

Chip-Chat: Challenges and Opportunities in Conversational Hardware Design

Jason Blocklove
New York University
New York, NY USA
jason.blocklove@nyu.edu

Siddharth Garg
New York University
New York, NY USA
siddharth.garg@nyu.edu

Ramesh Karri
New York University
New York, NY USA
rkarri@nyu.edu

Hammond Pearce
University of New South Wales
Sydney, Australia
hammond.pearce@unsw.edu.au

Abstract—Modern hardware design starts with specifications provided in natural language. These are then translated by hardware engineers into appropriate Hardware Description Languages (HDLs) such as Verilog before synthesizing circuit elements. Automating this translation could reduce sources of human error from the engineering process. But, it is only recently that artificial intelligence (AI) has demonstrated capabilities for machine-based end-to-end design translations. Commercially-available instruction-tuned Large Language Models (LLMs) such as OpenAI’s ChatGPT and Google’s Bard claim to be able to produce code in a variety of programming languages; but studies examining them for hardware are still lacking. In this work, we thus explore the challenges faced and opportunities presented when leveraging these recent advances in LLMs for hardware design. Given that these ‘conversational’ LLMs perform best when used interactively, we perform a case study where a hardware engineer co-architects a novel 8-bit accumulator-based microprocessor architecture with the LLM according to real-world hardware constraints. We then sent the processor to tapeout in a Skywater 130nm shuttle, meaning that this ‘Chip-Chat’ resulted in what we believe to be the world’s first wholly-AI-written HDL for tapeout.

Index Terms—Hardware Design, CAD, LLM

I. INTRODUCTION

A. Trends in hardware design

As digital designs continue to grow in capability and complexity, software components in Integrated Circuit (IC) Computer Aided Design (CAD) have adopted machine learning (ML) throughout the Electronic Design Automation flow (e.g. [1]–[3]). Where traditional approaches try to formally model each process, ML-based approaches focus on identifying and exploiting generalizable high-level features or patterns [3]—meaning ML can augment or even replace certain tools. Still, ML research in IC CAD tends to focus on the back-end processes such as logic synthesis, placement, routing, and property estimation. In this work, we instead explore the challenges and opportunities when applying an emerging type of ML model to the earliest stages of the hardware design processes: the writing of Hardware Description Language (HDL) itself.

B. Automating Hardware Description Languages (HDLs)

While hardware designs are expressed in formal languages (HDLs), they actually begin the design lifecycle as specifications provided in natural language (e.g. English-language requirements documents). The process of translating these into the appropriate HDL (e.g. Verilog) must be done by hardware engineers, which is both time-consuming and error-prone [4]. Alternative pathways, such as using high-level synthesis tools [5], can enable developers to specify functionality in higher-level languages like C, but these methods come at the expense of hardware efficiency. This motivates the exploration of Artificial Intelligence (AI) or ML-based tools as an alternative pathway for translating specifications to HDL.

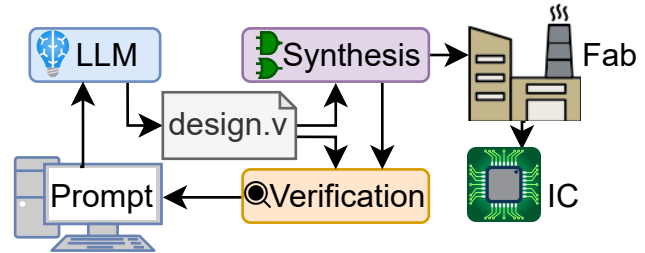


Fig. 1. Can conversational LLMs be used to iteratively design hardware?

The obvious candidate for this machine translation application comes from the Large Language Models (LLMs) [6] popularized by commercial offerings such as GitHub Copilot [7]. LLMs claim to produce code in a variety of languages and for a variety of purposes. Still, they focus on software, and benchmarks for these models evaluate them for languages such as Python, rather than on the needs present in the hardware domain. As such, adoption by the hardware design community continues to lag behind that in the software domain. Although steps for benchmarking the ‘autocomplete’ style models have begun to appear in the literature [8], the latest LLMs such as OpenAI’s ChatGPT [9] and Google’s Bard [10] instead provide a different ‘conversational’ chat-based interface to their capabilities.

Therefore, we pose the following question: What are the potential advantages and obstacles associated with integrating these tools into the HDL development process (Figure 1)? To understand this, we perform a directed but open-ended “free chat” where an LLM serves as a co-hardware architect during the development of a novel 8-bit processor (Section III).

In order to comprehend the significance of this emerging technology, it is crucial to conduct observational studies like this one. Similar studies are being carried out for ChatGPT in various domains, including software [11] and education [12], making this investigation into the impact of conversational LLMs on hardware design both relevant and timely.

