

TIGHT CLUSTERS MAKE SPECIALIZED EXPERTS

Anonymous authors

Paper under double-blind review

ABSTRACT

At the core of Sparse Mixture-of-Experts (MoE) models is the router that learns the clustering structure of the input distribution in order to direct tokens to suitable experts. However these latent clusters may be unidentifiable, causing slow convergence, vulnerability to contamination, and degraded representations. We examine the router through the lens of clustering optimization, deriving optimal feature weights that maximally distinguish these clusters. Using these weights, we compute token-expert assignments in an adaptively transformed space that better separates clusters, helping identify the best-matched expert for each token. In particular, for each expert cluster, we compute weights that scale features according to whether that expert clusters tightly along that feature. We term this novel router the Adaptive Clustering (AC) router. Our AC router confers three connected benefits: 1) faster convergence, 2) better robustness, and 3) overall performance improvement, as experts are specialized in semantically distinct regions of the input space. We empirically demonstrate the advantages of our AC router in language modeling and image classification in both clean and corrupted settings.

1 INTRODUCTION

Scaling up model capacity continues to yield performance gains across tasks, notably in visual representation learning and language modeling (Alexey, 2020; Bao et al., 2021; Raffel et al., 2020). However, larger models incur increasing computational cost, prompting research in Sparse Mixture-of-Experts (MoE) models (Shazeer et al., 2017; Fedus et al., 2022; Lepikhin et al., 2020), which activate only sub-modules, or *experts*, to reduce overhead. These models can outperform dense architectures with nearly constant computation on speech recognition, image recognition, machine translation, and language modeling (Riquelme et al., 2021; Kumatani et al., 2021).

At the core of the MoE layer is the learned router which segments the input space such that semantically similar input tokens are assigned to corresponding experts. Recent work explores various routing strategies, from linear programs (Lewis et al., 2021) and cosine similarity-based methods (Chi et al., 2022) to soft assignments (Puigcerver et al., 2023) and top-k routing (Shazeer et al., 2017; Zhou et al., 2022). These methods rely on dot-products between inputs and experts, which can be suboptimal when semantic regions are not easily identified in high-dimensional space. Typically, we expect that the true underlying clusters present in the data will cluster on different, potentially disjoint, subsets of features, and may not be discoverable when using the full feature set. This phenomenon can lead to slower convergence as experts are unable to specialize on semantically similar regions of the data, poor robustness as data contamination can spuriously assign inputs to unsuitable experts, and degraded overall downstream performance due to suboptimal input-expert matching.

Contribution. We introduce the Adaptive Clustering (AC) router and Adaptive Clustering Mixture-of-Experts (ACMoE), a novel MoE method in which the router computes token-expert assignments in a transformed space that maximally identifies latent clusters in the data and more easily discovers the best-matched expert for each token. This produces three benefits: 1) *faster convergence* as experts are able to specialize more quickly by being allocated semantically similar inputs, 2) *better robustness* as latent clusters are better separated, thereby minimizing the risk that data corruption erroneously assigns tokens to unsuitable experts, and 3) *better downstream performance* due to improved expert specialization. In order to discover the corresponding weights, we present a feature-weighted clustering optimization perspective on the MoE framework and demonstrate how the clustering solution obtains the required feature weights. We theoretically prove that our proposed routing mechanism learns the latent clustering structure of the data faster than standard routers and that our

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

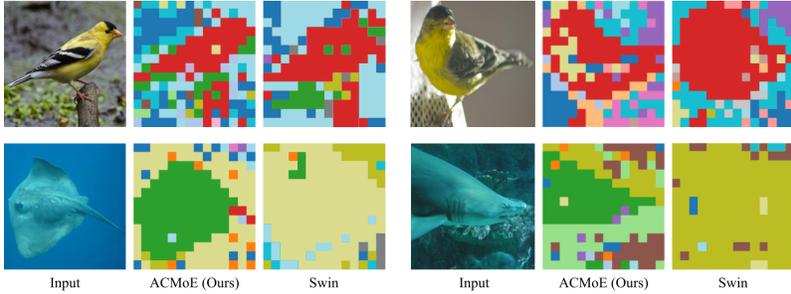


Figure 1: ACMoE discovers semantically distinct regions. We show 14x14 image reconstructions where patches are colored by assigned experts. **Top row:** Swin fails to segment the bird precisely while ACMoE accurately discovers the bird and relevant foreground. **Bottom row:** When the background and foreground are hard to distinguish, Swin fails to register the stingray (left) or shark (right) and allocates one expert for virtually the entire image. ACMoE, however, accurately discovers the semantically distinct regions and utilizes one expert (green) to specialize on the stingray and shark and different experts to specialize on the the background.

mechanism is more robust to data contamination. Furthermore, our proposed method involves no learnable parameters and can be computed highly efficiently. In summary, our contributions are:

1. We develop the novel Adaptive Clustering router for MoE architectures, which computes token-expert assignments in a transformed space that promotes separation of latent clusters in the data and more easily identifies the best-matched expert for each token.
2. We propose a feature-weighted clustering optimization perspective on token-expert assignment and derive the optimal feature weights for routing.
3. We provide a theoretical framework demonstrating how MoE robustness and convergence depend on the clustering structure of the input space.

We empirically demonstrate that 1) the AC router outperforms baseline routers in language modeling and downstream finetuning, and image classification in clean and contaminated settings, 2) the AC router exhibits faster convergence than baseline methods, and 3) the AC router attains these performance improvements with no learnable parameters and negligible computational overhead.

Preliminaries. We consider Transformer (Vaswani, 2017) based MoE architectures and follow the approach of previous work where the MoE layer is inserted after the self-attention layer, replacing the traditional feed-forward network (Fedus et al., 2022; Du et al., 2022; Liu et al., 2021). Let $\mathbf{h} \in \mathbb{R}^d$ be a hidden representation and $e_1, e_2, \dots, e_N \in \mathbb{R}^d$ be the N learnable expert embeddings for model hidden dimension d . The MoE layer selecting the top k experts is described by:

$$\mathcal{K} := \text{topk}_k(s_k) = \text{topk}_k(\mathbf{h}^\top \mathbf{e}_k) \tag{1}$$

$$f^{SMoE}(\mathbf{h}) = \mathbf{h} + \sum_{k \in \mathcal{K}} g(\mathbf{h}^\top \mathbf{e}_k) f_k^{\text{FFN}}(\mathbf{h}), \tag{2}$$

where f_k^{FFN} is the k^{th} expert feed-forward network, $s_k = \mathbf{h}^\top \mathbf{e}_k$ is the similarity score between token \mathbf{h} and the k^{th} expert \mathbf{e}_k and $g(\cdot)$ is a gating function often chosen as softmax. We refer to Eqn. 1 as the router, which learns the best matched experts per token, and Eqn. 2 as the overall MoE layer.

2 A CLUSTERING OPTIMIZATION PERSPECTIVE

We analyze the MoE router through the lens of feature-weighted clustering (Witten & Tibshirani, 2010; Friedman & Meulman, 2004; Brusco & Cradit, 2001). We explicitly model the router’s task as learning a token assignment that groups together similar tokens. We incorporate learnable feature weights in solving a clustering optimization problem to optimally reveal latent clusters and present an analytical solution for any given routing assignment. We discuss how this solution improves the MoE router before providing the full formulation of our AC router in the next section.

2.1 CLUSTERING OPTIMIZATION

Let $\mathbf{h}_i = [h_{i1}, \dots, h_{id}]^\top$ be the i^{th} hidden representation and $D_{ij}(\mathbf{w}) = \sum_{q \in [d]} w_q \rho_{ijq}$ be the distance between pairs of vectors \mathbf{h}_i and \mathbf{h}_j where $\mathbf{w} = [w_1, \dots, w_d]$ are nonnegative feature weights

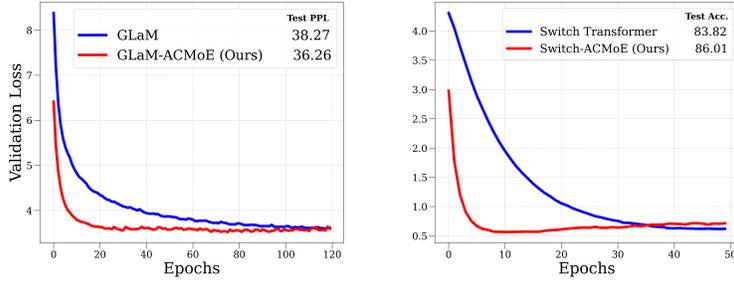


Figure 2: Fast Convergence of ACMoE. **Left:** Convergence speed on WikiText-103 pretraining using GLaM (Du et al., 2022) backbone. **Right:** Convergence speed on Banking-77 finetuning using Switch Transformer (Fedus et al., 2022) backbone. We see faster convergence and better final perplexity (PPL) and accuracy (Acc.).

summing to 1 and ρ_{ijq} is a chosen distance metric over the q^{th} feature. We wish to learn a classifier $r(i) = k$ assigning the i^{th} object to a group k over the input set of N objects, where objects within the same group are more similar to each other than to those in other groups. We furthermore wish to model the scenario that groupings exist in different latent subspaces with varying dependence on possibly disjoint subsets of features. We therefore use clustering criterion with cluster-dependent feature weights $\{\mathbf{w}_k\}_{k=1}^E$ for E groups, given by:

$$\begin{aligned} (r^*, \{\mathbf{w}_k^*\}_{k=1}^E) = \arg \min_{r, \{\mathbf{w}_k\}} & \sum_{k \in [E]} \frac{1}{N_k^2} \sum_{r(i)=k} \sum_{r(j)=k} D_{ij}^J(\mathbf{w}_k), \\ \text{such that} & \sum_{q \in [d]} w_{qk} = 1, \quad \forall k \in [E], \end{aligned} \quad (3)$$

where $D_{ij}^J(\mathbf{w}_k) = \sum_{l=1}^d w_{lk} \rho_{ijl} + \lambda J(\mathbf{w}_k)$ denotes the weighted distance between i and j combined with regularization J and regularization strength λ . We set the regularizer to the Kullback-Leibler divergence between w and the uniform distribution $\mathbf{u} = (1/d, \dots, 1/d) \in \mathbb{R}^d$, denoted by $J(\mathbf{w}_k) = D_{\text{KL}}(\mathbf{u} \parallel \mathbf{w}_k)$, where λ reflects our preference to maintain more or less features in the solution.

