# SGLang v0.2: Faster Interface and Runtime for LLM inference (confidential)

Ying Sheng, Liangsheng Yin, Zhiqiang Xie

Aug - Dec. 2023

Jan. - now 2024

now - 2024

Early Stage: the "programming LLM" paradigm
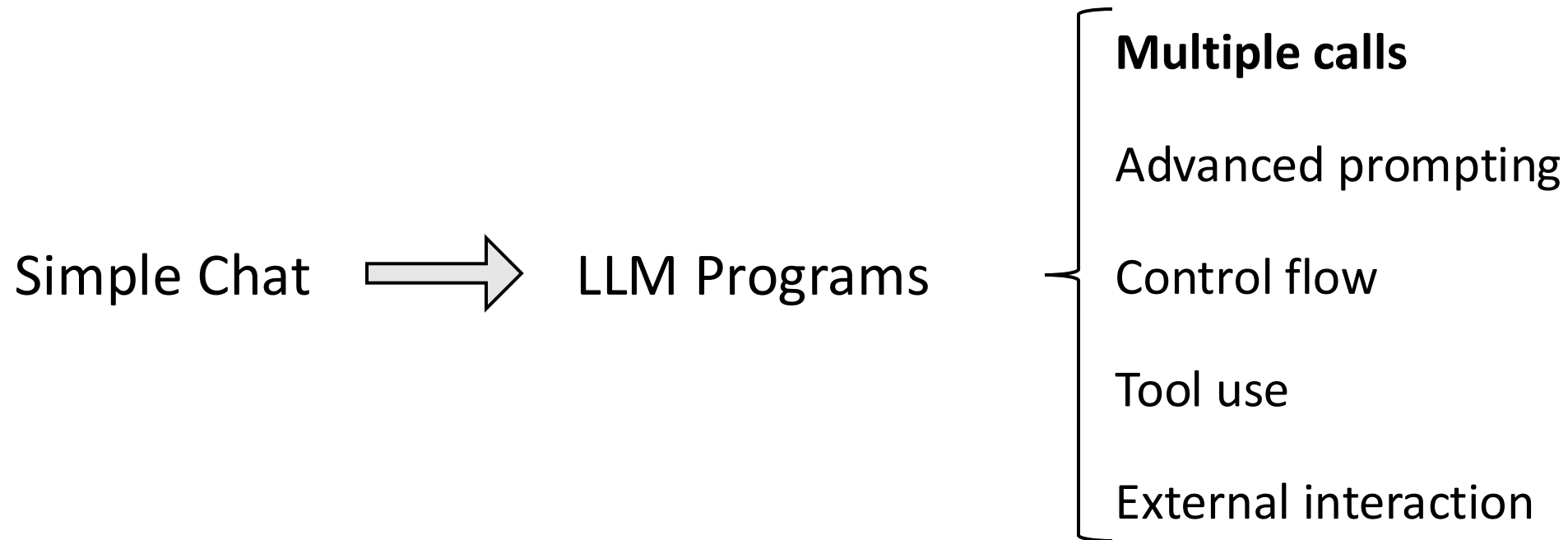
Middle Stage: innovative features and optimizations

Production Stage: research and industry use-cases

# Early Stage: the "Programming LLM" Paradigm

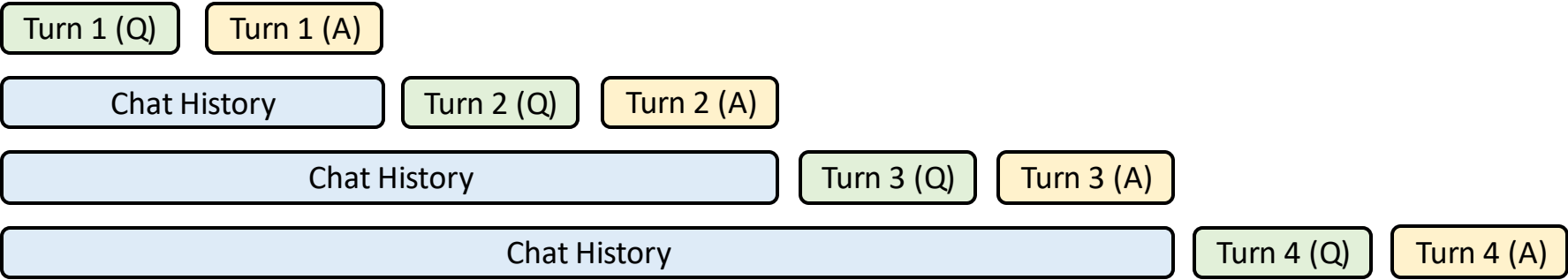From chat and simple prompting to programmatic usage of LLMs