C. Contributions

Our contributions include the following:

- Conducting the first investigation into the use of conversational LLMs in Hardware design.
- Conducting an observational study on the end-to-end co-design of a complex application in Hardware, utilizing ChatGPT-4.
- Achieving a significant milestone by using AI to write the complete HDL for a tapeout for the first time.
- Providing practical recommendations for the effective utilization of cutting-edge conversational LLMs in hardware-related tasks.

Open-source: All benchmarks, tapeout toolchain scripts, generated Verilog and LLM conversation logs are provided on Zenodo [13].

II. BACKGROUND AND RELATED WORK

A. Large Language Models (LLMs)

It wasn't until recently with the GPT-3 [14] family of models that the relative capabilities of these Large Language Models (LLMs) became apparent. These include Codex [6], which has billions of learned parameters and is trained on millions of open-source software repositories. In the state of the art there are dozens of LLMs, open-source, closed-source, and commercial, with options for general and task-specific applications.

Still, all LLMs share commonalities. They act as 'Scalable sequence prediction models' [6], meaning that given some 'input prompt' they will output the 'most likely' continuation of that prompt (think of them as a 'smart autocomplete'). For this I/O, they use *tokens*, which are common character sequences specified using byte pair encoding. This is efficient as LLMs have a fixed context size, meaning that they can ingest more text than they could by operating over characters. For OpenAI's models, each token represents about 4 characters, and their context windows range up to 8,000 tokens in size (meaning they can support about 16,000 characters of I/O).

B. Large Language Models for hardware design

The first work exploring LLMs for use in the hardware domain was by Pearce et al. [15]. They fine-tuned a GPT-2 model (that they termed DAVE) over synthetically generated Verilog snippets and evaluated the model outputs lexically for 'undergraduate-level' tasks. However, due to the limited training data, the model does not generalize to unfamiliar tasks. Thakur et al. [8] extended this idea, exploring both how model performance for generating Verilog could be evaluated rigorously and using different strategies for training Verilog-writing models. Other works have investigated the implications of such models: [16] examined the incidence rates of 6 types of hardware bugs in Verilog code by GitHub Copilot, and when [17] explored if automated bug repair could be achieved using the Codex models, they also included two hardware CWEs in Verilog.

In the industry, there is also increasing interest: new companies like RapidSilicon are promoting upcoming (but not yet released) tools like RapidGPT [18] which will work in this space.

C. Instruction-tuned 'conversational' models

Recently, a new kind of training methodology has been applied to LLMs which, when combined with labelled data for specific intents, can produce instruction-tuned models more capable of following a user's intent. Where previous LLMs focused on 'autocomplete', they can be instead trained to 'follow instructions'. Methodologies not requiring the (non-scalable) human feedback have followed [19]. These can then be fine-tuned to better focus on *conversational* style interactions. Models such as ChatGPT [9] (including ChatGPT-3.5 and ChatGPT-4 versions) were trained using these techniques. They provide an exciting new potential interface for works in the hardware domain. However, to the best of the authors knowledge, no such application has yet been explored.

III. LLM ASSISTED DESIGN SPACE EXPLORATION (CHIP CHAT)

A. Overview

Real-world hardware design can often have quite nuanced and complex requirements, as compared to standard HDL benchmarks. This is a challenge when considering the potential uses for conversational LLMs, such as rigidly scripting and constraining the ways that a user can interact with the LLMs, or letting a user communicate with the model how they see fit at any given time. We seek to investigate if unstructured conversations might allow for

```
Let us make a brand new microprocessor design together. We're severely constrained on space and I/O. We have to fit in 1000 standard cells of an ASIC, so I think we will need to restrict ourselves to an accumulator based 8-bit architecture with no multi-byte instructions. Given this, how do you think we should begin?
```

Fig. 2. 8-bit accumulator-based processor: Starting design prompt

greater levels of performance and mutual creativity. Investigating this would generally be done with a large-scale user-study, with hardware engineers paired with the tool during development. Such studies have been done in the software domain for LLMs, e.g. this example from Google which paired their proprietary LLM with >10,000 software developers [20] and found measurable, positive impacts on developer productivity (reduced their coding iteration duration by 6% and reduced number of context switches by 7%). We aim to motivate such a study for the hardware domain by performing a proof of concept experiment, where we pair a capable commercial LLM (OpenAI's ChatGPT-4) with an experienced hardware design engineer (one of the paper authors), and qualitatively examine the outcome when tasked with making a complex design, as outlined next.

B. Design Task: An 8-bit accumulator-based microprocessor

Constraints: With the goal of taping out this design on Tiny Tapeout, we must adhere to the constraints of that format. This restricts the design to 8 input bits, including the clock and reset, and 8 output bits, as well as only 1000 standard cells. We wish for ChatGPT-4 to write all the processor's Verilog (excluding the top-level Tiny Tapeout wrapper). To ensure we can load and unload data from the processor, we require all registers to be connected in a 'scan chain' of shift registers.

Overall goal: Design an 8-bit accumulator-based architecture alongside ChatGPT-4 in an unscripted manner. The initial prompt to ChatGPT-4 is provided in Figure 2. Given the space restriction, we aimed for a von Neumann type design with 32 bytes of memory (combined data and instruction).

