

EMBEDAGENT: Benchmarking Large Language Models in Embedded System Development

Ruiyang Xu^{1,2,*}, Jialun Cao^{3,*}, Mingyuan Wu³, Wenliang Zhong^{1,2}, Yaojie Lu¹,

Ben He^{1,2}, Xianpei Han¹, Shing-Chi Cheung³, Le Sun¹

¹Chinese Information Processing Laboratory, Institute of Software, Chinese Academy of Sciences, Beijing, China

²University of Chinese Academy of Sciences, Beijing, China

³The Hong Kong University of Science and Technology, Hong Kong, China

{xuruiyang2022,zhongwenliang2024,luyaojie,xianpei,sunle}@iscas.ac.cn,
{jcaoop, scc}@cse.ust.hk, 11849319mail.sustech.edu.cn, benhe@ucas.edu.cn

Abstract

Large Language Models (LLMs) have shown promise in various tasks, yet few benchmarks assess their capabilities in embedded system development. In this paper, we introduce EMBEDAGENT, a paradigm designed to simulate real-world roles in embedded system development, such as *Embedded System Programmer*, *Architect*, and *Integrator*. This paradigm enables LLMs to be tested in tasks that bridge the gap between digital and physical systems, allowing for a more comprehensive assessment of their capabilities. To evaluate LLMs on these tasks, we propose EMBEDBENCH, the first comprehensive benchmark for embedded system programming, circuit design, and cross-platform migration. EMBEDBENCH consists of 126 cases, covering 9 electronic components across 3 hardware platforms. Through extensive experiments on 10 mainstream LLMs, we uncover several key findings. Surprisingly, despite the simplicity of the cases, DeepSeek-R1 achieves only a 55.6% pass@1 rate when provided with schematic information, and 50.0% when tasked with generating the schematics itself. In the cross-platform migration tasks, LLMs show relatively strong performance with MicroPython on the Raspberry Pi Pico (with the top model achieving 73.8% pass@1), but perform poorly on ESP-IDF, where the best model reaches only 29.4% pass@1. Interestingly, we observe that general-purpose chat LLMs like DeepSeek-V3 often fail to utilize relevant pre-trained knowledge in this domain, while reasoning LLMs tend to overthink and overlook efficient knowledge during pretraining. Based on these insights, we propose two strategies—retrieval augmented generation and compiler feedback—to enhance LLM performance. These strategies result in significant improvements, with Deepseek-R1 reaching a 65.1% pass@1 with correct schematics, and 53.1% without. Additionally, the accuracy of the Arduino to ESP32 migration task improves from 21.4% to 27.8%.

CCS Concepts: • Computer systems organization → Embedded systems.

Keywords: Embedded System, Large Language Models

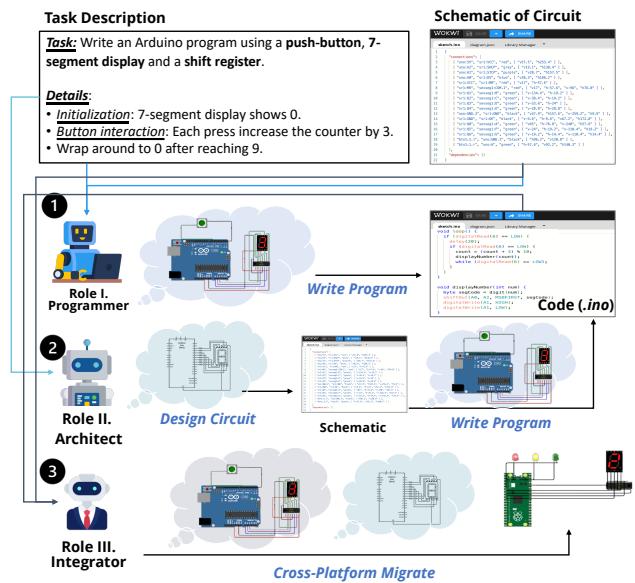


Figure 1. Three Settings of EMBEDAGENT: **❶ Embedded System Programmer**: Given the task description and schematic of circuit, LLMs are expected to write Arduino code. **❷ Embedded System Architect**: Given the task description, LLMs are expected to design the circuit and write the code. **❸ Embedded System Integrator**: Given the schematic of circuit and code of one hardware, LLMs are expected to migrate the design of circuit and code to another hardware platform.

1 Introduction

Large Language Models (LLMs) have demonstrated remarkable expertise across various software engineering tasks, such as code generation [13, 22], defect detection [40], program repair [17, 37, 43] and code translation [38]. To assess the these abilities of LLMs, benchmarks like HumanEval [8] and SWE-bench [18] have been introduced, offering valuable insights into the strengths and limitations of LLMs.

Research gap – However, as a key carrier connecting the digital LLMs to the physical world [19, 36], there are few benchmarks to evaluate the capabilities of LLMs in the

* Both authors contributed equally to this research.

development of embedded systems. Prior studies [14, 41] have primarily focused on code generation for embedding application, leaving their abilities in circuit design and cross-platform migration largely unexplored.

EMBEDBENCH – To address the research gap, we construct EMBEDBENCH, a new benchmark designed to evaluate the capabilities of LLMs on fundamental embedded system tasks. To facilitate the evaluation, we introduce **EMBEDAGENT**, which operates under three distinct settings that simulate real-world roles, as shown in Figure 1:

- ❶ **Embedded System Programmer:** Given a task description and a circuit schematic, LLMs are expected to write embedded system code (e.g., Arduino code).
- ❷ **Embedded System Architect:** Given a task description, LLMs are expected to design the circuit and write the corresponding embedded code.
- ❸ **Embedded System Integrator:** Given the code and the circuit schematic for one hardware platform, LLMs are expected to migrate both the circuit design and code to another platform (Arduino, ESP32, and Raspberry Pi Pico).

For each setting, there are manually constructed 126 cases in EMBEDBENCH, including task description, reference solution and automated correctness validation.

However, a comprehensive evaluation of LLMs’ capabilities in embedded system development presents several challenges. First, representing circuit schematics in a format that LLMs can understand and generate is challenging, since LLMs typically interact solely through natural language inputs and outputs. Second, verifying the correctness of the generated code is difficult. Previous work often evaluates correctness through serial output [14, 41], but this approach can be misleading, as LLMs may produce plausible outputs without ensuring correct hardware behavior. Finally, verifying LLM-generated solutions at scale remains costly and inefficient, as manual assembly and physical testing are time-consuming and resource-intensive.

Design for End-to-End Automated Evaluation – To overcome the above challenges, we design a comprehensive evaluation from **representation**, **evaluation**, and **pipeline** aspects. First, we design an **interpretable circuit representation** that allows LLMs to understand and generate circuit structures effectively, bridging the gap between natural language inputs and hardware representation. Second, to overcome the limitations of traditional evaluation based on serial output, we introduce a **hardware-driven evaluation approach**, where the behavior of virtual hardware in the Wokwi¹ environment is monitored in real-time to verify the correctness of the generated code. Third, to tackle the inefficiency of manual validation, we develop an **end-to-end automated evaluation pipeline**, which maps LLM-provided schematics to simulated connections within the

virtual environment, ensuring an efficient and scalable evaluation process.

Study – Through extensive experiments on 10 mainstream LLMs, we uncover several key findings. First, although our benchmark focuses on relatively basic tasks in embedded system development, these problems remain challenging for state-of-the-art LLMs. For example, even when provided with correct circuit schematics, the best-performing model, Deepseek-R1, achieves only 55.6% pass@1. When the correct schematics are not provided, its performance drops to 50% pass@1. In the cross-platform migration task (from Arduino to ESP32), the best result is achieved by Claude 3.7 Sonnet (Thinking), with a pass@1 score of only 29.4%. Second, we observed that different types of LLMs exhibit different weaknesses in this domain. For example, chat LLMs often struggle to flexibly apply their pretrained knowledge. In contrast, reasoning LLMs tend to approach problems from a very low-level perspective and often fail to effectively leverage the higher-level knowledge encoded during pretraining. Surprisingly, in some cases, a reasoning LLM that designs the circuit by itself can outperform the same model when given a predefined schematic.

Building on these insights, we explore two strategies to improve LLM performance in embedded system development tasks: 1) *Retrieval augmented generation*, which enhances LLMs by incorporating previous experience; 2) *Compiler feedback*, which helps correct syntax errors in the generated code. With these methods, Deepseek-R1 improves to 65.1% pass@1 when provided with correct circuit schematics, and 53.1% without them. For the Arduino to ESP32 migration task, accuracy improves from 21.4% pass@1 to 27.8%.

The contribution of this paper includes:

✓ **Benchmark EMBEDBENCH:** We introduce EMBEDBENCH, a benchmark for embedded system development. To the best of our knowledge, this is the first comprehensive benchmark designed to assess LLMs’ capability in embedded system programming, circuit design, and cross-platform migration. It includes of 126 cases, covering 9 electronic components across 3 hardware platforms.

✓ **Executable Context and automated validation mechanism:** We propose an efficient and automated evaluation pipeline for assessing embedded systems developed by LLMs. This framework streamlines the testing process, ensuring consistency and reliability.

✓ **Comprehensive Study and Vision:** We conduct a comprehensive evaluation of 10 mainstream LLMs on our benchmark and uncover several insightful findings regarding their strengths and limitations. Based on these findings, we propose two effective methods to improve LLM performance in embedded system development tasks.

¹<https://wokwi.com/>

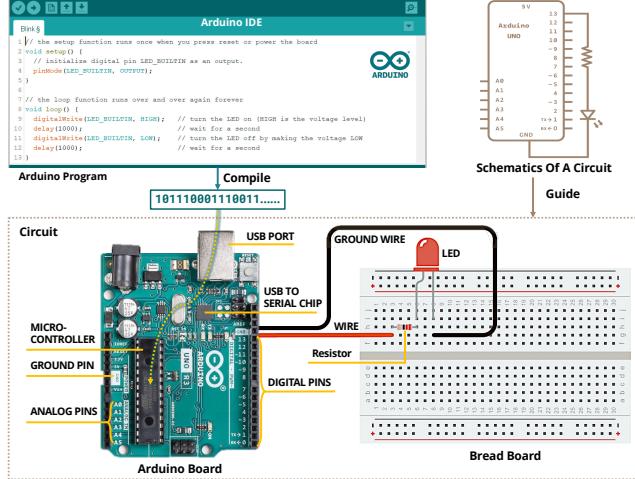


Figure 2. Arduino Workflow. The workflow requires a combination of hardware and software. On the hardware side, an Arduino board, electronic components, and a circuit schematic are needed. On the software side, code must be written based on the schematic and then uploaded to the board via USB port and UBS to serial chip, finally processed in microcontroller.

