

Mechanistic Interpretability for Vision Models Optimization



SAPIENZA
UNIVERSITÀ DI ROMA

Course: Computer Vision
Candidate: Alessio Olivieri
ID: 1973323
Date: 19/06/25

Tutti i diritti relativi al presente materiale didattico ed al suo contenuto sono riservati a Sapienza e ai suoi autori (o docenti che lo hanno prodotto). È consentito l'uso personale dello stesso da parte dello studente a fini di studio. Ne è vietata nel modo più assoluto la diffusione, duplicazione, cessione, trasmissione, distribuzione a terzi o al pubblico pena le sanzioni applicabili per legge

Outline:

- Motivation and problem statement
- Related work and gaps
- Method: Mechanistic Pruning with ACDC
- Dataset and Task Definition
- Experimental Setup
- Results and Analysis
- Conclusion and Future Work

Problem Statement:

ViTs deliver top-tier accuracy but their fully-connected attention graphs contain millions of redundant edges and heads.

This makes them computationally expensive...
Especially for edge devices.

We need an automated way to discover and remove superfluous internal edges ahead of deployment, yielding a sparse, hardware-friendly ViT that preserves task accuracy.

State of the art:

- Static Circuit Discovery:

ACDC | *Edge-Pruning* (NeurIPS 24)

- Dynamic Token Pruning:

MADTP, *Zero-TPrune*

- Efficient Architectures:

SpectFormer-H-L, *CRATE- α -L*, *EfficientViT*

- Quantisation & Binarisation:

BHViT, *DeepCompress-ViT*, *Joint Prune-Quant* (CVPR 25)

- Parameter-Efficient FT:

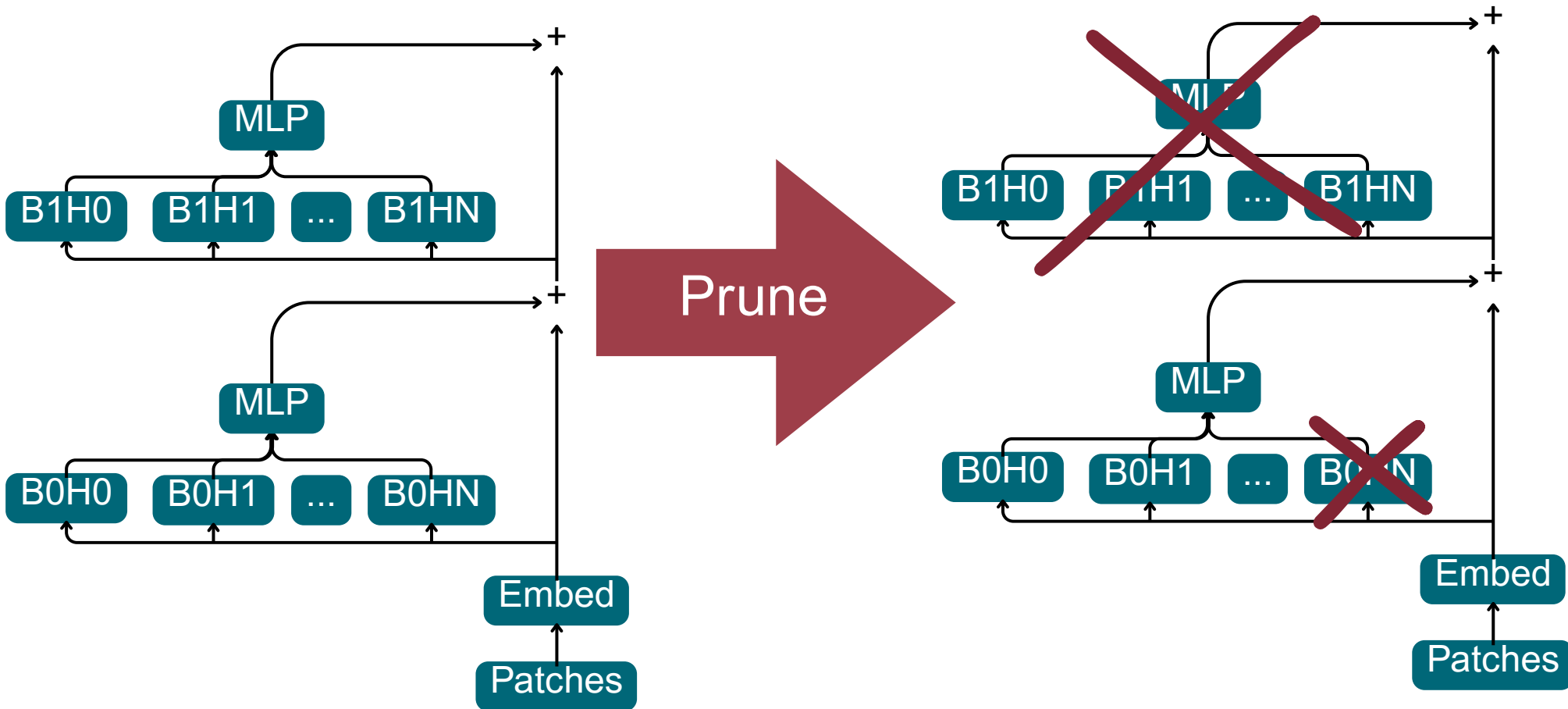
Serial-LoRA, *NOLA*

- Fast Attention Kernels:

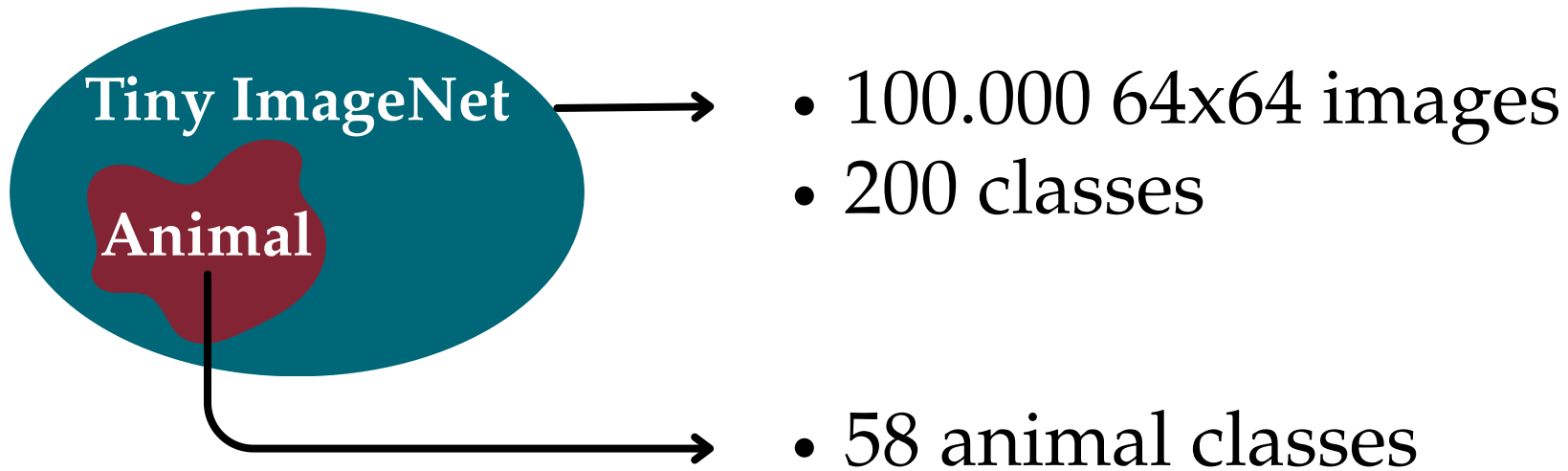
FlashAttention-3

Proposed method:

By pruning a ViT trained on a more complicated task
We can obtain a more efficient ViT



Dataset:



ACDC TASK: 6 coarse animal classes

Acquatic

Reptiles

Arthropods

Birds

Mammals

Marine Life

Experimental Setup (1)

Train the ViT to recognize the 58 fine classes

ViT architecture:

- Patch size: 8
- Hidden size: 64
- Blocks: 6
- Attention heads: 8

Data transformations:

- RandAugment
- Horizontal flip
- Random erasing: 25%

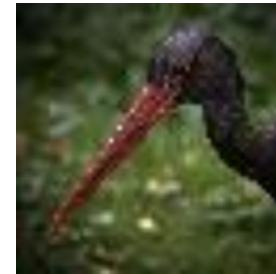
Trainer parameters:

- Dropout: 0.2
- Criterion: Soft CrossEntropy
- Optimizer: AdamW
- Mixup + CutMix
- LR: $3e-4$
- WarmUp: 20 epochs
- Cosine annealing
- Patience: 50 epochs

Experimental Setup (2)

ACDC (1):

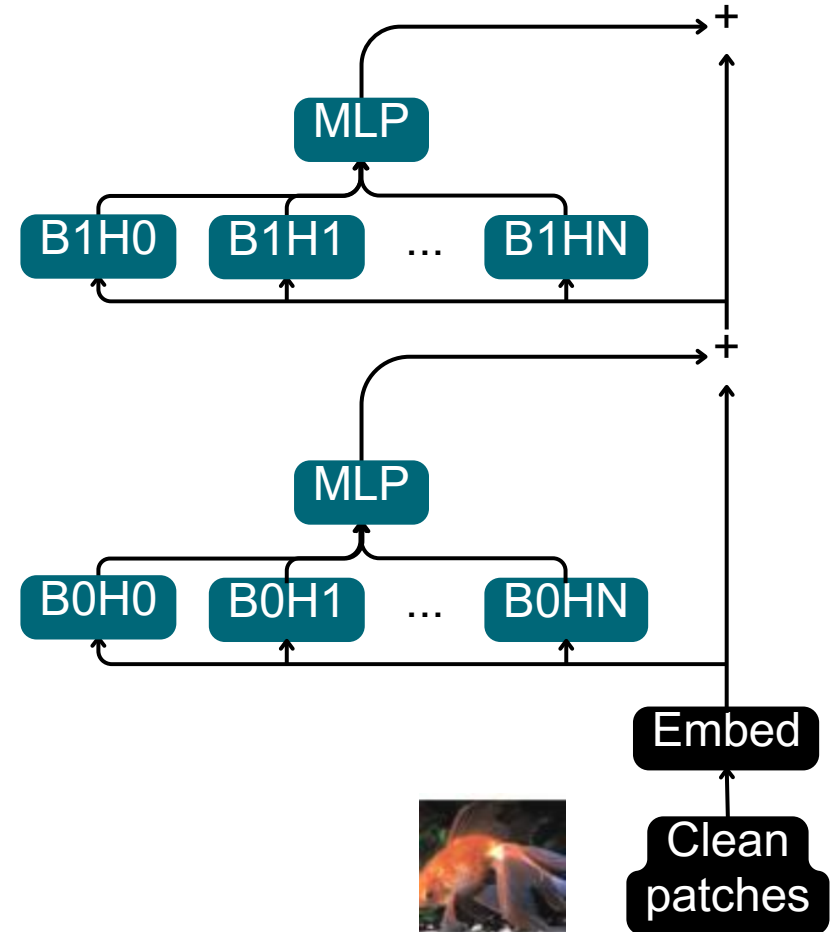
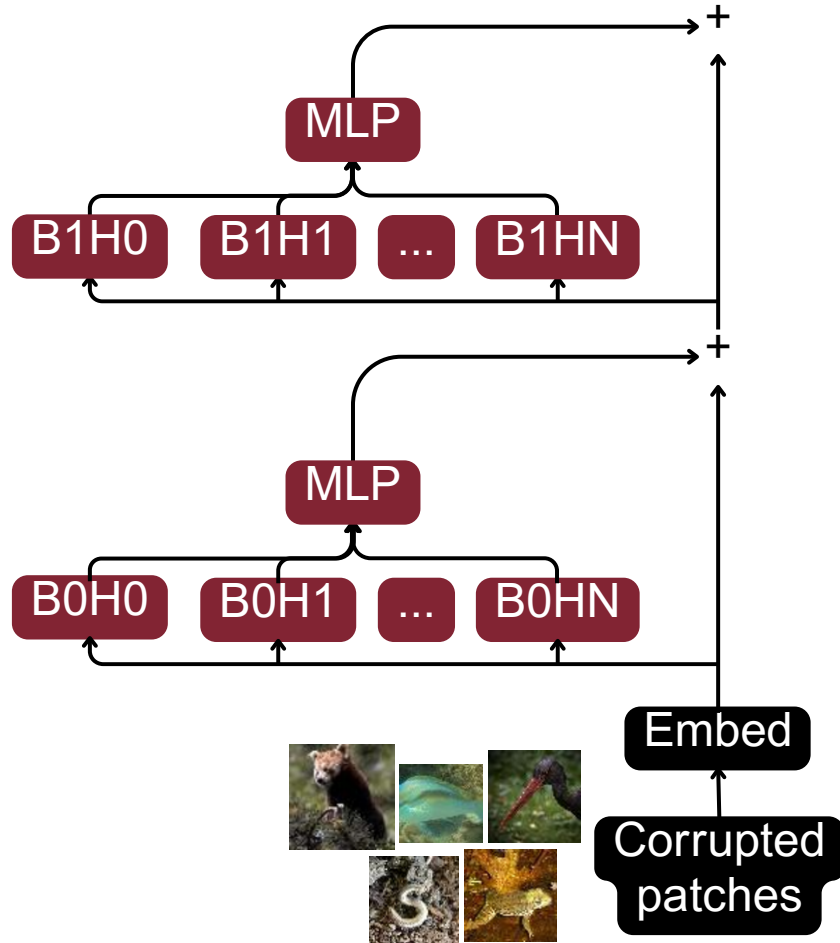
For each sample Get one "bad" sample (negative) from each of the other classes