Simple Chat ⟹ LLM Programs

**Multiple calls**

Advanced prompting

Control flow

Tool use

External interaction

# Existing Systems

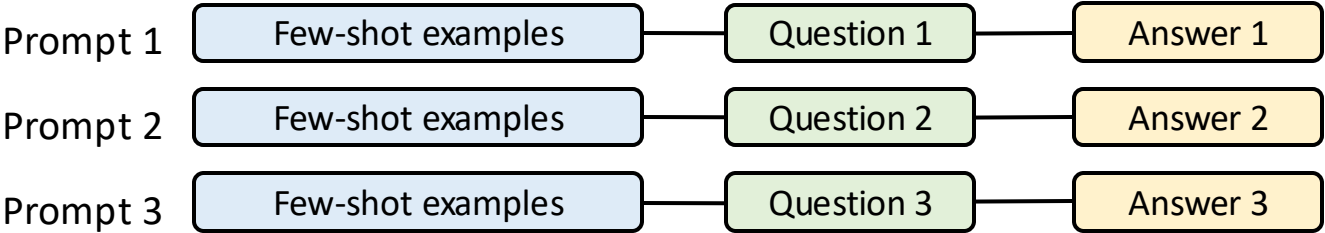**F**ront end language: ignored runtime optimizations
(Guidance, LMQL)

**B**ackend Inference engine: do not know program structure
(NVIDIA TensorRT-LLM, vLLM)
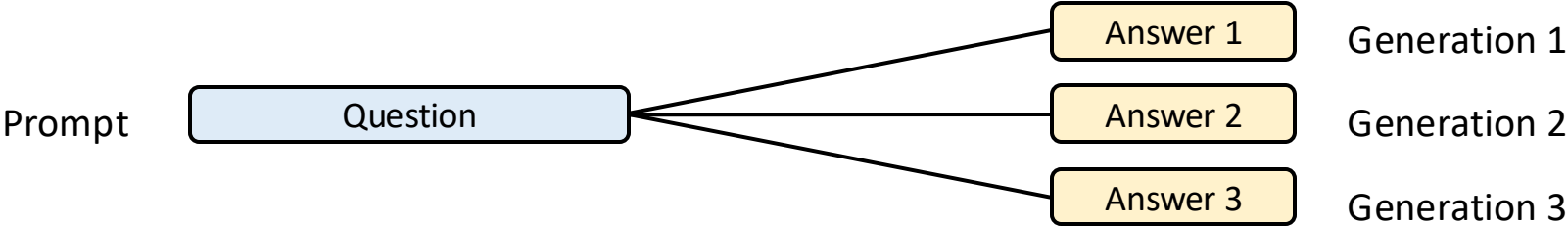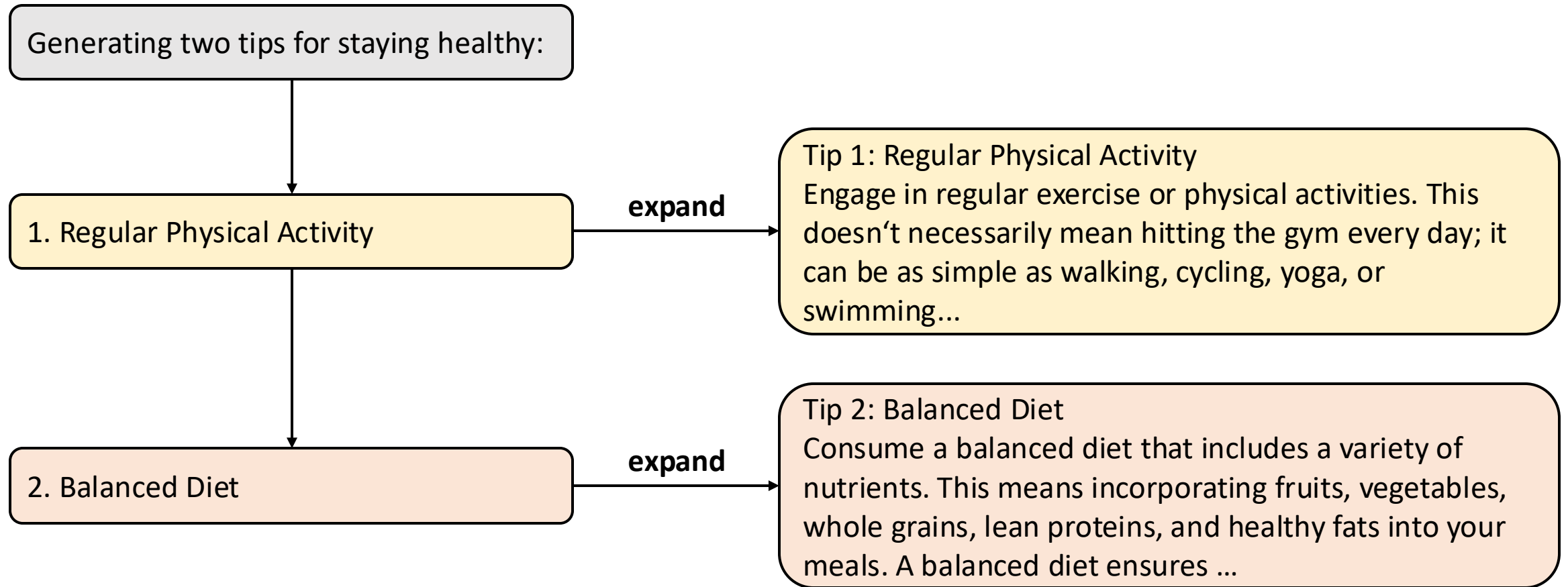
# Opportunity: KV Cache Reuse

