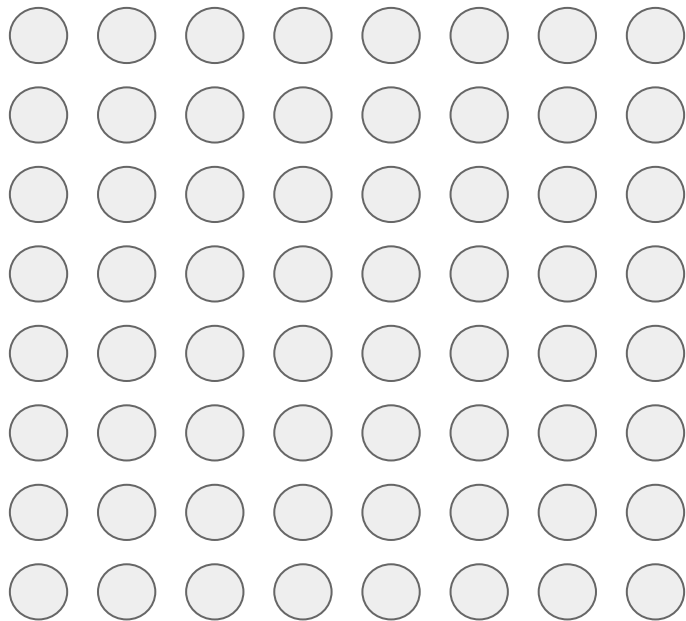


MIXTURE OF EXPERTS

Trelis Research

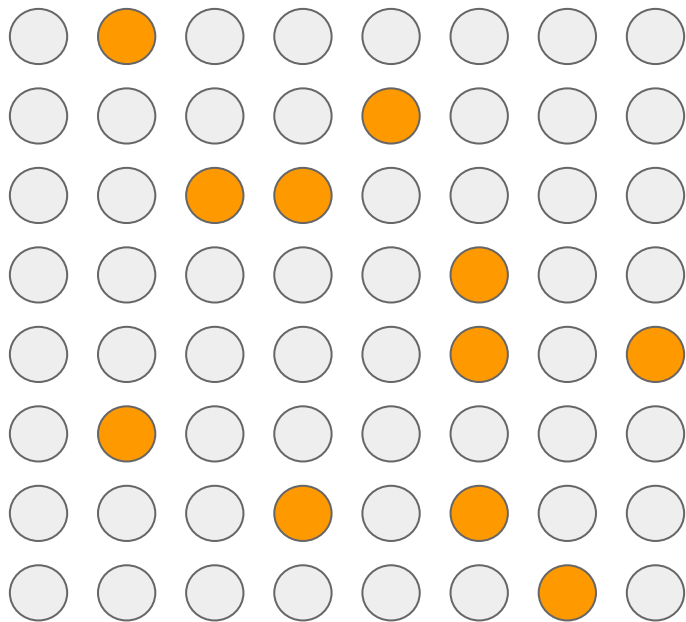
WHY MOE?



A TRADITIONAL GPT
USES ALL NEURONS IN
ALL MATRICES FOR
FORWARD PASSES



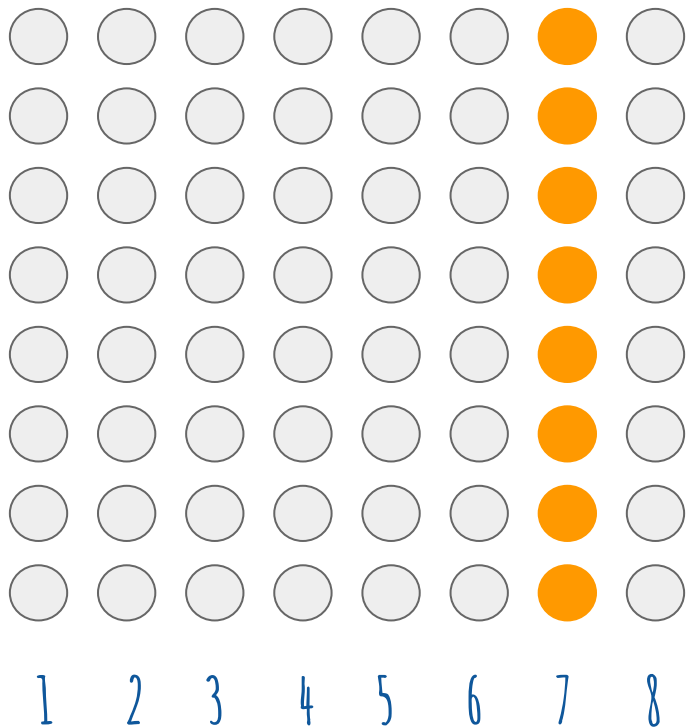
WHY MOE?



