

REFINING SPECS FOR LLM-BASED RTL AGILE DESIGN

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) are increasingly employed to assist in agile register-transfer-level (RTL) hardware design. This is a labor-intensive stage in developing FPGA-based acceleration services or prototyping ASICs, and successful automation can largely shorten the development cycle. However, benchmarks are reporting a relatively low functional correctness rate (sometimes called accuracy) when generating simple modules of less than 100 lines of code (LOC, in Verilog), questioning the practicality of current LLMs for real-world designs. This paper highlights that the low accuracy is attributed to the use of low-quality descriptions as prompts in both training datasets and benchmarks. First, the natural language descriptions (NLDs) do not contain all the semantics constrained by the testbenches (TB), causing false negatives during verification. Second, existing automatically generated NLDs are usually too detailed in implementation, which is not suitable for both training and benchmarking. We designed tools to quantify the clarity and simplicity of the cases, improve the quality of existing and future LLM-for-RTL datasets, and assist agile RTL designers in creating qualified specifications (specs, i.e., formatted and complete NLDs). We show by experiment that LLMs can create specs with high quality at a low cost. Additionally, when equipped with these specs, general-purpose LLMs can achieve a high pass@5 rate (up to 89% on RTLLM, 96% on VerilogEval-Human) without requiring expensive fine-tuning or post-generation self-fixing.

1 INTRODUCTION

In the recent past, large language models (LLMs) have emerged as general-purpose design assistants, and their capabilities have been examined in various aspects of software design, including code generation, debugging, and refactoring (Schmid et al., 2025; Esposito et al., 2025). Following this idea, researchers have begun to explore their potential in automatic, agile hardware design, particularly in register-transfer-level (RTL) design and verification, which facilitates the rapid development of FPGA accelerators or ASIC prototypes. The use case is usually defined as follows. The user prepares a natural language description (NLD) of the to-be-designed hardware, and a verification method such as an executable testbench (TB) or a formal verifier (FV). The LLMs are instructed to generate an RTL design according to the NLD, referred to as the device under test (DUT), that should pass the verification over the TB or the FV.

Benchmarks have been proposed to evaluate the quality of designs, especially the probability of correctness (Lu et al., 2024; Liu et al., 2023; Jin et al., 2025; Purini et al., 2025). Unfortunately, the results are not satisfactory. According to the OpenLLM paper (Liu et al., 2024), GPT-4, the leading general-purpose model of the time, managed only 55.8% pass@5 on the VerilogEval-Human benchmark (Liu et al., 2023), and 65.5% on the RTLLM benchmark (Lu et al., 2024). The tasks in these benchmarks are mostly simple components, so these low pass@5 scores suggest that LLMs are not yet practical for generating real-world designs directly.

Improvements are proposed for both the models and the generation strategies, i.e., the prompts and workflows. Works on model fine-tuning focus on how to automatically collect large-scale datasets and create better NLDs for training (Zhu et al., 2025; Zhang et al., 2024). [Code segments from real-world projects are serving as golden models \(GMs\), and NLDs and TBs are created from them.](#) Meanwhile, works on generation strategies focus on automatically correcting errors after an initial

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

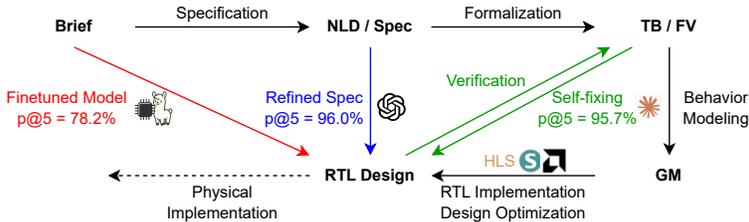


Figure 1: The workflow of traditional and agile RTL design. Black arrows mark the regular design workflow. Red (Li et al., 2025), blue (our work), and green (Zhao et al., 2024b) arrows mark the LLM-based agile design workflow. The percentages are the best reported pass rate on the VerilogEval-Human benchmark.

attempt (Li et al., 2025; Zhao et al., 2024b). Error reports created by simulators are fed back to the LLMs, and a correct design can be sampled within a few iterations. However, these efforts still have shortcomings: the accuracy improvement is limited (especially when fine-tuning from small models), scoring high only on line-by-line translation benchmarks (e.g., VerilogEval-Machine), or consuming several times more computing power.

Interestingly, through manual inspection of existing datasets and benchmarks, we found that many failures do not originate from inherent limitations of current LLMs. Instead, a substantial portion of the failing cases should be reinterpreted as false negatives (FNs), caused by two recurrent deficiencies in the input NLDs: (D1) insufficient or misleading semantic specifications, and (D2) overly complicated implementation-level instructions. Case studies and qualitative analyses are presented in Section 3. After correcting these deficiencies, LLMs are often able to generate functionally correct designs. This raises an important question: to what extent have existing benchmarks underestimated the true power of LLMs in RTL generation?

This paper proposes a paradigm for writing specifications (specs) in a fine format with clarity and simplicity, to prompt the LLM-based RTL design agents. Similar to the traditional workflow, human designers should provide a precise definition of functionalities, sequential behaviors, and semantics of the parameters, ports, and I/O signals before moving to the implementation phase. Since these documents must be produced (as user guides) regardless of whether LLMs design agents are involved, a minimal extra human workload is introduced. In fact, they can also be generated by LLMs when properly prompted. Given the improved specs, RTL generation accuracy can increase by approximately 10-20 percentage points. For example, GPT-4 achieves pass@5 rates of 89.0%, 96.0%, and 99.5% on RTLLM, VerilogEval-human, and VerilogEval-machine, respectively.

The main contributions of this work are:

- We demonstrated the importance of a clear yet simple spec for the accuracy of RTL generating LLMs. We designed templates and automated tools to create such specs, and improved existing benchmarks by replacing the original NLDs with the specs we created.
- We examined that simply using these specs can help LLMs achieve a score similar to (for weaker models) or much better than (for powerful models) existing works. This even holds when the spec is generated by an LLM other than the one that generates the RTL code.
- We show that enhancing the spec quality can also increase the accuracy of generating complex designs, such as FFT and Conv2d, which were previously reported to be challenging.

2 BACKGROUNDS

2.1 RTL DESIGN WORKFLOW

As illustrated in Figure 1, the traditional RTL design workflow can be divided into stages, each producing a solid outcome that is validated before proceeding to the next stage. These intermediate outcomes, including the specs, TBs / FVs, and GMs, are not only essential for subsequent stages but also serve as user guides and tools incorporated into the released product.

The brief is a short statement of what the module is designed for, typically appearing in the first paragraph of a user manual. The spec is a human-friendly yet complete semantic description of the design. It should be able to guide a user to use the design as a black box, without awareness of the internal details. The TB and the FV are machine-friendly versions of the spec. A design is functionally correct only if it passes verification. Of course, passing does not always imply correctness. But since orthogonal works have attempted to mitigate these false positives (FPs) by writing better TBs (Jin et al., 2025), we will not further discuss the FPs in our paper. The GM is a straightforward implementation, without performance or cost considerations, but simple enough to guarantee its correctness. The final design should be semantically equivalent to the GM, so that a system built using the GM can run as-is if the GM is replaced by it.

2.2 LLM-BASED AGILE DESIGN

LLMs have been adopted to assist human agile designers to leap from each stage directly to the last. Intuitively, the earlier the stage at which the LLMs start, the less human effort is required, but also the less information is provided, and the greater the difficulty for the LLM. Among these, the last stage (RTL implementation) is typically handled by high-level synthesis (HLS) tools, which lie outside the scope of this paper.

Metrics and Benchmarks. The accuracy of the designs is typically measured using the pass@ k metric (Chen et al., 2021), which indicates the probability of obtaining at least one passing candidate among k random samples. Each test case should be equipped with either a GM (for fuzzing-based differential tests), a TB (for constructed tests), or an FV (for formal verification over constraints) to determine whether a DUT passes. Although earlier benchmarks relied more on simulation-based TBs, FVs are increasingly employed (Fang et al., 2024; Jin et al., 2025).

RTLLM (Liu et al., 2024) and VerilogEval (Liu et al., 2023) are currently the most recognized benchmarks in this scenario. The latest version of RTLLM (-v2) contains 50 basic industrial function units, including arithmetic units, buffers, FSMs, and signal processing units such as edge detectors or parallel-serial converters. VerilogEval contains 156 crafted tasks collected from beginner exercises: implementing very simple logic with fewer than 10 lines of code (LOC), or translating truth tables or state graphs into code. Two versions of NLDs are provided in different styles. The -v2 (and the previous -Human) version has manually written high-level semantic descriptions, while the -Machine version has LLM-generated GM code summaries.