2.2 MOE AS CLUSTERING OPTIMIZATION

In the MoE setting, the router performs the role of the classifier $r : \mathbb{R}^d \rightarrow [E]$, which is learned via gradient descent on the output loss. Therefore, we fix r and optimize the criterion with respect to cluster-wise feature weights \mathbf{w}_k . Under this interpretation, the router learns via backpropagation to optimally allocate representations to experts, with representations adaptively transformed to maximally reveal the clustering structure of the input data. The objective in Eqn. 3 then becomes

$$\{\mathbf{w}_k^*\}_{k=1}^E = \arg \min_{\{\mathbf{w}_k\}} \sum_{k \in [E]} \frac{1}{N_k^2} \sum_{r(i)=k} \sum_{r(j)=k} D_{ij}^J(\mathbf{w}_k), \quad (4)$$

with the same summation to unity constraints as in Eqn. 3. The following theorem presents the optimal weights per feature q and cluster k :

Theorem 1 (Optimal feature weights). *Let $s_{qk} := N_k^{-2} \sum_{r(i)=k} \sum_{r(j)=k} \rho_{ijq}$ be a measure of dispersion on the q^{th} feature for the representations assigned to cluster k . Then, for a given router function $r : \mathbb{R}^d \rightarrow [E]$, the corresponding optimal weights $\{\mathbf{w}_k\}_{k \in [E]}$ that minimize the feature-weighted clustering optimization problem in Eqn. 4 are given by*

$$w_{qk} = \frac{\lambda/d}{s_{qk} + \alpha_k}, \quad (5)$$

for $(q, k) \in [d] \times [E]$, where $\{\alpha_k\}_{k \in [E]}$ are constants that for any $\lambda > 0$ satisfy $\sum_{q \in [d]} \frac{1}{s_{qk} + \alpha_k} = \frac{d}{\lambda}$. The existence α_k and the proof of Theorem 1 is provided in Appendix A.1. The optimal weights for a cluster k given in Eqn. 5 take an intuitive form in that they are inversely proportional to the measure of dispersion in cluster k along each dimension, $\mathbf{w}_k \propto [\frac{1}{s_{1k}}, \dots, \frac{1}{s_{dk}}]$. Hence, the optimal cluster-wise feature weights scale features according to their contribution to forming tight clusters.

This method enables the MoE router to perform better token-expert matching. The cluster-wise feature weights w_k scale each token according to the specialization of the experts, as large weights indicate those features are highly important to the identification of that expert cluster, thereby allowing the router to best identify the most suitable expert for each token. Note that this solution is local in that we learn the optimal weights adaptively *per cluster*, obtaining w_k for all $k \in [E]$, and so we compute a unique scaling of the feature space adaptively *per cluster* as well.

3 A TIGHT CLUSTER IS A SPECIALIZED EXPERT

In this section, we use the solution from the optimization problem in Eqn. 5 to obtain the AC router and present theoretical results on how AC routing promotes faster convergence and robustness.

3.1 FULL TECHNICAL FORMULATION

We first integrate the weights from Eqn. 5 into the transformation in Definition 1, which scales each dimension according to the k^{th} expert’s specialization:

Definition 1 (AC Router Transformation M_k). *Let $\mathcal{C}_k^\ell = \{\mathbf{h}_1^\ell, \dots, \mathbf{h}_{N_k}^\ell\}$ be the tokens assigned to expert k at layer ℓ . Let $s_{qk}^\ell \in \mathbb{R}$ be a measure of a spread in the q^{th} dimension for cluster k , such as mean absolute deviation $s_{qk}^\ell = \frac{1}{N_k} \sum_{i \in \mathcal{C}_k^\ell} |\mathbf{h}_{iq}^\ell - \bar{\mathbf{h}}_q^\ell|$. Then, the cluster-dependent router transformation for expert k at layer ℓ is given by a diagonal matrix $M_k^\ell := \text{diag}(1/s_{1k}^\ell, \dots, 1/s_{dk}^\ell)$.*

Using M_k^ℓ to adaptively scale the feature space according to the experts’ specialization yields our AC router and corresponding ACMoE layer:

Definition 2 (Adaptive Clustering Router and MoE Layer). *Let $\mathbf{h}^\ell \in \mathbb{R}^d$ be the hidden representation of an input, $\mathbf{e}_1^\ell, \dots, \mathbf{e}_N^\ell \in \mathbb{R}^d$ be expert embeddings at layer ℓ . Let $\mathbf{h}^{\ell-1} \in \mathbb{C}_{k^*}^{\ell-1}$ have been assigned to expert k^* in the previous layer. Let $M_{k^*}^{\ell-1} \in \mathbb{R}^{d \times d}$ be the Adaptive Clustering transformation (Definition 1) for input \mathbf{h} at layer $\ell - 1$. Let $g(\cdot)$ be the softmax function. Then the following equations describe the Adaptive Clustering router (Eqn. 6) and overall ACMoE layer (Eqn. 7):*

$$\mathcal{K} := \text{topk}_k(s_k) = \text{topk}_k(\mathbf{h}^{\ell\top} M_{k^*}^{\ell-1} \mathbf{e}_k^\ell) \quad (6)$$

$$f^{\text{ACMoE}}(\mathbf{h}^\ell) = \mathbf{h}^\ell + \sum_{k \in \mathcal{K}} g(\mathbf{h}^{\ell\top} M_{k^*}^{\ell-1} \mathbf{e}_k^\ell) f_k^{\text{FFN}, \ell}(\mathbf{h}^\ell). \quad (7)$$

Remark 1. *We see from Eqns. 6 and 7 that standard routers and MoE layer are recovered by setting the adaptive clustering router transformation to the identity matrix, $M_k = \mathbf{I}_d$ for all $k \in [E]$. Within our framework, standard routers assume all experts $k \in [E]$ depend equally on all dimensions.*

3.2 ADAPTIVE CLUSTERING PROMOTES ROBUSTNESS AND FAST CONVERGENCE

We now present theoretical propositions on how our method improves robustness and convergence speed. Robustness follows from the exponentially lower probability of erroneous expert assignment and faster convergence follows from improved Hessian conditioning with respect to expert embeddings. Proofs are deferred to Appendix A.3.

Robustness. Lemma 1 shows that our AC transformation (Def. 1) increases inter-cluster separation, and Lemma 2 provides a probability bound for incorrect assignments as a function of inter-cluster distance. Robustness of AC routing then follows as a direct combination of these two lemmas.

Lemma 1 (Adaptive Clustering Router Transformation Increases Cluster Separation). *Let the data be generated from a Gaussian mixture model with components, $g_c = \mathcal{N}(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$ for $c \in [E]$. Without loss of generality, consider two expert clusters $c \in \{a, b\}$ where a token representation $\mathbf{h} \sim g_a$ belongs to cluster a . Let $M_a = \text{diag}(1/s_{1a}, \dots, 1/s_{da})$ be the router transformation constructed from the feature-wise dispersions, s_{qa} , of cluster g_a for each feature $q \in [d]$ as given by Definition 1. Then the distance between cluster means in the M_a -transformed space, defined as $\|\boldsymbol{\mu}_k - \boldsymbol{\mu}_a\|_{M_a}^2 := (\boldsymbol{\mu}_k - \boldsymbol{\mu}_a)^\top M_a (\boldsymbol{\mu}_k - \boldsymbol{\mu}_a)$, is larger than in the original Euclidean space: $\|\boldsymbol{\mu}_k - \boldsymbol{\mu}_a\|_{M_a}^2 \geq \|\boldsymbol{\mu}_k - \boldsymbol{\mu}_a\|^2$.*

In Lemma 2, we derive the probability of mis-assignment as a function of inter-cluster distance, highlighting how cluster separation mitigates the effect of noise that can confuse the router.

Lemma 2 (Incorrect Assignment Probability). Let $\mathbf{h} \sim \mathcal{N}_{k^*}(\boldsymbol{\mu}_{k^*}, \boldsymbol{\Sigma}_{k^*})$ be a representation belonging to cluster k^* . Let $\mathbf{h}' = \mathbf{h} + \boldsymbol{\epsilon}$ be contaminated by some 0-mean noise $\boldsymbol{\epsilon} \sim (\mathbf{0}, \boldsymbol{\Sigma}_{\epsilon})$. Let k be the nearest, incorrect cluster to k^* . Let the inter-cluster mean distance between k^* and k be given by $\|\delta\boldsymbol{\mu}\| := \|\boldsymbol{\mu}_{k^*} - \boldsymbol{\mu}_k\|$. Let the routing assignment be given by $r : \mathbb{R}^d \rightarrow [E]$ and denote the cumulative density of a standard normal distribution by Φ . Then the probability of incorrect assignment is

$$\Pr(r(\mathbf{h}') \neq k^*) = 1 - \Phi\left(\frac{\|\delta\boldsymbol{\mu}\|^2}{2\sqrt{\delta\boldsymbol{\mu}^\top(\boldsymbol{\Sigma}_{k^*} + \boldsymbol{\Sigma}_{\epsilon})\delta\boldsymbol{\mu}}}\right). \quad (8)$$

It is worth noting that since $1 - \Phi(x) \sim (\sqrt{2\pi}x)^{-1}e^{-x^2/2}$ for large x and $\sqrt{\delta\boldsymbol{\mu}^\top(\boldsymbol{\Sigma}_{k^*} + \boldsymbol{\Sigma}_{\epsilon})\delta\boldsymbol{\mu}} = O(\|\boldsymbol{\mu}\|)$, we find that the probability of incorrect cluster assignment as given by Eqn. 8, $\Pr(r(\mathbf{h}') \neq k^*) = e^{-O(\|\delta\boldsymbol{\mu}\|^2)}$ is an exponentially decreasing function in $\|\delta\boldsymbol{\mu}\|$. We now combine the notions in Lemmas 1 and 2 to obtain that the probability of erroneous assignment using the AC router is exponentially smaller than under a standard routing scheme:

Proposition 1 (Robustness of ACMoE). Consider an expert assignment setting for the representation $\mathbf{h} \sim \mathcal{N}_{k^*}(\boldsymbol{\mu}_{k^*}, \boldsymbol{\Sigma}_{k^*})$ as in Lemma 2 with two routers given by $r : \mathbb{R}^d \rightarrow [E]$ and $r^{\text{AC}} : \mathbb{R}^d \rightarrow [E]$ for standard (Eqn. 2) and AC routers (Definition 2), respectively. Then the probabilities of incorrect assignments of routers r and r^{AC} satisfy $\Pr(r^{\text{AC}}(\mathbf{h}') \neq k^*) \leq \Pr(r(\mathbf{h}') \neq k^*)$.

