

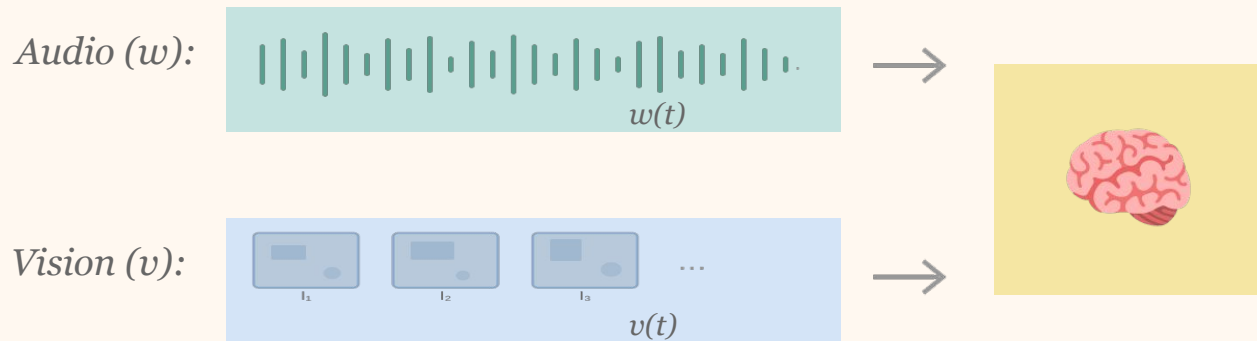
Stream Everything

Moving from Request input to Streaming input

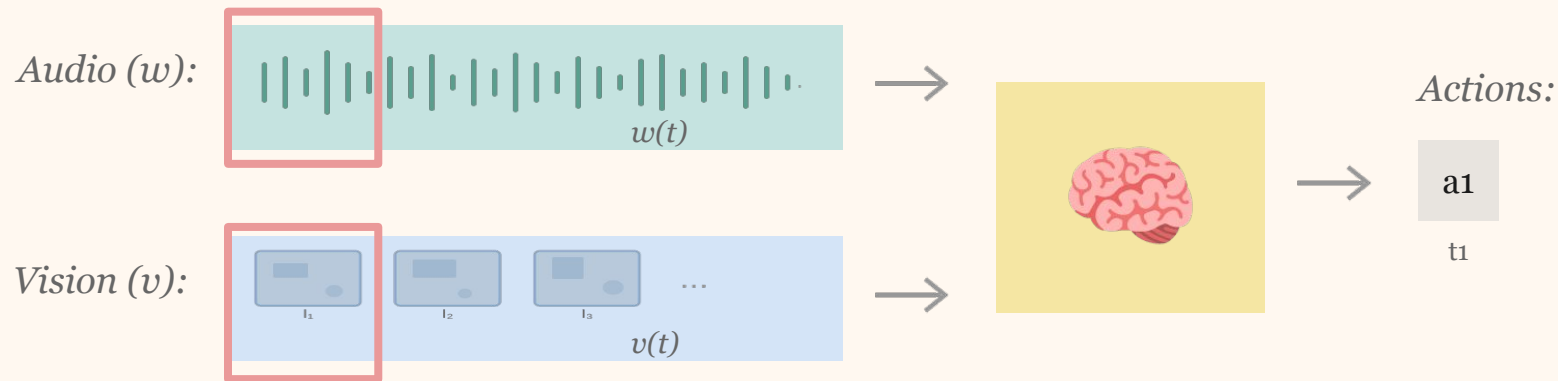
Patrick von Platen, Mistral AI



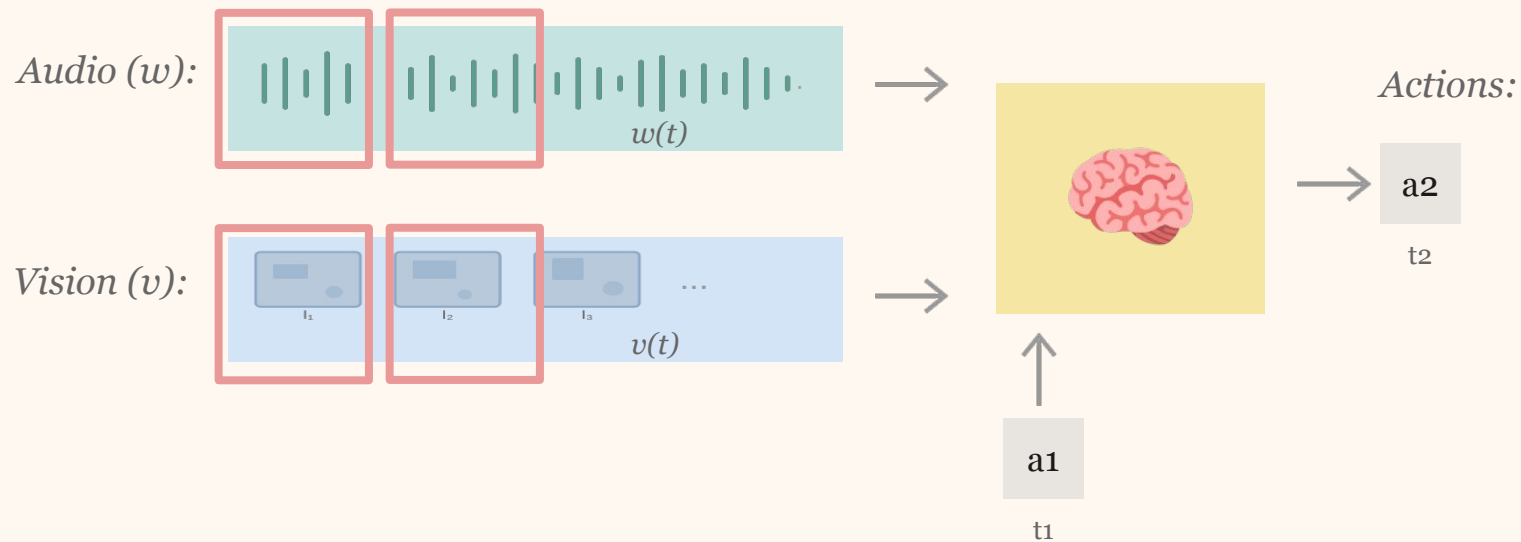
How humans process the world



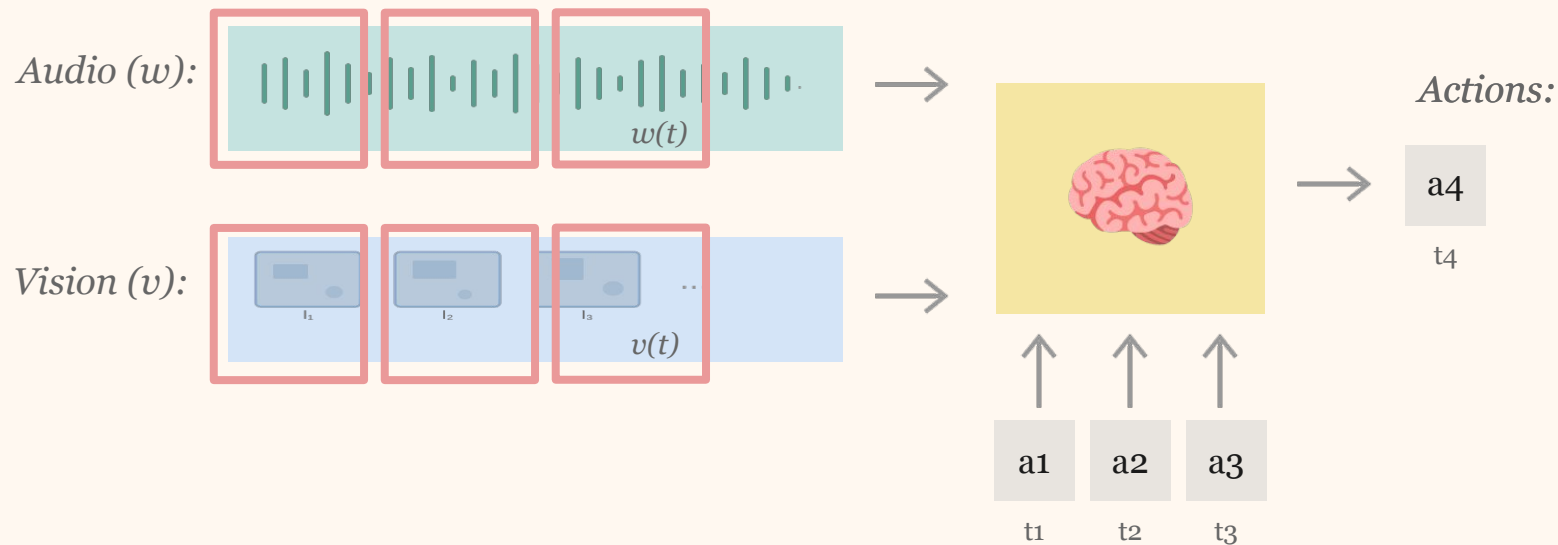
How humans process the world



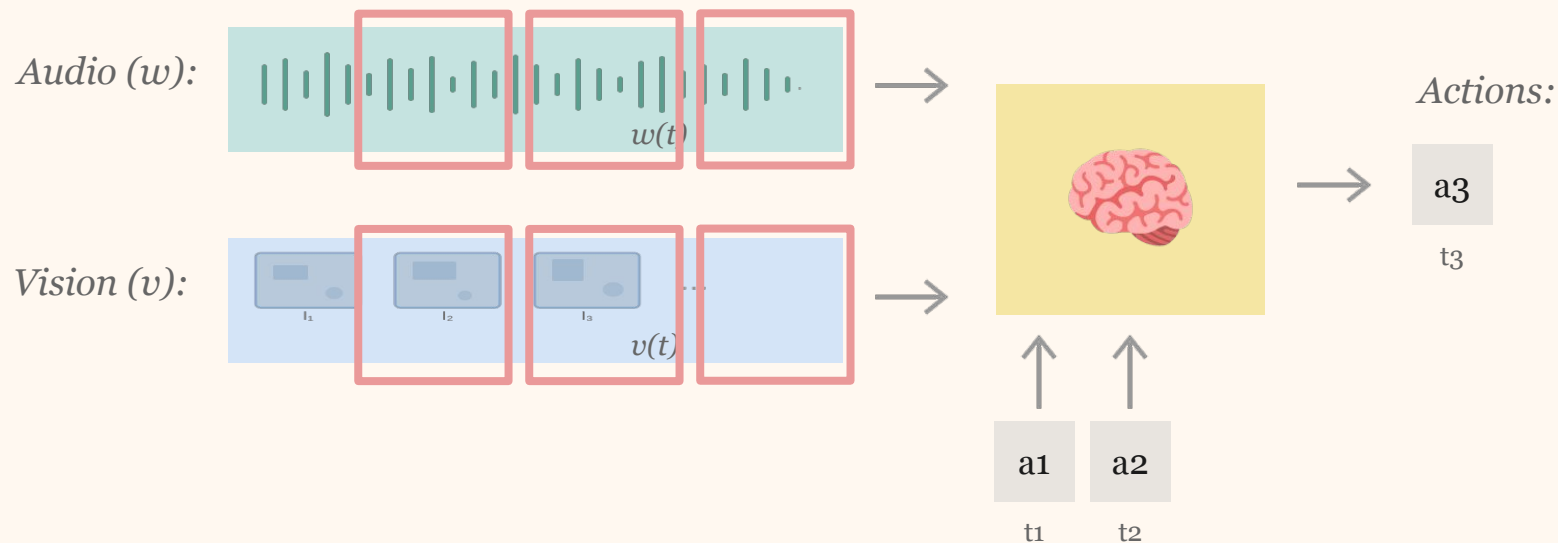
How humans process the world



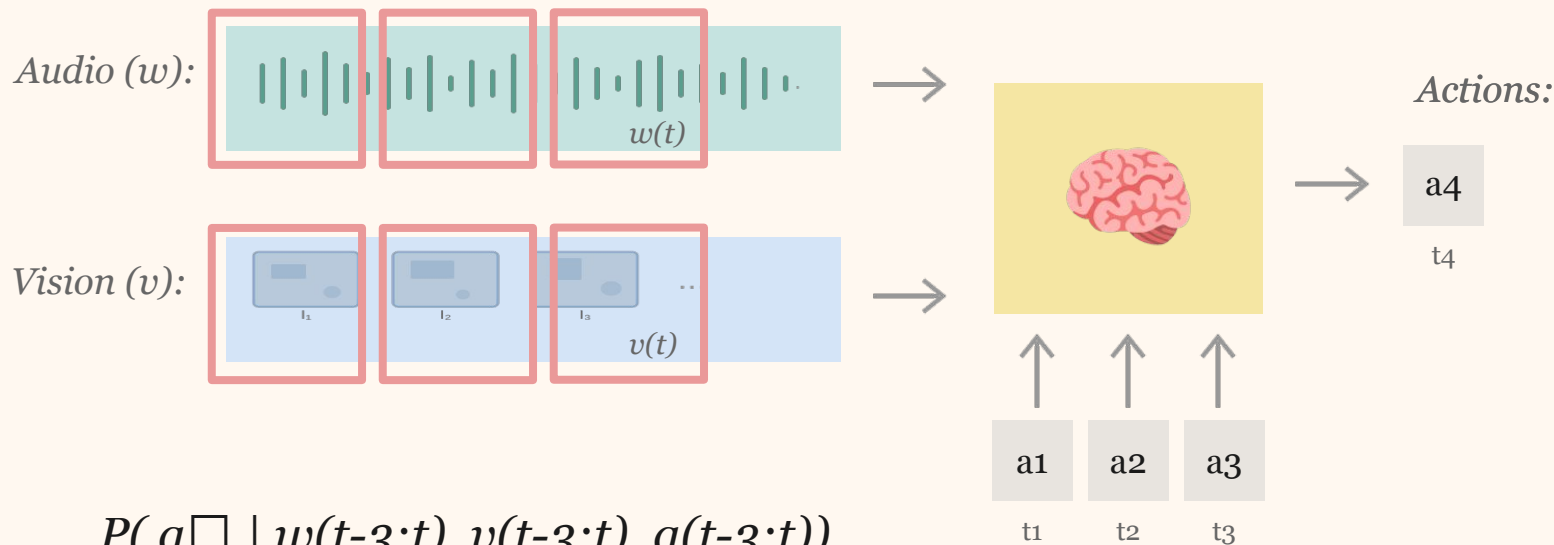
How humans process the world



How humans process the world



How humans process the world



$$P(a_t | w(t-3:t), v(t-3:t), a(t-3:t))$$

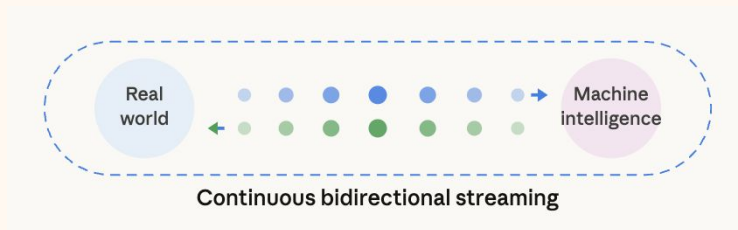
At every time step t , an action a is taken based on:

- Audio input stream
- Visual input stream
- Previous actions

■ Why does that matter for LLMs?

