Towards the Next-Generation Al Agent Infrastructure

Yiying Zhang - GenseeAl and UCSD





About Yiying Zhang

- Associate professor at UCSD CSE, leading the WukLab (wuklab.io)
- Founder and CEO of GenseeAl (gensee.ai), an agent infrastructure startup
- 15 years of experience in systems and infrastructure
- ~5 years in SysML, recent years focusing on LLM serving and agents

- Interested in joining GenseeAl as summer intern or founding engineer?
 - email yiying@gensee.ai



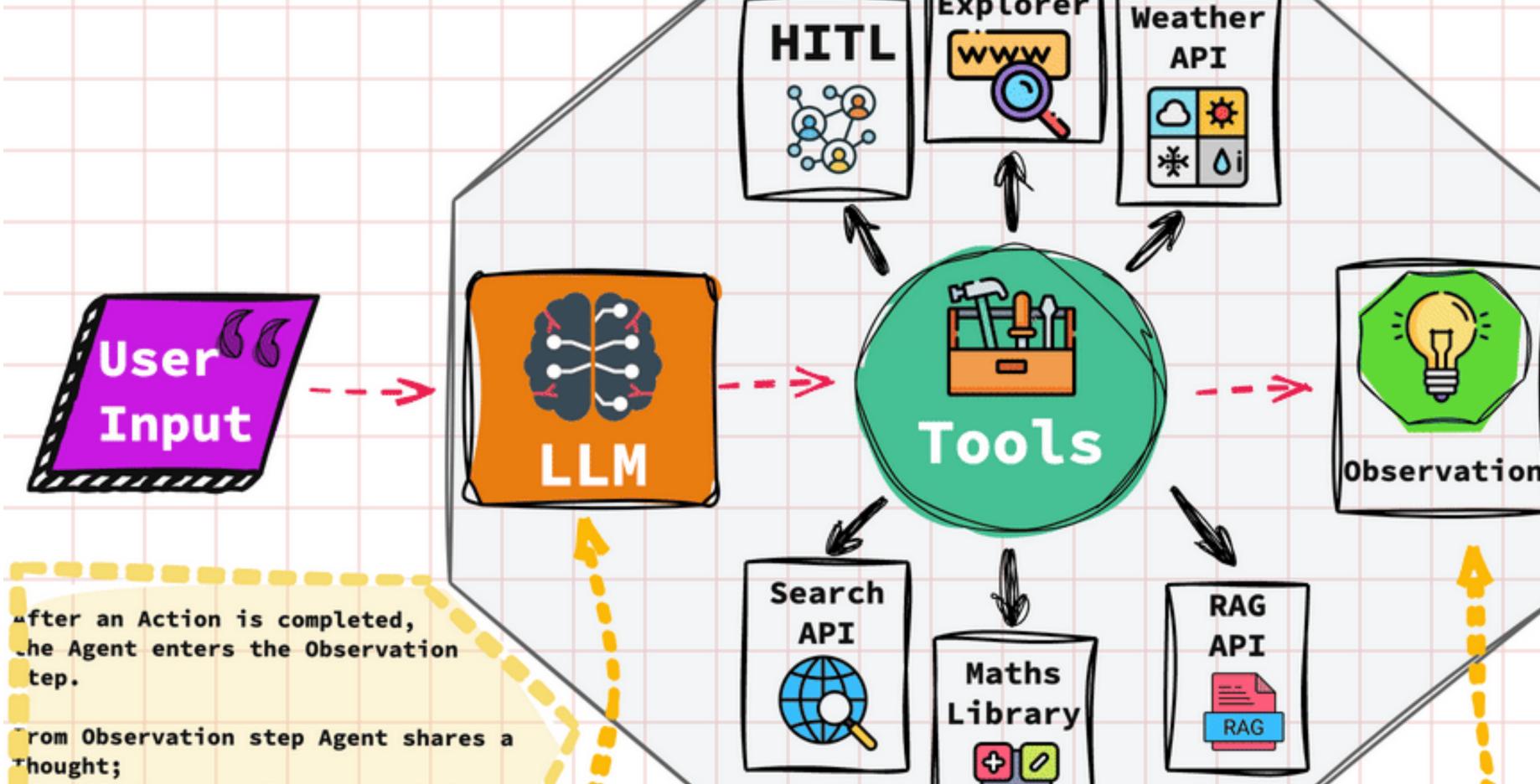
AI AGENTS - AGENTIC APPLICATIONS 1

Web

Explorer

<u>– [23]</u>







Final Answer Is Reached

COBUS GREYLING

& Al

of a final answer is not reached, he Agent cycles back to another ction in order

o move closer to a Final Answer.

The Status of Al Agent

Easy to demo, extremely hard to productionize

A large majority of organizations have deployed less than a third of their GenAl experiments into production

Deloitte.

Avg Time Nine Months — Prototype to Production

Time to Develop Al Prototype to Production Percentage of Respondents

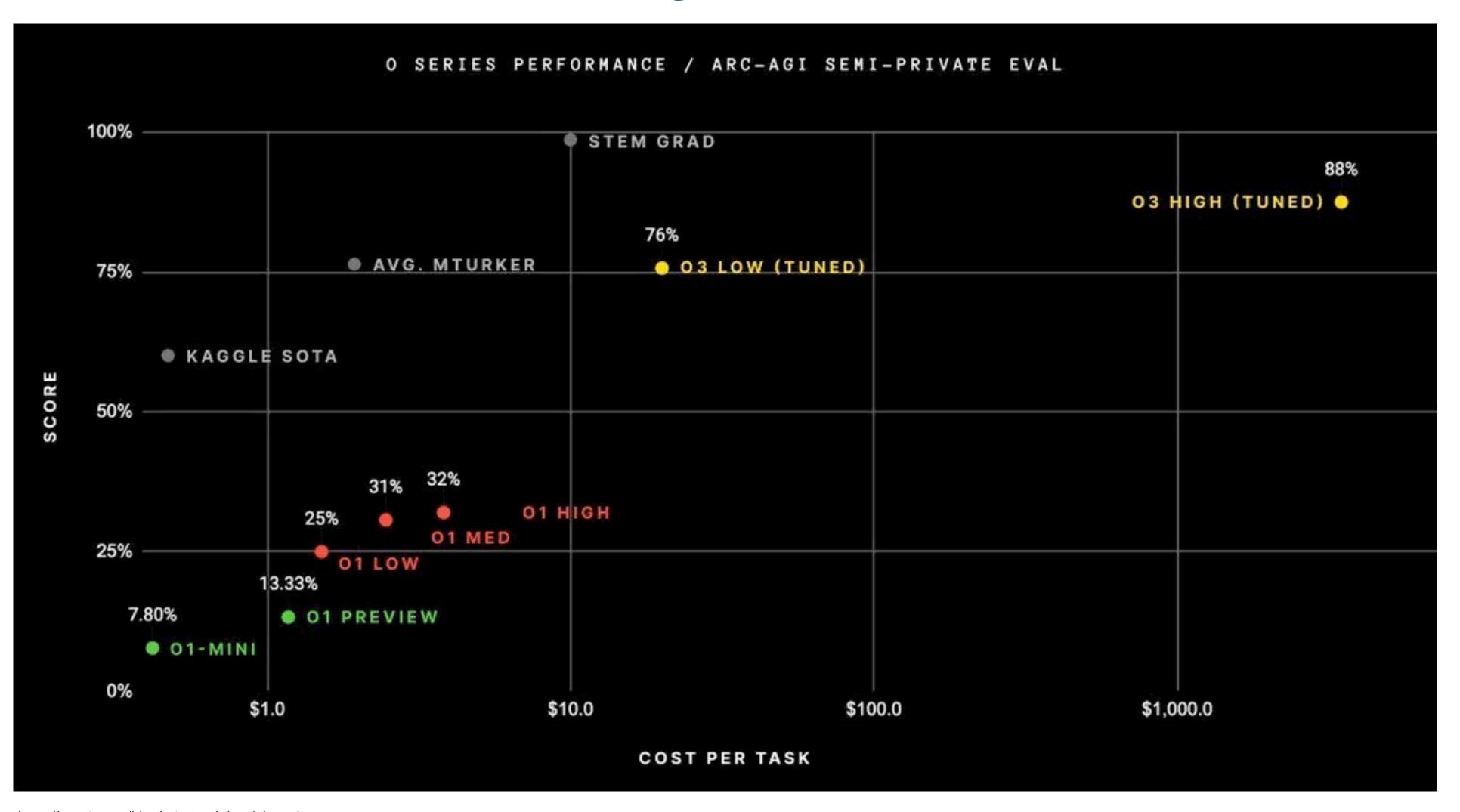
Gartner

5x more time to productionize [Al Refinery] than to demo

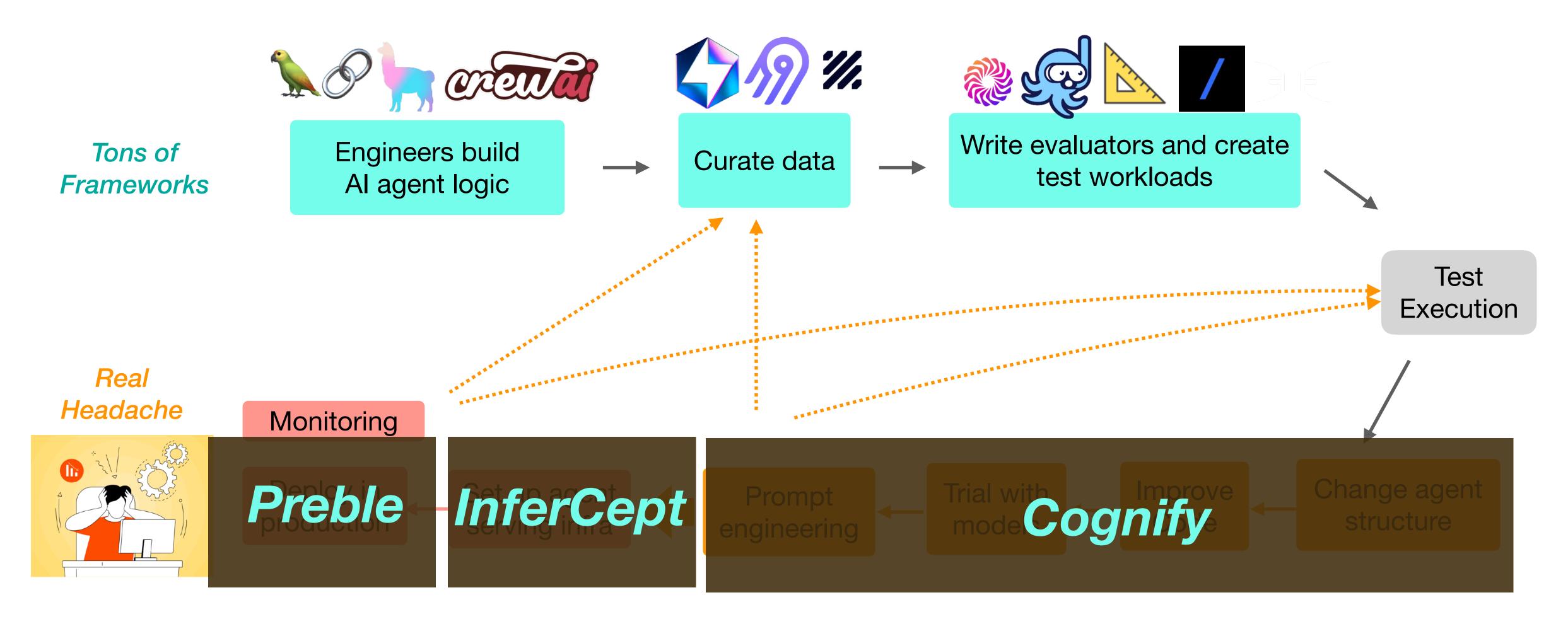


The Status of Al Agent

Rising inference cost, esp. as agent involves more model calls



Today's Al agent development and deployment life cycle



Today's Talk

Cognify [arxiv'25]: Al agent optimizer

Preble [ICLR'25]: long & shared prompt serving

InferCept [ICML'24]: compound LLM serving

Cognify: Multi-Facet Al Agent Optimization

Under submission, open source at https://github.com/GenseeAl/cognify

- Al Engineers spend a lot of effort tuning Al agents and workflows
- Cognify: aims to autotune Al agents and workflows with a small budget
- Key challenge: the search space is huge and proper searching needs \$\$
 - A simple 4-step workflow could need \$168K and weeks to search brute force
- AdaSeek: adaptive hierarchical BO-based search
 - Uses \$5 and 24 minutes to autotune the above workflow with 2.8x higher quality









Challenge: Large Search Space but Limited Budget

Assuming 3 LLM steps in an agent workflow

12 configurations in total, each with 3 options

Grid search requires: $(3)^12 = 531,441$ runs!!!

