

# Shared-Probe Priors for Diagnosing and Guarding Against Expert-Routing Collapse in Multilingual ASR

Anonymous ACL submission

## Abstract

Language-specific LoRA experts make multilingual ASR parameter-efficient, but they also turn language choice into a latent inference-time decision when labels are absent or unreliable. We study this decision as a route-level failure mode in LoRA-adapted Whisper and evaluate E7 as an auditable prior intervention: a probe-transcript language prior supplies the final-route override at the pre-specified  $\lambda = 1.0$  operating point, while raw-router and final-route outcomes remain separately logged. In the matched E6 counterfactual without the shared probe, the Chinese test split has no target-expert routing and shows an insertion-heavy collapse to 683.33 CER; applying the E7 prior override recovers 99.87% Chinese target routing and reduces CER to 7.03. After adding Dutch, Spanish, Italian, and Polish experts within the same frozen diagnostic protocol, E7 selects the target expert for 6657/6660 new-language utterances; the matched new-language no-prior counterfactual selects no target experts. Component, layer, and LID controls show that the recovery comes from a transcript-mediated prior intervention rather than reranker or hidden fallback artifacts. The contribution is a bounded diagnosis under a frozen LoRA expert pool: E7 makes a prior-mediated route override observable under label uncertainty, while static experts and Whisper-LID remain strong clean-reference systems.

## 1 Introduction

Modern multilingual ASR systems must support existing languages while adding new ones without retraining a full model for each deployment. Whisper and large multilingual representation learners provide strong cross-lingual initialization (Radford et al., 2023; Babu et al., 2022; Zhang et al., 2023). Still, efficient adaptation remains delicate: full fine-tuning is costly, and adding languages can perturb existing behavior. LoRA reduces this cost

by freezing the base model and learning low-rank updates (Hu et al., 2022). Recent ASR work uses shared LoRA experts or experts tied to language, speaker, and accent factors (Song et al., 2024; Li et al., 2025; Bagat et al., 2025). When language labels are absent or untrusted, however, the system must still infer which expert to use.

This selection step can fail more severely than the recognizer itself. In our matched counterfactual without a shared probe, Chinese utterances are routed to the wrong expert and produce an insertion-heavy CER collapse. Aggregate WER or CER alone cannot distinguish route choice, decoder looping, and expert quality. We therefore study expert selection as a latent decision and pair ASR error rates with route diagnostics: target-route distributions, prior sweeps, oracle gaps, and output-length checks. This follows a broader sparse-expert lesson: learned gates can develop collapse pathologies and often need stabilization (Shazeer et al., 2017; Chi et al., 2022; Fedus et al., 2022; Zoph et al., 2022; Zuo et al., 2022).

We study Shared-Probe Prior Intervention as a frozen diagnostic protocol for a fixed LoRA expert pool. E7 obtains an auxiliary shared-adapter probe transcript, estimates a language prior, and at the reported  $\lambda = 1.0$  operating point uses this prior as the final selected-expert route while retaining raw encoder-router scores for audit. Unlike a clean LID gate, the protocol separates the failed raw route, the auditable prior intervention, and the final expert choice; the object of study is utterance-level routing collapse, not a new full-capacity MoE architecture.

Our key contributions are as follows:

- We identify expert-routing collapse as a distinct failure mode in LoRA expert inference without trusted language labels; in a matched counterfactual, Chinese test utterances never select the Chinese expert.

084	• We introduce a route-audit protocol that pairs	2022). In ASR, LR-MoE uses shared LID routing	130
085	WER and CER with target-route rates, expert	to dispatch language-specific experts (Wang	131
086	distributions, oracle gaps, and output-length	et al., 2023). BLR-MoE separates router con-	132
087	checks to separate routing failure from ordi-	fusion from ASR confusion through router aug-	133
088	nary recognition error.	mentation and expert pruning (Ma et al., 2025),	134
		and Switch Conformer combines sparse language	135
089	• We bound the shared-probe prior with reten-	experts with universal phonetic experts (Mimura	136
090	tion and onboarding tests, prior sweeps, con-	et al., 2025). Selective-invocation work shows that	137
091	trols, and LID comparisons, reporting accu-	LID routing can add cost and suffer from misclas-	138
092	racy, latency, and auditability trade-offs.	sification (Xue et al., 2025).	139
093	The contribution is therefore a routing-collapse	<b>2.4 Diagnostics and Evaluation</b>	140
094	diagnosis and auditability study, not a leaderboard-	Layer-wise analyses show that acoustic and lin-	141
095	style ASR comparison.	guistic information changes with depth in speech	142
		representations (Pasad et al., 2021). Confidence	143
096	<b>2 Related Work</b>	and calibration diagnostics matter when ASR de-	144
		isions depend on posterior scores (Woodward	145
097	<b>2.1 Multilingual ASR</b>	et al., 2020; Guo et al., 2017), and NLP signifi-	146
098	Whisper, XLS-R, and USM show the value of mul-	cance guidance motivates paired tests with com-	147
099	tilingual pretraining at scale (Radford et al., 2023;	compact route-audit diagnostics (Dror et al., 2018).	148
100	Babu et al., 2022; Zhang et al., 2023). Work on	Recent work also highlights multilingual ASR	149
101	parameter-efficient extension asks how to add lan-	metric pitfalls, beyond-WER Whisper diagnostics,	150
102	guages without full retraining (Liu et al., 2024).	and group-fairness disparities in multilingual ASR	151
103	We focus on frozen Whisper-LoRA experts, where	(Manohar and Pillai, 2024; Liang et al., 2025; Zee	152
104	the central question is how to select language-	et al., 2024).	153
105	specific updates when the language label is un-		
106	available or untrusted.	<b>3 Setup and Route Audit</b>	154
107	<b>2.2 LoRA Experts</b>	<b>3.1 LoRA Expert Pool</b>	155
108	LoRA adapts a frozen model with compact low-	The experiments use Whisper-small as the back-	156
109	rank updates, avoiding full model copies (Hu	bone. Language-specific LoRA experts use rank	157
110	et al., 2022). In ASR, LoRA-Whisper uses one	32, alpha 64, dropout 0.05, and updates to $q_{proj}$	158
111	adapter per language (Song et al., 2024); LoRA	and $v_{proj}$ in encoder and decoder attention.	159
112	language experts and LoRA-MoLE fuse or dis-	Each adapter has 3,538,944 trainable parameters	160
113	still monolingual experts for inference without a	(3.54M in tables), verified from PEFT. The E7	161
114	trusted language label (Li et al., 2025). HLoRA	shared-probe adapter uses the same counted unit.	162
115	uses LID-posterior LoRA routing in a CTC ASR	E4 uses one shared multilingual LoRA adapter,	163
116	stack (Zheng et al., 2026), and mixture-of-LoRA	E5 uses static language-specific experts, and new-	164
117	work extends the idea to speaker or accent spe-	language onboarding inserts one added expert per	165
118	cialization (Zhao et al., 2024; Bagat et al., 2025).	new language.	166
119	Layer-aware and rank-level LoRA variants study		
120	where sharing and specialization should occur dur-	<b>3.2 Expert Selection</b>	167
121	ing adaptation (Xu et al., 2026; Mei et al., 2026);	E6 and E7 build on the E5 experts by adding	168
122	our focus is routing-time diagnosis under label un-	inference-time routing. Given an utterance, the	169
123	certainty.	system extracts encoder features, scores candi-	170
		date experts, and decodes with the selected LoRA	171
124	<b>2.3 Language Routing and MoE</b>	expert. E6 removes the shared-probe prior and	172
125	Sparse MoE systems route each input to a small	serves as the matched no-prior counterfactual; E7	173
126	set of experts (Shazeer et al., 2017; Fedus et al.,	applies the prior intervention introduced below.	174
127	2022), but representation collapse, training stabil-	Because E7 applies a text-based classifier to the	175
128	ity, and transfer stability remain central design is-	probe transcript and fixes $\lambda = 1.0$ , we describe	176
129	sues (Chi et al., 2022; Zoph et al., 2022; Zuo et al.,	it as an audit-path route override with an auxil-	177

