**Self-Supervised Fine-Tuning**

Input:

CS165

is

an

awesome

class

Decode:

Label:

is

an

awesome

class

EOS

**loss function**

pre-training

**Likelihood**

$$p(\text{is}, \text{an}, \ldots, \text{EOS})$$
$$= p(x_1, x_2, \ldots, x_N)$$
$$= \prod_{i=1}^{N} p(x_i \mid x_1, \ldots, x_{i-1})$$

**Cross-Entropy**

$$\mathscr{L} = -\frac{1}{N} \sum_{i=1}^{N} \log p(x_i)$$

**Perplexity**

$$e^{\mathscr{L}} = \exp\left( -\frac{1}{N} \sum_{i=1}^{N} \log p(x_i) \right)$$

DASlab
@ Harvard SEAS

Input:

Decode:    Label:

CS165 ➡️

is ➡️

an ➡️

awesome ➡️

class ➡️

📈 VS is

📈 VS an

📈 VS awesome

📈 VS class

📈 VS EOS

**loss function**

post-training

# Instruction Tuning

Input:

Decode: | Label:

prompt:
| CS165 |
| is |

answer:
| an |
| awesome |
| class |

Decode:
| 📈 | √§ |
| 📈 | √§ |
| 📈 | √§ |
| 📈 | √§ |
| 📈 | √§ |

Label:
| is |
| an |
| awesome |
| class |
| EOS |

**loss function**

post-training

# Instruction Tuning

Input:

prompt:

| CS165 |
| is |

answer:

| an |
| awesome |
| class |

Decode:

Label:

| null |
| null |
| awesome |
| class |
| EOS |

**loss function**

post-training

DASlab
@ Harvard SEAS

# Instruction Tuning

## Evaluation

**Perplexity**

Input:

CS165

is

an → 📈 √$ awesome

awesome → 📈 √$ class

class → 📈 √$ EOS

Decode:

Step 1:

loss function

**Accuracy**

Input:

CS165

is → an → awesome → class → EOS

Step 1: Step 2: Step 3:

**Generated Answer** class

**Golden Answer** class

post-training

DASlab
@ Harvard SEAS

**Trajectory**

**User:** Find the population of Tokyo and compare it to New York City.

**Thought:** I should search for the current population of Tokyo first.
**Action:** search("population of Tokyo 2024")
**Observation:** Tokyo metropolitan area population is approximately 37.4 million.

**Thought:** Now I need the population of New York City.
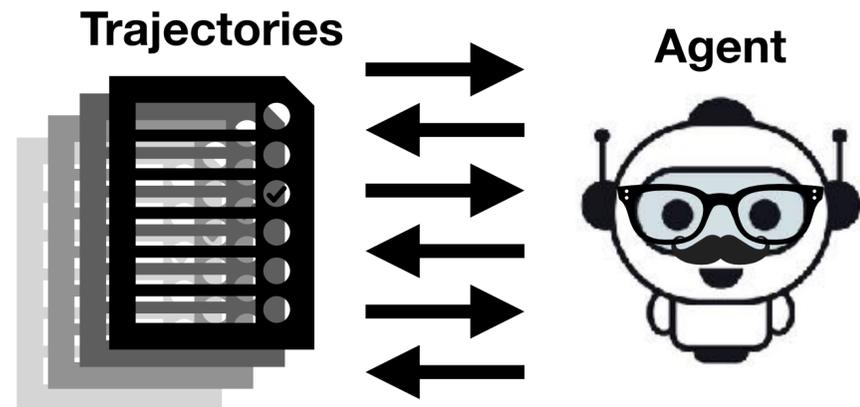**Action:** search("population of New York City 2024")
**Observation:** New York City population is approximately 8.3 million.

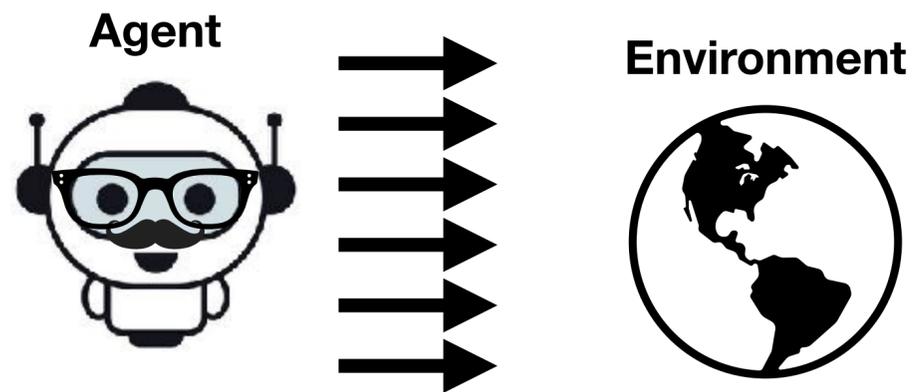**Thought:** I have both numbers. I can now compute the comparison.
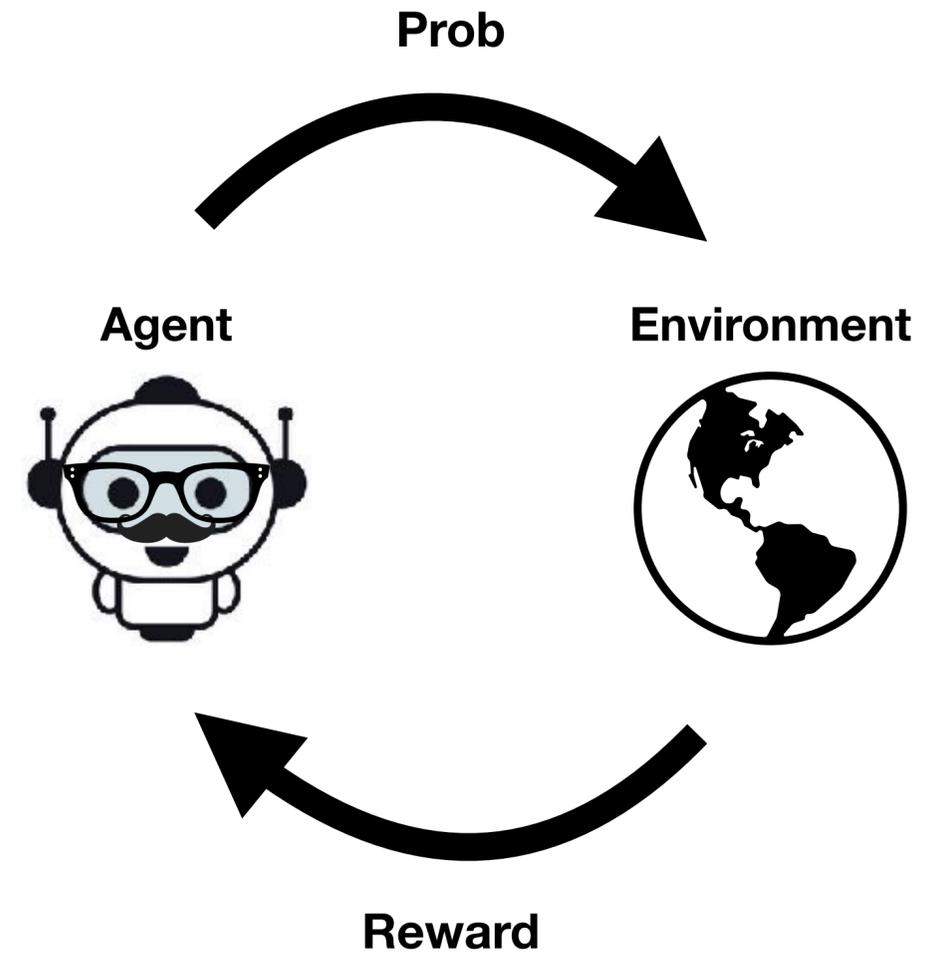**Final Answer:** Tokyo (37.4M) is approximately 4.5 times larger than New York City (8.3M).

