

Beginners Guide to Using Al for Hardware Engineers

Presented by Doug Smith



Where do we start?



Robotics

Deals with the design, construction, operation, and application of robots



Deep Learning

A subset of machine learning that uses neural networks with multiple layers to model and solve complex problems

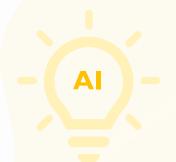


Knowledge Representation and Reasoning

Focuses on representing knowledge in a form that a computer system can utilize to solve complex tasks

Software systems that emulate

the decision-making ability of a human expert in a particular



Neural Networks

Computer systems modeled after the neural connections in the human brain



Cognitive Computing

Mimics human thought processes and augments human cognition



Focuses on the interaction between computers and human languages



Natural Language Processing (NLP)

How do I use AI in engineering?

Our focus

Only a subset of Al

How does Al work?

How do I work

with Al?



domain

Expert Systems

Computer Vision

Involves enabling machines to interpret and make decisions based on visual data

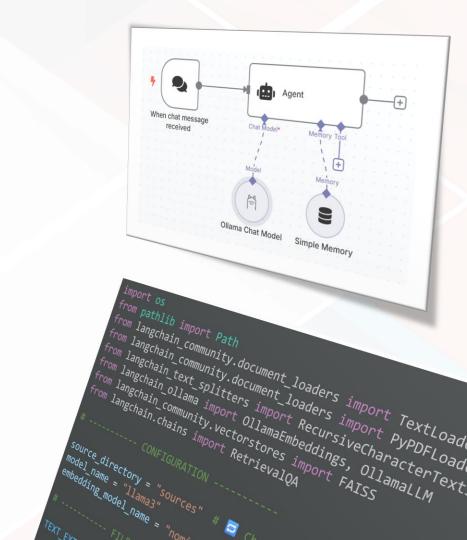


Machine Learning

The study of computer algorithms that improve automatically through experience and data



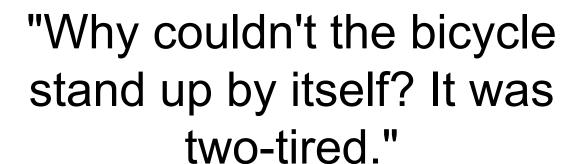
How does Al work?





+ Write a funny joke.





How can a computer write a joke?

The fundamental problem

1010 1010

Computers speak in numbers, not English.



Human language has nuance and context

Natural Language Processing (NLP)





[10445, 7846, 956, 279, 36086, 2559, 709]

Machine language has numbers and vectors

*

Chopping up language

Al parses words into tokens.

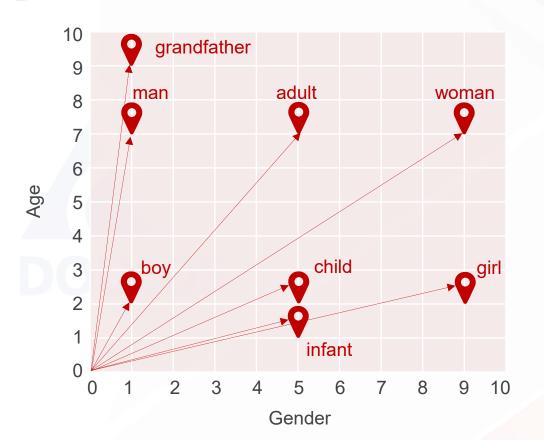


Tokens



A **token** is a unit of text, like an entry in a dictionary, but it has no **meaning**.

Embeddings

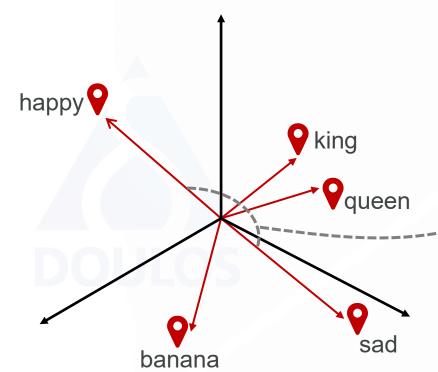


Gender Age

boy	(1,2)	woman	(9,7)
man	(1,7)	girl	(9,2)
adult	(5,7)	infant	(5,1)
child	(5.2)	grandfather	(1,9)

An embedding is a token's address into 'meaning space'.

Word relationships



Angle and direction between vectors encode similarity and relationships

Score	Meaning	Similarity		
1	Same direction	High	king & queen	
0	Unrelated	None	king & banana	
-1	Opposite direction	Low	happy & sad	



She rode the train.

She trained the dog.

(verb)

(noun)

She broke the record.

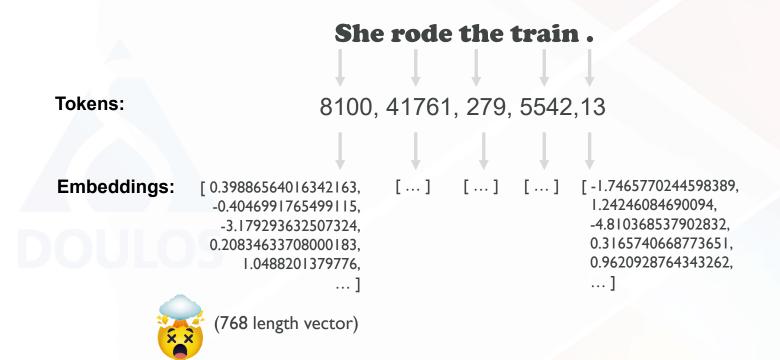
She broke the record player.

