# ReLU Strikes Back: Exploiting Activation Sparsity in Large Language Models

Iman Mirzadeh, Keivan Alizadeh, Sachin Mehta, Carlo C Del Mundo,

Oncel Tuzel, Golnoosh Samei, Mohammad Rastegari, Mehrdad Farajtabar
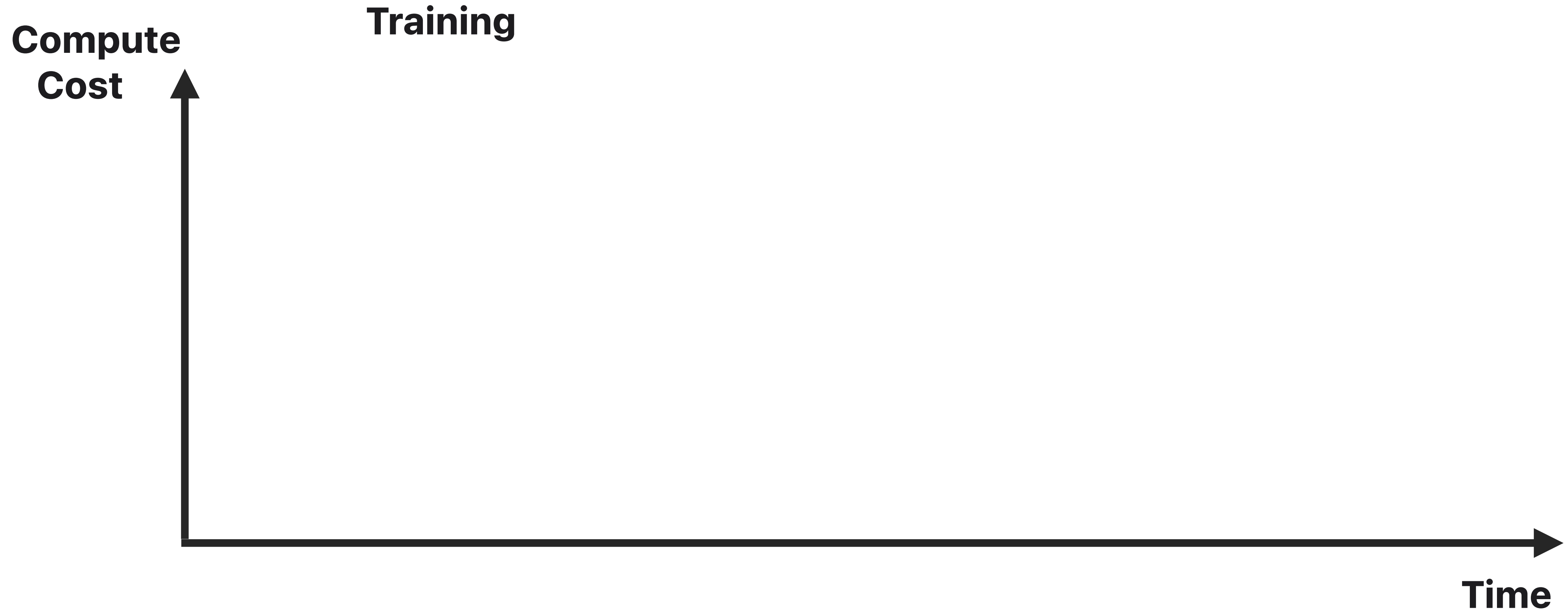
1

# Motivation

Training vs Inference
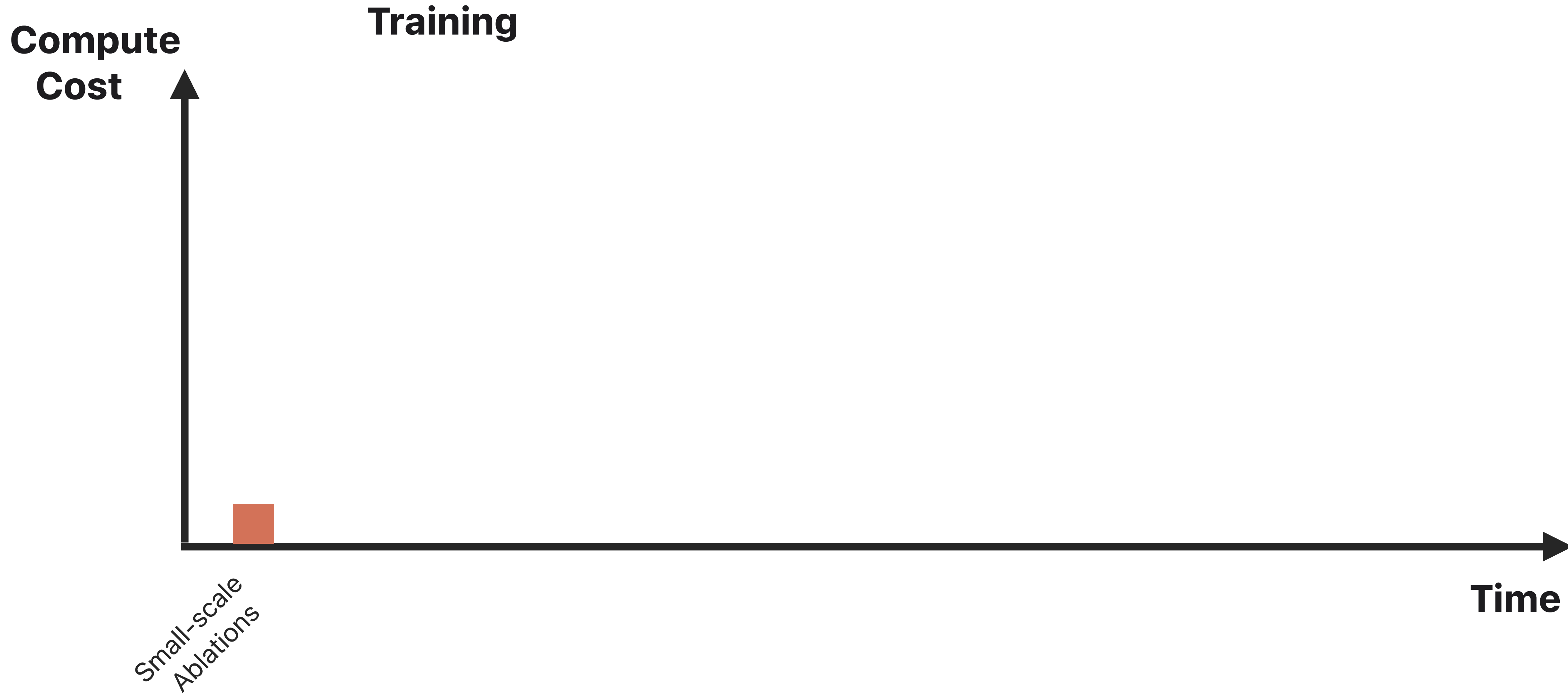
**Compute Cost**

**Time**

# Motivation
Training vs Inference

**Training**

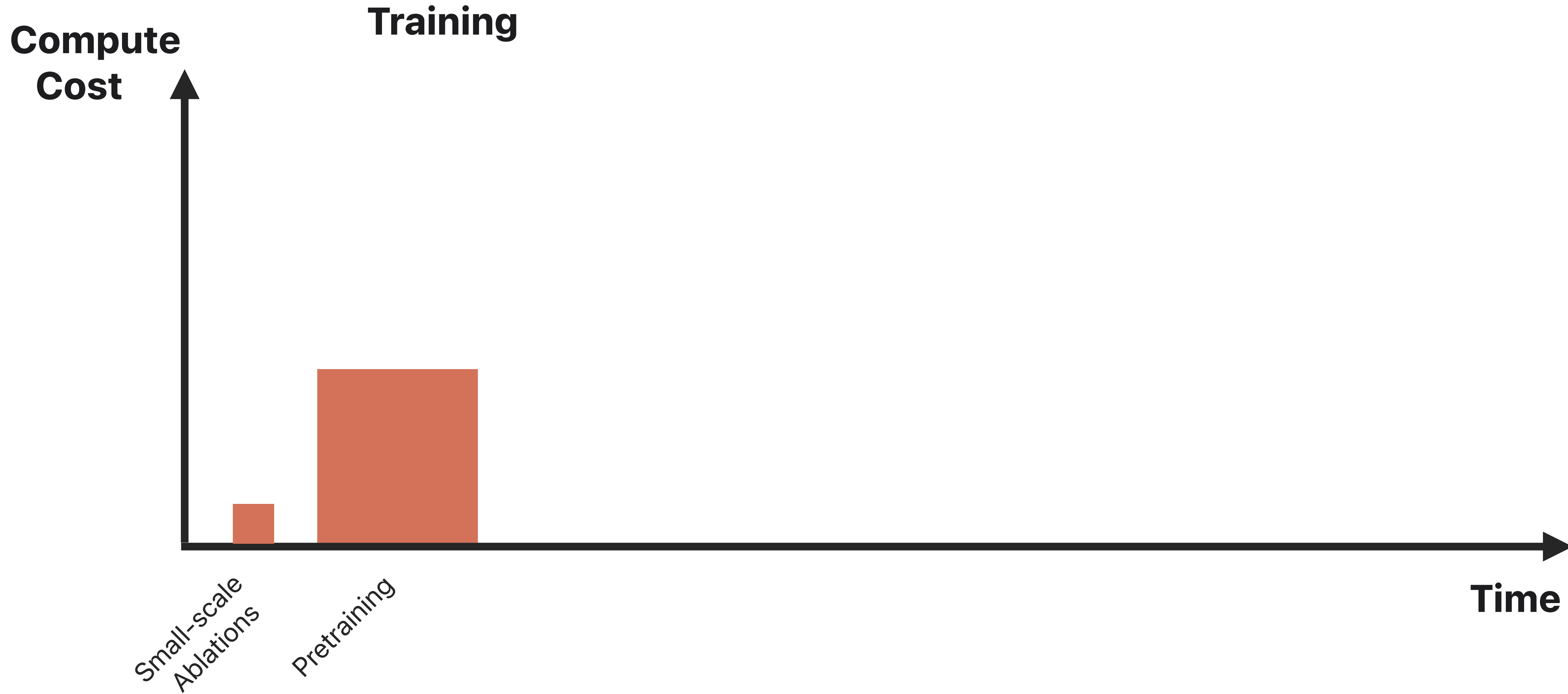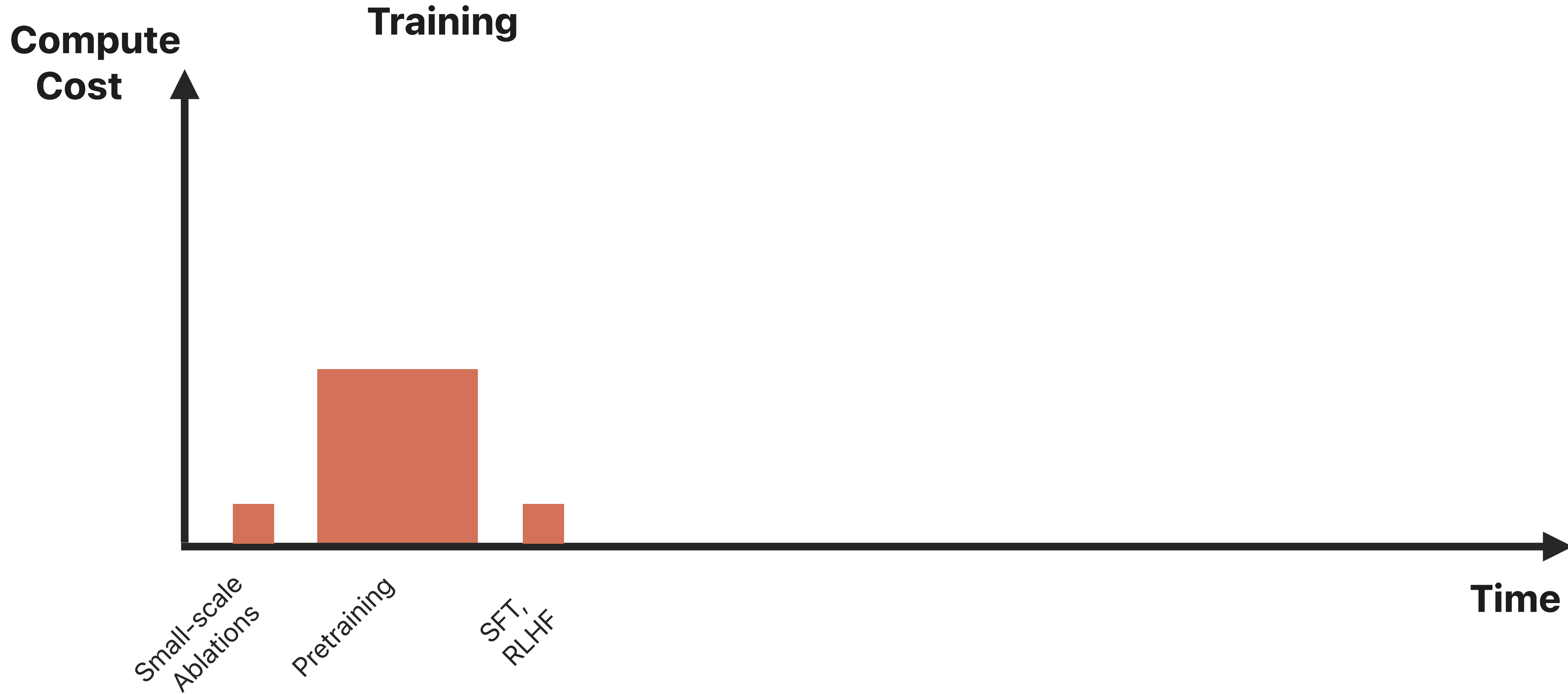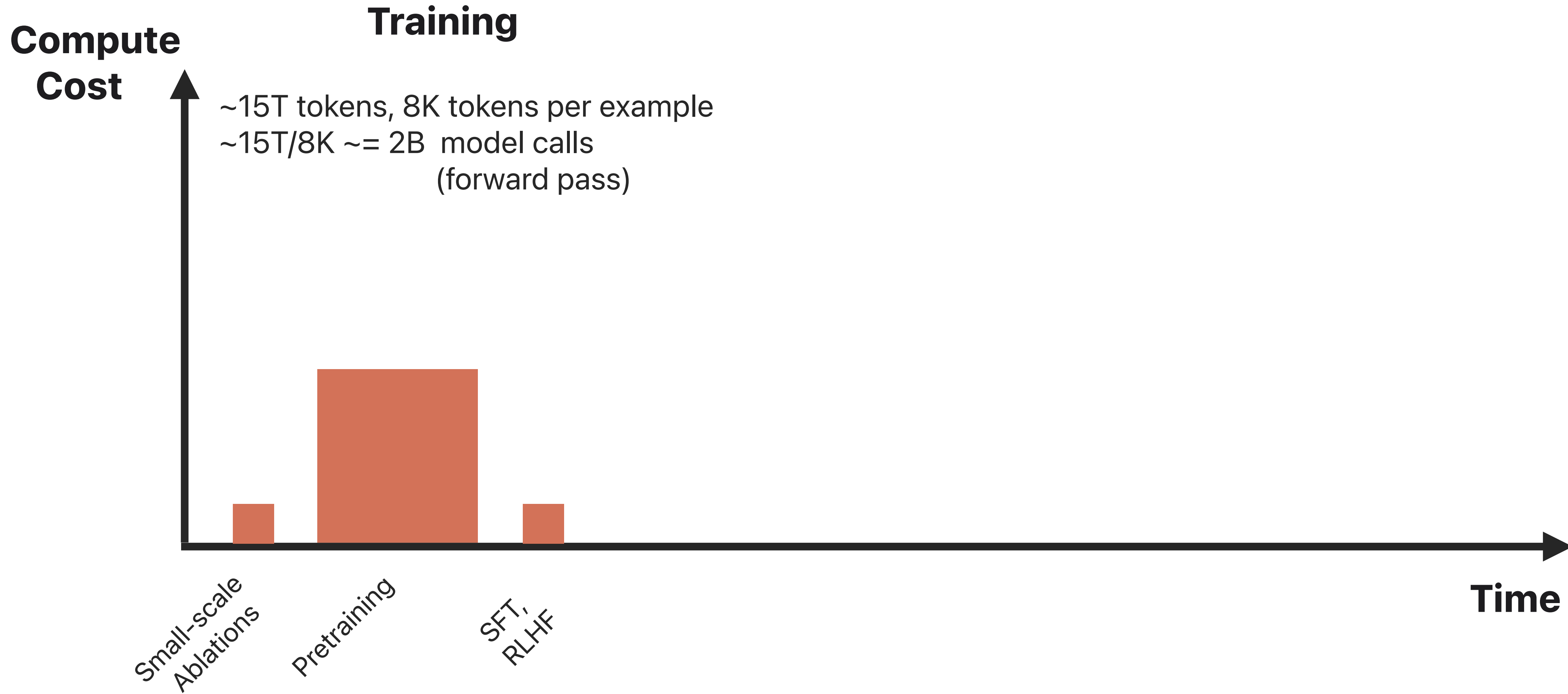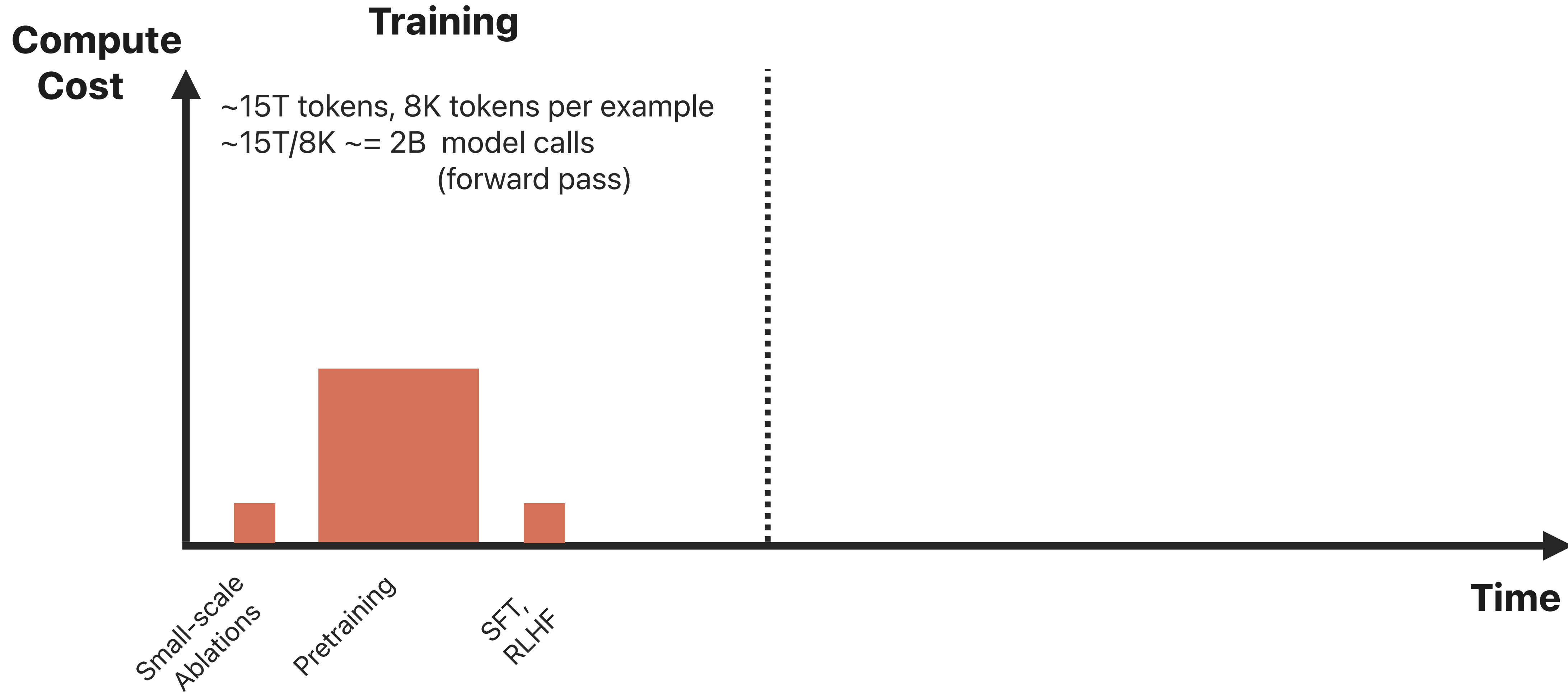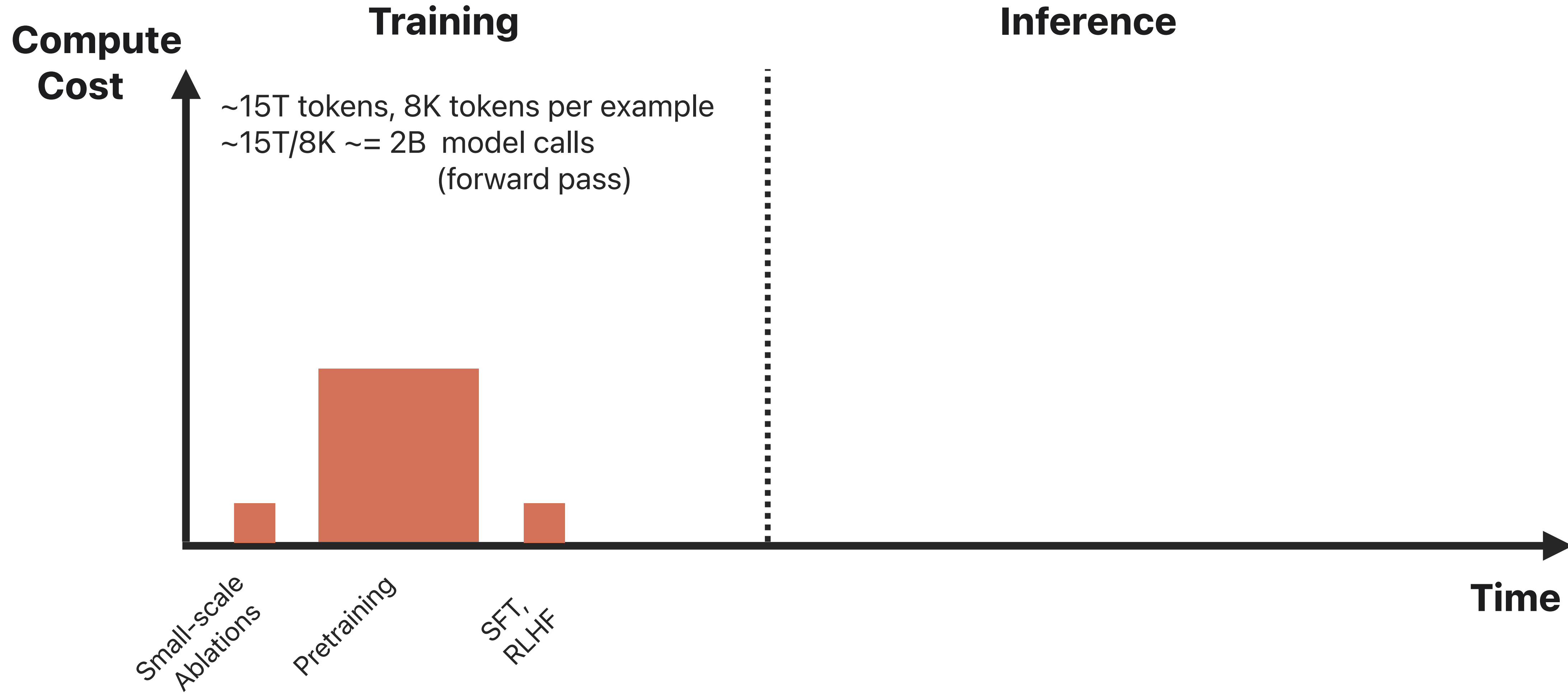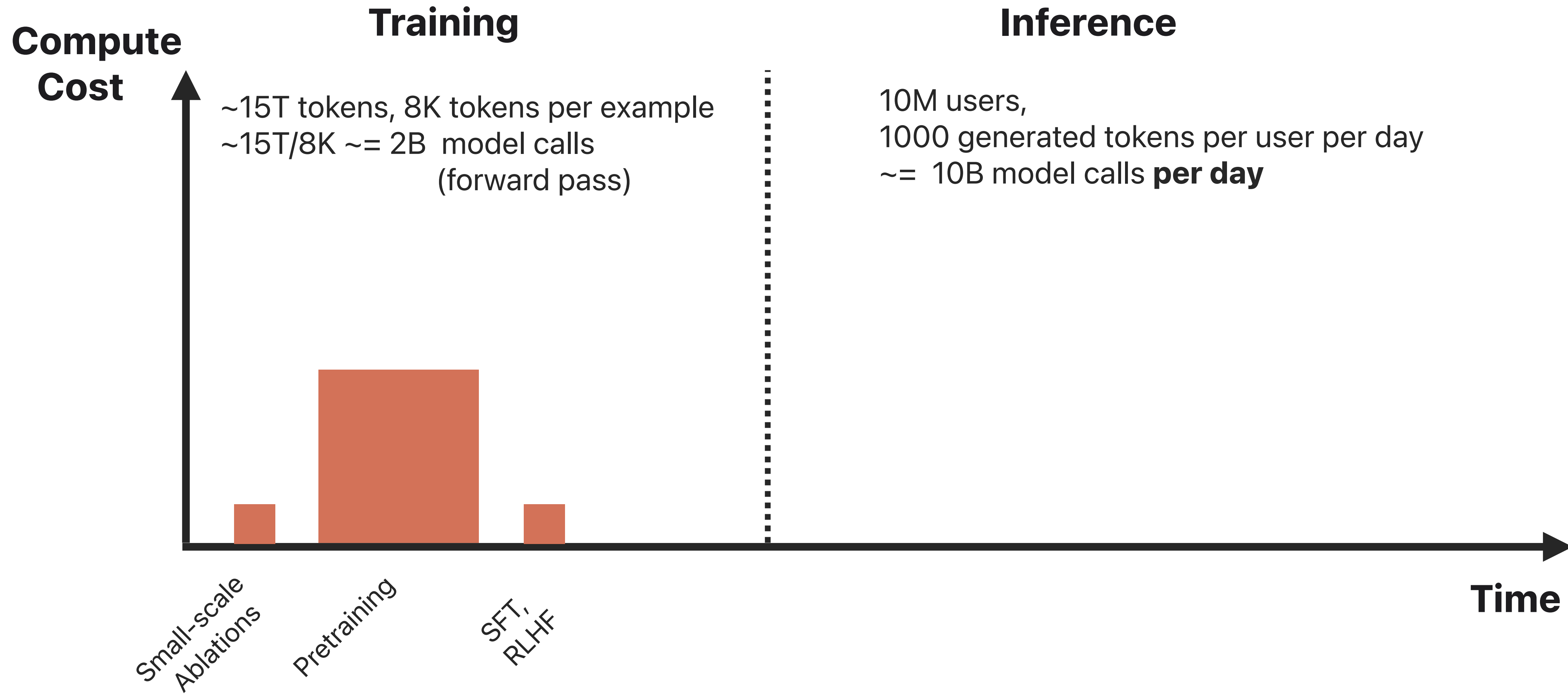**Compute Cost** ↑

