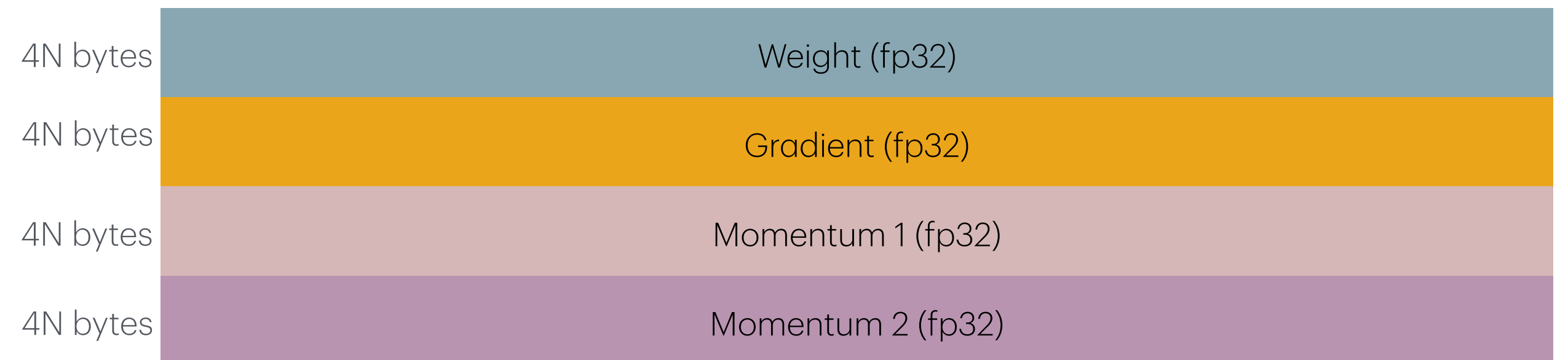# Low-rank adapters

Philipp Krähenbühl, UT Austin

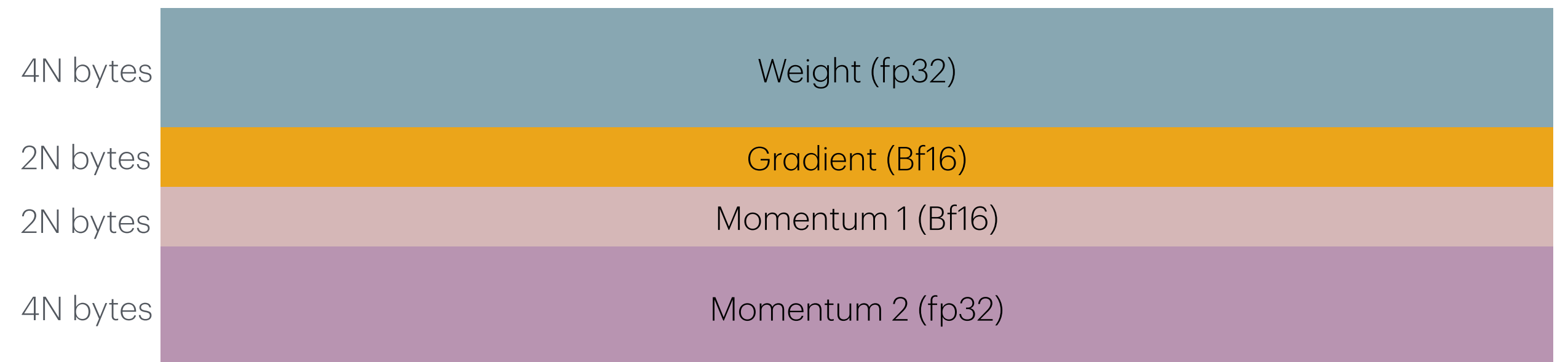# Training large models

## Memory requirements

- Without optimization:

  - Model parameters: N

  - Weights: N floats

  - Gradients: N floats

  - Momentum: N floats

  - 2nd momentum (ADAM): N floats

- 16N bytes without counting activations

| | |
|---|---|
| 4N bytes | Weight (fp32) |
| 4N bytes | Gradient (fp32) |
| 4N bytes | Momentum 1 (fp32) |
| 4N bytes | Momentum 2 (fp32) |

# Training large models

## Memory requirements

- Mixed precision

  - Model parameters: N

  - Weights: N floats

  - Gradients: N bfloat16

  - Momentum: N bfloat16

  - 2nd momentum (ADAM): N floats

- 12N bytes without counting activations

| | |
|---|---|
| 4N bytes | Weight (fp32) |
| 2N bytes | Gradient (Bf16) |
| 2N bytes | Momentum 1 (Bf16) |
| 4N bytes | Momentum 2 (fp32) |

# Training large models

## Memory requirements

FSDP

- Zero / FSDP

- 16N / M bytes without counting activations
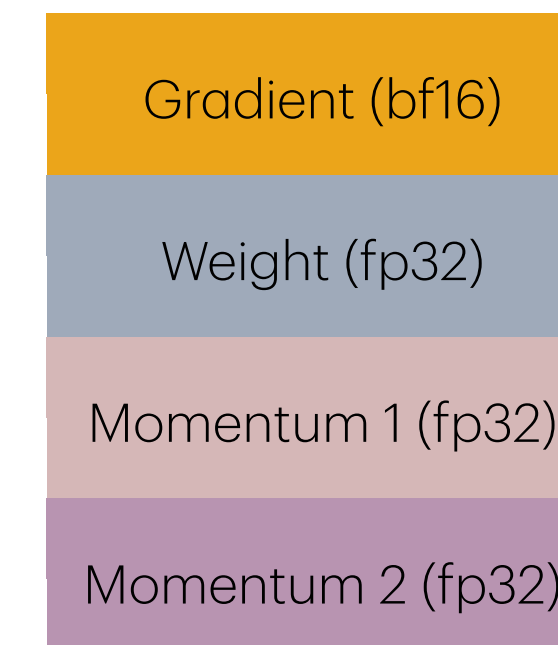
- For M GPUs