# PROBLEM
# Agent fine-tuning rests on benchmarks the quality of which has not been systematically assessed

*If the scores cannot be trusted — the conclusions built on top of them cannot be trusted either*

# 1. Quality Assessment of Agent Fine-Tuning Benchmarks

## SOLUTION
# Build a unified quality framework to characterize what agent benchmarks actually measure

*Propose new benchmarks to fill in the gaps*

# Benchmark Difficulty Curve

Each bar is a bucket of tasks. Height = number of tasks in that bucket. X-axis = empirical difficulty (fraction of model-run pairs that fail the task). Shape reveals whether a benchmark discriminates capability.

## Smooth unimodal spread  `GOOD`



- **Tasks span the full range.** Every difficulty tier is represented — easy, medium, hard.
- **Score differences are meaningful.** A better model shifts the boundary rightward; you can detect that.
- **Good for both SFT and RL.** Mid-range tasks provide demonstrations and reward signal.
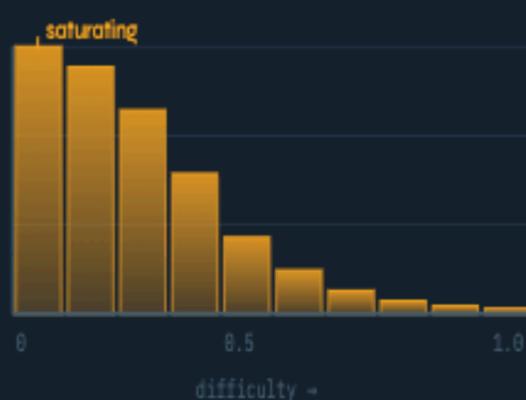
## Bimodal / U-shaped  `BAD`



- **Measures one threshold, not a range.** Tasks are trivial or impossible — nothing in between.
- **Training improvements are invisible.** A better model still fails the hard tasks and already solved the easy ones.
- **RL gets no gradient signal.** Reward variance near zero in both tails — nothing to optimize against.
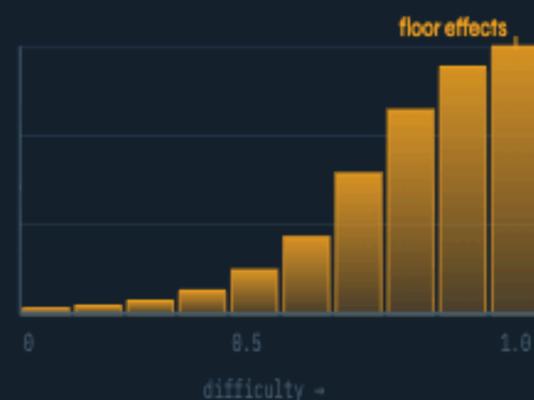
## Left-skewed — too easy  `BAD`



- **Scores cluster near the ceiling.** Most models solve most tasks — no room to improve. AgentGym shows this pattern.
- **Cannot distinguish strong models.** A fine-tuned model looks roughly the same as the baseline.
- **SFT overfits to easy patterns.** Training on these tasks doesn't generalize; it just memorizes solved trajectories.

