

Resource Efficient Artificial Intelligence

Dr Varun Ojha

School of Computing, Newcastle University

varun.ojha@newcastle.ac.uk

ojhavk.github.io



Newcastle
University



**National
Edge AI
Hub**

Resource Efficient Artificial Intelligence

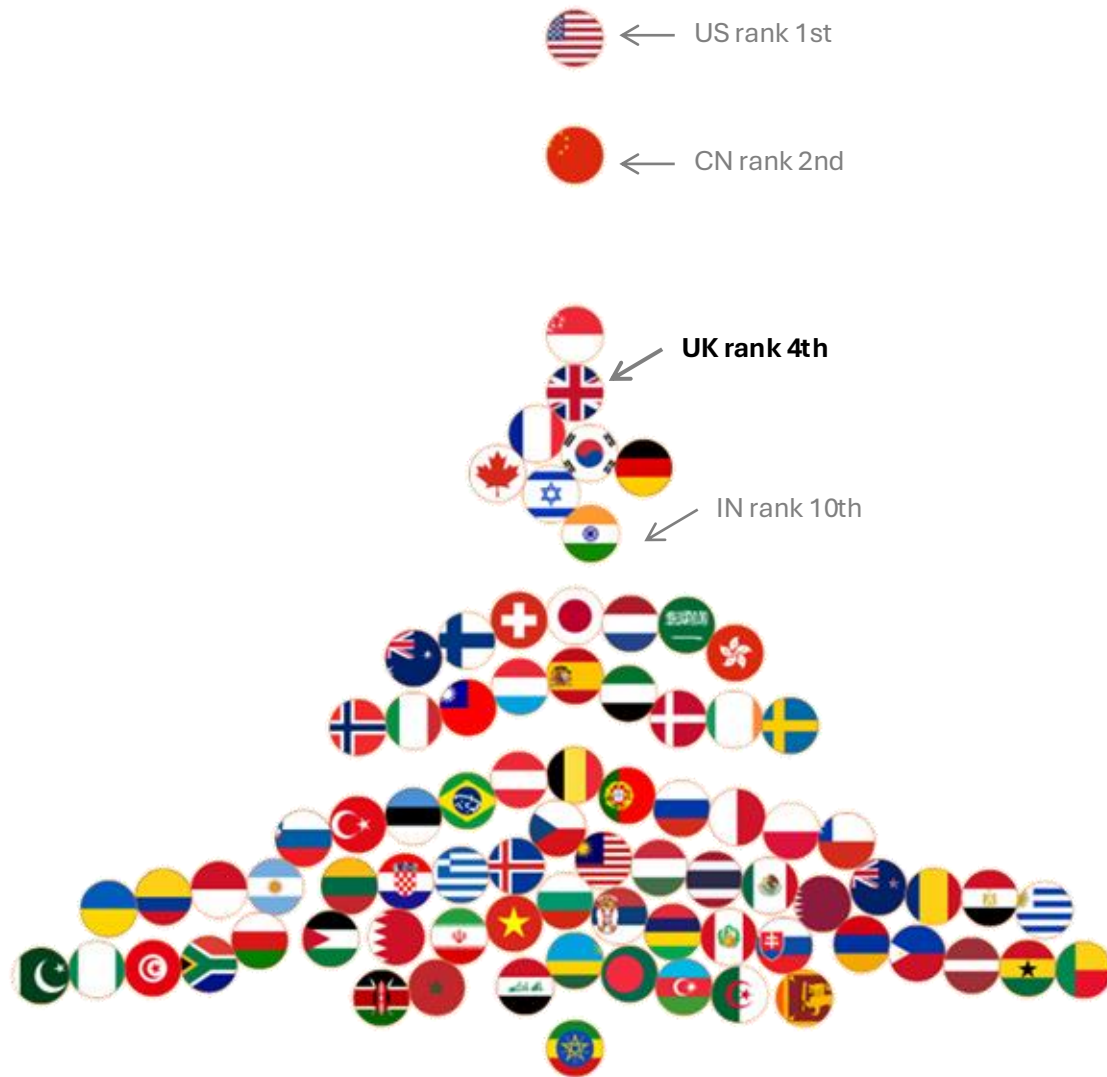
Agenda

- Why Resource Efficient AI
- Understanding AI Model Performance
- AI Model Compression Methods and Results
- Search for Efficient AI Models
- Tricks to make AI models resource efficient

Part 1

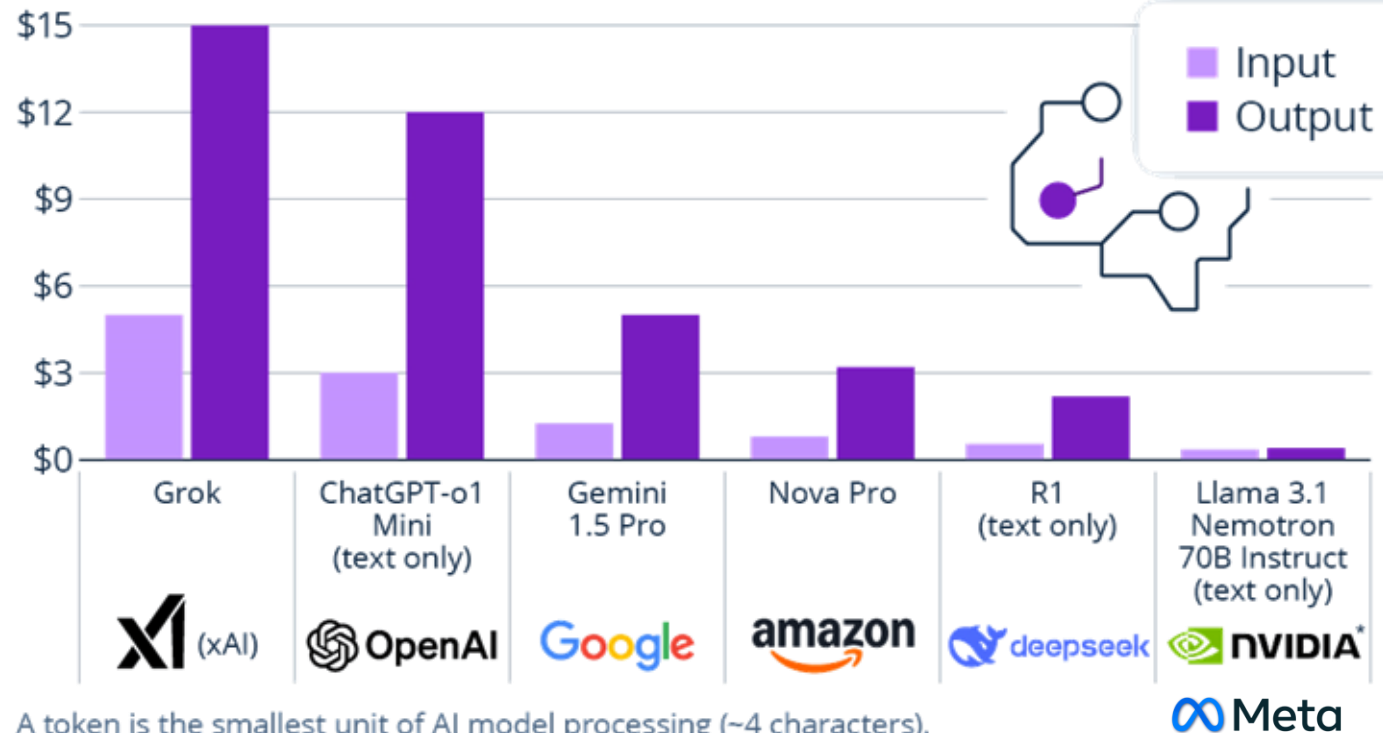
Why Resource Efficient AI

AI Arms Race



DeepSeek-R1 Upsets AI Market With Low Prices

Estimated price for processing one million input/output tokens on different AI models



A token is the smallest unit of AI model processing (~4 characters).
o1 is ChatGPT's latest model. List includes most comparable model per company

* Uses Meta's open-source Llama AI

Source: DocsBot



Advanced AI Models

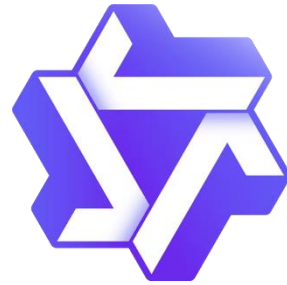
9:17 AM



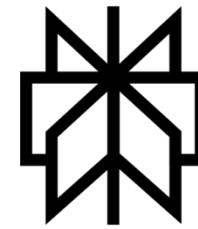
ChatGPT



Gemini



Qwen



perplexity



Mistral AI



Deepseek R1



Sora



midjourney

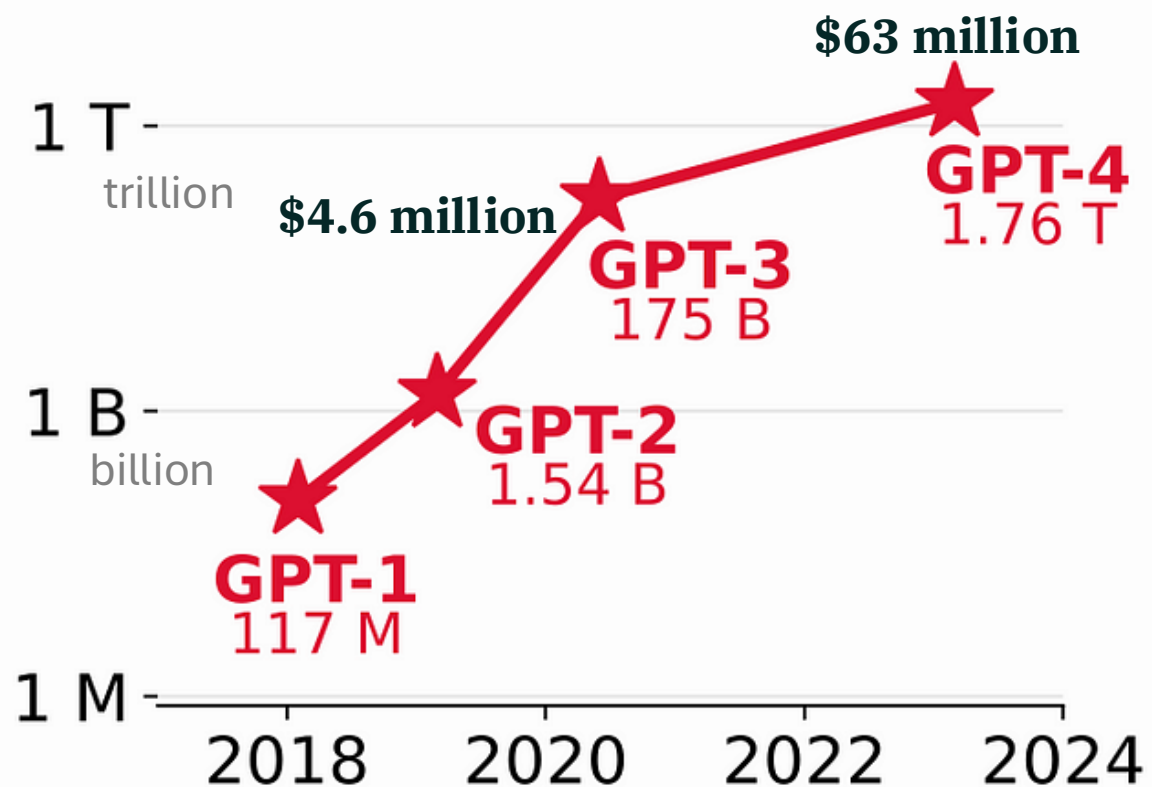


runway

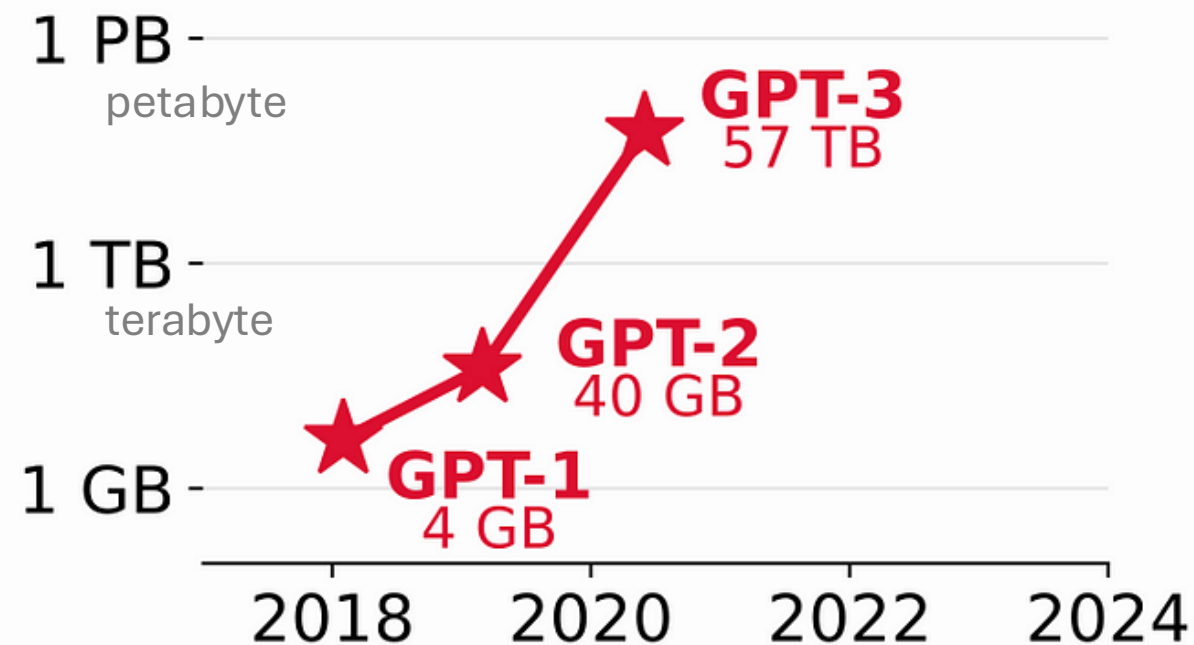
Google
Veo 3

Massive AI Models

Parameters Count



Training Data Size



AI Model Training Energy Consumption

Common carbon footprint benchmarks

in lbs of CO2 equivalent

Roundtrip flight b/w NY and SF (1 passenger)

1,984

Human life (avg. 1 year)

11,023

American life (avg. 1 year)

36,156

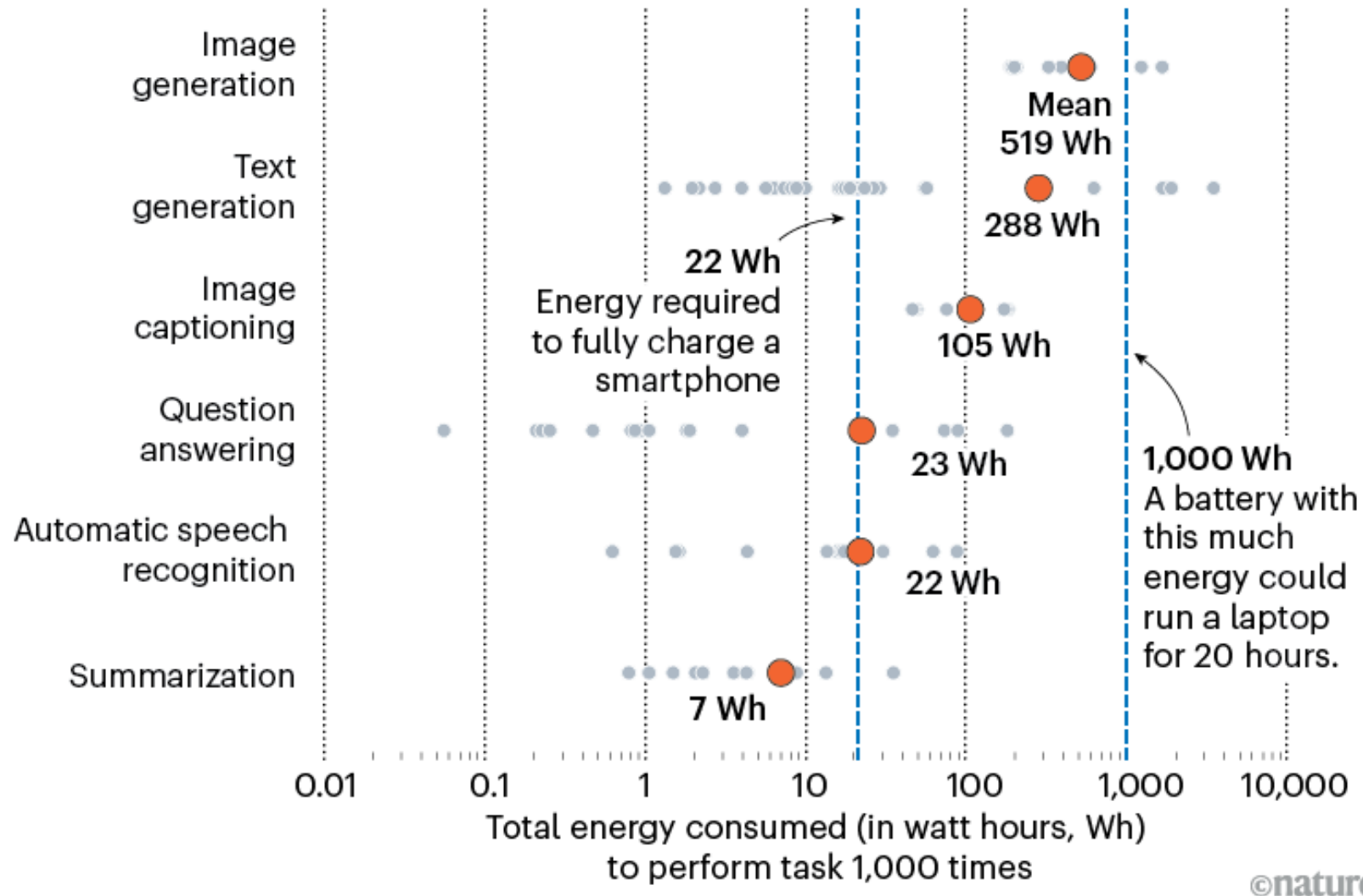
US car including fuel (avg. 1 lifetime)

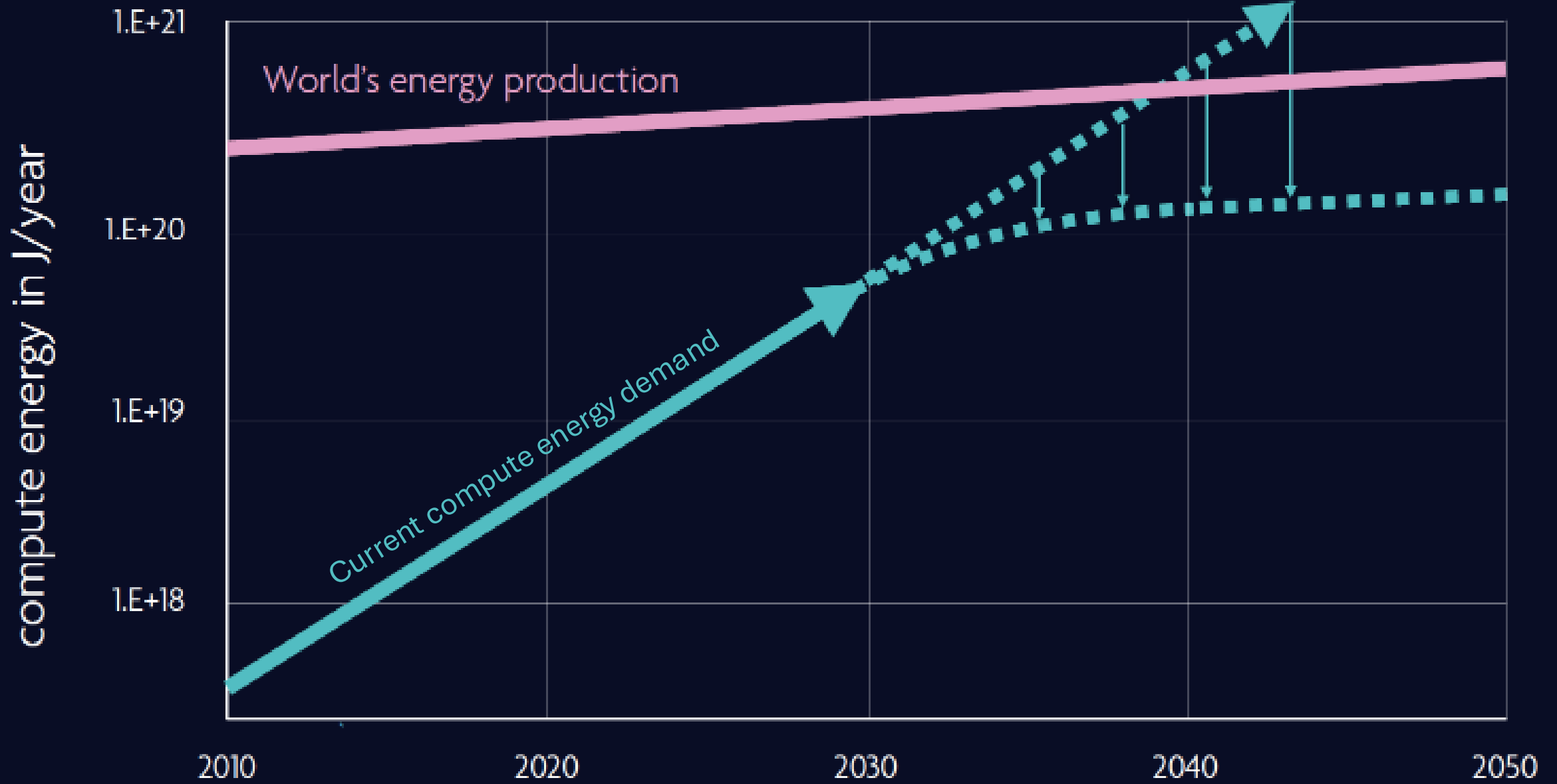
126,000

Transformer (213M parameters) w/
neural architecture search

626,155

AI Model Inference Energy Consumption

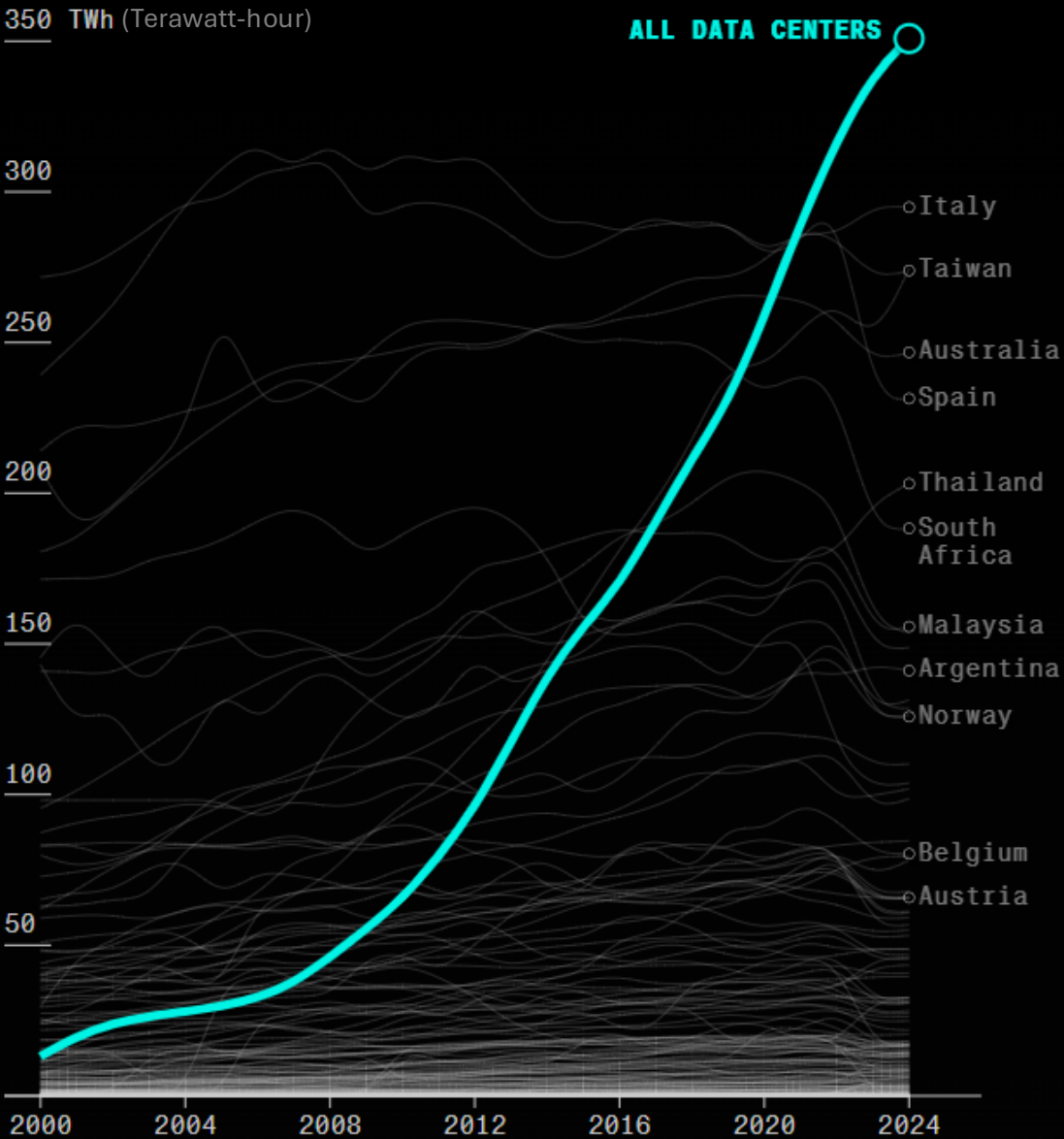




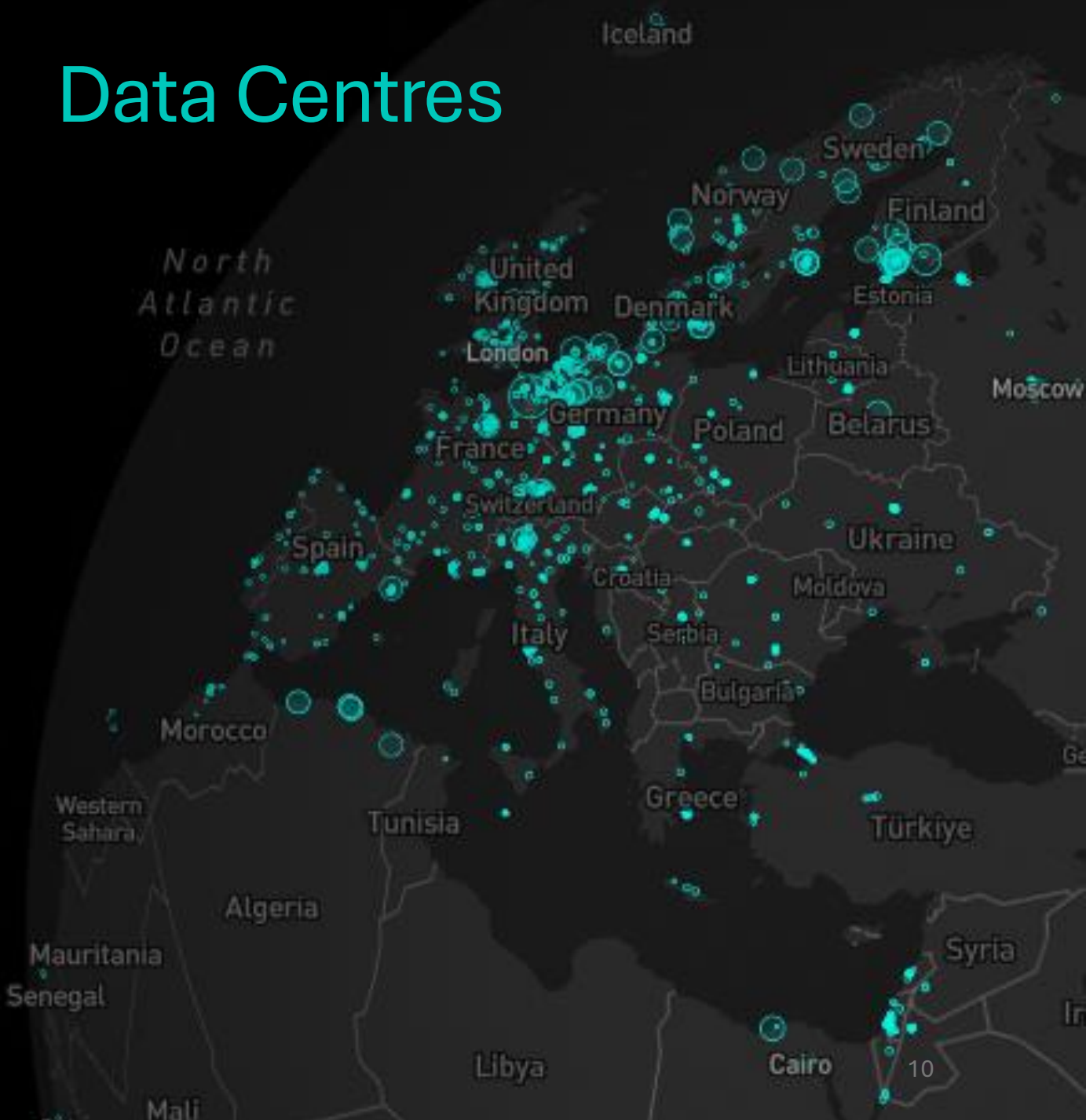
Altogether, data centers use more electricity than most countries

Only 16 nations, including the US and China, consume more

Source: Bloomberg



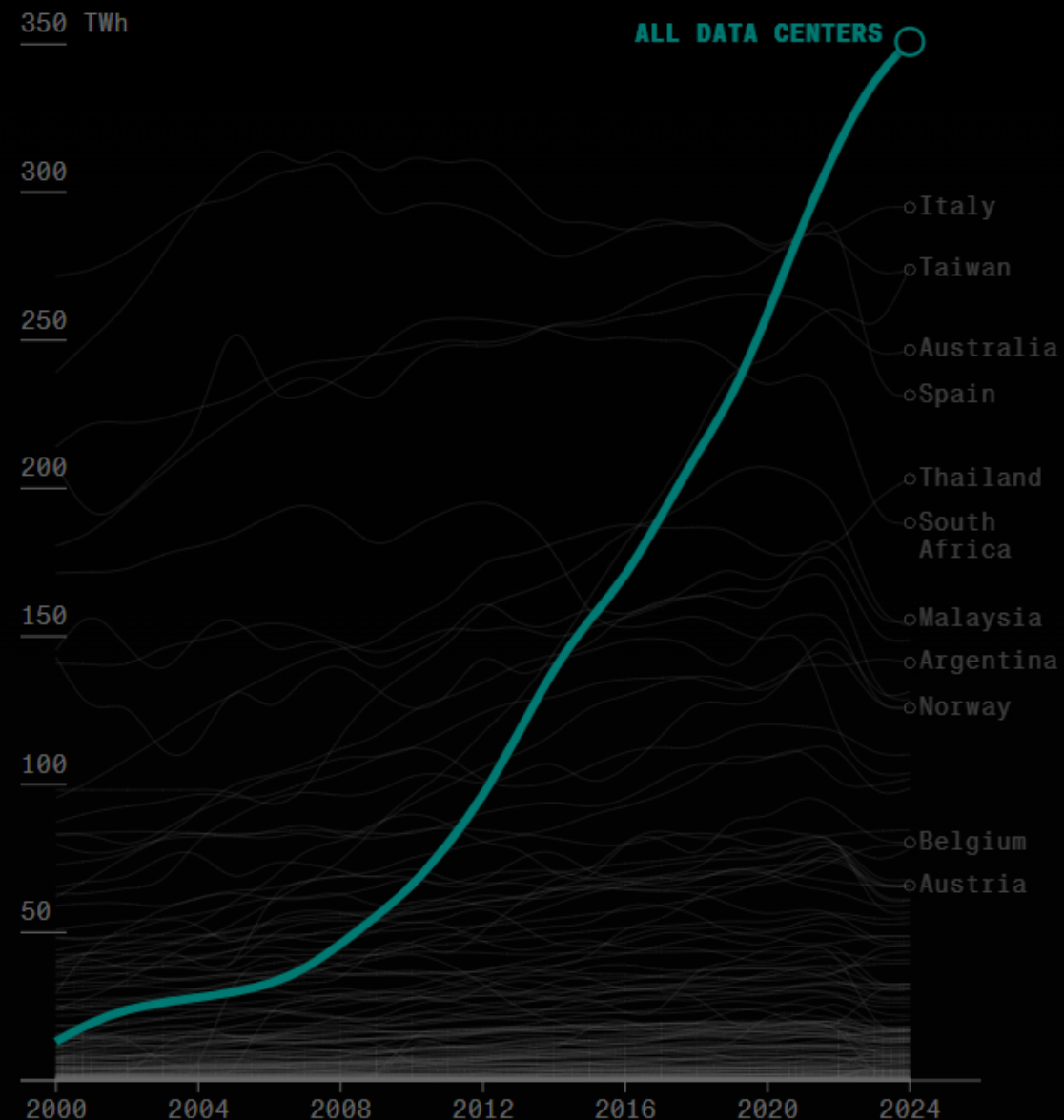
Data Centres



Altogether, data centers use more electricity than most countries

Only 16 nations, including the US and China, consume more

Only 16 nations, including the US and China, consume more



Potential solutions to energy problem:

Cloud AI → Edge AI

- Moving AI applications from the cloud to the edge
- AI model simplification techniques to reduce power consumption
- AI Model quantization, etc.
- On device (Edge) AI model training
- Federated Learning

Edge AI

- Unlike Cloud AI (e.g., ChatGPT that runs in data centers), edge AI runs at the edge computing devices such smartphones, cameras, cars, medical devices, ensuring quality of data for inference
- Reduces latency, cost, and power consumption
- Protects data privacy and reduce improve data security and cybersecurity
- Reduces risk of inference failure in critical systems (e.g., autonomous vehicles, healthcare devices) that may endanger lives

