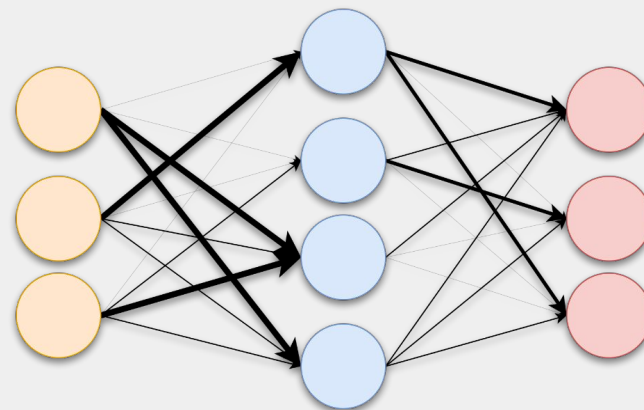
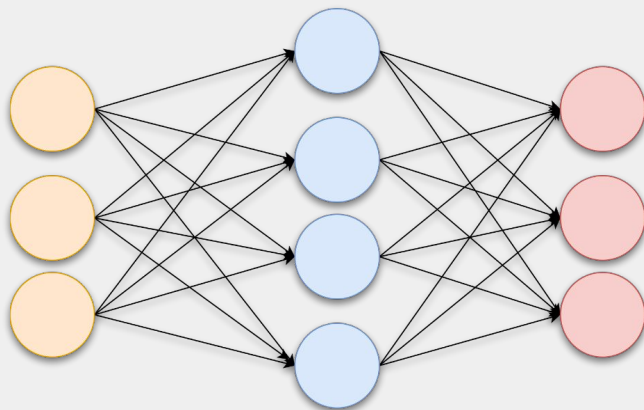


# Learning Dynamic Networks



 InstaDeep™

Kale-ab Tessera, Chiratidzo Matowe, Arnu Pretorius,  
Benjamin Rosman, Sara Hooker

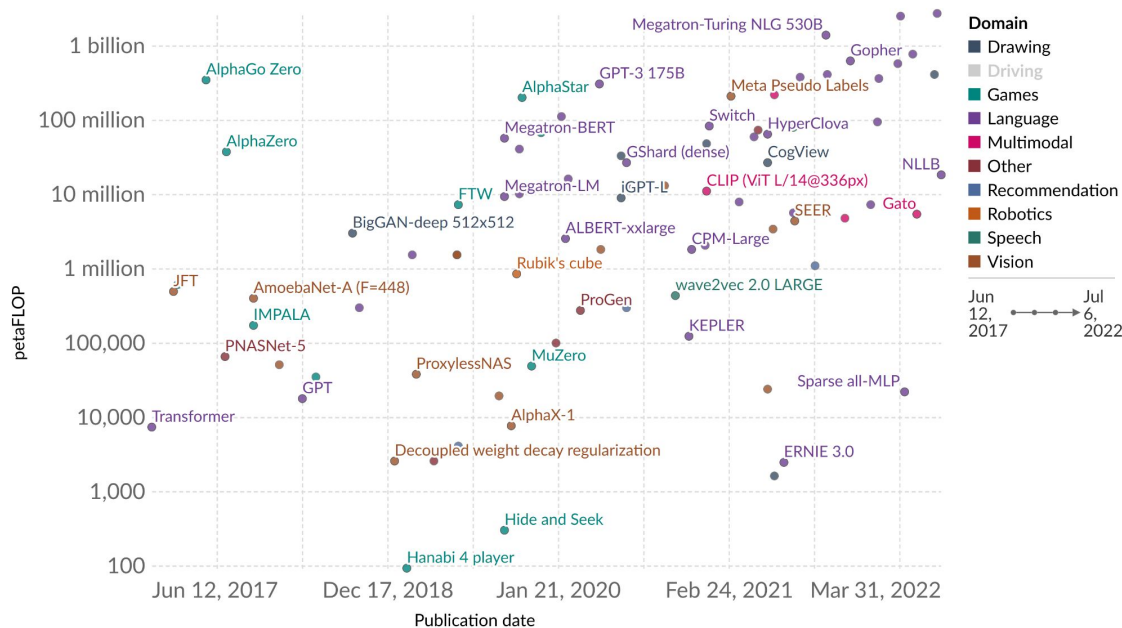


co:here

# Motivation - Era of Large Scale Models

## Estimated computation used to train notable AI systems

Computation is measured in petaFLOP, which is  $10^{15}$  floating-point operations.