A TRADITIONAL GPT
USES ALL NEURONS IN
ALL MATRICES FOR
FORWARD PASSES



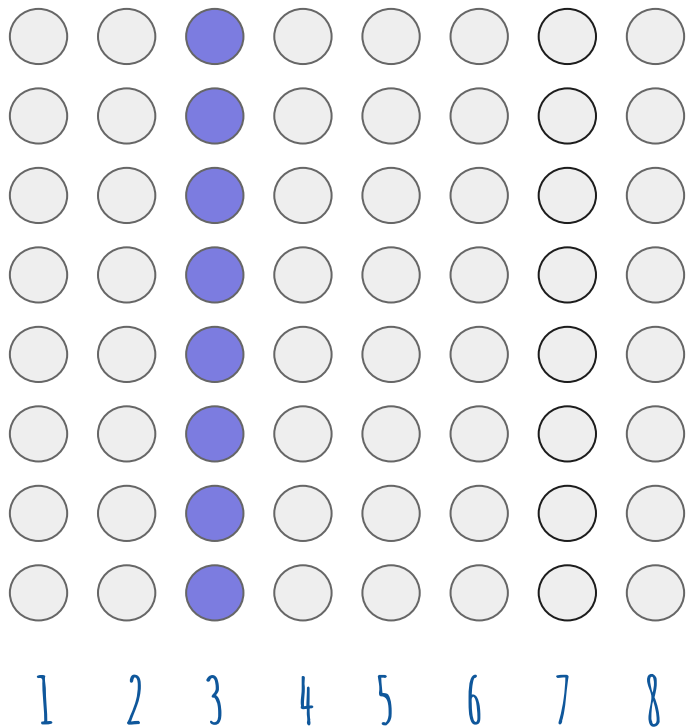
WHY MOE?



WHAT IF WE COULD SPLIT
THE NETWORK + CHOOSE
THE BEST COLUMN OF
WEIGHTS TO USE?



WHY MOE?



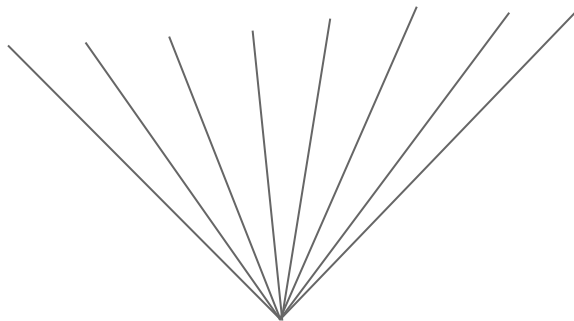
WHAT IF WE COULD SPLIT
THE NETWORK + CHOOSE
THE BEST COLUMN OF
WEIGHTS TO USE?



