

Code-Switching in End-to-End Automatic Speech Recognition: A Systematic Literature Review

Maha Tufail Agro¹ Atharva Kulkarni¹ Karima Kadaoui¹
Zeerak Talat² Hanan Aldarmaki¹

¹Mohamed bin Zayed University of Artificial Intelligence, UAE

²University of Edinburgh, UK

¹{maha.agro, hanan.alldarmaki}@mbzuai.ac.ae, ²z@zeerak.org

Abstract

Motivated by a growing research interest into automatic speech recognition (ASR), and the growing body of work for languages in which code-switching (CS) often occurs, we present a systematic literature review of code-switching in end-to-end ASR models. We collect and manually annotate papers published in peer-reviewed venues. We document the languages considered, datasets, metrics, model choices, and performance, and present a discussion of challenges in end-to-end ASR for code-switching. Our analysis thus provides insights on current research efforts and available resources as well as opportunities and gaps to guide future research.

1 Introduction

Automatic speech recognition (ASR) is one of the most widely used technologies for accessible communication and device interaction. A recent market analysis estimated the global market of speech and voice technologies at USD 20 billion in 2023, of which ASR accounts for over 60%, distributed across a range of sectors, including automotive, banking and financial, healthcare, and retail ([Grand View Research, 2024](#)). Such widespread adoption of ASR highlights the need and desire for high performing systems across different languages and speech patterns. While large-scale resources for developing robust ASR systems are common for monolingual settings, datasets and methods for code-switching (CS)—the practice of alternating between two or more languages in a single conversation or discourse—remains a challenge for current ASR models, despite that it is a widespread phenomenon across the globe ([UNESCO, 2024](#)). Multilingual models may be well suited for handling CS ([Peng et al., 2023](#); [Kadaoui et al., 2024](#)), but their performance across language pairs has yet to be fully explored. Studies that seek to address CS in speech are limited by a scarcity of datasets

that are designed to evaluate this phenomenon. In this paper, we present a comprehensive review of efforts towards end-to-end (E2E) ASR systems for CS. We seek to identify the current trends in research by answering the following research questions: (1) Which language pairs, modeling, and evaluation mechanisms have been pursued? (2) What is the current state-of-the-art for E2E ASR for CS? (3) What are the persisting challenges in ASR for CS research? We find that although the body of work in E2E ASR for code-switching is growing, a small subset of languages receive the most attention, and while the number of datasets is growing, there is a large disparity in the languages that are covered by the datasets. Moreover, we find a wide variety of training and evaluation methodologies, but the efforts are mostly sporadic, with no consistent benchmarking or clarity on directions for future research. To this end, we present a discussion of these challenges around three axes: resources, methodology, and fairness. Through this work, we hope to provide a road map for more inclusive, rigorous, and reproducible research in end-to-end ASR for code-switching.

1.1 Code-switching

People engage in CS in their communication for many and varied reasons, such as a speaker’s familiarity with the setting and language and how they wish to project themselves ([Myers-Scotton, 1993](#)); to express identity or tone ([Gumperz, 1982](#)); and due to diglossia, the presence of two languages or varieties used under different conditions, such as ‘high’ and ‘low’ (colloquial) varieties ([Ferguson, 1959](#)). CS can also be a consequence of ASR systems themselves; a recent study found that older African-Americans engaged in a form of code-switching while addressing ASR systems that aren’t trained to process their vernacular ([Harrington et al., 2022](#)). In code-switching research, there is a distinction made between the

matrix language—the language which provides the structure—and the *embedded* language—which provides the foreign words or phrases (Weller et al., 2022). Code-switching is often categorized into two types: *inter-sentential* (switching languages or dialects between sentences) and *intra-sentential* code-switching, where language switches happen within a single sentence. The latter is most challenging for speech systems and what is typically addressed in the surveyed literature.

1.2 Automatic Speech Recognition

ASR research and development has a long history dating back to the 1950s. Anusuya and Katti (2010) summarize the progress made until 2009; by then, the SOTA models were hybrid HMM-DNN models, combining the controlled statistical power of HMMs for lexical and language modeling with the adaptive nature of neural networks for acoustic modeling. Subsequent research led to the development of fully connected *end-to-end* (or E2E for short) ASR systems. As pointed out in Prabhavalkar et al. (2024), the term E2E hides various practical complexities in these systems, and it is an umbrella term used to describe various neural ASR systems that are characterized (with caveats) as *joint modeling and training of ASR components from scratch in a single computational graph and objective function*. Refer to Prabhavalkar et al. (2024) for a comprehensive review of these systems and their types. Relevant to our discussion here is that these systems are currently the state-of-the-art in speech recognition, enhanced with self-supervised learning (SSL) methods that utilize large unlabeled data for better generalization. A review of these SSL methods for speech is provided in Mohamed et al. (2022). These systems achieve remarkable recognition performance across speakers, noise conditions, and even across different languages (Yadav and Sitaram, 2022).

1.3 Related Work

While there are hundreds of reviews of “speech recognition” systems (e.g., Reddy, 1976; Prabhavalkar et al., 2024), we identify only one published survey on code-switching in ASR systems (Mustafa et al., 2022). In their study, Mustafa et al. (2022) analyze a sample of papers on bilingual and multilingual ASR, including 24 papers on CS. Their search methodology differs from ours in both focus and scope, as they analyze a relatively small subset of CS papers, many of which predate

the ones in our analysis. Meanwhile, in text processing Winata et al. (2023) present a comprehensive survey of code-switching, covering decades of research in NLP. Our work complements these efforts by focusing exclusively on E2E ASR, which currently represents the mainstream in ASR technology. Furthermore, our analysis includes all papers that fit our search and inclusion criteria, which is constrained only by the recency of the studied E2E systems; indeed, over half of the papers in our survey have been published since 2022. Our survey thus provides a comprehensive and up-to-date overview of the research on CS in E2E ASR.

1.4 Outline

We present a systematic literature review of published research on E2E ASR for code-switched speech. We analyze the coverage of languages, datasets, metrics, and modeling choices, thereby presenting a comprehensive overview of the field. We first describe the scope of the study and our search and annotation methodologies in Section 2, and describe the results of our annotation in the subsequent sections. Section 3 summarizes the languages and datasets covered in this literature. We then detail ASR modeling choices in Section 4, including architectures, language identification, and decoding strategies, among other dimensions identified in our annotation scheme. Section 5 describes the training and evaluation settings used in these papers, such as data augmentation and evaluation metrics. We summarize the state-of-the-art models in Section 6, and describe new and persisting challenges in Section 7.

2 Scope & Methodology

Data Collection We gathered related papers by querying the Semantic Scholar API¹ to retrieve papers that use terms related to code-switching, published from 2014 until February 27, 2025.² We initially retrieved 378 papers, from which we only kept papers that have been published in peer-reviewed venues and describe end-to-end code-switching systems, resulting in a set of 127 papers published between 2018 and 2024 for analysis.³

¹<https://www.semanticscholar.org/product/api>

²We used the following query for retrieving papers: (ASR | Speech recognition) & (code-switch* | codeswitch* | code switch* | code-mix* | codemix*).

³Based on our search results, the first paper related to ASR for CS using an E2E architecture was published in 2018.

Dimension	Attribute	Description
Problem Setup & Data	Language(s)	<i>Languages studied (see full list in Appendix)</i>
	Datasets	<i>Datasets used (see full list in Appendix)</i>
	Dataset Accessibility	Yes No
Model Design Choices	Monolingual Modeling	Yes No
	Multilingual Modeling	Yes No
	LID	Yes No
	Text Units	Words Subwords (BPE) Characters Phones
	Architecture	<i>Type of neural architecture used for ASR</i>
	Pretrained models	Yes No
	Language Model	Yes No
	Loss Function	CTC Cross Entropy Hybrid CTC/Attention Others
Training and Evaluation Settings	Decoding Strategy	Greedy Beam search Other
	Data augmentation	<i>Type of data augmentation used</i>
	Translation	Yes No
	Zero-Shot	Yes No
Performance	Evaluation metrics	WER CER MER TER <i>etc.</i>
	Best performance	<i>Best reported performance</i>
	Best model	<i>Model with best reported performance</i>

Table 1: List of annotated attributes in the survey.

Annotation The papers were annotated by five annotators after agreeing on the annotation dimensions. The annotations were carried out during scheduled annotation sessions, where any misunderstandings and questions were addressed during the session to ensure consistency in annotation procedures. Each paper was annotated by one annotator by manually going over the paper text. We extracted information from each paper across four major dimensions:

- **Problem Setup & Data:** We recorded the languages covered (e.g., Mandarin-English, Spanish-English), the datasets used (e.g., SEAME, Bangor Miami corpus) for training and evaluation, and dataset accessibility.⁴
- **Model Design Choices:** We annotated whether monolingual or multilingual components are explicitly modeled, and note whether code-switched data are used for training. We also annotated the text units used (e.g., characters, phones, BPE), the model architecture (e.g., hybrid CTC/attention-based model with Transformer encoder), use of large pretrained models (e.g., Whisper), external language model integration (e.g., shallow fusion during decoding), loss function and the decoding strategies employed (e.g., beam search).
- **Training and Evaluation Settings:** We annotated whether the study applied any form

of data augmentation (e.g., speed perturbation, synthetic code-switching, etc.), translation (e.g., translating code-switched utterances into monolingual text), and zero-shot evaluation (i.e., evaluating multilingual models on code-switched utterances without fine-tuning). We also annotated the evaluation metrics reported (e.g., WER, CER).

- **Performance:** We annotated the best reported result in each paper and the model that achieved that best performance.

Type	Count
Modeling	88
New Dataset	19
Data Augmentation	10
New Evaluation Metric	3
Shared Task	2
Other	6

Table 2: Number of papers by category.

We also coded the type of contributions in each paper (see Table 2). The majority of the surveyed papers describe empirical research focused on model design choices, such as architecture or training methodology. The ‘Other’ category includes papers that do not fit the other categories, such as papers that examine text encodings, data partitions, or analysis of existing models.

3 Languages & Datasets

We find a steady increase in publications over time with more than 38 datasets currently avail-

⁴We understand an ‘accessible’ dataset to be one where a link is made available, it can be obtained upon request, or it is used in subsequent research.