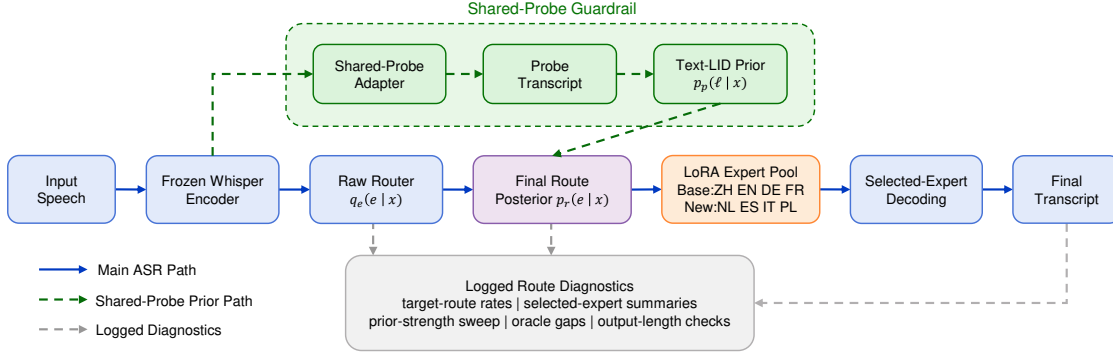


Figure 1: Shared-probe prior intervention and route diagnostics. The raw router remains separately logged, while the probe-derived prior supplies the final route at the audit operating point and the audit path logs route-level evidence.

178 iary transcript-LID prior. Unlike the Whisper-LID  
 179 comparator, E7 exposes raw route, prior interven-  
 180 tion, and final route traces; its value is diagnostic  
 181 visibility under label uncertainty, not replacing ex-  
 182 plicit LID in clean accuracy or latency.

### 183 3.3 Collapse Diagnostics

184 The diagnostics use selected-expert summaries,  
 185 target-route rates, oracle gaps, and output-length  
 186 checks, with cached route traces for collapse re-  
 187 runs. We organize them around the E1–E7 lad-  
 188 der. E1–E4 are non-routed or shared-adapter refer-  
 189 ences, E5 is the clean-label static-expert reference,  
 190 E6 isolates the learned encoder router without the  
 191 shared-probe prior, and E7 applies the prior while  
 192 keeping the expert pool fixed.

193 We label a case as routing collapse only when  
 194 route evidence and recognition symptoms agree:  
 195 target routing is missing or badly degraded, one  
 196 wrong expert dominates the selected-expert his-  
 197 togram, and oracle or length diagnostics show that  
 198 a better route was available. If target routing is re-  
 199 stored but error remains, we treat the residual as  
 200 post-routing ASR headroom, not another collapse  
 201 signal.

## 202 4 Shared-Probe Prior Intervention

### 203 4.1 Transcript Prior

204 The shared-probe path estimates the language ev-  
 205 idence used by the reported route override before  
 206 the final selected-expert decode. As shown in Fig-  
 207 ure 1, the raw router is trained from pooled Whisper  
 208 encoder states from layers 6 and 8. We use  
 209 these middle layers because speech cues vary with  
 210 depth (Pasad et al., 2021) and the layer-pair con-  
 211 trol in Appendix Table A2 gives the best displayed  
 212 Base Avg<sup>†</sup> for the tested layer-6 setting. E7 ob-  
 213 tains a shared-adapter probe transcript and uses

214 `langid.py 1.1.6 (Lui and Baldwin, 2012)` to form  
 215 an auditable text-LID prior.

216 At the reported  $\lambda = 1.0$  operating point, the  
 217 shared probe supplies the deciding transcript-level  
 218 evidence, preserving a clean boundary between  
 219 raw router, probe prior, and final selected route.

