# NEXUS: SPECIALIZATION MEETS ADAPTABILITY FOR EFFICIENTLY TRAINING MIXTURE OF EXPERTS

Anonymous authors

Paper under double-blind review

### ABSTRACT

Efficiency, specialization, and adaptability to new data distributions are qualities that are hard to combine in current Large Language Models. The Mixture of Experts (MoE) architecture has been the focus of significant research because its inherent conditional computation enables such desirable properties. In this work, we focus on "upcycling" dense expert models into an MoE, aiming to improve specialization while also adding the ability to adapt to new tasks easily. We introduce Nexus, an enhanced MoE architecture with *adaptive routing* where the model learns to project expert embeddings from domain representations. This approach allows Nexus to flexibly add new experts after the initial upcycling through separately trained dense models, without requiring large-scale MoE training for unseen data domains. Our experiments show that Nexus achieves a relative gain of up to 2.1% over the baseline for initial upcycling, and a 18.8% relative gain for extending the MoE with a new expert by using limited finetuning data. This flexibility of Nexus is crucial to enable an open-source ecosystem where every user continuously assembles their own MoE-mix according to their needs.

025 026 027

004

010 011

012

013

014

015

016

017

018

019

021

### 1 INTRODUCTION

In an era of bigger and bigger models (Canziani et al., 2016; Strubell et al., 2019; Rae et al., 2021; Raffel et al., 2020; Bommasani et al., 2022; Hooker, 2024), there are several key objectives driving state-of-art progress. Doing *more with less* by improving efficiency (Treviso et al., 2023) remains paramount, but in addition to efficiency, the deployment of these models in the wild means that the ability to adapt to new data (Pozzobon et al., 2023b; Gururangan et al., 2020a; Jang et al., 2022; Jin et al., 2022), and specialization of compute (Zadouri et al., 2024; Shazeer et al., 2018; Riquelme et al., 2021; Du et al., 2022; Fedus et al., 2022) have gained renewed focus. While all these properties are desirable, a formidable challenge is designing architectures that can fulfill *all* of these requirements.

The Mixture of Experts (MoE) approach gained prominence because of its efficiency properties. In contrast to dense models which require significant compute to deploy, MoE approaches only activate a subset of the parameters for every single token. Intuitively, not all parameters are necessary for each request, as some parameters will specialize on certain tasks, and those unrelated to the current request can be ignored. However, while MoEs greatly improved efficiency, the ability to induce meaningful specialization has been more limited, with observations that experts don't appear to exhibit dedicated expertise (Jiang et al., 2024; Zoph et al., 2022; Zadouri et al., 2023). Furthermore, MoEs tend to suffer from severe training instabilities (Zoph et al., 2022).

Recent work has attempted to address both the training instabilities and the lack of specialization.
These techniques often train completely separate experts and "upcycle" (combine) them into a single
unified MoE model *after* dense training (Sukhbaatar et al., 2024). This reduces the memory and
communication cost, and improves *efficiency* during training as computations are more local and
cross-device communication is reduced (Li et al., 2022; Gururangan et al., 2023). Notably, the other
major advantage of these approaches is the increase in *specialization* with separate experts that are
trained on specific domains, making them clearly responsible for their human-interpretable subset of
the data. On the other hand, MoEs with a standard router, which needs to be trained on a mix of all
training data, are not designed to maintain domain specialization (Jiang et al., 2024).

053 However, efficiently integrating new experts into upcycled MoE models - a setting that is of great interest for *adaptability* objectives - is far less studied. For most practitioners, given the scale of



Figure 1: Depiction of Nexus for a single Transformer block: A) In the initial training phase, each expert is trained separately. Its training data is embedded by an embedding model and stored. The experts are combined by initializing each block's MoE layer with the expert FFNs, and finetuning the model on a mix of all domains. During a forward pass, the seed model FFN is used as shared expert and always activated. For the other experts, we perform top-1 routing based on the similarity of the input data with the transformed expert embeddings, which is equivalent to viewing the learned projection as a hypernetwork whose output is the router weight matrix.
B) Later, we can add a new expert by appending its training data embedding to the existing domain embeddings. The router function is independent of the number of experts, and therefore adapts fast to the new one.

078

100

102

103

modern LLMs (Brown et al., 2020; Touvron et al., 2023; Kaplan et al., 2020; Anil et al., 2023)
training MoEs repeatedly is an infeasible computational cost. Furthermore, most model development
fails to take into account distribution drift in use cases, with limited flexibility and applicability across
different tasks and domains (Pozzobon et al., 2023a; Gururangan et al., 2020b). However, human
language is shaped by a cumulative culture, constantly building upon itself and evolving over time
(Silvey, 2016). Also, specialized use cases such as multilingual, code and math often require tailored
additional training.

In this work, we attempt to reconcile all three desirable properties: *efficiency, specialization, and adaptability.* We ask "*how can we adaptively combine separately trained specialized experts?*" To address this, we introduce **Nexus**, a novel MoE architecture that parameterizes the router based on domain-specific data by learning to project the embedding of each data domain to an expert embedding. This learnable projection for the router allows for the easy extension of the MoE model with new experts that are trained independently on new datasets of interest. This also avoids the difficulties of MoE training, as our learned router scales with the number of experts without needing to be trained from scratch, which enables adding or removing experts as desired.

Our experiments show that Nexus outperforms previous work when upscaling an MoE from separately trained specialized domain experts. Going beyond the single upscaling phase, Nexus can be efficiently extended with a new expert trained on a new domain, by finetuning it with much fewer tokens, compared to the finetuning after the initial upcycling.

- In summary, our contributions are as follows:
  - 1. We present Nexus, a novel MoE framework designed to enhance sparse upcycling of specialized dense experts, while reducing the training cost of MoEs by facilitating easy adaptation to unseen data distributions. In Nexus, the traditional linear router from vanilla MoE models is replaced with routing based on the similarity of layer inputs to an expert embedding vector, derived from the average embedding of the corresponding expert dataset.
- 2. Our method outperforms the existing approach for upcycling specialized models into MoE, leading to 2.1% and 1.6% relative increase over the upcycled MoE (linear router) in 470M and 2.8B scales respectively. This enables performance increase in general tasks with 5.8% and 7.4% relative gains over the dense seed model at 470M and 2.8B respectively.

- 3. Our method enables efficient adaptation to new domains by extending upcycled MoE with the new experts trained on unseen datasets. In this setting, Nexus outperforms the baseline MoE (linear router) when finetuning on the limited amount of data, leading 18.8% relative gain on the new domain with 1B finetuning tokens upon MoE extension.
  - 4. Finally, we show that our method is robust across different load balancing and data mixtures, and consistently outperforms the MoE with a linear router for specialized upcycling, confirming the benefits of the *adaptive* routing based on domain projections used in Nexus.

### 2 BACKGROUND

119 Sparse Mixture of Experts architectures (Shazeer et al., 2017; Fedus et al., 2022) replace the feed-120 forward network (FFN) with an MoE layer in the Transformer block (Vaswani et al., 2017). An 121 MoE layer consists of a router network R and a set of n experts,  $E_1, ..., E_n$ , where each expert  $E_i$ 122 corresponds to an independent dense feed-forward network. The router network R is commonly 123 parameterized by trainable weights  $W_r \in \mathbb{R}^{h \times n}$  where h is the model hidden dimension, and followed by a *softmax* function which takes an intermediate token representation x as input and combines the 124 output of each expert based on the gating scores  $s_1, ..., s_n$ . Sparse MoEs only use the top-k experts 125  $E_k$  based on experts gating scores  $s_i$ . 126

127

129 130 131

132 133

108

110

111

112

113

114

115 116 117

118

128

 $s_i = R(x) = \operatorname{softmax}(W_r^T x)$ (Router)

$$s_k = \text{TopK}(s_i)$$
 (Top-K Routing)

$$y = \sum_{i=1}^{k} s_k \cdot E_k(x) \tag{MoE}$$

134 Sparse Upcycling (Komatsuzaki et al., 2023) initializes an MoE model from a dense Transformer 135 model by copying FFN layers as MoE experts, and the router layer is trained from scratch. BTX 136 (Sukhbaatar et al., 2024) generalize this approach to initialize each MoE expert from the FFN layer 137 of a different dense model, and all other parameters are averaged over the dense models.