(related)

"Attention" means looking at a word's context in a sentence.

opyright © 2025 Doulos – All Rights Reserv

Meaning has many dimensions



1

Large language models

She rode the ...

Embeddings



(Transformer)

LLM

Predicted probability

train	0.08003556728363037
bus	0.07108145207166672
bike	0.024848591536283493
elevator	0.019790327176451683
back	0.01585303619503975



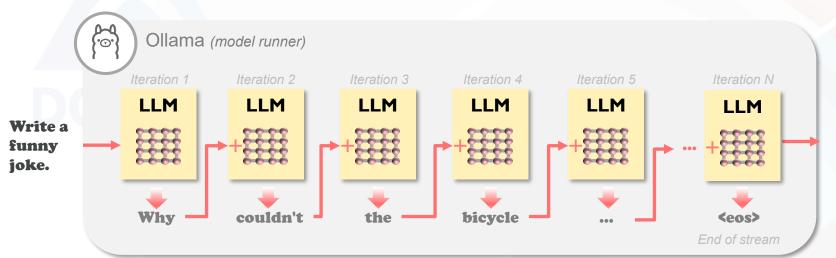
She rode the train

Achieving thought

+ Write a funny joke.

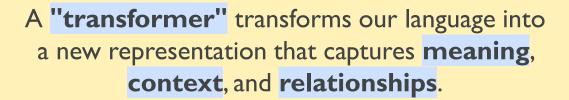
How?

"Why couldn't the bicycle stand up by itself? It was two-tired."

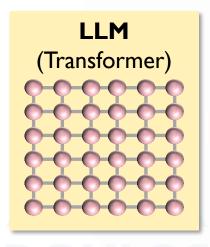


"Why couldn't the bicycle stand up by itself? It was two tired."

Generative Pre-trained = GPT Transformer



Hallucinations





train	0.08003556728363037
bus	0.07108145207166672
bike	0.024848591536283493
elevator	0.019790327176451683
back	0.01585303619503975



She rode the train

(temp = 0.2)



She rode the back

(temp = 0.8)

'Hallucination" is when an output is valid but factually incorrect based on its training data or input context

smaller → more accurate temperature larger → more creative picks from top k choices

picks from at least p





GPT4All Llama 3.1 8B Instruct 128k

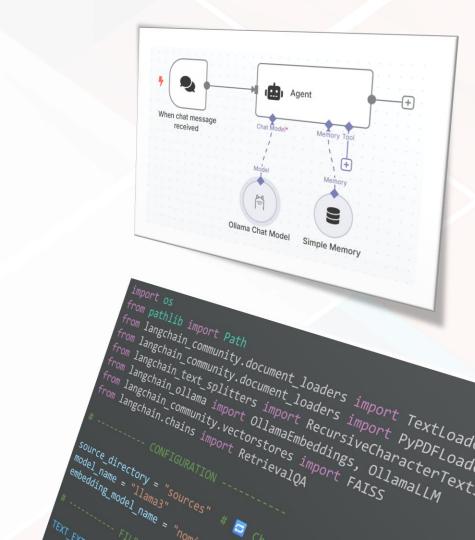
Here are two simple sentences that use the word "cloud" with different meanings:

- 1. The fluffy white cloud floated across the sky.
- 2. The company's profits were in the red, and they had to borrow money from a bank to cover their expenses.

In this example, "cloud" refers to a collection of water vapor or gas in the atmosphere (sentence 1), while it also means an amount owed that is not yet paid (sentence 2).



Working with Al





Key factors of an LLM:

Parameters
Licensing
Tool support
GPUs



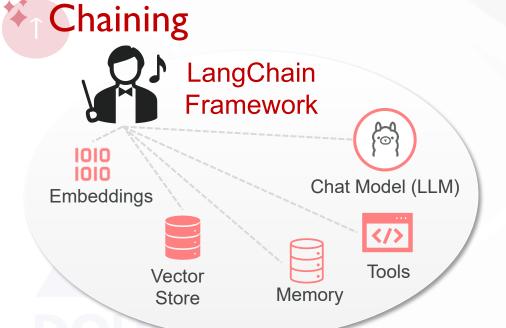
Parameters are weights

More weights ...

