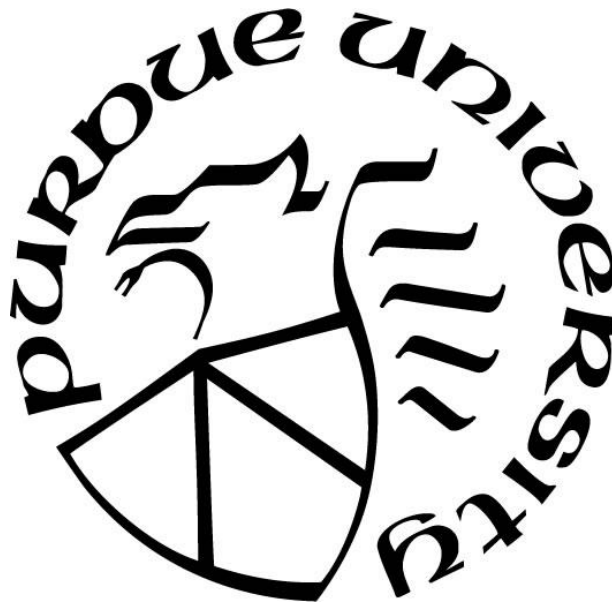# THE LANGUAGE OF ENGAGEMENT IN MATH INSTRUCTIONAL VIDEO TUTORIALS: A CORPUS-BASED STUDY

by

**Aleksandra Swatek**

**A Dissertation**

*Submitted to the Faculty of Purdue University*

*In Partial Fulfillment of the Requirements for the degree of*

**Doctor of Philosophy**

Department of English

West Lafayette, Indiana

May 2019

# THE PURDUE UNIVERSITY GRADUATE SCHOOL
# STATEMENT OF COMMITTEE APPROVAL

Dr. April Ginther, Co-Chair

    Department of English

Dr. Tony Silva, Co-Chair

    Department of English

Dr. Shelley Staples, Committee Member

    Department of English, University of Arizona

Dr. Bradley Dilger, Committee Member

    Department of English

**Approved by:**

    Dr. Manushag Powell

        Head of the Graduate Program

*To my family*

# ACKNOWLEDGMENTS

*I came to explore the wreck./ The words are purposes. / The words are maps./ I came to see the damage that was done / and the treasures that prevail. / I stroke the beam of my lamp / / slowly along the flank / of something more permanent / than fish or weed // the thing I came for: / the wreck and not the story of the wreck / the thing itself and not the myth*
from "Diving into the Wreck" by Adrienne Rich

Writing this dissertation was a tremendous undertaking, and its successful completion would have not happened without the support of my committee: Dr. April Ginther, Dr. Tony Silva, Dr. Shelley Staples, and Dr. Bradley Dilger. Thank you April for providing me with feedback and guidance on the project, as it developed over time and scope, conversations and support when I needed it. You inspire me to be a better person and a better scholar. I would also like to thank my other co-chair – Tony – whose wise advice and support made me a more confident writer and a researcher. Shelley, I am grateful for introducing me to the world of corpus linguistics and always being supportive and providing constructive feedback, showing me where I can grow. Bradley, I have learned so much from you in our Crow nest and beyond, including the importance of *poor quality* first drafts.

There is also my friend-family, coming from graduate programs at Purdue University and University of Maine. Ashley J. Velázquez, you provided me with strength and courage when I lost hope and energy. I will be forever grateful to the higher forces putting you in my way, with out friendship challenging us both to grow and become better. You are my Mer and Christina, and I am always counting myself lucky to see you persistently excel as a scholar, and be inspired. N. Claire Jackson, thank you for endless conversations about composition scholarship, new episodes of *Grey's Anatomy* and *Vampire Diaries*, seltzer, cats, and the daily adventures of writing teachers; thank you for believing that I can do it. We continue to be similar, but different. Sarah Cook, babe/mental sister, thank you for keeping me grounded and letting me feel I have a family here; thank you for all the little and big reminders of the world of outside of academia, being patient and kind, and supportive. Thank you for all the poetry and sending me lifelines (the opposite of dead-lines).  Words cannot describe how grateful I am to be counted as your friend

and your sister. Aleksandra Kasztalska, thank you for being the other A in the A-team. Thank you for all the tea-infused Skype conversations that fueled this project. Thank you also to Aleksandra Judejko, for the countless hours of discussion about life, courage, and mental health over all those years.

I am grateful to MIT OCW and Khan Academy for making their content available through creative commons licenses. This work would not have also come to be without the work of Python programmer: Adriana Picoral Scheidegger, Peter Collingridge, and Dmytro Pryimak. My family: my mom and dad, brother Karol and his wife Justyna, aunt Grażyna and uncle Krzysiek have always supported me in all the possible ways. I'd like to thank my grandma – Krystyna Dziedzic - for her warmth and support all these years. Thank you my Polish friends: Monika Krajewska-Woźniak (and family!), Dorota Buśko, Iza Krawczyk, Sylwia Koć and Olka Demolka for cheering for me from afar. Thank you to Dr. Magda Leszko, for your words of wisdom and your positive can-do energy (and for giving me my cat!).

I am grateful to the UMaine English Department faculty – Dr. Dylan Dryer and Dr. Charlsye Diaz for introducing me to the field of rhetoric and composition, and research methods. I kept learning from you long after I left the program. All of my UMaine friends - thank you! Marta and Mario Potoccy, you were my second Polish family, kept me using Polish and shared your home with me when I needed it. Thanks to my UMaine friends: Liz Maliga, Wes McMasters, James Brophy, Chad Van Buskirk, Erin Workman, Jesse Priest (by proxy), and Aaron Pinnix, Beata and Artur Palacz, Amam On / Stasiu, and countless others. My cohort at Purdue: Yiyang, Xiaorui, Hanyang, Kenny, Kai, Negin, for creating a vibrant academic mini-community. Thank you ESL GO! and all the SLS graduate students. Thank you, Tyler, for endless walks on the trail and endless conversations. I am also grateful for being part of the Corpus and Repository of Writing research group, which helped me grow as a researcher, writer, and a presenter.

Purdue English Department has been an intellectual home to me for the past five years. I am grateful to Dr. Margie Berns and Dr. Dwight Atkinson for sharing their knowledge in seminars and conferences. I am lucky to have been in those zones of proximal development. Thank you to the Oral English Proficiency Program: Dr. Nancy Kauper, Rochelle Hines, Kelley Farrell, Beth Lageveen, Dr. Mark Haugen, and all the tutors and instructors for teaching me what a good workplace looks like; the work I did in OEPP was the direct inspiration for this project.

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

# ABSTRACT

Author: Swatek, Aleksandra, M. MA
Institution: Purdue University
Degree Received: May 2019
Title: The Language of Engagement in Math Instructional Video Tutorials: a Corpus-Based
    Study.
Committee Chair: Ginther, April; Silva, Tony.

This dissertation investigates the linguistic features of engagement in spoken academic online
and face-to-face instruction in mathematics on two platforms: Khan Academy and MIT
OpenCourseWare. In particular, the study analyses involvement features (personal pronouns and
deixis) and interactional features (response elicitors, direct hypothetical reported speech). Using
corpus linguistics methodology and register analysis framework (Biber &Conrad, 2009), I
investigated normed frequency of occurrence for these features and multi-word expressions
which contain them to reveal patterns of use. Additionally, I investigated the function of these
features in concordance lines to reveal their use to engage audience in the learning process. The
findings of this study suggests that Khan Academy instruction in mathematics relies on using
conversational and academic spoken features similar to those found in the MIT lecture corpus,
including frequent use of personal pronouns (especially *we)*, and response elicitors (*right?)*. The
format of online video instruction also elicits more use of spatial deixis to draw attention to the
elements on the virtual board. The findings of this exploratory study add to the growing literature
on language used for educational purposes in online environments, especially the online
academic spoken discourse.

# CHAPTER 1: INTRODUCTION

**The Need for the Study of Spoken Discourse in Nontraditional E-Learning Environments**

Spoken academic discourse has been at the center of academic language research for over 20 years (Benson, 1989; Flowerdew, 1994; Johns, 1981; Richards, 1983; Waggoner, 1984). The centrality of the spoken discourse role in education is not surprising, since the history of education is inextricably linked to oral tradition (Jones, 2007; Ong, 1982). The format of a lecture has been chronicled as the basis of Ancient Greek education in the form of Socratic method, or in the Middle Ages as an alternative to expensive transmission of knowledge in the written form (Eisenstein, 1997). In the 21$^{st}$ century, however, the Internet technology has reshaped the genre even more than the radio and TV formats did. The development of Internet technologies allowed for on-demand, asynchronous access to educational content for anyone with adequate technology. It not only further removed the need for the speakers and listeners to be sharing the same time and space of the communicative event, just like TV and radio did, but it allowed anyone to become an instructor and produce their own content.

A new type of instructional form appeared – a video tutorial – concurrently with the rise of the YouTube platform as the central hub of participatory culture (Burgess & Green, 2009). This platform allows anyone to create and share video tutorials on a variety of topics. One of the people who decided to share his knowledge through the platform was Salman Khan, the founder of Khan Academy. Khan started uploading his math instructional video tutorials to YouTube in 2006 and has been producing more and more lessons ever since. As he joined the ranks of other YouTube educators who create educational materials accessible to anyone, other trends were on the rise, fueled by the same exponential growth in performance and utility of the Internet – the movements for open educational resources (OER), massive online open courses (MOOCs), and online instruction at higher education institutions.

The growth of online education ventures increased access to opportunities for students who otherwise might not be able to take time to travel to an institution or who may not afford college courses. Nowadays, a number of for-profit and non-profit organizations offer online courses to match the needs of traditional and non-traditional students as they prepare to enter or change their positions on the job market. Some of these courses are designed or endorsed by

higher education institutions and are offered for college credit, such as the ones on EDx platform. Others, such as Khan Academy, offer free self-paced courses that can supplement face to face and online instruction or support tutoring in AP (Advanced Placement) courses.

With the rise of online education offerings come questions about how online teaching compares to its traditional, face-to-face classroom counterparts. Such comparisons can be made by investigating a number of factors that have impact on learning. One such factor is language use, which can investigated by conducting a register analysis (Biber & Conrad, 2009). A register is defined as a language variety characterized by its situation of use. Analyzing a register requires describing its situational characteristics, linguistic features and their function in the particular context to arrive to conclusion. The analysis of online education requires an in depth examination of the non-traditional situational characteristics, which have not been explored at length in any previous register analysis studies. In fact, previous studies on spoken and written registers of university language have been conducted on corpora of texts collected in physical university spaces across the USA. For example, Biber, Conrad, Reppen, Byrd and Helt (2002) used TOEFL 2000 Spoken and Written Academic Language (T2K SWAL) corpus to compare spoken and written university language collected in a variety of instructional contexts. This corpus was composed of 674 texts from four universities: Northern Arizona University, Georgia State University, Iowa State University, and California State University – Sacramento. T2K SWAL, alongside Michigan Corpus of Academic Spoken English (MICASE) and British Academic Spoken English corpus (BASE) are the basis of hundreds of linguistic studies on oral academic English. All of them, however, were collected at the beginning of 21$^{st}$ century, and did not include language used in online academic instruction. Other university language studies rely on smaller, self-collected corpora (e.g. Fortanet, 2004; Crawford Camiciottoli, 2007) which are more manageable to analyze by one researcher or a small team, but do not represent large academic domains. However, despite the continued interest in the use of English for academic purposes, there has been no previous research on academic video tutorials made solely for the online environment for the general public.

**Motivation to Study Khan Academy Math Instruction Engagement Language**

Unlike the previous research on spoken academic discourse, my dissertation focuses on the under-researched register of spoken online video tutorials used for academic instruction.

Additionally, my corpus represents instruction done by Salman Khan, who is not a teacher by profession, but rather developed his teaching through interaction first with his family members and then with the wider audience on the Internet.  In this study, I explore Khan's use of linguistic features typical of engaged discourse in instructional videos. Engagement is a term that is used differently by researchers (Biber, 1989; Hyland, 2005), depending on their disciplinary affiliation. In my study, I use language of engagement to denote discourse that uses linguistic features of interaction, involvement and stance. Interactive discourse is typical of conversation, where two or more speakers take turns while discussing a topic (Morell, 2004). Involvement denotes the speaker's mode of participation, from a very casual to intense engagement in the topic or relationship between participants (Katriel & Dascal, 1989; Tannen, 1985).

My dataset is a self-collected corpus of transcripts extracted from Khan Academy[1] short video tutorials which, in sequence, provide instruction in various mathematical fields. As a comparable corpus, I used a self-collected MIT Mathematics (MITM) corpus of lecture transcripts. Since these lectures were recorded during the regular academic year at MIT, they provide a good comparable corpus of face to face instruction. Using these two corpora, both of which include instructional spoken academic discourse of high-school and college level mathematics, allows me to explore the differences and similarities in the use of engagement features in spoken academic discourse in both modalities.

Research on the spoken discourse of Khan Academy is important since it is one of the most popular open educational resource (OER) platforms online. Unlike other platforms, such as FutureLearn, EdX, Coursera, Khan Academy is not a Massive Online Open Course platform. On the contrary, the focus of Khan Academy is individual attainment and pacing of learning, rather than following a prescribed syllabus in a scheduled timeframe. This format is very popular among Internet users. In 2017, Khan Academy reported serving almost 12 million learners every month (Khan Academy Report, 2018). The impact of Khan's mathematical instruction exceeds the impact of any brick and mortar higher education institutions, or that of individual math instructors who only instruct the students that are physically present, because the viewers are not limited by time, space, or finances. Moreover, the effect of Khan's instructional videos, which are slowly being translated by volunteers into 27 languages, has global impact on the learners from all around the world. Indeed, Khan Academy's mission is to "provide a free, world-class

---

[1] "Note: All Khan Academy content is available for free at (www.khanacademy.org)".

education to anyone, anywhere" on the condition they have access to technology that allows them to access the platform.

Another important reason to conduct this study is to explore instructional speech of a non-academic affiliated instructor in an online environment. Salman Khan is not a math professor, but he did earn a B.A. degree in mathematics from MIT in the 1990s.Thus, he can be considered a proficient mathematician (see Chapter 3), but he has no formal training in education, or in teaching courses at high school or university level. Despite having benefited from prestigious higher education institution's offerings, his own teaching does not try to recreate that pedagogy. Instead of focusing on emulating a high-school or university classroom, Khan Academy is student-centered and tailored to individual students' learning goals and pace which is embodied in the design of the platform itself. This design, which allows each student to collect points, progress at their own pace and watch videos on particular academic concepts followed by exercises linked to these videos, is reflective of Khan's belief in *mastery learning* (Block & Burns, 1976). This educational psychology concept is based on the premise that students should move on from a simple concept to a more complex one only if they mastered the first one. As a result, the Khan Academy platform lets the students practice until they master the concept presented in the video. Thus, the videos are a crucial component of the pedagogy: they allow the student to focus their attention and practice on smaller units one at a time, unlike lectures, which require prolonged attention.

The need for this study lies in the scarcity of research on language used for Khan academy instruction in a technology-enabled student-centered pedagogy in math. The investigation of this computer mediated communication (CMC) register also provides unique insight into language that is used in visual-heavy teaching, because math instruction typically relies on using a board for writing. In the case of Khan academy, it is an interactive board on which Khan draws and writes and/or displays photos or videos for the lesson. His style of teaching, one in which he does not record or show his face, is specific to online teaching context. Research on how engaging different online videos are suggests that Khan-style videos are more engaging than voice-overed PowerPoint recordings or other type of screencasts (Guo, Kim & Rubin, 2014). Despite the fact that researchers from psychology and computer science have been calling for more personal and conversational style recordings (Clark & Mayer, 2011; Guo, Kim & Rubin, 2014), there has been little research on how such style is realized linguistically.

**A Register and Corpus-Based Approach to Discourse Analysis**

This dissertation project is based on the comparison of two large datasets – *corpora*, both of which are described in Chapter 4. It is important to first introduce the term *corpus*, as well as corpus-based research approach to discourse analysis to understand the framework in which this project is set in.

A *corpus* is a large, balanced and representative, collection of natural occurring texts (Biber, Conrad, Reppen, 1998). These texts can be written (e.g. college essays, newspaper articles, job ads) or spoken (e.g. transcripts of academic lectures, speeches, conversations). Apart from the modality of the represented register, the size of the corpus, representativeness, and balance are crucial for the validity of the results produced with the use of corpora (McCarthy, 1998). Because the particular advantage of corpus studies lies in the opportunity to explore very large language samples, the analysis of such large corpora overwhelmingly depends on computer analysis. Over the years, researchers either opted for creating their own programs for analysis of corpora, or relied on software such as AntConc (Anthony, 2019), WordSmith, SketchEngine (a platform), or more recently LancsBox.

Corpus linguistics methodology is particularly useful for empirical research of language variation and use (Biber, 2015). Researchers have used it to explore patterns of language in the corpora, noting if anything is typical or atypical, expected or unexpected, in relationship to intuitive or research-based predictions of previously explored linguistic forms (Conrad, 2002). An important element of corpus research is the recognition that different registers in one language, such as English, can be characterized by different patterns of language use – a variation in language use. This means that depending on the type of register, it could on average include more or less modal verbs, mental verbs, or any other type of linguistic feature (see *Longman Grammar of Spoken and Written English* [Biber et al. 1999]). Such language variation is seen from the perspective of corpus linguistic research as systematic, and thus can described using quantitative methods (Biber, 2015). The results of empirical analysis, then, can be investigated by researchers functionally in light of the situational characteristics of use. This type of research, in which linguistic variation is investigated in terms of function of use is defined as corpus-based research (Biber, Conrad, & Reppen, 1998).

Corpus-based and corpus-driven are two different approaches to conducting research on a corpus. Corpus-driven studies do not rely on pre-determined linguistic features and structures

from linguistic theories, rather they try to theorize about language based just on the data included in the corpus. Corpus-based approach to linguistic research relies on the premise that linguistic forms and structures derived from theory are valid (Biber, 2015). Further, the investigation of systematic patterns in corpus-based research can rely on using pre-defined linguistic features as a starting point of investigation. This dissertation project relies on corpus-based approach to systematically analyze such a set of pre-determined linguistic features, which were chosen based on previous research on interactivity and involvement in academic discourse.

No other research approach would lend itself better to this project than corpus-based research, since analyzing such large datasets would take enormous effort without the use of computer software. The use of computer software is essential in detecting patterns of language use in large datasets, and it also automates a more fine-grained analysis of co-occurring linguistic features or even concordance lines.

## Overview of the Study

The research used in this dissertation is based on a combination of quantitative and qualitative approaches, typical of register analysis studies (Biber & Conrad, 2009). Two corpora of high school to college level math instruction – one from Khan Academy, the second one from MIT – were created to investigate the linguistic features used engaging the audience. These corpora are both based on transcripts of lectures from MIT and video tutorials from Khan Academy instruction. For both corpora, I also conducted a detailed situational analysis, which is an obligatory component of register analysis. Such situational analysis allows the research to interpret the differences in functions of linguistic characteristics.

The study consists of three steps: situational analysis, analysis of linguistic features (normed frequency, patterns), and functional analysis of the features. First, I analyze the situational characteristics of Khan Academy math instruction and MIT lecture instructions. The next step of the study involves analyzing the frequencies of the linguistic features used for engaging the audience in both corpora and more fine grained analysis of the collocations of these features and their use in context. The last step, is the functional analysis of the results of my linguistic analysis in the context of the situation of use.

## Research Questions

The research questions put forward in my study draw on the register analysis framework (Biber & Conrad, 2009). Questions 1 and 2 refer to the situational characteristics, which are explored in a qualitative (i.e. descriptive) manner based on secondary sources and primary research (i.e. observation and description of primary sources). Questions 3-4 combine an corpus-based exploration linguistic features in both corpora and their functional analysis.

1. What are the situational characteristics of online video instruction in mathematics on Khan Academy (KA)?
2. What are the situational characteristics of recorded face to face lectures on the MIT Open Course Ware website (MIT)?
3. Are there differences in the involvement linguistic feature use in the two corpora?
4. Are there differences in the interactive linguistic feature use in the two corpora?

## Outline of the Study

The first chapter of this dissertation is an introduction. Chapter 2 of my study provides literature review of previous research on spoken academic discourse. In that chapter I also present the construct of engagement, which I define as a combination of interactive and involvement features. Chapter 3 contains description of the purpose of situational analysis in the register analysis framework (Biber & Conrad, 2009). In particular, I analyze the situational characteristics of both corpora – Khan Academy Math tutorials and MIT math lectures. Chapter 4 focuses on the methods used in the study, in particular corpus-based discourse and register analysis. Chapter 5 presents the results and discussion of my findings. The last part, Chapter 6, provides a summary of the study and paths for future research.

# CHAPTER 2: LITERATURE REVIEW

This review of literature focuses on three research areas that inform my study: monologic academic discourse in face-to-face and e-learning contexts, register analysis framework, and corpus-based studies on engagement linguistic features. Little research has been done on the specific register under investigation in this dissertation: video tutorials made specifically for e-learning environments.

## Monologic Academic Discourse

Monologic academic discourse is one of the oldest method of teaching in higher education. The term monologic refers to a communicative situation in which there is one person speaking to an audience, and there is very little to no response in the audience. The most common type of monologic academic discourse, the one with the longest – ancient – tradition, is the lecture. Instructors lecture to communicate foundational knowledge of the discipline to the students (Flowerdew, 1994; Flowerdew & Miller, 1997; Lee, 2009; Thompson, 1994). Depending on the lecturing style, the instructor might engage the audience in a dialogue, but discussion is not typical of lectures. The term *monologic* can be deceptive, as number of scholars have pointed out that such discourse still includes linguistic features that attempt to engage listeners and writers. Hoey (1994) defines monologue as "written or spoken, may be regarded as a dialogue in which the reader/listener's questions or comments have not been explicitly included but which retains clear indications of the assumed replies of the reader" (p. 29). The presence or absence of audience does not affect the dialogic potential of language use (Biber et al., 1999, p. 213). Since lectures are inextricably linked to traditional model of education, they came under criticism for the underlying assumption that students learn through passive listening (Crawford Camiciottoli & Querol-Julian, 2016). Still, lecture persists as the most often used method to deliver educational content in large classes.

Lecturing has been considered the key instructional method at universities for most of the existence of the institution (Brockliss, 1996; Flowerdew 1994). In the field of broadly defined applied linguistics and ESL, the research on lecture has often been motivated by the need to understand how second language learners are challenged and overcome challenges posed by that

instructional method (Barbieri, 2015). The research on listening comprehension focused in the past on discourse organization of lectures (Young, 1994), discourse organizing markers (Jung, 2003; Thompson, 2003), discourse signaling (Nesi & Basturkmen, 2006), and importance (Deroey, 2015). Research on lectures in other fields explores it as a transmedial pedagogical form that bridges oral and written communication through electronic and digital means (Friesen, 2011). Another line of inquiry explores the reimaging traditional lecturing to address the *active learning* pedagogy model, made known by the example of Eric Mazur (Lambert, 2012). Mazur, a physics professor, revised his traditional lecturing method into active learning, in which lectures are interrupted by comprehension checks, and students teach each other through planned and intentional inclusion of discussion, problem solving, collaboration and cooperation activities (Prince, 2004).

Despite being delivered in the spoken mode, lectures are not as spontaneous as everyday conversations (Flowerdew & Miller, 1995) as they require a degree of planning, and offer potential for rehearsing and repetition. Csomay (2000) categorizes lectures a "hybrid" register, since they are informational in purpose, but delivered on-line[2]. Lectures can have different levels of interactivness (Csomay, 2007) and styles (Dudley-Evans & Johns, 1991), with certain lecturers engaging their audience in minimal dialogue and others preferring purely monologic style (DeCarrico & Nattinger, 1988; Nattinger & DeCarrico, 1992). In the past, lectures were delivered using more formal language, but in recent years researchers started to notice that there are more informal and conversational markers used by instructors (Bamford, 2005; Morell, 2007). This less rigid style of lecturing allows for more interaction between the speaker and the audience (Chang, 2012; Crawford Camiciottoli, 2004, 2008; Csomay, 2007; Dudley-Evans, 1994).

The linguistic characteristics of lectures have been studied in the field of English for Academic Purposes (EAP), which emerged from the need to describe register/genre of education in English (Biber & Conrad, 2009; p. 3). Rather than a set of general English skills, the students who receive their education in English need to be familiarized with the specific text types, spoken and written, which they will encounter in high school and college. As the field of EAP

---

[2] Since the term on/-/line can be understood in two different ways in the context of my dissertation 1) live production 2) Internet environment, I will use different spellings two differentiate between the two. The term on-line to describe the method of production – meaning face to face, in real time. The term online will be used to discuss the Internet environment.

developed over the past three decades, the discourses that are used for education changed as more technology entered educational contexts. The short monologic academic speech, such as the one used in Khan Academy math tutorials, has emerged with the onset of online education.

Traditional lectures are based in a context in which students share space and time with the instructor. This can exclude non-traditional, lifelong learners from being able to participate in the educational environment (Hicks, Reid, & George, 2001). E-learning has emerged to challenge the traditional methods of instruction and offer access to education to a larger pool of prospective learners, utilizing Internet technologies to increase opportunity for learning. Lectures are being given a second life with the rise of massive online open courses (MOOCs) or OpenCourseWare platforms where universities share their course materials for free (such as MITs). With technology facilitating and preserving the live lecture performances of university instructors, the study of the spoken academic discourse across different contexts and with international audiences is gaining new importance. Despite the popularity of open education resources (OER), and especially the resources made for online courses, little research in EAP has been done on the linguistic characteristics of that register.

**Spoken Academic Discourse of Mathematics**

Mathematics register has been described in previous studies using systemic functional linguistics methodology (Lemke, 1989, 2003; O'Halloran, 1999, 2015; Veel, 1999, cited in Schleppegrell, 2007). Pimm (1987) examined the language used in mathematics classroom, specifically the oral and written modes of communication, using recordings of classroom sessions. His qualitative analysis which draws on Halliday's notion of register interrogates how language functions in mathematics between instructors and students. He explores the idea that mathematics is to some extent a foreign language, examines students' mathematical talk (talking to self vs. talking to others), and the power relations embodied in language used between teacher and students. Schleppegrell (2007) conducted the most recent synthesis on research in mathematics language. She described two main register features: use of multiple semiotic systems (symbolic notation, oral language, written language, graphs and visual displays), and grammatical patterns (technical vocabulary, dense noun phrases, *being* and *having* verbs, conjunctions with technical meanings, and implicit logical relationships).

The literature on language in mathematics teaching has been mostly focused on K-9 grades (Gerofsky, 2004, Morgan, 2006). Speer, Smith, and Horvath (2010) conducted a literature review on collegiate mathematics teaching which showed that mathematics education at post K-12 level is sparse, especially in terms of empirical research. In their analysis of previous literature, they pointed to the need for more research on: selecting and sequencing content in teaching, motivating specific content choice to help students orient themselves in the knowledge field and connect concepts, asking questions and using wait time, representing mathematical concepts and relationships. All of these elements are directly related to structuring lectures in which these elements come together and might pose a challenge to the students.

Artemeva and Fox (2011, 2012) explored the genre of teaching undergraduate mathematics. They defined the main pedagogical genre as *chalk talk* as "writing out a mathematical narrative on the board while talking out loud (Artemeva & Fox, 2011, p. 345). In their study of 50 math instructors from all over the world, they noticed that all teachers in their local context used the same pedagogical genre (*chalk talk*), which is characterized by: running commentary, metacommentary, gesturing for relationships, reference to textbook and notes, use of rhetorical questions for transition, asking questions to students (p. 356). The interviewed teachers emphasized how the *process* of mathematics needs to be narrated in order to introduce students into disciplinary thinking, with running commentary and metacommentary being key elements of showing the process.

**Instructional Videos as Computer-Mediated Discourse**

Instructional videos can be considered a type of computer-mediated discourse (Herring & Androutsopoulos, 2015) and register, since the underlying function of use is teaching via networked mobile devices. As a type of computer-mediated communication (CMD), instructional videos might be difficult to categorize in terms of modality, since they combine written and spoken elements. Swarts (2012) suggests that instructional videos can rely on design principles used for print (framing the introduction and goals), while also utilizing the modality affordances by spending time on demonstrating steps.

## Engagement in E-learning and Multimedia Academic Discourse

The ever-growing presence of technological tools supporting teaching and learning put in question the role of instructors in education. The growth of technological solutions completely reorganized distance education programs, which in the past relied on mail, television, or telephone (Open University, 2006), now are almost entirely Internet-based.

Discussion of engagement in e-learning, and specifically effective learning through multimedia materials has been explored in cognitive psychology and in computer science. Richard Mayer (2014) formulated cognitive theory of multimedia learning to describe the constraints or opportunities for engaging in learning online. These constraints are particularly important for analyzing both of my datasets which are examples of multimedia-based teaching materials. Mayer's theory is based on three assumptions about how people learn through multimedia. The first assumption – *dual processing* – states that humans process visual and auditory cues separately. The second assumption – *limited capacity* – suggests that humans have a limited capacity to process information in each channel at the same time (Baddeley, 1992; Chandler and Sweller, 1991; cited in Mayer, 2014). Finally, the *active processing* assumption reveals that "humans engage in active learning by attending to relevant incoming information, organizing selected information into coherent mental representations, and integrating mental representations with other knowledge" (Mayer, 2014; p.47). These assumptions are relevant in terms of understanding how cognitive research can provide evidence of engagement in multimedia learning (including video instruction). Mayer's research on e-learning revealed that when the content is presented in a narrative, conversational format (using first and second person pronouns, instructor turns aimed at the listeners) rather than a more formal style, the students remember more of the content. This phenomenon has been coined as *the personalization effect* (Moreno & Mayer, 2004). This effect is at the basis of the rule that an instructional message should be designed to promote the feeling of physical presence of the instructor (Moreno & Mayer, 2004). The personalization effect is often achieved by using personal pronouns and addressing the listeners directly in the instructional materials:

P version: *You* are about to start a journey where *you* will be visiting different planets. For each planet, *you* will need to design a plant. *Your* mission is to learn what type of

roots, stem, and leaves will allow *your* plant to survive in each environment. I will be guiding *you* through by giving out some hints. (Moreno & Mayer, 2004; p. 733 ).

The results of research on personalized (conversational and narrated) content explanations demonstrates that students who receive personalized input performed better on tests than students who received a more formal narration (Mayer, 2003). In a similar experiment with use of *your* possessive pronoun instead of *the,* students obtained significantly higher scores on transfer than the non-personalized group (Mayer, Fennell, Farmer, & Campbell, 2004). This research on personalized multimedia learning, especially in terms of using personal language, provides evidence for cognitive benefits of using engagement language, with an effect size of 0.65 for improving transfer performance (being able to use knowledge in a new problem).

Engaging in video instruction has been measured using user data in computer science. Guo, Kim, and Rubin (2014) explored student engagement in online educational videos using data collected in four edX courses. They analyzed 6.9 million video sessions to explore the reasons why students engaged or disengaged from watching the instructional videos. They categorized the online videos into four styles: 1) typical classroom lecture 2) "talking head" shot of an instructor at a desk 3) digital drawing format popularized by Khan Academy and 4) a slide presentation (p. 1). As measures of engagement, they analyzed a number of variables: length of viewing sessions, drop off instances, and answering questions displayed after the video finished. Their findings reveal that 1) shorter videos are more engaging 2) videos with a more personal feel (informal setting emulating one-on-one meeting) keep students involved more than big productions, 3) Khan-style videos with continuous visual flow and commentary are more engaging than PowerPoint slides with voiceover. This research provides big-data analysis evidence on engagement measured through patterns of students' actions. These results partially support the use of Khan-style video for increasing students engagement in their learning process.

