# From Words to Wisdom: Discourse Annotation and Baseline Models for Student Dialogue Understanding

Farjana Sultana Mim, Shuchin Aeron, Eric Miller and Kristen Wendell

*Abstract*—Identifying discourse features in student conversations is quite important for educational researchers to recognize the curricular and pedagogical variables that cause students to engage in constructing knowledge rather than merely completing tasks. The manual analysis of student conversations to identify these discourse features is time-consuming and labor-intensive, which limits the scale and scope of studies. Leveraging natural language processing (NLP) techniques can facilitate the automatic detection of these discourse features, offering educational researchers scalable and data-driven insights. However, existing studies in NLP that focus on discourse in dialogue rarely address educational data. In this work, we address this gap by introducing an annotated educational dialogue dataset of student conversations featuring knowledge construction and task production discourse. We also establish baseline models for automatically predicting these discourse properties for each turn of talk within conversations, using pre-trained large language models GPT-3.5 and Llama-3.1. Experimental results indicate that these state-of-the-art models perform suboptimally on this task, indicating the potential for future research.

*Index Terms*—Natural Language Processing, Large Language Model, Discourse, Dialogue, Education.

## I. INTRODUCTION

RESEARCH in classroom settings has shown that student learning outcomes are higher when students frame a classwork or homework activity as an opportunity for constructing knowledge rather than as a task to be produced for the instructor [1], [2]. In other words, two important features of student conversations are: *knowledge construction (KC)* discourse, which refers to the student talks focused on developing conceptual understanding, and *task production (TP)* discourse, where student talks are focused on completing an instructional task as expediently as possible [3].

Prior research in learning sciences has also demonstrated that when students frame their purpose within an instructional activity as constructing knowledge rather than just completing a task, they are more likely to develop expertise and be able to later transfer their expertise to new situations [4]. These

F. S. Mim is with the Department of Electrical and Computer Engineering, Tufts University, Medford, MA 02155, United States. (e-mail: farjana.mim59@gmail.com) (currently in the Department of Computer Science and Information Technology, Patuakhali Science and Technology University, Bangladesh)

S. Aeron is with the Department of Electrical and Computer Engineering, Tufts University, Medford, MA 02155, United States. (e-mail: shuchin.aeron@tufts.edu)

E. Miller is with the Department of Electrical and Computer Engineering, Computer Science and Biomedical Engineering, Tufts University, Medford, MA 02155, United States. (e-mail: eric.miller@tufts.edu)

K. Wendell is with the Department of Mechanical Engineering and Education, Tufts University, Medford, MA 02155, United States. (e-mail: kristen.wendell@tufts.edu)

**Homework Topic**

Design an experiment complete with instrumentation to determine the specific heats of a gas using a resistance heater. Discuss how the experiment will be conducted, what measurements need to be taken, and how the specific heats will be determined. What are the sources of error in your system? How can you minimize the experimental error?

**Task Production Discourse**

*Student X*: Although we just have to design the experiment. It's not like we have to actually do it.
*Student T*: No.
*Student A*: Just design and justify this will work.
*Student X*: How can you minimize the experimental error. That's one of the points there.

**Knowledge Construction Discourse**

*Student X*: Ok. So one thought I had too was that actually um whatever material the container is made out of when it heats up, it's going to expand -
*Student T*: Mhm.
*Student X*: - and that will change whatever the internal volume is. And I don't know if it makes it bigger or smaller actually. It um might make it bigger but if there were -

Fig. 1: Students' homework discussion's snippet of knowledge construction and task production discourse.

results have been found across several disciplines, including physics, chemistry, biology, and engineering education [1]–[3], [5]. However, the relationship between knowledge construction discourse and learning outcomes has yet to be translated into actionable principles for pedagogy and curriculum design. The major difficulty lies in pinpointing which particular aspects of the learning environment and instructional activity cue students into knowledge constructing discourse.

To address this gap, we aim to develop efficient methods for distinguishing students' knowledge construction discourse from their task production discourse so that researchers can more broadly investigate the conditions or contexts under which students tend to adopt a knowledge construction framing. Such findings would enable educators to design learning experiences and environments so that they cue students toward constructing knowledge.

Fig 1 shows examples of *knowledge construction* and *task production* discourse in an undergraduate engineering students' conversation. In the task production discourse of the example, the students remind each other that their homework task is to design an experiment and describe how they would minimize experimental error. These lines are focused on setting up the steps to complete their homework. In the

knowledge construction discourse from the same homework conversation, *student X* shares an idea about how the process of heating a gas will affect the material containing it. Rather than simply completing a pre-determined step of the homework, student X tries to envision the phenomena that will occur in the experiment the students are designing. At this moment, X's turn of talk is oriented toward understanding rather than expediency.

Traditional manual analysis of student dialogues to identify these discourse features is time-intensive, which limits the scope of studies. Leveraging natural language processing (NLP) techniques can facilitate the automatic detection of KC and TP discourse, providing educational researchers with valuable insights into how curricular and pedagogical variables influence students to engage in knowledge construction rather only task production.

Discourse in dialogue or conversations has been widely studied in NLP in different task settings such as dialogue act classification [6]–[9], dialogue topic segmentation and categorization [10]–[13], dialogue state tracking [14]–[19], and identifying dialogue system behaviors [20], [21]. However, although various discourse frameworks are being applied to different types of conversational data, hardly any of them are educational data [22], [23]. To address this gap, this study creates a novel educational dialogue dataset, annotated with knowledge construction (KC) and task production (TP) discourse[1]. We also formulate the NLP task of KCTP (Knowledge Construction and Task Production) prediction, aiming to automatically identify these discourse types within educational dialogues.

Lately, the NLP field has been revolutionized by pre-trained large language models (LLMs) such as GPT-3 [24], Llama [25], Gemini [26], Deepseek [27]. These models have demonstrated significant performance gains and yielded interesting findings across various NLP tasks, including the study of discourse in dialogues or conversations [14], [20]. Recently, a new paradigm called *"Pre-train, Prompt, and Predict"* [28] has gained popularity which leverages pre-trained LLMs through natural language prompts instead of fine-tuning them for specific tasks. By using such *"prompting"* method, one can probe task-specific knowledge from LLMs, which has shown remarkable performance in various tasks such as text classification and summarization [29], [30]. Another paradigm called *"instruction fine-tuning"* [31] which finetunes a model on a dataset via instructions, has significantly improved the performance of several tasks [32]. Therefore, we use GPT-3.5 with prompting techniques to establish a baseline for our Knowledge Construction vs. Task Production (KCTP) prediction task. However, as GPT-3.5 is not an open-source model, we also use the open-access LLaMA-3.1 (8B) model [33] and fine-tune it for the same task. Experimental results indicate that prompting and fine-tuning GPT-3.5 and LLaMA-3.1 yield suboptimal performance on KCTP prediction, suggesting the need for further research into models and methods better suited

---

[1]The human subjects protocol under which the data were generated does not allow for its public sharing. Readers interested in the data set may contact the authors for further information.

to educational discourse analysis. To summarize, the main contributions of this work are as follows:

- We create a novel educational dialogue dataset annotated with *Knowledge Construction (KC)* and *Task Production (TP)* discourse, addressing a gap in discourse-annotated educational data.
- We formulate the *Knowledge Construction vs. Task Production (KCTP)* classification as a natural language processing (NLP) task to automatically identify *KC* and *TP* discourse in student dialogues.
- We establish baseline models for the *KCTP* prediction task using GPT-3.5 and LLaMA-3.1 prompting as well as LLaMA-3.1 instruction fine-tuning, revealing current limitations of LLMs in modeling educational discourse and highlighting directions for future research.