  - Good solution for GPU-rich people

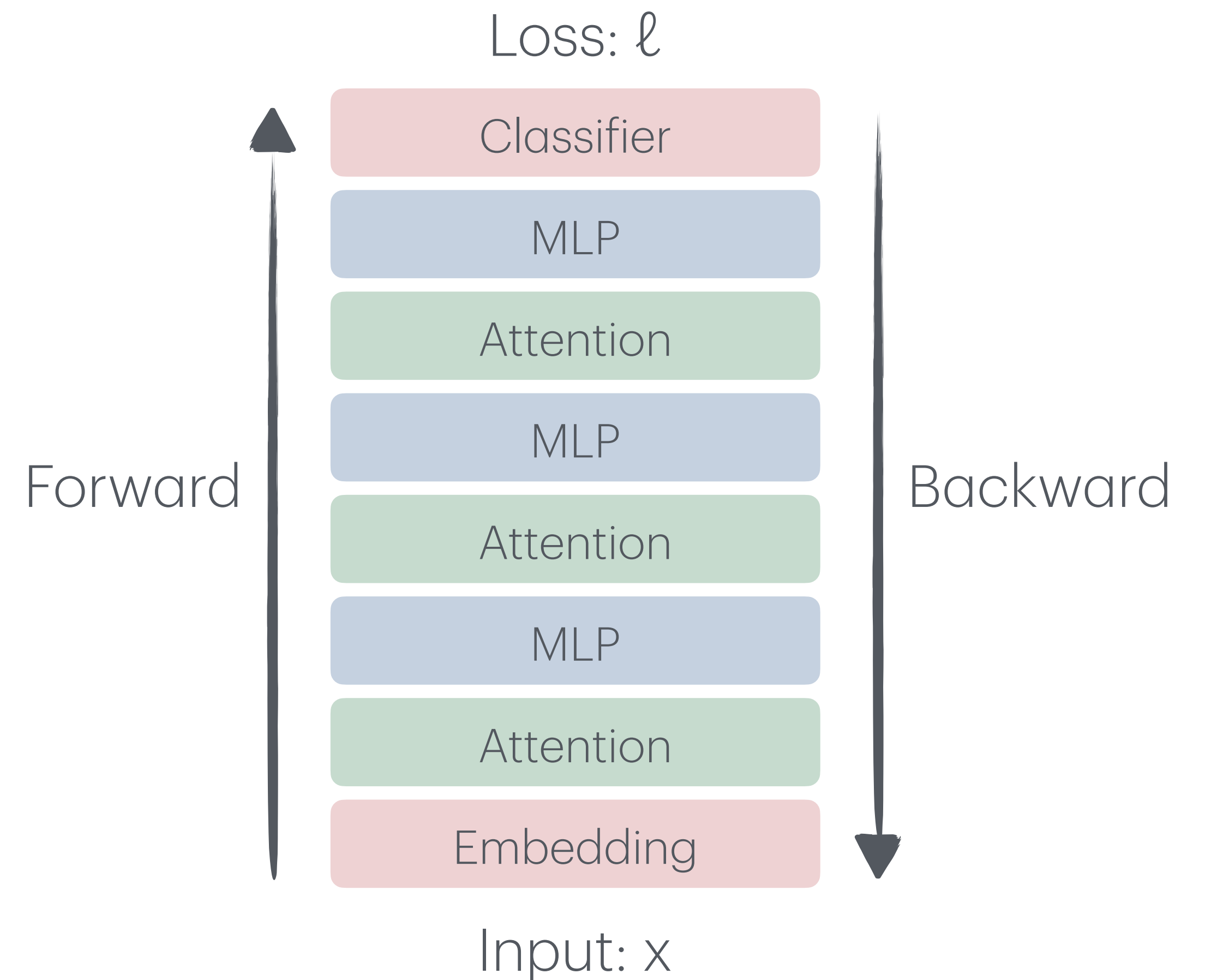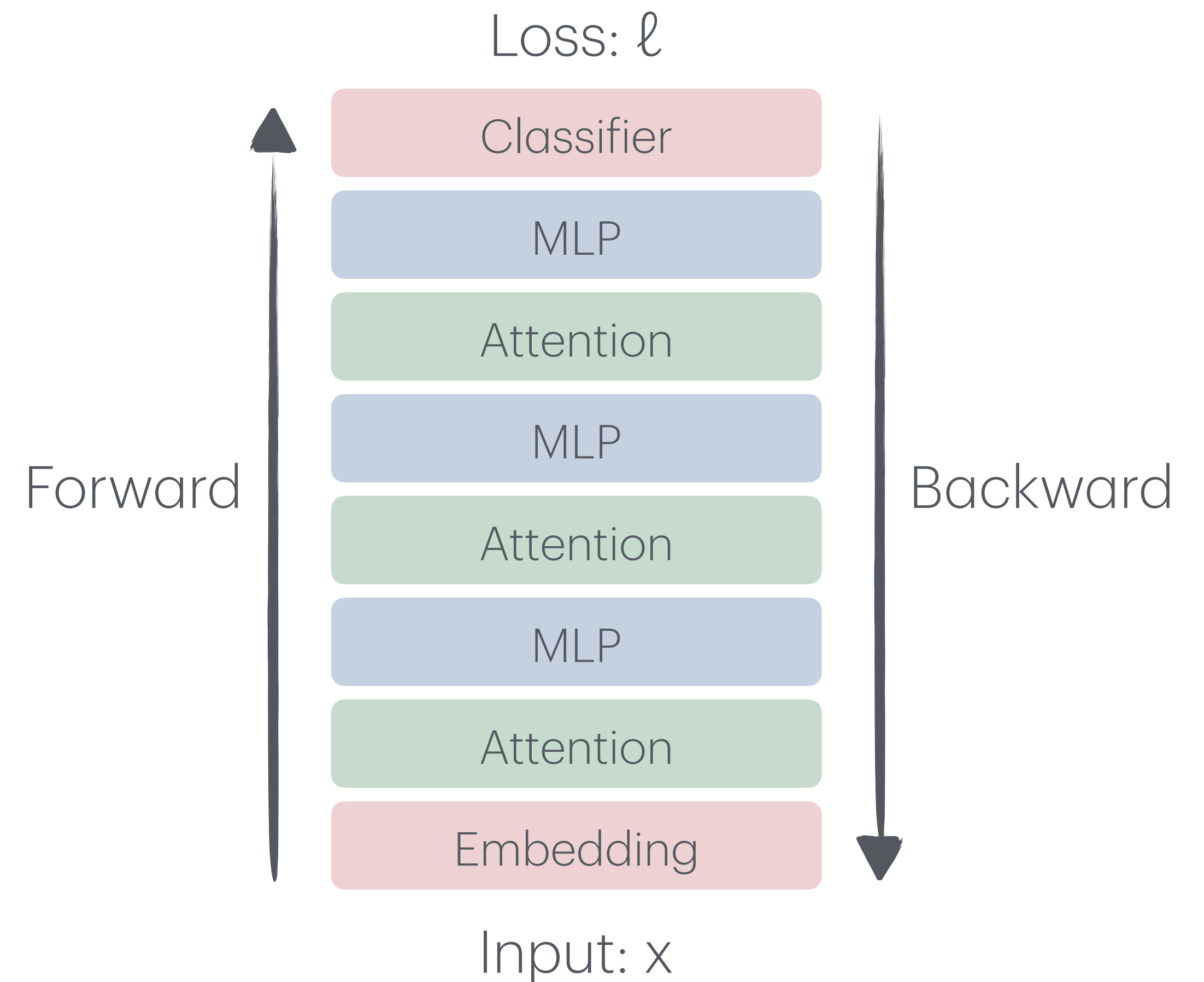| | |
|---|---|
| 4N bytes | Gradient (bf16) |
| 4N bytes | Weight (fp32) |
| 4N bytes | Momentum 1 (fp32) |
| 4N bytes | Momentum 2 (fp32) |

