

Cross-Domain Classification of Education Talk-Turns

Anonymous submission

Abstract

The study of classroom discourse is essential for enhancing child development and educational outcomes in academic settings. Prior research has focused on the annotation of conversational talk-turns within the classroom, offering a statistical analysis of the various types of discourse prevalent in these environments. In this work, we explore the generalizability and transferability of these discourse codes across different educational domains via automatic text classifiers. We examine two distinct English-language classroom datasets from the domains of literacy and mathematics. Our results show that models exhibit high accuracy and generalizability when the training and test datasets originate from the same or similar domains. However, as the distance between the training and test domains increases in terms of subject matter and teaching methodology, we observe a decline in model performance. We also observe that accompanying each talk turn with dialog-level context improves the accuracy of the generative models. We conclude by offering suggestions on how to enhance the generalization of these methods to novel domains, proposing directions for future studies to investigate new methods and techniques for boosting the model adaptability across varied educational domains.

1 Introduction

In recent years, computational approaches have increasingly demonstrated their potential in capturing and analyzing discourse-level features within educational settings (Ganesh et al., 2021). Previous research in this domain has provided valuable insights, particularly in the context of specific educational domains or settings (Wang et al., 2023). However, these studies often limit their focus to single domains, leaving a gap in understanding the adaptability and effectiveness of these models across varied educational contexts.

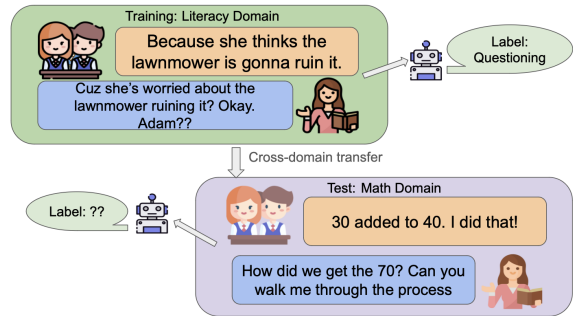


Figure 1: Cross Domain Training in an educational context where model trained on classroom discourse from the literacy domain is applied to the math domain to predict talk moves in a new context.

Addressing this gap, we aim to understand the divide between domain-specific and generalizable models in educational discourse analysis. We incorporate a varied set of models, from fine-tuned transformer-based models to in-context learning approaches using generative language models. Our focus is to evaluate how the performance of these models varies when the distance between the contexts of the training and test domains is increased. By 'distance,' we refer not only to the difference in academic domains (e.g., English vs. Mathematics) but also to differences in educational materials such as textbooks used and instructional variations among teachers. Our approach allows us to assess the generalizability of various computational methods across diverse educational settings. Figure 1 illustrates an example of cross-domain training in an educational context where a machine learning model trained on classroom discourse from the one domain is applied to an unseen domain to predict talk moves in a new context.

The primary contributions of this research are (1) an in-depth analysis of the generalizability of language models across diverse educational domains thus providing critical insights into the adaptability and limitations of these models when applied

068 in different educational contexts; (2) experiments
069 evaluating several types of models, ranging from
070 fine-tuned transformer-based encoder models to
071 in-context learning approaches, for classifying ed-
072 ucational talk-turns with discourse codes; and (3)
073 the creation of a ground truth dataset for testing
074 the generalizability and accuracy of various lan-
075 guage models for cross-application classification
076 of discourse codes.

077 Using data ranging from read-aloud discussions
078 in early education classrooms to interactions from
079 mathematics classes, we investigate how model
080 performance varies across contexts. We aim to
081 shed light on the strengths and limitations of cur-
082 rent computational approaches in educational dis-
083 course analysis. Our results demonstrate that mod-
084 els exhibit high accuracy of discourse codes when
085 the training and test datasets originate from the
086 same domains, but as expected, the effectiveness of
087 these models begins to decrease when the distance
088 between the training and testing data increases.
089 Despite this decline in performance, the explo-
090 ration into transformer-based and in-context learn-
091 ing models in cross-domain scenarios remains cru-
092 cial in this area of study in order to precisely quan-
093 tify the *extent* of this performance dip, especially
094 given recent advancements in large language mod-
095 els (LLMs). Further, we investigate how the choice
096 of model or the inclusion of additional information
097 such as conversational context might help mitigate
098 this drop-off, seeking to highlight which computa-
099 tional approaches might exhibit greater resilience
100 against the challenges posed by domain variance,
101 thereby contributing to the generalizability of these
102 models across diverse educational settings.

103 2 Related Work

104 The dynamics of student-teacher classroom dis-
105 course play a pivotal role in shaping the experience
106 and outcomes of students. Several papers have stud-
107 ied this phenomenon, particularly in the context of
108 K-12 mathematics education and other childhood
109 learning environments. For example, [Suresh et al. \(2022a\)](#)
110 found that sustained classroom discourse
111 is a critical component of equitable and a rich learn-
112 ing environment. Towards that goal, they built an
113 extensive collection of human-annotated transcripts
114 from K-12 classroom mathematics lessons as they
115 can be effective tools for understanding discourse
116 patterns in classroom instructions. Furthermore,
117 [Demszky et al. \(2021\)](#) argue that teachers' acknowl-

118 edgement, repetition and reformulation of students'
119 responses has been linked to higher student engage-
120 ment and achievement. The impact of building
121 upon student contributions in the classroom is ex-
122 plored in studies by [Brophy and Good \(1984\)](#) and
123 [Faculty and Michaels \(1993\)](#). They demonstrate
124 that acknowledgment, repetition, and elaboration
125 of student inputs can significantly enhance student
126 learning and academic achievement. [Wright \(2019\)](#)
127 delves into the significance of read-aloud activities
128 in nurturing children's reading skills and knowl-
129 edge. They deduce that engaging in interactive
130 read-alouds is beneficial for children in acquiring
131 new vocabulary, understanding textual functions,
132 and developing a diverse set of skills essential for
133 independent reading. [Giroir et al. \(2015\)](#) explore ef-
134 fective methodologies for implementing read-aloud
135 programs. Their research particularly focuses on
136 integrating aspects of second language acquisition
137 and culturally responsive teaching methods, outlin-
138 ing critical steps and applications for an effective
139 read-aloud strategy.

140 In the context of early childhood education,
141 [Christ et al. \(2023\)](#) study the interactions between
142 teacher and child talk-turns during read-aloud ses-
143 sions. The statistical discourse analysis conducted
144 in this study provides insights into how certain
145 talk-turns can influence children's comprehension
146 responses, thereby emphasizing the critical role of
147 teacher mediation in shaping learning outcomes.
148 Their findings also demonstrate that when chil-
149 dren's talk-turns mediate other children's actions,
150 they act as a predictor for those children's subse-
151 quent responses in terms of comprehension.

152 Given the breadth of actionable findings in this
153 area, a promising direction is to develop tools that
154 assist teachers in refining their instructional strate-
155 gies. [Suresh et al. \(2022b\)](#) outline the development
156 of the TalkBack application. Their tool leverages
157 deep learning capabilities to provide teachers with
158 automated feedback on their discourse strategies,
159 highlighting the importance of automated feedback
160 to enhance and enrich teacher learning. Specif-
161 ically, it aids in refining instructional strategies,
162 thereby enhancing the learning environment.

163 In recent years, advancements in NLP have
164 opened new and effective means of analyzing and
165 enhancing classroom discourse. [Ganesh et al. \(2021\)](#)
166 aims to enhance classroom learning and
167 engagement by developing a system to predict the
168 next talk move (an utterance strategy) in a class-

room discussion, based on the academically productive talk (APT) framework. They present a neural network model aimed at predicting the next talk move in a conversation based on its history and associated talk moves potentially leading to more interactive and personalized learning experiences. In this study (Suresh et al.) incorporate enriched contextual cues from previous and subsequent utterances using a transformer model, called RoBERTa to improve the automated classification of “talk moves” in educational settings. Similarly, (Alic et al., 2022) address the task of creating specific types of questions that promote responsive teaching. The authors created an annotated dataset and employed various supervised and unsupervised learning methods to demonstrate the importance of incorporating computational tools to assist teachers in refining their instructional techniques.

While existing work demonstrated the effectiveness of computational tools to assist in classroom settings, these have typically focused on a single domain. However, in order to use these tools broadly, they must generalize across topical areas and classroom contexts. Therefore, we set out to evaluate the extent to which current methods are able to accurately transfer to new contexts, in terms of both classrooms and teaching domains.

