



# Scaffolding Middle-School Mathematics Curricula With Large Language Models

Rizwaan Malik  
Stanford University

Dorna Abdi  
Stanford University

Rose Wang  
Stanford University

Dorottya Demszky  
Stanford University

Despite well-designed curriculum materials, teachers often face challenges in their implementation due to diverse classroom needs. This paper investigates whether Large Language Models (LLMs) can support middle-school math teachers by helping create high-quality curriculum scaffolds, which we define as the adaptations and supplements teachers employ to ensure all students can access and engage with the curriculum. Through Cognitive Task Analysis with expert teachers, we identify a three-stage process for curriculum scaffolding: observation, strategy formulation, and implementation. We incorporate these insights into three LLM approaches to create warmup tasks that activate background knowledge. The best-performing approach, which provides the model with the original curriculum materials and an expert-informed prompt, generates warmups that are rated significantly higher than warmups created by expert teachers in terms of alignment to learning objectives, accessibility to students working below grade level, and teacher preference. This research demonstrates the potential of LLMs to support teachers in creating effective scaffolds and provides a methodology for developing AI-driven educational tools.

VERSION: August 2024

Suggested citation: Malik, Rizwaan, Dorna Abdi, Rose Wang, and Dorottya Demszky. (2024). Scaffolding Middle-School Mathematics Curricula With Large Language Models. (EdWorkingPaper: 24-1028). Retrieved from Annenberg Institute at Brown University: <https://doi.org/10.26300/b47y-mh41>

# Scaffolding Middle-School Mathematics Curricula With Large Language Models

Rizwaan Malik\*

Dorna Abdi†

Rose Wang‡

Dorottya Demszky§

**ABSTRACT:** Despite well-designed curriculum materials, teachers often face challenges in their implementation due to diverse classroom needs. This paper investigates whether Large Language Models (LLMs) can support middle-school math teachers by helping create high-quality curriculum scaffolds, which we define as the adaptations and supplements teachers employ to ensure all students can access and engage with the curriculum. Through Cognitive Task Analysis with expert teachers, we identify a three-stage process for curriculum scaffolding: observation, strategy formulation, and implementation. We incorporate these insights into three LLM approaches to create warmup tasks that activate background knowledge. The best-performing approach, which provides the model with the original curriculum materials and an expert-informed prompt, generates warmups that are rated significantly higher than warmups created by expert teachers in terms of alignment to learning objectives, accessibility to students working below grade level, and teacher preference. This research demonstrates the potential of LLMs to support teachers in creating effective scaffolds and provides a methodology for developing AI-driven educational tools.

**KEYWORDS:** Curriculum Scaffolding; Middle School Mathematics; Human-Computer Interaction; Large Language Models; Cognitive Task Analysis; Large Language Model Evaluation

## Practitioner’s Notes

### What is already known about this topic

---

\*Rizwaan Malik (rizmalik@stanford.edu; corresponding author) is a Graduate Researcher in Educational Data Science at Stanford University, Stanford, CA, United States.

†Dorna Abdi (dabdi@stanford.edu) is a Graduate Researcher in Educational Data Science at Stanford University, Stanford, CA, United States.

‡Rose Wang (rewang@stanford.edu) is a Ph.D. student in Computer Science at Stanford University, Stanford, CA, United States.

§Dorottya (Dora) Demszky (ddemszky@stanford.edu; corresponding author) is an Assistant Professor in Educational Data Science at Stanford University, Stanford, CA, United States.

1. Scaffolding is essential for enabling students to access and engage with curriculum materials.
2. Large Language Models (LLMs) have shown promise in generating educational content and supporting teachers.
3. Teachers frequently need to adapt and supplement standardized curricula to meet the diverse needs of their students.

### **What this paper adds**

1. Identifies a three-stage curriculum scaffolding process (observation, strategy formulation, implementation) used by expert teachers.
2. Demonstrates that providing LLMs with additional context from the curriculum, such as the original warmup task, helps to ground the model and improve the quality of the generated warmup tasks.
3. When prompted well, LLMs can generate warmup tasks that are of similar or better quality to those created by expert teachers in terms of alignment to learning objectives, accessibility, and teacher preference.

### **Implications for practice and/or policy**

1. Provides practical suggestions for prompting LLMs to generate high-quality warmup tasks for middle-school math teachers, such as incorporating additional curriculum context and expert-informed prompts.
2. Demonstrates how CTA with expert teachers can be used to develop LLM-based tools for educators that align with their practices and preferences.
3. Additional research is needed to explore the potential for LLMs to support other types of curriculum adaptations, evaluate their effectiveness in real classroom settings, and investigate how they can be designed to effectively tailor to the specific needs and characteristics of individual students.

# 1 Introduction

Scaffolding is an instructional process that provides support tailored to the needs of students that enables them to “solve problems, complete tasks, or achieve educational goals that would be challenging without such support” (Wood et al., 1976). Research consistently shows that scaffolded instruction benefits all students, whether they have documented learning needs or not, leading to improved academic and non-academic outcomes (CAST, 2024; Goddard et al., 2015). The need for scaffolding has been magnified in the post-pandemic era. The national incidence of students with Individualized Education Plans (IEPs) and 504 plans has increased from 13% to 15%, with some districts reporting rates as high as 35% (Irwin et al., 2023). Furthermore, an alarming 49% of students are now below grade level in at least one subject (Irwin et al., 2023).

In recent years, the use of standardized curricula in math education has proliferated to support teachers in providing high-quality instruction (Remillard & Heck, 2014; Jackson & Makarin, 2018; Doan et al., 2022). However, the post-pandemic era has highlighted the limitations of these curricula in addressing the diverse needs of learners (Nerlino, 2022). With more students working below grade level or requiring special support, standardized curricula often lack the necessary scaffolding for accessibility and engagement. Curriculum scaffolding, a core aspect of instructional scaffolding, involves teachers adapting and supplementing the official curriculum to bridge this gap, particularly for struggling learners (Remillard & Heck, 2014; Remillard, 1999; Squire et al., 2003). This may include adding warmup tasks, formative assessments, or culturally relevant adaptations (Meidl & Meidl, 2011; Walkington, 2013).

Despite the documented effectiveness of and need for scaffolding, a majority of educators encounter significant challenges in implementing it, and there’s little work on how teachers approach the scaffolding process (Sherin & Drake, 2009). One of the primary challenges is the time-intensive nature of creating effective scaffolds, which requires teachers to carefully analyze the curriculum, identify potential barriers for learners, and design appropriate adap-

tations and supplements. This process can be particularly challenging for novice teachers who may lack the experience and pedagogical content knowledge to efficiently navigate the complexities of scaffolding (Lokey-Vega, 2015). Additionally, many teachers receive limited training on how to effectively scaffold instruction, highlighting the need for more targeted teacher education in this area (Yan & Goh, 2023; Galiatsos et al., 2019).

This research investigates the potential for Large Language Models (LLMs) to support curriculum scaffolding in middle-school mathematics. Effectively exploring this potential first requires a more granular understanding of how teachers approach the scaffolding process. Thus, we start with a qualitative research phase, conducting Cognitive Task Analysis (CTA) to gain insights into how expert teachers scaffold their curriculum. Building on our findings and existing frameworks (Remillard & Heck, 2014), we propose a three-stage process of curriculum scaffolding: observation, strategy formulation, and implementation. We apply this framework to a case study of using LLMs to generate warmup tasks, a specific type of instructional scaffold that activates and refreshes background knowledge. To assess LLMs' ability to produce warmup tasks of a similar quality as those made by experts, we develop a novel dataset of expert-created warmup tasks and prepared various LLM-based approaches informed by insights from our CTA. We then evaluate these LLM-based approaches against the expert-created tasks using considerations identified in the CTA as evaluation criteria.

Our research aims to explore the following interdependent questions:

1. How do expert teachers scaffold their curriculum, and which aspects of their process might AI facilitate?
2. How do instructional scaffolds created by LLMs compare in quality to those created by expert teachers?

## 2 Related Work

### 2.1 Conceptions of Curriculum Use

A large area of educational research centers on curriculum use: How do teachers interact with and draw upon curriculum resources designed to guide instruction (Remillard, 2005; Stein & Smith, 1998)? Recent frameworks have sought to provide a more nuanced understanding of the curriculum implementation process by distinguishing between the official curriculum and the operational curriculum (Remillard & Heck, 2014). The official curriculum, also known as the ‘formal curriculum’ (Gehrke et al., 1992), comprises curricular aims and objectives, content of consequential assessments, and the designated curriculum (Remillard & Heck, 2014). In contrast, the operational curriculum encompasses what actually happens in the classroom, including the teacher-intended curriculum, the enacted curriculum, and student outcomes (Remillard & Heck, 2014; Gehrke et al., 1992).

The transition from the official to the operational curriculum involves teachers interpreting and making decisions when designing instruction, drawing upon the designated curriculum and other resources. This process, termed “documentational genesis,” results in the generation of new documents tailored to specific students at a particular moment (Gueudet & Trouche, 2009). Sherin & Drake (2009) proposed a three-step process of reading, evaluating, and adapting, while Brown (2002) described three modes of engagement: offloading, adapting, and improvising. This process of curriculum enactment is highly dependent on teacher characteristics such as subject matter knowledge, pedagogical content knowledge, beliefs, and goals, as well as the local context (Brown, 2002; Remillard, 2005; Gonzalez Thompson, 1984; Sherin & Drake, 2009).

As noted by Sherin & Drake (2009), there have been relatively few empirically-based generalizations focused on the transition from the official curriculum to the operational, and even less from the designated to the teacher-intended. This is partly due to the challenge of accessing the teacher-intended curriculum, given that it typically exists within individual

teacher minds, files and presentations and is not centrally collected or stored (Remillard & Heck, 2014). Our study aims to contribute to this gap by examining the strategies and considerations expert teachers employ in scaffolding curriculum for diverse learners and exploring the potential for AI to support this process. By focusing on the transition from the designated curriculum to the teacher-intended curriculum and the specific adaptations and supplements teachers make to ensure accessibility and engagement for all students, we seek to provide a more nuanced understanding of curriculum enactment in contexts of learner variability.

