

From Child Adaptation to Adult Retention: Merging Specialized Arabic–English ASR Models Across Architectures

Anonymous ACL submission

Abstract

Automatic speech recognition (ASR) systems exhibit persistent performance disparities across age groups and speaker nativity, with children’s speech remaining a systematically underrepresented and challenging domain, even in high-resource languages. Existing adaptation strategies predominantly rely on fine-tuning, which often induces catastrophic forgetting and degrades performance on adult speech; these limitations are further amplified in bilingual children’s ASR, where robust cross-language generalization is required. In this work, we explore weight-space model merging as a principled framework for age-robust and language-inclusive speech modeling. Starting from a shared multilingual base model, we fine-tune complementary child-adapted checkpoints and merge them using balanced weighting to preserve adult representations while incorporating age- and language-specific adaptations. Across all benchmarks, model merging consistently improves recognition accuracy for children while retaining or in some cases improving—adult performance, outperforming fine-tuning and joint training baselines. These results demonstrate that model merging provides a scalable and data-efficient alternative to fine-tuning for inclusive ASR across age groups, speaker nativity (L1 and L2 Arabic speakers), and languages (Arabic and English).

1 Introduction

Automatic Speech Recognition (ASR) systems have achieved strong performance on adult speech across multiple languages, largely due to the availability of large-scale training data and robust model architectures. In contrast, children’s speech still remains challenging. Compared to adult speech, children’s speech exhibits distinct acoustic and linguistic characteristics, including higher pitch, shorter vocal tracts, increased pronunciation variability, and frequent disfluencies, all of which contribute

to degraded recognition performance (Wu et al., 2019; Bhardwaj et al., 2022). These challenges are further amplified by the limited availability of annotated children’s speech data.

To overcome, the widely used approach involves fine-tuning large, pre-trained adult models on available children’s data. While this often yields substantial gains on child-specific benchmarks like MyST (Pradhan et al., 2024; Linguistic Data Consortium, 2021) and CMU Kids (Jain et al., 2023), it introduces a critical representational trade-off. Standard fine-tuning often results in *catastrophic forgetting*, where the model’s adaptation to child-specific acoustic shifts causes a drift in the underlying representations that support adult speech recognition. This limitation is particularly severe in bilingual or multilingual settings, where the model must simultaneously encode age-related acoustic variance and distinct language-specific structures (e.g., Arabic and English).

Previous efforts to mitigate this forgetting have largely focused on continual learning frameworks, utilizing regularization techniques like Elastic Weight Consolidation (EWC) to preserve knowledge from the original domain (Sadhu and Hermansky, 2020; Ahadzi et al., 2025). While effective, these methods require modifying the training objective and maintaining complex protocols during the adaptation phase. In contrast, a parallel line of research in vision and NLP offers a more flexible paradigm: *weight-space model merging*. Techniques such as model soups and task arithmetic demonstrate that averaging or interpolating the weights of independently fine-tuned models can improve generalization and combine task-specific behaviors without additional training or inference costs (Wortsman et al., 2022; Ilharco et al., 2023).

In this study, we shift the focus from simple adaptation to a more fundamental question: “how can child-specific acoustic shifts be integrated into a model without much compromising the in-

084 *tegrity of its previously learned adult representa-*
085 *tions?”* Hence, rather than just finetuning with
086 child-specific ASR data, we propose to decou-
087 ple adaptation from retention using weight-space
088 model merging.

089 Starting from a bi/multilingual base ASR, we de-
090 rive multiple child-adapted checkpoints through tar-
091 geted fine-tuning to capture complementary induc-
092 tive biases including language-invariant adult pho-
093 netics, language-specific child acoustic patterns,
094 and cross-lingual age-related characteristics. We
095 then merge these checkpoints directly in weight
096 space using various techniques, Linear Interpo-
097 lation (LERP) , and TIES-merging (Yadav et al.,
098 2023). This process produces a unified model that
099 incorporates the nuances of children’s speech while
100 preserving high-fidelity adult performance across
101 both Arabic and English.

102 We demonstrate the broad utility and
103 architecture-agnostic nature of this approach
104 by experimenting with three distinct architectural
105 frameworks: (i) Multilingual Encoder-Decoder
106 using Whisper (Radford et al., 2023); (ii) Bilingual
107 Encoder-Decoder for Arabic-English ASR; and
108 (iii) Bilingual Encoder-CTC for Arabic-English
109 ASR; as the foundation model. Our evaluation
110 is uniquely comprehensive, spanning both native
111 and non-native speech across two languages.
112 We utilize the established English MyST corpus
113 (Pradhan et al., 2024) and introduce a new dataset
114 of Arabic native (**AraKidsL1**) and non-native
115 (**AraKidsL2**) children’s speech. Crucially, we
116 measure knowledge retention by benchmarking
117 against original adult domains, including MGB-2
118 (Arabic) and LibriSpeech (English). By evaluating
119 across these diverse architectures and datasets,
120 we show that weight-space merging can acts as
121 an implicit regularizer, aligning age-dependent
122 acoustic patterns while safeguarding the linguistic
123 structures required for robust, inclusive ASR.