Dataset	Matrix Language	Speech Type	Total (hrs)	Code-Switched Speech Duration (hr)		
				Train	Dev	Test
TALCS (Li et al., 2022)	zho	Spontaneous	587.0	555.9	8.0	23.6
ASCEND (Lovenia et al., 2022)	zho	Spontaneous	10.6	8.8	0.9	0.9
ASRU 2019 (Shi et al., 2020)	zho	Read	740.0	200.0	40.0	-
SEAME (Lyu et al., 2010)	zho	Spontaneous	192.0	101.1	11.4	-
KSC2 (Mussakhojayeva et al., 2022)	kaz	Read	1127.9	26.3	0.5	0.5
Hartanto CS (Roosadi and Lestari, 2023)	ind	-	7.15	6.0	-	1.15
VITB-HEBiC (Jain and Bhowmick, 2024)	hin	Read	7.5	-	-	-
MUCS 2021 (Diwan et al., 2021)	hin, ben	Spontaneous	600.0	135.9	-	12.2
IITG-HingCos (Sreeram et al., 2019)	hin	Read	25.0	-	-	-
Mixat (Ali and Aldarmaki, 2024)	ara	Spontaneous	15.0	-	-	-
ZAEBUC-Spoken (Hamed et al., 2024)	ara	Spontaneous	12.0	-	-	-
TunSwitch CS (Abdallah et al., 2024)	ara, fra	Spontaneous	163.62	8.5	0.15	0.25
HAC (Hamed et al., 2022b)	ara	Spontaneous	2.0	-	-	-
ESCWA-CS (Chowdhury et al., 2021)	ara	Spontaneous	2.8	-	-	-
QASR-CS (Mubarak et al., 2021)	ara	Read	5.9	-	-	-
ArzEn (Hamed et al., 2020)	ara	Spontaneous	12.0	-	-	-
FAME (Yilmaz et al., 2015)	fry-nld	Spontaneous	14	11.5	1.2	1.2
DECM (Ugan et al., 2024)	deu	Spontaneous	3.38	-	-	3.38
German Spoken Wikipedia Corpus (Khosravani et al., 2021)	deu	Read	34	-	-	-
Miami Bangor (Deuchar et al., 2014)	spa	Spontaneous	35.66	-	-	-
South African Soap Operas (van der Westhuizen and Niesler, 2018)	xho, sot, tsn, zul	Spontaneous	14.3	2.5	0.2	0.7
Two Sepedi SPCS Corpus (Modipa and Davel, 2022)	nso	Read	10	-	-	-

Table 3: Comprehensive list of publicly available datasets found in the surveyed literature. Embedded language is English unless otherwise specified.

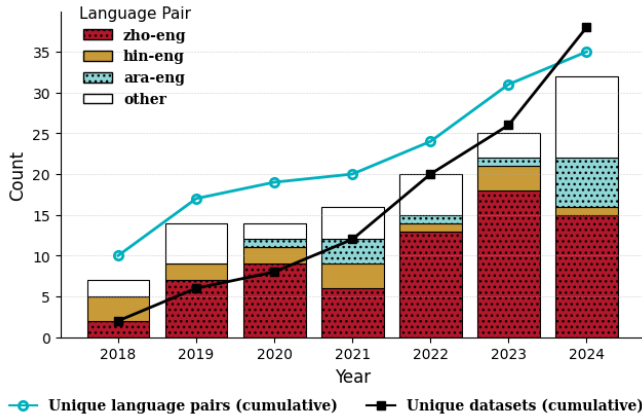


Figure 1: Total number of papers per year, and the number of papers on each of the top 3 language-pairs. The lines show the cumulative total of unique language pairs and unique datasets covered over time.

able over 35 unique language pairs; however, the majority of languages appear only in a single dataset. The research is dominated by three language pairs accounting for $\sim 76\%$ of all papers: Mandarin-English (zho-eng), Hindi-English (hin-eng), and Arabic-English⁵ (ara-eng) (see Figure 1). Mandarin-English is the most studied language pair and is covered by $\sim 55\%$ of all papers. Research on Mandarin-English is partially driven by the availability of large datasets such as SEAME (Lyu et al., 2010) and the ASRU 2019 challenge (Shi et al., 2020). The availability of datasets appears

⁵Here, we collate several dialects of Arabic.

to be one of the main drivers of research in this area; $\sim 77\%$ of papers make use of accessible datasets, while the rest use proprietary or unspecified datasets. For a full listing of languages and datasets, see Table 5 in the Appendix.

Below we briefly summarize frequently used accessible datasets in the literature grouped by the matrix language. Unless otherwise specified, the embedded language is English.⁶ See Table 3 for additional details on all datasets we identify as ‘accessible’.

3.1 Mandarin

Datasets for Mandarin code-switched data include: TALCS (Li et al., 2022), ASCEND (Lovenia et al., 2022), the ASRU 2019 Code-Switching Challenge dataset (Shi et al., 2020), and the Southeast Asia Mandarin-English (SEAME) corpus (Lyu et al., 2010), mainly covering the Mandarin-English pair. Among these, SEAME—which contains Mandarin- and Singaporean-accented speech—is the most frequently used ($\sim 24\%$) followed by ASRU 2019 ($\sim 11\%$). Despite being the largest dataset, TALCS is only used in $\sim 4\%$ of the papers.

⁶We use ISO-3 language codes and full language names interchangeably to represent language pairs. See Table 5 in the Appendix for a complete mapping.

3.2 Indic Languages

Datasets covering Indic languages include MUCS (Diwan et al., 2021) and ITTG-HingCos (Sreeram et al., 2019) which primarily cover Hindi-English (hin-eng) and Bengali-English (ben-eng). As the second most studied language group, $\sim 11\%$ of papers address hin-eng and $\sim 3\%$ cover ben-eng. Of papers studying code-switching with Indic languages, $\sim 38\%$ rely on the MUCS dataset—which includes monolingual speech in 7 Indic languages as well as CS speech hin-eng and ben-eng language pairs. In contrast, ITTG-HingCos focuses solely on the hin-eng pair and captures a range of dialectal variation, as its speakers are sourced from multiple states across India.

3.3 Arabic

Code-switching datasets for Arabic primarily embed English or French, with a few cases of intra-Arabic dialectal CS (i.e., code-switching between Arabic dialects).⁷ We identify 8 datasets for Arabic in the literature—most explicitly designed for CS research—more than any other language. Mixat (Ali and Aldarmaki, 2024) is the largest dataset, covering Emirati Arabic-English, followed by ArzEn (Hamed et al., 2020) which covers Egyptian Arabic-English code-switching.

3.4 African Languages

The most commonly used dataset for African languages in our survey is the Multilingual Code-switched Soap Opera Speech (van der Westhuizen and Niesler, 2018), which contains ~ 3.6 hours of code-switched speech of a total of 14.3 hours. The dataset contains code-switched speech for isiXhosa, isiZulu, Setswana, and Sesotho—four Bantu languages—paired with English, and occasional instances of code-switching between the Bantu languages.⁸

3.5 Japanese

Japanese-English (jpn-eng) is another common language pair in the survey, with about 5% of papers focusing on it. Interestingly, none of these works use or release a dedicated jpn-eng code-switched dataset. Instead, most rely on synthetic data, typically using the BTEC (Takezawa et al., 2007)

corpus—a parallel Japanese-English dataset for machine translation—to create code-switched segments. These datasets are proprietary, so no large-scale, publicly available jpn-eng code-switched dataset appears in the surveyed literature.

4 ASR Modeling Choices

We now summarize the findings related to how code-switching has been modeled in end-to-end ASR research. We present statistics and summarize findings related to: monolingual and multilingual modeling (Section 4.1), language identification, Section 4.2), text units (Section 4.3), architectures (Section 4.4), large pretrained models (Section 4.5), decoding schema (Section 4.6), and the use of language models (Section 4.7).

4.1 Monolingual vs. Multilingual Models

Monolingual modeling has become an important method for optimizing CS performance, even in medium and high-resource settings. We find that roughly 60% of surveyed papers included additional monolingual data for training their models—half of which include datasets with more than 100 hours of CS speech. Other efforts turned towards multilingual modeling to encode both languages in the latent space. More than a third of surveyed paper include multilingual components, where two or more languages are modeled as part of the same ASR architecture. These models can be appealing because they capture shared acoustic and lexical patterns across languages, and they obviate the need for separate language dependent acoustic models and pronunciation dictionaries.

Effect of monolingual training: Due to the limited availability of large-scale CS training data, which remains a major challenge in building effective CS ASR systems for most languages, many researchers model each language separately using monolingual data and introduce CS data at a later stage, effectively utilizing transfer learning (Yue et al., 2019; Yang et al., 2024; Wang et al., 2024). For example, Wang et al. (2024) propose a tri-stage training strategy for a two-pass E2E Mandarin-English ASR system. They use large monolingual Mandarin and English corpora to pretrain two symmetric, language-specific encoders. The pretrained representations are then combined using a feed-forward neural network and the combined system is retained using monolingual data, followed by fine-tuning on a CS corpus.

⁷We observe the following Arabic dialects in the literature: Egyptian, Emirati, MSA, Gulf, Jordanian, Mauritanian, Palestinian and Tunisian.

⁸Bantu is a branch of the Niger-Congo language family spoken in central, eastern, and southern Africa.

Effect of CS on monolingual performance:

While monolingual modeling can have positive impacts on CS, [Shah et al. \(2020\)](#) argue that the relationship is not necessarily isomorphic. They argue that fine-tuning pretrained monolingual models for CS may impair performance on monolingual datasets due to CS data contributing to catastrophic forgetting. They propose the Learning Without Forgetting (LWF) mechanism for when the monolingual model is available but its training data is not; and a regularization method when the monolingual training data and model are available. LWF applies a knowledge distillation loss to retain the original model’s predictions during CS training, which consistently outperforms simple fine-tuning on monolingual and CS test sets; while the regularization method minimizes the KL divergence between the output distribution of the pretrained and fine-tuned models.

Effect of multilingual training [Seki et al. \(2018, 2019\)](#) find that multilingual methods may outperform language-dependent models (i.e., models with language-specific ASR modules and a separate LID module) to optimize both CS and monolingual performance. They propose a monolithic multilingual ASR system that jointly performs language identification and transcription, which introduces dynamic language tracking within an utterance to handle intrasentential code-switching.

4.2 Language ID

We find that LID plays a supporting role in many E2E ASR systems for CS in our survey—approximately 33% of papers incorporate LID in some stage of their training process, either as a pre-processing step to separate the languages ([Lu et al., 2020](#)), as part of the loss function ([Zeng et al., 2019](#); [Shan et al., 2019](#); [Zhao et al., 2024](#)), as an auxiliary task ([Qiu et al., 2020](#)), or as part of the predicted output sequence ([Seki et al., 2018, 2019](#); [Zhang et al., 2021](#)). For example, [Lu et al. \(2020\)](#) propose a bi-encoder transformer network based Mixture-of-Experts (MoE) architecture, which uses an external LID to route language-specific data to the corresponding language-specific encoder to avoid cross-lingual contamination of the encoder modules. In contrast, [Shan et al. \(2019\)](#) propose a Multi-Task Learning (MTL) setup where they analyze the impact of LID loss at different points in the network and find that attention-related components (like the attention output and the attention hidden

state) yield the best results compared to adding the loss in decoder hidden state, suggesting that acoustic and alignment information is more useful for LID than the language modeling information in the decoder.

Other work has experimented with predicting LID at each time step as an auxiliary task alongside character predictions ([Qiu et al., 2020](#); [Shan et al., 2019](#)). [Shan et al. \(2019\)](#) find that this frame-level LID prediction, implemented via MTL, outperforms predicting LID token at the beginning of the sequence. They argue that the latter approach is more appropriate for inter-sentential switching than intra-sentential switching, as it does not capture spontaneous switching within sentences.

4.3 Text Units

The choice of text units in code-switching ASR, whether shared or distinct at phonetic levels or word/character levels, impacts the model’s ability to capture linguistic nuances and transitions between languages. A small subset of surveyed papers (~5%) used a common label set shared across languages. For instance, [Seki et al. \(2018, 2019\)](#) construct a common character set by taking the union of characters from all involved languages. While simple, this approach suffers from greater confusion between cross-language targets—due to similar acoustic realization of different symbols, and higher computational costs caused by the larger vocabulary. Seeking to address this issue, [Dhawan et al. \(2020\)](#) and [Sreeram and Sinha \(2020\)](#) propose a reduced set of target phones based on acoustic similarity using IITG-HingCos. Despite the reduction in the size of the target set, [Sreeram and Sinha \(2020\)](#) report an increase in WER which they attribute to ambiguities among homophones. [Li et al. \(2019\)](#) instead propose using Unicode bytes as output units which affords a fixed-size and language independent vocabulary, thereby avoiding the need to modify the softmax output layer due to new characters or languages. Nearly all papers on Mandarin-English CS adopt mixed units, where the former is encoded using characters and the latter using BPE (e.g., [Zhang et al., 2021](#); [Long et al., 2021](#); [Yan et al., 2023](#)).

