Upcycling Large Language Models into Mixture of Experts

Ethan He, Abhinav Khattar, Ryan Prenger, Vijay Korthikanti, Zijie Yan, Tong Liu, Shiqing Fan, Ashwath Aithal, Mohammad Shoeybi, Bryan Catanzaro

NVIDIA

Abstract

Upcycling pre-trained dense language models into sparse mixture-of-experts (MoE) models is an efficient approach to increase the model capacity of already trained models. However, optimal techniques for upcycling at scale remain unclear. In this work, we conduct an extensive study of upcycling methods and hyperparameters for billion-parameter scale language models. We propose a novel "virtual group" initialization scheme and weight scaling approach to enable upcycling into fine-grained MoE architectures. Through ablations, we find that upcycling outperforms continued dense model training. In addition, we show that softmax-then-topK expert routing improves over topK-then-softmax approach and higher granularity MoEs can help improve accuracy. Finally, we upcycled Nemotron-4 15B on 1T tokens and compared it to a continuously trained version of the same model on the same 1T tokens: the continuous trained model achieved 65.3% MMLU, whereas the upcycled model achieved 67.6%. Our results offer insights and best practices to effectively leverage upcycling for building MoE language models.

1 Introduction

Sparse Mixture of Experts (MoE) models [1] are becoming increasingly popular [2, 3, 4, 5, 6, 7] since they can help achieve better accuracy without a commensurate increase in model training compute. Most recently, state-of-the-art LLMs like Grok-1¹, DBRX², Phi-3.5³, Mixtral 8x22B [3], DeepSeek-V2 [8] and Qwen2 [9] are MoE models. However, an immense amount of compute has been spent on pre-training dense LLMs with only one MLP layer (one expert) [10, 11, 12, 13, 14]. These existing dense models may be able to achieve better accuracy for the same compute cost if they had access to more parameters through MoEs. Upcycling pre-trained dense language models into sparse mixture-of-experts models (referred to as simply 'upcycling' in this work) has emerged as an efficient approach to increase model capacity without the need for training from scratch [15, 16, 17, 3]. By leveraging the knowledge captured in existing dense checkpoints, upcycling enables the creation of large-scale mixture-of-experts (MoE) models while reducing the computational cost and time required for training.

Most previous work on upcycling either does not provide details into how their models were upcycled [3], or provides experiments only at a small scale [15]. We also find that the recommendations in [9] lead to sub-optimal models. To improve general knowledge on upcycling methods, we therefore publish this study of upcycling methods and best practices for billion-parameter scale language models.

In this work, we conduct an extensive study of upcycling techniques and hyperparameters. Our contributions are as follows:

- 1. We recommend training recipes to consistently upcycle billion-parameter scale LLMs.
- 2. We perform a comprehensive study to find the best hyperparameters for upcycling including learning rate, batch size, and load balancing loss.

 $^{\ ^*} equal\ contribution.\ Correspondence\ to:\ \{yihuih,akhattar,rprenger\}@nvidia.com$

¹https://x.ai/blog/grok-os

²https://www.databricks.com/blog/introducing-dbrx-new-state-art-open-llm

³https://huggingface.co/microsoft/Phi-3.5-MoE-instruct

- 3. We propose a novel "virtual group" initialization scheme to enable upcycling into fine-grained MoE architectures, along with a weight scaling approach which brings 1.5% better loss to both coarse-grained and fine-grained upcycled MoE models.
- 4. We compare softmax-then-topK expert routing with the topK-then-softmax approach.
- 5. We assess the benefits of higher granularity MoE models and using higher topK.

We demonstrate that our upcycling approach produces a better model than continued dense model training, softmax-then-topK routing improves over topK-then-softmax approach, and higher granularity MoEs can help boost model accuracy in certain training scenarios. Finally, we upcycle the Nemotron-4 15B model [13] into MoE on 1T tokens and show that it improves MMLU. To make the comparison fair, we train the Nemotron-4 model for an additional 1T tokens and achieve a 65.3% MMLU score, whereas our Nemotron-4 model which was upcycled on the same 1T tokens achieves 67.6. This shows that the improvement from upcycling is not just due to the additional tokens the model is trained on, but also due to the MoE architecture.

By sharing our upcycling approach and ablation results, we aim to contribute to the growing body of knowledge on efficient and scalable language model development, enabling researchers and practitioners to build upon our work and further advance the field of large-scale MoE models. We use Megatron-LM⁴ [18] to upcycle and train our MoE models.

2 Methodology

2.1 Sparse Mixture of Experts

In this work, we only investigate MoEs on the MLP layer of the transformer. These layers comprise the majority of compute and treat each token individually, avoiding issues with kv-cache consistency. A routing layer routes the tokens to a subset of multiple possible MLP layers. This increases the parameter count and presumably the model capacity without necessarily increasing the amount of compute required (measured in total training FLOPs).

2.1.1 Upcycling

Upcycling is the approach of converting a trained dense model into an MoE [15]. The most obvious way to convert a dense model into an MoE model without losing accuracy is to duplicate the MLP layer weights multiple times and use a randomized router, weighting the output of MLP layers by their probabilities

$$E_1(x) := E_2(x) \dots := E_N(x) := FFN(x)$$

$$MoE_activation = \sum_{i=1}^{T} P_i \times E_i(x)$$
(1)

where E_i are MLP experts with router probability P_i such that $P_i \in topK$. x is the output from attention, N is the total number of experts in the MoE layer, and T is the number of experts every token is routed to (topK).

2.1.2 Softmax - TopK Order

The standard MoE router formulation [1, 19] performs a softmax on router logits followed by a topK (softmax-then-topK). The activations from the experts are then multiplied by the softmax probability. In this case, the weight of each expert is found using:

$$TopK(Softmax(x \cdot W_r))$$

where x is the input to the MoE block and W_r is the router.

