



# Challenges in Automated Processing of Speech from Child Wearables: The Case of Voice Type Classifier

Tarek Kunze<sup>1</sup>, Marianne Métais<sup>1</sup>, Hadrien Titeux<sup>1</sup>, Lucas Elbert<sup>1</sup>, Joseph Coffey<sup>1</sup>, Emmanuel Dupoux<sup>1</sup>, Alejandrina Cristia<sup>1</sup>, Marvin Lavechin<sup>2</sup>

<sup>1</sup>LSCP, DEC, ENS, EHESS, CNRS, PSL University, France

<sup>2</sup>Computational Psycholinguistics Lab., Massachusetts Institute of Technology, United States

tarek.kunze@ens.psl.eu

## Abstract

Recordings gathered with child-worn devices promised to revolutionize both fundamental and applied speech sciences by allowing the effortless capture of children's naturalistic speech environment and language production. This promise hinges on speech technologies that can transform the sheer mounds of data thus collected into usable information. This paper demonstrates several obstacles blocking progress by summarizing three years' worth of experiments aimed at improving one fundamental task: Voice Type Classification. Our experiments suggest that improvements in representation features, architecture, and parameter search contribute to only marginal gains in performance. More progress is made by focusing on data relevance and quantity, which highlights the importance of collecting data with appropriate permissions to allow sharing.

**Index Terms:** wearables, in-the-wild audio, multi-label classification, voice type classification, LENA

## 1. Introduction

Child-worn recording devices have transformed both fundamental and applied speech sciences, capturing egocentric audio (or audio-video) that represents speech from the child's perspective. They enable effortless, long-duration recordings, offering unprecedented insights into children's language environments [1]. Unlike laboratory studies, which may misrepresent everyday speech, or short-form recordings, which are more prone to observer bias, long-form recordings provide more naturalistic data [2]. These recordings are key to reverse-engineering language acquisition [1, 3] and have also been used in applied research, such as evaluating interventions to enhance parental communication and even providing feedback to parents [4, 5].

The value of this technique hinges on the capabilities of advanced speech technologies, which are essential to transform the vast mounds of data collected via child-worn devices into intelligible descriptors usable by researchers and clinicians. Many users of this technology hope for interpretable metrics, such as how much parents talked to and around the child [6], how advanced the child's vocalizations were [7], and, in an ideal world, what children and caregivers were talking about [8, 9, 10]. From the speech technology viewpoint, this results in users' requests for algorithms that accurately identify, categorize, and transcribe verbal interactions assigned to different speakers. Where are we in delivering on this promise?

For the last two decades, the technology predominantly used to fulfill this need has been developed by the LENA Foundation, which has made advances in the identification and categorization of vocalizations by children and adults [11, 12]. Despite these successes, LENA faces significant obstacles to its future viability. Although the LENA software is proprietary

and not openly accessible, available descriptions indicate that it still relies on technology that, by 2025, appears to be outdated. Specifically, the software uses a Gaussian Mixture Model with 36 Mel-frequency cepstral coefficients as input, structured within a minimum duration framework and developed under a maximum likelihood estimation approach to categorize audio segments. These segments are classified into broad speaker and nonspeaker categories, which are often simplified into 4 key categories: Key Child, Other Child, Adult Male, and Adult Female, with everything else being treated as background noise or silence. Despite LENA's technological seniority, it remains the benchmark in the field, with anecdotal evidence suggesting that over 95% of scholarly articles rely on LENA technology.

It was only in 2020 that an open-source alternative emerged, matching or even exceeding LENA's accuracy in the fundamental task of segmenting audio into different speaker types. Capitalizing on the rise of deep learning, Lavechin et al. [13] introduced the Voice Type Classifier (VTC), a model that aimed to democratize access to cutting-edge speech processing technology by releasing as open-source the best-performing model, selected after quite extensive experimentation. The released VTC's architecture combined SincNet, to extract low-level features with a stack of LSTMs (Long Short-Term Memory) to aggregate them into context-dependent representations. The best-performing model was trained to return a score for each of Key Child, Other Child, Adult Male, Adult Female, and Speech (a class that was on when any of the others were on, as well as for live speech that human annotators had not been able to assign to a specific source). Unlike LENA, the VTC is freely accessible, updatable, and unconstrained by proprietary recording hardware, making it advantageous for many different use cases (see [14] for a recent performance comparison).

