

Towards Revolutionizing Chip Design with AI: The Integration of AI and Algorithm

Haoxing (Mark) Ren, Director of Design Automation Research, NVIDIA

Al works for many chip design problems

Research Projects @ NVIDIA



VerilogCoder

Analysis, Optimization, and Assistance are the main application scenarios

Assistance (LLM/Agent)							<i>ChipNeMe</i> (Engineerin	clusteringAg g) (Cell)	gent	(RTL) <i>FVAgent</i> (FV)
Know-how							<i>VerilogEv</i> (RTL)	al RTLFixer (RTL)	<i>FVEval</i> (FV)	OPCAssistant (Lithography)
Task automation					<i>VAESA</i> (Arch)	<i>Transsizer</i> (PD)	BufFormer (PD)	Clustering (Cell)		
Optimization						<i>TAG</i> (Analog)	Dream-GAN (PD)	<i>ILILT</i> (Lithography)		
Faster More scalable Better results			<i>NVCell-RL</i> (Cell)	<i>PrefixRL</i> (Synthesis)	G	raph Cluster (PD)	AutoCRAFT RL (Analog)			
		FIST (PD)	<i>ParaSize</i> (Analog)				AutoDMP (PD)			
Analysis (CNN/GNN)		ParaGraph (Parasitics)	<i>MAVIREC</i> (IR Drop)	5	[Lith	D <i>OINN</i> nography)				
Faster Cross-Stage	PRIMAL (Power)	<i>GRANI</i> (Switching	V <i>ITE</i> Activity)							
Detect Anomaly	HPGCN (Testability)	PowerNet (IR Drop)					Coverage Closure (Verification)			
	2019	2020	2021	20	022	20)23	2024		

3 📀 🔁 🛛 3

VerilogCoder

Al predicts cross-stage design metrics

Impact of layout parasitics on schematic design

Use AI to predict layout parasitics from schematic given a large dataset

Convert schematic to graph and learn with GNN

Predicted parasitics reduced simulation error to <10%



Al optimizes tools parameters

Macro placement quality is very important for physical design Placement tools parameters have a huge impact on macro placement Multi-objective Bayesian Optimization (sample and model) find best parameter set Find better macro placement with open-source GPU accelerated placement tools





Al optimizes datapath structure

Datapath synthesis important for GPU

Optimize prefix adder structure with Reinforcement Learning for better delay and area tradeoff Reward computed by open-source EDA tools







PrefixRL achieves better results than well known adder architectures

R. Roy et al, PrefixRL: Optimization of parallel prefix circuits using deep reinforcement learning



AI Generates Optimal Gate Sizes

Timing/power optimization such as gate sizing affects scalability of PD tools Model a path of gates as a sequence, generate optimized gate sizes using Transformer Trained on tool optimized netlist dataset

100X – 1000X speedup compared to traditional optimization with similar PPA



7

AI assists engineers in QnA, coding and analysis tasks

LLM is good at question answering, coding, extraction, rewriting, summarization, classification, reasoning, ...

- Answer questions about designs, infrastructures, tools, flows, HW domains, etc.
- Generate scripts for specific tasks (VLSI)
- Summarize bug reports, predict assignment Domain-adapted model have better performance







Will AI revolutionize Chip Design? Skeptics are saying :

Analysis

"Inaccurate predictions"

Optimization:

"Al is not as good as existing algorithms" "Al will never get better results than the data it trained on"

Assistance

"Hallucinations"

"Can not help real/complex design problems"



AI and Algorithms are good at different things

	Algorithm	AI
Categories	Placement, route, synthesis, CTS, etc	Supervised, unsupervised, reinforcement learnings, etc
Optimality	Known	Unknown
Robustness	Works on any data distribution	Do not work if training and inference distribution mismatch
Training	No	Require a lot data
Interpretability	Behavior explainable	Not explainable
Pros	Solve a known problem efficiently	Solve any problem by learning from complex data
Cons	Oversimplification of dynamic, complex problem	Difficult to leverage the mechanics of the problem



AI vs Algorithm is like System 1 vs System 2 We use both System 1 and System 2 thinkings for hard problems



11 📀 NVIDIA

Daniel Kahneman: Thinking, Fast and Slow

To revolutionize Chip Design with Al Integrate Al with Algorithms

Not a new Idea, Bayesian Optimization/Reinforcement Learning examples won't work without Algorithms,



Tighter integration!