This causes an issue with upcycling where the output of the upcycled model is not equivalent to the dense model right after upcycling when TopK < N. Even though the output of the upcycled model is

⁴https://github.com/NVIDIA/Megatron-LM

Upcycled Transformer

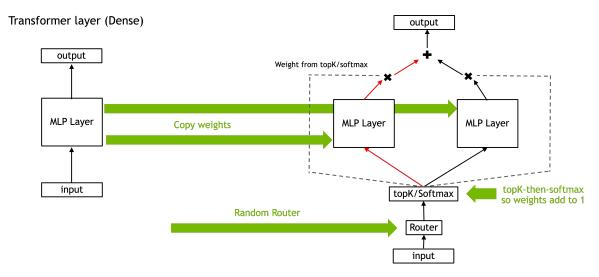


Figure 1: Converting pre-trained dense checkpoint into Mixture of Experts. MLP weights are duplicated to initialize the weights of the experts. The router is randomly initialized. Softmax is often applied after TopK to ensure the upcycled MoE is functionally the same as dense on the first iteration.

not equivalent to the dense model, training the model for a few steps might be able to adjust for this change.

Another way to fix this problem is to use the topK operator directly on the router logits and then only use logits of these topK experts for the softmax computation (topK-then-softmax). In this case, the weight of each expert is found using:

$$Softmax(TopK(x \cdot W_r))$$

The downside of taking this approach is that the information contained in the absolute magnitude of the router output is lost. Also, this approach only works for top K > 1 as softmax of a single element is a constant 1 which has no gradient w.r.t the input. This technique is used in Mixtral [3].

In this study, we compare upcycling with both topK-then-softmax and softmax-then-topK in topK> 1 regime to see which works better (section 3.4).

2.2 Granularity

Earlier work on MoE routed each token to a very small number of experts (topK = 1 or 2) [19, 5]. Routing to only one expert guarantees that the training FLOPs stay similar to the dense model, even though the MoE has more parameters. However, it has recently been suggested that increasing the number of experts to which a token is routed to, while shrinking the dimension of each expert might be a superior approach [20]. This approach is referred to as granular mixture of experts, shown in Figure 2.

Granularity introduces a new degree of freedom as every expert can be reduced in size. Since shrinking experts reduces FLOPs per expert, this approach allows us to increase topK by the same magnitude as the shrinking and still keep the overall FLOPs count the same. While FLOPs is only a proxy for the actual compute required to train or deploy a model, it is still useful and an easy metric to compare compute cost.

There are three hyperparameters that define a fine-grained MoE architecture. We use the nomenclature proposed in [20] and add another term T to refer to topK. The three hyperparameters convey the following:

• E: Expansion rate. How many times larger is the total number of parameters in the MoE layer as compared to the dense MLP counter part. $(N_{\text{MoE}}/N_{\text{dense_MLP}})$

- G: Granularity. How many times smaller is the expert hidden size compared with the original dense layer's FFN hidden size. $(d_{\text{ffn}}/d_{\text{expert}})$
- T: TopK. How many experts is a token being routed to.

For example, in Figure 2, from left to right are coarse-grained MoE E2G1T1 and fine-grained MoE E2G2T2.

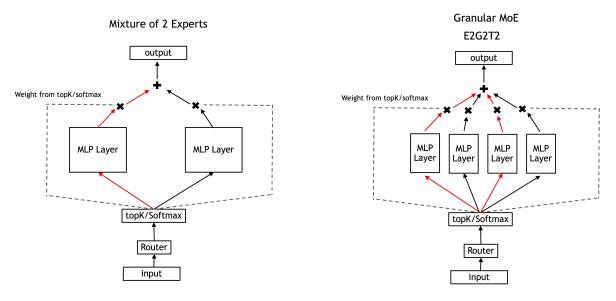


Figure 2: Finegrained Mixture of Experts reduces the size of each expert but activates more experts.

2.3 Granular Upcycling

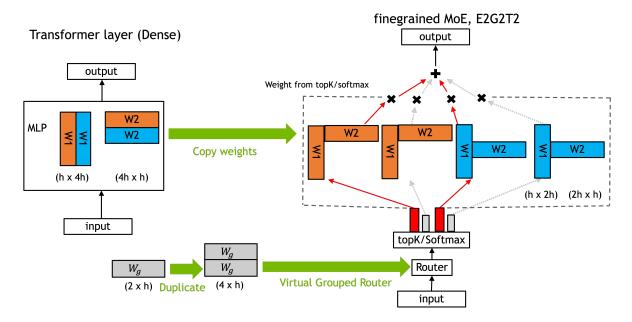


Figure 3: An example of granular upcycling a dense layer into E2G2T2 finegrained MoE. E2G2T2 denotes 4 experts, top 2, with half intermediate size. (1) We shard MLP weights in the intermediate dimension $(4h \to 2h)$ then duplicate the shards. (2) We initialize half the router weights then duplicate them. This ensures Top2 always selects one of each MLP shard so MoE output is the same as the dense model at the start of training.

Unlike standard MoE upcycling [15] where we can copy the dense MLP weights to MoE experts, granularity reduces the size of every MoE expert. This makes copying the dense MLP weights to MoE experts non-trivial. An intuitive way to upcycle a dense model into finegrained MoE would be to:

- 1. segment the dense layer into several shards (G) in the FFN dimension.
- 2. replicate each shard several times (E).
- 3. route to some experts (T) at training time.

For example, segmenting into 8 shards then replicating 8 times yields 64 experts. However, using this naive approach we found that the resulting model's loss was very high and the network did not converge to the original loss. There are two problems with this approach:

- 1. The expert output is scaled down by the router. For 8 expert top-2 model (E8G1T2), a straight forward way of upcycling is to use topK-then-softmax router which ensures that the topK probabilities sum up to 1. However, for fine-grained MoE, even when using topK-then-softmax strategy the outputs are still scaled down. For a 64 experts MoE top-8 model (E8G8T8), if using topK-then-softmax strategy, the output is scaled down by a factor of 8. If using softmax-then-topK strategy, the router probability is about 1/64 because of the random initialization. This implies that the expert outputs are scaled down by a factor of 64.
- 2. The MoE is no longer functionally similar to the dense model in the forward pass. For a coarse grained MoE, using the topK-then-softmax strategy, the upcycled model is functionally the same as dense model in the forward pass in the beginning of training. However, for fine-grained MoE, because the experts are segmented into smaller shards, the router needs to select exactly one replica from each segment in order to function the same as a dense MLP layer.

