

# Probabilistic Fusion and Calibration of Neural Speaker Diarization Models

Juan Ignacio Alvarez-Trejos, Sergio A. Balanya, Daniel Ramos, Alicia Lozano-Diez

**Abstract**—End-to-End Neural Diarization (EEND) systems produce frame-level probabilistic speaker activity estimates, yet since evaluation focuses primarily on Diarization Error Rate (DER), the reliability and calibration of these confidence scores have been largely neglected. When fusing multiple diarization systems, DOVER-Lap remains the only established approach, operating at the segment level with hard decisions. We propose working with continuous probability outputs, which enables more sophisticated fusion and calibration techniques that can leverage model uncertainty and complementary strengths across different architectures. This paper presents the first comprehensive framework for calibrating and fusing EEND models at the probability level. We investigate two output formulations—multilabel and powerset representations—and their impact on calibration and fusion effectiveness. Through extensive experiments on the CallHome two-speaker benchmark, we demonstrate that proper calibration provides substantial improvements even for individual models (up to 19% relative DER reduction), in some cases mitigating the absence of domain adaptation. We reveal that joint calibration in powerset space consistently outperforms independent per-speaker calibration, that fusion substantially improves over individual models, and that the Fuse-then-Calibrate ordering generally outperforms both calibrating before fusion and uncalibrated fusion while requiring calibration of only a single combined model. Our best configuration outperforms DOVER-Lap in terms of DER while providing reliable confidence estimates essential for downstream applications. This work proposes best practices for probability-level fusion of EEND systems and demonstrates the advantages of leveraging soft outputs over hard decisions.

**Index Terms**—End-to-End Neural Diarization (EEND), speaker diarization, probabilistic calibration, probability-level fusion.

## I. INTRODUCTION

Speaker diarization identifies and temporally localizes individual speakers in multi-speaker audio recordings. Traditional modular diarization systems rely on assigning speech segments to speakers through a pipeline of speaker embedding extraction and clustering algorithms [1]. While effective, these systems produce hard segment-level assignments where the distance to cluster centroids could provide some measure of uncertainty although it is not usually considered. To handle speaker overlap, external modules are often incorporated into the pipeline. The emergence of End-to-End Neural Diarization (EEND) [2] fundamentally changed this paradigm by directly predicting frame-level speaker activity probabilities for each time frame in an audio stream, enabling integrated modeling of overlapping speech and providing explicit confidence estimates at the frame level.

While recent EEND models achieve remarkably low error rates in speaker activity detection, both the reliability of their confidence estimates and the fusion of multiple diarization

systems remain largely unexplored. Despite the critical role that confidence estimates play in downstream applications and model fusion, research on calibration of neural diarization models is extremely limited. The only prior work [3] assesses calibration to select poorly calibrated samples for retraining, but does not explore post-hoc calibration techniques or model fusion. This gap becomes particularly critical when combining predictions from multiple models, as each model may exhibit distinct confidence patterns and biases that can significantly degrade fusion quality without proper calibration. The rapid evolution of EEND architectures has produced a diverse ecosystem of models with varying strengths: some excel at capturing long-range temporal dependencies [4], [5], others at handling overlapping speech [6], and still others at speaker discrimination in noisy conditions [7].

This diversity naturally raises the question of whether these complementary capabilities can be effectively combined. Existing fusion approaches like DOVER [8] and DOVER-Lap [9] operate on hard decisions at the segment level, limiting their ability to leverage the probabilistic nature of EEND outputs. We show that naive probabilistic model averaging is not able to effectively combine EEND models, and more advanced fusion techniques are required. We argue that this incompatibility arises from model *miscalibration*—that is, their confidence predictions cannot be reliably interpreted as actual class probabilities [10]–[12]. Modern neural networks (NNs), the backbone of EEND models, have shown to produce overconfident predictions [13], [14]: for instance, when a network assigns 0.8 confidence to its predictions, the actual accuracy may be substantially lower than 80%. Given that EEND systems provide frame-level probabilistic outputs, they are uniquely suited for calibration and fusion techniques that require continuous probability estimates—an approach not applicable to methods like DOVER-Lap that operate on hard segment-level decisions.

This paper addresses these fundamental limitations towards a framework for calibrating and fusing neural diarization models. While post-hoc calibration techniques have been successfully applied in speaker recognition tasks [15], [16], their application to speaker diarization presents unique challenges. Our key insight is that the multilabel nature of speaker diarization, where multiple speakers can be simultaneously active in a time-frame, requires rethinking traditional calibration and fusion approaches designed for single-label classification problems. Moreover, well-calibrated probabilities enable optimal decision-making under varying costs and priors through decision theory—an important avenue for future work.

We investigate two output formulations for EEND speaker diarization: the standard multilabel formulation that treats

speakers independently [2], and the *powerset* formulation that explicitly models one class for each of all speaker combinations [17], [18]. Hereafter we will refer to these as Mult (Multilabel) and Power (Powerset) formulations. These output representations lead to different fusion behaviors and calibration requirements that have not been previously explored in neural diarization.

Through extensive experiments on the CallHome two-speaker benchmark [19], we make several key contributions that represent, to our knowledge, the first systematic study of calibration and score-level fusion for EEND systems. First, we demonstrate that proper calibration provides substantial improvements even for individual EEND models, in some cases mitigating the absence of domain adaptation. Unlike prior work that operates on hard decisions, we show that calibrating frame-level probabilistic outputs enables more effective model combination. Second, we explore both unsupervised fusion methods and supervised approaches for combining EEND outputs, analyzing how each interacts with calibration strategies and domain adaptation. Third, we reveal that the choice between multilabel and powerset representations fundamentally impacts both calibration quality and fusion effectiveness, with certain methods benefiting from explicit modeling of speaker combinations while others performing similarly across both formulations. Importantly, we find that improvements in calibration quality (measured by Binary Cross-Entropy [20], [21]) generally translate to corresponding improvements in diarization performance (measured by Diarization Error Rate - DER), demonstrating that better-calibrated confidence estimates support more accurate multi-speaker segmentation. Finally, we propose best practices for score-level fusion of EEND speaker diarization systems, demonstrating that fusing first and then calibrating (Fuse-then-Calibrate) generally outperforms calibrating individual models before fusion, and that joint calibration in powerset space consistently outperforms independent per-speaker calibration.

The remainder of this paper is organized as follows. Section II provides necessary background on EEND models, probability formulations, calibration, and fusion. Section III presents our methodology for calibration and fusion of neural diarization models. Section IV describes the experimental setup. Section V presents comprehensive results analyzing calibration strategies, fusion methods, and their interactions. Section VI concludes with key findings and future work.

**Code Availability:** Our calibration and fusion framework is publicly available at [github.com/SergioAlvarezB/calibrated-fusion-diarization](https://github.com/SergioAlvarezB/calibrated-fusion-diarization), enabling researchers to apply these techniques to their own neural diarization systems.

## II. BACKGROUND AND PROBLEM FORMULATION

### A. End-to-End Neural Diarization

The EEND paradigm reformulates speaker diarization as a multilabel classification task in which NNs assign time frames to zero, one, or multiple speakers simultaneously [2]. Since its introduction, EEND has evolved substantially with developments including self-attention mechanisms [22], attractor-based approaches [23], computational efficiency improvements [24], and various architectural innovations [25].

Research has also explored diverse input feature representations for EEND models. Different feature types—including traditional acoustic features, speaker embeddings, and paralinguistic descriptors—have been shown to provide complementary information for speaker discrimination [26]–[28]. This diversity in feature representations creates natural opportunities for score-level model fusion, as EEND models trained with distinct input features may capture different aspects of speaker characteristics and conversational dynamics.

In this work, we utilize pre-trained EEND with Encoder-Decoder-based Attractors (EEND-EDA) [23] models with different input feature representations as the foundation for our calibration and fusion framework. Following our approach in [26], we leverage multiple EEND-EDA models trained with diverse feature sets to exploit their complementary speaker information.

### B. Output Formulations in Neural Diarization

While EEND established the multilabel classification paradigm as the foundation for neural diarization, recent research has explored alternative formulations that may better capture speaker interaction patterns by performing multi-class classification over the powerset of speakers [18]. The choice of a classification formulation has significant implications for both fusion strategies and calibration techniques.