(a) Multi-turn chat

Turn 1 (Q)  Turn 1 (A)

Chat History  Turn 2 (Q)  Turn 2 (A)

Chat History  Turn 3 (Q)  Turn 3 (A)

Chat History  Turn 4 (Q)  Turn 4 (A)

(b) Few-shot learning

Prompt 1  Few-shot examples — Question 1 — Answer 1

Prompt 2  Few-shot examples — Question 2 — Answer 2

Prompt 3  Few-shot examples — Question 3 — Answer 3

(c) Self-consistency

Prompt  Question

Answer 1  Generation 1

Answer 2  Generation 2

Answer 3  Generation 3

# Opportunity: Parallelism

Generating two tips for staying healthy:

1. Regular Physical Activity

**expand** →

Tip 1: Regular Physical Activity
Engage in regular exercise or physical activities. This doesn't necessarily mean hitting the gym every day; it can be as simple as walking, cycling, yoga, or swimming...

2. Balanced Diet

**expand** →

Tip 2: Balanced Diet
Consume a balanced diet that includes a variety of nutrients. This means incorporating fruits, vegetables, whole grains, lean proteins, and healthy fats into your meals. A balanced diet ensures ...

# System Challenges

- How to program these LLM applications?

- How to optimize across multiple LLM calls?

# Introducing SGLang: A <u>S</u>tructured <u>G</u>eneration <u>L</u>anguage

A "co-design" approach

**Front end**
- A new domain specific language embedded in Python
- Automatic parallelization and other compiler optimizations

**Back end**
- Automatic KV cache reuse with **RadixAttention**

# API example: A Multi-Dimensional Essay Judge

```python
dimensions = ["Clarity", "Originality", "Evidence"]

@function
def essay_judge(s, essay):
  s += "Please evaluate the following essay. " + essay

  # Evaluate an essay from multiple dimensions in parallel
  forks = s.fork(len(dimensions))                    # <------ Launch parallel prompts
  for f, dim in zip(forks, dimensions):
    f += (
      "Evaluate based on the following metric: " +
      dim + ". End your judgement with the word 'END'")
    f += "Judgment: " + f.gen("judgment", stop="END")  # <--- Non-blocking generation call

  # Merge judgments
  for f, dim in zip(forks, dimensions):
    s += dim + ": " + f["judgment"]                   # <------ Fetching generation results

  # Generate a summary and give a score
  s += "In summary," + s.gen("summary")
  s += "I give the essay a letter grade of " +
  s += s.gen("grade", choices=["A", "B", "C", "D"])   # <------ Constrained generation

ret = essay_judge.run(essay="A long essay ...")        # <------ Run the function
print(ret["grade"])
```