Assume workflow

- use gpt-4o: \$10/1M tokens
- each execution output 10K tokens

⇒ \$53K

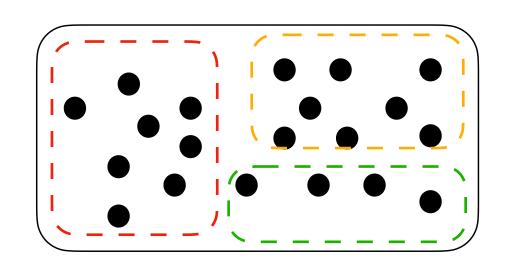
What's the best I can get with 128 trials?

Our Approach: Adaptive Hierarchical Search

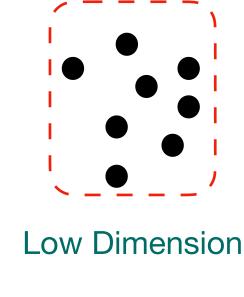
- 1. Hierarchical search with layered cog search space
- 2. Result-driven budget distribution
 - A. Search space partition
 - B. Search budget partition
- 3. Dynamic resource re-distribution

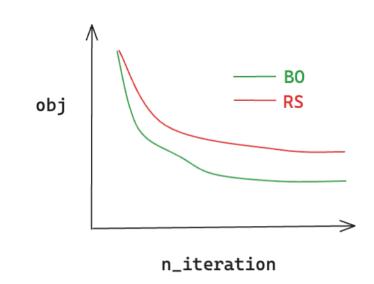
Cognify's Solution: Hierarchical Search

Flattened Search Space

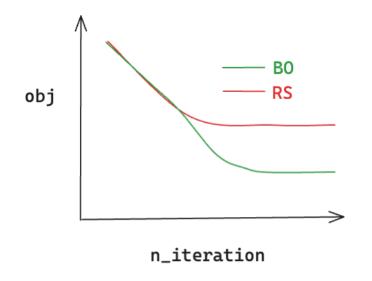


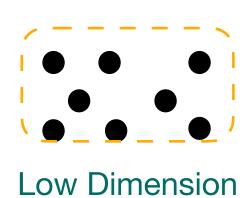
High Dimension

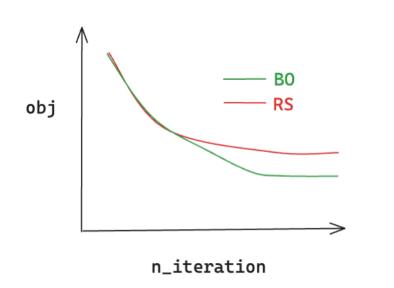




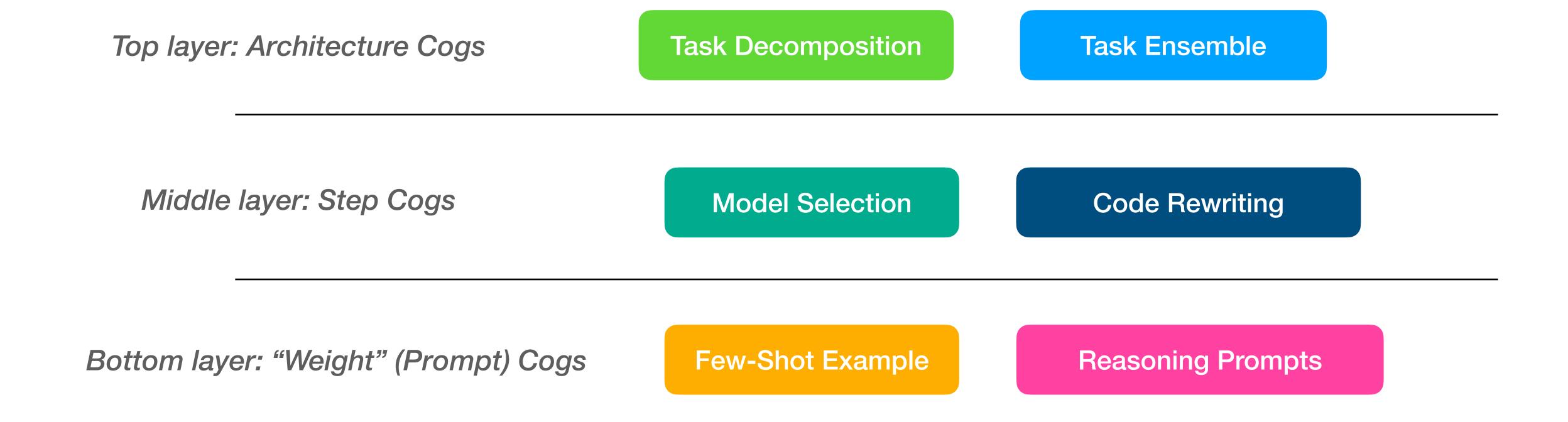








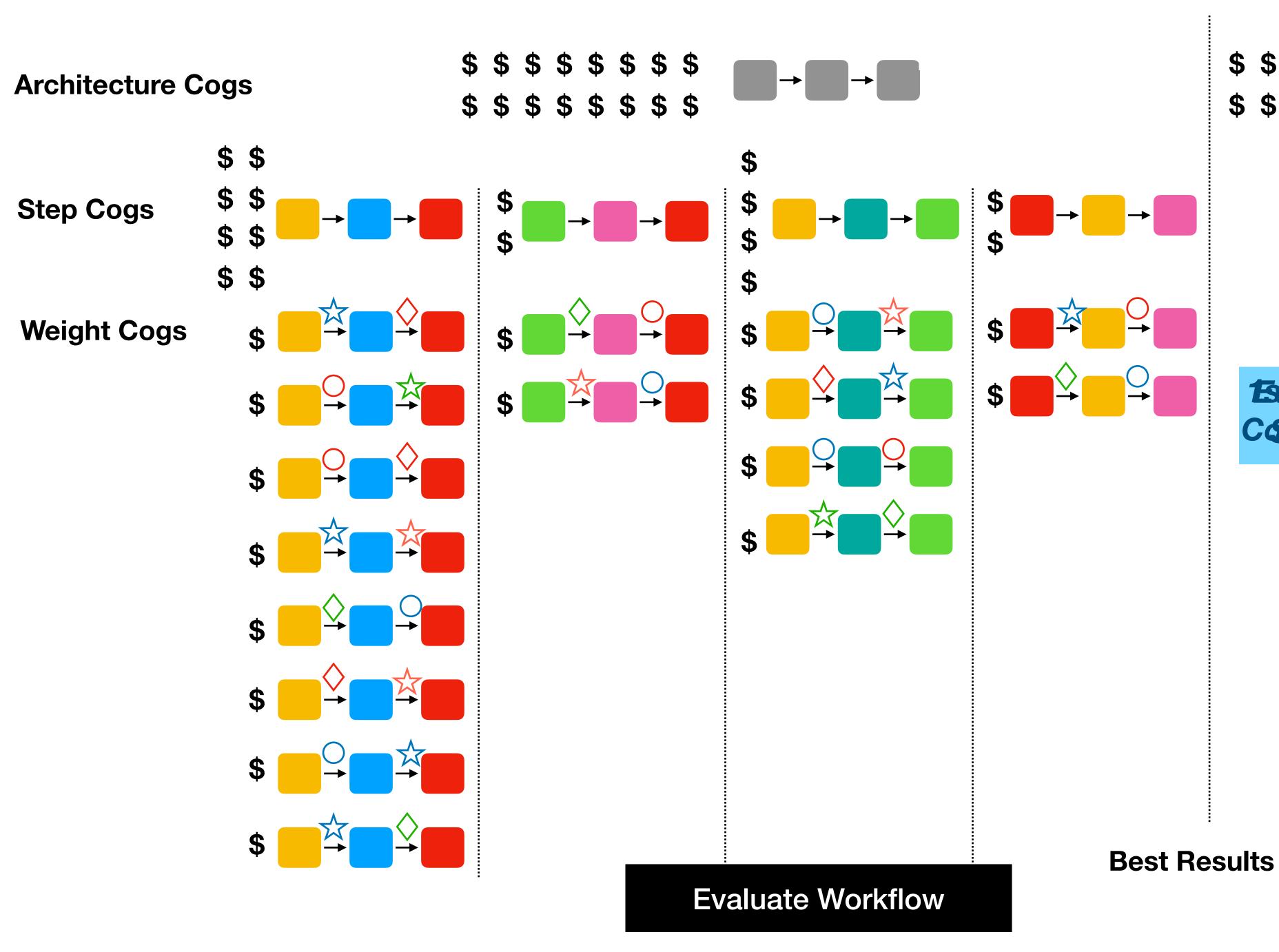
Cognify: Organize tuning knobs (cogs) hierarchically

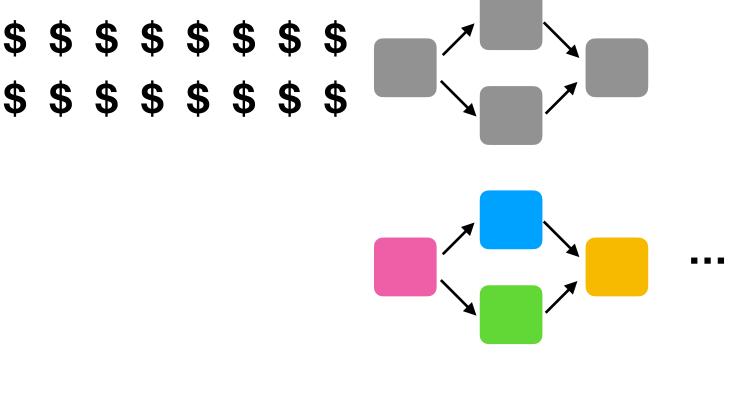


Cognify's Solution: Adaptive Search

• Partition search budget across hierarchies according to layer complexity

Direct search budget to more promising configurations using SH

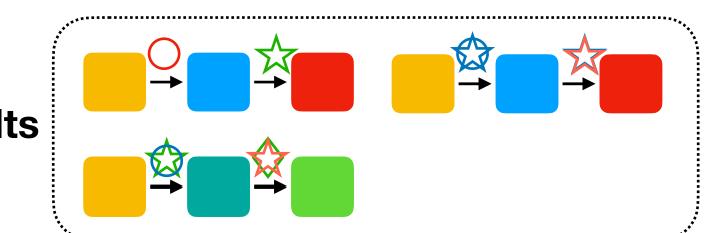




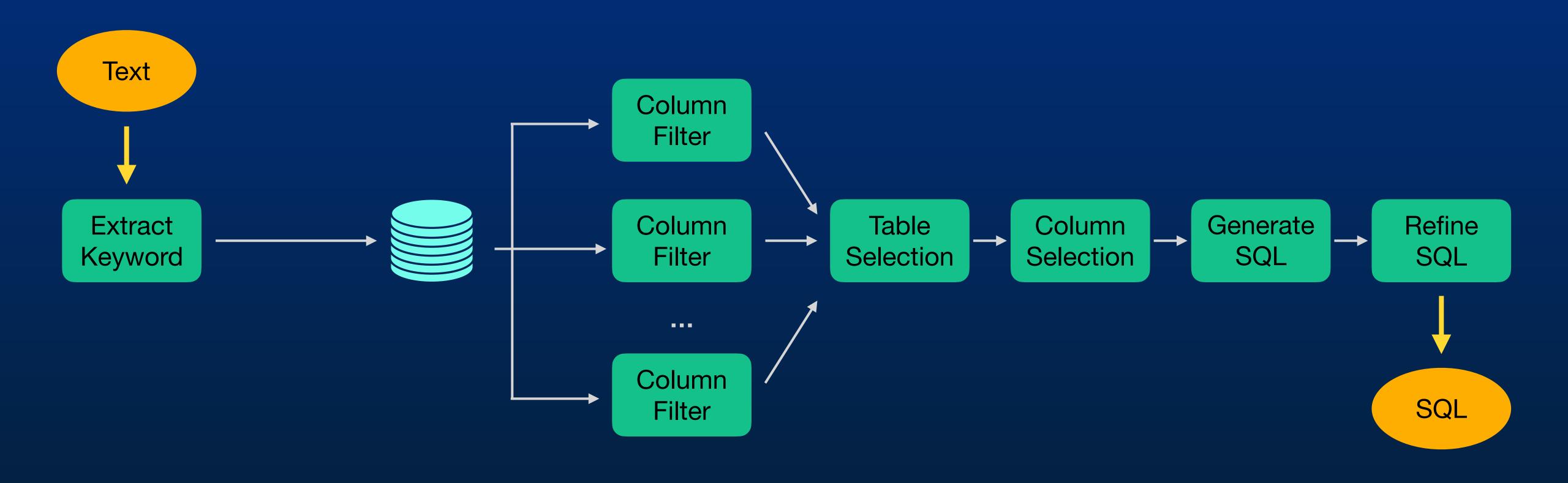


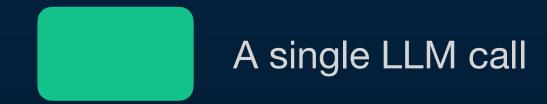
Search Budget (32 iterations)