Convergence. Our AC router reduces the conditioning number of the Hessian of the loss with respect to the expert e_k , improving the loss landscape and enabling faster convergence of the router. We find this result empirically supported, as shown in Fig. 2. Formally this is:

Proposition 2 (Faster convergence of ACMoE). Let $\mathcal{L}^{\text{MoE}} : \Theta \rightarrow \mathbb{R}_+$ and $\mathcal{L}^{\text{ACMoE}} : \Theta \rightarrow \mathbb{R}_+$ be the network loss functions over parameters Θ when employing the standard (Eqn. 2) and AC routers (Definition 2), respectively. Let $\kappa(\mathbf{A}) = \lambda_{\max}/\lambda_{\min}$ denote the conditioning number of a matrix \mathbf{A} with largest and smallest eigenvalues λ_{\max} and λ_{\min} respectively. Let the Hessian of an i^{th} expert be given by $\nabla_{e_i}^2$. Then for each $i \in [E]$ the following holds with high probability

$$\kappa(\nabla_{e_i}^2 \mathcal{L}^{\text{ACMoE}}) \leq \kappa(\nabla_{e_i}^2 \mathcal{L}^{\text{MoE}}) \quad (9)$$

4 EXPERIMENTAL RESULTS

We evaluate our method on Wikitext-103 (Merity et al., 2016) language modeling and ImageNet (Deng et al., 2009) image classification. We integrate our AC router into Switch Transformer (Fedus et al., 2022), Generalist Language Model (GLaM) (Du et al., 2022), and Swin Transformer (Liu et al., 2021) backbones. We show i) ACMoE obtains substantive improvements over baseline models across both language and vision tasks; ii) ACMoE offers robust improvements on contaminated samples. We additionally show in Appendix B that ACMoE attains these gains with negligible additional computational overhead. Results are averaged over 5 runs with different seeds.

4.1 LANGUAGE MODELING

Setup. Following Pham et al. (2024), we compare ACMoE against Switch Transformer and GLaM using 16 experts and top-2 routing with 220M parameters. We report pretraining test perplexity (PPL) for Wikitext-103 and top-1 accuracy for finetuning classification tasks on Stanford Sentiment Treebank-2 (SST2) (Socher et al., 2013), Stanford Sentiment Treebank-5 (SST5) (Socher et al., 2013), and Banking-77 (B77) (Casanueva et al., 2020). Full details are provided in Appendix C.

Language Modeling Results. Table 2 show test PPL on clean WikiText-103 and when contaminated by Text Attack, where words are randomly swapped with a token 'AAA'. We follow the setup of Han et al. (2024) and assess models by training them on clean data before attacking the test data. ACMoE outperforms baseline Switch and GLaM in both clean and contaminated settings with gains of up to 5.8%. Table 1 further shows ACMoE pretrained models surpass the performance of baselines in finetuning, with strong, consistent improvements of approximately 3%.

4.2 IMAGE CLASSIFICATION

Setup. Following Liu et al. (2021), we evaluate ACMoE against a 280M parameter Swin Transformer with 16 experts. We evaluate robustness under white box attacks fast gradient sign method

Table 1: WikiText-103 test PPL and top-1 test accuracy on SST2, SST5, and B77 finetuning.

Model	Test PPL (\downarrow)	SST2 (\uparrow)	SST5 (\uparrow)	B77 (\uparrow)
<i>Switch Transformer</i> (Fedus et al., 2022)	35.48	76.27	39.13	83.82
Switch-ACMoE (Ours)	34.42	77.32	40.04	86.01
<i>GLaM</i> (Du et al., 2022)	38.27	69.97	33.69	80.89
GLaM-ACMoE (Ours)	36.26	71.90	34.24	82.33

Table 2: Clean and Contaminated Test Perplexity (PPL) on WikiText-103

Model	Clean Test PPL (\downarrow)	Contaminated Test PPL (\downarrow)
<i>Switch Transformer</i> (Fedus et al., 2022)	35.48	48.12
Switch-ACMoE (Ours)	34.42	47.61
<i>GLaM</i> (Du et al., 2022)	38.27	50.84
GLaM-ACMoE (Ours)	36.26	47.91

Table 3: Test Accuracy on ImageNet corrupted PGD, FGSM, and SPSA.

Model	Clean Data		PGD		FGSM		SPSA	
	Top 1	Top 5						
<i>Swin</i> (Liu et al., 2021)	76.10	92.99	40.85	75.51	54.70	85.22	60.57	82.75
Swin-ACMoE (Ours)	76.31	93.14	43.74	78.55	55.78	85.80	63.47	86.05

(FGSM) (Goodfellow et al., 2014) and projected gradient descent (PGD) (Madry et al., 2017), and black box simultaneous perturbation stochastic approximation (SPSA) (Uesato et al., 2018).

Image Classification Results.. Table 3 shows performance on ImageNet against FGSM, PGD, and SPSA. Compared against Swin Transformer, ACMoE improves a noteworthy 7% against PGD.

5 RELATED WORK

Routing Methods. Recent studies propose routers based on reinforcement learning (Bengio et al., 2015), linear programs (Lewis et al., 2021; Nguyen et al., 2024), cosine similarity (Chi et al., 2022), greedy top-k experts per token (Shazeer et al., 2017) and greedy top-k tokens per expert (Zhou et al., 2022). These works have predominantly considered dot-products as a suitable similarity metric. This work continues with dot-product based learnable routing but computes the routing assignments in an adaptively transformed space to maximally identify the latent expert clusters.

MoE and Cluster Analysis. Recent studies on MoE show the router can recover the clustering structure of the input space and each expert specializes in a specific cluster (Dikkala et al., 2023; Chen et al., 2022). Our work considers transformations of the input space to identify expert clusters, and we learn these transformations via feature-weighted cluster analysis (Brusco & Cradit, 2001; Witten & Tibshirani, 2010; Gnanadesikan et al., 1995). Friedman & Meulman (2004) consider cluster-dependent feature weights to augment iterative clustering algorithms. Our approach similarly uses cluster-dependent feature weights but uses a different optimization problem to derive optimal weights that directly capture the importance of each feature to the clustering solution.

6 CONCLUSION AND FUTURE WORK

In this paper, we present the Adaptive Clustering (AC) router and ACMoE layer, a novel MoE routing method that computes token-expert assignments in a transformed space that maximally identifies latent clusters in the data and more easily discovers the best-matched expert for each token. We adaptively learn for each input which features are relevant to determining its latent cluster assignment and scale its features accordingly, where features that promote tight clustering are upweighted. Our AC routing method enables faster convergence by improving the Hessian conditioning of the router and better robustness by increasing the separation of latent clusters in the transformed space. For ongoing work, we are investigating improved methods for estimating the latent cluster memberships without reliance on previous layers and with provable consistency guarantees.

324 **Reproducibility Statement.** Source code for our experiments are provided in the supplementary
325 material. We provide the full details of our experimental setup – including datasets, model specifi-
326 cation, train regime, and evaluation protocol – for all experiments in Appendix C. All datasets are
327 publicly available.

328 **Ethics Statement.** Our work considers fundamental architectures, and in particular their robustness
329 and convergence properties. Given this, we foresee no issues regarding fairness, privacy, or security,
330 or any other harmful societal or ethical implications in general.

332 REFERENCES

334 Dosovitskiy Alexey. An image is worth 16x16 words: Transformers for image recognition at scale.
335 *arXiv preprint arXiv: 2010.11929*, 2020.

336 Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers.
337 *arXiv preprint arXiv:2106.08254*, 2021.

339 Emmanuel Bengio, Pierre-Luc Bacon, Joelle Pineau, and Doina Precup. Conditional computation
340 in neural networks for faster models. *arXiv preprint arXiv:1511.06297*, 2015.

341 Michael J Brusco and J Dennis Cradit. A variable-selection heuristic for k-means clustering. *Psy-*
342 *chometrika*, 66:249–270, 2001.

344 Iñigo Casanueva, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. Efficient
345 intent detection with dual sentence encoders. *arXiv preprint arXiv:2003.04807*, 2020.

346 Zixiang Chen, Yihe Deng, Yue Wu, Quanquan Gu, and Yuanzhi Li. Towards understanding the
347 mixture-of-experts layer in deep learning. *Advances in neural information processing systems*,
348 35:23049–23062, 2022.

350 Zewen Chi, Li Dong, Shaohan Huang, Damai Dai, Shuming Ma, Barun Patra, Saksham Singhal,
351 Payal Bajaj, Xia Song, Xian-Ling Mao, et al. On the representation collapse of sparse mixture of
352 experts. *Advances in Neural Information Processing Systems*, 35:34600–34613, 2022.

353 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-
354 erarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
355 pp. 248–255. Ieee, 2009.

356 Nishanth Dikkala, Nikhil Ghosh, Raghu Meka, Rina Panigrahy, Nikhil Vyas, and Xin Wang. On the
357 benefits of learning to route in mixture-of-experts models. In *Proceedings of the 2023 Conference*
358 *on Empirical Methods in Natural Language Processing*, pp. 9376–9396, 2023.

360 Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim
361 Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. Glam: Efficient scaling of language
362 models with mixture-of-experts. In *International Conference on Machine Learning*, pp. 5547–
363 5569. PMLR, 2022.

364 William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter
365 models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39,
366 2022.

367 Jerome H Friedman and Jacqueline J Meulman. Clustering objects on subsets of attributes (with
368 discussion). *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 66(4):
369 815–849, 2004.

371 Ram Gnanadesikan, Jon R Kettenring, and Shiao Li Tsao. Weighting and selection of variables for
372 cluster analysis. *Journal of classification*, 12:113–136, 1995.

373 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
374 examples. *arXiv preprint arXiv:1412.6572*, 2014.

376 Yongxin Guo, Zhenglin Cheng, Xiaoying Tang, Zhaopeng Tu, and Tao Lin. Dynamic mix-
377 ture of experts: An auto-tuning approach for efficient transformer models. *arXiv preprint*
arXiv:2405.14297, 2024.