- \rightarrow more brain power
- → more memory
- → more GPUs

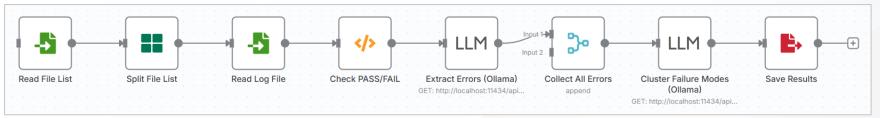
Model	Parameters	License	Originator	Tool Support	Commercial Use	GPU & RAM Requirements
OpenLLaMA 3.1 / 3.2 (Tool-enabled)	8B – 70B	OpenLLaMA Community	Meta Al / OpenLLaMA	✓ Yes	✓ Yes	1× NVIDIA A100 40GB or RTX 3090 24GB; 32–64GB RAM
Mistral / Mixtral (Instruct / Tool)	7B, 8×7B	Apache 2.0	Mistral Al	✓ Yes	✓ Yes	1× NVIDIA A100 40GB or RTX 3090 24GB; 32–64GB RAM
Falcon (Tool / Instruct)	7B, 40B, 180B	Apache 2.0	TII (UAE)	✓ Yes	✓ Yes	2× NVIDIA A100 40GB; 64– 128GB RAM
GPT-NeoX / GPT-J (Tool)	6B, 20B (NeoX), 6B (J)	Apache 2.0	EleutherAl	✓ Yes	✓ Yes	2× NVIDIA A100 40GB; 64– 128GB RAM
DeepSeek-LLM / V3 / R1	7B, 67B, 236B dense; 671B total (MoE, 37B active)	DeepSeek License	DeepSeek (China)	Yes (native tool reasoning)	✓ Yes	2× NVIDIA A100 40GB or RTX 3090 24GB; 64–128GB RAM
Gemma (1–3 series, Tool)	1B – 27B	Gemma License (open- weights)	Google DeepMind	✓ Yes	✓ Yes	1× NVIDIA A100 80GB or RTX 4090 24GB; 32–64GB RAM
Grok (1, 2.5)	~314B (dense est.), details vary	Apache 2.0	xAl (Elon Musk / X)	✓ Yes	✓ Yes	8× NVIDIA A100 40GB; 512GB RAM
OpenAl GPT-4 (API / Plugins)	175B+	OpenAl	OpenAl	✓ Yes	Yes (via API)	5× NVIDIA A100 80GB; 400GB+ RAM
OpenAl GPT-4-turbo (Function Calling)	175B+	OpenAl	OpenAl	✓ Yes	Yes (via API)	5× NVIDIA A100 80GB; 400GB+ RAM
GPT-OSS (20B / 120B)	20B, 120B	Apache 2.0	OpenAl	✓ Yes	✓ Yes	1× NVIDIA A100 80GB (120B); 1× RTX 4090 24GB (20B); 32–64GB RAM
Qwen (Tool-enabled variants)	3B, 7B, 14B, 72B	Apache 2.0	Alibaba Cloud	✓ Yes	Yes (with conditions)	1× NVIDIA A100 40GB or RTX 3090 24GB; 32–64GB RAM





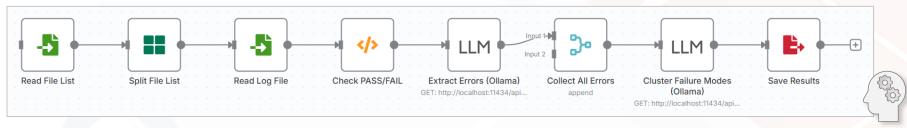
A "chain" is a sequence of steps that process input data through components like prompts, models, and memories.

An example chain of steps



Workflows and agents

Workflow

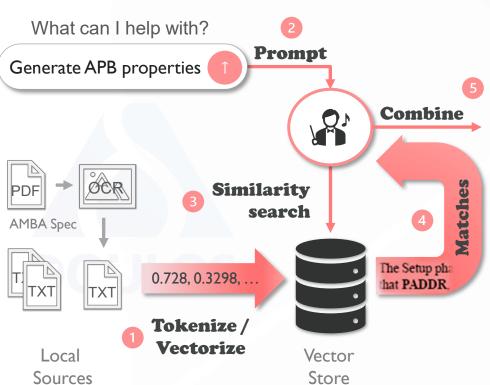


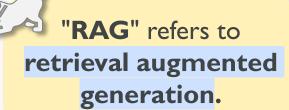
Agent (Agentic)

Al decides

An "agent"
independently decides
what it will do next.

Document stores (RAG)





Generate APB properties



The Setup phase of the write to that PADDR, PWRITE, and I

The Access phase of the write asserted by the Completer at the PWDATA, and any other cont

At the end of the transfer, PEN same peripheral.





Ouery LLM



New Prompt





Prompting



Task – always start a task with a generate verb



Context – ask yourself 3 questions:

What's the user's background? What does success look like? What environment are you in?



Exemplars – give examples of good and bad templates



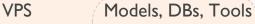
Persona – who do you want the Al to be?



Format – how does the result look like?



Tone – indicate the feeling – casual, formal, witty, etc.



OpenLLaMA 3.1 / 3.2 (Tool-enabled)

Mistral / Mixtral (Instruct / Tool)

http://

http://

Falcon (Tool / Instruct)

GPT-NeoX / GPT-J (Tool)

DeepSeek-LLM / V3 / R1

Gemma (1-3 series, Tool)

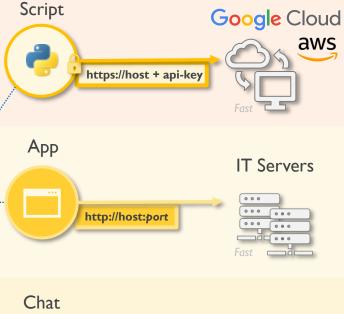
Grok (1, 2.5)

OpenAl GPT-4 (API / Plugins)

OpenAl GPT-4-turbo (Function Calling)

GPT-OSS (20B / 120B)

Qwen (Tool-enabled variants)