### 220 4.2 Prior-Guided Selection

221 The encoder router remains logged for diagnosis,  
 222 but at the reported  $\lambda = 1.0$  operating point the fi-  
 223 nal E7 route is a logged intervention by the probe-  
 224 derived prior. Let  $q_e(e|x)$  be the raw encoder-  
 225 router posterior for expert  $e$ ,  $m(e)$  map expert  $e$   
 226 to its language,  $p_p(\ell|x)$  be the probe-derived lan-  
 227 guage posterior for language  $\ell$ , and  $\mathcal{E}$  be the expert  
 228 set. E7 computes:

$$\begin{aligned}
 p_p(\ell|x) &= \text{LangID}(\text{ProbeASR}(x)), \\
 q_e(e|x) &= \text{softmax}_e(\text{Router}(x)), \\
 p_r(e|x) &\propto (1 - \lambda) q_e(e|x) + \lambda p_p(m(e)|x).
 \end{aligned}$$

229 The final line is normalized over  $e \in \mathcal{E}$ . Reported  
 230 E7 runs fix  $\lambda = 1.0$  as a pre-specified audit point  
 231 rather than a tuned interpolation: the final route is  
 232 intentionally determined by the probe prior, while  
 233 the raw router remains logged. Table 5 reports the  
 234 frozen full-test sweep: weak prior fusion leaves  
 235 Chinese routing unusable, whereas the collapsed  
 236 route is avoided only once the probe prior domi-  
 237 nates. Intermediate rows are diagnostics, not  
 238 model-selection candidates; near-identical strong-  
 239 prior rows indicate a stable override regime.

### 241 4.3 Controls

242 E6 and E7 use the same router-training recipe: ex-  
 243 pert supervision, class-balanced sampling, temper-  
 244 ature scaling, and load-balancing regularization  
 245 following standard sparse-expert practice (Zoph  
 246 et al., 2022). The recipe is fixed across the compar-  
 247 ison and is not an independent contribution.

Role	Lang	Training sources	Train	Eval	Test
Base	ZH	AISHHELL-1 150.9 + WenetSpeech S 100.1	250.9	18.1	10.0
	EN	LibriSpeech 100.6 + GigaSpeech S 250.0	350.6	5.4	5.4
	DE	MLS 150.0 + VoxPopuli 274.5	424.5	14.3	14.3
	FR	MLS 150.0 + VoxPopuli 215.6	365.6	10.1	10.1
New	NL	MLS 50.0 + VoxPopuli 46.4	96.4	12.8	12.8
	ES	MLS 49.9 + VoxPopuli 50.0	99.9	5.0	10.0
	IT	MLS 49.7 + VoxPopuli 50.0	99.7	5.0	2.5
	PL	MLS 50.6 + VoxPopuli 50.0	100.6	5.0	2.5

Table 1: Frozen data protocol after filtering. Durations are shown in hours and rounded to one decimal. Base languages support the collapse/diagnosis study; NL/ES/IT/PL are added only for the extension experiment.

controls rule out three alternative explanations for the observed recovery. First, the E6 and E7 boundary isolates the shared-probe prior while keeping the frozen experts and router recipe fixed. Second, prior sweeps, layer-pair checks, hard and soft routing checks, and no-reranker checks test whether the behavior comes from the logged prior intervention rather than a tuned endpoint, a special router layer, or hidden decoding fallback. Third, capacity, resource, and Whisper-LID controls bound the claim against stronger clean-label or explicit-LID references. Together, these controls support a narrow conclusion: E7 provides an auditable prior intervention for an observed routing collapse under the frozen protocol.

## 5 Experiments

### 5.1 Experimental Setup

Table 1 reports realized train, evaluation, and test hours for each language from frozen project splits. Chinese data come from AISHELL-1 (Bu et al., 2017) and WenetSpeech (Zhang et al., 2022); English from LibriSpeech (Panayotov et al., 2015) and GigaSpeech (Chen et al., 2021); German and French from MLS (Pratap et al., 2020) and VoxPopuli (Wang et al., 2021). The protocol fixes data manifests and sampling seeds before training. The main model runs use seed 42, as specified in Section 5.2. Sources differ in domain, speaker mix, and scale, so routing diagnostics are interpreted with WER and CER rather than as domain-independent behavior. For subset-based MLS and VoxPopuli data, speaker and duration distributions and source ratios are frozen before final evaluation. Main tables use only frozen test splits; development and pilot diagnostic runs are excluded. Leakage checks deduplicate splits by utterance ID, audio path, and speaker ID.

### 5.2 Implementation Details

All experiments used one NVIDIA A40 GPU. The stack is PyTorch 2.1, Hugging Face Transformers 4.36, PEFT 0.6, Datasets 2.16, and langid.py 1.1.6 (Lui and Baldwin, 2012). Main LoRA runs use effective batch size 32, learning rate 1e-4 with 500 warm-up steps and cosine decay to 1e-6, AdamW with weight decay 0.01, 20 epochs with patience-3 early stopping, rank 32, alpha 64, dropout 0.05, gradient clipping at 1.0, and seed 42 unless an ablation states otherwise.

For E7, the probe transcript is decoded with the shared LoRA adapter using one beam and automatic decode-language selection. langid is restricted to the candidate expert languages; at  $\lambda = 1.0$ , the final argmax is determined by the probe-derived prior, while the raw router posterior is logged only for audit. Decoder-side posterior-weighted language conditioning (weight 0.2) and confidence floors (0.60 base, 0.45 extension) were fixed before test decoding and are not the interpolation weights in Table 5. The independent frozen rerun that reproduced the E7 Base Avg<sup>†</sup> is reported only as a reproducibility check; all prediction records store raw-route probabilities, probe priors, final selected experts, and decoding hashes.

### 5.3 Systems

Table 2 defines the controlled base-language comparison. The rows separate system family, expert-selection rule, and stored-parameter cost: raw Whisper-small and full fine-tuning baselines, a shared multilingual LoRA, gold-language static LoRA experts (E5), Ours w/o shared-probe (E6), Ours with the shared-probe prior (E7), and a Whisper-LID gate that maps detect\_language to an E5 expert. Under fixed data, expert pools, and decoding settings, the table tests whether routing decisions—rather than corpus or decoder changes—explain the E6/E7 gap. E6 stores the four E5 language experts plus the 9.2K encoder router; E7 adds one 3.54M shared-probe LoRA adapter, so its stored count is 3.54M×5 + 9.2K rather than 3.54M×4 + 9.2K. Resource, capacity, component, layer, and stress controls are reported in the appendix.

