

# 000 MOE MEETS REPARAMETERIZATION: REPARAM- 001 ETERIZABLE MIXTURE-OF-EXPERTS MODEL EN- 002 HANCES DERMATOLOGY DIAGNOSIS VIA DENSE-TO- 003 EXPERTS DISTILLATION 004

005 **Anonymous authors**  
006 Paper under double-blind review  
007

## 013 1 RELATED WORK

014 **Mixture-of-Experts.** Early sparsely gated MoE layers demonstrated that conditional parameter  
015 utilisation can enlarge model capacity without proportional compute (Shazeer et al., 2017). GShard  
016 scaled this idea to  $\sim 600$  B parameters in multilingual translation by integrating automatic tensor  
017 sharding with expert routing (Lepikhin et al., 2020). Switch Transformer further simplified routing  
018 to one-expert selection and trained a trillion-parameter language model with reduced communica-  
019 tion overhead (Fedus et al., 2022). GLaM extended sparsity to multilingual web corpora and showed  
020 strong zero-shot generalisation (Du et al., 2022). Subsequent open-source efforts improved stability  
021 and transfer: ST-MoE introduced noise-based gating to equalise expert load (Zoph et al., 2022),  
022 OpenMoE provided modular implementations for continual pre-training (Xue et al., 2024), and  
023 MegaScale-MoE proposed pipeline parallelism for production-scale deployment (Jin et al., 2025).  
024 Auxiliary regularisers alleviate expert imbalance and collapse, for example expert-choice routing  
025 encourages uniform utilisation (Zhou et al., 2022), gating dropout reduces communication cost (Liu  
026 et al., 2022), and residual MoE blends dense and sparse paths (Wu et al., 2022a). Vision transform-  
027 ers later adopted sparse experts: V-MoE matched the accuracy of large dense ViT-G while saving  
028 FLOPs through token-level routing (Riquelme et al., 2021), M<sup>3</sup>ViT co-designed accelerator-aware  
029 MoE blocks for multi-task learning (Fan et al., 2022), ViMoE explored depth-wise expert placement  
030 (Han et al., 2024), and Mobile V-MoE delivered mobile-scale latency by weight pruning and sparsity  
031 (Daxberger et al., 2023).

032 **Reparameterisation.** Structural reparameterisation trains over-parameterised multi-branch blocks  
033 then algebraically merges them into single-branch kernels for inference efficiency. ACNet first  
034 showed asymmetric convolutions can be fused into  $3 \times 3$  filters for robustness (Ding et al., 2019).  
035 RepVGG stacked  $3 \times 3$  convolutions during deployment and achieved  $> 80\%$  ImageNet top-1  
036 with VGG-style speed (Ding et al., 2021c). Diverse Branch Block generalised fusion to heteroge-  
037 neous branches for improved generalisation (Ding et al., 2021b). ExpandNets applied linear over-  
038 parameterisation to compress student CNNs after training (Guo et al., 2020). ResRep decoupled  
039 “remembering” and “forgetting” filters to prune channels losslessly (Ding et al., 2021a). FastViT  
040 combined RepMixer token mixing with ViT blocks to reduce latency and improve accuracy (Vasu  
041 et al., 2023). RepFormer extended the idea to facial landmark transformers by fusing pyramid heads  
042 (Li et al., 2022).

043 **Dermatology Vision Models.** Esteva *et al.* trained an end-to-end CNN on  $\sim 130,000$  clinical and  
044 dermoscopic images and reached dermatologist-level melanoma detection (Esteva et al., 2017). Sub-  
045 sequent reader studies confirmed AI superiority over clinicians in large cohort benchmarks (Brinker  
046 et al., 2019). Human-computer collaboration improved overall sensitivity by combining visual  
047 saliency with expert opinion (Tschandl et al., 2020). Reinforcement learning further aligned pre-  
048 dictions with clinical risk preferences (Barata et al., 2023). Self-supervised pre-training on 75 M  
049 radiology and dermatology images advanced few-shot transfer (Azizi et al., 2021). Federated self-  
050 supervised contrastive learning preserved privacy across mobile dermatology assistants (Wu et al.,  
051 2022b). SwAVDerm mined unannotated community images via contrastive clustering and improved  
052 long-tail disease coverage (Shen et al., 2024). PanDerm combined slide-level and pixel-level super-  
053 vision over  $> 2$  M images plus 45 K pathology reports to set new state of the art across multiple  
benchmarks (Yan et al., 2025).