## Register Analysis Framework

Linguists have for a long time been interested in how the occasion of use evokes a certain language variety – a register. Speaking to an older person in a position of power in a formal meeting requires a different language than speaking to a child at a birthday party. How language varies between situations of use became the focus of first register variation studies (Ellis & Ure, 1969, cited in Biber & Finegan, 1994). The description of characteristics of such elements, or a

communication situation, for the purpose of sociolinguistic study had been informed by Malinowski's work on contextual theory of meaning (1923), followed by Firth's discussion of "context of situation" (1935, 1957). Hymes (1974) proposed the SPEAKING (setting and scene, participants, ends, acts sequence, key, instrumentalities, norms, & genre) framework for ethnography of communication, which became the basis or element of other frameworks used in sociolinguistics. With the development of the field of sociolinguistics, scholars were approaching register studies that used other characteristics for identifying registers (Basso, 1974; Biber, 1994; Brown & Fraser, 1979; Crystal & Davy, 1969; Duranti, 1985).

These frameworks are particularly important for understanding of how the situational characteristics shape a register. A person who uses language in certain situation repeatedly will with time develop linguistic resources similar to other people who use language in the same situation. The term *register* can be contrasted with other important sociolinguistic terms as *dialect*, a language variety dependent on the speaker's geographical (regional dialect) and social dialects (Biber et al. 1999, p. 17; Romaine, 2001). How language is shaped by the situation of use has also been explored in the field of rhetorical genre studies, with Miller's (1984) theorizing that genre "are dynamic rhetorical forms that develop from responses to recurrent situations and serve to stabilize experience and give it coherence and meaning" (p. 479).

Register analysis framework which is used in this study (Biber, 1994; Biber & Conrad, 2009), has its roots in the field of sociolinguistics whose focus centers on written and spoken language in use (Biber & Finegan, 1994). The work of pioneering scholars in sociolinguistics laid the foundation for recognizing the three component of register analysis: situational features, functions and conventions, and linguistic forms (Halliday, 1978; cited in Biber, 1994). Discussions around this model have centered on the questions of the nature of relationship between these elements - whether they were correlational or deterministic (Biber, 1994). Previous taxonomies of situational characteristics (Duranti, 1985; Hymes, 1974) and classification of registers (Biber, 1989), did not allow for comparison between levels of generality (e.g. formal vs. informal, novels vs. science articles, methodology sections of agronomy articles vs. linguistics articles). The question arose: how to decide at which level of generality a given register is and if these registers should be compared (Biber, 1994). The framework for register analysis was conceived as at attempt to describe the situation of register

use in order to allow for reliable comparisons across registers. This framework, fully developed in Biber and Conrad (2009), serves as the basis of my study.

The importance of situation of use is crucial for interpreting how a linguistic feature functions in a particular situation. Since the functions cannot be determined for the whole language in general, but need to be investigated in context of the register (Staples, Egbert, Biber & Conrad, 2015). Corpus-based register analysis can be helpful for lexical descriptions of registers, by considering how lexical items are used in two or more registers represented in corpora. The results of such studies can provide insight into how language differs based on the situation of use, avoiding generalizations associated with describing the language as a whole.

### Patterning: collocations and other multi-word expressions

Corpus linguistics can be used to investigate lexis, specifically a word and its co-occurrence with other words (Moon, 2010). Firth (1968) first drew attention to the patterns in which words co-occur with other words and how the language we use on the daily basis relies on formulaic expressions, rather than creative combinations (i.e. idiom principle). Baker (2006) defines a collocate is a word that occurs typically within the neighborhood, a span of another word (i.e. *good morning*, rather than *fine morning* or *splendid morning)*. For instance, the examination of the noun *chair* in British National Corpus reveals that it is often modified by adjectives *comfortable, swivel, high-backed, wooden*, while in the English Web corpus (*enTenTen08*), it is often modified by *folding, electric, comfy*. The difference across collocate lists depends on the corpora that are being used for research. Thus, even quite simple searches of collocations in corpora that represent two different registers can provide important information on how the same words or phrases are used in a similar or different manner, revealing differences in register.

The research on collocation has a long tradition in lexicographical and phraseological studies, where structure and co-occurrence are used in tandem to examine word meanings (for discussion see Sinclair, 1987; Stubbs, 2001). Greaves and Warren (2010) trace the development of collocation studies to the first corpus-driven studies of collocation patterns conducted by Sinclair's team with the use of a computer in 1960s. Their research revealed three findings which were a result of resistance to the primacy of grammar over lexis at that time. The first finding of their report was the concept of meaning creation through lexis, with grammar managing meaning

creation (Sinclair et al., 1970, 2004). To exemplify that concept, Sinclair used the example of *just a minute* as a phrase in which *just* and *a* are not used in because of their grammatical category, but because of the meaning imposed by the use of *minute, second, moment, sec* (Sinclair, 2004, xxvi). He compares the function of grammar to the function of intonation in speech (as presented by Brazil, 1995): a text without intonation does not have meaning. As you assemble intonation and text, the meaning is constructed in the act of articulation. The second finding was that meaning is created through co-selection of words and, third, the formulation of the idiom principle

> The principle of idiom is that a language user has available to him or her a large number of semi-preconstructed phrases that constitute single choices, even though they might appear to be analyzable into segments. (Sinclair 1991, p. 110)

Collocation patterns can be examined in concordance lines, where the right or left context of the collocation provides information on the discourse level function and meaning of collocation (see the analysis of now in Swales and Malczewski, 2001). Some concordance programmes, for example SketchEngine or WordSmith, allow the research to display and sort the list of words that appears within 1 to 5 words to left or right from the keyword (the span depends on the researcher's choice) to examine extended patterns. The ability to examine co-occurrence is crucial to the examination of the function of the collocation in the register under examination.

Examining collocations is just one aspect of pattern examination. Another framework of pattern analysis relies on examining extended collocations in context. The formulation of the idiom principle (Sinclair, 1991), coincided with research on fixed phrases. There is abundant empirical evidence that speakers and writers produce semi-preconstructed phrases. Pawley and Syder (1983) note that native speakers of language use "sentence stems", fixed phrases in terms of form-meaning pairings, to construct fluent speech. Rather than constructing novel form-meaning pairings each time they speak, native speakers rely on fixed phrases representing culturally recognized concepts (Ellis, 1996; Siyanova-Chanturia & Martinez, 2014; Wray & Perkins, 2000; Wray, 2005). Cortes (2004) found that the use of multi-word units (such as collocations and fixed expressions) can be examined as features of proficient language use by writers. Hyland (2008) provided evidence that writers in various discourse communities are proficient users of multi-word expressions typically found in the writings belonging to their

academic disciplines, showing that the proficient use of multi-word expressions is a marker of participating in a given discourse community. Formulaic expressions, in a cognitive framework, are said to be stored and retrieved by the speaker from memory, rather than being generated and analyzed with more mental effort (Wray, 2005). However, research is lacking on the effects of these phrases on listeners, particularly in terms of their processing by audience (Siyanova-Chanturia & Martinez, 2014).

Recurrent phrases, such as collocations, binomials (*thick and thin*), phrasal verbs (*keep up with)*, idioms (*it's raining cats and dogs),* and lexical bundles (*the thing is)* can all be categorized as multi-word expressions (Siyanova-Chanturia & Martinez, 2014). In corpus-based studies of multi-word expressions, researchers might set out with a predetermined lexical expressions, or decide to analyze all possible sequences in a corpus (Gray & Biber, 2015). Multi-word sequences can also be studied across two or more registers to explore register-dependent variation of their use (Gray & Biber, 2015). Corpus-driven studies of multiword expression depend on frequency and dispersion as the variables that determine which expressions need to be analyzed. Altenberg (1998) was one of the first studies to examine sequences of 3 consecutive words (i.e. *this is such,* or *a great example,* or *of an 3-gram;* so called 3-grams) which occurred at least ten times in the *London-Lund Corpus,* describing them in terms of grammar and function. A more robust study of multi-word expressions was part of the *Longman Grammar of Spoken and Written English* (Biber et al. 1999), in which the researchers analyzed 4 to 6 gram *lexical bundles* in conversation and academic writing.

Identifying multi-word expression is just the first step of phraseological studies of corpora. Quantitative-level frequency-driven analysis in this research study relies on detecting patterns of language use (Hunston, 2010) and analyzing their use in context (Sinclair, 2004). Another level of analysis is describing the multiword expression discourse functions (Biber & Gray, 2015). There are two widely used frameworks of multi-word expression, specifically lexical bundle, analysis – Biber, Conrad, and Cortes (2004) and Hyland (2008). The first framework (Biber, Conrad, & Cortes, 2004) divides lexical bundles into three main categories: referential, stance, and discourse organizing expressions. Hyland's (2008) framework divides the bundles into three groups research-oriented, text-oriented, and participant-oriented bundles. In both cases, the bundles were first selected based on the frequency and dispersion thresholds, and then analyzed in concordance lines to explore their functions.

**Concordance and Key Words in Context (KWIC)**

Concordance is defined as "a list of all of the occurrences of a particular search term in a corpus, presented within the context in which they occur" (Baker, 2006). Exploring concordance lines is an analytical technique that is particularly useful in generating and testing hypothesis about item use in their context (Evison, 2010). Patterns observed in even small number of lines can be tested out through corpus searches on a larger scale (Hunston, 2010). One caveat of this analytical technique is taking into consideration its limitation: the patterns that are used in natural language but are not used in the corpus, will not appear in results. Thus, there is limited generalizability of this concordance line examination, that, to some extent can be mitigated by use of large corpora. However, for features that occur very rarely in natural language, chances of recording the feature in a small or medium corpus are low.

Concordance lines can provide information about *typicality* of linguistic feature use (Evison, 2010), such as most frequent meaning and collocates (co-occurring items). Such relationship is more frequent that a chance occurrence. Analyzing multiple concordance lines is used to determine the discourse function of a pattern in a text. Concordance lines can have varying length, depending on the software use and the setting preference of the user.

**Linguistic Aspects of Engagement in Spoken Academic Discourse**

Academic lectures are marked by the presence of involvement and interaction features (Hyland, 2009) that work together to make the academic spoken discourse engaging. These features decrease the distance between the speaker and the audience, reducing the affective gap between the speaker and listener. Student engagement can be understood as a factor that promotes participation (Barbieri, 2015) and attention. Even though research has not established direct connection between students' engagement understood as participation, and involved production, in my dissertation I theorize that interactional and involvement features work together to engage audience in the learning experience in the online instructional video tutorials.

The division of the features into involvement and interaction in the literature is not uniform and defined in the same manner by scholars working in these areas. Hyland (2005), who works in the tradition of discourse studies, delineated the building blocks of interaction as *stance* (hedges, boosters, attitude markers, self-mentions) and *engagement* (reader pronouns, directives,

questions, shared knowledge, personal asides). While a data-driven multidimensional analysis on academic language conducted by Biber (1988) placed involved production on the opposite side of informational one. While there is not one framework that is widely used to investigate engagement in spoken academic discourse, there are overlaps in how researchers analyze single features or group of features to examine involvement and interaction in academic discourse.

**Involvement and Involvement Features in Spoken Academic Discourse**

The research on involvement as a feature of spoken discourse had been explored in the fields of discourse studies, sociolinguistics, and cognitive-functional approaches (Ädel, 2012). In the first of these disciplines involvement has been described as a feature of oral language (Tannen, 1982, p. 3), in contrast with the detachment, in terms of context and shared time-space, in writing. Tannen (1985, p.125) characterized involvement as a spectrum on which people express their need to be close to their audience, to feel connected, or are able to distance themselves. Tannen (2007) explored the relationship between linguistic form and involvement, by listing strategies that are markers of involvement, for example repetition of phonemic, syntactical and lexical elements within and across turns; the use of detail and imagery; use of reported speech or constructed speech. Katriel and Dascal (1989), who also work in the same field as Tannen, defined involvement as the speaker's emotional participation in a discussion or in the speaker's own monologic discourse. Chafe (1982, p. 45) used the term involvement in his discussion of speakers who communicate in a manner that shows their emotional and experiential participation. The term "experiential participation" is quite vague and built on the dichotomy between formal written language and speech, with only the latter including the speakers sharing the communicative experience at the same time.

In his work on discourse strategies from a sociolinguistic perspective, Gumperz (1982) defined *involvement* as the willingness and ability to initiate and sustain verbal interaction. In the cognitive-functional approaches to language use, the term *intersubjectivity* has been used to signal the awareness of an audience (Ädel, 2012). The term intersubjectivity stands in contrast to the term *subjectivity*, which denotes the speaker's attention to their own self (Traugott & Dasher, 2002). Intersubjectivity can be expressed explicitly and implicitly. The explicit expression can be achieved by the use of addressee honorifics and tag questions, i.e. *okay?*; implicit expression rely

on more subtle signals such as hedges *well, okay* (Traugottt & Dasher, 2002, p. 22; cited in Ädel, 2012).

Another approach to investigating involved production comes from multidimensional studies of linguistic variation in speech and writing (Biber, 1988). In analyzing groups of linguistic features which co-occur together, the feature on Dimension 1 were interpreted as marking involved production on one end of the spectrum and informational production on the opposite side. The involvement vs. informational dimension was named inductively, based on the analysis of features that co-occur together. Biber (1995) characterized involved production as a functional category that is typical of "content with an affective, interpersonal focus" (p. 145) produced in on-line conditions. The result of his multidimensional analysis, to some extent, supports Tannen's (1985) and Chafe's (1982) definitions of involvement as a spectrum that to some extent reflects the proximity of affective and emotional participation of spoken discourse and detachment typical of writing.

Involvement has also been explored in educational studies, in which involvement is related to positioning theory, which is used for discursive identity construction and interpersonal relation analysis (Hirvonen, 2016). Teachers, just like any person, act and speak to impact and organize social structures in order to position themselves against the others – in this case, students (Harre & van Lagenhove, 1999). Techers use language forms to relate to the students to create a social arrangements that might recreate the physical ones – with students facing the teacher who leads the class from under the blackboard, detached and commanding (Wagner & Herbel-Eisenmann, 2008; Mesa & Chang, 2010). However, teachers can also use language forms to align the students with the instructor's voice - or limit the choices students might have in expressing themselves (Mesa & Chang, 2010). Their positioning, then, serves as an involvement strategy by showing their own positioning and allegiance towards the students.

Involvement as an interpersonal dimension of academic discourse has not been studied in a systematic manner (Barbieri, 2015). Previous studies have focused on single linguistic features that are also typical of conversational discourse, such as imperatives (*let's, look)*, addressing the audience directly *you guys*, and discourse markers *okay* and *so* (Swales & Malczewski, 2001). Mauranen (2001) found similar markers in her study of discourse-reflexivity in the MICASE corpus: addressing the audience (*you guys)*, use of personal pronouns *(you)* and imperatives. Csomay's (2006) study on American classroom discourse presented evidence for academic

teaching containing a mix of linguistic features of academic prose (especially *content and information production)* and of conversational features (*involved production)*. Biber's (2006) multidimensional analysis of academic registers revealed that there are disciplines-specific differences in terms of narrative and non-narrative orientation of teaching, with humanities and social sciences relying more on narrative elements, and natural sciences relying on non-narrative features. Csomay (2005) found that based on the level of instruction (introductory to advanced course), the teaching in graduate courses was characterized by the most interpersonal involvement (i.e. directives and first and second person pronouns), stance, and orality. The most comprehensive analysis of academic lecturers' involved production was conducted by Barbieri (2015), who studied linguistic features of involvement across disciplinary domains, level of instruction and class size. Her operationalization of involvement encompassed features which included interactivity markers (questions, confirmation checks, response tokens) and features considered by other scholars as involvement markers (discourse markers, directives, stance, pronouns and other). She found that situational factors had a "relatively weak effect on linguistic marking of involvement in classroom discourse" (Barbieri, 2015, p. 168), but that involvement is pervasive irrespective of the context in American academic instruction.

### Personal Pronouns Spoken Academic Discourse

Personal pronouns are essential in intersubjective communication (Benveniste, 1966; Rounds, 1987). The use of personal pronouns in spoken academic discourse has been linked to how speakers and writers imagine their audience (Fortanet, 2004). In the context of a classroom, personal pronouns are markers of the student-teacher relationship, involvement, engagement, and allow the participants to locate themselves in the ongoing conversations in the classroom setting (Friginal et al., 2017) . The use of these pronouns also helps the speakers build shared knowledge of people, objects, and entities (Carter & McCarthy, 2006).

One of the widely cites studies of personal pronoun function is Kamio's (2001) theory of information territory. He linked the use of *I* and *we* to conversational space, marking proximity to speaker's territory (*I and we)* or positioning the hearer in the hearer's territory *(you)*. The distance between the speaker and hearer, as enacted by the uses of the personal pronouns, signals how the speaker conceptualizes the degree of closeness with the audience/hearer. Chafe (1985) pointed out that unlike academic writing which tends to detach the writer from the

audience/readers by relying on the passive forms, speakers tend to refer to themselves very often. This creates "experiential involvement" (p. 45), as the speaker and listener share time and space, and the speaker marks statements with his or her own mental processes.

The distribution of the pronouns in university lectures varies depending on the context of the study. The widely cited Rounds (1987a, 1987b) studies reported that *we* was the most frequently used pronoun in her corpus of math lectures led by TAs and ITAs at the University of Michigan. However, there was variability across instructors. The TAs who used *we* in their speech three times as often as others, on average, were also the ones who got higher teaching evaluations, fewer complaints, and good evaluations from supervisor (Rounds, 1987a). Fortanet (2004) used MICASE, a corpus of transcripts collected at the University of Michigan and made public in 2007, and came to contradictory conclusions. That study revealed that *we* was the least often used pronoun, which she attributed to changing discourse practices at university over time. Biber (2006) found that first person pronouns (*I, we*) are slightly less common in classroom teaching (~40 per 1000 words) than in other spoken academic discourses. Two studies explored the relationship between the class size and the use of personal pronouns. Cheng (2012) found that the use of pronouns is more frequent is small-class lecture endings, while Lee (2009) found that the use of pronouns *you* and *I* differed depending on the size of the class. In large classes the speakers tended to use *you* more, and in smaller classes they showed preference for *I.* Yeo and Ting (2014) studied the personal pronoun use in large-class introductions in lectures given in English at a Malaysian university. Similarly to Fortanet (2004) they found that pronoun *we* was used with least frequency, compared to *you* which was one of the most frequently used in lecture introductions. The instructors in that corpus also used *you* to refer to the audience more often than they used a general *you,* irrespective of discipline and class size. Lee and Subtirelu (2015) who examined two corpora of lectures, one from MICASE and one self-collected EAP classroom corpus, found that *you* was the most frequently used pronoun. The pronoun *I* was also more frequently used than the pronoun *we.*

In politeness theory, the use of personal pronoun *we* marks inclusion of the hearer in the speaker's discourse and is used to build and maintain rapport (Brown & Levinson, 1987). The use of *I* and *you* has the opposite effect, signaling detachment and distance. These two strategies of distancing and including are linked to negative and positive politeness strategies. The positive politeness strategy serves the purpose of building hearer's positive face linked to the need for

respect and being liked. The negative politeness strategies are speaker-centered, because they are used to minimize threats to the speaker's face (Brown & Levison, 1987; cited in Friginal et al., 2017).

The pronoun *we* can convey different types of reference, based on its context of use. It can function in two major roles: inclusive *I + you* and exclusive *I + they* (Haas, 1969; Spiegelberg, 1973). Rounds (1987a) noted that the exclusive use in mathematics discourse is especially common in the function of naming and defining terms, which students cannot participate in. For example,

*We say the function f of x . . . is differentiable . . at a point J . if its derivative exists here.*

In this situation, the *we* denotes the mathematic instructor and mathematicians as a discourse community. The same argument is made by Pimm (1989) and Rowland (1999), who approaches such use of *we* as a power imposition, in which the students are co-opted into an unspoken agreement (Muhlhausler & Harre, 1999) with the instructor's discourse community, without the opportunity to voice questions or signal misunderstanding. The use of *we*, in the mathematical discourse, shows teacher's power "masquerading" as solidarity (Rowland, 1999, p. 19). Pimm (1987) also suggested that mathematics instructors' use of *we* furthers their position as aligned with an imagined community of mathematicians, a powerful group. Such positioning might marginalize and exclude student voices. Pimm (1987) compared the use of *we* in the educational context of math classroom to echoing baby talk (i.e. we're getting you dressed, aren't we?), hospital discourse (i.e. we are going to take our temperature), editorial *we,* and school-context *we* (i.e. "Susan, we never bite our friends." Grenfell, 1977, p. 23, cited in Pimms, 1987, p. 109). All of these examples show unequal power-relations with the use of *we* patronizing and diminishing the voice of the less powerful party. Similarly, Rowland (1999) claimed that *we* traps students into agreeing with the instructor in a mathematics lesson, without allowing for questioning and disagreement. Other scholars (Bailey, 1982, 1984b; Rounds, 1985, 1987b) opposite effect of *we* on classroom environment and teaching effectiveness.

Other research points to the use of *we* as a unifying and motivating device, especially in politics. Steffens and Haslam (2013) found that politicians in Australian government who used collective pronouns *we* and *us* win more elections, than politicians who use more of *I* and *me.*

This suggests that the collective *we* can be a motivating device, inspiring to action, which could have a positive effect on any audience.

What is important to point out is that none of these studies looked at other dimensions of the pronoun use in conjunction with prosody, sociocultural context, and body movements to decipher the inclusive/exclusive role of personal pronoun *we.* It might be the case that the use of the pronoun *we* can be read by students as encouraging or discouraging, based on how the utterance and discourse is carried out by the instructor, including their pitch, prosody, body movement and other factors (their identity as perceived by the students).

<div align="center">**Deixis in Spoken Academic Discourse**</div>

Deictic markers connect the speakers and the audience to a shared time, space and personal context. Levinson (1983) provided a typology of deictic categories which include *personal deixis* (I, you, we), spatial deixis (this, that, here, there), temporal deixis (now, today, yesterday, etc.) , social deixis (i.e. polite pronouns, familiarity pronouns), and discourse deixis (i.e. *we are gonna stop here and pick up this topic tomorrow).*

Deictics, such as *this, that, here, there* can perform a pointing function in academic discourse and depend on the situational context of use (Friginal et al., 2017). They are used to connect the interaction to its context (Cairns, 1991). Without the context of use, space, and time location, they do not have any independent meaning (Levinson, 2004, p. 103). More importantly, deictics are considered an involvement feature, because they are used by speakers who want to draw the attention of the addressee to an element in space/context (Cairns, 1991). The speaker directs the addressee's attention either to proximal (close) elements by using *here, this,* and *these*, or points to the distal elements with the use of *there, that,* and *those* (Levinson, 1983). The object of the attention might not be present in the immediate context of the interaction but deixis can still be used to refer to such an element (Sidnell & Enfield, 2017).

Deictics, then, unify the attention of the speaker and the addressee. The use of deictics is necessary if there is more than one element present in the communicative situation. In the sentence: "Please bring *the* book", the addressee knows that there is only one book. When the sentence relies on deictics: "Please pick up *that* book, not *this* book", deictics help the addressee and the speaker shift the attention between elements, on condition that both of them are paying attention to the utterance. Krebs and Dawkins (1984) discussed how this state of being co-

present, allows the individuals to be manipulate each other's attention and gaze with the use of deictics. Biber et al. (1999) reported that deixis *that, here, there* are more common in conversation, whereas *this, these* and *those* are more often used in academic writing. Research on pragmatics in spoken corpora (O'Keeffe et al., 2011) revealed that the spatial deictic *that* was the most frequently used in a general conversation corpus, which means it plays an important part in interactive exchanges.

Deictics, despite their crucial role in manipulating attention and gaze, have not been explored at length by researchers who focus on educational contexts. There are two studies in the EFL/ESL contexts that discuss the use of deictics in teaching. Bamford (2004) explored the use of the deictic *here* in the discourse surrounding visual aids (such as graphs, diagram, table, etc.) in two academic corpora (MICASE and Siena corpus) and casual conversation corpus. Results showed the importance of gestural *here* in reference to visual aids and its importance for creating rapport with students. A more recent study of the use deixis in EAP classroom was conducted by Friginal et al. (2018), who found that spatial deixis is used very frequently in learning and teacher classroom talk, suggesting that EAP classroom register is similar to conversation register in that respect.

## Interaction and Interactional Features

Interaction is an essential component of academic contexts, since students also learn through dialogue with the instructors and with other students. In spoken discourse studies, interactivity is measured through pattern of turn-taking (Barbieri, 2015). The studies on interactivity in academic discourse have explored whether any interaction is present during a lecture (DeCarrico and Narringer, 1988; Nattinger and DeCarrico, 1992), the level of interactivity based on audience turns number per class session at different instruction levels (Csomay, 2007), the use of questions (Bamford, 2005) and discourse markers (Lee, 2009). Morell (2004) found that more interactive lectures include more first and second person pronoun as well as other discourse features such as elicitation, display, and referential questions.

### Response Elicitors: *Right? Okay?* in Spoken Academic Discourse

Response elicitors *right? okay?* are explored as a subset of the category of discourse markers. They are often called invariant, or fixed, tags. Response elicitors are defined as

generalized question tags and play an important role for the speaker who is seeking a response (verbal or non-verbal) that the message has been understood or accepted (Biber et al. 1999, p. 1089). In an academic setting, they function as a comprehension check. The question *right?* usually requires a verbal response from the hearer (Biber et al. 1999, p. 1089). In the earlier grammar (Quirk, 1985, p. 1481) it is considered a "comment clause" used to involve the addressee in the conversation. These questions tags usually follow a declarative clause, as in *This cat is beautiful, right?* and is considered an "exit technique" (Schlegoff, Jefferson, & Sacks, 1974) for the speaker who signals the end of their turn and shift of responsibility to the addressee. The use of question tags also helps the second speaker shape their response, by providing an option to agree or disagree (Heritage & Raymond, 2005). Their non-interrogative function of such markers as *right, okay* signals turn taking in a conversation (Biber et. al, 1999, p. 1046).

Research on tag questions pragmatics (Heritage, 2012) reveals their interactional function by allowing the speaker to 1) request information or 2) respond to support a point of view. Tag question meaning should also be explored with attention paid to intonation, which can reveal the difference in function (Cameron, et al., 1989). In one of the largest corpus-based study of question tags in British and American speech, Tottie and Hoffman (2006, p. 301) classified the question tags into six categories, by utilizing previous frameworks: 1) informational (genuine request) and 2) confirmatory (speaker is not sure about what they are saying), 3) affective: attitudinal (emphasis, but there is no need for response), 4) facilitating (i.e. teacher knows the answer but wants to involve a student), 5) preemptory, and 6) aggressive.

In academic settings the use of question tags *right? okay?* has been explored in lectures by native speakers (Othman, 2010; Thompson, 1998), seminar lectures by speakers of different genders (Schleef, 2005),and supervisor meetings with students (Bowker, 2012). Thompson (1998) explored the use of questions in a small corpus of academic lectures. In the two major question categories used in academic lectures: audience-oriented and content-oriented, question tags *okay?, right? all right? yeah?* belong to the audience -oriented category. While in conversations these tags are used as interactional signals, Hunston (1998) noted that in lectures they function more like a symbolic negotiation, rather than a real invitation for interaction. These audience-oriented questions might not evoke a verbal response but might involve a non-verbal response (i.e. head nod) from the audience, thus engaging students in the lecture. Hunston (1998)

posits that these questions are an overt demonstration of the speaker's concern over students' following the presentation.

Othman (2010) research on a corpus of four lectures given by 4 different speakers aimed to explore the use of *okay, right,* and *yeah.* She found that *okay?* plays a very important role in structuring involvement and interaction because it is a marker of progression or a confirmation check (p. 672), with sub-functions of response eliciting, assurance, and information. She contrasted the use of *okay?* with *right?* which revealed that the second question is used to mark a sense of shared knowledge (p. 674). Schleef (2005) analyzed MICASE lectures, finding that humanities instructors used fewer progression check question tags and modal question tags (asking for confirmation or information) than natural science instructors did. He commented that instructors in natural sciences are aware that students struggle with fact-oriented subject matter, therefore they check in comprehension with more frequency.

Bowker (2012) used conversational analysis methodology to explore the use of the tags *yeah? right? okay?* in conversations between advisors and students. He found that these question tags serve the purposes of addressing the supervisor's and students' face needs, with *okay?* serving as an affiliative tag to check for acceptability of a suggestion or an arrangement (p. 187). *Right?* functioned in those meetings to signal shared knowledge and avoidance of being patronizing, simultaneously signaling anticipation of affirmative response which might make raising questions difficult for the advisees. The literature on the response elicitors use in academic settings is not extensive; however, the current research shows their crucial function for building involvement and interaction in academic discourse.

**Direct Hypothetical Reported Speech: *You might say / think***

A typical element of interactions between people is reporting the discourse of others (Bakhtin, 1981). There is a particular type of reported discourse that does not represent any actual interaction between people, but is a result of fictious creation of the speaker. This type of discourse – hypothetical discourse (Haberland, 1986) – is used for different rhetorical purposes, based on the register it functions in. It is a feature present across various languages, such as Danish (Haberland, 1986), French (Fleischman & Yaguello, 2004) or Maya (Hanks, 1993). Golato (2012) analyzed hypothetical discourse function in German and English mundane speech, and found that there are three main functions of this type of discourse: 1) modeling discourse for

other 2) backing a claim 3) adding humor (p. 31). The modeling function has also been explored in other contexts in discussion of its grammatical iteration – hypothetical reported speech (HRS). For example, therapists use "hypothetical active voicing" to model a type of talk that the client could use in a future problematic situation (Simmons & LeCouteur, 2011). Koester (2014) found that in business contexts, HRS is used as a rhetorical persuasive and rapport-building device. Through the use of imaginary responses and quotes, the speaker can create affiliation, rapport, and understanding of their audience in business negotiation context.

Personal pronouns use in combination with communication or mental verbs, especially in the phrase "you might/could/can say/think" or "you are/ s/he was like" (Barbieri and Eckardt, 2007; Buchstaller, 2014) can suggest the presence of direct hypothetical reported speech (DHRS). In my corpora, this type of reported speech is an enactment of what the student might say or think as the lesson is progressing. An example from KAM corpus of a DHRS statement is

So let's multiply negative 7 times 3/49. So you might say, I don't see a fraction here. This looks like an integer. But you just have to remind yourself that the negative 7 can be rewritten as negative 7/1 times 3/49. [doc#25KAM]

The use of directed reported speech for creating hypothetical situations performs a very specific discourse function based on the type of register it is present in. Such a speech act is an example of direct reported speech, in which the speaker quotes what another person said, is an important and recurrent feature in conversation (Biber et al. 1999, p. 1118). Speakers rely on utterance openers such as *oh, well, look, okay* to mark the beginning of a quotation. Another way to signal quoting, especially among younger generations, is the use of particle *like,* preceded by form *be* (Biber et al. 1999, p. 1120).