## Right-skewed — too hard  `BAD`



- **Scores cluster near zero.** The benchmark was designed for a capability level the field hasn't reached yet.
- **RL starves.** Successes are too rare for stable reward signal — gradients are noisy and training diverges.
- **SFT has no demonstrations.** If the teacher model also fails, there are no successful trajectories to imitate.
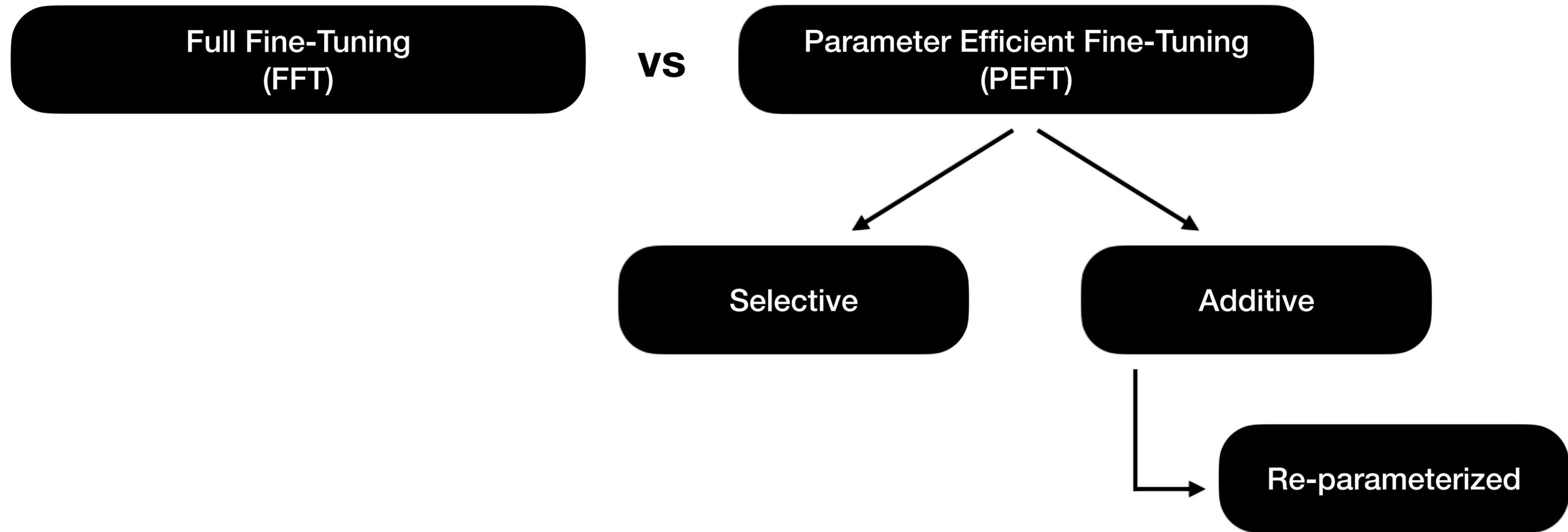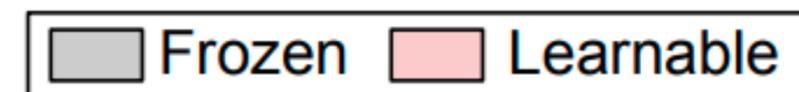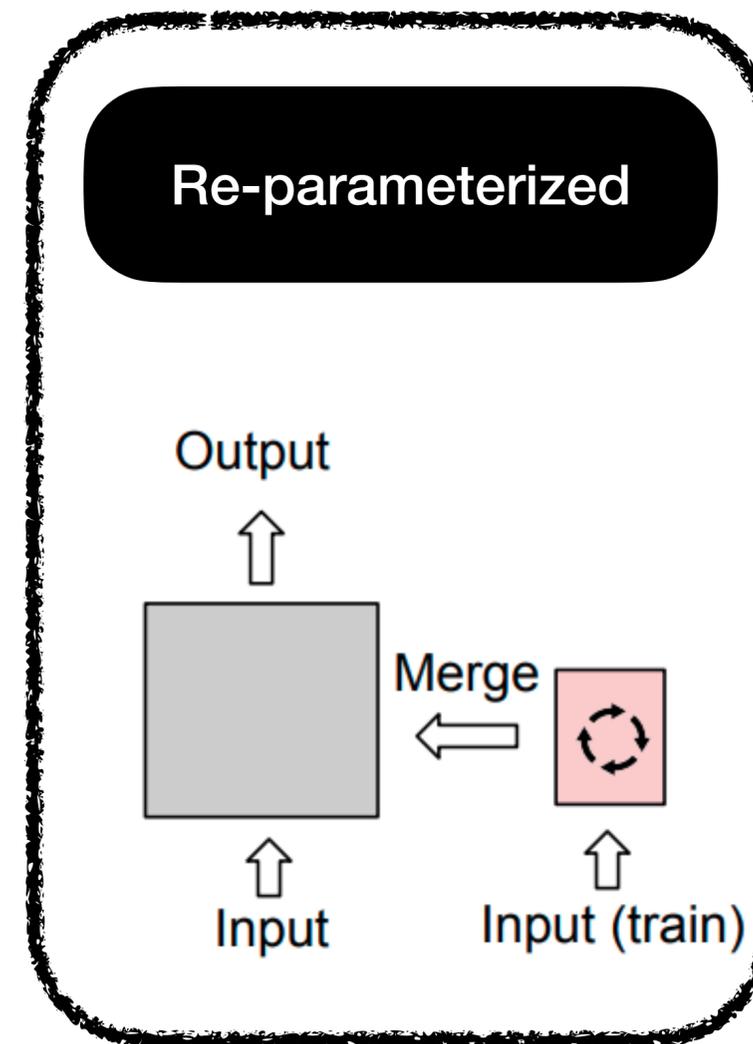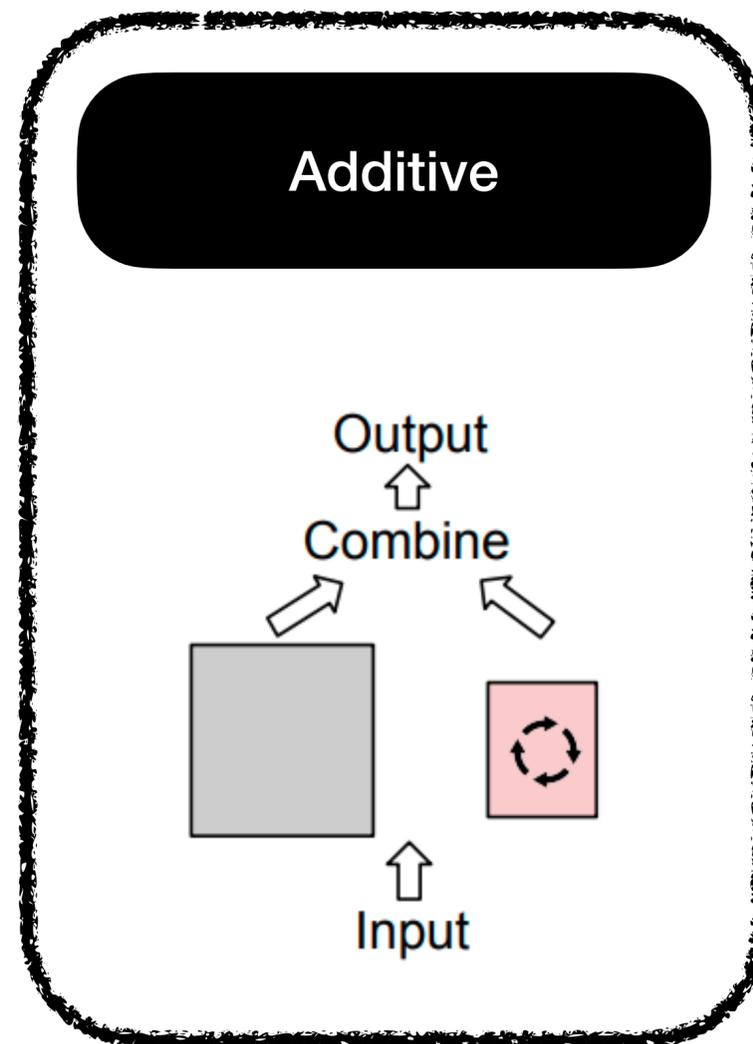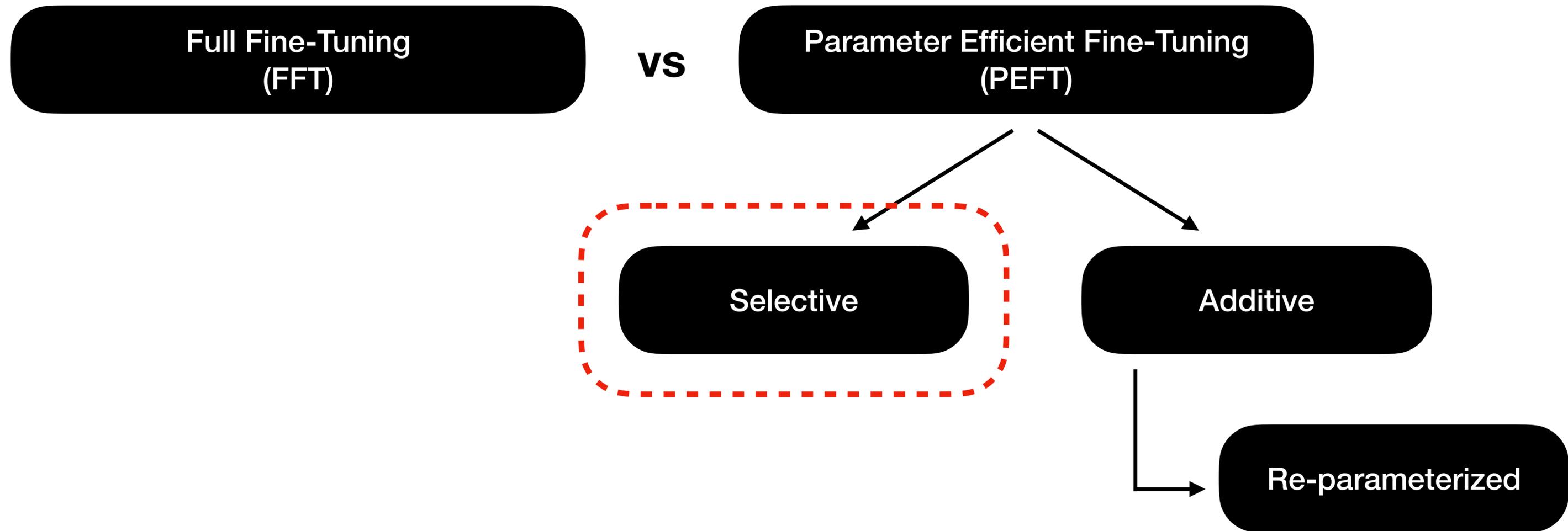
# Weight Update

Full Fine-Tuning (FFT)  **vs**  Parameter Efficient Fine-Tuning (PEFT)