**Time** →

# Motivation
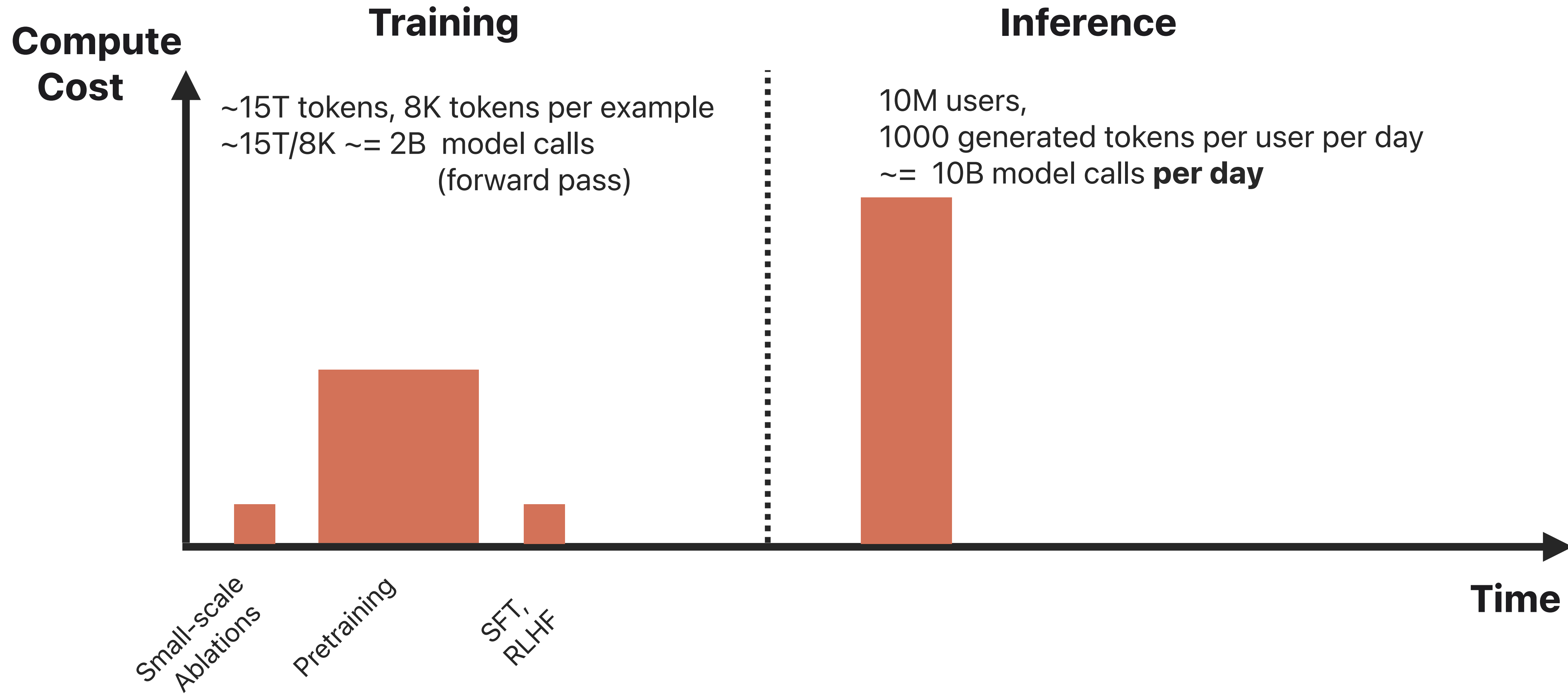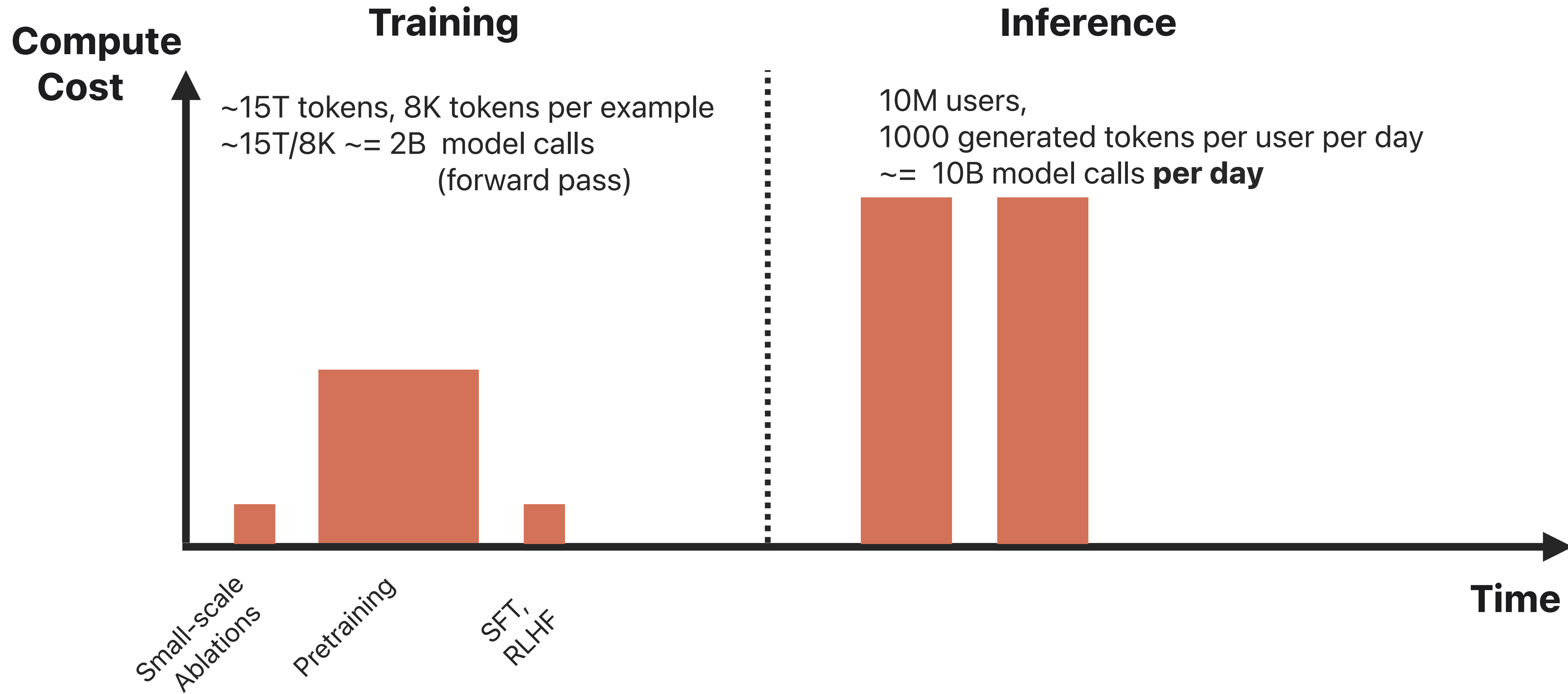Training vs Inference

# Motivation
Training vs Inference

# Motivation

Training vs Inference



**Compute Cost**

**Training**

Small-scale Ablations

Pretraining

SFT, RLHF

**Time**

# Motivation

Training vs Inference

**Compute Cost**

**Training**

~15T tokens, 8K tokens per example

~15T/8K ~= 2B  model calls

(forward pass)



Small-scale Ablations

Pretraining

SFT, RLHF

**Time**

# Motivation

Training vs Inference

**Compute Cost**

**Training**

~15T tokens, 8K tokens per example

~15T/8K ~= 2B  model calls

(forward pass)

Small-scale Ablations

Pretraining

SFT, RLHF

**Time**

# Motivation

Training vs Inference



**Compute Cost**

**Training**

**Inference**

~15T tokens, 8K tokens per example
~15T/8K ~= 2B model calls
(forward pass)

Small-scale Ablations

Pretraining

SFT, RLHF

**Time**

# Motivation

Training vs Inference



**Compute Cost**

**Training**

~15T tokens, 8K tokens per example
~15T/8K ~= 2B  model calls
                    (forward pass)

**Inference**

10M users,
1000 generated tokens per user per day
~=  10B model calls **per day**

Small-scale
Ablations

Pretraining

SFT,
RLHF

**Time**

# Motivation

Training vs Inference

**Compute Cost**

**Training**

~15T tokens, 8K tokens per example
~15T/8K ~= 2B  model calls
(forward pass)

**Inference**

10M users,
1000 generated tokens per user per day
~=  10B model calls **per day**

Small-scale Ablations

Pretraining

SFT, RLHF

**Time**

# Motivation

Training vs Inference



**Compute Cost**

**Training**

~15T tokens, 8K tokens per example
~15T/8K ~= 2B  model calls
(forward pass)

**Inference**

10M users,
1000 generated tokens per user per day
~=  10B model calls **per day**

Small-scale Ablations

Pretraining

SFT, RLHF

**Time**

# Motivation

Training vs Inference

**Compute Cost**

**Training**

~15T tokens, 8K tokens per example
~15T/8K ~= 2B  model calls
(forward pass)

**Inference**

10M users,
1000 generated tokens per user per day
~=  10B model calls **per day**



Small-scale Ablations

Pretraining

SFT, RLHF

**Time**

2

# Motivation

Training vs Inference



**Compute Cost**

## Training

~15T tokens, 8K tokens per example
~15T/8K ~= 2B  model calls
(forward pass)

## Inference

10M users,
1000 generated tokens per user per day
~=  10B model calls **per day**

Small-scale Ablations

Pretraining

SFT, RLHF

. . .

**Time**

[1] OpenAI's ChatGPT now has 100 million weekly active users
https://techcrunch.com/2023/11/06/openais-chatgpt-now-has-100-million-weekly-active-users/

# Inference costs almost always outweigh training costs.

# Motivation: Where are we spending our compute?



Transformer Block

# Motivation: Where are we spending our compute?



Transformer Block

# Motivation: Where are we spending our compute?



Simple FC Layers!

Feed Forward

Normalization

Attention

Normalization

~2/3 of FLOPS

~1/3

Transformer Block

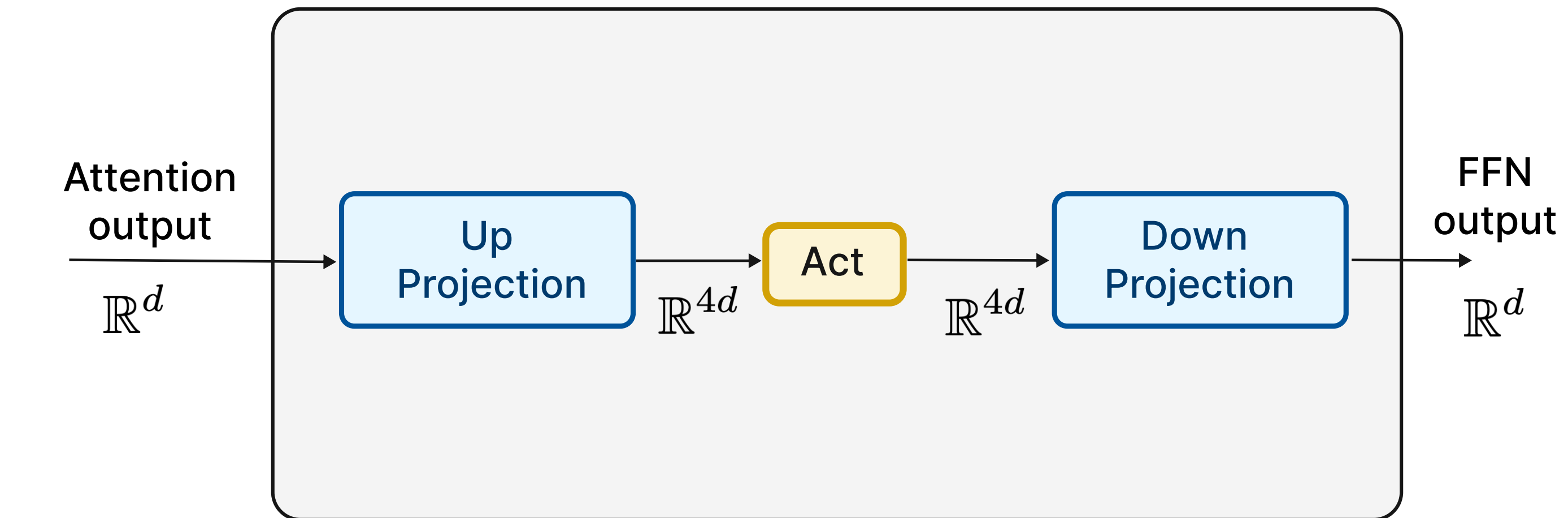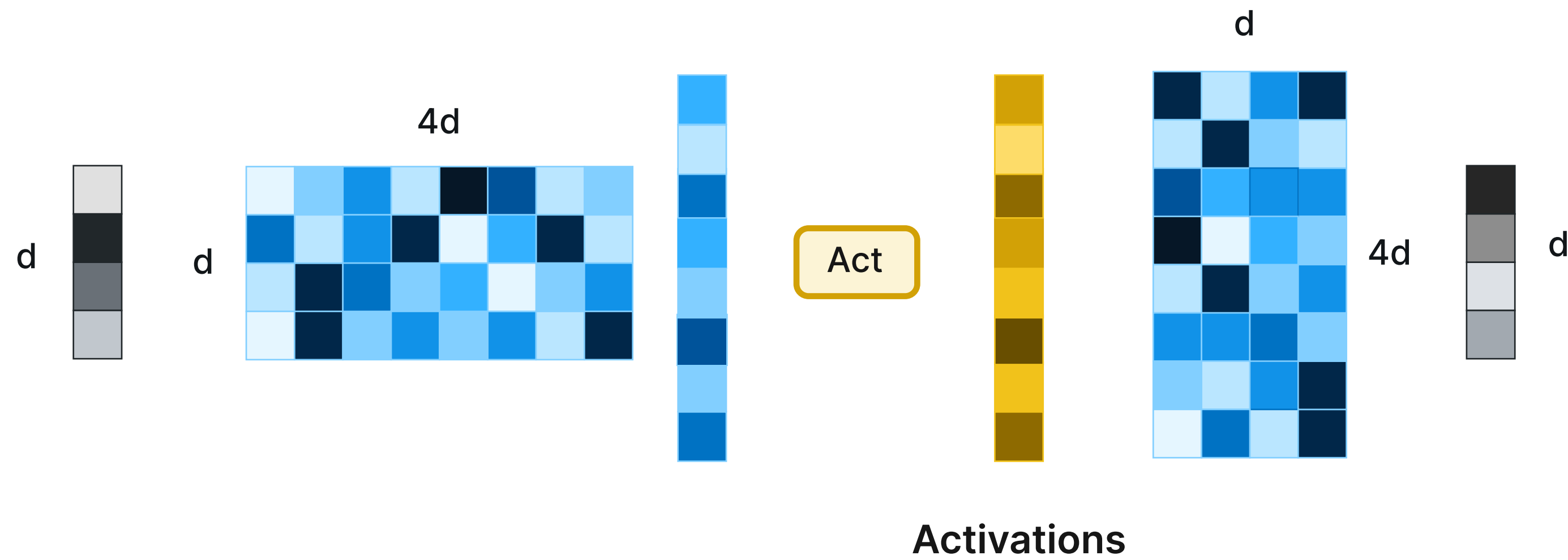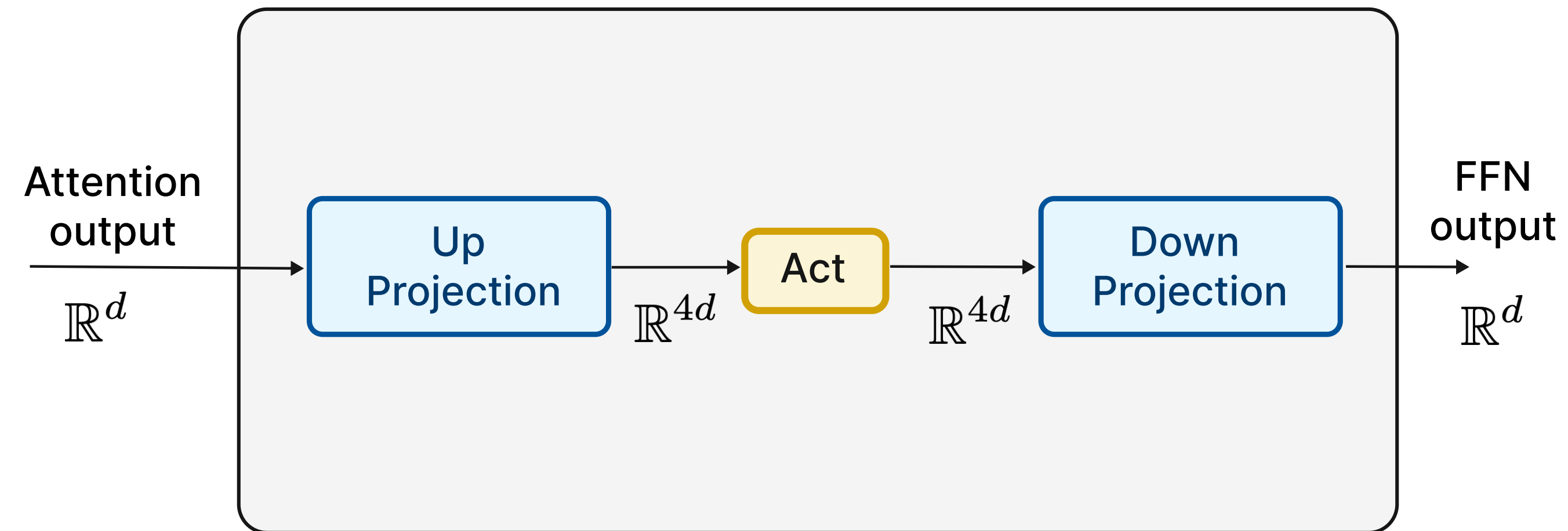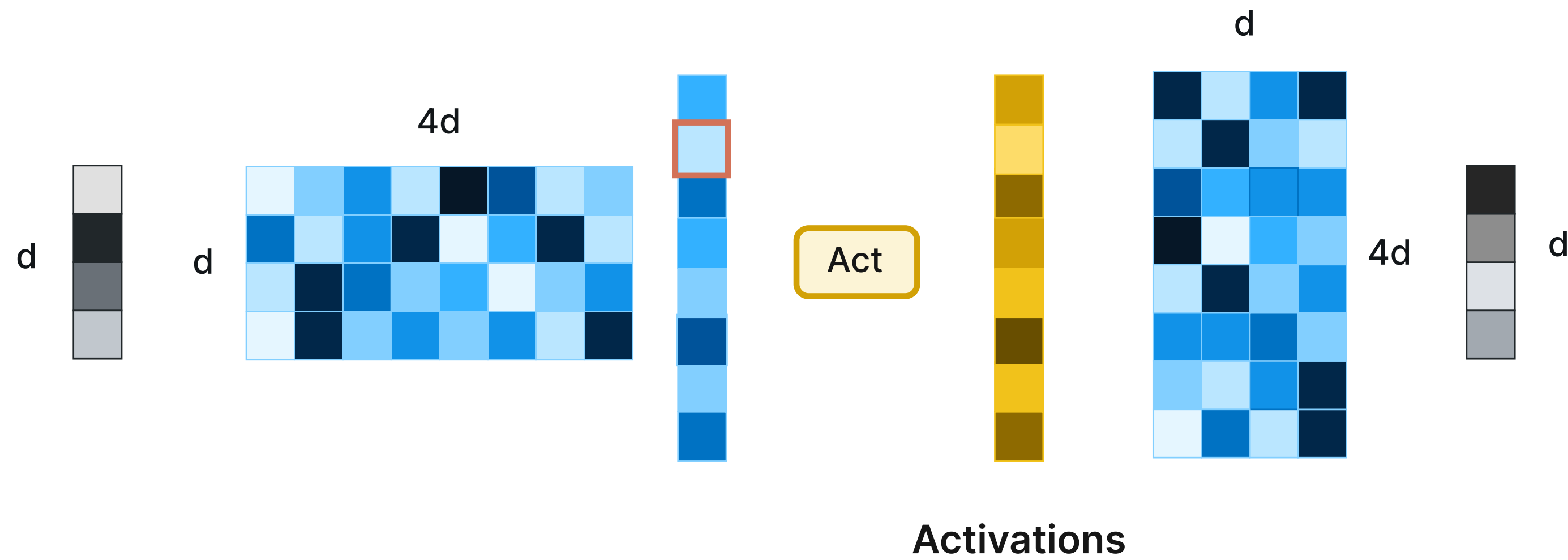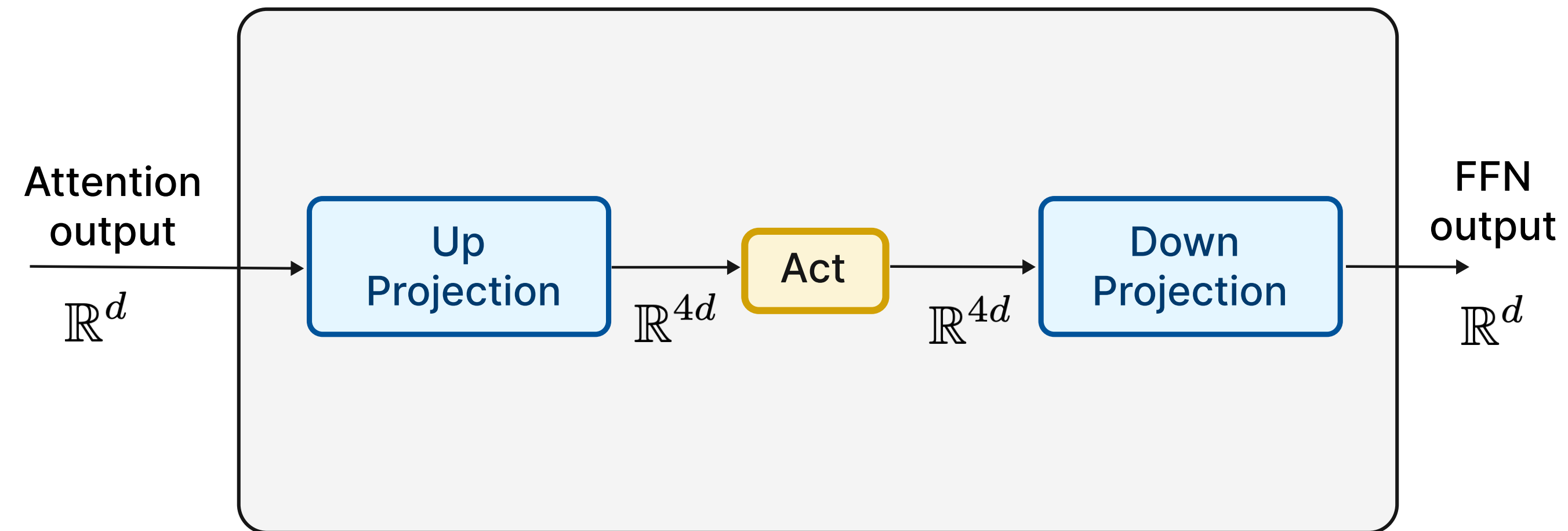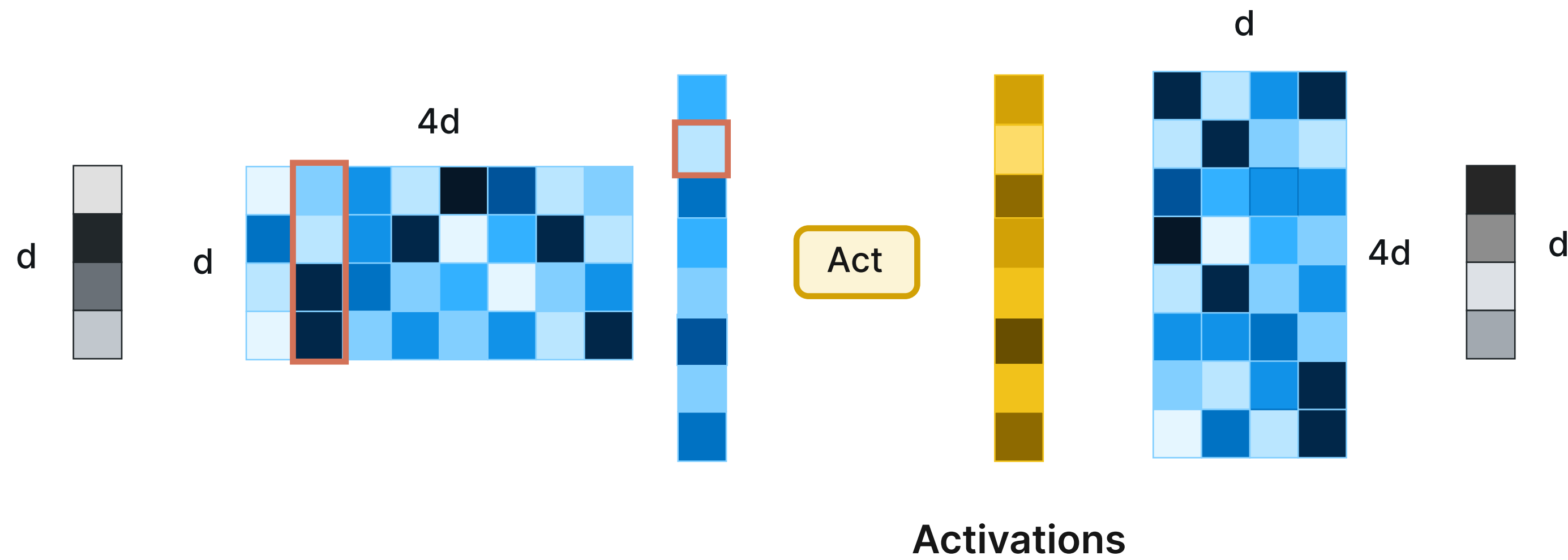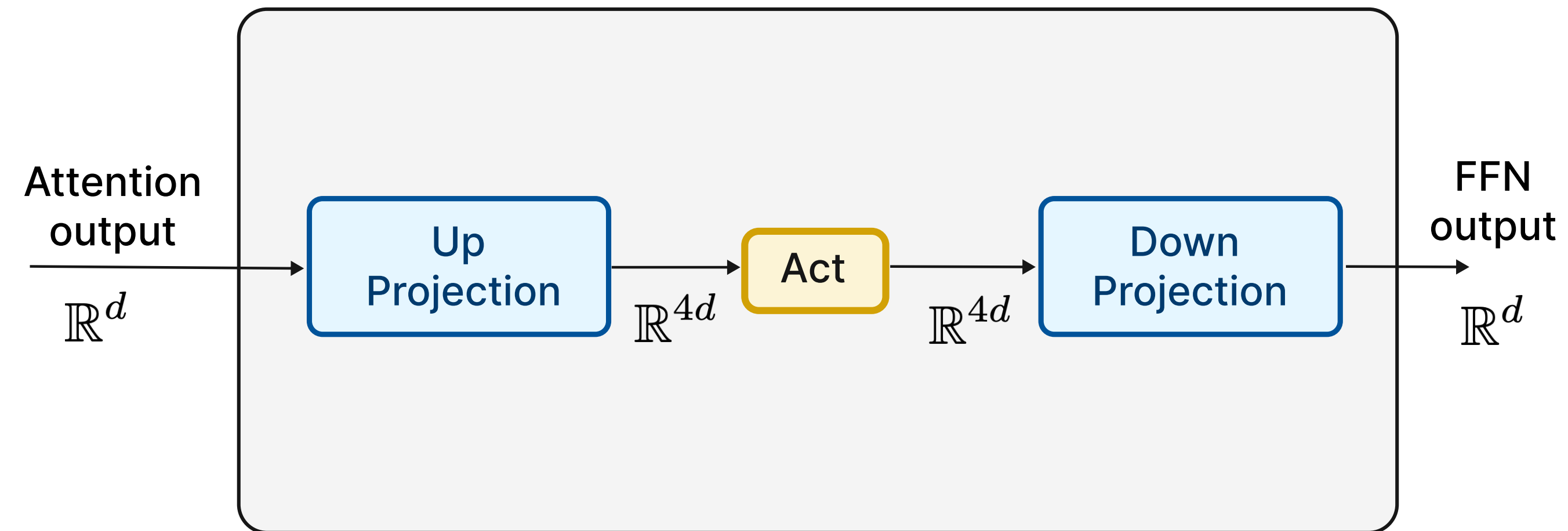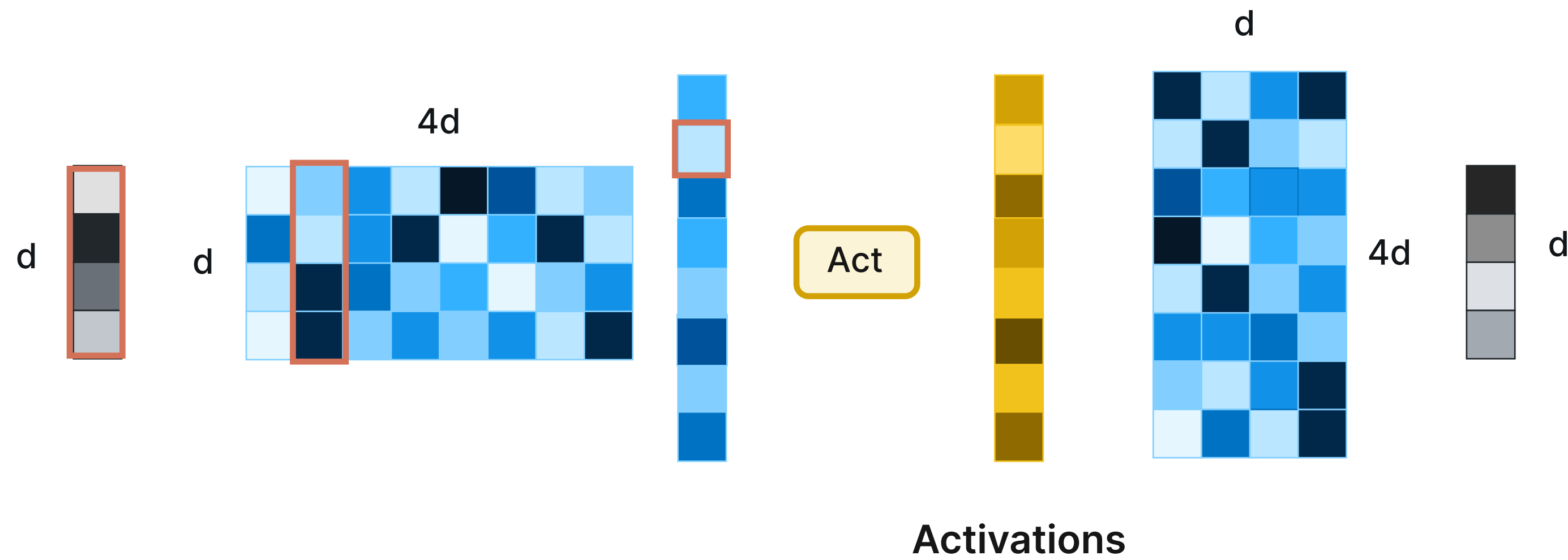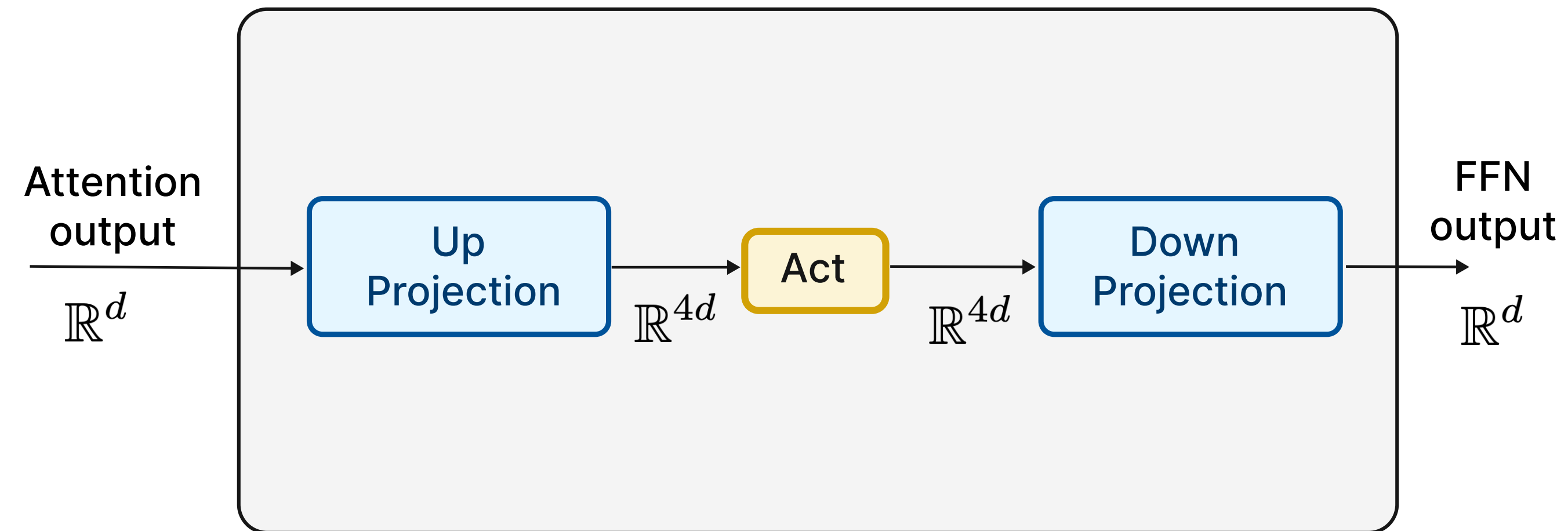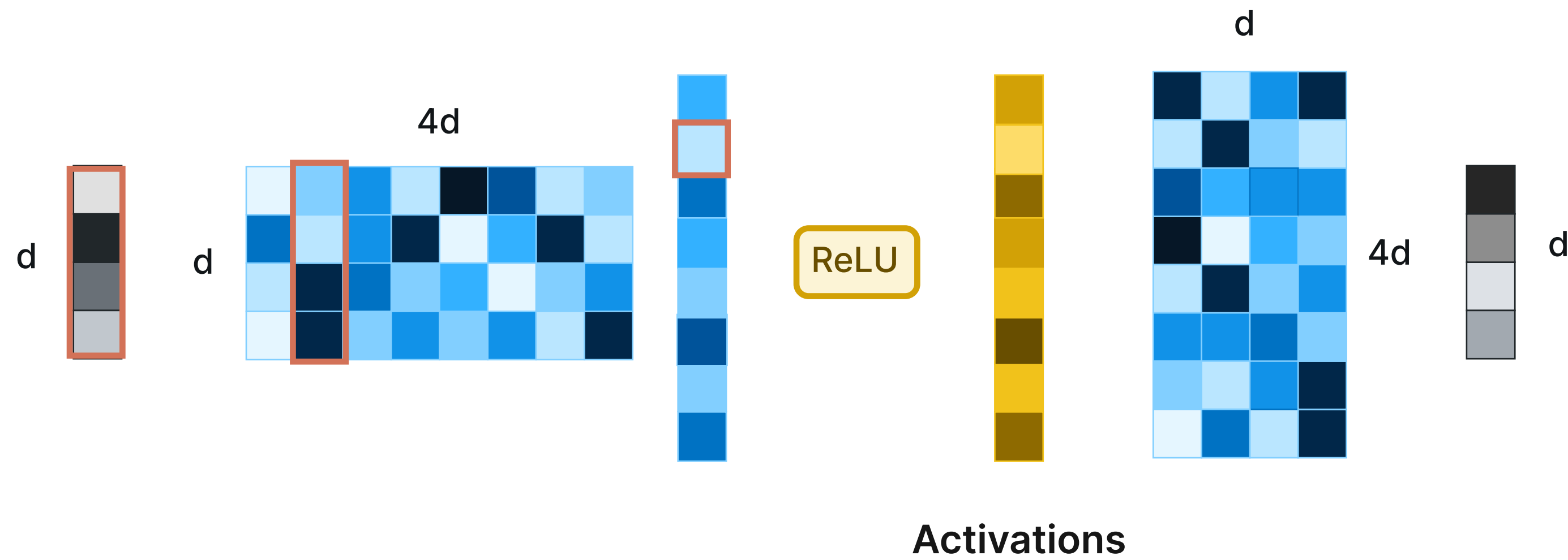# Motivation: Skipping FFN computation

