

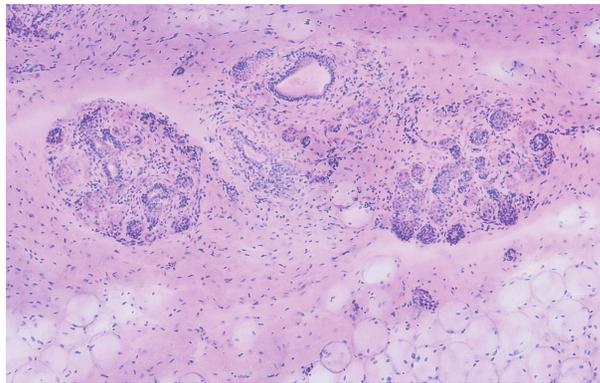
Scalable Image AI via Self-Designing Storage

*Utku
Sirin*  **DASlab**
@ Harvard SEAS

*work with Victoria Kauffman, Aadit Saluja, Florian Klein, Jeremy Hsu,
Vlad Cainamisir, Qitong Wang, Konstantinos Kopsinis, Stratos Idreos*

What if we can reason about image data?

Seeing is at center of AI



**Digital
Pathology**



**Textile
Recycling**

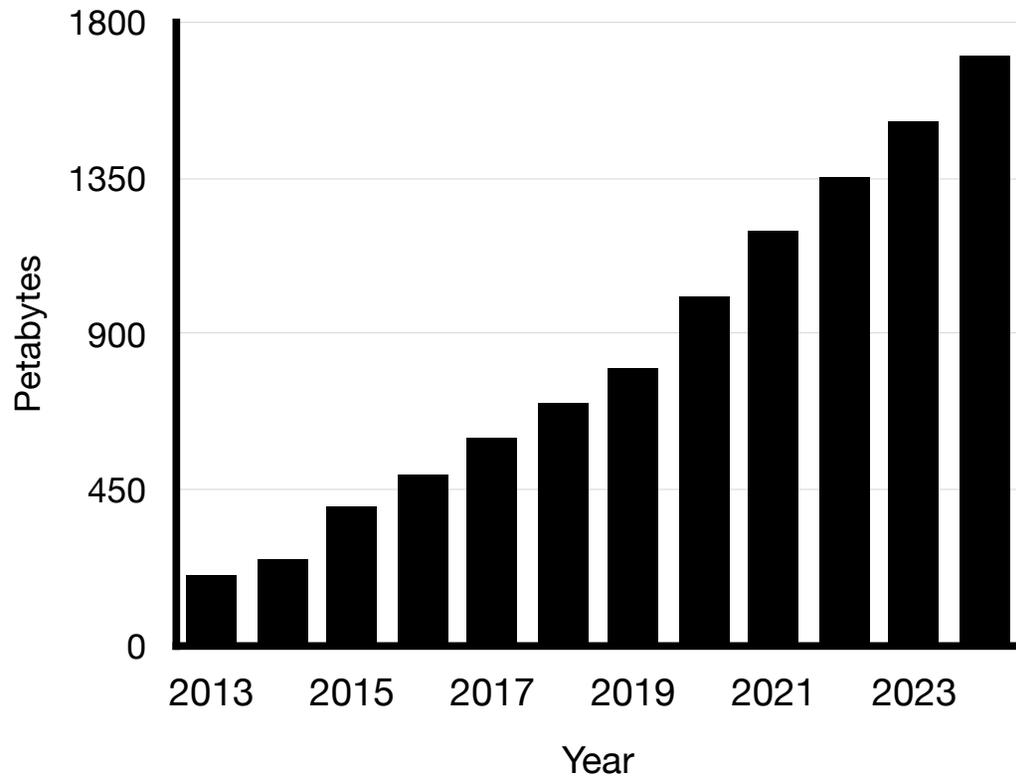


**Robot
Navigation**

Data size

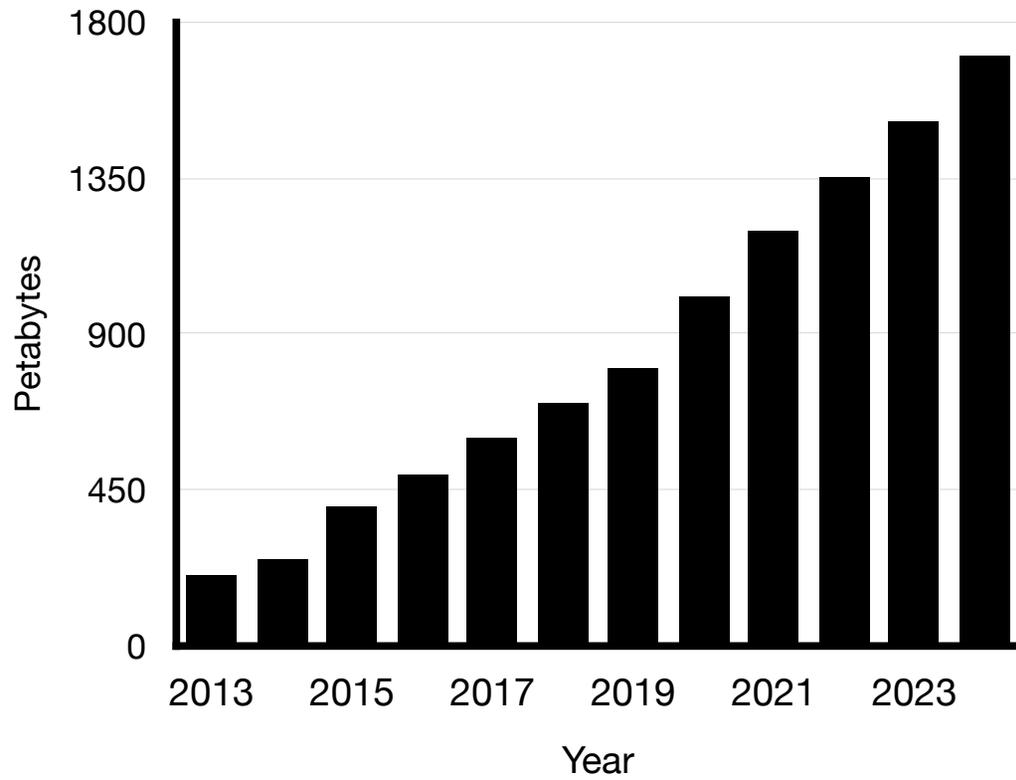
Model size

Data size

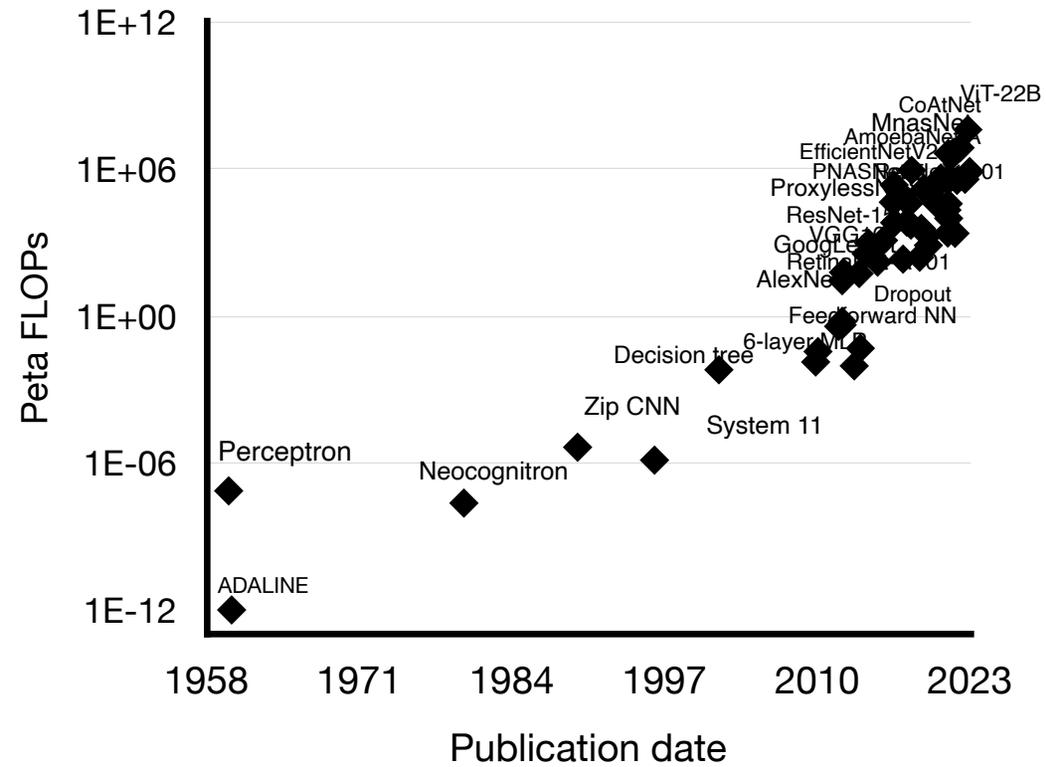


Model size

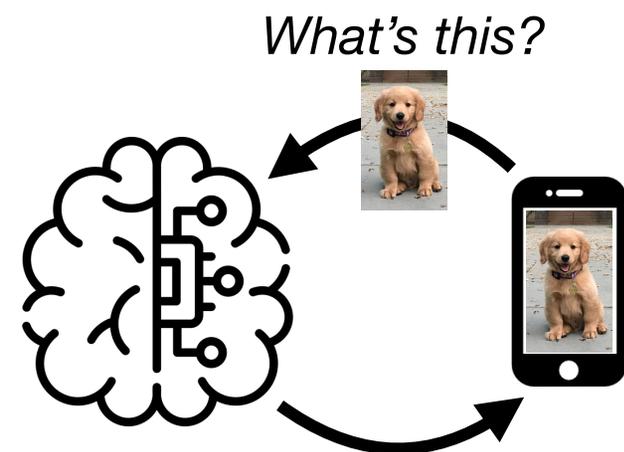
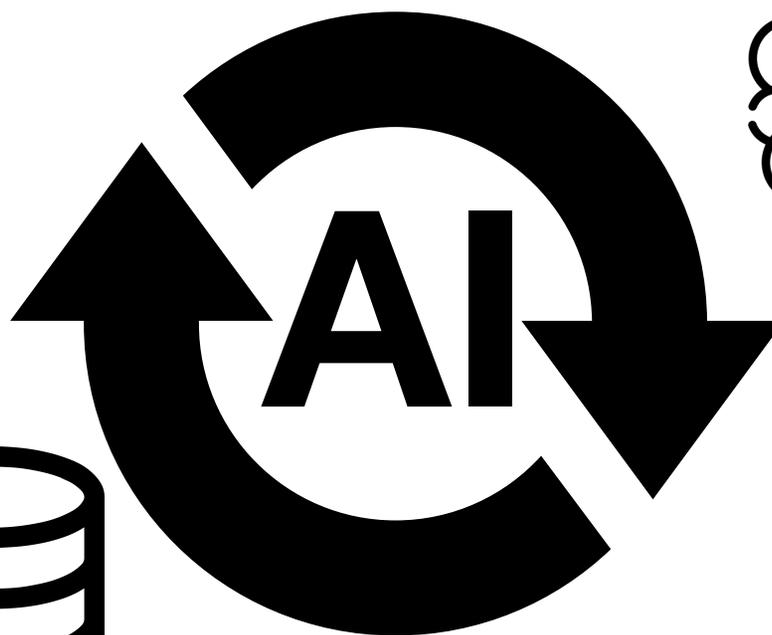
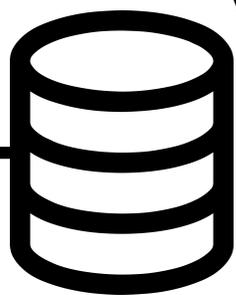
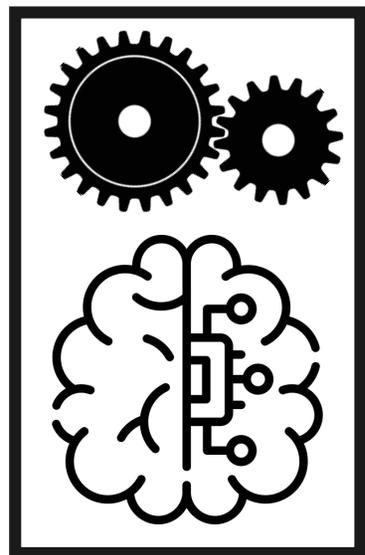
Data size



Model size

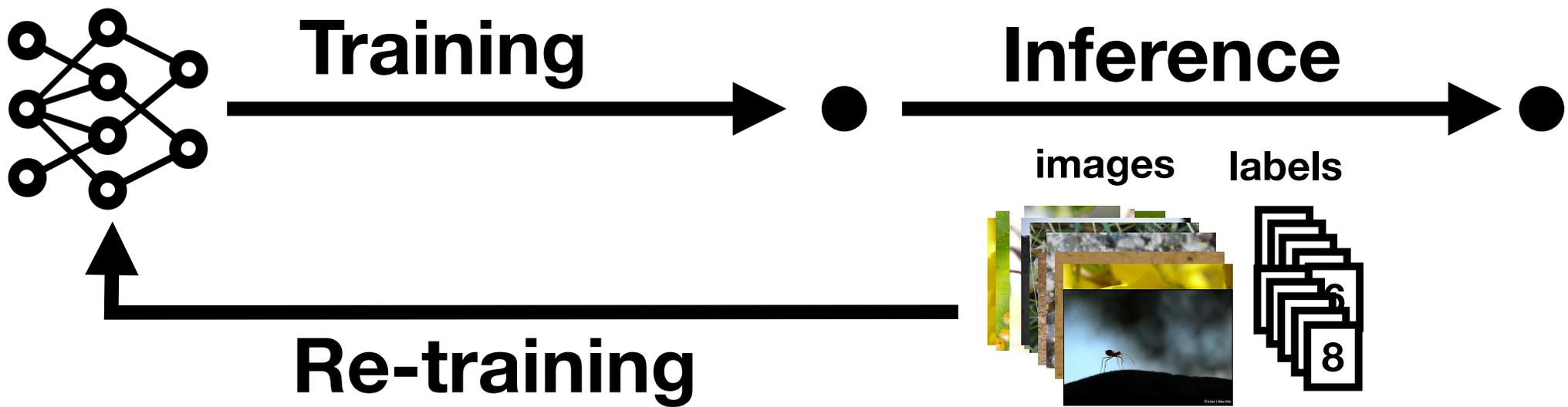


Training

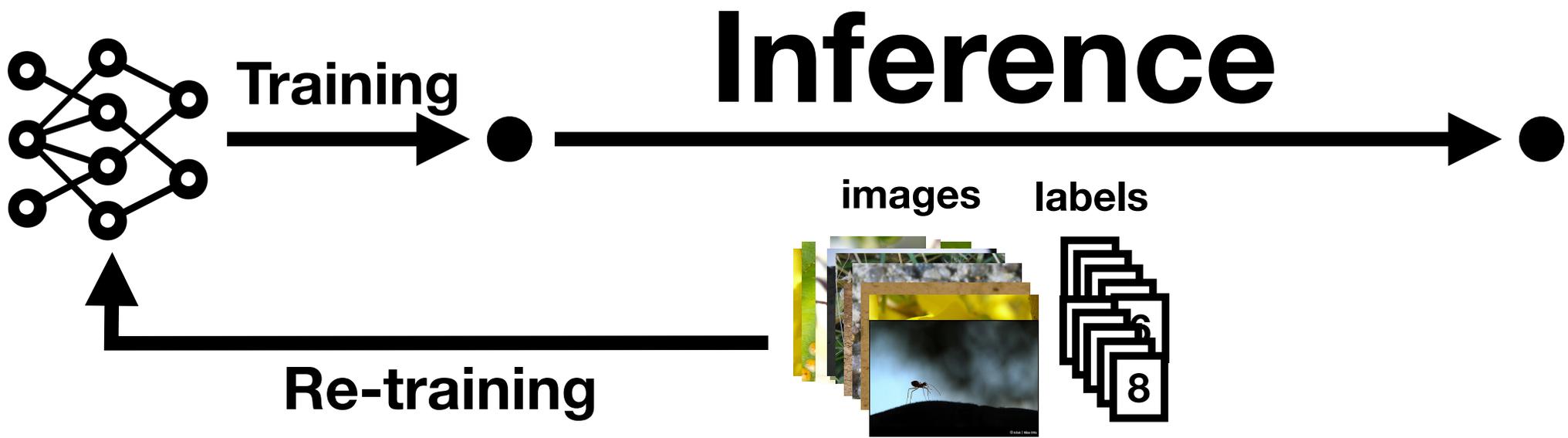


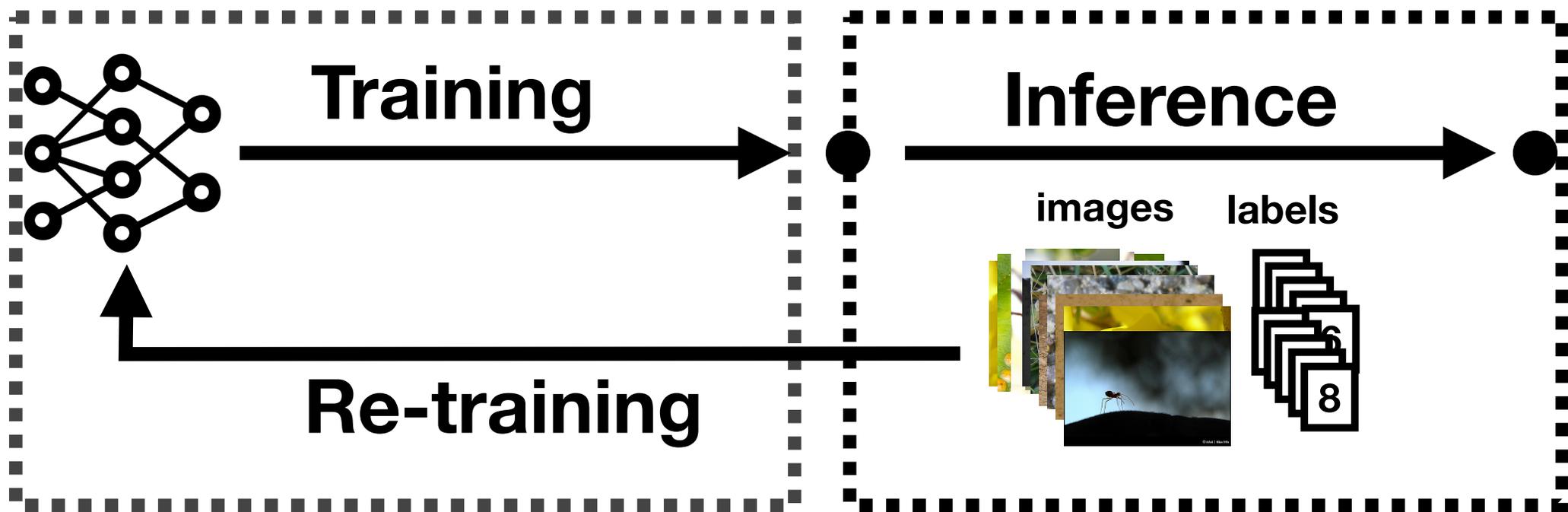
What's this?
It's a dog!

Inference

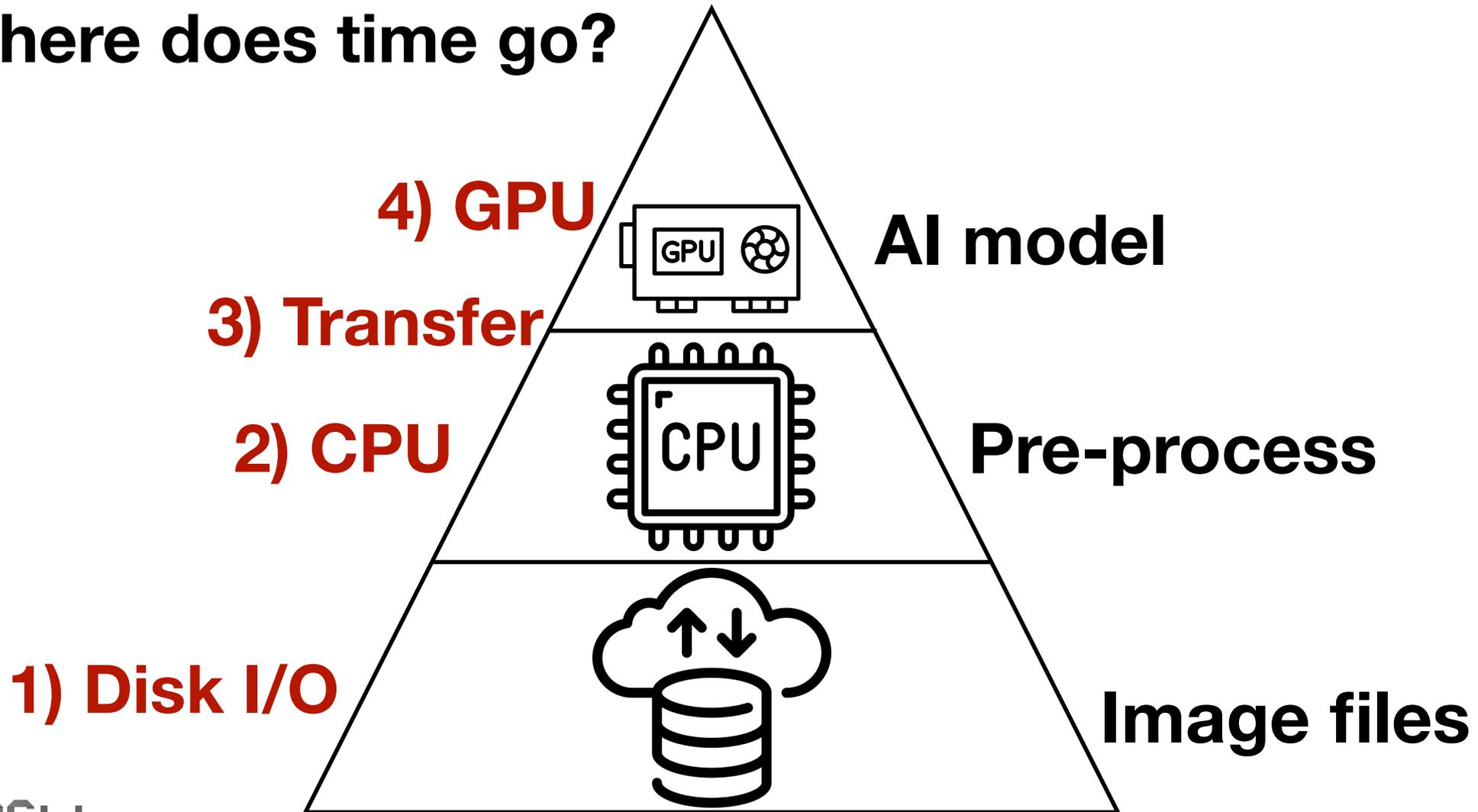






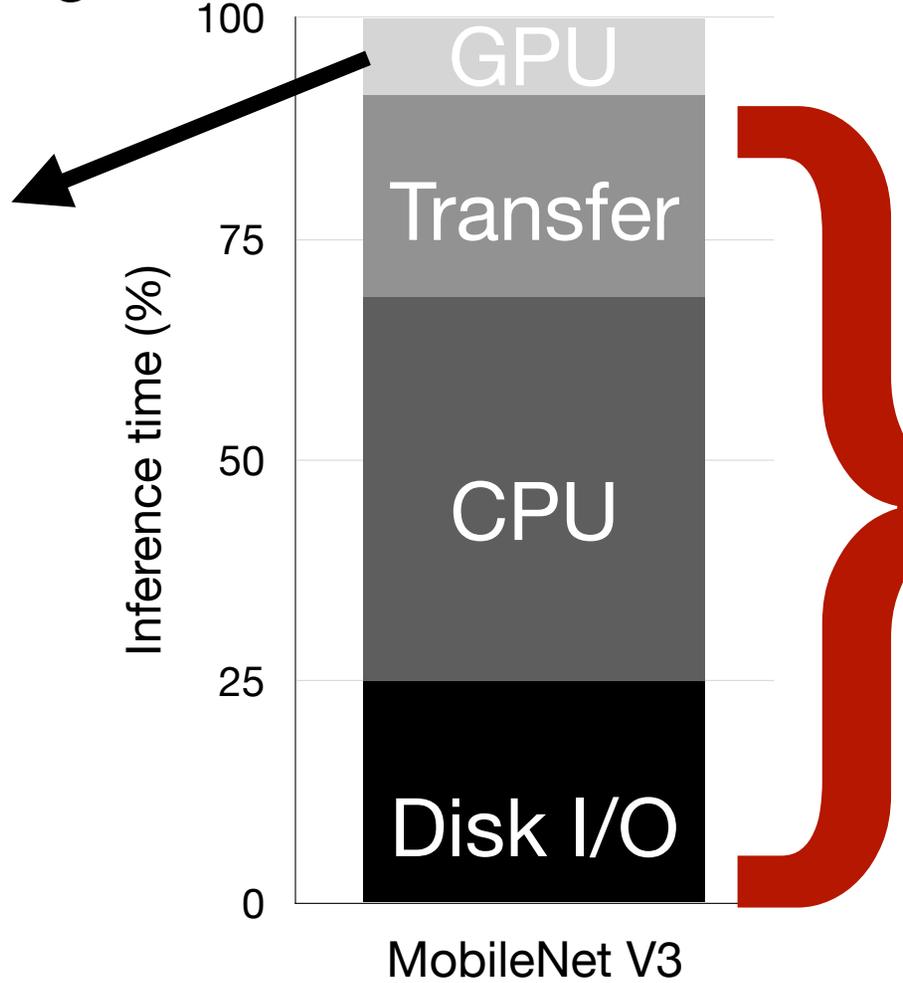


Where does time go?