Despite our excitement about this open-source technology, users and machine learning enthusiasts will be unimpressed by the consideration that VTC outperformed LENA only by a slight margin (F-score of 69% vs. 55% for key child; 33% vs. 29% for other children; 63% vs. 43% for female adult; and 43% vs. 37% for male adult, see [13] for more details). These results teach us a humbling lesson: detecting *who speaks when* remains vexingly difficult, even for modern deep learning approaches.

In this paper, we summarize three years of work attempting to improve VTC performance. Since only one previous study [15] investigated foundation models with English long-form data, we contribute novel data using a more diverse dataset and benchmark against the current open-source state-of-the-art [13]. First, we demonstrate that pre-trained representations learned by Whisper [16] are most effective for smaller training sets. Second, we demonstrate that these pre-trained representations benefit the rare male speech class the most. Third, we provide a detailed analysis of the types of errors made by our model.

Finally, we show that there is still room for improvement by benchmarking against human-human agreement. We conclude by reflecting on the numerous directions we explored to improve performance, too many to detail in a 4-page paper.

## 2. Methods

### 2.1. End-to-end voice type classification

As in [13], we framed the voice type classification problem as a multi-label classification problem, where the input is the audio stream divided into  $N$  frames  $S = \{s_1, s_2, \dots, s_N\}$  and the expected output is the corresponding sequence of labels  $\mathbf{y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N\}$  where each  $\mathbf{y}_i$  is of dimension  $K$  (the number of labels) with  $y_{i,j} = 1$  if the  $j^{\text{th}}$  class is activated,  $y_{i,j} = 0$  otherwise. In training, multiple sequences from the training set were sampled for the total number of frames across all sequences to sum to  $M$  in each mini-batch.

Our baseline model is PyanNet-VTC, PyanNet [13, 17] re-trained from scratch. The PyanNet architecture uses a SincNet [18] feature encoder to learn meaningful filter banks, followed by stacked bidirectional long short-term memory (LSTM) layers and feed-forward (FF) layers. In contrast, our proposed Whisper-VTC model replaces the SincNet encoder with frozen Whisper representations while introducing a different classification architecture. For Whisper-VTC, an 80-channel log-magnitude Mel spectrogram is computed using 25-millisecond windows with a 10-millisecond stride. These spectral features are processed through two convolutional layers before entering the frozen Whisper encoder. The encoder outputs are combined using a learnable weighted sum [19]. The key architectural difference is that these representations are then processed by shared bi-LSTM layers (size 256), whose output features are fed to  $K = 4$  independent binary classification heads, one for each voice type. Each head consists of simple FF layers that output a single scalar for binary classification, unlike PyanNet-VTC’s single multi-label output layer.

Due to Whisper’s fixed 30-second input requirement and the shorter duration of our audio samples, we truncate the encoder’s output representations accordingly. The network is trained to minimize the sum of  $K$  independent binary cross-entropy losses:

$$\mathcal{L} = -\frac{1}{M} \sum_{i=1}^M \sum_{j=1}^K y_{i,j} \log(\hat{y}_{i,j}) + (1 - y_{i,j}) \log(1 - \hat{y}_{i,j})$$

At test time, we use detection thresholds for each task, by default at 0.5, as performed during training. These can be tuned to improve performance or balance the precision-recall trade-off post-training.

For our use case, we consider 4 labels: 1) KCHI for key child vocalizations, 2) OCH for vocalizations produced by any other children, 3) FEM for female adult speech, and 4) MAL for male adult speech. Note that, due to the independent prediction heads, multiple labels can be activated simultaneously in the case of overlapping speech.

