

ADI-20: Arabic Dialect Identification dataset and models

Haroun Elleuch^{1,2}, Salima Mdhaffar¹, Yannick Estève¹, Fethi Bougares^{1,2}

¹LIA, Avignon Université, France

²Elyadata, France

haroun.elleuch@elyadata.com, salima.mdhaffar@univ-avignon.fr,
yannick.estève@univ-avignon.fr fethi.bougares@elyadata.com

Abstract

We present ADI-20, an extension of the previously published ADI-17 Arabic Dialect Identification (ADI) dataset. ADI-20 covers all Arabic-speaking countries' dialects. It comprises 3,556 hours from 19 Arabic dialects in addition to Modern Standard Arabic (MSA). We used this dataset to train and evaluate various state-of-the-art ADI systems. We explored fine-tuning pre-trained ECAPA-TDNN-based models, as well as Whisper encoder blocks coupled with an attention pooling layer and a classification dense layer. We investigated the effect of (i) training data size and (ii) the model's number of parameters on identification performance. Our results show a small decrease in F1 score while using only 30% of the original training data. We open-source our collected data and trained models to enable the reproduction of our work, as well as support further research in ADI.

Index Terms: Arabic Dialect Identification, Speech Processing, Low-resource Languages, corpus

1. Introduction

Dialect identification [1] is a crucial task in various natural language processing (NLP) and speech processing applications such as machine translation [2], automatic speech recognition (ASR) [3], sentiment analysis [4], Text-To-Speech [5], text normalization [6], etc. In everyday life, dialect is the most widely used form of communication. Dialects are geographical and social variants of a language, shaped by syntactic, lexical, and phonological differences, including pronunciation shifts (accent). For instance, languages such as Arabic, Chinese, English, and Spanish have numerous dialects that vary significantly across different regions and social groups, illustrating the complexity of dialect identification. Arabic dialects, in particular, are highly diverse, with significant variations in pronunciation, vocabulary, and grammar across regions. Unlike European dialects, which may differ mainly in accent or vocabulary, Arabic dialects can alter fundamental aspects of the language. Additionally, Arabic speakers frequently code-switch between local dialects, Modern Standard Arabic (MSA), and some foreign languages, making the task of dialect discrimination even more challenging. Arabic dialects can be categorized into country-level and city-level variations [7], reflecting both broad regional linguistic patterns and localized speech characteristics. Country-level dialects correspond to the general linguistic traits spoken across an entire nation. However, city-level dialects exist within each country, which exhibit even finer linguistic distinctions influenced by historical, social, and geographical factors.

Recently, the Arabic Dialect Identification (ADI) from speech task has gained significant attention, leading to the development of various datasets and models aimed at addressing its challenges [3, 8, 9, 10]. Several research initiatives have advanced ADI from speech, notably with the organization of multiple shared tasks. The ADI shared task was introduced in VarDial 2016¹ [11]. It consists of identifying a set of four regional Arabic dialects (Egyptian, Gulf, Levantine, North African) and MSA in a transcribed speech corpus. In its second edition, the task was further enriched with the inclusion of acoustic features and the release of audio files, providing richer resources for dialect classification of speech. Driven by the success of the VarDial challenge, ADI-5 [8] and ADI-17 [12, 13] shared tasks have been organized as part of the MGB challenge. ADI-5 contains 74 hours of audio segments from the Al Jazeera TV channel, classified into five dialect groups. ADI-17 consists of 3,033 hours of audio segments from YouTube programs covering a variety of different genres. The segments are split into 17 different dialectal categories, allowing for much finer grain dialectal analysis than the ADI-5 corpus. These efforts provided the research community with two valuable datasets for speech dialect identification, each offering different levels of granularity.

Initial approaches for ADI relied on acoustic features and classical machine learning methods like Support Vector Machines [14] and Deep Neural Networks [15]. More recently, Self-Supervised Learning (SSL) has emerged as a promising approach in the domain of speech processing. SSL models like Hubert [16] and wav2vec 2.0 [17], which leverage large amounts of unlabeled data to learn representations without the need for explicit annotations, have demonstrated considerable potential in improving performance for many speech processing tasks for Arabic [18, 19, 20], including ADI [3].

