

# LCMA-SRT: Language-Conditional Mixture-of-Experts Adapters for Joint Multilingual Speech Recognition and Translation

Anonymous ACL submission

## Abstract

Neural transducers provide an alignment-free framework for joint automatic speech recognition (ASR) and speech translation (ST). Hierarchical transducer architectures further improve multilingual speech-to-text modeling by stacking a translation-focused encoder on top of an ASR encoder to better handle reordering. However, scaling hierarchical transducers to multilingual many-to-many settings remains challenging: fully shared models often suffer from negative transfer and unstable target-language generation, while training separate models per direction is computationally prohibitive. We propose LCMA-SRT (Language-Conditional Mixture-of-Experts Adapters for Speech Recognition and Translation), which augments a hierarchical transducer with language-conditional Mixture-of-Experts (MoE) adapters. A source-conditioned MoE adapter (SRC-MoE) routes using the source-language embedding to improve acoustic-phonetic modeling and reduce cross-language interference for ASR. A target-conditioned MoE adapter (TGT-MoE) routes using the desired target language to guide reordering and lexical selection and to mitigate cross-target interference in many-to-many ST. Experiments on Europarl-ST (9 languages, 72 directions) show that LCMA-SRT improves both ASR and ST within a single joint model, reducing average WER and increasing BLEU and COMET over strong hierarchical transducer baselines. We release our codes and models at <https://anonymous.4open.science/r/LCMA-SRT>.

## 1 Introduction

Speech translation (ST) enables spoken content to be recognized and translated across languages, supporting cross-lingual communication in settings such as multilingual meetings, education, healthcare, and customer service (Köksal and Yürük, 2020; Al Shamsi et al., 2020). Conventional ST

systems are typically cascades: an ASR module first transcribes speech into text, which is then translated by an MT system (Matusov et al., 2005; Bertoldi and Federico, 2005). Despite the simplicity of this modular design, cascades suffer from error propagation, duplicated computation, and increased system complexity when ASR and MT are optimized separately (Sperber et al., 2017; Rabatin et al., 2024).

End-to-end speech translation has therefore gained momentum as a more streamlined alternative, aiming to map speech directly to target text while sharing representations across ASR and translation (Berard et al., 2016, 2018). Most progress has been driven by attention-based encoder-decoder (AED) models (Gaido et al., 2020; Hussein et al., 2023), but they can be sensitive to segmentation and complicate joint ASR and ST optimization (Anastasopoulos et al., 2022; Sinclair et al., 2014). In contrast, alignment-free objectives such as CTC and neural transducers provide a frame-synchronous modeling paradigm that has been highly effective for speech recognition (Graves et al., 2006; Graves, 2012), and can naturally combine recognition and translation within one architecture (Chiu et al., 2019).

A key difficulty is that ST is typically less monotonic than ASR and often requires substantial reordering (Yan et al., 2023). As a result, training ASR and ST in a single transducer pipeline with fully shared representations can degrade translation quality, especially for language pairs with strong word-order divergence (Wang et al., 2023; Tang et al., 2023).

Hierarchical transducer architectures mitigate this mismatch by explicitly separating recognition and translation within one model. A recent work (Hussein et al., 2025) first produces speech-aligned, approximately monotonic representations via an ASR encoder and then stacks a translation-specific encoder to better model non-monotonic re-

ordering, improving transducer-based ST (Dalmia et al., 2021).

Scaling such hierarchical transducers to multilingual many-to-many ASR and ST remains challenging. Fully shared multilingual models often suffer from negative transfer under language imbalance and cross-language interference, leading to unstable direction control and language mixing on the translation side (Xue et al., 2024), while training separate direction-specific models is computationally impractical. Meanwhile, large unified systems (e.g., Whisper (Radford et al., 2023), HuBERT (Hsu et al., 2021), Whistle (Yusuyin et al., 2024), Seamless (Barrault et al., 2023), Canary (Sekoyan et al., 2025)) show strong multilingual transfer via scaling, but often at high model/compute cost. TTA (Liu et al., 2025) demonstrates that ASR and ST multi-task training with speech and text alignment and language tokens improves cross-lingual transfer; however, such control signals alone may be insufficient to fully prevent cross-target interference in a fully shared many-to-many transducer, motivating language-conditional specialization.

To reconcile scalability with accuracy, we argue that a unified many-to-many model should preserve the benefits of hierarchical transduction while introducing language-conditional specialization, rather than relying on fully shared parameters everywhere. Mixture-of-experts (MoE) offers conditional capacity to capture both shared and language-specific structure (Shazeer et al., 2017; Fedus et al., 2022), and lightweight MoE adapters have been shown to improve multilingual ASR robustness under data imbalance with a small parameter overhead (Li et al., 2025; Mu et al., 2024).

In this work, we propose LCMA-SRT, a unified multilingual hierarchical neural transducer framework that combines a hierarchical transducer encoder (Hussein et al., 2025) with language-conditional MoE adapters to scale joint ASR and ST to many-to-many settings. LCMA-SRT targets the central difficulty of many-to-many multilingual transduction: fully shared models often suffer from negative transfer on the recognition side and cross-target interference on the translation side, while direction-specific training is computationally prohibitive. Our objectives are three-fold: (1) to extend hierarchical transducer modeling from many-to-one to a single many-to-many model for joint ASR and ST, (2) to improve multilingual ASR robustness under heterogeneous acoustic-phonetic

conditions without sacrificing cross-lingual sharing, and (3) to strengthen target-language control for multilingual ST by reducing cross-target interference and language mixing within a unified model. To this end, we introduce two complementary language-conditional adapter mechanisms: a source-conditioned MoE adapter (SRC-MoE) and a target-conditioned MoE adapter (TGT-MoE), integrated into the hierarchical transducer backbone. We evaluate LCMA-SRT on Europarl-ST (Iranzo-Sánchez et al., 2020) and perform controlled ablations to quantify the impact of each design component in many-to-many joint ASR and ST.

## 2 Proposed Approach

Given input speech features  $X = (x_1, \dots, x_T) \in \mathbb{R}^{T \times F}$ , we consider two token sequences: a source-language transcript  $\mathbf{y}^{(s)} = (y_1^s, \dots, y_{U_s}^s) \in \mathcal{V}_{U_s}^{(s)}$  (ASR) and a target-language translation  $\mathbf{y}^{(t)} = (y_1^t, \dots, y_{U_t}^t) \in \mathcal{V}_{U_t}^{(t)}$  (ST) (Graves, 2012). For each task  $k \in \{s, t\}$ , a neural transducer models the conditional probability  $P(\mathbf{y}^{(k)} | X)$  by marginalizing over all monotonic alignments  $\mathbf{a} \in \bar{\mathcal{V}}_{T+U_k}^{(k)}$ :

$$P(\mathbf{y}^{(k)} | X) = \sum_{\mathbf{a} \in \mathcal{B}^{-1}(\mathbf{y}^{(k)})} P(\mathbf{a} | X) \quad (1)$$

$$\bar{\mathcal{V}}^{(k)} = \mathcal{V}^{(k)} \cup \{\phi\} \quad (2)$$

where  $\phi$  denotes the blank symbol and  $\mathcal{B}$  deterministically removes blanks to map an alignment to its corresponding label sequence. The transducer parameterization uses an encoder, a predictor, and a joiner to produce a posterior distribution over  $\bar{\mathcal{V}}^{(k)}$  at each lattice state. To better handle the monotonicity mismatch between recognition and translation, a hierarchical transducer employs a stacked encoder hierarchy (Hussein et al., 2025). An ASR encoder first produces speech-aligned representations, which are then transformed by a translation-specific encoder:

$$\mathbf{F}^{(s)} = \text{Enc}_{\text{asr}}(X) \quad (3)$$

$$\mathbf{F}^{(t)} = \text{Enc}_{\text{st}}(\mathbf{F}^{(s)}) \quad (4)$$

Task-specific transducer heads are applied on top of  $\mathbf{F}^{(s)}$  (ASR) and  $\mathbf{F}^{(t)}$  (ST), respectively. We train the model with a multitask transducer objective:

$$\mathcal{L}_{\text{nt}} = \alpha_{\text{asr}} \mathcal{L}_{\text{nt}}^{(s)} + \alpha_{\text{st}} \mathcal{L}_{\text{nt}}^{(t)} \quad (5)$$

$$\mathcal{L}_{\text{nt}}^{(k)} = -\log P(\mathbf{y}^{(k)} | X) \quad (6)$$

where  $\alpha_{\text{asr}}$  and  $\alpha_{\text{st}}$  control the contribution of each task.

