

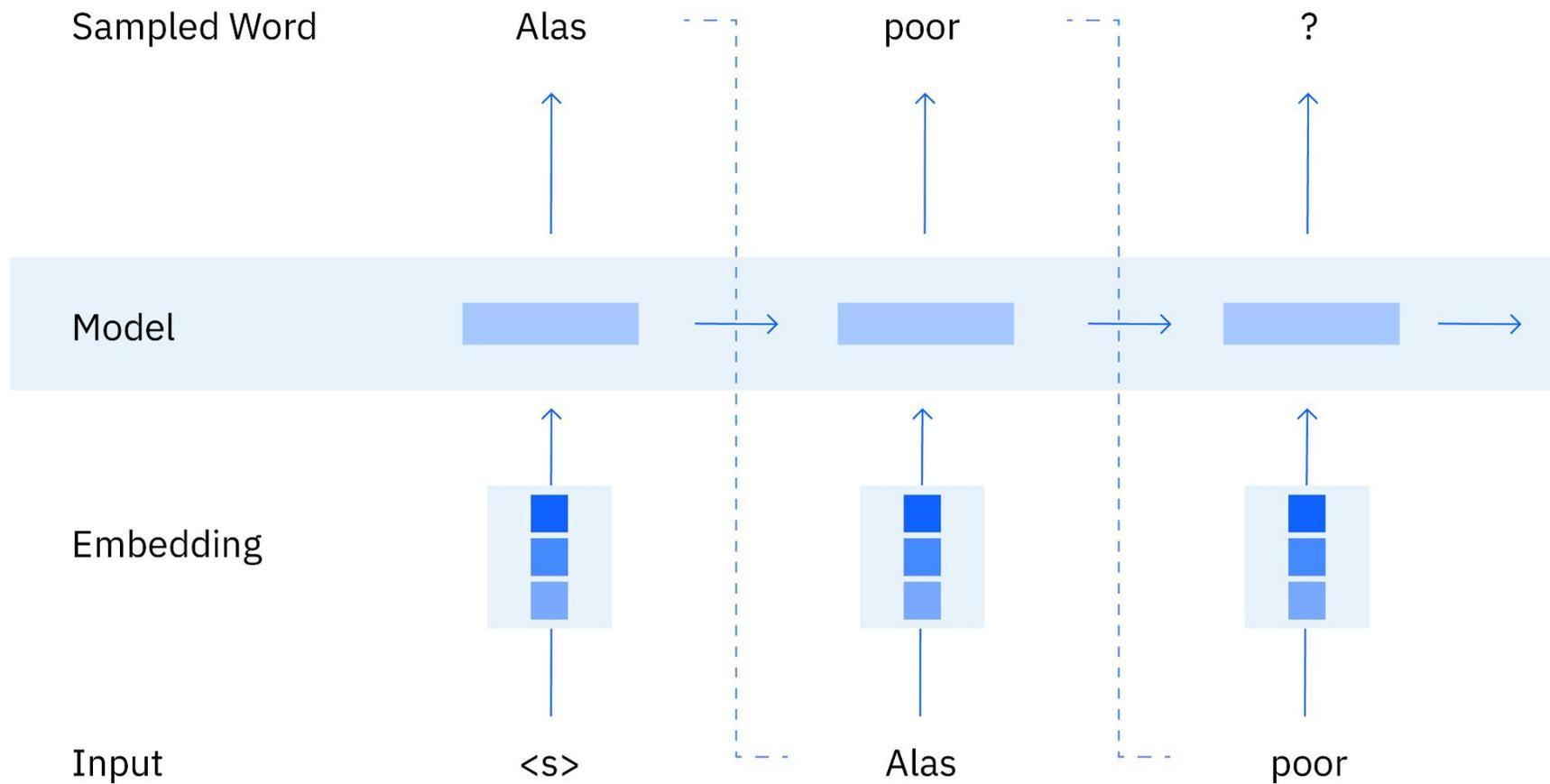


Flex Attention

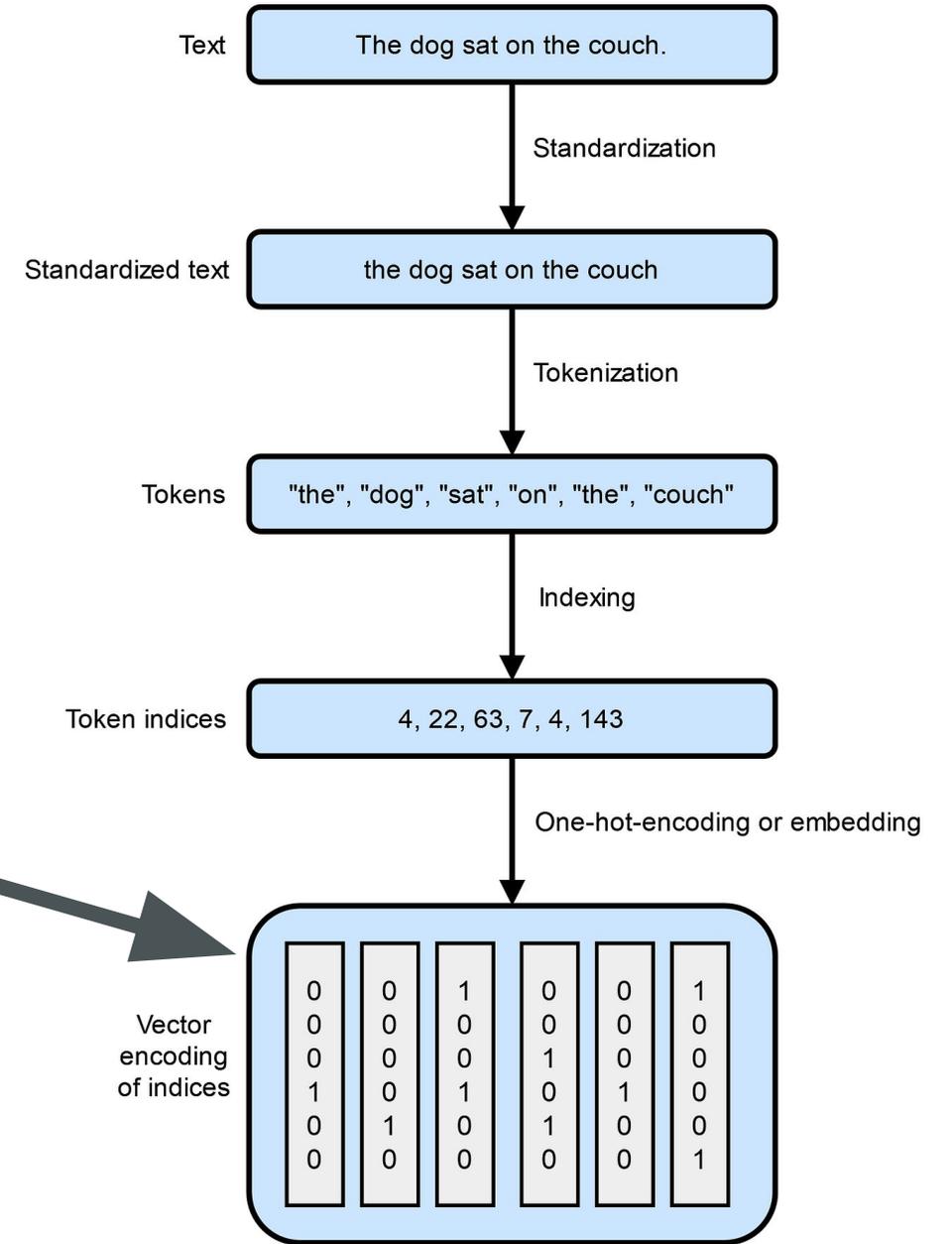
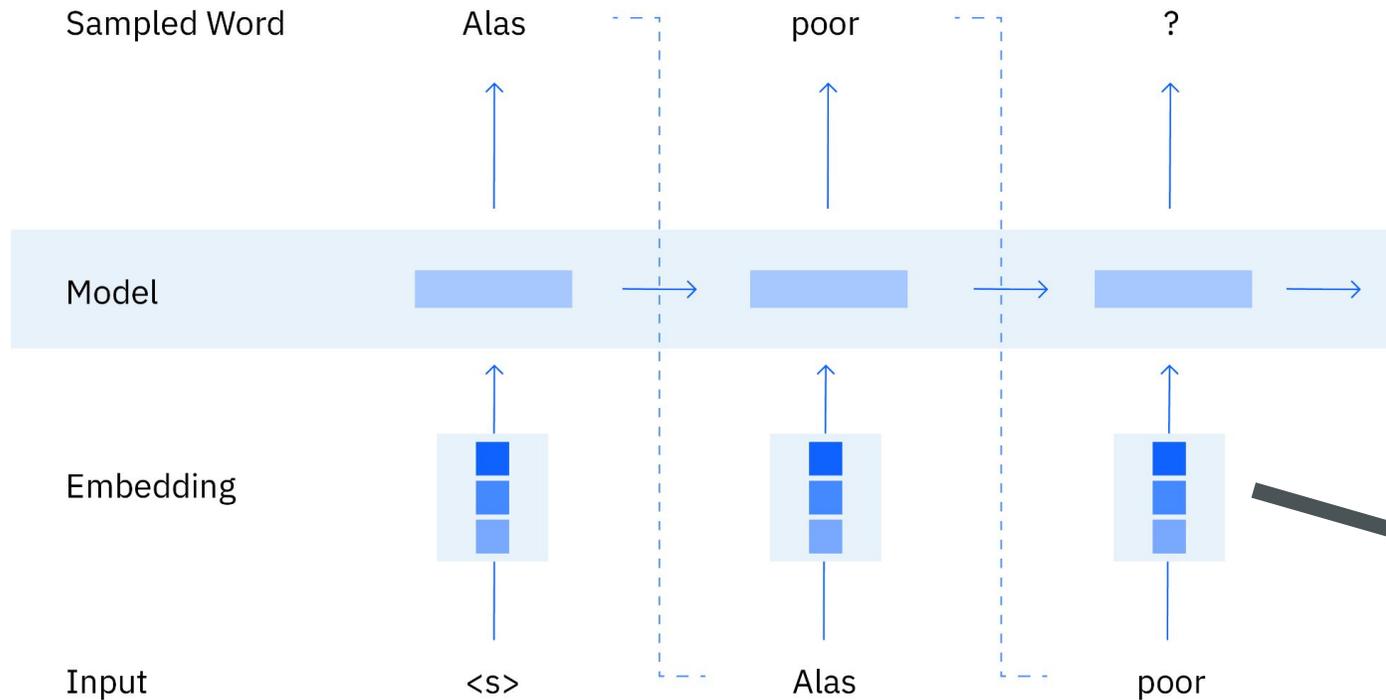
Automatic Optimised Attention Kernels