Motivated by these two observations, we propose weight scaling and virtual grouping for finegrained MoE upcycling. This approach is shown in Figure 3. We initialize the router using *virtual group initialization*. Virtual group initialization ensures that there is exactly one copy of every MLP shard within the router topK right after the dense model is converted into an MoE. Virtual group initialization initializes the experts and router weights such that:

- 1. each router group owns the duplicates of exactly one dense MLP shard.
- 2. all router groups have the same router weights.

A pseudo code snippet illustrating this can be found in appendix A.

2.3.1 Scaling the Weights

We found that with granular upcycling, the scaling of the network weights greatly influences the accuracy of the fine-tuned MoE model. While this scaling could be done entirely in the second linear projection of the MLP (W2 in Figure 3), we found empirically that this works worse than scaling both the linear projection weights (W1 and W2). Equation 2 calculates this scaling factor for the case of squared-relu activation which we use for our base 15B dense model.

$$MoE_activation = \sum_{i=1}^{T} P_i \times E_i(x)$$

where P_i is the probability for the top i^{th} expert, E_i is the corresponding expert layer, T is the topK, and x is the output from attention.

assuming approximately uniform distribution⁵ for iteration 0

$$P = P_1 = P_2 = \dots P_T = \frac{1}{E \times G}$$

 $^{^5}$ while this is a simplifying assumption, it holds true for our most important virtual grouping case where the shrinking factor is the same as topK

$$MoE_activation = P \times (\sum_{i=1}^{T} E_i(x))$$

So for virtual grouping,

$$MoE_activation = \frac{1}{E \times G} (\frac{T}{G} \times dense_activation) = \frac{T}{E \times G^2} \times dense_activation$$

assuming squared relu activation, we normalize W_1 and W_2 for each expert's MLP in case of virtual grouping by:

 $\sqrt[3]{\frac{E \times G^2}{T}} \tag{2}$

While for different activation functions a hyperparameter search for the optimal scaling of the input and output weights might be intuitively better, we find empirically that the equal weight distribution like above works just as well for our 2B models using swiglu activation. We also use this weight scaling for non-granular (a.k.a coarse-grained) MoE models and observe that it helps convergence.

3 Results

Model: We do all our ablation studies on our Nemotron $2B^6$ and Nemotron-4 15B [13] models followed by final results using a larger token count on Nemotron-4 15B. Nemotron 2B is a transformer-based decoder-only language model similar to GPT-2 and 3 [21]. This model was trained on 1.1T tokens using NeMo [22] and Megatron-LM [18]. The model uses the SwiGLU activation function [23], Rotary Positional Embeddings (RoPE) [24], maximum sequence length of 4096, no dropout, no bias terms in all linear layers, and untied embedding and output layer. Nemotron-4 15B is a 15-billion-parameter large multilingual language model trained on 8 trillion text tokens. Both these models use a vocab size of 256K.

Data: Upcycling can be performed on either pretraining data that the pretrained dense model has seen, new unseen data, or a combination of both. In our Nemotron 2B experiments, we upcycle on pretraining data the model has seen for simplicity. For all of the ablation studies, we use 110B tokens (about 10% of the pretraining 1.1T tokens). For Nemotron-4 15B experiments, we upcycle on continued training data so that we can compare upcycling against our existing dense continued training result [13, 25]. While the continued training data has still been seen in pretraining, it follows a different data blending distribution. For 15B ablations, we train on 0.1T tokens (10%) of continued training data blend. For our 15B final results, we train on the full continued training data blend of 1T tokens. Validation loss is measured on 1% held-out data.

3.1 Effectiveness of Upcycling

Upcycling vs. Dense Continued Training: Following previous works [15], we compare upcycling vs continued training the dense Nemotron 2B model with the same amount of tokens (0.1T) under a similar learning schedule. As shown in Figure 4a, continuous training plateaus quickly while upcycling keeps improving. From continued training to upcycling, LM loss improved by 1.1%.

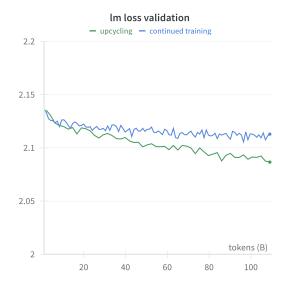
Upcycling vs. Training from Scratch: Figure 4b shows that upcycling achieves good improvement over training from scratch if one assumes a fixed compute budget. Upcycling is an efficient method to utilize pretrained dense model weights when the compute budget is much smaller than the pretraining compute budget. An interesting question that remains unanswered with our studies is whether it is still worth to upcycle a dense model instead of pretraining, assuming a larger compute budget. While some recent works like Skywork-MoE [26] try to answer this question, we leave investigating this as a potentially interesting future direction.

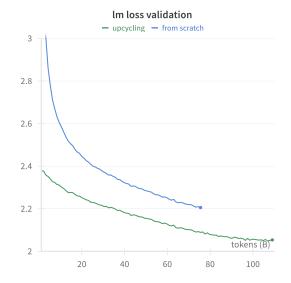
3.2 Learning Rate and Batch Size

3.2.1 Learning Rate Schedule

We found that the learning rate schedule plays an important role in upcycling. When upcycling a dense model, the model has usually been trained for a large number of steps already and the learning

 $^{^6 \}rm https://huggingface.co/nvidia/GPT-2B-001$





- (a) Upcycling Nemotron 2B outperforms continued training on loss
- (b) Upcycling Nemotron 2B outperforms training from scratch on loss under certain compute budget

Figure 4: Effectiveness of upcycling

rate is generally decayed during this training phase allowing the model to enter a local minimum. For example, our 2B model follows a warmup cosine decay learning rate schedule. The learning rate is warmed up to a 2e-4 peak and then decayed to 2e-5 at the end of 1.1T token training. During upcycling, it is unclear if increasing the learning rate above where pretraining ended (we call this resetting the learning rate) will improve or hurt model quality.

