

# DEMIX Layers: Disentangling Domains for Modular Language Modeling

Anonymous ACL submission

## Abstract

We introduce a new domain expert mixture (DEMIX) layer that enables conditioning a language model (LM) on the domain of the input text. A DEMIX layer is a collection of expert feedforward networks, each specialized to a domain, that makes the LM *modular*: experts can be mixed, added, or removed after initial training. Extensive experiments with autoregressive transformer LMs (up to 1.3B parameters) show that DEMIX layers reduce perplexity, increase training efficiency, and enable rapid adaptation. Mixing experts during inference, using a parameter-free weighted ensemble, enables better generalization to heterogeneous or unseen domains. Adding experts incorporates new domains without forgetting older ones, and removing experts restricts access to unwanted domains without additional training. Overall, these results demonstrate benefits of explicitly conditioning on textual domains during language modeling.

## 1 Introduction

Conventional language model (LM) training algorithms assume data homogeneity: all parameters are updated to minimize the loss on all of the data. We refer to this approach as *dense training*. Dense training leaves variation in the data, or *domains*, to be implicitly discovered (Aharoni and Goldberg, 2020), assuming that models will be able to fit all domains equally well.

While dense training is convenient, and densely trained LMs achieve impressive results (Brown et al., 2020), the approach has drawbacks with respect to generalization, efficiency, and flexibility. Even if training data is sourced from many domains, dense training can in practice emphasize subsets of the data in proportion to their ease of access (Oren et al., 2019; Fan et al., 2020), limiting generalization to less prevalent domains. Updating all parameters of the network gets substantially more expensive as model size grows (Strubell et al.,

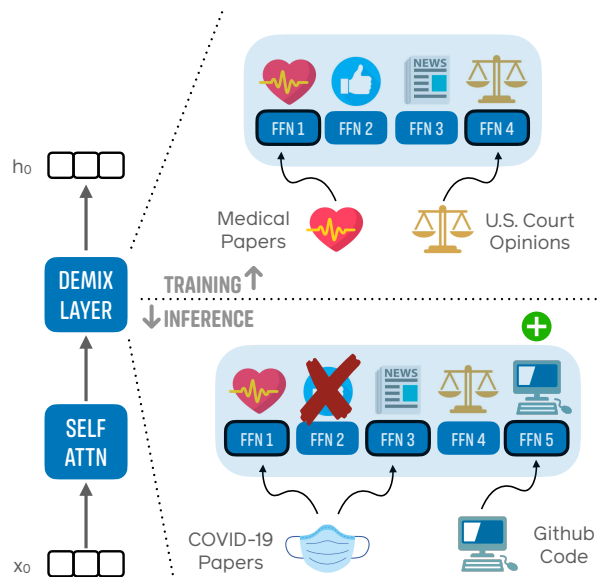


Figure 1: Illustration of a DEMIX layer in a single transformer block. During training, expert feedforward networks are conditionally activated based on the domain (i.e., scientific papers or court opinions) of the input sequence (i.e., scientific papers or court opinions). At inference time, the language model has new modular functions: domain experts can be mixed to handle heterogeneous domains (e.g., COVID-19 papers), added to adapt to novel domains (e.g., Github code), or removed to reduce the influence of unwanted domains (e.g., social media). Image attribution in §A.1.

2019), making fine-tuning or domain adaptation hard to perform with smaller computational budgets. It is also difficult to adapt to new domains without forgetting the original data (McCloskey and Cohen, 1989; Aghajanyan et al., 2021) or restrict access to certain domains the LM has been exposed to during training (e.g., those that contain hatespeech; Bender et al. 2021), leading to risks of unwanted behavior (Gehman et al., 2020).

To address these limitations of dense training, we argue that LMs should be designed with *modularity*. We propose a modular LM that has components specialized to distinct domains in the training data,

and can be customized at inference-time by mixing, adding, or removing these separated components as needed. This design principle emphasizes the ability to rapidly adapt the LM after training, a need that has been broadly advocated for language systems (Dinan et al., 2021; Lazaridou et al., 2021).

We introduce modularity into an LM with a new domain expert (DEMIX) layer that explicitly conditions the LM on the domain of the input text (when it is known), or estimates the input domain during inference (when it is not known). A DEMIX layer is a drop-in substitute for a feedforward layer in a transformer LM (e.g., GPT-3), creating a specialized version of the layer (or *expert*) per domain (see Figure 1; §3).<sup>1</sup>

We identify domains using coarse provenance categories (e.g., whether a document is a medical research paper or a Reddit post; §2). Training on data from eight different domains, we find that replacing every feedforward layer in the transformer with a DEMIX layer consistently improves in-domain performance (§4). To improve performance in settings in which the target data does not clearly align with a single domain, we introduce a parameter-free probabilistic approach to dynamically estimate a *weighted mixture* of domains during inference (§5). Mixing experts improves DEMIX performance not only on *novel* test-time domains, but also on test data from the *training* domains, which may themselves be heterogeneous. Our results suggest that introducing modularity into an LM need not come at a cost to generalization.

Because DEMIX forces experts to specialize to domains, the overall model can be (partially) disentangled after training. Beyond mixing, we can add (§6) or remove (§7) domain experts, predictably changing model behavior at inference time. Adding experts allows for model adaptation without updating all parameters (hence avoiding forgetting), and removing experts allows for simulating the removal of training domains without additional training. These results demonstrate benefits of moving away from treating data homogeneously during language modeling. Our code is publicly available.<sup>2</sup>

<sup>1</sup>This is an example of conditional computation (Fedus et al., 2021; Lepikhin et al., 2020; Lewis et al., 2021; Roller et al., 2021), which follows prior literature on mixture of experts (Jacobs et al., 1991; Shazeer et al., 2017). Unlike dense training, conditional computation activates different parameters for different inputs. Instead of learning how to route data to experts, the DEMIX layer routing mechanism follows from a natural, observable segmentation of the data.

<sup>2</sup>[anonymouse.com](https://anonymouse.com)

	Domain	Corpus
TRAINING	1B	(Chelba et al., 2014)
	CS	(Lo et al., 2020)
	LEGAL	(Caselaw Access Project)
	MED	(Lo et al., 2020)
	WEBTEXT <sup>†</sup>	(Gokaslan and Cohen, 2019)
	REALNEWS <sup>†</sup>	(Zellers et al., 2019)
	REDDIT	(Baumgartner et al., 2020)
	REVIEWS <sup>†</sup>	(Ni et al., 2019)
NOVEL	ACL PAPERS <sup>◇</sup>	(Dasigi et al., 2021)
	BREAKING NEWS <sup>†◇</sup>	(Baly et al., 2018)
	CONTRACTS <sup>†◇</sup>	(Hendrycks et al., 2021)
	CORD-19	(Wang et al., 2020)
	GITHUB	(Github Archive Project)
	GUTENBERG	(Project Gutenberg)
	TWEETS <sup>†◇</sup>	(Twitter Academic API)
	YELP REVIEWS <sup>†</sup>	(Yelp Reviews)

Table 1: Domains that make up our multi-domain training corpus. See §A.2 for specific token counts and details on how these data were collected. † indicates domains that we (partially) anonymize (§A.3). ◇ indicates smaller domains for which we use 1M (instead of 10M) tokens for evaluation.

## 2 Multi-Domain Corpus

We center this study around a large multi-domain corpus we constructed with document-level metadata that describe *provenance*, or the dataset used to access each document (Table 1). Defining domains in this way is intuitive and conveys a great deal about the language variation in a document. Other accounts of domains (e.g., Lucy and Bamman, 2021; Gururangan et al., 2020) may be studied in future work. While other multi-domain corpora (Koh et al., 2021; Gao et al., 2020) cover many more domains and tasks, our corpus contains substantial metadata-tagged text for language modeling, and datasets with friendly licensing to support reproducibility.

Our multi-domain corpus consists of two parts. The first is a collection of **training** domains: text from eight domains of largely English text (top of Table 1), each of which varies in complexity and coverage and has been the subject of study in NLP. The training domains consist of 73.8B whitespace-separated tokens in total (§A.2).

The second part is a collection of **novel** domains: text from eight domains also of largely English text (bottom of Table 1), which may or may not align with the training domains. The novel domains allow us to measure how models generalize to a more challenging data distribution shift, where domain boundaries may be less clear.