**Frontend**

- Launch parallel prompts
- Non-blocking generation call
- Fetching generation results
- Constrained generation
- Run the function

# Compiler Optimizations

- **Building a dataflow graph**
  - Remove Python Interpreter Overhead
  - Global scheduling optimization over the graph

- **Prefetching cached prefixes**
  - Insert prefetching nodes into the graph

- **Code movement for increasing sharable prefix length**
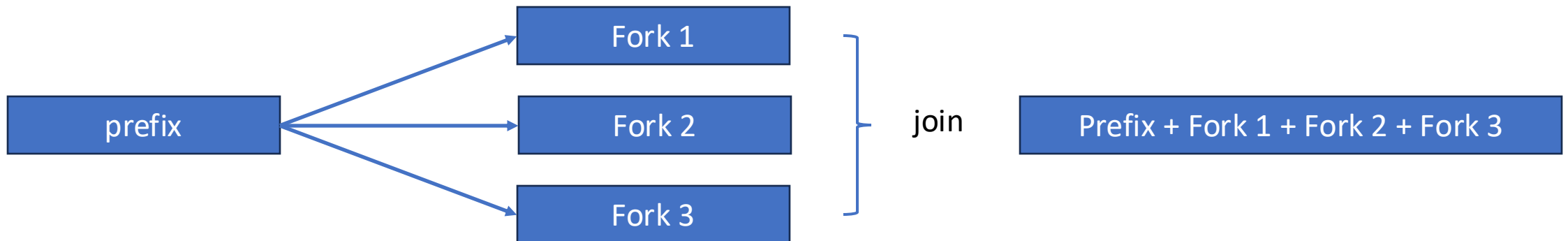  - Reorder some prompt elements with the help of GPT-4

Frontend

Backend

# Prefix caching from request tracking?

- In multi-turn chat, retrieval tasks, etc
  - The interpreter tracks the request id (rid) and caches the history before it ends.
  - Only needs to match the rid.
  - "pin" is a primitive of fixing a prefix to be cached.
  - "fork/join" primitives

Frontend

Backend



| prefix | → | Fork 1 |
| | | Fork 2 | join | Prefix + Fork 1 + Fork 2 + Fork 3 |
| | | Fork 3 |

Cannot reuse shared prefix across requests!

# Runtime (SRT) with RadixAttention

**Existing Systems**: Discard KV cache after a request finishes.

**Ours**: Maintain an LRU cache of the KV cache of all requests in a radix tree (compact prefix tree).

# Runtime (SRT) with RadixAttention

Maintain an LRU cache of the KV cache of all requests in a radix tree.

# Runtime (SRT) with RadixAttention

Maintain an LRU cache of the KV cache of all requests in a radix tree.

# Cache Aware Scheduling

- In the request queue, sort the requests according to the matched prefix length
  - Achieves good cache hit rate

- Future work
  - Distributed cache aware scheduling for multiple data parallel workers
  - Fairness to prevent starvation (https://arxiv.org/abs/2401.00588)

Aug -
Dec.
2023

Early Stage: the "programming LLM" paradigm

Jan. -
now
2024

Middle Stage: innovative features and optimizations

RadixAttention    Upper-level Scheduling

now -
2024

Production Stage: research and industry use-cases

# SGLang Structure: Pipeline

# SGLang Structure: Inside TP Worker



1.Decode Batch

2.Prefill a New Batch

3.Decode Batch

EOS token, finished this req
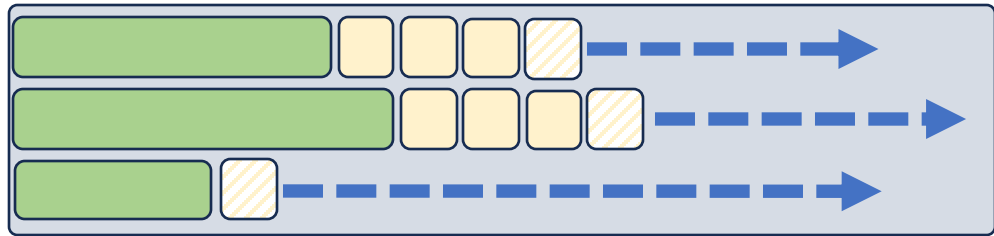
Merge into decode batch
ready for decode

