

A Declarative–Procedural Perspective on Expert Routing in Bilingual Mixture-of-Experts Language Models

Anonymous ACL submission

Abstract

We investigate whether Mixture-of-Experts (MoE) language models develop linguistically structured expert routing during bilingual language acquisition. Inspired by the Declarative–Procedural framework, we analyze lexical, grammatical, and syntactic processing in a decoder-only English–German MoE Transformer trained under sequential language exposure. We construct a probe-based validation set and extract token-level routing distributions to quantify category-dependent specialisation using mutual information, routing entropy, and Jensen–Shannon divergence. The curriculum-trained model exhibits a peak mutual information of 0.1148 at layer 5, indicating category-dependent differences in routing distributions across linguistic categories. Surprisingly, a no-curriculum baseline trained on mixed English–German data shows stronger aggregate specialisation, reaching a peak mutual information of 0.2599 at the same layer. These results suggest that interpretable linguistic organization emerges within MoE routing patterns even without sequential language exposure. Rather than uniformly increasing specialisation, staged bilingual exposure appears to redistribute specialisation across languages, yielding a more balanced bilingual routing profile. Code and data will be released upon acceptance.

1 Introduction

Mixture-of-Experts (MoE) language models route each token across a subset of expert networks instead of sending every token through the same feed-forward network (Shazeer et al., 2017; Fedus et al., 2022). This makes MoE models useful not only for scaling, but also for studying how neural language models choose experts internally. If different kinds of linguistic tokens are routed to different experts, the routing mechanism may exhibit an organized form of expert distribution inside the model.

In this work, we examine whether bilingual MoE language models develop routing patterns

that correspond to meaningful linguistic categories. We focus on three categories motivated by the Declarative–Procedural view of language: lexical knowledge, grammatical processing, and syntactic structure (Ullman, 2001b,a, 2020). In this paper, a probe is a token chosen from a sentence to be examined. Lexical probes are tokens that test word knowledge, such as irregular forms. Grammatical probes are tokens that test rule-based morphology and agreement. Syntactic probes are tokens that participate in clause structure and sentence organization. Instead of examining only the model’s final predictions, we directly analyse the expert-routing distributions associated with these probe tokens.

We train a sparse bilingual transformer using a staged English–German curriculum. The training curriculum starts with English and gradually transitions to German, approximating structured second-language exposure and curriculum learning in bilingual settings (Bengio et al., 2009; Platanios et al., 2019; Zhang et al., 2019). We compare this structured curriculum to an unstructured no-curriculum setting in which English and German are introduced without gradual cumulative exposure. This lets us examine whether MoE routing becomes organized by language and whether staged bilingual exposure affects how that organization is distributed across languages.

To measure the development of an organized routing schema, we use a held-out validation set with lexical, grammatical, and syntactic probes in both English and German. We examine the router’s probability distribution over experts at each MoE layer for each probe token. We then measure the relationship between expert routing and linguistic category using mutual information, entropy, Jensen–Shannon divergence, and sentence-level permutation testing (Shannon, 1948; Lin, 1991; Cover and Thomas, 2006). This complements recent work that studies behavioral declarative and procedural knowledge in datasets and large language models

(Li et al., 2024) by moving the analysis to the level of internal routing.

Our results show that MoE routing distributions contain measurable information about linguistic category membership: the category of a probe token can be partially inferred from its expert probability distribution because routing and category exhibit non-trivial mutual information. Lexical, grammatical, and syntactic probes are routed differently. The clearest specialisation appears in intermediate MoE layers, whereas later layers route more diffusely and show reduced mutual information. We also find that curriculum affects how specialisation is balanced between languages. The no-curriculum model develops a stronger English routing signature, while the staged L1–L2 curriculum yields a more balanced bilingual organization and supports clearer German expert structure in the middle layers.

In general, this study shows that bilingual MoE routers develop a language-sensitive quantifiable expert-routing structure that reflects linguistically significant differences. These findings suggest that expert routing provides a useful mechanistic window into how bilingual language models allocate computation across different types of language processing.

The central motivation of this paper is to move beyond asking whether a bilingual MoE model performs well on a diagnostic task and instead ask how its internal routing mechanism organizes linguistic computation. Our explicit contribution is a routing-level analysis framework for testing whether expert allocation is sensitive to lexical, grammatical, and syntactic probe categories under staged bilingual exposure. By comparing a forward English–German curriculum with a no-curriculum baseline, we isolate how training order changes the distribution of routing specialisation across languages. This makes the study a mechanistic investigation of expert routing in bilingual MoE models, rather than a general benchmark of language-model accuracy.

2 Related Work

This section places the study in the context of four connected areas: the Declarative–Procedural view of language, the linguistic analysis of transformer models, expert routing in MoE architectures, and curriculum learning for bilingual exposure. Together, these areas lead to the main question of the paper: does expert routing in a bilingual MoE

model become sensitive to lexical, grammatical, and syntactic differences? Also, does staged language exposure change that routing structure?

Declarative–Procedural theory. The Declarative–Procedural framework divides language knowledge into two main parts (Ullman, 2001b,a, 2020). Declarative knowledge relates to stored word knowledge, such as vocabulary and irregular forms. Procedural knowledge relates to rule-based grammar and compositional structure. In this paper, we use this framework only as a linguistic lens. We do not claim that MoE routers are the same as human memory systems.

Declarative and procedural knowledge in LLMs. Recent studies have looked at declarative and procedural knowledge in language models mainly through model outputs and task performance (Li et al., 2024). Our work explores this question within the model. Instead of just asking if the model provides the right answer, we investigate whether a sparse bilingual model allocates computation differently for varying types of probe tokens.

Transformer linguistic structure. Previous studies indicate that transformer layers do not all encode the same type of linguistic information (Vaswani et al., 2017; Tenney et al., 2019; Clark et al., 2019; Elhage et al., 2021). Intermediate layers often have clearer syntactic and semantic structures than very early or very late layers. This motivates our layer-wise analysis of MoE routing. We specifically question whether middle routed layers show stronger category-sensitive expert allocation.

Mixture-of-Experts routing. MoE models use routers to direct tokens to a subset of experts (Shazeer et al., 2017; Fedus et al., 2022). Earlier work has shown that experts can specialize based on token statistics, domains, routing design, and balancing constraints (Dai et al., 2024; Falke et al., 2026; Sun et al., 2026). Our research stands out by focusing on bilingual linguistic probe categories and exploring how curriculum structure impacts the distribution of routing specialization across English and German.

Curriculum learning and bilingual exposure. Curriculum learning investigates how the order of training examples influences learning (Bengio et al., 2009; Platanios et al., 2019; Zhang et al., 2019). In bilingual training, staged exposure is useful be-

cause it allows us to compare two scenarios: one where the model sees English first and German later, and another where both languages are mixed from the start. This comparison helps us assess whether the order of language exposure affects how routing specialization develops across languages.