138 In **Nexus**, we leverage upcycling specialized expert models similar to BTX, however, it diverges in 139 terms of MoE training, in particular with its novel MoE router, which enables to efficiently extend 140 the MoE in multiple rounds after the sparse upcycling. We describe our method in the next section.

- 141 142
- 143 144

### ADAPTIVE ROUTER FOR UPCYCLING SPECIALIZED EXPERTS AS MOE 3

145 The core component of an MoE model is the router, as it determines which experts to activate 146 for any given input. In vanilla MoEs, the router is a learned linear layer that takes the token intermediate representations as input and computes the expert probabilities. However, this router does 147 not necessarily learn specialization as MoEs are commonly trained using an auxiliary load balancing 148 loss to improve training stability (Fedus et al., 2022; Jiang et al., 2024). In Nexus, we propose a 149 novel MoE router where per MoE block we learn a projection layer from given pre-computed domain 150 embeddings to expert embeddings. We parametrize this projection layer  $P_r$  as a two-layer MLP with 151 a SwiGLU activation function (Shazeer, 2020): 152

153 154

156

$$e_i = P_r(d_i)$$
 (Domain to Expert Embeddings)  
=  $W_2 \cdot \text{SwiGLU}(W_1 \cdot d_i)$ 

where  $d_i \in \mathbb{R}^m$ , and  $e_i \in \mathbb{R}^h$  are the domain and expert embeddings for the *i*th domain respectively. 157 where m and h are the domain embedding and the model dimensions.  $W_1 \in \mathbb{R}^{2h \times d}, W_2 \in \mathbb{R}^{l \times l}$ 158 are linear layers, and SwiGLU is defined as  $\mathbb{R}^{2n} \to \mathbb{R}^n$ . Given the expert embeddings  $e_i$  and layer 159 inputs  $x \in \mathbb{R}^{s \times h}$ , we then compute routing probabilities  $s_i$  as: 160

161

$$s_i = \operatorname{softmax}(x \cdot e_i)$$
 (Routing Scores)

Unlike the standard router, Nexus's router includes a stronger inductive bias through pre-computed domain embeddings<sup>1</sup> that enables expert embedding to specialize. Thus,  $x \cdot e_i$  gives a high value for input tokens that are closer to the domain of the corresponding expert. Notably, this router is particularly suited for the sparse upcycling setting where the dense experts are separately trained on different domains.

167 168 169 169 169 169 169 169 170 170 171 Connection to hypernetworks. Our router parametrization is closely related to hypernetworks (Ha 170 et al., 2016) as the projection layer  $P_r$  generates parameters for the router during runtime for a given 171 in provide the projection layer, enabling efficient adaptation 172 and also a better cross-domain transfer based on the similarity between domain embeddings as shown 173 in previous work (Mahabadi et al., 2021; Üstün et al., 2022).

Upcycling dense experts as an MoE. After training dense expert models, we merge the individual experts into a unified MoE by appending their FFNs along a new dimension to create an MoE layer per Transformer block. Unlike Sukhbaatar et al. (2024), instead of using the original FFN of the seed model as one of the routed experts in an MoE layer, we use it as the "shared expert" FFNs (Rajbhandari et al., 2022; Dai et al., 2024) to better preserve the previous capabilities in the MoE model. For all non-FFN parameters including the attention weights, we merge expert parameters using simple weight averaging:

> $FFN_{moe} = FFN_s + [FFNe_1, FFNe_2, ..., FFNe_n]$ (MoE Layer FFNs)  $\phi_{moe} = \frac{\sum_{i=1}^{n} \phi_i}{n}$ (Merge Non-FFN params.)

**Efficient adaptation to new domains.** An important advantage of method is that when a new data domain is present after MoE training, we use the learned projection  $P_r$  to compute expert embedding of the new domain as  $e_{new} = P_r(d_{new})$ . This enables to enhance the trained MoE model with additional dense experts, which are trained in the same way as the initial experts. The FFN parameters of the new expert are simply appended to the array of existing experts.

To adequately preserve the non-FFN parameters of existing experts, we perform a weighted average  $\phi_f = (1 - \lambda) \cdot \phi_{moe} + \lambda \cdot \phi_{new}$  where  $\phi_f$ ,  $\phi_e$ , and  $\phi_{moe}$  are parameters of the final MoE, dense expert, and initial MoE model and  $\lambda = 1/(n + 1)$ . This enables *efficiently adapting* Nexus to new domain by extending it with the new dense expert trained independently. After extending the MoE with a new expert, we perform a lightweight finetuning with a limited number of tokens.

194 195

196 197

215

179

180

181 182 183

4 EXPERIMENTS

197 4.1 EXPERIMENTAL SETTING 198

199 Our experimental setup includes 3 phases. Figure 1 shows the architecture of Nexus and the corresponding experimental setting:

201 1. Training specialized expert LMs. For training the dense specialized experts, we use the 202 sub-datasets from the SlimPajama dataset (Soboleva et al., 2023), a 627B token English-language corpus assembled from web data of various sources. We initialize four dense experts from the 203 204 weights of the seed model and train them on the ARXIV, BOOKS, C4, GITHUB, STACKEXCHANGE, and WIKIPEDIA domains.<sup>2</sup> As the seed model, we use 470M and 2.8B parameters decoder-only 205 autoregressive Transformer models (Radford et al., 2019), each of them trained with a standard 206 language modeling objective for 750B tokens. We train dense experts for 20 and 40 billion tokens for 207 470M and 2.8B seed models respectively. We use parallel attention layers, (Anil et al., 2023; Wang, 208 2021), SwiGLU activation (Shazeer, 2020), no biases in dense layers, and a byte-pair-encoding (BPE) 209 tokenizer with a vocabulary size of 256,000. During training, we use a linear warmup (10% of total 210 steps) to a maximum learning rate of 1e-3 and a cosine decay schedule to 3e-4. 211

 <sup>&</sup>lt;sup>1</sup>We used Cohere Embed v3 (Cohere, 2023) as an external embedding model to compute domain embeddings based on individual data sources. However, similar to Gururangan et al. (2023), pre-training data can also be clustered and the centroids can be used for domain embeddings.

<sup>&</sup>lt;sup>2</sup>We exclude the Github and StackExchange datasets from SlimPajama in order to ablate adding a new expert model using the CODE domain



Figure 2: **Downstream performance at different scales:** Nexus consistently outperforms upcycled baselines on both the 470M and 2.8B parameters scale, showing the robustness of our method. We report the average performance on Knowledge, Science, Reasoning and MMLU.

216

217 218

219

220

222

224

225 226

227 228

229

**2. MoE training.** After the training of dense expert models, we merge them into a unified MoE by 233 appending their FFNs along a new dimension to create an MoE layer per Transformer block. For the 234 shared expert in our MoE layer, we use the original FFN layer of the seed model to better preserve 235 the previous capabilities in the MoE model. For all non-FFN parameters including the attention 236 weights, we merge expert parameters using simple weight averaging, following Sukhbaatar et al. 237 (2024). After the MoE model is created, we continually train it for an additional 25B and 40B tokens 238 respectively for the 470M and 2.8B experiments, on a mix of all domain and original pre-training 239 datasets, using the same training hyperparameters as in the single expert training. Finally, we train the MoE models using an additional 1B tokens by upweighting the original pre-training dataset as 240 it includes high-quality data sources such as instruction-style datasets using a cosine learning rate 241 decay to 3e-5 (Parmar et al., 2024). 242

3. Extending the MoE model with new experts. After adding a new expert as defined in Section
3, we finetune the extended MoE model for up to 1 billion tokens using a uniformly sampled data
mix consisting of 50% the previous domains and pre-training data and 50% the new domain. For the
new expert (CODE), we train a dense model using code documents from StarCoder (Li et al., 2023)
with the same settings as for the training of the initial experts. As the 470M scale MoE did not have
sufficient instruction following capabilities to attempt the code benchmarks, we only tested extending
the MoEs with a new expert on the 2.8B scale.

