

REASONING COMPILER: LLM-Guided Optimizations for Efficient Model Serving

Annabelle Sujun Tang, Christopher Priebe, Rohan Mahapatra, Lianhui Qin, Hadi Esmaeilzadeh

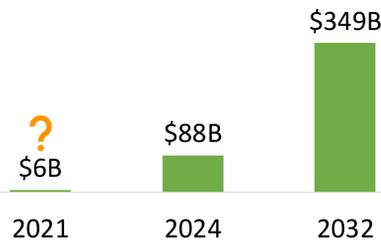
Alternative Computing Technologies (ACT) Lab, University of California San Diego



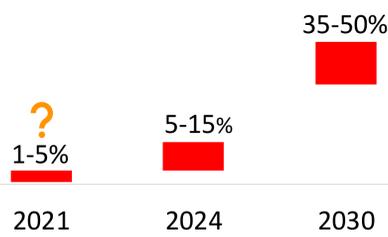
The Current Landscape

The Shift towards Model Serving

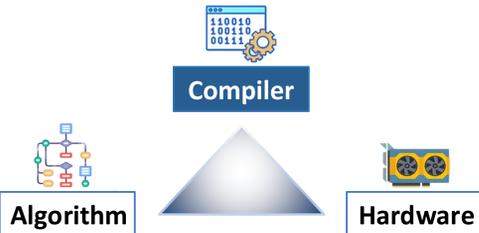
Global Inference Market



AI Data-Center Power Usage

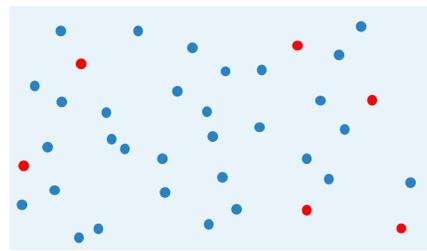


The Triangle of Improvement



Two Categories of Existing Compiler Optimizations

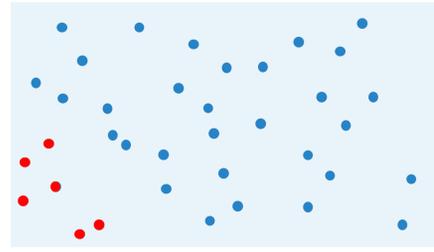
Category II: Stochastic Search for Compilation



- STOKE (Stochastic Super Optimization)
 - Markov Chain Monte Carlo (MCMC)
 - High quality programs often lie in regions separated by low-probability paths
- TVM, Anso, FlexTensor, Tensor Comprehensions
 - Genetic Algorithm
 - Simulated Annealing

- + Find Higher Quality Optimized Programs
- Sample Inefficient
- Cannot leverage context and interdependence

Category I: Rule-Based Compiler Optimization

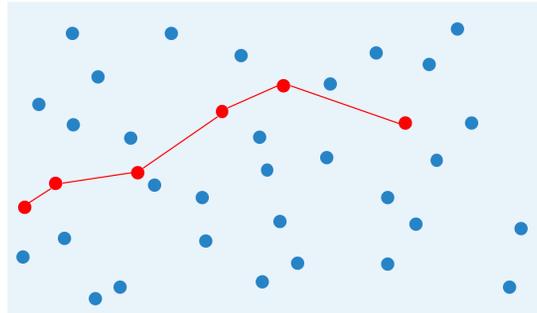


- Relies on hand-tuning or domain-specific heuristics
- Often overfit to a specific workload or hardware target