# Motivation: Skipping FFN computation

# Motivation: Skipping FFN computation

# Motivation: Skipping FFN computation

# Motivation: Skipping FFN computation

# Motivation: Skipping FFN computation
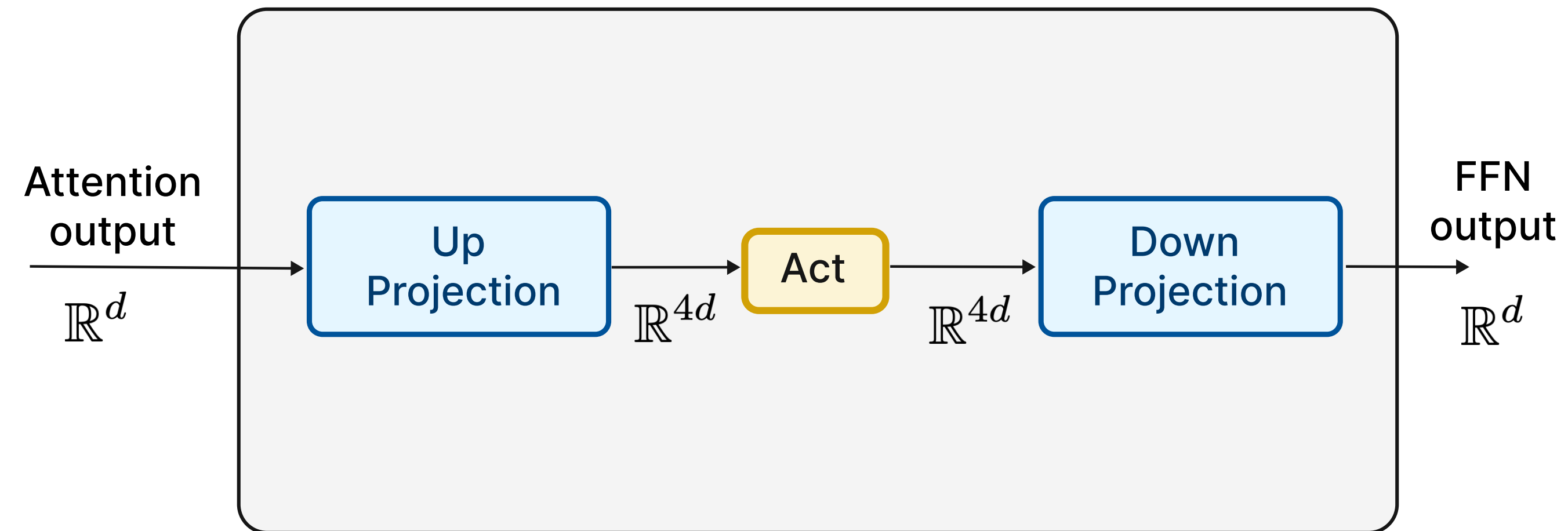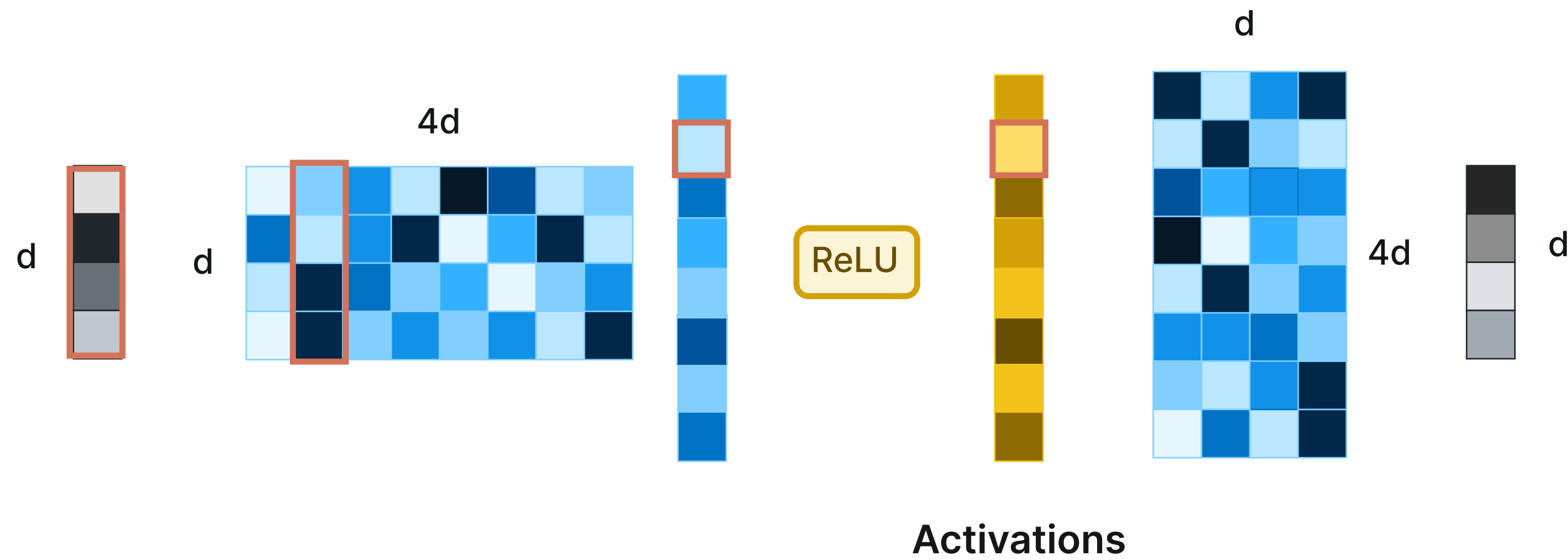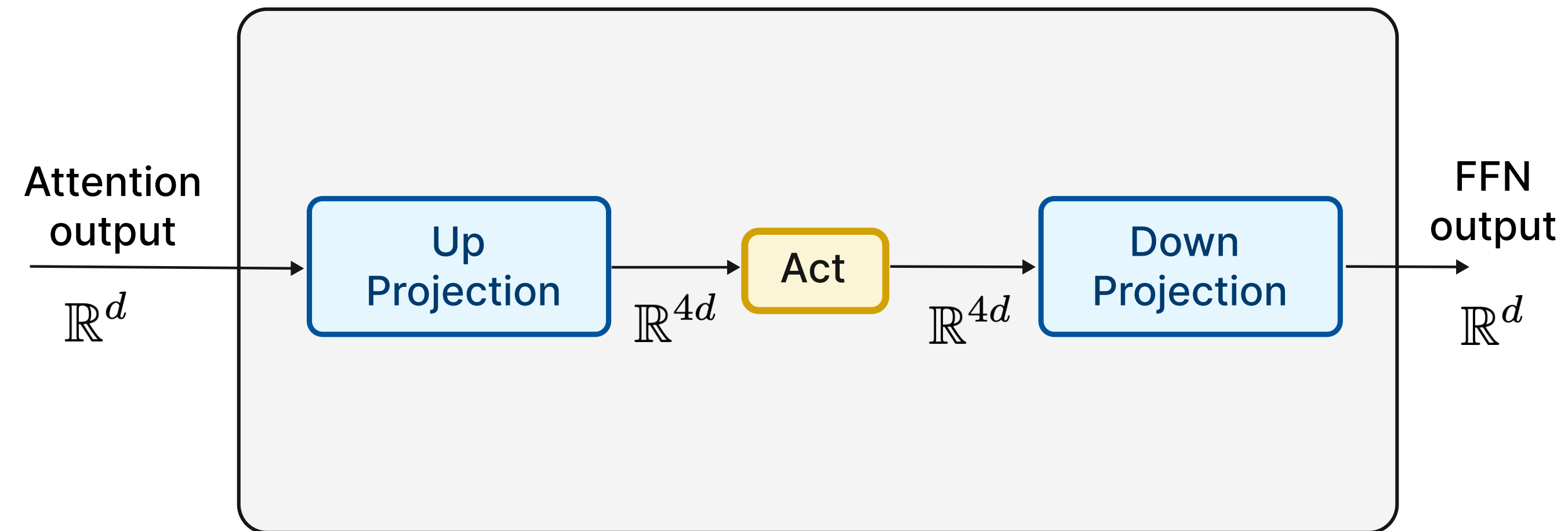


Attention output $\mathbb{R}^d$ → Up Projection → $\mathbb{R}^{4d}$ → Act → $\mathbb{R}^{4d}$ → Down Projection → FFN output $\mathbb{R}^d$

Activations

# Motivation: Skipping FFN computation



Attention output $\mathbb{R}^d$ → Up Projection $\mathbb{R}^{4d}$ → Act $\mathbb{R}^{4d}$ → Down Projection → FFN output $\mathbb{R}^d$
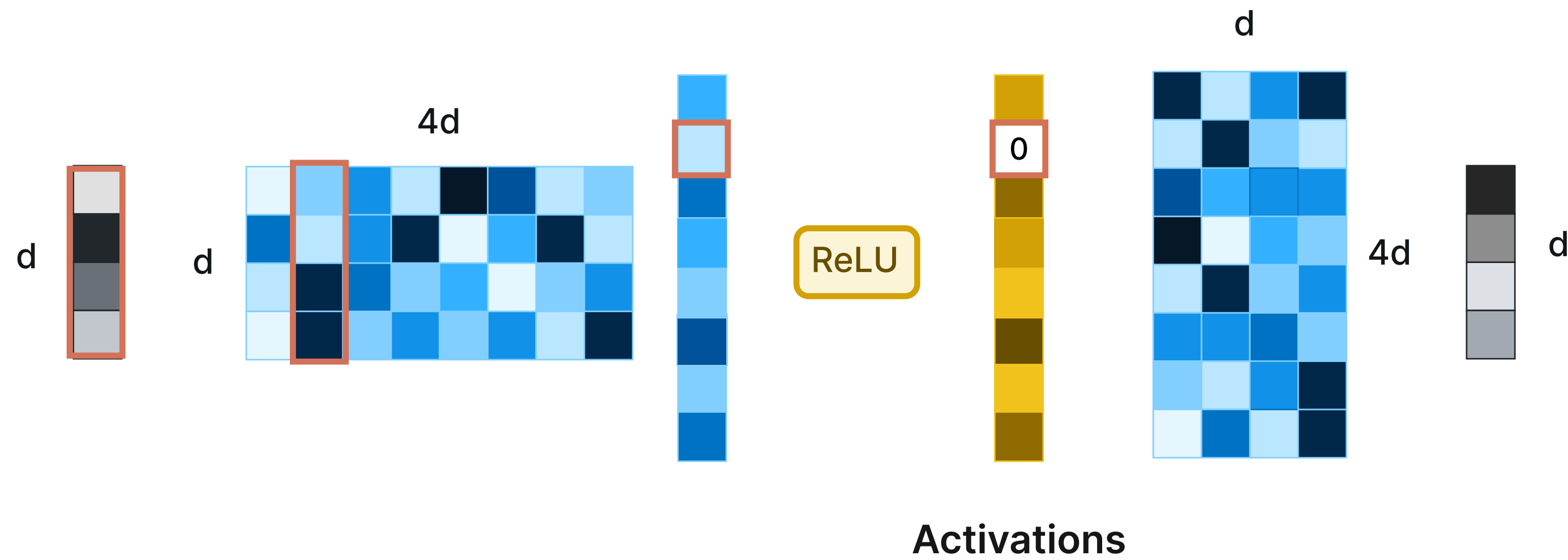
Activations