# Memory Use

- What takes up GPU memory during training?

  - Model weights

  - Gradients

  - Momentum

  - Activations (more later in class)

Loss: ℓ

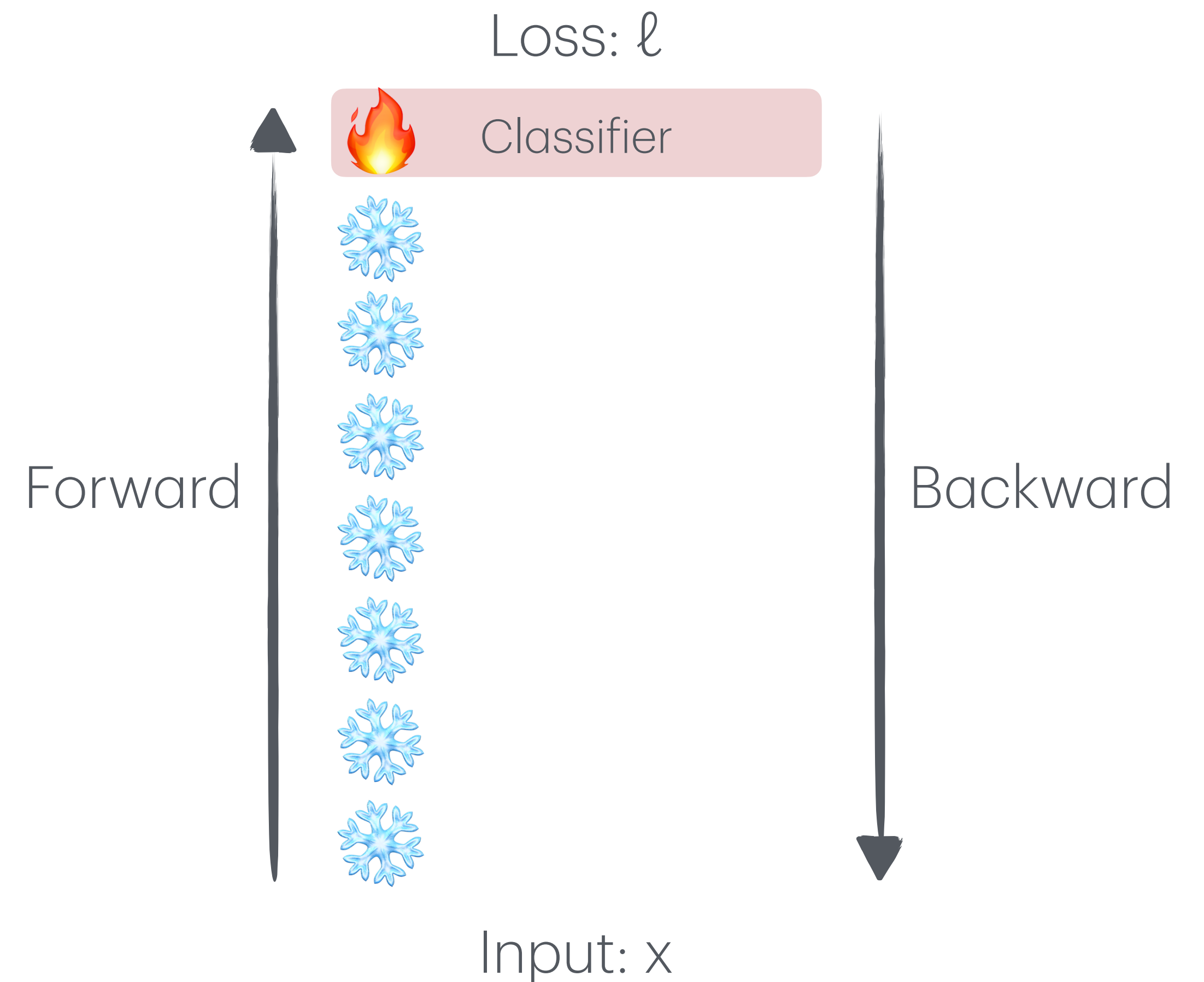| |
|---|
| Classifier |
| MLP |
| Attention |
| MLP |
| Attention |
| MLP |
| Attention |
| Embedding |

