# Tutorly: Turning Programming Videos Into Apprenticeship Learning Environments with LLMs

Wengxi Li
University of Michigan
USA
wengxili@umich.edu

Roy Pea
Stanford University
USA
roypea@stanford.edu

Nick Haber
Stanford University
USA
nhaber@stanford.edu

Hari Subramonyam
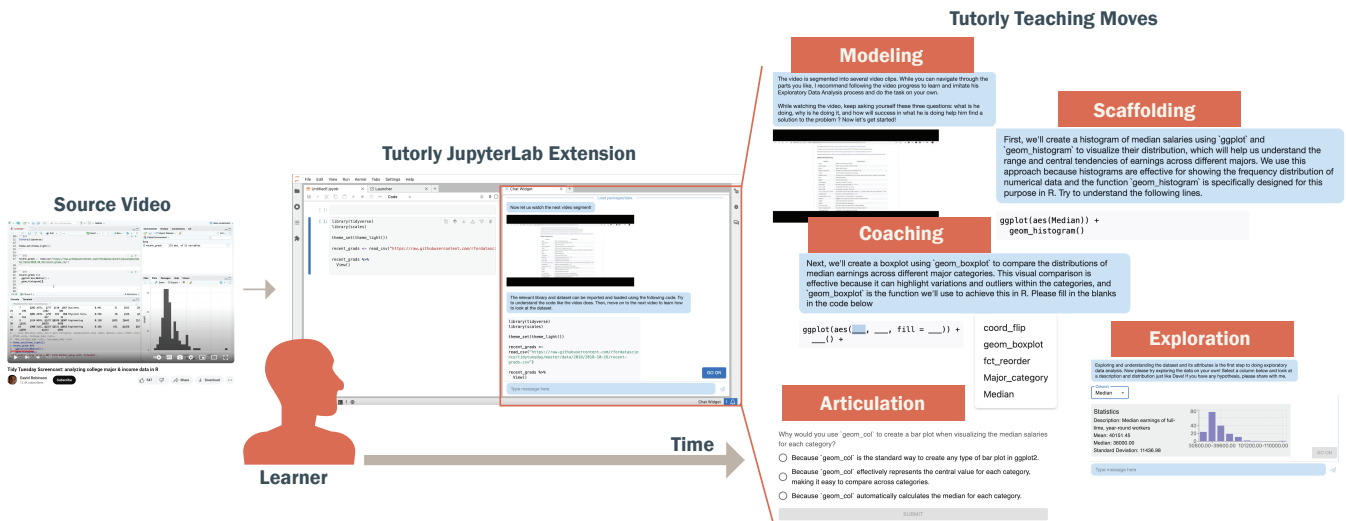Stanford University
USA
harihars@stanford.edu

Figure 1: Overview of Tutorly: As the learner engages with a programming video tutorial using our JupyterLab Extension, Tutorly guides the learners through rich multi-modal conversations that correspond to different moves in the Cognitive Apprenticeship Framework.

## ABSTRACT

Online programming videos, including tutorials and streamcasts, are widely popular and contain a wealth of expert knowledge. However, effectively utilizing these resources to achieve targeted learning goals can be challenging. Unlike direct tutoring, video content lacks tailored guidance based on individual learning paces, personalized feedback, and interactive engagement necessary for support and monitoring. Our work transforms programming videos into one-on-one tutoring experiences using the cognitive apprenticeship framework. Tutorly, developed as a JupyterLab Plugin, allows learners to (1) set personalized learning goals, (2) engage in learning-by-doing through a conversational LLM-based mentor agent, (3) receive guidance and feedback based on a student model that steers the mentor moves. In a within-subject study with 16 participants learning exploratory data analysis from a streamcast, Tutorly significantly improved their performance from 61.9% to 76.6% based on a post-test questionnaire. Tutorly demonstrates the potential for enhancing programming video learning experiences with LLM and learner modeling.

## 1 INTRODUCTION

Video has emerged as a preferred medium for learning programming. According to a recent survey by Stack Overflow, developers

who learned through "How-to-videos" and "Video-based Online Courses" accounted for 60.14% and 49.41%, respectively. Video learning has several advantages, including convenient and asynchronous access [65], flexibility to learn at different paces and focus on selective content [1], and several online content such as on YouTube are free or available at relatively low costs [29]. Regardless, the most effective learning paradigm for programming videos is *learning by doing* [13]. However, this is easier said than done. Videos and supplement materials are often static and spread out across different content formats and tooling (e.g., videos, slides, or GitHub repositories). Further, practicing coding tasks while watching videos is a difficult skill to learn, requiring clear instructions and immediate feedback [16]. While sometimes students are able to interact with teachers who record videos (e.g., online office hours), the process is time-consuming and burdensome[69].

Simulating instructor-led, supportive practice based on video content could enhance the learning experience and efficacy, enabling students to master key concepts presented in the videos more effectively. We propose that, given the generative power of large language models (LLMs), they offer the opportunity to assume the role of a teacher and provide guidance and immediate feedback for students. However, operationalizing LLMs as tutors for video learning requires addressing a number of current challenges. **First, current LLMs suffer from a verbosity bias,** which means that they favor lengthy answers even if they are of poorer quality and harder for people to understand compared to shorter answers [99]. **Second, due to the open-ended nature of LLM text generation, providing step-by-step and targeted guidance is challenging.** It is usually hard for LLMs to teach students step-by-step what they need to learn in a video, just through prompts [52]. In contrast, students benefit significantly from a one-to-one instructor who can choose and utilize pedagogy that suits specific students [90]. As a result, the instructor needs to summarize the knowledge and tasks in the video, teach the knowledge in concise sentences, and guide the students to work on their tasks [66]. Currently, LLM in a fully open-ended setting is not the answer to applying these strategies and sometimes can even be a significant burden for students [76].

To address the challenges of utilizing LLMs as an intelligent *instructor*, we adopt cognitive strategies that play a crucial role in knowledge acquisition and skill learning. Specifically, we draw on the *Cognitive Apprenticeship (CogApp)* framework from the learning sciences as the foundation for our approach [21]. The *CogApp* framework advocates six teaching methods (moves): *Modeling, Coaching, Scaffolding, Articulation, Reflection, and Exploration* that make the instructor's thinking visible to learners and the student's thinking visible to the instructor. For example, a teacher might use *Modeling* to demonstrate the process of solving mathematics problems by having students bring difficult new problems for her to solve in class. *Scaffolding* can take diverse forms, such as suggestions to facilitate writing or physical support in downhill skiing. *Articulation* involves encouraging students to articulate their thoughts as they carry out their problem-solving. The *CogApp* framework also advocates for the selection of teaching methods based on accurate assessments of a student's current skill level and then identifying an intermediate step of suitable difficulty to facilitate the target learning activity [72].

Leveraging these insights, we introduce Tutorly, an LLM-powered apprenticeship learning environment for programming videos. Tutorly is implemented as a JupyterLab [75] extension that supports conversational interactions for scaffolding and guidance in learning from videos. We developed a prompt generation pipeline that segments videos into distinct learning goals and, for each learning goal, selects the appropriate move from the *CogApp* framework by using an internal student model. Rather than allowing LLMs to autonomously choose the method for generating the conversation — a process that can lead to unpredictable outcomes due to the diversity of videos and student models — we adopt a more structured approach. Initially, we extract both procedural and conceptual knowledge from video segments that align with distinct learning goals. Subsequently, we select suitable methods for each type of knowledge, informed by insights from student models. Finally, leveraging these methods, along with corresponding actions and interactions, we craft targeted messages. This systematic process ensures a more controlled and relevant message generation tailored to the specific learning needs of students. As shown in Figure 1, the system progressively generates instructions for students to learn how to visualize the data, starting with MENTORING, then COACHING, and at a later point, SCAFFOLDING and REFLECTION.

We conducted a technical analysis of Tutorly prompting and a user study of Tutorly's end-to-end effectiveness as a system for helping students while watching videos. We found the guidance generated by our prompting pipeline is more concise, consistent with the learning goal and mentor move. Based on our user study, participants demonstrated better learning outcomes with Tutorly. Our approach for tightly integrating *CogApp* framework with LLMs **allows for targeted instruction** on knowledge from videos and **alleviates teaching constraints** through open-ended guidance generation. In summary, Tutorly highlights the potential of generative AI as a virtual instructor for online programming learning, combining pedagogy research on cognitive apprenticeship theory with research on the generative capabilities of LLMs. We contribute:

(1) Tutorly: an interactive conversational system as an extension in JupyterLab. Given the video, Tutorly generates student-specific guidance and lets students use the interactions to practice knowledge according to the learning goal.

(2) *CogApp-based Prompting*: A domain-specific language for organizing cognitive apprenticeship methods, allowing LLMs to generate direction and interactions consistently and stably. It also supports teachers to personalize the combination of methods, interactions, and learning goals.

(3) An evaluation of *CogApp prompting pipeline* and a user study of Tutorly with $N = 16$ participants. Our studies highlight Tutorly's effectiveness in applying cognitive apprenticing pedagogy compared to the status quo of teaching the same material.

## 2 RELATED WORK

The capabilities of LLMs offer novel perspectives to assist students in learning programming while watching tutorial videos. Here, we

synthesize relevant literature on learning by watching videos, cognitive apprenticeship framework, and LLMs in intelligent tutoring systems (ITS), which is the basis for the design of Tutorly.

## 2.1 Video-based Learning

Video-based learning is a pedagogical approach in which learners acquire knowledge and skills through video content, leveraging multimedia content to enhance student engagement and instructional delivery. In the context of programming and software development education, several studies have investigated the efficacy of video tutorials [30, 73, 81, 83]. Online learning platforms like MOOCs and video platforms like YouTube can deliver tutorial videos in a scalable way. However, a significant challenge is this form of learning does not encourage active engagement. Videos provide few opportunities for the viewer to practice and develop their skills unless they are intrinsically motivated [20, 26, 40]. Besides, many learning theories emphasize that in order to become responsible and autonomous learners, students need to take control of their learning [17], express their 'epistemic agency' [44, 82, 85], and participate in creativity-boosting activities alongside explicit instructions on computer programming [3, 102]. Empirical studies also show that the more students engage with the learning content, such as through extensive practice and making mistakes, the easier it is for them to retrieve it [4, 34, 63, 64].

Multiple research has delved into enhancing video-based learning through technology and human-centered design. For example, Guo et al. [33] highlight how video production quality and presentation style significantly affect student engagement in online courses, emphasizing the need for well-designed educational videos to maintain attention and facilitate learning. ToolScape [42] explores the intricacies of extracting step-by-step information from existing how-to videos, suggesting that enhanced metadata and interactive video components can significantly improve video content's usability and educational value. This is complemented by investigating in-video dropouts and interaction peaks in online lecture videos [41], which offers insights into student engagement patterns and learning behaviors. VT-Revolution [10], LV4LP [53], and ITSS [71] also incorporate interactive workflow at the time of video creation. Codemotion [39] and psc2code [9] demonstrated the extraction of executable code from programming tutorial videos, which facilitates active learning by allowing learners to interact directly with the code discussed in the videos.

To maximize student engagement, interactive programming video tutorial authoring and watching systems have been explored [14, 20, 32]. Zhao et al. [98] emphasized the benefits of video question-answering systems for screencast tutorials, which support learners in navigating through complex software learning processes by providing direct answers to specific queries raised during video playback. Soloist [92] uses audio processing to generate mixed-initiative tutorials from existing guitar instructional videos, showcasing the potential for automated educational content creation. Moreover, FlowMatic [96] points toward the future of educational environments on immersive authoring tools for interactive scenes in virtual reality, where learners can manipulate and experiment with the video content, enhancing the depth of interaction and engagement. To align the interactivity features with learning goals, some works enable users to find segments where a certain sub-topic or keyword is presented and allows the learner to click on them and to jump there directly [35, 45, 61, 95]. Our work builds on these prior approaches and leverages the video content to create interactive learning experiences aligned with learning goals.

## 2.2 Cognitive Apprenticeship in Computer Science Education

Knowledge integration, in which learners acquire, modify, and store learned information with what they already know, is central to a successful learning experience. In domains such as programming, it involves combining multiple *representations* into existing knowledge and fostering a unified understanding of a complex domain [54]. To facilitate such understanding, work in instructional design has identified patterns to help learners, such as eliciting, adding, distinguishing, and sorting our ideas through analysis and reflection [55]. A computational cognitive apprenticeship framework for embodying these patterns should focus not only on "how to do," but "how to think" [6]. This approach is in line with cognitive load theory, which suggests that learners manage intrinsic, extraneous, and germane cognitive loads during the learning process [62].

In the realm of computer programming education, cognitive apprenticeship has proven to be effective in enhancing learning outcomes. Studies have demonstrated the advantages of incorporating subgoals in programming education, as breaking down complex tasks into smaller subgoals can decrease cognitive load and improve comprehension [62]. Moreover, providing explicit programming strategies to adolescents has been shown to be beneficial, as it equips them with structured problem-solving approaches and encourages self-regulation throughout the programming process [43]. By engaging in authentic tasks, novices can not only acquire the explicit knowledge needed for a task but also the implicit knowledge that experts possess [57]. Cognitive apprenticeship framework highlights the significance of learning through guided experiences, where novices collaborate with experts to enhance their skills and establish habit in self-regulation [15, 21].

The utilization of technology, such as multimedia content creation tools, can facilitate collaborative learning experiences within the cognitive apprenticeship framework. For instance, studies have looked at anchored instruction [11] that integrates computation within disciplinary engineering practices coupled with scaffolding approaches [59, 60, 91]. Recent work has embedded the methods in computational notebooks to help students learn an unsupervised ML algorithm [86]. Our work implements the computational cognitive apprenticeship framework with subgoal setting, explicit programming instruction, and collaborative multimedia tools within computational notebook environment to facilitate student practice activities and integrate knowledge.

## 2.3 Conversational Intelligent Tutoring Systems

Conversational Intelligent Tutoring Systems (CITS) have emerged as a promising approach to personalized and interactive learning experiences, which guide students through scaffolding practice problems and provide correctness feedback, next-step hints, and adaptive feedback messages. Latham et al. [46] introduce Oscar, an Intelligent Adaptive Conversational Agent Tutoring System that

incorporates human-like natural language interfaces to support constructivist learning styles. My Science Tutor [93] is a conversational multimedia virtual tutor designed to provide individualized and adaptive instruction akin to human tutors, emphasizing the role of ITS in enhancing learning achievement. Furthermore, Ji et al. [37] investigate student and tutor perceptions of conversational ITS in online learning programs, aiming to understand the potential impact of chatbot-based tutoring systems on educational experiences. By leveraging advancements in natural language processing, machine learning, and human-computer interaction, CITS have the capacity to accommodate diverse learning styles and enhance knowledge acquisition in educational settings.