The architecture configuration for PyanNet-VTC was the exact same as used in [13]. Whisper-VTC uses bidirectional LSTM layers (size 256) with either 2 or 4 layers showing comparable performance. Table 1 presents results with 4 LSTM layers. The last FF stack uses 1 layer of size 256. The learning rate is set up by an AdamW optimizer with a plateau-based learning rate scheduler (for more details, see <https://github.com/LAAC-LSCP/VTC-IS-25>).

### 2.2. Datasets

We use the same dataset as in [13] that we augmented with three CHILDES [20] datasets, which were not collected with a long-form recording device or a wearable but were close enough in domain as to be found to improve performance during piloting (see Table 1). To increase the uncommon category OCH, we additionally segmented sections of audio from the same long-form datasets in BabyTrain-2020 and two new datasets, always following the DARCLE annotation scheme [21].

### 2.3. Evaluation metrics

We evaluate PyanNet-VTC and Whisper-VTC using the F-score between precision and recall, the identification error rate and percentage correct, as implemented in `pyannote.metrics` [17].

## 3. Results

### 3.1. Experimenting with different Whisper sizes

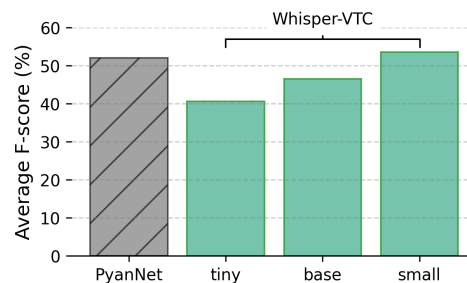


Figure 1: *F*-score (%) averaged across speaker categories for PyanNet (results from [13]) vs. Whisper-VTC (using frozen features from Whisper tiny, base, or small). Performance is computed on the hold-out set.

Figure 1 shows a clear trend: larger Whisper-VTC models achieve better performance, with the small model reaching approximately 50% F-score compared to the tiny model’s 40%. A comparison against PyanNet suggests that using pretrained Whisper representations yields only marginal performance improvement. Additionally, experiments with larger Whisper sizes (not shown here) yield similar performance despite massively longer training times.

### 3.2. Pretrained Whisper representations are most useful for smaller training sets

In this section, we investigate how the performance of Whisper-VTC and PyanNet-VTC varies as a function of training set size. As evident in Figure 2 (showing only Whisper-VTC base to facilitate inspection), Whisper-VTC yields the most significant advantages compared to PyanNet-VTC when training data is limited. While Whisper-VTC shows markedly superior performance in low-resource scenarios (42% F-score vs. 35% with only 10% of training data), this advantage diminishes as more training data becomes available. Beyond 70% data utilization (around 420h of audio), both approaches converge to similar average F-scores of approximately 48%. Interestingly, Whisper-VTC’s performance is much more stable across runs than PyanNet-VTC.

Table 1: Summary statistics of the BabyTrain-2025 dataset. For details about BabyTrain-2020 and the hold-out dataset, see [1]. Most datasets are available via CHILDES [20] or HomeBank [22], which are components of the TalkBank data-sharing platform [23]. BabyTrain-2021 highlighted here is a corrected version. UK: United Kingdom; NA: North America; PNG: Papua New Guinea