Forward

Backward

Input: x

# Reducing Memory Use

- Idea: Train fewer parameters

  - Keep most parameters frozen

    - No gradient, no momentum

  - Train a small subset

    - With gradient and momentum

Loss: $\ell$

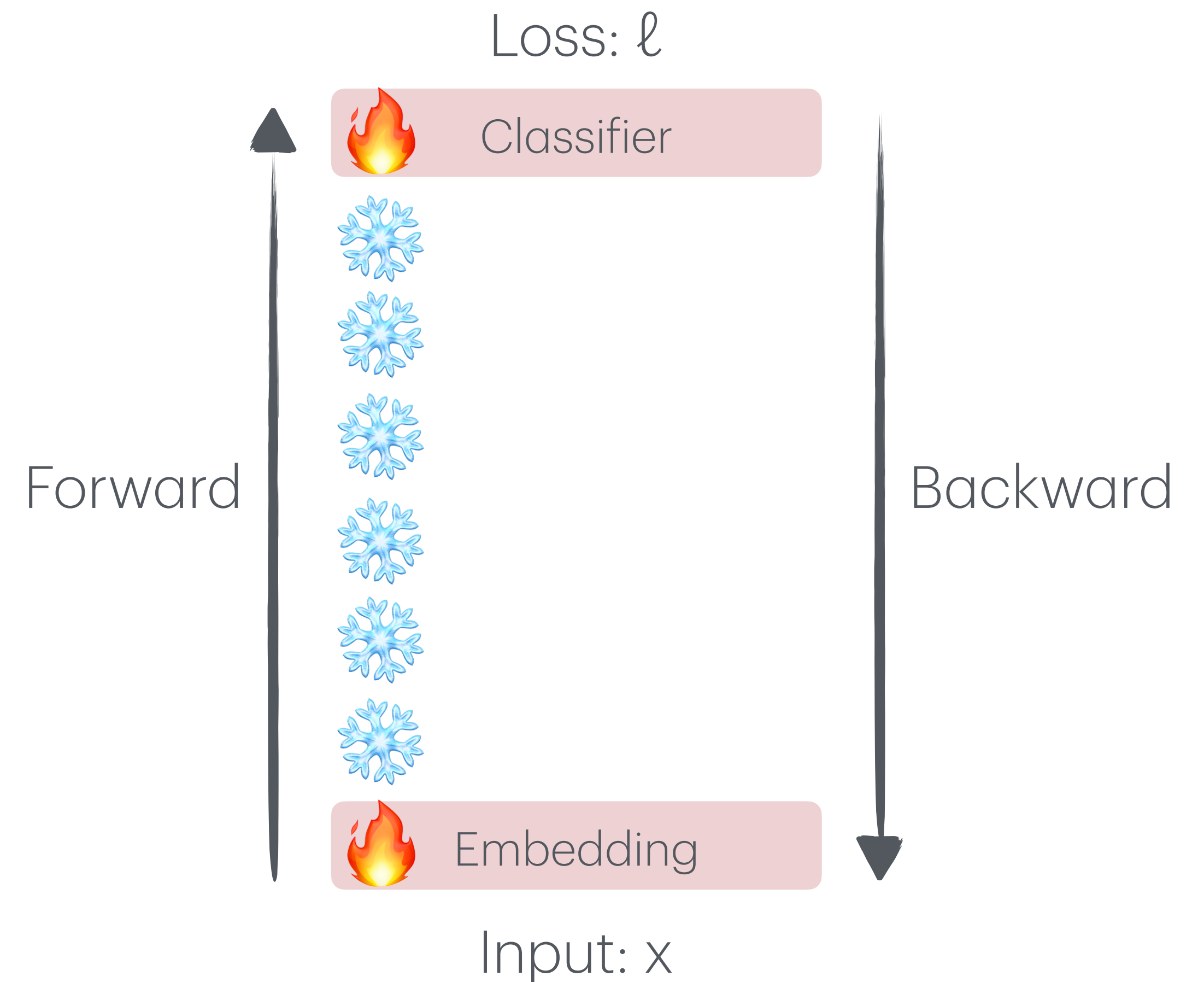| Classifier |
| MLP |
| Attention |
| MLP |
| Attention |
| MLP |
| Attention |
| Embedding |