The integration of conversational agents and natural language interfaces in ITS has been a key area of advancement. For example, AutoTutor [31, 68] and Watson Tutor [25] are tutors with dialogue in natural language, emphasizing the role of conversational interactions in enhancing learning experiences. PUMICE [51] is a multimodal agent that learns concepts through natural language and demonstrations, emphasizing the importance of integrating multiple modalities in intelligent tutoring systems for effective communication. Furthermore, SUGILITE [50] integrates GUI-grounded natural language instructions to develop a conversational assistant for Android, showcasing the fusion of interactive task learning with real-world applications. By integrating conversational breakdown repairs, CITS can provide users with a seamless learning experience by addressing communication challenges through multi-modal interfaces [49]. More recently, LLMs capacity to generate coherent responses, understand context, and adapt to user input has opened new avenues for educational applications[38, 58].

The effectiveness of CITS in supporting personalized learning has been a subject of interest in educational research. Zinn [101] discusses algorithmic debugging for intelligent tutoring, focusing on improving diagnosis and problem-solving support within ITS, underscoring the importance of multiple models for enhanced learning outcomes. Additionally, Aljameel et al. [5] evaluate the application of LANA CITS in supporting learning for autistic children, showcasing the potential of Conversational Intelligent Tutoring Systems to cater to diverse learning needs. It has been shown that CITS provides the capability to encourage mentees, especially when learning knowledge within the STEM domain [28]. Moreover, Akyuz [2] explores the impact of Intelligent Tutoring Systems on personalized learning, highlighting the positive contributions of ITS in enhancing student performance and time management. Tutorly builds upon these prior works, providing support for authoring and assessing skill mastery to build student models for adaptive teaching and leveraging advances in LLM applications in chatbots to build conversational capabilities.

## 3 USER EXPERIENCE

Tutorly is an integrated learning environment where learners can watch programming video segments, engage in learning conversations with the AI mentor, and practice writing code. As shown in Figure 2, Tutorly is implemented as an extension to *JupyterLab* [75] — an interactive computational environment for notebooks, code, and data. The extension primarily consists of a rich *multi-modal* conversational chat panel added to the main JupyterLab interface.

The default configuration consists of a notebook panel on the left and the Tutorly chat panel on the right. However, learners can adjust the position of the panels at will. The chat panel and the code cells in the notebook are linked and controlled by Tutorly. To better understand how Tutorly transforms programming videos into a one-on-one mentoring experience, let us follow Leon, a student in a Data Science undergraduate program interested in learning Visual Exploratory Data Analysis (EDA) using the 'R' programming language [36]. Leon has familiarity with R programming and has installed the Tutorly extension on his JupyterLab desktop application using a one-click installer.

### 3.1 Video Context

For this example, Leon will use the YouTube video "Tidy Tuesday Screencast: analyzing college major & income data in R [1]" by Dave Robinson. To provide context, Dave Robinson is a well-known figure in the data science community, particularly known for his involvement with the TidyTuesday project and his screencasts on YouTube [78]. TidyTuesday [22] is a weekly data project aimed at the R community that encourages participants to apply their data wrangling and visualization skills to a new dataset every week. While not intended as tutorials, Dave regularly publishes his live coding sessions on YouTube, which embodies his deep expertise in data science. Leon has attempted to learn by directly watching the video, but self-regulation [100] – in which the learner actively manages their own learning experiences through planning, monitoring, and evaluating their understanding – has been challenging. Our goal is for Leon to leverage the expert knowledge in the video to improve his EDA skills using the apprenticeship model approach.

### 3.2 Apprenticeship Learning with Tutorly

Upon loading the chat panel, the Tutorly chatbot prompts Leon to upload a configuration file for the topic he wishes to learn along with links to the video and code (we detail the file format and creation in Section 4.5). Based on the configuration, learning EDA involves the following *learning goals*: gaining familiarity with the dataset, data cleaning, forming hypotheses or visualization intent, visualization construction, visualization refinement, and visualization interpretation and insight generation [88]. Using the configuration Tutorly instantiates a student model for Leon with a skill level of 'novice' for all of the goals.

To get the tutoring session started, Tutorly identifies that gaining *familiarity with the data* is the first relevant segment in the video. Since this is also the first time Leon is encountering this goal, Tutorly takes the *MODELING* move and presents the video segment in which Dave explores the dataset using think-out-loud as a chat response. For goals involving *declarative knowledge* such as facts, **playing the video segment** in which the expert (Dave) is exploring the dataset is a suitable move. Tutorly instructs Leon to watch the short video clip. After learning about the dataset characteristics from the video, Leon clicks on the "Go On" button, indicating that he is ready to proceed. The next step in the video is creating a visualization to explore the relationships between *College Major* and *Income*. Since this step involves *procedural knowledge* by writing code to create a box plot, the modeling move involves Tutorly

---

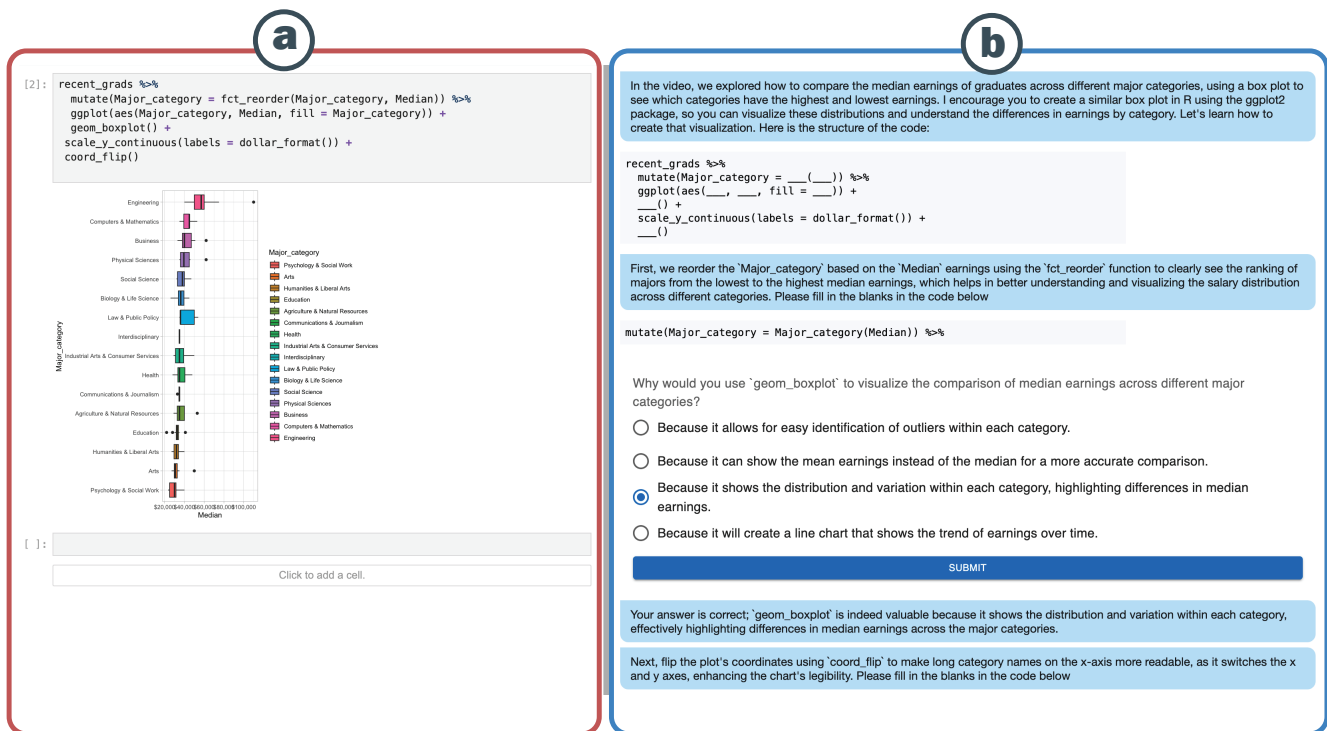[1]https://www.youtube.com/watch?v=nx5yhXAQLxw

**Figure 2: Tutorly user interface as a JupyterLab Extension. (a) Code panel with cell structure, (b) Tutorly Chat Interface with Rich Conversational Features**

presenting the learner with the code snippet and walking them through each of the functions and parameters. Tutorly has access to Dave's completed code along with the transcript, which is used to generate the correct code. Note that in the process, Tutorly also automatically adds the code to the notebook cell. At the end of the **code walkthrough**, Tutorly chatbot instructs Leon to execute the code and generate the visualization. Further, for visualization interpretation, this move again plays the relevant video segment.

As the video progresses, Leon encounters a second instance of a bar chart visualization. In this case, Tutorly takes the *COACHING* move and provides Leon with a **code template** for the bar chart with key functions and parameters left blank. Then, based on the visualization intent, Tutorly guides Leon in filling out the blanks through conversation and interactivity. For instance, the blanks in the code template are clickable, and Leon can select from the dropdown list of options. Tutorly provides feedback in case of incorrect selection. Once completed, Tutorly adds the code to the notebook, which can be executed. For visualization interpretation, the coaching involves specific guidance to Leon on what to attend to in the visualization. *SCAFFOLDING* and *ARTICULATION* moves for declarative knowledge involves **multiple choice questions**, asking Leon to pick the right answer, and providing feedback. These moves for procedural knowledge could involve Tutorly asking Leon to directly write the code or interactively explore the dataset to formulate hypotheses and provide **feedback**. As Leon is writing code in the notebook, Tutorly evaluates the code for errors and assists Leon in **fixing code errors**. Alternately, in *REFLECTION*

move, Tutorly can ask the student to assess the issues on their own. Lastly, *EXPLORATION* allows Leon to go beyond what is contained in the source video, such as generating and testing hypotheses that are not covered in the video. Throughout these moves, the Tutorly constantly updates the underlying student model to capture the content and learning goals Leon has encountered.

In summary, Tutorly takes the source video and applies different mentoring moves from the *CogApp* framework and interactively walks Leon through the video. In the learning conversations with Tutorly, Leon takes a learning-by-doing approach to develop his EDA and visualization skills and receives active feedback across different EDA tasks. In the process, Leon can also ask clarifying questions and work at his own pace to regulate his learning.

## 4  SYSTEM ARCHITECTURE

Tutorly implements a set of LLM prompting techniques that yield conversational utterances grounded in the *CogApp* framework. As shown in Figure 3, our system consists of (1) a video segmentation module that *slices* long videos based on learning goals using a prompt chaining approach, (2) a DSL generator that produces the *context* for the LLM to generate utterances in interaction with the learner, and (3) a student model to capture the student's learning progress and inform appropriate mentor *moves*. Here, we provide technical details about each of these modules along with challenges and prompt engineering strategies to make the generated conversations **consistent** and **controllable**:
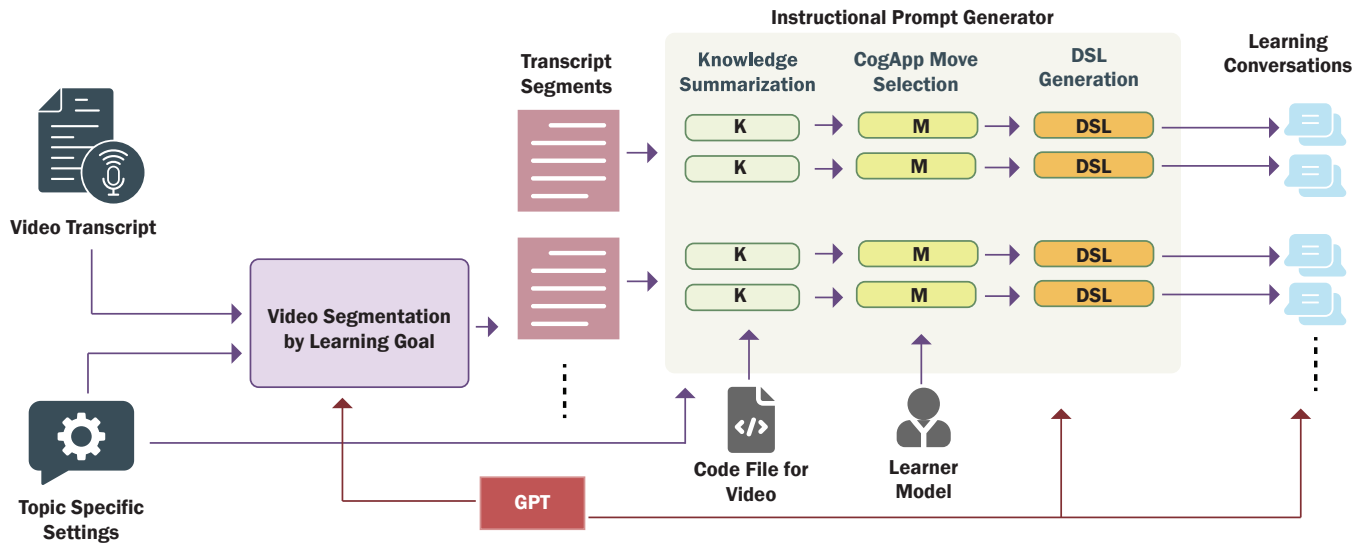
**Figure 3: Overview of Tutorly's System Architecture. From Left to Right: The video Transcript is segmented based on learning goals, and then for each segment, the knowledge is summarized based on knowledge type. Next the appropriate *CogApp* mentor move is selected for the knowledge and using the student model. Finally a DSL representation is produced to prompt GPT for generating the conversation.**

## 4.1 Video Segmentation by Learning Goals

To implement the *CogApp* framework, we need to be able to distinguish specific learning goals in the video and map them to mentoring strategies. While GPT-4 has strong summarization capabilities [97], at the time of this writing, GPT-4 could not yet directly perform long-context semantic classification. First, when segmenting the transcript by just providing the learning goals, we observed that even small references to other learning goals within a main segment (e.g., mention of data or hypothesis in segmentation about visualization construction) led GPT to distinguish it as having two learning goals, thus affecting the classification results. Additionally, current transcripts are segmented and timestamped based on pauses in speech or changes in speakers. While advanced ASR attempts to improve the readability of transcripts with complete sentences and the use of punctuation, it can be challenging and is not always accurate. These challenges in transcription impact the inductive capabilities of large language models.

To address these challenges, we start by defining each of the learning goals as a one-shot prompt using a single descriptive example as context. For instance, if the goal is *chart interpretation*, we provide a short description that it means instances in which "the expert interprets the visualizations and discusses the implications of the visualization, drawing conclusions, and theorizing about the underlying trends or patterns in data." In our iterative trials, GPT is able to reasonably generalize and infer goals based on these descriptions. In addition, we provide few-shot examples of the segmented videos to produce distinct summaries of multiple occurrences of the learning goals in the video. Given this context, we utilize the chain-of-thought (CoT) prompting technique [94] to define the *Video Segmentation by Learning Goal* (VSLG) algorithm defined as follows:

Let $V = \{v_1, v_2, ..., v_T\}$ be the set of all sentences in the video transcript, where $T$ is the total number of sentences, and $L = \{l_1, l_2, ..., l_M\}$ be the set of learning goals. The transcript is segmented into $S = \{s_1, s_2, ..., s_N\}$, where $s_n$ corresponds to one segment for one of the learning goals $l_m$. $N > M$, so that each learning goal has multiple segments. We define three functions executed sequentially in Algorithm 1, which mimics the cognitive steps an individual might take when classifying video clips semantically and provides a structured methodology for LLMs to segment video content by learning goals.