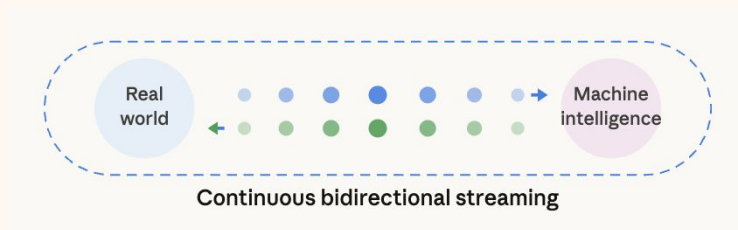
■ Why does that matter for LLMs?

Flawless real world <> machine interaction needs streaming!



■ Why does that matter for LLMs?

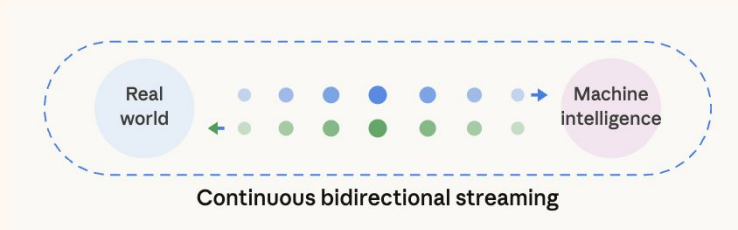
Flawless real world <> machine interaction needs streaming!



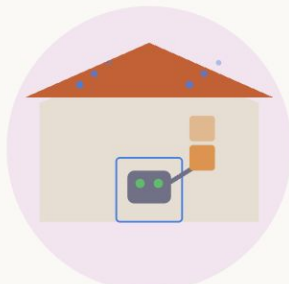
Autonomous driving

■ Why does that matter for LLMs?

Flawless real world <> machine interaction needs streaming!



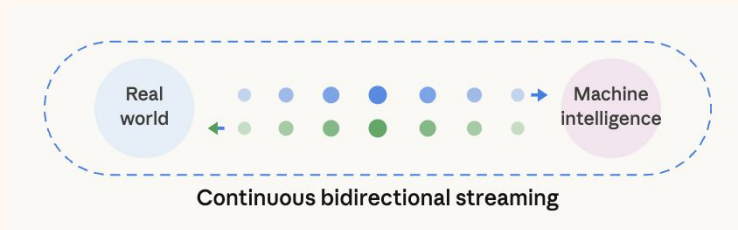
Autonomous driving



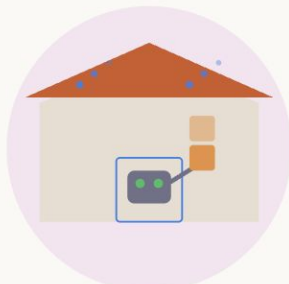
AI in warehouses

■ Why does that matter for LLMs?

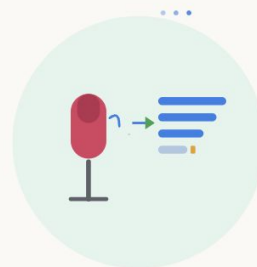
Flawless real world <> machine interaction needs streaming!



Autonomous driving



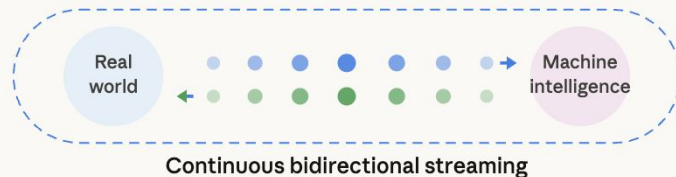
AI in warehouses



Live transcription

■ Why does that matter for LLMs?

Flawless real world <> machine interaction needs streaming!



Autonomous driving



AI in warehouses



Live transcription



Voxtral Realtime

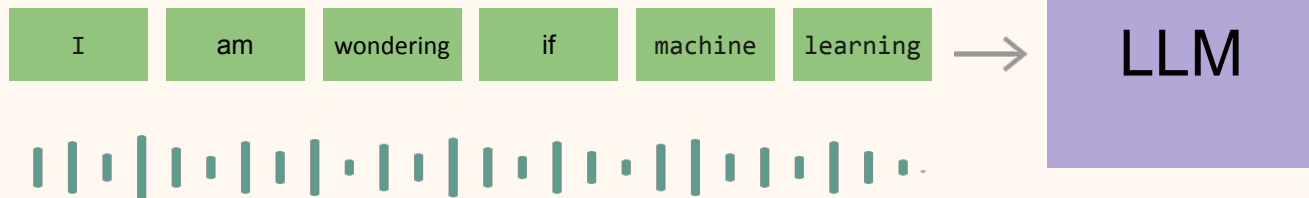
Mistral AI

Abstract

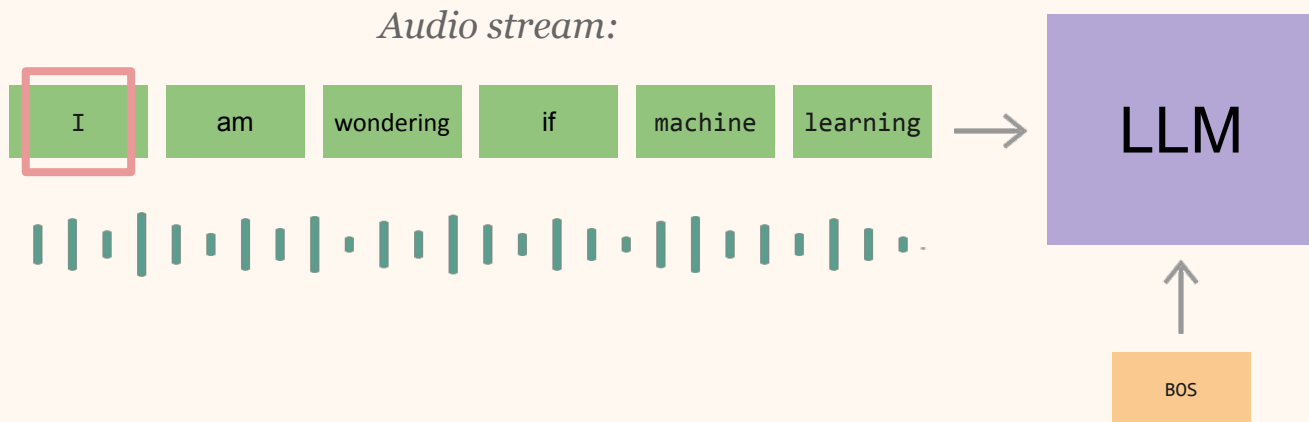
We introduce Voxtral Realtime, a natively streaming automatic speech recognition

■ Speech Recognition - Realtime

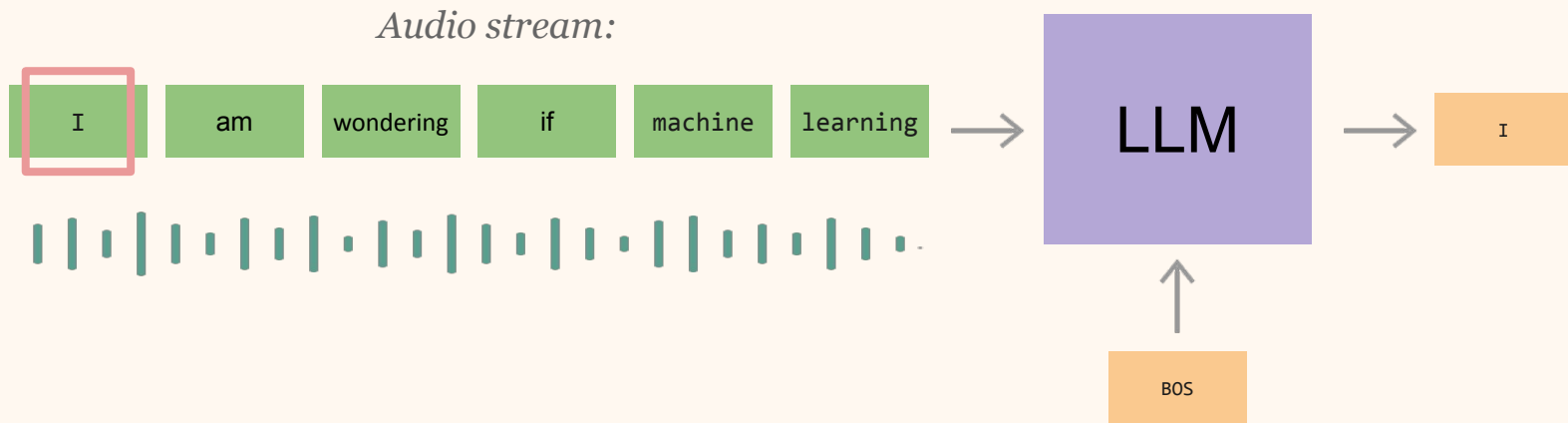
Audio stream:



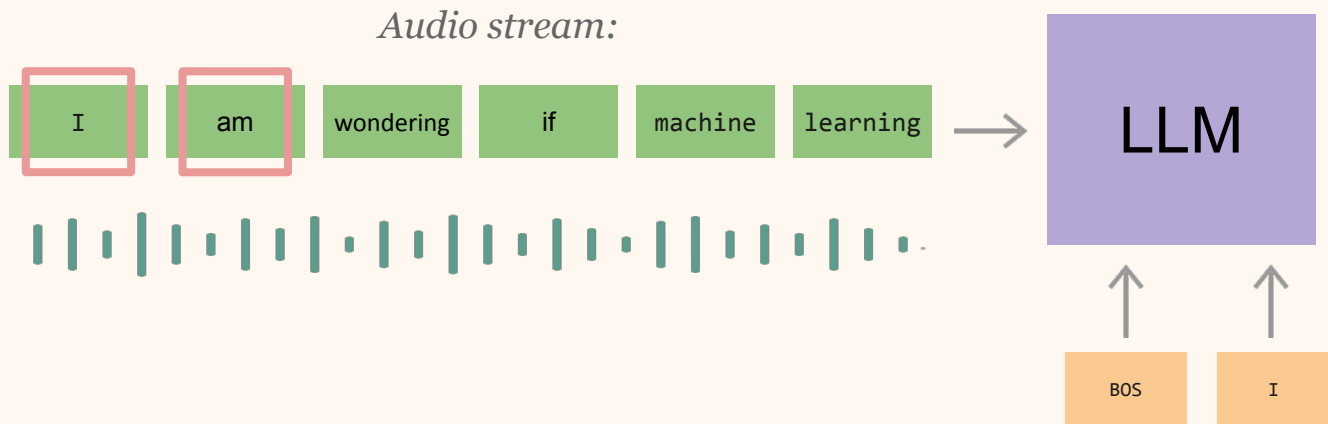
Speech Recognition - Realtime



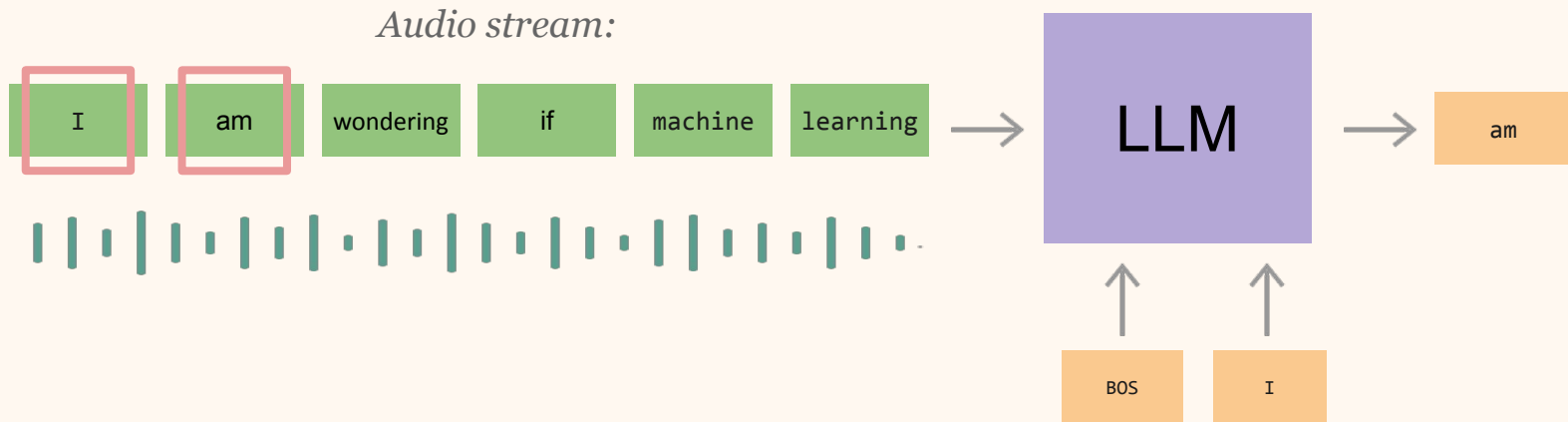
Speech Recognition - Realtime



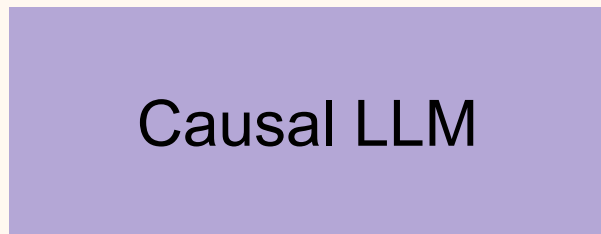
Speech Recognition - Realtime



Speech Recognition - Realtime



Speech Recognition - Realtime



Hid emb



I



BOS

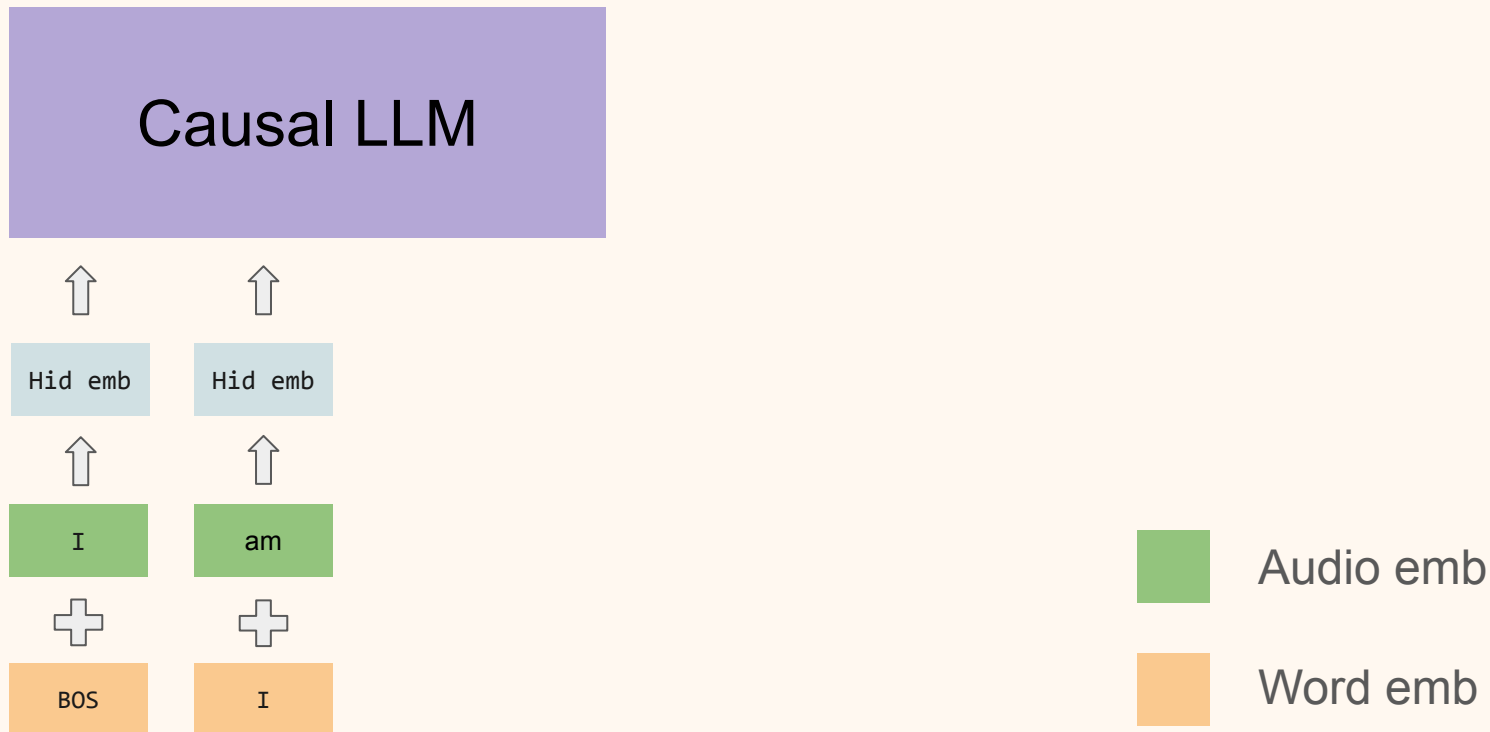


Audio emb



Word emb

Speech Recognition - Realtime

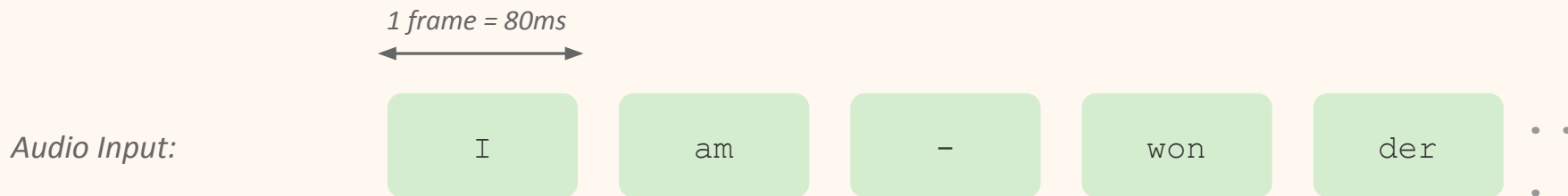


Speech Recognition - Realtime



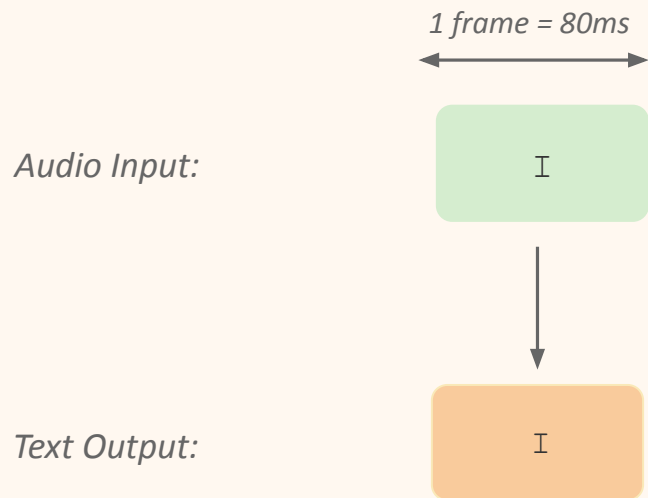
■ Speech Recognition - Realtime

Audio is moving faster than text... what to predict?



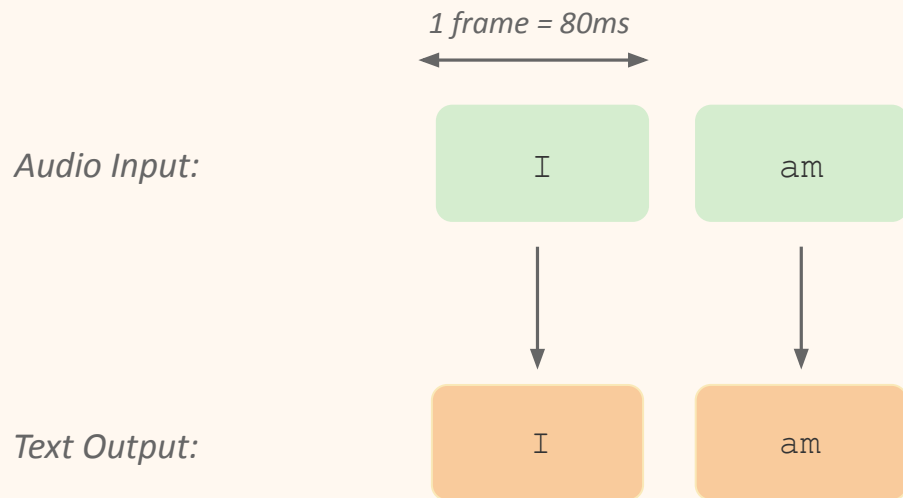
■ Speech Recognition - Realtime

Audio is moving faster than text... what to predict?



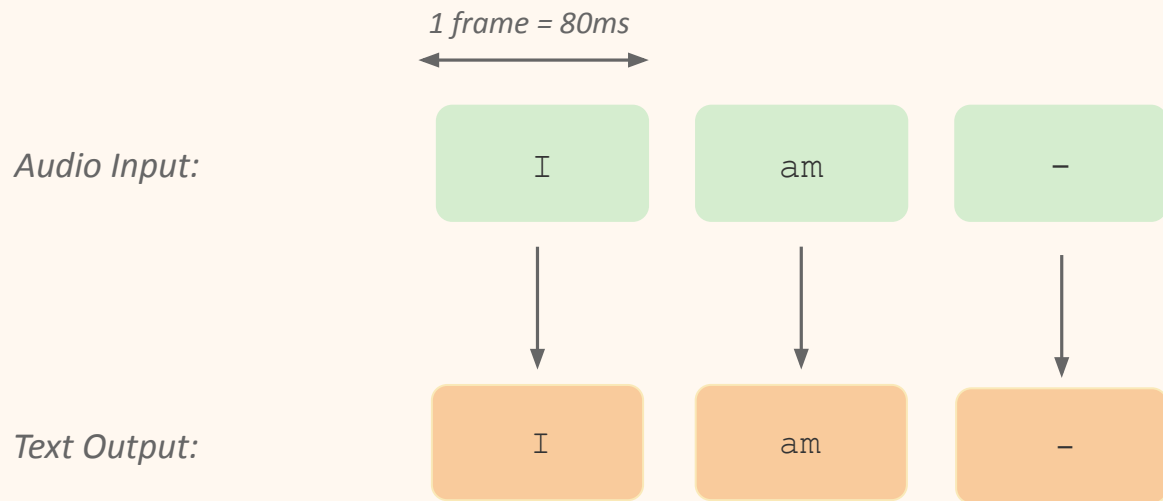
■ Speech Recognition - Realtime

Audio is moving faster than text... what to predict?



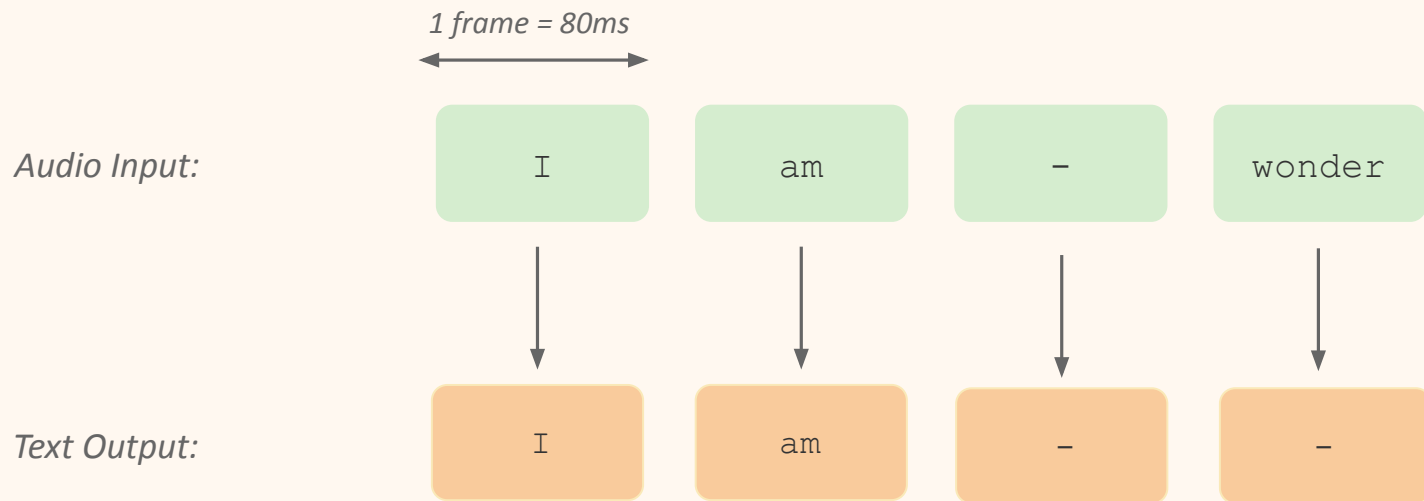
■ Speech Recognition - Realtime

Audio is moving faster than text... what to predict?