- 378 Xing Han, Tongzheng Ren, Tan Nguyen, Khai Nguyen, Joydeep Ghosh, and Nhat Ho. Design-
379 ing robust transformers using robust kernel density estimation. *Advances in Neural Information*
380 *Processing Systems*, 36, 2024.
- 381
382 Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul
383 Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical
384 analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF international*
385 *conference on computer vision*, pp. 8340–8349, 2021a.
- 386 Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial
387 examples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recogni-*
388 *tion*, pp. 15262–15271, 2021b.
- 389 Kenichi Kumatani, Robert Gmyr, Felipe Cruz Salinas, Linquan Liu, Wei Zuo, Devang Patel, Eric
390 Sun, and Yu Shi. Building a great multi-lingual teacher with sparsely-gated mixture of experts for
391 speech recognition. *arXiv preprint arXiv:2112.05820*, 2021.
- 392
393 D Lepikhin, H Lee, Y Xu, D Chen, O Firat, Y Huang, M Krikun, N Shazeer, and Z Gshard.
394 Scaling giant models with conditional computation and automatic sharding. *arXiv preprint*
395 *arXiv:2006.16668*, 2020.
- 396 Mike Lewis, Shruti Bhosale, Tim Dettmers, Naman Goyal, and Luke Zettlemoyer. Base layers:
397 Simplifying training of large, sparse models. In *International Conference on Machine Learning*,
398 pp. 6265–6274. PMLR, 2021.
- 399
400 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.
401 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the*
402 *IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- 403 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
404 Towards deep learning models resistant to adversarial attacks. *stat*, 1050(9), 2017.
- 405
406 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
407 models. *arXiv preprint arXiv:1609.07843*, 2016.
- 408 Huy Nguyen, Nhat Ho, and Alessandro Rinaldo. On least squares estimation in softmax gating
409 mixture of experts. *arXiv preprint arXiv:2402.02952*, 2024.
- 410
411 Quang Pham, Giang Do, Huy Nguyen, TrungTin Nguyen, Chenghao Liu, Mina Sartipi, Binh T
412 Nguyen, Savitha Ramasamy, Xiaoli Li, Steven Hoi, et al. Competesmoeeffective training of
413 sparse mixture of experts via competition. *arXiv preprint arXiv:2402.02526*, 2024.
- 414 Joan Puigcerver, Carlos Riquelme, Basil Mustafa, and Neil Houlsby. From sparse to soft mixtures
415 of experts. *arXiv preprint arXiv:2308.00951*, 2023.
- 416
417 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
418 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
419 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 420
421 Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André
422 Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts.
423 *Advances in Neural Information Processing Systems*, 34:8583–8595, 2021.
- 424
425 N Shazeer, A Mirhoseini, K Maziarz, A Davis, Q Le, G Hinton, and J Dean. The sparsely-gated
426 mixture-of-experts layer. *Outrageously large neural networks*, 2017.
- 427
428 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng,
429 and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment
430 treebank. In *Proceedings of the 2013 conference on empirical methods in natural language pro-*
431 *cessing*, pp. 1631–1642, 2013.
- Jonathan Uesato, Brendan O’donoghue, Pushmeet Kohli, and Aaron Oord. Adversarial risk and the
dangers of evaluating against weak attacks. In *International conference on machine learning*, pp.
5025–5034. PMLR, 2018.

432 A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
433
434 Daniela M Witten and Robert Tibshirani. A framework for feature selection in clustering. *Journal*
435 *of the American Statistical Association*, 105(490):713–726, 2010.
436
437 Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew M Dai, Quoc V
438 Le, James Laudon, et al. Mixture-of-experts with expert choice routing. *Advances in Neural*
439 *Information Processing Systems*, 35:7103–7114, 2022.
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

Supplement to “Tight Clusters Make Specialized Experts”

Table of Contents

486		
487		
488	Table of Contents	
489		
490	A Technical Proofs	10
491		
492	A.1 Proof of Theorem 1	10
493	A.2 Proof of Proposition 1	12
494	A.2.1 Proof of Lemma 1	12
495	A.2.2 Proof of Lemma 2	12
496	A.3 Proof of Proposition 2	13
497		
498		
499	B Implementation Procedure and Computational Efficiency	15
500		
501		
502	C Experimental Details and Additional Experiments	16
503		
504	C.1 Language Modeling	16
505	C.1.1 Datasets	16
506	C.1.2 Model, Optimizer, & Train Specification	16
507	C.2 Image Classification	17
508	C.2.1 Datasets and Attacks	17
509	C.2.2 Model, Optimizer, & Train Specification	17
510	C.3 Adversarial Attack At Higher Perturbation Budget	18
511	C.4 Cluster Visualization	18
512	C.5 Ablation Studies	19
513	C.5.1 Measures of Dispersion	19
514	C.5.2 Layer Placement	19
515	C.5.3 Random Ablation	20
516	C.6 Cluster Weight Mixing	20
517	C.7 Adaptive Clustering Integration into Soft Mixture of Experts	21
518	C.8 Image Classification in Swin Transformer Base Configuration	21
519	C.9 Router Stability	21
520	C.10 Dynamic Routing	22
521		
522		
523		
524		
525		
526		
527		
528	D Broader Impact	22
529		

A TECHNICAL PROOFS

A.1 PROOF OF THEOREM 1

To begin with, we present the following lemma to show the existence of constants α_k for $k \in [E]$ that satisfy Eqn. ??:

Lemma 3. *For any $\lambda > 0$, Eqn. ?? has exactly d real solutions with respect to α_k .*

540 *Proof of Lemma 3.* Without loss of generality, assume that $s_{1k} \geq s_{2k} \geq \dots \geq s_{dk}$. Denote

$$541 \varphi(\alpha) := \sum_{q \in [d]} \frac{1}{s_{qk} + \alpha} - \frac{d}{\lambda}. \quad (10)$$

542 Then, the existence of solutions to Eqn. ?? is equivalent to the condition $\varphi(\alpha_l) = 0$. Note that $\varphi(\alpha)$
543 is a strictly decreasing function in its connected continuity domains since

$$544 \varphi'(\alpha) = - \sum_{q \in [d]} \frac{1}{(s_{qk} + \alpha)^2} < 0 \quad (11)$$

545 for all $\alpha \in \mathbb{R} \setminus \{-s_{1k}, \dots, -s_{dk}\}$. Further, we observe that

$$546 \lim_{\alpha \rightarrow -s_{qk}^-} \varphi(\alpha) = -\infty, \quad \lim_{\alpha \rightarrow -s_{qk}^+} \varphi(\alpha) = +\infty \quad (12)$$

547 for all $q \in [d]$, and

$$548 \lim_{\alpha \rightarrow \pm\infty} \varphi(\alpha) = -\frac{d}{\lambda} < 0. \quad (13)$$

549 Now consider the domain of continuity of $\varphi(\alpha)$, namely $(-\infty, -s_{1k}) \cup (-s_{1k}, -s_{2k}) \cup \dots \cup (-s_{dk}, \infty)$.
550 Due to the monotonicity and limits 12 & 13, there exists a unique solution in each of the intervals
551 except for $(-\infty, -s_{1k})$ where the function is always strictly negative, thus, yielding d roots in total.
552 \square

553 Now we follow up with the main proof of this section.

554 *Proof of Theorem 1.* First, let $\mathcal{I}_k := \{i : r(i) = k\}$ for convenience. Now let us restate the clustering
555 optimization problem (3) here once again:

$$556 \min_{\mathbf{w}_k} Q(c, \{\mathbf{w}_k\}_{k \in [E]}) = \sum_{k \in [E]} \frac{1}{N_k^2} \sum_{i, j \in \mathcal{I}_k} \sum_{q \in [d]} \left(w_{qk} \rho_{ijq} + \frac{\lambda}{d} \log \frac{1}{dw_{qk}} \right),$$

$$557 \text{ such that } \sum_{q \in [d]} w_{qk} = 1, \quad \forall k \in [E], \quad (14)$$

558 where we have immediately used the fact that

$$559 D_{\text{KL}}(\mathbf{u} \parallel \mathbf{w}_k) = \sum_{q \in [d]} \frac{1}{d} \log \frac{1/d}{w_{qk}}. \quad (15)$$

560 Also, note that

$$561 \sum_{q \in [d]} \left(w_{qk} \rho_{ijq} + \lambda \frac{1}{d} \log \frac{1}{dw_{qk}} \right) = \sum_{q \in [d]} \left(w_{qk} \rho_{ijq} - \lambda \frac{1}{d} \log(dw_{qk}) \right)$$

$$562 = \sum_{q \in [d]} \left(w_{qk} \rho_{ijq} - \frac{\lambda}{d} \log w_{qk} \right) - \lambda \log d. \quad (16)$$

563 We can ignore the term $\lambda \log d$ since it does not depend on the optimization variable. Method of
564 Lagrange multipliers turns this constrained optimization problem into the following unconstrained
565 counterpart:

$$566 \min_{\mathbf{w}_k, \boldsymbol{\alpha}} \mathcal{L}(c, \{\mathbf{w}_k\}_{k \in [E]}, \boldsymbol{\alpha}) = \sum_{k \in [E]} \frac{1}{N_k^2} \sum_{i, j \in \mathcal{I}_k} \sum_{q \in [d]} \left(w_{qk} \rho_{ijq} - \frac{\lambda}{d} \log w_{qk} \right) + \sum_{k \in [E]} \alpha_k \left(\sum_{q \in [d]} w_{qk} - 1 \right),$$

567 where $\boldsymbol{\alpha} = [\alpha_1 \dots \alpha_L]^\top$ is the vector of Lagrange multipliers. Note that the last optimization
568 problem can be separated into the following L independent optimization subproblems:

$$569 \min_{\mathbf{w}_k, \alpha} \mathcal{L}_k(c, \mathbf{w}_k, \alpha) = \frac{1}{N_k^2} \sum_{i, j \in \mathcal{I}_k} \sum_{q \in [d]} \left(w_{qk} \rho_{ijq} - \frac{\lambda}{d} \log w_{qk} \right) + \alpha_k \left(\sum_{q \in [d]} w_{qk} - 1 \right),$$

for $k \in [E]$. Since the objective function is a positive combination of convex functions, the optimization problem is also convex. By setting the derivatives of \mathcal{L}_k with respect to both optimization variables to 0, we obtain the following system of equations:

$$\begin{cases} \frac{\partial \mathcal{L}_k}{\partial w_{qk}} = s_{qk} - \frac{\lambda}{d} \frac{1}{w_{qk}} + \alpha_k = 0, \\ \frac{\partial \mathcal{L}_k}{\partial \alpha_k} = \sum_{q \in [d]} w_{qk} - 1 = 0 \end{cases}$$

for all $k \in [E]$, where s_{qk} is the data dispersion measure defined in the theorem statement. The first equation yields

$$w_{qk} = \frac{\lambda}{d} \frac{1}{s_{qk} + \alpha_k}, \quad (17)$$

where α_k is found from $\sum_{q \in [d]} w_{qk} = 1$ which in fact gives

$$\sum_{q \in [d]} \frac{1}{s_{qk} + \alpha_k} = \frac{d}{\lambda} \quad (18)$$

for all $k \in [E]$ as desired. \square

A.2 PROOF OF PROPOSITION 1

To give the probability bound an exact form, we assume the clusters follow a Gaussian mixture model (GMM). We note that GMMs are a highly expressive and general framework, so this assumption does not place significant restrictions on our analysis. We further assume that though clusters may overlap, they are well-separated along the features for which they cluster tightly¹.

Since Proposition 1 is a composition of Lemma 1 and Lemma 2, we proceed by providing their proofs.

