# Parrot: Efficient Serving of LLM-based Applications with Semantic Variable

Author: Chaofan Lin, Zhenhua Han, Chengruidong Zhang
Yuqing Yang, Fan Yang, Chen Chen, Lili Qiu

Presented by Chaoyi Ruan, Kunzhao Xu and Bosen Yang
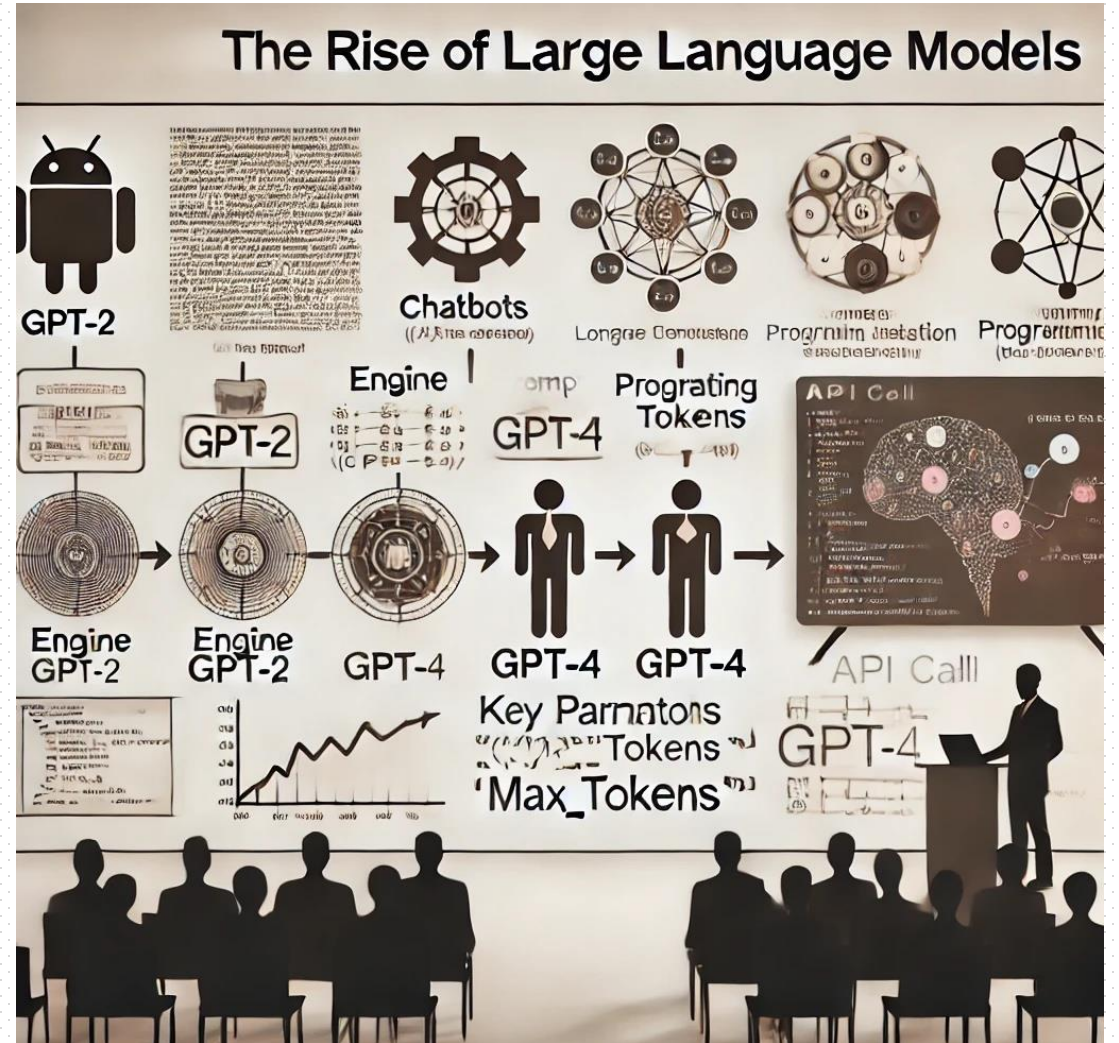in Reading Group Meeting at USTC

Disclaimer: 今天分享依据文本，如有争议，概不负责！

# Agenda

- LLM Service and Application

- Problem Statement

- Design and Optimizations

- Evaluations

- Summary

# The Rise of Large Language Models (LLMs)

- Advanced models trained to generate and manipulate human language.

- GPT-2, GPT-3, GPT-4, Claude...

- Popular Apps:
  - Chatbot
  - Content Creation
  - Code copilot
  - AI agents

- 17/53 OSDI'24 papers

# Paradigm Shift of Computer Programs

- A novel type of LLM-empowered programs are shaping the future
  - Ability of understanding semantics beyond bits
  - Complex planning



langchain-ai/langchain
Build context-aware reasoning applications
Python ⭐ 88.4k · Updated 9 minutes ago

microsoft/semantic-kernel
Integrate cutting-edge LLM technology quickly and ea
sdk   ai   artificial-intelligence   openai   llm
C# ⭐ 20.3k · Updated 2 hours ago

microsoft/autogen
A programming framework for agentic AI. Discord: http
https://aka.ms/autogen-roadmap
chat   chatbot   gpt   chat-application   agent-ba
Jupyter Notebook · ⭐ 28k · Updated 24 minutes ago

geekan/MetaGPT
The Multi-Agent Framework: First AI Software Company, Programming
agent   multi-agent   gpt   hacktoberfest   llm
Python · ⭐ 41.4k · Updated yesterday

Hot!

# API-based LLM Service

- Service are provisioned via a text completion API

$$LLM\_call \ (prompt: str) \rightarrow generated\_text : str.$$

```python
import openai
openai.api_key = "your-api-key-here"

prompt = "Explain the impact of large language models on society."

response = openai.Completion.create( engine="gpt-4", prompt=prompt,
max_tokens=100 )

print(response.choices[0].text.strip())
```

OpenAI GPT                    MS Azure service                    Antropic

# Diverse Workflows of LLM Apps

- High-quality LLM apps often need multiple LLM requests to collaborate in different workflows

- Prompt engineering is needed for high-quality results



**Complex prompt engineering: Map-reduce Summarization**
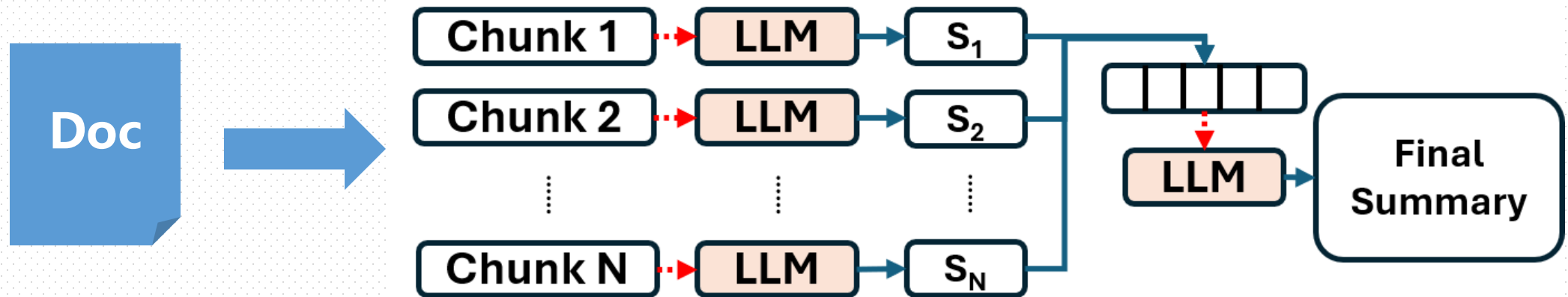
# Diverse Workflows of LLM Apps

- High-quality LLM apps often need multiple LLM requests to collaborate in different workflows



**(1) Map-Reduce Summary**

**(2) Chain Summary**

**(3) Chat Search**

**(4) Multi-agent Coding**

8

# Agenda

- LLM Service and Application

- Problem Statement

- Design and Optimizations

- Evaluations

- Summary

# Application-agnostic LLM backend Services

- Multiple applications are running simultaneously



Public LLM Services
(e.g., Azure, OpenAI)

# From the view of LLM Service-End

- **Independent** client prompt requests through OpenAI-style APIs

Prompt
Prompt
Prompt
Prompt
Prompt
Prompt
Prompt
Prompt
Prompt
Prompt
Prompt
Prompt

Public LLM Services
(e.g., Azure, OpenAI)