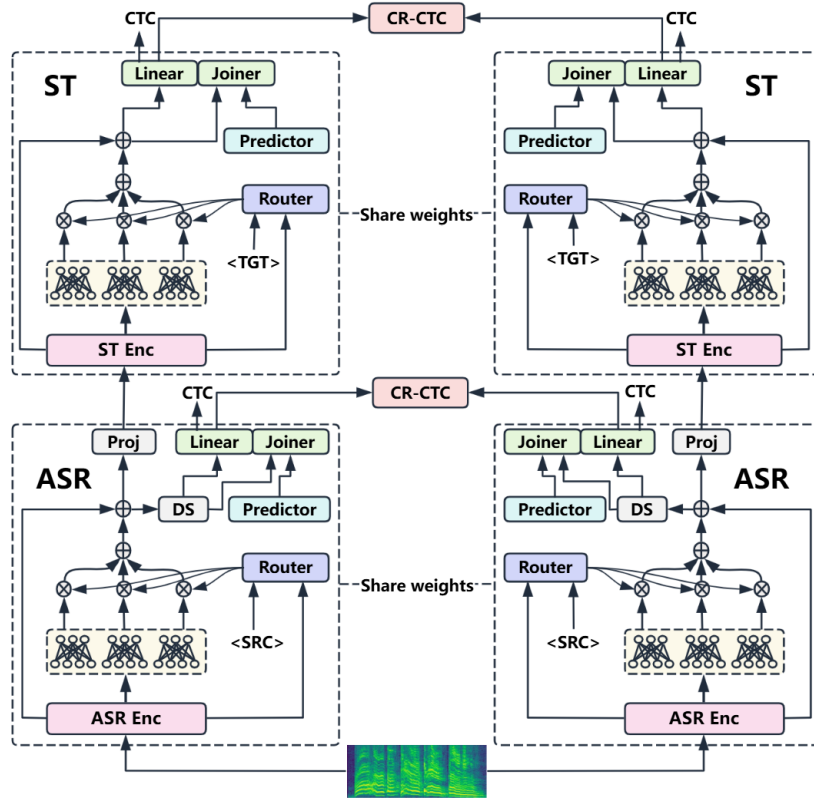


Figure 1: Proposed language-conditional mixture-of-experts adapters for unified multilingual speech recognition and translation (LCMA-SRT).

**LCMA-SRT Overview:** To improve multilingual ASR accuracy and extend hierarchical transducer modeling from many-to-one settings to a unified many-to-many setting, we propose LCMA-SRT. As illustrated in Figure 1, LCMA-SRT augments the hierarchical transducer backbone (Hussein et al., 2025) with two lightweight language-conditional MoE adapters inserted once at the encoder outputs: a SRC-MoE after the ASR encoder and TGT-MoE after the ST encoder.

## 2.1 ASR Encoder and SRC-MoE

As shown in Figure 1, the ASR branch takes the speech-aligned encoder representations  $\mathbf{F}^{(s)}$  (defined in Eq. 3) and applies a SRC-MoE to mitigate cross-language interference in multilingual recognition. Let  $\mathbf{F}^{(s)} = \{\mathbf{f}_t^{(s)}\}_{t=1}^T$  and let  $\mathbf{e}_s = \text{Embed}_{\text{src}}(\ell_s) \in \mathbb{R}^{d_\ell}$  denote a learned embedding of the source language identity (corresponding to <SRC> in Figure 1). SRC-MoE consists of a router  $g(\cdot)$  and  $E$  lightweight experts  $\{f_i(\cdot)\}_{i=1}^E$  (Shazeer et al., 2017; Fedus et al., 2022). For notational clarity, we describe the adapter in a generic form using an input hidden vector  $\mathbf{h}_t$ ; in the ASR branch we set  $\mathbf{h}_t = \mathbf{f}_t^{(s)}$ . The router conditions on the

hidden state and the source-language embedding to produce mixture weights:

$$\pi_t = g([\mathbf{h}_t; \mathbf{e}_s]) \in \mathbb{R}^E \quad (7)$$

$$\mathbf{w}_t = \text{softmax}(\pi_t) \in \mathbb{R}^E \quad (8)$$

where  $[\cdot; \cdot]$  denotes concatenation. Each expert transforms  $\mathbf{h}_t$ , and the adapter output is computed as a weighted mixture with a residual connection:

$$\mathbf{m}_t = \sum_{i=1}^E w_{t,i} f_i(\mathbf{h}_t) \quad (9)$$

$$\tilde{\mathbf{h}}_t = \mathbf{h}_t + \mathbf{m}_t \quad (10)$$

Stacking  $\tilde{\mathbf{h}}_t$  over time yields the adapted sequence  $\tilde{\mathbf{F}}^{(s)} = \{\tilde{\mathbf{f}}_t^{(s)}\}_{t=1}^T$ , where  $\tilde{\mathbf{f}}_t^{(s)}$  corresponds to  $\tilde{\mathbf{h}}_t$  under the substitution  $\mathbf{h}_t = \mathbf{f}_t^{(s)}$ . By conditioning routing on the source language, SRC-MoE allocates capacity to language-specific acoustic-phonetic variation while retaining sharing in the backbone, thereby improving multilingual ASR robustness.

Following the SRC-MoE adapter, we apply a  $2 \times$  downsampling module on the ASR branch before

the ASR transducer head:

$$\hat{\mathbf{F}}^{(s)} = \text{DS}\left(\tilde{\mathbf{F}}^{(s)}\right) \in \mathbb{R}^{T/2 \times d} \quad (11)$$

The ASR predictor–joiner head then defines  $P(\mathbf{y}^{(s)} | X)$  using  $\hat{\mathbf{F}}^{(s)}$  as encoder inputs. Meanwhile,  $\tilde{\mathbf{F}}^{(s)}$  is forwarded to the ST branch (via the projection module in Figure 1), enabling improvements in source-side modeling to benefit downstream translation.

## 2.2 ST Encoder and TGT-MoE

We adopt a two-stage training recipe: we first pre-train the ASR branch, and then perform a second-stage joint training of ASR and ST on top of the pretrained model. During the second stage, the ST encoder consumes the refined ASR representations  $\tilde{\mathbf{F}}^{(s)}$  produced by SRC-MoE (Section 2.1), which implicitly encode source-language characteristics learned from multilingual ASR training. These source-aware intermediate representations facilitate learning translation while maintaining effective cross-lingual sharing.

As shown in Figure 1, we first apply a lightweight projection to match the input dimension expected by the ST encoder:

$$\mathbf{G}^{(s)} = \text{Proj}\left(\tilde{\mathbf{F}}^{(s)}\right) \quad (12)$$

The ST encoder then produces translation-oriented representations:

$$\mathbf{F}^{(t)} = \text{Enc}_{\text{st}}\left(\mathbf{G}^{(s)}\right) \quad (13)$$

where  $\mathbf{F}^{(t)} = \{\mathbf{f}_t^{(t)}\}_{t=1}^T$ .

We apply a TGT-MoE once at the ST encoder output. Let  $\mathbf{e}_t = \text{Embed}_{\text{tgt}}(\ell_t) \in \mathbb{R}^{d_\ell}$  denote the learned embedding of the desired target language (corresponding to <TGT> in Figure 1). TGT-MoE follows Eqs. 7–10, where the hidden input is set to  $\mathbf{h}_t = \mathbf{f}_t^{(t)}$  and the conditioning embedding is  $\mathbf{e} = \mathbf{e}_t$ . The adapted ST representations are:

$$\tilde{\mathbf{F}}^{(t)} = \text{Adapter}_{\text{tgt}}\left(\mathbf{F}^{(t)}; \mathbf{e}_t\right) \quad (14)$$

Conditioning routing on the target language guides reordering and lexical selection toward the desired output language, reducing cross-target interference and stabilizing target-language generation in many-to-many ST. The ST predictor–joiner head defines the translation distribution using  $\tilde{\mathbf{F}}^{(t)}$  as encoder inputs, yielding  $P(\mathbf{y}^{(t)} | X, \ell_t)$ .

## 2.3 CR-CTC Self-Distillation

Following prior work on hierarchical transducer training for joint ASR and ST (Hussein et al., 2025), we adopt consistency-regularized CTC (CR-CTC) as an auxiliary self-distillation signal (Yao et al., 2024). CR-CTC constructs two stochastic views of the same utterance under shared parameters and encourages their CTC posteriors to be consistent.

We attach lightweight CTC heads to both branches. For ASR, the auxiliary losses are computed on the downsampled representations  $\hat{\mathbf{F}}^{(s)}$  (Eq. 11); for ST, they are computed on the TGT-MoE output  $\tilde{\mathbf{F}}^{(t)}$  (Eq. 14). We define the auxiliary loss for each task as a weighted sum of CTC and CR terms:

$$\mathcal{L}_{\text{aux}}^{(s)} = \lambda_{\text{ctc}}^{(s)} \mathcal{L}_{\text{ctc}}^{(s)} + \lambda_{\text{cr}}^{(s)} \mathcal{L}_{\text{cr}}^{(s)} \quad (15)$$

$$\mathcal{L}_{\text{aux}}^{(t)} = \lambda_{\text{ctc}}^{(t)} \mathcal{L}_{\text{ctc}}^{(t)} + \lambda_{\text{cr}}^{(t)} \mathcal{L}_{\text{cr}}^{(t)} \quad (16)$$

Here  $\mathcal{L}_{\text{ctc}}^{(k)}$  is the CTC negative log-likelihood, and  $\mathcal{L}_{\text{cr}}^{(k)}$  enforces view-consistent CTC posteriors via a symmetric KL divergence over valid frames. The coefficients  $\lambda_{\text{ctc}}^{(s)}$ ,  $\lambda_{\text{cr}}^{(s)}$  and  $\lambda_{\text{ctc}}^{(t)}$ ,  $\lambda_{\text{cr}}^{(t)}$  weight the CTC and CR-CTC terms for ASR and ST, respectively.