Selective · Additive · Re-parameterized

Frozen · Learnable

# PEFT Taxonomy



**Full Fine-Tuning (FFT)**

vs

**Parameter Efficient Fine-Tuning (PEFT)**

Selective

Additive

Re-parameterized

DASlab
@ Harvard SEAS

# LLM Block

# Selective PEFT

frozen                                                                trainable

# Selective PEFT

# Selective PEFT
## Two more ideas

$$\mathbf{Q}^{m,\ell}(\mathbf{x}) = \mathbf{W}_q^{m,\ell}\mathbf{x} + \mathbf{b}_q^{m,\ell}$$

$$\mathbf{K}^{m,\ell}(\mathbf{x}) = \mathbf{W}_k^{m,\ell}\mathbf{x} + \mathbf{b}_k^{m,\ell}$$

$$\mathbf{V}^{m,\ell}(\mathbf{x}) = \mathbf{W}_v^{m,\ell}\mathbf{x} + \mathbf{b}_v^{m,\ell}$$

$$\mathbf{h}_2^\ell = \text{Dropout}\big(\mathbf{W}_{m_1}^\ell \cdot \mathbf{h}_1^\ell + \mathbf{b}_{m_1}^\ell\big)$$

$$\mathbf{h}_3^\ell = \mathbf{g}_{LN_1}^\ell \odot \frac{(\mathbf{h}_2^\ell + \mathbf{x}) - \mu}{\sigma} + \mathbf{b}_{LN_1}^\ell$$

$$\mathbf{h}_4^\ell = \text{GELU}\big(\mathbf{W}_{m_2}^\ell \cdot \mathbf{h}_3^\ell + \mathbf{b}_{m_2}^\ell\big)$$

$$\mathbf{h}_5^\ell = \text{Dropout}\big(\mathbf{W}_{m_3}^\ell \cdot \mathbf{h}_4^\ell + \mathbf{b}_{m_3}^\ell\big)$$

$$\text{out}^\ell = \mathbf{g}_{LN_2}^\ell \odot \frac{(\mathbf{h}_5^\ell + \mathbf{h}_3^\ell) - \mu}{\sigma} + \mathbf{b}_{LN_2}^\ell$$

LayerNorm Tuning
(Bingchen Zhao et al., 2024, ICLR)

LLM Block

Add & Norm

FFN

Add & Norm

(Self/Cross) Attention

Q    K    V

x

DASlab
@ Harvard SEAS

# PEFT Taxonomy

# Additive PEFT
## The Adapter Architecture

# Additive PEFT
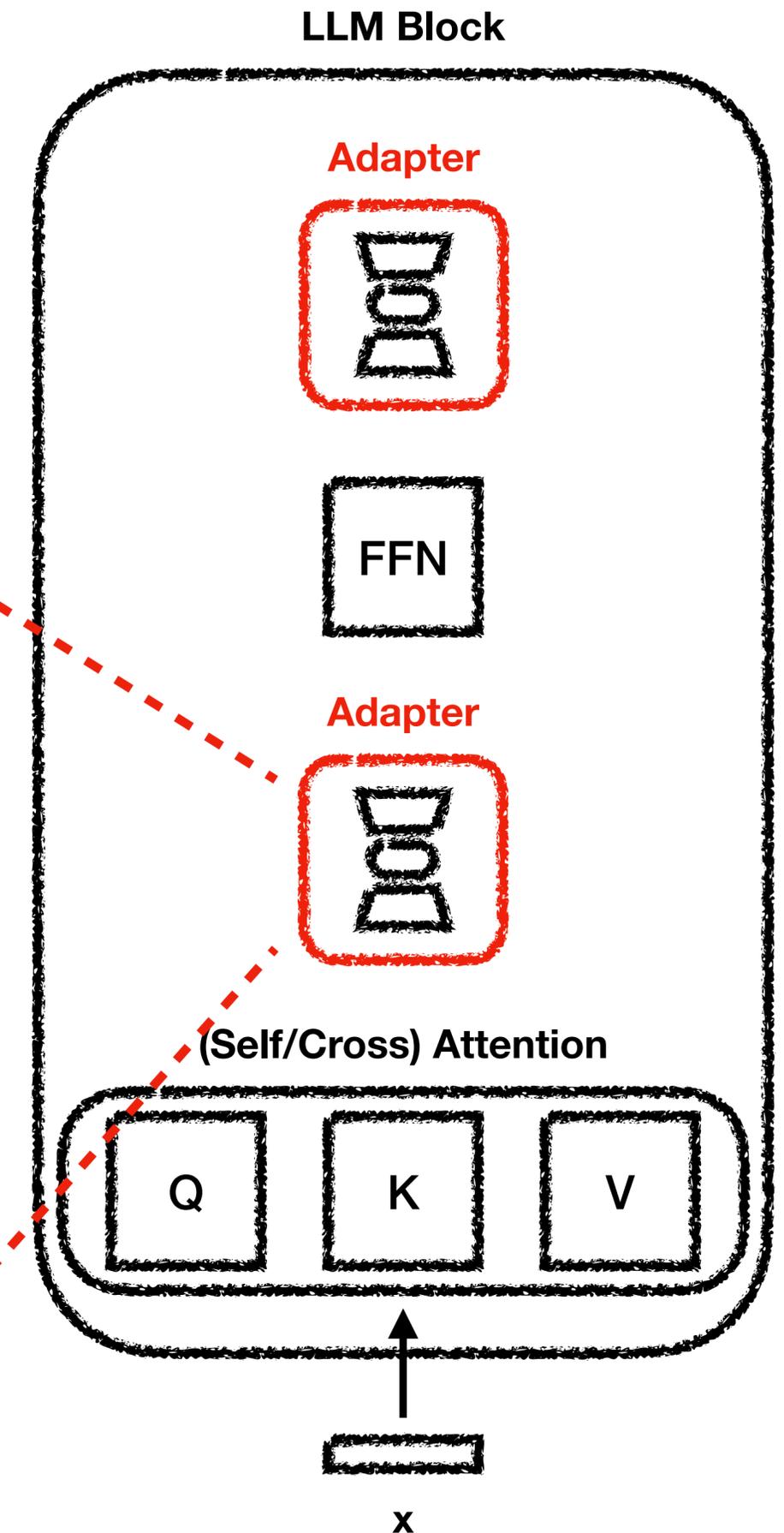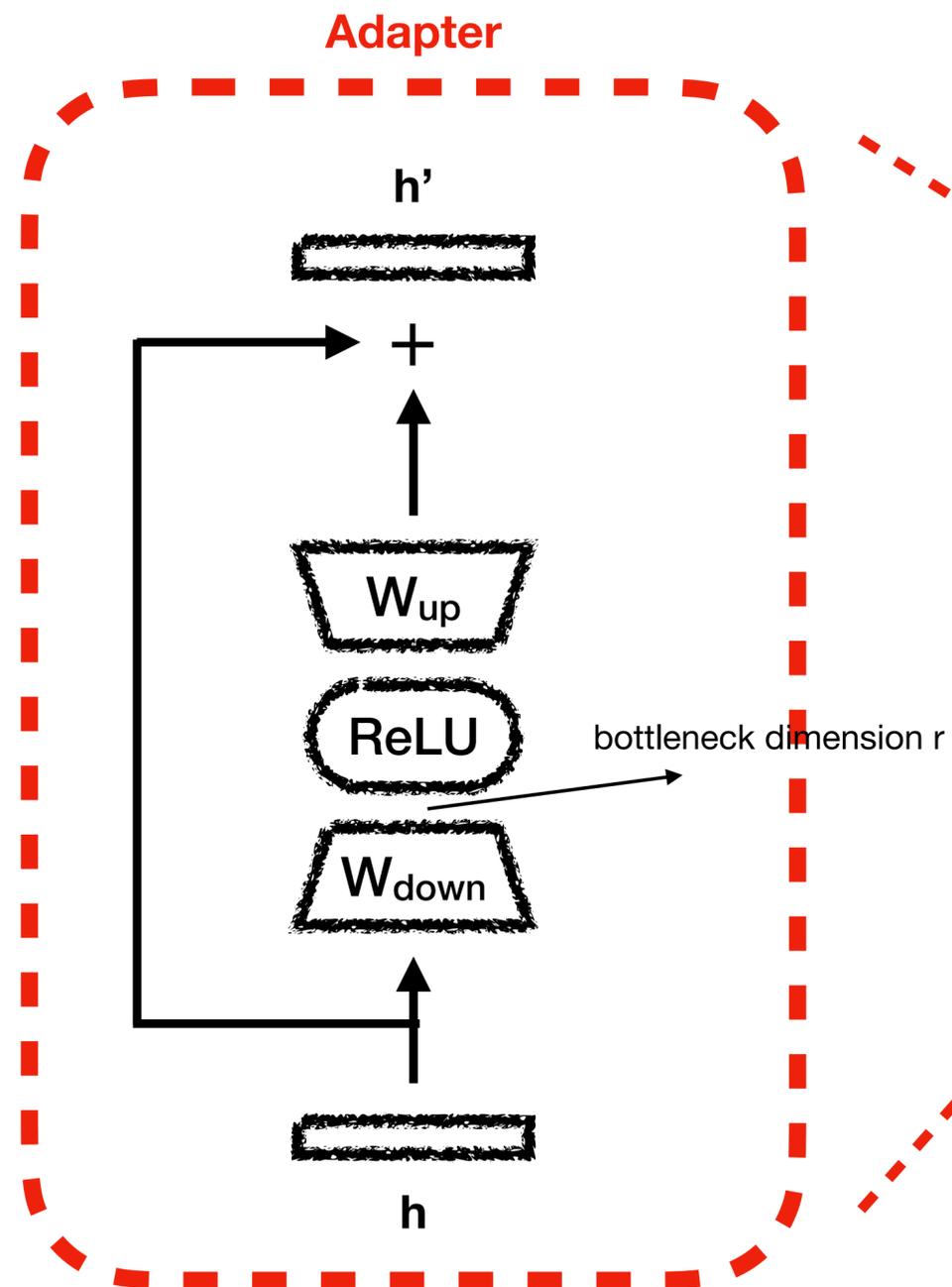## The Adapter Architecture