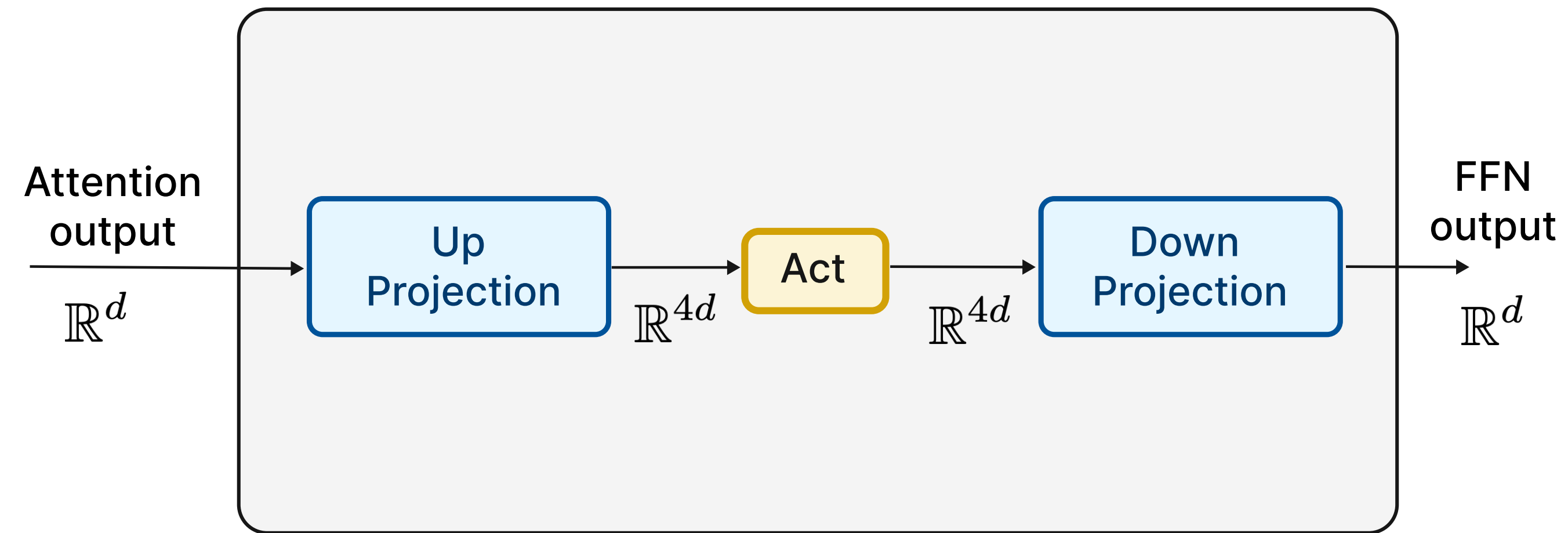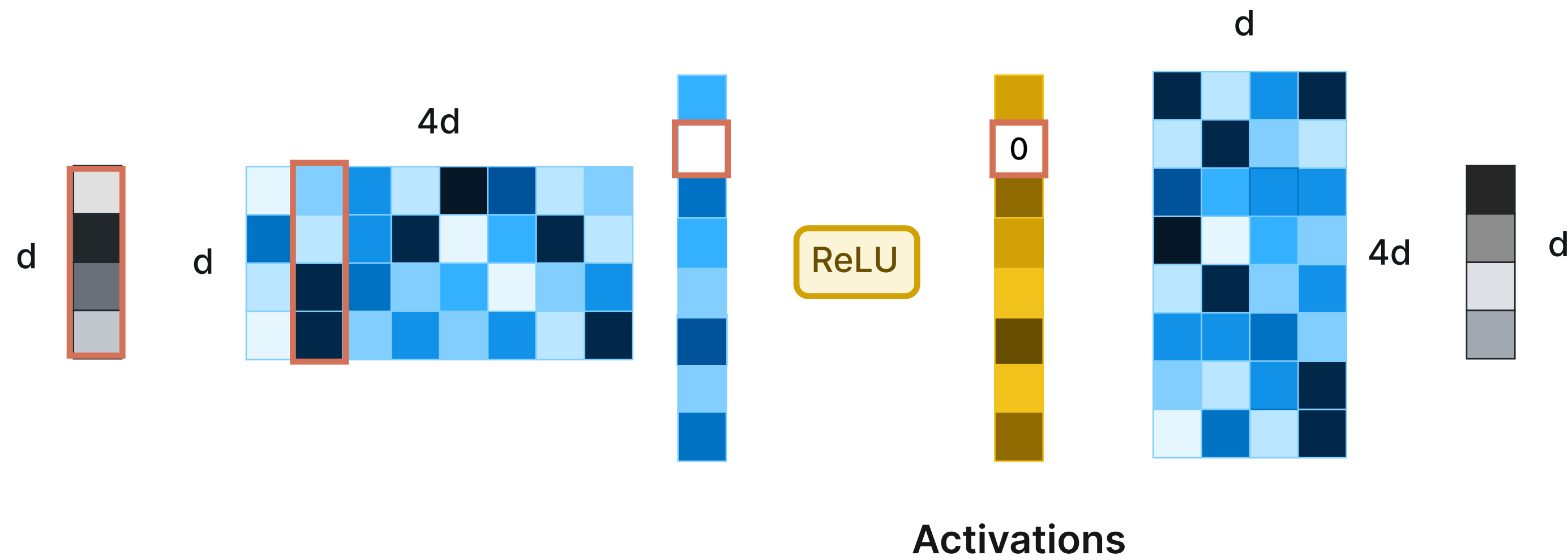# Motivation: Skipping FFN computation
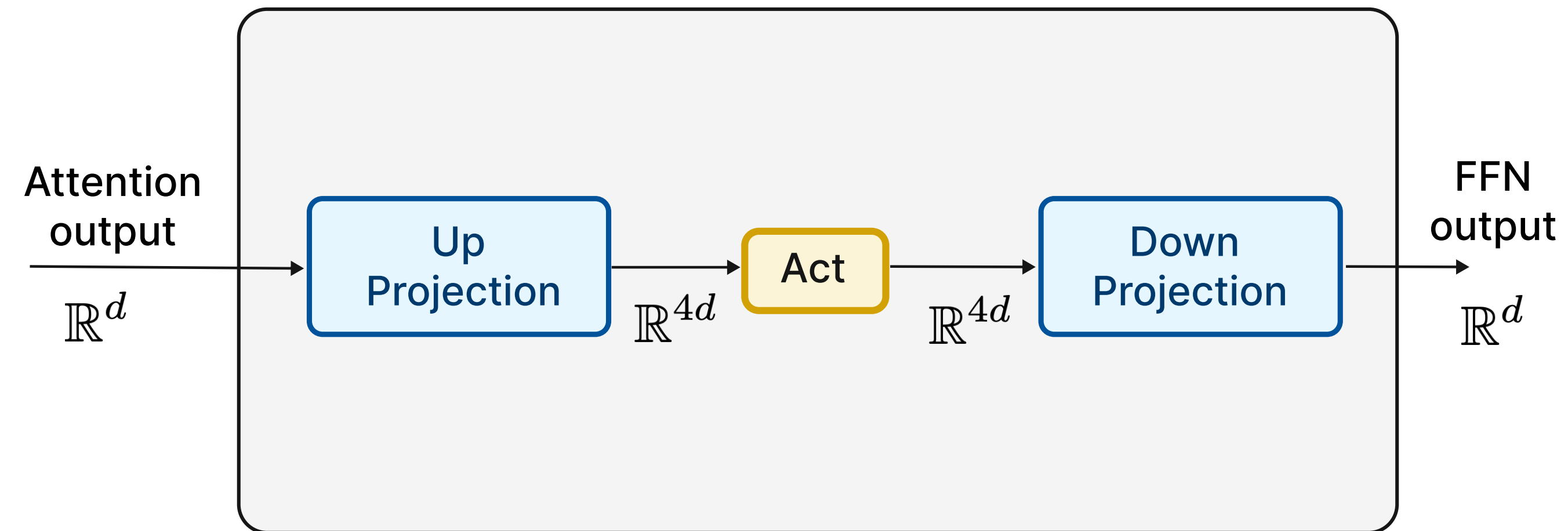
# Motivation: Skipping FFN computation



Attention output $\mathbb{R}^d$ → Up Projection → $\mathbb{R}^{4d}$ → Act → $\mathbb{R}^{4d}$ → Down Projection → FFN output $\mathbb{R}^d$
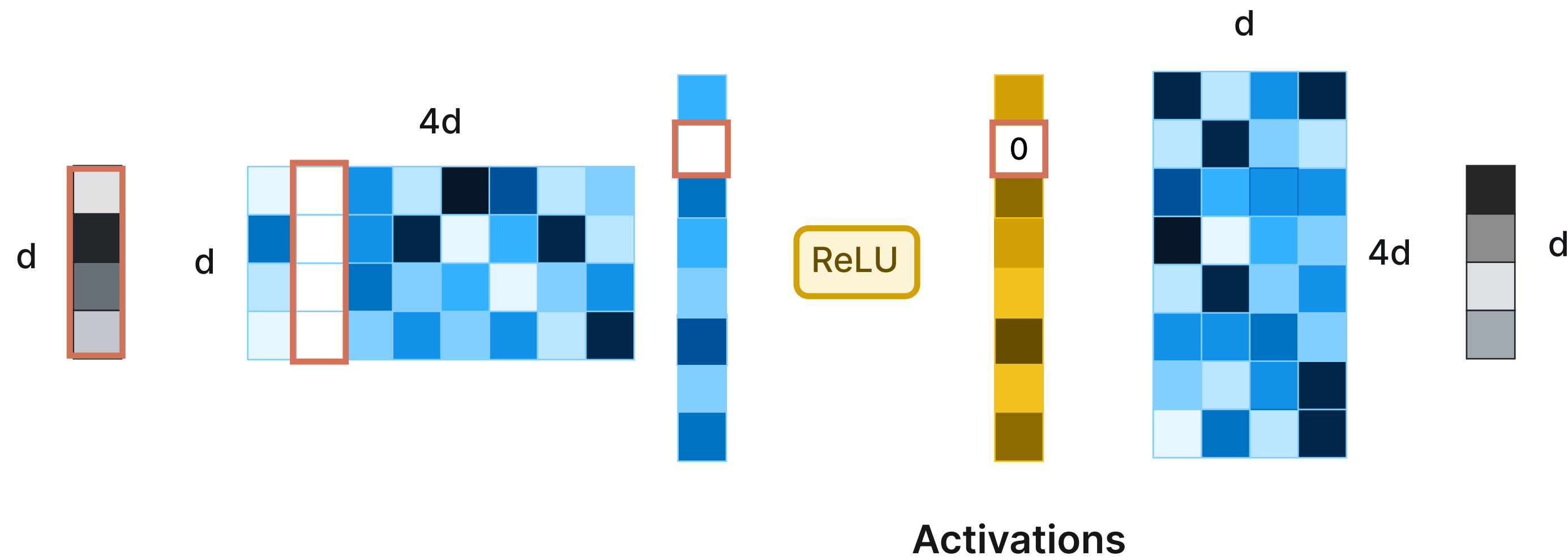
Activations

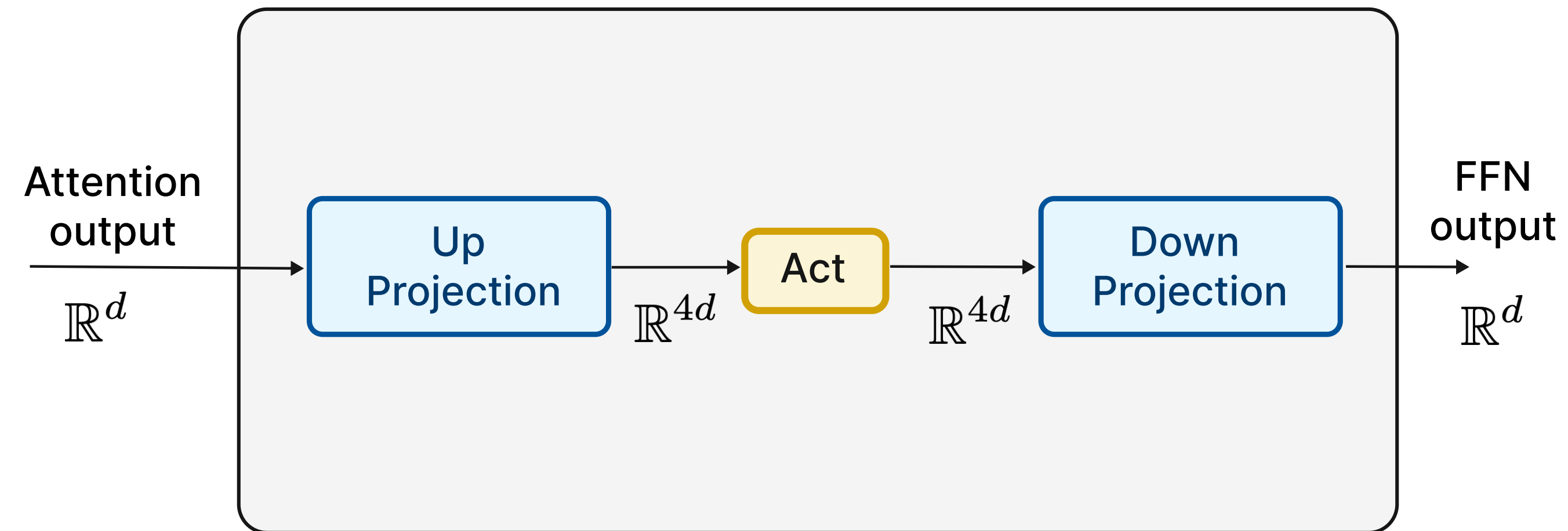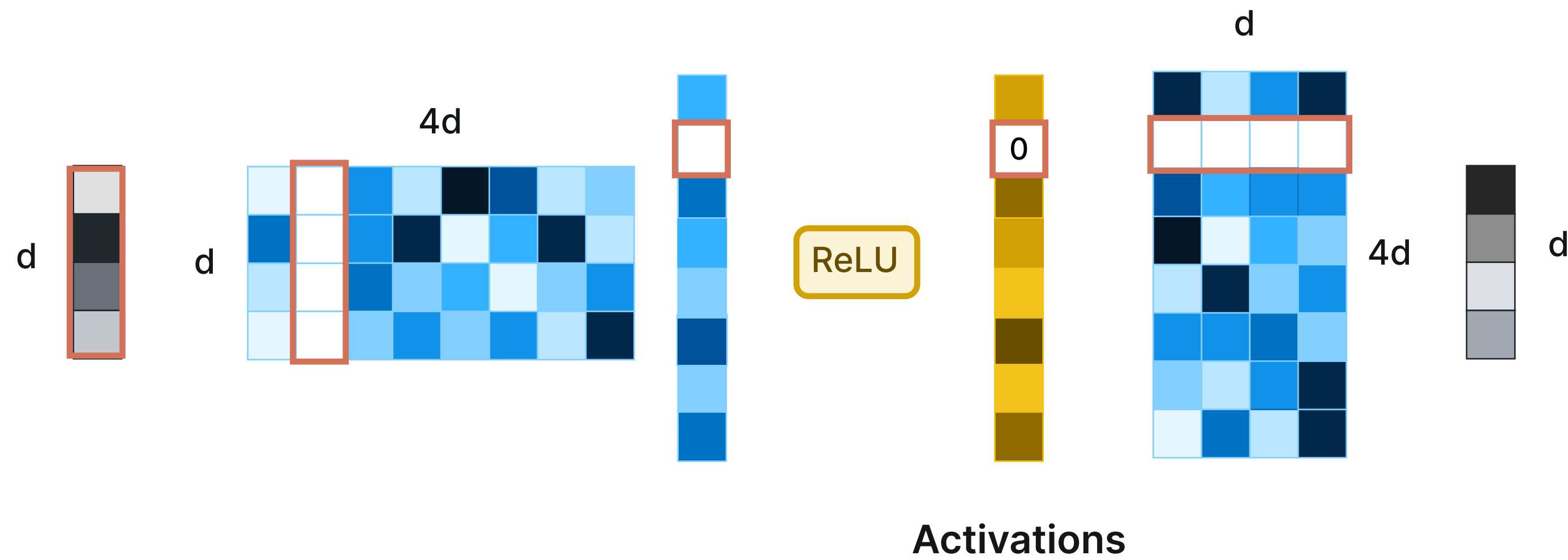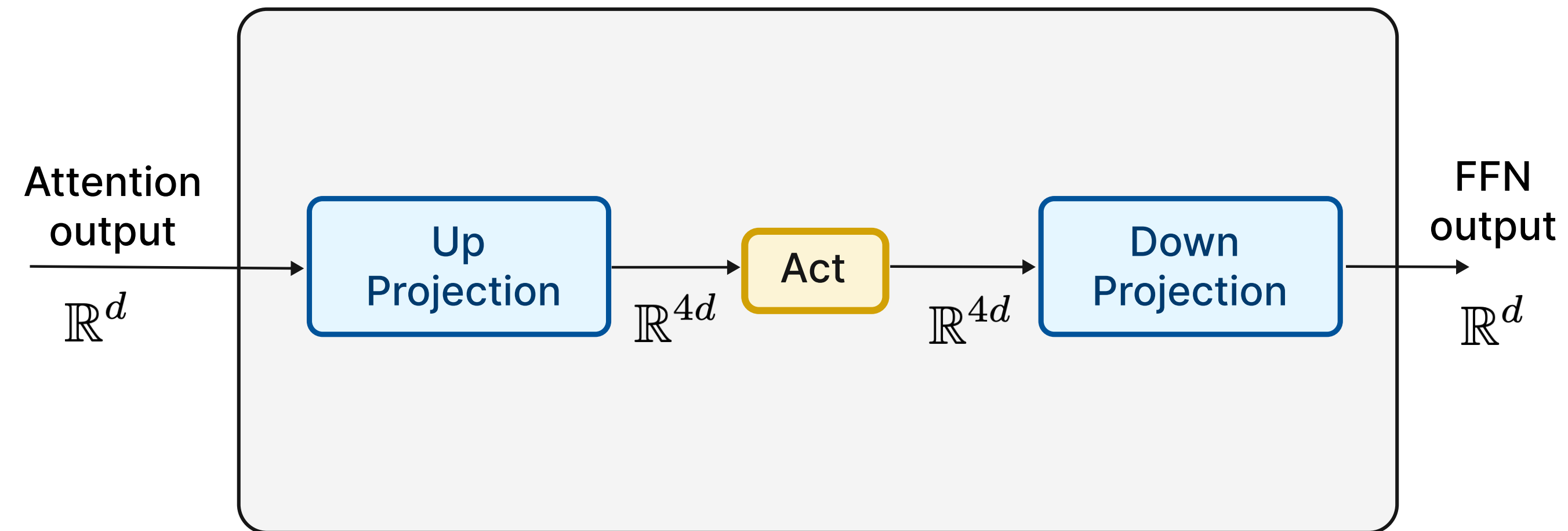# Motivation: Skipping FFN computation

# Motivation: Skipping FFN computation

# Motivation: Skipping FFN computation



Activations

# Motivation: Skipping FFN computation



Attention output $\mathbb{R}^d$

Up Projection $\mathbb{R}^{4d}$

Act $\mathbb{R}^{4d}$

Down Projection

FFN output $\mathbb{R}^d$
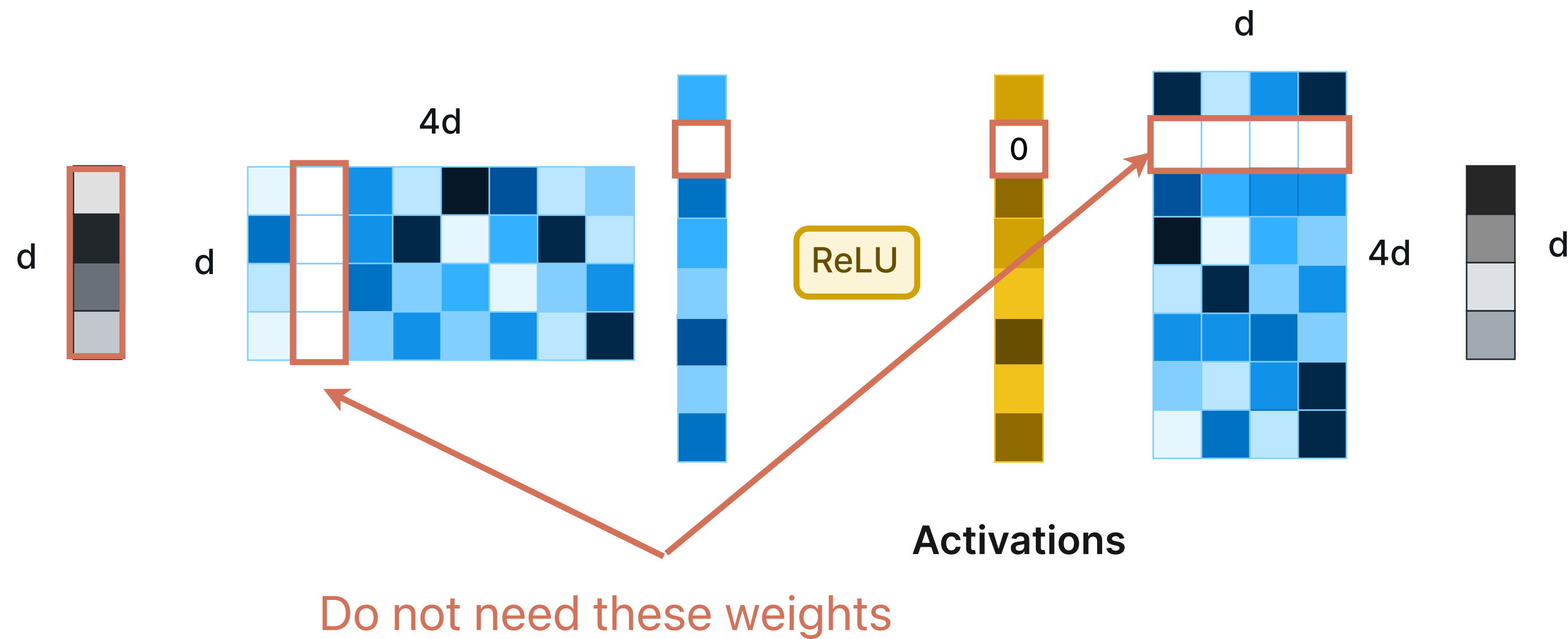
d

4d

d

ReLU

0

d

4d

d

Activations

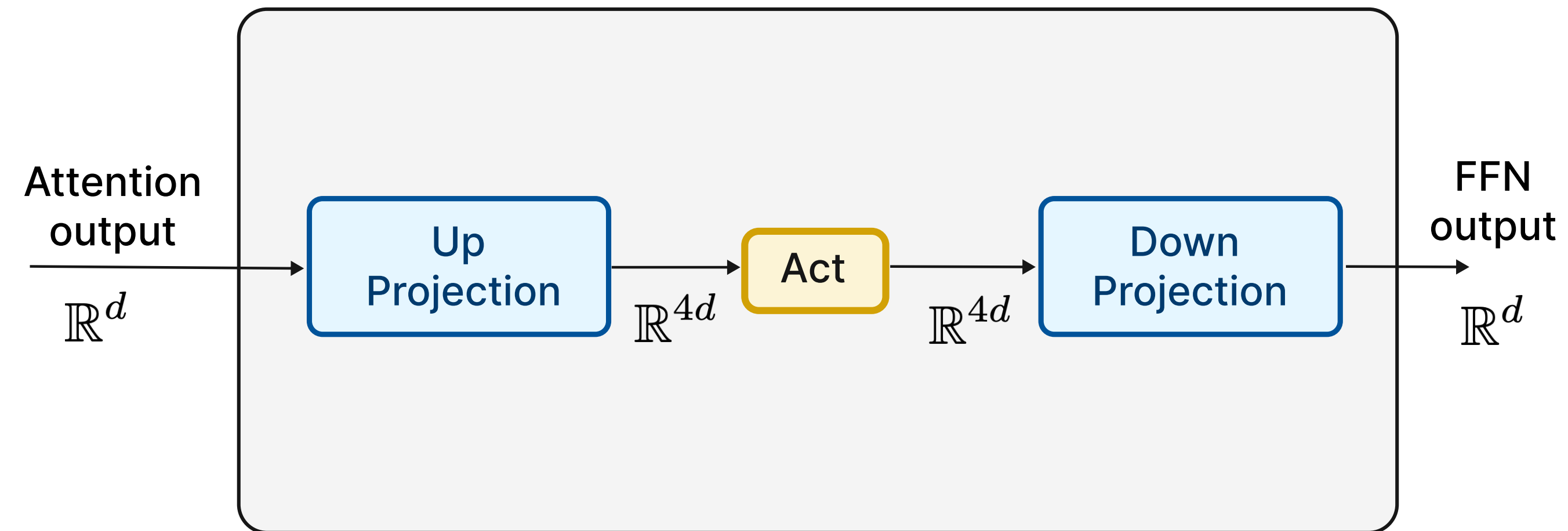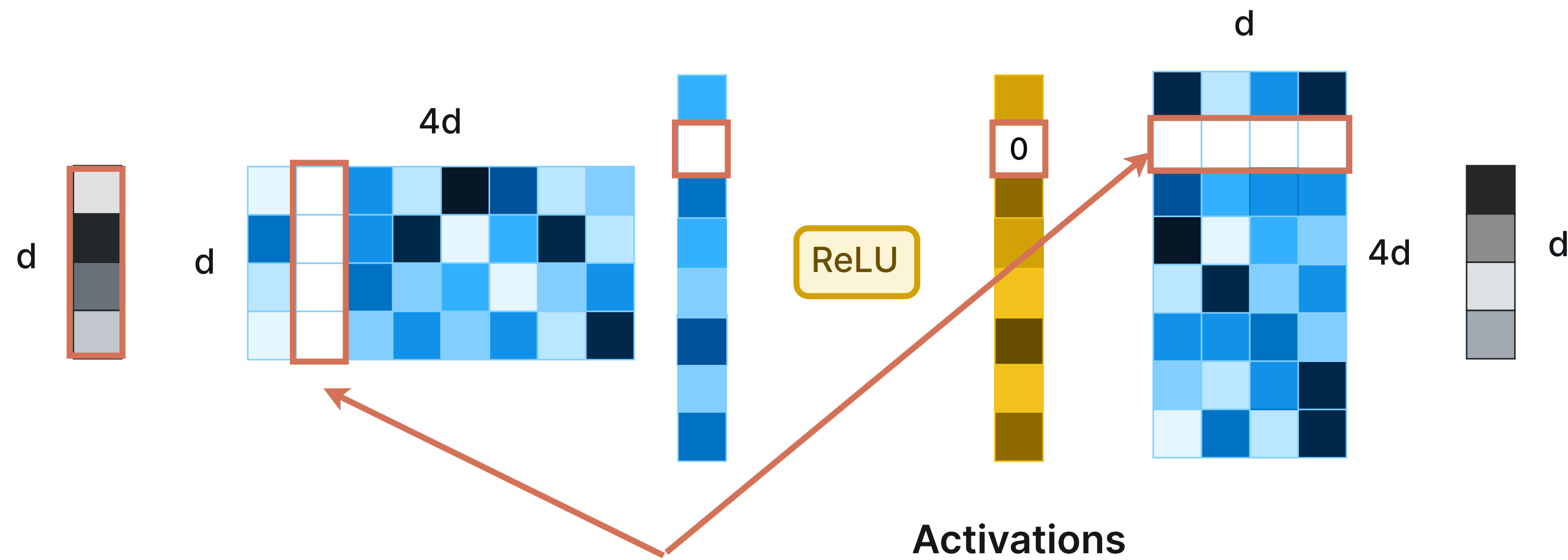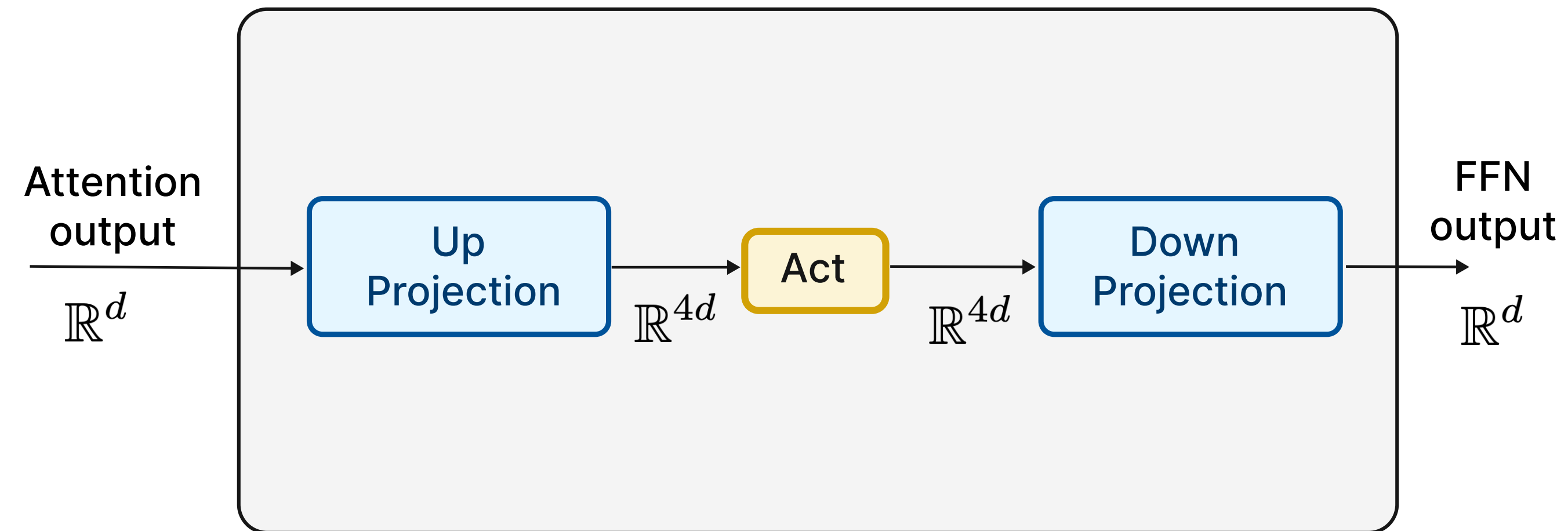# Motivation: Skipping FFN computation