Where does time go?

**Only 10%
is GPU!**

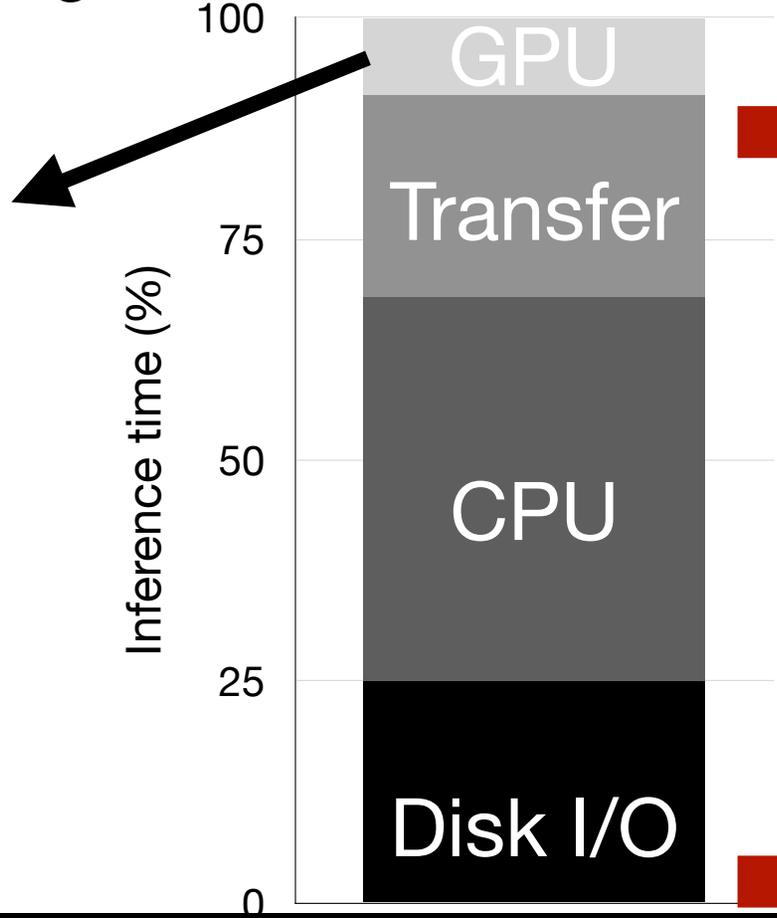


Data: ImageNet
AI Model: MobileNet V3
Machine: V100, PCIe Xeon, SSD
Framework: PyTorch v1

**Data movement/
pre-processing**

Where does time go?

**Only 10%
is GPU!**



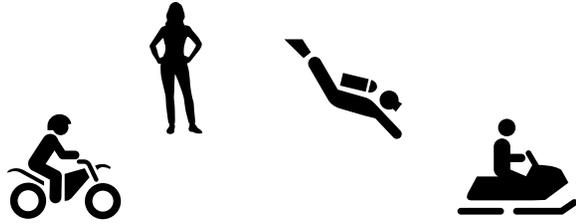
Data: ImageNet
AI Model: MobileNet V3
Machine: V100, PCIe Xeon, SSD
Framework: PyTorch v1

**Data movement/
pre-processing**

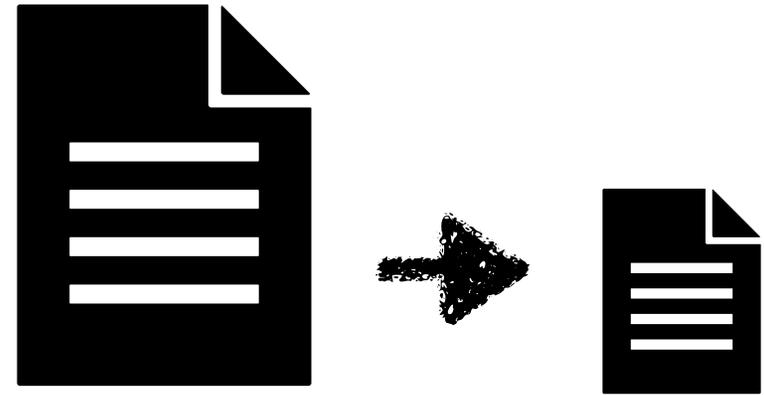
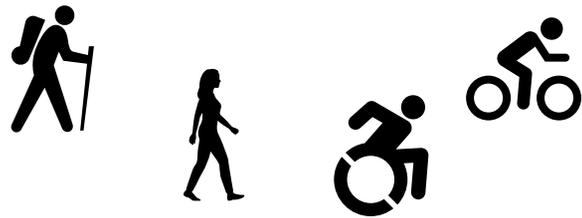
Re-consider Storage for AI

**How do machines
store images today?**

JPEG
Joint Photographic Experts Group



standard



compression

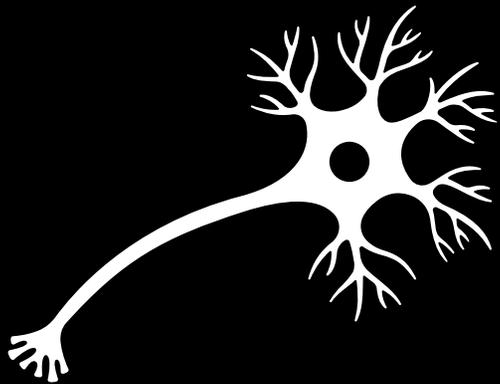
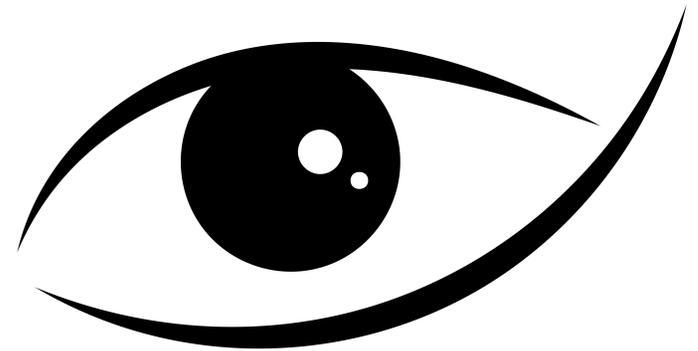
RAW



JPEG

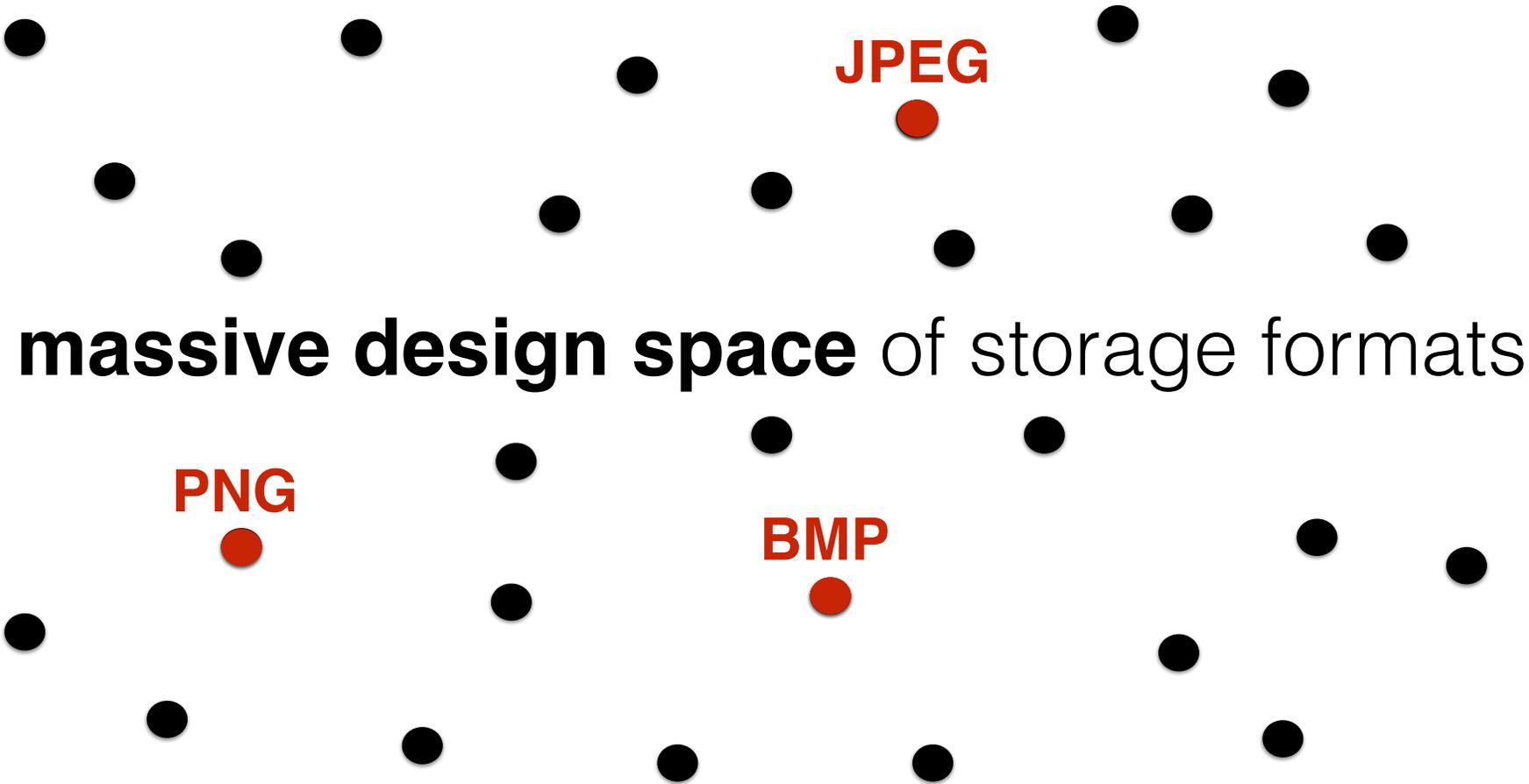


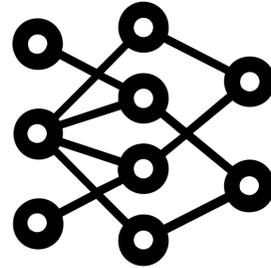
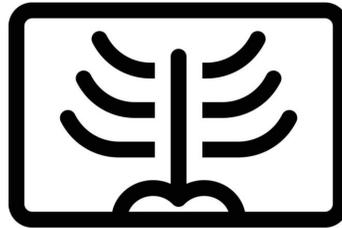
**JPEG is designed for
human eye**



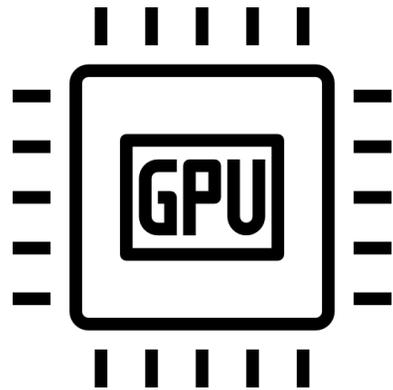
**AI algorithms “see”
images**

A few existing designs





massive space of AI problems



Thesis

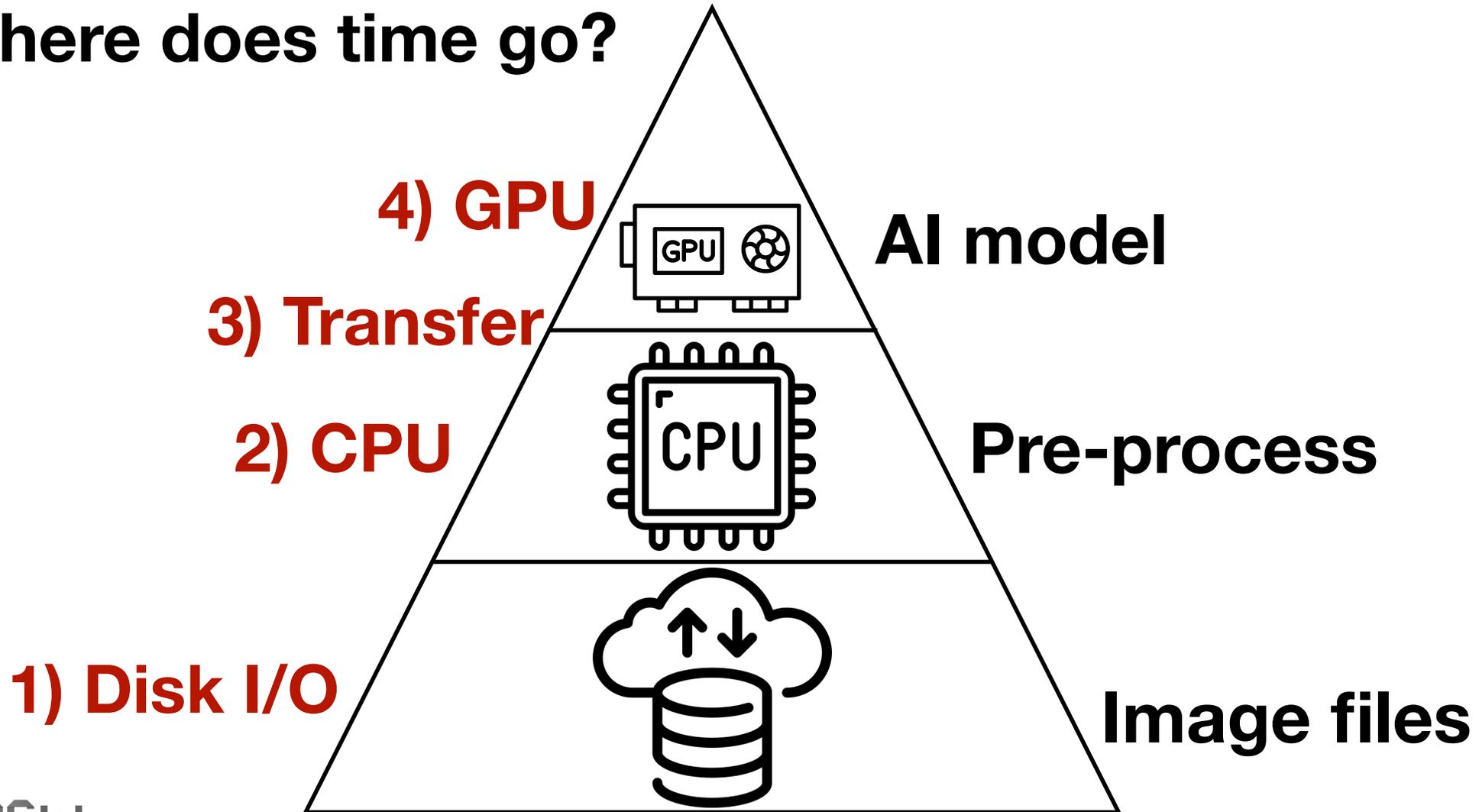
Storage determines end-to-end cost. Using a single storage for all problems results in inefficient AI systems. Orders of magnitude higher performance is only possible if the storage is tailored to the AI problem.

Image Calculator

A self-designing storage format that shapes itself for a given AI problem

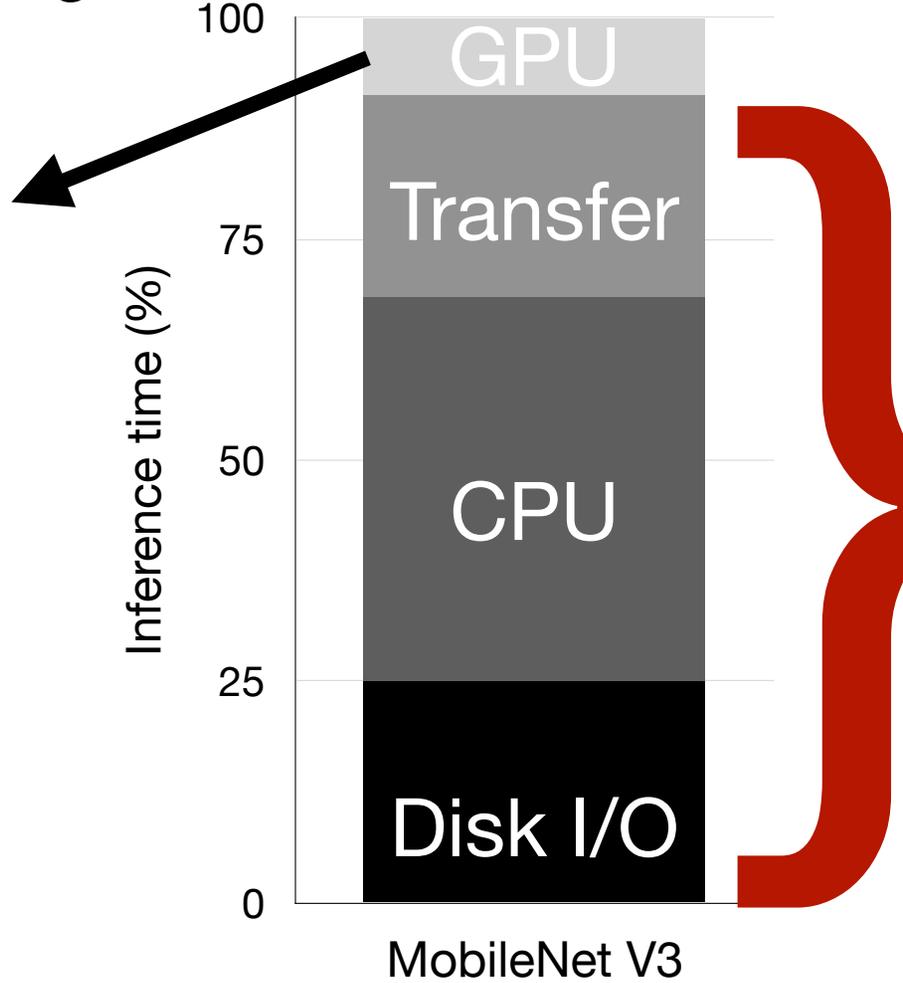
Recap from previous class

Where does time go?



Where does time go?

**Only 10%
is GPU!**

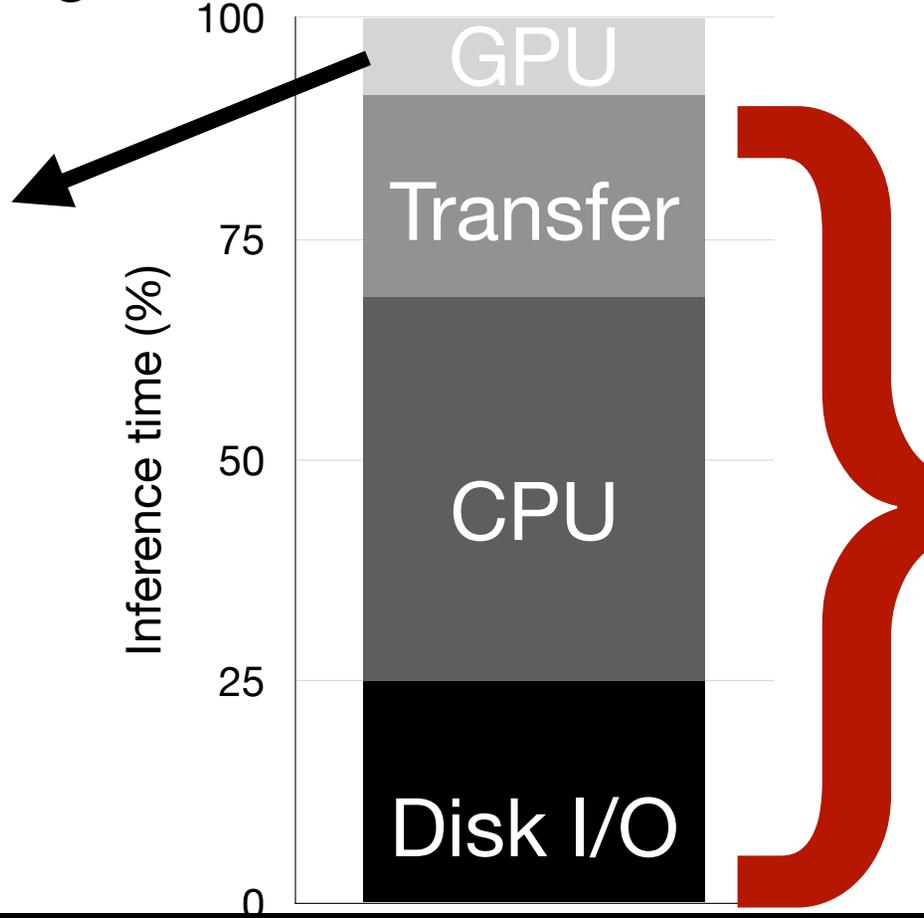


Data: ImageNet
AI Model: MobileNet V3
Machine: V100, PCIe Xeon, SSD
Framework: PyTorch v1

**Data movement/
pre-processing**

Where does time go?