Task partitioning: For this design task, the experienced human engineer was responsible (a) for shepherding ChatGPT-4, and (b) for verifying its output (e.g. syntax checks, authoring verification code / testbenches). Meanwhile, ChatGPT-4 was solely responsible for the Verilog code for the processor. It also produced the majority of the processor's specification.

C. Method: Conversation flow

General process: The microprocessor design process began by defining the Instruction Set Architecture (ISA), then implementing components that the ISA would require, before combining those components in a datapath with a control unit to manage them. Simulation and testing were used to find bugs which were then repaired.

Conversation threading: Given that ChatGPT-4, like other LLMs, has a fixed-size context window (see Section II-A), we assumed that the best way to prompt the model is by breaking up the larger design into subtasks which each had its own 'conversation thread' in the interface. This keeps the overall length below 16,000 characters. A proprietary back-end method performs some kind of text reduction when the length exceeds this, but details on its implementation are scarce. Since ChatGPT-4 does not share information between threads, the human engineer would copy the pertinent information from the previous thread into the new first message, growing a 'base specification' that slowly comes to define the processor. The

TABLE I

CONVERSATION FLOW MAP: THE PROCESSOR WAS BUILT THROUGH A LINEAR FLOW OF 125 USER MESSAGES ACROSS 18 TOPICS IN 11 ‘CONVERSATION THREADS’.

Cont. T. ID	T. ID	Topic	# User Msgs	# Restart	# User Lines	# User Chars	# LLM Lines	# LLM Chars
-	00	Specification	22	10	45	5025	498	44818
-	01	Register specification	6	2	59	4927	91	9961
-	02	Shift registers and memory	5	5	65	5444	269	9468
-	03	Multi-cycle planning and ALU	7	2	103	7284	243	10148
-	04	Control signal planning	13	21	216	9205	414	20364
-	05	Control Unit state logic	12	11	216	9898	742	21663
-	06	ISA to ALU opcode	4	0	72	4576	149	5624
-	07	Control unit output logic	11	6	266	8632	518	19180
-	08	Datapath components	12	0	144	5385	516	15646
-	09	Python assembler	3	4	127	4231	218	6270
00	10	Spec. branch update	1	1	14	1275	15	1635
07	11	Control Unit branch update	2	2	98	3743	101	3969
08	12	Datapath branch update	2	0	25	888	20	726
11	13	Control Unit bug fixing	6	1	190	5413	241	8001
-	14	Memory mapped components	7	0	79	3079	516	16237
-	15	Shift Register bug fix	2	0	38	985	85	2593
12	16	Datapath bug fixing & updates	6	0	116	2979	128	4613
14	17	Memory mapped constants	4	0	21	849	101	4655
03	18	ALU optimization	1	0	2	98	32	1368
TOTALS			125	65	1896	83916	4897	206939

base specification eventually included the ISA, a list of registers (Accumulator ‘ACC’, Program Counter ‘PC’, Instruction Register ‘IR’), the definitions for the memory bank, ALU, and control unit, and a high-level overview of what the processor should do in each cycle. Most of the information in this specification was produced by ChatGPT-4 and copy/pasted and lightly edited by the human.

Topics: One topic per thread worked well for the early design stages of the processor (with one exception, where the ALU was designed in the same thread as the multi-cycle processor clock cycle timing plan). However, once the processor got to the simulation stage and we ran programs on it, we found mistakes and bugs in the specification and implementation. Rather than starting new conversation threads and rebuilding the previous context, the design engineer instead chose to continue previous conversation threads where appropriate. We illustrate this in our flow map in Table I, where the ‘Cont. T. ID’ column indicates if they ‘Continued’ a previous thread (and if so, which thread).

Bug repair: After errors were encountered, we would use ChatGPT-4 to repair them. An example of this is presented in Figure 3. This is based on Conversation 15, and demonstrates how an error in the shift register primitive was resolved.

Restarts: Sometimes ChatGPT-4 outputs suboptimal responses. If so, the engineer has two options: (1) continue the conversation and nudge it to fix the response, or (2) use the interface to force ChatGPT-4 to ‘restart’ the response, i.e. regenerating the result by pretending the previous answer never occurred. Choosing between these has trade-offs and requires professional judgement: continuing the conversation allows for the user to specify which parts of the previous response are good or bad, but regeneration will keep the overall conversation shorter and more succinct (valuable considering the finite context window size). Still, as can be seen from the ‘# Restart’ column in Table I, the number of restarts tended to decrease as the engineer grew more experienced with using ChatGPT-4, with Topics 00-07 having 57 restarts compared to Topics 08-18 having just 8. The highest individual number of restarts on a single message was 10, in Topic 04 (Control signal planning) which has the message in Figure 4. This was a difficult prompt because it asks for a specific kind of output with a significant amount of detail, but eventually yielded a satisfactory answer as listed in Figure 5.

```

1 I have the following shift_register written in Verilog:
2
3 module shift_register #(parameter WIDTH = 8)(
4     // Ports omitted for space
5 );
6     reg [WIDTH-1:0] internal_data;
7     // Shift register operation
8     always @(posedge clk) begin
9         if (rst) begin
10            internal_data <= {WIDTH{1'b0}};
11        end else if (enable) begin
12            internal_data <= data_in;
13        end else if (scan_enable) begin
14            internal_data <= {internal_data[WIDTH-2:0],
15                scan_in};
16        end
17    end
18    // Output assignment
19    assign data_out = internal_data;
20    assign scan_out = internal_data[WIDTH-1];
21 endmodule