- + Relatively Fast
- Cannot explore the entirety of search space

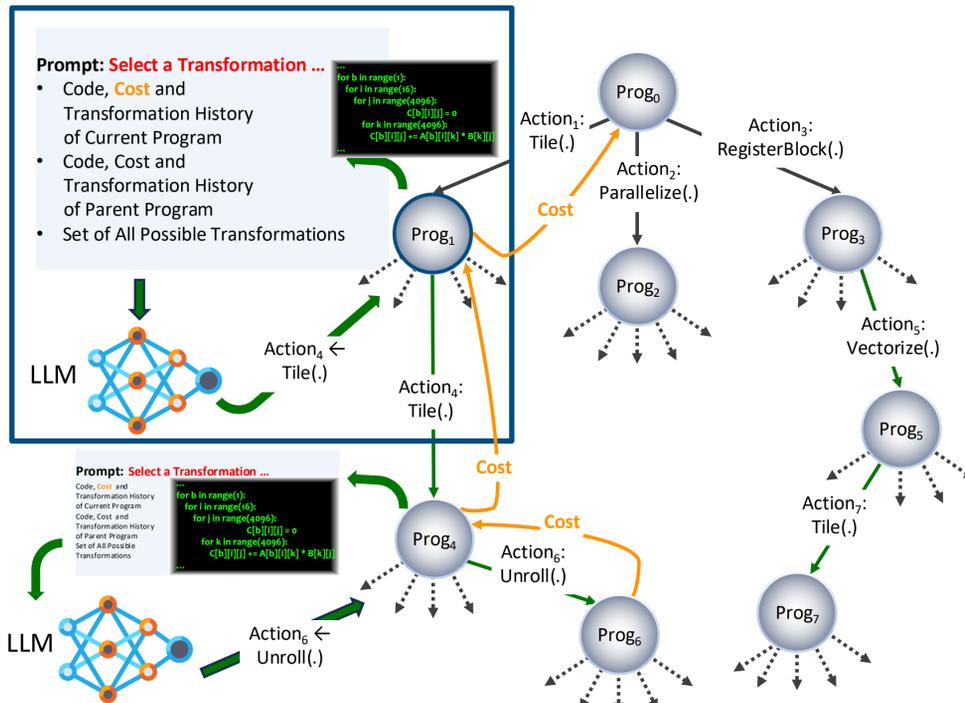
Our Innovation

REASONING COMPILER: Context-Informed, Guided, Structured Search



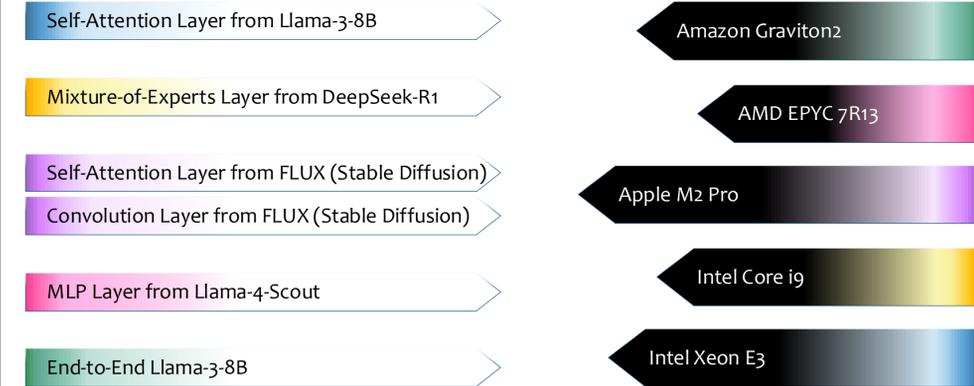
- Leaps from fully stochastic to structured optimizations
- Models optimization as Markov Decision Process
- Uses Monte Carlo Tree Search (MCTS) as a planner
- Utilizes LLM reasoning as a guide for Monte-Carlo Tree Search