## II. RELATED WORK

This study develops an educational dialogue dataset annotated with instances of *knowledge construction (KC)* and *task production (TP)* discourse. We also establish baseline models for the automatic prediction of KC and TP discourse, with the goal of enabling educational researchers to identify the curricular and pedagogical conditions that encourage students to engage in constructing knowledge rather than merely completing tasks. In this section, we briefly review prior work in three relevant areas: (1) discourse in learning sciences, (2) discourse analysis in Dialogue using NLP, and (3) use of pre-trained language models for discourse modeling in dialogue.

### A. Learning Sciences Approach to Educational Discourse Analysis

Discourse has been long studied in learning sciences to determine the nature of activity, understanding, and learning styles of students [1]–[5]. Gouvea et al. [1] presented a case study of a life-science major in a reformed physics course, showing how epistemological shifts in one discipline can transfer to another. Over a year, the student moved from rote learning to coherence-seeking reasoning in physics, integrating materials, peer discussion, and feedback. This reframing extended to biology, where the student began approaching the subject more conceptually. The study provides qualitative evidence that discourse-centered instructional strategies can foster cross-disciplinary epistemological development.

In another work, Scherr and Hammer [5] explored how students' collaborative behaviors such as posture, gaze, gestures, and vocal dynamics serve as observable indicators of their epistemological framing during active-learning physics activities. They analyze video recordings from introductory physics tutorial sessions and identify distinct behavioral clusters corresponding to different ways students frame the task: for instance, working through substance-based sensemaking versus perceiving it as a procedural worksheet exercise. The authors demonstrate that when students frame the activity as sensemaking, their behaviors align with deeper conceptual reasoning and engagement in discussing the substance of ideas. Their findings highlight the dynamic interplay between

observable behavior, framing, and the quality of students' scientific reasoning in small-group learning contexts.

Koretsky et al. [2] examined how the design of engineering tasks and instructional framing influence student team dynamics, balancing action ("doing") and reflection ("thinking") during collaborative open-ended projects. Through detailed cases of small-group engineering design work, they show that when tasks are meaningful, realistic, and properly scaffolded, teams display more equitable participation, distributed modeling and communication, and deeper conceptual reasoning rather than surface-level task execution alone. In particular, the interplay between material engagement (e.g., prototyping and sketching) and explicit discourse about design decisions fosters collective sense-making and shared agency. The study highlights how thoughtfully structured activities and facilitative framing can empower teams to engage in both productive action and epistemic dialogue, offering important implications for discourse-centric analyses and NLP applications in educational dialogue modeling.

### B. Discourse Analysis in Dialogue using NLP

Discourse in dialogue has been extensively studied in natural language processing (NLP) [34]–[36]. Raheja and Tetreault [7] proposed a hierarchical recurrent neural network and coupled it with a context-aware self-attention mechanism to model different levels of utterance and dialogue act semantics, achieving state-of-the-art performance on the Switchboard Dialogue Act Corpus. Liu et al. [10] introduced a joint model for dialogue segmentation and topic categorization, which was evaluated on a clinical spoken conversation dataset created by them. In another work, Xu et al. [15] developed a Dialogue State Distillation Network (DSDN), which leverages relevant information of previous dialogue states and employs an inter-slot contrastive learning loss to effectively capture the slot co-update relations from dialogue context. Their proposed method achieved state-of-the-art performance on the dialogue state tracking task. Sabour et al. [21] introduced a novel approach for empathetic response generation in dialogue, which leverages commonsense to draw more information about the user's situation and uses that to further enhance the empathy expression in generated responses. They showed that their approach outperforms the baseline models in both automatic and human evaluations.

### C. Use of Pre-trained Language Models for Discourse Modeling in Dialogue

Recent advancements of pre-trained large language models (LLMs) has significantly advanced the field of discourse modeling [37]–[40]. The importance of modeling speaker turns in dialogues was investigated by He et al. [6], where they incorporated turn changes in conversations among speakers for the dialogue act classification task. They introduced speaker turn embeddings and added them to utterance embeddings produced by the pretrained language model RoBERTa [41], which showed better performance for the dialogue act classification task. Xing and Carenini [11] utilized a neural utterance-pair coherence scoring model based on fine-tuning NSP BERT
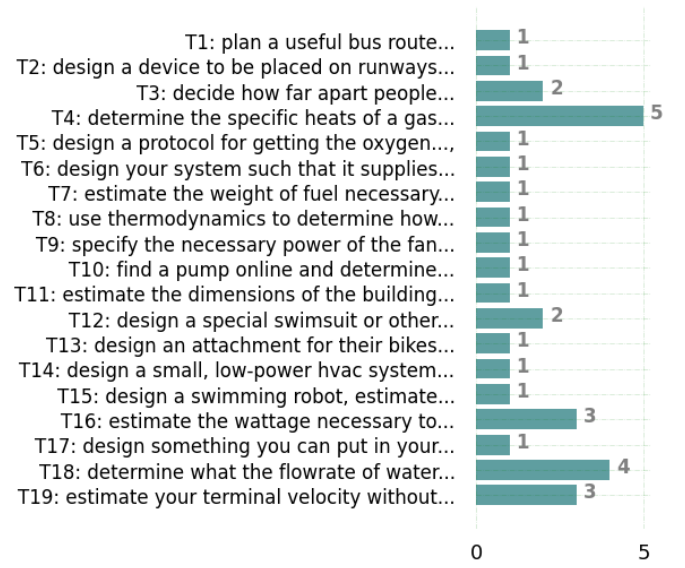


Fig. 2: Topic distribution across the dataset

[42] and achieved state-of-the-art results on the Dialogue topic segmentation task across three public datasets. Feng et al. [14] presented the first evaluation of ChatGPT on the dialogue state tracking task, highlighting its superior performance over prior methods. They also proposed an LLM-driven dialogue state tracking framework based on smaller, open-source foundation models and showed that it achieves comparable performance to ChatGPT. Finch et al. [20] investigated the ability of the state-of-the-art large language model (LLM),i.e., ChatGPT-3.5, to perform dialogue behavior detection for nine categories in real human-bot dialogues and showed that although ChatGPT performed promisingly, often outperforming specialized detection models, the result is still not up to human performance.

Few researches have been conducted that focus on discourse modeling on educational dialogue data. Jensen et al. [22] proposed a methodology for providing teachers with objective, automated feedback on the quality of their classroom discourse by comparing traditional open-vocabulary approaches using n-grams and Random Forest classifiers with a modern deep transfer learning method leveraging BERT. By modeling seven key features of teacher talk (such as questioning and elaborated evaluation) on 127 recordings of classroom talk, the authors demonstrated that while transfer learning with BERT offers a promising path for enhancing automated discourse analytics in education, its effectiveness hinges on the availability of sufficient annotated data to fine-tune the model effectively. Alic et al. [23] automatically distinguished between two pedagogically significant types of teacher questions: funneling questions, which guide students toward specific answers, and focusing questions, which encourage students to reflect on their reasoning. The authors create a labeled dataset of over 2,000 teacher questions annotated by experts and develop both supervised (fine-tuned RoBERTa) and unsupervised models to classify question types. Their supervised RoBERTa model showed strong alignment with expert judgments and correlates with key educational outcomes, such as instructional quality
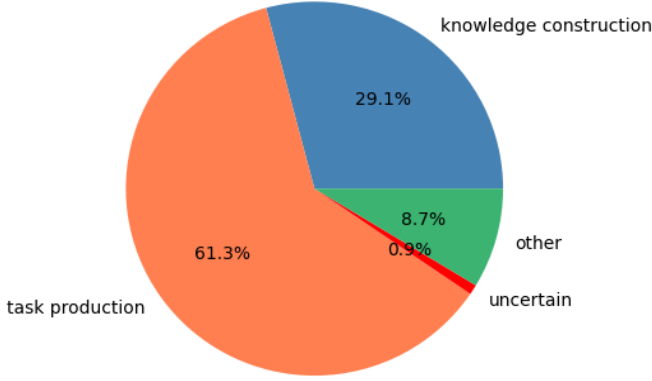
Fig. 3: Distribution of categories across the dataset.