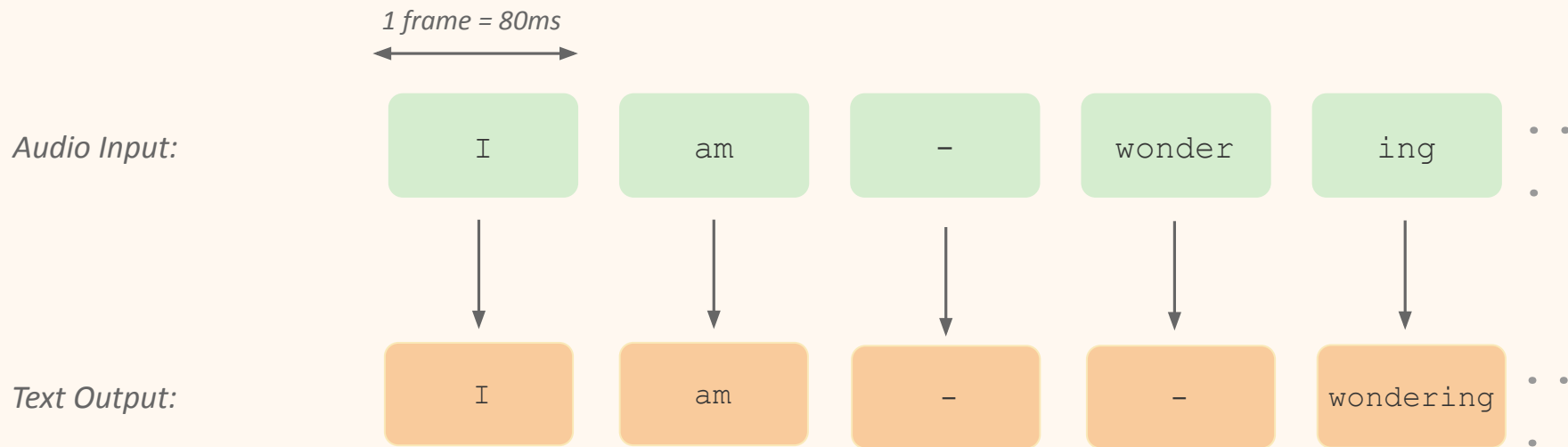
■ Speech Recognition - Realtime

Audio is moving faster than text... what to predict?

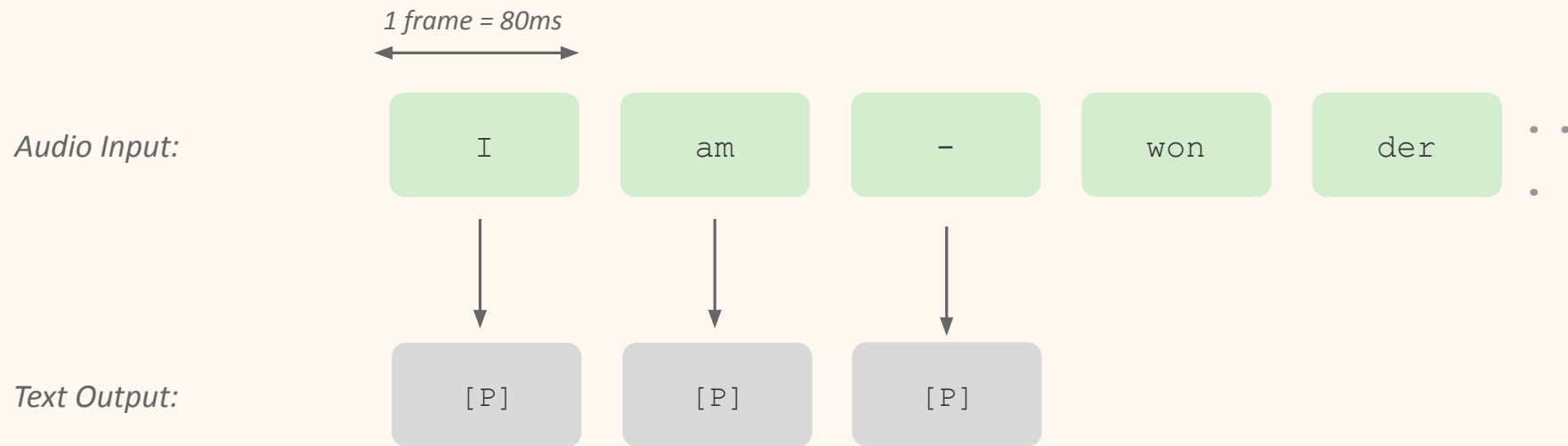


■ Speech Recognition - Realtime

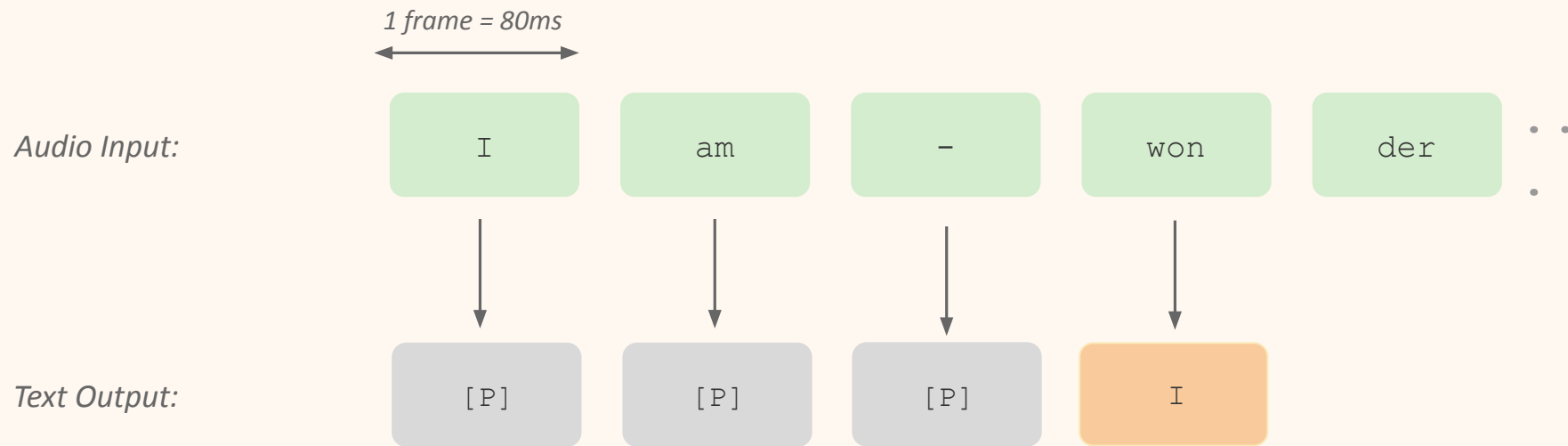
Audio is moving faster than text... what to predict?