---

**Algorithm 1** VSLG - Video Segmentation by Learning Goal

---

**Require:** Video transcript sentences $\mathcal{V}$, Learning goals $\mathcal{L}$
**Ensure:** Video segments by learning goals $\mathcal{S}$
1: $\mathcal{S}' = \{s'_1, s'_2, ..., s'_N\} \leftarrow$ **Summarize**$(\mathcal{V}, \mathcal{L})$
    ▷ summarizes transcript of each learning goal, yielding a set of summary points of all segments without timestamps.
2: $\mathcal{R} \leftarrow$ **Retrieve**$(\mathcal{S}', \mathcal{V})$
    ▷ retrieves sentences from the video transcript for each summary point $s'_n$.
3: $\mathcal{S} \leftarrow$ **Rearrange**$(\mathcal{R})$
    ▷ determines the range of segments based on the timestamps of sentences and sorts them by time.

---

Below is a snippet of the segmentation process using this approach ("start" and "end" are the start time and end time of this segment in seconds):

**Input: [video transcript]**
**[After the first step]**
```
[...
```

```
    {"category": "Visualize the data",
     "summary": "David decides to explore the
         median salaries by creating a
         histogram to understand the
         distribution of median earnings across
          majors.'},
 ...]
```
**[After the second step]**
```
[...
    {"category": "Visualize the data",
     "sentence": "alright so I take a look at
         every synchros now that it picked
         something I'm interested in...'},
 ...]
```
**[After the final step]**
```
[...
    {"category": "Visualize the data",
     "start": 435.23, "end": 461.93},
 ...]
```

## 4.2 Instructional Prompt Generator

The next step in the pipeline involves generating prompt *templates* and associated *context* for each video segment. The prompts need to align with different moves in the *CogApp* framework to drive the conversational interactions with the learner later. Applying *chain-of-thought*, if an instructor were teaching, they would begin by determining the specific knowledge or content that the goal should encompass [74]. Then, they would strategically determine pedagogical approaches for different learner needs. Similarly, our prompting strategy is first to use the segmented video transcripts and code to summarize the knowledge to be learned, then determine the teaching method based on the knowledge type for each of the moves in *CogApp* Framework. Finally, we map the method to predetermined actions and interaction types. By enforcing a step-by-step generation process, we emulate instructional practices and gain fine-grained control to carry out the conversation with learners progressively. Here, we detail each of these steps in generating instructional prompt templates:

*4.2.1 Summarize the knowledge in each video segment.* Grounded in the widely used revised Bloom's Taxonomy [7, 80], we focus on *declarative* Knowledge (i.e., propositional facts, information, and descriptions ) and *procedural* knowledge (i.e., knowledge about processes and sequences of actions) as they are most common in programming education. Through trial and error with different levels of knowledge abstractions, being too precise leads to overlapping information in the summaries, making downstream tasks more complex and less generalizable. Further, we find that in programming videos, declarative and procedural knowledge is spread across domain *concepts* (e.g., objects in object-oriented programming, dataset attributes) as well as *programming* constructs (visualization, sorting). Using this categorization, we use few-shot prompts to generate procedural and declarative knowledge summaries in each of the video segments. Additionally, as shown in Table 1, for each knowledge-content category, we define specific summary formats. First, this ensures that we are able to map the

knowledge to build and update the student model more accurately (Section 4.3). Second, the summary format contains the necessary information for GPT to generate conversational messages without making the knowledge content too long or too short. Third, the format allows us to link the transcript to the code. If the knowledge summary is vague, for instance, *"Filter the dataset to select the top 20 rows to focus on the highest-earning majors"*, it becomes difficult to map to functions and attributes in the code for downstream interactions. We designed the summary formats iteratively through design and testing.

*4.2.2 Select Pedagogical Strategy for each piece of knowledge.* Given the knowledge summary, the next step is for the LLM to determine the instructions grounded in the *CogApp* Framework. If LLMs have a high degree of freedom, i.e., zero-shot prompting by giving it the framework and letting it plan the instructions, it is challenging to maintain logical consistency and coherence [56]. In our trials, we observed that with this approach, the LLM might use the *REFLECTION* move even in the early stages and regardless of the knowledge to be learned. On the other hand, rigidly mapping specific teaching methods to particular knowledge can constrain Tutorly's ability to adapt to unforeseen content and interactions, and the monotony may lead to disengagement. Our goal is to support a dynamic instructional approach within the boundaries of *CogApp* frameworks.

To do this, our prompt engineering follows the following three principles in the *CogApp* framework to select an appropriate mentor move: **(1) Global before local skills.** The recommendation is to have learners build a conceptual map before attending to the details of the terrain. To comply with this principle, Tutorly chooses *Modeling* or *Coaching* as the first move to describe the tasks explored in the video segment or demonstrate the end-to-end process of solving the tasks. This move would then be followed by other methods to guide students in using more advanced problem-solving strategies. **(2) Increasing complexity.** With this principle, the difficulty of the task gradually increases, so each new task introduces additional layers of skills and concepts essential for mastery. We implement this through a combination of prompt engineering and student modeling, described later. **(3) Increasing Diversity.** This means providing students with a broad range of problems and contexts, which compels them to apply their skills and strategies in different ways [12]. To optimize the learning experience, it is essential to diversify teaching methods, even when students are engaging with different video segments that share the same learning goals. For example, if a student learned how to make a histogram by answering a multiple-choice question in the past video segment, next time, the Tutorly will ask them to fill in the code blanks to make a box plot. Varying the instructional approach could maintain student interest and cater to different learning styles.

The process is detailed in Algorithm 2. Here is one example in the prompt (additional details in the Appendix):

**Input: [Knowledge of a concept-related segment]**
**Output:**
```
{"knowledge": "To interpret the differences
     in median income by college major...",
 # Start with Modeling or Scaffolding
 "actions": [{"method": "Scaffolding"}]},
```

| Knowledge | Representation | Example |
|-----------|---------------|---------|
| *Concept Related* | | |
| Declarative | [Subject] + [verb phrase] + that + **[independent clause]**. | *The median income by college major shows that* **majors earn a median income of over $30K right out of college.** |
| Procedural | To achieve/understand + [specific goal/outcome] + one must + **[actions/processes]** + [additional details] + considering/using + [relevant factors/tools]. | *To understand the distribution of earnings by college major, one must* **examine the histogram and identify overall trend or extreme values**, *considering whether high earnings are due to the field's financial reward or influenced by factors such as low sample size and high variation.* |
| *Programming Related* | | |
| Declarative | The task is + [final goal] + using + [general method/tool] + and + [additional method/technique for enhancement]. | *The task is comparing the distribution of median earnings across different major categories using a box plot and adjusting the visualization for better readability and interpretation.* |
| Procedural | To achieve + **[specific goal]** + one must + **[action/verb] + [specific tool/method]** + on + [object/target] + because + [reason/purpose]. | *To* **achieve an ordered factor level** *based on the 'Median', one must* **use 'fct_reorder'** *on 'Major_category', making it easier to compare distributions.* |

**Table 1: Representations of the two types of knowledge in Tutorly. Each type of knowledge has concept-related and programming-related representations. The bold parts are later used in the student model.**

```
{"knowledge": "The chart on median income by
    college major contains...",
 # Increase diversity by using more methods
 "actions": [{"method": "Articulation"},
            {"method": "Reflection"}]},
{"knowledge": "To draw final conclusions from
     the chart...",
 # Increase complexity followed by Reflection
 "actions": [{"method": "Coaching"},
            {"method": "Reflection"}]}
```

*4.2.3 Generating a DSL for multi-turn conversation.* Given a knowledge summary and pedagogical move, we still need a reliable way to operationalize it in generating Tutorly chatbot conversations. Zero-shot prompting will not work. First, teaching moves are somewhat abstract principles and, without specificity and context, can lead LLMs to produce lengthy and unfocused text, hindering learning effectiveness. Second, the conversation should be generated more actively and in a step-by-step manner. Typically, chatbots respond reactively to user inputs rather than initiating a dialogue that steers the learning process. Further, the timing of messages is also critical: some should be delivered promptly to maintain the learning flow, while others should be contingent upon the learner's actions to ensure relevance and engagement.

To support these learning conversation needs, we developed a domain-specific language (DSL) that encompasses all the necessary information for directing chatbot interactions within a learning scenario. As a prerequisite, for each pedagogical move, we need experts to provide at least one action type and response interaction for programming-related and concept-related knowledge (see Section 4.5 for details). Our DSL incorporates information from earlier processing steps and the action information in the pipeline to orchestrate the conversation. As shown in Algorithm 3, we first map the expert-provided *actions* and *interactions* for each pedagogical move. Then, we used few-shot prompting to generate a *prompt* and *need-response* based on each pair of action-interaction, which is used to drive the LLM to generate an utterance and detect whether it needs the user to respond to generate the next utterance. The *parameters* is a list of parameters required in the "prompt", such as knowledge, code block, student's answer, etc. The sequence of messages depends on the knowledge and the actions that help students learn the knowledge. They are stored in a queue, and each time a message is generated, the head of the queue is removed. The full structure can be formed before uploading to the chatbot, thus avoiding calling the prompt pipeline to save time during the conversation. The algorithm and part of an example DSL structure is shown below:

```
"Visualize the data - 509": [
{
  "knowledge": "Use `geom_boxplot` on ...",
  "actions": [
  {
    "method": "Coaching",
```

---

**Algorithm 2** Knowledge Summary and Pedagogical Move

---

**Require:** Video segments $S$, Cognitive Apprenticeship framework $\mathcal{F}$, code blocks $C$, student model for each knowledge $p_k$
**Ensure:** Summarized knowledge $K$ and teaching methods $M$
1: Initialize knowledge list $K \leftarrow \emptyset$
2: Initialize teaching method list $M \leftarrow \emptyset$
3: **function** SUMMARIZEKNOWLEDGE($S$)
4:     **for** each segment $s \in S$ **do**
5:         $K_s \leftarrow$ Generate knowledge snippets for $s$ with $C_s$
6:         Sort $K_s$ in order of appearance within $s$
7:         $K \leftarrow K \cup K_s$
8:     **end for**
9:     **return** $K$
10: **end function**
11: **function** PLANTEACHINGMETHODS($K, C$)
12:     **for** each knowledge $k \in K$ **do**    ▷ Apply the principles
13:         Determine the index, complexity, diversity of $k$
14:         $M_k \leftarrow \emptyset$ ▷ Initialize the set of methods for knowledge $k$
15:         **if** $index <= 1$ **then**
16:             $M_k$.add(*Modeling/Scaffolding*)
                ▷ Start with global skills and task overviews
17:         **else if** $p_k <= 0.5$ **then**
18:             $M_k$.add(*Scaffolding/Coaching/Articulation*)
                ▷ Depends on which one has not been used
19:             $M_k$.add(*Reflection*)    ▷ Encourage self-assessment
20:         **else if** $p_k > 0.5$ **then**
21:             $M_k$.remove(*Scaffolding*)
22:         **end if**
23:         $M \leftarrow M \cup M_k$
24:     **end for**
25:     **return** $M$
26: **end function**
27: $K \leftarrow$ SUMMARIZEKNOWLEDGE($S$)
28: $M \leftarrow$ PLANTEACHINGMETHODS($K, C$)

---

**Algorithm 3** DSL Structure for Chatbot Message Generation

---

**Require:** Methods Set $\mathcal{M}$, Actions $\mathcal{A}$, Interactions $\mathcal{I}$, Parameters list $\mathcal{P}$, student behavior $\mathcal{B}$, knowledge $K$, student model $p_K$
**Ensure:** Chatbot interaction
1: **function** GETDSL($a, i$)
2:     **Description**: Generates DSL components for $\mathcal{M}$.
3:     $prompt \leftarrow$ Use $i$ to generate the prompt for action $a$
            ▷ "Use [interaction] to + definition of the method"
4:     $needResponse \leftarrow$ Determine if requires user response
5:     $\mathcal{P} \leftarrow$ Extract parameters required for the prompt
6:     **return** $\{prompt, needResponse, \mathcal{P}\}$
7: **end function**
8: Initialize message queue $Q$
9: **for** each $m \in \mathcal{M}$ **do**
10:     **for** each $(a, i) \in \mathcal{A} \times \mathcal{I}$ **do**
11:         $\{prompt, needResponse, \mathcal{P}\} \leftarrow$ GETDSL($a, i$)
12:         Enqueue $\{m, a, i, prompt, \mathcal{P}, needResponse\}$ into $Q$
13:     **end for**
14: **end for**
15: **while** not $Q$.empty **do**
16:     $message \leftarrow Q$.dequeue
17:     **Chatbot**.send($message.prompt, message.\mathcal{P}$)
18:     **if message**.$needResponse$ is **false then**
19:         **continue**
20:     **else**
21:         **Chatbot**.wait($\textbf{Signal}_{response}(t)$)
22:         $p_K \leftarrow$ **UpdateStudentModel**($\mathcal{B}, K, p_K$)
23:     **end if**
24: **end while**

---

```
    "action": "Prompt the student to use {
        interaction} to practice the
        knowledge",
    "interaction": "fill-in-blanks",
    "prompt": "[Use one sentence to prompt
        the student to fill in the {code-line
        -with-blanks} below]",
    "parameters": ["code-line-with-blanks"],
    "need-response": true
  },
  ...
  ]
},
...
]
```

## 4.3 Student Behavior Monitoring

A novel aspect of Tutorly is that the conversation is not only driven by the video context but also by modeling the learner. Here, we cover key details on student modeling, including the timing of

behavior evaluation, updating student models, and the influence of these models on the prompt pipeline:

*4.3.1 Monitor Trigger Conditions.* Based on literature [8, 89] and trial and error, we identify a set of signals that updates the knowledge state of the student model. The signals are intended to ensure optimal monitoring frequency; excessive instructor intervention can distract students and slow system responsiveness, while infrequent updates to student models may delay crucial learning feedback. We implement a signal function **Signal**($t$). When it is activated, Tutorly updates the student model. We define the four monitor trigger conditions as follows:

- *Post-viewing Guidance*: Initiated by $\textbf{Signal}_{video}(t)$, where the instructor provides guidance based on the video.
- *Progression Cue*: Triggered upon completion of an instruction $\textbf{Signal}_{response}(t)$, signaling the next step of guidance.
- *Corrective Feedback*: Activated by $\textbf{Signal}_{error}(t)$, when the code execution fails, prompting improvement suggestions.
- *Responsive Assistance*: In response to $\textbf{Signal}_{help}(t)$, the instructor answers questions or provides additional help.