# Additive PEFT
## [Serial] Adapter (Houlsby et al., 2019)

$$h' \leftarrow h + f(hW_{down})W_{up}$$

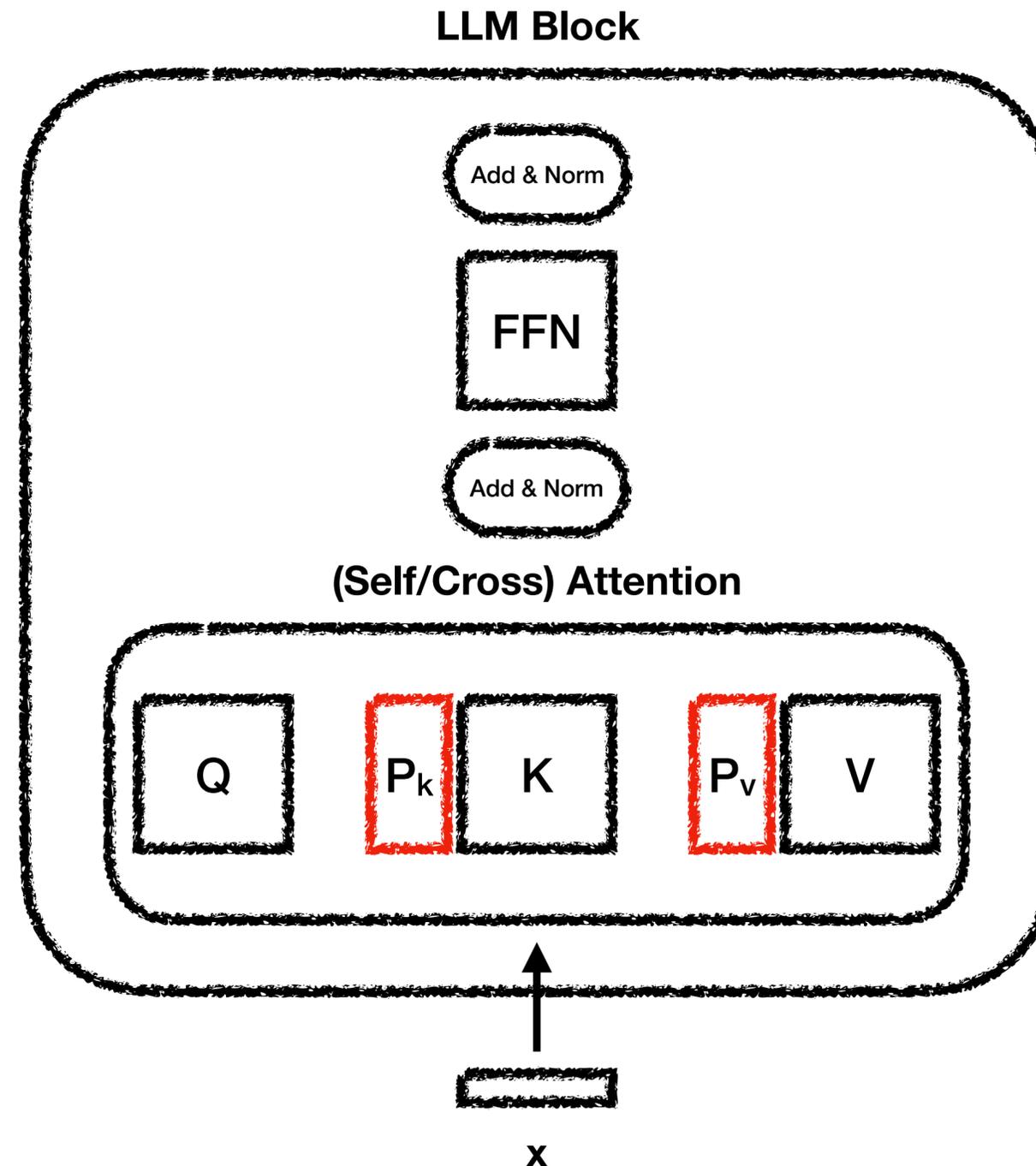$$W_{down} \in \mathbb{R}^{d \times r}$$

$$W_{up} \in \mathbb{R}^{r \times d}$$

Adapters: small trainable feed-forward networks inserted between the layers in the frozen pre-trained model.
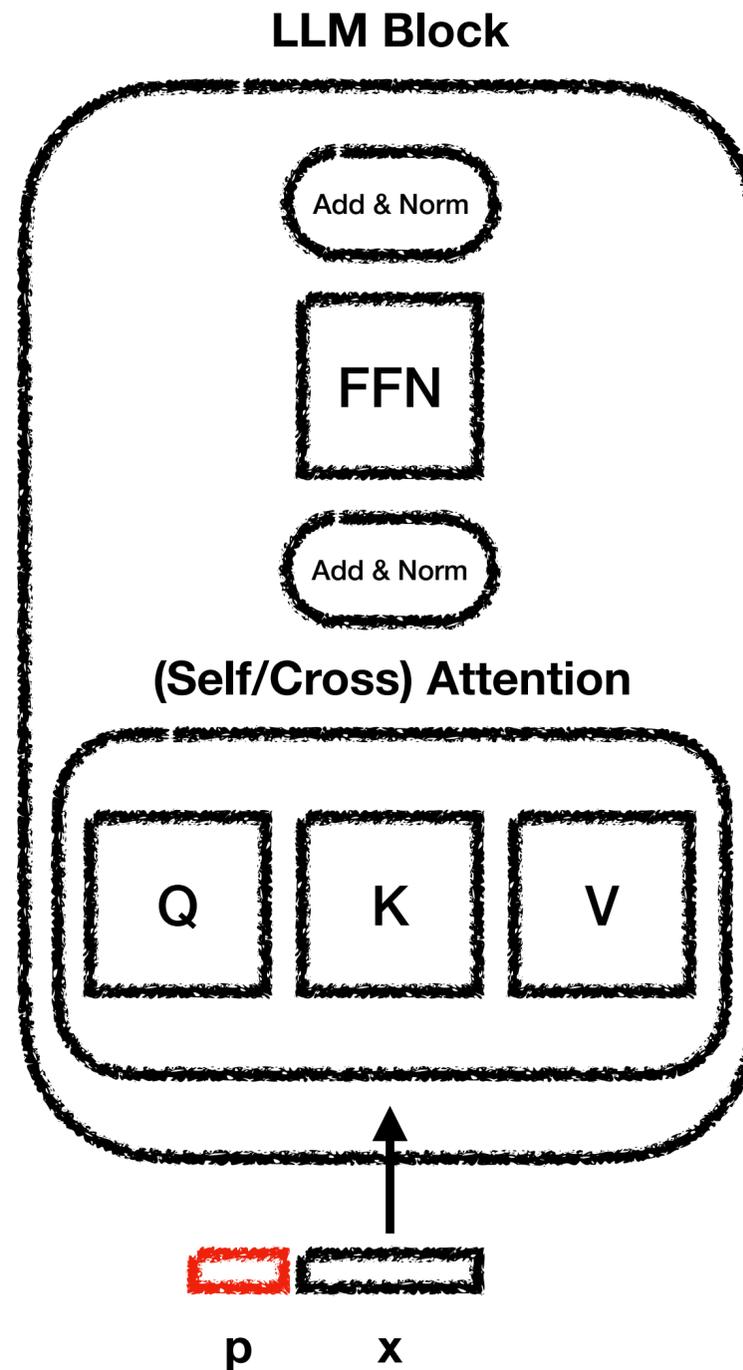
**Adapter**

