

Navigating the Reality Gap: Privacy-Preserving On-Device Continual Adaptation of ASR for Clinical Telephony

Anonymous ACL submission

Abstract

Automatic Speech Recognition (ASR) holds immense potential to assist in clinical documentation and patient report generation, particularly in resource-constrained regions. However, deployment is currently hindered by a technical deadlock: a severe “**Reality Gap**” between laboratory performance and noisy, real-world clinical audio, coupled with strict privacy and resource constraints. Such adaptation is essential for clinical telephony systems, where patient speech is highly variable and transcription errors can directly impact downstream clinical workflows. We quantify this gap, showing that a robust multilingual model (IndicWav2Vec) degrades up to a **40.94% WER** on rural clinical telephony speech from India, rendering it unusable. We demonstrate consistent improvements on these helpline interactions without transmitting raw patient data off-device via an on-device continual adaptation framework using Low-Rank Adaptation (LoRA). We conduct an investigative study of stabilization strategies, characterizing the trade-offs between data-driven and parameter-driven approaches. Our results demonstrate that multi-domain Experience Replay (ER) yields the primary performance gains, achieving a **17.1% relative improvement** in target WER and reducing catastrophic forgetting by **55%** compared to naive adaptation. Furthermore, we investigate a stabilized importance estimation strategy (Absolute Fisher) to ensure robust convergence against the high-variance gradients common in clinical telephony speech. Finally, we verify via a domain-specific spot check that acoustic adaptation is a fundamental prerequisite for usability in healthcare settings which cannot be bypassed by language models alone.

1 Introduction

The recent surge in Self-Supervised Learning (SSL) has propelled Automatic Speech Recognition (ASR) to near-human performance on standardized benchmarks. Foundational models like

Meta’s Wav2Vec 2.0 (Baevski et al., 2020) and OpenAI’s Whisper (Radford et al., 2023) promise a future where automated transcription can digitize patient reports, allowing clinics to improve throughput with reduced operational costs. However, for specialized, high-impact clinical domains such as rural health helplines and remote clinical services, this promise remains unfulfilled. The “reality gap,” representing the disparity between clean training corpora and the chaotic, noisy, and privacy-constrained environments of real-world clinical telephony, renders these state-of-the-art models practically unusable.

Our baseline analysis reveals that even robust multilingual models like IndicWav2Vec (Javed et al., 2022) degrade to a prohibitively high **40.94% Word Error Rate (WER)** when exposed to real-world clinical telephony audio from rural India (Bhanushali et al., 2022). This failure is compounded by a technical deadlock: patient data privacy laws preclude the use of cloud-based adaptation services (Leroy et al., 2019), while rural infrastructure lacks the high-end compute required for traditional model retraining. This creates a scenario where models cannot improve because data cannot leave the local environment, and localized compute is too constrained for standard fine-tuning.

To break this deadlock, we propose an on-device adaptation framework for localized, stream-based learning. We focus on data residency: ensuring that raw patient audio never leaves the local device. By leveraging Low-Rank Adaptation (LoRA) (Hu et al., 2022), we enable the model to fine-tune on incoming clinical data streams using a fraction of the trainable parameters. However, sequential adaptation in long-running clinical deployments introduces the risk of *catastrophic forgetting*, where the model loses its foundational linguistic capabilities (McCloskey and Cohen, 1989). While recent benchmarks on real-world Indian speech, such as NIRANTAR (Javed et al., 2025), indicate that

085 standard regularization methods (EWC, MAS) of 32
086 ten fail to prevent this forgetting in diverse do-
087 mains, we hypothesize that this instability is ex- 133
088 acerbated by the acoustic mismatch inherent in
089 low-bandwidth clinical telephony. To counteract
090 this, we integrate a multi-domain experience re-
091 play mechanism (Chaudhry et al., 2019b), inter-
092 leaving small buffers of general-domain data with
093 the incoming clinical stream to anchor the learning
094 trajectory.

095 Furthermore, we investigate parameter-
096 regularization strategies to stabilize this pipeline.
097 While standard quadratic formulations (EWC)
098 can be sensitive to the high-variance gradients
099 of noisy telephony speech, we observe that
100 a balanced importance estimation (Absolute
101 Fisher) (Benzing, 2021) provides a more resilient
102 regularization signal for this localized application.
103 Our results demonstrate that these strategies yield
104 a 17.1% relative reduction in WER on the target
105 clinical domain, effectively stabilizing the model’s
106 performance while respecting rigid data residency
107 constraints. The primary research contributions of
108 our work are:

- 109 • **On-Device Clinical Adaptation:** A localized
110 continual adaptation framework for clinical
111 telephony ASR that respects healthcare data
112 residency constraints while maintaining com-
113 putational efficiency.
- 114 • **Empirical Validation on Clinical Speech:**
115 Parameter-efficient adaptation (LoRA) com-
116 bined with multi-domain memory (Replay)
117 yields a 17.1% relative improvement in WER
118 on real-world rural clinical speech without
119 centralized retraining or raw data exfiltration.
- 120 • **Investigation of Optimization Robustness:**
121 The interplay between data-driven replay and
122 parameter-driven regularization, validating
123 that linearized Fisher importance promotes
124 more reliable convergence in noisy, multi-
125 speaker clinical deployments.
- 126 • **Quantification of the Clinical Reality Gap:**
127 Establish a rigorous benchmark for 8kHz clini-
128 cal telephony, providing insights into the tech-
129 nical prerequisites to close the gap between
130 foundational ASR models and frontline health-
131 care needs.

2 Related Work

Designing a privacy-preserving, adaptive ASR
pipeline for low-resource medical settings requires 134
synthesizing advancements in self-supervised 135
acoustic modeling, parameter-efficient adaptation, 136
and continual learning. We address critical gaps 137
regarding adaptation efficiency, data scarcity, and 138
optimization stability. 139

2.1 Self-Supervised Acoustic Foundations and the Domain Gap 140

Low-resource speech recognition increasingly re- 142
lies on Self-Supervised Learning (SSL) to lever- 143
age unlabeled audio. We build upon Wav2Vec 2.0 144
(Baevski et al., 2020), which achieves high data 145
efficiency through contrastive learning; Baevski 146
et al. (2020) showed that fine-tuning on just ten 147
minutes of data yields competitive performance 148
on benchmarks like Librispeech (Panayotov et al., 149
2015). To support linguistic diversity, we utilize 150
IndicWav2Vec (Javed et al., 2022), scaling the ar- 151
chitecture to cover 40 Indian languages. 152

The Gap: Despite these capabilities, a “usability 153
gap” exists between laboratory benchmarks and 154
real-world deployment. Clinical ASR errors are 155
particularly risky because they affect medically 156
salient information such as dosages, negations, and 157
symptom reporting (Chiu et al., 2018a,b). This 158
is exacerbated in rural healthcare by noisy tele- 159
phonic audio and diverse regional dialects. Recent 160
large-scale efforts such as NIRANTAR (Javed et al., 161
2025) have begun characterizing these challenges 162
across 22 Indian languages. We utilize the Gram 163
Vaani ASR dataset (Bhanushali et al., 2022) to rep- 164
resent these challenges; Bhanushali et al. (2022) 165
highlight that standard models struggle to general- 166
ize to such 8kHz telephony audio without targeted 167
adaptation. 168

While signal-level enhancement strategies, 169
such as bandwidth expansion and speech super- 170
resolution (Lin et al., 2023; Li et al., 2019), offer an 171
alternative by reconstructing high-frequency com- 172
ponents, they introduce significant computational 173
overhead. Deploying a separate BWE model (often 174
GAN or diffusion-based) alongside the ASR sys- 175
tem contradicts the strict low-latency requirements 176
of edge devices. Furthermore, privacy constraints 177
preclude cloud-based services (Leroy et al., 2019), 178
creating a deadlock where models cannot improve 179
because data cannot leave the local environment. 180

181 2.2 Parameter-Efficient Fine-Tuning (PEFT) 182 and the Efficiency Gap

183 Bridging the domain gap typically requires fine-
184 tuning on target data. However, for Large Au-
185 dio Models (LAMs), full fine-tuning is computa-
186 tionally prohibitive and prone to overfitting. This
187 presents an *efficiency gap*: high-end GPUs are un-
188 available on the resource-constrained edge devices
189 of rural hospitals.

190 To address this, we adopt Low-Rank Adapta-
191 tion (LoRA) (Hu et al., 2022), which enables local
192 fine-tuning on modest clinical hardware. This ap-
193 proach is essential for meeting healthcare **Data**
194 **Residency** requirements, as it allows for local-
195 ized, high-performance adaptation without need-
196 ing the centralized data storage often required by
197 federated or cloud-based approaches (Rieke et al.,
198 2020). Recent work by Song et al. (2024) on
199 LoRA-Whisper validates this approach for multilin-
200 gual ASR, demonstrating that PEFT can effectively
201 adapt large Transformers without catastrophic in-
202 terference. By constraining optimization to a low-
203 dimensional subspace, LoRA resolves deployment
204 challenges in privacy-sensitive clinics.