124 Our main contribution includes the following:

- 125 1. We presented a post-hoc, training-free framework
126 to balance “adaptation-retention” trade-off in
127 children’s ASR using weight-space model merg-
128 ing for knowledge integration.
- 129 2. We validate the effectiveness of such framework
130 across three distinct ASR architectures – Mul-
131 tilingual Encoder-Decoder (Whisper), Bilingual
132 Encoder-Decoder, and Bilingual Encoder-CTC,
133 showcasing the architecture-agnostic robustness.

3. We introduce and release **AraKids** – a new
dataset comprising both native and non-native
Arabic children’s speech. This addresses a sig-
nificant gap in the community and serves as a
challenging benchmark for cross-nativity ASR.

4. We evaluate and compare multiple existing merg-
ing strategies for their ability to align age-related
acoustic shifts in a bilingual (Arabic-English)
context.

2 Related Work

2.1 Children’s ASR

Children’s ASR remains difficult because child
speech differs from adult speech acoustically and
linguistically (e.g., higher pitch, shorter vocal
tracts, greater pronunciation variability, and more
disfluencies). A widely used approach is to adapt
adult-trained acoustic or end-to-end ASR models
with limited child data, and prior work shows that
fine-tuning large pretrained/foundation models can
yield substantial gains, especially with multilin-
gual pretraining (Wu et al., 2019; Bhardwaj et al.,
2022). Recent studies further demonstrate strong
improvements when adapting Whisper-style or self-
supervised models on common benchmarks such
as MyST and CMU Kids (Jain et al., 2023; Lin-
guistic Data Consortium, 2021). However, many
adaptation studies primarily report target-domain
gains and provide limited analysis of whether adult-
domain performance is retained after fine-tuning.

2.2 Forgetting, Continual Learning, and Training-Free Alternatives

Sequential domain adaptation can trigger *cata-*
trophic forgetting, motivating continual/lifelong
learning methods for ASR via regularization, re-
play, or constrained optimization (Sadhu and Her-
mansky, 2020; Houston and Kirchhoff, 2020). For
children’s ASR, online adaptation studies show
that regularization-based approaches (e.g., EWC,
Synaptic Intelligence) can better preserve earlier
capabilities than naive fine-tuning (Ahadzi et al.,
2025), but they add training-time complexity. In
parallel, weight-space *model merging* offers a
training-free route to combine fine-tuned behaviors:
model soups and task-arithmetic-style operations
can improve robustness when checkpoints share
an initialization and remain compatible (Wortsmann
et al., 2022; Ilharco et al., 2023; Ortiz-Jimenez
et al., 2023; Sung et al., 2023). Although less ex-
plored in speech than in vision/NLP, recent ASR

183 results indicate that merging independently fine-
 184 tuned checkpoints can improve robustness without
 185 extra inference cost (Ducorroy and Riad, 2025),
 186 making it a practical option when multiple adapta-
 187 tion runs are already available.

188 2.3 Positioning and Contributions

189 We position our work as a “training-free” paradigm
 190 adaption that shifts the burden of knowledge re-
 191 tention from the training phase to a post-hoc
 192 weight-space operation. While standard fine-tuning
 193 achieves high accuracy on children’s speech at the
 194 expense of adult performance, and regularization
 195 methods like EWC mitigate this by introducing
 196 training-time complexity, we show model merging
 197 method achieves a similar balance with negligi-
 198 ble computational overhead. Furthermore, we ex-
 199 pand the scope of current research by investigating
 200 merging across distinct architectures; specifically
 201 Encoder-Decoder and CTC-based ASR model and
 202 applying these techniques to a bilingual Arabic-
 203 English context. By introducing a new Arabic
 204 native and non-native children’s corpus, we are
 205 also able to evaluate the limits of cross-lingual
 206 and cross-age generalization in a way that remains
 207 largely unexplored in existing ASR literature.

208 3 Methodology

209 We aim to develop an inclusive ASR model that
 210 excels in the highly variable domain of Arabic and
 211 English children’s speech while maintaining the
 212 comparable performance benchmarks established
 213 for adult speakers. We define this objective as an
 214 optimization problem in the shared latent param-
 215 eter space such that the single model θ can accu-
 216 rately transcribe speech across distinct dimensions
 217 of variability: (i) **Speaker age** (child vs. adult);
 218 (ii) **Nativity** (native vs. non-native); and (iii) **Bilin-**
 219 **gualism** (Arabic vs English).

220 **Base and Naive Fine-tuned Models:** We define
 221 the base model (θ_{base}) as the original pre-trained
 222 ASR model, which is specialized in adult speech.
 223 Following, we designed monolingual and bilingual
 224 child-speech specialized models, θ_{child} , by fine-
 225 tuning the base with children training data \check{D}_{child}^L ,
 226 where $L \in \{\text{Ar, En, Ar+En}\}$.

$$227 \theta_{base} \xrightarrow{\text{FT}^L} \theta_{child}^L$$

228 The shift in knowledge during this adaptation can
 229 also be encoded by the task vector, $\tau^L = \theta_{child}^L -$

θ_{base} , which isolates the parameter-level updates
 230 necessary for child-specific acoustics. 231

232 **Catastrophic Forgetting:** We define the catas-
 233 trophic forgetting gap, $\Delta_{forgot}(\theta_*, \theta_{base})$, as the
 234 model performance degradation on the original
 235 source domain (Adult ASR) when evaluating a
 236 model adapted to the target domain (child ASR).

