

How to Teach Programming in the AI Era? Using LLMs as a Teachable Agent for Debugging

Qianou Ma^{1(⋈)}, Hua Shen², Kenneth Koedinger¹, and Sherry Tongshuang Wu¹

Carnegie Mellon University, Pittsburgh, PA, USA {qianoum,krk,sherryw}@cs.cmu.edu
University of Michigan, Ann Arbor, MI, USA huashen@umich.edu

Abstract. Large Language Models (LLMs) now excel at generative skills and can create content at impeccable speeds. However, they are imperfect and still make various mistakes. In a Computer Science education context, as these models are widely recognized as "AI pair programmers," it becomes increasingly important to train students on evaluating and debugging the LLM-generated code. In this work, we introduce HypoCompass, a novel system to facilitate deliberate practice on debugging, where human novices play the role of Teaching Assistants and help LLM-powered teachable agents debug code. We enable effective task delegation between students and LLMs in this learning-by-teaching environment: students focus on hypothesizing the cause of code errors, while adjacent skills like code completion are offloaded to LLM-agents. Our evaluations demonstrate that HypoCompass generates high-quality training materials (e.g., bugs and fixes), outperforming human counterparts fourfold in efficiency, and significantly improves student performance on debugging by 12% in the pre-to-post test.

Keywords: LLM · teachable agent · debugging · CS1

1 Introduction

LLMs are becoming an integral part of software development—commercialized tools like GitHub Copilot are now advertised as "your AI pair programmer" and generate up to 46% of users' code [6]. Despite their prevalence, LLMs often produce unpredictable mistakes [11], e.g., GPT-4 can still make mistakes 17% of the time in coding tasks for introductory and intermediate programming courses [22]. The impressive yet imperfect generative capabilities of LLMs, coupled with the associated risks of excessive reliance on these models, underscore the importance of teaching evaluation skills to students. In the context of programming, students must improve their debugging and testing skills [2].

However, debugging tends to be overlooked in formal educational curricula, especially in introductory Computer Science classes (i.e., CS1) [21]. Prior

[©] The Author(s), under exclusive license to Springer Nature Switzerland AG 2024 A. M. Olney et al. (Eds.): AIED 2024, LNAI 14829, pp. 265–279, 2024. $https://doi.org/10.1007/978-3-031-64302-6_19$

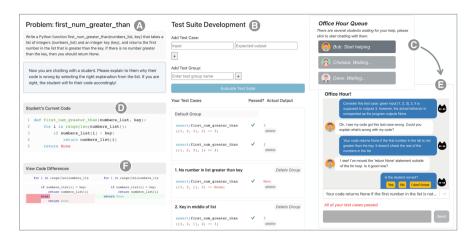


Fig. 1. In HYPOCOMPASS, given a programming problem description (A), a student user (in the role of a Teaching Assistant) needs to compile a test suite (B) and assist multiple LLM-simulated agents (e.g., Bob, Chelsea, Dave) in an Office Hour Queue (C) through a chat interface (E). Each LLM-agent acts as a novice seeking help with a buggy solution (D) and provides feedback to the user (F).

research has outlined various factors contributing to the absence of debugging instruction, such as instructors' limited time budget for developing specialized debugging materials and assessments [19]. Consequently, students primarily learn debugging from working on their own mistakes, which can be rather frustrating—they must invest substantial time and effort in *hypothesizing* the cause of bugs while grappling with other cognitively demanding tasks, such as understanding and writing code. These challenges prompt us to ask:

Research Question: Can we train students to improve debugging skills by providing *explicit* and *scaffolded* practice *with minimal cost to instructor time?*

In this work, we focus on training students' abilities in hypothesis construction, a critical step in debugging as established by prior work [29,30]. We introduce HYPOCOMPASS (Fig. 1, Sect. 3), an interactive, LLM-augmented intelligent tutoring system for debugging. Leveraging LLMs' material generation capability, we have these models imitate CS1 students who have written buggy code and require assistance from Teaching Assistants (TAs). Human novice students assume the role of the TA, who helps troubleshoot these bugs. This enables students to deliberately practice the skill of hypothesizing about the defects of LLM-generated code, delegating other tasks not core to hypothesis construction (e.g., code completion) to the LLM. As a result, HYPOCOMPASS fosters an engaging learning environment using the teachable agent framework [3] and provides students with guided exposure to LLM-generated bugs. We also employ prompting strategies such as focused task formation and over-generate-then-select to improve LLM generation quality in HYPOCOMPASS (Sect. 4).

We conducted two evaluation studies and found that HypoCompass saves instructors' time in material generation and is beneficial to student learning. In our LLM evaluation study (Sect. 5), expert inspections on six practice problems and 145 buggy programs showed that HypoCompass achieved a 90% success rate in generating and validating a complete set of materials, four times faster than human generation. Our learning evaluation study with 19 novices (Sect. 6) showed that HypoCompass significantly improved students' pre-to-post test performance by 12% and decreased their completion time by 14%.