In this work, we tackle the problem of the robustness of country-level Arabic speech dialect identification, making the following contributions: (1) we expand ADI-17's dialectal coverage by introducing the ADI-20 dataset, (2) we investigate the impact of training data quantity and model complexity on ADI model performance and (3) we release high-performing ADI models for the research community². According to our knowledge, this is the first publicly available ADI model.

2. Datasets

In this work, we use two datasets, namely ADI-17 and ADI-20, both of which are detailed in the following subsections.

¹Workshop on NLP for Similar Languages, Varieties, and Dialects

²Pre-trained models, dataset manifests, and complete recipes for ADI-17 and ADI-20 systems are available at github.com/elyadata/ADI-20

2.1. ADI-17

The official ADI-17 dataset [12] contains 3,033 hours of dialectal Arabic speech for training and around 2 hours per dialect for the test and validation splits. However, the training data is imbalanced, with substantial differences between some dialects. For example, the quantity of training data for the Iraqi dialect is 31 times higher than the quantity available for Jordanian.

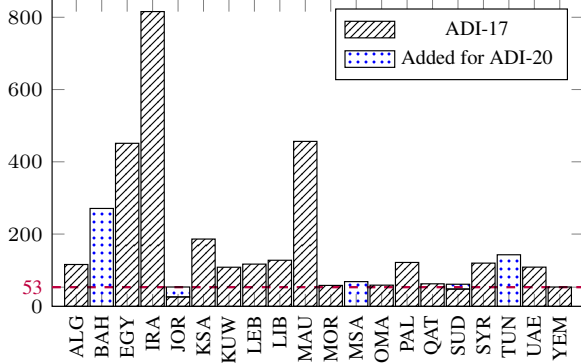


Figure 1: Distribution of ADI-17 and ADI-20 train data (hours) per country dialect. The horizontal line indicates the 53-hour threshold for ADI-17-53h and ADI-20-53h.

As shown in Figure 1, Iraqi, Egyptian, and Mauritanian (Hassaniya) dialects combined account for over half the full ADI-17 training set. Therefore, some dialects, such as Jordanian and Sudanese dialects, are poorly represented in the training set. Even more importantly, we noticed that MSA and some Arabic dialects were not considered during the ADI-17 data collection process, namely Tunisian and Bahraini. For all these reasons, we proposed ADI-20, described in section 2.2.

2.2. ADI-20

The ADI-20 data set is proposed to broaden the dialectal coverage of ADI-17 and include MSA. First, we incorporated MSA to differentiate between regional dialects and the more formal, standardized Arabic. Second, Tunisian and Bahraini dialects were added to the dataset. Obtaining audio data for dialects from countries like Somalia and Eritrea remains challenging and is left for future exploration. Figure 1 illustrates the differences in dialectal coverage between ADI-17 and ADI-20. For MSA, we used audio files from ADI-5 coupled with some scraped broadcast news data from YouTube. Overall, ADI-20 contains 68 hours of MSA speech. For evaluation sets, we have decided to use the same MSA dev and test sets used in ADI-5. Regarding the Tunisian dialect, the whole TunSwitch [21] dataset was used, including the code-switched subset, for a combined duration of 142.7 hours. The train and dev splits were also based on TunSwitch splits but were augmented with content scraped from YouTube to reach 2 hours for each, as is the case for the other dialects. Lastly, Bahraini content is composed exclusively of content sourced from YouTube, with approximately 271 hours for training and 2 hours for validation and testing each.

In addition to the entirely new dialects, we collected more data for the two most underrepresented ADI-17 dialects: Sudanese and Jordanian. Both have been increased to reach the 53-hour training data of the Yemeni dialect. In total, ADI-20 covers 20 Arabic dialects with at least 53 hours of training data

per dialect and a total duration of **3556.76 hours** of MSA and dialectal Arabic speech.