§A.2 has more details on how these data were collected, as well as per-domain token counts. For larger domains, we use an additional 10M tokens for the validation and test sets each. Smaller domains have 1M tokens in each (Table 1). To support future work with the data, we also release an API to download and preprocess it into a format compatible with Fairseq (Ott et al., 2019).<sup>3</sup>

### 3 DEMIX Layer

#### 3.1 Background: Mixture-of-Experts Transformers

The transformer architecture is comprised of interleaved multi-head self-attention, layer-norms, and feedforward networks (Vaswani et al., 2017). Each of these layers produces a vector representation for each of the input tokens. Our focus is on the feedforward component:

$$\mathbf{h}_{t,\ell} = \text{FFN}(\mathbf{h}_{t,\ell-1}), \quad (1)$$

where  $\mathbf{h}_{t,\ell}$  is the vector for the  $t$ th token produced by layer  $\ell$ .

Shazeer et al. (2017) propose to replace dense feedforward layers with an ensemble of  $n$  experts  $\text{FFN}_1, \dots, \text{FFN}_n$ , assigned weights respectively by functions  $g_1, \dots, g_n$ :

$$\text{FFN}(\mathbf{h}_{t,\ell-1}) = \sum_{j=1}^n g_j(\mathbf{h}_{t,\ell-1}) \cdot \text{FFN}_j(\mathbf{h}_{t,\ell-1}) \quad (2)$$

The  $g$  function routes tokens to different experts, usually each a separate dense feedforward network. If  $g$  routes to a single expert, then the computational cost (in floating-point operations; FLOPs) will be same as the original network, even though it has more than  $n$  times as many parameters.

#### 3.2 DEMIX Routing

Previous approaches *learn* the weighting functions  $g$  at a token-level, and either assign at most one (Fedus et al., 2021) or two (Lepikhin et al., 2020) experts per token. This necessitates load balancing to encourage the model to use all experts instead of relying on just a few (Lewis et al., 2021).

We instead use domain metadata provided with training documents to route data to experts at the *document* (i.e., sequence) level. During training,

<sup>3</sup>anonymous.com

		Parameters per GPU			
		125M	350M	760M	1.3B
DENSE	GPUs	32	64	128	128
	Total Experts	0	0	0	0
	GPUs/expert	0	0	0	0
	Total params	125M	350M	760M	1.3B
	TFLOPs/update TFLOPs/GPU	556 31	3279 37	13,637 45	23,250 51
DEMIX	GPUs	32	64	128	128
	Total Experts	8	8	8	8
	GPUs/expert	4	8	16	16
	Total params	512M	1.8B	3.8B	7.0B
	TFLOPs/update TFLOPs/GPU	556 31	3279 37	13,637 48	23,250 55

Table 2: Our specifications for training DENSE and DEMIX LMs. All models are trained for about 48 hours on V100 GPUs. DEMIX layers increase the total parameters of the LM while maintaining (or increasing) throughput, measured in TFLOPs/GPU. We use the formula described in Narayanan et al. (2021) to calculate these metrics. See §A.4 for more details.

every token in the same sequence is assigned to the same expert based on the domain label.

Let  $\mathcal{D}$  denote the set of domain labels (i.e., the eight labels in Table 8). If we index the experts by  $\mathcal{D}$  and  $d \in \mathcal{D}$  is the domain label for the current training instance, then

$$g_j(\mathbf{h}_{t,\ell}) = \begin{cases} 1 & \text{if } j = d \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

While we assume that each *training* document is associated with a single domain label, we relax this requirement at inference time (§5), which improves model performance in heterogeneous domains.

#### 3.3 DEMIX Architecture

Our design results in one expert in a DEMIX layer per domain (i.e., eight experts for eight training domains in our multi-domain corpus).

We replace *every* feedforward layer in the transformer with a separate DEMIX layer, in contrast to previous work (Fedus et al., 2021; Lepikhin et al., 2020) that interleaves shared and expert layers. Preliminary experiments showed that interleaving led to worse in-domain performance (see §A.5 for more details). Future work may comprehensively compare different architectural choices.

Each expert  $\text{FFN}_j$  is a two-layer MLP with the same dimensions as the original FFN layer of the transformer. This means that the effective number of parameters in the overall DEMIX LM increases (Table 2), despite using the equivalent FLOP count of a significantly smaller model.

### 3.4 DEMIX Training

To train an LM with DEMIX layers, we partition the GPUs among the domains, so that each GPU is assigned a single domain (along with its corresponding expert). We fill a mini-batch with  $k$  sequences from a particular domain, and we send each mini-batch to its dedicated expert. We use larger batch sizes with distributed data parallel between expert parameters on GPUs assigned to the same domain; we assign  $n/8$  GPUs to each domain (Table 2). To reduce overfitting, we ensure that each of these  $n/8$  GPUs is assigned to different shards of the training data for their domain.

Compared to DENSE LMs, DEMIX layers achieve the same or slightly higher throughput (measured in TFLOPs/GPU) for the same total FLOPs per update, despite adding significantly more parameters (Table 2). DEMIX achieves higher throughput because while we sync shared parameters across all GPUs, we only sync expert parameters allocated to the same domain. DENSE models sync all parameters across all GPUs. As we increase model size, this reduces latency costs between GPUs, and hence, faster training.

## 4 In-Domain Performance

Our first set of experiments considers the impact of replacing the feedforward layers in a transformer LM with DEMIX layers. We run all experiments in this section with the training domains (Table 1).

### 4.1 Experimental Setup

#### Architecture, Input and Hyperparameters

The model architecture is a randomly-initialized LM with the GPT-3 (Brown et al., 2020) architecture (i.e., small, medium, large, and XL) implemented in Fairseq (Ott et al., 2019). We use the GPT-2 (Radford et al., 2019) vocabulary of 50,264 BPE types, and train with 1,024-token sequences, across document boundaries. We prepend a beginning-of-sentence token to each document. See §A.6 for training hyperparameters.

**Evaluation** We follow previous work by using runtime as our computational budget (Lewis et al., 2021). We report test perplexities after 48 hours of training on NVIDIA V100 32GB GPUs. We display the number of GPUs used for each model size in Table 2. For all experiments, we report each result with respect to parameters per GPU.

	Parameters per GPU			
	125M	350M	760M	1.3B
DENSE	20.6	16.5	14.5	13.8
DENSE (balanced)	19.9	15.8	14.3	13.6
+DOMAIN-TOKEN	19.2	15.9	14.3	13.4
DEMIX (naive)	18.4	15.5	14.2	13.8
DEMIX (cached; §5.4)	17.8	14.7	13.9	13.4

Table 3: Average in-domain test-set perplexity across the 8 domains in the training data. We discuss the last row in §5.4. See §A.7 for per-domain results.

### 4.2 Compared Models

**DENSE** The first baseline is a DENSE LM, implemented with distributed data parallel (Li, 2021). There is no explicit conditioning on domain.

**DENSE (balanced)** Under this setting, we train densely but ensure that the model is exposed to an equal amount of data from each domain. While there is still no explicit conditioning on domain, the gradient updates that the model makes during training are an average of those computed across all domains represented in a batch.

**+DOMAIN-TOKEN** This model is trained identically to DENSE (balanced), but we prepend a token to every sequence indicating its domain (during training and test time). We ignore the domain token when computing perplexity during evaluation.

**DEMIX (naive)** We replace every feedforward layer in the transformer with a DEMIX layer, and employ DEMIX training (§3). Under this *naive* setting, the test data’s domain is known and revealed (e.g., the CS expert is used for CS test data). We relax this assumption in the next section.

### 4.3 Results

Table 3 shows test perplexities, averaged across the eight training domains. First, we observe that domain balancing is consistently helpful for DENSE training. Next, we observe that the benefits of additional domain information (i.e, domain tokens or DEMIX layers) are clearest for the smaller models; for larger models, the benefits are smaller but consistent. These results suggest that domain information enables the model to better specialize to different training domains. However, as the model size grows, the DENSE baseline improves, catching up to the DEMIX (naive) model, at least when considering the average perplexity across domains.

Domain	1.3B parameters per GPU		
	DENSE	DEMIX (naive)	DEMIX (cached prior; §5.4)
1B	11.8	11.5	<b>11.3</b>
CS	13.5	12.2	<b>12.1</b>
LEGAL	6.8	<b>6.7</b>	<b>6.7</b>
MED	9.5	9.2	<b>9.1</b>
WEBTEXT	<b>13.8</b>	14.6	14.3
REALNEWS	<b>12.5</b>	13.3	13.1
REDDIT	28.4	30.6	<b>28.1</b>
REVIEWS	14.0	12.6	<b>12.5</b>
<b>Average</b>	13.8	13.8	<b>13.4</b>

Table 4: Test perplexity by domain for largest models. We discuss the last column in §5.4.