A.2.1 PROOF OF LEMMA 1

Proof of Lemma 1. Notice that we can expand inequality (1) as

$$\sum_{i \in [d]} m_i \delta \mu_i^2 \geq \sum_{i \in [d]} \delta \mu_i^2,$$

where we let $\delta \mu := \mu_b - \mu_a$. Since M_a entries are mean-scaled, we can rewrite them as

$$m_i = \frac{dm'_i}{\sum_{j \in [d]} m'_j} \quad (19)$$

for some initial dispersion estimates $\{m'_j\}_{j \in [d]}$. Without loss of generality, assume that $[d']$ is the set of dimension indices for which the dispersions are relatively much smaller than those in the rest of the dimensions in the sense that $m'_i \gg m'_j$ for any $i \in [d']$ and $j \in [d] \setminus [d']$. Then, there exists a positive $\alpha \ll 1/2$ such that $\sum_{i \in [d']} m_i > d - \alpha$ and $\sum_{i \in [d] \setminus [d']} m_i < \alpha$. By the assumption that clusters are best-separated along the features for which they cluster tightly, this means that the weight matrix M_a maximizes the contribution of largest d' terms in $\sum_{i \in [d]} m_i \delta \mu_i^2$ corresponding to individual feature-wise distances in dimensions where the feature dispersions are the smallest instead of giving uniform weights to all dimensions, which leads to inequality (1). \square

A.2.2 PROOF OF LEMMA 2

Proof of Lemma 2. Since we use the \mathcal{L}_2 distance between the token \mathbf{h} and μ_c as a similarity metric, we assign cluster g_{k^*} to the token \mathbf{h}' iff $\|\mathbf{h}' - \mu_{k^*}\| \leq \|\mathbf{h}' - \mu_k\|$. Assume that the token \mathbf{h}' is a noisy observation of an underlying true token \mathbf{h} which actually originates from cluster g_{k^*} . Then, the token \mathbf{h}' can be decomposed as $\mathbf{h}' = \mathbf{h} + \epsilon$ for a random noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \Sigma_\epsilon)$. Now define

¹Intuitively, this assumption captures the natural property that the semantic regions of the input space are distinct along the dimensions that best identify them.

648 the decision variable $\mathcal{D}(\mathbf{h}') := \|\mathbf{h}' - \boldsymbol{\mu}_{k^*}\|^2 - \|\mathbf{h}' - \boldsymbol{\mu}_k\|^2$ which turns the clustering condition to
 649 $\mathcal{D}(\mathbf{h}') \leq 0$ for the cluster g_{k^*} . Let us analyze the decision variable \mathcal{D} as a random variable where
 650 randomness may come from the underlying sampling strategy and noise. Note that

$$\begin{aligned} 651 \mathcal{D}(\mathbf{h}') &= \|\mathbf{h} + \boldsymbol{\epsilon} - \boldsymbol{\mu}_{k^*}\|^2 - \|\mathbf{h} + \boldsymbol{\epsilon} - \boldsymbol{\mu}_k\|^2 \\ 652 &= \|\mathbf{h} - \boldsymbol{\mu}_{k^*}\|^2 - \|\mathbf{h} - \boldsymbol{\mu}_k\|^2 + 2(\boldsymbol{\mu}_k - \boldsymbol{\mu}_{k^*})^\top \boldsymbol{\epsilon} \\ 653 &= \mathcal{D}(\mathbf{h}) + 2\delta\boldsymbol{\mu}^\top \boldsymbol{\epsilon}, \end{aligned} \quad (20)$$

654 where $\delta\boldsymbol{\mu} := \boldsymbol{\mu}_k - \boldsymbol{\mu}_{k^*}$. Due to the assumption that \mathbf{h} is drawn from the distribution g_{k^*} , it can be
 655 rewritten as $\mathbf{h} = \boldsymbol{\mu}_{k^*} + \boldsymbol{\nu}$ with $\boldsymbol{\nu} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{k^*})$. Then for the first term in Eqn. 20, we have

$$\begin{aligned} 656 \mathcal{D}(\mathbf{h}) &= \|\mathbf{h} - \boldsymbol{\mu}_{k^*}\|^2 - \|\mathbf{h} - \boldsymbol{\mu}_k\|^2 \\ 657 &= \delta\boldsymbol{\mu}^\top (2\mathbf{h} - \boldsymbol{\mu}_{k^*} - \boldsymbol{\mu}_k) \\ 658 &= \delta\boldsymbol{\mu}^\top (2\boldsymbol{\nu} - \delta\boldsymbol{\mu}) \\ 659 &= 2\delta\boldsymbol{\mu}^\top \boldsymbol{\nu} - \|\delta\boldsymbol{\mu}\|^2. \end{aligned} \quad (21)$$

660 Substituting this back into Eqn. 20, we get

$$661 \mathcal{D}(\mathbf{h}') = 2\delta\boldsymbol{\mu}^\top (\boldsymbol{\nu} + \boldsymbol{\epsilon}) - \|\delta\boldsymbol{\mu}\|^2. \quad (22)$$

662 This shows that $\mathcal{D}(\mathbf{h}') \sim \mathcal{N}(-\|\delta\boldsymbol{\mu}\|^2, 4\delta\boldsymbol{\mu}^\top (\boldsymbol{\Sigma}_{k^*} + \boldsymbol{\Sigma}_\epsilon) \delta\boldsymbol{\mu})$. Since $\mathcal{D}(\mathbf{h}')$ follows a normal distri-
 663 bution with the derived parameters, the probability that \mathbf{h}' is assigned to cluster g_{k^*} is given by

$$664 \Pr(\text{correct cluster}) = \Pr(\mathcal{D}(\mathbf{h}') \leq 0) = \Phi\left(\frac{\|\delta\boldsymbol{\mu}\|^2}{2\sqrt{\delta\boldsymbol{\mu}^\top (\boldsymbol{\Sigma}_{k^*} + \boldsymbol{\Sigma}_\epsilon) \delta\boldsymbol{\mu}}}\right), \quad (23)$$

665 where Φ denotes the CDF of normal distribution as usual. Since Φ is an increasing function, the
 666 probability that the noisy token \mathbf{h} is assigned to the correct cluster is proportional to the distance
 667 between the cluster centroids and inverse proportional to the covariance matrices of the cluster and
 668 the additive noise. On the other hand, for the incorrect clustering probability, we have

$$669 \Pr(\text{incorrect cluster}) = 1 - \Phi\left(\frac{\|\delta\boldsymbol{\mu}\|^2}{2\sqrt{\delta\boldsymbol{\mu}^\top (\boldsymbol{\Sigma}_{k^*} + \boldsymbol{\Sigma}_\epsilon) \delta\boldsymbol{\mu}}}\right) \quad (24)$$

670 as claimed. \square

671 A.3 PROOF OF PROPOSITION 2

672 *Proof of Proposition 2.* Let the router be given by g and let the softmax function be given by
 673 $g_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d$, parameterized by expert embeddings $\{e_i\}_{i \in [E]}$. The network loss depends on expert
 674 embeddings only through the router function g . We shall explore the exclusive contribution of each
 675 expert embedding in minimizing $\mathcal{L}^{\text{ACMoE}}$. In order to do this, we look at the network loss as a scalar
 676 function of i^{th} expert embedding vector while treating all other network parameters as fixed. Then,
 677 we can write $\mathcal{L}^{\text{ACMoE}} : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\mathcal{L}^{\text{ACMoE}} = \mathcal{L}^{\text{ACMoE}}(g_\theta(e_i))$. For simplicity, we shall
 678 omit the subscript θ . The gradient that comes from back-propagation is then given by

$$679 \nabla_{e_i} \mathcal{L}^{\text{ACMoE}} = (\nabla_g \mathcal{L}^{\text{ACMoE}})^\top \nabla_{e_i} g, \quad (25)$$

680 where $\nabla_{e_i} g \in \mathbb{R}^{d \times d}$ denotes the Jacobian matrix of g since for $g_k := (g_\theta(e_i))_k$, we can write

$$681 \frac{\partial}{\partial e_{is}} \mathcal{L}^{\text{ACMoE}}(g_1, \dots, g_d) = \sum_k \frac{\partial \mathcal{L}^{\text{ACMoE}}}{\partial g_k} \frac{\partial g_k}{\partial e_{is}}. \quad (26)$$

682 Note that for $g_k = \text{softmax}(\mathbf{h}^\top \mathbf{M} e_k)$, we have

$$683 \frac{\partial g_k}{\partial e_{is}} = m_s h_s g_k (\delta_{ki} - g_i) = m_s h_s b_{ki}. \quad (27)$$

Then, the element of the Hessian matrix of the network loss at index $(s, t) \in [d] \times [d]$ can be written as

$$\begin{aligned}
\mathbf{H}_{st}^{(i)}(\mathcal{L}^{\text{ACMoE}}) &= \frac{\partial^2 \mathcal{L}^{\text{ACMoE}}}{\partial e_{is} \partial e_{it}} = \frac{\partial}{\partial e_{it}} \sum_k \frac{\partial \mathcal{L}^{\text{ACMoE}}}{\partial g_k} \frac{\partial g_k}{\partial e_{is}} \\
&= \sum_k \left(\sum_j \frac{\partial^2 \mathcal{L}^{\text{ACMoE}}}{\partial g_k \partial g_j} \frac{\partial g_j}{\partial e_{it}} \right) \frac{\partial g_k}{\partial e_{is}} + \frac{\partial \mathcal{L}^{\text{ACMoE}}}{\partial g_k} \frac{\partial^2 g_k}{\partial e_{is} \partial e_{it}} \\
&= m_s h_s m_t h_t \left[\sum_k \left(\sum_j \frac{\partial^2 \mathcal{L}^{\text{ACMoE}}}{\partial g_k \partial g_j} b_{ji} \right) b_{ki} + \frac{\partial \mathcal{L}^{\text{ACMoE}}}{\partial g_k} b'_{ki} \right] \\
&= m_s h_s m_t h_t B_i,
\end{aligned} \tag{28}$$

where B_i is some constant that depends only on index i . Due to Eqn. 28, the Hessian takes the following matrix form

$$\mathbf{H}^{(i)} = B_i (\mathbf{M} \mathbf{h}) (\mathbf{M} \mathbf{h})^\top. \tag{29}$$

Taking expectation from both sides, we obtain

$$\mathbb{E}_{\mathbf{h} \sim (\mu, \Sigma)} [\mathbf{H}^{(i)}] = B_i \mathbb{E}_{\mathbf{h} \sim (\mu, \Sigma)} [\mathbf{M} (\mathbf{h} \mathbf{h}^\top) \mathbf{M}] = B_i \mathbf{M} (\Sigma) \mathbf{M}, \tag{30}$$

where we assume \mathbf{h} is centered. Now recall that $\mathbf{M} = \text{diag}(m_1, \dots, m_d)$ where for each i , $m_i \sim 1/\sqrt{\Sigma_{ii}}$ holds. Assume that the covariance matrix Σ is symmetric positive definite. Then, it is diagonalizable as $\Sigma = \mathbf{U} \Lambda \mathbf{U}^\top$ with $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$, a diagonal matrix with eigenvalues of Σ . With the transformation \mathbf{M} , we get

$$\mathbf{M} \Sigma \mathbf{M} = \mathbf{M} \mathbf{U} \Lambda \mathbf{U}^\top \mathbf{M} = \mathbf{U} \mathbf{M} \Lambda \mathbf{M} \mathbf{U}^\top \tag{31}$$

$$= \mathbf{U} \begin{bmatrix} m_1^2 \lambda_1 & & \\ & \ddots & \\ & & m_d^2 \lambda_d \end{bmatrix} \mathbf{U}^\top. \tag{32}$$