By: Bhavye, Satyam, Zain, Charlie, Kushal

Intro to Language Models



Intro to Language Models



Current Day Language Models



Gemini



Claude



**Transformer
Architecture**

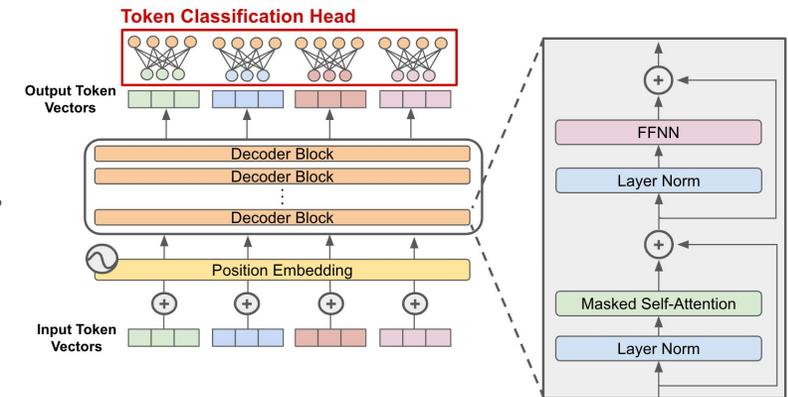