## 2.2 Scaffolding in Math Education

The effectiveness of scaffolding in math education is closely tied to the transition from the official to the operational curriculum. The concept of scaffolding refers to the temporary but essential support provided to learners to help them achieve tasks that would otherwise be beyond their reach (Wood et al., 1976). Grounded in Vygotsky’s sociocultural theory and the concept of the Zone of Proximal Development (ZPD), scaffolding emphasizes the importance of providing learners with the necessary guidance and support to move from assisted to independent performance (Vygotsky & Cole, 1978).

In the context of math education, scaffolding has been shown to be a highly effective instructional approach, with meta-analyses reporting significant effect sizes on student outcomes (Hattie, 2008; Zuo et al., 2023). The benefits of scaffolding extend beyond academic achievement, encompassing increased task effort, cognitive development, metacognitive awareness, independence, sensemaking, and self-confidence (Zuo et al., 2023).

Well-designed, standards-aligned instructional materials often remain inaccessible to students if the tasks fall outside their ZPD. Consequently, teachers often devote significant time and effort to adapting and supplementing their core curriculum materials to better suit their students’ learning needs (Philipp & Kunter, 2013). This process of curriculum scaffolding,

which we define as a subset of the teacher-intended curriculum, encompasses the various techniques and strategies teachers employ to ensure that all students can access and engage with the official curriculum.

These techniques include providing hints, modeling, asking probing questions, employing the gradual release of responsibility model, guided practice, visual aids, manipulatives, graphic organizers, supportive resources, and offering alternative representations (Drew, 2022). By strategically using these techniques, teachers can tailor their support to the diverse needs of their students, helping each learner navigate the path towards understanding and independence.

### **2.3 Automated Creation of Educational Materials**

The rapid advancement of Large Language Models (LLMs), such as GPT-4, has opened up new avenues for the automated creation of educational materials. These models, trained on vast amounts of text data, have demonstrated the ability to generate human-like text in response to given inputs, offering potential applications in various educational contexts. The proliferation of AI tools and chatbots has led to increased usage among educators, with nearly 50% now using ChatGPT at least once a week (Impact Research, 2024).

Recent studies have explored the use of LLMs in generating a wide range of educational materials. LLMs have been employed to author learning objectives, ensuring alignment with course content and desired outcomes (Sridhar et al., 2023). They have also been used to write worked examples and explanations, providing step-by-step guidance for learners (Jury et al., 2024; Prihar et al., 2023). In the domain of question generation, LLMs have demonstrated the ability to create question-answer pairs, multiple-choice questions, and open-ended questions across various subjects (Rodriguez-Torrealba et al., 2022; Elkins et al., 2023; Z. Wang et al., 2022; Shimmei et al., 2023; Bulathwela et al., 2023; Doughty et al., 2024). In the context of math education, LLMs have been utilized to generate problems at

varying levels of difficulty and to adapt existing problems for improved student understanding (Jiao et al., 2023; Norberg et al., 2023). LLMs have also been employed to generate entire course content, including syllabi, lectures, and assessments (Leiker et al., 2023; Diwan et al., 2023). This line of work suggests that LLMs could serve as powerful tools for creating instructional resources, potentially saving teachers time and effort in developing materials tailored to their students’ needs.

However, the use of LLMs in educational contexts is not without challenges. Researchers have noted that these models can sometimes generate unreliable solutions to math problems (Frieder et al., 2024) and may “hallucinate” information, producing content that appears plausible but is not actually accurate (Ji et al., 2023). There is also evidence that LLMs alone do not behave like expert instructors, such as when providing pedagogical feedback to teachers (R. Wang & Demszky, 2023), while remediating mathematical mistakes in tutoring (R. E. Wang et al., 2024), or providing pedagogical explanations (Jury et al., 2024; Prihar et al., 2023). These limitations highlight the need for careful consideration and human oversight when employing LLMs in the creation of educational materials.

One promising way to both ensure safety and pedagogically sound outputs is to create AI-in-the-loop systems where the teacher is in control of what is being sent to students (Ninaus & Sailer, 2022). For example, in a study exploring the use of LLMs for remediating student math mistakes — a type of scaffolding, R. E. Wang et al. (2024) demonstrated that LLM responses to students improve significantly when the teacher helps identify the cause of error and response strategy to use.

## **2.4 Modeling Expert Decision-Making**

Understanding the decisions instructional experts make while performing a task is critical to designing technological systems that seek to model such expertise and enhance teachers’ work. Cognitive Task Analysis (CTA) is a widely used structured, qualitative research

method for eliciting and formalizing the knowledge, thought processes, mental strategies, and goal structures that underlie expert performance (Clark et al., 2008). In recent years, researchers have increasingly applied CTA in educational settings to gain insights into the complex decision-making processes of expert teachers. For example, Lokey-Vega (2015) used CTA to detail a nine-step process that experts follow when designing and implementing technology-rich lessons, revealing that novice teachers were less familiar with distinct parts of this process. Similarly, R. E. Wang et al. (2024) employed CTA to develop a framework for expert teacher decision-making in the context of responding to student misconceptions, formalizing a process in which the expert identifies the student’s error, determines a remediation strategy, and articulates their instructional intention before generating a response.

CTA builds upon and extends traditional ethnographic research methods by incorporating not only observations but also verbal statements from experts as a primary source of information. Through the use of interview techniques, verbal reports, and the analysis of team communication, researchers can elicit valuable insights into the cognitive processes that guide expert performance. In the context of curriculum scaffolding, CTA offers a powerful tool for identifying the knowledge and strategies that expert teachers use to tailor curriculum materials and create effective scaffolds for diverse learners. By formalizing these often tacit processes, CTA can help to make expert teachers’ decision-making more explicit and accessible, potentially informing the design of professional development programs and instructional resources.

Moreover, recent research has demonstrated the potential for integrating insights from CTA into the development of Large Language Models (LLMs) to generate educational materials that align with expert teachers’ practices and preferences. R. E. Wang et al. (2024) found that incorporating expert decision-making models derived from CTA into LLMs led to outputs that were more favorably evaluated by teachers compared to those generated by LLMs without such models. This suggests that LLMs informed by CTA can become valuable tools for generating educational outputs that resonate with educators and reflect best

practices in instructional design.

### 3 Cognitive Task Analysis for Curriculum Scaffolding

To systematically uncover how expert middle-school math teachers scaffold their curriculum to meet the diverse needs of their students (RQ1), we conducted CTA with six expert teachers. The CTA allowed us formalize teachers’ curriculum scaffolding as a three-stage process, which we describe in Section 3.2.

#### 3.1 Methods

**Participants.** We recruited teachers from two public school districts in Washington and Chicago. Each teacher was selected by district administrators based on their extensive teaching experience (minimum of 10 years) and perceived expertise, a method that has been shown to effectively identify highly skilled teachers (Jacob & Lefgren, 2008). The participants had experience working in a range of public schools and student contexts, ensuring a diverse set of perspectives. All teachers were compensated for their time at a rate of \$45 per hour. The number of experts included in this study is comparable to other NLP studies that work closely with domain experts (R. E. Wang et al., 2024; Sharma et al., 2023).

**Data collection.** Data collection consisted of two main components: weekly surveys and CTA interviews. Each week, the expert teachers completed a short survey in which they provided (a) the official (designated) curriculum materials they had access to via their district for the lessons they were teaching that week, (b) the teacher-intended curriculum materials they specifically prepared to use in their classrooms, and (c) brief descriptions of the changes made to the original materials. This approach allowed us to capture the transition from the official curriculum to the teacher-intended curriculum, shedding light on the adaptations and scaffolds created by the teachers.

Prior to each CTA interview, the research team reviewed the official and teacher-intended curriculum materials shared by the teachers and formulated hypotheses about the adaptations made. The interviews were then structured around three main questions: (1) What student responses to the curriculum did the teachers expect? (2) How did they want to respond to these anticipated responses? (3) Why did they choose to respond in that particular manner? These questions were designed to elicit insights into the complex decision-making processes teachers employed when determining when, how, and why to adapt and scaffold their curriculum materials.

During the interviews, teachers were prompted to provide detailed explanations of their thought processes and reasoning behind the changes they made to the official curriculum materials. The interviewers used follow-up questions and probes to encourage teachers to elaborate on their responses and to clarify any ambiguities.

**Data analysis.** We analyzed the data collected through the weekly surveys and CTA interviews using a qualitative, iterative approach. First, we reviewed the original and adapted curriculum materials to identify patterns in the types of changes made by the expert teachers. We then compared these initial observations against teachers' own descriptions of their adaptations and refined based on the insights gained from the CTA interviews. Next, we coded the interview data using a combination of deductive and inductive coding techniques. Deductive codes were derived from existing literature on curriculum use and scaffolding, while inductive codes emerged from the data itself, capturing the unique strategies and considerations employed by the expert teachers. We organized the coded data into themes and sub-themes, representing the key components of the teachers' decision-making processes.

Throughout the analysis process, the research team engaged in regular discussions to ensure consistency in coding and interpretation. We resolved discrepancies through consensus, and iteratively refined the coding scheme to better capture the nuances of the data. We then used the final themes and sub-themes to construct a framework describing the expert

teachers' approach to curriculum scaffolding.

**Limitations.** The small sample size ( $n=6$ ) and the focus on two specific school districts may limit the generalizability of the findings to other contexts. Additionally, the reliance on self-reported data from the teachers may have introduced some degree of bias (Montibeller & Von Winterfeldt, 2015).

## 3.2 Findings

Based on our findings from the CTA and building on existing frameworks of curriculum use (Sherin & Drake, 2009; Brown, 2002), we propose a three-stage framework for understanding how expert teachers scaffold curriculum materials to meet the needs of students who are struggling to access and engage with the content. This framework comprises three key stages: observation, strategy formulation, and implementation.

