

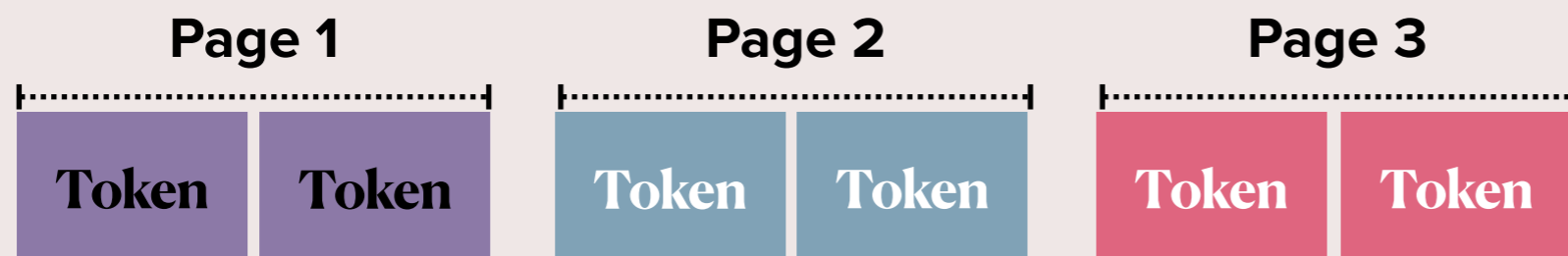
Span Queries

Nick Mitchell, Paul Castro, Nathan Ordonez, Thomas Parnell, Mudhakar Srivasta, Antoni Viros i Martin

IBM Research | May 2026

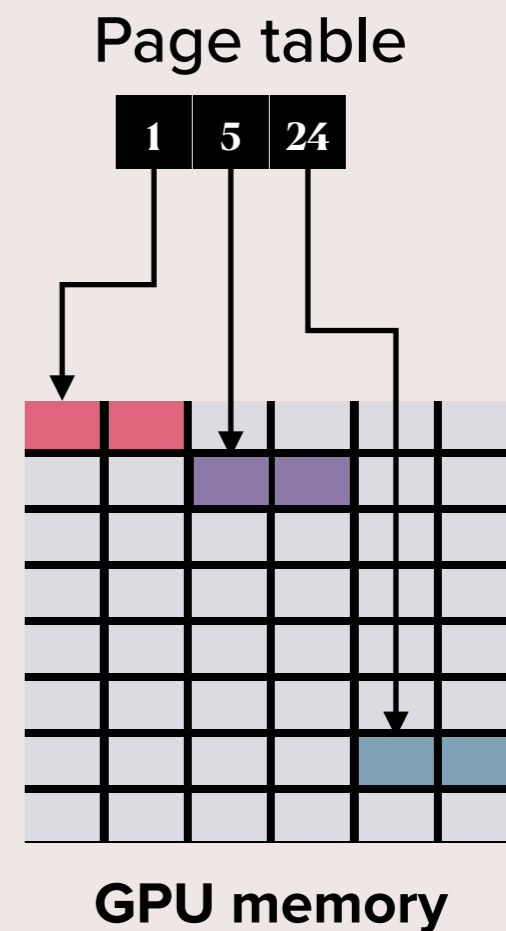
Pages (a.k.a. Blocks)

Contiguous sequence of tokens in a paged attention server



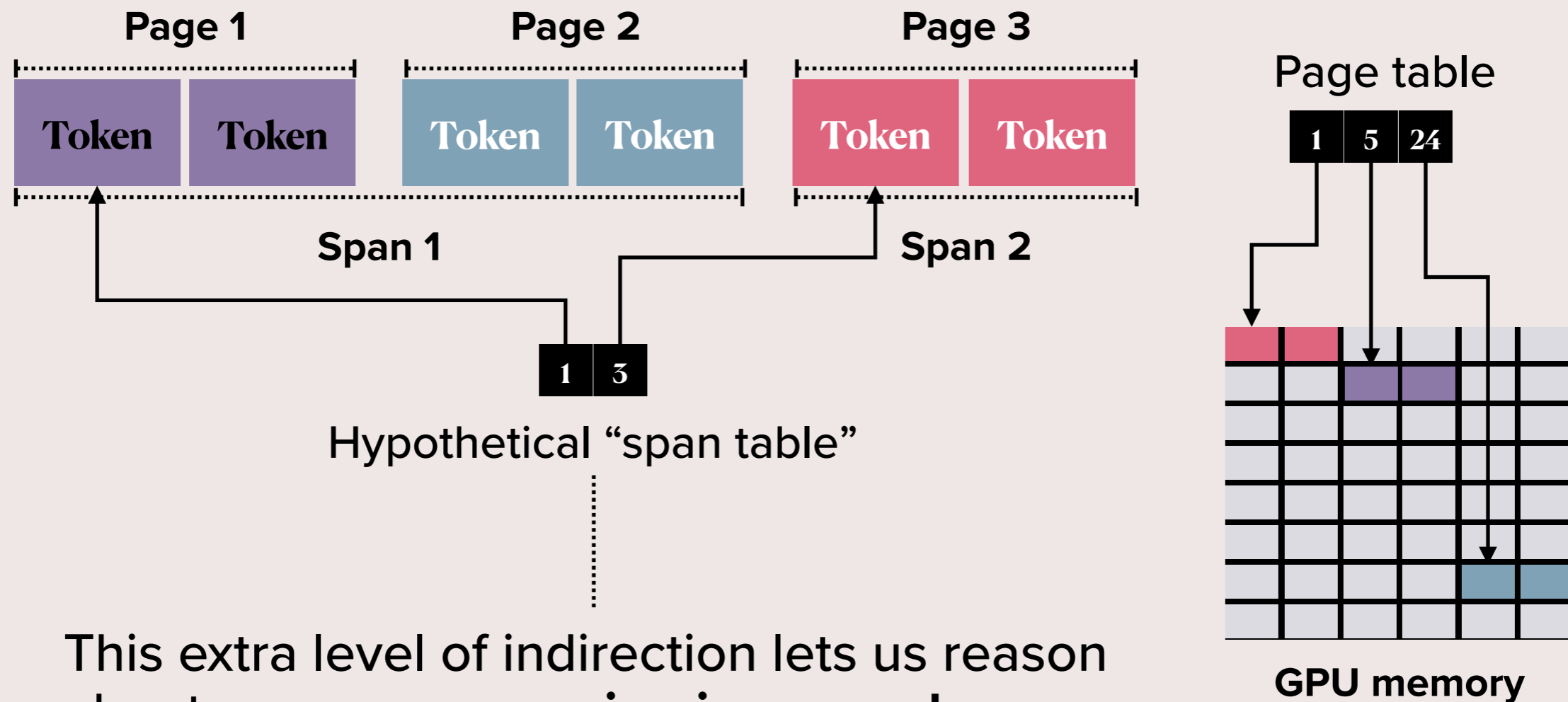
vLLM virtualizes **discontiguous physical** pages into **contiguous virtual** pages

So that they can be laid out in any physical order



Spans

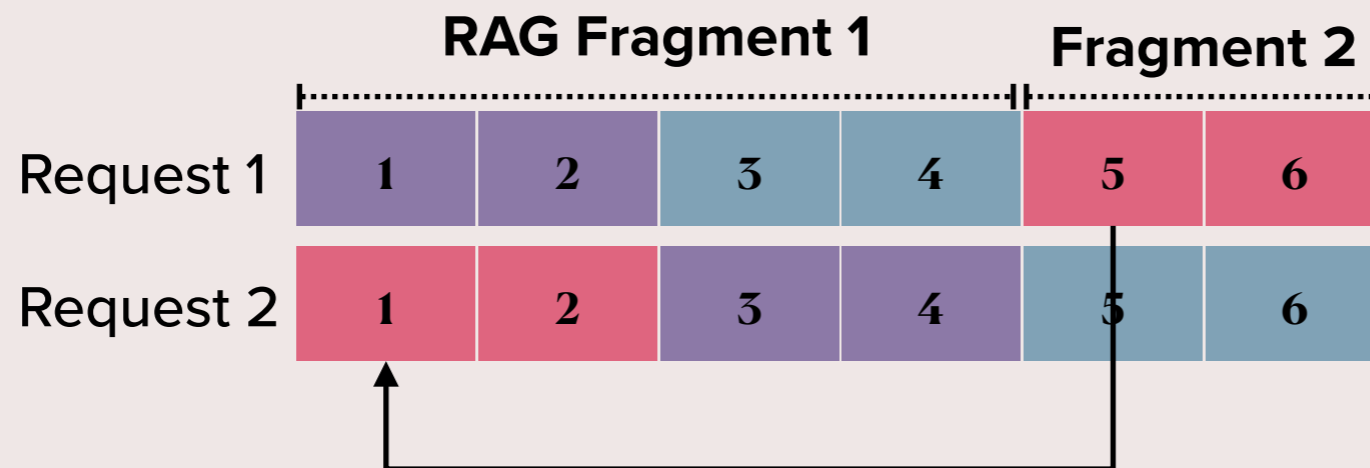
Contiguous sequence of pages in a paged attention server



This extra level of indirection lets us reason about **spans as occurring in any order**.

Commutativity Hypothesis

If order doesn't matter, we can optimize



$AB=BA \Rightarrow$ **Commutativity**

Q1: What does commutativity buy?



Cache Locality

Because we can **reuse pages** even if they were originally used in different context



"Attention Locality"

Because we can apply **divide and conquer** to lost-in-the-middle problems.

Q2: How do clients express a “span table”?