In summary, we contribute:

- A pragmatic solution that balances the benefits and risks of LLMs in learning. We use LLMs to prepare students to engage with imperfect LLMs, and we highlight the importance of *role-playing* for practical LLM application and *task delegation* to help students focus on essential skills.
- A theoretically grounded instructional design to enhance debugging skills. To
 the best of our knowledge, we are the first to provide aligned instruction and
 assessments on the hypothesis construction learning objectives, *i.e.*, forming
 hypotheses about the source of error, a core bottleneck in debugging [25].

2 Related Works

The Debugging Process. Debugging is a complicated process of various cognitively demanding tasks, including understanding the code, finding bugs, and fixing bugs, with the first two considered primary bottlenecks [19,25]. While many studies have attempted to improve students' code understanding [12], there is limited instruction on bug finding. Researchers characterize the cognitive model of bug finding as a hypothesis construction process, including initializing, modifying, selecting, and verifying hypotheses (Fig. 2B) [29]. This process is challenging: prior works show that novices struggle to systematically generate comprehensive hypotheses and identify the right hypothesis, in contrast to experts [7,8]. Hence, we emphasize teaching students to construct accurate hypotheses about bugs and develop comprehensive hypotheses about potential bugs.

Tutors and Tools for Debugging Training. Prior studies [19] and online discussions [21] indicate that teaching debugging is challenging and is rarely included in CS1 curricula, due to logistical challenges like the lack of instructional time and resources [5,10]. Existing tools demand instructor effort and often focus on the full debugging process, improving bug-fixing accuracy and efficiency [1,15]. In contrast, few studies emphasize accurate or comprehensive hypothesis construction (and they tend to be language-specific) [13,25]. To fill in the gap, we design HypoCompass to provide deliberate practice [9] on hypothesis construction, and use the LLM generation capability to provide easily adaptable and targeted exercises with immediate feedback.

LLM Capabilities for CS Learning. LLMs can perform well in a CS1 class-room [22], but concerns about misuse and LLM errors limit their use in education [2]. Therefore, current deployments tend to focus on generating instructional

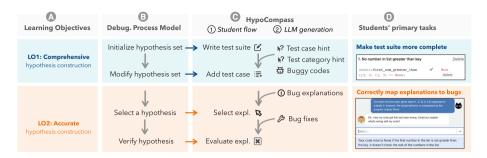


Fig. 2. To enable deliberate practice, we establish a close mapping between the (A) learning objectives, (B) the cognitive debugging process model, (C) the HYPOCOMPASS interaction flow, and (D) the primary tasks students perform in HYPOCOMPASS. We offload various material generation tasks to LLMs (C_2) .

materials (e.g., questions [24]). In our work, HYPOCOMPASS uses the LLM to generate inter-dependent materials in an integrated process and frame the LLM as a student asking for help [3], such that human novices can embrace imperfections in LLMs. Two unique capabilities of LLMs power this: (1) LLMs can simulate different personas and tutoring interactions [18]; (2) LLMs make common mistakes and natural bugs similar to humans [20], which can be used as buggy code practice examples. We adapt and develop various prompting methods [27] to enhance the quality of LLM generations.

3 The Design of HypoCompass

Grounded in the cognitive process [29] and the novice-expert difference in hypothesis-driven debugging (Sect. 2), we specify two crucial learning components for HypoCompass: comprehensive and accurate hypothesis construction. Prior work shows that hypothesis construction is closely connected with testing [30]: each additional test case should, ideally, be a hypothesis about what can go wrong in the program. In turn, a *comprehensive* test suite (*i.e.*, a set of test cases) should allow an effective debugger to construct a *accurate* hypothesis about why the program is wrong. We thus design toward two learning objectives (Fig. 2A,D):

- LO1. Comprehensive Hypothesis Construction: Construct a comprehensive test suite that well covers the possible errors for the given problem.
- LO2. Accurate Hypothesis Construction: Given the failed test cases, construct an accurate explanation of how the program is wrong.

Interface and Key Components. We designed HYPOCOMPASS through an iterative development process with 10 pilots, including CS1 students, TAs, and instructors. In the resulting interface (Fig. 1), a human student would be asked to play the role of a TA where they help an LLM-simulated student (LLM-agent)