newly decoded token     prompt tokens

EOS (newly decoded)     decoded token

## How to always keep the batch size large enough?

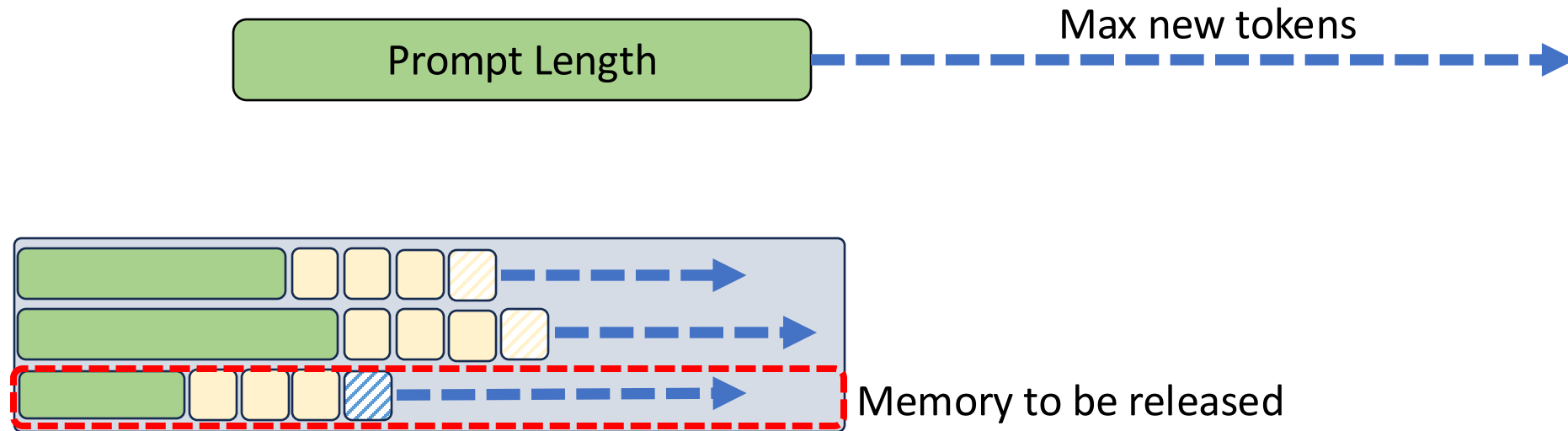# Dynamically Adjust the new token ratio estimation



The max context length
decided by max new tokens

- There is a lot of space left in the GPU memory
- We do not need to reserve every token in max new tokens

# Dynamically Adjust the new token ratio estimation

Max new tokens

Prompt Length

Memory to be released

1. The EOS would be earlier than the max new tokens.
2. There are always requests finished and release all the memory.

Only preserve $\beta \times$ max_new_token tokens in advance, and adjust $\beta$ dynamically.

Aug - Dec. 2023

Early Stage: the "programming LLM" paradigm

Jan. - now 2024

Middle Stage: innovative features and optimizations

RadixAttention    Upper-level Scheduling    Jump-forward decoding
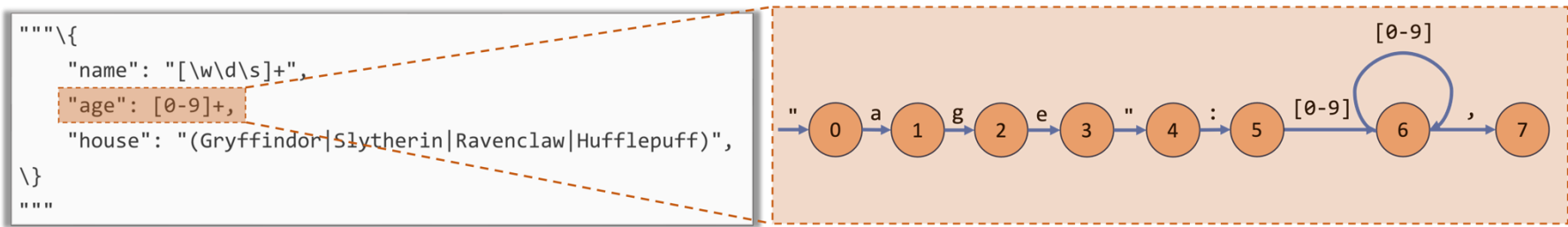
now - 2024