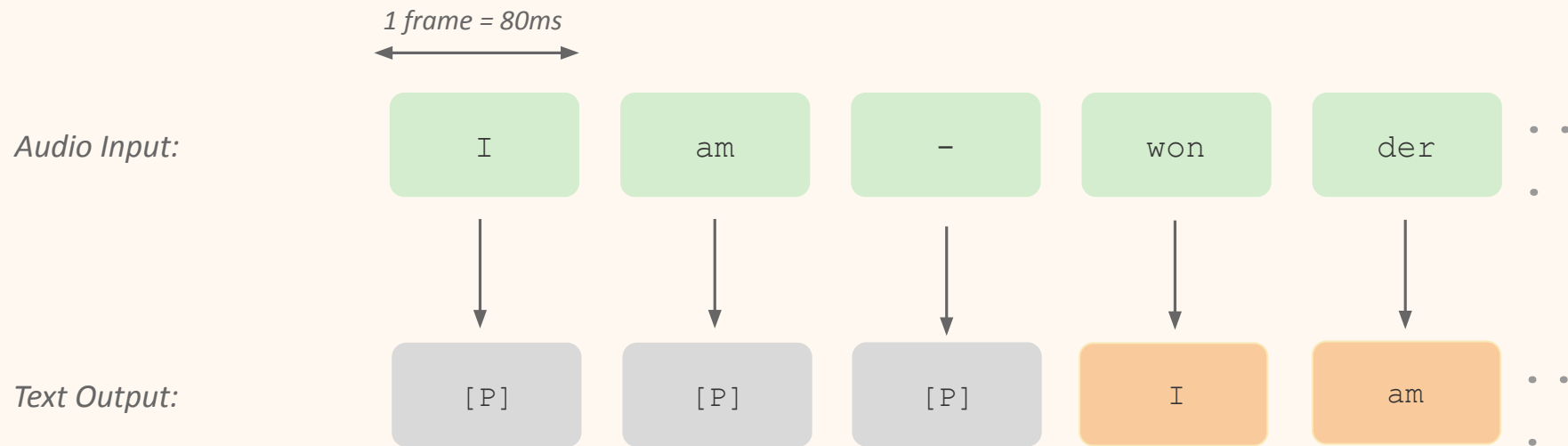
Speech Recognition - Realtime



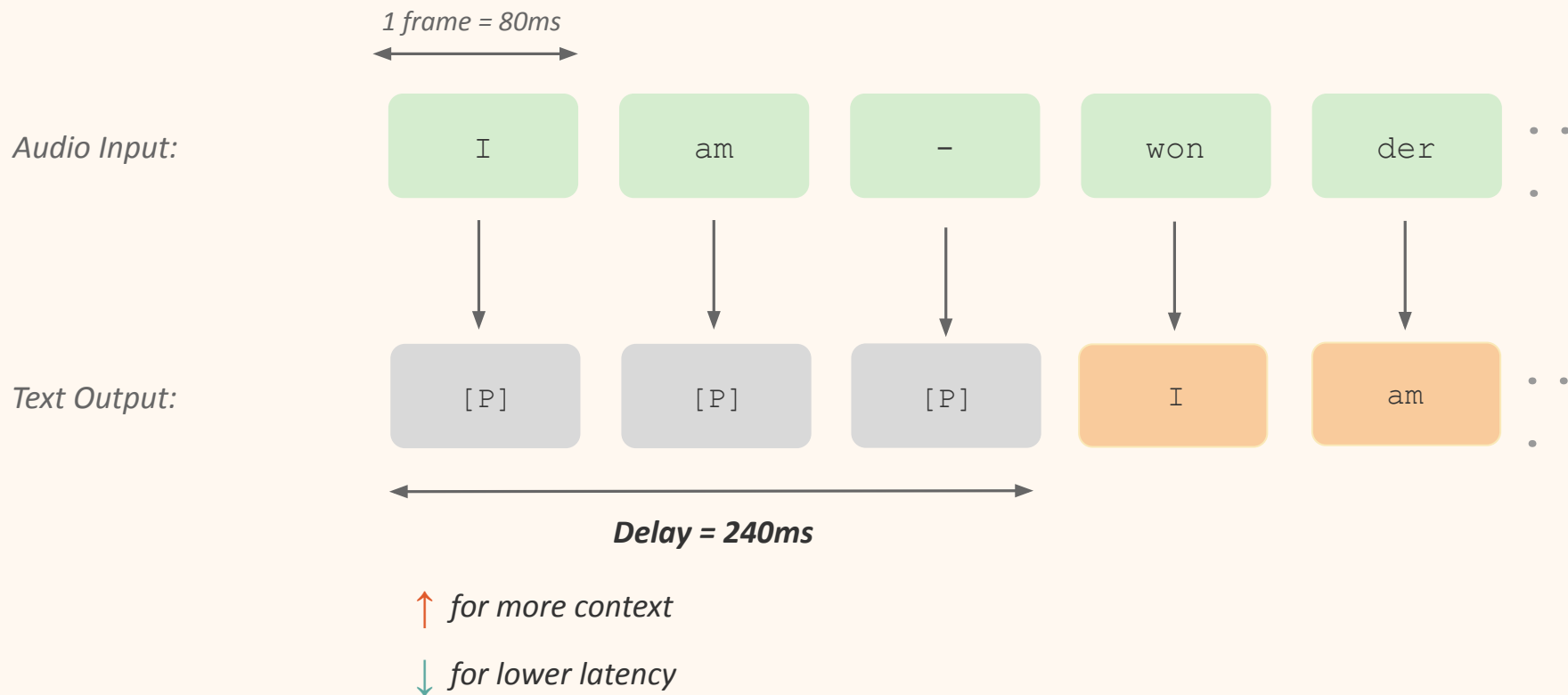
Speech Recognition - Realtime



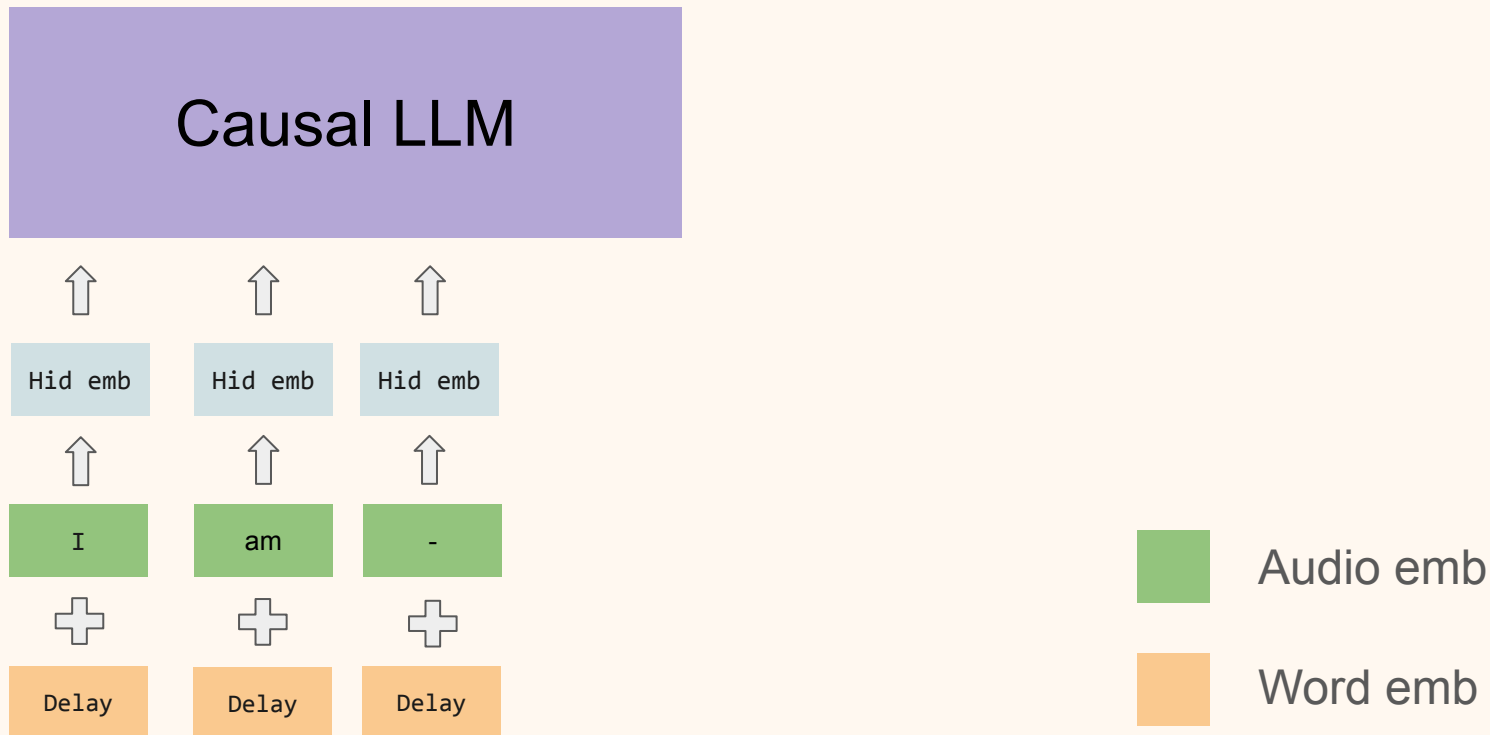
Speech Recognition - Realtime



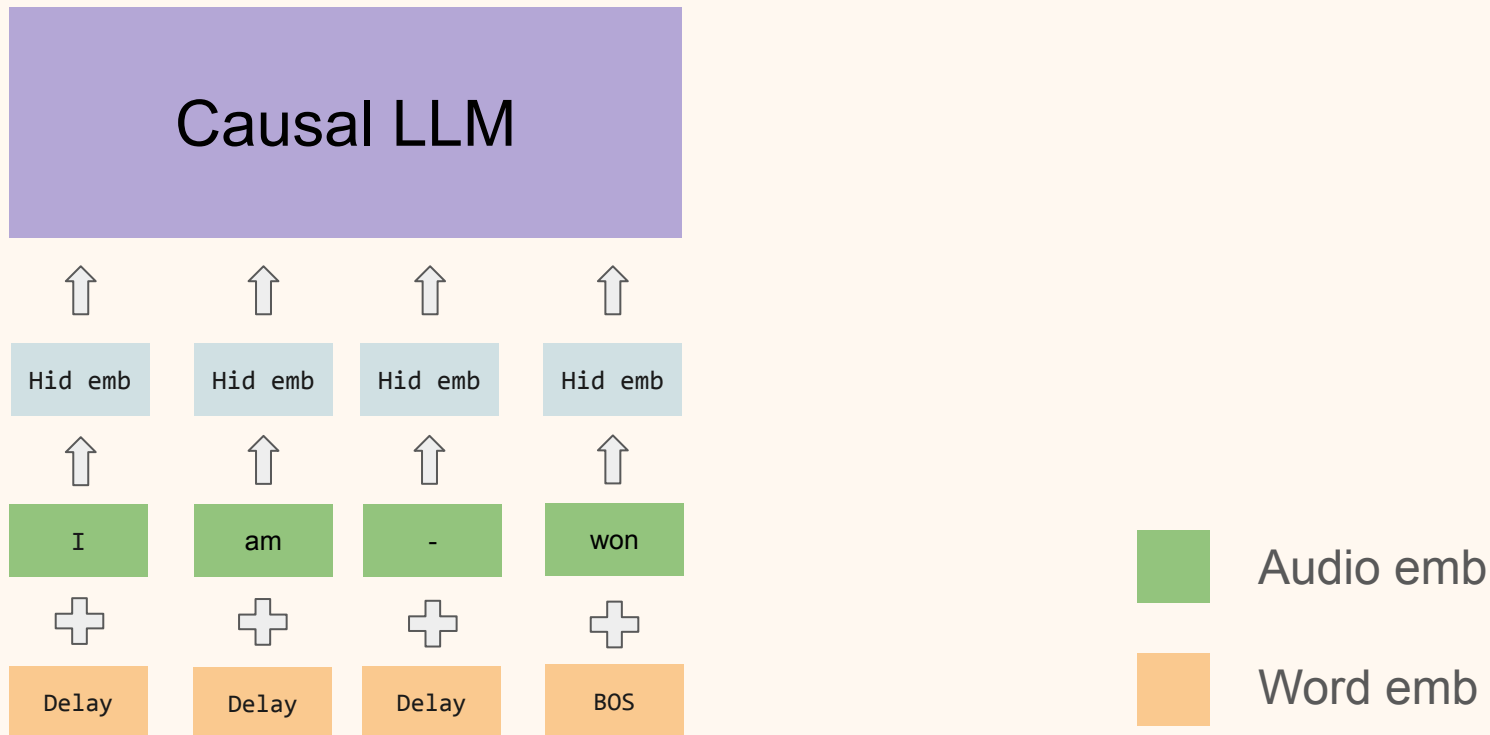
Speech Recognition - Realtime



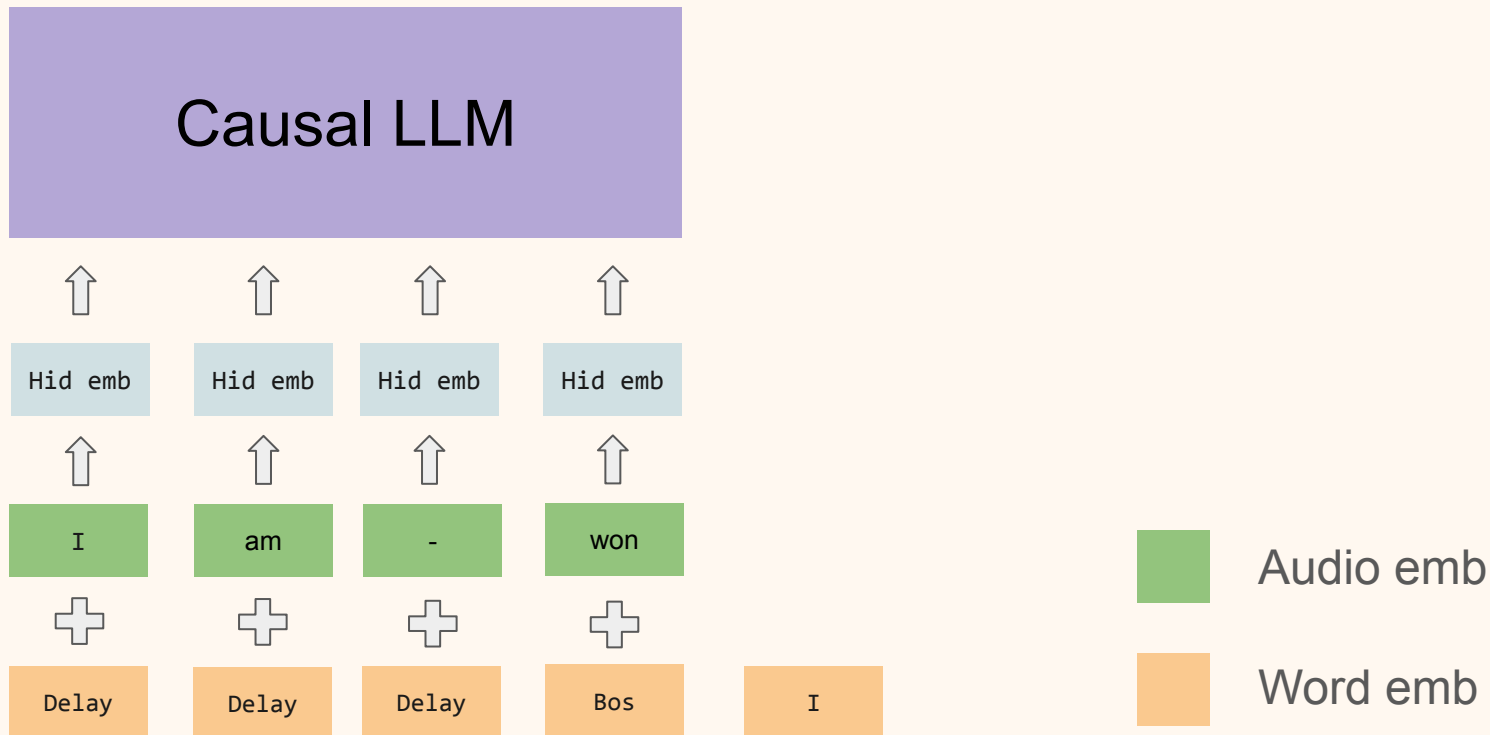
Speech Recognition - Realtime



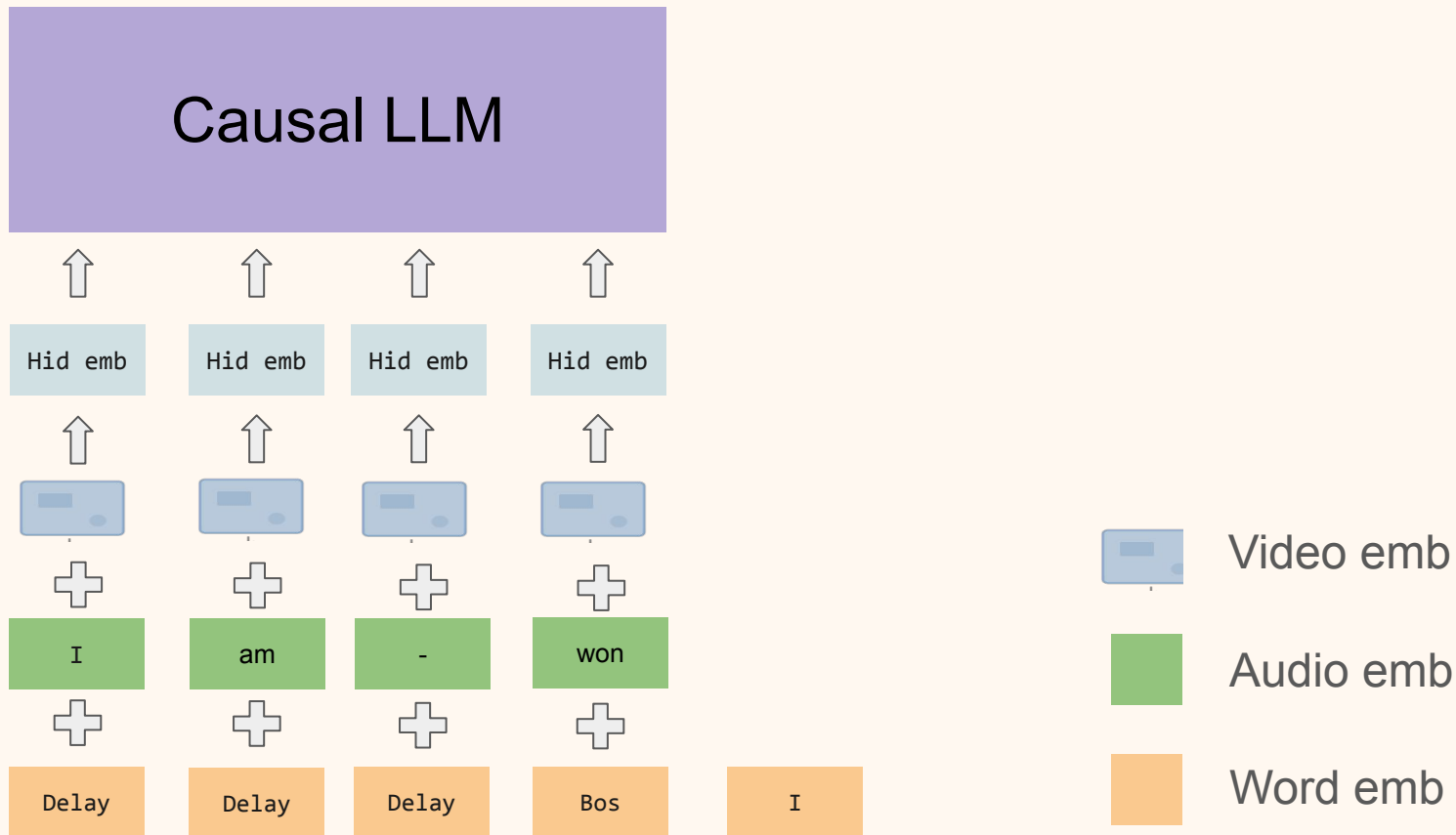
Speech Recognition - Realtime



Speech Recognition - Realtime

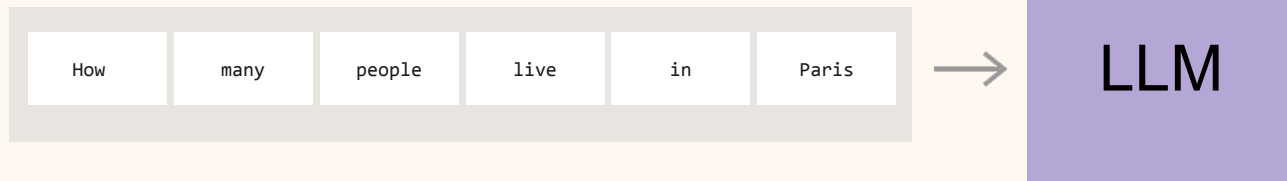


Speech & Video - Realtime

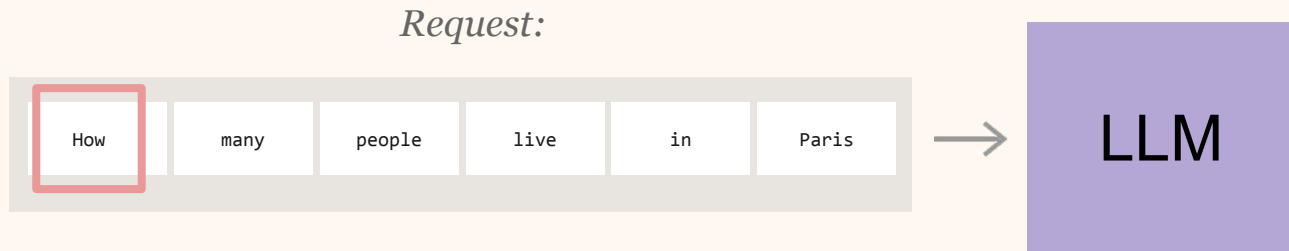


■ How LLMs currently work

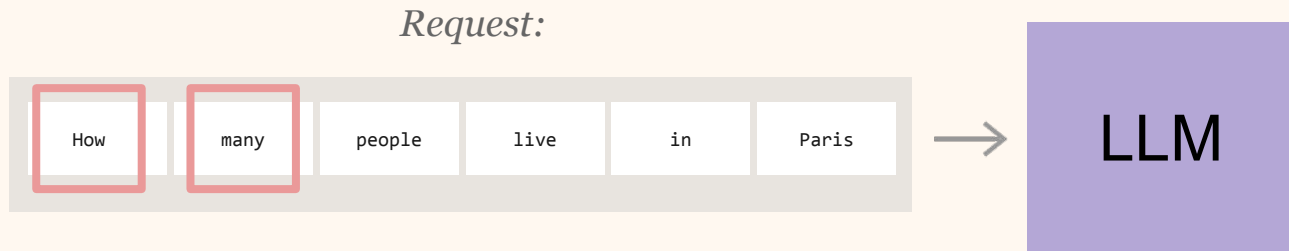
Request:



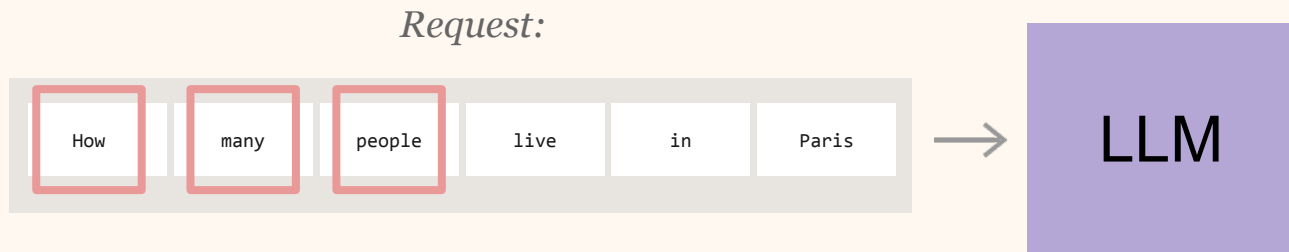
■ How LLMs currently work



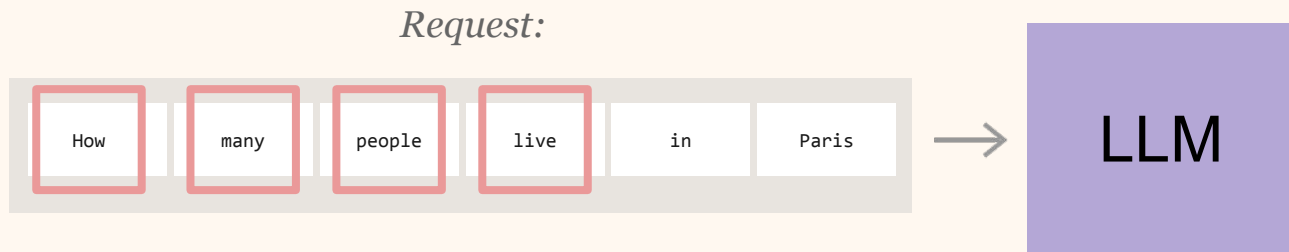
■ How LLMs currently work



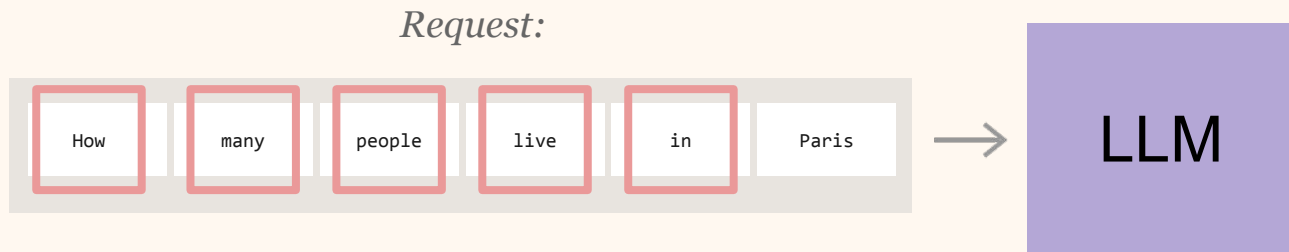
■ How LLMs currently work



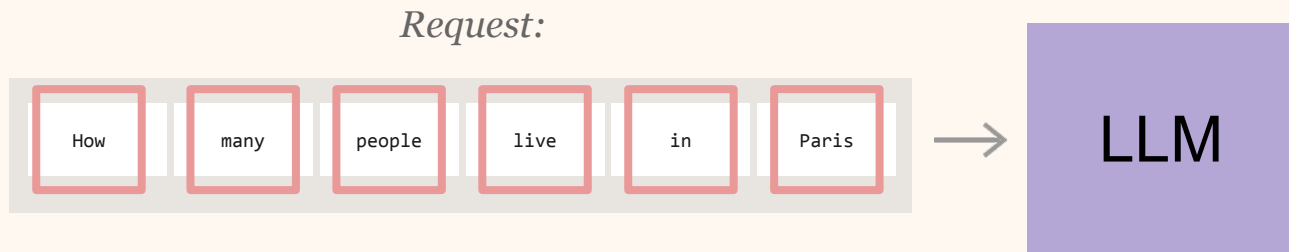
■ How LLMs currently work



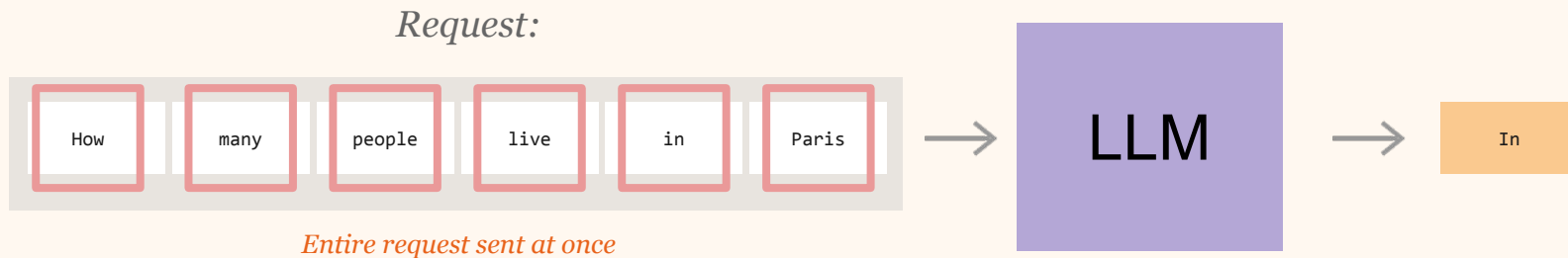
■ How LLMs currently work



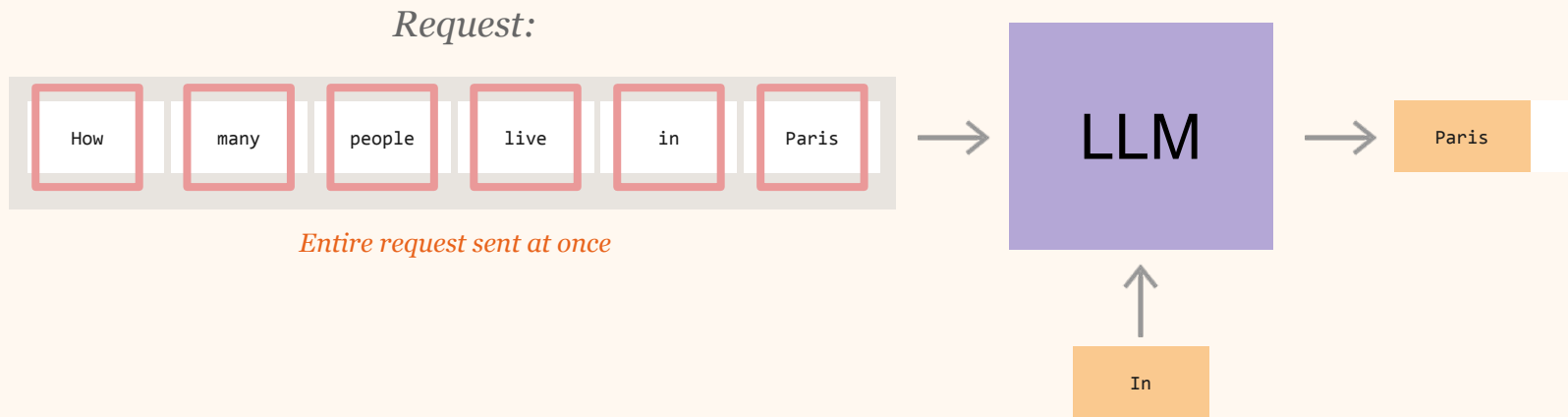
■ How LLMs currently work



■ How LLMs currently work

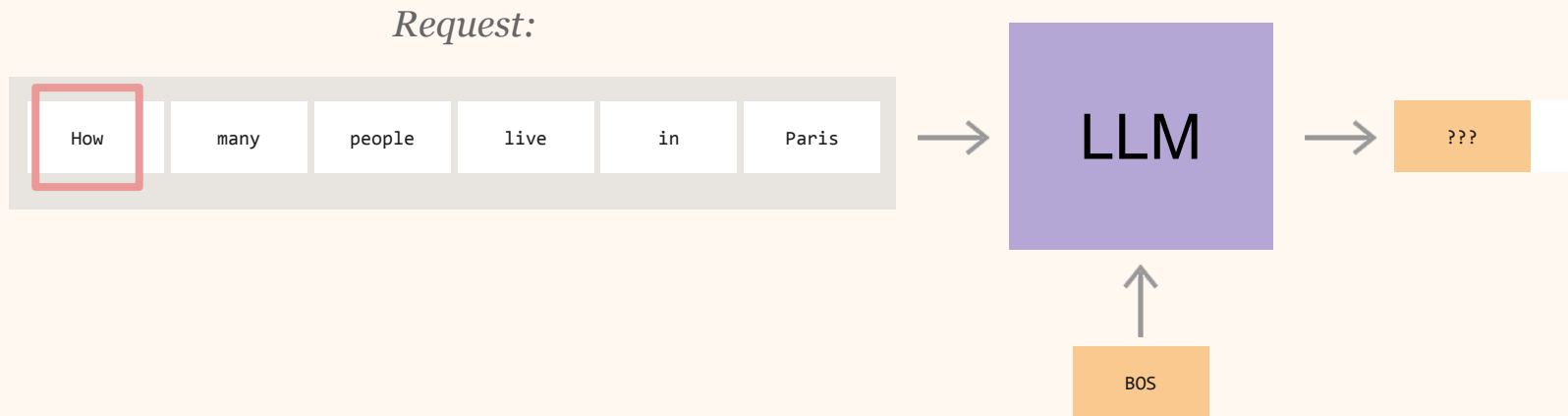


■ How LLMs currently work

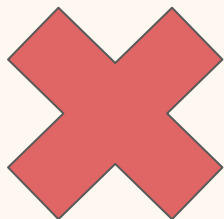
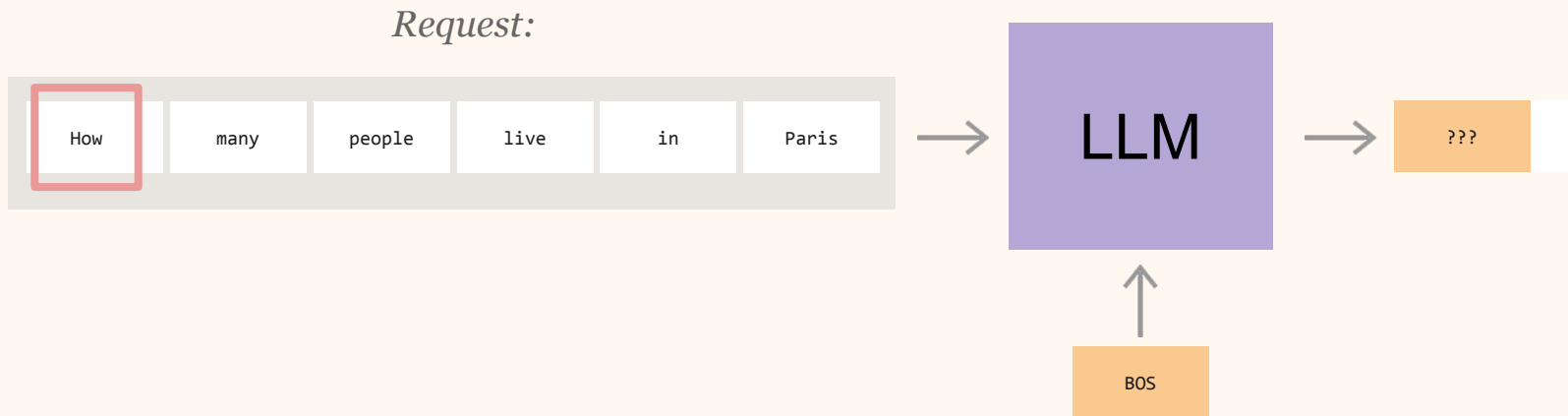


■ Can we stream LLMs?

Can we stream LLMs?



Can we stream LLMs?



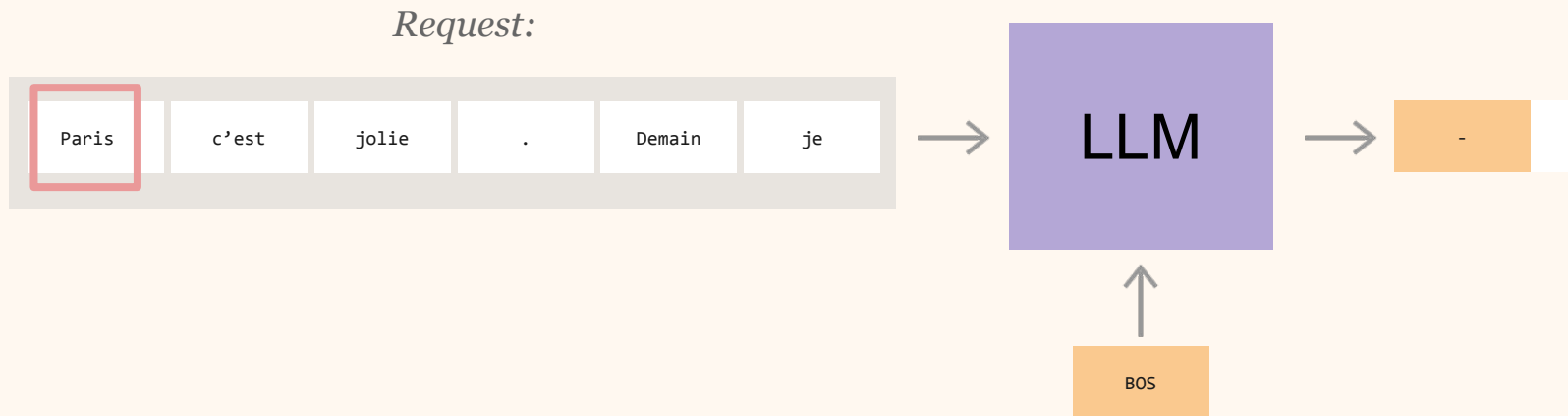
What should we predict from just “How”?

■ Live Translation

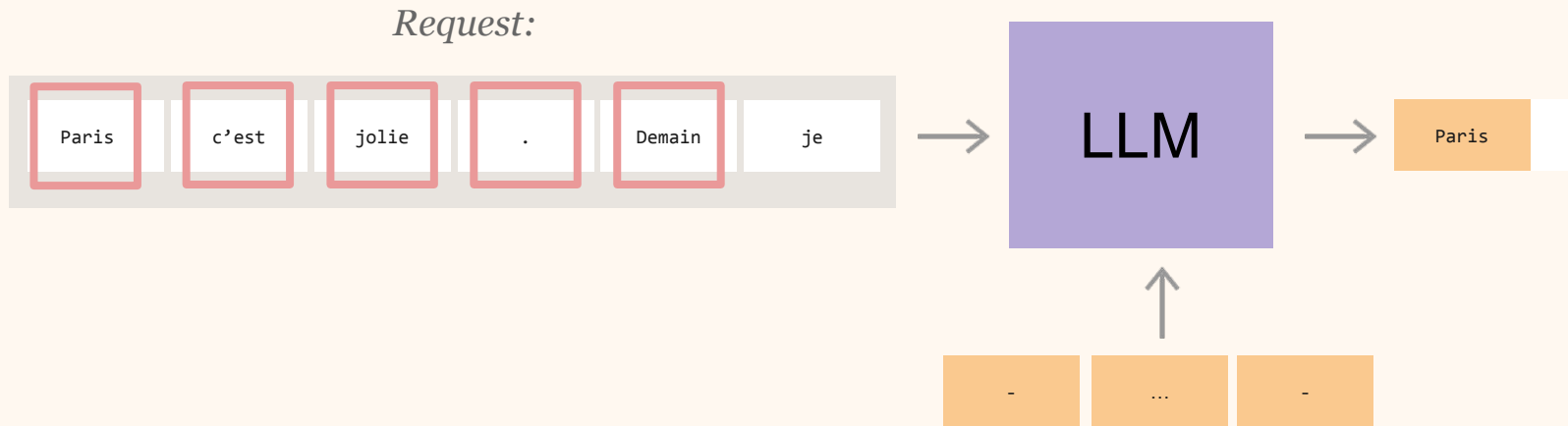
Request:



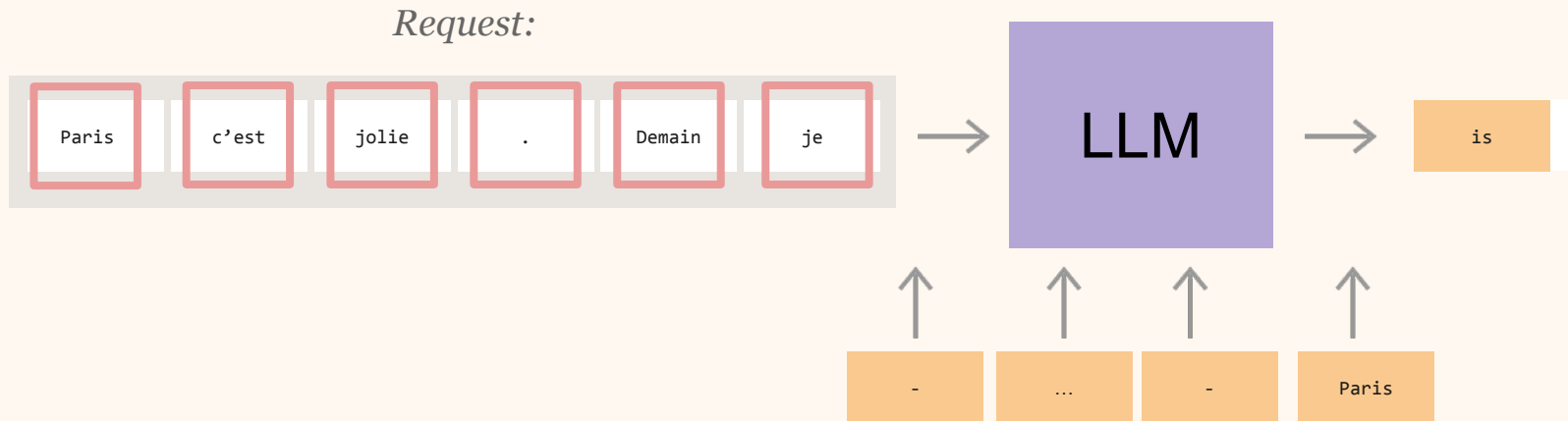