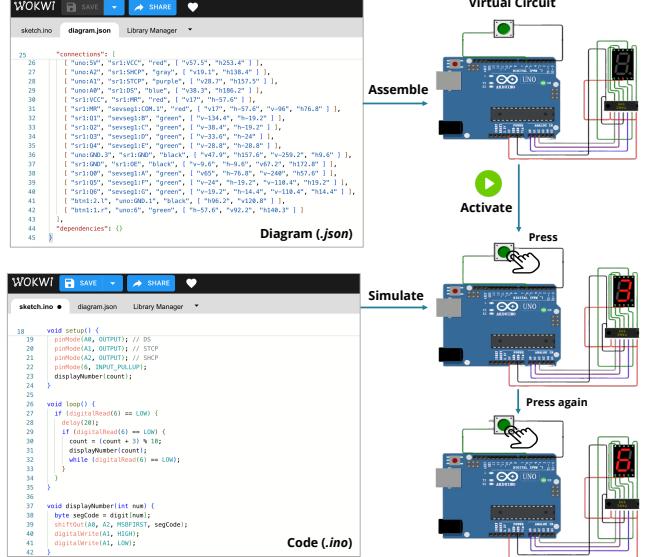


Figure 3. Illustration of Wokwi [35] - A Virtual Circuit Simulation Platform. The upper part (i.e., Diagram) shows how the virtual circuit is denoted in Wokwi, the lower part (i.e., Code) shows how the Arduino code simulates the virtual circuit. Once the virtual circuit is activated, it will be virtually powered. According to the Arduino code, the number shows in 7-segment display increase once the button has been pressed.

2 Background

In this section, we provide background information on the workflow of embedded systems and introduce Wokwi, a virtual circuit simulation platform.

2.1 Embedded System Workflow

Embedded systems are specialized computing systems that are based on hardware and driven by software. They are widely used in areas such as smart homes, robotics, and the Internet of Things (IoT). These systems serve as an important bridge between the virtual world and the physical world. Building a complete embedded system typically requires a combination of both hardware and software. Figure 2 illustrates the process of building a simple embedded system using an Arduino development board and an LED.

On the hardware side, it is necessary to prepare the development board and the required components. In this case, the components are an Arduino board and an LED. Then, based on the components and the serial port information of the board, the circuit schematic is constructed.

On the software side, code is written according to the schematic. In the figure, the Arduino IDE is used to write code that makes the LED blink every second. This code is compiled into a binary file and uploaded to the development board, enabling interaction between the virtual world and the physical world.

2.2 Virtual Circuit Simulation Platform Wokwi

Wokwi[35] is an online simulation platform for electronic circuits and embedded systems. It is primarily used to simulate the behavior of various microcontrollers such as Arduino, ESP32, and Raspberry Pi Pico, along with peripheral electronic components like sensors, LEDs, and displays. It provides a convenient virtual environment that allows developers to quickly build, test, and debug projects directly in the browser, enabling code and circuit design verification without the need for physical hardware.

Figure 3 shows an example of developing an embedded system using the Wokwi platform. Similar to real-world development, creating an embedded system in Wokwi also involves both software and hardware. On the hardware side, appropriate development boards and necessary components are selected, followed by schematics. The schematics are represented in the Diagram (.json) shown in Figure 3, where a structured file records the connections.

On the software side, the corresponding Arduino code is written, as shown in the Code (.ino) section of the figure. Once the code is complete, clicking the “Activate” button on the website compiles the program. Users can interact with the embedded system through their computer to determine whether the expected functionality has been achieved.

3 Benchmark Construction

In this section, we provide a detailed introduction to the construction process of our benchmark. It can be divided into three main steps. First, determining the electronic components required for each problem. Second, formulating appropriate problems based on the hardware usage. Finally, constructing test cases according to the problems. We also present statistics of our benchmark and analyze its quality.

3.1 Electronic Component Selection

To ensure balanced utilization of electronic components across problems, we selected seven primary hardware types as seed components: LED, RGB LED, Push Button, 7-Segment Display, LED Bar Graph, Slide Potentiometer, and Servo. Detailed descriptions of each component are provided in Appendix A. As illustrated in Figure 4 Step I, each problem is assigned a distinct set of electronic components. For each problem, we randomly select 1 to 3 components from the seven primary types. Each selected component is then randomly assigned a usage count between 1 and 3. We sample each unique hardware combination twice, resulting in a total of 126 cases, each with a unique distribution of electronic components. An example problem is shown in the Task Description of Figure 4 step II, where the sampled configuration includes two types of primary hardware (7-Segment Display and Push Button), each hardware use one instance.

3.2 Problem Construction

After determining the required electronic components for each problem, we formulate the problem descriptions through human annotation based on these components. Each problem must satisfy two key conditions. First, all selected electronic components must be actively utilized within the problem. Second, to ensure the stability of the evaluation results, any given system state must be maintained for at least 2 seconds.

Each problem description is composed of two parts. The first part provides a high-level overview, specifying the electronic components involved and the task’s general objective. The second part outlines the specific operations and the expected states of each electronic component following those operations.

A simplified version is shown in the Task Description of Figure 4, Step II. The problem description begins by stating that the electronic components include a button and a 7-segment display. The task is to implement a counter. It then specifies the detailed behavior: The initial state of the 7-segment display should show the number 0, and each button press should increment the counter by 3.

Based on the formulated problems, we develop corresponding test cases, with each problem containing 3 to 5 test cases. These test cases are carefully designed to thoroughly verify whether the code generated by the LLMs satisfies all the

requirements specified in the problem description. Additionally, we provide correct reference implementations to ensure that each problem is solvable.

For example, in Figure 4, Step II, to verify the correctness of the code, we need to check whether the output for the n th button click is equal to $(n * 3) \bmod 10$. Accordingly, the first three test cases check the behavior for 1 to 3 button clicks. A fourth test case simulates 20 button presses to verify whether the code handles extended sequences and edge conditions correctly, ensuring comprehensive evaluation.

3.3 Executable Scripts Preparation

Unlike previous work that relies on physical hardware and determines correctness through serial output

, we use Wokwi simulation platform to observe the real-time state of each electronic component at any given moment. This capability is essential for evaluating embedded systems, as the primary goal of embedded systems is to ensure that the electronic components perform the expected operations.

We use an automated submission bot to submit the code generated by LLMs and monitor the electronic components to ensure the behavior aligns with the specified requirements. Each problem is associated with a unique script containing all the test cases defined in Step II, which are used to verify correctness. For instance, in Step III of Figure 4, the validation process involves checking whether the 7-segment display shows the expected number after the push button has been pressed a specified number of times.

3.4 Quality Control

To ensure the correctness of the benchmark, each example was independently annotated by two annotators to minimize subjective bias. Upon completing their annotations, the two annotators cross-verified each other’s results. Once both annotators reached a consensus, the annotated content was reviewed by two additional volunteers for further validation. The details can be found in Appendix B (Data Annotation and Quality Control).

3.5 Dataset Statistics

In this section, we present statistical insights into our dataset, covering Problem Description Length, Reference Solution Length, the number of electronic components, and the number of Wires used. As shown in Figure 5, most problem descriptions and reference solutions fall within the 200–400 tokens range. The distribution of test cases is relatively uniform, with a majority of problems containing four test cases. Additionally, the number of electronic components is mostly concentrated between 4 and 8, while the number of wires predominantly falls in the 10–40 range.

3.6 Quality Analysis

Each problem, along with its corresponding solutions and test cases, is created by two volunteers holding bachelor’s

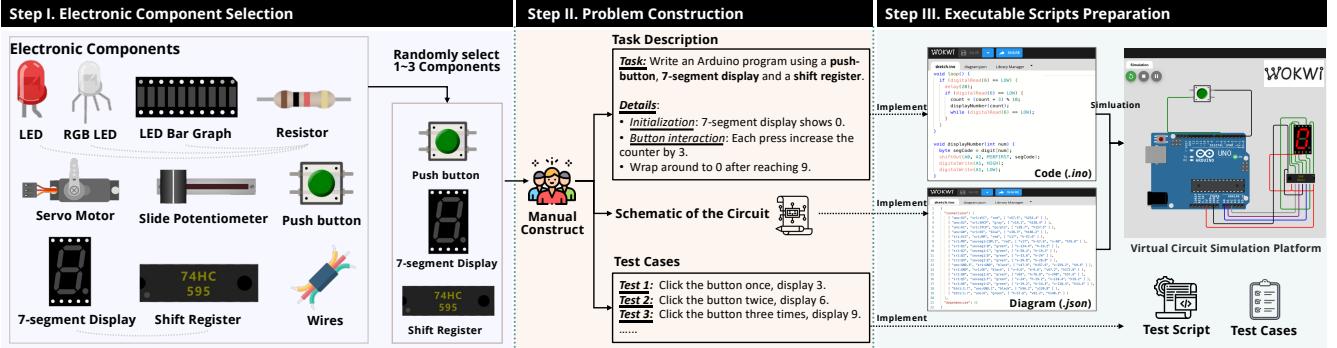


Figure 4. Data Preparation Pipeline

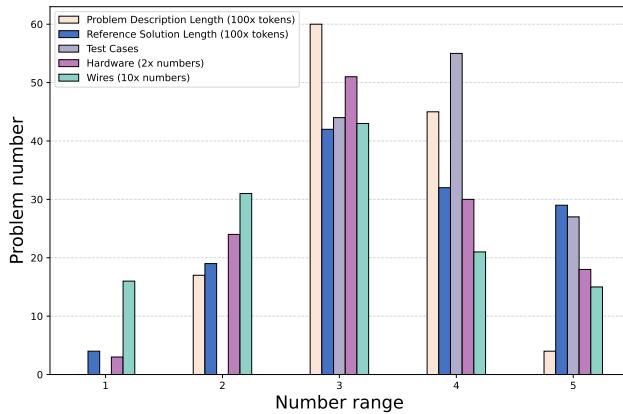


Figure 5. The dataset distribution about Problem Description Length, Reference Solution Length, Hardware numbers, Wires numbers. The horizontal axis number i means $[(i-1) * \text{interval}, i * \text{interval}]$, for example $(100x \text{ tokens})$ means interval is 100.

degrees. After the initial development, another two volunteers with bachelor’s degrees are responsible for reviewing the content to ensure its quality and accuracy. Annotators should follow these main guidelines: (1) The problem description must fully utilize all the specified electronic components. (2) The provided test cases must be sufficient to verify the correctness of the problem, ensuring that the reference code passes all of them. (3) Any states that require verification must be maintained for at least two seconds to avoid the impact of potential network delays. This multi-stage review process helps ensure the correctness, reliability, and overall quality of the questions. Further details regarding annotation rules can be found in Appendix C.

4 Experiments

4.1 Experiment Setup

4.1.1 Evaluation Tasks. To systematically evaluate embedded systems, we propose three task settings designed

to assess LLMs’ embedded system development capabilities under different conditions: Settings ①, ② and ③ which is shown in Figure 1. For the i -th problem, we define the problem description as q_i , the schematic configuration as w_i , the hardware specifications as h_i , the code solution as c_i , the test cases as u_i , the programming language as L , the development board as B , and the LLMs used for code generation as M . The LLMs’ predicted probability distribution is denoted by $P(\cdot)$, while $I(\cdot)$ represents a validation function to verify the correctness of the LLM’s output. The evaluation result r_i is an n -dimensional binary matrix, where n is the number of test cases, and each element is either 1 (indicating the test case is passed) or 0 (indicating the test case is failed).