**Overparameterized** models have led to many **breakthroughs** in machine learning (Chowdhery et al., 2022; Brown et al., 2020; Zhai et al., 2021; Reed et al., 2022).

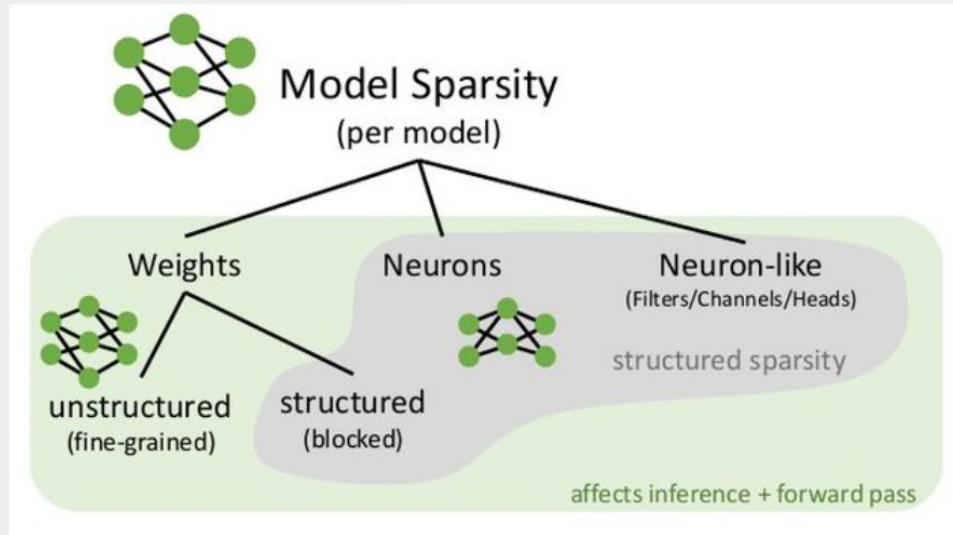
## Challenges:

- ❖ Larger Models
  - Efficient **storage** and inference.
  - Efficient **training** of large models.
  - **Overfitting**, regularization and generalization.
- ❖ Train Longer
  - Handle **temporal** dynamics of training.

# Motivation - Sparsity

Common method to handle challenges of overparameterization - **Sparsity/Pruning**.

Benefits - similar performance, with a **fraction of the weights** (Gale et al., 2019; Frankle & Carbin, 2019), **faster training** (Dettmers & Zettlemoyer, 2019; Luo et al., 2017) and more **robust to noise** (Ahmad & Scheinkman, 2019).



# Train Longer - Schedules

When we train longer -> more temporal decisions to make.

**Temporal Decisions** (examples include):

- **Learning Rate:**
  - Initial Learning Rate.
  - Function for the rate of change - [Need LR Schedule.](#)
- **Sparsity:**
  - Initial Sparsity (% active neurons, weights, filters, channels etc).
  - Function for the rate of change - [Need Sparsity Schedule.](#)

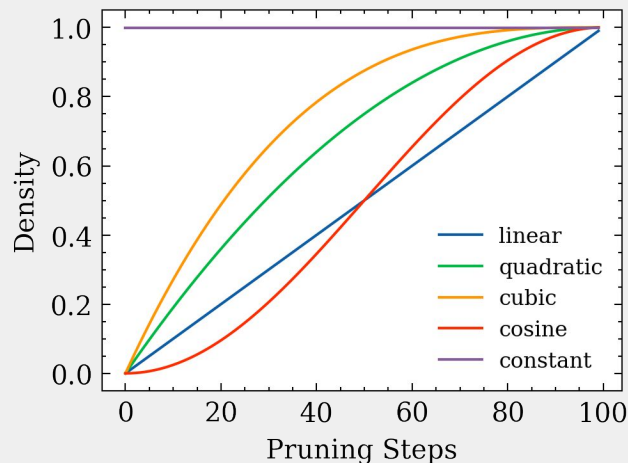
What about these choices per layer?! - **Layerwise Schedules.**