h'

+

$W_{up}$

ReLU → bottleneck dimension r

$W_{down}$

h

**LLM Block**

**Adapter**

**Adapter**

FFN

**(Self/Cross) Attention**

Q    K    V

x

DASlab
@ Harvard SEAS

# Additive PEFT
## Prefix Tuning (Li & Liang, 2021)

# Additive PEFT
## Prompt Tuning (Brian Lester et al., 2021)



LLM Block

Add & Norm

FFN

Add & Norm

(Self/Cross) Attention

Q    K    V

p    x

# PEFT Taxonomy

Full Fine-Tuning
(FFT)

**vs**

Parameter Efficient Fine-Tuning
(PEFT)

Selective

Additive

Re-parameterized

# Which one is better from a systems perspective ❓

# PEFT Taxonomy

Full Fine-Tuning
(FFT)

**vs**

Parameter Efficient Fine-Tuning
(PEFT)

Selective

Additive

Re-parameterized

# Re-parameterized PEFT

# Re-parameterized PEFT
## Low Rank Adaptation (Hu et al., 2021)

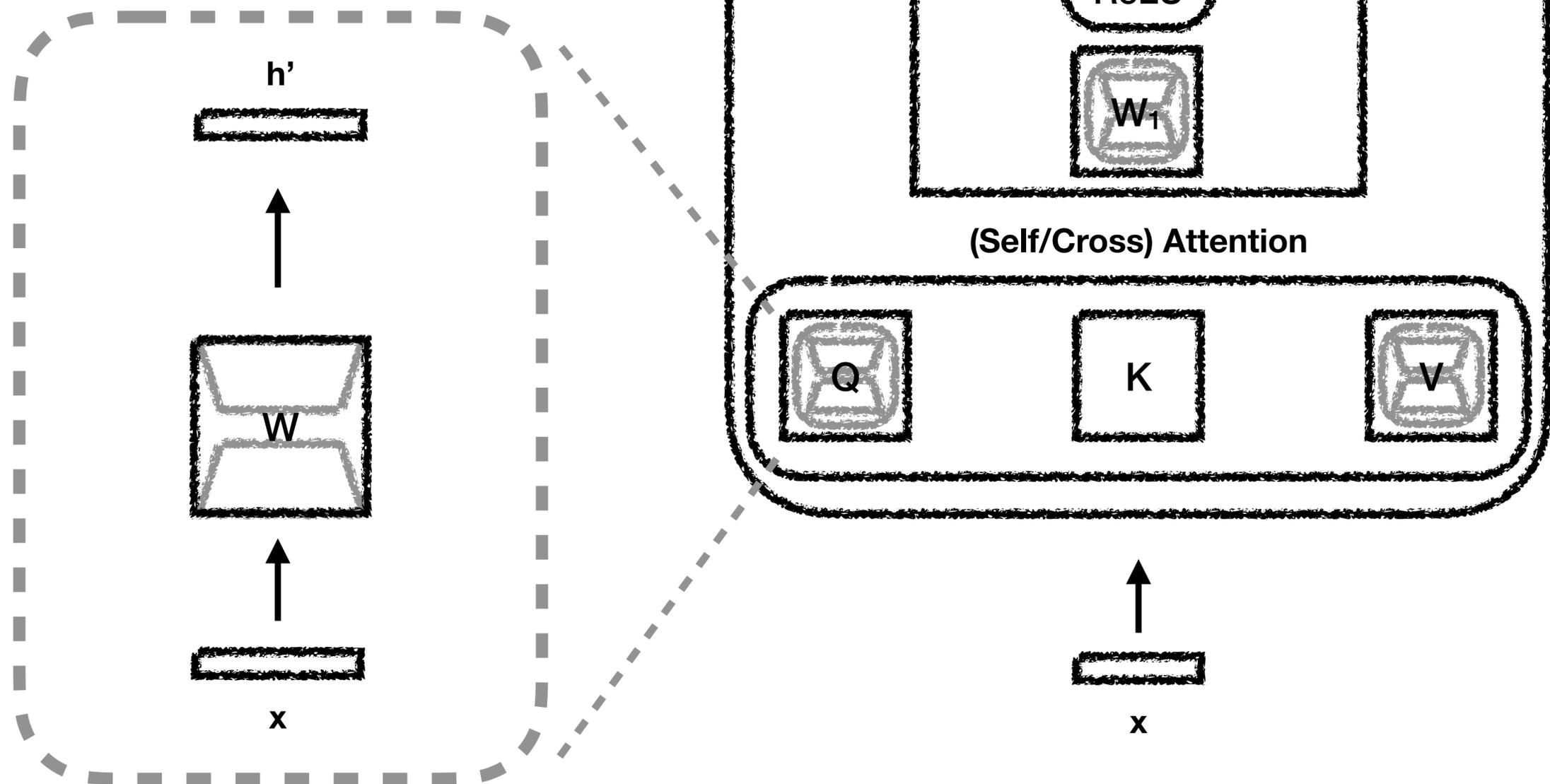$$h \leftarrow xW + \frac{a}{r} \cdot xW_{down}W_{up}$$

