

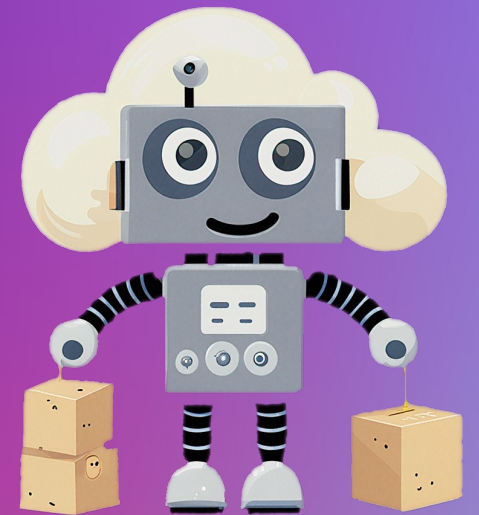


SES109

Parameter Efficient Finetuning (PEFT) with LoRA

Mariano Kamp

Principal Solutions Architect
Amazon Web Services



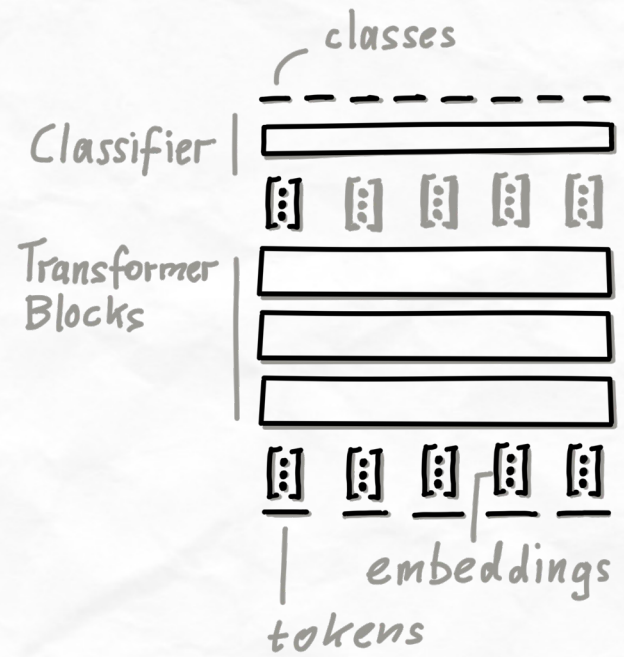
Model	BoolQ	CB	CoPA	MultiRC	Record	RTE	WiC	WSC
Few-shot	89.1	89.3	95	86.3/-	92.9/-	81.2	64.6	89.5
Finetuned	<u>92.2</u>	<u>100/100</u>	<u>100</u>	<u>90.1/69.2</u>	<u>94.0/94.6</u>	<u>95.7</u>	<u>78.8</u>	<u>100</u>

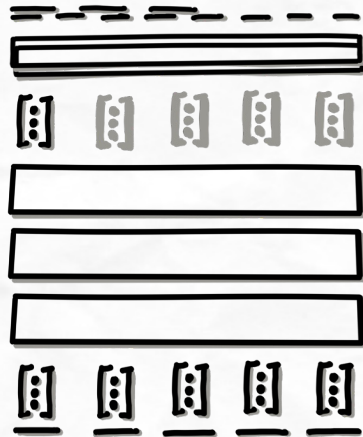
Table 8: Results on SuperGLUE dev set comparing PaLM-540B few-shot and finetuned.

PaLM: Scaling Language Modeling with Pathways
Chowdhery, Narang, Devlin et al.

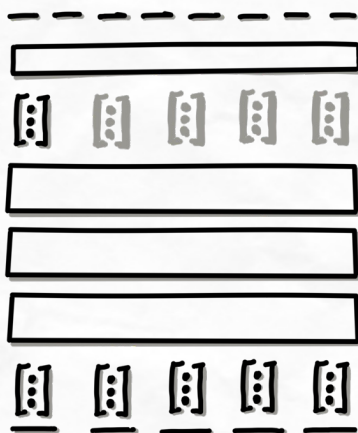
PRE-TRAINING

FINE-TUNING





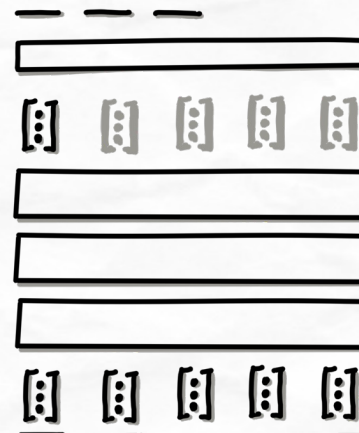
Next/masked
token prediction



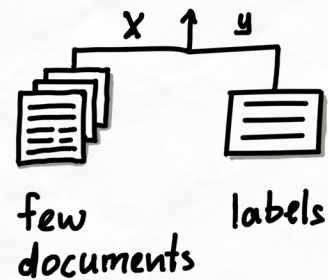
different
task

same
architecture

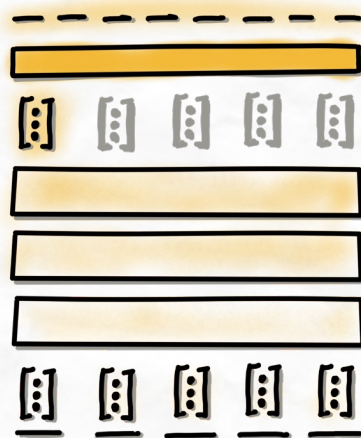
different
dataset



Sentiment
prediction



Next/masked
token prediction

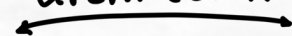


large
corpora

different
task



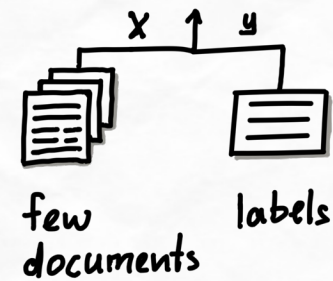
same
architecture



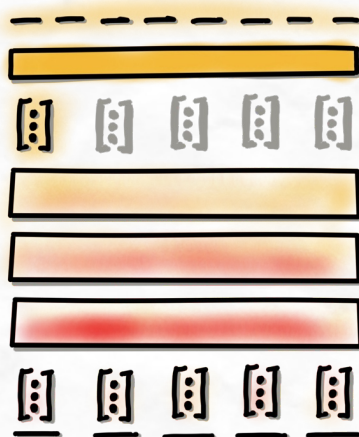
different
dataset



Sentiment
prediction

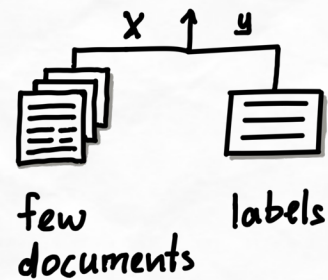


Next/masked
token prediction



Syntax
Grammar
Substitutions
...

Sentiment
prediction

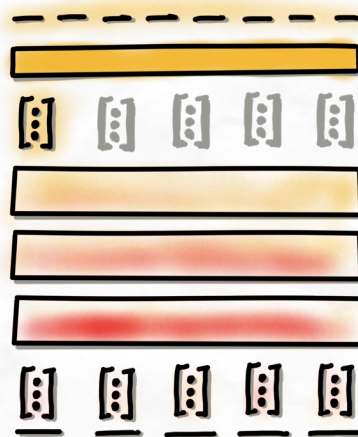


PRE-TRAINING „UPSTREAM-MODEL“

FINE-TUNING „DOWNSTREAM-MODEL“

Next/masked
token prediction

$h \sim N \rightarrow$

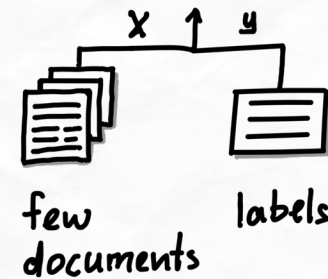


init $h \sim N \rightarrow$

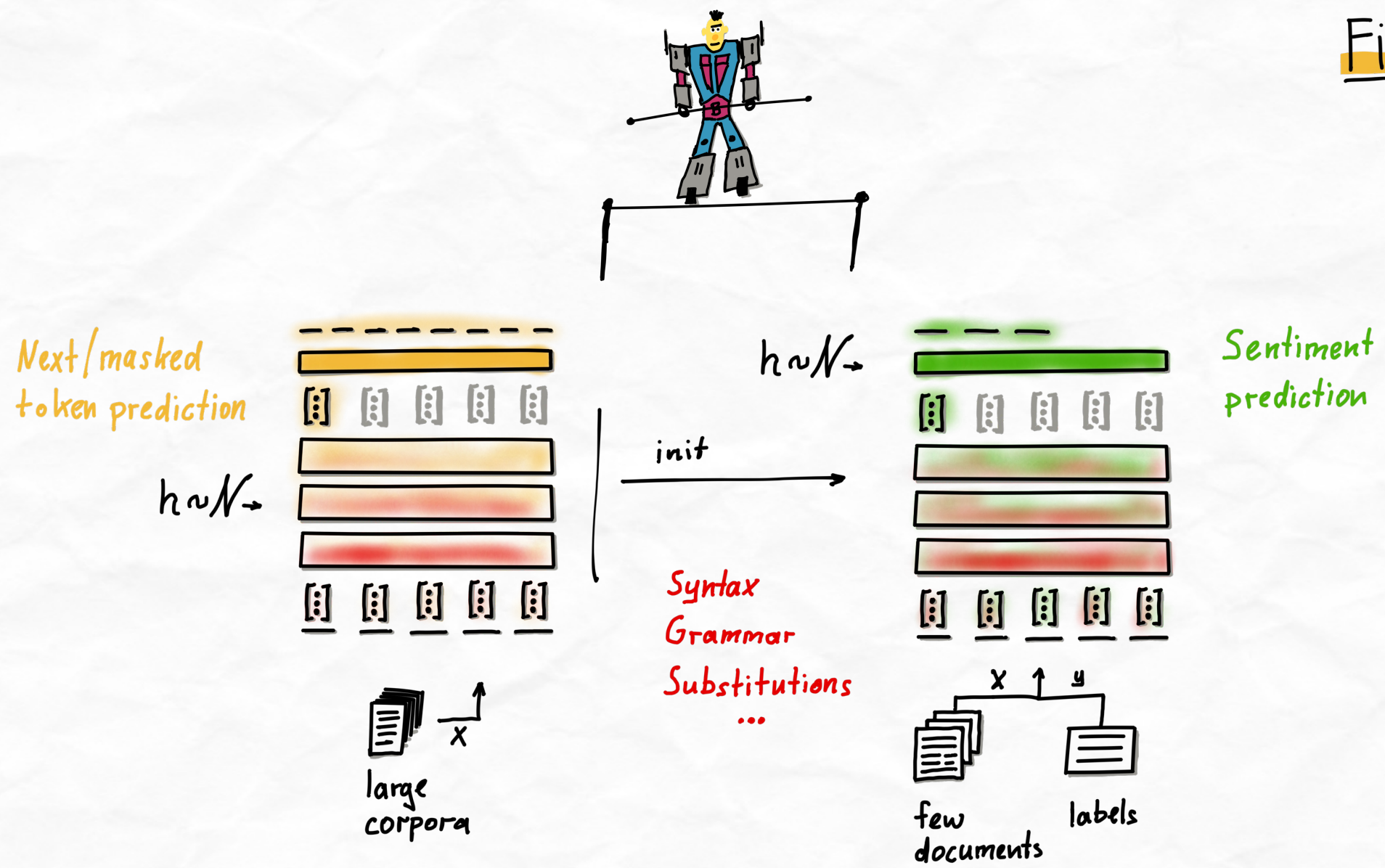
Syntax
Grammar
Substitutions
...

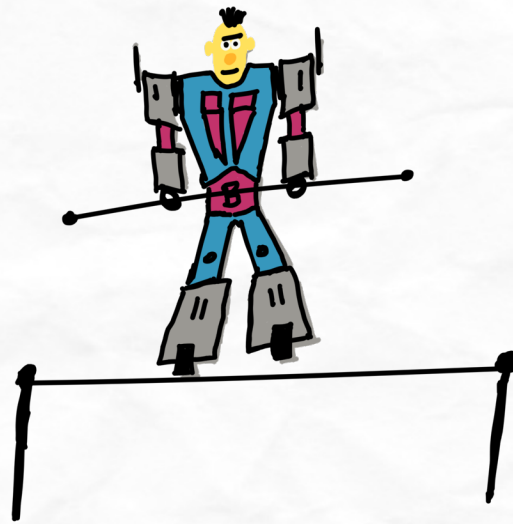
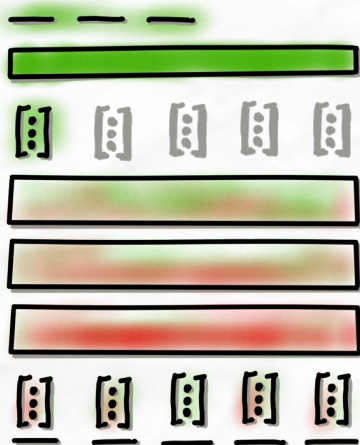


Sentiment
prediction



FINE-TUNING

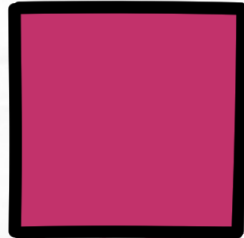




TIGHTROPE WALK

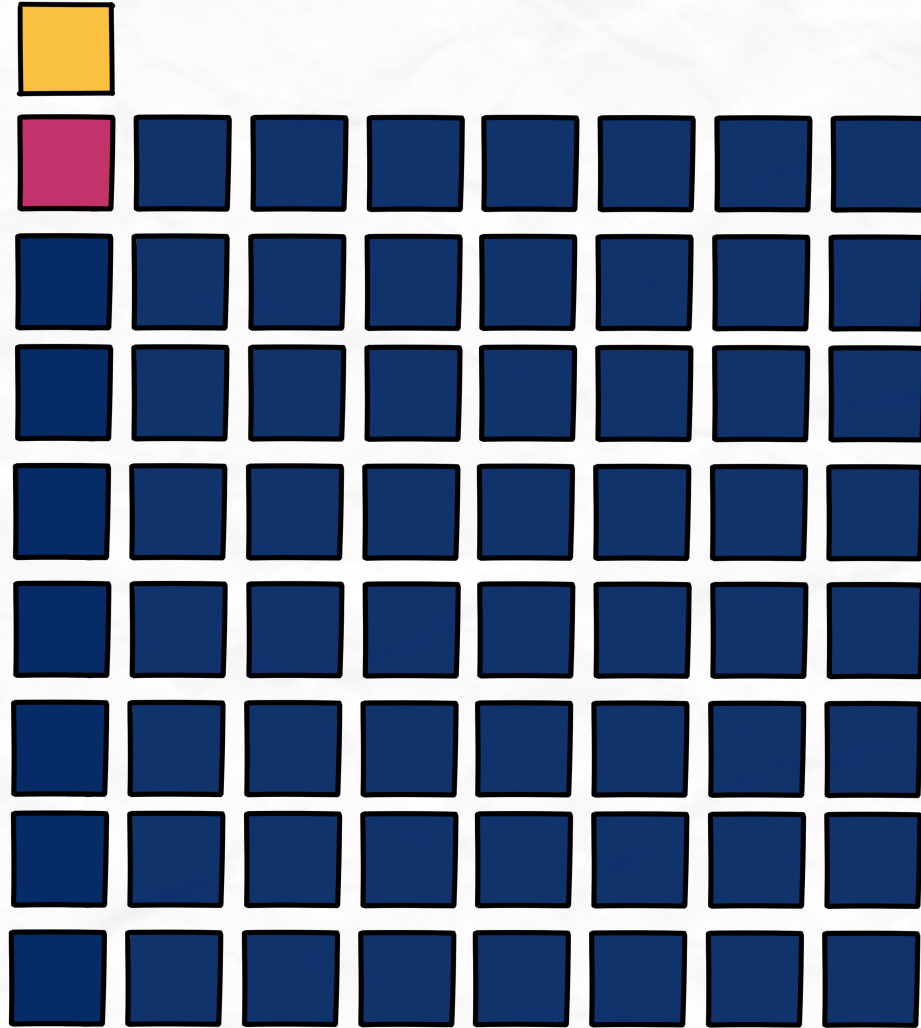
- Fewer epochs
- Cautious LR
- Per-Layer-LR
- Frozen layers