Standard approach - learn these schedules through **trial-and-error.**

## Related Work

**Handcrafted** schedules.

- **Constant** - SET [1], Deep Rewiring (DeepR) [2] and Neural Network Synthesis Tool (NEST) [3].
- **Cosine** - RigL [4] and Sparse Network From Scratch (SNFS) [5]
- **Cubic** - [6],[7].



# Our Approach - Can we learn these schedules using RL?

## Algorithm 1 Learning Sparsity Schedules using Reinforcement Learning

**Input:** train dataset  $X_{train\_set}$ , test dataset  $X_{test\_set}$ , train network  $f_{train}$ , eval network  $f_{eval}$ , agent  $a$ , number of episodes  $N$ , minimum density per layer  $min_d$  and maximum density per layer  $max_d$ .

$a \leftarrow \text{init}(min_d, max_d)$

▷ Initialize agent  $a$ .

$X_{train\_split}, X_{val\_split} \leftarrow \text{split}(X_{train\_set})$

▷ Split train dataset.

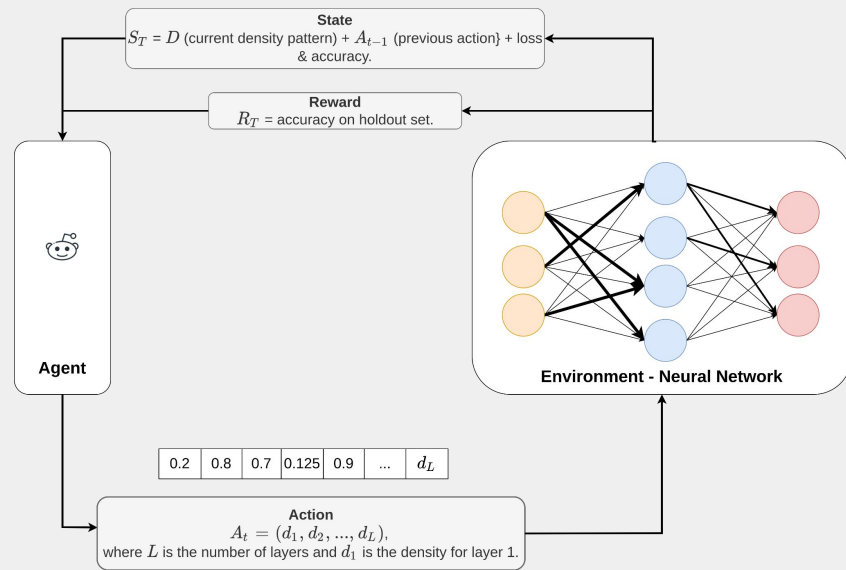
**for** episode=1, $N$  **do**

$a \leftarrow \text{train\_loop}(a, X_{train\_split}, X_{val\_split}, f_{train})$  ▷ Run train loop and retrieve trained agent  $a$ .

$\text{eval\_loop}(a, X_{train\_set}, X_{val\_set}, f_{eval})$  ▷ Run evaluation loop on unseen network  $f_{eval}$  using trained agent  $a$ .

**end for**

- Agent - PPO.
- Dataset - Cifar10.
- Sparsity:
  - Random Pruning, with Random Regrowth (**RP-RR**)
  - Magnitude Pruning, with Random Regrowth (**MP-RR**)