*4.3.2 Student Model Update.* Although many ITS use question-answer pairs to test student learning outcomes [24], evaluating students' performance on specific problems cannot directly represent students' knowledge models [84]. Utilizing a part of the **Knowledge** as outlined in Table 1, we construct a statistical model

to emulate the student's mental state using Bayesian Knowledge Tracing (BKT) based on a Hidden Markov Model. Each knowledge component is characterized by four parameters: $p_{mastery}^{(t)}$, $p_{transit}$, $p_{slip}$, and $p_{guess}$. These parameters are defined as follows:

- $p_{mastery}^{(t)}$: The probability that the student has mastered the knowledge at time $t$.
- $p_{transit}$: The probability of the student transitioning to a state of knowledge after an opportunity to apply it.
- $p_{slip}$: The probability the student makes an error when applying known knowledge.
- $p_{guess}$: The probability the student correctly applies unknown knowledge.

Our system updates the student model using Equation 1, based on the observed student behavior obs at time $t$ when **Signal**$(t) = True$

$$p_{mastery}^{(t|obs)} = \begin{cases} \dfrac{p_{mastery}^{(t)} \times (1-p_{slip})}{p_{mastery}^{(t)} \times (1-p_{slip}) + (1-p_{mastery}^{(t)}) \times p_{guess}} & \text{obs = correct} \\[2em] \dfrac{p_{mastery}^{(t)} \times p_{slip}}{p_{mastery}^{(t)} \times p_{slip} + (1-p_{mastery}^{(t)}) \times (1-p_{guess})} & \text{obs = incorrect} \end{cases}$$
$$\tag{1}$$

Tutorly uses this formula to update the student model throughout the learning process. Every time the student is practicing a piece of knowledge, it is recorded and initialized using the BKT parameters. If the student is practicing knowledge that has been practiced with similar knowledge before, the BKT parameters of the existing knowledge will be updated. We use semantic similarity to evaluate whether similar knowledge exists [67]. During the practicing process, if students respond correctly, such as choosing a correct choice in the multiple-choice question or writing the correct code line, the observation is *Correct* and Tutorly accordingly updates the parameters of that correctly practiced knowledge. On the contrary, if students respond incorrectly, Tutorly uses the formula for incorrect to update the parameters of the incorrectly practiced knowledge. Each time a student completes a learning session, BKT parameters are saved to the database. When students watch a new video, Tutorly will load previously saved parameters for the new learning course. The whole process is shown in Algorithm 4.

*4.3.3 Integration into the Prompt pipeline.* The student model is a direct reflection of the student's current knowledge and understanding. Based on prior work [23], our pipeline integrates the student model for the following three tasks:

- Selecting learning goals for students showing weak mastery, indicated by $p_{mastery}^{(t)} < 0.3$.
- Reducing the recurrence of knowledge components already mastered, denoted by $p_{mastery}^{(t)} > 0.7$.
- Adapting teaching methods according to the student's knowledge mastery; for instance, favoring *Coaching* over *Modeling* or *Scaffolding* if $p_{mastery}^{(t)} > 0.5$ for knowledge "plotting a histogram to see the distribution".

## 4.4 Implementation Details

Tutorly is implemented as a server-client architecture. The front end is a JupyterLab plugin [75] built using TypeScript and React [27].

---

**Algorithm 4** Update Student Model

---

**Require:** student behavior $\mathcal{B}$, knowledge $K$, student model $p_K$
**Ensure:** the Student model gets updated
1: **function** UPDATESTUDENTMODEL($\mathcal{B}, K, p_K$)
2:     **Description**: Updates the student model.
3:     Wait for **Signal**$(t)$ before proceeding
4:     **for** each knowledge $k \in K$ **do**
5:         **if** $k$ is semantically similar to an element in $K$ **then**
6:             $p_k \leftarrow$ the parameter for $k$ in $p_K$
7:             **if** $\mathcal{B}$ is correct **then**
8:                 Update $p_k$ based on correct behavior
9:             **else**
10:                Update $p_k$ based on incorrect behavior
11:             **end if**
12:         **else**
13:             Create new parameter $p_k$ for $k$
14:             **if** $\mathcal{B}$ is correct **then**
15:                Initialize $p_k$ based on correct behavior
16:             **else**
17:                Initialize $p_k$ based on incorrect behavior
18:             **end if**
19:         **end if**
20:     **end for**
21:     $p_K \leftarrow p_K \cup p_k$
22:     **return** updated $p_K$
23: **end function**
24: $p_K \leftarrow$ UPDATESTUDENTMODEL($\mathcal{B}, K, p_K$)

---

We make use of the provided APIs to integrate the plugin with the main computational notebook including @jupyterlab/notebook, @jupyterlab/cells, and @lumino/signaling. The chat UI design uses the Chat UI Kit from *chatscope* [19], an open-source UI toolkit for developing web chat applications based on React. In addition to basic incoming-outgoing message functions, it also provides timestamps, "is typing," and other functions in the interface. The backend is built using Python and uses OpenAI's API for all language model functions, specifically the GPT-4 and GPT-4-32k model [70], which are OpenAI's most capable language model at the time of writing. We use a temperature of 0.2 for mentioned prompts. We use the conversation buffer memory implemented by LangChain [2] to store messages during the chat. The video transcript and other necessary information are fetched using YouTube's API. The original R Markdown code files and dataset files of videos are fetched using GitHub's API. We use SQLite for the database design. We use the *Bilingual and Crosslingual Embedding for RAG* model [67] for text classification and semantic similarity. All server functions are implemented to ensure compatibility with the JupyterLab extension.

## 4.5 Expert Inputs

As mentioned above, Tutorly relies on expert inputs to define learning goals, actions, and interactions for each of the pedagogical moves. Since these are dependent on the topic content, we abstract them out from the base implementation. This allows for the generalizability of our approach. We developed a simple web-based utility

---

[2]https://python.langchain.com/docs/modules/memory/types/buffer

(Figure 4) to allow experts to input and configure these aspects of Tutorly. In our current implementation, interactions are specified using well-known interventions such as "multiple-choice" and "fill-in-blanks." Future work can look at developing more specialized interventions. Additionally, in the utility interface, the expert can turn on or off different learning goals and conversational behavior. For instance, in a classroom setting, the instructor can configure different learning goals for students based on their needs. Based on the user experience walkthrough, if Leon excels at chart interpretation but needs to focus on creating a visualization, the instruction can turn off the interpretation goals. During active learning, Tutorly will not generate conversations for chart interpretation. The configuration will be used in the pipeline to generate the DSL.
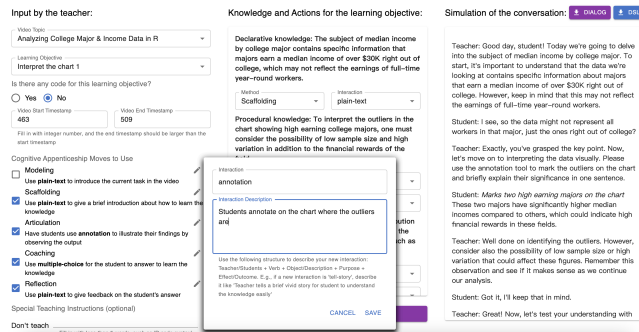


**Figure 4: Tutorly Prompt Configuration and Testing Utility**

## 5 TECHNICAL EVALUATION

We validate our approach for three different topics (EDA, Machine Learning, and Game Development) to measure both segmentation accuracy and the generated messages for different moves in the *CogApp* framework. To do this, we set up a prompting pipeline by extending the utility tool to simulate both the Tutorly chatbot and learner. Then, we ask experts to score the conversations for accuracy along the following two dimensions:

(1) **Segmentation Accuracy:** Prompting must correctly categorize the video transcript into segments that correspond to the learning goals.

(2) **Controllability:** Users can adjust the knowledge, methods, actions, and interactions in the prompting pipeline. Conversation must correctly follow the user's settings to demonstrate the controllability of the pipeline.

### 5.1 Procedure

The lead author, leveraging their expertise from reviewing numerous exploratory data analysis videos and familiarity with learning goals formulation, serves as the annotator for segmenting videos according to learning goals. Three additional experts recruited through our connections, proficient in each of the three video topics, served as evaluators to measure the accuracy of the intent.

*5.1.1 Evaluating Segmentation.* The lead author watched three 20-minute videos and manually segmented them using the same definitions of learning goals given to GPT. Then, we compared

each segment's timestamps between the labeled data and generated data. Considering that the transcript given to GPT is discrete and the information received by watching the video is continuous, we analyzed accuracy such that the timestamp differences are within five seconds error.

*5.1.2 Evaluating Controllability.* We created an initial set of dialog intent labels (Table 2) that have a hierarchical classification scheme. This scheme also represents the hierarchical prompt pipeline. The naming rules of labels are based on the previous Pair Programming Labels [77]. The ground truth data are extracted from the DSL files, and their expression is consistent with the (Table 2). The hierarchy begins with the root Knowledge classifier that categorizes knowledge into procedural and declarative knowledge. The second layer is teaching Methods defined in the cognitive apprenticeship theory [21]. We excluded *EXPLORATION* in the experiment because it would add open-ended dialogues that users could not control. In the third layer, we included all the interactions used to simulate conversations. Finally, we have four types of actions that represent the four most common behaviors in our programming teaching process [77].

Three expert evaluators label the data independently and separately for three topics of video (each topic has two videos): exploratory data analysis (124 utterances), machine learning (89 utterances), and game development (64 utterances). Each message is annotated with four labels. For example, the evaluator would annotate this message: "Now, let's delve deeper. I have a multiple-choice question for you: What could be the potential reason behind the pattern of high earnings for certain majors? A) The field's inherent financial reward, B) Low sample size, C) High variation in the data, or D) All of the above." with the following label: "Declarative - Articulation - multiple-choice - Comprehension". Using this method, the evaluators annotated a total of 277 dialogue messages from the instructor. Finally, we compared the labeled data with the ground truth data with the precision, recall, and F1-score metrics to evaluate the system's accuracy, providing an understanding of how well the system's output aligns with the expected results.

### 5.2 Results

*5.2.1 Tutorly has acceptable segmentation results.* Video transcript segmentation by learning goals algorithm achieved an overall accuracy of 73.7% within a five-second margin, which is an acceptable performance for educational content. However, the observed trend of decreasing accuracy with longer video duration suggests that the algorithm's effectiveness in identifying precise segmentation points diminishes over time. This observation may be attributed to the increasing complexity and variability of the content as the video progresses, making it challenging to maintain a consistent level of summarizing the video content and retrieving the relevant sentences. To alleviate this problem and improve segmentation accuracy, it could be helpful to split the video into smaller clips (e.g., every 10 to 12 minutes) and then segment by learning goals. This approach will limit the amount of content the algorithm needs to process at once, potentially reducing context load and allowing for more accurate identification of learning goals and corresponding segmentation points.

**Table 2: Prompt Pipeline Labels**

| Knowledge | Example |
| --- | --- |
| Procedural | To interpret the chart, one must analyze the histogram to identify outliers. |
| Declarative | Majors earn a median income of over $30k right out of college. |

| Methods | Example |
| --- | --- |
| Modeling | I encourage you to follow along with me to create a meaningful visualization. |
| Coaching | Now, please fill in the blanks to add the box plot layer to our ggplot. |
| Scaffolding | You'll need to sort the data by median earnings in descending order. |
| Articulation | Can you tell me why it's important to consider these factors? |
| Reflection | Please compare your code with the standard code block. |

| Interactions | Example |
| --- | --- |
| plain-text | I'd like you to tell me what you observe about the high-earning majors. |
| multiple-choice | What could be the potential reason behind the pattern? Here are the options. |
| fill-in-blanks | Now, can you fill in the blanks to reorder the Major_category based on Median? |
| show-code | Please execute the cell to see the visualization. |
| annotation | Please use annotation to mark the specific areas on the chart. |

| Intent | Example |
| --- | --- |
| Task Control | We're going to focus on interpreting a chart related to median income by major. |
| Comprehension | Visualizing the distribution of median salaries helps us to see the spread of earnings across different majors, which can be very insightful. |
| Code/Run Code | The correct line should be recent_grads %>% ggplot(aes(Median)). |
| Feedback | You've done well in understanding the importance of each element in the plot! |

**Table 3: Comparison of the intent performance across three video topics: exploratory data analysis (EDA), machine learning (ML), and game development (Game) using the precision, recall, and F1-score metrics.**

| Topic | Knowledge | | | Method | | | Action | | | Interaction | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Prec. | Recall | F1 | Prec. | Recall | F1 | Prec. | Recall | F1 | Prec. | Recall | F1 |
| Total | 0.791 | 0.787 | 0.789 | 0.814 | 0.809 | 0.807 | 0.902 | 0.895 | 0.896 | 0.970 | 0.968 | 0.968 |
| EDA | 0.810 | 0.806 | 0.808 | 0.849 | 0.815 | 0.818 | 0.896 | 0.871 | 0.872 | 0.984 | 0.984 | 0.984 |
| ML | 0.795 | 0.798 | 0.796 | 0.825 | 0.809 | 0.807 | 0.900 | 0.899 | 0.899 | 0.956 | 0.944 | 0.947 |
| Game | 0.827 | 0.781 | 0.792 | 0.808 | 0.813 | 0.798 | 0.954 | 0.953 | 0.953 | 0.973 | 0.969 | 0.969 |

*5.2.2 Tutorly is adept at interpreting user intent.* As shown in Table 3, the *Total* row provides an aggregate view of Tutorly's performance across all video topics while the other rows contain metrics of each video topic. The Knowledge column reflects Tutorly's capacity to extract procedural and declarative knowledge from the videos. This is a critical first step as it sets the foundation for subsequent instructional design. Precision and recall scores for different video topics are around 0.791 and 0.787, indicating that Tutorly is relatively accurate in generating messages according to the given knowledge, without containing misleading information (high precision), and without missing important content (high recall). The slightly lower performance (recall equals 0.781) in game development videos may suggest the presence of more diverse or ambiguous content in such videos, which poses a greater challenge for knowledge extraction.

The *Method* column represents Tutorly's ability to generate information that is well aligned with a defined pedagogical move. Precision and recall are also consistent for this category (difference between highest Recall to lowest is 0.006), with EDA videos showing better performance. This higher score suggests that Tutorly is more effective at following the given methods to generate appropriate messages in an EDA environment, which may be due to the inherently interactive and action-oriented characteristics of EDA tasks. The *Action* column metrics reflect how Tutorly translates interactions into specific messages or commands. This step also shows high accuracy (over 0.900), indicating that Tutorly can reliably trigger intended actions based on user's setting. The higher scores in game development (0.954) suggest Tutorly's action messaging is particularly effective at identifying factual content and responding to action-driven commands in game scenarios, which may involve a mix of dynamic and static content requiring precise system responses.

In the *Interaction* column, the high scores (all around 0.950) signify Tutorly's proficiency in generating messages based on the selected interactions. This step is crucial as it determines how students will engage with the material and how users can personalize

interactions based on their needs. The near-perfect performance on all three topics reflects Tutorly's robustness in seamlessly converting specified interactions into actual conversational utterances.