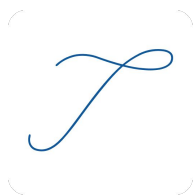
HOW IS MOE TRAINED? ROUTERS



1 2 3 4 5 6 7 8



ROUTER



$[0, 0, 0, 0, 0, 0, 1, 0]$

EXPERT CHOICE



ROUTER MATRIX WEIGHTS



$[0.34, 0.35, 0.73, 0.94]$

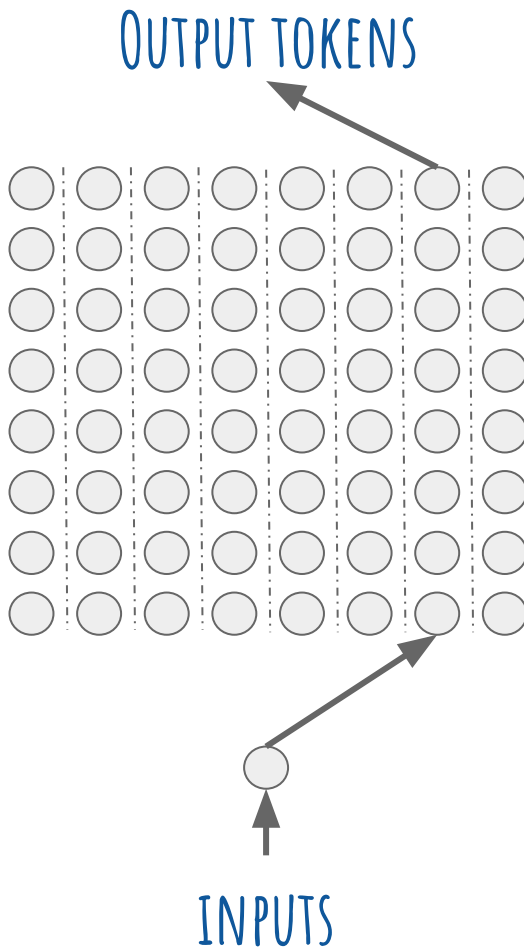
INPUT VECTOR

T

TRAINING

TRANSFORMER WEIGHTS

ROUTER WEIGHTS

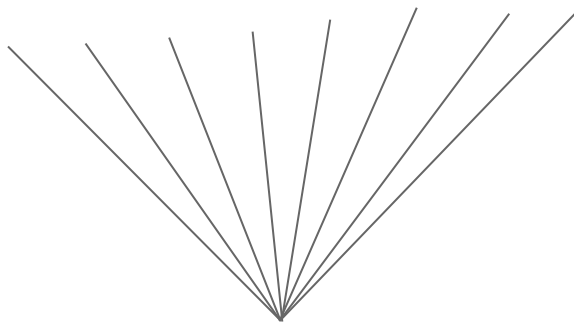
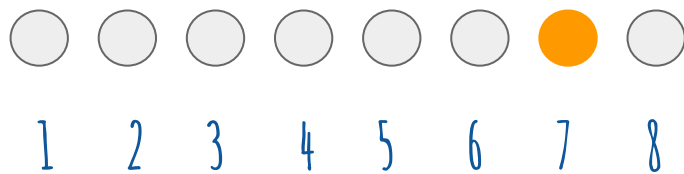


$$\text{LOSS} = \text{ACTUAL} - \text{PREDICTED}$$

BACK PROPAGATION
THROUGH GPT +
ROUTER

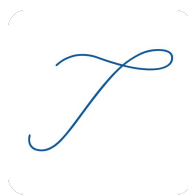
T

A PROBLEM TRAINING MIXTURE OF EXPERTS



ONE EXPERT CAN
DOMINATE!!!