Research on direct report speech complicates the idea that speakers report other's speech verbatim. Speakers can construct utterances that were never spoken or heard by them (Emmison, Butler, & Danby, 2011). Myers (1999) proposed a category of possible or conditional hypothetical represented discourse, in which the speaker "frames possible utterances as a way of provoking responses" (p. 576). In the case of hypothetical represented discourse, the speakers imagine what could be said or might be done in a given situation, such as in the example:

(2): you know we might see an image of malnourishment and disease and we think Gosh we've got to do something about that you know let's get organized or let's send parcels

out, let's lobby this or lobby that, er you know it's quite similar with Louise Woodward, for some people (Meyers, 1999, p. 576)

Tannen (2007) in her book *Talking Voices* coins the term "constructed dialogue" as a feature of involvement. This term is contrasted with the more common grammatical term "reported speech". Rather than being reported, Tannen claims an utterance is transformed by the speaker to communicate something else than was originally intended. In her discussion of involvement, Tannen (2007) points to the use of constructed dialogue as a source of emotion in discourse. The function of DHRS, then, is for the teacher to enact the student's thinking process, their challenges, and model responses to the posed problems. This is evidence of teacher's recognition students learning needs, since the instructor focuses on the needs of his audience.

### *Hypothetical reported speech and mathematics knowledge (Student and Teacher Components)*

In terms of mathematical education, the use of DHRS can be linked to the concept of Knowledge of Content and Students (KCS). KCS is defined as teacher's "content knowledge intertwined with knowledge of how students think about, know, and learn this particular content" (Hill, Loewenberg, Ball, Schilling, 2008). This type of knowledge also encompasses the understanding how, for example, students learn to add fractions and what challenges they typically encounter when they engaged in the process of learning that content. The construct of KCS is an important element of Shulman's (1986) concept of pedagogical content knowledge (PCK), which is a type of knowledge that allows the teacher to gauge whether a topic is easy or difficult for her/his students, given their age and their background. While the teacher's PCK can vary, the issues students encounter while solving certain tasks are shared across international contexts (Hadjidemetriou & Williams 2002; Leinhardt et al., 1990). Inexperienced instructors might overestimate solution rates among students on tasks, because they experience *expert blind spot* (Nathan & Petrosino, 2003), which is a term that denotes overreliance of expert educators on the advanced subject-matter knowledge, rather than focusing on their students' learning needs and developmental profiles (Koedinger & Nathan, 1997; Nathan, Koedinger, & Alibali, 2001).

Research on student development of mathematics knowledge rests on examining two foundational constructs: conceptual and procedural knowledge, two ideas which lie at the basis of learning in general, including mathematics education. Concepts are defined as "abstract and general principles such as cardinality and numeric magnitude" (Rittle-Johnson, 2017). The

knowledge of concepts can be implicit or explicit and is not tied to particular problem types (Rittle-Johnson & Schneider, 2015). There is also a difference in novice and expert conceptual knowledge: novice learners struggle with fragmentary, unintegrated conceptual knowledge, while experts' knowledge is organized and expands with time (diSessa, Gillespie, & Esterly, 2004; Schneider & Stern, 2009). Procedural knowledge can be defined as "'knowing how', or the knowledge of the steps required to attain various goals. Procedures have been characterized using such constructs as skills, strategies, productions, and interiorized actions" (Byrnes & Wasik, 1991, p. 777). There is debate on whether conceptual knowledge needs to precede procedural knowledge for learning to be effective (Baroody, 2003), or whether it is a bidirectional, iterative process (Rittle-Johnson, Siegler, & Wagner Alibali, 2001). Procedural flexibility is the ability to choose different procedures to solve the same problem, and being able to apply the adaptively.

Each type of mathematics knowledge can be assessed using sample tasks, as laid out by Rittle-Johnson (2017, p. 185).

Table 1. Types of knowledge in mathematics.

| Conceptual Knowledge | |
|---|---|
| a. Evaluate examples of concept | a. Decide whether the number sentence 3=3 makes sense |
| b. Translate quantities between representational systems | b. Place symbolic numbers on number lines |
| c. Compare quantities | c. Indicate which symbolic integer or fraction is larger |
| d. Generate or select definitions of concepts | d. Define the equal sign |
| Procedural Knowledge | |
| a. Solve problems in a familiar format | a. 8/10+6/10=__ |
| b. Solve problems with a new surface or problem feature | b. 2½+¼=__ |
| Procedural Flexibility | |
| a. Generate multiple methods | a. Solve this equation in two different ways: 4(x+2)=12b. |
| b. Evaluate nonconventional methods | b. Do you think this way of starting this problem is (a) a very good way; (b) OK to do, but not a very good way; (c) not OK to do? |

In my dissertation research I used these categories, to some extent, for identifying the uses of hypothetical reported speech into conceptual/factual or procedural/application category.

## Corpus Design: Representativeness, Balance, Sampling, Comparability

This dissertation project is grounded in corpus linguistics methodology, which allows for analysis of large, principled datasets of texts - corpora. Designing a corpus requires following a number of guidelines that the researcher in collecting, processing and analyzing the corpus. This section of my literature review presents a short synthesis of literature on issues in corpus design.

Representativeness is a characteristic of corpus design that raises the question of adequacy of the sampled language in the corpus to represent the characteristic of the register as a whole. Representativeness is defined by Biber as "the extent to which sample includes the full range of variability in a population" (1993, p. 243). My dissertation project consists of samples collected on Khan Academy platform (all high school and college level math instruction) and MIT Open Course Ware platform (all available instruction in math). My assessment of the sample representativeness is more problematic for MIT corpus, since it aims to represent the register of math instruction in face-to-face modality at MIT. Only a small percentage of these lectures are represented on the Open Course Ware platform, making the corpus a sample of the language used in teaching math at MIT.

The corpora used in my study are *sampled* corpora (Biber, 1993; Leech, 2007), because they reflect the language used on both platforms at the time of collection. In both cases, I used the total available number of transcripts for mathematics at high school and college level on Khan Academy and MIT at that point in time. Both platforms are continuing to grow, as more lessons or lectures are added to them. Both corpora can be referred to as *specialized* corpora (McEnery, Xiao,& Tono, 2006) since they are domain specific (i.e. math) and genre specific (i.e. monologic instruction).

Another important parameter in corpus is *balance*, which is a term that denotes the equal representation of texts types included in the corpus. The range of texts for a general corpus of English, then, should not include 90% of, for example radio talk shows and 10% of radio commercials, but rather it should strive for equal representation of both in a corpus of radio register. That is why proportional sampling is important, so that a balanced range of texts are represented and researched, rather than just one small subset. In the corpora used for my

dissertation, it was important to represent the instruction in different subfields of math at the high school and college level. However, in terms of text types, the corpora only consist of transcripts of monologic speech (Khan), or mostly monologic speech (MIT) with very short or inaudible audience interjections.

# CHAPTER 3: SITUATIONAL ANALYSIS

## Situational Analysis: Rationale and Summary

A vital part of any register analysis study is the analysis of its situation of use. This procedure – a situational analysis – serves the purpose of describing in detail the particular constraints that shape the use of a given variety. Rich description of circumstances, such as situational analysis, is used to account for differences in comparative analysis results of registers use. This is not a new approach, as scholars in the early years of applied linguistics research drew attention to the way that the nature of how communicative context shapes a speech act (e.g. SPEAKING model by Hymes, 1974; Halliday, 1978). The description of the context is grounded in my analysis of secondary resources that accompany the videos.

Each situational analysis needs to be preceded by the identification of the register under examination. Thus, examining the situation of use of a register involves describing the participants and the relations among them, the channel (mode and medium), production circumstances, setting, communicative purposes, and topic (Biber & Conrad, 2009; p. 40-47). Two registers that I examine in my study include Khan Academy tutorials in math and MIT math lectures. First, I present an overview of the situational characteristics (Table 2), followed by a detailed analysis of the characteristics.

Table 2. Overview of Situational Characteristics

| Situational Characteristics | Khan Academy Math (KAM) Corpus | MIT Math (MITM) Lectures Corpus |
|---|---|---|
| **Participants** | | |
| Addressor | Single addressor – Salman Khan. | Multiple addressors: Srini Devadas, David Jerison, Eric Demaine, Denis Auroux, Gilbert Strang, Peter Kempthorne, Choongbum Lee, Charles Leierson, Nancy Lynch, Haynes Miller, Vasily Strela, Jake Xia, Arthur Mattuck. |
| Social Characteristics | Male, started recording the tutorials at age 30 (in 2006) and continues until now, highly educated (graduate degrees from Ivy League schools), worked in finance and tech industry, worked in a hedge fund, social class: middle class background, first generation American. | 12 male addressees, 1 female addresses, varying in age and teaching experience (in terms of years teaching at MIT), highly educated (all of them have PhD). All teaching / doing research at MIT, one of the top universities in the world. |
| Addressees | Other; | Other; MIT undergraduate students, present in the classroom; |
| On-lookers | Yes. | Yes. |
| **Relations among participants** | | |
| Interactivness | No direct interaction. | Direct interaction or possibility of direct interaction. |
| Social roles | The addressor has authority but the power relations are not straightforward. The addressee can stop using service without consequences. | The addressors have authority and have power over the addressees who are physically present in their lectures (the students). |
| Personal relationships | Family (early on), strangers. | Professor-student relationship, or strangers (on-lookers). |
| Shared knowledge | Varies: part of the knowledge (background knowledge) is shared and part of it is new (the new concepts being taught). | Varies: part of the knowledge (background knowledge) is shared and part of it is new (the new concepts being taught). |

Table 2. continued

| Situational Characteristics | Khan Academy Math (KAM) Corpus | MIT Math (MITM) Lectures Corpus |
|---|---|---|
| **Channel** | | |
| Mode | Speech. | Speech. |
| Specific Medium | Taped / Monologue | Face-to-face / Taped |
| Production circumstances | Planned, real-time*, can be re-recorded, it is unedited, unscripted up to 10 minute long lesson. | Planned, real-time. |
| **Setting** | | |
| Time and place | Time and place are not shared by participants. The viewers can watch the videos on their devices at any time. | Time and place are shared by the participants. The onlookers can watch the recording at any time. |
| Place of communication | Public: Internet users. | Public, but exclusive: the classroom. |
| Time | Contemporary. | Contemporary. |
| **Communication purposes** | | |
| General Purposes | Exposit, inform, explain; how-to, procedural; teaching. | Exposit, inform, explain; how-to, procedural; teaching. |
| Specific purposes | Teaching math using virtual blackboard. | Lecturing / teaching math using examples written on the board. |
| Factuality | Factual. | Factual. |
| Expression of stance | Epistemic. | Epistemic. |
| **Topic** | | |
| General topic "domain" | Education/Academic. | Education/Academic. |
| Specific topic | Mathematics. | Mathematics. |
| Social status of person being referred to | Students-no special social status. | Students-no special social status. |

**Situational Analysis: Khan Academy**

Khan Academy is a non-profit company striving to provide "a free, world-class education for anyone, anywhere." (Khan Academy, nd.). This non-profit developed [www.khanacademy.org](www.khanacademy.org), an online educational platform that allows students, teachers, and parents to set up an account and track progress The data reported on the website shows that the platform is used by over 50 million registered users and many more who are not registered. The registered users come from different geographical locations, with majority coming from United States (70%) and the rest from India, Brazil, Mexico, South Africa and others. The Khan Academy contains instructional videos appropriate for elementary, middle and high school students, alongside college level and pre-professional content.

A substantial part of my situational analysis relies on two sources: my own observations and descriptions provided by Sal Khan in his non-fiction book-fiction book *The One World Schoolhouse: Education Reimagined* (2012), his speeches (for example TED talk or IBM talk) or the video "Behind the Scenes with Sal Khan" (Bank of America) in which he described his recording process. I also take a more qualitative approach to the situational analysis, with the consideration of the development of Khan's instruction over time, as presented by him through the various sources (talks, books, articles).

**Participants**

*Addressor*

Instructional videos in my dataset are all authored by Salman Khan, the original founder of Khan Academy. Khan was born in 1976 in Louisiana to immigrant parents. His father was a pediatrician from Bangladesh who moved to Louisiana for his medical residency at Louisiana State University (Khan, 2013; p. 15). His mother, who was born in India, married his father in Bangladesh and moved to USA in 1972 to start a family. Salman Khan grew up in New Orleans and went to school there. He graduated from MIT with two B.A. degrees in mathematics and computer science, followed by an M.Sc. degree in electrical engineering, and an MBA from Harvard Business School. In an IBM Think 2018 talk on the development of the non-profit, he mentioned that directly before starting to work full time for Khan Academy, he was worked as a hedge fund analyst. Khan did not have experience in teaching math or any other subject in a

formal setting. His degrees from MIT and HBS provide evidence that he has been successful in these highly competitive educational environments.

The addressors of the video developed over time, from the first tutorials that were addressed to his niece, then other young extended family members and friends. This speaker first taught the video math lessons to his niece in 2004, via telephone, in a form of conversation. Over time the niece requested video recordings that could be paused and rewind and did not require shared time and direct interaction. Khan uploaded his first math video to his YouTube channel in 2006 (NASA, 2014). Over time, Khan noticed that more people were watching his videos on YouTube (Khan, 2018) and decided to start a non-profit and develop an online platform which would present the videos in a more principled manner and added other features to it (exercises, badges, etc.). This trajectory of development, from addressing the lessons to a single person (his little cousin), to groups of family members, and then a bigger audience on YouTube might have left indelible marks on the way the mathematical concepts are presented to the listeners. The addressees of the videos shifted from single (the niece), to plural (other family members), to unenumerated.

### *The Characteristics of the Addressees*

Khan Academy now is approaching 60 million registered users, which includes students, teachers and parents. The tutorials are also available to unregistered users. The materials, including math tutorials, are extremely popular with general audience, whose main purpose is to learn or to teach. The videos have been viewed over 1,597,729,228 times on the YouTube platform. These videos are as popular as the lectures published by MIT, which has a recognizable brand as one of the top universities in the world. Unlike MIT, Khan Academy is not associated with any physical, brick-and-mortar educational institution.

The first addressee of the tutorials was Nadia, Khan's cousin who struggled with math. The story of the first tutorials with Nadia and the development of Khan's philosophy and practice of teaching is described in the chapter of his book titled "Teaching Nadia" (p. 16-25). Khan described his 12-year-old cousin Nadia as capable and destined getting a college degree (p. 16), but also lacking confidence in her math ability. Over the course of individual work with Sal, Nadia gained confidence in learning math and progressed to a more advanced math program. Following the success with Nadia, other family members and friends were tutored online, with

the use of different types of technology. The ultimate step for Khan was posting the videos on YouTube, and realizing that there is an audience for the videos that goes beyond his relatives and friends.

## Relationship among Participants

### *Interactivness*

At the beginning of the teaching history, Sal was able to interact with his tutees by phone or via the doodling platform. He was able to monitor and interact with his tutees in an interactive session. Khan's cousins expressed preference for less interactive learning saying that they preferred his YouTube lessons rather than in-person contact. In fact, Khan developed a philosophy of teaching which stipulates that video tutorials provide an opportunity for no face-threatening learning situation, in which students can pause and repeat the lessons as many times as they want. Such activity allows the students flexibility, without the pressure of immediate confirmation of understanding and response to comprehension checks (Khan, 2011). When Khan moved his instruction to YouTube, he started receiving feedback sent through comments under videos or via e-mail from strangers. As the popularity of his lessons grew, the venture turned into a non-profit with a bespoke platform, which gives the opportunity to interact with other users of the platform, and employees of Khan, which can drop in to check the comment section of the videos, to answer specific questions.

In speeches and interviews, Khan talks about the written feedback he got from his anonymous viewers that made him continue working on the *project*:

> "I started getting some comments and some letters and all sorts of feedback from random people around the world. These are just a few. This is from Zone of the original calculus videos. Someone wrote it on YouTube, it was a YouTube comment: 'First time I smiled doing a derivative.'" (Khan, TEDeX, 2011)

Clearly, the feedback is important to him and is a reason to continue developing the platform.

*Social Roles*

Another important aspect of the video tutoring, and Khan's approach to communication with his viewers, is building the sense of intimacy, familiarity, and unthreatening environment. There are two reasons that emerge as to why Khan decided to balance out the power differences between him and those who were listening. First, unlike most teachers who get degrees and start their careers as instructors in a formalized setting, his practice grew out from a personal and family relationship with his niece. Secondly, Khan also wanted to avoid the ineffective and disengaging aspects of traditional classroom instruction. In his description of the context of the teaching situation, Khan evokes an image of sitting with a student at "the kitchen table, elbow to elbow, working out problems together. I didn't want to appear as a talking head at a blackboard, lecturing from across the room" (Khan, 2012, p. 34). The context of the teaching he imagines is more informal and in direct opposition to the threatening figure of a lecturer. In his book he also describes the tone of his own school teachers as being "occasionally arrogant and even condescending" (Khan, 2012, p. 18). As his experiences of traditional model of teacher instruction embodied an unequal power relationship, his own teaching strives to be the opposite. He is aware that educational contexts can be challenging to the students in terms of they make them feel uncomfortable. The purpose of his tutorial session was to make it be "a safe, personal, comfortable, thought-provoking experience" (p. 18).

*Personal Relationships*

In the chapter "Learning how to teach" of his book, Khan explains what he wanted to achieve in terms of building the relationship with his audience "You talk *with* someone, not *at* someone." (Khan, 2012, p. 34). Despite talking to strangers, Khan wants to build a sense of familiarity and personal connection in his tutorials. Additionally, even though the tutorial, unlike his early phone conversations, does not provide a possibility of immediate response and conversation with the students, he still wants to be talking "with" them rather than "at" them. This suggests that his pedagogy relies on dialogic approach, recognizing that the students learn when they are involved and interacting, rather than through passive listening.

His choice to have a conversational approach in his tutorial did not influence his decision to not show his face in the tutorials. In the chapter "Focusing on the content" he mentions that

face time can help the teachers convey empathy and concern but is not helpful in teaching academic concepts (Khan, 2012, p. 35). That is how Khan decided to not include his face or body in the recording, just a voiceover and virtual blackboard, so as to not distract them with unnecessary component of facial expressions.

**Shared Knowledge**

Teaching is a special case when considering shared knowledge dimension. The teacher, in this case Khan, has the knowledge and skills he wants his audience to acquire in one tutorial. By the end of watching the video the audience is supposed to be able to demonstrate the same knowledge and skills as the instructor.

Many videos start with a statement of purpose of the video. This signals what the goal of the tutorial is. This is also the goal of shared knowledge development, which is being clearly defined at the beginning of the

You will hear me use the word abstract a lot so I thought I would actually give you an attempt at a definition [doc#81KAM] [3]

In this video, I want to familiarize you with the idea of a limit, which is a super important idea [doc#81KAM]

As Khan works through a variety of concepts in each video, he also refers to the previous content he recorded

Table 3. Collocations of the word "video" in the KAM corpus.

| Word | Co-occurrences (raw count) |
|------|---------------------------|
| last | 211 |
| next | 134 |
| previous | 68 |
| future | 58 |

---

[3] In my dissertation project I use excerpt from transcripts to provide examples of the linguistic feature use e.g., [doc#638KAM], The information in bracket is the number of transcript in the corpus: doc#638. KAM stands for excerpts from Khan Academy Math corpus; MITM signals a transcript from MIT Math corpus. All of the fragments used in the dissertation can be found on both platform with the use of a search engine (i.e. Google.com), because all of them are available for non-commercial purposes.

The cooccurrence of these adjectives with the lemma "video" suggests that Khan expects the students to be following the sequence of the videos as he prepared them. He refers to the knowledge that was shared before, so the students can build on the concepts as the teaching progresses. In the excerpt below, Khan reminds the audience about the rules he explained in the previous video

> But how do we evaluate that? Do we do the division first, or the multiplication first? And remember, I told you in the last video , when you have 2-- when you have multiple operations of the same level-- in this case, division and multiplication-- they're at the same level. You're safest going left to right. Or you should go left to right. [doc#21KAM]

Some videos begin with elements borrowed from pop culture, which can help build rapport with the listeners. For example, the video on matrices begins with a reference to the Wachowski movie trilogy *The Matrix*. Khan makes a reference to the use of matrices in computer graphics and simulation, relating it to the plot of the movie in which the protagonists are escaping a simulated world. This strategy of bringing a familiar object (the movie) to teach a new, more specialized concept relies on shared cultural knowledge that is being used for the purpose of learning academic content.

**Channel**

The mode of the in my corpus is speech, while the specific medium is permanent video recordings. Early in the history of the instruction, Sal and Nadia were using Yahoo online doodles and phone conversation for teacher-student talk, which was a more transient medium. With the transition to recording YouTube video tutorials, the format changed to up to 10 minute long video recordings. The length was dictated by the restrictions put on the content publishers to keep their videos shorter than 10 minutes (Khan, 2012, p. 28).

**Production Circumstances**

In a 2013 video ("Behind the Scene"), Sal Khan described his method of work as consisting of two parts: preparing for each individual video and then recording the lesson. He specifically mentioned that he does not write scripts for his lessons, rather he relies on his preparation work to produce a lesson. The mode of the register, then, can be listed as spoken.

Similarly, to face-to-face instruction (a lecture), when the speech is not scripted, but might have been practiced. Thus, the production circumstances can be described as real-time on the part of the speaker.

**Setting**

The setting of the tutorials is strictly a virtual blackboard. The listeners can use the website anywhere they have access to the Internet with a device that can display video.

The platform evolved over the years to accommodate the needs of three group of users: students, teachers, and parents. Each type of user can create an account on the platform, and depending on the choice, they will be presented with a different set of features. Parents and teachers can assign material for the students to study and to be quizzed on. They can also follow the progress and see which material is more problematic for the students.

The Khan Academy video lessons are unlike regular academic lectures because they typically focus on one concept from a discipline at a time. Exploring a concept happens over a series of shorter videos intertwined with interactive exercises. Students can see an overview of how the material is sequences in the menu of the page, after clicking on the

Figure 1. An example of sequencing of content on Khan Academy – Algebra foundations.

Another important feature of these instructional materials is that the instructors' face is typically never shown, instead the students follow a power-point like presentation or a "blackboard" on which the instructor writes the content of the lesson. The dependence of the spoken part of the lessons with the graphic representations being discussed perhaps creates a hybrid register that will display particular functions of language that are used to draw attention to the graphics involved. For example, in the example below, the word "tips" is color coded in green and associated with "30" and the color coding extends across multiple lines to allow the students to identify and track.



Figure 2. A typical tutorial on Khan Academy. A virtual blackboard.

**Communicative purpose**

The general purpose of the tutorials is to teach the listeners particular concepts or procedures in mathematics, fulfilling a descriptive and explanatory purpose. The tutorials include descriptions of the problems with a running commentary on procedural step-by-step which models actions for the students. Another type of tutorials provide definitions of key concepts in mathematics.

Because the tutorials are fairly short and focused, there are not many specific purposes that can be distinguished. The tutorials usually include introduction, with a statement of the purpose of the video, for example

**What we're going to do in this video is get** some practice evaluating exponents of decimals.[doc#17KAM]

The explanatory part focuses on a problem or a concept that is being taught in the tutorial. The conclusion section, which is usually a sentence or two can include a preview of the next video

And as **we'll see in the next video** , calculating by the inverse of a 3x3 matrix is even more fun. See you soon." [doc#1113KAM].

This basic macro structure of opening – instruction – closing has been discussed in the context of classroom discourse (Sinclair & Coulthard, 1975; Cazden, 1988).

Another key consideration in examining the communicative purpose is factuality (Biber & Conrad, 2009, p. 46). The tutorials in Khan Academy are set up to be factual, insofar as Khan is conveying factual information in the domain of mathematics. Many videos in which Khan made a mistake by misusing a term include a pop-up window showing a commentary on what the correct term should be in that instance.

Expression of stance is the final consideration in the communicative purpose analysis. A more detailed analysis of stance is not conducted in this dissertation. A broad overview of stance for instructional tutorials suggest that similar to classroom teaching (Biber, 2006), the use of modal verbs is the most overt grammatical device used in Khan's speech. Khan's style of teaching is marked by frequent use of modal verbs, especially by modals of permission and possibility. This can be interpreted as his teaching being more encouraging of showing different possibilities in math procedures:

Or **you could** flip both sides of this. Or **you could say** the ratio between 5 and 15 is going to be equal [doc#48 KAM]

So that's the same thing as 0.30. **Or I could just** write 0.3. I **could** ignore that zero if I like. [doc#33KAM]

As shown in the results chapter, Khan prefers the use of personal pronoun *we* as a referent of *I.* This also affects his personal stance presentation, since the shift to *we* as the most

frequent pronoun also de-personalizes his discourse, making it an expression of shared stance. Since he does not have to do any classroom management, there is less opportunities for using necessity modals that direct listeners.

The overall purpose of Khan's teaching – making it as personable as solving math together with the student at the kitchen table – suggests that he is not striving to present an objective, de-personalized communication. Since stance is a crucial element of making communication more personal and subjective, it is an important feature of Khan academy.

**Topic**

The general topical domain of Khan Academy tutorial register used in my corpus is mathematics. The specific topics were extracted from the website using the metadata which accompanies each transcript I collected. Using the information on Khan Academy website on the level of each of the courses, I selected topics to include in my corpus that align with high school and college-level mathematics. There was a number of tutorials that would be used in two different courses, such as Algebra Basics and Algebra I and II, or AP Calculus AB and AP Calculus BC. This shows how Khan Academy integrates foundational concepts wherever there is a need to use a foundational concept, common to many courses.

Table 4. Topic distribution in the Khan Academy Corpus.

| Course | Number of transcripts |
| --- | --- |
| Algebra I and II | 647 |
| AP Calculus AB & BC | 323 |
| Statistics Probability | 240 |
| Precalculus | 115 |
| Algebra basics | 80 |
| Trigonometry | 71 |
| TOTAL | 1476 |

**Situational Analysis: MITM Lectures**

The MIT Open CourseWare platform has been first published online in 2002, offering 32 courses. Since then it grew to be offering 2439 courses in September 2018 (MIT OCW Site Statistics) and boasts about 267 million visits per year, with over 166 million unique visitors. Out

of the 10 most popular courses in the last month (January, 2019), 7 were math courses that are included in my corpus.



Figure 3. MIT OpenCourseWare users location. Source: MIT OpenCourseWare *Site Statistics* https://ocw.mit.edu/about/site-statistics/

The MIT mathematics lectures were recorded during regular semester at MIT, often with MIT students being in the audience. The audience was informed about the recording process before the semester started, as the instructors often mention a leaflet with detailed information about the recording process being distributed. The students who did not want to be visible in the recording were advised to move back to the end of the classroom. The recordings were then made available online on the MIT Open CourseWare website for anyone to watch. The lectures were not specifically recorded with the online audience in mind. The recordings of the lectures often include a view of the instructors standing in front of the classroom, writing on the board, or walking around the room. Rarely is there a view of the students or their faces.

**Participants**

*Addressors*

The MIT Mathematics lecture corpus consists of transcripts from professors or visiting lecturers who taught mathematics courses at MIT. The full list of the instructors whose lectures were used in the corpus is listed in Table 5. There are 13 speakers represented in my corpus, sampled from all of the college-level math courses, with video transcripts available on the MIT OCW website. Only 4 speakers are middle-aged, with 9 who started their careers at the

beginning of 1980s or earlier. All of them are highly educated, often with Ivy League graduate degrees and extensive work experience at MIT, a leading global university. Nine professors who are represented in the corpus have been working at MIT for over 20 years. There are two instructors: Jake Xia and Vasily Strela, both with doctorates, who served a role of visiting lecturer in the Topics in Mathematics with Applications in Finance course. In terms of their educational background, all of the speakers have PhDs in a variety of fields, with the most representing mathematics. All but one of them are men. This is not surprising, as at MIT math department only 8% of senior faculty were women (Natanson, 2017).

Table 5. The list of the speakers in the MITM corpus.

| First & last name | Teaching experience at MIT | Affiliation | # words in the corpus | % of MITM corpus |
|---|---|---|---|---|
| Arthur Mattuck | At MIT since 1958, professor emeritus. | Department of Mathematics | ~ 200,520 | 14.1 |
| Charles Leiserson | At MIT since 1981. | Electrical Engineering and Computer Science | ~ 110,802 | 7.8 |
| Choongbum Lee | At MIT since 2012. Educated partly in Korea (B.Sc). | Department of Mathematics | ~ 26,799 | 1.9 |
| David Jerison | At MIT since 1981. | Department of Mathematics | ~ 148,255 | 10.4 |
| Denis Auroux | At MIT between 1999-2011. Educated in France. | Department of Mathematics | ~ 201,381 | 14.2 |
| Erik Demaine | At MIT since 2001. Born and educated in Canada. | Computer Science and Artificial Intelligence Lab | ~ 229,456 | 16.1 |
| Gilbert Strang | At MIT since 1959. | Department of Mathematics | ~ 213,180 | 15 |
| Haynes Miller | At MIT since 1986. | Department of Mathematics | ~ 69,078 | 4.9 |
| Jake Xia | Visiting Lecturer. Has a PhD from MIT in EE and CS. | - | ~ 9,449 | 0.7 |

Table 5. continued

| First & last name | Teaching experience at MIT | Affiliation | # words in the corpus | % of MITM corpus |
|---|---|---|---|---|
| Nancy Lynch | At MIT since ~1984. | Software Science and Engineering, Professor of Electrical Engineering and Computer Science, | ~21,739 | 1.5 |
| Peter Kempthorne | At MIT since 1986. | Associate Professor of Management Science | ~56,607 | 4 |
| Srini Devadas | At MIT since 1988. Educated (B.A) in India. | School of Engineering | ~ 148,255 | 10.4 |
| Vasily Strela | Visiting Lecturer | - | ~ 5,646 | 0.4 |

*Addressees*

The addressors of the instructor's speech are the students taking the courses at MIT. Secondary audience are the viewers who are watching the courses online. The addressors are aware of being recorded for the purpose of sharing the videos with the audience on the Internet. Unlike courses made specifically for the online mode of delivery, these recordings were meant to capture the courses as they are delivered at MIT.

MIT Admissions office provides a comprehensive description of its physical student body on their website (2019). The students are selected from a large pool of applicants, with only 6.7 % of students being admitted into MIT. 73% of accepted students score between 750-800 points on the SAT Math component, providing evidence for competence in mathematics. 10% of the undergraduate students at MIT come from abroad. The majority of MIT undergraduate students matriculate into MIT right after high school, which means they are 18 year old when they begin their studies and graduate when they are 22-23.

*On-lookers*

The audience of the video viewers can be considered on-lookers. However, this group can be described as individuals who are interested in advancing their mathematical knowledge. These videos are not usually watched for entertainment, since they contain very specific knowledge on advanced mathematics. The lectures are available through a few platforms: the MIT Open Course Ware, YouTube, and Internet Archive (archive.org). The MIT platform provides detailed information on the geographical location of the on-lookers Figure 4).



Figure 4. Self-declared educational status of the visitors to the MIT website.

The data provided by the MIT OpenCourseWare on the educational status of the on-lookers reveals that they are mostly students (presumably enrolled in courses at different institutions), or self-learners, who learn for the sake of learning (not to receive a credential).

**Relationship among participants**

*Interactivness*

There are 484 marked turns of the audience in the corpus. This means, that the lectures exhibit some degree of direct interactivness in terms of spoken exchanges between the instructors and the audience. This is despite the fact that many of these lectures have dedicated small group recitation sections that are geared towards interaction.