SIZE



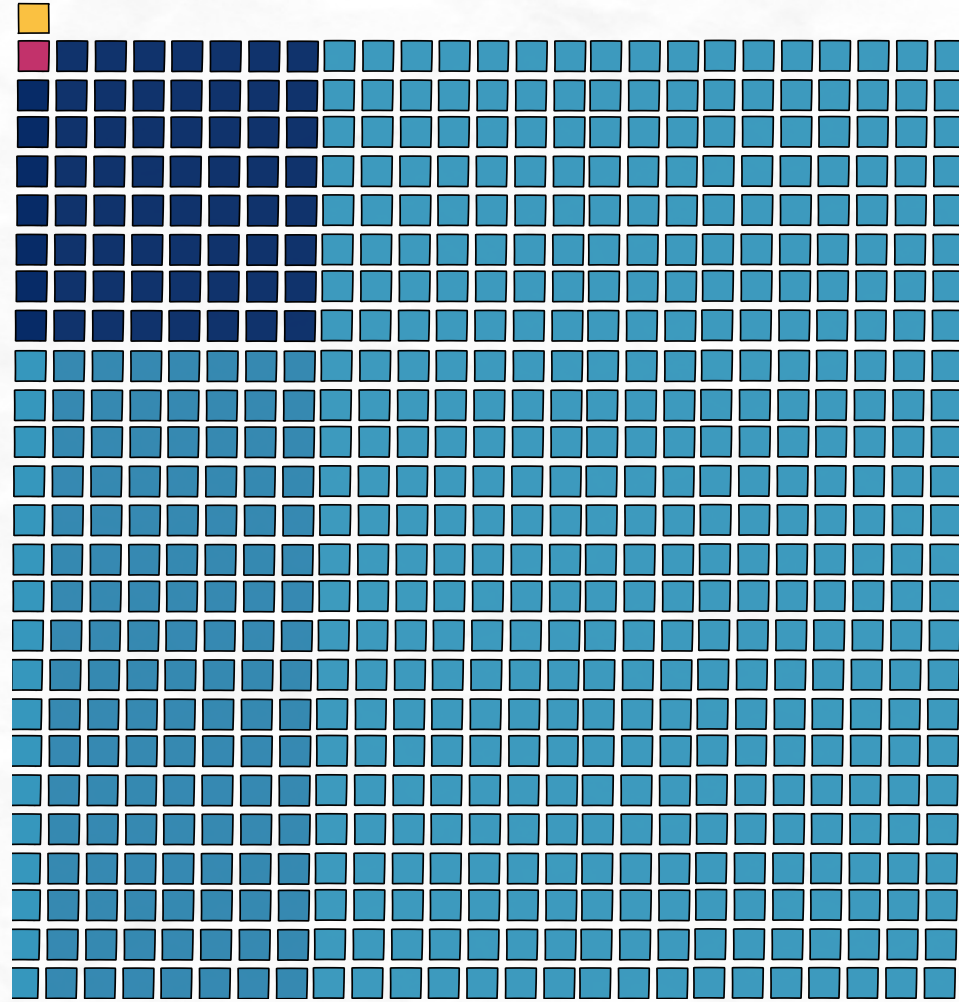
 LABELED DATA

 BERT BASE 110M



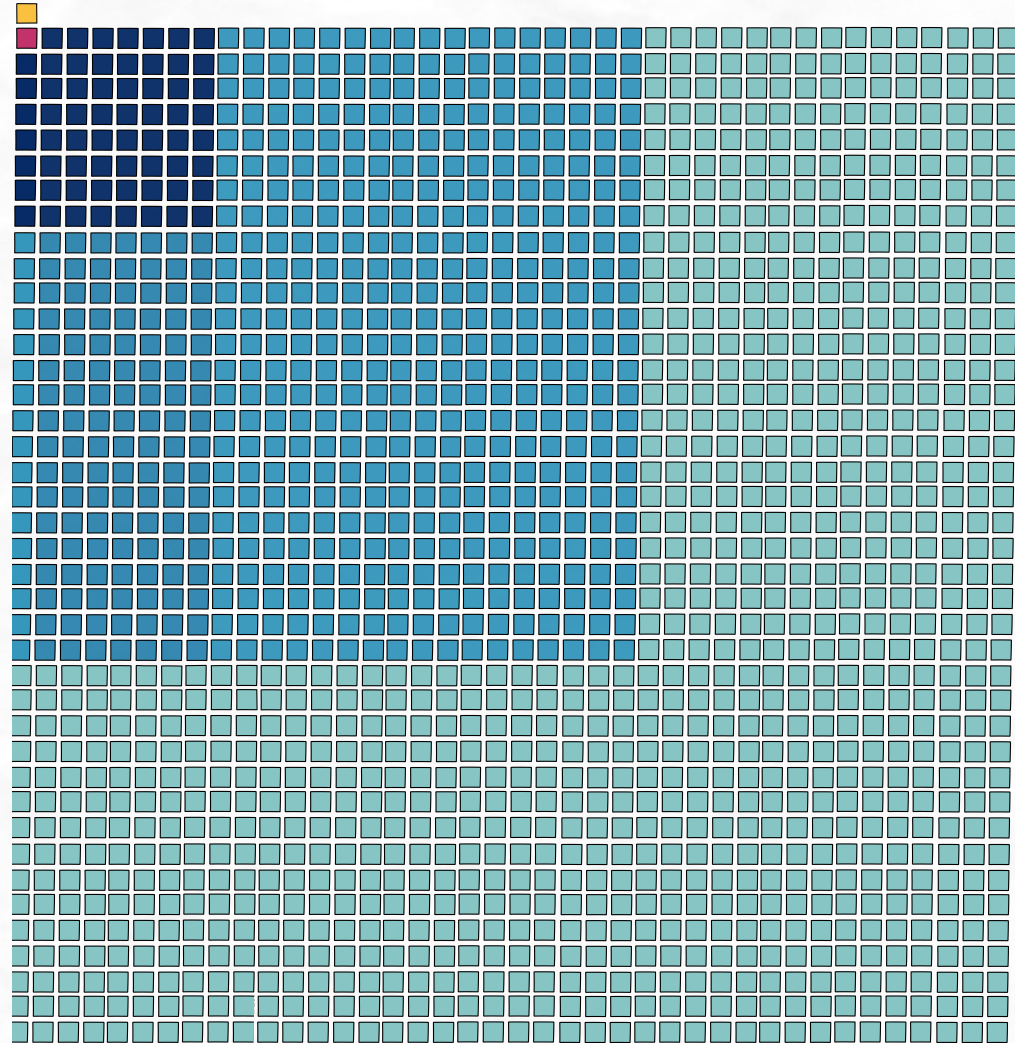
 LABELED DATA

  7B



 LABELED DATA

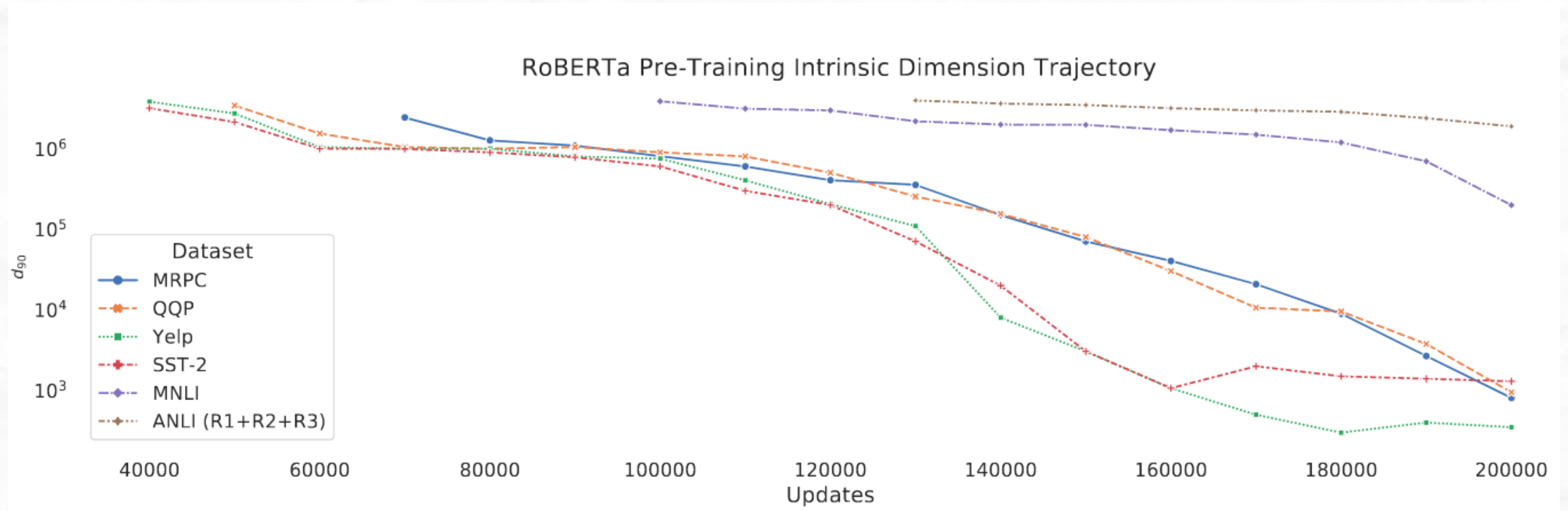
   70B



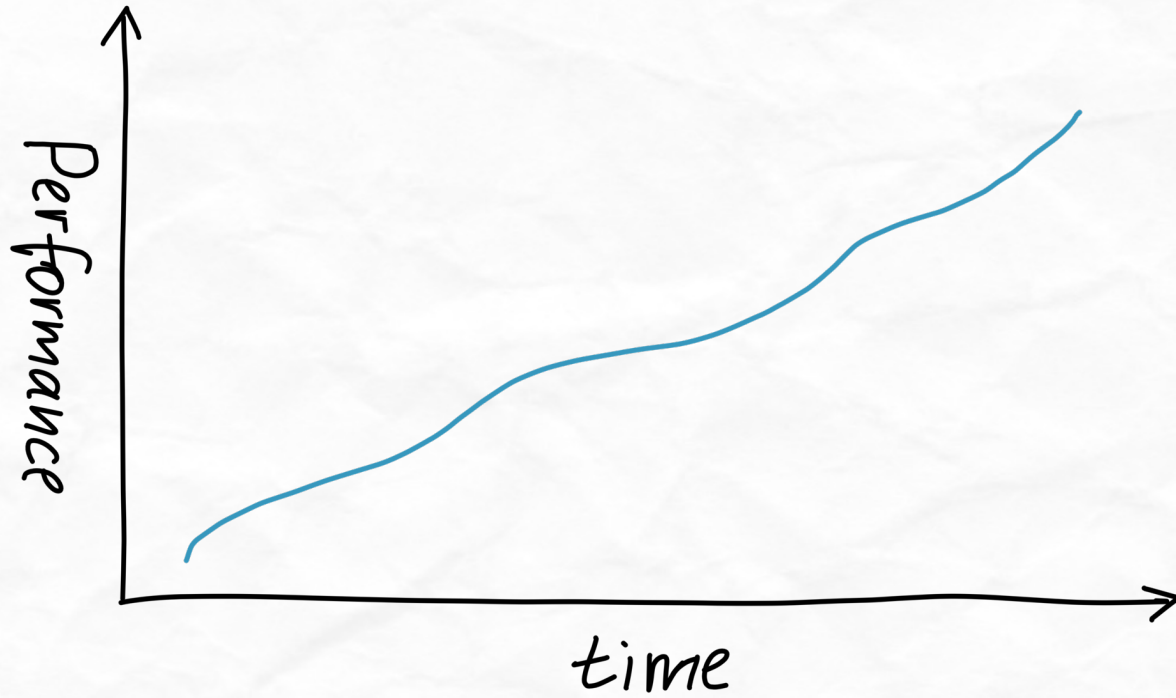
 LABELED DATA

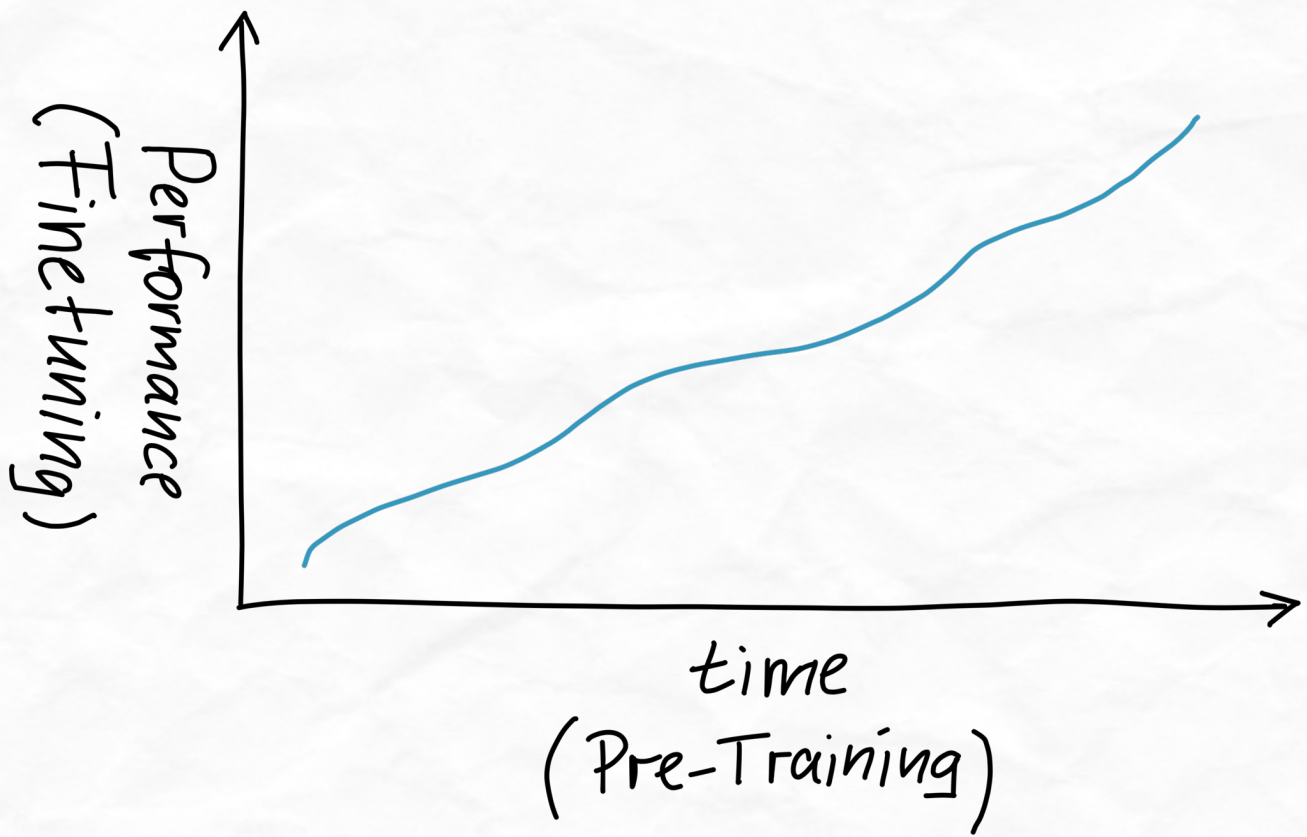
    180 B

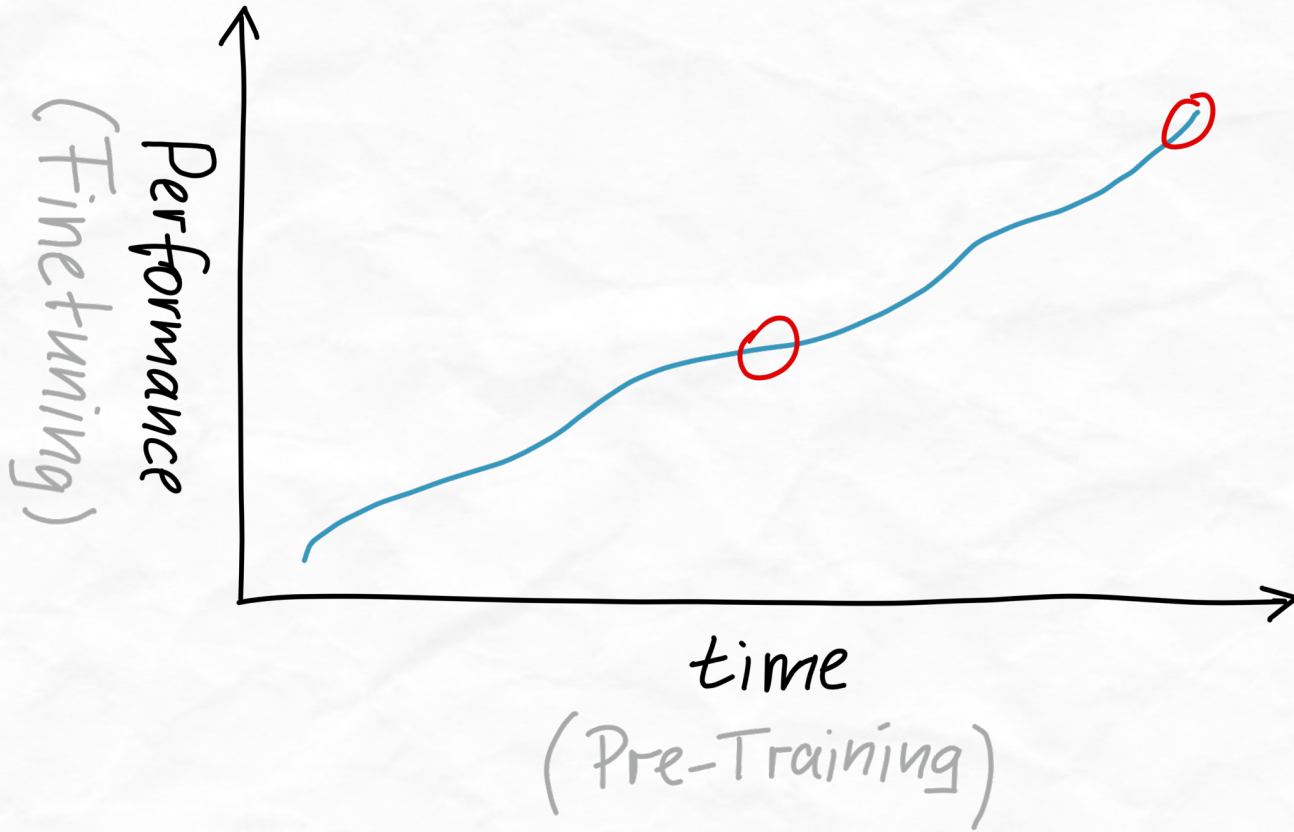
Intrinsic Dimensionality

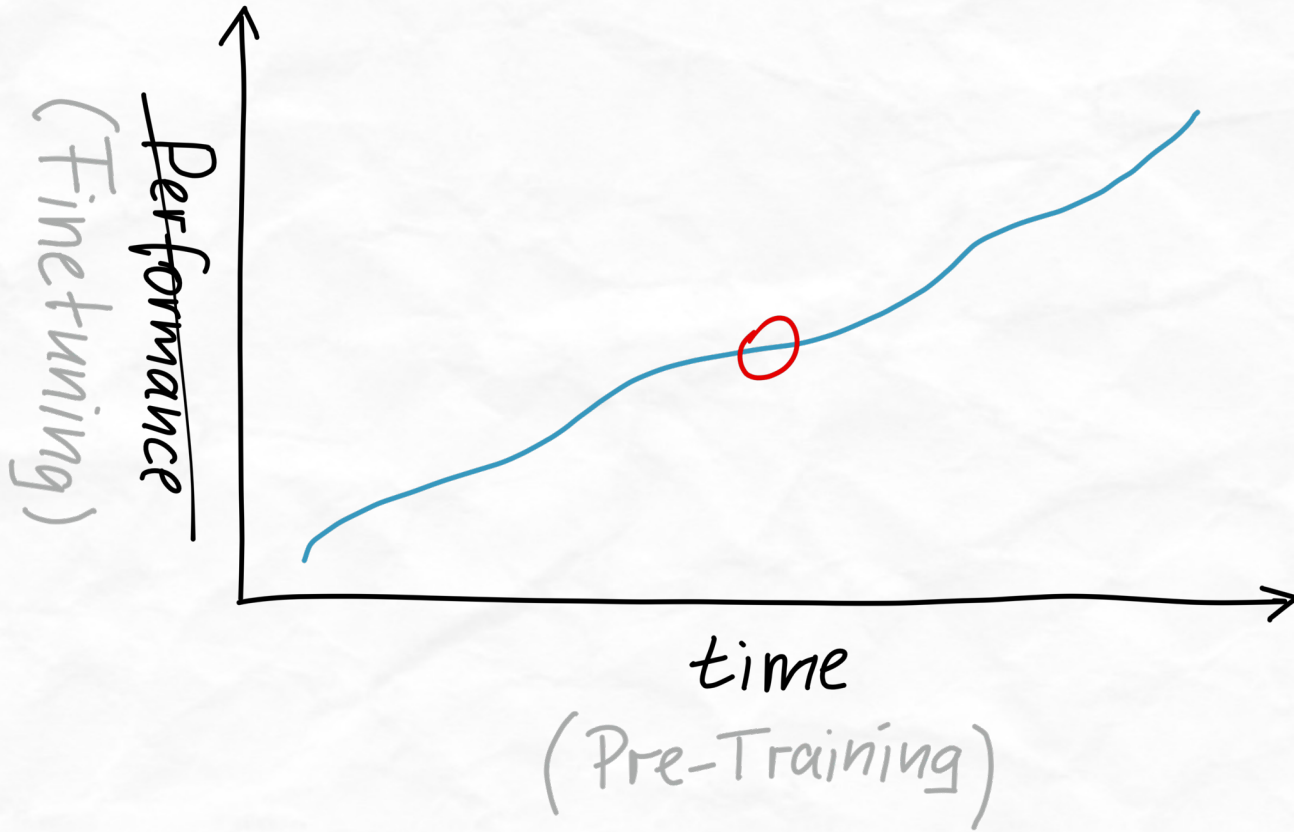


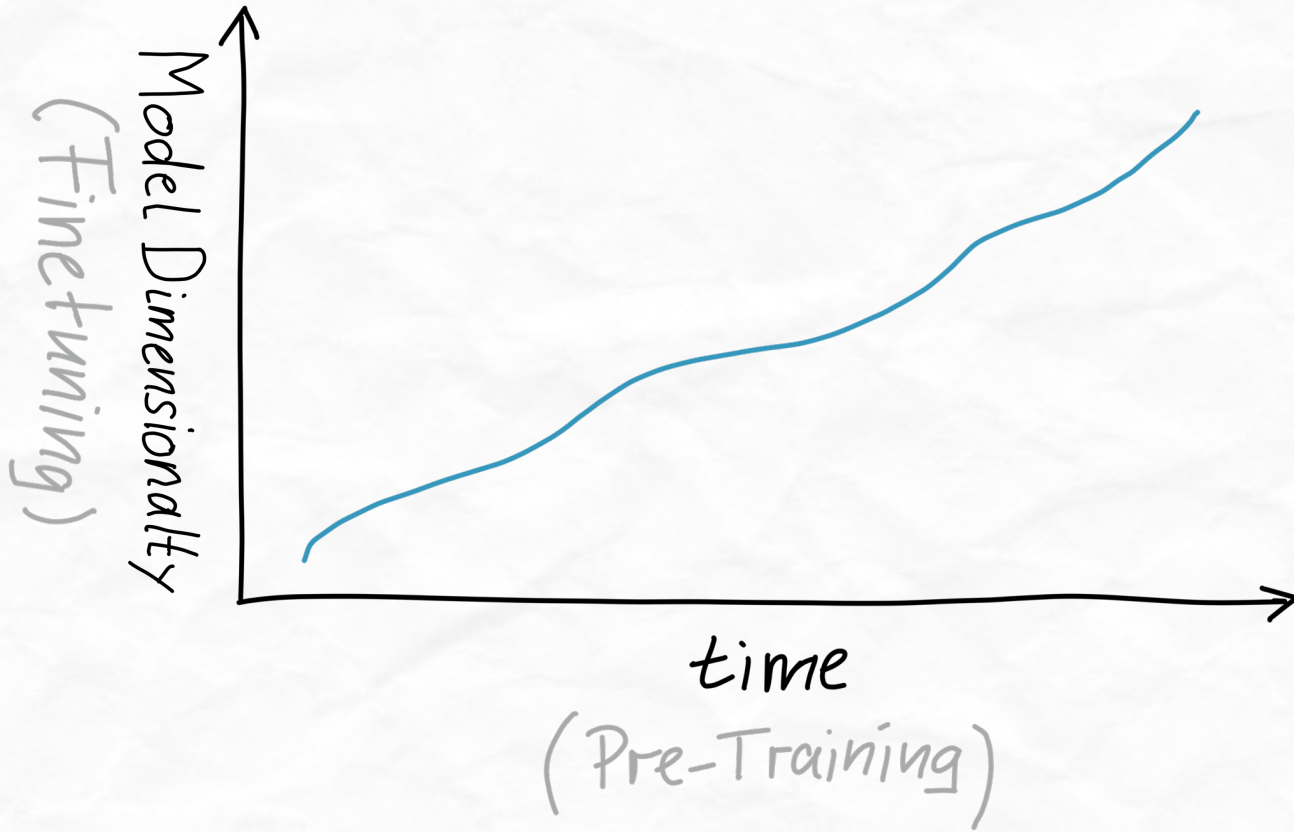
INTRINSIC DIMENSIONALITY EXPLAINS THE EFFECTIVENESS OF LANGUAGE MODEL FINE-TUNING
Aghajanyan, Zettlemoyer, Gupta, 2020