Starter test category hint	0	Explanations to choose from			
1. No number in list greater than key	Delete Group	Ok, I see my code got this test case wrong. Could you explain what's wrong with my code?			
2. Key in middle of list	Delete Group	Your code returns None if the first number in the list is not			
3. All numbers in list greater than key	Delete Group	Your function returns None as soon as it encounters a number that is smaller than the key, it only continue to check the loop if num is equal to key. But there might be numbers greater than the key later after a number smaller than key.			
Descriptive test case hint Hi Christinal Here is my code and I think problem with it. Can you walk me through		Your code returns None if the first number in the list is not greater than the key. It doesn't check the rest of the numbers in the list Correct explanation			
		Your function is printing the number instead of returning it, so it always returns None			
Write a test case to cover the scenario where seen in the list, but there are numbers (a)		Distracting explanations			
oresent in the list, but there are numbers (a)	greater than the key	Distracting explanations			
cresent in the list, but there are numbers (a)		Feedback to correct explanation selection			
coresent in the list, but there are numbers (a) Feedback to incorrect test cases add Test Case:	greater than the key	Distracting explanations			
resent in the list, but there are numbers (a) Feedback to incorrect test cases and Test Case: (3, 2, 1), 3	greater than the key	Feedback to correct explanation selection [2] [3] [4] [5] [6] [6] [7] [6] [7] [8] [8] [9] [9] [9] [9] [9] [9			
Feedback to incorrect test cases add Test Case: [3, 2, 1], 3 NONE The correct expected output is None, please to	greater than the key 1 1 ry again.	Feedback to correct explanation selection I seel I've moved the 'return None' statement outside of the for loop. Is it good now?			
Feedback to incorrect test cases and Test Case: (a, 2, 1], 3 NONE The correct expected output is None, please to the correct expected output is None.	try again. action agreeater than the key try again.	Feedback to correct explanation selection I seel I've moved the 'return None' statement outside of the for loop. Is it good now? Your code returns None if the first number in the list is not			
Feedback to incorrect test cases and Test Case: (3, 2, 1), 3 NONE The correct expected output is None, please to the	ty again. try again. try again. action agreater_than([1, ode would output 3.	Feedback to correct explanation selection I seel I've moved the 'return None' statement outside of the for loop. Is it good now? Your code returns None if the first number in the list is not All of your test cases passed. View Code Differences for i in range(len(numbers_lis for i in range(len(numbers_t));			
Feedback to incorrect test cases Add Test Case: [3, 2, 1], 3 NONE The correct expected output is None, please to But if that's the case, shouldn't first_nut	ty again. try again. try again. action agreater_than([1, ode would output 3.	Feedback to correct explanation selection I seel I've moved the 'return None' statement outside of the for loop. Is it good now? Your code returns None if the first number in the list is not All of your test cases passed. View Code Differences for i in range(len(numbers_lis for i in range(len(numbers_			

Fig. 3. (a) HYPOCOMPASS offers (1) test category hints to help write a comprehensive test suite systematically; (2) test case hints to help students add missing test scenarios; (3) candidate explanation pool to clarify misconceptions of alternative explanations. (b) HYPOCOMPASS provides immediate feedback to (1) incorrect test cases, ensuring students understand the code behavior; (2) correct explanations, as correct code fixes; (3) incorrect explanations, as confusion messages from the LLM-agent.

in debugging. They need to write and sort test cases into categories (Fig. 1B) that represent different hypotheses of what inputs may trigger errors in code.

Once the student is satisfied with their test suite, HYPOCOMPASS shows them an Office Hour Queue (OHQ) simulator (Fig. 1C). As the student interacts with each LLM-agent, the agent presents a buggy code snippet (Fig. 1D). The student guides the LLM-agent in debugging code through a dialog interface (Fig. 1E), selecting or creating test cases that reflect their hypotheses of the bug, and selecting explanations for the bug among a pool of candidate natural language explanations. These candidates each explain a different bug, representing alternative hypotheses that may confuse students (e.g., Fig. 3a₃).

The LLM-agent then uses the test case and explanation to revise the code, providing immediate feedback to the student (Fig. 3b). If the explanation is correct, the agent will conduct minimal code fixes, and present the color-coded edits

as feedback (Fig. 1F, a zoomed-in view is in Fig. 3b₂). Otherwise, the LLM-agent will ask the student to reflect on their hypothesis by responding with a confusion message that highlights the discrepancy between the student's explanation and the actual code behavior (Fig. 3b₃).

Once the student correctly confirms that all the bugs are fixed, they can move to help the next LLM-agent (Fig. 1C). Upon completion, HYPOCOMPASS will provide the next round of exercises with another programming problem. While the numbers are configurable, by default HYPOCOMPASS includes two programming exercises, each with three LLM-agents (buggy programs).

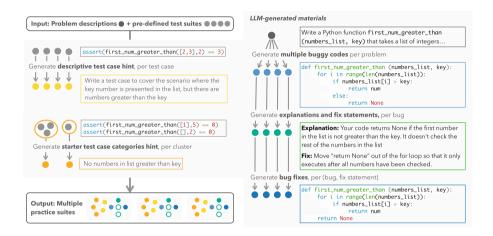


Fig. 4. Examples of inputs and outputs to the LLM material generation pipeline.