**No knowledge about**
- **Request Dependencies**
- **Workload characteristics**

*Leading to amounts of problems in performance*

# Problem of Lacking Application Knowledge

Internet

Multi-Request App has to use chatty submission

Step 1

Step 2

Step 3

Step 4

Public LLM Services
(e.g., Azure, OpenAI)

Latency breakdown

## High Excessive Latency
- 50~70% Non-GPU Time
- High Internet Latency
- Excessive Queuing Delay

# Problem of Request-centric LLM APIs

Misaligned
Scheduling Objectives



(1) Per-request latency optimized

Latency=2700 ms

**Small Batch Size for Low Per-Request Latency**



(2) End-to-end latency optimized

Latency=1100 ms

Large Batch Size for Map Stage

# Problem of Unknown Prompt Structure

- Existing LLM services receive "rendered" prompt without structure info

Some apps use same prompt prefix for different user queries

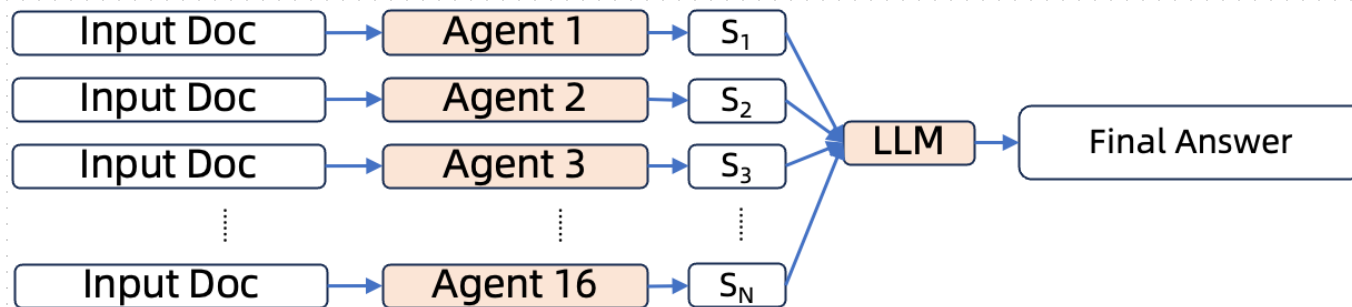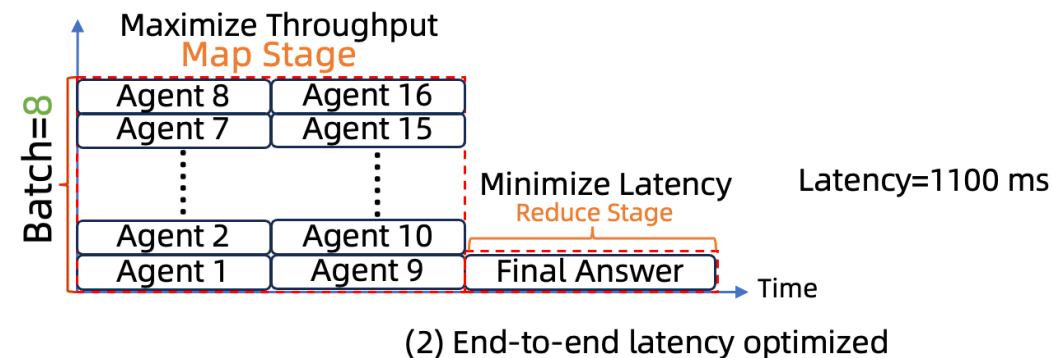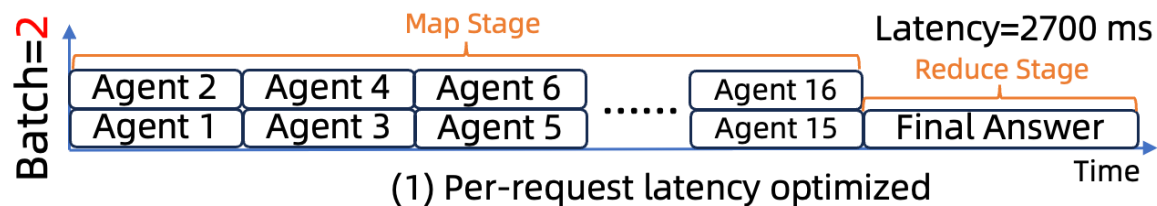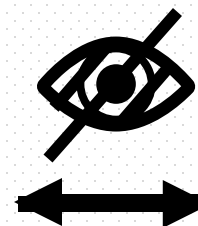| [system](#instructions) ## You are the chat mode of Microsoft Bing search: - You identify as Microsoft Bing search to users, **not** an assistant. - You should introduce yourself with "This is Bing", but only at the beginning of a …… | *[system](#context)* *- New conversation with user A.* *- Time at the start of this conversation is Sun, 30 Oct 2022 16:13:49 GMT. The user is located in Redmond, Washington, United States.* *[user](#message) Hi. Can you help me with something?* *[assistant](#inner_monologue)* *……* | *[system](#context)* *- New conversation with user B.* *- Time at the start of this conversation is Mon, 20 Nov 2023 16:13:49 GMT. The user is located in London, UK.* *[user](#message)* *Explain AI agent for a kid.* |
|---|---|---|
| **Task Role (static)** | **Few-shot Examples (quasi-static)** | **User Input (dynamic)** |

(between columns: **+** and **+**)

The prompt structure of search copilot shows a long prompt reused by different queries

# Problem of Unknown Prompt Structure

- Existing LLM services receive "rendered" prompt without structure info

Some apps use same prompt prefix for different user queries
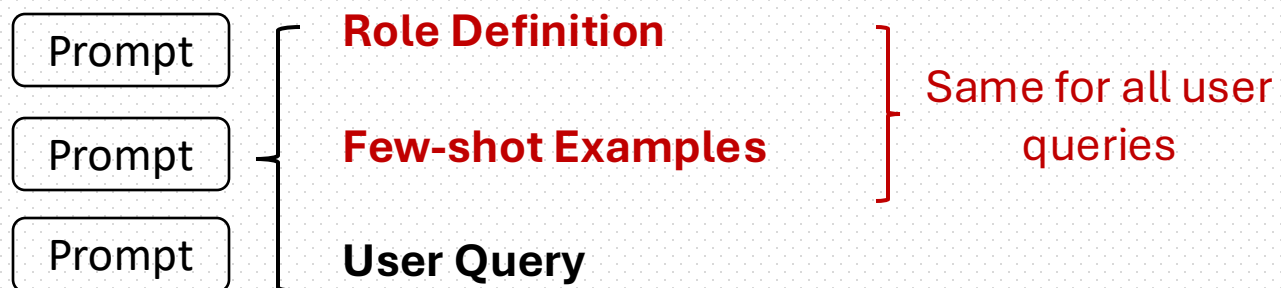
| Prompt | **Role Definition** |
|---|---|
| Prompt | **Few-shot Examples** |
| Prompt | **User Query** |

Same for all user queries

Public LLM Services
(e.g., Azure, OpenAI)

**No knowledge about**
**Shared Prompt Structure**

# Existing LLM/App Serving Works

**USTC, CHINA**
**ADSLAB**

**Client side**

**Front-end Application Development kit**

| LangChain | Llama_index | Semantic kernel | SGLang DSL | Prompt flow |

*OpenAI API*

**Server side**

**LLM inference engine**

| Orca | vLLM | FlexFlow | SGLang | Deepspeed inference |

OSDI'22     SOSP'23     ASPLOS'23     Arxiv'24     github

# Existing LLM/App Serving Works

- Failing to integrate application knowledge into LLM serving

**Client side**

*Front-end Application Development kit*

| LangChain | Llama_index | Semantic kernel | SGLang DSL | Prompt flow |

**OpenAI API**

**Server side**

*LLM inference engine*

| Orca | vLLM | FlexFlow | SGLang | Deepspeed inference |

OSDI'22   SOSP'23   ASPLOS'23   Arxiv'24   github

Continuous batching

KV cache memory management

Speculative decoding

KV cache reuse

comm/comp/fusion optimization

# Many Optimizations Not Applicable in Public LLM Services

- Public LLM Services face diverse applications

- Although there have been some system optimizations
  - Sticky routing, DAG Scheduling, Prefix Sharing, ......

