



Sparse MoEs meet efficient ensembles

Rodolphe Jenatton (rjenatton@google.com)

Credits

James Urquhart Allingham^{1,*} *jua23@cam.ac.uk*

Florian Wenzel^{2,†} *fln.wenzel@gmail.com*

Zelda E Mariet, Basil Mustafa *{zmariet,basilm}@google.com*

Joan Puigcerver, Neil Houlsby *{jpuigcerver,neilhoulby}@google.com*

Ghassen Jerfel^{3,†} *ghassen@google.com*

Vincent Fortuin^{1,4,*} *vbf21@cam.ac.uk*

Balaji Lakshminarayanan, Jasper Snoek *{balajiln,jsnoek}@google.com*

Dustin Tran, Carlos Riquelme, Rodolphe Jenatton *{trandustin,rikel,rjenatton}@google.com*

Google Research, Brain Team; ¹University of Cambridge; ²no affiliation; ³Waymo; ⁴ETH Zürich

Credits

James Urquhart Allingham ^{1,*}	<i>jua23@cam.ac.uk</i>
Florian Wenzel ^{2,†}	<i>fln.wenzel@gmail.com</i>
Zelda E Mariet, Basil Mustafa	<i>{zmariet,basilm}@google.com</i>
Joan Puigcerver, Neil Houlsby	<i>{jpuigcerver,neilhoulby}@google.com</i>
Ghassen Jerfel ^{3,†}	<i>ghassen@google.com</i>
Vincent Fortuin ^{1,4,*}	<i>vbf21@cam.ac.uk</i>
Balaji Lakshminarayanan, Jasper Snoek	<i>{balajiln,jsnoek}@google.com</i>
Dustin Tran, Carlos Riquelme, Rodolphe Jenatton	<i>{trandustin,rikel,rjenatton}@google.com</i>
<i>Google Research, Brain Team; ¹University of Cambridge; ²no affiliation; ³Waymo; ⁴ETH Zürich</i>	

Some parts of the talk reuse slides from:

- {balajiln, trandustin, jsnoek}@: NeurIPS tutorial 2020, [Uncertainty and Out-of-Distribution Robustness in Deep Learning](#)

Context and motivations

Project at the intersection of two topics:

- **Sparse mixture of experts (sparse MoEs)**
- **Reliability in deep learning**
 - **Why and how to measure it?**
 - **Ensembles**

Context and motivations

Project at the intersection of two topics:

- **Sparse mixture of experts (sparse MoEs)**
- **Reliability in deep learning**
 - **Why and how to measure it?**
 - **Ensembles**

Sparse mixture of experts (sparse MoEs)

In a nutshell:

- **Conditional computation** [Bengio et al., 2013]:
 - Only subpart of the model activated in an input-dependent fashion
 - (≠ standard “dense” models: all parameters used to process an input)

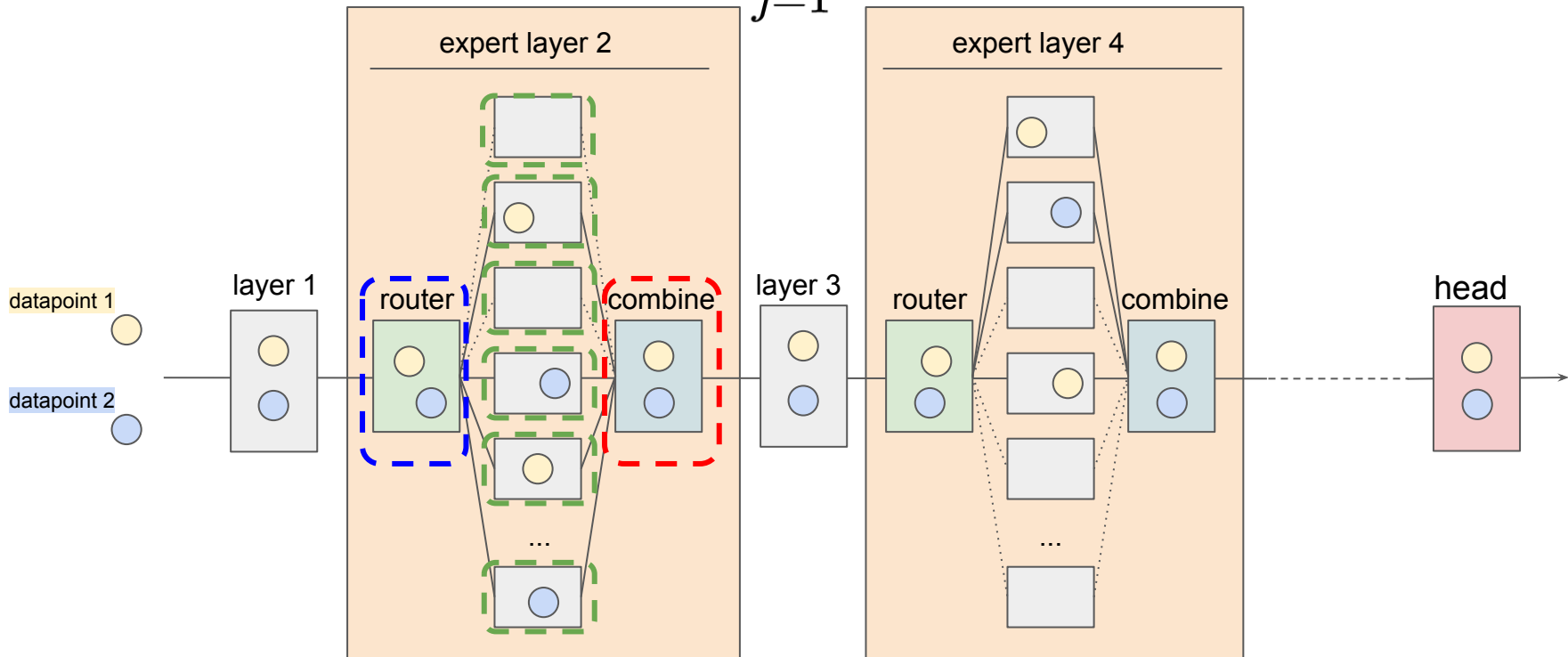
Sparse mixture of experts (sparse MoEs)

In a nutshell:

- **Conditional computation** [Bengio et al., 2013]:
 - Only subpart of the model activated in an input-dependent fashion
 - (\neq standard “dense” models: all parameters used to process an input)
- **Goal: Grow model size while keeping compute \sim constant**

Pictorial view of sparse MoEs

$$f(\mathbf{x}) = \sum_{j=1}^E g_j(\mathbf{x}) \cdot \text{expert}_j(\mathbf{x})$$



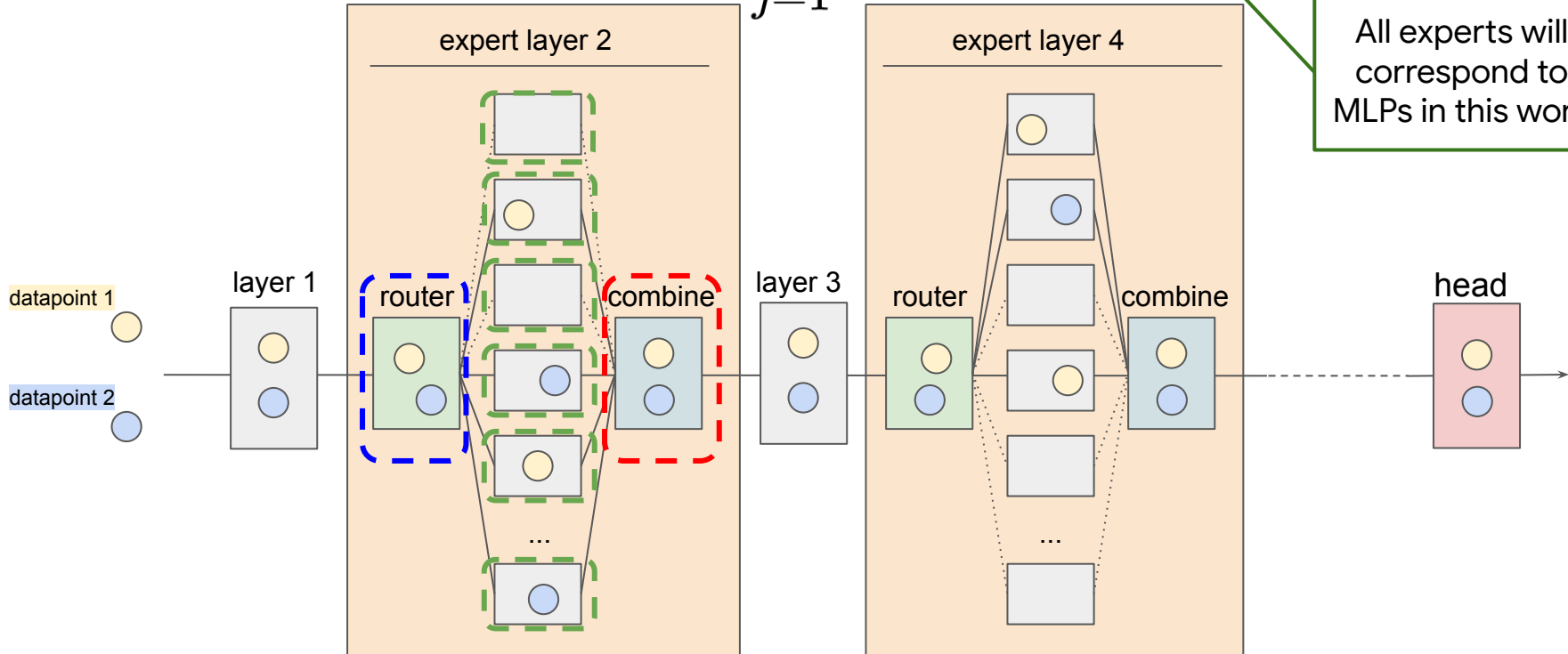
Pictorial view of sparse MoEs

$$f(\mathbf{x}) = \sum_{j=1}^E g_j(\mathbf{x}) \cdot \text{expert}_j(\mathbf{x})$$

#experts

Sparse gating weights
(only K non-zero)

All experts will correspond to MLPs in this work



Sparse MoEs successfully applied in NLP

OUTRAGEOUSLY LARGE NEURAL NETWORKS: THE SPARSELY-GATED MIXTURE-OF-EXPERTS LAYER

Noam Shazeer¹, Azalia Mirhoseini^{*†1}, Krzysztof Maziarz^{*2}, Andy Davis¹, Quoc Le¹, Geoffrey Hinton¹ and Jeff Dean¹

¹Google Brain, {noam,azalia,andydavis,qvl,geoffhinton,jeff}@google.com

²Jagiellonian University, Cracow, krzysztof.maziarz@student.uj.edu.pl

SWITCH TRANSFORMERS: SCALING TO TRILLION PARAMETER MODELS WITH SIMPLE AND EFFICIENT SPARSITY

William Fedus*
Google Brain

liamfedus@google.com

Barret Zoph*
Google Brain

barretzoph@google.com

Noam Shazeer
Google Brain

noam@google.com

PATHWAYS: ASYNCHRONOUS DISTRIBUTED DATAFLOW FOR ML

Paul Barham¹ Aakanksha Chowdhery¹ Jeff Dean¹ Sanjay Ghemawat¹ Steven Hand¹ Dan Hurt¹
Michael Isard¹ Hyeontaek Lim¹ Ruoming Pang¹ Sudip Roy¹ Brennan Saeta¹ Parker Schuh¹
Ryan Sepassi¹ Laurent El Shafey¹ Chandramohan A. Thekkath¹ Yonghui Wu¹

ABSTRACT

We present the design of a new large scale orchestration layer for accelerators. Our system, PATHWAYS, is explicitly designed to enable exploration of new systems and ML research ideas, while retaining state of the art performance for current models. PATHWAYS uses a *sharded* dataflow graph of *asynchronous* operators that consume and produce futures, and efficiently gang-schedules *heterogeneous* parallel computations on thousands of accelerators while coordinating data transfers over their dedicated interconnects. PATHWAYS makes use of a novel

GShard: Scaling Giant Models with Conditional Computation and Automatic Sharding

Dmitry Lepikhin
lepikhin@google.com

HyoukJoong Lee
hyouklee@google.com

Yuanzhong Xu
yuanzx@google.com

Dehao Chen
dehao@google.com

Orhan Firat
orhanf@google.com

Yanping Huang
huangyp@google.com

Maxim Krikun
krikun@google.com

Noam Shazeer
noam@google.com

Zhifeng Chen
zhifengc@google.com

Carbon Emissions and Large Neural Network Training

David Patterson^{1,2}, Joseph Gonzalez², Quoc Le¹, Chen Liang¹, Lluis-Miquel Munguia¹,
Daniel Rothchild², David So¹, Maud Texier¹, and Jeff Dean¹
{davidpatterson, qvl, crazydonkey, lmunguia, davidso, maudt, jeff}@google.com,
{pattsn, jegonzal, drothchild}@berkeley.edu

Abstract: The computation demand for machine learning (ML) [has grown rapidly](#) recently, which comes with a number of costs. Estimating the energy cost helps measure its environmental impact and finding greener strategies, yet it is [challenging without detailed information](#).

We calculate the energy use and carbon footprint of several recent large models—[T5](#), [Meena](#), [GShard](#), [Switch Transformer](#), and [GPT-3](#)—and refine earlier estimates for the neural architecture search that found [Evolved Transformer](#).

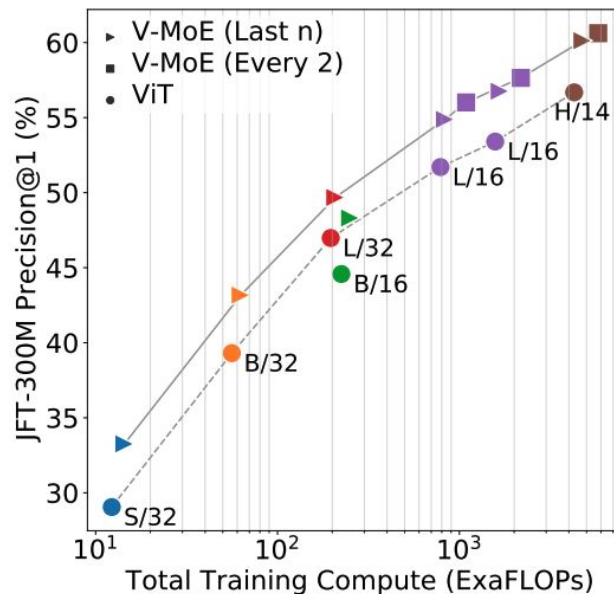
We highlight the following opportunities to improve energy efficiency and CO₂ *equivalent emissions* (CO₂e):

- Large but sparsely activated DNNs can consume <1/10th the energy of large, dense DNNs without sacrificing accuracy despite using as many or even more parameters.

Sparse MoEs for vision

[[Riquelme, Puigcerver, Mustafa et al., 2021](#)]:

- Inspired by successful applications of Sparse MoEs in NLP
- Strong performance vs. FLOPs trade-off



↑ dense to sparse MoEs

Image credit: [Riquelme et al., 2021](#)

Context and motivations

Project at the intersection of two topics:

- **Sparse mixture of experts (sparse MoEs)**
- **Reliability in deep learning**
 - **Why and how to measure it?**
 - **Ensembles**

Uncertainty & reliability in deep learning, why? (1/2)

ML systems are being deployed in many *safety-critical* applications, e.g.,

- **Health applications** [Miotto et al., 2016; Rajkomar et al., 2018; Liu et al., 2020; Mckinney et al., 2020;...]
- **Self-driving cars** [Levinson et al., 2011; Sun et al., 2018]
- **Benefit claims & welfare issues** (e.g., [the Guardian](#), 2019)
- **Dialog systems with LLMs**
- ...

Uncertainty & reliability in deep learning, why? (2/2)

In those applications, we need robust uncertainty estimation

- Knowing when to trust model's predictions, e.g., under dataset shift
- Better decision making, e.g., with asymmetric costs

Uncertainty & reliability in deep learning, why? (2/2)

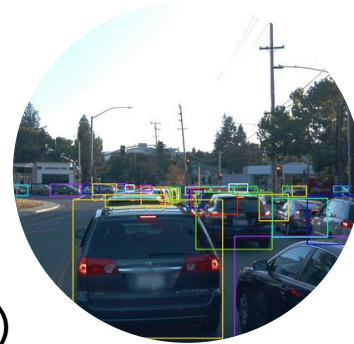
In those applications, we need robust uncertainty estimation

- Knowing when to trust model's predictions, e.g., under dataset shift
- Better decision making, e.g., with asymmetric costs
- Open set recognition
- Lifelong learning
- Active learning, RL, Bayesian optimization

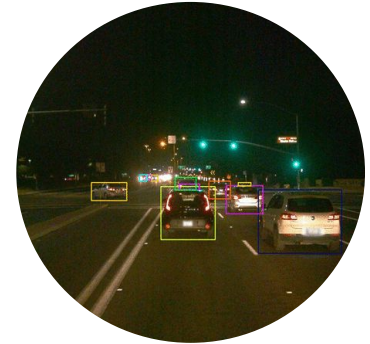
Example: Self-driving cars

Dataset shift:

- Time of day / Lighting
- Geographical location (City vs suburban)
- Changing conditions (Weather / Construction)