237 **Model Merging Objective:** The primary ob-
 238 jective is to derive a generalized (inclusive) θ_m
 239 model, that achieve performance parity with spe-
 240 cialized child model ($\theta_m \approx \theta_{child}^L$) while minimiz-
 241 ing $\Delta_{forgot}(\theta_m, \theta_{base})$. We achieve this by combin-
 242 ing the specialized models using a functional
 243 merging operator, ($Merge(\cdot)$).

$$244 \theta_m = \theta_{base} + Merge(\theta_1, \dots, \theta_N)$$

245 where N represents the number of specialized mod-
 246 els (e.g., Arabic-child, English-child, Bilingual-
 247 child) being integrated into the base model.

248 **Linear Interpolation (LERP):** Linear Interpo-
 249 lation (LERP) operator, $Merge_L$, allows a fine-
 250 grained balance between the base model’s stability
 251 and the specialists’ expertise. Unlike uniform merg-
 252 ing, LERP introduces a global scaling factor, λ , to
 253 control the overall intensity of the adaptation and
 254 individual mixing weights (α_i) to prioritize specific
 255 tasks.

$$256 Merge_L(\theta_1, \dots, \theta_N) = \lambda \sum_{i=1}^N \alpha_i \theta_i.$$

257 where $\alpha \in [0, 1]$ and $\sum \alpha = 1$; for $\lambda \in [0, 1]$
 258 and $\lambda \rightarrow 0$ values makes the θ_M bias to the θ_{base} ,
 259 decreasing Δ_{forgot} , whereas, $\lambda \rightarrow 1$ maximizes
 260 the child-speech adaptation.

261 **TIES-Merging:** To reduce the interference of
 262 different task vectors, τ^i to move same weight in
 263 conflicting direction, we opt for TIES, which in-
 264 clude three stages – Trim, Elect Sign and Disjoint
 265 Merge. Trim stage first retain the top $k\%$ weight
 266 values with highest magnitude, setting rest to zero
 267 to reduce the parameter noise. The Elect sign stage
 268 calculate a consensus sign vector γ , based on the di-
 269 rection that possesses the highest cumulative mag-
 270 nitude across all N task vectors. Following, only
 271 the selected parameters across task vectors that
 272 aligns with γ are averaged, filtering out conflicting
 273 updates.

4 AraKids for L1 and L2 Speakers

Overview: We present **AraKids**, an Arabic children’s speech corpus designed to benchmark ASR robustness across *nativity* (L1 vs. L2), *regional grouping*, and *age*. AraKids contains 64,705 utterances from 181 speakers (109.18 hours). It comprises two subsets: **AraKidsL1**, native Arabic children grouped by Arabic macro-regions (Egyptian, Gulf, Levantine, North African), and **AraKidsL2**, non-native Arabic children grouped by speakers’ broad origin regions (Africa, Asia, Americas, Europe) (Table 2). We plan to release AraKids publicly in future work. To the best of our knowledge, AraKids is among the first large-scale Arabic children’s speech dataset introduced specifically to support ASR research with explicit L1/L2 coverage and broad regional diversity, and we plan to release AraKids publicly in future work.

Prompts and text sources: Prompts are read and curated to be age-appropriate. For **AraKidsL1**, sentences are selected from school-grade materials to match speakers’ expected literacy level. For **AraKidsL2**, Arabic teachers adapt and calibrate prompts to span multiple proficiency levels while controlling lexical and syntactic complexity. All prompts are fully diacritized.

Speaker ages: AraKids train covers ages 7–18, split into 7–12 (71 speakers, 38.80 hours) and 13–18 (110 speakers, 70.38 hours). For a focused evaluation on younger children, the test set is restricted to ages 7–12.

Splits and leakage control: We enforce **speaker-disjoint** train/test partitioning and **text-disjointness** at the prompt level (no overlap in recorded prompt sentences between train and test), reducing memorization effects and measuring generalization to unseen speakers and unseen prompt texts. In addition, we report prompt inventory statistics: across **train** (L1+L2), AraKids contains 8 unique `ScriptIDs` (book source) and 3,194 unique `UtteranceIDs` (sentence), while **test** (L1+L2) contains 5 unique `ScriptIDs` and 380 unique `UtteranceIDs`.

Benchmark test set: We construct a region-balanced test set of 2,000 utterances by sampling 250 utterances per region across the 8 regions (4 AraKidsL1 macro-regions + 4 AraKidsL2 origin regions). This yields 1,000 utterances for AraKidsL1 (26 speakers, 94 minutes) and 1,000 utterances for

Split	#Spk	#Utt	Hours
Train (AraKidsL1)	85	29,222	42.27
Train (AraKidsL2)	96	35,483	66.91
Train (Total)	181	64,705	109.18
Test (AraKidsL1)	26	1,000	1.57
Test (AraKidsL2)	40	1,000	2.43
Test (Total)	66	2,000	3.99

Table 1: AraKids split overview. The test set is region-balanced by sampling 250 utterances per region across the 8 regions and is restricted to the 7–12 age group.