205 While PEFT solves efficiency, self-improving
206 systems face the *stability gap* of Continual Learn-
207 ing (CL), specifically *Catastrophic Forgetting*,
208 where new training erodes previously learned capa-
209 bilities (McCloskey and Cohen, 1989). In health-
210 care, ensuring model stability over time is a critical
211 safety requirement for clinical evaluation cycles
212 (De Lange et al., 2022). In ASR, naive updates on
213 recent telephony transcripts lead to overfitting on
214 specific speakers or acoustic conditions. Bench-
215 marks like NIRANTAR (Javed et al., 2025) further
216 underscore this, demonstrating that no single CL
217 strategy currently yields consistent performance
218 across the diverse shifts encountered in real-world
219 Indian speech.

220 To mitigate this, we employ a hybrid contin-
221 ual learning strategy that synergizes data-level and
222 parameter-level consolidation. First, to address the
223 distribution shift between segments, we implement
224 **Experience Replay (ER)** (Chaudhry et al., 2019a).
225 Building on recent regularized CL studies for Indic
226 ASR (T and Nirmala, 2025) and episodic memory
227 frameworks (Yang et al., 2022), we utilize a *multi-*
228 *domain experience replay* mechanism. While Yang
229 et al. (2022) focused on gradient projection, we
230 find that a direct data mixing strategy effectively
231 mitigates catastrophic forgetting when combined

with PEFT. Our buffer retains a mixture of “hard” 232
examples (high-loss) from the target domain (Gram 233
Vaani) and gender-balanced samples from a clean 234
auxiliary domain (Kathbath). 235

Complementing this, we constrain parameter 236
drift using **Online LoRA-EWC** (Xiang et al., 237
2023), utilizing the **Absolute Fisher** importance 238
estimation (Benzing, 2021). We estimate impor- 239
tance using the *absolute value* of accumulated gra- 240
dients ($|g|$) rather than the standard squared Fisher 241
Information. While this linear scaling draws inspi- 242
ration from the numerical stability of Memory 243
Aware Synapses (MAS) (Aljundi et al., 2018), Ab- 244
solute Fisher explicitly ties the importance to the 245
likelihood of the adaptation task rather than un- 246
supervised output sensitivity. By combining this 247
regularized signal with multi-domain replay, our 248
approach prevents catastrophic forgetting through 249
a dual mechanism: replay maintains the data distri- 250
bution, while the regularization penalty preserves 251
critical adapter parameters. 252

253 3 Methodology

We propose an on-device, adaptive ASR framework 254
designed to bridge the “reality gap” and ensure 255
deployment-ready reliability within the acoustic 256
constraints of rural clinical telephony. Our ap- 257
proach focuses on Continual Learning (CL), en- 258
abling a pre-trained model to adapt sequentially 259
to incoming clinical data streams (\mathcal{D}_{stream}) while 260
maintaining strict data residency and preventing 261
catastrophic forgetting. 262

263 3.1 Base Acoustic Backbone

We utilize **IndicWav2Vec** as our acoustic back- 264
bone. This model is built on the Wav2Vec 2.0 265
architecture and pre-trained on a massive corpus 266
of diverse Indian languages. The network consists 267
of a convolutional feature encoder $f(x)$ mapping 268
raw audio to latent representations, followed by a 269
Transformer context network optimized via Con- 270
nectionist Temporal Classification (CTC) loss: 271

$$272 \mathcal{L}_{CTC} = -\log P(y|x) \quad (1)$$

To enable efficient adaptation on edge devices, 273
we freeze the base model and inject trainable 274
Low-Rank Adaptation (LoRA) matrices into the 275
query and value projection layers. This parameter- 276
efficient approach serves as a stability safeguard, 277
preventing large-scale corruption of pre-trained 278

weights while allowing for specialized clinical adaptation. 280

3.2 Data-Driven Stability: Multi-Domain Experience Replay 281

To mitigate catastrophic forgetting, we employ Experience Replay (ER), explicitly grounding the model’s optimization trajectory with historical data. 282 We implement a **Multi-Domain Replay** strategy 283 that maintains a dual-source buffer \mathcal{B} : 284

- **General Domain Anchor** (\mathcal{B}_{gen}): A gender-balanced subset of high-resource, standard Hindi samples (sourced from Kathbath) to preserve foundational phonetic robustness, which is essential for handling diverse patient demographics in a remote clinical setting. 285
- **Target Domain History** (\mathcal{B}_{spec}): A sliding window of "hard" examples (high loss) and random samples from previous clinical segments, ensuring retention of recent domain-specific adaptations. 286

During training, the incoming clinical stream \mathcal{D}_{stream} is concatenated with samples from \mathcal{B} . The optimization objective for Experience Replay is: 287

$$\mathcal{L}_{ER} = \gamma \mathbb{E}_{x \sim \mathcal{D}_{stream}} [\mathcal{L}_{CTC}(x)] + (1 - \gamma) \mathbb{E}_{x \sim \mathcal{B}} [\mathcal{L}_{CTC}(x)] \quad (2)$$

where γ represents the mixing ratio between new data and replayed data. 288

3.3 Parameter-Driven Stability: Elastic Weight Consolidation (EWC) 289

As a standalone strategy (V4.5) and a component of our hybrid approach, we implement Elastic Weight Consolidation (EWC) to prevent drift in critical parameters. Unlike ER, which requires data storage, EWC regularizes the model by penalizing changes to weights that are important for previous tasks. 290

We compute the importance of each LoRA parameter θ_i using the diagonal of the Fisher Information Matrix (F). To ensure robust optimization against high-variance gradients, we approximate importance using accumulated **absolute gradients** rather than squared gradients. This ‘‘Absolute Fisher’’ approach (Benzing, 2021) functions as a reliability optimization. We distinguish this from Memory Aware Synapses (MAS) (Aljundi et al., 2018); while MAS assesses importance based 291

on the sensitivity of the output function, Absolute Fisher derives importance from a linearized measure of the likelihood’s Fisher Information, preventing transient acoustic artifacts from drowning out meaningful weights: 292

$$F_i = \frac{1}{N} \sum_{j=1}^N |\nabla_{\theta_i} \log P(y_j|x_j)| \quad (3)$$

This linear accumulation prevents outliers from dominating the importance metric. The EWC regularization loss is then applied as a quadratic penalty: 293

$$\mathcal{L}_{EWC}(\theta) = \frac{\lambda}{2} \sum_i F_i (\theta_i - \theta_i^*)^2 \quad (4)$$

where θ_i^* represents the frozen optimal parameters from the previous segment and λ controls the regularization strength. 294

3.4 Hybrid Optimization Framework (V5.1) 295

To leverage the synergistic effects of data-driven grounding and parameter-driven constraints, we propose a Hybrid ER + EWC optimization framework (Strategy V5.1). This approach combines the Multi-Domain replay buffer with the EWC regularization penalty. 296

The total optimization objective minimizes the transcription error over the mixed batch while simultaneously constraining parameter drift: 297

$$\mathcal{L}_{Total} = \mathcal{L}_{ER}(\mathcal{D}_{stream}, \mathcal{B}) + \mathcal{L}_{EWC}(\theta) \quad (5)$$

By using \mathcal{L}_{ER} to provide the necessary gradients for acoustic adaptation and \mathcal{L}_{EWC} to define a ‘‘safe’’ optimization region, this hybrid mechanism aims to maximize target domain accuracy while minimizing catastrophic forgetting. 298

4 Experiments and Results 299

We conducted a comprehensive evaluation to validate the effectiveness of our adaptive pipeline. The experimental design focuses on two key aspects: the ability to adapt to a specific clinical domain (Gram Vaani) and the ability to retain general linguistic knowledge (Kathbath) to prevent catastrophic forgetting. 300