3 Routing Analysis Framework

Let E denote the expert-routing random variable, C denote the linguistic category random variable, and $r^{(l)}(t)$ denote the routing probability vector of token t at layer l . Analysis is performed on the routed layers $L = \{1, 3, 5, 7\}$.

The category set is defined as

$$C = \{\textit{lexical}, \textit{grammatical}, \textit{syntactic}\}$$

Probing mechanism. Probe tokens are annotated in each sentence using a parser (details in Methodology and Appendix B). The routing vector corresponding to the probe token is extracted from the router during inference and used to compute the routing statistics described below.

Analysis is performed independently for each routed layer. All routing distributions and statistical measures are therefore computed separately for each layer.

To analyse the relationship between the model’s routing behaviour and different categories, we frame the null hypothesis as

$$H_0 : P(E|C) = P(E)$$

meaning that the expert routing distribution is independent of linguistic category.

The alternate hypothesis is

$$H_1 : P(E|C) \neq P(E)$$

meaning that the expert routing distribution depends on linguistic category.

3.0.1 Statistical Testing Framework

The primary statistical tool used is mutual information between expert-routing distributions and linguistic category. Higher mutual information indicates a stronger association between the router’s expert probability mass and linguistic category.

The marginal routing probability assigned to expert e is computed as

$$p(e) = \frac{\sum_{\tau \in T} r^{(l)}(\tau)[e]}{\sum_{e'} \sum_{\tau \in T} r^{(l)}(\tau)[e']} \quad (1)$$

and the joint probability between expert e and category c is computed as

$$p(e, c) = \frac{\sum_{\tau \in T_c} r^{(l)}(\tau)[e]}{\sum_{c'} \sum_{e'} \sum_{\tau \in T_{c'}} r^{(l)}(\tau)[e']} \quad (2)$$

where T denotes the complete set of probe tokens and T_c denotes the subset belonging to category c .

Mutual information is then computed as

$$MI(E; C) = \sum_e \sum_c p(e, c) \log \left(\frac{p(e, c)}{p(e)p(c)} \right) \quad (3)$$

where $p(e)$ is the marginal routing probability assigned to expert e and $p(c)$ is the marginal probability of the routed token belonging to category c .

Robustness of the mutual information is analysed with a permutation test where the p -level is set to be less than 0.01 in order to eliminate coincidence of mutual information.

We can now frame our hypothesis based on mutual information:

$$H_0 : MI(E; C) = 0$$

$$H_1 : MI(E; C) > 0$$

3.0.2 Entropy Analysis

Entropy is considered as the second statistical tool to analyse how distributed the routing is across experts. Higher entropy indicates a more diffuse expert-usage distribution across probes across experts, whereas lower entropy indicates concentration of probability mass on a smaller subset of experts.

Shannon entropy is given by

$$H(E) = - \sum_e p(e) \ln p(e) \quad (4)$$

where $p(e)$ is the marginal distribution of expert e being chosen. When entropy is computed for a specific linguistic category, $p(e)$ is estimated only from the probe tokens belonging to that category. Thus, category-wise entropy measures how concentrated or diffuse expert usage is for lexical, grammatical, or syntactic probes separately.

Other tools include Jensen–Shannon Divergence (JSD), which gives a measure of how routing between different linguistic categories is similar or different, which is given by

$$\begin{aligned} & JSD(P(E|C_i) \parallel P(E|C_j)) \\ &= \frac{1}{2} D_{KL}(P(E|C_i) \parallel M) \\ &+ \frac{1}{2} D_{KL}(P(E|C_j) \parallel M) \end{aligned} \quad (5)$$

where

$$M = \frac{1}{2} (P(E|C_i) + P(E|C_j)).$$

4 Methodology

4.1 Dataset Construction and Probe Annotation

The dataset is constructed from the FineWeb-Edu corpus (Penedo et al., 2024) for the English part and the German portion of mC4 (Xue et al., 2021) for the German part. The pipeline uses the Stanza module to categorise words in each sentence as lexical, grammatical, and syntactic. The taxonomy for this classification is given in Appendix A.

One probe token is annotated for each sentence-category pair. A sentence may contain probe tokens belonging to multiple linguistic categories and can therefore appear in more than one category-specific dataset. In such cases, the sentence is duplicated across the relevant categories, with each category using its corresponding probe token.

The following table provides the size of the English and German datasets:

Language	Lexical	Grammatical	Syntactic
English	460,250	650,000	650,000
German	138,127	175,000	175,000

Table 1: Dataset size by language and linguistic category.

Note: Lexical probe extraction was terminated before reaching the target dataset size because of computational constraints, resulting in smaller lexical subsets for both languages. This does not affect the reported analyses, which are conducted on a fixed held-out validation set.

4.2 Sequential L1–L2 Curriculum

The curriculum is designed to approximate staged second-language acquisition. Here, the first language is English (EN) and the second language is

German (DE). The first 1–8 epochs consist of an EN-only curriculum. From epochs 9–16, the curriculum slowly transitions from EN-only to bilingual (EN + DE). The choice of language was based on the differences in grammatical structure and the linguistic category taxonomy (more details in Appendix A).

$$\lambda_s = 0.20 + 0.05(s - 9), \quad 9 \leq s \leq 14 \quad (6)$$

where λ_s denotes the proportion of German samples presented during epoch s .

Within each epoch, the amount of lexical, grammatical, and syntactic probe tokens is equal and follows a repetition + new exposure system. This system can be given by

$$\begin{aligned} Cur_s &= 0.6 \cdot Cur_{s-1} \\ &+ 0.4 \cdot (New \text{ or } Cur_{s'} \text{ where } s' < s - 1) \end{aligned} \quad (7)$$

for $2 \leq s \leq 16$.

$$Cur_1 = 1 \cdot New \quad (8)$$

where Cur_s is the curriculum of epoch s and New is the set of samples that are not part of any previous curriculum.

4.3 Diagnostic Held-out Validation

The validation set consists of 12,000 samples (6,000 per language and 2,000 per category per language). The validation set is a minimal-pairs dataset consisting of a pair of sentences, S^+ and S^- , where S^+ is a linguistically valid and acceptable sentence and S^- is a linguistically invalid sentence.

The S^+ set was generated synthetically using OpenAI’s GPT-5.1 model in order to obtain controlled, diverse, and task-specific examples. The generation was guided by rigorous prompting (prompts provided in Appendix E), with valid probe tokens injected into the prompts. Automatic post-generation validation is applied to ensure that generated sentences follow the correct category subtype, are not duplicates, and contain the required metadata. The S^- sentence is then deterministically obtained by modifying the probe token to invalidate the sentence in a linguistically meaningful way (rules provided in Appendix E).