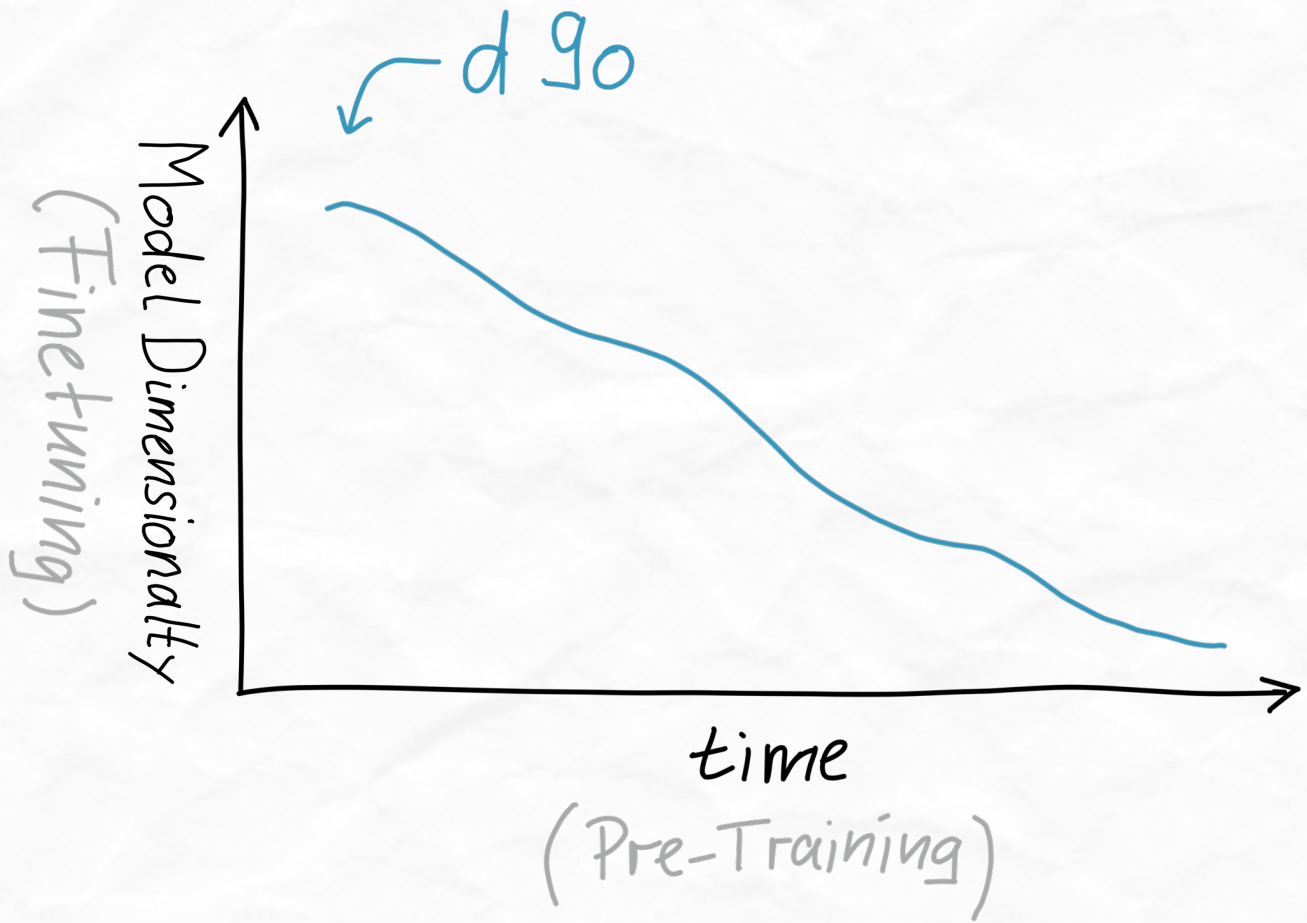


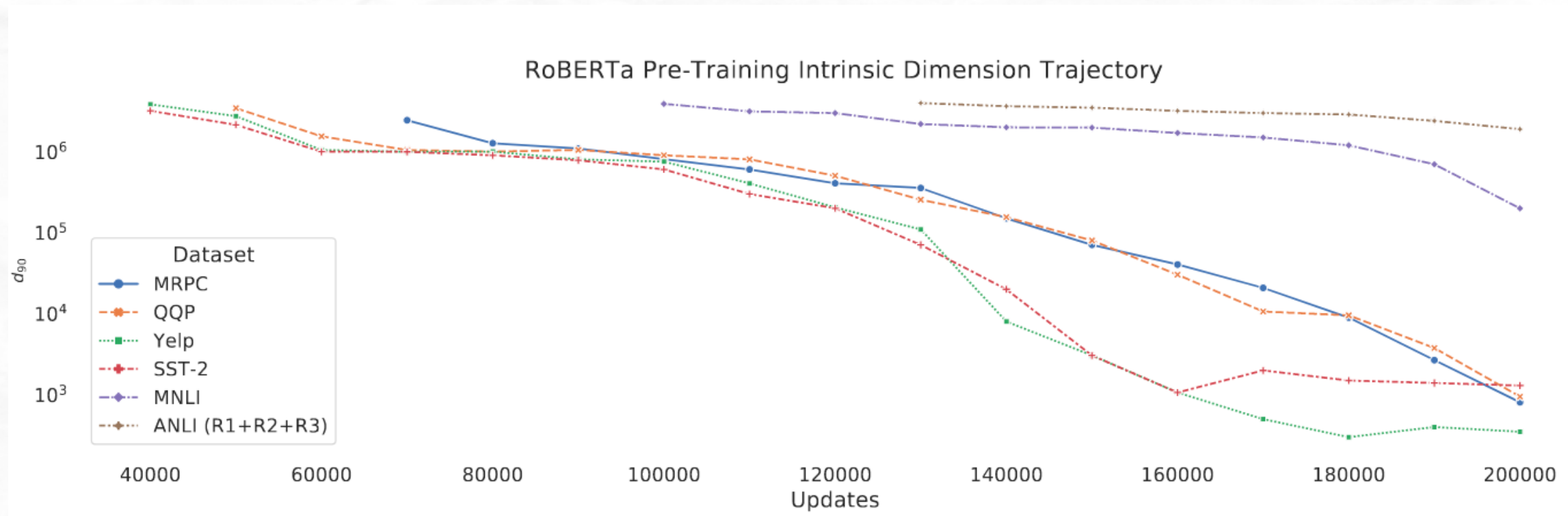




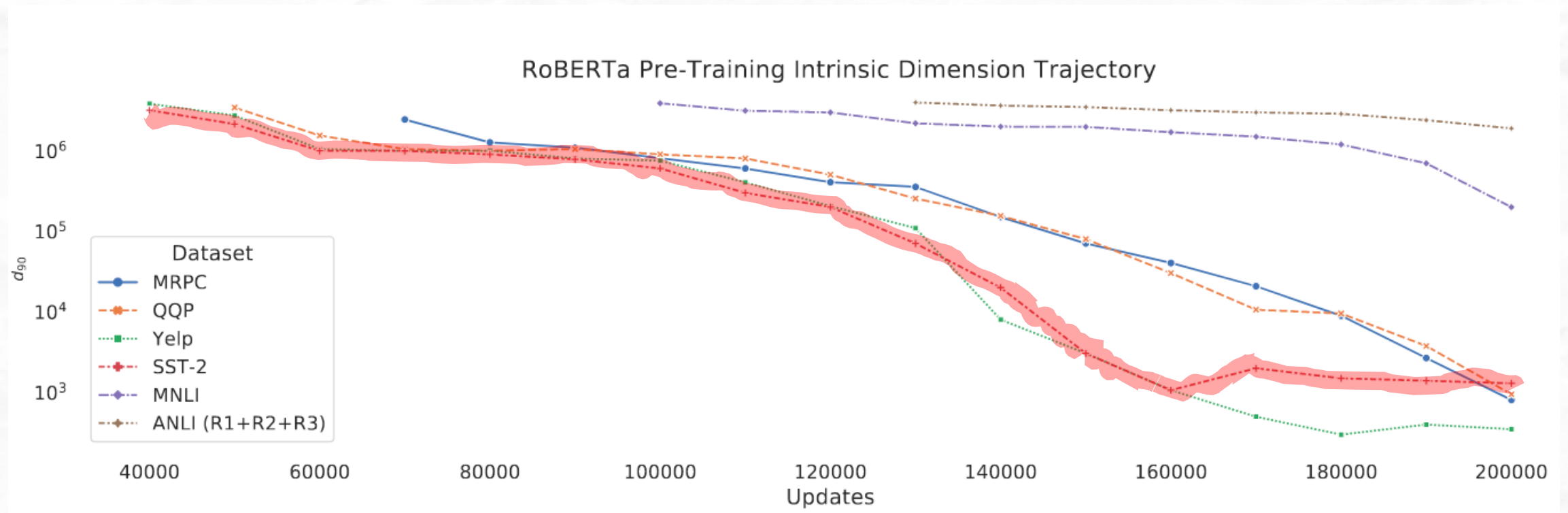




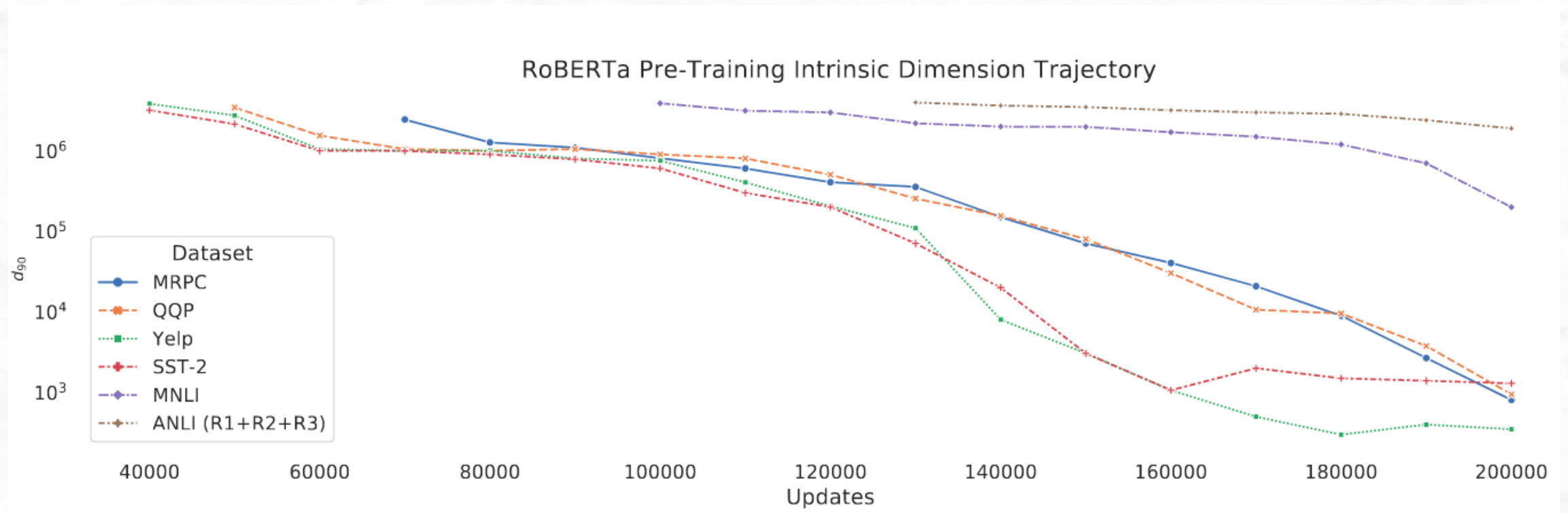




INTRINSIC DIMENSIONALITY EXPLAINS THE EFFECTIVENESS OF LANGUAGE MODEL FINE-TUNING
 Aghajanyan, Zettlemoyer, Gupta, 2020



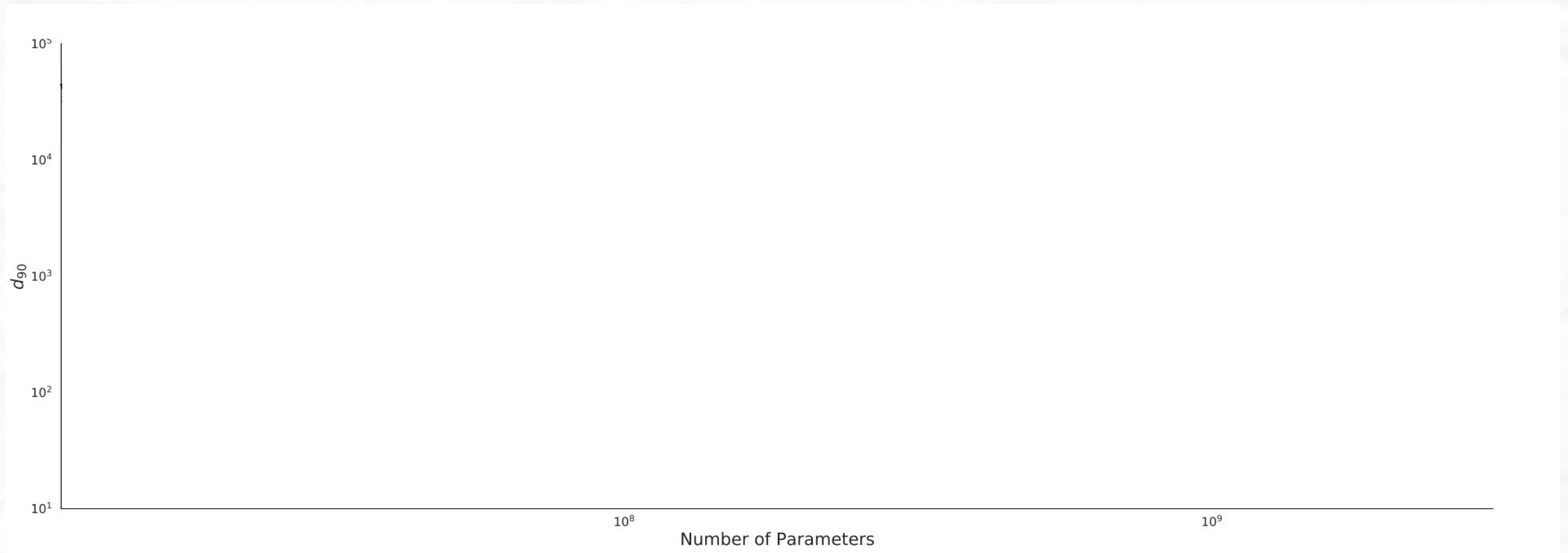
INTRINSIC DIMENSIONALITY EXPLAINS THE EFFECTIVENESS OF LANGUAGE MODEL FINE-TUNING
 Aghajanyan, Zettlemoyer, Gupta, 2020



INTRINSIC DIMENSIONALITY EXPLAINS THE EFFECTIVENESS OF LANGUAGE MODEL FINE-TUNING
 Aghajanyan, Zettlemoyer, Gupta, 2020



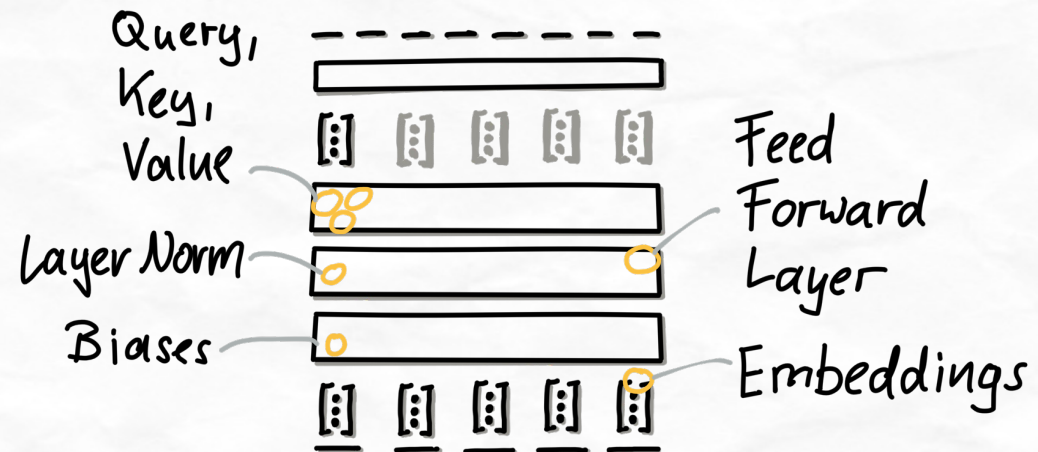
Intrinsic Dimensionality (by Model Size)