365	4.1 Experimental Setup	362	4.2.1 Paradigm 1: Naive Continual Fine-tuning (Baseline)	409
366	4.1.1 Metrics	362	4.2.1 Paradigm 1: Naive Continual Fine-tuning (Baseline)	410
367	We evaluate performance using two standard metrics: Word Error Rate (WER) and Character Error Rate (CER). WER measures transcription accuracy at the word level, while CER provides a finer-grained analysis of phonetic accuracy, particularly useful for agglutinative languages and dialectal variations.	365	This represents the simplest form of adaptation, where the model is fine-tuned sequentially on incoming clinical data streams without explicit regularization. We utilize a LoRA rank of 16 and $\alpha = 32$.	411
370	$\text{WER} = \frac{S_w + D_w + I_w}{N_w}$	368		412
371	$\text{CER} = \frac{S_c + D_c + I_c}{N_c}$	369		413
372		370		414
373	4.1.2 Datasets	371		415
374	• Gram Vaani (Clinical Helpline Proxy) (Bhanushali et al., 2022): This dataset consists of rural telephonic speech (originally 8kHz, upsampled to 16kHz) and serves as a proxy for audio encountered in frontline healthcare helplines. We evaluate on the official 3-hour Evaluation set (GV_Eval_3h). Note that prior work (Patel and Scharenborg, 2022) reported 30.3% WER on the 5-hour Development set using full Conformer models. Our result (33.94%) reflects performance on the unseen Evaluation partition under strictly constrained on-device conditions (LoRA vs. full tuning). To strictly simulate a continual clinical learning scenario, we partition the 103 hours of training data into sequential segments, processing them one by one to mimic a live data stream while ensuring patient data residency.	372	• Outcome (V1.1): While the model achieves a significant reduction in target WER (down to 34.00%), it suffers from severe catastrophic forgetting, with general domain error increasing by an absolute 5.93% .	416
375		373		417
376		374		418
377		375		419
378		376		420
379		377		421
380		378		422
381		379		423
382		380		424
383		381		425
384		382		426
385		383		427
386		384		428
387		385		429
388		386		430
389		387		431
390		388		432
391		389		433
392	• Kathbath (General Domain): A high-quality, read speech dataset representing standard Hindi (Javed et al., 2023). We utilize a subset of the training set (approx. 25,800 samples) to populate the experience replay buffer, ensuring the model retains knowledge of standard Hindi. The complete validation set (3,151 samples) is used exclusively to measure catastrophic forgetting after adaptation.	390		434
393		391		435
394		392		436
395		393		437
396		394		438
397		395		439
398		396		440
399		397		441
400		398		442
401	4.2 Continual Adaptation Paradigms	399		443
402	To characterize the optimal pathway for localized adaptation, we investigate four distinct continual learning paradigms. For all experiments, we utilize a base IndicWav2Vec model, which provides a pre-training baseline of 40.94% WER on the target rural clinical domain and 11.57% on the general Kathbath domain.	400		444
403		401		445
404		402		446
405		403		447
406		404		448
407		405		449
408		406		450
		407		451
		408		452
		409		453
		410		
		411		
		412		
		413		
		414		
		415		
		416		
		417		
		418		
		419		
		420		
		421		
		422		
		423		
		424		
		425		
		426		
		427		
		428		
		429		
		430		
		431		
		432		
		433		
		434		
		435		
		436		
		437		
		438		
		439		
		440		
		441		
		442		
		443		
		444		
		445		
		446		
		447		
		448		
		449		
		450		
		451		
		452		
		453		

Paradigm	Strategy	Final Target WER	Improvement (%)	Final General WER	Forgetting
Baseline	Pre-trained IndicWav2Vec	40.94%	-	11.57%	-
Naive	V1.1 Naive Fine-tuning	34.00%	+17.0%	17.50%	+5.93%
Experience Replay	V2.1 Single-Domain ER	33.98%	+17.0%	14.90%	+3.33%
Experience Replay	V3.1 Multi-Domain ER	33.94%	+17.1%	14.23%	+2.66%
EWC	V4.5 EWC ($\lambda = 1e1$)	33.94%	+17.1%	15.15%	+3.58%
Hybrid	V5.1 ER + EWC	34.51%	+15.7%	14.14%	+2.57%

Table 1: Comprehensive performance comparison of the investigated continual learning paradigms. Improvement is relative to the IndicWav2Vec baseline on rural clinical data. Forgetting is the absolute increase in WER on the general Kathbath domain.

4.2.4 Paradigm 4: Hybrid ER + EWC Optimization

Finally, we evaluate the hybrid framework (V5.1) described in Section 3.4, which combines the data-driven guidance of Multi-Domain ER with the parameter-driven Absolute Fisher constraints.

- **Outcome (V5.1):** This paradigm yielded our most stable model, achieving the lowest absolute forgetting rate with a final general WER of **14.14%** (+2.57% increase). However, this stability came at a slight cost to plasticity, with a final target WER of **34.51%**.

5 Analysis and Discussion

5.1 Comparative Performance Summary

Table 1 summarizes the final metrics for the five primary investigative strategies. All proposed paradigms successfully bridge the “Reality Gap,” reducing error on the clinical domain by at least 15% relative.

5.2 Stability-Plasticity Pareto Efficiency

The efficiency of our adaptation strategies is best characterized through a Pareto analysis of retention versus adaptation (Figure 3). We observe that with the current amount of continual learning, multiple successful strategies converge to a range around 34-35% WER. Below this threshold, further reductions in target WER begin to incur disproportionately high stability costs. For instance, while Multi-Domain ER (V3.1) successfully navigates this trade-off, attempts to achieve deeper adaptation in the Hybrid model (V5.1) encountered a stability bottleneck. This behavioral ceiling suggests that for a given model architecture and domain mismatch, there exists an optimal frontier of acoustic adaptation. Crucially, the V3.1 trajectory has not plateaued but has entered a logarithmic long-tail learning phase, suggesting that further gains are

possible but require significantly more data or linguistic context to unlock.

5.3 Mechanism Analysis: Data vs. Parameter Anchors

Our investigation reveals a fundamental difference between data-driven and parameter-driven regularization in lightweight models.

- **Data Anchors (ER):** Replaying Kathbath samples provides an explicit gradient signal toward a generalized linguistic center. This was the most efficient stabilizer, as shown by the V3.1 Pareto trajectory.
- **Parameter Anchors (EWC):** While EWC protects weights, it is “acoustically blind,” penalizing drift without reconciling new acoustic features. This is evidenced by V2.1 (ER) providing superior stabilization (+3.33% forgetting) compared to V4.5 (EWC, +3.58%) despite similar target performance.
- **The LoRA-Regularization Bottleneck:** We observe a conflict in the Hybrid model (V5.1), where combining LoRA’s low-rank space with EWC’s parameter-level penalties overly restricts optimization paths, leading to higher target WER.

5.4 Convergence Dynamics

While the Naive baseline (V1.1) exhibits erratic loss curves, the ER and Hybrid paradigms demonstrate controlled, monotonic trends (Figure 1). This stability is a critical requirement for clinical deployment, ensuring the system does not regress after extended use. However, we observe a “soft floor” around 34% WER, which we characterize as the **Acoustic Bottleneck of Clinical Telephony**. The mismatch between 8kHz telephony and the 16kHz model expectation creates a ceiling for acoustic-only adaptation, suggesting that further gains may

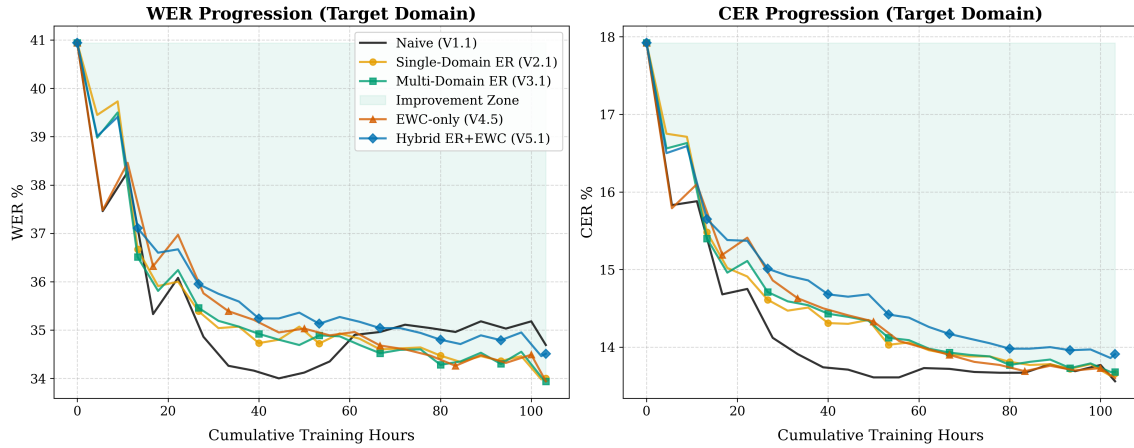


Figure 1: Progression of recognition performance on the target rural clinical data over 100 cumulative training hours. All investigated paradigms successfully bridge the “Reality Gap,” with Multi-Domain ER (V3.1) achieving the deepest adaptation.

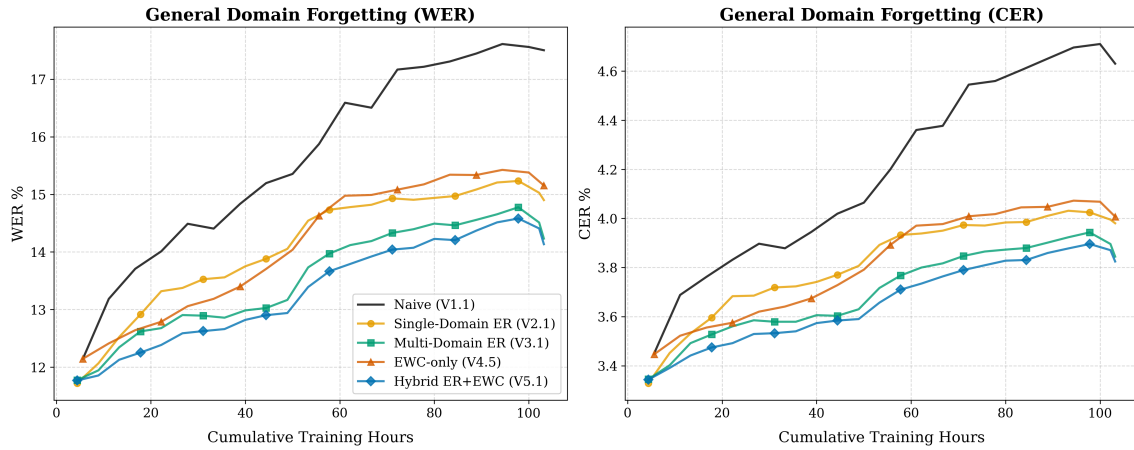


Figure 2: Catastrophic forgetting analysis on the general Kathbath domain. The plot illustrates how knowledge retention is preserved over time. Standard Naive adaptation (V1.1) and EWC-only (V4.5) show significant drift, whereas Multi-Domain ER and Hybrid models successfully flatten the error trajectory.