4.4 Experimental Setup

A custom 8-layer MoE language model is used, with a feed-forward dimension of 2048, an embedding size of 512, and 8 attention heads per layer. The model uses the multilingual mBERT tokenizer with a vocabulary size of 119,547. The total parameter count of the model is 86.4 million.

A relatively small model size is chosen to facilitate controlled experimentation and clearer analysis of expert routing behaviour. The MoE layers are placed in alternate transformer layers (1, 3, 5, and 7). The primary training setup uses the proposed sequential L1–L2 curriculum, with an additional No-Curriculum setting used as an ablation baseline.

The no-curriculum baseline uses a fixed 80:20 English-German mixture at every epoch. This keeps the overall bilingual composition broadly similar to the forward curriculum, although the aggregate exposure is not exactly identical: the forward curriculum yields 81.56% English and 18.44% German across the full schedule, whereas the no-curriculum condition remains fixed at 80% English and 20% German.

Both conditions receive the same total number of training samples. The difference between them is the order in which the languages are introduced. In the forward curriculum, the model is trained on English first and German is introduced gradually. In the no-curriculum condition, English and German are mixed from the start using the same overall English-German proportions.

The primary runs use top- k routing with $k = 3$ and a router load-balancing coefficient of 0.01. Additional hyperparameters are provided in Appendix C.

5 Results and Analysis

All the results, mathematical analyses, and metrics reported below are obtained through inference on the forward-curriculum model and the no-curriculum ablation model over the held-out validation set. Unless otherwise stated, all MI values reported in Sections 5.1–5.5 are computed on the pooled bilingual validation set combining English and German probes. Language-specific analyses are introduced separately in Section 5.6.

5.1 Behavioural Performance

Table 2 reports validation performance for the forward-curriculum and no-curriculum conditions. The two setups achieve broadly comparable be-

havioural outcomes on the held-out diagnostic validation set. The Accuracy is computed as token-level next-token accuracy on the grammatical sentence strings only, averaged over all non-padding tokens; it is not a pairwise sentence-choice accuracy between S^+ and S^- . The forward-curriculum model attains lower perplexity (111.19) than the no-curriculum model (117.01), whereas the no-curriculum condition achieves marginally higher accuracy (29.65% versus 28.99%).

Because the two training conditions achieve similar behavioural performance, the routing patterns reported in the following sections are less likely to reflect simple differences in overall model competence.

The relatively low absolute scores reflect the challenging evaluation setup. The purpose of the diagnostic set is to expose routing behaviour under controlled linguistic contrasts rather than to serve as a benchmark of language-model capability. These scores are not the main outcome of the paper. They just indicate that both models can handle the diagnostic examples. The primary analysis focuses on the router. We examine whether different token types are sent to different experts during validation.

Condition	Perplexity	Accuracy (%)
Forward Curriculum	111.19	28.99
No Curriculum	117.01	29.65

Table 2: Validation performance on the held-out diagnostic probe set.

5.2 Emergence of Category-Dependent Routing

The primary analysis investigates whether routing specialisation occurs across different linguistic categories. Mutual Information (MI) is computed from the inference routing logs to quantify this phenomenon. The MI values vary across layers, indicating that the degree of specialisation differs throughout the network. Layer 5 shows the highest mutual information value (0.1148), suggesting the strongest category-dependent routing behaviour. In contrast, the final routed layer (layer 7) shows an MI value substantially lower than the other routed layers.

To determine whether the obtained MI values (routing-category associations) are statistically significant, a permutation test was conducted using 1,000 sentence-level permutations. Across all

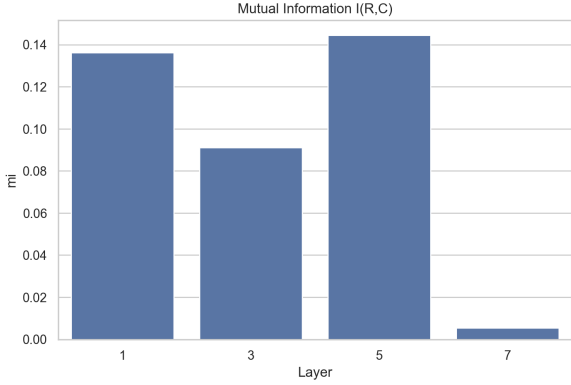


Figure 1: Layer-wise mutual information between expert-routing distributions and linguistic categories.

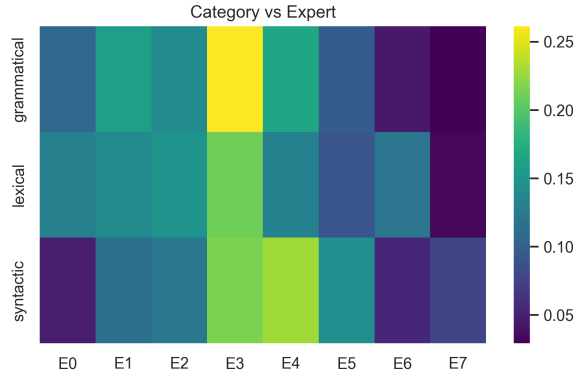


Figure 2: Category-conditioned expert allocation patterns.

433 routed layers, the observed MI values were sig-
 434 nificantly higher than the corresponding layer-wise
 435 null distributions, with statistical significance of
 436 $p < 0.001$ for all layers. This result makes a
 437 shuffled-label explanation unlikely and supports
 438 the presence of category-dependent routing be-
 439 haviour.

Layer	Observed MI	Null Mean	Null Std	p -value
1	0.0520	0.0004	0.0002	0.000999
3	0.0384	0.0003	0.0001	0.000999
5	0.1148	0.0003	0.0001	0.000999
7	0.0207	0.0001	0.0000	0.000999

Table 3: Permutation test results for routing-category mutual information. All routed layers exhibit statistically significant routing-category associations. All reported p -values correspond to the minimum attainable empirical value under 1000 permutations, indicating that no permuted sample exceeded the observed MI.

440 5.3 Expert Allocation Patterns

441 The secondary analysis supporting the expert spe-
 442 cialisation claim is performed through category-
 443 wise and layer-wise expert allocation patterns. This
 444 analysis also explains how experts have been uti-
 445 lized across different linguistic categories. Al-
 446 though the probability mass is distributed across
 447 different experts for each category, certain experts
 448 consistently receive higher routing probabilities for
 449 specific categories. This indicates the presence of
 450 preferential expert allocation rather than uniform
 451 routing.

452 We can also observe that there is no collapse to
 453 a single expert for any particular linguistic cate-
 454 gory. Instead, each category exhibits a distribu-
 455 tion of expert preferences, where some experts
 456 are selected more frequently than others. For ex-

457 ample, grammatical probes show stronger prefer-
 458 ences towards Experts E1 and E3, while syntactic
 459 probes exhibit higher utilisation of Experts E3–E5.
 460 This shows that category-dependent specialisation
 461 emerges through differences in routing preference
 462 rather than strict expert exclusivity.