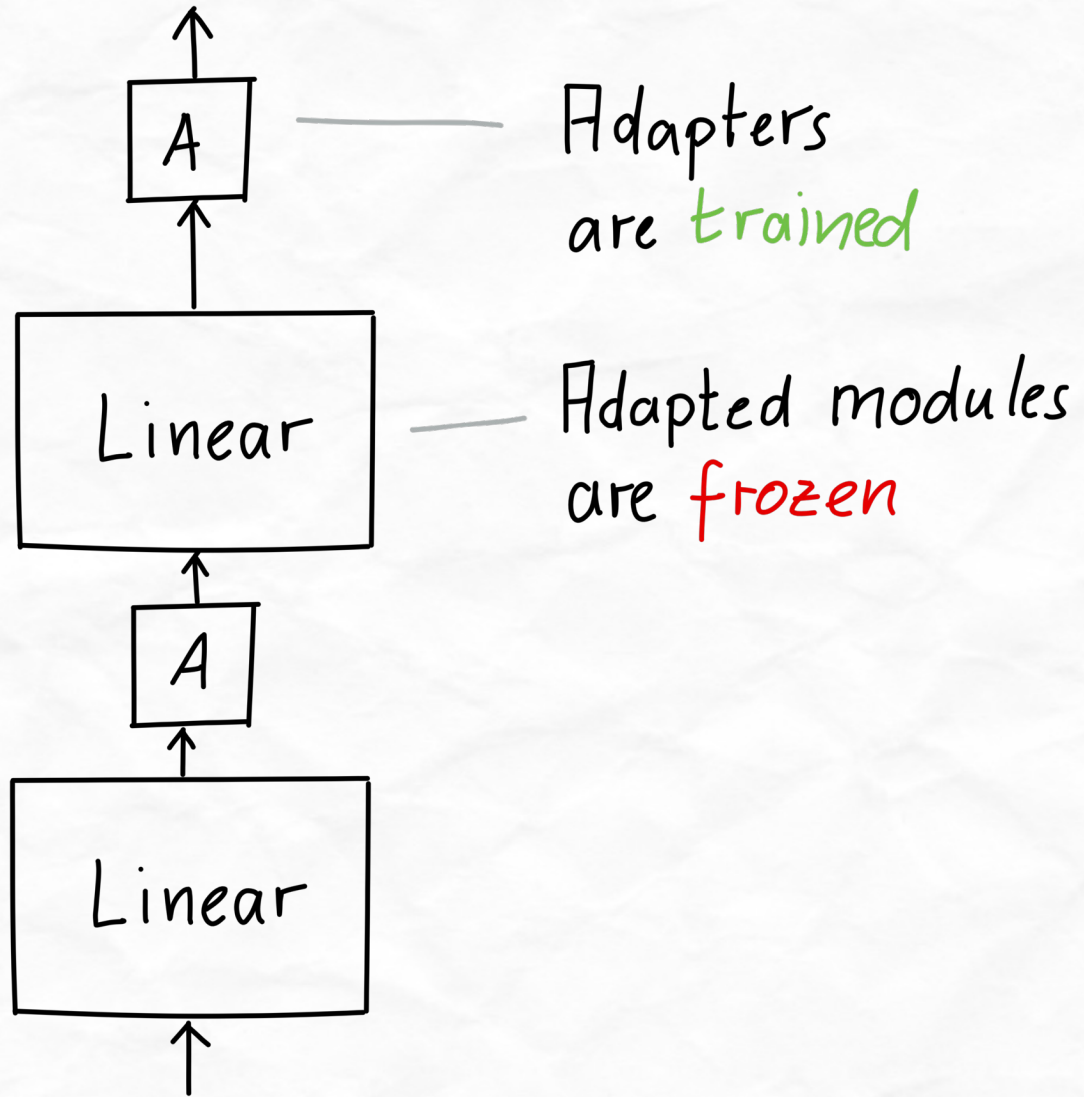
INTRINSIC DIMENSIONALITY EXPLAINS THE EFFECTIVENESS OF LANGUAGE MODEL FINE-TUNING
Aghajanyan, Zettlemoyer, Gupta, 2020

PARAMETER EFFICIENT FINE-TUNING
(PEFT)

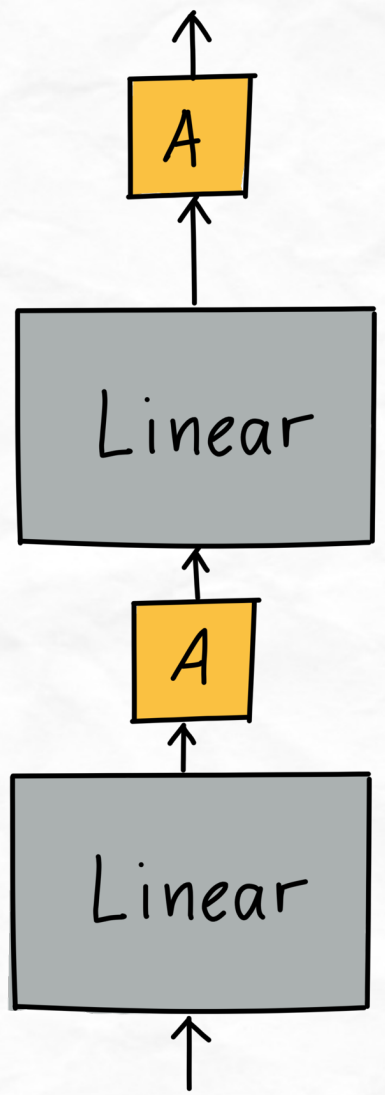
WHERE TO APPLY?