528 require specialized front-end processing or linguist-
529 tic context.

530 5.5 The Necessity of Acoustic Adaptation

531 To verify that acoustic adaptation is a *prerequisite*⁵⁴⁰
532 *sic* for clinical deployment, we paired the base-⁵⁴¹
533 line and adapted models with a domain-specific⁵⁴²
534 4-gram LM (Table 2). The unadapted baseline+LM
535 achieves only **34.96% WER**, while V3.1+LM
536 reaches **30.26% WER**. This confirms that **Acous-**
537 **tic Data Residency** is the primary hurdle; solving
538 for local acoustic mismatch is fundamental before
539 linguistic correction can be effectively applied.

Model	WER (%)	CER (%)
Baseline + LM	34.96	17.59
V3.1 Adapted + LM	30.26	15.08

Table 2: LM Spot Check: Acoustic adaptation remains essential even with LM assistance.

5.6 Clinical Implications and Deployment Feasibility

The **17.1% relative improvement** in WER demon-
strated by our Multi-Domain ER strategy (V3.1)
represents a significant bridge between “raw” and
“usable” clinical ASR. In a clinical dictation set-
ting, this reduction minimizes the manual correc-
tion burden on healthcare workers, whose time
is a critical resource in rural India. By enabling

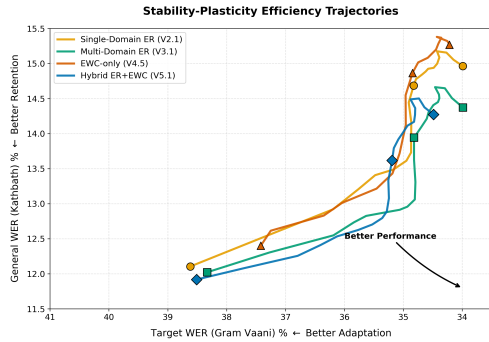


Figure 3: Stability-Plasticity Pareto Analysis: The trajectories illustrate the “cost” of adaptation. Notably, Multi-Domain ER (V3.1) achieves the best overall balance, outperforming other paradigms by maintaining the lowest error rates on both stability and plasticity metrics, whereas Hybrid and EWC approaches incur higher stability costs for similar plasticity gains.

549 the system to stabilize on local dialects and noisy
 550 line conditions without transmitting patient data to
 551 central servers, we provide a technical blueprint
 552 for scalable, privacy-compliant medical transcrip-
 553 tion. We note that prior work on the Gram Vaani
 554 dataset (Patel and Scharenborg, 2022) achieved a
 555 lower WER (30.3%) using heavy-duty architec-
 556 tures (Conformer/TDNN) and external Language
 557 Models. However, our framework (33.94%) oper-
 558 ates under strict on-device constraints without ex-
 559 ternal LMs during training, reducing computational
 560 footprint while ensuring adherence to rigid data re-
 561 sidency requirements. Furthermore, the ability to
 562 run these adaptations on mobile GPUs ensures that
 563 hospitals with limited digital infrastructure can still
 564 leverage state-of-the-art ASR foundational models.

565 5.7 Ablation Studies

566 5.7.1 Sensitivity to EWC Regularization (λ)

567 A critical component of our investigation was iden-
 568 tifying the optimal regularization strength for EWC.
 569 We explored a spectrum of λ values to character-
 570 ize the trade-offs between model freezing and forget-
 571 ting. Initial configurations with high regularization
 572 ($\lambda \geq 10^3$) resulted in extreme gradient dominance
 573 by the EWC constraint, effectively “freezing” the
 574 LoRA parameters and preventing any meaningful
 575 adaptation to the noisy clinical domain. Conversely,
 576 we found that a lower value of $\lambda = 10$ (V4.5) pro-
 577 vided the optimal balance, allowing the model to
 578 bridge the reality gap while providing sufficient
 579 parameter protection.

580 5.7.2 LR Warmup and Convergence Speed

581 We evaluated the impact of learning rate warmup
 582 schedules on the stability of continual adaptation.
 583 By comparing a conservative 100-step warmup
 584 against an aggressive 10-step schedule, we ob-
 585 served that the latter significantly accelerated con-
 586 vergence in the early phases of each training seg-
 587 ment. Crucially, this aggressive schedule did not
 588 result in an increase in catastrophic forgetting
 589 ($< 0.05\%$ difference across all runs). This suggests
 590 that in the context of LoRA-based ASR adaptation,
 591 the model can safely utilize high initial learning
 592 rates to rapidly escape local minima from previous
 593 segments without overwriting foundational linguis-
 594 tic features.

595 6 Conclusion

596 In this work, we investigated the “Reality Gap” that
 597 hinders the deployment of state-of-the-art ASR in
 598 rural clinical settings. By quantifying this disparity
 599 at 40.94% WER, we demonstrated the necessity for
 600 localized, privacy-preserving adaptation. Through
 601 a rigorous comparative study of four continual
 602 learning paradigms, we established that a Multi-
 603 Domain Experience Replay strategy provides the
 604 most efficient balance of plasticity and stability.
 605 Our results show a 17.1% relative improvement in
 606 target accuracy and a 55% reduction in catastrophic
 607 forgetting compared to naive baselines. Further-
 608 more, our characterization of the trade-off between
 609 low-rank optimization and acoustic mismatch pro-
 610 vides a technical roadmap for understanding the
 611 interplay between stability and plasticity. These
 612 findings establish a viable blueprint for building
 613 self-improving ASR systems that remain robust and
 614 reliable in high-impact, real-world environments.

615 7 Privacy, Ethics, and Clinical Safety 616 Considerations

617 The deployment of ASR in frontline healthcare set-
 618 tings introduces unique ethical and safety consid-
 619 erations which our framework explicitly addresses.

Data Governance and Residency: Our primary
 620 focus on on-device adaptation is a direct response
 621 to healthcare data residency requirements. By en-
 622 suring that raw patient audio never leaves the local
 623 environment, we mitigate the risk of centralized
 624 data breaches and maintain compliance with emerg-
 625 ing data protection laws in various jurisdictions
 626 (Leroy et al., 2019).
 627

628 **Clinical Safety and Error Risks:** While our
 629 framework significantly reduces WER, automated
 630 transcription in medical contexts is never risk-free.
 631 ASR errors in medically salient keywords (e.g.,
 632 negations or dosage numbers) can lead to serious
 633 diagnostic errors. We advocate for a **Clinician-in-**
 634 **the-Loop** model, where our adaptive system serves
 635 as a draft-generation tool that reduces manual work-
 636 load while maintaining human oversight for final
 637 clinical validation.

638 **Algorithmic Bias and Equity:** By prioritizing
 639 adaptation on rural dialects and low-bandwidth tele-
 640 phony, this work contributes to healthcare equity.
 641 Standard ASR models often fail on marginalized
 642 populations due to accent mismatch; our frame-
 643 work allows localized clinics to "correct" these bi-
 644 ases in real-time by adapting specifically to their
 645 local community's speech patterns.

646 8 Limitations and Future Directions

647 While our framework demonstrates significant po-
 648 tential for clinical deployment, several limitations
 649 remain:

- 650 1. **Metrics as a Clinical Proxy:** Our evaluation
 651 relies on WER and CER, which are standard
 652 in ASR but do not explicitly weight the clin-
 653 ical significance of errors. An error in a pa-
 654 tient's name may be less critical than an error
 655 in a medication dosage; developing medical-
 656 concept-aware metrics remains an essential
 657 future direction.
- 658 2. **Reliance on Clinician Supervision:** Our con-
 659 tinual learning assumption assumes a stream
 660 of corrected transcripts (e.g., from medical
 661 professionals). In under-staffed rural clinics,
 662 this supervision may be sparse or delayed, ne-
 663 cessitating future work into uncertainty-aware
 664 or semi-supervised adaptation.
- 665 3. **Acoustic-Only Training:** While we validated
 666 the necessity of acoustic adaptation, we did
 667 not integrate Language Models (LMs) into
 668 the training loop (e.g., via shallow fusion dur-
 669 ing replay). Capturing the synergy between
 670 acoustic and linguistic adaptation is critical
 671 for handling specialized medical terminology.
- 672 4. **Language and Dialect Scope:** While we fo-
 673 cus on rural Hindi dialects, the efficacy of
 674 our replay strategy on tonal or Dravidian lan-
 675 guages remains to be validated.