2.3. Data preparation

Since this work aims to assess the impact of training data quantity on ADI model performance, we created three size-based variations of the original ADI-17 dataset described below:

1. **ADI-17-10h**: Limiting training data to 10 hours per dialect. In this configuration, we are using only 5.6% of the full ADI-17 train set.
2. **ADI-17-25h**: 25 hours per dialect. In this setup, only 14% of the full ADI-17 train set will be used for model training.
3. **ADI-17-53h**: 53 hours per dialect. This is equivalent to the usage of about 30% of the total duration of the full training set. Note that we still use only 17 dialects except for Jordanian and Sudanese dialects, which we completed with data from ADI-20 to reach the target size of 53 hours per dialect.
4. **ADI-17-full**: refers to the setup in which the entire ADI-17 dataset is used as it is and without any additional data.

As a result of this data preparation step, we are able to explore different ADI models' performance with respect to the quantity of training data. It is worth mentioning that all three subsets are sampled using the same seed. We would also like to emphasize that the stratification considers the speech duration per dialect irrespective of the number of segments, as those can still differ between dialects. Moreover, all segments exceeding 30 seconds are split into smaller segments of 30 seconds or shorter, and those shorter than 3 seconds are discarded for all train datasets. For the scraped content, we rely on YouTube captions timestamps for audio segmentation.

On top of the 4 subsets described above, we created two additional subsets that integrated the newly introduced dialects into ADI-20. We refer to the first subset as **ADI-20-53h**: Created as an extension of **ADI-17-53h** with an additional 53 hours of speech segments from MSA, Bahraini, and Tunisian dialects. **ADI-20-full** refers simply to the full ADI-20 data set.

3. Models

In this section, we discuss the various trained ADI models. We experimented with two model architectures: (1) ECAPA-TDNN and (2) Whisper-based identification systems. All our models are trained using the open-source SpeechBrain [22] toolkit.

3.1. ECAPA-TDNN

Building on the Time-Delay Neural Network (TDNN) framework, ECAPA-TDNN (Emphasized Channel Attention, Propagation, and Aggregation in TDNN) [23] introduces several enhancements, such as residual connections, squeeze-and-excitation modules, and channel-dependent attentive statistical pooling. The model also leverages Res2Blocks (residual blocks with dilated convolutions) to capture multiscale speech features while integrating position-aware attention to refine feature representations. These improvements enable ECAPA-TDNN to extract highly discriminative speech features, making it a strong candidate for classification tasks like dialect identification. ECAPA-TDNN has achieved state-of-the-art (SoTA) results on the ADI-17 dataset [9], which motivates its choice as a baseline system in this study. All our ECAPA-TDNN systems

are implemented following the CommonLanguage recipe³. We also used the pre-trained model available on HuggingFace⁴ since it has been shown to give better results³. A key difference from the original implementation is that 80-dimensional Mel filterbanks are used instead of MFCC features.

3.2. Whisper encoder

Whisper [24] is a weakly supervised family of encoder-decoder transformer models trained on 680,000 hours of audio crawled from the Internet, using a multitask objective that includes ASR, speech translation, and language identification. These models achieve near SoTA performance across multiple benchmarks and cover 96 languages beyond English, including Arabic. Notably, studies such as [25] have demonstrated Whisper’s ability to deliver competitive results in ASR for MSA and particular Arabic dialects, even outperforming fine-tuned XLS-R models [26] in some cases. Moreover, Whisper is available in multiple size variants (tiny, base, small, medium, large), which makes it a good choice to study model size scalability. Since ADI is fundamentally a classification task, we kept only the Whisper encoder block, and we replaced the decoder with an attention pooling mechanism followed by a classification layer.

4. Experimental setup

We consider three experimental configurations to address our research questions. All our models were trained on 4 GPUs with a batch size of 16.

4.1. ECAPA-TDNN vs Whisper

This setup is designed to compare ECAPA-TDNN and Whisper-based model performances for the ADI task. In order to ensure a fair comparison, we selected the Whisper-base model alongside ECAPA-TDNN. Both models are similar in size: 20.6 and 20.8 million parameters for Whisper-base and ECAPA-TDNN, respectively. Both models were trained using the same **ADI-17-full** without data augmentation.

4.2. Data quantity vs model complexity

The main purpose of this setup is to assess the impact of training data quantity and model size on model performance. For this purpose, we used three pre-trained Whisper models: Whisper small with 88.2M parameters, Whisper medium with 307.2M parameters, and Whisper large with 637 million parameters. All these models are trained on the **ADI-17-full** dataset and its three subsets, described in section 2.3.