Experimental Setup (2)

ACDC (2):

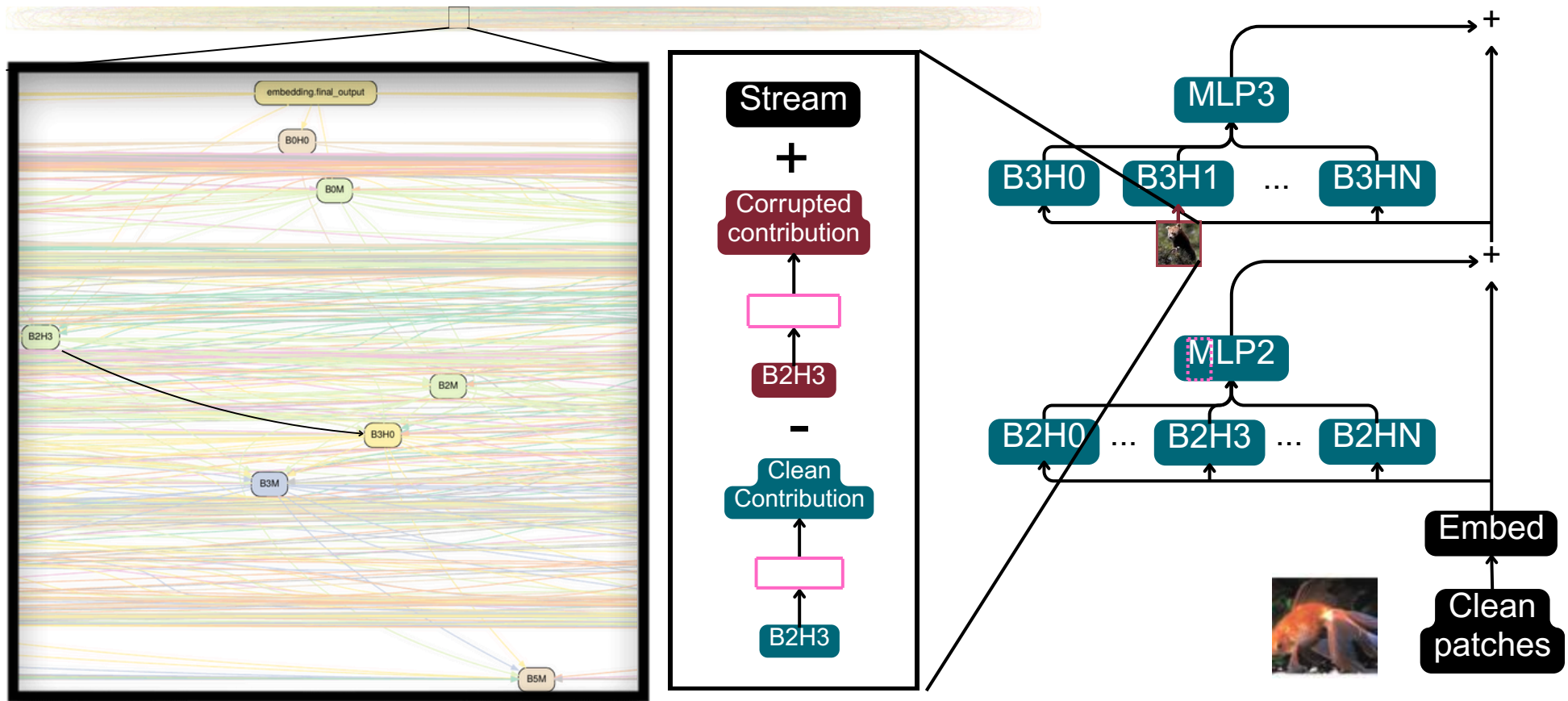
Cache all activations



Experimental Setup (2)

ACDC (3):

- Test each edge changing the input stream in input to the dest node
- If the average KL-Divergence over train dataset between clean logits is less than τ , Prune the edge