Setting ① Embedded System Programmer Prompt
Problem Description
You are tasked with programming an Arduino (mega) to create an interactive game using three LEDs ...
Components Usage
{'type': 'arduino-mega', 'id': 'mega'}
...
Components Information
type: led
pins:
'A': The positive pin of the LED.
'C': The negative pin of the LED. ...
Connections
['mega:GND.1', 'led_red:C']
['led_red:A', 'mega:2']
['mega:GND.1', 'led_yellow:C']
...

Table 1. The prompt of Setting ① task

Setting ① As illustrated in Figure 1, this task requires LLMs to generate embedded system code based on a given problem description and a schematic configuration. Successfully completing the task demands a comprehensive understanding of the problem requirements, the functionality of each

Setting ② Embedded System Architect Prompt	Setting ③ Embedded System Integrator Prompt
<p>### Problem Description You are tasked with programming an Arduino (mega) to create an interactive game using three LEDs ...</p> <p>### Components Usage {`type': 'arduino-mega', 'id': 'mega'} {`type': 'rgb-led', 'id': 'rgb1'} ...</p> <p>### Components Information type: rgb-led pins: 'R': Controls the brightness of the red LED. ... ### Example Connections Here is an example of a connection format: [CONNECTIONS] [{"uno:A2", "sr1:SHCP"}] [{"sr1:MR", "sevseg1:COM.1"}] [{"rgb1:G", "uno:10"}] ... [{"r1:1", "bargraph1:C1"}] [CONNECTIONS] ...</p>	<p>### Arduino Code void setup() { ... } ...</p> <p>### Components Usage ... ### Components Information type: pi-pico pins: 'GP0' to 'GP22': Digital GPIO pins for input/output operations. ... ### Arduino Connection ...</p>

Table 2. The prompt of Setting ② task

hardware component, and the roles of individual pins. This task evaluates the LLM’s capability as a programmer. The process can be formalized as:

$$r_i = I(P(c_i | q_i, w_i, h_i, L, B; M); u_i) \quad (1)$$

The prompt for the “Setting ①” tasks is shown in Tabel 1. First, the LLM is provided with a problem description. Next, a list of components is given, where each component is represented as a dictionary containing two keys: the component’s type and its unique identifier. Following this, detailed specifications for each component type are presented, including the available pins and their respective functions. Lastly, the schematics are described. Each connection is specified as a pair of strings in the format “id:pin”, denoting the connection between specific pins of two components. For example, [“uno:GND.1”, “led1:COM”] represents a connection between pin GND.1 of the uno component and the COM pin of the led1 component. The LLM is expected to generate Arduino code that utilizes the given components and schematics configuration to solve the specified problem.

Setting ② As shown in Figure 1, in this task, in addition to generating code for the embedded system, LLMs are also required to design circuits. This requires a deeper understanding of hardware components, familiarity with the functionalities of various serial ports, and a certain level of logical reasoning and spatial capabilities. This task evaluates the

Table 3. The prompt of Setting ③ task

LLM’s capability as an expert in embedded system development. The process can be formalized as:

$$r_i = I(P(w'_i, c'_i | q'_i, h'_i, L', B'; M); u_i) \quad (2)$$

As shown in Table 2, the prompt is constructed similarly to the “Setting ①” task. LLMs are given a task description without schematics of circuits. To format the output, we include example connections to guide the LLM in constructing an appropriate schematic in the required format. The LLM is expected to return both the schematic for the given components and the Arduino code necessary to control the hardware components and complete the task.

Setting ③ As shown in Figure 1, the “Setting ③” task requires LLMs to migrate both the schematics and code from a given source platform to a target platform. The LLM must consider the similarities and differences between platforms, including serial ports, programming languages, and the syntax of various function calls. This process can be formalized as:

$$r_i = I(P(w'_i, c'_i | h'_i, L', B', w_i, c_i; M); u_i) \quad (3)$$

Here, the variables with superscripts (w'_i, c'_i, h'_i, L', B') denote the corresponding elements in the target platform. Detailed information of each platform used in this task is provided in Appendix B.

Unlike the “Setting ①” and “Setting ②” tasks. As shown in Table 3, the “Setting ③” task does not include an explicit problem description. Instead, because it involves platform migration, the input includes the Arduino code and schematics connections. The component usage and specifications remain consistent with the previous two tasks, with the addition of the target platform’s details. The LLM must first translate the schematics from Arduino to the target platform and then generate equivalent code that replicates the same functionality.

4.1.2 Evaluation Method. pass@ k is a metric commonly used in test-driven development, where software requirements are first converted into test cases, and success is defined by a program that passes these tests. The formula is given in Equation 4, where n represents the number of generated samples per problem, k denotes the number of samples evaluated and c is the number of correct samples:

$$\text{pass}@k := \mathbb{E}_{\text{Problems}} \left[1 - \frac{\binom{n-c}{k}}{\binom{n}{k}} \right] \quad (4)$$

We use pass@1 as the evaluation metric following prior work [8, 21]. The temperature is set to be 0 and employ greedy decoding for generation. A problem is considered successfully solved if the LLM passes all test cases.

4.1.3 LLMs for evaluation. We select 10 state-of-the-art LLMs spanning 6 different LLM families, and categorize them into three types for evaluation: **Reasoning LLMs** (O3-mini[28], Claude 3.7 Sonnet (Thinking)[9], DeepSeek-R1[15], QWQ[32]), **SFT-Distilled Reasoning LLMs** (DeepSeek-R1-Distill-Llama[15], DeepSeek-R1-Distill-Qwen[15]), and **Chat LLMs** (DeepSeek-V3[11], Llama-3.3-Instruct[3], Starcoder2[24], Qwen2.5-Instruct[39]).

4.2 Overall Result

Table 4 presents the pass@1 evaluation results for different types of LLMs, arranged in descending order based on their parameter size. The key observations are as follows:

1) EMBEDBENCH is challenging for studied LLMs. In both the “Setting ①” and “Setting ②” tasks, the highest-performing LLM, DeepSeek-R1, achieves only 55.6 and 50.0 pass@1, respectively. For the “cross-platform” task, most LLMs struggle to transfer Arduino code to ESP32’s ESP-IDF effectively. Even the best-performing LLM, Claude 3.7 Sonnet (thinking), reaches only 29.4 pass@1. Transferring code to the Raspberry Pi Pico using MicroPython is slightly more successful, with a pass@1 of 73.8 at best.

2) Reasoning LLMs outperform Chat LLMs. However, SFT-distilled reasoning LLMs show mixed results compared to their base versions. For example, QwQ surpasses its base LLM, Qwen2.5-Instruct, across all tasks. However, DeepSeek-R1-Distill-Qwen does not demonstrate a significant improvement over its base LLM. Additionally, DeepSeek-R1 consistently outperforms DeepSeek-V3 in all tasks.

5 Analysis

We first analyze why LLMs struggle to construct embedded systems by investigating concrete cases. Next, we unpack the chain of thought generated by different reasoning LLMs when tackling embedded system-related tasks. Finally, we try some methods to explore LLMs’ upper bound of capabilities in constructing embedded systems.

Models	Size	Task			
		Setting ①		Setting ②	
		Arduino C++	Arduino C++	ESP32 ESP-IDF	Pi Pico MicroPy
<i>Reasoning LLM</i>					
Claude 3.7	\	48.4	48.4	29.4	73.8
O3-mini	\	47.6	38.9	16.7	61.9
DS-R1	671B/31B	55.6	50.0	21.4	60.3
QWQ	32B	39.7	50.0	4.0	52.4
<i>SFT-Distilled Reasoning LLM</i>					
DS-R1-Dist.	70B	35.7	23.8	0.0	11.1
DS-R1-Dist.	32B	27.0	11.1	3.2	27.8
<i>Chat LLM</i>					
DS-V3	671B/31B	37.3	25.4	2.4	46.8
Llama-3.3-Ins.	70B	32.5	21.4	0.0	23.0
Qwen2.5-Ins.	32B	25.4	23.8	3.2	33.3
Starcoder2	15B	11.1	2.4	0.0	4.0

Table 4. pass@1(%) results of each LLMs. Abbreviations: sch (schematics), DS(DeepSeek), Claude 3.7(Claude 3.7 Sonnet (thinking)), DS-R1-Dist. 70B (DeepSeek-R1-Distill-Llama), DS-R1-Dist. 32B (DeepSeek-R1-Distill-Qwen), Ins. (Instruct). Ds-R1 and DS-V3 are MOE LLMs with 32B parameters activated during inference.

5.1 Case Study

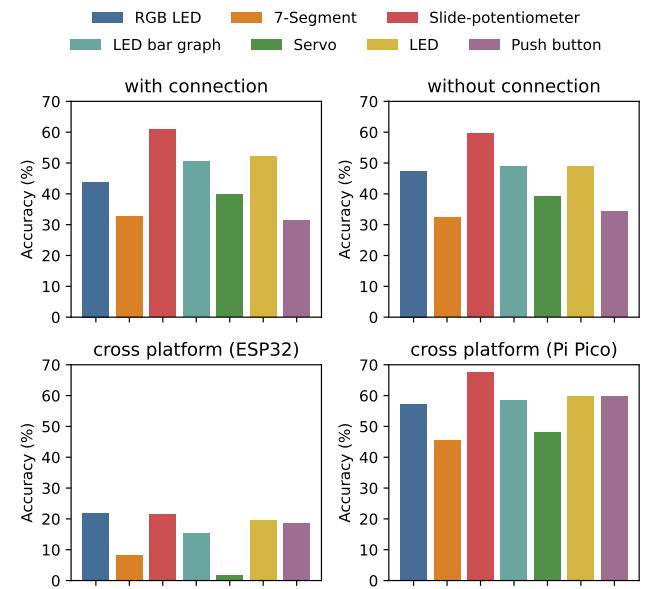


Figure 6. The average accuracy of reasoning LLMs (QwQ, DeepSeek-R1, O3-mini, Claude 3.7 Sonnet (Thinking)) on problems involving specific electronic components.

To investigate why LLMs struggle with constructing embedded systems, we categorize problems based on the specific electronic components they involve. For instance, if a problem requires the use of both a “Servo” and an “LED”, we classify it under both “Servo” and “LED” categories. We then measure the average accuracy of various reasoning LLMs (QwQ, DeepSeek-R1, O3-mini, Claude 3.7 Sonnet (Thinking)) across seven major electronic components in four different tasks. The results are presented in Figure 6. From the figure, we observe that LLMs struggle particularly with problems involving “7-Segment displays” and “Push Buttons” in both the “Setting ①” and “Setting ②” tasks. In the “Setting ② (Pi Pico)” task, LLMs generally perform well. However, in the “Setting ② (ESP32)” task, accuracy is consistently low across all electronic components. In the following content, we provide a detailed analysis of these challenges, focusing on the “7-Segment”, “Push Button” and the ESP32 platform.