LLawin

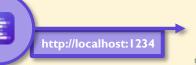
https://url

- api-key

*****Claude









Desktop



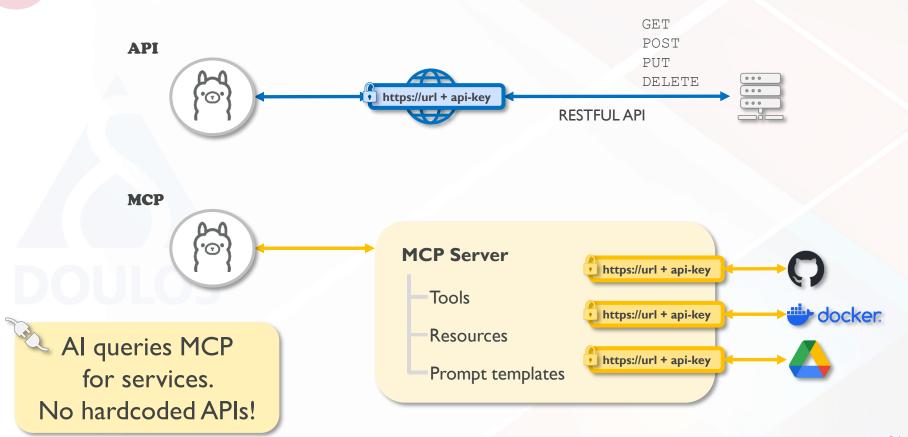
Beware! Data is not private!





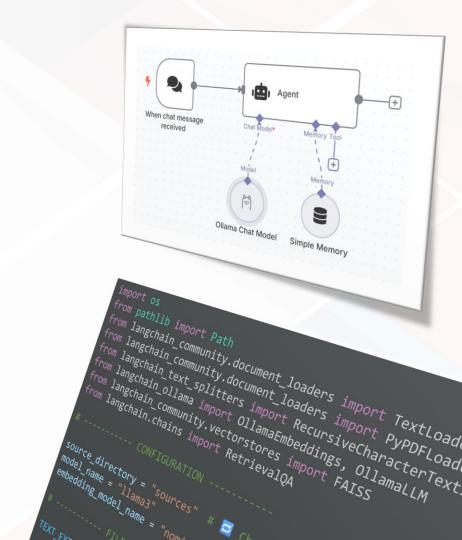
http://

Model Context Protocol (MCP)





Al for Engineering







Analytical / Assistive AI

	Ge	nerative AI	Bug / Issue Detection	Analyze RTL/testbench code to find		
	RTL/HDL Generation	Generate synthesizable RTL code from high-level descriptions + specifications	Performance / Timing Estimation	bugs or inconsistencies Predict timing paths or bottlenecks based on RTL patterns		
	Testbench Generation	Create verification environments, stimulus, and test vectors	Design Rule Checking	Check RTL against coding standards or design rules		
ביים אמם אמם	SVA/Assertion Generation	Create SystemVerilog Assertions (SVA) to cover corner cases	Test Coverage Analysis	Evaluate coverage metrics and suggest missing tests		
	Specification Summarization	Read long design specs and produce concise summaries or design outlines	Debug Assistance	Analyze simulation logs and highlight probable failure causes		
11 @ 2023 Dou	Documentation / Commenting	Generate design documentation, comments, or explanations for existing RTL	Specification Consistency Checking	Compare design spec vs RTL/testbench to find gaps or contradictions		

Where do I begin?

Ask Chat

What's on the agenda today?

+ I want to use AI on my local machine. Where do I begin?





Chat with AI to get started using AI.

+ How do I install a local LLM?

1

+ How do I use it in docker?



+ Write a script using Ollama



+ How do I enable my GPU?



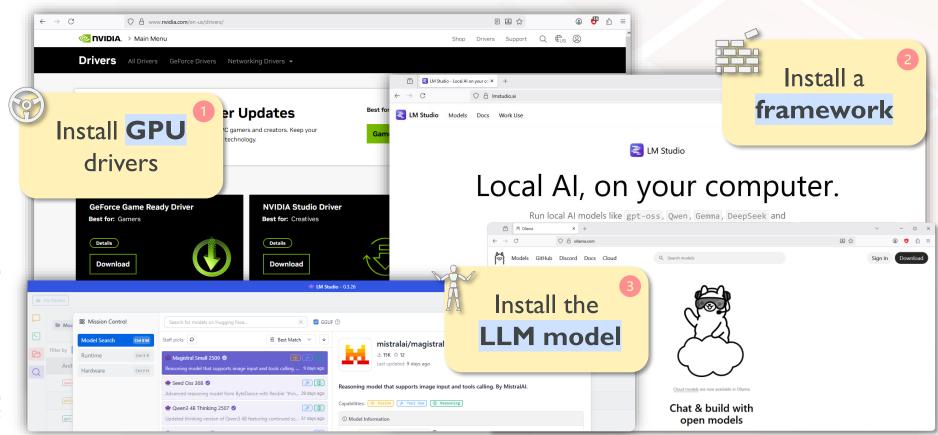
+ Add RAG to my script



+ How do I install an AI model?



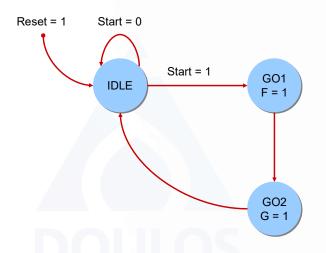
Local LLM install



opyright © 2025 Doulos - All Rights Reserved

*

Experiment with LLMs





I am a hardware engineer and I need to generate a finite state machine written in SystemVerilog. There are 3 states: IDLE, GO1, and GO2. The state machine transitions to IDLE on 'Reset' = 1. When the 'Start' = 1, the state machine transitions to GO1 and the output 'F' = 1. When in the GO1 state, the state machine transitions to GO2 and the output 'G' = 1. The defaults for 'F' and 'G' are 0. When in the GO2 state, the state machine transitions back to IDLE. Generate this state machine as a SystemVerilog module.





gpt-oss:20b





Thought for 91.4 seconds

Below is a compact, synthesizable SystemVerilog implementation of the requested FSM. It uses a two-bit enumerated state type, a synchronous state register, and combinational logic for the next state and the two outputs **F** and **G**.