### 3.2.1 Observation: Assessing Curriculum Materials and Student Needs

The curriculum scaffolding process begins with teachers making observations about the existing curriculum materials and assessing how well they align with their professional intentions and the needs of their students. Consistent with prior research (Sherin & Drake, 2009; Collopy, 2003), our CTA revealed that teachers often read curriculum materials with their students in mind, making determinations about how to use and adapt suggested activities based on their perceptions of students' needs and deficits.

The most common observation made by the expert teachers in our study was the significant gap between the curriculum materials and the current level of their students, particularly in low-income communities with high percentages of students working below grade level. Teachers noted that despite wanting to use the curriculum, they were acutely aware that students would not be able to access the materials as presented, potentially leading to disengagement and behavioral issues.

Other observations included concerns about the curriculum’s presentation of content, even when the content itself was deemed accessible. These concerns often intersected with teachers’ professional and pedagogical beliefs about effective instruction. For example, some teachers noted that their curriculum relied heavily on open-ended questions and student discussions, with relatively little time spent on explicitly modeling processes and algorithms. In these cases, teachers opted to increase opportunities for direct instruction within their lessons.

Some observations were deeply rooted in teachers’ understanding of their students’ abilities, interests, and learning needs. Teachers explained that certain phrasing, topics, or cultural references in the curriculum might not be relevant to their students, or that certain activities would not be effective for students with specific learning needs. In other cases, teachers’ observations were driven more by their professional judgments than their understanding of their students. For instance, if a teacher believed there was a more effective way to introduce a topic than what was provided in the curriculum, they would be inclined to implement that method.

School-level dynamics and priorities also influenced teachers’ observations about the curriculum materials. Efforts to increase the use of formative assessment or teach in a culturally relevant manner, for example, led some teachers to note that the curriculum did not fully reflect these priorities.

### **3.2.2 Strategy Formulation: Developing Approaches to Address Observations**

After making observations about the curriculum materials, teachers formulate strategies to address each of these observations. These strategies, typically developed at the lesson level, reflect teachers’ professional experience and knowledge of what has worked well for them in the past.

Our CTA identified several common strategies employed by expert teachers, including:

- **Activating and refreshing students' background knowledge:** For example, if a lesson on fractions is planned, a teacher might include a warmup task that reviews basic concepts of division and multiplication, ensuring students have the necessary foundational knowledge.
- **Supporting the decoding of text, mathematical notation, and symbols:** Teachers might add glossaries or visual aids to help students understand complex terminology or symbols, making the content more accessible.
- **Incorporating additional formative assessment:** Teachers often include formative assessment activities, such as quick quizzes or exit tickets, to gauge student understanding and adjust instruction accordingly.
- **Explicitly modeling necessary skills and procedures:** For instance, if students struggle with solving equations, teachers might provide step-by-step demonstrations and think-alouds to model problem-solving processes.
- **Providing opportunities for additional practice:** Teachers might include fluency drills or extra practice problems to reinforce key skills and concepts.
- **Guiding student information processing:** This could involve creating structured note-taking guides or graphic organizers to help students capture and organize information effectively.

These strategies were often grounded in learning sciences and reflected professional development within the school, emphasizing evidence-based practices to enhance student learning outcomes. The strategy formulation stage aligns with previous research suggesting that teachers engage in a process of evaluation and adaptation when using curriculum materials (Sherin & Drake, 2009). By developing strategies to address the gaps and challenges identified during the observation stage, teachers demonstrate their pedagogical design ca-

capacity (Brown, 2002) and their ability to craft instructional approaches that support student learning.

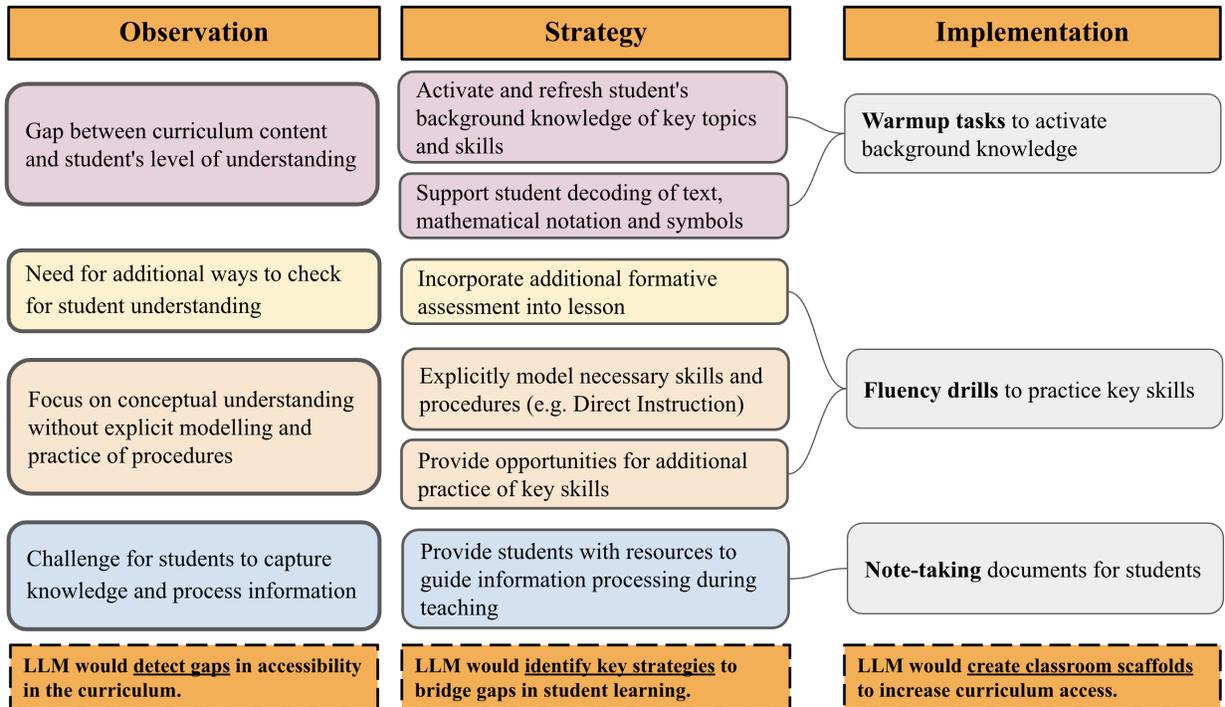
### **3.2.3 Implementation**

The final stage of the curriculum scaffolding process involves implementing the strategies developed in the previous stage. This implementation occurs primarily at the resource level, with teachers describing two main approaches: adaptation and supplementation of existing curriculum materials.

Adaptation often involved light-touch edits and modifications to the original materials, such as adding definitions, reducing cognitive demand, and making other minor tweaks to improve accessibility and relevance for students. In some cases, teachers described more substantial adaptations, such as restructuring an instructional task or changing its format. Adaptation was more common when the original curriculum included resources that aligned with the teacher’s desired strategy, such as a glossary of key terms that could be extended to support student decoding of text and symbols.

Supplementation, on the other hand, usually occurred in response to identified gaps or missing elements in the original materials. Teachers would create or source additional resources to fill these gaps, such as developing warmup tasks to review relevant background knowledge, creating note-taking guides for students, or crafting formative assessment questions. Supplementation often involved drawing upon third-party repositories like Teachers Pay Teachers to find suitable materials.

Figure 1 summarizes the most common observations made by the experts, the associated strategies they employed, and illustrative examples of how they implemented this strategy in the lesson. We note that multiple strategies can be implemented via a single classroom scaffold. For example, a warmup task could be used to activate relevant background knowledge and to build student understanding of important language that will be used in the lesson. Appendix A also provides additional examples of curriculum scaffolding strategies



**Figure 1:** Framework showing the levels of curriculum scaffolding found in the CTA. Expert teachers described first making observations about the existing curriculum materials, then crafting strategies to address these, and finally implementing these strategies through adaptations or modifications at the resource level. LLMs can potentially support teachers at each of these steps.

and implementation approaches.

### 3.3 Opportunities for AI Augmentation

The findings from our CTA highlight several aspects of the curriculum scaffolding process that could potentially be augmented through AI approaches. At each stage of the framework, we see opportunities for LLMs to support teachers in their efforts to create accessible and engaging instructional materials.

At the observation stage, LLMs could assist teachers by analyzing curriculum materials to identify potential gaps and areas of challenge. Similar approaches have been used to identify student misconceptions in programming (Fwa, 2024) and math education (Smart et al., 2024). By providing teachers with preliminary observations based on student performance

data and patterns in student learning over time, LLMs could save teachers time and offer a foundation upon which to build their instructional strategies.

At the strategy formulation stage, LLMs could be designed to suggest strategies that address the observations made by teachers. For example, if a gap in background knowledge is identified, the LLM could propose strategies to activate prior knowledge. These suggestions could be grounded in learning sciences and professional development best practices, ensuring pedagogical soundness. LLMs could also potentially personalize strategy recommendations based on specific student profiles and learning needs, similar to how they have been used to interpret and explain students' wrong answers (Smart et al., 2024).

At the implementation stage, LLMs could generate scaffolded instructional materials, such as modified tasks or supplementary worksheets, aligned with the teacher's desired strategies. By incorporating expert decision-making models, LLMs could provide suggestions that mimic the nuanced adaptations and supplements typically employed by expert teachers. For example, LLMs could create warmup tasks that review relevant prior content or generate formative assessment questions tailored to the lesson objectives. This aligns with recent work to use LLMs to generate a wide range of educational materials, as outlined in 2.3.

These opportunities for AI augmentation resonate with Brown (2002)'s suggestions for supporting teachers' pedagogical design capacity. Brown proposes that teachers should receive help in evaluating the features and affordances of curriculum materials and identifying necessary modifications to align these materials with instructional goals. The use of LLMs to provide this kind of support could not only save teachers time and effort but also potentially serve as a form of embedded professional development, helping teachers to refine their understanding of instruction and student learning.