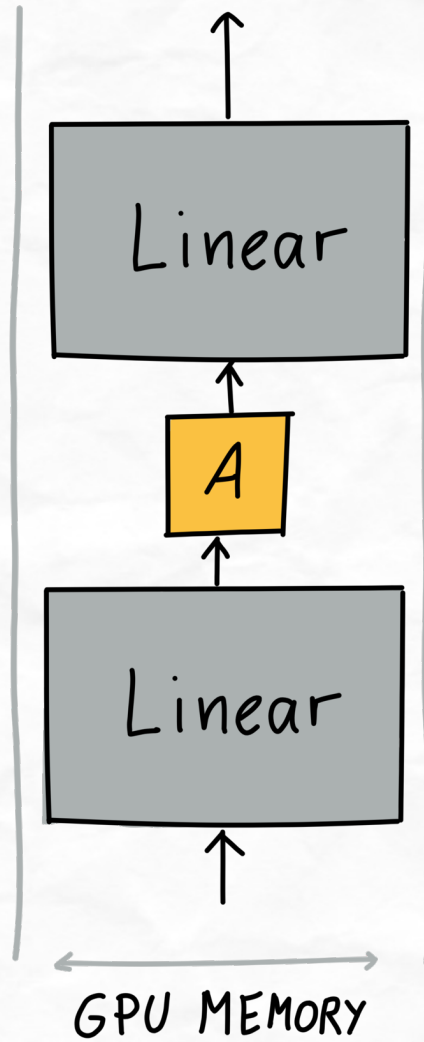


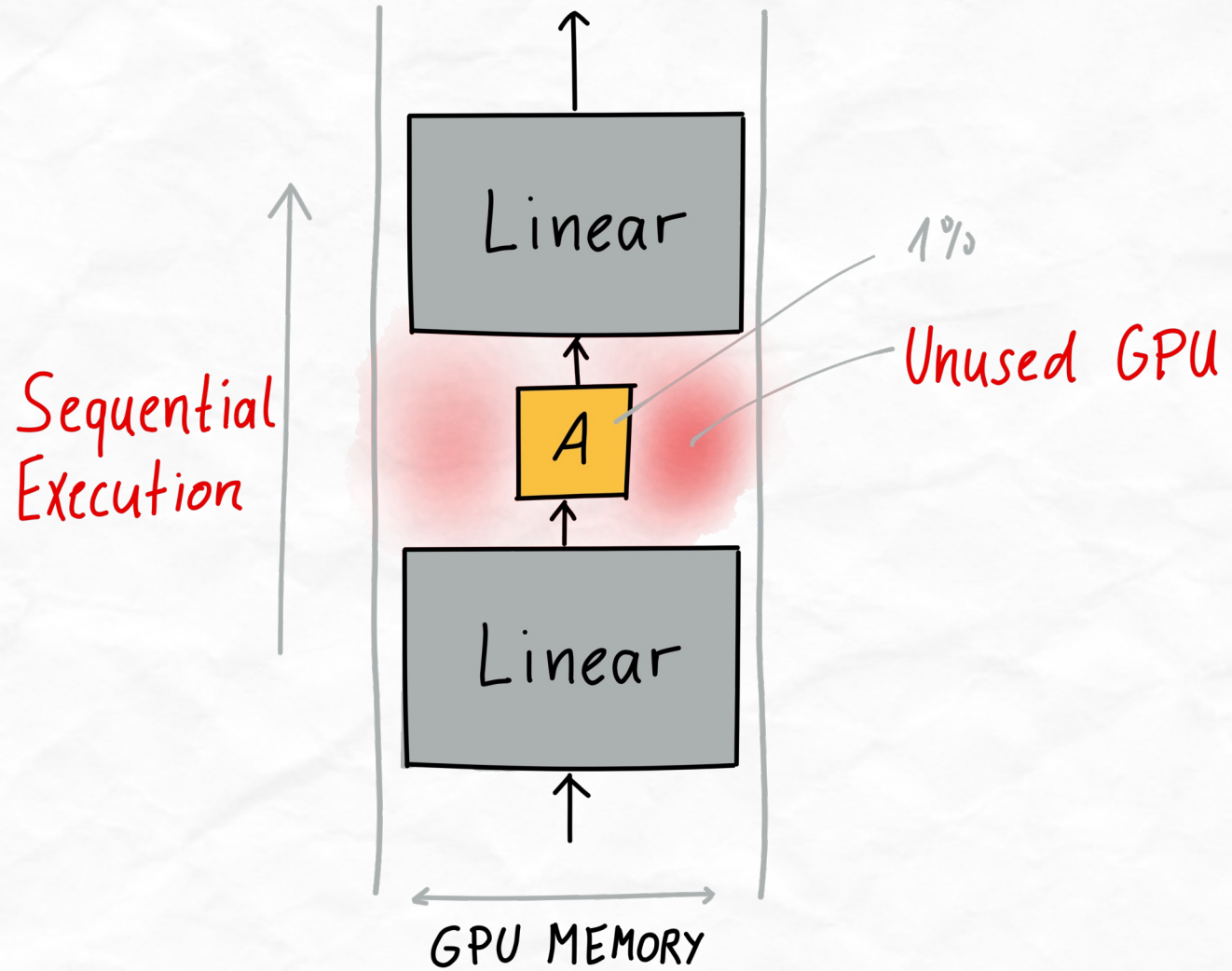
ADAPTER

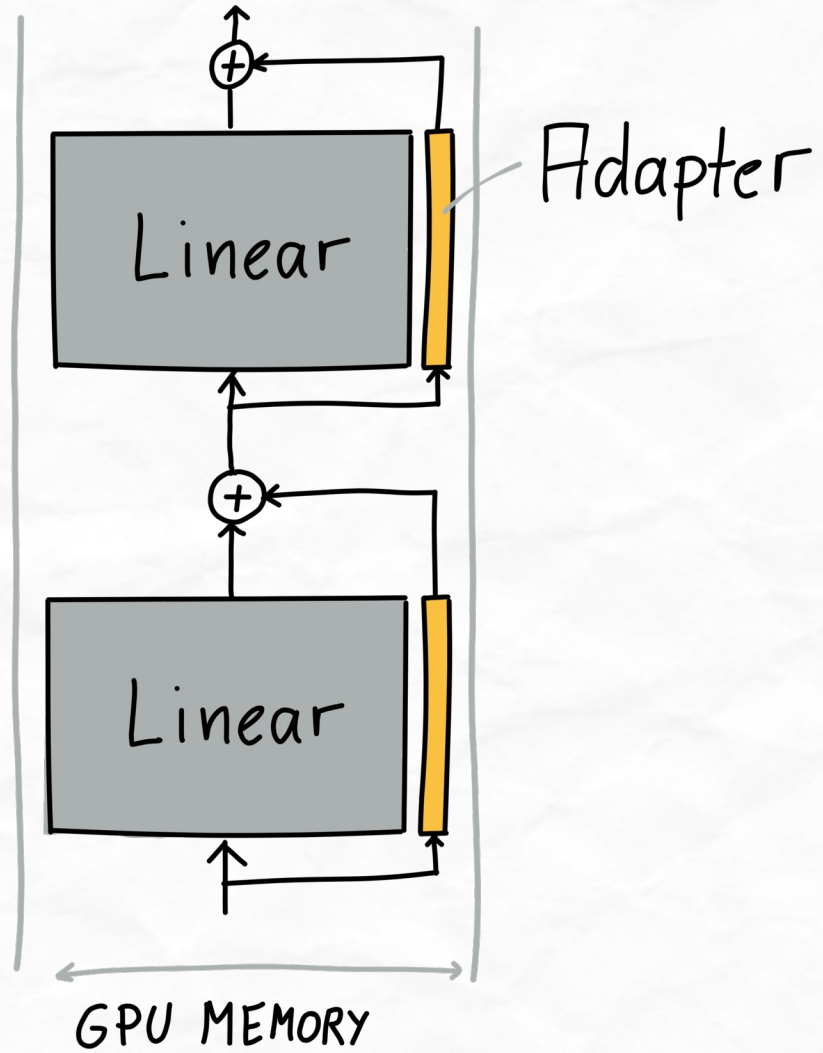


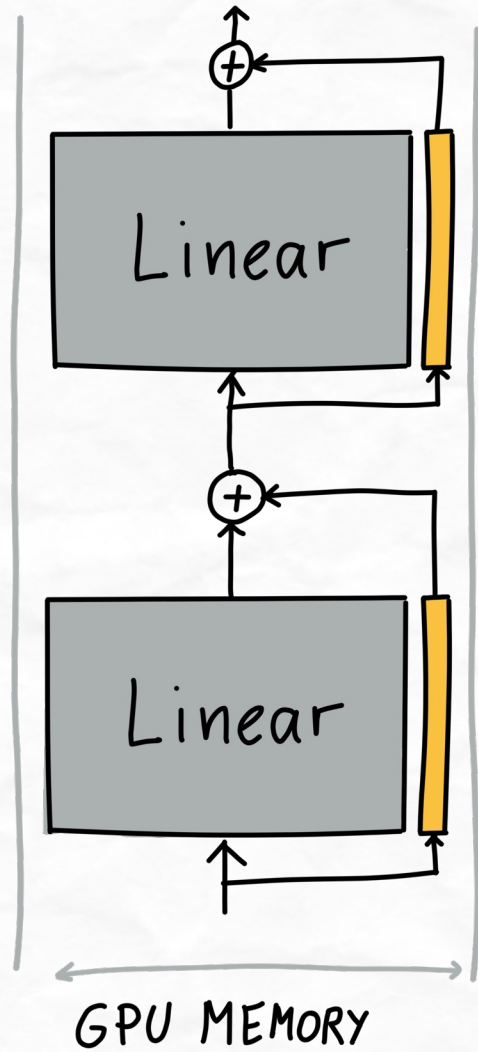
ADAPTER



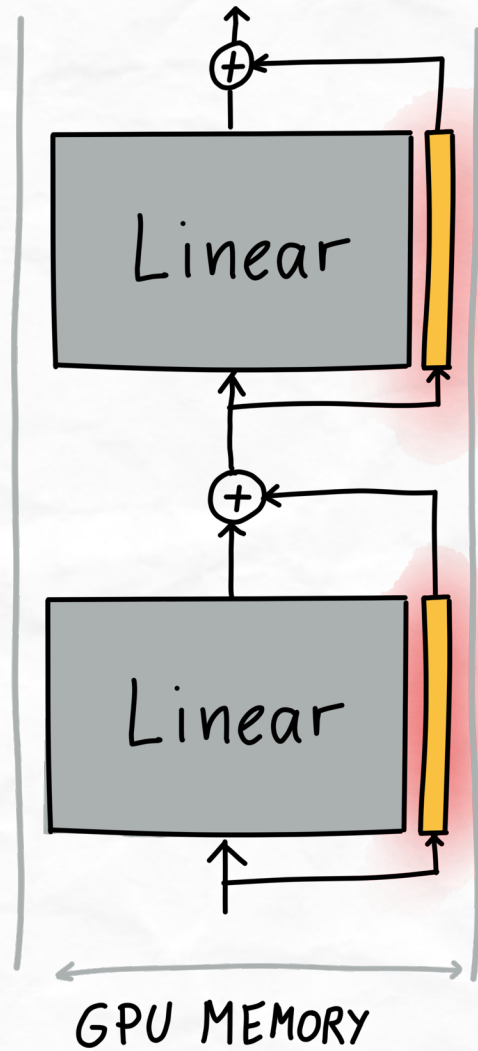




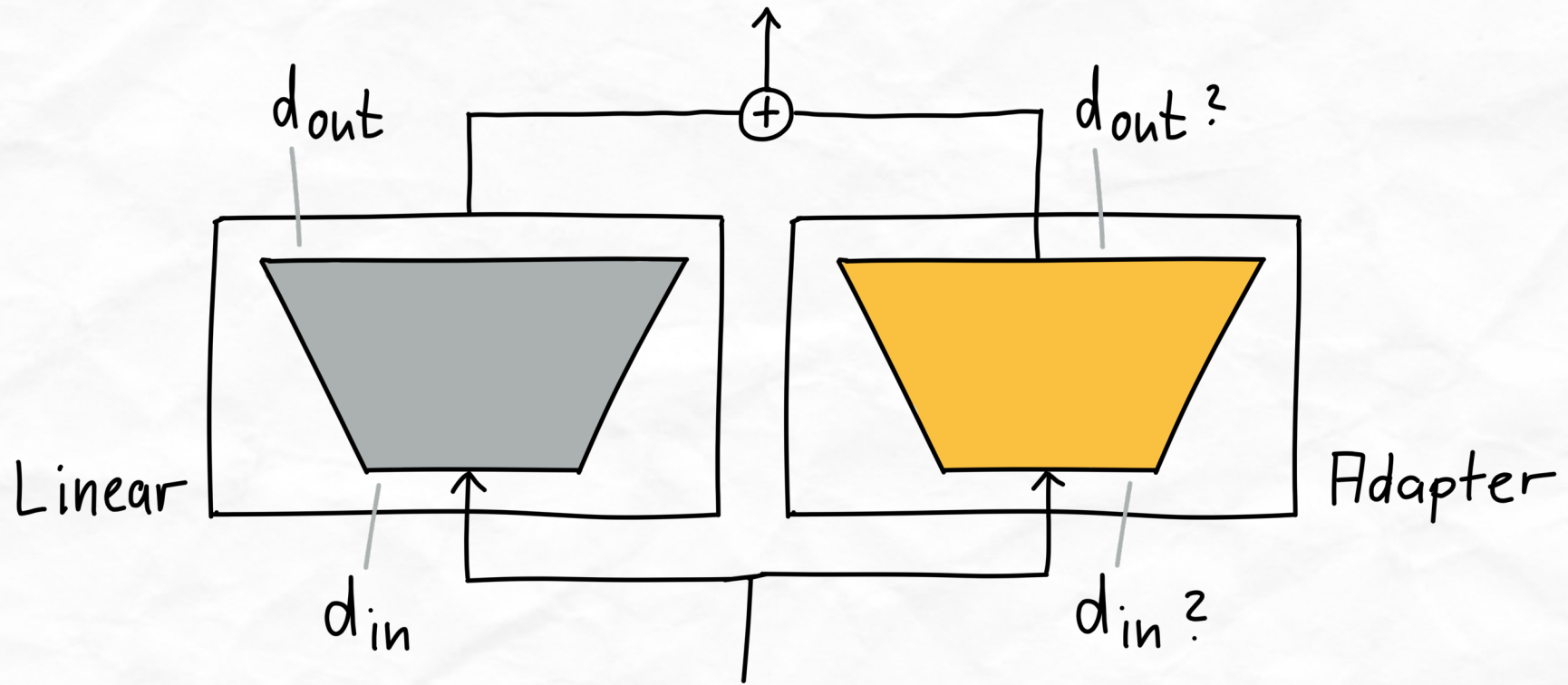


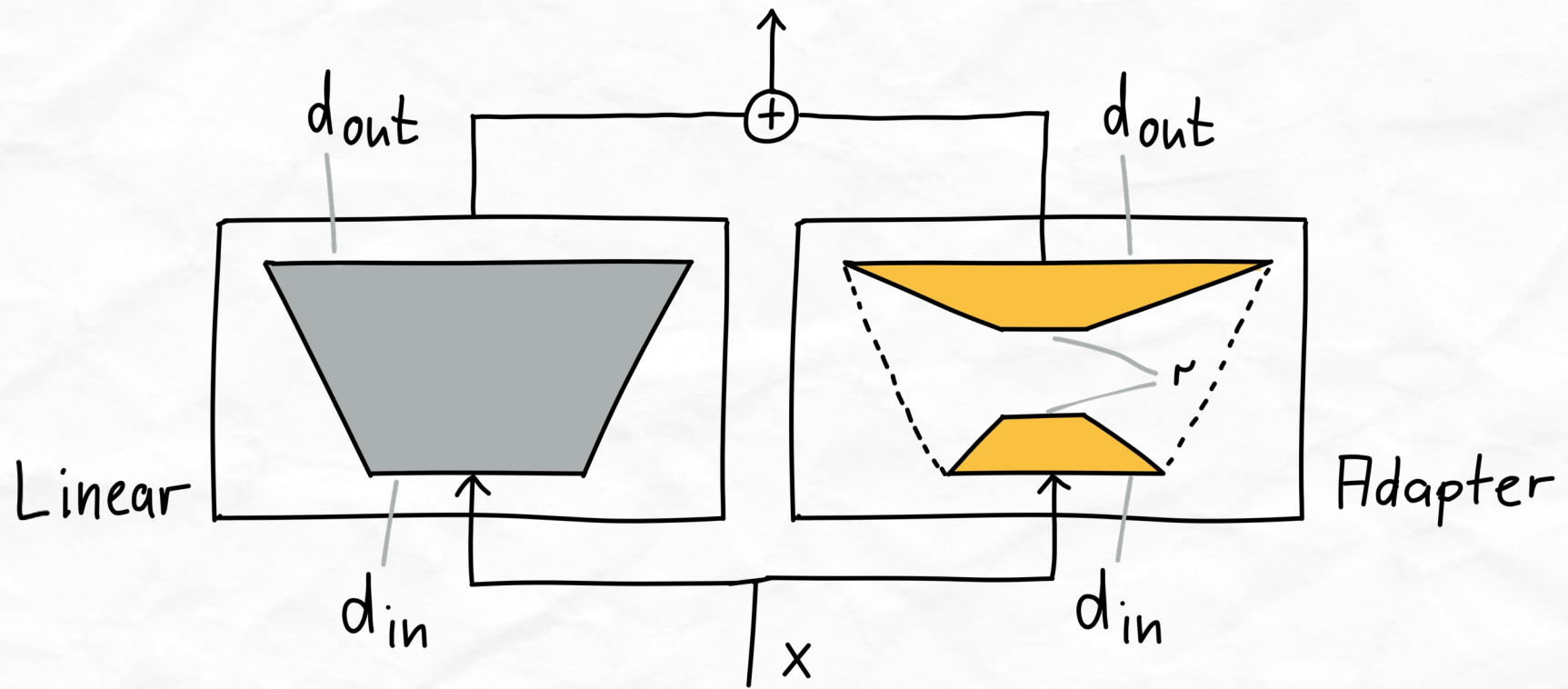


- Parallel Execution
- No Memory Bottleneck



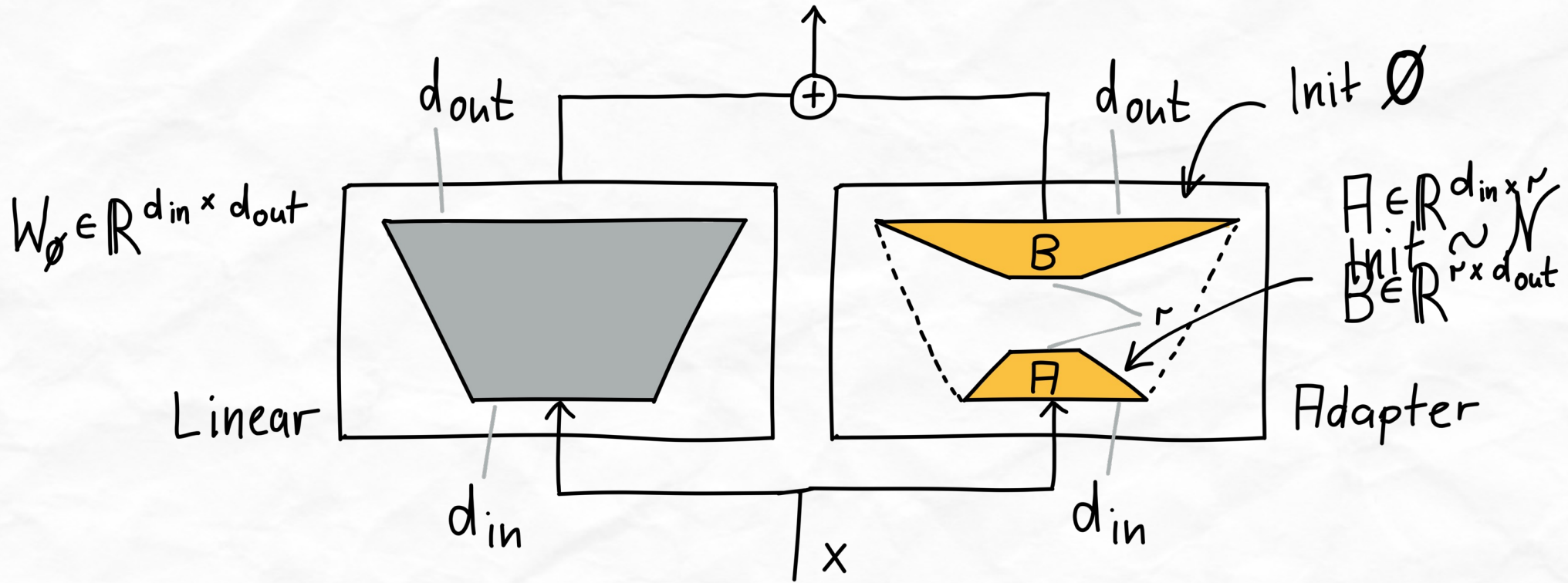
- Parallel Execution
- No Memory Bottleneck
- Memory & Execution Overhead (Inference)
- Adapter Size ?





LoRA

$$x' = x W_{\emptyset} + x AB$$



Hu, Chen et al. 2021

```

class LoRAAdapter(nn.Module):
    def __init__(self, adaptee, r):
        super().__init__()

        self.adaptee = adaptee
        self.r = r

        self.orig_forward = adaptee.forward
        adaptee.forward = self.forward

        adaptee.lora_A = nn.Parameter(
            torch.randn(adaptee.in_features, r) / math.sqrt(adaptee.in_features)
        )
        adaptee.lora_B = nn.Parameter(torch.zeros(r, adaptee.out_features))

    def forward(self, x, *args, **kwargs):
        return (
            self.orig_forward(x, *args, **kwargs)
            + x @ self.adaptee.lora_A @ self.adaptee.lora_B
        )

```

```

class LoRAAdapter(nn.Module):
    def __init__(self, adaptee, r):
        super().__init__()

        self.adaptee = adaptee
        self.r = r

        self.orig_forward = adaptee.forward
        adaptee.forward = self.forward

        adaptee.lora_A = nn.Parameter(
            torch.randn(adaptee.in_features, r) / math.sqrt(adaptee.in_features)
        )
        adaptee.lora_B = nn.Parameter(torch.zeros(r, adaptee.out_features))

    def forward(self, x, *args, **kwargs):
        return (
            self.orig_forward(x, *args, **kwargs)
            + x @ self.adaptee.lora_A @ self.adaptee.lora_B
        )

```

```

class LoRAAdapter(nn.Module):
    def __init__(self, adaptee, r):
        super().__init__()

        self.adaptee = adaptee
        self.r = r

        self.orig_forward = adaptee.forward
        adaptee.forward = self.forward

        adaptee.lora_A = nn.Parameter(
            torch.randn(adaptee.in_features, r) / math.sqrt(adaptee.in_features)
        )
        adaptee.lora_B = nn.Parameter(torch.zeros(r, adaptee.out_features))

    def forward(self, x, *args, **kwargs):
        return (
            self.orig_forward(x, *args, **kwargs)
            + x @ self.adaptee.lora_A @ self.adaptee.lora_B
        )

```

```

class LoRAAdapter(nn.Module):
    def __init__(self, adaptee, r):
        super().__init__()

        self.adaptee = adaptee
        self.r = r

        self.orig_forward = adaptee.forward
        adaptee.forward = self.forward

        adaptee.lora_A = nn.Parameter(
            torch.randn(adaptee.in_features, r) / math.sqrt(adaptee.in_features)
        )
        adaptee.lora_B = nn.Parameter(torch.zeros(r, adaptee.out_features))

    def forward(self, x, *args, **kwargs):
        return (
            self.orig_forward(x, *args, **kwargs)
            + x @ self.adaptee.lora_A @ self.adaptee.lora_B
        )

```

```

class LoRAAdapter(nn.Module):
    def __init__(self, adaptee, r):
        super().__init__()

        self.adaptee = adaptee
        self.r = r

        self.orig_forward = adaptee.forward
        adaptee.forward = self.forward

        adaptee.lora_A = nn.Parameter(
            torch.randn(adaptee.in_features, r) / math.sqrt(adaptee.in_features)
        )
        adaptee.lora_B = nn.Parameter(torch.zeros(r, adaptee.out_features))

    def forward(self, x, *args, **kwargs):
        return (
            self.orig_forward(x, *args, **kwargs)
            + x @ self.adaptee.lora_A @ self.adaptee.lora_B
        )

```

```

class LoRAAdapter(nn.Module):
    def __init__(self, adaptee, r):
        super().__init__()

        self.adaptee = adaptee
        self.r = r

        self.orig_forward = adaptee.forward
        adaptee.forward = self.forward

        adaptee.lora_A = nn.Parameter(
            torch.randn(adaptee.in_features, r) / math.sqrt(adaptee.in_features)
        )
        adaptee.lora_B = nn.Parameter(torch.zeros(r, adaptee.out_features))

def forward(self, x, *args, **kwargs):
    return (
        self.orig_forward(x, *args, **kwargs)
        + x @ self.adaptee.lora_A @ self.adaptee.lora_B
    )

```

```
class LoRAAdapter(nn.Module):
    def __init__(self, adaptee, r):
        super().__init__()

        self.adaptee = adaptee
        self.r = r

        self.orig_forward = adaptee.forward
        adaptee.forward = self.forward

        adaptee.lora_A = nn.Parameter(
            torch.randn(adaptee.in_features, r) / math.sqrt(adaptee.in_features)
        )
        adaptee.lora_B = nn.Parameter(torch.zeros(r, adaptee.out_features))

    def forward(self, x, *args, **kwargs):
        return (
            self.orig_forward(x, *args, **kwargs)
            + x @ self.adaptee.lora_A @ self.adaptee.lora_B
        )
```

QLoRA:
LoRA + Quantization



```
model = AutoModelForSequenceClassification.from_pretrained(
    "tiiuae/falcon-7b-instruct",
    num_labels=2,
    load_in_4bit=True
)

model = prepare_model_for_kbit_training(model)

peft_config = LoraConfig(
    task_type=TaskType.SEQ_CLS,
    r=args.lora_r,
    target_modules=["query_key_value", "dense_h_to_4h", "dense_4h_to_h"]
)

model = get_peft_model(model, peft_config)
```

QLoRA:
LoRA + Quantization



```
model = AutoModelForSequenceClassification.from_pretrained(
    "tiiuae/falcon-7b-instruct",
    num_labels=2,
    load_in_4bit=True
)

model = prepare_model_for_kbit_training(model)

peft_config = LoraConfig(
    task_type=TaskType.SEQ_CLS,
    r=args.lora_r,
    target_modules=["query_key_value", "dense_h_to_4h", "dense_4h_to_h"]
)

model = get_peft_model(model, peft_config)
```

QLoRA:
LoRA + Quantization



```
model = AutoModelForSequenceClassification.from_pretrained(  
    "tiiuae/falcon-7b-instruct",  
    num_labels=2,  
    load_in_4bit=True  
)
```

```
model = prepare_model_for_kbit_training(model)
```

```
peft_config = LoraConfig(  
    task_type=TaskType.SEQ_CLS,  
    r=args.lora_r,  
    target_modules=["query_key_value", "dense_h_to_4h", "dense_4h_to_h"]  
)
```

```
model = get_peft_model(model, peft_config)
```

QLoRA:
LoRA + Quantization



```
model = AutoModelForSequenceClassification.from_pretrained(
    "tiiuae/falcon-7b-instruct",
    num_labels=2,
    load_in_4bit=True
)
```

```
model = prepare_model_for_kbit_training(model)
```

```
peft_config = LoraConfig(
    task_type=TaskType.SEQ_CLS,
    r=8,
    target_modules=["query_key_value", "dense_h_to_4h", "dense_4h_to_h"]
)
```

```
model = get_peft_model(model, peft_config)
```

QLoRA:
LoRA + Quantization



```
model = AutoModelForSequenceClassification.from_pretrained(
    "tiiuae/falcon-7b-instruct",
    num_labels=2,
    load_in_4bit=True
)

model = prepare_model_for_kbit_training(model)

peft_config = LoraConfig(
    task_type=TaskType.SEQ_CLS,
    r=8,
    target_modules=["query_key_value", "dense_h_to_4h", "dense_4h_to_h"]
)

model = get_peft_model(model, peft_config)
```



QLoRA:
LoRA + Quantization



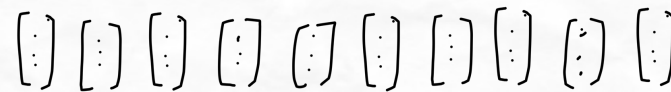
WHERE TO APPLY?



WHERE TO APPLY?