References

- 676
- Rahaf Aljundi, Francesca Babiloni, Mohamed Elho-
 677 seiny, Marcus Rohrbach, and Tinne Tuytelaars. 2018.
 678 **Memory aware synapses: Learning what (not) to**
 679 **forget.** In *Computer Vision – ECCV 2018: 15th Eu-*
 680 *ropean Conference, Munich, Germany, September*
 681 *8–14, 2018, Proceedings, Part III*, page 144–161,
 682 Berlin, Heidelberg. Springer-Verlag. 683
- Alexei Baevski, Henry Zhou, Abdelrahman Mohamed,
 684 and Michael Auli. 2020. wav2vec 2.0: a framework
 685 for self-supervised learning of speech representations.
 686 In *Proceedings of the 34th International Conference*
 687 *on Neural Information Processing Systems, NIPS '20*,
 688 Red Hook, NY, USA. Curran Associates Inc. 689
- Frederik Benzing. 2021. **Unifying regularisation meth-**
 690 **ods for continual learning.** 691
- Anish Bhanushali, Grant Bridgman, Deekshitha G,
 692 Prasanta Ghosh, Pratik Kumar, Saurabh Kumar,
 693 Adithya Raj Kolladath, Nithya Ravi, Aaditeshwar
 694 Seth, Ashish Seth, Abhayjeet Singh, Vrunda Sukha-
 695 dia, Umesh S, Sathvik Udupa, and Lodagala V. S.
 696 V. Durga Prasad. 2022. **Gram vaani asr challenge**
 697 **on spontaneous telephone speech recordings in re-**
 698 **gional variations of hindi.** In *Interspeech 2022*, pages
 699 3548–3552. 700
- Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus
 701 Rohrbach, and Mohamed Elhoseiny. 2019a. **Effi-**
 702 **cient lifelong learning with a-GEM.** In *International*
 703 *Conference on Learning Representations.* 704
- Arslan Chaudhry, Marcus Rohrbach, Mohamed El-
 705 hoseiny, Thalaiyasingam Ajanthan, Puneet Kumar
 706 Dokania, Philip H. S. Torr, and Marc'Aurelio Ran-
 707 zato. 2019b. **Continual learning with tiny episodic**
 708 **memories.** *CoRR*, abs/1902.10486. 709
- Chung-Cheng Chiu, Anshuman Tripathi, Katherine
 710 Chou, Chris Co, Navdeep Jaitly, Diana Jaunzeikare,
 711 Anjuli Kannan, Patrick Nguyen, Hasim Sak, Ananth
 712 Sankar, Justin Tansuwan, Nathan Wan, Yonghui Wu,
 713 and Xuedong Zhang. 2018a. **Speech recognition for**
 714 **medical conversations.** In *Interspeech 2018*, pages
 715 2972–2976. 716
- Chung-Cheng Chiu, Anshuman Tripathi, Katherine
 717 Chou, Chris Co, Navdeep Jaitly, Diana Jaunzeikare,
 718 Anjuli Kannan, Patrick Nguyen, Hasim Sak, Ananth
 719 Sankar, Justin Tansuwan, Nathan Wan, Yonghui Wu,
 720 and Xuedong Zhang. 2018b. **Speech recognition for**
 721 **medical conversations.** In *Interspeech 2018*, pages
 722 2972–2976. 723
- Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah
 724 Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh,
 725 and Tinne Tuytelaars. 2022. **A continual learning sur-**
 726 **vey: Defying forgetting in classification tasks.** *IEEE*
 727 *Transactions on Pattern Analysis and Machine Intel-*
 728 *ligence*, 44(7):3366–3385. 729
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-
 730 Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu 731

- Chen. 2022. [LoRA: Low-rank adaptation of large language models](#). In *International Conference on Learning Representations*. 732
- 733
- 734
- 735 Tahir Javed, Kaushal Bhogale, and Mitesh M. Khapra. 2025. [NIRANTAR: Continual Learning with New Languages and Domains on Real-world Speech Data](#). In *Interspeech 2025*, pages 918–922. 736
- 737
- 738
- 739 Tahir Javed, Kaushal Bhogale, Abhigyan Raman, Pratyush Kumar, Anoop Kunchukuttan, and Mitesh Khapra. 2023. [Indicsuperb: A speech processing universal performance benchmark for indian languages](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 37:12942–12950. 740
- 741
- 742
- 743
- 744
- 745 Tahir Javed, Sumanth Doddapaneni, Abhigyan Raman, Kaushal Santosh Bhogale, Gowtham Ramesh, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh M. Khapra. 2022. [Towards building asr systems for the next billion users](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(10):10813–10821. 746
- 747
- 748
- 749
- 750
- 751 David Leroy, Alice Coucke, Thibaut Lavril, Thibault Gisselbrecht, and Joseph Dureau. 2019. [Federated learning for keyword spotting](#). In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6341–6345. 752
- 753
- 754
- 755
- 756
- 757 Xinyu Li, Venkata Chebiyyam, and Katrin Kirchhoff. 2019. [Speech audio super-resolution for speech recognition](#). In *Interspeech 2019*, pages 3416–3420. 758
- 759
- 760 Yin-Tse Lin, Bo-Hao Su, Chi-Han Lin, Shih-Chan Kuo, Jyh-Shing Roger Jang, and Chi-Chun Lee. 2023. [Noise-robust bandwidth expansion for 8k speech recordings](#). In *Interspeech 2023*, pages 5107–5111. 761
- 762
- 763
- 764 Michael McCloskey and Neal J. Cohen. 1989. [Catastrophic interference in connectionist networks: The sequential learning problem](#). volume 24 of *Psychology of Learning and Motivation*, pages 109–165. Academic Press. 765
- 766
- 767
- 768
- 769 Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. [Librispeech: An asr corpus based on public domain audio books](#). In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5206–5210. 770
- 771
- 772
- 773
- 774 Tanvina Patel and Odette Scharenborg. 2022. [Using cross-model learnings for the gram vaani asr challenge 2022](#). In *Interspeech 2022*, pages 4880–4884. 775
- 776
- 777 Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. [Robust speech recognition via large-scale weak supervision](#). In *Proceedings of the 40th International Conference on Machine Learning, ICML’23*. JMLR.org. 778
- 779
- 780
- 781
- 782 Nicola Rieke, Jonny Hancox, Wenqi Li, Fausto Milletari, Holger R. Roth, Shadi Albarqouni, Spyridon Bakas, Mathieu N. Galtier, Bennett A. Landman, Klaus H. Maier-Hein, Sébastien Ourselin, Micah J. Sheller, Ronald M. Summers, Andrew Trask, Daguang Xu, Maximilian Baust, and M. Jorge Cardoso. 2020. [The future of digital health with federated learning](#). *npj Digit. Medicine*, 3. 783
- 784
- 785
- 786
- 787
- 788
- 789
- 790 Zhesu Song, Jianheng Zhuo, Yifan Yang, Ziyang Ma, Shixiong Zhang, and Xie Chen. 2024. [LoRA-Whisper: Parameter-Efficient and Extensible Multilingual ASR](#). In *Interspeech 2024*, pages 3934–3938. 791
- 792
- 793
- 794 Gokul Adethya T and S. Jaya Nirmala. 2025. [A study on regularization-based continual learning methods for indic asr](#). *Preprint*, arXiv:2508.06280. 795
- 796
- 797 Jiannan Xiang, Tianhua Tao, Yi Gu, Tianmin Shu, Zirui Wang, Zichao Yang, and Zhiting Hu. 2023. [Language models meet world models: Embodied experiences enhance language models](#). In *Thirty-seventh Conference on Neural Information Processing Systems*. 798
- 799
- 800
- 801
- 802 Muqiao Yang, Ian Lane, and Shinji Watanabe. 2022. [Online continual learning of end-to-end speech recognition models](#). In *Interspeech 2022*, pages 2668–2672. 803
- 804

805 A Experimental Setup and 806 Reproducibility

807 To ensure the reproducibility of our findings, we
808 provide the full configuration details for all experi-
809 mental paradigms. All models were trained for **3**
810 **epochs per data segment** to ensure local conver-
811 gence.

812 A.1 Hyperparameter Configurations

813 Table 4 details the parameters used, utilizing a split-
814 table format to distinguish between universal set-
815 tings and those tailored to specific continual learn-
816 ing strategies.

817 A.2 Algorithm and Buffer Management

818 The technical contribution of our framework lies in
819 the efficient integration of LoRA, linearized EWC,
820 and prioritized experience replay. Algorithm 1 be-
821 low details the localized adaptation loop, and Al-
822 gorithm 2 describes our multi-domain buffer man-
823 agement strategy.