Recently, benchmarks for generating complex designs or large systems have also been proposed. The authors of AutoSilicon (Li et al., 2025) added generation tasks for FFT, I/O ports, and RISC-V CPU, using the same format as RTLLM. ArchXBench (Purini et al., 2025) provides commonly used functional primitives, such as signal filters and cryptographic algorithms, along with their floating-point and pipelined variants. RealBench (Jin et al., 2025) provides structured design tasks with high-quality manual written specs. Until the submission of this article, these harder tasks have never been solved by any LLM design agent.

The naive approach. The earliest attempts directly fed the briefs into general-purpose LLMs or coding models trained with software codes. Minimal human effort is required, but the pass rate is low: only 20%-40% of samples are correct on simple tasks (Thakur et al., 2023). While LLMs are becoming larger and more powerful, this score increased only to 50% (Lu et al., 2024). Further improvements were mainly proposed from three perspectives, as also illustrated in Figure 1.

Dataset construction and model fine-tuning. Some believe that the low accuracy is due to a lack of training materials for hardware design. So, they build large-scale datasets (Zhu et al., 2025), including cases involving new input formats (e.g., truth tables and state diagrams (Liu et al., 2025a)) or new hardware description languages (HDLs, e.g., the emerging Chisel (Zhao et al., 2024a)). Then, they fine-tune the models, either supervised or via reinforcement learning. These works usually use small backbone models (with ≤ 32 B parameters) because training larger models is too expensive. This category is exemplified by CodeV-R1 (Zhu et al., 2025).

Pre-generation prompt engineering. In parallel with fine-tuning, some works turn to teaching the LLMs through prompts. They highlight the difference between HDL and software codes, decompose states or submodules, or instruct the LLMs to do so by themselves (referred to as self-planning in the RTLLM paper (Lu et al., 2024)). This category is usually employed as part of other approaches.

162 However, if the backbone model is powerful, using it alone can also yield significant performance
 163 improvements. Our work falls into this category.

164 **Post-generation self-fixing.** In addition to improving the single-shot performance, post-generation
 165 mechanisms are proposed to fix minor errors reported during synthesis or simulation. Post-
 166 generation mechanisms can achieve a high pass rate of 95.7%, but incur significantly higher infer-
 167 ence costs. For example, Mage sampled 20 code candidates, selected the best two, and debugged
 168 each for 5 iterations (Zhao et al., 2024b). This is much more expensive than sampling 5 independent
 169 candidates, which is the practice of a regular pass@5. This category is exemplified by Mage (Zhao
 170 et al., 2024b). Models that are fine-tuned with a self-fixing corpus (given erroneous code and error
 171 report), such as CraftRTL (Liu et al., 2025a), can be regarded as merging the self-fixing step into
 172 generation, and can also be considered in this category.

174 3 THE REFINED SPEC

175 In this section, we first highlight the two major deficiencies with existing datasets and benchmarks
 176 for LLM-based RTL design tasks. Then, we outline a writing paradigm for high-quality specs and
 177 present our LLM-based tools that can help dataset creators or hardware designers to improve their
 178 specs, either automatically or interactively. We introduce the clarity and simplicity metrics and use
 179 their joint distribution to evaluate the quality of the datasets before and after refinement. Finally, we
 180 discuss related works on benchmark improvements and future directions.

183 3.1 CASE STUDY 1: THE NLD-GM MISMATCH

184 Deficiency (D1) is the semantic mismatch between the NLDs and GMs, TBs, or FVs. Consider the
 185 serial2parallel task from the RTLLM v2 benchmark (Liu et al., 2024), as shown in Figure 2. This module is normally ex-
 186 pected to have eight states, each processes one bit, to accept bitstreams continuously. But the
 187 GM creates an extra state, in which `cnt == 4'd8`, skipping every ninth bit received. While
 188 this feature is not explicitly described in the NLD, the TB requires the DUTs to have strictly
 189 consistent behavior with the GM. Therefore, it is almost impossible for the DUTs (whether
 190 written by humans or LLMs) to pass this test unless the designer is allowed to infer unwritten
 191 requirements from simulation error reports.

192 (D1) is harmful for both benchmarks and datasets. Mismatches in benchmarks lead to
 193 an underestimation of LLMs’ ability. What is worse, mismatches in datasets should be regarded as a form of mislabeling, which, in theory, re-
 194 duces the accuracy of the trained model.

205 (D1) occurs in both manually-written NLDs and LLM-generated summaries. Recall that in the tradi-
 206 tional workflow, a human designer must first present the NLD, then formalize it as a TB or FV, and
 207 finally present a semantically equivalent GM. During this process, constraints are applied to down-
 208 stream outcomes, but upstream documents are not always updated accordingly. Thus, mismatches
 209 occur. On the other hand, existing works that automatically create large-scale datasets often prompt
 210 the LLMs to “summarize” the functions of modules to generate NLDs. However, the generated
 211 summaries may omit the semantic details that are supposed to be in the context. In other words, the
 212 LLMs think they are writing an introduction instead of a complete spec.

213 As the MG-Verilog paper concluded, “although high-level global summaries are the most user-
 214 friendly data format, their ambiguity often results in a lack of detailed information necessary for
 215 precise code generation” (Zhang et al., 2024). However, their approach of creating multi-grained
 NLDs leads to another problem, as demonstrated below.

Implement a series-parallel conversion circuit. It receives a serial input signal "din_serial" along with a control signal "din_valid" indicating the validity of the input data. The module operates on the rising edge of the clock signal "clk" and uses a synchronous design. The input din_serial is a single-bit data, and when the module receives 8 input data, the output dout_parallel outputs the 8-bit data (The serial input values are sequentially placed in dout_parallel from the most significant bit to the least significant bit), and the dout_valid is set to 1.

```

always@(posedge clk or negedge rst_n)begin
  // ...
  else if(din_valid)
    cnt <= (cnt == 4'd8)?0:cnt+1'b1;
  else
    cnt <= 0;
  end
always@(posedge clk or negedge rst_n)begin
  // ...
  else if(cnt == 4'd8)begin
    dout_valid <= 1'b1;
    dout_parallel <= din_tmp;
  end
  // ...
end

```

Figure 2: The original NLD and GM of the RTLLM serial2parallel task (Lu et al., 2024). The red font marks the mismatch between the NLD and the GM.

3.2 CASE STUDY 2: OVER-SPECIFIED IMPLEMENTATION GUIDES

Deficiency (**D2**) is the over-detailed implementation guidance in automatically generated NLDs. Figure 3 shows the case #2456 of the MG-Verilog training dataset. This paper proposed a multi-grained training set, claiming that “high-level descriptions can facilitate user-friendly LLM interactions, while detailed descriptions are crucial for enabling LLMs to create complex designs” (Zhang et al., 2024). However, probably to solve (**D1**), the authors prompts the LLMs to “be very specific”, and the resulting NLDs turned out to be a line-by-line translation of the GM code without high-level semantics. Another typical example is the VerilogEval-Machine benchmark (Liu et al., 2023). We speculate that the LLMs may have misunderstood this instruction and became specific on *how* (instead of *what*) to design. Also, the excessively fine-grained block partitioning disperses the semantic context, making it difficult for the model to produce meaningful abstractions.

(**D2**) is also harmful for both benchmarks and datasets. Low-level benchmarks cannot assess LLMs’ pattern recombination skill because they focus on translating statements rather than selecting them to build functionality. Moreover, in a real-world use case, it is impractical for designers to write low-level NLDs, as this requires the same effort as manually designing the module. On the other hand, although multi-grained curriculum learning is theoretically beneficial, the actual accuracy improvement of these fine-tuned models is significant only on low-level benchmarks (Zhao et al., 2024a; Zhang et al., 2024). The excessive, unintended low-level training data may have led to overfitting.

It is worth noting that the RTL design corpus of existing LLMs is small, thus they are unfamiliar with microarchitecture-level instructions. Especially, they tend to confuse registers and wires. For this reason, retaining a small portion of low-level data in the datasets is still essential. We speculate that, after models are properly trained, implementation details will no longer be detrimental.

3.3 SPEC EVALUATION, REFINEMENT, AND CREATION

Despite these deficiencies, the need for automated NLD generation remains inevitable as the demand for large-scale benchmarks or training datasets grows. So, we designed three LLM-based tools, namely the *evaluator*, the *converter*, and the *creator*. These tools helped us to identify which aspects of the specification are essential for enabling LLMs to produce correct outputs.

The *evaluator* is prompted to score and criticize specs from a user’s perspective. It asks an LLM to measure NLDs using two metrics, clarity and simplicity, defined as follows.

- **Clarity:** Can you learn the module’s exact behavior from the spec without correcting any misunderstanding by reading the Verilog code?
- **Simplicity:** Can you learn the module’s high-level semantics easily without awareness of its internal implementation?