463 Figure 3 presents the layer-wise expert usage
 464 distribution. It is interesting to note that the dif-
 465 ferent routed layers exhibit different expert usage
 466 distributions, with different subsets of experts do-
 467 minating at different stages of the processing. This
 468 observation is consistent with the layer-wise MI
 469 analysis, suggesting that the routing distribution
 470 differs across different parts of the model.

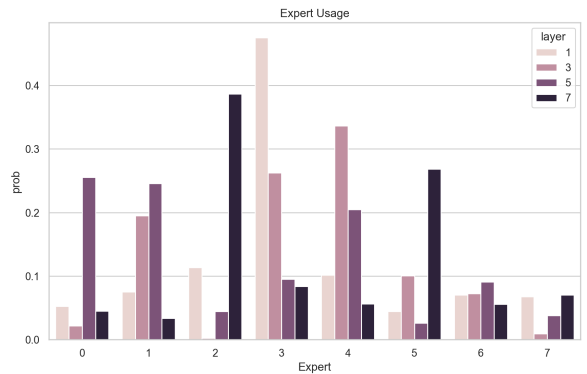


Figure 3: Layer-wise expert utilisation across routed layers.

471 While the heatmap and expert allocation reveal
 472 how the experts are allocated layer-wise as well as
 473 category-wise, they do not indicate how diverse or
 474 concentrated the routing distribution is. To examine
 475 this aspect, we analyse routing entropy across the
 476 routed layers.

5.4 Routing Entropy Analysis

The entropy values vary across different layers and among the different categories as well. Grammatical probes exhibit the lowest entropy consistently in all the layers (with the strongest effect in layer 3), which suggests that the expert utilisation distribution is more concentrated with a smaller subset of experts consistently utilized for this category.

In contrast, the syntactic probes exhibit the highest entropy values across different layers (or equal to lexical in layer 1). This indicates that expert usage is spread across more experts for this particular category across all layers. Together, these observations suggest that different linguistic categories rely on distinct routing strategies, with grammatical processing exhibiting more concentrated expert utilisation and syntactic processing exhibiting more distributed expert usage.

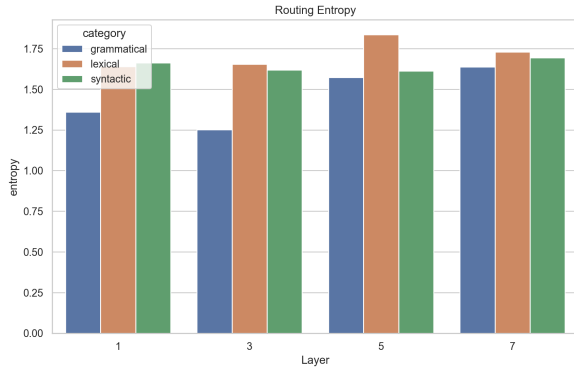


Figure 4: Routing entropy across routed layers and linguistic categories. Lower entropy indicates more concentrated expert utilisation.

5.5 Category Separation Analysis

While entropy explains how concentrated or distributed the expert utilisation is, it does not explain how the routing distribution differs between different linguistic categories. To examine this aspect, we compute pairwise Jensen–Shannon Divergence (JSD) between category-wise routing distributions.

Left Category	Right Category	JSD
Grammatical	Lexical	0.0788
Grammatical	Syntactic	0.1147
Lexical	Syntactic	0.1218

Table 4: Pairwise Jensen–Shannon Divergence between category-conditioned routing distributions.

Table 4 shows the pairwise JSD values between grammatical, lexical and syntactic routing distri-

butions. All category pairs exhibit non-zero divergence, suggesting that the router assigns different expert utilisation distributions to different linguistic categories. We can observe that the highest divergence in the routing distribution is between lexical and syntactic processing, whereas the lowest divergence is between grammatical and lexical processing. These observations are consistent with the mutual information and expert allocation analyses, providing additional evidence for category-dependent routing.

5.6 Effect of Curriculum on Specialisation

In the no-curriculum setting, the mutual information peaks at 0.2599 in layer 5 and reaches 0.1945 in layer 1. Similar to the forward curriculum setting, the strongest routing-category association is observed in the intermediate routed layers. These observations indicate that routing specialisation emerges under both training conditions, suggesting that curriculum learning is not strictly necessary for category-dependent routing to develop.

Layer	No-Curr MI
1	0.1945
3	0.0674
5	0.2599
7	0.0739

Table 5: Pooled bilingual routing-category MI under the no-curriculum condition.

The following table shows language-wise MI values computed separately for English and German.

Layer	Fwd EN	Fwd DE	No-Curr EN	No-Curr DE
1	0.0508	0.1361	0.3810	0.1357
3	0.0366	0.0911	0.1364	0.0804
5	0.1025	0.1444	0.6457	0.1116
7	0.0571	0.0053	0.1141	0.0649

Table 6: Language-wise routing-category mutual information for the forward curriculum and no-curriculum settings.

To verify that the observed category-sensitive routing is not reducible to language identity alone, we additionally compute conditional mutual information $I(R; C \mid L)$, where L denotes language identity. This analysis measures whether routing still contains information about lexical, grammatical, and syntactic category membership after conditioning on whether the probe is English or German.

As reported in Appendix Table 11, conditional mutual information remains non-zero across routed layers in both training conditions, indicating that the routing-category relationship is not explained solely by English-German separation.

Table 6 shows the language-wise mutual information for each routed layer under both the forward curriculum and no-curriculum settings. The language-wise MI values are drastically different for English. In particular, at layer 5, the MI value increases from 0.1025 in the forward curriculum setting to 0.6457 in the no-curriculum setting, representing an increase of approximately 6.3 times. A similar trend can be observed across all routed layers.

In contrast, the difference between the category-dependent routing effect of the two setups is much less visible in German, where MI remains around 0.1 at layers 1 and 5 for both setups (0.1444 in the forward curriculum setting and 0.1116 in the no-curriculum setting at layer 5 for German language). A similar effect is observed in the other layers, where the differences remain comparatively small.

From the above observations, curriculum does not appear to uniformly increase routing specialisation. Instead, the effect of curriculum differs across languages. While the no-curriculum condition produces substantially stronger specialisation in English, the forward curriculum yields a more balanced EN-DE specialisation profile. In contrast, the no-curriculum specialisation pattern is strongly dominated by English routing behaviour.

6 Discussion

Routing distributions provide significant insights into linguistic category membership, even with held-out validation data. Specialisation shows up in both the forward-curriculum and no-curriculum conditions. This indicates that curriculum learning is not necessary for category-sensitive routing to develop. The key question is how curriculum changes how this specialisation is distributed across different languages and layers.