We highlight the two most essential components of the interaction:

- Frame imperfect LLMs through role-play. We use the LLM to simulate students who wrote bugs and have human novices offer help. This teachable agent setup supports learning, helping students reflect on their knowledge and reason through diverse bugs [23]. Having students work through "other people's errors" also boosts their motivation and protects their self-efficacy [3]. More importantly, it actively involves novices in identifying bugs in LLM-generated code, enabling guided exposure to LLM imperfectness.
- Task delegation between students and LLMs. To ensure deliberate practice on comprehensive and accurate hypothesis construction, students primarily engage in two tasks corresponding to each learning objective (Fig. 2D): (1) making the test suite more complete (LO1); and (2) correctly mapping explanations to bugs (LO2). We align student interaction flow (Fig. 2C₁) with the cognitive model of debugging [29] (Fig. 2B). LLMs take over other tasks that are *indirectly related* to the core learning goals, including generating diverse bugs and fixes, which frees students from code writing. We also use

LLMs to support scaffolding, generate hints (Fig. 3a), and provide immediate feedback throughout the practice (Fig. 3b).

4 LLM Integration

As shown in Fig. 2C₂, we use LLM to generate five types of materials: (1) test case category hints, (2) test case hints, (3) buggy programs, (4) explanations of bugs, and (5) programs with bugs fixed. We reduce instructor workload by generating practices using just a problem description, a reference solution, and a reference test suite with about 10 inputs, and we further minimize human verification overhead with optimized prompts and automated algorithms. Our generation process is detailed in Fig. 4, example prompts are in Table 1, and full prompts are in Table 3 in Supplements¹. OpenAI's gpt-3.5-turbo is used for all materials, except for explanation generation, which uses gpt-4 for enhanced reasoning capabilities. Below are key factors to the success of generation:

Task Formation and Decomposition. We iterate on our prompts according to the nature of the task. First, as LLMs behave inconsistently when the user tasks conflict with LLMs inherent training objectives [28], we carefully formulate the task to avoid introducing competing tasks. Take Local Bug Fix (Table 1) as an example: when we directly ask the LLM to fix a bug according to an explanation, we observe that the model almost always over-fix all bugs irrespective of the provided instructions. This is because LLMs can be biased towards generating fully correct code (part of the LLM pre-training) and away from local bug fixing (changing only the buggy snippet described by the instruction, the desired task). Hence, we re-frame it as a translation task, converting bug-fixing instructions to its code format old \rightarrow new code snippet. This task re-framing mitigates the model's inherent bias, reducing over-fixing errors by 70%.

Second, for multi-step tasks (e.g., Local Bug Fix), we adopt LLM-chains [27], decomposing tasks into sub-tasks handled by separate steps, such that each step contributes to stable performance. Third, we also address prompt complexity by explicitly prioritizing essential requirements. For tasks like generating Bug Explanations and Fix Instructions (Table 1), we prioritize precise bug extraction, instructing the model to list all unique bugs upfront. Secondary requirements (e.g., word limits) are specified only in the output format. This hierarchical disentanglement significantly improves success rates by over 40%.

Over-Generate-then-Select. While LLMs can easily generate random materials, it is nontrivial to ensure that their generations have pedagogical values. For example, behaviorally distinct bugs help students practice with varied instances, but it is hard to enforce through prompting as it requires LLMs to "know" bug behaviors. Nonetheless, we can configure the non-deterministic LLMs to overgenerate multiple solutions with mixed qualities [17], and then select a subset of desired ones (Fig. 5). We apply this strategy in multiple places:

¹ Supplemental materials are at: http://tinyurl.com/hypocompass-sup.

Table 1. Prompts and temperatures (Temp.) for generating bugs, explanations, and fixes. The temperature is set higher for more diverse and random outputs.

	Generation goal	Temp.
Buggy code	To over-generate bugs with mixed quality for further selection.	0.7
	are a novice student in intro CS , you make mistakes and wrigy code.	te
[User] Wri	<pre>blem Description: {problem.description} te different buggy solutions with common mistakes like novic dents:</pre>	e
	To describe each unique bug, and write a corresponding fix instruction. If there are multiple bugs in the code, generate their explanations and fixes separately.	
[Sys.] You	are a helpful and experienced ${f TA}$ of an introductory programss.	ming
[User] prol	I'm a student in your class. I'm having trouble with this blem in the programming assignment: {problem_description} Her	
unio in s wha fix	buggy code: {buggy_code} What's wrong with my code? List all que bugs included, but do not make up bugs. For each point, the format of: {explanation: accurate and concise explanation the code does and what the bug is, for a novice, fix: how the bug, within 30 words} yreturn the bullet list. Do not write any other text or code.	put n of to
unio in wha fix Only	que bugs included, but do not make up bugs. For each point, the format of: {explanation: accurate and concise explanation the code does and what the bug is, for a novice, fix: how the bug, within 30 words}	put n of to e.
unio in what fix Only	que bugs included, but do not make up bugs. For each point, the format of: {explanation: accurate and concise explanation the code does and what the bug is, for a novice, fix: how the bug, within 30 words} y return the bullet list. Do not write any other text or cod To edit the buggy code according to the fix instruction, w/o over- or	put n of to e.
unid in some what fix only Bug fix [Sys.] You [User] Origination for the fore edited and the fore edited	que bugs included, but do not make up bugs. For each point, the format of: {explanation: accurate and concise explanation the code does and what the bug is, for a novice, fix: how the bug, within 30 words} y return the bullet list. Do not write any other text or cod To edit the buggy code according to the fix instruction, w/o over-or under-fix.	put n of to e. 0.3
unid in what fix Only Bug fix [Sys.] You [User] Orig for for edi -> o	que bugs included, but do not make up bugs. For each point, the format of: {explanation: accurate and concise explanation: the code does and what the bug is, for a novice, fix: how the bug, within 30 words} yreturn the bullet list. Do not write any other text or cod To edit the buggy code according to the fix instruction, w/o over-or under-fix. fix bugs in Python code closely following the instructions. ginal code: {buggy_code}; Code modification: {explanation} nslate the statement into actual, minimal code change in thimat: iginal code snippet: ""copy the lines of code that need ting""	put n of to e. 0.3