Fig. 4: Confusion matrix of dual annotations

and student learning gains.

## III. DATASET CURATION

### A. Data Collection

We recorded homework discussions among undergraduate mechanical engineering students, focusing on topics from their thermal fluid systems course. Between 2 and 5 students participated in each conversation. Then, we transcribed the conversations ensuring that all data were de-identified. As part of the consent process, students were asked if their de-identified transcripts could be used in future research. Only transcripts from students who consented were included in the dataset.

### B. Dataset Statistics

The dataset consists of 32 small-group conversations covering 19 homework topics, each topic corresponding to a distinct *task description* that students were required to complete collaboratively through discussion. Fig 2 shows the topic distribution across the dataset. The utterances in the conversations are segmented based on the fact that one "turn," or utterance, consists of everything a single person utters until another person speaks (either because the first person has finished or because they interrupt the first person). The average token per conversation is 6404, and the average turns of talk is 321. Please see the appendix for the details of each topic.

### C. Annotation Study

*1) Setup:* Two expert annotators, including one co-author of this paper, participated in the annotation study. We developed a comprehensive annotation guideline and instructed the annotators to label each conversational turn as *knowledge construction*, *task production*, *uncertain*, or *other*. We created the label *uncertain* for the turns of talk where there is insufficient evidence to determine whether the speaker is continuing the current framing of either knowledge construction or task production. If a single utterance includes indicators for both
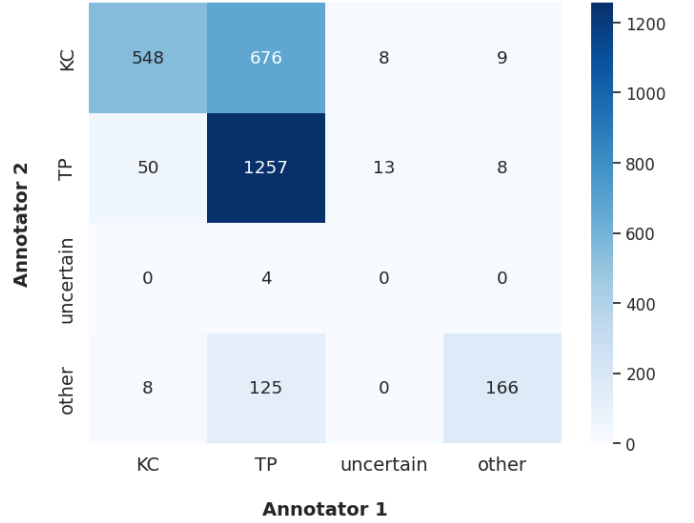
KC and TP classification, and the annotator cannot determine which category is the predominant framing for the student during the utterance, the utterance should be classified as *uncertain*. The label *other* refers to the turns of talk where students discuss a topic other than the assigned problem, such as the purpose of participating in the research study, or other academic classes or social events.

We trained the annotators in a pilot annotation phase where they were asked to annotate 5 conversations. After the pilot annotation, we discussed the disagreements and, if needed, adjusted the annotation guidelines. In our main annotation study, 6 conversations were annotated by two annotators and 21 conversations were annotated by a single annotator. For inter-annotator agreement (IAA) and the analysis of annotations, we report the results of dual annotations. Fig 3 illustrates the distribution of annotated labels across the dataset of 32 conversations. As anticipated, we see that the dataset is imbalanced, with the majority of annotated labels falling into the *knowledge construction* and *task production* categories.

*2) Inter-annotator agreement (IAA):* We computed IAA using Cohen's ($\kappa$) [43] for the dual annotations across four annotated discourse labels (i.e., knowledge construction, task production, uncertain and other) . We obtained Cohen's ($\kappa$) of 0.45 which indicates a moderate agreement [44], [45].

Discerning undergraduate students' aims and purposes based on their spoken word is notoriously difficult for learning sciences researchers [5]. The difficulty in quickly determining whether students are in task production or knowledge construction mode (or when those modes are co-occurring) is one reason researchers are interested in exploring algorithm-assisted annotation. This also means that it is not surprising that the agreement between the two annotators was only moderate.

*3) Analysis of Annotations:* The confusion matrix of the dual annotations of 6 conversations is shown in Figure 4. We see that both annotators mostly agree with each other during the labeling of *task production (TP)* discourse and the most disagreement happens when one annotator thinks a turn of

talk is *knowledge construction (KC)* while other thinks that it's *task production.*

Figure 4 reveals that Annotator 1 leaned toward classifying discourse as TP, while Annotator 2 leaned toward classifying discourse as KC. Of the utterances on which there was TP vs KC disagreement between the two human annotators, Annotator 1 chose KC for only 7% (50/726) of the disagreements while Annotator 2 chose KC in 93% (676/726) of the cases. Besides, where there was TP vs "other" disagreement, 94% (125/133) times Annotator 1 chose TP over other, while just 0.06% (8/133) times Annotator 2 chose TP over "other."

We also found that the disagreement mostly happens under two conditions: (i) when students discuss the details of their problem-solving steps, and (ii) when students ask each other questions. For example, Consider the conversation snippet below (from Topic 4, "determine the specific heats").

> T: Right. Also how long are we doing it for. It's for like -
> X: Yeah. Do it for ten hours. Do we need another you know microsecond.
> T: Yeah. Um ok. So then graph um V I versus time and take the area under the curve. Um. Ok. That area under the curve is just gonna be equal to Q, right?

The students here are discussing the details of an experimental design. Their homework task is to specify the design set-up. They consider the duration of the experiment, the plot they will produce from the data, and the physical quantity represented on that plot. Annotators 1 and 2 disagree on whether this portion of the discussion is aimed toward deeper understanding or toward making progress on the assignment. On one hand, discussion of experiment timescales and of the meaning of a graph might help students build knowledge about the physical quantity to be measured in the experiment. On the other hand, the students' statements about the length of the experiment and of the graphs it will generate could simply comprise another step forward in specifying an experimental design, which is completing the homework task.

The correct label in the case is knowledge construction. When student X discusses the experiment's timescale and student T discusses the meaning of the area under the curve, they are talking about concepts that they contributed anew to this homework discussion; these were not concepts mentioned in the homework problem statement, notes, or textbook for this course. Therefore, the students were calling up other intellectual resources to construct new ideas for this homework activity.

## IV. EXPERIMENTS

### A. Task setting

We consider the prediction task of KCTP discourse in a conversation as a label-generation task for each turn of talk in the conversation, where the model is instructed to generate one label out of the four annotated labels i.e., *knowledge construction (KC), task production (TP), uncertain, other* for each turn of talk.

To create a strong baseline, we assume that in cases where such KCTP discourse-specific resources are unavailable, pre-trained large language models (LLMs) could be the most effective means of generating KCTP discourse labels for each turn in the conversation. We evaluate our task in three settings: (i) **Zero-shot prompting setting**: Zero-shot prompting is a technique used with large language models in which a task is defined using only natural language instructions, without providing any examples of how the task should be performed. This method relies on the model's pre-trained knowledge and ability to generalize in order to accurately interpret and carry out the given instruction. (ii) **Few-shot prompting setting**: Few-shot prompting is a technique where a language model is given a few input-output examples along with a natural language instruction to guide its response to new, similar inputs. In contrast to zero-shot prompting, which relies only on instructions, few-shot prompting uses these examples to establish a pattern or context that the model can mimic. This method exploits the in-context learning ability of large language models, allowing them to generalize from a small number of examples without the need for task-specific fine-tuning [24]. (iii) **Fine-tuning setting**: Fine-tuning refers to the process of taking a pre-trained large language model (which is generally trained on a large, general-purpose corpus) and further training it on a smaller, task-specific dataset to improve its performance on a particular task. This transfer learning strategy [42] allows models to leverage the rich representations learned during pre-training and adapt them to specialized tasks.