Everyone is a walking AI computer

Any random mobile configuration/specification these days

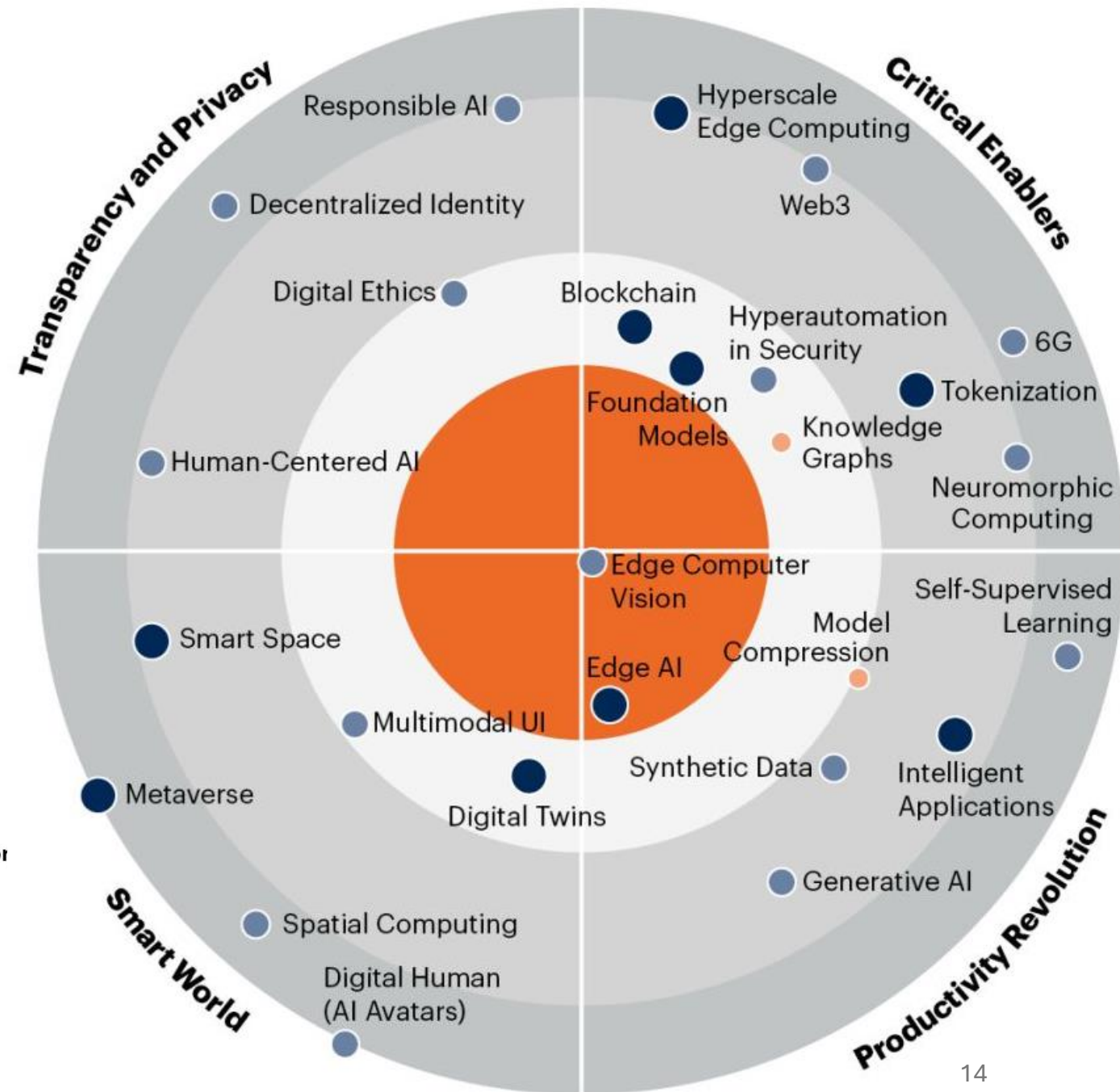
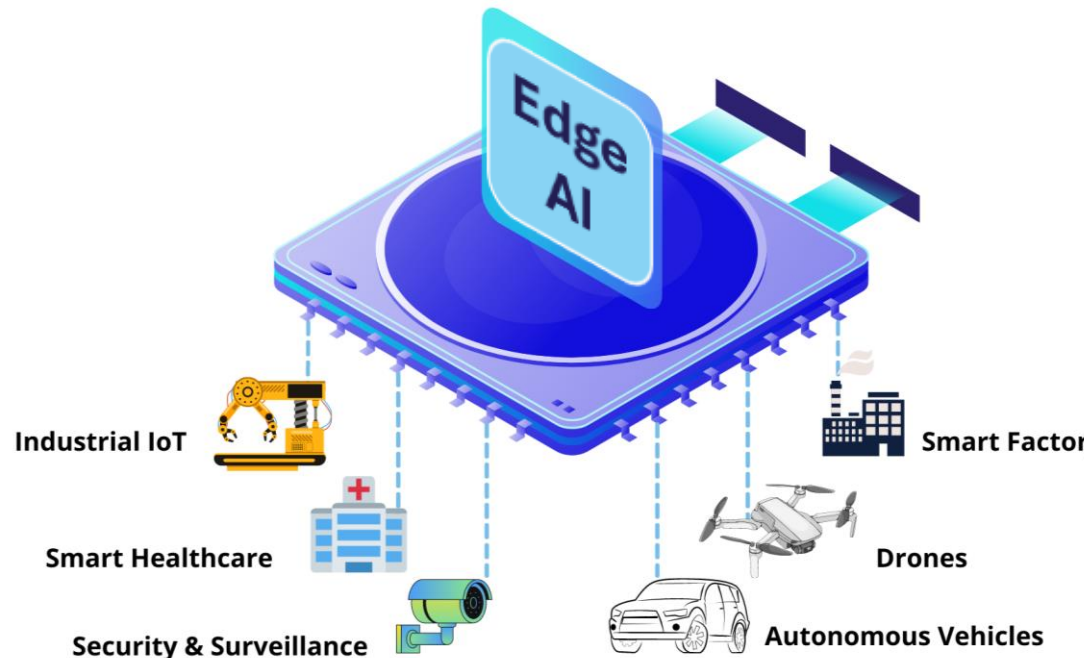


- **Chip**
- A18 Bionic chip
- 6-core CPU with 2 performance
- 4 efficiency cores
- 5-core GPU
- 16-core Neural Engine
- **Capacity**
- 512GB
- **Multiple Sensors**
- 48MP camera
- Satellite and GPS



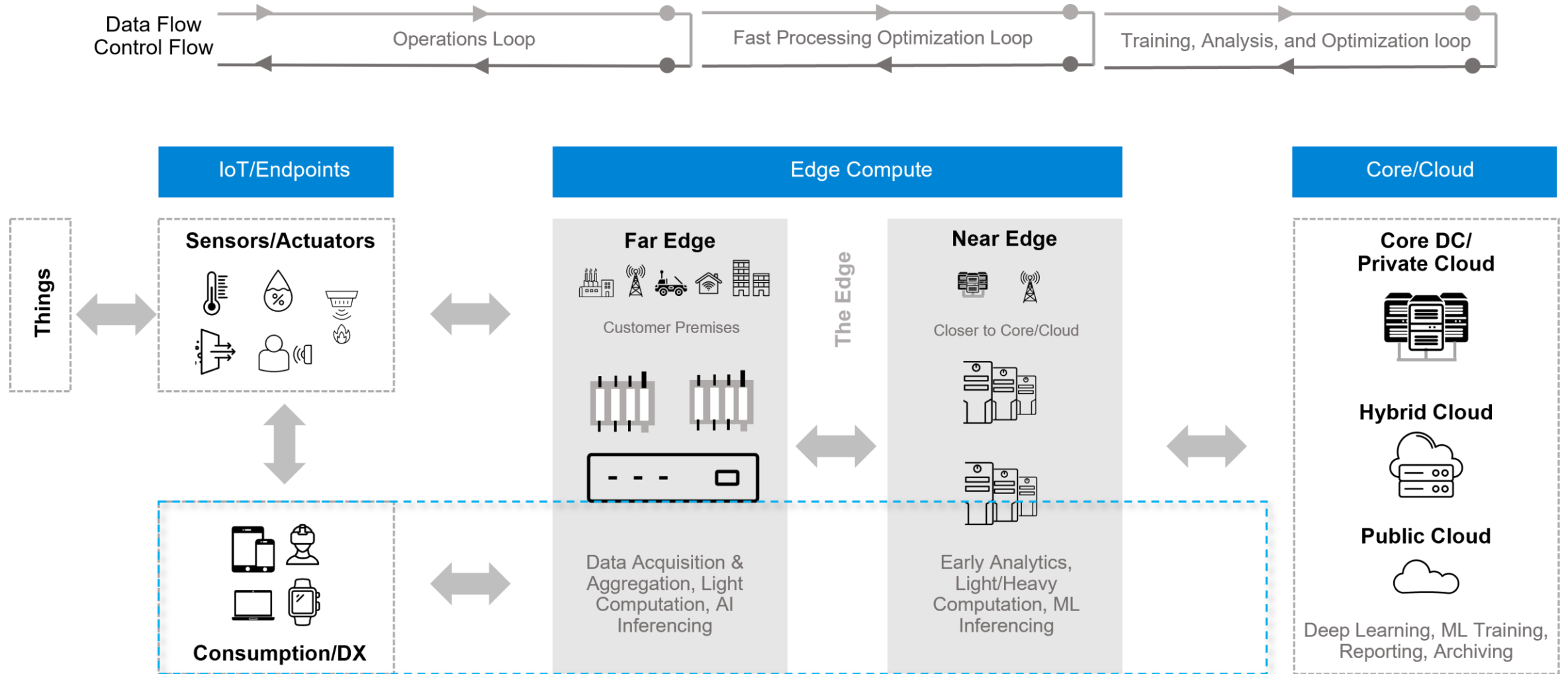
The Emergence of Edge AI: a game changer for industries

(Gartner 2023)



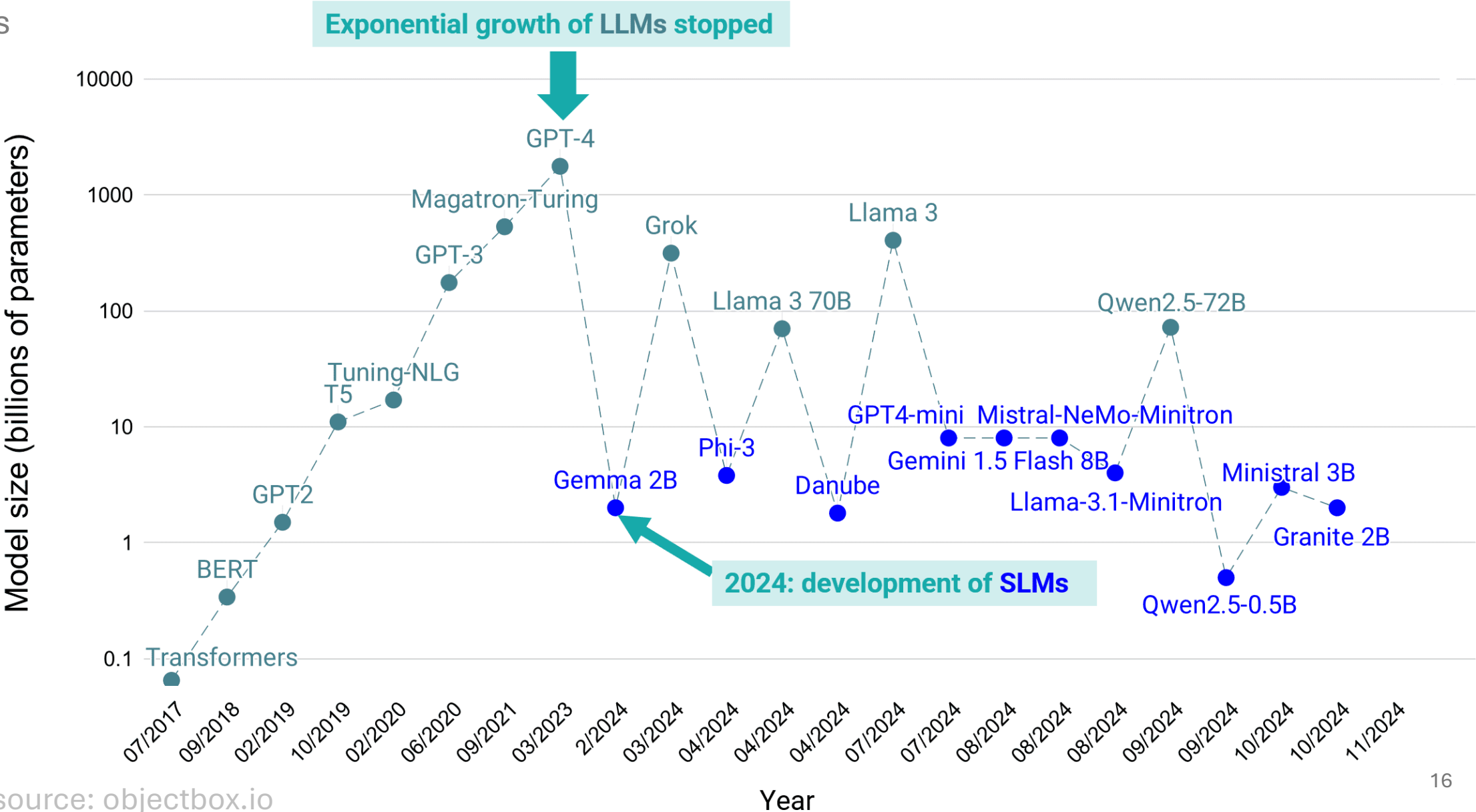
The Rise of Generative AI at the Edge

image source: Arrow Intelligent Solutions



Small Language Models Could Redefine The AI Race,

Forbes



AI Model Size



Performance Trade-offs

Larger models, such as GPT-3 with 175 billion parameters, can handle complex tasks but require significant computational resources and incur higher operational costs.



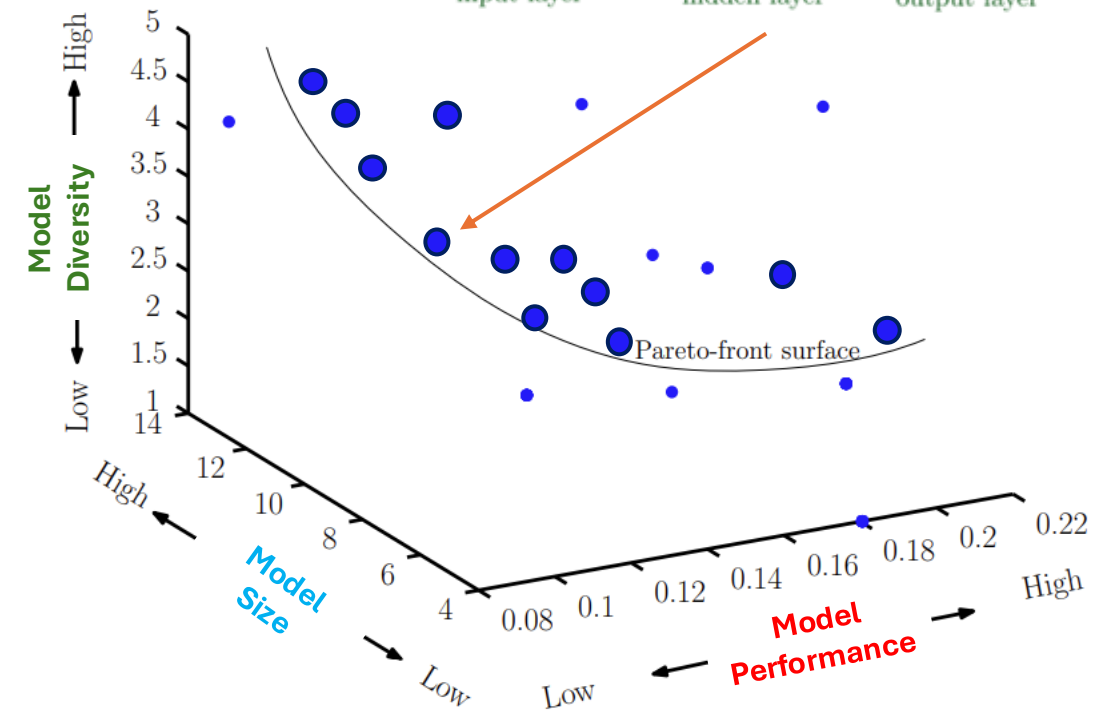
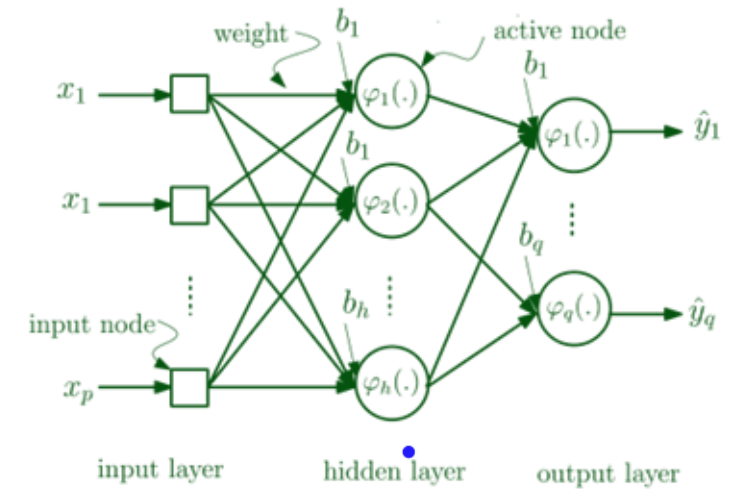
Efficiency of Smaller Models

Smaller models offer faster performance and lower costs, making them suitable for real-time applications, though they may compromise on accuracy.



Architectural Impact

The structure of a model significantly affects how effectively its parameters are utilized, with optimized architectures enhancing performance metrics.

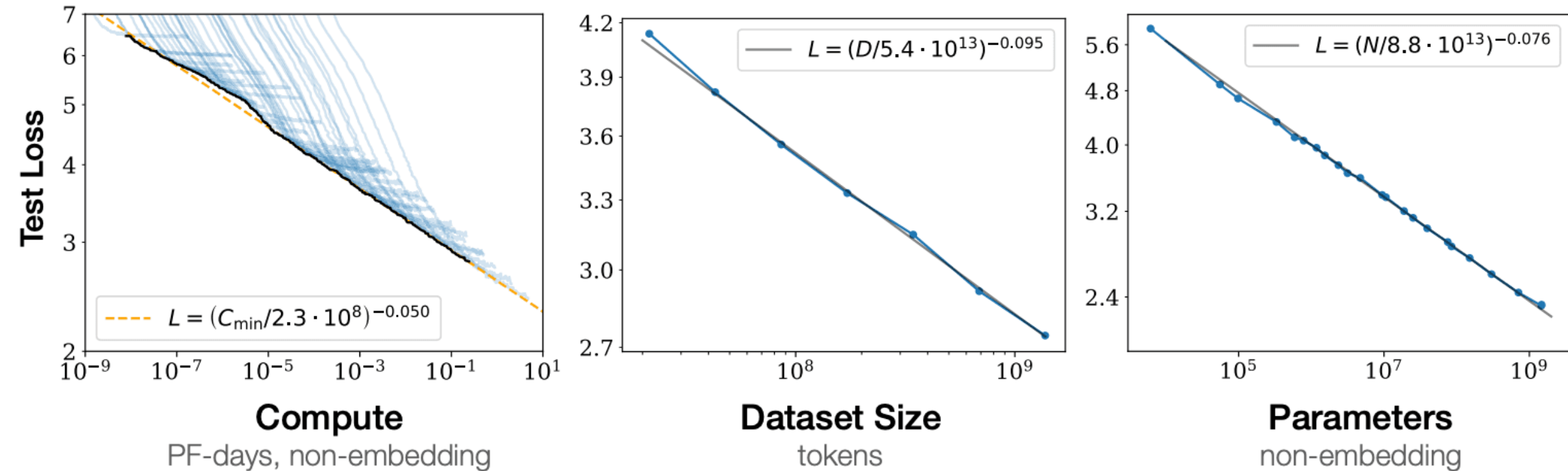


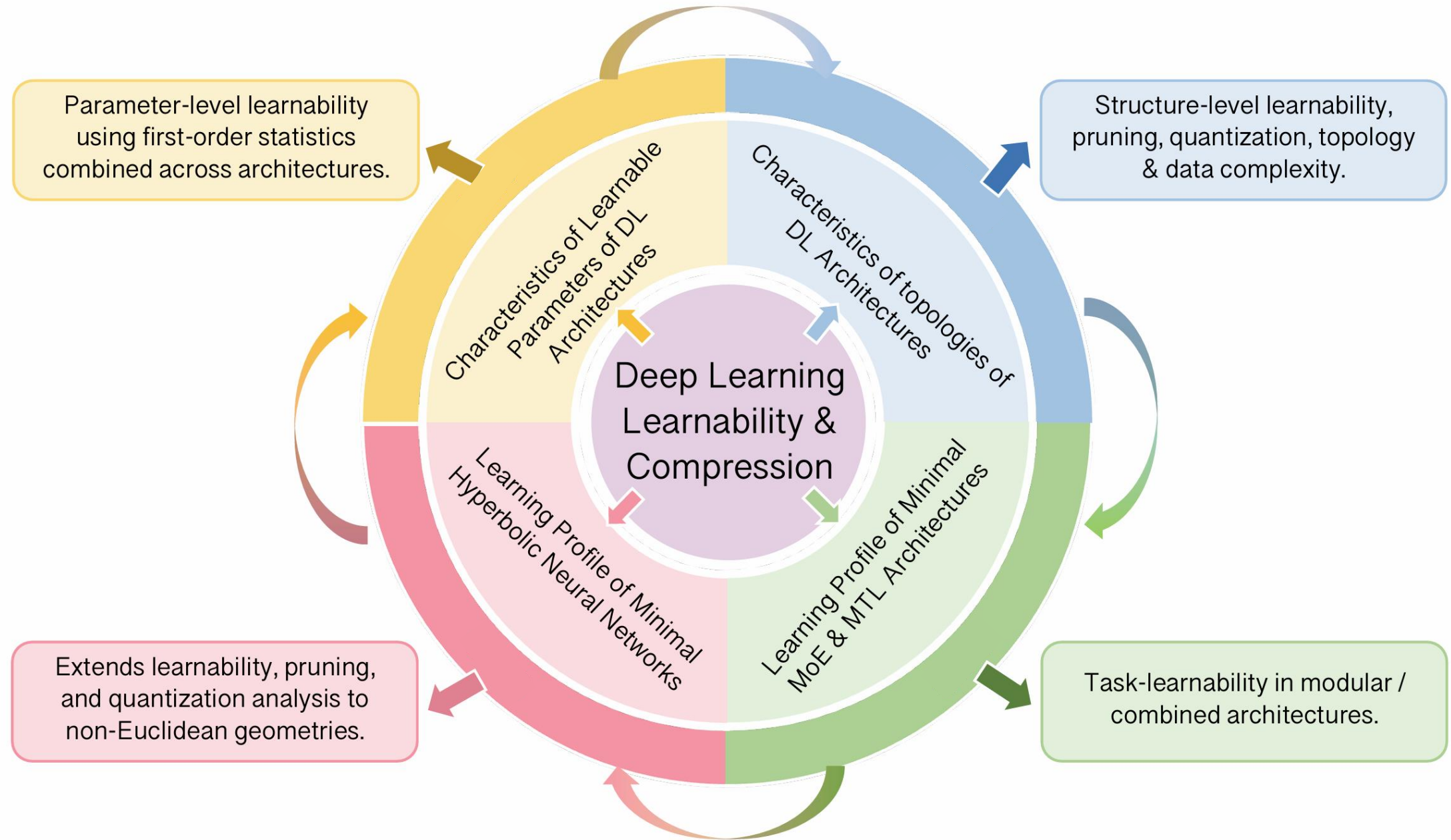
Part 2

Understanding AI Model Performance

AI Model Performance

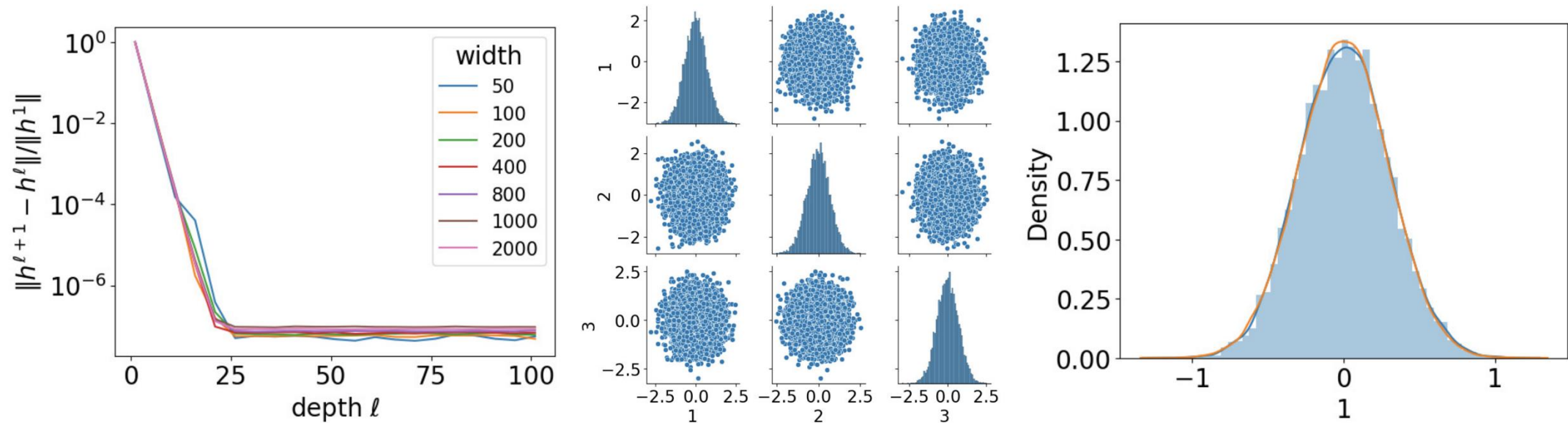
AI model's performance predictably improves as you increase resources like model size (parameters), dataset size, and compute power, often following power-law relationships where gains diminish but remain consistent.





Trained Network Exhibits Gaussian distribution

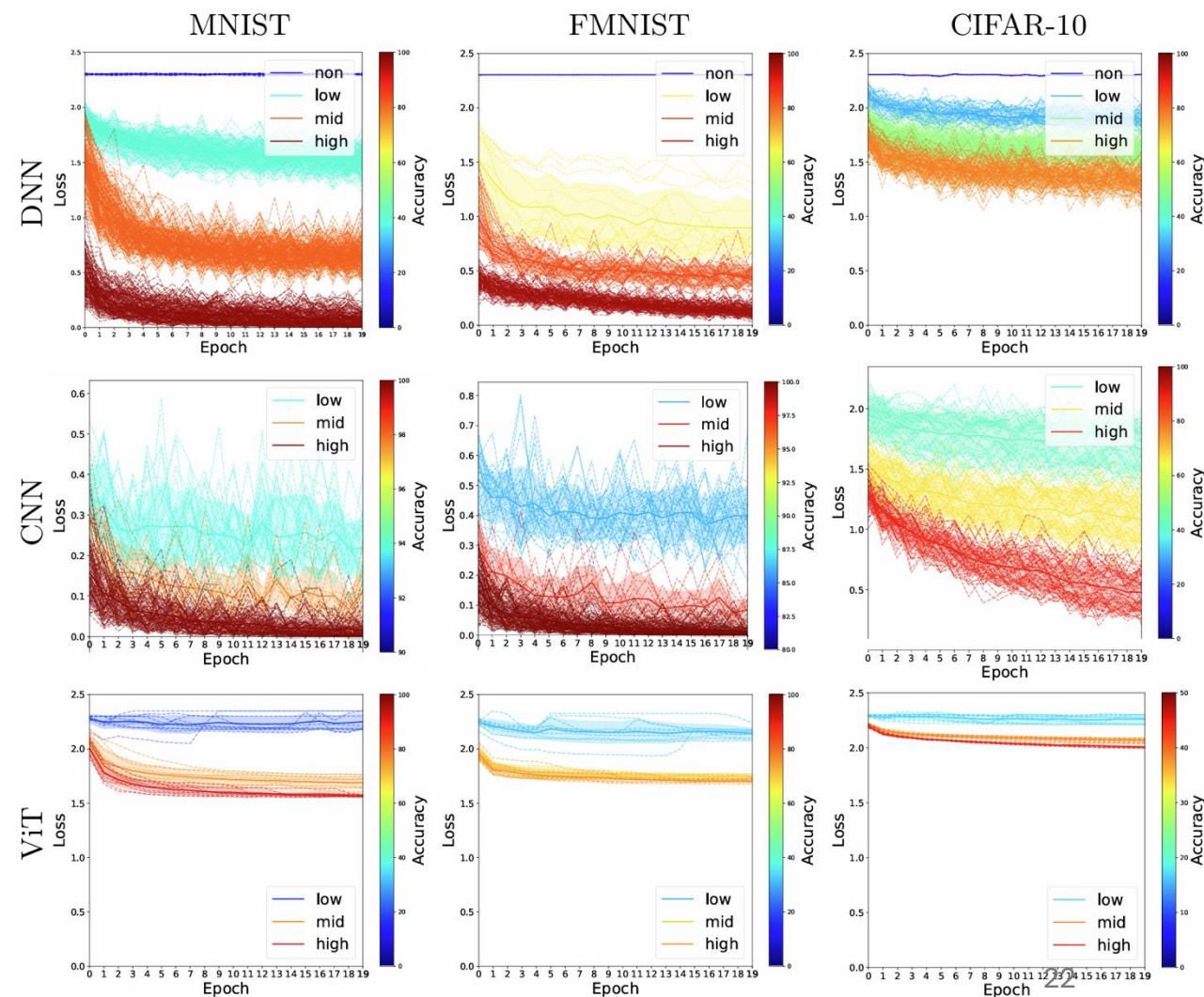
Convergence to the fixed point (left); joint distributions for the first neuron for three outputs for 10, 000 neuron networks, with orange curve denotes the Gaussian distribution (right)



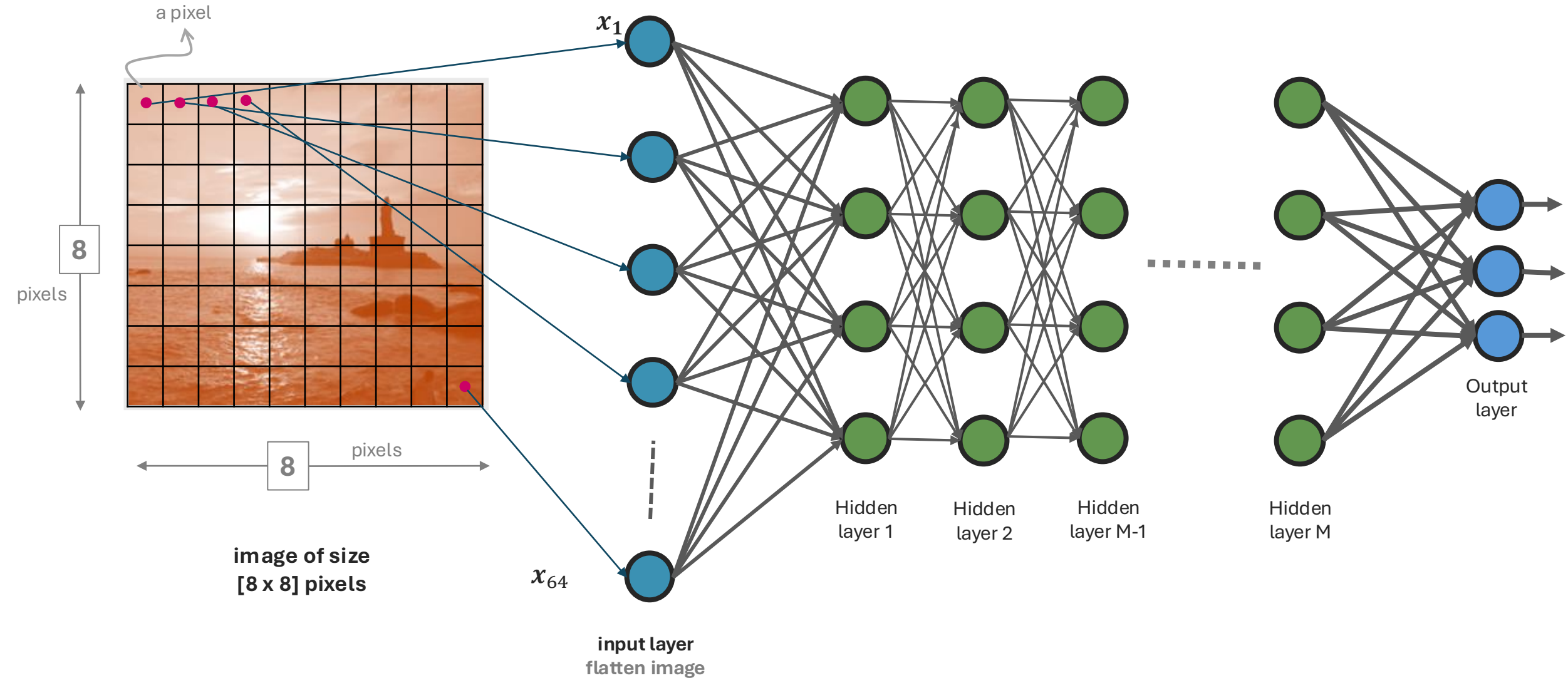
Performance Characterization of Models

Irrespective of model parameter size the model's convergence cluster into optimal and suboptimal patterns

Network	Data	Layer	Input	Output
DNN	MNIST/ F-MNIST	FC1	28×28	5-200
		FC2	5-200	5-200
		FC3	5-200	10
	CIFAR-10	FC1	$3 \times 32 \times 32$	5-1000
		FC2	5-1000	5-1000
		FC3	5-1000	10
CNN	MNIST	Conv1	$28 \times 28 \times 1$	$26 \times 26 \times C$
	F-MNIST	FC	$26 \times 26 \times C$	10
	CIFAR-10	Conv1	$3 \times 32 \times 32$	$30 \times 30 \times C$
		FC	$30 \times 30 \times C$	10
ViT	MNIST	Encoder	$d_model=784$	-
	F-MNIST	FC	$nhead=2-16$	10
	CIFAR-10	-	784 (input)	-

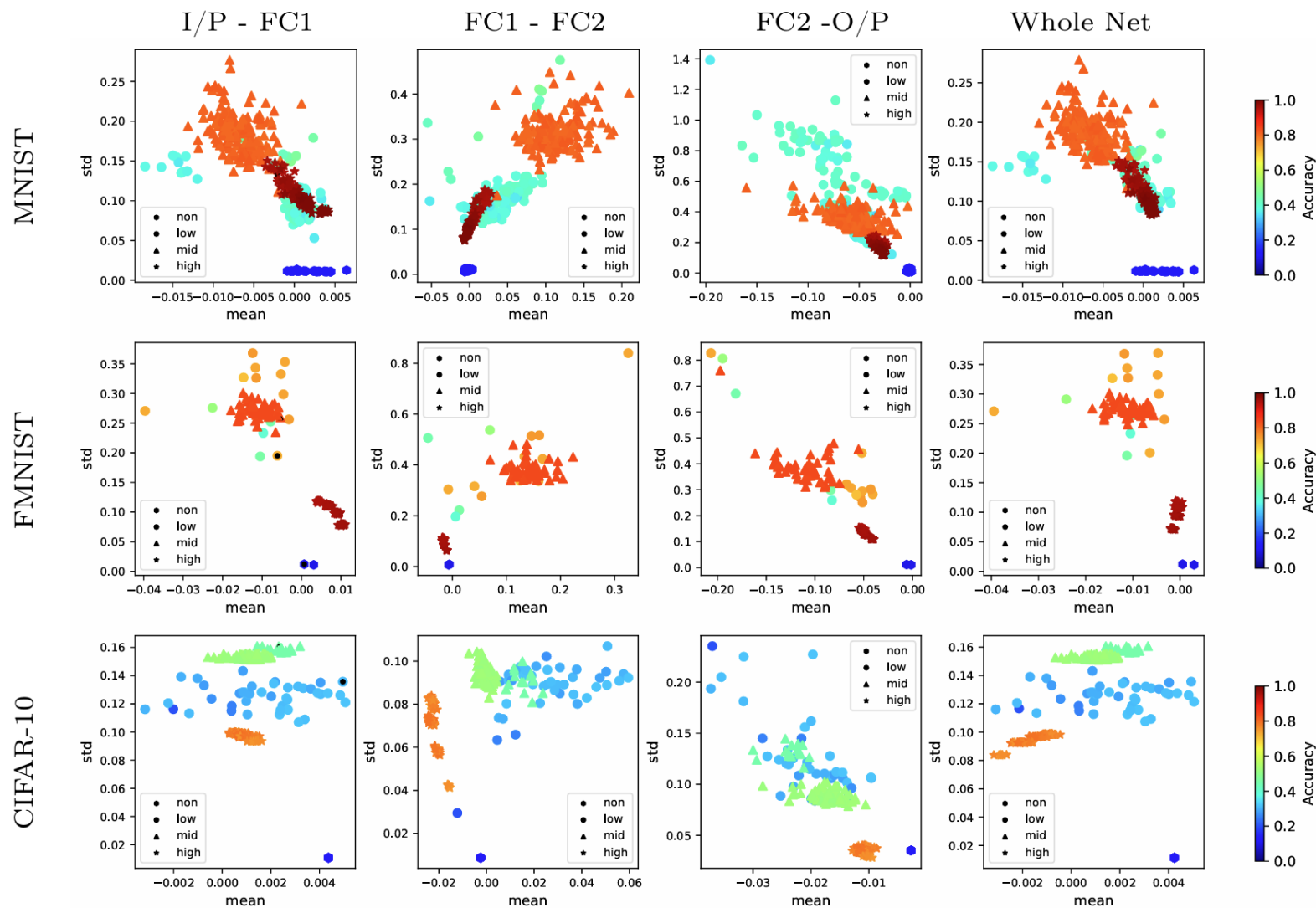


Deep Neural Networks



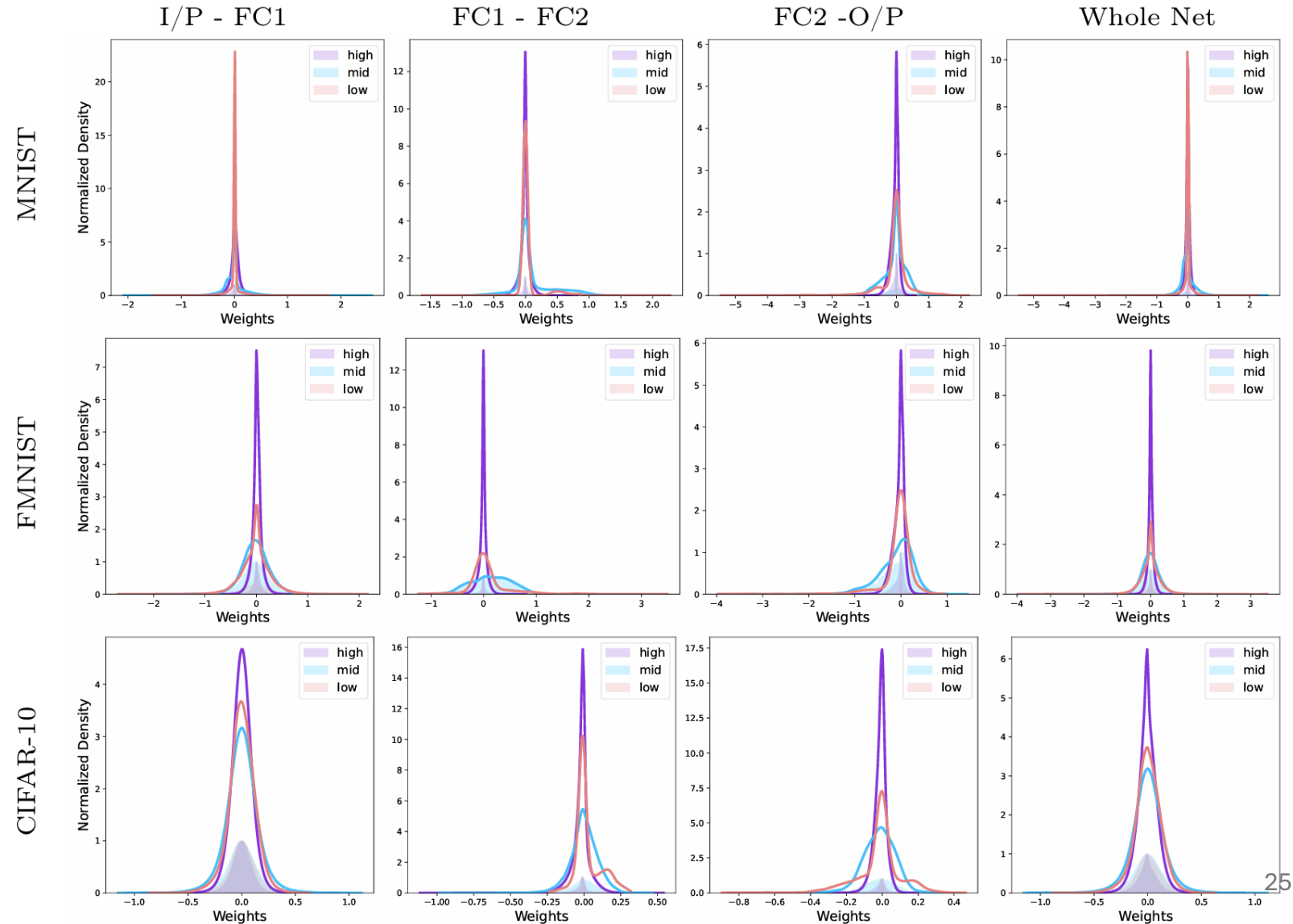
DNN's Learnable Parameters

Learned parameter of all **failed models**, **suboptimal models** and **successes full models** clearly show similarity of their statistical properties and weight distribution irrespective of their model size. i.e. models of any size converge to a similar weight distribution when trained to a fixed number of epochs on a dataset