The *evaluator* scores each metric on a scale of 1 to 5 and provides suggestions for improvement. The prompt template and an example response are available in Appendix A. We ran the *evaluator* on the NLDs of both benchmarks and training datasets. Although these scores are produced by the LLMs and are not rigorous, they still reflect the difficulty level that the LLMs themselves experience during training and evaluation. The distributions of the scores are shown in Figure 4. Intuitively, a low

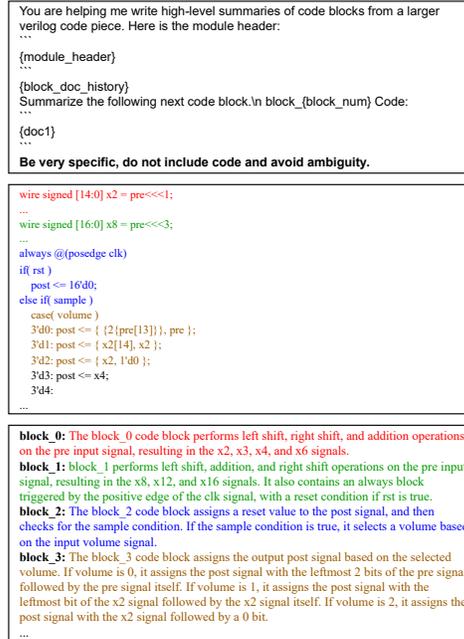


Figure 3: The prompt for block-level summary generation, the generated NLD, and the GM code of MG-Verilog training set case #2456. The colors show the mapping between code lines and description sentences.

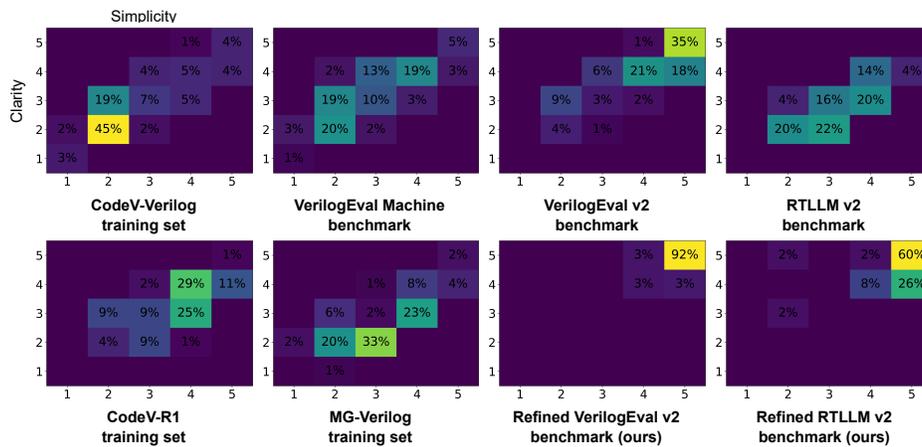


Figure 4: The joint distribution of clarity and simplicity ratings of benchmarks and datasets, scored by our *evaluator* based on GPT-4. 200 cases (or all, if insufficient) are sampled from each dataset.

clarity score means **(D1)**, and a low simplicity score means **(D2)**. We can see that among the training datasets, the NLDs of MG-Verilog and CodeV-Verilog (both using multi-grained summaries) are of relatively low quality. Meanwhile, the relatively higher quality of the CodeV-R1 dataset may explain why the fine-tuned model performs better.

From the resulting suggestions, we concluded 8 key features to create an ideal spec. First, it should include complete **black-box semantics** to prevent **(D1)**, so that if the user replaces the module with another implementation that strictly adheres to the spec, the system should work as-is without any further changes. This includes:

1. **Declaration.** Module name, design signature, and parameters.
2. **Behavior details.** Describes the functional relationship between the output and the input. For arithmetic primitives, this should include the computation formula. When LLMs fail to understand technical terms in a specific area, the spec must explain them.
3. **Sequential features.** Describes reset behaviors, I/O handshaking protocols, and whether the module is pipelined (i.e., can accept another input when processing one). Although intuitively more like implementation details, some FSM states should be regarded as semantics. For example, in a serial-to-parallel module, it is essential to specify whether a frame includes a start bit or parity bits, which introduces extra states.
4. **I/O details.** Describes the semantics and formats of each I/O port and its values, especially control signals (e.g., the function selection signal of an ALU) or complex data structures (e.g., floating point values).

Additionally, in cases where some functionalities are hard to interpret, the spec can include some explanations of the semantics, such as:

5. **Examples** (optional). Describes normal cases and corner cases. The normal cases are examples of how the module will be instantiated and interacted with. The corner cases clarify ambiguous behaviors.
6. **Implementation overview** (optional). Describes the block-level semantics and how blocks collaborate. This follows the idea of many existing works (Zhang et al., 2024; Zhao et al., 2024a; Liu et al., 2025b); however, the blocks should be partitioned by functionality (e.g., FSM states, pipeline stages, or submodules) rather than by syntax (e.g., code lines).

Meanwhile, to prevent **(D2)**, restrictions are imposed on the above parts:

7. **Limited RTL details.** The spec should avoid describing the microarchitecture directly; otherwise, the design synthesis task degenerates into a line-by-line translation problem.

This restriction should apply to benchmarks, because real-world designers will not provide these details as mentioned. However, when generating the lowest-level (i.e., RTL) cases for multi-grained datasets, this restriction can be disabled.

8. **Prevent trivial corner case behaviors.** For example, if overflows do not exist or are not expected to be handled, the spec may simply say “assume there are no overflows”.

In accordance with this paradigm, we designed the *converter* and the *creator* to refine the specs. We believe this process can be automated by LLMs, because they “excel in summarizing Verilog code rather than generating it” (Zhao et al., 2024a). [The prompts used are listed in Appendices B and C.](#)

The *converter* converts existing projects (containing code, documents, tests, etc.) into datasets. It automatically completes or corrects the NLDs using high-level semantics extracted from the GM, the TB, or the FV. [Theoretically, when creating datasets without a GM, using a TB or an FV is also fine.](#) However, in the benchmarks we are using, the TBs are too simple for the LLMs to infer the expected behavior. Figure 4 also shows the score distribution of the RTLLM and VerilogEval specs refined by the *converter*. The refined specs have higher scores on both metrics, and we will show in Section 4.2 that they really lead to higher RTL generation accuracy.

The *creator* helps real-world users to create specs for new designs from scratch (i.e., the user cannot provide any of GM, TB, or FV) interactively. It advises possible improvements to the spec based on the comments given by the *evaluator*. Conceptually, the *creator* implements RTLLM’s self-planning idea (Lu et al., 2024), guiding the LLMs to reason step by step. We further explicitly instruct it to carefully consider corner cases and sequential behaviors, and to ask users to clarify any ambiguity rather than trust its own hallucinations.

3.4 CODE GENERATION

After the specs are created, they can be fed to the LLMs for code generation. Comparing with existing workflows, we replace the NLDs with refined specs and add hints to prevent syntax errors.

As current LLMs still struggle to distinguish between combinational and sequential logic, we add prompts to strengthen LLMs’ understanding. We instructed the LLMs to “use `wire` assignments or `always @ (*)` to implement the combinational computations within each clock cycle, and add `_next` suffix to these signal names”. This increased the functional correctness rate of FSM tasks, such as RTLLM `radix2_div`.

3.5 DISCUSSIONS OF RELATED WORKS

Works have proposed other means of NLD refinement or attempted to achieve related goals, such as optimizing the performance, power, and area (PPA), or designing more reliable verification schemes. Here, we discuss some future directions inspired by them.

NLD clarification. The NLDs of newer benchmarks, e.g., ArchXBench and RealBench, have better quality but are manually designed (Purini et al., 2025; Jin et al., 2025). We validated our LLM-created specs and found that they have a very similar structure and format to RealBench. [While RealBench aims to explore the abilities of the LLMs with harder tasks, we point the way for future benchmark \(and training dataset\) creators.](#) We believe that our work can overcome “the input format shortcomings in existing benchmarks” (Jin et al. (2025)) with minimal human efforts, while retaining comparability with existing results evaluated on the original version.

PPA optimization. RTLLM and AutoSilicon evaluated the PPA of LLM-generated codes (Lu et al., 2024; Li et al., 2025). However, given that the functional correctness rate remains low, we think it is too early to require the LLMs to optimize the PPA of their design. We observed that existing benchmarks include cases that aim to do so but end up being functionally incorrect. An example is given in Section 4.4. Only after removing the PPA optimization statements may the LLMs present test-passing HDL codes “in software style”, e.g., using address decoders to access arrays instead of shift registers, although they are more expensive.

[Despite this, LLM-generated designs with unsatisfactory PPA can still be used for agile deployment of FPGA-based services, provided their correctness is verified.](#) Accelerators, such as the systolic arrays and spatial accelerators mentioned by some benchmarks (Chang et al., 2024), leverage hard-

ware only for parallelization. A better PPA is usually not the user’s main concern if the design can be deployed onto the FPGA without timing or power violations. Meanwhile, trusted computing services may rely on hardware-based security features (e.g., physical isolation (Zhao et al., 2022)) and primitives (e.g., oblivious memory access (Wang et al., 2015)), which need to be deployed on FPGAs. The sooner the security features are deployed, the lower the risk to the sensitive applications.