Forward

Backward

Input: x

# Fine-tuning classifier

- Freeze backbone

- Train classifier

- Most memory efficient

  - No backprop

  - Very few learnable parameters

- Not very expressive

Loss: $\ell$

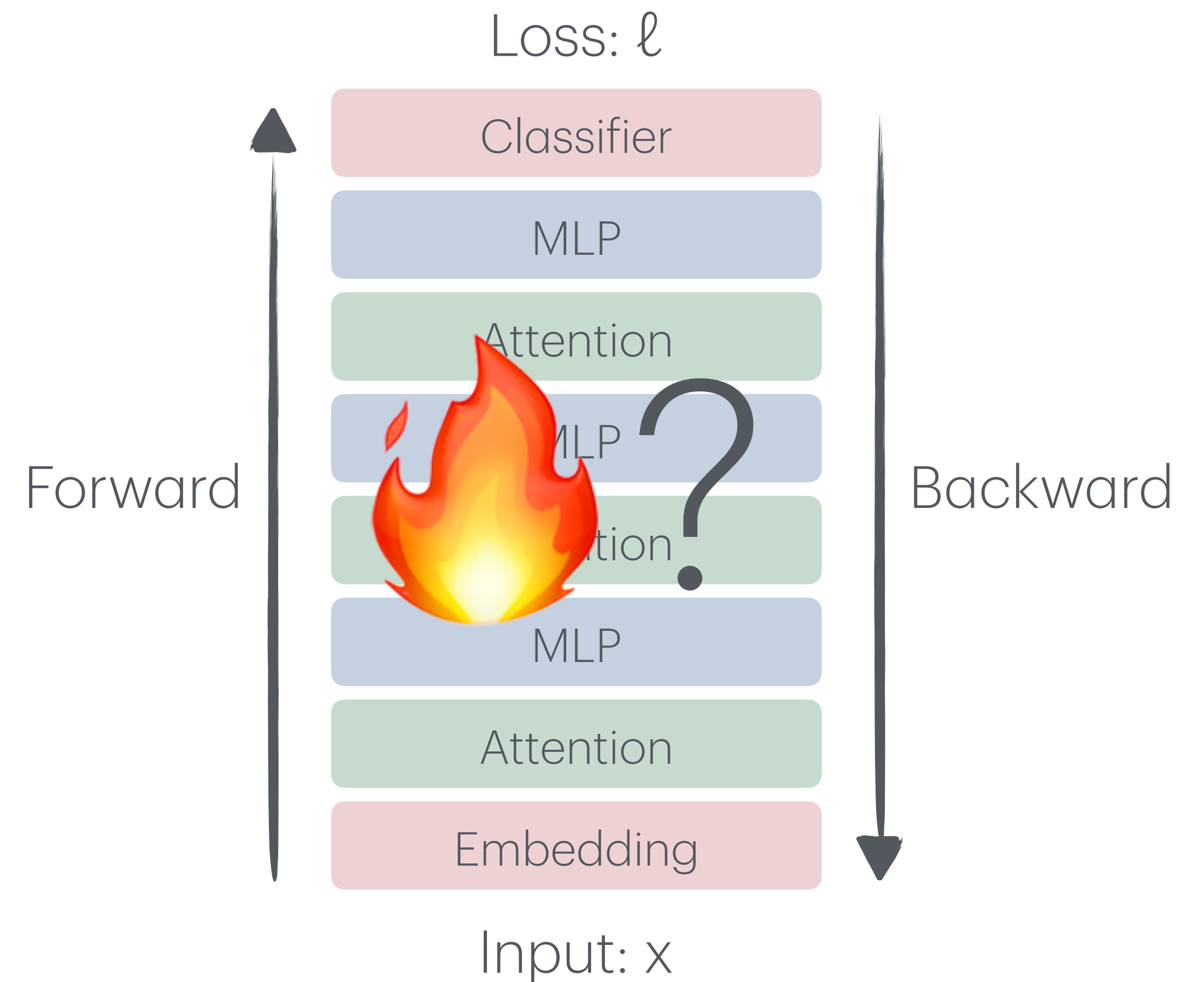🔥 Classifier

Forward

Backward

Input: x

# Input adapters

- Freeze backbone

- Train input embedding (maybe classifier)