LLMs are already trained in hardware description languages.

```
С Сору
* Module: simple fsm
* Purpose: 3-state FSM (IDLE → GO1 → GO2 → IDLE)
                                                   * Synchronous state register (reset is asynchronous)
                                                   * -----*/
* Inputs:
      clk : Clock (positive-edge synchronous)
                                                  always_ff @(posedge clk or posedge reset) begin
      reset : Asynchronous active-high reset
                                                     if (reset)
      start : Pulse that moves the FSM out of IDLE
                                                        state <= IDLE:
                                                                           // Return to IDLE on reset
                                                     else
  Outputs:
                                                        state <= next state; // Update state on clock edge
           : 1 when the ESM is in state GO1
                                                  end
           : 1 when the FSM is in state GO2
*********************

    Next-state logic (combinational)

module simple fsm (
                                                                                                               end
   input logic clk,
                    // Clock
                                                   * ------ */
                                                  always comb begin
                   // Asynchronous active-high reset
   input logic reset,
                                                                                                            endmodule
                                                                            // Default: hold current state
                                                     next state = state;
   input logic start.
                   // Start pulse
   output logic F,
                    // Output asserted in GO1
                                                     case (state)
                                                        IDLE: begin
   output logic G
                    // Output asserted in GO2
                                                                                                       How it works
                                                           if (start)
                                                               next state = GO1:
                                                                              // Start -> GO1
                                                        end
   * State encoding
                                                        GO1: next state = GO2; // GO1 -> GO2
   * _____
   typedef enum logic [1:0] {
                                                        GO2: next state = IDLE; // GO2 -> IDLE
      IDLE = 2'b00,
      G01 = 2'b01,
                                                        default: next state = IDLE; // Safety
      G02 = 2'b10
                                                     endcase
   } state_t;
```

- 1. Reset (reset=1) forces the FSM to IDLE .
- 2. While in IDLE, a high start pulse moves the FSM to GO1 and sets F=1.
- 3. In GO1, the FSM automatically goes to GO2 on the next clock and sets G=1.
- 4. In GO2, it goes back to IDLE on the next clock.
- 5. Both F and G default to 0 in all other cases.

The module is fully synthesizable, uses only standard constructs, and can be instantiated in any SystemVerilog design.

povriaht © 2025 Doulos – All Rights Reserved

Look, UVM!

1. Package – common types

```
systemverilog
// file: fsm pkg.sv
package fsm pkg;
    import uvm pkg::*;
    `include "uvm macros.svh"
    // Sequence item that the driver will push into the DUT
    class fsm seg item extends uvm seguence item;
       rand bit start:
                           // stimulus to the DUT
       bit F exp;
                               // expected value of output F
       bit G exp;
                                // expected value of output G
        `uvm object utils(fsm seq item)
       function new(string name = "fsm seq item");
            super.new(name);
        endfunction
    endclass : fsm seg item
    // Sequence - just a simple generator
    class fsm sequence extends uvm sequence #(fsm seq item);
        `uvm object utils(fsm sequence)
       function new(string name = "fsm sequence");
            super.new(name);
```

- 2. Environment agent, driver & monitor
- 3. Scoreboard
- 4. Environment glue all the pieces together
- 5. Test
- 6. Top-level testbench (instantiates DUT and UVM)

How it works

- 1. Test (fsm_test) creates the environment and starts a sequence that toggles s
- The driver applies start to the DUT, waits for the next clock edge, records wheexpected values back to the scoreboard via add_expected().
- 3. The monitor simply raises an event whenever F or G changes in this simple ca

ovright © 2025 Doulos – All Rights Reserve

Summarization



You are an expert hardware engineer. Analyze the spi_clgen, spi_defines, spi_shift, and spi_top modules. Create a list of properties describing the functionality of the spi_top module.

Below is a concise "properties" list that captures what the <code>spi_top</code> wrapper brings together, based on the public interfaces and typical design intent of an SPI controller implemented with the <code>spi_defines</code>, <code>spi_shift</code>, and <code>spi_clgen</code> sub-modules.

Property	Description
Wishbone-compliant interface	wb_clk_i, wb_rst_i, wb_adr_i, wb_dat_i, wb_sel_i, wb_we_i, wb_stb_i, wb_cyc_i are the standard Wishbone signals. The module returns wb_dat_o, wb_ack_o, wb_err_o, and wb_int_o. This allows the controller to be addressed, read from, and written to by any wishbone master (CPU, DMA, etc.).
SPI I/O pins	External SPI pins are exposed as: ss_pad_o (slave-select), sclk_pad_o (serial clock), mosi_pad_o (master-out-to-slave), and miso_pad_i (master-in-from-slave). These correspond to the physical pins on an FPGA or ASIC.
Parameterised clock period	The module has a parameter $T_p = 1$; While not used directly in the snippet, it is typically leveraged by timing constraints or internal calculations for generating the SPI bit-rate.
SDI clock goneration	Internally and allege produces the SCIK salk had a from wholk it based on a programmable divider (dividen)