vLLM

500 lines changed

1. RoPE on read, not on write
2. Prefix scans must selectively disable block hash chaining for “spanned” regions

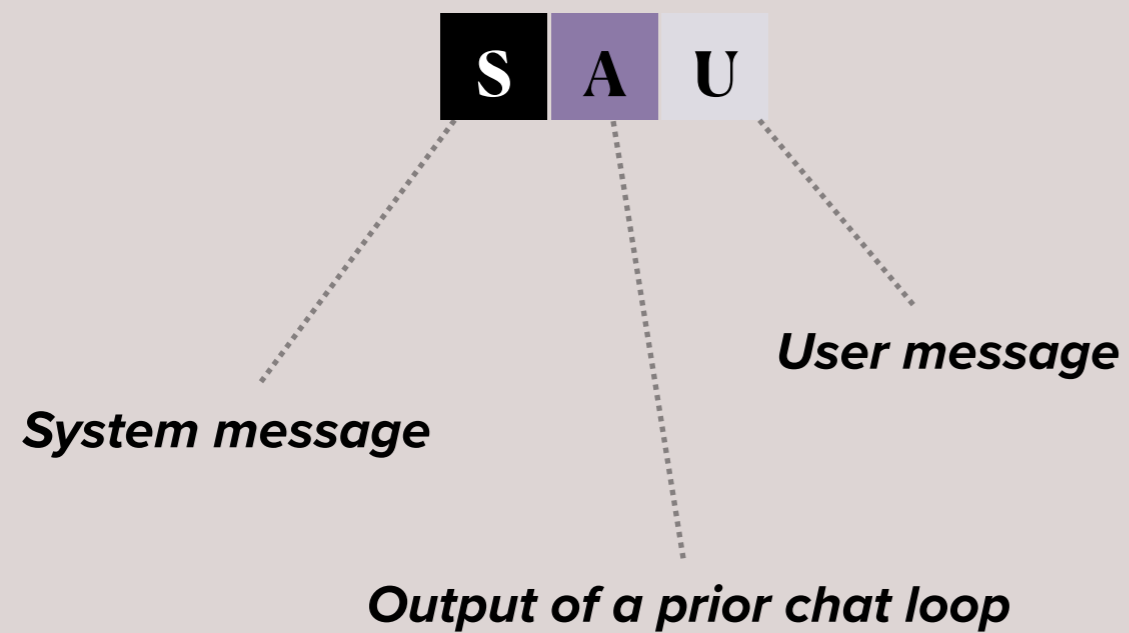
spn1

Clients issue span queries, vLLM parses, plans, optimizes

The Gist

Generalize token sequences

Chat Completion

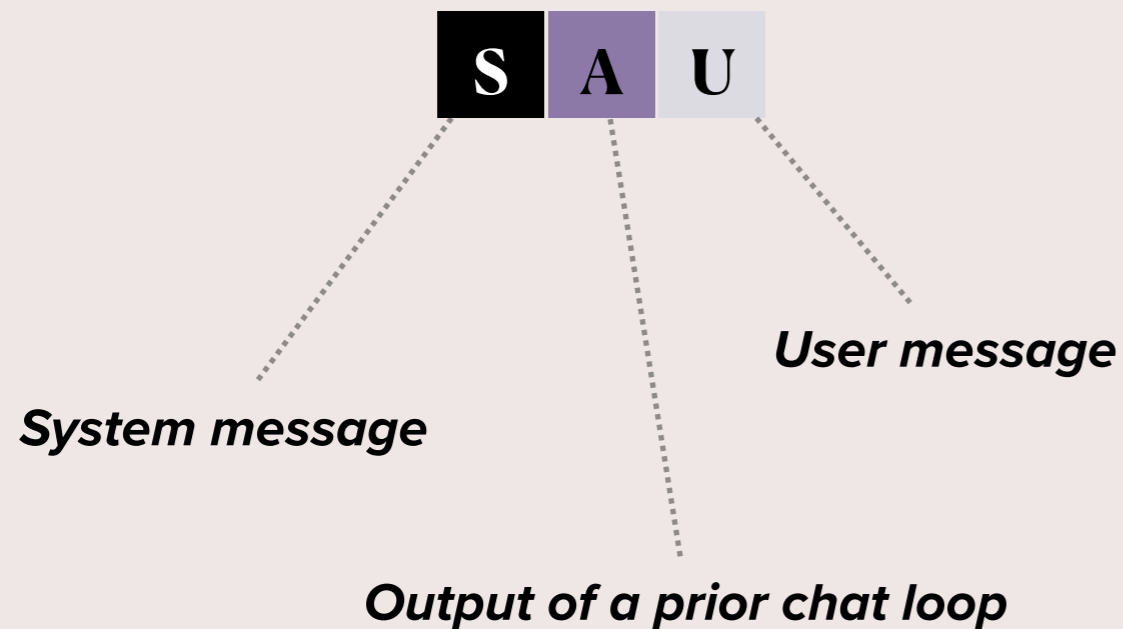


API of Message Lists

The Gist

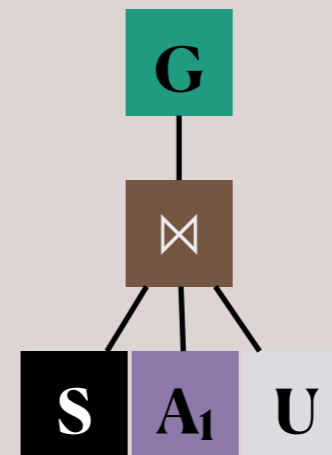
Generalize token sequences

Chat Completion



API of Message Lists