### B. Models

We employ state-of-the-art LLMs namely GPT-3.5-turbo [24] and Llama-3.1-8B-Instruct [33] models for the KCTP discourse prediction task while we use GPT-4-1106-preview [46] for our prompt engineering [28]. A GPT (Generative Pre-trained Transformer) model [47] is an auto-regressive large language model developed by OpenAI that uses transformer [48] architecture to generate and understand human-like text. GPT models use a transformer decoder architecture, which is trained to predict the next word in a sequence, followed by fine-tuning on labeled datasets for specific applications. GPT-3.5 Turbo is optimized for speed and cost-efficiency, making it ideal for high-volume tasks. In contrast, GPT-4 offers superior reasoning, accuracy, and contextual understanding for more complex applications while costing more as well. Like GPT, the LLaMA (Large Language Model Meta AI) series [25] developed by Meta is also an auto-regressive language model based on the transformer architecture. Its key advantage lies in being open-source, enabling cost-free use while still delivering competitive performance.

### C. Prompt Design

We create 5 prompts for the KCTP prediction task and use the GPT-4-1106-preview model to optimize our created prompts. We report results for both types of prompts, i.e., our curated prompts and GPT-4 optimized prompts. We also use the optimized prompts for instruction fine-tuning of Llama-3.1 8B model. Among the 5 prompts, prompts 1 and 2 consist of the previous dialogue context along with the current turn of talk. Prompts 3 and 4 include the task description and the definitions of the labels respectively in addition to the previous

| Prompts | Zero-Shot | | | | Few-Shot | | | | Fine-Tuned |
|---|---|---|---|---|---|---|---|---|---|
| | Curated prompt | | Optimized prompt | | Curated prompt | | Optimized prompt | | Optimized prompt |
| | GPT-3.5 | Llama-3.1 | GPT-3.5 | Llama-3.1 | GPT-3.5 | Llama-3.1 | GPT-3.5 | Llama-3.1 | Llama-3.1 |
| Prompt 1 | 0.34 | 0.43 | 0.29 | 0.49 | 0.35 | 0.48 | 0.26 | 0.50 | 0.54 |
| Prompt 2 | 0.28 | 0.38 | 0.47 | 0.40 | 0.35 | 0.45 | 0.37 | 0.48 | 0.51 |
| Prompt 3 | 0.49 | 0.47 | 0.55 | 0.44 | 0.39 | 0.51 | 0.40 | 0.47 | 0.45 |
| Prompt 4 | 0.32 | 0.52 | 0.46 | 0.50 | 0.38 | 0.54 | 0.44 | 0.57 | 0.49 |
| Prompt 5 | 0.27 | 0.39 | 0.33 | 0.46 | 0.27 | 0.44 | 0.27 | 0.49 | 0.55 |

TABLE I: Performance of GPT-3.5 and Llama-3.1 in the label prediction task under zero-shot, few-shot, and fine-tuned settings using different prompt types.

dialogue context and current turn of talk. Prompt 5 consists of both the previous and afterward dialogue context and the current turn of talk. Please see the details of these prompts in Appendix.

### D. Setup

We conduct experiments in zero-shot and eight-shot (two examples for each of the four labels) prompt settings where the number of shots reflects the number of examples provided in the prompt. Few-shot examples were sampled from two conversations and topics not included in the dataset.

We use OpenAI's API for GPT models and set the temperature of the model at 0.0 and maximum tokens at 10. To use Llama-3.1 (8B) for our task, we use *Unsloth*, an open-source AI library that enables us to train an LLM faster and efficiently with less GPU memory by applying techniques like quantization [49] and low-rank adaptation (LoRA) [50]. In our zero-shot and few-shot experiments, we set the Llama-3.1 model with a temperature of 1.5 and a maximum of 64 new tokens. Fine-tuning is performed for 5 epochs using a learning rate of 1e-4 with the AdamW 8-bit optimizer. We use a batch size of 8, gradient accumulation of 16, and a weight decay of 0.01. All experiments are conducted with a fixed random seed of 3407.

### E. Evaluation Procedure

We use the weighted F1 score to evaluate model performance. During fine-tuning Llama-3.1, we employ five-fold cross-validation to obtain results across the entire dataset and enable comparison with zero-shot and few-shot prompting.

## V. RESULTS AND DISCUSSIONS

Table I presents the performance of GPT-3.5 and LLaMA-3.1 models across three experimental settings: Zero-Shot, Few-Shot, and Fine-Tuned, using both curated and GPT-4 optimized prompts on the KCTP discourse label prediction task. Five prompts (Prompt 1–5) developed by us are evaluated, and the scores represent the average F1 score.

### A. Zero-shot and Few-shot effectiveness

The results indicate that overall, Llama-3.1 performs better than GPT-3.5 for both curated and optimized prompts across zero-shot and few-shot settings. However, the overall performance remains suboptimal, with the highest F1-score reaching only 0.57.

In the zero-shot setting, the best result is obtained by GPT-3.5 with optimized prompt 3, which includes a topic description in addition to the previous dialogue contexts. The result suggests that explicitly providing the communicative goal of the student conversation helps the model infer appropriate labels without prior examples. However, the GPT-3.5 performance drops in the few-shot setting. We assume that one of the reasons the few-shot prompting did not perform better here could be attributed to the fact that the examples we used didn't generalize well with the dataset, or the model had too much information to process. Moreover, performance degradation can sometimes happen in some LLMs for adding too domain-specific examples [51]. In contrast, LLaMA-3.1 attains its highest score in the few-shot setting with optimized prompt 4, which incorporates label definitions alongside the preceding dialogue context. It means that when we include examples, LLaMA-3.1 can better generalize than GPT-3.5 on this task by leveraging explicit label information. Table I also shows that optimizing our curated prompt with GPT-4 improves the overall performance.

### B. Fine-tuning effectiveness

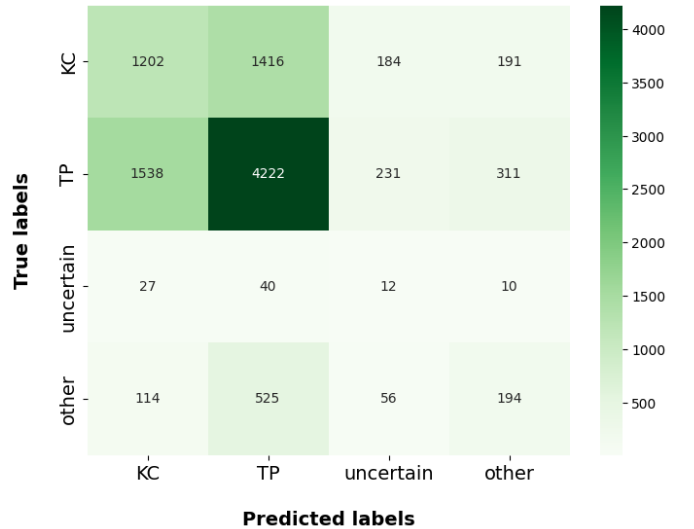The fine-tuning results presented in Table I show that fine-tuning Llama-3.1 does not exceed its zero-shot and few-



Fig. 5: Confusion matrix of model prediction vs. true annotated labels

| Prompts | Same topic | Single different topic | Three different topics |
|---|---|---|---|
| Prompt 1 | 0.64 | 0.61 | 0.65 |
| Prompt 2 | 0.58 | 0.63 | 0.61 |
| Prompt 3 | 0.60 | 0.63 | 0.61 |
| Prompt 4 | 0.53 | 0.58 | 0.62 |
| Prompt 5 | 0.64 | 0.62 | 0.61 |

TABLE II: Performance of finetuned Llama-3.1 when trained on the same topic, a single different topic, and three different topics.

shot performance. We hypothesize that since our data is very domain-specific, it may have reduced the model's generalization ability by overfitting to narrow linguistic or conceptual patterns. Recent studies also showed that, in certain cases, fine-tuning often yields limited or even negative performance gains for large language models (LLMs) [52].