### 5.4 Metrics

Decoding is controlled separately from expert selection. All main comparisons use the recorded offline Whisper-small snapshot, beam size 5, fixed

ID	System	Selector	Stored Params	ZH	EN	DE	FR	Avg <sup>†</sup>
E1	Whisper-small raw	no	–	20.57	34.77	13.26	16.83	21.36
E2	Full-FT multilingual	shared model	240.6M	8.58	6.52	11.04	10.77	9.23
E3	Full-FT monolingual	gold language	240.6M×4	8.58	5.23	<b>10.81</b>	<b>10.10</b>	<b>8.68</b>
E4	Shared multilingual LoRA	shared adapter	3.54M	8.82	<b>4.67</b>	12.75	12.46	9.67
E5	Static LoRA experts	gold language	3.54M×4	<b>6.85</b>	5.55	11.34	11.40	8.79
E6 <sup>‡</sup>	Ours w/o shared-probe	encoder router	3.54M×4 + 9.2K	683.33	5.54	33.91	11.49	183.57
E7	Ours: shared-probe prior	shared-probe prior	3.54M×5 + 9.2K	7.03	5.55	11.73	11.38	8.92
LID	Whisper-LID + E5 expert	Whisper-LID gate	3.54M×4	6.86	5.55	11.79	11.38	8.90

Table 2: Base-language clean-test results under the frozen protocol. ZH uses CER; EN/DE/FR use WER; Base Avg<sup>†</sup> is an unweighted descriptive average. Bold marks the best value among leaderboard/reference rows, excluding the E6 diagnostic row. <sup>‡</sup> marks the seed-42 full-test no-shared-probe collapse; full-test seed-43/44 checks are discussed in Section 8.

task setting, language-token policy, text normalization, generation temperature, and decoding configuration hashes stored with the result JSONs. Static experts, routed experts, and the Whisper-LID comparator therefore differ in adapter/expert selection, not in ASR decoding hyperparameters or oracle prompt information. All reported scores use frozen manifests and splits unless stated otherwise. ZH uses CER; all other languages use WER. Within a language, WER, CER, and oracle gaps are corpus-level micro averages over the displayed split. Base Avg<sup>†</sup> and New-language Avg are unweighted macro summaries across languages, not pooled cross-script ASR metrics. Paired bootstrap tests, where reported, resample aligned utterance-level predictions. Oracle gap is a route-audit margin, not an additional deployment setting: for each utterance, we compare the transcript produced by the selected expert with the lowest-error transcript among the same per-utterance candidate set of available experts. The reported Gap is the selected-route corpus error minus this per-utterance best-available corpus error, using the same metric scale as the table row (CER points for ZH and WER points otherwise). The target-language expert is therefore an operational proxy for the intended route rather than an absolute gold route, and we call collapse only when target routing is missed together with a large oracle gap or output-length pathology.

## 6 Results

### 6.1 Collapse and Prior Recovery

Table 2 reports clean base-language results including the E6 routing-collapse diagnostic. E6 shows an extreme Chinese failure and elevated German error, while the E7 prior path recovers recognition quality near strong clean references. E5 remains slightly better on Base Avg<sup>†</sup> (8.79 vs. 8.92), and Whisper-LID is also highly competitive (8.90).

Lang	Metric	E6: no probe			E7: probe			Diagnosis
		Err	Target	Gap	Err	Target	Gap	
ZH	CER	683.33	0.00	409.65	7.03	99.87	1.00	collapse→prior fix
EN	WER	5.54	100.00	2.76	5.55	99.96	2.77	stable
DE	WER	33.91	69.42	25.27	11.73	99.53	3.19	partial→prior fix
FR	WER	11.49	99.51	1.89	11.38	100.00	1.80	stable

Table 3: Route-level collapse diagnostics. Target is the percentage of utterances assigned to the target-language expert, and Gap is the selected-route corpus error minus the per-utterance best-available candidate corpus error in points; low E6 Target with a large Gap indicates routing failure, while high E7 Target with a small Gap indicates prior-mediated recovery.

Per-language paired tests show that E7 has higher error than E5 on ZH and DE, no reliable EN difference, and a small FR gain. Thus the clean-test result is best read as a successful prior-mediated route recovery near strong label-aware references, not as an accuracy win over static E5. We therefore do not make a formal non-inferiority claim against E5; the claim is prior-mediated route recovery with a small clean-label cost. An independent frozen rerun reproduced the E7 Base Avg<sup>†</sup> value.

Table 3 is included to make the mechanism claim explicit rather than to add another leaderboard. For each language, it asks whether E6’s error is accompanied by missed target-expert routing and a large oracle gap, and whether the E7 prior path recovers target routing before the remaining recognition error is interpreted. This connects Table 2’s aggregate scores to a route-level diagnosis.

Table 3 shows why the Chinese failure should be read as routing collapse rather than ordinary ASR degradation: the no-probe route misses the target expert and leaves a large oracle gap, while the shared-probe prior recovers target selection and leaves only small residual gaps. German shows a milder version of the same mechanism. Thus the table turns the aggregate error pattern into a route-level diagnosis: the main E6 failure is wrong expert selection, while residual E7 errors mainly reflect post-routing ASR limits.

## 6.2 New-Language Onboarding

Table 4 reports the eight-language onboarding view with both base- and new-language blocks. The base-language entries are rerun or retained references for the expanded onboarding protocol: E4 is retrained as the shared LoRA reference for the eight-language setting, E5 is the gold-language expert reference, and E7 reruns the expanded routed pool. Thus the block is a retention and protocol-alignment check rather than a replacement for the four-language leaderboard in Table 2. Under the expanded expert pool, E7 keeps Base Avg<sup>†</sup> at 8.92, indicating no measurable base-language degradation in this frozen protocol.