3 Data

In this study, we leveraged existing classroom discourse datasets comprising of turn-level student-teacher interactions that were broken down into dialogue discourse moves and annotated by educational experts. To investigate the generalizability of these codes across different academic domains, we re-annotated a small subset from each dataset using the discourse codes that were originally developed for the other datasets.

3.1 Datasets Used

We used four existing English-language datasets: The MuMo Talk moves dataset (Christ et al., 2023), the National Center for Teacher Effectiveness (NCTE) Transcripts dataset (Demszky and Hill, 2023), and two additional datasets referred to by the pseudonymous teachers of their respective classrooms: Mason (Christ and Cho, 2023) and Newman (Cho and Christ, 2022).

The MuMo Talk moves dataset includes three kindergarten teachers’ interactive read-alouds comprising of 736 talk-turns across six video recorded

Variable	MuMo	Mason	Newman
Response Evaluation	67	345	656
Providing Information	115	290	344
Revoicing	113	330	335
Strategy Related	62	206	302
Questioning	219	563	503
Behavior Management	55	326	354
Turn Management	207	283	325
Total	736	2550	2467

Table 1: Number of talk turns with each label across MuMo, Mason, and Newman datasets. **Total** indicates the number of talk turns in the dataset. Note that each talk turn may have more than one label.

and transcribed sessions. The talk-turns were coded using *a priori* and emergent codes.

The NCTE transcripts consists of the largest dataset of mathematics classroom transcripts available. The dataset consists of 2348 anonymous transcripts of whole lessons collected as part of the National Center for Teacher Effectiveness, NCTE main study, spanning across the K-12 math classrooms across four districts serving largely historically marginalized students.

In the Mason dataset (Christ and Cho, 2023), the authors investigated the engagement of four second-grade emergent bilingual students and their teacher with listening comprehension during interactive read-aloud sessions. These sessions used books with varying levels of cultural relevance. The study aimed to understand how this engagement related to the teacher’s implementation of culturally relevant and sustaining pedagogical practices. To conduct the analysis, the researchers collected data through cultural relevance ratings of the books, video recordings, and transcripts of nine 20-minute lessons, resulting in a total of 2781 talk-turns.

The Newman dataset investigates how two emergent bilingual students from refugee families interacted with texts of varying cultural relevance in a third-grade class in the Midwest U.S by using video recordings and transcripts of 12 read-aloud discussions, interviews, and cultural relevance ratings. In this paper, they analyzed the students’ inference-making processes and examined students’ use of text information, background knowledge, and the coherence of their inferences. This dataset had a total of 2470 talk-turns. Table 1 shows class distribution of variables belonging to Class 1 across MuMo, Mason and Newman datasets and Table 2 shows class distribution of variables belonging to Class 1 across the NCTE dataset.

In assembling our datasets for this study, we

257 began with the established MuMo talk moves code-
 258 book. Recognizing the similarities across vari-
 259 ous datasets in the literacy domain, we extended
 260 this codebook to include the Mason and New-
 261 man datasets. This decision was taken due to
 262 the fact that, despite originating from different ini-
 263 tial codebooks, the discourse codes and categories
 264 across these datasets (MuMo, Mason and New-
 265 man) shared enough similarities to justify a unified
 266 codebook. This also allows for a comprehensive
 267 cross-dataset analysis. Originally, the discourse
 268 codes within these three datasets had over over
 269 100 unique codes across the MuMo, Mason, and
 270 Newman datasets, capturing the subtleties of class-
 271 room interactions between teachers and students.
 272 However, for the purposes of our research, we we
 273 decided to streamline this approach by consolidat-
 274 ing these codes into 15 broader categories. By
 275 implementing this approach we not only simplify
 276 the analysis but also enhances the classification per-
 277 formance of our models, including both generative
 278 and transformer architectures.

279 3.2 Annotation Process

280 To investigate the cross-domain generalizability of
 281 large language models, we sampled data points
 282 from our datasets. Specifically, we chose 140
 283 data points from the MuMo, Mason, and Newman
 284 datasets combined, selecting 10 talk-turns from
 285 each session. MuMo contributed data from 6 ses-
 286 sions, while Mason and Newman each had 4 ses-
 287 sions, collectively providing the 140 data points.
 288 In contrast, the NCTE dataset, which is comprised
 289 of a single extensive session, contributed a total of
 290 100 data points. This selection was made to ensure
 291 a balanced representation of interactions across dif-
 292 ferent educational settings. Once the annotation
 293 guidelines were established, five trained annotators
 294 re-annotated these subsets from each dataset. With
 295 the goal of creating a ground truth for evaluating
 296 the language models, the annotators applied the
 297 discourse codes from the math domain dataset (i.e.,
 298 NCTE) to the literacy domain datasets (i.e., MuMo,
 299 Mason, and Newman) and vis-versa. Each talk-
 300 turn was annotated by at least three annotators to
 301 ensure reliable accuracy and consistency. The inter-
 302 annotator agreement was quantitatively measured
 303 using Krippendorff’s alpha (Krippendorff, 2011)
 304 which yielded the following scores. Please refer
 305 to Table 3 for the results. In case of discrepancies
 306 among annotators, we employed a majority vote

Variable	Count
Student on Task	1964
Teacher on Task	2004
High Uptake	813
Focusing Question	359
Total	2348

Table 2: Number of talk turns with each label in the NCTE dataset. **Total** indicates the number of talk turns in the dataset. Note that each talk turn may have more than one label.

Label	Alpha	Source Dataset	Target Dataset
Response Evaluation	0.849	MMN	NCTE
Providing Information	0.958		
Revoicing	0.899		
Strategy-related	0.818		
Questioning	0.909		
Behavior Management	0.801		
Turn Management	0.936		
Misinformation	N/A	NCTE	MMN
Student on Task	0.912		
Teacher on Task	0.943		
High Uptake	0.847		
Focusing Question	0.851		

Table 3: Krippendorff’s α intercoder agreement scores for the combined datasets of MuMo, Mason, and Newman (MMN) using labels created for the NCTE dataset and vice versa. The label “Misinformation” was never assigned to any text.

307 system, where the label that received the majority
 308 consensus among the three annotators was chosen
 309 as the final label for each talk-turn in our test set.

310 4 Experimental Methodology

311 In our experimental methodology, specifically for
 312 experiments within the same domain, the NCTE
 313 dataset was partitioned using an 80-10-10 split,
 314 allocating 80% of the data for training, 10% for
 315 validation, and the remaining 10% for testing. For
 316 cross-domain experiments, we adopted a distinct
 317 approach where one dataset served as both the train-
 318 ing and validation set, while a dataset from a dif-
 319 ferent domain was designated as the test set. This
 320 strategy was applied to explore the adaptability of
 321 models across varied educational contexts. In the
 322 case of the MuMo, Mason, and Newman datasets,
 323 our experimental design included two setups.

324 First, for intra-domain experiments where the
 325 test set originated from the same dataset, we se-
 326 lected one entire session to function as the test set.
 327 Secondly, for experiments within the domain but
 328 across different datasets, the same session that was
 329 used as the test set in the previous experimental
 330 setup was held out to serve as the test set. This

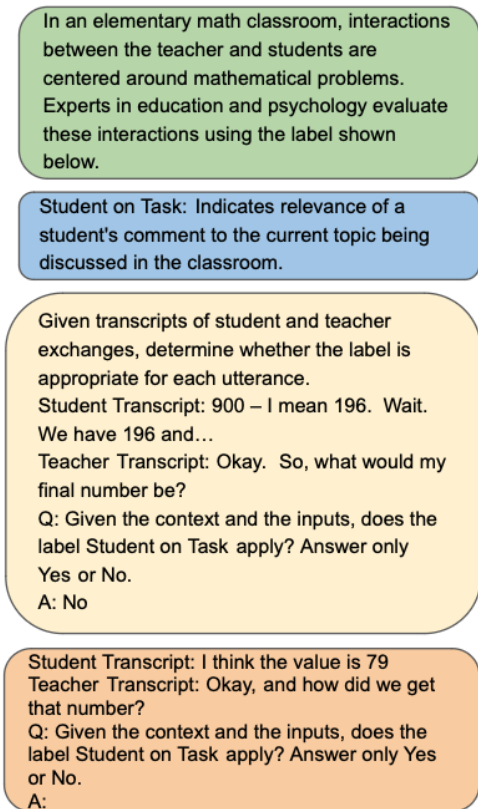


Figure 2: Prompt components for generative models for the math setting (used for NCTE). From top to bottom, the blocks display the context (green), labels (blue), few-shot examples (yellow).

approach allowed us to examine both the domain-specific and cross-domain efficacy of our models. We investigated both fine-tuned transformer encoder models, and auto-regressive generative models focusing on in-context learning.¹

For transformer-based deep learning models, we chose BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019) models, using weights and fine-tuning code from the HuggingFace transformers Wolf et al. (2019) library. We utilized the bert-base-uncased and roberta-base checkpoints along with their default tokenizers. The output from the [CLS] input token was then used as the input for a trainable classification layer. The training hyperparameters for these models are also specified in the Appendix.