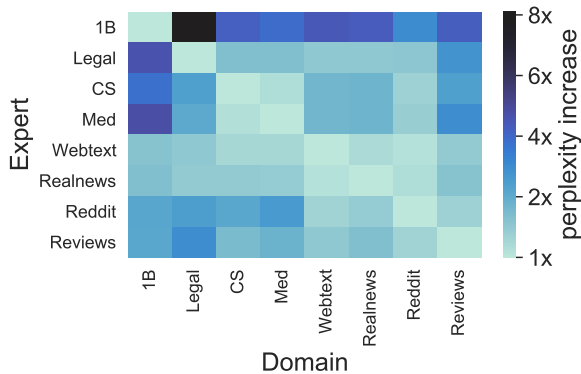


Figure 2: Heatmap of expert performance ratios, using the largest DEMIX LM (1.3B parameters per GPU). The diagonal indicates that expert specialization to their own domain. While some experts (e.g., 1B, MED) do not transfer well to most domains in the training corpus, WEBTEXT and REALNEWS experts transfer much better, confirming the heterogeneity of those domains.

#### 4.4 Domain Heterogeneity

However, a more complete view of the experiments with the largest model is shown in Table 4. We see that even at scale, most training domains benefit from DEMIX layers in a naive setting (where the domain label is revealed at test time), but some do not; WEBTEXT, REALNEWS, and REDDIT fare worse than the DENSE baseline. We hypothesize that DENSE training is advantageous for heterogeneous domains. Heterogeneous domains have a higher overlap with other training domains, and therefore, benefit from parameter sharing.

We support this explanation with a matrix of performance ratios across all domain experts (Figure 2), comparing the performance of all experts against the expert explicitly trained for each domain. We generally find that experts perform best on their assigned domain. However, the experts assigned to domains that benefit from DENSE training

perform relatively well on many training domains. Expert affinity to a target domain correlates positively with the bigram overlap between the expert and target domains ( $r=0.40$ ,  $t=3.45$ ,  $p=0.001$ ).

These findings suggest that a discrete notion of domain, while helpful on average, is too rigid. In the next section, we mitigate this issue by softening Equation 3 into a mixture of domain experts.

## 5 Mixing Experts at Inference Time

The previous section establishes that incorporating DEMIX layers improves LM performance on test data from *known* training domains. In practice, however, text may not come with a domain label, may straddle multiple domains, or may not belong to any of the domains constructed at training time.

In these cases, rather than a hard choice among experts (Equation 3), we propose to treat  $g_1, \dots, g_n$  as mixture coefficients, transforming the domain membership of an input text into a matter of probabilistic belief. Unlike previously proposed mixture-of-experts formulations (Shazeer et al., 2017; Lepikhin et al., 2020), this approach is parameter-free and computed only at test time.

### 5.1 Dynamically Estimating Domain Membership

Consider the probabilistic view of language modeling, where we estimate  $p(\mathbf{x}_t | \mathbf{x}_{<t})$ . We introduce a domain variable,  $D_t$ , alongside each word. We assume that this hidden variable depends on the history,  $\mathbf{x}_{<t}$ , so that:

$$p(\mathbf{x}_t | \mathbf{x}_{<t}) = \sum_{j=1}^n p(\mathbf{x}_t | \mathbf{x}_{<t}, D_t = j) \cdot \underbrace{p(D_t = j | \mathbf{x}_{<t})}_{g_j} \quad (4)$$

This model is reminiscent of class-based  $n$ -gram LMs (Brown et al., 1992; Saul and Pereira, 1997).

We have already designed the DEMIX LM to condition on a domain label, giving a form for  $p(X_t | \mathbf{x}_{<t}, D_t = j)$ . The modification is to treat  $g_1, \dots, g_n$  as a posterior probability over domains, calculated at each timestep, given the history so far.

To do this, we apply Bayes' rule:

$$\begin{aligned} p(D_t = j | \mathbf{x}_{<t}) &= \frac{p(\mathbf{x}_{<t} | D_t = j) \cdot p(D_t = j)}{p(\mathbf{x}_{<t})} \quad (5) \\ &= \frac{p(\mathbf{x}_{<t} | D_t = j) \cdot p(D_t = j)}{\sum_{j'=1}^n p(\mathbf{x}_{<t} | D_t = j') \cdot p(D_t = j')} \quad (6) \end{aligned}$$

The conditional probabilities of word sequences given a domain label, as noted above, are already

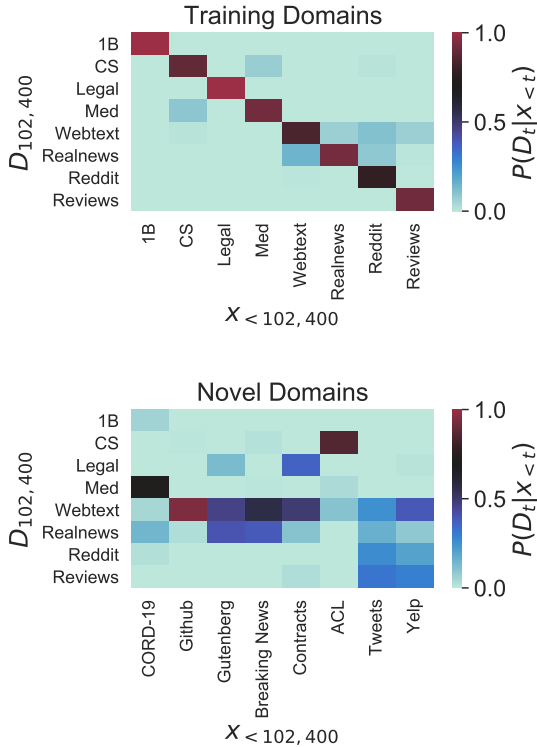


Figure 3: Estimates of posteriors  $p(D_t | \mathbf{x}_{<t})$  with a DEMIX LM (1.3B parameters per GPU), after 100 sequences (i.e., 102,400 tokens) of data in training (top heatmap) and novel domains (bottom heatmap).

defined by the DEMIX LM. For the prior over domain labels, we consider three alternatives:

**Uniform** Set a uniform prior across domains.

**Updating** Set the prior at timestep  $t$  to be an exponentially-weighted moving average of the posteriors from previous timesteps:

$$p(D_t = j) \propto \sum_{t'=1}^{t-1} \lambda^{t-t'} \cdot p(D_{t'} = j | \mathbf{x}_{<t'}) \quad (7)$$

During evaluation, this moving average is calculated over the posterior at the end of each sequence. The decay factor avoids putting too much weight on calculations made early in the dataset, when posterior calculations are noisier (§A.8). We performed a small grid search to set the value  $\lambda$ , and found that  $\lambda = 0.3$  worked well.

**Cached** We calculate the posterior over domain labels from additional data from the test distribution, and fix the prior to that estimate. We use 100 sequences from the validation set to estimate the prior, which we found to result in stable posterior probabilities. See §A.8 for more details, and Figure 7 for an illustration of expert mixing.

	Parameters per GPU			
	125M	350M	760M	1.3B
<b>DENSE</b>	25.9	21.4	18.4	17.8
<b>DENSE (balanced)</b>	25.3	19.6	18.3	17.1
<b>+DOMAIN-TOKEN</b>	24.8	20.4	18.4	18.0
<b>DEMIX (naive)</b>	28.8	23.8	21.8	21.1
<b>DEMIX (average)</b>	27.2	22.4	21.5	20.1
<b>DEMIX (uniform)</b>	24.5	20.5	19.6	18.7
<b>DEMIX (updating)</b>	21.9	18.7	17.6	17.1
<b>DEMIX (cached)</b>	<b>21.4</b>	<b>18.3</b>	<b>17.4</b>	<b>17.0</b>

Table 5: Average perplexity on novel domains. Mixing domain experts with a prior estimated using a small amount of data in the target domain outperforms all other baselines. See §A.7 for per-domain results.

## 5.2 Visualizing Domain Membership

In Figure 3, we plot domain posteriors calculated using the largest DEMIX LM from §4 and the updating prior, after 100 sequences of validation data. For training domains, the associated domain label has the highest probability, but some of the domains are more heterogeneous than we assumed. More variation is observed for the novel domains.