On new languages, E7 selects the target expert for 6657/6660 utterances. In WER, it is descriptively lower than E4 on all four added languages, reducing macro WER from 15.64 to 12.76 and nearly matching static experts (12.73). Paired-bootstrap checks show that the E7–E4 gains are reliable for IT and PL ( $p < 0.001$ ) and small for NL and ES; the E7–E5 comparison is mostly tied except for a small Spanish gap favoring E5 ( $p = 0.022$ ). Appendix Table A1 gives the matched no-prior counterfactual: without shared-probe, no added language selects its target expert and macro WER rises to 79.05. Raw Whisper-small remains strongest on NL and ES, so the claim is targeted routing stability and retention rather than uniform low-resource superiority.

## 7 Analysis and Ablations

This section keeps two reviewer-facing diagnostics in the body: the prior sweep and the route visualization. The appendix collects paired tests, new-language counterfactuals, layer and capacity checks, component checks, latency, and Whisper-LID stress results.

### 7.1 Prior Strength

Table 5 makes the prior boundary explicit on the frozen full-test split. The target-route block uses one column per language, so the override threshold is visible directly rather than through an ordered language list. Weak interpolation leaves Chinese routing unusable, whereas usable routing returns once the probe prior dominates. We keep  $\lambda = 1.0$  as the named E7 point because it was the pre-specified, probe-prior-dominated operating point: after diagnosing raw-router collapse, the final route is chosen from the probe-derived

prior rather than from a test-set-tuned mixture. The intermediate rows are frozen diagnostics, not model-selection candidates; the strong-prior rows only show that the override is not sensitive to the endpoint choice.

### 7.2 Probe Visualization

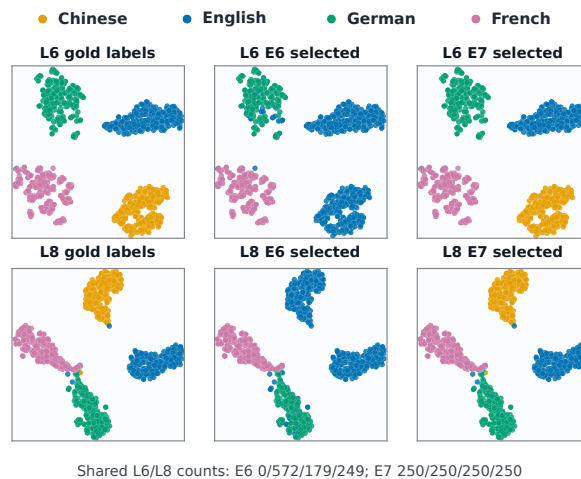


Figure 2:  $2 \times 3$  t-SNE diagnostic over 250 samples per base language from Whisper layers 6 and 8. Columns show gold labels, E6 selected experts, and E7 selected experts; the shared count line reports selected-expert counts in legend order.

Figure 2 visualizes the shared-probe prior intervention as a qualitative route audit: gold-label panels show language structure, while E6/E7 panels contrast route decisions before and after the probe prior. It supports routing-collapse diagnosis and prior-mediated recovery, not perfect language separability.

### 7.3 Cost and Control Boundary

Appendix Table A3 quantifies the cost side of the trade-off. Whisper-LID is the stronger clean deployment reference when only accuracy and latency are considered: Table 2 gives Base Avg<sup>†</sup> 8.90 for Whisper-LID and 8.92 for E7, while Table A3 reports 904.3 versus 3538.2 ms/utterance. The distinction is diagnostic scope. Whisper-LID directly maps an utterance to an E5 expert, whereas E7 keeps the raw-router failure, probe-prior intervention, final route, and route-level diagnosis in the same auditable trace. We therefore do not present E7 as a faster or stronger LID gate; its role is to make the collapse and the prior-mediated repair observable under the frozen protocol.

## 8 Error Analysis

The error analysis separates three mechanisms that a single aggregate score would obscure. Table 3

ID	System	Selector	Stored Params	Base languages					New languages				
				ZH	EN	DE	FR	Avg <sup>†</sup>	NL	ES	IT	PL	Avg
E1	Whisper-small raw	no	–	20.57	34.77	13.26	16.83	21.36	<b>16.80</b>	<b>7.56</b>	17.36	17.77	14.87
E4	Shared multilingual LoRA	shared adapter	3.54M	7.60	5.83	14.21	12.54	10.05	17.70	8.04	17.66	19.15	15.64
E5	Static LoRA experts	gold language	3.54M×8	<b>6.85</b>	<b>5.55</b>	<b>11.34</b>	11.40	<b>8.79</b>	17.56	7.85	<b>11.76</b>	13.75	<b>12.73</b>
E7	Ours: shared-probe prior	shared-probe prior	3.54M×9 + 9.2K	7.02	<b>5.55</b>	11.73	<b>11.38</b>	8.92	17.61	7.94	11.78	<b>13.70</b>	12.76

Table 4: Eight-language onboarding results. ZH uses CER; all other language columns use WER. The E4 base-language cells are from the shared-LoRA rerun in the eight-language onboarding protocol, so they are not intended to duplicate the four-language E4 row in Table 2. E5 is the gold-language expert reference, and E7 is rerun with the expanded expert pool. Base Avg<sup>†</sup> is a descriptive mixed average over ZH CER and EN/DE/FR WER; New Avg is macro WER over NL/ES/IT/PL. Stored Params counts retained LoRA adapters plus the router; – marks no additional stored adapter/router.

Prior Setting	$\lambda$	Recognition error					Target route (%)			
		ZH	EN	DE	FR	Avg <sup>†</sup>	ZH	EN	DE	FR
Raw router	0.00	683.33	5.54	33.91	11.49	183.57	0.0	100.0	69.4	99.5
Fusion	0.25	683.33	5.54	13.58	11.47	178.48	0.0	100.0	96.1	99.6
Fusion	0.50	615.48	5.55	11.77	11.41	161.05	17.0	100.0	99.5	99.8
Fusion	0.75	7.03	5.55	11.75	11.38	8.93	99.9	100.0	99.5	100.0
E7	1.00	7.03	5.55	11.73	11.38	8.92	99.9	100.0	99.5	100.0

Table 5: Shared-probe prior sweep on the frozen full-test split. ZH uses CER; EN/DE/FR use WER. Splitting target-route rates by language shows the prior-strength boundary directly: Chinese remains collapsed through  $\lambda = 0.50$  and enters the strong-prior recovery regime at  $\lambda = 0.75$ . E7 fixes  $\lambda = 1.0$  as a probe-prior-dominated operating point rather than selecting a tuned interpolation.