Scripting a chat bot

```
<u>-</u>
                                                                        Use framework
import ollama
# Choose your model
                                                                          Specify LLM
model_name = "llama3" # change if you have another ollama model
print("=== Simple Ollama Chatbot ===")
print("Type 'quit' to exit.\n")
while True:
    user_input = input("You: ")
                                                                          Get prompt
    if user_input.lower() == "quit":
        break
    # Generate response from Ollama
    response = ollama.generate(
                                                                            Ask LLM
       model=model_name,
        prompt=user_input,
        options={"temperature": 0.7, "max_tokens": 4000}
                                                                         Print response
    print("Bot:", response["response"])
```

Equipping your chat bot

```
RAG
      ----- LOAD DOCUMENTS ---
def load_documents_from_directory(directory):
  ----- SPLIT DOCUMENTS -----
                                      llm = ollamaLLM(model=model_name)
text_splitter = RecursiveCharacterTextS
docs = text_splitter.split_documents(do
                                       lef search_docs(query: str):
                                                                        ----- AGENT ----
print(f"聚Split into {len(docs)} chunks
                                          """Search the local knowledgeagent = initialize_agent(
                                          ga_chain = RetrievaloA.from
                                                                          tools,
 ---- EMBEDDINGS & VECTORSTORE
                                               11m=11m
embeddings = OllamaEmbeddings(model=embe
                                              retriever=retriever.
                                                                          agent=AgentType.ZERO_SHOT_REACT_DESCRIPTION,
faiss_path = "faiss_index"
                                              return_source_documents=1
                                                                          memory=memory,
                                              chain_type_kwargs={"outpu
                                                                          handle_parsing_errors=True,
                                      r memory
                                                                          verbose=True
 ----- MEMORY -----
def load_memory(path):
                                          result = qa_chain.invoke(quer
                                          sources = [doc.metadata.get(\frac{1}{#} ------ INTERFACE LOOP -----
def save_memory(path, messages):
                                      ocuments"]]
                                                                      print("@ Ready! Ask something (type 'exit' to quit):")
                                          return f"{result['result']}\
                                                                      while True:
previous_messages = load_memory(memory_tools = [
                                                                          query = input("You: ")
memory = ConversationBufferMemory(memor
                                          Tool (
                                                                          if query.lower() in ["exit", "quit"]:
                                              name="KnowledgeBaseSearch
                                                                              print("   Exiting. Saving memory...")
memory.chat_memory.messages = previous_r
                                              func=search docs.
                                                                              save_memory(memory_path, memory.chat_memory.messages)
                                              description="Use this to
                                                                               break
Memory
                                                                          response = agent.invoke(guery)
                                                                          print("AI:", response)
                                                Tools
                                                                          save_memory(memory_path, memory.chat_memory.messages)
```

Agentic flow

```
-- Agentic Flow ---
def agentic_flow(prompt):
    print(f"\nPrompt: {prompt}\n")
    # Step 1: Generate initial response
   response = llm_generate(prompt)
print(f"Initial Response: {response}\n"
    # Step 2: Critique the response
    critique = llm_critique(response)
    print(f"Critique: {critique}\n")
    # Step 3: Generate refined response
    refined_response = llm_refine(prompt, c
    print(f"Refined Response: {refined_response
    return refined_response
   -- Main program -
   __name__ == "__main__":
    prompt_text = input("Enter your prompt
    final_output = agentic_flow(prompt_text
```

- Enter your prompt: Write a funny joke
- Prompt: Write a funny joke
- Initial Response: I asked the librarian if the library had books on paranoia. She whispered, "They're right behind you."
 - Critique: Critique of the response
- 3 | "I asked the librarian if the library had books on paranoia. She whispered, They're right behind you.""

Universal

1. Clarity

Aspect	Asses		Refined Response: New Joke			
Language	Very c ambig		"I asked the librarian for a book on reverse psychology. She leaned in and whispered, "Just stay away from it."			
Narrative flow	Seam! natura	Why this or	ne works better			
Tone	Light,	Element	Why it's an improvement			
Potential	None.	Clear setup	The request for a "book on reverse psychology" immediately signals that the answer will be a subversion.			
confusion Verdict: The response	clear,	Tighter punchline	The librarian's reply is a literal echo of reverse-psychology: by telling the person to stay away, she nudges them toward the book.			
2. Complete	eness	Visual humor	The image of a librarian whispering "Just stay away from it" in a library is vivid and unexpected.			
~ -						

A generative adversarial network

(GAN) is where a <u>generator</u> is critiqued by a <u>discriminator</u> to make a better response.

Anyone familiar with the phrase "reverse psychology" (or the idea of a paradoxical

An engineering example

UVM ERROR @ 7165: uvm test top.m env.m wb spi sb [wb spi sb] WB write to an invalid address! Address = 16647

UVM INFO @ 7186: uvm test top.m env.m wb agent.m wb drv [WB DRIVER] wb read task: addr = 0000000c. data = 00000000

UVM INFO @ 7175: uvm test top.m env.m wb agent.m wb mon [wb monitor] Detected a WB transaction ...