```

1 Unfortunately, it doesn't work if the WIDTH is set to 1.
2 Can you fix it?

(a) User request to fix an error in the shift register

```

1 Sure! The issue arises from the part where you define the
2 concatenation of the internal_data register when
3 scan_enable is high. When the WIDTH is set to 1, the
4 range internal_data[WIDTH-2:0] will cause an out-of-
5 bounds access to the internal_data register.
6
7 To fix this, you can use a conditional expression to handle
8 the case where WIDTH is set to 1. Here's the updated
9 code:

```

```

1 module shift_register #(parameter WIDTH = 8)(
2     // Ports omitted for space
3 );
4     reg [WIDTH-1:0] internal_data;
5     // Shift register operation
6     always @(posedge clk) begin
7         if (rst) begin
8            internal_data <= {WIDTH{1'b0}};
9        end else if (enable) begin
10            internal_data <= data_in;
11        end else if (scan_enable) begin
12            if (WIDTH == 1) begin
13                internal_data <= scan_in;
14            end else begin
15                internal_data <= {internal_data[WIDTH-2:0],
16                    scan_in};
17            end
18        end
19    end
20    // Output assignment
21    assign data_out = internal_data;
22    assign scan_out = internal_data[WIDTH-1];
23 endmodule

```

1 Now, the code will properly handle the case where WIDTH is
2 set to 1.

(b) ChatGPT response to fixing the shift register

Fig. 3. Portion of a conversation (Topic ID 15) asking ChatGPT to fix a bug in one of the Verilog modules that it had previously authored.

D. ISA Development

All chat logs are provided in the data repository [13]. The ISA co-generated with ChatGPT-4 in Conversation 00 (and updated in 10) is presented in Table II. It is a relatively straightforward accumulator-based design with some notable features: (1) given the size constraints, the memory-access ‘Instructions with Variable-Data Operands’ use just five bits to specify the memory address, meaning the processor would be limited to an absolute maximum of 32 bytes of memory. (2) There is just one instruction with an immediate-data encoding. (3) The instructions use the full 256 possible byte encodings. (4) The JSR instruction makes it possible to implement

This looks excellent. According to this list, please produce the module definition for a control unit in Verilog which could operate the processor datapath. Please comment the purpose of each I/O. If a signal is for controlling a multiplexer, please also comment what each possible value should correspond to in the datapath.

Fig. 4. The most difficult prompt (10 restarts), which was provided in Topic 04 after ChatGPT-4 produced a list of datapath control signals and definitions.

```

1 module control_unit (
2     // Clock input
3     input wire clk ,
4     // Reset input
5     input wire rst ,
6     // Input from the Instruction Register (IR)
7     input wire [7:0] instruction ,
8     // Enables writing to the PC
9     output wire PC_write_enable ,
10    // Selects the input for the PC multiplexer
11    // 00: PC + 1 (FETCH cycle, BEQ_FWD, BNE_FWD)
12    // 01: ACC (JMP, JSR)
13    // 10: PC - 3 (BEQ_BWD, BNE_BWD)
14    output wire [1:0] PC_mux_select ,
15    // Enables writing to the ACC
16    output wire ACC_write_enable ,
17    // Selects the input for the ACC multiplexer
18    // 00: ALU output
19    // 01: Memory contents (LDA, LDAR)
20    // 10: PC (JSR)
21    output wire [1:0] ACC_mux_select ,
22    // Enables loading new instruction into IR from memory
23    output wire IR_load_enable ,
24    // Control signal specifying the ALU operation
25    output wire [3:0] ALU_opcode ,
26    // Selects input B for the ALU multiplexer
27    // 0: Memory contents (ADD, SUB, AND, OR, XOR)
28    // 1: Immediate (ADDI)
29    output wire ALU_inputB_mux_select ,
30    // Enables writing to memory (STA)
31    output wire Memory_write_enable ,
32    // Selects input for memory address multiplexer
33    // 00: IR[4:0] (LDA, STA, ADD, SUB, AND, OR, XOR)
34    // 01: ACC (LDAR)
35    // 10: PC (Instruction fetching)
36    output wire [1:0] Memory_address_mux_select
37 );

```

Fig. 5. Code produced by ChatGPT-4 for difficult prompt (11th attempt). It is still missing some I/O, corrected by later messages.