250 251

253

262

- 4.2 BASELINES
- We compare our experiments against two baselines:

Dense Merging where all separately pre-trained experts and the seed model merged into a dense
 Transformer via equal weight averaging similar to BTM (Li et al., 2022). This allows us to ask *What are the benefits of routing MoE over simple averaging*?

MoE (Linear Router) which is an MoE with a standard linear router that is upcycled from dense
 experts, to evaluate Nexus's novel router for upcycling. Here, we ask *how does our specialized routing compare to conventional learned linear routing?* For a fair comparison, we also train this
 MoE model on the same datasets and for the same number of tokens as our method, and use the same
 architectural modifications such as shared experts.

263 4.3 EVALUATION

For the downstream evaluation, we measure the performance of each model on 15 tasks from five evaluation categories that reflect different capabilities based on the tasks and the datasets used in the benchmarks.

These task categories are (1) Knowledge, to measure question-answering capabilities based on world knowledge and web documents such as Wikipedia, we report the performance on OpenBookQA (Mihaylov et al., 2018), Natural Questions (Kwiatkowski et al., 2019), TriviaQA (Joshi et al., 2017),

	Know.	Science	Reason.	MMLU	Code (excl. in upcyc.)	Avg. (w/o Code)
SEED MODEL (2.8B)	27.1	62.0	63.8	35.4	8.4	47.1
Upcycled Models						
Dense Merging	17.6	60.3	59.2	36.0	3.4	43.3
MOE (LINEAR ROUTER)	31.5	66.5	62.9	38.6	2.6	49.8
NEXUS	33.2	67.3	62.6	39.4	2.7	50.6

Table 1: Downstream task results for Nexus with a 2.8B parameter seed model. Our approach outperforms the baselines in 3 out of 4 evaluation categories. Dense merging corresponds a dense
model with 2.8B parameters, while both Nexus and MoE (linear router) have 4.3B active
and 9.1B total parameters. Note that the trained models show severe forgetting on code benchmarks,
as we exclude CODE data on purpose during the upcycling phase to simulate extending models with a
new dataset in Section 5.2.

285 QUAC (Choi et al., 2018) (all 0-shot) and SQuAD (4-shot) (Rajpurkar et al., 2016). (2) Science, to 286 measure knowledge in science-oriented academic benchmarks, we use ARC-Easy, ARC-Challenge 287 (Clark et al., 2018), SciQ (Welbl et al., 2017) (all 0-shot). (3) Reasoning, we use CommonSenseQA 288 (Talmor et al., 2019), SIQA (Sap et al., 2019), PIQA (Bisk et al., 2020), WinoGrande (Sakaguchi et al., 289 2019), and HellaSwag (Zellers et al., 2019) (all 0-shot). (4) General Language Understanding, 290 we use MMLU (5-shot) (Hendrycks et al., 2021). (5) Code, we evaluate models on MBPP (Austin 291 et al., 2021), LBPP (Matton et al., 2024) and HumanEval-Pack (Chen et al., 2021) that includes Cpp, 292 Javascript, Java, Go, Python, and Rust (all 0-shot). 293

293 294 295

296 297

298

284

### 5 RESULTS AND DISCUSSION

5.1 MAIN RESULTS FOR UPCYCLED MODELS

We first compare Nexus to the upcycled baselines MoE with linear router and dense merging. Here, we ask "*How does our MoE upcycling recipe with adaptive routing compare against baseline upcycling approaches?*"

470M parameter seed model. Table 4 (Appendix D) shows performances of upcycled models
 including Nexus where a 470M seed model is used to train dense experts. Both Nexus and the
 upcycled MoE (linear router) consist of 1 shared and 6 routed experts, corresponding to a
 total number of 1.3B parameters where 605M parameters are activated per input for top-2 routing
 (1 expert always activated, 1 chosen by the router). The dense merging baseline is created by
 averaging the weights of all dense experts and the seed model, and therefore has the same number of
 parameters as the seed model.

Compared to the seed model, Nexus performs better in all evaluation categories with a 5.8% relative gain on average (38.5 vs 36.4). Compared to upcycled models, Nexus outperforms MoE (linear router) in 3 out of 4 categories with 3.2% relative gain (38.5 vs 37.3) on average, and beats dense merging by 8.5% overall relative increase (38.5 vs 35.5). Notably, while both upcycled MoEs outperform the seed model, dense merging underperforms on average, showing the benefits of MoE upcycling over parameter averaging.

2.8B parameter seed model. Next, we experiment by upcycling dense models with 2.7B parameters to validate if the results from the 470M seed model hold at a larger scale. Table 1 compares Nexus with MoE (linear router) and dense merging. Both Nexus and MoE (linear router) use 1 shared expert and 4 routed experts in these experiments, corresponding to 4.3B active parameters per input (top-2) out of 9.1B total parameters.

Our results show that Nexus leads to higher upcycling results compared to the baselines at the 2.8B
 scale, confirming the findings from smaller scale experiments. Nexus enables a 7.4% relative gain
 over the seed model and outperforms the MoE (linear router) with a 1.6% relative increase
 (50.6 vs. 49.8). Nexus outperforms the best baseline in 3 out of 4 task categories and achieves
 the highest increase in *knowledge* tasks with 22.5% and 5.6% relative to the seed model and the

325 326

327

328

330

331

332

333

334

335

336



Figure 3: Extending upcycled MoE models with the Code experts: After initial upcycling, we extended MoEs (both Nexus and MoE with linear router) using an independently trained dense Code expert and finetuned the resulting models small number of tokens (200M, 500M, and 1B finetuning tokens) as described in 3. Nexus consistently outperforms the baseline in Code performance after extension without losing general performance. General tasks is the macro average of the knowledge, science, reasoning, and general knowledge categories reported in section 5.1. Note that the dense Code expert achieves scores of 42.1 and 14.3 for general and code tasks respectively.

MoE (linear router) respectively. These tasks include knowledge retrieval from Wikipedia in
 which one of our specialized experts is trained for.

Similar to the 470M experiments, both Nexus and MoE (linear router) outperform the
dense merging baseline. We relate this to potential cross-task interference between diverse
specialized experts (including the seed model as an additional expert), leading to poor performance
by applying a simple weight averaging.

352<br/>3535.2EXTENDING THE UPCYCLED MOE MODEL WITH A NEW EXPERT

354 To support fully modular and efficient training of MoEs, besides upcycling the existing expert models, 355 it is crucial for an adaptive method to have the ability to continuously extend the upcycled MoE with 356 new experts trained using previously unseen data domains. To evaluate this, we train a dense CODE 357 expert and extend the upcycled MoEs (both Nexus and MoE (linear router)) as described 358 in Section 3. We perform a small-scale finetuning of up to 1B tokens after extending the models. Figure 3 shows both the general performance and the target code performance at 200M, 500M, and 359 1B finetuning tokens. Here, we ask "Can we continuously upcycle dense models into an MoE without 360 requiring large-scale MoE training each time?" 361

Performance on the new domain. As shown in Figure 3 (right), Nexus outperforms the MoE (linear router) for 200M, 500M and 1B finetuning tokens with 18.4%, 6.2% and 18.8% relative gains respectively. Unlike MoE (linear router), where the router weights are reset after extending the MoE layers, Nexus uses the information that is available about the new domain by mapping the domain embedding to a new expert embedding for the router, and therefore finetunes the router weights without a restart.

368 **Comparison with the dense models.** Nexus reaches the code performance of the seed model while 369 retaining superior performance on general tasks. In comparison to the seed model and the dense code expert (trained for 8B code-only tokens on top of the seed model), although the dense code expert 370 still performs higher than both upcycled MoEs with a score of 14.3, its performance on general tasks 371 is far inferior (42.1). Our method also achieves up to 18.8% relative gains over the MOE (linear 372 router). These results show that with a fraction of the original upcycling budget (1B vs 40B 373 tokens for initial upcycling, and 1B vs 8B tokens for code expert training), Nexus can acquire a new 374 capability. 375

Performance on general tasks. As a proxy for the knowledge for previously learned domains,
 Figure 3 (left) shows the average performance of Nexus and MoE (linear router) in general tasks. Although there is a slight drop on the general tasks for Nexus compared to initial upcycling

Router decision ArXiv Books C4 Wiki 0.7 0.6 0.5 0.4 0.3 0.2 0.1 0.0 ArXiv Books C4 Wiki

Figure 4: Average routing probabilities for each expert per domain in Nexus: We compute the average routing probabilities across Transformer blocks for 512 samples per domain (from the 2.8B experiment). The x-axis denotes the samples' domain and the colored bars show the routing probabilities for the corresponding expert. We show the domains that are used to train specialized experts.