The strongest relationship between routing and category occurs in the intermediate routed layers, particularly in layer 5. This aligns with previous research that suggests intermediate transformer layers often feature clearer linguistic structures. Conversely, the final routed layer exhibits lower mutual information and higher entropy. This implies that later routing becomes more evenly spread across

experts instead of being distinctly separated by category.

The comparison between the two training conditions is crucial. The no-curriculum model shows stronger routing-category dependence in English, which boosts the overall bilingual mutual information. However, this stronger pooled signal mainly originates from English and does not distribute evenly between English and German. The forward curriculum creates a more balanced bilingual routing profile and maintains the German category structure more robustly in the middle layers. Therefore, sequential L1-L2 exposure does not simply boost specialisation everywhere; it alters how specialisation spreads across languages.

This is important because pooled statistics can be misleading. When combining English and German into one score, the no-curriculum model seems stronger. Yet, when analyzing the two languages separately, the staged curriculum appears to lessen English’s dominance and foster a more balanced bilingual routing structure. This provides a clearer understanding of the effects of the curriculum than merely assessing which model has higher total mutual information.

Interpreting behavioral accuracy values requires considering the study design. The model is intentionally small, and the validation set is challenging since it uses controlled minimal-pair contrasts. The aim is not to achieve high benchmark accuracy but to compare internal routing patterns across controlled linguistic probe families. In this context, the behavioral results primarily indicate that both models are sufficiently functional for meaningful routing analysis.

7 Conclusion

Bilingual MoE routing develops measurable category-dependent organization across lexical, grammatical, and syntactic probe families. The strongest routing specialisation is observed in intermediate routed layers and remains robust on held-out validation data. Curriculum learning does not uniformly increase specialisation; instead, it redistributes it across languages. While the no-curriculum condition produces stronger English routing-category dependence, staged L1-L2 exposure yields a more balanced bilingual specialisation profile and stronger German middle-layer routing structure.

635 Limitations

636 This study uses a single EN–DE language pair
637 and a single primary sparse architecture. While
638 the observed routing patterns are consistent across
639 multiple analyses, including mutual information,
640 entropy, Jensen–Shannon divergence, permutation
641 testing, and conditional mutual information con-
642 trols, the extent to which the findings generalize to
643 other language families, larger-scale MoE architec-
644 tures, or different curriculum schedules remains an
645 open question.

646 The diagnostic validation set is generated
647 through a controlled LLM-assisted pipeline and
648 subsequently validated using parser-based checks
649 and deterministic probe construction rules. Al-
650 though these procedures improve consistency and
651 coverage, the evaluation set does not undergo man-
652 ual human validation. Future work could incorpo-
653 rate expert human review and naturally occurring
654 linguistic examples to further verify that the ob-
655 served routing patterns generalize beyond synthetic
656 diagnostic probes. The validation set is dominated
657 by single-piece probes under the mBERT tokenizer,
658 but full subword-aggregation robustness was not
659 recomputed for the final archived routing logs used
660 in this draft

661 The present study uses distinct source corpora
662 for English and German. Although language-
663 conditioned analyses remain significant, future
664 work should examine matched-domain multilin-
665 gual corpora.

666 A reverse curriculum is not included because the
667 forward curriculum is explicitly designed around
668 a fixed cumulative English-German exposure ratio.
669 Simply reversing the schedule would alter not only
670 the temporal order of language presentation but
671 also the total exposure received by each language
672 across training. Consequently, a naive DE→EN
673 reversal would confound sequencing effects with
674 differences in cumulative language exposure. A
675 fair reverse-curriculum comparison would there-
676 fore require a separately constructed schedule that
677 preserves overall language proportions while re-
678 versing the order of introduction.

679 Several potential confounds were explicitly ex-
680 amined. Routing-category associations remained
681 stable across alternative subword aggregation
682 strategies and remained non-zero after condition-
683 ing on language identity and coarse lexical fre-
684 quency controls. Nevertheless, the present study
685 does not fully disentangle all possible interactions

686 between linguistic category, token identity, lexi-
687 cal frequency, morphological complexity, and part-
688 of-speech information. As a result, the reported
689 routing effects should be interpreted as category-
690 sensitive routing behaviour rather than evidence of
691 perfectly isolated category-specific mechanisms.

692 A stricter frequency-balanced lexical control was
693 explored but was not included in the final analy-
694 sis. The curriculum-generated lexical distributions
695 were highly skewed, making it difficult to construct
696 well-matched irregular and regular subsets while
697 maintaining sufficient probe coverage and sample
698 diversity.

699 We do not yet report extensive architectural
700 sweeps, multiple language pairs, or broad-scale
701 hyperparameter sensitivity analyses. In addition,
702 expert intervention experiments (e.g., expert mask-
703 ing or routing interventions) are not included, lim-
704 iting our ability to make strong causal claims about
705 the functional role of individual experts.

706 Finally, the work focuses specifically on sparse
707 MoE architectures because expert-routing distri-
708 butions constitute the primary object of analysis.
709 Dense transformers do not expose an explicit rout-
710 ing mechanism, making direct comparisons of rout-
711 ing specialisation impossible. While dense base-
712 lines remain useful for behavioural benchmarking,
713 they cannot provide the routing-level signals stud-
714 ied in this work.

715 While the present study focuses on identifying
716 category-dependent routing behaviour, future work
717 could perform expert masking or routing interven-
718 tions to determine whether the identified expert
719 preferences play a causal role in linguistic process-
720 ing.

721 References

- 722 Yoshua Bengio, Jérôme Louradour, Ronan Collobert,
723 and Jason Weston. 2009. [Curriculum learning](#). In
724 *Proceedings of the 26th International Conference on*
725 *Machine Learning*, pages 41–48.
- 726 Kevin Clark, Urvashi Khandelwal, Omer Levy, and
727 Christopher D. Manning. 2019. [What does BERT](#)
728 [look at? an analysis of BERT’s attention](#). *Proceed-*
729 *ings of the 2019 ACL Workshop BlackboxNLP: An-*
730 *alyzing and Interpreting Neural Networks for NLP*,
731 pages 276–286.
- 732 Thomas M. Cover and Joy A. Thomas. 2006. *Elements*
733 *of Information Theory*, 2 edition. Wiley.
- 734 Damai Dai, Chengqi Deng, Chenggang Zhao, R. X. Xu,
735 Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng,