subroutine calls, albeit a little awkwardly (there’s no stack pointer). (5) The branch instructions are restrictive but useful. Skipping two instructions backwards allows for efficient polling (e.g. load an input, mask it for relevant bit, then check if 0 or not). Skipping 3 instructions forwards allows to skip over the instructions needed for a JMP or JSR. These were designed over a number of iterations, including a later modification (Conversations 10-12, the ‘branch update’) which increased the jump forwards from 2 instructions to 3 after during simulation we realized that we could not easily encode JMP/JSR in just 2 instructions. (5) The LDAR instruction allows for pointer-like dereferences for memory loads. This enabled us to efficiently use a table of constants in our memory map (added in Conversation 17) to convert binary values into LED patterns for a 7-segment display.

IV. RESULTS: PROCESSOR IMPLEMENTATION

The processor datapath was assembled in Conversation 08 and is illustrated in Figure 6. The von Neumann design (shared memory for data and instructions) necessitated a 2-state multi-cycle control unit (‘FETCH’ and ‘EXECUTE’). A third ‘HALT’ state is entered after reaching a HLT instruction (reset to exit). ‘HALT’ also sets a processor_halted output flag. Notably, because the ‘FETCH’ state also increments the PC register, the branch instructions in the ISA require ‘-3’ and ‘+2’ modifiers. The Memory Bank is

TABLE II
ISA CO-CREATED WITH CHATGPT-4 (USES ALL 256 ENCODINGS).

Instruction	Description	Opcode
Instructions with Immediate Operands		
ADDI	Add immediate to Accumulator	1110XXXX
Instructions with Variable-Data Operands		
LDA	Load Accumulator with memory contents	000MMMMM
STA	Store Accumulator to memory	001MMMMM
ADD	Add memory contents to Accumulator	010MMMMM
SUB	Subtract memory contents from Accumulator	011MMMMM
AND	AND memory contents with Accumulator	100MMMMM
OR	OR memory contents with Accumulator	101MMMMM
XOR	XOR memory contents with Accumulator	110MMMMM
Control and Branching Instructions		
JMP	Jump to memory address	11110000
JSR	Jump to Subroutine (save address to ACC)	11110001
BEQ_FWD	Branch if ACC==0, forward (PC = PC + 3)	11110010
BEQ_BWD	Branch if ACC==0, backward (PC = PC - 2)	11110011
BNE_FWD	Branch if ACC!=0, forward (PC = PC + 3)	11110100
BNE_BWD	Branch if ACC!=0, backward (PC = PC - 2)	11110101
HLT	Halt the processor until reset	11111111
Data Manipulation Instructions		
SHL	Shift Accumulator left	11110110
SHR	Shift Accumulator right	11110111
SHL4	Shift Accumulator left by 4 bits	11111000
ROL	Rotate Accumulator left	11111001
ROR	Rotate Accumulator right	11111010
LDAR	Load Accumulator via indirect mem. access M[ACC]	11111011
DEC	Decrement Accumulator	11111100
CLR	Clear (Zero) Accumulator	11111101
INV	Invert (NOT) Accumulator	11111110

parameterized globally, allowing the human engineer to change the memory size from inside the Tiny Tapeout wrapper (the only file they authored, which was used to perform non-processor-related wiring). The processor was eventually synthesized with 17 bytes of register memory, with the 17th byte used for I/O (7-segment LED outputs, one button input). A look-up constant memory table of 10 bytes used for segment patterns was concatenated. After synthesis, the processor results in the ‘GDSII’ in Figure 7.

The synthesis for Tiny Tapeout was done using OpenLane [21], which produced the completed GDSII to be added to the full chip, also provided static timing and power analysis through OpenSTA [22]. The initial timing constraint for the design was set to use a 50kHz clock, and the reported timing slack allowed us to calculate the expected maximum clock frequency. This potential clock frequency only applies to this design, though, as the full Tiny Tapeout chip is restricted by its use of a scan-chain based method to select each design on the chip, so it was ultimately taped out with the slower clock frequency.

We also synthesized this processor for a Digilent Cmod A7 (XC7A35T) FPGA development board using Xilinx Vivado and

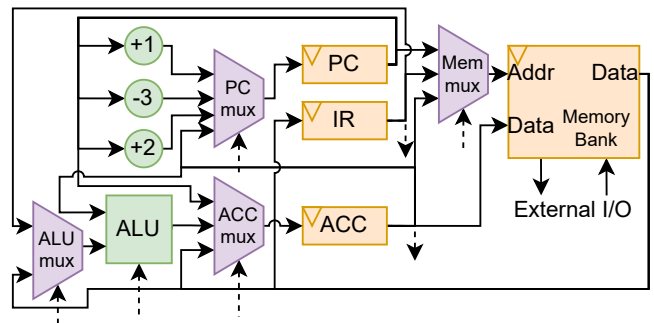
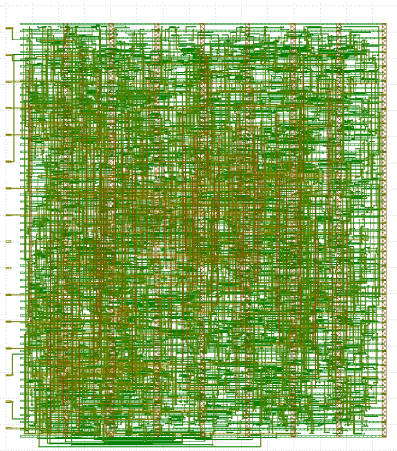


Fig. 6. Accumulator-based datapath designed by GPT-4 (illustration by human). Control signals indicated with dotted lines.