## 5.3 Experimental Setup

Here, we experiment with the corpus of novel domains (Table 1). We evaluate the three mixture treatments of DEMIX layers (§5.1) against five baselines. No new models are trained for this experiment beyond those used in §4.

**DENSE and DENSE (balanced)** These are the basic baselines from §4.

**+DOMAIN-TOKEN** Here test data is evaluated using each domain label token, and we choose the lowest among these perplexity values per test set.

**DEMIX (naive)** Similar to +DOMAIN-TOKEN, we evaluate the data separately with each of the eight experts, and report the lowest among these perplexity values per test set.

**DEMIX (average)** At every timestep, we take a simple average of the eight experts’ predictions.

## 5.4 Results

**Novel Domain Performance** Ensembling DEMIX experts outperforms DENSE baselines and using experts individually (i.e., the “naive” baseline), and caching a prior before evaluation results in the best average performance (Table 5). Ensembling DEMIX experts with a cached prior allows smaller models to match or outperform

much larger DENSE models. Weighted ensembling outperforms simple averaging, confirming the importance of specificity in the expert mixture. These results demonstrate that modularity need not come at a cost to generalization to new domains.

**In-Domain Performance** We can also apply the expert mixture variant of inference (using a cached prior) to the training domains; see the last line of Table 3. We see performance improvements across all training domains for every scale, though the largest gains come from heterogeneous domains (Table 4 and §A.7; across all model sizes, REDDIT improves on average 10.7%, WEBTEXT 2.4%, REALNEWS 1.9%), confirming that domain labels may not align with the most effective boundaries.

## 6 Domain Adaptation with New Experts

Domain-adaptive pretraining (DAPT)<sup>4</sup> is an effective method for adapting LMs to new domains (Gururangan et al., 2020). However, DAPT with DENSE training (DENSE-DAPT) can be expensive and susceptible to forgetting older domains.

DEMIX layers allow for cheap adaptation without forgetting through a technique we call DEMIX-DAPT (see Figure 8 for an illustration). To adapt to a new domain, we initialize a new expert in each DEMIX layer using the parameters of the nearest pretrained expert, which we identify using the posterior calculations from §5 on a held-out sample. We then train the added expert on target data, *updating only the new expert parameters*. For inference, we mix experts with a cached prior (§5).

### 6.1 Experimental Setup

We compare DEMIX-DAPT to DENSE-DAPT on the novel domains. We report test perplexity after adapting to each domain for 1 hour with 8 NVIDIA V100 32GB GPUs, tracking validation perplexity every 10 minutes for early stopping. We adapt to each novel domain with the same hyperparameters as §4, except with a 10x smaller learning rate. DEMIX-DAPT updates about 10% of the total parameters in the DEMIX LM, while DENSE-DAPT updates all parameters of the DENSE LM.

### 6.2 Results

**Adding One Expert** We display examples of DEMIX-DAPT and DENSE-DAPT on a single do-

<sup>4</sup>This approach involves continued training on data from the new domain, and typically precedes supervised fine-tuning on task data, hence *pretraining*.

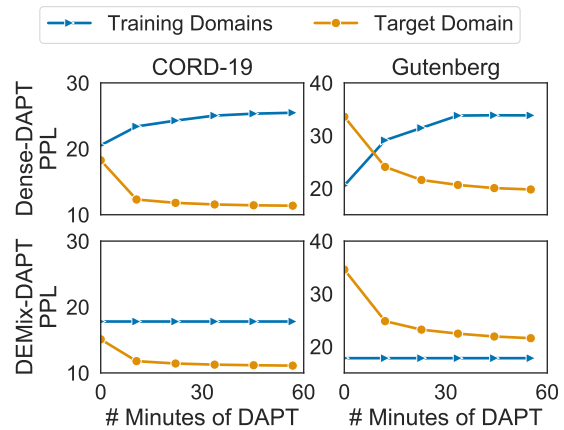


Figure 4: Adapting an LM (125M parameters per GPU) to CORD-19 or GUTENBERG. Top row: with DENSE-DAPT, average perplexity on all training domains degrades. Bottom row: DEMIX-DAPT avoids forgetting while achieving close (for GUTENBERG) or better (for CORD-19) performance on the target domain.

Domains	# Experts	Parameters per GPU			
		125M	350M	760M	1.3B
TRAINING	8	17.8	14.7	13.9	<b>13.4</b>
	16	<b>17.7</b>	<b>14.6</b>	<b>13.7</b>	<b>13.4</b>
NOVEL	8	21.4	18.3	17.4	17.0
	16	<b>16.0</b>	<b>14.0</b>	<b>13.5</b>	<b>12.5</b>

Table 6: Average perplexity in training and novel domains before and after adding 8 experts adapted to the novel domains (via DEMIX-DAPT). Adding experts reduces perplexity on novel and training domains.

main in Figure 4. As DENSE-DAPT proceeds, its performance on the training domains progressively worsens (see §A.11 for results with larger LMs). In contrast, DEMIX-DAPT reduces perplexity on the novel domain *without* forgetting.

**Adding Eight Experts** We find that adding *all* eight experts adapted to novel domains to the DEMIX model from §4 reduces perplexity on novel and previously seen domains (Table 6). For example, across all model sizes, on average, we see a 2.4% reduction on MED, 1.8% reduction on REALNEWS, and 2% reduction on REDDIT (§A.7). These gains are small, which is not surprising given our small budget for adaptation. Even so, these results suggest DEMIX-DAPT can support continual learning (Chen et al., 2018).

Domain	125M Parameters per GPU		
	+EXPERT	-EXPERT	-DOMAIN
1B	13.7	25.5	<b>30.4</b>
CS	15.7	22.4	<b>25.4</b>
LEGAL	8.9	20.9	<b>22.7</b>
MED	12.4	18.6	<b>21.9</b>
WEBTEXT	20.9	<b>27.3</b>	25.4
REALNEWS	18.9	<b>26.7</b>	25.0
REDDIT	34.4	47.8	<b>51.3</b>
REVIEWS	20.5	39.0	<b>43.0</b>
<b>Average</b>	18.2	28.5	<b>30.6</b>

Table 7: Removing a domain expert (−EXPERT) degrades perplexity on the corresponding domain, approaching the performance of an LM that has not been exposed to that domain (−DOMAIN). Here we bold the *worst* performing model for each domain.

## 7 Language Models with Removable Parts

LM training data contain undesirable content such as hatespeech (Gehman et al., 2020). DENSE training precludes the ability to restrict model access to these domains during inference, as may be desired for user-facing tasks (Xu et al., 2020).

DEMIX layers offer a simple solution: since domain experts specialize (Figure 2), experts assigned to unwanted domains can be disabled at test-time.<sup>5</sup>

### 7.1 Experimental Setup

Does disabling an expert simulate a model that has not been exposed to a particular training domain? To answer this question, we compare three settings: **+EXPERT**, a DEMIX LM with all experts active, **−EXPERT**, a DEMIX LM with a domain expert deactivated, and **−DOMAIN**, a DEMIX LM trained from scratch without a particular domain.<sup>6</sup>

For all settings, we use a DEMIX LM (125M parameters per GPU) from §4 and expert mixing with a cached prior (§5) for inference.

### 7.2 Results

Removing a domain expert harms model performance on the associated domain, in most cases approaching the performance of a model that has not been exposed to data from that domain (Table 7). In some cases (e.g., WEBTEXT and REALNEWS), −EXPERT even underperforms −DOMAIN. This

<sup>5</sup>Removing an expert offers no guarantee of having fully forgotten content from the removed domain, since there are shared parameters in the model.

<sup>6</sup>We replace the removed domain with GUTENBERG, since our cluster allocates training jobs via 8-GPU nodes.

leads us to conjecture that most domain-specific learning happens within the DEMIX layer.

## 8 Related Work

Document metadata has been used to improve topic models (Mimno and McCallum, 2012), adapt RNN-based LMs (Jaech and Ostendorf, 2018), learn document representations (Card et al., 2018), and improve text generation control (Zellers et al., 2019; Keskar et al., 2019). Other inference-time methods (Dathathri et al., 2020; Liu et al., 2021) may be used to steer text generation with DEMIX experts.

Related to variation across domains is crosslingual variation, and multilingual models benefit from language-specific parameters (Fan et al., 2020; Pfeiffer et al., 2020; Chau et al., 2020). Future work might explore a compositional approach with language and domain experts.