ROUTER



TRICK #1

$[0, 0, 0, 0, 0, 0, 1, 0]$

EXPERT CHOICE

↑ +RANDOMNESS

ROUTER MATRIX WEIGHTS

↑

$[0.34, 0.35, 0.73, 0.94]$

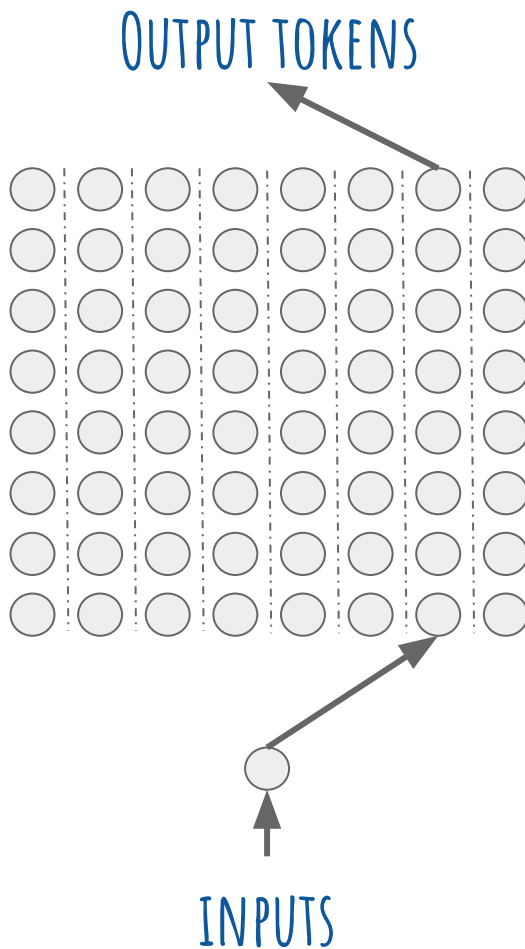
INPUT VECTOR

T

TRICK #2

TRANSFORMER WEIGHTS

ROUTER WEIGHTS



LOSS = ACTUAL - PREDICTED
TOKENS

+ PENALTY FOR UNEVEN
ROUTER CHOICE

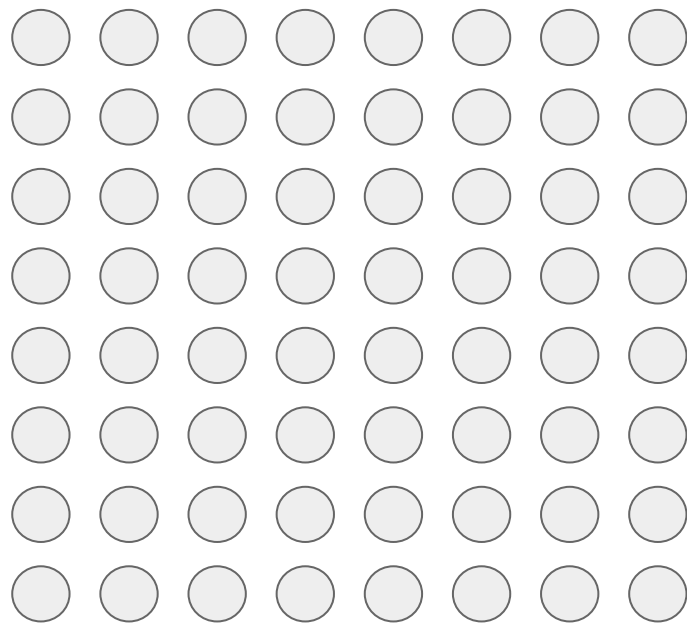
BACK PROPAGATION
THROUGH GPT +
ROUTER

T

IS MOE USEFUL FOR LAPTOP INFERENCE?