■ Live Translation



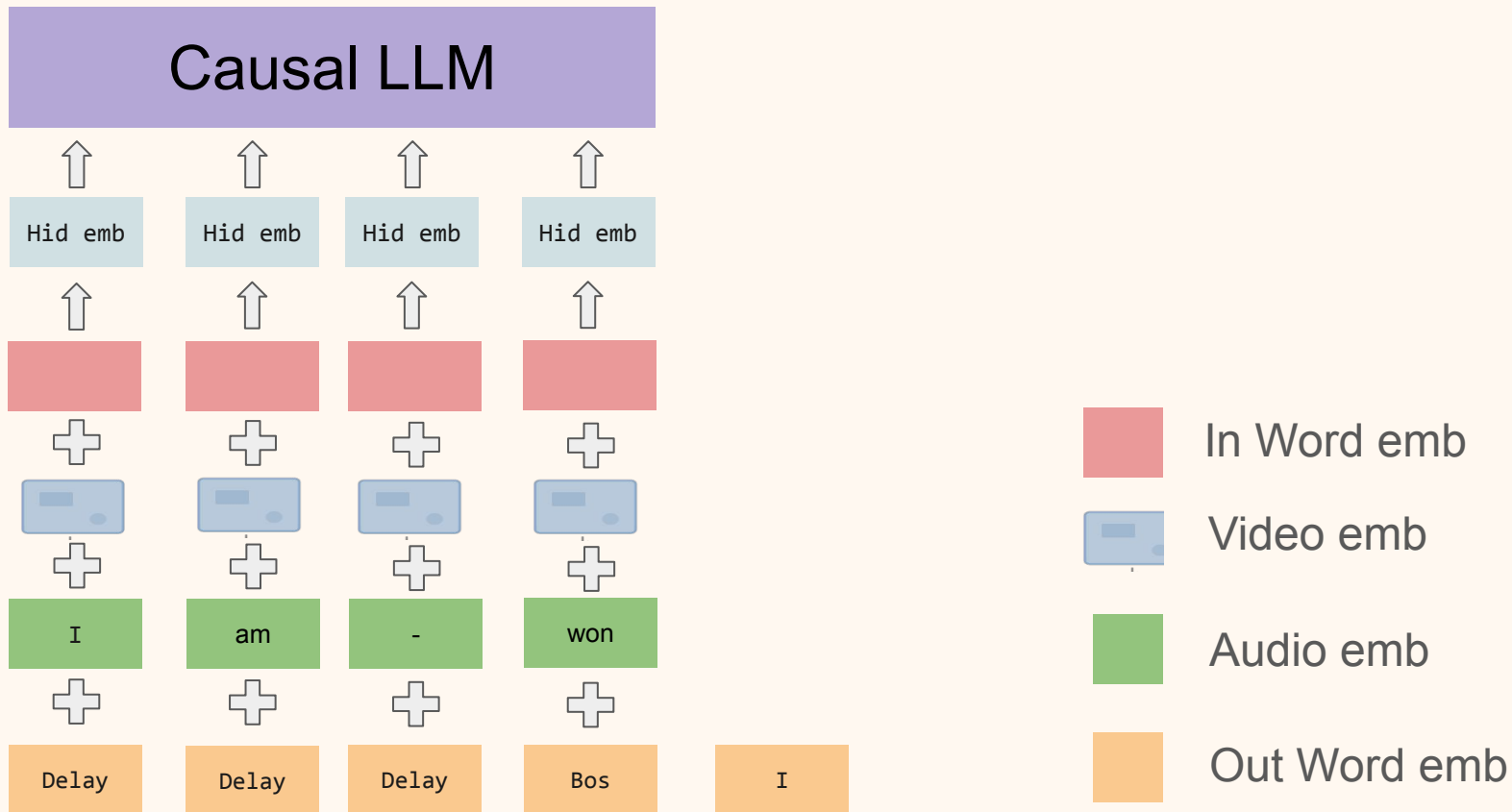
■ Live Translation



■ Live Translation



Speech & Video & Text - Realtime



■ Practical considerations

■ Practical considerations

Difficult in practice:

- How to pretrain? Need input & output pair
- Input & output pair is not necessarily aligned (log parsing)
- ...

■ Practical considerations

Difficult in practice:

- How to pretrain? Need input & output pair
- Input & output pair is not necessarily aligned (log parsing)
- ...

BUT:

- Promising approach to combine modalities
- Enables early stopping / termination
- Many industry task need “realtime” systems
- Pretrained LLMs are adaptable

■ Streaming Input for current LLMs still?

■ Streaming Input for current LLMs still?

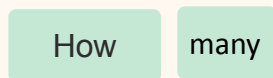
Yes, Faster prefill!

Long input has to be chunked anyways, let's start early!