Table 6. Frequency of audience turns in the MITM corpus, by instructor.

| Speaker | Occurrences of Audience Response | Audience turns per 10 thousand tokens |
|---|---|---|
| Srini Devadas | 156 | 13.3 |
| Eric Demaine | 108 | 4.7 |
| Choongbum Lee | 56 | 20.89 |
| Peter Kempthorne | 45 | 7.94 |
| Nancy Lynch | 42 | 19.32 |
| Co-taught course | 26 | 18.25 |
| Jake Xia | 24 | 25 |
| David Jerison | 14 | 0.9 |
| Haynes Miller | 9 | 1.3 |
| Vasily Strela | 4 | 10 |
| Gilbert Strang | 3 | 0.14 |

Since the instructors may be represented more or less in the corpus, the normed frequency shows the number of turns audience takes per 10 thousand tokens. The instructors who elicited the most audience turns are Choongbum Lee and Nancy Lynch, with a course that was lead by two or more instructors at the same time coming close after. It is important to note that all of the lectures are paired with recitation sessions, which are specifically designed for interaction. Thus, the primary function of these lectures is not to interact, but to deliver content to a large number of students.

*Social roles*

The MIT lectures naturally divide the participants into the social roles of professors and students. There are important social differences between these two groups, as often professors are in position of power because they design assessments and can decide directly or indirectly whether a student passes or fails a course. The unequal status of power relations can also stem from age, experiential, and professional background between the students and the professors. Since only a handful of the instructors are middle-aged, the age difference might contribute to the sense of social distance. Another aspect that is worth noticing is the gender imbalance, as most of the instructors are male. Research on student evaluation of male instructors show that they tend to ascribe them more power, in comparison to female instructors (Johnson, 2006).

*Personal relationships*

Typically, the students and professors might develop a friendly professional relationship, but personal relationship of romantic nature are prohibited. It can also be assumed that over the course of the semester the students and the professors become more familiar with each other, as they meet weekly for the lectures. The size of lecture groups might still prevent from more personal and familiar relationships forming, since there is very little opportunity for learning names or discussion time. There are, however, very few instances of the instructors calling out students by their name, which suggests that at least some of the students have developed more personal relationship with their instructors.

*Shared knowledge*

Similarly to KAM corpus, the participants in the MIT Math lectures have unequal knowledge level. This is typical of a university setting, in which the students work towards increasing their knowledge in a specific field. The students' progress through a series of courses in order to increase their knowledge in mathematics. In fact, in five different lectures in the MITM corpus, the instructors mention prerequisite courses, which allows the instructors to assume what previous knowledge the students have. Another layer of shared knowledge is the skills that students bring from their previous education institution, represented through the SAT math scores. There are also other non-content shared knowledge elements, such as how to behave in an educational setting or during a lecture.

**Channel**

The MIT lecture mode of delivery is speech, specifically a university-level lecture. However, the kind of oral discourse used in lecture is somewhat different from general conversation discourse, because it requires preparation, with notes and outlines guiding the speech of the instructor.

**Production circumstances**

As typical of other lecturing situations, the production circumstances of MITM corpus are planned. The lecturers follow a course structure and the instructors know ahead of time what topic they will be discussing. The OCW platform provides view of the lecture notes (plans) that

are prepared and shared by the lecturers on the course pages (Figure 5). These lecture notes, whether typed or hand-written, show that there is a certain amount of planning that goes into preparing for a lecture. These notes can also serve as a plan for the lecture, guiding the lecturer back to the main point when they lose track or go off-topic.



Figure 5. MIT Instructor notes.

**Setting**

The setting of the recorded MIT lectures is always an MIT classrooms in Boston. The set up of the classrooms is very typical for traditional teaching: the professors are standing or walking between their desk and the blackboard, and they are always in front of the students who are all facing the professors. The students are seated next to each other, in fixed rows, stadium-seat style.

In terms of the recording time, the lectures took place anywhere between 2005 and 2015. The immediate audience shared the time and place with the lecturers. However, the secondary audience of onlookers are removed from the lectures, being able to watch them any time and any place on their devices connected to the Internet.

The videos include an opening slide, which states the title of the course, the name of the instructor and the time of recording. The OCW also includes more materials from the courses and detailed information about the syllabus, time of meeting of the course, deadlines, etc.

The communication place is public, since the lectures are held at a university. However, the on-lookers can view it in any setting, such as home, which takes the lecture into a more private setting.

## Communicative purposes

The communicative purpose of MIT lectures is educational, since the instructors are trying to describe, explain, and model the use of mathematical concepts and procedures. Additionally, because the courses are taking place on a university campus, during a regular academic year, the other important part of the communicative purpose is course administration. Each course starts with explanation how the course works, what are its components, who is teaching each section, communication practices, etc. The instructor mention assessment procedures, requirements for participation and other elements typical of non-content related aspects of teaching.

## Topic

The lecture topics are centered on various mathematics subfields. Table 7 includes information about course topics represented in the corpus, the instructors who taught the courses, as well as the year and semester the recordings took place.

Table 7. Courses represented in the MITM corpus, instructors and semester of teaching.

| Course Topic | Instructors | Semester and year |
| --- | --- | --- |
| Topics in Mathematics with Applications in Finance | Dr. Peter Kempthorne<br>Dr. Choongbum Lee<br>Dr. Vasily Strela<br>Dr. Jake Xia | Fall 2013 |
| Single Variable Calculus | Prof. David Jerison | Fall 2006, Fall 2010 |
| Design and Analysis of Algorithms | Prof. Erik Demaine<br>Prof. Srini Devadas<br>Prof. Nancy Lynch | Spring 2015 |
| Differential Equations | Prof. Haynes Miller<br>Prof. Arthur Mattuck | Spring 2010 |
| Introduction to Algorithms | Prof. Charles Leiserson,<br>Prof. Erik Demaine | Fall 2005 |
| Linear Algebra | Prof. Gilbert Strang | Spring 2010, Fall 2011 |
| Multivariable Calculus | Prof. Denis Auroux | Fall 2007, Fall 2010 |

In collecting data for the corpus, I used the video transcripts listed alongside mathematics courses. There are three courses that are interdisciplinary in nature *Topics in Mathematics with Applications in Finance, Introduction to Algorithms* and *Design and Analysis of Algorithms*. These courses rely on mathematics as the foundation for knowledge in another discipline - finance or computer science. These courses are taught in part by mathematics professors and are listed under Mathematics category on the MIT OCW website.

# CHAPTER 4: METHODS

## Motivation of the research design

This exploratory research study investigates differences and similarities in the engagement language used in two large corpora of mathematics instruction: Khan Academy and MIT. The purpose of the analysis is to investigate potential differences in the way engagement is realized in these corpora representing two teaching contexts: online video instruction and face-to-face instruction.

The combined size of both corpora – over 2 million words – allows for a two-step research process: frequency-based analysis of linguistic features present in the corpus and qualitative analysis of their use in context through concordance line analysis. The frequency analysis is conducted using three groups of linguistic variables: words (e.g. personal pronouns *I, we, you*; modal and semi-modal verbs *will, should, ought to*), multi-word expressions which include these lexemes (i.e. 2-6 gram bundles), and frequency-based collocations (e.g. which verbs follow the bundle *I am going to* most often ). The results of frequency analysis are then investigated qualitatively in concordance line. The discussion of linguistic variables analysis is conducted in the context of their situational use, concluding in functional analysis.

## Corpora

### Corpus data collection

The two corpora used in this research represent two math instructional settings: Khan Academy (KAM) and MIT (MITM) lectures. Two sets of transcripts were collected for corpora used in this dissertation study. The transcripts for Khan Academy were collected in September 2018, while the MIT transcripts were collected in November 2018. Transcripts in both corpora are annotated with a header, which includes the information about the speaker and the subject of each lecture or video tutorial. KAM speaker data was based on description of the videos available through the website. In the case of the MITM corpus, I added the information to the header manually, as I was collecting the transcripts from the website. The description of the

speaker in the video was additionally confirmed by watching fragments of the lectures. Unlike in Khan Academy, a viewer can see the lecturer's faces in the MIT videos.

The Khan Academy transcripts were downloaded using the Application Programming Interface (API) made available by Khan Academy organization. Khan Academy shares their educational materials through Creative Commons Attribution-NonCommercial-ShareAlike 3.0 United States License (https://creativecommons.org/licenses/by-nc-sa/3.0/us/). This license allows the user to share and adapt the materials for non-commercial purposes.

The API allows anyone to obtain various data (including anonymous user data) from the Khan Academy platform, provided they know programming basics. Three programmers helped me with various parts of the script, until it was functioning as I intended. I was able to use the script to retrieve all of the available transcripts from the website and save them as .txt files in UTF-8 encoding. The script automatically fills in the headers with metadata from the Khan Academy website. That metadata includes information on discipline, sub-discipline, topic, author (speaker), a number of the video and a link to the original transcript.

```
<doc discipline="economicsfinancedomain" sub_discipline="macroeconomics"
topic="pikettycapital" author="SalKhan" video_title="12241"
transcript_url="http://www.khanacademy.org/api/internal/videos/o5-T52bh-eQ/transcript">
Sal: Before talking more about inequality I think it's worth talking about the
difference between wealth and income. Wealth and income, because I think they often get
confused in conversations about wealth and income, and also about inequality. As you
can imagine these two things move together. You tend to associate someone who has more
```

Figure 6. A screenshot of a sample text file with a header.

These headers are used by the corpus hosting platform (SketchEngine, described in a following section) for providing additional options for data analysis. This allows me to perform searches which are related to a particular speaker, discipline, and field. The original files are stored on my computer, while SketchEngine has a copy which is used to perform the searches.

The MITM corpus was collected manually, via the OpenCourseWare website. The materials on the website are available under Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0) creative commons license. First, I located all of the courses in mathematics, which also have video lectures and transcripts of audio. I manually collected the transcripts, by copying and pasting the text of the transcript from the website into a plain text file. As I was copying the text, I was also manually filling in the information in the header: the

course name and the lecturer. Some courses were taught interchangeably by two or more instructors, and the headers reflect the instructors present in the class on the day. I used the video feature to ensure that the data about the speakers in the headers, matched the speakers in the transcript.

For the purpose of my dissertation research, I collected only the transcripts of mathematics instruction for high school, AP (advanced placement) and college level from Khan Academy and MIT. In order to select appropriate transcripts from Khan Academy, I used the Common Core Map available on the website. I also researched college-level mathematics courses and selected the content from Khan Academy. All MIT courses are by definition college-level courses. The selection process was also guided by the classification provided on the MIT OCW website.

## Corpus data description

Using these procedures resulted in creating two corpora of similar sizes. Since the data was collected on the Internet, I used SPSS and SketchEngine to create data descriptions of each corpora (counting speakers, final number of transcripts, number of transcripts per course). The two corpora differ in terms of the number of speakers, number of courses included, and length of transcript. However, their overall sizes are similar.

KAM corpus consists of transcripts of monologic narration of tutorials, which total over 1.4 million words. These tutorials are taught and narrated by Salman Khan, who is the founder of Khan Academy. The transcription is solely orthographic, with some variety in terms of transcribing repetitions and restarts (see the section on transcript reliability).

The second corpus I collected consists of 1.4 million words and includes 7 undergraduate mathematics courses from MIT Open Course Ware website (MITM). A few of the courses are recordings of the same course taught at two different times in academic year or in different years. The MIT courses are taught by 15 instructors, with only 2 courses being taught just by one person and the remaining ones taught by two or more instructors. The lectures typically were taught by one instructor at a time, with certain instructors taking over a few consecutive lectures in a sequence.

The transcripts in KAM corpus are monologic, while the transcripts of the MIT lectures have minimal audience participation. In the MITM corpus there are 484 audience turns, with 84 inaudible and the remaining ones are short confirmation responses.

Table 8. Corpora sizes: Khan Academy Math and MIT Math.

|  | Khan Academy Math | MIT Math |
|---|---|---|
| Tokens | 1,774,532 | 1,736,215 |
| Words | 1,440,133 | 1,421,497 |
| Speakers | 1 | 15 |
| # of transcripts | 1477 | 217 |

Table 8 presents the difference between the two corpora – KAM corpus has almost seven times more transcripts, with an average transcript length of 975 words. MITM corpus has fewer transcripts, but their average length is 6550 words. This represents the difference in modality, with short concept-focused video tutorials being shorter than lectures that were recorded in university settings.

Token in the SketchEngine platform is the smallest unit that the corpora divide to. In my corpora, a token is a word or a punctuation mark. Contracted forms, such as *don't*, are counted as two tokens. The corpora were divided into tokens by a tokenizer, a program that was written specifically for English language and made available through SketchEngine platform. In my research, I used norming to 10 000 tokens, since the frequency count shows in SketchEngine in that unit by default. However, since I compare some of my results to previous studies which use word norming, this is a limitation of my study.

Table 9. Instructors represented and word count for each of them.

| Instructor | Token Count (words + punctuation) | Word count | Share of the corpus (%) |
|---|---|---|---|
| Demaine | 280,258 | 229,456 | 16.1 |
| Strang | 260,379 | 213,180 | 15 |
| Auroux | 245,967 | 201,381 | 14.2 |
| Mattuck | 244,916 | 200,520 | 14.1 |
| Jerison | 181,079 | 148,255 | 10.4 |
| Devadas | 139,610 | 114,303 | 8 |
| Leiserson | 135,334 | 110,802 | 7.8 |

Table 9. continued

| Instructor | Token Count (words + punctuation) | Word count | Share of the corpus (%) |
|---|---|---|---|
| Miller | 84,373 | 69,078 | 4.9 |
| Kempthorne | 69,140 | 56,607 | 4 |
| Lee | 32,733 | 26,799 | 1.9 |
| Lynch | 26,552 | 21,739 | 1.5 |
| Xia | 11,541 | 9,449 | 0.7 |

I used SketchEngine to automatically tag my corpus using CLAWS 7 tagger, which has reported average 96-97% accuracy, with certain features being less accurately tagged than others. Since my research does not rely heavily on analysis of grammatical structures, the use of the tagger and its accuracy is not cause validity issues in my research results. When the analysis relies on the use of the taggers, I perform error rate adjustment to the normed frequency count.

## Transcript reliability

My research project relies on data (transcripts) made available to the public as a means of making the videos accessible to a variety of audiences. Because the data was not transcribed by me, nor was it created for research purposes, checking the reliability of transcription became important before the analysis process began. Specifically, my goal was to assess how reliable the transcripts are for the purpose of my research of analyzing the frequency of linguistic features in both corpora. Since my quantitative analysis relies on comparing general trends in normed frequencies of linguistic feature occurrence, without any testing procedures, minor issues in transcriptions (i.e. irregular conventions of repetition transcription, occasional omission of discourse markers) does not have an impact on my results. Furthermore, I analyze randomly selected concordance line samples in terms of their function, which are not affected by minor issues of transcription (i.e. missing repetitions or single discourse markers). In my analysis it is important that, for example, for a phrase *I think that* to be noted at least once when the speaker utters it twice. A major issue in transcription would be an omission of this phrase, rather than unreliably transcribing it once or twice when the speaker repeats it, as he or she writes on the board. The repetition would be important for a study on stance or hesitation, but this is beyond

the scope of this project. In the analysis of transcripts, I did not observe major issues in conventions that would affect the overall frequency patterns that I report on.

The transcripts of the online tutorials and lectures was already available on the websites without information on who performed the transcription and what convention was used. In order to check the quality of the orthographic transcription of speech, I analyzed three randomly selected transcripts from both corpora. The analysis included playing 10 minute fragments of videos and following the transcription, coding any discrepancies between the speech and the transcript. The low number if transcript used for reliability analysis is a limitation of this study. A more representative sample of transcripts will be used in future studies to prevent negative impact on results.

**Khan Academy**

In the Khan Academy transcripts, there are discrepancies in terms of how closely the transcript follows repetitions, restarts, and omission of discourse markers. In one of the earlier videos from December 10th 2010 titled "IIT JEE trigonometry problem 1" the transcription follows the speech very closely. This transcript is free of omissions and includes unedited transcription of speech. In the example below the phrase *I'll do it* is repeated twice, as Khan is deciding which color to choose.

> So let me just write it down so we have i'll do it.. i'll do it in blue a over c a over c sine of two times capital C [doc#1473KAM]

Khan also repeats certain phrases of words as he writes them on the virtual board:

> So c squared is equal to a squared plus b squared minus 2 ab cosine cosine of C of capital C [doc#1473KAM]

This habit of repeating whatever Khan is writing on the board is present across videos recorded at different points in time. Example 1. transcript is marked with high fidelity in terms of representing orthographically repetitions and restarts. In the 10 minute transcription of speech, I found only 3 major omissions: missing discourse marker "*so*"; the word *expression* was substituted with the word *equation*; and a missing interjection *just in case it's useful*.

In two other videos I examined – "Focio of a hyperbola" and "Frequency table independent events", with a date of publication in the same year (2010), the transcription conventions differ. The person transcribing the text used punctuation marks to make the transcript follow written conventions. The two transcripts also do not represent the speech with the same level of fidelity, with repetitions and restarts not being reflected in the transcript.

Other evidence of editing is writing the full form of the word *because,* rather than *cause* which is used by Khan. There are also instances of not transcribing recasting (repeating the error in a corrected form), instead of skipping the unfinished clauses and transcribing only the fully formed clauses. There were 20 editing instances in the "Frequency" transcript and 18 in the "Focio" text. Despite the fact that the conventions differ,

Both types of transcripts lend themselves well to my analysis, which does not focus on repetitions or restarts, but rather on the frequency of lexical items. The discrepancies I detected do not influence the results of my research to the extent of invalidation.

## MIT

Three transcripts I examined from the MIT had between 3 to 6 minor transcription issues. These transcripts were the first lectures in the series: "Topics in Mathematics with Applications in Finance," (Fall 2013), "Differential Equations" (Spring 2010), and "Introduction to Algorithms" (Fall 2005).

The most often coded issue was omission of discourse markers and pause fillers*: so, you know*. There was one instance of missed "and" and "that a". The MIT transcripts also were unedited in terms of L2 speaker discourse issues, preserving his syntax and grammar mistakes. The number of missed or not transcribed items was 3,5,6 per each transcript respectively. Thus, the MIT transcripts are on average more accurate than Khan Academy transcripts, but in both cases the errors and editing interventions are minor and should not have impact on the type of analysis I performed.

**Corpus tools used in the study**

**SketchEngine**

The quantitative linguistic comparison of the two registers represented by the corpora was conducted using the SketchEngine platform. After uploading, the texts are compiled into a corpus, which involves the process of parsing, tagging, and computation of data. Even though SketchEngine cannot completely substitute writing one's own script to perform searches, it allows the researcher to use regular expressions, CQL (Corpus Query Language), and wild cards in simple searches.

I used CQL to look for hypothetical reported speech patterns. Since reported speech is usually introduced with the communication verb *say*, and that hypothetical reported speech utilizes utterance-launchers *hey,* I can look for these two items in proximity using CQL notation:

[lemma="you"][]{0,3}[tag="MD"]?[lemma="say"]

This notation asks SketchEngine to find instances of "you" followed by anything between no to 3 words, then an optional modal verb (signaled by the question mark) and lemma "say". Using the search term of lemma, allows me to look for different forms of the verb *say*: *said saying says* etc. Tag MD stands for all modal verbs. The results of this search were exported into Excel spreadsheet.

**Concordance and Key Words in Context (KWIC)**

Concordance is defined as "a list of all of the occurrences of a particular search term in a corpus, presented within the context in which they occur" (Baker, 2006). Exploring concordance lines is an analytical technique that is particularly useful in generating and testing hypothesis about item use in their context (Evison, 2010). Patterns observed in even small number of lines can be tested out through corpus searches on a larger scale (Hunston, 2010). One caveat of this analytical technique is taking into consideration its limitation: the patterns that are used in natural language, but are not used in the corpus, will not appear in results. Thus, there is limited generatability of this concordance line examination, that, to some extent, can be mitigated by use of large corpora. However, for features that occur very rarely in natural language, chances of recording the feature in a small or medium corpus are low.

Concordance lines can provide information about *typicality* of linguistic feature use, such as most frequent meaning and collocates (co-occurring items) (Evison, 2010). A collocate of a word is simply another word that tends to occur in the company of the other one (Moon, 2010). Such relationship is more frequent than a chance occurrence. The research on collocation already has a long tradition in lexicographical and phraseological studies, where structure and co-occurrence are used in tandem to examine meaning of words (for discussion see Sinclair, 1987; Stubbs, 2001). Collocations can be analyzed using simple frequency count procedures or use more sophisticated statistical analysis of mutual information scores (MI) or t-scores. This research study does not utilize statistical scores.

In the Sketch Engine platform, the key word might be a single word or a phrase that the researcher queries (marked in red). Sorting concordance lines in the right or left context can reveal patterns of use on a discourse level (see the analysis of *now* in Swales and Malczewski, 2001). The same patterns can be revealed using *Frequency* tool on the SketchEngine platform.



Figure 7. KWIC view in SketchEngine.

Search in SketchEngine can be performed using a number of functions, depending on the need of the researcher. Because my corpora are tagged, SketchEngine has the ability to either look for particular lemma in the canonical form (e.g will find only the instances of *goes*, when you type it in) or can search for all word forms (e.g. search for *go* will find *goes, went, gone*). I can also input a tag, such as V.*, to find all verbs in the corpus. The search functions also allow

the researcher to use regular expression (regex). Regex is a set of pre-determined text strings, which are used for making search patterns more efficient.

In my dissertation research searches, I often use * sign to look for a wildcard (like in Figure 7 example) or ? for an optional tag, lemma, or character in the DHRS pattern search.

**English TreeTagger PoS Tagset with modifications**

In order to be able to perform more advanced analysis in corpora, a tagger can be used to attach information about part of speech category to each word in the corpus. This is done automatically with the use of a script. Research on the English TreeTagger (Schmid, 1995) shows that it has 88% accuracy in tagging speech transcripts (Horsmann, Erbs, Zesch, 2015). Because my analysis is minimally dependent on using tags, the inaccuracy does not impact my results. The tagged search function was used in case of:

- Demonstrative determiners and demonstrative pronouns: I used the DT tag (determiner) which yields a result of concordances with both demonstrative pronouns and determiners.
- Hypothetical reported speech: looking for optional modal verbs in a phrase: *personal pronoun + modal verb (MD) + say*

**Automatic tagging error adjustment**

When a more sophisticated method of analysis is required involving tagged parts of speech, I follow the procedure of error analysis presented in Biber, Conrad, Reppen (1998, p. 91). There are four steps in the process:

1) Extraction of a random sample of the linguistic feature that is automatically tagged, but needs to be examined manually.
2) Analysis of the accuracy of the tags in the sample.
3) Computing the proportional use of each grammatical category in the sample.
4) Multiplying the total number of occurrences obtained before the hand-editing procedure by the percentages obtained in the error-analysis process.

This procedure is used for discerning *this* in the role of demonstrative, rather than other grammatical function (i.e. relative pronoun) in the analysis of deixis.

**Analytical framework**

My study adopts the register analysis framework outlined by Biber and Conrad (2009). The overarching goal of register analysis is to analyze language variation – how English language is use differs depending on the context of use, for example speech vs. writing, e-mails vs. singing telegram. My study focuses on a variation in a small subset of linguistic features that are used for engaging audience in two registers: online (KAM) and face-to-face (MITM) instruction in math. The analysis begins with predetermined set of linguistic characteristics that emerge from literature review on engagement features used in spoken academic discourse. These features are used in performing a register analysis in a very narrow scope.

The use of register analysis framework requires three major steps. The first step is describing the situational characteristics of the register, which allows a for more robust description of the participants and the situation of the register use. This step also helps to identify "distinctive aspects of the context and communicative purpose" (p. 50). The second step is analyzing idiosyncratic linguistic characteristics of the register, which in my case focuses on selected linguistic features that have particular importance for the study of engagement, defined as involvement and interactivity. The final step is interpreting how the relationship between the situational and linguistic descriptions work to engage the audience in this particular communicative situation.

In their discussion of selection of register features, Biber and Conrad (2009) point to the need of comparative approach, in combination with quantitative analysis of representative samples (p. 51). Other scholars (Sinclair, 2001; Haarman et al. 2002, Crawford Camiciottoli, 2008) also posit that in comparing registers, the distinguishing features of each registers can become apparent. There are three key aspects of the comparative approach used in my study. First, I use two corpora which represent a very similar domain - mathematical instruction at an advanced high school and college level. The similarity of instructional purpose makes the case for comparison between these two corpora. Since the topic of instruction is quite similar, the differences that emerge in the analyses may be attributed to the modality with more confidence (face-to-face or online) than if the domains were different. In my analysis, however, I acknowledge that the unequal comparison between the number of speakers represented in the corpora has an impact on the results of my research.

In this study I also rely on the results of previous comparisons between different spoken academic registers to discuss the use of engagement features that are particularly pertinent to spoken academic register. Selecting a few of the same linguistic features as the ones used in previous studies allows me to contrast my results from two very particular registers - Khan tutorials and MIT lectures - to results obtained in previous studies on the same features. Thus, my research is positioned to add context the previous research done on these features with the use of corpus-based methodology.

**Corpus-based study of register variation: investigating association patterns in language use**

As described in the introduction chapter, my research methods are based on the corpus-based approach to register analysis (Biber & Conrad, 2009). Using a set of pre-selected linguistic features used for engagement in academic discourse, I analyze how speakers in both corpora use these features to engage their audience. The size of the corpora, which exceed 1 million words, provide a good sample of the language from both register and have the potential to be mined for variation in engagement feature use.

This corpus-based study of register variation relies on investigating association patterns in language use (Biber, Conrad, & Reppen, 1998), by examining the frequency of association of the particular linguistic features (i.e. personal pronouns, modals and semi-modals) with other linguistic and non-linguistic features.

Any analysis of a complex association pattern relies on the premise that linguistic features are not isolated entities, but rather they tend to associate with other linguistic or non-linguistic features (Firth, 1957; Biber, Conrad, and Reppen, 1998, p. 5). There are two main approaches to researching association of patterns in language use. First, examining the use of a linguistic feature (lexical or grammatical) by investigating its linguistic (lexical or grammatical) association or non-linguistic association (across registers, dialects, time periods (p. 6)). Second, investigating a variety (or varieties of texts) by examining linguistic association patterns. These patterns can either be examined by looking at one feature or classes of features, or by examining co-occurring patterns of features (p. 6).

In my dissertation project, the second approach is used to investigate the registers of math instruction delivered online (KAM corpus) or face-to-face (MITM corpus). In particular, I am investigating two types of lexical linguistic associations, with a minor use of grammatical

associations. Using predetermined involvement and interaction  I examine how these features are systematically associated with particular words (*lexical association*) or multi-word units. The primary analysis, however, focuses on lexical association. Because patterns of language use vary between registers, it is important to compare a pattern across register to determine the difference and provide analysis of the possible explanations in terms of the functions behind the differences.

**Top-down and Bottom-up Linguistic Feature Functional Analysis**

*Top-down Analysis*

The top-down analysis of linguistic features involved selecting the items for investigation. In the group of features used for investigation, there were two main categories: unambiguous and ambiguous group. The first group included features that can be searched for reliably without the need for grammatical tagging (i.e. personal pronouns, deixis – excluding *that*, and response elicitors). The second category included one item – *that* – which can function as a demonstrative or a relative pronoun. Using error adjustment methods, I excluded the use of *that* as a relative pronoun, because it would impact my results (i.e. overestimating the use of *that* for deictic function).

Figure 8. Top-down frequency-based analysis of linguistic patterns in the corpora.

After obtaining an overall frequency count for each feature, I performed lexical analysis (analyzed randomly selected 100 concordance lines) and phraseological analysis, looking for patterns that are used for engagement. Patterns that stood out specifically in comparison between two corpora were investigated in depth.

Functional analysis involved analyzing concordance lines (10% of the concordance lines for particular phraseological pattern) and interpreting the results in light of the description of situational characteristics.

### *Bottom-up Analysis*

A bottom-up approach was used to look for linguistic patterns used by speakers to engage audience. This approach begun with close reading and annotating linguistic patterns that are used for audience. 10 transcripts were closely read and annotated in this process. There was one pattern that was noticed and selected for investigation – direct hypothetical reported speech.

Other patterns were not selected to be included in this study, because of limitations on the scale of the research for this dissertation.

A corpus search was performed to elicit data for concordance line instances of direct reported hypothetical speech (DHRS) Using iterative process of searching for previously established lexico-grammatical patterns typical of DRHS, I analyzed a variety of phrasal combinations to generate the most comprehensive list of concordance lines that could potentially include DRHS. I appended the phrasal searches based on any items that signaled DRHS but were not represented in previous literature. I combined all of the concordance lines that potentially included DRHS in both corpora and manually marked each instance as including or excluding the presence of DHRS in an Excel spreadsheet. As I manually categorized these instances, I also made notes on their function. I considered a response elicitor to be formed with one word or two word question.

**Frequency-based lexical analysis**

My dissertation study explores register variation in terms of engagement language. There are two types of analysis I conduct: a lexical description of linguistic feature use and phraseological description (Staples et al. 2015). The selection of the features to investigate is based on their relative frequency: the similarity or difference in normed frequency count of the feature in both corpora. It is important to note that the word *frequent* in corpus linguistics is always dependent on the corpora that are being used in the study (Hunston, 2010) There are, however, benchmarks established in previous research for word frequencies or multi-word units that I describe below. In the *University Language* (2006) study, Biber presented a three-tier band for word types by frequency level (p. 36):

Very common words – occurring 200 times per 100 000 words (or 20 per 10 thousand words)

Moderately common words – between 21 and 200 per 100 000 words (or 2-20 per 10 K)

Rare words – fewer than 20 times per million words (or fewer than 2 per 10 thousand words).

Multi-word expression cut-off frequency in previous studies (Biber, Conrad, & Cortes, 2004; Biber, 2006) is set to occurrence of 40 times per 1 million words (or 0.4 per 10 000), although there are expressions that are much more frequent that than in spoken academic corpora. In my corpora, lexical bundles have to be used by at least two speakers in the MITM corpus or in at least 10 tutorials to prevent from speaker or topic idiosyncrasy to skew the results.

**Lexical description of words**

The overarching method of selecting linguistic features for qualitative investigation is based on a difference between normed frequency count of each linguistic feature used for engagement. This feature, a single word such as a personal pronoun or a modal verb, in both corpora. The procedure for lexical description I used is as follows:

1. Single word item normed frequencies are analyzed and reported for each corpus.
2. Normed frequencies are compared between the two corpora and to the results from previous studies on spoken academic discourse (if such results exist).
3. The linguistic items are analyzed in context – 100 or 10% of the results (whichever is fewer) randomly selected concordance lines for each linguistic feature were analyzed for function in context.
4. The functional analysis is concluded with contextualizing how linguistic items function similarly or differently based on the situational characteristics of the register.

**Phraseological descriptions**

Research on phraseology in register studies centers on the frequency, form, and function of multi-word units (Staples et al. 2015), which is also the focus of my analysis. There are three steps in my multi-word unit analysis:

*Multi-word unit search.*

For each category of engagement linguistic features - personal pronouns, deictics, demonstratives – I performed a search of multi-word units which contained the lexical item. For example, in order to understand how the use of personal pronoun *I* is used in both corpora, I performed a search of 4-5 gram. This shows the use the personal pronoun *I* in the context of 3-4 other words in a unit that is occurring frequently, more than 40 times per 1 million words (or 4

times per 10 thousand words). An example of 4 gram with personal pronoun *I* is *I encourage you to*.