Daylight



Night



Weather



Construction



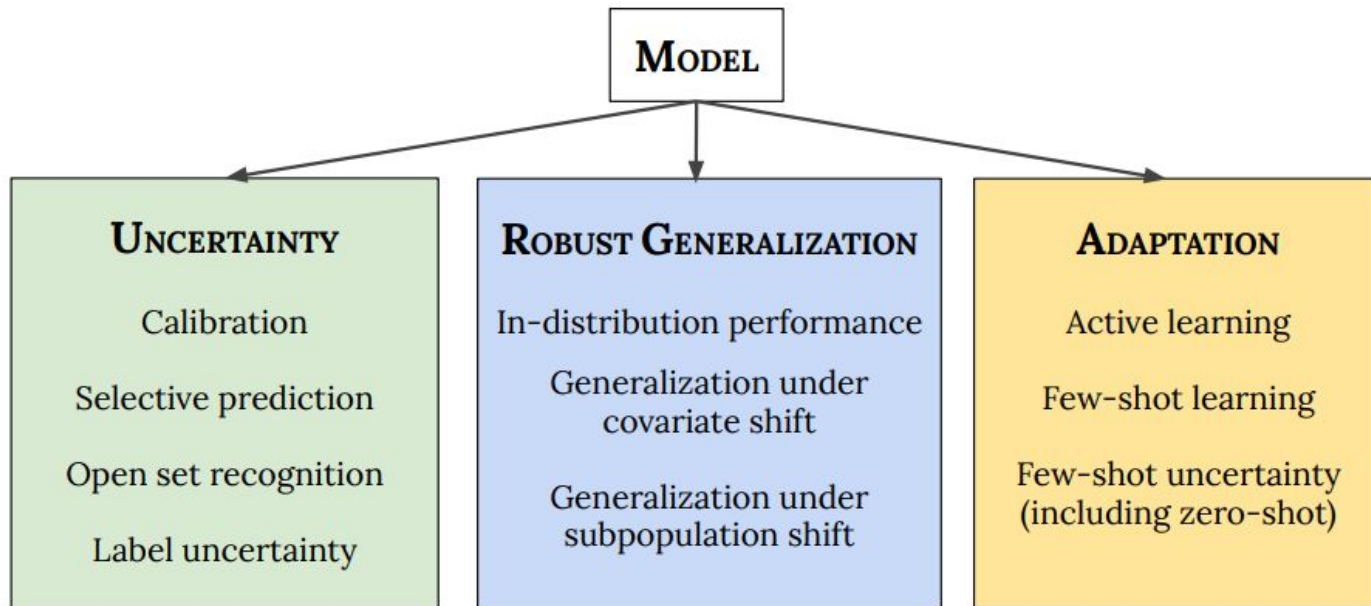
Downtown



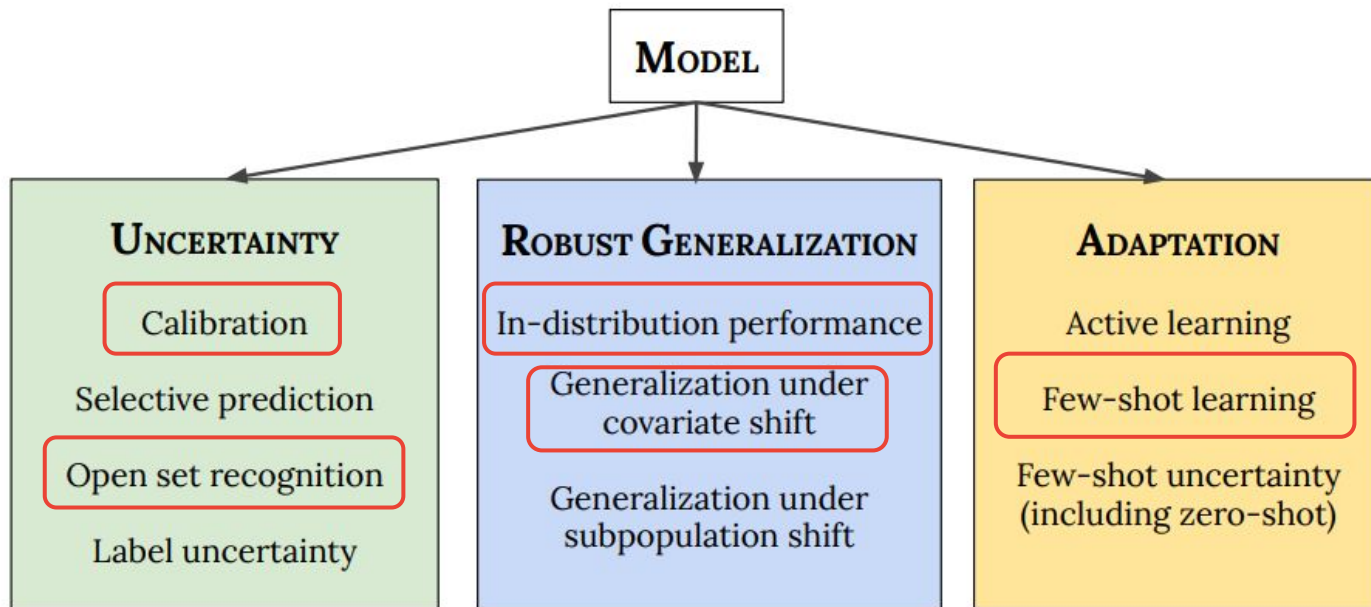
Suburban

Image credit: Sun et al, [Waymo Open Dataset](#)

How can we measure reliability, robustness & uncertainty?



How can we measure reliability, robustness & uncertainty?



Calibration

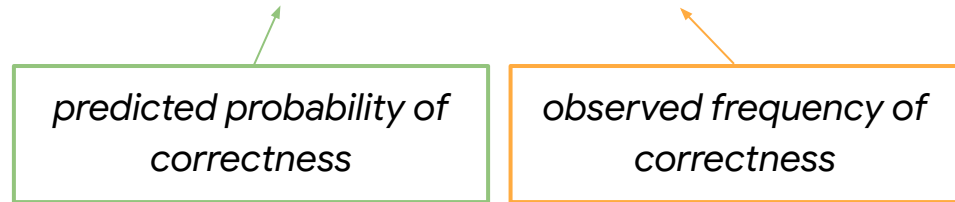
$$\text{Calibration Error} = |\text{Confidence} - \text{Accuracy}|$$

*predicted probability of
correctness*

*observed frequency of
correctness*

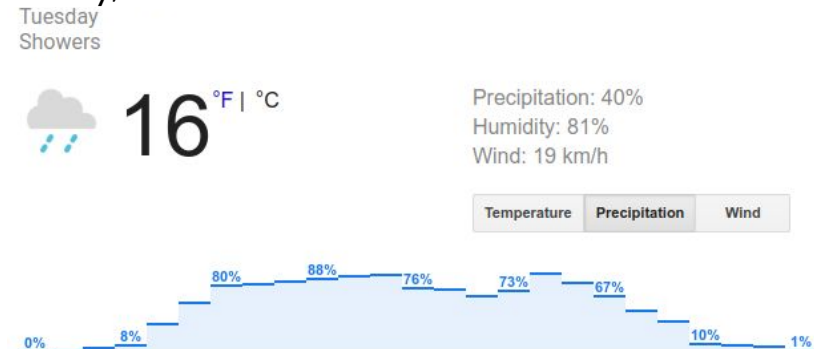
Calibration

$$\text{Calibration Error} = |\text{Confidence} - \text{Accuracy}|$$



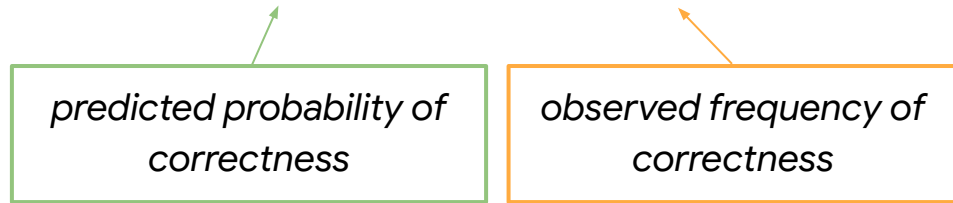
Of all the days where the model predicted rain with 80% probability, what fraction did we observe rain?

- 80% implies perfect calibration
- Less than 80% implies model is *overconfident*
- Greater than 80% implies model is *under-confident*



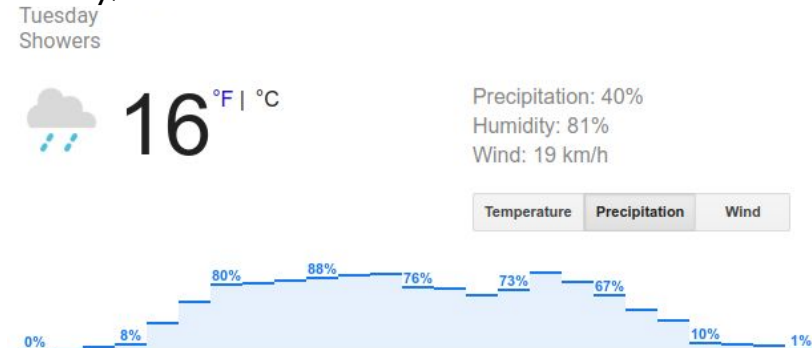
Calibration

$$\text{Calibration Error} = |\text{Confidence} - \text{Accuracy}|$$



Of all the days where the model predicted rain with 80% probability, what fraction did we observe rain?

- 80% implies perfect calibration
- Less than 80% implies model is *overconfident*
- Greater than 80% implies model is *under-confident*
- In practice, use some binning [Naeini et al., 2015]:



$$\text{ECE} = \sum_{b=1}^B \frac{n_b}{N} |\text{acc}(b) - \text{conf}(b)|$$

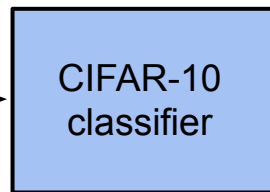
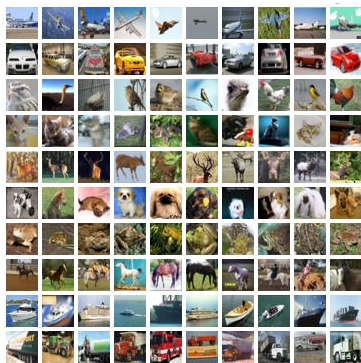
Other ways to quantify reliability & uncertainty?

- **Proper scoring rules** [Gneiting & Raftery, 2007]: Log-likelihood, Brier score, ...

Other ways to quantify reliability & uncertainty?

- **Proper scoring rules** [Gneiting & Raftery, 2007]: Log-likelihood, Brier score, ...
- **Out-of-distribution (OOD) detection:**

CIFAR-10 (test inputs)

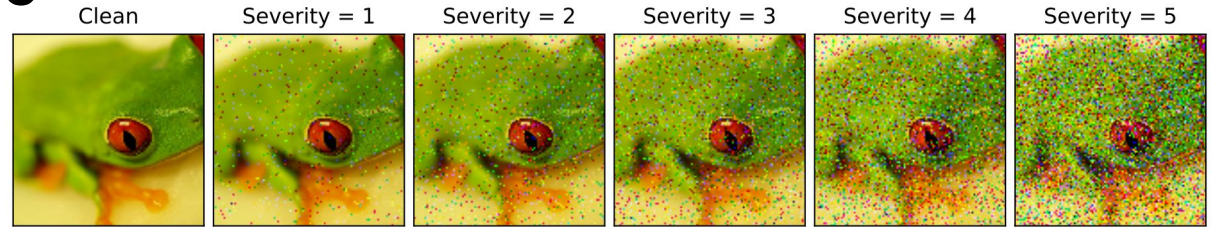


SVHN (OOD inputs)



Confidence on test inputs > Confidence on OOD inputs ? (e.g., via AUC)

Models accuracy degrades under dataset shift



- **Accuracy drops** with increasing shift on Imagenet-C

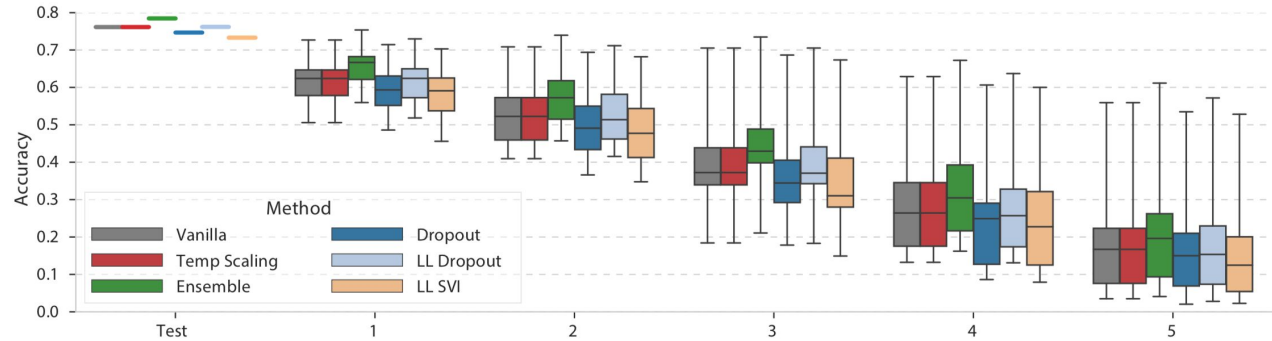
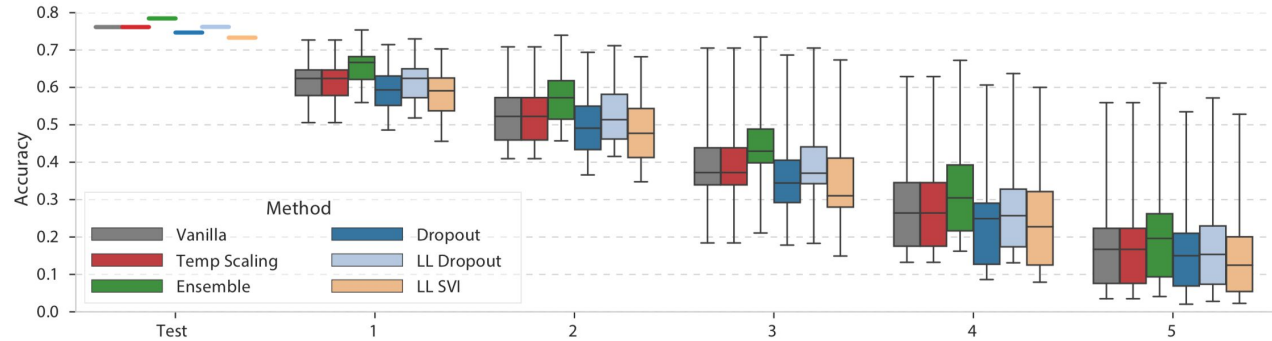
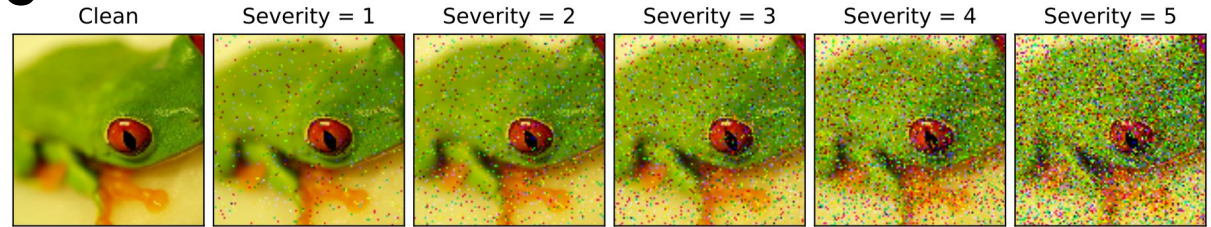


Image source: [Can You Trust Your Model's Uncertainty? Evaluating Predictive Uncertainty Under Dataset Shift?](#), Ovadia et al. 2019

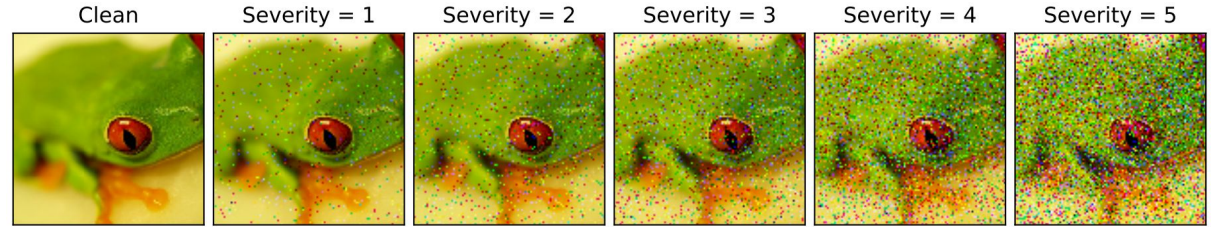
Models accuracy degrades under dataset shift



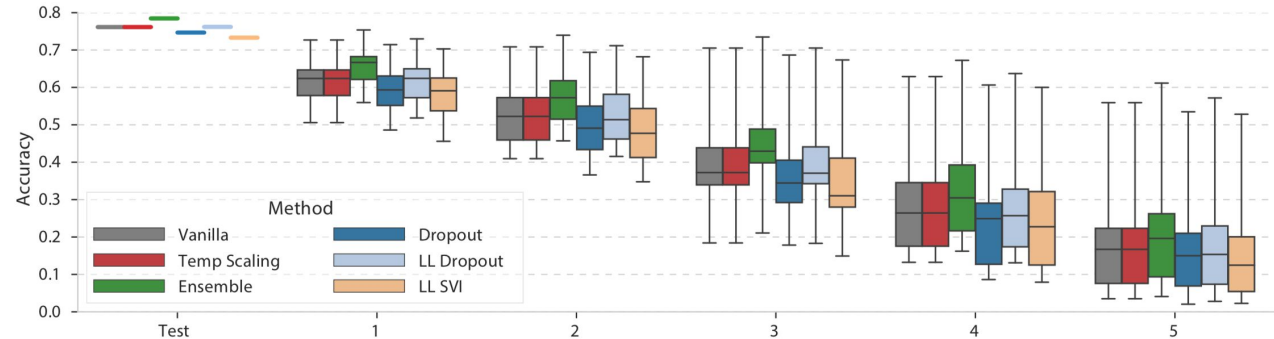
- Accuracy drops with increasing shift on Imagenet-C

- But do the models know that they are less accurate?