# Motivation: Skipping FFN computation

# Motivation: Skipping FFN computation

# Motivation: Skipping FFN computation



Attention output $\mathbb{R}^d$ → Up Projection → $\mathbb{R}^{4d}$ → Act → $\mathbb{R}^{4d}$ → Down Projection → FFN output $\mathbb{R}^d$

ReLU

Activations

Do not need these weights

# Motivation: Skipping FFN computation



Attention output $\mathbb{R}^d$ → Up Projection $\mathbb{R}^{4d}$ → Act $\mathbb{R}^{4d}$ → Down Projection → FFN output $\mathbb{R}^d$
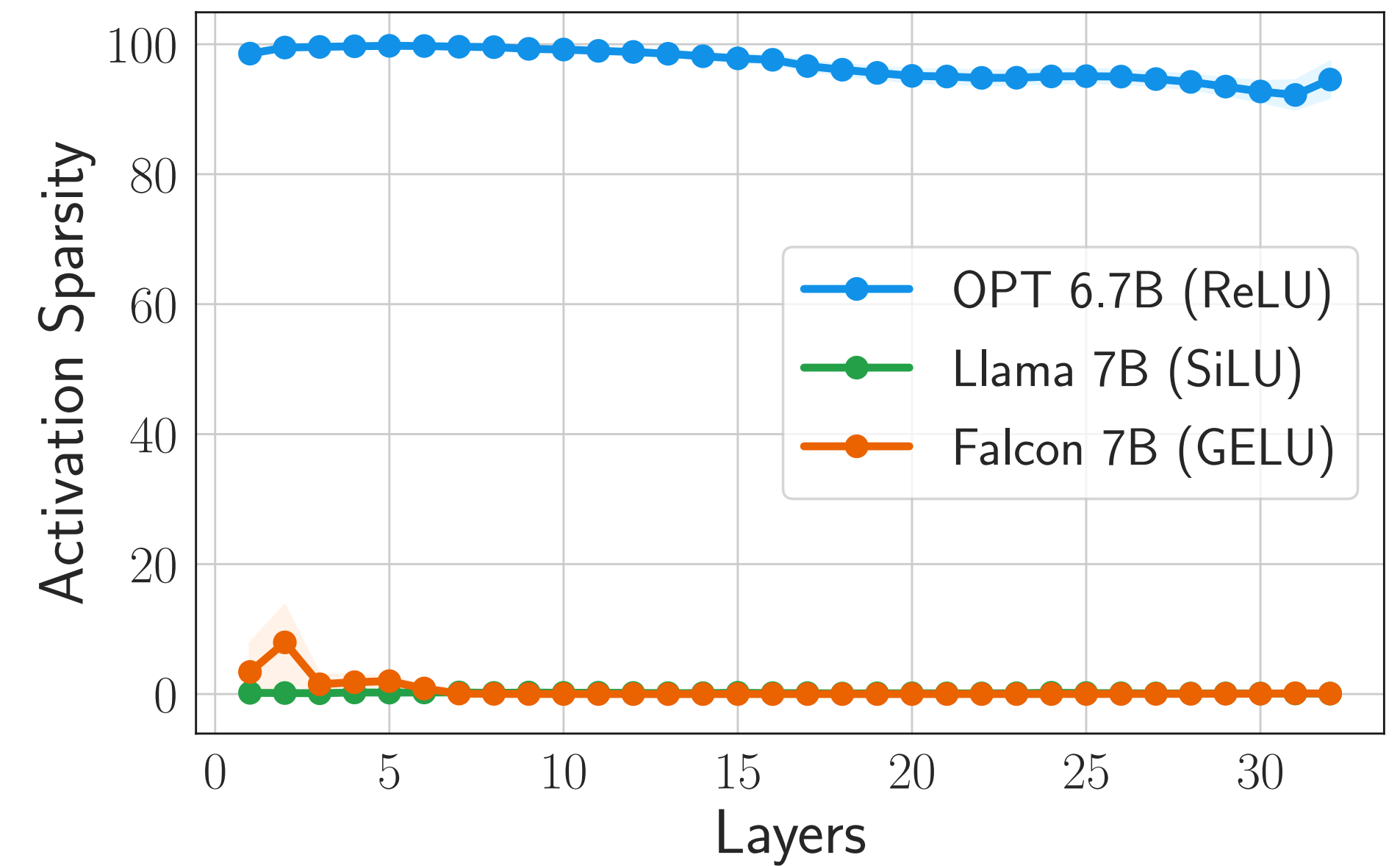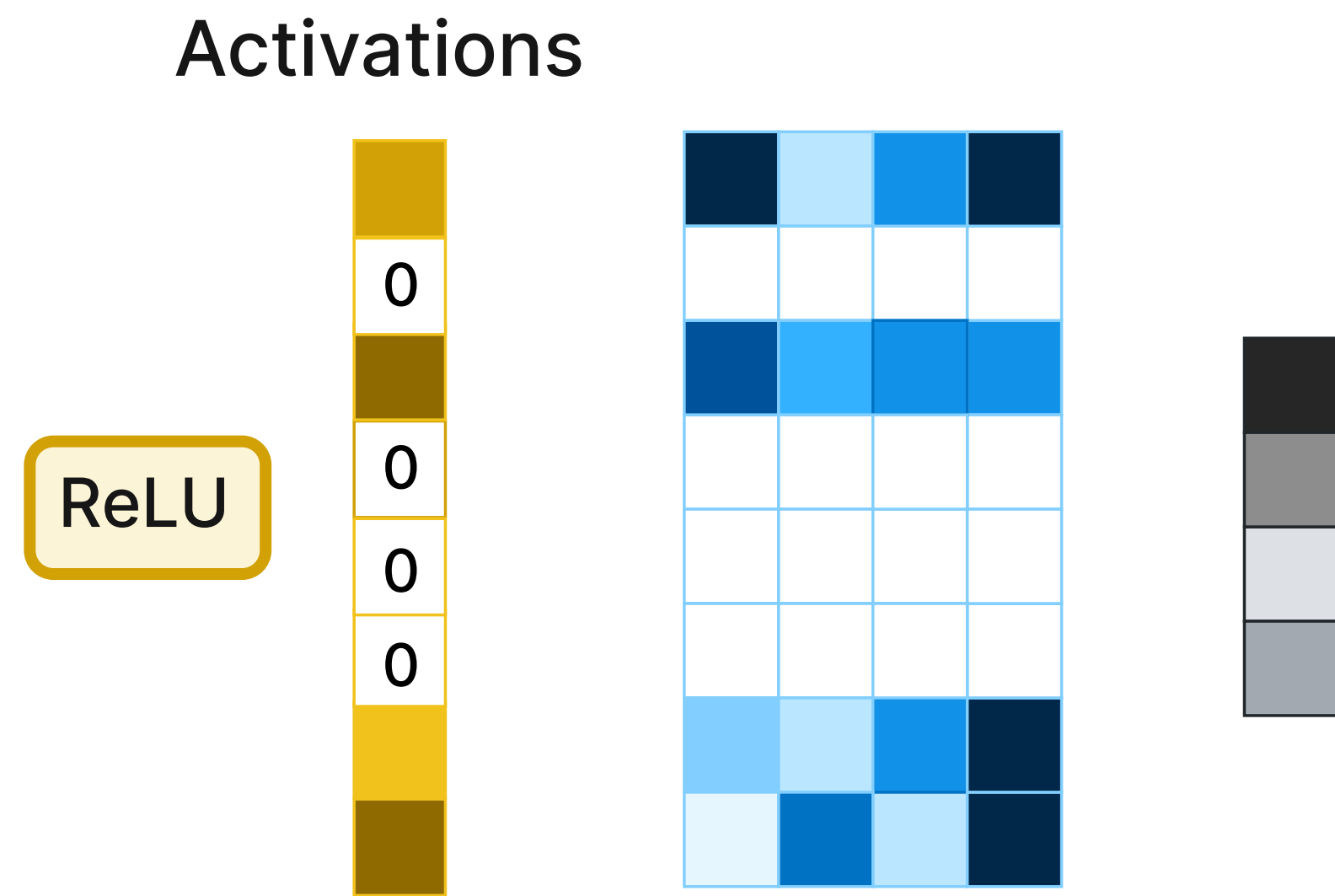
4d

d

d

ReLU

0

Activations

Do not need these weights **for this token only**

# Activation Sparsity: The Benefits of ReLU

Activations



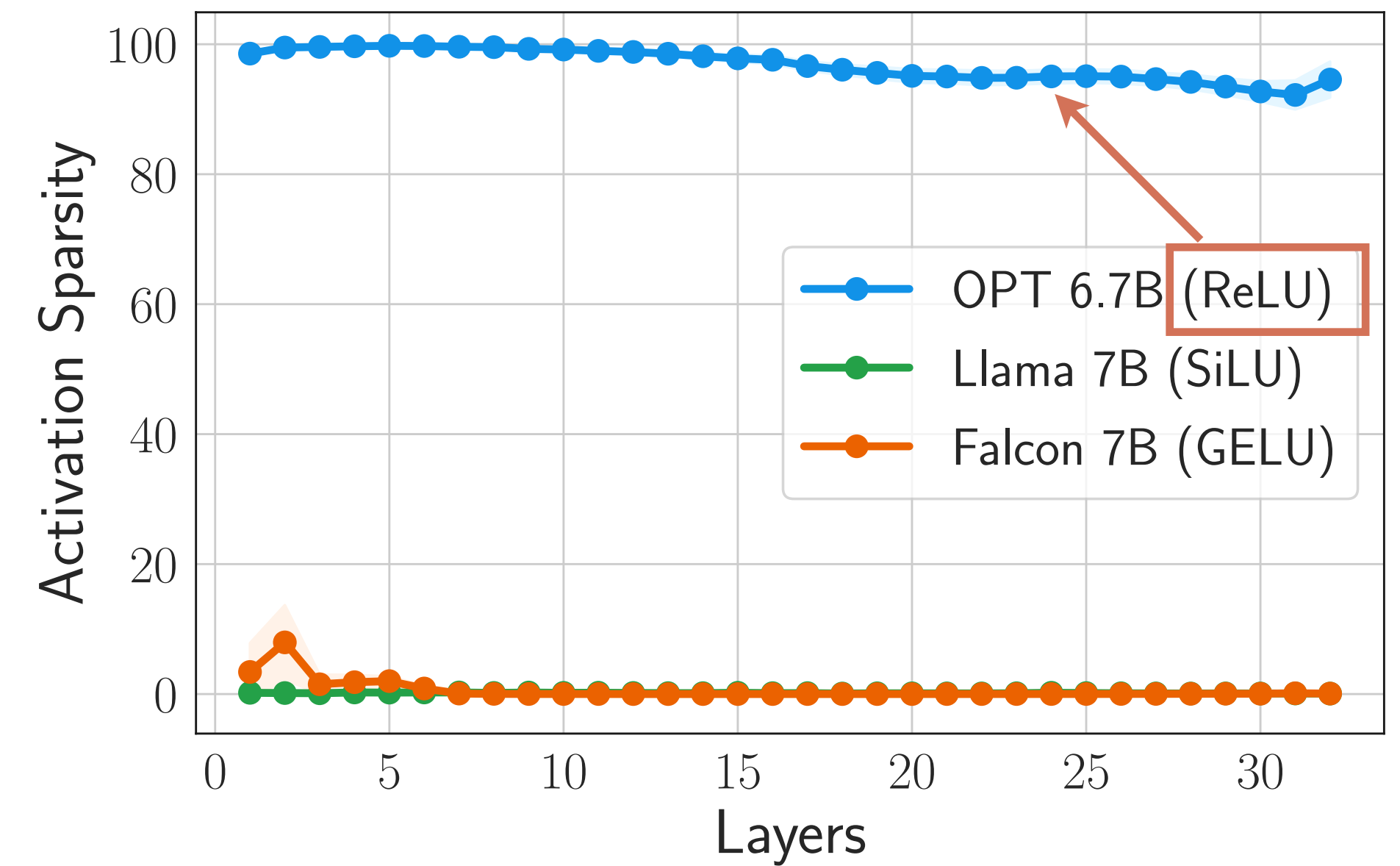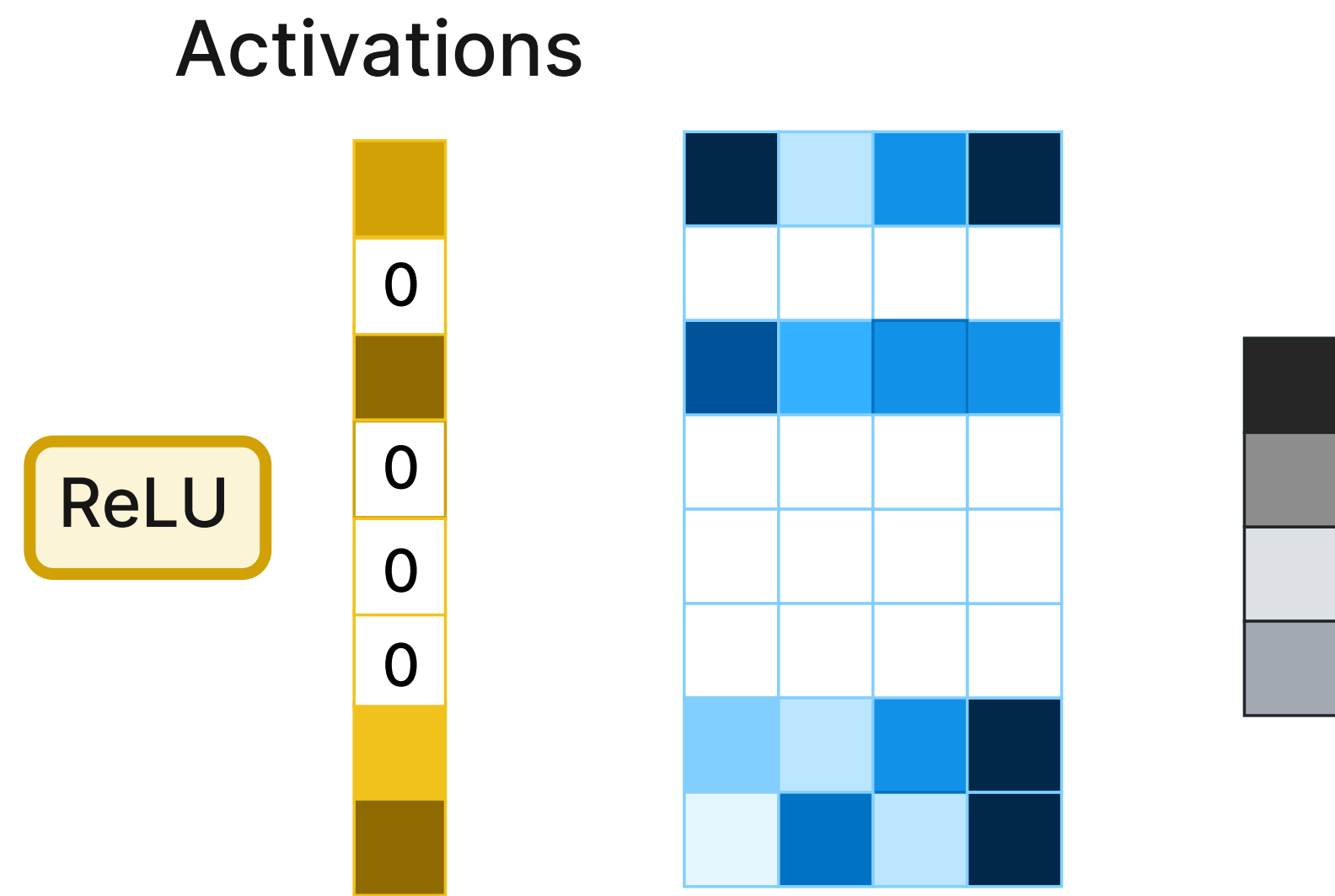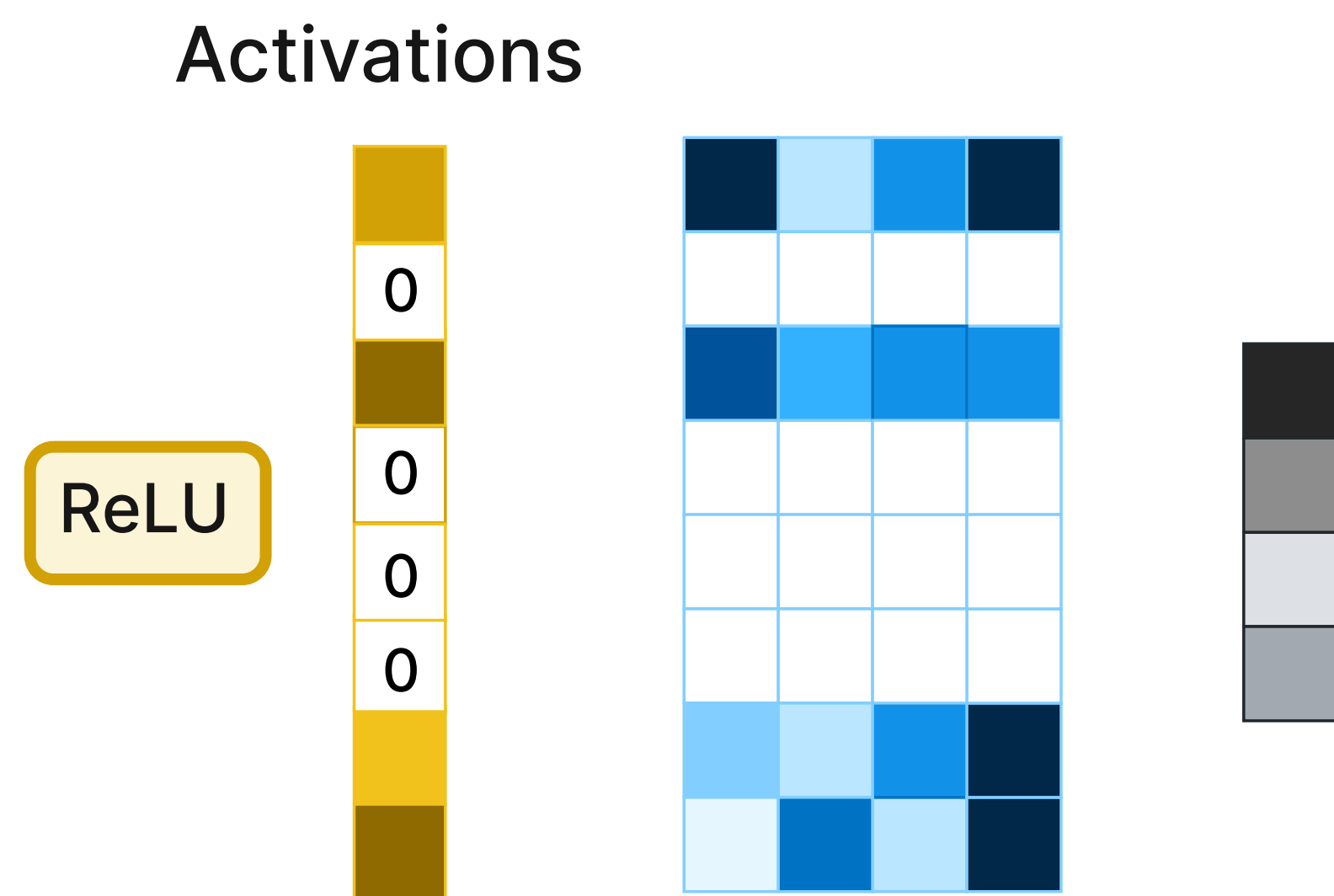ReLU