## 2.4 MoE Entropy Regularization

MoE routing may collapse to a small subset of experts. To encourage balanced expert utilization, we regularize the router to have high entropy (Shazeer et al., 2017). Let  $\mathbf{w} \in \mathbb{R}^{T \times B \times E}$  denote the routing weights returned by an MoE module, where  $\mathbf{w}_{t,b}$  is a probability simplex over  $E$  experts. The per-frame entropy is

$$H(\mathbf{w}_{t,b}) = - \sum_{i=1}^E w_{t,b,i} \log w_{t,b,i} \quad (17)$$

With an optional padding mask  $\mathbf{M} \in \{0, 1\}^{T \times B}$  ( $M_{t,b} = 1$  indicates padded frames), we compute the mean entropy over valid frames as

$$\bar{H}(\mathbf{w}) = \frac{\sum_{t,b}(1 - M_{t,b}) H(\mathbf{w}_{t,b})}{\sum_{t,b}(1 - M_{t,b})} \quad (18)$$

We apply entropy regularization to both SRC-MoE and TGT-MoE, and average the two weighted terms:

$$\mathcal{L}_{\text{ent}} = -\frac{1}{2} \sum_{k \in \{s,t\}} \lambda_{\text{ent}}^{(k)} \cdot \bar{H}\left(\mathbf{w}^{(k)}\right) \quad (19)$$

Here  $\lambda_{\text{ent}}^{(s)}$  and  $\lambda_{\text{ent}}^{(t)}$  control the strength of entropy regularization for SRC-MoE and TGT-MoE, respectively.

Finally, we optimize the overall multitask transducer loss together with the auxiliary CR-CTC losses and the MoE entropy regularizer:

$$\mathcal{L} = \mathcal{L}_{\text{nt}} + \alpha_{\text{asr}} \mathcal{L}_{\text{aux}}^{(s)} + \alpha_{\text{st}} \mathcal{L}_{\text{aux}}^{(t)} + \mathcal{L}_{\text{ent}} \quad (20)$$

## 2.5 Decoding for ASR and ST

We decode both ASR and ST with a modified transducer beam search using their task-specific predictor-joiner heads (Graves, 2012; Jain et al., 2019). During inference, we condition the model on the source and target language identities (i.e., the <SRC> and <TGT> signals in Figure 1), so that SRC-MoE and TGT-MoE routing remains active in the ASR and ST branches.

For ST, we enforce the translation direction with a target-language prefix token  $\tau(\ell_t)$  (e.g., <2xx>) (Johnson et al., 2016) for all systems in this work. ST training labels are prepended with  $\tau(\ell_t)$ , and at inference we force the first emitted token to be the same prefix. For the direction-specific HENT-SRT-M2O $\times$ 9 baseline,  $\tau(\ell_t)$  is fixed within each model (one target language per model). The prefix token is removed for evaluation. Following prior transducer-based ST practice, we further apply a blank penalty during ST decoding to discourage excessive blank emissions (Hussein et al., 2025). However, a prefix-only constraint does not prevent target-language drift in fully shared models (Sec. 3.2).

## 2.6 Computation Efficiency

LCMA-SRT is designed to improve many-to-many multilingual ASR and ST without sacrificing the compactness of hierarchical transducer modeling: instead of training and deploying separate direction-specific systems, it supports all translation directions within a single shared backbone, avoiding costs that grow with the number of language pairs. The extra computation is confined to two lightweight language-conditional MoE adapters inserted once after  $\text{Enc}_{\text{ASR}}$  and once after  $\text{Enc}_{\text{ST}}$  (Figure 1); each adapter comprises a small router and  $E$  shallow experts (two-layer MLPs) with a residual connection, so the dominant cost still lies in the shared encoders and the adapter overhead scales mainly with  $E$ , not with the number of translation directions. In our implementation, both branches

use the Zipformer encoder and we train the transducer with the pruned transducer loss (Yao et al., 2023; Kuang et al., 2022), preserving the efficiency of the hierarchical transducer backbone while enabling conditional capacity allocation via SRC-MoE and TGT-MoE.<sup>1</sup>

## 3 Experiments

### 3.1 Experimental setup

We conduct experiments on Europarl-ST (Iranzo-Sánchez et al., 2020), a multilingual speech translation benchmark with paired speech, source-language transcripts, and target-language translations. We consider nine languages (en, de, es, fr, it, nl, pl, pt, ro) and cover all ordered cross-lingual pairs, yielding 72 translation directions. We follow the official train/dev/test splits and report test-set results for both ASR and ST.

**Data processing:** Our experiments utilize the icefall framework, with the Lhotse toolkit (Želasko et al., 2021) for speech data preparation and dataset construction. ASR pretraining uses 276 hours of unique audio. For joint ASR and ST training, we reuse the same segments and pair each utterance with its available translations into the remaining target languages, yielding 1,642 pair-hours. All audio recordings are resampled to 16 kHz. We extract 80-dimensional log-Mel filterbank features using a 25 ms window with a 10 ms frame shift. Additionally, we apply on-the-fly SpecAugment (Park et al., 2019), incorporating time warping (maximum factor of 80), frequency masking (two regions, max width of 27 bins), and time masking (ten regions, max width of 100 frames). For CR-CTC self-distillation, we strengthen the augmentation of the second view by increasing both the number of time-masked regions and the maximum masking fraction by a factor of 2.5, following prior practice (Yao et al., 2024). For text, we apply Whisper-style normalization rules (Radford et al., 2023) prior to subword tokenization. For subword tokenization, we train SentencePiece BPE models (Kudo and Richardson, 2018) with a vocabulary size of 1k for ASR and 1k for ST. On the ASR side, we use a single multilingual BPE model shared across all source-language transcripts. On the ST side, we include a target-language tag token (e.g., <2xx>) in the ST vocabulary and prepend the corresponding tag to every ST training label sequence to indicate

<sup>1</sup><https://github.com/k2-fsa/icefall>

the desired output language. For all multilingual many-to-many systems, we use a shared 1k ST BPE model. For the direction-specific baseline trained as nine separate systems (one per target language), we instead train an independent 1k ST BPE model for each target language, while keeping all other preprocessing settings identical.

**Models:** We reproduce a multitask hierarchical neural transducer for joint ASR and ST, employing the Zipformer encoder architecture and task-specific transducer heads. The encoder is partitioned into an ASR encoder stacked below an ST encoder, and we keep the Zipformer configuration identical to HENT-SRT for fair comparison. For Stage 1 ASR pretraining, we train three ASR-only variants: CR-CTC, CR-CTC+S-Bias (with a source-identity bias added to the encoder outputs), and CR-CTC+SRC-MoE (with the source-conditioned MoE adapter applied after the ASR encoder). The first two models have 30M parameters; adding SRC-MoE increases the parameter count to 36M. For Stage 2 joint ASR and ST training, we compare (i) direction-specific many-to-one training as nine separate models HENT-SRT-M2O $\times$ 9 (each 61M), (ii) HENT-SRT-M2M, a fully shared many-to-many extension of the original many-to-one HENT-SRT (61M; Sec. 2.5), and (iii) LCMA-SRT and its ablations, which examine the SRC-MoE and dissect the TGT-MoE into an unconditioned MoE and a simple target-identity bias. Unless otherwise specified, we use  $E=8$  experts for SRC-MoE and  $E=16$  experts for TGT-MoE, resulting in model sizes ranging from 66M to 77M parameters. Both tasks use a stateless transducer predictor implemented as a single Conv1D layer with kernel size 2, together with a standard joiner.

Our training uses the ScaledAdam optimizer (Yao et al., 2023) with a learning rate of 0.02 and a 5k-step warmup. We train in two stages on 4 NVIDIA A800 GPUs, each for 50 epochs, using duration-based batching with max-duration=900 s in Stage 1 and 450 s in Stage 2. Stage 2 is initialized from the best checkpoint of Stage 1. The multitask weights  $\alpha_{\text{asr}}$  and  $\alpha_{\text{st}}$  in Eq. (5) are both set to 1. For the auxiliary CR-CTC objective (Sec. 2.3), we set  $\lambda_{\text{ctc}}^{(s)} = \lambda_{\text{ctc}}^{(t)} = 0.1$  and  $\lambda_{\text{cr}}^{(s)} = \lambda_{\text{cr}}^{(t)} = 0.05$  for the CTC and CR terms, respectively. When MoE adapters are enabled, we apply router entropy regularization with  $\lambda_{\text{ent}}^{(s)} = \lambda_{\text{ent}}^{(t)} = 0.015$ . All decoding results are obtained with beam search using a beam size of 20. For ST, we apply a blank penalty of 2.