Assertion-based verification and generation. Most TBs of RTLLM and VerilogEval are based on cycle-accurate diffests (Lu et al., 2024; Liu et al., 2023). So, to reduce FNs, our current *converter* has to require the LLMs to strictly follow the GM. However, removing constraints that are not semantically necessary can also reduce FNs. For example, when designing an accelerator, the microarchitecture and consequently the sequential behavior are not determined until the final implementation stage. So the TB should only verify the correctness of the output value, not the number of cycles required for the outputs to appear.

The authors of RealBench considered the opposite type of error, i.e., the FPs, caused by the low coverage of diffest-based TBs (Jin et al., 2025). They designed FVs for system-level design tasks. In the meantime, AssertLLM (Fang et al., 2024) served as a good step towards using LLMs to translate NLDs into formal SystemVerilog assertions (SVAs). **We believe we can further develop this idea by adding assertions to the FV step by step, until it contains sufficiently comprehensive constraints to formally generate the design, e.g., until the logic between each pair of registers is formalized. Because every step can be verified, this approach is theoretically more stable than directly completing the HDL code. We leave this as future work.**

4 EXPERIMENTS

We conducted experiments to evaluate the refined specs, aiming to answer three questions: **(Q1)** How does the quality of the spec affect the correctness of design generation? **(Q2)** How do the non-basic parts of the spec affect its quality? **(Q3)** Are the better specs helpful in complex tasks?

4.1 SETUP

Benchmarks. We applied the widely used RTLLM v2, VerilogEval-v2, and VerilogEval-Machine benchmarks (Liu et al., 2024; Pinckney et al., 2024; Liu et al., 2023) to compare with existing works. Additionally, we challenged some more complex tasks from levels 4 and 5 of ArchXBench (Purini et al., 2025). We used the original TBs provided with the benchmarks. We used Synopsys VCS as the simulator, which supports SystemVerilog syntax required by the TBs.

LLMs. For spec generation, we chose GPT-4 and Claude-3.5-sonnet (C3.5), the two commonly used general-purpose LLMs. For code generation, we chose Qwen2.5-Coder-32B (QC) as the smaller model, CodeV-R1 (CVR1) as the distilled model, and GPT-4 and C3.5 as the larger models. QC is used as a base model for CodeV (Zhao et al., 2024a) and CVR1 (Zhu et al., 2025), with CVR1 currently the most powerful finetuned model to our knowledge. We fixed the temperature to 0.3.

4.2 MAIN RESULTS

Power of the refined specs. Table 1 shows the pass@1 and pass@5 scores of LLM-based RTL design agents, given the original NLDs and the full specs (having all eight configurations in Section 3.3 enabled) generated by GPT-4. **The specs increased the pass scores of the CVR1, C3.5, and GPT-4 agents on both RTLLM-v2 and VerilogEval-v2 (Human).** Remarkably, GPT-4 succeeded on 47 (out of 50) RTLLM tasks and 151 (out of 156) VerilogEval tasks. **An in-depth analysis of the failing cases is available in Appendix D, showing that the specs can fix not only the FN cases but also true negative ones (too hard to generate even when the NLD is clear).** On VerilogEval-Machine, our spec also helped CVR1 and GPT-4 achieve an almost 100% pass@5 rate.

Compared with fine-tuning. We noticed that our specs had a negative effect on QC, especially on the VerilogEval benchmarks. QC is a model trained purely on software codes, so without the implementation guides in the original NLDs, it tends to produce syntax errors. For the CVR1 model, which is also trained from QC, we observe a 16% increase in its pass@5 score. This indicates that the refined specs can work for models powerful enough to interpret them, regardless of their size.

Table 1: The score of LLM-based RTL design agents given the original NLDs and GPT-4 specs, on RTLLM and VerilogEval. Scores in yellow cells are evaluated using RTLLM v1.1 or VerilogEval-Human (v1), which have fewer, easier cases and may yield higher scores. Scores of QC are cited from CodeV-R1 (Zhu et al., 2025). Scores of GPT-4 are cited from CodeV (Zhao et al., 2024a). Other systems marked with (*) are cited from their original papers. SP stands for self-planning. RID means Distilled R1.

Agent	Description			RTLLM		VE-H / VE-v2		VE-M	
	Base Model	FT Dataset	Strategy	pass@1	pass@5	pass@1	pass@5	pass@1	pass@5
CraftRTL-SDG-CC-Repair*	SartCoder2	Non-textual	Self-fixing	49.0	65.8	68.0	72.4	81.9	86.9
QC*	QC-32B	-	-	47.8	63.9	47.5	60.7	66.6	76.6
CodeV-All-QC*	QC-7B	Multi-layer	-	-	55.2	56.6	67.9	81.9	89.9
QC-Spec-full (Ours)	QC-32B	-	Spec	53.6	63.1	24.7	31.3	37.1	47.8
CVR1*	QC-7B	RID + Solvable ²	-	68.0	78.2	68.8	78.2	76.5	84.1
CVR1+Spec-full (Ours)	QC-7B	RID ³	Spec	79.6	83.7	81.8	82.0	99.6	100.0
C3.5-SP	C3.5	-	SP	50.0	66.1	62.4	76.8	-	-
MAGE*	C3.5	-	Self-fixing	-	-	72.4	95.7 ⁴	-	-
C3.5-SP-Spec-full (Ours)	C3.5	-	SP + spec	68.4	82.1	65.8	77.4	-	-
GPT-4*	GPT-4	-	-	-	65.5	43.5	55.8	60.0	70.6
GPT-4-SP ¹	GPT-4	-	SP	57.6	72.6	74.5	85.6	98.6	100.0
GPT-4-SP-Spec-full (Ours)	GPT-4	-	SP + spec	77.2	89.0	84.6	96.0	93.7	99.5

Compared with self-fixing. On VerilogEval-v2, the state-of-the-art self-fixing design agent (Zhao et al., 2024b) achieves a similar pass@5 score as GPT-4 using our refined spec, while their one-shot pass rate is approximately consistent with earlier works⁴. Although a high score is achieved, this self-fixing process is computationally expensive, as discussed in Section 2.2. While we acknowledge the power of self-fixing, we suggest that a better spec can reduce the required number of iterations.

Generalizability. We additionally used C3.5 to generate specs for RTLLM. To prevent a model from embedding information into the specs that only itself can interpret, we performed cross-evaluation: using one model to generate the specs, and another to generate the codes. We measured all the scores with a global temperature of 0.3, so the results may differ from those reported in previous papers.

The scores are shown in Table 2. We can see that both C3.5 and GPT-4 specs can improve the accuracy of distilled and larger models. Among the three NLDs, specs generated by GPT-4 yield the highest pass rates regardless of which model performs code generation. The score of C3.5 using GPT-generated specs is even higher than using specs generated by itself. This shows that one refined-spec version can be beneficial for multiple models.

Table 2: Cross evaluation results on the RTLLM v2 benchmark. All data comes from our experiments.

	Original NLD		C3.5 Spec		GPT-4 Spec	
	pass@1	pass@5	pass@1	pass@5	pass@1	pass@5
QC	44.8	52.6	43.6	50.7	53.6	63.1
CVR1	51.2	59.0	70.4	77.7	79.6	83.7
CVR1-SP	59.2	73.2	56.0	66.0	72.0	82.1
C3.5-SP	50.0	66.1	58.8	74.2	68.4	82.1
GPT-4-SP	57.6	72.6	66.4	79.5	77.2	89.0

We observed a performance decline compared with the original NLD when giving the C3.5 Spec to QC and CVR1-SP (i.e., with self-planning). While we have discussed the limitations of QC, comparing CVR1 with or without SP brings another insight. Although of lower quality, the length of C3.5 specs is similar to GPT-4 specs, making it more difficult for smaller models with limited context length to process. Manual inspection discovered that 14.4% and 15.2% of the CVR1-SP samples have truncated outputs (causing syntax errors) on the original NLD and C3.5, respectively. Without SP, CVR1 achieves a normal pass rate. We conclude that self-planning is unsuitable for smaller models.

Answer for Q1. The deficiencies we pointed out are limiting the power of LLM hardware design agents. Given refined specs with high clarity and simplicity, LLMs such as GPT-4 and CVR1 are practical to handle basic RTL design tasks.

¹We implemented our own self-planning prompt for improved accuracy.

²Solvable means can be solved by a powerful model, which usually implies “no NLD-GM mismatch”.

³We directly used the released CodeV-R1 model, and no extra training is applied.

⁴MAGE’s score is given as pass@1, but this “1 attempt” actually combines 20 samples and 10 debugging iterations. We regard the one-shot correct rate as pass@1, while that after 5 debugging iterations as pass@5.