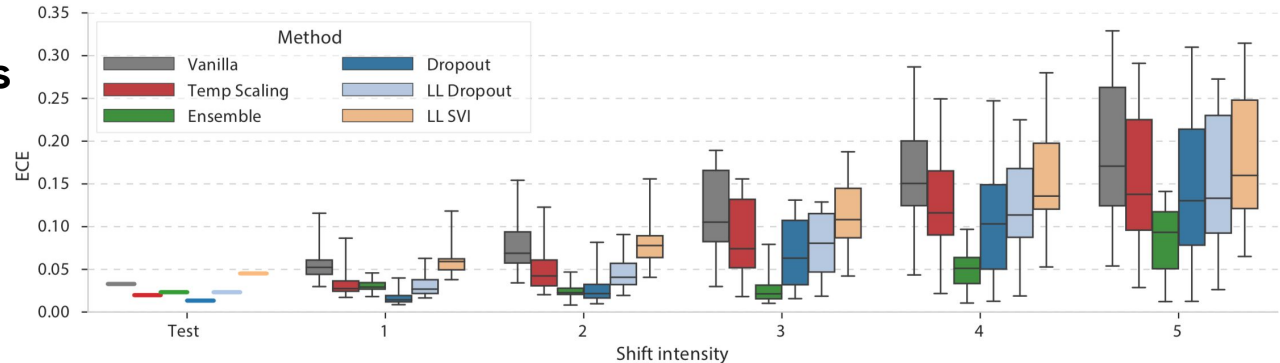
Models are not calibrated under dataset shift



- Accuracy drops with increasing shift on Imagenet-C



- Calibration degrades with shift: “overconfident mistakes”

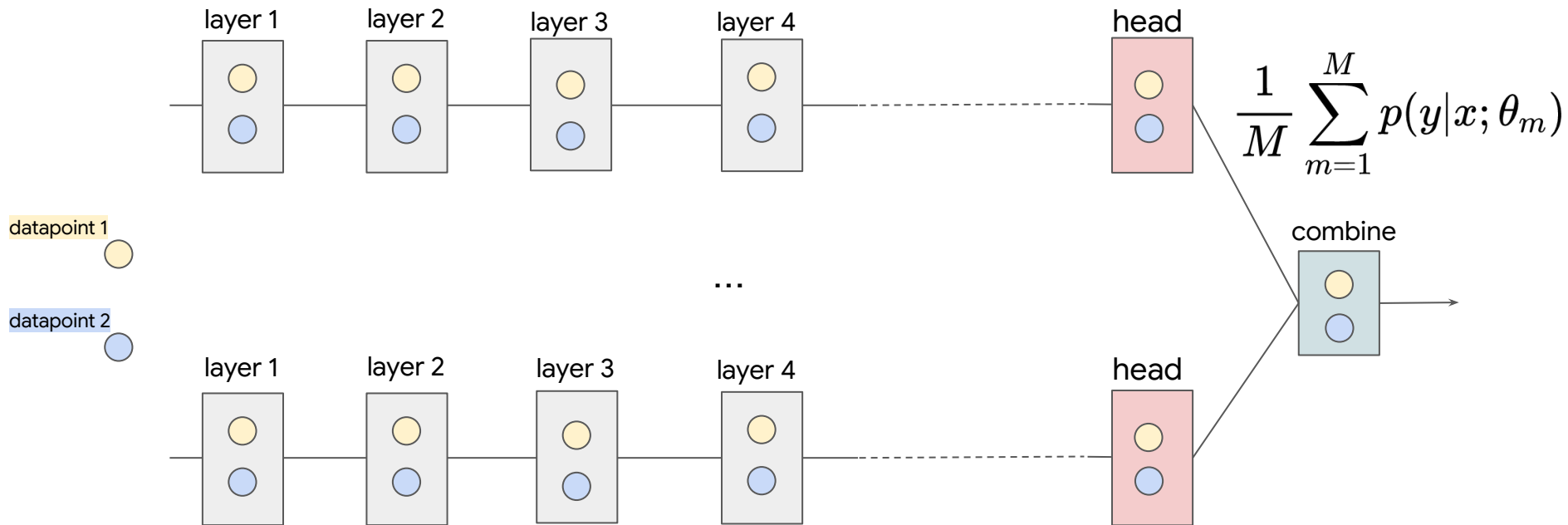


Context and motivations

Project at the intersection of two topics:

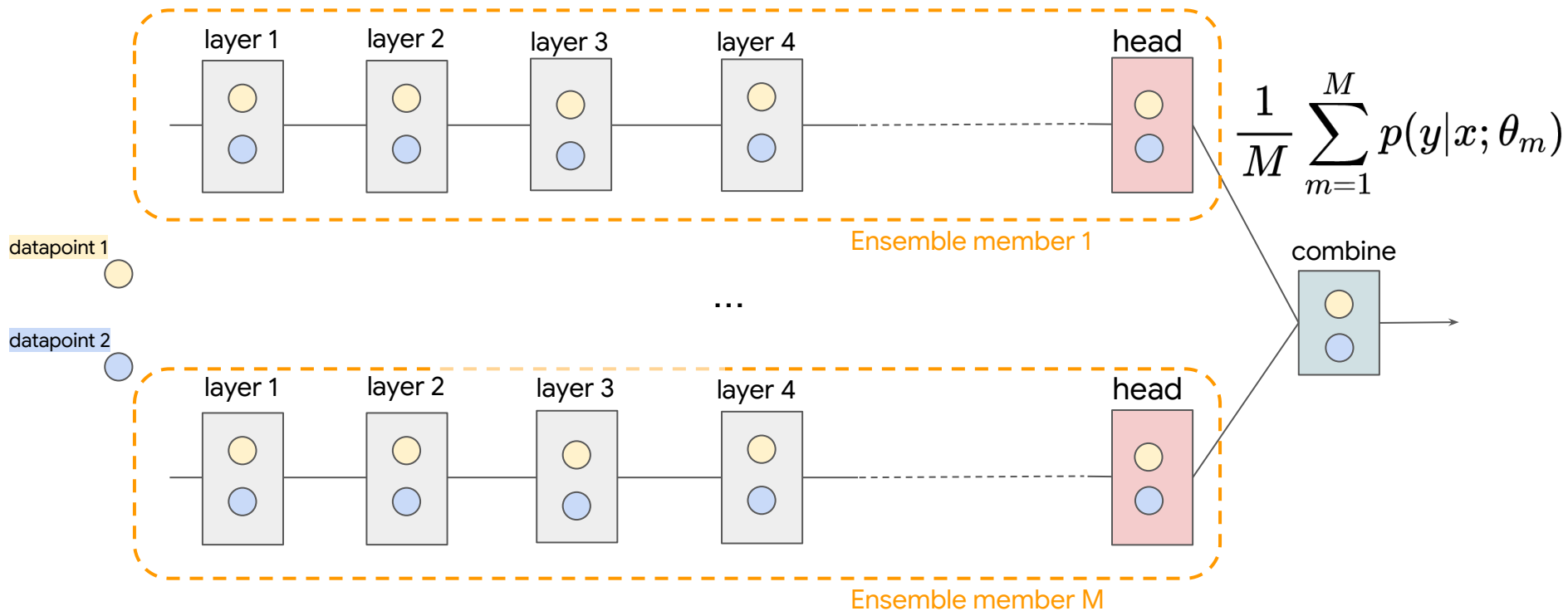
- **Sparse mixture of experts (sparse MoEs)**
- **Reliability in deep learning**
 - **Why and how to measure it?**
 - **Ensembles**

Primer on deep ensembles [Lakshminarayanan et al. 2017, ..., Hansen et al., 1990]



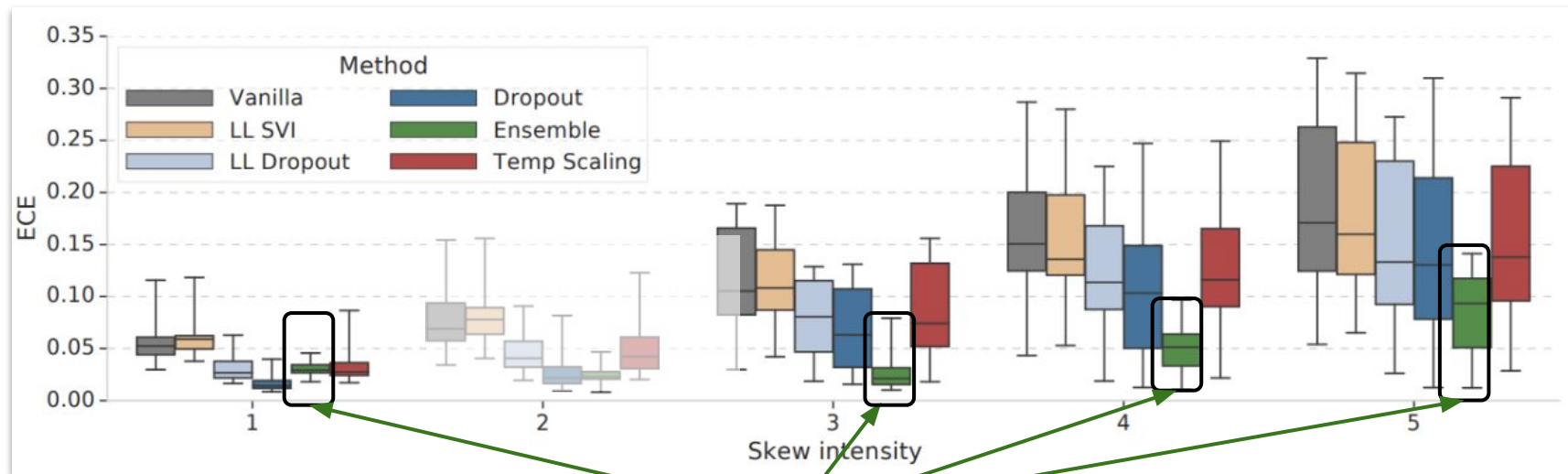
- Multiple trainings from different seeds
- Average the predictions
- Simple...but expensive

Primer on deep ensembles [Lakshminarayanan et al. 2017, ..., Hansen et al., 1990]



- Multiple trainings from different seeds (**M ensemble members**)
- Average the predictions
- Simple...but expensive

Deep ensembles work surprisingly well in practice



Deep ensembles are consistently among the best performing methods, especially under dataset shift

Sparse MoEs vs. ensembles: Overview of properties

Sparse MoEs	Ensembles
Single prediction	Multiple predictions

Sparse MoEs vs. ensembles: Overview of properties

Sparse MoEs	Ensembles
Single prediction	Multiple predictions
Per-example adaptivity	Static combination

Sparse MoEs vs. ensembles: Overview of properties

Sparse MoEs	Ensembles
Single prediction	Multiple predictions
Per-example adaptivity	Static combination
Combination at activation level	Combination at prediction level

Sparse MoEs vs. ensembles: Overview of properties

Sparse MoEs	Ensembles
Single prediction	Multiple predictions
Per-example adaptivity	Static combination
Combination at activation level	Combination at prediction level
???	Robust to distribution shift Strong OOD detection Well-calibrated uncertainty
Strong few-shot performance	???

Sparse MoEs vs. ensembles: Overview of properties

Sparse MoEs	Ensembles
Single prediction	Multiple predictions
Per-example adaptivity	Static combination
Combination at activation level	Combination at prediction level
???	Robust to distribution shift Strong OOD detection Well-calibrated uncertainty
Strong few-shot performance	???
Compute \approx standard NN	Compute \gg standard NN

Sparse MoEs vs. ensembles: Overview of properties

Sparse MoEs	Ensembles
Single prediction	Multiple predictions
Per-example adaptivity	Static combination
Combination at activation level	Combination at prediction level
???	Robust to distribution shift Strong OOD detection Well-calibrated uncertainty
Strong few-shot performance	???
Compute \approx standard NN	Compute \gg standard NN

Goals:

- Understanding the interplay between those two classes of models
- Design approaches “taking the best of both worlds”

Interplay between sparse MoEs & ensembles

Some questions of interest:

- **Can we combine ensembles and sparse MoEs?**

Interplay between sparse MoEs & ensembles

Some questions of interest:

- Can we combine ensembles and sparse MoEs?
- What are the most important factors?

- **M**: The number of ensemble members.
- **K**: The sparsity, i.e., the number of selected experts.
- (**E**: The total number of experts.)

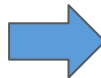
} “Static” combination

} “Adaptive” combination

(Small digression: Experiment setup)

“**Upstream**”, assume checkpoints:

- Pretrained on JFT-300M (~18k classes)
- Vision transformers (ViT) [Dosovitskiy et al., 2020]
- Sparse MoEs (V-MoE) [Riquelme et al., 2021]



“**Downstream**”, fine-tuning:

- ImageNet, Cifar10, Cifar100,...
- ViT, V-MoE
- (all other models assumed compatible with checkpoints)

(Small digression: Experiment setup)



[Image credit](#)

Vision transformers with different backbone sizes: “scale/patch_size”

(Small digression: Experiment setup)



[Image credit](#)

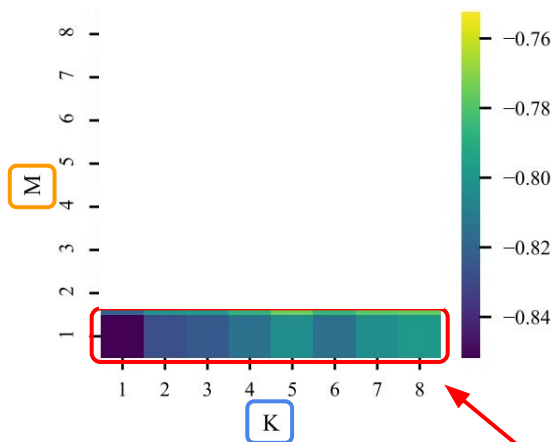
Vision transformers with different backbone sizes: “scale/patch_size”

- S/32, B/32, L/32, B/16, L/16, H/14 (from 36.5M to 2.7B params)

	HIDDEN DIMENSION	MLP DIMENSION	# LAYERS
Small	512	2048	8
Base	768	3072	12
Large	1024	4096	24
Huge	1280	5144	32

“Static” vs. “adaptive” combination

- Deep ensembles of M x V-MoEs with sparsity K
- ImageNet fine-tuned performance

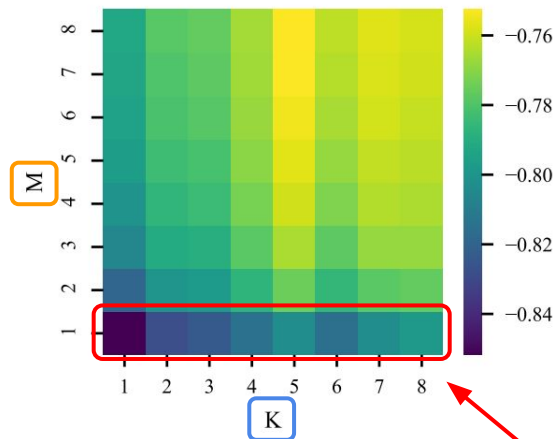


(a) Log-Likelihood

Standard V-MoE (no ensemble)
More selected experts help [Riquelme et al., 2021]

“Static” vs. “adaptive” combination

- Deep ensembles of M x V-MoEs with sparsity K
- ImageNet fine-tuned performance



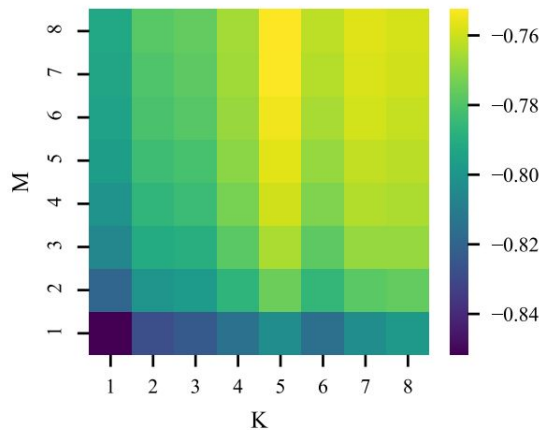
(a) Log-Likelihood

Ensembling beneficial at all sparsity levels

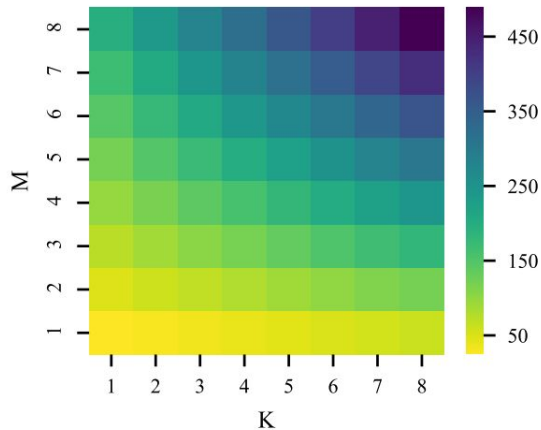
Standard V-MoE (no ensemble)
More selected experts help [Riquelme et al., 2021]

“Static” vs. “adaptive” combination

- What about the cost? (GFLOPs = 10^9 FLOPs)



(a) Log-Likelihood

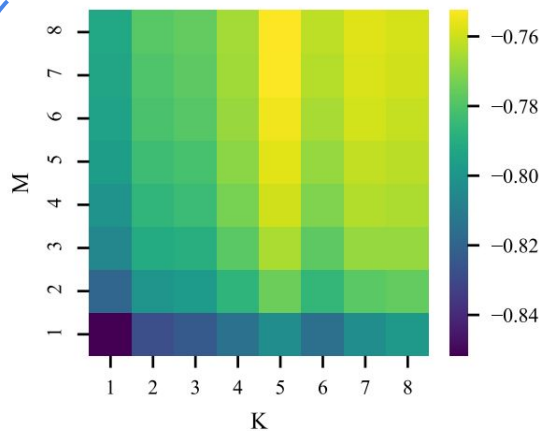


(b) (downstream) GFLOPs

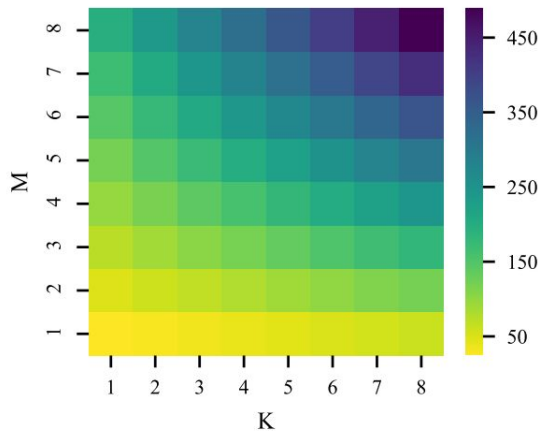
Ensembling more expensive

“Static” vs. “adaptive” combination

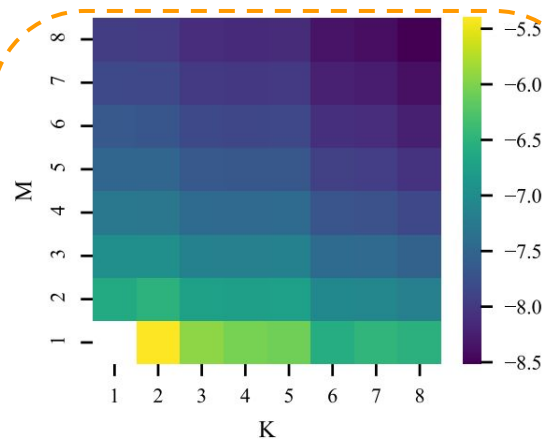
- More cost-effective to use “adaptive” combination
- ...but best absolute performance requires both types of combinations



(a) Log-Likelihood

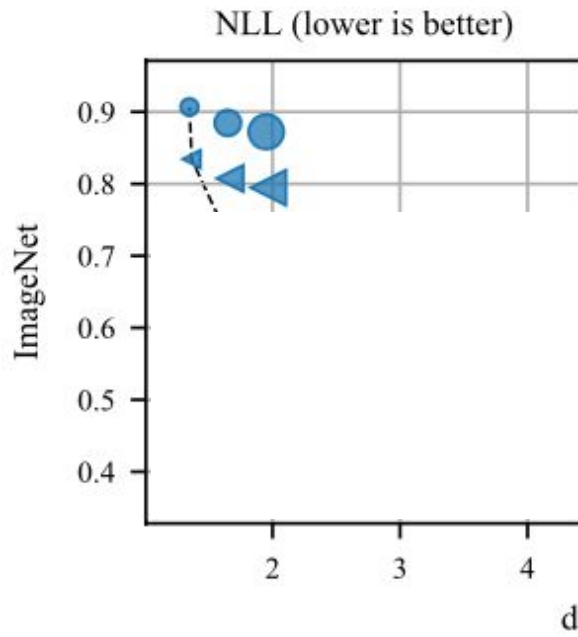


(b) (downstream) GFLOPs



(c) $\log \left[\frac{LL_{(K,M)} - LL_{(1,1)}}{GFLOPs_{(K,M)} - GFLOPs_{(1,1)}} \right]$