Production Stage: research and industry use-cases

# Jump-forward JSON Decoding

**Method**

- Analyze the regular expression
- Compress the finite state machine
- Decode multiple tokens at the same time



Regular Expression

Finite State Machine

Constrained Decoding With Logits Mask

# Speedup Regex Guided Generation



Jump-Forward Decoding With Compressed FSM

Normal Decoding With FSM

Generated JSONs

# Jump-forward JSON Decoding

**Results:**

3x faster latency

2.5x higher throughput

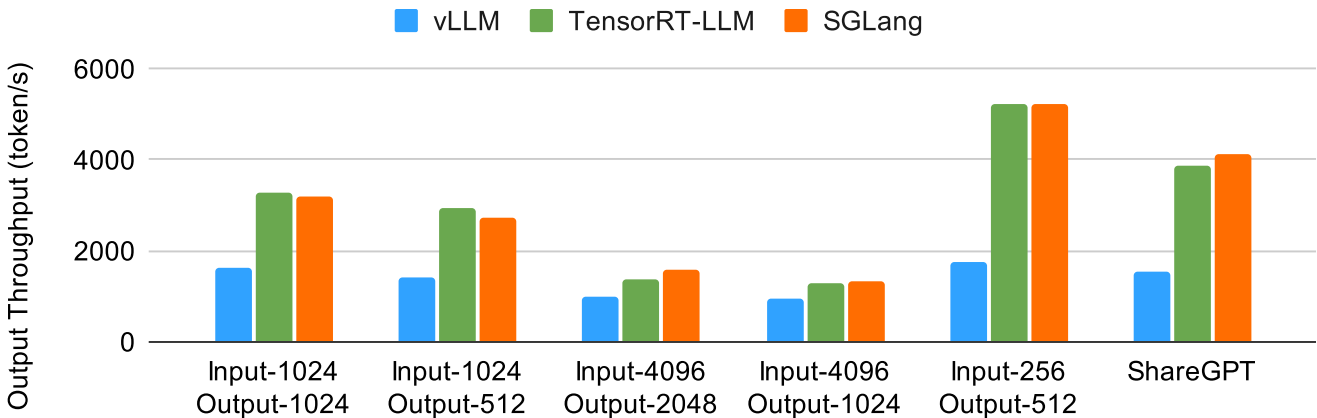| SGLang | Outlines + vLLM |
|---|---|
| | |

# Summary: techniques in SGLang

- RadixAttention
- Jump-forward JSON Decoding
- Torch Compile
- Flashinfer Kernels
- Chunked Prefill
- Continuous Batching
- Token Attention(Paged Attention with page_size = 1)
- CUDA Graph
- Interleave window attention

# SGLang v0.2 Results

Llama-8B (bf16) on 1 x A100. Higher Throughput is Better.



Llama-70B (fp8) on 8 x H100. Higher Throughput is Better.

Aug -
Dec.
2023

Early Stage: the "programming LLM" paradigm

Jan. -
now
2024

Middle Stage: innovative features and optimizations

now -
2024

Production Stage: research and industry use-cases

# Research and industry use cases

x.ai: Production serving of grok-2 and grok-2-mini on X

Databricks: accelerate research workflow by 3.3x

LMSys Chatbot Arena: serving vision language models

LLaVA OneVision: serving multi-modal image and video models

ByteDance

# Future work

- multi-level cache
- distributed radix attention for replicas
- long-context sequence parallelism
- speculative decoding
- communication overlapping
- ……

# Do the serving engines come to converge on performance?

YES and NO

Basic performance eventually converge

But there are more sophisticated workloads from different scenarios:
RAG systems, agent systems, …

We never forget about the origin of SGLang!
 Structured inputs, interactions with different resources, multi-modality, …