Corpus	Access	Language	Tot. Dur.	Cumulated utterance duration			
				KCHI	OCH	MAL	FEM
BabyTrain-2020 [13]	Mixture	Mixture	159h 25m	39h 58m	3h 57m	1h 41m	64h 45m
BabyTrain-2021	Mixture	Mixture	50h 7m	12h 7m	4h 8m	3h 29m	12h 13m
Forrester [24]	CHILDES	English (UK)	11h 47m	4h 28m	0h 31m	4h 19m	2h 29m
Thomas [25]	CHILDES	English (UK)	403h 18m	92h 35m	0h 4m	0h 57m	164h 53m
Soderstrom [26]	CHILDES	English (NA)	25h 59m	2h 30m	1h 00m	0h 19m	12h 30m
Png2019 [27]	Private	Yélf Dnye (PNG)	0h 23m	0h 3m	0h 2m	0h 1m	0h 3m
Solomon	Private	Mixture	5h 23m	0h 36m	1h 3m	0h 32m	1h 8m
Cougar [28]	HomeBank	English (NA)	13h 0m	5h 47m	0h 54m	1h 25m	3h 57m
BabyTrain-2025 (total)	Mixture	Mixture	669h 25m	158h 8m	11h 41m	12h 45m	262h 0m
Hold-out [13]	Mixture	Mixture	20h 0m	1h 39m	0h 45m	0h 43m	2h 48m

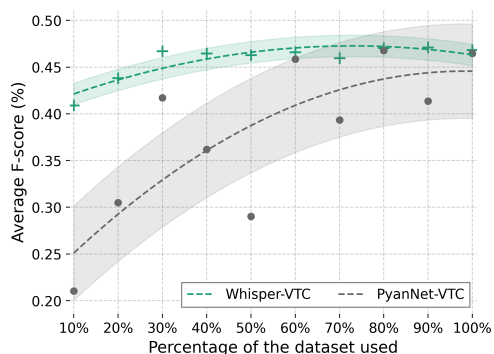


Figure 2:  $F$ -score (%) averaged across speaker categories as a function of training set size for PyanNet-VTC and Whisper-VTC (using frozen features from Whisper base). Performance is computed on the hold-out set.

### 3.3. Segmentation errors: Similarities and differences

Figure 3 (showing only PyanNet-VTC and Whisper-VTC base to facilitate inspection) suggests that, by and large, Whisper-VTC and PyanNet-VTC behave similarly, presenting higher percent correct and lower miss rates when SNR is higher. Thus, it does not seem like Whisper’s extensive pretraining has rendered it more robust to noise than what can be achieved through training with in-domain data. These figures also bring forth a much greater tendency of Whisper-VTC than PyanNet-VTC to miss speech, and for the former to false alarm more than the latter. This difference in behavior is far from obvious, suggesting that perhaps one of the challenges facing systems using Whisper is the ubiquitous presence of far-field speech in long-form recordings. On the other hand, the higher rate of false alarms for PyanNet-VTC is part of the design choice of favoring recall over precision, which could be revisited in the future.

### 3.4. Towards reaching human-level performance

Numerous unsuccessful attempts (see Section 4) at improving the performance led us to question what performance level was realistically achievable on the voice type classification task. Table 2 shows how our automated systems compare to human performance by presenting  $F$ -scores from PyanNet-VTC, Whisper-VTC, and a second human annotator on our hold-out set [21]. Note that these recordings are so challenging that even two well-

Table 2:  $F$ -scores (%) obtained on the hold-out set by PyanNet-VTC [13], Whisper-VTC (using frozen features from Whisper tiny, base or small), and a second human annotator (human 2). Best performances are shown in bold, second best are underlined.

System	KCHI	OCH	MAL	FEM	Ave.
Human 2	<b>79.7</b>	<b>60.4</b>	<b>67.6</b>	<b>71.5</b>	<b>69.8</b>
PyanNet-VTC [13]	68.2	<u>30.5</u>	41.2	63.7	50.9
Whisper-VTC (tiny)	62.6	1.34	39.0	59.5	40.6
Whisper-VTC (base)	63.7	6.7	49.9	66.0	46.6
Whisper-VTC (small)	<u>68.4</u>	20.6	<u>56.7</u>	<u>68.9</u>	<u>53.6</u>

trained human annotators can’t agree perfectly, with  $F$ -scores as low as 60% (OCH) and maximally 80% (KCHI). While PyanNet-VTC and Whisper-VTC achieve similar results in the overall averages, both systems underperform compared to human annotation. Focusing on Whisper-VTC, in the best case scenario, the performance gap is just 2% for FEM, but still massive for OCH. This being the most challenging class is reasonable since (a) children are likely under-represented in Whisper’s original pretraining data and (b) other children can be far-field in long-form recordings.