Current Day Language Models



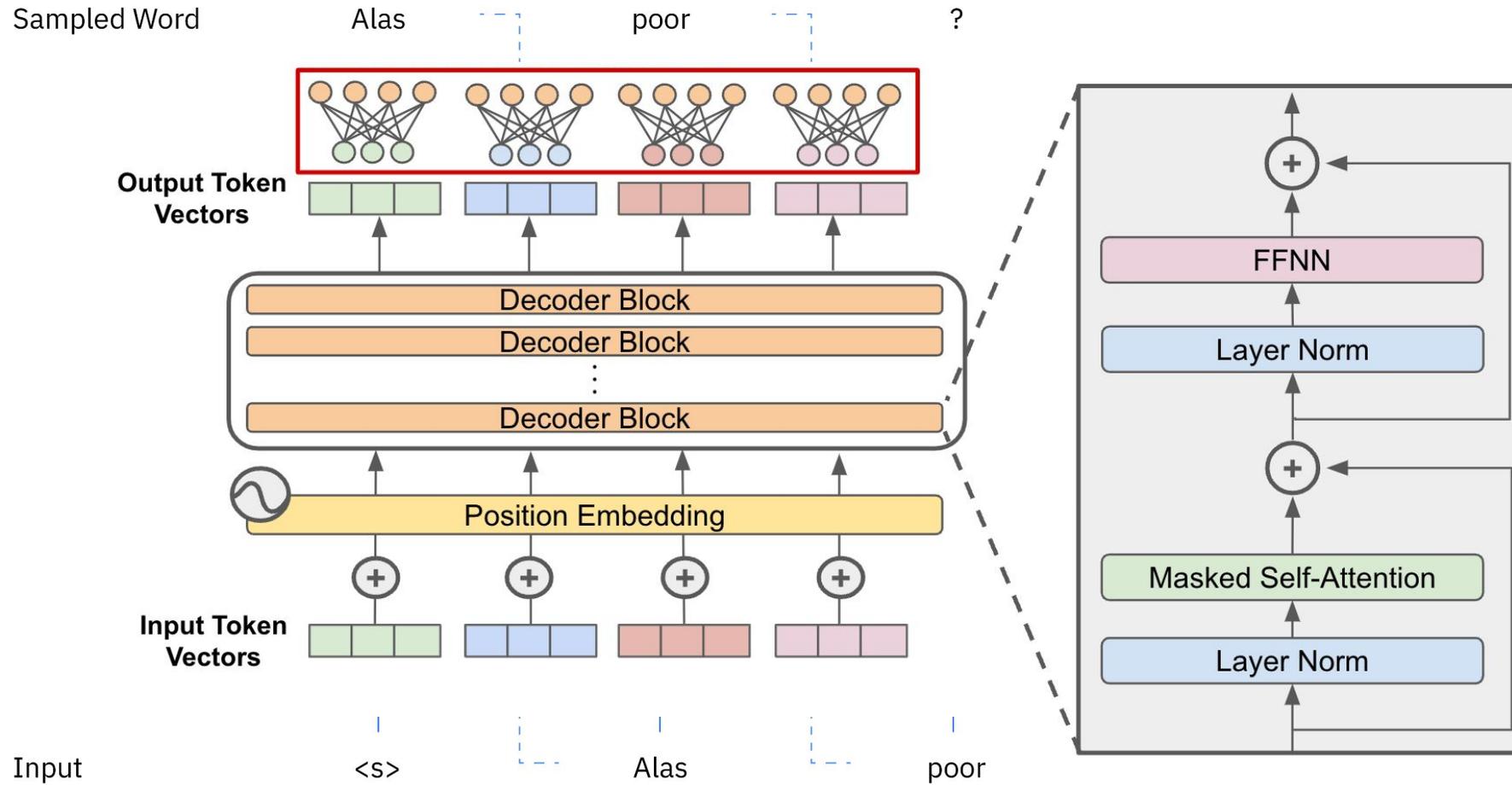
Gemini



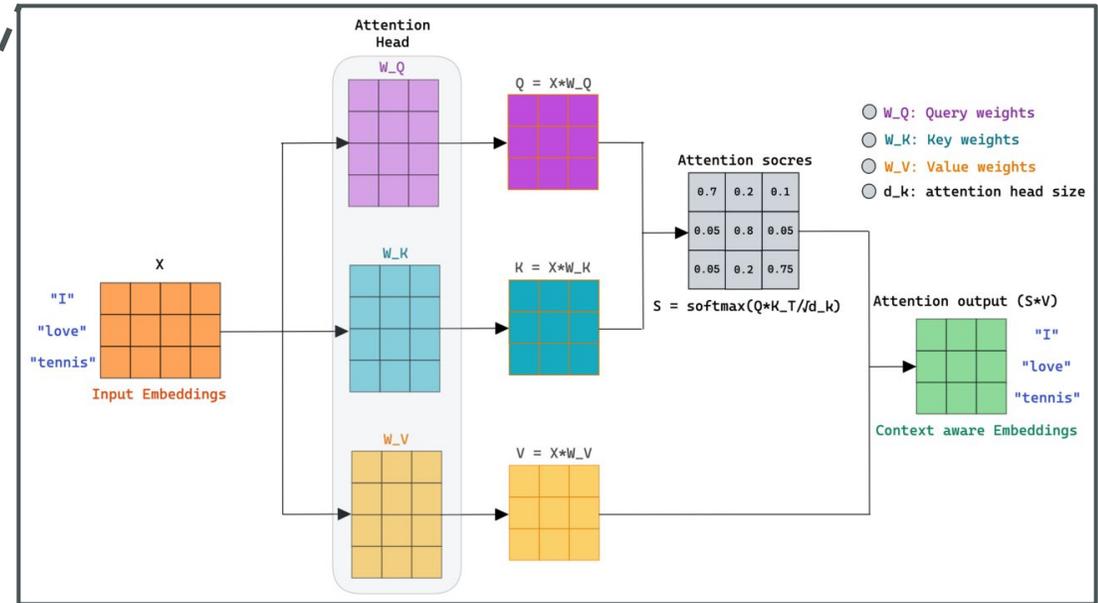
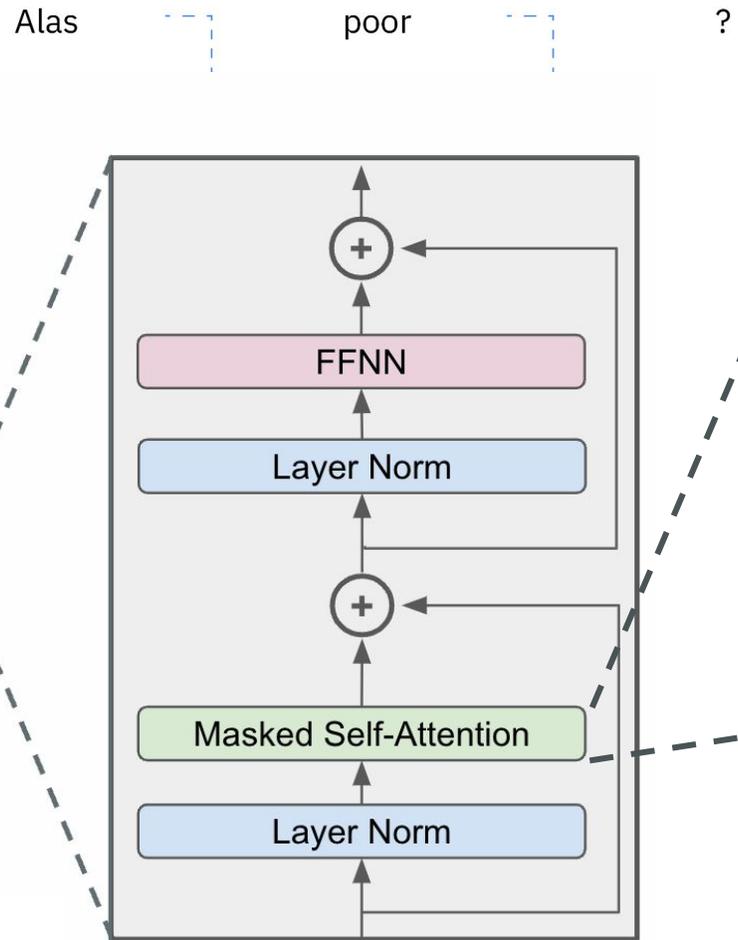
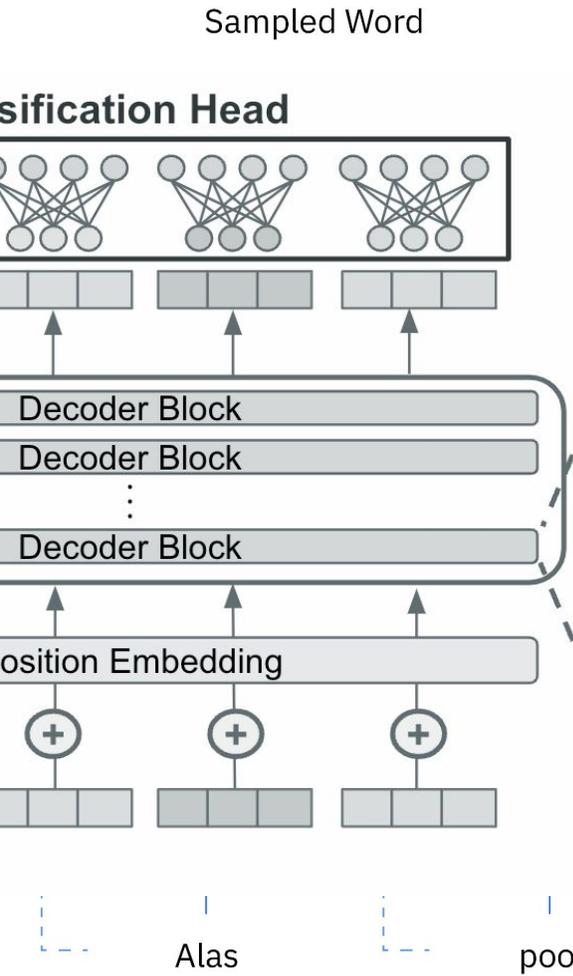
Claude