4.3. ADI-20 systems

Based on the results from the previously described setup, the third one aims to find the best possible model configuration to obtain the best ADI system. Several ECAPA-TDNN and Whisper models (small, medium, and large) were trained and evaluated on ADI-20 evaluation sets. Therefore, we used data augmentation methods and experimented with different layer-freezing approaches.

5. Results

Table 2 presents the weighted F1-scores achieved by our ADI systems. The first half of the table reports results for models trained on ADI-17 subsets and evaluated on ADI-17 test sets. The second half shows results for ADI-20 systems evaluated on both ADI-17 and ADI-20 test splits. Table 1 reports evaluation results of ECAPA-TDNN systems and the Whisper-base system.

5.1. ECAPA-TDNN vs Whisper

Table 1 presents the dialect identification results on ADI-17 test sets using both ECAPA-TDNN and the Whisper-base models. Obtained results demonstrate that ECAPA-TDNN outperforms the Whisper-based system when they have a similar number of parameters: A weighted F1-score of 93.16% is obtained with ECAPA-TDNN, while the Whisper-based model obtains only 91.34%. We believe this could be explained by the fact that ECAPA-TDNN, unlike the Whisper-base model, was pre-trained on a speaker recognition task, which is somehow similar to dialect identification.

Table 1: *Weighted F1-scores of ADI using Whisper-base and ECAPA-TDNN models trained using ADI-17-full and testing on ADI-17 test sets.*

	Weighted F1-score
ECAPA-TDNN	93.16
Whisper-base	91.34

5.2. Data quantity vs model complexity

We carried out several experiments in order to evaluate the impact of training dataset quantity and model size on dialect identification model performance. As previously mentioned, Whisper-based systems were chosen for their multiple-size option, which is aligned with our goals to investigate model scaling.

First, we experimented with using a bigger Whisper model to replace the Whisper-base used in the previous section. As shown in the results presented in the first column of Table 2, the F1-score improves with increasing the Whisper model size. A better F1-score of 95.66% is achieved compared to the 91.34% obtained with Whisper-base and reported in Table 1. This shows that, as expected, larger models outperform smaller ones for our ADI task.

Next, we focused on training different ADI models with smaller amounts of training data. The results of this experiment are summarized in Table 2. As we can see from the table, using only 10 hours of training data per dialect, only 5.6% of the total ADI-17 training data, significantly decreases the ADI system performance irrespective of the model size. For instance, for the Whisper-large model, we obtained an F1-score of 92.96% with ADI-17-10h while ADI-17-Full gives 95.66%. In the same vein, Table 2 gives the F1-score of models trained using ADI-17-25h and ADI-17-53h. As shown in the table, the gap gets narrower, and we were even able to get a better F1-score for the Whisper-small model using ADI-17-53h of 94.42%, compared to 93.80% when using ADI-17-Full (i.e., with only 30% of the original training data).

Finally, it is worth mentioning that Whisper-medium achieves a similar F1-score compared to Whisper-large (95.29%

³ github.com/speechbrain/speechbrain/tree/develop/recipes/CommonLanguage

⁴ huggingface.co/speechbrain/spkrec-ecapa-voxceleb

Table 2: *Weighted F1-scores of ADI-17 and ADI-20 models on their respective test sets. ADI-20 models are also evaluated on the ADI-17 test split. Bold indicates our best overall system.*

Train dataset	ADI-17 Full	ADI-17 53h	ADI-17 25h	ADI-17 10h	ADI-20 53h		ADI-20-53h + frz. + aug.	
Test dataset	ADI-17				ADI-17	ADI-20	ADI-17	ADI-20
HuBERT-17 [27]	92.12	-	-	-	-	-	-	-
Whisper-small	93.80	94.42	92.86	89.52	93.52	90.88	93.45	91.54
Whisper-med.	95.46	95.29	94.10	92.88	94.96	93.59	95.03	93.39
Whisper-large	95.66	<u>95.16</u>	93.91	92.96	<u>95.41</u>	94.79	95.82	94.83

vs 95.66%) model despite having half as many parameters and only using 30% of the training data.