## 4 Case Study: Using LLMs to Generate Curriculum Scaffolds

Building on the three-stage framework obtained through our CTA, we conducted a case study to understand how LLMs might be used at the *implementation* stage of curriculum scaffolding, and how these scaffolds compare to those created by experts (RQ2). We restrict our case study to only one stage in order to develop a proof-of-concept with a simpler scenario that does not require jointly evaluating modeling decisions for the three stages at once. Among the three stages, we chose to prioritize implementation because this step tends to be the most time intensive for teachers (Grossman & Thompson, 2008) and because this stage produces observable outputs that are directly comparable to expert implementations. Thus, we draw on a pre-defined observation and scaffolding strategy for the implementation task: (1) We focus on the *observation* that the designated curriculum is often inaccessible to a significant proportion of students, and that there exists a gap between the curriculum content and students' present level of understanding, (2) we use a common *strategy* identified during the CTA: the intention to activate and refresh students' background knowledge of key topics and skills. We leave the task of evaluating LLMs performance at producing observations and scaffolding strategies for future work.

The specific implementation example we focus on is a warmup task, which not only emerged as a common strategy in the CTA but also has a strong basis in learning science research. The Universal Design for Learning framework and the National Council of Teachers of Mathematics (NCTM) emphasize the importance of building on students' prior knowledge and experiences (CAST, 2024; National Council of Teachers of Mathematics, 2016). NCTM recommends starting lessons with a review of previously learned skills, using a 5-10 minute warmup activity, such as a discussion prompt or a mathematical task. For example, a lesson on the Pythagorean Theorem might begin with a review of square numbers, surds, and the properties of right-angled triangles. Given its prevalence and importance, warm up tasks

serve as an ideal first case study for research-driven approaches to developing and testing algorithms that create specific classroom scaffolds.

## 4.1 Methods

### 4.1.1 Data Collection

Creating a comparison set for LLM outputs required us to collect example implementations of the teacher-intended curriculum. As noted in 2.1, this data is challenging to collect as it is not centrally stored and often exists within individual teacher minds, files and presentations. As such, we collaborated with two expert teachers from school districts in Washington and Illinois to collect this novel dataset. Each expert has a minimum of 10 years of teaching experience, deep familiarity with the Illustrative Mathematics curriculum, and roles as subject and curriculum leaders within their schools, where they mentor other teachers. These qualifications help to ensure that the dataset is both robust and reflective of high-quality instructional practices. The teachers were compensated at a rate of \$50 per hour for their contributions.

The data collection process began with selecting ten lesson plans from Units 2-7 of the Illustrative Mathematics 6th-grade curriculum. These units target common misconceptions in algebra, categorized into four primary areas: ratios and proportional relationships, the number system, expressions and equations, and functions (Bush & Karp, 2013). The lessons were evenly distributed across these categories, covering 21 of the 26 unique standards within these units. This selection was deliberate, focusing on critical areas of need and spanning a broad range of skills. It is noteworthy that while there is a warmup in every Illustrative Mathematics lesson, the CTA revealed that these often assumed students were on grade level and sometimes did not effectively activate background knowledge. Therefore, modification was necessary to better meet the needs of students working below grade level.

Each expert was provided with the context of these ten lesson plans. The context included

comprehensive details such as the lesson narrative, learning goals, standards, and specific instructional routines and activities. Importantly, the experts were instructed to assume they were teaching the lesson to a class where 50% of the students were working below grade level. They were tasked with designing a warmup task that activates background knowledge, concise enough to fit on a single Google Slide, aligning with common practices for creating instructional materials. Additionally, they were asked to provide a brief commentary on the rationale for their created task, including how, if at all, they used the original curriculum warmup in their version. This commentary captures the teacher’s resource-level decision-making process.

The dataset comprises 20 items in total, with each lesson plan having two modified warmup tasks (one from each expert). Figure 2 illustrates an example of a lesson plan context, the expert-created warmup task, and the commentary provided by the teachers.

#### 4.1.2 Model Development

With expert examples of the teacher-intended curriculum for this specific classroom artifact, we began to develop LLM-based approaches to generate similar warmup tasks. We prepared three distinct approaches, each informed extensively by the findings from the CTA:

1. **Expert Informed Prompt:** In this condition, each model was provided with the learning goals and lesson narrative. The prompt included general best practices for creating effective warmups, incorporating insights gathered from the CTA. For example, the model was encouraged to consider the skills and knowledge that would be important for that particular lesson, mimicking the thought process that expert teachers described.
2. **Additional Curriculum Context:** In addition to the learning goals and lesson narrative provided in the Expert Informed Prompt condition, the model was given the original warmup task from the Illustrative Mathematics curriculum. This condition

Illustrative Mathematics  
LEARNING WITH PURPOSE

Grade 6 Unit 3 Lesson 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17

LESSON 6  
**Interpreting Rates**

PREPARATION LESSON PRACTICE [View Student Lesson](#)

**Lesson Narrative**  
In previous lessons students have calculated and worked with rates per 1. The purpose of this lesson is to introduce the two **unit rates**,  $\frac{a}{b}$  and  $\frac{b}{a}$ , associated with a ratio  $a : b$ . Each unit rate tells us how many of one quantity in the ratio there is per unit of the other quantity. An important goal is to give students the opportunity to see that both unit rates describe the same situation, but that one or the other might be preferable for answering a given question about the situation. Another goal is for students to recognize that they can just divide one number in a ratio by another to find a unit rate, rather than using a table or another representation as an intermediate step. The development of such fluency begins in this section and continues over time. In the Cooking Oatmeal activity, students have explicit opportunities to justify their reasoning and critique the reasoning of others (MP3).

**Learning Goals**  
Teacher Facing  
Calculate and interpret the two unit rates associated with a ratio, i.e.,  $\frac{a}{b}$  and  $\frac{b}{a}$  for the ratio  $a : b$ .  
Choose which unit rate to use to solve a given problem and explain the choice (orally and in writing).  
Comprehend the term "unit rate" (in spoken and written language) refers to a rate per 1.

**Student Facing**  
Let's explore unit rates.

**Learning Targets**  
**Student Facing**  
I can choose which unit rate to use based on how I plan to solve the problem.  
When I have a ratio, I can calculate its two unit rates and explain what each of them means in the situation.

Reformat the [original] warm-up as a 'Think, Pair, Share' and add a lower level first question that supports all students to access the task.

**INTERPRETING RATES Warm-Up:**

THINK-

1. Think of things that you can measure, record, or collect information about with numbers.
2. Now use those ideas to think of two things you have heard described in terms of "something per something."

PAIR-

3. Share your ideas with your group, and listen to everyone else's idea. Make a group list of all unique ideas.

SHARE-

4. Pick a reporter to share these with the class.

**Figure 2:** For each lesson context (top-left), the expert teachers were given a pre-defined observation (curriculum inaccessibility) and scaffolding strategy (activating background knowledge). They implemented this strategy by creating a student-facing warmup resource (bottom) and provided commentary on their rationale and use of the original curriculum warmup, capturing their resource-level decision-making process (top-right).

tested the model's ability to utilize existing curriculum components to generate a related but distinct educational task.

**3. Expert Implementation Guidance Provided:** As detailed in 4.1.1, the expert teachers involved in data collection provided commentary focused on the implementation stage, justifying how and why they either adapted or supplemented the original

warmup in the curriculum, in line with the process outlined in Section 3.2.3. In this third condition, the model received the original warmup task and the expert teacher’s commentary on implementation. This condition tested the model’s capability to integrate human expert insights to refine or alter educational content meaningfully. The rationale for this condition was to explore whether providing the model with explicit guidance on how to implement the scaffolding strategy, as opposed to just the original task itself, would lead to higher-quality outputs.

The full prompts used for each condition are documented in Appendix C for reference.

### 4.1.3 Evaluation Setup

Our evaluation methodology employs the comparative judgment technique, widely recognized for its reliability in educational research (Pollitt, 2012).

**Evaluators.** To facilitate a comprehensive assessment, we engaged evaluators from two distinct sources: a group of 19 current or former middle-school math teachers recruited through partnerships with two school districts in Washington and Illinois and through personal and institutional networks, and an additional 31 active teachers recruited through the Prolific platform. Prolific participants were screened to ensure they were based in the US or UK, were current teachers, had at least an undergraduate degree, and had their degree in mathematics or statistics. Although it was not possible to filter directly for math teachers, we assumed it was likely they had experience teaching math if they met these criteria.

Table 1 summarizes the sample of evaluators based on self-reported information we collected at the beginning of the task. The sample consisted primarily of current math teachers, with 74% actively teaching math. An additional 18% were former math teachers, while the remaining 8% were teachers of another subject with a math degree. In terms of teaching experience, 52% had more than 10 years of experience, 28% had 5-10 years, and 20% had less than 5 years. The majority of evaluators taught at the high-school level (78%), followed

by middle-school teachers (61%), and a smaller proportion teaching at the elementary level (13%).

**Table 1:** Demographic Information of Evaluators

Category	Percentage
Current Math Teacher	74%
Former Math Teacher	18%
Never Taught Math	8%
More than 10 years of teaching experience	52%
5 - 10 years of teaching experience	28%
Less than 5 years of teaching experience	20%
High-School Level	78%
Middle-School Level	61%
Elementary Level	13%

**Evaluation task.** Each evaluator completed 10 paired comparisons, one for each of 10 distinct lessons. These sixth grade mathematics lessons were derived from our dataset (Section 4.1.1). In each comparison, they evaluated the expert-made warmup against one of four warmups: the original curriculum warmup and three LLM generated warmups, one from the Expert-Informed Prompt condition, one from the Additional Curriculum Context condition and one from the Expert Implementation Guidance Condition. The original curriculum warmup was included as a baseline comparison, to assess whether evaluators preferred the unmodified materials from Illustrative Mathematics to the versions modified by the expert teacher. To mitigate potential biases, we randomized the order of the lessons, the selection of the expert warmup (expert 1 or expert 2) for each lesson, and the comparison condition presented for each lesson.<sup>1</sup> In total, 500 comparisons were conducted, 125 for the Expert-Informed Prompt condition, 132 for the Additional Curriculum Context condition, 122 for the Expert Implementation Guidance condition, and 121 for the Original Curriculum Warmup condition.