(a relative decrease of 1.9%), the competitive performance is maintained across different numbers of finetuning tokens. We relate this to the composition of the finetuning mix where we use a high percentage of the code data (50% of the code and 50% of the previous domains).

### 395 396 397 398

378

379

380

381

382

384

385

386

387

388

389

390

391 392

393

394

5.3 EXPERT SPECIALIZATION

To measure the specialization in our MoE, we take a closer look at how the MoE experts are activated 399 for samples of separate domains. We compute average routing frequencies across all Transformer 400 layers in Figure 4, where the labels on the x-axis represent which domain the tokens are coming 401 from, and the colored bars show the routing frequencies for each of the experts trained on one 402 of the domains. Since we select only one routed expert per token in each MoE layer, and expert 403 FFN layers are inherited from dense experts, average routing frequencies present a good proxy for 404 specialization of each of the experts. Here, we ask "can Nexus retain a high degree of specialization 405 after upcycling?" 406

**Routing for the upcycled experts.** As shown in Figure 4, we find that the expert trained on the cor-407 responding domain always receives the highest share of the tokens from that domain, confirming that 408 Nexus retains the specialization from the specialized dense models. Concretely, this specialization 409 is higher for ArXiv, Books, and Wikipedia with 63.0%, 64.7%, and 69.8% respectively. Interestingly, 410 tokens from C4 are routed only 40.9% of the time to the C4 expert and distributed to the other experts 411 approximately 20% for each one. We relate this to the broad coverage of the C4 dataset, which 412 potentially includes samples closer to other domains and also a large percentage of the C4 used in the 413 MoE training phase (proportional to its size in the SlimPiama dataset). Especially the latter factor 414 pushes tokens from C4 to be distributed to the other experts due to the load balancing factor.

Specialized routing for the new expert. Next, we measure expert specialization for the newly added expert on the new code domain. Figure 5 shows the average routing probability per expert for sampled code tokens. We compute routing probabilities on the Nexus model with the code expert after 1B finetuning tokens (See Section 5.2 for details). Here, we see clearly that code tokens are routed to the code expert 69.1% of the time on average. This shows that Nexus not only retains the specialization for the initial upcycling but also exhibits a high degree of specialization for a newly added expert for its own domain.

422 423

424

5.4 Ablations

Mixture-of-expert models are known to be sensitive to the choice of load balancing loss factor (Fedus et al., 2022; Zoph et al., 2022) and sampling weights for each data domains during training. As additional ablations, we run two new sets of experiments at 470M scale, one with a lower load balancing factor and the other one with equal weighting of each domain during training (whereas originally the weights were proportional to the share of tokens of that domain in SlimPajama). Figure 6 compares Nexus and MoE (linear router) in terms of their downstream performances for these ablations. Finally, in this section, we also visualize domain and projected expert embeddings to see if the relationship between embeddings is preserved after the learned projection.





Figure 5: Average routing probabilities per expert for the new domain in extended Nexus: We show the routing probabilities for code tokens after extending MoE (1B finetuning).



Lowering the load balancing loss factor. In Figure 6 (baseline vs low load-bal.), we compare two Nexus models with the corresponding MoE (linear router) baselines where we use load balancing loss factor of 0.05 and 0.0005 for each set of experiments. We find that using a significantly lower factor for the load balancing loss hurts MoE (linear router) performance by approximately 2% relative drop while Nexus shows a robust performance across both load balancing factors. We hypothesize that because the expert embeddings in our router are always based on the domain representations, we achieve more stable distribution of tokens even if the load balancing loss is weighted extremely low.

459 Changing the training data composition. Next, we compare our default of sampling specialized 460 domain data proportional to the size of the domain (total amount of tokens in SlimPajama), with 461 a uniform sampling over all domains. Figure 6 (baseline vs equal data) shows the downstream 462 performances for both Nexus and MOE (linear router). Although sampling uniform sampling 463 domains' data does not significantly impact the downstream performance for both models, we find that 464 it helps Nexus to improve specialization for all the domains in terms of expert routing probabilities 465 (Figure 11, Appendix I). In particular, compared to the size proportional sampling, tokens from the C4 domain are routed more accurately (27.6% vs 71.1%) when data is equally sampled. 466

467 **Domain embeddings before and after projection.** Finally, in Figure 8, we visualize cosine 468 similarities between domains and the projected expert embeddings from the last Transformer block, 469 in our main upcycling experiments at the 470M scale. Comparing the embeddings before and after 470 mapping, we find that the router's learned projection preserves the main relationship between domains. 471 For instance, relatively high cosine similarity between Books & C4, and StackExchange & GitHub exist both between their domain embeddings and the projected expert embeddings. Interestingly, 472 while preserving the main relationships, we also find that the learned projection pushes expert 473 embeddings further away from each other, potentially due to our choice of only activating a single 474 expert per token besides the shared expert. 475

476 477

478

445

446

447

448

449 450 451

### 6 RELATED WORK

Routing Variants of MoEs. The most common MoE architecture (Shazeer et al., 2017; Lepikhin et al., 2020; Fedus et al., 2022) employs a linear router with a top-k routing scheme, where k typically equals 1 or 2. In this standard routing schema, only the k experts with the highest router gate values are activated. There is substantial research proposing alternatives to top-k expert assignments (Hazimeh et al., 2021; Lewis et al., 2021; Roller et al., 2021; Zhou et al., 2022; Zuo et al., 2022). DeepSeek-MoE (Dai et al., 2024) introduces a routing variant where a number of experts are "shared"

<sup>485</sup> 

<sup>&</sup>lt;sup>3</sup>We report the average performance on Knowledge, Science, Reasoning, and MMLU.

and always assigned to all tokens. Our work also adopts this approach for our general base expert.
 However, these efforts primarily focus on improving the general performance and/or training stability of MoEs. In contrast, our work puts emphasis adaptability and extensibility.

489 Efficient MoE Training by Re-Using Existing Dense Models. Training MoEs from scratch is 490 computationally expensive (Gale et al., 2023; Fedus et al., 2022) and often challenging due to training 491 instabilities (Zoph et al., 2022). Alternatively, recent works have explored re-using existing dense 492 models to initialize MoEs. Sparse Upcycling (Komatsuzaki et al., 2023) re-uses a single dense model 493 to initialize the MoE by replicating the FFN weights in an MoE layer. The router is initialized 494 randomly, and all other parameters are copied directly from the dense model. BTX (Sukhbaatar 495 et al., 2024) extends this approach by upcycling not from a single dense model, but from multiple 496 specialized dense expert models. Furthermore, BAM (Zhang et al., 2024) expands BTX to upcycle not only FFN experts but also attention experts. Our work also leverages this approach by reusing 497 specialized dense experts for an MoE, while extending it further to facilitate on-the-fly adaptations 498 for new experts specialized in unseen data domains. 499

500 Efficient MoE Architectures. Zadouri et al. (2024) proposes replacing traditional MoE's 501 computation-heavy feed-forward network (FFN) experts with more efficient experts comprised of smaller vectors and adapters, which are activated in parallel to a single dense FFN. This lightweight 502 503 architecture necessitates only a limited number of parameter updates when finetuning, offering efficiency advantages. However, unlike our approach, it does not leverage existing specialized dense 504 models and lacks a notion of specialized experts, which are central to our method. Similar to our 505 work, Muqeeth et al. (2024) and Ostapenko et al. (2024) study combining separately trained experts 506 into a unified model. However, they focus on parameter-efficient adapters such as LoRA (Hu et al., 507 2021) and supervised finetuning. In this work, we focus on efficiently pre-training fully-fledged MoE 508 models via upcycling. 509