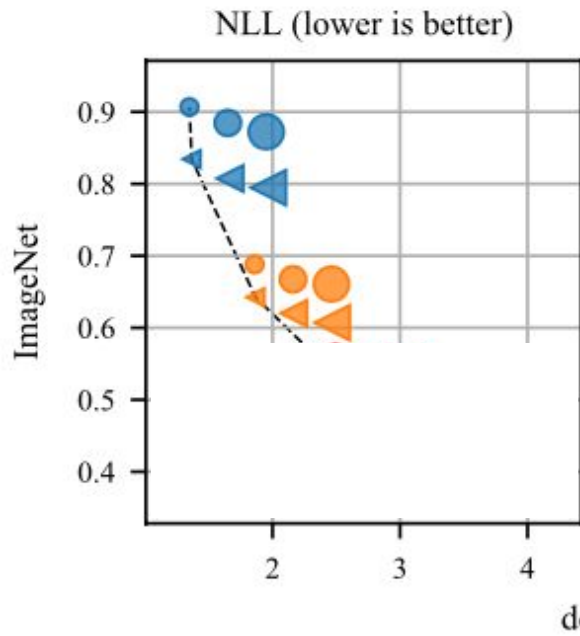
Is ensembling equally helpful for ViT & V-MoE? (ImageNet)



Ensembling helps ViT as much as V-MoE
(Sparsity $K=1$, M in $\{1, 2, 4\}$)



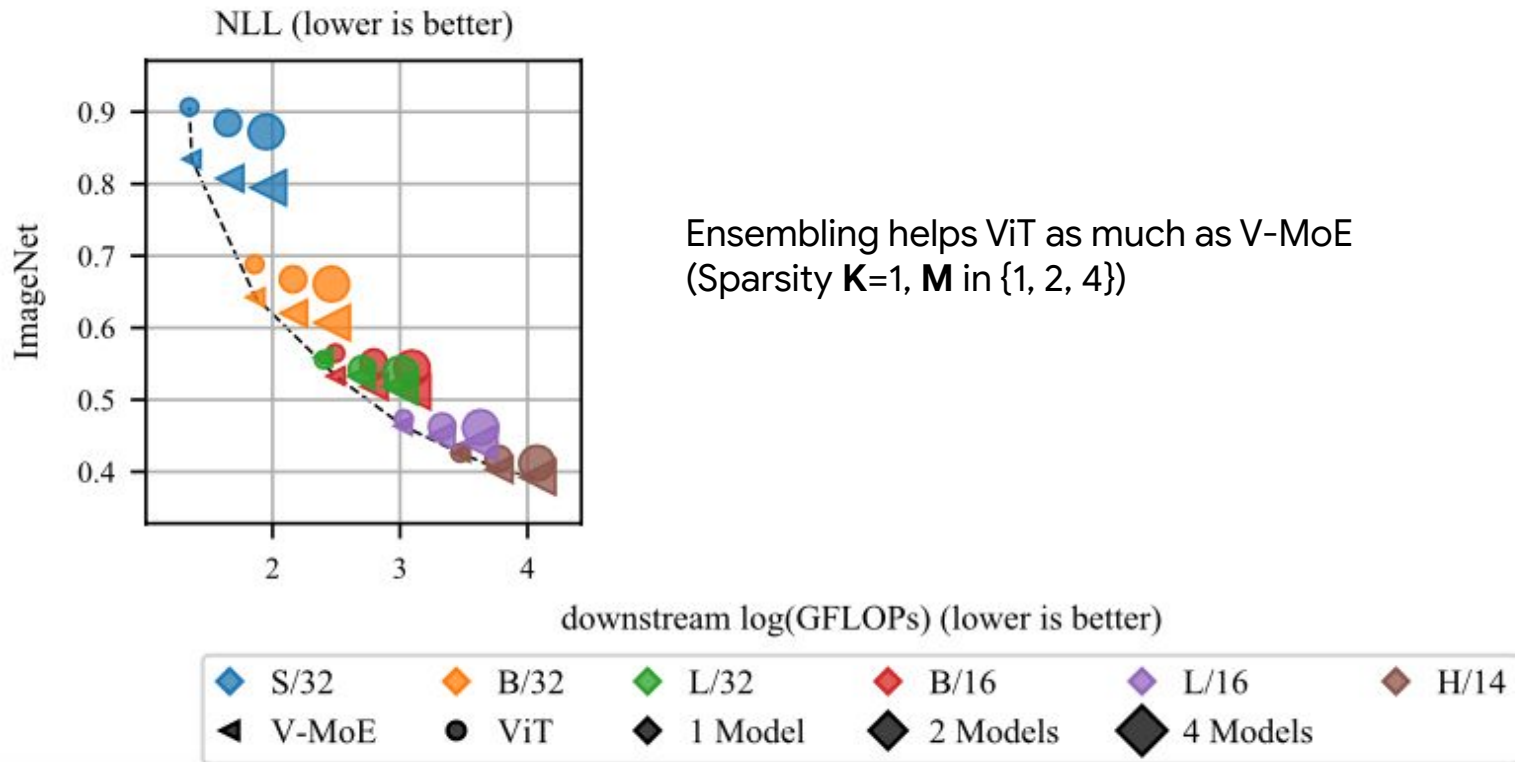
Is ensembling equally helpful for ViT & V-MoE? (ImageNet)



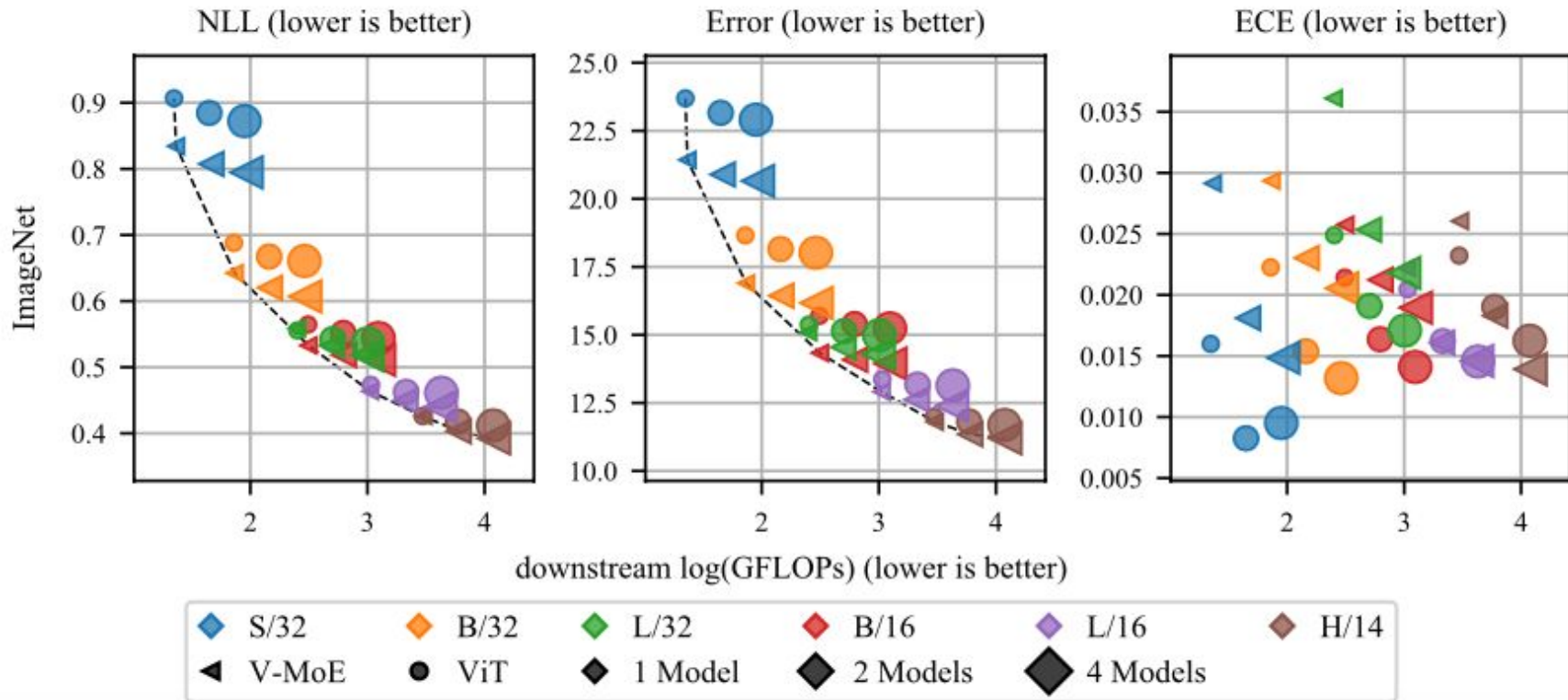
Ensembling helps ViT as much as V-MoE
(Sparsity $K=1$, M in $\{1, 2, 4\}$)



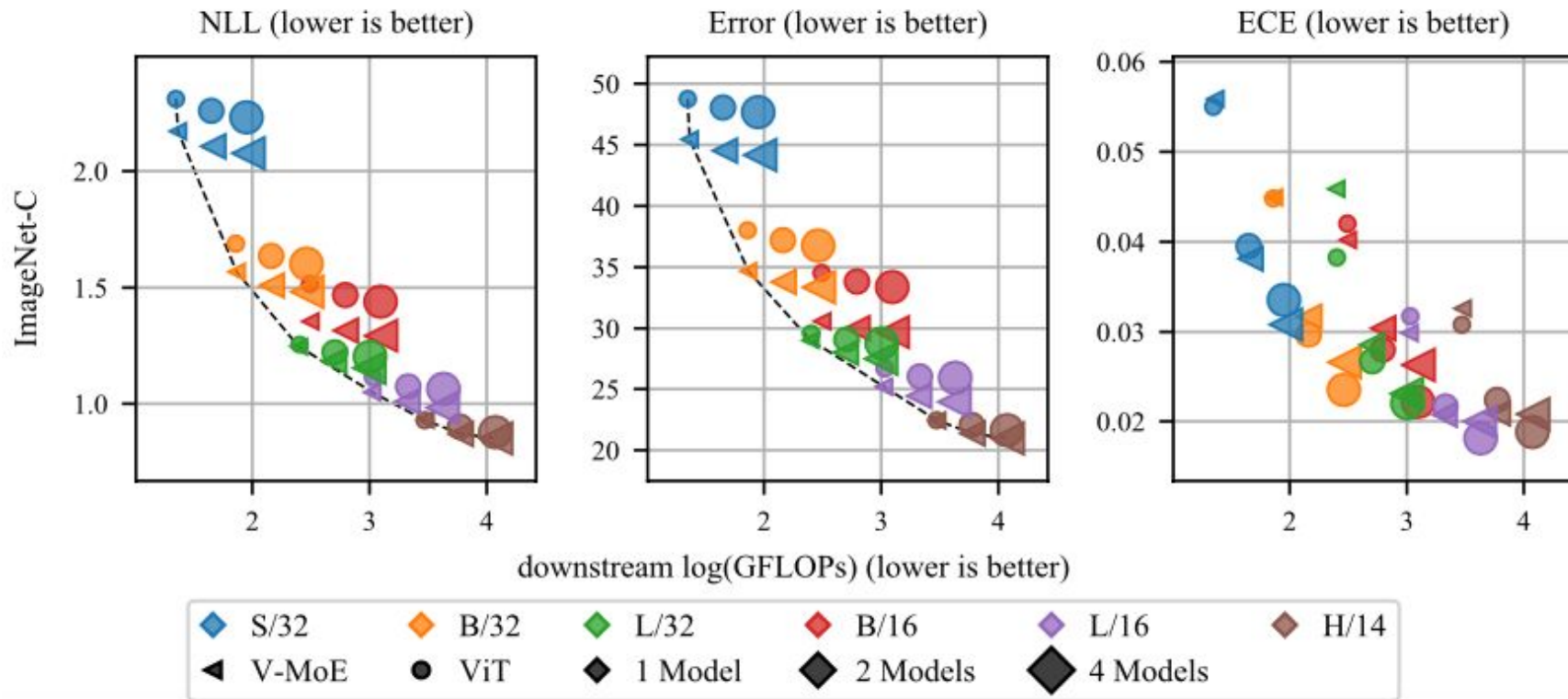
Is ensembling equally helpful for ViT & V-MoE? (ImageNet)



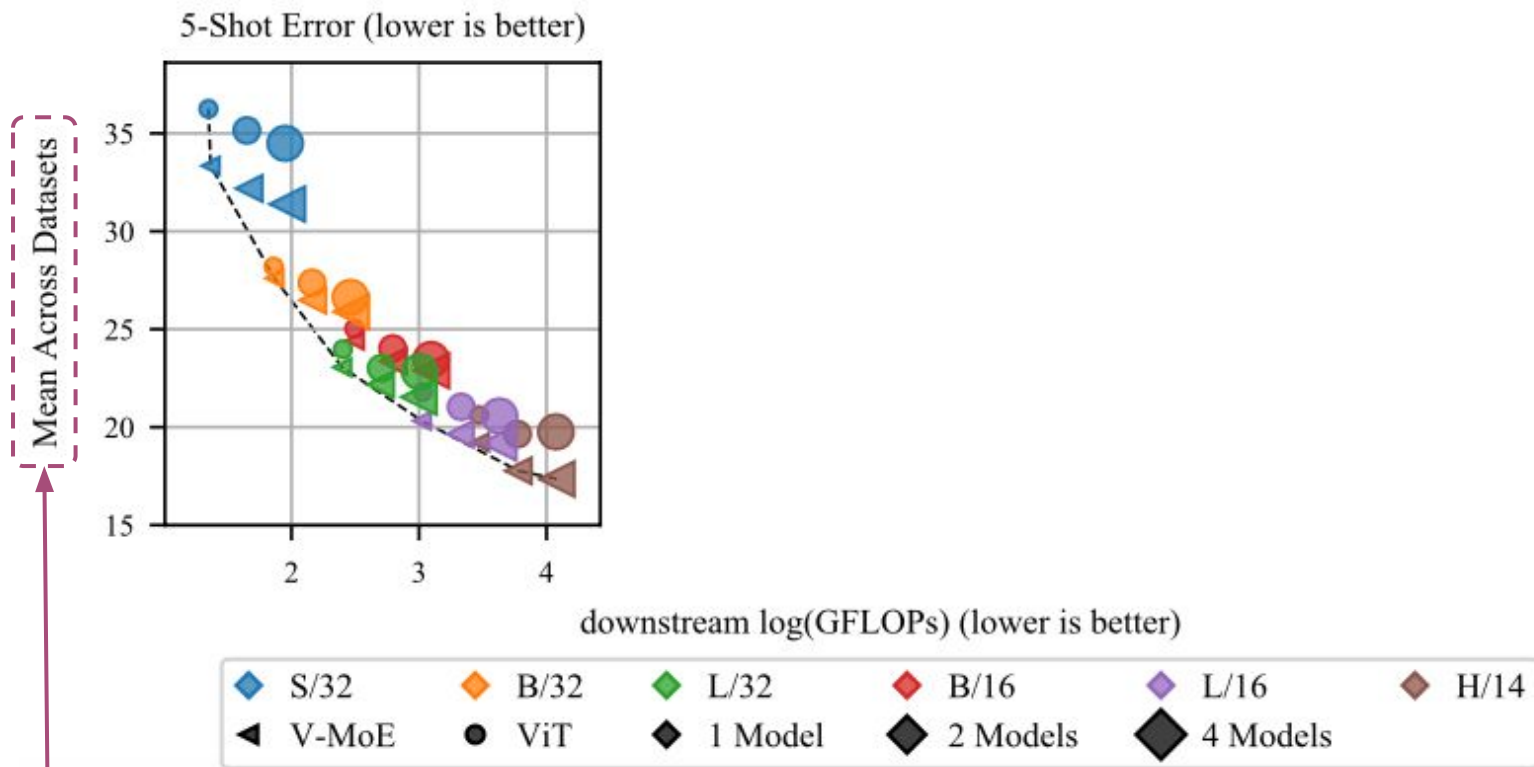
Is ensembling equally helpful for ViT & V-MoE? (ImageNet)



Is ensembling equally helpful for ViT & V-MoE? (ImageNet-C)

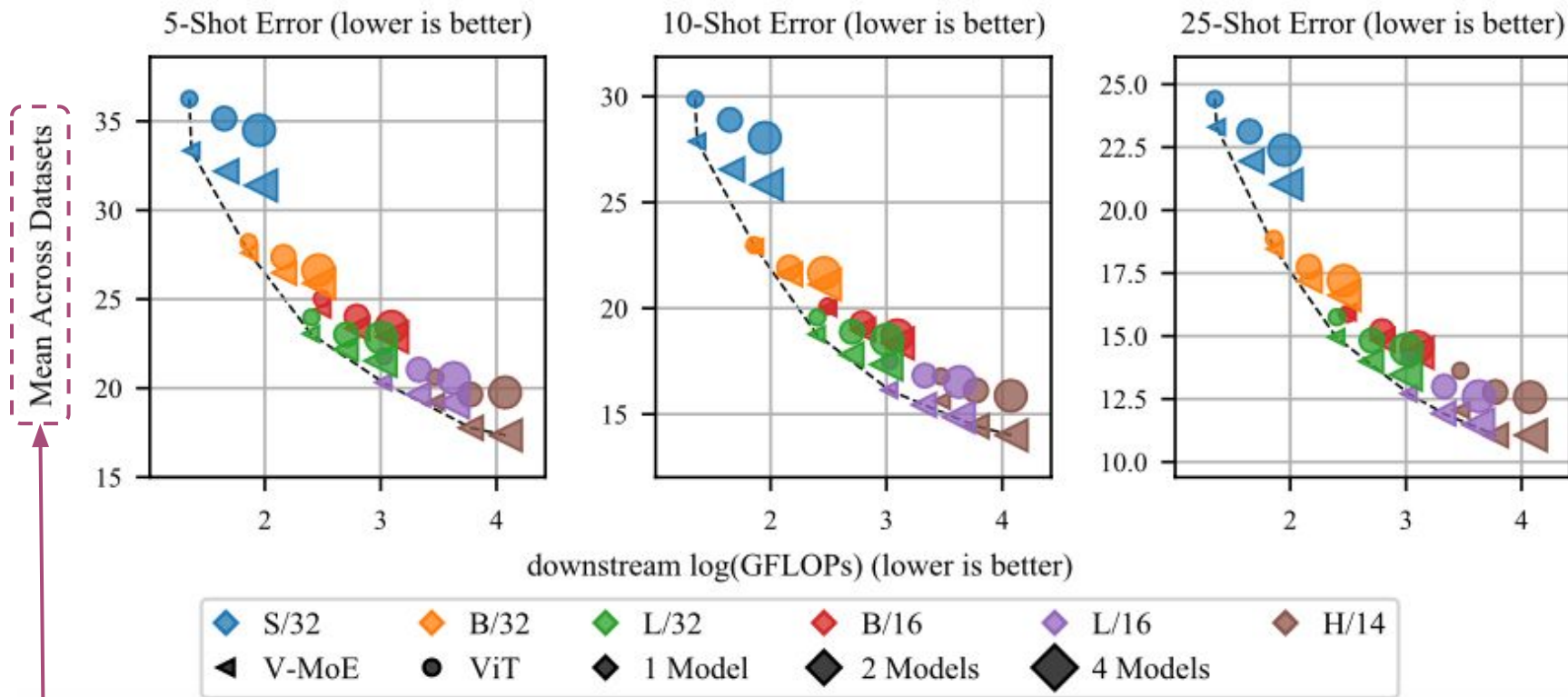


Is ensembling equally helpful for ViT & V-MoE? (few-shot)