4.4 Architectures

While using CTC is attractive for CS due to the output independence assumption, auto-regressive models with encoder decoder architectures are shown to perform better in [Peng et al. \(2022\)](#). Almost half

of the surveyed papers (47%) employ an encoder-decoder architecture, with RNN/LSTM (Zeng et al., 2019; Ma et al., 2019) being most common before 2020, and Transformer architectures being more common recently (J et al., 2020; Nga et al., 2023; Sailor et al., 2021; Kronis et al., 2024). Around 17% of the papers relied on pretrained models, such as Wav2Letter2+ (Naowarat et al., 2021), XLS-R (Ogunremi et al., 2023; Akhi and Arefin, 2024), Wav2Vec2 (Chen et al., 2024), and most commonly Whisper (Ugan et al., 2024; Kim et al., 2024; Alharbi et al., 2024); these are discussed in more details in the next section. The remaining papers used CTC-based encoder only models (Chuang et al., 2021; Tian et al., 2022), RNN or Transformer based transducers (Dalmia et al., 2021; Zhang et al., 2021), or Mixture of Experts (Lu et al., 2020; Ma et al., 2023; Yang et al., 2024).

4.5 Large Multi-Lingual Pretrained Models

The use of pretrained language models is another common feature across our sample. In particular, ~17% of papers use Whisper (Radford et al., 2023), Wav2Vec 2.0 (Baeviski et al., 2020), or XLS-R (Conneau et al., 2021). We identify three notable directions in work that use pretrained models: (i) methods that focus on distillation and parameter efficient fine-tuning (PeFT), such as Tseng et al. (2024) who use a filtering method to distill Whisper Large V2 by 50% and speed up generation five-fold, while outperforming the teacher model by 30% in some settings, and Kim et al. (2024) present Gated Low Rank Adaption, a weight separation-based PeFT method to enable the use of pretrained models on low-spec devices, e.g., mobile phones. (ii) Methods that perform multilingual fine-tuning of pretrained models—specifically XLSR—for low-resource languages, e.g., Southern Bantu languages from South Africa (Ogunremi et al., 2023), Kichwa (Taguchi et al., 2024) from the South America and Bangla (Akhi and Arefin, 2024) from Bangladesh. Finally, (iii) simple prompting techniques have been explored. For example, Penagarikano et al. (2023) show that by concatenating LID tags with the input prompt, Whisper can be used for zero-shot CS detection for Mandarin-English despite not officially supporting code-switching.

4.6 Decoding Strategies

Most papers do not explicitly mention the decoding strategy. As greedy decoding is the simplest strategy and is commonly used in end-to-end models,

we assume that most papers follow this strategy (only seven papers explicitly mention greedy decoding). Around 30% of papers explicitly mention a strategy other than greedy decoding, approximately 70% of which used some form of beam search. Other papers have used multi-graph decoding (Yue et al., 2019; Ali et al., 2021), or other techniques. For instance, Yue et al. (2019) introduce a multi-graph decoding strategy to address data imbalance between high-resourced and low-resourced languages. They construct parallel Weighted Finite-State Transducer (WFST) search graphs, one for each monolingual language and one for the bilingual setup (Frisian-Dutch), while sharing the acoustic and lexicon models, which allows the decoder to dynamically select the most appropriate path during inference. However, this setup does not support intra-sentence switching, as decoding is limited to a fixed monolingual or bilingual path. In contrast, Ali et al. (2021) propose a similar multi-graph decoding approach but incorporate a Kleene-closure, enabling transitions between language graphs and allowing intra-sentence switching during decoding.

4.7 Language Model

Around 45% papers report using a language model (LM) in their ASR pipeline to improve decoding accuracy through shallow fusion or rescore. While 20% use some form of n-gram LMs (Tian et al., 2022; Ma et al., 2023; Srivastava and Sitaram, 2018), the remainder use RNN-based LMs (Yue et al., 2019; Sharma et al., 2020; Zhou et al., 2020) or transformer-based LMs (Liu et al., 2023, 2024; Chen et al., 2023a). Li and Vu (2020) experiment with a word-based LM (one-layer LSTM) and a subword-based LM (two-layer RNN) for Mandarin-English code-switched ASR. The models are trained on natural and synthetic ASR and integrated into the ASR using shallow fusion. They find that the subword-based LM trained on SEAME and augmented with CycleGAN-generated (Zhu et al., 2017) text achieves the highest performance. Chen et al. (2023b) investigate the capability of T5 (Raffel et al., 2020), MT5 (Xue et al., 2021) and PaLM (Chowdhery et al., 2023) to rescore hypotheses generated by a ASR model for long-form speech recognition in US English and Indian English. The LLMs compute log-likelihoods of ASR outputs and improve performance through rescore with optional segment-level context. The paper finds that MT5 rescore—when fine-tuned on CS

data—improves WER over neural and maximum entropy based baselines.

5 Training & Evaluation

5.1 Data Augmentation

Code-switching is a low-resource setting due to the small size of datasets for code-switching and the limited coverage of language pairs. Unsurprisingly, $\sim 30\%$ of papers explicitly mention data augmentation in their methodology. SpecAugment (Park et al., 2019), speed perturbation, and TTS synthesis (Sharma et al., 2020) constitute the most commonly used data augmentation techniques.

While SpecAugment and speed perturbation remain the main techniques used, some methods have been developed specifically for CS speech. Sharma et al. (2020) synthesize Hindi-English CS speech and apply Mixup regularization (Zhang et al., 2018) to bridge the distribution gap between synthetic and real data. Results on proprietary data show reductions up to 5% absolute WER. Du et al. (2021) find mixed results in an exploration of TTS-based augmentation. They experiment with audio splicing, random noun/verb translation, and random English word insertion into Mandarin sentences followed by TTS. While the combination of all three methods along with SpecAugment yielded the best overall performance, SpecAugment still outperformed them individually. Liang et al. (2022) compare numerous data augmentation techniques, namely pitch shifting, speed perturbation, audio codec augmentation, SpecAugment, and a TTS approach combining transcripts from one dataset and style IDs from another. Experiments on a Conformer favor TTS and speed perturbation, while audio codec and pitch shifting can have an adverse effect on performance. Other methods (Chi et al., 2023; Hussein et al., 2024) explore CS speech generation from monolingual data, using grid beam search (Hokamp and Liu, 2017) and audio splicing with energy normalization.

5.2 Translation

17 papers use translation in their methodology to augment their language model data (Penagarikano et al., 2023), to improve model evaluation (Taguchi et al., 2024; Kadaoui et al., 2024), or to use TTS synthesis to augment their speech data (Yu et al., 2023; Tazakka et al., 2024).

For example, Yu et al. (2023) use machine translation to generate a parallel dataset of sentences,

then perform word alignment and substitution to ensure the synthetic data mirrors statistics of CS in the natural data. They convert their synthesized text into speech using a multilingual TTS system to augment their training data. Their method obtains a 16% relative reduction in error rate for Mandarin-English. Tazakka et al. (2024) take a similar approach to perform semi-supervised training for Indonesian-English.

5.3 Zero-Shot Evaluation

Given the scarcity of CS speech data and the prevalence of multilingual ASR systems, zero-shot evaluation appears to be an attractive choice for CS ASR. While using large multilingual pretrained models, such as Whisper presents as an obvious choice, Whisper is also only trained with monolingual data. However, Peng et al. (2023) find that concatenating the language tokens for the languages in a CS scenario improves Whisper performance, despite it not being explicitly trained to receive such combined tokens. Ugan et al. (2024) conduct a zero-shot evaluation of different ASR models on their German-English DECM dataset. They use the massively multilingual speech dataset (MMS, Pratap et al., 2024) and two large Whisper models, one of which is used with a German decoding prefix (WhisperDe). Whisper performs best overall with WhisperDe close behind, while MMS has almost a 12% higher WER. Further analysis suggests that WhisperDe’s lower overall performance stems from its superior German monolingual performance, rather than CS in particular. Zhou et al. (2024b) adapt kNN-CTC (Zhou et al., 2024a) to enable zero-shot Chinese-English ASR through the use of dual monolingual datastores of labeled examples. This setup ensures that retrieved examples with similar audio frames come from the appropriate language datastore during decoding, with the help of a gating mechanism. This approach outperforms CTC fine-tuning and even small Whisper variants (Radford et al., 2023; Peng et al., 2023) that are almost double the parameter size. In Yan et al. (2023), two monolingual modules transcribe all segments using their respective vocabularies, resulting in transliterations for foreign segments. This delays the CS boundary decision to a textual bilingual module that generates a bilingual sequence from the previous outputs, thereby mitigating error-propagation compared to earlier methods (Tian et al., 2022; Song et al., 2022).

Dataset	Best Model	Year	Monolingual	Multilingual	LM	LID	Translation	Augmentation	Zero-shot	Architecture	Languages	Metric	Best Result
SEAME	(Aditya et al., 2024)	2024	✓							Whisper	zho-eng	MER	14.2
ASRU	(Wang et al., 2023)	2023	✓			✓		✓		MoE	sgp-eng	MER	20.8
TALCS	(Wang et al., 2024)	2024	✓	✓						Enc	zho-eng	MER	8.2
ASCEND	(Tseng et al., 2024)	2024	✓					✓		Whisper	zho-eng	MER	6.17
MUCS	(Kumar et al., 2021)	2021		✓	✓	✓	✓			Enc-Dec	hin-eng	WER	17.86
											ben-eng	WER	22.0
ArzEn	(Hamed et al., 2022a)	2022	✓	✓	✓			✓		Non E2E	ara-eng	WER	27.8
Mixat	(Kadaoui et al., 2024)	2024					✓			Whisper	ara-eng	WER	32.1
Miami Bangor	(Hillah et al., 2024)	2024		✓		✓				Whisper	spa-eng	WER	24.8
ESCWA	(Chowdhury et al., 2021)	2021		✓						Enc-Dec	ara-eng	WER	48.38
												WER	37.7

Table 4: Best performing models on popular datasets.

5.4 Evaluation Metrics

We now turn to discussing notable metrics used in the surveyed work (see Table 6 in the Appendix for a full listing of metrics).

Standard Metrics: The standard metric for evaluating ASR accuracy is the Word Error Rate (WER). WER computes the edit distance between a reference and prediction and divides the total number of substitutions (S), insertions (I), and deletions (D) by the total number of words in the reference (N): $WER = (S + I + D)/N$.

Character Error Rate (CER) offers a higher level of granularity by employing the same principle of WER but on a character-level. It is particularly useful for languages that lack word boundaries such as Chinese, Japanese and Thai.

Mixed Metrics: For code-switching, other metrics such as Mixed Error Rate (MER) and Token Error Rate (TER) (Zhou et al., 2024b; Shen and Guo, 2022) have been proposed due to the nature of code-switching that often involves different writing systems, e.g., where one language employs word boundaries but another does not (such as Japanese-English). MER combines WER for word-based language segments and CER for character-based segments within the same transcription, while TER evaluates prediction at the token-level, be it a character for a language like Chinese or a sub-word unit for languages like English.