Alt: Q K V O FF: Up Down





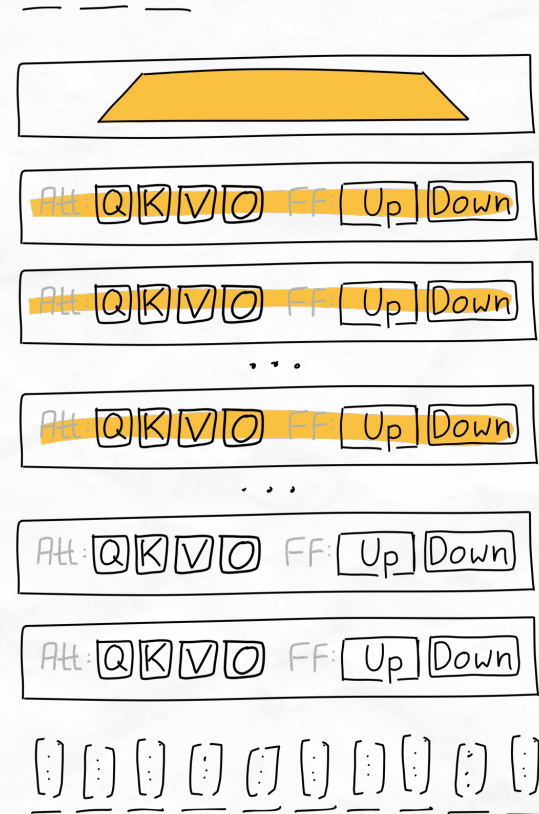
BY LAYERS



Even

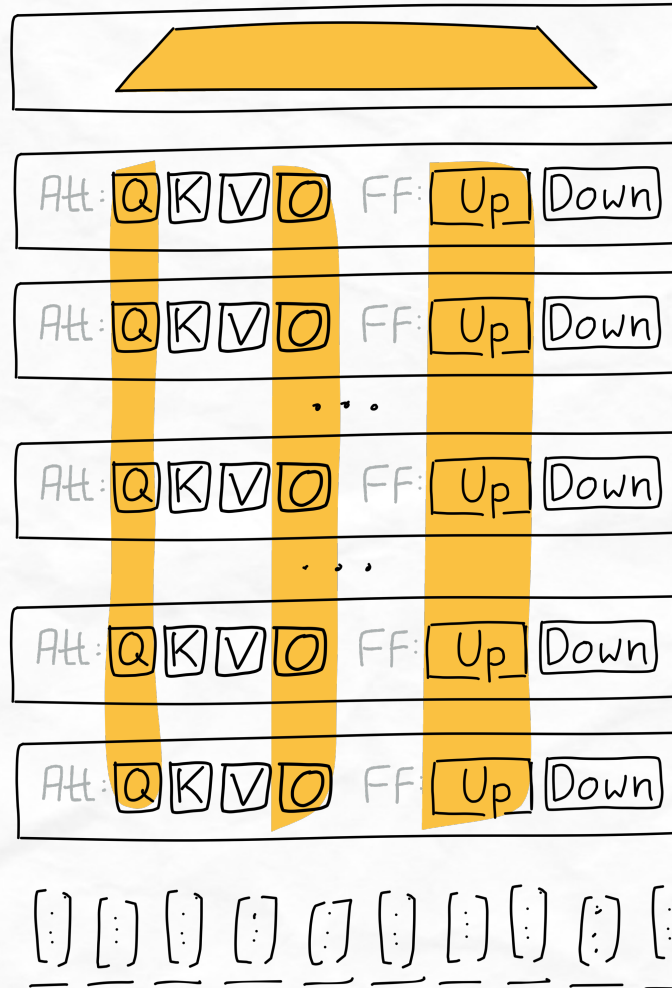


Lower

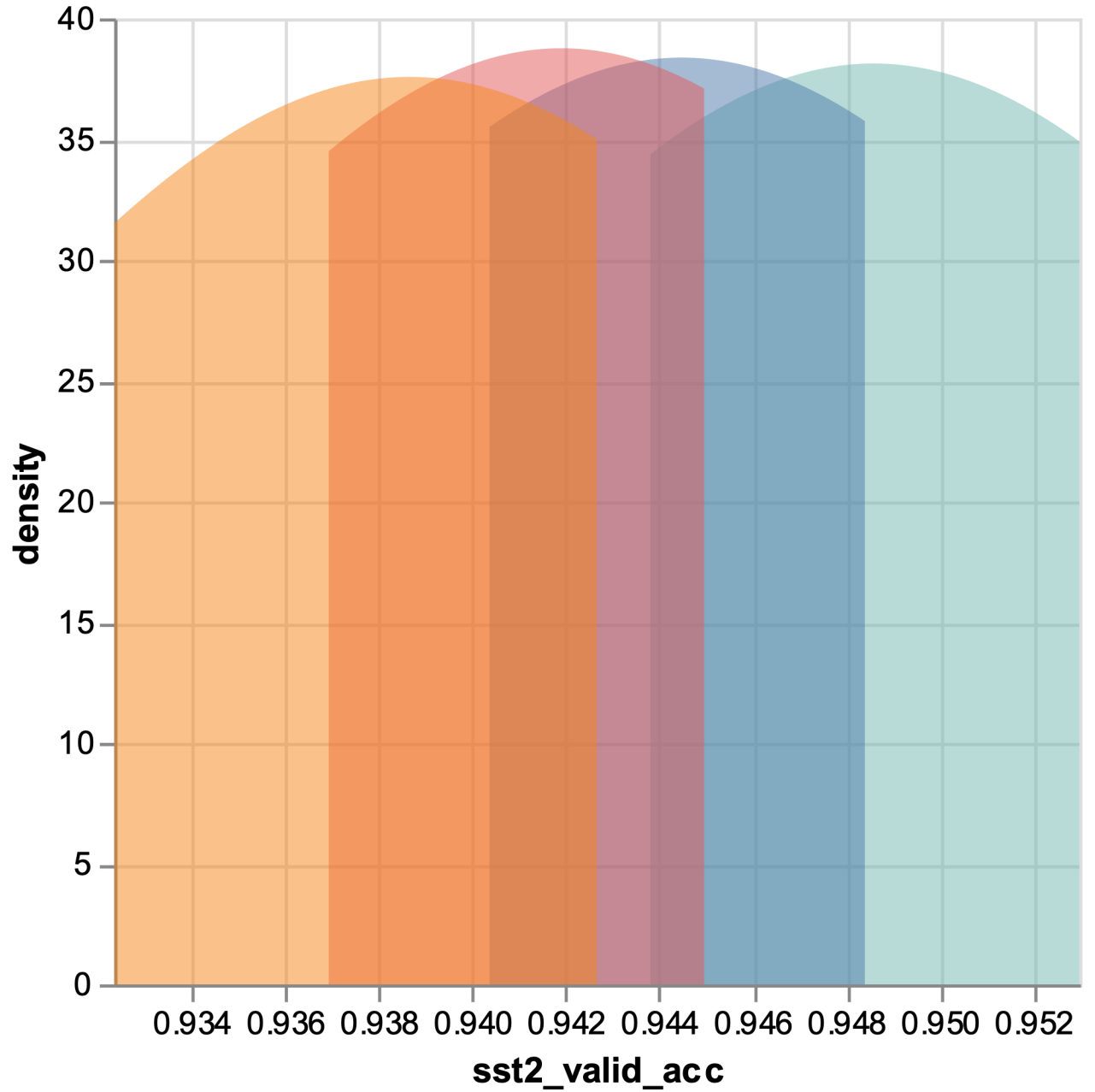
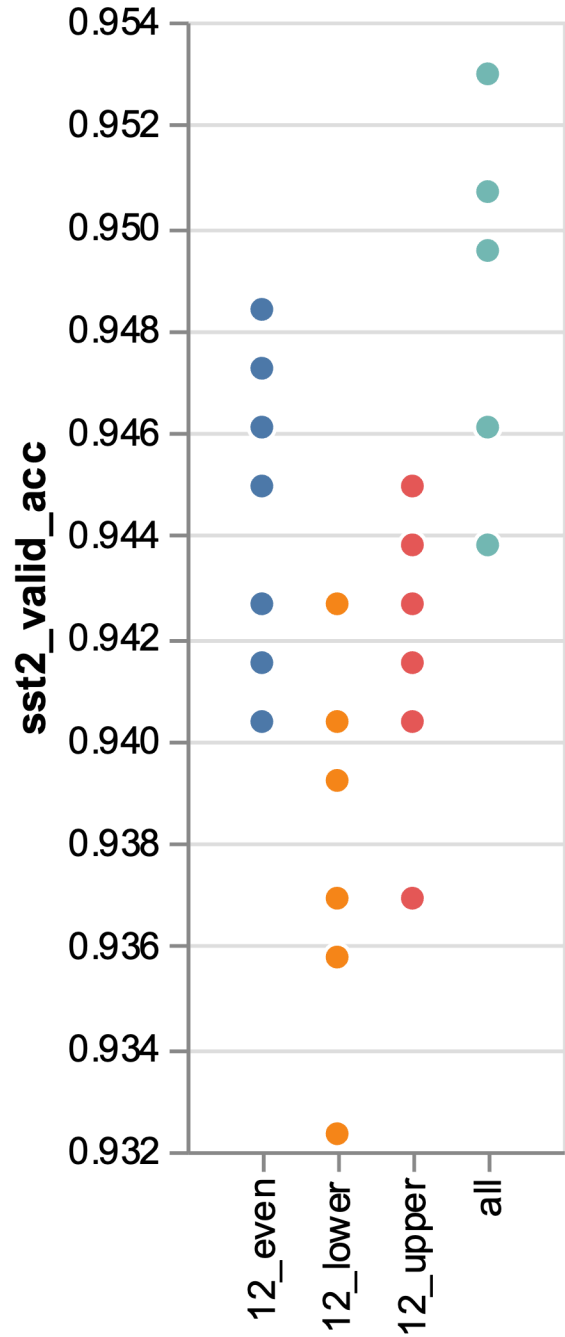


Upper

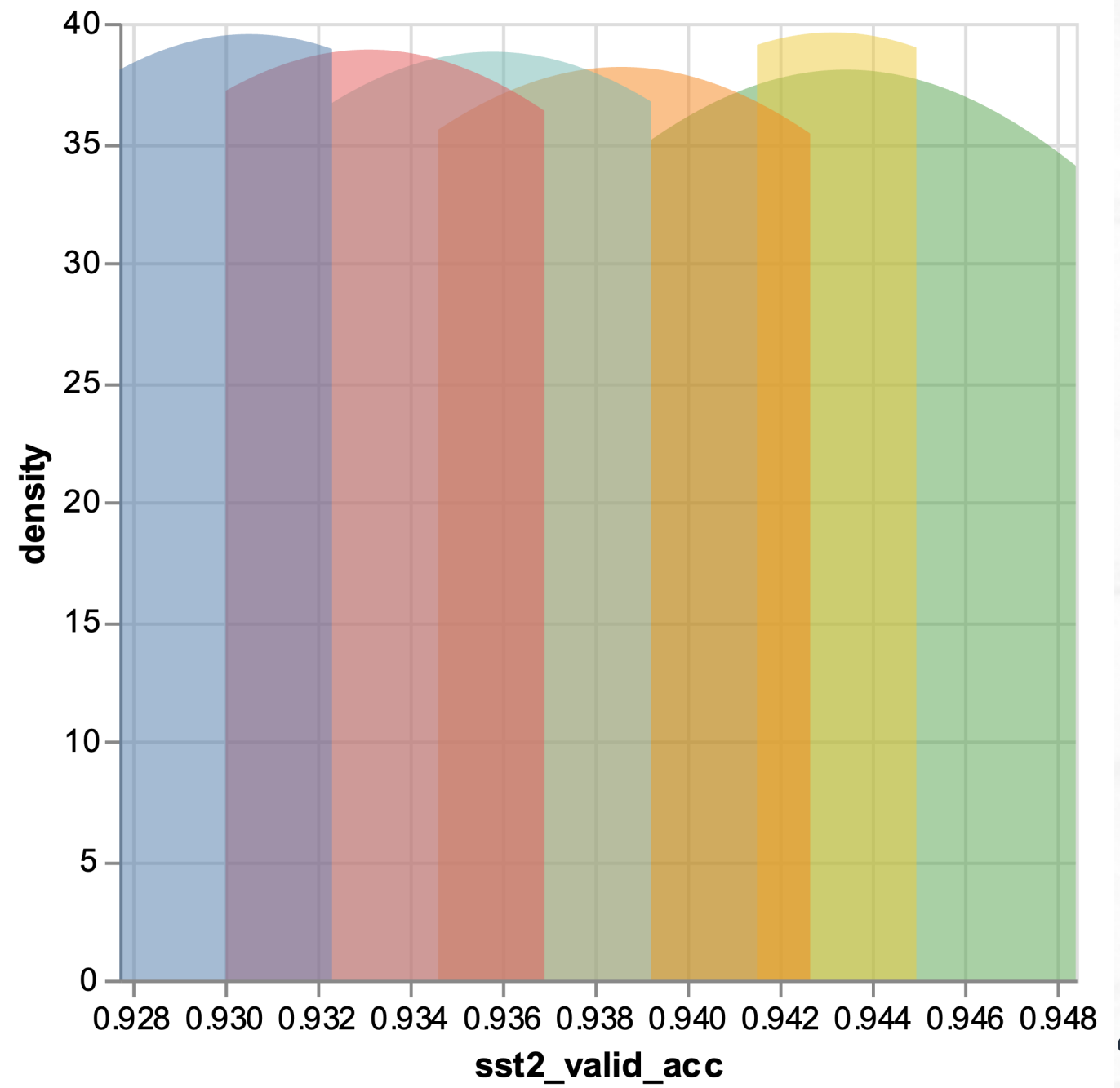
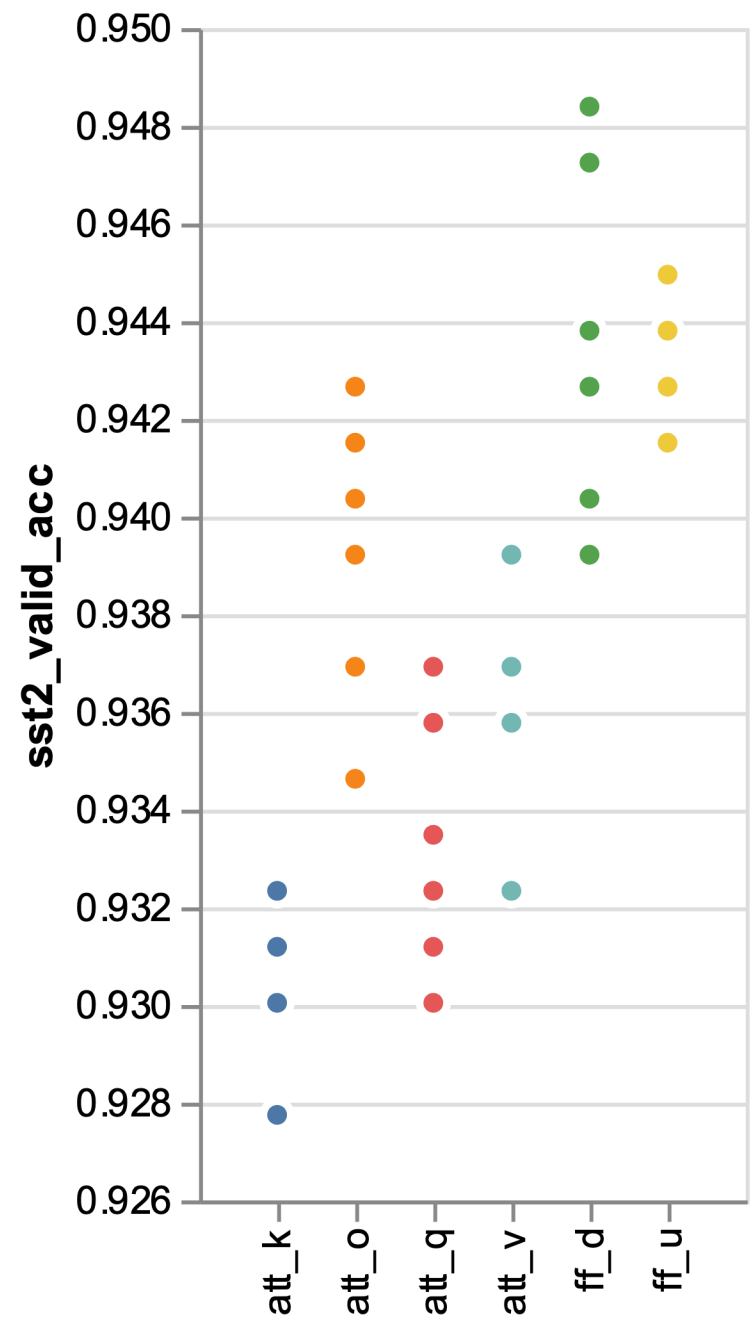
BY COMPONENT