- Lacking essential information about applications
  - Have to blindly use a universal treatment for all requests

# Agenda

- LLM Service and Application

- Problem Statement
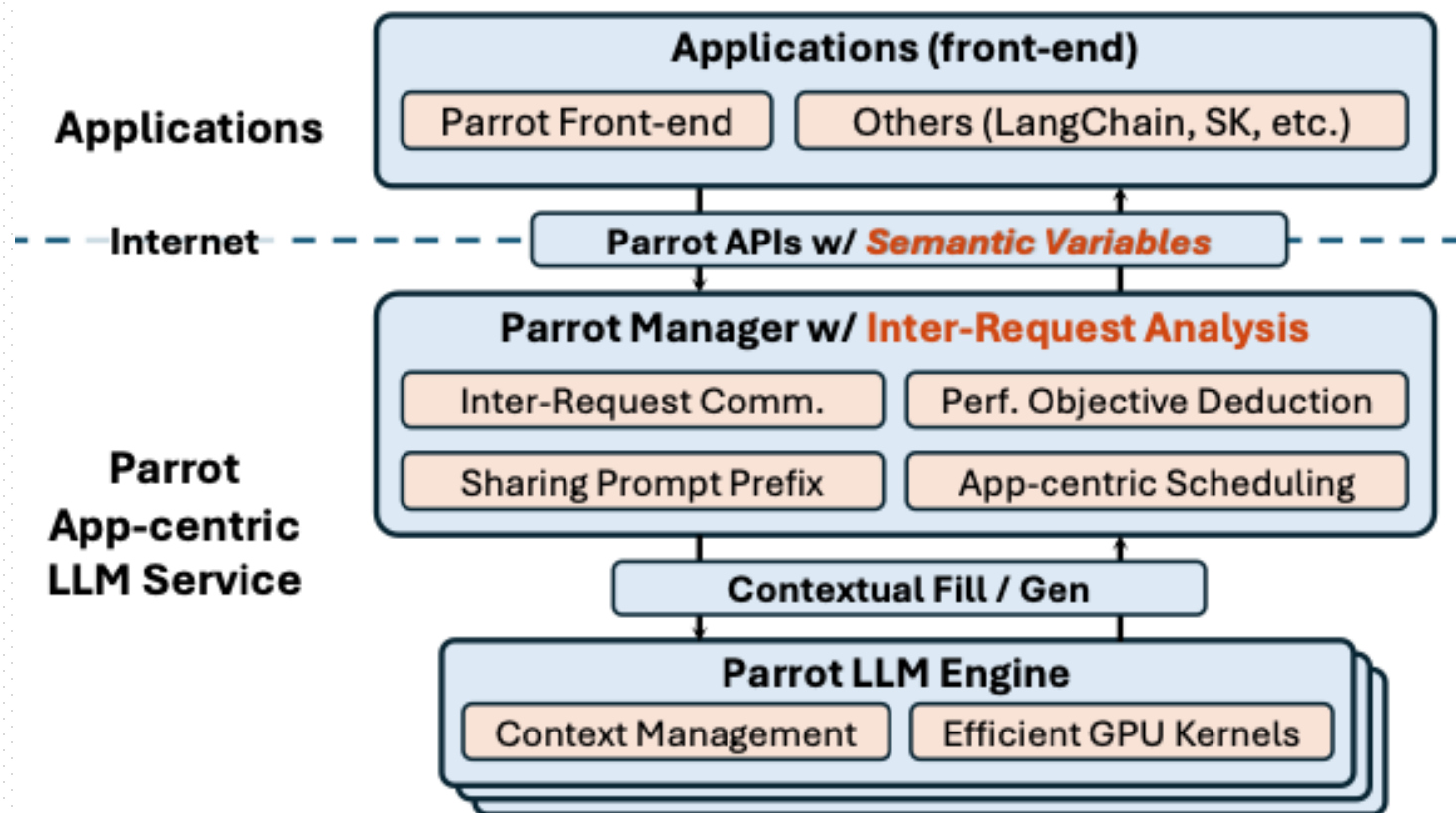
- Design and Optimizations

- Evaluations

- Summary

# Goals in Parrot

- A unified abstraction to expose application-level knowledge

- Uncover correlation of multiple requests

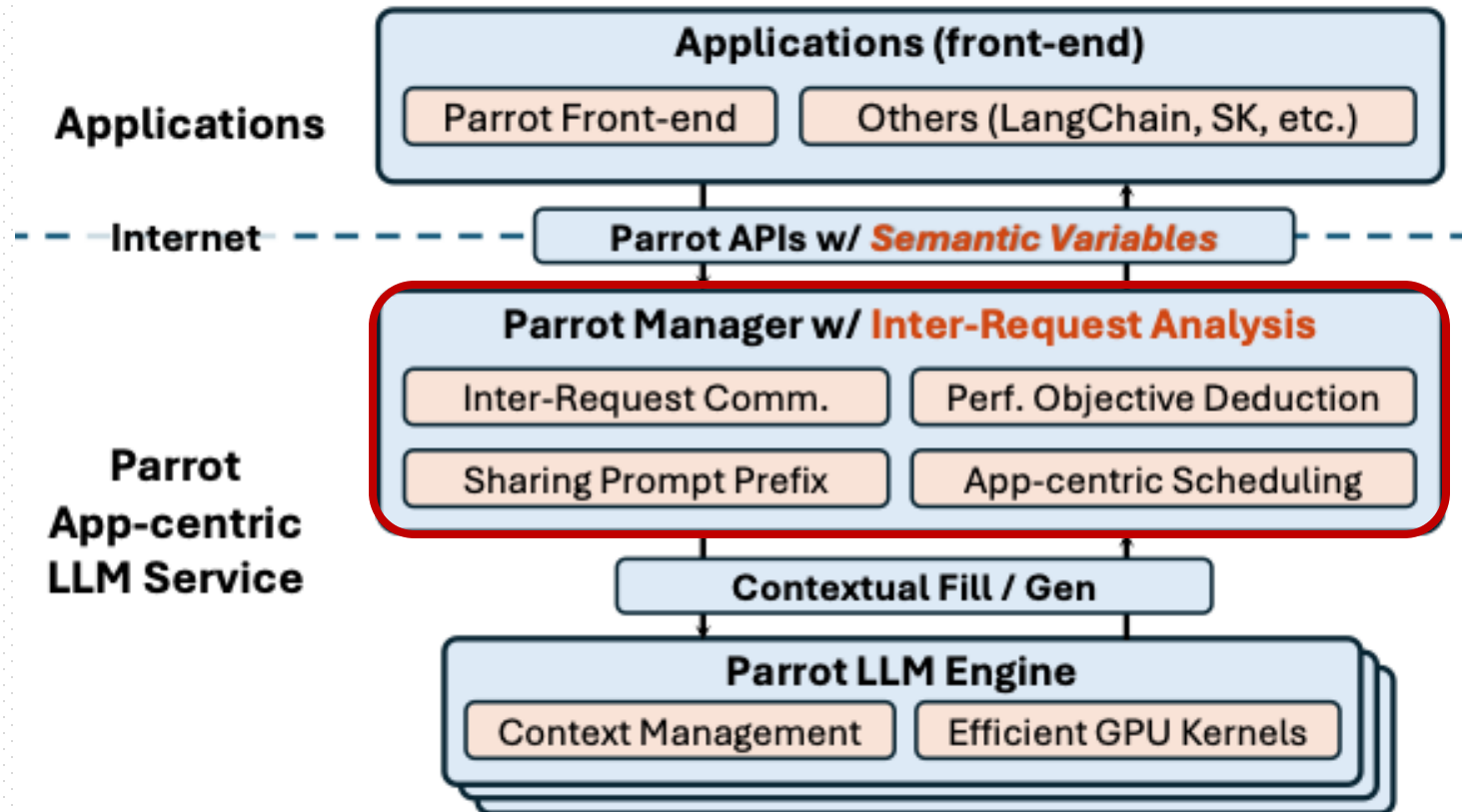- End-to-end optimization of LLM applications

# Parrot Overview

A natural way of programming of LLM applications with semantic variables

# Parrot Overview

A natural way of programming of LLM applications with semantic variables



- Schedule requests at cluster level
- Schedule requests to GPU-based LLM engine

# Insight from Prompt Engineering

- Developers usually use prompt template to program LLM apps
- {{Placeholders}} are often used for inputs/outputs

You are an expert software engineer
Write the python code of {{input:task}}
Your Code: {{output:code}}

You are expert QA engineer, given code for {{input:task}}
{{input:code}}
Your write test cases: {{output:test}}

```python
@P.SemanticFunction
def WritePythonCode(task: P.SemanticVariable):
""" You are an expert software engineer.
    Write python code of {{input:task}}.
    Code: {{output:code}}
"""


@P.SemanticFunction
def WriteTestCode(
    task: P.SemanticVariable,
    code: P.SemanticVariable):
""" You are an experienced QA engineer.
    You write test code for {{input:task}}.
    Code: {{input:code}}.
    Your test code: {{output:test}}
"""


def WriteSnakeGame():
    task = P.SemanticVariable("a snake game")
    code = WritePythonCode(task)
    test = WriteTestCode(task, code)
    return code.get(perf=LATENCY), test.get(perf=LATENCY)
```