Since the eigenvalues capture the variances along the principal components of the covariance matrix, m_i^2 , as a reciprocal of a measure of dimension-wise dispersion, is reasonably correlated with $1/\lambda_i$, as demonstrated by Lemma 4, implying $\lambda_j \leq \lambda_i \implies m_j \geq m_i$ with high probability. Therefore, we obtain that

$$\kappa(\mathbf{M} \Sigma \mathbf{M}) = \frac{\lambda_{\max}(\mathbf{M} \Sigma \mathbf{M})}{\lambda_{\min}(\mathbf{M} \Sigma \mathbf{M})} \approx \frac{m_{\min}^2 \lambda_{\max}(\Sigma)}{m_{\max}^2 \lambda_{\min}(\Sigma)} \leq \kappa(\Sigma), \tag{33}$$

which implies the claim. \square

Lemma 4 (Correlation between dimension-wise variances and covariance eigenvalues). *Let $\{\mathbf{b}_i\}_{i \in [d]}$ be the set of normalized basis vectors of \mathbb{R}^d . Consider a symmetric positive definite covariance matrix Σ and its unit eigenvectors $\{\mathbf{v}_i\}_{i \in [d]}$. Assume that the eigenvector \mathbf{v}_i is a reasonably small perturbation of the basis vector \mathbf{b}_i such that $\mathbf{v}_i^\top \mathbf{b}_i \geq 1 - \epsilon$ for all $i \in [d]$ and a small constant $\epsilon > 0$. Then, for all $i \in [d]$, we have*

$$|\lambda_i - \Sigma_{ii}| \leq \epsilon \cdot \max_{j \neq i} |\lambda_i - \lambda_j|, \tag{34}$$

where $\{\lambda_i\}_{i \in [d]}$ is the set of ordered eigenvalues of Σ corresponding to eigenvectors $\{\mathbf{v}_i\}_{i \in [d]}$.

Proof of Lemma 4. Note that each diagonal element of the SPD covariance matrix Σ can be written as

$$\Sigma_{ii} = \mathbf{b}_i^\top \Sigma \mathbf{b}_i = \mathbf{b}_i^\top \left(\sum_{j \in [d]} \lambda_j \mathbf{v}_j \mathbf{v}_j^\top \right) \mathbf{b}_i = \sum_{j \in [d]} \lambda_j (\mathbf{v}_j^\top \mathbf{b}_i)^2. \tag{35}$$

Table 4: Efficiency Comparison between ACMoE and baseline MoE models

Model	Compute Speed (ms/it)	Max Memory (K)	#Params (M)
<i>GLaM</i> (Du et al., 2022)	422.62	25.69	220
GLaM-ACMoE (Ours)	425.15	25.72	220
<i>Switch Transformer</i> (Fedus et al., 2022)	391.93	34.64	216
Switch-ACMoE (Ours)	393.29	34.68	216
<i>Swin Transformer</i> (Liu et al., 2021)	403.36	22.00	280
Swin-ACMoE (Ours)	408.56	22.19	280

Then, the difference on the left hand side of Eqn. 34 can be bounded as

$$\begin{aligned}
|\lambda_i - \Sigma_{ii}| &= \left| \lambda_i - \sum_{j \in [d]} \lambda_j (\mathbf{v}_j^\top \mathbf{b}_i)^2 \right| = \left| \lambda_i (1 - (\mathbf{v}_i \mathbf{e}_i)^2) - \sum_{j \neq i} \lambda_j (\mathbf{v}_j^\top \mathbf{b}_i)^2 \right| \\
&= \left| \lambda_i \sum_{j \neq i} (\mathbf{v}_j^\top \mathbf{b}_i)^2 - \sum_{j \neq i} \lambda_j (\mathbf{v}_j^\top \mathbf{b}_i)^2 \right| \tag{36}
\end{aligned}$$

$$\begin{aligned}
&= \left| \sum_{j \neq i} (\lambda_i - \lambda_j) (\mathbf{v}_j^\top \mathbf{b}_i)^2 \right| \\
&\leq \max_{j \neq i} |\lambda_i - \lambda_j| \sum_{j \neq i} (\mathbf{v}_j^\top \mathbf{b}_i)^2 \\
&= \max_{j \neq i} |\lambda_i - \lambda_j| (1 - (\mathbf{v}_i \mathbf{b}_i)^2) \tag{37} \\
&\leq \epsilon \max_{j \neq i} |\lambda_i - \lambda_j|,
\end{aligned}$$

where we used the fact that

$$\sum_{j \in [d]} (\mathbf{v}_j^\top \mathbf{b}_i)^2 = \left(\sum_{j=1}^n (\mathbf{v}_j^\top \mathbf{b}_i) \mathbf{v}_j \right)^\top \left(\sum_{k=1}^n (\mathbf{v}_k^\top \mathbf{b}_i) \mathbf{v}_k \right) = \mathbf{b}^\top \mathbf{b} = 1$$

to obtain Eqn. 36 and Eqn. 37 since the eigenvectors of Σ are orthonormal. \square

B IMPLEMENTATION PROCEDURE AND COMPUTATIONAL EFFICIENCY

Training and Inference. Given the AC routing scheme requires the expert assignment per token from the previous layer, we can only implement AC routing from the second layer on. We incorporate AC routing into both training and inference stages. This is because, firstly, AC routing is designed to offer improvements to both clean and contaminated data, and so even in the presence of completely clean train and test data, it is advantageous to incorporate the AC method into both stages. Secondly, it is commonplace to encounter data contamination only at the test stage and indeed highly possible to encounter it in train as well. Therefore, in the interest of robustness as well, AC routing is incorporated into both stages.

Computational Efficiency. Computing the required $\{\mathbf{w}_k\}_{k \in [E]}$ for number of experts E requires no learnable parameters and is obtained simply by computing the mean absolute deviation for each set of tokens assigned to the k^{th} expert. This can be computed using just two computations of the mean – once for the mean per cluster and once again for the mean of the absolute deviations per cluster – done in parallel over all clusters using `torch.index_reduce()` and is of the order $\mathcal{O}(2nd) = \mathcal{O}(n)$ for n tokens. Hence the upper-bound time complexity of the MoE layer is unaffected. We provide in Table 4 additional efficiency analysis in terms of throughput, max GPU memory allocated, and parameters which shows no significant efficiency loss compared to baseline MoE architectures.

C EXPERIMENTAL DETAILS AND ADDITIONAL EXPERIMENTS

C.1 LANGUAGE MODELING

C.1.1 DATASETS

WikiText-103. The WikiText-103² dataset contains around 268K words and its training set consists of about 28K articles with 103M tokens. This corresponds to text blocks of about 3600 words. The validation set and test sets consist of 60 articles with 218K and 246K tokens respectively.

EnWik-8. The EnWik-8 dataset is a byte-level dataset of 100 million bytes derived from Wikipedia that, in addition to English text, also includes markup, special characters, and text in other languages. EnWik-8 contains 90M characters for training, 5M for validation, and 5M for testing.

Stanford Sentiment Treebank-2. The Stanford Sentiment Treebank-2 (SST2) (Socher et al., 2013) is a 2 class corpus with fully labeled parse trees for analysis of the compositional effects of sentiment in language. The dataset consists of 11,855 single sentences extracted from movie reviews. It was parsed with the Stanford parser and includes 215,154 unique phrases from the parse trees, each annotated by 3 human judges.

Stanford Sentiment Treebank-5. Stanford Sentiment Treebank-5 (SST5) (Socher et al., 2013) is a 5 class dataset used for sentiment analysis. It consists of 11,855 single sentences extracted from movie reviews. It includes 215,154 unique phrases from parse trees, each annotated by 3 human judges. Phrases are classified as negative, somewhat negative, neutral, somewhat positive, or positive.

Banking-77. Banking-77 (B77) (Casanueva et al., 2020) is a highly fine-grained 77 class classification dataset comprising 13083 customer service queries labelled with 77 intents.

C.1.2 MODEL, OPTIMIZER, & TRAIN SPECIFICATION

Models. We use as backbones the Switch Transformer (Fedus et al., 2022) and Generalist Language Model (Du et al., 2022). Table 5 contains the specification over self-attention (SA) layers, feed-forward network (FFN) layers, Mixture-of-Experts (MoE) layers, attention span (Att. Span), embedding size and parameter count for both backbones at small and medium configurations for each pretraining task. All backbones use 16 experts with top-2 expert routing.

Table 5: Language Modeling Backbone Specifications

Model	SA Layers	FFN Layers	MoE Layers	Att. Span	Embed Size	Params
<i>WikiText-103 Pretrain</i>						
Switch-small	3	-	3	256	128	70M
Switch-medium	6	-	6	1024	352	216M
GLaM-small	6	3	3	2048	144	79M
GLaM-medium	12	6	6	2048	352	220M
<i>EnWik-8 Pretrain</i>						
Switch	8	-	8	2048	352	36M

Optimizer. All experiments use Adam with a base learning rate of 0.0007. Small configurations use 3000 iterations of learning rate warmup while medium configurations use 4000 iterations.

²www.salesforce.com/products/einstein/ai-research/the-wikitext-dependency-language-modeling-dataset/

Pretrain Specification. For WikiText-103 pretraining, small Switch backbones are trained for 40 epochs with a batch size of 96 and medium Switch backbones are trained for 80 epochs with a batch size of 48. Small GLaM backbones are trained for 60 epochs with a batch size of 48 and medium GLaM backbones are trained for 120 epochs with a batch size of 48. We use 0.01 auxiliary load balancing loss.

For EnWik-8 pretraining, both Switch and GLaM backbones are trained for 80 epochs with batch size 48. We use 0.01 auxiliary load balancing loss.

Finetune Specification. For SST2 and SST5 finetuning, we finetune for 5 epochs using Adam and a base learning rate of 0.001 without warmup and a batch size of 16. For B77 we finetune for 50 epochs using Adam and a base learning rate of 0.00001 without warmup and a batch size of 16.

Compute Resources. All models are trained, evaluated, and finetuned on four NVIDIA A100 SXM4 40GB GPUs.