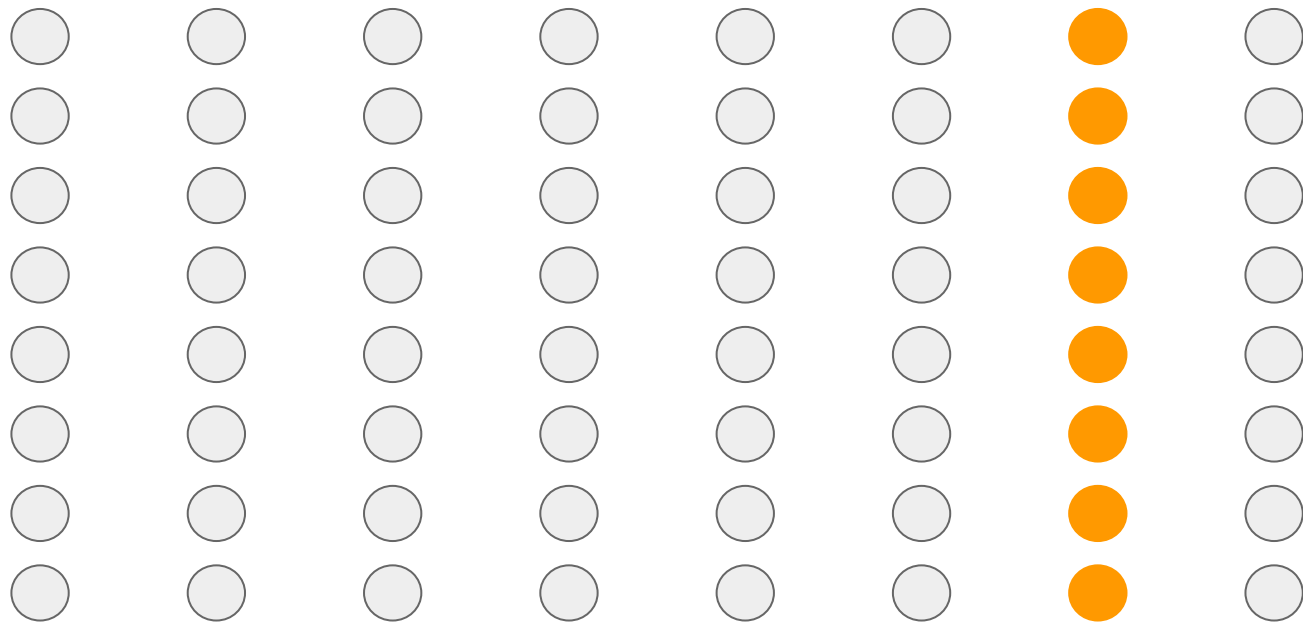
MOE



STANDARD GPT



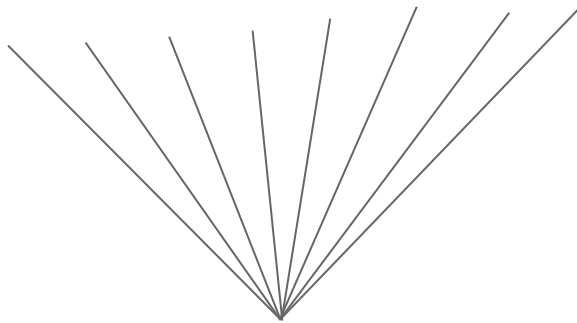
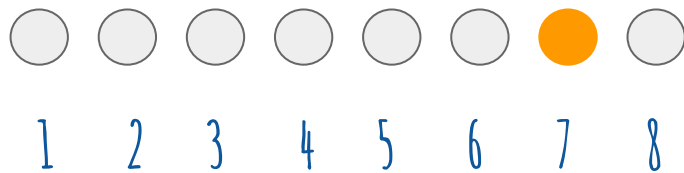
DOES MOE REDUCE COSTS AT SCALE?



EACH EXPERT GETS ITS OWN GPU AND QUERIES ARE ROUTED AND BATCHED

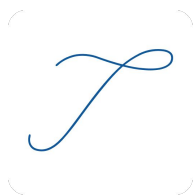


WHAT ARE THE PROBLEMS WITH MOE

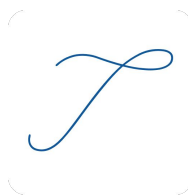
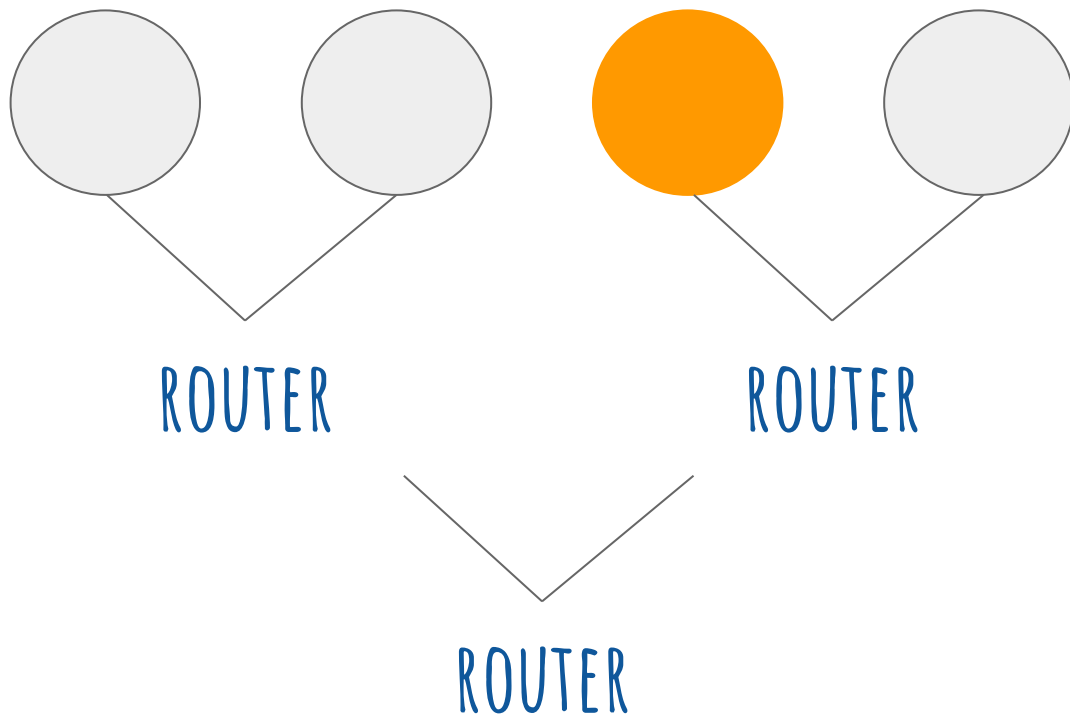