Furthermore, the results suggest that incorporating dialogue context alone (as in Prompts 1 and 3), without additional task-specific contextual signals, leads to relatively better outcomes during fine-tuning. This observation implies that minimal yet coherent contextual grounding may help the model retain its pretrained reasoning and discourse capabilities.

### C. Analysis

Figure 5 presents the confusion matrix of the true annotated labels versus the model-predicted labels for the best result (i.e., Llama 3.1 performance in the few-shot settings with optimized prompt 4). It shows that the model is good at predicting the *task production* labels. However, it struggles with predicting the *knowledge construction* label and misclassifies about half of these instances as *task production*, indicating notable room for improvement. Moreover, the model also kind of struggles with predicting the *"other"* label, only correctly predicting it 22% of the time. These findings suggest that while the model captures dominant discourse functions well, it struggles with more nuanced or less frequent categories, highlighting opportunities for improving label representation and contextual modeling.

Since our dataset consists a diverse set of 19 topics (see Appendix for the details of each topic), we investigated how topic similarity between training and testing data influences model performance. For this experiment, we kept the test data same, and the training data size fixed, but varied the topical composition of the training data across three configurations: (i) training and test data drawn from the same topic, (ii) training data from a different single topic than the test set, and (iii) training data from three different topics distinct from the test set.

The results, summarized in Table II, reveal that, the model performs better when trained on data from different topics, except for prompt 5. This finding suggests that exposure to a broader range of linguistic and conceptual patterns enhances the model's generalization ability, whereas training on a single, homogeneous topic may lead to overfitting for most prompts. Notably, the model exhibits slightly better generalization within the same-topic setting for Prompt 5, which incorporates succeeding dialogue context. This indicates that the model might learn topic-bounded discourse dynamics

i.e., how participants introduce, elaborate on, and shift ideas within a coherent topical space in this prompt setting whereas such recurring patterns are less transferable across topics.

## VI. LIMITATIONS

In this study, we only consider two types of discourse features in educational conversations, namely *knowledge construction (KC)* and *task production (TP)*. Besides, the dataset is limited to topics from a thermal fluid systems course in mechanical engineering Moreover, we use a limited set of prompt templates because of the resource and time constraints.

## VII. CONCLUSIONS AND FUTURE WORK

This work presents a novel educational conversational dataset, annotated with *knowledge construction (KC)* and *task production (TP)* discourse. Such discourse properties are crucial for framing student learning activities to develop more effective pedagogical settings that emphasize knowledge construction over mere task completion. In this work, we establish baselines for the KCTP discourse prediction task using state-of-the-art language models with prompting techniques and fine-tuning. Our results demonstrate that state-of-the-art LLMs struggle with this task, both under prompting-based and fine-tuning settings.

For future work, we plan to create reasoning chains that will help the model better understand the definitions of the labels in the prompts. We also intend to annotate low-level discourse structure for these student dialogues so that looking at the lower levels might help to see how the higher-order concepts emerge from a particular interaction of dialogue moves. Furthermore, we aim to expand the dataset by including a broader range of undergraduate subjects, thereby capturing more diverse discourse patterns across academic domains. This increased topical diversity will support more robust fine-tuning and facilitate the development of models capable of domain-general discourse understanding.

APPENDIX

| Topic | Description |
|---|---|
| Topic 1 (T1) | about a decade ago, stanford university successfully tried using waste vegetable oil from the dining halls as fuel for campus shuttles (https://news.stanford.edu/news/2006/january25/biodiesel-012506.html). what if [institution] tried to do this? plan a useful bus route around [institution] and specify the volume of fuel needed for the bus to travel this route without having to refuel. you may assume the energy density of vegetable oil is 42.20 mj/kg or 30.53 mj/l. |
| Topic 2 (T2) | contrails are giant vortices left by airplanes on the runway and in the sky. if other planes pass through these, it can cause problems because it is like going through a mini tornado, and the planes are not equipped to handle such a pressure gradient. boeing has hired you to design a device to be placed on runways to help get rid of contrails there. this could be done by moving the contrails out of the way or by stopping them altogether. justify your design using fluid mechanics. |
| Topic 3 (T3) | covid-19 has drastically changed how people live their daily lives. guidelines have been created for how far apart people should stay when talking normally to each other. however, if people are doing something like singing, which takes more effort results in air (and droplets potentially carrying the coronavirus) being expelled from the lungs more forcefully, the guidelines for simply talking may not be adequate. if six feet apart is the recommendation for talking, use fluid mechanics argument to decide how far apart people who are singing should stand in order to be far enough away from any particles that may be expelled into the air by their singing. |
| Topic 4 (T4) | design an experiment complete with instrumentation to determine the specific heats of a gas using a resistance heater. discuss how the experiment will be conducted, what measurements need to be taken, and how the specific heats will be determined. what are the sources of error in your system? how can you minimize the experimental error? |
| Topic 5 (T5) | gas turbine engines used in airplanes consist of a fan followed by a compressor, diffuser, combustor, turbine, and sometimes an afterburner. you are designing the engine for a high- altitude airplane. normally, commercial planes operate best around 35,000 ft above sea level, but your plane should operate optimally at around 100,000 ft. because of the high altitude, there will be a lower concentration of oxygen than normal, and the air entering the engine will be colder. design a protocol for getting the oxygen up to the appropriate temperature and pressure needed for combustion. keep in mind your solution has to be relatively light. |
| Topic 6 (T6) | geothermal heat pumps harness renewable geothermal energy by using thermal reservoirs of water deep within the earth for heating. such reservoirs have temperatures up to around 370 degrees celsius. geothermal heat pumps use this energy by transporting room-temperature or cold liquid deep into the ground via pipes, exposing it to the hot reservoir, and carrying it back up to the surface. imagine one of these reservoirs is discovered beneath the building where you live and design a heat exchanger system that uses the reservoir to heat your building. sketch your system and specify the diameter, length, and material of the pipe, the flow rate, and the working fluid. design your system such that it supplies a significant portion of the energy required for your building to operate normally. |
| Topic 7 (T7) | hybrid rockets use a combination of solid and liquid or gaseous propellants. in hybrid rockets, a stable oxidizer is used with a solid fuel. in order to be used, the fuel needs to be vaporized. the primary difficulty with hybrids is with mixing the propellants during the combustion process. in a hybrid rocket, the mixing happens at the melting or evaporating surface. the mixing is not well-controlled and generally, a lot of propellant is left unburned, limiting the motor's efficiency. on the other hand, liquid propellants are generally mixed with oxidizer by an injector at the top of the combustion chamber which directs many small streams of fuel and oxidizer into one another. based on reasonable efficiencies of both liquid fuel and hybrid fuel processes, estimate the weight of fuel necessary to get a specific rocket of your choice to low earth orbit if the fuel is liquid vs. hybrid. |
| Topic 8 (T8) | most ski resorts in the u.s. use snow guns to make additional snow to supplement natural snow. these machines use water and compressed air. the air forces the water to form tiny droplets, which are then expelled from the nozzle and form ice crystals, which then fall to the ground as snow. compressed air cools as it expands, which assists with converting the water droplets into snow. choose your favorite ski resort and the desired depth of snow for the best skiing, and use thermodynamics to determine how long it will take to cover the ski trails in that amount of snow. you may assume that one snow gun uses about 100 gallons of water per minute and that the compressor can produce 50 cfm (cubic feet per minute) of air. |