Span Query

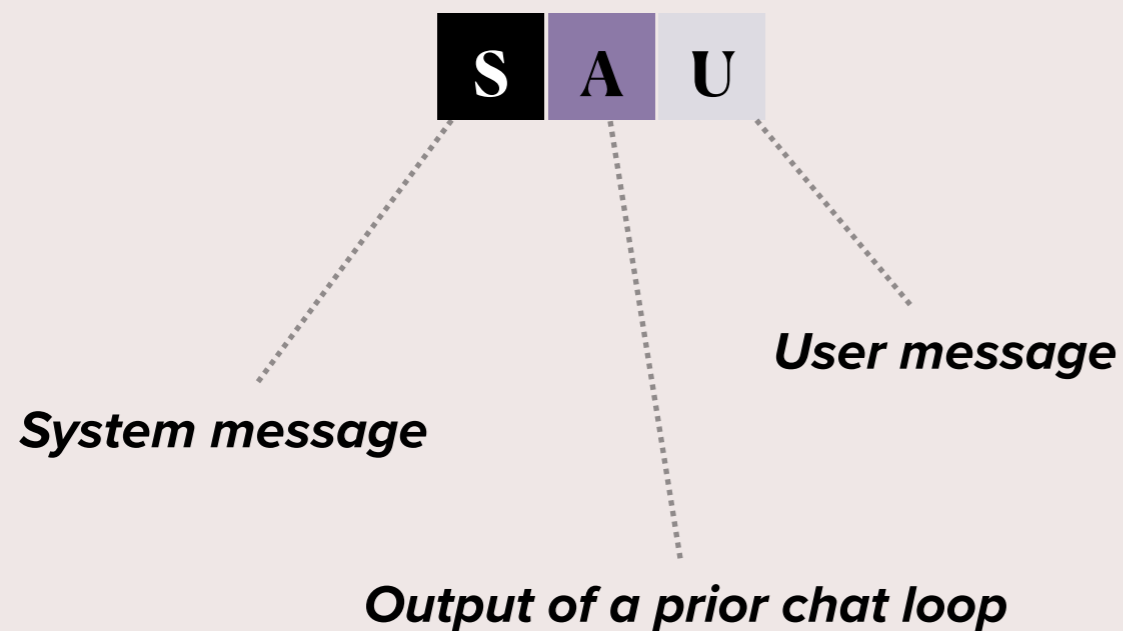


API of Message Trees

The Gist

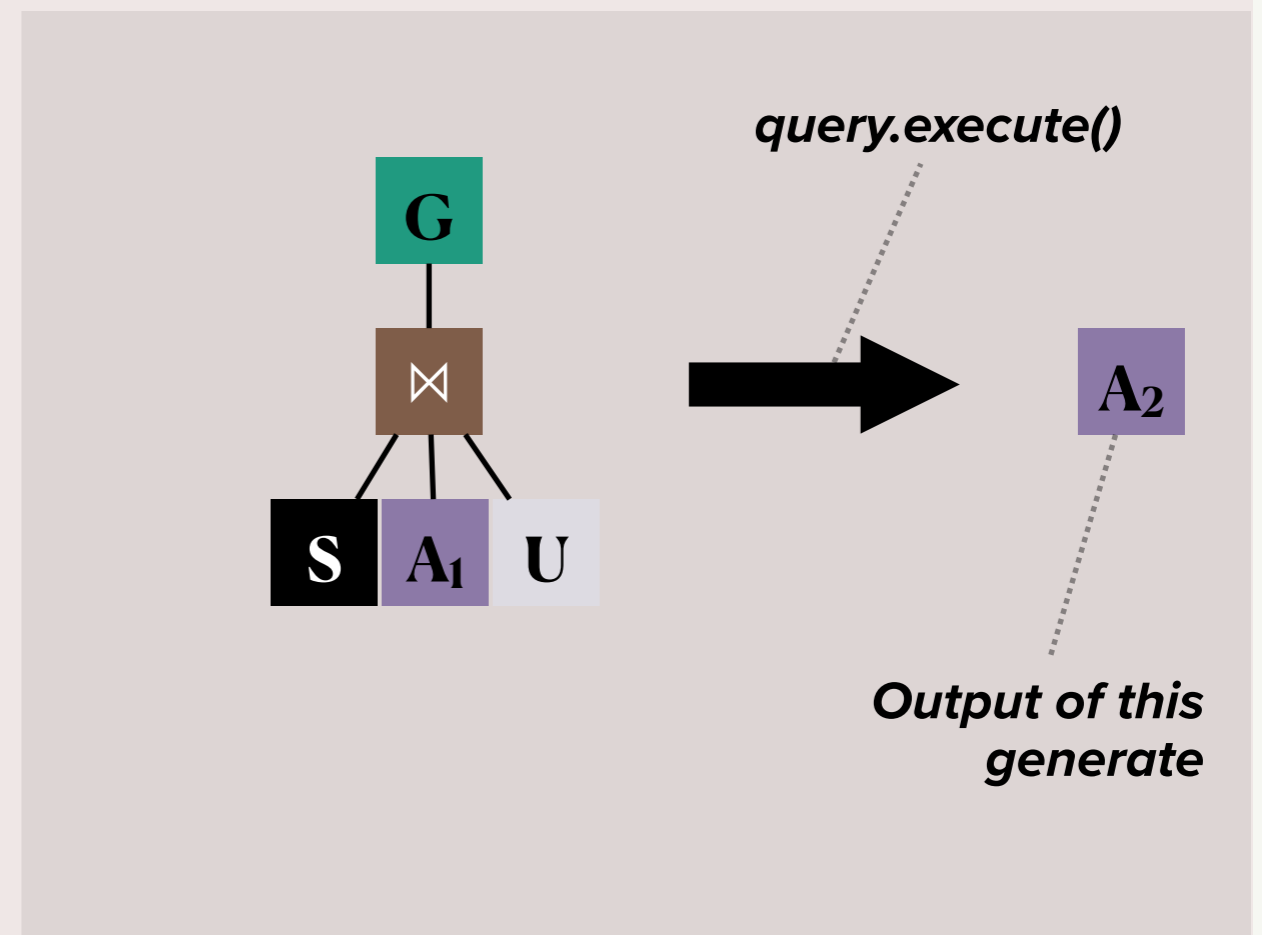
Generalize token sequences

Chat Completion



API of Message Lists

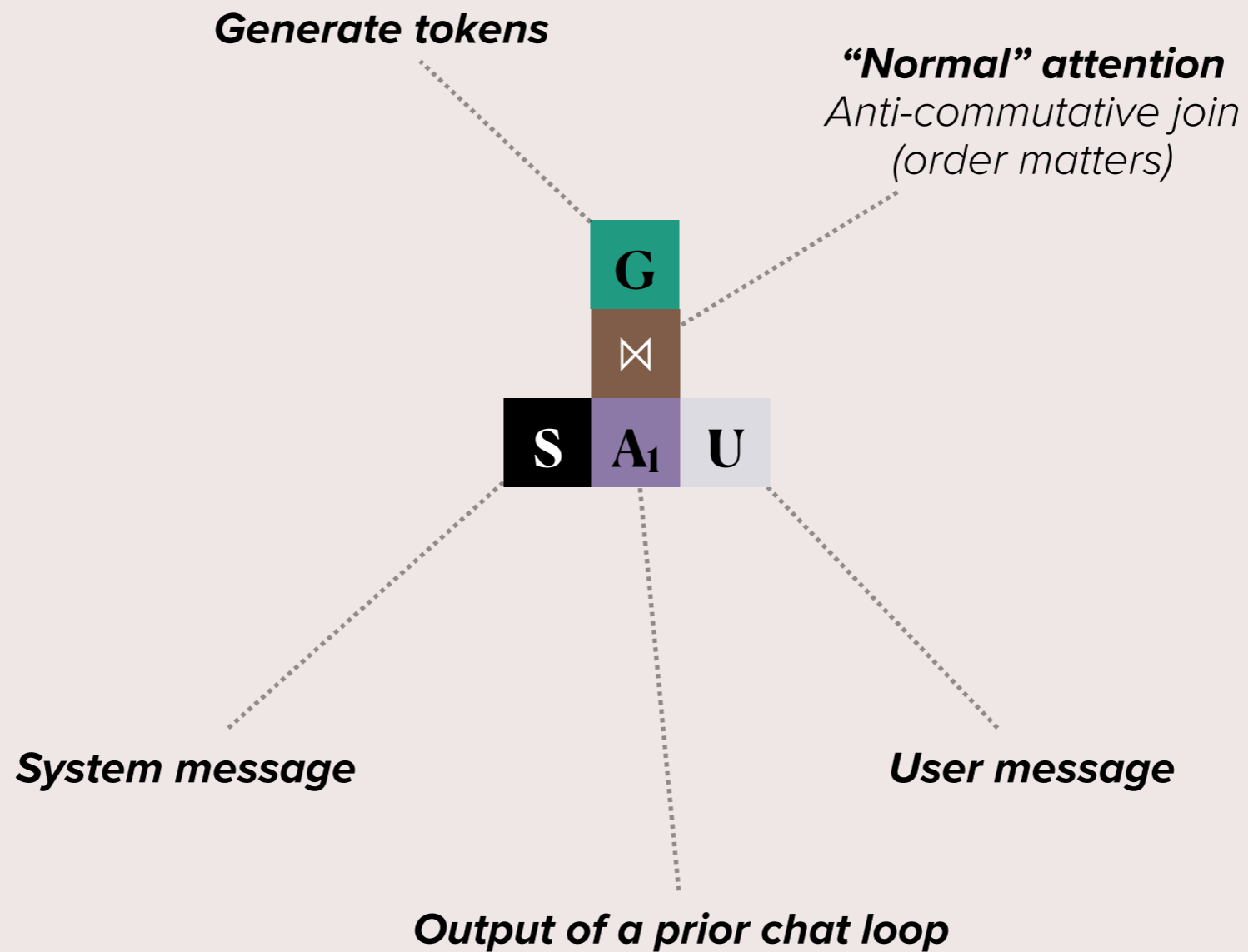
Span Query



API of Message Trees

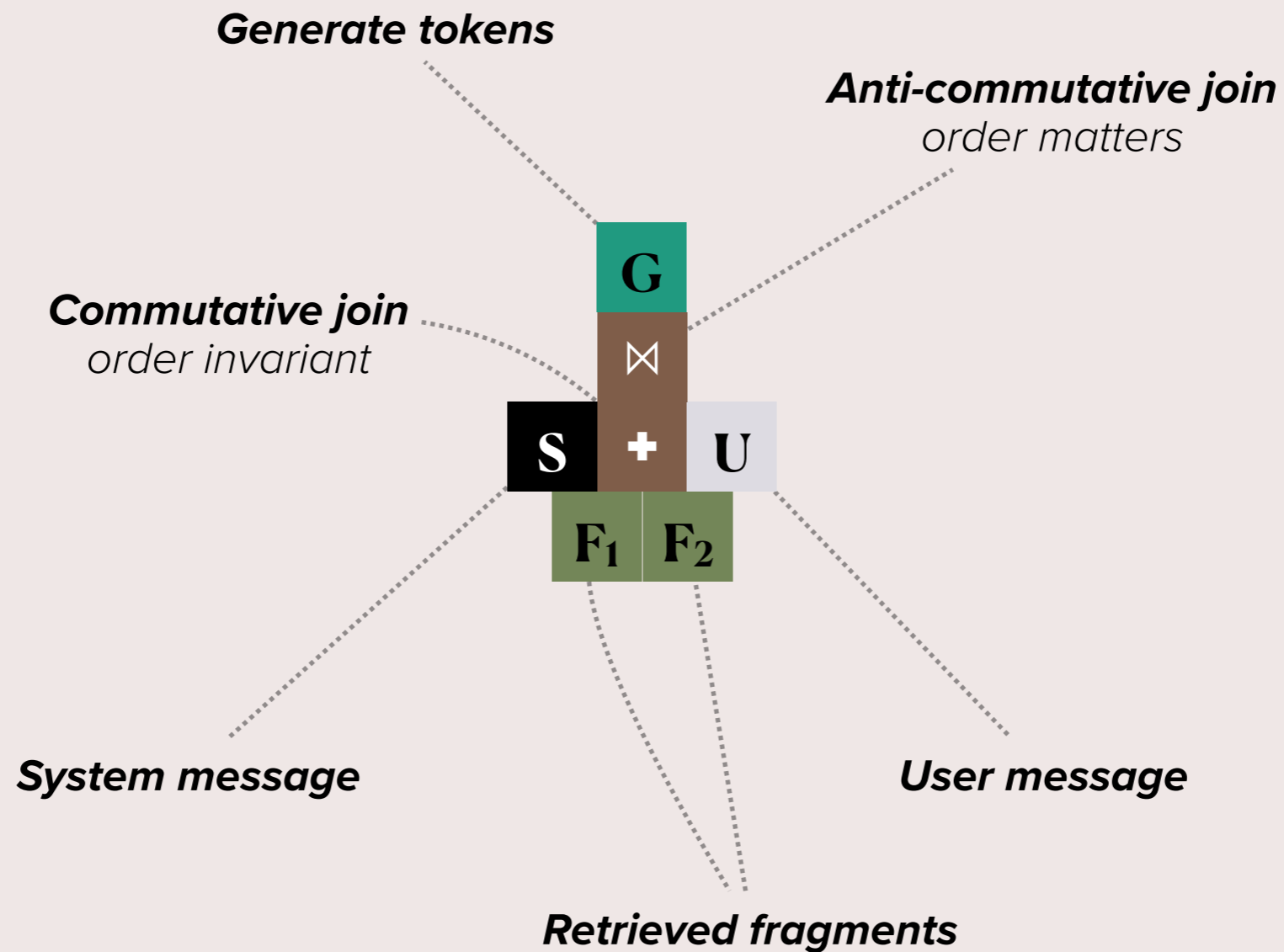
Span Query

Chat use case



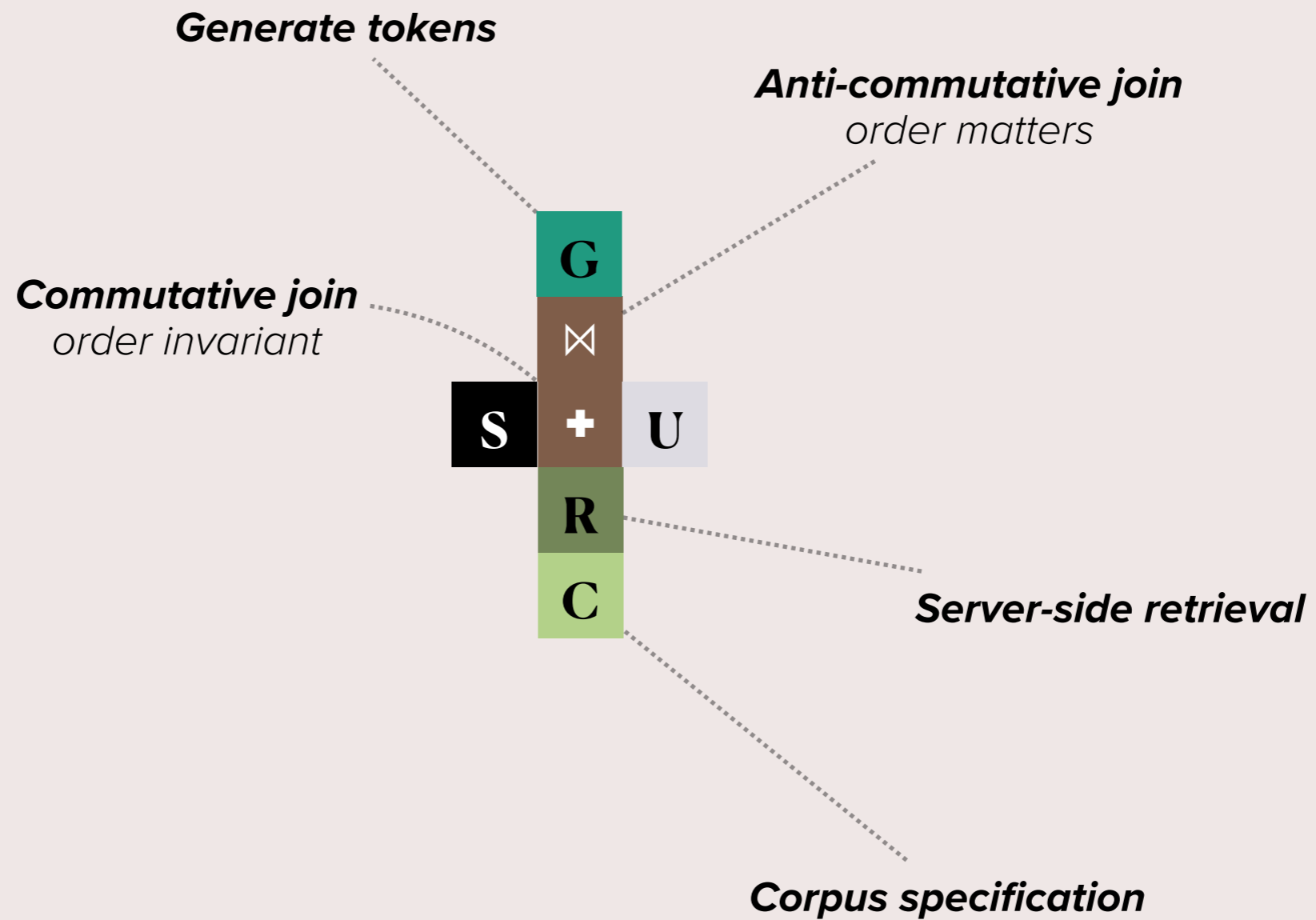
Span Query

RAG use case



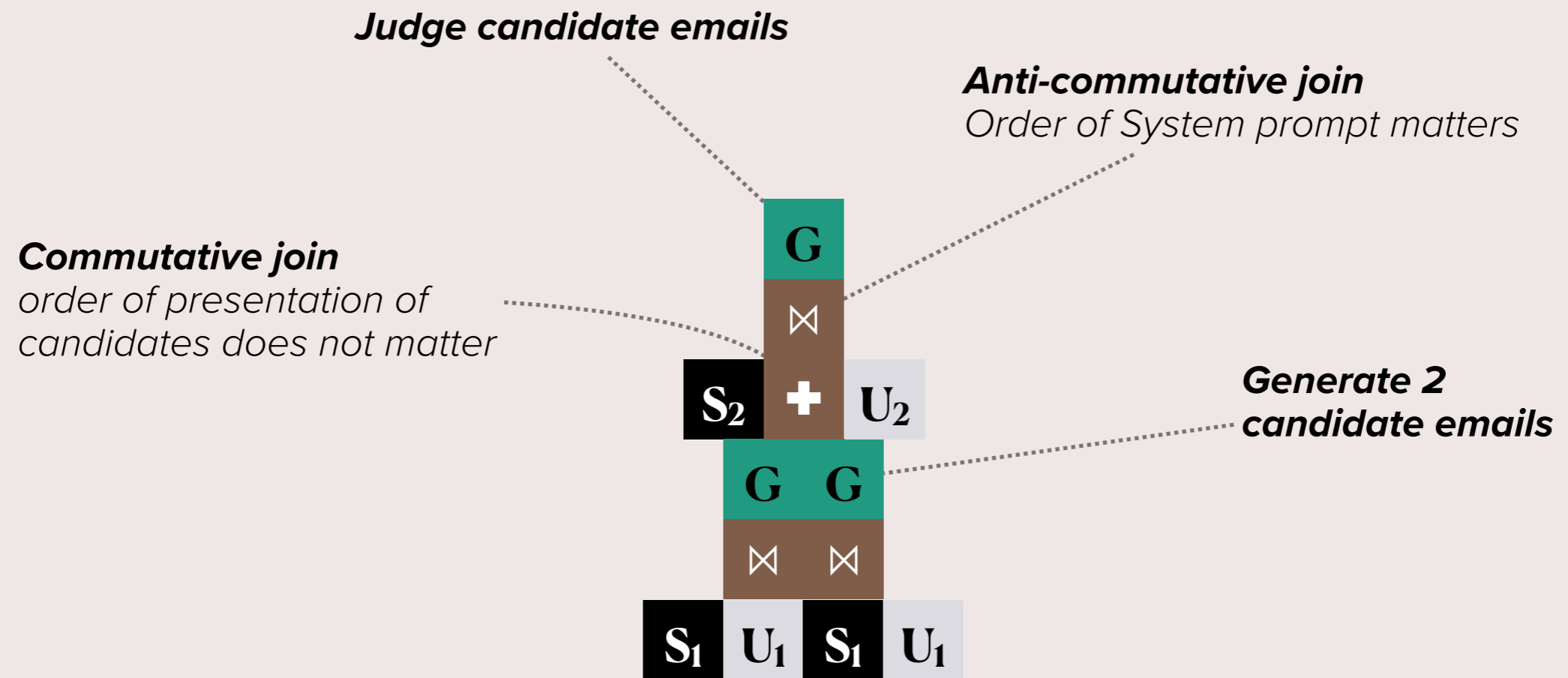
Span Query

RAG use case with retrieval