4.3 ABLATION STUDY ON CONFIGURATIONS

We evaluated different configurations of our spec on RTLLM v2. We began with the Spec-basic configuration, which contains features 1-4 as declared in Section 3.3. This forms the essential sections of a spec. Then, we added examples by feature 5 (-cases), and after that, the block-level implementation overview by features 6-8 (-full, where features 7 and 8 are to prevent feature 6 from being too detailed). The results are displayed in Table 3.

The -case configuration gets the lowest pass rates among GPT-4 specs. We speculate that this is because the examples required by part 5 are presented in multi-modal formats, such as tables or timing diagrams, which are hard for current LLMs to interpret. Meanwhile, the pass rate of the -full configuration is only 2% lower than -basic while still 19% higher than the original NLD. This implies that the accuracy would still be acceptable if the implementation overview is essential (e.g., for large pipelined architectures).

Answer for Q2. Examples and implementation overviews may, to some extent, degrade the performance of design agents. However, if the contents are properly restricted, the degree of performance degradation can be controlled.

Table 3: The score of GPT-4 given specs generated with different configurations.

	RTLLM	
	pass@1	pass@5
Original NLD	51.6	66.2
SP-Spec-basic	83.6	91.8
SP-Spec-cases	75.6	85.0
SP-Spec-full	77.2	89.0

4.4 CASE STUDY 3: ARCHXBENCH LEVEL 4-5

After spec refinement, the accuracy of existing agents over simple tasks is acceptable. So, we explored the limits of our spec *creator* by letting GPT-4 challenge the recently proposed harder benchmark ArchXBench (Purini et al. (2025)). We selected the fixed-point FFT and IFFT cases in level 4 and the Conv2d case in level 5.

When creating specs for FFT and IFFT, we are asked to specify the format of the twiddle factors and decide whether the inputs should be reorganized into bit-reversed order. Then, using this spec as input, 2 out of 5 generation attempts passed the test. This is a plausible progress since the original paper reported a complete failure.

We made similar attempts on the Conv2d case. The original spec requires using shift registers to implement the sliding window to reduce the cost of address decoders. Although the convolution part was implemented correctly, GPT failed to shift the data along the correct dimension or at the correct time. Following our observation that details are harmful, we removed this requirement from the spec. This time, GPT generated an offset-pointer-based implementation. Unfortunately, it failed to build a correct zero-padding mask due to a misunderstanding of the registers’ cycle-delay behavior. This is a common issue already observed when generating dividers or floating-point arithmetic units. Nevertheless, although we consider this case a failure, we found the code highly readable and “almost correct”, such that a human engineer could fix it in about 10 minutes.

Answer for Q3. A better spec can slightly reduce the human workloads on complex designs. However, fine-tuning on tasks that involve multiple submodules or sequential logic still seems essential. We plan as future work to train the models using our refined specs.

5 CONCLUSION

A refined spec can help LLMs to produce RTL design more accurately. It should explicitly specify the semantics to avoid mismatches between requirements and test cases. Also, before LLMs are properly trained, users should not guide them to write complex implementations.

LLM has sufficient capabilities (about 90% pass@5 accuracy) to complete simple RTL designs when provided with high-quality specs generated using our *converter* or *creator* tools. Future works that create datasets or benchmarks can use our tools to automatically generate NLDs of higher quality.

REPRODUCIBILITY STATEMENT

The LLMs used in our experiments (QwenCoder, MG-Verilog, CodeV-R1, GPT-4, Claude 3.5 sonnet) are either open or provide a purchasable API. The prompts of the LLM-based tools (i.e., the *evaluator*, *converter*, and *creator*) are shown in the appendices. The original datasets and benchmarks used in our experiments (i.e., RTLLM, MG-Verilog, VerilogEval, and CodeV-R1) are publicly available, and our refined specs can be reproduced from the provided descriptions.

We plan to release the refined specs under the MIT license.

ETHICS STATEMENT

This work focuses on improving the quality of natural language descriptions for LLM-based RTL design generation. The experiments use publicly available datasets and do not involve sensitive personal data or harmful model outputs. We do not foresee significant ethical concerns associated with this research.

LLM USAGE

Large language models (LLMs) are the primary subject of this study. We used LLMs in two ways: (1) as evaluation agents to rate the clarity and simplicity of natural language descriptions (NLDs), and (2) as code-generation agents to test how NLD quality affects design correctness. LLMs were also used to provide suggestions during the exploration of what information an ideal RTL spec should include. However, all decisions regarding methodology design, analysis, and writing were made by the authors.

REFERENCES

- Kaiyan Chang, Zhirong Chen, Yunhao Zhou, Wenlong Zhu, Haobo Xu, Cangyuan Li, Mengdi Wang, Shengwen Liang, Huawei Li, Yinhe Han, et al. Natural language is not enough: Benchmarking multi-modal generative ai for verilog generation. *arXiv preprint arXiv:2407.08473*, 2024.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgren Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Matteo Esposito, Xiaozhou Li, Sergio Moreschini, Noman Ahmad, Tomas Cerny, Karthik Vaidyanathan, Valentina Lenarduzzi, and Davide Taibi. Generative ai for software architecture. applications, challenges, and future directions. *arXiv preprint arXiv:2503.13310*, 2025.
- Wenji Fang, Mengming Li, Min Li, Zhiyuan Yan, Shang Liu, Zhiyao Xie, and Hongce Zhang. AsserTLM: Generating and evaluating hardware verification assertions from design specifications via multi-llms. *arXiv preprint arXiv:2402.00386*, 2024.
- Pengwei Jin, Di Huang, Chongxiao Li, Shuyao Cheng, Yang Zhao, Xinyao Zheng, Jiaguo Zhu, Shuyi Xing, Bohan Dou, Rui Zhang, Zidong Du, Qi Guo, and Xing Hu. Realbench: Benchmarking verilog generation models with real-world ip designs. *arXiv preprint arXiv:2507.16200*, 2025.
- Cangyuan Li, Chujie Chen, Yudong Pan, Wenjun Xu, Yiqi Liu, Kaiyan Chang, Yujie Wang, Mengdi Wang, Ying Wang, Huawei Li, et al. Autosilicon: Scaling up rtl design generation capability of large language models. *ACM Transactions on Design Automation of Electronic Systems*, 2025.

- 594 Mingjie Liu, Nathaniel Pinckney, Brucec Khailany, and Haoxing Ren. Verilogeval: Evaluating large
595 language models for verilog code generation. In *2023 IEEE/ACM International Conference on*
596 *Computer Aided Design (ICCAD)*, pp. 1–8. IEEE, 2023.
- 597
598 Mingjie Liu, Yun-Da Tsai, Wenfei Zhou, and Haoxing Ren. Craftrtl: High-quality synthetic data
599 generation for verilog code models with correct-by-construction non-textual representations and
600 targeted code repair. *arXiv preprint arXiv:2409.12993*, 2025a.
- 601
602 Shang Liu, Yao Lu, Wenji Fang, Mengming Li, and Zhiyao Xie. Openllm-rtl: Open dataset and
603 benchmark for llm-aided design rtl generation. In *Proceedings of the 43rd IEEE/ACM Interna-*
604 *tional Conference on Computer-Aided Design*, pp. 1–9, 2024.
- 605
606 Yi Liu, Changran Xu, Yunhao Zhou, Zeju Li, and Qiang Xu. Deeprtl: Bridging verilog under-
607 standing and generation with a unified representation model. *arXiv preprint arXiv:2502.15832*,
2025b.
- 608
609 Yao Lu, Shang Liu, Qijun Zhang, and Zhiyao Xie. Rtlm: An open-source benchmark for design rtl
610 generation with large language model. In *2024 29th Asia and South Pacific Design Automation*
611 *Conference (ASP-DAC)*, pp. 722–727. IEEE, 2024.
- 612
613 Nathaniel Pinckney, Christopher Batten, Mingjie Liu, Haoxing Ren, and Brucec Khailany. Revis-
614 iting verilogeval: Newer llms, in-context learning, and specification-to-rtl tasks. *arXiv preprint*
arXiv:2408.11053, 2024.
- 615
616 Suresh Purini, Siddhant Garg, Mudit Gaur, Sankalp Bhat, Sohan Mupparapu, and Arun Ravindran.
617 Archxbench: A complex digital systems benchmark suite for llm driven rtl synthesis. *arXiv*
preprint arXiv:2508.06047, 2025.
- 618
619 Larissa Schmid, Tobias Hey, Martin Armbruster, Sophie Corallo, Dominik Fuchß, Jan Keim, Haoyu
620 Liu, and Anne Koziolok. Software architecture meets llms: A systematic literature review. *arXiv*
621 *preprint arXiv:2505.16697*, 2025.
- 622
623 Shailja Thakur, Baleegh Ahmad, Zhenxing Fan, Hammond Pearce, Benjamin Tan, Ramesh Karri,
624 Brendan Dolan-Gavitt, and Siddharth Garg. Benchmarking large language models for automated
625 verilog rtl code generation. In *2023 Design, Automation & Test in Europe Conference & Exhibi-*
tion (DATE), pp. 1–6. IEEE, 2023.
- 626
627 Xiao Wang, Hubert Chan, and Elaine Shi. Circuit oram: On tightness of the goldreich-ostrovsky
628 lower bound. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communi-*
cations Security, 2015.
- 629
630 Yongan Zhang, Zhongzhi Yu, Yonggan Fu, Cheng Wan, and Yingyan Celine Lin. Mg-verilog:
631 Multi-grained dataset towards enhanced llm-assisted verilog generation. In *2024 IEEE LLM Aided*
632 *Design Workshop (LAD)*, pp. 1–5. IEEE, 2024.
- 633
634 Mark Zhao, Mingyu Gao, and Christos Kozyrakis. Shef: shielded enclaves for cloud fpgas. In *ACM*
ASPLOS, 2022.
- 635
636 Yang Zhao, Di Huang, Chongxiao Li, Pengwei Jin, Muxin Song, Yanan Xu, Ziyuan Nan, Mingju
637 Gao, Tianyun Ma, Lei Qi, et al. Codev: Empowering llms with hdl generation through multi-level
638 summarization. *arXiv preprint arXiv:2407.10424*, 2024a.
- 639
640 Yujie Zhao, Hejia Zhang, Hanxian Huang, Zhongming Yu, and Jishen Zhao. Mage: A multi-agent
641 engine for automated rtl code generation. *arXiv preprint arXiv:2412.07822*, 2024b.
- 642
643 Yaoyu Zhu, Di Huang, Hanqi Lyu, Xiaoyun Zhang, Chongxiao Li, Wenxuan Shi, Yutong Wu, Jianan
644 Mu, Jinghua Wang, Yang Zhao, Pengwei Jin, Shuyao Cheng, Shengwen Liang, Xishan Zhang,
645 Rui Zhang, Zidong Du, Qi Guo, Xing Hu, and Yunji Chen. Codev-r1: Reasoning-enhanced
646 verilog generation. *arXiv preprint arXiv:2505.24183*, 2025.
- 647