- Fairly memory efficient

  - Very few learnable parameters

- Popular with LLMs

  - Soft-prompting

  - Adapters for new inputs

Loss: $\ell$

🔥 Classifier

❄️
❄️
❄️
❄️
❄️
❄️

Forward

Backward

🔥 Embedding

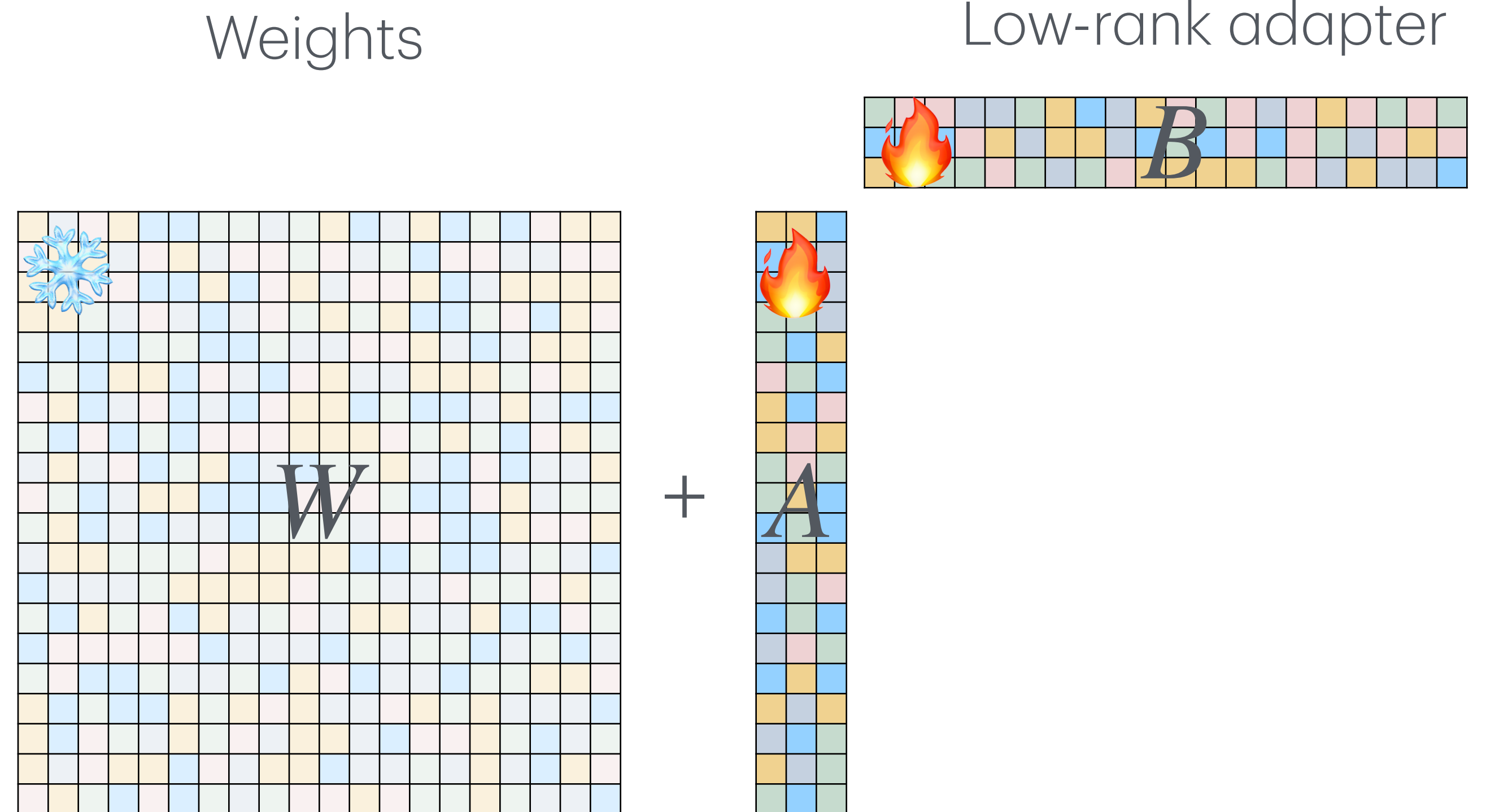Input: x

# Intermediate layers

- Fine-tuning input and output does not change computation inside network significantly

  - Cannot learn new "internal computation"

- Can we learn a subset of intermediate layer parameters?