Region	#Spk	#Utt	Hours
<i>AraKidsL1</i>			
Egyptian	29	9,478	14.49
Gulf	18	6,796	10.11
Levantine	22	6,873	8.75
North African	16	6,075	8.92
<i>AraKidsL2</i>			
Africa	14	5,339	11.03
Americas	16	6,000	12.44
Asia	32	11,796	22.01
Europe	34	12,348	21.43
Total	181	64,705	109.18

Table 2: AraKids train distribution by region. AraKidsL1 regions are Arabic macro-regions; AraKidsL2 regions are speakers’ origin regions.

AraKidsL2 (40 speakers, 2.43 hours). Overall, the test set totals 3.99 hours, with average utterance durations of 5.64s (L1) and 8.74s (L2) (Table 1).

5 Experimental Setup

5.1 ASR Model Architectures

We evaluate weight-space merging across three representative ASR architectures to verify that our approach is not tied to a specific modeling or decoding paradigm. Throughout the paper, we denote the encoder–decoder family as **ED**, the encoder–CTC family as **ECTC**, and Whisper Large v3 as **WL**.

5.2 ASR Model Architectures

Bilingual Encoder-CTC (E-CTC): Our encoder–CTC backbone is an encoder-only bilingual (Arabic–English) ASR model adapted from the OWSM-CTC architecture. It consists of (i) a speech encoder built on top of a frozen Arabic-centric self-supervised front-end whose embeddings are processed by a stack of L transformer encoder layers, and (ii) a transformer text-history encoder whose representations are injected into selected speech-encoder layers via cross-attention to condition recognition on previous context.

347 Training uses self-conditioned CTC by applying
348 auxiliary CTC losses at intermediate layers in
349 addition to the final-layer CTC loss, averaged
350 to improve convergence. For long-form audio,
351 we employ parallel chunk-wise greedy decoding
352 with overlapping 30-second windows, followed by
353 overlap alignment and confidence-based stitching
354 to produce the final transcript. We tokenize
355 transcripts using a 20,002 BPE vocabulary, and
356 apply SpecAugment during training.

357 **Bilingual Encoder–Decoder (ED):** Our bilin-
358 gual encoder–decoder (ED) baseline is an ESPnet
359 Conformer–Transformer model trained on Arabic–
360 English speech only. It uses an XLSR-53 self-
361 supervised front-end (via S3PRL) with multi-layer
362 feature extraction, followed by a 12-block Con-
363 former encoder and a 6-block Transformer decoder,
364 jointly optimized with a hybrid CTC/attention ob-
365 jective ($\lambda_{\text{CTC}} = 0.3$). We tokenize transcripts
366 using a 10k BPE vocabulary, and apply SpecAug-
367 ment during training.

368 **Multilingual Encoder-Decoder (WL):** For the
369 multilingual encoder–decoder setting, we use
370 **Whisper Large v3 (WL)** (Radford et al., 2023), a
371 large-scale ASR model pretrained on diverse mul-
372 tilingual and multitask speech data. WL serves
373 as a strong architecture-independent baseline for
374 our merging experiments due to its robust multilin-
375 gual representations and widely adopted encoder–
376 decoder formulation.

377 5.3 Datasets and Benchmarks

378 We evaluate adaptation to children’s speech and
379 retention on adult speech using a consistent bench-
380 mark suite spanning both domains.

381 **Child Specialized Model Training Data:** We
382 build child-specialized models by fine-tuning on
383 AraKids (Arabic children; Table 1), and MyST
384 (Pradhan et al., 2024) with 171 hours of training
385 audio. For all child-training datasets and all model
386 architectures, we apply SpecAugment (Park et al.,
387 2019) and speed perturbation (Ko et al., 2015) to
388 improve robustness.

389 **Children’s Benchmarks:** We measure in-
390 domain performance on two children-focused test
391 sets: AraKids Test (Table 1) and MyST clean test
392 set (Pradhan et al., 2024).

393 **Adult Retention Benchmarks:** To quantify re-
394 tention of adult ASR capability, we evaluate on

395 LibriSpeech test-clean (Panayotov et al., 2015) and
396 the MGB-2 Arabic multi-dialect broadcast bench-
397 mark (Ali et al., 2016). These adult benchmarks
398 serve as out-of-domain anchors, enabling us to mea-
399 sure how well merging preserves general-purpose
400 performance while improving children’s ASR.

401 5.4 Naive Fine-tuning Parameters

402 As a reference adaptation strategy, we perform
403 *naive fine-tuning* of each baseline model on the chil-
404 dren training data using standard supervised ASR
405 objectives. To limit overfitting and to keep adapta-
406 tion budgets comparable across architectures, we
407 fine-tune the ESPnet models (ED and ECTC) for **5**
408 **epochs** each, while WL is fine-tuned for **1 epoch**.