<sup>1</sup>We verified that randomization worked correctly, as each unique model/original warmup item was compared 6.25 times with a standard deviation of 1.196, which is in line with expectations.

To minimize bias in the presentation of the warmups, we performed minimal post-processing of the model-generated warmups for comparison with the expert versions. This involved formatting the model output on a Google Slide with the same template and theme as the expert version, without making any changes to the text. We also included up to one image per slide to maintain aesthetic consistency. For example, if the model output involved questions about paint, we inserted an image of paints. These images served purely aesthetic purposes and had no pedagogical or mathematical relevance, encouraging evaluators to focus on the content rather than the presentation. Figure 3 provides an example of the comparison screen presented to evaluators.

Evaluators assessed the warmups based on the following criteria, chosen for their relevance to practical and effective classroom instruction:

1. **Alignment to Learning Objectives:** Evaluators rated each scaffold based on its alignment to the specific learning objectives, ensuring that the task directly contributes to achieving these goals.
2. **Accessibility for Below-Level Students:** This criterion examines the ease with which below-level students can engage with the task, considering the clarity, simplicity, and support embedded within the scaffold.
3. **Readiness for Classroom Use:** This measures the extent to which a scaffold can be utilized in the classroom without further modification. Evaluators assessed how much additional preparation, adaptation, or modification each scaffold would require before it could be effectively used. Readiness for use is important as it reduces preparation time, allowing teachers to focus more on instruction and interaction with students.
4. **Preference:** Evaluators were asked which warmup they would prefer to use if they were teaching the specific learning objective. This captures the overall practical preference of educators based on their professional judgment.

You are a 6th grade math teacher selecting a warmup to prepare students for the following lesson. Use the lesson title and learning objectives below to answer the questions below.

**Lesson Title:** Color Mixtures

**Learning Objectives:**

- Comprehend and respond (orally and in writing) to questions asking whether two ratios are equivalent, in the context of color mixtures.
- Draw and label a discrete diagram with circled groups to represent multiple batches of a color mixture.
- Explain equivalent ratios (orally and in writing) in terms of the amounts of each color in a mixture being multiplied by the same number to create another mixture that is the same shade.

**Warmup A:**

**Color Mixtures Warm-Up:**

Think about creating a mixture using blue and yellow water. If you use 2 cups of blue water and 3 cups of yellow water:

1. How many cups of each color would you need if you want to make double the amount of the mixture?
2. How many cups of each color would you need if you want to make triple the amount of the mixture?
3. Turn to a partner and discuss: If you use 4 cups of blue water and 6 cups of yellow water, will the shade of green be the same as the original mixture? Why or why not?
4. Share your discussion points with the class.

**Warmup B:**

**Color Mixtures Warm-Up:**

Micah was selling cakes to raise money for graduation. Over last month he sold 34 cakes and charged \$25 for each.

	20	+	5	
30	60		150	210
+				+100
4	80		20	310

1. What does the work Micah did tell us about his situation?
2. Do you agree with Micah's work? What did he do wrong?
3. Would you have solved this problem using the same strategy? If not, what strategy or steps would you use?

**Figure 3:** For each of the ten lessons, evaluators were shown two warmups, one created by an expert and one created in one of our model conditions. The order of expert and model were randomized each time.

Each of these dimensions was assessed using a direct comparison between the expert-created and model-generated warmups. Evaluators indicated their preference on a 5-point Likert scale ranging from “Warmup A is much more [aligned to learning objectives/accessible/ready-

Metric	Expert Informed Prompt	Additional Curriculum Context	Implementation Guidance	Original Curriculum Warmup
...aligned with learning objectives.	0.528***	0.636***	0.467***	-0.331**
...accessible.	0.448***	0.447***	0.410**	-0.099
...ready-to-go.	0.440***	0.371***	0.270*	0.017
...preferred.	0.496***	0.545***	0.410**	-0.124

**Table 2:** Results of Comparative Judgments Between Model-Generated Warmups and Expert-Created Benchmarks

to-go/preferred]” to “Warmup A and B are equally [aligned to learning objectives/accessible/ready-to-go/preferred]” to “Warmup B is much more [aligned to learning objectives/accessible/ready-to-go/preferred]”. These ratings were then coded on a scale from -2 to +2, with -2 and -1 corresponding to Warmup A being much more or somewhat more preferred, 0 indicating equal preference, and +1 and +2 corresponding to Warmup B being somewhat more or much more preferred.

## 4.2 Results

Table 2 summarizes the results of the comparative judgments between the model-generated warmups and the expert-created benchmarks across four dimensions: alignment to learning objectives, accessibility for below-level students, readiness for classroom use, and preference. The results are also visually represented in Appendix B, which shows violin plots for each of the dimensions assessed.

Overall, the Additional Curriculum Context condition performed best across all evaluation criteria. These outputs were significantly preferred to the expert outputs in terms of alignment to learning objectives ( $M=0.636$ ,  $p<0.001$ ), accessibility ( $M=0.447$ ,  $p<0.001$ ), readiness for classroom use ( $M=0.371$ ,  $p<0.001$ ), and overall preference ( $M=0.545$ ,  $p<0.001$ ). However, the magnitude of the difference is moderate, with an average of 0.50 across the criteria on a -2 to 2 point scale.

The Expert Informed Prompt condition also performed well, with outputs significantly

preferred to the expert versions in alignment to learning objectives ( $M = 0.528, p < 0.001$ ), accessibility ( $M = 0.448, p < 0.001$ ), readiness for classroom use ( $M = 0.440, p < 0.001$ ), and overall preference ( $M = 0.496, p < 0.001$ ). The magnitude of the difference is similar to the Additional Curriculum Context condition. A paired t-test indicates that the scores for the Expert Informed Prompt condition are not significantly different from the Additional Curriculum Context condition for alignment to learning objectives ( $t = -0.649, p = 0.517$ ), accessibility ( $t = -0.238, p = 0.812$ ), readiness for classroom use ( $t = 0.163, p = 0.871$ ), and overall preference ( $t = -0.342, p = 0.733$ ).

The Implementation Guidance Provided condition also performed well, with outputs significantly preferred to the expert versions in terms of alignment to learning objectives ( $M = 0.467, p < 0.001$ ), accessibility ( $M = 0.410, p < 0.01$ ), readiness for classroom use ( $M = 0.270, p < 0.05$ ), and overall preference ( $M = 0.410, p < 0.001$ ). A paired t-test indicates that the scores for the Implementation Guidance Provided condition are not significantly different from the top-performing Additional Curriculum Context condition for alignment to learning objectives ( $t = 1.051, p = 0.295$ ), accessibility ( $t = 0.533, p = 0.595$ ), readiness for classroom use ( $t = 1.063, p = 0.290$ ), and overall preference ( $t = 1.037, p = 0.302$ ). We provide hypotheses for why this condition did not outperform the other conditions in the discussion and future work sections.

The Original Curriculum Warmup condition was not significantly preferred to the expert outputs on most of the evaluation criteria, with mean scores close to zero for accessibility ( $M = -0.099, p > 0.05$ ), readiness for classroom use ( $M = 0.017, p > 0.05$ ), and overall preference ( $M = -0.124, p > 0.05$ ). However, for alignment to learning objectives, the expert-created warmups were significantly preferred to the original warmups in the curriculum ( $M = -0.331, p < 0.01$ ). This suggests that while the original curriculum warmups were perceived as comparable to the expert-created ones in terms of accessibility, readiness for use, and overall preference, the expert-created warmups were considered better aligned with the learning objectives of the lessons.

## 5 Discussion

Our research aimed to explore the potential of LLMs in supporting teachers with scaffolding their curricula, specifically focusing on middle-school mathematics. We were particularly interested in the transition from the designated curriculum to the teacher-intended curriculum, and within that, the subset of adaptations and supplementations that teachers employ to ensure that all students can access and engage with the curriculum, which we define as curriculum scaffolding.

Building on CTA with expert teachers, we proposed a three-step process for curriculum scaffolding: observation, strategy formulation, and implementation. This framework aligns with existing literature, such as Sherin & Drake (2009)’s conceptualization of teachers reading, evaluating, and adapting the curriculum and Brown (2002)’s three modes of engagement with curriculum: offloading, adapting and improvising. We extend prior work by providing additional granularity specific to the context of preparing curriculum for middle-school mathematics classrooms where students are struggling to access and engage with the content. Our framework offers insights into the day-to-day implementations of various adaptation strategies in this specific subject and context, some of which are described in Appendix A. For example, we provide additional color to what ‘adapting’ means in Sherin & Drake (2009)’s conceptualization, detailing specific techniques such as activating background knowledge, supporting the decoding of text and symbols, and explicitly modeling necessary skills and procedures. These insights contribute to a more nuanced understanding of curriculum scaffolding in contexts of learner variability.

To apply this framework and investigate the potential of LLMs to produce high-quality curriculum scaffolds, we focused on a specific scenario and strategy identified in our CTA with expert teachers: creating warmup tasks that activate and refresh students’ background knowledge. We presented a novel dataset of expert-created warmup tasks for this purpose and developed several LLM-based approaches informed by the insights gained from the

CTA. Our evaluation, which involved 500 comparisons across three model conditions and the original curriculum warmup by math teachers, found that the LLM-generated warmup tasks performed better than those created by expert teachers. The best-performing approach was the Additional Curriculum Context condition, where the model was provided with the original warmup task from the curriculum in addition to an expert-informed prompt. This finding resonates with the expert teacher approach identified in the CTA, whereby teachers review the original curriculum materials before deciding how to adapt them for their specific context.

Our results highlight the potential for LLMs to support teachers in the curriculum scaffolding process, particularly when provided with relevant context and guidance informed by expert practices. By automating the creation of high-quality instructional scaffolds, such as warmup tasks, LLMs could help alleviate the burden on teachers, allowing them to focus more on direct instruction and interaction with students. Furthermore, our work demonstrates an approach for closely involving teachers in the development and evaluation of LLM-based approaches for education, which the research community has called for (Nazaretsky et al., 2022). By replicating this approach for other specific curriculum scaffolding strategies and implementation examples, researchers can continue to explore the potential of AI to support teachers in creating effective, engaging, and inclusive learning experiences for all students.