7-Segment Display *LLMs struggle to determine the correct voltage levels for each port when displaying numbers (0 to 9) or characters such as “A” and “P”.*

A significant portion of the errors in the “Setting ①” task related to 7-segment displays comes from this issue, accounting for 48 out of 116 errors (approximately 41.4%). Similarly, in the “Setting ②” task, 45 out of 103 errors (approximately 43.7%) fall into the same category.

Listing 1 Subject-101 (QWQ)

```

1 // Lookup table for 7-segment display
2 // (common anode)
3 byte seg_code[10] = {
4     0x90, // 0
5     0xF9, // 1
6     0xA2, // 2
7     0xB0, // 3
8     0x99, // 4
9     0x92, // 5
10    0x82, // 6
11    0xF8, // 7
12    0x80, // 8
13    0x90 // 9
14 };

```

An example is shown in Listing 1. The LLM uses hexadecimal values to represent the voltage levels of each segment in the 7-segment display. For instance, to represent the number 1, it uses the value 0xF9, which corresponds to the binary pattern 0b1111001. This means that two segments are set to low voltage (i.e., they are activated in a common-anode display), correctly displaying the digit 1. However, the LLM struggles to consistently generate correct values for each digit. In this example, both the digits 0 and 9 are represented using the same value 0x90, which indicates an error, as they should have different segment configurations.

Push Button *For push buttons, LLMs often struggle with handling button debounce.*

In the “Setting ①” task, nearly half of the errors related to push buttons, 43 out of 108 (approximately 39.8%), are caused by improper debounce handling. Similarly, in the “Setting ②” task, 51 out of 106 errors (approximately 48.1%) result from the same issue.

Listing 2 Subject-76 (Claude 3.7 sonnet (thinking))

```

1 // debounceDelay = 150ms
2 if (digitalRead(btn1Pin) == LOW && (
3     currentTime - lastBtn1Press) >
4     debounceDelay) {
5     lastBtn1Press = currentTime;
6     incrementNumber(1);
7 }

```

As shown in Listing 2, the LLM implements a basic debounce mechanism by checking whether the button is pressed (LOW) and ensuring that at least 150 ms have passed since the last press. However, even though the button press duration is intended to be 150 ms, various factors can make it difficult to maintain this timing precisely. As a result, this approach may lead to false triggers, where a single press is misinterpreted as multiple presses. A more reliable method would involve tracking the previous button state and updating the debounce timer only when a state change occurs (for example, from HIGH to LOW). This helps reduce the likelihood of erroneous detections.

ESP32 *The primary issue observed in ESP32 code generated by LLMs is the high frequency of syntax errors.*

A substantial portion of the generated code, 211 out of 504 cases (approximately 41.9%), contains such errors. Common problems include missing or incorrect usage of relevant header files, as well as the use of functions that are incompatible with the specified version of the ESP-IDF, even when the version is clearly provided in the prompt. When comparing Python and ESP-IDF code generation, it becomes evident that current LLMs lack sufficient training in domain-specific embedded programming languages such as ESP-IDF.

Finding 1: The main challenges for LLMs are: (1) Using a 7-Segment Display to display numbers or letters. (2) Debouncing button inputs. (3) Syntax of ESP-IDF.

5.2 Deep thinking for Embedded System

After conducting case studies on hardwares and platforms where LLMs are prone to make errors, we further investigate the underlying causes of these errors. Specifically, we analyze the reasoning chains produced by the reasoning LLMs to better understand how and why such errors occur in the development of embedded system.

5.2.1 The thinking Format of each LLMs. Given that current LLMs are trained on limited data related to embedded system development, we begin by analyzing the differences

in the chains of thought produced by each reasoning LLM when confronted with this unfamiliar task.

Claude 3.7 Sonnet (Thinking) devotes a significant portion of its reasoning process to debugging. This frequent and thorough debugging may help explain why Claude 3.7 Sonnet (Thinking) achieves the best performance on Setting ❸ tasks (e.g., ESP32), where LLMs are particularly prone to syntax errors. The reasoning process of Claude 3.7 Sonnet (Thinking) begins with a brief preparation phase that involves understanding the question, analyzing schematics diagrams, and outlining a plan for code generation. Following this initial stage, the model spends an average of 4,440 tokens (out of a total 5,703, or approximately 77.8% of its reasoning chain) identifying and resolving issues in the generated code. As illustrated in Table 5, the model first uses 955 tokens to generate an initial answer, then proceeds to debug the code multiple times

Claude 3.7 Sonnet Thinking Process	
Generate the answer use 955 tokens	Wait , the connections say that ...
Let's analyze what we have and what we need to do to migrate from Arduino Mega to a Pi Pico running MicroPython ...	Wait , let me think again. Let me list the segments: ...
Here's the original Arduino code: ...	Wait , the DP isn't connected here, so we can ignore it. ...
Repair the code use 2025 tokens	Now, moving on to the code structure.
Wait, I need to refine this a bit ...	Wait , but the LEDs should alternate: ...
Let me further check and correct this code ...	Wait , the STCP is the latch pin. ...
Let me fix the code ...	Wait , the standard procedure is: ...
This logic is a bit confusing. Let me simplify ...	Testing the code:
Wait, I need to clarify the LED logic once more ...	Wait , looking back at the connections:
I think this code should work correctly, but let's make one final review to ensure it's accurate ...	Wait , the connections say:...
I believe the code is correct, but I just realized we need to fix the GPIO numbers ...	The only possible issue is ...
	...

Table 5. An example of Claude 3.7 Sonnet (thinking) Thinking Process

DeepSeek-R1 and QWQ. exhibit similar reasoning patterns, marked by frequent reflection throughout each stage of the problem-solving process. A typical indicator of this reflective behavior is the word “wait”[15, 27], which often signals moments of uncertainty or perceived inconsistency. This tendency toward continual self-evaluation may contribute to their strong performance on both code generation tasks which require advanced reasoning skills.

Compared to LLMs like Claude 3.7 Sonnet, DeepSeek-R1, and QWQ demonstrate a more introspective approach. From understanding the initial problem to analyzing potential flaws in the generated code, these LLMs repeatedly assess and revise their reasoning. The keyword “wait” frequently

appears in their outputs as a marker of this process. On average, DeepSeek-R1 uses “wait” approximately 20.4 times per problem, while QWQ uses it about 40.7 times. An example is shown in 6. From the example, we can see that Deepseek-R1 reflects on all stages of problem understanding, code structure construction, and code testing in a balanced manner.

Deepseek-R1 Thinking Process
Okay, let's tackle this problem step by step. First, I need to understand all the requirements and how the components are connected .
...
Wait , the connections say that ...
Wait , let me think again. Let me list the segments: ...
Wait , the DP isn't connected here, so we can ignore it. ...
Now, moving on to the code structure .
Wait , but the LEDs should alternate: ...
Wait , the STCP is the latch pin. ...
Wait , the standard procedure is: ...
Testing the code:
Wait , looking back at the connections:
Wait , the connections say:...
The only possible issue is ...
...

Table 6. An example of Deepseek-R1 Thinking Process

SFT-distilled reasoning LLMs show significant instability when tackling embedded system development tasks. This instability suggest that SFT-distilled reasoning LLMs struggle with unfamiliar tasks, and may explain that in some cases these LLMs underperform compared to their base LLMs.

Compared to their teacher LLMs, both DeepSeek-R1-Distill-Llama and DeepSeek-R1-Distill-Qwen exhibit significantly higher variance in token usage during reasoning. For the “Setting ❶” and “Setting ❷” tasks, DeepSeek-R1-Distill-Llama shows variances of 40,552,828 and 35,707,662 tokens, respectively, while DeepSeek-R1-Distill-Qwen reaches 118,677,362 and 89,303,221 tokens—over an order of magnitude greater than the teacher model’s 11,435,790 and 7,748,574.

On Setting ❸ tasks (ESP32 and Pi Pico), DeepSeek-R1-Distill-Llama averages only 696 and 569 thinking tokens, and DeepSeek-R1-Distill-Qwen averages 540 and 849, compared to DeepSeek-R1’s 3,760 and 3,848. These findings reinforce the idea that SFT-distilled reasoning models are more brittle and less consistent when faced with unfamiliar problem domains.

Finding 2: (1) Claude 3.7 Sonnet (Thinking) tends to focus more on code repair, which gives it an advantage in ensuring syntactic correctness. (2) DeepSeek-R1 and QWQ engage in reflective thinking throughout the reasoning process, which helps ensure the accuracy of their logical reasoning. (3) SFT-distilled reasoning LLMs exhibit significant instability in their thought processes when faced with this entirely new task.

5.2.2 Why struggle to display 7-Segment? To further explore why LLMs encounter difficulties in embedded system development, we analyze their performance on a specific hardware task: controlling a 7-segment display. This task is widely used in embedded systems and is likely represented in the training data of many LLMs. However, it also requires a multi-step reasoning process that includes hardware comprehension, pin mapping, and display logic. Based on the case study, we select three LLMs with varying capabilities, ranging from weak to strong, and evaluate how each performs when solving problems related to 7-segment displays:

Listing 3 Typical Digit table for 7-segment display

```

1  /* Digit table for the 7-segment display */
2  byte digits[10] = {
3      0b00111111, // 0
4      0b00000110, // 1
5      0b01011011, // 2
6      0b01001111, // 3
7      0b01100110, // 4
8      0b01101101, // 5
9      0b01111101, // 6
10     0b00000111, // 7
11     0b01111111, // 8
12     0b01101111 // 9
13 };

```

① **Starcoder2** is a code LLM trained on Stack v2 (67.5TB of GitHub code), which is publicly available. Using the provided tool [25], we verified that the digit table shown in Listing 3 is indeed present in the Stack v2 training set. Despite this, evaluation results reveal that StarCoder2 achieved 0% accuracy (0/45) on 7-segment display tasks, even in “Setting ①” task where the schematics are included. Moreover, only 22.2% (10/45) of its generated solutions incorporated a digit table similar to the one in Listing 3. These findings indicate that **LLMs with small parameter sizes struggle to retrieve and apply pretrained knowledge effectively** without sufficient exposure to task-specific training.

② **DeepSeek-V3** is a state-of-the-art chat LLM that performs well on a variety of general tasks. However, in the “Setting ①” task, which uses the 7-segment display, it achieves only 6.7% accuracy (3/45), despite including a digit table in 93.3% (42/45) of its outputs. While DeepSeek-V3 successfully activates knowledge related to 7-segment display data present

in its pretraining corpus, only 15.5% (7/45) of its generated solutions correctly adapted the digit table to a common anode configuration as required by the task. Even fewer solutions incorporated more complex changes in response to the problem’s requirements. These results indicate that while larger models trained on more extensive datasets **demonstrate an increased ability to recall relevant pretraining knowledge**, they still face challenges in flexibly applying that knowledge based on the task-specific contexts.