DEMIX-DAPT is closely related to model expansion techniques used to incorporate new reinforcement learning or visual tasks (Rusu et al., 2016; Draelos et al., 2017) as well as adapter modules for pretrained LMs (Houlsby et al., 2019; Pfeiffer et al., 2020). Other continual learning methods, include regularization (Kirkpatrick et al., 2017), meta-learning (Munkhdalai and Yu, 2017), episodic memory (de Masson d’Autume et al., 2019), and data replay (Sun et al., 2019) may be combined with DEMIX-DAPT.

Multi-domain models have been studied in machine translation (Pham et al., 2021), or with smaller LMs using explicit supervision (e.g., Wright and Augenstein, 2020) or dense training (e.g., Maronikolakis and Schütze, 2021). Previous studies have shown the importance considering domains when adapting LMs (Ramponi and Plank, 2020; Gururangan et al., 2020). Our study establishes the importance of considering domains when training LMs from scratch.

## 9 Conclusion

We introduce DEMIX layers, which provide modularity to an LM at inference time, addressing limitations of dense training by providing a rapidly adaptable system. DEMIX layers experts can be mixed to handle heterogeneous or unseen domains, added to iteratively incorporate new domains, and removed to restrict unwanted domains. Future work may combine domain and token-level routing, discover domains automatically with unsupervised learning, or scale the number of training domains.



535  
536  
537  
538  
539  
540  
541  
  
542  
543  
544  
545  
546  
547  
  
548  
549  
550  
551  
552  
553  
554  
  
555  
556  
557  
  
558  
559  
560  
561  
562  
563  
564  
  
565  
566  
567  
568  
  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
  
581  
582  
583  
584  
585  
  
586  
  
587  
588  
589  
590

## References

- Armen Aghajanyan, Akshat Shrivastava, Anchit Gupta, Naman Goyal, Luke Zettlemoyer, and Sonal Gupta. 2021. [Better fine-tuning by reducing representational collapse](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Roei Aharoni and Yoav Goldberg. 2020. [Unsupervised domain clusters in pretrained language models](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7747–7763, Online. Association for Computational Linguistics.
- Ramy Baly, Georgi Karadzhov, Dimitar Alexandrov, James Glass, and Preslav Nakov. 2018. [Predicting factuality of reporting and bias of news media sources](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3528–3539, Brussels, Belgium. Association for Computational Linguistics.
- Jason Baumgartner, Savvas Zannettou, Brian Keegan, Megan Squire, and Jeremy Blackburn. 2020. [The pushshift reddit dataset](#).
- Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. [On the dangers of stochastic parrots: Can language models be too big?](#) In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT '21*, page 610–623, New York, NY, USA. Association for Computing Machinery.
- Peter F. Brown, Vincent J. Della Pietra, Peter V. deSouza, Jenifer C. Lai, and Robert L. Mercer. 1992. [Class-based  \$n\$ -gram models of natural language](#). *Computational Linguistics*, 18(4):467–480.
- 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. 2020. [Language models are few-shot learners](#).
- Dallas Card, Chenhao Tan, and Noah A. Smith. 2018. [Neural models for documents with metadata](#). *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.
- Caselaw Access Project. [Caselaw access project](#).
- Ethan C. Chau, Lucy H. Lin, and Noah A. Smith. 2020. [Parsing with multilingual BERT, a small corpus, and a small treebank](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1324–1334, Online. Association for Computational Linguistics.
- Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and Tony Robinson. 2014. [One billion word benchmark for measuring progress in statistical language modeling](#).
- Zhiyuan Chen, Bing Liu, Ronald Brachman, Peter Stone, and Francesca Rossi. 2018. *Lifelong Machine Learning: Second Edition*.
- Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A. Smith, and Matt Gardner. 2021. [A dataset of information-seeking questions and answers anchored in research papers](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4599–4610, Online. Association for Computational Linguistics.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2020. [Plug and play language models: A simple approach to controlled text generation](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Cyprien de Masson d’Autume, Sebastian Ruder, Lingpeng Kong, and Dani Yogatama. 2019. [Episodic memory in lifelong language learning](#).
- Emily Dinan, Gavin Abercrombie, A. Stevie Bergman, Shannon Spruit, Dirk Hovy, Y-Lan Boureau, and Verena Rieser. 2021. [Anticipating safety issues in e2e conversational ai: Framework and tooling](#).
- T. Draelos, N. Miner, Christopher C. Lamb, Jonathan A. Cox, Craig M. Vineyard, Kristofor D. Carlson, William M. Severa, C. James, and J. Aimone. 2017. [Neurogenesis deep learning: Extending deep networks to accommodate new classes](#). *2017 International Joint Conference on Neural Networks (IJCNN)*, pages 526–533.
- Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Man-deep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. 2020. [Beyond english-centric multilingual machine translation](#).
- William Fedus, Barret Zoph, and Noam Shazeer. 2021. [Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity](#).
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. [The pile: An 800gb dataset of diverse text for language modeling](#).



752	Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019.	Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz,	809
753	<a href="#">Justifying recommendations using distantly-labeled</a>	Andy Davis, Quoc V. Le, Geoffrey E. Hinton, and	810
754	<a href="#">reviews and fine-grained aspects</a> . In <i>Proceedings</i>	Jeff Dean. 2017. <a href="#">Outrageously large neural net-</a>	811
755	<i>of the 2019 Conference on Empirical Methods in</i>	<a href="#">works: The sparsely-gated mixture-of-experts layer</a> .	812
756	<i>Natural Language Processing and the 9th Interna-</i>	In <i>5th International Conference on Learning Rep-</i>	813
757	<i>tional Joint Conference on Natural Language Pro-</i>	<i>resentations, ICLR 2017, Toulon, France, April 24-</i>	814
758	<i>cessing (EMNLP-IJCNLP)</i> , pages 188–197, Hong	<i>26, 2017, Conference Track Proceedings</i> . OpenRe-	815
759	Kong, China. Association for Computational Lin-	<a href="#">view.net</a> .	816
760	guistics.		
761	Yonatan Oren, Shiori Sagawa, Tatsunori Hashimoto,	Emma Strubell, Ananya Ganesh, and Andrew McCal-	817
762	and Percy Liang. 2019. <a href="#">Distributionally robust lan-</a>	lum. 2019. <a href="#">Energy and policy considerations for</a>	818
763	<a href="#">guage modeling</a> . In <i>Proceedings of the 2019 Con-</i>	<a href="#">deep learning in NLP</a> . In <i>Proceedings of the 57th</i>	819
764	<i>ference on Empirical Methods in Natural Language</i>	<i>Annual Meeting of the Association for Computa-</i>	820
765	<i>Processing and the 9th International Joint Confer-</i>	<i>tional Linguistics</i> , pages 3645–3650, Florence, Italy.	821
766	<i>ence on Natural Language Processing (EMNLP-</i>	Association for Computational Linguistics.	822
767	<i>IJCNLP)</i> , pages 4227–4237, Hong Kong, China. As-		
768	sociation for Computational Linguistics.	Fan-Keng Sun, Cheng-Hao Ho, and Hung-Yi Lee.	823
769	Myle Ott, Sergey Edunov, Alexei Baevski, Angela	2019. <a href="#">Lamol: Language modeling for lifelong lan-</a>	824
770	Fan, Sam Gross, Nathan Ng, David Grangier, and	<a href="#">guage learning</a> .	825
771	Michael Auli. 2019. <a href="#">fairseq: A fast, extensible</a>	Twitter Academic API. <a href="#">Twitter academic api</a> .	826
772	<a href="#">toolkit for sequence modeling</a> . In <i>Proceedings of</i>	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	827
773	<i>the 2019 Conference of the North American Chap-</i>	Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz	828
774	<i>ter of the Association for Computational Linguistics</i>	Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all</a>	829
775	<i>(Demonstrations)</i> , pages 48–53, Minneapolis, Min-	<a href="#">you need</a> .	830
776	nesota. Association for Computational Linguistics.		
777	Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Se-	Lucy Lu Wang, Kyle Lo, Yoganand Chandrasekhar,	831
778	bastian Ruder. 2020. <a href="#">MAD-X: An Adapter-Based</a>	Russell Reas, Jiangjiang Yang, Doug Burdick, Dar-	832
779	<a href="#">Framework for Multi-Task Cross-Lingual Transfer</a> .	rarin Eide, Kathryn Funk, Yannis Katsis, Rodney	833
780	In <i>Proceedings of the 2020 Conference on Empirical</i>	Kinney, Yunyao Li, Ziyang Liu, William Merrill,	834
781	<i>Methods in Natural Language Processing (EMNLP)</i> ,	Paul Mooney, Dewey Murdick, Devvret Rishi, Jerry	835
782	pages 7654–7673, Online. Association for Computa-	Sheehan, Zhihong Shen, Brandon Stilson, Alex	836
783	tional Linguistics.	Wade, Kuansan Wang, Nancy Xin Ru Wang, Chris	837
784	MinhQuang Pham, Josep Maria Crego, and François	Wilhelm, Boya Xie, Douglas Raymond, Daniel S.	838
785	Yvon. 2021. <a href="#">Revisiting multi-domain machine trans-</a>	Weld, Oren Etzioni, and Sebastian Kohlmeier. 2020.	839
786	<a href="#">lation</a> . <i>Transactions of the Association for Computa-</i>	<a href="#">Cord-19: The covid-19 open research dataset</a> .	840
787	<i>tional Linguistics</i> , 9:17–35.	Dustin Wright and Isabelle Augenstein. 2020. <a href="#">Trans-</a>	841
788	Project Gutenberg. <a href="#">Project gutenberg</a> .	<a href="#">former based multi-source domain adaptation</a> . In	842
789	Alec Radford, Jeff Wu, Rewon Child, David Luan,	<i>Proceedings of the 2020 Conference on Empirical</i>	843
790	Dario Amodei, and Ilya Sutskever. 2019. <a href="#">Language</a>	<i>Methods in Natural Language Processing (EMNLP)</i> ,	844
791	<a href="#">models are unsupervised multitask learners</a> .	pages 7963–7974, Online. Association for Computa-	845
792	Alan Ramponi and Barbara Plank. 2020. <a href="#">Neural unsu-</a>	tional Linguistics.	846
793	<a href="#">pervised domain adaptation in NLP—A survey</a> . In	Jing Xu, Da Ju, Margaret Li, Y-Lan Boureau, Jason	847
794	<i>Proceedings of the 28th International Conference</i>	Weston, and Emily Dinan. 2020. <a href="#">Recipes for safety</a>	848
795	<i>on Computational Linguistics</i> , pages 6838–6855,	<a href="#">in open-domain chatbots</a> .	849
796	Barcelona, Spain (Online). International Committee	Yelp Reviews. <a href="#">Yelp reviews</a> .	850
797	on Computational Linguistics.	Rowan Zellers, Ari Holtzman, Hannah Rashkin,	851
798	Stephen Roller, Sainbayar Sukhbaatar, Arthur Szlam,	Yonatan Bisk, Ali Farhadi, Franziska Roesner, and	852
799	and Jason Weston. 2021. <a href="#">Hash layers for large</a>	Yejin Choi. 2019. <a href="#">Defending against neural fake</a>	853
800	<a href="#">sparse models</a> .	<a href="#">news</a> . In <i>NeurIPS</i> .	854
801	Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Des-		
802	jeardins, Hubert Soyer, James Kirkpatrick, Koray		
803	Kavukcuoglu, Razvan Pascanu, and Raia Hadsell.		
804	2016. <a href="#">Progressive neural networks</a> .		
805	Lawrence Saul and Fernando Pereira. 1997. <a href="#">Aggregate</a>		
806	<a href="#">and mixed-order Markov models for statistical lan-</a>		
807	<a href="#">guage processing</a> . In <i>Second Conference on Empiri-</i>		
808	<i>cal Methods in Natural Language Processing</i> .		