## 4. Discussion

As alluded to in the Introduction, we summarize here some of the extensive experiments we did with null improvements. In one, we found no improvement when replacing the SincNet architecture (PyanNet-VTC) processing raw waveforms with a stack of convolutional layers processing spectrograms. Later, our colleagues raised valid concerns about class imbalance, as Other Children (OCH) and Male Adult (MAL) speakers represent only 4.7% and 1.2% of the total speech/vocalization time, respectively – a well-studied issue [29, 30]. To address this, we explored multiple approaches: 1) oversampling the least frequent classes, 2) undersampling the most frequent classes, 3) augmenting the training set with MAL speech from CHiME5 [31], and 4) implementing the widely used Grad-Norm technique to balance gradient norms associated to each class [30]. We also explored powerset transformation of target classes, an approach that reframes multi-label classification into powerset multi-class classification [32]. None of these techniques improved performance [33]. Given our incredibly noisy

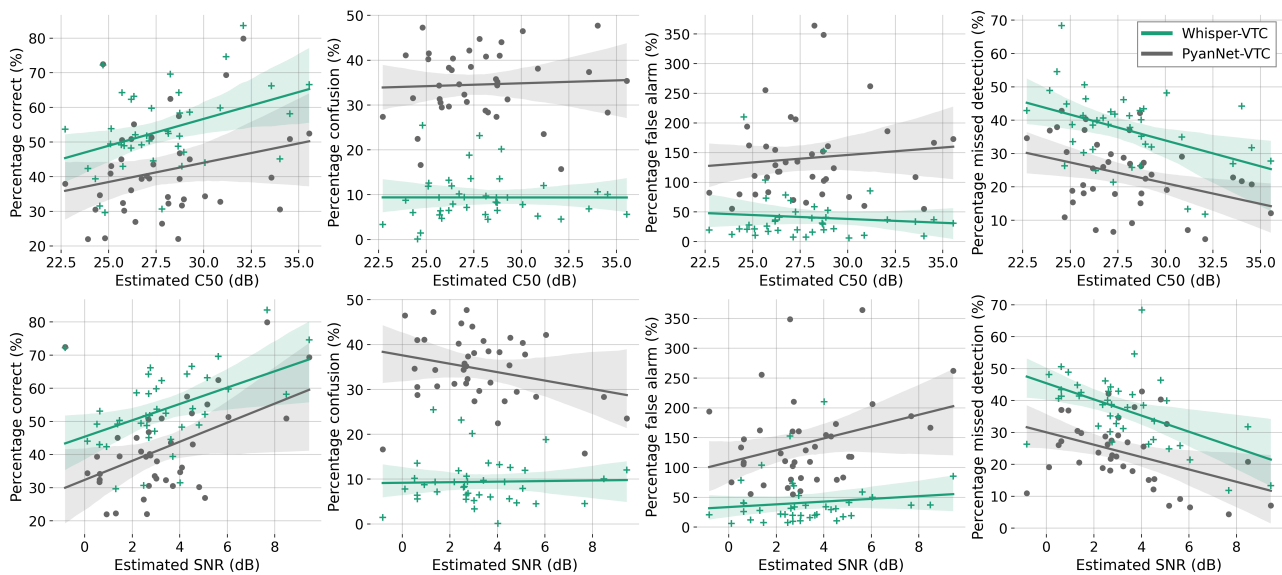


Figure 3: Miss (%), false alarm (%), confusion (%), and correct (%) obtained by PyanNet-VTC and Whisper-VTC (base) as a function of speech-to-noise ratio (SNR) and  $C_{50}$  (estimated by Brouhaha [8]). Each point represents the average performance over all audio from a given child. Solid lines show linear regression fits with shaded 95% confidence intervals.

data, we reasoned that the model could benefit from robust representations learned from massive datasets – similar to other reports [15]. Our exploration began with wav2vec 2.0 pretrained on LibriLight and our own child-centered long-form recordings before moving to the Whisper experiments reported on in this paper.