To find a good learning rate schedule for upcycling, we experimented with three different settings:

- Constant learning rate. We use the minimum pretraining learning rate of 2e-5 which is a typical learning rate schedule for finetuning.
- Peak learning rate 2e-4, cosine decay to 2e-5. This is the exact same learning rate schedule as pretraining, except that it decays to 2e-5 at the end of 0.1T tokens upcycling, as we only have 10% of pretraining tokens.
- Peak learning rate 1e-4, cosine decay to 2e-5. Since it's possible that using learning rate as high as pretraining can lead to catastrophic forgetting, using a lower peak learning rate might be a good option.

As shown in Figure 5, we found that while constant learning rate schedule starts off with much lower loss than the reset learning rate schedule, it eventually plateaus. The reset learning schedule gradually catches up and eventually outperforms the constant learning rate schedule. When resetting, peak learning rates of 1e-4 and 2e-4 seem to perform similarly. Interestingly, using a peak learning rate as high as pretraining (2e-4) does not lead to catastrophic forgetting.

Weights Similarity: Typically, finetuned model weights have very high cosine similarity with the base model weights. For example, Llama 2 chat model has cosine similarity close to 1 with Llama 2.

We compute cosine similarity between an upcycled model and the original dense model layer-by-layer and since MoE layer does not match with MLP, we compute cosine similarity between each expert and the original MLP and take the average. We finally average the cosine similarity across all layers to get a single number.

Shown in Figure 6, like in most finetuning tasks, upcycling with the minimum learning rate does not change the weights much. The cosine similarity between the MoE and base model is close to 1. By applying the reset learning schedule method, we observed the cosine similarity reduced to 0.6-0.7. This might imply that the higher learning rate helps the model escape from the dense model's local minimum, and find a superior minimum. Additionally, the high learning rate also helps the experts

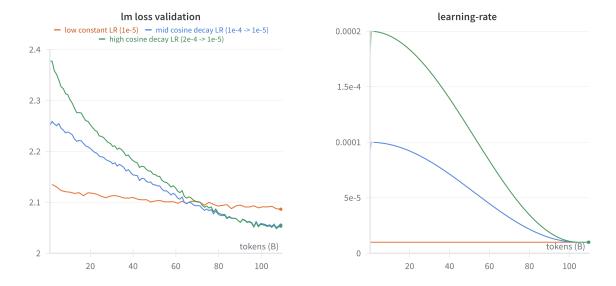


Figure 5: Resetting LR to peak pre-training LR when upcycling Nemotron 2B improves accuracy.

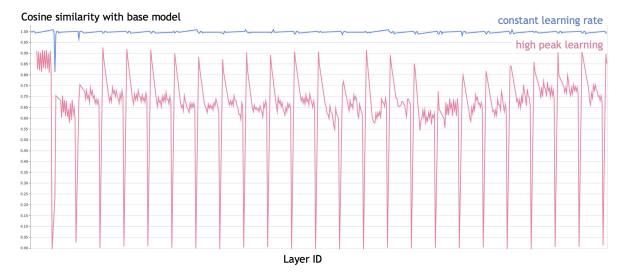


Figure 6: Cosine similarity between upcycled MoE and base Nemotron 2B model weights. When higher learning rate is used, similarity is lower and upcycling validation loss improves.

diversify. On the other hand, small constant learning rate leads to experts being very similar to each other, which makes the model not too different from the base dense model.

3.2.2 Batch Size

Aside from the learning rate, we observed that batch size also heavily affects MoE training and upcycling. We hypothesize that MoEs benefit from larger batch size than dense equivalents for two reasons:

- Since each expert receives only a portion of tokens, gradients are noisier than the dense model.
- The load balancing loss is noisier if there are fewer tokens to balance.

As shown in Figure 7, we compared batch size of 512, 1024, and 8192 (2M, 4M, and 32M tokens respectively) for upcycling the Nemotron 2B model. While the batch size of 32M tokens performs the worst, batch size of 4M tokens converges faster than 2M tokens. The training efficiency (per-GPU

throughput) is also much better with larger batch sizes. Recently, Deepseek-V2 [8] also used a large batch size of 9216 (more than 37M tokens) and learning rate of 2.4e-4 to pretrain MoE models.

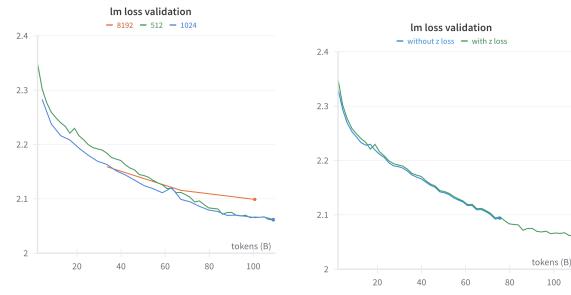


Figure 7: Effect of batch size on validation loss when upcycling Nemotron 2B. Increasing batch size to 1024 (4M tokens) does not degrade accuracy while improving model FLOP utilization.

Figure 8: Upcycling Nemotron 2B with or without z loss does not have noticeable difference

100

3.3 Load Balancing and Regularization Loss

Load Balancing Auxiliary Loss: We used the same aux loss as described in ST-MoE [5] and Switch Transformer [19] and experimented with different aux loss coefficients. We found that while not using aux loss at all leads to dead experts and causes the training loss to plateau early, aux loss coefficients set too high leads to aux loss overwhelming the language modeling loss. We found aux loss coefficients between 1e-2 to 1e-3 to give the best language model loss.

Z Loss: We used the same z loss described in ST-MoE [5]. As shown in Figure 8, we compared upcycling with and without z loss (with a coefficient of 1e-3) and found that z loss has no impact on the final model quality.