1) *Multilabel Formulation:* The traditional EEND approach treats each speaker independently at the output level, producing speaker activity probabilities:

$$p_s^{\text{Mult}} \in [0, 1], \quad s \in \{1, \dots, S\} \quad (1)$$

where  $p_s^{\text{Mult}}$  represents the probability of speaker  $s$  being active. Each  $p_s^{\text{Mult}}$  is a probability from a one-vs-all binary classification problem, where  $1 - p_s^{\text{Mult}}$  represents the probability that speaker  $s$  is not active. These probabilities can be collected in a vector  $\mathbf{p}^{\text{Mult}} = (p_1^{\text{Mult}}, \dots, p_S^{\text{Mult}}) \in [0, 1]^S$  for convenience.

While neural architectures such as self-attention mechanisms can capture dependencies between speakers within the model's internal representations, the final output layer treats each speaker prediction independently. Any speaker interaction patterns learned by the network remain implicit in the hidden representations rather than being explicitly modeled in the output formulation. This approach naturally handles overlapping speech but may not fully exploit speaker interaction information.

2) *Powerset Formulation:* An alternative approach considers all possible combinations of active speakers as mutually exclusive classes. For  $S$  speakers, this creates  $K = 2^S$  possible classes:

$$\mathcal{C} = \{\emptyset, \{1\}, \{2\}, \dots, \{S\}, \{1, 2\}, \dots, \{1, 2, \dots, S\}\} \quad (2)$$

The neural network outputs a probability for each class  $c_k \in \mathcal{C}$ :

$$p_{c_k}^{\text{Power}} \in [0, 1], \quad k \in \{1, \dots, K\} \quad (3)$$

where these probabilities form a valid distribution:  $\sum_{k=1}^K p_{c_k}^{\text{Power}} = 1$ . These probabilities can be collected in a vector  $\mathbf{p}^{\text{Power}} = (p_{c_1}^{\text{Power}}, \dots, p_{c_K}^{\text{Power}}) \in \mathbb{S}^K$ , where  $\mathbb{S}^K$

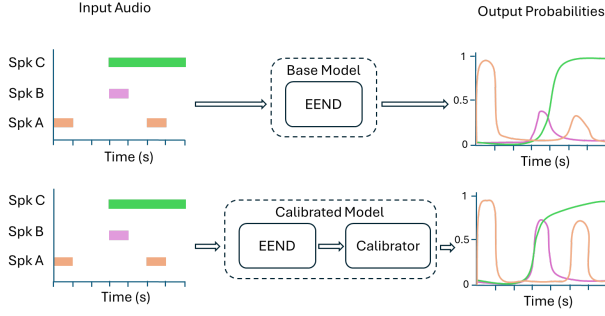


Fig. 1. Post-hoc Calibration Process

denotes the  $(K - 1)$ -dimensional probability simplex, the geometric locus where a  $K$ -dimensional probability vector can lie, i.e.:

$$\mathbb{S}^K = \left\{ \mathbf{x} \in \mathbb{R}^K : \sum_{i=1}^K x_i = 1, x_i \geq 0 \right\} \quad (4)$$

### C. Post-hoc Calibration of Probabilistic Classifiers

We use the term *probabilistic classifier* to denote a classification system that outputs posterior probabilities or confidence scores. For a given input  $\mathbf{x}$ , the system outputs a prediction  $\hat{\mathbf{p}}$  for the posterior class probabilities, where the  $i$ -th component represents the probability, or confidence, assigned to class  $C_i$  such that  $\hat{p}_i \simeq p(C_i|\mathbf{x})$ . Let  $\mathbf{y}$  denote the one-hot encoded target variable, where  $y_i \in \{0, 1\}$  indicates whether  $\mathbf{x}$  belongs to class  $C_i$ .

A probabilistic classifier is said to be well calibrated when its predicted probabilities accurately reflect the likelihood of correct predictions in terms of long-run observed frequencies. Formally, a classifier  $f$  that produces confidence scores  $f(\mathbf{x}) = \hat{\mathbf{p}}$  over  $K$  possible classes, where  $\hat{\mathbf{p}} \in \mathbb{S}^K$ , is perfectly calibrated if:

$$P(\mathbf{y}|\hat{\mathbf{p}} = \mathbf{p}) = \mathbf{p}, \quad \forall \mathbf{p} \in \mathbb{S}^K \quad (5)$$

Modern neural networks are typically not well calibrated out of the box [13], [14]. A commonly used approach to address this issue is *post-hoc calibration*, where the outputs of a miscalibrated classifier are recalibrated. This process involves training a secondary model (the calibration method) on the outputs of the primary classifier, as illustrated in Figure 1.

In this work, we employ Platt scaling [29] and its multi-class extension [30], which essentially applies flavours of logistic regression to the classifier outputs. The specific implementation of this regression varies depending on the calibration strategy, with detailed methods presented in the next section.

The general approach of Platt scaling, also known as matrix scaling, minimizes the cross-entropy loss over a calibration set with respect to learnable parameters. Let  $\mathbf{z} \in \mathbb{R}^K$  denote the pre-softmax output vector, or logit vector, of a classifier over  $K$  classes. The method is implemented as:

$$\tilde{\mathbf{p}} = \text{softmax}(\mathbf{W}\mathbf{z} + \mathbf{b}) \quad (6)$$

where  $\mathbf{W} \in \mathbb{R}^{K \times K}$  and  $\mathbf{b} \in \mathbb{R}^K$  are the learnable parameters, and  $\tilde{\mathbf{p}}$  is the calibrated counterpart of  $\hat{\mathbf{p}} = \text{softmax}(\mathbf{z}) = f(\mathbf{x})$ .

Despite the importance of reliable confidence estimates in speaker diarization systems, research on calibration of neural diarization models remains extremely limited. To date, only one study has specifically addressed this topic, examining the calibration properties of powerset-based diarization models [3]. This lack of research represents a significant gap, particularly considering the multilabel nature of the diarization task and the critical role that confidence estimates play in both downstream applications and model fusion scenarios.

### D. Model Fusion in Neural Diarization

The diversity of neural architectures and feature representations in modern diarization systems creates opportunities for performance improvements through model fusion. Different models may excel in complementary scenarios: acoustic feature-based models might better capture phonetic distinctions, while speaker embedding models could provide superior speaker discrimination capabilities, and models with different temporal receptive fields may capture distinct speaker patterns.

Consider  $M$  diarization models  $f_1, f_2, \dots, f_M$  producing predictions  $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_M$  for the same input, where each  $\mathbf{p}_m$  can be either a multilabel vector  $\mathbf{p}_m \in [0, 1]^S$  or a powerset distribution  $\mathbf{p}_m \in \mathbb{S}^K$ . The ensemble prediction can be formulated as:

$$\mathbf{p}_{\text{fused}} = g(\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_M; \boldsymbol{\theta}) \quad (7)$$

where  $g$  is the fusion function parameterized by  $\boldsymbol{\theta}$ .

Despite the potential benefits, model fusion in neural diarization remains largely underexplored compared to other speech processing tasks. The primary contribution to date is DOVER-Lap [9], which extends the original DOVER framework [8] to handle overlapping segments through modified label mapping and voting mechanisms. DOVER-Lap operates by first aligning speaker labels across different system outputs using a global cost tensor that considers all pairwise overlaps simultaneously, followed by overlap-aware weighted majority voting that can assign multiple speakers to temporal regions.

However, DOVER-Lap still operates on hard decisions rather than probabilistic outputs, discarding valuable uncertainty information that could improve fusion quality.

## III. METHODOLOGY

This section presents our complete framework for fusing and calibrating multilabel neural diarization models. While we focus on EEND-EDA architectures [23], our methodology is applicable to any end-to-end neural diarization system that produces frame-level speaker activity probabilities in multilabel format.