Linear few-shot [Dosovitskiy et al., 2020, Riquelme et al., 2021], aggregated over 8 datasets

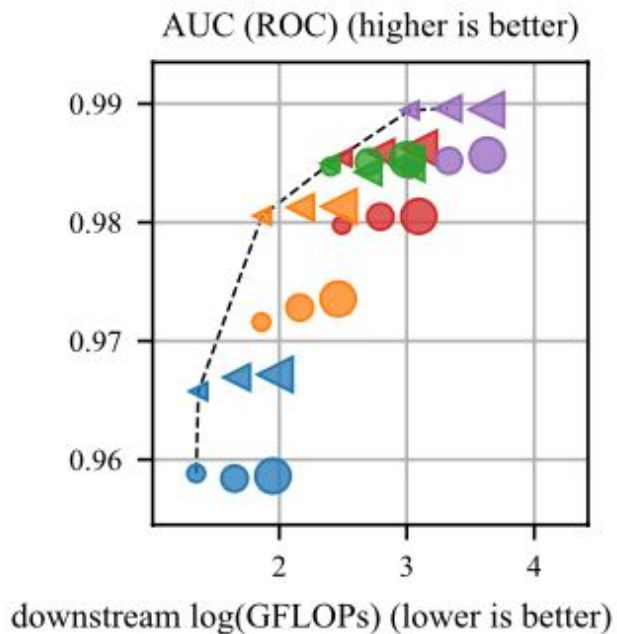
Is ensembling equally helpful for ViT & V-MoE? (few-shot)



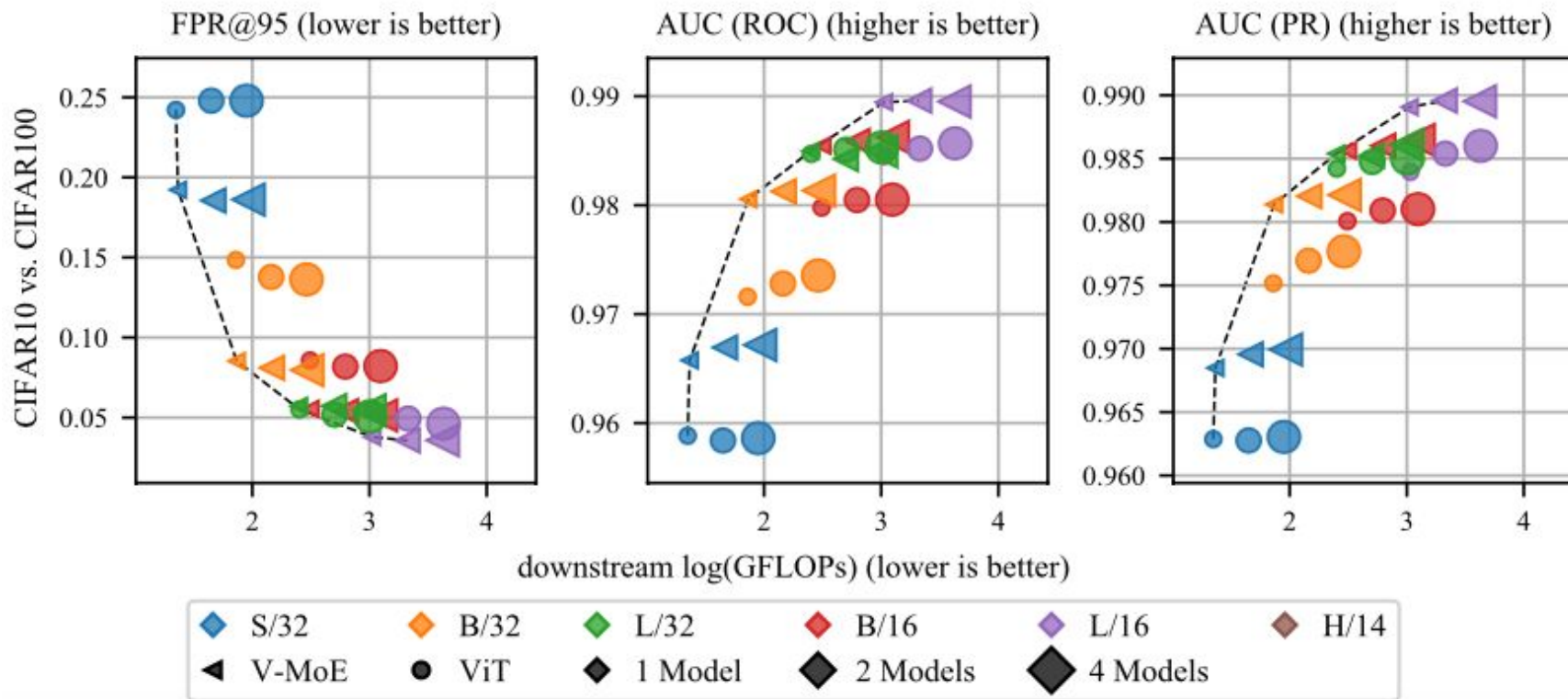
Linear few-shot [Dosovitskiy et al., 2020, Riquelme et al., 2021], aggregated over 8 datasets

Is ensembling equally helpful for ViT & V-MoE? (OOD detection)

CIFAR10 vs. CIFAR100



Is ensembling equally helpful for ViT & V-MoE? (OOD detection)



Summary so far:

Some questions of interest:

- **Can we combine ensembles and sparse MoEs?** →  ... but costly 

Summary so far:

Some questions of interest:

- **Can we combine ensembles and sparse MoEs?** →



... but costly



- **What are the most important factors?**

- **M**: The number of ensemble members.
- **K**: The sparsity, i.e., the number of selected experts.

} “**Adaptive**” combination more cost-effective but both “**static**” & “**adaptive**” are needed.

Summary so far:

Some questions of interest:

- **Can we combine ensembles and sparse MoEs?** →



... but costly



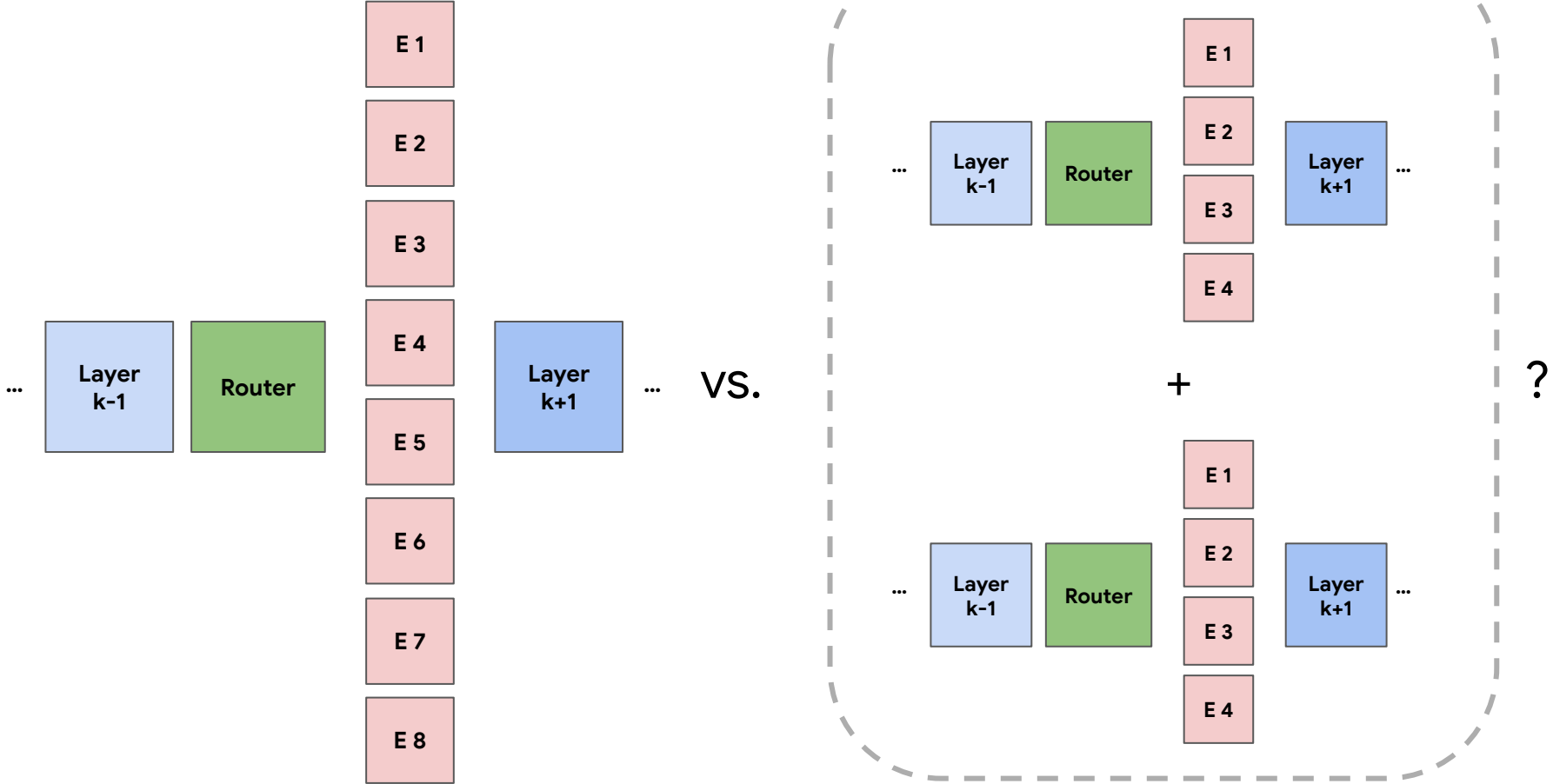
- **What are the most important factors?**

- **M**: The number of ensemble members.
- **K**: The sparsity, i.e., the number of selected experts.

} “**Adaptive**” combination more cost-effective but both “**static**” & “**adaptive**” are needed.

- **Can we do better than a naive ensemble of sparse MoEs?**

Motivating experiment



Motivating experiment: V-MoE-B/32 (K=1, E=32) on ImageNet

M	E	NLL ↓	ERROR ↓
1	32	0.642 ± 0.002	16.90 ± 0.05
2	32	0.588	15.74
4	32	0.561	15.10

Already observed:
Ensembling V-MoEs helps

Motivating experiment: V-MoE-B/32 (K=1, E=32) on ImageNet

M	E	NLL ↓	ERROR ↓
1	32	0.642 ± 0.002	16.90 ± 0.05
2	32	0.588	15.74
4	32	0.561	15.10
2	16	0.588	15.97
4	8	0.577	15.82

Main observation:

Ensemble of smaller sparse MoEs > One single larger sparse MoE

Efficient Ensemble of Experts (E^3)

Constructing an efficient ensembles of sparse MoEs:

- **Idea:** End-to-end training of M x sparse MoEs with E/M experts

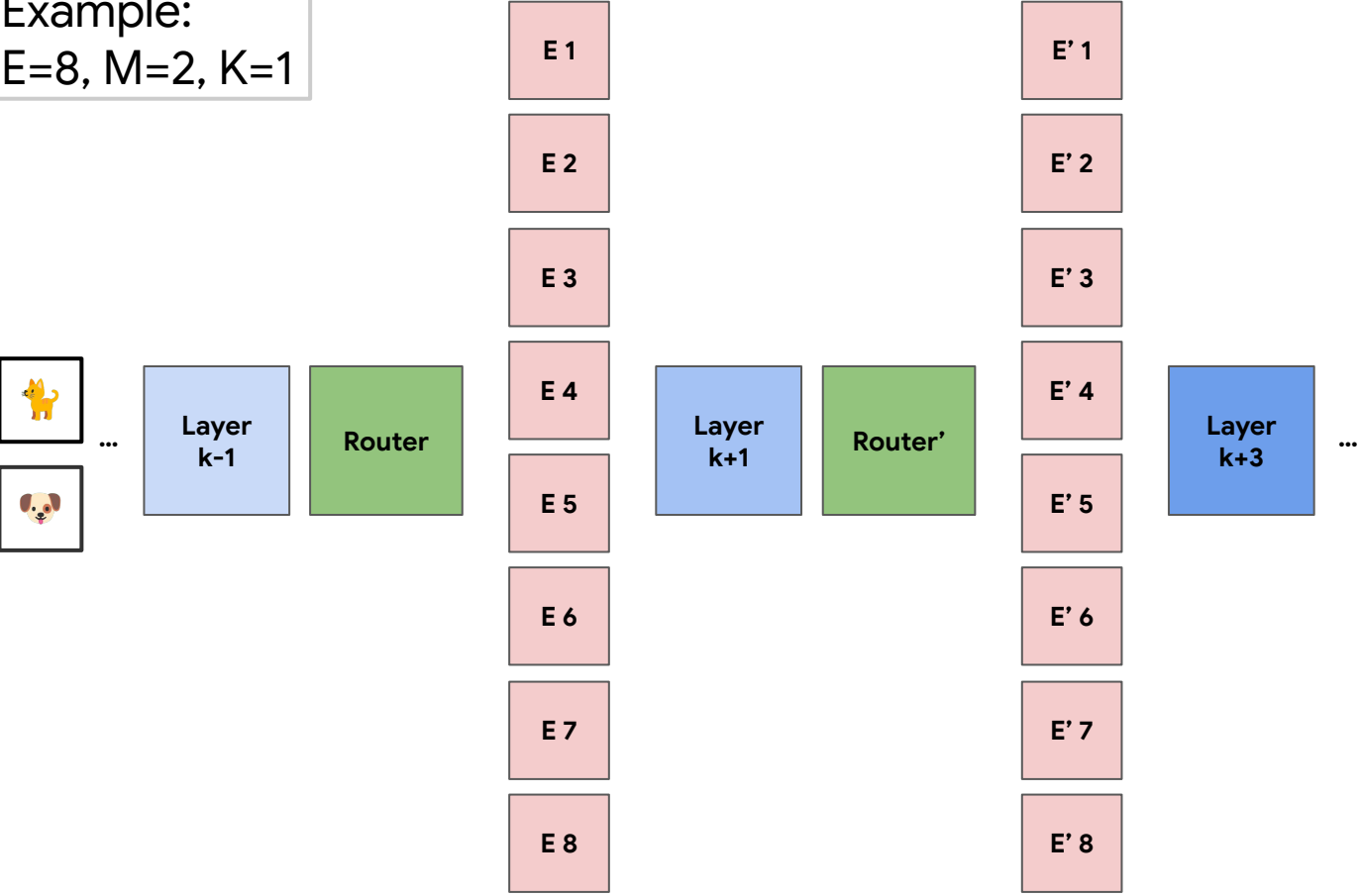
Efficient Ensemble of Experts (E^3)

Constructing an efficient ensembles of sparse MoEs:

- **Idea:** End-to-end training of M x sparse MoEs with E/M experts
 - An ensemble member = one sparse MoEs with E/M experts
 - Efficient simultaneous training via a tiled representation
 - Sharing of parameters in non-expert layers
- Inspired by batch ensemble [Wen et al., 2019]

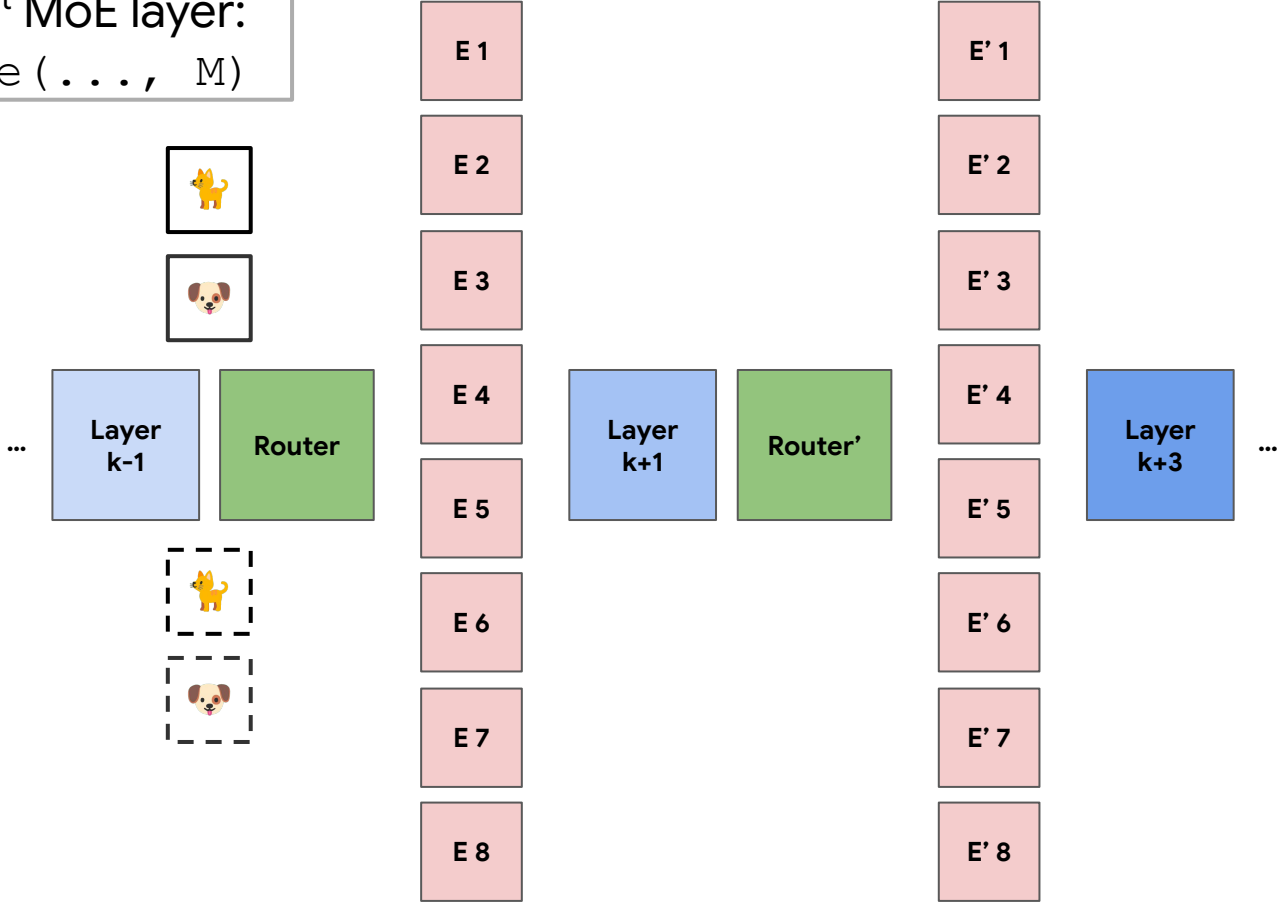
Pictorial view of Efficient Ensemble of Experts (E³)