Transliteration-Based Metrics: Other metrics first transliterate to a single writing system, then compute performance to address the ambiguous distinction between code-switching and loan words in cases where the matrix and embedded languages use different scripts and predictions may be unfairly penalized for being written in the wrong

script. For example, Transliteration-Optimized WER (toWER, Emond et al., 2018) first transliterates using a weighted finite state transducer then computes WER on the “transliterated space”, and the Pronunciation-Optimized Word Error Rate (poWER, Srivastava and Sitaram, 2018) process the reference and prediction with a grapheme-to-phoneme conversion for each word to ensure that predicted words that are faithful to the audio are not penalized. Kadaoui et al. (polyWER, 2024) instead propose a modification to the edit distance algorithm to enable multi-reference evaluation. They apply PolyWER across two dimensions: transliteration with a CER threshold for allowed variations in orthography, or translation with a cosine similarity threshold to allow synonyms.

6 Best Performing Models

We summarize the modeling choices corresponding to the state-of-the-art approaches on popular datasets in Table 4. Note that there is no single methodology shared across all datasets. Roughly half use monolingual modeling, and half use multilingual modeling, with only two SOTA systems utilizing both approaches. The use of LM, LID, and data augmentation is relatively small, in spite of efforts showing the potential of these approaches. None of the SOTA models rely on zero-shot evaluation, which underscores the importance of dedicated CS training. The choice of architecture also varies, but Encoder-Decoder architectures (including Whisper) is most common. In one of the datasets (ArzEn), the explored E2E models do not even outperform the sequential baseline, which is a hybrid CNN-TDNN model.

7 Discussion: Challenges & Opportunities

Research in code-switching is undergoing rapid methodological innovation and resource development. Here, we present a discussion of challenges and opportunities as drawn from our analysis of the surveyed literature.

Data Scarcity: Research in this area is largely driven by data availability, which has led to a stronger focus on higher-resourced languages. For instance, Mandarin-English code-switching has attracted the most attention in end-to-end ASR due to the pre-existing datasets and well-defined benchmarks. Moreover, while there is an upward trend in the publication of new resources, they remain concentrated around a small subset of languages, thereby privileging some languages while leaving many others under-resourced and under-studied. Adding to the challenge of resource availability is that a large proportion of datasets used in the surveyed literature are proprietary, and therefore are not available to the wider research community, e.g., for further studies or replicating results. The development of publicly available training and evaluation data thus presents an opportunity for researchers to help steer the future of research towards under-represented languages.

Disparities in coverage: While ASR research is moving towards inclusive and fair representation with coverage of hundreds of languages, research in code-switching in ASR remains confined to a small subset of languages. Moreover, dialectal variations are only partially represented for the Arabic and Indic language families. Considering the colloquial nature of code switching, dialectal variations also need to be taken into account, with a focus on marginalized communities who often need to code-switch in order to utilize such systems. As ASR technology continues to be deployed in critical domains such as accessibility, healthcare, and education, this disparity in coverage risks exacerbating existing inequalities.

Evaluation: The recent introduction of new evaluation metrics provides further avenues for precise evaluation of end-to-end ASR systems for code-switching, but they also underscore the existing conceptual challenges in ASR evaluation and the shortcomings of existing metrics. These challenges arise due to differences between scripts (e.g., dif-

ferent levels of tokenization) and norms (e.g., standardization of spelling and transliteration). This calls for additional efforts to examine current evaluation practices and propose metrics that are valid (i.e., consistent, interpretable, and meaningful) (Jacobs and Wallach, 2021; Delobelle et al., 2024) for a grounded assessment of progress.

Systematic studies & Reproducibility: The rapid methodological innovation in the field has led to a lack of insight into (i) the robustness of proposed methods, (ii) their applicability across language pairs, and (iii) the interaction of different methods, e.g., the impact of data augmentation methods on different architectures or language pairs. Furthermore, with the exception of a small set of languages, the efforts appear sporadic and lack a consistent evaluation framework, benchmarking, and continuity to compare proposed methods with past efforts. As a result, positive findings from earlier studies have often not been replicated on newer datasets, leading to performance gaps and uncertainty about the robustness and validity of prior work. These gaps present opportunities for future research, such as the development of benchmarks to standardize experimental methodology and studies focused on replication and comparison.

Data Augmentation: Data augmentation lies at the intersection of resource and methodological challenges. In the surveyed literature, data augmentation, e.g., through code-switched speech synthesis and audio perturbation, has been frequently employed as a redress to concerns of data sparsity. While increasing data resources for a particular language may mitigate the data sparsity issue, data augmentation can still prove useful for other purposes, e.g., to ensure model robustness and generalization. Furthermore, the availability of monolingual data will likely continue to outpace that of code-switching resources, which remain limited in part due to the challenges of annotation. The field could benefit from further innovation in data augmentation techniques and comprehensive studies in the utility of data augmentation across different language pairs and settings.

8 Conclusion

Research on code-switching in end-to-end ASR is experiencing a boom with an increasing number of papers published in recent years. In this paper, we conducted an extensive literature review

of 127 papers published on the topic of E2E ASR for code-switched speech. By annotating and analyzing these papers, we sought to identify and explicate developments and trends in this body of research across languages and datasets, modeling choices, and training and evaluation methods, and presented a discussion of existing and potential challenges identified through our literature review. We note a significant bias towards language pairs with established datasets and benchmarks, underscoring the importance of resources and standardization in driving research. Through this study, we hope to encourage future research towards more inclusive and replicable research.

References

- Ahmed Amine Ben Abdallah, Ata Kabboudi, Amir Kannon, and Salah Zaiem. 2024. [Leveraging data collection and unsupervised learning for code-switched Tunisian Arabic automatic speech recognition](#). In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2024, Seoul, Republic of Korea, April 14-19, 2024*, pages 12607–12611. IEEE.
- Bobbi Aditya, Mahdin Rohmatillah, Liang-Hsuan Tai, and Jen-Tzung Chien. 2024. [Attention-guided adaptation for code-switching speech recognition](#). In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2024, Seoul, Republic of Korea, April 14-19, 2024*, pages 10256–10260. IEEE.
- Fatema Tuz Zohra Akhi and Mohammad Shamsul Arfin. 2024. [Transformer based Bangla-English code-switching speech recognition and language identification model](#). In *2024 IEEE International Conference on Computing, Applications and Systems (COMPAS)*, pages 1–5.
- Sadeen Alharbi, Reem Binmuqbil, Ahmed Ali, Raghad Aloraini, Saiful Bari, Areeb Alowisheq, and Yaser Alonaizan. 2024. [Leveraging LLM for augmenting textual data in code-switching ASR: Arabic as an example](#). *Proceedings of SynData4GenAI*.
- Ahmed Ali, Shammur Absar Chowdhury, Amir Hussein, and Yasser Hifny. 2021. [Arabic Code-Switching Speech Recognition Using Monolingual Data](#). pages 3475–3479.
- Maryam Al Ali and Hanan Aldarmaki. 2024. [Mixat: A data set of bilingual Emirati-English speech](#). *CoRR*, abs/2405.02578.
- M. A. Anusuya and S. K. Katti. 2010. [Speech recognition by machine, A review](#). *CoRR*, abs/1001.2267.
- Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. [wav2vec 2.0: A framework for self-supervised learning of speech representations](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Peikun Chen, Fan Yu, Yuhao Liang, Hongfei Xue, Xucheng Wan, Naijun Zheng, Huan Zhou, and Lei Xie. 2023a. [BA-MoE: Boundary-Aware Mixture-of-Experts Adapter for Code-Switching Speech Recognition](#). *IEEE Conference Publication | IEEE Xplore*. In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 1–7.
- Po-Kai Chen, Li-Yeh Fu, Cheng-Kai Chen, Yi-Xing Lin, Chih-Ping Chen, Chien-Lin Huang, and Jia-Ching Wang. 2024. [Self-supervised learning and masked language model for code-switching automatic speech recognition](#). In *2024 Tenth International Conference on Communications and Electronics (ICCE)*, pages 387–391.
- Tongzhou Chen, Cyril Allauzen, Yinghui Huang, Daniel Park, David Rybach, W. Ronny Huang, Rodrigo Cabrera, Kartik Audhkhasi, Bhuvana Ramabhadran, Pedro J. Moreno, and Michael Riley. 2023b. [Large-scale language model rescoring on long-form data](#). In *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5.
- Jie Chi, Brian Lu, Jason Eisner, Peter Bell, Preethi Jyothi, and Ahmed M. Ali. 2023. [Unsupervised code-switched text generation from parallel text](#). In *24th Annual Conference of the International Speech Communication Association, Interspeech 2023, Dublin, Ireland, August 20-24, 2023*, pages 1419–1423. ISCA.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2023. [PaLM: Scaling language modeling with pathways](#). *J. Mach. Learn. Res.*, 24:240:1–240:113.