For the generative models, we opted to use the Llama2-7B (Touvron et al., 2023) and Mistral-7B-Instruct-v0.2 (Jiang et al., 2023)² models specifi-

cally due to their open weights³ availability. The decision to use open-weights models align with our commitment to transparency and reproducibility allowing others to replicate and extend our experiments without the barriers often associated with proprietary models, and also helps avoid leakage of datasets without the consent of the participants of the original studies (Balloccu et al., 2024). These models operate by receiving an instruction or a prompt as input and generating a response that aligns with the given context or question. The primary objective of employing these models in our experiment was to determine whether auto-regressive models are capable of accurately predicting talk-turn labels particularly in scenarios where there is limited data availability. This focus also aligns with the broader objective of exploring the potential of generative language models to adapt in data-constrained environments, a common challenge in the field of educational discourse analysis.

The experimental setup for the auto-regressive models was conducted in both a zero-shot and a few-shot learning context. In the zero shot setup, the context and label description are prepended to the prompt followed by a transcript from the test set. Finally we ask the model if the given label is appropriate for the transcript. The model is parametrically constrained to answer only in a Yes or No format for calculating accuracy, F1 score and other metrics. We repeated the experiments with same prompt three times to check for any variability in the model’s outputs.

In the few-shot setup, several example interactions between teachers and students, along with their correct labels, were prepended to the prompt in a question-answer format along with the definitions of each output label. Additional instructions developed by the annotators as part of the annotation guidelines were also provided to the model as the input. Each experiment was conducted three times to account for any variance in the model’s results. Similar to the zero-shot setup, experiments were conducted using three different prompts for both Llama2 7B and Mixtral 8x7B. In case of the generative models, the average of the three turns of the best performing prompt was reported. A summary of the various components utilized in this setup can be found in Figures 2 and 3

¹We also explored classical machine learning approaches using bag-of-words features, but found these to always underperform the transformer-based approaches.

²Henceforth referred to simply as “Mixtral.”

³We distinguish between “open source” and “open weights”, where the former includes cases where all code to fully reproduce the model is available, while the latter refers to the open availability of the trained model’s parameters.

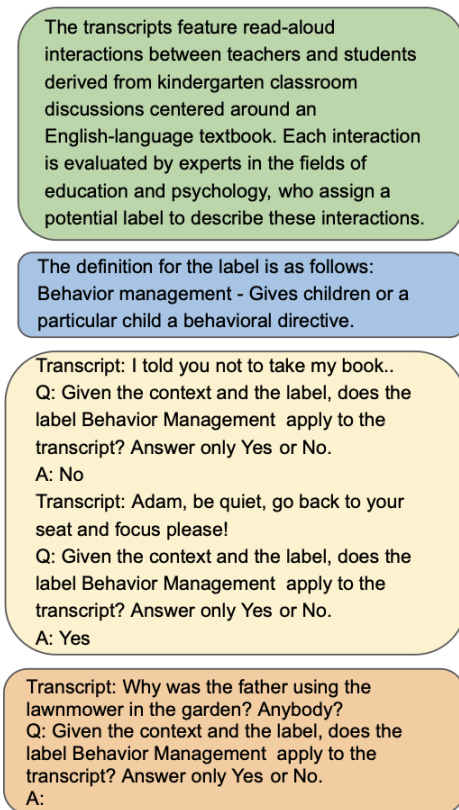


Figure 3: Prompt components for generative models for the read-aloud setting (used for MuMo/Mason/Newman). From top to bottom, the blocks display the context (green), labels (blue), few-shot examples (yellow).

We also experimented with varying the number of prior talk turns provided as context. This setup was only applied to the generative models due to the input size limitations of 512 tokens for the BERT-based models. In this setup, each talk turn was accompanied by the preceding interaction(s) and the speaker tag (whether the talk turn was uttered by a teacher or a student) to provide additional context (where context = 1, 3, and 5 prior interactions) to the model.

5 Results

Tables 5, 6, 7, and 8 present the performance of different classes of models averaged across all the output labels for the domains of literacy and mathematics. Please refer to appendix E for a full breakdown of results for each label and for each model. Among these, fine-tuned models, *i.e.*, BERT and RoBERTa, demonstrated superior performance in most of the experiments over generative models with in-context learning. This was consistent across a majority of test scenarios and across domains.

Among the generative models, Mixtral outperformed Llama2 in most scenarios. Interestingly, Mixtral, when prompted with prior interactions, outperformed BERT and RoBERTa models for certain variables in the binary classification tasks. Despite their overall lower performance compared to fine-tuned models, the fact that these generative models utilized far fewer training data (few shots with $n = 3$ and context $c = \{3, 5\}$) while learning highlights their potential in specific contexts.

The performance metrics, *i.e.*, the F1 scores, showcased a common trend across all models: a higher degree of accuracy when both the training and test sets originated from the same domain. However, we observed a decline as the contextual distance between training and testing data increased. This reflects the challenge of applying machine learning models to diverse educational content due to their varying subject matter and teaching methodology.

A critical observation from our experiments is the distribution of classes within our datasets. Notably, several variables have only a very small number of positive examples in the data. The lack of sufficient support for the majority class in these datasets likely contributed to the lower F1 scores observed for those variables in many scenarios.

5.1 Error Analysis

In this section we investigate the discrepancies between the ground truth labels and the model predictions on the test set. We used the best performing model, *i.e.* BERT out of all the various experiments on a specific test set for this analysis. Table 4 shows paraphrased examples of interactions where the model failed to accurately predict the output label. We found that the model’s predictions were mostly incorrect due to the lack of prior interactions as context. Classroom discourse is inherently continuous and time-series, meaning that understanding any given talk-turn often depends on the preceding turns, and the output labels rely heavily on the interactional context to be accurately classified. For example, in the MuMo dataset, the model mislabels the talk turn “Rocks?” as Class 1 (Questioning) because the context was insufficient and the ground truth wasn’t Class 1 in this specific scenario. The teacher wanted the students to simply repeat the utterance. Similarly, according to the codebook used for this paper, a compliment given by a teacher can be considered Providing Informa-

Test Set	Category	Speaker	Transcript	Actual Label	Pred Label
MuMo	Questioning	Student	Rocks?	0	1
	Turn Management	Teacher	Very good. And let me see what this is down here. A mallow.	1	0
	Questioning	Teacher	Rocks?	1	0
	Providing Information	Teacher	Very good. And a home for everyone.	1	0
Mason	Providing Information	Teacher	We're going to listen to that story now, talk to each other if we notice certain things. We can discuss more tomorrow as well.	1	0
	Revoicing	Teacher	I understand, I get that you are sleepy. But that also means you need to go to bed earlier.	1	0
	Literal Responses	Student	Then we can have two weddings	0	1
	Behavior Management	Teacher	Can you sit down please? [T reading: After the cake was served... We are doing the flower girl] They are pretending to be a flower girl while dancing.	1	0
	Questioning	Teacher	Can you sit down please? [T reading: After the cake was served... We are doing the flower girl] They are pretending to be a flower girl while dancing.	0	1
Newman	Literal Responses	Student	I was attached last winter, everybody hit me. James is their boss. I was very upset. I threw Quincy on accident.	0	1
	Questioning	Student	You know the paper? I mixed the two up. I was gonna write sad and the word said how they feel and then how they feel. How they feel	0	1
	Reading	Teacher	[reads from the book]. So was it okay for Jack to go to the library since there was no book to read from?	1	0
	Questioning	Student	Can I see? I cannot see the book.	0	1
NCTE	Focusing Question	Teacher	Oops, I bet, you know what? I made something, you know what happens?	0	1
	Focusing Question	Teacher	Where would you line up those Xs?	1	0
	Student on Task	Student	I am going to do a shout out.	0	1

Table 4: Error Analysis of Model Predictions Across Different Test Sets

tion, but the model struggled to label some of those interactions accurately.