Along the process of Knowledge-Method-Action-Interaction, the scores show an upward trend (e.g., precision is 0.791-0.814-0.902-0.970 respectively). The ascending scores across these stages can be understood as a natural outcome of the instructional guidance generation process, which moves from the broad and general to the specific and interactive. At each step, Tutorly leverages the work done in the previous phase, refining and sharpening its outputs.

# 6 USER STUDY

We conducted a within-subject study to evaluate the effectiveness of Tutorly's active guidance and feedback. Concretely, we gathered user feedback on Tutorly's support for (1) knowledge improvement in several learning goals pertaining to EDA and (2) the overall usability and cognitive load while using Tutorly.

## 6.1 Method

We recruited 16 participants who were proficient in the R programming language but novices in EDA through our personal networks and social media messaging. The individual sessions were conducted via Zoom and lasted approximately 90 minutes. We compensated participants with $50 for the time. At the start of the session, participants were given instructions to install the extension on their own system or remotely control Tutorly deployed on the lead author's laptop. Most participants opted for the latter as they also had to install JuypterLab, pip, and R on their machines, which can be time-consuming.

In each session, participants were asked to first fill in a pre-test questionnaire to assess their prior knowledge of EDA topics. The questionnaire had five questions, including proposing a hypothesis, visualizing the data, and interpreting the chart. For example, a question for visualizing the data is "For the problem 'What are the countries with the highest and lowest student/teacher ratios,' which visualization technique is most appropriate?". The questions were made by the first author based on a YouTube video "Tidy Tuesday Screencast: analyzing student/teacher ratios and other country statistics [3]". Next, participants were randomly assigned to one of two conditions: (1) watch the video directly on YouTube and use an IDE of their choice and ChatGPT to support a learning-by-doing approach, or (2) use Tutorly environment. Each condition was assigned a different video from Dave Robinson's YouTube channel [4] [5]. For the Tutorly condition, participants also received a 15-minute tutorial on the main features and could try out the system on their own before proceeding to the test tasks. In the second half of the study, participants worked on the other condition, allowing for comparative feedback between watching videos with and without Tutorly.

Finally, participants were given a post-test quiz to fill in, in which there were five questions to test their knowledge of the same learning goals of EDA after learning. The questions were based on a new YouTube video "Tidy Tuesday Screencast: Analyzing

Horror Movies in R [6]". Each participant then completed a usability questionnaire [48], a cognitive load questionnaire [47], and left optional open-ended comments on their learning and potential improvements. The study materials are provided in the supplement.

## 6.2 Results

*6.2.1 Practice with Tutorly results in improvements in EDA tasks.* Because the questions of each learning goal are independent (i.e., no overlap), we conducted a paired *t*-test (Table 4) for the questions in each learning goal based on the results from the pre-test quiz and post-test quiz. Using Tutorly has certain improvements in EDA tasks, but there are varying degrees of improvements for different learning goals. In the case of "Interpret the chart," the mean score increased from 20.69 to 28.44, which was statistically significant, indicating that Tutorly was particularly effective for this learning goal ($p < 0.05$). This may mean the teaching methods were well suited to enhance participants' competency in interpreting graphs. For "Visualize the data," the mean score increased from pre- to post-test (10 to 14.38), although this improvement did not reach statistical significance ($p = 0.150$). "Propose a hypothesis" showed a slight increase in the mean score from 31.25 to 33.75, but it did not reach statistical significance ($p = 0.544$). Overall results indicate a positive trend in these two learning goals that may become significant with an increased quiz size or longer learning sessions.

**Table 4: Paired *t*-test results across the learning goals**

| Learning Goal | Mean Diff | t-statistic | p-value |
|---|---|---|---|
| Visualize the Data | 4.38 | −1.52 | 0.150 |
| Interpret the Chart | 7.75 | −2.26 | 0.039 |
| Propose a Hypothesis | 2.50 | −0.62 | 0.544 |

*6.2.2 Practice with Tutorly shows higher engagement in watching videos.* Across all sessions, participants responded positively to using Tutorly to learn programming through watching videos. Although 75% of the participants expressed that they typically practice programming while watching videos (i.e., learning by doing), more than half felt there was not enough explanation of the code or how to debug it (66.7%), and there were problems applying the concepts in practice due to a lack of detailed examples (58.3%). By contrast, 87.5% agreed or strongly agreed that they like using Tutorly's interface, and 81.25% agreed or strongly agreed that it was easy to learn to use the system (Figure 5). According to P2, when watching a long instructional video, the hardest thing is to stay focused. Instead, using Tutorly encourages them to become more immersed in watching videos and practicing: *"if I want to answer the questions or fill in the code blanks correctly, I should watch the video clip carefully to avoid multiple viewings."*(P2) Similarly, P16 commented: *"I feel like I can actually practice instead of just watching videos numbly."*

*6.2.3 Tutorly demonstrates a better alternative to ChatGPT..* Based on feedback from using both Tutorly and their typical IDE alongside ChatGPT, although using ChatGPT to practice is quite common among participants (93.8%), over half (68.8%) point out it is hard

---

[3]https://youtu.be/NoUHdrailxA?si=8fziI_chqWRWFwnv
[4]https://youtu.be/nx5yhXAQLxw?si=g6YmEUDeL3-vZ3Cw
[5]https://www.youtube.com/live/Kd9BNI6QMmQ?si=JqXX5AbZ3r2fnLiG

[6]https://www.youtube.com/live/Eucfn-KY-t0?si=sEHUyGtZ9E2Vr0eJ

to express their needs clearly, which led to frustration. In contrast, 93.75% of participants believed that the information provided by Tutorly effectively helped them complete tasks and scenarios, and 68.75% of participants thought that the information provided by Tutorly was clear. For example, P10 commented that Tutorly does not generate long and tedious messages, which means it is easier to understand: *"[Tutorly ] doesn't give super long messages like ChatGPT, but it also has enough stuff so I won't get lost in it."* Besides, P3 found it very convenient not to describe the video content to Tutorly: *"I always need to describe a specific situation before asking how to do it [with ChatGPT]."* Given the inconvenience of practicing directly with ChatGPT, it is necessary to develop a tool that could provide video-specific guidance like Tutorly.

*6.2.4 Tutorly reveals the potential for personalized learning.* Among participants who do not practice while watching, a common reason is "Uncertain about how to start or what to practice". With Tutorly, learners could learn the video content guided by learning goals. According to P8: *"it's better if I can choose learning goals I wanna learn."* Another participant (P14) pointed out, *"I dislike reading texts about code".* P9 said, *"I prefer writing code than answering questions."* Using our action configurator, we can flexibly support these learner preferences. Participants also proposed additional features that could be added to make the code cells more integrated with Tutorly, for example, P16: *"Add code function explanation next to each line of the code instead of in the chat interface".*

*6.2.5 Cognitive load of Tutorly's interface.* Tutorly presented a moderate level of intrinsic load, minimized extraneous cognitive demands, and maximized germane cognitive engagement for learning. In the cognitive load survey, we use a 10-point scale where 0 represents "not at all the case" and 10 represents "completely the case". The Intrinsic Load (IL), which represents task complexity and the learner's prior knowledge, averaged at 5.08 ($SD$ = 0.42). This score was in the middle of our scale, indicating that participants found the task complexity to be moderately challenging. Extraneous Load (EL) reflects unnecessary cognitive load imposed by instructional features, with a mean score of 2.75 ($SD$ = 0.66). This significantly lower score indicates that instructional elements are largely perceived as not unduly adding to cognitive load, suggesting that they are carefully designed to avoid unnecessary complexity. Germane Load (GL), associated with cognitive processes beneficial for learning, stood out with a mean of 8.05 ($SD$ = 0.37). This indicates that the vast majority of participants found the instructional features to significantly enhance their learning and understanding, suggesting an effective instructional design aimed at promoting meaningful cognitive processing. The relatively low standard deviations indicate a general consensus among participants regarding this assessment.

*6.2.6 Overall usability of Tutorly's interface.* In general, participants responded positively to the learning process with Tutorly. In the PSSUQ survey, we anchored at the endpoints with the terms "Strongly agree" for 7 and "Strongly disagree" for 1. We found the overall usability score for the system was 6.16 ($SD$ = 1.27). In terms of sub-scales, the System Usefulness (SYSUSE) score was 6.25 ($SD$ = 1.23), indicating that the system was perceived as relatively easy to use and learn by the respondents. The Information Quality

(INFOQUAL) score is 5.95 ($SD$ = 1.43), suggesting that the information provided by the system was generally clear and helpful, but that there may be room for improvement in specific areas such as line-by-line code explanation and debugging assistance. The Interface Quality (INTERQUAL) score is 6.28 ($SD$ = 1.06), which means that participants found the system's interface to be reasonably pleasant and visually appealing. These results suggest a good level of Tutorly's usability overall but highlight potential areas for enhancement, especially in making information more accessible and understandable to further improve the user experience.

## 6.3 Tutorly's Failure Modes

Our user study also highlighted a range of failure patterns that the videos, the LLMs, or the system itself could cause. For instance, at times, the video clips corresponding to the learning goals are too short or too long, which affects the learning experience. We noticed that because video creators solved problems of varying difficulty at different speeds, a video clip can be as short as 12 seconds or as long as 348 seconds. One participant (P15) noted that it would be better if she could watch the video more coherently: *"Just from my point of view, the Tutorials sometimes were cut a little suddenly. So, for people who don't know data analysis so much, it's a challenge for them to give a reaction immediately."* Similarly, one participant (P4) expressed that it is a waste of time to watch the video author correct his own errors (common in unedited screencasts). Merging adjacent short video clips or removing unnecessary parts from long video clips is a venue for future work. Additionally, waiting for LLMs to generate a message can be frustrating. Several participants observed that waiting for an incoming message from Tutorly takes too long. According to our observations, the waiting time is generally between ten seconds to fifteen seconds for one message, which considerably slows down learning efficiency. When we use a self-built chatbot framework instead of the pre-built implementation by LangChain, the waiting time is reduced to about five seconds, which appears acceptable. However, in order to achieve faster response speed, LLMs should be iterated out for faster models with small to moderate context window to cope with chat mode.

## 7 DISCUSSION

In Tutorly, we employ the *CogApp* framework to realize a conversational LLM tutor system to assist in learning from programming videos. We implement a novel pipeline and DSL to generate high-quality conversations while allowing domain and content experts to specific key pedagogical strategies. Our user study findings highlight the benefits, limitations, and opportunities of using Tutorly. Here, we elaborate on these limitations and discuss future opportunities for improvements.

## 7.1 Broader Utility

Beyond a mentorship learning context, a compelling application of Tutorly is in pair programming. Pair programming is a technique in which two individuals share a single computer as they work together to develop software. In traditional pair programming, one person actively writes code (the "driver"), and the other person provides guidance, feedback, and suggestions (the "navigator"). We can imagine our Tutorly's prompting techniques can be extended to
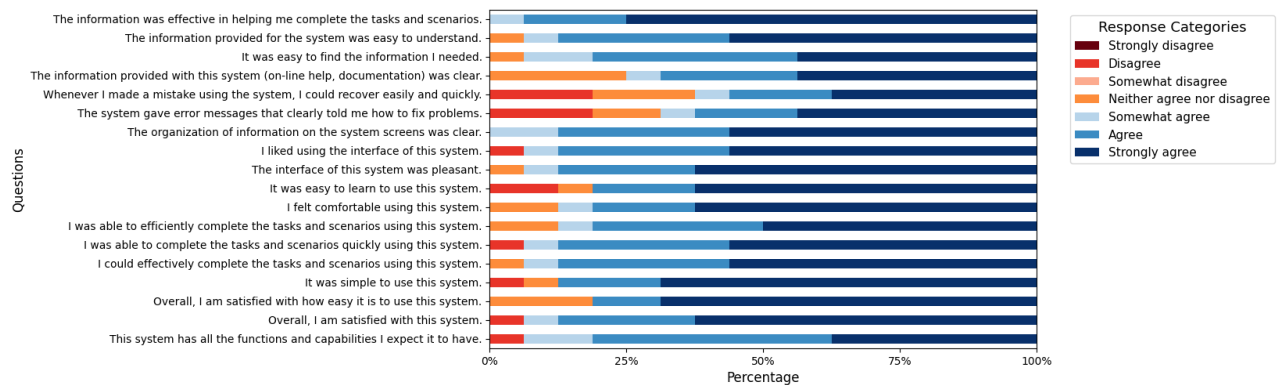
**Figure 5: Participant Responses to Post-Study System Usability Questionnaire.**

support collaborative brainstorming and problem-solving in open-ended domains. Rather than relying on video and transcript, the system might build on code examples and datasets to generate productive interactions. Such approaches are shown to provide a variety of benefits, including increased code quality, productivity, creativity, knowledge management, and self-efficacy [18, 79].

Second, we imagine content creators can leverage Tutorly to provide more engaging viewing and learning experiences within their own platforms to scale the utility of their content. As domain and topic experts, they can take a human-in-the-loop approach to configure the prompt behavior and contexts. Unlike straightforward factual information, their expertise involves deep understanding, the ability to navigate nuances based on experiential and tacit knowledge, and the application of knowledge in varied, often unpredictable, real-world situations. Tutorly can be used to formulate novel learning goals and teaching moves to help learners acquire these skills. Third, we imagine Tutorly being useful in formal classrooms. Since the pandemic, many courses have been offered in a hybrid format with recorded lectures for students to watch asynchronously. We imagine Tutorly can be configured by instructors to integrate and interleave watching and completing specific assignments set by the instructors. Further, teaching assistants might be able to personalize the prompt parameters in Tutorly to provide better guidance and practice on hard-to-grasp concepts and for low-achieving students. Rather than accompanying students to practice their ability to write code, Tutorly can also offset the burden on teaching assistants by helping students to find problems, think about problems, and solve problems while watching videos.

## 7.2 Limitations and Future Work

*7.2.1 Adding Novel Interventions.* In Tutorly's conversational interactions, we allow experts to specify well-known intervention strategies such as multiple-choice-questions (MCQs) or fill in the blank. There is an opportunity to support more expressive interactions, such as allowing students to directly annotate over visualizations to engage in diagram construction, etc. However, our current implementation favors generalizability over integrating topic-specific techniques. Future research can investigate ways for LLMs to transform users' natural language description into a domain-specific language for domain-specific interaction formats.

Alternately, the DSL can be extended to include techniques for generating interactive widgets that are interpretable for the student model.

*7.2.2 Cross-Platforms and Cross-Language Learning.* We deploy Tutorly as an extension in JupyterLab that provides a convenient platform for interactive data science and scientific computing. However, users may have preferences for an integrated development environment (IDE), such as using R Studio for R, PyCharm for Python, XCode for C++, or Visual Studio Code. Therefore, future work can extend Tutorly to multiple IDEs as plug-ins. A different scenario is when students are learning game development; they not only need IDE to write C# code but also need Unity to watch and manipulate. This will require cross-platform integration. More importantly, we imagine future iterations of Tutorly can take an input video in one language, such as R, and generate code and instructions in a different language, such as Python while keeping the underlying learning content the same.