- + Finds Higher Quality Optimized Programs
- + Sample Efficient

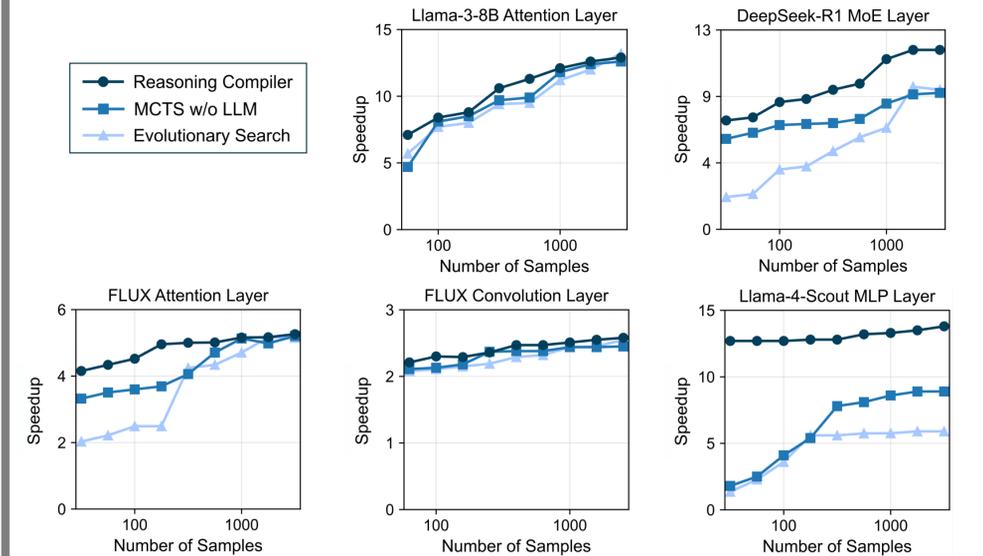


Empirical Evaluation

Diversity of Workloads and Hardware Platforms



Higher Speedup with Fewer Samples



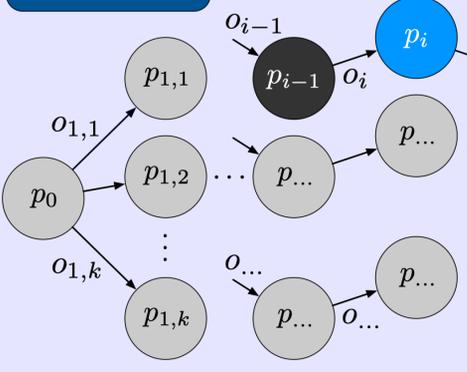
Across all 25 layer-platform pairs, **REASONING COMPILER** achieves a **5.0x speedup** with **5.8x fewer samples**: **10.8x** improvement in **sample efficiency** over Evolutionary Search.

Consistent End-To-End Improvements

Hardware Platforms	Evolutionary Search		Reasoning Compiler		Improvement	
	# Samples	Speedup	# Samples	Speedup	Sample Reduction	Sample Efficiency Gain
Amazon Graviton2	4,560	3.7x	1,440	5.1x	3.2x	4.4x
AMD EPYC 7R13	410	2.0x	140	2.2x	2.9x	3.2x
Apple M2 Pro	4,820	2.2x	1,770	3.9x	2.7x	4.8x
Intel Core i9	3,800	2.2x	720	4.9x	5.3x	11.8x
Intel Xeon E3	4,640	5.0x	670	5.0x	6.9x	6.9x
Geomean	-	2.8x	-	4.0x	3.9x	5.6x

For Llama-3-8B, **REASONING COMPILER** achieves a **4.0x speedup** using **3.9x fewer samples**, achieving **5.6x sample efficiency** over Evolutionary Search.

MCTS Tree



MCTS Tree Update



Chain-of-Thought Prompt

You are an AI compiler optimization assistant. We are performing a Monte Carlo tree search (MCTS) to find an optimal transformation sequence for a given node. In the MCTS tree, the "current node" is the leaf we are expanding, while "immediate parent" and "grandparent" refer to the ancestors in the tree. You are given the following information.

Programs

- p_{i-2}
- p_{i-1}
- p_i

Costs

- $\kappa(p_{i-2})$
- $\kappa(p_{i-1})$
- $\kappa(p_i)$

The set of all transformations O

Transformation Traces

- $S_{i-2} = \langle o_{1,\dots}, \dots, o_{i-2} \rangle$
- $S_{i-1} = \langle o_{1,\dots}, \dots, o_{i-1} \rangle$
- $S_i = \langle o_{1,\dots}, \dots, o_{i-1}, o_i \rangle$

Please compare the code, their transformation trace, and scores of these programs to see what changes might improve the current node's performance. Then, propose a *sequence* of optimizations (one or more) from the provided list. Output your chain-of-thought reasoning as well as the final list of transformations.

■ Grandparent ■ Parent ■ Selected Leaf Node