Our novel contribution was not only a new voice-type classifier with careful benchmarking against the most commonly used open-source alternative (PyanNet-VTC), but also the comparison against inter-human agreement to establish the most meaningful performance ceiling. As shown in Section 3.4, although Whisper-VTC (small) outperforms PyanNet-VTC in three classes, and results are particularly encouraging for some classes (notably MAL), we are still far from human-level performance (particularly for OCH). The small differences in average F-scores, together with the apparent plateau in Whisper-VTC’s performance (Fig. 2), leave us with little hope for significant improvement in the near future.

Despite sacrificing countless interns’ and engineers’ time to the altar of performance optimization, all our sophisticated attempts have fallen short of declaring the task of voice type classification “solved”. These attempts echo Sutton’s “bitter lesson”<sup>1</sup>: incorporating domain knowledge through architectural choices and data balancing techniques do not hold a candle to scaling up the amount of data. While it may be tempting to address our performance gap through increasingly sophisticated models or clever data manipulation, our experience suggests that the path forward likely lies in expanding our annotation campaign to build a larger, more diverse training set. This aligns with the historical pattern in AI, where leveraging computation and data quantity has often proven more effective than hand-engineered solutions based on human intuition about the task.

We believe obtaining high-quality, manually annotated data will be crucial. One approach is increased data sharing across research laboratories. The darcle.org mailing list, compris-

ing over 200 researchers working with daylong recordings, exemplifies existing collaborative infrastructure with established data-sharing practices. When we distributed a request for data contributions several months ago, two researchers offered 60 hours of human-segmented audio each, while five others expressed interest but had incompatible sampling or annotation formats. Despite this positive community response, the available contributions would increase our annotated audio corpus by only 14%, highlighting both the potential and limitations of current data-sharing approaches.

An alternative we hope to explore in future work involves collecting additional data. Messinger et al. [34, 35] have equipped each child and teacher in a classroom with both an audio recording device and RFID-based location tracking (Ubisense Dimension4). This system captures distance and relative angle (i.e., whether individuals are facing each other), offering a promising avenue for automatically generating high-quality labels for key child, female adult, and other child categories. Expanding this approach to home settings could further enable fully automated voice type labeling for all four categories, potentially obviating the need for human annotation.

## 5. Conclusions

We report on over three years of work improving Voice Type Classification in long-form, child-centered recordings—a challenging task due to the realistic recording conditions. Despite these challenges, this technology holds transformative potential for educational interventions and insights into children’s language learning mechanisms. Given that errors in voice classification cascade through subsequent analyses, we echo Soderstrom et al. [21]’s call for continued collaboration and better-annotated datasets to enhance performance.

## 6. Acknowledgements

This work was performed using HPC resources from GENCI- IDRIS (Grant 2024-AD011015450).

<sup>1</sup>The Bitter Lesson of Rich Sutton: <http://incompleteideas.net/IncIdeas/BitterLesson.html>