5.3. ADI-20 systems results

Based on these previous observations, we moved forward to the ADI-20 setup presented in section 4.3. In this context, we trained various ADI systems using the ADI-20-53h training subset. The results are presented in the last column of Table 2. Note that the obtained models are evaluated on both ADI-17 and ADI-20 test sets. This is done for comparability reasons with models trained and evaluated using only ADI-17 (i.e. column 2 in Table 2).

As we can see from the reported results, we obtained slightly better results on the ADI-17 test set by training Whisper-large with ADI-20-53h rather than ADI-17-53h: 95.41% of F1-score for the former versus 95.16% of F1-score for the latter (see underlined scores in Table 2). Note that this improvement only manifests with Whisper-large architecture: Whisper-small and Whisper-medium systems trained on ADI-20-53h don’t outperform their counterparts trained on ADI-17-53h. We also carried out other sets of experiments, in order to push further the ADI models trained using ADI-20-53h data set. Mainly we experimented with data augmentation and encoder lower layer freezing since it has been shown to be helpful for other speech processing tasks [28, 29]. Overall, the best results are obtained using Whisper-large system by freezing the first half of its encoder layers and using data augmentation as implemented in the Speechbrain toolkit [22].

Finally, we also trained the ECAPA-TDNN model using the ADI-20-53h dataset, and the obtained model gives an F1-score of respectively 92.89% and 90.49% on ADI-17 and ADI-20 test sets.

5.4. Zero-shot evaluation

We further evaluate our best systems on the recently released Casablanca dataset [27] in a zero-shot fashion. Results presented in Table 3, show that Whisper large trained with ADI-20-53h+frz+aug achieves an F1-score of 62.74%, significantly outperforming the 39.24% reported in [27] for the system described in [10], which was also evaluated in a zero-shot manner.

According to the results reported in Table 3, the best results are obtained with ADI-20-53h with layer freezing and data augmentation. We hypothesize that augmentations enhance robustness to the unseen conditions of TV dramas, which differ from the domains of our training datasets consisting of YouTube-sourced videos.

Table 3: *Weighted F1-scores of ADI-17 and ADI-20 systems zero-shot evaluation on the Casablanca test set. Bold indicates the best overall system.*

Models	Train data	ADI-17 Full	ADI-20 53h frz+aug
HuBERT-17 [27]		39.24	-
Whisper med.		53.84	58.11
Whisper large		58.89	62.74

We also evaluated our ADI-17-full system, which is trained using the same training data as the HuBERT-17 model reported in [27]. Both Whisper-medium and Whisper-large models clearly outperform HuBERT-17 with an F1-score of 53.84% and 58.89%, respectively. An error analysis of the results reveals that most errors happen when classifying Jordanian, which is often mis-classified as Egyptian or Syrian. Moreover, a considerable amount of misclassification happens between Algerian, Moroccan, and Libyan (which are Maghrebi Arabic dialects).

5.5. Error analysis

Error analysis reveals frequent misclassifications between geographically or linguistically close dialects, e.g., Jordanian as Lebanese/Palestinian/Syrian and Bahraini as Emirati/Qatari. MSA, though not a dialect, is primarily confused with North African varieties (Algerian, Libyan, Egyptian), challenging expectations given Arabic’s origins in the Arabian Peninsula.

6. Conclusion

This paper presents ADI-20, an extension of ADI-17 designed for country-level Arabic dialect identification, offering comprehensive coverage of dialects from almost all Arabic-speaking countries. Using this dataset and the ADI-17 dataset, we investigate the impact of training data quantity and model complexity on ADI model performance, and we demonstrate that competitive ADI models can be trained with as little as 53 hours of data per dialect. We release our data, models, and training recipes to the community to facilitate further research and reproducibility. For future work, we plan to investigate ECAPA-TDNN models further, which demonstrated promising results but were not explored in depth due to time constraints and the availability of Whisper models in various sizes. Additionally, we aim to extend our research to city-level Arabic dialect identification, capturing finer linguistic variations within countries to improve dialect classification granularity.