Component	Count
Comb. Logic	999
Diode	4
Flip Flops	168
Buffer	126
Tap	300

Above: (a) Components.

Left: (b) Final processor GDSII render by 'klayout', I/O ports on left side, grid lines = 0.001 um.

Fig. 7. ASIC processor synthesis information.

TABLE III
ESTIMATED MAXIMUM CLOCK FREQUENCIES AND POWER CONSUMPTIONS FOR ASIC AND FPGA IMPLEMENTATIONS

	Approx. F_{max}	Power Consumption
TT03 ASIC	125kHz ^a	7×10^{-7} W
Cmod A7	114MHz	8.9×10^{-3} W

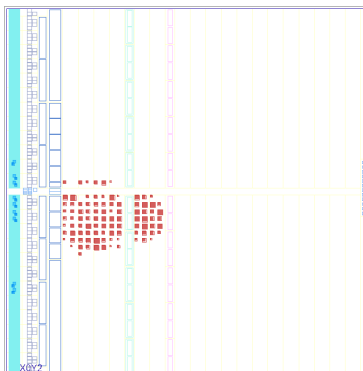
^aThis number is much lower than one might expect. This is due to Tiny Tapeout 3, which places all I/O in a relatively slowly clocked global scan chain and as such considerably constrains the operating clock rates.

recorded the potential maximum clock frequency and power estimates as implemented for that device.

The power and timing estimates for both devices are given in Table III. We found that processor is able to be clocked several orders of magnitude faster on the FPGA than with the ASIC implementation. This is due to the additional constraints used to ensure the tile is compatible with Tiny Tapeout's scan-chain. The FPGA has significantly more power overhead, with 0.072W of the total estimated 0.089W of power consumption coming from the static device power. Figure 8 shows the FPGA utilization taking a small part of a single clock region of the Artix7 device, as well as the components used to implement the processor, as compared to the layout and component usage shown with the GDSII image in Figure 7.

A. Observations

In general, ChatGPT-4 produced relatively high-quality code, as can be seen by the short verification turnaround. Once the Python



Component	Count
LUTs	268
Registers	173
Muxes	31
Bonded IOB	12

Above: (a) Components.

Left: (b) FPGA component layout on a single clock region of the XC7A35T device, highlighted in red.

Fig. 8. FPGA processor synthesis information.

assembler was written (Conversation 09), the bug-fixing Conversations (10-13, 15-16) used just 19 of the total 125 messages. Given the 25 messages per 3 hours rate limit by ChatGPT-4, the total time budget for this design was 22.8 hours of ChatGPT-4 (including the restarts). The actual generation averaged around 30 seconds per message: with no rate limit the whole design *could* have been completed in <100 minutes, subject to the human engineer. Although ChatGPT-4 produced the Python assembler with relative ease, it struggled to author programs written for our design, and no non-trivial test programs were written by ChatGPT. Overall, we evaluated all 24 instructions across a series of comprehensive human-authored assembly programs both in simulation and FPGA emulation.

V. EVALUATION

A. Discussion

Steps for practical adoption: Ideally with the rise of conversational LLMs it would be possible to go from idea to functional design with minimal effort. Although much emphasis has been placed on their single-shot performance (i.e., making a design in a single step) we found here that for hardware applications that they may function better as a *design assistant*, rather than a designer in and of itself. Where they work in lock-step with an experienced engineer, they may serve as an effort 'force multiplier', providing 'first-pass' designs which may then be tweaked and quickly iterated over.

A notable observation is that the outcome of a conversation depends heavily on the early interactions: the response to the initial prompt and first few instances of feedback. As a result, we recommend evaluating the responses to the early prompts, and if they are unsatisfactory, consider 'restarting' the conversation from an earlier point. This was done several times throughout the conversations creating the processor, as errors became evident several messages after initially emitted.

Security analysis: As a relatively feature-limited 8-bit microprocessor lacking modern security features, such as secure enclaves or memory protection, there are not many aspects of security encoded into the design of the processor itself. Still, no matter the design, security-relevant weak design patterns may be present in the Verilog. As using LLMs to produce this Verilog was the main exercise of this work, we perform an analysis of the HDL using the CWEAT tool [23], which can scan for 6 common 'hardware CWEs' [24]. It found no weaknesses, indicating that the Verilog has reached at least some basic level of quality. This is pertinent given [16] which found that language models like GitHub Copilot emit these kinds of CWEs.