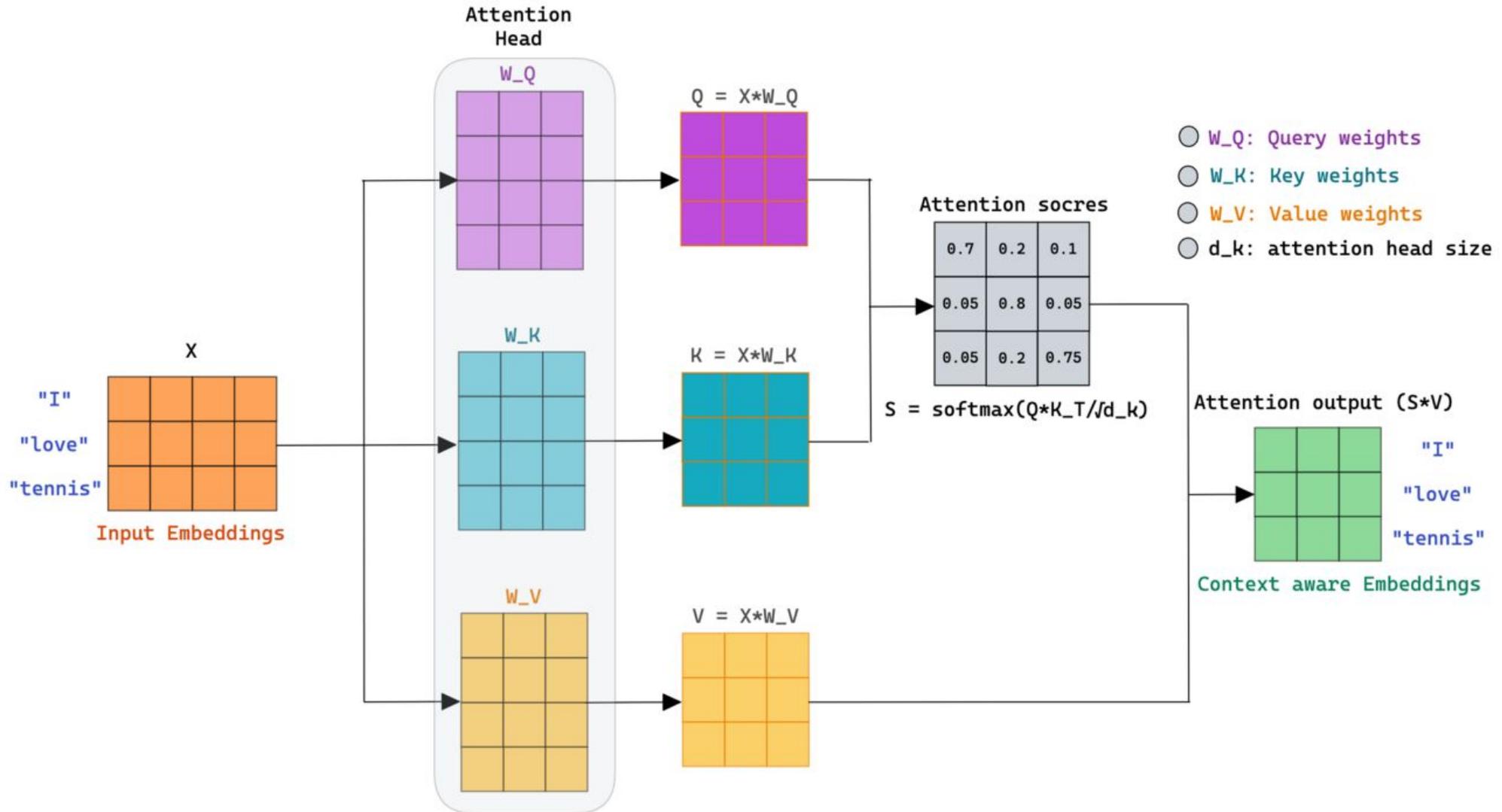
Transformer Block



Transformer Block

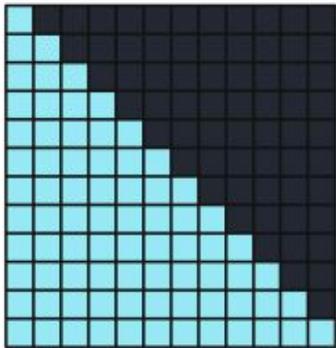


Attention

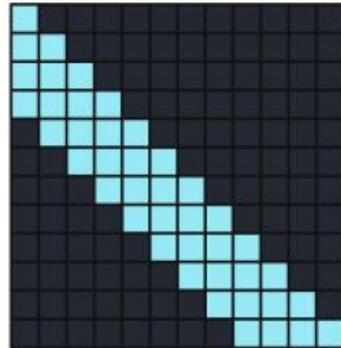


Attention Variants

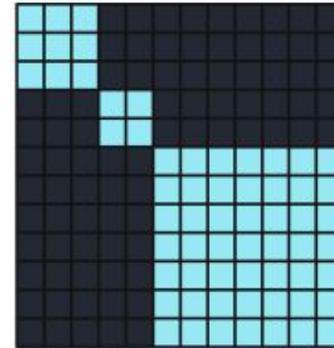
- Masks are applied to the attention matrix.



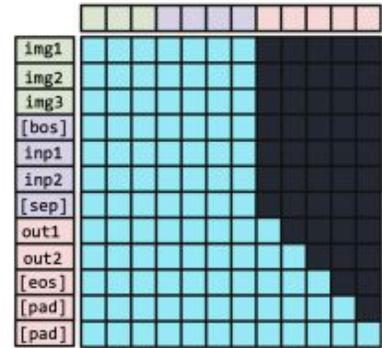
```
def causal_mask(b, h, q, kv):  
    return q >= kv
```



```
def sliding_window_mask(b, h, q, kv):  
    return q <= kv + window and q >= kv
```



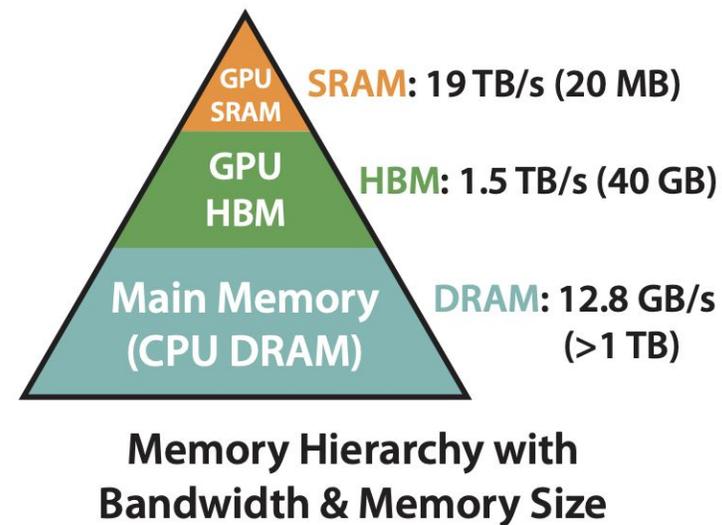
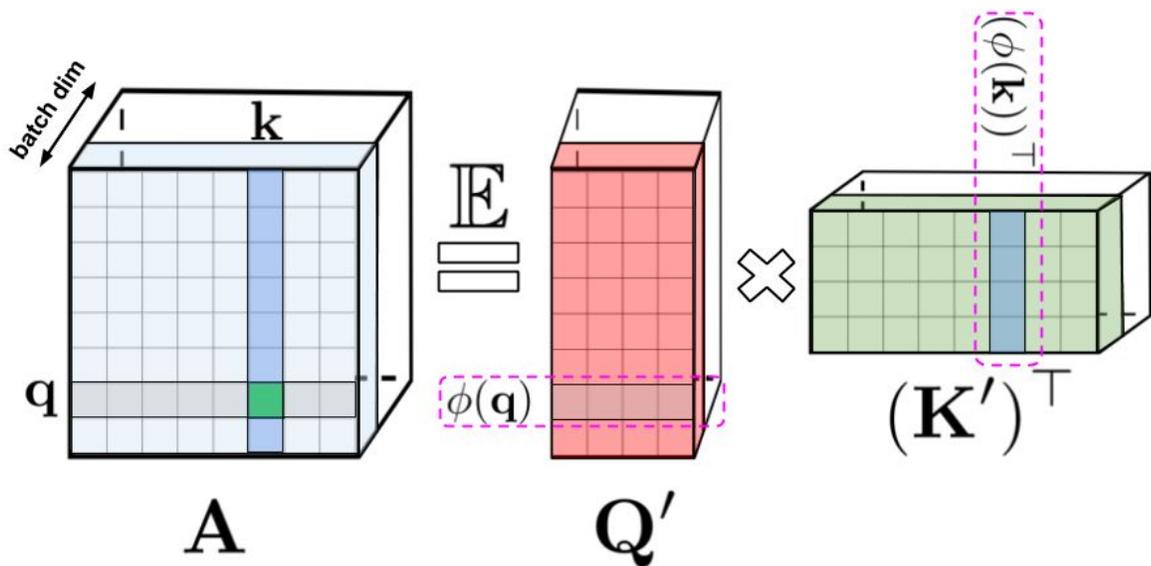
```
def document_mask(b, h, q, kv):  
    return doc_id[q] == doc_id[kv]
```



```
def prefix(b, h, q_idx, kv_idx):  
    return kv_idx <= prefix_length[b]  
prefixLM = or_mask(prefix, causal_mask)
```