(1) To expose students to behaviorally distinct bugs, we over-generate buggy code (Table 1). We filter out correct code, and we vectorize buggy code's behavior based on the reference test suite (Fig. 5A, 0 being failed tests). We then greedily choose a diverse subset of buggy programs with the maximum pairwise distance, using Euclidean distance on the error vectors (Fig. 5B).



Fig. 5. Over-generate and automatically select materials with pedagogical values.

- (2) To help students clarify misconceptions (Fig. 5C), we want distracting explanations that look similar to the actual explanation for each practice buggy code. We choose from the over-generated buggy code pool, find two with the smallest Euclidean distance to the target code, and use their corresponding explanations as distractors. The mapping also helps generate the confusion messages (Fig. 3b₃)—when a student selects the distractor explanation, we use its corresponding buggy code to find test cases to present to students.
- (3) To capture key testing aspects in our test category hints (Fig. 5D), we cluster reference test cases into semantically meaningful groups. We build dendrograms from test case vectors with Agglomerative Hierarchical Clustering [14], which guide the selection of test category hints from the overgenerated pool.

Human-in-the-Loop Verification. As shown in Fig. 4, while the hints for test cases and categories are generated separately, the materials relevant to bugs are generated in sequential order. We perform human verification per step to mitigate the risk of cascading errors in subsequent steps. We provide more details on human verification and editing times in Sect. 5.

5 LLM Evaluation: Generation Efficiency and Quality

We evaluated the generations on six different problems from prior work [4] and our own problems (detailed in Table 4 in Supplements). On average, for each problem, we generated 3 test category hints, 10 test case hints, 24 buggy programs, explanation and fix instructions, and 33 bug fixes. The total number and the success rates are summarized in Table 2. We provided the success criteria for all types of materials in Table 5 in Supplements.

Material	Raw LLM outputs			Human verification		
	# Generation	Avg. gen time	Success%	Avg. edit time	IRR%	κ
Test case description hint	61	0:00:37	98.36%	0:00:08	100%	-
Test case category hint	18	0:00:10	94.44%	0:00:10	100%	-
Buggy code	145	0:01:30	57.93%	0:00:02	n/a	n/a
Bug explanation and fix	145	0:03:36	91.72%	0:00:52	90%	0.875
Bug fix	195	0:02:45	86.15%	0:00:37	92%	0.752

Table 2. LLM Evaluation: Time, Success rate, and Inter-Rater Reliability scores (*i.e.*, IRR% = #agreements / #total labels, κ is Cohen's Kappa coefficient).⁴

Method. Two authors annotated 10% of the generations at each step individually, and discussed to resolve the disagreement and update the codebook. An external instructor annotated the same 10% of LLM-generated materials, using the updated codebook. We calculated the inter-rater reliability (IRR) between the external instructor and the resolved annotation among the two authors using percent IRR and Cohen's Kappa. As shown in Table 2, the agreements are satisfactory across different model generations (IRR% > 90% and $\kappa > 0.75$)². One author annotated the rest of the materials to calculate the success rates. We log the verification and editing *time*, as proxies to the instructor overhead.

To compare LLM and human generations, we recruited two experienced CS TAs to each create practice materials for a specific problem. Each TA received the same input as LLMs, was asked to produce one set of materials matching the amount of content LLMs produced, and was compensated for their time.

Result: Efficient and High-Quality Generation. We achieve high-quality generation: a complete set of practice materials with 9 buggy programs (3 for practice and 6 more as distractors), 9 bug explanations, 9 bug fixes, 10 test case hints, and 3 test category hints can be generated with a 90% success rate and only takes 15 min to label and edit. As we *over-generate* and automatically select buggy code, a success rate over 50% is reasonable for practical use.

Employing LLMs can also be significantly more efficient. In total, a TA spent around 60 min to generate one set of practice materials for HypoCompass. One TA noted the difficulty in consistently creating unique and high-quality materials after 30 min, saying that "the importance of the bug I create would start to decline." The same author evaluated the TAs' generations using the annotation codebook, which had a 100% success rate and took 11 min. The time invested in generating and editing instructional materials for HypoCompass using LLMs was 4.67 times less than that of the human TAs.