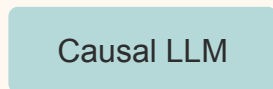
Streaming Input for faster prefilling

Input stream:

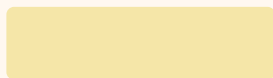
Chunk 1:



Model:

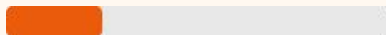


Output:



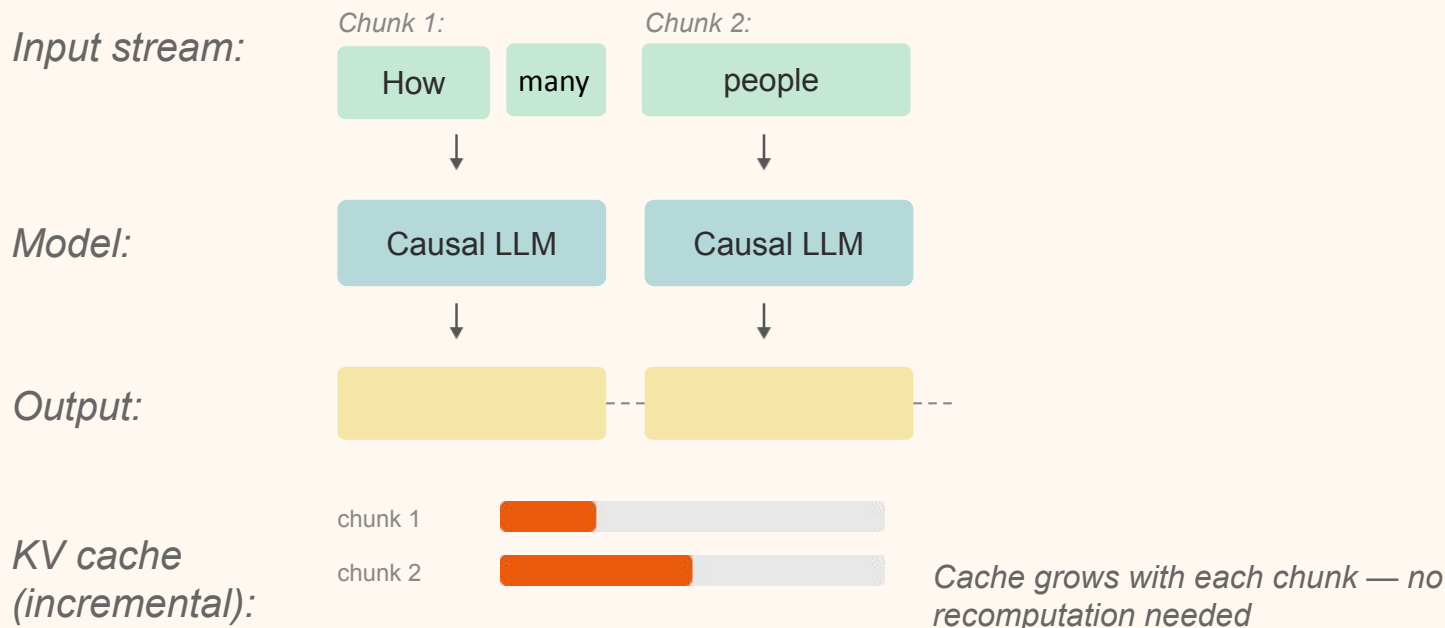
*KV cache
(incremental):*

chunk 1

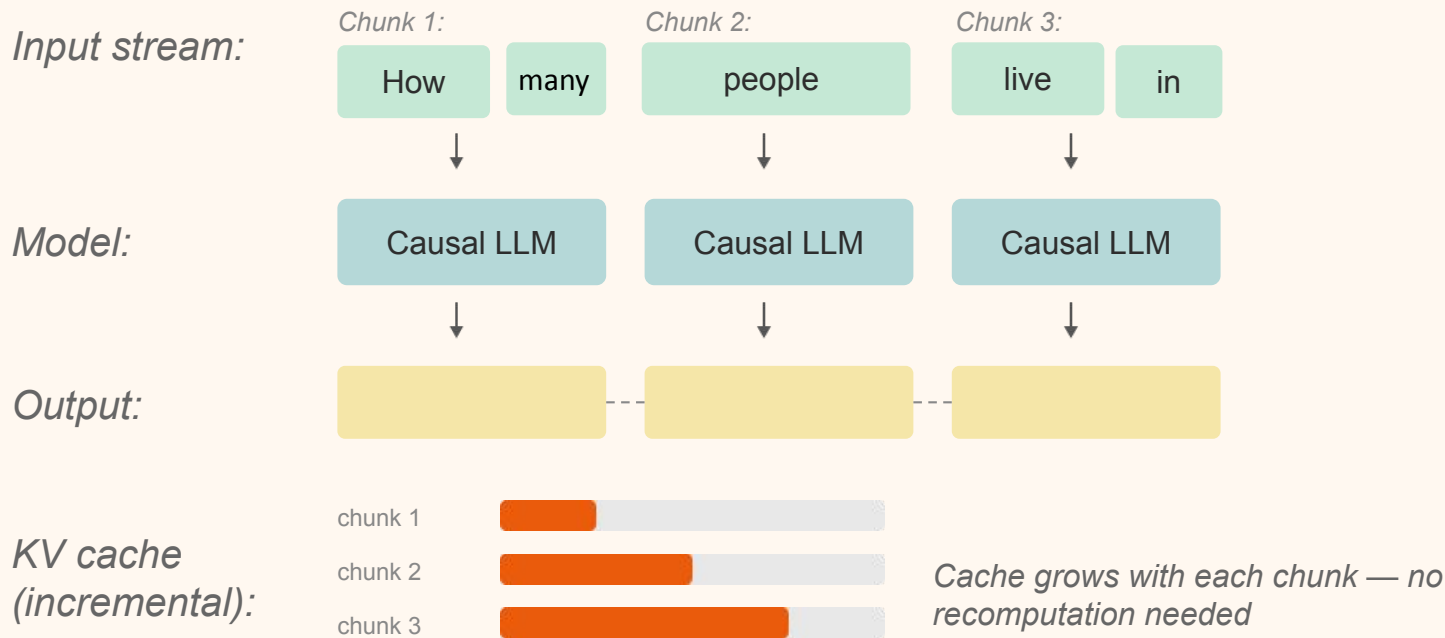


*Cache grows with each chunk — no
recomputation needed*

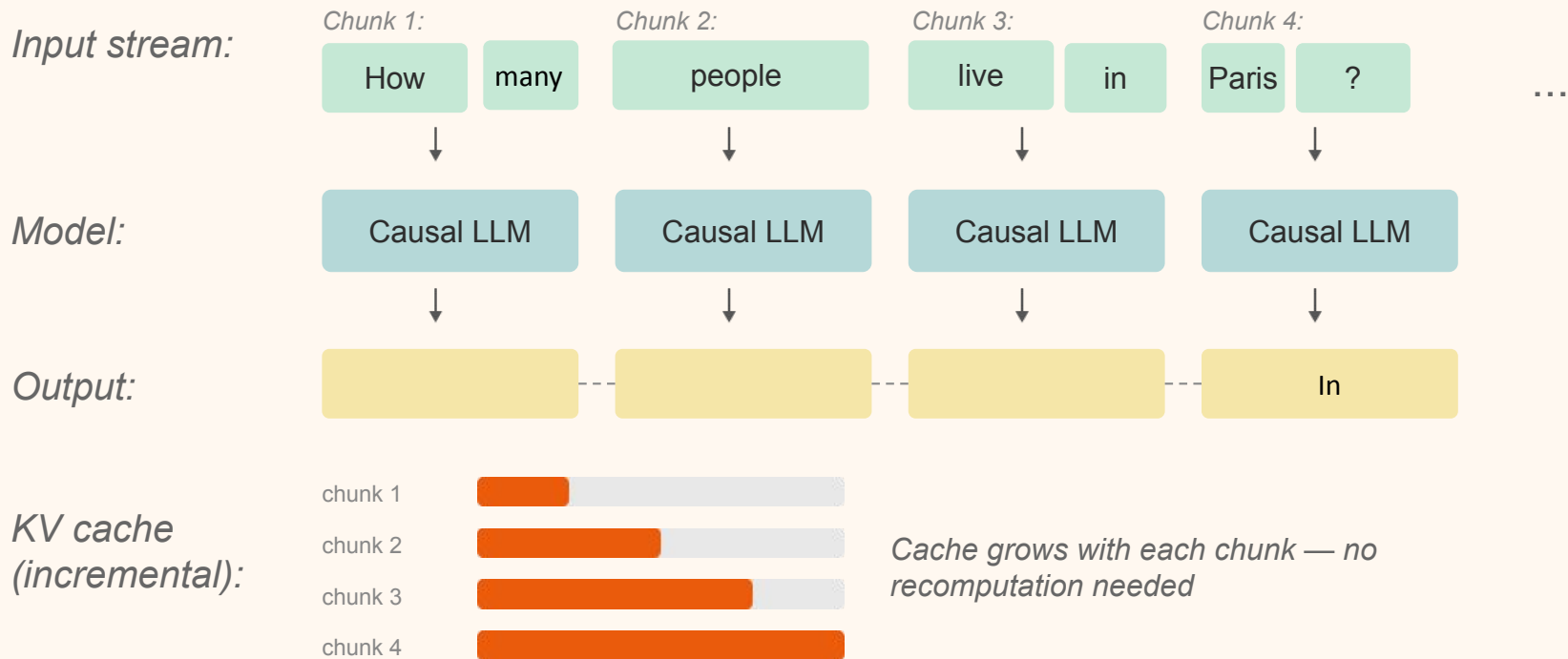
Streaming Input for faster prefilling



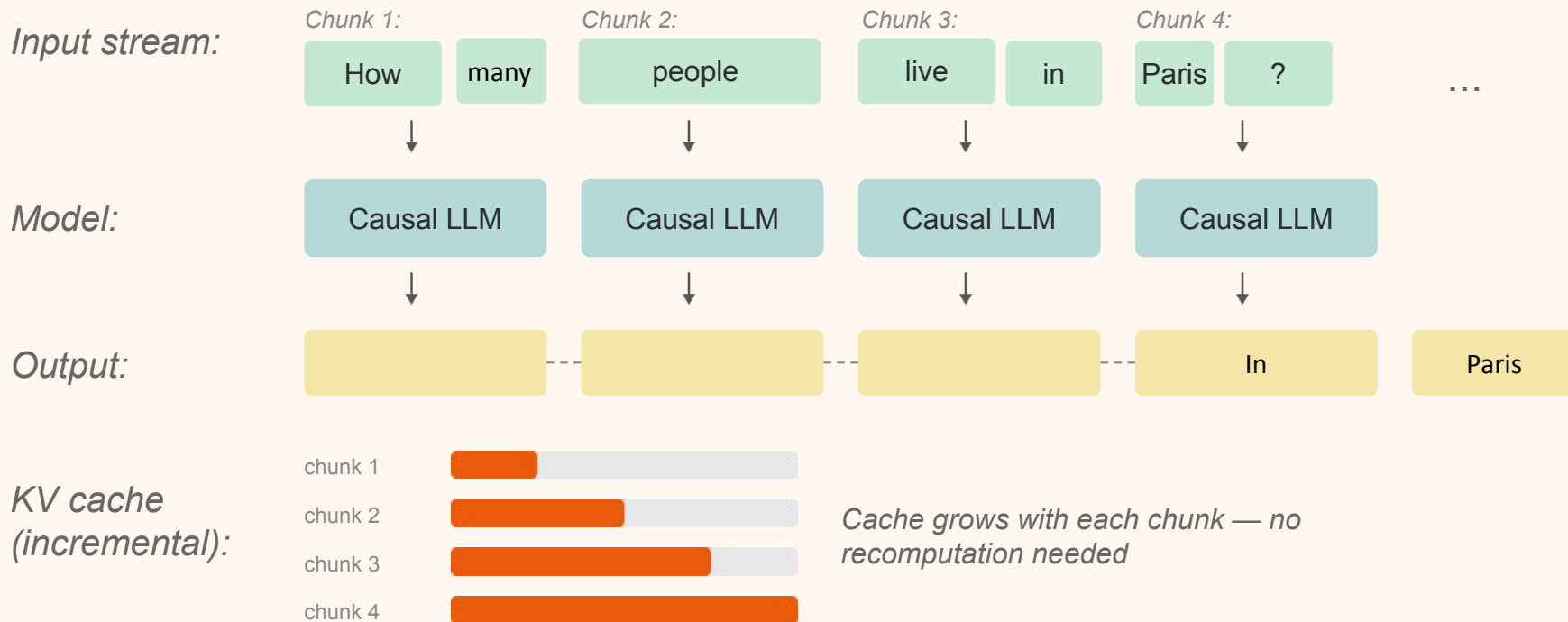
Streaming Input for faster prefilling



Streaming Input for faster prefilling



Streaming Input for faster prefilling



Support for input streaming in vLLM

- Blog:
<https://vllm-project-github-a2gx9w6hp-simon-mos-projects.vercel.app/2026/01/31/streaming-realtime.html>
- Docs:
https://docs.vllm.ai/en/latest/serving/openai_compatible_server/?h=realtime#realtime-api
- Voxtral:
<https://mistral.ai/news/voxtral-transcribe-2>

```
from dataclasses import dataclass
from vllm.inputs import PromptType
from vllm.sampling_params import SamplingParams

@dataclass
class StreamingInput:
    prompt: PromptType
    sampling_params: SamplingParams | None = None
```