$$x \in \mathbb{R}^{1 \times d}$$
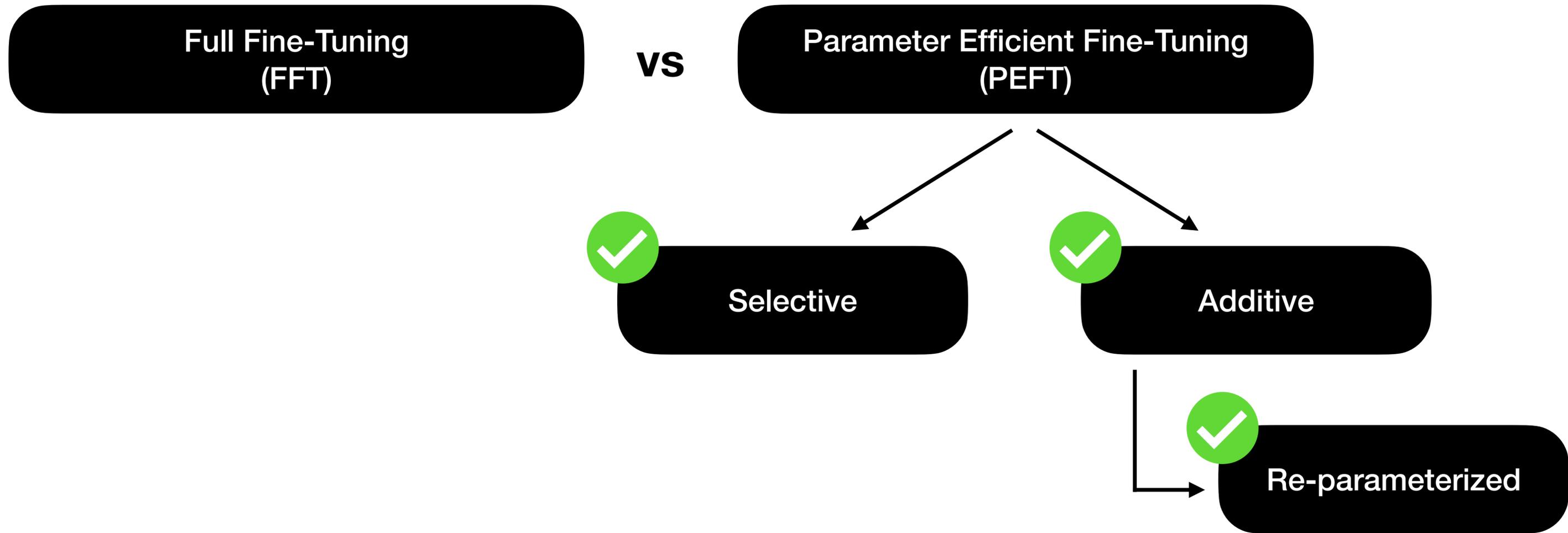
$$h, h' \in \mathbb{R}^{1 \times k}$$

$$W \in \mathbb{R}^{d \times k}$$
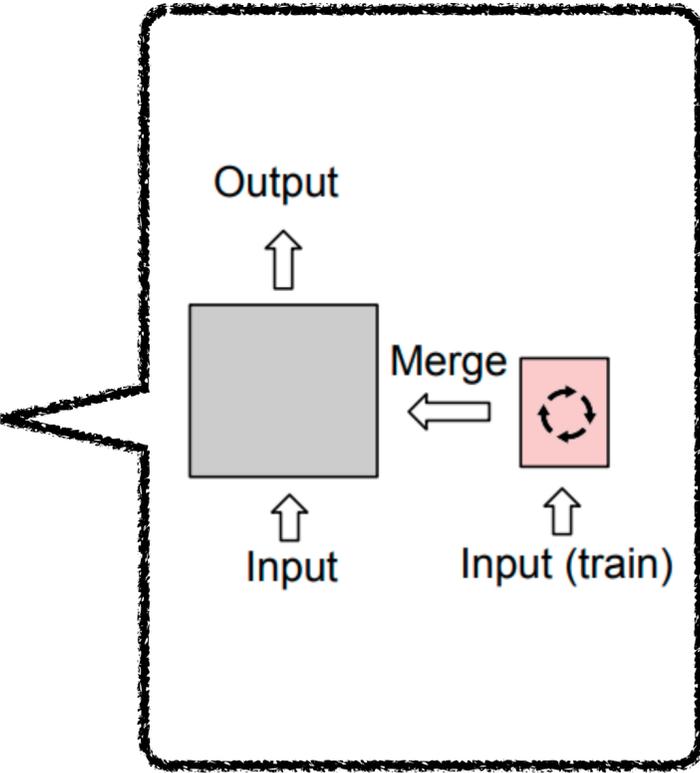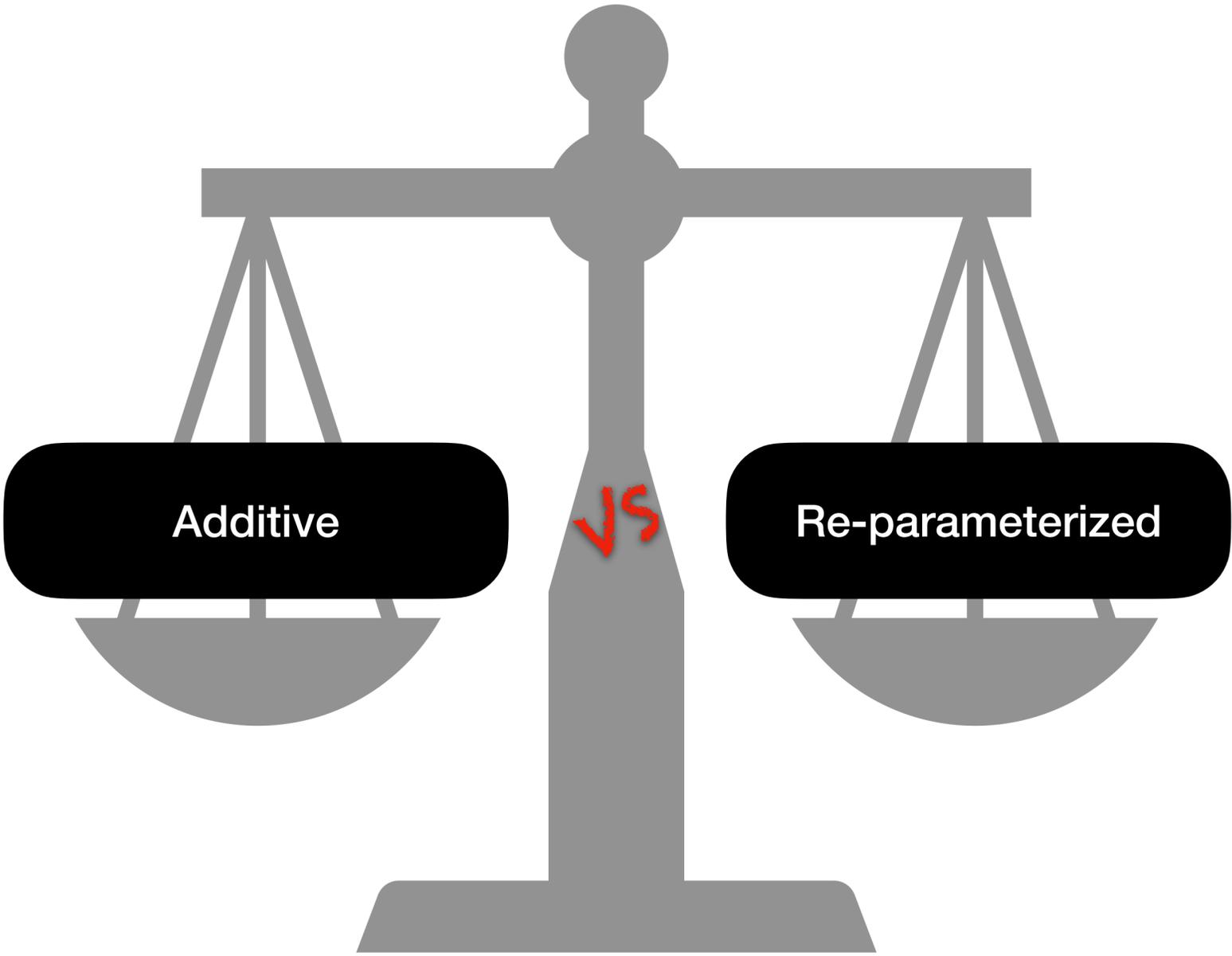
$$W_{down} \in \mathbb{R}^{d \times r}$$

$$W_{up} \in \mathbb{R}^{r \times k}$$

# Re-parameterized PEFT
## Low Rank Adaptation (Hu et al., 2021)

$$h \leftarrow xW + \frac{a}{r} \cdot xW_{down}W_{up}$$

$x \in \mathbb{R}^{1 \times d}$

$h, h' \in \mathbb{R}^{1 \times k}$

$W \in \mathbb{R}^{d \times k}$

$W_{down} \in \mathbb{R}^{d \times r}$

$W_{up} \in \mathbb{R}^{r \times k}$



LoRA

h'