Adaptive MoEs and Ensemble Models. ModuleFormer (Shen et al., 2023) also aims to produce 510 adaptable MoEs. The authors achieve adaptability by freezing existing MoE parameters while 511 only training newly added modules with optimization constraints to the router. Unlike our work, 512 ModuleFormer does not leverage existing expert dense seed models for efficiency gains, nor does it 513 have a notion of specialization which is central to our work. Similar to our work, DEMix (Gururangan 514 et al., 2021) independently trains different FFN experts on specialized data domains, with each expert 515 functioning as a domain-specific module. Modules can be added on-the-fly for adaptability. Followup 516 works BTM and C-BTM (Li et al., 2022; Gururangan et al., 2023) extend DEMix to create adaptive 517 ensemble models. However, all three works use a router requiring a forward pass for every expert at 518 inference instead of sparsely activating them, which significantly increases inference costs, especially 519 with a large number of experts. Unlike these approaches, our router cost is approximately the same 520 as standard top-k routing during both training and inference, offering a more scalable solution for adaptability. 521

522 523

524

# 7 CONCLUSION

We propose Nexus, a new LLM framework that enables efficient upcycling of specialized dense experts into a sparsely activated MoE model. We show that individual experts in our method retain their specialization after upcycling, and that our router based on expert embeddings outperforms previous approaches for combining the dense experts. Furthermore, the model can be extended efficiently with new dense experts after the initial training phase, saving much compute compared to re-training the upcycled model or training from scratch.

531 532

533

# 8 LIMITATIONS

The MoE architecture is often employed for larger models in the multi-billion parameter range, where efficiency is paramount. However, to facilitate a broader set of experiments, we limit our setup to using 2.8B parameter seed models for the main results and 470M parameter seed models for ablations. Furthermore, our dense experts are based on existing data sources in the SlimPajama dataset which is pre-defined. Future work could extend our method by discovering specialized data domains through unsupervised clustering similar to Gururangan et al. (2023).

# 540 REFERENCES

542 Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, 543 Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark 544 Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan Botha, James Bradbury, 546 Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. 547 Choquette-Choo, Aakanksha Chowdhery, Clément Crepy, Shachi Dave, Mostafa Dehghani, Sunipa 548 Dev, Jacob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu Feng, Vlad 549 Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-Ari, 550 Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, Andrea Hu, Jeffrey Hui, Jeremy Hurwitz, 551 Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun, 552 Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li, Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu, Frederick Liu, Marcello Maggioni, 553 Aroma Mahendru, Joshua Maynez, Vedant Misra, Maysam Moussalem, Zachary Nado, John 554 Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie Pellat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan Richter, Parker Riley, Alex Castro Ros, Aurko Roy, 556 Brennan Saeta, Rajkumar Samuel, Renee Shelby, Ambrose Slone, Daniel Smilkov, David R. So, Daniel Sohn, Simon Tokumine, Dasha Valter, Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang, 558 Pidong Wang, Zirui Wang, Tao Wang, John Wieting, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting 559 Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven Zheng, Ce Zheng, Weikang Zhou, Denny 560 Zhou, Slav Petrov, and Yonghui Wu. Palm 2 technical report, 2023. 561

- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,
   Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large
   language models, 2021.
- 565 566

567

568 569 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020.

Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, 570 Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, 571 Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel, 572 Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano 573 Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren 574 Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter 575 Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil 576 Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar 577 Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal 578 Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Munyikwa, 579 Suraj Nair, Avanika Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles, 580 Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu Ren, Frieda 582 Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh, Shiori Sagawa, 583 Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori, Armin W. 584 Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, 585 Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi Zhang, 586 Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. On the opportunities and risks of foundation models, 2022. URL https://arxiv.org/abs/2108.07258. 588

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler,
Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott
Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya
Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.

594 595 596	Alfredo Canziani, Adam Paszke, and Eugenio Culurciello. An Analysis of Deep Neural Network Models for Practical Applications. <i>arXiv e-prints</i> , pp. arXiv:1605.07678, May 2016.
597 598 599 600 601 602 603 604 605 606	Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Ka- plan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021. URL https://arxiv.org/abs/2107.03374.
607 608 609	Eunsol Choi, He He, Mohit Iyyer, Mark Yatskar, Wen-tau Yih, Yejin Choi, Percy Liang, and Luke Zettlemoyer. Quac: Question answering in context. <i>arXiv preprint arXiv:1808.07036</i> , 2018.
610 611 612 613	<ul> <li>Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. <i>arXiv preprint arXiv:1803.05457</i>, 2018.</li> </ul>
614 615	introducing-embed-v3.
616 617 618 619	Damai Dai, Chengqi Deng, Chenggang Zhao, R. X. Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai Yu, Y. Wu, Zhenda Xie, Y. K. Li, Panpan Huang, Fuli Luo, Chong Ruan, Zhifang Sui, and Wenfeng Liang. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models, 2024. URL https://arxiv.org/abs/2401.06066.
620 621 622 623 624	Nan Du, Yanping Huang, Andrew M. Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, Barret Zoph, Liam Fedus, Maarten Bosma, Zongwei Zhou, Tao Wang, Yu Emma Wang, Kellie Webster, Marie Pellat, Kevin Robinson, Kathleen Meier-Hellstern, Toju Duke, Lucas Dixon, Kun Zhang, Quoc V Le, Yonghui Wu, Zhifeng Chen, and Claire Cui. Glam: Efficient scaling of language models with mixture-of-experts, 2022.
625 626	William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity, 2022.
627 628 629	Trevor Gale, Deepak Narayanan, Cliff Young, and Matei Zaharia. Megablocks: Efficient sparse training with mixture-of-experts. <i>Proceedings of Machine Learning and Systems</i> , 5:288–304, 2023.
630 631 632 633 634	Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. Don't stop pretraining: Adapt language models to domains and tasks. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pp. 8342–8360, Online, July 2020a. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.740. URL https://aclanthology.org/2020.acl-main.740.
635 636 637	Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. Don't stop pretraining: Adapt language models to domains and tasks. <i>arXiv</i> preprint arXiv:2004.10964, 2020b.
638 639 640	Suchin Gururangan, Mike Lewis, Ari Holtzman, Noah A Smith, and Luke Zettlemoyer. Demix layers: Disentangling domains for modular language modeling. <i>arXiv preprint arXiv:2108.05036</i> , 2021.
641 642 643	Suchin Gururangan, Margaret Li, Mike Lewis, Weijia Shi, Tim Althoff, Noah A. Smith, and Luke Zettlemoyer. Scaling expert language models with unsupervised domain discovery, 2023. URL https://arxiv.org/abs/2303.14177.
644 645	David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. arXiv preprint arXiv:1609.09106, 2016.
646 647	Hussein Hazimeh, Zhe Zhao, Aakanksha Chowdhery, Maheswaran Sathiamoorthy, Yihua Chen, Rahul Mazumder, Lichan Hong, and Ed H. Chi. Dselect-k: Differentiable selection in the mixture of experts with applications to multi-task learning, 2021.