| | |
|---|---|
| Topic 9 (T9) | race cars need to be as aerodynamic as possible. in many cases, to test the aerodynamics of a car, a wind-tunnel is used. you have been hired by chevrolet to analyze the air flow around their race cars. the wind tunnel you will be using to do this is an open circuit wind tunnel, where air is drawn from the laboratory environment, rather than being recirculated in the wind tunnel itself. such wind-tunnels consist of a nozzle to accelerate the air, the test section in which the car sits, and a diffuser which decelerates the air. based on reasonable values for air speed around the vehicle being tested, design a wind-tunnel for testing a race car. include all necessary specifications of the different parts of the wind tunnel, such as dimensions and air speeds. also specify the necessary power of the fan and estimate the head loss due to the vehicle. use fluid mechanics to justify your response. |
| Topic 10 (T10) | since they know you are a mechanical engineer, your neighbors have asked you to help them design a waterfall for their garden similar to the one in the image below. you need to devise a way to get water from the pool at the bottom up to the top of the waterfall, and there needs to be enough water so that the waterfall actually looks good. design a system to do this. include a diagram of how the pump system will work, and include any important specifications such as flow rates and dimensions. then find a pump online and determine approximately how much power the waterfall pump will use per day. you may make as many assumptions as needed, just specify what assumptions you are making and why. |
| Topic 11 (T11) | the building of farfar's danish ice cream shop in duxbury, ma is somewhat old and thus does not seem to have a great cooling system. as a result, sometimes the ice cream gets a bit melty even when it's still in the freezer. the temperature in the ice cream shop is to be maintained at 55°f. estimate the dimensions of the building, use thermodynamics principles to determine the maximum heat loss the shop can have, and suggest a method for minimizing this heat loss. |
| Topic 12 (T12) | the butterfly swimming stroke is considered by many to be one of the most difficult strokes. it is also one of the fastest. when used over longer distances, the butterfly stroke is slightly slower than freestyle, partly due to the greater physical exertion required by the butterfly. however, butterfly has the fastest peak speed. explain why you think this stroke has the fastest peak speed. then, design a special swimsuit or other (non-motorized) device for a swimmer to further increase the speed of the butterfly stroke so that it will always be faster than freestyle no matter the distance over which the stroke is used. include a diagram of your design, and use fluid mechanics principles to prove that it will work. |
| Topic 13 (T13) | trek bikes has contracted you to design an attachment for their bikes to help make the bike and rider more streamlined. this attachment should effectively reduce the bike and rider's air resistance without impeding the cyclist's ability to ride their bike as usual. also specify what material this should be made of, and include a diagram of your design. justify your design using fluid mechanics. |
| Topic 14 (T14) | you are designing a tiny home that can be used for camping adventures. you want to be able to take your tiny home on camping trips in vermont and new hampshire during the fall to see the foliage, but you are worried that it might get a bit too cold for comfort, as that time of year, the temperatures at night can get down to 30°F. design a small, low-power hvac system to keep the inside of your tiny home at a temperature no lower than 45°F. specify what parts will be needed and how this system will be compatible with your tiny home. use fluid mechanics and heat transfer principles to justify that this system will indeed keep the temperature at 45°F or higher. |
| Topic 15 (T15) | you are working at a robotics company to design a robot that can swim in water to collect data on sharks. this robot needs to be as hydrodynamic as possible so that it is efficient, and you need to be able to control how fast the robot will go so it can keep up with the sharks, as well as be able to make it turn while swimming. design a swimming robot, estimate its drag coefficient and the drag on the robot when it is moving at three different speeds (so you should have three different values for drag). then determine how much power will be needed to make the robot move forward at each of the three speeds. include a diagram of your robot in your response. |
| Topic 16 (T16) | you have been contracted by [institution] to design a system to get hot water to different parts of the science and engineering complex (sec). in particular, this system needs to work well during winter, when it is colder outside and most likely slightly colder than usual within the outer walls of the building and in the building in general. estimate the wattage necessary to keep the water at a reasonably hot temperature, and determine the flow rates and pressures necessary to get the hot water to various parts of the building. include a labeled sketch of your design, and be sure to use fluid mechanics to justify that your design will work. |

| Topic 17 (T17) | you have been doing a lot of baking recently and wish that you had a convection oven. convection ovens have one or more fans that help circulate the air in the oven, whereas in regular ovens, the only thing moving the air is natural convection. therefore, you want a convection oven so that you can bake everything faster and more evenly. however, you don't want to spend the money on an entirely new oven since convection ovens are expensive, and you don't want to have to get rid of the regular oven you already have. design something you can put in your regular oven that will make it function similarly to a convection oven. specify air flow rates and estimate the power needed for any components. also draw a diagram of your design and specify where any proposed components will go in the oven. use fluid mechanics to justify that your design will make your oven work similarly to a convection oven. |
|---|---|
| Topic 18 (T18) | you have been hired by firefighters to design a tripod to hold a large hose when fighting fires. the stream of water that comes out of the hose is 5 cm in diameter. determine what the flowrate of water out of the hose should be in order to work well for fighting a fire that is 9 meters away. then calculate how much reaction force will be needed at the base of the tripod to keep it from moving when the hose is being used. use fluid mechanics to support your response. |
| Topic 19 (T19) | you have recently gotten into skydiving. when you are skydiving, once you get close enough to the earth, you have to deploy a parachute. the skydiving part is exciting, but once you deploy the parachute, you have been getting bored since when you're falling through the air, you eventually reach one constant speed (the terminal velocity). you want to design an attachment that enables you to increase and decrease your terminal velocity as you are falling. estimate your terminal velocity without this attachment, and then estimate the maximum and minimum terminal velocities when the attachment is being used. use lift and drag calculations to justify your answer. |

TABLE III: Details of the dataset topics, where each topic corresponds to a distinct *task description* that students were required to complete collaboratively through discussion.

| Prompts | Author curated prompt template | GPT-4 optimized prompt template |
|---|---|---|
| **Prompt 1** (Previous dialogue context) | You will be provided with a dialogue and its context. The context is the previous dialogue lines of the given dialogue and each line in context is separated by a newline character. Classify the given dialogue considering its context into one of the four categories: knowledge construction, task production, uncertain, other. Output only one of the categories and do not provide any explanation. #### Context: Dialogue: | Classify the provided dialogue into the correct category based on its context. Choose one of these categories: knowledge construction, task production, uncertain, or other. Only provide the category name as your response. Context: Dialogue: |
| **Prompt 2** (Previous dialogue context) | You will be provided with a current dialogue line and its previous dialogue lines. Each line in the previous dialogue lines is separated by a newline \n character. Classify the current dialogue line considering its previous dialogue lines into one of the four categories: knowledge construction, task production, uncertain, other Output only one of the categories and do not provide any explanation. #### Previous dialogue lines: Current dialogue line: | Classify the current dialogue line into one of the following categories based on its context within the preceding dialogue lines: knowledge construction, task production, uncertain, or other. Provide the category without any explanation. Previous dialogue context: Current dialogue line: |