- Shammur Absar Chowdhury, Amir Hussein, Ahmed Abdelali, and Ahmed Ali. 2021. [Towards one model to rule all: Multilingual strategy for dialectal code-switching Arabic ASR](#). In *22nd Annual Conference of the International Speech Communication Association, Interspeech 2021, Brno, Czechia, August 30 - September 3, 2021*, pages 2466–2470. ISCA.
- Shun-Po Chuang, Heng-Jui Chang, Sung-Feng Huang, and Hung-yi Lee. 2021. [Non-autoregressive Mandarin-English code-switching speech recognition](#). In *IEEE Automatic Speech Recognition and Understanding Workshop, ASRU 2021, Cartagena, Colombia, December 13-17, 2021*, pages 465–472. IEEE.
- Alexis Conneau, Alexei Baevski, Ronan Collobert, Abdelrahman Mohamed, and Michael Auli. 2021. [Un-supervised cross-lingual representation learning for speech recognition](#). In *22nd Annual Conference of the International Speech Communication Association, Interspeech 2021, Brno, Czechia, August 30 - September 3, 2021*, pages 2426–2430. ISCA.
- Siddharth Dalmia, Yuzong Liu, Srikanth Ronanki, and Katrin Kirchhoff. 2021. [Transformer-transducers for code-switched speech recognition](#). In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2021, Toronto, ON, Canada, June 6-11, 2021*, pages 5859–5863. IEEE.
- Pieter Delobelle, Giuseppe Attanasio, Debora Nozza, Su Lin Blodgett, and Zeerak Talat. 2024. [Metrics for What, Metrics for Whom: Assessing Actionability of Bias Evaluation Metrics in NLP](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 21669–21691, Miami, Florida, USA. Association for Computational Linguistics.
- Margaret Deuchar, Peredur Davies, Jon Russell Herring, M. Carmen Parafita Couto, and Diana Carter. 2014. [5. Building Bilingual Corpora](#), pages 93–110. Multilingual Matters, Bristol, Blue Ridge Summit.
- Kunal Dhawan, Ganji Sreeram, Kumar Priyadarshi, and Rohit Sinha. 2020. [Investigating target set reduction for end-to-end speech recognition of Hindi-English code-switching data](#). In *2020 National Conference on Communications, NCC 2020, Kharagpur, India, February 21-23, 2020*, pages 1–5. IEEE.
- Anuj Diwan, Rakesh Vaideeswaran, Sanket Shah, Ankita Singh, Srinivasa Raghavan K. M., Shreya Khare, Vinit Unni, Saurabh Vyas, Akash Rajpuria, Chiranjeevi Yarra, Ashish R. Mittal, Prasanta Kumar Ghosh, Preethi Jyothi, Kalika Bali, Vivek Shashtri, Sunayana Sitaram, Samarth Bharadwaj, Jai Nanavati, Raoul Nanavati, and Karthik Sankaranarayanan. 2021. [MUCS 2021: Multilingual and code-switching ASR challenges for low resource Indian languages](#). In *22nd Annual Conference of the International Speech Communication Association, Interspeech 2021, Brno, Czechia, August 30 - September 3, 2021*, pages 2446–2450. ISCA.
- Chenpeng Du, Hao Li, Yizhou Lu, Lan Wang, and Yanmin Qian. 2021. [Data augmentation for end-to-end code-switching speech recognition](#). In *IEEE Spoken Language Technology Workshop, SLT 2021, Shenzhen, China, January 19-22, 2021*, pages 194–200. IEEE.
- Jesse Emond, Bhuvana Ramabhadran, Brian Roark, Pedro J. Moreno, and Min Ma. 2018. [Transliteration based approaches to improve code-switched speech recognition performance](#). In *2018 IEEE Spoken Language Technology Workshop, SLT 2018, Athens, Greece, December 18-21, 2018*, pages 448–455. IEEE.
- Charles A. Ferguson. 1959. [Diglossia](#). *WORD*, 15(2):325–340.
- Grand View Research. 2024. Voice and speech recognition market size, share & trends analysis report by function (speech recognition, voice recognition), by technology, by vertical, by region, and segment forecasts, 2024 - 2030. <https://www.grandviewresearch.com/industry-analysis/voice-recognition-market>. Accessed: 2025-05-25.
- John J. Gumperz. 1982. *Discourse Strategies*. Studies in Interactional Sociolinguistics. Cambridge University Press.
- Injy Hamed, Pavel Denisov, Chia-Yu Li, Mohamed Elmahdy, Slim Abdennadher, and Ngoc Thang Vu. 2022a. Investigations on speech recognition systems for low-resource dialectal arabic–english code-switching speech. *Computer Speech & Language*, 72:101278.
- Injy Hamed, Fadhil Eryani, David Palfreyman, and Nizar Habash. 2024. [ZAEBUC-spoken: A multilingual multidialectal Arabic-English speech corpus](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/COLING 2024, 20-25 May, 2024, Torino, Italy*, pages 17770–17782. ELRA and ICCL.
- Injy Hamed, Amir Hussein, Oumnia Chellah, Shammur Absar Chowdhury, Hamdy Mubarak, Sunayana Sitaram, Nizar Habash, and Ahmed Ali. 2022b. [Benchmarking evaluation metrics for code-switching automatic speech recognition](#). In *IEEE Spoken Language Technology Workshop, SLT 2022, Doha, Qatar, January 9-12, 2023*, pages 999–1005. IEEE.
- Injy Hamed, Ngoc Thang Vu, and Slim Abdennadher. 2020. [ArzEn: A speech corpus for code-switched Egyptian Arabic-English](#). In *Proceedings of The 12th Language Resources and Evaluation Conference, LREC 2020, Marseille, France, May 11-16, 2020*, pages 4237–4246. European Language Resources Association.
- Christina N. Harrington, Radhika Garg, Amanda T. Woodward, and Dimitri Williams. 2022. [“It’s Kind](#)

- of Like Code-Switching”: Black Older Adults’ Experiences with a Voice Assistant for Health Information Seeking. In *CHI ’22: CHI Conference on Human Factors in Computing Systems, New Orleans, LA, USA, 29 April 2022 - 5 May 2022*, pages 604:1–604:15. ACM.
- Leopold Hillah, Mateusz Dubiel, and Luis A Leiva. 2024. “¿te vienes? sure!” joint fine-tuning of language detection and transcription improves automatic recognition of code-switching speech. In *Proceedings of the 6th ACM Conference on Conversational User Interfaces*, pages 1–7.
- Chris Hokamp and Qun Liu. 2017. [Lexically constrained decoding for sequence generation using grid beam search](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1535–1546, Vancouver, Canada. Association for Computational Linguistics.
- Amir Hussein, Dorsa Zeinali, Ondrej Klejch, Matthew Wiesner, Brian Yan, Shammur Absar Chowdhury, Ahmed Ali, Shinji Watanabe, and Sanjeev Khudanpur. 2024. [Speech collage: Code-switched audio generation by collaging monolingual corpora](#). In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2024, Seoul, Republic of Korea, April 14-19, 2024*, pages 12006–12010. IEEE.
- Metilda Sagaya Mary N. J, Vishwas M. Shetty, and Srinivasan Umesh. 2020. [Investigation of methods to improve the recognition performance of Tamil-English code-switched data in transformer framework](#). In *2020 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2020, Barcelona, Spain, May 4-8, 2020*, pages 7889–7893. IEEE.
- Abigail Z. Jacobs and Hanna Wallach. 2021. [Measurement and Fairness](#). In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 375–385, Virtual Event Canada. ACM.
- Palash Jain and Anirban Bhowmick. 2024. Vitb-hebic: A bilingual corpus for evaluating asr in diverse indian code-switching scenarios. *Applied Acoustics*, 224:110119.
- Karima Kadaoui, Maryam Al Ali, Hawau Olamide Toyin, Ibrahim Mohammed, and Hanan Aldarmaki. 2024. [PolyWER: A holistic evaluation framework for code-switched speech recognition](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 6144–6153, Miami, Florida, USA. Association for Computational Linguistics.
- Abbas Khosravani, Philip N. Garner, and Alexandros Lararidis. 2021. [An evaluation benchmark for automatic speech recognition of German-English code-switching](#). In *IEEE Automatic Speech Recognition and Understanding Workshop, ASRU 2021, Cartagena, Colombia, December 13-17, 2021*, pages 811–816. IEEE.
- Gwantae Kim, Bokyeung Lee, Donghyeon Kim, and Hanseok Ko. 2024. [Gated low-rank adaptation for personalized code-switching automatic speech recognition on the low-spec devices](#). In *IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP 2024 - Workshops, Seoul, Republic of Korea, April 14-19, 2024*, pages 760–764. IEEE.
- Martins Kronis, Askars Salimbajevs, and Mārcis Pinnis. 2024. [Code-mixed text augmentation for Latvian ASR](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 3469–3479, Torino, Italia. ELRA and ICCL.
- Mari Ganesh Kumar, Jom Kuriakose, Anand Thyagachandran, Ashish Seth, Lodagala Durga Prasad, Saish Jaiswal, Anusha Prakash, Hema Murthy, et al. 2021. Dual script e2e framework for multilingual and code-switching asr. *arXiv preprint arXiv:2106.01400*.
- Bo Li, Yu Zhang, Tara N. Sainath, Yonghui Wu, and William Chan. 2019. [Bytes are all you need: End-to-end multilingual speech recognition and synthesis with bytes](#). In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2019, Brighton, United Kingdom, May 12-17, 2019*, pages 5621–5625. IEEE.
- Chengfei Li, Shuhao Deng, Yaoping Wang, Guangjing Wang, Yaguang Gong, Changbin Chen, and Jinfeng Bai. 2022. [TALCS: an open-source Mandarin-English code-switching corpus and a speech recognition baseline](#). In *23rd Annual Conference of the International Speech Communication Association, Interspeech 2022, Incheon, Korea, September 18-22, 2022*, pages 1741–1745. ISCA.
- Chia-Yu Li and Ngoc Thang Vu. 2020. [Improving code-switching language modeling with artificially generated texts using cycle-consistent adversarial networks](#). In *21st Annual Conference of the International Speech Communication Association, Interspeech 2020, Virtual Event, Shanghai, China, October 25-29, 2020*, pages 1057–1061. ISCA.
- Yuhao Liang, Peikun Chen, Fan Yu, Xinfu Zhu, Tianyi Xu, Yingying Gao, and Lei Xie. 2022. [The NPU-ASLP system for the ISCSLP 2022 Magichub code-switching ASR challenge](#). In *13th International Symposium on Chinese Spoken Language Processing, ISCSLP 2022, Singapore, December 11-14, 2022*, pages 532–536. IEEE.
- Hexin Liu, Leibny Paola García, Xiangyu Zhang, Andy W. H. Khong, and Sanjeev Khudanpur. 2024. [Enhancing code-switching speech recognition with interactive language biases](#). In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2024, Seoul, Republic of Korea, April 14-19, 2024*, pages 10886–10890. IEEE.

- Hexin Liu, Haihua Xu, Leibny Paola García, Andy W. H. Khong, Yi He, and Sanjeev Khudanpur. 2023. [Reducing language confusion for code-switching speech recognition with token-level language diarization](#). In *IEEE International Conference on Acoustics, Speech and Signal Processing ICASSP 2023, Rhodes Island, Greece, June 4-10, 2023*, pages 1–5. IEEE.
- Yanhua Long, Shuang Wei, Jie Lian, and Yijie Li. 2021. [Pronunciation augmentation for Mandarin-English code-switching speech recognition](#). *EURASIP J. Audio Speech Music. Process.*, 2021(1):34.
- Holy Lovenia, Samuel Cahyawijaya, Genta Indra Winata, Peng Xu, Yan Xu, Zihan Liu, Rita Frieske, Tiezheng Yu, Wenliang Dai, Elham J. Barezi, Qifeng Chen, Xiaojuan Ma, Bertram E. Shi, and Pascale Fung. 2022. [ASCEND: A spontaneous Chinese-English dataset for code-switching in multi-turn conversation](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference, LREC 2022, Marseille, France, 20-25 June 2022*, pages 7259–7268. European Language Resources Association.
- Yizhou Lu, Mingkun Huang, Hao Li, Jiaqi Guo, and Yanmin Qian. 2020. [Bi-encoder transformer network for Mandarin-English code-switching speech recognition using mixture of experts](#). In *21st Annual Conference of the International Speech Communication Association, Interspeech 2020, Virtual Event, Shanghai, China, October 25-29, 2020*, pages 4766–4770. ISCA.
- Dau-Cheng Lyu, Tien Ping Tan, Engsiong Chng, and Haizhou Li. 2010. [SEAME: a Mandarin-English code-switching speech corpus in south-east asia](#). In *11th Annual Conference of the International Speech Communication Association, INTERSPEECH 2010, Makuhari, Chiba, Japan, September 26-30, 2010*, pages 1986–1989. ISCA.
- Guodong Ma, Wenxuan Wang, Yuke Li, Yuting Yang, Binbin Du, and Haoran Fu. 2023. [LAE-ST-MOE: boosted language-aware encoder using speech translation auxiliary task for E2E code-switching ASR](#). In *IEEE Automatic Speech Recognition and Understanding Workshop, ASRU 2023, Taipei, Taiwan, December 16-20, 2023*, pages 1–8. IEEE.
- Min Ma, Bhuvana Ramabhadran, Jesse Emond, Andrew Rosenberg, and Fadi Biadsy. 2019. [Comparison of data augmentation and adaptation strategies for code-switched automatic speech recognition](#). In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2019, Brighton, United Kingdom, May 12-17, 2019*, pages 6081–6085. IEEE.
- Thipe I. Modipa and Marelle H. Davel. 2022. [Two Sepedi-English code-switched speech corpora](#). *Lang. Resour. Evaluation*, 56(3):703–727.
- Abdelrahman Mohamed, Hung-yi Lee, Lasse Borgholt, Jakob D. Havtorn, Joakim Edin, Christian Igel, Katrin Kirchhoff, Shang-Wen Li, Karen Livescu, Lars Maaløe, Tara N. Sainath, and Shinji Watanabe. 2022. [Self-supervised speech representation learning: A review](#). *IEEE J. Sel. Top. Signal Process.*, 16(6):1179–1210.
- Hamdy Mubarak, Amir Hussein, Shammur Absar Chowdhury, and Ahmed Ali. 2021. [QASR: QCRI Aljazeera speech resource A large scale annotated Arabic speech corpus](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 2274–2285. Association for Computational Linguistics.
- Saida Mussakhojayeva, Yerbolat Khassanov, and Huseyin Atakan Varol. 2022. [KSC2: an industrial-scale open-source Kazakh speech corpus](#). In *23rd Annual Conference of the International Speech Communication Association, Interspeech 2022, Incheon, Korea, September 18-22, 2022*, pages 1367–1371. ISCA.
- Mumtaz Begum Mustafa, Mansoor Ali Yusoo, Hasan Kahtan Khalaf, Ahmad Abdel Rahman Mahmoud Abushariah, Miss Laiha Mat Kiah, Hua Nong Ting, and Saravanan Muthaiyah. 2022. [Code-switching in automatic speech recognition: The issues and future directions](#). *Applied Sciences*, 12(19).
- Carol Myers-Scotton. 1993. *Social Motivations For Codeswitching: Evidence from Africa*. Oxford University Press.
- Burin Naowarat, Thananchai Kongthaworn, Korrawe Karunratanakul, Sheng Hui Wu, and Ekapol Chuangsuwanich. 2021. [Reducing spelling inconsistencies in code-switching ASR using contextualized CTC loss](#). In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2021, Toronto, ON, Canada, June 6-11, 2021*, pages 6239–6243. IEEE.
- Cao Hong Nga, Duc-Quang Vu, Huong Hoang Luong, Chien-Lin Huang, and Jia-Ching Wang. 2023. [Cyclic transfer learning for Mandarin-English code-switching speech recognition](#). *IEEE Signal Process. Lett.*, 30:1387–1391.
- Tolulope Ogunremi, Christopher Manning, and Dan Jurafsky. 2023. [Multilingual self-supervised speech representations improve the speech recognition of low-resource African languages with codeswitching](#). In *Proceedings of the 6th Workshop on Computational Approaches to Linguistic Code-Switching*, pages 83–88, Singapore. Association for Computational Linguistics.
- Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le. 2019. [SpecAugment: A simple data augmentation method for automatic speech recognition](#). In *20th Annual Conference of the International Speech Communication Association, Interspeech 2019, Graz, Austria, September 15-19, 2019*, pages 2613–2617. ISCA.