A PROMPT OF THE EVALUATOR

Below is the prompt we used to evaluate the original NLDs.

```

648
649
650
651
652 1 Please act as a professional Verilog designer. We have designed a
653   submodule of our project, and tried to write a spec for it.
654 2
655 3 Here is the verilog code of the module.
656 4
657 5   ```
658 6   {{code}}
659 7   ```
660 8
661 9 Here is the spec we wrote.
662 10
663 11   ```
664 12   {{spec}}
665 13   ```
666 14
667 15 Please help us evaluate, if you want to use our module,
668 16 a) Can you learn its exact behavior from the spec without correcting any
669 17   misunderstanding by reading the Verilog code?
670 18 b) Can you learn its high-level semantics easily without awareness of its
671 19   internal implementation?
672 20
673 21 Please score (range from 1 to 5) for each of the questions, and give some
674 22   advice for us to improve. For both questions, a higher score means
675 23   that the spec is easier to read. Please give your response in JSON
676 24   format:
677 25
678 26 {
679 27     "clarity score": X,
680 28     "simplicity score": Y,
681 29     "suggestions": "..."/>

```

Below is an example of the response of the RTLLM freq_divbyodd case.

```

680 1 {
681 2     "clarity score": 2,
682 3     "simplicity score": 2,
683 4     "suggestions": "The current spec is too brief and lacks important
684   details. It does not describe the timing, duty cycle, phase, or
685   waveform characteristics of clk_div. It also omits the behavior
686   for different values of NUM_DIV, especially how the output clock
687   waveform looks for arbitrary odd divisors. There is no mention of
688   the output's duty cycle (e.g., is it 50%), phase relationship
689   to clk, or reset behavior in terms of output state. To improve:
690   1) Explicitly describe the output waveform, including duty cycle
691   and phase. 2) State how the divider works for any odd NUM_DIV,
692   and what happens at reset. 3) Add a timing diagram or example
693   waveforms. 4) Specify any constraints or limitations (e.g., only
694   odd divisors, minimum/maximum values). 5) Clarify whether NUM_DIV
695   is static or can be changed at runtime. This will make the
696   module much easier to use as a black box and understand its
697   behavior without reading the code."
698 5 }

```

B PROMPT OF THE CONVERTER

Below is the prompt we used to convert RTLLM and VerilogEval descriptions into our specs.

701

702

703 1 Please act as a professional Verilog designer. We have designed a {{name
704 }} module, and we wish to write a spec for it.

705 2

706 3 Here is the verilog code of the module.

707 4

```
707 5 ```
708 6 {{code}}
709 7 ```
```

710 8

711 9 Here is the spec we wrote.

712 10

```
712 11 {{spec}}
```

713 12

714 13 Please help us improve the spec, such that a third-party developer can
715 learn the exact behavior of this module from the spec without reading
716 the verilog code. Please be cautious that the behavior of the code
717 might be unnatural, and the description might be unclear or mistaken.
718 You should strictly follow the code (especially, the output
sometimes should be advanced or delayed for several cycles).

719 14

720 15 What should be contained:

721 16

- 722 17 - Module name.
- 723 18 - Behavior details. Please specify the function of this module, i.e., how
724 the output is computed from the input. If it is a sequential logic (e.g.,
725 an FSM or pipeline), please specify the timing, i.e., the handshaking
726 protocols, the number of pipeline stages, or the number of cycles before
the valid output comes out.
- 727 19 - Parameters. If the module has configurable parameters, please specify
their meanings and default values.
- 728 20 - I/O details. Please specify the semantics of each input or output port.
729 Please detail the meaning of each possible value of the inputs,
730 especially for control signals.
- 731 21 - Normal cases and corner cases. Please design some cases (input-output
732 pairs) according to the code. We have assured that the code is correct,
733 however its corner case behavior might not have been written in the
original spec.
- 734 22 - Implementation overview. Please include some high-level hints of the
735 module's structure. For example, you may describe the semantics of each
736 state in an FSM, each stage in a pipeline, or submodules that might be
instantiated.

737 23

738 24 What should NOT be contained:

739 25

- 740 26 - Complicated white-box details. Third-party users wish to treat the
741 module as a black box, and details might be too complex.
- 742 27 - Trivial corner case behaviors. For example, if the code did not handle
743 overflows, we may just say "assume there are no overflows".

743 28

744 29 Please give us the spec in Markdown format.

745

746

747

Below is the improved spec for the RTLLM serial2parallel case as shown in Figure 2.

748

```
748 1 # Specification: `serial2parallel` Module
```

749 2

```
749 3 ## 1. Module Name
```

750 4

```
750 5 **serial2parallel**
```

751 6

```
751 7 ---
```

752 8

```
752 9 ## 2. Functional Overview
```

753 10

```

756 11 The `serial2parallel` module converts a stream of serial input bits into
757 12 an 8-bit parallel output. It collects 8 consecutive serial bits, in
758 13 order, and then outputs them as an 8-bit parallel value. The module
759 14 uses a clocked, synchronous design with an active-low reset. The
760 15 output data is accompanied by a validity signal to indicate when the
761 16 parallel output is valid.
762 17
763 18 **Important Note:**
764 19 The module's output (`dout_parallel` and `dout_valid`) is only asserted
765 20 after the **8th valid serial bit is received**. The output is not
766 21 immediately available after the 8th input; instead, it is presented
767 22 on the next clock cycle, and the valid signal is asserted for one
768 23 clock cycle only.
769 24
770 25 ---
771 26
772 27 ## 3. Parameters
773 28
774 29 - **No configurable parameters.**
775 30 The module is fixed to collect 8 serial bits per parallel output.
776 31
777 32 ---
778 33
779 34 ## 4. I/O Ports and Semantics
780 35
781 36 | Port Name          | Direction | Width | Description |
782 37 |-----|-----|-----|-----|
783 38 | clk                | Input    | 1     | Clock signal. All operations are
784 39 |                    |          |       | synchronous to the rising edge of `clk`.
785 40 | rst_n              | Input    | 1     | Asynchronous reset, active low.
786 41 |                    |          |       | When `rst_n` is 0, all internal state and outputs are reset.
787 42 | din_serial         | Input    | 1     | Serial data input. The value is
788 43 |                    |          |       | sampled on the rising edge of `clk` when `din_valid` is high.
789 44 | din_valid          | Input    | 1     | Input data valid signal. When high,
790 45 |                    |          |       | indicates that `din_serial` holds a valid bit to be sampled.
791 46 | dout_parallel      | Output   | 8     | Parallel output data. Valid only
792 47 |                    |          |       | when `dout_valid` is high.
793 48 | dout_valid         | Output   | 1     | Output valid signal. High for one
794 49 |                    |          |       | clock cycle when `dout_parallel` contains valid data.
795 50
796 51 ### Input Signal Details
797 52
798 53 - **clk**: All state transitions and data sampling occur on the rising
799 54 edge.
800 55 - **rst_n**: When low, the module resets all internal registers and
801 56 outputs.
802 57 - **din_serial**: Carries the serial data bit to be collected.
803 58 - **din_valid**:
804 59 - `1`: `din_serial` is sampled and counted toward the 8-bit output.
805 60 - `0`: No data is sampled; the internal counter and buffer are reset.
806 61
807 62 ### Output Signal Details
808 63
809 64 - **dout_parallel**: Contains the 8 collected serial bits, with the first
810 65 received bit as the most significant bit (MSB), and the last as the
811 66 least significant bit (LSB).
812 67 - **dout_valid**:
813 68 - `1`: Indicates `dout_parallel` contains valid data. Asserted for **
814 69 one clock cycle** after 8 valid serial bits are received.
815 70 - `0`: Output is not valid.
816 71