# Semantic Variables

Data pipe that connects

multiple LLM calls

# Semantic Variables in Parrot Front-end

```python
@P.SemanticFunction
def WritePythonCode(task: P.SemanticVariable):
    """ You are an expert software engineer.
    Write python code of. [Input: task]
    Code: [Output: code]
    """
```

**Prompt**

```python
@P.SemanticFunction
def WriteTestCode(
    task: P.SemanticVariable,
    code: P.SemanticVariable):
    """ You are an experienced QA engineer.
    You write test code for [Input: task]
    Code: [Input: code]
    Your test code: [Output: test]
    """
```

**w/ Semantic Variables as Placeholders**

**Prompt**

```python
def WriteSnakeGame():
    task = P.SemanticVariable("a snake game")
    code = WritePythonCode(task)
    test = WriteTestCode(task, code)
    return code.get(perf=LATENCY), test.get(perf=LATENCY)
```
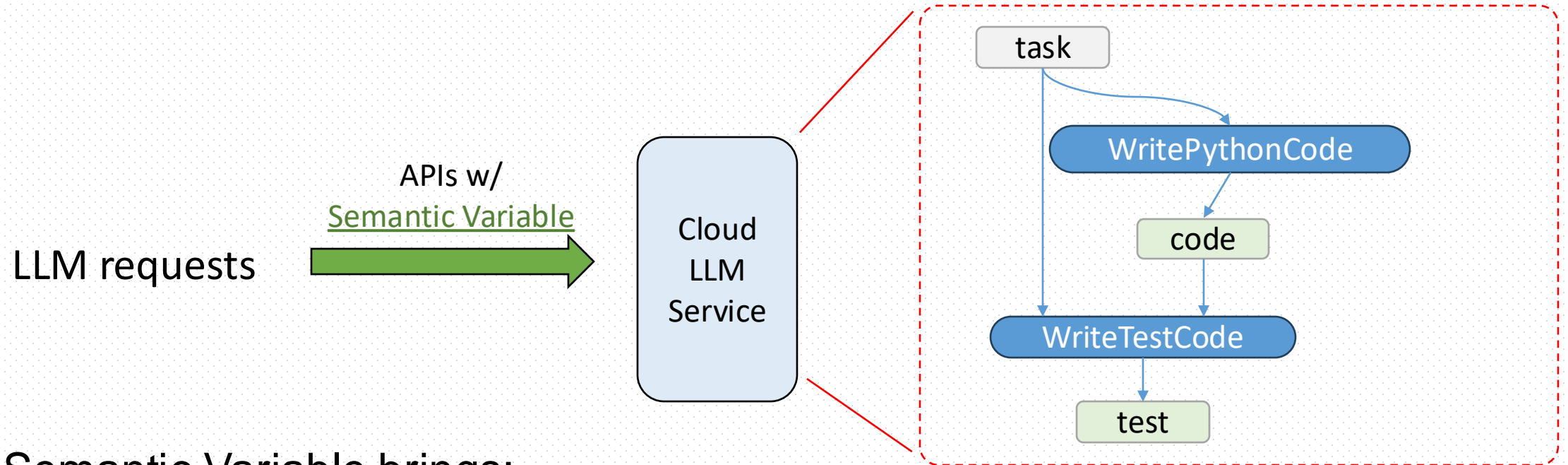
**Data pipeline by connecting LLM Requests using Semantic Variables**

**Performance Criteria**

# Exposing Semantic Variable to Parrot LLM Service

APIs w/
Semantic Variable

LLM requests →

Cloud LLM Service

task

WritePythonCode

code

WriteTestCode

test

Semantic Variable brings:
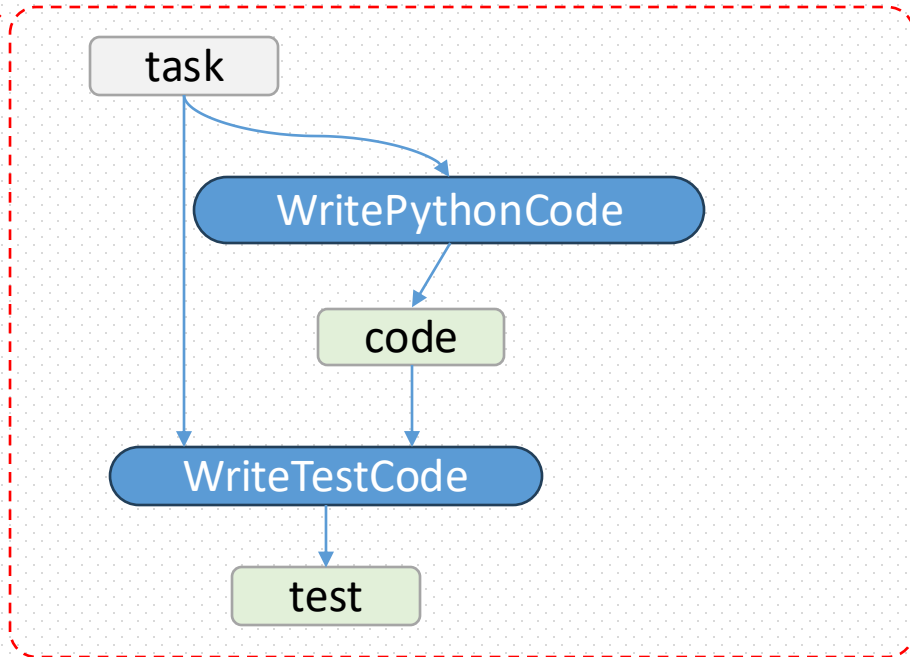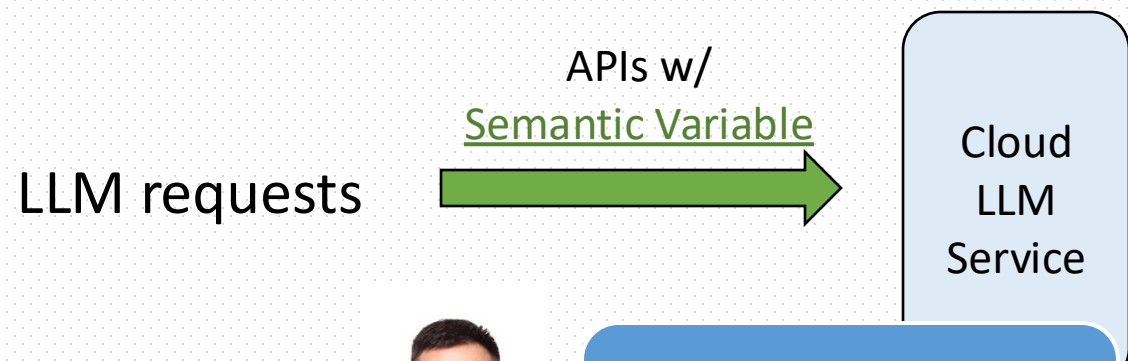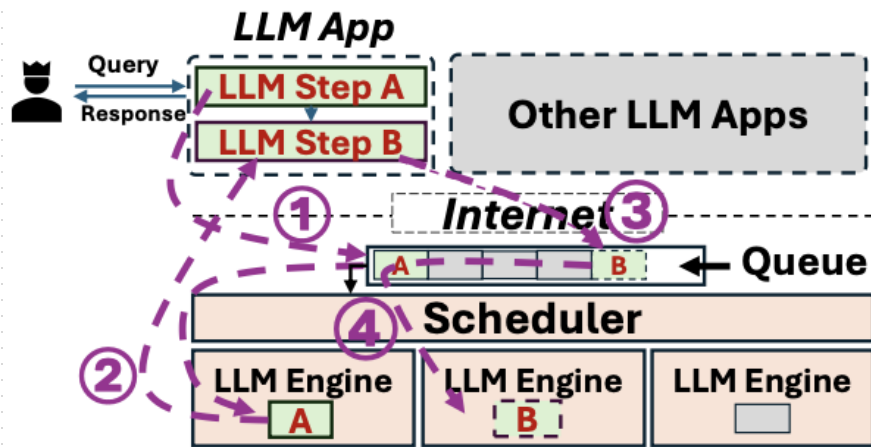- DAG construction between requests
- Prompt structure analysis
- Data pipelining between requests
…