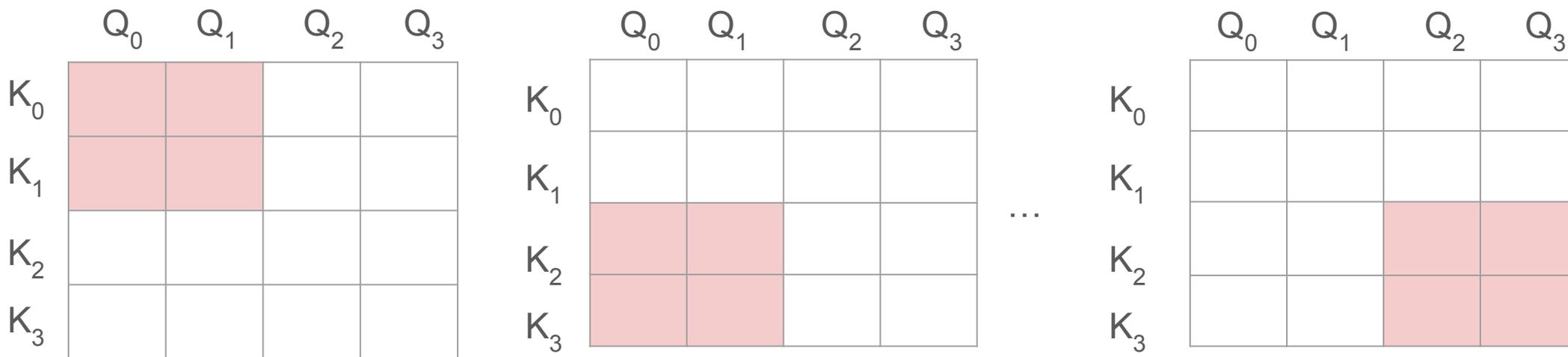
Flash Attention: Context

Memory is main bottleneck in attention computation



Flash Attention: Solution

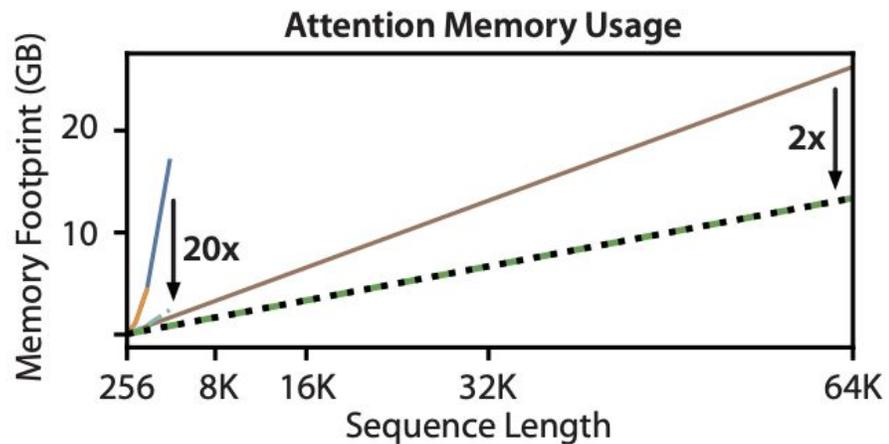
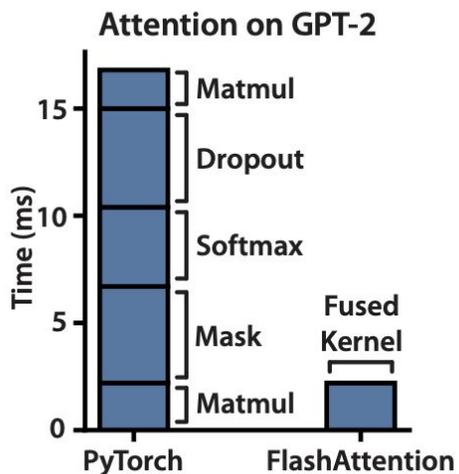
Key: Avoid HBM access by tiling



$$\begin{aligned} \text{Col_Max}_i &= \max(Q_i K_j) \\ \text{Col_Sum_Exp}_i &= \sum(\exp(Q_i K_j - \text{Col_Max}_i)) \end{aligned} \implies O_{ij} = \frac{\exp(Q_i K - \text{Col_Max}_i) * V_j}{\text{Col_Sum_Exp}_i}$$

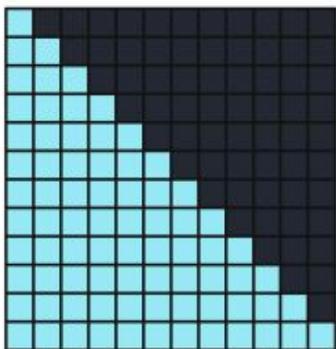
Flash Attention: Results

Noticeable speedups in runtime

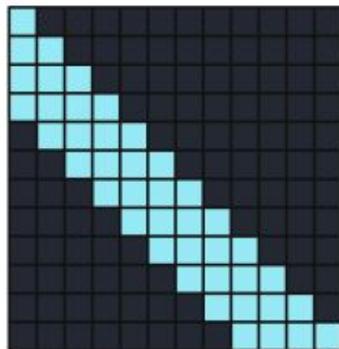


Flash Attention: Limitation

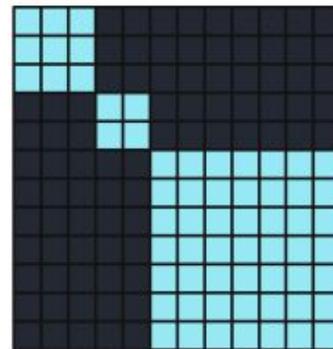
- Only supports standard, rigid attention
- If researchers need to experiment with a attention variant, they need to develop new CUDA kernels entirely from scratch



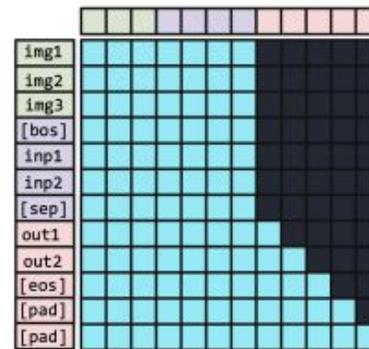
```
def causal_mask(b, h, q, kv):  
    return q >= kv
```



```
def sliding_window_mask(b, h, q, kv):  
    return q <= kv + window and q >= kv
```

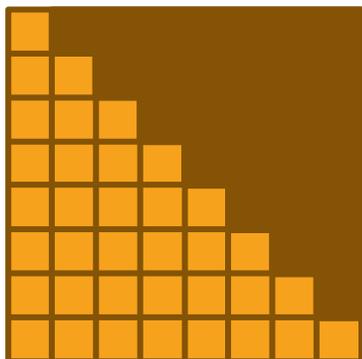


```
def document_mask(b, h, q, kv):  
    return doc_id[q] == doc_id[kv]
```

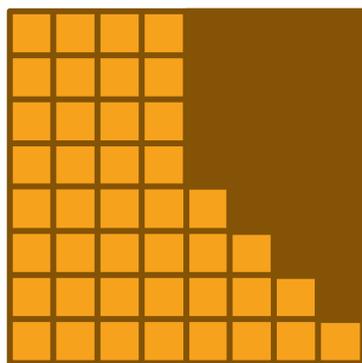


```
def prefix(b, h, q_idx, kv_idx):  
    return kv_idx <= prefix_length[b]  
prefixLM = or_mask(prefix, causal_mask)
```

Block Masks + Logical Fusions



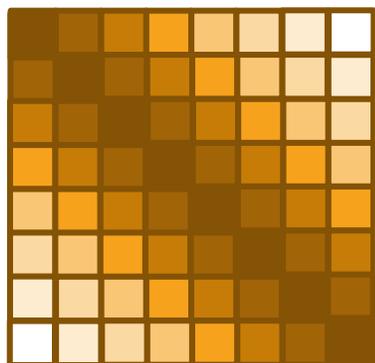
```
def causal_mask(batch: int, head: int, q_idx: int, kv_idx: int):  
    return q_idx >= kv_idx
```