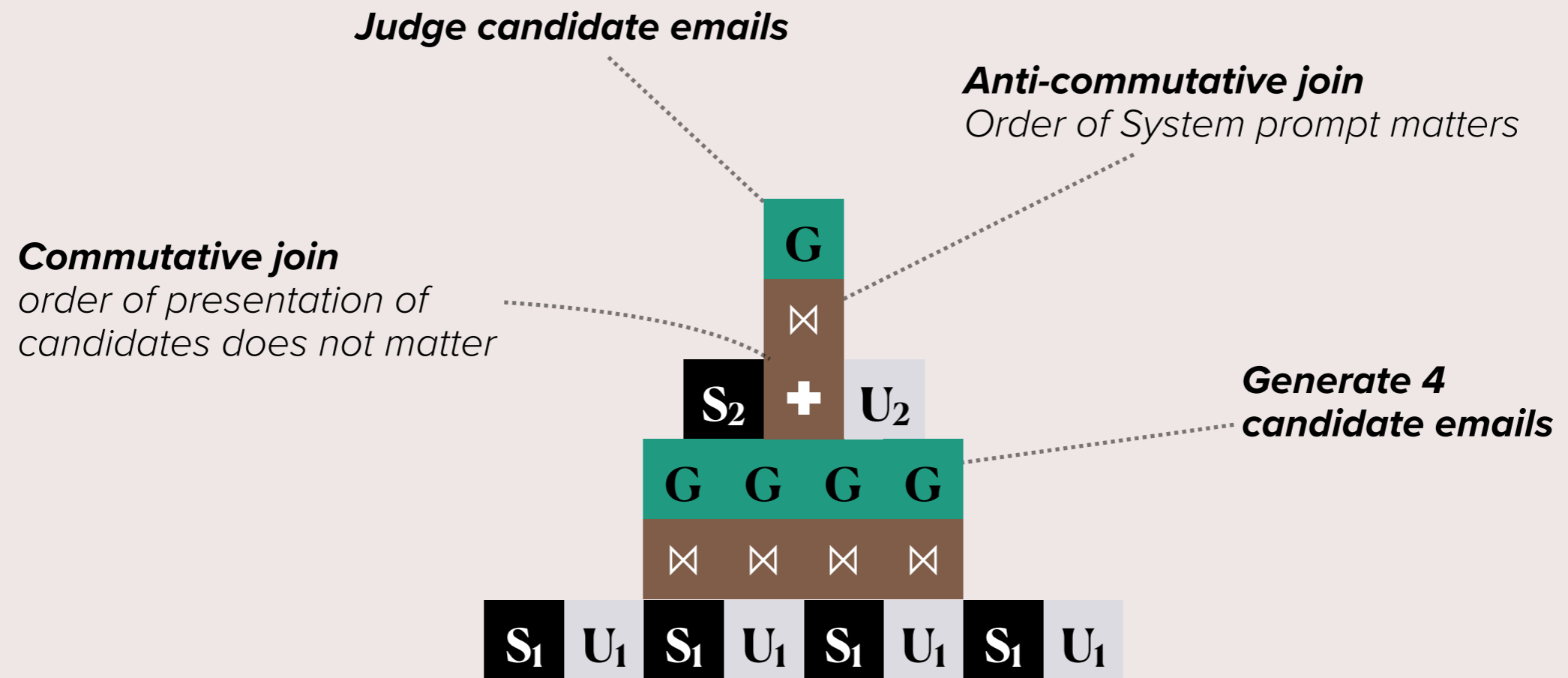
Span Query

2-way Judge-Generator use case



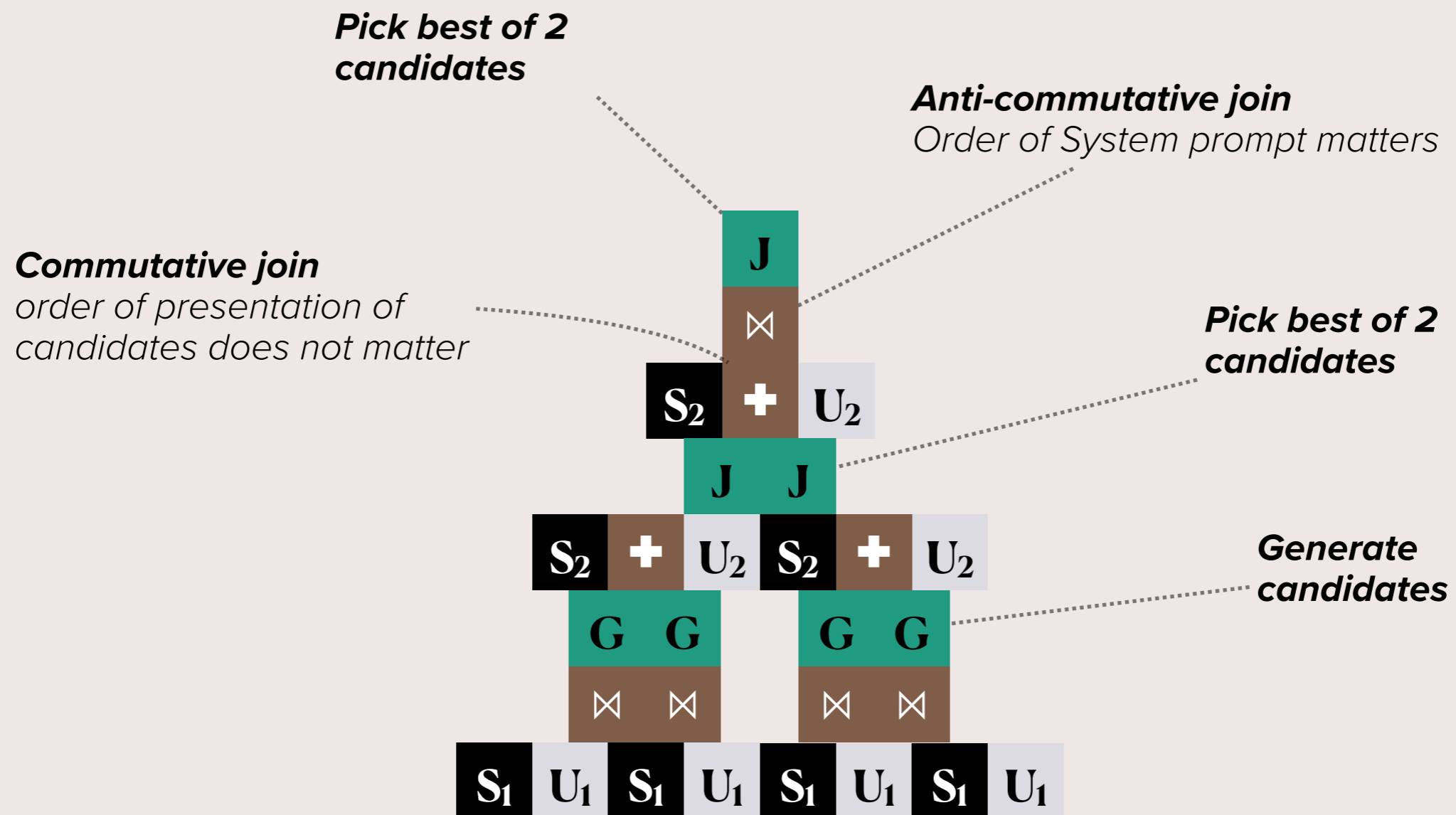
Span Query

4-way Judge-Generator use case



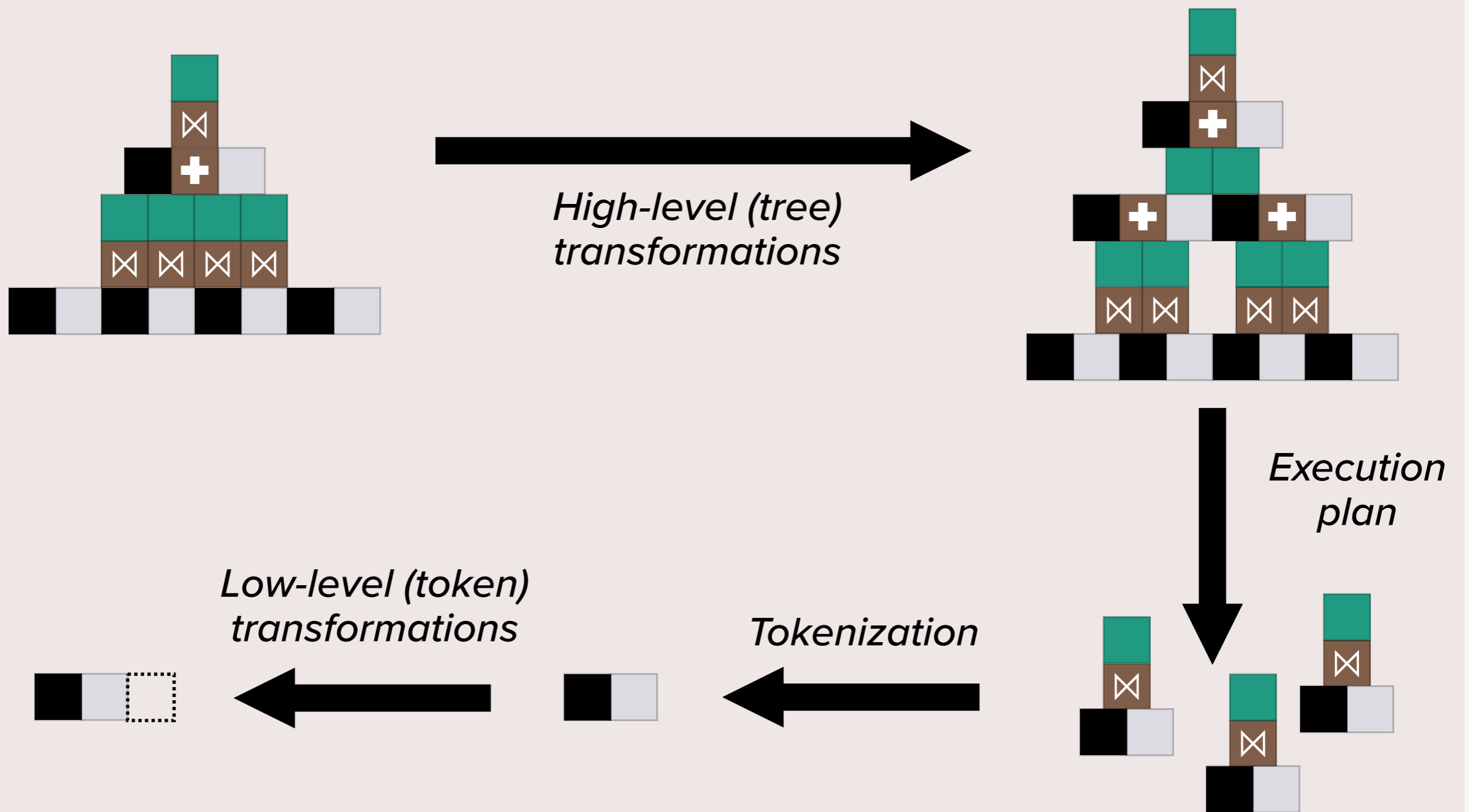
Span Query

4-way Judge-Generator query, optimized to avoid LITM



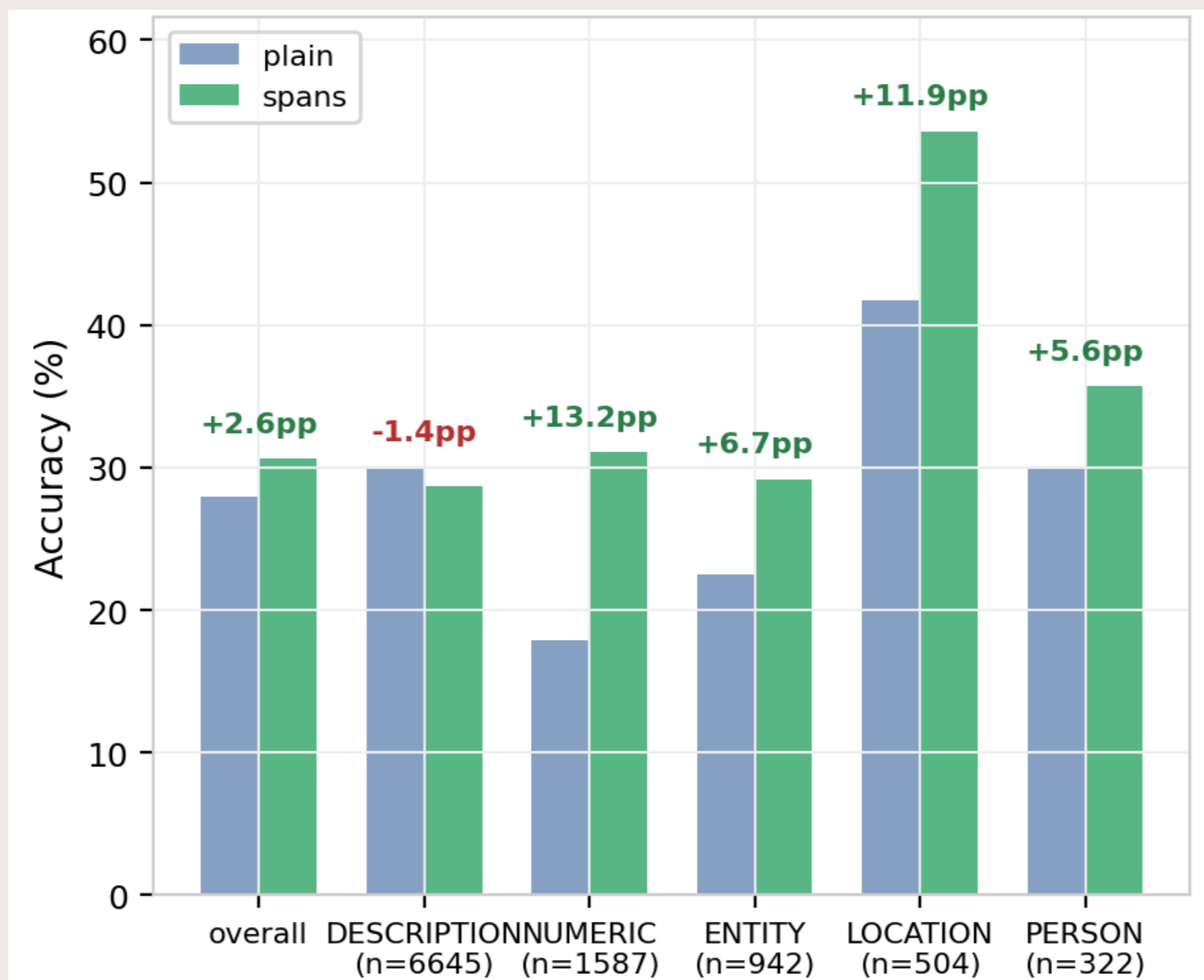
Span Query Operation

Optimization and Execution



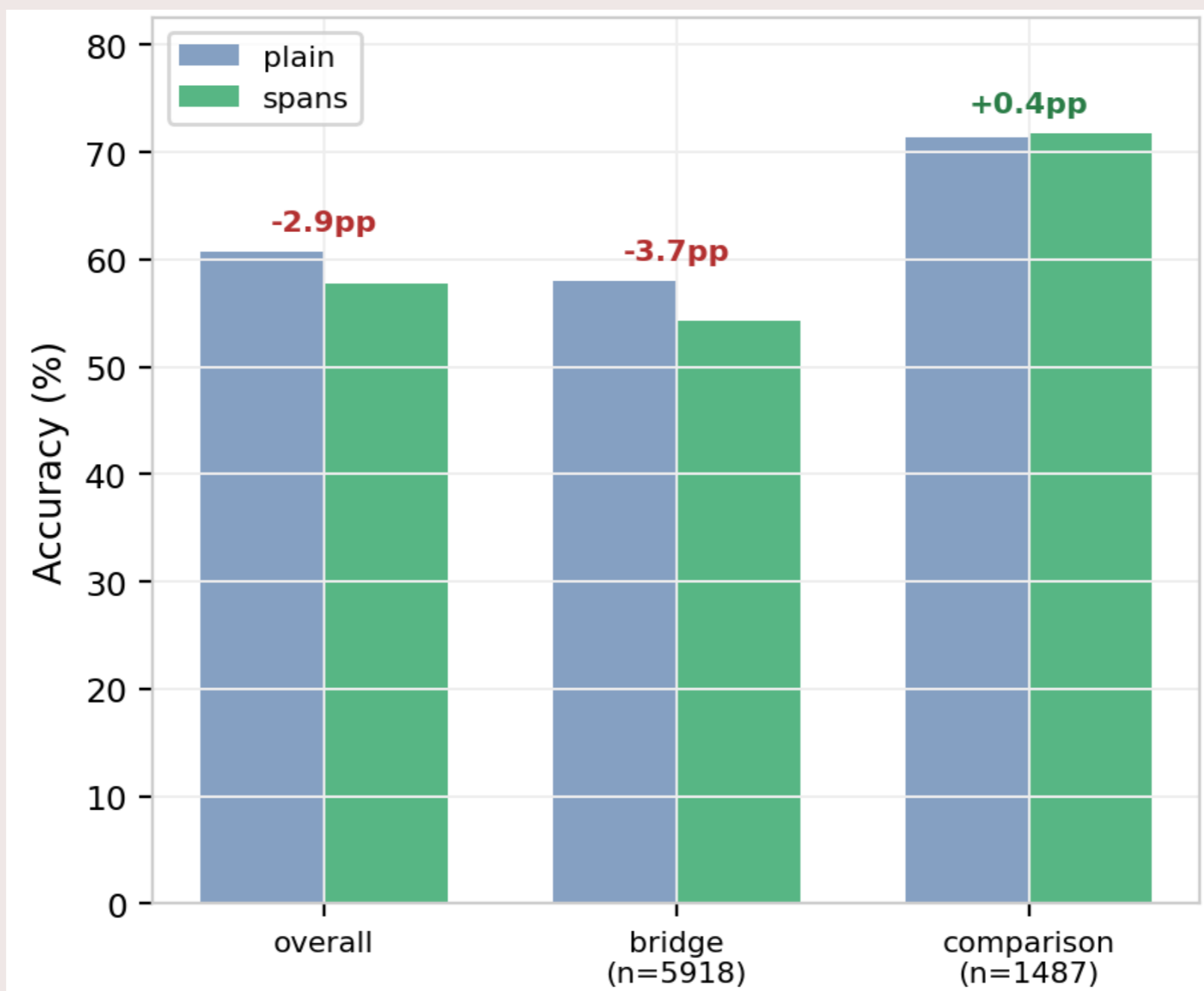
Accuracy

With span table implementation (MSMARCO)