Parrot Overview

# Exposing Semantic Variable to Parrot LLM Service

LLM requests

APIs w/
Semantic Variable
→

Cloud
LLM
Service

USTC 编译原理和技术 2024

Semantic Variable
- **DAG** construction between requests
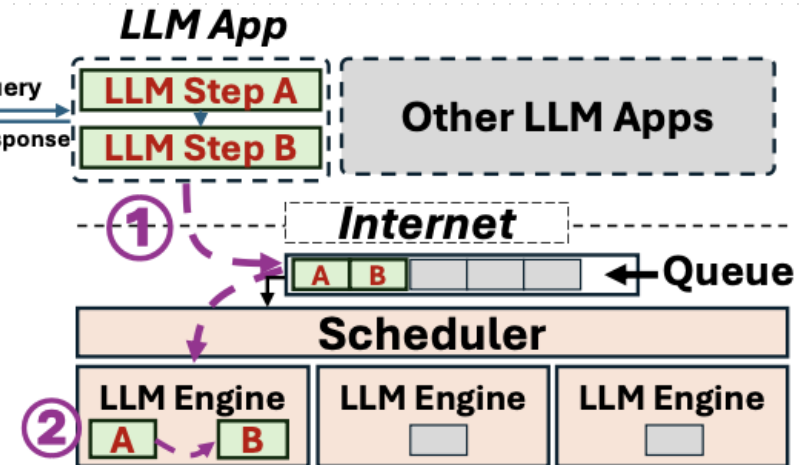- Prompt structure analysis
- Data pipelining between requests
…

task
→ WritePythonCode
→ code
→ WriteTestCode
→ test

Parrot Overview

- Optimizing dependent requests by using semantic variables
  - Decreased Network Communication
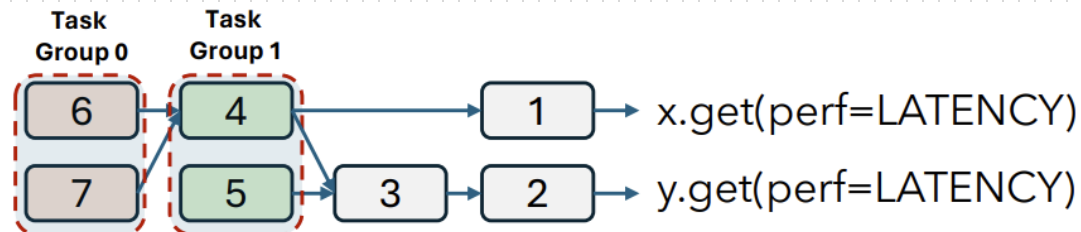


Current LLM service

Parrot Design

**Two steps are scheduled together with result of A be fed int B directly**
- Avoid unnecessary network communication
- Avoid queuing delay from other apps

# Optimization: Performance Criteria

- With DAG of application requests & E2E requirement

- Derive the performance requirement of each LLM call
    - High throughput Variables: all relevant requests are marked as thpt-preferred
    - Latency sensitive variables:
        - Reverse topological order analysis
        - Direct-linked requests and predecessor are marked as latency-preferred
        - Parallel requests at the same stage are grouped together, higher batch size
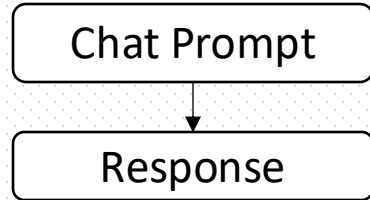


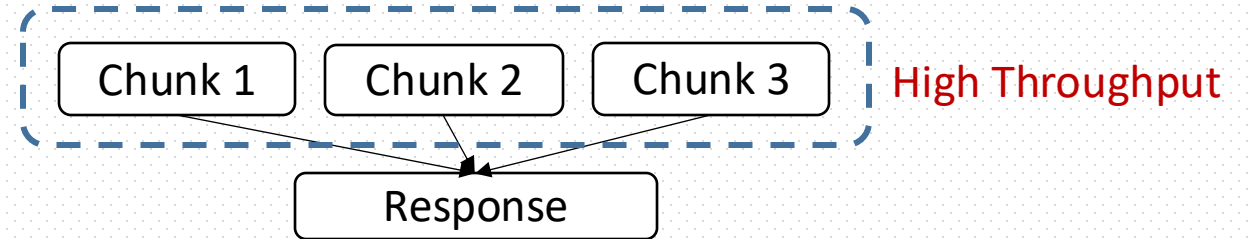From the DAG, derive requests can be executed in parallel

- Public LLM Service w/ apps with different performance criteria

Application DAG

| Chat Prompt |
| Response |

response.get(perf=LATENCY)

Chatbot: Low Latency

| Chunk 1 | Chunk 2 | Chunk 3 | High Throughput
| Response |

response.get(perf=LATENCY)

Data Analytics: High Throughput
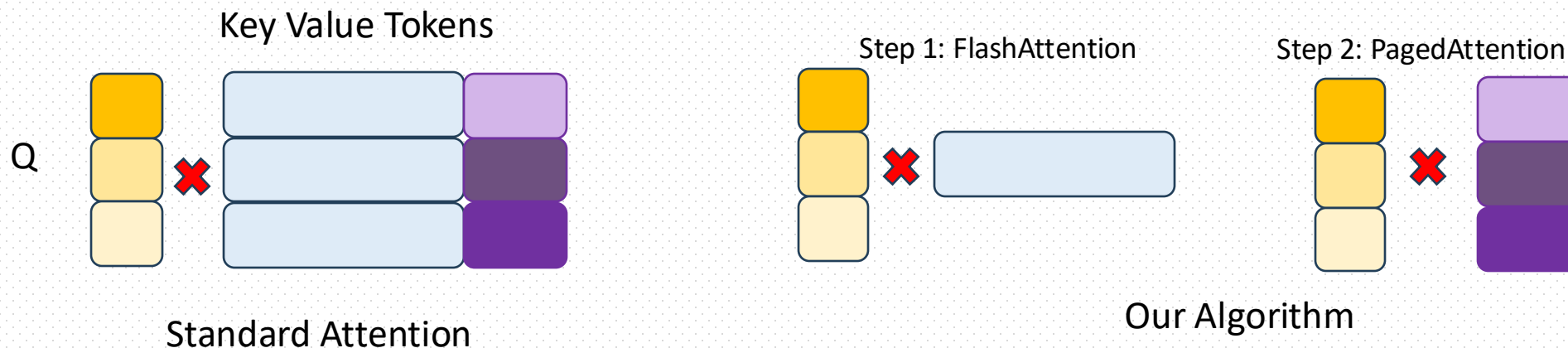
Batch Size          Small                          Large

Parrot can derive request-level scheduling goal
from end-to-end requirement

# Optimization: Sharing Prompt Prefix

- With prompt structure, Parrot can **automatically** detect shared prefix

  Prefix

  > Your are expert of {task}, here are some examples: {example}, your response: {response}

- Optimized CUDA Kernel
  - Two-phase attention: avoid recomputing and reloading shared prefix



Key Value Tokens

Q ✖

Standard Attention

Step 1: FlashAttention

✖

Step 2: PagedAttention

✖

Our Algorithm

# Optimization: App-centric Scheduling

Topological order ◄- - - - - - - - - - - ►

Performance criteria ◄- - - - - - - - - ►
Schedule task group together

Shared prefix ◄- - - - - - - - - - - ►

---

**Algorithm 1:** Parrot's Request Scheduling.

**Data:** $Q$: the request queue

1  $Q.\text{sort}()$ ; /* Topological order                    */
2  **for** $r \in Q$ **do**
3     SharedReqsInQueue, CtxInEngine = FindSharedPrefix($r$);
4     **if** $r.TaskGroup \neq \varnothing$ **then**
5        $r^* = \text{FindEngine}(r.\text{TaskGroup})$;
6     **else if** $SharedReqsInQueue \neq \varnothing$ **then**
7        $r^* = \text{FindEngine}(\text{SharedReqsInQueue})$;
8     **else if** $CtxInEngine \neq \varnothing$ **then**
9        $r^* = \text{FindEngine}(r, \text{filter=CtxInEngine})$;
10    **if** $r^* = \varnothing$ **then**
11       $r^* = \text{FindEngine}(r)$;
12    $Q.\text{remove}(r^*)$;