SLOW TO TRAIN EXPERTS
EVENLY DUE TO NOISE

ROUTER



FAST FEED FORWARD - BETTER THAN MOE?



[1] OR [0] WITH PROBABILITY P

EXPERT CHOICE



ROUTER MATRIX WEIGHTS



[0.34, 0.35, 0.73, 0.94]

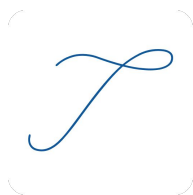
INPUT VECTOR

T

FAST FEED FORWARD - BETTER THAN MOE?

BINARY TREE NETWORKS ALLOW:

- SAME TRAINING TIME AS GPT
- FASTER INFERENCE



FAST FEED FORWARD - BETTER THAN MOE?

Width	Model											
	feedforward				mixture-of-experts ($e=16, k=2$)				fast feedforward ($\ell=32$)			
	M_A	ETT	G_A	ETT	M_A	ETT	G_A	ETT	M_A	ETT	G_A	ETT
$w = 64$	87.2	307	49.3	55	57.8	5354	29.4	4880	85.8	302	45.9	22
$w = 128$	95.5	200	51.5	46	62.0	6074	33.6	938	90.1	305	45.5	22
$w = 256$	99.9	105	52.0	48	62.4	2001	33.9	372	91.2	244	44.4	17
$w = 512$	99.9	85	52.4	31	65.4	3834	34.5	315	96.2	175	43.7	10
$w = 1024$	99.9	82	53.0	21	65.3	1575	35.2	327	96.0	180	41.3	9

Table 2: The results of the comparison of feedforward, mixture-of-experts, and fast feedforward networks, for various training widths. The inference width is fixed to 32 for mixture-of-experts and fast feedforward networks. The ETT columns to the right of metric columns list the “epochs to train”, i.e. the number of training epochs that have elapsed until the score to the left was observed.

ETT = EPOCHES TO TRAIN. M_A =MEMORISATION, G_A =GENERALISATION



FAST FEED FORWARD - INFERENCE

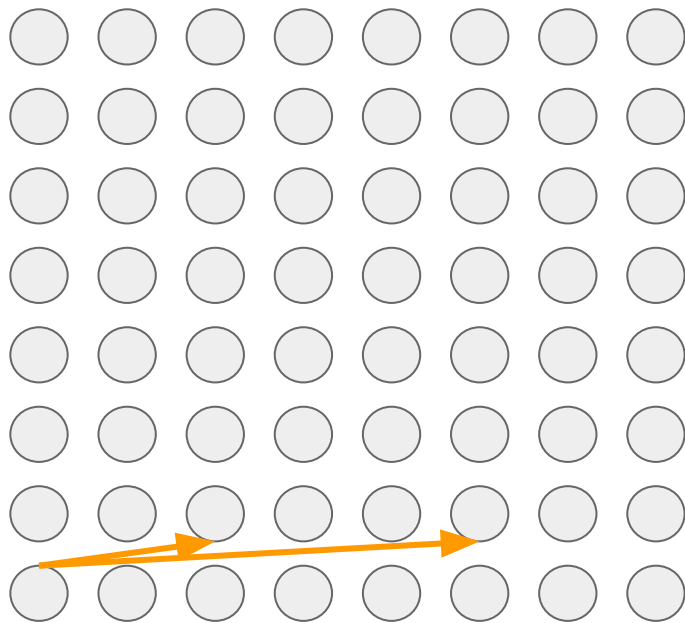
Model		Property						
		depth	training width	training size	inference width	inference size	speedup	G_A
FF	$w = 128$	–	128	128 (100%)	128 (100%)	128 (100%)	1.00x	84.7
fast FF	$\ell = 32$	2	128	131 (102%)	32 (25%)	34 (27%)	2.44x	83.6
	$\ell = 16$	3	128	135 (105%)	16 (13%)	19 (15%)	2.80x	83.2
	$\ell = 8$	4	128	143 (112%)	8 (6%)	12 (9%)	3.29x	82.8
	$\ell = 4$	5	128	159 (124%)	4 (3%)	9 (7%)	3.39x	81.6
	$\ell = 2$	6	128	191 (149%)	2 (1%)	8 (6%)	3.47x	80.1
	$\ell = 1$	7	128	255 (199%)	1 (1%)	8 (6%)	3.93x	79.8