| | | |
|---|---|---|
| **Prompt 3** (Previous dialogue context & Topic description) | You will be provided with a dialogue, its context and a task description.<br><br>The context is the previous dialogue lines of the given dialogue and each line in context is separated by a newline character. The dialogue and context are about completing the task details in the task description.<br><br>Classify the given dialogue considering its context and task description into one of the four categories: knowledge construction, task production, uncertain, other.<br><br>Output only one of the categories and do not provide any explanation.<br>####<br>Task description:<br>Context:<br>Dialogue: | Given a dialogue along with its preceding context and a specific task description, classify the provided dialogue into one of four categories (knowledge construction, task production, uncertain, other). Provide only the category without any further explanation.<br>Task Description:<br>Context:<br>Dialogue: |
| **Prompt 4** (Previous dialogue context & Label definitions) | You will be provided with a dialogue and its context.<br><br>The context is the previous dialogue lines of the given dialogue, and each line in context is separated by a newline \n character. Classify the given dialogue considering its context into one of the four categories: knowledge construction, task production, uncertain, other.<br>The definition of each of the categories is given below:<br><br>knowledge construction: means the dialogue is focused on expressing understandings of concepts, phenomena, or technologies. Simply stating a definition from textbook or notes does not count as knowledge construction.<br><br>task production: means the dialogue is focused on completing the assigned task to satisfy the instructor, without verbalizing regard for understanding the bigger picture. For example, the dialogue is stating an equation, or asksing for a numerical answer, or calculating a number out loud, or discussing what to do next.<br><br>uncertain: means there is insufficient evidence for classifying the dialogue either as a knowledge construction or task production. It is the default category for one word conversational fillers such as 'yeah', 'okay'.<br><br>other: means the dialogue is about a topic other than the assigned task.<br><br>Output only one of the categories and do not provide any explanation.<br>####<br>Context:<br>Dialogue: | Given the dialogue and its preceding context, classify the dialogue into one of the following four categories: knowledge construction, task production, uncertain, or other.<br><br>- Knowledge Construction: The dialogue focuses on deep understanding of concepts or phenomena, going beyond mere definitions.<br><br>- Task Production: The dialogue aims at completing a task or assignment, primarily focusing on procedural steps.<br><br>- Uncertain: The dialogue does not provide enough information for classification into the above categories or includes filler words like 'yeah', 'okay'.<br><br>- Other: The dialogue discusses topics unrelated to the assigned task.<br><br>Provide only the category without any explanation.<br>Context:<br>Dialogue: |
| **Prompt 5** (Previous and afterward dialogue context) | You will be provided with a dialogue, its before context and its after context<br><br>The before context is the previous dialogue lines and after context is the succeeding dialogue lines of the given dialogue. Each line in before and after context is separated by a newline \n character<br><br>Classify the given dialogue considering its before and after context into one of the four categories: knowledge construction, task production, uncertain, other.<br><br>Output only one of the categories and do not provide any explanation.<br><br>####<br>Before Context:<br>Dialogue:<br>After Context: | Your task is to classify a specific dialogue based on the surrounding context into one of the following categories: knowledge construction, task production, uncertain, other. You will be given the dialogue, as well as the lines of conversation that precede it (Before Context) and follow it (After Context). Each dialogue line in the contexts is separated by a newline character.<br><br>Please provide only the category as your response without any explanation.<br>Before Context:<br>Dialogue:<br>After Context: |

TABLE IV: Details of the prompts used in modeling for knowledge construction and task production discourse prediction.

REFERENCES

[1] J. Gouvea, V. Sawtelle, and A. Nair, "Epistemological progress in physics and its impact on biology," *Physical Review Physics Education Research*, vol. 15, no. 1, p. 010107, 2019.

[2] M. D. Koretsky, D. M. Gilbuena, S. B. Nolen, G. Tierney, and S. E. Volet, "Productively engaging student teams in engineering: The interplay between doing and thinking," in *2014 IEEE Frontiers in Education Conference (FIE) Proceedings*. IEEE, 2014, pp. 1–8.

[3] J. E. S. Swenson, "Developing knowledge in engineering science courses: Sense-making and epistemologies in undergraduate mechanical engineering homework sessions," Ph.D. dissertation, Tufts University, 2018.

[4] D. L. Schwartz, J. M. Tsang, and K. P. Blair, *The ABCs of how we learn: 26 scientifically proven approaches, how they work, and when to use them*. WW Norton & Company, 2016.

[5] R. E. Scherr and D. Hammer, "Student behavior and epistemological framing: Examples from collaborative active-learning activities in physics," *Cognition and Instruction*, vol. 27, no. 2, pp. 147–174, 2009.

[6] Z. He, L. Tavabi, K. Lerman, and M. Soleymani, "Speaker turn modeling for dialogue act classification," *arXiv preprint arXiv:2109.05056*, 2021.

[7] V. Raheja and J. Tetreault, "Dialogue act classification with context-aware self-attention," *arXiv preprint arXiv:1904.02594*, 2019.

[8] R. Li, C. Lin, M. Collinson, X. Li, and G. Chen, "A dual-attention hierarchical recurrent neural network for dialogue act classification," *arXiv preprint arXiv:1810.09154*, 2018.

[9] H. Kumar, A. Agarwal, R. Dasgupta, and S. Joshi, "Dialogue act sequence labeling using hierarchical encoder with crf," in *Proceedings of the aaai conference on artificial intelligence*, 2018.

[10] Z. Liu, S. U. M. Salleh, H. C. Oh, P. Krishnaswamy, and N. Chen, "Joint dialogue topic segmentation and categorization: A case study on clinical spoken conversations," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, 2023, pp. 185–193.

[11] L. Xing and G. Carenini, "Improving unsupervised dialogue topic segmentation with utterance-pair coherence scoring," *arXiv preprint arXiv:2106.06719*, 2021.

[12] S. Somasundaran *et al.*, "Two-level transformer and auxiliary coherence modeling for improved text segmentation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 05, 2020, pp. 7797–7804.

[13] S. Kim, R. E. Banchs, and H. Li, "Towards improving dialogue topic tracking performances with wikification of concept mentions," in *Proceedings of the 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, 2015, pp. 124–128.

[14] Y. Feng, Z. Lu, B. Liu, L. Zhan, and X.-M. Wu, "Towards llm-driven dialogue state tracking," *arXiv preprint arXiv:2310.14970*, 2023.

[15] J. Xu, D. Song, C. Liu, S. C. Hui, F. Li, Q. Ju, X. He, and J. Xie, "Dialogue state distillation network with inter-slot contrastive learning for dialogue state tracking," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 11, 2023, pp. 13 834–13 842.

[16] M. D. Ma, J.-Y. Kao, S. Gao, A. Gupta, D. Jin, T. Chung, and N. Peng, "Parameter-efficient low-resource dialogue state tracking by prompt tuning," *arXiv preprint arXiv:2301.10915*, 2023.

[17] J. Guo, K. Shuang, J. Li, Z. Wang, and Y. Liu, "Beyond the granularity: Multi-perspective dialogue collaborative selection for dialogue state tracking," *arXiv preprint arXiv:2205.10059*, 2022.

[18] Y. Zhou, G. Zhao, and X. Qian, "Dialogue state tracking based on hierarchical slot attention and contrastive learning," in *Proceedings of the 31st ACM international conference on information & knowledge management*, 2022, pp. 4737–4741.

[19] G. Qixiang, G. Dong, Y. Mou, L. Wang, C. Zeng, D. Guo, M. Sun, and W. Xu, "Exploiting domain-slot related keywords description for few-shot cross-domain dialogue state tracking," in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2022, pp. 2460–2465.

[20] S. E. Finch, E. S. Paek, and J. D. Choi, "Leveraging large language models for automated dialogue analysis," *arXiv preprint arXiv:2309.06490*, 2023.

[21] S. Sabour, C. Zheng, and M. Huang, "Cem: Commonsense-aware empathetic response generation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 10, 2022, pp. 11 229–11 237.

[22] E. Jensen, S. L. Pugh, and S. K. D'Mello, "A deep transfer learning approach to modeling teacher discourse in the classroom," in *LAK21: 11th international learning analytics and knowledge conference*, 2021, pp. 302–312.

[23] S. Alic, D. Demszky, Z. Mancenido, J. Liu, H. Hill, and D. Jurafsky, "Computationally identifying funneling and focusing questions in classroom discourse," *arXiv preprint arXiv:2208.04715*, 2022.