***Examining collocates of the unit and their function in terms of engagement.***

I have examined each multi-word unit by randomly selecting 100 concordance lines, using random selection function in SketchEngine. I examined its function in concordance lines, making notes on potential engagement function. Additionally, I used SketchEngine frequency feature to analyze first word to the right to see what are typical word forms that follow a multi-word expression. For example, multi-word expression *I encourage you to*, is frequently followed by verbs *pause, watch, try.* These follow up procedures usually yielded smaller number of concordance lines for examinations, such as *watch* and *try* in Figure 9, which were easier to examine in terms of their functions. In case of multi-word expressions that yielded frequency larger than 100, I examined 100 randomly selected concordance lines.

| | Word | ↓ Frequency | Frequency per million | |
|---|---|---|---|---|
| 1 | pause | 204 | 114.96 | ... |
| 2 | watch | 22 | 12.40 | ... |
| 3 | try | 18 | 10.14 | ... |

Figure 9. Frequency feature of SketchEngine.

These results are analyzed for their potential uses in engaging the audience. I investigate these function by examining the function of word units in the concordance lines.

**Dispersion**

For each feature, I also check the dispersion rate, to avoid one speaker (in case of MIT) or one course skewing the results of the analysis. Dispersion (Gries, 2008) is the measure of the spread of a linguistic variable across texts in a corpus. In my dissertation study, I use a SketchEngine feature *Visualise* to check whether my search term is evenly distributed in the corpus. In the two figures below, the terms *already* and *hey* are visualized. The visualization process depends on dividing each corpus into parts - in the case of both figure below it is 100 parts. Then, then relative frequency of the search words in each part is displayed in the column.

Figure 10. Distribution of the adverbial *already* in the MITM corpus.



Figure 11. Distribution of *hey* in the MITM corpus.

The adverbial *already* is relatively evenly dispersed in the corpus, while the attention signal *hey* is dispersed unevenly. This uneven dispersion would prompt further investigation on speaker or course-related reasons of the differences.

More sophisticated measures of dispersion (Julliand et al., 1970) are not used in the study, since the more rudimentary visual approach is sufficient for the purpose of this research.

**Selection of Linguistic Features Used for Engagement**

I define engagement in this study as the use of features in academic instruction that invite were previously explored in studies of interaction and involvement. Linguistic features typical of involved and interactive academic speech were selected from the oral and involved discourse dimension (Biber, 1988; Biber, 2006) as well as other studies on engaging academic audiences.

I rely on the definition and linguistic features used for involved production identified through a type of statistical factor analysis – multidimensional analysis (MDA) - by Biber (1988; 2006). Using a multi-register academic corpus, which includes spoken and written registers, overall patterns of register variations were analyzed (Biber, 2006; p. 177) to discern dimensions through factors analysis. A dimension describes a relationship between co-occurring patterns of pre-selected linguistic features. The registers under investigation were classroom teaching, class management, labs, office hours, study programs, service encounters, textbooks, course packs, syllabi, institutional writing. In the case of the 2006 study, there were 90 linguistics features, out of initial 129, that were used in the MDA to show meaningful differences between the registers. In Biber (2006) study, the first dimension – involved versus informational production – grouped features that correlated with each other highly to create a set of features used for involved production, and features that were negatively correlated used for informational production. Based on the situational analysis and the focus of this study, I have selected the linguistic features that have been identified by Biber (2006) as relevant for analyzing involvement. Following the design of Csomay (2000; 2002), I selected feature groups relevant for involved aspect of the register I am investigating, rather than investigating every feature in that dimension. Alongside the selected features that are associated with involvement, I also selected feature groups associated with interaction (response elicitors) and stance. These linguistic feature groups considered together constitute the construct under investigation in my study: math instructors' engagement with their audience.

Previous multidimensional studies aim to provide a much needed broad characterization of academic registers, but are not used for extensive in-depth analysis for single linguistic feature use in context. My study is an exploration of how such linguistic features which are marked as functioning together to signal involvement or oral production, are used in context to engage audience in the KAM and MITM corpora.

**Personal pronouns (I, we, you)**

I conducted lexical and phraseological analysis and description for personal pronoun in both corpora. I used previous research on the use of personal pronouns in conversation and academic prose (Biber et al., 1999) and Biber (2006) as a reference for determining normed frequency of personal pronoun occurrence.

Table 10. Personal pronoun frequency across registers.

| Personal pronouns (reported per 10 thousand words) | Conversation (Biber et al. 1999) | Academic prose (Biber et al. 1999) | Classroom Teaching (Biber, 2006) |
|---|---|---|---|
| *I* | 370 | 20 | - |
| *we* | 70 | 30 | - |
| *I+we* | 440 | 50 | ~460 |
| *you* | 300 | 10 | ~447 |

For pronoun *you* and *we*, I additionally extracted 100 concordance lines and coded them based on the type of use as direct reference to audience or generalized you as an indefinite or impersonal subject (see Friginal et al., 2017). Both pronouns can have literal meaning, when the instructor is addressing his or her particular audience, or have a general meaning, in which both pronouns could be substituted with the noun *people*, *a person* (i.e. *You can eat broccoli every day and feel fine.*). In borderline cases, I coded the use of the pronoun as general reference.

**Deictics**

My analysis of deictics focuses on frequency analysis of adverbs of location (*here* and *there),* demonstrative determiners (*this/that* book), and demonstrative pronouns (*this/that).* I report on the normed frequency of all of the above features. I exclude the use of existential there from the analysis (i.e. *There was a cat, F.D.C Willard, which co-authored a physics paper in 1975.*) as well as using *that* in function other than deixis. *That* was only analyzed in the function of a determiner and a pronoun, excluding the relative pronoun use and the use in complement clauses. Manual error adjustment procedure (Biber, Conrad, & Reppen, 1998, p. 91). was used to obtain frequency count and excluded instances of *that* as a relative pronoun from concordance line samples.

A 10% sample or 100 concordance lines (whichever was smaller) of the results are analyzed to distinguish between spatial deictic, discourse deictic or non-deictic function. Functional analysis is performed in terms of special uses for engaging audience.

**Response Elicitors**

Response elicitors are clearly a group of linguistic features used for engagement. In order to find response elicitors in my two corpora, I performed the query searching for one or two words followed by a question mark

[word="\."][ ]{1,2}[word="\?"]

Most of the response elicitors present in spoken discourse are discourse markers followed by a question mark, such as *right, okay*. The quantitative analysis includes frequency count per 10 000 tokens. In the results, I included only the elicitors with frequency of 0.5 or more per 10 000 tokens. The second step in analysis, was functional coding for patterns of engagement of 100 randomly selected concordance lines.

Additionally, in the MITM corpus, I looked specifically for discourse markers that elicited responses in the lecture audience. The audience turns are marked by the use of tags: <audience>, <inaudible>, or <student>. I did not analyze the audience turns, but analyzed which response elicitors prompted anyone from the audience to respond to the instructor. The same function could not be performed, for obvious reasons, in the KAM math corpus since it does not have live audience.

**"You might say, hey:" Direct Hypothetical Reported Speech**

The only feature that is explored with a bottom-up approach is direct hypothetical reported speech. As I was examining personal pronoun use, I noticed the phrases that invoke reported speech to paraphrase what . In case of Sal Khan, the use of the discourse marker/attention signal *hey*, is invariably linked with the hypothetical reported speech, as is addressing himself in the third person (*hey, Sal*). Using these utterance-launchers (Biber et al. 1999, p. 1118-19), I performed a series of searches to compile the most complete list of KWIC lines that were suspected of containing hypothetical reported speech. My aim was not to extract

all of the instances of DHRS, but a most representative sample, that is why I relied on iterative process that would yield most often used phrases that initiated DHRS.

First, I searched for sentence fragments that contained both the pronoun *you* and the discourse marker *hey.* This first search revealed that Khan uses the verbs *say, think, wonder* to introduce the hypothetical thoughts and statement of his students. This search resulted in adding other utterance-launchers, such as *well, okay, look* (as listed in Biber et al. 1999, p. 1118-19). The features were almost always in proximity of personal pronouns *you* and *we*, and very often contained DHRS. Instances of DHRS are also present when Khan talks about himself in third person (i.e. "So you're probably saying, hey, Sal. What is all this opposite, hypotenuse.." KAMDoc#1406). Following up on the results of sociolinguistic research in direct speech use a more recent trend in using utterance-launchers is the construction *be + like*, which is also present in the KAM:

> then you just look at this expression right here, **you're like , hey, what's wrong with x?** [doc#638KAM]

Qualitative analysis focused on the right context of the pattern *pronoun + optional modal verb / other verb + say / think / wonder / ask + discourse marker + message* (Koester, 2018). The next step was refining and compiling all of the results from both corpora into a spreadsheet and analyzing the purpose of the DHRS in the math instruction, with aim of involving the audience. I analyzed both spreadsheets to exclude instances of concordance lines that did not include hypothetical reported speech.

In case of searching for DHRS in KAM corpus, I relied on fact that there is just one speaker represented in it. Khan reuses certain phrases, marking his idiolect (Dittimar, 1996) and making searches more effective. The research on idiolects using corpus linguistics methodology show that individual lexical and syntactic patterns appear to be stable over a period of 1-2 years and are maintained even when they differ from the communities (Barlow, 2013). However, there is limited research especially on spoken idiolects, with much of research focusing on written corpora, especially in authorship analysis (Barlow, 2013). The same search was performed in the MITM corpus.

The qualitative analysis of the KWIC was performed to search for patterns that marked the DHRS as related to process/application (apply/analyze) narration or knowledge/factual

(understand) narration. scheme was based to limited extent on the revised Bloom's taxonomy of Knowledge Dimensions (Krathwol, 2002). An example of knowledge/factual (KF) hypothetical reported speech is:

> is going to be equal to the square of the longer side, or the square of the hypotenuse. And if you're not sure about that, **you're probably thinking** **, hey Sal, how do I know that a is shorter than this side over here? How do I know it's not 15 or 16?** And the way to tell is that [doc#73KAM]

In the case of KF-DHRS Khan is pointing out basic facts or concepts in math, often explaining a definition or in the case of the example above an observation that pertains to basic facts that are used for more advanced processing or calculations. An example of process/application (PA) direct hypothetical reported speech is:

> to this one right over here. If we were to make them consistent, if you were to make this definition consistent with this, **you would say** **hey , let's start with a 1**, and then multiply it by 1 eight times. And you're still going to get a 1 right over here. [doc#15KAM]

I do not report the frequency of patterns, since it is not the goal of this study. Rather, these patterns are used to provide an example of register-specific engagement strategies.

# CHAPTER 5: RESULTS AND DISCUSSION

## Linguistic features of involvement

The analysis of linguistic features used for involvement included personal pronouns (*I, we, you)* and deixis (demonstrative pronouns, determiners and location adverbs).

## Personal Pronouns

The search results conducted with SketchEngine are presented in Table 11. The results were normed per 10000 tokens to represent the difference in the frequency of the pronouns reliably. The use of "i" that often appears in both corpora as a letter used to represent an unknown number was excluded manually and the occurrence rate was recounted.

Table 11. Occurrences and Frequency per 10 Thousand Tokens of Personal Pronouns

|  | Frequency (per 10 thousand tokens) | |
| --- | --- | --- |
|  | **KAM** | **MITM** |
| First Person Pronouns | | |
| I | 106 | 185 |
| me | 25 | 17 |
| my | 9.7 | 10 |
| Total | 140 | 212 |
| we | 204 | 126 |
| let's | 6.5 | 3.3 |
| our | 30.2 | 10.1 |
| us | 13 | 4.6 |
| Total | 253.7 | 144 |
| Second Person Pronouns | | |
| you * | 131 | 148 |
| your | 9.2 | 9.9 |
| Total | 140 | 159 |
| Pronouns Total | 534.6 | 513.9 |

Data in Table 11 shows that personal pronouns are a frequently used in math instruction, irrespective of the mode. In both corpora, the pronouns constitute slightly more than 5% of the total word count. There is a difference in the frequency of personal pronoun use by the speakers in these two corpora. While *I* is used most frequently in the MITM corpus (185 times per 10 thousand tokens), we is used most frequently in KAM (204 times per 10 thousand tokens). The

rate of occurrence of the personal pronoun *you* is similar in both corpora (131 and 148 times per 10 thousand tokens). In comparison to conversation register (Biber et al., 1999) and classroom teaching (Biber, 2006), both corpora use fewer personal pronouns. In comparison to frequency of *I* and *we* use of in the Biber (2006) study (460 occurrences per 10K words), KAM corpus includes fewer instances (306 instances per 10K tokens), as does MITM corpus (311). The reason for this difference are situational characteristics since T2KSWAL sampled three types of interactive classroom sessions (low, medium, and high interactivity), while my corpora are focused more on monologic academic discourse which would fall in the low interactivity category (MITM) or not interactive (KAM) in terms of turn takes. However, a detailed analysis of the pronoun *you* and *we* is required to discuss how these personal pronouns function in my corpora, since they are used very often in both.



Figure 12. Normed frequency of personal pronoun occurrence per 10 000 tokens in KAM and MITM corpora.

Contrary to results from previous studies on the language of lectures (Fortanet, 2004; Lee, 2009; Cheng, 2012; Friginal et al., 2017) and similarly to a small scale study in math instruction by Rounds (1987a, 1987b), the pronoun *we* is the most frequently occurring in the KAM corpus. Lee (2009) proposed that using *we* more often in his dataset (large-class lecture

introductions) serves the purpose of reducing the affective and physical distance between the lecturer and the listeners. In the context of Khan Academy tutorials face even greater threat in terms of affective and physical distance between the instructor and students, thus the overwhelming use of *we* can be read at the direct reaction to the challenge of detached online instruction. Khan needs to overcome the depersonalized feel of online instruction, in which the students are interacting with technology interface, rather than physical instructor presence. Such frequent use of *we* can be read as an attempt to mitigate these situational characteristics.

Pronoun *we* can be particularly important in showing teacher involvement and creating student engagement in mathematics. Through the use of *we* the speaker is building a feeling of joint work between the listeners and the speaker (Flowerdew & Miller, 1997). Since the issues of *self-efficacy* - student's belief of ability to complete task on their own (Bandura, 1987) - and general confidence are linked to the performance on math tests (Hackett, 1985), taking off the responsibility from the *you* as the student and shifting it to the collective *we* can potentially decrease the anxiety of the student and increase sense of confidence in completing math tasks.

**Personal pronoun *I:* lexical description**

The high frequency of the use of pronoun *I* in the MITM lecture corpus is in line with previous lecture discourse studies (Fortanet, 2004; Lee, 2009; Cheng, 2012; Friginal et al. 2017). There are context-specific and mode-specific reasons for the frequency of use of pronoun *I* in the face-to-face modality. The physical presence of audience and clear division of roles when discussing course organization contributes to positioning of the speaker who refers to themselves with the use of *I* and the audience – *you.*

Khan, the sole speaker in the corpus, prefers to use *we* more often in his speech. However, *we* is often used as an *I* referent, which means that *I* could be substituted for *we* without substantial change in meaning. A functional analysis included later in this chapter reveals that the personal pronoun *we* in the KAM corpus is very frequently used in the narration of mathematical reasoning. Each step in the Khan's reasoning is referred as if it is done by the inclusive *we* – the instructor and the audience. Since in MITM there is a physical audience present, the distinction between *I* and *we* is related to the on-line production and circumstances, in which *we* means a delimited number of people. In one study, Dafouz, Nunez, and Sancho (2008) noticed discrepancies in preferences in personal pronoun use among three lecturers,

suggesting that personal pronoun use might be dependent on teaching style. In analyzing these results, the speaker representation also needs to be taken into account. KAM corpus represents just one speaker's teaching language, so it is likely to be biased towards his teaching and communication style.

In comparison to the results from Biber (2006) who reported that 1st person pronouns are used with the frequency of 40 per 1000, the results from both corpora are similar, 31 personal pronouns per 1000 words. However, what is not accounted for is the difference in the distribution of the pronouns *I* and *we.* The differences in the frequency of personal pronoun *I* use in both corpora prompted a functional analysis. One method of exploring how this pronoun is used in both corpora is the analysis of bigrams.

Table 12. Personal pronoun I bigram normed frequencies (per 10 thousand tokens)

| KAM | Normed Frequency | MITM | Normed Frequency |
|---|---|---|---|
| I am/ I'm | 17 | I am/I'm | 29 |
| I will | 10 | I have / I've | 13 |
| I have | 8 | I want | 10 |
| I could | 5 | I will/I'll | 9 |
| I want | 5 | I do(n't) | 8 |

The bigram analysis provides more context on the use of *I* in both corpora, specifically the differences in *I* lemma context. In both corpora, the most frequent collocate is the verb be in the present tense: I am/I'm. However, in the MITM corpus, the collocation *I am* occurs almost twice as often, which provides evidence for stronger focus on the instructor actions in the MITM corpus. The main uses of I am relate to *chalk talk:*

> For the null space, **I'm looking at combinations of** columns to get the zero column. [doc#1MITM]

> If I have three triples-- if I have a triple xi, xj, xk,     **I'm going to convert that** into a number that looks like this where the one positions are [doc#56MITM]

> So, **I'm going to define** Y_n to be two to the power of X_n. [doc#188MITM]

The MITM instructors also use I am to reflect on conceptual aspects of problem solving:

If the right-hand side is the sum of the functions, well, so is the left. But **I'm saying it the other** way around. If the left is an even function, why does the right-hand side have [doc#52MITM]

then I allowed to say -- what's the matter with this argument? That gave us the constant term eight. It's wrong. **What I'm going to write up is wrong. I'm going to say** Bx is alpha x. Add those up, and you get A plus B x equals lambda plus alpha x. [doc#84MITM]

The use of the personal pronoun *I* in MITM corpus is used mostly for chalk talk – the think aloud the math instructors do as they solve the problem. The preference for *I* positions the instructors at the center of informational space (Kamio, 2001), in the position of authority and action. Other bi-grams, *I have, I want, I will* when qualitatively analyzed in a random selection of 100 concordance lines, are almost exclusively used to narrate problem solving:

So capital S is going to grow to include that node. **I've extracted** it from the queue. [doc#20MITM]

**And I want to produce** out of that q 1 and q2, I want to produce orthonormal vectors. [doc#55 MITM]

Unlike the results from O'Boyle (2014), in which she found *I think* to be the most frequent bi-gram in multi-party classroom context, in my monologic corpora *I'm/I am* is the most frequent. In O'Boyle (2014) research, the students and instructors interacted with one another, creating opportunities for expressing stance and seeking to reach mutual understanding through the use of *I think* (Fung & Carter, 2007). MITM lectures do not serve the purpose of teaching through interaction, but rather modelling problem solving methods. Since MIT is also a well-established, top-tier STEM school, the status of the professors as experts in the fields is embodied in the language they use, focusing on their expertise.

The second most frequent bi-gram in KAM *I will,* even though has a similar frequency to MITM corpus bi-gram, shows that apart from narrating problem solving, Khan uses *I* to signal next steps in the problem solving with *I will do, I will just write/rewrite.* Such use of this modal verb prepares the audience for the next step of mathematical procedure. As Khan uses the virtual board, he signals plans with *I am going to* (5 times in 10 thousand tokens) or *I will* (10 times in

10 thousand tokens). Both in MITM and KAM corpus the use of the semi-modal *I'm going to* is evidence of involvement and recognition of the needs of the audience. The speakers prepare the audience on what is going to happen next in mathematical narration.

**Personal pronoun *I*: phraseological description**

Description of the frequency of occurrence of personal pronouns and bigram analysis does not provide a comprehensive information on how these pronouns are used by the speakers to engage the audience in the video or the lecture. In order to investigate the difference between the use of these pronouns in both corpora for engagement, I generated a list of 4-5 multiword expressions (with 0.5 normed frequency as threshold) and analyzed their engagement functions in concordance line context.

Table 13. Personal pronoun I in multi-word expressions.

| KAM | Normed frequency (per 10 thousand tokens) | MITM | Normed frequency (per 10 thousand tokens) |
|---|---|---|---|
| I encourage you to… <br> - Pause <br> - Watch <br> - Try | 1.8 | I am going to | 13 |
| if I were to <br> - take <br> - draw <br> - say <br> - ask <br> - multiply | 1.5 | I don't know <br> - Whether <br> - Why <br> - How <br> - Exactly <br> - If | 1.7 |
| I want to do | 1.4 | I want to do | 0.8 |
| I don't know | 1.4 | I don't have | 0.7 |
| what I want to <br> - do | 1.4 | I don't want | 0.7 |
| I don't want to | 1.4 | what I want to | 0.6 |
| I guess you could | 0.6 | that I want to | 0.6 |
| see if I can | 0.6 | I have to do | 0.5 |

The most frequent multiword expression containing the pronoun *I* in KAM "I encourage you to pause…" [the video], which is a direct reference to the audience that suggest an action - pausing the video. Because of the modality of the instruction, a video recording which can be controlled by the audience, the listeners are encouraged to perform a task on their own or to consider the material or a problem that has just been presented:

**I encourage you to** pause this video and *think* about it.[doc#1107KAM]

**I encourage you to** pause the video and try to *work* through it on your own. [doc#463KAM]

And **I encourage you to** pause this video and *give it a go* on your own. [doc#1023KAM[

The frequent use of this multi-word expression in the corpus suggests that Khan designs his videos with his audience and their learning process in mind. Similarly, to a class in which an instructor might intertwine moments of lecturing and demonstration with student practice, Khan wants the students to respond to his request and perform an action or a learning task. The use of this lexical item suggests that Khan is aware of his audience and their needs, as well as their thought and learning processes. Lee (2016) pointed out that teachers use the pronoun *you* when they give instructions, just as is the case with the lexical phrase "*I encourage you",* which in turn can keep the students involved and engaged.

Another important multi-word expression is *if I were to + verb* which is used to introduce a problem or a scenario that will be solved by Khan:

**If I were to tell you,** that the absolute value of x is less than 10 what does that mean? That means that the distance from x to 0 has to be less than 10. So **if I were to draw** a number line and put 0 here we can only go up 10 away and even that's too far. It has to be less than 10. [doc#366KAM]

And so **if I were to ask you**, what is the probability-- I'm going to flip a coin. And I want to know what is the probability of getting heads.[doc#1277KAM]

The multi-word expression *if I were to + verb* introduces subjunctive mood, a hypothetical scenario that engages the audience (i.e. if I were to ask *you)* by offering a proposition. The personal pronoun *you* follows the phrase 1 in 10 times, making the audience as

an object in the discourse. Even when Khan does not directly address the audience, using a hypothetical scenario is an engaging device. Conditional clauses require a more active participation from the audience, as they ask them to imagine a scenario. Biber et al. (1999) suggest that conditional clauses make commands or suggestions less forceful (p. 821), which can increase positive politeness and build more engagement. Rather than telling his students "What is the probability of…?", Khan asks "If I were to ask you what is the probability" (#1277), the question becomes less direct and threatening. Another use of this conditional phrase is to construct a discourse around attempts and trials:

> The easiest way to think about it is: If I want one X on this left-hand side, that is a third of the total X's here. So what **if I were to multiply** the left-hand side by one-third -- -- but if I want to keep the scale balanced, I have to multiply the right-hand side by one-third. If we can do that mathematically….[doc#132KAM]

As Khan solves the mathematical problems, he presents them as an exploration of possibility, modeling to his students a less commanding and authoritative identity. The same modelling of openness shows in the use of phrase [*let's*/*let me*] *see if I can*, which again shows Khan's own involvement. By narrating his learning process, modelling of his thinking as he solves the problems, he reveals his own persona and voice in teaching.

In comparison to KAM corpus, the multi-word expressions in MITM corpus follow the patterns found in Biber (2006) study. The most frequently used pattern is *I don't know + wh/if clause* which is a personal epistemic stance phrase. The use of that phrase reveals the instructor's stance in the teaching discourse, expressing uncertainty. There are instances when the instructor expresses hesitation in reference to the needs of the audience

> That's our question. Let me, **I don't know if you** see the answer. Whether there's -- so let's see. I guess we could do it properly. [doc#19MITM]

> So today we will embrace our inner Cookie Monster and eat as many-- eat the largest cookie first, would be the standard algorithm for Cookie Monster. **I don't know if you** learned that in Sesame Street, but-- all right [doc#20MITM]

Because the audience is physically present, even indirect questions about visibility of material or previous knowledge can be answered by non-verbal responses from the audience. The use of

personal pronoun *you* is a direct reference to the audience present in the classroom, with some of the instances

## Pronoun you: lexical description

The normed frequency of occurrence for the pronoun *you* shows quite similar use levels in MITM and KAM corpus. Since *you* can serve as a direct reference to the listener (*Do you have any questions*?) or be used in a general sense (*You would take this a and add it to the b to get the c.* ) in mathematical narration, to mean "a person" Qualitative analysis of the function of *you* in a randomly selected 100 concordance lines in terms of their reference to audience or indefinite / impersonal object. In the process of coding, I also included coding for hypothetical reported speech. The indefinite use of personal pronoun *you* was strictly related to mathematical reasoning narration.

Table 14. Functions of you in both corpora in 100 randomly selected concordance lines. Raw count.

|                                      | KAM | MITM |
| ------------------------------------ | --- | ---- |
| Direct reference to audience         | 15  | 23   |
| Hypothetical reported speech         | 3   | 0    |
| Idiomatic expression "you could say" | 5   | 1    |
| Mathematical narration               | 77  | 74   |
| Invalid transcription (you -.your)   | -   | 2    |

The use of *you* in direct hypothetical reported speech is discussed at length in section 5.3. In case of the MITM corpus, two concordance lines had typos (you was substituted for yours), and the selection was invalid. I also singled out *you could say* as an idiomatic phrase that introduces example, but is more fixed lexically than any other uses and should be counted separately.

## Direct reference to audience

Both corpora include direct turns aimed at the audience to a similar degree. The examples in MITM corpus refer to the students' circumstances – being present in the previous lecture, registering for other courses, remembering what was said in previous lectures (i.e. *Do you remember remember Perron-Frobenius theorem*?). These direct references could either elicit a

verbal or a non-verbal response from the audience. However, it is important to note that Khan also acknowledges audience as actual people who are interacting with the videos:

> Or are these two quantities equal? Or is there not enough information to tell? So like always, pause this video **and see if you can work through** it on your own and now I will work through it with you. All right, so let's just write down the information [doc#295KAM]

In this instance, Khan asks the audience to work on the problem, after he introduced the question. This is a direct request to an actual student, using *you* as a referent to a person listening to the video. This pronoun is used when Khan addresses the listeners so they engage in performing a task that requires them to use time off the video. He also uses *you* in reference to his viewers' longitudinal learning development

> of two numbers that, when I multiply them I get 50, and when I add them, I get 15. And this is going to be a bit of **an art that you're going to develop**, but the more practice you do, you're going to see that it'll start to come naturally. So what could a [doc#471KAM]

Such use of pronoun *you* is used in the most straightforward manner of engage the audience. Since *you* belongs to the hearer's territory, it implies the distance between hearer and the speaker (Kamio, 2001). As the speaker directly refers to *you* as the audience, he or she reaches out across the speaker space and acknowledges, and invites the audience to respond. The higher the use of *you*, the more interactive is the classroom (Hyland, 2009). This is not surprising, since the basis of interaction is communication in which all the parties acknowledge each other's positions addressing each other.

### *Chalk talk: mathematical narration*

The most frequent use of *you* in both corpora relates to an impersonal narration of mathematical reasoning. As the instructors go through problem sets, they might choose *I, we,* or *you* to walk the audience through a set of procedural steps. When the instructor is using *you* they do not directly refer to a particular student or an action that is happening at the moment of speech. Rather, they hypothesize about the problem solving procedure as done by undefined person.

here you have some big O of 1. **You** are probably doubling the constant in there every time **you** do this relation. If **you** have a finite number of doubling of constants, no big deal, it is just a constant, two the power number of doublings. But **[doc#123MITM]**

So, let's see, we get 2 minus 12 is negative, is negative 10 plus c equals negative 10. So you add 10 to both sides you get c in this case is equal to 0. So we figured out what our position function is as well. The c right over here is just going **[doc#972KAM]**

Despite the fact that this is a generalized use of *you* as a subject, the fact that *chalk talk* can also be realized with the use of *I* (see the previous section), suggest that *you* might be functioning as a less authoritative and more engaging manner of chalk talk. In both corpora the use of *you* in problem solving narration to some extent leads the audience through the problem. In KAM, the phrase *you're going to get* (1 occurrence per 10 thousand tokens) signals result of the problem. The speakers are making a choice to use *you*, rather than authoritative *I* or collective *we.*

The analysis of bi-gram reveals, that Khan has a preference to use modal of possibility: *you could* rather than *you can* which signals ability in the MITM corpus. Just as it was the case with the pronoun *I*, Khan shows preference for modal verb *could* with the pronoun *you.* Additionally, the phrase *you are going to* (3 occurrences per 10 thousand tokens), shows that Khan recognizes the audience's presence and signals to them the results or steps they should have gotten to.

Table 15. Bi-gram with personal pronoun *you* in KAM and MITM corpora.

| KAM | Normed frequency (occurrence per 10 thousand tokens) | MITM | Normed frequency (occurrence per 10 thousand tokens) |
|---|---|---|---|
| you could | 13 | you can | 13 |
| you're | 13 | you have | 12 |
| you have | 11 | you're | 9 |
| you can | 9 | you don't | 6 |
| you get | 6 | you want | 4 |

**Pronoun *you*: phraseological description**

The analysis of multi-word expressions with personal pronoun *you* reveal that there is a more diverse set of fixed phrases used in KAM corpus (15) than in MITM corpus (8). The reason for this diversity of fixed phrases in KAM corpus can be the function of having one speaker with a limited repertoire of phrases being represented in a large corpus.

Table 16. Pronoun *you*: multi-word expressions with frequency 0.4. or greater per 10 thousand tokens.

| KAM | Normed frequency (occurrence per 10 thousand tokens) | MITM | Normed frequency (occurrence per 10 thousand tokens) |
|---|---|---|---|
| you're going to | 3.2 | you're going to | 1.8 |
| I encourage you to pause | 1.8 | if you have a | 0.84 |
| if you were to + specialized verb (plot, graph, expand, calculate, multiply) | 1.2 | if you want to | 0.76 |
| see if you can | 0.9 | if you look at | 0.76 |
| and see if you | 0.8 | you don't have | 0.69 |
| and then you have | 0.8 | if you don't | 0.51 |
| I guess you could | 0.6 | you have to do | 0.48 |
| Or you could say | 0.6 | you can think of | 0.42 |
| You could view this | 0.5 | | |
| If you wanted to | 0.5 | | |
| you don't have to | 0.5 | | |
| if you look at | 0.5 | | |
| I guess you could say I guess you could guess you could say | 0.5 | | |
| if you have a | 0.4 | | |
| there you have it | 0.4 | | |

The results of the analysis of frequency of multiword expressions with the personal pronoun *you* reveal that in KAM and MITM the stance expression of personal intention/prediction (Biber, 2006) *you're going to* is the most used phrase. This is also in live

with previous research on expressions used in classroom teaching. Biber (2006) describes this category of expressions as conveying prediction that the instructor is making. In case of the situational use in my corpora the instructor's predictions refer specifically to the results of mathematical problem solving.