Table 3: The results of the testing of vision transformers leveraging feedforward and fast feedforward layers. All sizes are given in neurons. Bracketed percentages describe quantities relative to their counterparts in the vanilla feedforward layers. G_A is the generalization accuracy of the fully trained vision transformer and “speedup” gives the performance improvement over vanilla feedforward layers in our testing setup.

CAVEAT - THE ABOVE IS A VERY SMALL NEURAL NET FOR VISION.



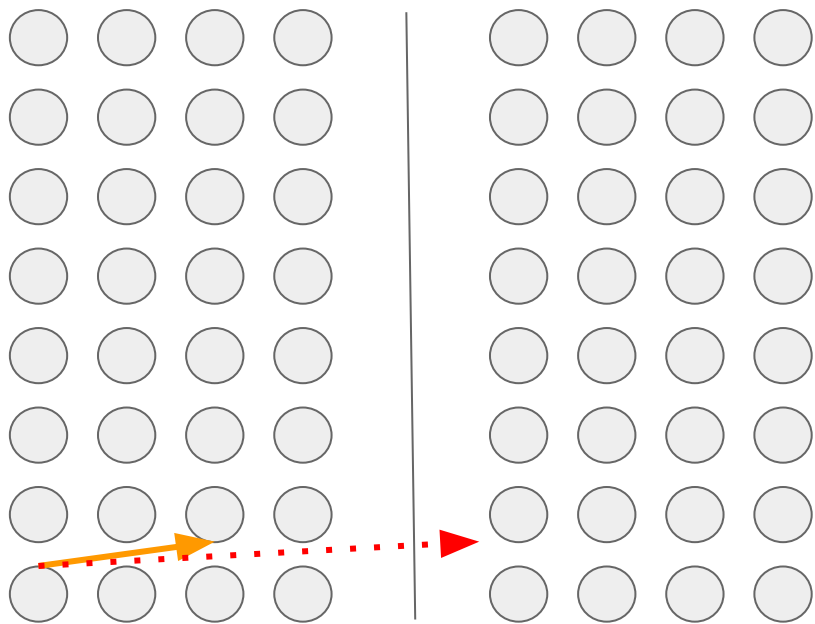
WHY DOES MOE (OR FFF) WORK?



STANDARD GPT



WHY DOES MOE (OR FFF) MAKE SENSE?



MOE/FFF

SAME MODEL SIZE
FEWER CONNECTIONS
SMALL ACCURACY LOSS
 $\sim 1/N$ SPEED GAIN



PAPERS AND LINKS

BINARY-TREE/FFF PAPER: [HTTPS://ARXIV.ORG/PDF/2308.14711.PDF](https://arxiv.org/pdf/2308.14711.pdf)

MOE PAPERS: [HTTPS://ARXIV.ORG/PDF/2208.02813.PDF](https://arxiv.org/pdf/2208.02813.pdf);

[HTTPS://ARXIV.ORG/PDF/1701.06538.PDF](https://arxiv.org/pdf/1701.06538.pdf)

REDDIT THREAD: [HTTPS://TINYURL.COM/YTHSU2ND](https://tinyurl.com/ythsu2nd)

YOUTUBE VIDEO: [HTTPS://YOUTU.BE/OU_65FLOTQ0](https://youtu.be/OU_65fLOTQ0)

TRELIS.COM