C.2 IMAGE CLASSIFICATION

C.2.1 DATASETS AND ATTACKS

ImageNet-1K. We use the full ImageNet dataset that contains 1.28M training images and 50K validation images. The model learns to predict the class of the input image among 1000 categories. We report the top-1 and top-5 accuracy on all experiments.

ImageNet-A/O/R. ImageNet-A (Hendrycks et al., 2021b) contains real-world adversarially filtered images that fool current ImageNet classifiers. A 200-class subset of the original ImageNet-1K’s 1000 classes is selected so that errors among these 200 classes would be considered egregious, which cover most broad categories spanned by ImageNet-1K.

ImageNet-O (Hendrycks et al., 2021b) contains adversarially filtered examples for ImageNet out-of-distribution detectors. The dataset contains samples from ImageNet-22K but not from ImageNet-1K, where samples that are wrongly classified as an ImageNet-1K class with high confidence by a ResNet-50 are selected.

ImageNet-R (Hendrycks et al., 2021a) contains various artistic renditions of object classes from the original ImageNet dataset, which is discouraged by the original ImageNet. ImageNet-R contains 30,000 image renditions for 200 ImageNet classes, where a subset of the ImageNet-1K classes is chosen.

Adversarial Attacks. We use produce corrupted ImageNet samples using white box attacks fast gradient sign method (FGSM) (Goodfellow et al., 2014) and projected gradient descent (PGD) (Madry et al., 2017), and black box simultaneous perturbation stochastic approximation (SPSA) (Uesato et al., 2018). FGSM and PGD use a perturbation budget of 1/255 while SPSA uses a perturbation budget 1. All attacks perturb under l_∞ norm. PGD and uses 20 steps with step size of 0.15 and SPSA uses 20 iterations.

C.2.2 MODEL, OPTIMIZER, & TRAIN SPECIFICATION

Models. Our results are based off of the Swin Transformer (Liu et al., 2021) architecture. This backbone uses 4 base layers of depth 2, 2, 18, and 2. The first two base layers each contain 2 self-attention layers and 2 feed-forward layers. The third base layer contains 18 self-attention layers with alternating feed-forward and MoE layers. The final base layer contains 2 self-attention layers with one feed-forward and one MoE layer. The embedding dimension is 96 and the heads per base layer are 3, 6, 12, and 24. We use 16 total experts and present results for both top-1 and top-2 expert routing. The total parameter count is 280M.

Optimizer. We use AdamW with a base learning rate of 1.25e-4, minimum learning rate of 1.25e-7, 0.1 weight decay and cosine scheduling.

Train Specification. We train for 60 epochs with a batch size of 128 and 0.1 auxiliary balancing loss.

Compute Resources. All models are trained and evaluated on four NVIDIA A100 SXM4 40GB GPUs.

C.3 ADVERSARIAL ATTACK AT HIGHER PERTURBATION BUDGET

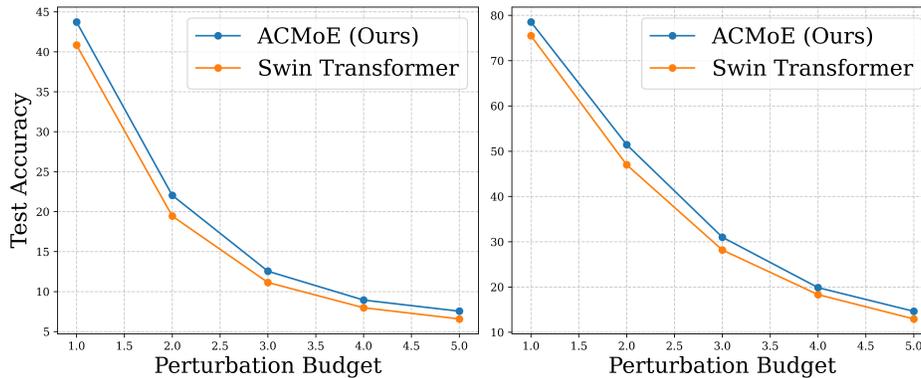


Figure 3: ACMoE and Swin Transformer under PGD attack at increasing perturbation budgets. ACMoE widens its performance gain over Swin at increasingly severe attacks in both top-1 test accuracy (**left**) and top-5 test accuracy (**right**), starting at approximately 7% improvement at 1/255 and ending at just over 10% at 5/255.

Figure 3 shows that for PGD perturbation budgets 1/255 through to 5/255, ACMoE widens its already substantive robust performance gain over Swin, with top-1 and top-5 test accuracy improvements increasing from 7% to approximately 10%.

C.4 CLUSTER VISUALIZATION

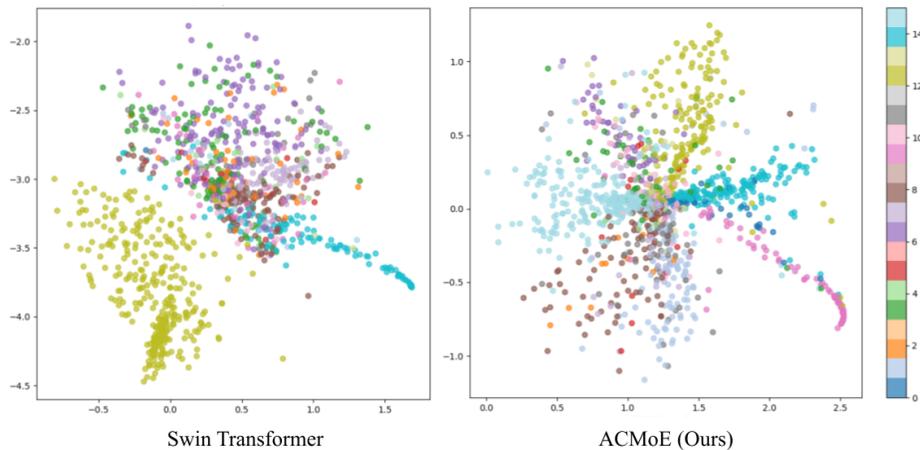


Figure 4: Cluster Visualization on ImageNet. Each token is represented as a point and colored by its assigned expert. **Left:** Swin identifies one cluster clearly (yellow/gold) but otherwise fails to distinguish remaining clusters **Right:** ACMoE learns better-defined expert clusters.

We pass random ImageNet batches through Swin and ACMoE and plot the representations along with their assigned experts, using t-sne to represent the high dimensional data in 2 dimensions. The result is shown in Fig. 4, where we see Swin learns overlapping and indistinguishable expert clusters. ACMoE, on the other hand, performs better in learning the clusters, producing much clearer and better-distinguished clusters.

Table 6: Ablation on Measure of Spread in Switch Transformer (Fedus et al., 2022)

Measure of Spread	Test PPL (\downarrow)
Variance	34.87
MAD	34.42

Table 7: Ablation on Layer Placement in Switch Transformer (Fedus et al., 2022)

Layer Placement	Test PPL (\downarrow)
Back Half	34.95
Alternating	34.80
Skip 1	34.42
Full	34.88

C.5 ABLATION STUDIES

C.5.1 MEASURES OF DISPERSION

We present in Tables 6 and 8 results for Switch-ACMoE and Swin-ACMoE when changing the measure of dispersion used in the AC routing transformation (Definition 1) from mean absolute deviation (MAD) to variance. We see mean absolute deviation outperforms variance as a measure of spread. This is an intuitive finding given that squared distances, as used in variance computations, are highly sensitive to outliers. Using mean absolute deviation as an alternative measure of spread reduces this issue and produces a more robust estimate of dispersion. We note that MAD is not the only robust measure of spread. We conjecture that taking interquartile range as an additionally robust measure of spread may produce good results in both clean and contaminated data. We, however, leave this interesting direction to future research as interquartile range poses implementation challenges as it requires designing concurrent linear scans over the expert clusters. MAD, by contrast, requires just two computations of the mean which is easily parallelizable using `torch.index_reduce()`.

C.5.2 LAYER PLACEMENT

We consider the effect of layer placement in the Switch-medium configuration and in the Swin Transformer (see Sections C.1.2 and C.2.2 for the full model specifications). In particular, Switch is a 6 layer model and Swin is a 24 layer model. With regard to Swin, we focus on the deepest block of depth 18 to implement our ACMoE layers. This is due to the change in embedding size between base layers, meaning we are restricted to this base layer of depth 18. Note further that Swin only uses MoE layers in an alternating pattern with feed-forward networks between each MoE layer. For example, for Switch, a full ACMoE specification would mean placing ACMoE on layers 2,3,4,5,6. For Swin, a full specification means placing ACMoE on layers 4,6,8,10,12,14,16,18. To examine the effect of layer placement we consider the following models:

- *Alternating*: For Switch this means we place ACMoE on layers 2,4,6. For Swin this means we place ACMoE on layers 4,8,12,16.
- *Back Half*: For Switch this means we place ACMoE on just the last 3 layers of the network. For Swin this means we place ACMoE on just the last 5 layers of the network.
- *Skip 2*: For Swin this means we place ACMoE on layers 8,10,12,14,16,18.
- *Skip 1*: For Switch this means we place ACMoE on layers 3,4,5,6. For Swin this means we place ACMoE on layers 6,8,10,12,14,16,18.
- *Full*: We place ACMoE on every possible layer.

We present in Table 7 results for Switch and Swin ACMoE models when changing the positions of the ACMoE layers throughout the network. The results agree with our expectation that, generally speaking, more ACMoE layers improve performance, but a in some circumstances a threshold is met at the point where ACMoE layers are used too early in the network such that the model has not been able to learn reasonably good approximations of the cluster membership of the tokens yet.

We find that in the Switch backbone, performance improves the more ACMoE layers we add, which agrees with our expectation that more ACMoE layers improve performance. However, we find that top performance is attained when allowing two standard MoE layers to go before the first ACMoE, as opposed to the minimum of 1 standard MoE layer. We conjecture this is because we need to give the model a few layers before the first ACMoE in order to learn decent representations such that we

Table 8: Ablation on Measure of Spread in Swin Transformer

Measure of Spread	Test Acc.	
	Top 1	Top 5
<i>Swin-Top1</i> (Liu et al., 2021)		
Variance	75.06	92.49
MAD	75.39	92.56
<i>Swin-Top2</i> (Liu et al., 2021)		
Variance	76.11	93.08
MAD	76.31	93.14

Table 9: Ablation on Layer Placement in Swin Transformer

Layer Placement	Test Acc.	
	Top 1	Top 5
<i>Swin-Top1</i> (Liu et al., 2021)		
Back Half	75.16	92.46
Skip 2	75.34	92.42
Skip 1	75.35	92.45
Full	75.39	92.56
<i>Swin-Top2</i> (Liu et al., 2021)		
Back Half	76.16	93.02
Skip 2	76.10	92.93
Skip 1	76.29	92.98
Full	76.31	93.14

have good enough estimated cluster assignments for use in the ACMoE layer. Encouragingly, we find just one additional standard MoE layer is sufficient for the benefits of ACMoE to be obtained.