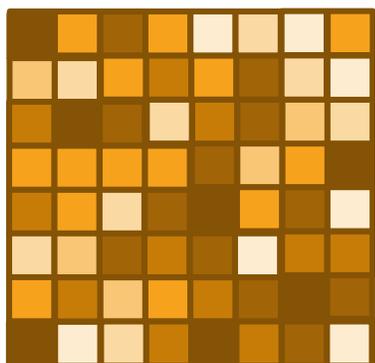
```
def prefix(batch: int, head: int, q_idx: int, kv_idx: int):  
    return kv_idx <= prefix_length[batch]
```

```
prefix_mask = or_mask(prefix, causal_mask)
```

Score Modification Functions



```
def alibi_bias(score: double, batch: int, head: int, q_idx: int, kv_idx: int):  
    return score + alibi_bias[head] * (q_idx - kv_idx)
```

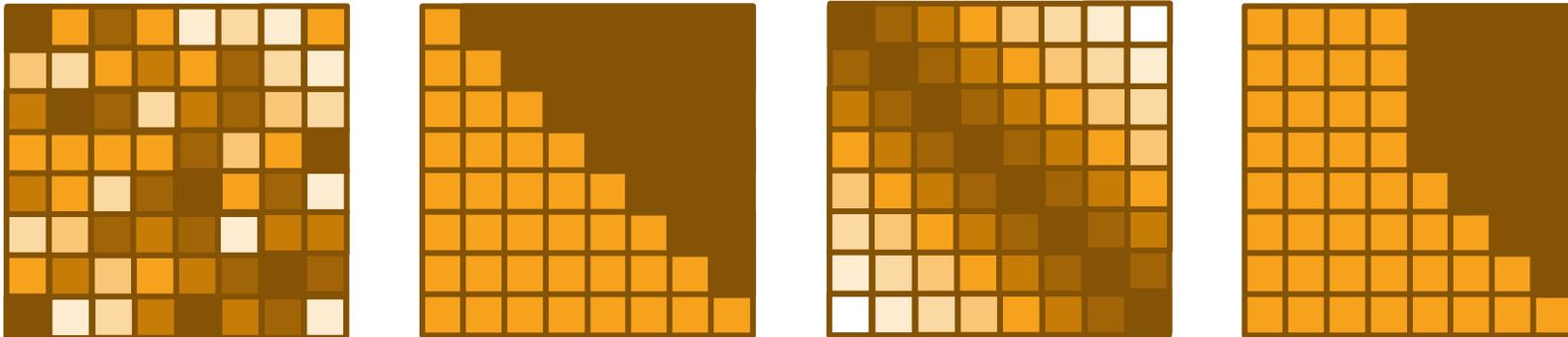


```
def tanh_mod(score: double, batch: int, head: int, q_idx: int, kv_idx: int):  
    return tanh(score)
```

Unified Attention Abstraction

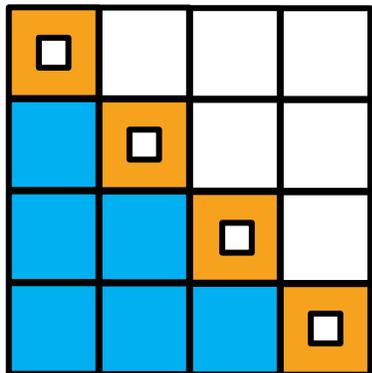
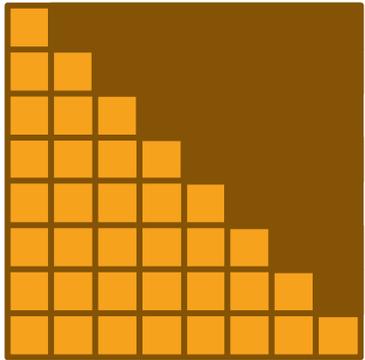
```
def mask_mod(batch: int, head: int, q_idx: int, kv_idx: int)  
    → bool
```

```
def score_mod(score: double, batch: int, head: int, q_idx: int, kv_idx: int)  
    → double
```



Why Mask Modifications?

“attention variants usually show high sparsity such as 50%”



full blocks → ignore mask

partial blocks → apply mask element-wise

empty blocks → skip computation

autotuned + configurable block sizes

Flex Attention Overview

Implemented new attention variant via 2 simple functions

Abstracted into

- Mask Mod: specifies if scalar value should be masked out
- Score Mod: applies arb., fine-grained update to score scalars

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

$$\text{FlexAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\text{mod}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\right)\mathbf{V}$$

Mod Signatures

```
def mask_mod(batch_idx: int, head_idx: int,  
             q_idx: int, kv_idx: int) -> bool
```

```
def score_mod(score: T, batch_idx: int,  
             head_idx: int, q_idx: int, kv_idx: int)  
             -> T
```

Code Compilation

```
User defined <mask_mod>  
def causal_mask(b, h, q_idx, kv_idx):  
    return q_idx >= kv_idx
```

```
User defined <mask_mod>  
def rel_pos(score, b, h, q_idx, kv_idx):  
    return score + (q_idx - kv_idx)
```

↓ torch.compile

```
@triton.jit
```

```
mask = tl.where(q_idx - kv_idx, 1, 0)
```

```
@triton.jit
```

```
score = score + (q_idx - kv_idx)
```

```
@triton.jit
```

Attention Template

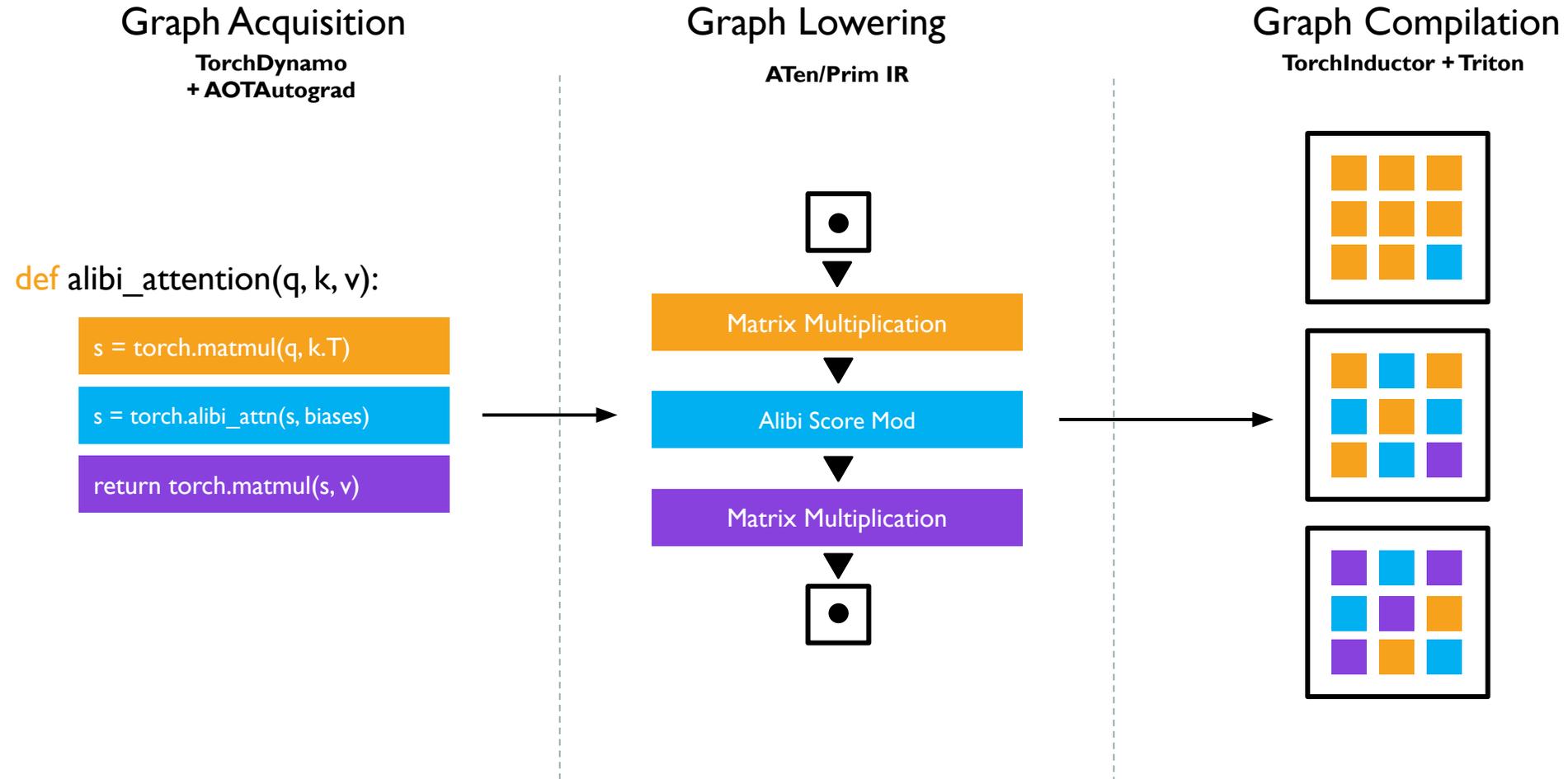
```
...  
score = q @ k  
q_idx = tl.arange(BLOCK_M) + q_offset  
kv_idx = tl.arange(BLOCK_N) + kv_offset  
q_idx = q_idx[:, None]  
kv_idx = kv_idx[None, :]
```

```
mask = tl.where(q_idx - kv_idx, 1, 0)  
score = tl.where(mask, score, float("-inf"))
```

```
score = score + (q_idx - kv_idx)
```

```
out = score @ v  
...
```