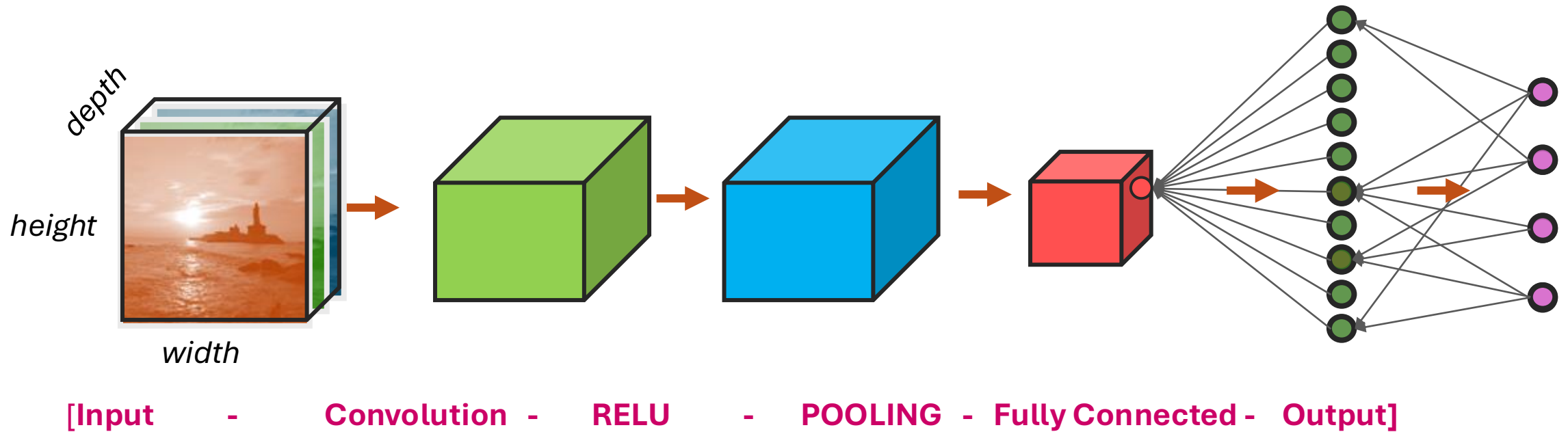


DNN's Learnable Parameters Distribution

Learned parameter of all failed models, suboptimal models and successes full models clearly show similarity of their statistical properties and weight distribution irrespective of their model size. i.e. models of any size converge to a similar weight distribution when trained to a fixed number of epochs on a dataset



Convolutional Neural Nets

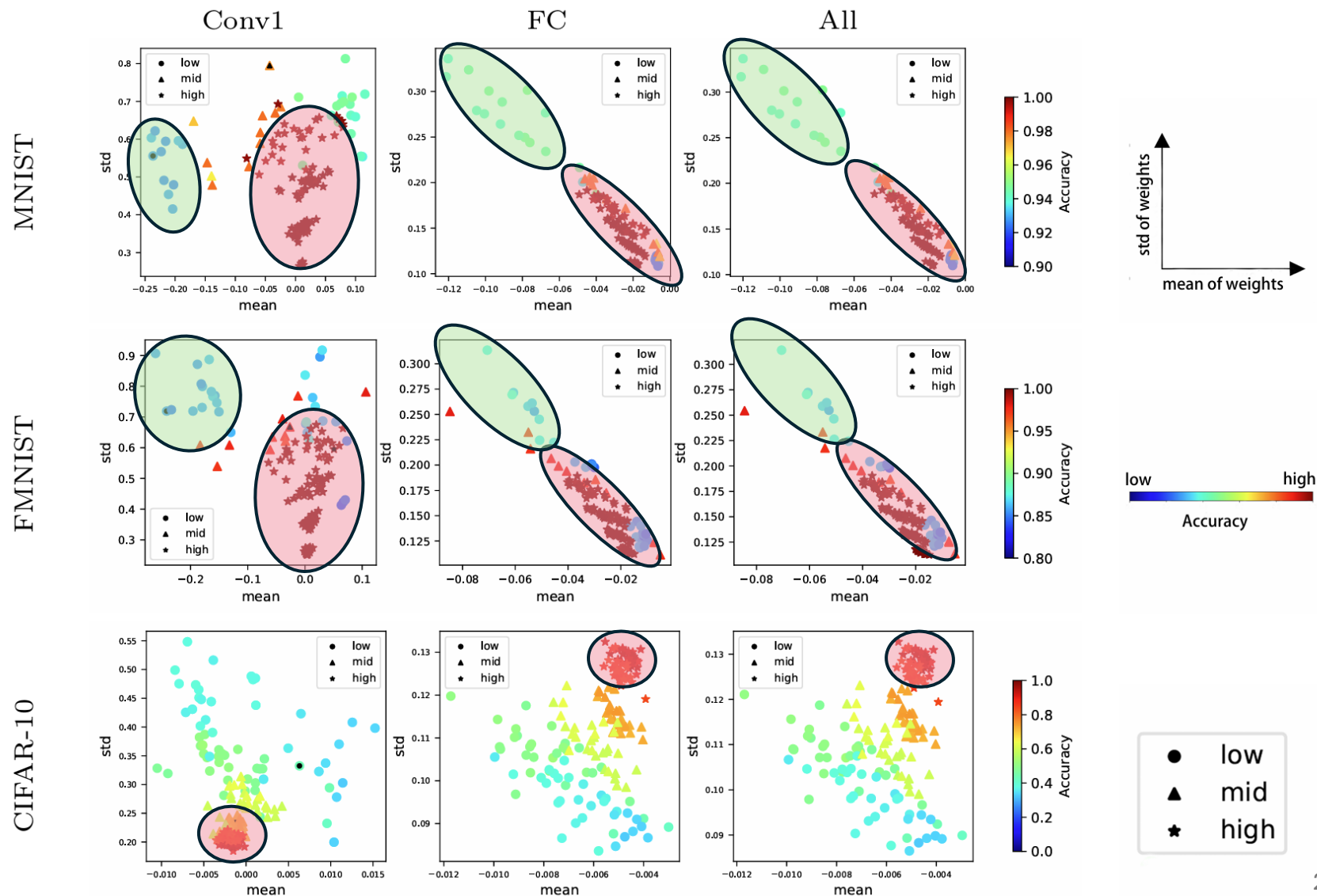


CNN's Learnable Parameters

Similar complexity data shows similar trainable parameter **statistical property pattern** (weight distribution), irrespective of model size

High performing model weights show that models of any size converge to parameter with **similar statistics** (weight distribution)

High complexity data shows that statistical property of successful models has very **specific statistical property** or specific weight distribution

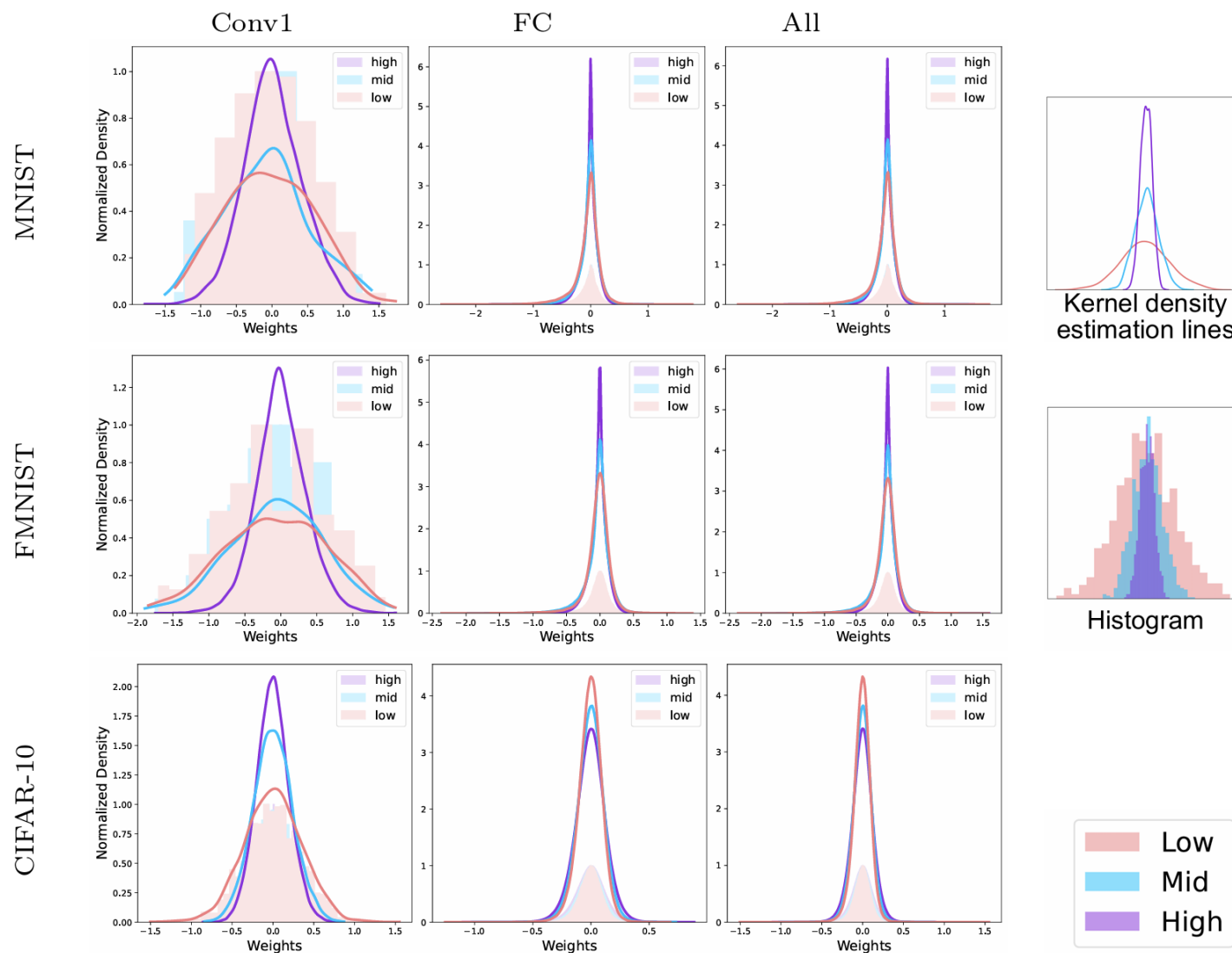


CNN's Learnable Parameters Distribution

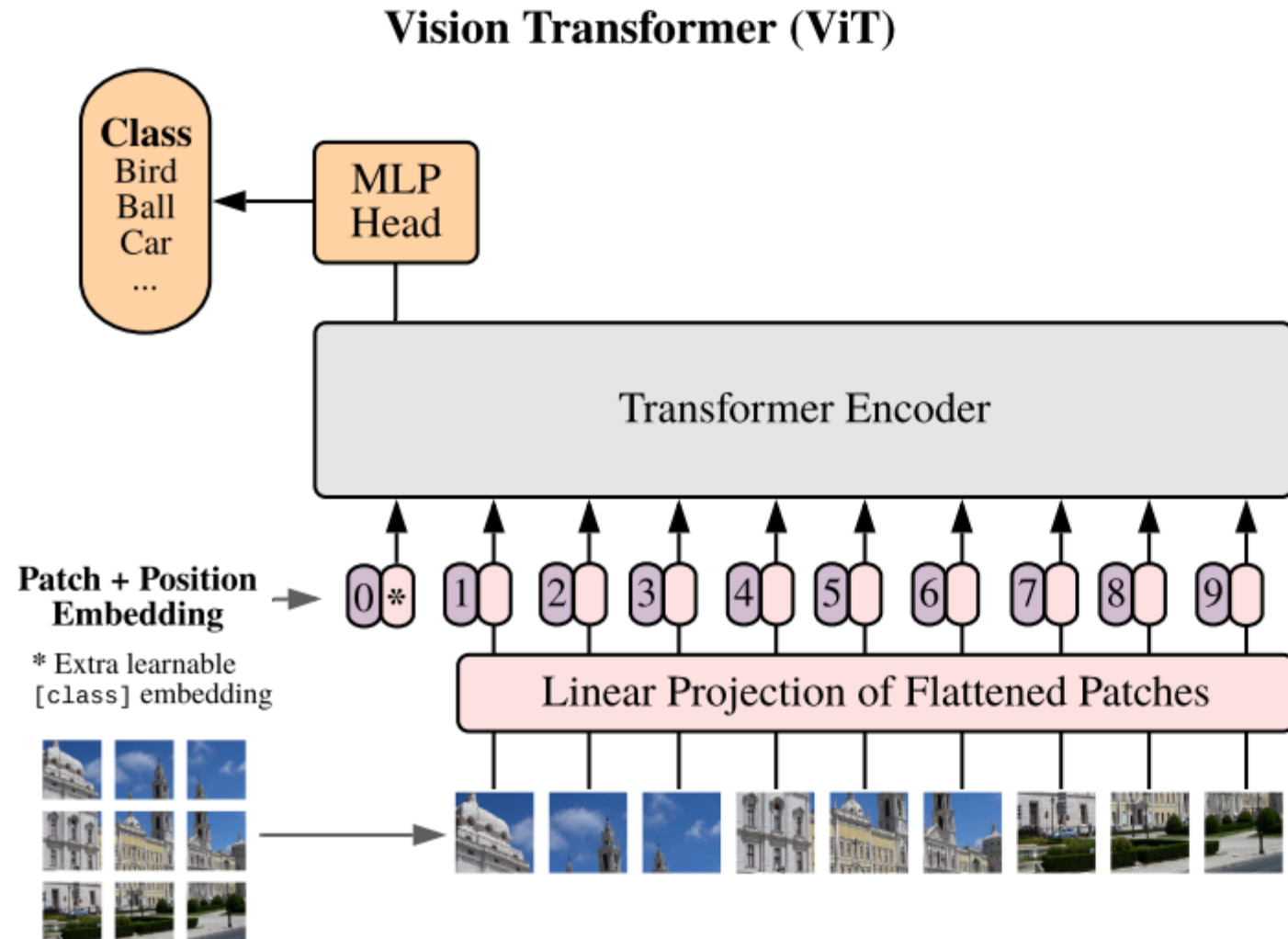
Similar complexity data shows similar trainable parameter statistical property pattern (**weight distribution**), irrespective of model size

High performing model weights show that models of any size converge to parameter with similar statistics (**weight distribution**)

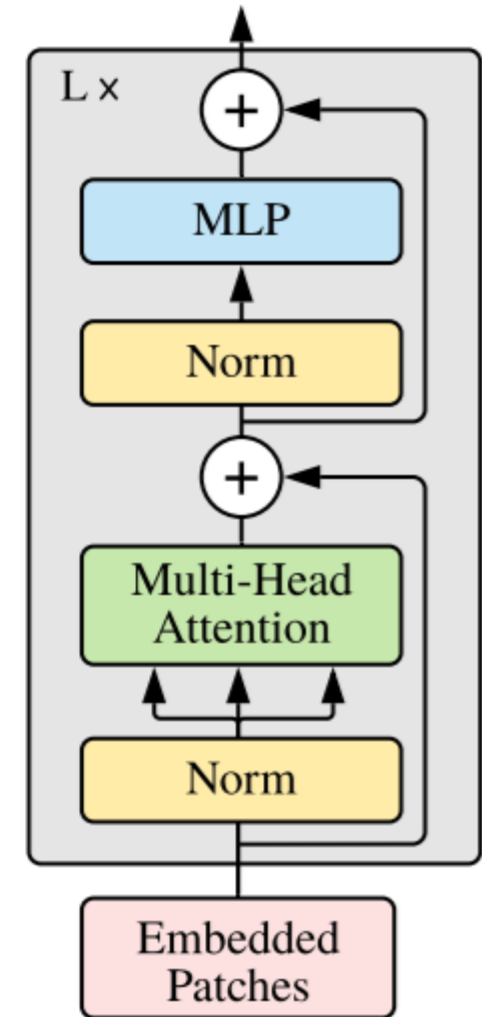
High complexity data shows that statistical property of successful models has very specific statistical property or **specific weight distribution**



Vision Transformer Models

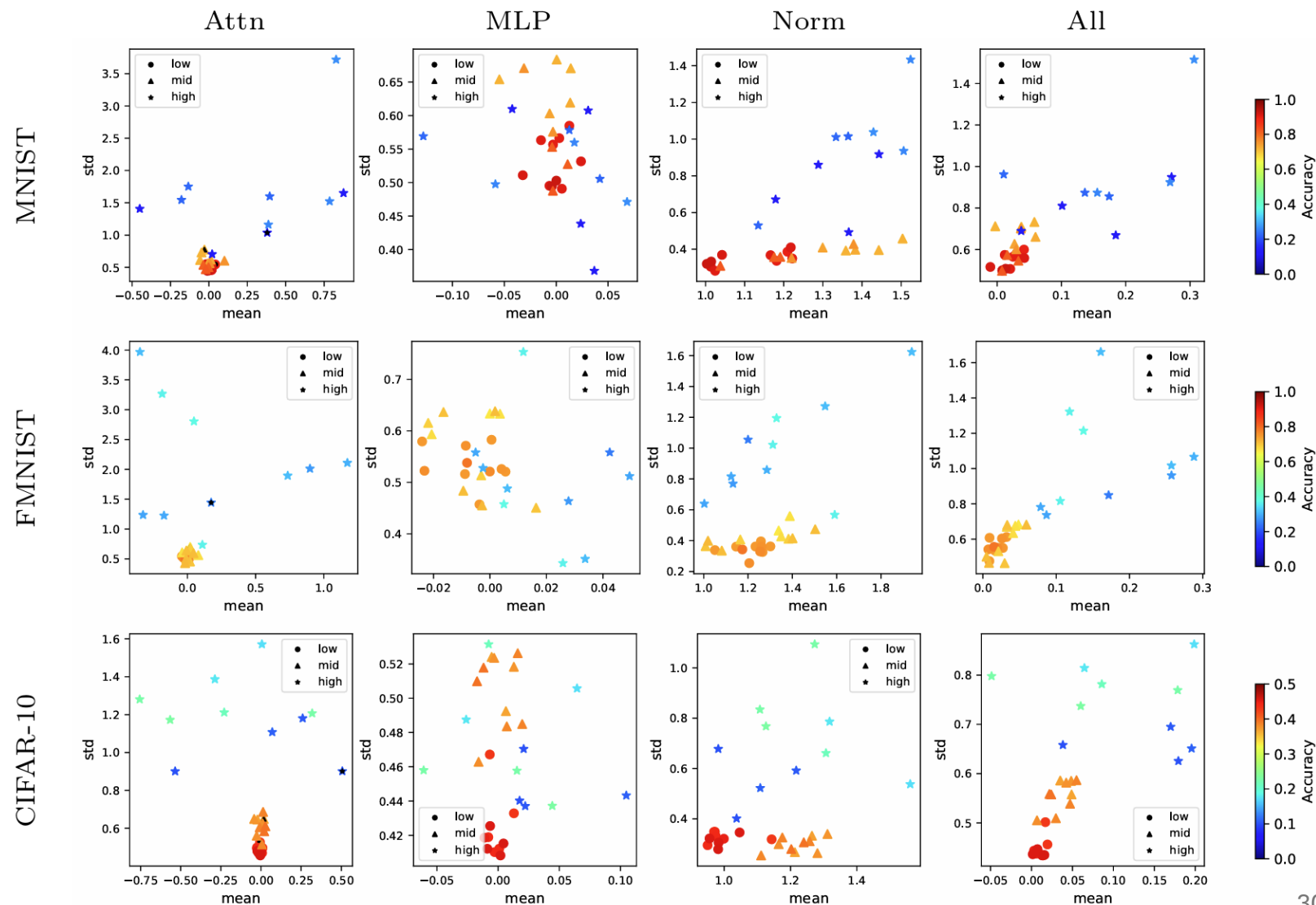


Transformer Encoder



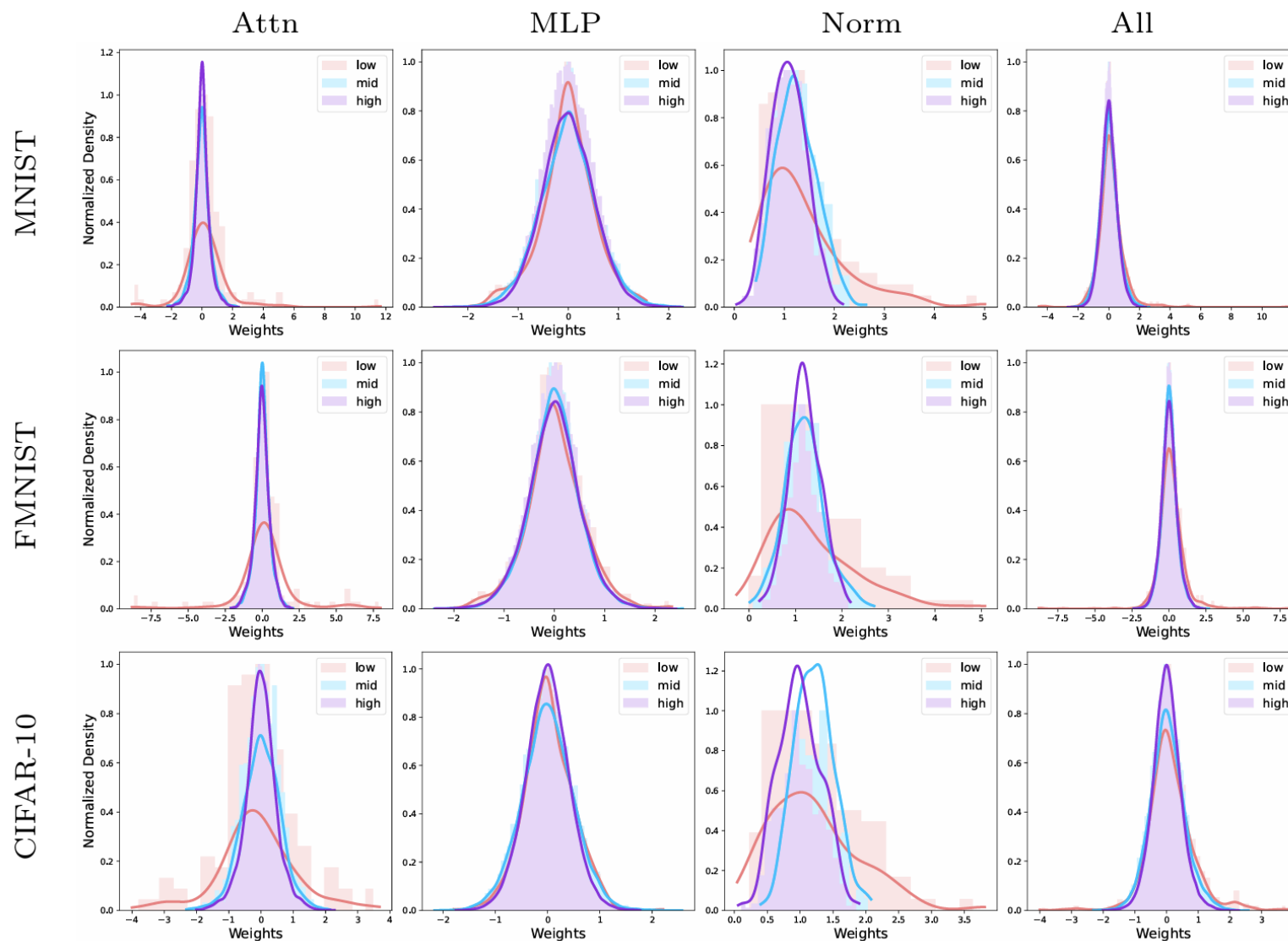
ViT's Learnable Parameters

All failed models, suboptimal models and successes full model's learned parameter clearly show similarity of their statistical properties (weight distribution) **irrespective of Model Architecture**



ViT's Learnable Parameters Distribution

All failed models, suboptimal models and successes full model's learned parameter clearly show similarity of their statistical properties (weight distribution) **irrespective of Model Architecture**



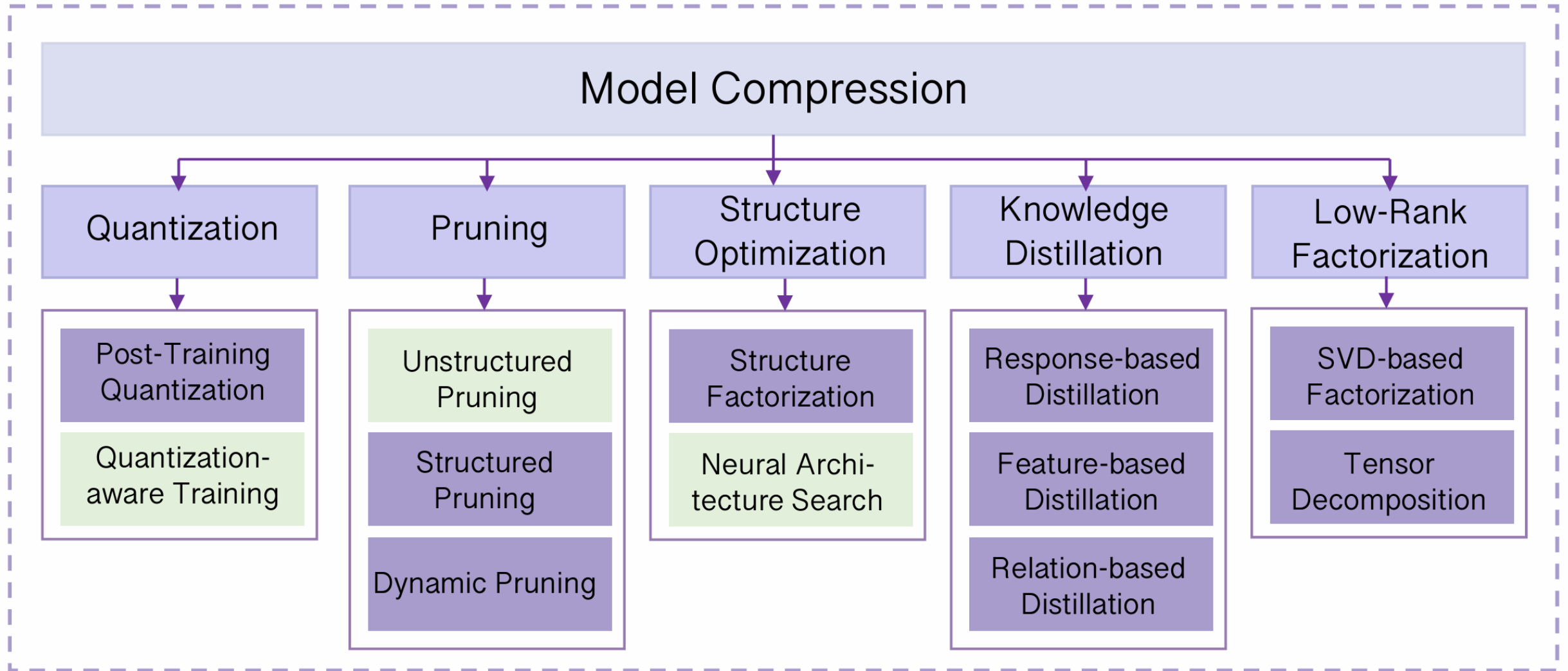
Part 3

AI Model

Compression

Methods and Results

AI Model "Slimming" Techniques



AI Model "Slimming" Techniques



Quantization

Reduces model precision (32-bit to 8-bit), cutting memory by ~75% with minimal accuracy loss.



Pruning

Eliminates less important weights to create smaller, sparser models while maintaining accuracy.



Knowledge Distillation

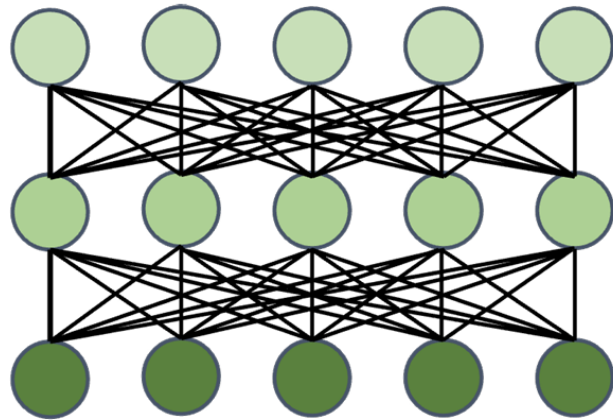
Trains a compact student model to replicate behaviors of a large teacher model.



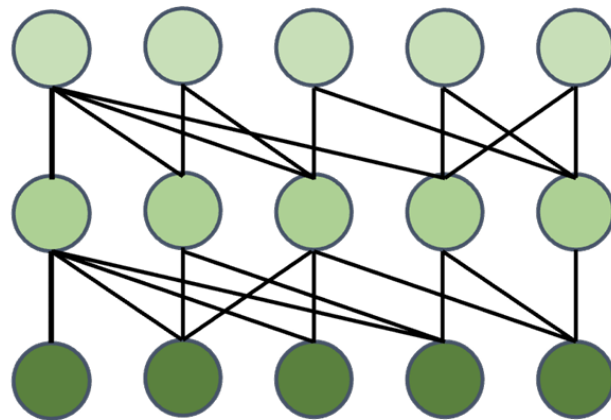
Model Splitting

Partitions processing between edge device and cloud for co-inference.