Log file triage

```
# UVM ERROR verilog src/uvm-1.1c/src/seq/uvm sequencer base.svh(1236) @ 7096: uvm test top.m env.m wb agent.m wb seqr [SE0FINERR] Parent sequence
'uvm test top.m env.m v sequencer.vseq.dut config' should not finish before locks from itself and descedent sequences are removed. The lock held by the child sequence
'uvm test top.m env.m v sequencer.vseq.dut config' is being removed.
# UVM INFO @ 7096: uvm test top.m env.m wb agent.m wb dry [WB DRIVER] Performing a RX ...
# UVM INFO @ 7115: uvm test top.m env.m wb agent.m wb mon [wb monitor] Detected a WB transaction ...
# UVM INFO @ 7126: uvm test top.m env.m wb agent.m wb drv [WB DRIVER] wb read task: addr = 00000004, data = 00000000
# UVM INFO @ 7127: uvm test top.m env.m wb agent.m wb seqr@@vseq.wb read rx regs.wb read [wb read] Generated sequence:
# UVM INFO @ 7127: uvm test top.m env.m wb agent.m wb drv [WB DRIVER] Performing a RX ...
# UVM INFO @ 7135: uvm test top.m env.m wb agent.m wb mon [wb monitor] Sending transaction to the scoreboard
# uvm sequence item
                            wb trans -
# addr
                            integral 5
                            integral 32
# UVM ERROR @ 7135: uvm test top.m env.m wb spi sb [wb spi sb] SPI error callback. Error selecting the bit to twiddle!"
# UVM INFO @ 7145: uvm test top.m env.m wb agent.m wb mon [wb monitor] Detected a WB transaction ...
# UVM INFO @ 7156: uvm test top.m env.m wb agent.m wb drv [WB DRIVER] wb read task: addr = 00000000, data = 0000050e
# UVM INFO @ 7157: uvm test top.m env.m wb agent.m wb seqr@@vseq.wb read rx regs.wb read [wb read] Generated sequence:
# UVM INFO @ 7157: uvm test top.m env.m wb agent.m wb drv [WB DRIVER] Performing a RX ...
# UVM INFO @ 7165: uvm test top.m env.m wb agent.m wb mon [wb monitor] Sending transaction to the scoreboard
# ------
# Name
# uvm sequence item
                            wb trans -
                                           @1315
                            integral 5
   data
                            integral 32
                                                                                                              Can we use AI to group
```

Can we use AI to group common failure modes?

* Al limitations





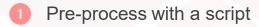




LLM limitations:

- Bad at scale (too many log files)
- Bad at long noisy text
- Distracted by random, test-specific values
- RAG helps retrieval, but ... messy raw data is messy retrieval
- Better at semantic than symbolic analysis

Solution





- Filter random, test-specific values
- Extract error signatures
- Normalize into a compact form:

[FAIL] parity_mismatch in module A

[FAIL] addr_decode_error in module B

[FAIL] timeout waiting for handshake in module C

Cluster and categorize messages



- Use embeddings to group similar errors
- Use sentence level instead of word level
- Return signature groupings
- 3 Summarize with an LLM





Simple reporting script code





Results

```
def export_report(clustered, summaries, out_csv="failure_report.csv"
   rows = []
   for label, sigs in clustered.items():
       summary = summaries.get(label, "Uncategorized")
       rows.append({
            "cluster": label,
            "Summary": summary,
"Num_Failures": len(sigs),
            "Examples": "; ".join(set(sigs[:3]))
   df = pd.DataFrame(rows)
   df.to_csv(out_csv, index=False)
   return df
   <u>__name__</u> == "__main__":
   logs_dir = "./workflow/log_parsing"
   sigs = collect_signatures_from_logs(logs_dir)
   clustered = cluster_signatures(sigs)
   summaries = summarize_clusters(clustered)
   report = export_report(clustered, summaries)
   print(report)
                                                    172,1
                                                                   Bot
```



Results OK, but ... easier (and faster) with a regular script

Output

```
[cluster_signatures] device = cuda
tried eps=0.200 -> clusters=4, labels_set={np.int64(0), np.int64(1), np.int64(2), np.int64(3)}
tried eps=0.300 -> clusters=2, labels_set={np.int64(0), np.int64(1)}
tried eps=0.400 -> clusters=2, labels_set={np.int64(0), np.int64(1)}
tried eps=0.500 -> clusters=1, labels_set={np.int64(0)}
tried eps=0.600 -> clusters=1, labels_set={np.int64(0)}

Chosen clusters: 4 (eps=0.2)
Cluster ... Examples

0  0 ... UVM_ERROR @ <NUM>: uvm_test_top.m_env.m_wb_spi...

1  1 ... UVM_ERROR verilog_src/uvm-<NUM>.<NUM>c/src/seq...
2  2 ... UVM_ERROR : <NUM>; UVM_FATAL : <NUM>; UV...
3  3 ... UVM_ERROR @ <NUM>: uvm_test_top.m_env.m_wb_spi...
[4 rows x 4 columns]
```

CSV

Cluster	Num_Failures	Examples			
0	70	UVM_ERROR @ < NUM>: uvm_test_top.m_env.m_wb_spi_sb [wb_spi_sb] WB write to an	invalid add	dress! Addre	ess = <num></num>
1	20	UVM_ERROR verilog_src/uvm- <num>.<num>c/src/seq/uvm_sequencer_base.svh(<nu< td=""><td>M>) @ <nl< td=""><td>JM>: uvm_te</td><td>st_top.m_env.m_v</td></nl<></td></nu<></num></num>	M>) @ <nl< td=""><td>JM>: uvm_te</td><td>st_top.m_env.m_v</td></nl<>	JM>: uvm_te	st_top.m_env.m_v
2	60	UVM_ERROR: <num>; UVM_FATAL: <num>; UVM_WARNING: <num></num></num></num>			
3	10	UVM_ERROR @ < NUM>: uvm_test_top.m_env.m_wb_spi_sb [wb_spi_sb] SPI error callba	ck. Error s	electing the	bit to twiddle!"