Automatic Kernel Compilation



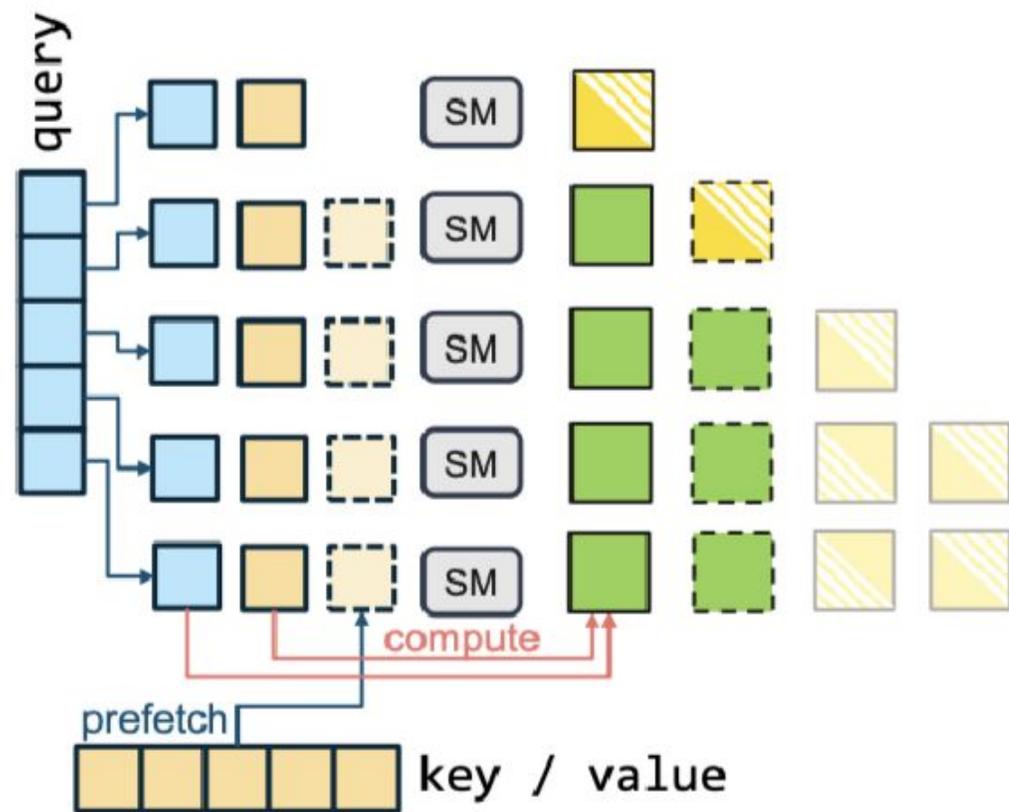
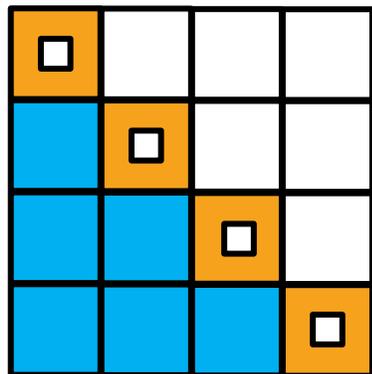
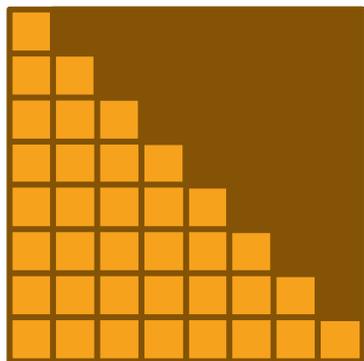


Figure 4. Scheduling full and partial blocks to SM.

Why Mask Modifications?

“attention variants usually show high sparsity such as 50%”



full blocks → ignore mask

partial blocks → apply mask element-wise

empty blocks → skip computation

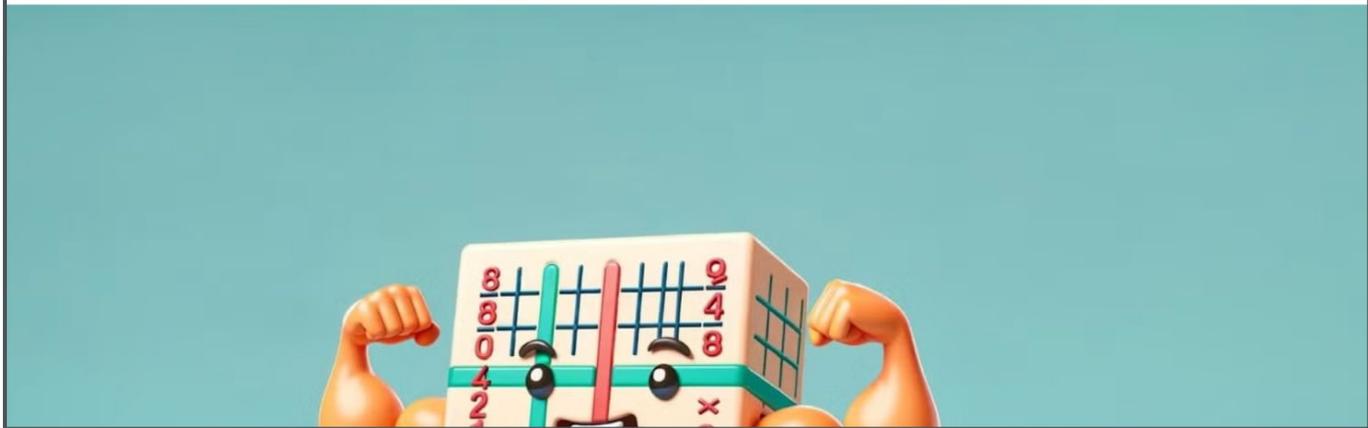
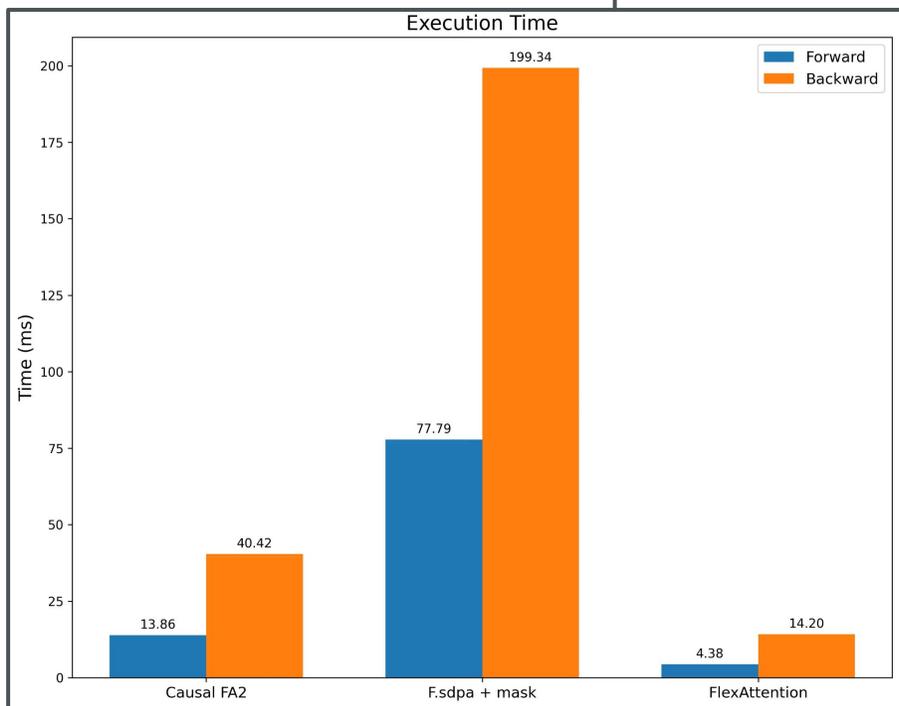
autotuned + configurable block sizes