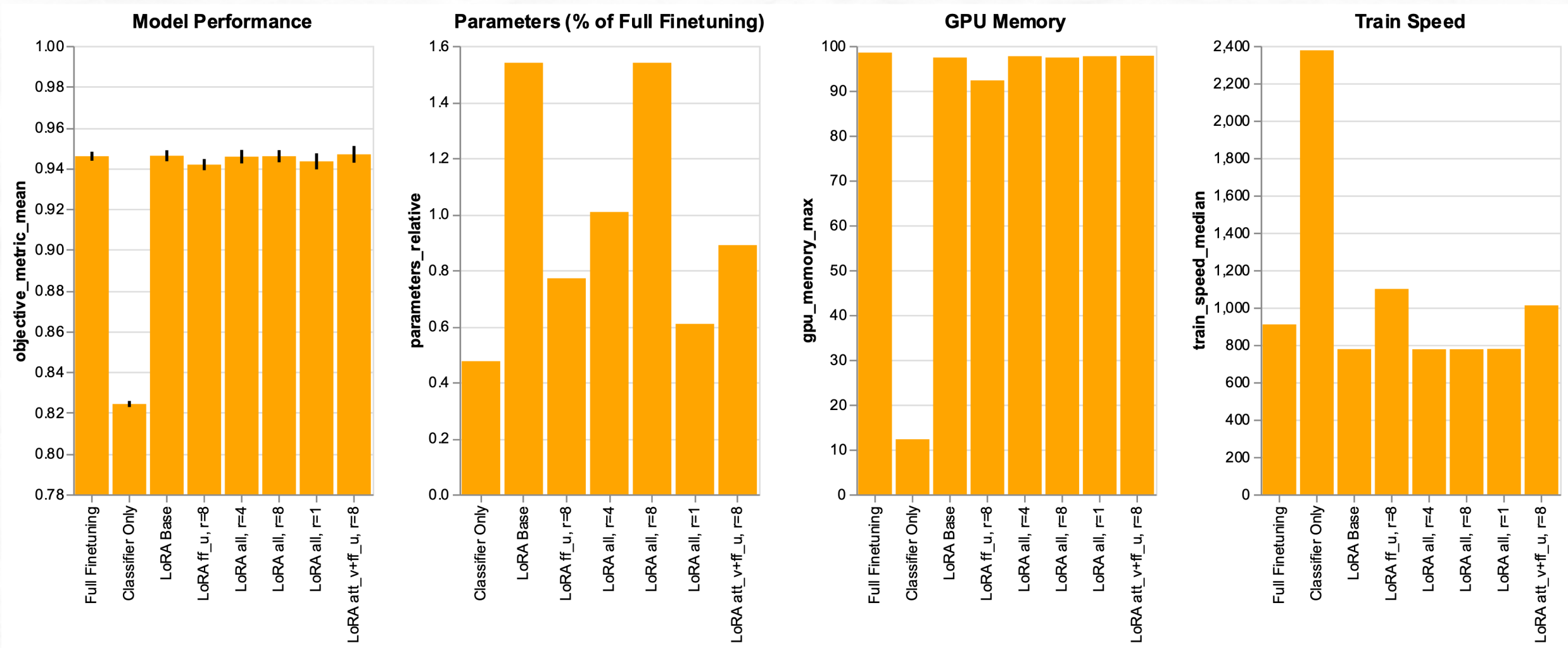


sst2-lora-config

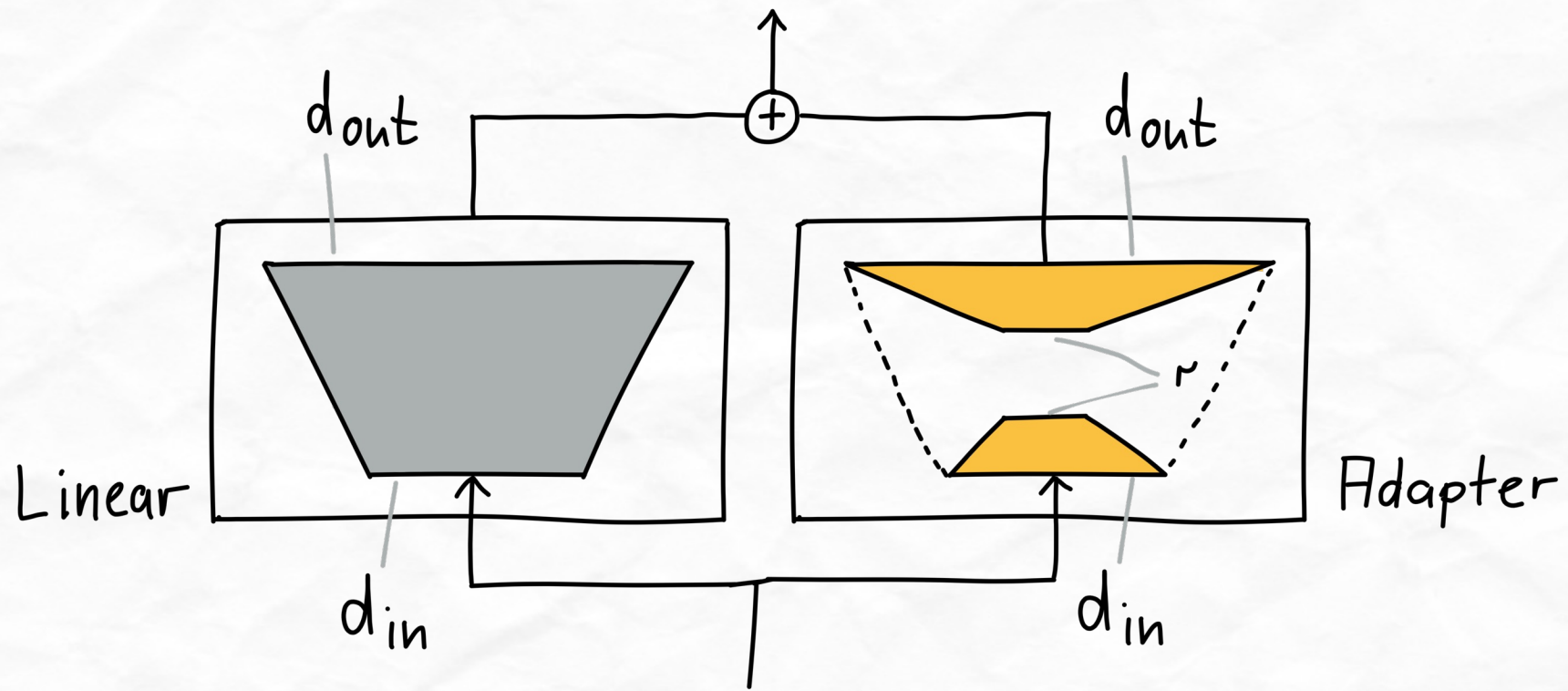


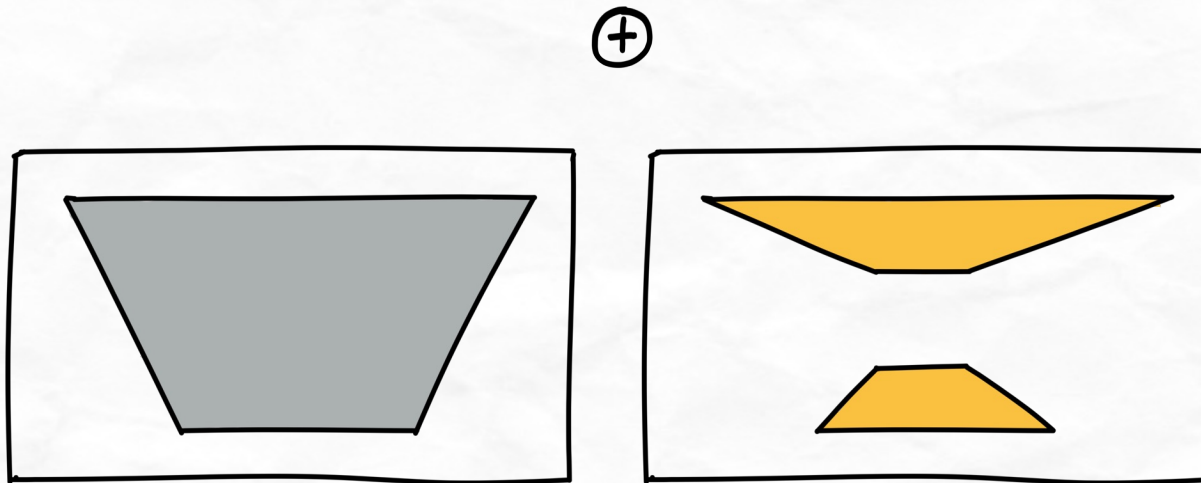
sst2-lora-config

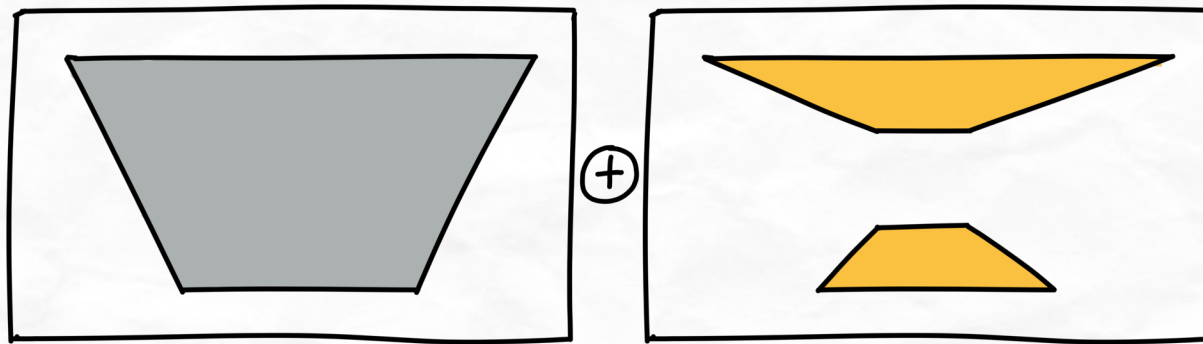


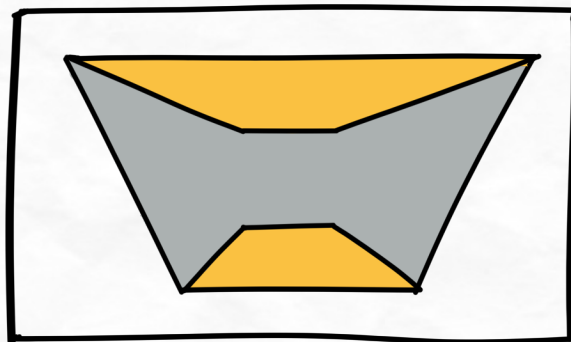


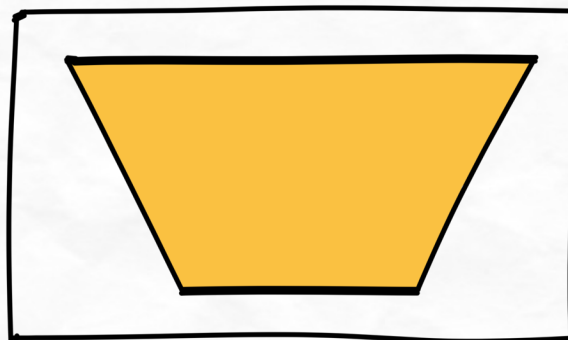
INFERENCE

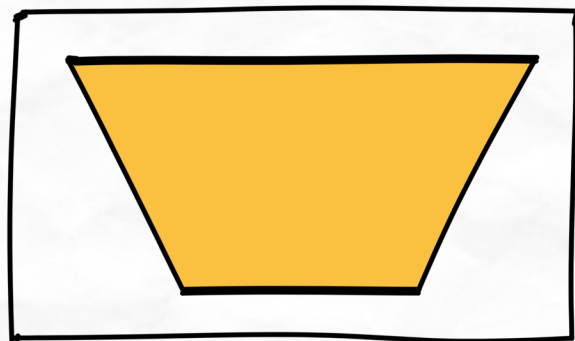








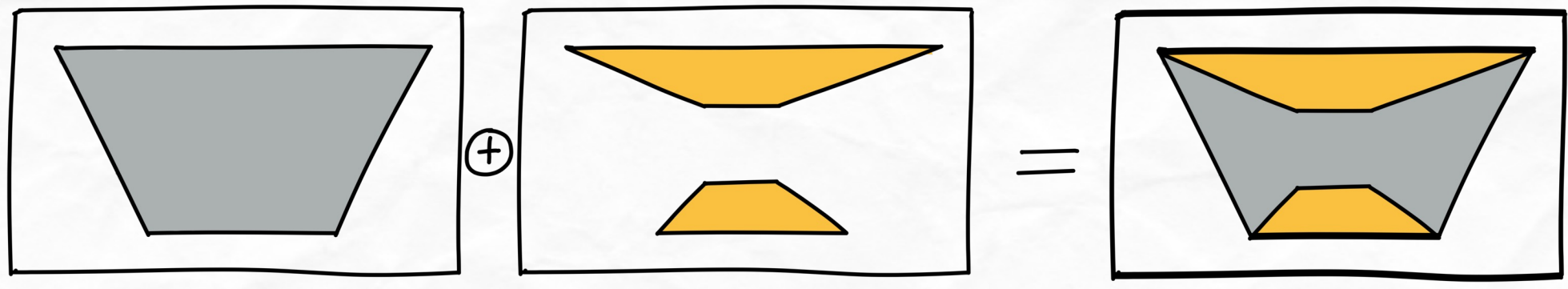




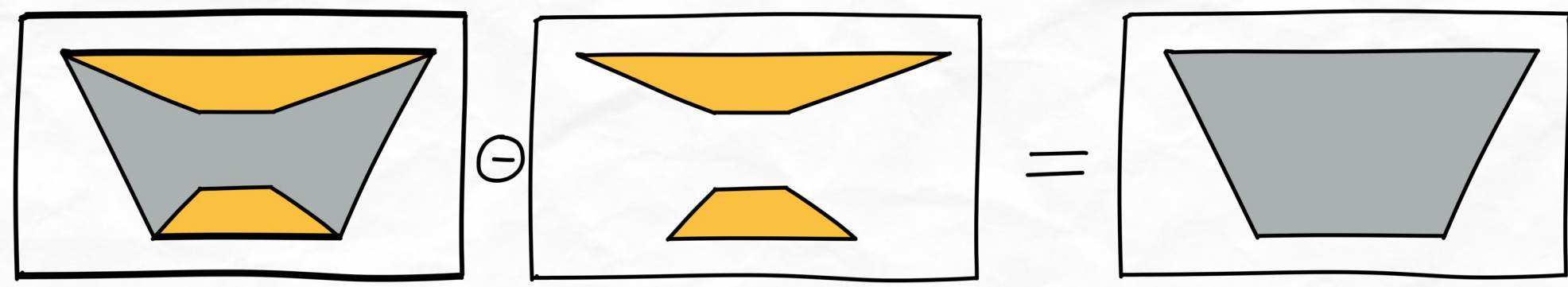
Baking in the product
of A and B:

- No special handling at inference
- No overhead in compute or memory
- Small footprint of finetuned weights