736	Xingkai Yu, Y. Wu, Zhenda Xie, Y. K. Li, Panpan Huang, Fuli Luo, Chong Ruan, Zhifang Sui, and Wenfeng Liang. 2024. DeepSeekMoE: Towards ultimate expert specialization in mixture-of-experts language models . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1280–1297.	
737		
738		
739		
740		
741		
742		
743	Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, and 6 others. 2021. A mathematical framework for transformer circuits . <i>Transformer Circuits Thread</i> .	
744		
745		
746		
747		
748		
749		
750		
751	Tobias Falke, Nicolas Anastassacos, Samson Tan, Chankrisna Richy Meas, Chandana Satya Prakash, Nitesh Sekhar, M. Saiful Bari, Krishna Kompella, and Gamaleldin F. Elsayed. 2026. Moe routing testbed: Studying expert specialization and routing behavior at small scale .	
752		
753		
754		
755		
756		
757	William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity . <i>Journal of Machine Learning Research</i> , 23(120):1–39.	
758		
759		
760		
761	Zhuoqun Li, Hongyu Lin, Yaojie Lu, Hao Xiang, Xianpei Han, and Le Sun. 2024. Meta-cognitive analysis: Evaluating declarative and procedural knowledge in datasets and large language models . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 11222–11228.	
762		
763		
764		
765		
766		
767		
768	Jianhua Lin. 1991. Divergence measures based on the shannon entropy . <i>IEEE Transactions on Information Theory</i> , 37(1):145–151.	
769		
770		
771	Guilherme Penedo, Hynek Kydliček, Loubna Ben Allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro von Werra, and Thomas Wolf. 2024. The fineweb datasets: Decanting the web for the finest text data at scale . In <i>Advances in Neural Information Processing Systems Datasets and Benchmarks Track</i> .	
772		
773		
774		
775		
776		
777	Emmanouil Antonios Platanios, Otilia Stretcu, Graham Neubig, Barnabás Póczos, and Tom Mitchell. 2019. Competence-based curriculum learning for neural machine translation . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1</i> , pages 1162–1172.	
778		
779		
780		
781		
782		
783		
784	Claude E. Shannon. 1948. A mathematical theory of communication . <i>Bell System Technical Journal</i> , 27(3):379–423.	
785		
786		
787	Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer . In <i>Proceedings of the International Conference on Learning Representations</i> .	
788		
789		
790		
791		
792		
	Hanchi Sun, Yixin Liu, Yonghui Wu, and Lichao Sun. 2026. Expert threshold routing for autoregressive language modeling with dynamic computation allocation and load balancing .	793 794 795 796
	Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. BERT rediscovers the classical NLP pipeline . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4593–4601.	797 798 799 800 801
	Michael T. Ullman. 2001a. The neural basis of lexicon and grammar in first and second language: The declarative/procedural model . <i>Bilingualism: Language and Cognition</i> , 4(2):105–122.	802 803 804 805
	Michael T. Ullman. 2001b. A neurocognitive perspective on language: The declarative/procedural model . <i>Nature Reviews Neuroscience</i> , 2(10):717–726.	806 807 808
	Michael T. Ullman. 2020. The declarative/procedural model . In Bill VanPatten, Jessica Williams, Gregory D. Keating, and Stefanie Wulff, editors, <i>Theories in Second Language Acquisition: An Introduction</i> , 3 edition, pages 128–161. Routledge.	809 810 811 812 813
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need . In <i>Advances in Neural Information Processing Systems 30</i> , pages 5998–6008.	814 815 816 817 818
	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mt5: A massively multilingual pre-trained text-to-text transformer . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 483–498, Online. Association for Computational Linguistics.	819 820 821 822 823 824 825 826
	Xuan Zhang, Pamela Shapiro, Gaurav Kumar, Paul McNamee, Marine Carpuat, and Kevin Duh. 2019. Curriculum learning for domain adaptation in neural machine translation . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1</i> , pages 1903–1915.	827 828 829 830 831 832 833
	A Probe Token Taxonomy	834
	English lexical probes	835
	A token τ is classified as English lexical if either condition holds:	836 837
	1. Irregular past-tense verb. $\text{UPOS}(\tau) = \text{VERB}$, $\text{Tense}=\text{Past} \in \text{FEATS}(\tau)$, $\text{VerbForm}=\text{Fin} \in \text{FEATS}(\tau)$, and $\text{LEMMA}(\tau) \in \mathcal{V}_{\text{irr}}^{\text{EN}}$.	838 839 840
	2. Irregular plural noun. $\text{UPOS}(\tau) = \text{NOUN}$ and $\text{TEXT}(\tau) \in \mathcal{N}_{\text{irr}}^{\text{EN}}$.	841 842

German lexical probes

German lexical probes are operationalized as item-specific or lexically exceptional forms:

1. **Strong past-tense verbs.** $\text{UPOS}(\tau) = \text{VERB}$, $\text{Tense}=\text{Past} \in \text{FEATS}(\tau)$, $\text{VerbForm}=\text{Fin} \in \text{FEATS}(\tau)$, and $\text{LEMMA}(\tau) \in \mathcal{V}_{\text{irr}}^{\text{DE}}$.
2. **Suppletive or highly irregular finite paradigm forms.** $\text{TEXT}(\tau) \in \mathcal{S}_{\text{DE}}$ and $\text{UPOS}(\tau) \in \{\text{VERB}, \text{AUX}\}$.
3. **Irregular plural nouns.** $\text{UPOS}(\tau) = \text{NOUN}$, $\text{Number}=\text{Plur} \in \text{FEATS}(\tau)$, and $\text{TEXT}(\tau) \in \mathcal{N}_{\text{irr}}^{\text{DE}}$.
4. **Gender-marked determiners.** $\text{UPOS}(\tau) = \text{DET}$ with nominative singular definite or indefinite gender-marked forms. These probes are grouped with lexical probes because they depend on lexically specified noun-gender distinctions, even though they also involve morphosyntactic marking.

English grammatical probes

1. **Regular past-tense verb.** $\text{UPOS}(\tau) = \text{VERB}$, $\text{Tense}=\text{Past} \in \text{FEATS}(\tau)$, $\text{VerbForm}=\text{Fin} \in \text{FEATS}(\tau)$, and $\text{LEMMA}(\tau) \notin \mathcal{V}_{\text{irr}}^{\text{EN}}$.
2. **Auxiliary agreement.** $\text{UPOS}(\tau) = \text{AUX}$ with number-marked auxiliary agreement.

German grammatical probes

1. **Weak past-tense verb.** $\text{UPOS}(\tau) = \text{VERB}$, $\text{Tense}=\text{Past} \in \text{FEATS}(\tau)$, $\text{VerbForm}=\text{Fin} \in \text{FEATS}(\tau)$, $\text{LEMMA}(\tau) \notin \mathcal{V}_{\text{irr}}^{\text{DE}}$, and the form is not part of the suppletive set.
2. **Auxiliary agreement.** $\text{UPOS}(\tau) = \text{AUX}$ with number-marked auxiliary agreement.

Syntactic probes

1. $\text{UPOS}(\tau) = \text{SCONJ}$, or
2. $\text{DEPREL}(\tau) \in \{\text{advcl}, \text{ccomp}, \text{xcomp}, \text{acl}, \text{csubj}\}$.

These labels should be interpreted as operational probe families rather than perfectly discrete theoretical classes. Their purpose is to test whether MoE routing is sensitive to broad linguistically motivated contrasts under a consistent parser-based annotation scheme.