**Only 10%
is GPU!**

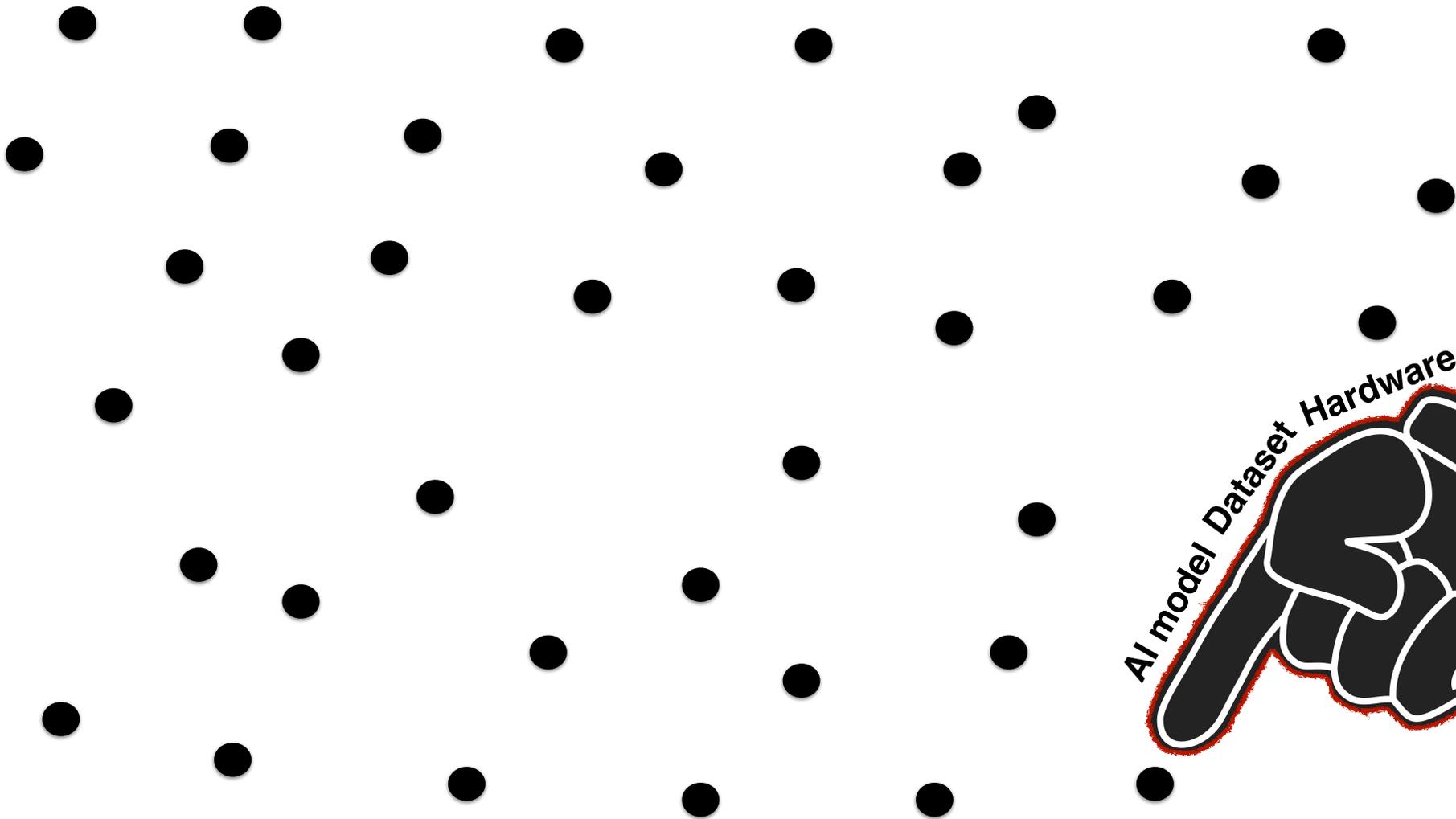


Data: ImageNet
AI Model: MobileNet V3
Machine: V100, PCIe Xeon, SSD
Framework: PyTorch v1

**Data movement/
pre-processing**

Re-consider Storage for AI

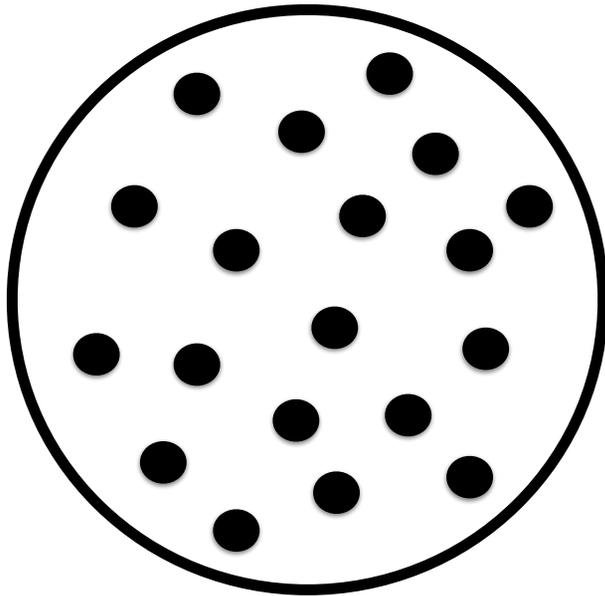
A self-designing File format



AI model Dataset Hardware Budget

A stylized hand icon pointing to the left. The hand is black with white outlines. The index finger is extended and points towards the left. The text "AI model Dataset Hardware Budget" is written along the length of the index finger, following its curve. The text is in a bold, sans-serif font.

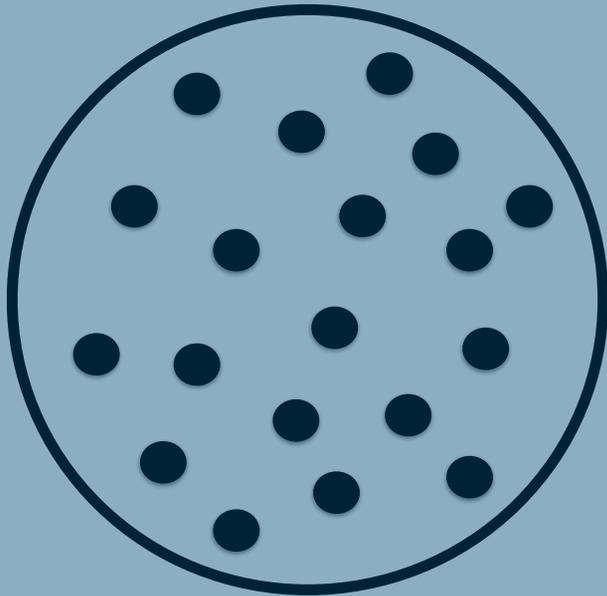
Design space



Search

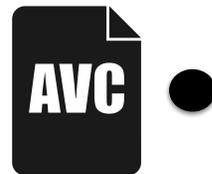
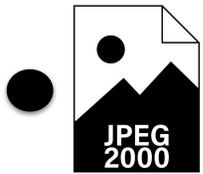
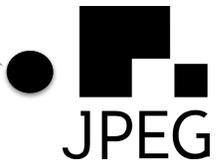
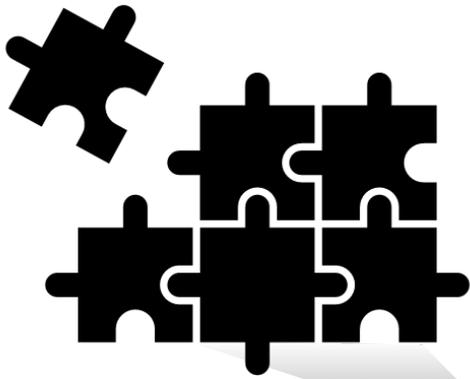


Design space

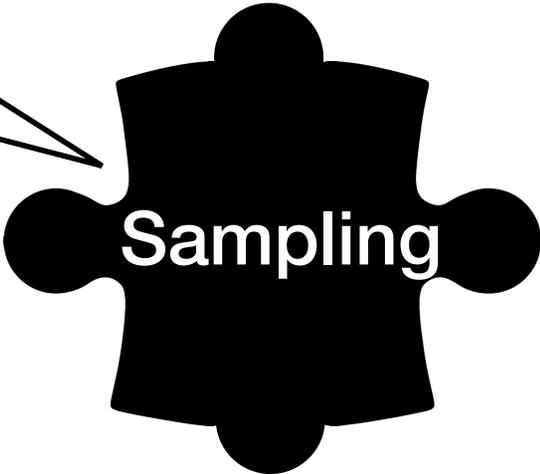


Search

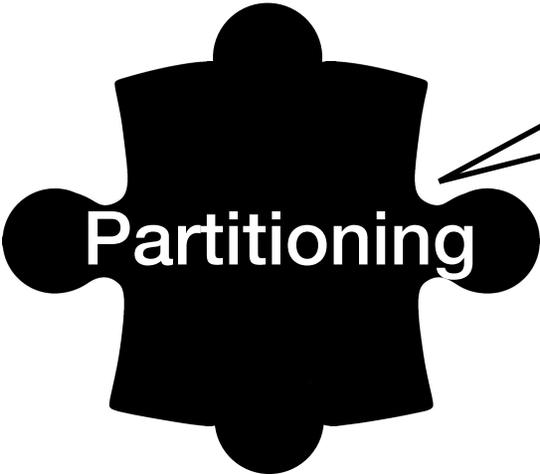




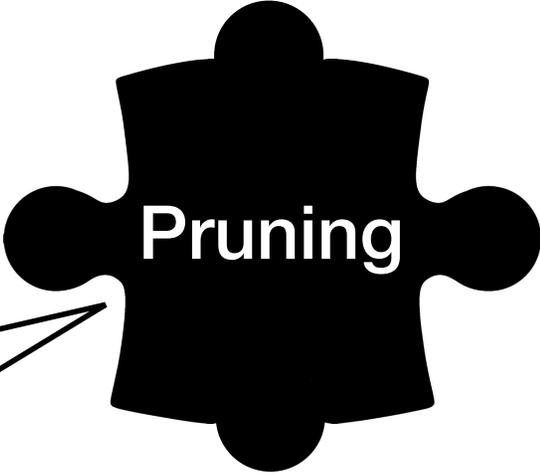
Remove rows/
columns



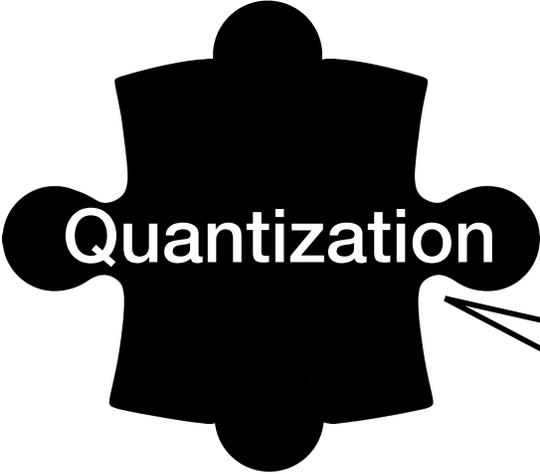
Processing
granularity



Remove
unuseful data



Magnitude
reduction



Remove rows/
columns

Processing
parallelity

10^{150K}

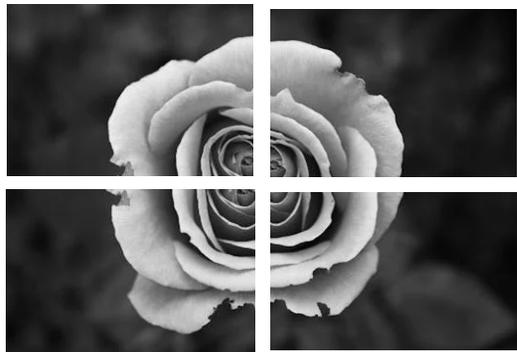
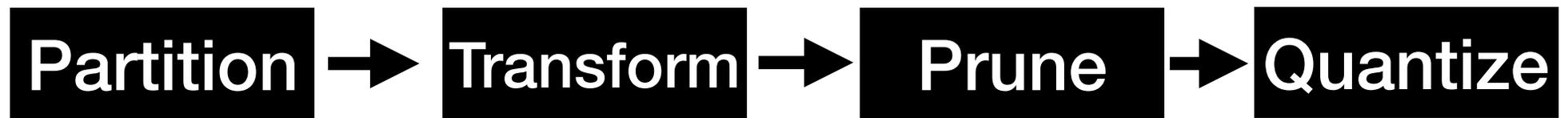
possibilities

Remove
unuseful data

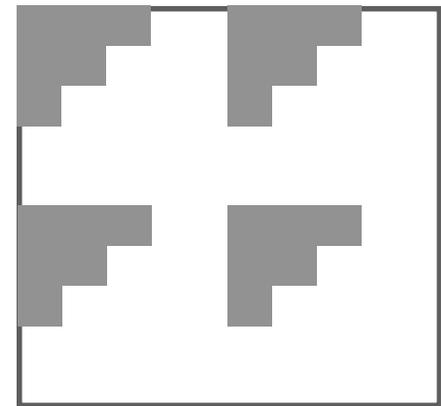
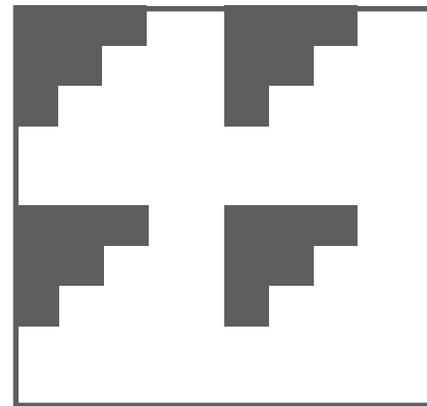
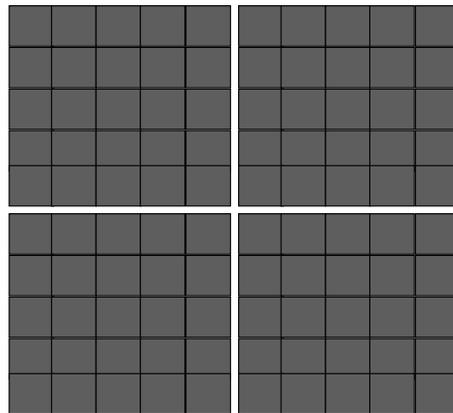
magnitude
reduction