③ **DeepSeek-R1** is a reasoning LLM that differs from chat LLMs by employing *a longer chain-of-thought* before output. This reasoning ability enables it to outperform the previous models, achieving 44.4% accuracy (20/45) on the “Setting ①” task. However, reasoning-based LLMs also face challenges in 7-segment tasks. Unlike its base model, DeepSeek-V3, which only reuses pretraining knowledge with minimal adaptation, DeepSeek-R1 tends to over-analyze. On average, it uses approximately 4429 tokens ((4429/10342) 42.8% of the total chain of thought) focusing on binary representations of digits—an unnecessary detour for this relatively simple task. This excessive reasoning often leads to errors, with some chains of thought even stating, “But this is tedious.” (See Appendix D). A more effective approach would be to retrieve a known digit table (e.g., Listing 3), and then adapt it based on the task requirements. These results show that reasoning LLMs may **entangle in unnecessarily complex thought processes**, limiting their ability to apply pretrained knowledge efficiently in domain-specific tasks.

Finding 3: For 7-segment display tasks, different type of LLMs face different difficulties: (1) For chat LLMs, the primary challenge lies in retrieving relevant pretraining knowledge and adapting it appropriately to the task-specific constraints. (2) For reasoning LLMs, the issue stems from overcomplicating the problem, leading to inefficient use of known information and reduced task performance.

5.2.3 The gap between Setting ① and Setting ②. From the results in Table 4, we observe an interesting phenomenon. For chat LLMs and SFT-distilled reasoning LLMs, performance on the “Setting ①” task is significantly better than on the “Setting ②” task. However, for reasoning LLMs, the “Setting ①” task does not show a clear advantage over the “Setting ②” task. In fact, for the QWQ model, the “Setting ②” performance is even better than the “Setting ①” result. In this section, we aim to analyze why the performance gap between these two tasks varies across different types of LLMs.

As discussed in Section 4.2.1, reasoning LLMs such as DeepSeek-R1 and QwQ often use the keyword “wait” to indicate moments of uncertainty or inconsistencies in their reasoning process. Therefore, the frequency of the word “wait” can serve as a rough proxy for the LLM’s perceived difficulty in handling a task. Interestingly, DeepSeek-R1 used

QwQ Thinking Process in Different Tasks	
Setting ① thinking process	
Okay, let's tackle this problem step by step ...	
Looking at the connections:	
- The RGB LED's R, G, B pins are connected to Mega's digital pins 2,3,4. Since these are PWM pins (pins 2-13 have PWM), we can use analogWrite for color control ...	
- The 7-segment is driven via the shift register (sr1). The shift register's DS, SHCP, STCP are connected to Mega's A0, A1, A2. Wait, no, checking the connections again: the connections list says: ...	
Looking at the connections for the 7-segment: ...	
Think 9276 tokens	
So in the code: ...	
The bar graph and RGB LED parts should work as intended.	
Setting ② thinking process	
Okay, let's tackle this problem step by step.	
...	
Now, let's outline the connections:	
For the shift register (sr1):	
- DS (data) to Arduino pin 2	
- SHCP (shift clock) to pin 3	
- STCP (latch clock) to pin 4	
- OE to GND (to enable)	
- MR to 5V (to not reset)	
VCC to 5V, GND to GND...	
For the 7-segment display,	
Think 1565 tokens	
Now, compiling all the connections into the required format.	
Now, the code: ...	

Table 7. An example of QwQ Thinking Process in “Setting ①” and “Setting ②”

the word “wait” 3,369 times in the “Setting ②” task, but this number increased to 4,503 in the “Setting ①” task. The discrepancy is even more pronounced for QwQ, which used “wait” 5,779 times in the “Setting ②” setting and 9,782 times in the “Setting ①” setting. Contrary to expectations, these results suggest that for LLMs with strong reasoning capabilities, it may actually be more effective to allow them to perform the schematics themselves, rather than supplying the correct schematics in advance.

To better understand this phenomenon, we examined the cases that QwQ successfully solved in the “Setting ②” task but failed to solve in the “Setting ①” setting. We found that the model often discovered alternative schematic schemes that were easier for it to reason about than the official, correct schematics. A representative example is shown in Table 7. In this example, the main difference lies in how the

74HC595’s DS, SHCP, and STCP pins are connected. In the correct schematics, since some tasks involve more complex hardware setups with limited digital ports, we standardize the pin connections to A0, A1, and A2 on the Arduino. However, for this specific task, connecting to pins 2, 3, and 4 would be a more natural and efficient choice. The LLM selected this alternative schematics in the “Setting ②” setting, leading to a smoother reasoning process. This small change had a significant impact. In the “Setting ②” task, QwQ used only 1,565 tokens to construct a digital logic table. In contrast, for the “Setting ①” task—despite being given the correct schematics—it used 9,276 tokens, indicating a much more complex and less efficient reasoning path.

Finding 4: For reasoning LLMs, rigidly enforcing a predefined reasoning path Instead, allowing LLMs the flexibility to determine their own schematic strategy can result in more effective problem-solving.

5.3 LLMs’ Upper-bound for Embedded Systems

Given the subpar performance of LLMs in our benchmark, this section explores methods to unlock their potential for embedded system development based on the analysis above. **R1-Retrieval** As discussed in Section 4.2.2, we observe that LLMs encounter different challenges in the code generation task (Setting ① and ②). Chat LLMs can leverage pretrained knowledge but often struggle to adapt flexibly to specific problem contexts. In contrast, reasoning LLMs tend to approach tasks from an overly low-level perspective, missing higher-level abstractions. Ideally, an effective LLM should combine pretrained knowledge with adaptive reasoning capabilities, enabling it to adjust and apply its understanding to the specific problem dynamically.

To address this, we propose a method called R1-Retrieval, which combines the strengths of both DeepSeek-V3 and DeepSeek-R1. The approach is as follows: we first use LLM DeepSeek-V3 to generate code for embedded system tasks that can be addressed using pretrained knowledge alone. This generated code is then provided to DeepSeek-R1 as a form of retrieval-augmented input within its prompt. With this setup, DeepSeek-R1 can leverage the pretrained output of DeepSeek-V3 as contextual knowledge, enabling it to reason more effectively and solve problems that require both learned knowledge and adaptive reasoning ability.

For evaluation, we selected a simple task for each of the seven main electronic components—tasks that DeepSeek-V3 can handle with ease. For each target task, if it involves a particular electronic component, the corresponding code from DeepSeek-V3 is included in the retrieval input. These code and the associated prompts are detailed in the Appendix E. The results are summarized in Table 8. With the R1-Retrieval method, the pass@1 score of DeepSeek-R1 improved from

Methods	Setting ①	Setting ②	Setting ③ (ESP32)
DeepSeek-R1	55.6	50.0	21.4
R1-Retrieval	65.1	53.1	\
R1-Compiler	\	\	27.8

Table 8. The result of R1-Retrieval and R1-Compiler

55.6 to 65.1 on the “Setting ①” task, and from 50.0 to 53.1 on the “Setting ②” task.

R1-Compiler. For Setting ③ task development using the ESP32 and the ESP-IDF framework, LLMs often generate code with syntax errors, making deployment infeasible. To mitigate this, we introduce R1-Compiler, a method that incorporates the compiler as an external feedback mechanism to guide code repair. Specifically, we capture compiler error messages and feed them back into DeepSeek-R1, prompting it to revise its output accordingly. The detailed prompting strategy is provided in the Appendix F.

The incorporation of compiler feedback significantly reduces syntax errors from 34.1% (43/126) to 3.2% (4/126) after just one repair iteration. The final result is shown in Table 8, the pass@1 score increases from 21.4 to 27.8.

5.4 Implications

Although LLMs perform poorly on our fundamental embedded systems design benchmark, they exhibit remarkable potential. For example, even without being given the correct schematics, they are able to rely on their inherent reasoning abilities to find better solutions. When activated with pre-trained knowledge and provided with the correct schematics, the reasoning LLM can achieve a pass@1 score of 65.1%.

We believe that with more data and improved training methods, such as reinforcement learning, multi-agent collaboration, or multimodal LLM, it is possible to enable interaction between the virtual and physical worlds solely through LLMs. This could represent the first step toward achieving embodied intelligence.

6 Related Work

Benchmark for Large Language Models.

As LLMs continue to advance, there is an increasing need to evaluate their performance across a wide variety of tasks. Consequently, more and more benchmarks have been developed to assess LLMs. These benchmarks and LLMs mutually drive progress: benchmarks highlight limitations in specific areas, guiding improvements, while more capable LLMs motivate the development of more challenging and higher-quality benchmarks. Existing benchmarks mainly cover areas such as natural language processing[1, 2, 6], logical reasoning[33, 42], mathematics[4, 10, 34], and code[31, 38, 44]. However,

most benchmarks are limited to tasks in virtual environments, lacking evaluation of LLMs’ capabilities to interact with the physical world.

Benchmark for Embedded System. There are also benchmarks related to embedded systems [14, 23, 41], but they mainly focus on evaluating LLMs from a coding perspective rather than assessing their ability to design embedded systems from scratch. Also, these evaluation methods often rely on serial output embedded within the generated code, which introduces a clear gap between the simulated results and the actual performance on hardware. Evaluating LLMs on real hardware is also significantly more costly and resource-intensive. Our benchmark leverages the Wokwi simulation platform, enabling automated evaluation that closely reflects real hardware behavior, without the time and resource costs associated with using physical hardware components.

Large Language Models for Embedded Systems.

Existing approaches that use LLMs for hardware platform programming primarily rely on integrated development environments or programming frameworks. However, they rarely enable genuine interaction between the LLM and the hardware[5, 7, 12, 20, 29]. Some studies investigated the use of LLMs for hardware-related work. These include generating code for Field Programmable Gate Arrays (FPGAs)[16, 26, 30], using LLMs to write Arduino code but focusing on human-in-the-loop testing[14], or leveraging serial output to provide additional information for LLMs debugging[41]. Despite these advancements, most of the work remains limited to code generation and does not involve complete embedded system design. Our benchmark is the first to allow an LLM to independently implement the full process of embedded system development.

7 Threat to Validity

Our study identifies three potential threats to validity, along with the steps we have taken to mitigate them.

First, the manual construction of our benchmark dataset introduces the possibility of subjective bias in labeling. To reduce this risk, two independent annotators labeled the data separately, and any differences were discussed and resolved through consensus.

Second, the LLMs used in our experiments inherently involve randomness, which can cause variability in the observed results. To mitigate this issue, we set the decoding temperature to 0 and employed greedy decoding during inference, ensuring deterministic and consistent results across multiple runs.

Third, our evaluation relies on the correctness and stability of the Wokwi simulation platform. Wokwi supports a large and active user community, hosting more than one million Arduino projects, 800,000 ESP32 projects, and 100,000 Raspberry Pi Pico projects. During the process of creating our benchmark and associated test cases, we encountered

no unexpected behaviors, which enhances our confidence in the platform's reliability.

8 Conclusion

In this work, we propose a benchmark to evaluate the ability of LLMs in embedded system design. Our benchmark consists of four types of tasks, each with 126 cases. These tasks assess the LLMs' capabilities in embedded programming with given wiring, autonomous circuit wiring, end-to-end embedded system design, and cross-platform migration. We also provide an automated and accurate evaluation method and conduct experiments on 10 state-of-the-art LLMs. Through our analysis, we find that although the LLMs face challenges in solving these tasks, they also demonstrate significant potential. For instance, discovering wiring schemes that are even better than the reference solutions. Based on our findings, we propose several simple strategies that can effectively improve the performance of LLMs in embedded system development.