# Activation Sparsity: The Benefits of ReLU
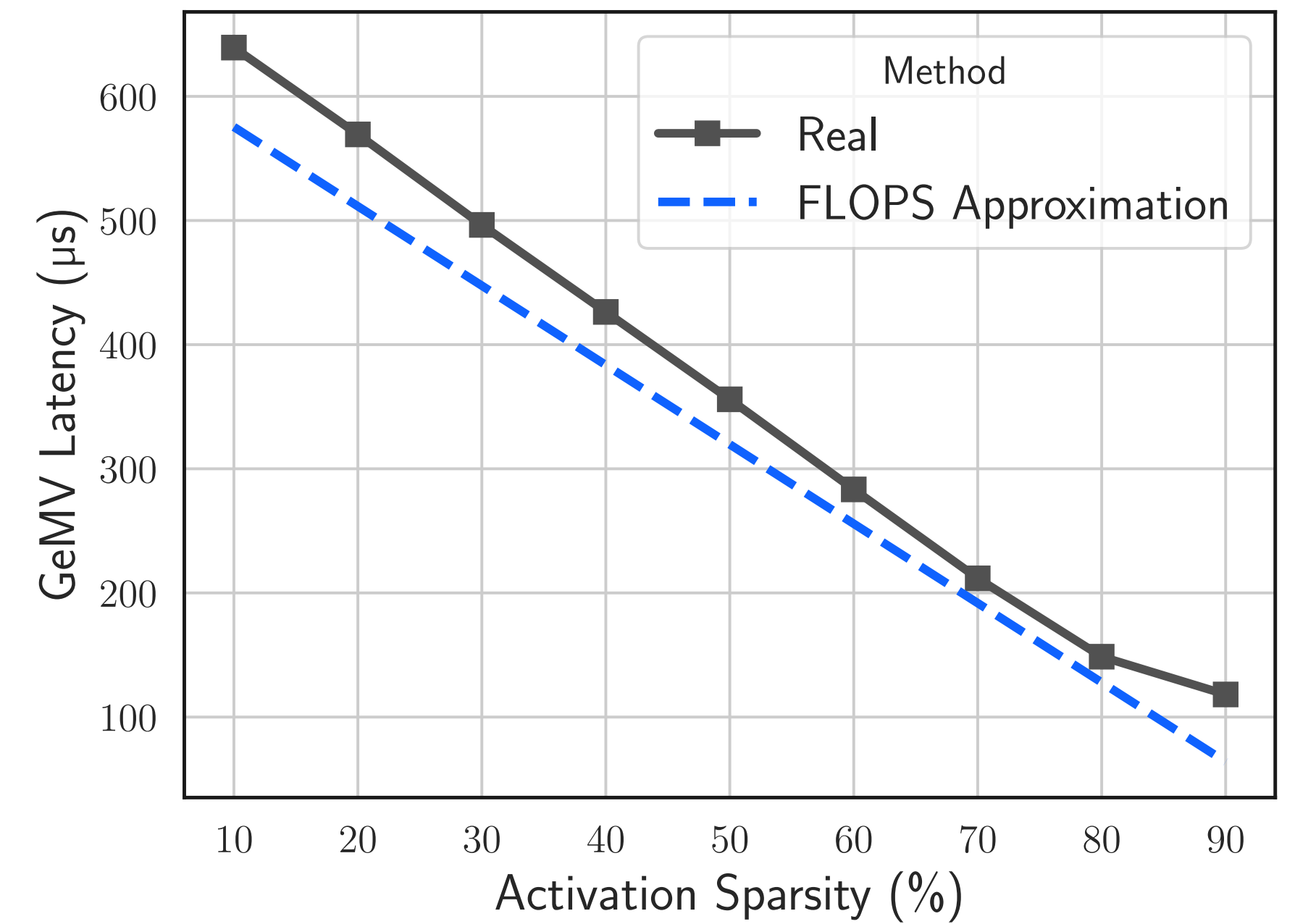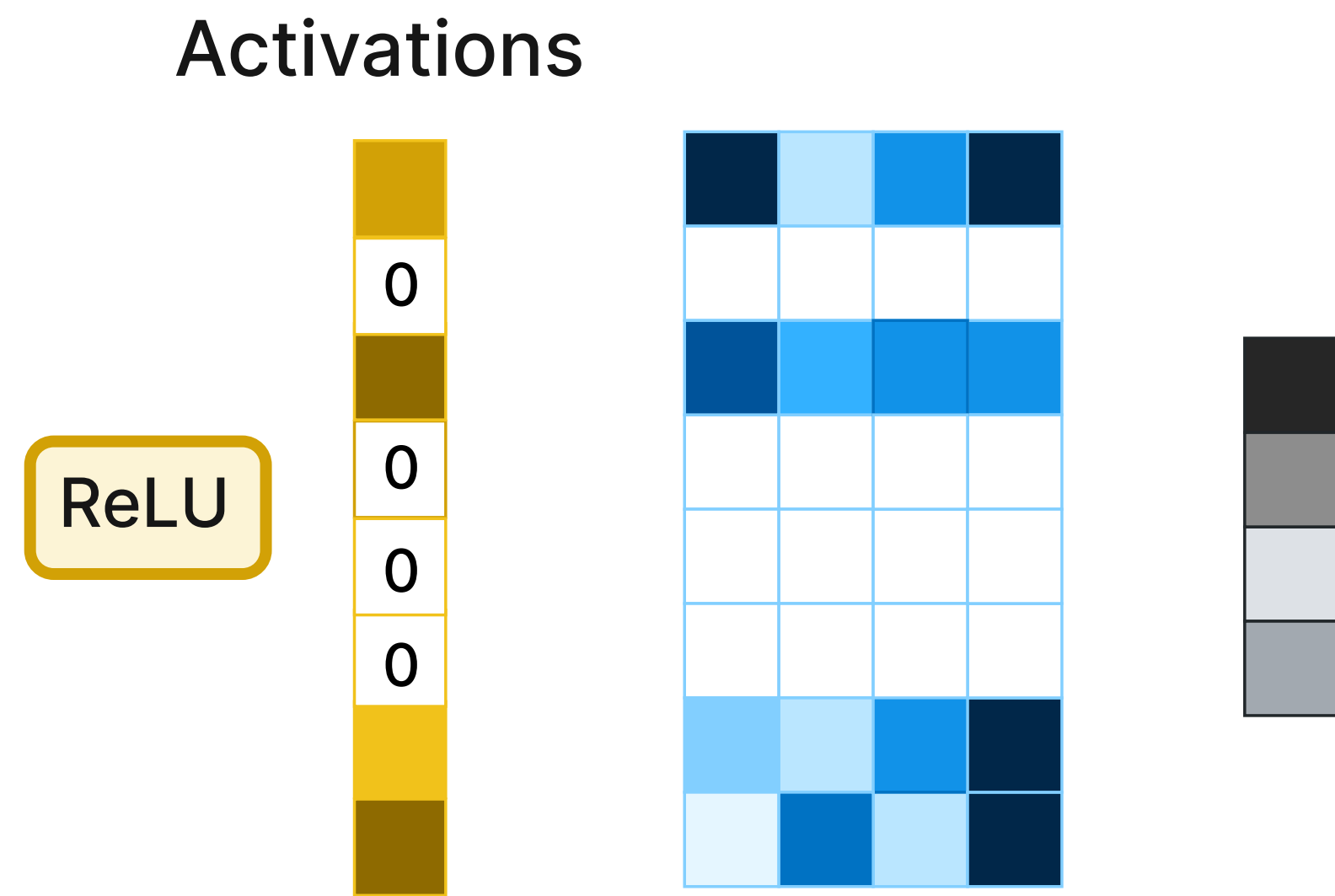


Activations

ReLU

# Activation Sparsity: The Benefits of ReLU

# Activation Sparsity: Efficient Implementation

Activations

# Activation Sparsity: Efficient Implementation
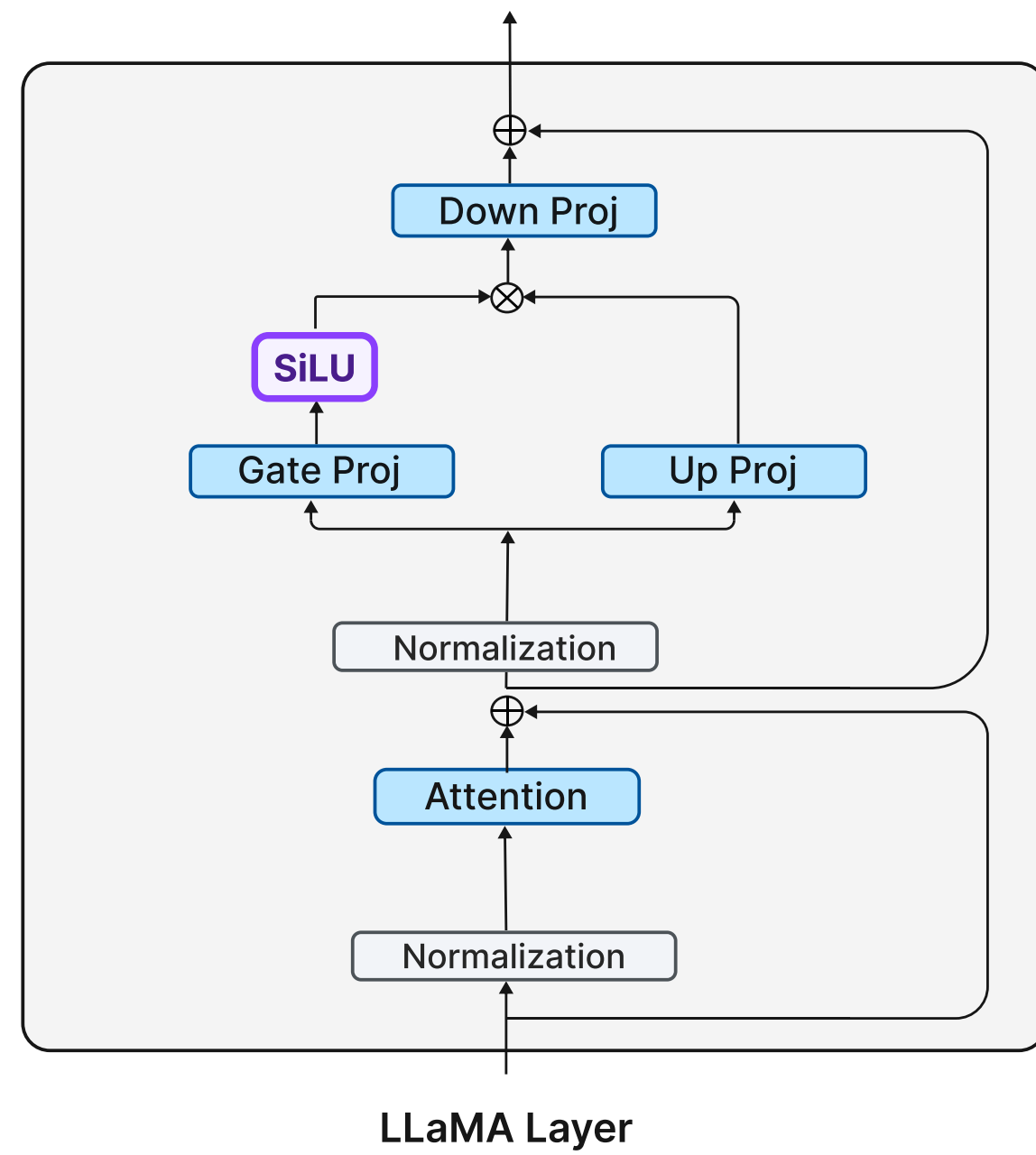
Activations



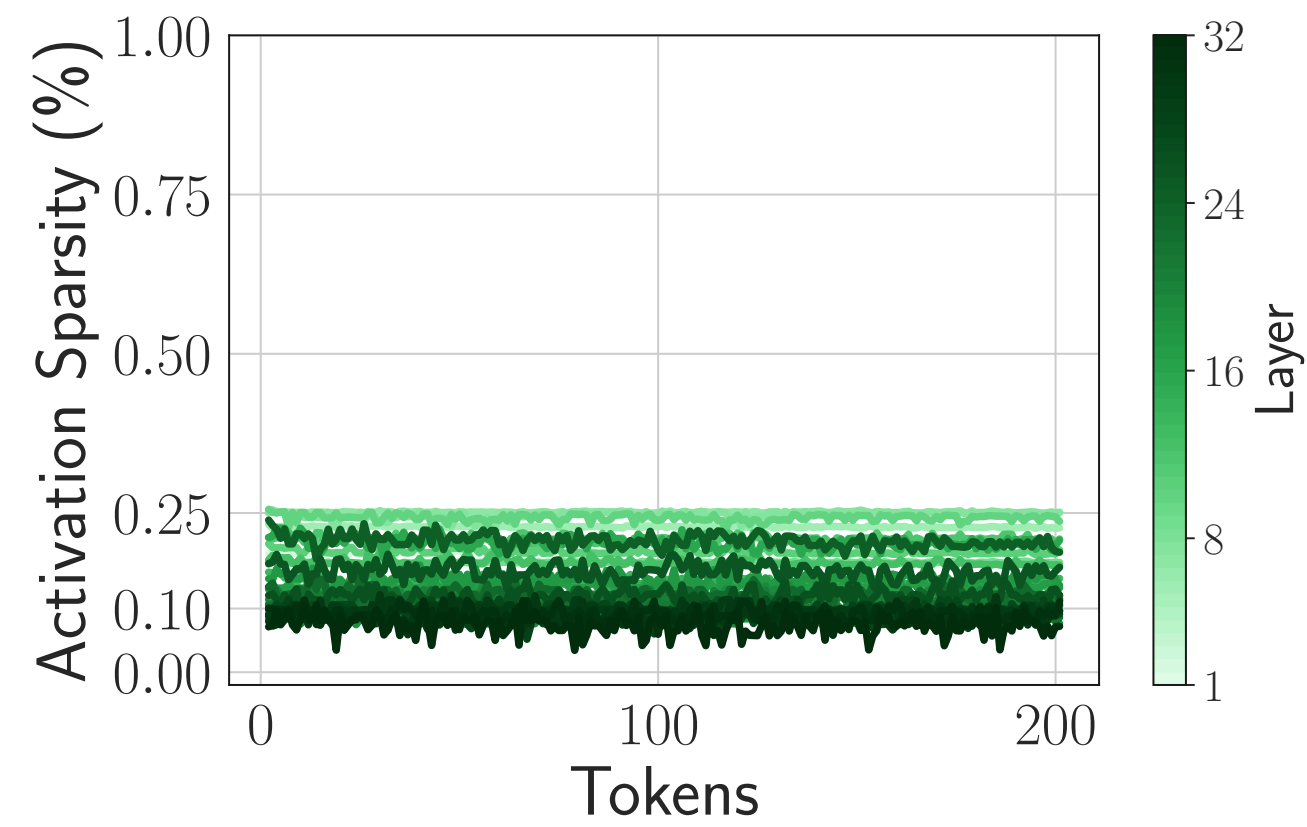Efficient Matrix-Vector product Metal
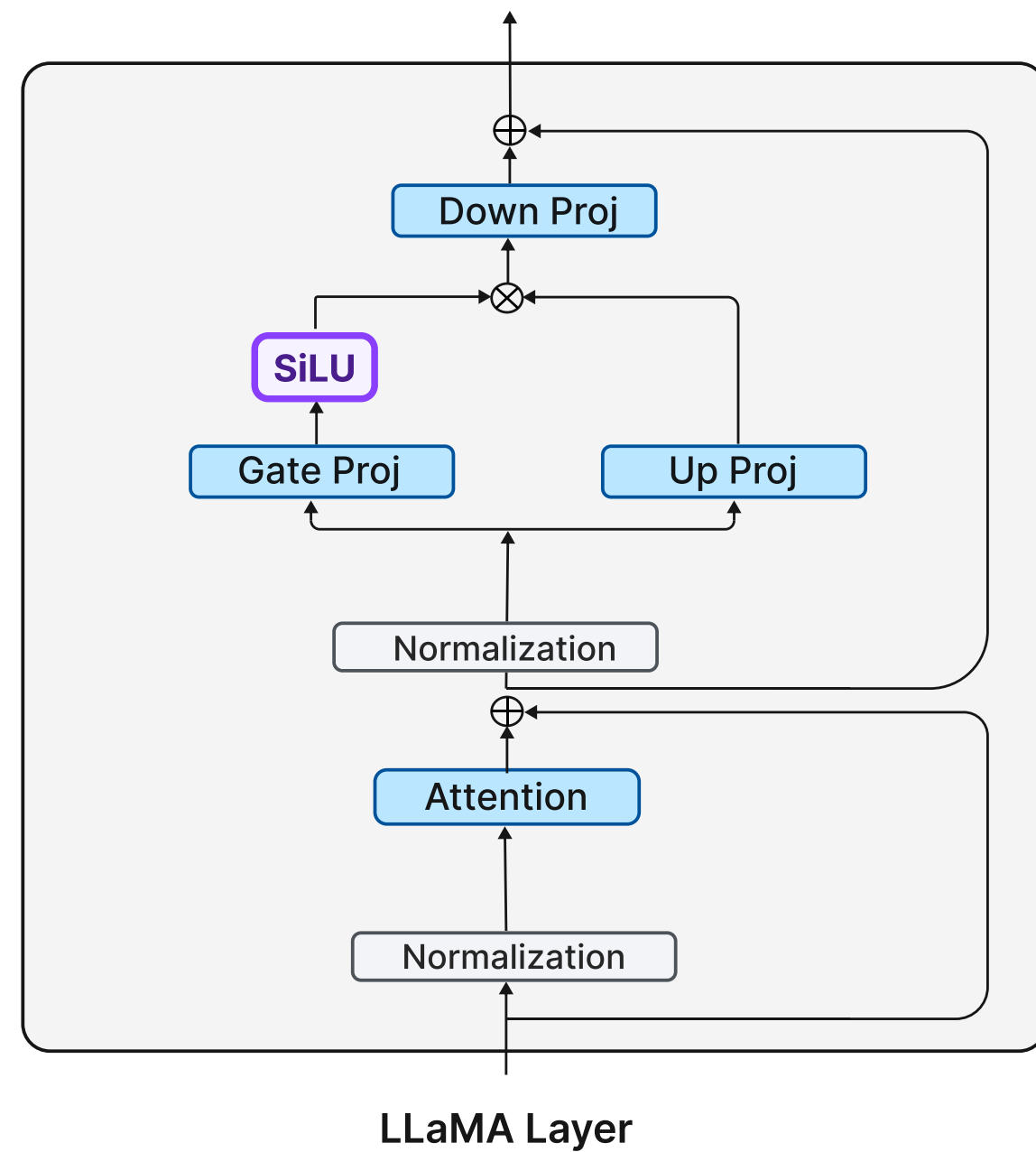Kernel on M2-Macbook Pro

# But most of the LLMs are already pre-trained without ReLU

# ReLUfication: Bringing ReLU Back to non-ReLU Pre-Trained Models

# ReLUfication: Bringing ReLU Back to non-ReLU Pre-Trained Models



**LLaMA Layer**

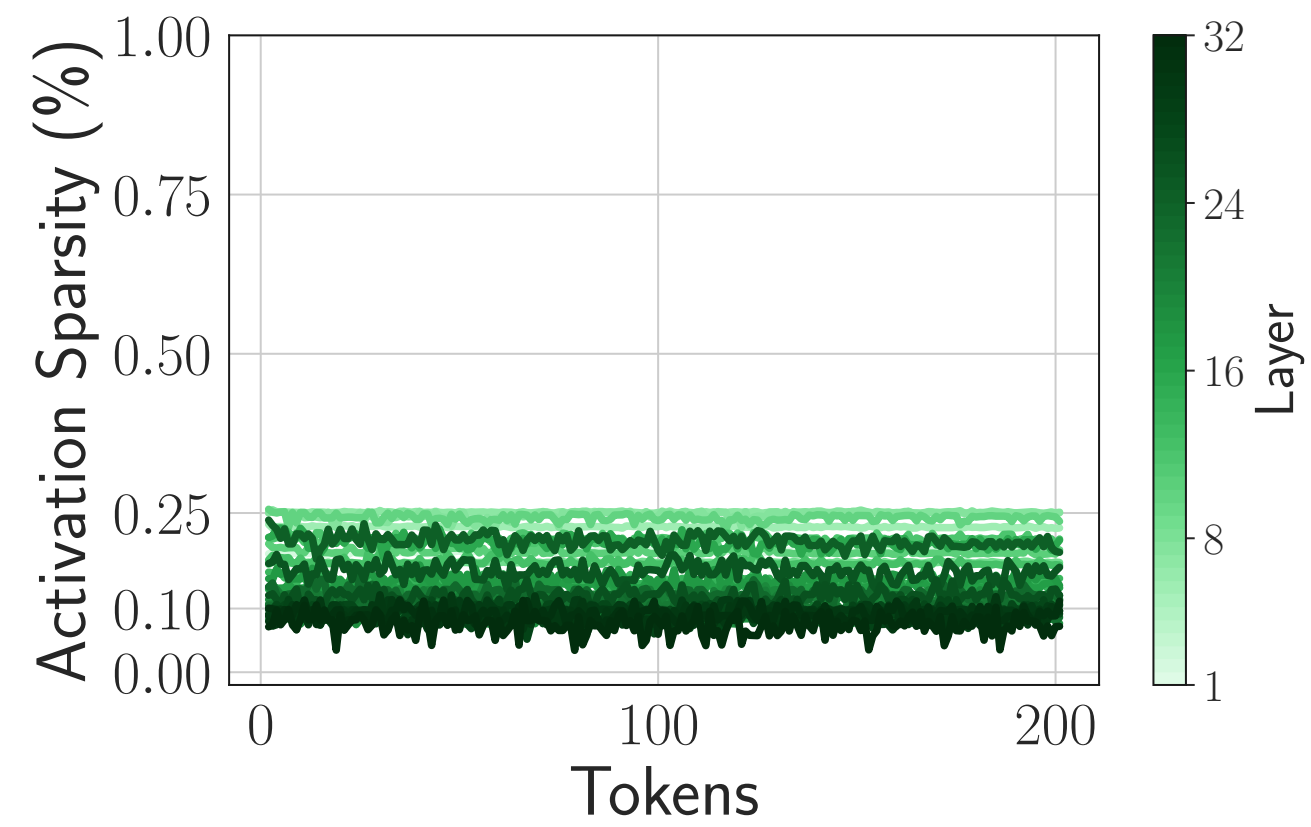# ReLUfication: Bringing ReLU Back to non-ReLU Pre-Trained Models