Algorithm 1: LoRA-based Hybrid Adaptation Loop
Require: Base model θ_{base} , Replay buffer \mathcal{B} , Regular-
 ization λ
1: Initialize LoRA adapters $\theta_0 \subset \theta_{base}$ and importance
 $F \leftarrow 0$
2: for segment $k = 1, \dots, K$ **do**
3: Receive clinical stream \mathcal{D}_k
4: $\mathcal{D}_{train} \leftarrow \text{Sample}(\mathcal{D}_k) \cup \text{Sample}(\mathcal{B})$
5: Adaptation Step:
 $\theta_k \leftarrow \text{argmin}_{\theta} \mathcal{L}_{CTC}(\mathcal{D}_{train})$
 $+ \frac{\lambda}{2} \sum F_i (\theta_k - \theta_{k-1}^*)^2$
6: Importance Estimation:
 $F_{new} \leftarrow \frac{1}{N} \sum_{x \in \mathcal{D}_{train}} |\nabla_{\theta} \mathcal{L}_{CTC}(x)|$
7: Asynchronous Consolidation:
 $F \leftarrow \frac{F \times (k-1) + F_{new}}{k}$
8: $\theta_k^* \leftarrow \text{detach}(\theta_k)$ *Checkpointing*
9: $\mathcal{B} \leftarrow \text{UpdateBuffer}(\mathcal{D}_k, \mathcal{L}_{inst})$
10: end for

Algorithm 2: Buffer Management
Require: Segment data \mathcal{D}_k , Buffer \mathcal{B} , Threshold τ
1: Compute $\mathcal{L}_{inst}(x)$ for all instances $x \in \mathcal{D}_k$
2: Hard Example Mining:
 Identify $\mathcal{S}_{hard} \leftarrow \{x \in \mathcal{D}_k \mid \mathcal{L}_{inst} > \tau \bar{\mathcal{L}}\}$
3: Buffer Update:
 $\mathcal{B}_{gv} \leftarrow \text{Sample}(60\% \text{ from } \mathcal{S}_{hard}$
 $\cup 40\% \text{ Random})$
 $\mathcal{B}_{gen} \leftarrow \text{Sample}(300 \text{ Balanced from Kathbath})$
4: $\mathcal{B} \leftarrow \mathcal{B}_{gv} \cup \mathcal{B}_{gen}$

824 A.3 Computational Efficiency and Hardware 825 Setup

A core goal of our work is to prove that high- 826
 performance, privacy-preserving adaptation is possi- 827
 ble on standard mobile workstations. We evalu- 828
 ated our pipeline on the configurations detailed in 829
 Table 5. These benchmarks demonstrate the techni- 830
 cal feasibility of localized, self-improving ASR 831
 systems in environments where high-end compute 832
 clusters are unavailable. 833

B Detailed Experimental Results 834

This appendix provides the detailed training pro- 835
 gression for all experimental configurations. Each 836
 plot summarizes the recognition performance (WER/CER) 837
 on the target Gram Vaani domain alongside the CTC loss convergence trend. 838
 839

Gradient Stability: A critical challenge encoun- 840
 tered during the development of the EWC-based 841
 paradigms was gradient explosion. As shown later 842
 in Figure 13, our proposed Linearized EWC (L- 843
 EWC) strategy successfully stabilizes the gradient 844
 norm throughout the 24 adaptation segments com- 845
 pared to standard quadratic formulations. 846

Forgetting Analysis: The preservation of foun- 847
 dational general-domain knowledge (Kathbath) is 848
 evaluated across all experimental paradigms, with 849
 detailed comparisons provided in Figures 9 through 850
 12. 851

C Dataset Characteristics 852

The Gram Vaani dataset serves as a rigorous proxy 853
 for rural clinical environments due to its telephonic 854
 acquisition (originally 8kHz upsampled to 16kHz) 855
 and focus on medical/agricultural discussions. Ta- 856
 ble 3 summarizes the key characteristics. 857

Table 3: Characteristics of the partitioned Gram Vaani dataset used for continual adaptation.

Metric	Value
Total Duration	103.2 Hours
Number of Segments (k)	24
Samples per Segment	$\approx 1,600$
Avg. Duration per Sample	9.4 s
Acoustic Condition	Noisy
Primary Dialect	Rural Hindi

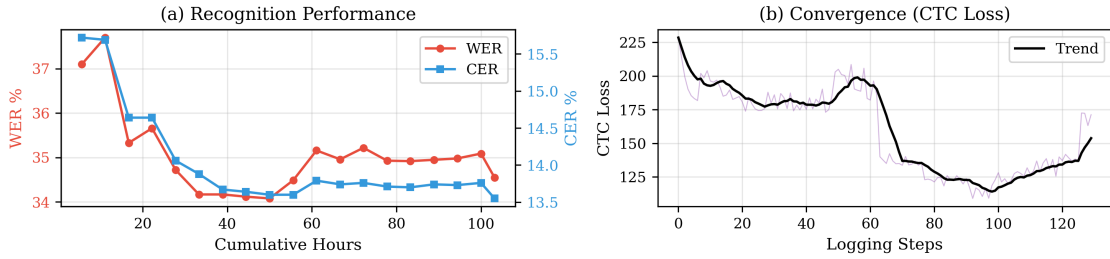
Parameter	Value
Optimizer	AdamW
Base Learning Rate	3×10^{-4}
Weight Decay	0.01
Warmup Steps	10 (Aggressive)
LoRA Target Modules	query, value (Attention blocks)
Batch Size (Effective)	64
Max Audio Duration	30.0 seconds
Training Epochs	3 per segment
Strategy-Specific Settings	
V1.1 Naive	$r = 16, \alpha = 32$
V2.1 Single-Domain ER	$r = 24, \alpha = 48$, Buffer: 400 Target
V3.1 Multi-Domain ER	$r = 24, \alpha = 48$, Buffer: 300 Target + 300 General
V4.5 EWC	$r = 24, \alpha = 48, \lambda = 10$
V5.1 Hybrid	$r = 24, \alpha = 48, \lambda = 100$, Buffer: 300 Target + 300 General

Table 4: Comprehensive hyperparameter settings. r and α represent the LoRA rank and scaling factor, respectively. λ denotes the EWC regularization strength.

Configuration	CPU	GPU	Training Time (per segment)
Config A (High-End)	Intel i7-13700H (14 cores)	NVIDIA RTX 4050 (35W)	25–30 minutes
Config B (Mid-Range)	Intel i5-12500H (12 cores)	NVIDIA RTX 3050 (70W)	50–60 minutes

Table 5: Hardware configurations and training time benchmarks for on-device adaptation.

Detailed Training Progression - V1: Naive (Conservative Warmup)



Detailed Training Progression - V1.1: Naive (Aggressive Warmup)

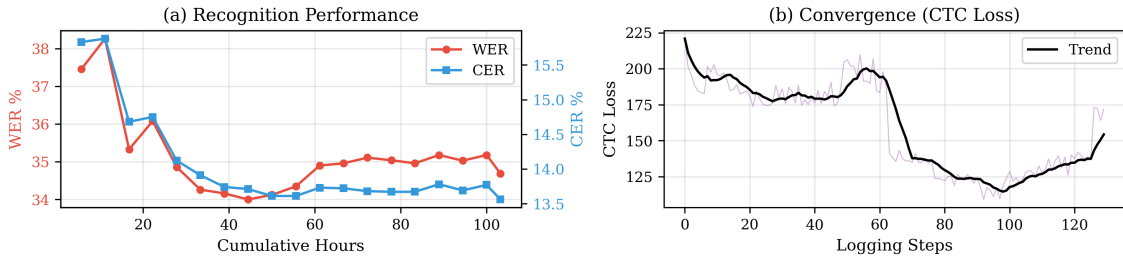
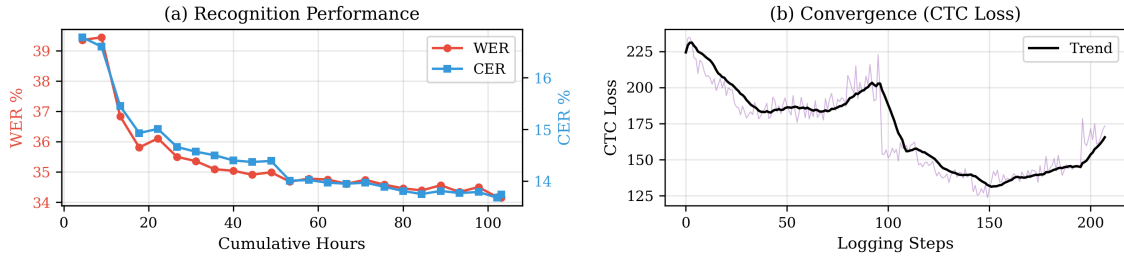


Figure 4: **Naive Baseline Dynamics (V1, V1.1):** Showing rapid initial adaptation but significant volatility, highlighting the need for regularization.