Furthermore, in the Mason dataset, we noticed that the model failed to label the talk turn “I understand, I get that you are sleepy, but that also means you need to go to bed earlier” as Revoicing. In the prior interaction, the student told the teacher that they were sleepy, and missing this context led to a misclassification. Another issue arose when the teacher’s instructions were mixed with reading sections of a book. For example, the teacher’s phrase “Can you sit down please?” accompanied by the teacher reading a small paragraph from the textbook in the same interaction led to prediction errors. This problem was observed in both the Mason and Newman datasets. Additionally, quite a few errors in the student’s talk turns were a result of short sentences and the lack of context provided to the models. These short, isolated statements were often misclassified because the model couldn’t access the surrounding interactions that would clarify their meaning. While the fine-tuned models gave the best performance, they are limited in their ability to incorporate the necessary context for precise predictions in certain scenarios. Therefore future work might explore the fine-tuning of models with large context windows to reap the benefits of both additional context and fine-tuning.

6 Conclusion

Understanding classroom discourse is pivotal for improving educational outcomes and child development. In this study, we assess the generalizability of discourse codes across distinct educational domains of literacy and mathematics using automatic text classifiers such as transformer based models and in context learning based open weights generative models. We utilized several datasets from prior studies both from literacy and mathematics disciplines; annotated a subset of those data sets to generate ground truths for cross domain classification of educational classifiers. Our findings suggest show that transformer-based models, particularly BERT, and RoBERTa were better at classifying classroom discourse compared to open weights generative models. We further noticed a performance drop when transitioning between different educational domains highlighting the challenges of using large language models in the field of education.

In addition to these findings, we conducted a comprehensive error analysis using the best performing model, providing a fresh perspective on the model failures. We also experimented with providing context to the generative models in the form of prior interactions, and found out that such context could significantly impact the models’ ability to understand and classify discourse accurately.

Train Set	BERT	RoBERTa	Mixtral				Llama2			
			c=0	c=1	c=3	c=5	c=0	c=1	c=3	c=5
<i>MuMo</i>	0.375	0.367	0.326	0.355	0.375	0.363	0.313	0.332	0.351	0.348
<i>Mason</i>	0.497	0.475	0.390	0.426	0.466	0.499	0.357	0.398	0.439	0.495
<i>Newman</i>	0.473	0.466	0.379	0.409	0.448	0.474	0.348	0.387	0.425	0.460
<i>Ms+Nw</i>	0.473	0.453	0.348	0.395	0.430	0.470	0.334	0.378	0.421	0.450
<i>Mu+Ms</i>	0.471	0.463	0.351	0.400	0.422	0.427	0.331	0.376	0.413	0.418
<i>Mu+Nw</i>	0.448	0.438	0.361	0.384	0.419	0.409	0.329	0.362	0.401	0.407

Table 5: Average F1-score for each model across different training sets for test set Mason. **c** denotes the number of prior interactions provided as context to the generative models during classification.

Train Set	BERT	RoBERTa	Mixtral				Llama2			
			c=0	c=1	c=3	c=5	c=0	c=1	c=3	c=5
<i>MuMo</i>	0.350	0.336	0.325	0.321	0.341	0.345	0.306	0.306	0.326	0.316
<i>Mason</i>	0.527	0.471	0.384	0.407	0.423	0.435	0.337	0.364	0.392	0.400
<i>Newman</i>	0.494	0.462	0.387	0.407	0.439	0.447	0.341	0.373	0.399	0.415
<i>Ms+Nw</i>	0.486	0.468	0.379	0.401	0.420	0.421	0.337	0.368	0.392	0.386
<i>Mu+Ms</i>	0.464	0.453	0.342	0.361	0.399	0.395	0.320	0.351	0.378	0.371
<i>Mu+Nw</i>	0.469	0.431	0.343	0.357	0.388	0.407	0.318	0.338	0.388	0.401

Table 6: Average F1-score for each model across different training sets for test set Newman. **c** denotes the number of prior interactions provided as context to the generative models during classification.

Train Set	BERT	RoBERTa	Mixtral				Llama2			
			c=0	c=1	c=3	c=5	c=0	c=1	c=3	c=5
<i>MuMo</i>	0.542	0.517	0.419	0.480	0.485	0.484	0.348	0.420	0.451	0.466
<i>Mason</i>	0.541	0.536	0.455	0.519	0.532	0.570	0.404	0.463	0.516	0.523
<i>Newman</i>	0.515	0.529	0.393	0.455	0.470	0.497	0.342	0.407	0.448	0.471
<i>Ms+Nw</i>	0.488	0.488	0.335	0.367	0.396	0.436	0.306	0.348	0.390	0.406
<i>Mu+Ms</i>	0.601	0.582	0.405	0.448	0.479	0.473	0.352	0.413	0.447	0.445
<i>Mu+Nw</i>	0.577	0.561	0.406	0.475	0.518	0.506	0.338	0.413	0.444	0.454

Table 7: Average F1-score for each model across different training sets for test set MuMo. **c** denotes the number of prior interactions provided as context to the generative models during classification.

Train Set	BERT	RoBERTa	Mixtral				Llama2			
			c=0	c=1	c=3	c=5	c=0	c=1	c=3	c=5
<i>NCTE</i>	0.437	0.488	0.358	0.395	0.401	0.382	0.326	0.351	0.361	0.339

Table 8: Average F1-score for each model across different training sets for NCTE data. **c** denotes the number of prior interactions provided as context to the generative models during classification.

525 The cross-domain experiments involving the Ma-
526 son, Momo, and Newman datasets, labeled with the
527 NCTE labels, achieved decent scores, except for
528 the high uptake variable. This indicates a potential
529 for these models to understand and classify dis-
530 course in educational settings to some extent. How-
531 ever, the experiments relating to NCTE data labeled
532 for the literacy discourse codes did rather poorly,
533 highlighting the difficulties in accurately captur-
534 ing and generalizing discourse patterns within this
535 domain. Given these challenges, we recommend
536 future directions in this area of study to enhance
537 the effectiveness of these models in the field of ed-
538 ucation. Enhancing the collection and annotation

of classroom discourse data across a wider range
of educational settings could improve the represen-
tation within training datasets. This could help the
models learn more generalized features of class-
room discourse that are not specific to any single
domain. Implementing novel cross-domain tech-
niques could help with better transfer learning and
adaptation. Using better architectures and state-
of-the-art (SOTA) models to help generalize the
discourse codes across domains more effectively.

Limitations

In case of generative models, we used only open
weights models and local data processing strictly

539
540
541
542
543
544
545
546
547
548
549
550
551

552 adhering to our data privacy and ethical standards.
 553 While this approach aligns with our ethical stance
 554 and ensures data confidentiality, it also narrows
 555 our selection of computational tools. Potentially
 556 more sophisticated and proprietary models with
 557 higher performance metrics were not considered
 558 in this study due to these constraints. The nature
 559 of our datasets presents another potential limita-
 560 tion. The datasets utilized in our analysis were
 561 shared by the original authors of the work on the
 562 condition that we do not make them publicly avail-
 563 able. This restriction could impose a barrier to
 564 the reproducibility of our study for future research.
 565 Our research primarily concentrates on specific
 566 subject domains like mathematics in English liter-
 567 eracy. These subjects represent only a fraction of
 568 the diverse disciplines within the educational field,
 569 which our current paper does not account for.

570 Ethics Statement

571 The study utilized existing datasets derived from
 572 prior research that were shared with us by the au-
 573 thors of that work. In alignment with our com-
 574 mitment to confidentiality, we have anonymized
 575 all personal information. Names and other iden-
 576 tifying details of students and teachers have been
 577 replaced with pseudonyms, thereby protecting their
 578 identities. Furthermore, the tools and models ap-
 579 plied in our research such as Mixtral and Llama2
 580 7B are open-weights generative models. The de-
 581 cision to use open-weights models supports trans-
 582 parency of our methods and further protects pri-
 583 vacy by eliminating the need to transfer sensitive
 584 data to external servers. The use of open-weights
 585 models can also facilitate reproducibility in the re-
 586 search, allowing other researchers to validate and
 587 build upon the findings in our paper. The primary
 588 goal of the research was to investigate the efficacy
 589 of cross domain classification of educational dis-
 590 course, particularly doctors within the classroom
 591 setting. We recognize the implications of applying
 592 AI in analyzing children’s classroom interactions.
 593 It is important to approach the application of our
 594 research with the understanding of the potential im-
 595 pacts of AI application, making sure that it serves
 596 to enhance the educational experience rather than
 597 compromising it.