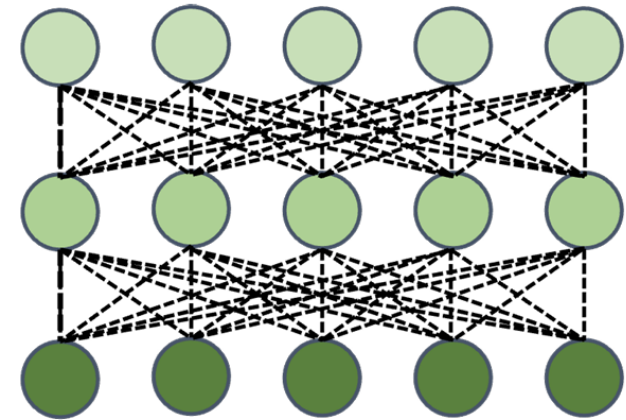
Pruning and Quantization Aware Neural Network Training



32-bit network



pruned network



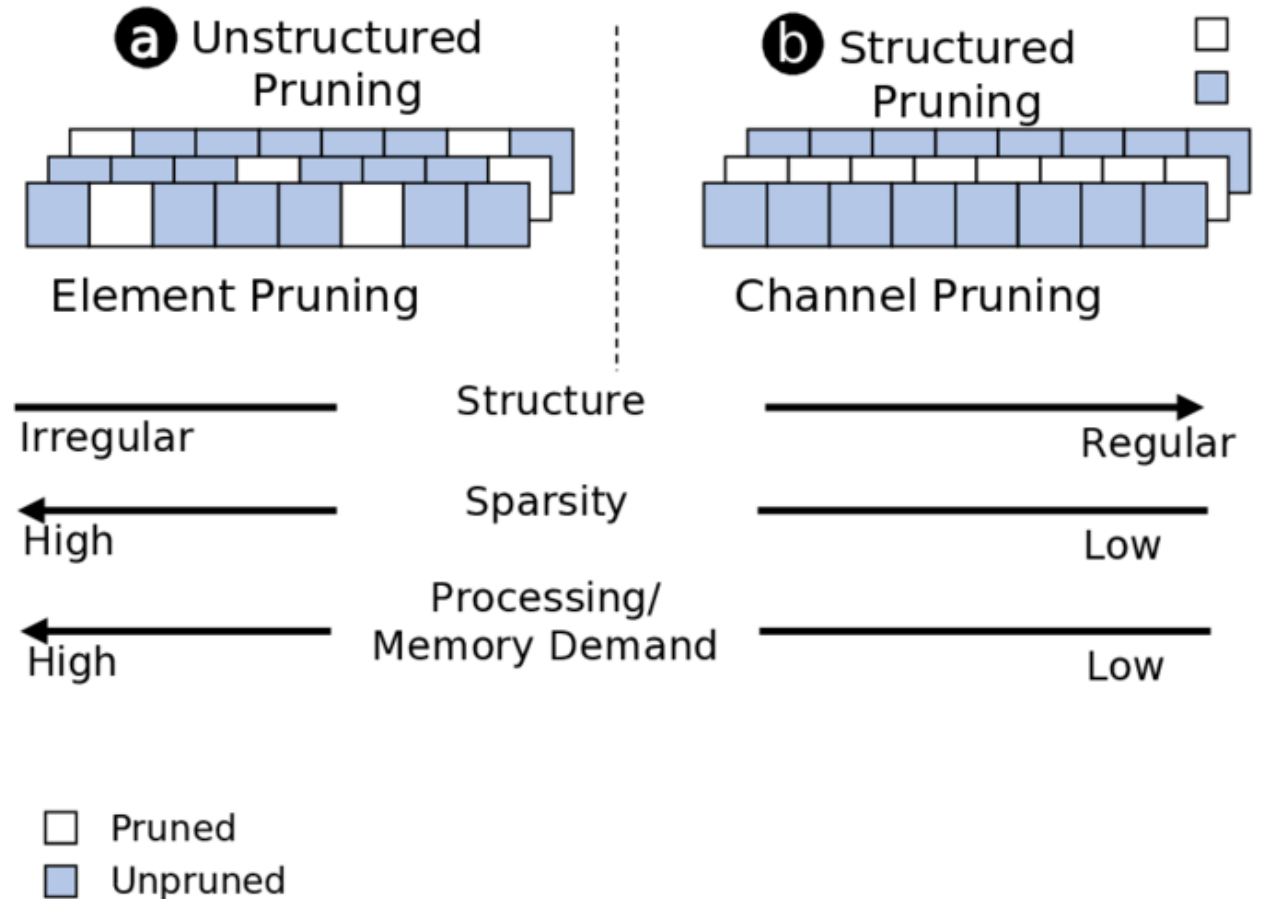
8-bit network

Pruning

- **Train-time (dynamic) pruning** involves integrating the pruning process directly into the training phase of the neural network. During training, the model is trained in a way that encourages sparsity or removes less important connections or neurons as part of the optimization process
- **Unstructured pruning** involves zeroing individual weights within the weight matrices
- **Structured pruning** removing entire structured groups of weights, the method reduces the scale of calculations that would have to be made in the forward pass through the model's weights graph

Pruning

- **Magnitude-based Pruning:**
Simple, small weights/channels are pruned
- **Gradient-based Pruning:**
method prunes weights that show smaller gradients over time.
- **Importance-based Pruning:**
Weights with lower “importance scores” are pruned first
- **Training-time (dynamic) pruning:** dynamic pruning adjusts the network structure during training based on real-time performance metrics



Quantization

The main task of model quantization in deep learning is to convert high-precision floating-point numbers in neural networks into low-precision numbers

0.34	3.75	5.64
1.12	2.7	-0.9
-4.7	0.68	1.43

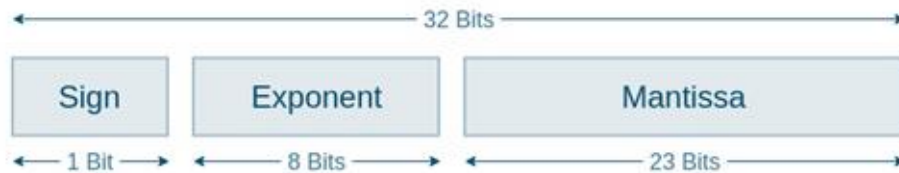
FP32



Quantization

64	134	217
76	119	21
3	81	99

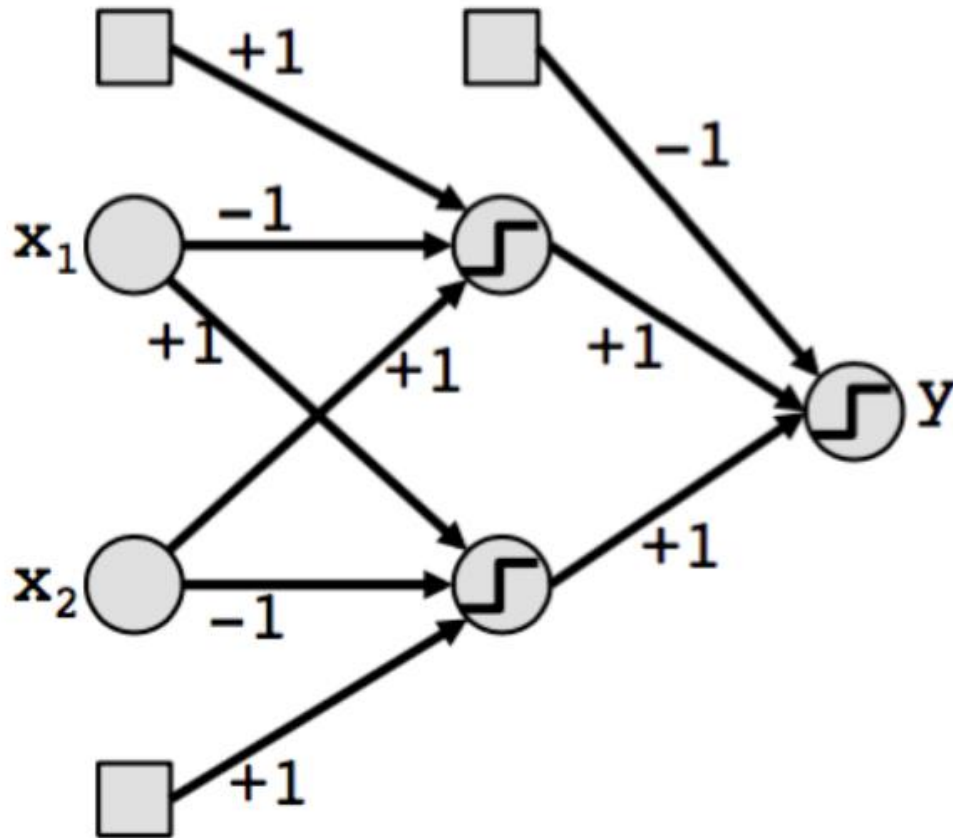
INT8



Single Precision
IEEE 754 Floating-Point Standard



Binarized Neural Networks

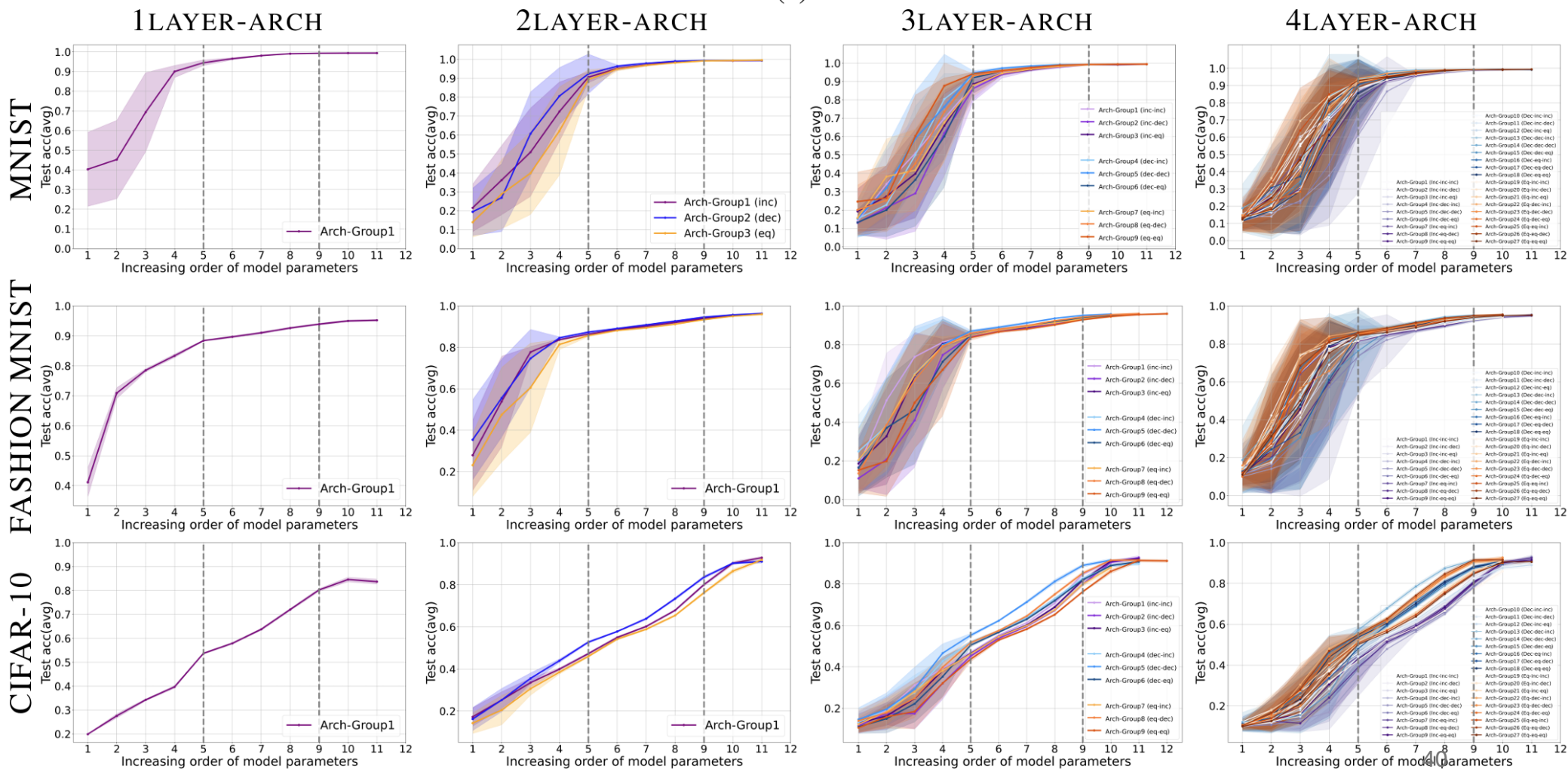
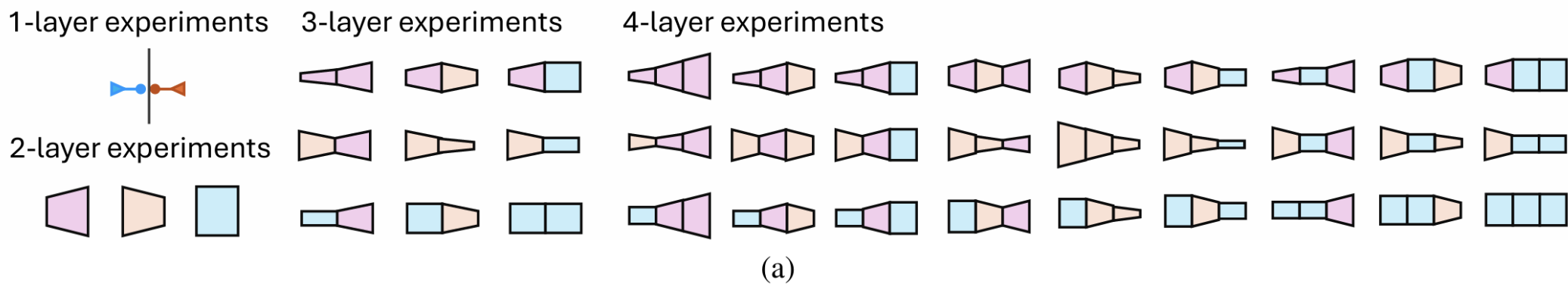


Weights and activation of Neural Networks:

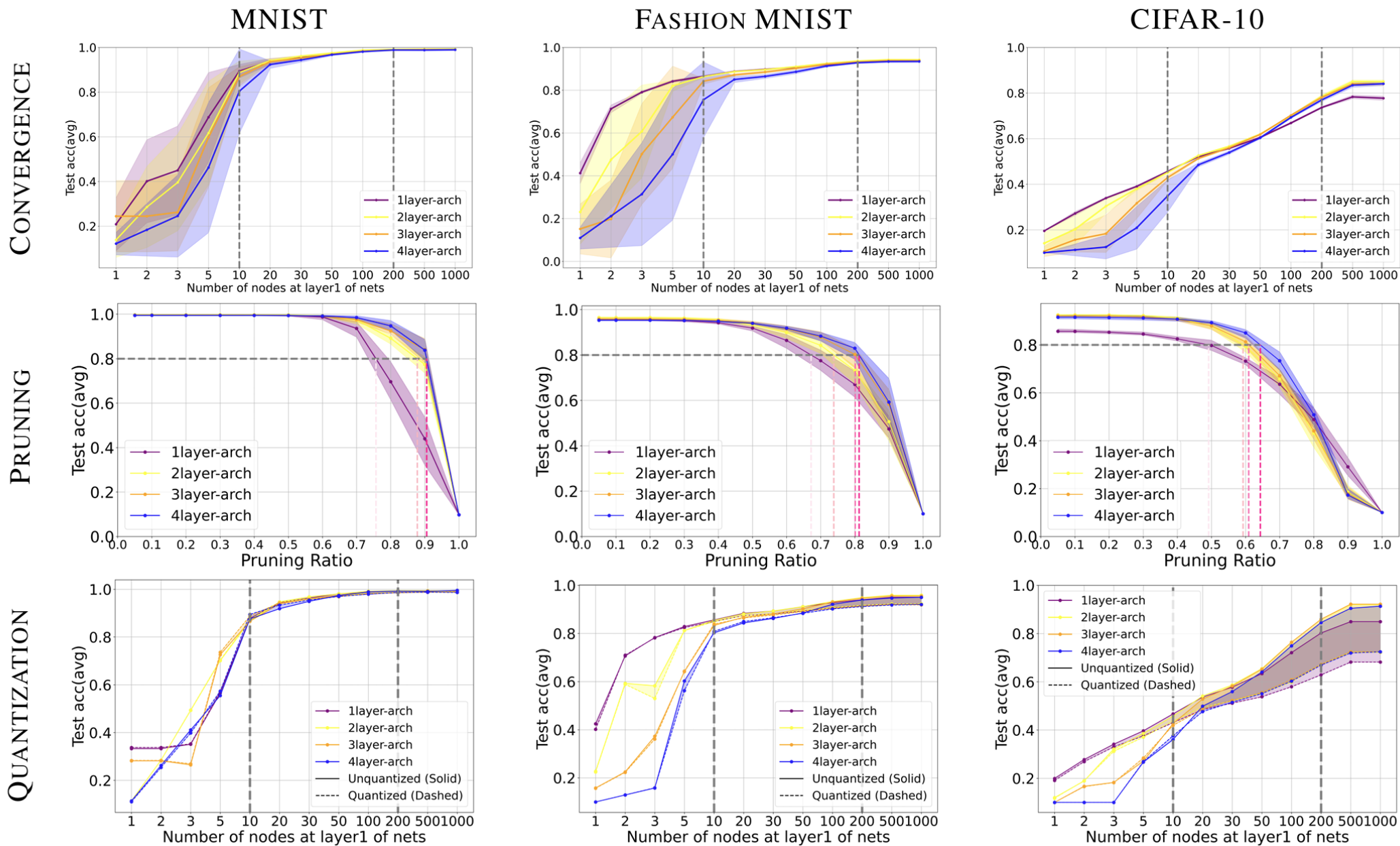
$$x^b = \text{Sign}(x) = \begin{cases} +1 & \text{if } x \geq 0, \\ -1 & \text{otherwise} \end{cases}$$

$$x^b = \begin{cases} +1 & \text{with probability } p = \sigma(x), \\ -1 & \text{with probability } 1 - p, \end{cases}$$

Architectural Invariance

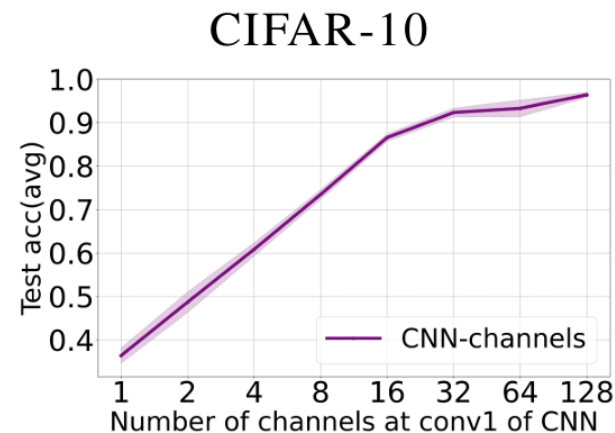
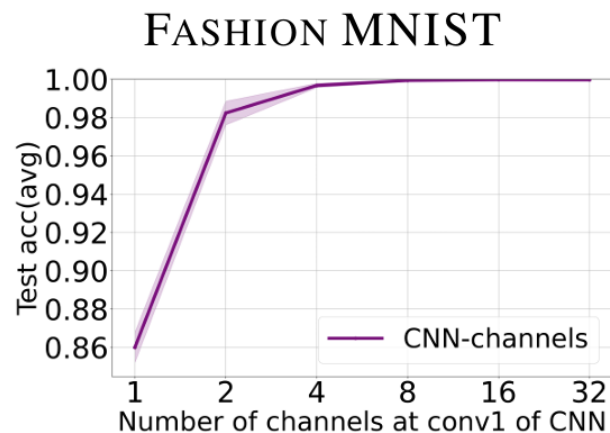
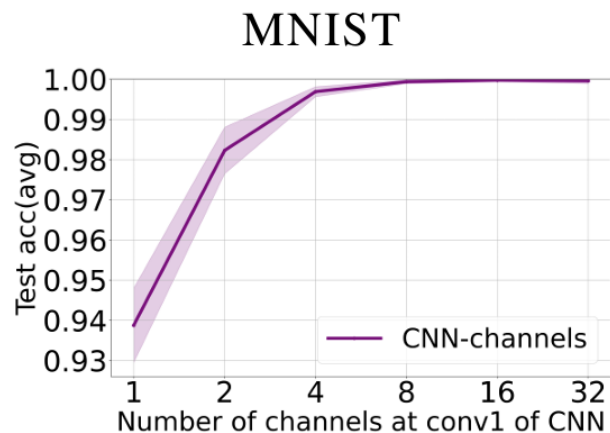


DNN Architecture

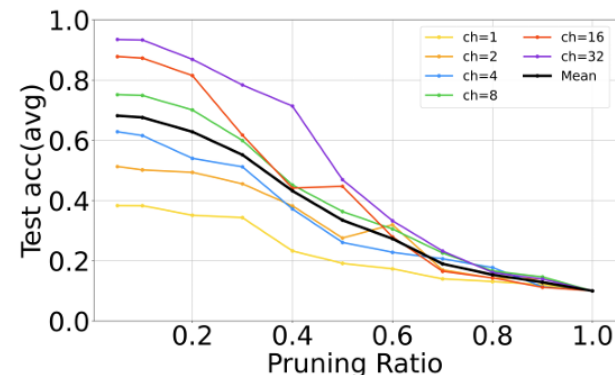
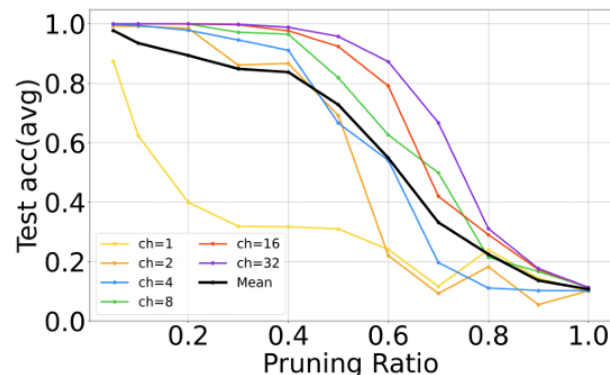
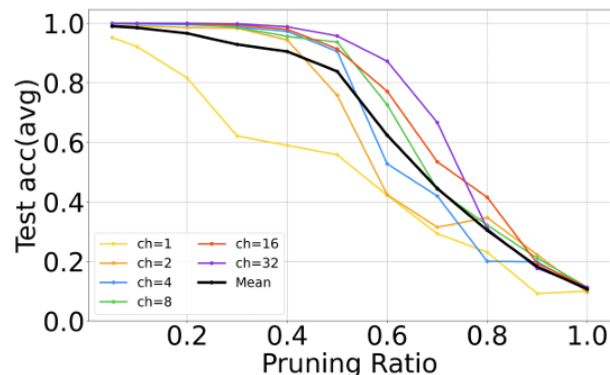


CNN Architecture

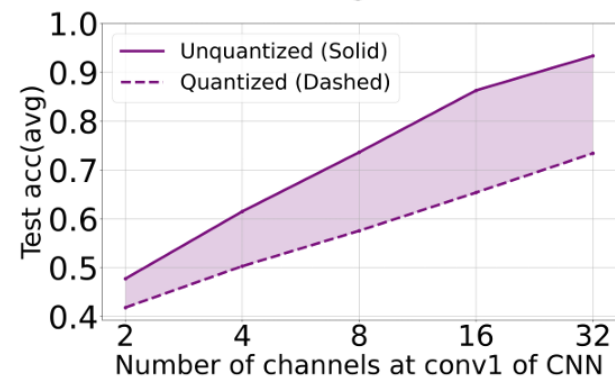
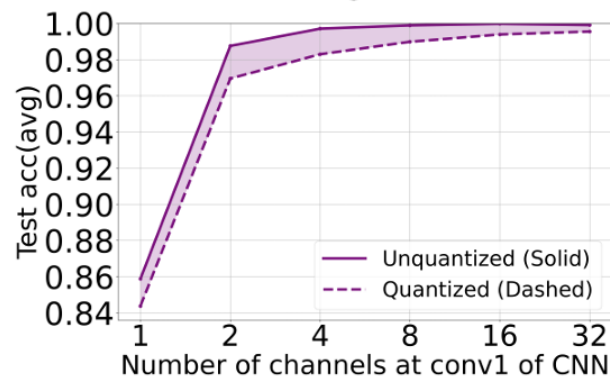
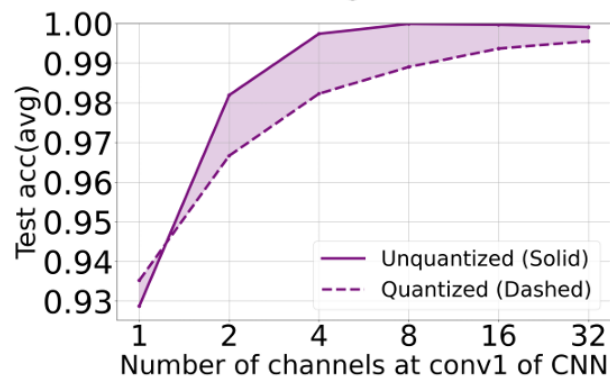
CONVERGENCE



PRUNING



QUANTIZATION



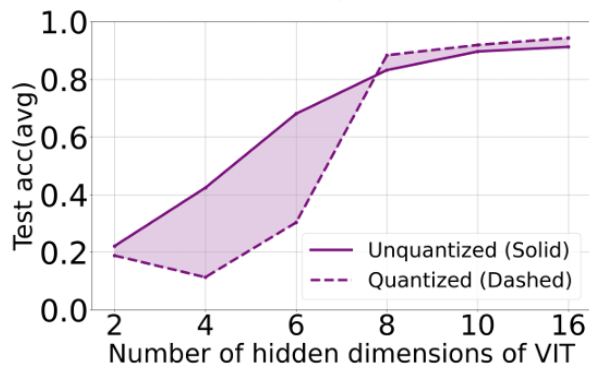
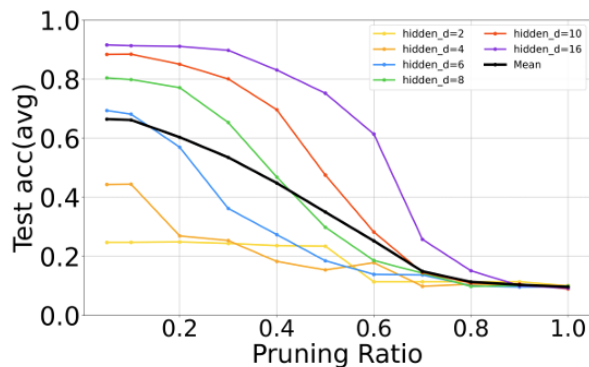
ViT Architecture

CONVERGENCE

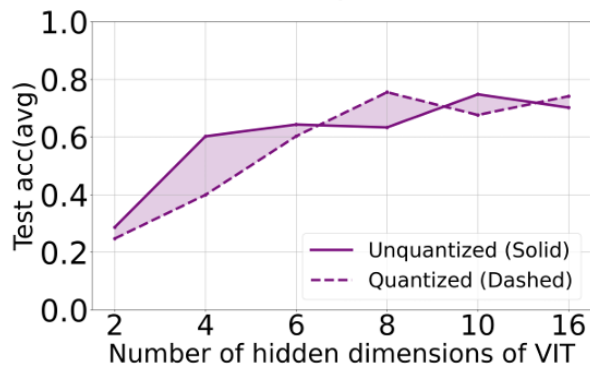
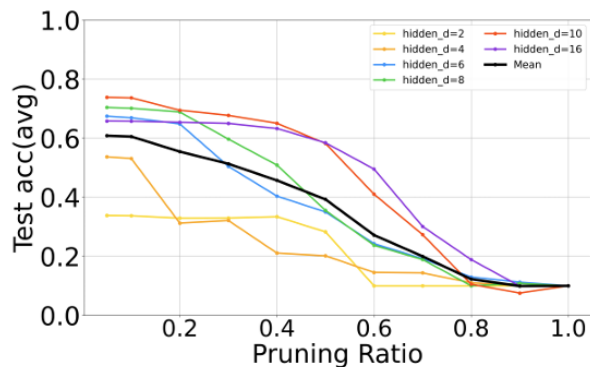
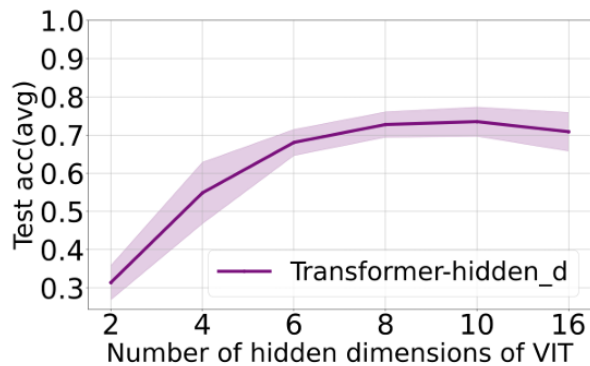
PRUNING

QUANTIZATION

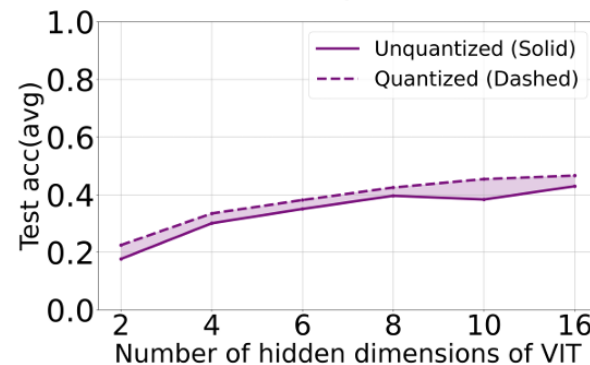
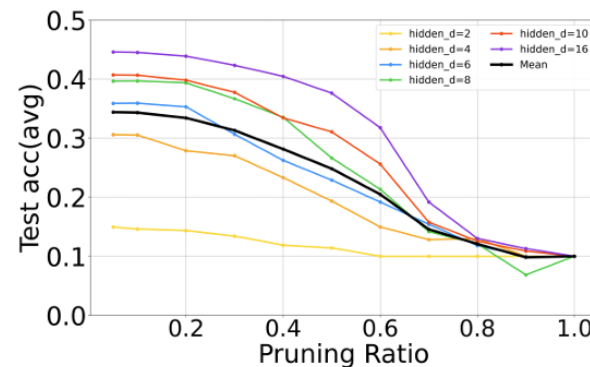
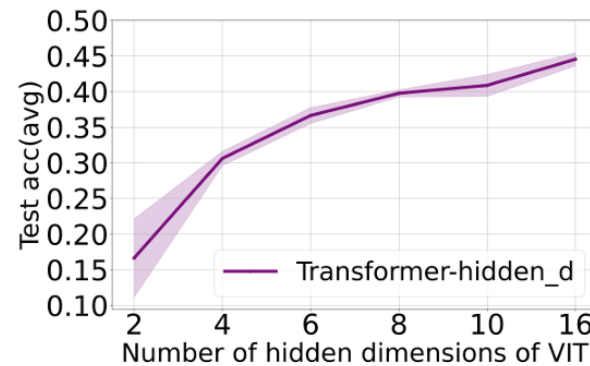
MNIST



FASHION MNIST



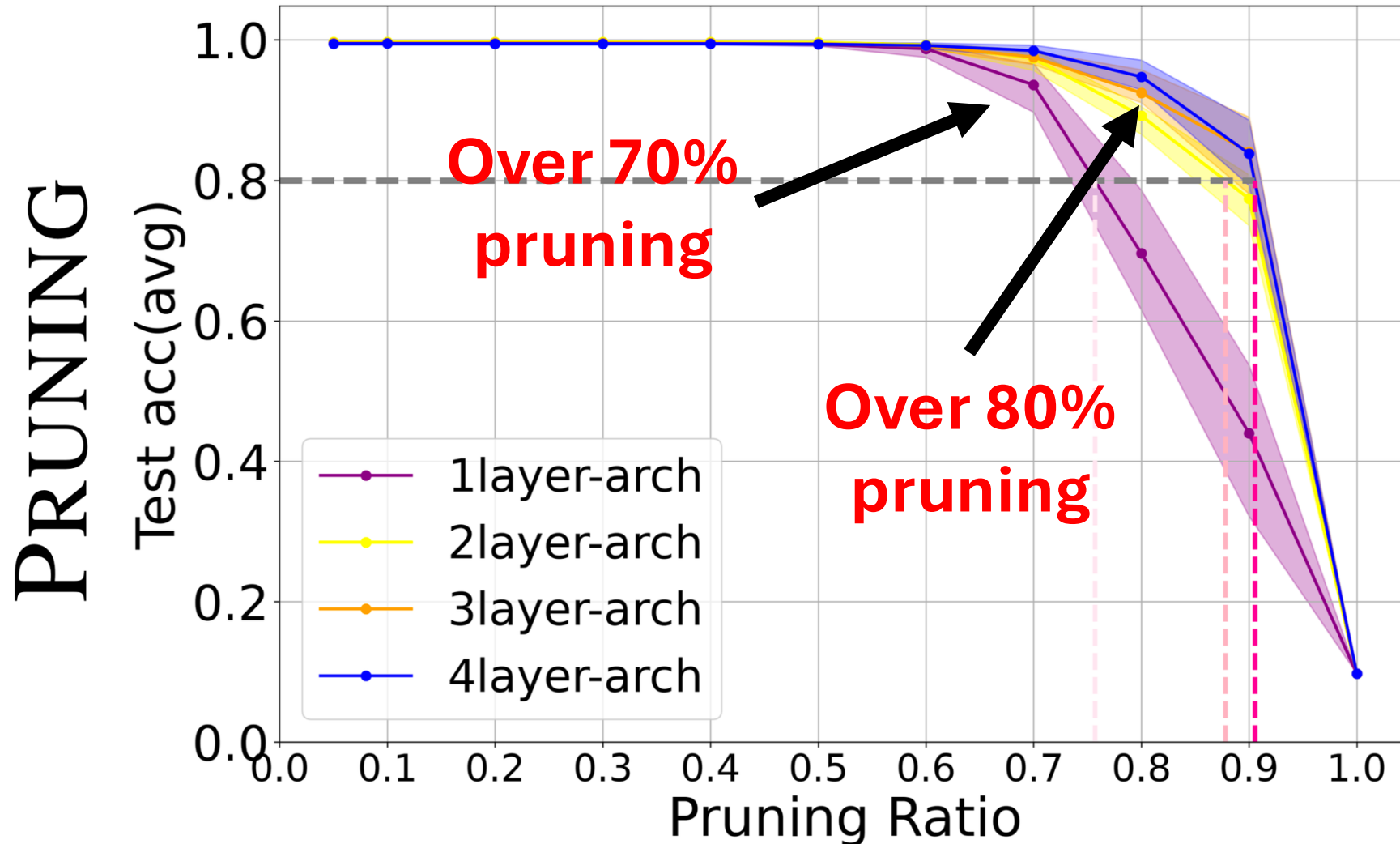
CIFAR-10



Maximal Pruning for Minimal Network

test results on MNIST

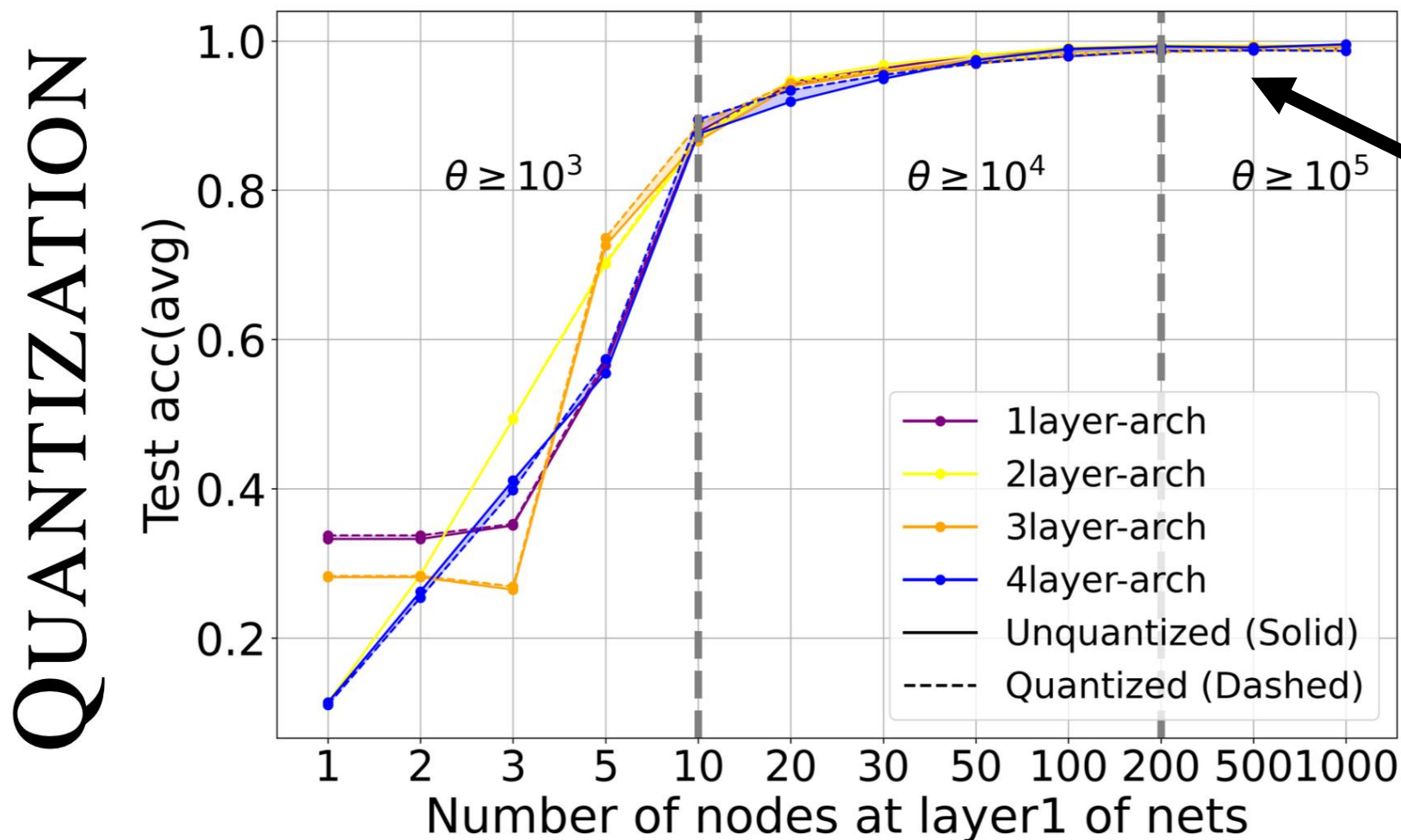
similar results of CIFAR 100 and other 8 datasets