**Evaluation:** We evaluate ASR with word error rate (WER). For speech translation, we report BLEU (Papineni et al., 2002), COMET (Rei et al., 2020), and sentence-level target-language mismatch rate (LMR) using an off-the-shelf language identification model (Joulin et al., 2016): a hypothesis is matched only if it is classified as the specified target language with confidence  $\geq 0.7$ . To ensure consistent comparisons, we normalize ASR transcripts and ST outputs by removing punctuation and evaluating case-insensitively.

### 3.2 Main Results

We evaluate our proposed framework on Europarl-ST following the two-stage training recipe. In this section, we analyze the effectiveness of the proposed adapters in ASR pretraining, followed by a comprehensive comparison against strong baselines in the joint ASR and ST setting.

**Multilingual ASR Pretraining:** Table 1 reports pretraining results. Relative to the CR-CTC baseline (22.35% Avg. WER), an unconditioned MoE yields a modest gain (22.06%), while adding a source-identity bias reduces WER to 21.08%. The proposed SRC-MoE achieves the best performance (20.88%), obtaining the lowest WER on 8/9 languages, suggesting that source-conditioned routing learns more robust speech-aligned representations for subsequent training.

**Many-to-Many Joint Training:** We present the many-to-many joint training results in Table 2 and Table 3. We evaluate LCMA-SRT against the direction-specific baseline HENT-SRT-M2O $\times$ 9. Notably, LCMA-SRT significantly outperforms this strong baseline, yielding an average improvement of +5.2 BLEU (15.3  $\rightarrow$  20.5) and +0.076 COMET (0.575  $\rightarrow$  0.651). We observe consistent gains across diverse data conditions: high-resource directions achieve substantial absolute improvements (e.g., Spanish: 19.1  $\rightarrow$  25.8 BLEU), while lower-resource directions see remarkable relative gains (e.g., Polish: 7.2  $\rightarrow$  10.7 BLEU, +48% relative improvement). This suggests that LCMA-SRT enables effective positive transfer from high-resource to low-resource languages within the shared MoE framework. Comparison with the fully shared baseline HENT-SRT-M2M further highlights the advantage of explicit target conditioning. Despite forced target-prefix decoding (Sec. 2.5), the fully shared HENT-SRT-M2M collapses under 72 competing directions due to severe

Model	WER (%) ↓									
	de	en	es	fr	it	nl	pl	pt	ro	Avg
CR-CTC	24.57	18.59	20.76	19.24	17.33	36.75	25.28	19.82	18.77	22.35
+ MoE	24.39	18.41	20.16	18.61	17.28	36.83	24.36	19.70	18.79	22.06
+ S-Bias	23.89	17.60	19.58	17.41	16.73	<b>34.72</b>	23.63	18.21	17.97	21.08
+ SRC-MoE	<b>23.34</b>	<b>17.45</b>	<b>19.41</b>	<b>17.34</b>	<b>16.27</b>	35.20	<b>23.28</b>	<b>18.16</b>	<b>17.48</b>	<b>20.88</b>

Table 1: Multilingual ASR results on Europarl-ST. WER is reported per source language, and Avg denotes the overall average. We report the CR-CTC baseline and its variants with an unconditioned MoE adapter (+MoE), a source-identity bias (+S-Bias), and the proposed source-conditioned MoE adapter inserted after the ASR encoder (+SRC-MoE).

Model	WER (%) ↓	Average BLEU ↑									
		de	en	es	fr	it	nl	pl	pt	ro	Avg
HENT-SRT-M2O×9	23.28	10.7	21.2	19.1	18.2	14.2	16.5	7.2	18.4	12.1	15.3
HENT-SRT-M2M	16.65	2.6	12.8	5.5	4.0	1.8	3.5	1.2	4.9	2.5	4.3
LCMA-SRT	<b>15.71</b>	<b>15.2</b>	<b>25.9</b>	<b>25.8</b>	<b>24.7</b>	<b>20.0</b>	<b>20.5</b>	<b>10.7</b>	<b>23.9</b>	<b>17.6</b>	<b>20.5</b>
TGT-MoE→MoE	16.42	2.3	14.7	4.7	3.3	1.7	2.7	1.1	4.5	2.0	4.1
TGT-MoE→T-Bias	15.84	13.1	22.7	23.5	22.3	17.7	18.1	8.3	21.8	14.5	18.0
w/o TGT-MoE	16.48	2.0	12.8	5.9	3.9	1.6	3.0	1.3	5.0	2.2	4.2
w/o SRC-MoE	16.11	14.5	24.9	25.0	24.6	19.6	20.0	10.5	23.7	17.5	20.0

Table 2: Joint ASR and ST results on Europarl-ST. WER is averaged over all 72 translation directions. BLEU is averaged over directions grouped by their target language, and Avg denotes the overall average across all directions. We compare HENT-SRT-M2O×9 and HENT-SRT-M2M against LCMA-SRT and ablations that replace TGT-MoE with an unconditioned MoE (TGT-MoE→MoE) or a target-identity bias (TGT-MoE→T-Bias), or remove TGT-MoE / SRC-MoE (w/o TGT-MoE, w/o SRC-MoE).

Model	LMR (%) ↓	Average COMET ↑									
		de	en	es	fr	it	nl	pl	pt	ro	Avg
HENT-SRT-M2O×9	<b>0.65</b>	0.507	0.656	0.587	0.542	0.565	0.558	0.550	0.609	0.598	0.575
HENT-SRT-M2M	84.95	0.380	0.543	0.478	0.427	0.435	0.401	0.385	0.471	0.406	0.436
LCMA-SRT	<u>0.75</u>	<b>0.574</b>	<b>0.715</b>	<b>0.682</b>	<b>0.627</b>	<b>0.656</b>	<b>0.613</b>	<b>0.616</b>	<b>0.693</b>	<b>0.678</b>	<b>0.651</b>
TGT-MoE→MoE	85.23	0.380	0.559	0.476	0.426	0.438	0.395	0.386	0.472	0.408	0.438
TGT-MoE→T-Bias	0.78	0.529	0.675	0.642	0.583	0.612	0.563	0.562	0.651	0.621	0.604
w/o TGT-MoE	85.19	0.376	0.545	0.480	0.427	0.434	0.398	0.387	0.473	0.407	0.436
w/o SRC-MoE	0.81	0.568	0.708	0.671	0.621	0.646	0.606	0.605	0.685	0.675	0.643

Table 3: Average COMET and LMR scores for the same models and ablations as Table 2.

cross-direction interference and target-language drift (LMR 84.95%), resulting in poor translation quality (4.3 BLEU). In contrast, LCMA-SRT maintains a low LMR (0.75%), comparable to the direction-specific HENT-SRT-M2O×9 (0.65%), while substantially improving translation quality (+5.2 BLEU, +0.076 COMET) and reducing WER to 15.71%. These findings demonstrate that LCMA-SRT effectively scales hierarchical transducer modeling to the many-to-many setting within a single unified model, achieving performance that outperforms the strong hierarchical transducer baselines reproduced in this work on both tasks simultaneously. Direction-wise ASR

WER results are reported in Appendix A Table 4; direction-wise ST BLEU/COMET/LMR results are reported in Appendix A Tables 5, 6, and 7.

### 3.3 Ablation Analysis

To assess the contribution of each component in LCMA-SRT, we conduct ablations in Table 2 and Table 3. We first examine the critical role of Target-Conditioned MoE (TGT-MoE). Removing TGT-MoE (w/o TGT-MoE) leads to a collapse in translation quality (4.2 BLEU) and an extremely high LMR (85.19%), showing that a monolithic ST encoder fails to preserve target-language fidelity under 72 competing directions. Replacing condi-

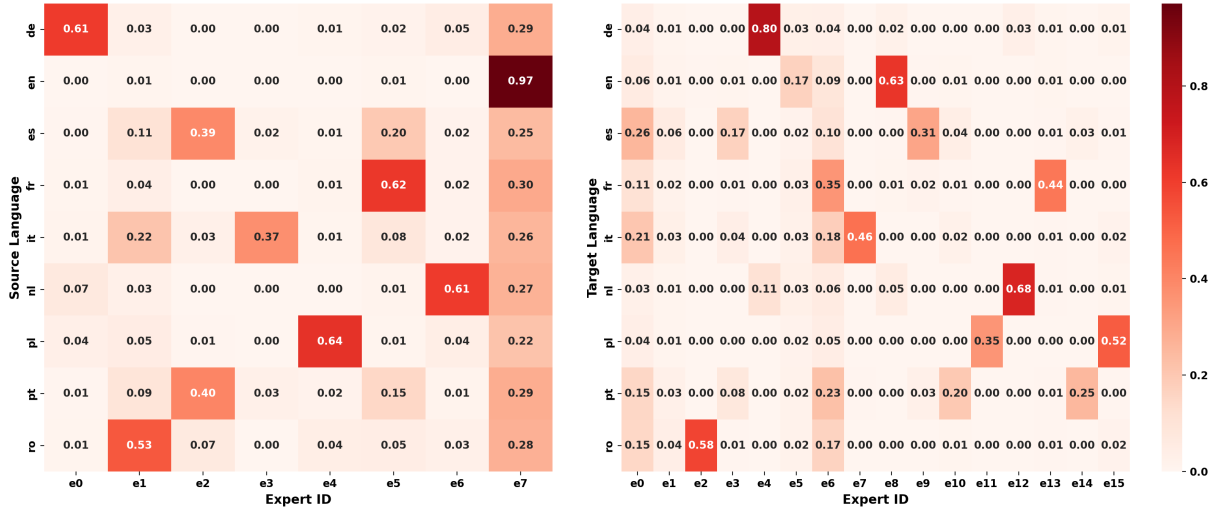


Figure 2: Heatmaps of language-aggregated expert routing weights in the language-conditional MoE adapters. Left: SRC-MoE on the ASR side, where each row is a source language and each column is an expert (e0–e7). Right: TGT-MoE on the ST side, where each row is a target language and each column is an expert (e0–e15).