# Results - Simple CNN - 5 Layers

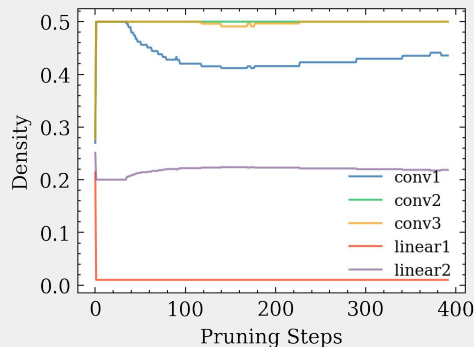
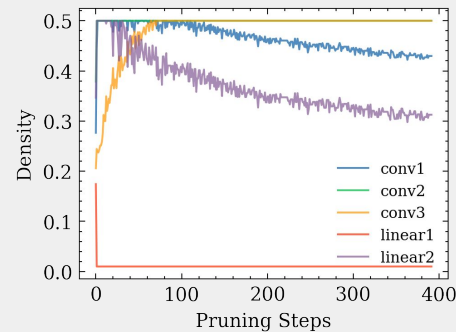
## 1. Learned Schedules are Competitive

Table 1: Test Accuracy (mean and standard deviation) of different schedules on CIFAR-10, using Simple-CNN.

Target Density (%)	Schedule	Random Pruning with Random Regrowth (RP-RR)	Magnitude Pruning with Random Regrowth (MP-RR)
10	Linear	20.365 $\pm$ 17.952	60.418 $\pm$ 1.362
	Quadratic	23.15 $\pm$ 20.043	61.259 $\pm$ 1.485
	Cubic	33.721 $\pm$ 21.693	60.396 $\pm$ 0.832
	Cosine	18.302 $\pm$ 14.379	59.807 $\pm$ 0.384
	Constant	<b>61.475 <math>\pm</math> 0.731</b>	62.536 $\pm$ 0.314
	Learned (Ours)	61.071 $\pm$ 1.574	<b>63.191 <math>\pm</math> 0.810</b>
50	Linear	64.54 $\pm$ 0.477	64.78 $\pm$ 0.464
	Quadratic	64.987 $\pm$ 0.86	63.933 $\pm$ 0.431
	Cubic	65.31 $\pm$ 0.49	64.315 $\pm$ 0.437
	Cosine	64.672 $\pm$ 0.771	64.737 $\pm$ 0.345
	Constant	65.1 $\pm$ 0.283	65.388 $\pm$ 0.375
	Learned (Ours)	<b>65.655 <math>\pm</math> 0.515</b>	<b>65.686 <math>\pm</math> 0.284</b>
100	Linear	66.228 $\pm$ 0.691	66.711 $\pm$ 0.423
	Quadratic	66.947 $\pm$ 0.749	67.25 $\pm$ 0.578
	Cubic	66.857 $\pm$ 0.627	67.395 $\pm$ 0.547
	Cosine	66.074 $\pm$ 0.282	66.18 $\pm$ 1.027
	Full Dense	<b>67.815 <math>\pm</math> 0.146</b>	67.878 $\pm$ 0.482
	Learned (Ours)	67.534 $\pm$ 0.174	<b>67.908 <math>\pm</math> 0.162</b>

## 3. Learned a Handcrafted Schedule! Piecewise Schedule for Random Pruning- [8].

## 2. Learned Schedules are Layerwise Diverse



# Results - ResNet18 & Conclusion

More **challenging** environment  
for our agent.

Schedule	Test Accuracy
Linear	93.019 +- 0.024
Quadratic	93.106 +- 0.107
<b>Cubic</b>	<b>93.148 +- 0.156</b>
Cosine	92.916 +- 0.105
Constant (Fully Dense)	92.481 +- 0.641
Learned (Ours)	92.818 +- 0.048

## Challenges:

1. **Non-stationarity** environment.
  - a. Our environment (the network we are learning a schedule for) is learning and adapting while our agent is learning to model the environment.
  - b. Worse for challenging networks - use techniques like **data augmentation** and **learning rate decay** (e.g. ResNet-18).
2. **High dimension** action and (possibly) state space.
3. **Slow convergence** - 25-50 episodes.

## Conclusion:

In this work, we demonstrate that it is **possible** to learn **well-performing dynamic sparsity schedules** using reinforcement learning. The schedules learned are not arbitrary and are distinct per layer and pruning method.

[ICML Workshop Paper](#)- Workshop on Dynamic  
Neural Networks.