409 For ECTC, we optimize the E-CTC model with
410 self-conditioned intermediate CTC supervision (in-
411 termediate layers at {6, 9, 11}) and a text-history
412 prompt encoder. We use AdamW with a peak
413 learning rate of 2.5×10^{-4} and a piecewise linear
414 warmup schedule, batch size 32 with gradient ac-
415 cumulation of 2, mixed-precision training (AMP),
416 and SpecAugment; the self-supervised HArNESS
417 front-end is kept frozen.

418 For ED, we fine-tune the ED under a hybrid
419 CTC/attention objective ($\lambda_{\text{CTC}} = 0.3$). Optimiza-
420 tion uses Adam with learning rate 2×10^{-4} and
421 WarmupLR (15k warmup steps), SpecAugment,
422 and gradient accumulation of 16 with numel-based
423 batching. The XLSR-53 front-end is frozen during
424 training.

425 For WL, we fine-tune the WL using the Hugging-
426 Face Seq2SeqTrainer setup with learning rate
427 10^{-5} , warmup of 500 steps, per-device batch size
428 32 with gradient accumulation of 2, FP16 training,
429 and checkpoint selection based on validation WER.

430 5.5 Model Merging Parameters

431 We instantiate the merging operator $Merge(\cdot)$ us-
432 ing two weight-space methods: LERP and TIES
433 (Yadav et al., 2023). In all cases, we merge *task-*
434 *specialized* checkpoints into a single inclusive
435 model while controlling the retention–adaptation
436 trade-off through mixing weights (LERP) or sparse,
437 sign-consensus task-vector aggregation (TIES). Un-
438 less stated otherwise, merges follow the check-
439 point precision of each model family (ESPnet in
440 float32, Whisper in float16). For each merging
441 technique, we manually swept the main hyperpa-
442 rameters (e.g., interpolation weights, density, and
443 update weights) and selected the reported config-
444 uration based on the best average WER computed

445 across *all* evaluation test sets.

446 **LERP.** For LERP, we use convex weights $\{\alpha_i\}$
447 over the participating checkpoints (and set the
448 global scaling $\lambda = 1$ as per our implementation),
449 yielding a weighted average in parameter space.
450 For **ED**, we merge the encoder and decoder *sep-*
451 *arately* using the same weights: $\alpha_{base} = 0.6$ and
452 $\alpha_{Ar} = \alpha_{Ar+En} = 0.2$. For **WL**, we merge the en-
453 coder with $\alpha_{base} = 0.6$ and $\alpha_{Ar+En} = 0.4$. For
454 **ECTC**, we merge the encoder with $\alpha_{base} = 0.6$,
455 $\alpha_{Ar} = 0.2$, and $\alpha_{Ar+En} = 0.2$.

456 **TIES.** For TIES, we use the baseline checkpoint
457 as θ_{base} and merge one or more specialists by ap-
458 plying (i) trimming via a density parameter (re-
459 taining the top-magnitude fraction of each task
460 update), (ii) electing a consensus sign, and (iii)
461 disjointly averaging only sign-aligned updates (Ya-
462 dav et al., 2023). We enable `normalize=true` and
463 `int8_mask=true` in all TIES runs. For **ED**, we ap-
464 ply TIES on the encoder (baseline as `base_model`)
465 with per-expert settings `density` $\in \{0.8, 0.6\}$ and
466 `weight` $\in \{0.6, 0.4\}$. For **WL**, we apply TIES
467 between the baseline and the 1-epoch fine-tuned
468 checkpoint with `density=0.7` and `weight=0.7`
469 (baseline as `base_model`). For **ECTC**, we apply
470 TIES using the baseline encoder as `base_model`
471 and two specialists (θ_{child}^{Ar+En} and $A\theta_{child}^{Ar}$), with
472 `density` $\in \{0.8, 0.6\}$ and `weight` $\in \{0.6, 0.4\}$.

473 **ED merging scope (what is merged vs. kept**
474 **from baseline).** For the **ED** architecture, merg-
475 ing is applied *only* to the main encoder and de-
476 coder tensors. We explicitly do *not* merge aux-
477 iliary components and statistics: the frozen self-
478 supervised front-end is kept from the baseline,
479 the **CTC head** is inherited from the baseline
480 (`ctc.safetensors`), and **batch-normalization**
481 **statistics** are also kept from the baseline
482 (`bn_stats.safetensors`). In addition, we re-
483 tain a baseline `rest.safetensors` block (i.e., the
484 remaining non-encoder/decoder parameters, in-
485 cluding feature-extraction/extractor-side and other
486 auxiliary modules) and exclude small tensors
487 (`small_tensors.safetensors`) from merging.

488 5.6 Evaluation Measures

489 We use Word Error Rate ($WER(\theta, \mathcal{D}_*)$) as the pri-
490 mary metric to measure ASR performances. To
491 ensure fair scoring under Arabic orthographic vari-
492 ability, we apply a lightweight normalization to
493 both references and hypotheses before computing

494 WER: we remove diacritics, normalize common
495 letter variants (e.g., *alef* forms, *hamza* variants, *alef*
496 *maksura* \rightarrow *ya*, and *teh marbuta*), and lowercase
497 Latin text. For **MyST**, we additionally exclude
498 from WER computation segments whose reference
499 is empty or contains only non-speech markers (e.g.,
500 `<silence>`, `<noise>`, (`silence`)), as well as unin-
501 telligible placeholders (e.g., `(())`); segments with
502 empty hypotheses are also skipped. This prevents
503 the metric from penalizing models for orthographic
504 conventions or no-speech regions.