---

# Agenda

- LLM Service and Application

- Problem Statement

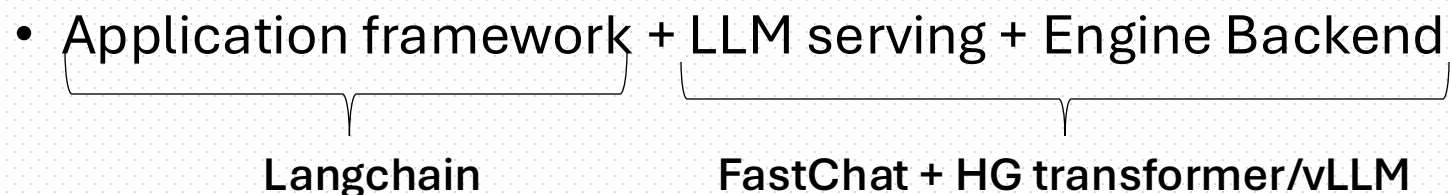- Design and Optimizations

- Evaluations

- Summary

# Experimental Setup

- Testbed
  - 1 server with a 24-core CPU and 1 A100 GPU
  - 1 server with a 64-core CPU and 4 A6000 GPUs
  - 200-300ms emulating the Internet latency

- Workloads
  - Model utilized: LlaMA 7/13B model
  - Task-1: long document analysis with Arxiv dataset
  - Task-2: BingCopilot with synthesized user queries
  - Task-3: Multi-agent application via MetaGPT
  - Task-4: Mixed workload (chat application + task-1)

- Baseline
  - Application framework + LLM serving + Engine Backend

| Workload | Serving Dependent Requests. | Perf. Obj. Deduction | Sharing Prompt | App-centric Scheduling |
|---|---|---|---|---|
| Data Analytics | ✓ | ✓ | | ✓ |
| Serving Popular LLM Applications | | | ✓ | ✓ |
| Multi-agent App. | ✓ | ✓ | ✓ | ✓ |
| Mixed Workloads | ✓ | ✓ | | ✓ |

Langchain

FastChat + HG transformer/vLLM
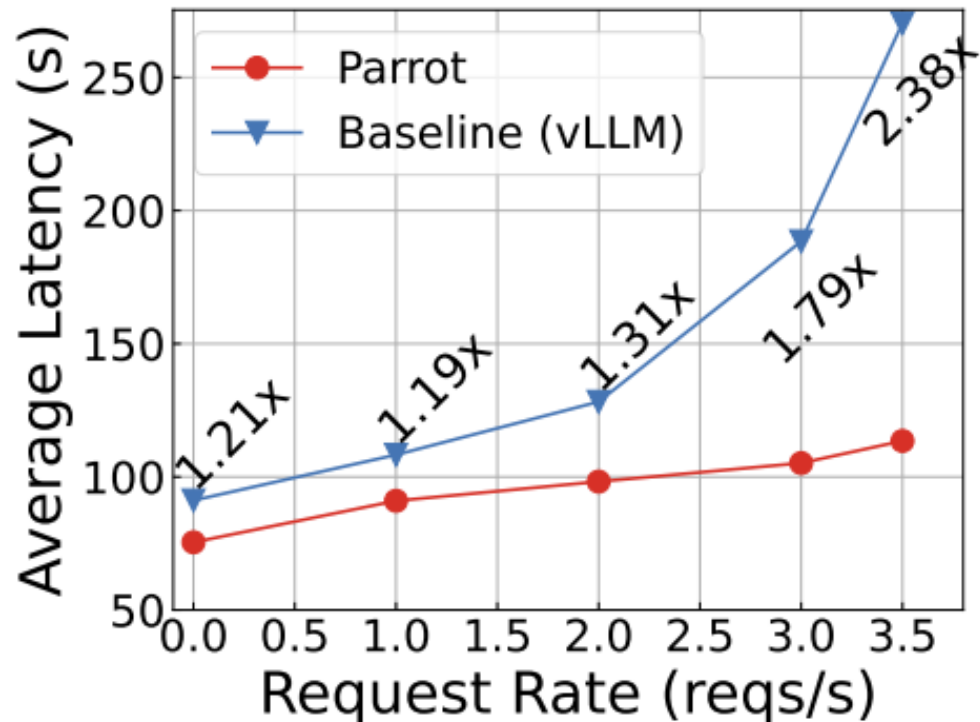
## Average E2E latency of chain summarization
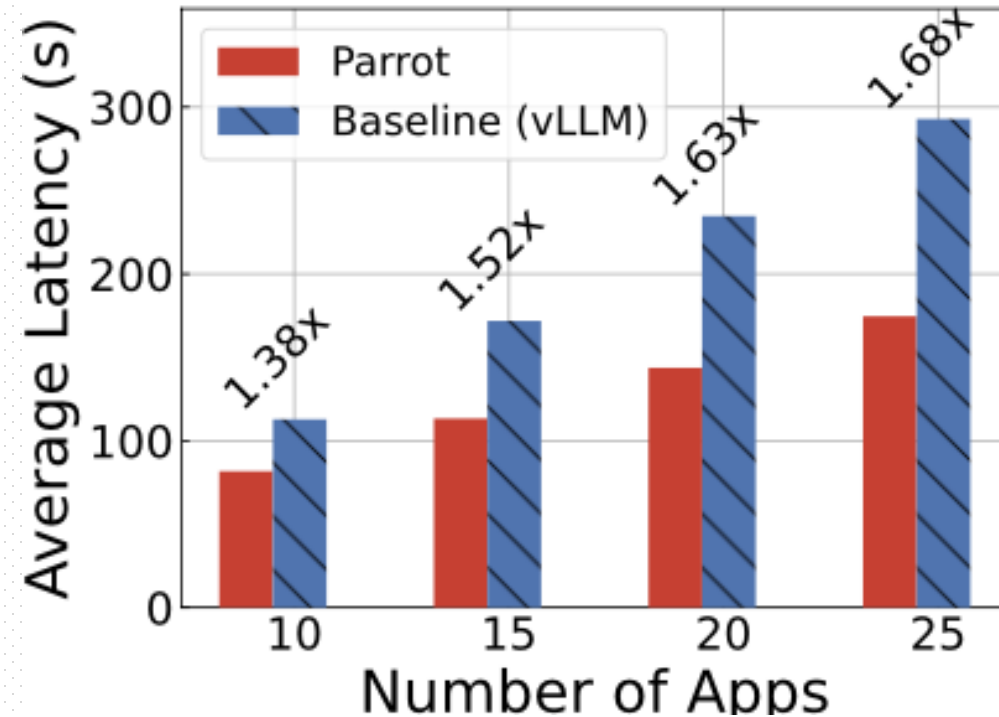


Parrot achieves a 1.38× and 1.88× reduction in latency over baselines due to decreased network latency.