Thus, we used an aux loss coefficient of 1e-2 and no z loss for all our experiments.

3.4 Softmax TopK order

We consistently found that using softmax-then-topK works better than topK-then-softmax for upcycling. We hypothesize this is because the softmax-then-topK approach preserves the information contained in the absolute value of the router output. However, keeping the output of the original upcycled model similar to the output of the dense model is more difficult with this approach because the outputs no longer sum to one. We overcome this issue with our weight scaling method.

3.5 Fixing Output Scale

We tried multiple approaches to compensate for the issue of expert outputs being scaled down. Apart from weight scaling described in Equation 2, we also experimented with the following:

Scaling the MoE output: instead of scaling the weights, we tried scaling the output of the MoE layer by a constant factor or a learnable scalar. We empirically found that neither of them work well. Post Expert Layernorm: Work on finegrained MoE scaling laws [20] recommended adding a layernorm at the end of MoE layer. Typically, dense models do not have this layernorm. We tried adding the post expert layernorm during upcycling and found that while it can achieve the same effect and stop the loss from exploding, it takes a lot longer to adapt.

We compared different methods for upcycling Nemotron 2B into 64 experts top-8 fine-grained MoE (E8G8T8). The expert intermediate hidden size is 1/8 of the original FFN intermediate hidden size so that they are iso-FLOP. Our proposed weight scaling method performed the best and is also the easiest to implement, since it does not modify the model architecture.

Using equation 2, for the finegrained MoE with 64 experts top-8, the scaling factor should be 4. We tried multiple scaling factors (2x, 2.5x, 4x) and found that the scaling factor of 4 also performed the best empirically. So our weight scaling function, while not exact, helps convergence.

We also discovered that weight scaling helps upcycling standard coarse-grained MoEs as well. As shown in Figure 9, we upcycled Nemotron-4 15B into 8 experts top-1. With weight scaling, the training loss is 1.5% better than when not using weight scaling.

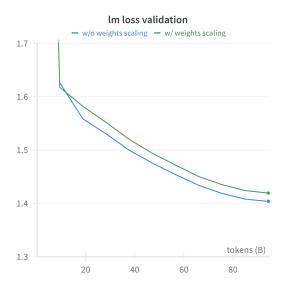


Figure 9: Using weight scaling helps achieve better loss when upcycling Nemotron-4 15B into 8 experts top-1 MoE

3.6 Increasing Granularity

We found that increasing granularity improves loss when training with small token counts (0.1T tokens for our 15B example) since virtual grouped granular upcycling is able to achieve a better loss more quickly than the non-granular version. However, when training on much larger token regimes (\geq 1T tokens), we saw that granularity did not help proportionally and both granular and non-granular runs converged to similar loss values. Since these large token horizon runs required significant compute, we did not perform ablations on them.

We tried scaling up the number of experts from 8 to 256 without increasing FLOPs. On upcycling Nemotron 2B and Nemotron-4 15B, we compared 8, 64, 128, and 256 experts. We kept all these networks iso-FLOP by scaling down the expert hidden size proportionally with respect to the topK. Shown in Figure 10, on Nemotron 2B, 64 experts performed better than 8 experts. However, scaling up further to 128 or 256 brought in little benefit. Similarly, on Nemotron-4 15B, the improvement maxed out at 64 experts. Surprisingly, on upcycling Nemotron-4 15B, 256 experts performed slightly worse than 64 or 128 experts. Too many experts seemed to hurt accuracy when upcycling. We hypothesize that this is because the experts were all copies and the larger the number of experts get, it becomes difficult for the network to find new superior minima. While increasing granularity is promising, it is important to note that it comes with higher MoE permute/unpermute cost and smaller GEMM sizes. Owing to these system-level factors, we worked with both granular and non-granular recipes.

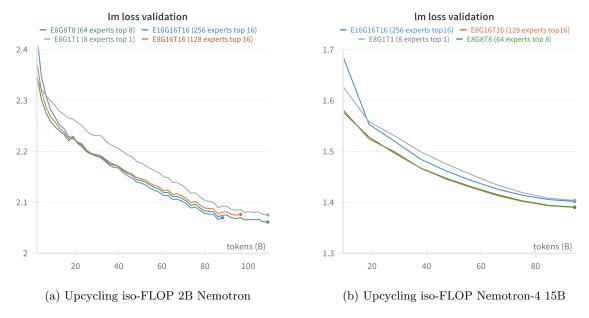


Figure 10: Increasing granularity provides diminishing benefits. We found granularity of 8 to be good for both Nemotron 2B and Nemotron-4 15B.

3.7 Increasing TopK

Top-2 routing is often used with MoE models [3, 5]. While this increases the amount of compute required to run the model, it helps in achieving better accuracy. We compared increasing topK for both coarse-grained and fine-grained MoE upcycling. Figure 11 shows that 8 experts top-2 (E8G1T2) performed better than top-1 (E8G1T1) on both upcycling Nemotron 2B and Nemotron-4 15B. On 15B, the top-2 achieved lower training loss than top-1 (1.35757 vs 1.38031). Previous works [5] have shown that a tradeoff with wall clock time rather than compute is a better metric and in such cases topK greater than granularity might make more sense. This claim heavily depends on the actual implementation and we leave it as a future systems study.