² Buggy programs undergo automatic testing, so human verification is unnecessary (n/a). If both raters unanimously agree in one category, kappa is undefined (-), so κ is only noted when there's less than 100% IRR agreement on a single label.

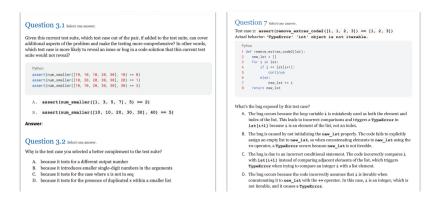


Fig. 6. Pre-post test question examples for LO1 comprehensive (Q3.1 and Q3.2) and LO2 accurate hypothesis construction (Q7).

6 Learning Evaluation: Pre- / Post-Test Study

Can novices better formulate hypotheses after engaging with HYPOCOMPASS? We conducted a learning evaluation with 19 students and compared the difference in speed and performance from the pre-test to the post-test.

Assessment. To best capture student learning gains on our learning objectives, we took a backward design method [26] to create an aligned assessment for the comprehensive LO1 and accurate LO2 hypothesis construction skills. We conducted multiple rounds of pilots to refine our intervention and pre-post tests. Our final tests are based on two programming exercises with comparable difficulties. We counterbalanced pre-post tests' problems to control for problem sequence influence. Each test consists of seven questions, with three assessing LO1 and four for LO2. Figure 6 provides a sample for each. For instance, Question 3.1 asks students to identify the more suitable test case to add to an existing test suite, evaluating their ability to construct comprehensive hypotheses (LO1). We measure students' performance using their test scores based on a standard rubric. We also log the pre-post tests' completion time as a proxy for proficiency.

Method: Study Procedure and Participants. Our hour-long user study constituted a pre-survey, pre-test, interaction with HYPOCOMPASS, post-test, and a post-survey. Participants began with a pre-survey, which asked demographic information and 7-level Likert Scale questions on their debugging experiences. Then, participants had up to 20 min for the pre-test. The system interaction consisted of two problems, where participants needed to write a test suite and explain bugs in three different buggy programs for each problem. The first problem was the same as in the pre-test, and the second problem matched the screening survey's exercise. By reusing problems that students have seen, we isolate our learning objectives from the program comprehension skills. After a subsequent 20-min post-test, participants filled out a post-survey with Lik-

ert Scale and open-ended questions on their experience and perceptions using HypoCompass. Participants received a \$15 Gift Card for their time.

We recruited a diverse group of undergraduate and graduate students from four public or private US institutions. Interested participants completed a screening survey, which included a programming exercise that also served as the second exercise in our study. To ensure a suitable skill range, we excluded those with extensive programming experience or who quickly solved the exercise. After filtering, 19 participants (S1-19) were included in the study—12 females, 6 males, 1 non-binary, and 8 non-native English speakers, with an average age of 20.7.

Quantitative Result: Learning Gains. A two-tailed paired t-test showed that students' pre-test to post-test scores significantly improved by 11.7% (p=0.033<0.05), and the time of completion significantly reduced by 13.6% (p=0.003), indicating success in learning through HypoCompass interaction. Note that the bugs used in pre-post tests are generated by humans and are not the same as in HypoCompass. As such, the significant learning gains indicate that students could learn debugging skills transferable to real-world bugs.

Where does the learning gain come from? We break down the analyses by learning objectives. We found a small 6.1% improvement in the score and a large 23.6% time reduction for *comprehensive hypothesis construction* (LO1), and a large 15.8% improvement in the score and a small 9.0% time reduction for *accurate hypothesis construction* (LO2). Therefore, students showed more efficiency enhancement in LO1, and more learning gains in LO2. Note that these improvements may confound with problem difficulty, as the items corresponding to LO1 (pre-test $\mu = 54\%$) seem easier than the ones for LO2 (pre-test $\mu = 38\%$).

Qualitative Result: Student Perceptions. We further unpack how HYPOCOMPASS contributed to learning by analyzing the survey responses. Students valued being able to offload some debugging subtasks to HYPOCOMPASS, such as writing code and explanations. For example, S1 said "looking at the test behavior and the explanation options really helps relieve that burden." Students also generally felt that the LLM-generated bugs and fixes were authentic. Most participants could not tell if their practiced programs were written by students or AI because of their experiences making or seeing similar mistakes from peers.

Moreover, students reported that HYPOCOMPASS was engaging, fun, not frustrating, and helped build confidence in debugging. A Wilcoxon signed-rank test shows a significant increase in self-rated confidence in debugging by 15% (p=0.007). Students rated HYPOCOMPASS as significantly more engaging (6.0 out of 7), fun (6.0), and less frustrating (2.5) than their conventional way of learning debugging and testing (p < 0.005) for each). S8 especially liked the teachable agent setup: "the role play just feels more natural because it feels like explaining to a rubber duck instead of to talking to myself".