538 tional routing with an unconditioned MoE (TGT-  
 539 MoE→MoE) yields similarly high LMR (85.23%)  
 540 and poor BLEU, demonstrating that merely in-  
 541 creasing parameters is insufficient—explicit tar-  
 542 get conditioning is essential to disentangle lan-  
 543 guages. Interestingly, a target-identity bias (TGT-  
 544 MoE→T-Bias) restores low LMR (0.78%), indicat-  
 545 ing that directional cues help avoid language leak-  
 546 age; however, it still lags behind full TGT-MoE  
 547 in BLEU/COMET, suggesting that expert-based  
 548 non-linear capacity is necessary beyond simply se-  
 549 lecting the right language. For the ASR encoder,  
 550 removing the SRC-MoE (w/o SRC-MoE) leads to  
 551 a degradation in ASR WER (15.71% → 16.11%)  
 552 and a consequent drop in ST BLEU (20.5 → 20.0),  
 553 validating that mitigating source-side acoustic in-  
 554 terference via conditional adapters yields better  
 555 speech-aligned representations that benefit down-  
 556 stream translation.

### 557 3.4 Analysis of MoE Routing Behavior

558 Figure 2 visualizes the language-aggregated  
 559 routing distributions learned by our language-  
 560 conditional MoE adapters under entropy regulariza-  
 561 tion. On the ASR side (SRC-MoE, left), each ex-  
 562 pert is biased toward certain source languages (e.g.,  
 563 e<sub>4</sub> is strongly activated by pl), while several ex-  
 564 perts are concurrently used by multiple languages (e.g.,  
 565 overlapping usage between es and pt). This sug-  
 566 gests that SRC-MoE separates language-specific  
 567 acoustic–phonetic variation while retaining shared  
 568 capacity for transferable structure. On the ST side  
 569 (TGT-MoE, right), we observe an analogous but

570 more target-oriented behavior: routing is target-  
 571 dependent (e.g., de and en have different dominant  
 572 experts), yet remains soft with noticeable weight on  
 573 secondary experts, enabling partial sharing when  
 574 beneficial. Such soft target-conditioned specializa-  
 575 tion provides a direct mechanism to reduce inter-  
 576 direction interference in many-to-many translation  
 577 without discarding sharing among related targets.  
 578 In contrast, without entropy regularization, routing  
 579 collapses into near-disjoint partitions: each expert  
 580 becomes predominantly focused on a single lan-  
 581 guage with little to no cross-language sharing, re-  
 582 ducing the intended benefit of conditional capacity  
 583 allocation across languages.

## 584 4 Conclusion

585 We propose LCMA-SRT, a unified many-to-many  
 586 multilingual speech translation extension of hi-  
 587 erarchical neural transducers. By inserting two  
 588 lightweight language-conditional MoE adapters  
 589 (SRC-MoE after the ASR encoder; TGT-MoE af-  
 590 ter the ST encoder), it allocates capacity across  
 591 source and target languages while keeping a  
 592 shared backbone. On Europarl-ST, LCMA-SRT  
 593 markedly improves many-to-many translation over  
 594 fully shared baselines and matches or surpasses  
 595 direction-specific systems within a single model.  
 596 Ablations and routing analysis show SRC-MoE  
 597 mainly strengthens multilingual ASR representa-  
 598 tions, while TGT-MoE is crucial for stable target-  
 599 controlled translation; together they provide com-  
 600plementary gains.

## 601 Limitations

602 This work evaluates LCMA-SRT only in the of-  
603 fline setting and does not benchmark streaming  
604 performance; moreover, experiments are primar-  
605 ily conducted on Europarl-ST with limited domain  
606 and language-family coverage. LCMA-SRT relies  
607 on explicit source/target language signals (and a  
608 target prefix for ST) to drive conditional routing,  
609 which may be affected by imperfect language iden-  
610 tification or missing target specifications in practice.  
611 Finally, we do not systematically study accuracy-  
612 efficiency trade-offs or hyperparameter sensitivity  
613 (e.g., number of experts and entropy coefficient),  
614 leaving broader scalability and deployment consid-  
615 erations for future work.

## 616 References

617 Hilal Al Shamsi, Abdullah G Almutairi, Sulaiman  
618 Al Mashrafi, and Talib Al Kalbani. 2020. Implica-  
619 tions of language barriers for healthcare: a systematic  
620 review. *Oman medical journal*, 35(2):e122.

621 Antonios Anastasopoulos and 1 others. 2022. Findings  
622 of the iwslt 2022 evaluation campaign. In *Proceed-*  
623 *ings of the 19th International Conference on Spoken*  
624 *Language Translation*, pages 98–157.

625 Loïc Barrault, Yu-An Chung, Mariano Coria Megli-  
626 oli, David Dale, Ning Dong, Mark Duppenhaler,  
627 Paul-Ambroise Duquenne, Brian Ellis, Hady Elsahar,  
628 Justin Haaheim, and 1 others. 2023. *Seamless: Mul-*  
629 *tilingual expressive and streaming speech translation.*  
630 *Preprint*, arXiv:2312.05187.

631 Alexandre Berard, Laurent Besacier, Ali Can Ko-  
632 cabiyikoglu, and Olivier Pietquin. 2018. End-to-  
633 end automatic speech translation of audiobooks. In  
634 *ICASSP 2018*, pages 6224–6228.

635 Alexandre Berard, Olivier Pietquin, Christophe Servan,  
636 and Laurent Besacier. 2016. Listen and translate: A  
637 proof of concept for end-to-end speech-to-text trans-  
638 lation. In *NIPS Workshop on end-to-end learning for*  
639 *speech and audio processing*.

640 Nicola Bertoldi and Marcello Federico. 2005. A new  
641 decoder for spoken language translation based on  
642 confusion networks. In *IEEE Workshop on Auto-*  
643 *matic Speech Recognition and Understanding*, pages  
644 86–91.

645 Chung-Cheng Chiu and 1 others. 2019. A comparison  
646 of end-to-end models for long-form speech recogni-  
647 tion. In *ASRU 2019*, pages 889–896.

648 Siddharth Dalmia, Brian Yan, Vikas Raunak, Florian  
649 Metze, and Shinji Watanabe. 2021. Searchable hid-  
650 den intermediates for end-to-end models of decom-  
651 posable sequence tasks. In *NAACL 2021*, pages 1882–  
652 1896.

653 William Fedus, Jeff Dean, and Barret Zoph. 2022. A re-  
654 view of sparse expert models in deep learning. *arXiv*  
655 *preprint arXiv:2209.01667*.

656 Marco Gaido, Mattia A Di Gangi, Matteo Negri, and  
657 Marco Turchi. 2020. End-to-end speech-translation  
658 with knowledge distillation: Fbk@ iwslt2020. In  
659 *Proceedings of the 17th International Conference on*  
660 *Spoken Language Translation*, pages 80–88.

661 Alex Graves. 2012. Sequence transduction with  
662 recurrent neural networks. *arXiv preprint*  
663 *arXiv:1211.3711*.

664 Alex Graves, Santiago Fernández, Faustino Gomez, and  
665 Jürgen Schmidhuber. 2006. Connectionist temporal  
666 classification: labelling unsegmented sequence data  
667 with recurrent neural networks. In *Proceedings of the*  
668 *23rd international conference on Machine learning*,  
669 pages 369–376.

670 Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai,  
671 Kushal Lakhotia, Ruslan Salakhutdinov, and Abdel-  
672 rahman Mohamed. 2021. *Hubert: Self-supervised*  
673 *speech representation learning by masked prediction*  
674 *of hidden units*. *IEEE/ACM Transactions on Audio,*  
675 *Speech, and Language Processing*, 29:3451–3460.

676 Amir Hussein, Dan Povey, Cihan Xiao, Leibny Paola  
677 Garcia, Matthew Wiesner, and Sanjeev Khu-  
678 danpur. 2025. Hent-srt: Hierarchical efficient  
679 neural transducer with self-distillation for joint  
680 speech recognition and translation. *arXiv preprint*  
681 *arXiv:2506.02157*.