Designing A Domain: Pruning

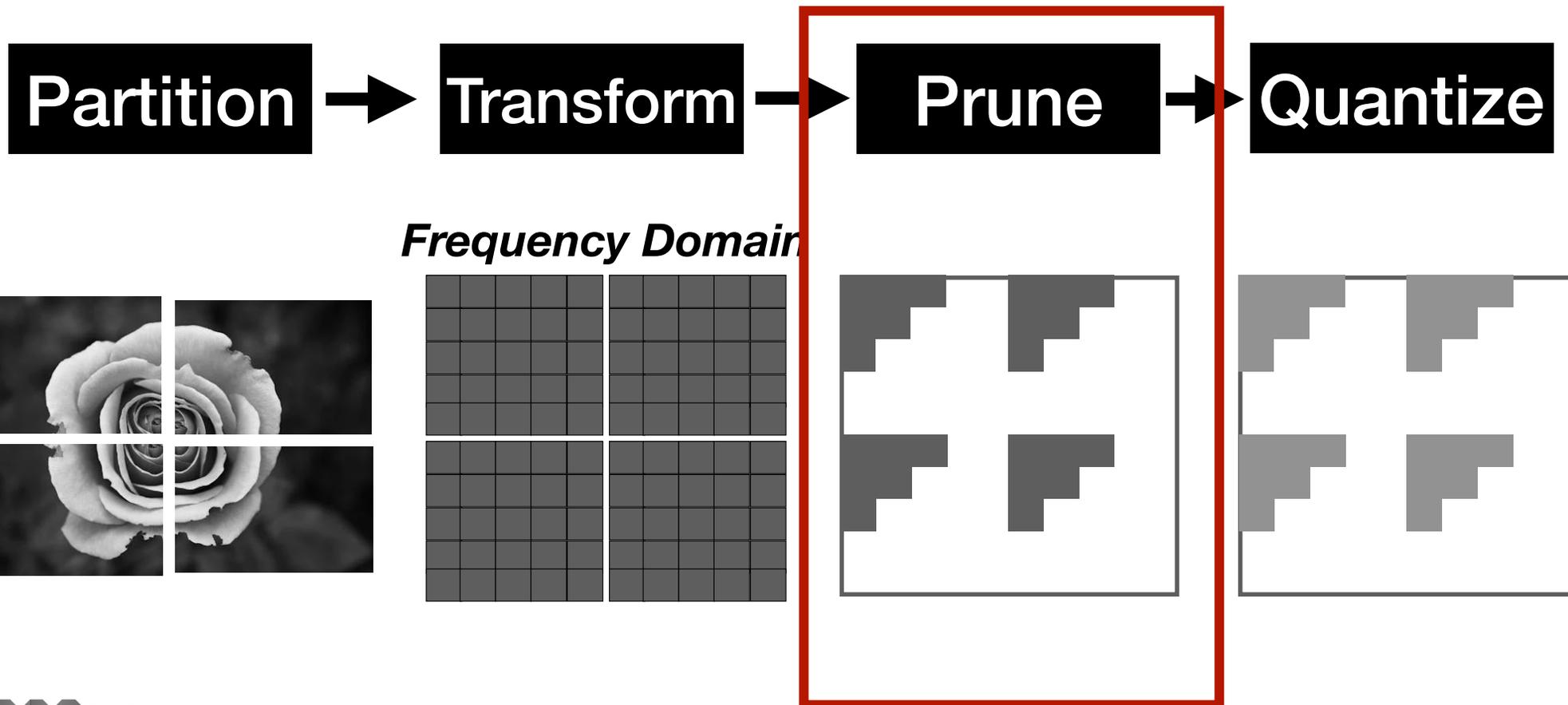
Background: Image Storage

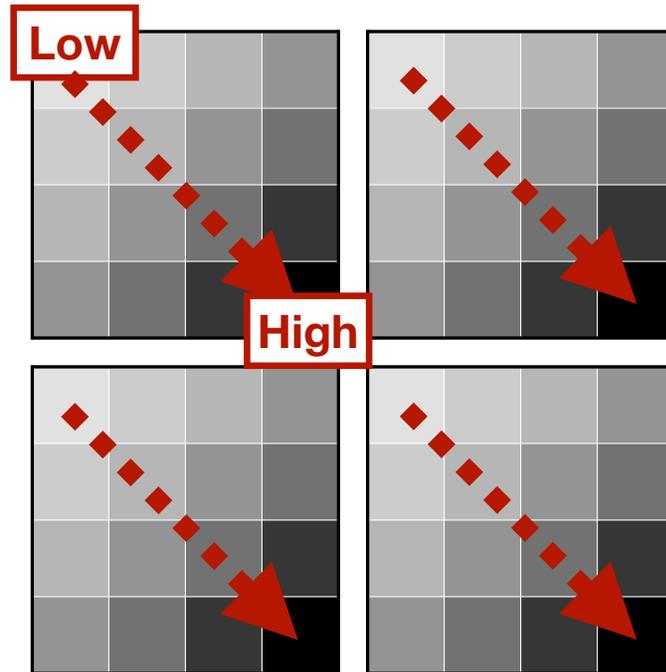


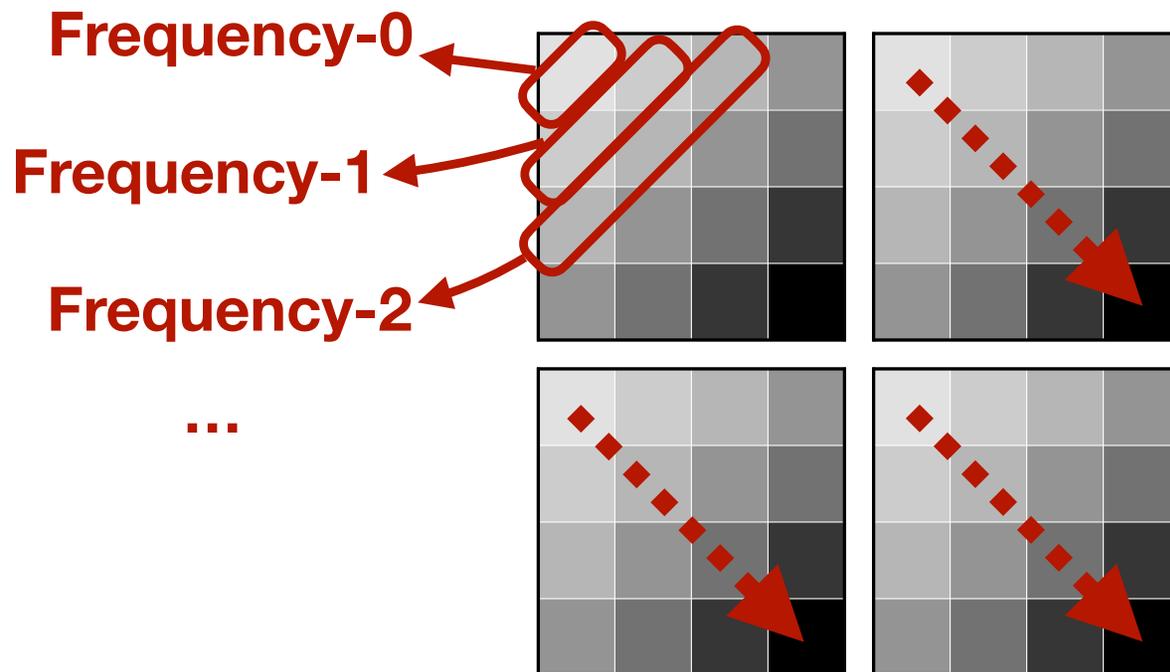
Frequency Domain



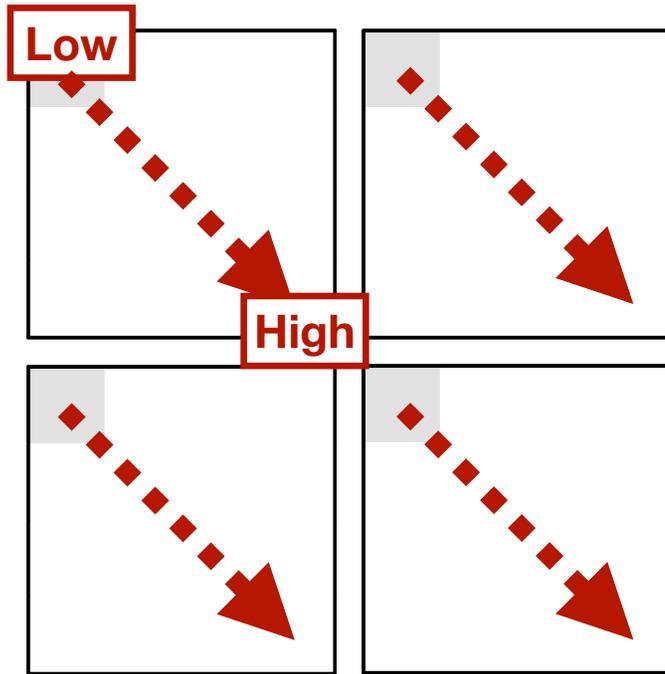
Background: Image Storage



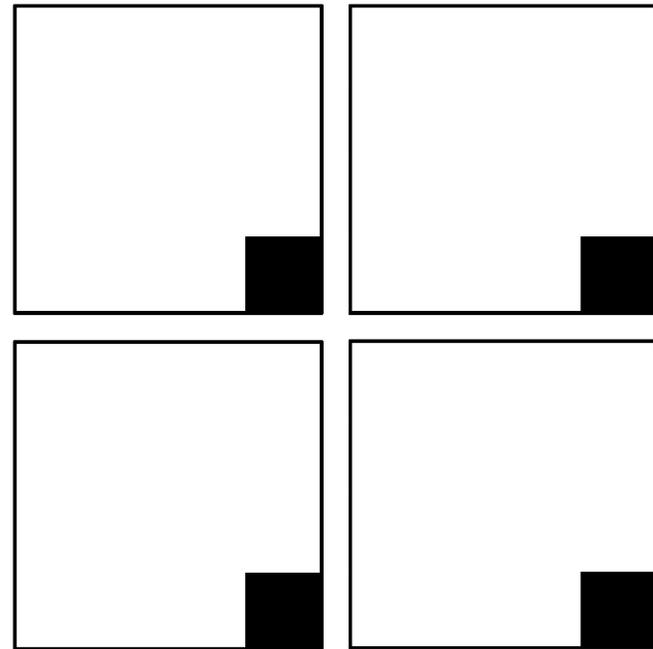




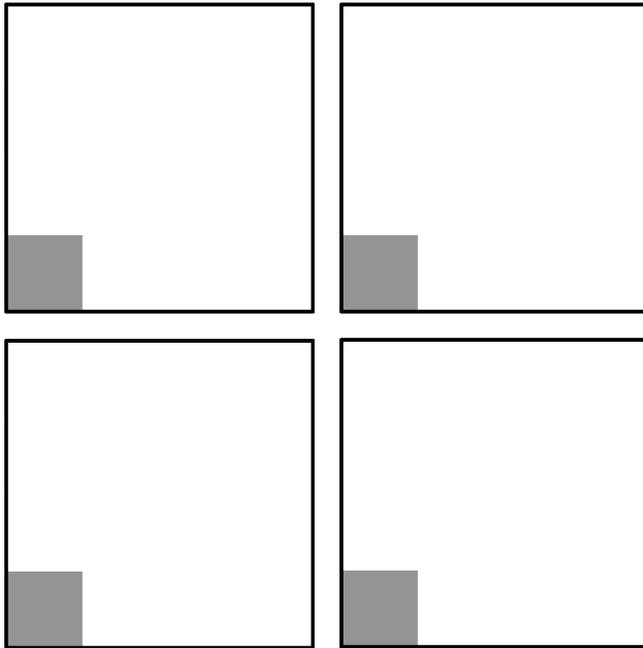
(1) High to low



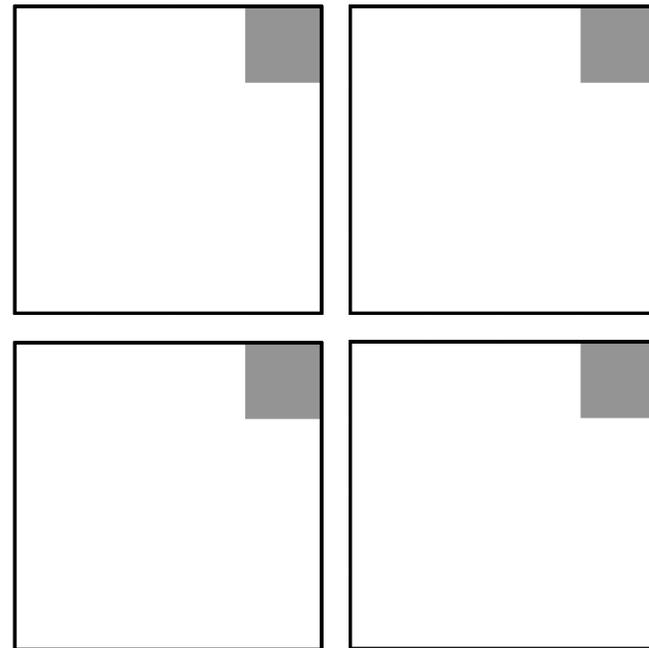
(2) Low to high

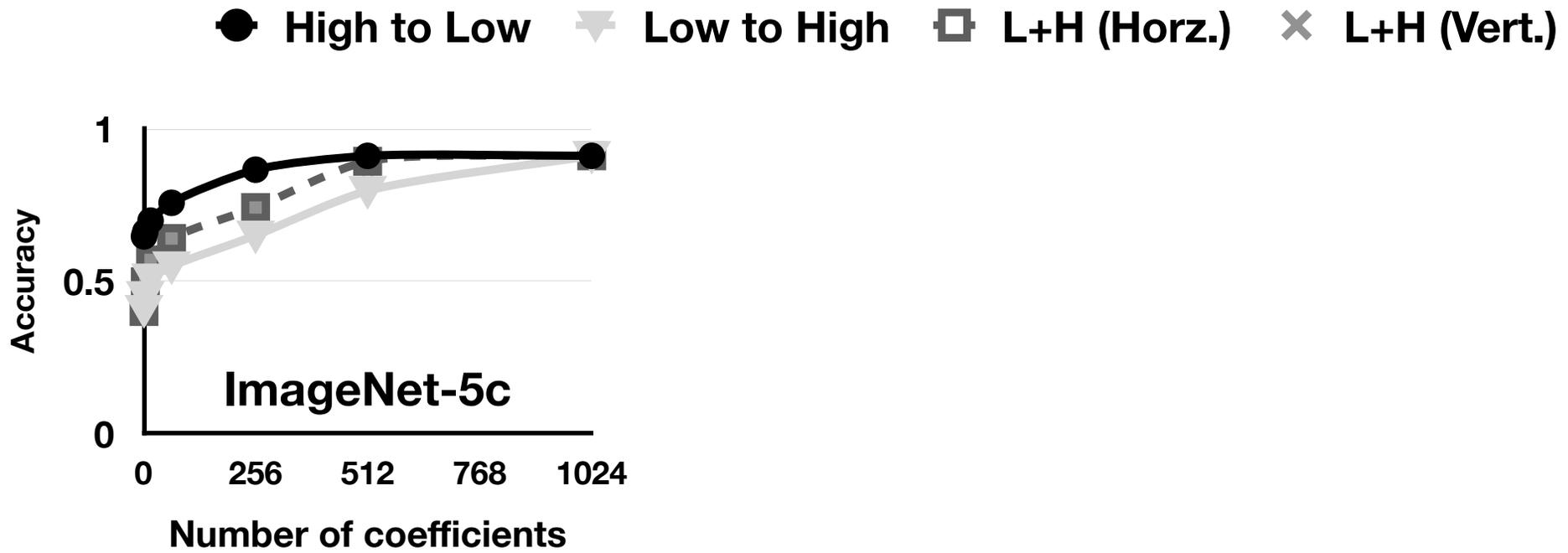


(3) Low & high — Horizontal

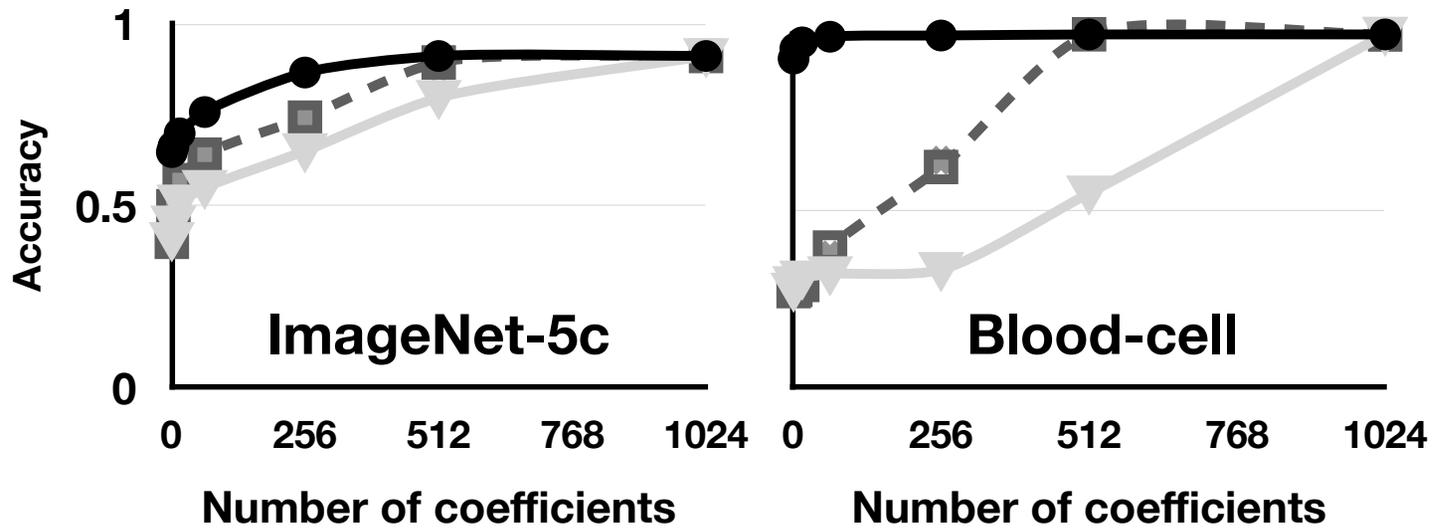


(4) Low & high — Vertical

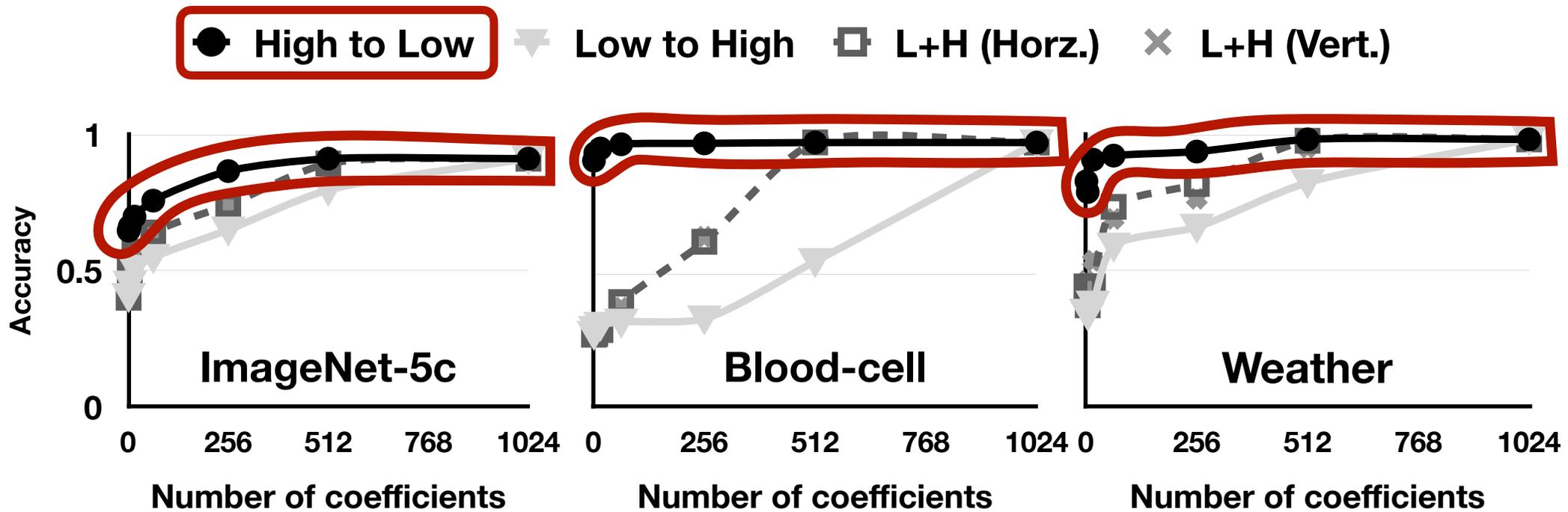




● High to Low ▼ Low to High □ L+H (Horz.) × L+H (Vert.)



ResNet50, A100, PyTorch v1

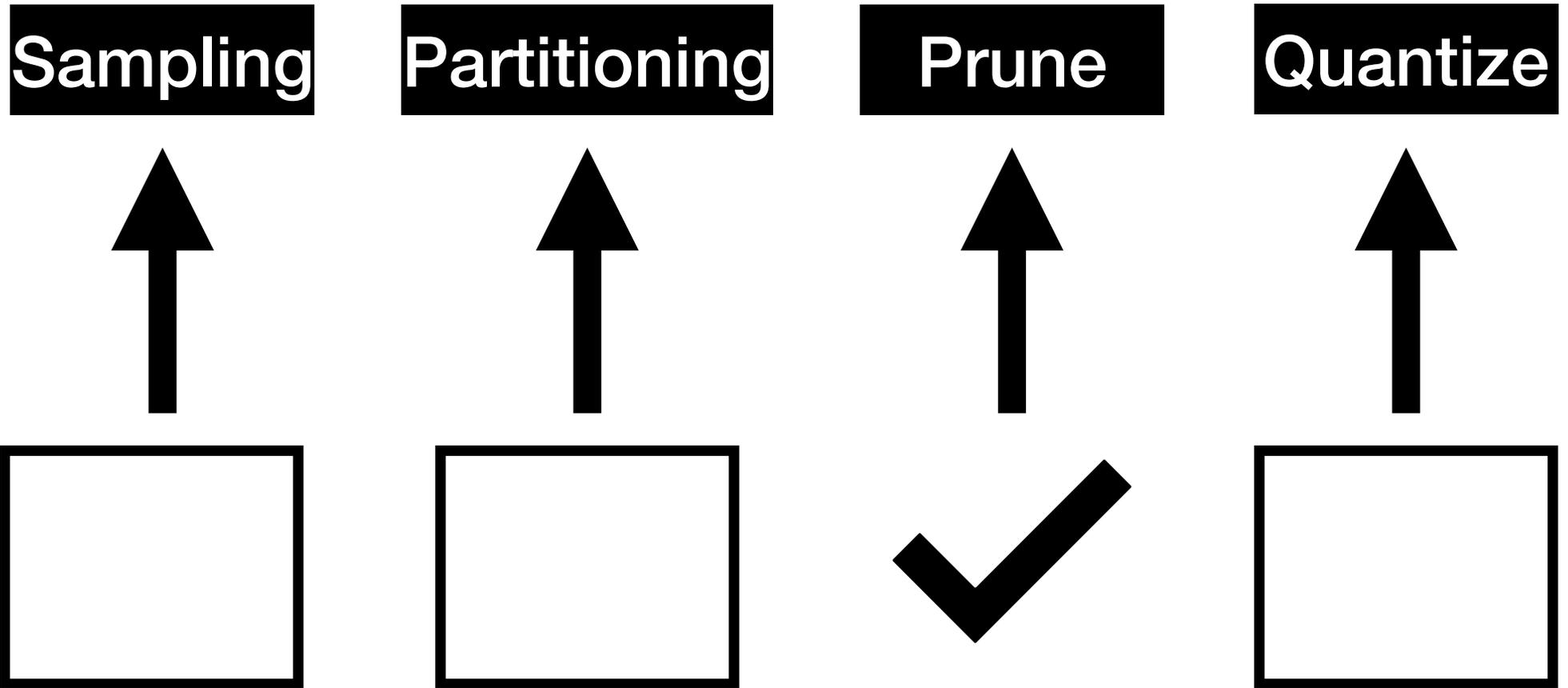


***Pruning is from
high to low frequency***

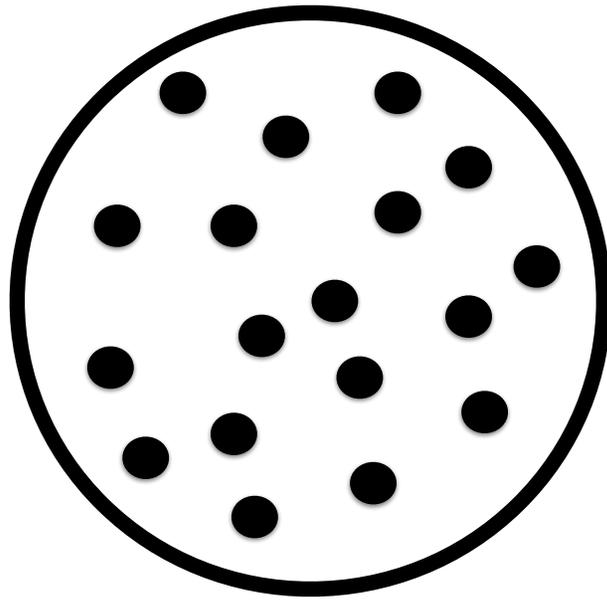
Pruning Domain

$2B \times B \rightarrow 2B-1$

Other primitives

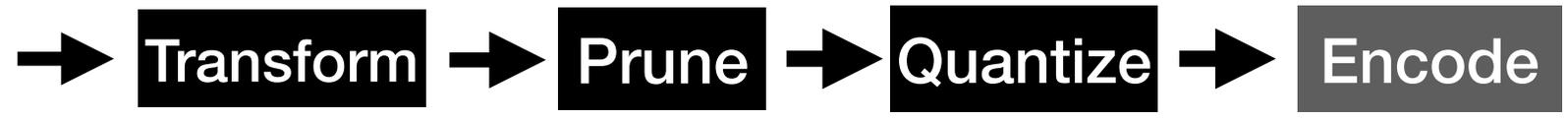


4104 new designs



Reading & Writing Algorithms

Sample & Partition



Frequency Domain

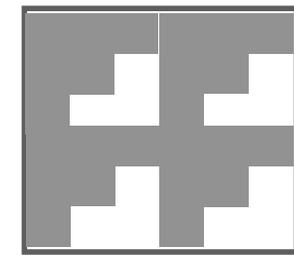
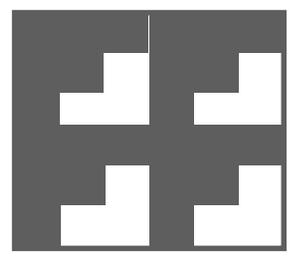
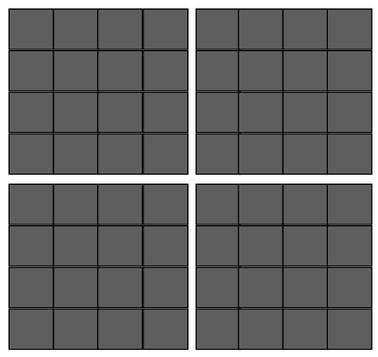
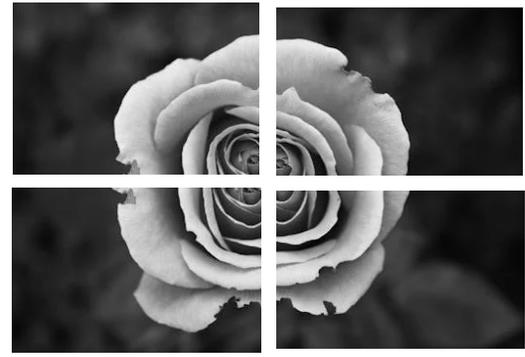
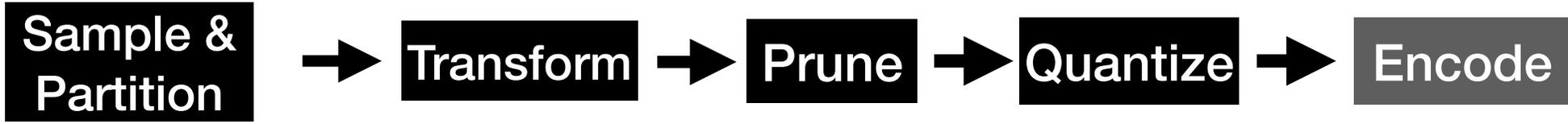


image.ic

110010
101101



Frequency Domain

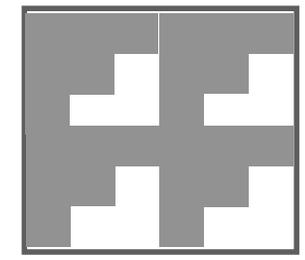
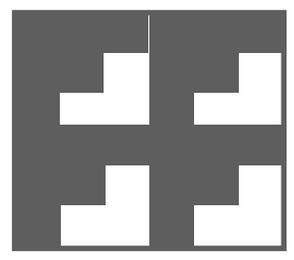
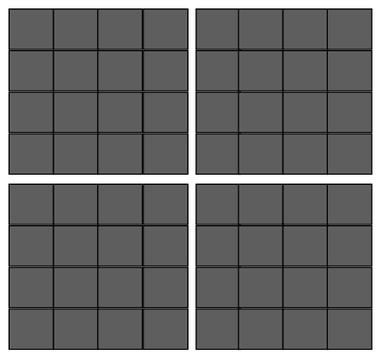
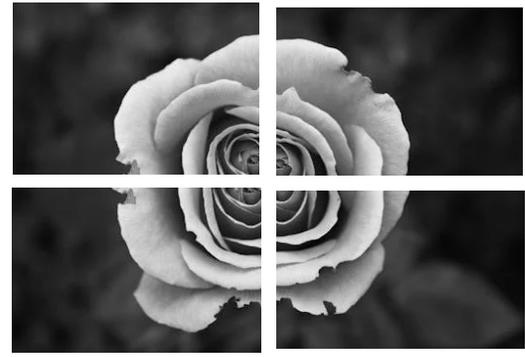
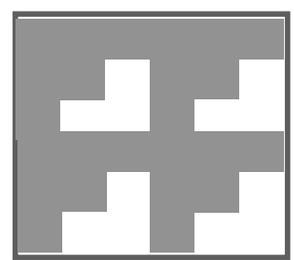
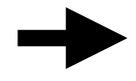


image.ic
110010
101101

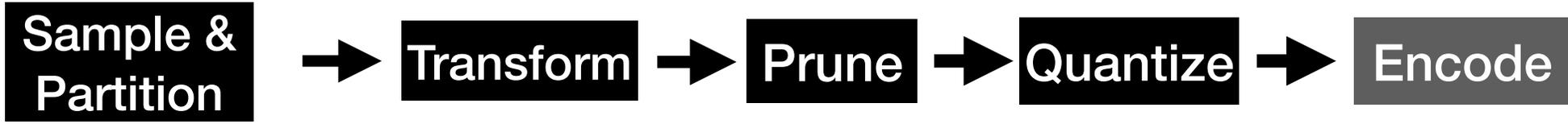
Decode

image.ic
110010
101101



“Learning in Frequency Domain”