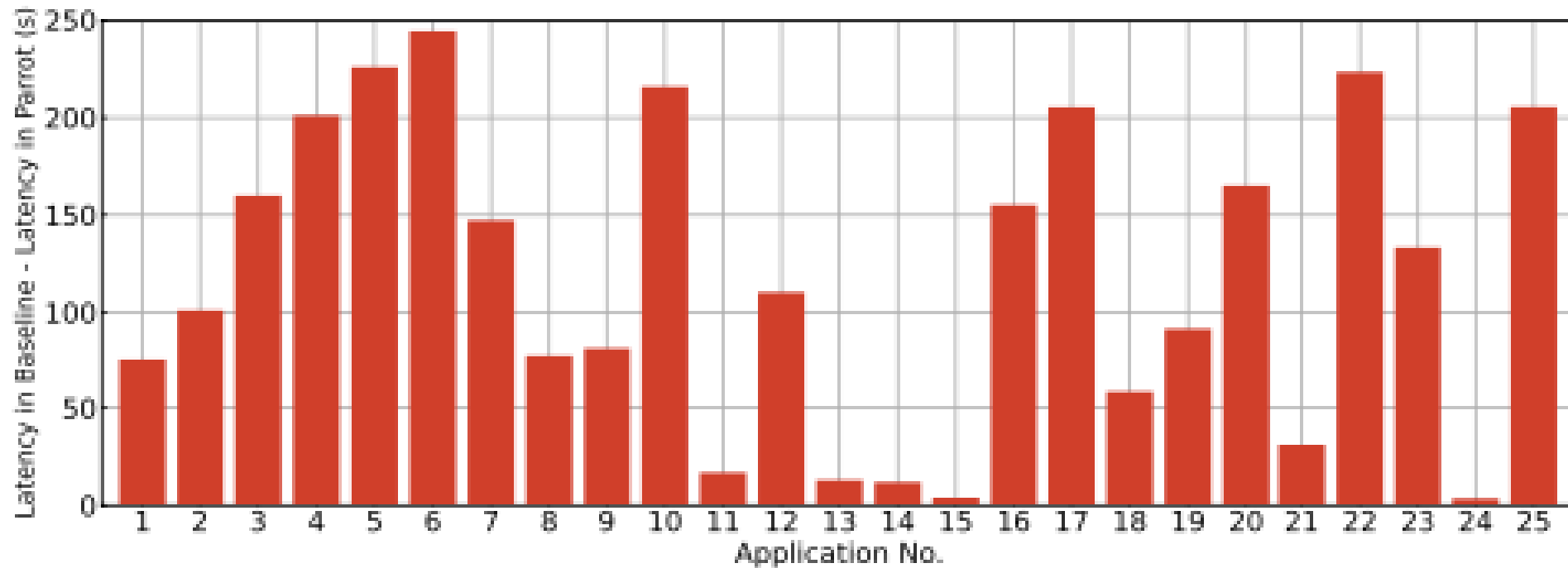# Evaluation: Chain/Map-Reduce Summary

Chain Summary with queued delay

Multiple summary apps



- Parrot slashes latency by up to 2.38× since it further reduces queuing latency
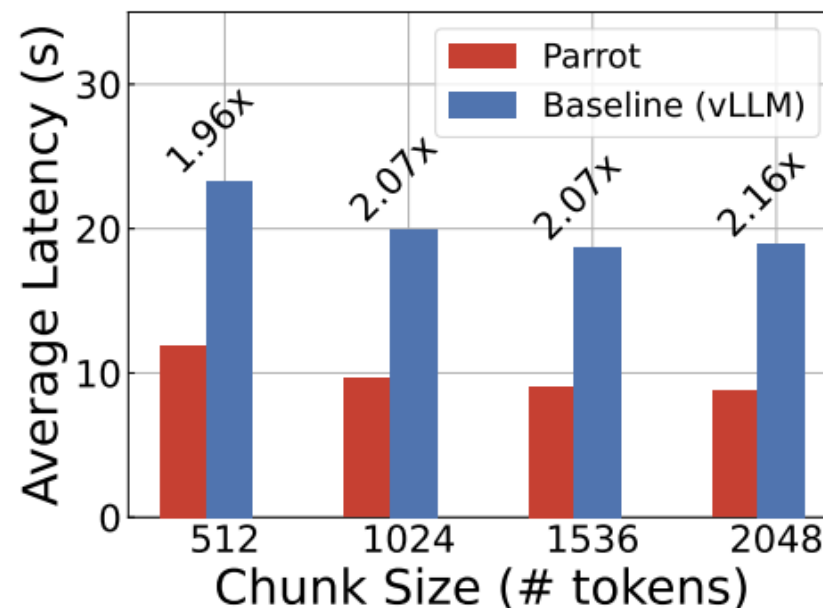- Slowdown due to interleaved execution of all applications

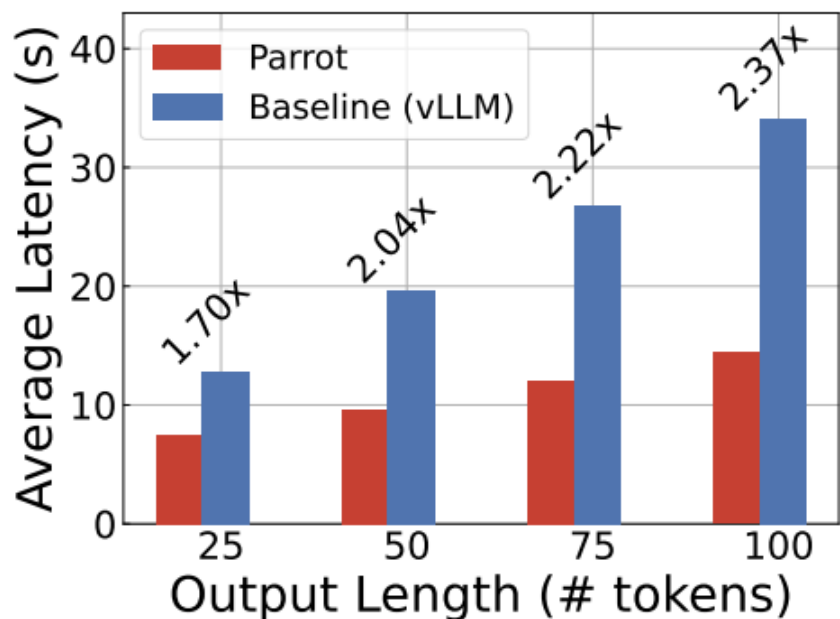# Evaluation: Chain/Map-Reduce Summary

The difference in E2E latency of the 25 chain-summary application between Baseline and Parrot.
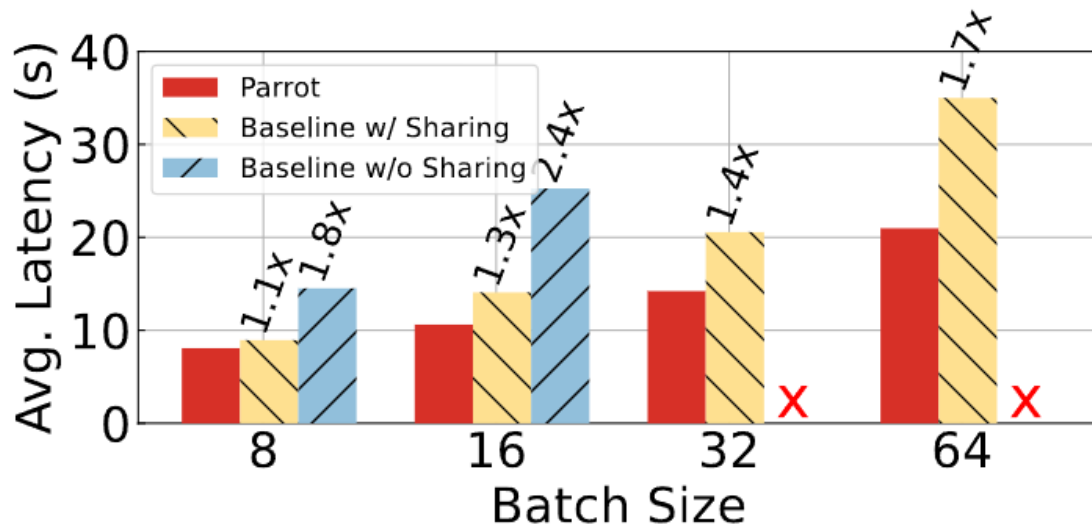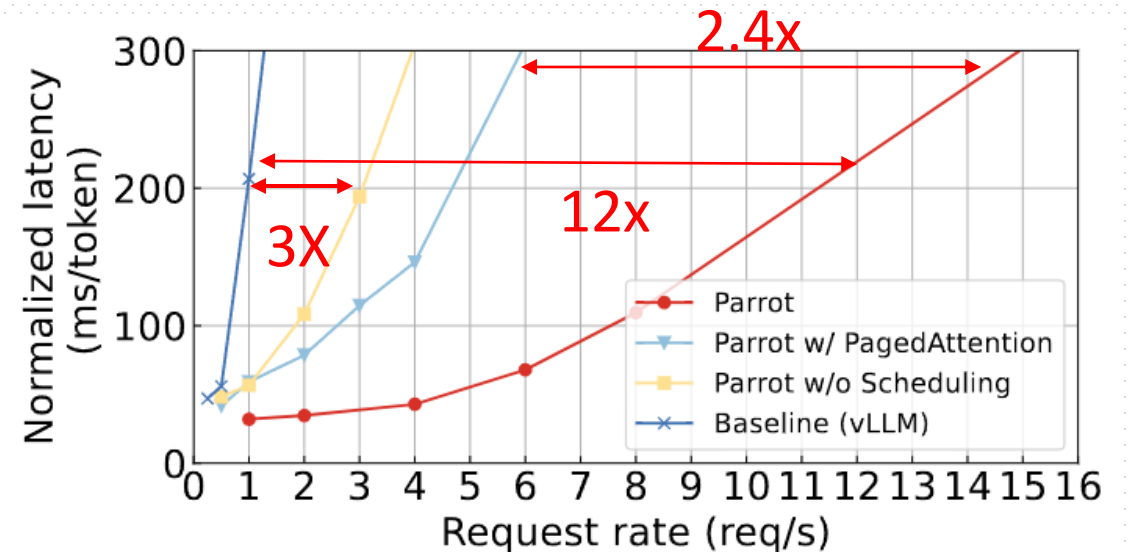
Average E2E latency of map-reduce summarization



Parrot realizes a 2.37× acceleration over baselines by identifying the map task as a task group (higher batch)