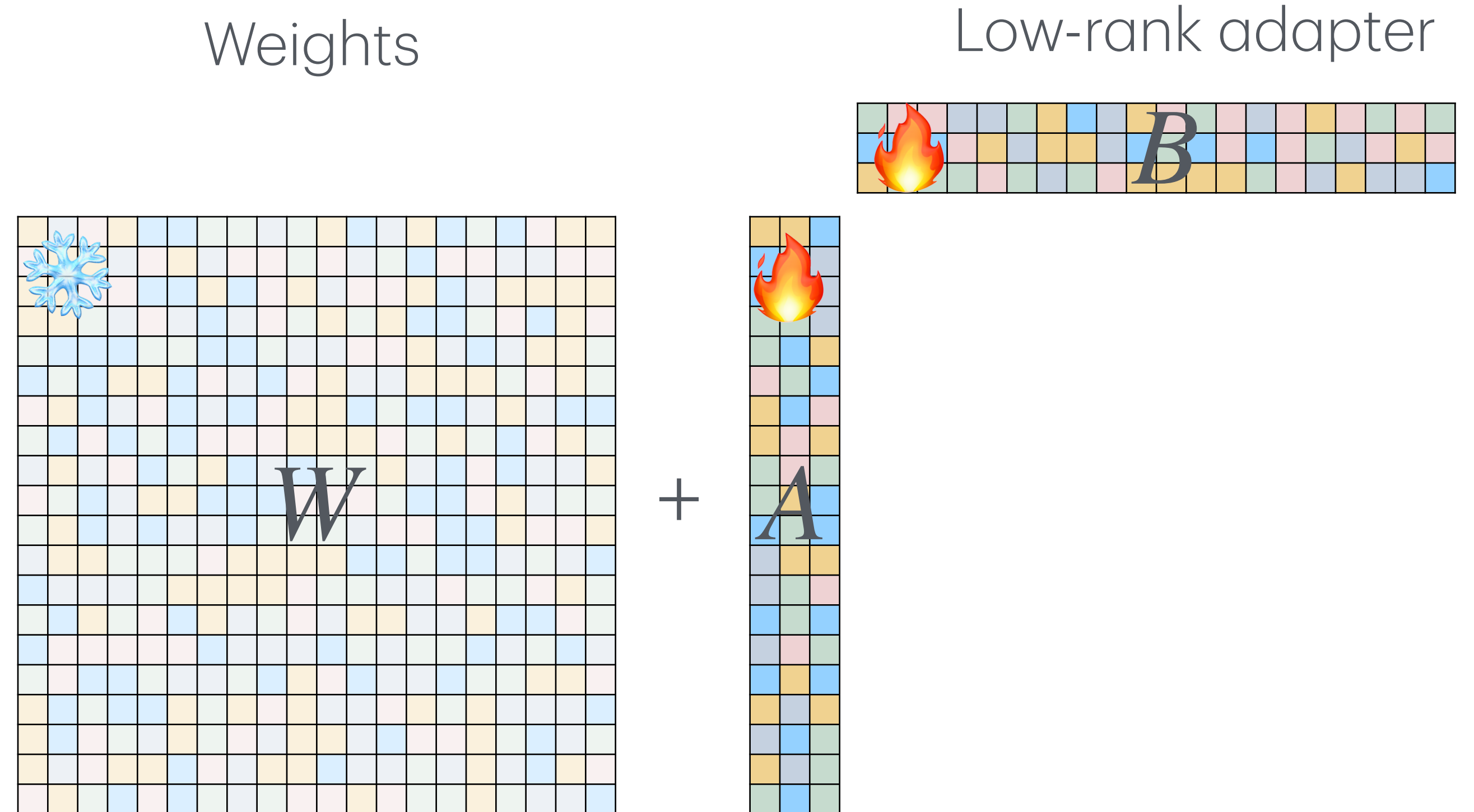# Low Rank Adapters

## LoRA

Weights

Low-rank adapter

- Keep weights $W \in \mathbb{R}^{N \times M}$ frozen

- Learn adapter $AB$

  - $A \in \mathbb{R}^{N \times R}, B \in \mathbb{R}^{R \times M}$

  - Rank $R \ll min(M, N)$

- Total parameters: $R(N + M)$



[1] Edward J. Hu, et al. Lora: Low-rank adaptation of large language models. 2021

# Low Rank Adapters

Weights
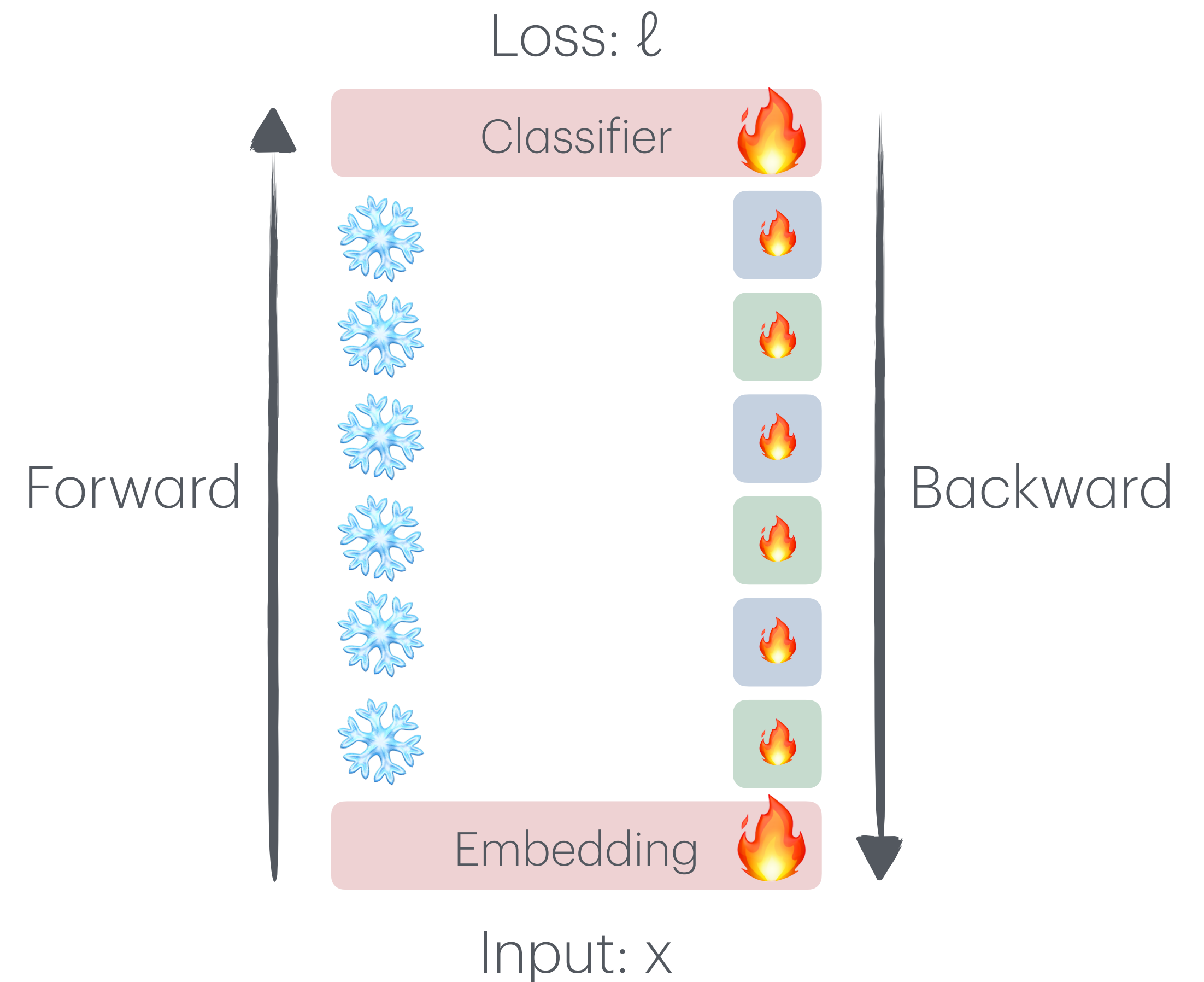
Low-rank adapter

- How do we train $W$?

  - We use a pre-trained model

- How do we initialize $A$ and $B$?