Results

Blog

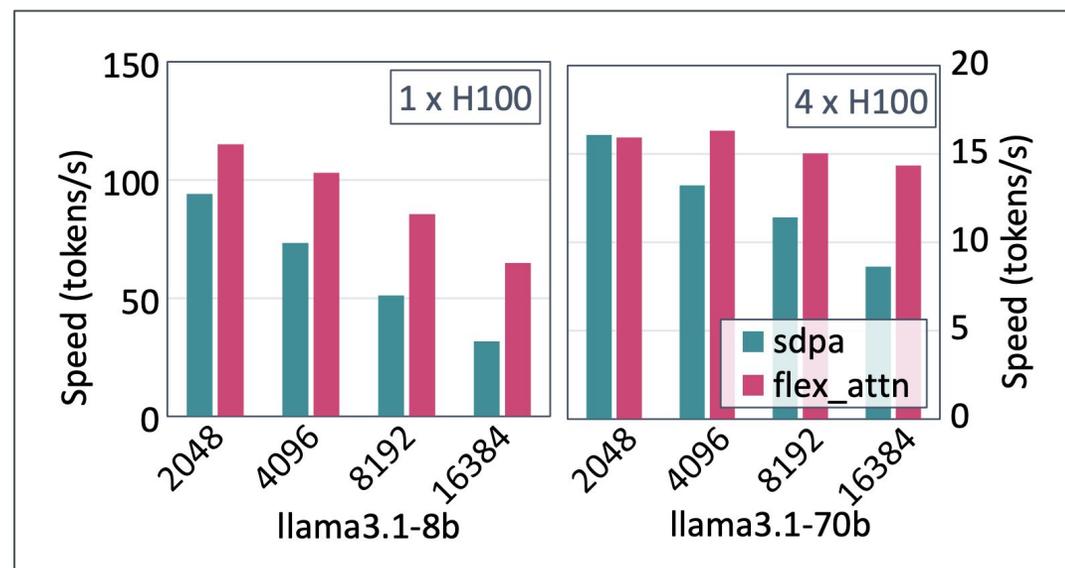
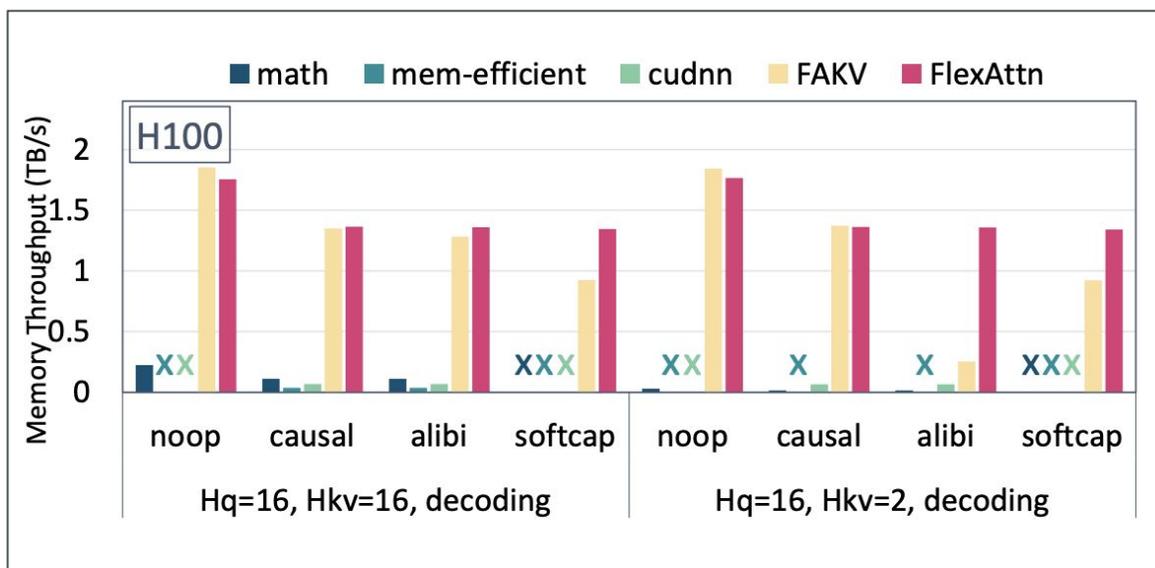
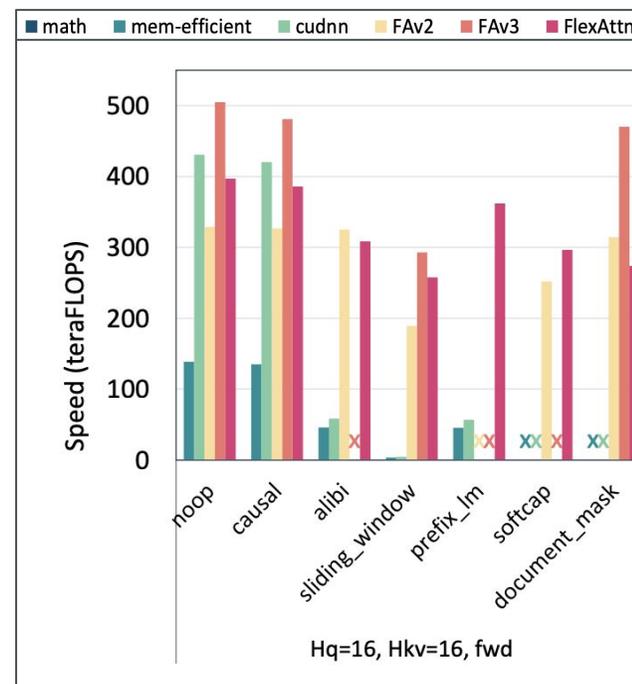
FlexAttention: The Flexibility of PyTorch with the Performance of FlashAttention

By Team PyTorch: Driss Guessous, Yanbo Liang, Joy Dong, Horace He | August 7, 2024



Results

- FlexAttention is an improvement on FlashAttention that enables more attention variants and often turns out to be more efficient than eager implementations.



Limitations / Future Scope

- FlexAttention is an improvement that allows for more attention variants than FlashAttention does, but on common variants (causal, noop, etc) it is actually slower than FlashAttention or hand-written cuda implementations.
- The paper loosely concludes that it hopes to help “explore new attention variants without being held back by what hand-written kernels support.” A future scope could include not only adapting new attention variants, but optimizing for the most common variants as well.
- Similarly, the paper could continue to show performance on more new variants and show a stronger use case for FlexAttention over FlashAttention, given widespread adoption is the goal.