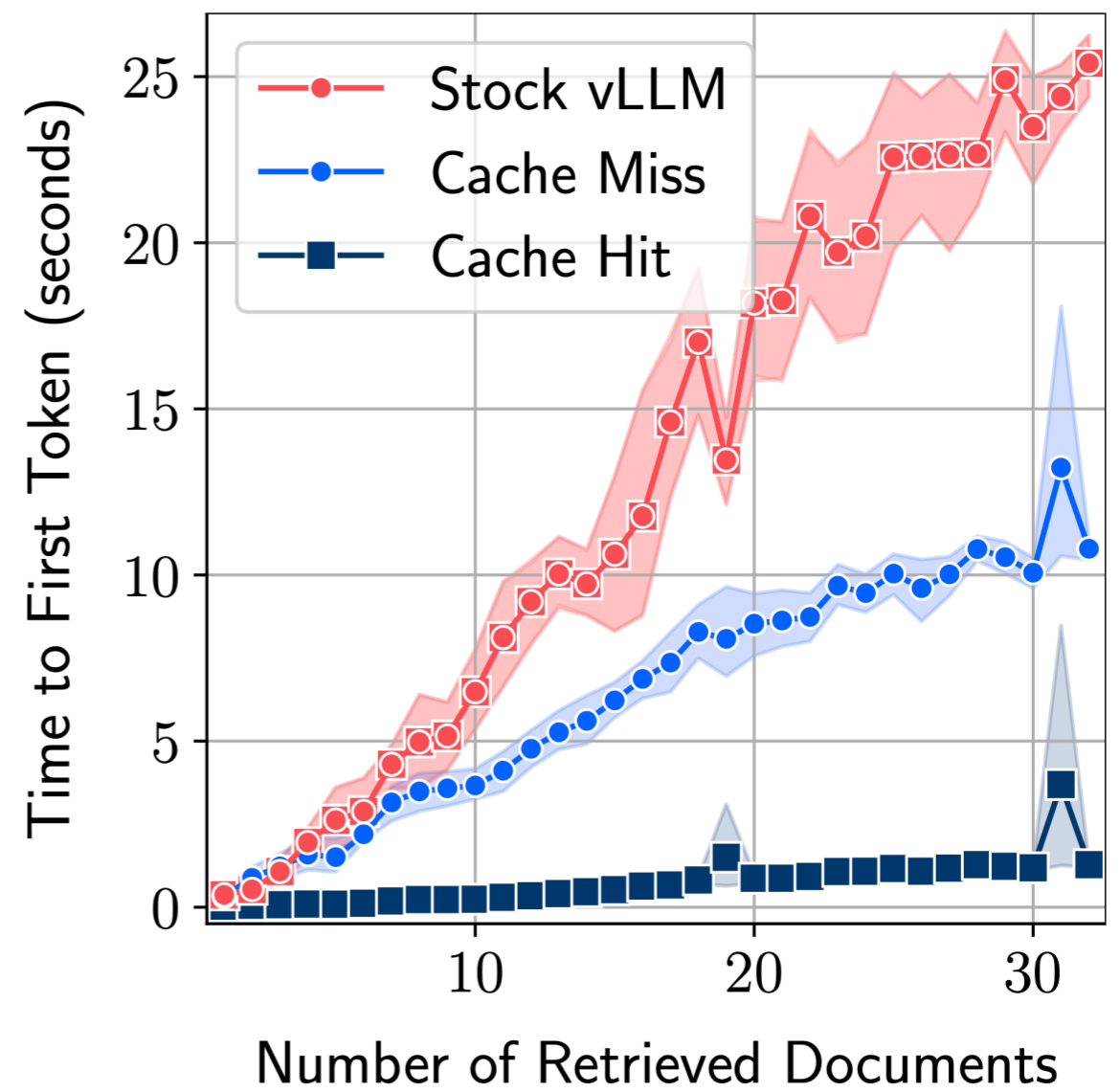
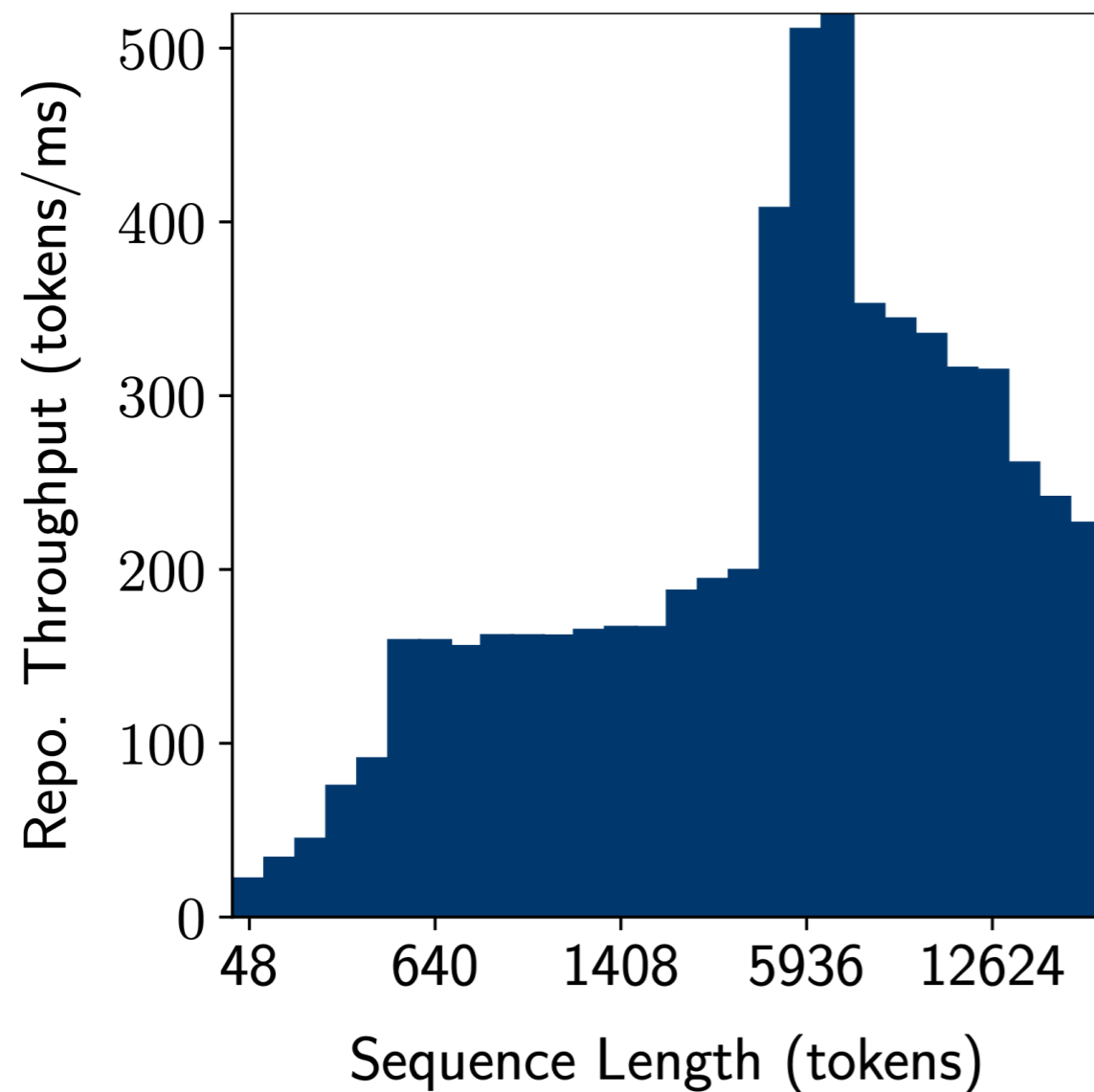
Accuracy

With span table implementation (HotpotQA)



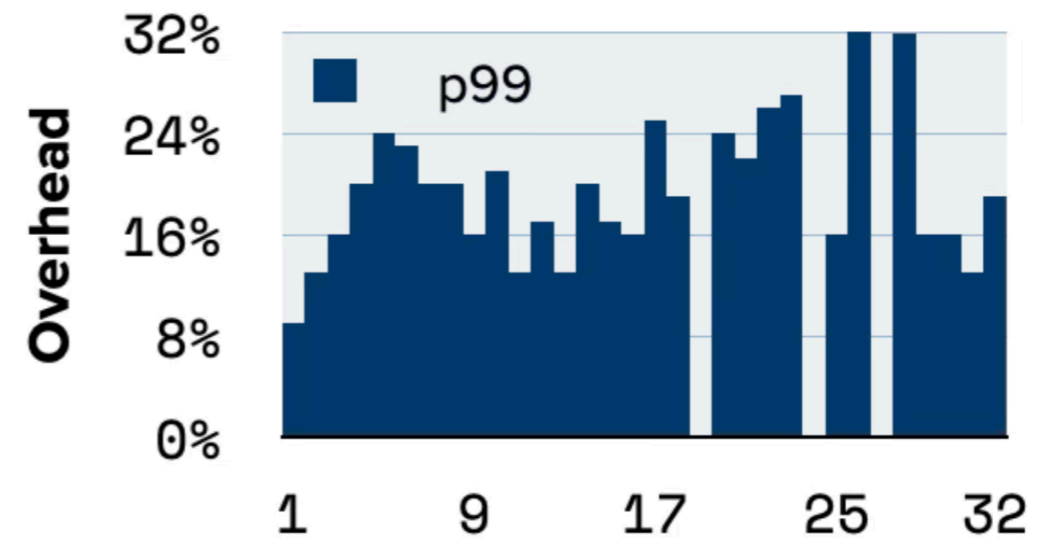
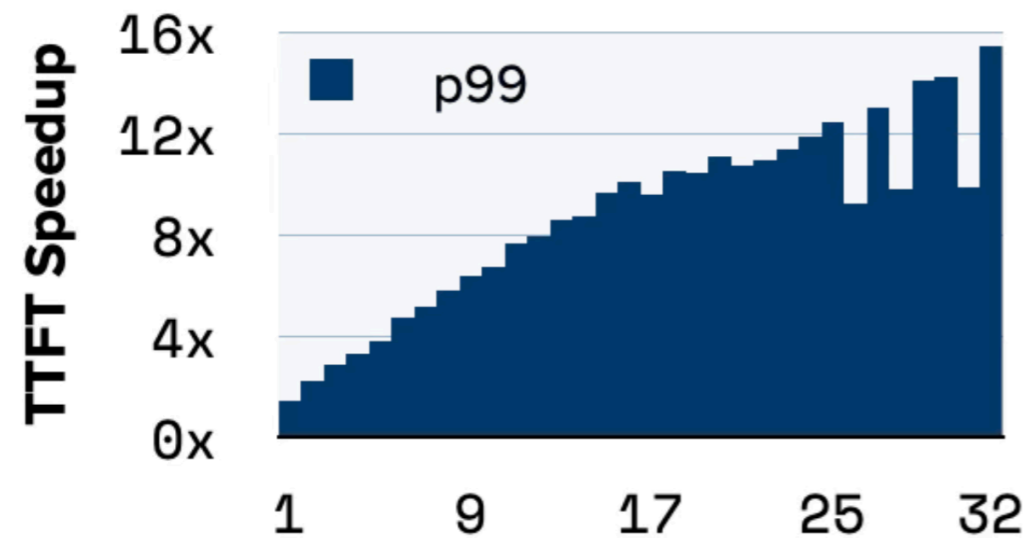
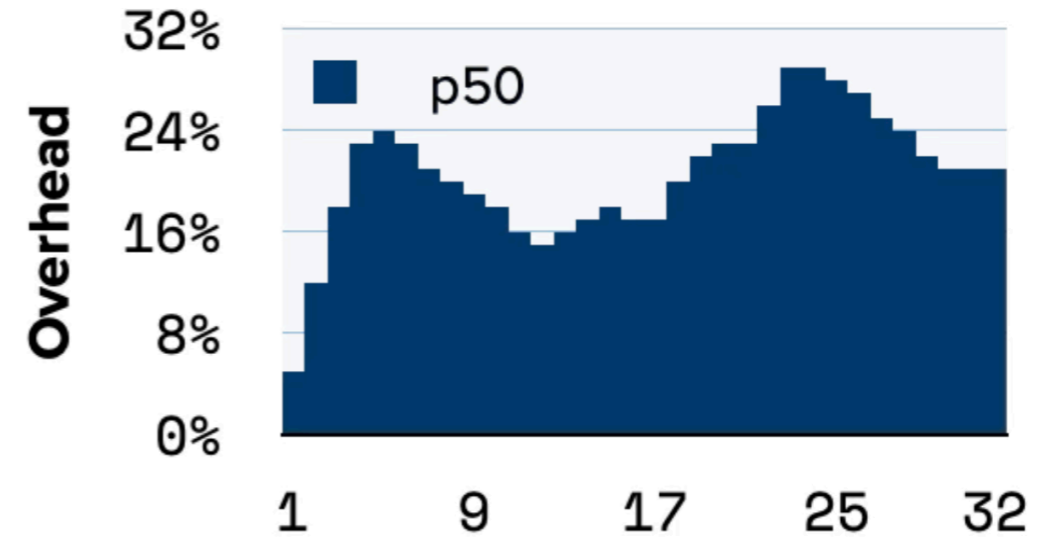
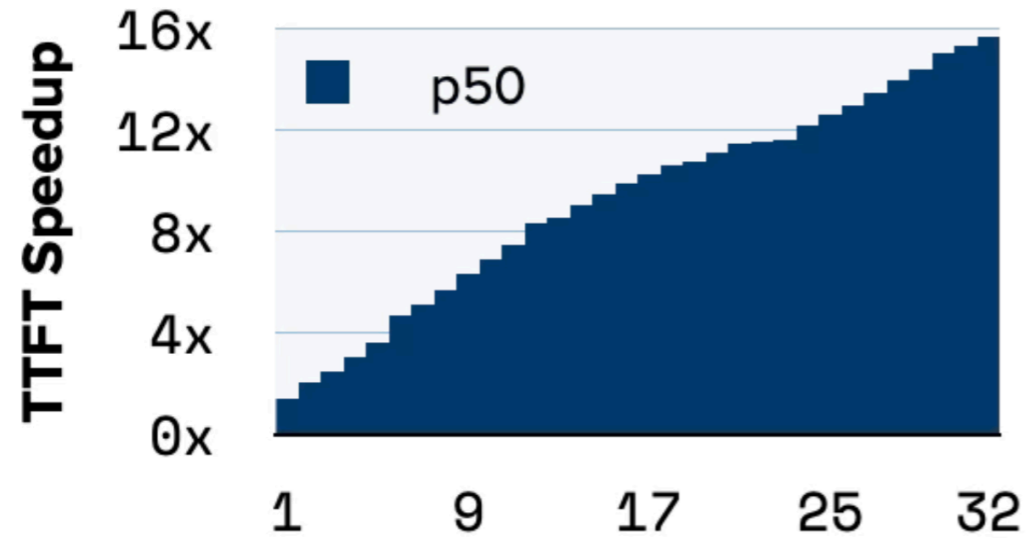
TTFT Speedup