References

- [1] Ahmed Abdelali, Hamdy Mubarak, Shammur Absar Chowdhury, Maram Hasanain, Basel Mousi, Sabri Boughorbel, Yassine El Kheir, Daniel Izham, Fahim Dalvi, Majd Hawasly, et al. 2023. LAraBench: Benchmarking Arabic AI with large language models. *arXiv preprint arXiv:2305.14982* (2023).
- [2] Kabir Ahuja, Harshita Diddee, Rishav Hada, Millicent Ochieng, Krithika Ramesh, Prachi Jain, Akshay Nambi, Tanuja Ganu, Sameer Segal, Maxamed Axmed, et al. 2023. Mega: Multilingual evaluation of generative ai. *arXiv preprint arXiv:2303.12528* (2023).
- [3] Meta AI. 2024. Introducing meta llama 3: The most capable openly available llm to date. Blog. <https://ai.meta.com/blog/meta-llama-3/> Online; accessed 15-January-2024.
- [4] Daman Arora, Himanshu Gaurav Singh, et al. 2023. Have LLMs advanced enough? A challenging problem solving benchmark for large language models. *arXiv preprint arXiv:2305.15074* (2023).
- [5] Thomas Ball, Abhijith Chatra, Peli de Halleux, Steve Hodges, Michal Moskal, and Jacqueline Russell. 2019. Microsoft MakeCode: embedded programming for education, in blocks and TypeScript. In *Proceedings of the 2019 ACM SIGPLAN symposium on SPLASH-E*. 7–12.
- [6] Y Bang, S Cahyawijaya, N Lee, W Dai, D Su, B Wilie, H Lovenia, Z Ji, T Yu, W Chung, et al. 2023. A multitask, multilingual, multimodal evaluation of ChatGPT on reasoning, hallucination, and interactivity (arXiv: 2302.04023). arXiv.
- [7] Robert W Brennan and Jonathan Lesage. 2022. Exploring the implications of OpenAI codex on education for industry 4.0. In *International Workshop on Service Orientation in Holonic and Multi-Agent Manufacturing*. Springer, 254–266.
- [8] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374* (2021).
- [9] Claude. 2025. Claude 3.7 Sonnet. <https://www.anthropic.com/clause/sonnet>
- [10] Katherine M Collins, Albert Q Jiang, Simon Frieder, Lionel Wong, Miri Zilka, Umang Bhatt, Thomas Lukasiewicz, Yuhuai Wu, Joshua B Tenenbaum, William Hart, et al. 2024. Evaluating language models for mathematics through interactions. *Proceedings of the National Academy of Sciences* 121, 24 (2024), e2318124121.
- [11] DeepSeek-AI. 2024. DeepSeek-V3 Technical Report. arXiv:2412.19437 [cs.CL] <https://arxiv.org/abs/2412.19437>
- [12] James Devine, Joe Finney, Peli de Halleux, Michal Moskal, Thomas Ball, and Steve Hodges. 2019. MakeCode and CODAL: Intuitive and efficient embedded systems programming for education. *Journal of Systems Architecture* 98 (2019), 468–483.
- [13] Xueying Du, Mingwei Liu, Kaixin Wang, Hanlin Wang, Junwei Liu, Yixuan Chen, Jiayi Feng, Chaofeng Sha, Xin Peng, and Yiling Lou. 2024. Evaluating large language models in class-level code generation. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*. 1–13.
- [14] Zachary Englhardt, Richard Li, Dilini Nissanka, Zhihan Zhang, Girish Narayanswamy, Joseph Breda, Xin Liu, Shwetak Patel, and Vikram Iyer. 2024. Exploring and characterizing large language models for embedded system development and debugging. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. 1–9.
- [15] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025).
- [16] Zhi Guo, Walid Najjar, and Betul Buyukkurt. 2008. Efficient hardware code generation for FPGAs. *ACM Transactions on Architecture and Code Optimization (TACO)* 5, 1 (2008), 1–26.
- [17] Xueyu Hu, Kun Kuang, Jiankai Sun, Hongxia Yang, and Fei Wu. 2024. Leveraging print debugging to improve code generation in large language models. *arXiv preprint arXiv:2401.05319* (2024).
- [18] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2023. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770* (2023).
- [19] Manuel Jiménez, Rogelio Palomera, and Isidoro Couvertier. 2013. *Introduction to embedded systems*. Springer.
- [20] Philip Koopman, Howie Choset, Rajeev Gandhi, Bruce Krogh, Diana Marculescu, Priya Narasimhan, Joann M Paul, Ragunathan Rajkumar, Daniel Siewiorek, Asim Smailagic, et al. 2005. Undergraduate embedded system education at Carnegie Mellon. *ACM Transactions on Embedded Computing Systems (TECS)* 4, 3 (2005), 500–528.
- [21] Sumith Kulal, Panupong Pasupat, Kartik Chandra, Mina Lee, Oded Padon, Alex Aiken, and Percy S Liang. 2019. Spoc: Search-based pseudocode to code. *Advances in Neural Information Processing Systems* 32 (2019).
- [22] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2024. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances in Neural Information Processing Systems* 36 (2024).
- [23] Mingjie Liu, Nathaniel Pinckney, Brucek Khailany, and Haoxing Ren. 2023. Verilogeval: Evaluating large language models for verilog code generation. In *2023 IEEE/ACM International Conference on Computer Aided Design (ICCAD)*. IEEE, 1–8.
- [24] Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtiar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov, Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul, Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii, Nii Osae Osae Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muh-tasham Oblokulov, Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, Julian McAuley, Han Hu, Torsten Scholak, Sébastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostafa Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. 2024. StarCoder

- 2 and The Stack v2: The Next Generation. arXiv:2402.19173 [cs.SE]
- [25] Marc Marone and Benjamin Van Durme. 2023. DataPortraits. <https://dataportraits.org/>
- [26] Tomas G Moreira, Marco A Wehrmeister, Carlos E Pereira, Jean-Francois Petin, and Eric Levrat. 2010. Automatic code generation for embedded systems: From UML specifications to VHDL code. In *2010 8th IEEE International Conference on Industrial Informatics*. IEEE, 1085–1090.
- [27] Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. 2025. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393* (2025).
- [28] OpenAI. 2025. OpenAI o3-mini. <https://openai.com/index/openai-o3-mini/>
- [29] Sudeep Pasricha. 2022. Embedded systems education in the 2020s: Challenges, reflections, and future directions. In *Proceedings of the Great Lakes Symposium on VLSI 2022*. 519–524.
- [30] Abhinav Podili, Chi Zhang, and Viktor Prasanna. 2017. Fast and efficient implementation of convolutional neural networks on FPGAs. In *2017 IEEE 28Th international conference on application-specific systems, architectures and processors (ASAP)*. IEEE, 11–18.
- [31] Shanghaoran Quan, Jiaxi Yang, Bowen Yu, Bo Zheng, Dayiheng Liu, An Yang, Xuancheng Ren, Bofei Gao, Yibo Miao, Yunlong Feng, et al. 2025. CodeElo: Benchmarking Competition-level Code Generation of LLMs with Human-comparable Elo Ratings. *arXiv preprint arXiv:2501.01257* (2025).
- [32] Qwen Team. 2025. QwQ-32B: Embracing the Power of Reinforcement Learning. <https://qwenlm.github.io/blog/qwq-32b/>
- [33] Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. MuSiQue: Multihop Questions via Single-hop Question Composition. *Transactions of the Association for Computational Linguistics* 10 (2022), 539–554.
- [34] Tianwen Wei, Jian Luan, Wei Liu, Shuang Dong, and Bin Wang. 2023. Cmath: Can your language model pass chinese elementary school math test? *arXiv preprint arXiv:2306.16636* (2023).
- [35] Wokwi. 2019. Wokwi: World’s most advanced ESP32 simulator. <https://wokwi.com/>
- [36] Wayne Wolf. 2010. *High-performance embedded computing: architectures, applications, and methodologies*. Elsevier.
- [37] Chunqiu Steven Xia and Lingming Zhang. 2023. Keep the Conversation Going: Fixing 162 out of 337 bugs for \$0.42 each using ChatGPT. *arXiv preprint arXiv:2304.00385* (2023).
- [38] Ruiyang Xu, Jialun Cao, Yaojie Lu, Hongyu Lin, Xianpei Han, Ben He, Shing-Chi Cheung, and Le Sun. 2024. CRUXEval-X: A Benchmark for Multilingual Code Reasoning, Understanding and Execution. *arXiv preprint arXiv:2408.13001* (2024).
- [39] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2024. Qwen2.5 Technical Report. *arXiv preprint arXiv:2412.15115* (2024).
- [40] Aidan ZH Yang, Claire Le Goues, Ruben Martins, and Vincent Hellen-dorff. 2024. Large language models for test-free fault localization. In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*. 1–12.
- [41] Huanqi Yang, Mingzhe Li, Mingda Han, Zhenjiang Li, and Weitao Xu. 2024. EmbedGenius: Towards Automated Software Development for Generic Embedded IoT Systems. *arXiv preprint arXiv:2412.09058* (2024).
- [42] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600* (2018).
- [43] Li Zhong, Zilong Wang, and Jingbo Shang. 2024. Ldb: A large language model debugger via verifying runtime execution step-by-step. *arXiv preprint arXiv:2402.16906* (2024).
- [44] Qiming Zhu, Jialun Cao, Yaojie Lu, Hongyu Lin, Xianpei Han, Le Sun, and Shing-Chi Cheung. 2025. Domaineval: An auto-constructed benchmark for multi-domain code generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 39. 26148–26156.

A Hardware Descriptions

This section provides an overview of the basic hardware components used in our benchmark, including their functions and pin configurations.

LED

A Light Emitting Diode (LED) is a basic component used for visual indication

Pins:

- A – Anode (positive terminal)
- C – Cathode (negative terminal)

Pushbutton

A pushbutton is a momentary switch used for digital input.

Pins:

- 1.l, 1.r – First contact terminals
- 2.l, 2.r – Second contact terminals

RGB LED

An RGB LED combines red, green, and blue LEDs in a single package to produce different colors.

Pins:

- R, G, B – Controls red, green, and blue channels respectively
- COM – Common cathode (connected to ground)

Seven-Segment Display

A seven-segment display is used for numerical output.

Pins:

- A–G – Individual segment controls
- DP – Decimal point
- COM.1, COM.2 – Common anode pins

LED Bar Graph

An LED bar graph displays signal levels or statuses with multiple LEDs in a row.

Pins:

- A1–A10 – Anodes of each LED
- C1–C10 – Cathodes of each LED

Slide Potentiometer

A slide potentiometer provides analog output based on position.

Pins:

- GND – Ground
- SIG – Analog output
- VCC – Supply voltage

Attribute: The output value ranges from 0 to 1023.

Resistor

A resistor limits current flow in a circuit.