There are three phrases in KAM corpus that are used together in a variety of ways to engage with the audience: *I encourage you to pause, see if you can, and see if you can:*

So I encourage you to pause the video and see if you can do that. [doc#318KAM]

Based on this, **pause the video and see if you can figure out** what the inverse of g is. [doc#319KAM]

The phrase *pause the video*, actually, occurs with the rate of 1 per 10 thousand tokens. It is always directed towards *you,* the physical audience. I analyzed the phrase *I encourage you to* in the previous section. However, it is important to notice how it co-occurs with *see if you can*, as a direct encouragement of the audience to attempt to solve the problem on their own. The personal pronoun *you* is used twice *I encourage you* and *see if you can*, to acknowledge the presence of the listener and instruct them.

KAM corpus also features the lexical bundle or *you could say, I guess you could (say)* which serves the function of translating one representation of the concept into another one, modelling for the audience different representation of the same concept

It's 0.65. And if you want to write it as a percentage you essentially multiply this number by 100. Or another **way        you    could say** is you shift the decimal point over two spots to the right. So this is going to be equal to 65%. Now, there's [doc#34KAM]

5, taking that product. Then you have negative 2 times 1 times 3. Well that's negative 6. So we'll have negative 6. Or **you     could say** plus negative 6 there. And then you have 2 times 2 times 4. Well that's just 4 times 4, which is just 16. So we have [doc#726KAM]

Research in mathematics education shows that one of the key sources of error in students mathematical reasoning is "translation among numeric, graphical and algebraic representations of associated mathematical relations" (Bosse, Adu-Gyamfi, & Cheetham, 2011). The use of the phrase *you could say*, is directly related to Khan's attempt to show students how else the same concept could be represented. This phrase plays an important instructional role, as it signals to

the audience the possibility of representing the same result or variable using different notation. Using idiomatic phrasing, *you could say*, Khan invites the audience to verbalize another option. The same gesture towards the audience is realized in the phrase *you could view this*

> Let's do a few more of these examples. So then we have a 2 by 1, **you could view this is** as a 2 by 1 matrix or **you could view this** as a column vector. This is another 2 by 1 matrix, or a column vector. [doc#1101KAM]

In both cases, *you* is used not as a direct reference and idiomatically, but in phrases that invite the audience to see or name a mathematical concept in another way. Both phrases are used in order to recognize the need of the students to be able to see the same concept from multiple angles, under multiple names or graphical representations. Since Khan started his teaching experience with quite basic math instruction when his niece was very young, he might be continuing to use the strategies that are key for young learners even at more complex levels of math instruction.

The only other expression that functions similarly in MITM corpus is *you can think of*, in which the instructors invite the audience to understand a new concept in a particular way

> But, **you can think of** L, the way to think of it is as a black box, a function of what goes into the black box, well, if this were a function box, what would go it would be a number, and what would come out with the number [doc#14MITM]

> So the way this is usually said is that we'd like to reuse the existing blocks in the cache as much as possible. And this **you can think of** as temporal locality. When I access a particular block, I'm going to access it again fairly soon. That way it's actually useful to bring it into my cache, and then I use it many times. [doc#93KAM]

This is a more specialized use of an expression that serves to either to introduce a metaphor for a concept (i.g. L as a black box) or introduction of a new term (i.e. temporal locality). In both cases the instructor recognizes the gap in shared knowledge, either in terms of translational skills or knowledge of specialized concepts. The MITM instructors use this phrase to signal to the audience that a new concept is introduced or an old one is presented in a different manner.

The MITM most frequent multi-word expressions align with Biber (2006) results, with discourse organizing bundle: *if you look at, if you have a*; stance - attitudinal/desire : *if you want to*, obligation / directive expressions: *if you don't, you have to do.* These bundles are typical of

classroom teaching discourse, in which the students' attentions and actions are being managed by the instructor.

**Pronoun *we*: lexical description**

In order to analyze the use of personal pronoun *we*, I examined how it functions in the registers. In the sample of 200 randomly selected concordance lines, I examined the function of these pronouns. I used previous research on the uses of *we* as exclusive (general) or inclusive (direct reference to the audience) to code the lines for function of the pronoun.

Table 17. Function of personal pronoun *we* in both corpora. Analysis of 200 randomly selected concordance lines.

| Function | KAM | MITM |
|---|---|---|
| Mathematical Procedure Narration (general) | 190 (95%) | 171 (85%) |
| Direct reference to audience | 10 (5%) | 29 (15%) |

Not surprisingly, the personal pronoun *we* is used in narrating the mathematical procedures, so called "chalk talk," which is a central genre of mathematical instruction around the world (Artemeva and Fox, 2011). Writing on the board while narrating the process of problem solving provides a means of introducing students to "disciplinary thinking, practices, and procedures of 'doing mathematics'" (Artemeva and Fox, 2011, p. 356). The use of personal pronoun *we* for this type of narration involves the audience in the process of mathematical reasoning:

> 's l of x. And let's say this, right over here, this right over here, is the point (4.36, f of 4.36), and **the way that we're gonna approximate this value** is to figure out what this value is right over here. And what is this [doc#797KAM]

Even though Khan in the only active participant and the instructor in the tutorial, he uses *we* to involve the listeners in the organization of the mathematical problem solving steps. The personal pronoun *we* could easily be changed to *I* and the meaning of the sentence would not change fundamentally. Since it is Khan who is approximating the value and is "figuring out" the meaning of the other value, using *we* serves the purpose of including the audience. The same use of *we* is present in the MITM corpus, albeit to slightly smaller degree (85%).

(100 cm)^3. And with this unit notation **we** really do want to cube the centimeters and cube the 100. So **we** get pi / 2. **And here we get 10 ^ 6 cm^3.** And that comes out to pi / 2 * 1000 liters. Or, in other words, about 1600 liters. So I'd like to ask you first {doc#187MITM]

The MIT instructor uses the personal pronoun *we* in a sequence that provides a step-by-step guidance through a mathematical procedure. In the case of MIT lectures, the audience is physically present in front of the speaker, and presumably is taking notes as the lecturer writes on the board. This could explain, to some extent, a reference that is more directly pointing to the students as participants in the problem solving process. In case of the KAM, the audience which is removed from the speaker in space and time, is still involved in the process of solving the problem, as if they were physically sharing space and time with Khan.

Other uses of personal pronoun *we* in the MITM corpus serves as referent in discourse organizing function. The use of *we* relates to the instructor discussing the course organization, schedule, plans, etc. performing a metadiscoursive function. In the example below, the instructor refers to the focus of the lecture:

So, again welcome to 18.01. **We're** getting started today with what **we're calling** Unit One, a highly imaginative title. And it's differentiation. So, [doc#208MITM]

Other uses of *we* in the MITM corpus focus on organization of teaching segments, which can be understood as fragment of the lecture focusing on one concept:

for cache oblivious algorithms. All right, so that's sort of review of why this model is reasonable. LRU is good. So now **we're going to talk about** two basic problems-- searching for stuff in array, sorting an array in both of these models. [doc#48MITM]

Because MITM lectures are six to ten times longer than Khan Academy tutorials, the instructors need to organize the progression of the concepts they are introducing. The speakers use the phrases "we're going to talk about" or "what we'll do is just focus on" to signal a new section in teaching. The use of *we* points again to the inclusive nature of the instruction, in which the instructor brings the listeners to his proximal space (Kamio, 1999). Fortanet (2004), also found this frequent substitute of *we* for *I* in the MICASE corpus. There. She found such use of the pronouns *we* as attempting to involve the audience in the instructional activity.

The category of other functions in Khan Academy includes organizational and reflective commentary on where students should be at, given the sequence of the video instruction

So in the other videos, **we** looked at it in terms of breaking it down to its simplest parts, but **I think we have enough practice now** to be able to do a little bit more of it in our heads. So what is the largest number that divides [doc#462KAM]

Such use of *we* which points to the grouping of the speaker and audience, who are both engaged in the sequential process of teaching and learning the concepts through the series of videos. Similarly to the discourse organizing use of *we* in the MITM (doc#48) fragment, Khan refers to the past experiences of the audience using personal pronoun *we*, to signal inclusiveness and involvement in the teaching and learning process with the students. Unlike in lecture, where metadiscourse is essential to help students overcome the challenge of comprehending content delivered at relatively high speed (Thompson, 1994), the short tutorials do not require discourse organizing expressions since they are already focused on one concept at a time. In fact, the Khan Academy itself, its infrastructure serves as an organization map that guides the students through the learning process. That work of organizing discourse, introducing topics, elaborating, and transitioning to another topic from exercises back to the lecture is the responsibility of the instructor in the physical classroom.

The bi-gram results presented in Table 18 show how the pronoun *we* is used differently in both corpora. The functional analysis revealed that Khan has preference to use *we* in his mathematical narration of problems. This is why the bi-gram *we're* is used so frequently (7 occurrences per 10 thousand tokens) in the phrase *we're going to,* a key phrase in his virtual *chalk talk.*

Table 18. Bi-grams with *we* as the personal pronoun.

| KAM | Normed Frequency | MITM | Normed Frequency |
|---|---|---|---|
| we're | 37 | we're | 23 |
| we have/we've | 29 | we have | 19 |
| we can | 19 | we'll | 9 |
| we could | 14 | we can | 7 |
| we know | 8 | we do / we don't | 5.6 |

**Pronoun *we*: phraseological description**

Salman Khan's preference to use *we* as the narrator in his *chalk talk* impacts the diversity of multiword expressions with 40 occurrences per one million tokens, or 0.4 per 10 000 tokens (used often as a cutoff score in multi-word expression studies). Since he is the only speaker in corpus, the multiword units he relies on might be a part of his idiolect. Indeed, there are 5 expressions in the MITM corpus that include *we* and occur at least 0.4 times per 10 thousand tokens. There is 36 multiword expression in the KAM corpus, with some expression overlapping with each other to some extent.

Table 19. Multi-word expressions with *we* pronoun. Normed frequency.

| KAM | Normed frequency (per 10 thousand tokens) | MITM | Normed frequency (per 10 thousand tokens) |
|---|---|---|---|
| we're going to | 7.8 | We're going to | 2.13 |
| we just have to | 1.6 | we don't know | 0.68 |
| see if we can | 1.6 | we know how to | 0.60 |
| and then we have | 1.6 | we don't have | 0.60 |
| if we were to | 1.3 | we want to do | 0.41 |
| if we want to | 1.1 | | |
| we don't know | 1.0 | | |
| we could say that | 0.9 | | |
| we know that the | 0.9 | | |
| if we wanted to | 0.8 | | |
| so we could say | 0.7 | | |
| and then we can | 0.7 | | |
| we are going to | 0.7 | | |
| we don't have | 0.7 | | |
| we know that this | 0.6 | | |
| and now we can | 0.6 | | |
| so we know that | 0.6 | | |
| and we are done | 0.6 | | |
| or we could say | 0.5 | | |
| and we know that | 0.5 | | |
| so we have a | 0.5 | | |
| we want to do | 0.5 | | |

Table 19. continued

| KAM | Normed frequency (per 10 thousand tokens) | MITM | Normed frequency (per 10 thousand tokens) |
|---|---|---|---|
| and we want to | 0.5 | | |
| we could write this | 0.5 | | |
| we were able to | 0.5 | | |
| what we want to | 0.4 | | |
| we are left with | 0.4 | | |
| if we look at | 0.4 | | |
| so we could write | 0.4 | | |
| it looks like we | 0.4 | | |
| we can rewrite this | 0.4 | | |
| we want to figure | 0.4 | | |
| we want to figure out | 0.4 | | |
| we can figure out | 0.4 | | |
| and then we could | 0.4 | | |
| now we just have | 0.4 | | |

Overall, the use of the *we* in these expressions does not fall into direct reference to Khan and audience, but is used in a more generalized manner in the mathematical instruction register. There are a number of phrases that are very particular to the situational context of use. The expression, *see if we can*, engages the audience in a translational exercise that is expressed conditionally, without a predetermined result:

So this is equal to negative 8 plus or minus the square root of 60. All of that over negative 2. And let's **see if we can** simplify the radical expression here, the square root of 60. Let's see, 60 is equal to 2 times 30. 30 is equal to 2 times 15. And then 15 is 3 times 5. So we do have a perfect square here. [doc#393KAM]

The verb *simplify* is the most frequently used (0.2 per 10 thousand tokens) verb that follows the phrase *see if we can*. The other verbs: *figure out, solve, find, factor* suggest that Khan invites the audience into a narration of an attempt to solve the mathematical problem at hand. The use of conditional *if* leaves the act of solving as uncertain, and the process of problem solving as more open to mistakes and mishaps. In fact, other multi-word expression also include

the collective *we* into an uncertain trial-and-error problem situation: *we want to figure out, if we were to, if we wanted to*

> So we've got an equation here, it says five times x minus three is equal to four times x plus three. So, what we want to do is **we want to figure out** an x that satisfies this, so there's some number that if I take five, multiply it by that number, subtract three from it, that's going to be the same thing as if I take four times that number and add three to it. [doc#121KAM]

> Well, a, we're now saying we can represent that as some integer k times p. So **we can rewrite that** as some integer k times p. And so, let's see**, if we were to multiply this out.** So we get b squared times p-- and you probably see where this is going-- is equal to k squared times p squared. We can divide both sides by p, and we get b squared is equal to p times k squared [doc#115KAM]

> So here's another identity. Another way to write my cosine of 2a. We're discovering a lot of ways to write our cosine of 2a. Now **if we wanted to solve** for sine squared of 2a we could add it to both sides of the equation.[doc#1461KAM]

The conditional phrasing of the *chalk talk* can build a more engaging and open discourse. Rather than being authoritative, Khan opens up the process of problem solving to simulated discovery alongside the audience, rather than modelling the process to the audience from a more certain position. In fact, KAM corpus is characterized by triple the use of possibility modal *could* in comparison to MITM corpus. The modal *could* is used in a number of multi-word expressions with the pronoun *we*, which suggests that Khan wants to voice possibility for the collective *we*

Table 20. Modal verb frequency in both corpora. Normed frequency per 10 thousand tokens.

| Possibility/Ability/Permission Modals | KAM | MITM |
| --- | --- | --- |
| can | 39 | 36 |
| could | 39 | 12 |
| may | 0.1 | 2 |
| might | 10 | 4 |

Previous research on *we* in mathematical instruction (Pimm, 1987; Rowland, 1999) considered it as an imposing pronoun, which in *chalk talk* cuts off the individual opinions and voice of the student. However, in my study, Khan uses *we* with modal verbs of ability (*can*) and possibility (*could*), presenting the mathematical problem solving as full of possibilities and uncertain ends.

**Pronoun *our*: shared process solving**

I analyzed the use of the possessive personal pronoun *our*, since the frequency of use in both corpora were notable. Khan also uses the possessive pronoun *our*, more than twice as much as often as the speakers in the MIT corpus. A more detailed analysis of the collocations of *our,* reveals that *our* is often used in preceding nouns *change* and *function*, especially in the lessons in algebra:

And what's **our change** in X? [doc#211KAM]

So for g of x, if we were to write **our change** in y over **our change** in x [doc#223KAM]

The use of *our* instead of definite article the exemplifies how Khan builds rapport and shared process solving with the audience. While the same phrase *change in* is preceded in KAM with the possessive pronoun *our* twice as often as with the definite article *the*, in MITM corpus the same phrase is preceded only by definite article *the*. In the MITM corpus there are no instances of the use of *our* preceding *change*, but there are 30 instances of *our function*, 11 of which belong to one speaker – Dennis Auroux. The most common collocate of the pronoun our in the MITM corpus are *first* and *goal*, which point to the topic introduction function in the corpus. The use of the pronoun *our* in KAM corpus, however, might be caused by increased overall frequency of functions and change discussion in algebraic problems. More importantly, Khan involves the audience in the process of solving the mathematical problem by putting responsibility on both himself and his audience.

**Spatial Deixis Use for Engagement**

*Locative Adverbs: Here and There*

The analysis of frequency of demonstratives *here* and *there,* with manual correction for automatic tagging error, shows that in KAM corpus these adverbs are used more often than in the MITM corpus.



Figure 13. Demonstratives here and there in KAM and MITM corpora. Normed frequency.

This is not surprising, since KAM tutorials rely solely on writing on the virtual board. Bamford (2004) found that gestural *here* frequently appears in reference to visual aids. Since Khan's instruction relies almost exclusively on visual aids, the high frequency of *here* is not surprising. MITM instructors lectures potentially include other phases of the lecture that might not require blackboard, such as transition between major parts of the lecture, classroom organization, short exchanges between the audience, or discussion of materials students have read. KAM language is marked by more proximal feature *here*, which presents the object of discussion as being closer both to the speaker and the audience.

There are two elements to Khan's teaching: his voice and his writing on the board, he needs to rely heavily on deixis to help the audience follow his instruction. Bamford (2004) emphasizes how deixis fundamentally is tied to gesturing and the situational context of pointing in space. Khan's disembodied instruction makes him rely on deixis, rather than just body

movements in conjunction with mouse pointer to involve the audience and guide their attention. Khan points to elements on the virtual board constantly which is aided with the use of *here* and *there* alongside of the mouse pointer movement:

> Let me draw a number line **here** [doc#0 KAM]

> This is going to land us right over **here** on the number line [doc#1KAM]

> 2 times 4 is 8. It goes into the ones', or the constants' place. 2 times negative 6x is negative 12x. **And we'll put a plus there.** That was a plus 8. 2 times 9x squared is 18x squared, so we'll put that in the x squared place. [doc#450KAM]

In the above examples, Khan uses plural pronouns to mark movements together with the audience (i.e. land us right over here, we'll put a plus there). The actions are not only done together, but the location/destination of the action is shared across the speaker and the audience. Khan's description of procedures located on the virtual board make the audience involved and engaged in the process of solving the problems. The situational context of Khan imagining the student sitting next to him at the kitchen table, working out and talking through a problem is further linked to the deictic space (Hanks, 1992) both speaker and viewer are interacting with. The use of gestural *there* and *here* is described by Bamford (2005) as speaker-centered in which the speaker and the audience are in close physical proximity (p. 125).

Khan also uses *there* and *here* symbolically to refer to a part of the mathematical reasoning process or he is addressing the audience:

> So we could say that the set of polynomials, polynomials closed, closed... I won't even put it in quotes. Closed under subtraction. And I didn't prove it **here** , I just did one example where I subtracted two polynomials and I got another one, and there's clearly more rigorous proofs that you can do. [doc#422KAM]

*There* is also used symbolically to refer back not to a physical point on the board, but also a phase in the lesson (*i.e.* end of a tutorial) or a crucial point in a mathematical procedure. Quirk et al. (1985, p. 1453) noted that lecturers have a tendency to use *here* and *there* as adverbials of time. *There* becomes a point in a lecture, a time deixis (Fillmore, 1997):

is actually a fairly large number relative to the percentage that do take the drugs and test positive. **So I will leave you  there**. This is fascinating not just for this particular case, but you will see analysis like this all the time [doc#1312KAM]

I have to add the same amount to both sides or subtract the same amount again. Now, the reason why **I was careful there**       is I didn't just add 4 to the right hand side of the equation. Remember, the 4 is getting multiplied by 5. [doc#397KAM]

Such use of symbolic deixis as a referent to a symbolic object requires the listener to be more cognitively engaged than in the case of gesture use (Bamford, 2005).
The MITM corpus also includes instances of direct reference to a point on the blackboard:

And, say that because we are on the, let's see, sorry, we are going from the origin to (1,0). Well, we know we are on the x-axis. So, y **there** is actually just zero. And, the variable will be x from zero to one. [doc#78MITM]

### *Right Over Here*

The analysis of *here* in both corpora revealed that the adverb *here* is very often used in the KAM corpus in the phrase *right over here* (23 times per 10 thousand tokens) or simply *over here* (33 times per 10 tokens). This is not the case for MITM corpus, in which *over here* is occurring with 0.3 times per 10 thousand tokens, and the phrase *right over here* is used only 4 times in the whole corpus. The examples from KAM corpus further reveal that often its context included another demonstrative *this*, adding to the indexical function of the whole phrase:

This **right over here** is a positive six. [doc#4KAM]

This point **right over here** on the function would be x comma f of x. [doc#791KAM]

At this point **right over here** , the x value is b, and the y value, of course, [doc#865KAM]

The use of this phrase in KAM corpus stems from reliance on the use of a mouse cursor and voice to move the audience's attention. The movement of the cursor and the verbal cue "right over here" often involve the audience by moving their attention from one point of the board into

another. In fact, 621 transcripts include the phrase *right over here* at least once, while 443 transcripts featuring it 4 or more times.



Figure 14. Frequency of occurrence of *right over here.* Y axis number of transcripts *right over here* occurs. X axis is number of instances of *right over here*.

The multiple use of the phrase *right over here* per transcripts signals how many times audiences are asked to shift their gaze to a different point of the virtual board. These shifts require the audience to engage visually with the board and to pay attention to the instruction. Such high occurrence of the phrase *right over here* in the KAM corpus might additionally be explained by Khan's idiolect.

### Demonstrative Pronouns: This and That

The results for demonstrative pronouns *this* and *that* normed occurrence frequency, are similar to the results of *here* and *there* in that KAM corpus includes more proximal features. The use of proximal item *this* and *here* in KAM corpus suggests that Khan discusses objects as if they are physically or symbolically closer to him and the viewers. The distal demonstrative pronoun *that* is more frequent in the MITM corpus. The results exclude use of *this* and *that* in functions other than spatial or discourse deixis. Similarly to results on EAP teacher and student discourse in EAP (Friginal et al., 2017), the speakers in my corpora use proximal deixis more often than distal deixis.

Figure 15. Demonstrative pronouns. Normed frequency per 10 thousand tokens.

Deixis is key element of instruction, as it supports the instructors goal of shifting the students' attention from one object, problem, item, to another in order to guide the ongoing mathematical reasoning process.

> So then we're going to move 29 over to the right. **That's** the 29 part. Now **this** is a positive 29, and so how do we figure out what **this** is? **This** is 29 right **here** that we're adding. **This** is going to land us right over **here** on the number line. So how do we figure out what number **that** is? Well, once again, we can just visualize it. [doc#1339KAM]

Demonstrative pronouns allow the audience to track the progress in the mathematical procedure, with the use of gestural deixis (i.e. this is 29 right here). In the KAM corpus, the audience needs to be following the mouse pointer on the black background of the virtual board. The reliance on the proximal deixis can be interpreted as egocentric positioning (Kamio, 2001; Friginal et al., 2017), since the audience is focusing on the speaker's territory. The less substantial difference between *this* and *that* in the MITM corpus, suggests that instructors are present in a classroom, interacting with other people and objects. Research in the EAP physical space showed a more balanced used of spatial deixis (Friginal et al., 2017).

The reliance on proximal deixis in KAM can further contribute to Khan's vision of sitting next to a student at a table, solving the problem that is right in front of both of them. In that particular imagined context, the reliance on the proximal *this* brings the student closer to the speaker's perspective, and allows them to see problem solving from Khan's point of view. The frequent use of demonstratives supports the audience's tracing of spatial movement on the board conducted by the instructor. Since the mouse pointer can occassionally be difficult to distinguish from the background, the verbal cues – including changes in prosody which are beyond the scope of this analysis – can support the audiences reading and following of the movements of the most, especially its change in position. The positioning of the audience in the proximal space of Khan involves them in the virtual context by overlaying the audiences' visual attention with Khan's point of view, further engaging them through use of deixis for attention control.

### *Demonstrative Determiners: This / That + Noun*

Another key aspect of analysis is investigating objects of attention in both corpora. By looking at the patterns of *this/that+noun*, I analyzed what the attention objects across different mathematical subjects/ courses. The use of proximal *this*+noun is more frequent in the KAM corpus, while the distal that *that* + noun occurs somewhat more frequently in the MITM corpus.
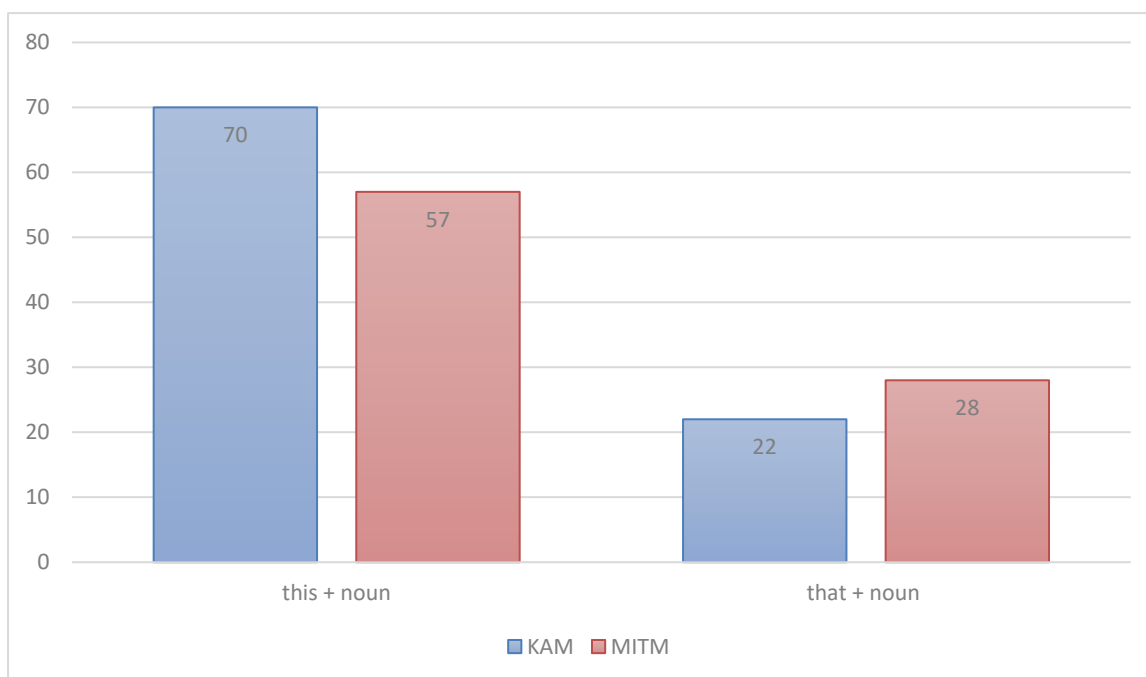


Figure 16. Demonstrative Determiners. Normed frequency per 10 thousand tokens.

There is a bigger difference between proximal and distal demonstrative adjectives in these corpora. Using situational characteristics can provide some context as to why these differences occur. The teaching space in the KAM is the virtual blackboard, which is set up as a space that is in front of the viewer and the speaker. Moreover, since the majority of the teaching happens through *chalk talk,* there are fewer opportunities to point to any part of the writing that is more distant in the imagined physical space of the board.

The analysis of *this*+noun patterns in both corpora reveal that the most frequent object of attention are two nouns: *thing* and *guy.* Previous research on the use of *thing* noun in spoken academic discourse shows it plays a number of key functions: refers to actions or ideas, important information (i.e. *the thing is),* to identify points of information (Biber, 2006, p. 56). However, five speakers in the MITM account for the frequent use of the bigram *this guy*, a surprising finding.

Table 21. *This* + noun combinations. Normed frequency per 10 thousand tokens.

| KAM | Normed frequency per 10 thousand tokens | MITM | Normed frequency per 10 thousand tokens |
|---|---|---|---|
| this thing | 3.8 | this guy | 2.2 |
| this equation | 3.7 | this point | 2 |
| this video | 3.4 | this thing | 2 |
| this point | 2.6 | this matrix | 1.1 |
| this expression | 2.2 | this formula | 1.1 |
| this term | 1.5 | this equation | 1 |
| this function | 1.5 | this function | 1 |
| this angle | 1.5 | this problem | 1 |
| this part | 1.4 | | |
| this line | 1.4 | | |
| this number | 1.2 | | |
| this side | 1.2 | | |
| this distance | 1.1 | | |
| this business | 1 | | |

The frequency of use of *thing* as a noun in KAM corpus is in line with research on spoken academic discourse. *This thing* in the mathematical register can function as a general description of parts of equations or group of variables, etc. Using *this thing*, rather than a more precise

phrase, with the support of mouse movement allows for an easier reference to an object on the board

> And we want to solve for n. Well the easiest way to solve for n is maybe multiply both-- **this thing** on the left is equal **to this thing** on the right. So we can multiply them both by the same thing. [doc#46KAM]

The use of *this thing* focuses the audience's attention on objects on the virtual board, without the need of using technical vocabulary that requires more cognitive processing on the part of the audience. The use of quite vague phrase *this thing* by Khan (and to lesser extent MITM) allows for more flexibility and ease in directing audience's attention.

However, in the MITM corpus, the attention of the students is being shifted to *this guy*, which refers to an unknown or a part of an equation

> Can you tell me an eigenvector here for **this guy**? x equal -- what is -- actually, can you tell me one vector that which is lambda x, and I have a three x, [doc#84MITM]

> There will be a square root and some squares and some stuff. What is the normal vector? Well, you take this vector and you scale it down to unit length. Just to emphasize it, **this guy** here is not n and this guy here is not dS. Each of them is more complicated than that, but the combination somehow simplifies nicely. [doc#112MITM]

In both cases, the use of nouns that are not precise and can stand in for a variety of mathematical concepts allow the speaker to exert less effort in terms of generating more precise labels.

The frequency of other *this* + noun combinations shows that the speakers do refer to specific elements of mathematical procedures: equation, point, function, formula, matrix, etc.. There are more pairings of *this* + noun in the KAM corpus, which suggests that he refers to objects on the board more often than the speakers in MITM corpus. Since his level of instruction is geared to high school students or students early in their college career, he draws the audience's attention to more fundamental elements of mathematical reasoning: line, number, side, angle distance. In the MITM corpus, the deictics point to more complex concepts, without referring to the basics of math.

One item that frequently appears in KAM corpus *this video*, has no presence in the MITM corpus. The use of *this video* is clearly the function of the situational characteristics – the

video tutorial format – of the KAM corpus. *This video* functions as a discursive deictic, most frequently used in the phrases *what is want to do in this video is* or *I encourage you to pause this video*. *This video* serves as a unit of teaching/learning, and the phrase *this video* draws the students' attention to the video tutorial and its goals. The use of *this video* allows the students to follow a sequence of videos, since Khan often juxtaposes what the goal in *this video* is with what will happen in the *next video* (0.7 occurrences per 10 thousand) or was taught in the *last video* (1 occurrence per 10 thousand).

In contrast with the proximal use of demonstrative determiners, the distal use is less frequent. In both corpora, the bigram *that point* is used most frequently. Both corpora include deixis pointing to unknown x*: that x* or a function *that f*, which are written on the board.

Table 22. Demonstrative determiners: distal deixis. Normed frequency per 10 thousand tokens.

| KAM | Normed frequency | MITM | Normed frequency |
|---|---|---|---|
| that point | 2.0 | that point | 1.17 |
| that x | 1.6 | that matrix | 1.13 |
| that right | 1.0 | that way | 0.97 |
| that f | 0.6 | that f | 0.48 |
| that way | 0.6 | that equation | 0.47 |
| that number | 0.6 | that x | 0.42 |
| that line | 0.5 | that formula | 0.42 |
| that interval | 0.4 | that vector | 0.41 |
| that distance | 0.4 | | |

The most frequently used combination *that point* function both as a spatial deixis and as a discourse deixis. Since mathematical reasoning, especially in graphic algebraic functions, *point* as a noun plays a central role in that discourse:

Here is delta x, and here is the height. It's because **that point is** 0, 0 [doc#157MITM]

We can figure it out numerically. I'll in a second draw it graphically. So what's our change in y? Our change in y is literally how much did our y values change going **from this point to that point** ? [doc#199KAM]

In KAM corpus especially, both *this point* and *that point* are used more frequently, suggesting that there is more discussion of functions in general. The example from document #199 shows especially how Khan uses proximal and distal signaling when he discusses differentiate between the two points on the virtual board. The MITM corpus includes the use of *that point* as a discourse marker more frequently than KAM corpus.