Seed	$N$	Target (%)	Route hist.	$H$	CER	Gap
42	7176	0.00	EN:7176	0.5583	683.33	409.65
43	7176	0.00	DE:7176	0.5104	2096.19	1554.36
44	7176	0.00	EN:2649/FR:4527	0.5856	1693.67	1147.79

Table 6: Full-test E6 Chinese seed replication. Table 3 establishes the seed-42 route-level mechanism; this table tests whether the zero-target Chinese collapse persists across seeds. Target is the percentage routed to the Chinese expert; Route hist. is the selected wrong-expert histogram; Gap is the oracle CER gap in points. The replicated finding is zero target routing, while the wrong expert and CER magnitude are seed-dependent symptoms.

diagnosed the seed-42 route-level mechanism; Table 6 below asks whether the zero-target Chinese collapse is specific to that seed. The remaining subsections then separate residual post-routing error from language/source limits.

## 8.1 Routing Collapse

Table 6 shows that all three E6 Chinese runs assign 0.00% of 7176 utterances to the Chinese expert. The replicated fact is the absent target route; the dominant wrong expert and CER/gap magnitude are seed-dependent symptoms. This pattern is not consistent with a global label-map permutation: EN and FR remain stable in Table 3, E5 static experts decode the same manifests under trusted labels, and E7 restores Chinese target routing using the same expert indices. Seed-42 logs instead suggest a miscalibrated boundary: ZH is sent to EN despite  $H = 0.558$ , top-1 probability 0.752, top-2 margin 0.504, and coarse ZH posterior mass 0.018. This suggests, but does not prove, calibration pressure from class balance plus script/domain mismatch. The seed-42 Chi-

nese trace is also insertion-heavy: the mean predicted/reference character-length ratio is 6.70, and 711/7176 utterances exceed a ratio of 6. After E7 applies the prior override, the ratio returns to 1.00 with no utterance above 6, so E6 is a controlled no-prior diagnostic rather than a competing ASR system.

## 8.2 Post-Routing Error

The E7–E5 gap is modest but systematic on some languages: +0.18 CER on ZH, -0.00 WER on EN, +0.40 WER on DE, and -0.02 WER on FR. The ZH and DE increases are reliable under paired bootstrap, EN is indistinguishable, and FR is slightly lower for E7. This is not coarse language-ID failure: E7 selects the German expert for 3378/3394 German utterances and the French expert for all 2426 French utterances. The remaining gap is therefore residual recognition cost after inferred expert selection; E5 remains the trusted-label upper baseline.

## 8.3 Language Limits

Stable expert selection is not uniform ASR improvement. Raw Whisper-small is strongest on Dutch and Spanish (16.80/7.56 WER versus E7 17.61/7.94), while Italian and Polish benefit more from expert onboarding (E7 11.78/13.70 versus raw 17.36/17.77 and E4 17.66/19.15). We therefore report per-language scores and retention summaries rather than claiming SOTA multilingual ASR or uniform gains.

## 9 Conclusion

This paper identifies expert-routing collapse as a distinct failure mode in LoRA expert ASR and provides a frozen diagnostic protocol for measuring it. E7 uses a pre-specified probe-transcript prior override to recover stable target routing in the observed collapse setting and after four-language expert expansion, while logging raw-router behavior, prior-sweep behavior, selected-expert traces, and oracle/length diagnostics. The evidence shows that the main failure is not weak expert capacity or ordinary ASR variance, but a latent expert-selection decision that can silently dominate recognition quality. By formulating the raw router, transcript prior, and final route as separate quantities, the method makes this decision observable under the evaluated protocol.

## Limitations

The claim remains deliberately constrained. E7 is a prior-dominated audit path, not evidence that the raw learned router alone is reliable or that auxiliary LID evidence can be removed. Clean-label static experts and Whisper-LID remain strong references, new-language gains are uneven, and the audit path adds substantial latency (Table A3: 3538.2 ms/utterance versus 867.3 ms for E5 and 904.3 ms for Whisper-LID). The hard-vs-soft and no-reranker controls also show that aggregate improvements can come from less interpretable fallback behavior, so route transparency is part of the evaluation target rather than an optional diagnostic.

Future work should reduce probe cost through caching or selective invocation, learn confidence-aware prior weighting, and test learned audio-LID and matched-hour controls before extending the claim beyond the evaluated languages and corpora.

## Ethical Considerations

This work uses public ASR corpora under their released research licenses and collects no new speech. AISHELL-1, WenetSpeech, LibriSpeech, GigaSpeech, MLS, and VoxPopuli are used only through frozen manifests and documented splits; no attempt is made to identify speakers. Per-language scores, routing histograms, stress behavior, and latency limits are reported because recognition quality can vary across languages, domains, speakers, and acoustic conditions. These diag-

nostics are not fairness guarantees: multilingual ASR can show group-level and intersectional performance disparities, so deployment requires additional group-aware auditing (Zee et al., 2024). The experiments are limited to Whisper-small adaptation and small LoRA experts, and we report parameter counts, storage requirements, and batch-1 latency so that readers can assess the method’s resource footprint.