```

```

810 ---
811 52
812 53
813 54 ## 5. Behavioral and Timing Details
814 55
815 56 ### Data Collection and Output Protocol
816 57
817 58 - The module begins collecting serial bits when `din_valid` is high.
818 59 - Each rising edge of `clk` with `din_valid = 1`:
819 60   - The value of `din_serial` is shifted into an internal buffer.
820 61   - An internal counter increments.
821 62 - When 8 valid bits have been collected:
822 63   - On the next clock cycle, `dout_parallel` outputs the 8-bit value,
823 64     and `dout_valid` is asserted high for one clock cycle.
824 65   - The output bits are ordered such that the first serial bit received
825 66     becomes the MSB (`dout_parallel[7]`), and the last bit becomes
826 67     the LSB (`dout_parallel[0]`).
827 68 - If `din_valid` is low at any time, the internal counter and buffer are
828 69   reset to zero, and the collection process restarts.
829 70
830 71 ### Reset Behavior
831 72
832 73 - When `rst_n` is low:
833 74   - All internal state (counter, buffer) and outputs (`dout_parallel`, `
834 75     dout_valid`) are reset to zero immediately.
835 76
836 77 ---
837 78 ## 6. Normal and Corner Case Examples
838 79
839 80 ### Example 1: Normal Operation
840 81
841 82 | Cycle | din_valid | din_serial | Internal Counter | Buffer (MSB..LSB) |
842 83 |-----|-----|-----|-----|-----|
843 84 | 0     | 0        | X          | 0                | 00000000         |
844 85 | 1     | 1        | 1          | 1                | 00000001         |
845 86 | 2     | 1        | 0          | 2                | 00000010         |
846 87 | 3     | 1        | 1          | 3                | 00000101         |
847 88 | 4     | 1        | 1          | 4                | 00001011         |
848 89 | 5     | 1        | 0          | 5                | 00010110         |
849 90 | 6     | 1        | 1          | 6                | 00101101         |
850 91 | 7     | 1        | 0          | 7                | 01011010         |
851 92 | 8     | 1        | 1          | 8                | 10110101         |
852 93 | 9     | 1        | X          | 1                | X                 |
853 94 | 10    | 1        | 10110101  | 1                | X                 |
854 95
855 96 - On cycle 9, `dout_valid` is asserted and `dout_parallel` outputs the
856 97   collected bits: `10110101`.
857 98
858 99 ### Example 2: din_valid Drops Before 8 Bits
859 100
860 101 | Cycle | din_valid | din_serial | Internal Counter | dout_valid |
861 102 |-----|-----|-----|-----|-----|
862 103 | 0     | 1        | X          | 0                | 0           |
863 104 | 1     | 1        | 1          | 1                | 1           |
864 105 | 2     | 1        | 0          | 2                | 1           |
865 106 | 3     | 1        | 1          | 3                | 1           |
866 107 | 4     | 1        | 1          | 4                | 1           |
867 108 | 5     | 1        | 0          | 5                | 1           |
868 109 | 6     | 1        | 1          | 6                | 1           |
869 110 | 7     | 1        | 0          | 7                | 1           |
870 111 | 8     | 1        | 1          | 8                | 1           |
871 112 | 9     | 0        | X          | 9                | 0           |
872 113 | 10    | 0        | X          | 10               | 0           |

```

```

864 97 | 1      | 1      | 1      | 1      | 0      |
865      | 00000000 |      |      |      |      |
866 98 | 2      | 1      | 0      | 2      | 0      |
867      | 00000000 |      |      |      |      |
868 99 | 3      | 0      | X      | 0      | 0      |
869      | 00000000 |      |      |      |      |
870 100 | 4      | 1      | 1      | 1      | 0      |
871      | 00000000 |      |      |      |      |
872 101 - If `din_valid` goes low before 8 bits are collected, the counter and
873 102     buffer reset. The collection process restarts from the next valid bit
874      .
875 103 ### Example 3: Continuous din_valid High
876 104
877 105 - If `din_valid` remains high, the module will output a valid parallel
878 106     word every 8 cycles, with no gaps, as long as `din_serial` continues
879      to provide data.
880 107 ### Example 4: Reset Behavior
881 108
882 109 - If `rst_n` is asserted low at any time, all outputs and internal states
883 110     are cleared immediately, regardless of other inputs.
884 111 ---
885 112
886 113 ## 7. Implementation Briefs
887 114
888 115 - The module uses:
889 116     - An 8-bit shift register to collect serial bits.
890 117     - A 4-bit counter to track the number of valid bits received.
891 118 - Data is shifted in on each clock cycle where `din_valid` is high.
892 119 - When the counter reaches 8, the parallel output is presented and the
893 120     valid signal is asserted for one cycle.
894 121 - If `din_valid` is ever low, the counter and buffer reset, and the
895 122     collection process restarts.
896 123 ---
897 124 ## 8. Summary
898 125
899 126 - **Collects 8 serial bits (with `din_valid` high) and outputs them as an
900 127     8-bit parallel word.**
901 128 - **Output is valid for one clock cycle after the 8th bit is received.**
902 129 - **If `din_valid` drops before 8 bits are collected, the process
903 130     restarts.**
904 131 - **Reset (`rst_n` low) clears all state and outputs immediately.**
905 132 ---
906 133 **End of Specification**
907

```

The 9th cycle of Example 1 in the "Normal and Corner Case Examples" section pointed out that every 9th bit will be discarded. 12 out of 20 attempts passed using this spec (pass@5 = 0.99), while 0 out of 20 attempts passed using the original description.

C PROMPT OF THE CREATOR

The prompt for the *creator* is similar to the *converter*'s, only appended with the following sentence.

Please add placeholders wherever you wish us to add more details.

918 This allows the creator to seek completion for undefined behaviors. Below is the response for the
 919 ArchXBench FFT/IFFT task as shown in 4.4, added with placeholders.
 920

```

921 1 Here is the improved, professional-grade specification for your module `
922   fft16_iterative`, rewritten in clear Markdown format. It captures
923   both user-facing semantics and timing behavior without exposing
924   excessive implementation detail.
925 2
926 3 ---
927 4 # Specification: `fft16_iterative`
928 5
929 6 ## Title
930 7
931 8 **16-point Iterative Fixed-Point Radix-2 FFT/IFFT Module**
932 9
933 10 ---
934 11
935 12 ## Summary
936 13
937 14 `fft16_iterative` is a resource-efficient, fixed-size, radix-2 FFT/IFFT
938 15 processor that computes the 16-point Discrete Fourier Transform (DFT)
939   or its inverse (IDFT) over complex fixed-point input vectors. It
940   uses an iterative, decimation-in-time (DIT) architecture to
941   minimize area and twiddle storage by reusing a single butterfly unit
942   across stages.
943 16
944 17 This module is designed for use in DSP pipelines and embedded signal-
945 18 processing systems with limited hardware resources, and supports 
946 19 fixed-point input/output, IFFT mode, and automatic output
947 20 scaling.
948 21
949 22 ---
950 23 ## Functional Behavior
951 24
952 25 ### Transform Type
953 26
954 27 The module supports:
955 28
956 29 * Forward FFT (mode = 0): Compute standard DFT with radix-2 DIT.
957 30 * Inverse FFT (IFFT) (mode = 1): Use conjugated twiddle factors and
958 31   output normalization by 1/N.
959 32
960 33 ### Timing Protocol
961 34
962 35 * Computation begins on the rising edge of `start` (a 1-cycle
963 36   pulse).
964 37 * Input data must be stable at the moment of `start`.
965 38 * Output will become valid after a fixed number of cycles (see below).
966 39 * `done` is asserted high for one clock cycle to indicate that `
967 40   data_real_out` and `data_imag_out` are valid.
968 41 * The module is sequential, not pipelined. It processes each FFT in
969 42   isolation.
970 43
971 44 ### Latency
972 45
973 46 * Latency = `log2(N) * (N/2)` cycles (for N=16, that is 4 * 8 = 32
974 47   cycles).
975 48 * `done` is asserted exactly 32 cycles after `start`.
976 49
977 50 ---
978 51 ## Module Declaration
  