855	<b>A Appendix</b>		
856	<b>A.1 Image attribution</b>		
857	Images retrieved from <a href="http://emojipedia.org">emojipedia.org</a> or		
858	<a href="http://istockphoto.com">istockphoto.com</a> .		
859	<b>A.2 Collecting Domains</b>		
860	For most domains, we use the associated		
861	sources, listed in Table 8, without modifica-		
862	tion. For GUTENBERG, we use the scrap-		
863	ing tool provided in <a href="https://github.com/aparrish/gutenberg-dammit">https://github.com/</a>		
864	<a href="https://github.com/aparrish/gutenberg-dammit">aparrish/gutenberg-dammit</a> . For BREAKING		
865	NEWS, we identify a list of factually reli-		
866	able English news sources, using the list cu-		
867	rated by Baly et al. (2018). Specifically,		
868	we filter on "high" factuality in the data		
869	provided in this repository: <a href="https://github.com/ramybaly/News-Media-Reliability">https://github.</a>		
870	<a href="https://github.com/ramybaly/News-Media-Reliability">com/ramybaly/News-Media-Reliability</a> . We		
871	then use Newspaper3K ( <a href="https://newspaper.readthedocs.io/en/latest/">https://newspaper.</a>		
872	<a href="https://newspaper.readthedocs.io/en/latest/">readthedocs.io/en/latest/</a> ) to scrape the lat-		
873	est 1000 articles from each site. After dropping		
874	duplicates, we arrive at about 20K articles from		
875	400 news sources. We provide downloading links		
876	and general instructions at <a href="http://anonymous.com">anonymous.com</a> .		
877	<b>A.3 Dataset Anonymization</b>		
878	To anonymize certain datasets, we apply a suite		
879	of regexes that aim to identify common patterns		
880	of user-identifiable data and substitute them with		
881	dummy tokens. We display anonymization regexes		
882	and associated dummy tokens in Table 9.		
883	<b>A.4 Calculating TFLOPs/GPU</b>		
884	We use the formula presented in Narayanan		
885	et al. (2021) to calculate TFLOPs/GPU and		
886	TFLOPs/update. The spreadsheet that contains		
887	the calculations and formula can be accessed here:		
888	<a href="http://anonymous.com">anonymous.com</a>		
889	<b>A.5 Interleaving Experiments</b>		
890	We hypothesize that shared layers may serve as a		
891	bottleneck to find shared features between domains,		
892	and may impact performance adversely when train-		
893	ing domains are highly different from one another.		
894	Indeed, preliminary experiments suggest that in-		
895	terleaving expert layers causes large performance		
896	hits in the most distinct domains, i.e., those with		
897	lower vocabulary overlap with other domains in the		
898	corpus.		
	<b>A.6 Hyperparameter Assignments</b>		899
	We display hyperparameter assignments for LM		900
	pretraining in Tables 11, 12, 13, and 14. We set the		901
	total number of training steps based on this allo-		902
	cated runtime, set 8% of these steps to be warm-		903
	up, and use the Adam optimizer (Kingma and Ba,		904
	2017) with a polynomial learning rate decay. Learn-		905
	ing rates are tuned for each model separately over		906
	{0.0001, 0.0003, 0.0005}, taking the fastest learn-		907
	ing rate that avoids divergence. Each worker pro-		908
	cesses two sequences of length 1,024, and gradients		909
	are accumulated over 8 updates. We clip gradients		910
	if their $L_2$ norm exceeds 0.1. These settings are		911
	inspired by Lewis et al. (2021).		912
	<b>A.7 Per-Domain Results</b>		913
	We display the rest of the per-domain test results in		914
	the spreadsheets at the following link: <a href="http://anonymous.com">anonymous.</a>		915
	<a href="http://anonymous.com">com</a>		916
	<b>A.8 Domain Posterior Calculations</b>		917
	We track calculated domain posteriors over se-		918
	quences of development data in Figure 5 (train-		919
	ing domains) and Figure 6 (novel domains). The		920
	calculate domain posteriors are noisier for earlier		921
	sequences, stabilizing usually after around 50 se-		922
	quences. For all experiments, we conservatively		923
	use 100 sequences of data to compute the domain		924
	posterior, though one may be able to accurately		925
	calculate the domain posterior for some domains		926
	with less data.		927
	<b>A.9 Illustration of Expert Mixing</b>		928
	See Figure 7 for an illustration of expert mixing.		929
	<b>A.10 Illustration of DEMIX-DAPT</b>		930
	See Figure 8 for an illustration of DEMIX-DAPT.		931
	<b>A.11 Perplexity changes after DENSE-DAPT</b>		932
	In Table 10, we display the average perplexity		933
	change after performing DENSE-DAPT on a new		934
	domain. We observe that across all model sizes,		935
	DENSE-DAPT improves performance in the novel		936
	domain, at the cost of a large performance hit in		937
	the training domains.		938