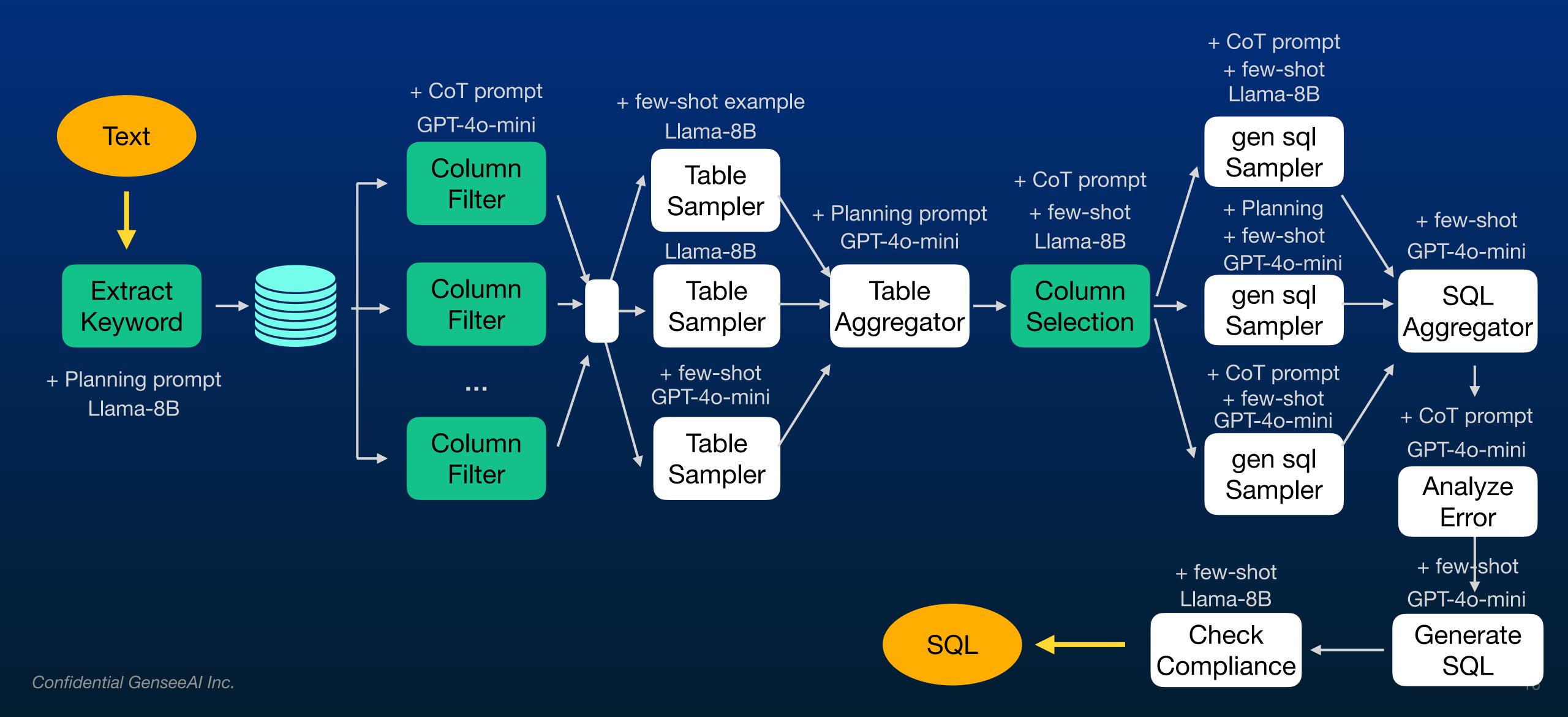
A real gen-Al workflow: text-to-SQL



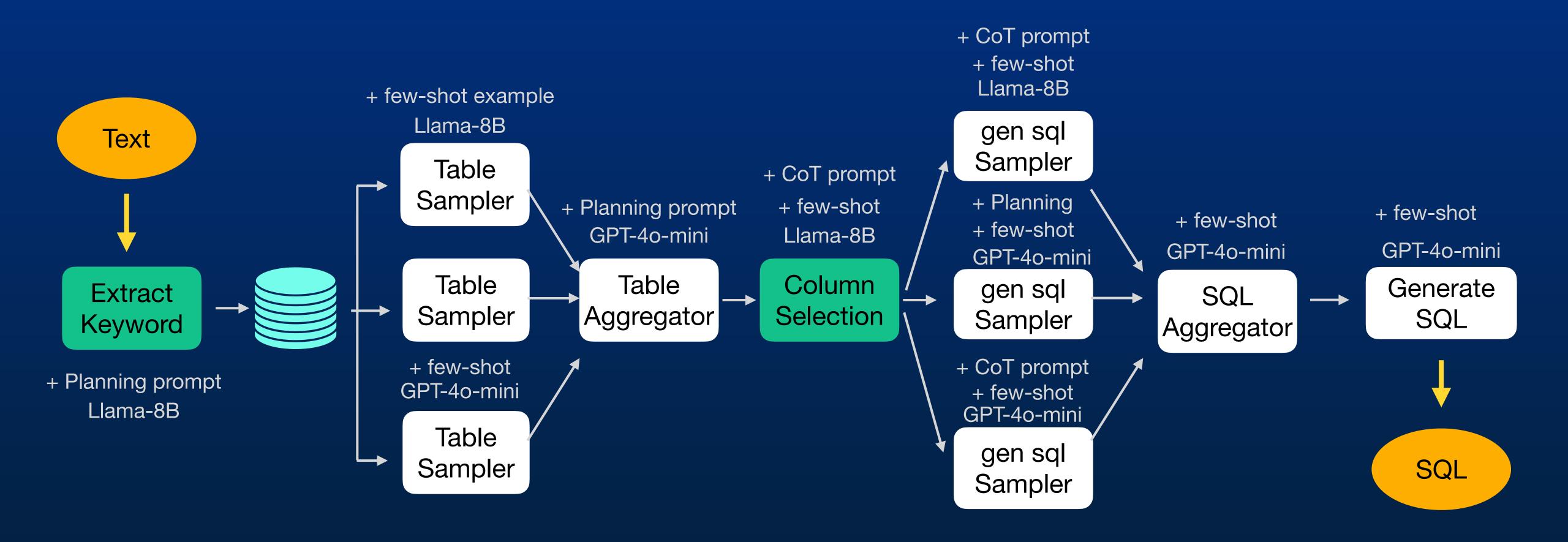


Confidential GenseeAl Inc.

Gensee's optimized text-to-SQL (training+AutoML)

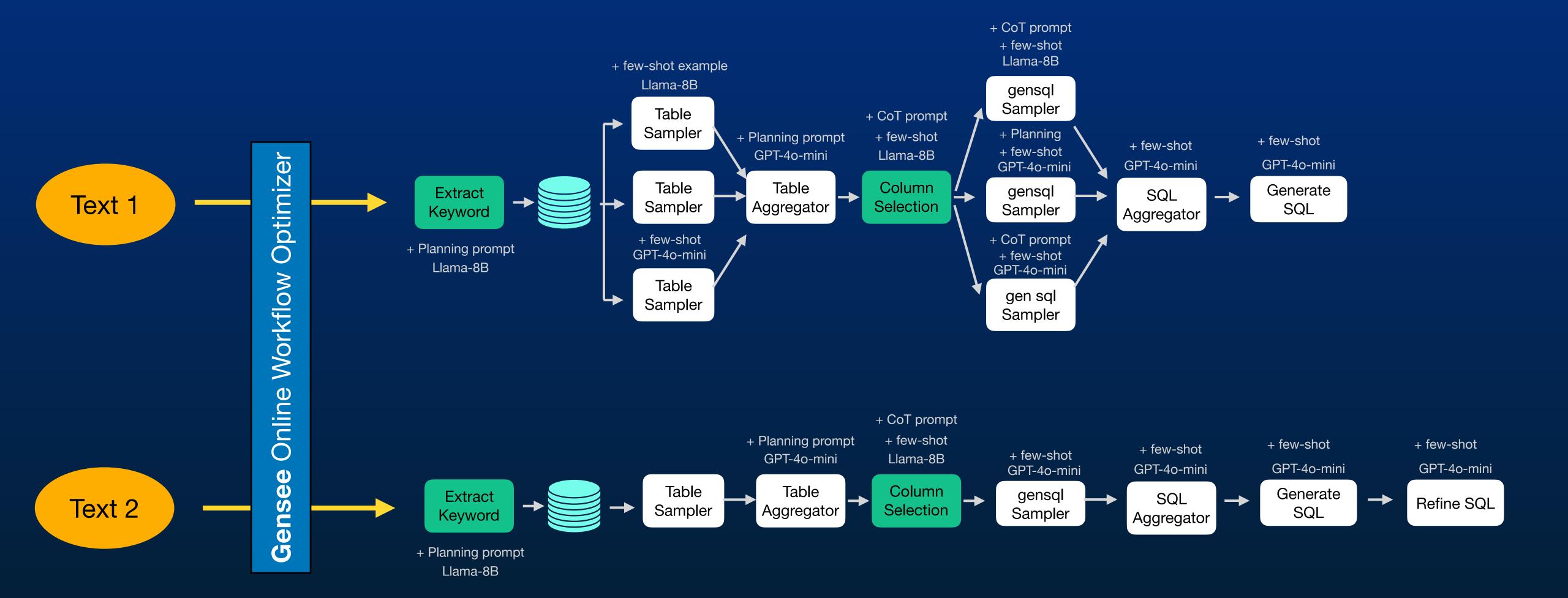


Another Gensee's optimized text-to-SQL (latency-oriented)



Confidential GenseeAl Inc.

Gensee's online per-request optimization (serving)



Confidential GenseeAl Inc.