Gueguen et al., NIPS’18; Xu et al., CVPR’20; Tancik et al., ICLR’20; Ning et al., ICLR’25



Frequency Domain

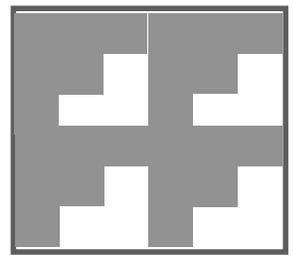
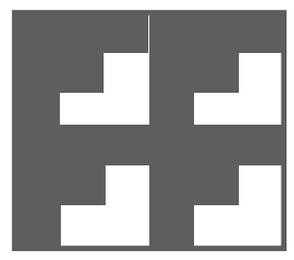
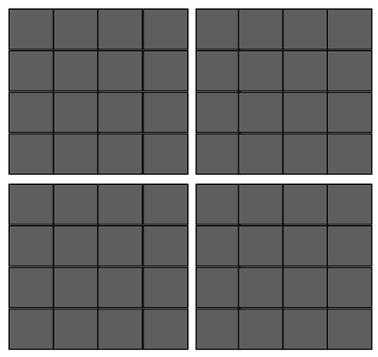
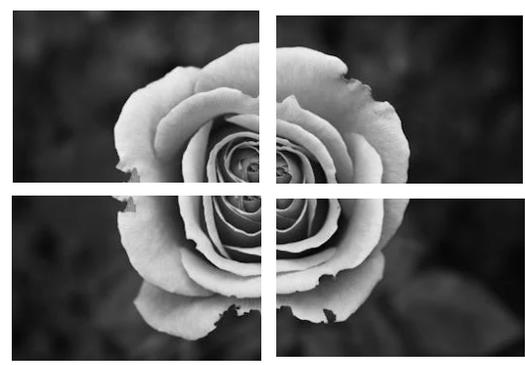
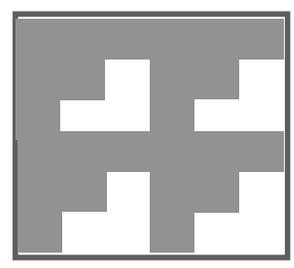
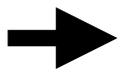


image.ic
110010
101101

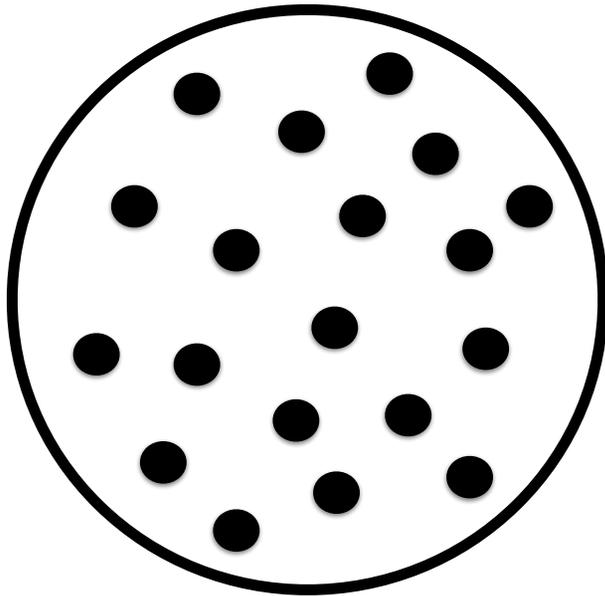
Decode

image.ic
110010
101101



- ✓ **Fast & Simple Reconstruction**
- ✓ **Low GPU Time**

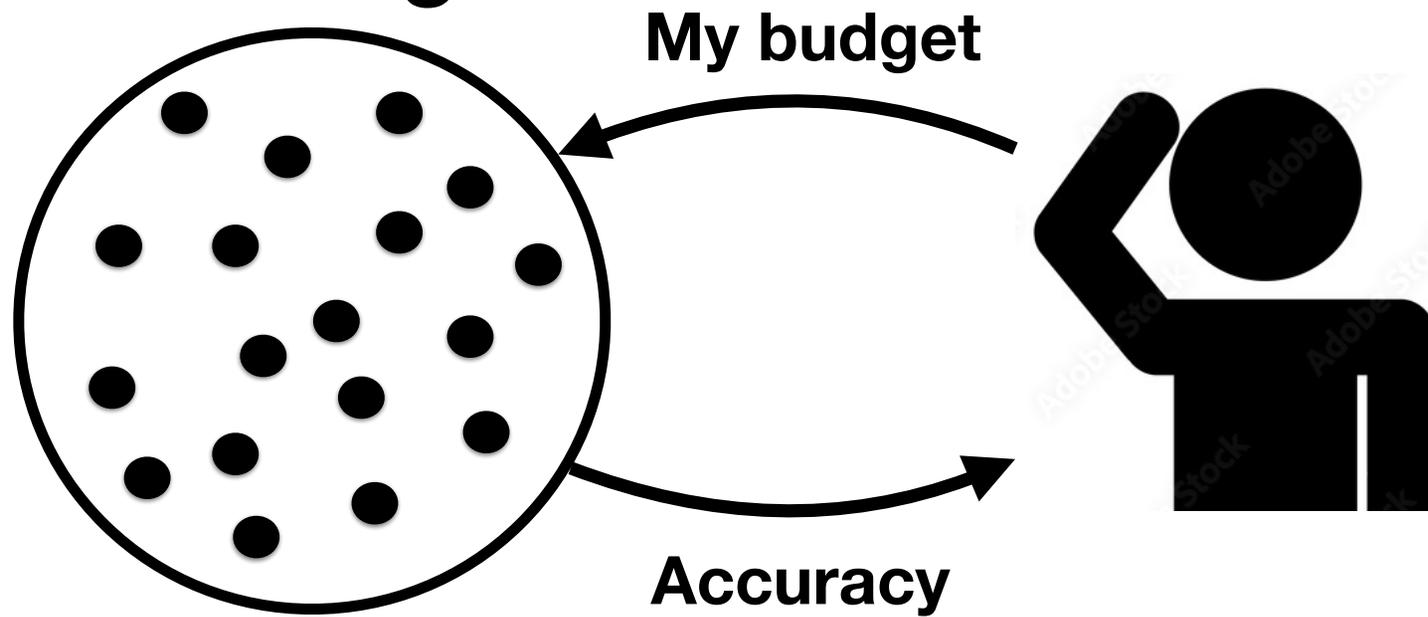
Design space



Search

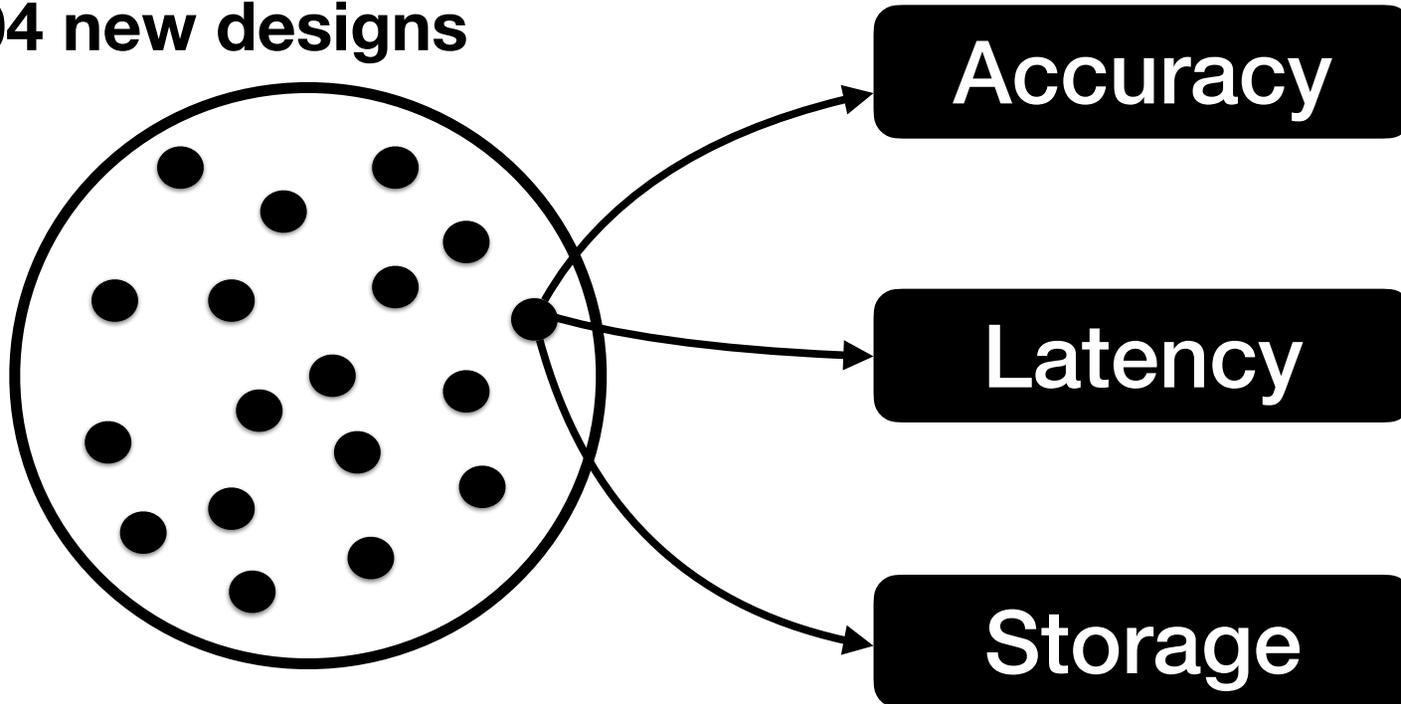


4104 new designs



Performance models

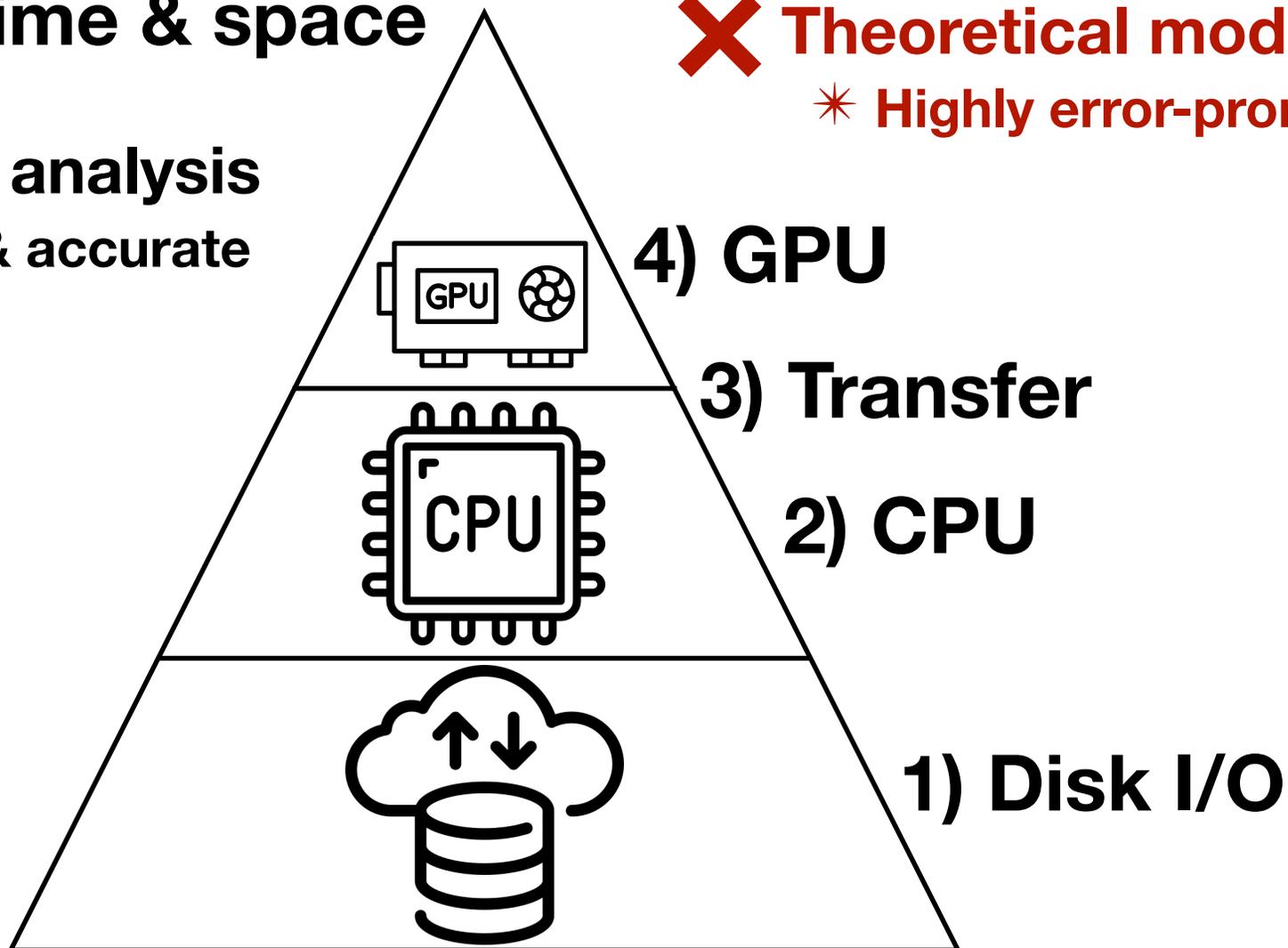
4104 new designs



Inf./training time & space

✓ **Empirical analysis**
* Cheap & accurate

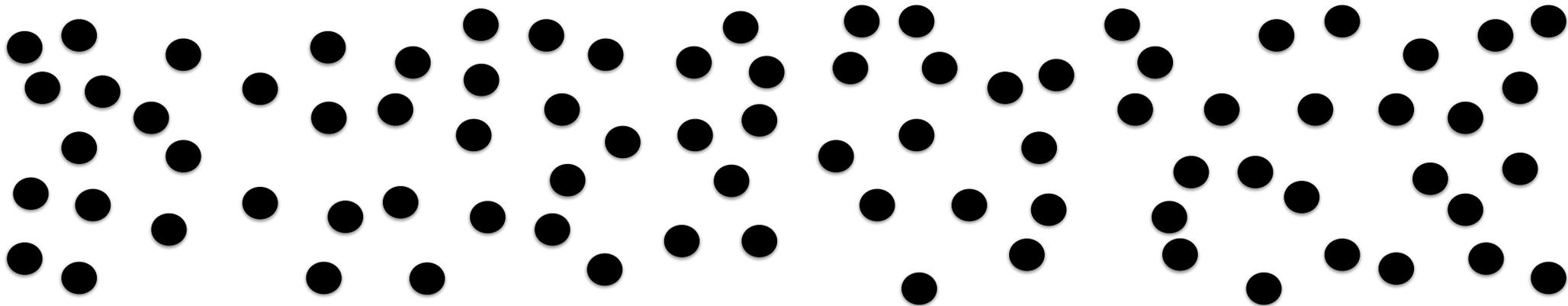
✗ **Theoretical model**
* Highly error-prone