7 Discussion

Teachable Agent for Appropriate Reliance with Imperfect AIs. Our work illustrates a scenario in which LLM-generated bugs are not seen as problems

but rather as features. HYPOCOMPASS's teachable agent setup provides students with moderated exposure to imperfect LLMs, and may help them learn that LLMs are fallible and calibrate trust accordingly. Future iterations could remove material validation and allow direct exposure to unfiltered LLM mistakes in real-time interactions, taking full advantage of the teachable agent framework. Students will naturally expect that the LLM-agent seeking help may make mistakes (e.g., fail to follow bug-fixing explanations). This approach, however, requires a more sophisticated design for scaffolding students in recognizing LLM errors.

Task Delegation for Shifting Learning Focus. Our exploration lays the foundation for a paradigm shift toward cultivating higher-order evaluation skills in the generative AI era. Essentially, we asked: what skills should we offload, and what should we learn? Most students in our study appreciated offloading subtasks to LLM (Sect. 6); however, some need more scaffolds, while others prefer less. Future research can investigate more personalized task delegation. For example, students who need more help can use LLMs to facilitate code tracing, and students can also write their own explanations for bugs based on their proficiency. Deciding the bare minimum programming skills and human-AI collaboration skills to teach also warrants further exploration [16].

Modularize to Adapt to Different Needs. Though most students and instructors found HypoCompass engaging, some expressed concerns about the deployment and maintenance cost of a new tool. To maximize utility to diverse users, we can modularize different components in HypoCompass. Instructors who prefer to distribute training materials as handouts can rely entirely on the material generation module. In contrast, instructors who want to experiment with TA training can employ HypoCompass with practice generated using their training questions. Future studies may perform ablation studies to evaluate different HypoCompass components with more extensive classroom deployment.

Limitation. We primarily evaluated whether HYPOCOMPASS can bring learning and efficiency gains through small in-lab experiments. With this prerequisite, we plan to conduct future classroom deployment with controlled comparisons. There is also a limitation regarding the reported efficiency of the LLM-assisted instructional material development, as the instructors need some familiarization time with the tool and the process.

8 Conclusion

In an attempt to answer how LLMs can reshape programming education's focus, we introduce a novel system, HYPOCOMPASS, and new instructional designs for hypothesis construction skills. We aim to provide engaging and deliberate practice on debugging to novices, using our theoretically motivated and empirically tested teachable agent augmented by LLM. Our evaluations show that HYPOCOMPASS can efficiently help instructors create high-quality instructional materials, effectively train novices on comprehensive and accurate hypothesis construction, and facilitate students' confidence and engagement in debugging.

Acknowledgment. Thanks to the participants, reviewers, Kelly Rivers, Michael Taylor, Michael Hilton, Michael Xieyang Liu, Vicky Zhou, Kexin Yang, Jionghao Lin, Erik Harpstead, and other Ken's lab members for the insights and help. This study was supported by gift funds from Adobe, Oracle, and Google.

References

- Ardimento, P., Bernardi, M.L., Cimitile, M., Ruvo, G.D.: Reusing bugged source code to support novice programmers in debugging tasks. ACM Trans. Comput. Educ. 20(1), 1–24 (2019). https://doi.org/10.1145/3355616
- Becker, B.A., Denny, P., Finnie-Ansley, J., Luxton-Reilly, A., Prather, J., Santos, E.A.: Programming is hard-or at least it used to be: educational opportunities and challenges of AI code generation. In: Proceedings of the 54th ACM Technical Symposium on Computer Science Education V, vol. 1. pp. 500–506 (2023)
- 3. Blair, K., Schwartz, D.L., Biswas, G., Leelawong, K.: Pedagogical agents for learning by teaching: teachable agents. Educ. Technol. Res. Dev. 47(1), 56–61 (2007)
- Dakhel, A.M., Majdinasab, V., Nikanjam, A., Khomh, F., Desmarais, M.C., Jiang, Z.M.J.: Github copilot AI pair programmer: asset or liability? J. Syst. Softw. 203, 111734 (2023). https://doi.org/10.48550/ARXIV.2206.15331
- Desai, C., Janzen, D.S., Clements, J.: Implications of integrating test-driven development into CS1/CS2 curricula. SIGCSE Bull. 41(1), 148–152 (2009)
- Dohmke, T.: GitHub copilot x: the AI-powered developer experience (2023). https://github.blog/2023-03-22-github-copilot-x-the-ai-powered-developer-experience/. Accessed 5 Sept 2023
- Edwards, S.H., Shams, Z.: Comparing test quality measures for assessing studentwritten tests. In: Companion Proceedings of the 36th International Conference on Software Engineering. ICSE Companion 2014, pp. 354–363. Association for Computing Machinery, New York, NY, USA (2014)
- 8. Edwards, S.H., Shams, Z.: Do student programmers all tend to write the same software tests? In: Proceedings of the 2014 Conference on Innovation and Technology in Computer Science Education. ITiCSE '14, pp. 171–176. Association for Computing Machinery, New York, NY, USA (2014)
- 9. Ericsson, A., Pool, R.: Peak: secrets from the new science of expertise. Random House (2016)
- Fitzgerald, S., McCauley, R., Hanks, B., Murphy, L., Simon, B., Zander, C.: Debugging from the student perspective. IEEE Trans. Educ. 53(3), 390–396 (2010)
- Ganguli, D., et al.: Predictability and surprise in large generative models. In: Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency, pp. 1747–1764 (2022)
- Kallia, M.: The search for meaning: Inferential strategic reading comprehension in programming. In: Proceedings of the 2023 ACM Conference on International Computing Education Research (ICER '23). ACM (2023)
- Ko, A.J., Myers, B.A.: Debugging reinvented: asking and answering why and why not questions about program behavior. In: Proceedings of the 30th International Conference on Software Engineering. ICSE '08, pp. 301–310. Association for Computing Machinery, New York, NY, USA (2008)
- Lukasová, A.: Hierarchical agglomerative clustering procedure. Pattern Recogn. 11(5-6), 365-381 (1979)