Our findings contribute to the growing body of literature on curriculum use and the application of AI in education. By providing a more granular understanding of the curriculum scaffolding process and demonstrating the potential of LLMs to generate high-quality instructional materials, this study lays the groundwork for the development of AI-driven tools that can support teachers in their efforts to meet the diverse needs of their students, particularly in contexts where many learners are struggling to access and engage with the curriculum.

## 6 Limitations and Future Work

While our study provides valuable insights, it is only a first step towards understanding the potential of LLMs to support curriculum scaffolding. One limitation is the reliance on a small set of expert-created warmups for evaluation. Although we collaborated with experienced teachers to develop the dataset, having only two examples of teacher-created warmups per lesson may not fully capture the diverse ways in which warmups can be implemented. Evaluators' preferences might have been influenced by factors beyond the quality of the warmups themselves. To address this limitation, future studies should aim to expand the dataset to include a broader range of expert-created warmups, ensuring a more comprehensive and accurate evaluation of the LLM-generated scaffolds.

Another limitation is that the strategies used in the Expert Strategy condition were derived from the dataset creators rather than the evaluators themselves. It is possible that the evaluators might have had different ideas about how to improve the warmups, and the strategies infused into the model might not fully reflect their approaches. This discrepancy could have affected the model's performance in this condition. To mitigate this issue, future research should consider a more dynamic approach, where teachers review the original warmup, provide their own strategies for improvement, and then evaluate the automatically generated warmup based on their input. Such an interactive AI-in-the-loop system would help ensure that the LLM-generated scaffolds align more closely with the evaluators' perspectives and preferences.

Our study lays the groundwork for several exciting avenues of research. One priority should be to expand the dataset to include a wider variety of instructional scaffolds beyond warmup tasks, as well as observations and strategies to capture all three stages of the scaffolding process. By applying the framework and methodology developed in this study to other types of curriculum adaptations and supplements, researchers can gain a more comprehensive understanding of the potential for LLMs to support teachers in the curriculum

scaffolding process. However, this type of data collection process would need to be dynamic to address the limitations mentioned above.

Additionally, piloting the LLM-generated warmups in real classroom settings would provide valuable insights into their practical effectiveness and help identify areas for further refinement. Collaborative research with teachers, involving the use of AI-generated scaffolds in their day-to-day practice, could shed light on the challenges and opportunities associated with integrating these tools into existing instructional workflows. Moreover, conducting focus groups and interviews with students to qualitatively understand their perception of the materials, and measuring outcomes such as their relationship to math and academic performance, could help assess the impact of these scaffolds on learners, similarly to what other studies on scaffolding have done.

Another promising direction for future research is the exploration of more nuanced and personalized scaffolding strategies. By leveraging real-time feedback from teachers and data on student performance and engagement, LLMs could potentially generate instructional scaffolds that are tailored to the specific needs and characteristics of individual classrooms and learners (Lim et al., 2024). The development of adaptive, data-driven scaffolding tools could represent a significant step forward in supporting teachers' efforts to create inclusive and effective learning experiences for all students.

## **7 Acknowledgments**

We would like to thank Stanford's Center for Human-Centered AI and the Stanford Accelerator for Learning for the funding that has supported this work.

## References

- Brown, M. W. (2002). *Teaching by design: Understanding the intersection between teacher practice and the design of curricular innovations*. Northwestern University.
- Bulathwela, S., Muse, H., & Yilmaz, E. (2023). Scalable educational question generation with pre-trained language models. In N. Wang, G. Rebolledo-Mendez, N. Matsuda, O. C. Santos, & V. Dimitrova (Eds.), *Artificial intelligence in education* (pp. 327–339). Cham: Springer Nature Switzerland.
- Bush, S. B., & Karp, K. S. (2013). Prerequisite algebra skills and associated misconceptions of middle grade students: A review. *The Journal of Mathematical Behavior*, 32(3), 613–632. doi: 10.1016/j.jmathb.2013.07.002
- CAST. (2024). *Universal design for learning guidelines*. <https://udlguidelines.cast.org/>. (Accessed 15-April-2024)
- Clark, R., Feldon, D., Van Merriënboer, J. J. G., Yates, K., & Early, S. (2008, 01). Cognitive task analysis. *Handbook of Research on Educational Communications and Technology*, 577-593.
- Collopy, R. (2003). Curriculum materials as a professional development tool: How a mathematics textbook affected two teachers' learning. *The elementary school journal*, 103(3), 287–311.
- Diwan, C., Srinivasa, S., Suri, G., Agarwal, S., & Ram, P. (2023). Ai-based learning content generation and learning pathway augmentation to increase learner engagement. *Computers and Education: Artificial Intelligence*, 4, 100110. Retrieved from <https://www.sciencedirect.com/science/article/pii/S2666920X22000650> doi: <https://doi.org/10.1016/j.caeai.2022.100110>

- Doan, S., Kaufman, J. H., Woo, A., Tuma, A. P., Diliberti, M. K., & Lee, S. (2022). How states are creating conditions for use of high-quality instructional materials in k–12 classrooms.
- Doughty, J., Wan, Z., Bompelli, A., Qayum, J., Wang, T., Zhang, J., ... Sakr, M. (2024). A comparative study of ai-generated (gpt-4) and human-crafted mcqs in programming education. New York, NY, USA: Association for Computing Machinery. Retrieved from <https://doi.org/10.1145/3636243.3636256> doi: 10.1145/3636243.3636256
- Drew, C. (2022). *18 scaffolding examples in education*. Helpful Professor. Retrieved from <https://helpfulprofessor.com/scaffolding-examples-in-education/> (Accessed: 2024-05-24)
- Elkins, S., Kochmar, E., Cheung, J. C. K., & Serban, I. (2023). *How useful are educational questions generated by large language models?*
- Frieder, S., Pinchetti, L., Griffiths, R.-R., Salvatori, T., Lukasiewicz, T., Petersen, P., & Berner, J. (2024). Mathematical capabilities of chatgpt. *Advances in Neural Information Processing Systems*, 36.
- Fwa, H. L. (2024). Experience report: Identifying common misconceptions and errors of novice programmers with chatgpt. In *Proceedings of the 46th international conference on software engineering: Software engineering education and training* (p. 233–241). New York, NY, USA: Association for Computing Machinery. Retrieved from <https://doi.org/10.1145/3639474.3640059> doi: 10.1145/3639474.3640059
- Galiatsos, S., Kruse, L., & Whittaker, M. (2019). Forward together: Helping educators unlock the power of students who learn differently. *National Center for Learning Disabilities*.
- Gehrke, N. J., Knapp, M. S., & Sirotnik, K. A. (1992). Chapter 2: In search of the school curriculum. *Review of research in education*, 18(1), 51–110.

- Goddard, Y., Goddard, R., & Kim, M. (2015). School instructional climate and student achievement: An examination of group norms for differentiated instruction. *American Journal of Education*, *122*(1), 111–131.
- Gonzalez Thompson, A. (1984). The relationship of teachers' conceptions of mathematics and mathematics teaching to instructional practice. *Educational studies in mathematics*, *15*(2), 105–127.
- Grossman, P., & Thompson, C. (2008). Learning from curriculum materials: Scaffolds for new teachers? *Teaching and Teacher Education*, *24*(8), 2014–2026. Retrieved from <https://www.sciencedirect.com/science/article/pii/S0742051X0800084X> doi: <https://doi.org/10.1016/j.tate.2008.05.002>
- Gueudet, G., & Trouche, L. (2009). Towards new documentation systems for mathematics teachers? *Educational studies in mathematics*, *71*, 199–218.
- Hattie, J. (2008). *Visible learning: A synthesis of over 800 meta-analyses relating to achievement*. routledge.
- Impact Research. (2024). *Ai chatbots in schools: Findings from a poll of k-12 teachers, students, parents, and college undergraduates*. (Accessed at: <https://8ce82b94a8c4fdc3ea6d-b1d233e3bc3cb10858bea65ff05e18f2.ssl.cf2.rackcdn.com/bf/24/cd3646584af89e7c668c7705a006/deck-impact-analysis-national-schools-tech-tracker-may-2024-1.pdf>)
- Irwin, V., Wang, K., Tezil, T., Zhang, J., Filbey, A., Jung, J., ... Parker, S. (2023). *Report on the condition of education 2023. nces 2023-144* (Tech. Rep.). National Center for Education Statistics.
- Jackson, K., & Makarin, A. (2018, August). Can online off-the-shelf lessons improve student outcomes? evidence from a field experiment. *American Economic Journal: Economic Pol-*

- icy*, 10(3), 226–54. Retrieved from <https://www.aeaweb.org/articles?id=10.1257/pol.20170211> doi: 10.1257/pol.20170211
- Jacob, B., & Lefgren, L. (2008). Can principals identify effective teachers? evidence on subjective performance evaluation in education. *Journal of Labor Economics*, 26(1), 101–136. Retrieved 2024-05-24, from <http://www.jstor.org/stable/10.1086/522974>
- Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., . . . Fung, P. (2023). Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12), 1–38.
- Jiao, Y., Shridhar, K., Cui, P., Zhou, W., & Sachan, M. (2023). Automatic educational question generation with difficulty level controls. In N. Wang, G. Rebolledo-Mendez, N. Matsuda, O. C. Santos, & V. Dimitrova (Eds.), *Artificial intelligence in education* (pp. 476–488). Cham: Springer Nature Switzerland.
- Jury, B., Lorusso, A., Leinonen, J., Denny, P., & Luxton-Reilly, A. (2024). Evaluating llm-generated worked examples in an introductory programming course. In *Proceedings of the 26th australasian computing education conference* (p. 77–86). New York, NY, USA: Association for Computing Machinery. Retrieved from <https://doi.org/10.1145/3636243.3636252> doi: 10.1145/3636243.3636252
- Leiker, D., Finnigan, S., Gyllen, A. R., & Cukurova, M. (2023). *Prototyping the use of large language models (llms) for adult learning content creation at scale*.
- Lim, L., Bannert, M., van der Graaf, J., Fan, Y., Rakovic, M., Singh, S., . . . Gašević, D. (2024). How do students learn with real-time personalized scaffolds? *British Journal of Educational Technology*, 55(4), 1309-1327. Retrieved from <https://bera-journals.onlinelibrary.wiley.com/doi/abs/10.1111/bjet.13414> doi: <https://doi.org/10.1111/bjet.13414>
- Lokey-Vega, A. (2015, March). Expert as the tpack misfit: A cognitive task analysis to map expert-teacher cognitive processes of technology-rich lesson planning. In D. Rutledge &