Cognify Results

Six workloads: RAG-based QA, text-to-SQL, data visualization, financial analysis, code generation, BigBench

Up to 2.8x quality improvement

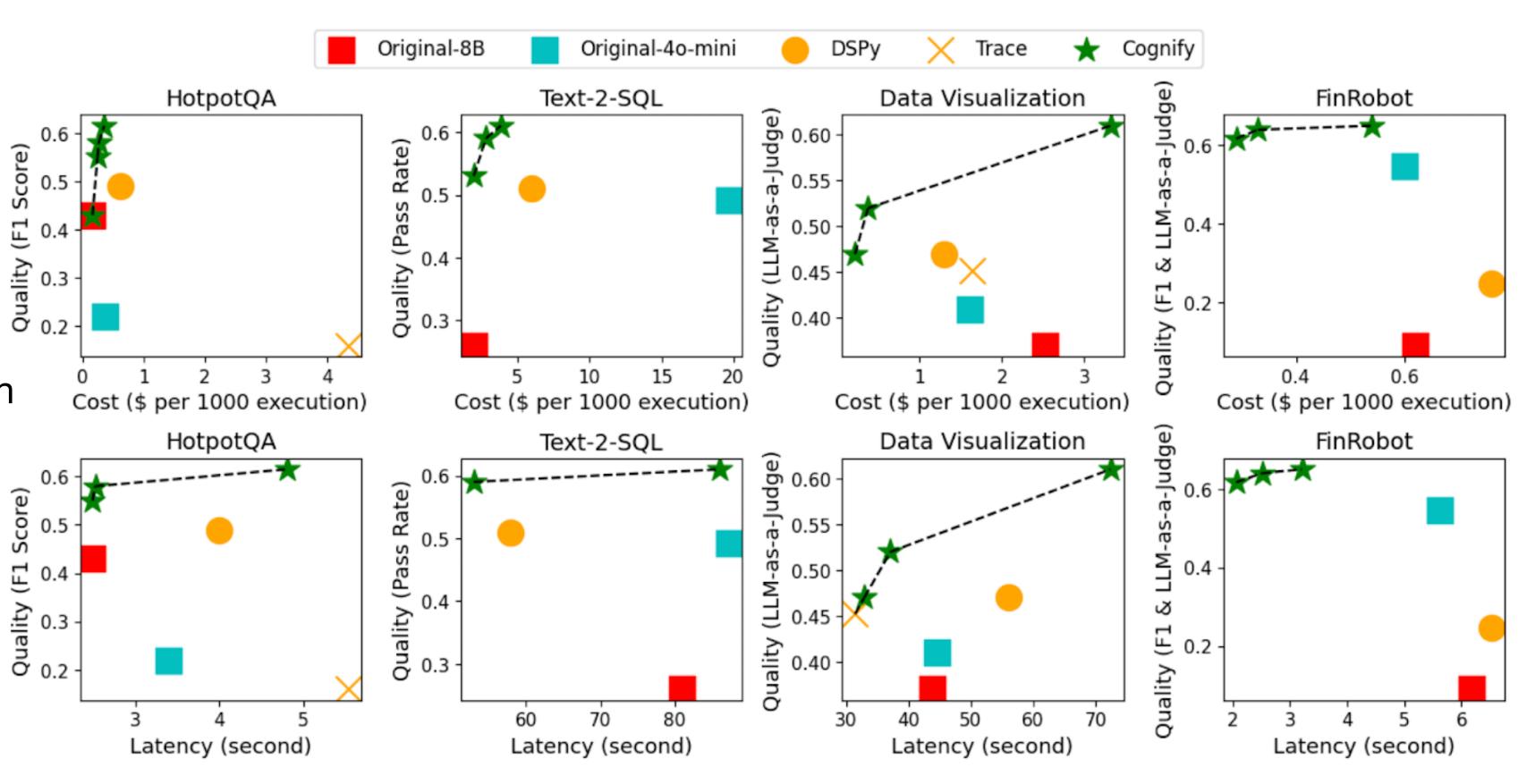
Up to 2.7x latency reduction

Up to 10x cost saving

Multiple choices on Pareto frontier

Low optimization costs: \$1.5-\$10, 8-42min

(roughly half of DSPy)



Cognify Takeaways

- Production-grade agents require manual input to incorporate business logic
- But manual efforts for tuning production agents can be avoided
- Cognify is the first production-ready autotuning tool for Al agents
- Please join our discord and star our repo to learn more!
- https://github.com/GenseeAl/cognify https://discord.gg/8TSFeZA3V6

Today's Talk

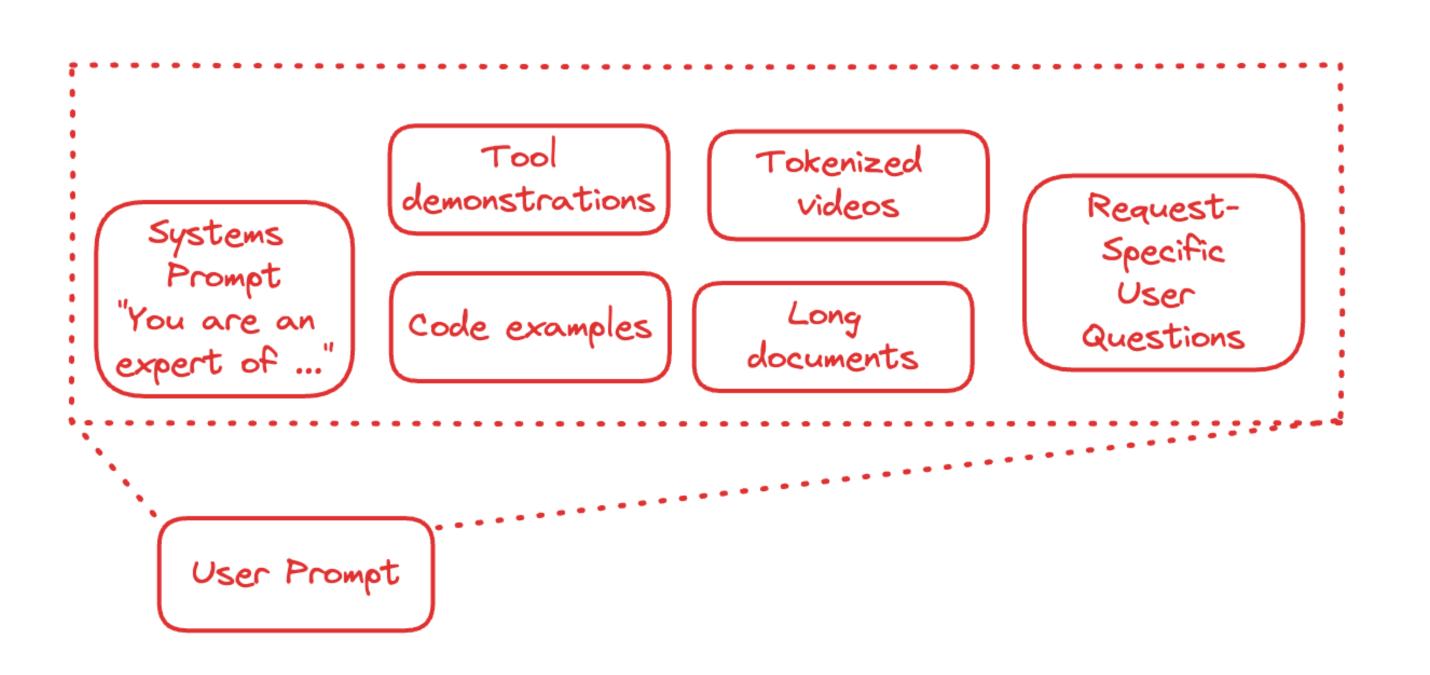
Cognify [arxiv'25]: Al agent optimizer

Preble [ICLR'25]: long & shared prompt serving

InferCept [ICML'24]: compound LLM serving

It's all about prompting

- Agent prompts are more than just a simple question





https://towardsdatascience.com/how-i-won-singapores-gpt-4-prompt-engineering-competition-34c195a93d41

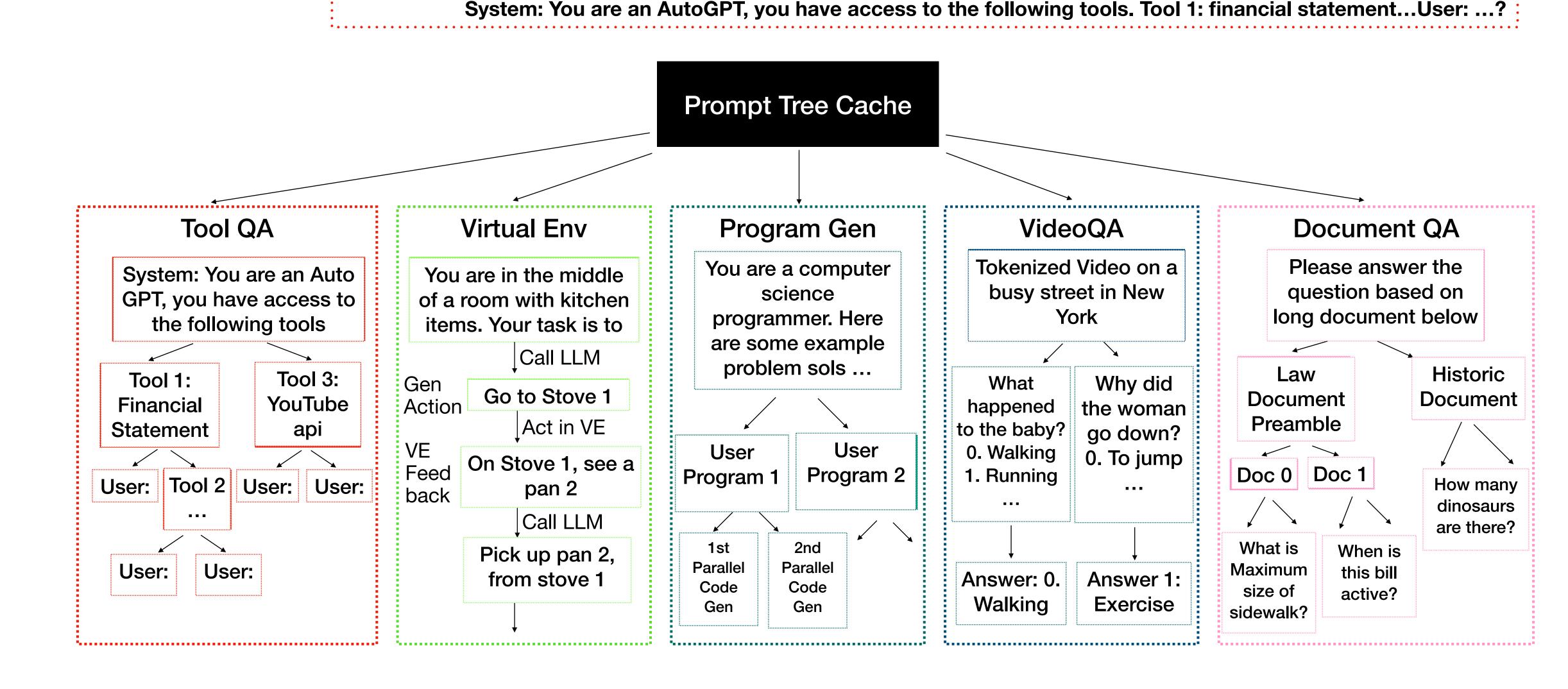
Prompts

Request Please answer the question based on the long document below ...

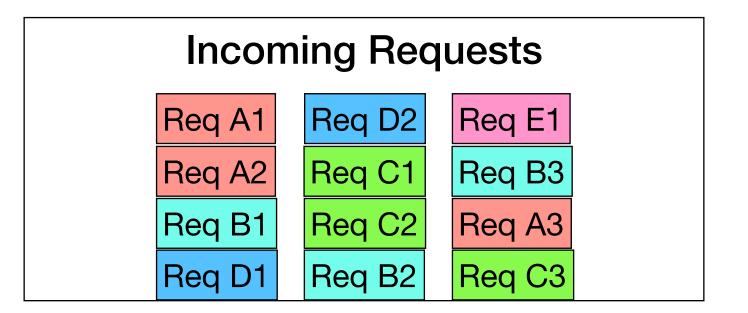
Queue Tokenized Video ... What happened to the baby? 0. Walking 1. Running
You are a computer science programmer. Here are some example problem sols ...
You are in the middle of a room with kitchen items. Your task is to

System: You are an AutoGPT, you have access to the following tools. Tool 3: ... Parameters...

System: You are an AutoGPT, you have access to the following tools. Tool 1: financial. Tool 2: ...