598 References

599 Sterling Alic, Dorottya Demszky, Zid Mancenido, Jing
 600 Liu, Heather Hill, and Dan Jurafsky. 2022. Computa-

tionally identifying funneling and focusing questions 601
 in classroom discourse. 602

Simone Balloccu, Patrícia Schmidtová, Mateusz Lango, 603
 and Ondrej Dusek. 2024. Leak, cheat, repeat: Data 604
 contamination and evaluation malpractices in closed- 605
 source LLMs. In *Proceedings of the 18th Confer- 606*
ence of the European Chapter of the Association 607
for Computational Linguistics (Volume 1: Long Pa- 608
pers), pages 67–93, St. Julian’s, Malta. Association 609
 for Computational Linguistics. 610

Jere E Brophy and Thomas L. Good. 1984. *Teacher 611*
behavior and student achievement microform jere 612
brophy and thomas l. good. 613

Hyonsuk Cho and Tanya Christ. 2022. How two emer- 614
 gent bilingual students from refugee families make 615
 inferences with more and less culturally relevant texts 616
 during read-alouds. *TESOL Quarterly*, 56(4):1112– 617
 1135. 618

Tanya Christ, Iman Bakhoda, Ming Ming Chiu, 619
 X. Christine Wang, Alexa Schindel, and Yu Liu. 620
 2023. Mediating learning in the zones of develop- 621
 ment: Role of teacher and kindergartner talk-turns 622
 during read-aloud discussions. *Journal of Research 623*
in Childhood Education, 37(4):519–549. 624

Tanya Christ and Hyonsuk Cho. 2023. Emergent bilin- 625
 gual students’ small group read-aloud discussions. 626
Literacy Research and Instruction, 62(3):203–232. 627

Dorottya Demszky and Heather Hill. 2023. *The NCTE 628*
transcripts: A dataset of elementary math classroom 629
transcripts. In *Proceedings of the 18th Workshop 630*
on Innovative Use of NLP for Building Educational 631
Applications (BEA 2023), pages 528–538, Toronto, 632
 Canada. Association for Computational Linguistics. 633

Dorottya Demszky, Jing Liu, Zid Mancenido, Julie 634
 Cohen, Heather Hill, Dan Jurafsky, and Tatsunori 635
 Hashimoto. 2021. *Measuring conversational uptake: 636*
A case study on student-teacher interactions. 637

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 638
 Kristina Toutanova. 2019. *BERT: Pre-training of 639*
deep bidirectional transformers for language under- 640
standing. In *Proceedings of the 2019 Conference of 641*
the North American Chapter of the Association for 642
Computational Linguistics: Human Language Tech- 643
nologies, Volume 1 (Long and Short Papers), pages 644
 4171–4186, Minneapolis, Minnesota. Association for 645
 Computational Linguistics. 646

Mary Faculty and Sarah Michaels. 1993. *Aligning aca- 647*
ademic task and participation status through revoicing: 648
Analysis of a classroom discourse strategy. *Anthro- 649*
pology & Education Quarterly, 24:318 – 335. 650

Ananya Ganesh, Martha Palmer, and Katharina Kann. 651
 2021. *What would a teacher do? predicting future 652*
talk moves. 653

654	Shannon Giroir, Leticia Grimaldo, Sharon Vaughn, and Greg Roberts. 2015. Interactive read-alouds for english learners in the elementary grades . <i>The Reading Teacher</i> , 68.	712
655		713
656		714
657		715
658	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L�lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth�e Lacroix, and William El Sayed. 2023. Mistral 7b .	716
665	Klaus Krippendorff. 2011. Computing krippendorff’s alpha-reliability .	717
666		718
667	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	719
668		720
669		721
670		722
671		723
672	Abhijit Suresh, Jennifer Jacobs, Charis Harty, Margaret Perkoff, James H. Martin, and Tamara Sumner. 2022a. The TalkMoves dataset: K-12 mathematics lesson transcripts annotated for teacher and student discursive moves . In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 4654–4662, Marseille, France. European Language Resources Association.	724
673		725
674		726
675		727
676		728
677		729
678		730
679		731
680	Abhijit Suresh, Jennifer Jacobs, Margaret Perkoff, James H. Martin, and Tamara Sumner. Fine-tuning transformers with additional context to classify discursive moves in mathematics classrooms . <i>Proceedings of the 17th Workshop on Innovative Use of NLP for Building Educational Applications</i> .	732
681		733
682		734
683		735
684		736
686	Abhijit Suresh, Jennifer Jacobs, Margaret Perkoff, James H. Martin, and Tamara Sumner. 2022b. Fine-tuning transformers with additional context to classify discursive moves in mathematics classrooms . In <i>Proceedings of the 17th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2022)</i> , pages 71–81, Seattle, Washington. Association for Computational Linguistics.	737
687		738
688		739
689		740
690		741
691		742
692		743
693		744
694	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	745
695		746
696		747
697		748
698		749
699		750
700		751
701		752
702		753
703		754
704		755
705		756
706		757
707		758
708		759
709		760
710		761
711		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

765	A.1 Labels and Definitions		
766	1. Student on Task: This label indicates		
767	whether a student's utterance is relevant to		
768	the current topic being discussed in the class-		
769	room.		
770	2. Teacher on Task: This label reflects whether		
771	the teacher's utterance pertains to the topic of		
772	the current classroom session.		
773	3. High Uptake: This label identifies instances		
774	where a speaker (teacher or student) builds		
775	upon what their interlocutor has said, demon-		
776	strating an understanding and extension of the		
777	conversation.		
778	4. Focusing Question: This label is used when a		
779	teacher asks a question that prompts students		
780	to articulate, clarify, or reflect upon their own		
781	thoughts or those of their classmates.		
782	A.2 Labeling Process		
783	1. Student on Task		
784	(a) Label as 1 (On Task): If a student's ut-		
785	terance directly relates to the topic of		
786	the lecture or session. For example, dis-		
787	cussing a specific math problem when		
788	the topic is math. Or if the classroom		
789	session is discussing the NLP textbook,		
790	then the topic would be NLP or anything		
791	related to it.		
792	(b) Label as 0 (Off Task): If a student's ut-		
793	terance is unrelated to the topic of the lec-		
794	ture. Such as talking about the weather		
795	or making a joke unrelated to the topic at		
796	hand.		
797	2. Teacher on Task		
798	(a) Label as 1 (On Task): If the teacher's		
799	utterance is directly related to the subject		
800	matter of the current session similar to		
801	student on task label.		
802	(b) Label as 0 (Off Task): If the teacher's		
803	utterance is not related to the topic of the		
804	session.		
805	3. High Uptake		
806	(a) Label as 1 (High Uptake): When a		
807	teacher acknowledges, repeats, or re-		
808	formulates what the student has said,		
809	thereby extending the conversation.		
	(b) Label as 0 (Low Uptake): When the re-	810	
	sponse does not build upon the previous	811	
	speaker's (student's) contribution.	812	
	4. Focusing Question	813	
	(a) Label as 1 (Focusing Question): If a	814	
	teacher's question prompts the student to	815	
	think deeply, articulate their understand-	816	
	ing, or engage in reflection about their	817	
	own thoughts or those of other students.	818	
	(b) Label as 0 (Funneling Question): If the	819	
	teacher's question or teacher's set of	820	
	questions to lead students to a desired	821	
	procedure or conclusion, while giving	822	
	limited attention to student responses	823	
	that veer from the desired path	824	
	A.3 Examples for Each Category	825	
	1. Student on Task / Teacher on Task	826	
	(a) Example (Label 1): Topic is English	827	
	textbook "My friend Jamal". S: "It is	828	
	because Jamal was a friend of Joseph	829	
	and they lived nearby." T: "yes! They	830	
	were friends and what does that mean for	831	
	Joseph?"	832	
	(b) Example (Label 0): Topic is English text-	833	
	book "My friend Jamal". S: "I played	834	
	soccer yesterday." T: "Shhh! Sit down	835	
	quietly. We have 15 minutes left."	836	
	2. High Uptake	837	
	(a) Example (Label 1): S: "Cause you took	838	
	away 10 and 70 minus 10 is 60". T:	839	
	"Why did we take away 10?"	840	
	(b) Example (Label 1): S: "There's not	841	
	enough seeds". T: "There's not enough	842	
	seeds. How do you know right away that	843	
	128 or 132 or whatever it was you got	844	
	doesn't make sense?"	845	
	(c) Example (Label 0): S: "Because the base	846	
	of it is a hexagon". T: "Student K?"	847	
	3. Focusing Question	848	
	(a) Example (Label 1): S: "I disagree with	849	
	Student A because if you skip count by	850	
	100 ten times, that will get you to 1,000".	851	
	T: "Let's try it. You ready? Let's start	852	
	right here with Student F". S: "A hun-	853	
	dred."	854	

- 855 (b) Example (Label 1): S: I first got 32 and
856 then I got 48. T: And how did you find
857 that? S: “Because I did 16 times two is
858 32”.
- 859 (c) Example (Label 0): S: “Do we eat pizza
860 today”. T: “Student K? What are you
861 doing there???”.