*7.2.3 Supporting Applications in Diverse Domains.* Although our evaluation study provides valuable insights into Tutorly's usage in exploratory data analysis videos, this topic does not encompass all the programming videos. Based on user feedback, many people also want to learn topics such as machine learning, game development, etc. Our technical evaluation has proven that the current DSL structure can support conversation for other topics. More broadly, how to support not only programming learning but also more everyday life teaching videos, such as cooking, makeup, and fitness, opens up problems for inquiry and solutions in the future. There are already some works that support the learning processes of life-teaching videos [87]. However, students may need to be provided with feedback rather than just learning. Multi-modal LLMs may be a potential solution to monitor and analyze student performance in the future.

## 8 CONCLUSION

We developed Tutorly, a system for assisting students learn programming while watching online tutorial videos. While using Tutorly, learners interact with an LLM-powered tutor through rich conversational interactions and follow the guidance and feedback to meet their learning goals. Tutorly design is informed by the *CogApp*

framework in learning sciences, and backed by a student model to generate targeted pedagogical strategies to support effective learning. In other words, Tutorly dynamically adapts its teaching strategies based on the observed learning progress of students. Our prompting pipeline is generalizable across different topics, actions, and learning goals. Our user evaluation reveals that Tutorly augments video-based learning experiences, is aligned with learning goals, and does not increase the learner's cognitive workload. Using Tutorly, learners can effectively engage in self-directed learning to acquire and apply the concepts and techniques contained in programming videos.

## REFERENCES

[1] Shamim Akhter, Muhammad Kashan Javed, Syed Qasim Shah, and A Javaid. 2021. Highlighting the advantages and disadvantages of E-learning. *Psychol. Educ* 58, 5 (2021), 1607–14.

[2] Yasar Akyuz. 2020. Effects of Intelligent Tutoring Systems (ITS) on Personalized Learning (PL). *Creative Education* 11 (2020), 953–978. https://api.semanticscholar.org/CorpusID:225735520

[3] Carlos Alario-Hoyos, Carlos Delgado-Kloos, Iria Estévez-Ayres, Carmen Fernández Panadero, Jorge Blasco, Sergio Pastrana, Guillermo Suarez-Tangil, and Julio Villena. 2016. Interactive activities: the key to learning programming with MOOCs.

[4] Carlos ALARIO-HOYOS, Carlos DELGADO KLOOS, Iria ESTÉVEZ-AYRES, Carmen FERNÁNDEZ-PANADERO, Jorge BLASCO, Sergio PASTRANA, Guillero SUÁREZ-TANGIL, and Julio VILLENA-ROMÁN. 2016. Interactive activities: the key to learning programming with MOOCs. *Proceedings of the European Stakeholder Summit on experiences and best practices in and around MOOCs (EMOOCS 2016)* (2016), 319.

[5] Sumayh S. Aljameel, James D. O'Shea, Keeley A. Crockett, and Annabel Latham. 2018. Can LANA CITS Support Learning in Autistic Children? A Case Study Evaluation. In *Intelligent Systems with Applications*. https://api.semanticscholar.org/CorpusID:62929662

[6] JosÉ Luis Ambite, Lily Fierro, Jonathan Gordon, Gully A. P. C. Burns, Florian Geigl, Kristina Lerman, and John D. Van Horn. 2021. BD2K Training Coordinating Center's ERuDIte: The Educational Resource Discovery Index for Data Science. *IEEE Transactions on Emerging Topics in Computing* 9, 1 (2021), 316–328. https://doi.org/10.1109/TETC.2019.2903466

[7] Lorin W Anderson and David R Krathwohl. 2001. *A taxonomy for learning, teaching, and assessing: A revision of Bloom's taxonomy of educational objectives: complete edition*. Addison Wesley Longman, Inc.

[8] Ryan SJd Baker. 2007. Modeling and understanding students' off-task behavior in intelligent tutoring systems. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. 1059–1068.

[9] Lingfeng Bao, Shengyi Pan, Zhenchang Xing, Xin Xia, D. Lo, and Xiaohu Yang. 2020. Enhancing developer interactions with programming screencasts through accurate code extraction. *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering* (2020). https://api.semanticscholar.org/CorpusID:226274152

[10] Lingfeng Bao, Zhenchang Xing, Xin Xia, and D. Lo. 2019. VT-Revolution: Interactive Programming Video Tutorial Authoring and Watching System. *IEEE Transactions on Software Engineering* 45 (2019), 823–838. https://api.semanticscholar.org/CorpusID:49581547

[11] Brigid JS Barron, Daniel L Schwartz, Nancy J Vye, Allison Moore, Anthony Petrosino, Linda Zech, and John D Bransford. 2014. Doing with understanding: Lessons from research on problem-and project-based learning. In *Learning through problem solving*. Psychology Press, 271–311.

[12] John D Bransford and Daniel L Schwartz. 1999. Chapter 3: Rethinking transfer: A simple proposal with multiple implications. *Review of research in education* 24, 1 (1999), 61–100.

[13] John Seely Brown, Allan M. Collins, and Paul Duguid. 1989. Situated Cognition and the Culture of Learning. *Educational Researcher* 18 (1989), 32 – 42. https://api.semanticscholar.org/CorpusID:9824073

[14] Yining Cao, Hariharan Subramonyam, and Eytan Adar. 2022. VideoSticker: A Tool for Active Viewing and Visual Note-taking from Videos. In *27th International Conference on Intelligent User Interfaces*. 672–690.

[15] Carl Casey. 1996. Incorporating cognitive apprenticeship in multi-media. *Educational Technology Research and Development* 44 (1996), 71–84. https://api.semanticscholar.org/CorpusID:62657577

[16] Victor Cassone. 2022. How to Use Deliberate Practice to Learn Programming More Efficiently. https://www.freecodecamp.org/news/how-to-use-deliberate-practice-to-learn-programming-fast/

[17] Mehmet Celepkolu and Kristy Elizabeth Boyer. 2018. The Importance of Producing Shared Code Through Pair Programming. In *Proceedings of the 49th ACM Technical Symposium on Computer Science Education* (Baltimore, Maryland, USA) *(SIGCSE '18)*. Association for Computing Machinery, New York, NY, USA, 765–770. https://doi.org/10.1145/3159450.3159506

[18] Mehmet Celepkolu and Kristy Elizabeth Boyer. 2018. Thematic Analysis of Students' Reflections on Pair Programming in CS1. In *Proceedings of the 49th ACM Technical Symposium on Computer Science Education* (Baltimore, Maryland, USA) *(SIGCSE '18)*. Association for Computing Machinery, New York, NY, USA, 771–776. https://doi.org/10.1145/3159450.3159516

[19] ChatScope Development Team. 2024. ChatScope UI Kit. https://chatscope.io Accessed: 2024-04-03.

[20] Yuanchun Chen, Walter S. Lasecki, and Tao Dong. 2020. Towards Supporting Programming Education at Scale via Live Streaming. *Proceedings of the ACM on Human-Computer Interaction* 4 (2020), 1 – 19. https://api.semanticscholar.org/CorpusID:225094301

[21] Allan Collins, John Seely Brown, Ann Holum, et al. 1991. Cognitive apprenticeship: Making thinking visible. *American educator* 15, 3 (1991), 6–11.

[22] R4DS Online Learning Community. 2023. Tidy Tuesday: A weekly social data project. https://github.com/rfordatascience/tidytuesday

[23] Yixiang Dai and Zihao Chen. 2019. Research on MOOC's Ways of Teaching Students According to Their Aptitude. https://api.semanticscholar.org/CorpusID:239321757

[24] Bidyut Das, Mukta Majumder, Santanu Phadikar, and Arif Ahmed Sekh. 2021. Automatic question generation and answer assessment: a survey. *Research and Practice in Technology Enhanced Learning* 16 (2021). https://api.semanticscholar.org/CorpusID:232301997

[25] Cassius D'Helon, Jae wook Ahn, Vinay Kumar Reddy Kasireddy, and Nirmal K. Mukhi. 2019. INTERACTIVE LEARNING IN A CONVERSATIONAL INTELLIGENT TUTORING SYSTEM USING STUDENT FEEDBACK, CONCEPT GROUPING AND TEXT LINKING. *INTED2019 Proceedings* (2019). https://api.semanticscholar.org/CorpusID:145967789

[26] Travis Faas, Lynn Dombrowski, Alyson Young, and Andrew D. Miller. 2018. Watch Me Code: Programming Mentorship Communities on Twitch.Tv. *Proc. ACM Hum.-Comput. Interact.* 2, CSCW, Article 50 (nov 2018), 18 pages. https://doi.org/10.1145/3274319

[27] Facebook Inc. 2024. React: A JavaScript library for building user interfaces. React Official Website. https://reactjs.org/ Accessed: 2024-04-03.

[28] Shi Feng, Alejandra J. Magana, and Dominic Kao. 2021. A Systematic Review of Literature on the Effectiveness of Intelligent Tutoring Systems in STEM. *2021 IEEE Frontiers in Education Conference (FIE)* (2021), 1–9. https://api.semanticscholar.org/CorpusID:245388778

[29] Christine Geith and Karen Vignare. 2008. ACCESS TO EDUCATION WITH ONLINE LEARNING AND OPEN EDUCATIONAL RESOURCES: CAN THEY CLOSE THE GAP? *Online Learning* (2008). https://api.semanticscholar.org/CorpusID:166615002

[30] Michail N. Giannakos. 2013. Exploring the video-based learning research: A review of the literature. *Br. J. Educ. Technol.* 44 (2013), 191–. https://api.semanticscholar.org/CorpusID:27007535

[31] Arthur C. Graesser, Shulan Lu, G. Tanner Jackson, Heather H. Mitchell, Mathew Ventura, Andrew M. Olney, and Max M. Louwerse. 2004. AutoTutor: A tutor with dialogue in natural language. *Behavior Research Methods, Instruments, & Computers* 36 (2004), 180–192. https://api.semanticscholar.org/CorpusID:62771305

[32] Philip J. Guo. 2015. Codeopticon: Real-Time, One-To-Many Human Tutoring for Computer Programming. *Proceedings of the 28th Annual ACM Symposium on User Interface Software & Technology* (2015). https://api.semanticscholar.org/CorpusID:11709892

[33] Philip J. Guo, Juho Kim, and Rob Rubin. 2014. How video production affects student engagement: an empirical study of MOOC videos. *Proceedings of the first ACM conference on Learning @ scale conference* (2014). https://api.semanticscholar.org/CorpusID:14910737

[34] Jana M. Hackathorn, Erin D. Solomon, Kate L. Blankmeyer, Rachel E. Tennial, and Amy M. Garczynski. 2011. Learning by Doing: An Empirical Study of Active Teaching Techniques. *The Journal of Effective Teaching* 11 (2011), 40–54. https://api.semanticscholar.org/CorpusID:15580253

[35] Moula Husain and S. M. Meena. 2019. Multimodal Fusion of Speech and Text using Semi-supervised LDA for Indexing Lecture Videos. *2019 National Conference on Communications (NCC)* (2019), 1–6. https://api.semanticscholar.org/CorpusID:174819414

[36] Ross Ihaka and Robert Gentleman. 2024. R: A Language and Environment for Statistical Computing. https://www.r-project.org/. Accessed: 2024-04-02.

[37] Siyuan Ji and Tangming Yuan. 2022. Conversational Intelligent Tutoring Systems for Online Learning: What do Students and Tutors Say? *2022 IEEE Global Engineering Education Conference (EDUCON)* (2022), 292–298. https://api.semanticscholar.org/CorpusID:248698848

[38] Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, George Louis Groh, Stephan Günnemann, Eyke Hüllermeier, Stephan Krusche, Gitta Kutyniok, Tilman Michaeli, Claudia Nerdel, Jürgen Pfeffer, Oleksandra Poquet, Michael Sailer, Albrecht Schmidt, Tina Seidel, Matthias Stadler, Jochen Weller, Jochen Kuhn, and Gjergji Kasneci. 2023. ChatGPT for good? On opportunities and challenges of large language models for education. *Learning and Individual Differences* (2023). https://api.semanticscholar.org/CorpusID:257445349

[39] Kandarp Khandwala and Philip J. Guo. 2018. Codemotion: expanding the design space of learner interactions with computer programming tutorial videos. *Proceedings of the Fifth Annual ACM Conference on Learning at Scale* (2018). https://api.semanticscholar.org/CorpusID:49304895

[40] Ada S. Kim and Amy J. Ko. 2017. A Pedagogical Analysis of Online Coding Tutorials. *Proceedings of the 2017 ACM SIGCSE Technical Symposium on Computer Science Education* (2017). https://api.semanticscholar.org/CorpusID:6717766

[41] Juho Kim, Philip J. Guo, Daniel T. Seaton, Piotr Mitros, Krzysztof Z Gajos, and Rob Miller. 2014. Understanding in-video dropouts and interaction peaks inonline lecture videos. *Proceedings of the first ACM conference on Learning @ scale conference* (2014). https://api.semanticscholar.org/CorpusID:220600751

[42] Juho Kim, Phu Tran Nguyen, Sarah A. Weir, Philip J. Guo, Rob Miller, and Krzysztof Z Gajos. 2014. Crowdsourcing step-by-step information extraction to enhance existing how-to videos. *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (2014). https://api.semanticscholar.org/CorpusID:1470846

[43] Amy J. Ko, Thomas D. Latoza, Stephen Hull, Ellen A. Ko, William Kwok, Jane Quichocho, Harshitha Akkaraju, and Rishin Pandit. 2019. Teaching Explicit Programming Strategies to Adolescents. *Proceedings of the 50th ACM Technical Symposium on Computer Science Education* (2019). https://api.semanticscholar.org/CorpusID:67866527

[44] Monlin Ko and Christina Krist. 2019. Opening up curricula to redistribute epistemic agency: A framework for supporting science teaching. *Science Education* (2019). https://api.semanticscholar.org/CorpusID:155416039

[45] Dimitrios Kravvaris and Katia Lida Kermanidis. 2018. Automatic point of interest detection for open online educational video lectures. *Multimedia Tools and Applications* 78 (2018), 2465–2479. https://api.semanticscholar.org/CorpusID:49671017

[46] Annabel Latham, Keeley A. Crockett, David Mclean, and Bruce Edmonds. 2011. Oscar: An Intelligent Adaptive Conversational Agent Tutoring System. In *Agent and Multi-Agent Systems: Technologies and Applications*. https://api.semanticscholar.org/CorpusID:32673366

[47] Jimmie Leppink, Fred Paas, Cees Van der Vleuten, Tamara Gog, and Jeroen J. G. Van Merrienboer. 2013. Development of an instrument for measuring different types of cognitive load. *Behavior research methods* 45 (04 2013). https://doi.org/10.3758/s13428-013-0334-1

[48] James R. Lewis. 1992. Psychometric Evaluation of the Post-Study System Usability Questionnaire: The PSSUQ. *Proceedings of the Human Factors Society Annual Meeting* 36, 16 (1992), 1259–1260. https://doi.org/10.1177/154193129203601617 arXiv:https://doi.org/10.1177/154193129203601617