Detailed Training Progression - V2: Single-Domain ER (Conservative)



Detailed Training Progression - V2.1: Single-Domain ER (Aggressive)

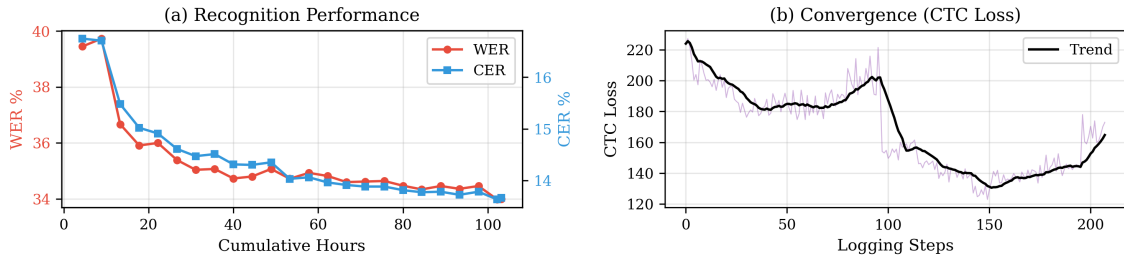
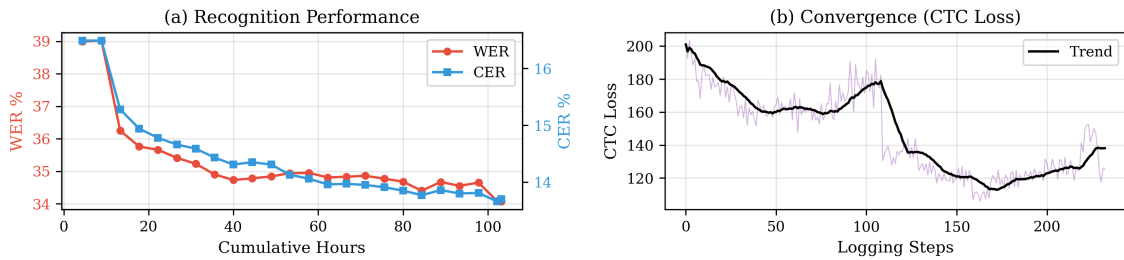


Figure 5: **Single-Domain Replay Dynamics (V2, V2.1):** Incorporating target-domain replay buffers stabilizes the learning trajectory compared to naive fine-tuning.

Detailed Training Progression - V3: Multi-Domain ER (Conservative)



Detailed Training Progression - V3.1: Multi-Domain ER (Aggressive)

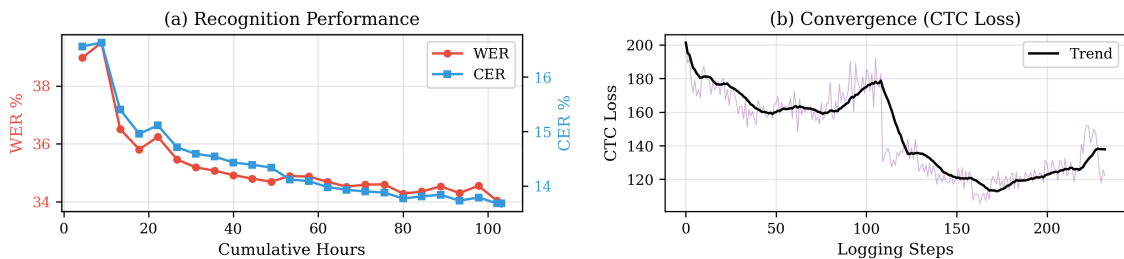
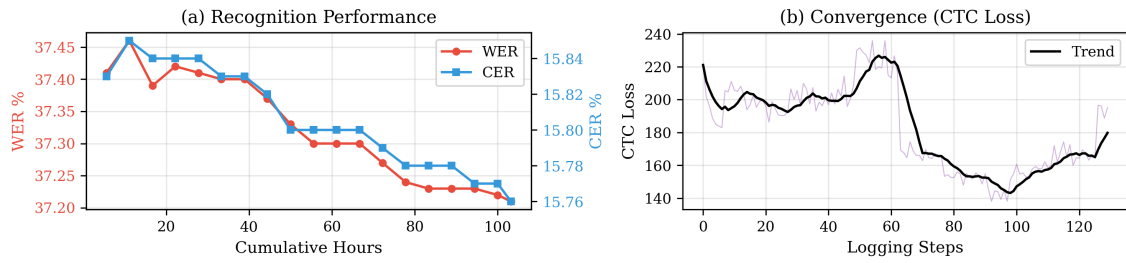


Figure 6: **Multi-Domain Replay Dynamics (V3, V3.1):** Balancing target and general domain samples in the replay buffer offers a tradeoff between plasticity and stability.

Detailed Training Progression - V4.3: EWC ($\lambda=1e3$)



Detailed Training Progression - V4.4: EWC ($\lambda=1e2$)

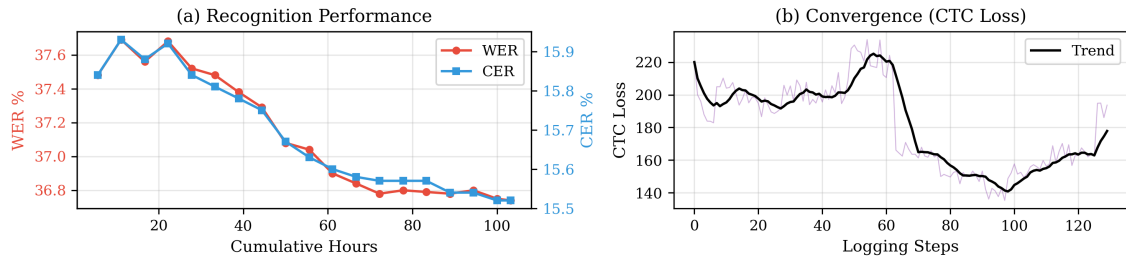
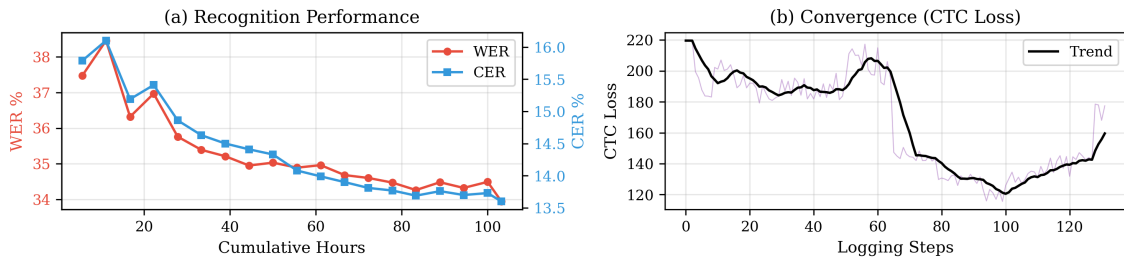
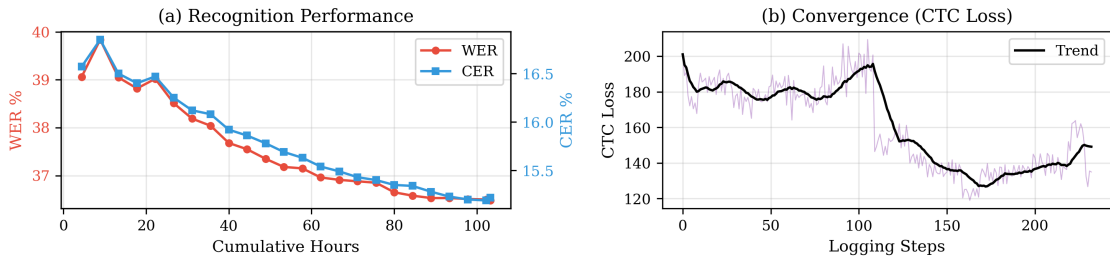


Figure 7: **EWC Regularization (V4.3, V4.4):** High regularization (λ) dampens plasticity, while moderate values allow for adaptation, though standard EWC shows instability.

Detailed Training Progression - V4.5: EWC ($\lambda=1e1$)



Detailed Training Progression - V5: Hybrid ER+EWC ($\lambda=1e2$)



Detailed Training Progression - V5.1: Hybrid ER+EWC ($\lambda=1e1$)

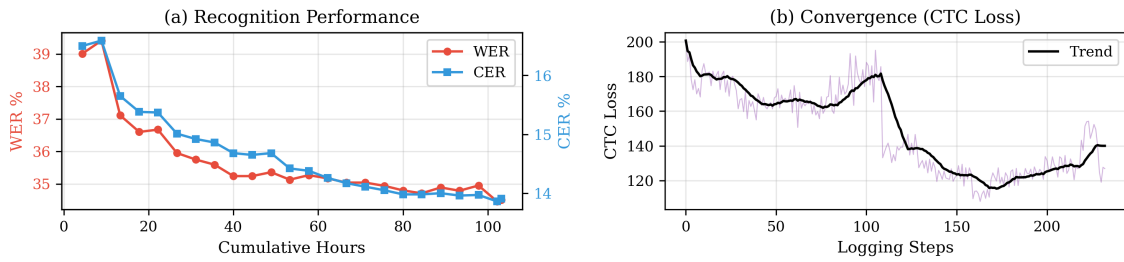


Figure 8: **Optimal Hybrid Strategies (V5, V5.1):** Combining Linearized EWC with Replay (V5.1) achieves the lowest final WER with stable convergence.