862 B Annotation Guidelines to label NCTE 863 dataset

864 These guidelines are designed to assist annotators
865 in labeling the classroom interactions between stu-
866 dents and the teacher based on the categories de-
867 fined in the research conducted by Christ, T., et al
868 (2022). Note that this is a multi-label classification
869 task and each individual interaction can have one
870 or more possible output labels.

871 B.1 Labels and Definitions

- 872 1. **Response Evaluation:** When the teacher ei-
873 ther compliments a child or expresses uncer-
874 tainty about an incorrect response.
- 875 2. **Providing Information:** Extending or elabo-
876 rating what was said either by the teacher or
877 the student, building background knowledge,
878 defining, using target words that are utilized
879 in the context.
- 880 3. **Misinformation:** Either by providing misin-
881 formation or verifying an incorrect response.
- 882 4. **Revoicing:** When the teacher acknowledges
883 and repeats what the student has said earlier.
- 884 5. **Strategy related:** Teacher directs a child to
885 look at or think about text clues, or asks chil-
886 dren to check their prediction.
- 887 6. **Questioning:** When a teacher questions a
888 child to get a more detailed response, or elicit
889 noticing text clues, or to define a target vocabu-
890 lary etc.
- 891 7. **Behavior Management:** Gives children or a
892 particular child a behavioral directive.
- 893 8. **Turn Management:** Teacher calls on particu-
894 lar child to respond or acknowledges or rejects
895 a child’s initiative to talk.

B.2 Labeling Process 896

To annotate this dataset: 897

- 898 (a) We read a transcript, and identify the pos-
899 sible codes that apply to that utterance
900 using the codebook provided in the orig-
901 inal paper.
- 902 (b) Look up the category that those particu-
903 lar codes fall under, and label either 1 or
904 0 on the spreadsheet. NOTE: The idea is
905 to map the codes back to their categories
906 and use them as labels instead.
- 907 (c) For example, when the teacher says “He
908 is mowing the grass. Good!! What will
909 the mower do to the flower if the dad
910 gets closer? What would the mower
911 do? Student K??”, the authors of the
912 paper identified that the teacher repeated
913 the student’s response. Then proceeded
914 to compliment the child, acknowledging
915 that the child has given the correct an-
916 swer. Then proceeds to ask a question
917 while directing that question to a particu-
918 lar child. Therefore we ended up with 4
919 possible codes for that one teacher utter-
920 ance. Now we map those codes back to
921 their categories.

C Model Hyperparameters 922

923 For the transformer-based deep learning models,
924 we initialize each from the model checkpoint and
925 fine-tune on our training data for 5 epochs with
926 a batch size of 16, weight decay of 0.01, and a
927 learning rate of $2e-5$.

928 Details of the hyperparameter tuning for the gener-
929 ative models, Mixtral and Llama2:

- 930 1. **Do sample:** Set to false, this parameter en-
931 sures deterministic outputs by selecting tokens
932 based on their probability distribution rather
933 than introducing variability. This aligns with
934 the experiment’s objective of restricting out-
935 puts to only "yes" and "no" tokens.
- 936 2. **Max new tokens:** With a value of 1, this pa-
937 rameter confines the model to generate exactly
938 one token after the input prompt. Given our
939 experiment’s focus on producing either "yes"
940 or "no," a value of one facilitates the desired
941 output format of one token per response.
- 942 3. **Temperature:** Set to 0, indicating no ran-
943 domness in output selection. By eliminating

944 randomness, the model consistently chooses
945 the same sequence of tokens from the input
946 prompt, thereby ensuring deterministic out-
947 put.

- 948 4. **Top k:** Set to 2, this parameter limits consid-
949 eration to the top two tokens with the highest
950 probabilities. Since the objective of this exper-
951 iment is binary output ("yes" or "no"), setting
952 Top k to 2 effectively restricts the model’s
953 outputs to these two options.

- 954 5. **Num_return_sequence:** set to 1

955 D Dataset Statistics

956 In order to evaluate the generalization performance
957 of the models, we annotated new data as described
958 in subsection 3.2. The number of datapoints as-
959 signed each label are presented in Tables 9 and
960 10.

Variable	Class 0	Class 1
Student on Task	21	119
Teacher on Task	13	127
High Uptake	74	66
Focusing Question	83	57

Table 9: Class distribution of MuMo/Mason/Newman data annotated with NCTE labels. Class 0 indicates the label does not apply and Class 1 indicates that it does.

Variable	Class 0	Class 1
Response Evaluation	55	45
Providing Information	58	42
Revoicing	74	26
Strategy Related	69	31
Questioning Behavior	32	68
Management	90	10
Turn Management	70	30

Table 10: Class distribution of NCTE data when annotated with MuMo/Mason/Newman label set. Class 0 indicates the label does not apply and Class 1 indicates that it does.

961 E Experimental Details

Train Set	Models	Questioning	Response Evaluation	Providing Information	Revoicing	Strategy Related	Behavior Management	Turn Management
MuMo	<i>Baseline</i>	0.241	0.414	0.151	0.178	0.230	0.194	0.198
	BERT	0.641	0.331	0.320	0.444	0.233	0.323	0.269
	RoBERTa	0.612	0.266	0.319	0.415	0.279	0.303	0.253
	Mixtral	0.450	0.285	0.294	0.371	0.320	0.307	0.250
	Mixtral (c=1)	0.462	0.270	0.299	0.364	0.303	0.310	0.248
	Mixtral (c=3)	0.525	0.299	0.316	0.386	0.288	0.313	0.252
	Mixtral (c=5)	0.537	0.283	0.333	0.405	0.284	0.322	0.255
	Llama2	0.438	0.276	0.239	0.310	0.244	0.279	0.259
	Llama2 (c=1)	0.445	0.240	0.240	0.313	0.241	0.290	0.250
	Llama2 (c=3)	0.484	0.279	0.282	0.333	0.274	0.300	0.251
Llama2 (c=5)	0.487	0.270	0.311	0.347	0.266	0.304	0.251	
Mason	BERT	0.699	0.460	0.324	0.538	0.333	0.460	0.458
	RoBERTa	0.699	0.422	0.315	0.530	0.329	0.459	0.457
	Mixtral	0.535	0.340	0.330	0.397	0.332	0.350	0.420
	Mixtral (c=1)	0.601	0.360	0.315	0.455	0.327	0.365	0.421
	Mixtral (c=3)	0.637	0.407	0.315	0.473	0.331	0.384	0.443
	Mixtral (c=5)	0.644	0.436	0.318	0.470	0.333	0.401	0.444
	Llama2	0.488	0.281	0.284	0.395	0.275	0.334	0.351
	Llama2 (c=1)	0.503	0.312	0.281	0.423	0.292	0.365	0.369
	Llama2 (c=3)	0.585	0.303	0.293	0.488	0.307	0.372	0.397
Llama2 (c=5)	0.599	0.326	0.305	0.544	0.330	0.351	0.405	
Newman	BERT	0.710	0.489	0.339	0.535	0.396	0.473	0.510
	RoBERTa	0.702	0.460	0.325	0.530	0.383	0.472	0.497
	Mixtral	0.472	0.384	0.349	0.420	0.313	0.370	0.400
	Mixtral (c=1)	0.510	0.412	0.352	0.429	0.327	0.384	0.428
	Mixtral (c=3)	0.666	0.440	0.343	0.481	0.367	0.386	0.459
	Mixtral (c=5)	0.669	0.490	0.336	0.493	0.365	0.368	0.450
	Llama2	0.444	0.340	0.287	0.359	0.304	0.344	0.401
	Llama2 (c=1)	0.592	0.365	0.295	0.382	0.324	0.359	0.403
	Llama2 (c=3)	0.594	0.425	0.316	0.436	0.315	0.387	0.418
Llama2 (c=5)	0.571	0.484	0.344	0.443	0.339	0.368	0.435	
Ms+Nw	BERT	0.689	0.479	0.353	0.530	0.388	0.462	0.495
	RoBERTa	0.676	0.472	0.353	0.492	0.364	0.444	0.453
	Mixtral	0.478	0.367	0.311	0.400	0.345	0.369	0.381
	Mixtral (c=1)	0.614	0.384	0.325	0.427	0.353	0.392	0.409
	Mixtral (c=3)	0.617	0.427	0.335	0.438	0.370	0.429	0.403
	Mixtral (c=5)	0.582	0.401	0.352	0.427	0.374	0.457	0.419
	Llama2	0.413	0.352	0.279	0.333	0.326	0.388	0.325
	Llama2 (c=1)	0.497	0.336	0.291	0.369	0.344	0.402	0.334
	Llama2 (c=3)	0.500	0.361	0.330	0.437	0.360	0.440	0.353
Llama2 (c=5)	0.490	0.379	0.339	0.428	0.341	0.460	0.333	
Mu+Ms	BERT	0.666	0.457	0.350	0.466	0.378	0.446	0.494
	RoBERTa	0.641	0.446	0.327	0.428	0.354	0.429	0.507
	Mixtral	0.430	0.360	0.241	0.279	0.292	0.314	0.377
	Mixtral (c=1)	0.473	0.384	0.252	0.315	0.315	0.317	0.401
	Mixtral (c=3)	0.577	0.412	0.304	0.379	0.345	0.344	0.402
	Mixtral (c=5)	0.654	0.449	0.300	0.353	0.373	0.354	0.396
	Llama2	0.414	0.345	0.248	0.245	0.250	0.279	0.383
	Llama2 (c=1)	0.479	0.349	0.270	0.286	0.271	0.307	0.401
	Llama2 (c=3)	0.567	0.404	0.308	0.379	0.327	0.336	0.444
Llama2 (c=5)	0.562	0.392	0.299	0.369	0.375	0.343	0.419	
Mu+Nw	BERT	0.691	0.463	0.347	0.419	0.396	0.430	0.500
	RoBERTa	0.676	0.460	0.324	0.412	0.376	0.423	0.471
	Mixtral	0.468	0.301	0.350	0.344	0.319	0.355	0.306
	Mixtral (c=1)	0.504	0.333	0.356	0.366	0.329	0.366	0.342
	Mixtral (c=3)	0.612	0.404	0.353	0.387	0.360	0.405	0.346
	Mixtral (c=5)	0.694	0.472	0.349	0.411	0.401	0.427	0.343
	Llama2	0.384	0.300	0.275	0.286	0.325	0.287	0.295
	Llama2 (c=1)	0.453	0.328	0.285	0.308	0.340	0.313	0.337
	Llama2 (c=3)	0.573	0.393	0.320	0.362	0.367	0.379	0.324
Llama2 (c=5)	0.680	0.462	0.344	0.411	0.400	0.432	0.310	