RAG use case



TTFT Speedup

Judge/generator use case



Number of Inner Generates

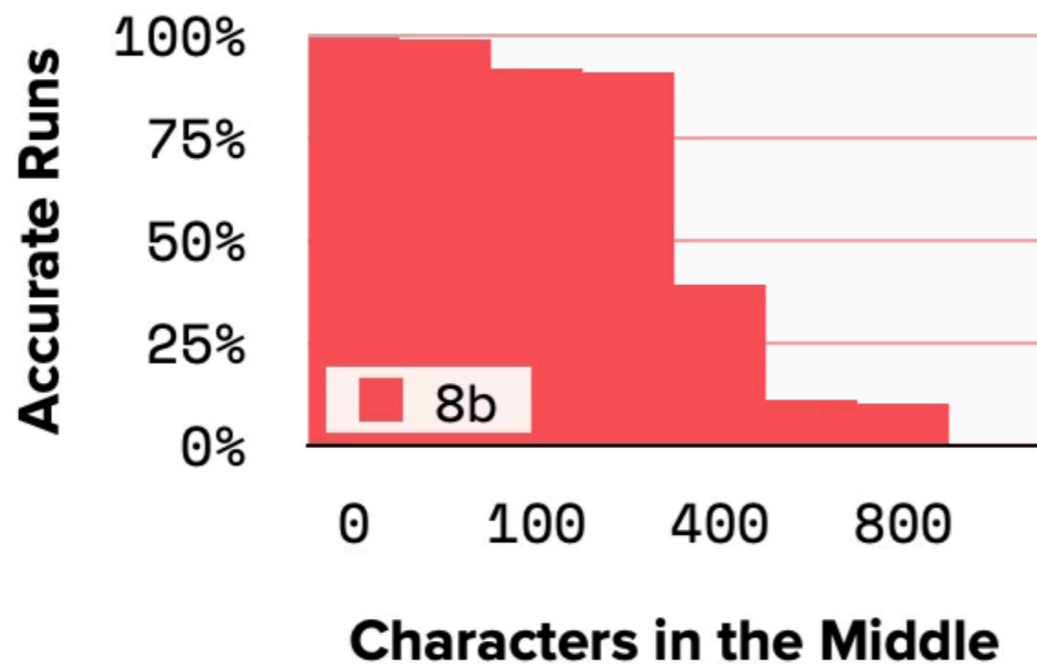
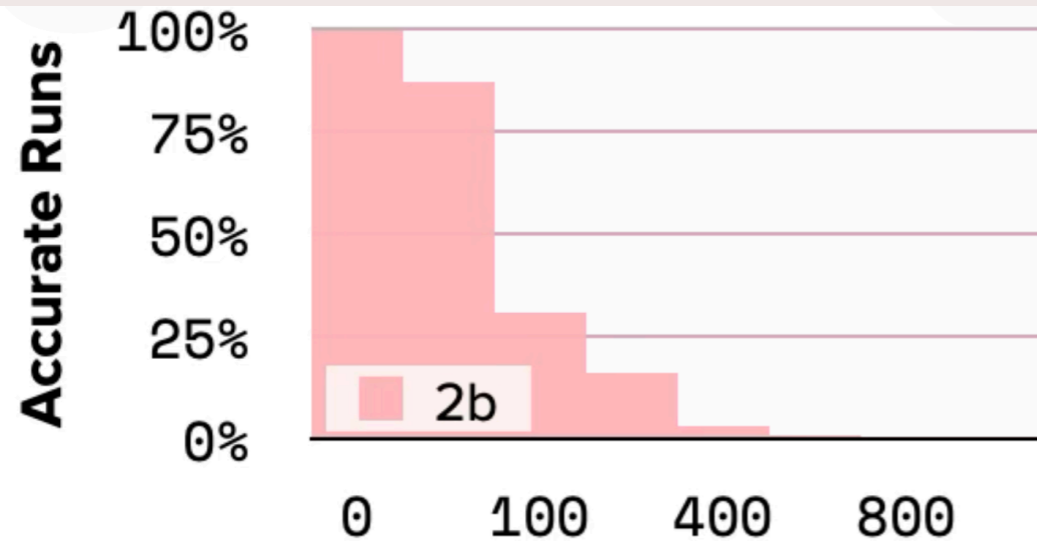
Number of Inner Generates

(a) TTFT speedup

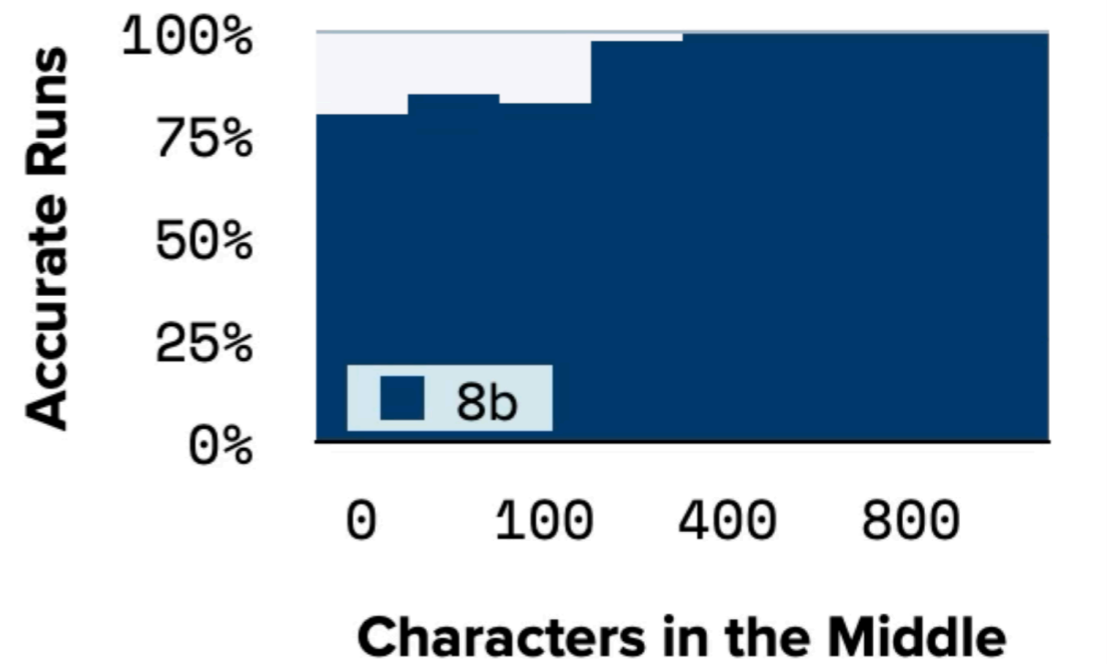
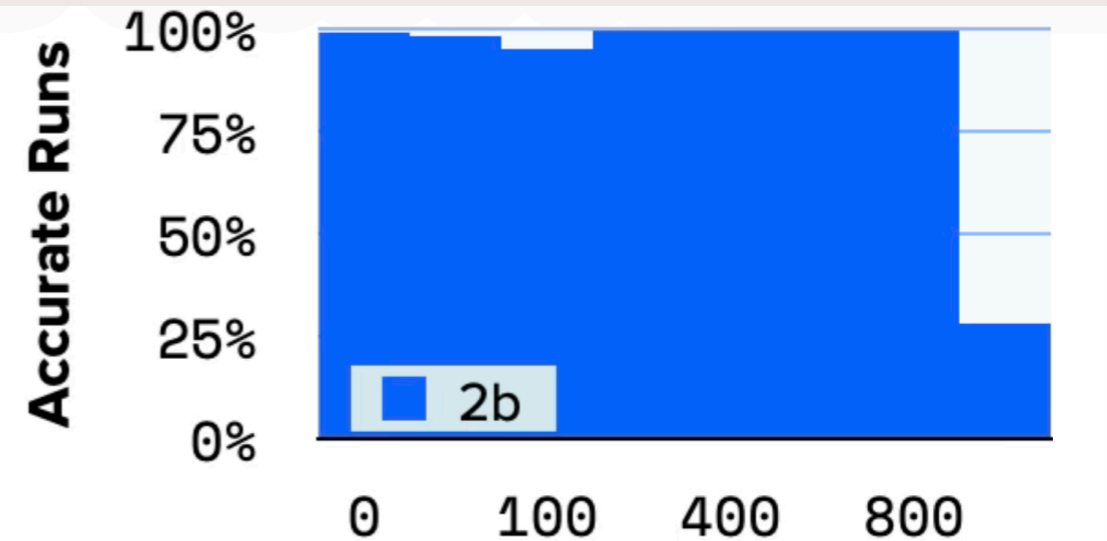
(b) Repositioning overhead

Attention Locality

Needle in haystack experiments



(a) Unoptimized query

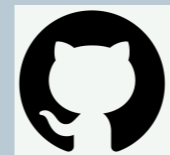


(b) Attention-optimized $k = 2$

Thanks!

Nick Mitchell, Paul Castro, Nathan
Ordonez, Thomas Parnell, Mudhakar
Srivasta, Antoni Viros i Martin