- D. Slykhuis (Eds.), *Proceedings of society for information technology & teacher education international conference 2015* (pp. 3322–3330). Las Vegas, NV, United States: Association for the Advancement of Computing in Education (AACE). Retrieved from <https://www.learnstechlib.org/p/150461>
- Meidl, T., & Meidl, C. (2011, Apr.). Curriculum integration and adaptation: Individualizing pedagogy for linguistically and culturally diverse students. *Current Issues in Education*, *14*(1). Retrieved from <https://cie.asu.edu/ojs/index.php/cieatasu/article/view/579>
- Montibeller, G., & Von Winterfeldt, D. (2015). Cognitive and motivational biases in decision and risk analysis. *Risk analysis*, *35*(7), 1230–1251.
- National Council of Teachers of Mathematics. (2016). *High expectations*. <https://www.nctm.org/Standards-and-Positions/Position-Statements/High-Expectations/>. ([Online; accessed 15-April-2024])
- Nazaretsky, T., Bar, C., Walter, M., & Alexandron, G. (2022). Empowering teachers with ai: Co-designing a learning analytics tool for personalized instruction in the science classroom. In *Lak22: 12th international learning analytics and knowledge conference* (p. 1–12). New York, NY, USA: Association for Computing Machinery. Retrieved from <https://doi.org/10.1145/3506860.3506861> doi: 10.1145/3506860.3506861
- Nerlino, E. (2022, January). Navigating “the chaos”: teacher considerations while adapting curriculum and instruction during the COVID-19 pandemic. *Qualitative Research Journal*, *22*(4), 433–447. Retrieved 2024-06-19, from <https://doi.org/10.1108/QRJ-02-2022-0026> (Publisher: Emerald Publishing Limited) doi: 10.1108/QRJ-02-2022-0026
- Ninaus, M., & Sailer, M. (2022). Closing the loop—the human role in artificial intelligence for education. *Frontiers in psychology*, *13*, 956798.

- Norberg, K., Almoubayyed, H., Fancsali, S. E., DeLey, L., Weldon, K., Murphy, A., & Ritter, S. (2023, July). Rewriting math word problems with large language models. In *Aied23: Artificial intelligence in education, empowering education with llms workshop*. Tokyo, Japan.
- Philipp, A., & Kunter, M. (2013). How do teachers spend their time? a study on teachers' strategies of selection, optimisation, and compensation over their career cycle. *Teaching and Teacher Education*, *35*, 1–12.
- Pollitt, A. (2012). The method of adaptive comparative judgement. *Assessment in Education: principles, policy & practice*, *19*(3), 281–300.
- Prihar, E., Lee, M., Hopman, M., Kalai, A. T., Vempala, S., Wang, A., ... Heffernan, N. (2023). Comparing different approaches to generating mathematics explanations using large language models. In N. Wang, G. Rebolledo-Mendez, V. Dimitrova, N. Matsuda, & O. C. Santos (Eds.), *Artificial intelligence in education. posters and late breaking results, workshops and tutorials, industry and innovation tracks, practitioners, doctoral consortium and blue sky* (pp. 290–295). Cham: Springer Nature Switzerland.
- Remillard, J. T. (1999). Curriculum materials in mathematics education reform: A framework for examining teachers' curriculum development. *Curriculum Inquiry*, *29*(3), 315–342. Retrieved from <https://onlinelibrary.wiley.com/doi/abs/10.1111/0362-6784.00130> doi: <https://doi.org/10.1111/0362-6784.00130>
- Remillard, J. T. (2005). Examining key concepts in research on teachers' use of mathematics curricula. *Review of educational research*, *75*(2), 211–246.
- Remillard, J. T., & Heck, D. J. (2014). Conceptualizing the curriculum enactment process in mathematics education. *Zdm*, *46*, 705–718.

- Rodriguez-Torrealba, R., Garcia-Lopez, E., & Garcia-Cabot, A. (2022). End-to-end generation of multiple-choice questions using text-to-text transfer transformer models. *Expert Systems with Applications*, *208*, 118258.
- Sharma, A., Rushton, K., Lin, I., Wadden, D., Lucas, K., Miner, A., ... Althoff, T. (2023, July). Cognitive reframing of negative thoughts through human-language model interaction. In A. Rogers, J. Boyd-Graber, & N. Okazaki (Eds.), *Proceedings of the 61st annual meeting of the association for computational linguistics (volume 1: Long papers)* (pp. 9977–10000). Toronto, Canada: Association for Computational Linguistics. Retrieved from <https://aclanthology.org/2023.acl-long.555> doi: 10.18653/v1/2023.acl-long.555
- Sherin, M. G., & Drake, C. (2009). Curriculum strategy framework: investigating patterns in teachers' use of a reform-based elementary mathematics curriculum. *Journal of Curriculum Studies*, *41*(4), 467–500.
- Shimpei, M., Bier, N., & Matsuda, N. (2023). Machine-generated questions attract instructors when acquainted with learning objectives. In N. Wang, G. Rebolledo-Mendez, N. Matsuda, O. C. Santos, & V. Dimitrova (Eds.), *Artificial intelligence in education* (pp. 3–15). Cham: Springer Nature Switzerland.
- Smart, F., Bos, N. D., & Bos, J. T. (2024). Can large language models recognize and respond to student misconceptions? In R. A. Sottolare & J. Schwarz (Eds.), *Adaptive instructional systems* (pp. 288–299). Cham: Springer Nature Switzerland.
- Squire, K. D., MaKinster, J. G., Barnett, M., Luehmann, A. L., & Barab, S. L. (2003). Designed curriculum and local culture: Acknowledging the primacy of classroom culture. *Science Education*, *87*(4), 468-489. Retrieved from <https://onlinelibrary.wiley.com/doi/abs/10.1002/sce.10084> doi: <https://doi.org/10.1002/sce.10084>
- Sridhar, P., Doyle, A., Agarwal, A., Bogart, C., Savelka, J., & Sakr, M. (2023). *Harnessing llms in curricular design: Using gpt-4 to support authoring of learning objectives*.

- Stein, M. K., & Smith, M. S. (1998). Mathematical tasks as a framework for reflection: From research to practice. *Mathematics teaching in the middle school*, 3(4), 268–275.
- Vygotsky, L. S., & Cole, M. (1978). *Mind in society: Development of higher psychological processes*. Harvard university press.
- Walkington, C. A. (2013). Using adaptive learning technologies to personalize instruction to student interests: The impact of relevant contexts on performance and learning outcomes. *Journal of Educational Psychology*, 105(4), 932–945. doi: 10.1037/a0031882
- Wang, R., & Demszky, D. (2023, June). Is chatgpt a good teacher coach? measuring zero-shot performance for scoring and providing actionable insights on classroom instruction. In *18th workshop on innovative use of nlp for building educational applications*.
- Wang, R. E., Zhang, Q., Robinson, C., Loeb, S., & Demszky, D. (2024, July). Bridging the novice-expert gap via models of decision-making: A case study on remediating math mistakes. In *Proceedings of the 2024 conference of the north american chapter of the association for computational linguistics*. Association for Computational Linguistics.
- Wang, Z., Valdez, J., Basu Mallick, D., & Baraniuk, R. G. (2022). Towards human-like educational question generation with large language models. In M. M. Rodrigo, N. Matsuda, A. I. Cristea, & V. Dimitrova (Eds.), *Artificial intelligence in education* (pp. 153–166). Cham: Springer International Publishing.
- Wood, D., Bruner, J. S., & Ross, G. (1976). The role of tutoring in problem solving. *Journal of Child Psychology and Psychiatry*, 17(2), 89–100.
- Yan, J., & Goh, H. H. (2023). Exploring the cognitive processes in teacher candidates' collaborative task-based lesson planning. *Teaching and Teacher Education*, 136, 104365.
- Zuo, M., Kong, S., Ma, Y., Hu, Y., & Xiao, M. (2023). The effects of using scaffolding in online learning: A meta-analysis. *Education Sciences*, 13(7), 705.

# Appendix A Curriculum Scaffolding Strategies and Implementation Approaches

---

## Curriculum Scaffolding Implementation Examples

### Strategy

---

**Increase relevance, value, authenticity**

1. Use motivating questions to engage students from the start of the lesson.
2. Introduce the lesson with a hook to capture students' interest.
3. Incorporate students' personal interests and experiences into the lesson content.

---

**Activate background knowledge**

1. Start with a task, such as a warmup or quiz, to review and activate prerequisite topics and skills.
  2. Provide a handout with relevant background knowledge, including worked examples, helpful formulas, and definitions.
  3. Use graphic organizers to help students structure and visualize information.
  4. Embed notes within the lesson materials to reinforce key concepts and provide guidance.
- 
-

---

## Curriculum Scaffolding Implementation Examples

### Strategy

---

**Build student fluency with targeted direct instruction and graduated levels of practice.**

1. Conduct fluency drills that focus on the specific skill being taught in the lesson.
2. Ask probing questions that vary in difficulty and form to deepen understanding.
3. Use quizzes and drills to reinforce prerequisite skills.
4. Implement digital practice activities that provide immediate feedback to students.
5. Provide direct instruction to address common misconceptions and errors.
6. Model key skills through direct instruction to ensure students understand the process.