054 REFERENCES  
055

056 Shekoofeh Azizi, Basil Mustafa, Fiona Ryan, Zachary Beaver, Jan Freyberg, Jonathan Deaton,  
057 Aaron Loh, Alan Karthikesalingam, Simon Kornblith, Ting Chen, et al. Big self-supervised  
058 models advance medical image classification. In *Proceedings of the IEEE/CVF international*  
059 *conference on computer vision*, pp. 3478–3488, 2021.

060 Catarina Barata, Veronica Rotemberg, Noel CF Codella, Philipp Tschandl, Christoph Rinner,  
061 Bengu Nisa Akay, Zoe Apalla, Giuseppe Argenziano, Allan Halpern, Aimilios Lallas, et al. A  
062 reinforcement learning model for ai-based decision support in skin cancer. *Nature Medicine*, 29  
063 (8):1941–1946, 2023.

064 Titus J Brinker, Achim Hekler, Alexander H Enk, Joachim Klode, Axel Hauschild, Carola Berking,  
065 Bastian Schilling, Sebastian Haferkamp, Dirk Schadendorf, Tim Holland-Letz, et al. Deep learning  
066 outperformed 136 of 157 dermatologists in a head-to-head dermoscopic melanoma image  
067 classification task. *European Journal of Cancer*, 113:47–54, 2019.

068 Erik Daxberger, Floris Weers, Bowen Zhang, Tom Gunter, Ruoming Pang, Marcin Eichner, Michael  
069 Emmersberger, Yinfai Yang, Alexander Toshev, and Xianzhi Du. Mobile v-moes: Scaling down  
070 vision transformers via sparse mixture-of-experts. *arXiv preprint arXiv:2309.04354*, 2023.

071 Xiaohan Ding, Yuchen Guo, Guiguang Ding, and Jungong Han. Acnet: Strengthening the kernel  
072 skeletons for powerful cnn via asymmetric convolution blocks. In *Proceedings of the IEEE/CVF*  
073 *international conference on computer vision*, pp. 1911–1920, 2019.

074 Xiaohan Ding, Tianxiang Hao, Jianchao Tan, Ji Liu, Jungong Han, Yuchen Guo, and Guiguang  
075 Ding. Resrep: Lossless cnn pruning via decoupling remembering and forgetting. In *Proceedings*  
076 *of the IEEE/CVF international conference on computer vision*, pp. 4510–4520, 2021a.

077 Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Diverse branch block: Building  
078 a convolution as an inception-like unit. In *Proceedings of the IEEE/CVF conference on computer*  
079 *vision and pattern recognition*, pp. 10886–10895, 2021b.

080 Xiaohan Ding, Xiangyu Zhang, Ningning Ma, Jungong Han, Guiguang Ding, and Jian Sun. Repvgg:  
081 Making vgg-style convnets great again. In *Proceedings of the IEEE/CVF conference on computer*  
082 *vision and pattern recognition*, pp. 13733–13742, 2021c.

083 Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim  
084 Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. Glam: Efficient scaling of language  
085 models with mixture-of-experts. In *International conference on machine learning*, pp. 5547–  
086 5569. PMLR, 2022.

087 Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and  
088 Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks.  
089 *nature*, 542(7639):115–118, 2017.

090 Zhiwen Fan, Rishov Sarkar, Ziyu Jiang, Tianlong Chen, Kai Zou, Yu Cheng, Cong Hao, Zhangyang  
091 Wang, et al. M<sup>3</sup>vit: Mixture-of-experts vision transformer for efficient multi-task learning with  
092 model-accelerator co-design. *Advances in Neural Information Processing Systems*, 35:28441–  
093 28457, 2022.

094 William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter  
095 models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39,  
096 2022.

097 Shuxuan Guo, Jose M Alvarez, and Mathieu Salzmann. Expandnets: Linear over-parameterization  
098 to train compact convolutional networks. *Advances in Neural Information Processing Systems*,  
099 33:1298–1310, 2020.

100 Xumeng Han, Longhui Wei, Zhiyang Dou, Zipeng Wang, Chenhui Qiang, Xin He, Yingfei Sun,  
101 Zhenjun Han, and Qi Tian. Vimoe: An empirical study of designing vision mixture-of-experts.  
102 *arXiv preprint arXiv:2410.15732*, 2024.