> Here it was a plane. Okay. I'm going to **stop at that point.** That's the central idea of -- the great example of how to create a subspace from a matrix. [doc#1MITM]

Actually, in a random sample of 100 concordance lines, there were no instances of *that point* in a discourse deixis function. Since the KAM tutorials are short and focused, there might not be a need to use *that point,* in reference to the point in a tutorial because their organization and duration is quite simple in comparison to lectures.

## Interactional features

### Hypothetical reported speech in Mathematics Instruction

The results of pattern searches revealed that Khan uses direct hypothetical reported speech almost ten times more often than speakers in the MITM corpus. Table 23 presents results for the DHRS pattern normed occurrence frequency in both corpora. Since my methods rely on examining features that have been previously described in the literature, there is a possibility that not all of the instances of direct hypothetical reported speech were found in the corpus. The patterns were to some extent developed using an iterative process of the search terms appending (i.e. I added *Sal,* and *wait* as additional lemma signaling potential presence of DRHS). There is an overall tendency of Khan to rely on DHRS in his teaching.

Table 23. Hypothetical reported speech - occurrence in KAM and MITM corpora.

| | KAM | | MITM | |
|---|---|---|---|---|
| | **Occurrences** | **Frequency (per 10 thousand tokens)** | **Occurrences** | **Frequency (Per 10 thousand tokens)** |
| Direct Hypothetical Reported Speech (total) | 657 | 3.7 | 72 | 0.38 |
| *you* | 491 | 2.8 | 72 | 0.38 |
| *we* | 162 | 0.97 | 0 | 0 |

In previous research, the frequency of hypothetical reported speech in business meetings (Koester & Handford, 2018) were 2.7 occurrences per 10 thousand tokens. The difference, then, between MITM and KAM corpus becomes even more stark, since the use of DRHS is even higher in KAM corpus than in a business meeting corpus.

The analysis of the patterns reveals that the personal pronoun *you* was used more frequently to launch DHRS than *we*. The communicative verb *say* is used in combination with modal verbs to launch DHRS. The frequency of *say* as the main verb for introducing DHRS is not surprising, since it is also the most frequently used verb for reported speech in general (Clift & Holt, 2007). The use of modals: *might*, *could*, and *would* signals Khan's assumptions or predictions of what could or might be said in this situation by the students. These modal verbs accompany the verb *say* and pronouns *you* and *we* in signaling hypothetical reported speech in KAM corpus. In terms of its function, DHRS was used to recreate responses or questions of the imagined students.

Table 24. Patterns used to signal DHRS in KAM and MITM corpus.

| KAM | Normed frequency (per 10 thousand tokens) | MITM | Normed frequency (per 10 thousand tokens) |
|---|---|---|---|
| you might say | 0.63 | you say | 0.11 |
| you say | | you might say | 0.04 |
| | 0.55 | | |
| you could say | 0.45 | you could say | 0.02 |
| we could say | 0.17 | you can say | 0.02 |
| we say | 0.14 | | |
| you're saying | 0.13 | | |
| you'd say | 0.10 | | |
| we said | 0.10 | | |
| you might be saying | 0.10 | | |
| you're like | 0.08 | | |

Each structure was examined functionally in its discursive context. The grammatical structure of DHRS does not correlate with any particular teaching function - knowledge-factual or procedural/application – but can be used to express either of these functions. DHRS in general is most often used to enact the imagined student's procedural/application knowledge, in what could be called a form of imagined students' *chalk talk*:

> The absolute value of negative 3 is essentially saying, how far are you away from 0**? How far is negative 3 from 0? And you say, well, it's 1, 2, 3 away from 0**. So you'd say that the absolute value of negative 3 is equal to positive 3. Now that's really [doc#9KAM]

In this case the imagined student is describing the steps of the procedure of arriving at the absolute value of negative three. Khan not only gives the imaginary student a voice, but also uses it to answer his own (instructor's) procedural question. This allows him to ask and answer his own question, without excluding the audience from the process. Instead, he connects with the viewers by trying to guess and enact their thought processes and their possible correct responses. This use of direct hypothetical reported speech is similar to results of researchers description of problem-solution pattern in which the speaker creates a possible problem and provides a solution to that problem (Hoey, 1983; 1994; Koester & Handford, 2018).

Another use of the DHRS in KAM corpus is in enacting a student question as the narration of procedure progresses, highlighting moments that might pose challenges to the listeners' learning process:

> 7 plus 2 times 5. And then, 2 times 5. And then close brackets. Minus 25. Now, this thing-- we want to do multiplication. **you could say, hey , wait.** I still have a parentheses here. Why don't I do that first? But when you just evaluate what's inside of this parentheses, you just get a negative 7. It doesn't really change anything. So we can just leave this [doc#22KAM]

In this case, Khan commented on the order of steps to solve the equation. He is imagining that the student might be wanting to do the part in parenthesis first and adds a commentary that would enhance the students' procedural knowledge. Such anticipation of the audience's need shows to the listeners that he recognizes their presence and engages with their thought process. This example also includes the use of the attention signal *wait,* which is often used to mark such a point where the audience might have a question:

> triangle. And so we could use the Pythagorean theorem in order to figure out what is this distance right over here. **And you might say, well, wait. How do we figure that out?** Well, we already know that this length of the triangle is 36 feet. And we know that this base right over here is 1/2 of the width of the goal. [doc#74KAM]

This example shows how Khan pauses his procedural narration to ask a question about explanation in the student's voice. The question posed to Khan serves to prompt himself to provide a more detailed analysis of the solution, drawing on the previous knowledge explained earlier in the video (i.e. we know that…). The same question pattern is used in the MITM corpus

> becomes the square root of x squared over x squared, which is one, plus y squared over x squared. It's homogeneous. Now, you might say , hey, this looks like you had to be rather clever to figure out if an equation is homogeneous. Is there some other way? [doc#83MITM]

Khan also uses DHRS to engage with students' conceptual knowledge, by posing questions about definition of new concepts in the student voice and answering them

are equations of this form over here. And in a traditional algebra curriculum, they're called linear equations. And you might be saying, well , OK, this is an equation. I see that this is equal to that. But what's so linear about them? What makes them look like a line? And to realize why they are linear, you have to make this jump that [doc#83KAM]

In this case, Khan imagines what the student could be saying when new term - *linear equation* - is presented. The shift in the subject of speaking from Khan to an imaginary student is marked by transition to the pronoun *I*, after the direct reported speech signal *you might be saying*. Taking the personal pronoun *I* again shows how Khan needs to engage with his imagined audience and create a dialogue that otherwise would not be present in the online environment. The introduction of the audience's voice in form of the question is very brief and prompts the instructor to elaborate on conceptual knowledge. Both instructors put a key question in the voice of the student and proceed to answer it, shifting the focus on the students addressing them with the pronoun *you*.

The same strategy of using DHRS to ask and answer an imagined student question is present in the MITM corpus:

of these rows. The columns in the answer are coming as combinations of those columns. And so that's three ways. **Now you can say, okay, what's the fourth way?** The fourth way -- so that's -- now we've got, like, the regular way, the column way, the row [doc#133MITM]

Despite the situational characteristics – the students being physically present to ask questions – the instructor still uses a hypothetical voice to move the discourse forward and to imagine the question in audience's line of thinking. Using questions in teaching is the fundamental way to engage and interact with students (Bamford, 2005; Goody, 1978) and stimulate teaching / learning process (Crawford Camiciotolli, 2004; Long & Sato, 1983). However, asking questions to live audience and waiting for a response can take time. In case of MITM lectures the goal is to transmit knowledge, since interaction happens in dedicated recitation sessions. Thus, the use of the questions can still function, similarly to KAM corpus, to create an interactional exchange without the actual participation of the audience.

The use of DRHS in two corpora differs in the preference of personal pronouns used to launch the patterns. The instructors in MITM did not use *we* to initiate hypothetical reported

speech, which is the case in KAM. The use of *we* can be explained by Khan's tendency to use that personal pronoun more often in general than other pronouns, especially in relation to narrating mathematical procedures.

> , 4 times 9 is 36. But let's just think of the cards as being 1 through 36, and we're going to pick nine of them. **So at first we'll say, well look, I have** nine slots in my hand, right? 1, 2, 3, 4, 5, 6, 7, 8, 9. Right? I'm going to pick nine cards for my hand. And [doc#1325KAM]

> interval. And this would be the probability of success in that smaller interval. And in the last video we tried it out. **we said , oh, well , what if we** make this interval a minute and this is the probability of success per minute? We'd have maybe a [doc#1365KAM]

In these examples, direct DHRS is initiated with the use of *we*, but in one case Khan still transitions to first person narration in the audience's voice. The other example presents an imagined collective voice (the students and the instructor), in which Khan refers back to a point made previously in the reasoning process (anaphoric reference). The pronoun *we* is also used in the *be+like* structure

> "s probably bothering you, because it's bothering me, is these x's that we have in the denominators right over here. **We' re like, well , how do we deal with that?** Well, whenever we see an x in the denominator, the temptation is to multiply it by x. [doc#143KAM]

> There are 21 instances of *personal pronoun + be + like* in the KAM corpus, but only 1 instance in the MITM corpus

> So this is called a satisfying assignment. Satisfying just means make true. **And you're like , no, I don't believe you.** And your friend says no, no, no, really, it's true. And here's how I can prove it. [doc#67MITM]

> Since the non-traditional quotative *be + like* is a feature of younger generation of speakers (Macaulay, 2002), it might not be used for the older generation of instructors at MITM. In this particular case, the instructor – Dr. Devadas – takes on a voice of a students' friend

(presumably young) to develop a scenario, thus allowing for the imagined conversation between two young students.

> And so what you do is, you take each of the equation and you like, when you graph it, you like to put it in kind of the y-intercept form or slope intercept form and so you do that, **so you say, well let me solve both of these for b** so if you want to solve this first equation for b you just subtract 2a from both sides if you subtract 2a from both sides of this first equation you get b is = to -2a + 3. [doc#242KAM]

In this case, Khan narrates mathematical calculation, guiding the listener as if they are doing the calculation themselves. This is a particular case of *chalk talk* (Artemeva & Fox, 2011), which is a central pedagogical genre of undergraduate math courses, in which writing out mathematical calculations is narrated by the instructor. Transition of the authority enacting the procedures from Khan to the imagined student creates a form of hypothetical engagement between Khan and his audience. The audience gains a voice, through direct hypothetical speech, and can feel engaged in the learning process.

This section of results presented the function of DRHS in both corpora. The overwhelming presence of this discourse pattern in KAM corpus suggests that it is used to create engaging lessons by giving voice to the imagined audience. Khan makes the experience of learning interactive by predicting how a student might struggle with the material, what questions they might ask or how the steps they might follow while solving a problem.

**Response elicitors**

The results show that the most frequently occurring response elicitor is *right/alright*? which is present in both corpora. In MITM corpus speakers also use *ok/okay?*, *yes / yeah*?, and *all right?/alright?* as response elicitors, which are not present in the KAM corpus.

Table 25. Response elicitors normed frequency per 10 thousand tokens in both corpora.

| Response elicitors | Khan Academy | | MIT | |
|---|---|---|---|---|
| | Raw count | Frequency (Per 10 thousand tokens) | Raw count | Frequency (Per 10 thousand tokens) |
| Right? | 322 | 1.8 | 890 | 5.4 |
| OK/okay? | - | - | 487 | 2.8 |
| Yes/yeah? | - | - | 192 | 1.1 |

The presence of *right?* in the KAM provides evidence that despite not having physical audience in front of him, Khan still uses discourse marker *right?,* which is typical for conversation register (Biber et al. 1999, p. 1086). The occurrence of *right*? is not 2.5 less frequent than in MITM corpus. A more detailed look at the MITM corpus distribution of the response elicitor *right?* reveals that there are speaker-dependent frequency differences. While Demaine and Devadas taught one course together, their frequency of using response elicitors differs greatly.

Table 26. Top five speakers in MITM by response elicitor use (normed frequency by 10 thousand tokens in their own subcorpora).

| | Raw occurrence | Normed frequency per 10 thousand tokens in the author sub-corpus |
|---|---|---|
| Strang | 242 | 9.11 |
| Devadas | 237 | 17.2 |
| Xia | 57 | 49.7 |
| Leiserson | 57 | 5 |
| Demaine | 52 | 2.08 |

Comparison of Khan's results to speakers in MITM corpus shows that his frequency of response elicitor *right?* use is slightly less than some of the speakers in MITM corpus. Actually, the speaker with the largest representation in the MITM corpus – Eric Demaine (280K tokens) – uses *right?* with just slight higher frequency to Khan. Biber et al. (1999, p. 1089) comment that *right?* typically acts as a prompt for a verbal response from the audience, which in the face-to-face modality is possible immediately. Actually, there are just five instance of audience reactions following the response elicitor *right?* in the MITM corpus. However, the presence of this conversational discourse marker provides evidence that Khan is signaling to the audience that he

recognizes their presence and is seeking verbal responses, irrespective of the fact that he is not able to hear it.

Functional analysis reveals that in the KAM corpus, in which each video is focused on a distinct mathematical concepts or procedure/exercise, response elicitors are often used to check comprehension of the imagined audience in response to what is presented:

> all the way down to n minus k plus 1, **what's it going to look like? It's going to be a polynomial, right?** We're multiplying a bunch of-- well really, we're multiplying a bunch of binomials and we're doing it k times. [doc#1365KAM]

This is an example of the use of response elicitor, right after a naming a concept – polynomial – Khan seeks confirmation from the listeners. After this comprehension check, he continues to explain why the preceding part can be called a polynomial. Khan also seeks agreement of the listeners after providing an explanation of algebraic representation:

> algebra. So that equals b sine-- b squared sine squared of theta**. Sine squared of theta just means sine of theta squared, right?** Plus, and we just foiled this out, although I don't like using foil. I just multiply it out. c squared minus 2cb cosine [doc#1424KAM]

> mean squared over 2 sigma squared. This is the standard deviation. **Standard deviation squared is just the variance, right?** And just so you know how to use this-- you're like, oh, wow, there's so many Greek letters here. What do I do? This tells you ,[doc#1232KAM]

Since there are multiple ways to represent the same concept in mathematics, it is important for the students to be able to recognize and transform the notation to fit the purpose of the exercises (Bosse, Adu-Gyamfi, & Cheetham, 2011). Khan shows not only how concepts in statistics are equivalent, but also represents different notations:

> the inequality. And now, what is this? **You know, this doesn't look like a rational expression. But I can rewrite minus 2, right?** Minus 2 is the same thing as minus 2 times x plus 4 over x plus 4. This is minus 2 times 1. This is the same thing. And this is the [doc#682KAM]

is going to shrink by 25$. **So if something shrinks by 25%, that means it's just going to be 0.75 or 75% of what it was before,   right?** 1 minus 25%. 0.75 times $125. So let's work that out here. $125 times 0.75. And just in case you're confused, I don't want [doc#161KAM]

In example #161, Khan is showing two different types of notation 0.75/75% of the simple concept of fractions. In the process of teaching mathematics, the *chalk talk* narration, he involves the listeners in the crucial moments of translation between quantities or concepts, drawing attention of the listeners in that moment. However, more often than the use of *right?* for representing and translating between concepts, Khan uses *right?* in simple *chalk talk*, narrating problem solving in his tutorials, for example:

approximately equal to the plus or minus. We can take the square root of this. Plus or minus the square root of 4/9 is 2/3**,   right?** Square root of 4 over square root of 9, times x. So, these are the asymptotes. There's two lines here. There's y is equal to [doc#1063KAM]

for t, let's subtract 10 from both sides. You get x minus 10 is equal to 5t. Divide both sides by 5. You get x over 5 minus 2, **right?** 10 divided by 5 is equal to t. And now, we can take this and substitute it. This is t, right? So we can take this and [doc#687KAM]

As Khan narrates solving problems step by step, he adds the response elicitors in between the steps in the procedure. He uses *right?* in particular after getting to a partial result that will be the stepping stone for the next problem solving fragment. In the case of KAM register, the response elicitor *right?* serves the purpose of building and signaling shared knowledge between Khan and his students. It also serves as attention device, since not only is *right?* a discourse marker, but it also functions as a question adding an interactive dimension to its function. Thompson (1998) suggested that the use of *right?* in monologic academic speech functions to seek agreement and symbolic negotiation. Even though *right?* might be seen as persuading the listener to agree and does not open up much space for questions, Khan does not just move directly to the next fragment, but provides an explanation of why the listener should agree:

Well, that's just positive 2**, right?** A negative divided by a negative is a positive. [doc#141KAM]

It makes sense because x can't be negative 2, **right? Because** negative 2 has an absolute value less than 3.[ doc#351KAM]

Well, we're going to add 2 to this one, **right**? **So** it's going to be x plus 4 [doc#151KAM]

In fact, the response elicitor *right?* is followed by linking adverbial *so* (45), which is a conversational feature which helps to make logical connections between points that are being made (Biber et al., 1999, P. 866). This response elicitor is also followed twenty times by a *because*-clause, which provides an explanation in case the audience did not understand, or have that knowledge. These patterns suggest that Khan is aware of the presence and needs of his audience, just as the speakers in physical lecture are aware. He is using conversation grammatical markers in his monologue to involve and interact with his online audience, even though they are not physically present to respond to his questions.

Similar patterns and uses of *right?* can be found in MITM corpus. The MIT instructors are using *right?* to draw attention to conceptual information. For example, Strang presents a description of square matrix qualities:

So I've got a square matrix A. And it may or may not have an inverse, **right**? Not all matrices have inverses. In fact, that's the most important question you can ask about the matrix, is if it's -- if you know it's square, is it invertible or not? If it is invertible, then there is some other matrix, shall I call it A inverse? And what's the -- if A inverse exists -- there's a big "if" here. If this matrix exists, and it'll be really central to figure out when does it exist? [doc#67MITM]

Just like Khan presents a key concept (e.g. polynomials) and follows it with response elicitors and explanation, so do the speakers in MIT. In this instance of teaching about matrices, Strang seeks confirmation of shared knowledge with the audience on its characteristics (i.e. "may or may not have an inverse"), but follows it up with an explanation for those who do not share the information. In the other examples, response elicitors are positioned after the correct answer, and signal the explanation.

What is all R of P? It's R**, right?** Because the weighting doesn't matter. [doc#485MITM]

What's log of 1,000? Ten, approximately, **right**? Log base two of 1,000 is about ten, so that's 10^2. [doc#435MITM]

There are also very few instances of the audience verbally responding to the instructors after hearing *right?* in the MITM corpus.

The MIT instructors also use *okay?* as a response elicitor in their instruction. The analysis of concordance lines reveals that one of the function of *okay?* is similar to *right?,* as it functions as an attention device in a math procedure, followed by a new clause starting with linking adverbial *so*

> the answer, OK, which is the same as the serial code, not surprisingly. That's what we want. **And that's equal to order n, OK? So**, the parallelism is then p bar equals T_1 of n over T infinity of n is equal to theta of log n. Is that a lot of parallelism? [doc#86MITM]

> things. So, that means that we've written was one over 4a times plus something squared plus something else squared, **OK?** So, these guys have the same sign, and that means that this term here will always be greater than or equal to zero. And that [doc#163MITM]

In the process of narrating a mathematical problem solving, the speakers check with the audience if they are following. Othman (2010)'s analysis of *okay?* use revealed that it can function as a progress check in a variety of other academic lectures. The instructors will often use *okay?* following an explanation of a step in a calculation procedure. Using a question in, otherwise monologic discourse allows the speakers to control the attention of the students in a crucial moment that might impact how they follow the rest of the procedure. Actually, in one case of the speaker follows the audience check *okay?* with a more precise question

> be Pythagoras, but I don't have a right angle. So, I have a third term which is twice length A, length B, cosine theta, **OK? Has everyone seen this formula sometime? I** hear some yeah's. I hear some no's. Well, it's a fact about, I mean, you [doc#106MITM]

This use of the response elicitor shows how it is used as an audience-oriented question that helps to gauge audience's knowledge. The use of *okay?* also comes after a crucial piece of information that needs to be understood in order for students to follow along and proceed to the next step:

> that has the maximum yij, that is going to be my upper tangent. **Because only for that will I have no points ahead of that, OK?** **So** yij is upper tangent. This is going to be maximum. And I'm not going to write this down, but it makes sense that the lower [doc#142 MITM]

, OK, and so we have, actually, completed the proof for this part. Now, well, that's only for a vertically simple region, **OK?**   So, if D is not vertically simple, what do we do? Well, we cut it into vertically simple pieces. OK so, concretely, I [doc#102 MITM]

There are also 4 instances of my random sample of the audience reacting almost directly after hearing *okay?*. In fact, there are about 4 instances in the whole corpus or verbal audience reaction to *OK?*

a merge process that's even faster. And we obviously tried to cook up a theta one merge process. **But that didn't work out, OK?**     <STUDENT>: But are there algorithms that [INAUDIBLE] ? <PROFESSOR>: First-- if you assume certain things about the input, [doc#122MITM]

As Thompson (1998) theorizes, the use of these response elicitors might be only symbolic, as there in very little expectation of verbal response. Rather, these discourse markers be serving the purpose of gauging a non-verbal response, in terms of gaze. In Othman's (2012) study, the use of *okay?* coincides with lecturer's gaze moving towards the audience. The quick question *okay?* might be used as an opportunity to read student's faces to see if the point was explained adequately so the lecturer can move on with the content.

The use of *yes?* in MITM corpus creates the conditions for the most interaction between lecturers and the audience. 20% of the response elicitor *yes?* in MITM corpus if followed by an audience turn. This suggests that this discourse marker functions differently from *right?* and *okay?* The questions *yes?* is used in a response to a spontaneous audience question, as is the

plus a multiple of the return per risk of the market portfolio. This term here is called the price of market risk.   **Yes?**   <AUDIENCE >: Is that the same as the Sharpe ratio? [#doc30MITM]

Frequency analysis of these turns reveals that there are two courses *Application of Math in Finance* and *Multivariate Calculus* that are the most interactive in terms of audience turns after *Yes?* questions. In the process of qualitative analysis of the concordance lines for other cases revealed that these questions were still used for direct interaction with the audience. Since the transcription conventions were not uniform, some of the turns of the audience were not

marked in the transcripts, but video analysis revealed that the instructor was interacting with the students.

The analysis of response elicitors *right? okay?* and *yes?* revealed that they serve the purpose of involving audience in mathematical reasoning instances. Since they function as questions, they engage the audience, even symbolically in the mathematical reasoning process. In fact, they often mark an important point in the narration: they direct the listener's attention to a crucial component of the procedure that is essential for understanding the next steps in the calculation.

Despite the difference in modalities – online or face-to-face – response elicitor *right*? is present in both corpora. It is not as frequent in KAM corpus, but plays an important function in symbolically asking students for agreement, even if they are not physically present in the same time and space to answer. It is not surprising that instructor speech represented in MITM course includes a more diverse repertoire of response elicitors, especially that they have a function that relates to assessing non-verbal behavior of the audience or guiding verbal interactions. The discourse marker *right?*, the most frequent one in both corpora, serves the purpose of involving the audience through directing their attention.

# CHAPTER 6: CONCLUSION

## Chapter Overview

This dissertation study is an exploratory research into the presence and function of engagement linguistic features in two corpora of video instructional materials sourced from Khan Academy and MIT OpenCourseWare. Limited research has examined the linguistic features of video tutorials made for online teaching including for popular platforms like Khan Academy. The current research provides a valuable first step into examining the language students encounter when learning online in non-institutional settings. The key findings are summarized below.

## Key Findings

In response to research question 3 and 4 about the difference in the use of involvement and interactive linguistic features, there are differences in terms of frequency and function of these features in the corpora. Although, KAM corpus represents a register aimed at virtual audiences (removed in time, space, and ability to respond), still exhibits frequent use of involvement and interactive features. These differences are summarized next.

### Speaking in one voice: using we in KAM corpus

KAM corpus represents a very large sample of Salman Khan's language used for teaching mathematics. The analysis of pronoun frequency revealed that he prefers to use the personal pronoun *we* in his *chalk talk* narration, involving the audience in solving the mathematical problems. There is general discrepancy in previous literature on the function of *we* in instructional language, interpreting it either as a signal of joint work (Flowerdew & Miller, 1997) or stifling and co-opting student voices into the constructed disciplinary authority of mathematicians (Pimm, 1987; Rowland, 1999). The use of *we* may also reduce the affective distance between Khan and his audience, as he aims to talk *with* his audience rather than *at* them (Khan, 2012). In structuring his *chalk talk,* the use of *we* takes up his promise of structuring tutoring in a manner that makes the student feel like they are "working out problems together" (Khan, 2012, p. 34). The insistence on *we* as structuring the work as a shared activity is intended

to be inclusive and empowering, rather than limiting and condescending. In contrast to MITM corpus, in which instructors use more references to self through *I,* Khan blends his authority with the audience's voice by using *we* as a reference for *I* often. This shared ownership of the work also shows through the use of possessive pronoun *our* in reference to objects on the board (i.e. *our function, our change).* On the other hand, MITM instructors with credentials and hired by one of the top STEM universities in the world, can position themselves as a voice of authority, modeling instruction and talking about their own thinking processes.

More research needs to examine how students perceive such language of engagement and inclusion: as imposing and restricting, or as inclusive or empowering. Another dimension of the engagement aspect can be brought by analyzing prosodic features, rather than only focusing on lexical aspects of engaging students in learning. How instructors sound when they try to engage students in learning is as important as which words they use.

**Engaging by discovering the unknown: let's see if we can…**

Khan's use of involvement and interactive features ask the audience to participate in imagined scenarios. There are two aspects of activating hypothetical scenarios in the KAM corpus to engage students: conditional phrases and direct hypothetical reported speech.

Khan relies on conditional phrases: *if I were to + take, draw, say, ask…, (let's/let me) see if I can, see if you can, if you wanted to, see if we can, if we were to* to model uncertainty and engage the student in the hypothetical situations. Khan signals to the audience that he processes each problem set as he goes through it, without knowing the answer. The use of the phrase *let me see if I can* models for students not only a method, but an attitude towards problems. Just as the audience is not sure if they can be successful in the problem solving, so is Khan as an expert uncertain about his outcomes. Revealing his uncertainty shows his involvement in the mathematical reasoning process. More often than modelling his own uncertainty, his language builds engagement through inviting the students to participate in the attempt to solve the mathematical problem with the use of such multi-word expression as *see if we can*.

The use of directed hypothetical reported speech in KAM corpus is ten times as frequent as in the MITM corpus. One interpretation of such persistent use of this feature is enacting an imagined interaction between the tutor (Khan) and his students, since there is no physical audience present in the time and space of his recording. The use of DHRS combines a direct

reference to the audience with the phrase *you might say*, signaling to them Khan's recognition of their presence.

The presence of this feature in KAM corpus is one of the most interesting findings revealed by this study. The use of DHRS for modelling to audiences the possible questions they could ask at various stages of a mathematical problem solving procedure engages the audience by enacting their voice and suggesting ways of thinking about the problem. Khan anticipates what a student might ask or struggle with and voices that concern through direct hypothetical reported speech. More research on that particular feature should explore in detail its particular grammatical structures including (a) the use of statement and questions in the imagined student voice, (b) the function of the DRHS for conceptual and procedural knowledge building, and (c) the effects of using DRHS on student learning.

**Working at the kitchen table: gesturing to draw attention in KAM**

Language used in KAM corpus is marked with a frequent use of spatial deixis, locative adverbs and demonstrative pronouns. Situational characteristics could explain the difference between the use of these devices in KAM and MITM corpora. While the lectures in MITM include classroom management language and use the opportunity for other types of communication, the Khan videos are restricted to explain one concept at a time. The infrastructure of the Khan Academy platform organizes the learning experience as the videos are embedded in learning paths that guide the audience through short readings, the videos, and exercise space. MIT lecturers need to orchestrate the lecture and its various components on their own (or in teams), while Khan Academy's audience controls their own learning experience and their learning trajectory.

The videos may be more engaging for the audience since they require more attention to space and movement on the board. The use of *right over here* multiple times during one video is evidence for repeated call for attention to the virtual board. Because there is no face or body to gesture, it might be more difficult to follow the movement on the board in comparison to face to face lecture.

**Interacting with virtual audience is not difficult, right?**

Despite the lack of an audience present in the same space and time, KAM corpus still contains 180 occurrences of *right?* per 1 million words. Even though Khan is addressing a virtual audience, he still uses *right* quite frequently to prompt the audience to agree with him as he is presenting his mathematical solutions. The use of *right?*, which is more frequent in the MITM corpus, adds a more conversational aspect to his teaching. It can be also understood as engaging the audience in a shared activity of learning. If we consider that most of *chalk talk* is done with the personal pronoun *we* as the subject of action (i.e. *we are adding x to y, right?)*, the response elicitor further cements engagement understood as shared action performance of Khan and the audience.

## Limitations and future directions

There is a number of limitations of this exploratory study. First, my research was limited in terms of the focus on mostly lexical analysis of linguistic features, with minor use of grammatical analysis. Secondly, using open source data is a double-edged sword – it allows for otherwise unparalleled analysis of a particular register (mathematical instruction), but it also means that more quality control needs to be conducted in terms of transcription conventions. Another limitation of this study is the fact that it reveals differences in virtual and face-to-face instruction between Khan Academy platform and MIT lectures, but more studies need to be done on other video instruction sources (and instructors) to understand the difference in language use for this particular academic purpose. Also, the use of two large corpora presented challenges and opportunities for research: while they can be used to analyze patterns of linguistic feature use, the qualitative analysis becomes a challenge in terms of coding discourse functions and uses

This exploratory study could serve as a starting point for future research on spoken academic discourse in online instruction. Despite the fact that made-for-online video instruction materials have been widely used in academic instruction, little is known about the features of this discourse. Understanding what kind of language students might encounter in online environment can help in developing English for Academic Purposes pedagogical materials to prepare students better for such experience.

It is also important to note that educational studies on the effectiveness of video online instruction only covered a very basic understanding of personalization as use of the personal pronoun *you / your* (Moreno & Mayer, 2004). Future research should examine how particular linguistic features facilitate learning and transfer, especially when the design accounts for linguistic variation.

This exploratory study will hopefully motivate additional research on online spoken academic discourse and other aspects of open educational resources. The reach and impact of the instructional materials offered through Khan Academy is immense yet little is known about its linguistic features in comparison to more traditional face-to-face academic lectures. The current study provides a first step towards understanding the functions of some of these features in one domain (mathematics). As more and more students use online video instruction, more research is needed to understand how to engage them using academic language. Khan Academy is a perfect site for conducting such research through the use of data that they make available through their infrastructure.

The current study revealed that Khan as an online tutor achieves to some extent the goals he sets for his pedagogy: to talk with the audience (a conversation) rather than talking at them. Khan knows that reducing the distance between him and his students can encourage them to try again and again when they fail to solve problems. It is a challenge to be close to someone when you talk into the dark void of the virtual blackboard, and have to use imagination to picture a student sitting somewhere else in the world, listening to your voice, as they prepare for challenging classes and exams.