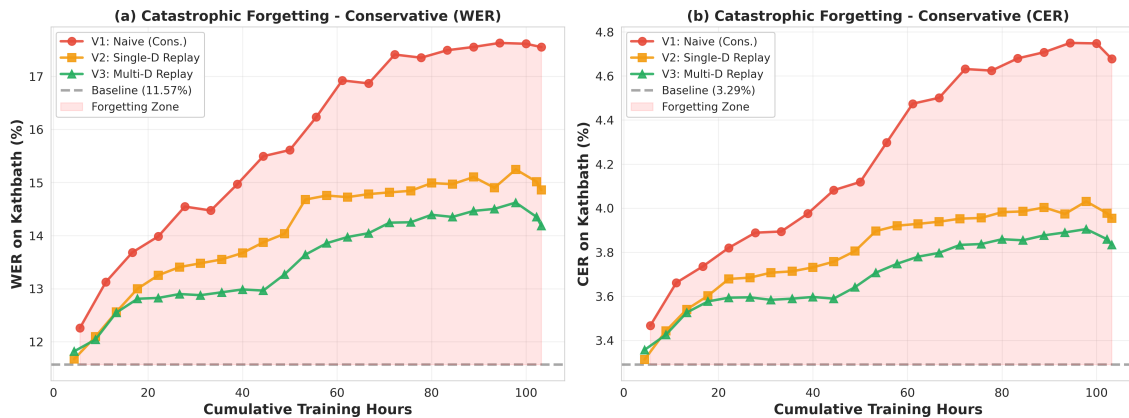


Figure 9: **Forgetting Analysis (Conservative):** Baseline strategies (V1-V3) show varying degrees of knowledge retention on the Kathbath dataset.

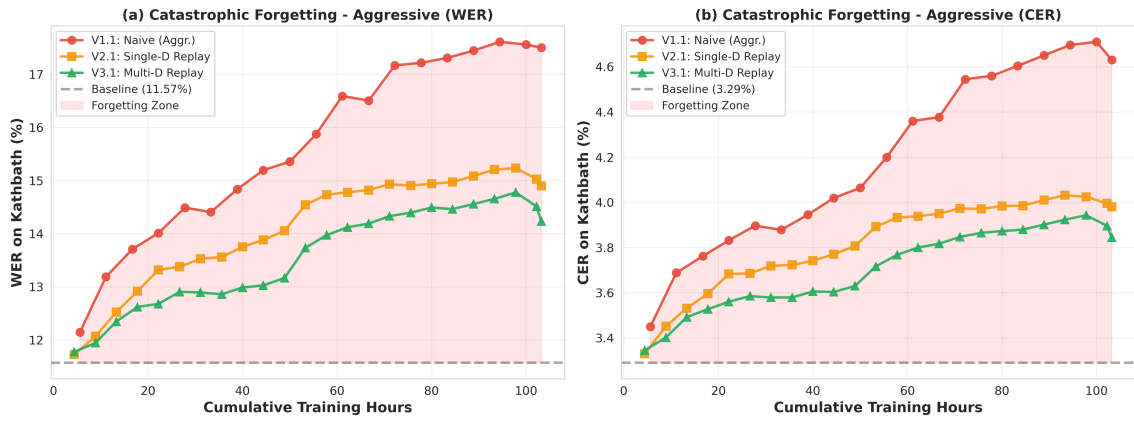


Figure 10: **Forgetting Analysis (Aggressive)**: Increased learning rates in V1.1-V3.1 accelerate adaptation but risk higher catastrophic forgetting.

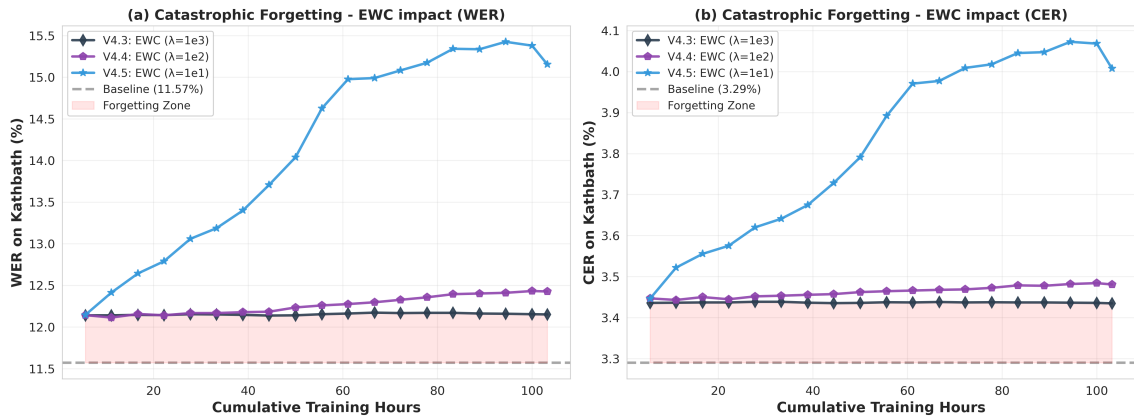


Figure 11: **EWC Forgetting Analysis**: Stronger regularization (λ) effectively reduces forgetting but may hinder adaptation speed.

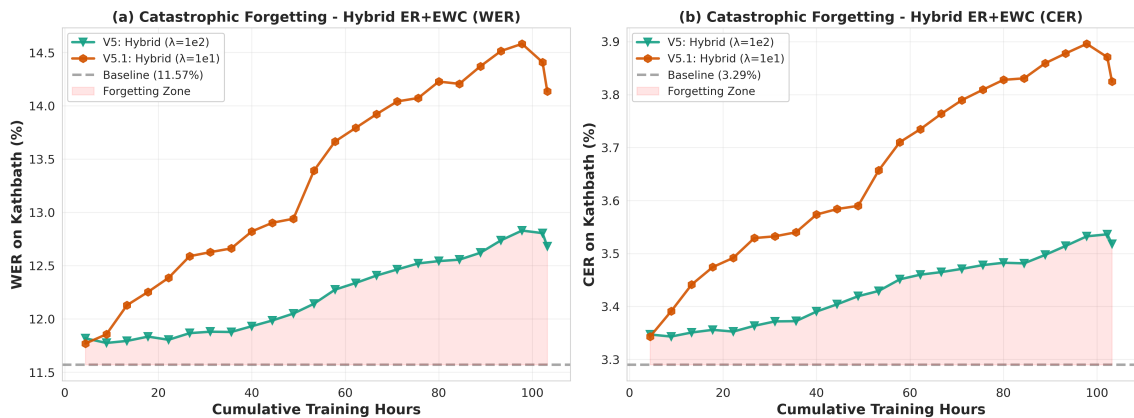


Figure 12: **Hybrid Forgetting Analysis (V5, V5.1)**: The hybrid approach demonstrates the best balance, minimizing forgetting while maintaining plasticity.

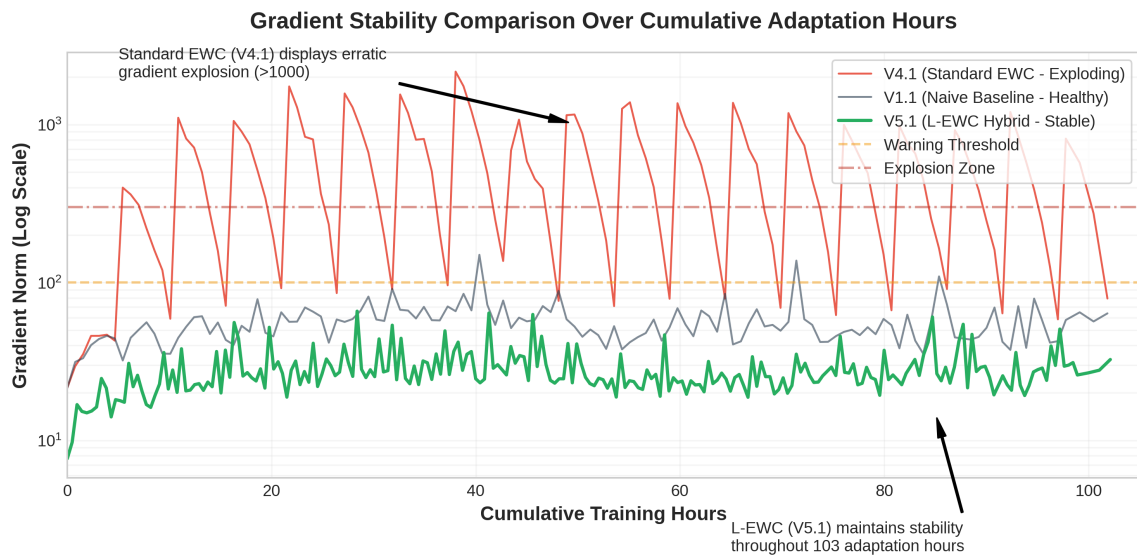


Figure 13: **Gradient Stability Analysis:** Comparing Standard EWC (V4.1), Naive Baseline (V1.1), and our proposed Hybrid L-EWC (V5.1). The log-scale plot demonstrates how the linearized importance estimation in L-EWC prevents the gradient explosion seen in standard quadratic formulations.