## 7. References

- [1] M. Lavechin, M. De Seyssel, L. Gautheron, E. Dupoux, and A. Cristia, "Reverse engineering language acquisition with child-centered long-form recordings," *Annual Review of Linguistics*, vol. 8, no. 1, pp. 389–407, 2022.
- [2] E. Bergelson, A. Amatuni, S. Dailey, S. Koorathota, and S. Tor, "Day by day, hour by hour: Naturalistic language input to infants," *Developmental science*, vol. 22, no. 1, p. e12715, 2019.
- [3] E. Dupoux, "Cognitive science in the era of artificial intelligence: A roadmap for reverse-engineering the infant language-learner," *Cognition*, vol. 173, pp. 43–59, 2018.
- [4] A. Weber, A. Fernald, and Y. Diop, "When cultural norms discourage talking to babies: Effectiveness of a parenting program in rural senegal," *Child Development*, vol. 88, no. 5, pp. 1513–1526, 2017.
- [5] J. A. List, J. Pernaudet, and D. L. Suskind, "Shifting parental beliefs about child development to foster parental investments and improve school readiness outcomes," *Nature communications*, vol. 12, no. 1, p. 5765, 2021.
- [6] O. Räsänen, S. Seshadri, M. Lavechin, A. Cristia, and M. Casillas, "Alice: An open-source tool for automatic measurement of phoneme, syllable, and word counts from child-centered daylong recordings," *Behavior Research Methods*, vol. 53, pp. 818–835, 2021.
- [7] M. Cychosz, A. Cristia, E. Bergelson, M. Casillas, G. Baudet, A. S. Warlaumont, C. Scaff, L. Yankowitz, and A. Seidl, "Vocal development in a large-scale crosslinguistic corpus," *Developmental science*, vol. 24, no. 5, p. e13090, 2021.
- [8] M. Lavechin, M. Métais, H. Titeux, A. Boissonnet, J. Copet, M. Rivière, E. Bergelson, A. Cristia, E. Dupoux, and H. Bredin, "Brouhaha: multi-task training for voice activity detection, speech-to-noise ratio, and c50 room acoustics estimation," in *Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2023, pp. 1–7.
- [9] A. Sun, J. J. Londono, B. Elbaum, L. Estrada, R. J. Lazo, L. Vitale, H. G. Villasanti, R. Fusaroli, L. K. Perry, and D. S. Messinger, "Who said what? an automated approach to analyzing speech in preschool classrooms," in *International Conference on Development and Learning (ICDL)*, 2024, pp. 1–8.
- [10] B. Long, V. Xiang, S. Stojanov, R. Z. Sparks, Z. Yin, G. E. Keene, A. W. Tan, S. Y. Feng, C. Zhuang, V. A. Marchman *et al.*, "The babyview dataset: High-resolution egocentric videos of infants' and young children's everyday experiences," *arXiv preprint arXiv:2406.10447*, 2024.
- [11] D. Xu, U. H. Yapanel, S. S. Gray, J. Gilkerson, J. A. Richards, and J. H. Hansen, "Signal processing for young child speech language development," in *WOCCI*, 2008, p. 20.
- [12] J. Gilkerson, J. A. Richards, S. F. Warren, J. K. Montgomery, C. R. Greenwood, D. Kimbrough Oller, J. H. Hansen, and T. D. Paul, "Mapping the early language environment using all-day recordings and automated analysis," *American journal of speech-language pathology*, vol. 26, no. 2, pp. 248–265, 2017.
- [13] M. Lavechin, R. Bousbib, H. Bredin, E. Dupoux, and A. Cristia, "An open-source voice type classifier for child-centered daylong recordings," in *Interspeech*, 2020. [Online]. Available: <https://api.semanticscholar.org/CorpusID:218889736>
- [14] M. Lavechin, L. R. Hamrick, B. Kelleher, and A. Seidl, "Performance and biases of the lena@ and aclew algorithms in analyzing language environments in down, fragile x, angelman syndromes, and populations at elevated likelihood for autism," 2025.
- [15] J. Li, M. Hasegawa-Johnson, and K. Karahalios, "Enhancing child vocalization classification with phonetically-tuned embeddings for assisting autism diagnosis," in *Interspeech*, 2024, pp. 5163–5167.