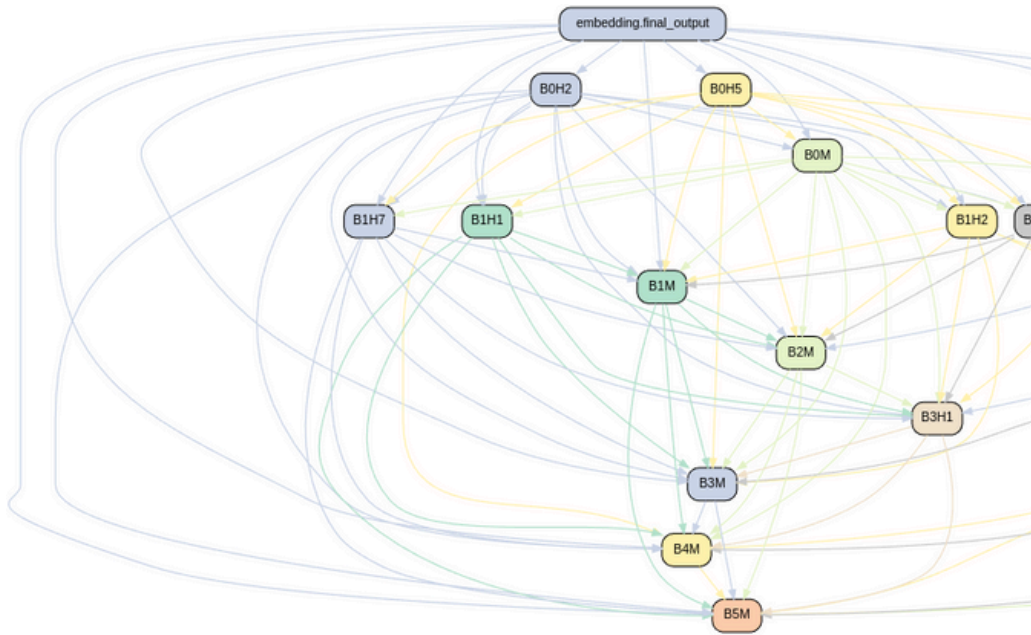
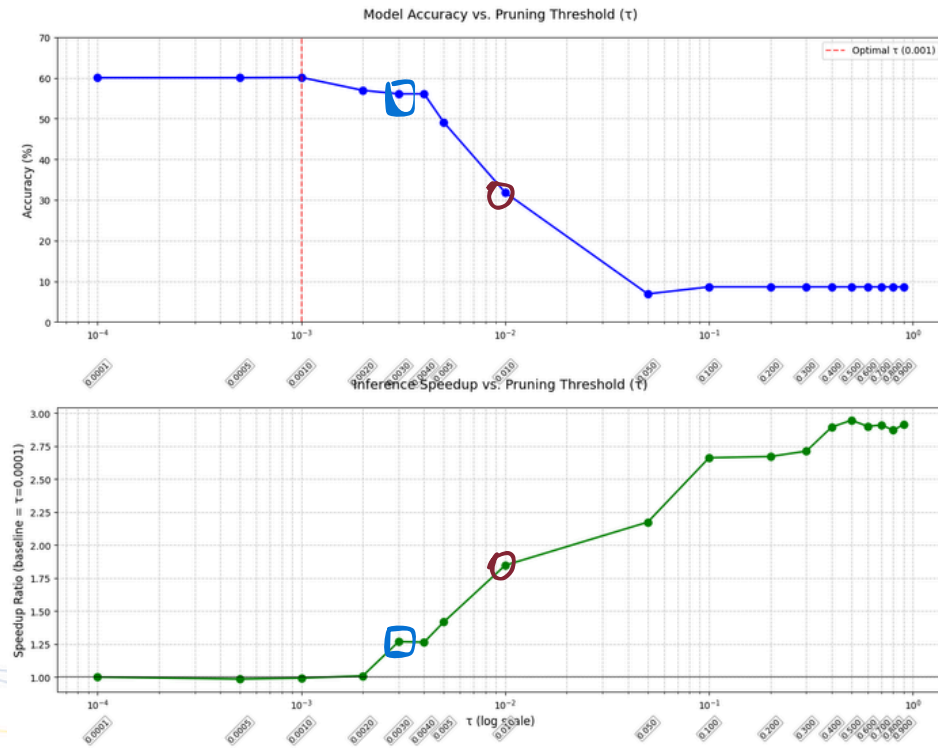


Model Evaluation (1)

Inference time and accuracy over τ

$\tau = 0.01$:

- ACDC pruned 94% of edges and 41 unused attention heads
- 50% accuracy decrease
- Roughly 2 times faster on CPU