- Mikel Penagarikano, Amparo Varona, Germán Bordel, and Luis Javier Rodriguez-Fuentes. 2023. [Semisupervised speech data extraction from basque parliament sessions and validation on fully bilingual Basque–Spanish ASR](#). *Applied Sciences*, 13(14).
- Puyuan Peng, Brian Yan, Shinji Watanabe, and David Harwath. 2023. [Prompting the hidden talent of web-scale speech models for zero-shot task generalization](#). In *24th Annual Conference of the International Speech Communication Association, Interspeech 2023, Dublin, Ireland, August 20-24, 2023*, pages 396–400. ISCA.
- Yizhou Peng, Jicheng Zhang, Haihua Xu, Hao Huang, and Eng Siong Chng. 2022. [Minimum word error training for non-autoregressive transformer-based code-switching ASR](#). In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2022, Virtual and Singapore, 23-27 May 2022*, pages 7807–7811. IEEE.
- Rohit Prabhavalkar, Takaaki Hori, Tara N. Sainath, Ralf Schlüter, and Shinji Watanabe. 2024. [End-to-end speech recognition: A survey](#). *IEEE ACM Trans. Audio Speech Lang. Process.*, 32:325–351.
- Vineel Pratap, Andros Tjandra, Bowen Shi, Paden Tomasello, Arun Babu, Sayani Kundu, Ali Elkahky, Zhaoheng Ni, Apoorv Vyas, Maryam Fazel-Zarandi, Alexei Baevski, Yossi Adi, Xiaohui Zhang, Wei-Ning Hsu, Alexis Conneau, and Michael Auli. 2024. [Scaling speech technology to 1,000+ languages](#). *Journal of Machine Learning Research*, 25(97):1–52.
- Zimeng Qiu, Yiyuan Li, Xinjian Li, Florian Metze, and William M. Campbell. 2020. [Towards context-aware end-to-end code-switching speech recognition](#). In *21st Annual Conference of the International Speech Communication Association, Interspeech 2020, Virtual Event, Shanghai, China, October 25-29, 2020*, pages 4776–4780. ISCA.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. [Robust speech recognition via large-scale weak supervision](#). In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 28492–28518. PMLR.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *J. Mach. Learn. Res.*, 21:140:1–140:67.
- Dabbala Rajagopal Reddy. 1976. Speech recognition by machine: A review. *Proceedings of the IEEE*, 64(4):501–531.
- Hizkia Raditya Pratama Roosadi and Dessi Puji Lestari. 2023. [Handling of Indonesian-English codeswitching speech in end-to-end indonesian speech recognition system using connectionist temporal classification model](#). In *International Conference on Electrical Engineering and Informatics, ICEEI 2023, Bandung, Indonesia, October 10-11, 2023*, pages 1–6. IEEE.
- Hardik B. Sailor, Kiran Praveen, Vikas Agrawal, Abhinav Jain, and Abhishek Pandey. 2021. [SRI-B end-to-end system for multilingual and code-switching ASR challenges for low resource Indian languages](#). In *22nd Annual Conference of the International Speech Communication Association, Interspeech 2021, Brno, Czechia, August 30 - September 3, 2021*, pages 2456–2460. ISCA.
- Hiroshi Seki, Takaaki Hori, Shinji Watanabe, Jonathan Le Roux, and John R. Hershey. 2019. [End-to-end multilingual multi-speaker speech recognition](#). In *20th Annual Conference of the International Speech Communication Association, Interspeech 2019, Graz, Austria, September 15-19, 2019*, pages 3755–3759. ISCA.
- Hiroshi Seki, Shinji Watanabe, Takaaki Hori, Jonathan Le Roux, and John R. Hershey. 2018. [An end-to-end language-tracking speech recognizer for mixed-language speech](#). In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2018, Calgary, AB, Canada, April 15-20, 2018*, pages 4919–4923. IEEE.
- Sanket Shah, Basil Abraham, Gurunath Reddy M, Sunayana Sitaram, and Vikas Joshi. 2020. [Learning to recognize code-switched speech without forgetting monolingual speech recognition](#). *CoRR*, abs/2006.00782.
- Changhao Shan, Chao Weng, Guangsen Wang, Dan Su, Min Luo, Dong Yu, and Lei Xie. 2019. [Investigating end-to-end speech recognition for mandarin-english code-switching](#). In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2019, Brighton, United Kingdom, May 12-17, 2019*, pages 6056–6060. IEEE.
- Yash Sharma, Basil Abraham, Karan Taneja, and Preethi Jyothi. 2020. [Improving low resource code-switched ASR using augmented code-switched TTS](#). In *21st Annual Conference of the International Speech Communication Association, Interspeech 2020, Virtual Event, Shanghai, China, October 25-29, 2020*, pages 4771–4775. ISCA.
- Zhijie Shen and Wu Guo. 2022. [An improved deliberation network with text pre-training for code-switching automatic speech recognition](#). In *23rd Annual Conference of the International Speech Communication Association, Interspeech 2022, Incheon, Korea, September 18-22, 2022*, pages 3854–3858. ISCA.
- Xian Shi, Qiangze Feng, and Lei Xie. 2020. [The ASRU 2019 Mandarin-English code-switching speech recognition challenge: Open datasets, tracks, methods and results](#). *CoRR*, abs/2007.05916.

- Tongtong Song, Qiang Xu, Meng Ge, Longbiao Wang, Hao Shi, Yongjie Lv, Yuqin Lin, and Jianwu Dang. 2022. [Language-specific characteristic assistance for code-switching speech recognition](#). In *Interspeech 2022*, pages 3924–3928.
- Ganji Sreeram, Kunal Dhawan, and Rohit Sinha. 2019. [IITG-HingCoS corpus: A Hinglish code-switching database for automatic speech recognition](#). *Speech Commun.*, 110:76–89.
- Ganji Sreeram and Rohit Sinha. 2020. [Exploration of end-to-end framework for code-switching speech recognition task: Challenges and enhancements](#). *IEEE Access*, 8:68146–68157.
- Brij Mohan Lal Srivastava and Sunayana Sitaram. 2018. [Homophone identification and merging for code-switched speech recognition](#). In *19th Annual Conference of the International Speech Communication Association, Interspeech 2018, Hyderabad, India, September 2-6, 2018*, pages 1943–1947. ISCA.
- Chihiro Taguchi, Jefferson Saransig, Dayana Velásquez, and David Chiang. 2024. [Killkan: The automatic speech recognition dataset for Kichwa with morphosyntactic information](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/COLING 2024, 20-25 May, 2024, Torino, Italy*, pages 9753–9763. ELRA and ICCL.
- Toshiyuki Takezawa, Gen-ichiro Kikui, Masahide Mizushima, and Eiichiro Sumita. 2007. [Multilingual spoken language corpus development for communication research](#). *Int. J. Comput. Linguistics Chin. Lang. Process.*, 12(3).
- Rais Vaza Man Tazakka, Dessi Lestari, Ayu Purwarianti, Dipta Tanaya, Kurniawati Azizah, and Sakriani Sakti. 2024. [Indonesian-English code-switching speech recognition using the machine speech chain based semi-supervised learning](#). In *Proceedings of the 3rd Annual Meeting of the Special Interest Group on Under-resourced Languages @ LREC-COLING 2024*, pages 143–148, Torino, Italia. ELRA and ICCL.
- Jinchuan Tian, Jianwei Yu, Chunlei Zhang, Yuexian Zou, and Dong Yu. 2022. [LAE: language-aware encoder for monolingual and multilingual ASR](#). In *23rd Annual Conference of the International Speech Communication Association, Interspeech 2022, Incheon, Korea, September 18-22, 2022*, pages 3178–3182. ISCA.
- Liang-Hsuan Tseng, Zih-Ching Chen, Wei-Shun Chang, Cheng-Kuang Lee, Tsung-Ren Huang, and Hung-Yi Lee. 2024. [Leave no knowledge behind during knowledge distillation: Towards practical and effective knowledge distillation for code-switching ASR using realistic data](#). In *IEEE Spoken Language Technology Workshop, SLT 2024, Macao, December 2-5, 2024*, pages 118–125. IEEE.
- Enes Yavuz Ugan, Ngoc-Quan Pham, and Alexander Waibel. 2024. [DECM: Evaluating bilingual ASR performance on a code-switching/mixing benchmark](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 4468–4475, Torino, Italia. ELRA and ICCL.
- UNESCO. 2024. [Multilingual education: A key to quality and inclusive learning](#).
- Ewald van der Westhuizen and Thomas Niesler. 2018. [A first South African corpus of multilingual code-switched soap opera speech](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Wenxuan Wang, Guodong Ma, Yuke Li, and Binbin Du. 2023. [Language-routing mixture of experts for multilingual and code-switching speech recognition](#). pages 1389–1393.
- Xuefei Wang, Yuan Jin, Fenglong Xie, and Yanhua Long. 2024. [Tri-stage training with language-specific encoder and bilingual acoustic learner for code-switching speech recognition](#). *Applied Acoustics*, 218:109883.
- Orion Weller, Matthias Sperber, Telmo Pires, Hendra Setiawan, Christian Gollan, Dominic Telaar, and Matthias Paulik. 2022. [End-to-end speech translation for code switched speech](#). In *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 1435–1448. Association for Computational Linguistics.
- Genta Winata, Alham Fikri Aji, Zheng Xin Yong, and Tamar Solorio. 2023. [The decades progress on code-switching research in NLP: A systematic survey on trends and challenges](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 2936–2978, Toronto, Canada. Association for Computational Linguistics.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. [mT5: A massively multilingual pre-trained text-to-text transformer](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 483–498. Association for Computational Linguistics.
- Hemant Yadav and Sunayana Sitaram. 2022. [A survey of multilingual models for automatic speech recognition](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference, LREC 2022, Marseille, France, 20-25 June 2022*, pages 5071–5079. European Language Resources Association.
- Brian Yan, Matthew Wiesner, Ondrej Klejch, Preethi Jyothi, and Shinji Watanabe. 2023. [Towards zero-](#)