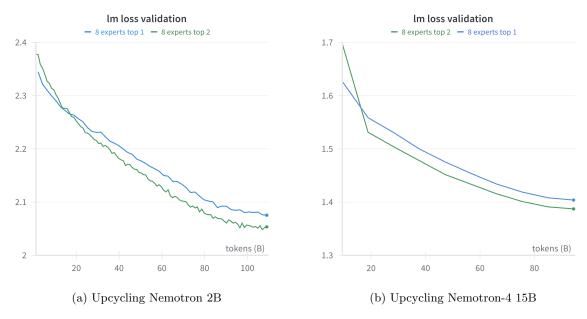


Figure 11: Increasing topK from 1 to 2 improves accuracy but requires extra compute

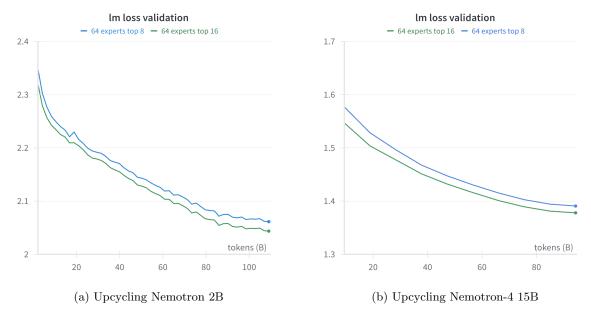


Figure 12: E8G8T8 vs. E8G8T16 (64 experts top-8 vs top-16, 1/8 expert intermediate size). Increasing TopK also helps granular models.

3.8 Promoting Expert Diversity: Weight Permutation and Reinitialization

We experimented with the weight permutation and reinitialization methods proposed in Qwen 2 [17]. Weight permutation permutes the FFN weights before copying it into each expert. Weights reinitialization randomly reinitializes 50% of expert weights. In our experiments, we did not find any improvement with these two methods on upcycling Nemotron 2B. Due to limited compute, we did not experiment with these techniques on a larger network.

3.9 Shared Experts

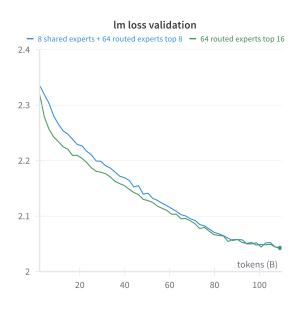


Figure 13: MoE with shared expert achieves similar accuracy as iso-FLOP fine-grained MoE in upcycling Nemotron 2B

We experimented with the shared experts approach proposed in Deepseek-MoE [16]. Shared experts are always 'on' i.e. every token is routed to the shared expert. They act like a dense layer in parallel with the MoE experts. As shown in Figure 13, we compared 8 shared experts + 64 experts top-8 which is iso-FLOP with 64 experts top-16 on Nemotron 2B. We found that shared expert performed on par with the iso-FLOP, no shared expert counterpart. We did not switch to using a shared expert since we did not see any accuracy improvement by using it.

4 Large Scale Upcycling

To ensure our proposed recipes work on large scale (model size and token count) regimes, we compared upcycling against continued training Nemotron-4 15B base model on 1T tokens. We decided to test only 2 variants: (1) E8G8T8 and (2) E8G1T2, due to the high training compute requirement (0.3 yottaFLOPs). We chose E8G8T8 with virtual group initialization since it worked better than E8G1T1 in our ablations. We chose E8G1T2 to show the effect of increase in compute and upcycling FLOPs. We chose E8G1T2 over E8G8T16 to show that our recommendations, including weight scaling along with softmax-then-topK router, work well for non-granular cases well. Non-granular use-cases are important to study since they practically achieve a better GPU FLOP utilization than their iso-FLOP granular counterparts.

For 1T upcycling, we initialized lr to min-lr of pretraining (4.5e-4) and used a peak-lr of 3e-4. We used cosine decay and decayed the lr to 1/100-th of the pretraining min-lr as done for the original continued training of Nemotron-4 15B [25]. We also used a higher batch size for the E8G8T8 model than E8G1T2 since E8G1T2 receives more tokens per expert on average. As shown in Table 1, upcycled E8G8T8 (64 experts top-8 1/8 expert hidden size) achieved a 4.1% lower validation loss than the dense continued training model while being iso-FLOP. It also achieved a better MMLU score of 66.2 (vs 65.3 for dense). We observed that the percentage difference in MMLU is sensitive to the continuous training data and the difference increased along with the token count - so longer token horizons favor MoE models. With increased training FLOPs, we also upcycled E8G1T2 (8 experts top-2) which achieved 5.2% lower validation loss and an even better MMLU of 67.6.

Model	val loss	MMLU (5 shots)
Nemotron-4 15B [13]	1.623	59.3
Nemotron-4 15B continued training [25]	1.377	65.3
Nemotron-4 15B upcycling E8G8T8	1.320	66.2
Nemotron-4 15B upcycling E8G1T2	1.306	67.6

Table 1: Upcycling Nemotron-4 15B on 1T tokens

5 Related Work

5.1 Mixture of Experts (MoE) Models

Mixture of Experts (MoE) models [27, 1] have gained significant attention in the field of large language models due to their ability to scale model capacity while maintaining computational efficiency. The MoE architecture employs a gating mechanism to selectively activate a subset of expert networks for each input token. This approach allows for an increased model capacity without a proportional increase in computational cost during inference and training.

Recent work has focused on improving the scalability and efficiency of MoE models. Switch Transformer [19] simplified the MoE architecture by using a top-1 routing mechanism and demonstrated the ability to scale to trillion-parameter models. The GShard framework [2] addressed challenges in training large-scale MoE models, introducing techniques such as expert capacity thresholds and local group dispatching to improve load balancing and training stability.

Research has explored various aspects of expert specialization and routing mechanisms in MoE models [28, 29]. Studies have investigated the impact of the number of experts on model performance, finding that increasing the number of experts leads to improved sample efficiency and faster training, albeit with diminishing returns beyond certain thresholds [20]. The choice of routing algorithm (e.g.,

top-1 vs top-2) and gating function (e.g., softmax vs sigmoid) has also been examined [30, 31, 32], with Mixtral 8x7B switching softmax and topK in the router [3].

5.2 Upcycling and Model Expansion

The concept of upcycling in the context of MoE models refers to the practice of leveraging pre-trained dense models to initialize MoE architectures. This approach has gained traction as a way to efficiently create large-scale MoE models while benefiting from the knowledge captured in existing pre-trained checkpoints. Notable work in this area includes:

Sparse Upcycling: [15] proposed a method for training MoE models from dense checkpoints, demonstrating the ability to expand model capacity while maintaining or improving performance. Recently, Qwen [17] and Deepseek [14] have adopted this approach. However, the recipe to scale upcycling beyond 1B parameters is not well known.