- [16] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," in *International conference on machine learning*. PMLR, 2023, pp. 28 492–28 518.
- [17] H. Bredin, "pyannote. metrics: A toolkit for reproducible evaluation, diagnostic, and error analysis of speaker diarization systems," in *Interspeech*, 2017, pp. 3587–3591.
- [18] M. Ravanelli and Y. Bengio, "Speaker recognition from raw waveform with sincnet," in *Spoken Language Technology workshop (SLT)*. IEEE, 2018, pp. 1021–1028.
- [19] A. Xu, K. Huang, T. Feng, L. Shen, H. Tager-Flusberg, and S. Narayanan, "Exploring speech foundation models for speaker diarization in child-adult dyadic interactions," in *Interspeech*, 2024, pp. 5193–5197.
- [20] B. MacWhinney, "The chldes system," *Handbook of child language acquisition*, pp. 457–494, 1998.
- [21] M. Soderstrom, M. Casillas, E. Bergelson, C. Rosemberg, F. Alam, A. S. Warlaumont, and J. Bunce, "Developing a cross-cultural annotation system and metacorpus for studying infants' real world language experience," *Collabra: Psychology*, vol. 7, no. 1, p. 23445, 2021.
- [22] M. VanDam, A. S. Warlaumont, E. Bergelson, A. Cristia, M. Soderstrom, P. De Palma, and B. MacWhinney, "Homebank: An online repository of daylong child-centered audio recordings," in *Seminars in speech and language*, vol. 37, no. 02. Thieme Medical Publishers, 2016, pp. 128–142.
- [23] B. MacWhinney, "The talkbank project," in *Creating and digitizing language corpora: Volume 1: Synchronic databases*. Springer, 2007, pp. 163–180.
- [24] M. A. Forrester, "Appropriating cultural conceptions of childhood: Participation in conversation," *Childhood*, vol. 9, no. 3, pp. 255–276, 2002.
- [25] E. Lieven, D. Salomo, and M. Tomasello, "Two-year-old children's production of multiword utterances: A usage-based analysis," 2009.
- [26] M. Soderstrom, M. Blossom, R. Foygel, and J. L. Morgan, "Acoustical cues and grammatical units in speech to two preverbal infants," *Journal of Child Language*, vol. 35, no. 4, pp. 869–902, 2008.
- [27] A. Cristia and M. Casillas, "Lena recordings gathered from children growing up in rossel island. osf," 2023.
- [28] M. VanDam, "Vandam cougar homebank corpus," 2018.
- [29] J. M. Johnson and T. M. Khoshgoftaar, "Survey on deep learning with class imbalance," *Journal of big data*, vol. 6, no. 1, pp. 1–54, 2019.
- [30] Z. Chen, V. Badrinarayanan, C.-Y. Lee, and A. Rabinovich, "Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks," in *International conference on machine learning*. PMLR, 2018, pp. 794–803.
- [31] J. Barker, S. Watanabe, E. Vincent, and J. Trmal, "The fifth 'chime' speech separation and recognition challenge: Dataset, task and baselines," in *Interspeech*, 2018. [Online]. Available: <https://api.semanticscholar.org/CorpusID:4397499>
- [32] A. Plaquet and H. Bredin, "Powerset multi-class cross entropy loss for neural speaker diarization," in *Interspeech*, 2023, pp. 3222–3226.
- [33] L. Elbert, "Explorations with a voice-type classifier for child-centered audio recordings," *Cognitive Machine Learning*, INRIA Paris (M2 report), Tech. Rep., 2020.
- [34] S. G. Mitsven, L. K. Perry, Y. Tao, B. E. Elbaum, N. F. Johnson, and D. S. Messinger, "Objectively measured teacher and preschooler vocalizations: Phonemic diversity is associated with language abilities," *Developmental science*, vol. 25, no. 2, p. e13177, 2022.
- [35] L. K. Perry, S. G. Mitsven, S. Custode, L. Vitale, B. Laursen, C. Song, and D. S. Messinger, "Reciprocal patterns of peer speech in preschoolers with and without hearing loss," *Early childhood research quarterly*, vol. 60, pp. 201–213, 2022.