# Principles in future development

Simplism          Minimalism          Modularity          Ease of use

Development velocity          Performance

# Future work

- <span style="color:red">multi-level cache</span>

- distributed radix attention for replicas

- long-context sequence parallelism

- speculative decoding

- communication overlapping

- ……

# Context KV caching is huge for serving!

"Our system achieves a **95% cache rate**, further reducing inference cost."
-- from Character AI

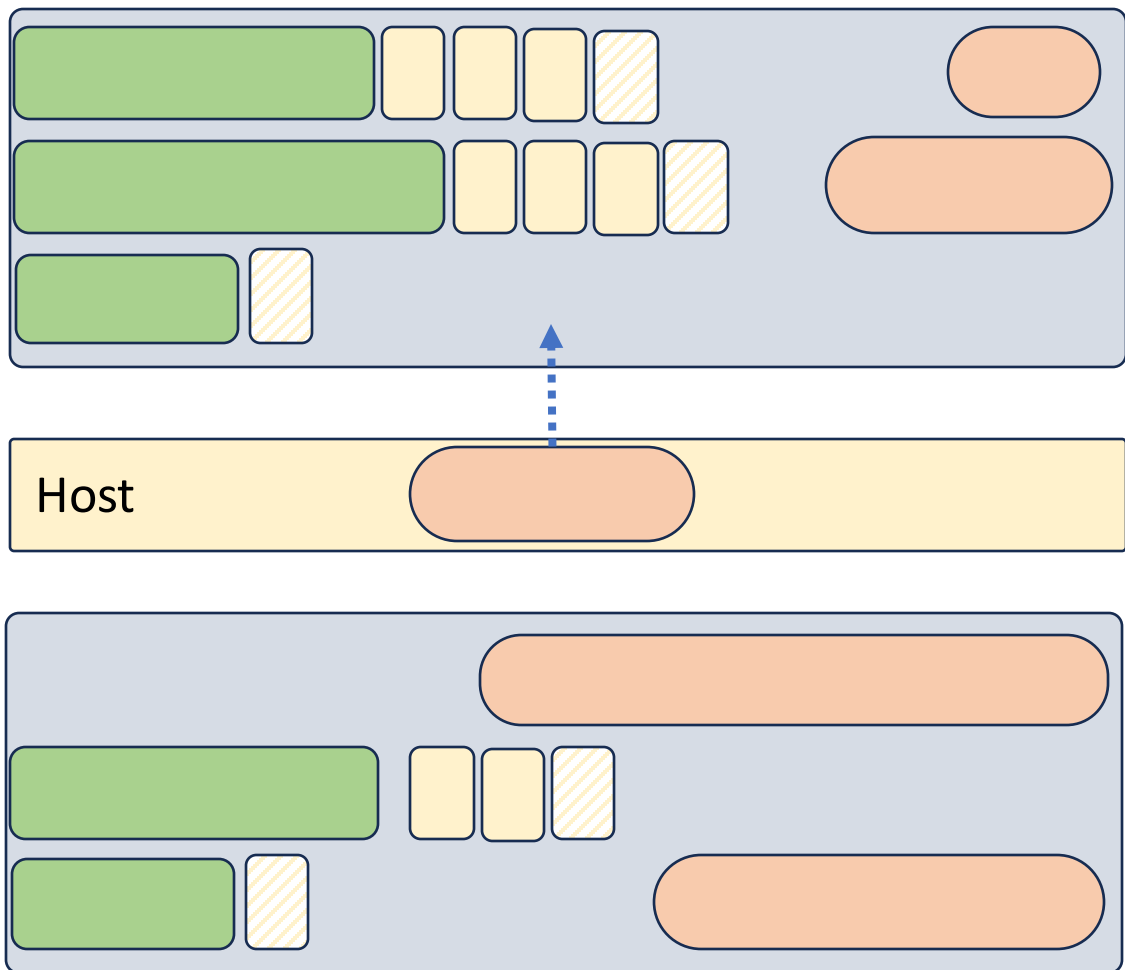"Even without any optimization, historical data shows that users **save over 50% on average**."
-- from DeepSeek