IBM Research | May 2026



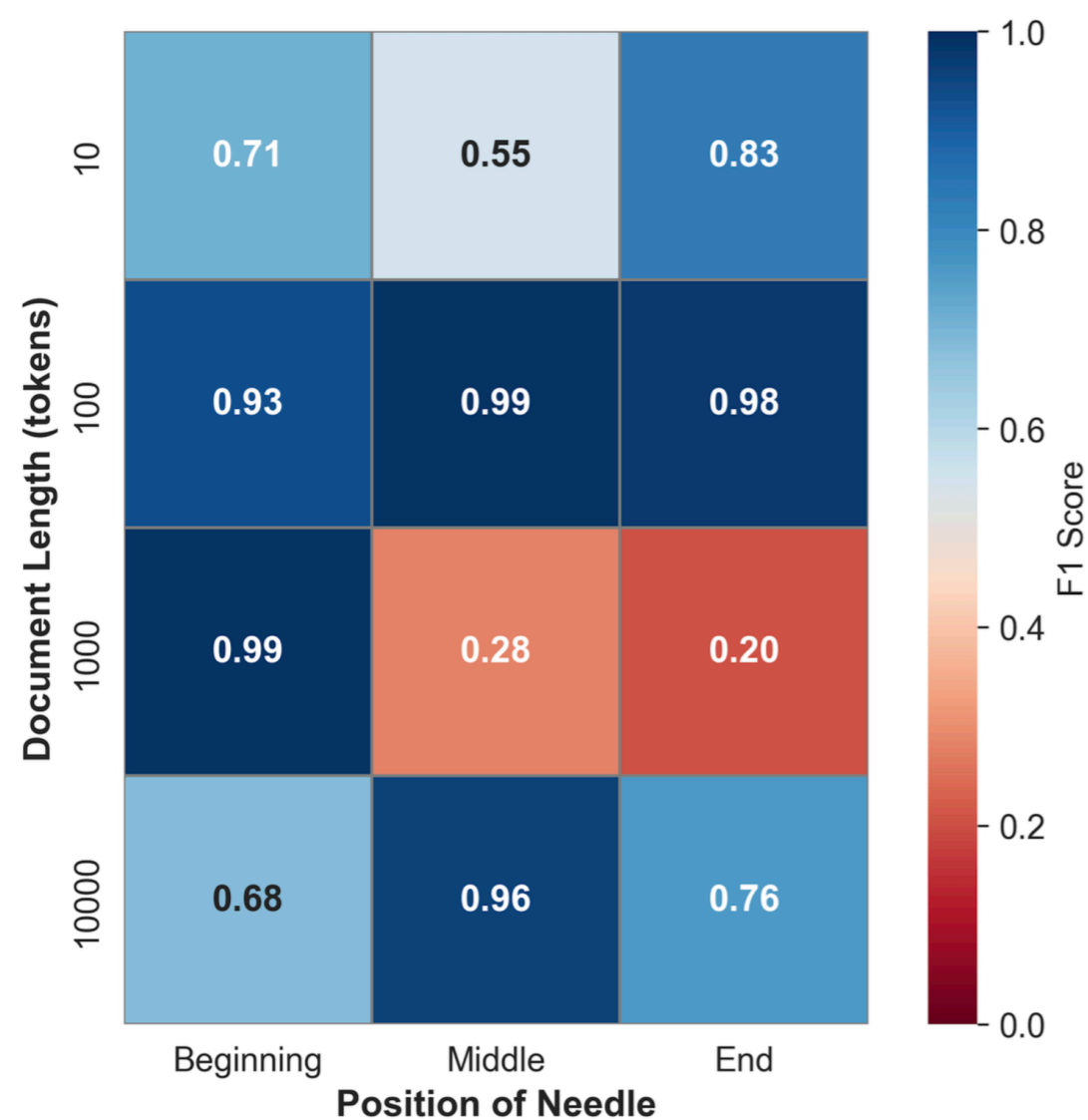
<https://github.com/IBM/spnl>



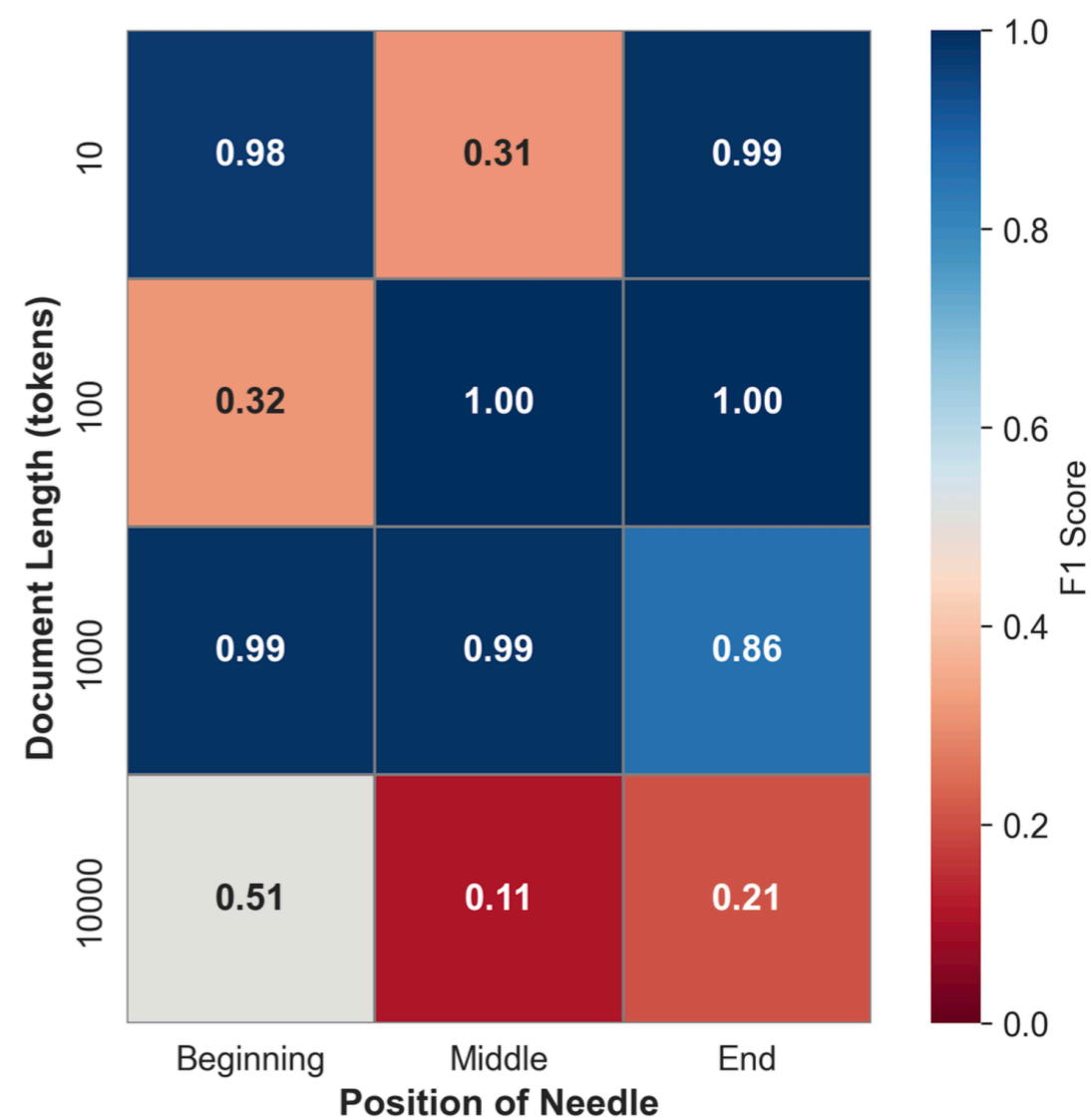
[@starpit.bsky.social](https://bsky.app/profile/starpit.bsky.social)

Accuracy

With cropping implementation



(a) Without cropping

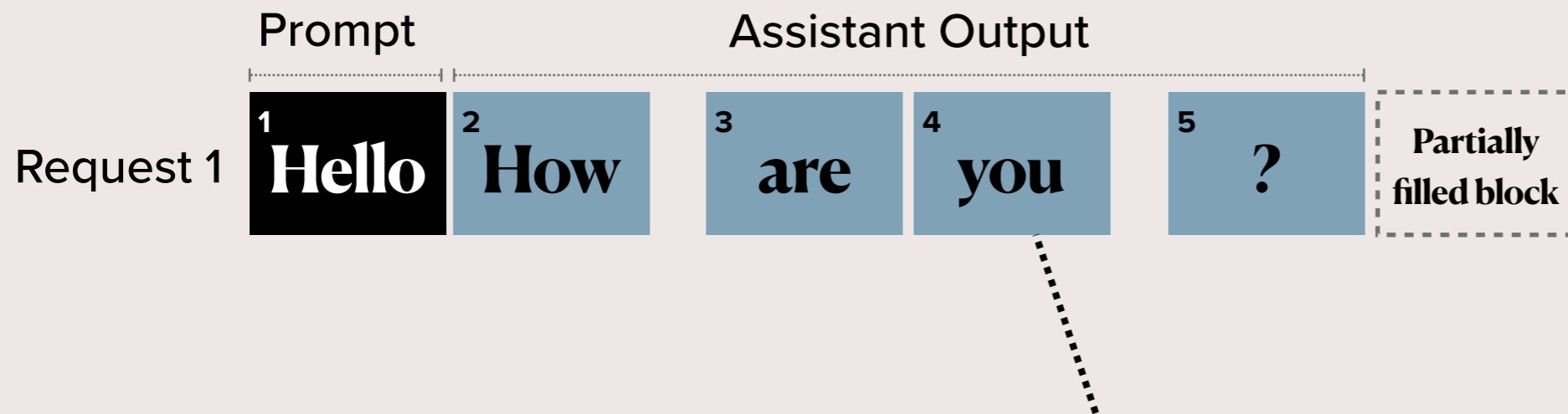


(b) With cropping

Backup Material on Tokenization

Chat: Token Sequences

After first request



Blocks in a paged attention inference server such as vLLM

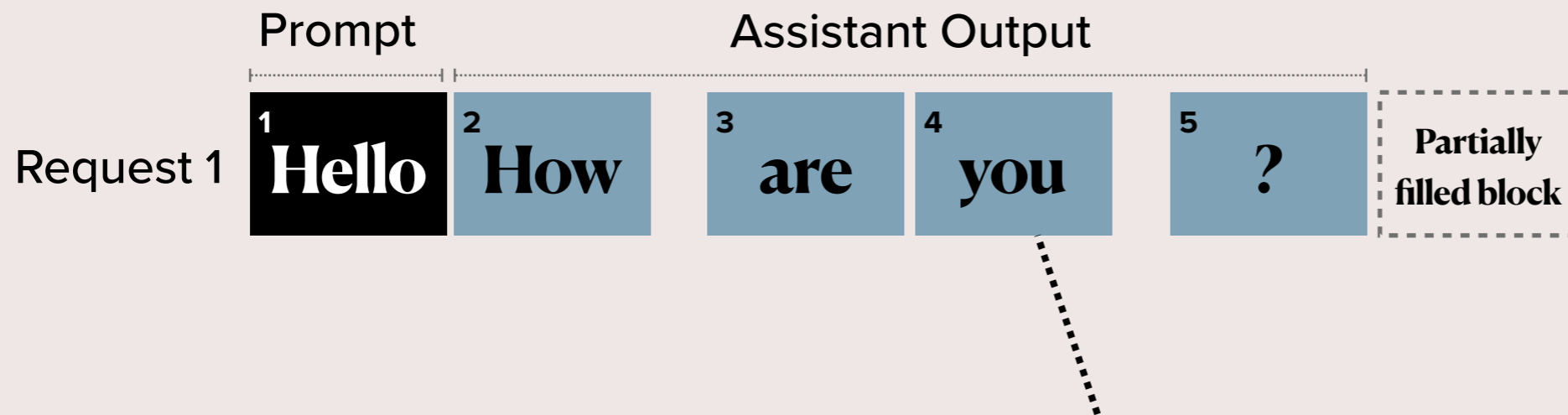
e.g. 2 tokens per block

Chat: Token Sequences

After first request

TODOs

Tally of features
we need to support
with span queries



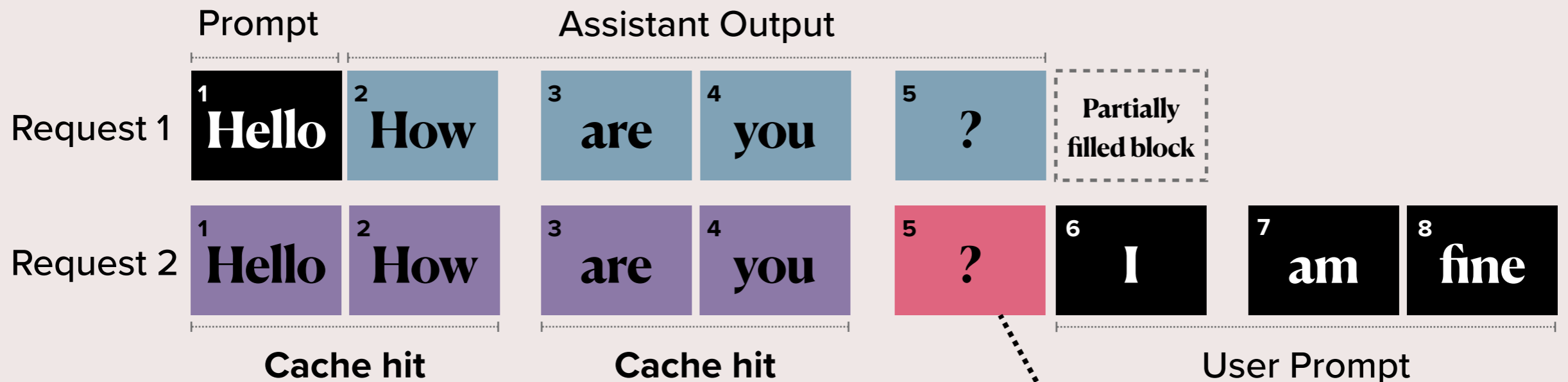
Blocks in a paged attention
inference server such as vLLM

Chat: Token Sequences

After second request

TODOs

- Prefix caching
- ~~Partial blocks~~
(1/N situation)



Hits because of prefix caching

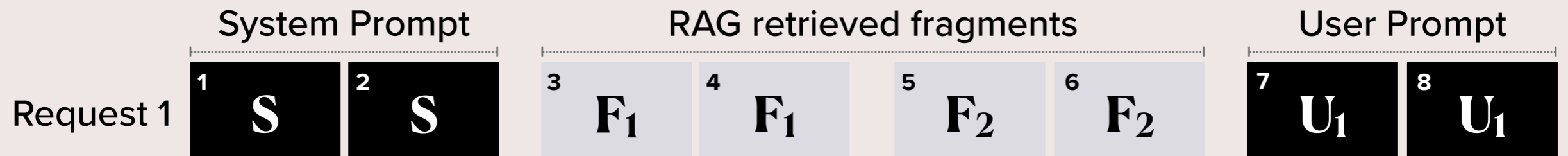
vLLM doesn't cache partial blocks,
who cares about 1/N misses