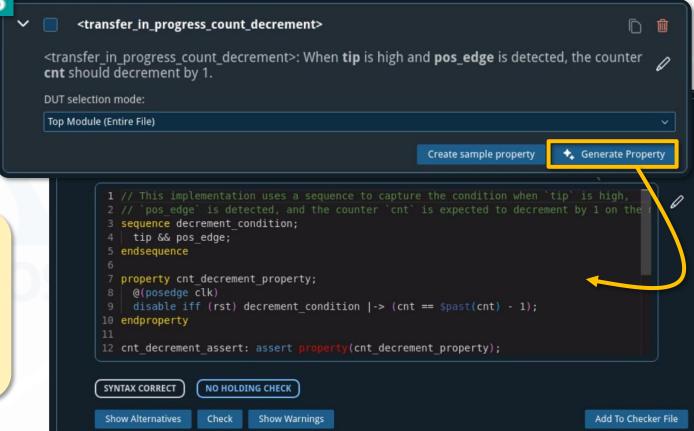


EDA tools

call LLMs

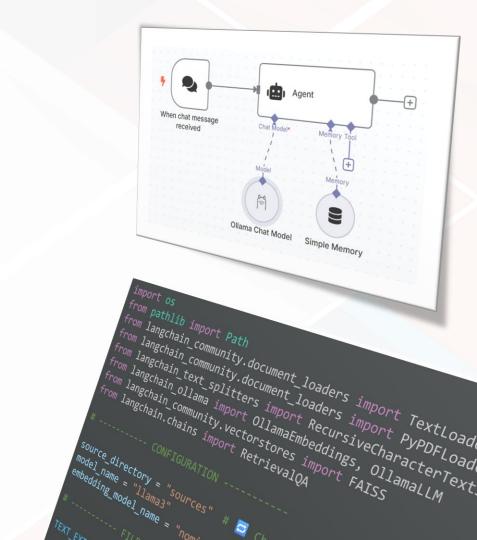
to help

with tasks





Al or not Al?



"Chat, summarize Al"

What it's good at ...

♠ What it's bad at ...

Understanding language

Understanding context deeply

Summarization

- ▲ Common sense reasoning
- Generating content quickly from its training
- Creativity and originality

Automating repetitive tasks

1

Recommendations with Al

X Avoid	i FYI
X Blind trust in its outputs	More resources → better performance
★ AI on non-GPUs & small amounts of GPU memory	Reasoning models mainly help with math & coding
X Lack of human oversight	Quantized models may give faster results
X Non-private or insecure Al	Use Vision Large Language Model (VLM) for text + images
X Over-reliance on Al	LLMs are not great for parsing

Practical uses of Al

Continuous integration support

* Task
Automate code generation
Chat-driven design assistant
Verification automation
Design review and refactoring



How Al Can Help

Generate RTL, testbenches, and assertions from specs.

Convert design intents into RTL and testbenches via chat.

Summarize simulation results and flag anomalies.

Detect issues and suggest improvements in HDL code.

Auto-comment or review on check-in.

References

- I. https://poloclub.github.io/transformer-explainer/
- Vaswani, Ashish; Noam M. Shazeer; et ali. "Attention is All you Need." Neural Information Processing Systems, 2017.
- 3. Kumar, Aman and Deepak Narayan Gadde. <u>Generative Al Augmented Induction-based Formal Verification</u>. 37th IEEE International System-on-Chip Conference, Sep. 16-19, Dresden Germany, 2024.
- 4. Truong, Andy; Daniel Helström; Harry Duque; and Lars Viklund. "Clustering and Classification of UVM Test Failures Using Machine Learning Techniques." DVCon Europe, 2018.
- 5. Ismail, K.A.; Ghany, M.A.A.E. <u>Survey on Machine Learning Algorithms Enhancing the Functional Verification Process</u>. *Electronics* 2021, *10*, 2688. <u>https://doi.org/10.3390/electronics10212688</u>
- 6. Thakur, Shailja; Baleegh Ahmad; et ali. "Benchmarking Large Language Models for Automated Verilog RTL Code Generation." https://arxiv.org/abs/2212.11140.2022.





SoC Design & Verification

FPGA & Hardware Design

Embedded Software & Arm

Python, AI & Machine Learning

» SystemVerilog » UVM » Formal

» SystemC » TLM-2.0

» VHDL » Verilog » SystemVerilog

» Tcl » AMD

» Emb C/C++ » Emb Linux » Yocto » RTOS

» Security » Android » Arm » Rust » Zephyr

» Python » Edge AI » Deep Learning





















opyright © 2025 Doulos - All Rights Reserv

New Design and Verification courses

Advanced Formal Verification

- equips you to tackle complex verification challenges by taking full advantage of formal verification in your engineering projects.
- forms a complete learning path with Essential Formal Verification course, which gives you a solid, practical grounding in formal verification.

SystemVerilog for New Designers: Self-Paced course

- get project-ready for FPGA or ASIC design, including RTL synthesis, block-level test benches, and FPGA design flows.
- high-quality content developed by expert instructors now in self-paced format.

IC Verification with Python and cocotb

➤ learn cocotb principles, constructing different aspects of a verification environment using cocotb and Verification strategies and tactics.







+ Any questions?

1