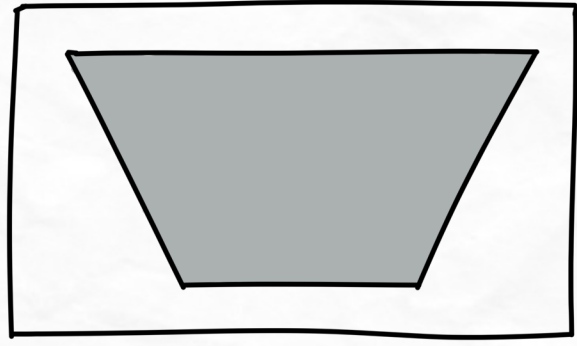
Adding Adapters



Removing Adapters



HOT SWAPPING

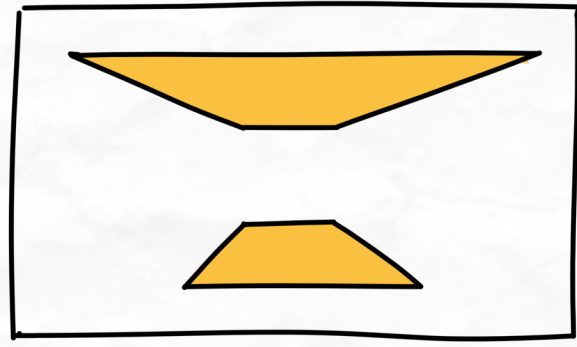


Pre-trained

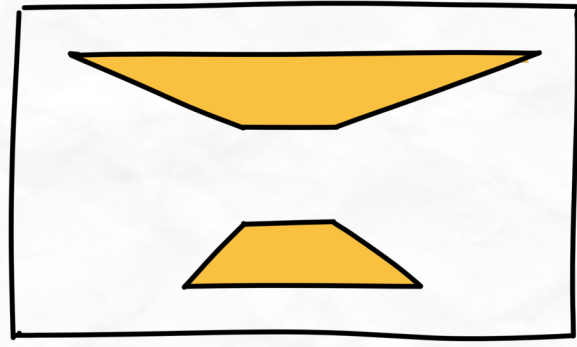
+

+

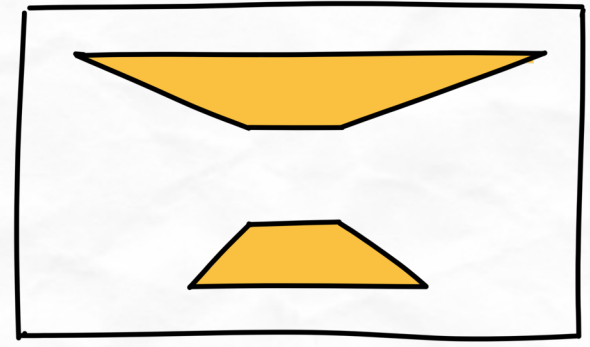
+



TASK A



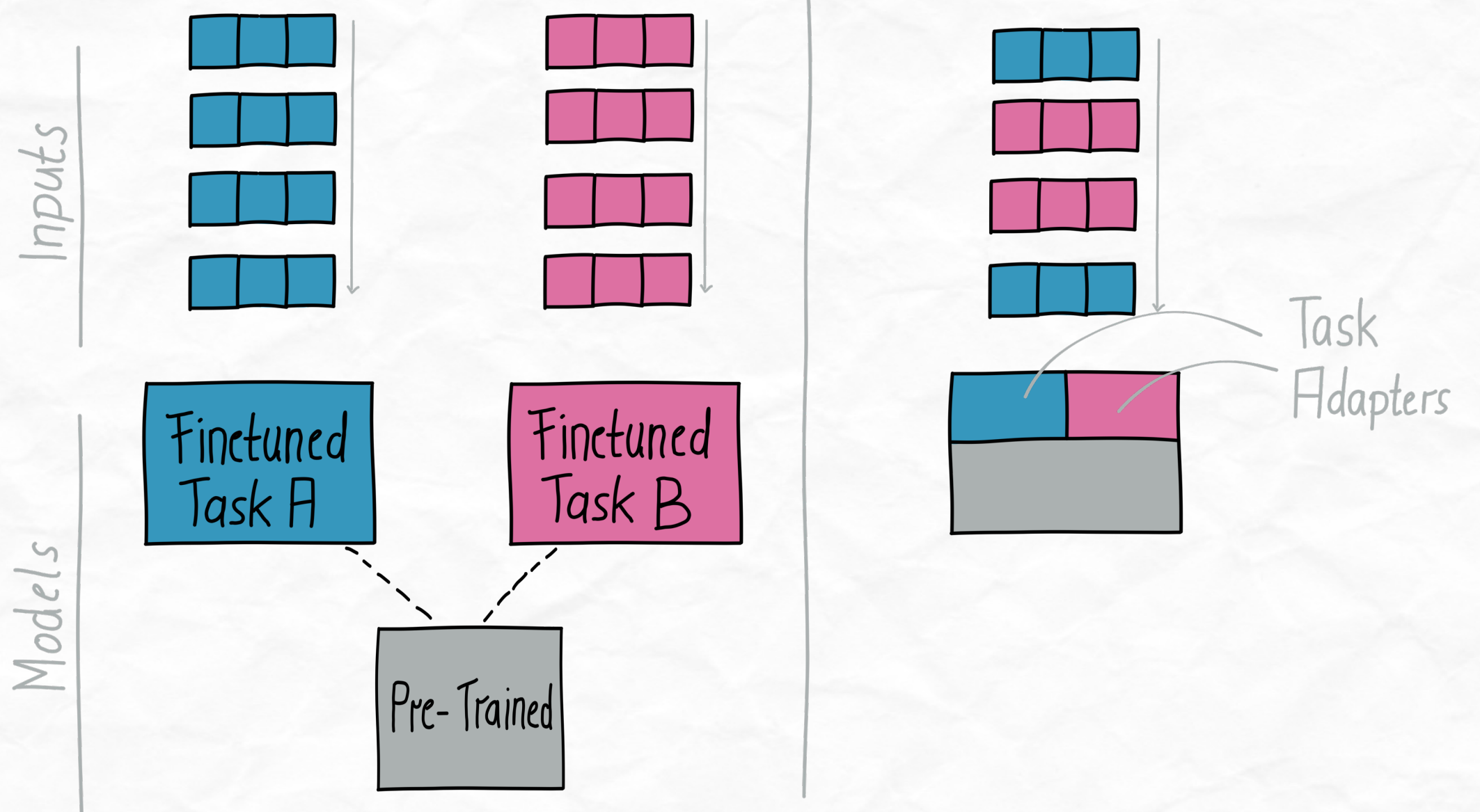
TASK B



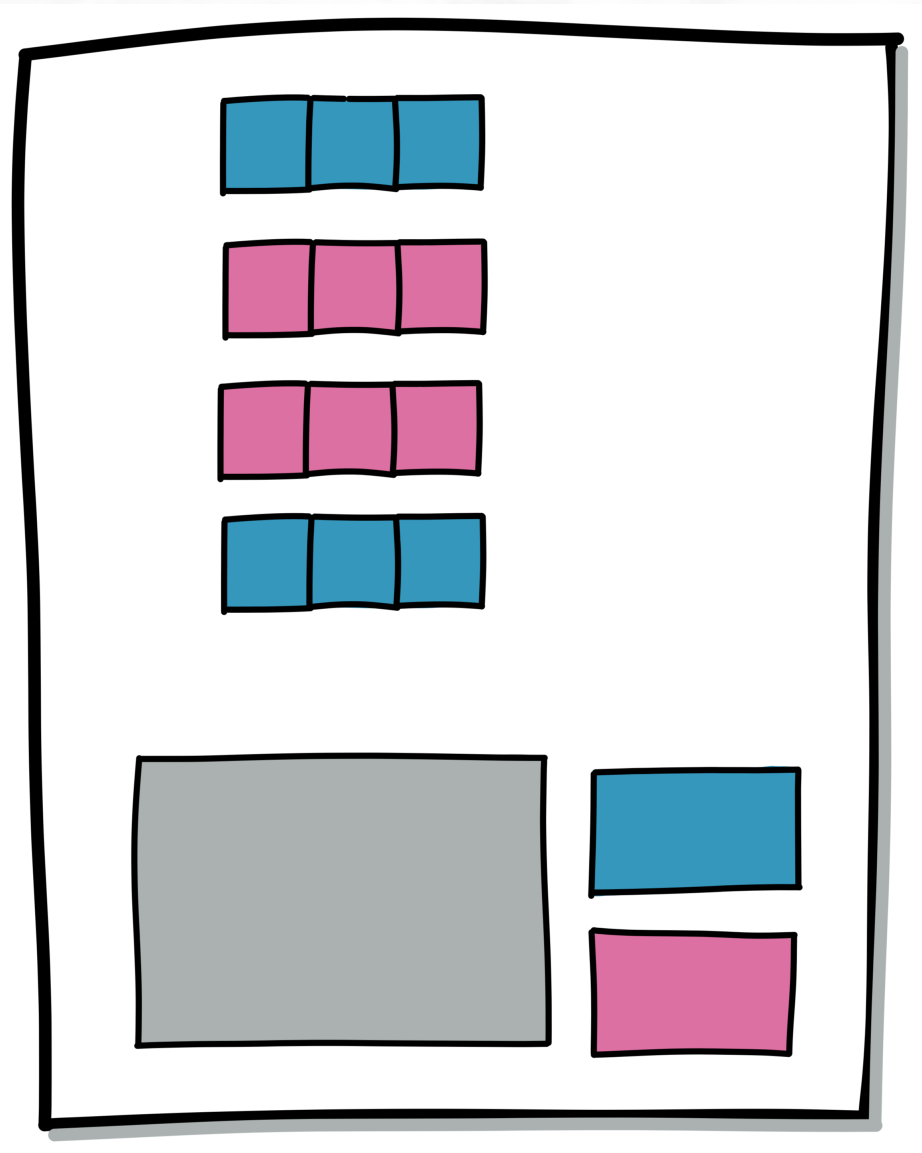
TASK C

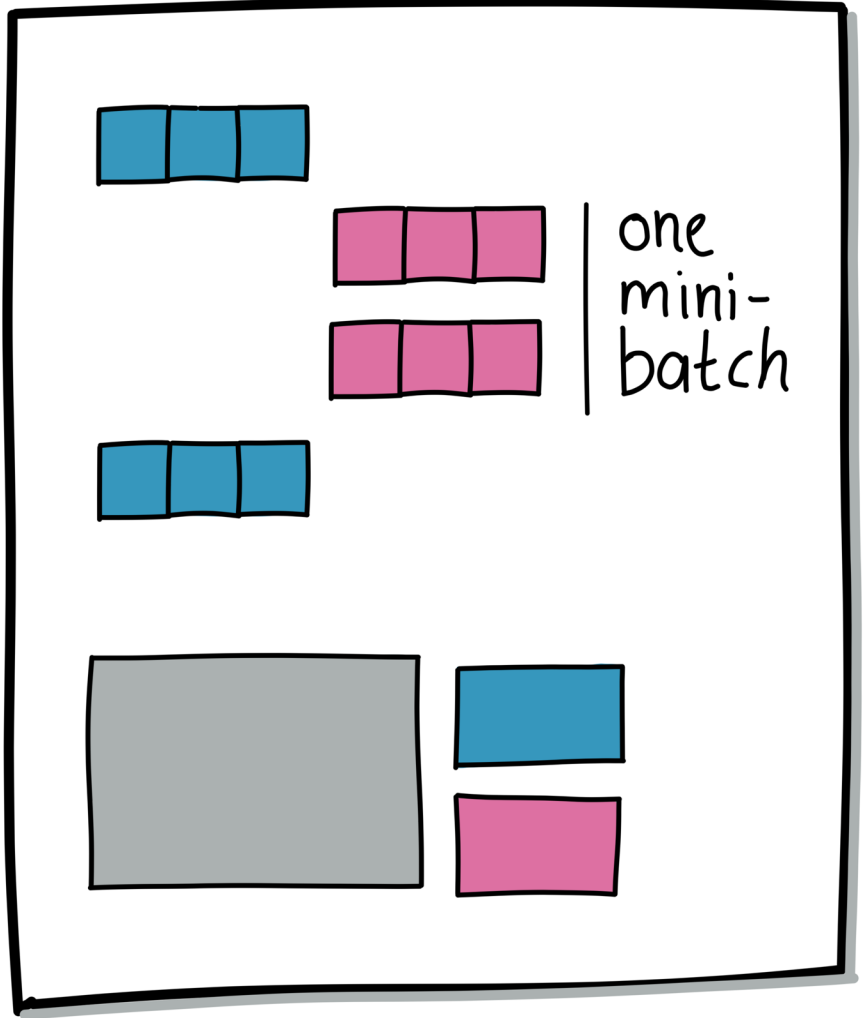
Fine-tuned

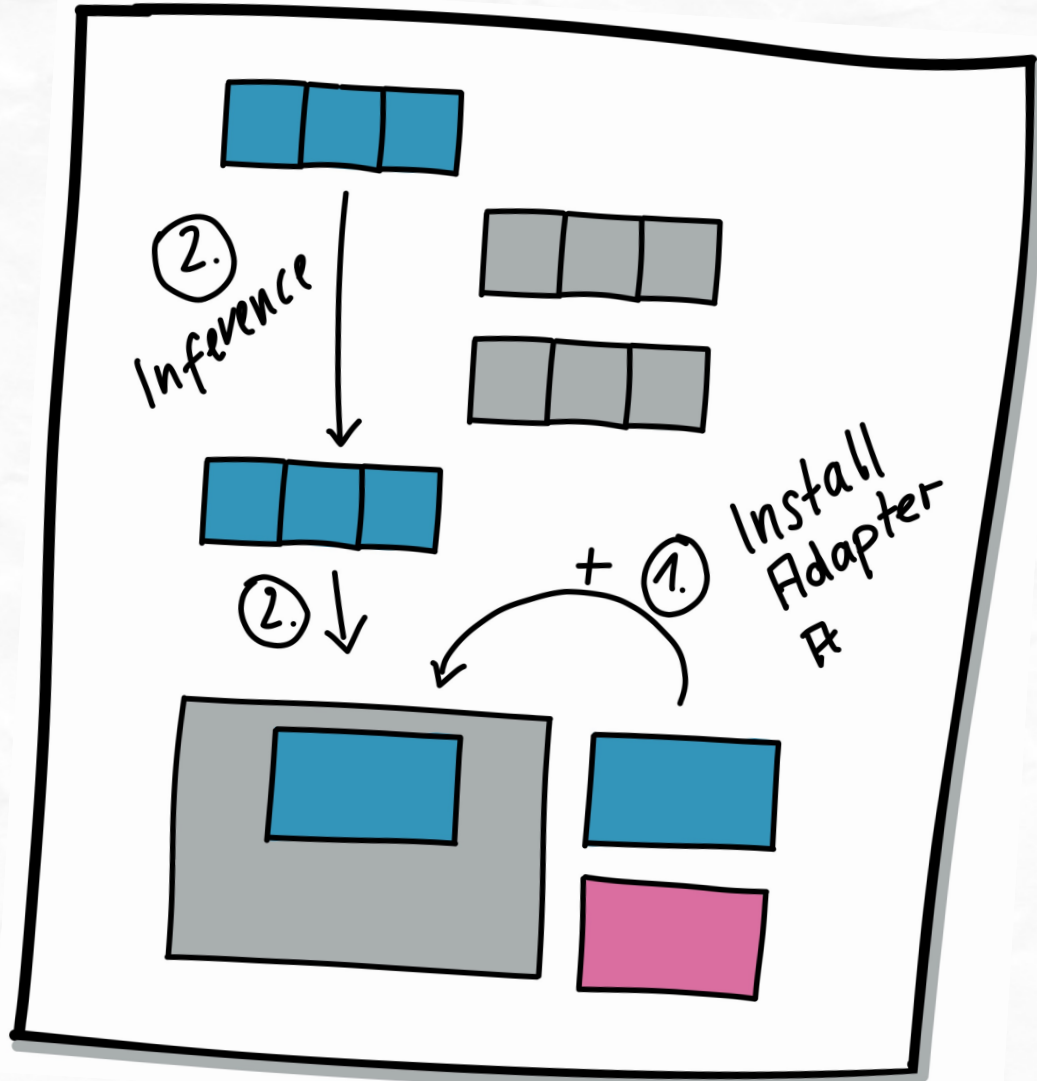
MULTI-TASK

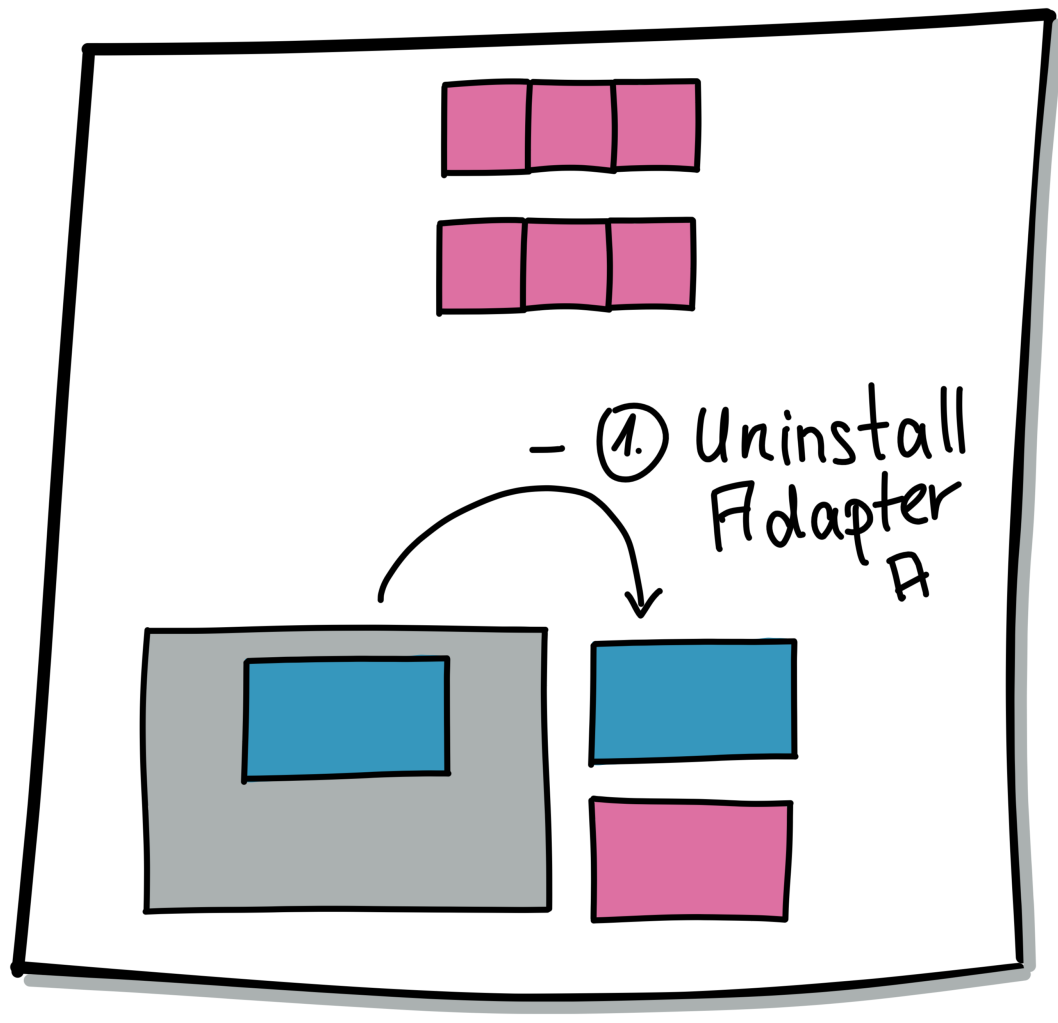


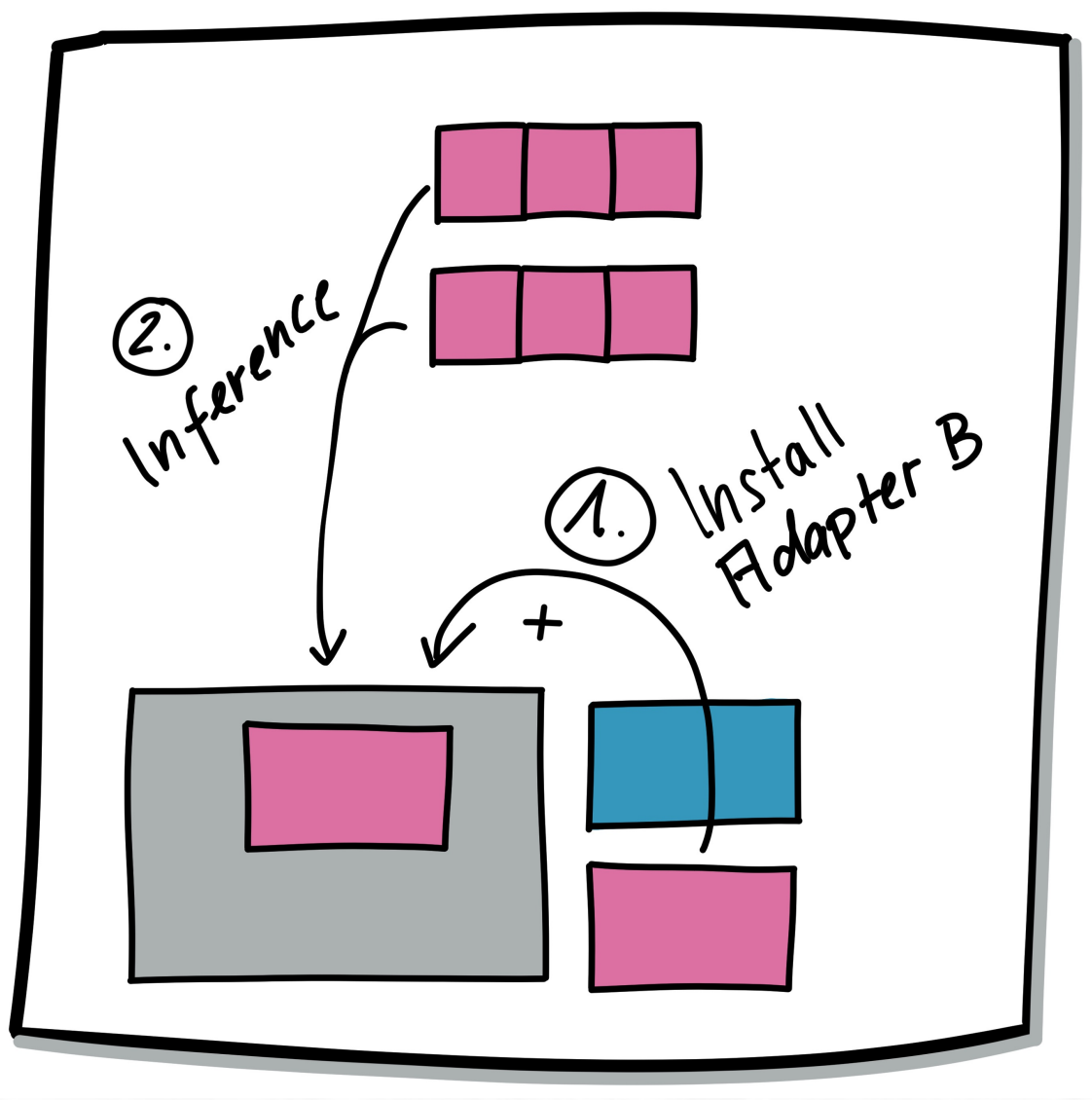
MULTI-TASK



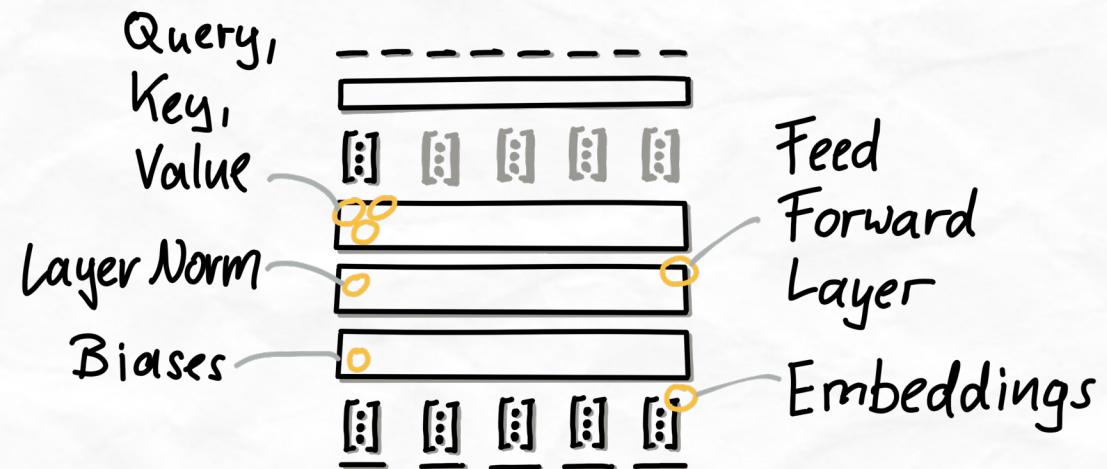




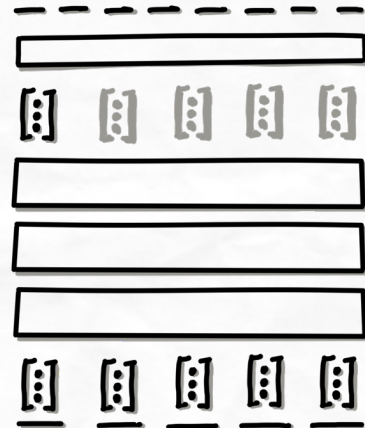


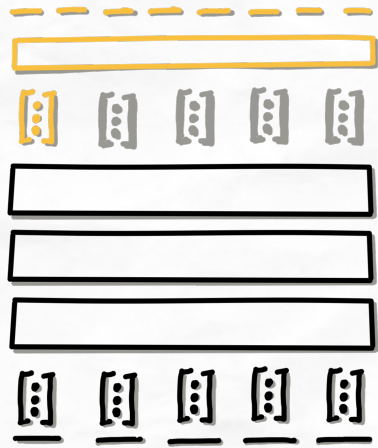


PEFT Approaches

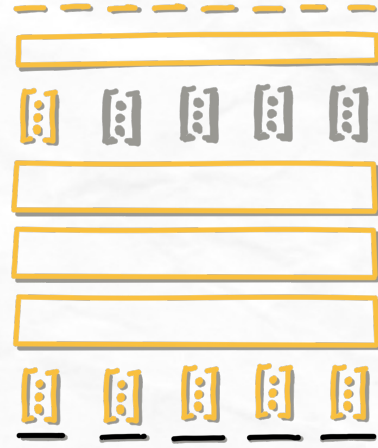


PEFT Approaches



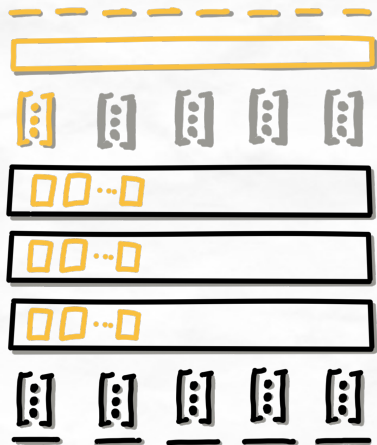


Classifier Only

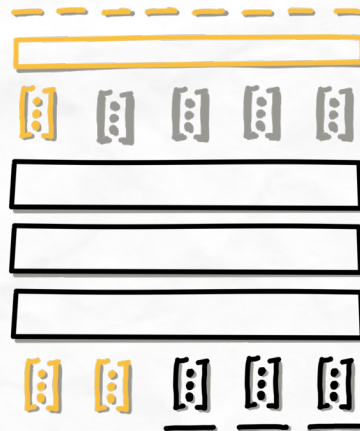


Full Finetuning

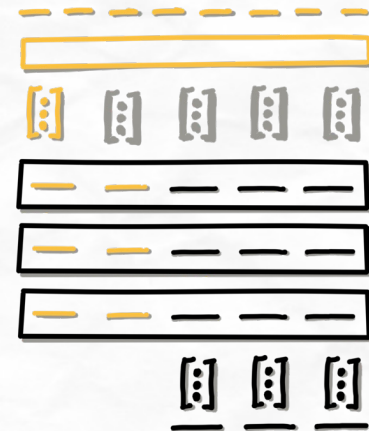
PEFT
Approaches



LoRA Finetuning

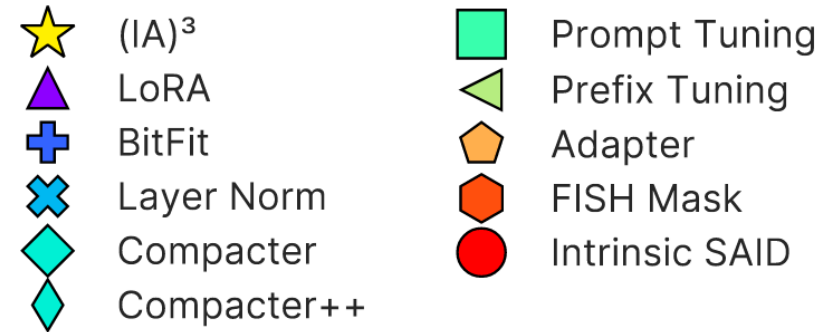
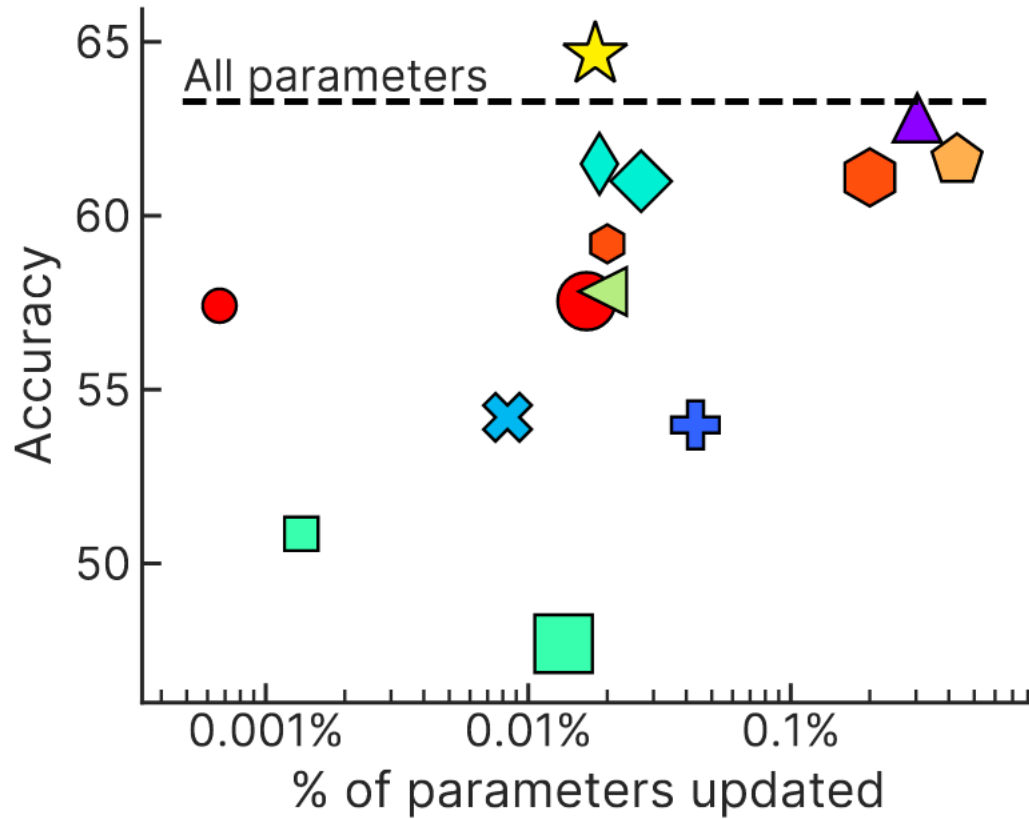


Soft Prompt Tuning



Prefix Tuning

PEFT Approaches



Liu, Tam, Muqeeth et al. 2022

Ressources



**Dive Into LoRA Adapters:
Intuitively Understanding
Finetuning Using LoRA**



**Fine-tune
LLaMA 2 (7-70B)
on Amazon SageMaker**



**Hugging Face
PEFT documentation**



Thank you!

Mariano Kamp

 <https://www.linkedin.com/in/mariano-kamp/>



Please complete the session survey.