505 Following, to quantify the trade-offs between
506 the inherent in model adaptation, we introduce two
507 specific ratios: the *Retention Index* (\mathcal{R}_{adult}) and
508 the *Adaptation Recovery* (\mathcal{R}_{child}).

509 **Retention Index** The Retention Index, $\mathcal{R}_{adult} =$
510 $\frac{WER(\theta_{base}, D_{adult}^L)}{WER(\theta_m, D_{adult}^L)}$ measures the extent to which the
511 merged model θ_m preserves the original capabili-
512 ties of the base model θ_{base} on adult benchmarks
513 D_{adult}^L ($L \in Ar, En$). with $\mathcal{R}_{adult} < 1.0$ represent
514 the forgetting cost occurred during the adaptation.

515 **Adaptation Recovery** The Adaptation Recovery,
516 $\mathcal{R}_{child} = \frac{WER(\theta_{ft}, D_{child}^L)}{WER(\theta_m, D_{child}^L)}$, measures the perfor-
517 mance gap between a unified merged model θ_m
518 and a child-specific ASR.

519 6 Results and Discussion

520 **Child Fine-Tuning vs. Adult Pretraining**
521 Across all three architectures – Encoder-Decoder
522 (ED), Encoder-CTC (ECTC), and WhisperLargev3
523 (WL) – fine-tuning the adult pretrained ASR model
524 on child speech consistently improves performance
525 on child-centric benchmarks (see Table 3). Rela-
526 tive to the adult-only baseline (θ_{base}), fine-tuning
527 with *Arabic child data* yields clear WER reductions
528 on AraKidsL1 and AraKidsL2, indicating effec-
529 tive adaptation to child acoustics and pronunciation
530 patterns. Incorporating *bilingual Arabic+English*
531 *child data* further strengthens this trend and leads
532 to more balanced improvements across AraKids
533 and MyST, particularly for encoder-decoder based
534 architectures: ED and EL. However, these gains are
535 accompanied by a generalization trade-off. While
536 child-domain WER improves substantially, fine-
537 tuned models may degrade on adult or out-of-
538 domain benchmarks, especially on MGB2 and, to a
539 lesser extent, Libri-C. This suggests that naive fine-
540 tuning on niche data such as children’s speech can
541 partially overwrite representations learned from
542 large-scale adult corpora.

Models	AraKidsL1	AraKidsL2	MyST	Libri-C	MGB2	Avg. WER
<i>Encoder-Decoder</i>						
θ_{base}	9.14	44.11	36.68	9.86	10.92	22.142
θ_{child}^{Ar}	6.41	26.76	44.23	12.21	10.7	20.062
θ_{child}^{En}	25.75	106.98	22.68	6.66	14.16	35.246
θ_{child}^{Ar+En}	6.37	22.71	22.52	6.66	12.97	14.312
<i>Encoder-CTC</i>						
θ_{base}	14.72	63.37	26.63	3.83	13.18	24.346
θ_{child}^{Ar}	13.11	16.32	77.32	18.64	23.63	29.804
θ_{child}^{En}	51.44	92.73	15.11	6.8	26.26	38.468
θ_{child}^{Ar+En}	12.45	24.24	15.9	6.66	21.2	16.09
<i>Whisper Large v3</i>						
θ_{base}	9.68	40.68	21.96	2.47	15.15	17.988
θ_{child}^{Ar+En}	8.7	11.84	14.47	3.11	23.52	12.328

Table 3: Reported WER for the base adult-specific θ_{base} model and specialized model with naive fine-tuning for 3 different ASR architecture.

Model Arch	MOperation	AraKidsL1	AraKidsL2	MyST	Libri-C	MGB2	Avg. WER
<i>Encoder-Decoder (ED)</i>							
ED	LERP	7.71	31.88	33.59	9.33	10.84	18.67
ED	LERP (Enc)	7.86	32.04	33.86	9.16	10.86	18.756
ED	TIES	7.14	26.76	28.67	9.23	11.26	16.612
ED	TIES (Enc)	7.35	27.01	28.25	8.16	10.86	16.326
<i>Encoder-CTC (ECTC)</i>							
ECTC	LERP	10.38	39.56	21.91	3.91	13.13	17.778
ECTC	TIES	11.07	18.94	18.68	5.66	18.71	14.612
<i>Whisper-Large v3 (WL)</i>							
WL	LERP	6.55	20.01	17.78	2.43	14.49	12.252
WL	LERP (Enc)	7.64	26.38	20.45	2.41	15.00	14.376
WL	TIES	8.48	11.37	14.36	2.91	22.68	11.96
WL	TIES (Enc)	7.16	19.83	17.55	2.61	17.01	12.832

Table 4: WER (%) across benchmarks for different merge operations. Within each architecture block (ED/ECTC/WL), the best (lowest) score in each column is highlighted in blue.