Maximal Quntization for Minimal Network

test results on MNIST

similar results of CIFAR 100 and other 8 datasets

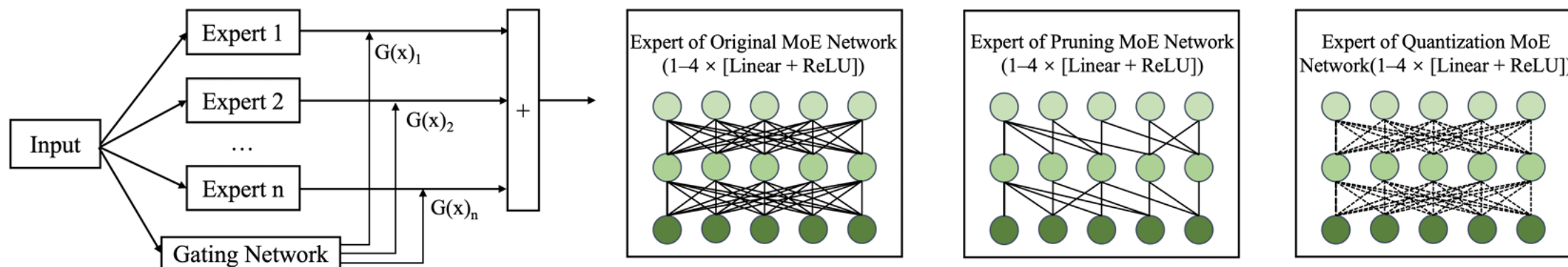
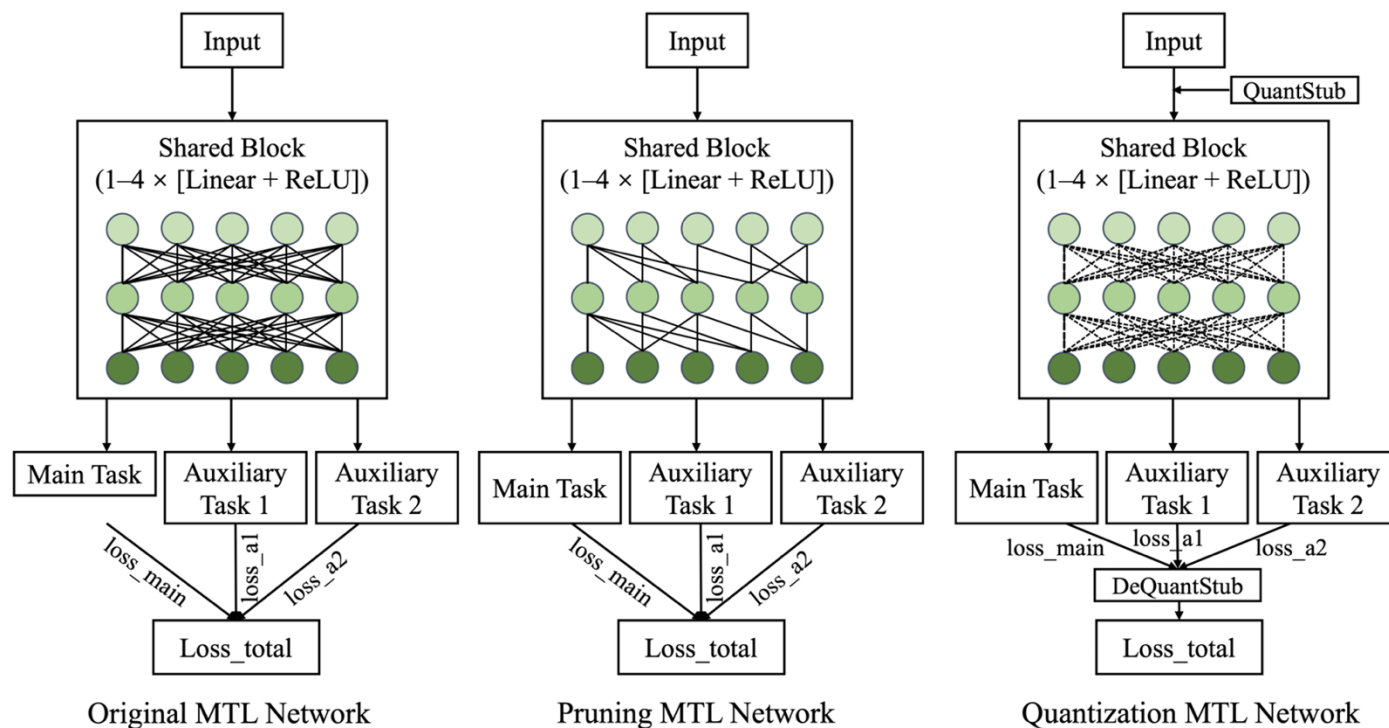


**Minimal
degradation
on 8-bit Nets**

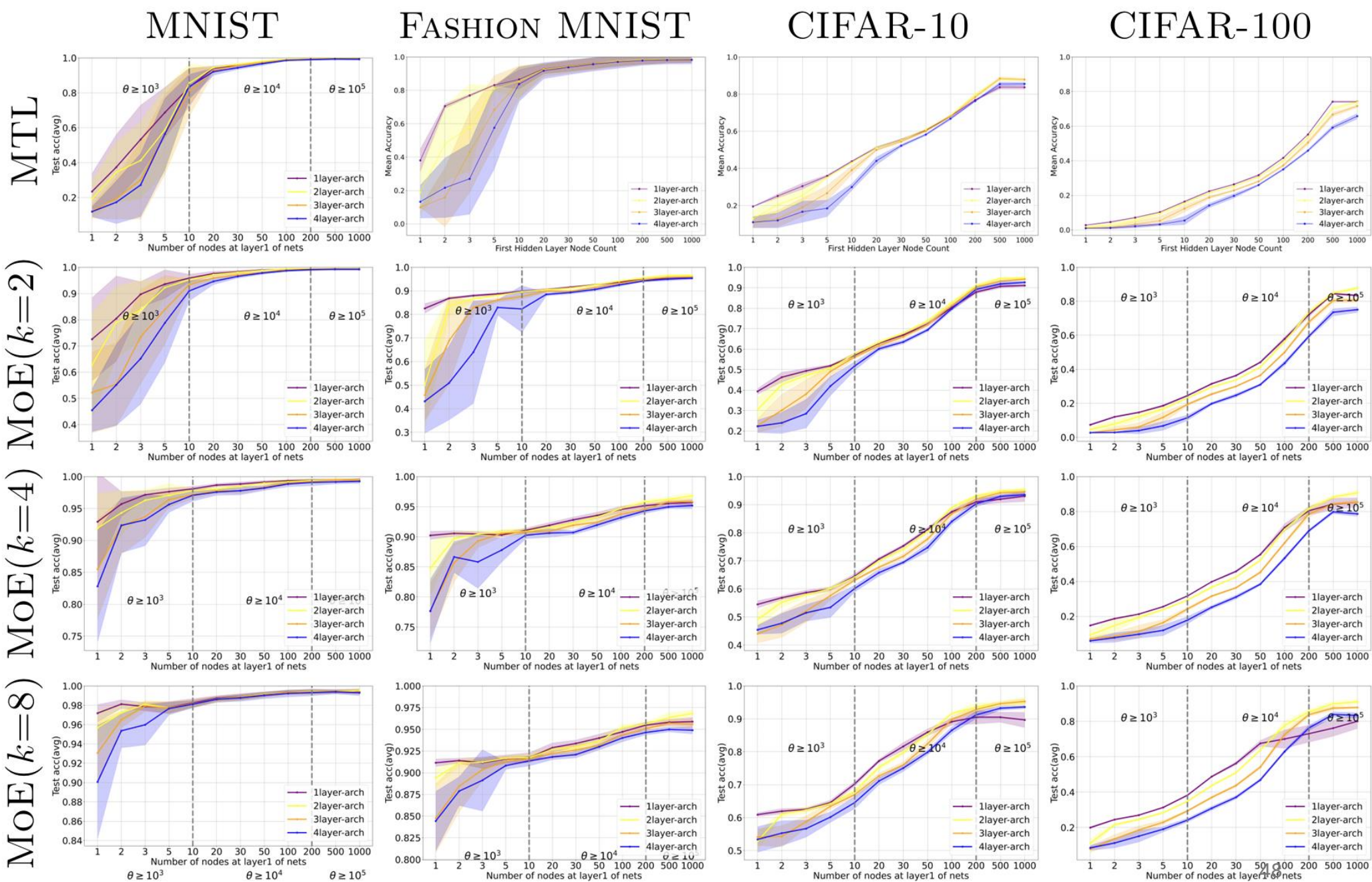
Summary of Pruning and Quantization on Models

Dataset	Architecture	Min. Params (Stability)	Safe Pruning %	8-bit Gap %
MNIST	DNN	2×10^4	60	0.15
	CNN	3×10^4	20	0.74
	ViT	5×10^3	20	10.66
F-MNIST	DNN	2×10^4	40	1.42
	CNN	3×10^4	20	1.10
	ViT	5×10^3	20	3.22
CIFAR-10	DNN	2×10^6	40	7.50
	CNN	3×10^6	5	14.80
	ViT	6×10^3	10	4.17

Multitask Learning and Mixture of Experts

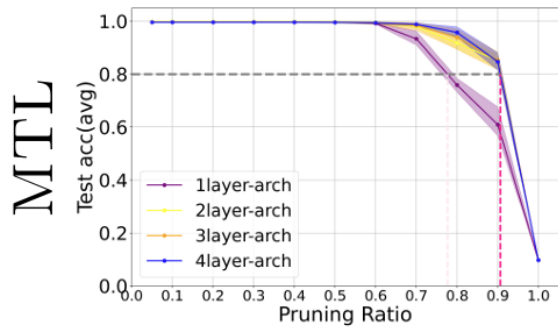


MTL and MoE Convergence

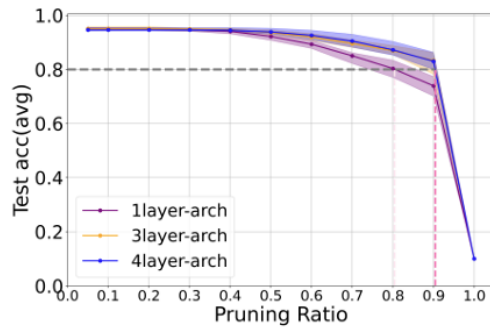


MTL and MoE Pruning

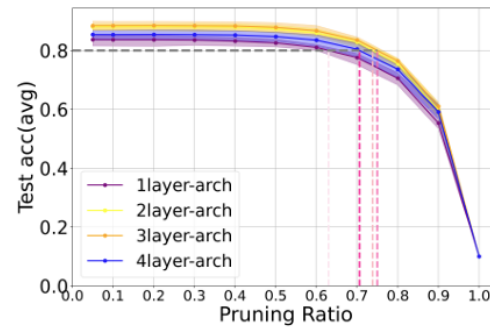
MNIST



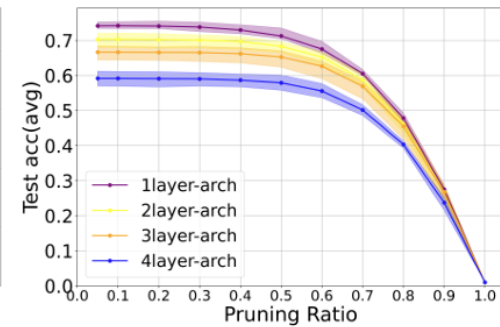
FASHION MNIST



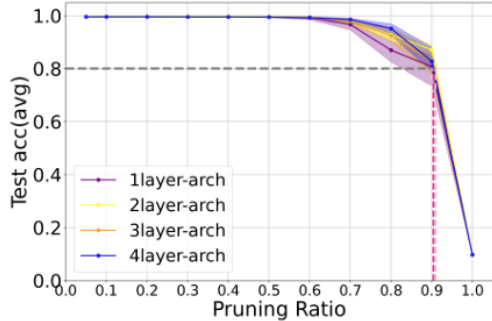
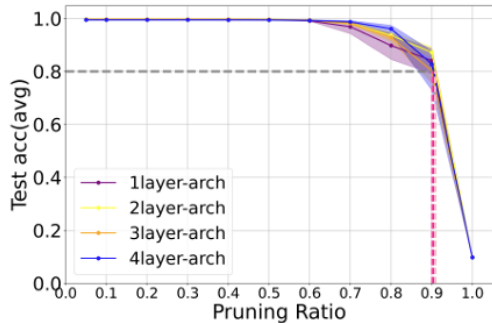
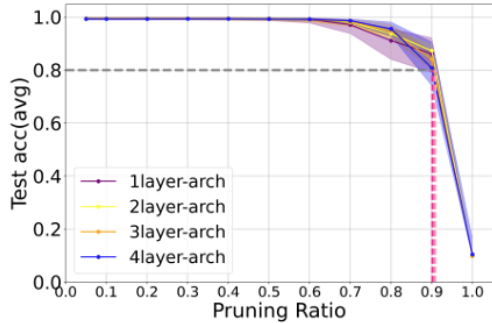
CIFAR-10



CIFAR-100



MTL

 $\text{MoE}(k=2)$ 
$$\text{MoE}(k=4)$$
 $\text{MoE}(k=8)$ 

Number of layers	Ours (w/ 1000)	Ours (w/ 100)	Ours (w/ 10)	Ours (w/ 1)	Ours (w/ 0.1)	Baseline
1	0.85	0.80	0.75	0.70	0.65	0.60
2	0.88	0.83	0.78	0.73	0.68	0.63
3	0.90	0.85	0.80	0.75	0.70	0.65
4	0.92	0.87	0.82	0.77	0.72	0.67
5	0.93	0.88	0.83	0.78	0.73	0.68
6	0.94	0.89	0.84	0.79	0.74	0.69
7	0.94	0.90	0.85	0.80	0.75	0.70
8	0.95	0.90	0.85	0.80	0.75	0.70
9	0.95	0.90	0.85	0.80	0.75	0.70
10	0.95	0.90	0.85	0.80	0.75	0.70

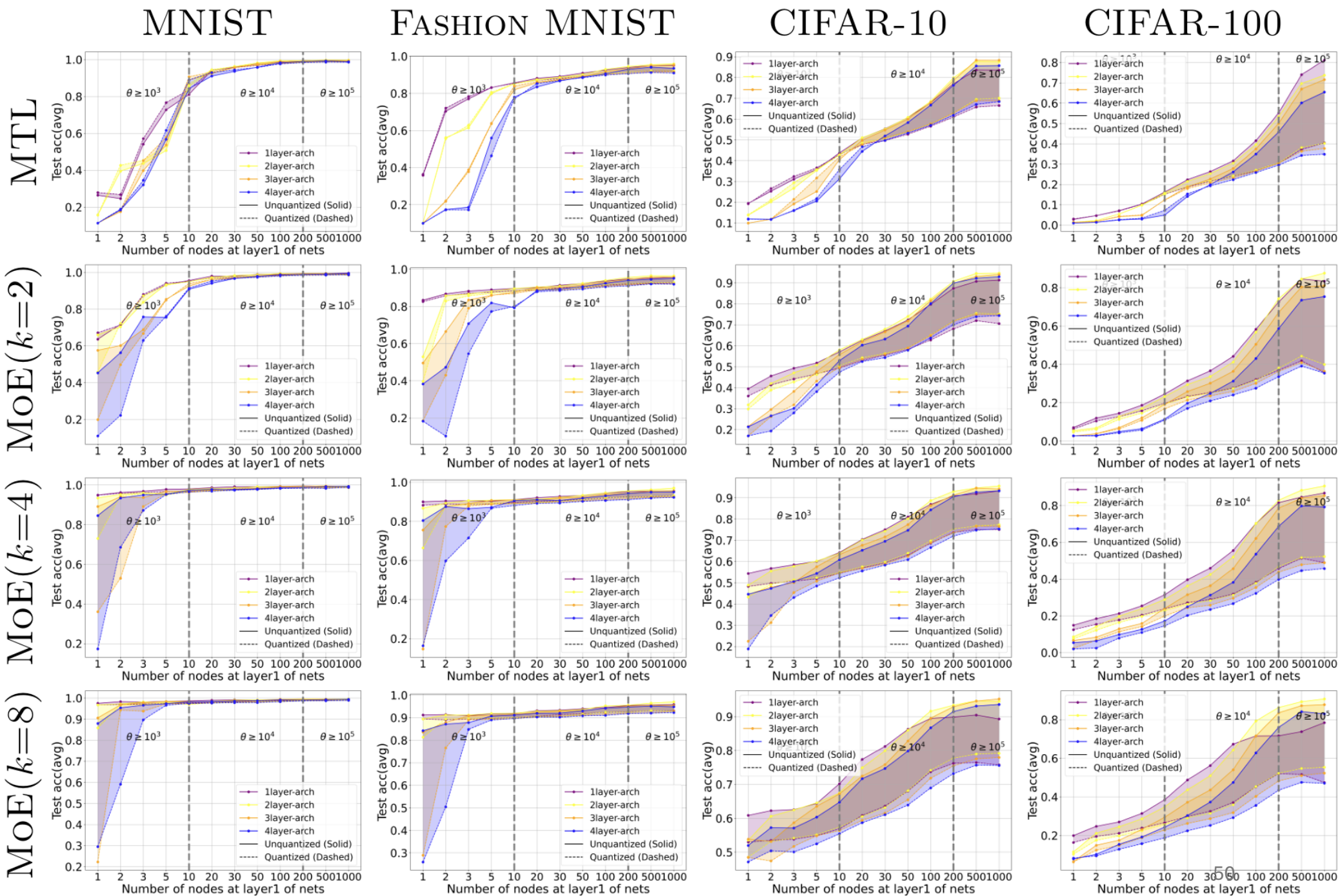
[illegible][illegible]

Test acc(avg)

[illegible]

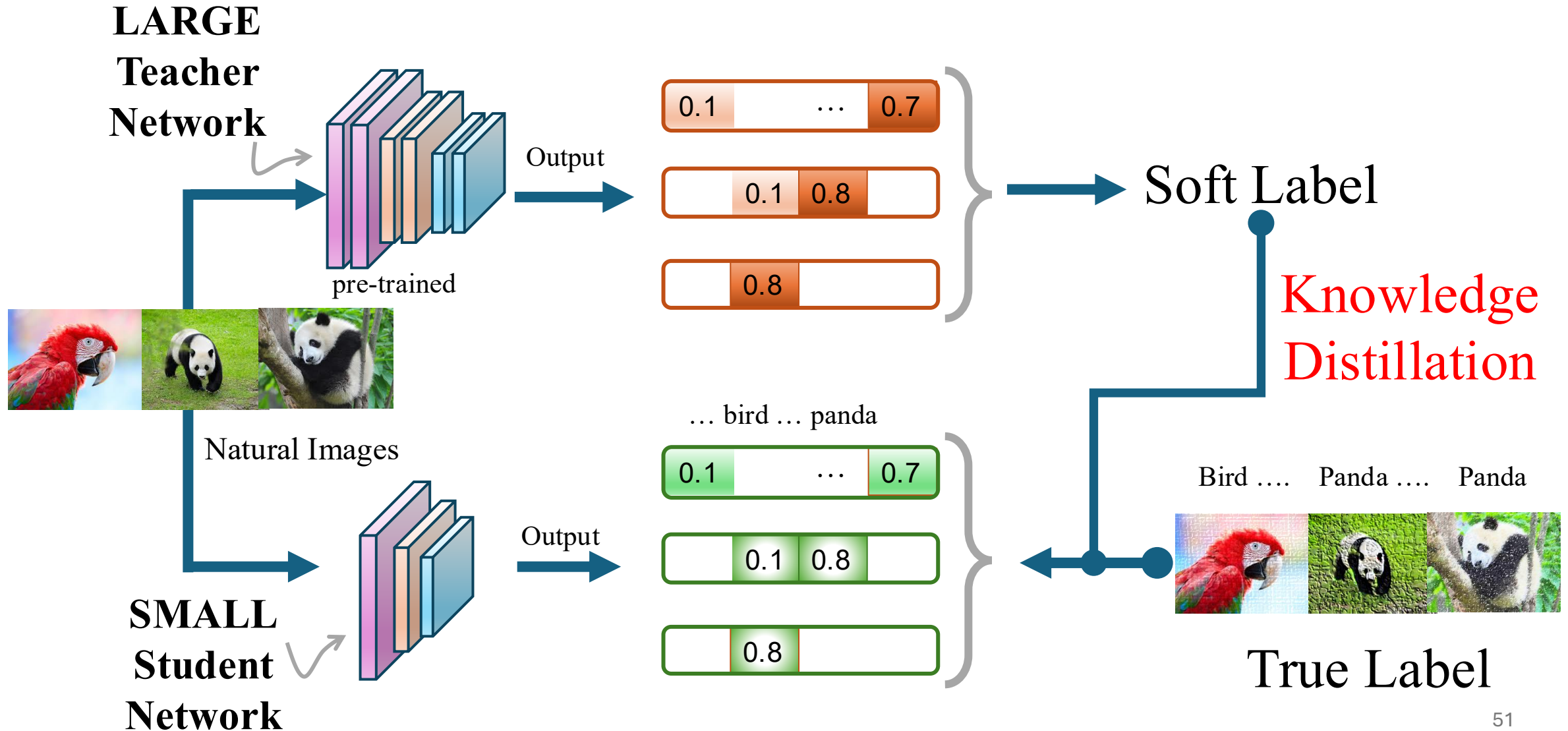
Pruning Ratio	arch1 (Yellow)	arch2 (Orange)	arch3 (Blue)	arch4 (Green)
0.3	~0.85	~0.85	~0.75	~0.75
0.4	~0.85	~0.85	~0.75	~0.75
0.5	~0.85	~0.85	-	-

MTL and MoE Quantization

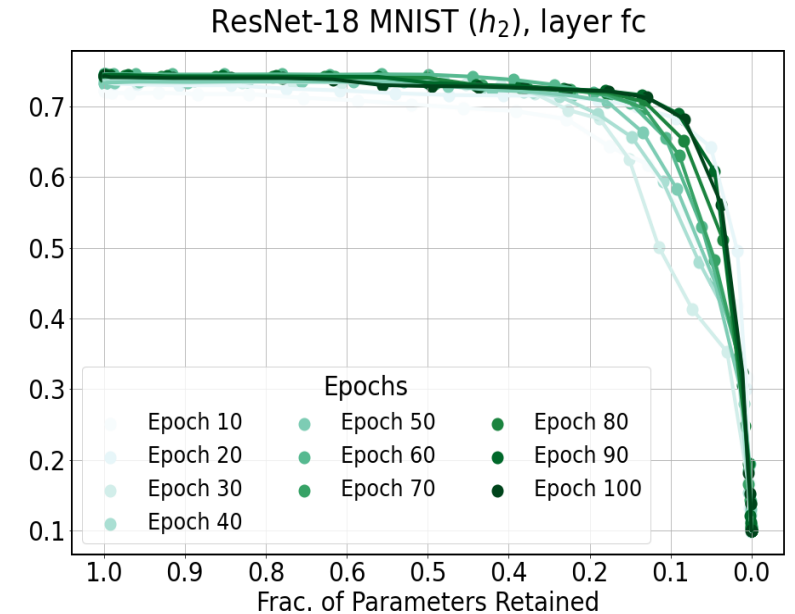
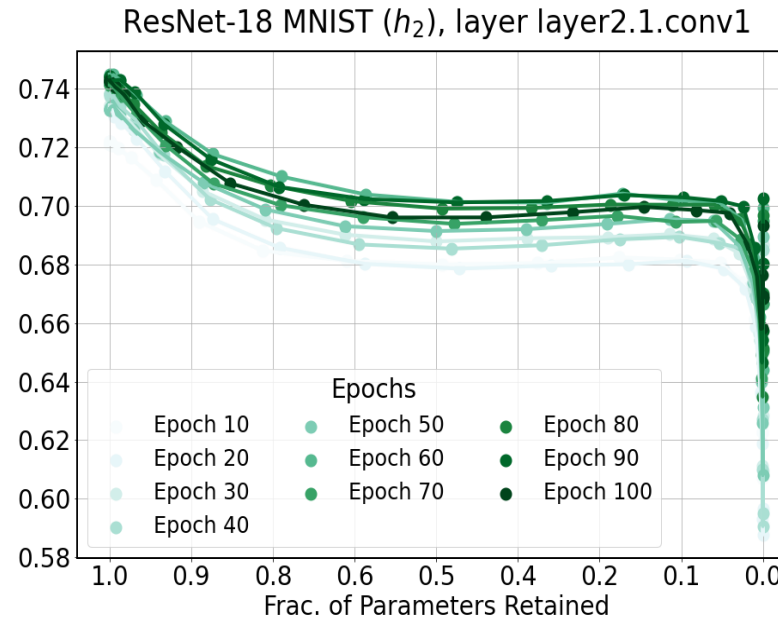
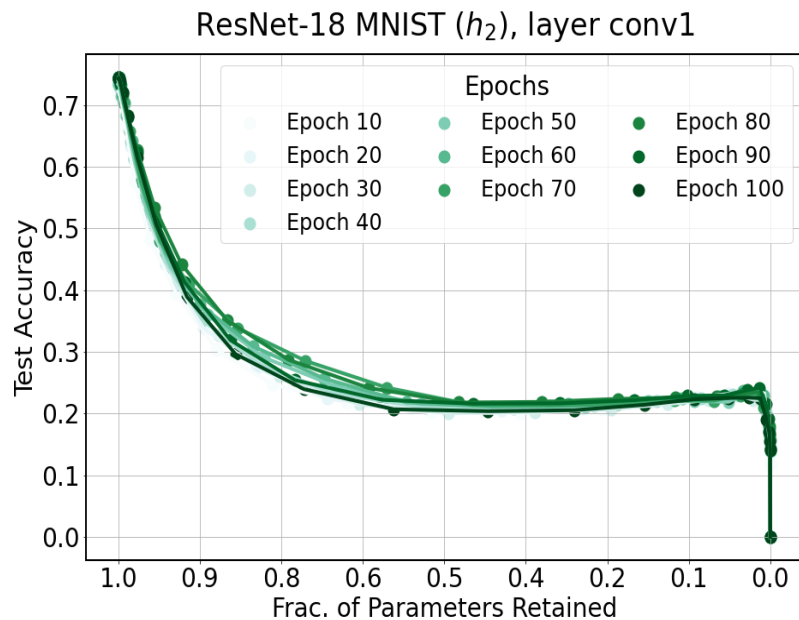
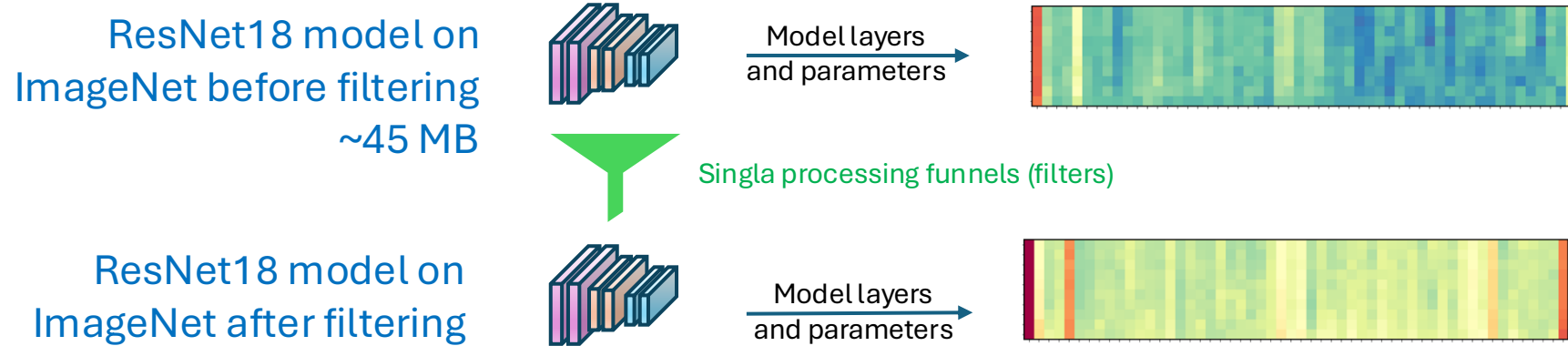


Knowledge Distillation and Model Robustness

Zhenyu, Ojha et al (2024), ICONIP



Compression Using “Parameter Score” -> “Importance Score”



Optimization Techniques for Edge Deployment

Beyond quantization, pruning, and distillation, several other methods can enhance the deployment of large models on edge and IoT devices:



Neural Architecture Search (NAS)

Automates the design of efficient neural networks tailored for specific hardware constraints.



Hardware-Aware Training (HAT)

Incorporates hardware characteristics during training to optimize models for particular devices.



Sparsity-Inducing Regularization

Encourages models to develop sparse weight matrices during training, leading to smaller and faster models.



Winograd Transformations

Optimizes convolutional operations to reduce the number of multiplications, improving inference speed.



Fast Fourier Transforms (FFTs)

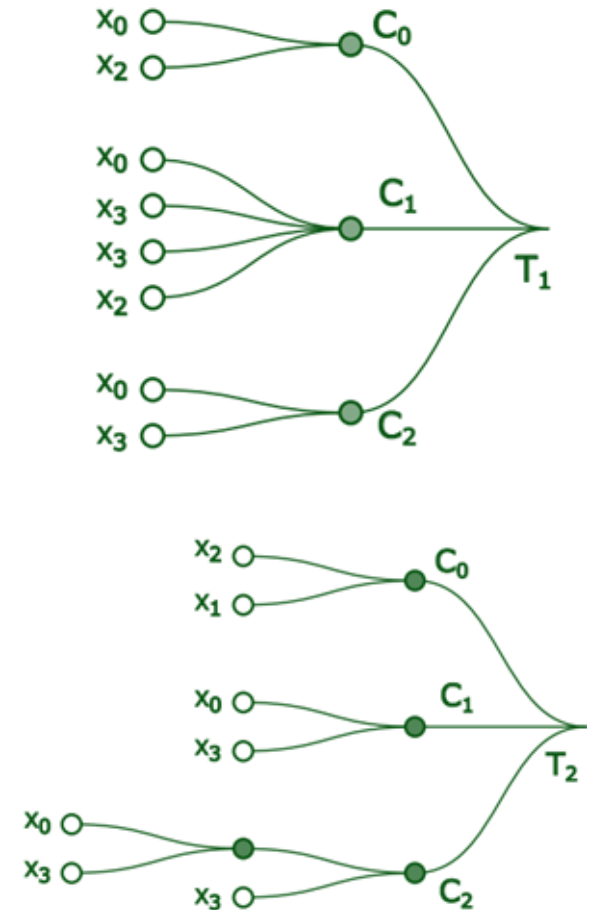
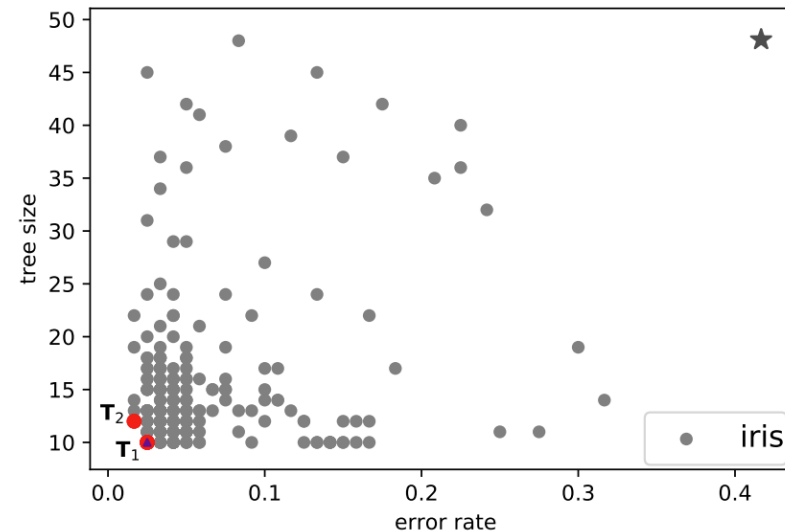
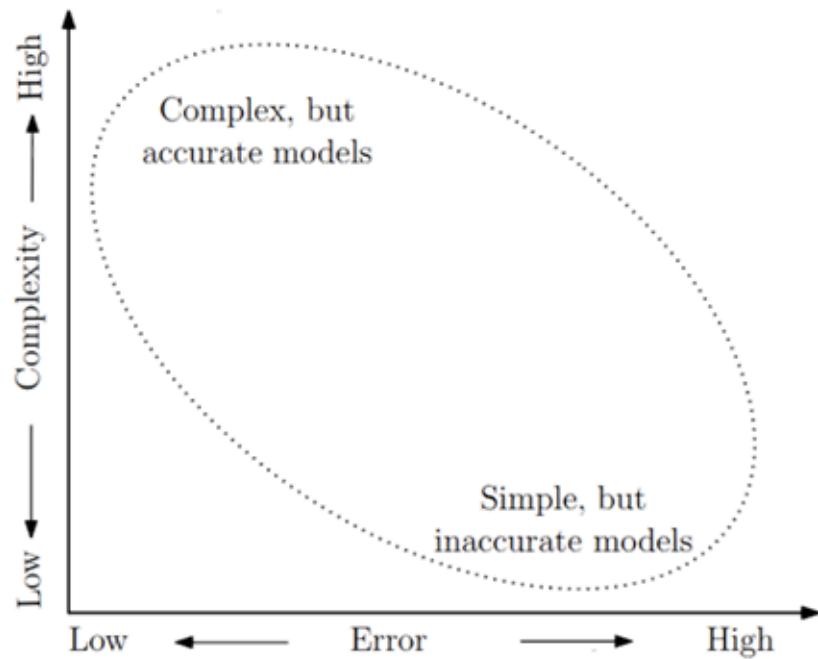
Efficiently performs convolutions in the frequency domain, particularly useful for large kernel sizes.

Part 4

Search for Efficient AI Models

Neural Architecture Search (NAS)

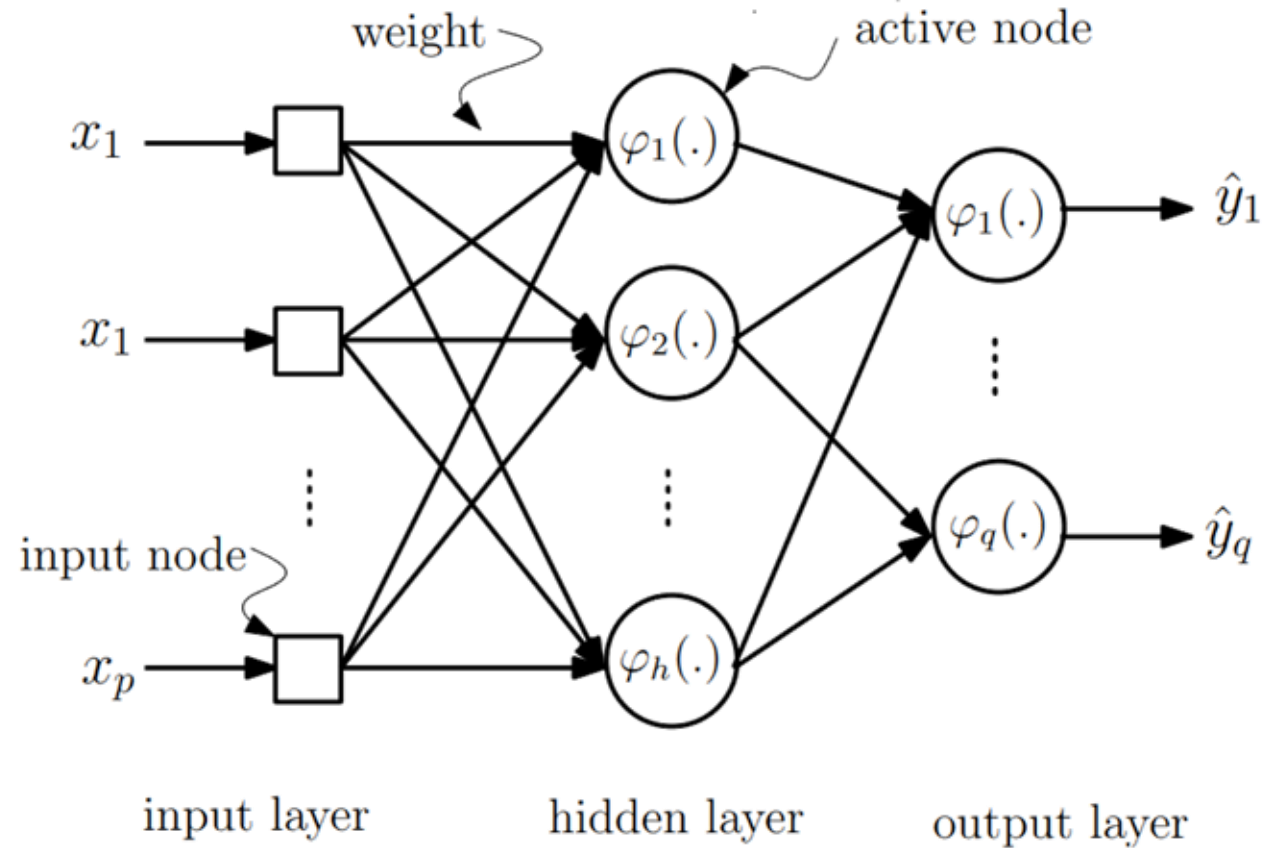
Catering the need of specialized problem and hardware



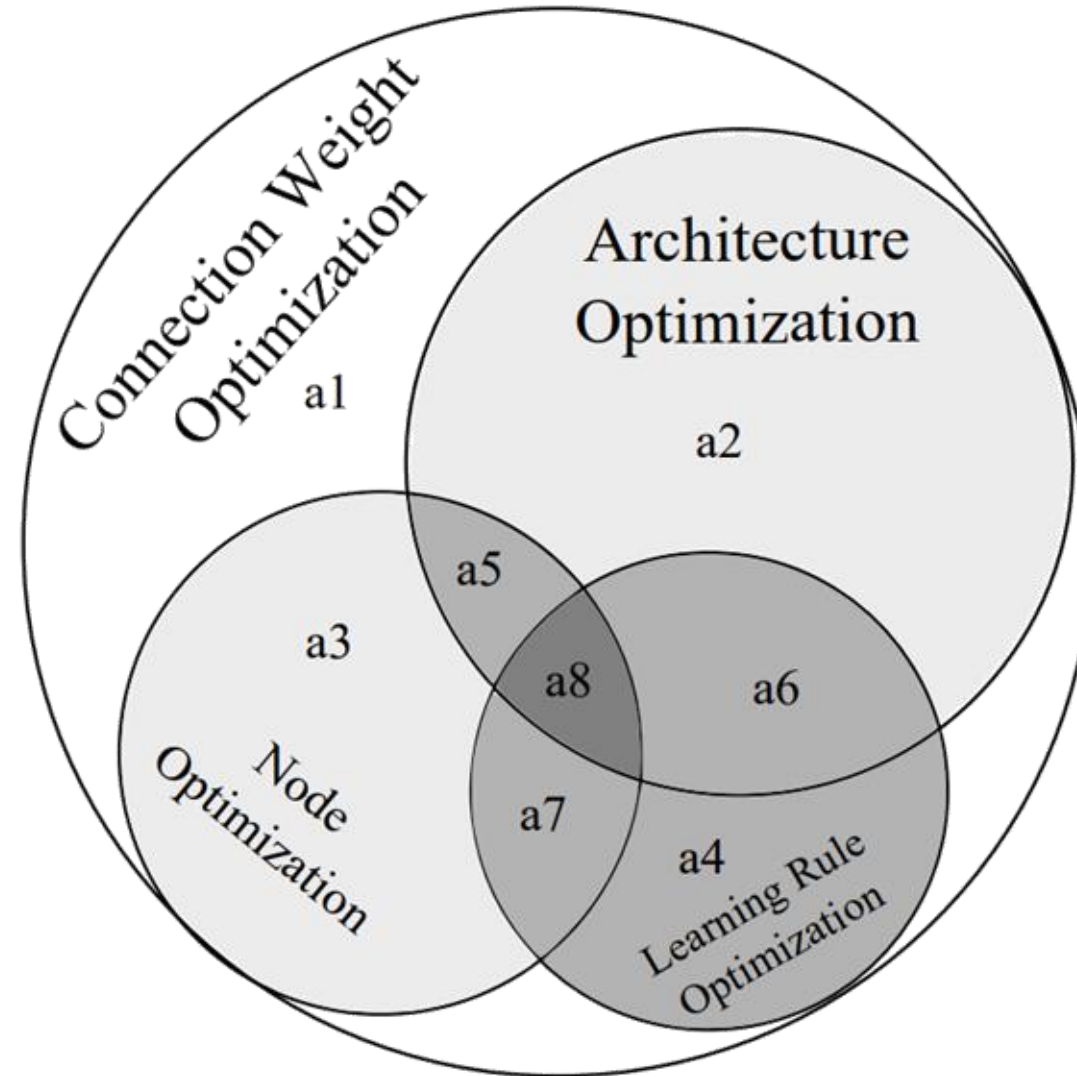
Neural Networks

NN components:

- Inputs
- Weights
- Architecture
- Activation functions
- Learning algorithms



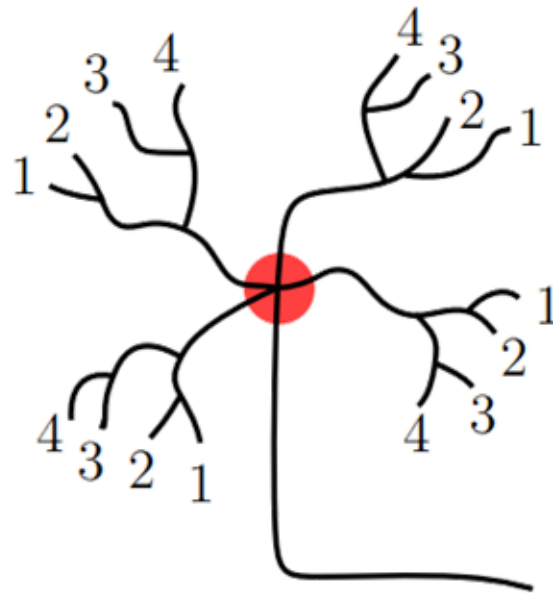
What Could be optimized?