[49] Toby Jia-Jun Li, Jingya Chen, Haijun Xia, Tom Michael Mitchell, and Brad A. Myers. 2020. Multi-Modal Repairs of Conversational Breakdowns in Task-Oriented Dialogs. *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology* (2020). https://api.semanticscholar.org/CorpusID:221733725

[50] Toby Jia-Jun Li, Tom Michael Mitchell, and Brad A. Myers. 2020. Interactive Task Learning from GUI-Grounded Natural Language Instructions and Demonstrations. In *Annual Meeting of the Association for Computational Linguistics*. https://api.semanticscholar.org/CorpusID:220057195

[51] Toby Jia-Jun Li, Marissa Radensky, Justin Jia, Kirielle Singarajah, Tom Michael Mitchell, and Brad A. Myers. 2019. PUMICE: A Multi-Modal Agent that Learns Concepts and Conditionals from Natural Language and Demonstrations. *Proceedings of the 32nd Annual ACM Symposium on User Interface Software and Technology* (2019). https://api.semanticscholar.org/CorpusID:202120459

[52] Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let's Verify Step by Step. arXiv:2305.20050 [cs.LG]

[53] Yu-Tzu Lin, Martin K.-C. Yeh, and Shenglong Tan. 2022. Teaching Programming by Revealing Thinking Process: Watching Experts' Live Coding Videos With Reflection Annotations. *IEEE Transactions on Education* 65 (2022), 617–627. https://api.semanticscholar.org/CorpusID:247469461

[54] Marcia Linn. 1995. Designing computer learning environments for engineering and computer science: The scaffolded knowledge integration framework. *Journal of Science Education and Technology* 4 (06 1995), 103–126. https://doi.org/10.1007/BF02214052

[55] MarciaC. Linn, Philip Bell, and ElizabethA. Davis. 2013. *3 Specific Design Principles: Elaborating the Scaffolded Knowledge Integration Framework*. 343–368. https://doi.org/10.4324/9781410610393-23

[56] Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying Zhang, Ruocheng Guo, Hao Cheng, Yegor Klochkov, Muhammad Faaiz Taufiq, and Hang Li. 2024. Trustworthy LLMs: a Survey and Guideline for Evaluating Large Language Models' Alignment. arXiv:2308.05374 [cs.AI]

[57] Dastyni Loksa and Amy J. Ko. 2016. The Role of Self-Regulation in Programming Problem Solving Process and Success. *Proceedings of the 2016 ACM Conference on International Computing Education Research* (2016). https://api.semanticscholar.org/CorpusID:15818535

[58] Romina Soledad Albornoz-De Luise, Miguel Arevalillo-Herráez, and David Arnau. 2023. On Using Conversational Frameworks to Support Natural Language Interaction in Intelligent Tutoring Systems. *IEEE Transactions on Learning Technologies* 16 (2023), 722–735. https://api.semanticscholar.org/CorpusID:257167610

[59] Alejandra J. Magana, Michael L. Falk, and Michael J. Reese. 2013. Introducing Discipline-Based Computing in Undergraduate Engineering Education. *ACM Trans. Comput. Educ.* 13, 4, Article 16 (nov 2013), 22 pages. https://doi.org/10.1145/2534971

[60] Alejandra J. Magana, Michael L. Falk, Camilo Vieira, and Michael J. Reese. 2016. A Case Study of Undergraduate Engineering Students' Computational Literacy and Self-Beliefs about Computing in the Context of Authentic Practices. *Comput. Hum. Behav.* 61, C (aug 2016), 427–442. https://doi.org/10.1016/j.chb.2016.03.025

[61] Debabrata Mahapatra, Ragunathan Mariappan, Vaibhav Rajan, Kuldeep Yadav, A. Seby, and Sudeshna Roy. 2018. VideoKen: Automatic Video Summarization and Course Curation to Support Learning. *Companion Proceedings of the The Web Conference 2018* (2018). https://api.semanticscholar.org/CorpusID:13796682

[62] Lauren Elizabeth Margulieux, Richard Catrambone, and Mark Guzdial. 2016. Employing subgoals in computer programming education. *Computer Science Education* 26 (2016), 44 – 67. https://api.semanticscholar.org/CorpusID:19782935

[63] Angela McGlynn. 2005. Teaching Millennials, Our Newest Cultural Cohort. *Education Digest: Essential Readings Condensed for Quick Review* 71 (01 2005).

[64] Janet Metcalfe, Judy Xu, Matti Vuorre, Robert Siegler, Dylan Wiliam, and Robert A Bjork. 2024. Learning from errors versus explicit instruction in preparation for a test that counts. *British Journal of Educational Psychology* (2024).

[65] K. Mukhtar, K. Javed, M. Arooj, and A. Sethi. 2020. Advantages, Limitations and Recommendations for online learning during COVID-19 pandemic era. *Pak J Med Sci* 36, COVID19-S4 (May 2020), S27–S31. https://doi.org/10.12669/pjms.36.COVID19-S4.2785

[66] Medhini Narasimhan, Arsha Nagrani, Chen Sun, Michael Rubinstein, Trevor Darrell, Anna Rohrbach, and Cordelia Schmid. 2022. TL;DW? Summarizing Instructional Videos with Task Relevance & Cross-Modal Saliency. In *European Conference on Computer Vision*. https://api.semanticscholar.org/CorpusID:251564419

[67] Inc. NetEase Youdao. 2023. BCEmbedding: Bilingual and Crosslingual Embedding for RAG. https://github.com/netease-youdao/BCEmbedding.

[68] Benjamin D. Nye, Arthur C. Graesser, and Xiangen Hu. 2014. AutoTutor and Family: A Review of 17 Years of Natural Language Tutoring. *International Journal of Artificial Intelligence in Education* 24 (2014), 427–469. https://api.semanticscholar.org/CorpusID:14424486

[69] Sharmaine Gek Teng Ong and Gwendoline Choon Lang Quek. 2023. Enhancing teacher–student interactions and student online engagement in an online learning environment. *Learning Environments Research* (2023), 1 – 27. https://api.semanticscholar.org/CorpusID:255843640

[70] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel

Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. GPT-4 Technical Report. arXiv:2303.08774 [cs.CL]

[71] Eng Lieh Ouh, Benjamin Kok Siew Gan, and D. Lo. 2022. ITSS: Interactive Web-Based Authoring and Playback Integrated Environment for Programming Tutorials. *2022 IEEE/ACM 44th International Conference on Software Engineering: Software Engineering Education and Training (ICSE-SEET)* (2022), 158–164. https://api.semanticscholar.org/CorpusID:248239730

[72] Roy D. Pea. 2004. The Social and Technological Dimensions of Scaffolding and Related Theoretical Concepts for Learning, Education, and Human Activity. *Journal of the Learning Sciences* 13 (2004), 423 – 451. https://api.semanticscholar.org/CorpusID:10481973

[73] Oleksandra Poquet, Lisa-Angelique Lim, Negin Mirriahi, and Shane Dawson. 2018. Video and learning: a systematic review (2007–2017). *Proceedings of the 8th International Conference on Learning Analytics and Knowledge* (2018). https://api.semanticscholar.org/CorpusID:3706765

[74] Tine S Prøitz. 2010. Learning outcomes: What are they? Who defines them? When and where are they defined? *Educational assessment, evaluation and accountability* 22, 2 (2010), 119–137.

[75] Project Jupyter. 2024. JupyterLab. https://jupyter.org/.

[76] Cheng Ren, Zachary Pardos, and Zhi Li. 2024. Human-AI Collaboration Increases Skill Tagging Speed but Degrades Accuracy. *arXiv preprint arXiv:2403.02259* (2024).

[77] Peter Robe, Sandeep K. Kuttal, Jake AuBuchon, and Jacob Hart. 2022. Pair programming conversations with agents vs. developers: challenges and opportunities for SE community. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering* (<conf-loc>, <city>Singapore</city>, <country>Singapore</country>, </conf-loc>) *(ESEC/FSE 2022)*. Association for Computing Machinery, New York, NY, USA, 319–331. https://doi.org/10.1145/3540250.3549127

[78] Dave Robinson. 2024. Educational Content on Data Science: TidyTuesday and YouTube Tutorials. YouTube Channel and TidyTuesday Community Project. https://www.youtube.com/user/DaveRobinson Accessed: 2024-04-02.

[79] Fernando J. Rodríguez, Kimberly Michelle Price, and Kristy Elizabeth Boyer. 2017. Exploring the Pair Programming Process: Characteristics of Effective Collaboration. In *Proceedings of the 2017 ACM SIGCSE Technical Symposium on Computer Science Education* (Seattle, Washington, USA) *(SIGCSE '17)*. Association for Computing Machinery, New York, NY, USA, 507–512. https://doi.org/10.1145/3017680.3017748

[80] Gilbert Ryle. 1945. Knowing how and knowing that: The presidential address. In *Proceedings of the Aristotelian society*, Vol. 46. JSTOR, 1–16.

[81] Marija Sablić, Ana Mirosavljević, and Alma Škugor. 2020. Video-Based Learning (VBL)—Past, Present and Future: an Overview of the Research Published from 2008 to 2019. *Technology, Knowledge and Learning* 26 (2020), 1061 – 1077. https://api.semanticscholar.org/CorpusID:220511137

[82] Marlene Scardamalia. 2002. Collective cognitive responsibility for the advancement of knowledge. https://api.semanticscholar.org/CorpusID:142643453

[83] Daniel L Schwartz and Kevin Hartman. 2014. It's not television anymore: Designing digital video for learning and assessment. In *Video research in the learning sciences*. Routledge, 335–348.

[84] Ulker Shafiyeva. 2021. Assessing Students' Minds: Developing Critical Thinking or Fitting into Procrustean Bed. *European Journal of Education* 4 (2021), 78 – 91. https://api.semanticscholar.org/CorpusID:243842555

[85] David Stroupe. 2014. Examining Classroom Science Practice Communities: How Teachers and Students Negotiate Epistemic Agency and Learn Science-as-Practice. *Science Education* 98 (2014), 487–516. https://api.semanticscholar.org/CorpusID:145116373

[86] Matilde Sánchez-Peña, Camilo Vieira, and Alejandra J. Magana. 2023. Data science knowledge integration: Affordances of a computational cognitive apprenticeship on student conceptual understanding. *Computer Applications in Engineering Education* 31, 2 (2023), 239–259. https://doi.org/10.1002/cae.22580 arXiv:https://onlinelibrary.wiley.com/doi/pdf/10.1002/cae.22580

[87] Anh Tuan Truong, Peggy Chi, D. Salesin, Irfan Essa, and Maneesh Agrawala. 2021. Automatic Generation of Two-Level Hierarchical Tutorials from Instructional Makeup Videos. *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (2021). https://api.semanticscholar.org/CorpusID:232045282

[88] John W Tukey et al. 1977. *Exploratory data analysis*. Vol. 2. Reading, MA.

[89] Kurt VanLehn. 2006. The behavior of tutoring systems. *International journal of artificial intelligence in education* 16, 3 (2006), 227–265.

[90] Kurt VanLehn. 2011. The relative effectiveness of human tutoring, intelligent tutoring systems, and other tutoring systems. *Educational psychologist* 46, 4 (2011), 197–221.

[91] Camilo Vieira, Alejandra Magana, Anindya Roy, and Michael Falk. 2021. Providing students with agency to self- scaffold in a computational science and engineering course. *Journal of Computing in Higher Education* 33 (08 2021). https://doi.org/10.1007/s12528-020-09267-7

[92] Bryan Wang, Mengyu Yang, and Tovi Grossman. 2021. Soloist: Generating Mixed-Initiative Tutorials from Existing Guitar Instructional Videos Through Audio Processing. *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (2021). https://api.semanticscholar.org/CorpusID:231693070

[93] Wayne H. Ward, Ronald A. Cole, Daniel Bolaños, Cindy Buchenroth-Martin, Edward Svirsky, and Timothy B. Weston. 2013. My Science Tutor: A Conversational Multimedia Virtual Tutor. *Journal of Educational Psychology* 105 (2013), 1115–1125. https://api.semanticscholar.org/CorpusID:262737077

[94] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. arXiv:2201.11903 [cs.CL]

[95] Kuldeep Yadav, Ankit Gandhi, Arijit Biswas, Kundan Srivastava, Saurabh Srivastava, and Om Deshmukh. 2016. ViZig: Anchor Points based Non-Linear Navigation and Summarization in Educational Videos. *Proceedings of the 21st International Conference on Intelligent User Interfaces* (2016). https://api.semanticscholar.org/CorpusID:14086148

[96] Lei Zhang and Steve Oney. 2020. FlowMatic: An Immersive Authoring Tool for Creating Interactive Scenes in Virtual Reality. *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology* (2020). https://api.semanticscholar.org/CorpusID:222799772

[97] Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, Kathleen McKeown, and Tatsunori B Hashimoto. 2024. Benchmarking large language models for news summarization. *Transactions of the Association for Computational Linguistics* 12 (2024), 39–57.

[98] Wentian Zhao, Seokhwan Kim, Ning Xu, and Hailin Jin. 2020. Video Question Answering on Screencast Tutorials. *ArXiv* abs/2008.00544 (2020). https://api.semanticscholar.org/CorpusID:220484710

[99] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A Survey of Large Language Models. arXiv:2303.18223 [cs.CL]

[100] Barry J Zimmerman and Dale H Schunk. 2011. Self-regulated learning and performance: An introduction and an overview. *Handbook of self-regulation of learning and performance* (2011), 15–26.

[101] Claus Zinn. 2013. Algorithmic Debugging for Intelligent Tutoring: How to Use Multiple Models and Improve Diagnosis. In *Deutsche Jahrestagung für Künstliche Intelligenz*. https://api.semanticscholar.org/CorpusID:2161547

[102] Serkan Şendağ, İlker Yakin, and Nuray Gedik. 2023. Fostering creative thinking skills through computer programming: Explicit or integrated teaching? *Education and Information Technologies* 28 (02 2023), 1–20. https://doi.org/10.1007/s10639-023-11629-4

# A APPENDIX

## A.1 Cognitive Apprenticeship Methods

Figure 6 shows the definition of the six cognitive apprenticeship methods and their corresponding messages.