7. Acknowledgements

This work was partially funded by the ESPERANTO project. The ESPERANTO project has received funding from the European Union's Horizon 2020 (H2020) research and innovation program under the Marie Skłodowska-Curie grant agreement No 101007666. This work was granted access to the HPC resources of IDRIS under the allocations AD011015051, AD011012551R3, and AD011012108R made by GENCI.

8. References

- [1] A. G. Pawar and N. V. Patil, "Comparative study of techniques for spoken language dialect identification," *International Journal of Computer Applications*, vol. 975, p. 8887, 2024.
- [2] W. Salloum, H. Elfardy, L. Alamir-Salloum, N. Habash, and M. Diab, "Sentence level dialect identification for machine translation system selection," in *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, 2014.
- [3] A. Waheed, B. Talafha, P. Sullivan, A. Elmadany, and M. Abdul-Mageed, "VoxArabic: A robust dialect-aware Arabic speech recognition system," in *Proceedings of ArabicNLP 2023*. Association for Computational Linguistics, Dec. 2023, pp. 441–449.
- [4] A. Kaseb and M. Farouk, "SAIDS: A novel approach for sentiment analysis informed of dialect and sarcasm," in *Proceedings of the Seventh Arabic Natural Language Processing Workshop (WANLP)*. Association for Computational Linguistics, Dec. 2022, pp. 22–30.
- [5] K. Doan, A. Waheed, and M. Abdul-Mageed, "Towards zero-shot text-to-speech for Arabic dialects," in *Proceedings of the Second Arabic Natural Language Processing Conference*. Association for Computational Linguistics, Aug. 2024, pp. 123–129.
- [6] B. Alhafni, S. Al-Towaity, Z. Fawzy, F. Nassar, F. Eryani, H. Bouamor, and N. Habash, "Exploiting dialect identification in automatic dialectal text normalization," in *Proceedings of the Second Arabic Natural Language Processing Conference*. Association for Computational Linguistics, Aug. 2024, pp. 42–54.
- [7] C. A. Ferguson, "Diglossia," in *The bilingualism reader*. Routledge, 2003, pp. 71–86.
- [8] A. Ali, S. Vogel, and S. Renals, "Speech recognition challenge in the wild: Arabic MGB-3," in *2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, Dec. 2017, pp. 316–322.
- [9] A. Kulkarni and H. Aldarmaki, "Yet Another Model for Arabic Dialect Identification," in *Proceedings of ArabicNLP 2023*. Association for Computational Linguistics, Dec. 2023, pp. 435–440.
- [10] P. Sullivan, A. Elmadany, and M. Abdul-Mageed, "On the Robustness of Arabic Speech Dialect Identification," in *INTERSPEECH 2023*. ISCA, Aug. 2023, pp. 5326–5330.
- [11] S. Malmasi, M. Zampieri, N. Ljubešić, P. Nakov, A. Ali, and J. Tiedemann, "Discriminating between similar languages and arabic dialect identification: A report on the third dsl shared task," in *Proceedings of the third workshop on NLP for similar languages, varieties and dialects (VarDial3)*, 2016, pp. 1–14.
- [12] S. Shon, A. Ali, Y. Samih, H. Mubarak, and J. Glass, "ADI17: A Fine-Grained Arabic Dialect Identification Dataset," in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 8244–8248.
- [13] A. Ali, S. Shon, Y. Samih, H. Mubarak, A. Abdelali, J. Glass, S. Renals, and K. Choukri, "The MGB-5 Challenge: Recognition and Dialect Identification of Dialectal Arabic Speech," in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, Dec. 2019, pp. 1026–1033.
- [14] S. Wray, "Classification of closely related sub-dialects of arabic using support-vector machines," in *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, 2018.
- [15] L. Lulu and A. Elnagar, "Automatic arabic dialect classification using deep learning models," *Procedia computer science*, vol. 142, pp. 262–269, 2018.