648 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob 649 Steinhardt. Measuring massive multitask language understanding, 2021. URL https://arxiv. 650 org/abs/2009.03300. 651 Sara Hooker. On the limitations of compute thresholds as a governance strategy, 2024. URL 652 https://arxiv.org/abs/2407.05694. 653 654 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 655 and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. 656 Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun Kim, Stan-657 ley Jungkyu Choi, and Minjoon Seo. Towards continual knowledge learning of language models, 658 2022. URL https://arxiv.org/abs/2110.03215. 659 660 Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris 661 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, 662 Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le 663 Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 664 Mixtral of experts, 2024. URL https://arxiv.org/abs/2401.04088. 665 666 Xisen Jin, Bill Yuchen Lin, Mohammad Rostami, and Xiang Ren. Learn continually, generalize 667 rapidly: Lifelong knowledge accumulation for few-shot learning, 2022. URL https://arxiv. 668 org/abs/2104.08808. 669 Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. TriviaQA: A large scale distantly 670 supervised challenge dataset for reading comprehension. In Proceedings of the 55th Annual 671 Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 1601– 672 1611, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/ 673 v1/P17-1147. URL https://aclanthology.org/P17-1147. 674 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, 675 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models, 676 2020. 677 678 Aran Komatsuzaki, Joan Puigcerver, James Lee-Thorp, Carlos Riquelme Ruiz, Basil Mustafa, Joshua 679 Ainslie, Yi Tay, Mostafa Dehghani, and Neil Houlsby. Sparse upcycling: Training mixture-of-680 experts from dense checkpoints, 2023. 681 Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris 682 Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. 683 Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 684 Natural questions: a benchmark for question answering research. Transactions of the Association 685 of Computational Linguistics, 2019. 686 Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, 687 Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional 688 computation and automatic sharding, 2020. 689 690 Mike Lewis, Shruti Bhosale, Tim Dettmers, Naman Goyal, and Luke Zettlemoyer. Base layers: 691 Simplifying training of large, sparse models. In Proceedings of the 38th International Conference 692 on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pp. 6265–6274. 693 PMLR, 18-24 Jul 2021. URL https://proceedings.mlr.press/v139/lewis21a. 694 html. Margaret Li, Suchin Gururangan, Tim Dettmers, Mike Lewis, Tim Althoff, Noah A. Smith, and Luke 696 Zettlemoyer. Branch-train-merge: Embarrassingly parallel training of expert language models, 697 2022. URL https://arxiv.org/abs/2208.03306. 698 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, 699 Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue 700 Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro, 701 Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Logesh Kumar

702 Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy, Jason 703 Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, 704 Nour Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luccioni, Paulo Villegas, 705 Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor, Jennifer Ding, Claire 706 Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun 708 Guha, Leandro von Werra, and Harm de Vries. Starcoder: may the source be with you!, 2023. 710 Rabeeh Karimi Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. Parameter-711 efficient multi-task fine-tuning for transformers via shared hypernetworks, 2021. 712 Alexandre Matton, Tom Sherborne, Dennis Aumiller, Elena Tommasone, Milad Alizadeh, Jingyi He, 713 Raymond Ma, Maxime Voisin, Ellen Gilsenan-McMahon, and Matthias Gallé. On leakage of code 714 generation evaluation datasets. arXiv preprint arXiv:2407.07565, 2024. 715 716 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct 717 electricity? a new dataset for open book question answering. In Proceedings of the 2018 Conference 718 on Empirical Methods in Natural Language Processing, pp. 2381–2391, Brussels, Belgium, 719 October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1260. 720 URL https://aclanthology.org/D18-1260. 721 Mohammed Muqeeth, Haokun Liu, Yufan Liu, and Colin Raffel. Learning to route among specialized 722 experts for zero-shot generalization. arXiv preprint arXiv:2402.05859, 2024. 723 724 Oleksiy Ostapenko, Zhan Su, Edoardo Maria Ponti, Laurent Charlin, Nicolas Le Roux, Matheus Pereira, Lucas Caccia, and Alessandro Sordoni. Towards modular llms by building and reusing a 725 library of loras. arXiv preprint arXiv:2405.11157, 2024. 726 727 Jupinder Parmar, Shrimai Prabhumoye, Joseph Jennings, Mostofa Patwary, Sandeep Subramanian, 728 Dan Su, Chen Zhu, Deepak Narayanan, Aastha Jhunjhunwala, Ayush Dattagupta, Vibhu Jawa, Jiwei 729 Liu, Ameya Mahabaleshwarkar, Osvald Nitski, Annika Brundyn, James Maki, Miguel Martinez, 730 Jiaxuan You, John Kamalu, Patrick LeGresley, Denys Fridman, Jared Casper, Ashwath Aithal, 731 Oleksii Kuchaiev, Mohammad Shoeybi, Jonathan Cohen, and Bryan Catanzaro. Nemotron-4 15b 732 technical report, 2024. URL https://arxiv.org/abs/2402.16819. 733 Luiza Pozzobon, Beyza Ermis, Patrick Lewis, and Sara Hooker. Goodtriever: Adaptive toxi-734 city mitigation with retrieval-augmented models. In Findings of the Association for Com-735 putational Linguistics: EMNLP 2023, pp. 5108-5125, Singapore, December 2023a. Asso-736 ciation for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.339. URL 737 https://aclanthology.org/2023.findings-emnlp.339. Luiza Pozzobon, Beyza Ermis, Patrick Lewis, and Sara Hooker. Goodtriever: Adaptive toxicity 739 mitigation with retrieval-augmented models, 2023b. URL https://arxiv.org/abs/2310. 740 07589. 741 742 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language 743 models are unsupervised multitask learners, 2019. URL https://api.semanticscholar. 744 org/CorpusID:160025533. 745 Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John 746 Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, 747 Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, 748 Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron 749 Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, 750 Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen 751 Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, 752 Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, 753 Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d'Autume, 754 Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, 755 Aurelia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake Hechtman, Laura Weidinger,

792

796

797

798

- Iason Gabriel, William Isaac, Ed Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol
  Vinyals, Kareem Ayoub, Jeff Stanway, Lorrayne Bennett, Demis Hassabis, Koray Kavukcuoglu,
  and Geoffrey Irving. Scaling Language Models: Methods, Analysis & Insights from Training
  Gopher, 2021.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
   Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text
   transformer, 2020.
- Samyam Rajbhandari, Conglong Li, Zhewei Yao, Minjia Zhang, Reza Yazdani Aminabadi, Ammar Ahmad Awan, Jeff Rasley, and Yuxiong He. DeepSpeed-MoE: Advancing mixture-of-experts inference and training to power next-generation AI scale. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 18332–18346. PMLR, 17–23 Jul 2022. URL https://proceedings.mlr.press/v162/rajbhandari22a.html.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2383–2392, Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1264. URL https://aclanthology.org/D16-1264.
- Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André
   Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts.
   *Advances in Neural Information Processing Systems*, 34:8583–8595, 2021.
- 779
   780
   781
   781
   782
   784
   785
   786
   786
   787
   787
   788
   788
   789
   789
   780
   780
   781
   781
   781
   781
   781
   782
   782
   783
   784
   784
   785
   784
   785
   785
   786
   786
   787
   786
   786
   787
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
   788
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale, 2019.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialiqa: Commonsense reasoning about social interactions, 2019. URL https://arxiv.org/abs/1904.09728.
- 787
   788
   789
   Noam Shazeer. Glu variants improve transformer, 2020. URL https://arxiv.org/abs/ 2002.05202.
- 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, 2017.
- Noam Shazeer, Youlong Cheng, Niki Parmar, Dustin Tran, Ashish Vaswani, Penporn Koanantakool,
   Peter Hawkins, HyoukJoong Lee, Mingsheng Hong, Cliff Young, Ryan Sepassi, and Blake
   Hechtman. Mesh-tensorflow: Deep learning for supercomputers, 2018.
  - Yikang Shen, Zheyu Zhang, Tianyou Cao, Shawn Tan, Zhenfang Chen, and Chuang Gan. Moduleformer: Modularity emerges from mixture-of-experts. *arXiv e-prints*, pp. arXiv–2306, 2023.
- Catriona Silvey. Speaking our minds: Why human communication is different, and how languageevolved to make it special, by thom scott-phillips, 2016.
- Daria Soboleva, Faisal Al-Khateeb, Robert Myers, Jacob R Steeves, Joel Hestness, and Nolan Dey.
   SlimPajama: A 627B token cleaned and deduplicated version of RedPajama, June 2023. URL https://huggingface.co/datasets/cerebras/SlimPajama-627B.
- Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for deep learning in nlp, 2019. URL https://arxiv.org/abs/1906.02243.
- 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.