### A. Probability Space Transformation

Given the multilabel and powerset formulations presented above, we can transform between these probability spaces to enable flexible fusion and calibration strategies. Starting from multilabel predictions  $\mathbf{p}^{\text{Mult}} = [p_{s_1}, p_{s_2}]$ , where  $p_{s_i}$  represents the probability of speaker  $s_i$  being active, the transformation to powerset space assumes independence between speakers:

$P(s_1, s_2) = P(s_1)P(s_2)$ . Under this assumption, the probability of each possible speaker activity combination is given by:

$$\mathbf{p}^{Power} = \begin{bmatrix} (1-p_{s_1})(1-p_{s_2}) \\ p_{s_1}(1-p_{s_2}) \\ (1-p_{s_1})p_{s_2} \\ p_{s_1}p_{s_2} \end{bmatrix} \quad (8)$$

While the independence assumption ignores any speaker dependencies learned by the neural network—as these are implicit in the multilabel predictions and cannot be recovered during the transformation—calibration methods applied in the powerset space may partially compensate for this limitation by learning to model speaker interactions through the calibration parameters.

### B. Model Fusion Strategies

We investigate multiple fusion approaches that combine predictions from several diarization models, ranging from simple averaging to sophisticated confidence-weighted schemes. Let  $\mathbf{p}_m$  denote the probability predictions from model  $m \in \{1, 2, \dots, M\}$ , where  $M = 3$  in our experiments, and  $\mathbf{z}_m$  denote the corresponding logits. These predictions can be represented in either multilabel space ( $\mathbf{p}_m \in [0, 1]^S$ ,  $\mathbf{z}_m \in \mathbb{R}^S$ ) or powerset space ( $\mathbf{p}_m \in [0, 1]^K$ ,  $\mathbf{z}_m \in \mathbb{R}^K$ ), depending on the output formulation. These fusion strategies can be broadly categorized into two groups: unsupervised methods that directly combine model outputs without requiring additional training data, and supervised methods that learn combination weights from labeled examples.

1) *Unsupervised Fusion Methods*: We explore four unsupervised fusion strategies to combine model outputs based on different principles:

- **Average Probabilities** represents the most straightforward fusion approach, computing the arithmetic mean of model predictions:

$$\mathbf{p}_{fused} = \frac{1}{M} \sum_{m=1}^M \mathbf{p}_m \quad (9)$$

This method assumes that all models contribute equally to the decision and provides a stable baseline for comparison.

- **Average Logits** performs fusion in the logit space before applying the appropriate activation function (also known as the inverse-link function [31]):

$$\mathbf{p}_{fused} = \phi \left( \frac{1}{M} \sum_{m=1}^M \mathbf{z}_m \right) \quad (10)$$

where  $\phi$  is the activation function: sigmoid  $\phi(\mathbf{z})_s = \sigma(z_s) = \frac{1}{1+\exp(-z_s)}$  for multilabel, or softmax  $\phi(\mathbf{z})_k = \frac{\exp(z_k)}{\sum_{k'=1}^K \exp(z_{k'})}$  for powerset. This approach handles extreme probability values more effectively and preserves the relative confidence differences between models.

- **Dynamic Logits Fusion** introduces adaptive weighting based on prediction confidence:

$$\mathbf{p}_{fused} = \phi \left( \sum_{m=1}^M w_m \cdot \mathbf{z}_m \right) \quad (11)$$

where weights  $w_m$  are computed based on the absolute sum of logits as a confidence measure:

$$w_m = \frac{\sum_i |z_{m,i}|}{\sum_{k=1}^M \sum_i |z_{k,i}|} \quad (12)$$

where  $i$  indexes speakers (multilabel) or classes (powerset). This approach addresses the fact that different models may produce logits at different scales. By normalizing based on logit magnitudes, models with higher confidence (larger absolute logits) receive greater weight in the fusion, allowing the system to dynamically prioritize more certain predictions while accounting for scale differences across models.

- **Entropy Fusion** employs inverse entropy weighting, where models with lower predictive entropy contribute more to the final decision:

$$\mathbf{p}_{fused} = \sum_{m=1}^M w_m \cdot \mathbf{p}_m \quad (13)$$

where the weights are computed as:

$$w_m = \frac{H_{\max} - H(\mathbf{p}_m)}{\sum_{k=1}^M (H_{\max} - H(\mathbf{p}_k))} \quad (14)$$

and  $H(\mathbf{p}_m) = -\sum_i p_{m,i} \log p_{m,i}$  is the predictive entropy, with  $i$  indexing speakers (multilabel) or classes (powerset), and  $H_{\max} = \log S$  (multilabel) or  $H_{\max} = \log K$  (powerset) denoting the maximum possible entropy for uniform distributions.

2) *Supervised Fusion Method*: In contrast to the unsupervised approaches, we investigate a data-driven fusion strategy using a so-called **MetaLearner**—a logistic regression model that learns to combine system outputs. The MetaLearner learns optimal combination weights from labeled data, taking the concatenated logits from all models as input:

$$\mathbf{z}_{fused} = \mathbf{W}[\mathbf{z}_1; \mathbf{z}_2; \dots; \mathbf{z}_M] + \mathbf{b} \quad (15)$$

where  $[\cdot; \cdot]$  denotes concatenation, and  $\mathbf{W} \in \mathbb{R}^{D \times (M \cdot D)}$  and  $\mathbf{b} \in \mathbb{R}^D$  are learned parameters optimized to minimize cross-entropy loss on the training set, with  $D = S$  for multilabel or  $D = K$  for powerset. This supervised approach can capture complex relationships between model outputs that may not be apparent in unsupervised methods.

### C. Calibration Framework

Neural models often produce poorly calibrated predictions [13], where the predicted probabilities do not accurately reflect the empirical proportions of events [10], [11]. This miscalibration can significantly impact fusion performance, as methods that rely on probability magnitudes may be misled by systematic biases in individual models.

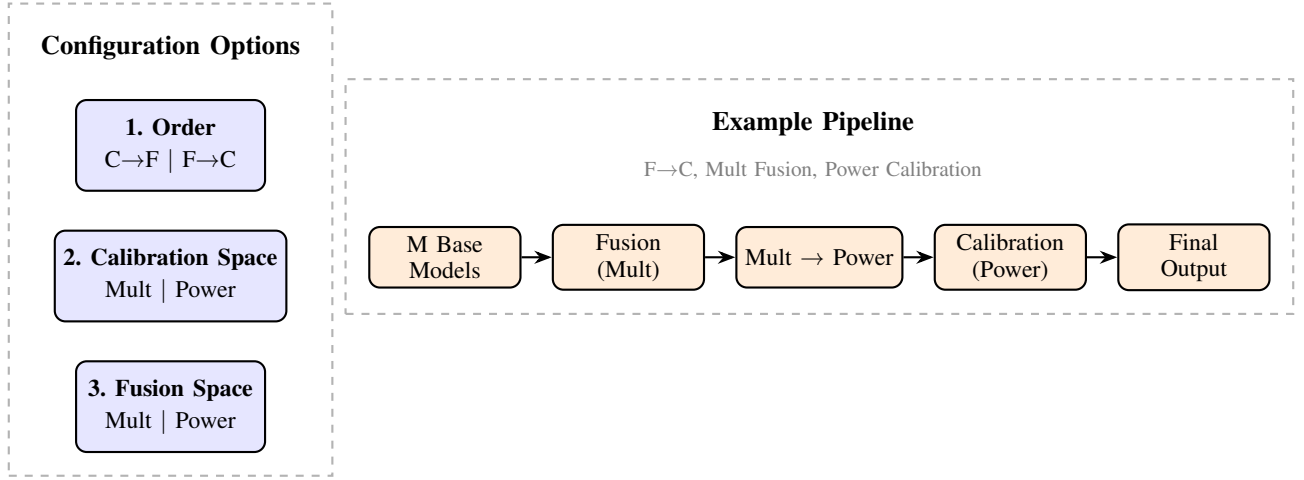


Fig. 2. Modular framework design with three independent configuration decisions (left) and an example horizontal processing pipeline (right). The framework allows independent selection of: (1) calibration-fusion ordering (C→F: Calibrate-then-Fuse or F→C: Fuse-then-Calibrate), (2) calibration probability space (Mult: Multilabel or Power: Powerset), and (3) fusion probability space. The example pipeline demonstrates the F→C strategy with fusion in multilabel space followed by calibration in powerset space.