W

W_up

W_down

bottleneck dimension r

x

LLM Block

FFN

LoRA

W₂

ReLU

LoRA

W₁

(Self/Cross) Attention

Q

K

V

LoRA

W_q

LoRA

W_k

LoRA

W_v

x

DASlab
@ Harvard SEAS

# Re-parameterized PEFT
## Low Rank Adaptation (Hu et al., 2021)

$$h \leftarrow xW + \frac{a}{r} \cdot xW_{down}W_{up}$$

$$x \in \mathbb{R}^{1 \times d}$$

$$h, h' \in \mathbb{R}^{1 \times k}$$

$$W \in \mathbb{R}^{d \times k}$$

$$W_{down} \in \mathbb{R}^{d \times r}$$

$$W_{up} \in \mathbb{R}^{r \times k}$$

# Re-parameterized PEFT
## Low Rank Adaptation (Hu et al., 2021)

$$h \leftarrow xW + \frac{a}{r} \cdot xW_{down}W_{up}$$

$$x \in \mathbb{R}^{1 \times d}$$

$$h, h' \in \mathbb{R}^{1 \times k}$$

$$W \in \mathbb{R}^{d \times k}$$

$$W_{down} \in \mathbb{R}^{d \times r}$$

$$W_{up} \in \mathbb{R}^{r \times k}$$



**LLM Block**

**FFN**

$W_2$

ReLU

$W_1$

**(Self/Cross) Attention**

Q  K  V

h'

W

x

x

# PEFT Taxonomy

Full Fine-Tuning
(FFT)

**vs**

Parameter Efficient Fine-Tuning
(PEFT)

✅ Selective

✅ Additive

✅ Re-parameterized

# Which one is better from a systems perspective ❓



Additive **VS** Re-parameterized

# PEFT Taxonomy



**LoRA Variants**

Fig. 3: Taxonomy of Parameter-Efficient Fine-Tuning Methods for Large Models.

(2024, arXiv) Parameter-Efficient Fine-Tuning (PEFT) for Large Models: A Comprehensive Survey

**LoRA Variants**



Fig. 1. Hierarchical taxonomy of LoRA variants based on four core principle operational axes.

## PROBLEM

Given a fixed parameter budget, how should LoRA adapters be allocated across the layers and modules of an LLM to maximize task performance?

*The design space is combinatorially large — methods distribute them uniformly or dynamically*

DASlab
@ Harvard SEAS

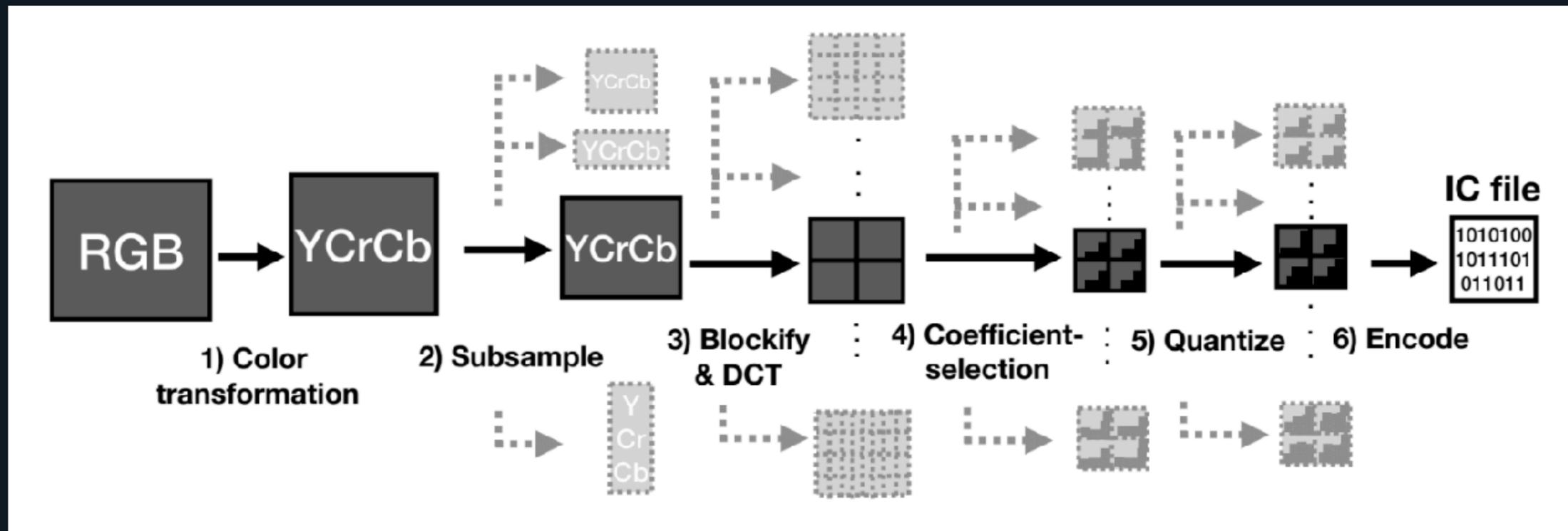# 2. Parameter Budget Allocation in Parameter Efficient Fine-Tuning

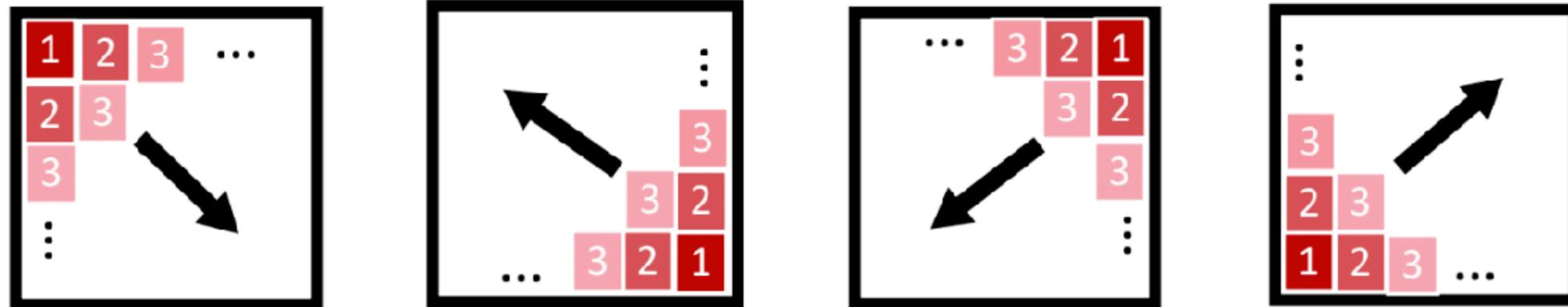# 2. Parameter Budget Allocation in Parameter Efficient Fine-Tuning

## SOLUTION
# Formally define the primitives of the LoRA placement design space and identify good domains to restrict the otherwise infinite design space

*Generalize across models and domains — find simple allocation heuristics*

DASlab
@ Harvard SEAS

# Image Calculator Throwback



(a) Strategy 1: Increasing-sized upper-left triangle (b) Strategy 2: Increasing-sized lower-right triangle (c) Strategy 3: Increasing-sized upper-right triangle (d) Strategy 4: Increasing-sized lower-left triangle