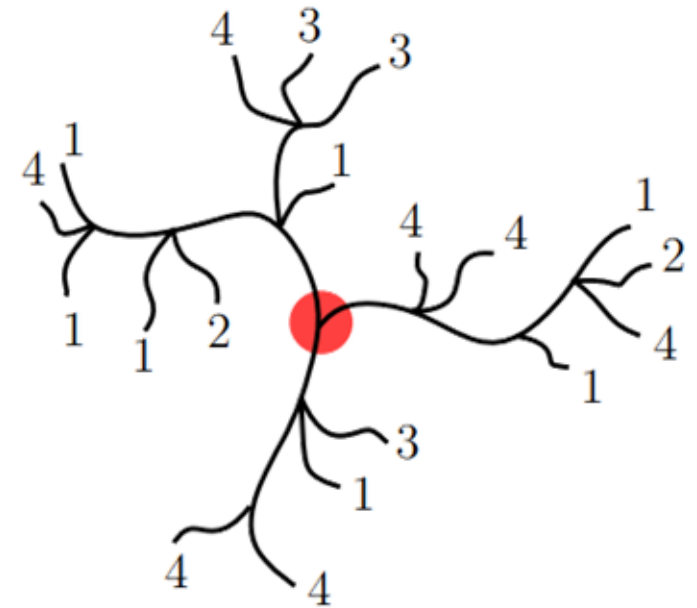
Plausible Biological Inspiration



Travis et al. (2005)



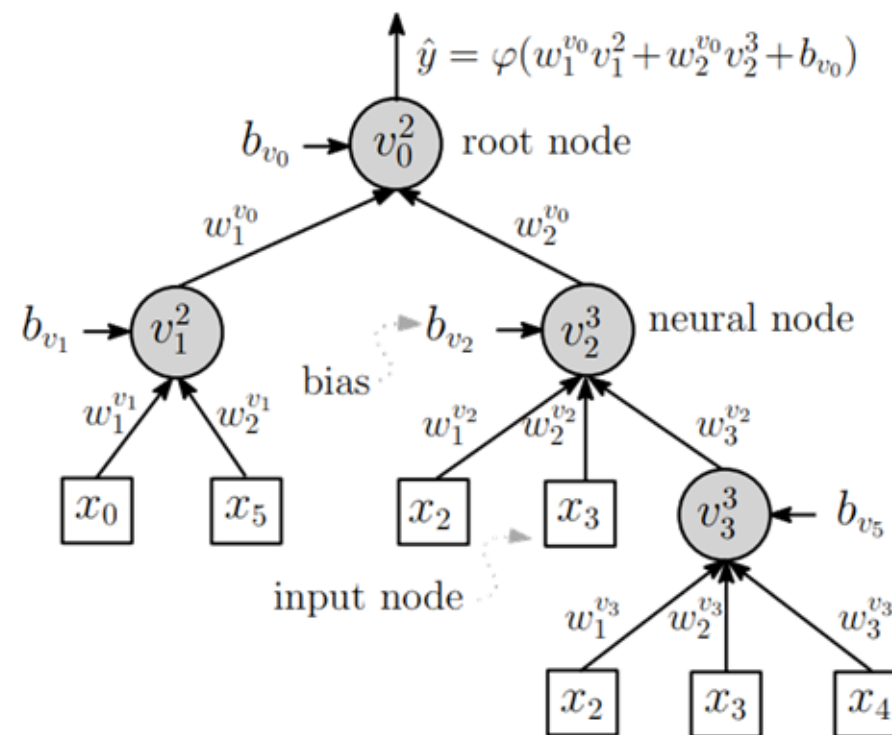
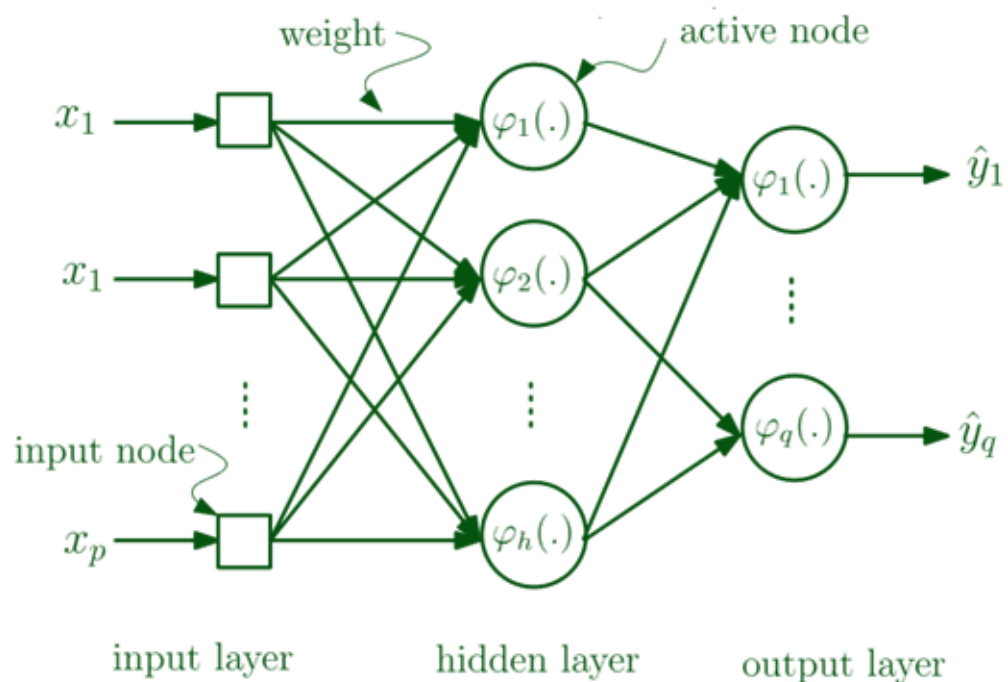
Jones and Kording (2021)



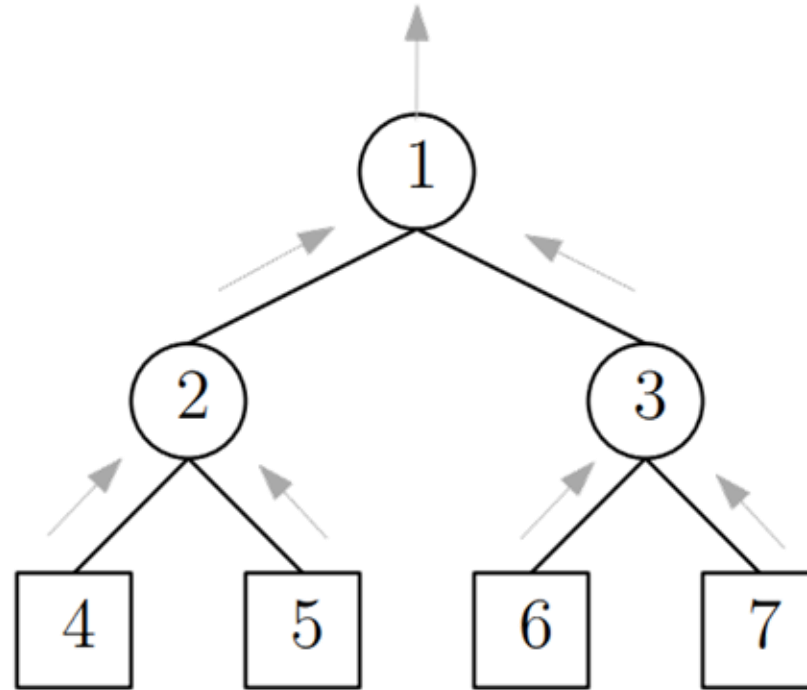
Ojha and Nicosia (2022)

Neural Tree

Neural Networks Architecture Search



Neural Computation

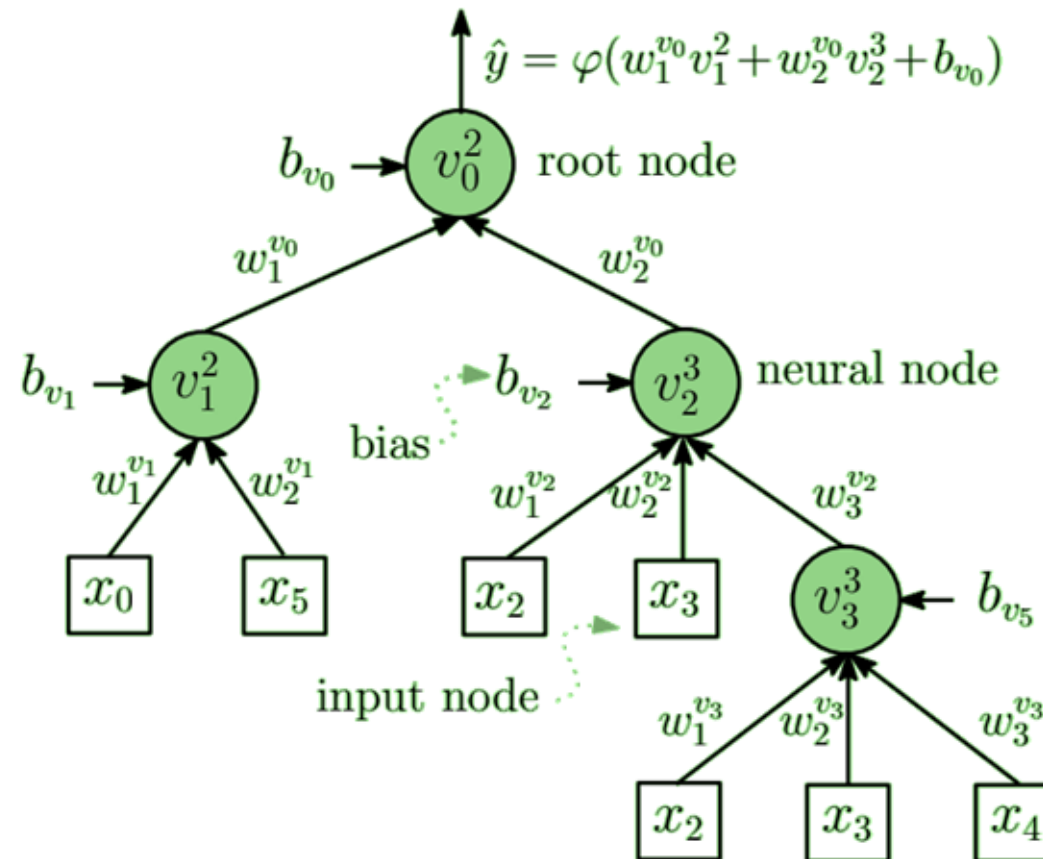


$[((4\ 5) \rightarrow 2) \quad ((6\ 7) \rightarrow 3)] \rightarrow 1$

forward pass: post-order

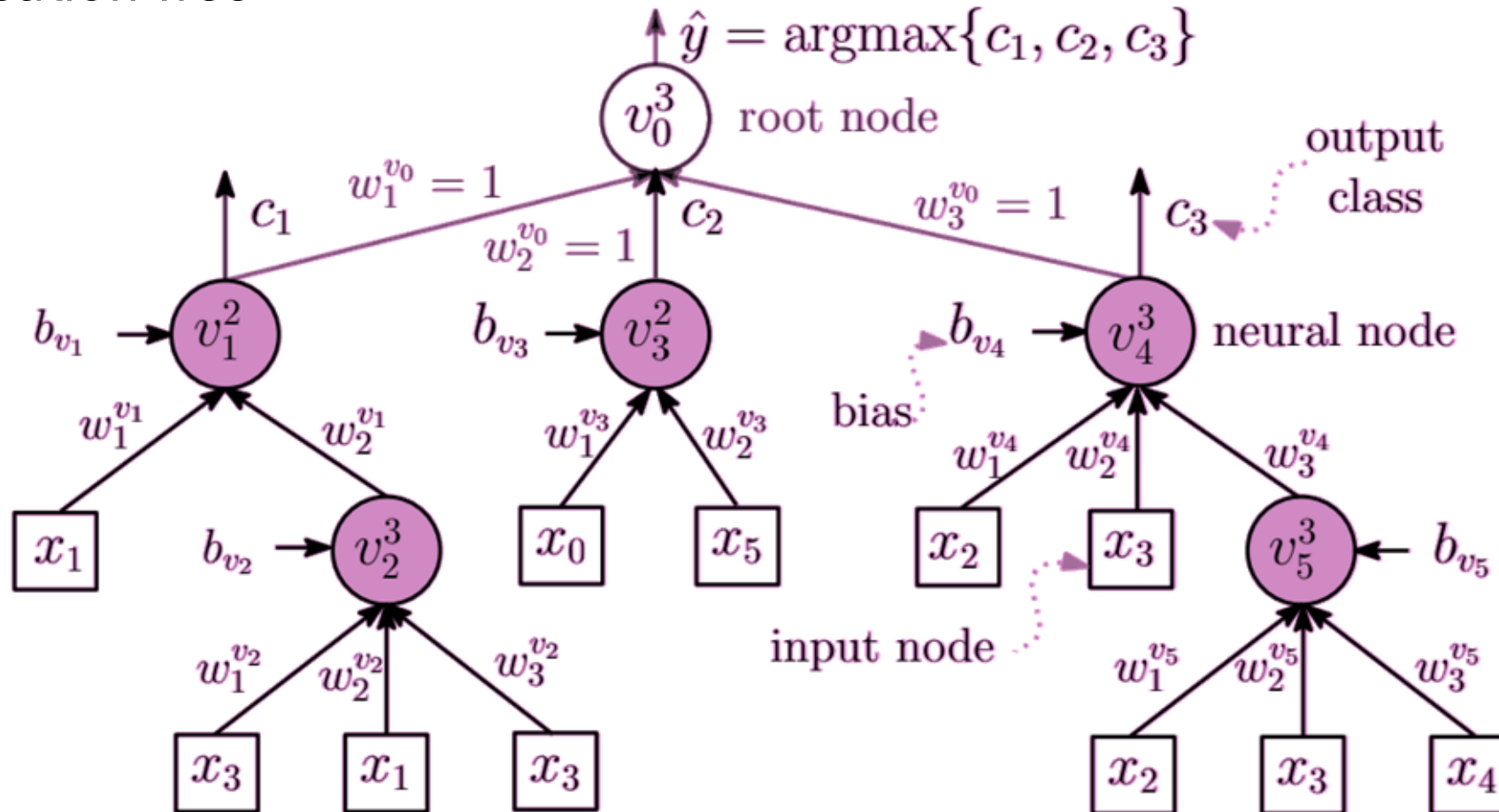
Types of Neural Tree

Regression Tree



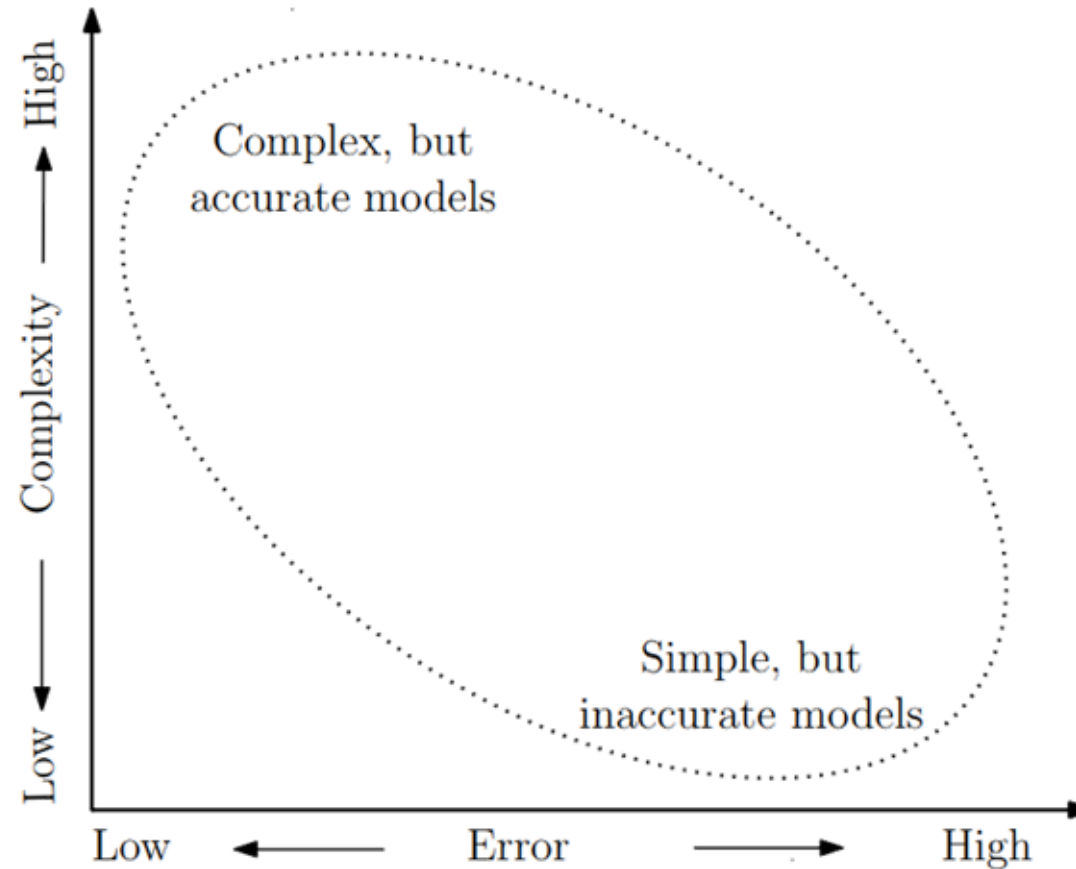
Types of Neural Tree

Classification Tree



Neural Architecture Search

Trade-offs

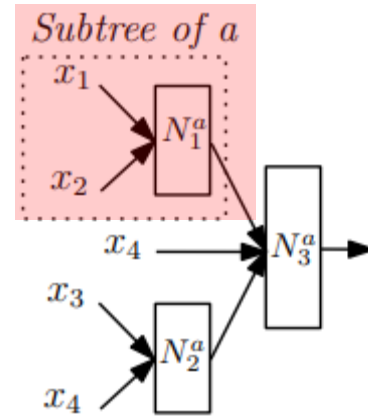


Neural Architecture Search

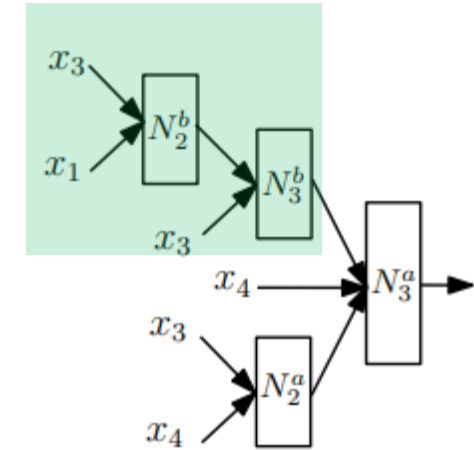
Trade-offs

Multiobjective
Genetic Programming
Crossover

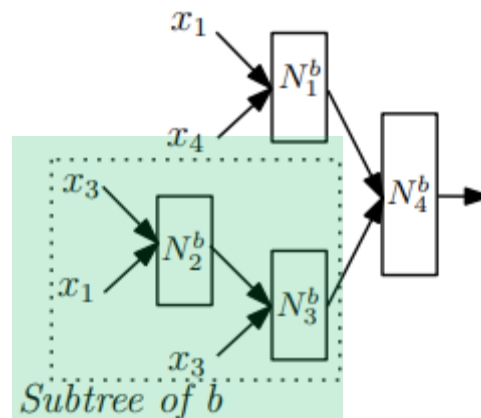
Ojha et al (2017), *IEEE Trans. Fuzzy Systems*



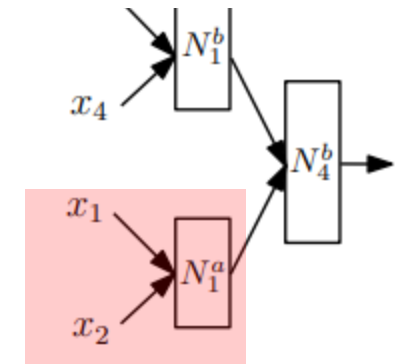
Parent tree: a



Child tree: c



Parent tree: b



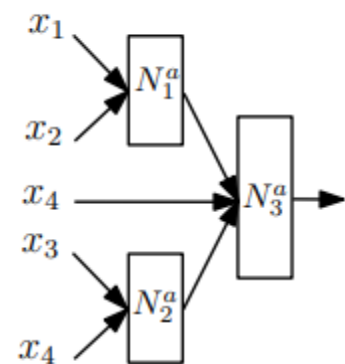
Child tree: d

Neural Architecture Search

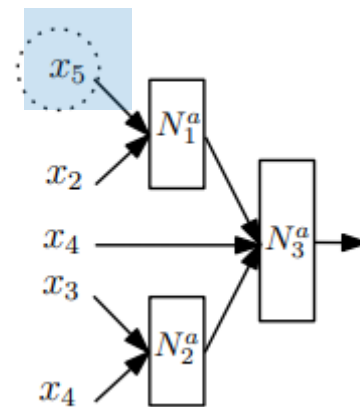
Trade-offs

Multiobjective
Genetic Programming
Mutation

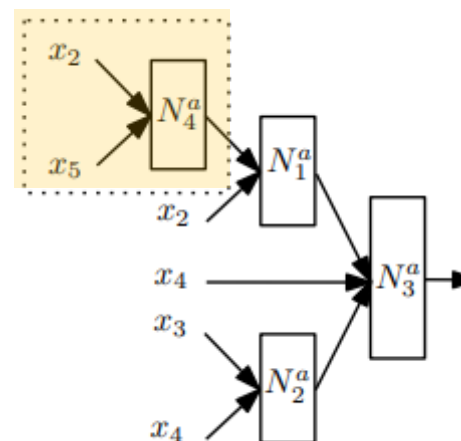
Ojha et al (2017), *IEEE Trans. Fuzzy Systems*



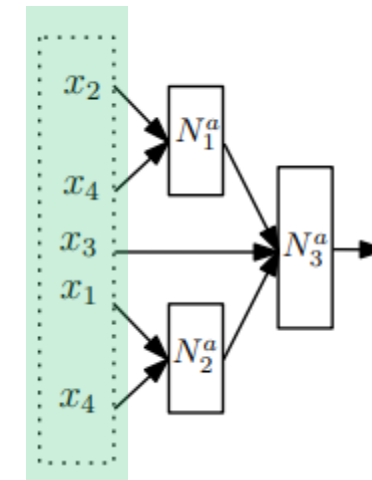
Parent tree



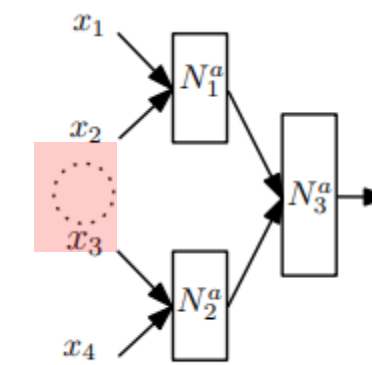
Single leaf mutation



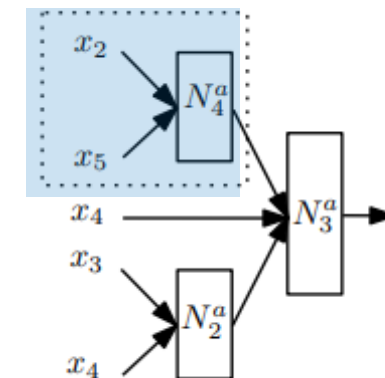
A subtree insertion



All leaves mutation



A subtree deletion

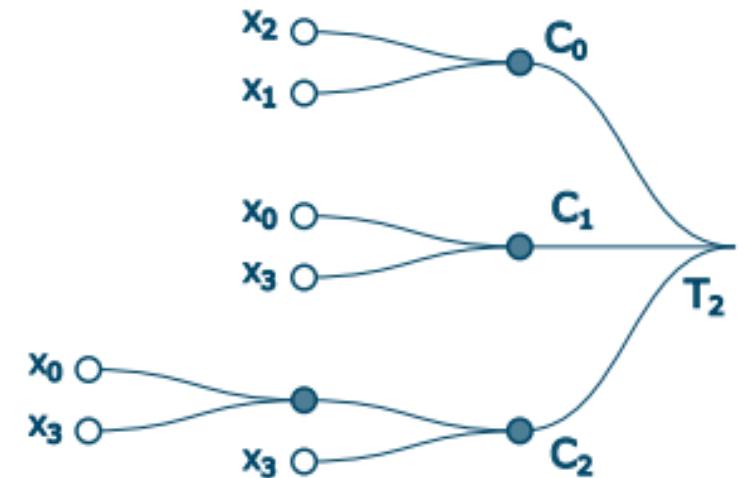
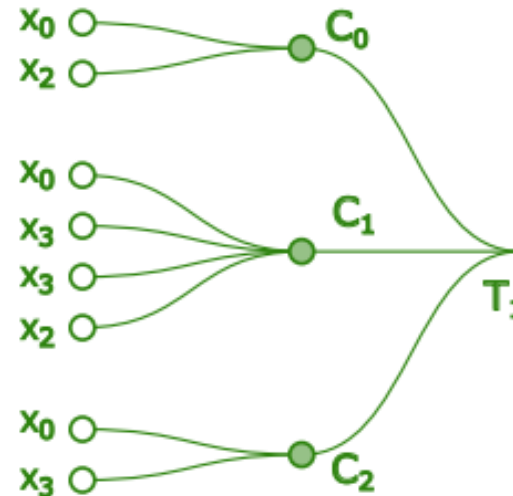
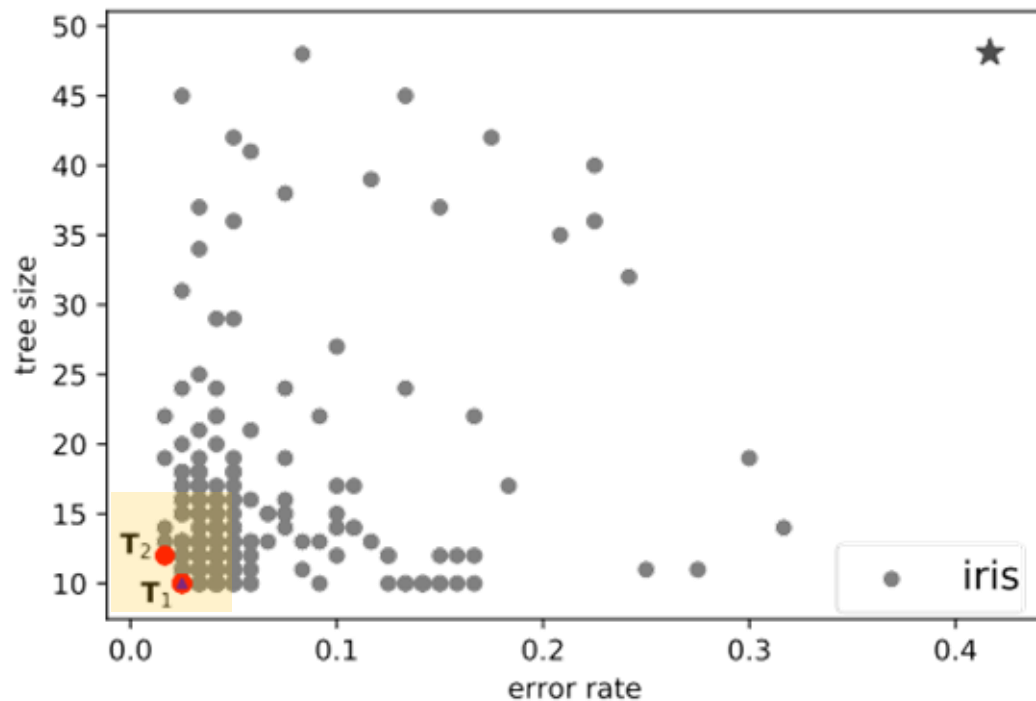


A subtree replacement

Architecture Search Trade-offs

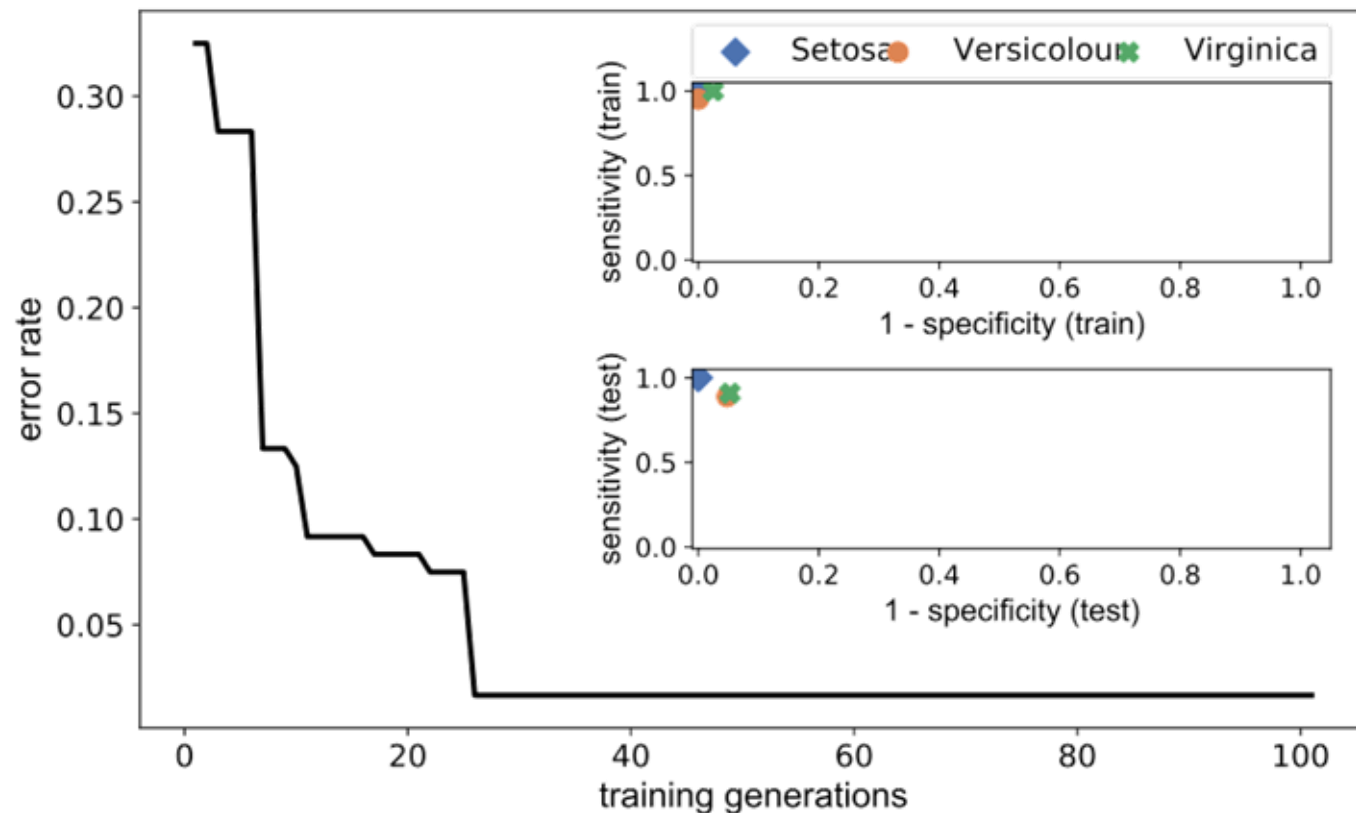
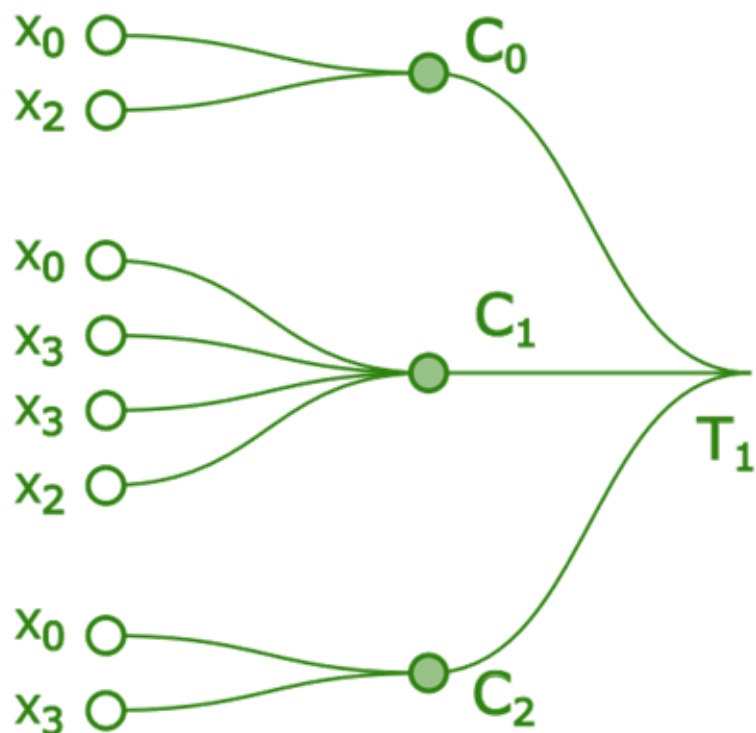
Multiobjective Genetic Programming