Prompt Agnostic Serving



Round Robin Scheduler

GPU 0 Llama Model

Processing Queue

recompute

GPU 1 Llama Model

Processing Queue

GPU 2 Llama Model

Processing Queue

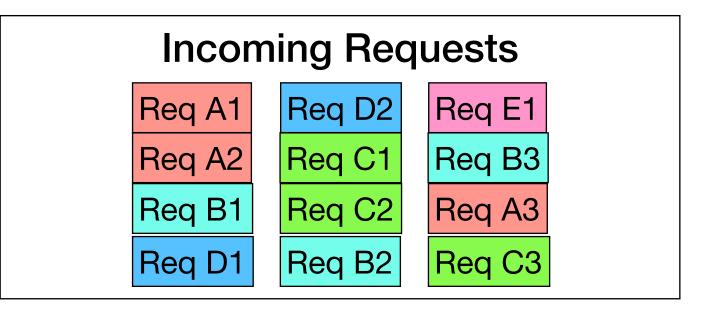
GPU 3 Llama Model

Processing Queue

recompute recompute recompute

All Requests are recomputed

Prompt Aware Scheduling



Prompt Aware
Distributed Scheduler

GPU 0 Llama Model

Processing Queue

GPU 1 Llama Model

Processing Queue

GPU 2 Llama Model

Processing Queue

GPU 3 Llama Model

Processing Queue

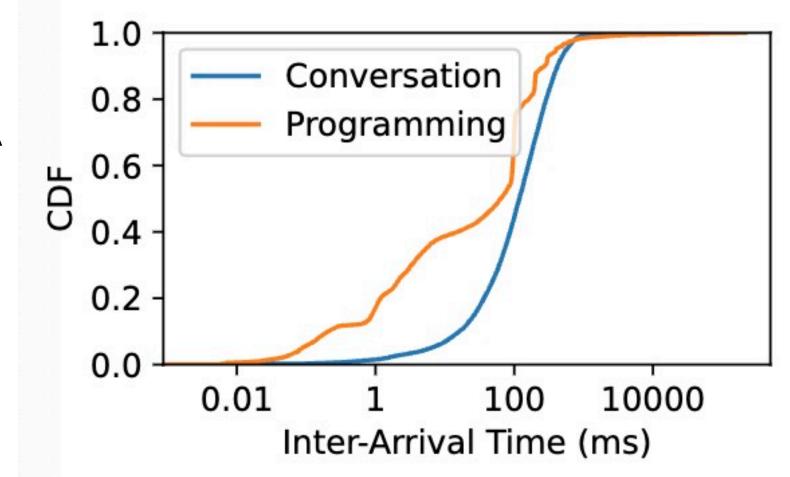
rearenneente rearenneente rearenneente rearenneente rearenneente rearenneente

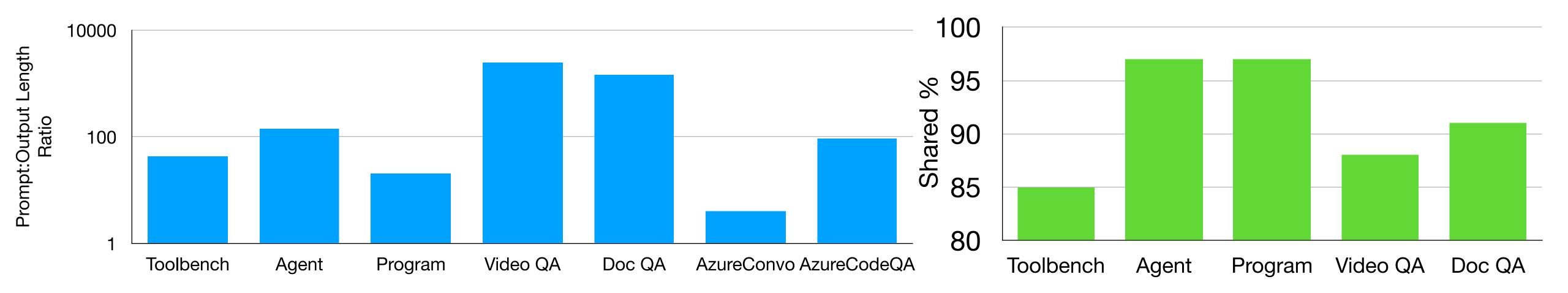
Faster completion due to prefix reusing

How are real-world prompts like?

- A study of prompts from systems perspective

- Studied 5 workloads and 1 real LLM request trace
 - tool use, embodied agents, program generation, video QA, and long document QA
 - 2023 Azure LLM Inference Trace
- Long prompts followed by short output
- High sharing degree
- Variation in request load





Preble: Distributed LLM Serving for Long and Shared Prompt ICLR'25, open source at https://github.com/WukLab/preble

- A distributed serving system targeting long and shared prompts
- Co-designs prefix sharing and load balancing
- Centers around a new E2 distributed scheduling algorithm
- Two-level scheduler for scalability
- Optimizes avg and p99 latency (up to 14.5X and 10X improvement over SoTA)









Preble's E2 Scheduling:

Exploration + Exploitation

prefix3 uniqque

E2 Scheduler

shared parterinc exculorage

>>>

explored lie presty bad GPU4

load = 1
recompute

load = 2
recompute

load = 4
reuse

load = 1
reuse

GPU 1

prefix1

GPU 2

prfx2

req

GPU 3

prefix3

req req

p4
prefix3

Exploit and Explore HeuristicGreedy Exploit

- Length of the prefix is proportional to compute
- Shared prefix > rest of prefix then exploit

Intuition on the computation:

- New Server: Shared + unique portion of prefix cost
- Existing Server: unique portion of prefix cost + eviction cost(unique # tkns)

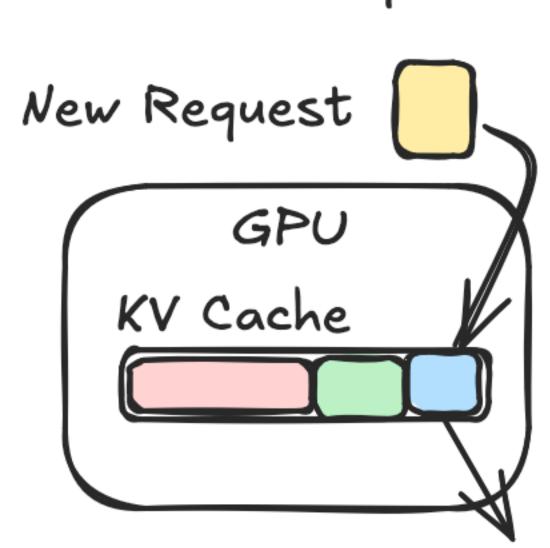
The amount of recomputation saved is larger than new computation

How do you Explore?

Prefix aware exploration

- 1. Total Request Load in Window(H)
- Running Batch

2. Load to be evicted to run request



New Request

3. Cost to run request

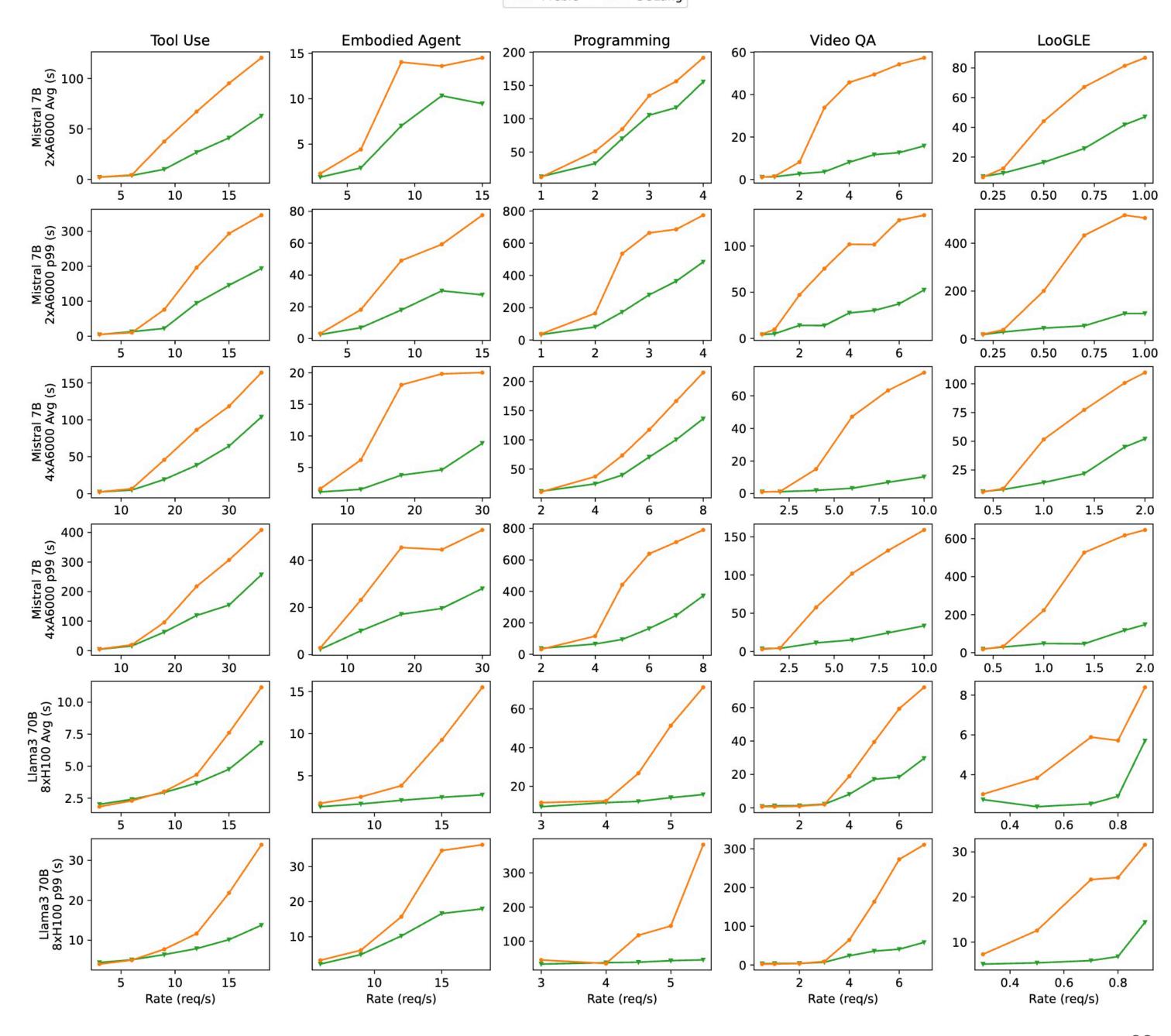
Cost = 1. + 2. + 3.

Calculate a load cost for each GPU and pick the min one

Completed Request **Preble Architecture Tokenizer** Queue Queue Req 0 Req 0 Global Scheduler Req 1 Req 1 Req 2 **Prefix Tree** Req 2 (Global and local) Req 3 Req 3 Rebalancer **Metrics Analysis** Engine Finished/ **Eviction Req** E2 Scheduler **GPU N** GPU 0 Tree Matching Tree Matching **Fair Waiting** Fair Waiting Running Running Queue Queue

Initial Results

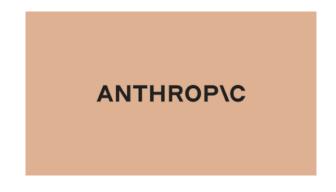
- Mistral-7B on 2 A6000
- Mistral-7B on 4 A6000
- Llama3-70B on 8 H100
- Five workloads
- cmp SGLang, Oracle partition
- avg & p99 req latency, RPS
- 1.5X to 14.5X on avg
- 2X to 10X on p99



Impact in Real World

Paper released May 2024

Cheaper Prefix Sharing in Real World Production systems



August 2024



October 2024



October 2024



MoonShot Al

 Directly inspired by Preble



Heuristic based

Preble Takeaways

- LLM Serving is getting more expensive using more complex prompting
- Workloads are longer and shared
- Preble(ICLR '25) enables cache and load to be effectively utilized for performance
 - Utilizing E2 scheduler and fair waiting queue