Pins:

- 1, 2 – Connection terminals

Attribute: Resistance value: $1000\Omega(1k\Omega)$

74HC595 Shift Register

A shift register expands the number of digital outputs using serial data input.

Pins:

- DS – Serial data input
- SHCP – Shift register clock input
- STCP – Storage register clock input (latch)
- OE – Output enable (active low)
- Q0–Q7 – Parallel outputs
- MR – Master reset (active low)
- GND – Ground
- VCC – Supply voltage (typically 5V)

Servo Motor

A servo motor allows precise angular positioning.

Pins:

- PWM – Control signal
- V+ – Power supply (5V)
- GND – Ground

B Platform Descriptions

This section outlines the hardware platforms used in our benchmark, including general-purpose input/output (GPIO) pins and power configurations.

Raspberry Pi Pico

A microcontroller board based on the RP2040 chip.

Pins:

- GP0–GP22 – Digital GPIO pins
- GP26–GP28 – Analog input pins
- GND.1–GND.8 – Ground connections
- VSYS, VBUS, 3V3 – Power supply inputs and outputs

Arduino Uno Rev3

A popular microcontroller board based on the ATmega328P.

Pins:

- 0–13 – Digital GPIO pins
- A0–A5 – Analog input pins
- GND.1–GND.3 – Ground pins
- VIN – External power input (7–12V)
- 5V – 5V power output

Arduino Mega

A microcontroller board with expanded I/O capabilities, based on the ATmega2560.

Pins:

- 0–53 – Digital GPIO pins
- A0–A15 – Analog input pins
- GND.1–GND.5 – Ground pins

- VIN – External power input (7–12V)
- 5V – 5V power output

Note: Pins 2–13, 44–46 support hardware PWM output (15 channels total).

ESP32-DevKitC V4

A development board based on the ESP32 microcontroller with integrated Wi-Fi and Bluetooth.

Pins:

- 0, 2, 4, 5, 12–19, 21–23, 25–27, 32, 33, 34, 35 – GPIO pins
- GND.1–GND.3 – Ground pins
- 3V3 – 3.3V power output

C Data Annotation and Quality Control

Human Annotation

In this section, we provide the detailed process for human annotation of our benchmark.

Initially, we recruited two bachelor’s degree holders as annotators. After undergoing rigorous training, each annotator independently performed the annotations. To ensure quality control, an additional pair of bachelor’s degree holders reviewed all annotated content. The guidelines for our annotation primarily cover the following aspects:

(1) **Optimal Hardware Resource Utilization Principle:** the problem description must ensure the complete and efficient utilization of all available hardware resources.

(2) **State Duration Stability Specification:** with the exception of instantaneous interaction scenarios (e.g., button-triggered interactions), all device states must maintain a stable duration of no less than two seconds.

(3) **The Test Validity Criteria:** test cases must satisfy completeness criteria, meaning they should comprehensively verify the functional characteristics and correctness of the target problem. Each test case should be designed with a specific validation objective in mind.

These guidelines aim to ensure the quality and effectiveness of the test case repository while enhancing the stability and reproducibility of the testing process.

Quality Control

To ensure the correctness of the benchmark, each example was independently annotated by two annotators to minimize subjective bias. Upon completing their annotations, the two annotators cross-verified each other’s results. Once both annotators reached a consensus, the annotated content was reviewed by two additional volunteers for further validation.

Volunteers meticulously examined the annotations and corrected errors. Through discussion, they selected the most accurate versions. The finalized annotations include the problem statements, corresponding solutions, and test cases. During the review phase, volunteers adhered to the following standardized quality assessment protocol:

(1) **Problem Description Evaluation:** Assess the completeness and clarity of the problem formulation. Problem Description Evaluation: The problem formulation must be complete and clear, with unambiguous requirements covering all specifications.

(2) **Test Case Validation:** Verify the logical soundness of the test case design, ensuring alignment with the problem’s functional and edge-case requirements. Evaluate whether the test suite provides sufficient validation for algorithmic correctness and robustness.

(3) **Stability Testing:** Perform a minimum of five independent test executions under controlled conditions. We mitigated external interference (e.g., computational resource contention) to ensure reproducible results. Document response times and output consistency for all runs to assess system reliability.

D Thinking Process of 7-Segment Display

Part of Thinking Process in 7-Segment Display

...
 For digit 1, the bits would be:
 Bit7 (Q7): 0 Bit6 (Q6): 1 Bit5 (Q5): 1 Bit4 (Q4): 1 Bit3 (Q3): 1 Bit2 (Q2): 0 Bit1 (Q1): 0 Bit0 (Q0): 1
 So binary 0 1 1 1 1 0 0 1 → which is 0b01111001 (0x79). Wait, but that would set Q0 to 1 (A off), Q1 to 0 (B on), Q2 to 0 (C on), and the rest (Q3-6 to 1, Q7 0). So segments B and C are on. That's correct. So the byte for digit 1 is 0x79.
This process is tedious, but necessary. Let's do this for each digit:
 Digit 0: 0x40 (binary 01000000) → A-F on, G off.
 Digit 1: 0x79 (0b01111001) → B and C on.
 Digit 2: Segments A, B, G, E, D → A (Q0=0), B (Q1=0), G (Q6=0), E (Q4=0), D (Q3=0). So which bits are 0? Q0, Q1, Q3, Q4, Q6. So the byte bits:
 Q7:0 Q6:0 (G on) Q5:1 (F off) Q4:0 (E on) Q3:0 (D on) Q2:1 (C off)
 Q1:0 (B on) Q0:0 (A on)
 So bits: 0 0 1 0 0 1 0 0 → binary 00100100 → 0x24. Wait, let's write it as 8 bits: 0 0 1 0 0 1 0 0 → that's 0x24 (36). Wait:
 Q7:0 → bit7:0 Q6:0 → bit6:0 Q5:1 → bit5:1 Q4:0 → bit4:0 Q3:0 → bit3:0 (Wait, D is on, so Q3 is 0 → bit3 is 0) Wait, no. Wait, for digit 2, D is on (Q3=0). So Q3 is 0 → bit3 is 0. Let me re-calculate: Segments for digit 2: A, B, G, E, D → A (Q0=0), B (Q1=0), G (Q6=0), E (Q4=0), D (Q3=0). Other segments (C and F) are off → Q2=1, Q5=1.
 So the bits:
 Q7:0 Q6:0 Q5:1 Q4:0 Q3:0 Q2:1 Q1:0 Q0:0
 Which is 0 0 1 0 0 1 0 0 → binary 00100100 → 0x24. So yes, 0x24 for digit 2.
 Hmm, **this is time-consuming**. Maybe there's a better way to approach this. Alternatively, perhaps we can create a lookup table using a helper function that for each digit, defines which segments are on (LOW), and then build the byte accordingly. Alternatively, perhaps using a standard 7-segment code table for common anode displays. For example, the standard byte codes for a common anode 7-segment display connected to a shift register. Wait, the standard segment codes for common anode (assuming segments a to g are connected to Q0 to Q6):
 For example, the standard code for 0 is 0x03 (binary 00000011) but that depends on the wiring. Wait, perhaps this approach is error-prone. Let me look for a standard 7-segment code table. Upon checking, the standard byte patterns for a common anode 7-segment display (assuming segments are a, b, c, d, e, f, g, dp connected to Q0-Q7), but in our case, the segments are connected as:
 Q0 → A Q1 → B Q2 → C Q3 → D Q4 → E Q5 → F Q6 → G
 So each segment corresponds to a specific Q output. The standard codes for the digits would then need to be mapped to these Q outputs. Let's list the segments for each digit and then map to the corresponding bits.
 ...

E Prompt and Example of R1-Retrieval

R1-Retrieval Prompt

{ prompt same as "w/ schematics" or "w/o schematics" }

Retrieve Code

Here are some retrieved code snippets that may be helpful for solving the problem:

LED

```

1 // Define the pin where the LED is connected
2 const int ledPin = 2;
3
4 void setup() {
5     // Initialize the LED pin as an output
6     pinMode(ledPin, OUTPUT);
7     // Ensure the LED is off initially
8     digitalWrite(ledPin, LOW);
9 }
10
11 void loop() {
12     // Turn the LED on for 1 second
13     digitalWrite(ledPin, HIGH);
14     delay(1000);
15
16     // Turn the LED off for 2 seconds
17     digitalWrite(ledPin, LOW);
18     delay(2000);
19
20     // Turn the LED on for 3 seconds
21     digitalWrite(ledPin, HIGH);
22     delay(3000);
23
24     // Turn the LED off for 3 seconds
25     digitalWrite(ledPin, LOW);
26     delay(3000);
27
28     // The sequence will repeat indefinitely
29     // due to the loop() function
30 }
```

PushButton

```

1 // Define the pins for the RGB LED and the
2 // push button
3 const int redPin = 2;      // Red pin of the
4 // RGB LED
5 const int greenPin = 3;    // Green pin of the
6 // RGB LED
7 const int bluePin = 4;     // Blue pin of the
8 // RGB LED
9 const int buttonPin = 5;   // Push button pin
10
11 // Variables to manage the button state and
12 // debouncing
```

```

8 int buttonState = HIGH;      // Current state
   of the button
9 int lastButtonState = HIGH; // Previous
   state of the button
10 unsigned long lastDebounceTime = 0; // Timestamp of the last button state
   change
11 const unsigned long debounceDelay = 50; // Debounce delay in milliseconds
12
13 // Variable to track the current color state
14 int colorState = 0; // 0: Off, 1: Red, 2:
   Green, 3: Blue
15
16 void setup() {
17   // Set the RGB LED pins as OUTPUT
18   pinMode(redPin, OUTPUT);
19   pinMode(greenPin, OUTPUT);
20   pinMode(bluePin, OUTPUT);
21
22   // Set the button pin as INPUT_PULLUP
23   pinMode(buttonPin, INPUT_PULLUP);
24
25   // Initialize the RGB LED to be off
26   digitalWrite(redPin, LOW);
27   digitalWrite(greenPin, LOW);
28   digitalWrite(bluePin, LOW);
29 }
30
31 void loop() {
32   // Read the current state of the button
33   int reading = digitalRead(buttonPin);
34
35   // Check if the button state has changed (due to noise or pressing)
36   if (reading != lastButtonState) {
37     // Reset the debounce timer
38     lastDebounceTime = millis();
39   }
40
41   // Check if the debounce delay has passed
42   if ((millis() - lastDebounceTime) >
        debounceDelay) {
43     // If the button state has changed
44     if (reading != buttonState) {
45       buttonState = reading;
46
47       // If the button is pressed (LOW state)
48       if (buttonState == LOW) {
49         // Cycle through the color states
50         colorState = (colorState + 1) % 4;
51
52         // Update the RGB LED based on the
           current color state
53         switch (colorState) {
54           case 0: // Off
55             digitalWrite(redPin, LOW);