810 Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A question 811 answering challenge targeting commonsense knowledge. In Proceedings of the 2019 Conference of 812 the North American Chapter of the Association for Computational Linguistics: Human Language 813 Technologies, Volume 1 (Long and Short Papers), pp. 4149–4158, Minneapolis, Minnesota, June 814 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1421. URL https: //aclanthology.org/N19-1421. 815

- 816 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée 817 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand 818 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language 819 models, 2023. 820
- Marcos Treviso, Ji-Ung Lee, Tianchu Ji, Betty van Aken, Qingqing Cao, Manuel R. Ciosici, Michael 821 Hassid, Kenneth Heafield, Sara Hooker, Colin Raffel, Pedro H. Martins, André F. T. Martins, 822 Jessica Zosa Forde, Peter Milder, Edwin Simpson, Noam Slonim, Jesse Dodge, Emma Strubell, 823 Niranjan Balasubramanian, Leon Derczynski, Iryna Gurevych, and Roy Schwartz. Efficient 824 Methods for Natural Language Processing: A Survey. Transactions of the Association for Compu-825 tational Linguistics, 11:826-860, 07 2023. ISSN 2307-387X. doi: 10.1162/tacl\_a\_00577. URL 826 https://doi.org/10.1162/tacl\_a\_00577. 827
- Ahmet Üstün, Arianna Bisazza, Gosse Bouma, Gertjan van Noord, and Sebastian Ruder. Hyper-X: A 828 unified hypernetwork for multi-task multilingual transfer. In Proceedings of the 2022 Conference 829 on Empirical Methods in Natural Language Processing, pp. 7934–7949, Abu Dhabi, United Arab 830 Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022. 831 emnlp-main.541. URL https://aclanthology.org/2022.emnlp-main.541. 832
- 833 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz 834 Kaiser, and Illia Polosukhin. Attention is all you need, 2017.
- Ben Wang. Mesh-Transformer-JAX: Model-Parallel Implementation of Transformer Language 836 Model with JAX. https://github.com/kingoflolz/mesh-transformer-jax, 837 May 2021. 838
  - Johannes Welbl, Nelson F Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. arXiv preprint arXiv:1707.06209, 2017.
  - Ted Zadouri, Ahmet Üstün, Arash Ahmadian, Beyza Ermiş, Acyr Locatelli, and Sara Hooker. Pushing mixture of experts to the limit: Extremely parameter efficient moe for instruction tuning, 2023. URL https://arxiv.org/abs/2309.05444.
- Ted Zadouri, Ahmet Üstün, Arash Ahmadian, Beyza Ermis, Acyr Locatelli, and Sara Hooker. 845 Pushing mixture of experts to the limit: Extremely parameter efficient moe for instruction tuning. In The Twelfth International Conference on Learning Representations, 2024. URL https: //openreview.net/forum?id=EvDeiLv7qc. 848
- 849 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine 850 really finish your sentence?, 2019. 851
- Qizhen Zhang, Nikolas Gritsch, Dwaraknath Gnaneshwar, Simon Guo, David Cairuz, Bharat 852 Venkitesh, Jakob Foerster, Phil Blunsom, Sebastian Ruder, Ahmet Ustun, and Acyr Locatelli. 853 Bam! just like that: Simple and efficient parameter upcycling for mixture of experts, 2024. URL 854 https://arxiv.org/abs/2408.08274. 855
- 856 Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew Dai, Zhifeng Chen, Quoc Le, and James Laudon. Mixture-of-experts with expert choice routing, 2022.
- Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and 859 William Fedus. St-moe: Designing stable and transferable sparse expert models, 2022. 860
- 861 Simiao Zuo, Xiaodong Liu, Jian Jiao, Young Jin Kim, Hany Hassan, Ruofei Zhang, Tuo Zhao, and 862 Jianfeng Gao. Taming sparsely activated transformer with stochastic experts, 2022.

857

858

835

839

840

841

842

843

844

846

847

### NEXUS ROUTING ALGORITHM А

Figure 7 outlines the code for the Nexus router, which consists of (1) a 2-layer MLP network (domain\_to\_expert\_ffn) to project domain embeddings to expert embeddings, (2) shared and routed expert FFNs, and (3) sparse Top-k gating. Note that the expert embeddings are independent of the input and could be precomputed once and stored as long as the weights of the model do not change. This means that the routing layer during inference closely resembles a vanilla MoE router, with the difference being that the router matrix in Nexus is not learnt during training but computed using the domain embeddings as an informative prior.

```
875
876
       3
877
       4
```

```
1 def router(self, inputs, domain_embeddings):
           # domain_to_expert_ffn learns projection domain to expert embeddings
           # domain_embeddings: [e_dim x n_experts]
           # expert_embeddings: [h_dim x n_experts]
878
           expert_embeddings = self.domain_to_expert_ffn(self.domain_embeddings)
     5
879
     6
880
           # router probs: [batch, seq, n_experts]
     7
           router_probs = nn.softmax(inputs @ expert_embeddings)
     8
     9
882
    10
           # Top-1 gate for routed experts
883
           index, gate = nn.topk(1, router_probs)
884
    12
           # routed_experts_ffns: An MoE layer with FFN experts
885
           # routed_expert_out: [batch, seq, h_dim]
    14
           # shared_expert_out: [batch, seq, h_dim]
    15
887
    16
           routed_expert_out = self.routed_expert_ffns[index](input)
           shared_expert_out = self.shared_expert_ffn(input)
889
    18
890
    19
```

return shared\_expert\_out + gate \* routed\_expert\_out

Figure 7: Router layer in Nexus: PyTorch-like pseudo-code illustrating the routing mechanism, situated before the expertized MLP layer in each transformer block.

### В COMPARISON OF EXISTING APPROACHES WITH NEXUS

Table 2 compares Nexus with previous approaches in the field of efficient MoE training. Unlike the vanilla MoE architecture (Shazeer et al., 2017; Fedus et al., 2022), the Branch-Train-Merge (BTM; Li et al., 2022) and the Branch-Train-Mix (BTX; Sukhbaatar et al., 2024) approaches train experts separately in different domains, reducing training cost and improving specialization. However, they either merge the experts during inference (BTM) or learn an MoE router layer from scratch, where prior domain information is not used (BTX). Our approach trains the MoE router based on domain information, maintaining the specialization and enabling efficient extension of the MoE with a new expert after training.

	MOE	BTM	BTX	NEXUS
	(Vanilla)	(Merge)	(Linear router)	(Ours)
Dense experts are trained independently (upcycling)	×	✓	<b>v</b>	✓
Experts are specialized in different domains	×	<ul> <li>✓</li> </ul>	<ul> <li>✓</li> </ul>	<ul> <li>✓</li> </ul>
Experts are chosen by a learned router per input token	<ul> <li>✓</li> </ul>	×	<ul> <li>✓</li> </ul>	<ul> <li>✓</li> </ul>
Router is adaptive via learned projection for new domains	×	×	×	<ul> <li>✓</li> </ul>

Table 2: A comparison of existing approaches with Nexus: We choose the vanilla MoE architecture (Shazeer et al., 2017; Fedus et al., 2022), Branch-Train-Merge (BTM; Li et al., 2022), and Branch-Train-Mix (BTX; Sukhbaatar et al., 2024) for comparison. Nexus combines the advantages of the existing MoE extensions while also allowing easy adaptation to new domains.

Furthermore, Table 3 shows parameter counts during training and inference of Nexus vs. the baselines.
 From this, we can infer that Nexus has the same memory and compute complexity as a vanilla MoE model during inference, and a slight overhead of 1% additional trainable parameters during training.

	Total Parameters	Active Parameters (Training)	Active Parameters (Inference)
470M Models			
SEED MODEL (470M)	467,682,304	467,682,304	467,682,304
MOE (LINEAR ROUTING)	1,298,252,800	606,110,720	606,110,720
NEXUS	1,312,834,560	620,692,480	606,110,720
2.8B Models			
SEED MODEL	2,752,565,760	2,752,565,760	2,752,565,760
MOE (LINEAR ROUTING)	9,044,226,560	4,325,429,760	4,325,429,760
NEXUS	9,129,218,560	4,410,421,760	4,325,429,760

Table 3: Total and active parameter counts. Comparison of the seed model, linear MoE, and Nexus architectures for both 470M and 2.8B parameter models. During inference, the router weights of Nexus can be precomputed once by the learned MLP hypernetworks, making it exactly equal to the vanilla MoE in terms of memory and compute complexity. During training, we also observe exactly the same step time for the vanilla MoE and Nexus, as the overhead of the additional MLP is negligible. In the 470M category, the MoE/Nexus models use 6 routed and 1 shared expert. In the 2.8B category, the MoE/Nexus models use 4 routed and 1 shared expert. In both categories, the models activate the shared expert and the top-1 of the routed experts during inference.