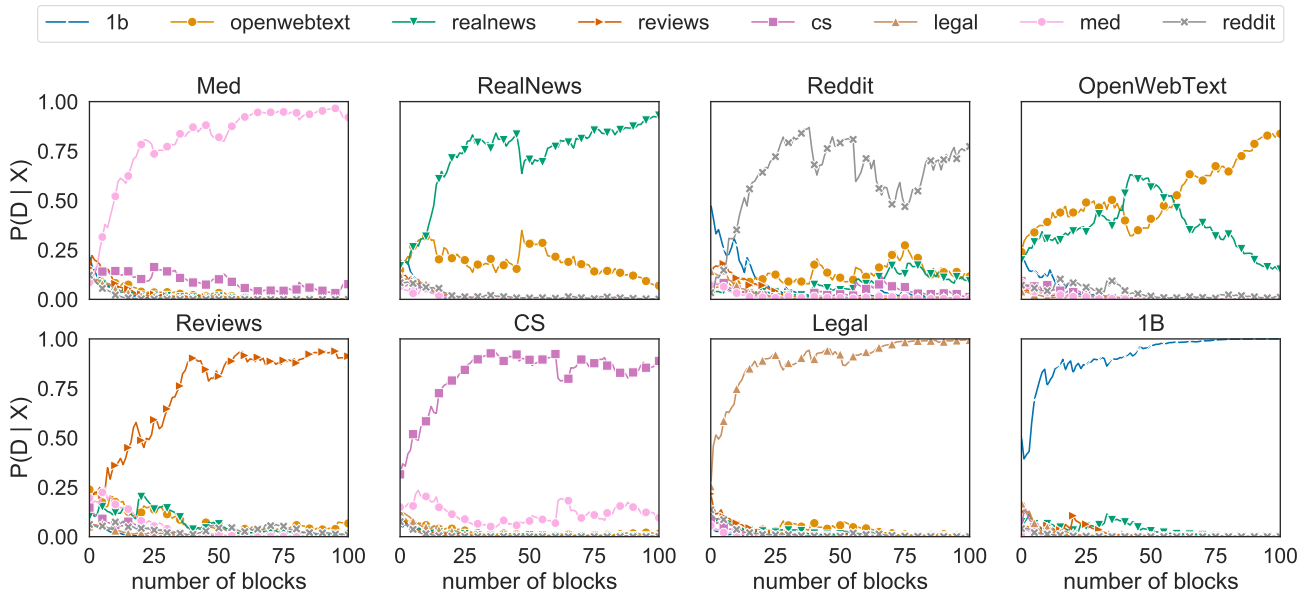


Figure 5: Calculated domain posteriors for 8 training domains.

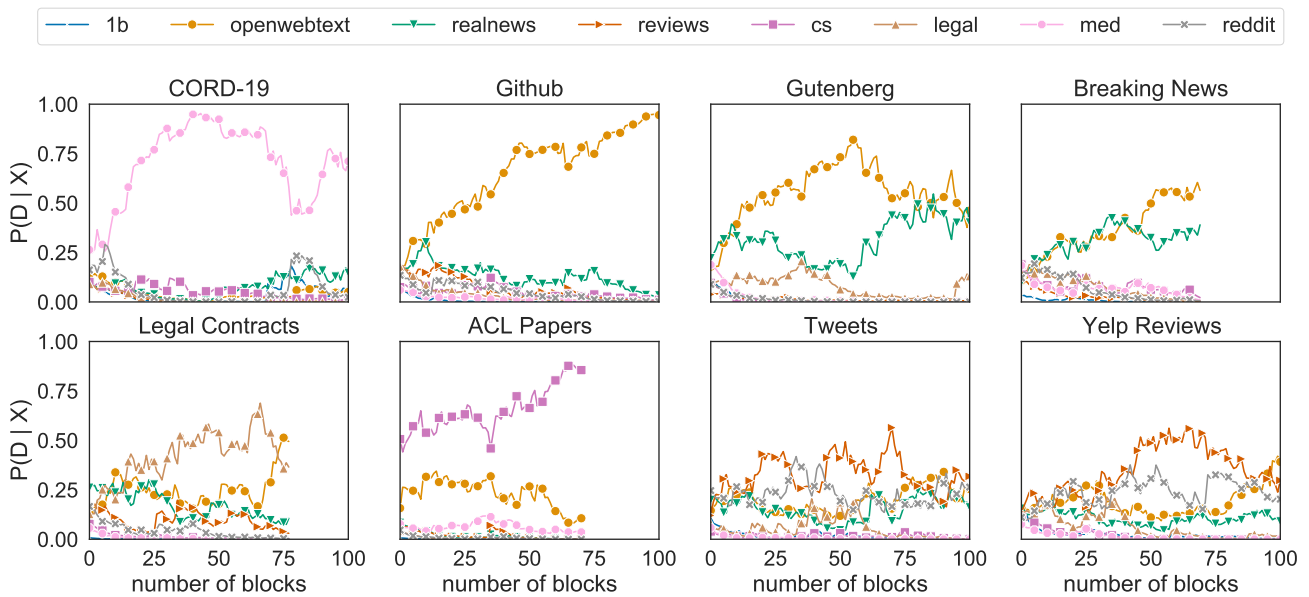


Figure 6: Calculated domain posteriors for 8 novel domains.

	Domain	Corpus	# Train (Eval.) Tokens
TRAINING	1B	30M NewsWire sentences (Chelba et al., 2014)	700M (10M)
	CS	1.89M full-text CS papers from S2ORC (Lo et al., 2020)	4.5B (10M)
	LEGAL	2.22M U.S. court opinions, 1658 to 2018 (Caselaw Access Project)	10.5B (10M)
	MED	3.2M full-text medical papers from S2ORC (Lo et al., 2020)	9.5B (10M)
	WEBTEXT <sup>†</sup>	8M Web documents (Gokaslan and Cohen, 2019)	6.5B (10M)
	REALNEWS <sup>†</sup>	35M articles from REALNEWS (Zellers et al., 2019)	15B (10M)
	REDDIT	Reddit comments from pushshift.io (Baumgartner et al., 2020)	25B (10M)
	REVIEWS <sup>†</sup>	30M Amazon product reviews (Ni et al., 2019)	2.1B (10M)
<b>Total</b>			<b>73.8B (80M)</b>
	Domain	Corpus	# Train (Eval.) Tokens
NOVEL	ACL PAPERS	1.5K NLP papers from ACL (Dasigi et al., 2021)	1M (1M)
	BREAKING NEWS <sup>†</sup>	20K latest articles from 400 English news sites (Baly et al., 2018)	11M (1M)
	CONTRACTS <sup>†</sup>	500 commercial legal contracts (Hendrycks et al., 2021)	1.5M (1M)
	CORD-19	400K excerpts from COVID-19 research papers (Wang et al., 2020)	60M (10M)
	GITHUB	230K public Github repository contents (Github Archive Project)	200M (10M)
	GUTENBERG	3.2M copyright-expired books (Project Gutenberg)	3B (10M)
	TWEETS <sup>†</sup>	1M English tweets from 2013-2018	8M (1M)
	YELP REVIEWS <sup>†</sup>	6M Yelp restaurant reviews (Yelp Reviews)	600M (10M)

Table 8: Domains that make up our multi-domain training corpus, including the size of our training and evaluation (i.e. validation and test) data, in whitespace-separated tokens. † indicates datasets that we (partially) anonymize (§2). REDDIT was extracted and obtained by a third party and made available on [pushshift.io](https://pushshift.io), and was anonymized by Xu et al. (2020); we use their version. See Appendix §A.2 for more details on how these data were collected.

Category	Link to Regex	Dummy Token
Email	<a href="https://regex101.com/r/ZqsF9x/1">https://regex101.com/r/ZqsF9x/1</a>	<EMAIL>
DART	<a href="https://regex101.com/r/0tQ6EN/1">https://regex101.com/r/0tQ6EN/1</a>	<DART>
FB User ID	<a href="https://regex101.com/r/GZ15EZ/1">https://regex101.com/r/GZ15EZ/1</a>	<FB_USERID>
Phone Number	<a href="https://regex101.com/r/YrDpPD/1">https://regex101.com/r/YrDpPD/1</a>	<PHONE_NUMBER>
Credit Card Number	<a href="https://regex101.com/r/9NTO6W/1">https://regex101.com/r/9NTO6W/1</a>	<CREDIT_CARD_NUMBER>
Social Security Number	<a href="https://regex101.com/r/V5GPNL/1">https://regex101.com/r/V5GPNL/1</a>	<SSN>
User handles	<a href="https://regex101.com/r/vpey04/1">https://regex101.com/r/vpey04/1</a>	<USER>

Table 9: Anonymization schema. We anonymize text using the regexes provided in the above links for the categories listed.