- Luxton-Reilly, A., McMillan, E., Stevenson, E., Tempero, E., Denny, P.: Ladebug: an online tool to help novice programmers improve their debugging skills. In: Proceedings of the 23rd Annual ACM Conference on Innovation and Technology in Computer Science Education. ITiCSE 2018, pp. 159–164. Association for Computing Machinery (2018)
- 16. Ma, Q., Wu, T., Koedinger, K.: Is AI the better programming partner? Humanhuman pair programming vs. human-AI pAIr programming. arXiv preprint arXiv:2306.05153 (2023). http://arxiv.org/abs/2306.05153
- 17. MacNeil, S., Tran, A., Mogil, D., Bernstein, S., Ross, E., Huang, Z.: Generating diverse code explanations using the GPT-3 large language model. In: Proceedings of the 2022 ACM Conference on International Computing Education Research-Volume 2, pp. 37–39 (2022)
- Markel, J.M., Opferman, S.G., Landay, J.A., Piech, C.: GPTeach: interactive TA training with GPT-based students. In: Proceedings of the Tenth ACM Conference on Learning @ Scale. L@S '23, pp. 226–236. Association for Computing Machinery, New York, NY, USA (2023). https://doi.org/10.1145/3573051.3593393
- McCauley, R., et al.: Debugging: a review of the literature from an educational perspective. Comput. Sci. Educ. 18(2), 67–92 (2008)
- Mozannar, H., Bansal, G., Fourney, A., Horvitz, E.: Reading between the lines: modeling user behavior and costs in AI-assisted programming. arXiv preprint arXiv:2210.14306 (2022)
- 21. News, Y.H.: Why don't schools teach debugging? (2014). https://news.ycombinator.com/item?id=7215870. Accessed 8 Sept 2023
- Savelka, J., Agarwal, A., An, M., Bogart, C., Sakr, M.: Thrilled by your progress!
 Large language models (GPT-4) no longer struggle to pass assessments in higher education programming courses. arXiv preprint arXiv:2306.10073 (2023)
- 23. Shahriar, T., Matsuda, N.: What and how you explain matters: inquisitive teachable agent scaffolds knowledge-building for tutor learning. In: Wang, N., Rebolledo-Mendez, G., Matsuda, N., Santos, O.C., Dimitrova, V. (eds.) AIED 2023. LNCS, vol. 13916, pp. 126–138. Springer, Cham (2023). https://doi.org/10.1007/978-3-031-36272-9_11
- 24. Wang, Z., Valdez, J., Basu Mallick, D., Baraniuk, R.G.: Towards human-like educational question generation with large language models. In: Rodrigo, M.M., Matsuda, N., Cristea, A.I., Dimitrova, V. (eds.) AIED 2022. LNCS, vol. 13355, pp. 153–166. Springer, Cham (2022). https://doi.org/10.1007/978-3-031-11644-5_13
- 25. Whalley, J., Settle, A., Luxton-Reilly, A.: Analysis of a process for introductory debugging. In: Proceedings of the 23rd Australasian Computing Education Conference. ACE '21, pp. 11–20. Association for Computing Machinery (2021)
- 26. Wiggins, G.P., McTighe, J.: Understanding by Design. ASCD (2005)
- 27. Wu, T., Terry, M., Cai, C.J.: AI chains: transparent and controllable human-AI interaction by chaining large language model prompts. In: Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems. No. Article 385 in CHI '22, pp. 1–22. Association for Computing Machinery, New York, NY, USA (2022)
- 28. Xie, J., Zhang, K., Chen, J., Lou, R., Su, Y.: Adaptive chameleon or stubborn sloth: unraveling the behavior of large language models in knowledge clashes (2023)
- Xu, S., Rajlich, V.: Cognitive process during program debugging. In: Proceedings of the Third IEEE International Conference on Cognitive Informatics, pp. 176–182. IEEE (2004)
- Zeller, A.: Why Programs Fail: A Guide to Systematic Debugging, 2nd edn. Morgan Kaufmann, Oxford (2009)