Network Growth: Research on model expansion techniques, such as those explored in the Gopher model [33], has shown that it's possible to significantly increase model size while maintaining performance comparable to or better than models trained from scratch.

Progressive Expansion: Approaches like LLAMA PRO [10] have investigated progressive expansion techniques, where model size is increased gradually during training [34, 35, 36, 37].

5.3 Challenges and Ongoing Research

Despite the promising results in upcycling and MoE pretraining, several challenges remain active areas of research:

Expert Collapse: The phenomenon of expert collapse, where certain experts become underutilized or inactive, has been observed in MoE training. While some studies suggest that expert collapse may not necessarily harm model accuracy, addressing this issue remains an important consideration in MoE design [38, 31].

Load Balancing: Ensuring an even distribution of work across experts continues to be a challenge, with various approaches proposed to improve load balancing, including auxiliary losses and specialized routing mechanisms. Skywork-MoE [26] used aux loss to promote expert diversity during upcycling. BTX [39, 40] trained experts on different tasks and then mixed them. DBRX ⁷ added weight norm on experts' weights.

Efficient Training and Inference: Ongoing work focuses on optimizing the training and inference processes for MoE models, including techniques for reducing communication costs and improving parallelization strategies. Deepseek-MoE [16] and Snowflake Arctic ⁸ parallelized shared expert to increase utilization. MegaBlocks [41] reformulated MoE computation in terms of block-sparse operations. Scattermoe [42] optimized memory for finegrained MoE.

6 Conclusion

In this work, we conducted an extensive study of upcycling techniques and best practices for billion-parameter scale language models. We proposed a "virtual group" initialization scheme and weight scaling approach to successfully enable both coarse-grained and fine-grained MoE upcycling at scale. We found that upcycling outperforms continued dense model training for the same amount of compute on both 2B and 15B parameter models. Based on the target inference and available upcycling FLOPs, an architecture that uses more FLOPs like E8G1T2 can achieve better accuracy than dense iso-FLOP MoE models. On our 2B ablations, we found that upcycling needs a different set of hyper-parameters than fine-tuning. Softmax-then-topK expert routing performs better than topK-then-softmax in the MoE router during upcycling. Higher granularity MoEs boost upcycling accuracy but require a more careful weight scaling and sharding strategy and also lead to a lower GPU FLOP utilization. In a purely FLOP-bound scenario, using virtual group init with granular model upcycling seems to be the best strategy. We hope this work illuminates the details of upcycling billion-parameter scale MoE models. Future directions include studying upcycling in larger models, improving expert diversity and utilization, and co-optimizing model architectures and system designs.

⁷https://www.databricks.com/blog/introducing-dbrx-new-state-art-open-llm

⁸ https://www.snowflake.com/blog/arctic-open-efficient-foundation-language-models-snowflake/

7 Acknowledgement

We would like to thank Aditya Vavre, Deepak Narayanan, Dennis Liu, Helen Ngo, Jakub Krajewski, June Yang, Jupinder Parmar, Marcin Chochowski, Mostofa Patwary, Ray Wang, Xin Yao, and Zeeshan Patel for their helpful contributions and discussions.

References

- [1] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. arXiv preprint arXiv:1701.06538, 2017.
- [2] Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional computation and automatic sharding. arXiv preprint arXiv:2006.16668, 2020.
- [3] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. arXiv preprint arXiv:2401.04088, 2024.
- [4] William Fedus, Jeff Dean, and Barret Zoph. A review of sparse expert models in deep learning. arXiv preprint arXiv:2209.01667, 2022.
- [5] Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and William Fedus. St-moe: Designing stable and transferable sparse expert models. arXiv preprint arXiv:2202.08906, 2022.
- [6] Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew M Dai, Quoc V Le, James Laudon, et al. Mixture-of-experts with expert choice routing. Advances in Neural Information Processing Systems, 35:7103-7114, 2022.
- [7] Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. Glam: Efficient scaling of language models with mixture-of-experts. In *International Conference on Machine Learning*, pages 5547–5569. PMLR, 2022.
- [8] DeepSeek-AI et al. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model, 2024.
- [9] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. arXiv preprint arXiv:2407.10671, 2024.
- [10] Chengyue Wu, Yukang Gan, Yixiao Ge, Zeyu Lu, Jiahao Wang, Ye Feng, Ping Luo, and Ying Shan. Llama pro: Progressive llama with block expansion. arXiv preprint arXiv:2401.02415, 2024.
- [11] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024.
- [12] Bo Adler, Niket Agarwal, Ashwath Aithal, Dong H Anh, Pallab Bhattacharya, Annika Brundyn, Jared Casper, Bryan Catanzaro, Sharon Clay, Jonathan Cohen, et al. Nemotron-4 340b technical report. arXiv preprint arXiv:2406.11704, 2024.
- [13] Jupinder Parmar, Shrimai Prabhumoye, Joseph Jennings, Mostofa Patwary, Sandeep Subramanian, Dan Su, Chen Zhu, Deepak Narayanan, Aastha Jhunjhunwala, Ayush Dattagupta, et al. Nemotron-4 15b technical report. arXiv preprint arXiv:2402.16819, 2024.
- [14] Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with longtermism. arXiv preprint arXiv:2401.02954, 2024.