Example:
E=8, M=2, K=1



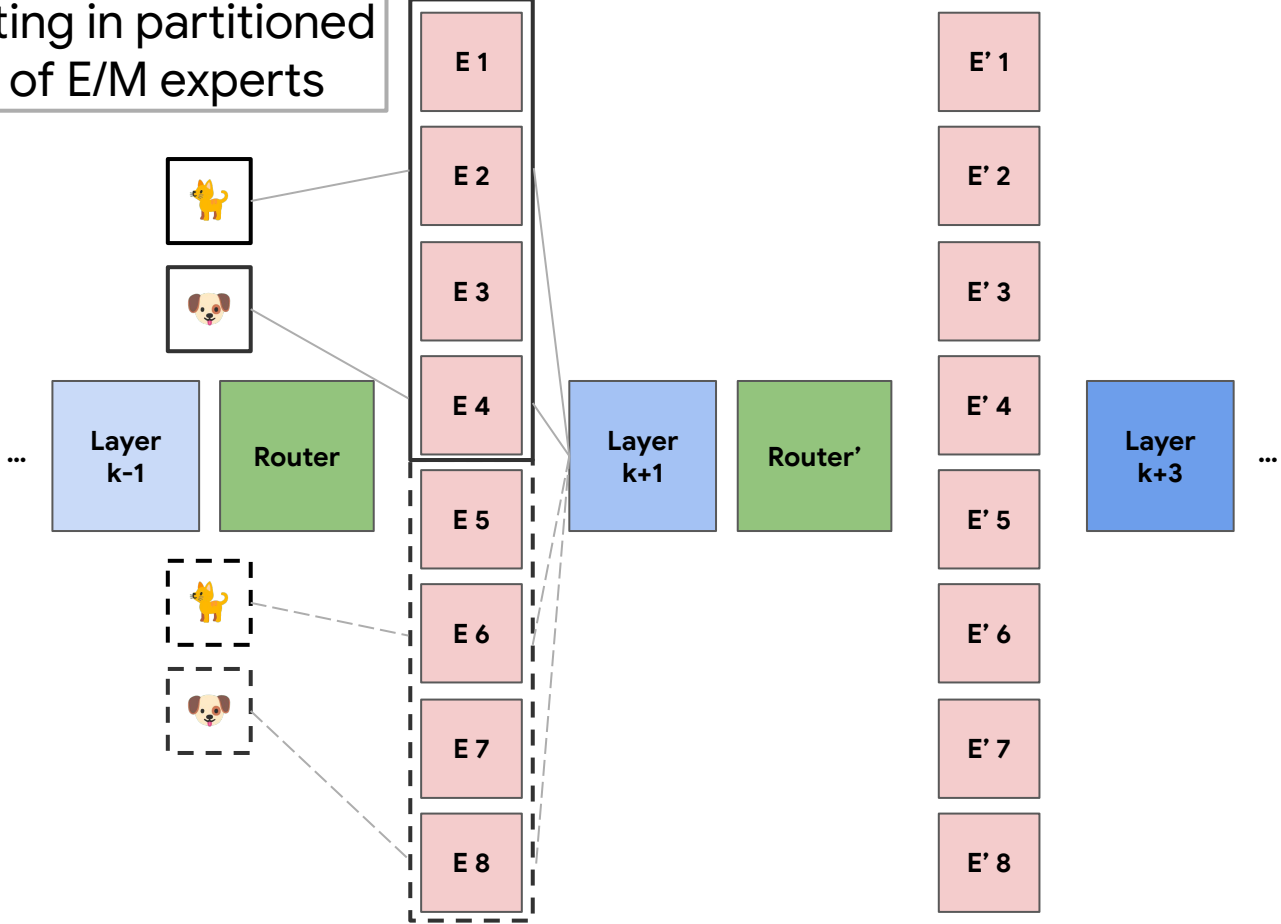
Pictorial view of Efficient Ensemble of Experts (E³)

At 1st MoE layer:
`tile(..., M)`



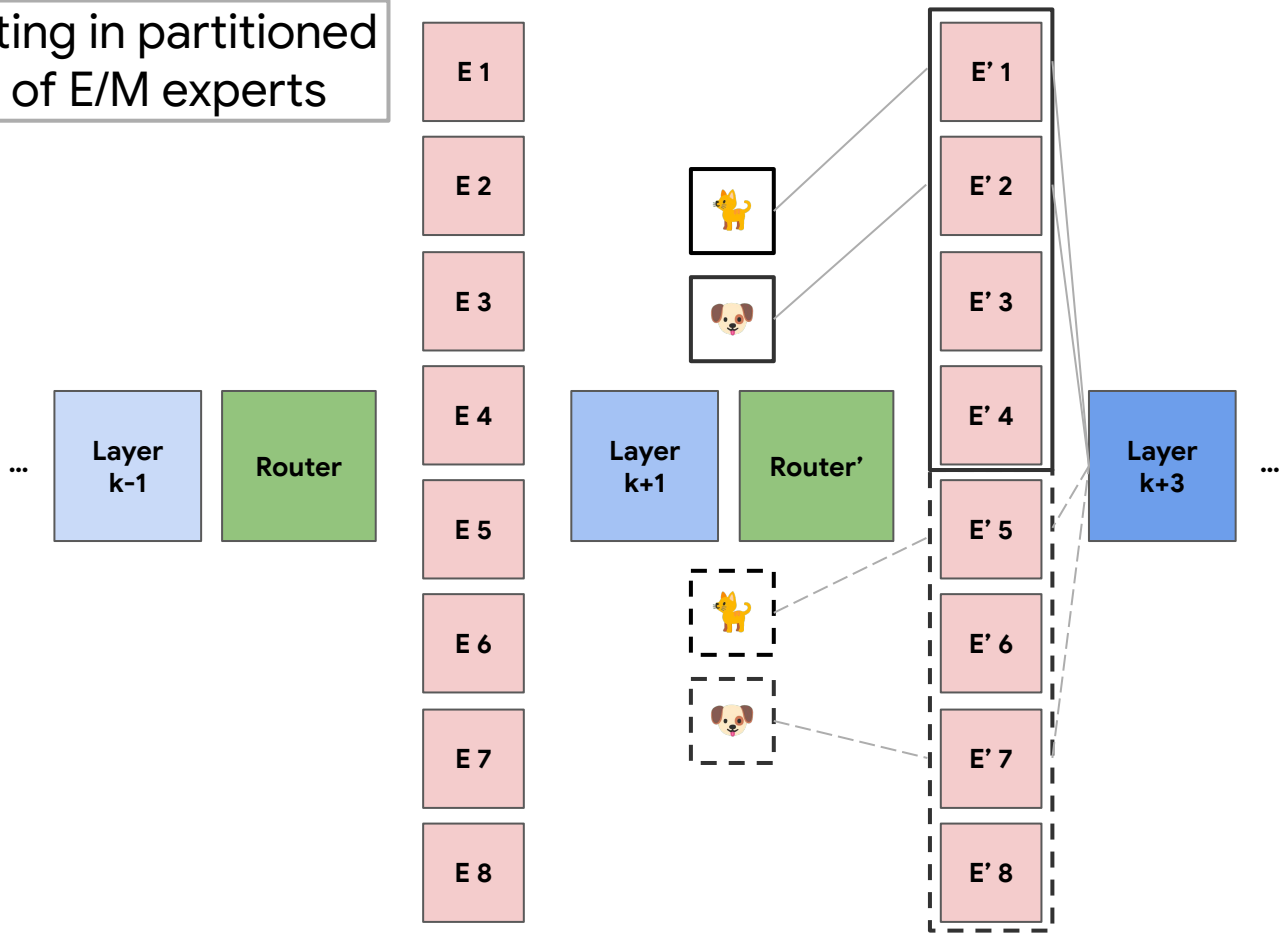
Pictorial view of Efficient Ensemble of Experts (E³)

Routing in partitioned sets of E/M experts

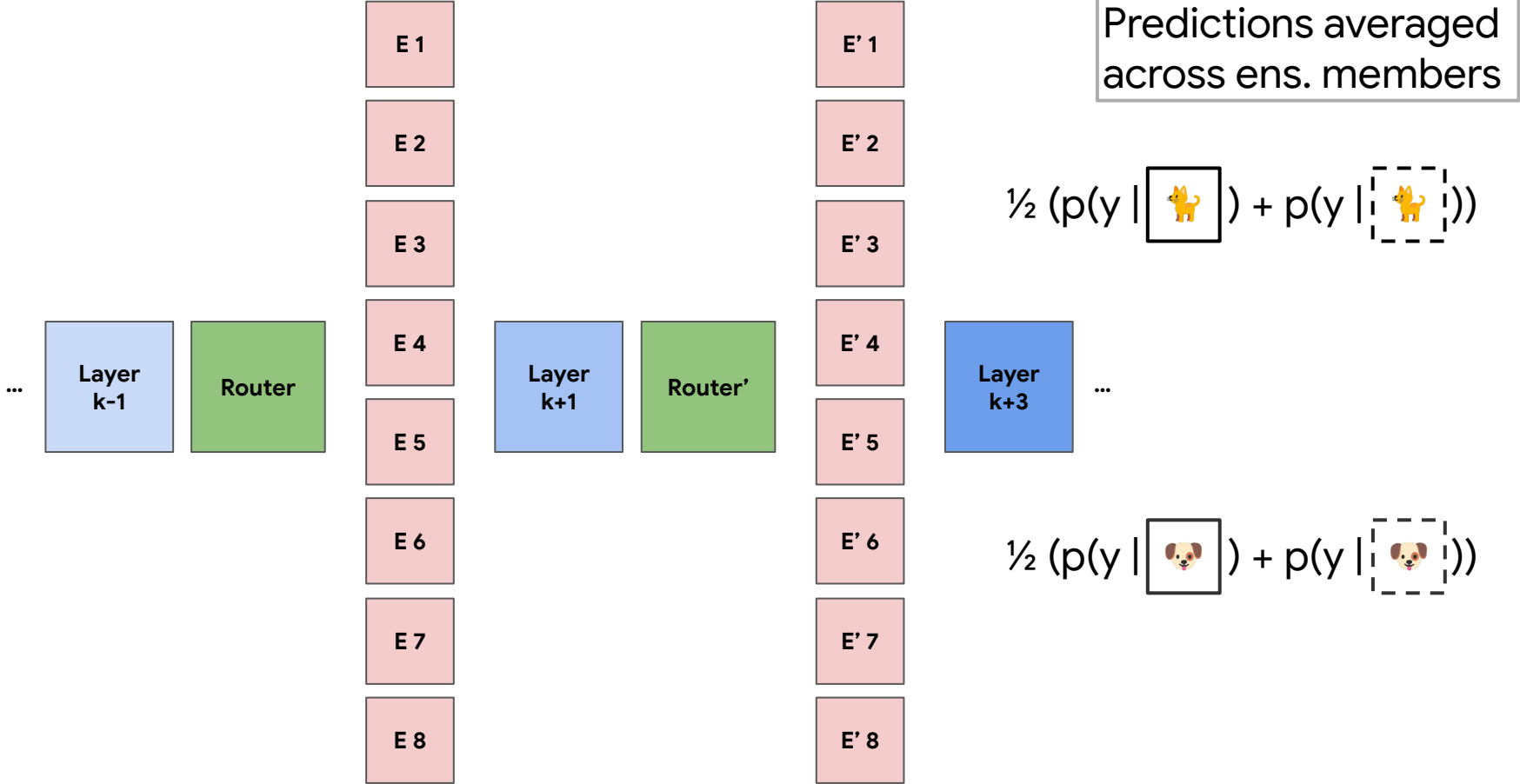


Pictorial view of Efficient Ensemble of Experts (E^3)

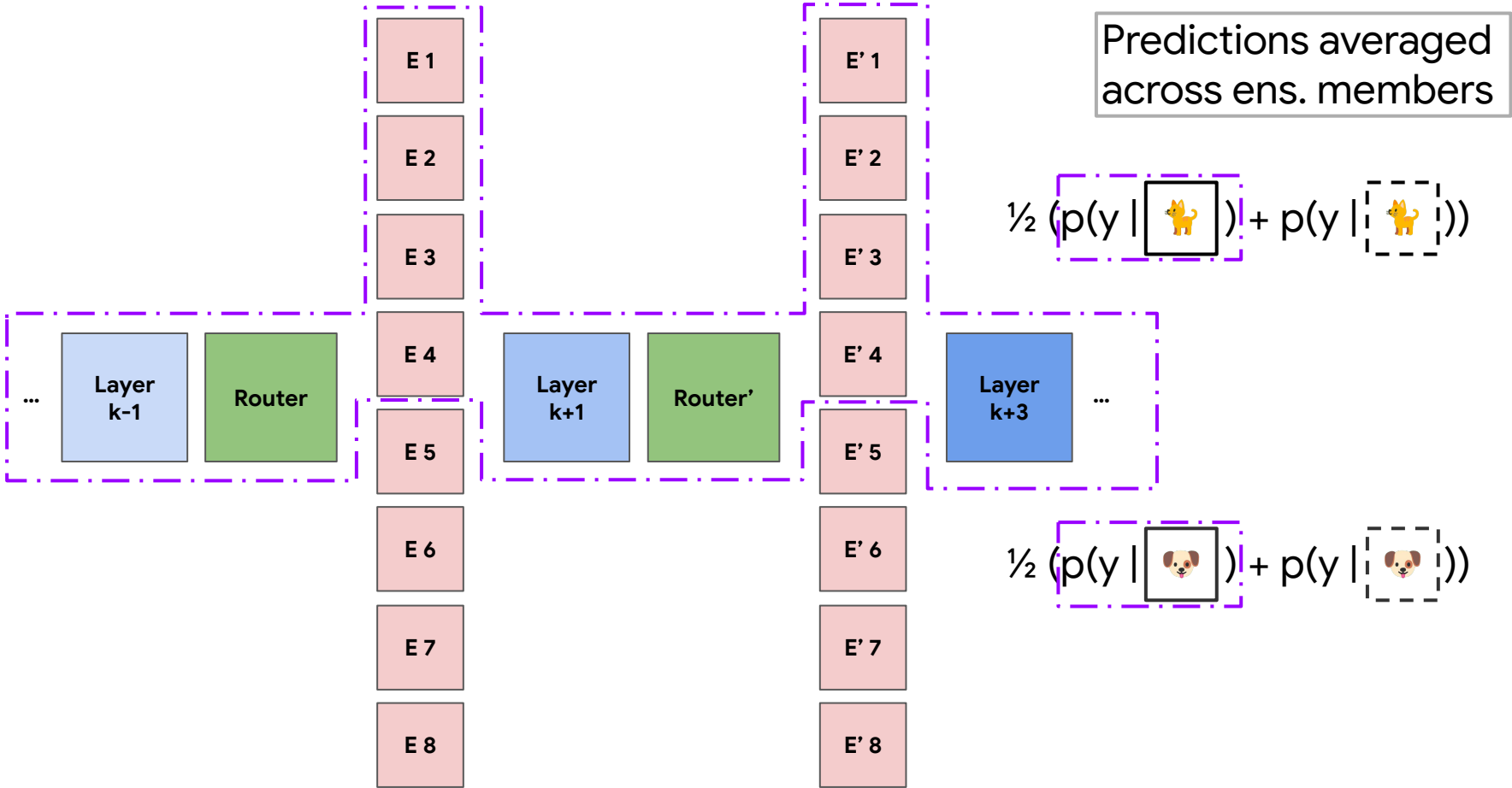
Routing in partitioned sets of E/M experts