LLaMA Layer

# ReLUfication: Bringing ReLU Back to non-ReLU Pre-Trained Models
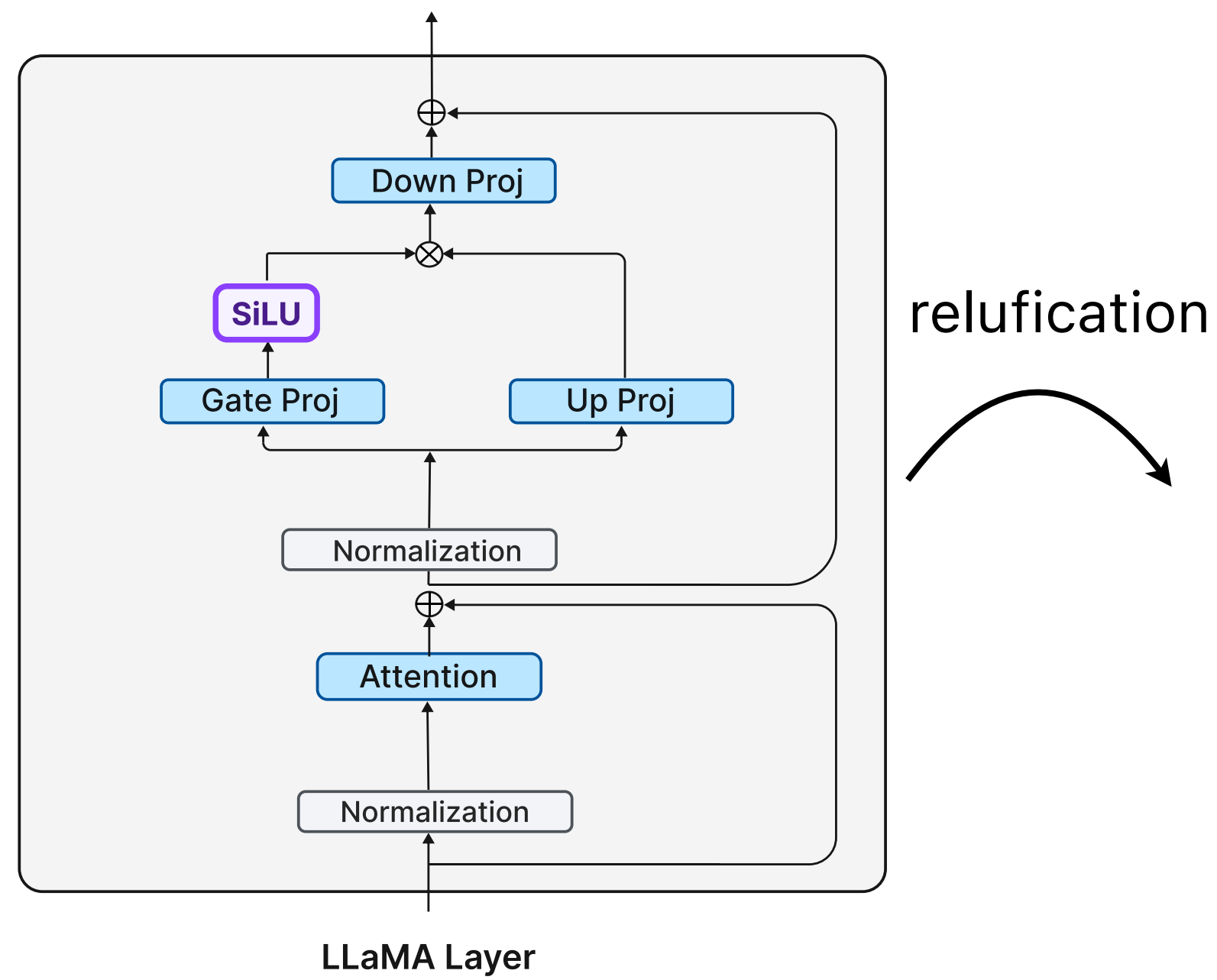


LLaMA Layer

relufication

# ReLUfication: Bringing ReLU Back to non-ReLU Pre-Trained Models

# ReLUfication: Bringing ReLU Back to non-ReLU Pre-Trained Models

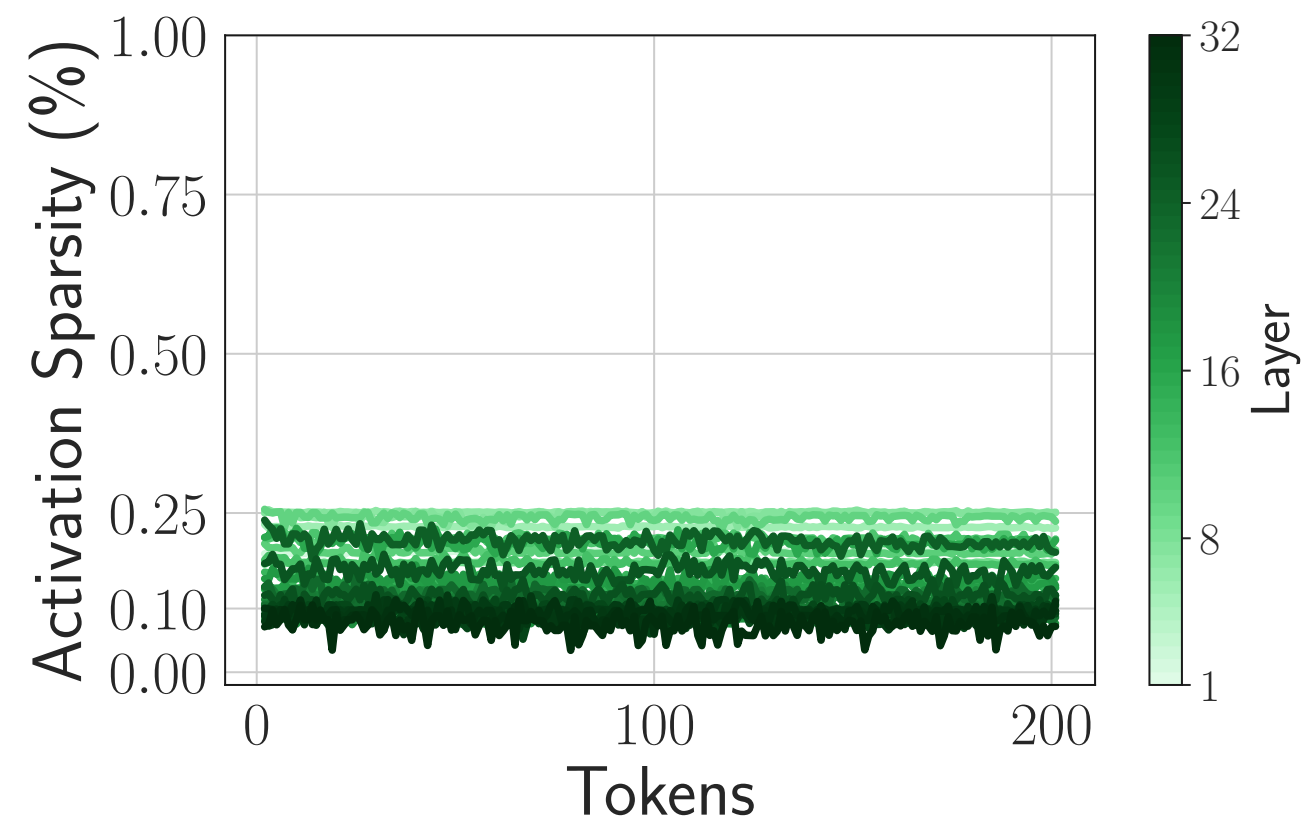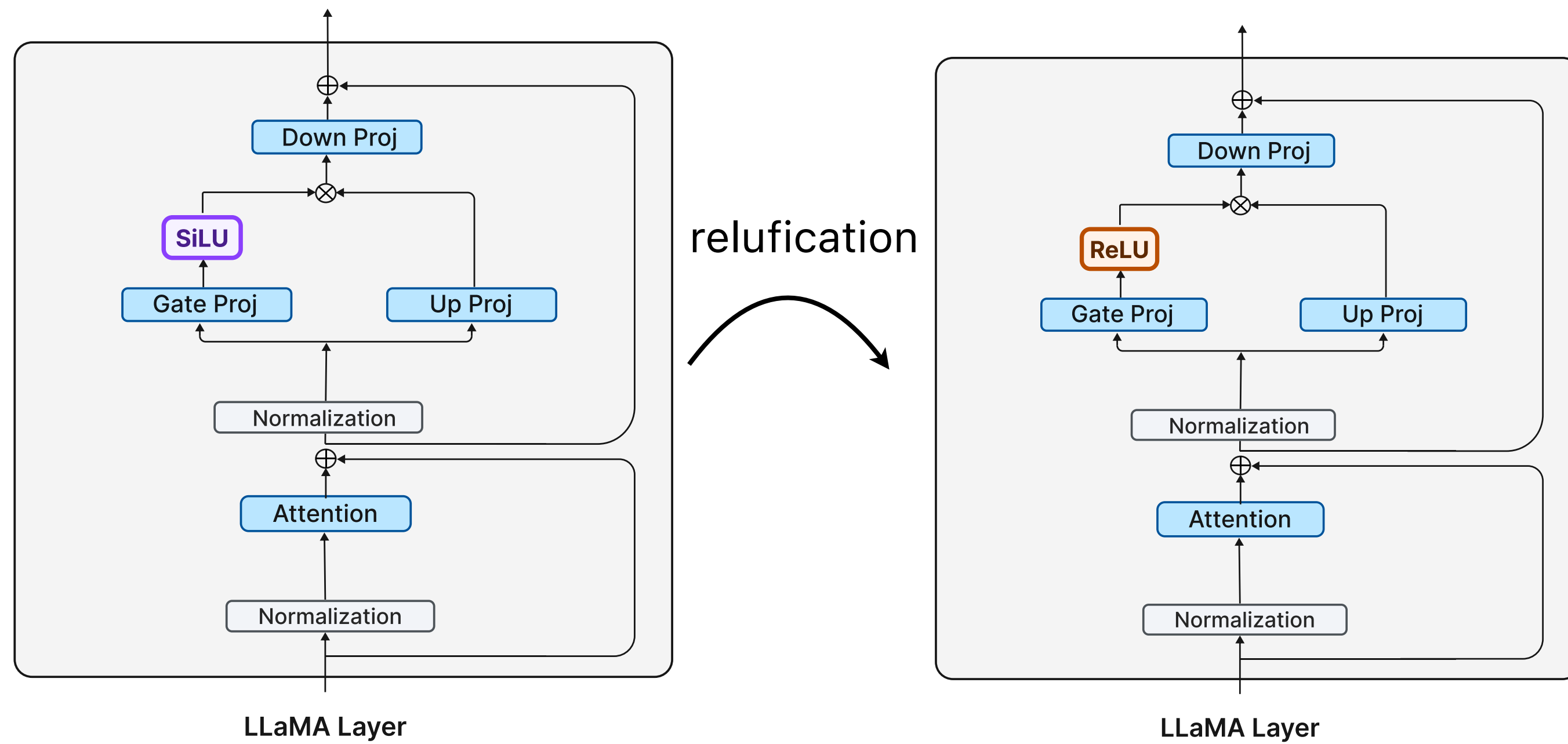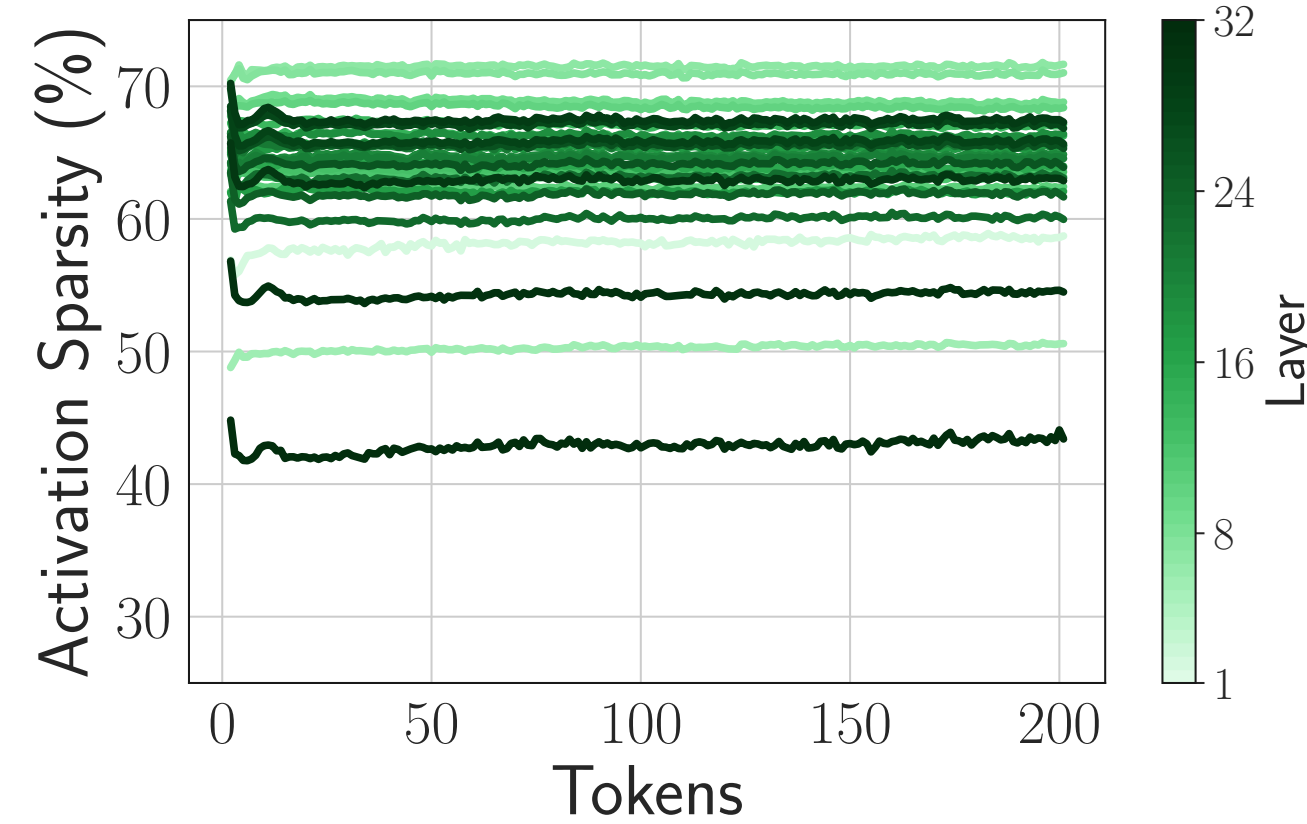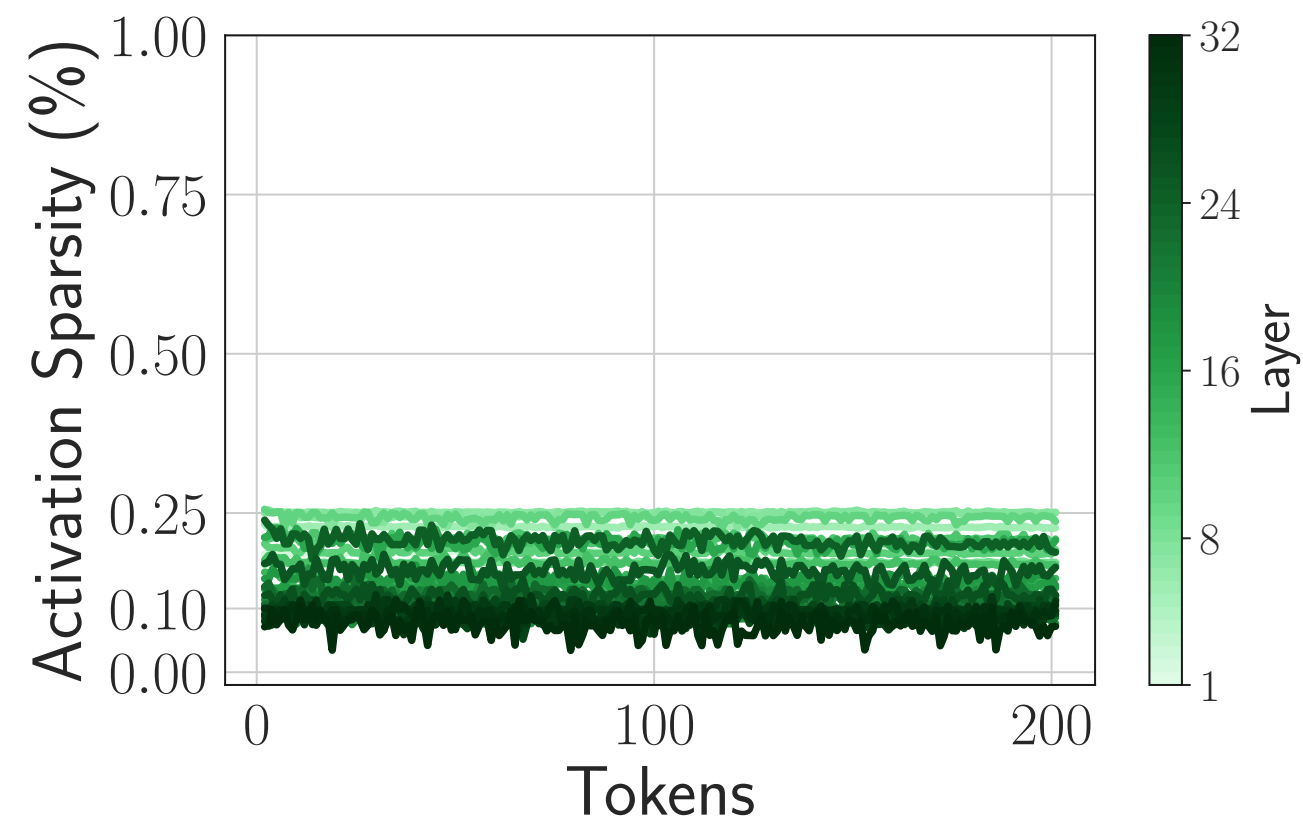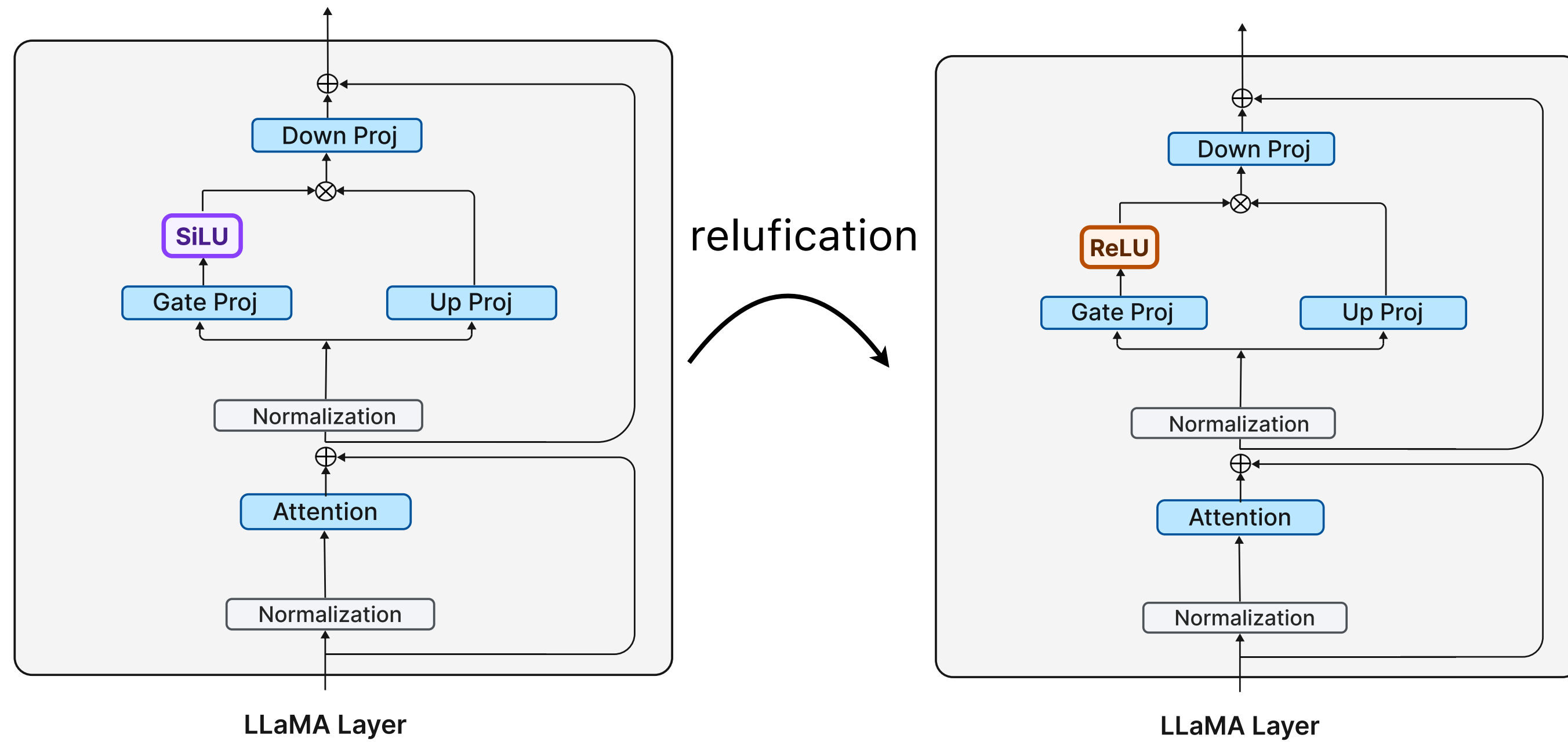# ReLUfication: Bringing ReLU Back to non-ReLU Pre-Trained Models

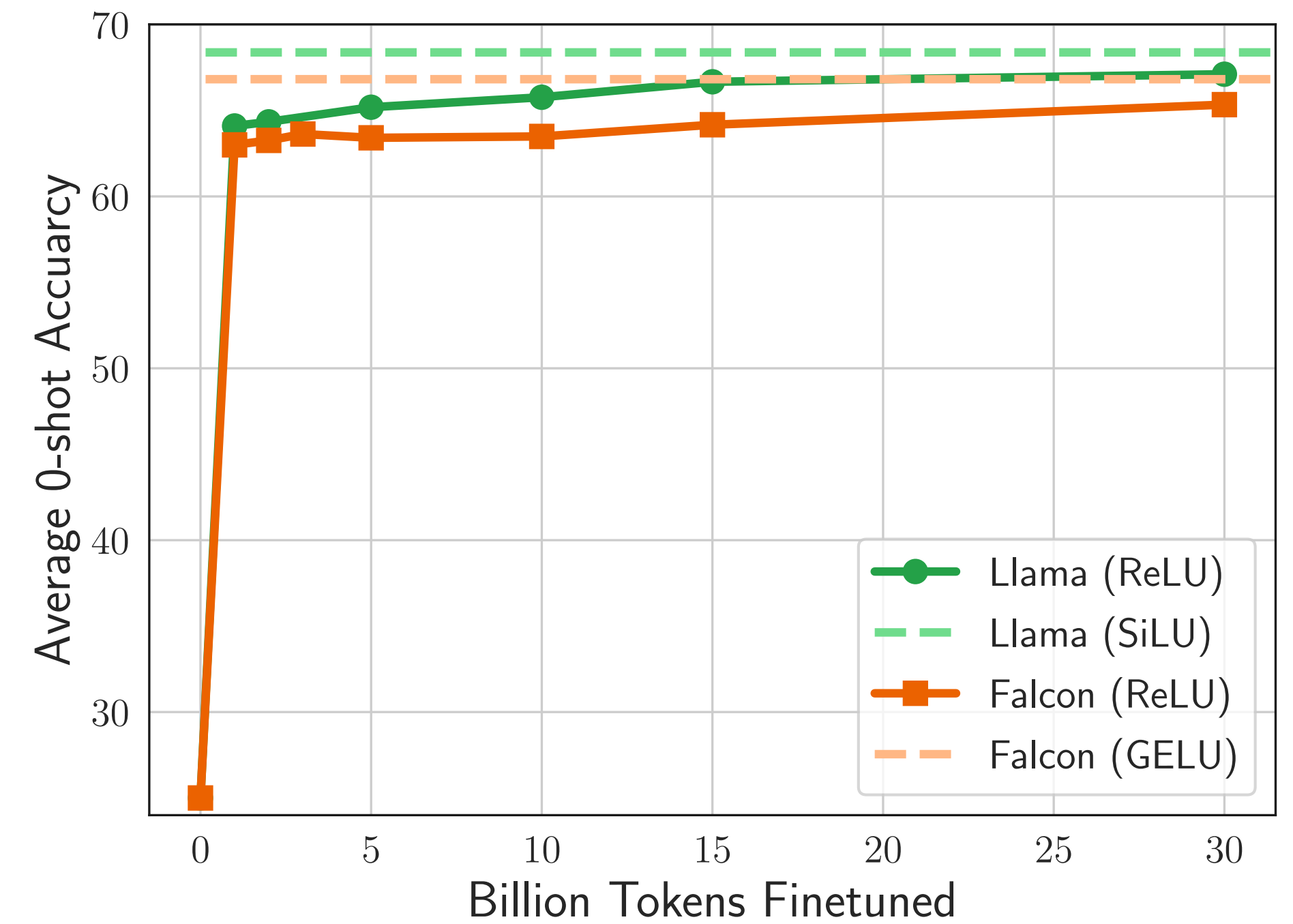# ReLUfication Stages



no relufication

relufication - stage 1

# ReLUfication Stages



**no relufication**

**relufication - stage 1**

sparse input

# ReLUfication Stages



no relufication

relufication - stage 1

relufication - stage 2

# ReLUfication Stages



no relufication

relufication - stage 1

relufication - stage 2

# ReLUfication Stages



no relufication

relufication - stage 1

relufication - stage 2

# ReLUfication Stages



no relufication

relufication - stage 1

relufication - stage 2

# ReLUfication Stages



no relufication

relufication - stage 1

relufication - stage 2

10

# ReLUfication: Results

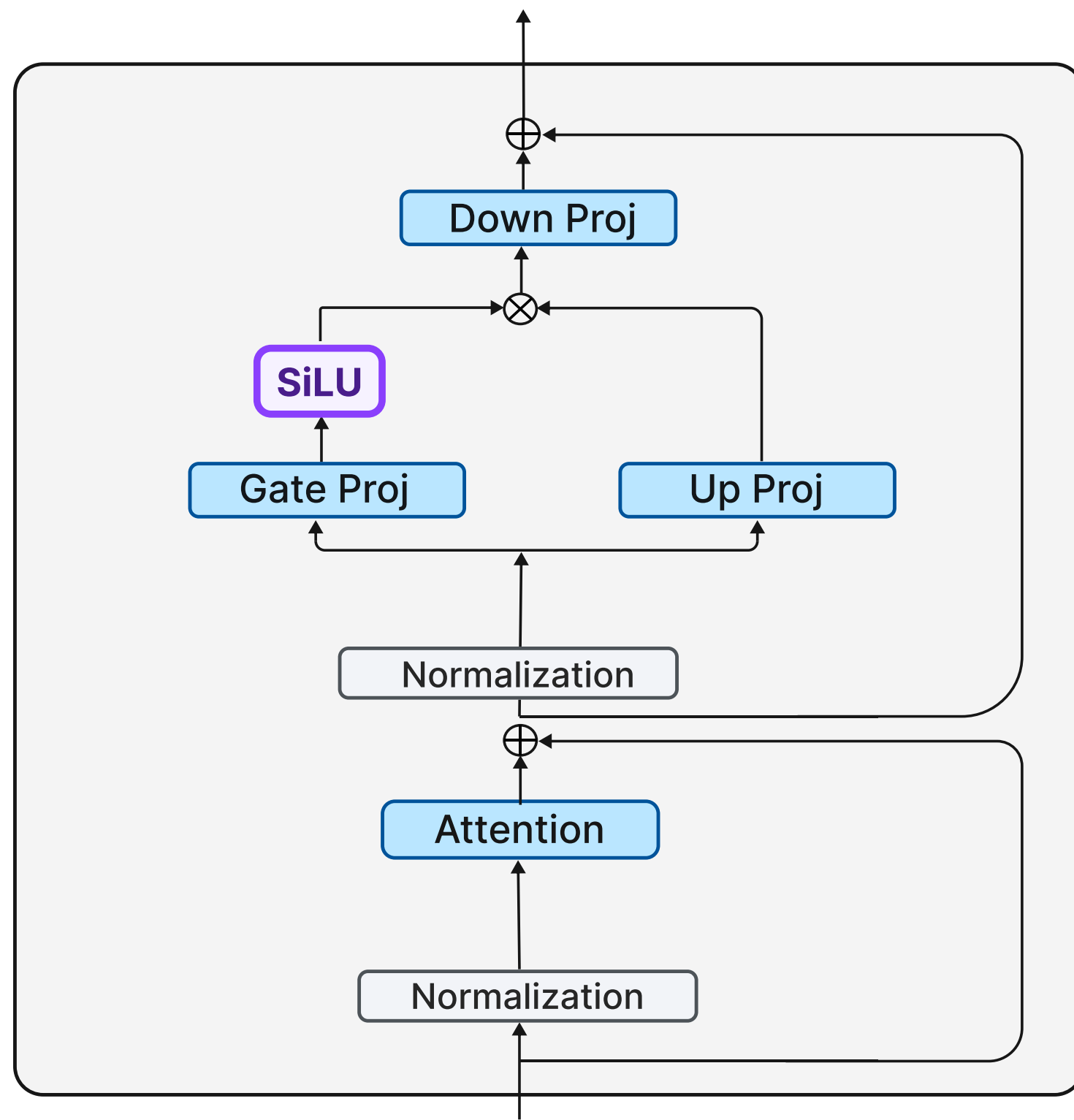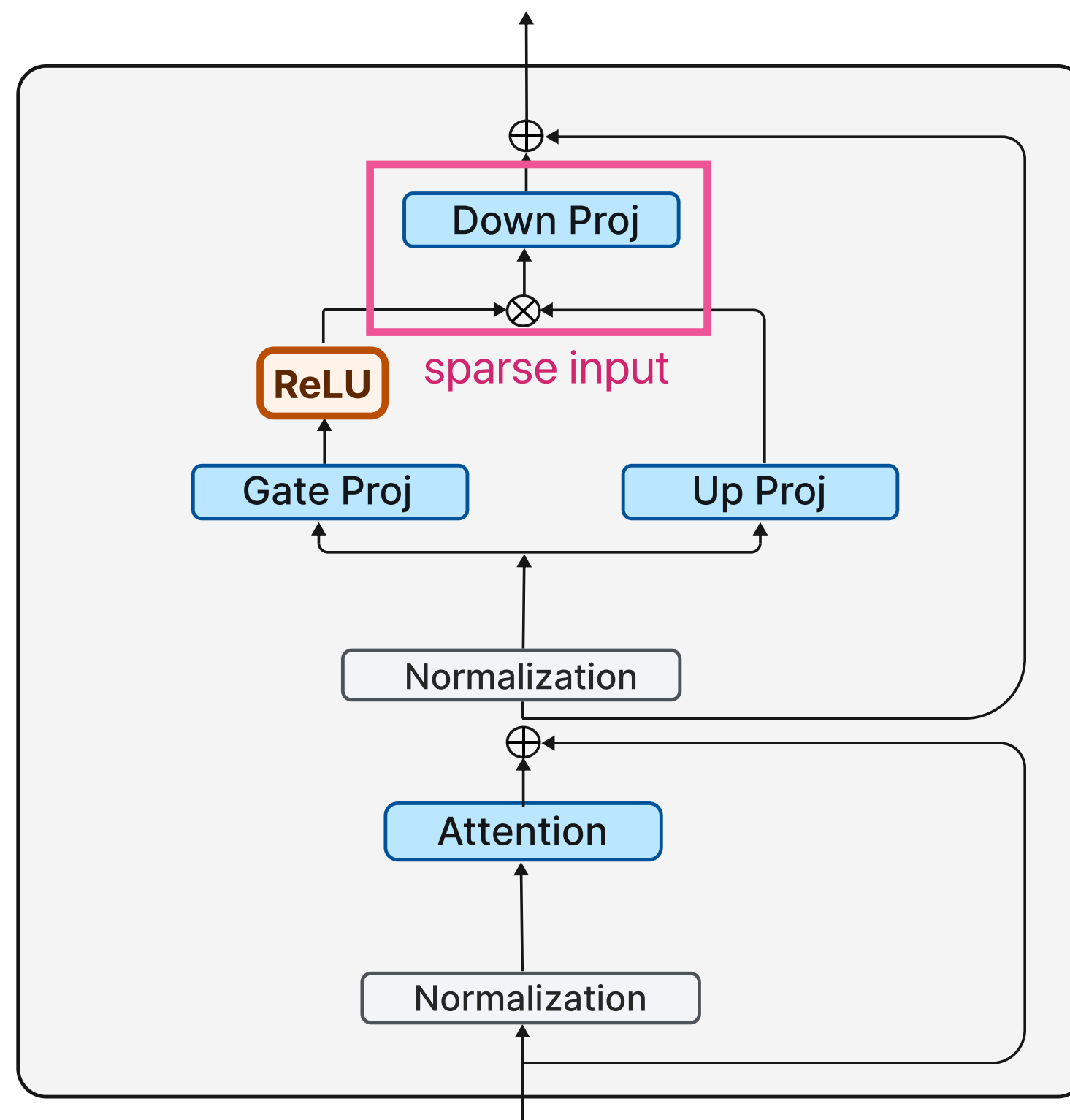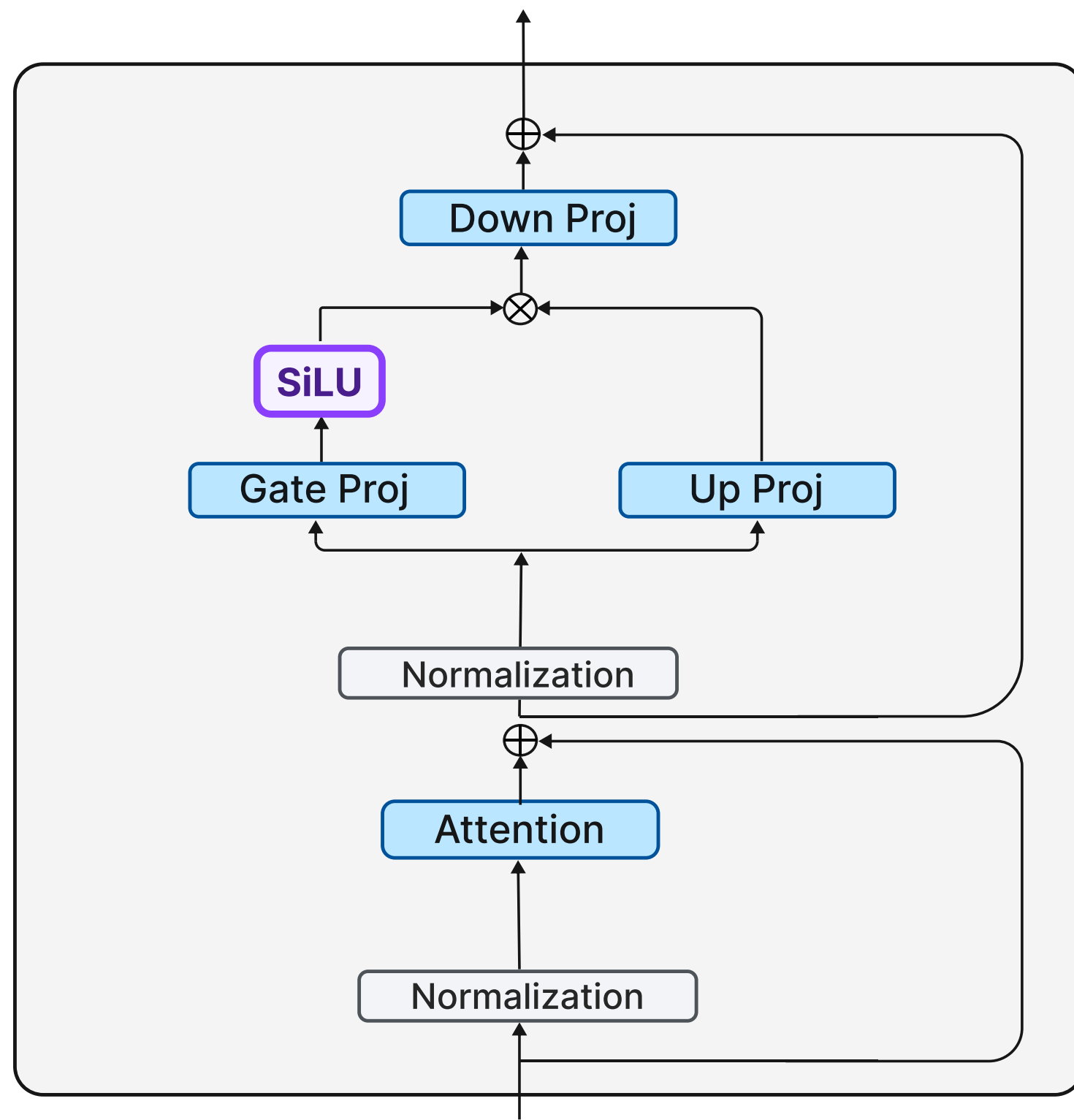| Model | Input Sparsity (%) | | | FLOPS (G) | Avg 0-shot Acc % |
|---|---|---|---|---|---|
| | QKV | UpProj | DownProj | | |
| Llama 7B | 0 | 0 | 0 | 6.6 | 68.4 |
| Llama 7B (relufied-s1) | 0 | 0 | 62 | 4.8 | 67.1 |
| Llama 7B (relufied-s2) | 51 | 67 | 65 | 2.9 | 66.4 |
| Falcon 7B | 0 | 1 | 0 | 6.6 | 66.8 |
| Falcon 7B (relufied-s1) | 0 | 0 | 94 | 4.1 | 65.2 |
| Falcon 7B (relufied-s2) | 56 | 56 | 95 | 2.2 | 64.8 |
| OPT 6.7B | 0 | 0 | 97 | 4.5 | 59.8 |
| OPT 6.7B (relufied-s2) | 50 | 40 | 97 | 2.8 | 58.6 |