## References

- Arun Babu, Changan Wang, Andros Tjandra, Kushal Lakhotia, Qiantong Xu, Naman Goyal, Kritika Singh, Patrick von Platen, Yatharth Saraf, Juan Pino, Alexei Baevski, Alexis Conneau, and Michael Auli. 2022. *XLS-R: Self-supervised cross-lingual speech representation learning at scale*. In *Proceedings of Interspeech*, pages 2278–2282.
- Raphaël Bagat, Irina Illina, and Emmanuel Vincent. 2025. *Mixture of LoRA experts for low-resourced multi-accent automatic speech recognition*. In *Proceedings of Interspeech*, pages 1143–1147.
- Hui Bu, Jiayu Du, Xingyu Na, Bengu Wu, and Hao Zheng. 2017. *AISHELL-1: An open-source Mandarin speech corpus and a speech recognition baseline*. In *Proceedings of the 20th Conference of the Oriental Chapter of the International Coordinating Committee on Speech Databases and Speech I/O Systems and Assessment*, pages 1–5.
- Guoguo Chen, Shuzhou Chai, Guan-Bo Wang, Jiayu Du, Wei-Qiang Zhang, Chao Weng, Dan Su, Daniel Povey, Jan Trmal, Junbo Zhang, Mingjie Jin, Sanjeev Khudanpur, Shinji Watanabe, Shuaijiang Zhao, Wei Zou, Xiangang Li, Xuchen Yao, Yongqing Wang, Zhao You, and Zhiyong Yan. 2021. *GigaSpeech: An evolving, multi-domain ASR corpus with 10,000 hours of transcribed audio*. In *Proceedings of Interspeech*, pages 3670–3674.
- Zewen Chi, Li Dong, Shaohan Huang, Damai Dai, Shuming Ma, Barun Patra, Saksham Singhal, Payal Bajaj, Xia Song, Xian-Ling Mao, Heyan Huang, and Furu Wei. 2022. *On the representation collapse of sparse mixture of experts*. In *Advances in Neural Information Processing Systems*, volume 35, pages 34600–34613.
- Rotem Dror, Gili Baumer, Segev Shlomov, and Roi Reichart. 2018. *The hitchhiker’s guide to testing statistical significance in natural language processing*. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pages 1383–1392.
- William Fedus, Barret Zoph, and Noam Shazeer. 2022. *Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity*. *Journal of Machine Learning Research*, 23(120):1–39.

639	Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. 2017. <a href="#">On calibration of modern neural networks</a> . In <i>Proceedings of the 34th International Conference on Machine Learning</i> , volume 70 of <i>Proceedings of Machine Learning Research</i> , pages 1321–1330.	695
640		696
641		697
642		698
643		
644		
645	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. <a href="#">LoRA: Low-rank adaptation of large language models</a> . In <i>International Conference on Learning Representations</i> .	699
646		700
647		701
648		702
649		703
650	Jiahong Li, Yiwen Shao, Jianheng Zhuo, Chenda Li, Liliang Tang, Dong Yu, and Yanmin Qian. 2025. <a href="#">Efficient multilingual ASR finetuning via LoRA language experts</a> . In <i>Proceedings of Interspeech</i> , pages 1138–1142.	704
651		705
652		706
653		707
654		
655	Siyu Liang, Nicolas Ballier, Gina-Anne Levow, and Richard Wright. 2025. <a href="#">Beyond WER: Probing Whisper’s sub-token decoder across diverse language resource levels</a> . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 31237–31247, Suzhou, China. Association for Computational Linguistics.	708
656		709
657		710
658		711
659		712
660		713
661		714
662	Wei Liu, Jingyong Hou, Dong Yang, Muyong Cao, and Tan Lee. 2024. <a href="#">A parameter-efficient language extension framework for multilingual ASR</a> . In <i>Proceedings of Interspeech</i> , pages 3929–3933.	715
663		716
664		717
665		718
666		719
667		720
668	Marco Lui and Timothy Baldwin. 2012. <a href="#">langid.py: An off-the-shelf language identification tool</a> . In <i>Proceedings of the ACL 2012 System Demonstrations</i> , pages 25–30, Jeju Island, Korea. Association for Computational Linguistics.	721
669		722
670		723
671		724
672		725
673		
674	Guodong Ma, Wenxuan Wang, Lifeng Zhou, Yuting Yang, Yuke Li, and Binbin Du. 2025. <a href="#">BLR-MoE: Boosted language-routing mixture of experts for domain-robust multilingual E2E ASR</a> . In <i>Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing</i> , pages 1–5.	726
675		727
676		728
677		729
678		730
679		731
680		732
681		733
682		
683		
684	Kavya Manohar and Leena G. Pillai. 2024. <a href="#">What is lost in Normalization? Exploring pitfalls in multilingual ASR model evaluations</a> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 10864–10869, Miami, Florida, USA. Association for Computational Linguistics.	734
685		735
686		736
687		737
688		
689	Yuxiang Mei, Delai Qiu, Shengping Liu, Jiaen Liang, and Yanhua Long. 2026. <a href="#">Zipper-LoRA: Dynamic parameter decoupling for speech-LLM based multilingual speech recognition</a> . <i>arXiv preprint arXiv:2603.17558</i> .	738
690		739
691		740
692		741
693		742
694		
695		
696		
697		
698		
699		
700		
701		
702		
703		
704		
705		
706		
707		
708		
709		
710		
711		
712		
713		
714		
715		
716		
717		
718		
719		
720		
721		
722		
723		
724		
725		
726		
727		
728		
729		
730		
731		
732		
733		
734		
735		
736		
737		
738		
739		
740		
741		
742		
743		
744		
745		
746		
747		
748		
749		
750		

751 [speech recognition difficulty](#). In *Proceedings of In-*  
752 *terspeech*, pages 2580–2584.

753 Anna Zee, Marc Zee, and Anders Søgaard. 2024.  
754 [Group fairness in multilingual speech recognition](#)  
755 [models](#). In *Findings of the Association for Com-*  
756 *putational Linguistics: NAACL 2024*, pages 2213–  
757 2226, Mexico City, Mexico. Association for Com-  
758 putational Linguistics.

759 Binbin Zhang, Hang Lv, Pengcheng Guo, Qijie Shao,  
760 Chao Yang, Lei Xie, Xin Xu, Hui Bu, Xiaoyu Chen,  
761 Chenchen Zeng, Di Wu, and Zhendong Peng. 2022.  
762 [WenetSpeech: A 10000+ hours multi-domain Man-](#)  
763 [darin corpus for speech recognition](#). In *Proceedings*  
764 *of the IEEE International Conference on Acoustics,*  
765 *Speech and Signal Processing*, pages 6182–6186.