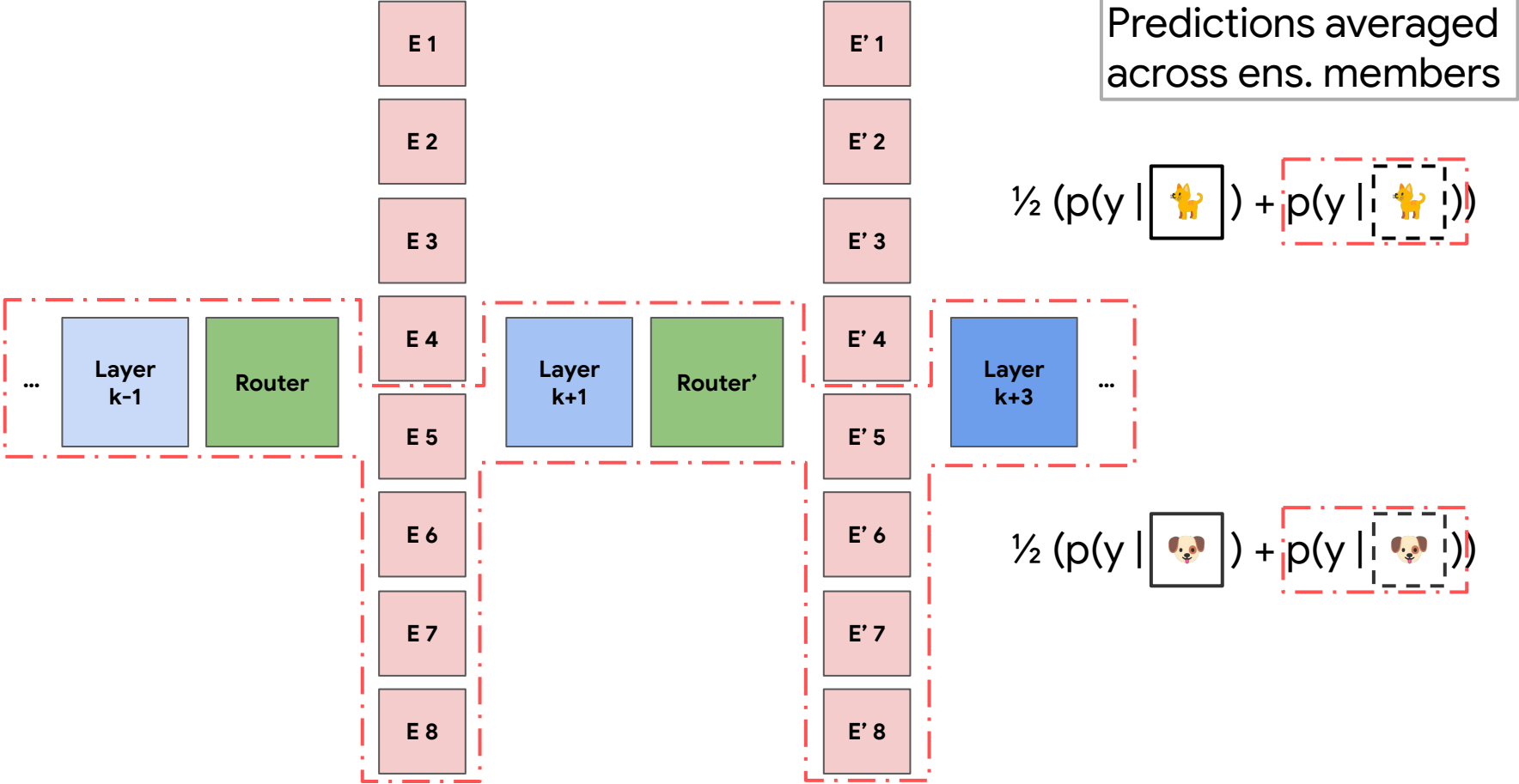
Pictorial view of Efficient Ensemble of Experts (E³)



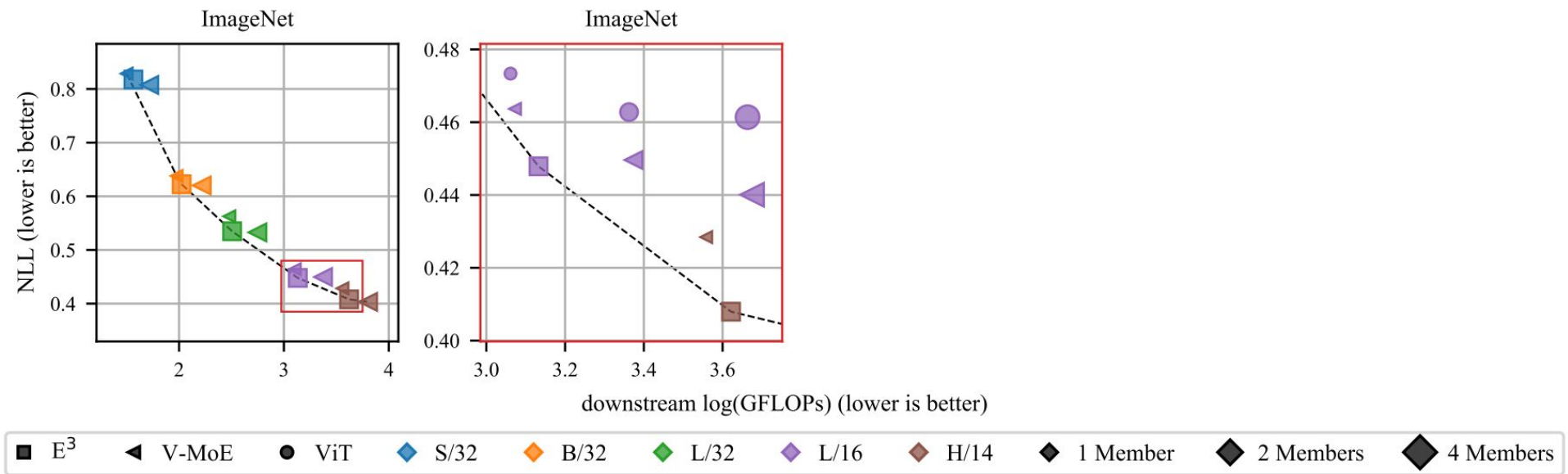
Pictorial view of Efficient Ensemble of Experts (E³)



Pictorial view of Efficient Ensemble of Experts (E³)

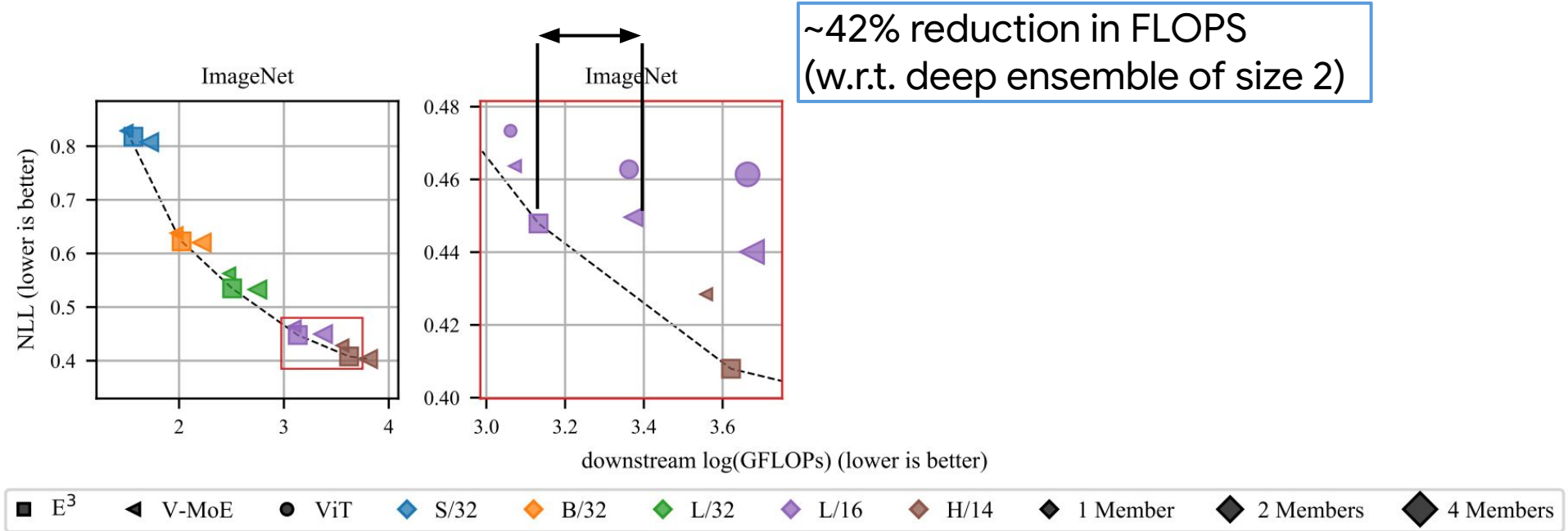


Evaluation of E³ (ImageNet & few-shot)



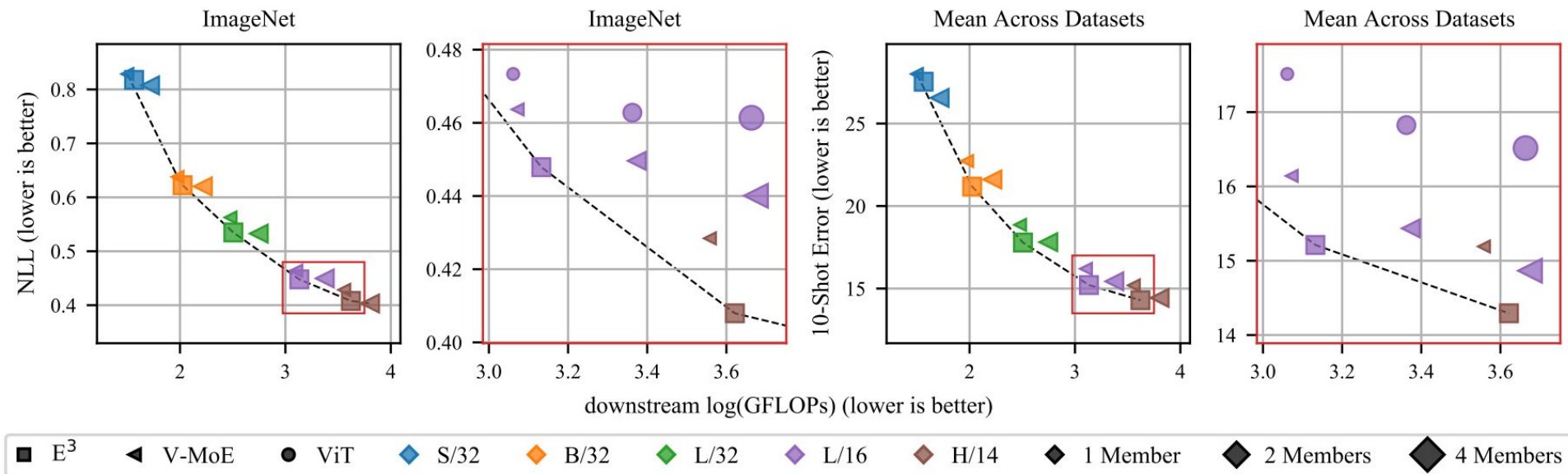
- E³ tends to be on the frontier performance vs. FLOPs

Evaluation of E³ (ImageNet & few-shot)



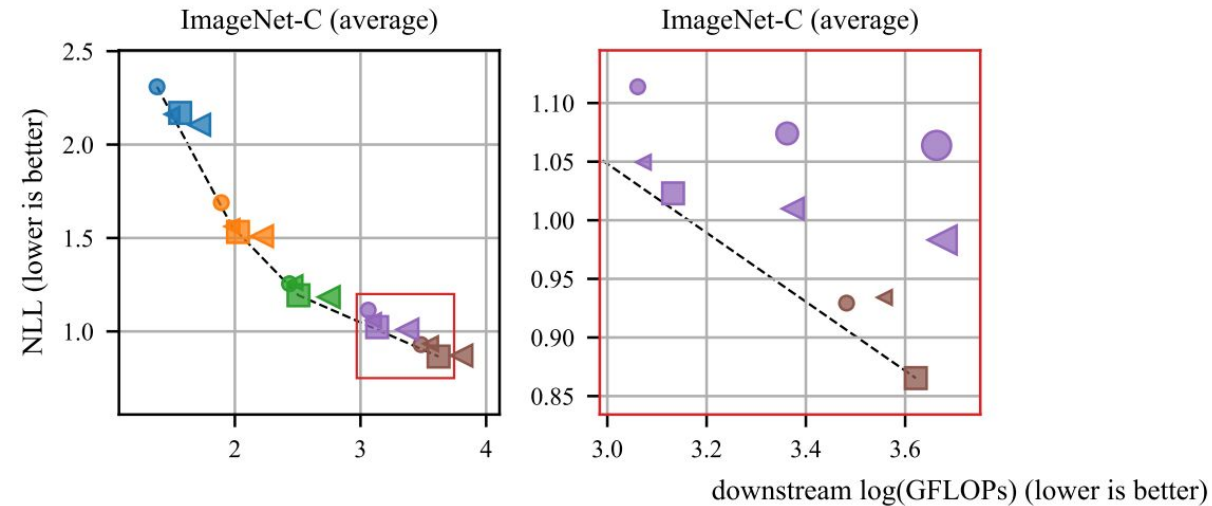
- E³ tends to be on the frontier performance vs. FLOPs

Evaluation of E³ (ImageNet & few-shot)



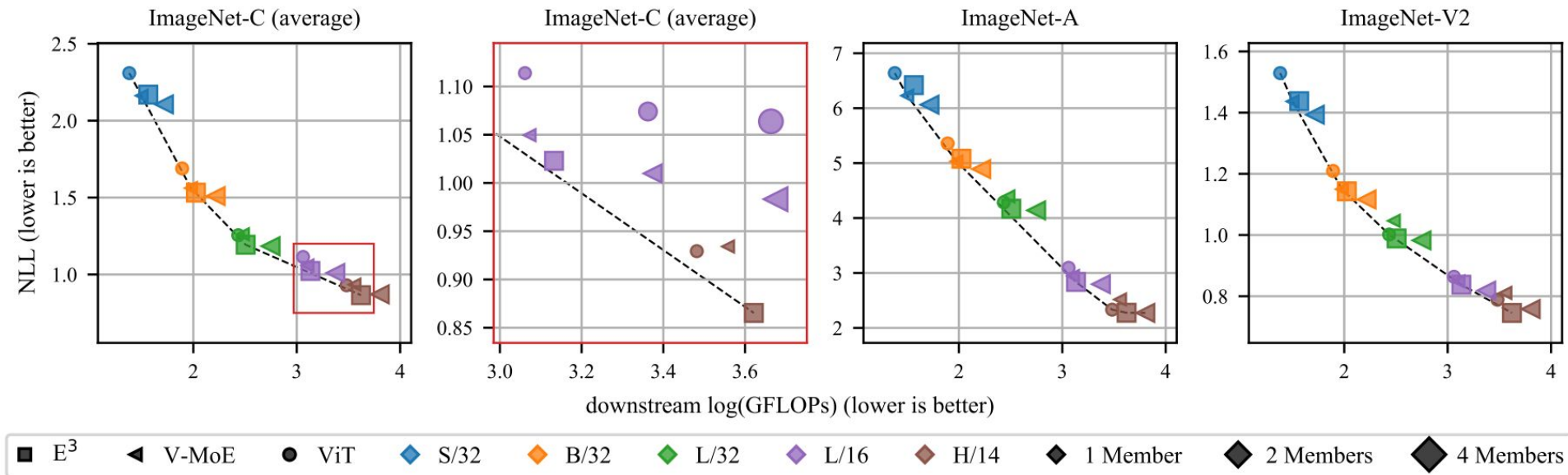
- E³ tends to be on the frontier performance vs. FLOPs
- E³ works well on benchmark where either ensembles or sparse MoEs are known to perform well

Evaluation of E³ (ImageNet under dataset shifts)



- E³ tends to be on the frontier performance vs. FLOPs

Evaluation of E³ (ImageNet under dataset shifts)



- E³ tends to be on the frontier performance vs. FLOPs

What about other efficient ensemble approaches?

Table : ImageNet performance (means \pm standard errors over 8 seeds) of different efficient ensemble approaches based on a ViT-B/32 architecture.

	K	M	NLL \downarrow	ERROR \downarrow	ECE \downarrow	GFLOPs \downarrow
	–	–	0.688 \pm 0.003	18.65 \pm 0.08	0.022 \pm 0.000	78.0
[Wen et al., 2019] BE ViT	–	2	0.682 \pm 0.003	18.47 \pm 0.05	0.021 \pm 0.000	97.1
	2	–	0.638 \pm 0.001	16.76 \pm 0.05	0.033 \pm 0.001	94.9
[Gal et al., 2015] MC Dropout V-MoE	1	2	0.648 \pm 0.002	17.10 \pm 0.05	0.019 \pm 0.001	97.2
	2	2	0.636 \pm 0.002	16.97 \pm 0.04	0.028 \pm 0.001	96.3
[Havasi et al., 2020] MIMO V-MoE	2	4	0.672 \pm 0.001	17.72 \pm 0.04	0.037 \pm 0.000	99.0
	1	2	0.622 \pm 0.001	16.70 \pm 0.03	0.018 \pm 0.000	105.9

What about other efficient ensemble approaches?

Table : ImageNet performance (means \pm standard errors over 8 seeds) of different efficient ensemble approaches based on a ViT-B/32 architecture.

	K	M	NLL \downarrow	ERROR \downarrow	ECE \downarrow	KL \uparrow	GFLOPs \downarrow
ViT	–	–	0.688 \pm 0.003	18.65 \pm 0.08	0.022 \pm 0.000	–	78.0
[Wen et al., 2019] BE ViT	–	2	0.682 \pm 0.003	18.47 \pm 0.05	0.021 \pm 0.000	0.040 \pm 0.001	97.1
V-MoE	2	–	0.638 \pm 0.001	16.76 \pm 0.05	0.033 \pm 0.001	–	94.9
[Gal et al., 2015] MC Dropout V-MoE	1	2	0.648 \pm 0.002	17.10 \pm 0.05	0.019 \pm 0.001	0.046 \pm 0.000	97.2
[Havasi et al., 2020] MIMO V-MoE	2	2	0.636 \pm 0.002	16.97 \pm 0.04	0.028 \pm 0.001	0.000 \pm 0.000	96.3
	2	4	0.672 \pm 0.001	17.72 \pm 0.04	0.037 \pm 0.000	0.001 \pm 0.000	99.0
E ³	1	2	0.622 \pm 0.001	16.70 \pm 0.03	0.018 \pm 0.000	0.217 \pm 0.003	105.9

- E³ tends to have diverse predictions

$$\frac{1}{M(M-1)} \sum_{m \neq m'}^M D_{\text{KL}}(p_m(y|x) || p_{m'}(y|x))$$

Some other results

In a nutshell:

- **Cifar10, Cifar10-C and Cifar100:** Even cleaner conclusions
- **OOD detection:**
 - In general, E^3 performs worse than V-MoE
 - Especially, “far” OOD detection task: Cifar* vs. {SVHN, Places365, DTD}
 - As the scale increases, E^3 performs better
- **Broader trend:** E^3 tends to perform better as the scale increases

Conclusions

- Sparse MoEs growing prevalence, e.g., in NLP [Patterson et al., 2021]
- Important, but also challenging, to study robustness at scale
- E³ often on the “performance vs. FLOPs” frontier
- E³ **simple & convenient**
 - Little code change
 - Can finetune standard checkpoints
- TMLR paper: <https://openreview.net/pdf?id=i0ZM36d2qU>
- Code: <https://github.com/google-research/vmoe>