We implement post-hoc calibration using logistic regression, also known as Platt Scaling [29], which has proven effective in other speech multiclass tasks [30]. The rationale behind logistic regression calibration relies on its objective function, the cross-entropy, which is a proper scoring rule and its minimization improves calibration [30]. For the case when we have independent scores per speaker, i.e. the multilabel domain, we explore two calibration strategies: Independent and Joint Calibration. For the powerset domain, only joint calibration is applicable, as the output represents a probability distribution over mutually exclusive classes.

**Independent Calibration** treats each class separately, learning individual transformation parameters for each speaker:

$$p_i^{cal} = \sigma(\alpha_i \log(p_i) + \beta_i) \quad (16)$$

where  $\alpha_i$  and  $\beta_i$  are class-specific parameters learned independently for each class  $i$ . This approach allows for class-specific correction of calibration biases but ignores potential dependencies between classes.

**Joint Calibration** learns a unified transformation that calibrates all classes simultaneously. For multilabel scenarios, the sigmoid is applied element-wise:

$$\mathbf{p}^{cal} = \sigma(\mathbf{A} \log(\mathbf{p}) + \mathbf{b}) \quad (17)$$

where  $\mathbf{A} \in \mathbb{R}^{K \times K}$  is the learned transformation matrix and  $\mathbf{b} \in \mathbb{R}^K$  is a bias vector. For multiclass (powerset) scenarios:

$$\mathbf{p}^{cal} = \text{softmax}(\mathbf{W} \log(\mathbf{p}) + \mathbf{b}) \quad (18)$$

where the softmax ensures valid probability distributions, i.e.  $\mathbf{p}^{cal} \in \mathbb{S}^K$ .

#### D. Modular Framework Design

Our framework provides a modular approach to combining calibration and fusion for neural diarization models. The design encompasses three key decisions that can be independently configured:

- 1) **Calibration-Fusion Order:** Whether to first calibrate individual base models and then fuse their predictions, or to first fuse the base model predictions and then calibrate the combined output.
- 2) **Calibration Space:** Whether to perform calibration in the multilabel or powerset probability space.
- 3) **Fusion Space:** Whether to perform fusion in the multilabel or powerset probability space.

This modular design allows for flexible exploration of different configurations. Probability space transformations enable calibration and fusion to be performed in different spaces, allowing us to identify the optimal combination of ordering, representation, and methodology for each component.

The framework also enables investigation of whether supervised fusion methods like the MetaLearner inherently perform calibration during training, potentially reducing the need for explicit calibration steps.

## IV. EXPERIMENTAL SETUP

### A. Base Diarization Models and Conditions

We evaluate our framework using three EEND-EDA models previously developed for two-speaker diarization [26]. All models share the same EEND-EDA architecture: a 4-block encoder with 256-dimensional outputs, trained on 50-second segments with a maximum of 15 attractors (though only 2 are actively trained for two-speaker scenarios). Complete implementation details are provided in the cited work; here we highlight the key differences in their input features:

- **MFB** utilizes 23-dimensional Mel-filterbank features with frame stacking. Features are extracted using a 25ms window with 10ms shift, then contextualized by concatenating 7 previous and 7 subsequent frames, resulting in a 345-dimensional input vector ( $23 \times 15$  frames).
- **ECAPA-TDNN  $\oplus$  MFB** enhances the MFB model by concatenating 512-dimensional ECAPA-TDNN speaker embeddings extracted with 1-second windows and 100ms

shift with the 345-dimensional MFB features, creating a 857-dimensional input vector. This configuration leverages specialized speaker-discriminative representations alongside traditional acoustic features.

- **GeMAPS  $\oplus$  MFB** incorporates paralinguistic features by combining a reduced 52-parameter GeMAPS feature set with Mel-filterbank features. The GeMAPS features are extracted using 60ms windows with 10ms shift and contextualized with 2 frames on each side (5 total frames), resulting in a 260-dimensional vector ( $52 \times 5$  frames) that is concatenated with the 345-dimensional MFB features for a final 605-dimensional input.

Each model is evaluated under two training conditions to assess the impact of domain adaptation:

- **No Fine-tuning (No FT)** represents models trained exclusively on simulated conversations, without adaptation to the target domain. This condition allows us to investigate whether calibration can compensate for domain mismatch.
- **Fine-tuned (FT)** includes models that underwent domain adaptation to CallHome using the CH1 two-speaker subset with Adam optimizer (learning rate  $1e-4$ ) following the original EEND-EDA methodology [23].

### B. Dataset and Evaluation Protocol

**Training Data:** All base models are initially trained on simulated conversations (SC) generated following established methodologies [32]. The dataset comprises 2,480 hours of two-speaker conversations created using recordings from multiple corpora: Switchboard-2 (Phases I, II, and III), Switchboard Cellular (Parts 1 and 2) [33], and NIST Speaker Recognition Evaluation datasets (2004, 2005, 2006, and 2008) [34]–[36]. All source material is standardized to 8 kHz sampling rate to match telephone speech conditions. The simulated conversations include background noise from the MUSAN dataset [37] and room impulse responses from the RIR dataset [38] to enhance robustness.

**Evaluation Data:** All experiments are conducted on the CallHome corpus [19], a widely-used benchmark for conversational telephone speech diarization. We focus on the two-speaker subset of the CH2 test set, which contains 148 recordings with approximately 3 hours of total audio.

**Training and Calibration Protocol:** Calibration parameters and MetaLearner weights are estimated using model predictions on the CallHome two-speaker subset of the CH1 adaptation set, which serves as the training data for both calibration methods and supervised fusion. Note that this is the same dataset in which the base models are fine-tuned.

### C. Implementation Details

**Model Processing:** All models process input sequences with a subsampling factor of 10, reducing the temporal resolution from 10ms to 100ms frames for computational efficiency. During inference, predictions are upsampled back to the original temporal resolution for accurate boundary detection.

**Post-processing:** Model outputs undergo median filtering with an 11-frame window (corresponding to 110ms at the

upsampled resolution) to smooth predictions and reduce spurious activations. Final speaker decisions are made using a fixed threshold of 0.5 applied to the filtered probabilities. While calibration enables optimal decision-making in a richer decision-theoretic framework with application-specific costs and utilities, exploring such frameworks is beyond the scope of this work. Here, calibration primarily serves to improve fusion quality by ensuring that probability estimates from different models are comparable and can be meaningfully combined.

**Platt Scaling:** The training procedure is identical across all implementations of Platt Scaling, whether used for calibration or fusion. Parameters are optimized by minimizing cross-entropy loss using the L-BFGS (Limited-memory Broyden-Fletcher-Goldfarb-Shanno) algorithm with default L2 regularization ( $C=1.0$ ) and a maximum of 1000 iterations.

**Fusion Processing:** All fusion experiments are conducted at the 100ms frame level before applying post-processing steps. This ensures that fusion decisions are made on the same temporal granularity across all methods.

### D. Evaluation Metrics

We employ two complementary metrics that assess both diarization performance and prediction quality:

**Diarization Error Rate (DER)** serves as our primary task-specific metric and the standard evaluation measure in speaker diarization. It measures the percentage of time frames with incorrect speaker attribution, including missed speech, false alarms, and speaker confusion errors.

**Cross-Entropy (CE)** is a proper scoring rule that measures the overall quality of probabilistic predictions, considering both calibration and discriminative power. For multilabel classification, CE is computed as Binary Cross-Entropy (BCE) per speaker. For powerset representations, we transform predictions back to multilabel space and compute BCE to maintain a consistent evaluation metric across both formulations. Well-calibrated models do not necessarily achieve lower CE scores, as a model can be perfectly calibrated while having poor discrimination (e.g., by predicting the prior distribution). Therefore, it is important to measure overall prediction quality with a proper scoring rule like cross-entropy [20]. However, post-hoc calibration improves prediction quality by mainly reducing calibration error [11].

## V. FRAMEWORK ANALYSIS AND RESULTS

This section presents a comprehensive analysis of our calibration and fusion framework across multiple dimensions. First, we examine the impact of calibration on the performance of individual base models and their fusion. Second, we compare joint calibration of all speakers versus independent per-speaker calibration. Third, we investigate how the choice of probability space affects both calibration and fusion performance. Fourth, we analyze the effects of different calibration-fusion ordering strategies. Finally, we provide a visual analysis of the impact that fusion and calibration have on each DER component and on the BCE.