766 Yu Zhang, Wei Han, James Qin, Yongqiang Wang,  
767 Ankur Bapna, Zhehuai Chen, Nanxin Chen, Bo Li,  
768 Vera Axelrod, Gary Wang, Zhong Meng, Ke Hu,  
769 Andrew Rosenberg, Rohit Prabhavalkar, Daniel S.  
770 Park, Parisa Haghani, Jason Riesa, Ginger Perng,  
771 Hagen Soltau, Trevor Strohman, Bhuvana Ramab-  
772 hadran, Tara Sainath, Pedro Moreno, Chung-Cheng  
773 Chiu, Johan Schalkwyk, Françoise Beaufays, and  
774 Yonghui Wu. 2023. [Google USM: Scaling au-](#)  
775 [tomatic speech recognition beyond 100 languages](#).  
776 *arXiv preprint arXiv:2303.01037*.

777 Qiuming Zhao, Guangzhi Sun, Chao Zhang, Mingx-  
778 ing Xu, and Thomas Fang Zheng. 2024. [SAML:](#)  
779 [Speaker adaptive mixture of LoRA experts for end-](#)  
780 [to-end ASR](#). In *Proceedings of Interspeech*, pages  
781 777–781.

782 Yuang Zheng, Yuxiang Mei, Dongxing Xu, Jie  
783 Chen, and Yanhua Long. 2026. [A language-](#)  
784 [agnostic hierarchical LoRA-MoE architecture for](#)  
785 [CTC-based multilingual ASR](#). *arXiv preprint*  
786 *arXiv:2601.00557*.

787 Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du,  
788 Yanping Huang, Jeff Dean, Noam Shazeer, and  
789 William Fedus. 2022. [ST-MoE: Designing sta-](#)  
790 [ble and transferable sparse expert models](#). *arXiv*  
791 *preprint arXiv:2202.08906*.

792 Simiao Zuo, Xiaodong Liu, Jian Jiao, Young Jin Kim,  
793 Hany Hassan, Ruofei Zhang, Jianfeng Gao, and Tuo  
794 Zhao. 2022. [Taming sparsely activated transformer](#)  
795 [with stochastic experts](#). In *International Conference*  
796 *on Learning Representations*.

## 797 A Supplementary Checks

798 This appendix answers likely reviewer questions  
799 while keeping the main text focused on the routing-  
800 collapse story. Table A1 gives the matched no-  
801 prior counterfactual for the new languages, show-  
802 ing whether the new-language gains come from  
803 target-expert routing rather than expert quality  
804 alone. Table A2 compactly reports the key con-  
805 trol evidence for likely artifact explanations: layer

choice, reranker behavior, explicit LID substitu- 806  
tion, expansion to eight experts, and stored capaci- 807  
ty. Table A3 states the resource trade-off so the 808  
audit path is not mistaken for a low-latency de- 809  
ployment claim. Table A4 stress-tests the clean 810  
Whisper-LID reference under degraded audio, ad- 811  
dressing whether an explicit LID gate makes the 812  
routing study unnecessary. 813

## 814 B New-Language Counterfactual

815 Table A1 isolates routing from expert quality for  
816 the added languages. It reports the matched no-  
817 prior counterfactual that the main text summarizes  
818 but does not need to tabulate in the body.

Lang	N	No prior	No-prior route	E7	E7 route
NL	3075	94.14	en:2982/de:93	17.61	nl:3074/en:1
ES	2385	56.22	en:1784/de:536/fr:65	7.94	es:2384/en:1
IT	594	53.40	en:441/fr:146/de:7	11.78	it:593/en:1
PL	606	112.46	en:549/de:57	13.70	pl:606
Macro	–	79.05	0 target	12.76	6657/6660 target

819 Table A1: New-language no-prior counterfactual. Without  
820 the shared-probe prior, no added language selects its target  
821 expert; E7 restores 6657/6660 target routes.

## 822 C Controls, Cost, and Stress

823 Table A2 summarizes the non-core controls by  
824 concern, check, and route-audit takeaway, so the  
825 appendix supports the method-boundary claims  
826 without expanding each check into a separate ta-  
827 ble.

Concern	Check	Takeaway
Layer pair	6/8=6/10 Avg <sup>†</sup> 8.92	Later pairs degrade (≈9.9)
Hard gate	Avg <sup>†</sup> 8.93; min route 99.50%	Not a soft-routing artifact
No reranker	Avg <sup>†</sup> 8.71; default ≤21.80%	Fallback-heavy; audit routes
Explicit LID	Avg <sup>†</sup> 8.84; min route 99.50%	Strong clean reference
8-expert E7	Avg <sup>†</sup> 8.92; min route 99.53%	Base routing is preserved
Capacity	Rank-128 Avg <sup>†</sup> 9.37	Capacity alone insufficient

828 Table A2: Non-core controls on the frozen base-language  
829 split. Each row changes only the named factor and states the  
830 reviewer concern it addresses; Avg<sup>†</sup> follows Section 5.4.

831 Table A3 exists to make the cost boundary ex- 825  
plicit: E7 is an audit-time prior intervention with 826  
added probe computation, not a speed or storage- 827  
efficiency claim. 828

System	Stored	Active params		Runtime	
		Fast	Audit	Fast	Audit
E2 full FT	240.6M	240.6M	240.6M	–	–
E4 shared LoRA	3.54M	3.54M	3.54M	–	–
E5 experts	14.16M	3.54M	3.54M	867.3	867.3
Whisper-LID	14.16M	3.54M	3.54M	904.3	904.3
E7 routed	17.70M+9.2K	3.54M+9.2K	7.08M+9.2K	2173.0	3538.2

832 Table A3: Resource and latency costs. Stored and Active are  
833 parameter counts; Runtime is batch-1 ms/utterance.

834 Table A4 qualifies the strong clean Whisper- 829  
LID reference under short/noisy speech. It is a 830

831  
832

stress diagnostic, not a claim that E7 replaces explicit LID on clean accuracy or latency.

Stress	W-LID Avg	E7 Avg	W-LID acc.	E7 target route
short1s+SNR0	122.03	122.89	70.95	50.83
short2s+SNR10	85.69	84.99	95.36	91.06
short5s+SNR20	52.17	52.51	99.95	99.76

Table A4: Whisper-LID stress diagnostic under degraded audio. Scores use Base Avg<sup>†</sup>; W-LID acc. and E7 target route are macro percentages.