### C EVALUATION DETAILS

For the downstream evaluation, we measure the performance of each model on 15 tasks<sup>4</sup> from five evaluation categories that reflect different capabilities based on the tasks and the datasets used in the benchmarks:

- **Knowledge**: To measure question-answering capabilities based on world knowledge and web documents such as Wikipedia, we report the performance on OpenBookQA (Mihaylov et al., 2018), Natural Questions (Kwiatkowski et al., 2019), TriviaQA (Joshi et al., 2017), QUAC (Choi et al., 2018) (all 0-shot) and SQuAD (4-shot) (Rajpurkar et al., 2016).
  - Science: For measuring knowledge in science-oriented academic benchmarks, we use ARC-Easy, ARC-Challenge (Clark et al., 2018), SciQ (Welbl et al., 2017) (all 0-shot).
  - **Reasoning**: For reasoning abilities, we use CommonSenseQA (Talmor et al., 2019), SIQA (Sap et al., 2019), PIQA (Bisk et al., 2020), WinoGrande (Sakaguchi et al., 2019), and HellaSwag (Zellers et al., 2019) (all 0-shot).
  - General Language Understanding: We use MMLU (5-shot) (Hendrycks et al., 2021) to test general language understanding.
  - Code: For code generation, we evaluate models on MBPP (Austin et al., 2021), LBPP (Matton et al., 2024), and HumanEval-Pack (Chen et al., 2021) which includes Cpp, Javascript, Java, Go, Python, and Rust (all 0-shot).

# D RESULTS FOR THE 470M PARAMETER MODEL

Table 4 shows the downstream task results for Nexus with a 470M parameter seed model. Our approach outperforms the baselines in all downstream benchmarks. Dense merging corresponds a dense model with 470M parameters, while both Nexus and MoE (linear router) consist of 605M active and 1.3B total parameters.

<sup>&</sup>lt;sup>4</sup>We did not include ARC-Challenge and Natural Questions in 470M experiments as some model variants were unable to achieve non-random performance.

0 - 0						
972		Know.	Science	Reason.	MMLU	Avg.
973		1110.00	Science	neusoni		
974	SEED MODEL (470M)	14.0	51.4	50.5	29.8	36.4
975	Upcycled Models					
976	DENSE MERGING	10.9	52.0	50.3	27.8	35.5
977	MOE (LINEAR ROUTER)	13.4	55.0	51.3	29.6	37.3
978	NEXUS	16.7	55.0	52.3	29.8	38.5

Table 4: Downstream task results for Nexus with a 470M parameter seed model. Dense merging merges all separately pretrained experts, while both Nexus and MoE (linear router) upcycle them and are evaluated with top-2 routing.

### 

### 

### Е **RESULTS FOR INDIVIDUAL EXPERTS**

To further contextualize the performance of the Nexus models, we report the performance of each individual expert in Table 5. The experts initialized from the 470M seed model are trained for 20B tokens on their domains, while the experts initialized from the 2.8B seed model are trained for 40B tokens.

993						
994		Know.	Science	Reason.	MMLU	Avg.
995	470M Experts					
996	ARXIV	9.5	47.8	44.3	31.2	33.2
997	BOOKS	9.0	51.8	51.4	32.0	36.1
000	C4	3.9	52.6	51.5	27.6	33.9
998	GITHUB	11.3	44.8	45.2	30.2	32.9
999	STACKEXCHANGE	9.9	45.4	44.9	29.2	32.4
1000	WIKIPEDIA	15.3	46.4	44.1	25.4	32.8
1001	2.8B Experts					
1002	ARXIV	13.4	57.3	51.3	36.2	39.5
1003	BOOKS	19.4	62.5	60.0	39.6	45.4
1004	C4	11.0	64.5	61.9	37.8	43.8
1005	WIKIPEDIA	22.6	60.3	55.3	37.2	43.9
1006	CODE	13.4	59.9	52.4	37.8	40.9
1007	Upcycled Models					
1008	NEXUS (470M)	16.7	55.0	52.3	29.8	38.5
1009	NEXUS (2.8B)	33.2	67.3	62.6	39.4	50.6

Table 5: Downstream task performance of individual experts. We report the separate performance of all experts used during the upcycling and extension stages. Note that the Nexus models beat every individual expert used for their upcycling, with one exception.

F **RESULTS FOR CONTINUAL TRAINING OF THE SEED MODEL** 

To compare Nexus to another dense baseline, for Table 6 we continually train the 470M and 2.8B seed models in a data matched setting. This means the 470M model has seen a total of 750B pretraining tokens (general pretraining data mix), 120B tokens from SlimPajama domains (the shuffled training tokens of all 6 experts), and 25B tokens from SlimPajama to match the Nexus finetuning phase. The 2.8B model has seen a total of 750B pretraining tokens, 160B tokens from SlimPajama (the shuffled training tokens of all 4 experts), and 40B additional tokens from SlimPajama to match the Nexus finetuning phase.

1026						
1020		Know.	Science	Reason.	MMLU	Avg.
1027						
1028	470M Models					
1000	SEED MODEL	14.0	51.4	50.5	29.8	36.4
1029	SEED MODEL + 145B TOKENS	19.9	53.8	50.8	29.6	38.5
1030	NEXUS	16.7	55.0	52.3	29.8	38.5
1031						
1032	2.8B Models					
1033	SEED MODEL	27.1	62.0	63.8	35.4	47.1
1000	SEED MODEL + 200B TOKENS	28.8	66.4	62.7	41.4	49.8
1034	NEXUS	33.2	67.3	62.6	39.4	50.6
1025						

Table 6: Downstream task results for data-matched continued pretraining of the 470M and 2.8B seed models. Both seed models are data matched to the Nexus/Linear MoE variants, including all expert training and finetuning. For the 2.8B parameter models, we also train for 1B tokens on instruction-style datasets from the original pretraining data before measuring the performance on downstream tasks (see Section 4.2). Note that the seed model training takes a lot more wallclock-time compared to our method, as the Nexus experts can all be trained in parallel, which is not possible with a single model.

# G COMPARISON OF DOMAIN EMBEDDINGS AND EXPERT EMBEDDINGS

Nexus maps the *domain embeddings* which are computed from each domain's training dataset to *expert embeddings* which represent experts. Figure 8 shows that trends in the similarity matrix are preserved after the mapping, but the embeddings are pushed away from each other a bit (lower similarity).



1063 1064 1065 <del>-</del>

1043 1044

1046

1051

1052

1053

1054

1055

1056

1057 1058

1061

1062

Figure 8: **Domain and the projected expert embeddings for Nexus:** We visualize cosine similarities between domains and the projected expert embeddings from the last Transformer block. The similarities are obtained from the 470M experiments. Our projected router maintains the relative similarity between the original domains (e.g. Books & C4, Github & StackExchange) after the router's projection.

1070 1071 1072

1073

### H ROUTING PROBABILITIES FOR THE LINEAR MOE MODEL

To investigate how specialized individual experts are in the Nexus approach vs. the vanilla MoE baseline, we also compute the routing distributions for the MoE (Linear Router) baseline with 2.8B parameters. Figure 9 shows the router distribution of this model with 4 experts. Figure 10 shows the router distribution for code data after adding the new code expert to the baseline. Although the specialization of the linear MoE model (Figure 9) matches that of Nexus for pretraining (Figure 4), it adapts much worse to the new expert, as fewer tokens from the code domain actually get routed to the code expert (Figure 10) than with Nexus (Figure 5).



Figure 9: Average routing probabilities for each expert per domain in the MoE (Linear router) baseline: We compute the average routing probabilities across Transformer blocks for 512 samples per domain (from the 2.8B experiment). The x-axis denotes the samples' domain and the colored bars show the routing probabilities for the corresponding expert. We show the domains that are used to train specialized experts.





### I ROUTING PROBABILITIES FOR UPCYCLING ABLATIONS

Figure 11 shows the expert routing probabilities for Nexus for all three settings described in Section 5.4.