- [16] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhota, R. Salakhutdinov, and A. Mohamed, "Hubert: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE/ACM transactions on audio, speech, and language processing*, vol. 29, pp. 3451–3460, 2021.
- [17] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in neural information processing systems*, vol. 33, pp. 12 449–12 460, 2020.
- [18] A. Djanibekov, H. O. Toyin, R. Alshalan, A. Alitr, and H. Aldarmaki, "Dialectal coverage and generalization in arabic speech recognition," *arXiv preprint arXiv:2411.05872*, 2024.
- [19] S. Mdhaffar, F. Bougares, R. de Mori, S. Zaiem, M. Ravanelli, and Y. Estève, "TARIC-SLU: A Tunisian benchmark dataset for spoken language understanding," in *LREC-COLING 2024*. Torino, Italia: ELRA and ICCL, May 2024, pp. 15 606–15 616. [Online]. Available: <https://aclanthology.org/2024.lrec-main.1357/>
- [20] S. Mdhaffar, H. Elleuch, F. Bougares, and Y. Estève, "Performance analysis of speech encoders for low-resource slt and asr in tunisian dialect," in *Proceedings of The Second Arabic Natural Language Processing Conference*, 2024, pp. 130–139.
- [21] A. A. B. Abdallah, A. Kabboudi, A. Kanoun, and S. Zaiem, "Leveraging Data Collection and Unsupervised Learning for Code-Switched Tunisian Arabic Automatic Speech Recognition," in *ICASSP*, Apr. 2024, pp. 12 607–12 611.
- [22] M. Ravanelli, T. Parcollet, A. Moumen, S. de Langen, C. Subakan, P. Plantinga, Y. Wang, P. Mousavi, L. D. Libera, A. Ploujnikov, F. Paissan, D. Borra, S. Zaiem, Z. Zhao, S. Zhang, G. Karakasidis, S.-L. Yeh, P. Champion, A. Rouhe, R. Braun, F. Mai, J. Zuluaga-Gomez, S. M. Mousavi, A. Nautsch, H. Nguyen, X. Liu, S. Sagar, J. Duret, S. Mdhaffar, G. Laperrière, M. Rouvier, R. D. Mori, and Y. Estève, "Open-source conversational AI with SpeechBrain 1.0," *Journal of Machine Learning Research*, 2024.
- [23] B. Desplanques, J. Thienpondt, and K. Demuyne, "ECAPA-TDNN: Emphasized Channel Attention, Propagation and Aggregation in TDNN Based Speaker Verification," in *Interspeech 2020*. ISCA, Oct. 2020, pp. 3830–3834.
- [24] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. Mcleavey, and I. Sutskever, "Robust Speech Recognition via Large-Scale Weak Supervision," in *Proceedings of the 40th International Conference on Machine Learning*. PMLR, Jul. 2023, pp. 28 492–28 518.
- [25] B. Talafha, A. Waheed, and M. Abdul-Mageed, "N-Shot Benchmarking of Whisper on Diverse Arabic Speech Recognition," in *INTERSPEECH 2023*. ISCA, Aug. 2023, pp. 5092–5096.
- [26] A. Babu, C. Wang, A. Tjandra, K. Lakhota, Q. Xu, N. Goyal, K. Singh, P. Von Platen, Y. Saraf, J. Pino, A. Baevski, A. Conneau, and M. Auli, "XLS-R: Self-supervised Cross-lingual Speech Representation Learning at Scale," in *Interspeech 2022*. ISCA, Sep. 2022, pp. 2278–2282.
- [27] B. Talafha, K. Kadaoui, S. M. Magdy, M. Abdul-Mageed *et al.*, "Casablanca: Data and Models for Multidialectal Arabic Speech Recognition," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Nov. 2024, pp. 21 745–21 758.
- [28] O. Eberhard and T. Zesch, "Effects of layer freezing on transferring a speech recognition system to under-resourced languages," in *KONVENS 2021*, K. Evang, L. Kallmeyer, R. Osswald, J. Waszczuk, and T. Zesch, Eds. KONVENS 2021 Organizers, 6–9 Sep. 2021, pp. 208–212.
- [29] M. Zanon Boito, J. Ortega, H. Riguidel, A. Laurent, L. Barrault, F. Bougares, F. Chaabani, H. Nguyen, F. Barbier, S. Gabbiche, and Y. Estève, "ON-TRAC consortium systems for the IWSLT 2022 dialect and low-resource speech translation tasks," in *IWSLT 2022*, E. Salesky, M. Federico, and M. Costa-jussà, Eds. Association for Computational Linguistics, May 2022, pp. 308–318.