For gaming, AI + Algorithm is much better than AI only

Tight integration of AI with search algorithm helps AI game model perform much better in Poker and Go



Noam Brown: Parables on the Power of Planning in AI: From Poker to Diplomacy



Three ways AI works with Combinatorial Optimization algorithms



Yoshua Bengio et al, Machine Learning for Combinatorial Optimization: a Methodological Tour d'Horizon



An inaccurate AI prediction model still works by acting as an Oracle in an algorithm

Tackle the test point insertion problem to improve design testability

GCN model predicts node testability: Difficult-to-test nodes (DTN) prediction

Model acts as a predictor/oracle for a test point insertion (TPI) algorithm

Although model accuracy is only 90%, reduced test points by 11% and test patterns by 6% to achieve similar coverage as a commercial test point insertion algorithm



Y. Ma et al, High Performance Graph Convolutional Networks with Applications in Testability Analysis

15 📀 🔍 🔂 15

Ground the AI with an algorithm helps AI get better results than the algorithm that generates the training data

Mask optimization is a timing consuming process for every design Learn the optimized mask from ILT (inverse lithography Technology) solver Use Lithography forward simulation algorithm to ground the model Increase inference time compute with multi-step inference

Better quality than ILT (generated the training data) with 10X speedup



H. Yang et al, ILILT: Implicit Learning of Inverse Lithography Technologies



AIDIV I 😒

Al as part of a meta algorithm to solve challenges difficult for algorithms to solve

Layouts of thousands cells per library designed manually, difficult to automate

Layout/Routing challenges: satisfying all the DRC rules

Use an algorithm to route, but RL to fix DRC

Generate std cell layout of entire industrial libraries with better quality tl





RAGs are algorithms to ground LLM with facts to reduce hallucination

LLMs tend to hallucination

Incorrect, outdated, conflicting, scarcity of the training data

Retrieval Augmented Generation (RAG) helps to ground LLM with curated documents

RAG leverages the search algorithm (embedding similarity "Hash") to retrieve the relevant chunks as context for the question

Advanced RAGs use even more complicated algorithms



RAG Illustration



Agentic systems are algorithms to decompose complex tasks and use tools

Agent handles complex tasks by problem decomposition

Planning, memory and tools form algorithms

These algorithms use LLM to compute at each step





Planning/Reasoning methods are algorithms



Multi-agent communication patterns are algorithms

Multiple agents can talk to each other form an agentic system

Planning for multi-agent is the design of communication patterns

Communication patterns between agents are algorithms: static and dynamic patterns



Report Agent: compare and summarize timing reports with CoT and ReACT

Reading many timing reports are tedious work.

Tools

timing_metric_calculation_tool: Calculate changes in WNS, TNS, or FEP slack_distribution_calculation_tool: Calculate changes in the slack distribution



Chain-of-Thought

Let's think step by step.

Firstly, analyze the change of WNS, and TNS for all designs in all PVT corners tables of these two "datecode" settings and explain with numbers using tools.

Secondly, comparing "FEP" of two datecode settings in all PVT corners tables.

Thirdly, analyze the slack distribution to identify the distribution of "slack less than 0" paths.

Finally, summarize the analysis in the following aspects:

1. Provide key takeaways, comparison, and suggestions with bullet points of the two "datecode" settings based on previous steps.

2. Identify the corner which still suffers many timing violations if any.

You need to use the provided tools to analyze the timing metrics! You are not good at math!



RTLFixer: fix RTL syntax errors generated by LLM with RAG

55% of GPT-3.5 Verilog generation errors are Syntax errors

Agent can fix 99% of Syntax Errors

Get feedback/error messages from Verilog compiler

Retrieve human guidance for each syntax error with **RAG**

Compiler Logs:

Object 'clk' is not declared. Verify the object name is correct. If the name is correct, declare the object.

Human Expert Guidance:

Check if 'clk' is an input. If not, and if 'clk' is used within the module, make sure the name is correct. If it's meant to trigger an 'always' block, replace 'posedge clk' with '*'.



Tools



FVAgent : Generate SVA from natural language (NL) input with self-learned rules and multi-agents

SVA generation is even harder because of less data available

Articulated 'rules' in the context improves the accuracy

Self-Learning: Generate rule knowledge base from training data

Multi-Agent task flow : rule generation/retrieval/fix syntax

FVAgent improves Syntax/Function corrections on all models



Tools

💿 NVIDIA

24

J. Wan et al, FVAgent: Bridging Natural Language Specifications and Formal Verification Assertions with Adaptive Multi-Agent Learning

VerilogCoder Agent significantly outperforms baseline LLMs in Verilog benchmark: VerilogEval



C.-T. Ho et al, VerilogCoder: Autonomous Verilog Coding Agents with Graph-based Planning and Abstract Syntax Tree (AST)-based Waveform Tracing Tool

25 📀 nvidia

Task planning and Code Implementation are major stages of VerilogCoder

The task planning stage generates a task plan

The code implementation stage writes code for each planned task

Leverage special knowledge base and tools

Task Planning

Task-Driven Circuit Relation Graph (TCRG)