# REFERENCES

Ädel, A. (2012). "What I want you to remember is…": Audience orientation in monologic academic discourse. *English Text Construction*, *5*(1), 101-127.

Anderson, J. R.(1993). *Rules of the Mind*. Hillsdale, NJ: Erlbaum.

Anthony, L. (2019). *AntConc* (Version 3.5.8) [Computer Software]. Tokyo, Japan: Waseda University. Available from http://www.laurenceanthony.net/software

Artemeva, N., & Fox, J. (2011). The writing's on the board: The global and the local in teaching undergraduate mathematics through chalk talk. *Written Communication*, *28*(4), 345-379.

Baker, P. (2006). *Glossary of corpus linguistics*. Edinburgh University Press.

Bakhtin, M. M. (1981). The dialogic imagination: Four essays, ed. *Michael Holquist, trans. Caryl Emerson and Michael Holquist (Austin: University of Texas Press, 1981)*, *84*(8), 80-2.

Bamford, J. (2005). Interactivity in academic lectures: The role of questions and answers. *Dialogue within discourse communities: Metadiscursive perspectives on academic genres*, 123-145.

Barbieri, F. (2015). Involvement in university classroom discourse: Register variation and interactivity. *Applied Linguistics*, *36*(2), 151–173. https://doi.org/10.1093/applin/amt030

Barbieri, F., & Eckhardt, S. E. (2007). Applying corpus-based findings to form-focused instruction: The case of reported speech. *Language Teaching Research*, *11*(3), 319-346.

Barlow, M. (2013). Individual usage: a corpus-based study of idiolects. *International Journal of Corpus Linguistics*, *18*(4).

Baroody, A. J. (2003). The development of adaptive expertise and flexibility: The integration of conceptual and procedural knowledge. *The development of arithmetic concepts and skills: Constructing adaptive expertise*, 1-33.

Basso, K. (1974). The ethnography of writing. *Explorations in the ethnography of speaking*, *425*, 432.

Benson, M. J. (1989). The academic listening task: A case study. *Tesol Quarterly*, *23*(3), 421-445.

Benveniste, E. (1971). Problems in General Linguistics. 1966. *Trans. Mary Elizabeth Meek. Coral Gables: U of Miami P*, 5-23.

Besnier, N. (1994). Involvement in linguistic practice: An ethnographic appraisal. *Journal of Pragmatics*, *22*(3–4), 279–299. https://doi.org/10.1016/0378-2166(94)90113-9

Bhatia, V. K. (1991). A genre-based approach to ESP materials. *World Englishes 10* https://doi.org/10.1111/j.1467-971X.1991.tb00148.x

Brown, P. and C. Fraser. (1979). Speech as a marker of situation. In Klaus R. Scherer and Howard Giles, eds., *Social Markers in Speech*. Cambridge: Cambridge University Press, pp. 33–62.

Buchstaller, I., & Van Alphen, I. (Eds.). (2012). *Quotatives: Cross-linguistic and cross-disciplinary perspectives* (Vol. 15). John Benjamins Publishing.

Biber, D. (1988). *Variation across speech and writing.* New York: Cambridge University Press.

Biber, D. (1994). An analytical framework for register studies. *Sociolinguistic perspectives on register*, 31-56.

Biber, D. (1995). *Dimensions of register variation.* New York: Cambridge University Press.

Biber, D. (2015). Corpus-Based and Corpus-Driven Analyses of Language Variation and Use. In B. Heine & H. Narrog (Eds.), *The Oxford Handbook of Linguistic Analysis* (pp. 1–30). Oxford University Press. https://doi.org/10.1093/oxfordhb/9780199677078.013.0008

Biber, D., & Barbieri, F. (2007). Lexical bundles in university spoken and written registers. *English for Specific Purposes 26*(3), 263-286. https://doi.org/10.1016/j.esp.2006.08.003

Biber, D., & Conrad, S. (2009). *Register, genre, and style*. Cambridge University Press.

Biber, D., Conrad, S., & Cortes, V. (2004). If you look at…: Lexical bundles in university teaching and textbooks. *Applied linguistics*, *25*(3), 371-405.

Biber, D., Conrad, S., & Reppen, R. (1998). *Corpus linguistics: Investigating language structure and use*. Cambridge University Press.

Biber, D., Conrad, S., Reppen, R., Byrd, P., & Helt, M. (2002). Speaking and writing in the university: A multidimensional comparison. *TESOL Quarterly, 36*(1), 9–48.

Biber, D., Conrad, S., Reppen, R., Byrd, P., Helt, M., Clark, V., Cortes, V., Csomay, E., & Urzua, A. (2004). *Representing language use in the university: Analysis of the TOEFL 2000 spoken and written academic language corpus* (ETS TOEFL Mongraph series, ms-25). Princeton, NJ: Educational Testing Service.

Biber, D., Johansson, S., Leech, G., Conrad, S., & Finegan, E. (1999). *Longman Grammar Of Spoken And Written English*. New York: Longman.

Biber, D., & Finegan, E. (1989). Styles of stance in English: Lexical and grammatical marking of evidentiality and affect. *Text-interdisciplinary journal for the study of discourse*, *9*(1), 93-124.

Biber, D., & Finegan, E. (Eds.). (1994). *Sociolinguistic perspectives on register*. Oxford University Press on Demand.

Block, J. H., & Burns, R. B. (1976). 1: Mastery learning. *Review of research in education*, *4*(1), 3-49.

Bowker, D. (2012). Okay? Yeah? Right?: Negotiating understanding and agreement in master's supervision meetings with international students. Unpublished dissertation.

Brockliss, L. (1996). *Curricula. A History of the University in Europe*, ed de Ridder-Symoens H.

Brown, P., & Levinson, S. C. (1987). *Politeness: Some universals in language usage* (Vol. 4). Cambridge University Press.

Byrd, P., & Constantinides, J. C. (1992). The language of teaching mathematics: Implications for training ITAs. *TESOL Quarterly*, *26*(1), 163-167.

Byrnes, J. P., & Wasik, B. A. (1991). Role of conceptual knowledge in mathematical procedural learning. *Developmental psychology, 27*(5), 777.

Cairns, B. (1991). Spatial Deixis-The Use of Spatial Co-ordinates in Spoken Language. *Working Papers in Linguistics*, *38*, 19-28.

Carter, R., & McCarthy, M. (2006). *Cambridge grammar of English: a comprehensive guide; spoken and written English grammar and usage*. Ernst Klett Sprachen.

Chafe, W. (1982). Integration and involvement in speaking, writing, and oral literature. *Spoken and written language: Exploring orality and literacy*, 35-54.

Chang, Y. Y. (2012). The use of questions by professors in lectures given in English: Influences of disciplinary cultures. *English for Specific Purposes*, *31*(2), 103–116. https://doi.org/10.1016/j.esp.2011.08.002

Cheng, S. W. (2012). "That's it for today": Academic lecture closings and the impact of class size. *English for Specific Purposes, 31*, 234– 248.

Clark, R. C., & Mayer, R. E. (2016). *E-learning and the science of instruction: Proven guidelines for consumers and designers of multimedia learning*. John Wiley & Sons.

Cortes, V. (2004). Lexical bundles in published and student disciplinary writing: Examples from history and biology. *English for specific purposes*, *23*(4), 397-423.

Crawford Camiciottoli, B. (2004). Interactive discourse structuring in L2 guest lectures: Some insights from a comparative corpus-based study. *Journal of English for Academic Purposes*, *3*(1), 39–54. https://doi.org/10.1016/S1475-1585(03)00044-4

Crawford Camiciottoli, B. (2008). Interaction in academic lectures vs. written text materials: The case of questions. *Journal of Pragmatics*, *40*(7), 1216–1231. https://doi.org/10.1016/j.pragma.2007.08.007

Crawford Camiciottoli, B., & Querol-Julián, M. (2016). Lectures. In K. Hyland & P. Shaw (Eds.), *The Routledge Handbook of English for Academic Purposes* (pp. 309–319).

Crystal, D., & Davy, D. (1969). *Investigating English style*. Routledge.

Csomay, E. (2000). Academic lectures: An interface of an oral and literate continuum. *NovELTy*, *7*(3), 30-48.

Csomay, E. (2002). Variation in academic lectures: Interactivity and level of instruction. *Using corpora to explore linguistic variation*, 205-224.

Csomay, E. (2004). Linguistic variation within university classroom talk: A corpus-based perspective. *Linguistics and Education*, *15*(3), 243-274.

Csomay, E. (2006). 'Academic talk in American university classrooms: Crossing the boundaries of oral-literate discourse?,' *Journal of English for Academic Purposes 5*. 117–35

Csomay, E. (2007). A corpus-based look at linguistic variation in classroom interaction: Teacher talk versus student talk in American University classes. *Journal of English for Academic Purposes*, *6*(4), 336-355.

Csomay, E. (2012). Lexical bundles in discourse structure: A corpus-based study of classroom discourse. *Applied linguistics*, *34*(3), 369-388.

Csomay, E., & Cortes, V. (2010). Lexical bundle distribution in university classroom talk. In *Corpus-linguistic applications* (pp. 153-168). Brill Rodopi.

DeCarrico, J., & Nattinger, J. R. (1988). Lexical phrases for the comprehension of academic lectures. *English for specific purposes*, *7*(2), 91-102.

Deroey, K. L. (2015). Marking importance in lectures: Interactive and textual orientation. *Applied linguistics*, *36*(1), 51-72.

Dittmar, Norbert. 1996. Explorations in 'Idiolects'. In Robin Sackmann and Monika Budde (eds). Theoretical Linguistics and Grammatical Description: Papers in honour of Hans-Heinrich Lieb. Amsterdam: Benjamins.

diSessa, A. A., Gillespie, N. M., & Esterly, J. B.(2004). Coherence versus fragmentation in the development of the concept of force. *Cognitive Science, 28*, 843–900.

Dudley-Evans, T. (1994). Variations in the discourse patterns favoured by different disciplines and their pedagogical implications. *Academic listening: Research perspectives*, 146-158.

Duranti, A. (1985). Sociocultural dimensions of discourse. *Handbook of discourse analysis*, *1*, 193-230.

Eisenstein, E. L. (1997). *The printing press as an agent of change: Communications and cultural transformations in early-modern Europe*. Cambridge, UK: Cambridge University Press.

Ellis, J., & Ure, J. (1969). Language varieties: register. *Encyclopedia of linguistics: Information and control*, 251-259.

Ellis, N. C. (1996). Sequencing in SLA: Phonological memory, chunking, and points of order. *Studies in second language acquisition*, *18*(1), 91-126.

Emmison, M., Butler, C. W., & Danby, S. (2011). Script proposals: A device for empowering clients in counselling. *Discourse studies*, *13*(1), 3-26.

Enfield, N. J., & Sidnell, J. (2017). *The concept of action*. Cambridge University Press.

Evison, J. (2010). What are the basics of analysing a corpus. *The Routledge handbook of corpus linguistics*, 122-135.

Firth, J. R. (1968). *Selected papers of JR Firth, 1952-59*. Indiana University Press.

Fleischman, S., & Yaguello, M. (2004). Discourse markers across languages. *Discourse across languages and cultures*, 129-147.

Flowerdew, J. (1994). Research of relevance to L2 lecture comprehension: An overview. In J. Flowerdew (Ed.), *Academic listening* (pp. 7–29). Cambridge: Cambridge University Press.

Flowerdew, J., & Miller, L. (1995). On the notion of culture in L2 lectures. *TESOL quarterly*, *29*(2), 345-373.

Flowerdew, J., & Miller, L. (1997). The teaching of academic listening comprehension and the question of authenticity. *English for Specific Purposes,16*(1), 27–46.

Fortanet, I. (2004). The use of "we" in university lectures: Reference and function. *English for Specific Purposes, 23*(1), 45–66. https://doi.org/10.1016/S0889-4906(03)00018-8

Fox, J., & Artemeva, N. (2012). The cinematic art of teaching university mathematics: chalk talk as embodied practice. *Journal Multimodal Communication*, *1*(1), 83-103.

Friesen, N. (2011). The Lecture as a Transmedial Pedagogical Form. *Educational Researcher*, *40*(3), 95–102. https://doi.org/10.3102/0013189x11404603

Friginal, E., Lee, J. J., Polat, B., & Roberson, A. (2017). *Exploring spoken English learner language using corpora: Learner talk*. Springer.

Gerofsky, S. (2004). *A man left Albuquerque heading east: Word problems as genre in mathematics education*. New York, NY: Peter Lang.

Gómez, I. F. (2006). Interaction in Academic Spoken English: The Use of 'I' and 'You' in the MICASE, 35–51. https://doi.org/10.1007/978-0-387-28624-2_3

Gray, B., & Biber, D. E. (2015). Phraseology. In *The Cambridge handbook of English corpus linguistics* (pp. 125-145). Cambridge University Press.

Greaves, C., & Warren, M. (2010). What can a corpus tell us about multi-word units. *The Routledge handbook of corpus linguistics. London: Routledge*, 212-226.

Gries, S. T. (2008). Dispersions and adjusted frequencies in corpora. *International Journal of Corpus Linguistics*, *13*(4), 403-437.

Gumperz, J. J. (1982). *Discourse strategies* (Vol. 1). Cambridge University Press.

Guo, P. J., Kim, J., & Rubin, R. (2014, March). How video production affects student engagement: An empirical study of MOOC videos. In *Proceedings of the first ACM conference on Learning@ scale conference* (pp. 41-50). ACM.

Haas, M. R. (1969). "Exclusive" and "inclusive": A look at early usage. *International Journal of American Linguistics, 35*, 1 - 6.

Haberland, H. (1986). Reported speech in Danish. *Direct and indirect speech*, 219-254.

Hadjidemetriou, C., & Williams, J. (2002). Children's graphical conceptions. *Research in Mathematics Education*, *4*(1), 69-87.

Halliday, M.A.K. (1978) *Language as a Social Semiotic*. London: Edward Arnold.

Harré, R., Brockmeier, J., & Mühlhäusler, P. (1999). *Greenspeak: A study of environmental discourse*. Thousand Oaks, CA: Sage Publications.

Heritage, J. (2012). Epistemics in action: Action formation and territories of knowledge. *Research on Language & Social Interaction*, *45*(1), 1-29.

Heritage, J., & Raymond, G. (2005). The terms of agreement: Indexing epistemic authority and subordination in talk-in-interaction. *Social psychology quarterly*, *68*(1), 15-38.

Herring, S. C., & Androutsopoulos, J. (2015). Computer-mediated discourse 2.0. *The handbook of discourse analysis*, *2*, 127-151.

Holt, E., & Clift, R. (Eds.). (2006). *Reporting talk: Reported speech in interaction* (Vol. 24). Cambridge University Press.

Hoey, M. (2002). Signalling in discourse: a functional analysis of a common discourse pattern in written and spoken English. In *Advances in written text analysis* (pp. 40-59). Routledge.

Holt, Elizabeth, 1996. Reporting on talk: the use of direct reported speech in conversation. Res. Lang. Soc. Interact. 29 (3), 219e245.

Holt, Elizabeth, 2000. Reporting and reacting: concurrent responses to reported speech. *Res. Lang. Soc. Interact. 33* (4), 425e454.

Horsmann, T., Erbs, N., & Zesch, T. (2015). Fast or Accurate?-A Comparative Evaluation of PoS Tagging Models. In *GSCL* (pp. 22-30).

Holt, Elizabeth, 2007. "I'm eyeing your chop up mind": reporting and enacting. In: Holt, E., Clift, R. (Eds.), *Reporting Talk: Reported Speech in Interaction.* Cambridge University Press.

Hicks, M., Reid, I., & George, R. (2001). Enhancing on-line teaching: Designing responsive learning environments. *International Journal for Academic Development*, *6*(2), 143-151.

Hill, H. C., Ball, D. L., & Schilling, S. G. (2008). Unpacking pedagogical content knowledge: Conceptualizing and measuring teachers' topic-specific knowledge of students. *Journal for research in mathematics education*, 372-400.

Hirvonen, P. (2016). Positioning theory and small-group interaction: Social and task positioning in the context of joint decision-making. *Sage Open*, *6*(3), 2158244016655584.

Hunston, S. (2010). How can a corpus be used to explore patterns. *The Routledge handbook of corpus linguistics*, 152-166.

Hymes, D. (1974). Ways of speaking. *Explorations in the ethnography of speaking*, *1*, 433-451.

Hyland, K. (2005). Stance and engagement: A model of interaction in academic discourse. *Discourse studies*, *7*(2), 173-192.moon

Hyland, K. (2008). As can be seen: Lexical bundles and disciplinary variation. *English for specific purposes*, *27*(1), 4-21.

Hyland, K. (2009). *Academic discourse: English in a global context*. A&C Black.

Johns, A. M. (1981). Necessary English: A faculty survey. *TESOL Quarterly*, *15*(1), 51-57.

Johns, A. M., & Dudley-Evans, T. (1991). English for specific purposes: International in scope, specific in purpose. *TESOL Quarterly, 25*(2), 297-314.

Johnson, A. (2006). *Power, privilege, and difference*. Boston, MA: McGraw-Hill.

Jones, S. E. (2007). Reflections on the lecture: Outmoded medium or instrument of inspiration? *Journal of Further and Higher Education, 31,* 397–406.

Jung, E. H. (2003). The role of discourse signaling cues in second language listening comprehension. *The Modern Language Journal*, *87*(4), 562-577.

Kamio, A. (2001). English generic we, you, and they: An analysis in terms of territory of information. *Journal of Pragmatics*, *33*(7), 1111-1124.

Katriel, T., & Dascal, M. (1989). Speaker's commitment and involvement in discourse. *From Sign to Text/Ed. Y. Tobin.–Amsterdam (Philadelphia): John Benjamins Publishing Company*, 275-295.

Kena, G., Hussar, W., McFarland, J., de Brey, C., Musu-Gillette, L., Wang, X., ... & Barmer, A. (2016). The Condition of Education 2016. NCES 2016-144. *National Center for Education Statistics*.

Khan, S. (2012). *The one world schoolhouse: Education reimagined*. Twelve.

Koedinger, K. R., & Nathan, M. J. (2004). The real story behind story problems: Effects of representations on quantitative reasoning. *The Journal of the Learning Sciences*, *13*(2), 129-164.

Koester, A. (2014). "We'd be prepared to do something, like if you say…" hypothetical reported speech in business negotiations. *English for Specific Purposes*, *36*, 35-46.

Koester, A., & Handford, M. (2018). 'It's not good saying "Well it it might do that or it might not"': Hypothetical reported speech in business meetings. *Journal of Pragmatics*, *130*, 67-80.

Krebs, J. R., Dawkins, R., & Davies, N. B. (1984). Behavioral ecology: An integrated approach.

Krathwohl, D. R. (2002). A revision of Bloom's taxonomy: An overview. *Theory into practice*, *41*(4), 212-218.

Lagenhove, L. V., & Harré, R. (1999). Positioning theory: Moral contexts of intentional action. *Maiden, Mass.: Blackwell*.

Lambert, C. (2012, March/April). Twilight of the Lecture. *Harvard Magazine*, 23-27.

Laurillard, D. (2002). *Rethinking university teaching: A framework for the effective use of educational technology* (2nd ed.). London: Routledge.

Leech, G. (2007). New resources, or just better old ones? The Holy Grail of representativeness. In *Corpus linguistics and the web* (pp. 133-149). Brill Rodopi.

Lee, J. J. (2009). Size matters: An exploratory comparison of smalland large- class university lecture introductions. *English for specific purposes, 29*, 42– 57.

Lee, J. J., & Subtirelu, N. C. (2015). Metadiscourse in the classroom: A comparative analysis of EAP lessons and university lectures. *English for Specific Purposes*, *37*, 52-62.

Leinhardt, G., Zaslavsky, O., & Stein, M. K. (1990). Functions, graphs, and graphing: Tasks, learning, and teaching. *Review of educational research*, *60*(1), 1-64.

Lemke, J. L. (1989). *Using language in the classroom*. Oxford University Press, USA.

Lemke, J. L. (2003). Mathematics in the middle: Measure, picture, gesture, sign, and word. *Educational perspectives on mathematics as semiosis: From thinking to interpreting to knowing*, *1*, 215-234.

Levinson, Stephen C. 1983. *Pragmatics.* Cambridge, England: Cambridge University.

Malinowski, B., Ogden, C. K., & Richards, I. A. (1923). The meaning of meaning. *New York & London: Harcourt Brace Jovanovich*.

Mauranen, A. 2000. Strange Strings in Translated Language: A Study on Corpora. In M. Olohan (ed.) *Intercultural Faultlines. Research Models in Translation Studies 1: Textual and CognitiveAspects*. Manchester: St. Jerome Publishing, 119 -141.

Mayer, R. E. (2014). Cognitive Theory of Multimedia Learning. *The Cambridge Handbook of Multimedia Learning*, 43.

Mayer, R. E., Fennell, S., Farmer, L., & Campbell, J. (2004). A personalization effect in multimedia learning: Students learn better when words are in conversational style rather than formal style. *Journal of Educational Psychology*, *96*(2), 389.

Mayer, R. E., & Moreno, R. (2003). Nine ways to reduce cognitive load in multimedia learning. *Educational psychologist*, *38*(1), 43-52.

Mesa, V., & Chang, P. (2010). The language of engagement in two highly interactive undergraduate mathematics classrooms. *Linguistics and Education*, *21*(2), 83-100.

McEnery, T., Xiao, R., & Tono, Y. (2006). *Corpus-based language studies: An advanced resource book*. Taylor & Francis.

Miller, C. R. (1984). Genre as social action. *Quarterly journal of speech*, *70*(2), 151-167.

Moon, R. (2010). What can a corpus tell us about lexis. *The Routledge handbook of corpus linguistics*, 197-211.

Morell, T. (2004). Interactive lecture discourse for university EFL students. *English for specific purposes*, *23*(3), 325-338.

Moreno, R., & Mayer, R. E. (2000). Engaging students in active learning: The case for personalized multimedia messages. *Journal of Educational Psychology*, *92*(4), 724.

Moreno, R., & Mayer, R. E. (2004). Personalized messages that promote science learning in virtual environments. *Journal of Educational Psychology*, *96*(1), 165.

Morgan, C. (2006). What does social semiotics have to offer mathematics education research? *Educational Studies in Mathematics, 61*, 219-245.

Morell, T. (2004). Interactive lecture discourse for university EFL students. *English for Specific Purposes*, *23*(3), 325–338. https://doi.org/10.1016/S0889-4906(03)00029-2

Morell, T. (2007). What enhances EFL students' participation in lecture discourse? Student, lecturer and discourse perspectives. *Journal of English for Academic Purposes*. https://doi.org/10.1016/j.jeap.2007.07.002

Myers, G. (1999). Unspoken speech: Hypothetical reported discourse and the rhetoric of everyday talk. *Text-Interdisciplinary Journal for the Study of Discourse*, *19*(4), 571-590.

NASA. (2018). Salman Khan - Khan Academy: Education Reimagined. Retrieved from https://www.nasa.gov/ames/ocs/2014-summer-series/salman-khan

Nathan, M. J., Koedinger, K. R., & Alibali, M. W. (2001, April). Expert blind spot: When content knowledge eclipses pedagogical content knowledge. In *Proceedings of the third international conference on cognitive science* (pp. 644-648).

Nathan, M. J., & Petrosino, A. (2003). Expert blind spot among preservice teachers. *American Educational Research Journal*, *40*(4), 905-928.

Nattinger, J. R., & DeCarrico, J. S. (1992). *Lexical phrases and language teaching*. Oxford University Press.

Nesi, H., & Basturkmen, H. (2006). Lexical bundles and discourse signalling in academic lectures. *International Journal of Corpus Linguistics*, *11*(3), 283-304.

O'Keeffe, A., Clancy, B., & Adolphs, A. (2011). *Introducing pragmatics in use*. London: Routledge.

O'Halloran, K. L. (1999). Interdependence, interaction and metaphor in multisemiotic texts. *Social Semiotics*, *9*(3), 317-354.

O'Halloran, K. L. (2015). The language of learning mathematics: A multimodal perspective. *The Journal of Mathematical Behavior*, *40*, 63-74.

Othman, Z. (2010). The use of okay, right and yeah in academic lectures by native speaker lecturers: Their 'anticipated'and 'real'meanings. *Discourse Studies*, *12*(5), 665-681.

Ong, W. (1982). *Orality and literacy: The technologizing of the word.* London: Routledge.

Pawley, A., & Syder, F. H. (1983). Natural selection in syntax: Notes on adaptive variation and change in vernacular and literary grammar. *Journal of Pragmatics*, *7*(5), 551-579.

Pimm, D. (1989). Speaking mathematically: Communication in mathematics classrooms.

Prince, M. (2004). Does active learning work? A review of the research. *Journal of Engineering Education*, *93*(3), 223-231.

Quirk, R. (2010). *A comprehensive grammar of the English language*. Pearson Education India.

Richards, J. C. (1983). Listening comprehension: Approach, design, procedure. *TESOL Quarterly*, *17*(2), 219-240.

Richards, J. C. 1983. "Listening comprehension: approach, design, procedure". *TESOL Quarterly* 17 (2): 219-239.

Rittle-Johnson, B. (2017). Developing mathematics knowledge. *Child Development Perspectives*, *11*(3), 184-190.

Rittle-Johnson, B., & Schneider, M. (2015). Developing conceptual and procedural knowledge of mathematics. *Oxford handbook of numerical cognition*, 1118-1134.

Rittle-Johnson, B., Siegler, R. S., & Alibali, M. W. (2001). Developing conceptual understanding and procedural skill in mathematics: An iterative process. *Journal of educational psychology*, *93*(2), 346.

Romaine, S. (2001). Dialect and dialectology. In R. Mesthrie, & R. E. Asher (Eds.), *Concise encyclopedia of sociolinguistics*. Oxford, UK: Elsevier Science & Technology. Retrieved from https://search-credoreference-com.ezproxy.lib.purdue.edu/content/entry/estsocioling/dialect_and_dialectology/0

Rounds, P. (1987a). Characterizing successful classroom discourse for NNS teaching assistant training. *TESOL Quarterly, 21*(4), 643–671.

Rounds, P. (1987b). Multifunctional personal pronoun use in educational setting. *English for Specific Purposes, 6*(1), 13–29.

Rowland, T. (1999). Pronouns in mathematics talk: Power, vagueness and generalisation. *For the Learning of Mathematics*, *19*(2), 19-26.

Sacks, H., Schegloff, E. A., & Jefferson, G. (1978). A simplest systematics for the organization of turn taking for conversation. In *Studies in the organization of conversational interaction* (pp. 7-55). Academic Press.

Schleppegrell, M. J. (2007). The linguistic challenges of mathematics teaching and learning: A research review. *Reading & Writing Quarterly*, *23*(2), 139-159.

Schleef, E. (2008). The "lecturer's OK" revisited: changing discourse conventions and the influence of academic division. *American Speech*, *83*(1), 62-84.

Schmid, H. (1999). Improvements in part-of-speech tagging with an application to German. In *Natural language processing using very large corpora* (pp. 13-25). Springer, Dordrecht.

Schneider, M., & Stern, E. (2009). The inverse relation of addition and subtraction: A knowledge integration perspective. *Mathematical Thinking and Learning*, *11*(1-2), 92-101.

Shulman, L. S. (1986). Those who understand: Knowledge growth in teaching. *Educational researcher*, *15*(2), 4-14.

Simmons, K., & LeCouteur, A. (2011). 'Hypothetical active-voicing': Therapists 'modelling'of clients' future conversations in CBT interactions. *Journal of Pragmatics, 43*(13), 3177-3192.

Sinclair, J. (2004). Trust the text. In *Trust the text* (pp. 19-33). Routledge.

Sinclair, J., & M. Coulthard. (1975). *Towards an Analysis of Discourse.* Oxford University Press.

Sinclair, J. M. (Ed.). (1987). *Looking up: An account of the COBUILD project in lexical computing and the development of the Collins COBUILD English language dictionary*. Collins Elt.

Siyanova-Chanturia, A., & Martinez, R. (2014). The idiom principle revisited. *Applied Linguistics*, *36*(5), 549-569.

Spiegelberg, H. (1973). On the right to say "we": A linguistic and phenomenological analysis. In G. Psathas (Ed.), *Phenomenological sociology* (pp. 129- 156). New York: Wiley.

Speer, N. M., Smith III, J. P., & Horvath, A. (2010). Collegiate mathematics teaching: An unexamined practice. *The Journal of Mathematical Behavior*, *29*(2), 99-114.

Steffens, N. K., & Haslam, S. A. (2013). Power through 'us': Leaders' use of we-referencing language predicts election victory. *PLoS One*, *8*(10), e77952.

Staples, S., Egbert, J., Biber, D., & Conrad, S. (2015). 24 Register Variation A Corpus Approach. *Discourse Analysis*, 505.

Stubbs, M. (2001). *Words and phrases: Corpus studies of lexical semantics*. Oxford: Blackwell Publishers.

Swarts, J. (2012). New modes of help: Best practices for instructional video. *Technical Communication*, *59*(3), 195-206.

Swales, J. M., & Malczewski, B. (2001). *Discourse management and new-episode flags in MICASE* (pp. 145-164). na.

Tannen, D. (1982). Oral and literate strategies in spoken and written narratives. *Language*, 1-21.

Tannen, D. (1985). Relative focus on involvement in oral and written discourse. *Literacy, language, and learning: The nature and consequences of reading and writing*, 124-147.

Thompson, S. (1994). Frameworks and contexts: A genre-based approach to analysing lecture introductions. *English for Specific Purposes, 13*(2), 171–186.

Thompson, S. (1998). 11 Why Ask Questions in Monologue? Language. In *Language at Work: Selected Papers from the Annual Meeting of the British Association for Applied Linguistics Held at the University of Birmingham, September 1997* (Vol. 13, p. 137). Multilingual Matters.

Thompson, S. E. (2003). Text-structuring metadiscourse, intonation and the signalling of organisation in academic lectures. *Journal of English for Academic Purposes*, *2*(1), 5-20.Shulman, L. S. (1986). Those who understand: Knowledge growth in teaching. *Educational researcher*, *15*(2), 4-14.

Tottie, G., & Hoffmann, S. (2006). Tag questions in British and American English. *Journal of English Linguistics*, *34*(4), 283-311.

Traugott, E. C. (2003). From subjectification to intersubjectification. *Motives for language change*, *124*, 139.

Veel, R. (1999). Language, knowledge and authority in school mathematics. *Pedagogy and the shaping of consciousness: Linguistic and social processes*, 185-216.

Waggoner, M. (1984). The new technologies versus the lecture tradition in higher education: Is change possible? *Educational Technology, 24* (3),7-12.

Wagner, D., & Herbel-Eisenmann, B. (2008). "Just don't": The suppression and invitation of dialogue in the mathematics classroom. *Educational Studies in Mathematics*, *67*(2), 143.

Wray, A. (2005). *Formulaic language and the lexicon*. Cambridge University Press.

Wray, A., & Perkins, M. R. (2000). The functions of formulaic language: An integrated model. *Language & Communication*, *20*(1), 1-28.

Yeo, J. Y., & Ting, S. H. (2014). Personal pronouns for student engagement in arts and science lecture introductions. *English for Specific Purposes*, *34*, 26-37.

Young, L. (1994). University lectures–macro-structure and micro-features. *Academic listening: Research perspectives*, 159-176.