[24] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.

[25] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.

[26] M. Reid, N. Savinov, D. Teplyashin, D. Lepikhin, T. Lillicrap, J.-b. Alayrac, R. Soricut, A. Lazaridou, O. Firat, J. Schrittwieser *et al.*, "Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context," *arXiv preprint arXiv:2403.05530*, 2024.

[27] D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi *et al.*, "Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning," *arXiv preprint arXiv:2501.12948*, 2025.

[28] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig, "Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing," *ACM Computing Surveys*, vol. 55, no. 9, pp. 1–35, 2023.

[29] T. Gao, A. Fisch, and D. Chen, "Making pre-trained language models better few-shot learners," *arXiv preprint arXiv:2012.15723*, 2020.

[30] X. L. Li and P. Liang, "Prefix-tuning: Optimizing continuous prompts for generation," *arXiv preprint arXiv:2101.00190*, 2021.

[31] J. Wei, M. Bosma, V. Y. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le, "Finetuned language models are zero-shot learners," *arXiv preprint arXiv:2109.01652*, 2021.

[32] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, Y. Li, X. Wang, M. Dehghani, S. Brahma *et al.*, "Scaling instruction-finetuned language models," *Journal of Machine Learning Research*, vol. 25, no. 70, pp. 1–53, 2024.

[33] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan *et al.*, "The llama 3 herd of models," *arXiv e-prints*, pp. arXiv–2407, 2024.

[34] W. Li, L. Zhu, W. Shao, Z. Yang, and E. Cambria, "Task-aware self-supervised framework for dialogue discourse parsing," in *2023 Conference on Empirical Methods in Natural Language Processing (EMNLP 2023)*. Association for Computational Linguistics, 2023, pp. 14 162–14 173.

[35] Y. Tulpan and O. Tsur, "A deeper (autoregressive) approach to non-convergent discourse parsing," *arXiv preprint arXiv:2305.12510*, 2023.

[36] F. S. Mim, N. Inoue, S. Naito, K. Singh, and K. Inui, "LPAttack: A feasible annotation scheme for capturing logic pattern of attacks in arguments," in *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, N. Calzolari, F. Béchet, P. Blache, K. Choukri, C. Cieri, T. Declerck, S. Goggi, H. Isahara, B. Maegaard, J. Mariani, H. Mazo, J. Odijk, and S. Piperidis, Eds. Marseille, France: European Language Resources Association, Jun. 2022, pp. 2446–2459. [Online]. Available: https://aclanthology.org/2022.lrec-1.261/

[37] C. Li, Y. Yin, and G. Carenini, "Dialogue discourse parsing as generation: a sequence-to-sequence llm-based approach," in *Proceedings of the 25th annual meeting of the special interest group on discourse and dialogue*, 2024, pp. 1–14.

[38] G. Cimino, C. Li, G. Carenini, and V. Deufemia, "Coherence-based dialogue discourse structure extraction using open-source large language models," in *Proceedings of the 25th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, 2024, pp. 297–316.

[39] X. Gu, K. M. Yoo, and J.-W. Ha, "Dialogbert: Discourse-aware response generation via learning to recover and rank utterances," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 14, 2021, pp. 12 911–12 919.

[40] F. S. Mim, N. Inoue, P. Reisert, H. Ouchi, and K. Inui, "Corruption is not all bad: Incorporating discourse structure into pre-training via corruption for essay scoring," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 2202–2215, 2021.

[41] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized bert pretraining approach," *arXiv preprint arXiv:1907.11692*, 2019.

[42] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, 2019, pp. 4171–4186.

[43] J. Cohen, "A coefficient of agreement for nominal scales," *Educational and psychological measurement*, vol. 20, no. 1, pp. 37–46, 1960.

[44] R. Artstein and M. Poesio, "Inter-coder agreement for computational linguistics," *Computational linguistics*, vol. 34, no. 4, pp. 555–596, 2008.

[45] W. Spooren and L. Degand, "Coding coherence relations: Reliability and validity," 2010.

[46] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat *et al.*, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.

[47] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever *et al.*, "Improving language understanding by generative pre-training," 2018.

[48] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[49] B. Jacob, S. Kligys, B. Chen, M. Zhu, M. Tang, A. Howard, H. Adam, and D. Kalenichenko, "Quantization and training of neural networks for efficient integer-arithmetic-only inference," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2704–2713.

[50] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen *et al.*, "Lora: Low-rank adaptation of large language models." *ICLR*, vol. 1, no. 2, p. 3, 2022.

[51] Y. Tang, D. Tuncel, C. Koerner, and T. Runkler, "The few-shot dilemma: Over-prompting large language models," *arXiv preprint arXiv:2509.13196*, 2025.

[52] S. Barnett, Z. Brannelly, S. Kurniawan, and S. Wong, "Fine-tuning or fine-failing? debunking performance myths in large language models," *arXiv preprint arXiv:2406.11201*, 2024.

**Eric Miller** is a Professor in the Department of Electrical and Computer Engineering and an adjunct Professor in the Departments of Computer Science and Biomedical Engineering at Tufts University. He previously served in the Department of Electrical and Computer Engineering at Northeastern University from 1994 to 2006. He is also a Senior Scientist at the Jean Meyer Human Nutrition Research Center on Aging at Tufts University and currently serves as the Director of the Engineering Education and Centers Division in the Directorate for Engineering at the U.S. National Science Foundation. Dr. Miller received National Science Foundation CAREER Award in 1996 and the Outstanding Research Award from the Northeastern University College of Engineering in 2002. From 2014 to 2018, he served on the Technical Liaison Committee for the IEEE Transactions on Computational Imaging and chaired the SIAM Imaging Sciences Special Interest Group from 2015 to 2017. He was an Associate Editor for the IEEE Transactions on Geoscience and Remote Sensing from 2003 to 2015 and for the IEEE Transactions on Image Processing from 1999 to 2003.

**Kristen Wendell** is an Associate Professor in the Department of Mechanical Engineering and Education at Tufts University. She earned her B.S.E. from Princeton University, her M.S. from the Massachusetts Institute of Technology, and her Ph.D. from Tufts University in 2003, 2005, and 2011, respectively. She currently serves as a CEEO Fellow in the Center for Engineering Education Outreach and as the Co-Director of the Institute for Research on Learning and Instruction at Tufts University. Her research work focuses on characterizing and supporting inclusive, sophisticated disciplinary practices during engineering learning experiences in undergraduate course, K-8 classrooms, and teacher education contexts.

**Farjana Sultana Mim** received her B.Sc. in Computer Science and Engineering from Patuakhali Science and Technology University, Bangladesh, in 2016, and her M.S. and Ph.D. in System Information Sciences from Tohoku University, Japan, in 2019 and 2022. She was a postdoctoral scholar in the Department of Electrical and Computer Engineering at Tufts University, USA, in 2023. She is currently a lecturer in the Department of Computer Science and Information Technology at Patuakhali Science and Technology University. Her research interests include NLP in education, large language models, discourse analysis, unsupervised learning, argumentation, and commonsense reasoning.

**Shuchin Aeron** is a professor in the Department of Electrical and Computer Engineering at Tufts School of Engineering. He received his Ph.D. from Boston University in 2009 and was awarded the best PhD thesis award from both the School Of Engineering and from the Department of Electrical and Computer Engineering. From 2009-2011, he was a postdoctoral research fellow at Schlumberger-Doll Research (SDR), where he worked on signal processing solution products for borehole acoustics resulting in a number of patents. In 2016, he received the NSF CAREER award for his work on multidimensional signals and systems. Shuchin Aeron is presently a senior member of the Institute of Electrical and Electronics Engineers (IEEE) and an associate editor for the ACM transactions on Theory of Probababilistic Machine Learning.