Today's Talk

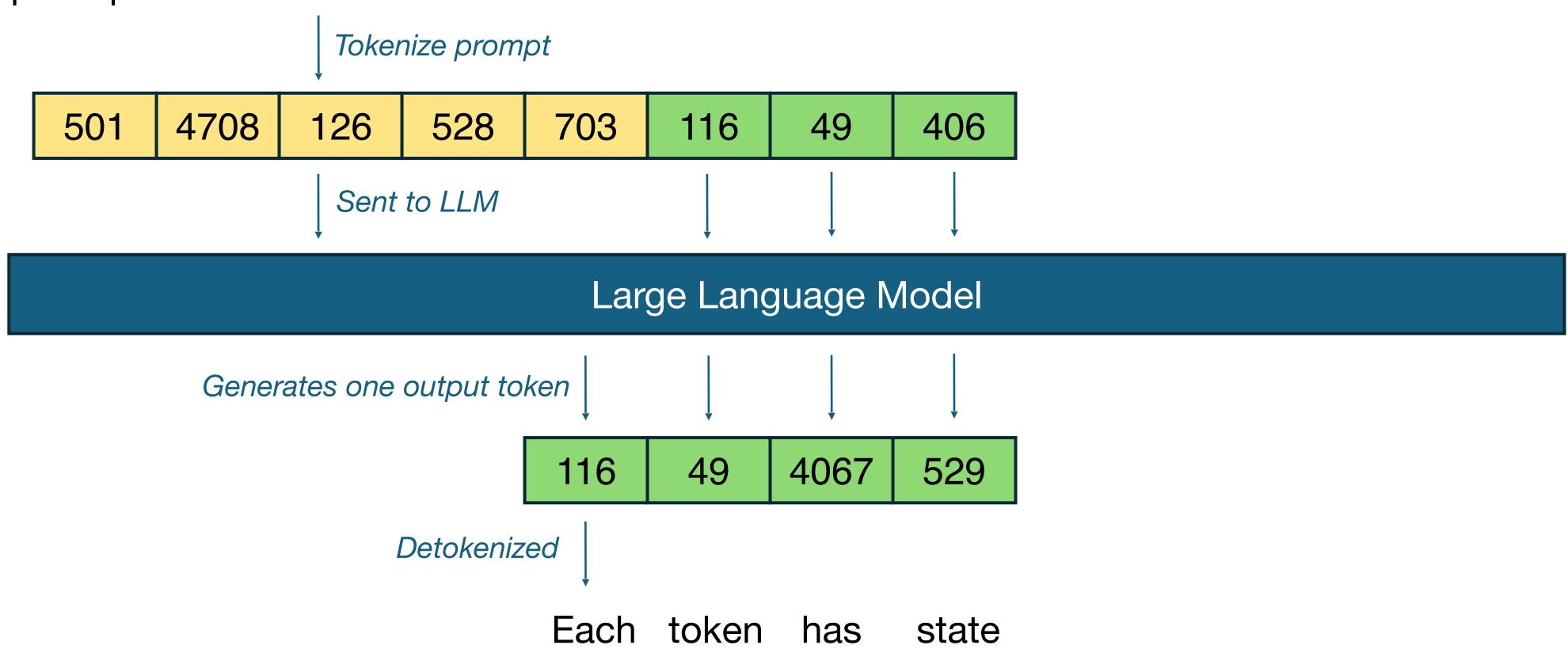
Cognify [arxiv'25]: Al agent optimizer

Preble [ICLR'25]: long & shared prompt serving

InferCept [ICML'24]: compound LLM serving

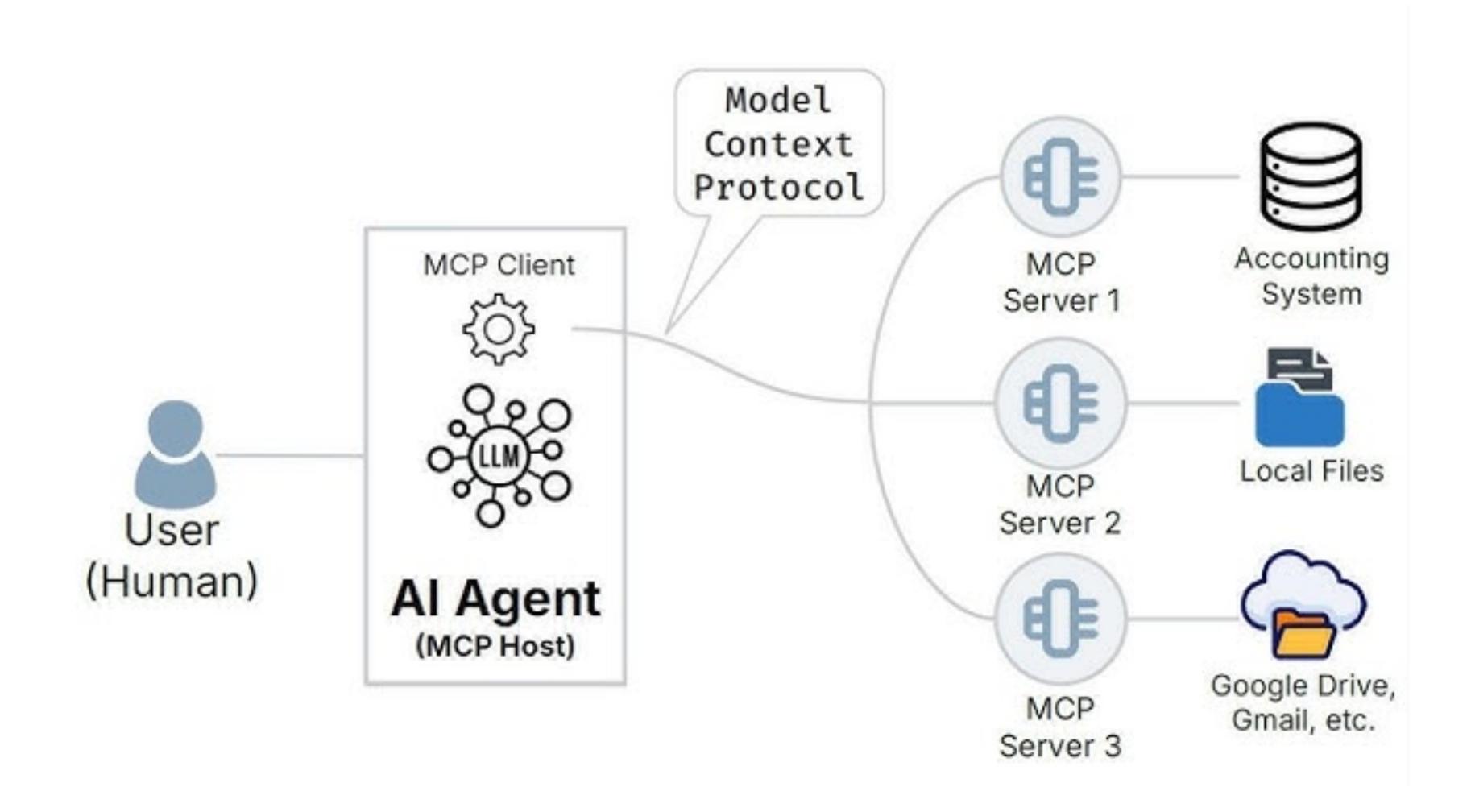
Typical LLM Inference: A closer look

Prompt: Explain the attention KV cache.

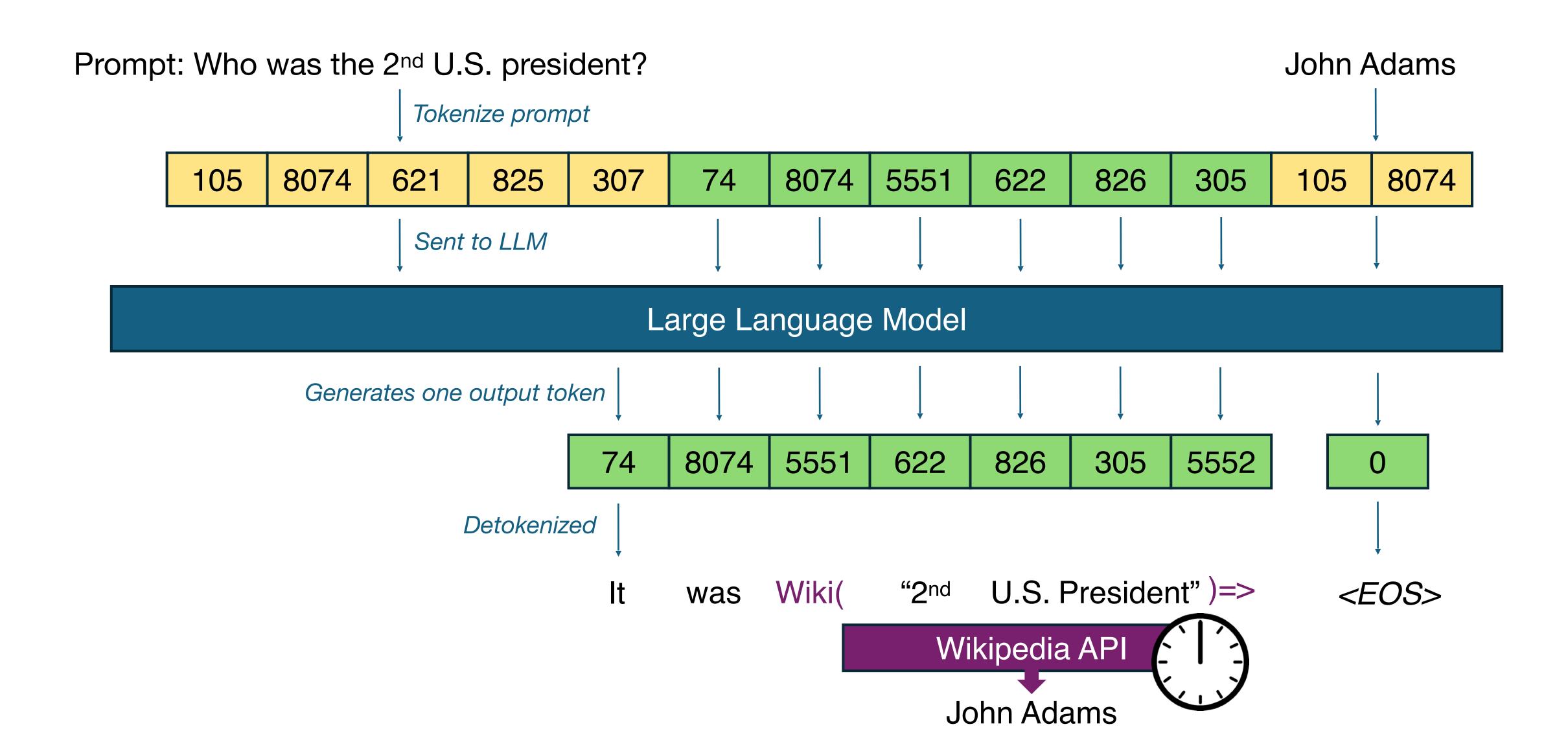


Agent Inference with Tool/Data Calling

Model Context Protocol (MCP)