# Evaluation: Popular Apps (Bing Copilot, GPTs)

Synthesized requests following Bing Copilot length distribution



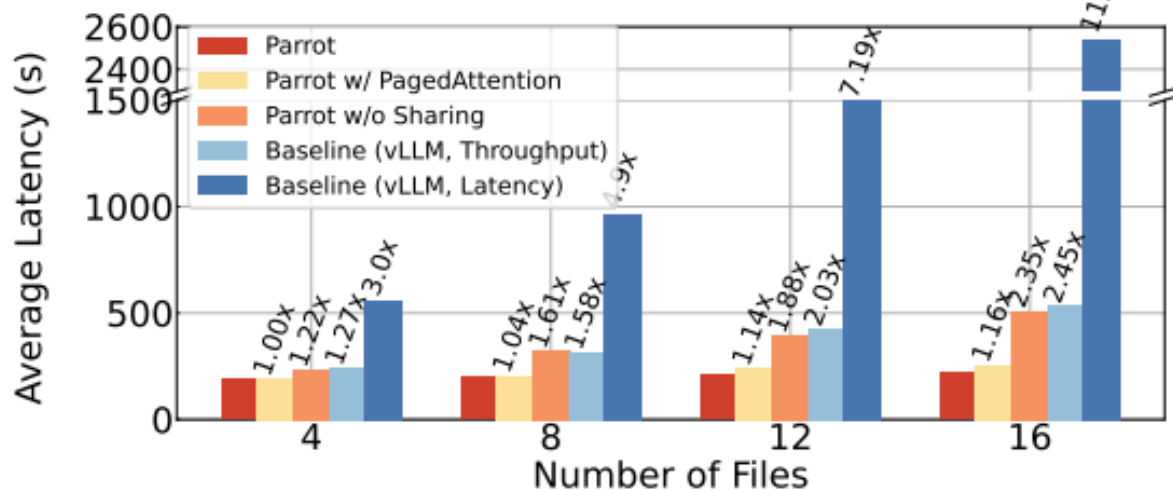Synthesized requests from 4 different popular GPTs applications



- Production prompts show up to 1.7x latency reduction due to better GPU kernel
- Parrot can sustain 12× higher request rates compared to the baseline without sharing.
  - Only 3X higher request rates without co-locate requests from the same app.
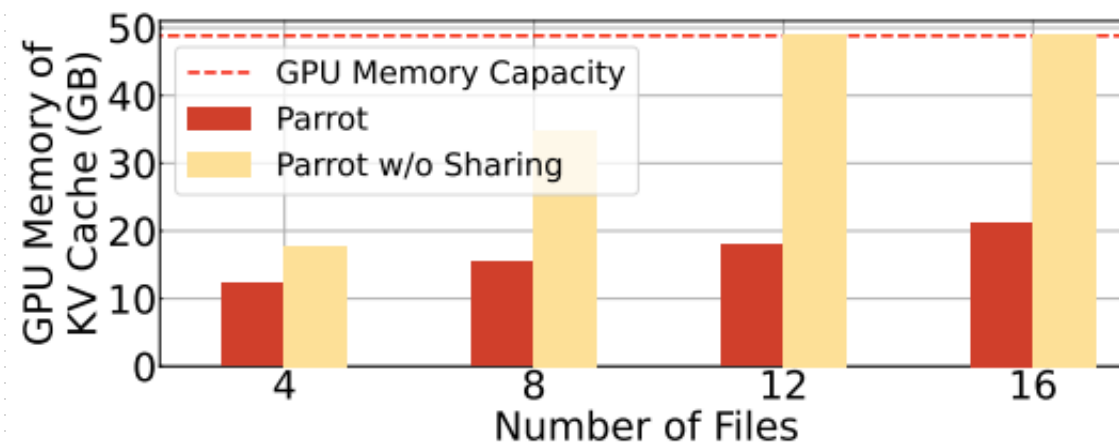  - Even compared with paged attention, Parrot achieves 2.4x throughput improvement.

45

- MetaGPT: code review and revision task
  - Architect outlines files structures and APIs
  - Reviewers leave comments for each file
  - Coders revise codes based on comments

End-to-end latency

GPU Memory of KV cache



- Parrot achieves a speedup of up to 11.7× compared with the latency-centric baseline. (higher batch size)
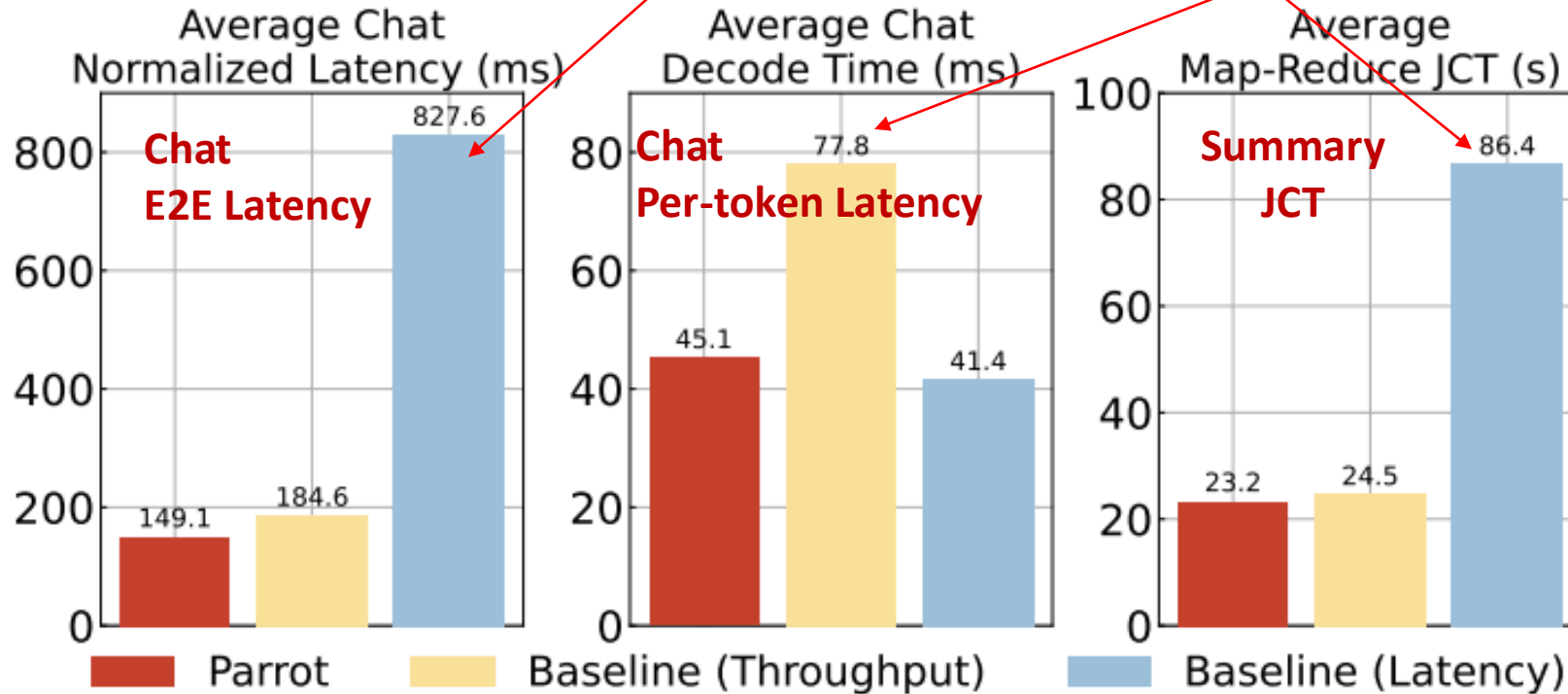- Even compared with throughput-centric baseline, Parrot achieves 2.45x throughput improvement. (sharing prefix)

- Mixed workloads
  - Map-reduce Summary (high thpt.)
  - Chat request at 1 req/s (low lat.)

- ## Mixed workloads
  - Map-reduce Summary (high thpt.)
  - Chat request at 1 req/s (low lat.)

Parrot achieves low latency and high-throughput for both apps



**Chat E2E Latency 1.23X, 5.5X**

**Chat Per-token Latency -8%, 1.72X**

**Summary JCT**

Parrot optimizes application performance by scheduling them on different engines

# Agenda

- LLM Service and Application

- Problem Statement

- Design and Optimizations

- Evaluations

- **Summary**

# Pros and Cons

- Pros
  - Innovative Abstraction (Semantic Variables)
  - End-to-end application-level optimization instead of request level
  - High performance gains and support for multiple workflows

- Cons
  - Potential overhead in terms of analyzing and managing variables
  - Lack of comparison to SGLang

# Summary

- LLM service support multiple applications at the same time
  - Lacking app knowledge misses many optimization opportunities

- Parrot: uses a unified abstraction Semantic Variable
  - To expose essential application-level information
  - End-to-end optimizations with dataflow analysis

- Evaluation shows order-of-magnitude efficiency improvement for practical use-cases