Supplementary material

Detailed view of E^3

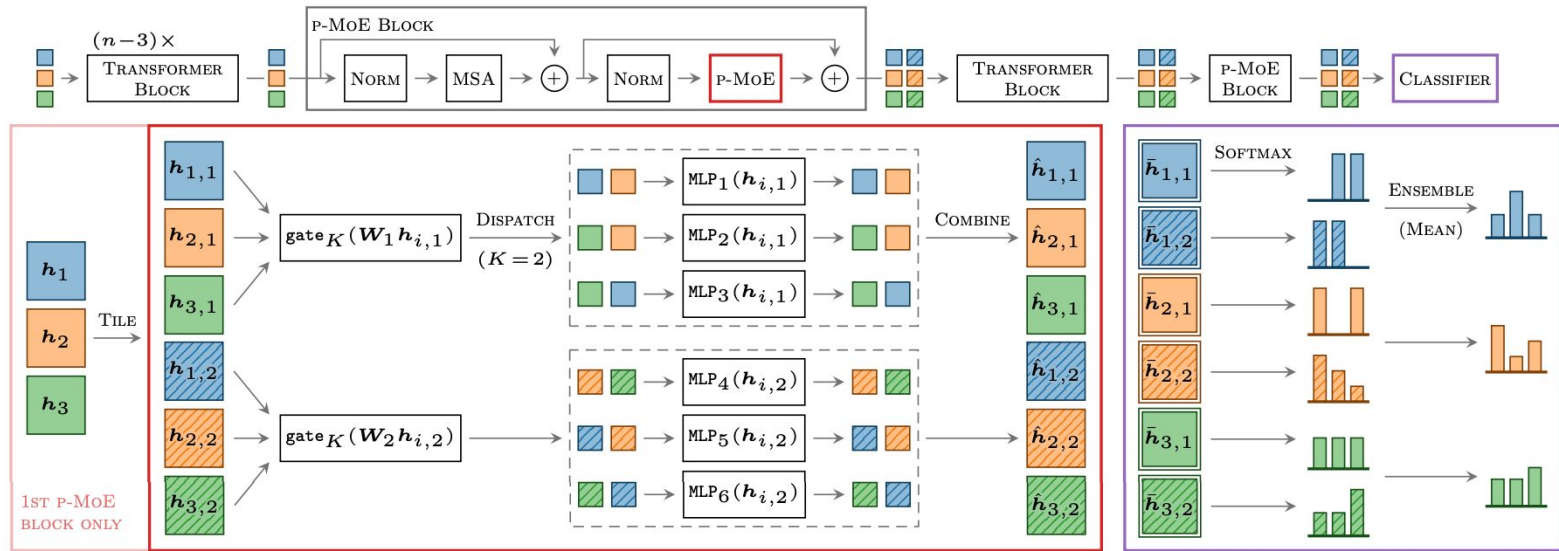
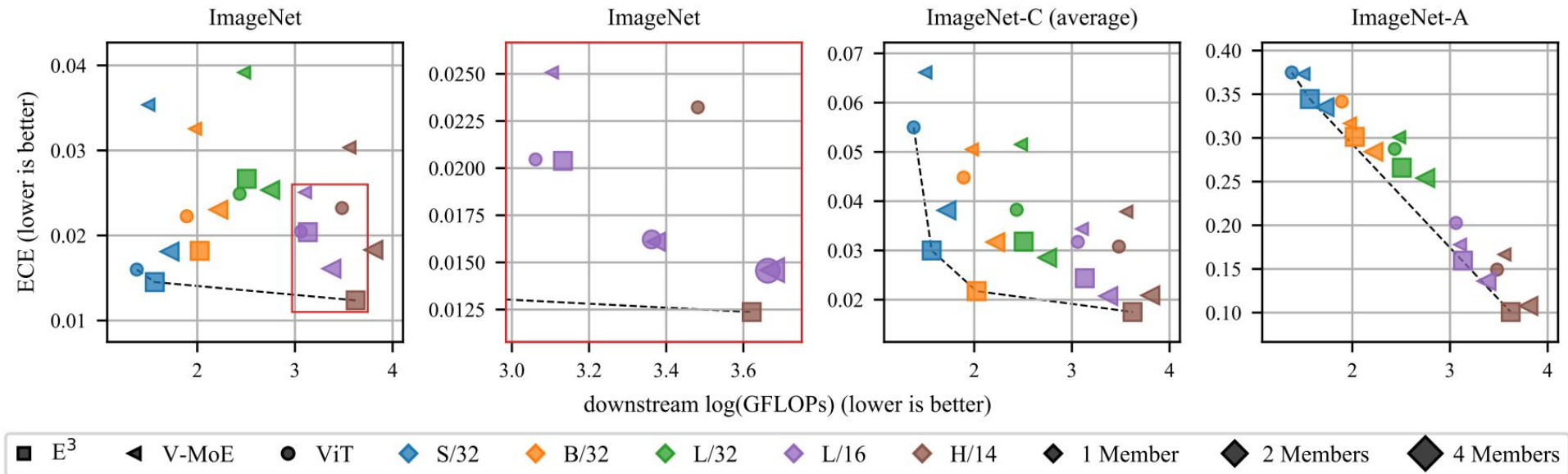


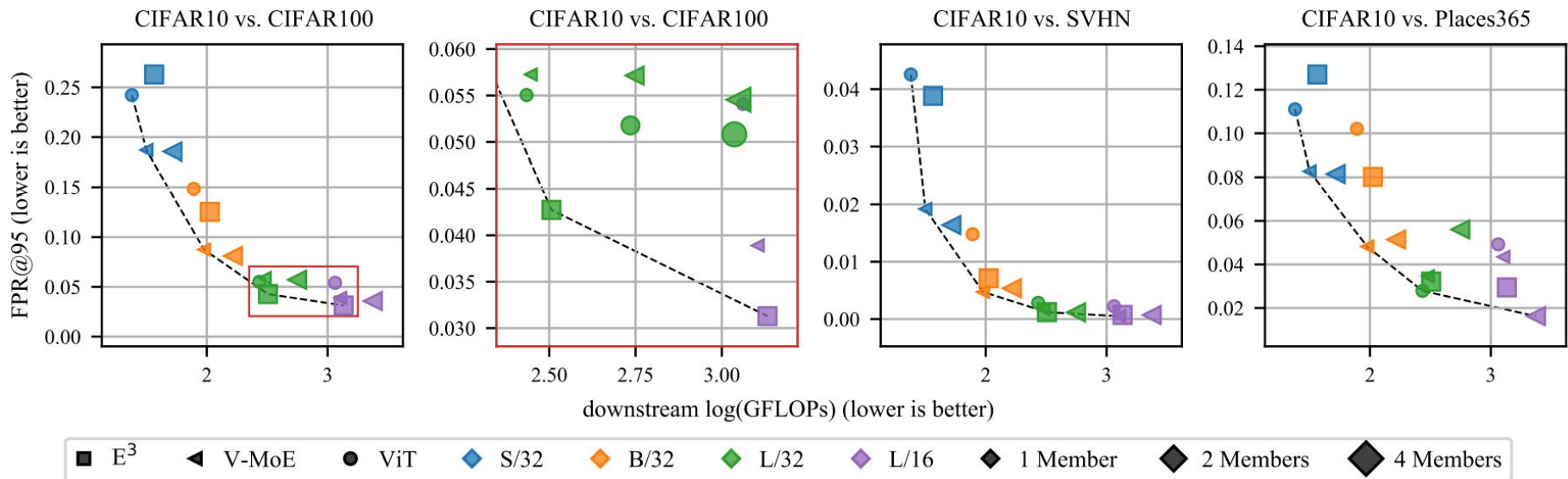
Figure 1: End-to-end overview of E^3 with $E=6$ experts, partitioned into $M=2$ groups, with sparsity of $K=2$, and a “last-2” configuration. **Top:** E^3 contains a sequence of transformer blocks, followed by alternating transformer and p(artitioned)-MoE blocks. As in ViT, images are split into patches whose embeddings are processed by each block. Here, we show 1 embedding for each of three images (■, ■, ■). **Bottom left:** In a p-MoE block, we replace the transformer block’s MLP with parallel partitioned expert MLPs, see (2). The effect of the routing weights is not depicted. Embeddings are tiled (▨) in the first p-MoE block only. **Bottom right:** The classifier averages predictions from the final tiled representations (▣).

Evaluation of E³: ECE



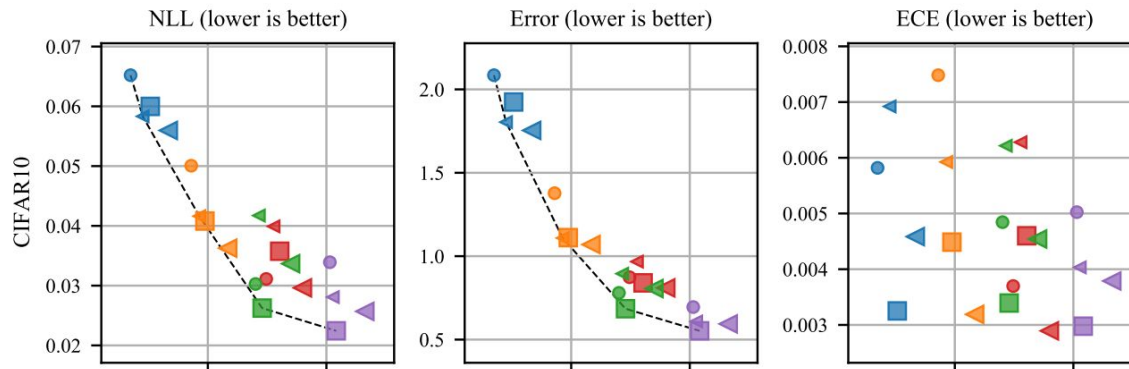
- ECE is not consistent for different ViT/V-MoE families
- E³ improves ECE over ViT and V-MoE (& V-MoE provides poor ECE)

Evaluation of E³: OOD detection



- E³ does not provide consistent OOD detection performance
- Improvement at larger scales

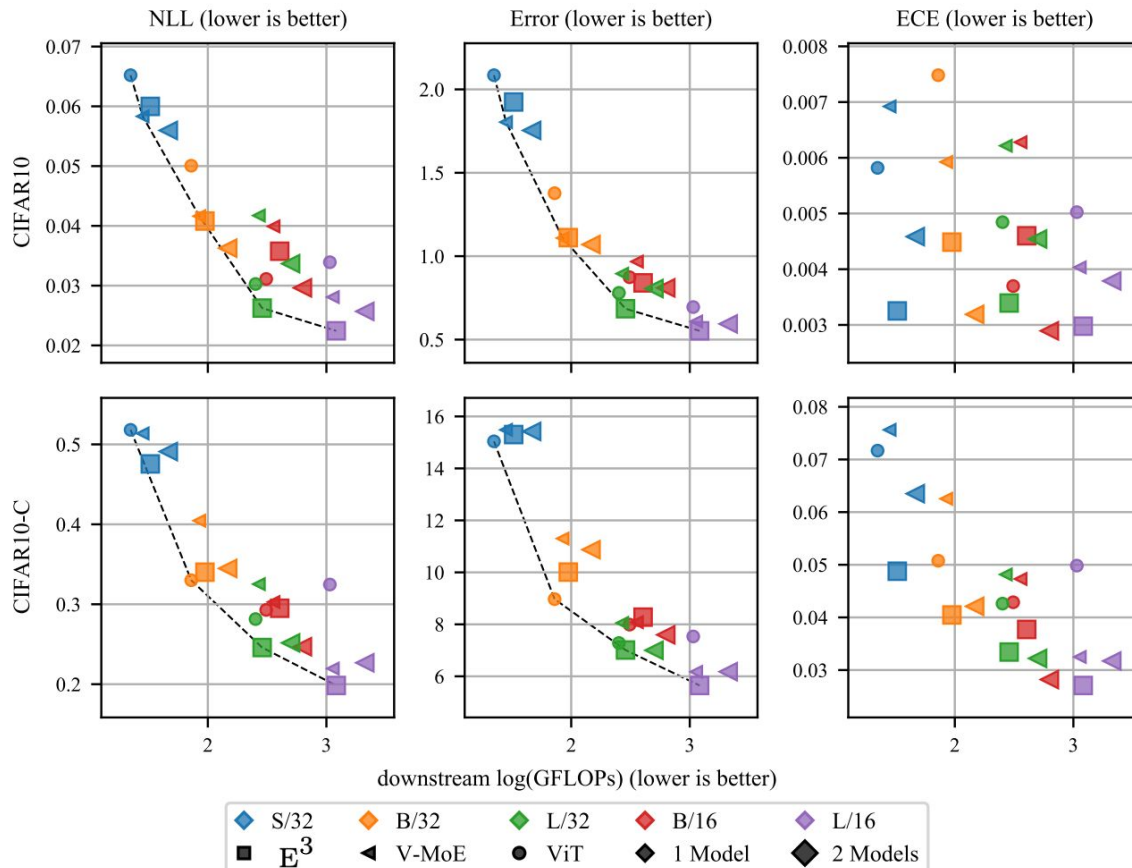
Evaluation of E³: Cifar10 & Cifar10-C



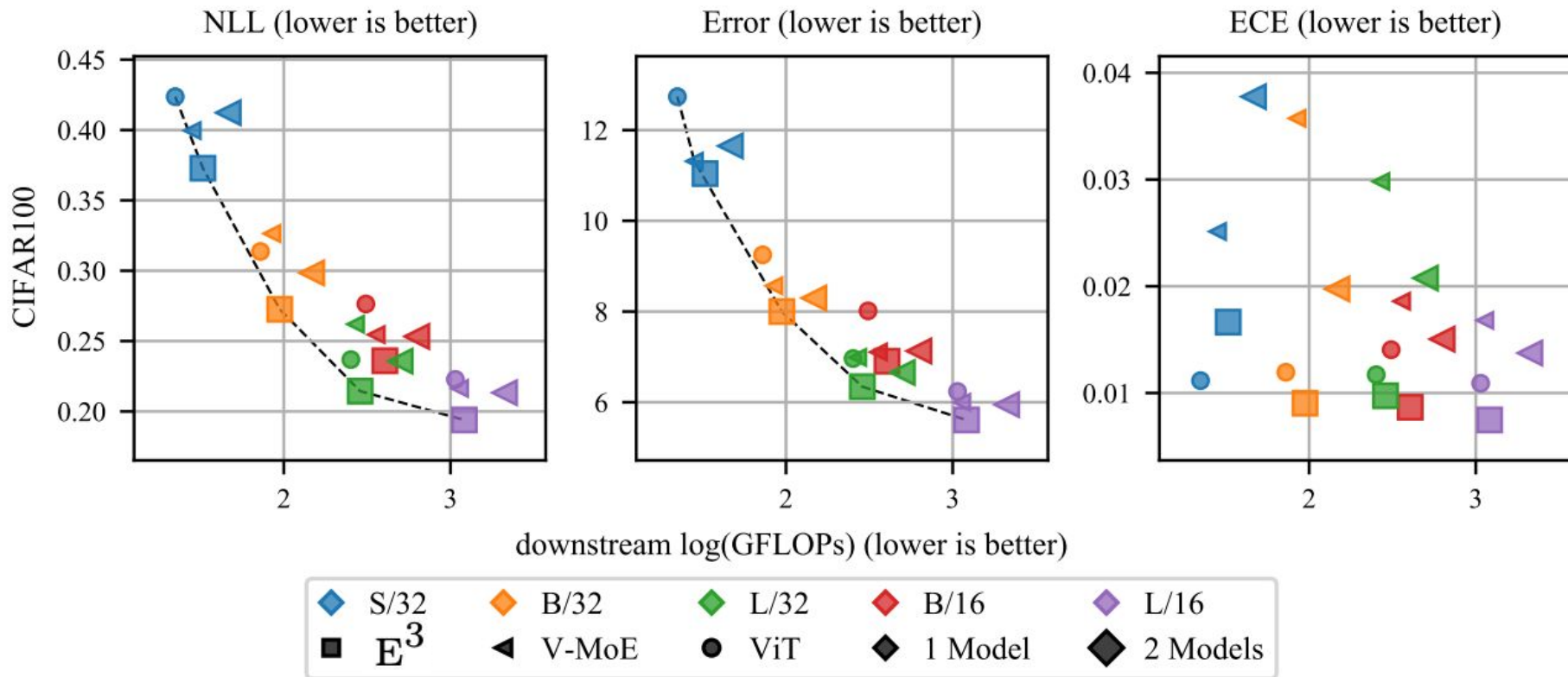
downstream log(GFLOPs) (lower is better)



Evaluation of E³: Cifar10 & Cifar10-C



Evaluation of E³: Cifar100

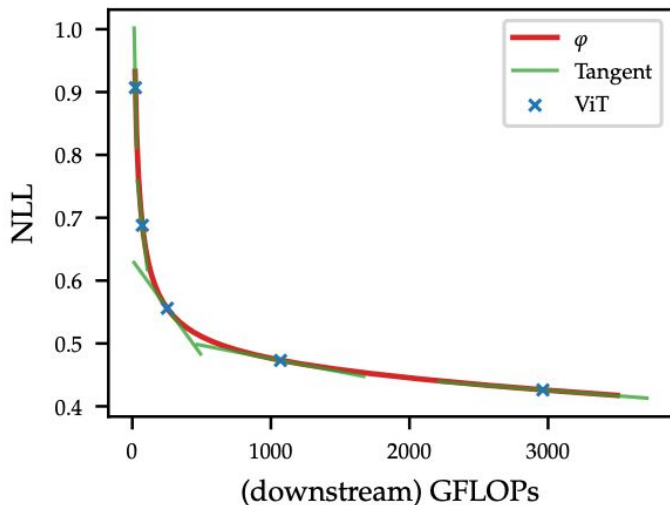


Ablation E³: Tiling and partitioning

Table 2: ImageNet performance (means \pm standard errors over 8 seeds) of E³-B/32 ($K = M = 2$), V-MoE ($K = 4$), and two ablations: *only* tiling and *only* partitioning. The noise in gate_K is denoted by σ .

	NLL \downarrow	ERROR \downarrow	ECE \downarrow	KL \uparrow
V-MoE	0.636 \pm 0.001	16.70 \pm 0.04	0.034 \pm 0.001	—
E ³	0.612 \pm 0.001	16.49 \pm 0.02	0.013 \pm 0.000	0.198 \pm 0.003
Tiling	0.637 \pm 0.002	16.74 \pm 0.06	0.028 \pm 0.001	0.000 \pm 0.000
Tiling ($\sigma \times 2$)	0.638 \pm 0.001	16.72 \pm 0.03	0.033 \pm 0.001	0.001 \pm 0.000
Tiling ($\sigma \times 4$)	0.638 \pm 0.001	16.74 \pm 0.03	0.033 \pm 0.001	0.002 \pm 0.000
Partitioning	0.640 \pm 0.001	16.72 \pm 0.05	0.034 \pm 0.001	—

E³ performs better as the scale increases



$$\text{Normalised improvement}(v) = \text{improvement}(v) \times \frac{\varphi'(\text{FLOPs}_{H/14})}{\varphi'(\text{FLOPs}_v)}$$

		S/32	B/32	L/32	L/16	H/14
Normalised	E ³ vs. ViT	0.02%	0.09%	0.24%	2.35%	4.27%
	V-MoE vs. ViT	0.01%	0.06%	-0.04%	0.89%	0.02%
Not normalised	E ³ vs. ViT	9.82%	9.53%	3.76%	5.38%	4.27%
	V-MoE vs. ViT	7.98%	6.62%	-0.60%	2.05%	0.02%