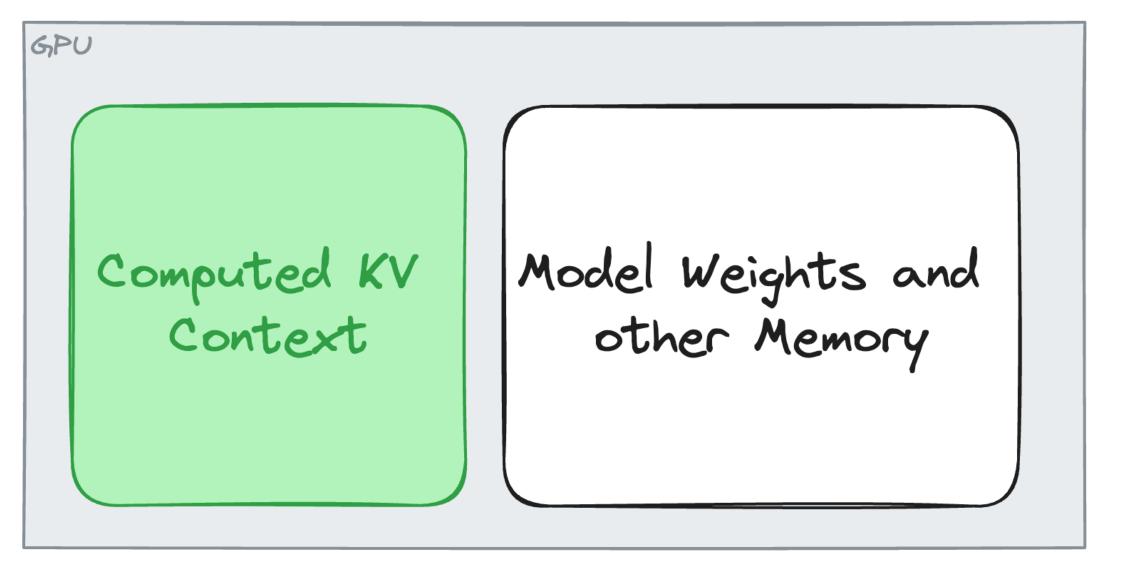


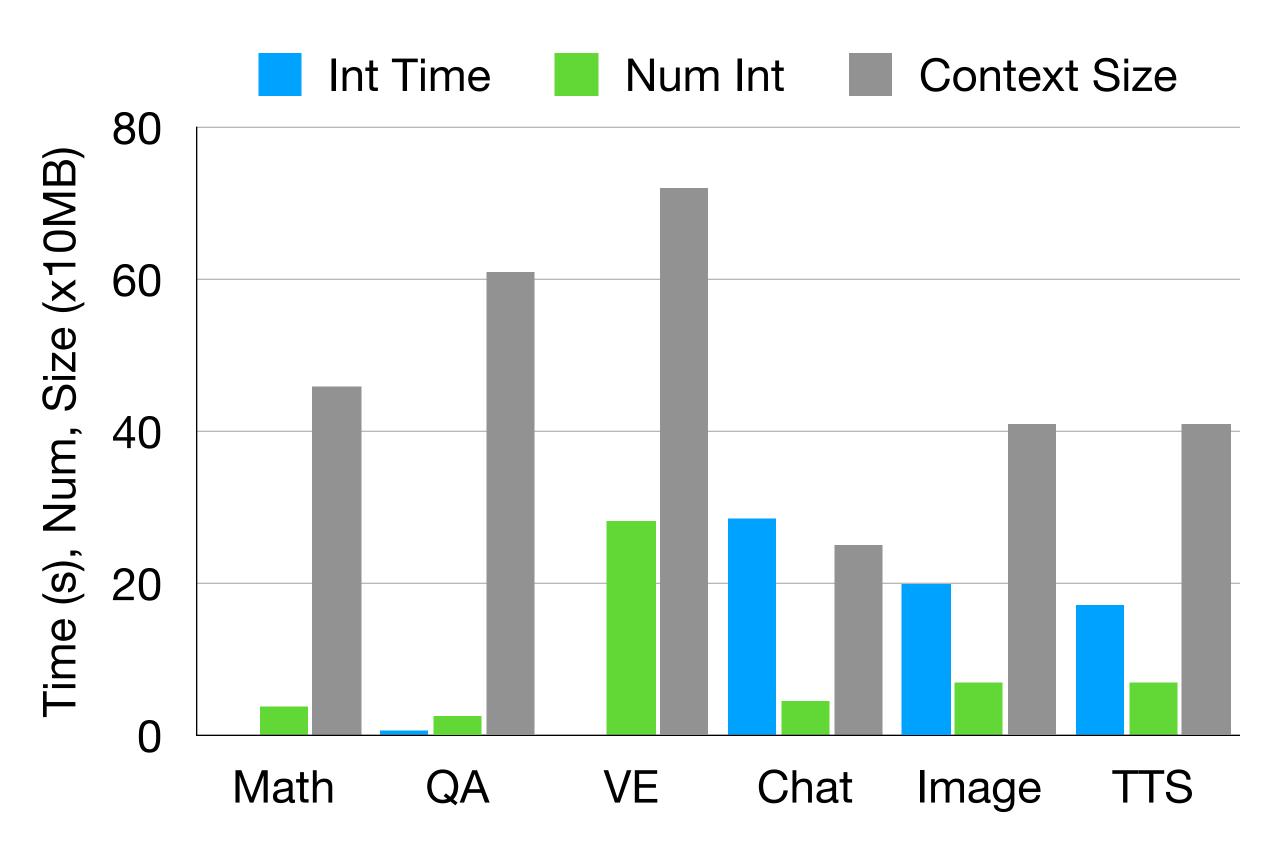
Agent Inference by Augmenting LLM Calls



How much is LLM intercepted?

- A study of six sets of real-life compound LLM workloads
- Intercepting for microseconds to minutes
- 2.5 to 28 interceptions per request
- Context 248MB to 720MB per req*





^{*} Assuming a 70B model

How are LLM interceptions handled now?

- They are not!
- SoTA LLM serving systems treat LLM interceptions as end of requests
 - Discard all KV context
 - (Re)compute KVs for tokens in context when interception ends
 - 37% 40% e2e request latency spent on recomputation
 - Wastes 40% GPU resources

InferCept: Adaptive LLM-Centric Workflow Inference ICML'24, open source at https://github.com/WukLab/infercept

- Pause a request upon intercepting
- Adaptively choose strategies for dealing with KV context
- Efficient implementation of intercept strategies
- Multiple intercepting endpoints supported (tool, other model, human, ...)
- 1.6x to 10x improvement over vLLM (SoTA LLM serving system)









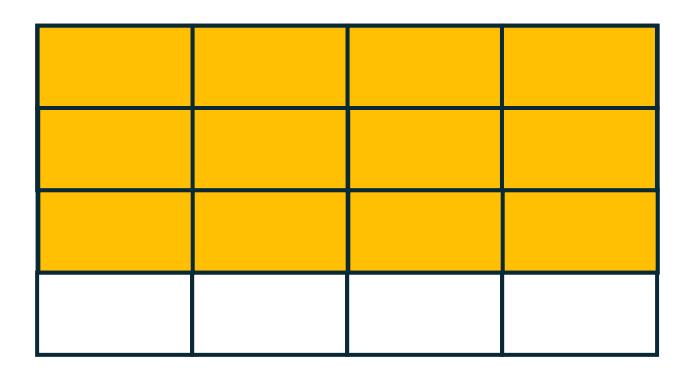


Three Intercepting Strategies

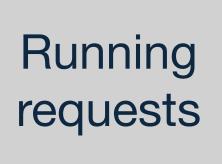
- when dealing with KV context
- Discard KV context and recompute upon return
- Preserve KV in GPU memory during interception
- Swap KV to CPU memory during interception

Strategy 1: Discard + Recompute

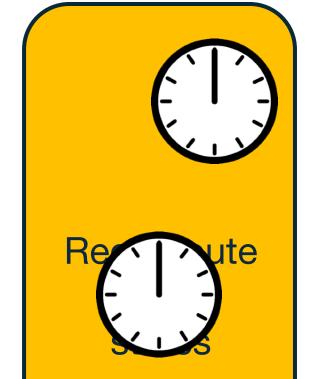
GPU memory



Bieccamobute all indetembrates the indetembrates the indetembrate in the indeterminant in the inde



Typical iteration time: 40ms

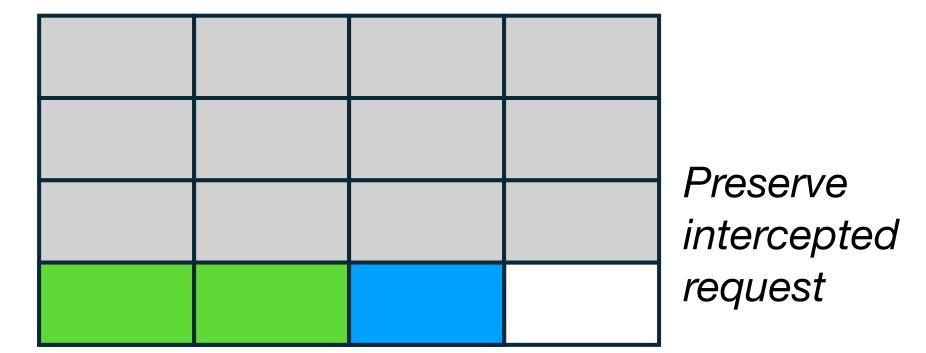


Other running requests waiting for 4x+ more time!

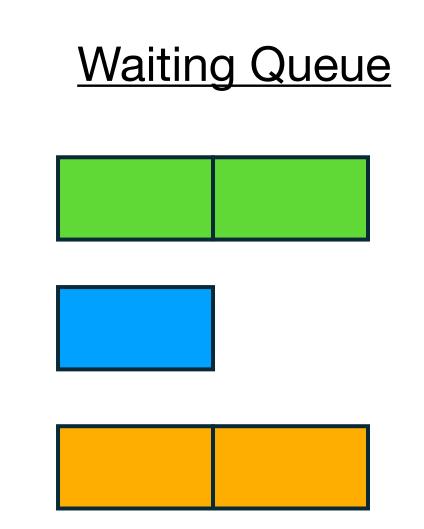
Typical recomputation time: 200+ms

Strategy 2: Preserve



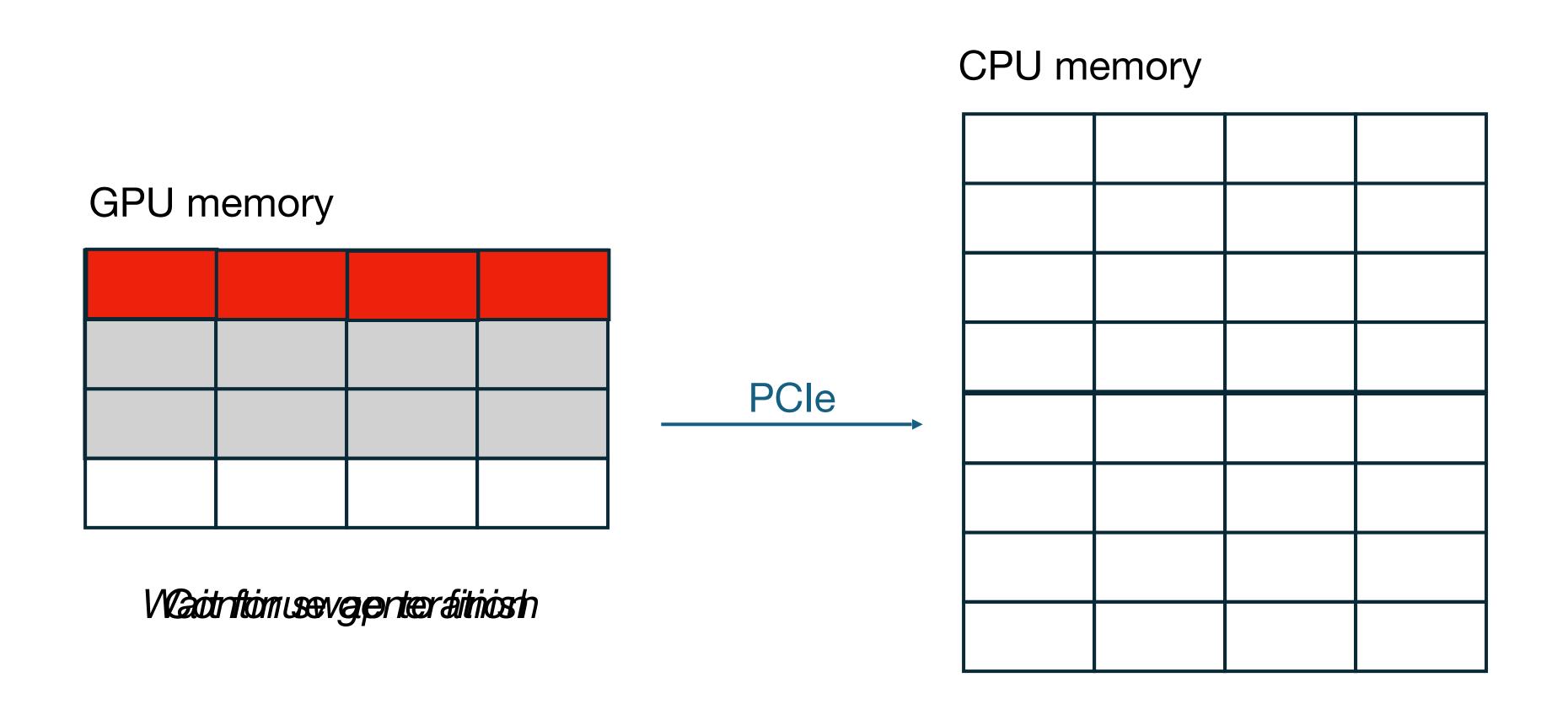


Prefill other requests



NOT ENOUGH MEMORY

Strategy 3: Swap to CPU memory



Each strategy has different tradeoffs

• Discard + recompute KV

Recomputing cost + stalls running requests

Preserve KV in GPU

Memory unused during interception

Swap KV to CPU memory

Swap bandwidth limited

Which is the best strategy for a given request?

How do we determine this?