682 Amir Hussein, Brian Yan, Antonios Anastasopoulos,  
683 Shinji Watanabe, and Sanjeev Khudanpur. 2023. En-  
684 hancing end-to-end conversational speech translation  
685 through target language context utilization. *arXiv*  
686 *preprint arXiv:2309.15686*.

687 Javier Iranzo-Sánchez, Joan Albert Silvestre-Cerda,  
688 Javier Jorge, Nahuel Roselló, Adria Giménez, Al-  
689 bert Sanchis, Jorge Civera, and Alfons Juan. 2020.  
690 Europarl-st: A multilingual corpus for speech transla-  
691 tion of parliamentary debates. In *ICASSP 2020-2020*  
692 *IEEE International Conference on Acoustics, Speech*  
693 *and Signal Processing (ICASSP)*, pages 8229–8233.  
694 IEEE.

695 Mahaveer Jain, Kjell Schubert, Jay Mahadeokar, Ching-  
696 Feng Yeh, Kaustubh Kalgaonkar, Anuroop Sriram,  
697 Christian Fuegen, and Michael L. Seltzer. 2019. Rnn-  
698 t for latency controlled asr with improved beam  
699 search. *arXiv preprint arXiv:1911.01629*.

700 Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim  
701 Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thor-  
702 rat, Fernanda Viégas, Martin Wattenberg, Greg Cor-  
703 rado, Macduff Hughes, and Jeffrey Dean. 2016.  
704 Google’s multilingual neural machine translation sys-  
705 tem: Enabling zero-shot translation. *arXiv preprint*  
706 *arXiv:1611.04558*.

707 Armand Joulin, Edouard Grave, Piotr Bojanowski, and  
708 Tomas Mikolov. 2016. Bag of tricks for efficient text  
709 classification. *arXiv preprint arXiv:1607.01759*.

710	Onur Köksal and Nurcihan Yürük. 2020. The role of translator in intercultural communication. <i>International Journal of Curriculum and Instruction</i> , 12(1):327–338.	764
711		765
712		766
713		767
714	Fangjun Kuang, Liyong Guo, Wei Kang, Long Lin, Mingshuang Luo, Zengwei Yao, and Daniel Povey. 2022. Pruned rnn-t for fast, memory-efficient asr training. <i>arXiv preprint arXiv:2206.13236</i> .	768
715		769
716		
717		
718	Taku Kudo and John Richardson. 2018. <a href="#">Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing</a> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 66–71.	770
719		771
720		772
721		773
722		774
723		
724	Junjie Li, Jing Peng, Yangui Fang, Shuai Wang, and Kai Yu. 2025. Mosa: Mixtures of simple adapters outperform monolithic approaches in llm-based multilingual asr. <i>arXiv preprint arXiv:2508.18998</i> .	775
725		776
726		777
727		778
728	Wei Liu, Jiahong Li, Yiwen Shao, and Dong Yu. 2025. Tta: Transcribe, translate and alignment for cross-lingual speech representation. <i>arXiv preprint arXiv:2511.14410</i> .	779
729		780
730		781
731		782
732	Evgeny Matusov, Stephan Kanthak, and Hermann Ney. 2005. On the integration of speech recognition and statistical machine translation. In <i>Interspeech</i> , pages 3177–3180.	783
733		784
734		785
735		786
736	B Mu, K Wei, Q Shao, Y Xu, and L Xie. 2024. Hd-mole: Mixture of lora experts with hierarchical routing and dynamic thresholds for fine-tuning llm-based asr models. <i>arXiv preprint arXiv:2409.19878</i> .	787
737		788
738		789
739		790
740	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. <a href="#">Bleu: a method for automatic evaluation of machine translation</a> . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL)</i> , pages 311–318.	791
741		792
742		793
743		794
744		795
745	Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le. 2019. <a href="#">SpecAugment: A simple data augmentation method for automatic speech recognition</a> . In <i>Proc. Interspeech</i> , pages 2613–2617.	796
746		797
747		798
748		799
749		800
750	Rastislav Rabatin, Frank Seide, and Ernie Chang. 2024. Navigating the minefield of mt beam search in cascaded streaming speech translation. <i>arXiv preprint arXiv:2407.11010</i> .	801
751		802
752		803
753		804
754	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In <i>International conference on machine learning</i> , pages 28492–28518. PMLR.	805
755		806
756		807
757		808
758		809
759	Ricardo Rei, Craig Stewart, Ana C. Farinha, and Alon Lavie. 2020. <a href="#">Comet: A neural framework for MT evaluation</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2685–2702.	810
760		811
761		812
762		813
763		814
	Monica Sekoyan, Nithin Rao Koluguri, Nune Tadevosyan, Piotr Zelasko, Travis Bartley, Nikolay Karpov, Jagadeesh Balam, and Boris Ginsburg. 2025. <a href="#">Canary-1b-v2 &amp; parakeet-tdt-0.6b-v3: Efficient and high-performance models for multilingual asr and ast</a> . <i>Preprint</i> , arXiv:2509.14128.	815
		816
		817
		818
	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. <i>arXiv preprint arXiv:1701.06538</i> .	
	Mark Sinclair, Peter Bell, Alexandra Birch, and Fergus McInnes. 2014. A semi-markov model for speech segmentation with an utterance-break prior. In <i>Interspeech</i> , pages 2351–2355.	
	Matthias Sperber, Jan Niehues, and Alex Waibel. 2017. Toward robust neural machine translation for noisy input sequences. In <i>Proceedings of the 14th International Conference on Spoken Language Translation</i> , pages 90–96.	
	Yun Tang and 1 others. 2023. Hybrid transducer and attention based encoder-decoder modeling for speech-to-text tasks. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics</i> , pages 12441–12455.	
	Peidong Wang, Eric Sun, Jian Xue, Yu Wu, Long Zhou, Yashesh Gaur, Shujie Liu, and Jinyu Li. 2023. Lamassu: A streaming language-agnostic multilingual speech recognition and translation model using neural transducers. In <i>Interspeech</i> , pages 57–61.	
	H Xue, W Ren, X Geng, K Wei, L Li, Q Shao, L Yang, K Diao, and L Xie. 2024. Ideal-llm: Integrating dual encoders and language-adapted llm for multilingual speech-to-text. <i>arXiv preprint arXiv:2409.11214</i> .	
	Brian Yan, Siddharth Dalmia, Yosuke Higuchi, Graham Neubig, Florian Metzger, Alan W Black, and Shinji Watanabe. 2023. Ctc alignments improve autoregressive translation. In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 1623–1639.	
	Zengwei Yao, Liyong Guo, Xiaoyu Yang, Wei Kang, Fangjun Kuang, Yifan Yang, Zengrui Jin, Long Lin, and Daniel Povey. 2023. Zipformer: A faster and better encoder for automatic speech recognition. <i>arXiv preprint arXiv:2310.11230</i> .	
	Zengwei Yao, Wei Kang, Xiaoyu Yang, Fangjun Kuang, Liyong Guo, Han Zhu, Zengrui Jin, Zhaoqing Li, Long Lin, and Daniel Povey. 2024. Cr-ctc: Consistency regularization on ctc for improved speech recognition. <i>arXiv preprint arXiv:2410.05101</i> .	
	Saierdaer Yusuyin, Te Ma, Hao Huang, Wenbo Zhao, and Zhijian Ou. 2024. <a href="#">Whistle: Data-efficient multilingual and crosslingual speech recognition via weakly phonetic supervision</a> . <i>CoRR</i> , abs/2406.02166.	

819 Piotr Żelasko, Daniel Povey, Jan Trmal, and Sanjeev  
820 Khudanpur. 2021. Lhotse: a speech data representa-  
821 tion library for the modern deep learning ecosystem.  
822 *arXiv preprint arXiv:2110.12561*.