Code Implementation

AST-guided waveform debugging tool

Dynamic task plan and knowledge base

C.-T. Ho et al, VerilogCoder: Autonomous Verilog Coding Agents with Graph-based Planning and Abstract Syntax Tree (AST)-based Waveform Tracing Tool



Multi-agent planner builds Knowledge Graph to extract detailed information for task planning



Multi-agent coder leverages syntax checker and waveform tracing tools to write and debug code

Code Agent: Write partial Verilog code



Debug Agent: Check and Correct the functionality



Writing code for each task

Debug final code

28 📀 **NVIDIA**

```
Mate Terminal
                                                                                    X Mate Terminal
                                                                                                                                                                           X Mate Terminal
(/home/scratch.chiatungh nvresearch/nvcell env/llm env) bash-4.2$ python hardware agent/examples/verilog eval agent/Verilog eval task flow ver1.py
reading /home/scratch.chiatungh nvresearch/hardware-agent-marco/hardware agent/examples/verilog testcases/verilog-eval-v2/dataset dumpall/ Prob139 2013 q2bfsm prompt.txt
reading /home/scratch.chiatungh_nvresearch/hardware-agent-marco/hardware_agent/examples/verilog_testcases/verilog-eval-v2/dataset_dumpall/ Prob139_2013_q2bfsm_ref.sv
reading /home/scratch.chiatungh_nvresearch/hardware-agent-marco/hardware_agent/examples/verilog_testcases/verilog-eval-v2/dataset_dumpall/ Prob139_2013_q2bfsm_test.sv
('llm class': <class 'adlrchat.langchain.LLMGatewayChat'>, 'llm kwargs': {'model name': 'gpt-4', 'temperature': 0.0, 'top p': 1.0}}
case manager length = 1
current test is 2013 q2bfsm
current task id is 2013 q2bfsm
current test is 2013 q2bfsm
register tool sequential flipflop latch identify tool < function sequential flipflop latch identify tool at 0x7f65abdac040> caller: <autogen.agentchat.assistant agent.AssistantAgent object at 0x7f65a93ad1b0> executer: <autogen.agentc
hat.user_proxy_agent.UserProxyAgent object at 0x7f65a93ad390>
Hardware Agent Initialized 1 proxy, 0 rag proxy, 2 assistants
Hardware Agent Initialized 1 proxy, 0 rag proxy, 1 assistants
user (to chat manager):
You are a Verilog RTL designer that can break down complicated implementation into subtasks implementation plans.
[Example Begin]
### Problem
I would like you to implement a module named TopModule with the following
interface. All input and output ports are one bit unless otherwise
specified.

    input clk

    input reset

    input in (8 bits)

                                                                                                                                                                I

    output out (8 bits)

The module should implement an 8-bit registered incrementer. The 8-bit
input is first registered and then incremented by one on the next cycle.
Assume all sequential logic is triggered on the positive edge of the
clock. The reset input is active high synchronous and should reset the
output to zero.
### Solution
module TopModule
  input logic
                      clk,
  input logic
                      reset,
  input logic [7:0] in,
                                                                       output logic [7:0] out
 // Sequential logic
  logic [7:0] reg_out;
  always @( posedge clk ) begin
    if ( reset )
      reg out <= \theta;
    else
      reg out <= in;
  end
  // Combinational logic
  logic [7:0] temp_wire;
  always @(*) begin
    temp_wire = reg_out + 1;
```

AI learns to generate a problem-specific algorithm

Static algorithms (predefined CoT, taskflow) work but not flexible

Dynamic algorithms ('think step by step', VerilogCoder task plan) are more desirable

O1 improves this capability by learning – learning to reason

Future: Make LLM learn to reason in Chip Design problem

Write a bash script that takes a matrix represented as a string with format '[1,2], [3,4],[5,6]' and prints the transpose in the same format.

> Formulating the problem Crafting the script Reading input and parsing Breaking down the matrix Determining matrix dimensions Mapping out the matrix Constructing the output



"o1 learns to hone its chain of thought and refine the strategies it uses. It learns to recognize and correct its mistakes. It learns to break down tricky steps into simpler ones. It learns to try a different approach when the current one isn't working. "

OpenAl, Learning to Reason with LLMs



Conclusions

Al alone can not revolutionize how we design chip today

We should deeply integrate AI with algorithms to leverage the advantages and avoid the pitfalls of either one

Explore the vast potential of integrating AI with algorithms for EDA problems: design analysis and optimization.

Exploit the revolutionary advancement of the Agentic system (an algorithm using LLM as a computing node) for design assistance

Train AI to learn to generate its own algorithms for chip design problems.