Accuracy model

**Insight 1: Inference Cost and Accuracy
are Correlated**

Sort

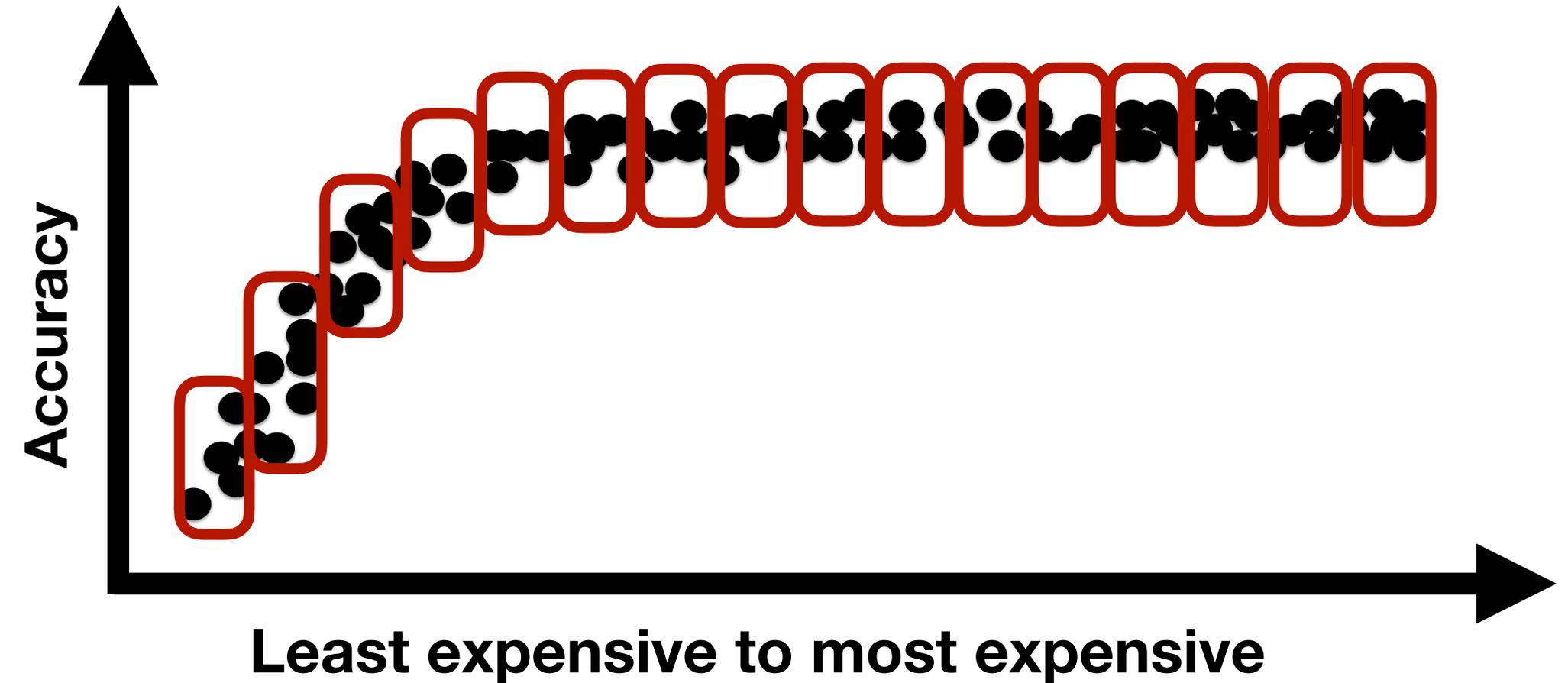


Least expensive to most expensive

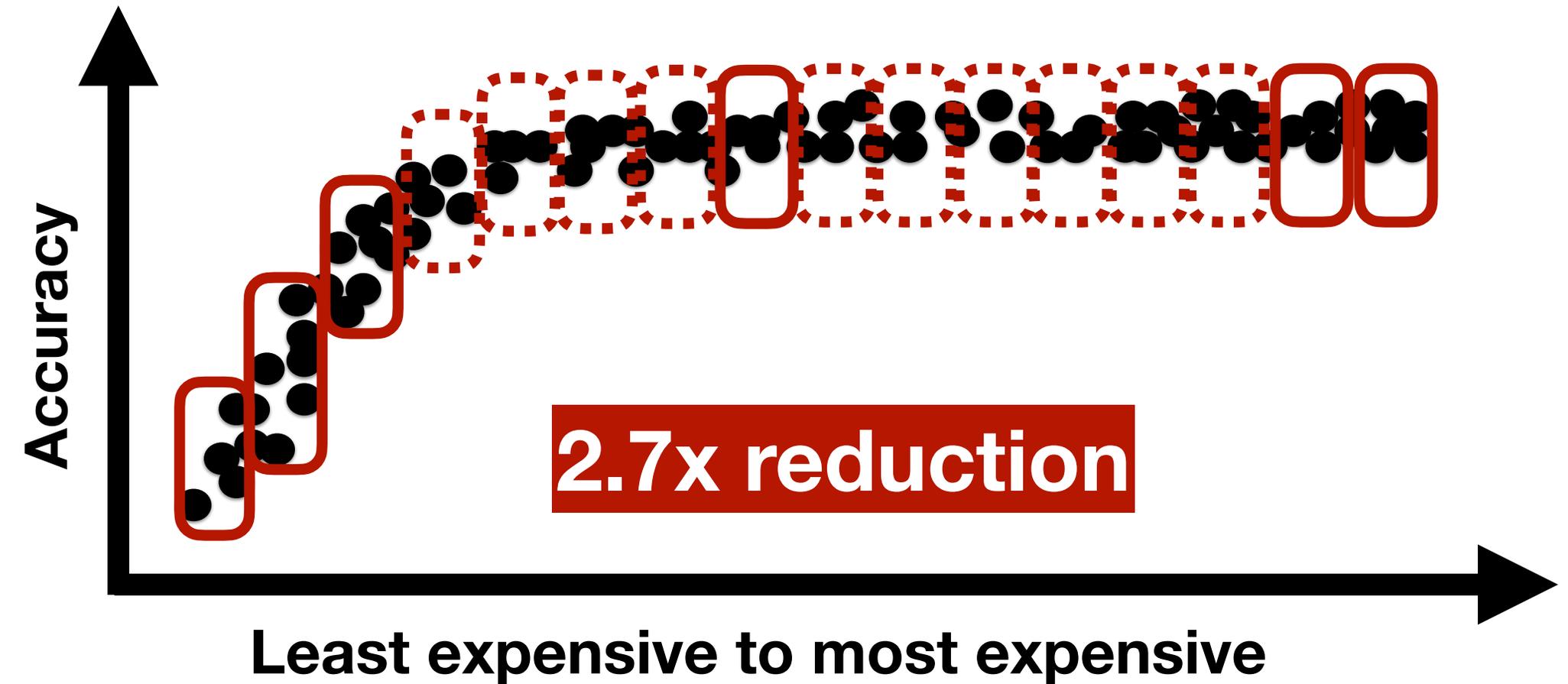
ImageNet-5c, ResNet50
A100, PyTorch v1



Bucket-sampling & interpolation

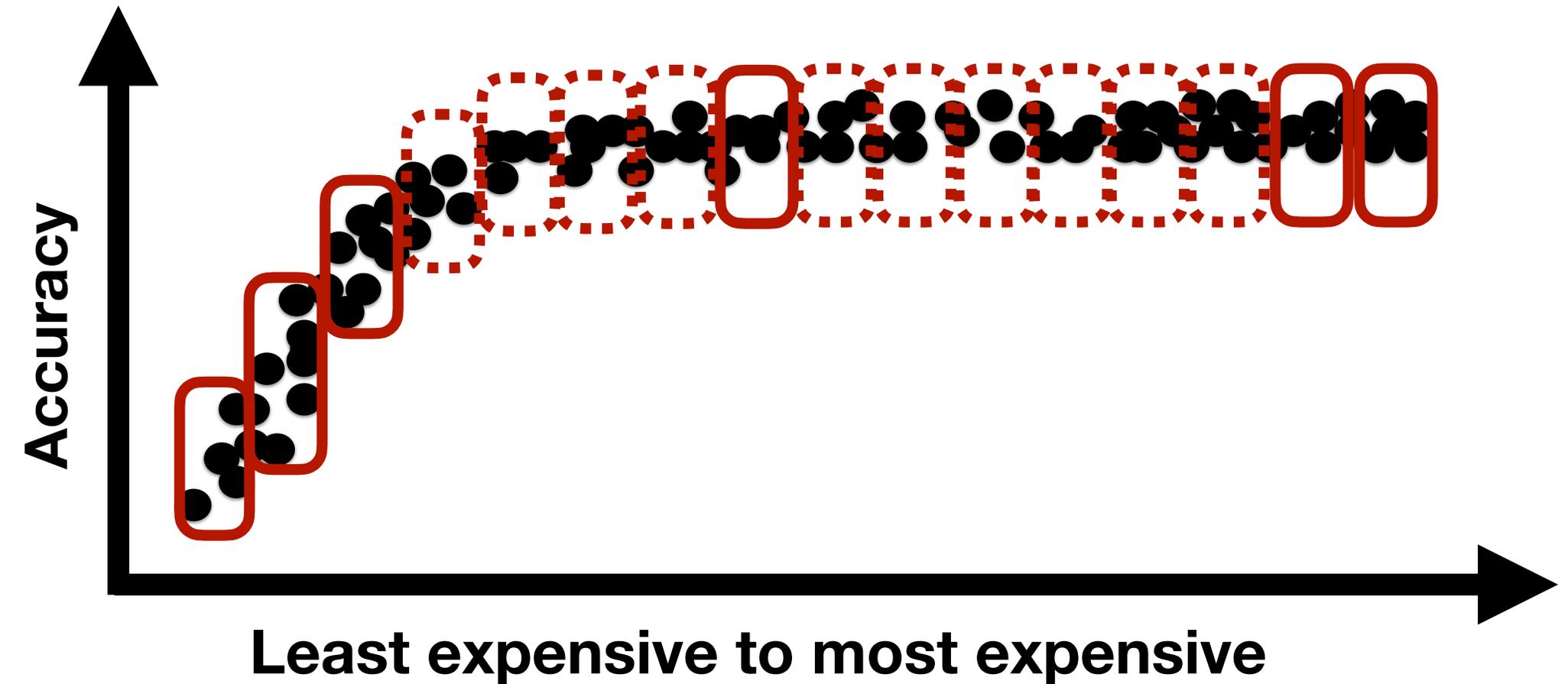


Bucket-sampling & interpolation

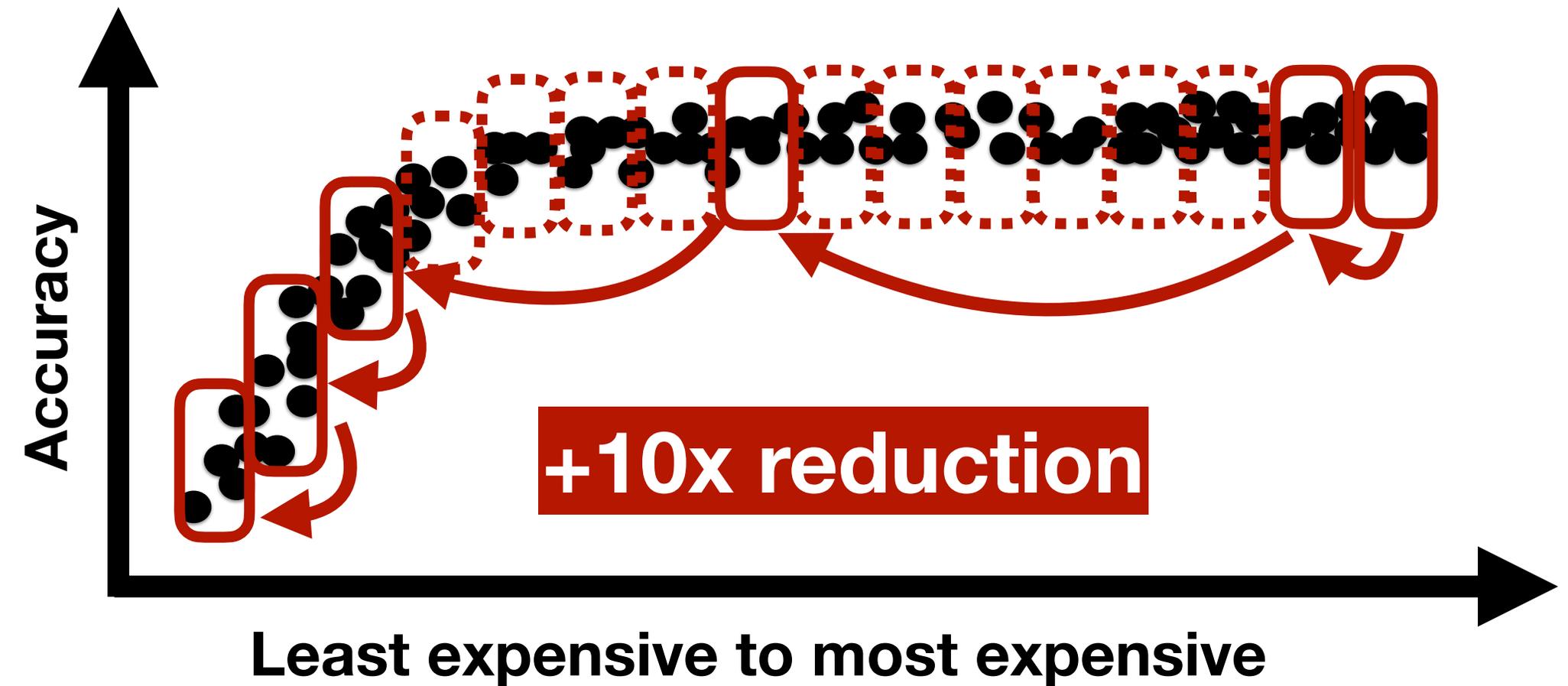


Insight 2: Storage Formats are Correlated

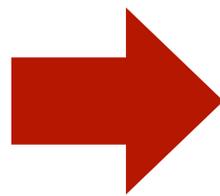
Bucket-sampling & interpolation



Transfer learning

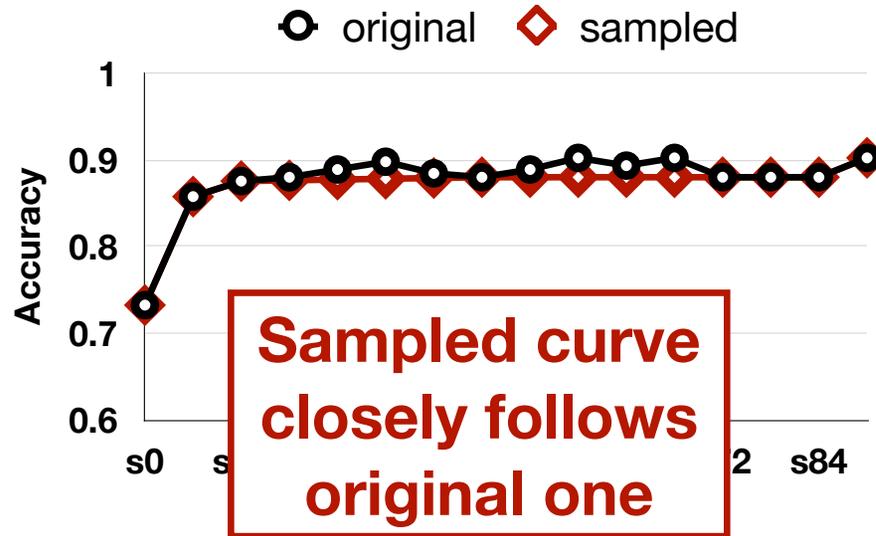
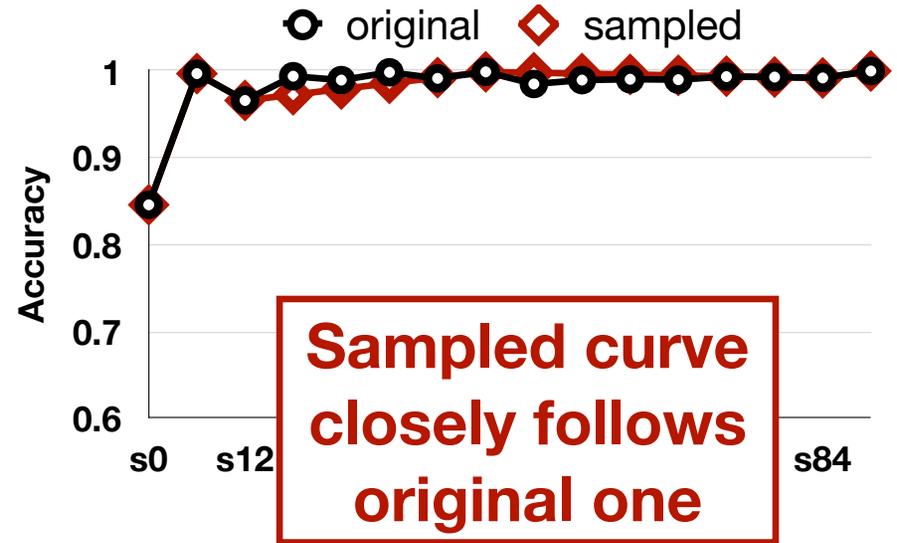
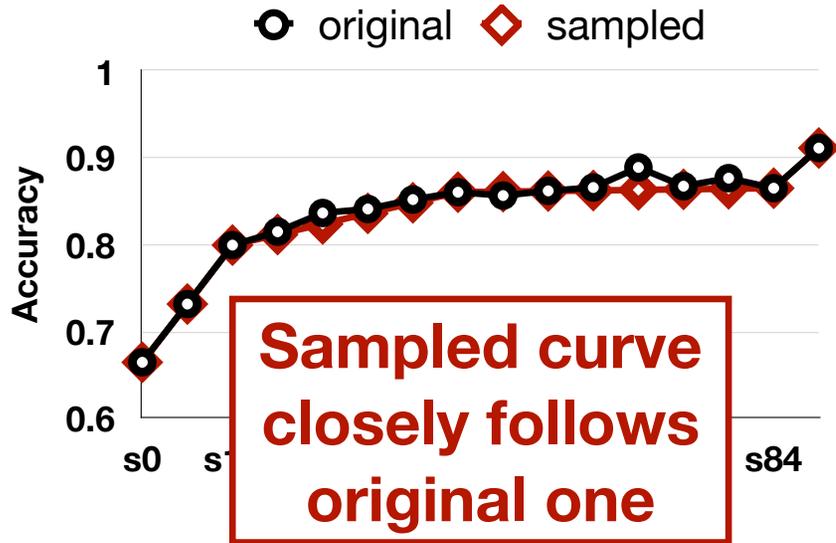


**~4K AI model
trainings**

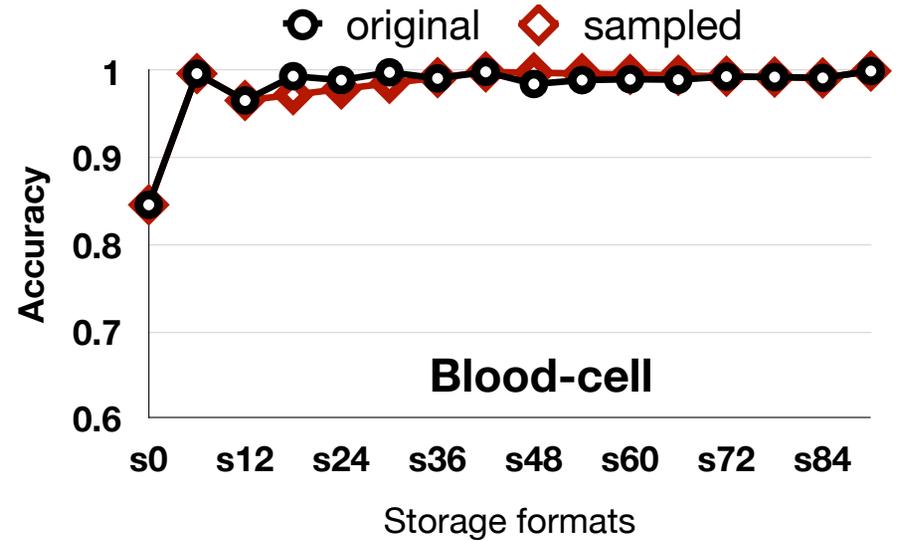
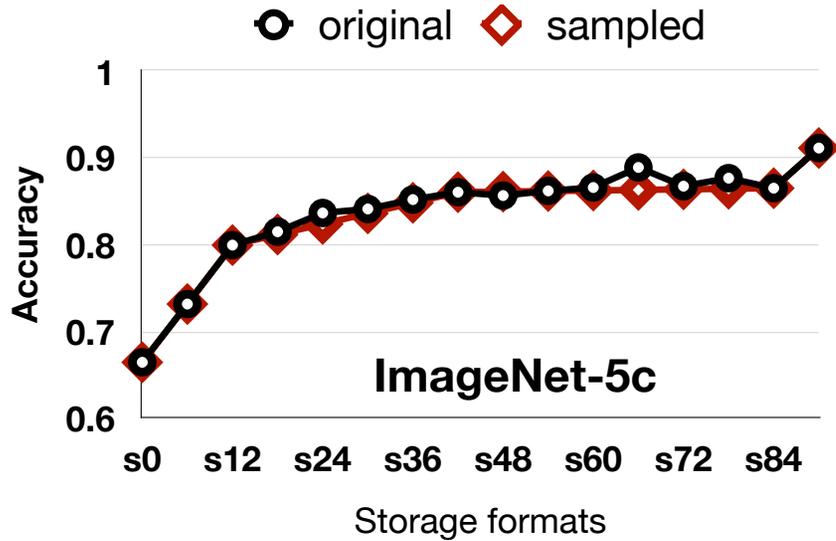


**150 AI model
trainings**

Verification



ResNet50, A100, PyTorch v1



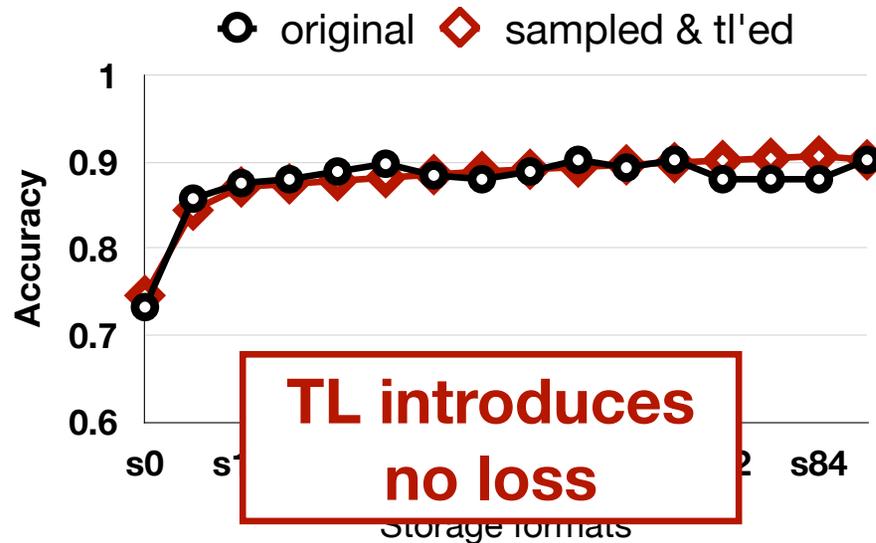
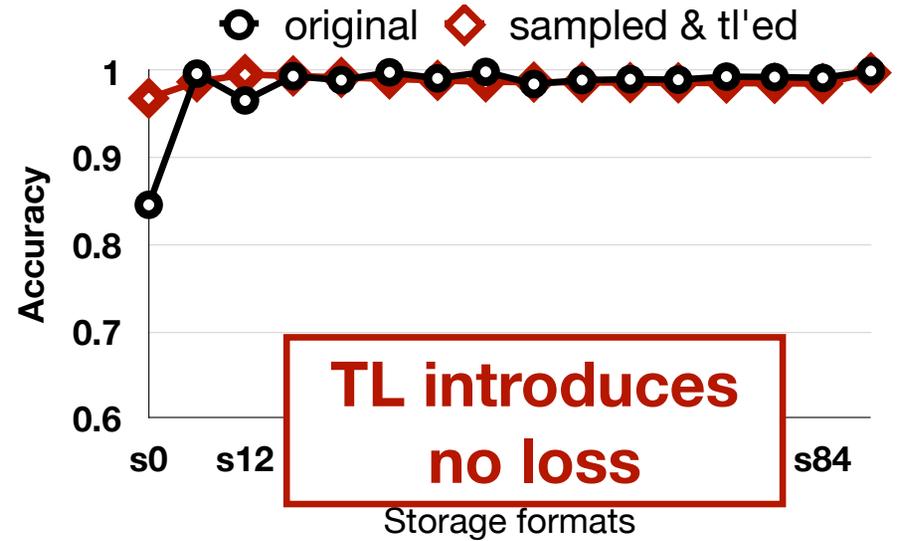
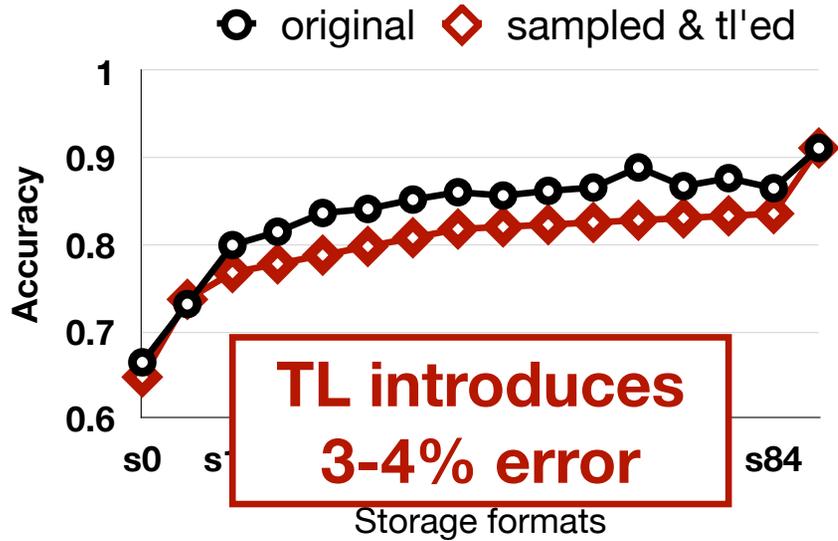
Legend: original (circle), sampled (diamond)