"We observed a **52.4% cache hit** rate for LLaVA-Next-34B and **74.1%** for Vicuna-33B."
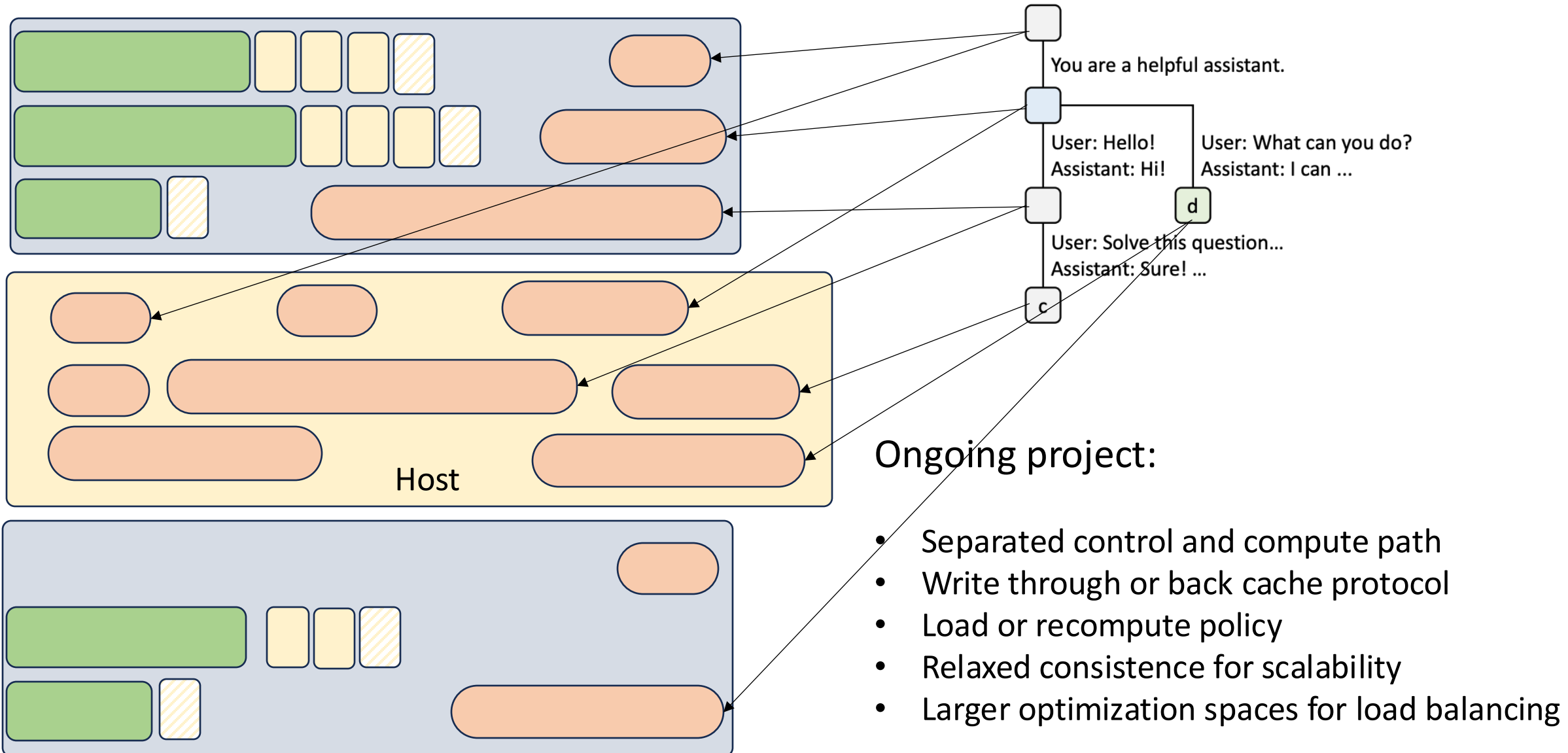--  from Chatbot Arena

…

# Challenges for caching



Historical KV caches compete GPU memory with ongoing requests

Host

Offloading and reloading requires more user effort and potentially increase latency

Trade-off between load balancing and cache hit rate for multi-GPU scenarios

# Towards a transparent multi-level caching



You are a helpful assistant.

User: Hello!
Assistant: Hi!

User: What can you do?
Assistant: I can …

d

User: Solve this question…
Assistant: Sure! …

c

Host

Ongoing project:

- Separated control and compute path
- Write through or back cache protocol
- Load or recompute policy
- Relaxed consistence for scalability
- Larger optimization spaces for load balancing

# Design principle: separated control and compute path