Selection of trees using Hypervolume indicator from a Pareto Front



Learnability of Classes

Competition between classes: TPR/Recall/Sensitivity vs FPR/(1 - Specificity)

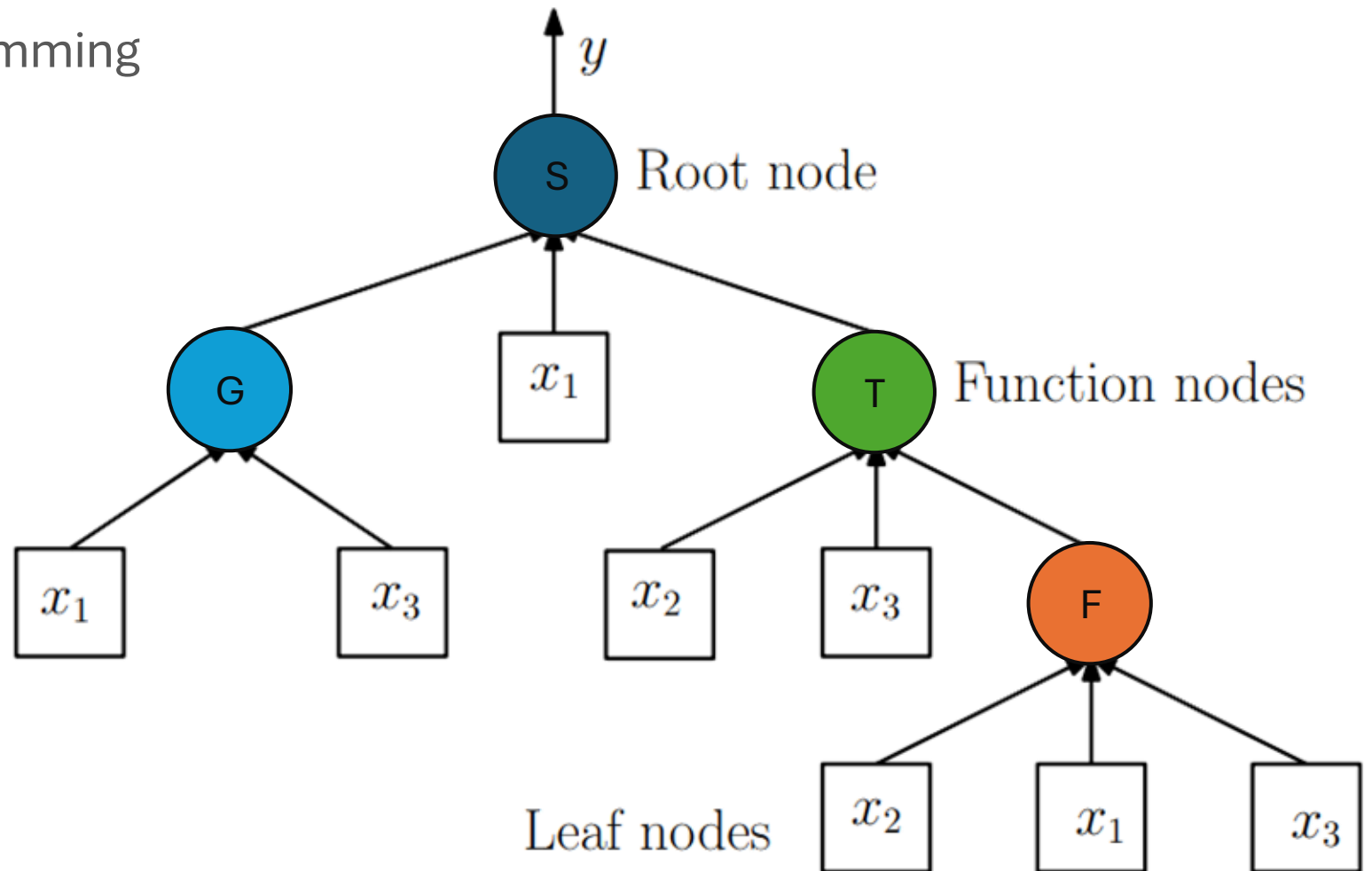


Heterogeneous Neural Tree

Multiobjective Genetic Programming

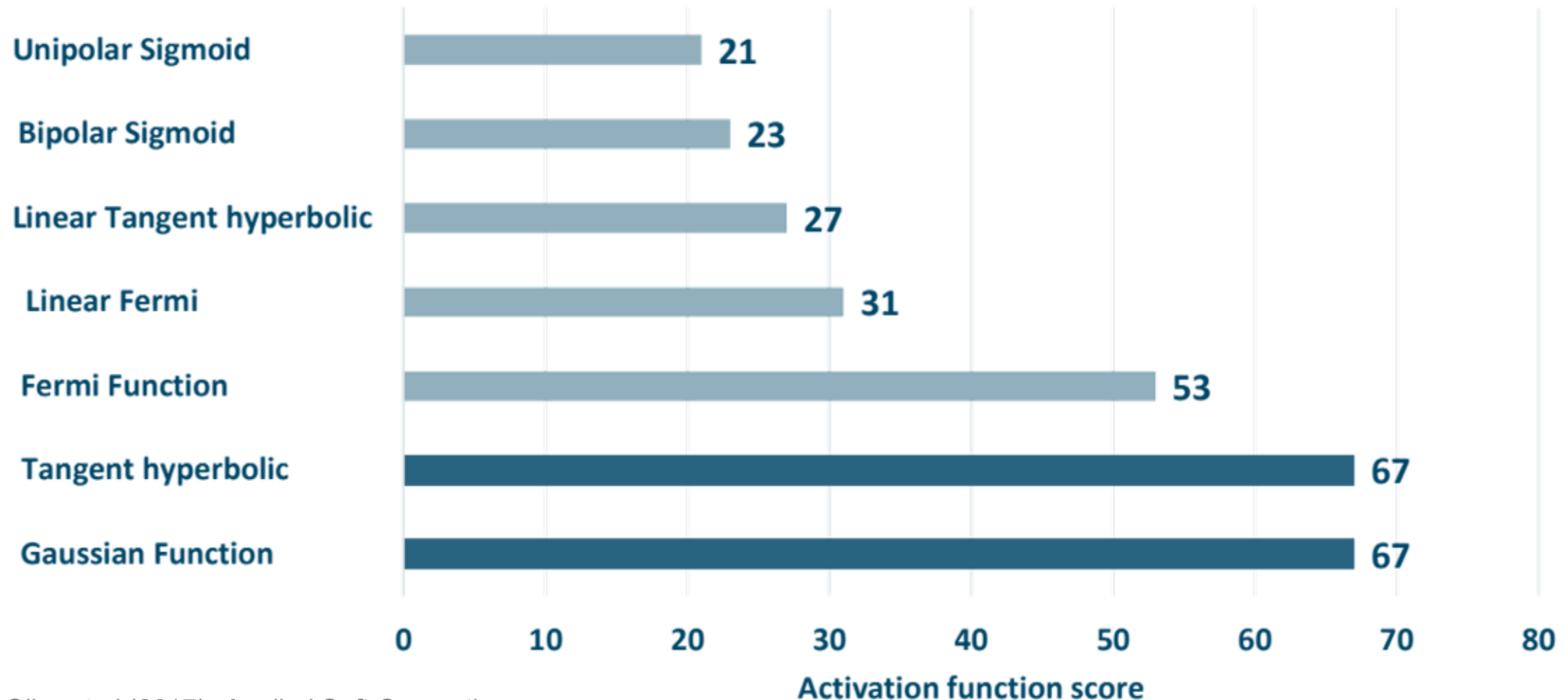
Activation Function Search

- S – Sigmoid
- G – Gaussian
- T – Tanh
- F – Fermi



Activation Function Performance

Higher values are better



Neural Modelling

Sparse Neural Tree

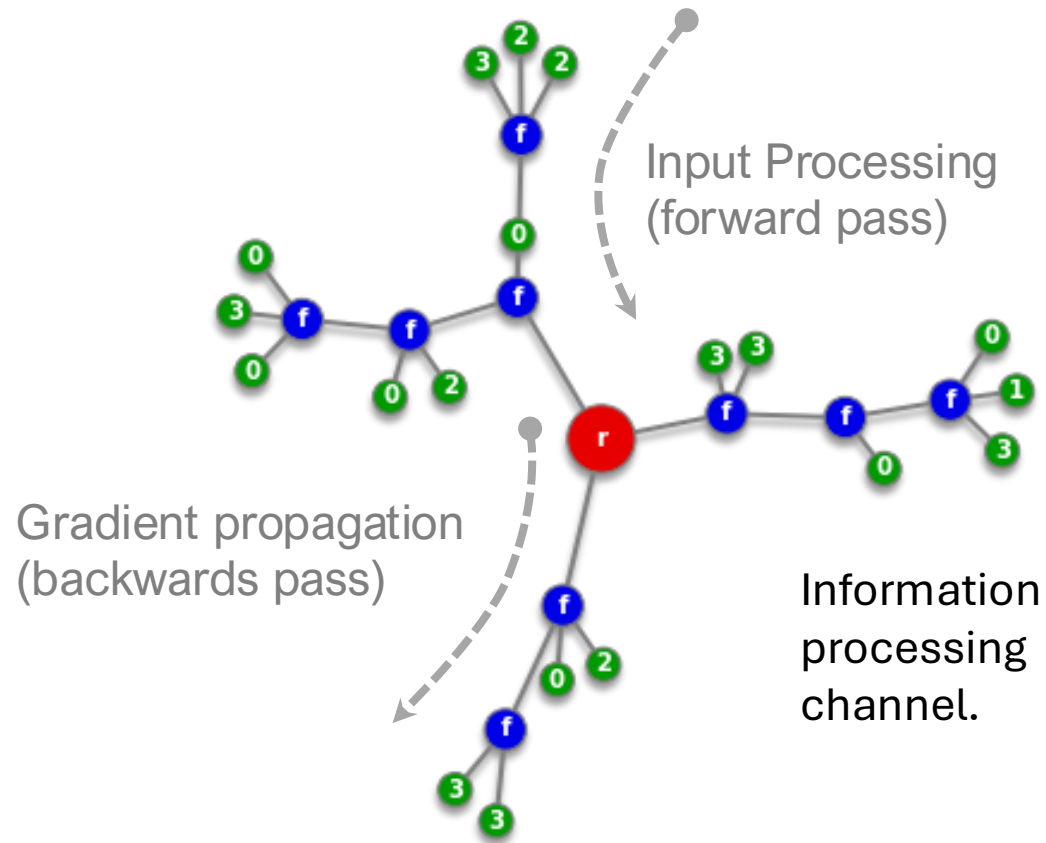
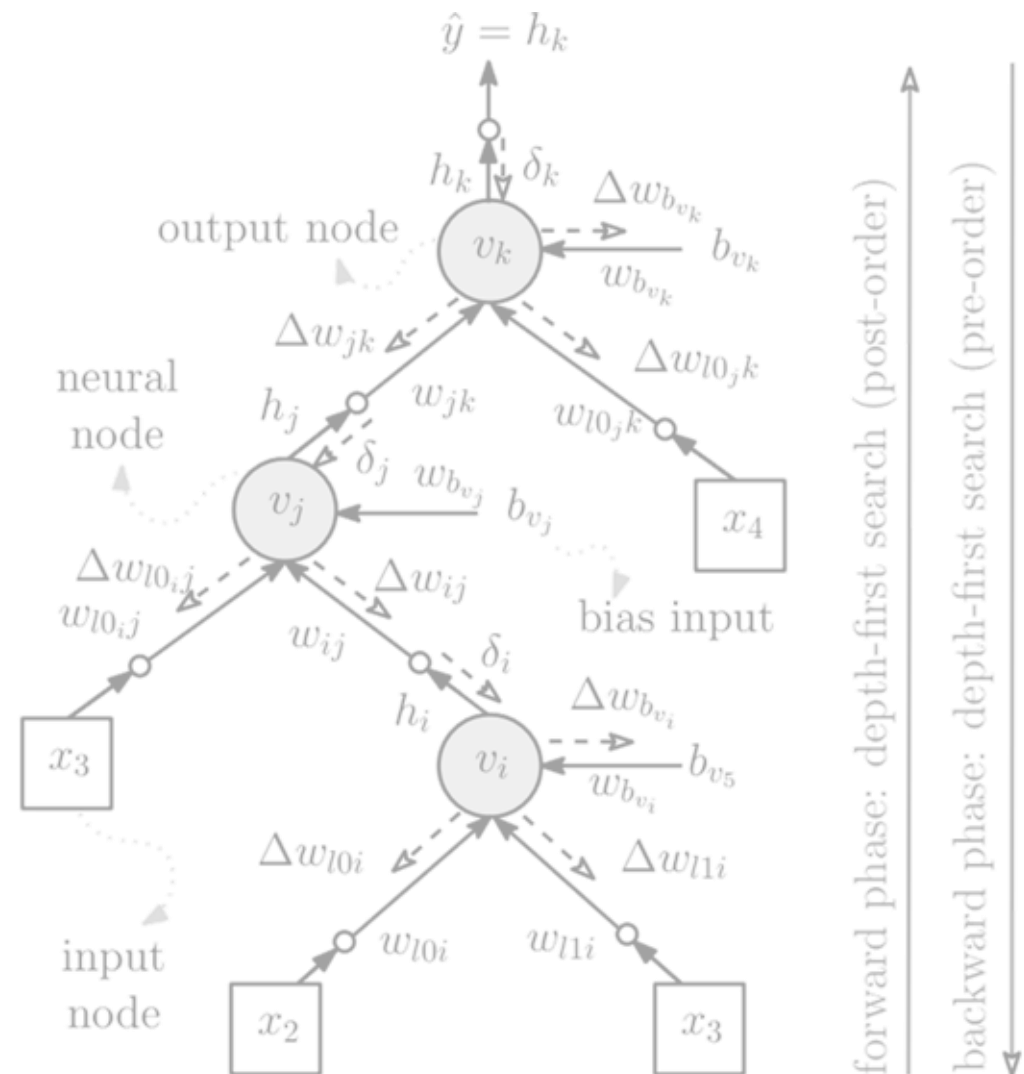
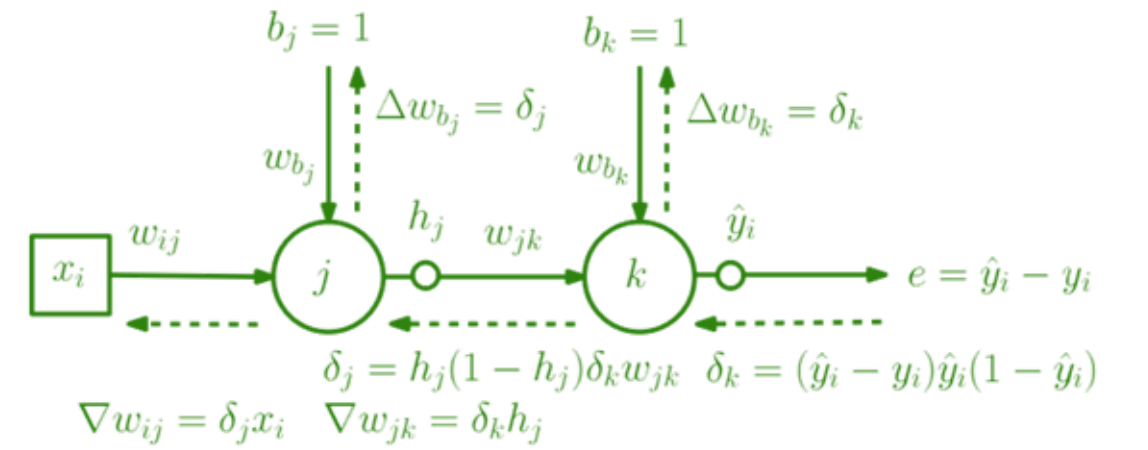
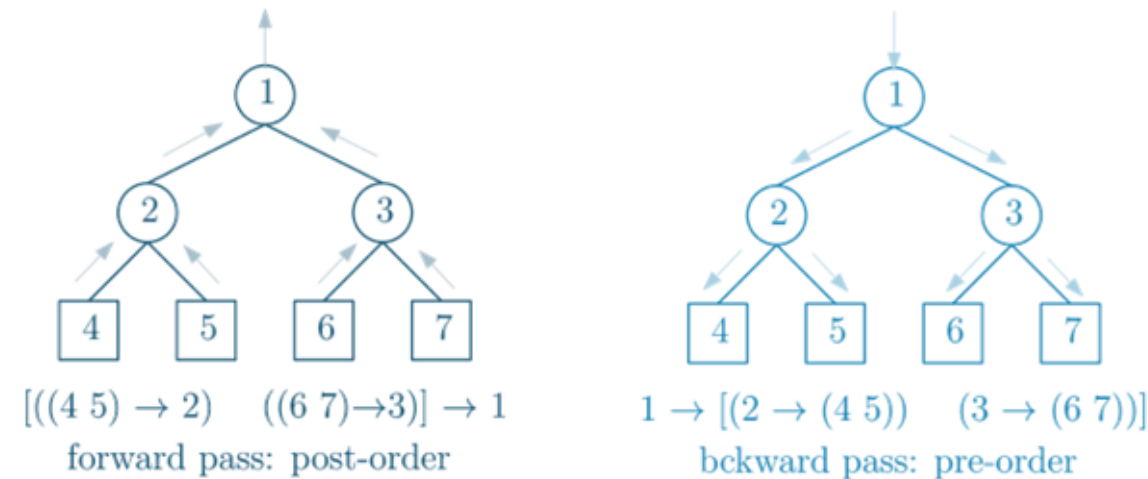
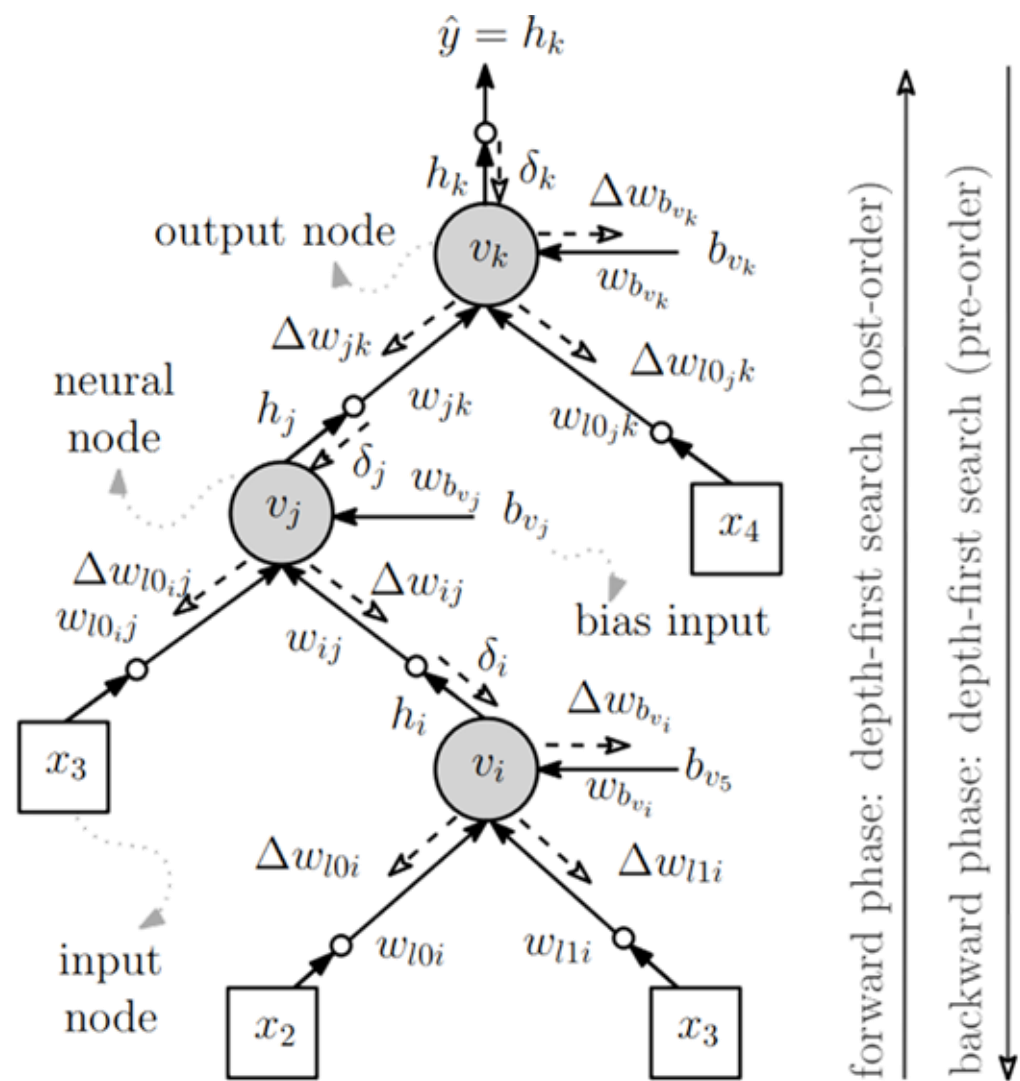


Fig A. Forward pass and gradient backpropagation

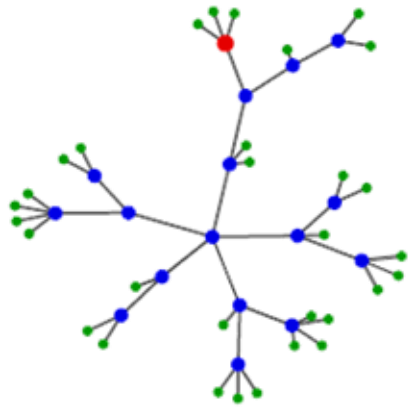


Backpropagation Neural Tree

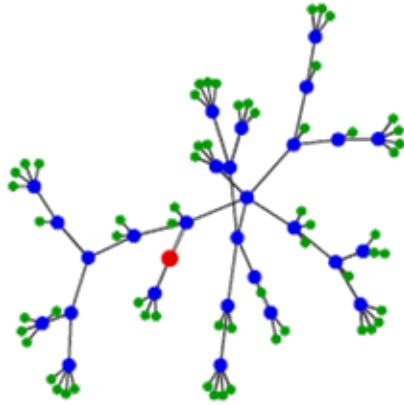


Backpropagation Neural Tree

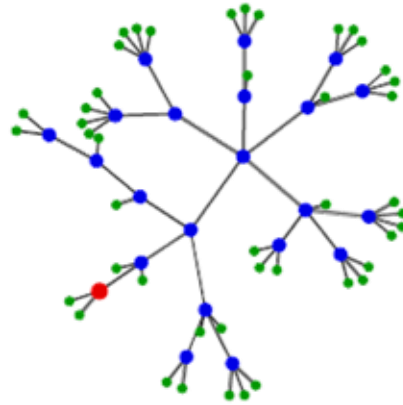
Regression results



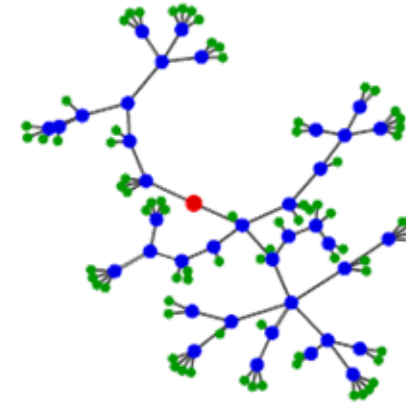
(a) baseball (.85, 48)



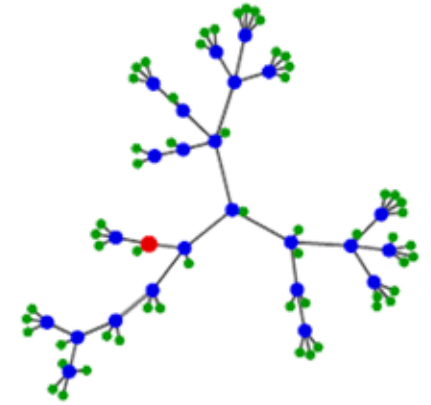
(b) dee (.89, 89)



(c) diabetes (.63, 67)



(d) friedman (.95, 116)



(e) mpg6 (.9, 82)

Algorithm	Bas	Dee	Dia	Frd	Mpg	Avg Acc	Avg Weights
BNeuralT	0.665	0.837	0.492	0.776	0.867	0.727	152
MLP	0.721	0.829	0.49	0.943	0.874	0.772	1041

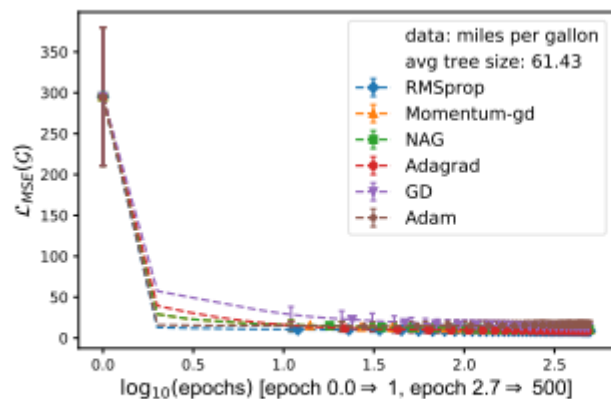
Backpropagation Neural Tree

Regression results

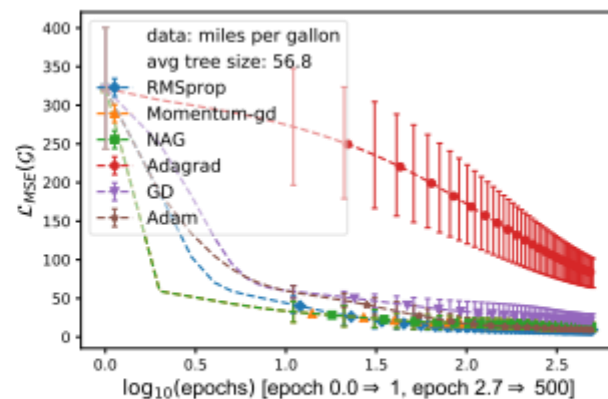
- BNeuralT used only 14.6% of MLP
- Accuracy differs only 5.8% lower than the best MLP result

Neural Tree vs Neural Networks

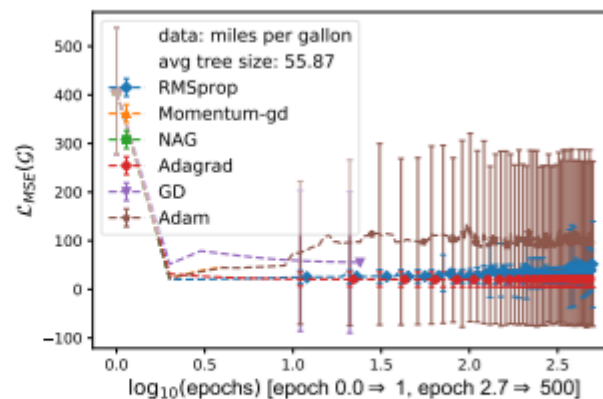
Regression Problems



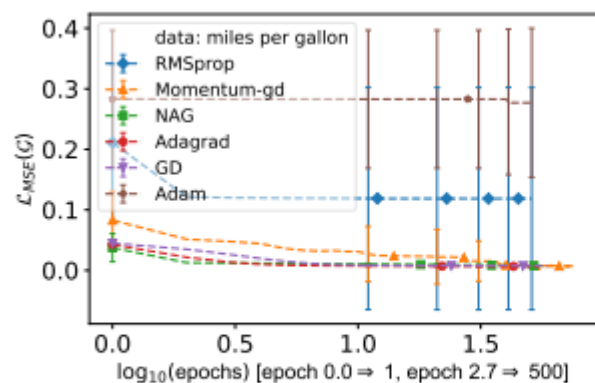
(g) BNeuralT: Sigmoid, $\eta = 0.1$



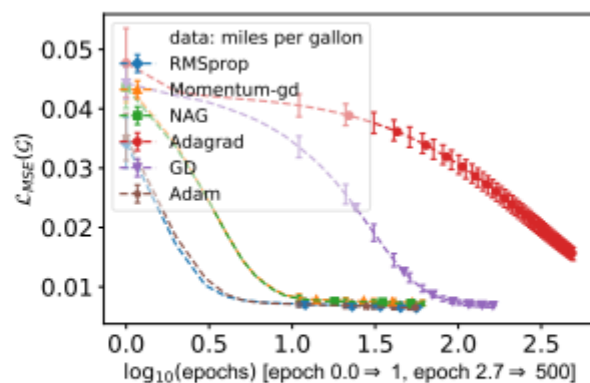
(h) BNeuralT: Sigmoid, $\eta = \text{default}$



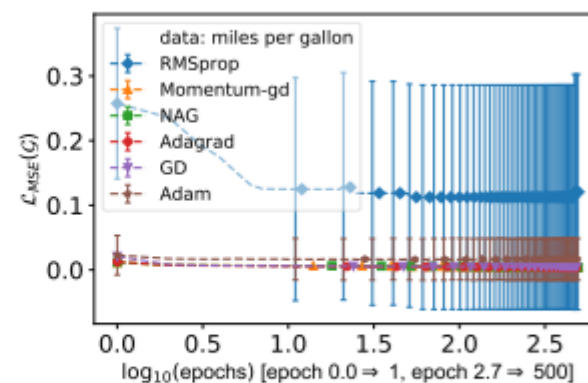
(i) BNeuralT: ReLU, $\eta = 0.1$



(j) MLP: Sigmoid, $\eta = 0.1$



(k) MLP: Sigmoid, $\eta = \text{default}$

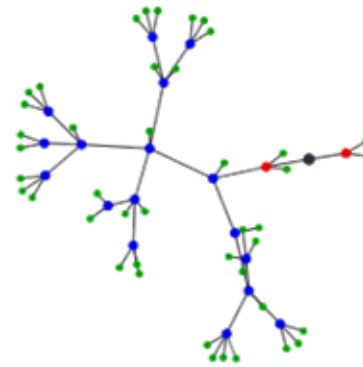


(l) MLP: ReLU, $\eta = 0.1$

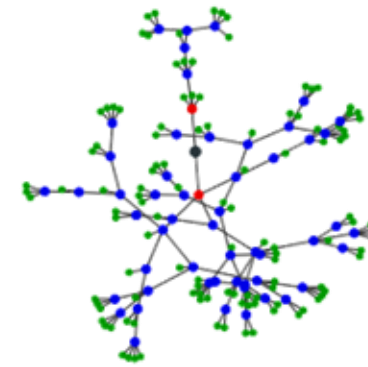
Backpropagation Neural Tree

Classification results.

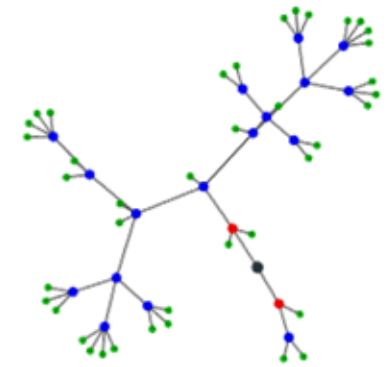
Data	BNeuralT	MLP
Aus	0.895	0.876
Hrt	0.897	0.833
Ion	0.952	0.882
Pma	0.822	0.774
Wis	0.986	0.984
Irs	0.992	0.972
Win	0.991	0.991
Vhl	0.75	0.826
Gls	0.732	0.635
Avg. Accuracy	0.891	0.863
Avg. Weights	261	1969



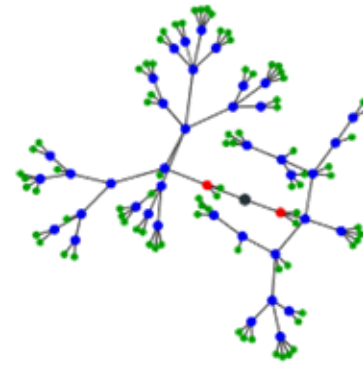
(a) australian (92%, 63)



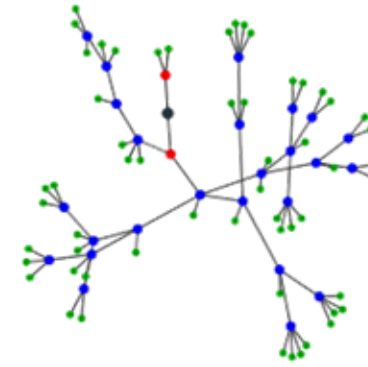
(b) heart (96%, 173)



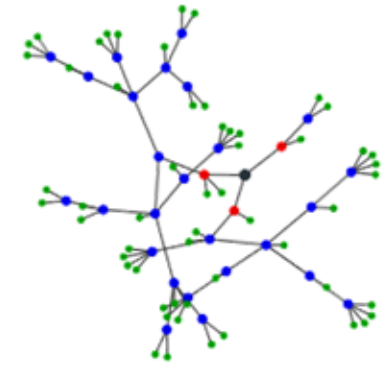
(c) ionosphere (99%, 60)



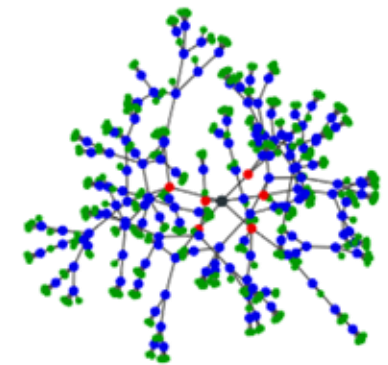
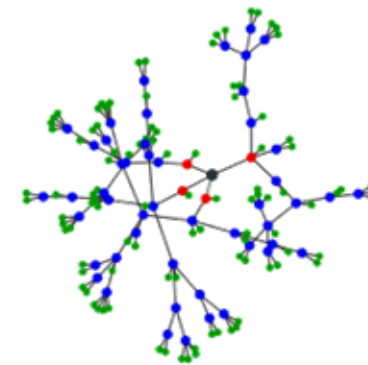
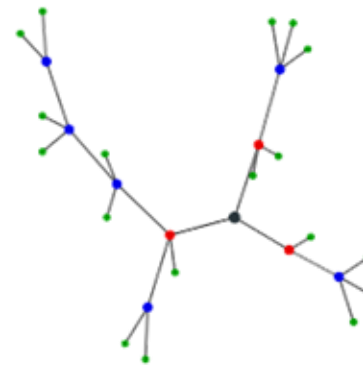
(d) pima (87%, 125)



(e) wiscosin (100%, 85)



(f) iris (100%, 86)



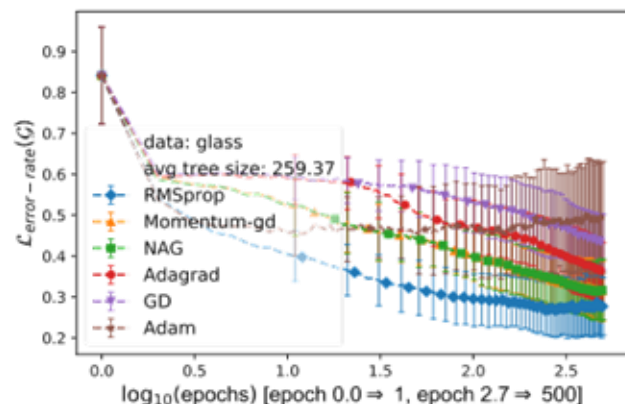
Novel type of Neural Modelling

Sparse Neural Tree

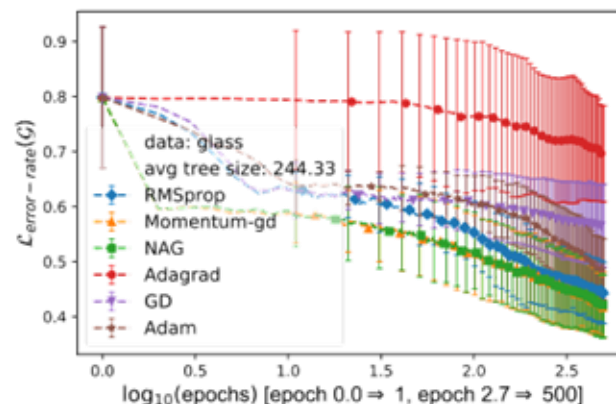
- Neural Tree used **only 13.25% parameters** of standard MLP
- Accuracy is **2.65% better than the best MLP** result