Table 11: F1-score when training on various train sets and evaluating on the test set from *Newman*. **Bold** indicates the best score for each column for each training set, underline indicates the best overall score for each column.

Train Set	Model	Questioning	Response Evaluation	Providing Information	Revoicing	Strategy Related	Behavior Management	Turn Management
	<i>Baseline</i>	<i>0.706</i>	<i>0.308</i>	<i>0.625</i>	<i>0.625</i>	<i>0.308</i>	<i>0.533</i>	<i>0.308</i>
	BERT	0.821	0.480	0.525	0.000	0.333	0.000	0.800
NCTE	RoBERTa	0.7724	0.666	0.649	0.000	0.495	0.000	0.813
	Mixtral	0.692	0.389	0.363	0.241	0.337	0.000	0.525
	Llama2	0.614	0.321	0.381	0.219	0.238	0.190	0.316
	Mixtral c=1	0.686	0.359	0.370	0.258	0.322	0.200	0.569
	Mixtral c=3	0.721	0.378	0.365	0.255	0.317	0.214	0.555
	Mixtral c=5	0.714	0.377	0.333	0.263	0.309	0.179	0.565
	Llama c=1	0.604	0.338	0.383	0.222	0.281	0.222	0.407
	Llama c=3	0.628	0.334	0.385	0.246	0.325	0.222	0.411
	Llama c=5	0.624	0.313	0.342	0.210	0.287	0.213	0.405

Table 12: Generalization performance on NCTE data labeled with MuMo, Mason and Newman dataset labels.

Train Set	Model	Student on Task	Teacher on Task	High Uptake	Focusing Question
	<i>Baseline</i>	<i>1.000</i>	<i>1.000</i>	<i>0.929</i>	<i>0.636</i>
	BERT	0.962	0.800	0.333	0.694
	RoBERTa	0.941	0.785	0.369	0.656
	Mixtral	0.784	0.601	0.303	0.666
	Mixtral (c=1)	0.740	0.600	0.300	0.661
	Mixtral (c=3)	0.767	0.637	0.297	0.628
	Mixtral (c=5)	0.749	0.615	0.270	0.612
	Llama2	0.678	0.610	0.263	0.537
	Mixtral (c=1)	0.684	0.580	0.284	0.550
	Mixtral (c=3)	0.647	0.613	0.289	0.523
	Mixtral (c=5)	0.666	0.600	0.251	0.501

Table 13: Generalization performance on subset of MuMo Data labeled using NCTE’s labels.

Train Set	Models	Questioning	Response Evaluation	Providing Information	Revoicing	Strategy Related	Behavior Management	Turn Management
<i>MuMo</i>	<i>Baseline</i>	0.470	0.230	0.200	0.364	0.105	0.036	0.364
	BERT	0.904	0.375	0.564	0.522	0.362	0.533	0.737
	RoBERTa	0.871	0.344	0.555	0.501	0.370	0.555	0.747
	Mixtral	0.790	0.318	0.231	0.444	0.333	0.378	0.478
	Mixtral (c=1)	0.881	0.320	0.324	0.451	0.333	0.415	0.538
	Mixtral (c=3)	0.888	0.337	0.362	0.484	0.345	0.467	0.594
	Mixtral (c=5)	0.865	0.330	0.363	0.476	0.358	0.525	0.575
	Llama2	0.643	0.275	0.286	0.310	0.219	0.344	0.404
	Llama2 (c=1)	0.726	0.294	0.323	0.371	0.252	0.387	0.480
	Llama2 (c=3)	0.810	0.297	0.404	0.375	0.303	0.465	0.511
	Llama2 (c=5)	0.807	0.280	0.417	0.380	0.306	0.541	0.518
Mason	BERT	0.722	0.345	0.542	0.500	0.500	0.557	0.688
	RoBERTa	0.718	0.333	0.514	0.500	0.500	0.548	0.680
	Mixtral	0.523	0.327	0.289	0.346	0.334	0.694	0.675
	Mixtral (c=1)	0.657	0.335	0.439	0.451	0.402	0.690	0.689
	Mixtral (c=3)	0.685	0.335	0.447	0.447	0.430	0.734	0.680
	Mixtral (c=5)	0.636	0.347	0.452	0.492	0.439	0.697	0.701
	Llama2	0.551	0.287	0.282	0.322	0.321	0.474	0.595
	Llama2 (c=1)	0.580	0.294	0.353	0.381	0.354	0.493	0.614
	Llama2 (c=3)	0.651	0.323	0.434	0.422	0.420	0.533	0.652
	Llama2 (c=5)	0.658	0.348	0.436	0.417	0.419	0.521	0.653
Newman	BERT	0.741	0.460	0.461	0.470	0.330	0.500	0.595
	RoBERTa	0.742	0.440	0.562	0.476	0.388	0.528	0.555
	Mixtral	0.464	0.349	0.387	0.341	0.307	0.444	0.463
	Mixtral (c=1)	0.696	0.376	0.405	0.364	0.309	0.454	0.487
	Mixtral (c=3)	0.695	0.409	0.430	0.426	0.323	0.485	0.545
	Mixtral (c=5)	0.735	0.461	0.452	0.466	0.329	0.502	0.599
	Llama2	0.400	0.350	0.378	0.330	0.240	0.295	0.301
	Llama2 (c=1)	0.568	0.371	0.390	0.358	0.262	0.340	0.355
	Llama2 (c=3)	0.655	0.416	0.430	0.411	0.298	0.422	0.481
	Llama2 (c=5)	0.663	0.454	0.452	0.465	0.330	0.507	0.585
Ms+Nw	BERT	0.620	0.333	0.499	0.458	0.333	0.430	0.467
	RoBERTa	0.611	0.333	0.504	0.443	0.327	0.442	0.476
	Mixtral	0.525	0.298	0.395	0.354	0.311	0.380	0.294
	Mixtral (c=1)	0.542	0.307	0.411	0.378	0.315	0.386	0.331
	Mixtral (c=3)	0.570	0.322	0.462	0.410	0.323	0.405	0.405
	Mixtral (c=5)	0.620	0.327	0.494	0.464	0.332	0.435	0.470
	Llama2	0.485	0.290	0.371	0.350	0.279	0.320	0.344
	Llama2 (c=1)	0.522	0.300	0.404	0.377	0.291	0.345	0.369
	Llama2 (c=3)	0.573	0.316	0.455	0.419	0.311	0.382	0.414
	Llama2 (c=5)	0.623	0.337	0.499	0.460	0.330	0.423	0.463
Mu+Ms	BERT	0.840	0.369	0.542	0.492	0.351	0.511	0.651
	RoBERTa	0.822	0.338	0.530	0.464	0.351	0.500	0.666
	Mixtral	0.622	0.295	0.389	0.400	0.317	0.381	0.471
	Mixtral (c=1)	0.776	0.308	0.420	0.412	0.320	0.450	0.571
	Mixtral (c=3)	0.750	0.378	0.422	0.474	0.333	0.458	0.588
	Mixtral (c=5)	0.743	0.342	0.450	0.457	0.309	0.501	0.547
	Llama2	0.595	0.244	0.352	0.373	0.308	0.331	0.386
	Llama2 (c=1)	0.651	0.273	0.397	0.396	0.330	0.408	0.441
	Llama2 (c=3)	0.734	0.295	0.465	0.436	0.332	0.442	0.443
	Llama2 (c=5)	0.729	0.269	0.434	0.414	0.341	0.469	0.413
Mu+Nw	BERT	0.832	0.364	0.518	0.490	0.323	0.509	0.668
	RoBERTa	0.833	0.349	0.516	0.489	0.330	0.510	0.640
	Mixtral	0.643	0.313	0.334	0.387	0.307	0.363	0.444
	Mixtral (c=1)	0.858	0.325	0.422	0.447	0.315	0.430	0.563
	Mixtral (c=3)	0.801	0.347	0.442	0.493	0.321	0.515	0.609
	Mixtral (c=5)	0.786	0.345	0.428	0.487	0.325	0.486	0.622
	Llama2	0.594	0.238	0.248	0.278	0.231	0.321	0.367
	Llama2 (c=1)	0.748	0.268	0.298	0.374	0.246	0.353	0.430
	Llama2 (c=3)	0.739	0.360	0.413	0.410	0.282	0.436	0.473
	Llama2 (c=5)	0.819	0.326	0.408	0.402	0.317	0.405	0.573