---

**Support decoding of text, mathematical notation, and symbols to build math literacy.**

1. Create vocabulary lists, word walls, and use strategic word substitutions to support language comprehension.
  2. Design practice activities that target students' use of mathematical language in both large group and small group tasks.
- 
-

---

## Curriculum Scaffolding Implementation Examples

### Strategy

---

#### **Add formative assessment**

1. Use exit tickets and diagnostic assessment questions to gauge student understanding at the end of lessons.
2. Implement think, pair, share activities to encourage student discussion and reflection.
3. Conduct low-stakes quizzes to monitor ongoing progress without adding pressure.
4. Provide digital practice with immediate feedback to help students correct mistakes in real-time.
5. Include mid-lesson checks for understanding to adjust instruction as needed.
6. Encourage student self-assessment to promote self-reflection and goal setting.

---

#### **Provide opportunities for deeper conceptual understanding to optimize challenge.**

1. Offer more challenging practice that covers the same core skill to push student thinking.
  2. Provide practice with immediate feedback to help students self-correct.
  3. Design challenge practice activities that build new knowledge based on current skills.
- 
-

---

## Curriculum Scaffolding Implementation Examples

### Strategy

---

**Adapt curriculum to include appropriate goal-setting, planning, and strategy development and managing information to support executive function and SEL.**

1. Prompt students to plan their problem-solving approach before starting tasks.
2. Encourage students to break down tasks themselves to develop planning skills.
3. Break up complex tasks or lessons into smaller, manageable parts for students.
4. Use flowcharts to guide students through decision-making processes.
5. Provide knowledge organizers with key notes and information to aid in learning.

---

## Curriculum Scaffolding Implementation Examples

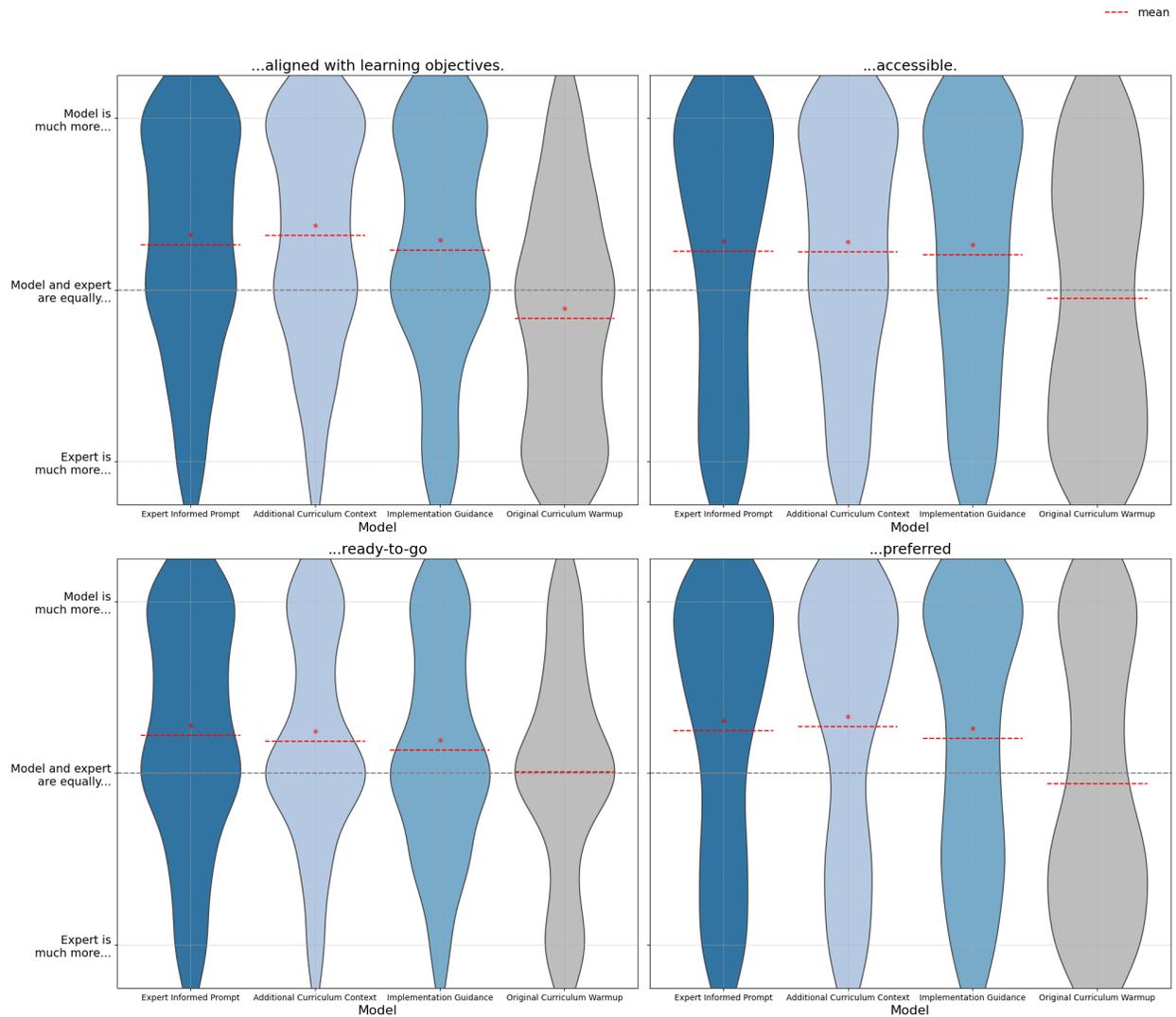
### Strategy

---

**Adapt curriculum to guide information processing and lighten the cognitive load on learners.**

1. Simplify numbers and data in problems to make them more accessible.
2. Use worked examples to demonstrate step-by-step solutions.
3. Include hints and suggestions in handouts, especially guidance for the first step in tasks.
4. Provide lesson notes in handouts to support learning.
5. Break down student-facing material into multiple artifacts, such as warm-ups, in-class practice, individual practice, and exit tickets.
6. Ensure the visual design of modified practice materials is similar to general practice to avoid confusion.
7. Reformat student-facing work to make it more intuitive and easier to follow.

# Appendix B Model Evaluation



**Figure B1:** Evaluators compared warmup tasks created under four model conditions (Expert-Informed Prompt, Additional Curriculum Context, Implementation Guidance, Original Curriculum Warmup) to those crafted by experts. The violin plots show the distribution of scores across four dimensions: alignment with learning objectives, accessibility for below-level students, readiness for classroom use, and overall preference. Red dashed lines indicate the mean scores.

# Appendix C Model Prompts

```
Expert-Informed Prompt

Your goal is to create a short warmup task that activates and refreshes relevant background knowledge for the lesson information provided.
-----
Lesson Title: {row['lesson_title']}
Lesson Narrative: {row['lesson_narrative']}
Learning Objectives: {row['learning_goals']}
-----
Class Context: 50% of students are performing below grade level and find it challenging to access the curriculum.
-----
Here are tips for creating a great warmup task:
- Think about the specific skills and knowledge students will need to access today's lesson
- If they're provided, you can draw on the original curriculum warmup and teacher response to original warmup
- The warmup should be focused on a single short scenario with up to 3 subparts
- Some subparts should be accessible for all students
- The warmup should take no more than 5 minutes to complete in class
- The task should be student-facing and not need anything else for students to complete it
- Do not provide a title. Do not provide anything apart from the task. Provide response in natural language, do not use latex
- The output should be short, concise and clear. Good warmups usually have 40-100 words
- Think step-by-step
-----
Now, create the warmup task.
WARMUP TASK:"""
```

**Figure C2:** This prompt was iteratively designed with input from expert teachers. It pipes in relevant information from the curriculum.

```
Expert-Informed Prompt + Additional Curriculum Context

Your goal is to create a short warmup task that activates and refreshes relevant background knowledge for the lesson information provided.
-----
Lesson Title: {row['lesson_title']}
Lesson Narrative: {row['lesson_narrative']}
Learning Objectives: {row['learning_goals']}
Original Curriculum Warmup: {row['warmup']}
-----
Class Context: 50% of students are performing below grade level and find it challenging to access the curriculum.
-----
Here are tips for creating a great warmup task:
- Think about the specific skills and knowledge students will need to access today's lesson
- If they're provided, you can draw on the original curriculum warmup and teacher response to original warmup
- The warmup should be focused on a single short scenario with up to 3 subparts
- Some subparts should be accessible for all students
- The warmup should take no more than 5 minutes to complete in class
- The task should be student-facing and not need anything else for students to complete it
- Do not provide a title. Do not provide anything apart from the task. Provide response in natural language, do not use latex
- The output should be short, concise and clear. Good warmups usually have 40-100 words
- Think step-by-step
-----
Now, create the warmup task.
WARMUP TASK:"""
```

**Figure C3:** This prompt included the original curriculum warmup from Illustrative Mathematics.

```
Expert-Informed Prompt + Additional Curriculum Context + Implementation Guidance

Your goal is to create a short warmup task that activates and refreshes relevant background knowledge for the lesson information provided.
-----
Lesson Title: {row['lesson_title']}
Lesson Narrative: {row['lesson_narrative']}
Learning Objectives: {row['learning_goals']}
Original Curriculum Warmup: {row['warmup']}
Teacher Response to Original Warmup: {row['teacher_feedback_to_warmup']}
-----
Class Context: 50% of students are performing below grade level and find it challenging to access the curriculum.
-----
Here are tips for creating a great warmup task:
- Think about the specific skills and knowledge students will need to access today's lesson
- If they're provided, you can draw on the original curriculum warmup and teacher response to original warmup
- The warmup should be focused on a single short scenario with up to 3 subparts
- Some subparts should be accessible for all students
- The warmup should take no more than 5 minutes to complete in class
- The task should be student-facing and not need anything else for students to complete it
- Do not provide a title. Do not provide anything apart from the task. Provide response in natural language, do not use latex
- The output should be short, concise and clear. Good warmups usually have 40-100 words
- Think step-by-step
-----
Now, create the warmup task.
WARMUP TASK:""
```

**Figure C4:** This prompt also includes the expert teacher’s strategy for modification of the original curriculum warmup.