Challenges for functional verification: During the course of this case study, we performed both specification design and functional implementation with the aid of ChatGPT. However, when we attempted to perform verification of the design, ChatGPT was repeatedly unable to produce plausible (or in many cases, even compilable) testbenches, test scripts, and test programs. We hypothesize that this is due to a lack of suitable open-source training data for well written Verilog testbenches. As a testbench is often very tightly coupled to the design being tested, it could be especially challenging for a LLM to provide salient testcases for a novel design of this nature.

B. Threats to Validity

Reproducibility: As the conversational LLMs tested are non-deterministic and generative, the outputs are not consistently reproducible. ChatGPT is closed-source and run remotely, so we are unable to examine the parameters of the models and analyze the method for generating outputs. The conversational nature of these tests hamper the reproducibility, as each user response in the conversation depends

on the previous model response, so slight variations can create substantial changes in the final design. Regardless, we do provide the full conversation logs for result reconstruction [13].

Statistical validity: As the goal of this work was to design hardware conversationally, we did not automate any part of this process, and each conversation needed to be done manually. This limited the scale of the experiments that could be performed, which were also hampered by rate limits and model availability (OpenAI’s ChatGPT-4 still has limited access at time of writing). As a result, the processor designed here may not provide enough data to draw formal statistical conclusions.

VI. CONCLUSIONS

Challenges: While using a conversational LLM to assist in designing and implementing a hardware device can be beneficial overall, it is clear that the technology still needs improvement. The ChatGPT LLM produced errors in aspects of both the specification and implementation, requiring intervention by the experienced hardware designer. It seems unlikely, then, that the model could produce designs without assistance (i.e. in the zero-shot setting). Further, we observed deficiencies when attempting to use the model for producing verification code.

Opportunities: Still, when a human is paired with ChatGPT-4, the language model seems to be a ‘force multiplier’, allowing for rapid design space exploration and iteration. We demonstrated this in our case study, where it helped to architect and implement a novel processor. In general, we observed that ChatGPT-4 could produce functionally correct code, which in general could free up designer time when implementing common modules and thus improving developer productivity. Potential future work would involve a larger user study to investigate this potential, as well as the development of conversational LLMs specific to hardware design to improve upon the results.