## A.2 Domain Specific Language Structure

Below is an example of the DSL structure for the video segment "Visualize the data" when the student has a median skills level.

```
[
    {
        "knowledge": "Declarative knowledge:
            The task is visualizing the
            distribution of median earnings
            across major categories using a
            box plot and enhancing
            readability by reordering
            categories and formatting axis
            labels.",
        "actions": [
            {
                "method": "Scaffolding",
                "action": "Demonstrate the
                    current task and provide
                    explanations of the
                    concepts underlying the
                    current step of the task
                    using plain-text",
                "prompt": "[Use one sentence
                    to explain the
                    Declarative knowledge:
                    The task is visualizing
                    the distribution of
                    median earnings across
                    major categories using a
                    box plot and enhancing
                    readability by reordering
                     categories and
                    formatting axis labels.
                    at this step, such as
                    what effect we want to
                    achieve, why we do it,
                    and what function we use
                    to do it]",
                "interaction": "plain-text",
                "parameters": ["knowledge"],
                "need-response": false
            }
        ]
    },
    {
        "knowledge": "Procedural knowledge:
            To achieve a clear visualization
            of categorical data distributions
            , one must use 'geom_boxplot' on
            'ggplot' in R because it
            effectively displays the spread
            and central tendency of the data.
            ",
        "actions": [
            {
                "method": "Scaffolding",
                "action": "Demonstrate the
                    current task and provide
                    explanations of the
                    concepts underlying the
                    current step of the task
                    using plain-text.",
                "prompt": "[Use one sentence
                    to explain the Procedural
                     knowledge: To achieve a
                    clear visualization of
                    categorical data
                    distributions, one must
                    use 'geom_boxplot' on '
                    ggplot' in R because it
                    effectively displays the
                    spread and central
                    tendency of the data. at
                    this step, such as what
                    effect we want to achieve
                    , why we do it, and what
                    function we use to do it]
                    ",
                "interaction": "plain-text",
                "parameters": ["knowledge"],
                "need-response": false
            },
            {
                "method": "Coaching",
                "action": "Use fill-in-blanks
                     to guide the student
                    through practice
                    exercises, offering
                    targeted hints and
                    feedback.",
                "prompt": "[Use one sentence
                    to prompt the student to
                    fill in the {code-line-
                    with-blanks} below][
                    Provide a brief hint to
                    help them through it]",
                "interaction": "fill-in-
                    blanks",
                "parameters": ["code-line-
                    with-blanks"],
                "need-response": true
            }
        ]
    },
    {
        "knowledge": "Procedural knowledge:
            To achieve an ordered factor
            level based on the 'Median', one
            must use 'fct_reorder' on '
            Major_category', because it
            facilitates easier comparison
            across categories by sorting them
             from lowest to highest median
            earnings.",
        "actions": [
```
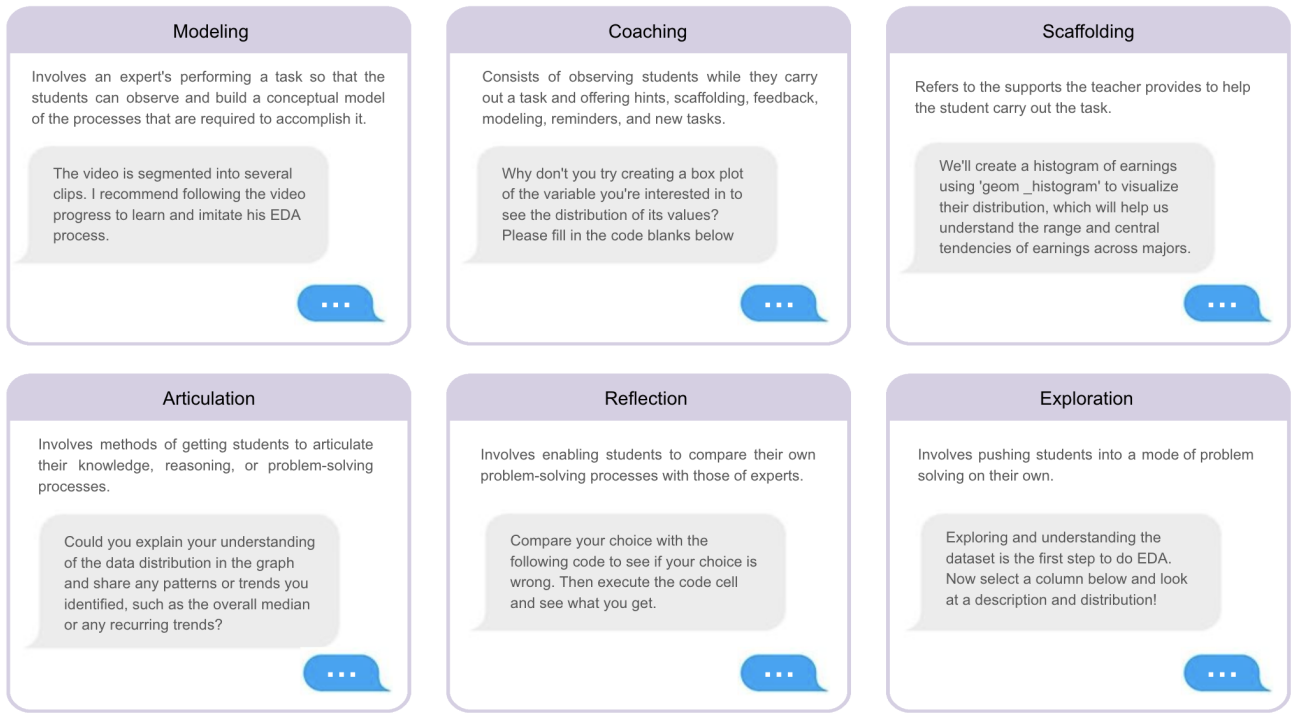
**Figure 6: The six cognitive apprenticeship methods and an example message of each method.**

```
{
    "method": "Scaffolding",
    "action": "Demonstrate the
        current task and provide
        explanations of the
        concepts underlying the
        current step of the task
        using plain-text.",
    "prompt": "[Use one sentence
        to explain the Procedural
         knowledge: To achieve an
         ordered factor level
        based on the 'Median',
        one must use 'fct_reorder
        ' on 'Major_category',
        because it facilitates
        easier comparison across
        categories by sorting
        them from lowest to
        highest median earnings.
        at this step, such as
        what effect we want to
        achieve, why we do it,
        and what function we use
        to do it]",
    "interaction": "plain-text",
    "parameters": ["knowledge"],
    "need-response": false
}
```

```
    ]
},
{
    "knowledge": "Procedural knowledge:
        To achieve improved readability
        of axis labels, one must use '
        coord_flip' and '
        scale_y_continuous' with '
        dollar_format' on the plot
        because flipping the coordinates
        helps in reading long category
        names and dollar formatting makes
         the earnings data more
        interpretable.",
    "actions": [
        {
            "method": "Scaffolding",
            "action": "Demonstrate the
                current task and provide
                explanations of the
                concepts underlying the
                current step of the task
                using plain-text.",
```

```
                "prompt": "[Use one sentence
                    to explain the Procedural
                     knowledge: To achieve
                    improved readability of
                    axis labels, one must use
                     'coord_flip' and '
                    scale_y_continuous' with
                    'dollar_format' on the
                    plot because flipping the
                     coordinates helps in
                    reading long category
                    names and dollar
                    formatting makes the
                    earnings data more
                    interpretable. at this
                    step, such as what effect
                     we want to achieve, why
                    we do it, and what
                    function we use to do it]
                    ",
                "interaction": "plain-text",
                "parameters": ["knowledge"],
                "need-response": false
        },
        {
                "method": "Reflection",
                "action": "Encourage students
                    to review and debug
                    their code using show-
                    code, and to reflect on
                    the learning process by
                    executing the complete
                    code block to verify
                    their understanding.",
                "prompt": "[Use one sentence
                    to let the student
                    compare his answer with
                    the standard {code-block
                    }][Use one sentence to
                    encourage the student to
                    execute the complete code
                     block to verify his
                    understanding]",
                "interaction": "show-code",
                "parameters": ["code-block"],
                "need-response": true
            }
        ]
    }
]
```

## A.3 Action Set

Table 5 shows the action set that generates prompts.

## A.4 Prompts for LLMs

Below are the prompts we fed into the ChatGPT API to perform each of Tutorly's key algorithms.

*A.4.1 Video Segmentation.* Below is the prompt doing **Summarize** in the video segmentation algorithm.

**Summarize:** Here is a video transcript about {video_topic}. Summarize the video content that corresponds to each given learning goal. The transcript is not necessarily arranged in the order in which the learning goals are defined and can contain multiple segments with the same learning goal. The script may contain only some of the learning goals. Please do not include summary of learning goals that do not exist in the transcript. Increase the granularity. For example, if the video author creates two different visualizations, they should be summarized into two points.

Response only in a list in the order of their appearance in the video without any explanations, for example:

```
[
    ("Introduction", summary),
    ("Load data/packages", summary),
    ("Understand the dataset", summary),
    ("Visualize the data", summary),
    ("Interpret the chart", summary),
    ("Visualize the data", summary),
    ("Interpret the chart", summary),
    ("Preprocess the data", summary),
    ...
]
```

*A.4.2 Knowledge Summarize.* Below is the prompt summarizing the declarative and procedural knowledge in a concept-related video segment. The prompts used for the programming-related clips differed only in the definition of the knowledge structure.

The following {video_type} video transcript is about a learning goal: {learning_goal}. Summarize the declarative and procedural knowledge in the video transcript.

The result should be summarized in one sentence of procedural knowledge and no more than {num_1} sentences of declarative knowledge in the order in which it should be learned.

Each piece of knowledge should follow this format:

Procedural knowledge: "To achieve/understand + [specific goal/outcome] + one need to + [general actions/processes] + [additional details] + and consider/use + [relevant factors/tools]." The [general actions/processes] should be quoted in a pair of & sign.

For example, "To understand the distribution of earnings by college major, one need to &examine the histogram and identify overall trend or extreme values&, and consider whether high earnings are due to the field's financial reward or influenced by factors such as low sample size and high variation."

Declarative knowledge: "[Subject] + [verb phrase] + that + [independent clause]".

For example, "The median income by college major shows that majors earn a median income of over $30K right out of college." And sort the output knowledge order according to the correct cognitive order. For example, students need to first learn how to interpret the chart then find out the facts in the chart.

Your response should be in a list format without any explanations:

```
[
    'knowledge_1',
    'knowledge_2',
    ...
]
```

*A.4.3 Teaching Methods Arrangement.* Below is the prompt for arranging the teaching methods according to the cognitive apprenticeship framework.

You are an expert mentor who is good at arranging teaching methods to help students learn from a video about {video_type}. You are teaching students to learn knowledge for {learning_obj} using the

| Move | Action | Prompt |
|------|--------|--------|
| *Programming-Related* | | |
| Scaffolding | Demonstrate the current task and provide explanations of the concepts underlying the current step of the task using **plain-text**. | [Use one sentence to explain the {knowledge} at this step, such as what effect to achieve, why we do it, and what function we use to do it] |
| Coaching | Use **fill-in-blanks** to guide the student through practice exercises, offering targeted hints and feedback. | [Use one sentence to prompt the student to fill in the {code-line-with-blanks} below to practice the {knowledge}][Provide a brief hint to help them through it] |
| Articulation | Use **plain-text** to allow students to articulate their understanding of knowledge. | [Use one sentence to ask the student to explain their understanding and reasoning about {knowledge}, such as articulate why make this kind of visualization rather than others] |
| Reflection | Encourage students to review and debug their code using **show-code** and to reflect on the learning process by executing the complete code block to verify their understanding. | [Use one sentence to let the student compare their answer with the standard {code-block}][Use one sentence to encourage the student to execute the complete code block to verify their understanding] |
| *Concept-Related* | | |
| Scaffolding | Provide structured guidance through **plain-text** as the student works on the task to learn the {knowledge}. | [Use no more than three sentences to guide the student step by step on how to learn and apply the {knowledge}, the student has made the visualization] |
| Coaching | Use **multiple-choice** to observe the student's approach to tasks, offering feedback to guide learning. | Propose a multiple-choice question for the student to understand the {knowledge}, such as what could be the potential reason behind the pattern |
| Articulation | Encourage students to use **interaction** to verbally explain their thought process and reasoning behind their observations and conclusions. | [Use one sentence to ask the student to explain their understanding and reasoning about {knowledge}, such as articulate what patterns they found in the chart] |
| Reflection | Encourage students to use **plain-text** to self-evaluate their performance, identifying strengths and areas for improvement. | [Use one sentence to give feedback on the {student-answer}][Use one sentence to tell the student if any additional steps could confirm their choice][Ask the student to remember the choice and see if it makes sense as they watch the rest of the video] |

**Table 5: Programming and concept-related cognitive apprentice moves in teaching programming and concepts. The bold terms are interactions used in prompts. Parameters are quoted in curly brackets.**

Cognitive Apprenticeship framework. Definition of Cognitive Apprenticeship framework methods:

Coaching: mentor observes mentee's activities along with provision of guidance and feedback

Scaffolding: mentor supports mentee while they work through the task with gradual fading of such supports

Articulation: mentor encourage mentees to verbalize their knowledge and thinking

Reflection: mentor enable mentees to self-assesses own performance

Your task is to choose proper Cognitive Apprenticeship methods to teach the student the given knowledge. The input knowledge list contains the declarative and procedural knowledge in the video. The input student mastery level list has a one-to-one correspondence with the knowledge, representing the student's mastery of each knowledge.

Teaching method arrangement rules:

1. Global before local skills: use Scaffolding as the first move to teach the first knowledge.

2.1 Increasing complexity for concept-related video: Choose one method from Scaffolding, Coaching, or Articulation to teach each knowledge. If the student's mastery level of the corresponding knowledge exceeds 0.5, Scaffolding should fade out. Coaching should be followed by Reflection. Reflection should only be used after Coaching.

2.2 Increasing complexity for programming-related video: For each declarative knowledge, choose between Scaffolding and Articulation. For each procedural knowledge, if the student's mastery level of the corresponding knowledge is lower than 0.3, use Scaffolding only. If the student's mastery level of the corresponding knowledge is between 0.3 and 0.7, use Scaffolding and Coaching. If the student's mastery level of the corresponding knowledge is higher than 0.7, Scaffolding fades out, and Coaching is used only. Reflection should be used once as the last method for the last knowledge.

3. Increasing diversity: diversify the selection of teaching methods based on the first two conditions.

You should use no more than three methods for each knowledge. Please include your choice in a structure in the same format like the following. Response Example:

```
[
    {{
        "knowledge": ...,
        "method": [...]
    }},
    ...
]
```

### A.4.4 Conversation.
The prompt used for generating dialog messages is as follows:

You are an expert in Data Science, specializing in {video_type}. Your task is to use the Cognitive Apprenticeship approach to assist a student in learning {video_type} through David Robinson's Tidy Tuesday tutorial series.

You will be provided with one or more of the following inputs:
- knowledge: the knowledge that will be learned by the student
- pedagogy: the specific cognitive apprenticeship move you must follow to guide students.
- student's code or question or choice: the student's current performance, encompassing either the code in the student's notebook or the student's query sent to you or the student's choice in the multiple-choice question.
- other parameters or requirements: additional information or requirements you must follow to guide the student.

Notes for Response:
- Don't answer or say anything irrelevant to the topic ({video_type}) or the programming language ({kernel_type}).
- Use natural language to communicate in the first person as a teaching assistant.
- You must strictly follow the pedagogy to provide guidance.
- Tailor your advice to the programming language the student uses: {kernel_type}.
- Don't tell the student your response is based on the transcript or code.
- You can find out the full list of conversation history below.