108 Chao Jin, Ziheng Jiang, Zhihao Bai, Zheng Zhong, Juncai Liu, Xiang Li, Ningxin Zheng, Xi Wang,  
 109 Cong Xie, Wen Heng, et al. Megascale-moe: Large-scale communication-efficient training of  
 110 mixture-of-experts models in production. *arXiv preprint arXiv:2505.11432*, 2025.

111 Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang,  
 112 Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional  
 113 computation and automatic sharding. *arXiv preprint arXiv:2006.16668*, 2020.

114 Jinpeng Li, Haibo Jin, Shengcai Liao, Ling Shao, and Pheng-Ann Heng. Refformer: Refinement  
 115 pyramid transformer for robust facial landmark detection. *arXiv preprint arXiv:2207.03917*,  
 116 2022.

117 Rui Liu, Young Jin Kim, Alexandre Muzio, and Hany Hassan. Gating dropout: Communication-  
 118 efficient regularization for sparsely activated transformers. In *International Conference on Ma-  
 119 chine Learning*, pp. 13782–13792. PMLR, 2022.

120 Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André  
 121 Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts.  
 122 *Advances in Neural Information Processing Systems*, 34:8583–8595, 2021.

123 Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton,  
 124 and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer.  
 125 *arXiv preprint arXiv:1701.06538*, 2017.

126 Yue Shen, Huanyu Li, Can Sun, Hongtao Ji, Daojun Zhang, Kun Hu, Yiqi Tang, Yu Chen, Zikun  
 127 Wei, and Junwei Lv. Optimizing skin disease diagnosis: harnessing online community data with  
 128 contrastive learning and clustering techniques. *NPJ Digital Medicine*, 7(1):28, 2024.

129 Philipp Tschandl, Christoph Rinner, Zoe Apalla, Giuseppe Argenziano, Noel Codella, Allan  
 130 Halpern, Monika Janda, Aimilios Lallas, Caterina Longo, Josep Malvehy, et al. Human–computer  
 131 collaboration for skin cancer recognition. *Nature medicine*, 26(8):1229–1234, 2020.

132 Pavan Kumar Anasosalu Vasu, James Gabriel, Jeff Zhu, Oncel Tuzel, and Anurag Ranjan. Fastvit:  
 133 A fast hybrid vision transformer using structural reparameterization. In *Proceedings of the  
 134 IEEE/CVF international conference on computer vision*, pp. 5785–5795, 2023.

135 Lemeng Wu, Mengchen Liu, Yinpeng Chen, Dongdong Chen, Xiyang Dai, and Lu Yuan. Residual  
 136 mixture of experts. *arXiv preprint arXiv:2204.09636*, 2022a.

137 Yawen Wu, Dewen Zeng, Zhepeng Wang, Yi Sheng, Lei Yang, Alaina J James, Yiyu Shi, and  
 138 Jingtong Hu. Federated self-supervised contrastive learning and masked autoencoder for derma-  
 139 tological disease diagnosis. *arXiv preprint arXiv:2208.11278*, 2022b.

140 Fuzhao Xue, Zian Zheng, Yao Fu, Jinjie Ni, Zangwei Zheng, Wangchunshu Zhou, and Yang  
 141 You. Openmoe: An early effort on open mixture-of-experts language models. *arXiv preprint  
 142 arXiv:2402.01739*, 2024.

143 Siyuan Yan, Zhen Yu, Clare Primiero, Cristina Vico-Alonso, Zhonghua Wang, Litao Yang, Philipp  
 144 Tschandl, Ming Hu, Lie Ju, Gin Tan, et al. A multimodal vision foundation model for clinical  
 145 dermatology. *Nature Medicine*, pp. 1–12, 2025.

146 Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew M Dai, Quoc V  
 147 Le, James Laudon, et al. Mixture-of-experts with expert choice routing. *Advances in Neural  
 148 Information Processing Systems*, 35:7103–7114, 2022.

149 Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and  
 150 William Fedus. St-moe: Designing stable and transferable sparse expert models. *arXiv preprint  
 151 arXiv:2202.08906*, 2022.

152

153

154

155

156

157

158

159

160

161