Table 14: F1-score when training on various train sets and evaluating on the test set from *MuMo*. **Bold** indicates the best score for each column for each training set, underline indicates the best overall score for each column.

Train Set	Models	Questioning	Response Evaluation	Providing Information	Revoicing	Strategy Related	Behavior Management	Turn Management
MuMo	Baseline	0.420	0.120	0.033	0.275	0.160	0.065	0.128
	BERT	0.644	0.305	0.327	0.458	0.255	0.341	0.294
	RoBERTa	0.638	0.284	0.306	0.443	0.306	0.322	0.273
	Mixtral	0.472	0.295	0.277	0.380	0.303	0.288	0.269
	Mixtral (c=1)	0.556	0.303	0.289	0.389	0.296	0.300	0.357
	Mixtral (c=3)	0.581	0.344	0.308	0.420	0.280	0.323	0.370
	Mixtral (c=5)	0.552	0.303	0.327	0.455	0.259	0.345	0.303
	Llama2	0.444	0.300	0.264	0.339	0.263	0.306	0.279
	Llama2 (c=1)	0.493	0.303	0.280	0.366	0.261	0.316	0.310
	Llama2 (c=3)	0.477	0.327	0.302	0.450	0.258	0.330	0.324
Llama2 (c=5)	0.453	0.309	0.322	0.466	0.252	0.341	0.295	
Mason	BERT	0.729	0.462	0.334	0.610	0.365	0.480	0.501
	RoBERTa	0.700	0.445	0.337	0.608	0.325	0.434	0.478
	Mixtral	0.560	0.380	0.309	0.399	0.316	0.367	0.400
	Mixtral (c=1)	0.704	0.389	0.318	0.436	0.322	0.394	0.418
	Mixtral (c=3)	0.748	0.433	0.321	0.515	0.352	0.428	0.465
	Mixtral (c=5)	0.727	0.467	0.337	0.614	0.365	0.486	0.493
	Llama2	0.562	0.291	0.275	0.401	0.269	0.318	0.383
	Llama2 (c=1)	0.689	0.328	0.285	0.444	0.288	0.351	0.405
	Llama2 (c=3)	0.669	0.391	0.313	0.522	0.323	0.409	0.448
	Llama2 (c=5)	0.718	0.463	0.338	0.602	0.361	0.480	0.506
Newman	BERT	0.711	0.447	0.339	0.541	0.328	0.449	0.495
	RoBERTa	0.693	0.444	0.309	0.517	0.326	0.461	0.512
	Mixtral	0.500	0.422	0.300	0.381	0.279	0.344	0.425
	Mixtral (c=1)	0.643	0.427	0.305	0.406	0.285	0.361	0.433
	Mixtral (c=3)	0.714	0.438	0.320	0.486	0.309	0.410	0.460
	Mixtral (c=5)	0.721	0.447	0.344	0.543	0.328	0.448	0.489
	Llama2	0.469	0.375	0.238	0.366	0.278	0.343	0.366
	Llama2 (c=1)	0.608	0.396	0.258	0.406	0.287	0.358	0.394
	Llama2 (c=3)	0.633	0.412	0.299	0.467	0.314	0.400	0.448
	Llama2 (c=5)	0.614	0.452	0.335	0.541	0.327	0.448	0.502
Ms+Nw	BERT	0.716	0.434	0.329	0.563	0.348	0.444	0.478
	RoBERTa	0.707	0.432	0.313	0.520	0.341	0.400	0.463
	Mixtral	0.475	0.353	0.279	0.342	0.268	0.332	0.384
	Mixtral (c=1)	0.633	0.366	0.290	0.388	0.287	0.361	0.438
	Mixtral (c=3)	0.692	0.404	0.305	0.468	0.311	0.391	0.440
	Mixtral (c=5)	0.719	0.400	0.325	0.568	0.351	0.443	0.482
	Llama2	0.527	0.301	0.233	0.344	0.287	0.300	0.348
	Llama2 (c=1)	0.654	0.341	0.253	0.390	0.304	0.328	0.375
	Llama2 (c=3)	0.718	0.337	0.294	0.467	0.319	0.383	0.430
	Llama2 (c=5)	0.680	0.326	0.324	0.567	0.351	0.452	0.448
Mu+Ms	BERT	0.703	0.440	0.323	0.587	0.349	0.444	0.455
	RoBERTa	0.694	0.438	0.319	0.560	0.347	0.429	0.454
	Mixtral	0.512	0.373	0.304	0.331	0.296	0.259	0.382
	Mixtral (c=1)	0.650	0.380	0.302	0.481	0.304	0.293	0.390
	Mixtral (c=3)	0.698	0.419	0.316	0.493	0.334	0.273	0.422
	Mixtral (c=5)	0.630	0.435	0.323	0.501	0.355	0.338	0.407
	Llama2	0.499	0.334	0.292	0.306	0.297	0.251	0.340
	Llama2 (c=1)	0.593	0.355	0.293	0.427	0.305	0.292	0.368
	Llama2 (c=3)	0.617	0.401	0.308	0.453	0.332	0.368	0.408
	Llama2 (c=5)	0.628	0.437	0.325	0.445	0.349	0.353	0.388
Mu+Nw	BERT	0.683	0.436	0.310	0.540	0.316	0.405	0.444
	RoBERTa	0.665	0.428	0.307	0.527	0.300	0.421	0.421
	Mixtral	0.489	0.366	0.264	0.460	0.247	0.301	0.399
	Mixtral (c=1)	0.569	0.379	0.270	0.473	0.263	0.323	0.412
	Mixtral (c=3)	0.598	0.409	0.297	0.547	0.268	0.380	0.434
	Mixtral (c=5)	0.616	0.401	0.267	0.521	0.261	0.354	0.443
	Llama2	0.420	0.341	0.251	0.414	0.233	0.258	0.384
	Llama2 (c=1)	0.540	0.359	0.259	0.445	0.250	0.290	0.396
	Llama2 (c=3)	0.573	0.403	0.292	0.485	0.283	0.345	0.427
	Llama2 (c=5)	0.543	0.403	0.305	0.535	0.249	0.400	0.412

Table 15: F1-score when training on various train sets and evaluating on the test set from *Mason*. **Bold** indicates the best score for each column for each training set, underline indicates the best overall score for each column.