# ReLUfication: Results

| Model | Input Sparsity (%) | | | FLOPS (G) | Avg 0-shot Acc % |
|---|---|---|---|---|---|
| | QKV | UpProj | DownProj | | |
| Llama 7B | 0 | 0 | 0 | 6.6 | 68.4 |
| Llama 7B (relufied-s1) | 0 | 0 | 62 | 4.8 | 67.1 |
| Llama 7B (relufied-s2) | 51 | 67 | 65 | 2.9 | 66.4 |
| Falcon 7B | 0 | 1 | 0 | 6.6 | 66.8 |
| Falcon 7B (relufied-s1) | 0 | 0 | 94 | 4.1 | 65.2 |
| Falcon 7B (relufied-s2) | 56 | 56 | 95 | 2.2 | 64.8 |
| OPT 6.7B | 0 | 0 | 97 | 4.5 | 59.8 |
| OPT 6.7B (relufied-s2) | 50 | 40 | 97 | 2.8 | 58.6 |

# Why aren't we training our LLMs with ReLU?

# The Impact of Activation Function on Performance

When trained from scratch...

# The Impact of Activation Function on Performance

When trained from scratch…

$$\mathrm{act}(x) = x\sigma(\beta x)$$

# The Impact of Activation Function on Performance

When trained from scratch...

$$\text{act}(x) = x\sigma(\beta x)$$
$$\text{SiLU}(x) = x\sigma(x)$$
$$\text{GELU}(x) \approx x\sigma(1.7x)$$

# The Impact of Activation Function on Performance

When trained from scratch…

Activation Function



$$\text{act}(x) = x\sigma(\beta x)$$
$$\text{SiLU}(x) = x\sigma(x)$$
$$\text{GELU}(x) \approx x\sigma(1.7x)$$

# The Impact of Activation Function on Performance

When trained from scratch...

Activation Function

Activation Sparsity



$$\mathrm{act}(x) = x\sigma(\beta x)$$
$$\mathrm{SiLU}(x) = x\sigma(x)$$
$$\mathrm{GELU}(x) \approx x\sigma(1.7x)$$

# The Impact of Activation Function on Performance

When trained from scratch...



Activation Function

Activation Sparsity

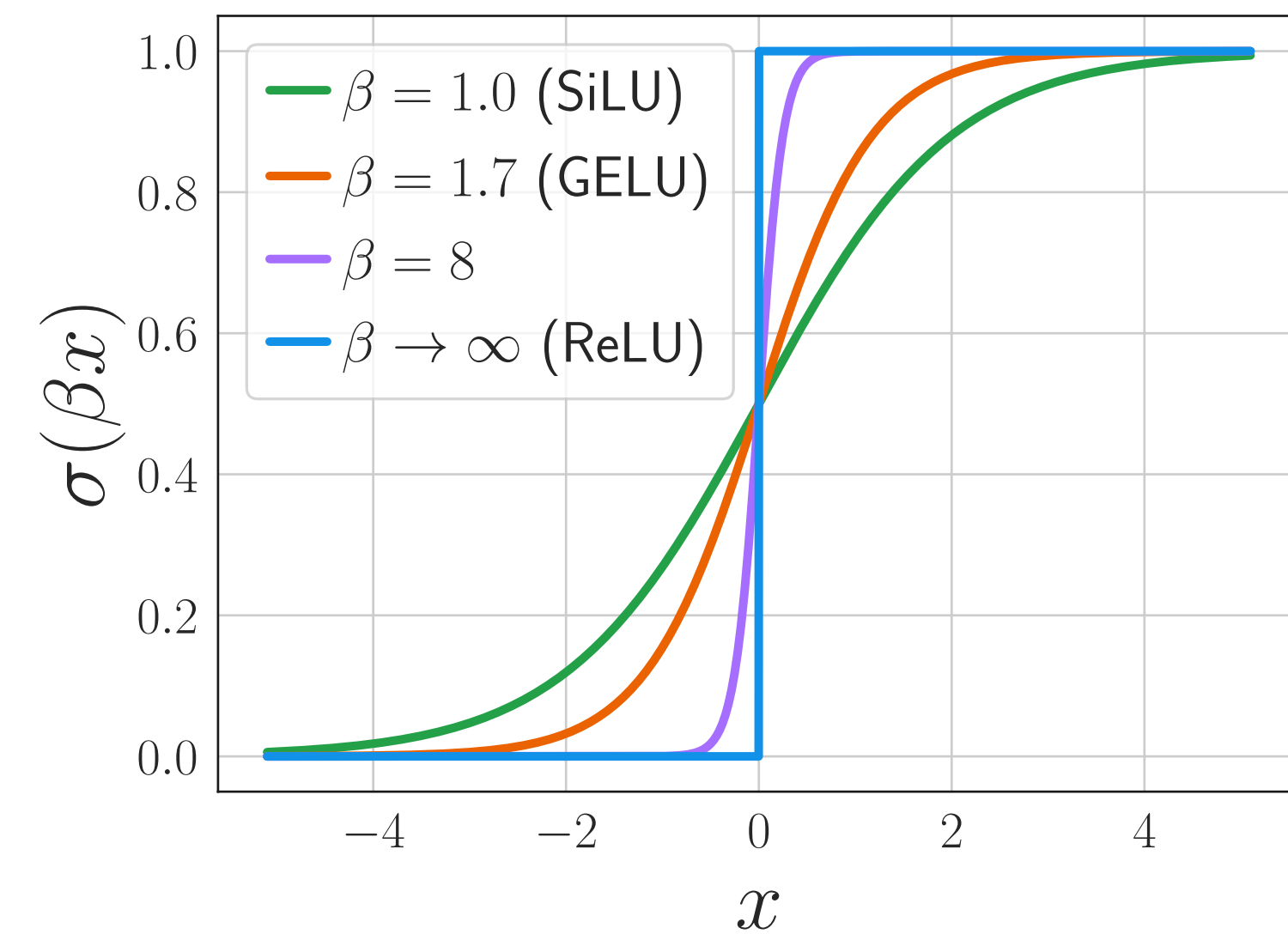Performance

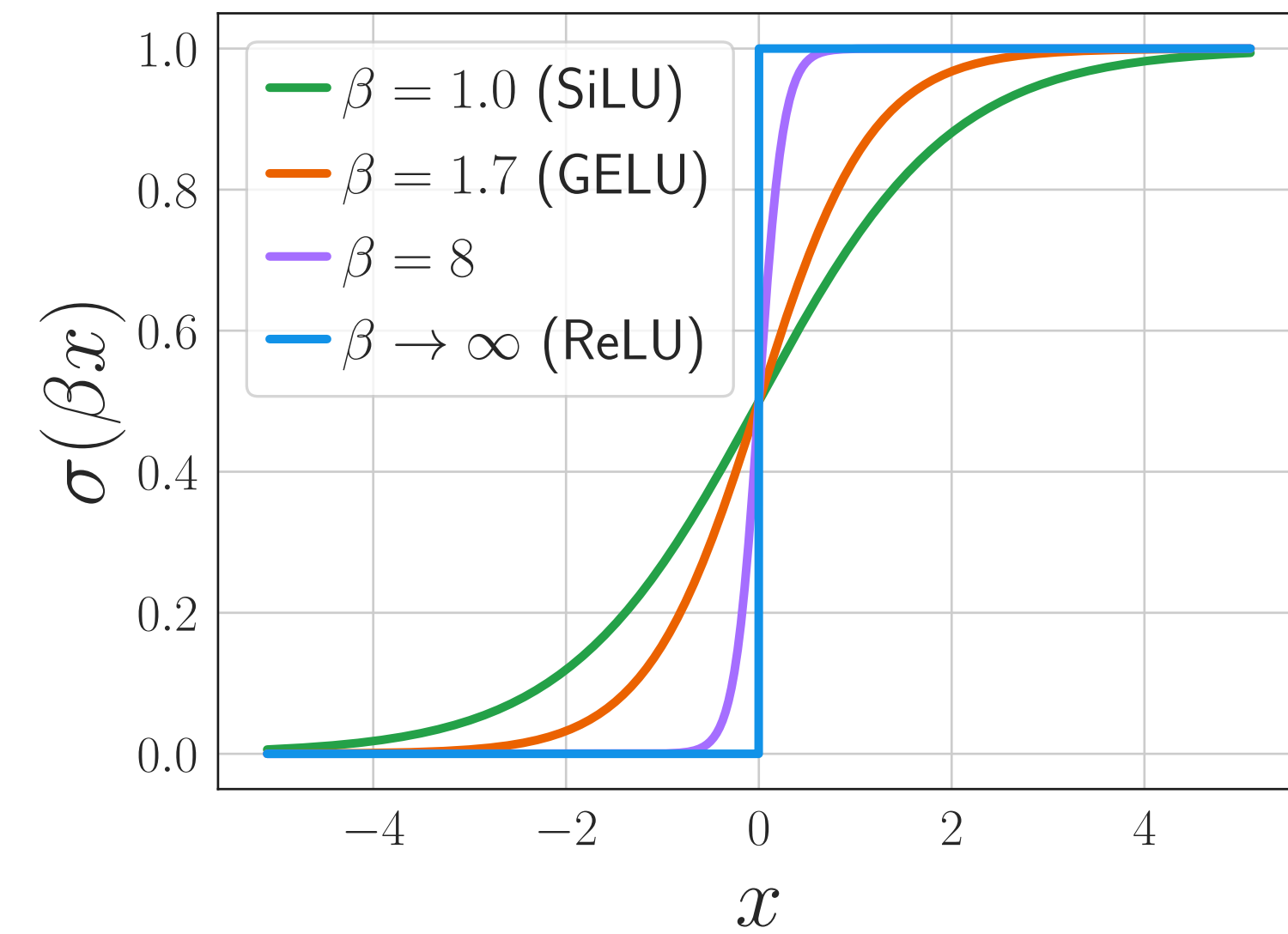$$\mathrm{act}(x) = x\sigma(\beta x)$$

$$\mathrm{SiLU}(x) = x\sigma(x)$$

$$\mathrm{GELU}(x) \approx x\sigma(1.7x)$$

# The Impact of Activation Function on Performance

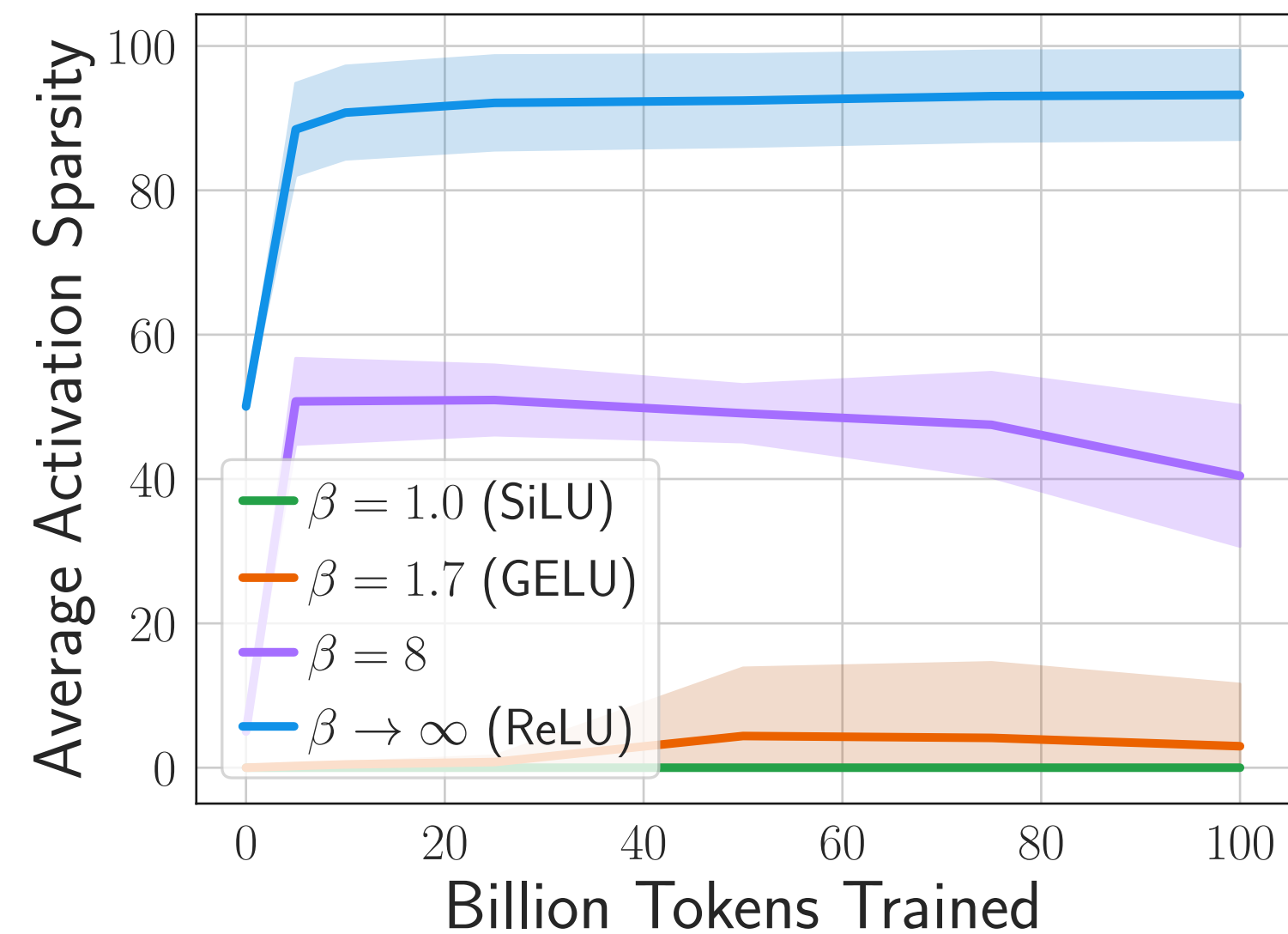When trained from scratch...

Activation Function

Activation Sparsity

Performance



$$\mathrm{act}(x) = x\sigma(\beta x)$$

$$\mathrm{SiLU}(x) = x\sigma(x)$$

$$\mathrm{GELU}(x) \approx x\sigma(1.7x)$$

# Concluding Remarks

# Recap

- ReLU activations are sparse. We can use this property for faster inference.

- The majority of LLMs are trained without ReLU.

- We can change their activation function to ReLU, and fine-tune them for a short time (ReLUfication).

- But do these alternative activation functions (e.g., GELU, SwiGLU) really improve performance?
  Maybe we should reconsider the choice of activation function.

# Future Directions: Multiple Tokens

# Future Directions: Multiple Tokens

Limitation of our work:
- The assumption that bsz = 1,
- We are generating one token at a time.

# Future Directions: Multiple Tokens

Limitation of our work:
- The assumption that bsz = 1,
- We are generating one token at a time.

It is difficult to exploit activation sparsity when we have beam search, speculative decoding, batched inference

# Future Directions: Multiple Tokens

Limitation of our work:
- The assumption that bsz = 1,
- We are generating one token at a time.

It is difficult to exploit activation sparsity when we have beam search, speculative decoding, batched inference

**Aggregated Sparsity**
How much of total available neurons have **not been used** for the first **t** tokens?
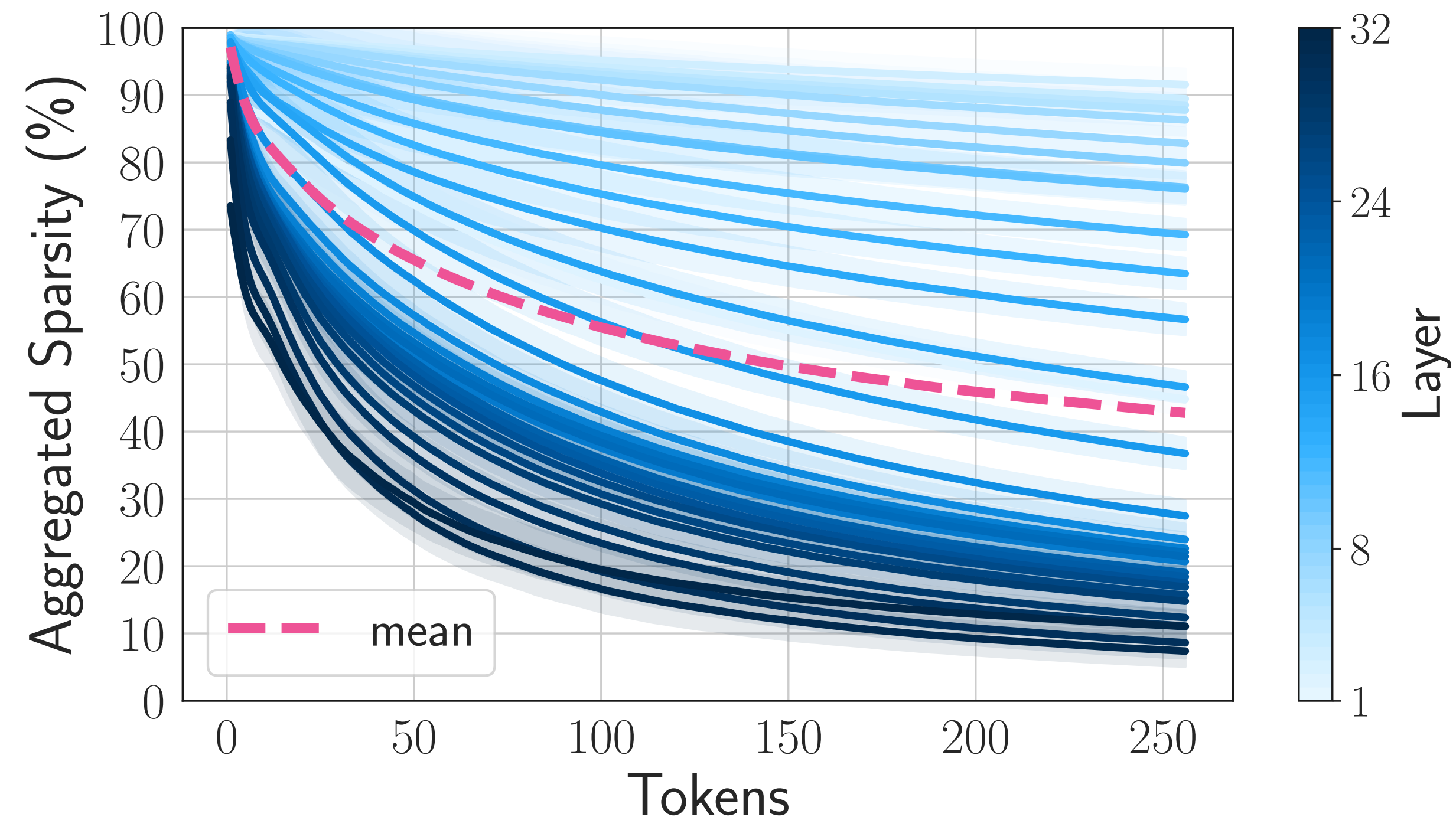
# Future Directions: Multiple Tokens

Limitation of our work:
- The assumption that bsz = 1,
- We are generating one token at a time.

It is difficult to exploit activation sparsity when we have beam search, speculative decoding, batched inference

**Aggregated Sparsity**
How much of total available neurons have **not been used** for the first **t** tokens?

# Future Directions: Improved Sparsity

For the GLU-based activation functions (e.g., SwiGLU), we notice lower activation sparsity (~60%) even after ReLUfication.

Tackling this issue:

[1] Song, Chenyang, et al. "ProSparse: Introducing and Enhancing Intrinsic Activation Sparsity within Large Language Models." arXiv preprint arXiv:2402.13516.

[2] Lee, Je-Yong, et al. "CATS: Contextually-Aware Thresholding for Sparsity in Large Language Models." arXiv:2404.08763.

Poster Session: Halle B #89