Before presenting our results, it is important to note that the base models employed in this work do not represent the



current state-of-the-art in speaker diarization. More advanced architectures such as SortFormer-Hybrid-Loss with 123M parameters achieve 5.87% DER on CallHome 2-speaker [39], and recent EEND variants like AED-EEND-EE with 11.6M parameters reach 6.93% DER [40]. However, the goal of this work is not to achieve state-of-the-art diarization performance, but rather to systematically investigate calibration and fusion strategies that can be applied to any neural diarization system. The performance of our base models provides a clearer view of how calibration and fusion techniques contribute to overall system performance.

### A. Calibration of Base Models

In this subsection, we present a baseline case that demonstrates the effect of calibrating base models before combining their predictions. First, base models are calibrated with Platt Scaling in the powerset domain following the procedure detailed in Section III-C. Then their predictions are combined in the multilabel space using the different fusion strategies described in the previous section. Finally, performance of each combination is measured in terms of DER and BCE.

Table I shows performance metrics for models without and with fine-tuning, before and after calibration. The fused models showcased here were obtained using the simple Average Probs strategy, a common approach for combining predictions.

TABLE I  
DER (%) AND BCE PERFORMANCE OF BASE MODELS AND THEIR FUSION BEFORE AND AFTER CALIBRATION. RESULTS SHOWN FOR MODELS WITHOUT (NO FT) AND WITH FINE-TUNING (FT).

Training Condition	Model	DER (%)		BCE	
		Before	After	Before	After
No FT	MFB	10.381	8.397	0.419	0.271
	ECAPA-TDNN	11.136	10.370	0.551	0.322
	GeMAPS	10.780	9.467	0.416	0.297
	Fused (Avg Probs)	8.425	7.764	0.257	0.252
FT	MFB	8.236	7.834	0.310	0.258
	ECAPA-TDNN	9.054	9.036	0.436	0.290
	GeMAPS	9.068	8.988	0.288	0.263
	Fused (Avg Probs)	7.067	7.031	0.214	0.213

The results show that individual base models benefit from calibration, especially those without fine-tuning. Since the calibration set coincides with the fine-tuning set, calibration provides not only better-calibrated probabilities but also a degree of domain adaptation. For example, MFB improves from 10.381% to 8.397% DER (19.1% relative reduction), while BCE improves from 0.419 to 0.271. After calibration, non-fine-tuned models achieve performance comparable to fine-tuned models: calibrated MFB (8.397% DER) approaches fine-tuned uncalibrated MFB (8.236% DER). This effect becomes more pronounced when fusing predictions: the fusion of calibrated models provides a 7.8% relative improvement in DER over raw fusion for non-fine-tuned models (8.425% to 7.764%), but only 0.5% for fine-tuned ones (7.067% to 7.031%). Nonetheless, calibration consistently improves performance across all model configurations, with BCE improve-

ments ranging from 0.025 (FT GeMAPS) to 0.229 (No FT ECAPA-TDNN).

Figure 3 presents a more comprehensive comparison including all fusion strategies. The observation that non-fine-tuned models benefit more from calibration holds across all fusion methods. The figure also demonstrates that improvements in BCE resulting from calibration are almost always accompanied by corresponding improvements in DER. The only notable exception is the fine-tuned ECAPA model, which exhibits peculiar behavior likely attributable to the use of oracle VAD during training [26]. By setting feature vectors to zero during non-speech segments, oracle VAD alters the distribution of the model's output scores, potentially affecting the relationship between calibration and task performance. However, a detailed analysis of this phenomenon is beyond the scope of this work.

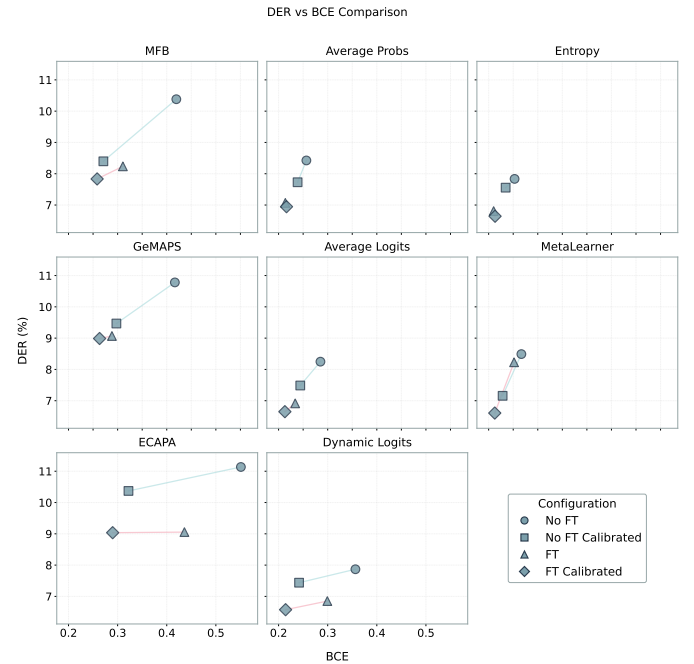


Fig. 3. DER (%) vs BCE comparison on Powerset space for different models and fusion methods.

### B. Joint versus Independent Calibration Strategies

This section compares the two calibration strategies applicable to the multilabel probability space: joint calibration versus independent per-speaker calibration. Following the calibration-then-fusion approach, we first calibrate the non-fine-tuned base models in the multilabel domain, then fuse their predictions and compare performance across both calibration strategies.

Table II shows that joint calibration almost consistently outperforms independent per-speaker calibration across both base models and fusion methods. The improvements are substantial for individual models: MFB improves from 12.059% to 10.874% DER (9.8% relative reduction), ECAPA-TDNN from 17.724% to 12.363% DER (30.2% relative reduction), and GeMAPS from 14.037% to 11.883% DER (15.3% relative reduction). Fusion methods also benefit considerably, with DER improvements ranging from 0.920% (Entropy: 8.453%

TABLE II

DER (%) AND BCE PERFORMANCE OF BASE MODELS AND THEIR FUSION AFTER INDEPENDENT (INDEP) AND JOINT CALIBRATION IN THE MULTILABEL SPACE. RESULTS SHOWN FOR MODELS WITHOUT FINE-TUNING.

Model	DER (%)		BCE	
	Indep	Joint	Indep	Joint
MFB	12.059	10.874	0.351	0.329
ECAPA-TDNN	17.724	12.363	0.514	0.360
GeMAPS	14.037	11.883	0.378	0.354
Average Probs	9.674	8.251	0.315	0.284
Average Logits	8.672	7.698	0.309	0.266
Dynamic Logits	8.313	7.342	0.365	0.263
Entropy	8.453	7.533	0.291	0.264
MetaLearner	7.535	7.534	0.252	0.252
DOVER-Lap*	11.440	10.610	-	-

\*DOVER-Lap operates at segment-level fusion using MFB as primary system.

to 7.533%) to 1.423% (Average Probs: 9.674% to 8.251%). BCE improvements follow similar patterns, with reductions of 0.022 to 0.154 for individual models and 0.027 to 0.102 for fusion methods. The MetaLearner shows minimal change, performing similarly under both calibration strategies (7.535% vs 7.534% DER, 0.252 BCE for both).

This finding is particularly significant because it reveals that speaker dependencies exist even in the multilabel formulation, where speakers are treated independently at the output level. While the multilabel representation assumes independence between speakers, the joint calibration strategy can exploit underlying dependencies in the prediction errors to improve calibration quality. This suggests that calibration should account for inter-speaker relationships, and that joint calibration strategies are preferable when operating in the multilabel space, as they can capture and correct for systematic biases that affect multiple speakers simultaneously. Based on this finding, all subsequent experiments in this work employ joint calibration.

### C. Comparison of Fusion Methods Across Probability Spaces

Table III presents a comparison of fusion methods across both probability spaces and training conditions, including comparison with the existing DOVER-Lap [9] fusion approach.

All fusion methods provide substantial improvements in DER over the best individual model (8.236% DER for fine-tuned MFB). The existing DOVER-Lap method achieves 7.030% DER with fine-tuned models, which serves as a competitive baseline but is outperformed by most of our proposed probability-level fusion approaches.