Neural Tree vs Neural Networks

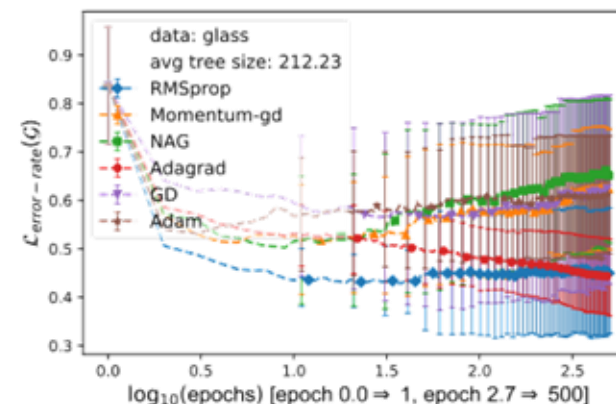
Classification Problems



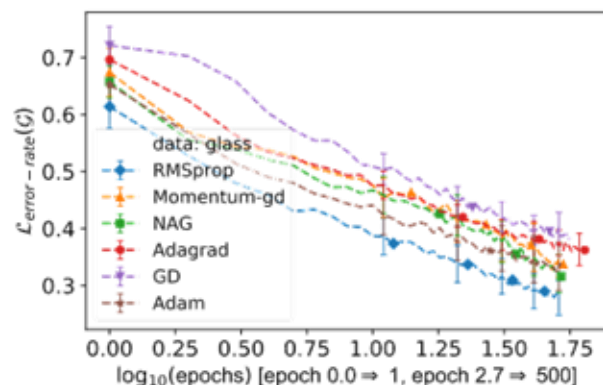
(a) BNeuralT: Sigmoid, $\eta = 0.1$



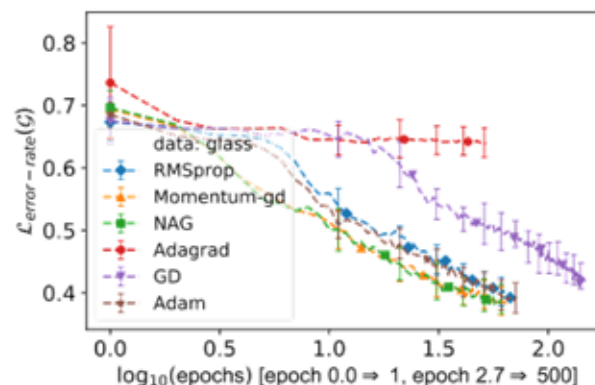
(b) BNeuralT: Sigmoid, $\eta = \text{default}$



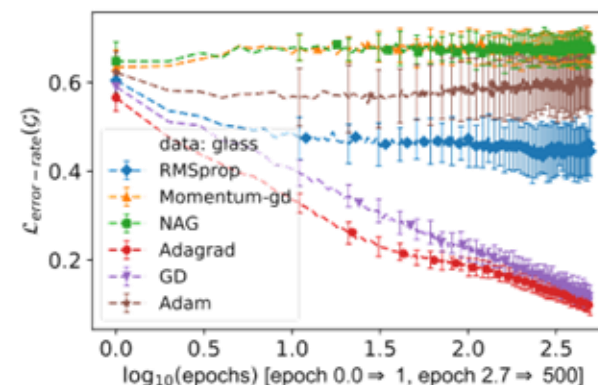
(c) BNeuralT: ReLU, $\eta = 0.1$



(d) MLP: Sigmoid, $\eta = 0.1$

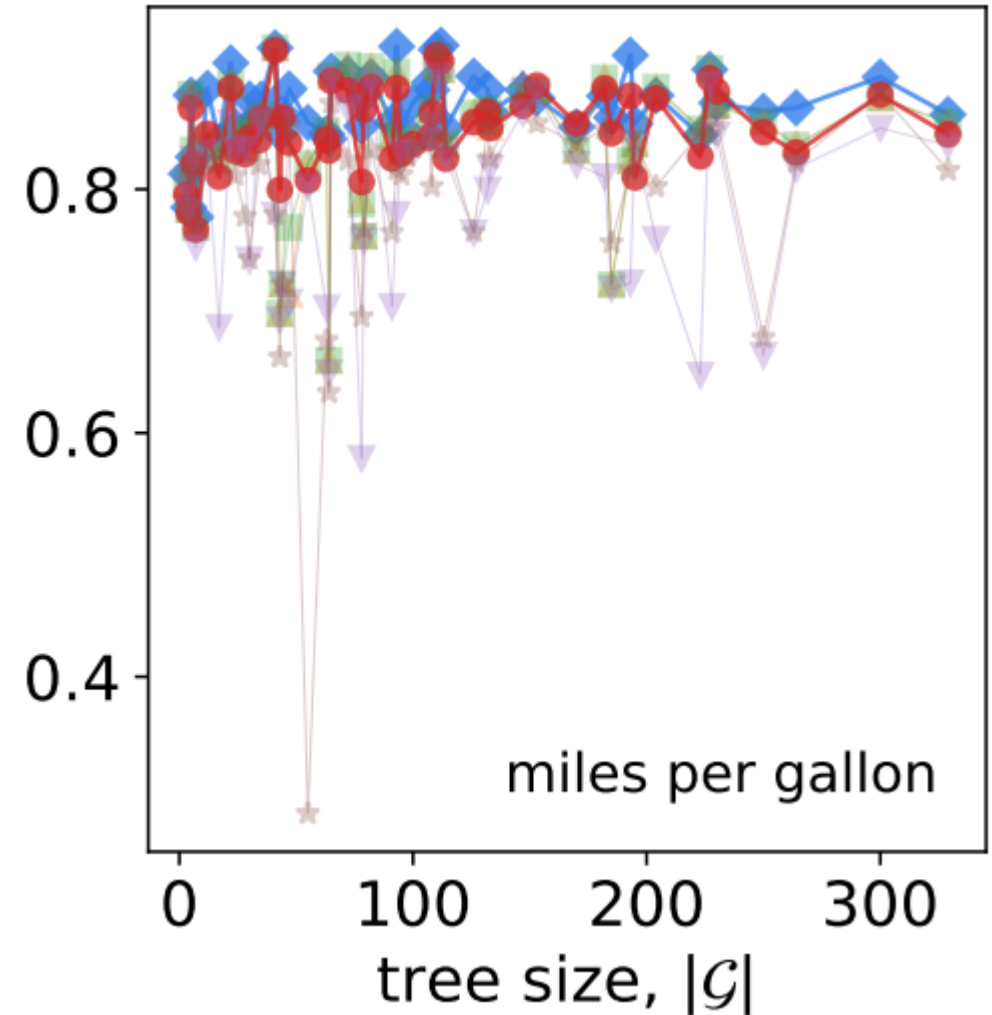
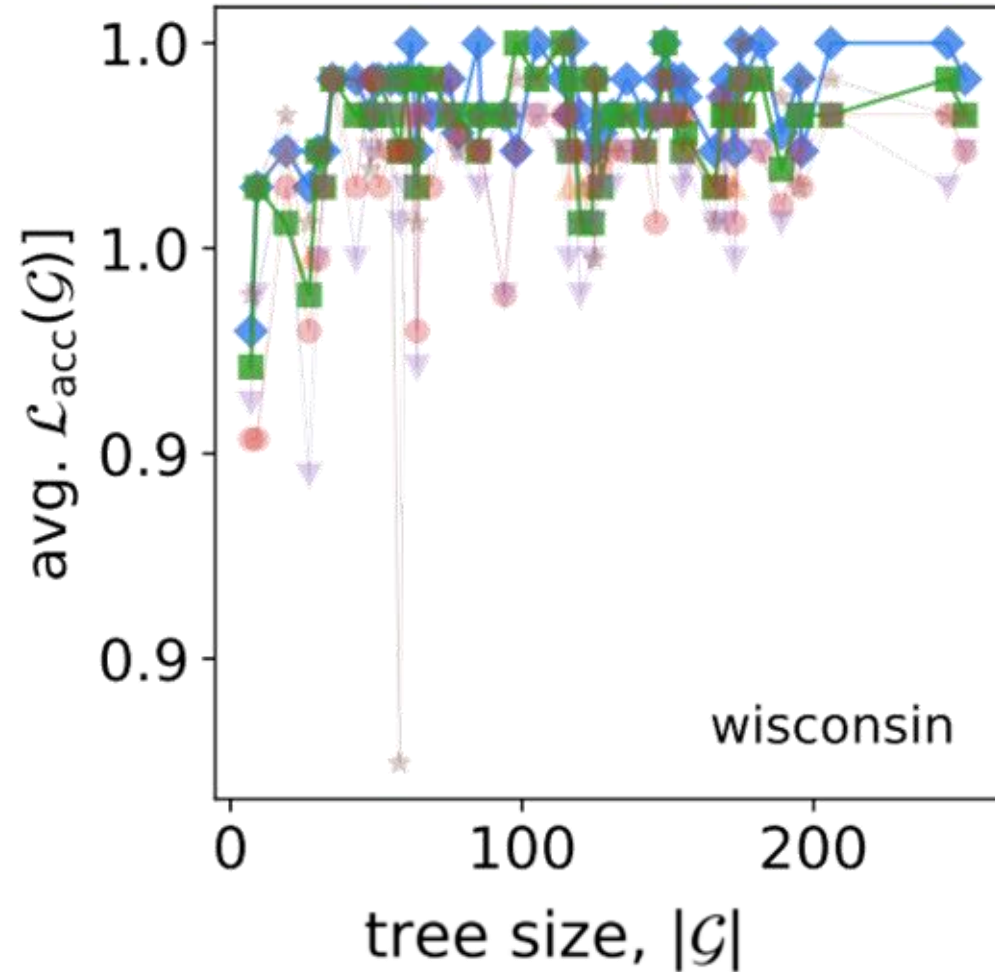


(e) MLP: Sigmoid, $\eta = \text{default}$

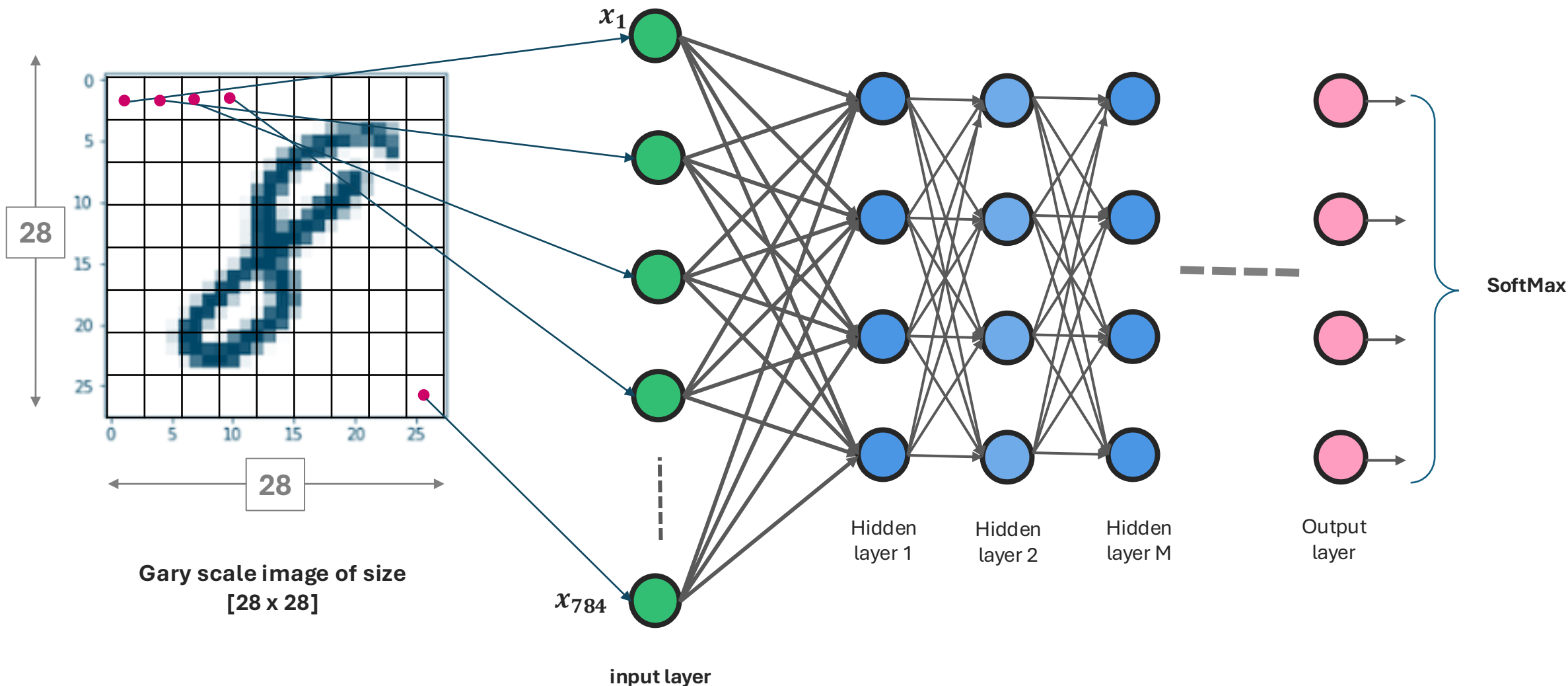


(f) MLP: ReLU, $\eta = 0.1$

Architectural Stochasticity



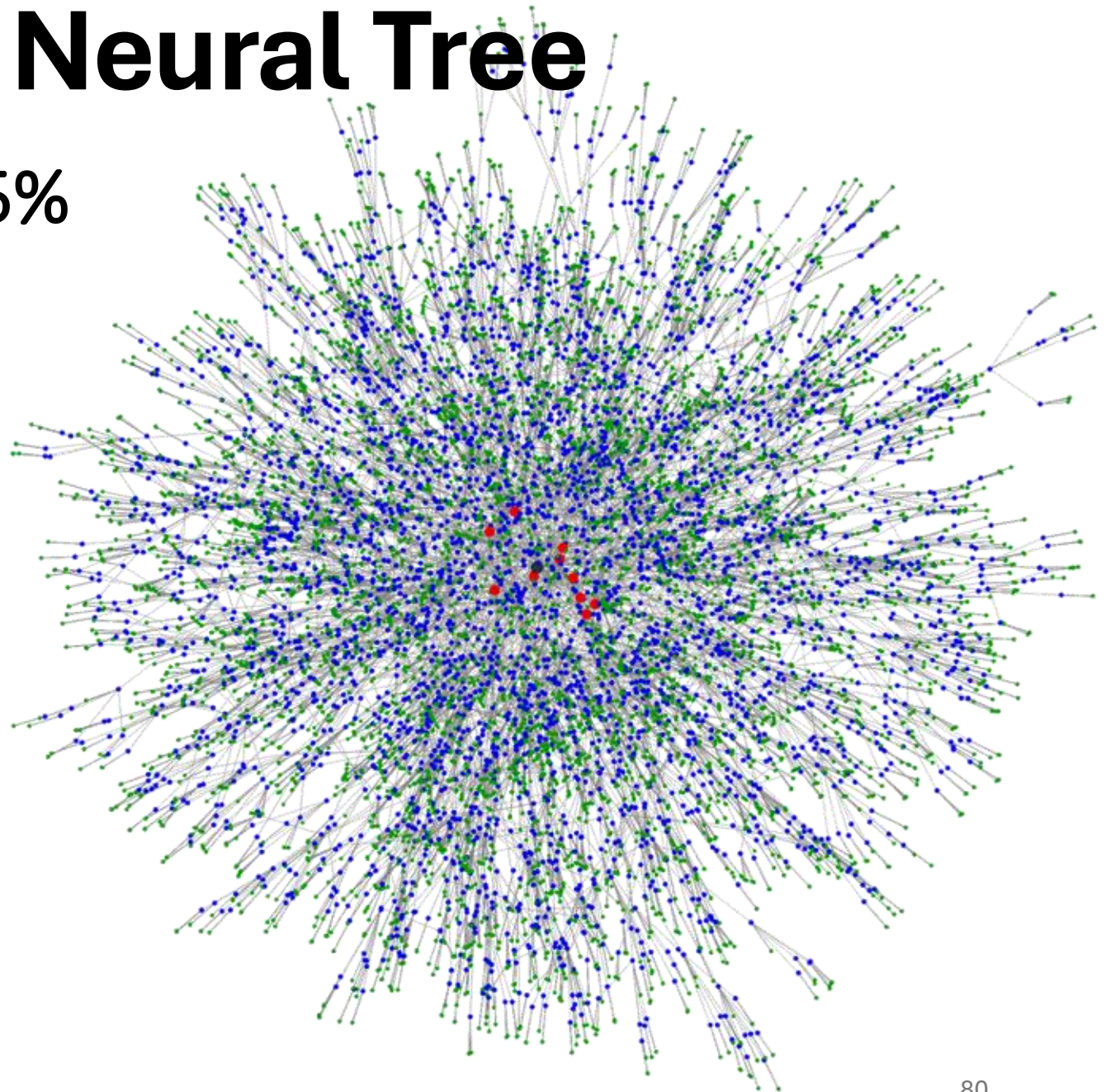
Deep Neural Networks



Backpropagation Neural Tree

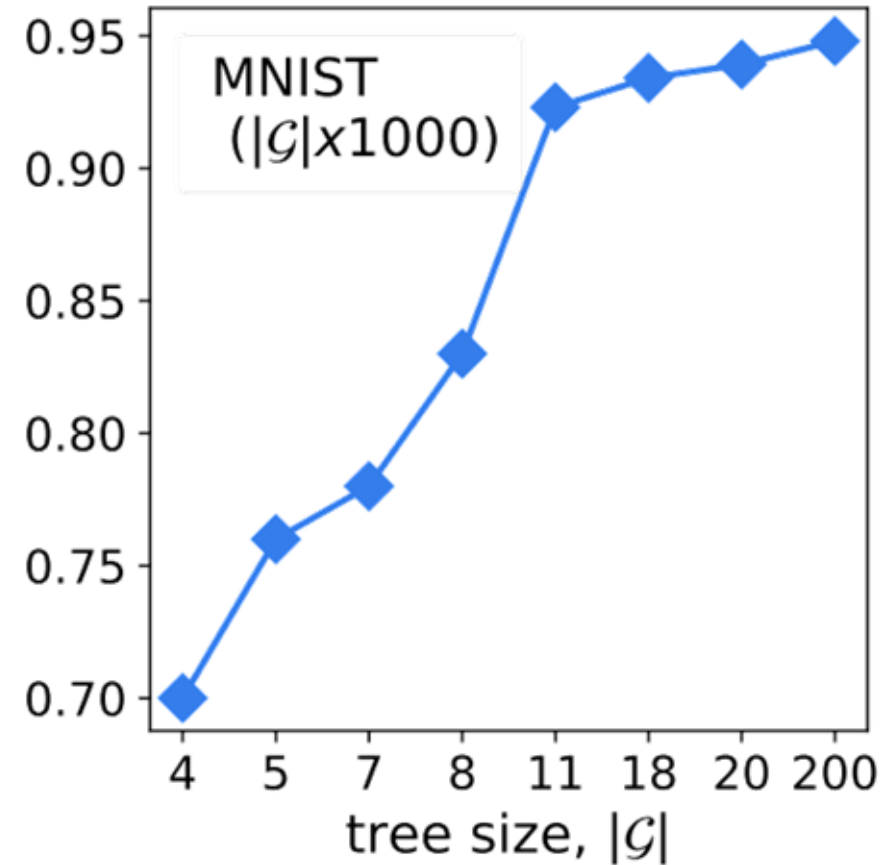
MNIST Model Accuracy ~95%

	Algorithms	Error(%)
BNeuralTs	BNeuralT-10K (pixels)	7.74
	BNeuralT-18K (pixels)	6.58
	BNeuralT-20K (pixels)	6.08
	BNeuralT-200K [†] (pixels)	5.19
Classification Trees	GUIDE (pixels, oblique split)	26.21
	OC1 (pixels, oblique split)	25.66
	GUIDE (pixels)	21.48
	CART-R (pixels)	11.97
	CART-P (pixels)	11.95
	C5.0 (pixels)	11.69
	TAO (pixels)	11.48
	TAO (pixels, oblique split)	5.26



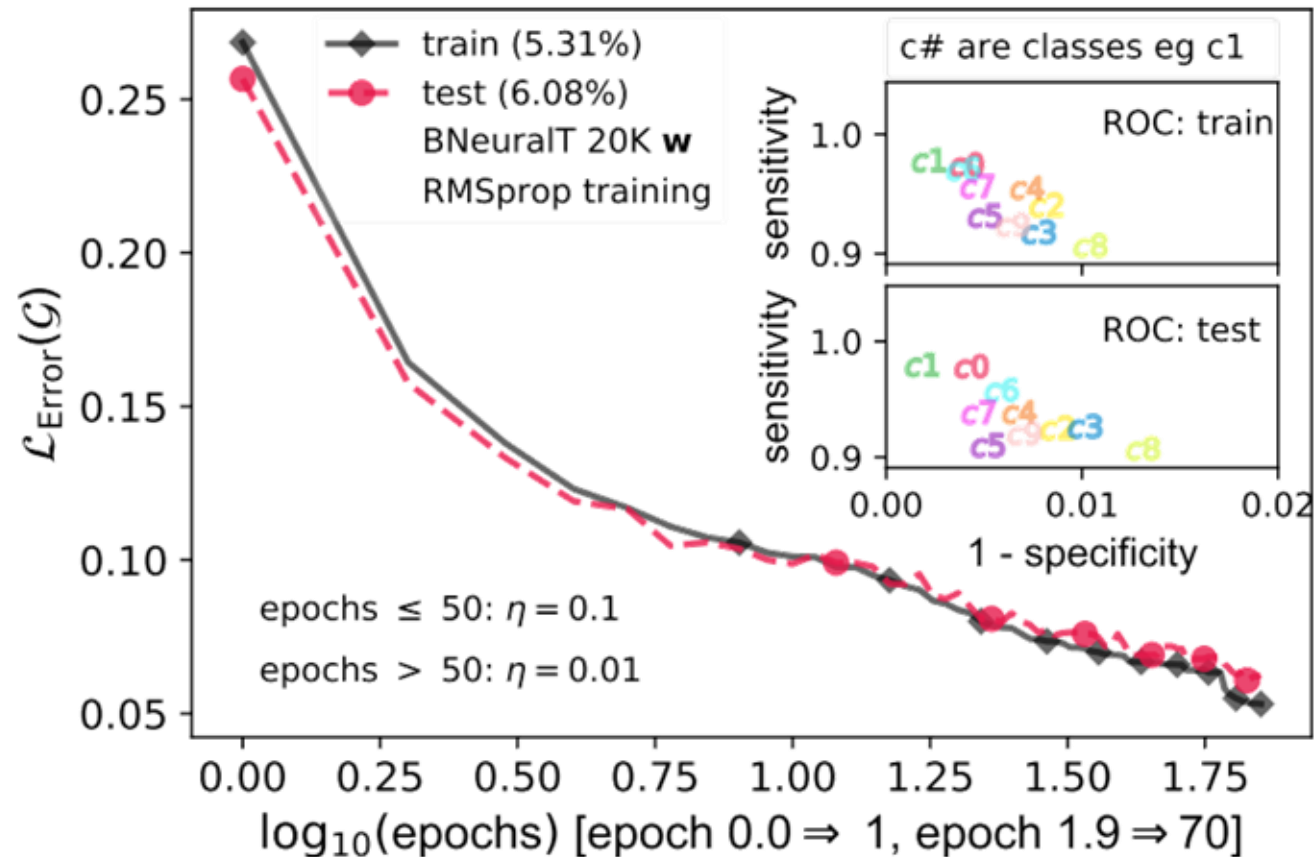
Model Size vs Accuracy

	Algorithms	Error(%)
BNeuralTs	BNeuralT-10K (pixels)	7.74
	BNeuralT-18K (pixels)	6.58
	BNeuralT-20K (pixels)	6.08
	BNeuralT-200K [†] (pixels)	5.19
Classification Trees	GUIDE (pixels, oblique split)	26.21
	OC1 (pixels, oblique split)	25.66
	GUIDE (pixels)	21.48
	CART-R (pixels)	11.97
	CART-P (pixels)	11.95
	C5.0 (pixels)	11.69
	TAO (pixels)	11.48
	TAO (pixels, oblique split)	5.26



Learnability of Different Classes

Competition between classes: TPR/Recall/Sensitivity vs FPR/(1 - Specificity)

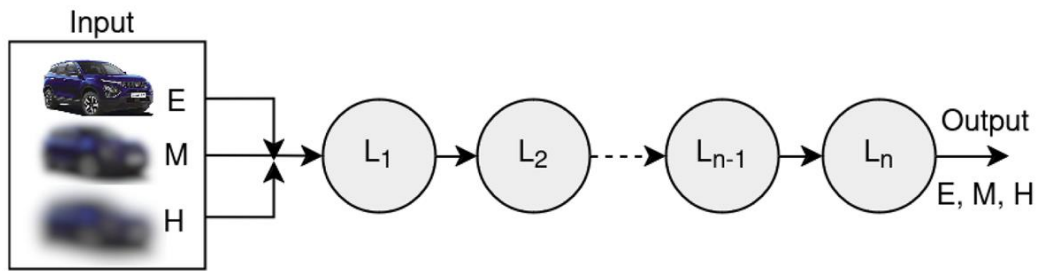


Part 5

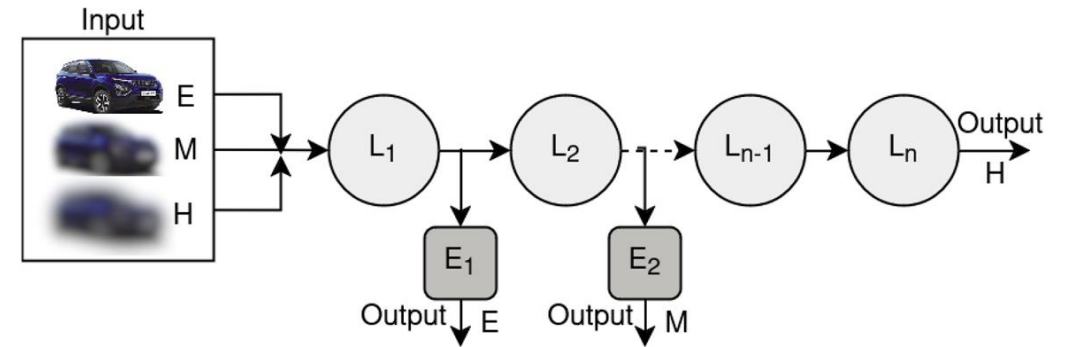
Clever tricks to Make AI Models Resource Efficient

Early Exit Neural Nets

(E, M, and H represent samples with easy, medium, and hard complexity data inputs)

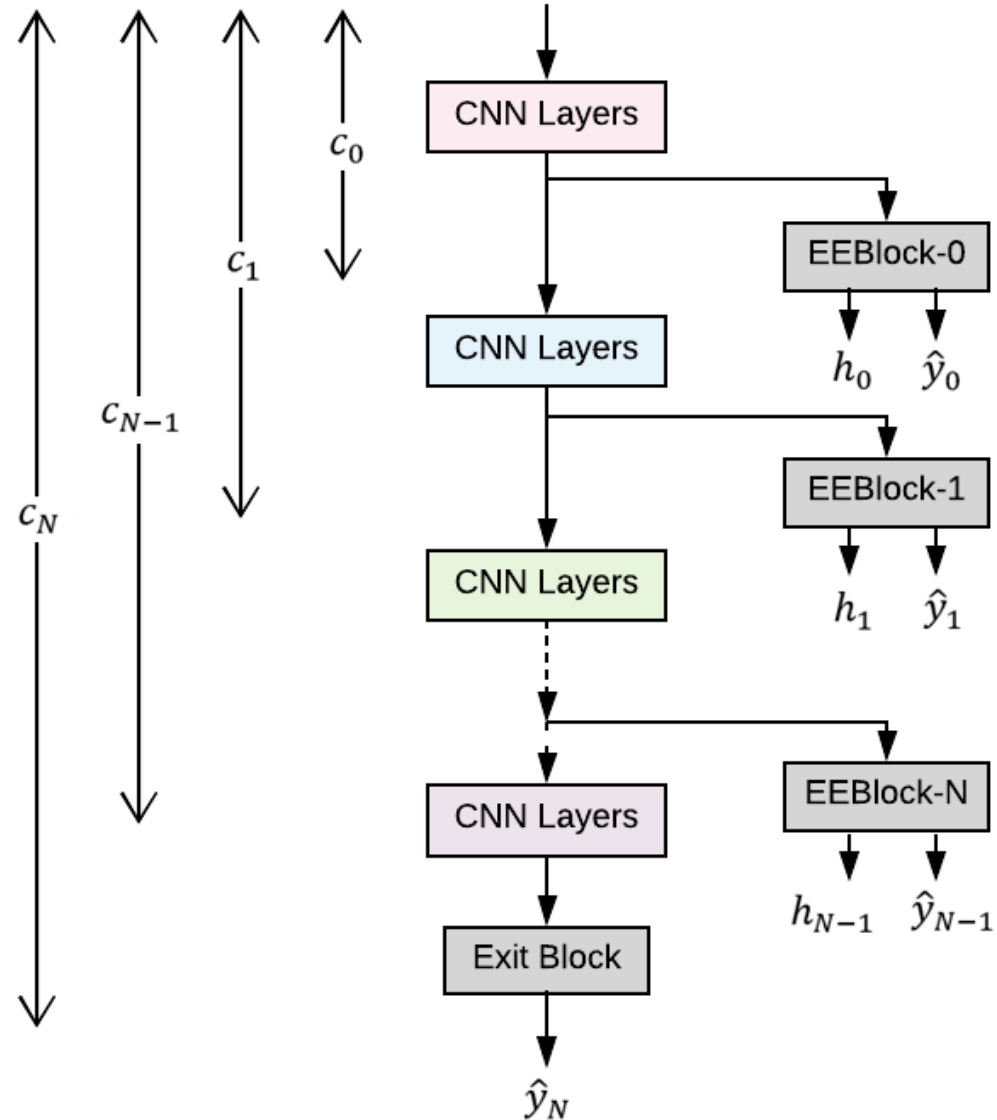
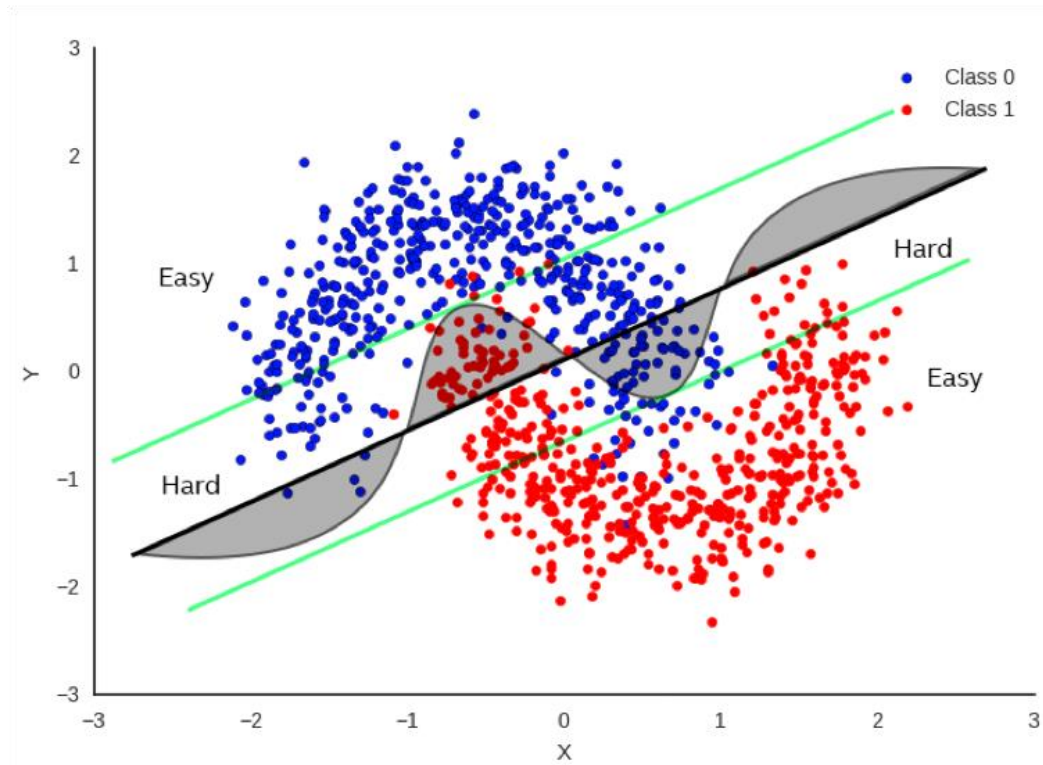


(a) Conventional DNN



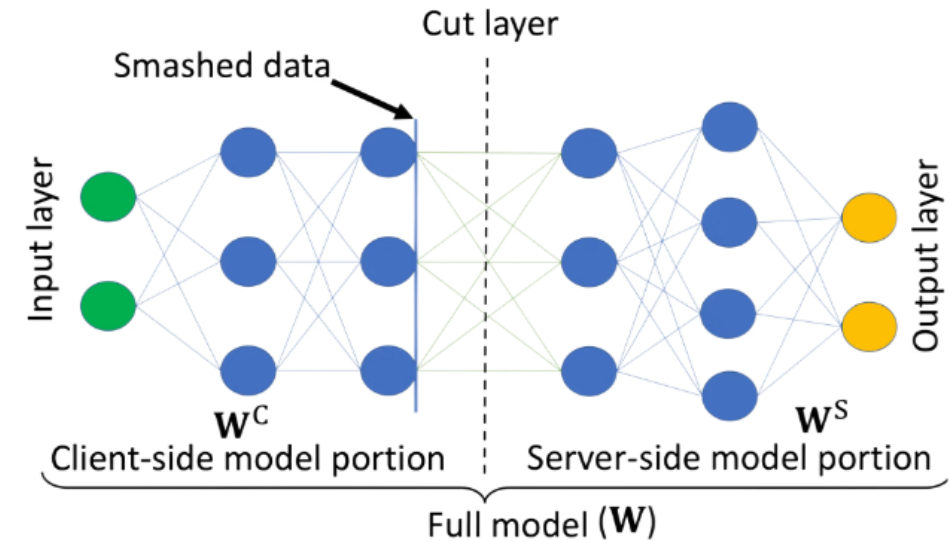
(b) Early-exit DNN

Early Exit Neural Nets



Model Splitting

- Model splitting is a deployment technique used to divide a neural network into two or more parts, which are then executed on different devices or platforms.
- This is particularly useful when a model is too large or computationally intensive to run on a single edge or IoT device.
- There are several strategies for model splitting:
 - Layer-wise splitting, where different layers of the network are executed on different devices
 - Feature-wise splitting, where the model is split based on feature maps
- The choice of splitting strategy depends on the specific architecture of the model and the capabilities of the devices involved.
- Challenges associated with model splitting include:
 - Minimizing communication overhead between devices
 - Ensuring that the split model maintains acceptable accuracy and latency
- Techniques such as knowledge distillation and quantization can be used to mitigate these challenges.



References

- [Backpropagation Neural Tree](#)
Neural Networks, Elsevier. (2022)
Ojha V, Nicosia G
- [Fragility, Robustness and Antifragility in Deep Learning](#)
Artificial Intelligence, Elsevier. (2024)
Pravin C, Martino I, Nicosia G, Ojha V
- [Dynamic Label Adversarial Training for Deep Learning Robustness Against Adversarial Attacks](#)
31st International Conference on Neural Information Processing (ICONIP). (2024)
Liu Z, Duan H, Liang H, Long Y, Snasel V, Nicosia G, Ranjan R and Ojha V
- [On Learnable Parameters of Optimal and Suboptimal Deep Learning Models](#)
31st International Conference on Neural Information Processing (ICONIP). (2024)
Zheng Z, Liang H, Snasel V, Latora V, Pardalos P, Nicosia G, and Ojha V
- [Sensitivity analysis for deep learning: Ranking hyper-parameter influence](#)
33rd IEEE Int. Conf. on Tools with Artificial Intelligence, ICTAI. IEEE (2021)
Taylor R, Martino I, Nicosia G, Ojha V
- [Multi-objective optimization of multi-output neural tree](#)
IEEE Congress on Evolutionary Computation (CEC), Glasgow, UK (pp 1–8). IEEE. (2020)
Ojha V, Nicosia G
- [Multiobjective programming for type-2 hierarchical fuzzy trees \(Code\)](#)
IEEE Transaction on Fuzzy Systems, 26(2), 915–936. (2018)
Ojha VK, Snášel V, Abraham A
- [Metaheuristic design of feedforward neural networks: a review of two decades of research](#)
Engineering Applications in Artificial Intelligence, 60, 97-116. Elsevier (2017)
Ojha VK, Snášel V, Abraham A
- [Ensemble of heterogeneous flexible neural trees using multiobjective genetic programming](#)
Applied Soft Computing, 52, 909-924. Elsevier (2017)
Ojha VK, Abraham A, Snášel V

Resource Efficient Artificial Intelligence

Varun Ojha

Senior Lecturer in Artificial Intelligence

AI Theme Lead National Edge AI Hub

School of Computing, Newcastle University

varun.ojha@newcastle.ac.uk

<https://ojhavk.github.io/>



Newcastle
University



**National
Edge AI
Hub**