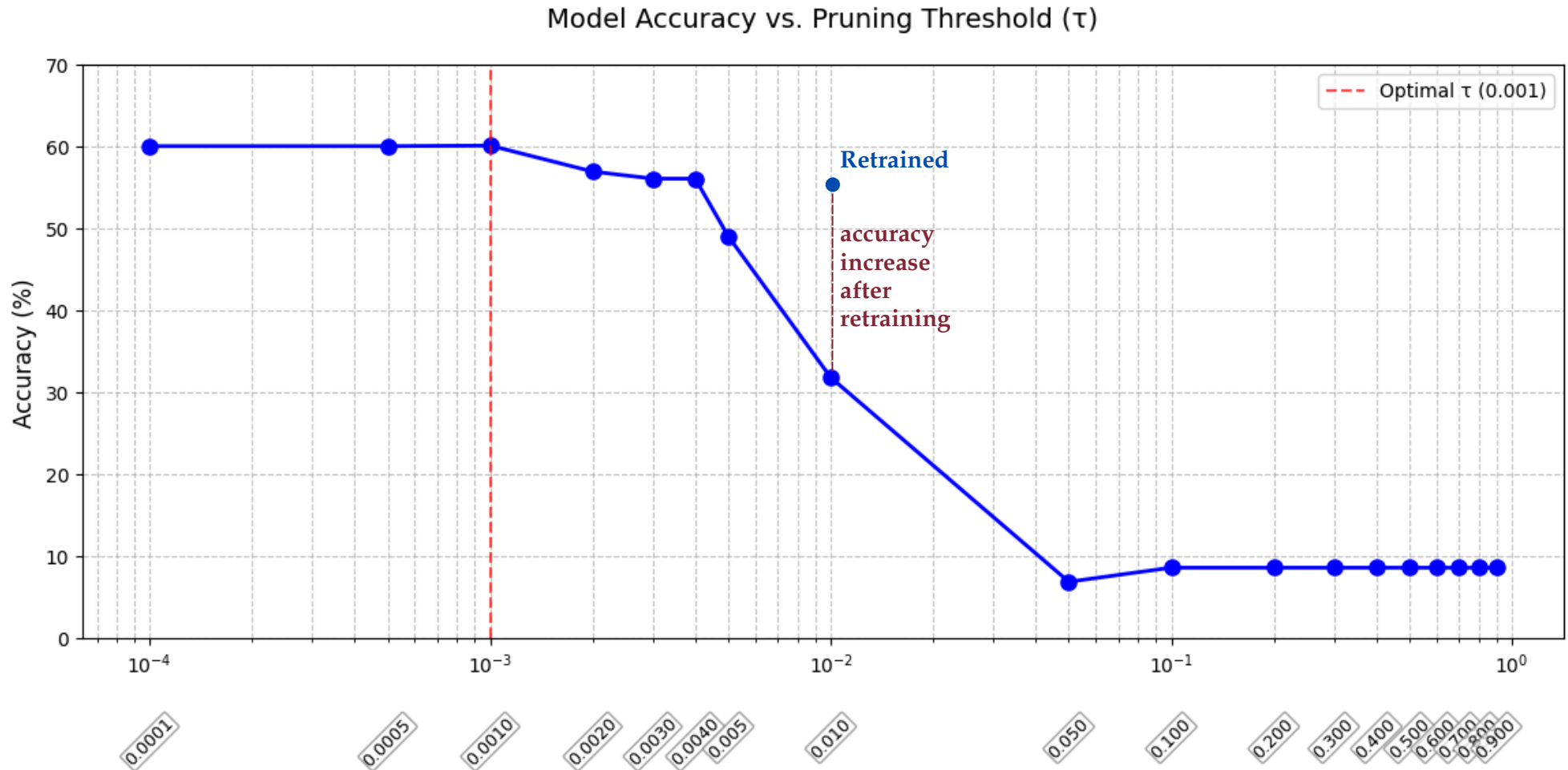


$\tau = 0.003$:

- ACDC pruned 57% of edges
- pruned 18 unused heads
- 9% accuracy decrease
- 25% speedup

Model Evaluation (2)

Accuracy on coarse classes after re initialization and training of the pruned ViT on the 58 fine labels



Conclusions:

- Mechanistic interpretability allows to squeeze performances from existing models
- 2 times speed increase and only 9% decrease of performances on CPU

Future work

- Create a prunable ViT architecture that scales well on GPU
- Benchmark performances on GPU



References

Miller et al. (2024). **Transformer Circuit Evaluation Metrics Are Not Robust.**

Conference on Language Modeling.

Conmy et al. (2023). **Towards Automated Circuit Discovery for Mechanistic Interpretability.**

NeurIPS. arXiv:2304.14997 [cs.LG].

Elhage, et al. (2021) "A Mathematical Framework for Transformer Circuits", Transformer Circuits Thread.

Matt Nguyen. (2024). "Building a Vision Transformer Model from Scratch." Medium.

Yuechun (Ethan) Gu (2024) "Fine-Tuning Vision Transformer (ViT) on Tiny ImageNet Dataset"