Minimizing Waste

A unified measurement for all strategies

Waste = unused GPU memory * time

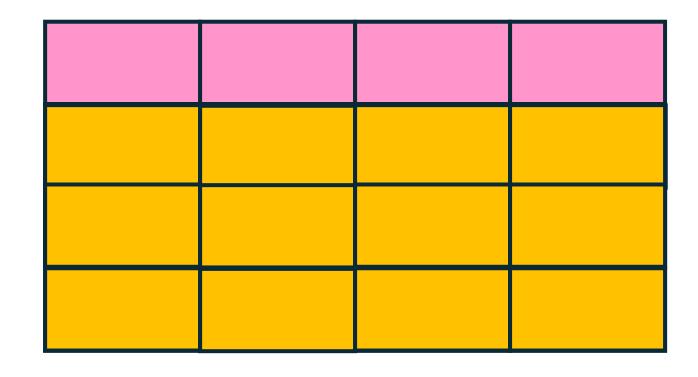
Accounting for intercepted request and remaining request

For each intercepted request, choose the minimal-waste strategy

Can we improve the existing strategies further?

MinWaste Discard: Chunk Recomputation

GPU memory



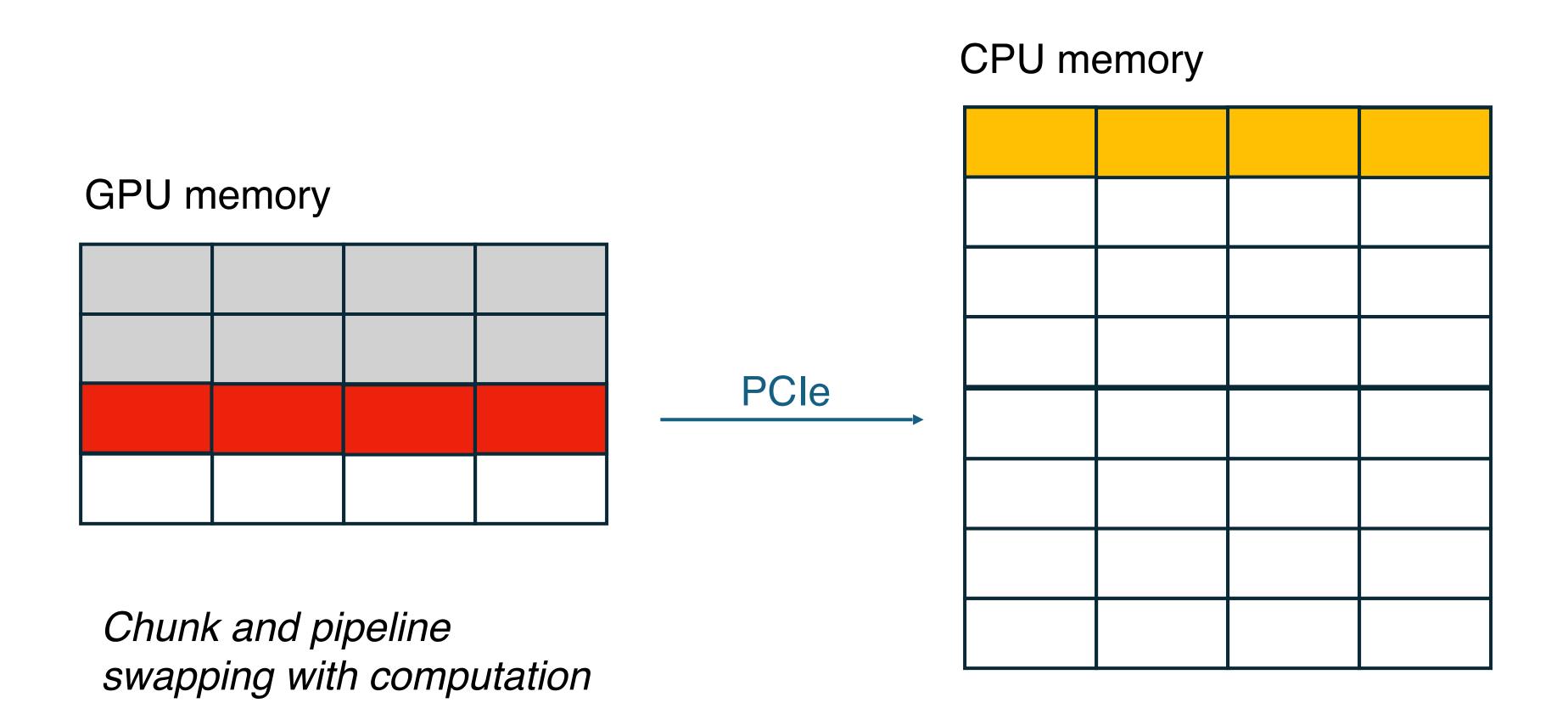
Discarded memory reused by other requests

Running requests

Recompute token states (prefill)

Recompute one chunk at a time to not stall other running requests

MinWaste Swap: Hide Swap Latency

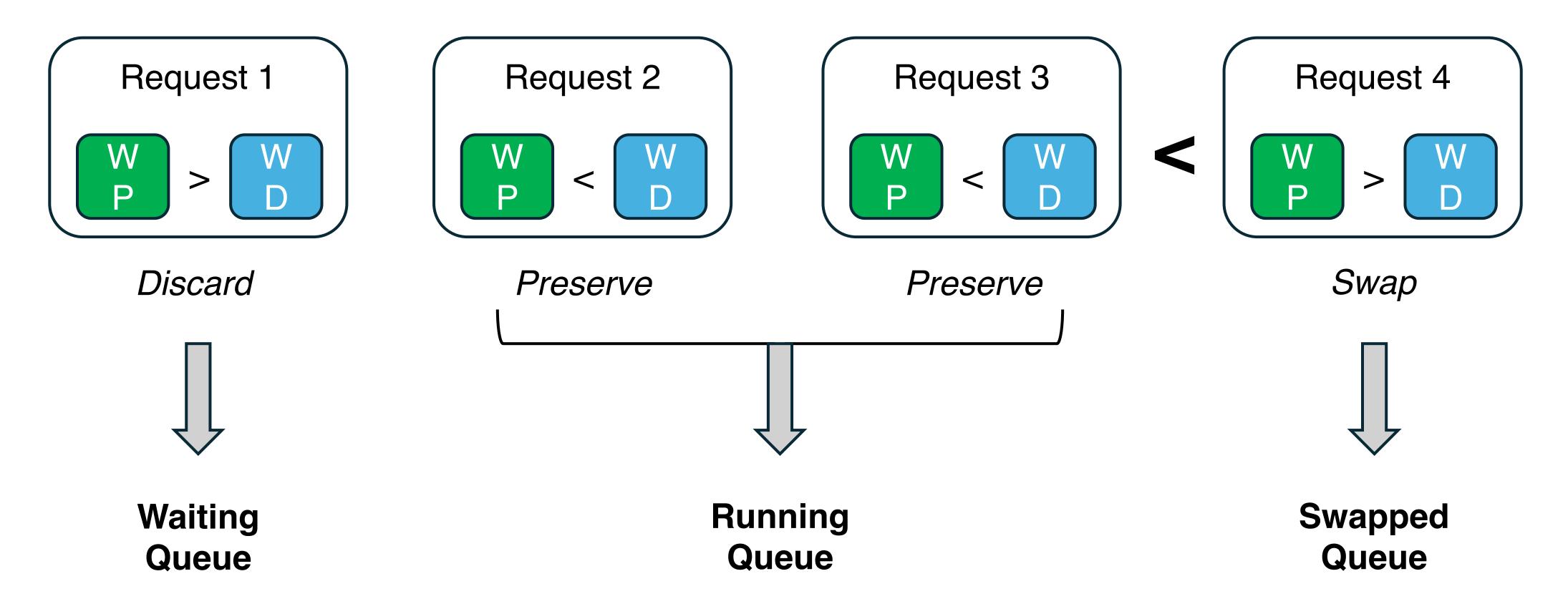


Swap (within a limit) is completely free!

Scheduling Across Requests

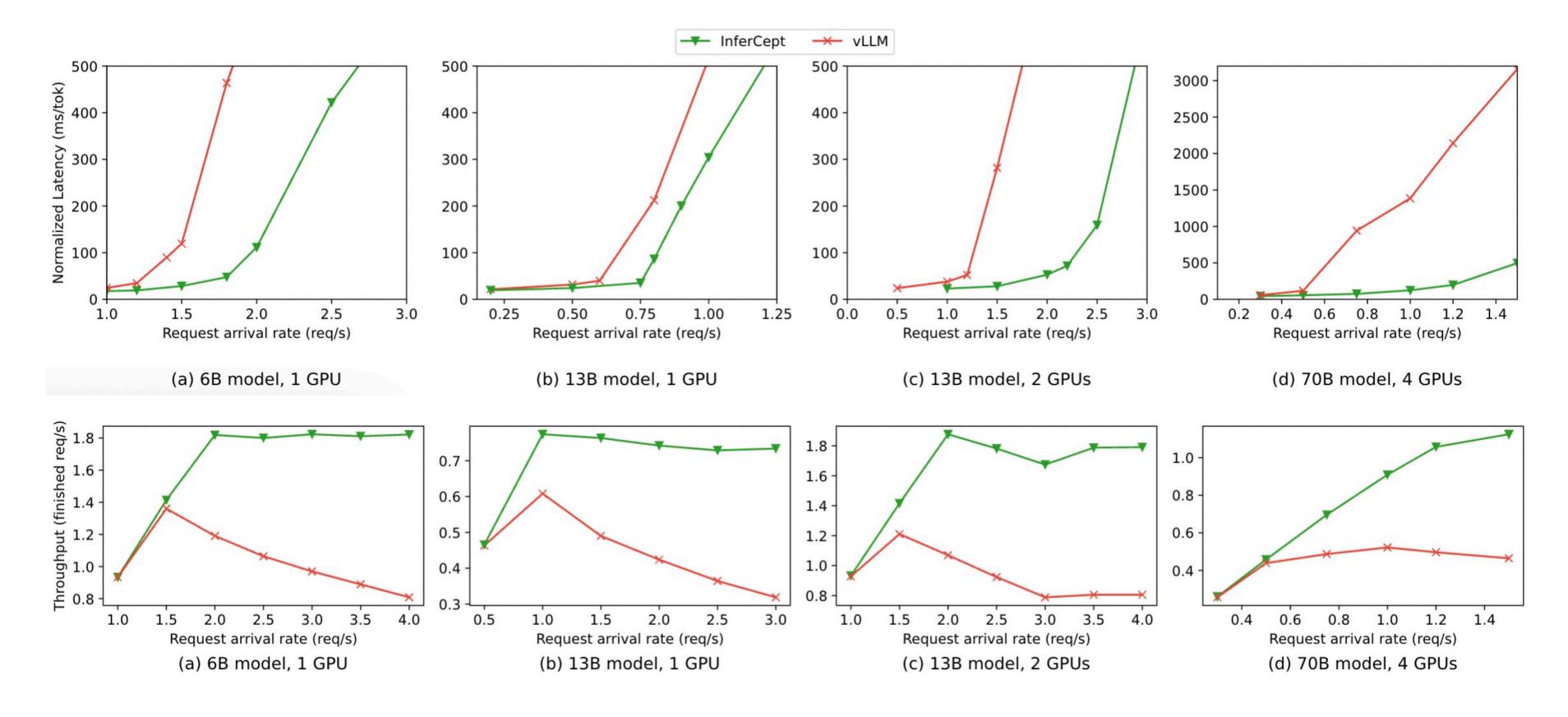
Out of all intercepted requests in an iteration

- Use the swap budget for otherwise most wasteful requests
- Choose the smaller waste of preserve and discard for the remaining requests



Results

- 6B GPT-J and 13B Vicuna calling six different tools
- Sustains 1.6x to 2x higher request load and 1.3x to 12x lower latency



InferCept Takeaways

- Model calls are increasingly accompanied by external tool and data calling
- KVs need to be properly managed when external entities intercept model calls
- Three basic strategies, each with pros and cons
- InferCept: first work to manage model/non-model interactions at the system level

Conclusion

- 2025 is the year of Al agents
- Many efforts in agent development and development frameworks and tools
- Al agent infrastructure is largely an unexplored area
- WukLab and GenseeAl are building a cross-stack Al agent platform
- Stay tuned and follow gensee.ai and mlsys.wuklab.io