823 **A Per-direction Results for Joint ASR**  
824 **and ST**

SRC \ TGT	Model	WER (%)↓									
		de	en	es	fr	it	nl	pl	pt	ro	
de	HENT-SRT-M2O×9	-	21.80	26.64	27.12	27.44	26.43	26.51	26.51	26.62	
	HENT-SRT-M2M	-	19.09	18.77	18.82	19.09	18.86	18.79	18.99	18.86	
	LCMA-SRT	-	<b>18.01</b>	<b>17.84</b>	<b>17.85</b>	<b>18.23</b>	<b>17.92</b>	<b>17.75</b>	<b>17.93</b>	<b>17.99</b>	
en	HENT-SRT-M2O×9	16.42	-	17.32	17.56	17.08	17.45	17.29	17.18	17.21	
	HENT-SRT-M2M	13.92	-	13.94	14.02	13.79	13.84	13.99	13.90	13.53	
	LCMA-SRT	<b>12.94</b>	-	<b>12.93</b>	<b>13.02</b>	<b>12.84</b>	<b>12.87</b>	<b>12.92</b>	<b>12.97</b>	<b>12.64</b>	
es	HENT-SRT-M2O×9	21.29	17.77	-	22.25	22.83	22.68	21.93	22.82	22.66	
	HENT-SRT-M2M	15.96	15.80	-	15.91	15.69	15.97	15.85	15.89	15.75	
	LCMA-SRT	<b>15.30</b>	<b>15.14</b>	-	<b>15.27</b>	<b>15.02</b>	<b>15.31</b>	<b>15.26</b>	<b>15.25</b>	<b>15.14</b>	
fr	HENT-SRT-M2O×9	19.37	16.07	19.82	-	20.41	19.30	19.45	19.86	20.80	
	HENT-SRT-M2M	13.40	13.38	13.28	-	13.42	13.36	13.39	13.38	13.37	
	LCMA-SRT	<b>12.58</b>	<b>12.51</b>	<b>12.53</b>	-	<b>12.51</b>	<b>12.50</b>	<b>12.56</b>	<b>12.55</b>	<b>12.65</b>	
it	HENT-SRT-M2O×9	18.18	15.05	19.19	19.32	-	19.06	18.60	19.00	19.91	
	HENT-SRT-M2M	13.10	13.19	13.13	13.24	-	13.17	12.98	13.18	13.27	
	LCMA-SRT	<b>12.50</b>	<b>12.41</b>	<b>12.52</b>	<b>12.63</b>	-	<b>12.59</b>	<b>12.42</b>	<b>12.62</b>	<b>12.66</b>	
nl	HENT-SRT-M2O×9	38.99	32.95	38.85	38.85	39.52	-	38.99	39.32	39.26	
	HENT-SRT-M2M	28.59	28.65	28.73	28.46	28.62	-	28.46	28.46	28.47	
	LCMA-SRT	<b>27.01</b>	<b>27.23</b>	<b>26.89</b>	<b>26.91</b>	<b>27.20</b>	-	<b>26.93</b>	<b>27.07</b>	<b>26.82</b>	
pl	HENT-SRT-M2O×9	25.89	22.01	26.33	27.19	25.99	26.47	-	27.13	27.36	
	HENT-SRT-M2M	18.26	18.27	18.14	18.21	17.87	18.27	-	18.29	18.00	
	LCMA-SRT	<b>17.54</b>	<b>17.39</b>	<b>17.32</b>	<b>17.36</b>	<b>17.01</b>	<b>17.43</b>	-	<b>17.57</b>	<b>17.11</b>	
pt	HENT-SRT-M2O×9	19.90	16.27	21.74	20.82	20.77	20.99	20.48	-	20.53	
	HENT-SRT-M2M	13.60	13.59	13.52	13.59	13.38	13.58	13.57	-	13.34	
	LCMA-SRT	<b>12.37</b>	<b>12.72</b>	<b>12.28</b>	<b>12.37</b>	<b>12.08</b>	<b>12.38</b>	<b>12.40</b>	-	<b>12.19</b>	
ro	HENT-SRT-M2O×9	22.32	15.85	21.87	22.04	23.97	22.88	23.63	22.82	-	
	HENT-SRT-M2M	14.59	14.20	14.52	14.42	14.17	14.59	14.49	14.65	-	
	LCMA-SRT	<b>13.64</b>	<b>13.29</b>	<b>13.51</b>	<b>13.46</b>	<b>13.38</b>	<b>13.61</b>	<b>13.54</b>	<b>13.72</b>	-	

Table 4: Direction-wise ASR WER (% , ↓) on the Europarl-ST test set under the same (SRC,TGT)-conditioned decoding used in joint ASR and ST. Rows denote source speech languages and columns denote target translation languages. We compare HENT-SRT-M2O×9, HENT-SRT-M2M, and LCMA-SRT. Bold indicates the lowest WER for each direction.

SRC	TGT	Model	BLEU $\uparrow$								
			de	en	es	fr	it	nl	pl	pt	ro
de		HENT-SRT-M2O $\times$ 9	-	17.5	13.3	12.1	8.7	16.2	5.9	12.4	8.3
		HENT-SRT-M2M	-	11.0	3.7	3.3	1.1	4.1	1.6	4.0	2.2
		LCMA-SRT	-	<b>22.0</b>	<b>19.7</b>	<b>20.2</b>	<b>14.5</b>	<b>19.0</b>	<b>8.9</b>	<b>18.7</b>	<b>13.5</b>
en		HENT-SRT-M2O $\times$ 9	15.4	-	26.0	24.6	19.0	21.9	9.7	23.1	19.8
		HENT-SRT-M2M	4.0	-	9.7	6.5	3.1	5.3	1.6	7.1	4.5
		LCMA-SRT	<b>20.1</b>	-	<b>33.4</b>	<b>30.7</b>	<b>25.0</b>	<b>25.4</b>	<b>14.7</b>	<b>29.4</b>	<b>26.3</b>
es		HENT-SRT-M2O $\times$ 9	9.9	22.1	-	20.2	15.7	15.1	6.9	22.4	12.2
		HENT-SRT-M2M	2.1	13.4	-	3.9	1.5	3.1	0.9	5.4	2.2
		LCMA-SRT	<b>13.7</b>	<b>26.1</b>	-	<b>26.3</b>	<b>21.0</b>	<b>19.4</b>	<b>10.3</b>	<b>26.6</b>	<b>17.7</b>
fr		HENT-SRT-M2O $\times$ 9	11.0	23.5	20.3	-	17.6	16.9	7.4	23.3	13.0
		HENT-SRT-M2M	2.9	11.9	6.4	-	2.2	4.0	1.3	6.5	2.4
		LCMA-SRT	<b>14.9</b>	<b>28.6</b>	<b>27.0</b>	-	<b>22.5</b>	<b>21.3</b>	<b>11.1</b>	<b>27.5</b>	<b>18.3</b>
it		HENT-SRT-M2O $\times$ 9	11.3	23.0	21.3	20.3	-	16.1	8.3	22.4	13.4
		HENT-SRT-M2M	2.9	14.7	5.1	4.0	-	3.2	1.7	5.6	2.0
		LCMA-SRT	<b>14.8</b>	<b>27.0</b>	<b>27.3</b>	<b>25.3</b>	-	<b>20.2</b>	<b>11.0</b>	<b>26.1</b>	<b>17.8</b>
nl		HENT-SRT-M2O $\times$ 9	7.1	15.6	11.3	10.4	7.3	-	3.7	10.4	6.3
		HENT-SRT-M2M	2.3	9.8	3.1	2.6	1.2	-	0.9	2.9	1.9
		LCMA-SRT	<b>12.1</b>	<b>21.0</b>	<b>17.6</b>	<b>16.5</b>	<b>13.6</b>	-	<b>7.0</b>	<b>16.9</b>	<b>11.6</b>
pl		HENT-SRT-M2O $\times$ 9	9.5	19.3	17.1	15.7	11.9	14.3	-	14.6	10.0
		HENT-SRT-M2M	2.4	12.1	4.6	3.8	1.6	3.4	-	3.8	2.1
		LCMA-SRT	<b>14.3</b>	<b>23.9</b>	<b>24.1</b>	<b>22.9</b>	<b>18.6</b>	<b>19.5</b>	-	<b>20.8</b>	<b>16.5</b>
pt		HENT-SRT-M2O $\times$ 9	10.9	23.7	22.1	21.3	17.3	15.6	7.5	-	13.9
		HENT-SRT-M2M	2.3	13.3	6.6	4.0	1.9	3.0	1.1	-	2.5
		LCMA-SRT	<b>15.4</b>	<b>28.1</b>	<b>28.3</b>	<b>27.0</b>	<b>22.8</b>	<b>19.7</b>	<b>10.5</b>	-	<b>19.0</b>
ro		HENT-SRT-M2O $\times$ 9	10.9	25.3	21.4	21.4	15.8	16.0	7.9	18.8	-
		HENT-SRT-M2M	1.9	16.4	4.6	3.7	1.5	2.2	0.7	3.9	-
		LCMA-SRT	<b>15.8</b>	<b>30.1</b>	<b>28.9</b>	<b>28.4</b>	<b>22.1</b>	<b>19.7</b>	<b>12.2</b>	<b>25.3</b>	-

Table 5: Direction-wise speech translation BLEU scores on Europarl-ST test set (9 languages, 72 directions). Rows denote source speech languages and columns denote target translation languages. We compare the direction-specific hierarchical transducer baseline HENT-SRT-M2O $\times$ 9, the fully shared many-to-many baseline HENT-SRT-M2M, and our unified LCMA-SRT. Bold indicates the best score for each direction.