- [15] Aran Komatsuzaki, Joan Puigcerver, James Lee-Thorp, Carlos Riquelme Ruiz, Basil Mustafa, Joshua Ainslie, Yi Tay, Mostafa Dehghani, and Neil Houlsby. Sparse upcycling: Training mixtureof-experts from dense checkpoints. arXiv preprint arXiv:2212.05055, 2022.
- [16] Damai Dai, Chengqi Deng, Chenggang Zhao, RX Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai Yu, Y Wu, et al. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models. arXiv preprint arXiv:2401.06066, 2024.
- [17] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. arXiv preprint arXiv:2309.16609, 2023.
- [18] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. arXiv preprint arXiv:1909.08053, 2019.
- [19] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022.
- [20] Jakub Krajewski, Jan Ludziejewski, Kamil Adamczewski, Maciej Pióro, Michał Krutul, Szymon Antoniak, Kamil Ciebiera, Krystian Król, Tomasz Odrzygóźdź, Piotr Sankowski, et al. Scaling laws for fine-grained mixture of experts. arXiv preprint arXiv:2402.07871, 2024.
- [21] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- [22] Oleksii Kuchaiev, Jason Li, Huyen Nguyen, Oleksii Hrinchuk, Ryan Leary, Boris Ginsburg, Samuel Kriman, Stanislav Beliaev, Vitaly Lavrukhin, Jack Cook, et al. Nemo: a toolkit for building ai applications using neural modules. arXiv preprint arXiv:1909.09577, 2019.
- [23] Noam Shazeer. Glu variants improve transformer, 2020.
- [24] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- [25] Jupinder Parmar, Sanjev Satheesh, Mostofa Patwary, Mohammad Shoeybi, and Bryan Catanzaro. Reuse, don't retrain: A recipe for continued pretraining of language models, 2024.
- [26] Tianwen Wei, Bo Zhu, Liang Zhao, Cheng Cheng, Biye Li, Weiwei Lü, Peng Cheng, Jianhao Zhang, Xiaoyu Zhang, Liang Zeng, et al. Skywork-moe: A deep dive into training techniques for mixture-of-experts language models. arXiv preprint arXiv:2406.06563, 2024.
- [27] Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991.
- [28] Fuzhao Xue, Zian Zheng, Yao Fu, Jinjie Ni, Zangwei Zheng, Wangchunshu Zhou, and Yang You. Openmoe: An early effort on open mixture-of-experts language models. arXiv preprint arXiv:2402.01739, 2024.
- [29] Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi, Pete Walsh, Oyvind Tafjord, Nathan Lambert, et al. Olmoe: Open mixture-of-experts language models. arXiv preprint arXiv:2409.02060, 2024.
- [30] Róbert Csordás, Kazuki Irie, and Jürgen Schmidhuber. Approximating two-layer feedforward networks for efficient transformers. arXiv preprint arXiv:2310.10837, 2023.
- [31] Zewen Chi, Li Dong, Shaohan Huang, Damai Dai, Shuming Ma, Barun Patra, Saksham Singhal, Payal Bajaj, Xia Song, Xian-Ling Mao, et al. On the representation collapse of sparse mixture of experts. *Advances in Neural Information Processing Systems*, 35:34600–34613, 2022.
- [32] Huy Nguyen, Nhat Ho, and Alessandro Rinaldo. Sigmoid gating is more sample efficient than softmax gating in mixture of experts. arXiv preprint arXiv:2405.13997, 2024.

- [33] Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models: Methods, analysis & insights from training gopher. arXiv preprint arXiv:2112.11446, 2021.
- [34] Yu Pan, Ye Yuan, Yichun Yin, Zenglin Xu, Lifeng Shang, Xin Jiang, and Qun Liu. Reusing pretrained models by multi-linear operators for efficient training. *Advances in Neural Information Processing Systems*, 36:3248–3262, 2023.
- [35] Andrea Gesmundo and Kaitlin Maile. Composable function-preserving expansions for transformer architectures. arXiv preprint arXiv:2308.06103, 2023.
- [36] Abhishek Panigrahi, Nikunj Saunshi, Kaifeng Lyu, Sobhan Miryoosefi, Sashank Reddi, Satyen Kale, and Sanjiv Kumar. Efficient stagewise pretraining via progressive subnetworks. arXiv preprint arXiv:2402.05913, 2024.
- [37] Peihao Wang, Rameswar Panda, Lucas Torroba Hennigen, Philip Greengard, Leonid Karlinsky, Rogerio Feris, David Daniel Cox, Zhangyang Wang, and Yoon Kim. Learning to grow pretrained models for efficient transformer training. arXiv preprint arXiv:2303.00980, 2023.
- [38] Dongyang Fan, Bettina Messmer, and Martin Jaggi. Towards an empirical understanding of moe design choices. arXiv preprint arXiv:2402.13089, 2024.
- [39] Sainbayar Sukhbaatar, Olga Golovneva, Vasu Sharma, Hu Xu, Xi Victoria Lin, Baptiste Rozière, Jacob Kahn, Daniel Li, Wen-tau Yih, Jason Weston, et al. Branch-train-mix: Mixing expert llms into a mixture-of-experts llm. arXiv preprint arXiv:2403.07816, 2024.
- [40] Qizhen Zhang, Nikolas Gritsch, Dwaraknath Gnaneshwar, Simon Guo, David Cairuz, Bharat Venkitesh, Jakob Foerster, Phil Blunsom, Sebastian Ruder, Ahmet Ustun, et al. Bam! just like that: Simple and efficient parameter upcycling for mixture of experts. arXiv preprint arXiv:2408.08274, 2024.
- [41] Trevor Gale, Deepak Narayanan, Cliff Young, and Matei Zaharia. Megablocks: Efficient sparse training with mixture-of-experts. *Proceedings of Machine Learning and Systems*, 5:288–304, 2023.
- [42] Shawn Tan, Yikang Shen, Rameswar Panda, and Aaron Courville. Scattered mixture-of-experts implementation. arXiv preprint arXiv:2403.08245, 2024.

A Virtual Grouped Router

The following pseudo code illustrates our virtual grouped router approach for upcycling into E2G2T2:

```
# virtual group - initialize every group with same weights router_prob = tensor([0.3, 0.3, 0.2, 0.2]) router_top2 = tensor([0.3, 0.3, 0.0, 0.0]) # one of each FFN shard is guaranteed to be selected FFN_moe(x) = 0.3 \ FFN_0 + 0.3 \ FFN_1 \approx FFN(x) / 4
```