B Full Curriculum Schedule

Category sampling

Within every epoch, lexical, grammatical, and syntactic sentences are sampled uniformly and shuffled randomly.

Forward curriculum

Epochs	λ_{DE}	Description
1–8	0.00	EN only
9	0.20	L2 introduction
10	0.25	
11	0.30	
12	0.35	
13	0.40	
14	0.45	
15–16	0.50	Stable bilingual

Table 7: Forward curriculum L2 exposure schedule.

C Hyperparameter Details

Hyperparameter	Value
Layers	8
Embedding size	512
Attention heads	8
Feed-forward size	2048
MoE layers	1, 3, 5, 7
Experts N	8
Top- k	3
Load-balance α	0.01
Vocabulary	mBERT WordPiece, $\sim 119\text{k}$
Optimizer	AdamW
Peak LR	3×10^{-4}
Warmup steps	500
Batch size	8 sequences
Gradient clip	1.0
Epochs	16

Computing Infrastructure and Budget

All training, validation, and routing-analysis experiments were conducted on NVIDIA T4 GPUs and NVIDIA RTX A6000 GPUs. The primary bilingual MoE model used in the study contains 86.4M parameters. Each training condition was run for 16 epochs with 300,000 sentence instances per epoch, corresponding to 4.8M sentence instances per condition. Validation, routing extraction, and downstream statistical analyses were performed from saved checkpoints and held-out routing logs on the same hardware class. Depending on hardware availability, experiments were conducted on

906	either NVIDIA T4 or NVIDIA RTX A6000 GPUs.	number forms, German strong verbs were weak-	955
907	Training time therefore varied across runs, ranging	regularized, German suppletive forms were re-	956
908	from approximately 14–28 GPU-hours per condi-	placed with present/agreement-incompatible alter-	957
909	tion	natives, and subordinating conjunction probes were	958
910	D Additional Routing Analysis Details	removed to create a syntactically degraded variant.	959
911	Permutation testing	This yielded minimal-pair contrasts in which S^+	960
912	Category labels are shuffled at the sentence level for	remained parser-validated and S^- differed by a	961
913	1,000 permutations per layer. Empirical p -values	controlled probe-level manipulation.	962
914	are computed as the fraction of permutations with	Below we provide two illustrative prompt exam-	963
915	$I(R; C) \geq I_{\text{observed}}$. This tests whether the ob-	ples adapted directly from the generation templates	964
916	erved routing-category dependence is larger than	used in the pipeline.	965
917	would be expected under category-independent	Prompt Example 1: English irregular verb	966
918	routing.	probe.	967
919	E OpenAI Generation Details, Validation	Write exactly 10 EN sentences.	968
920	Parsing, and Prompt Examples	Subtype: irregular_verb	969
921	The held-out validation sentences were generated	Preferred lemmas: become, draw, awake, speak,	970
922	with the OpenAI Responses API using the Python	throw, choose, write, drive, sing	971
923	OpenAI client. In the generation script, the de-	Avoid recent probes: none	972
924	fault model was GPT-5.1, configurable through	Rules: - Output exactly 10 lines, no more and	973
925	the OPENAI_MODEL environment variable. Re-	no fewer. - One sentence per line. - 6–12 words	974
926	quests were made with max_output_tokens=420,	per sentence. - Prefer a different listed lemma	975
927	reasoning.effort="none", store=false, re-	on each line. - Mark target as [PROBE: word]. -	976
928	quest timeout 40s, and max_retries=0. The API	Prefer the correct irregular simple-past form of	977
929	call also included a fixed instruction string: “ <i>Fol-</i>	one listed lemma. - Probe must be the main verb.	978
930	<i>low the output format exactly. Return only the</i>	- Use simple past only. - No auxiliaries or partici-	979
931	<i>requested lines and nothing else.</i> ”	ples with the probe. - Use the real irregular past	980
932	Generated candidates were then automatically	form, not the base form. - Examples: become →	981
933	parsed and validated with Stanza using the proces-	became, awake → awoke, draw → drew. - Prefer	982
934	sors tokenize, pos, lemma, and depparse. For	the listed lemmas strongly. - If one listed lemma	983
935	each generated S^+ sentence, the pipeline checked	feels awkward, use another listed lemma. - No	984
936	that the marked probe token occurred exactly once,	numbering, bullets, quotes, parentheses, or extra	985
937	that the output matched the required subtype in-	text. - Bare sentence text only.	986
938	ventory, and that the parsed probe token satisfied	Prompt Example 2: German subordinating con-	987
939	the subtype-specific constraints used in dataset con-	junction probe.	988
940	struction (e.g., finite past-tense verb, plural noun,	Write exactly 10 DE sentences.	989
941	auxiliary, or subordinating conjunction, depending	Subtype: sconj	990
942	on the subtype). Additional filters removed mal-	Preferred subordinating conjunctions: weil, ob-	991
943	formed outputs, duplicates, and subtype violations	wohl, bevor, nachdem, falls, sobald, damit, dass	992
944	before acceptance into the held-out set.	Avoid recent probes: none	993
945	For the final validation pair construction, the un-	Rules: - Output exactly 10 lines, no more and no	994
946	grammatical sentence S^- was obtained determinis-	fewer. - One sentence per line. - 6–12 words per	995
947	tically from the validated S^+ sentence by editing	sentence. - Prefer a different listed conjunction	996
948	only the probe token or deleting the probe in the	on each line. - Mark target as [PROBE: word]. - Pre-	997
949	syntactic conjunction cases. The transformation	fer one listed subordinating conjunction exactly	998
950	rule depended on subtype. For example, English	as shown. - Use the probe as a true subordinat-	999
951	irregular past forms were replaced with regular-	ing conjunction with UPOS=SCONJ. - Build a	1000
952	ized forms, English regular past-tense verbs were	clear subordinate-clause construction. - Keep the	1001
953	replaced with their lemma/base form, auxiliary	sentence fully grammatical as written. - Avoid	1002
954	agreement probes were swapped to mismatching	coordinators, adverbs, or discourse markers. - No	1003
		numbering, bullets, quotes, parentheses, or extra	1004
		text. - Bare sentence text only.	1005
		Parsing and Validation Tools	1006
		Parser-based probe annotation and validation were	1007
		performed with Stanza. We used the Stanza	1008
		pipelines for English and German with the proces-	1009
		sors tokenize, pos, lemma, and depparse. These	1010
		pipeline outputs were used to identify probe tokens,	1011

verify subtype constraints, and validate generated held-out examples.

F Dataset Diversity Statistics

To address concerns regarding dataset health and template collapse, we provide diversity statistics for the diagnostic set. All reported probe sets maintain 100% sentence-pair uniqueness.