```

```

56             digitalWrite(greenPin, LOW);
57             digitalWrite(bluePin, LOW);
58             break;
59           case 1: // Red
60             digitalWrite(redPin, HIGH);
61             digitalWrite(greenPin, LOW);
62             digitalWrite(bluePin, LOW);
63             break;
64           case 2: // Green
65             digitalWrite(redPin, LOW);
66             digitalWrite(greenPin, HIGH);
67             digitalWrite(bluePin, LOW);
68             break;
69           case 3: // Blue
70             digitalWrite(redPin, LOW);
71             digitalWrite(greenPin, LOW);
72             digitalWrite(bluePin, HIGH);
73             break;
74         }
75     }
76   }
77 }
78
79 // Save the current reading as the last
   button state
80 lastButtonState = reading;
81 }

```

RGB LED

```

1
2   // Define the pins for the RGB LED
3   const int redPin = 2;    // Red LED connected
   to pin 2
4   const int greenPin = 3; // Green LED
   connected to pin 3
5   const int bluePin = 4;  // Blue LED
   connected to pin 4
6
7   void setup() {
8     // Set the RGB LED pins as output
9     pinMode(redPin, OUTPUT);
10    pinMode(greenPin, OUTPUT);
11    pinMode(bluePin, OUTPUT);
12
13    // Initially turn off the RGB LED
14    digitalWrite(redPin, LOW);
15    digitalWrite(greenPin, LOW);
16    digitalWrite(bluePin, LOW);
17  }
18
19  void loop() {
20    // Red for 3 seconds
21    digitalWrite(redPin, HIGH);
22    digitalWrite(greenPin, LOW);
23    digitalWrite(bluePin, LOW);
24    delay(3000);

```

```

25
26 // Yellow (Red + Green) for 3 seconds
27 digitalWrite(redPin, HIGH);
28 digitalWrite(greenPin, HIGH);
29 digitalWrite(bluePin, LOW);
30 delay(3000);
31
32 // White (Red + Green + Blue) for 3 seconds
33 digitalWrite(redPin, HIGH);
34 digitalWrite(greenPin, HIGH);
35 digitalWrite(bluePin, HIGH);
36 delay(3000);
37
38 // Turn off the LED for 3 seconds
39 digitalWrite(redPin, LOW);
40 digitalWrite(greenPin, LOW);
41 digitalWrite(bluePin, LOW);
42 delay(3000);
43 }
```

Slide Potentiometer

```

1
2 // Define the pins for the RGB LED
3 const int redPin = 2;      // Red LED connected
   to digital pin 2
4 const int greenPin = 3;    // Green LED
   connected to digital pin 3
5 const int bluePin = 4;     // Blue LED
   connected to digital pin 4
6
7 // Define the pin for the slide potentiometer
8 const int potPin = A3;    // Slide
   potentiometer connected to analog pin A3
9
10 void setup() {
11   // Set the RGB LED pins as output
12   pinMode(redPin, OUTPUT);
13   pinMode(greenPin, OUTPUT);
14   pinMode(bluePin, OUTPUT);
15
16   // Initialize the RGB LED to be off
17   digitalWrite(redPin, LOW);
18   digitalWrite(greenPin, LOW);
19   digitalWrite(bluePin, LOW);
20 }
21
22 void loop() {
23   // Read the value from the slide
   potentiometer (0 to 1023)
24   int potValue = analogRead(potPin);
25
26   // Determine the color based on the
   potentiometer value
27   if (potValue < 341) {
28     // First third of the range: Red
     digitalWrite(redPin, HIGH);
```

```

30     digitalWrite(greenPin, LOW);
31     digitalWrite(bluePin, LOW);
32 } else if (potValue >= 341 && potValue <
   682) {
33   // Second third of the range: Green
34   digitalWrite(redPin, LOW);
35   digitalWrite(greenPin, HIGH);
36   digitalWrite(bluePin, LOW);
37 } else {
38   // Final third of the range: Blue
39   digitalWrite(redPin, LOW);
40   digitalWrite(greenPin, LOW);
41   digitalWrite(bluePin, HIGH);
42 }
43 }
```

7 Segment Display

```

1 #include <Servo.h>
2
3 // Pin definitions
4 const int dataPin = A0;      // DS pin of 74
   HC595
5 const int latchPin = A1;    // STCP pin of 74
   HC595
6 const int clockPin = A2;    // SHCP pin of 74
   HC595
7
8 // 7-segment display digit patterns (common
   anode)
9 const byte digitPatterns[10] = {
10   B11000000, // 0
11   B1111001, // 1
12   B10100100, // 2
13   B10110000, // 3
14   B10011001, // 4
15   B10010010, // 5
16   B10000010, // 6
17   B11111000, // 7
18   B10000000, // 8
19   B10010000 // 9
20 };
21
22 unsigned long previousMillis = 0; // Stores
   the last time the timer was updated
23 const long interval = 1000;        // Interval
   at which to update the timer (1 second)
24 int timer = 0;                  // Timer
   value
25
26 void setup() {
27   // Set the pin modes
28   pinMode(dataPin, OUTPUT);
29   pinMode(latchPin, OUTPUT);
30   pinMode(clockPin, OUTPUT);
31 }
```

```

32 // Initialize the 7-segment display to be
33 // off
34 digitalWrite(latchPin, LOW);
35 shiftOut(dataPin, clockPin, MSBFIRST,
36         B11111111); // Turn off all segments
37 digitalWrite(latchPin, HIGH);
38 }
39
40 void loop() {
41     unsigned long currentMillis = millis();
42
43     // Check if 1 second has passed
44     if (currentMillis - previousMillis >=
45         interval) {
46         previousMillis = currentMillis;
47
48         // Increment the timer by 2 seconds
49         timer += 2;
50
51         // Determine what to display on the 7-
52         // segment display
53         int displayValue;
54         if (timer > 10) {
55             displayValue = timer % 10; // Display
56             // only the last digit
57         } else {
58             displayValue = timer; // Display the
59             // full timer value
60         }
61     }
62 }

Servo

1 #include <Servo.h>
2
3 // Define the pins
4 const int btn1Pin = 2; // Button connected
5 // to digital pin 2
6 const int servo1Pin = 3; // Servo1 connected
7 // to digital pin 3
8 const int servo2Pin = 4; // Servo2 connected
9 // to digital pin 4
10
11 // Create servo objects
12 Servo servo1;
13 Servo servo2;
14
15 // Variables to store the servo positions
16 int servo1Pos = 0;
```

```

14 int servo2Pos = 180;
15
16 // Variables for button debouncing
17 int buttonState = HIGH;
18 int lastButtonState = HIGH;
19 unsigned long lastDebounceTime = 0;
20 const unsigned long debounceDelay = 50; // //
21 // Debounce delay in milliseconds
22
23 void setup() {
24     // Initialize the button pin as input with
25     // pull-up resistor
26     pinMode(btn1Pin, INPUT_PULLUP);
27
28     // Attach the servos to their respective
29     // pins
30     servo1.attach(servo1Pin);
31     servo2.attach(servo2Pin);
32
33     // Set initial positions of the servos
34     servo1.write(servo1Pos);
35     servo2.write(servo2Pos);
36 }
37
38 void loop() {
39     // Read the current state of the button
40     int reading = digitalRead(btn1Pin);
41
42     // Check if the button state has changed (
43     // due to noise or pressing)
44     if (reading != lastButtonState) {
45         // Reset the debouncing timer
46         lastDebounceTime = millis();
47     }
48
49     // Check if the debounce delay has passed
50     if ((millis() - lastDebounceTime) >
51         debounceDelay) {
52         // If the button state has changed
53         if (reading != buttonState) {
54             buttonState = reading;
55
56             // Only toggle the servo positions if
57             // the button is pressed (LOW)
58             if (buttonState == LOW) {
59                 // Toggle the servo positions
60                 if (servo1Pos == 0) {
61                     servo1Pos = 180;
62                     servo2Pos = 0;
63                 } else {
64                     servo1Pos = 0;
65                     servo2Pos = 180;
66                 }
67
68                 // Move the servos to their new
69                 // positions
70                 servo1.write(servo1Pos);
71             }
72         }
73     }
74 }
```

```

64         servo2.write(servo2Pos);
65     }
66   }
67 }
68
69 // Save the current reading as the last
// state for the next loop
70 lastButtonState = reading;
71 }

LED Bar Graph

1
2 // Define the pins for bargraph1 and
// bargraph2
3 int bargraph1Pins[] = {2, 3, 4, 5, 6, 7, 8,
9, 10, 11};
4 int bargraph2Pins[] = {3, 4, 5, 6, 7, 8, 9,
10, 11, 12};
5
6 void setup() {
7   // Initialize all LED pins as OUTPUT
8   for (int i = 0; i < 10; i++) {
9     pinMode(bargraph1Pins[i], OUTPUT);
10    pinMode(bargraph2Pins[i], OUTPUT);
11  }
12
13 // Turn off all LEDs initially
14 for (int i = 0; i < 10; i++) {
15   digitalWrite(bargraph1Pins[i], LOW);
16   digitalWrite(bargraph2Pins[i], LOW);
17 }
18 }
19 void loop() {
20   // Left to right sequence
21   for (int i = 0; i < 10; i++) {
22     // Turn on the current LED in bargraph1
23     digitalWrite(bargraph1Pins[i], HIGH);
24     // Turn on the previous LED in bargraph2
     // (if not the first LED)
25     if (i > 0) {
26       digitalWrite(bargraph2Pins[i - 1], HIGH
         );
27     }
28     // Turn off the LED in bargraph2 if
     // bargraph1 is at the first LED
29     if (i == 0) {
30       digitalWrite(bargraph2Pins[9], LOW);
31     }
32     // Wait for 2 seconds
33     delay(2000);
34     // Turn off the current LED in bargraph1
35     digitalWrite(bargraph1Pins[i], LOW);
36     // Turn off the previous LED in bargraph2
     // (if not the first LED)
37     if (i > 0) {

```

```

38       digitalWrite(bargraph2Pins[i - 1], LOW)
         ;
39     }
40   }
41   // Right to left sequence
42   for (int i = 9; i >= 0; i--) {
43     // Turn on the current LED in bargraph1
44     digitalWrite(bargraph1Pins[i], HIGH);
45     // Turn on the previous LED in bargraph2
     // (if not the first LED)
46     if (i > 0) {
47       digitalWrite(bargraph2Pins[i - 1], HIGH
         );
48     }
49     // Turn off the LED in bargraph2 if
     // bargraph1 is at the first LED
50     if (i == 0) {
51       digitalWrite(bargraph2Pins[9], LOW);
52     }
53     // Wait for 2 seconds
54     delay(2000);
55     // Turn off the current LED in bargraph1
56     digitalWrite(bargraph1Pins[i], LOW);
57
58     // Turn off the previous LED in bargraph2
     // (if not the first LED)
59     if (i > 0) {
60       digitalWrite(bargraph2Pins[i - 1], LOW)
         ;
61     }
62   }
63 }


```

F Prompt of R1-Compiler

R1-Compiler Prompt

The code is build failed because of the following error:

{the error message from compiler}

Please fix the error

Note:

1. Your connections should in [CONNECTIONS] and [/CONNECTIONS]
2. Your code should in [ESP32 Code] and [/ESP32 Code]
3. Even if the origin connection or code is correct, you should still rewrite it.
4. First give the connection and then give the code.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009