  - $A$ small random (normal)

  - $B$ zero

$$W + A$$



[1] Edward J. Hu, et al. Lora: Low-rank adaptation of large language models. 2021

# LoRA models

- Most weights frozen

- Train adapter for

  - All linear layers

  - Just MLPs

  - Just Attention

- Optionally train full input and output embedding

Loss: ℓ

Classifier 🔥

Forward

Backward

Embedding 🔥

Input: x

# LoRA in practice

- Download weights of a pre-trained model

- Define LoRALinear layer

$$Wx + b \rightarrow Wx + b + \alpha \frac{AB}{R} x$$

- Rewrite model with LoRALinear layers

  - If careful, `load_state_dict` just works

```python
class LoRALinear(nn.Linear):
    def __init__(self, in_features: int, out_features: int, rank: int,
                 alpha: float, bias: bool = True, device=None,
                 dtype=None):
        super().__init__(in_features, out_features, bias, device, dtype)

        self.lora_a = nn.Linear(in_features, rank, bias=False,
                                device=device, dtype=dtype)
        self.lora_b = nn.Linear(rank, out_features, bias=False,
                                device=device, dtype=dtype)
        self.alpha_div_rank = alpha / rank

        nn.init.kaiming_uniform_(self.lora_a.weight)
        nn.init.zeros_(self.lora_b.weight)

    def forward(self, x: Tensor) -> Tensor:
        return super().forward(x) + \
            self.alpha_div_rank / self.lora_b(self.lora_a(x))
```
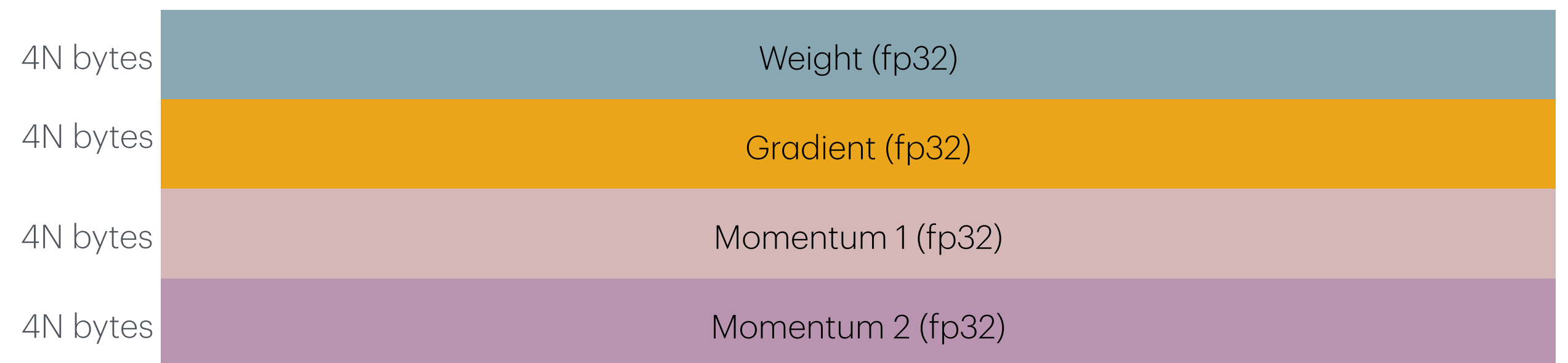
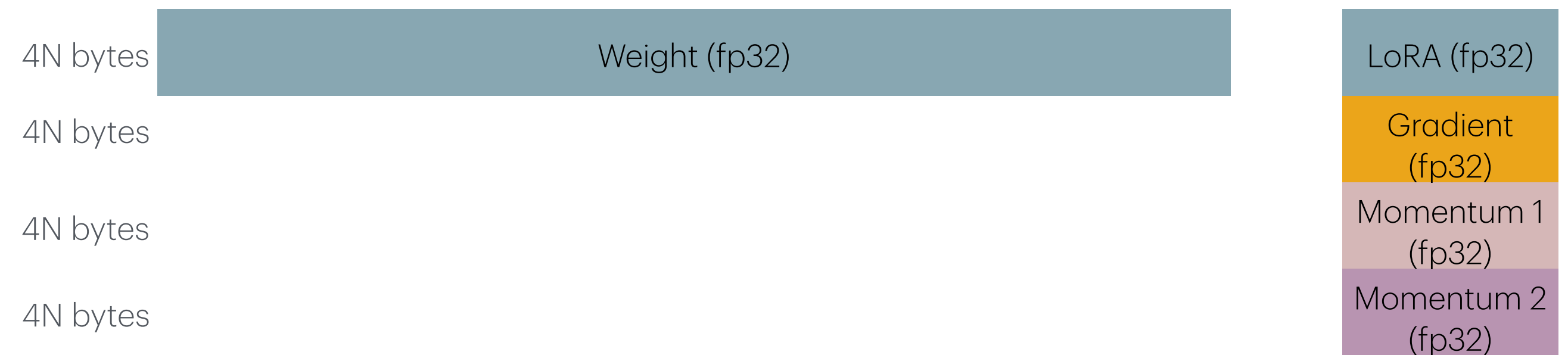# Training large models

## Memory requirements

- Without optimization:

  - Model parameters: N

  - Weights: N floats

  - Gradients: N floats

  - Momentum: N floats

  - 2nd momentum (ADAM): N floats

- 16N bytes without counting activations

| 4N bytes | Weight (fp32) |
|---|---|
| 4N bytes | Gradient (fp32) |
| 4N bytes | Momentum 1 (fp32) |
| 4N bytes | Momentum 2 (fp32) |

# Training LoRA models

## Memory requirements

- LoRA

  - Model parameters: N, LoRA param M

  - Weights: N+M floats

  - Gradients: M floats

  - Momentum: M floats

  - 2nd momentum (ADAM): M floats

- 4N+16M bytes without activations

- M often ~1-5% of N

| | |
|---|---|
| 4N bytes | Weight (fp32) |
| 4N bytes | |
| 4N bytes | |
| 4N bytes | |

LoRA (fp32)

Gradient (fp32)

Momentum 1 (fp32)

Momentum 2 (fp32)

# References

- [1] Edward J. Hu, et al. Lora: Low-rank adaptation of large language models. 2021 ([link](link))