Sampling reduces cost with no loss

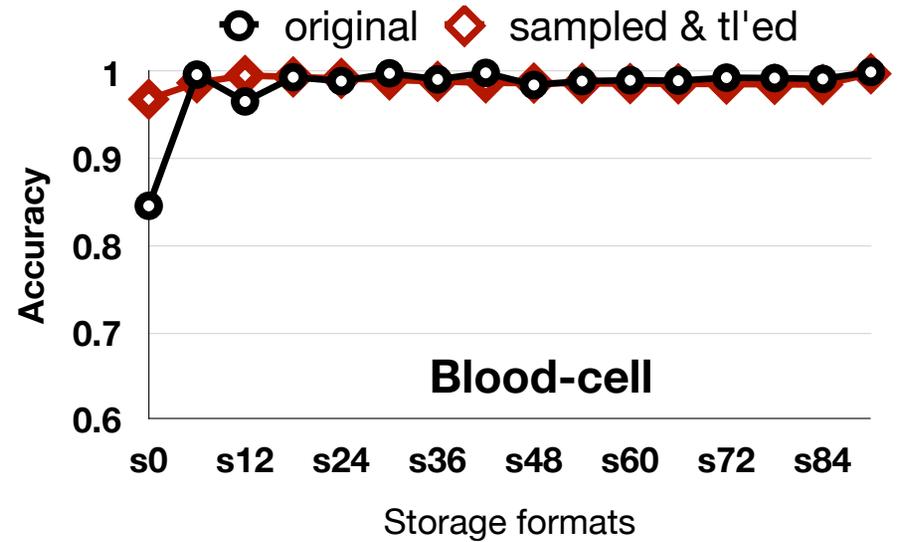
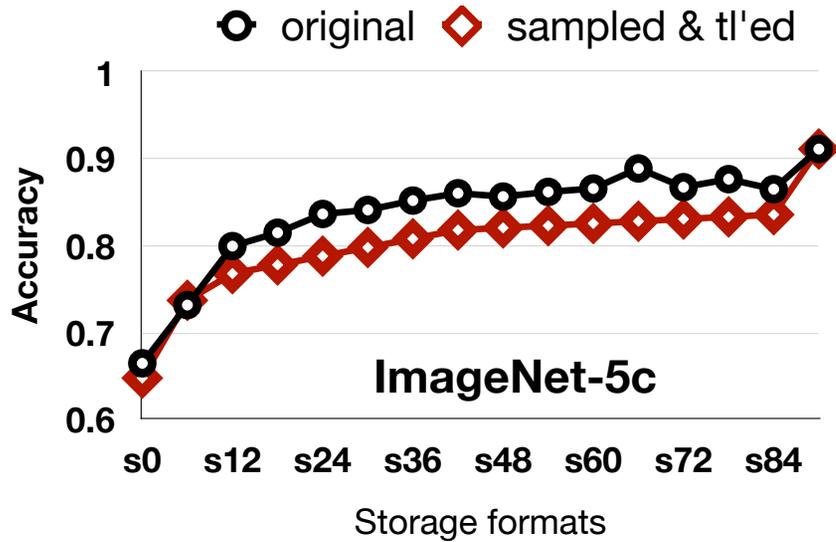
0.6 1

s0 s12 s24 s36 s48 s60 s72 s84

Storage formats



ResNet50, A100, PyTorch v1

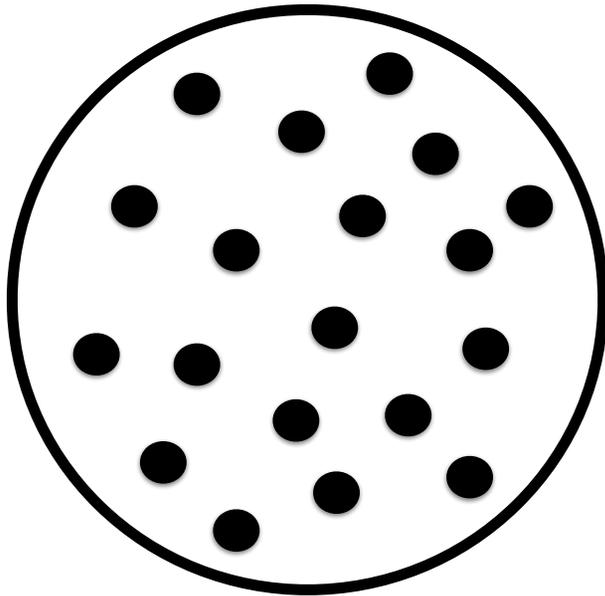


original sampled & tl'ed

Sampling & TL reduces cost by 30x with little or no loss



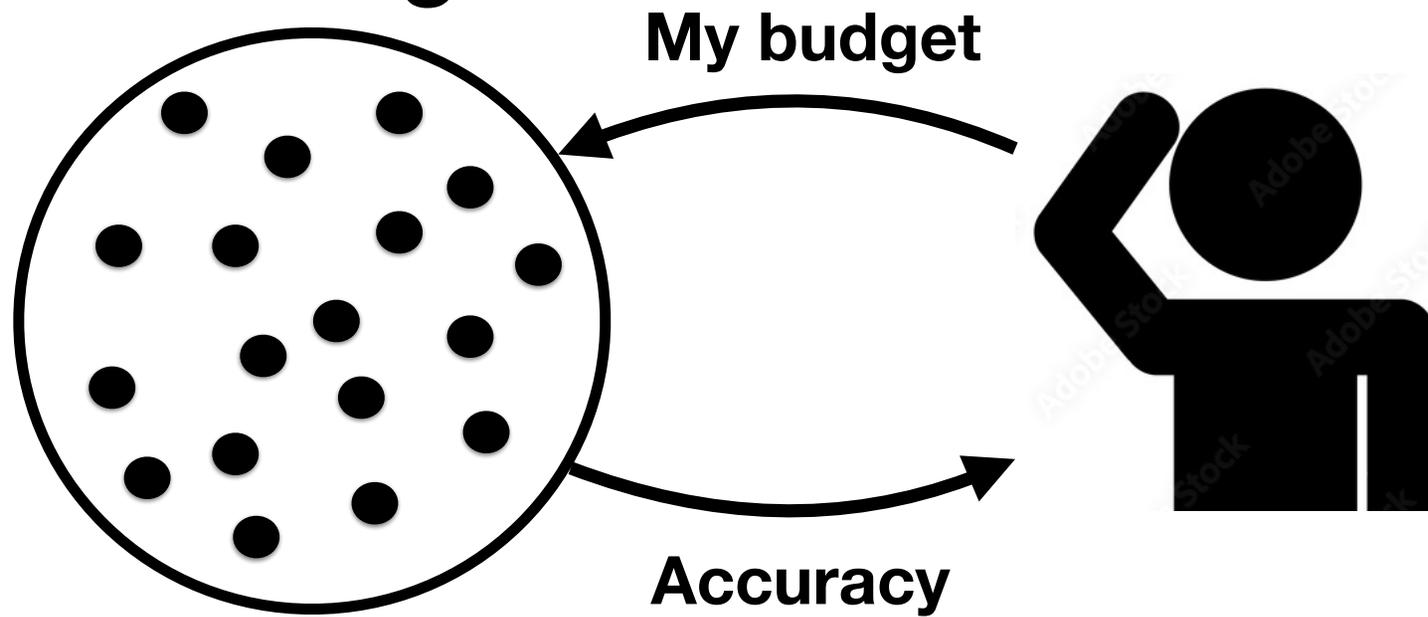
Design space



Search



4104 new designs



Machine: Nvidia A100 w/ 4 GPUs

Datasets:

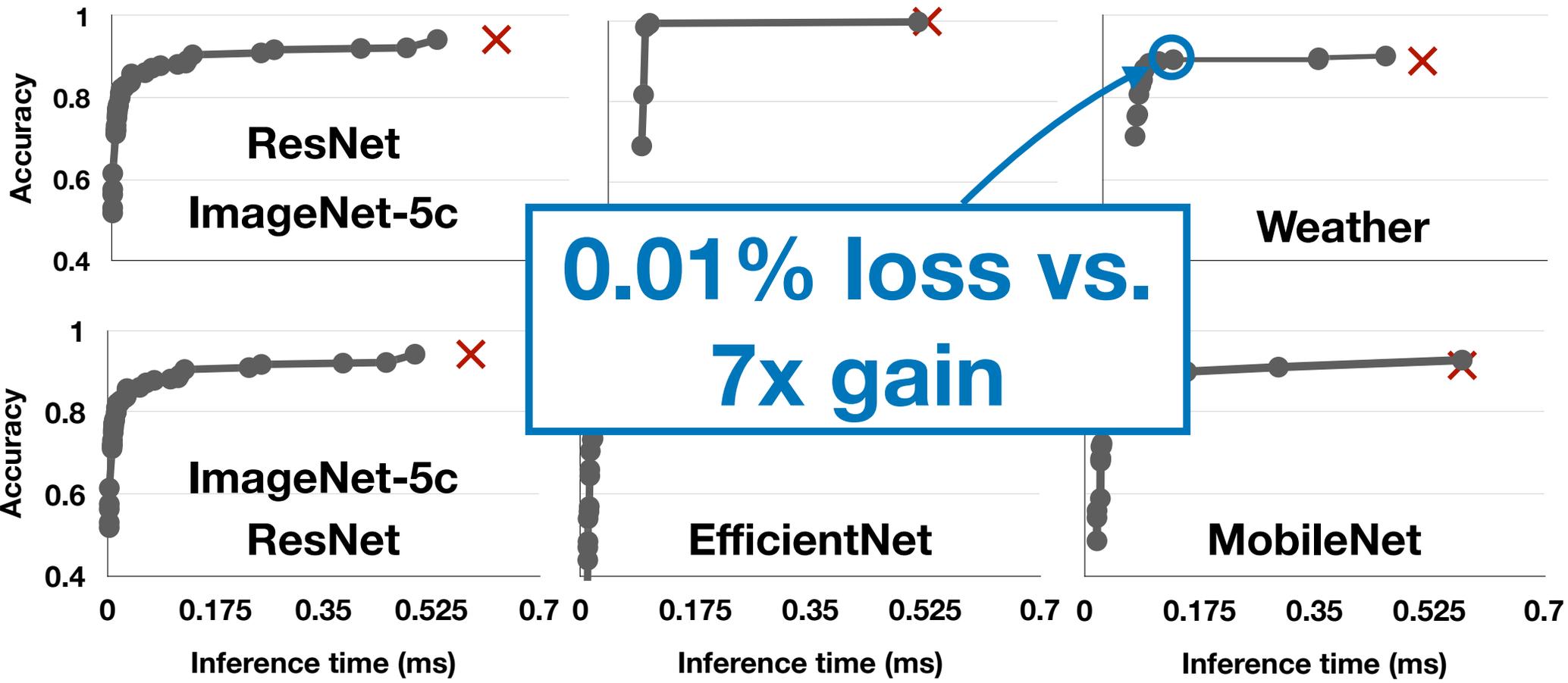
- **5-class ImageNet**
- **Blood-cell images**
- **Weather prediction**
- **50-class + full ImageNet**