We find in Table 9 that with Swin, best performance is obtained using ACMoE on every possible layer, again agreeing with our expectation that more ACMoE layers improve performance. With Swin, however, we do not face any drop in performance from placing ACMoE too early in the network, and indeed we see *Full* attaining top performance. We conjecture that Swin does not encounter this issue since Swin uses four layers of feed forward networks before the first MoE layer, and so by the first MoE layer the representations are of reasonably good quality to produce good estimates of the cluster membership.

C.5.3 RANDOM ABLATION

We show the efficacy of the adaptive clustering transformation M (Definition 1) in our AC router at capturing meaningful feature-wise information by ablating it against an alternate $d \times d$ diagonal matrix made up of normal random variables with mean 1 and standard deviation 0.5 (where we clip any negative values to prevent negative weights). We present in Tables 10 and 11 results for language modeling (using Switch) and image classification (using Swin), which show fairly substantial drops in performance in both backbones. This offers evidence to the claim that our AC routing transformation is meaningfully weighting features to improve routing, and that performance gains of our proposed method do not flow from a kind of implicit regularization of introducing noise into the router.

Table 10: Random Ablation in Switch (Fedus et al., 2022)

Model	Test PPL (\downarrow)
<i>Switch-Random</i> (Fedus et al., 2022)	38.17
Switch-ACMoE	34.42

Table 11: Random Ablation in Swin (Liu et al., 2021)

Model	Top 1 Acc.	Top 5 Acc.
<i>Swin-Random</i>	74.22	91.87
Swin-ACMoE	76.31	93.14

C.6 CLUSTER WEIGHT MIXING

The AC routing scheme estimates the cluster membership of each token based on its highest affinity cluster assigned in the previous layer. We could also further leverage the top-k structure of the MoE models by mixing the cluster-wise feature weights with weights corresponding to the affinities in the top-k routing. For example, if h has affinity scores α and $1 - \alpha$ to clusters k and k' respectively, then we could also obtain the required AC routing transformation for h as $M_{k^*} = \alpha M_k + (1 - \alpha) M_{k'}$. This approach therefore factors in the confidence with which we believe h belongs to cluster k or k' , and can be used for integrating ACMoE into higher expert granularity backbones (i.e higher

top-k settings). Tables 12 and 13 show results for computing M_{k^*} by mixing the top-affinity cluster weights (Mix 2) in Switch and GLaM with top-2 routing, versus our presented results which compute M_{k^*} just based off of the highest affinity cluster (Mix 1). We see that GLaM-ACMoE benefits substantially from cluster weight mixing whereas Switch-ACMoE prefers just using its top affinity cluster weights. For consistency across models, we present in our main body the Mix 1 results, as GLaM-ACMoE already performs extremely strongly using Mix 1 and so we prefer to opt for the added performance gain in the Switch backbone.

Table 12: Results on Cluster Weight Mixing in Switch (Fedus et al., 2022)

Clusters Mixed	Test PPL (\downarrow)
Mix 2	34.66
Mix 1	34.42

Table 13: Results on Cluster Weight Mixing in GLaM (Du et al., 2022)

Clusters Mixed	Test PPL (\downarrow)
Mix 2	35.29
Mix 1	36.26

C.7 ADAPTIVE CLUSTERING INTEGRATION INTO SOFT MIXTURE OF EXPERTS

We present here results for integrating ACMoE into SoftMoE (Puigcerver et al., 2023). To use ACMoE in the SoftMoe setting, which can be understood as a top-E routing setting where all experts are active for every token, we compute M_{k^*} using cluster weight mixing (Section C.6) over the top-8 highest affinity clusters. We present the performance of Soft-ACMoE on clean data, adversarially attacked data, and ImageNet-A/O/R in the following Tables 14 and 15.

Table 14: Test Accuracy on ImageNet corrupted PGD, FGSM, and SPSA using SoftMoE (Puigcerver et al., 2023) backbone

Model	Clean Data		PGD		FGSM		SPSA	
	Top 1	Top 5						
<i>SoftMoE</i> (Puigcerver et al., 2023)	72.86	90.92	45.29	78.91	56.95	85.60	66.59	88.70
Soft-ACMoE (Ours)	73.21	91.23	48.25	80.49	59.01	86.69	70.63	93.22

Table 15: Test Accuracy on Image Classification in Imagenet-A/O/R using SoftMoE (Puigcerver et al., 2023) backbone

Model	Im-A	Im-R	Im-O
	Top-1 Acc. (\uparrow)	Top-1 Acc. (\uparrow)	AUPR (\uparrow)
<i>SoftMoE</i> (Puigcerver et al., 2023)	6.69	31.63	17.97
Soft-ACMoE (Ours)	6.93	32.18	18.35

We see in Tables 14 and 15 the efficacy of ACMoE in the SoftMoE backbone, offering evidence of the adaptability of our framework into further MoE setups. In particular, the SoftMoE framework models a setting in which expert clusters are highly overlapping, as each token is soft assigned to all experts. Therefore, the performance gains shown in clean and contaminated data of Soft-ACMoE demonstrates that our AC router is well-suited to modeling such a clustering structure.

C.8 IMAGE CLASSIFICATION IN SWIN TRANSFORMER BASE CONFIGURATION

We further evaluate the performance ACMoE when scaling up model size in Table 16. We integrate ACMoE into the Base configuration of Swin (0.5B parameters) and evaluate on clean ImageNet-1K as well as under adversarial attacks.

C.9 ROUTER STABILITY

We present in Fig. 5 the routing stability of ACMoE, SMoE, XMoE, and StableMoE in the Switch backbone evaluated on WikiText-103. Routing instability computes over adjacent layers the proportion of tokens that are assigned to different experts across the two layers. Specifically, for n tokens $[h_1, \dots, h_n]$, we compute at layer ℓ the matrix $S^\ell \in \mathbb{R}^{n \times n}$ such that $S_{ij}^\ell = 1$ if the i^{th} and j^{th} tokens are assigned to the same expert in layer ℓ and is 0 otherwise. The router instability at layer ℓ can

Table 16: Test Accuracy on ImageNet corrupted PGD, FGSM, and SPSA using Swin Base (Liu et al., 2021) backbone

Model	Clean Data		PGD		FGSM		SPSA	
	Top 1	Top 5						
<i>Swin-Base</i> (Liu et al., 2021)	79.06	94.37	44.61	79.20	59.91	87.72	68.94	89.00
Swin-ACMoE-Base (Ours)	79.25	94.42	46.28	80.24	61.78	87.55	70.18	89.33

then be calculated as $r^\ell = \text{mean}(|S^{\ell-1} - S^\ell|)$. This metric therefore captures the degree to which tokens that are assigned to the same experts remain together through the model. A high r^ℓ indicates the router doesn't maintain consistent expert assignments, as tokens that it considers semantically similar at one layer it considers different at the next.

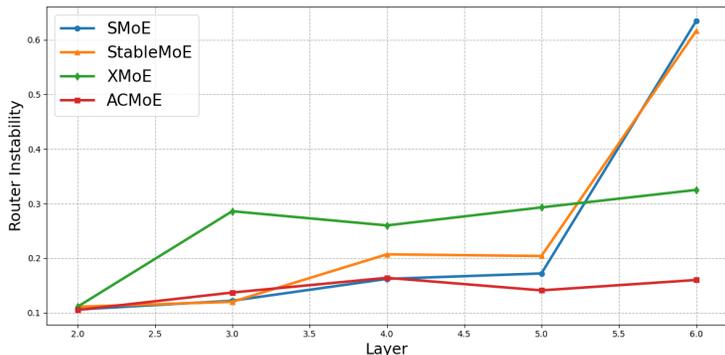


Figure 5: Router Instability of ACMoE, SMoE, XMoE, and StableMoE. ACMoE maintains consistent routing, while baseline routers more frequently change the expert assignments of tokens.

In Fig. 5, we see that baseline routers reach high levels of instability, where in the case of SMoE and StableMoE, at the last layer over 60% of tokens are assigned to a different expert. ACMoE, by contrast, maintains a more consistent, stable assignment through the model, with no more than 20% of tokens changing expert assignment across any layer.

C.10 DYNAMIC ROUTING

We further test the compatibility of our Adaptive Clustering routing scheme in dynamic top-p routing. In this setting, rather than routing each token to its top-k highest affinity experts in each MoE layer, we route each token to all experts that have affinity over a certain threshold p . This setting permits activating more or less experts for different tokens at different layers throughout the model, therefore dynamically assigning experts to tokens. We integrate our AC routing directly into this setting using the same setup as in Section 3, where the AC routing transformation is computed based on the estimated cluster membership of each token using the top affinity assignment of the previous layer. We present the results for Switch transformer on WikiText-103 language modeling in the following Table 17.

For fixed p , we set $p = 0.05$. For learnable p , we initialize the parameter to 0.05. We select this initialization as it reproduces approximately similar performance in the Switch backbone under default top-2 routing, thereby aiding direct comparison between fixed top-k and dynamic top-p routing. We see in the dynamic routing setting, ACMoE maintains the same consistent improvement over the Switch baseline of roughly 1 full PPL. These results suggest ACMoE is well-suited to the dynamic routing setting.

D BROADER IMPACT

Our research offers benefits to Mixture-of-Expert (MoE) architectures in both clean and contaminated settings. In particular, our work offers socially beneficial outcomes with regard to defense

Table 17: Results on Top- p Dynamic Routing in Switch Backbone (Fedus et al., 2022)

Model	Test PPL (\downarrow)
<i>Fixed top-k routing (Shazeer et al., 2017)</i>	
<i>Switch-medium (Fedus et al., 2022)</i>	35.48
ACMoE-medium (Ours)	34.42
<i>Dynamic top-p routing (Guo et al., 2024)</i>	
<i>Switch-Fixed p</i>	35.20
Switch-ACMoE-Fixed p (Ours)	34.14
<i>Switch-Learnable p</i>	34.29
Switch-ACMoE-Learnable p (Ours)	33.49

against adversarial attack, which we hope can be used to protect important AI systems from malicious actors. Furthermore, as large language models, many of which are built on MoE backbones, continue to profligate and be used in important societal settings, we hope our improved robustness to data contamination can aid this promising technology to continue to grow and improve in realistic settings of noisy training and evaluation data. Our research also shows substantially faster convergence than comparative baselines. We believe this faster convergence can deliver significant social benefit in terms of reducing the energy requirements of large model training, thereby helping to ease the growing environmental burden of AI training runs. We recognize there will always be risk of misuse with AI systems, however we hope that our work can be used to enhance and protect socially beneficial AI while also decreasing the environmental impact of this technology. We furthermore hope that our research can spur others on to continue building on robust and efficient AI for social good.