- shot code-switched speech recognition. In *IEEE International Conference on Acoustics, Speech and Signal Processing ICASSP 2023, Rhodes Island, Greece, June 4-10, 2023*, pages 1–5. IEEE.
- Tzu-Ting Yang, Hsin-Wei Wang, Yi-Cheng Wang, Chi-Han Lin, and Berlin Chen. 2024. [An effective mixture-of-experts approach for code-switching speech recognition leveraging encoder disentanglement](#). In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2024, Seoul, Republic of Korea, April 14-19, 2024*, pages 11226–11230. IEEE.
- Emre Yılmaz, Maaïke Andringa, Sigrid Kingma, Frits van der Kuip, Hans Van de Velde, Frederik Kampstra, Joke Algra, Henk van den Heuvel, and David van Leeuwen. 2015. FAME! - the Frisian audio mining enterprise.
- Haibin Yu, Yuxuan Hu, Yao Qian, Ma Jin, Linqun Liu, Shujie Liu, Yu Shi, Yanmin Qian, Edward Lin, and Michael Zeng. 2023. [Code-switching text generation and injection in Mandarin-English ASR](#). In *IEEE International Conference on Acoustics, Speech and Signal Processing ICASSP 2023, Rhodes Island, Greece, June 4-10, 2023*, pages 1–5. IEEE.
- Xianghu Yue, Grandee Lee, Emre Yılmaz, Fang Deng, and Haizhou Li. 2019. [End-to-end code-switching ASR for low-resourced language pairs](#). In *IEEE Automatic Speech Recognition and Understanding Workshop, ASRU 2019, Singapore, December 14-18, 2019*, pages 972–979. IEEE.
- Zhiping Zeng, Yerbolat Khassanov, Van Tung Pham, Haihua Xu, Eng Siong Chng, and Haizhou Li. 2019. [On the end-to-end solution to Mandarin-English code-switching speech recognition](#). In *20th Annual Conference of the International Speech Communication Association, Interspeech 2019, Graz, Austria, September 15-19, 2019*, pages 2165–2169. ISCA.
- Hongyi Zhang, Moustapha Cissé, Yann N. Dauphin, and David Lopez-Paz. 2018. [mixup: Beyond empirical risk minimization](#). In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.
- Shuai Zhang, Jiangyan Yi, Zhengkun Tian, Jianhua Tao, and Ye Bai. 2021. [Rnn-transducer with language bias for end-to-end Mandarin-English code-switching speech recognition](#). In *12th International Symposium on Chinese Spoken Language Processing, ISCSLP 2021, Hong Kong, January 24-27, 2021*, pages 1–5. IEEE.
- Jiahui Zhao, Hao Shi, Chenrui Cui, Tianrui Wang, Hexin Liu, Zhaocheng Ni, Lingxuan Ye, and Longbiao Wang. 2024. [Adapting Whisper for code-switching through encoding refining and language-aware decoding](#). *CoRR*, abs/2412.16507.
- Jiaming Zhou, Shiwan Zhao, Yaqi Liu, Wenjia Zeng, Yong Chen, and Yong Qin. 2024a. [KNN-CTC: enhancing ASR via retrieval of CTC pseudo labels](#). In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2024, Seoul, Republic of Korea, April 14-19, 2024*, pages 11006–11010. IEEE.
- Jiaming Zhou, Shiwan Zhao, Hui Wang, Tian-Hao Zhang, Haoqin Sun, Xuechen Wang, and Yong Qin. 2024b. [Improving zero-shot chinese-english code-switching ASR with kNN-CTC and gated monolingual datastores](#). *CoRR*, abs/2406.03814.
- Xinyuan Zhou, Emre Yılmaz, Yanhua Long, Yijie Li, and Haizhou Li. 2020. [Multi-encoder-decoder transformer for code-switching speech recognition](#). In *21st Annual Conference of the International Speech Communication Association, Interspeech 2020, Virtual Event, Shanghai, China, October 25-29, 2020*, pages 1042–1046. ISCA.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. 2017. [Unpaired image-to-image translation using cycle-consistent adversarial networks](#). In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2242–2251. IEEE Computer Society.

A Lists of Datasets, Languages, Metrics and Papers


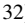

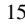

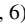

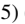

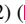

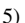

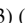


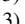

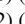


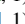

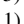


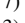

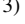


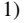


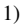
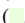





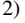

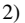







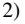






































Language Pair	ISO-3 Code	# Papers	Dataset Breakdown
Mandarin–English	zho-eng	70	( , 32) ( , 15) ( , 6) ( , 5) ( , 2) ( , 1) ( , 1) ( , 1) ( , 1) ( , 11)
Hindi–English	hin-eng	15	( , 5) ( , 3) ( , 1) ( , 1) ( , 5)
Arabic–English	ara-eng	12	( , 3) ( , 2) ( , 2) ( , 1) ( , 1) ( , 1) ( , 1) ( , 1) ( , 1)
Japanese–English	jpn-eng	7	( , 7)
Bengali–English	ben-eng	5	( , 3) ( , 2)
German–English	deu-eng	4	( , 1) ( , 1) ( , 2)
Korean–English	kor-eng	3	( , 1) ( , 1) ( , 1)
Spanish–English	spa-eng	3	( , 2) ( , 1)
French–English	fra-eng	3	( , 3)
Italian–English	ita-eng	2	( , 2)
Dutch–English	nld-eng	2	( , 2)
Portuguese–English	por-eng	2	( , 2)
Russian–English	rus-eng	2	( , 2)
Thai–English	tha-eng	2	( , 1) ( , 1)
Arabic–French	ara-fra	2	( , 1) ( , 1) ( , 1)
Indonesian–English	ind-eng	2	( , 2)
isiZulu–English	zul-eng	2	( , 2)
isiXhosa–English	xho-eng	2	( , 2)
Sesotho–English	sot-eng	2	( , 2)
Setswana–English	tsn-eng	2	( , 2)
Sepedi–English	nso-eng	2	( , 2) ( , 1)
Frisian–Dutch	fry-nld	1	FAME!
French–Mandarin	fra-zho	1	( , 1)
Japanese–Mandarin	jpn-zho	1	( , 1)
Tamil–English	tam-eng	1	( , 1)
Cantonese–English	yue-eng	1	( , 1)
Marathi–English	mar-eng	1	( , 1)
Kazakh–English	kaz-eng	1	KSC2
Basque–Spanish	eus-spa	1	Bilignual Basque-Spanish Dataset
Manipuri–English	mni-eng	1	MECOS
Quechua–Spanish	que-spa	1	Killkan
Urdu–English	urd-eng	1	Roman Urdu Code-Mixed Dataset
Latvian–English	lav-eng	1	( , 1)

Table 5: Number of papers per language pair and dataset usage breakdown.

Dataset Legend:

	SEAME		ASRU 2019		TALCS		ASCEND		MagicData-RAMC		ASRU 2021				
	AICUBES Competition Dataset		ISCSLP 2022 CSASR challenge		NTUT-AB01										
	MUCS		IITG-HingCoS		IIITH-HE-CM		VITB-HEBiC								
	ArzEn		ESCWA		Mixat		QASR.CS		HAC		TunSwitch		ZAEBUC-Spoken		Casablanca
 Saudilang Code-switch Corpus (SCC)															
	German Spoken Wikipedia Corpus		DECM		MEDREC		KECS		Bangor Miami		LOTUS-BI				
	Hartanto CS		South African Soap Operas		Two Sepedi-English Code-Switched Speech Corpora		Proprietary Data								

Metric	Languages	Description
WER	spa, kor, hin, fry-nld, nld, cmn, tha	Word-level ED
CER	jpn, cmn, deu, spa, fra, ita, nld, rus, por	Character-level ED
TER	jpn, cmn	Token-level ED
poWER	hin	WER + g2p conversion
toWER	hin, ben	WER + WFST transliteration
MER	cmn	WER (segmental) + CER (syllabic)
PPL	cmn	Token-level perplexity
WpER	jpn	Wordpiece ED
UER	kor	Unit-level ED (e.g. jamo, syllable)
SER	kor	Sentence-level ED
LER	cmn, eng, jpn, deu, spa, fra, ita, nld, rus, por	LID error rate

Table 6: Evaluation metrics and the languages they were employed for in our survey. Embedded language is English unless specified otherwise. ED: Edit distance.