		Parameters			
		125M	350M	760M	1.3B
<b>DENSE-</b>	T	+70.1%	+21.4%	+16.7%	+20.6%
<b>DAPT</b>	N	-55.1%	-46.6%	-38.3%	-44.4%

Table 10: Average change in perplexity in training (T) and novel (N) domains after DENSE-DAPT. Negative values indicate better performance relative to the original DENSE LM. While average perplexity in the novel domains decreases more for DENSE-DAPT, this comes at the cost of a significant deterioration in performance in training domains.

<b>Computing Infrastructure</b> 32 Volta 32GB GPUs	
<b>Hyperparameter</b>	<b>Assignment</b>
architecture	GPT-3 small
tokens per sample	1024
batch size	2
number of workers	2
learning rate	[5e-4, 3e-4, 1e-4]
clip norm	0.1
gradient accumulation steps	8
number of steps	300,000
save interval updates	6,000
validation interval	3,000
number of warmup steps	24,000
learning rate scheduler	polynomial decay
learning rate optimizer	Adam
Adam beta weights	(0.9, 0.95)
Adam epsilon	10e-8
weight decay	0.1

Table 11: Hyperparameters for pretraining the LM with 125M parameters per GPU. All hyperparameters are the same for DEMIX and DENSE training.

<b>Computing Infrastructure</b> 64 Volta 32GB GPUs	
<b>Hyperparameter</b>	<b>Assignment</b>
architecture	GPT-3 medium
tokens per sample	1024
batch size	2
number of workers	2
learning rate	[5e-4, 3e-4, 1e-4]
clip norm	0.1
gradient accumulation steps	8
number of steps	120,000
save interval updates	3,000
validation interval	2,000
number of warmup steps	9,600
learning rate scheduler	polynomial decay
learning rate optimizer	Adam
Adam beta weights	(0.9, 0.95)
Adam epsilon	10e-8
weight decay	0.1

Table 12: Hyperparameters for pretraining the LM with 350M parameters per GPU. All hyperparameters are the same for DEMIX and DENSE training.

<b>Computing Infrastructure</b> 128 Volta 32GB GPUs	
<b>Hyperparameter</b>	<b>Assignment</b>
architecture	GPT-3 large
tokens per sample	1024
batch size	2
number of workers	2
learning rate	[5e-4, 3e-4, 1e-4]
clip norm	0.1
gradient accumulation steps	8
number of steps	65,000
save interval updates	2,000
validation interval	1,000
number of warmup steps	5,200
learning rate scheduler	polynomial decay
learning rate optimizer	Adam
Adam beta weights	(0.9, 0.95)
Adam epsilon	10e-8
weight decay	0.1

Table 13: Hyperparameters for pretraining the LM with 760M parameters per GPU. All hyperparameters are the same for DEMIX and DENSE training.

<b>Computing Infrastructure</b> 128 Volta 32GB GPUs	
<b>Hyperparameter</b>	<b>Assignment</b>
architecture	GPT-3 XL
tokens per sample	1024
batch size	2
number of workers	2
learning rate	[5e-4, 3e-4, 1e-4]
clip norm	0.1
gradient accumulation steps	8
number of steps	50000
save interval updates	2,000
validation interval	500
number of warmup steps	4000
learning rate scheduler	polynomial decay
learning rate optimizer	Adam
Adam beta weights	(0.9, 0.95)
Adam epsilon	10e-8
weight decay	0.1

Table 14: Hyperparameters for pretraining the LM with 1.3B parameters per GPU. All hyperparameters are the same for DEMIX and DENSE training.



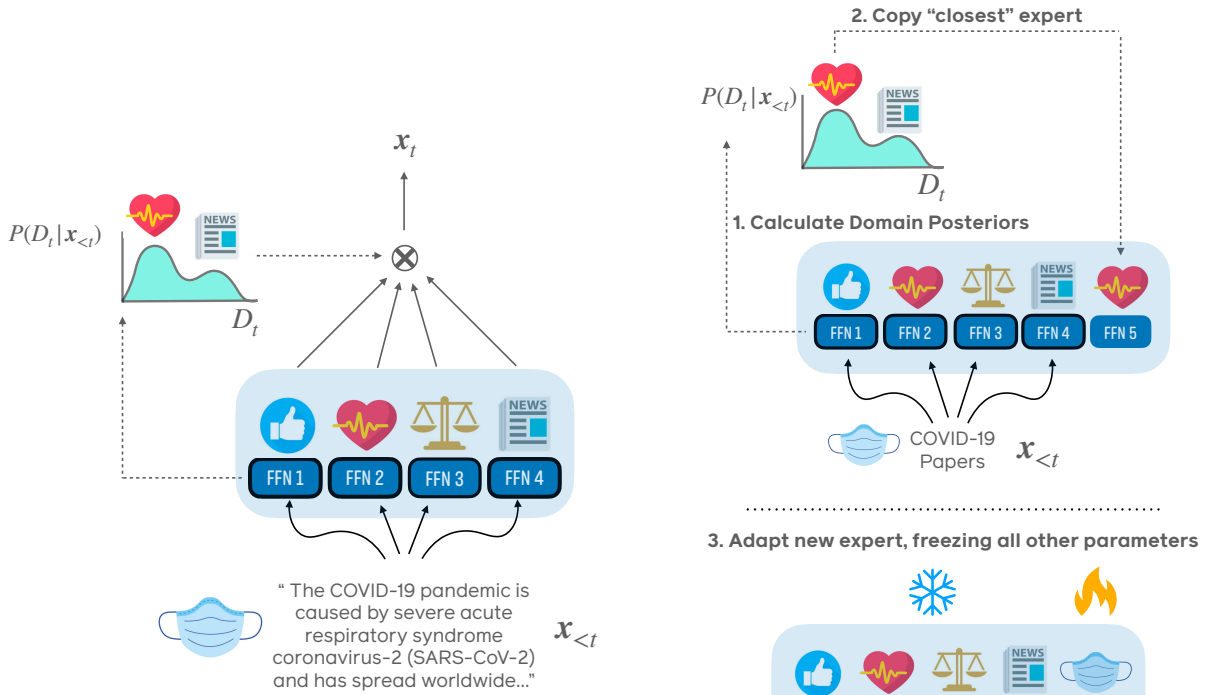


Figure 7: Illustration of inference with domain expert mixing. For a given input text  $x_{<t}$  from CORD-19, we estimate a posterior domain probabilities  $p(D_t | x_{<t})$ , informed by a prior that is either iteratively updated during inference, or is precomputed and cached on held-out data. In this example, the model assigns highest domain probabilities to the medical and news domains. We use these probabilities in a weighted mixture of expert outputs to compute the output  $x_t$ .

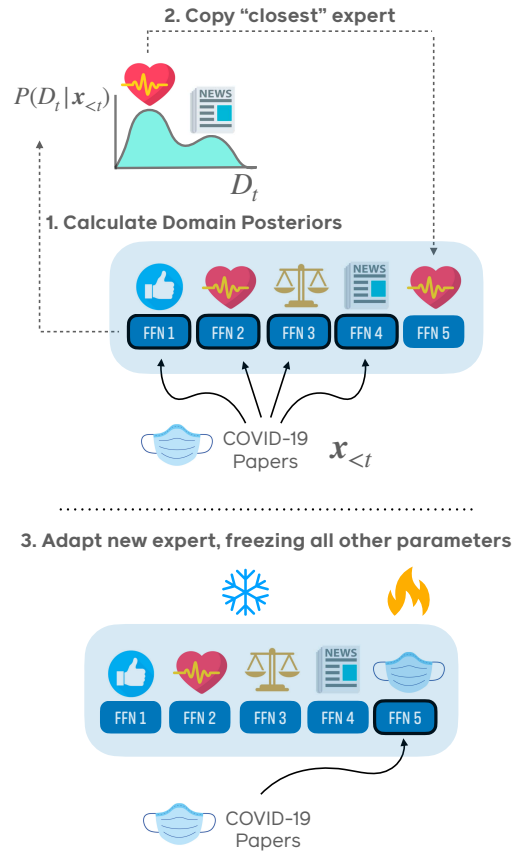


Figure 8: Illustration of DEMix-DAPT. First, we estimate domain posteriors on a held out sample of the target domain (in this case, CORD-19). We then initialize a new expert with the parameters of the most probable expert under the domain posterior distribution. Finally, we adapt the parameters of the newly initialized expert to the target domain, keeping all other parameters in the LM frozen.