REFERENCES

- [1] A. B. Kahng, “Machine Learning Applications in Physical Design: Recent Results and Directions,” in *Proceedings of the 2018 International Symposium on Physical Design*, ser. ISPD ’18. New York, NY, USA: Association for Computing Machinery, Mar. 2018, pp. 68–73. [Online]. Available: <https://doi.org/10.1145/3177540.3177554>
- [2] C. Yu, H. Xiao, and G. De Micheli, “Developing synthesis flows without human knowledge,” in *Proceedings of the 55th Annual Design Automation Conference*, ser. DAC ’18. New York, NY, USA: Association for Computing Machinery, Jun. 2018, pp. 1–6. [Online]. Available: <https://doi.org/10.1145/3195970.3196026>
- [3] G. Huang, J. Hu, Y. He, J. Liu, M. Ma, Z. Shen, J. Wu, Y. Xu, H. Zhang, K. Zhong, X. Ning, Y. Ma, H. Yang, B. Yu, H. Yang, and Y. Wang, “Machine Learning for Electronic Design Automation: A Survey,” *ACM Transactions on Design Automation of Electronic Systems*, vol. 26, no. 5, pp. 40:1–40:46, Jun. 2021. [Online]. Available: <https://doi.org/10.1145/3451179>
- [4] G. Dessouky, D. Gens, P. Haney, G. Persyn, A. Kanuparthi, H. Khattri, J. M. Fung, A.-R. Sadeghi, and J. Rajendran, “HardFails: Insights into Software-Exploitable Hardware Bugs,” 2019, pp. 213–230. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity19/presentation/dessouky>
- [5] P. Coussy and A. Morawiec, *High-level synthesis*. Springer, 2010, vol. 1.
- [6] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. d. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, A. Ray, R. Puri, G. Krueger, M. Petrov, H. Khlaaf, G. Sastry, P. Mishkin, B. Chan, S. Gray, N. Ryder, M. Pavlov, A. Power, L. Kaiser, M. Bavarian, C. Winter, P. Tillet, F. P. Such, D. Cummings, M. Plappert, F. Chantzis, E. Barnes, A. Herbert-Voss, W. H. Guss, A. Nichol, A. Paino, N. Tezak, J. Tang, I. Babuschkin, S. Balaji, S. Jain, W. Saunders, C. Hesse, A. N. Carr, J. Leike, J. Achiam, V. Misra, E. Morikawa, A. Radford, M. Knight, M. Brundage, M. Murati, K. Mayer, P. Welinder, B. McGrew, D. Amodei, S. McCandlish, I. Sutskever, and W. Zaremba, “Evaluating Large Language Models Trained on Code,” Jul. 2021, arXiv:2107.03374 [cs]. [Online]. Available: <http://arxiv.org/abs/2107.03374>
- [7] GitHub, “GitHub Copilot · Your AI pair programmer,” 2021. [Online]. Available: <https://copilot.github.com/>
- [8] S. Thakur, B. Ahmad, Z. Fan, H. Pearce, B. Tan, R. Karri, B. Dolan-Gavitt, and S. Garg, “Benchmarking Large Language Models for Automated Verilog RTL Code Generation,” Dec. 2022, arXiv:2212.11140 [cs]. [Online]. Available: <http://arxiv.org/abs/2212.11140>
- [9] OpenAI, “Introducing ChatGPT,” Nov. 2022. [Online]. Available: <https://openai.com/blog/chatgpt>
- [10] S. Pichai, “An important next step on our AI journey,” Feb. 2023. [Online]. Available: <https://blog.google/technology/ai/bard-google-ai-search-updates/>
- [11] A. Ahmad, M. Waseem, P. Liang, M. Fehmideh, M. S. Aktar, and T. Mikkonen, “Towards Human-Bot Collaborative Software Architecting with ChatGPT,” Feb. 2023, arXiv:2302.14600 [cs]. [Online]. Available: <http://arxiv.org/abs/2302.14600>
- [12] M. R. King and chatGPT, “A Conversation on Artificial Intelligence, Chatbots, and Plagiarism in Higher Education,” *Cellular and Molecular Bioengineering*, vol. 16, no. 1, pp. 1–2, Feb. 2023. [Online]. Available: <https://doi.org/10.1007/s12195-022-00754-8>
- [13] J. Blocklove, S. Garg, R. Karri, and H. Pearce, “Data Repository for Chip-Chat: Challenges and Opportunities in Conversational Hardware Design,” May 2023. [Online]. Available: <https://zenodo.org/records/7953725>
- [14] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, “Language Models are Few-Shot Learners,” in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 1877–1901. [Online]. Available: <https://proceedings.neurips.cc/paper/2020/file/1457c0d6b6fcb4967418bfb8ac142f64a-Paper.pdf>
- [15] H. Pearce, B. Tan, and R. Karri, “DAVE: Deriving Automatically Verilog from English,” in *2020 ACM/IEEE 2nd Workshop on Machine Learning for CAD (MLCAD)*, Nov. 2020, pp. 27–32.
- [16] H. Pearce, B. Ahmad, B. Tan, B. Dolan-Gavitt, and R. Karri, “Asleep at the Keyboard? Assessing the Security of GitHub Copilot’s Code Contributions,” in *2022 IEEE Symposium on Security and Privacy (SP)*, May 2022, pp. 754–768, iSSN: 2375-1207.
- [17] H. Pearce, B. Tan, B. Ahmad, R. Karri, and B. Dolan-Gavitt, “Examining Zero-Shot Vulnerability Repair with Large Language Models,” in *2023 IEEE Symposium on Security and Privacy (SP)*, May 2023, pp. 2339–2356, iSSN: 2375-1207. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/10179324>
- [18] RapidSilicon, “RapidGPT,” 2023. [Online]. Available: <https://rapidsilicon.com/rapidgpt/>
- [19] Y. Wang, Y. Kordi, S. Mishra, A. Liu, N. A. Smith, D. Khashabi, and H. Hajishirzi, “Self-Instruct: Aligning Language Model with Self Generated Instructions,” Dec. 2022, arXiv:2212.10560 [cs]. [Online]. Available: <http://arxiv.org/abs/2212.10560>
- [20] M. Tabachnyk and S. Nikolov, “ML-Enhanced Code Completion Improves Developer Productivity,” Jul. 2022. [Online]. Available: <http://ai.googleblog.com/2022/07/ml-enhanced-code-completion-improves.html>
- [21] “OpenLane,” May 2023, original-date: 2020-07-20T19:35:02Z. [Online]. Available: <https://github.com/The-OpenROAD-Project/OpenLane>
- [22] “Parallax Static Timing Analyzer,” Jul. 2023, original-date: 2018-09-28T15:45:57Z. [Online]. Available: <https://github.com/The-OpenROAD-Project/OpenSTA>
- [23] B. Ahmad, W.-K. Liu, L. Collini, H. Pearce, J. M. Fung, J. Valamehr, M. Bidmeshki, P. Sapiecha, S. Brown, K. Chakrabarty, R. Karri, and B. Tan, “Don’t CWEAT It: Toward (CWE) (A)nalysis (T)echniques in Early Stages of Hardware Design,” in *IEEE/ACM 2022 International Conference on Computer-Aided Design*, San Diego, CA, Oct. 2022, (Accepted).
- [24] T. M. Corporation, “CWE - CWE-1194: Hardware Design (4.1),” 2022. [Online]. Available: <https://cwe.mitre.org/data/definitions/1194.html>