DOI	Year	Type	Dataset	Languages
10.21437/interspeech.2019-1429	2018	Modeling	SEAME	zho-eng
10.1109/ICASSP.2019.8682674	2018	Modeling	Proprietary Data	jpn-eng
10.21437/Interspeech.2018-1171	2018	Modeling	Proprietary Data	hin-eng
10.1109/ICSDA.2018.8693044	2018	Dataset	Proprietary Data	jpn-eng
10.1109/SLT.2018.8639699	2018	Evaluation Metric	Proprietary Data	hin-eng
10.1109/ICASSP.2018.8462180	2018	Modeling	Proprietary Data	zho-eng, jpn-eng, deu-eng, spa-eng, fra-eng, ita-eng, nld-eng, rus-eng, por-eng
10.21437/SLTU.2018-23	2018	Dataset	IIITH-HE-CM	hin-eng
10.1109/ASRU46091.2019.9004035	2019	Modeling	FAME	fry-nld
10.1109/ICASSP.2019.8682824	2019	Other	Proprietary Data	ben-eng
10.21437/Interspeech.2019-1867	2019	Modeling	SEAME	zho-eng
10.1109/ASRU46091.2019.9003926	2019	Modeling	Proprietary Data	jpn-eng, jpn-zho, zho-eng, fra-eng, fra-zho
10.1109/IALP48816.2019.9037688	2019	Modeling	SEAME	zho-eng
10.18653/v1/K19-1026	2019	Data Augmentation	SEAME	zho-eng
10.1109/O-COCOSDA46868.2019.9060847	2019	Modeling	Proprietary Data	jpn-eng
10.1109/ICASSP.2019.8683223	2019	Modeling	Proprietary Data	zho-eng
10.1109/NCC48643.2020.9056083	2019	Modeling	IITG-HingCoS	hin-eng
10.1109/O-COCOSDA46868.2019.9041195	2019	Dataset	LOTUS-BI	tha-eng
10.1016/J.SPECOM.2019.04.007	2019	Dataset	IITG-HingCoS	hin-eng
10.21437/Interspeech.2020-2440	2019	Modeling	MEDREC	kor-eng
10.1109/ICASSP.2019.8682850	2019	Modeling	Proprietary Data	zho-eng
10.21437/INTERSPEECH.2019-3038	2019	Modeling	Proprietary Data	jpn-eng, zho-eng, deu-eng, fra-eng, ita-eng, nld-eng, por-eng, rus-eng
10.18653/v1/2020.acl-main.348	2020	Modeling	SEAME	zho-eng
10.1109/ICASSP39728.2021.9413806	2020	Modeling	Proprietary Data	tha-eng
10.1109/ICASSP39728.2021.9413562	2020	Modeling	SEAME	zho-eng
10.21437/Interspeech.2020-2402	2020	Data Augmentation	Proprietary Data	hin-eng
10.21437/Interspeech.2020-2177	2020	Data Augmentation	SEAME	zho-eng
10.21437/Interspeech.2020-2485	2020	Modeling	ASRU 2019	zho-eng
10.1109/ICASSP40776.2020.9054138	2020	Modeling	Proprietary Data	tam-eng
10.1109/SLT48900.2021.9383620	2020	Data Augmentation	ASRU 2019	zho-eng
10.1109/ISCSLP49672.2021.9362075	2020	Modeling	SEAME	zho-eng
aclanthology.org/2020.lrec-1.523/	2020	Dataset	ArzEn	ara(egyptian)-eng
10.21437/Interspeech.2020-2488	2020	Modeling	SEAME	zho-eng
10.1109/ACCESS.2020.2986255	2020	Modeling	IITG-HingCoS	hin-eng
10.1109/ICASSP39728.2021.9414428	2020	Modeling	ASRU 2019	zho-eng
10.21437/Interspeech.2020-1980	2020	Modeling	SEAME	zho-eng
10.1109/ISCSLP49672.2021.9362080	2021	Modeling	SEAME	zho-eng
10.21437/Interspeech.2021-978	2021	Modeling	MUCS	hin-eng, ben-eng
10.1109/ASRU51503.2021.9687941	2021	Dataset	German Spoken Wikipedia Corpus	deu-eng
10.3390/app11199106	2021	Modeling	SEAME	zho-eng
10.1109/icassp43922.2022.9747537	2021	Modeling	ASRU 2019	zho-eng
10.1186/s13636-021-00222-7	2021	Other	Proprietary Data	zho-eng
10.21437/Interspeech.2021-1809	2021	Modeling	ESCWA, QASR.CS	ara-eng, ara-fra
10.21437/Interspeech.2021-1578	2021	Modeling	MUCS	hin-eng, ben-eng
10.1016/j.csl.2021.101278	2021	Modeling	ArzEn	ara-eng
10.21437/Interspeech.2021-2231	2021	Modeling	ESCWA	ara-eng
10.1109/ASRU51503.2021.9688174	2021	Modeling	SEAME	zho-eng
10.1587/transinf.2021edp7005	2021	Modeling	Proprietary Data	jpn-eng
10.21437/Interspeech.2021-1339	2021	Other	MUCS	hin-eng, ben-eng
10.3390/APP11062866	2021	Modeling	Proprietary Data	kor-eng
10.1109/icassp43922.2022.9746830	2021	Modeling	SEAME	zho-eng
10.1587/transinf.2022edl8036	2022	Modeling	SEAME	zho-eng
10.21437/Interspeech.2022-221	2022	Modeling	Proprietary Data	zho-eng
10.1109/ICCT56141.2022.10072473	2022	Modeling	Proprietary Data	yue-eng
10.1007/s11042-022-12136-3	2022	Modeling	Proprietary Data	zho-eng
10.21437/Interspeech.2022-877	2022	Dataset	TALCS corpus	zho-eng
10.1109/ICASSP49357.2023.10095878	2022	Modeling	SEAME	zho-eng
10.1109/ISCSLP57327.2022.10037962	2022	Modeling	TALCS corpus and MagicData-RAMC	zho-eng
10.1080/09720529.2021.2014134	2022	Modeling	N.A	mar-eng
10.1109/ISCSLP57327.2022.10037997	2022	Modeling	Proprietary Data	zho-en
10.1007/s10579-022-09592-6	2022	Dataset	Two Sepedi Corpus	nso-eng
10.1109/ISCSLP57327.2022.10038051	2022	Other	TALCS corpus and MagicData-RAMC	zho-eng
10.21437/Interspeech.2022-923	2022	Modeling	ASRU 2019	zho-eng
10.1109/ICASSP49357.2023.10097151	2022	Modeling	SEAME	zho-eng
10.21437/Interspeech.2022-10286	2022	Modeling	ASRU 2021	zho-eng
10.21437/Interspeech.2022-10763	2022	Modeling	MUCS	hin-eng
10.21437/Interspeech.2022-11426	2022	Modeling	Alcubes competition dataset	zho-eng
10.1109/ISCSLP57327.2022.10038194	2022	Modeling	ISCSLP 2022 Code-Switching Challenge	zho-eng
10.21437/Interspeech.2022-719	2022	Modeling	Proprietary Data	zho-eng
10.1109/SLT54892.2023.10023181	2022	Evaluation Metric	HAC	ara-eng
10.21437/Interspeech.2022-421	2022	Dataset	KSC2	kaz-eng
10.21437/Interspeech.2023-2292	2023	Modeling	ASRU 2019	zho-eng
10.21437/Interspeech.2023-2032	2023	Other	ASCEND, SEAME	zho-eng
10.18653/v1/2023.findings-emnlp.543	2023	Data Augmentation	SEAME, ASCEND	zho-eng
10.1109/APSIPAASC58517.2023.10317410	2023	Modeling	ASCEND, NTUT-AB01	zho-eng
10.1109/ASRU57964.2023.10389644	2023	Data Augmentation	SEAME	zho-eng
10.21437/Interspeech.2023-1050	2023	Data Augmentation	SEAME	zho-eng
10.1109/ICASSP48485.2024.10448335	2023	Modeling	ASRU 2019	zho-eng
10.1109/ICEEI59426.2023.10346710	2023	Modeling	Hartanto CS	ind-eng
10.1109/ICASSP49357.2023.10096429	2023	Modeling	Proprietary Data	hin-eng

10.1109/ICASSP49357.2023.10096317	2023	Modeling	■ Proprietary Data	zho-eng
10.3390/app13148492	2023	Dataset	Bilingual Basque-Spanish Dataset	eus-spa
10.1109/ICASSP48485.2024.10446258	2023	Modeling	■ SEAME	zho-eng
10.18653/v1/2023.calcs-1.7	2023	Modeling	■ MUCS	hin-eng, spa-eng
10.21437/interspeech.2023-923	2023	Data Augmentation	■ ASRU 2019	zho-eng
aclanthology.org/2023.calcs-1.8/	2023	Modeling	■ South African Soap Operas	eng-(zul, xho, sot, tsn)
10.1109/LSP.2023.3307350	2023	Modeling	■ SEAME	zho-eng
10.1109/ICASSP48485.2024.10445734	2023	Modeling, Dataset	■ TunSwitch	ara(tunisian)-fra-eng
10.21437/interspeech.2023-262	2023	Modeling	■ SEAME	zho-eng
10.18653/v1/2023.calcs-1.4	2023	Modeling	■ SEAME	zho-eng
10.1109/ASRU57964.2023.10389662	2023	Modeling	■ ASRU 2019	zho-eng
10.1109/SLT54892.2023.10022475	2023	Modeling	■ Proprietary Data	hin-eng
10.21437/interspeech.2023-1465	2023	Modeling	■ ASRU 2019	zho-eng
10.1109/ASRU57964.2023.10389798	2023	Modeling	■ ASRU 2019	zho-eng
10.1109/ICASSP48485.2024.10446857	2023	Data Augmentation	■ SEAME, ■ ESCWA	zho-eng, ara-eng
10.1109/ICEIB57887.2023.10170166	2023	Modeling	■ TALCS corpus	zho-eng
10.1016/j.apacoust.2024.110119	2024	Dataset	■ VITB-HEBiC	hin-eng
10.1109/ICASSP48485.2024.10446652	2024	Modeling	■ SEAME	zho-eng
10.1016/j.csl.2024.101627	2024	Dataset	MECOS	mni-eng
aclanthology.org/2024.lrec-main.174/	2024	Other	■ South African Soap Operas	eng - (zul, xho, sot, tsn, nso)
aclanthology.org/2024.sigul-1.26/	2024	Dataset	■ Mixat	ara(emirati)-eng
aclanthology.org/2024.lrec-main.400/	2024	Dataset	■ DECM	deu-eng
aclanthology.org/2024.lrec-main.852/	2024	Dataset	Kilkarn	que-spa
aclanthology.org/2024.sigul-1.18.pdf	2024	Modeling	■ Hartanto CS	ind-eng
aclanthology.org/2024.lrec-main.1546/	2024	Dataset	■ ZAEBUC-Spoken	ara-eng, ara-ara
10.1016/j.apacoust.2024.109883	2024	Modeling	■ TALCS corpus	zho-eng
10.1145/3640794.3665579	2024	Modeling	■ Bangor Miami	spa-eng
aclanthology.org/2024.lrec-main.262/	2024	Modeling	■ Proprietary Data	zho-eng
10.1109/SLT61566.2024.10832233	2024	Other	■ SEAME	zho-eng
10.1109/SLT61566.2024.10832265	2024	Modeling	■ ASRU 2019	zho-eng
10.1109/ICASSP49660.2025.10889634	2024	Modeling	■ SEAME	zho-eng
10.1109/apsipaasc63619.2025.10849279	2024	Modeling	■ TALCS corpus	zho-eng
10.1109/ACCESS.2024.3496617	2024	Modeling	Roman Urdu code-mixed dataset	urd-eng
10.18653/v1/2024.findings-emnlp.356	2024	Evaluation Metric	■ Mixat	ara(emirati)-eng
10.21437/Interspeech.2024-259	2024	Modeling	■ ASRU 2019	zho-eng
10.18653/v1/2024.emnlp-main.1211	2024	Dataset	■ Casablanca	ara-fra, ara-eng
10.1109/SLT61566.2024.10832290	2024	Modeling	■ ASCEND	zho-eng
aclanthology.org/2024.lrec-main.308/	2024	Data Augmentation	■ Proprietary Data	lav-eng
10.1109/SLT61566.2024.10832173	2024	Other	■ ASRU 2019	zho-eng
10.21437/interspeech.2024-1418	2024	Modeling	■ SEAME	zho-eng
10.1109/COMPAS60761.2024.10796602	2024	Modeling	■ Proprietary Data	ben-eng
10.1016/j.procs.2024.10.184	2024	Modeling	■ ArzEn	ara(egyptian)-eng
10.1109/SLT61566.2024.10832326	2024	Modeling	■ SEAME	zho-eng
10.21437/syndata4genai.2024-6	2024	Dataset	■ Saudi Lang Corpus (SCC) (Eval only)	ara(saudi)-en
10.1109/ICCE62051.2024.10634607	2024	Modeling	■ SEAME	zho-eng
10.1109/ICASSP49660.2025.10890805	2024	Modeling	■ ASRU 2019	zho-eng
10.1109/ICASSP49660.2025.10890024	2024	Modeling	■ ASCEND	zho-eng
10.1109/ICASSPW62465.2024.10627333	2024	Modeling	■ KECS	kor-eng

Table 7: Full list of surveyed papers. **Year** corresponds the earliest available version, which can be a preprint.