Probe Set	Records	Tokens	Uniqueness
EN Irregular Verb	1500	46	100%
EN Irregular Plural	1500	30	100%
DE Strong Verb	500	35	100%
EN SCONJ	3000	21	100%
EN Aux Agreement	1500	17	100%

Table 8: Illustrative diagnostic probe diversity statistics.

G Ablation Controls

The appendix currently includes four additional controls beyond the main forward and no-curriculum MoE comparisons: random routing, frozen routing, top- k , and load-balancing ($\alpha = 0.005$) ablations under the forward curriculum. These controls are intended to distinguish curriculum effects from generic sparse-routing effects and to test how sensitive the observed specialisation is to router flexibility and routing sparsity.

Ablation	Combined	EN	DE
Random routing	0.0018	0.0029	0.0057
Frozen routing	0.0274	0.0647	0.0140
Load balancing ($\alpha = 0.005$)	0.1482	0.3288	0.1042
Forward top- $k = 2$	0.1100	0.1767	0.1349
Forward top- $k = 4$	0.2382	0.4107	0.1144

Table 9: Routing-control and sparsity ablations evaluated on the held-out validation set (MI@L5). Random routing collapses routing-category dependence toward near-null values, indicating that specialisation does not arise from architectural sparsity alone. Freezing router parameters preserves weaker but non-zero specialisation. Lower load balancing increases routing-category dependence, while the routing-sparsity ablations show that top- $k = 2$ yields a similar but slightly lower pooled MI than the main model, whereas top- $k = 4$ produces substantially stronger category-conditioned routing, especially for English probes.

The random-routing condition replaces learned router assignments with uniformly sampled expert selection while preserving the underlying expert parameters. As expected, routing-category mutual information collapses to near-zero values, indicating that the observed specialisation patterns require

adaptive routing rather than arising from architectural sparsity alone.

The frozen-routing condition initializes routing normally but prevents subsequent router optimization during training. Although specialisation is substantially weaker than in the fully trainable model, non-trivial routing-category dependence remains, suggesting that expert differentiation can emerge through changes in token representations interacting with fixed routing boundaries.

Load-balancing strength additionally exerts a strong influence on specialisation structure. Reducing the auxiliary load-balancing coefficient to $\alpha = 0.005$ markedly increases routing-category dependence, indicating that weaker balancing constraints permit stronger expert partitioning. However, this increase coincides with increasingly skewed expert utilisation, suggesting a trade-off between specialisation strength and balanced expert participation.

Routing sparsity also affects specialisation, but not in the initially expected direction. Reducing expert selection from top- $k = 3$ to top- $k = 2$ leaves pooled routing-category dependence at a similar level, with a slight decrease overall. In contrast, increasing routing breadth to top- $k = 4$ yields a substantially stronger routing-category association, especially for English probes. In the present runs, broader routing therefore coincides with stronger category-conditioned separation rather than weaker specialisation.

G.1 Frequency-Binned Lexical Analysis

To examine whether lexical routing specialisation is driven primarily by token frequency, English lexical probes were grouped into frequency bins using curriculum-weighted lemma frequencies. Mutual information was then recomputed using only lexical irregular-versus-regular contrasts within each frequency bin.

Frequency Bin	Forward MI@L5	No-Curr MI@L5
Frequent	0.0149	0.0055
Medium	0.0206	0.0976
Rare	0.2261	0.2494

Table 10: English lexical irregular-versus-regular routing mutual information at layer 5 after binning probe lemmas by curriculum-weighted frequency.

Mutual information remains non-zero across all frequency bins in both training conditions. Notably, the strongest routing-category association is observed for rare lexical items rather than the

1077
1078
1079

most frequent items. This suggests that the lexical routing effect cannot be explained solely by the highest-frequency lexical forms.

Layer	Forward $I(R; C L)$	No-Curr $I(R; C L)$
1	0.0935	0.2584
3	0.0639	0.1084
5	0.1235	0.3787
7	0.0312	0.0895

Table 11: Conditional mutual information between routing and linguistic category given language, $I(R; C | L)$, across routed layers. In both the forward-curriculum and no-curriculum conditions, conditional MI remains clearly non-zero and permutation-significant at all routed layers ($p = 0.000999$), indicating that category-sensitive routing is not reducible to language separation alone.

We additionally examined the evolution of routing specialisation during training using epoch-wise routing logs. These analyses are provided here as supplementary evidence.

G.2 Routing Specialisation Dynamics During Training

To analyse routing specialisation throughout training, mutual information and entropy were computed from the routing logs at checkpoints from each training epoch. Figure 5 shows that category-dependent routing emerges rapidly during the early stages of training and stabilises over the following epochs.

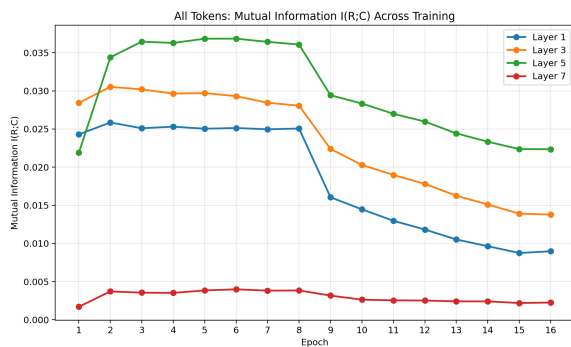


Figure 5: Evolution of routing-category mutual information across training epochs for the routed layers.

The relative ordering of the routed layers remains largely consistent throughout training, with layer 5 exhibiting the highest mutual information across all epochs. The figure also shows a decline in category-dependent routing after epoch 8. Following this transition, the mutual information de-

1080
1081
1082
1083

1084
1085

1086
1087
1088
1089
1090
1091
1092

creases across all routed layers before stabilising again during the later stages of training.

1099
1100

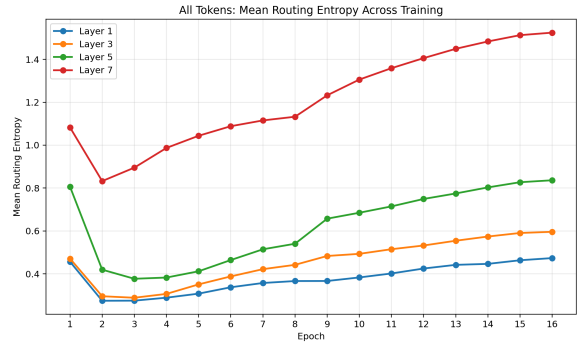


Figure 6: Evolution of routing entropy across training epochs for the routed layers.

Figure 6 presents the routing entropy across the routed layers throughout training. A complementary trend is observed between mutual information and entropy. As mutual information decreases after epoch 8, routing entropy increases across all routed layers during the same period. This indicates that expert utilisation becomes progressively more distributed while category-dependent routing behaviour remains present.

1101
1102
1103
1104
1105
1106
1107
1108
1109