Comparing performance across probability spaces, the results show no clear systematic advantage of fusing in the powerset domain versus the multilabel domain. While some methods (e.g., Entropy fusion) show slight improvements in powerset space, others (e.g., MetaLearner) perform substantially better in multilabel space. Given this lack of consistent advantage and the fact that the multilabel space has significantly lower dimensionality ( $S$  versus  $K = 2^S$  classes), fusion

TABLE III

PERFORMANCE COMPARISON OF FUSION METHODS IN MULTILABEL (MULT) VERSUS POWERSET (POWER) PROBABILITY SPACES. RESULTS SHOW DER (%) AND BCE FOR BOTH NON-FINE-TUNED (NO FT) AND FINE-TUNED (FT) MODELS.

Training Condition	Method	DER (%)		BCE	
		Mult	Power	Mult	Power
No FT	Average Probs	8.425	8.425	0.257	0.257
	Average Logits	8.248	8.248	0.286	0.285
	Dynamic Logits	7.811	7.867	0.375	0.356
	Entropy	8.428	7.835	0.269	0.253
	MetaLearner	<b>7.326</b>	8.490	<b>0.234</b>	0.267
	DOVER-Lap	9.030	-	-	-
FT	Average Probs	7.067	7.067	0.214	0.214
	Average Logits	6.915	6.915	0.234	0.234
	Dynamic Logits	6.830	6.850	0.314	0.299
	Entropy	6.910	<b>6.806</b>	0.220	<b>0.210</b>
	MetaLearner	6.963	8.227	0.217	0.252
	DOVER-Lap	7.030	-	-	-

in the multilabel domain is preferable for better computational efficiency and comparable performance.

### D. Impact of Calibration Probability Space

This section examines the impact of the probability space on calibration performance. Base models are first calibrated in either the multilabel or powerset space and then fused in the multilabel space. To account for possible interactions with the fusion strategy, we report results for all fusion methods.

TABLE IV

IMPACT OF CALIBRATION SPACE ON NON-FINE-TUNED MODEL PERFORMANCE. DER (%) AND BCE ARE COMPARED AFTER CALIBRATION IN MULTILABEL (MULT) VERSUS POWERSET (POWER) SPACES FOR BOTH INDIVIDUAL MODELS AND FUSION METHODS.

Method	DER (%)		BCE	
	Mult	Power	Mult	Power
MFB	10.874	8.397	0.329	0.271
ECAPA-TDNN	12.363	10.370	0.360	0.322
GeMAPS	11.883	9.467	0.354	0.297
Average Probs	8.251	7.764	0.284	0.252
Average Logits	7.698	7.660	0.266	0.240
Dynamic Logits	7.342	7.691	0.263	0.237
Entropy	7.533	7.758	0.264	0.242
MetaLearner	7.534	7.234	0.252	0.229
DOVER-Lap	10.610	7.940	-	-

For individual models, the results show a striking pattern: calibration in multilabel space consistently degrades performance compared to powerset calibration. Looking at Table IV, individual models calibrated in Mult space show substantially worse DER than when calibrated in powerset space (MFB: 10.874% vs 8.397%, ECAPA-TDNN: 12.363% vs 10.370%, GeMAPS: 11.883% vs 9.467%). This degradation becomes even clearer when examining Table I, which shows performance before and after calibration: for non-fine-tuned models, calibrating in powerset space improves DER (e.g., MFB: 10.381% to 8.397%), but the same models calibrated



in multilabel space (Table IV) actually perform worse than their uncalibrated versions (MFB: 10.874% vs 10.381% uncalibrated). This indicates that multilabel calibration not only fails to improve individual model performance but can actively harm it.

In contrast, fusion methods demonstrate more robustness to the calibration space choice. Interestingly, despite degrading individual model performance, multilabel calibration can still improve fusion results compared to using uncalibrated models. Comparing Table IV with Table III, fusion with multilabel-calibrated models (e.g., Average Probs: 8.251% DER) outperforms fusion without calibration (8.425% DER), suggesting that calibration normalizes prediction scales across models even in suboptimal spaces.

However, the choice of calibration space affects DER and BCE differently for fusion methods. While BCE consistently improves with powerset calibration across all fusion methods without exception (e.g., Average Probs: 0.252 vs 0.284, Dynamic Logits: 0.237 vs 0.263, MetaLearner: 0.229 vs 0.252), the DER results show no clear pattern, with some methods favoring multilabel (Dynamic Logits: 7.342% vs 7.691%, Entropy: 7.533% vs 7.758%) and others favoring Power (Average Probs: 7.764% vs 8.251%, MetaLearner: 7.234% vs 7.534%). This discrepancy reinforces our earlier finding that optimizing calibration quality (BCE) does not necessarily align with optimizing diarization performance (DER).

#### E. Effect of Calibration-Fusion Ordering

This section explores whether calibration can be performed after combining base model predictions. Table V compares the performance of both ordering strategies for each fusion method: Calibrate-then-Fuse (C→F) and Fuse-then-Calibrate (F→C). Results are shown for non-fine-tuned and fine-tuned base models and compared with DOVER-Lap. Note that the F→C approach cannot be applied to DOVER-Lap, as it performs fusion at the hard decision level where frame-level speaker scores are no longer available. All configurations are fused in the multilabel space and calibrated in the powerset space.

The results demonstrate that calibration after fusion (F→C) is not only feasible but often yields superior performance compared to the C→F approach. For fine-tuned models, all fusion methods except Average Probs show improvements with the F→C strategy. This F→C ordering also offers a significant computational advantage: only the single fused model requires calibration, rather than all  $M$  base models in the C→F approach.

Comparing against the DOVER-Lap baseline, several fusion strategies already outperform it with the C→F approach. With the F→C strategy, performance improvements become more consistent—all methods except Average Probs surpass DOVER-Lap’s 6.910% DER for fine-tuned models. The best F→C result (6.543% with Dynamic Logits) represents a substantial improvement over DOVER-Lap while requiring calibration of only a single combined model.

TABLE V  
DER (%) COMPARISON FOR CALIBRATION AND FUSION ORDER ON POWERSSET SPACE. RESULTS SHOWN FOR CALIBRATE-THEN-FUSE (C→F) AND FUSE-THEN-CALIBRATE (F→C) STRATEGIES, WITH AND WITHOUT FINE-TUNING (FT).

Method	Strategy	DER (%)		BCE	
		No FT	FT	No FT	FT
Average Probs	C→F	7.764	7.031	0.252	0.223
	F→C	7.792	6.975	0.242	0.218
Average Logits	C→F	7.660	7.030	0.240	0.210
	F→C	7.512	6.664	0.240	0.213
Dynamic Logits	C→F	7.691	6.910	0.237	0.213
	F→C	7.458	6.543	0.239	0.217
Entropy	C→F	7.758	7.030	0.242	0.217
	F→C	7.498	6.989	0.238	0.218
MetaLearner	C→F	7.234	6.762	0.229	0.213
	F→C	7.202	6.651	0.233	0.214
DOVER-Lap (Baseline)	C→F	7.940	6.910	–	–

#### F. DER Components and Calibration Quality Analysis

This section provides a visual analysis of how calibration, fine-tuning, and fusion affect DER error components (miss, false alarm, and confusion) and BCE at the system level. We compare individual base models against the best fusion method (Dynamic Logits) across four processing stages: no fine-tuning without calibration (No FT No Calib), no fine-tuning with calibration (No FT Calib), fine-tuning without calibration (FT No Calib), and fine-tuning with calibration (FT Calib). All results use joint calibration in the powerset space.

Figure 4 reveals how calibration and fine-tuning affect error distribution across different DER components. The effect of calibration is consistent across all systems: calibration reduces false alarms while increasing miss errors, but the magnitude of false alarm reduction substantially exceeds the miss error increase, resulting in overall DER improvement. For example, MFB shows false alarms decreasing from 7.161% to 2.275% (a reduction of 4.886 percentage points) while miss errors increase from 2.351% to 5.117% (an increase of only 2.766 percentage points), yielding a net improvement in total DER. This asymmetric behavior suggests that calibration relaxes overconfident predictions, shifting decision boundaries toward more conservative thresholds that reduce spurious detections more effectively than they increase missed speech. Fine-tuning further enhances this effect by providing domain-adapted representations that work synergistically with calibration. The Dynamic Logits fusion consistently outperforms individual models across all conditions, demonstrating the benefits of combining calibrated predictions—particularly evident in the substantial reduction of confusion errors, where the fusion effectively resolves speaker identity ambiguities that individual models struggle with.