RAG: Token Sequences

After first request

TODOs

- Prefix caching

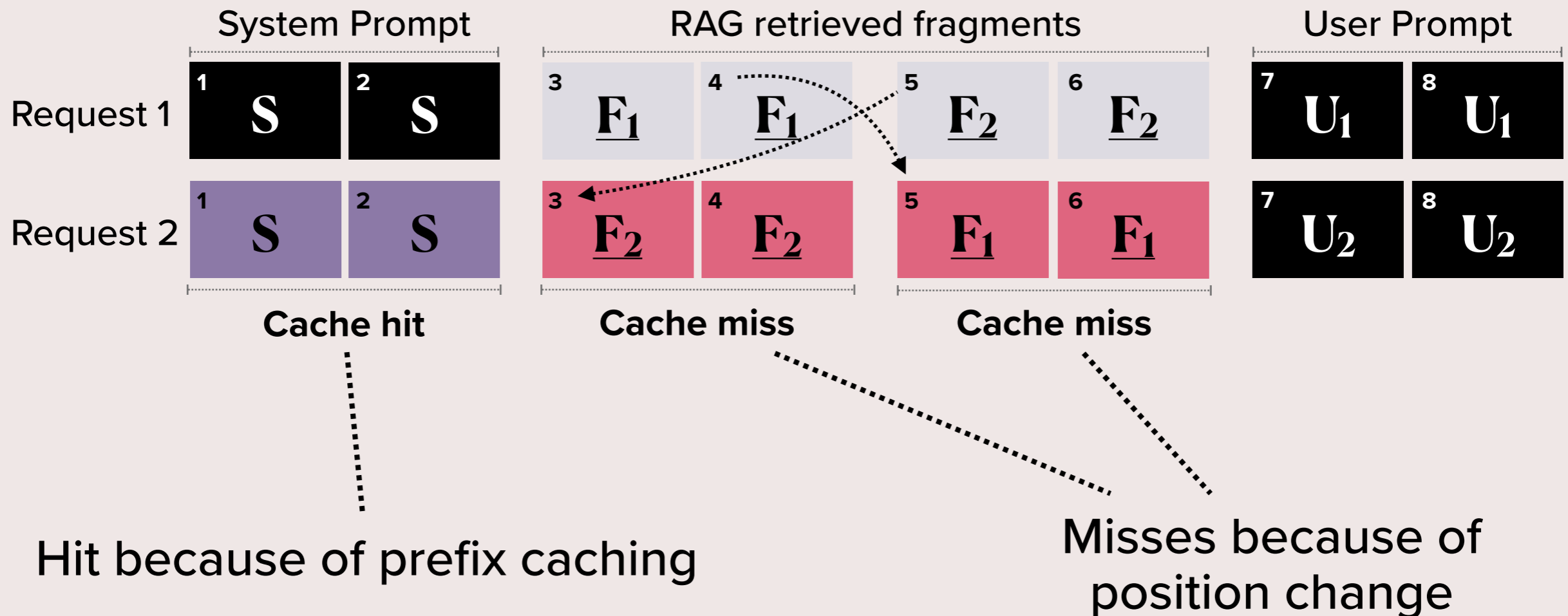


RAG: Token Sequences

After second request

TODOs

- Prefix caching
- Position change

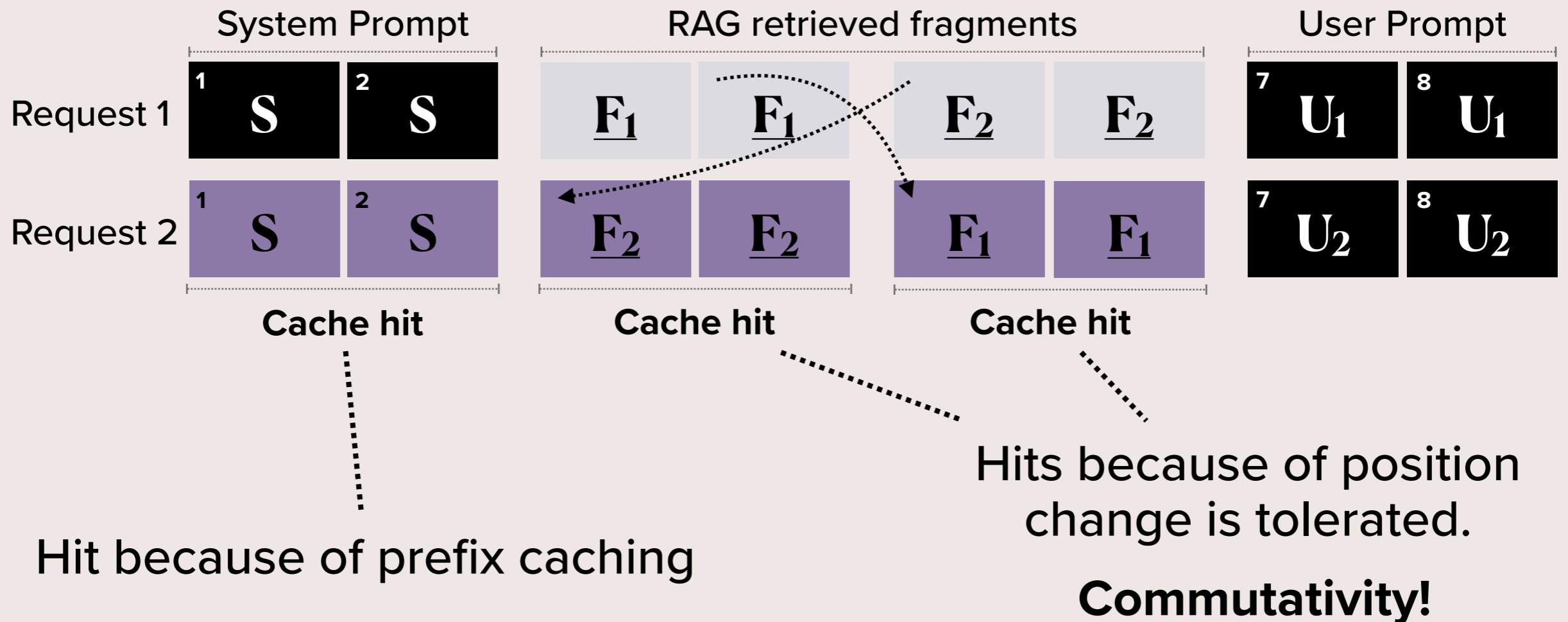


RAG: Token Sequences

After selectively deleting positions

TODOs

- Prefix caching
- Position change

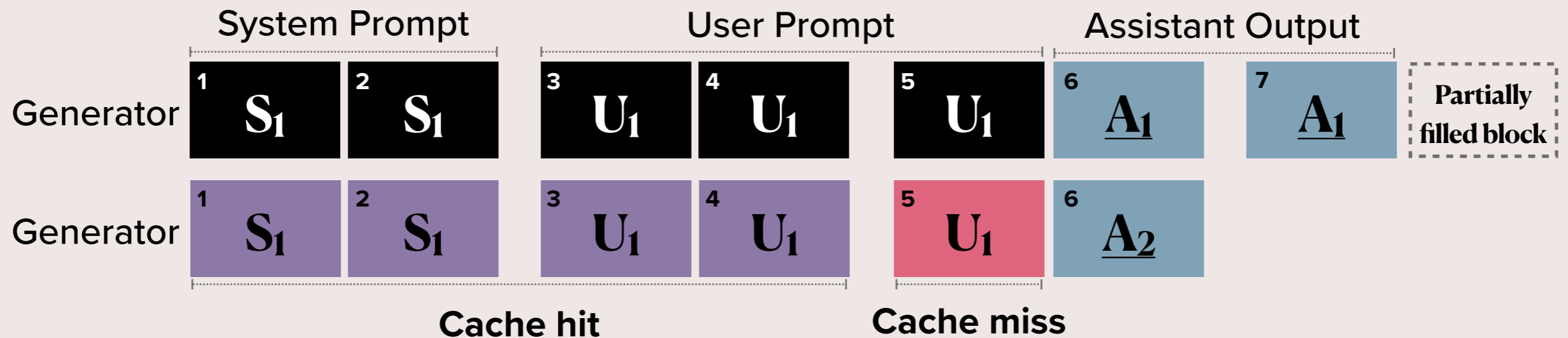


Judge/Generator: Token Sequences

After two generators

TODOs

- Prefix caching
- Partial blocks



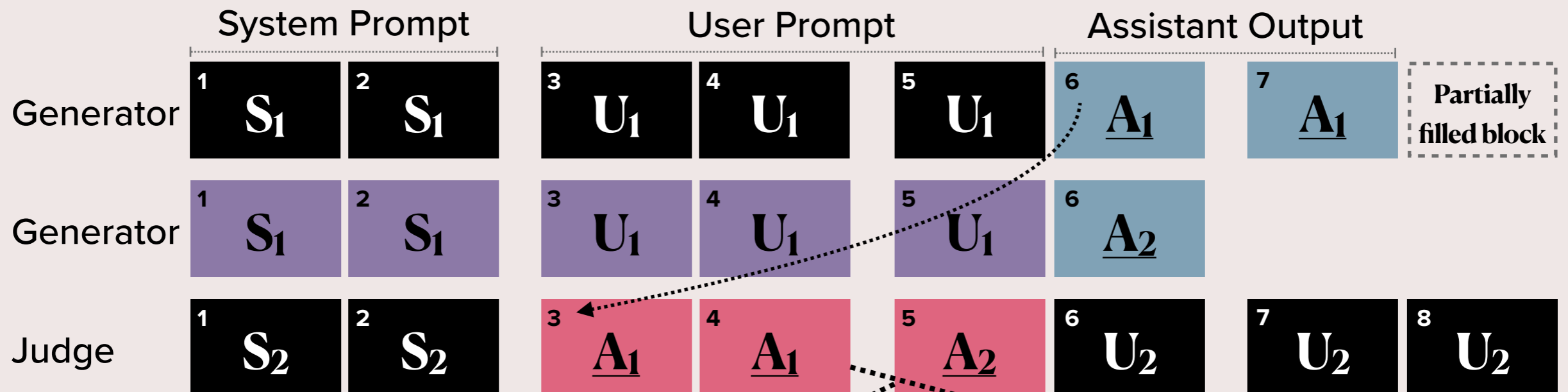
Partially filled block, resulting in mixed block content, but maybe who cares because 1/N...

Judge/Generator: Token Sequences

Two generators and one judge

TODOs

- Prefix caching
- Partial blocks
- Position change



Misses because **prefix** and **position** both have changed

Misses because **partial block**

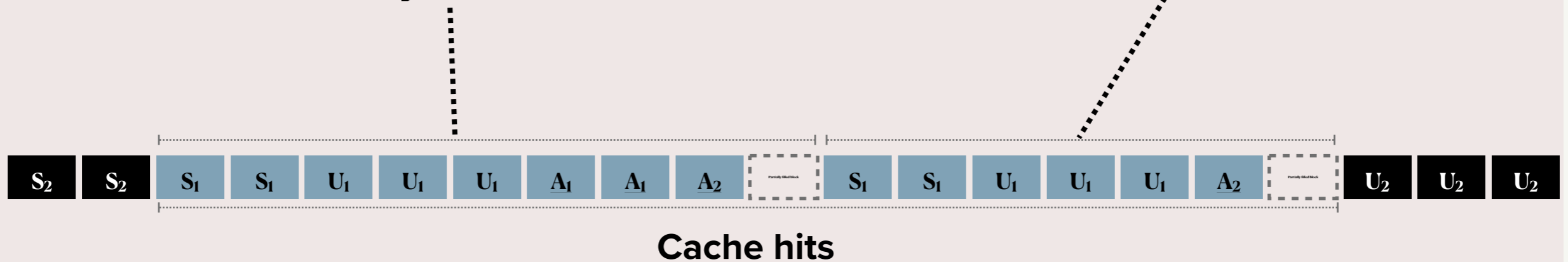
Judge/Generator: Token Sequences

What would give us 100% cache hits TODOs

- Prefix caching
- Partial blocks
- Position change

Hits, because, with **padding or cropping**, and commutativity, the this is literally Generate 1

Same, from Generate 2



Low-level Optimizations

Token-level transformations required to work with vLLM

We may need to reintroduce prefixes for agentic workloads

We may need to pad or crop assistant output

We need to allow for the expression of commutativity, and be able to selectively delete positions in commutative intervals.

TODOs

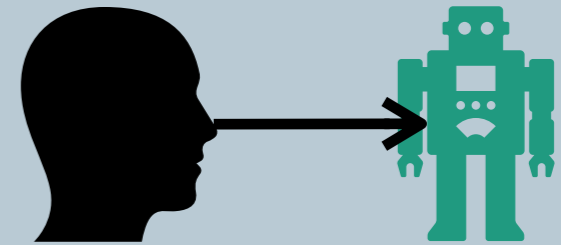
- Prefix caching
- Partial blocks
- Position change

The Use Cases

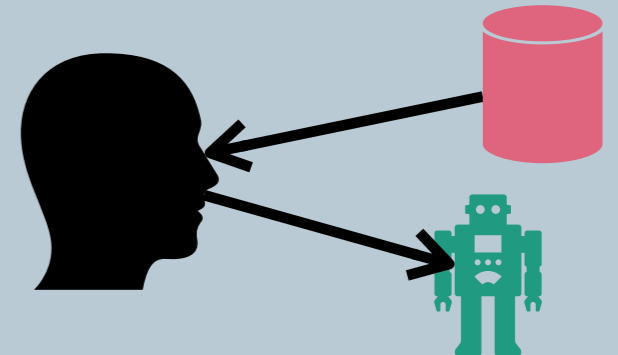
Use cases are evolving.

The **inference API needs to evolve** along with them.

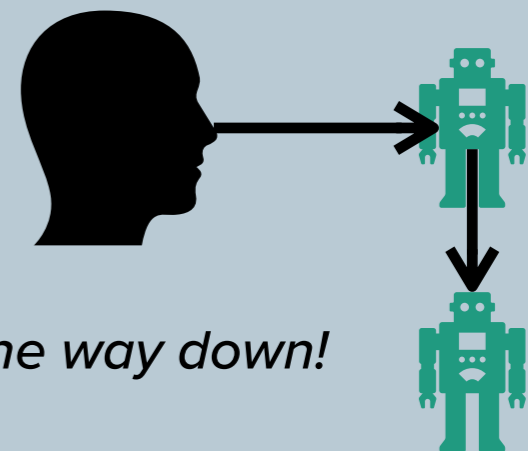
Chat



RAG



BOTS



It's bots all the way down!