SRC \ TGT	Model	COMET $\uparrow$								
		de	en	es	fr	it	nl	pl	pt	ro
de	HENT-SRT-M2O $\times$ 9	-	0.615	0.531	0.479	0.504	0.549	0.521	0.544	0.545
	HENT-SRT-M2M	-	0.522	0.453	0.407	0.409	0.397	0.383	0.447	0.391
	LCMA-SRT	-	<b>0.683</b>	<b>0.624</b>	<b>0.572</b>	<b>0.591</b>	<b>0.604</b>	<b>0.591</b>	<b>0.636</b>	<b>0.627</b>
en	HENT-SRT-M2O $\times$ 9	0.571	-	0.641	0.606	0.625	0.620	0.584	0.668	0.680
	HENT-SRT-M2M	0.421	-	0.533	0.470	0.487	0.430	0.419	0.524	0.458
	LCMA-SRT	<b>0.638</b>	-	<b>0.741</b>	<b>0.690</b>	<b>0.714</b>	<b>0.674</b>	<b>0.663</b>	<b>0.749</b>	<b>0.765</b>
es	HENT-SRT-M2O $\times$ 9	0.488	0.652	-	0.548	0.571	0.534	0.546	0.636	0.589
	HENT-SRT-M2M	0.357	0.536	-	0.416	0.424	0.385	0.374	0.464	0.396
	LCMA-SRT	<b>0.544</b>	<b>0.708</b>	-	<b>0.627</b>	<b>0.657</b>	<b>0.584</b>	<b>0.609</b>	<b>0.709</b>	<b>0.663</b>
fr	HENT-SRT-M2O $\times$ 9	0.499	0.685	0.603	-	0.603	0.551	0.555	0.650	0.618
	HENT-SRT-M2M	0.373	0.535	0.484	-	0.440	0.396	0.385	0.481	0.408
	LCMA-SRT	<b>0.561</b>	<b>0.737</b>	<b>0.700</b>	-	<b>0.685</b>	<b>0.603</b>	<b>0.616</b>	<b>0.723</b>	<b>0.701</b>
it	HENT-SRT-M2O $\times$ 9	0.507	0.679	0.614	0.569	-	0.551	0.568	0.650	0.623
	HENT-SRT-M2M	0.372	0.560	0.477	0.425	-	0.393	0.380	0.472	0.404
	LCMA-SRT	<b>0.560</b>	<b>0.728</b>	<b>0.698</b>	<b>0.640</b>	-	<b>0.600</b>	<b>0.619</b>	<b>0.717</b>	<b>0.686</b>
nl	HENT-SRT-M2O $\times$ 9	0.444	0.581	0.500	0.460	0.467	-	0.486	0.509	0.509
	HENT-SRT-M2M	0.367	0.508	0.435	0.397	0.402	-	0.365	0.435	0.380
	LCMA-SRT	<b>0.538</b>	<b>0.660</b>	<b>0.595</b>	<b>0.544</b>	<b>0.561</b>	-	<b>0.556</b>	<b>0.604</b>	<b>0.593</b>
pl	HENT-SRT-M2O $\times$ 9	0.515	0.643	0.568	0.518	0.545	0.543	-	0.584	0.583
	HENT-SRT-M2M	0.385	0.539	0.469	0.424	0.429	0.397	-	0.462	0.401
	LCMA-SRT	<b>0.584</b>	<b>0.709</b>	<b>0.667</b>	<b>0.612</b>	<b>0.651</b>	<b>0.608</b>	-	<b>0.683</b>	<b>0.677</b>
pt	HENT-SRT-M2O $\times$ 9	0.522	0.692	0.631	0.584	0.605	0.556	0.576	-	0.636
	HENT-SRT-M2M	0.381	0.557	0.491	0.439	0.444	0.402	0.390	-	0.409
	LCMA-SRT	<b>0.581</b>	<b>0.744</b>	<b>0.722</b>	<b>0.662</b>	<b>0.695</b>	<b>0.609</b>	<b>0.632</b>	-	<b>0.710</b>
ro	HENT-SRT-M2O $\times$ 9	0.514	0.697	0.606	0.575	0.596	0.563	0.569	0.627	-
	HENT-SRT-M2M	0.381	0.585	0.487	0.443	0.446	0.408	0.386	0.480	-
	LCMA-SRT	<b>0.587</b>	<b>0.753</b>	<b>0.711</b>	<b>0.667</b>	<b>0.696</b>	<b>0.624</b>	<b>0.642</b>	<b>0.724</b>	-

Table 6: Direction-wise speech translation COMET scores on Europarl-ST test set (9 languages, 72 directions). Rows denote source speech languages and columns denote target translation languages. We compare the direction-specific hierarchical transducer baseline HENT-SRT-M2O $\times$ 9, the fully shared many-to-many baseline HENT-SRT-M2M, and our unified LCMA-SRT. Bold indicates the best score for each direction.

SRC \ TGT	Model	LMR (%)↓								
		de	en	es	fr	it	nl	pl	pt	ro
de	HENT-SRT-M2O×9	-	<b>0.08</b>	0.70	0.64	0.66	1.00	<b>0.00</b>	3.39	<b>1.70</b>
	HENT-SRT-M2M	-	56.60	87.83	90.14	94.98	77.09	83.11	83.75	82.55
	LCMA-SRT	-	0.38	<b>0.49</b>	<b>0.43</b>	<b>0.58</b>	<b>0.84</b>	0.73	<b>2.38</b>	1.79
en	HENT-SRT-M2O×9	<b>0.00</b>	-	<b>0.79</b>	0.25	0.35	0.65	<b>0.16</b>	<b>1.98</b>	<b>1.64</b>
	HENT-SRT-M2M	78.21	-	78.93	86.08	95.93	78.54	85.14	81.62	80.55
	LCMA-SRT	0.08	-	0.95	<b>0.16</b>	<b>0.09</b>	<b>0.24</b>	0.40	<b>1.98</b>	<b>1.64</b>
es	HENT-SRT-M2O×9	<b>0.18</b>	<b>0.22</b>	-	<b>0.18</b>	<b>0.46</b>	1.01	<b>0.09</b>	<b>1.19</b>	<b>0.88</b>
	HENT-SRT-M2M	88.51	58.54	-	92.98	96.76	89.58	90.93	80.07	86.81
	LCMA-SRT	0.54	0.61	-	0.37	0.65	<b>0.92</b>	0.38	1.29	1.54
fr	HENT-SRT-M2O×9	<b>0.00</b>	<b>0.11</b>	<b>0.55</b>	-	<b>0.00</b>	<b>0.26</b>	<b>0.18</b>	<b>1.55</b>	<b>1.16</b>
	HENT-SRT-M2M	85.18	71.51	87.89	-	95.60	86.66	89.85	77.71	84.93
	LCMA-SRT	<b>0.00</b>	0.33	0.73	-	0.38	0.87	0.36	1.91	<b>1.16</b>
it	HENT-SRT-M2O×9	<b>0.11</b>	<b>0.00</b>	<b>0.57</b>	<b>0.23</b>	-	<b>0.36</b>	<b>0.25</b>	2.20	<b>0.54</b>
	HENT-SRT-M2M	88.54	57.23	92.05	94.25	-	91.61	92.77	86.59	90.38
	LCMA-SRT	<b>0.11</b>	0.42	0.91	<b>0.23</b>	-	0.84	0.37	<b>1.62</b>	0.95
nl	HENT-SRT-M2O×9	<b>0.09</b>	<b>0.34</b>	<b>0.49</b>	<b>0.49</b>	0.90	-	<b>0.21</b>	3.18	<b>1.60</b>
	HENT-SRT-M2M	78.25	56.62	88.85	89.02	93.71	-	86.96	82.91	81.98
	LCMA-SRT	0.19	0.86	1.19	0.79	<b>0.79</b>	-	0.62	<b>1.80</b>	1.94
pl	HENT-SRT-M2O×9	<b>0.00</b>	<b>0.22</b>	<b>0.72</b>	<b>0.16</b>	<b>0.42</b>	1.80	-	2.08	1.72
	HENT-SRT-M2M	82.54	60.55	89.70	91.34	96.01	86.68	-	83.45	87.29
	LCMA-SRT	<b>0.00</b>	0.40	0.80	0.48	0.85	<b>1.14</b>	-	<b>1.60</b>	<b>1.01</b>
pt	HENT-SRT-M2O×9	<b>0.00</b>	<b>0.17</b>	0.16	<b>0.16</b>	<b>0.08</b>	<b>0.41</b>	<b>0.08</b>	-	<b>0.81</b>
	HENT-SRT-M2M	88.36	65.79	86.39	91.99	97.51	90.31	90.80	-	87.09
	LCMA-SRT	0.16	0.39	<b>0.08</b>	0.24	<b>0.08</b>	0.81	0.50	-	1.17
ro	HENT-SRT-M2O×9	<b>0.00</b>	<b>0.15</b>	<b>0.42</b>	0.26	<b>0.77</b>	<b>0.91</b>	<b>0.00</b>	1.58	-
	HENT-SRT-M2M	90.98	50.74	93.27	95.33	98.46	93.14	93.04	89.42	-
	LCMA-SRT	0.16	0.56	0.50	<b>0.17</b>	<b>0.77</b>	1.24	0.43	<b>1.33</b>	-

Table 7: Direction-wise ST LMR (% , ↓) on the Europarl-ST test set (9 languages, 72 directions). Rows denote source speech languages and columns denote target translation languages. We compare HENT-SRT-M2O×9, HENT-SRT-M2M, and LCMA-SRT. Bold indicates the lowest LMR for each direction.