```

```

972 47 ``verilog
973 48 module fft16_iterative #(
974 49     parameter N = 16,
975 50     parameter DATA_W = 12,
976 51     parameter COEFF_W = 16,
977 52     parameter GAIN_W = 4
978 53 ) (
979 54     input clk,
980 55     input rst,
981 56     input start,
982 57     input mode, // 0: FFT, 1: IFFT
983 58     input signed [DATA_W-1:0] data_real_in [0:N-1],
984 59     input signed [DATA_W-1:0] data_imag_in [0:N-1],
985 60     output signed [DATA_W+GAIN_W-1:0] data_real_out [0:N-1],
986 61     output signed [DATA_W+GAIN_W-1:0] data_imag_out [0:N-1],
987 62     output done
988 63 );
989 64 ``
990 65 ---
991 66
992 67 ## Parameters
993 68
994 69 | Name          | Default | Description
995 70 |-----|-----|-----
996 71 | `N`           | 16      | Transform size (only 16 is currently supported).
997 72 | `DATA_W`      | 12      | Bit-width of real/imag input.
998 73 | `COEFF_W`     | 16      | Bit-width of twiddle coefficients.
999 74 | `GAIN_W`      | 4       | Maximum bit-growth from internal computation;
1000 75 |               |         | output width is `DATA_W + GAIN_W`.
1001 76 |-----|-----|-----
1002 77 ---
1003 78
1004 79 ## Inputs
1005 80
1006 81 | Signal          | Width  | Description
1007 82 |-----|-----|-----
1008 83 | `clk`          | 1      | Clock input.
1009 84 | `rst`          | 1      | Synchronous reset (active high).
1010 85 |               |         | Aborts current computation and resets FSM.
1011 86 | `start`        | 1      | One-cycle pulse to trigger a new
1012 87 |               |         | computation. Input data must be stable during this pulse.
1013 88 | `mode`         | 1      | `0`: FFT; `1`: IFFT.
1014 89 |-----|-----|-----
1015 90 | `data_real_in[N]` | `DATA_W` signed | Real parts of input vector.
1016 91 |               |         | Indexed from 0 to N-1.
1017 92 | `data_imag_in[N]` | `DATA_W` signed | Imaginary parts of input vector.
1018 93 |               |         | Indexed from 0 to N-1.
1019 94 |-----|-----|-----
1020 95 ## Outputs

```

```

1026 | Signal          | Width          | Description
1027 |
1028 | |
102995 | ----- | ----- | -----
1030 |
103196 | `data_real_out[N]` | `DATA_W + GAIN_W` signed | Real parts of the
1032 | output FFT/IFFT vector.
103397 | `data_imag_out[N]` | `DATA_W + GAIN_W` signed | Imaginary parts of the
1034 | output FFT/IFFT vector.
103598 | `done`          | 1              | Pulses high **exactly
1036 | 32 cycles after `start`**, for one cycle, to signal valid output. |
103799 |
1038100 | ---
103901 | ## Output Semantics
104002 |
104103 | * **IFFT scaling**:  
104204 | Output is normalized by 1/N (for N=16) using  
104305 | arithmetic shift (`>> 4`).
104406 | * **Output precision**:  
104507 | The result allows `GAIN_W` bits of headroom to  
104608 | avoid overflow due to butterfly additions/multiplications.
104709 | * **Output format**:  
104810 | The output order is **bit-reversed** compared to  
104911 | natural frequency ordering. You may apply a post-reordering step if  
105012 | needed.
105113 |
105214 | ---
105315 | ## Corner Case Behavior
105416 |
105517 | * `start` must not be asserted while a previous FFT is ongoing. If  
105618 | asserted early, behavior is undefined.
105719 | * If `rst` is asserted during a transform, computation is aborted and  
105820 | internal state is reset. No output will be produced for that round.
105921 | * If `mode` changes during a transform, the behavior is undefined. `mode`  
106022 | must be stable from `start` until `done`.
106123 |
106224 | ---
106325 | ## Example Timing (N = 16)
106426 |
106527 | | Cycle | `start` | Action |
106628 | | ----- | ----- | ----- |
106729 | | 0 | 1 | Input latched, transform starts |
106830 | | 1-31 | 0 | Internal butterfly computation |
106931 | | 32 | 0 | `done = 1`, output valid |
107032 |
107133 | ---
107234 | ## Implementation Overview
107335 |
107436 | * **Butterfly unit**:  
107537 | Computes radix-2 DIT butterfly:  
107638 | * Inputs: two complex numbers + twiddle factor  
107739 | * Computes: `X[k] = A + BW`, `X[k+N/2] = A - BW`
107840 | * **Control FSM**:  
107941 | * Manages stage iteration (0 to log(N)-1)  
108042 | * Selects butterfly pairs and ROM twiddles
108143 | * **Twiddle ROM**:  
108244 | * Precomputed cos/sin values (for FFT)  
108345 | * For IFFT, twiddles are complex conjugated
108446 | * **Data memory**:  
108547 | * Intermediate results stored in dual-ported BRAM or registers

```

```

1080      * Reused across stages
1081      * **Output scaling**:
```

```

1082      * Optional normalization for IFFT via arithmetic right shift by `log(N)`
1083
1084      ---
1085
1086      ## Test Cases
1087
1088      ### Case 1: Impulse Input (FFT)
1089
1090      * Input: `data_real_in[0]=1`, rest zeros; `mode=0`
1091      * Expected Output: `data_real_out[i] = 1`, `data_imag_out[i] = 0` for all
1092        i
1093
1094      ### Case 2: Constant Input (IFFT)
1095
1096      * Input: All `data_real_in[i] = 1`, `data_imag_in[i] = 0`; `mode=1`
1097      * Expected Output: `data_real_out[0]=16`, rest zeros, normalized to `1` (
1098        after >>4)
1099
1100      ### Case 3: Overflow Check
1101
1102      * Input: Mix of +2047 and -2048 at `DATA_W=12`
1103      * Expected Output: No overflow; values within `DATA_W + GAIN_W` signed
1104        range
1105
1106      ---
1107
1108      ## To Be Added
1109
1110      * [ ] **Exact internal FSM state list and transitions** (optional for doc
1111        readers)
1112      * [ ] **Twiddle factor format (e.g., fixed-point Q1.15?)**
1113      * [ ] **Clarify output order (natural vs bit-reversed)**
1114      * [ ] **Specify rounding/truncation method (e.g., truncation vs rounding-
1115        to-nearest)**
1116      * [ ] **Saturation behavior (if any)**
1117
1118      ---
1119
1120      Let me know if you'd like to generate a diagram (FSM or butterfly flow),
1121      testbench scaffolding, or a verification matrix from this spec.
1122
1123
1124
1125

```

The implementation overview part indicates the function of the butterfly unit. Also, the creator noticed that the output might be bit-reversed, so it added a placeholder and asked us to fill it in.

D FAILED CASES ANALYSIS OF GPT-4-SP-SPEC-FULL ON RTLLM

We manually examined the designs of the 50 RTLLM cases generated by GPT-4, and listed our findings in Table 4.

GPT-4 achieves higher accuracy on our spec than on the original NLDs in 20 cases, whereas the original NLDs completely failed in 11 cases. We confirmed that 9 of the completely failed cases and 4 of the less accurate cases are caused by **(D1)**. These cases can be fixed or improved with our spec.

Only 1 case failed on our spec but succeeded on the original NLD. We found that the full spec we generated unexpectedly included too many implementation details, and causes **(D2)**. When given the basic spec, i.e., removing the optional “implementation overview”, GPT-4 passed the test.

1134 Table 4: RTLLM cases in which GPT-4 performed differently given the original NLD and Spec, or
 1135 both failed. The numbers are the test-passing count among 10 samples.
 1136

Type	Name	Original NLD	Spec-full	FP
Fixed	freq_divbyeven	0	10	✓
	sequence_detector	0	6	
	signal_generator	0	10	✓
	asyn_fifo	0	2	
	radix2_div	0	2	✓
	alu	0	10	✓
	serial2parallel	0	6	✓
	parallel2serial	0	10	✓
	pulse_detect	0	4	✓
	clkgenerator	0	10	✓
multi_8bit	0	10	✓	
Improved	barrel_shifter	2	10	✓
	LFSR	2	10	
	traffic_light	6	10	✓
	freq_divbyfrac	2	8	
	fixed_point_subtractor	4	8	
	JC_counter	4	10	✓
	multi_booth_8bit	2	10	
	freq_div	4	10	✓
multi_pipe_4bit	8	10		
Both failed	float_multi	0	0	
	freq_divbyodd	0	0	
Weakened	pe	0	0	
	div_16bit	2	0	

1153
 1154 In the other 3 cases, GPT-4 failed on both specs and original NLDs. This indicates that, although
 1155 a better spec indeed yields a large improvement, enhancing LLMs’ understanding of sequential
 1156 behavior remains necessary.
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187