AI Models: ResNet, EfficientNet, MobileNet

Hyper-parameters: Standard from PyTorch

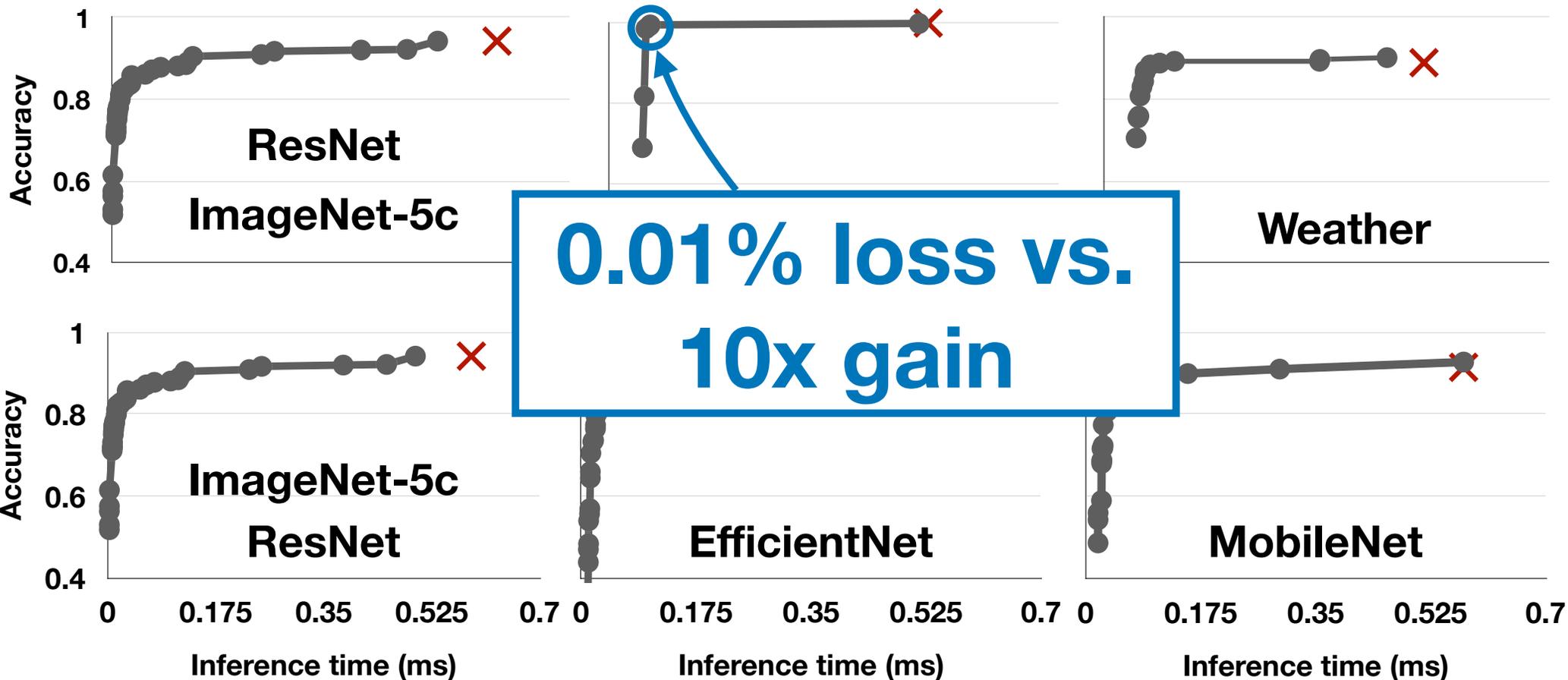
Baseline: JPEG

● Image Calculator ✕ JPEG

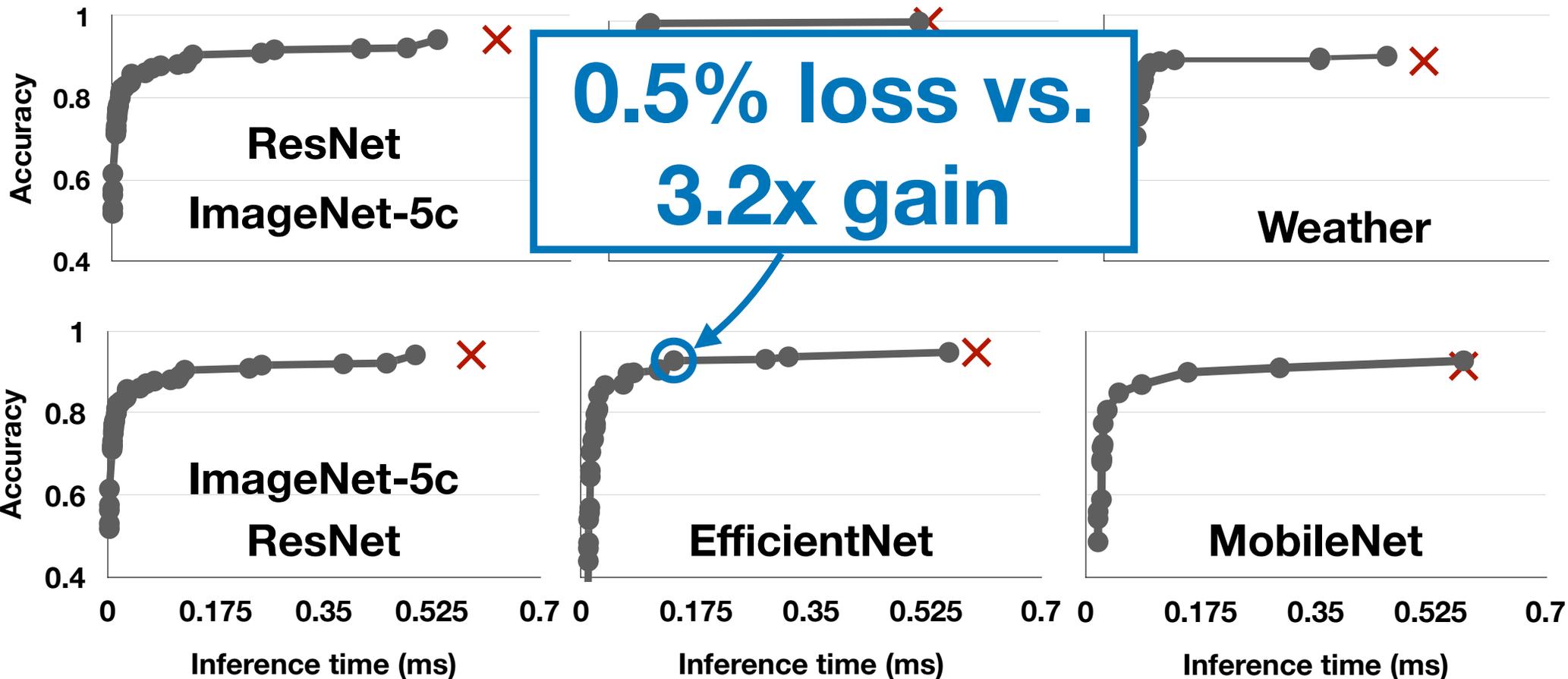


**0.01% loss vs.
7x gain**

● Image Calculator ✕ JPEG

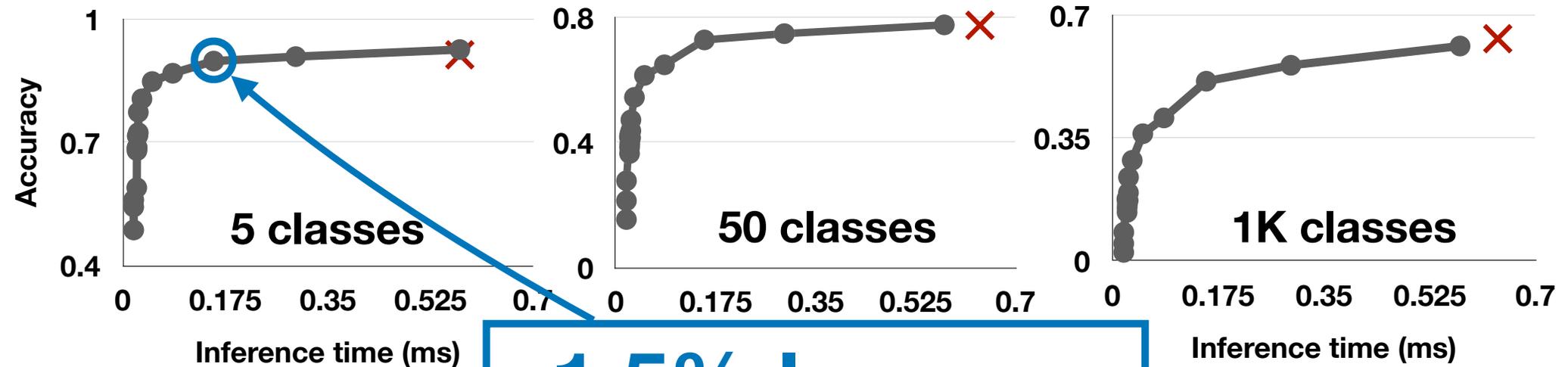


● Image Calculator ✕ JPEG



AI Model: MobileNet

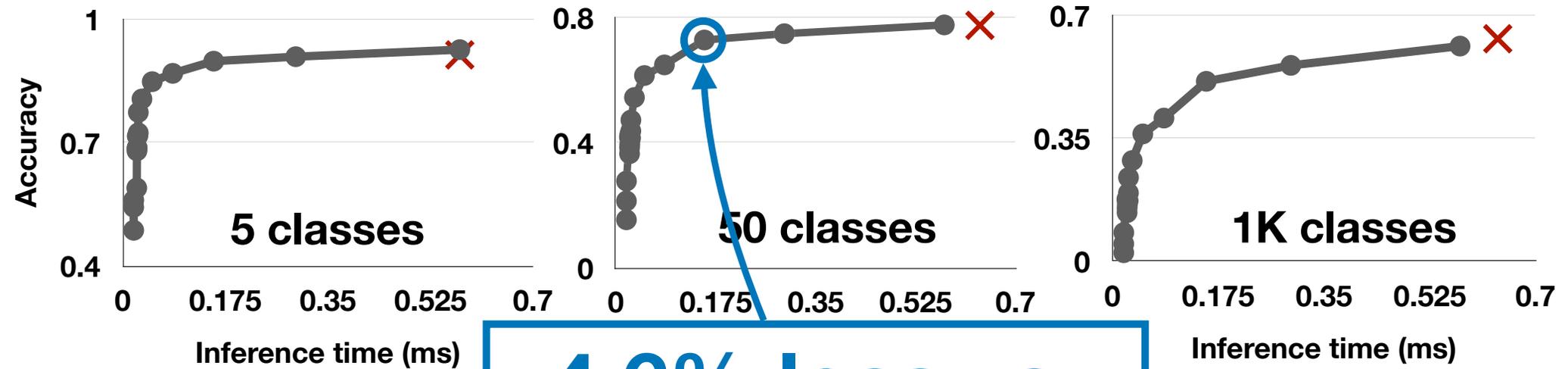
● Image Calculator ✕ JPEG



**1.5% loss vs.
3.7x gain**

AI Model: MobileNet

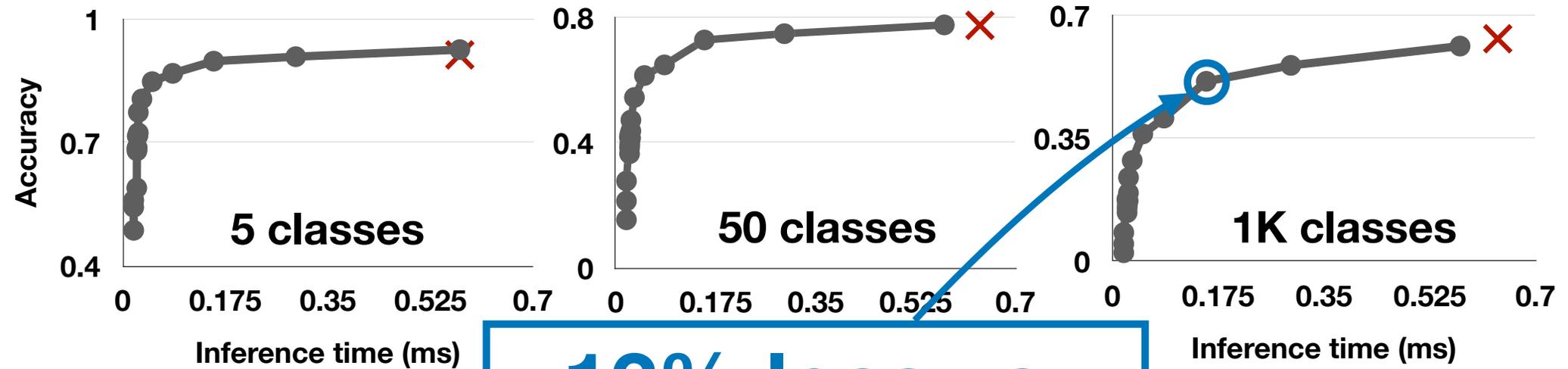
● Image Calculator ✕ JPEG



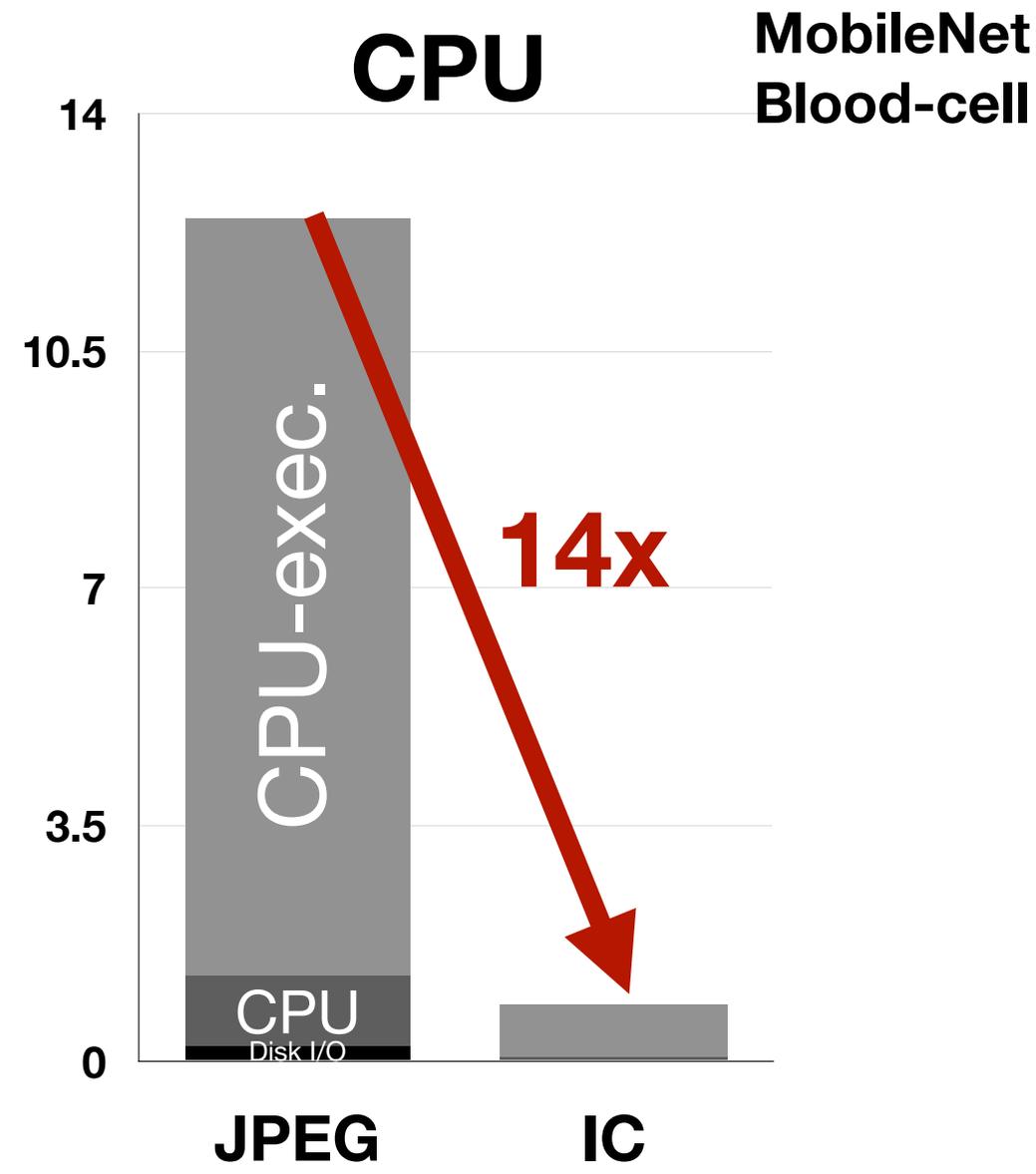
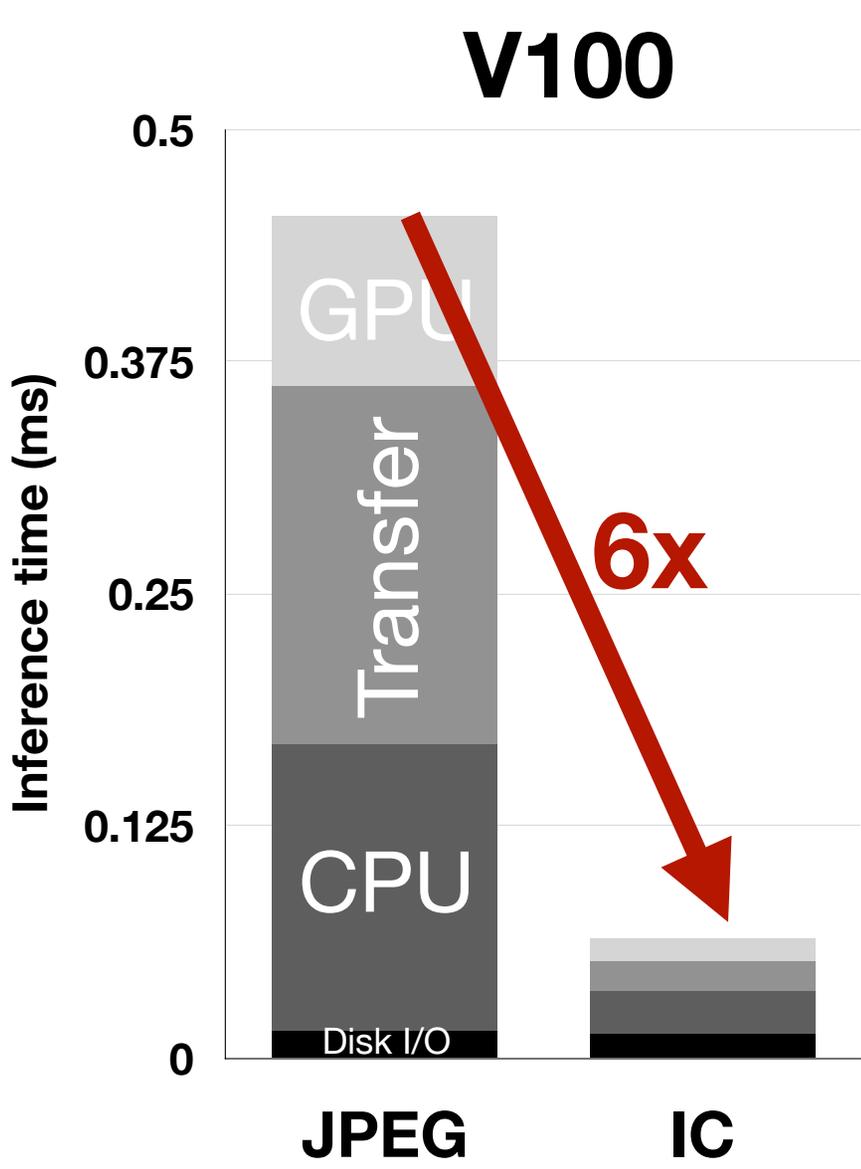
**4.6% loss vs.
3.7x gain**

AI Model: MobileNet

● Image Calculator ✕ JPEG

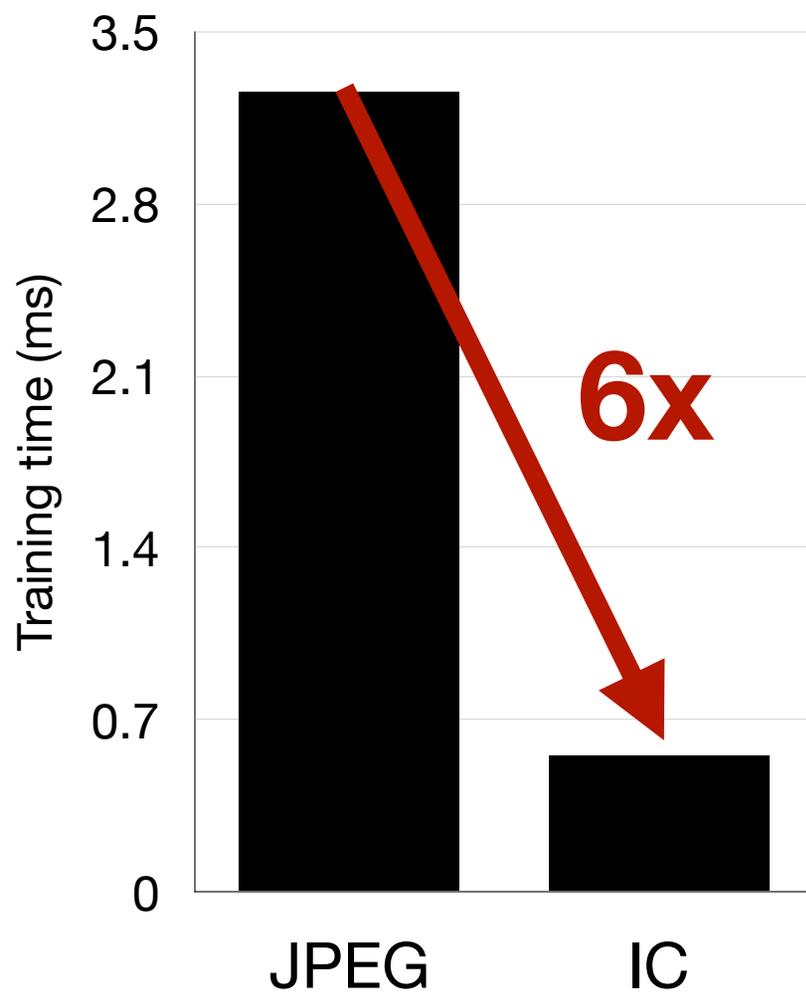


**12% loss vs.
3.7x gain**

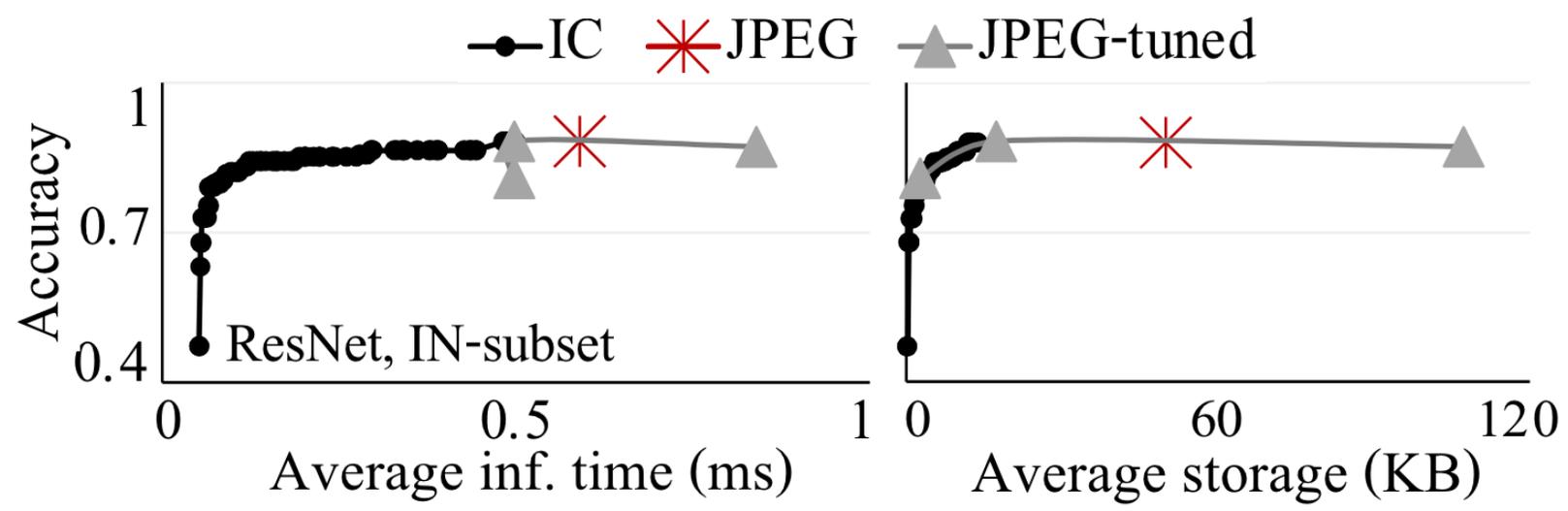


IC reduces training time

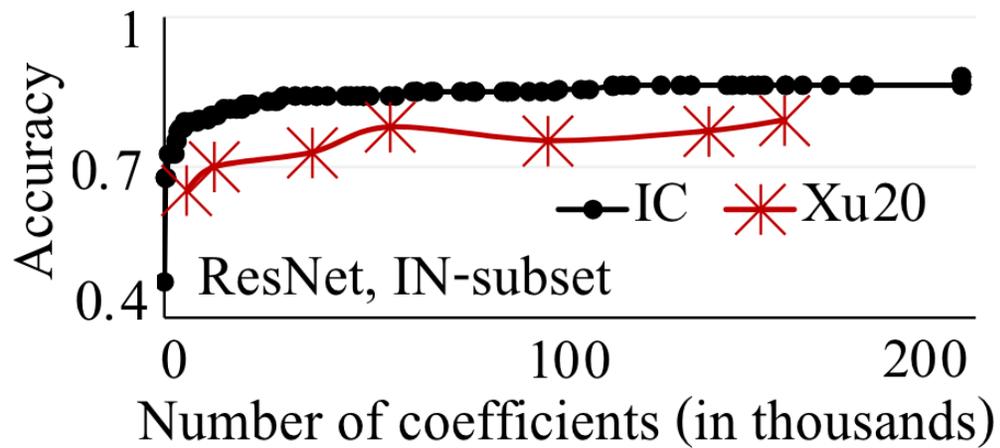
*Blood-cell, ResNet50
A100, PyTorch v1*



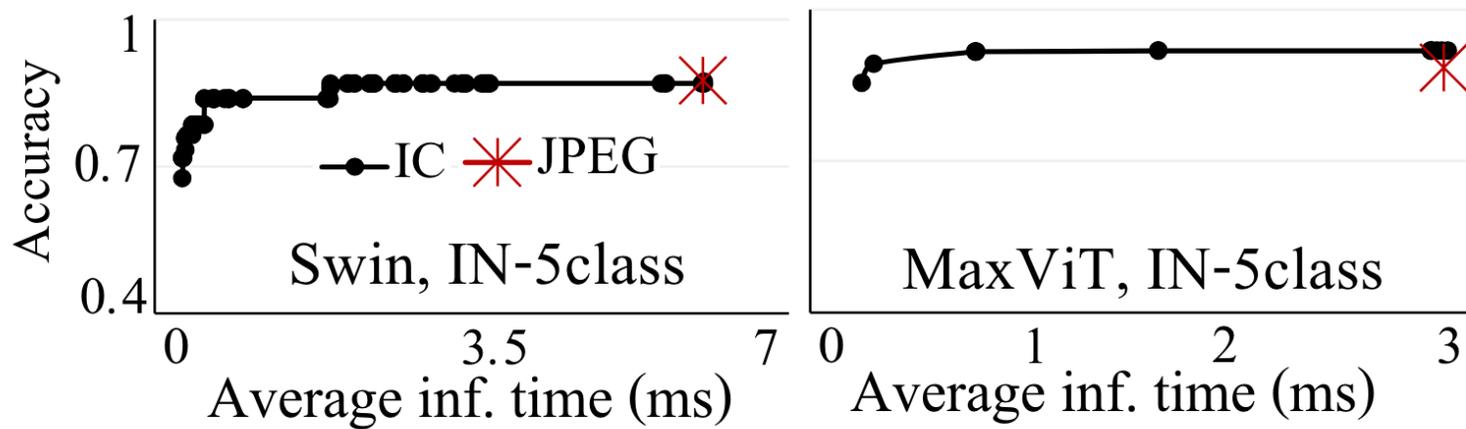
JPEG vs. JPEG-tuned



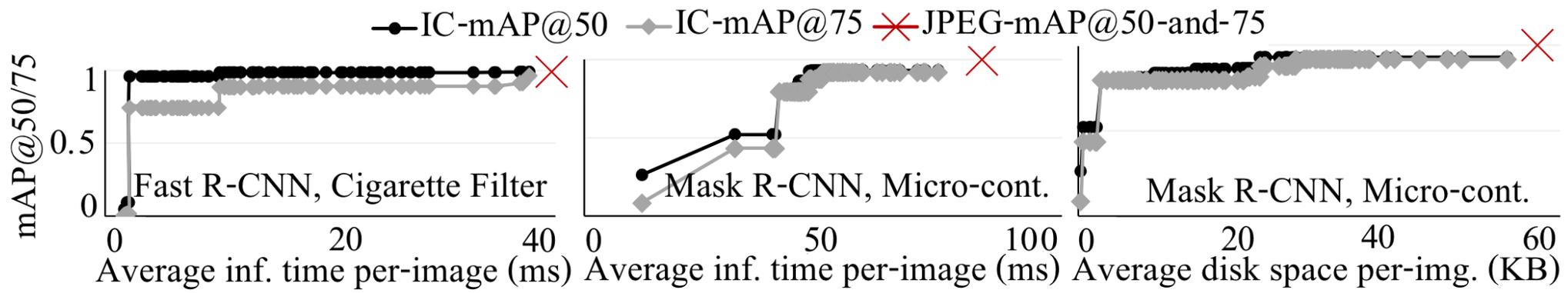
Related work, “Learned JPEG” (CVPR’20)



Visual Transformers



Object detection



Latency-accuracy

Latency-accuracy

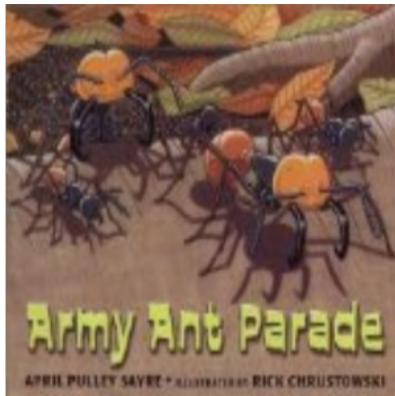
Storage-accuracy

Object detection & instance segmentation — examples

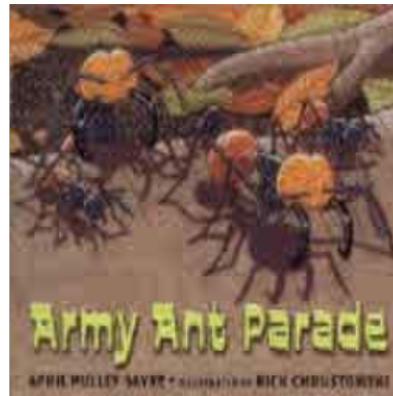


IC maintains high visual quality

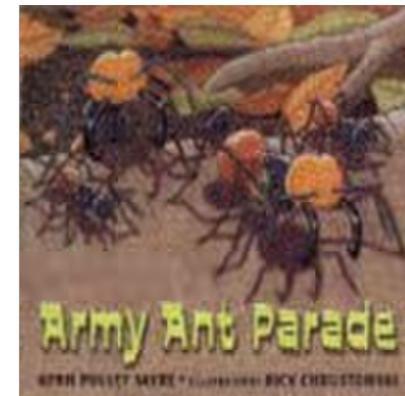
ImageNet-5c
PyTorch v1



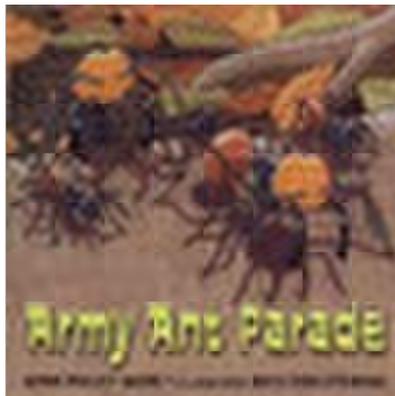
Original



60x



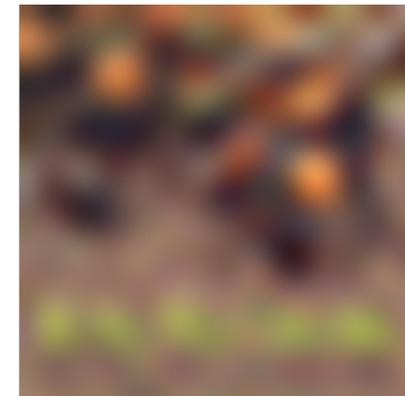
100x



200x

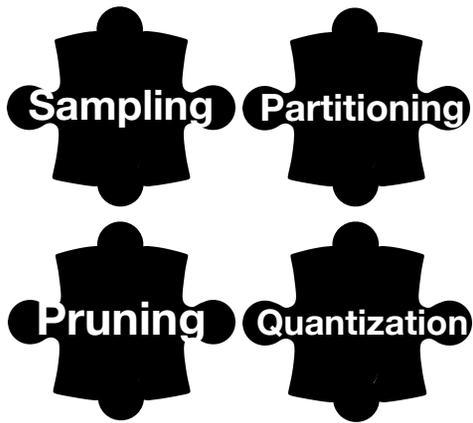


500x

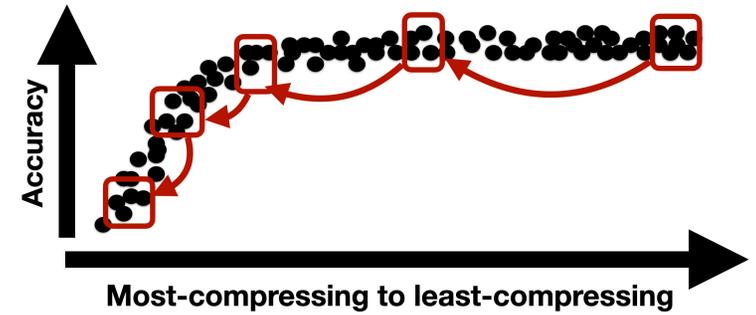
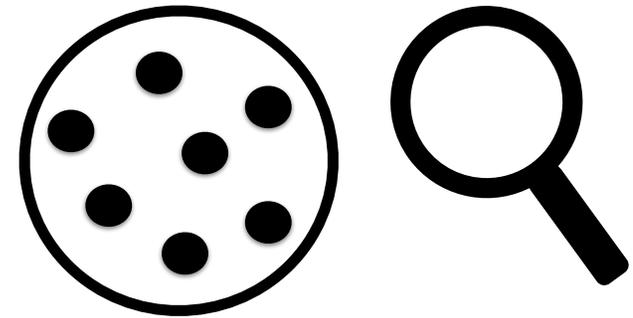


1200x

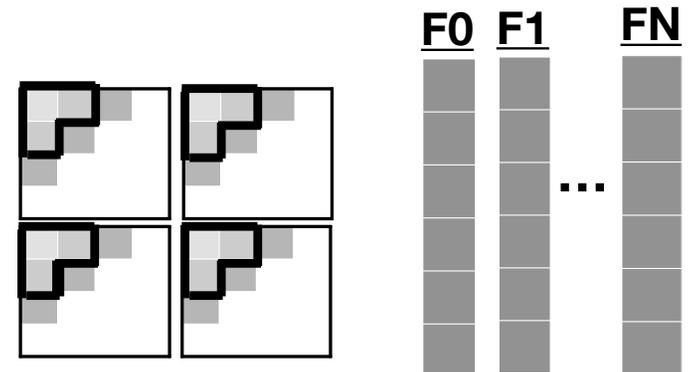
Self-Designing AI Storage



10^{150K} \rightarrow 4104



Efficiency vs. Scalability





DASlab

@ Harvard SEAS

daslab.seas.harvard.edu

THANKS!