Figure 5 shows BCE evolution across processing stages, revealing complementary effects of calibration and fusion. Calibration provides the largest BCE improvements for non-fine-tuned models. Notably, the Dynamic Logits fusion consistently achieves lower BCE than individual base models across nearly all processing stages, with the sole exception

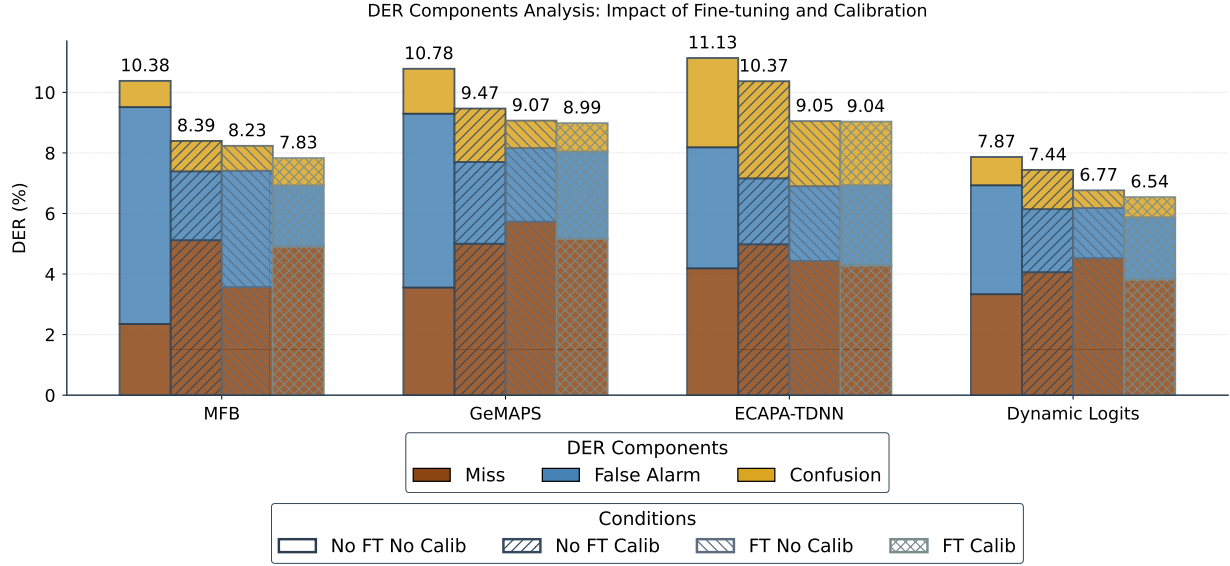


Fig. 4. DER (%) components for individual models and the best fusion method at different processing stages.

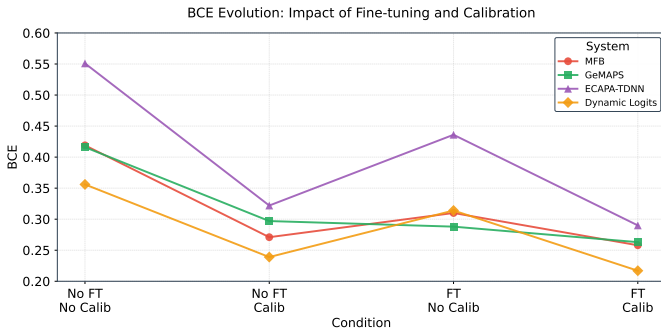


Fig. 5. BCE for individual models and the best fusion method at different processing stages.

of the FT No Calib condition. This pattern demonstrates that combining predictions from diverse models improves prediction quality even without explicit calibration, as the fusion process naturally benefits from averaging complementary and partially uncorrelated model uncertainties. The combination of fine-tuning and calibration yields the best overall results, with Dynamic Logits achieving 0.217 BCE and 6.543% DER, underscoring the synergistic benefits of domain adaptation, model diversity through fusion, and proper calibration.

The relationship between BCE and DER improvements reveals that better calibration typically leads to improved task performance. Across most systems and conditions, reducing BCE through calibration corresponds to DER reductions, demonstrating that well-calibrated confidence estimates support better diarization decisions. For instance, MFB achieves both substantial BCE improvement (0.419 to 0.271) and significant DER reduction (10.381% to 8.397%) through calibration alone. However, some interesting exceptions exist: ECAPA-TDNN shows minimal DER improvement (9.054% to 9.036%) despite substantial BCE reduction (0.436 to 0.290), suggesting that while its predictions become better calibrated,

the improved confidence estimates do not translate to better speaker boundary decisions in this particular case. Similarly, comparing Dynamic Logits configurations shows that the non-fine-tuned calibrated version achieves better BCE (0.239) than the fine-tuned uncalibrated version (0.314), yet the latter yields superior DER (6.77% vs 7.44%), highlighting that domain adaptation provides complementary benefits beyond calibration. These observations indicate that while calibration quality and task performance are generally well-aligned, optimal diarization performance requires both well-calibrated predictions and domain-adapted representations. Dynamic Logits with both fine-tuning and calibration demonstrates this synergy, achieving the best overall results.

## VI. CONCLUSIONS

This work presents the first comprehensive framework for calibrating and fusing EEND systems at the probability level. Through experiments on CallHome, we demonstrate that proper calibration provides substantial improvements even for individual models (up to 19% relative DER reduction), in some cases mitigating the absence of domain adaptation.

Our key findings proposes critical best practices for neural diarization: (1) Powerset representations with joint calibration consistently outperform independent per-speaker calibration (8.26% vs 9.18% DER), highlighting the importance of modeling speaker dependencies explicitly; (2) The Fuse-then-Calibrate strategy achieves superior performance (6.543% DER, 5.2% relative improvement over DOVER-Lap) while requiring calibration of only a single combined model; (3) Dynamic Logits fusion demonstrates the best performance across experimental conditions, effectively combining models at the logit level before softmax transformation.

We observe that calibration quality and task performance are not always aligned—improved BCE does not necessarily translate to improved DER. This finding highlights the need

for task-specific calibration objectives in practical deployments. While we use a fixed threshold of 0.5 for speaker decisions in this work, the application of decision theory to optimize thresholds based on application-specific costs and well-calibrated probabilities remains an important direction for future work. Furthermore, calibration and fine-tuning prove complementary: calibration alone provides significant gains when domain data is unavailable, while their combination yields optimal results.

Future work should extend this framework to scenarios with more speakers, explore alternative proper scoring rules that may better correlate with DER optimization, and investigate applications in downstream tasks such as speaker-attributed automatic speech recognition. By demonstrating that probability-level techniques outperform hard-decision approaches while providing reliable confidence estimates, this work establishes the foundation for more effective neural diarization systems that fully exploit their probabilistic outputs.

#### ACKNOWLEDGMENTS

This research was supported by project PID2021-125943OB-I00 funded by MCIN/AEI/10.13039/501100011033 FEDER, UE and project PID2024-160789OB-I00 funded by MICIU/AEI/10.13039/501100011033 FEDER, UE and project SI4/PJI/2024-00237 (COSER-IA), Comunidad de Madrid.