Arabic-only vs. bilingual child models and Model Merging Across merging setups, using the bilingual child model θ_{child}^{Ar+En} tends to outperform using Arabic-only child models θ_{child}^{Ar} (see Table 6). This pattern indicates that bilingual fine-tuning induces more robust intermediate representations that generalize better when merged with the adult anchor, reducing over-specialization to a single child-language distribution. At the same time, θ_{child}^{Ar} encodes Arabic-centric representations that remain useful to disentangle Ar-only vs. Ar+En directions in weight space. Based on the empirical comparison in Table 3, we therefore select θ_{child}^{Ar+En} and θ_{child}^{Ar} as primary merging candidates anchored with θ_{base} .

LERP vs. TIES. For our reported results, in Table 4, we observed LERP provides moderate improvements on children benchmark, but its recovery on adult benchmarks is remarkable. This is because the interpolation model gives much priority to the base (as anchor). In contrast, TIES consistently reduces average WER across architectures

for children ASR benchmark while lacking in the retention capabilities, suggesting the operation is more effective for alignment child-adapted parameters, while retaining major portion of the results. The advantage of TIES is most pronounced for ED and WL, where it yields a stronger adaptation.

Architectural Considerations in Model Merging: For ED and WL architectures, we compare merging *encoder-only* parameters against the *full model* (encoder and decoder) in Table 4. For TIES, Encoder-only merging generally yields superior adult-domain retention – most notably on the MGB2 benchmark, though it may slightly reduce the competitive edge in child-domain specialization relative to full-model merging. This disparity in stability arises because decoders adapted to niche domains, such as children’s speech, may collapse toward domain-specific token distributions, thereby reducing their capacity to retain the linguistic breadth acquired during large-scale adult pre-training. By preserving the adult decoder as a stable language model anchor, encoder-only merg-

ing avoids this failure mode while still benefiting from child-adapted acoustic representations in the encoder.

In contrast, for CTC-based architectures, merging the CTC linear projection layer alongside the encoder tends to reduce overall adaptation capability. Due to the limited capacity of the CTC head and its lack of an internal autoregressive language model, the weights are highly sensitive to distribution shifts. Consequently, merging these layers with the base model often dilutes the specialized child-speech distributions learned during fine-tuning, leading to a more significant performance degradation than that observed in ED decoders. This sensitivity is reflected in a simple CTC-head ablation: using the adult CTC head yields 14.61 WER, while swapping to child-specialized CTC heads results in 15.41 WER (CTC head of θ_{child}^{Ar+En} , +0.80) and 18.79 WER (CTC head of θ_{child}^{Ar} , +4.18).

Key Findings: From our exploratory experiments, we observed the following:

•**Adaptation-Generalization Trade-off:** While fine-tuning adult pretrained models consistently improves performance on child-centric benchmarks, it introduces a generalization trade-off where representations learned from large-scale adult corpora are partially overwritten by niche child-speech data.

•**Bilingual Robustness:** Bilingual child models (θ_{child}^{Ar+En}) produce more robust intermediate representations than monolingual experts. This reduces over-specialization and allows for better alignment when merged with adult anchors.

•**Operator Specialization:** The choice of merging operator determines the performance profile:

–**LERP** excels at adult-domain retention by prioritizing the base anchor, making it ideal for stability.

–**TIES** excels at child-domain adaptation, effectively aligning child-specific parameters to achieve lower WERs in complex architectures like ED and WL.

•**Decoder Stability vs. CTC Constraints:**

–In **Encoder-Decoder** and **Whisper** architectures, *encoder-only merging* is the optimal strategy to prevent "decoder collapse," preserving the adult linguistic breadth while adapting to child acoustics.

–In **CTC** architectures, the linear projection layer lacks the capacity to handle distribution merging, often leading to performance dilution when merged with the base model.

7 Conclusion

Weight-space model merging consistently improves children’s ASR while retaining strong adult performance, offering a practical balance between adaptation and retention without additional training cost. Across architectures and benchmarks, merging reduces child-domain WER and mitigates the degradation typically observed after naive fine-tuning, supporting a single inclusive ASR model across age, nativity (L1/L2), and language (Arabic/English).

8 Limitations

While our results consistently show that weight-space model merging improves children’s ASR without sacrificing adult retention, our hyperparameter exploration is limited. For each merging technique (e.g., LERP and TIES), we performed only a small manual sweep over the main parameters (interpolation weights, density, and update weights) and selected a single configuration based on the best average WER across all evaluation sets. We did not conduct an exhaustive or automated search (e.g., grid/random/Bayesian optimization), nor did we systematically study sensitivity to these parameters across architectures and datasets. As a result, the reported numbers may underestimate the best achievable performance, and a more comprehensive tuning study could further improve both adaptation and retention, or reveal more stable parameter regimes.