#### REFERENCES

- [1] X. Anguera, S. Bozonnet, N. Evans, C. Fredouille, G. Friedland, and O. Vinyals, "Speaker diarization: A review of recent research," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 20, no. 2, pp. 356–370, 2012.
- [2] Y. Fujita, N. Kanda, S. Horiguchi, K. Nagamatsu, and S. Watanabe, "End-to-end neural speaker diarization with permutation-free objectives," *arXiv preprint arXiv:1909.05952*, 2019.
- [3] A. Plaquet and H. Bredin, "On the calibration of powerset speaker diarization models," in *Interspeech 2024*, 2024, pp. 3764–3768.
- [4] K.-W. Huang and C.-P. Chen, "Long audio file speaker diarization with feasible end-to-end models," in *2024 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, 2024, pp. 1–6.
- [5] K. Kinoshita, M. Delcroix, and N. Tawara, "Integrating end-to-end neural and clustering-based diarization: Getting the best of both worlds," in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 7198–7202.
- [6] J. M. Coria, H. Bredin, S. Ghannay, and S. Rosset, "Overlap-aware low-latency online speaker diarization based on end-to-end local segmentation," in *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2021, pp. 1139–1146.
- [7] S. Maiti, Y. Ueda, S. Watanabe, C. Zhang, M. Yu, S.-X. Zhang, and Y. Xu, "EEND-SS: Joint end-to-end neural speaker diarization and speech separation for flexible number of speakers," in *2022 IEEE Spoken Language Technology Workshop (SLT)*, 2023, pp. 480–487.
- [8] A. Stolcke and T. Yoshioka, "DOVER: A method for combining diarization outputs," in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2019, pp. 757–763.
- [9] D. Raj, L. Paola Garcia-Perera, Z. Huang, S. Watanabe, D. Povey, A. Stolcke, and S. Khudanpur, "DOVER-lap: A method for combining overlap-aware diarization outputs," in *2021 IEEE Spoken Language Technology Workshop (SLT)*, 2021, pp. 881–888.
- [10] A. P. Dawid, "The well-calibrated Bayesian," *Journal of the American Statistical Association*, vol. 77, no. 379, pp. 605–610, 1982.
- [11] J. Bröcker, "Reliability, sufficiency, and the decomposition of proper scores," *Quarterly Journal of the Royal Meteorological Society*, vol. 135, no. 643, pp. 1512–1519, 2009.
- [12] M. H. DeGroot and S. E. Fienberg, "The comparison and evaluation of forecasters," *Journal of the Royal Statistical Society: Series D (The Statistician)*, vol. 32, no. 1–2, pp. 12–22, 1983.
- [13] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, "On calibration of modern neural networks," in *International conference on machine learning*. PMLR, 2017, pp. 1321–1330.
- [14] Y. Ovidia, E. Fertig, J. Ren, Z. Nado, D. Sculley, S. Nowozin, J. V. Dillon, B. Lakshminarayanan, and J. Snoek, "Can you trust your model's uncertainty? Evaluating predictive uncertainty under dataset shift," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [15] D. A. van Leeuwen and N. Brümmer, "An introduction to application-independent evaluation of speaker recognition systems," in *Speaker Classification I*, ser. Lecture Notes in Computer Science, C. Müller, Ed. Berlin, Heidelberg: Springer, 2007, vol. 4343. [Online]. Available: [https://doi.org/10.1007/978-3-540-74200-5\\_19](https://doi.org/10.1007/978-3-540-74200-5_19)
- [16] N. Brummer, L. Burget, J. Cernocky, O. Glembek, F. Grezl, M. Karafiat, D. A. van Leeuwen, P. Matejka, P. Schwarz, and A. Strasheim, "Fusion of heterogeneous speaker recognition systems in the stbu submission for the nist speaker recognition evaluation 2006," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 7, pp. 2072–2084, 2007.
- [17] Z. Du, S. Zhang, S. Zheng, and Z.-J. Yan, "Speaker overlap-aware neural diarization for multi-party meeting analysis," in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2022, pp. 7458–7469.
- [18] A. Plaquet and H. Bredin, "Powerset multi-class cross entropy loss for neural speaker diarization," in *Proc. INTERSPEECH 2023*, 2023.
- [19] A. Martin and M. Przybicki, "The nist 1999 speaker recognition evaluation—an overview," *Digital Signal Processing*, vol. 10, no. 1, pp. 1–18, 2000.
- [20] L. Ferrer and D. Ramos, "Evaluating posterior probabilities: Decision theory, proper scoring rules, and calibration," *arXiv preprint arXiv:2408.02841*, 2024.
- [21] D. Ramos, J. Franco-Pedroso, A. Lozano-Diez, and J. Gonzalez-Rodriguez, "Deconstructing cross-entropy for probabilistic binary classifiers," *Entropy*, vol. 20, no. 3, 2018. [Online]. Available: <https://www.mdpi.com/1099-4300/20/3/208>
- [22] Y. Fujita, N. Kanda, S. Horiguchi, Y. Xue, K. Nagamatsu, and S. Watanabe, "End-to-end neural speaker diarization with self-attention," in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2019, pp. 296–303.
- [23] S. Horiguchi, Y. Fujita, S. Watanabe, Y. Xue, and K. Nagamatsu, "End-to-End Speaker Diarization for an Unknown Number of Speakers with Encoder-Decoder Based Attractors," in *Proc. Interspeech 2020*, 2020, pp. 269–273.
- [24] F. Landini, M. Diez, T. Stafylakis, and L. Burget, "DiaPer: End-to-end neural diarization with perceiver-based attractors," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, pp. 1–16, 2024.
- [25] Z. Chen, B. Han, S. Wang, and Y. Qian, "Attention-based encoder-decoder network for end-to-end neural speaker diarization with target speaker attractor," in *Proc. Interspeech 2023*, 2023, pp. 3552–3556.
- [26] J. I. Alvarez-Trejos, A. Lozano-Diez, and D. Ramos, "Feature integration strategies for neural speaker diarization in conversational telephone speech," *Applied Sciences*, vol. 15, no. 9, 2025.
- [27] J. Han, F. Landini, J. Rohdin, A. Silnova, M. Diez, and L. Burget, "Leveraging self-supervised learning for speaker diarization," in *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2025, pp. 1–5.
- [28] A. Khare, E. Han, Y. Yang, and A. Stolcke, "ASR-aware end-to-end neural diarization," in *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 8092–8096.
- [29] J. Platt, "Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods," *Advances in large margin classifiers*, vol. 10, no. 3, pp. 61–74, 1999.
- [30] N. Brummer and D. A. Van Leeuwen, "On calibration of language recognition scores," in *2006 IEEE Odyssey - The Speaker and Language Recognition Workshop*, 2006, pp. 1–8.
- [31] P. McCullagh, "Generalized linear models," *European Journal of Operational Research*, vol. 16, no. 3, pp. 285–292, 1984.
- [32] F. Landini, A. Lozano-Diez, M. Diez, and L. Burget, "From Simulated Mixtures to Simulated Conversations as Training Data for End-to-End Neural Diarization," in *Proc. Interspeech 2022*, 2022, pp. 5095–5099.
- [33] J. Godfrey, E. Holliman, and J. McDaniel, "Switchboard: telephone speech corpus for research and development," in *ICASSP-92: 1992 IEEE*

- International Conference on Acoustics, Speech, and Signal Processing*, vol. 1, 1992, pp. 517–520 vol.1.
- [34] O. Sadjadi, C. Greenberg, E. Singer, L. Mason, and D. Reynolds, “NIST 2021 speaker recognition evaluation plan,” 2021.
  - [35] M. A. Przybocki, A. F. Martin, and A. N. Le, “NIST speaker recognition evaluation chronicles - part 2,” in *2006 IEEE Odyssey - The Speaker and Language Recognition Workshop*, 2006, pp. 1–6.
  - [36] A. F. Martin and C. S. Greenberg, “NIST 2008 speaker recognition evaluation: Performance across telephone and room microphone channels,” in *Tenth Annual Conference of the International Speech Communication Association*, 2009.
  - [37] D. Snyder, G. Chen, and D. Povey, “Musn: A music, speech, and noise corpus,” *arXiv preprint arXiv:1510.08484*, 2015.
  - [38] T. Ko, V. Peddinti, D. Povey, M. L. Seltzer, and S. Khudanpur, “A study on data augmentation of reverberant speech for robust speech recognition,” in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 5220–5224.
  - [39] T. Park, I. Medennikov, K. Dhawan, W. Wang, H. Huang, N. R. Koluguri, K. C. Puvvada, J. Balam, and B. Ginsburg, “Sortformer: A novel approach for permutation-resolved speaker supervision in speech-to-text systems,” in *Forty-second International Conference on Machine Learning*, 2025. [Online]. Available: <https://openreview.net/forum?id=AyYjRvrbDx>
  - [40] Z. Chen, B. Han, S. Wang, and Y. Qian, “Attention-based encoder-decoder end-to-end neural diarization with embedding enhancer,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 32, pp. 1636–1649, 2024.