References

- Edem Ahadzi, Vishwanath Pratap Singh, Tomi Kinnunen, and Ville Hautamäki. 2025. [Continuous learning for children’s asr: Overcoming catastrophic forgetting with elastic weight consolidation and synaptic intelligence](#). In *Proc. Interspeech*.
- Ahmed Ali, Peter Bell, James Glass, Yacine Messaoui, Hamdy Mubarak, Steve Renals, and Yifan Zhang. 2016. The MGB-2 challenge: Arabic multi-dialect broadcast media recognition. In *2016 IEEE Spoken Language Technology Workshop (SLT)*, pages 279–284. IEEE.
- V. Bhardwaj and 1 others. 2022. [Automatic speech recognition \(asr\) systems for children](#). *Applied Sciences*.

684	Alexandre Ducorroy and Rachid Riad. 2025. Robust fine-tuning of speech recognition models via model merging: application to disordered speech . In <i>Proc. Interspeech</i> .	Mitchell Wortsman, Gabriel Ilharco, Samir Yitzhak Gadre, and 1 others. 2022. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time . <i>arXiv preprint arXiv:2203.05482</i> .	737
685			738
686			739
687			740
688	Brady Houston and Katrin Kirchhoff. 2020. Continual learning for multi-dialect acoustic models . In <i>Proc. Interspeech</i> .	F. Wu and 1 others. 2019. Advances in automatic speech recognition for child speech . In <i>Proc. Interspeech</i> .	742
689			743
690			
691	Gabriel Ilharco and 1 others. 2023. Editing models with task arithmetic . In <i>ICLR (OpenReview)</i> .	Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. 2023. Ties-merging: Resolving interference when merging models . <i>arXiv preprint arXiv:2306.01708</i> .	744
692			745
693	R. Jain, L. Barcovschi, and 1 others. 2023. Adaptation of whisper models to child speech recognition . In <i>arXiv preprint arXiv:2307.13008</i> .		746
694			747
695			
696	Tom Ko, Vijayaditya Peddinti, Daniel Povey, and Sanjeev Khudanpur. 2015. Audio augmentation for speech recognition . In <i>Interspeech 2015</i> , pages 3586–3589.		
697			
698			
699			
700	Linguistic Data Consortium. 2021. Myst (my science tutor) children’s conversational speech . LDC Catalog.		
701			
702			
703	Guillermo Ortiz-Jimenez, Alessandro Favero, and Pascal Frossard. 2023. Task arithmetic in the tangent space: Improved editing of pre-trained models . In <i>NeurIPS</i> .		
704			
705			
706			
707	Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. Librispeech: An ASR corpus based on public domain audio books . In <i>Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 5206–5210.		
708			
709			
710			
711			
712			
713	Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le. 2019. SpecAugment: A simple data augmentation method for automatic speech recognition . In <i>Interspeech 2019</i> , pages 2613–2617.		
714			
715			
716			
717			
718	Sameer Pradhan, Ronald A. Cole, and Wayne H. Ward. 2024. My science tutor (MyST)—a large corpus of children’s conversational speech . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 12040–12045, Torino, Italia. ELRA and ICCL.		
719			
720			
721			
722			
723			
724			
725	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine Mcleavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision . In <i>Proceedings of the 40th International Conference on Machine Learning</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 28492–28518. PMLR.		
726			
727			
728			
729			
730			
731			
732	Samik Sadhu and Hynek Hermansky. 2020. Continual learning in automatic speech recognition . In <i>Proc. Interspeech</i> .		
733			
734			
735	Yi-Lin Sung and 1 others. 2023. An empirical study of multimodal model merging . In <i>Findings of EMNLP</i> .		
736			
		A Appendix	748

Arch.	Merge (*)	Adaptation Recovery (AR)				Retention Index (RI)		
		AraKidsL1	AraKidsL2	MyST	Avg.	Libri-C	MGB2	Avg.
ED	LERP	82,62%	71,24%	67,04%	73,63%	100%↑	100% ↑	100%↑
	TIES	89,22%	84,87%	78,55%	84,21%	100%↑	96,98%	98,49%
ECTC	LERP	61,37%	57,41%	100%↑	72,92%	100%↑	83,17%	91,58%
	TIES	57,54%	100%↑	100% ↑	85,84%	100%↑	58,36%	79,18%
WL	LERP	97,25%	100%↑	100%↑	99,03%	100%↑	75,36%	87,68%
	TIES	75,12%	100%↑	100%↑	91,70%	100%↑	48,15%	74,07%

Table 5: AR and RI summary across architectures and merge methods. Values $\geq 100\%$ are capped at 100% and marked with \uparrow .

Technique	Setup	Children			Adult		Overall
		AraKidsL1	AraKidsL2	MyST	Libri-C	MGB2	Avg. WER
LERP	M0 + B	10,27	34,92	27,8	4,01	12,98	17,996
	M0 + C	11,65	46,72	19,18	3,97	13,47	18,998
	M0 + A + B	11,45	52,34	21,52	3,9	13,1	20,462
	M0 + B + C	10,38	39,56	21,91	3,91	13,13	17,778
	M0 + A + B + C	11,72	48,32	19,21	3,92	13,24	19,282

Table 6: Impact of Arabic-only vs. bilingual child specialists in merging. M0 denotes the adult baseline; A/B/C denote child-specialized checkpoints used as merging candidates (A: MyST-only, B: AraKids-only, C: AraKids+MyST bilingual). Average WER is computed over all reported test sets.