

CL-MASR: A Continual Learning Benchmark for Multilingual ASR

Luca Della Libera*, Pooneh Mousavi*, Salah Zaiem, Cem Subakan, Mirco Ravanelli

Abstract—Modern multilingual automatic speech recognition (ASR) systems like Whisper have made it possible to transcribe audio in multiple languages with a single model. However, current state-of-the-art ASR models are typically evaluated on individual languages or in a multi-task setting, overlooking the challenge of continually learning new languages. There is insufficient research on how to add new languages without losing valuable information from previous data. Furthermore, existing continual learning benchmarks focus mostly on vision and language tasks, leaving continual learning for multilingual ASR largely unexplored. To bridge this gap, we propose CL-MASR, a benchmark designed for studying multilingual ASR in a continual learning setting. CL-MASR provides a diverse set of continual learning methods implemented on top of large-scale pretrained ASR models, along with common metrics to assess the effectiveness of learning new languages while addressing the issue of catastrophic forgetting. To the best of our knowledge, CL-MASR is the first continual learning benchmark for the multilingual ASR task. The code is available at <https://github.com/speechbrain/benchmarks>.

Index Terms—Continual learning, multilingual ASR, benchmark.

I. INTRODUCTION

AUTOMATIC speech recognition (ASR) has traditionally focused on converting speech into written text for individual languages. However, with the increasing need for cross-lingual communication and the availability of large-scale multilingual datasets, the focus has recently shifted towards the development of massively multilingual ASR models. These models leverage the well-known advantages of scaling [1] and effectively use similarity across languages to improve performance. Notably, the emergence of models like M-CTC-T [2], Whisper [3], USM [4], and MMS [5] has enabled the automatic transcription of audio of hundreds of languages using a single shared model. Despite these remarkable achievements, with over 7,000 languages existing worldwide [6], the problem of multilingual ASR is far from fully solved. Indeed, due to the dynamic nature of language, even the most advanced ASR system would need regular updates to effectively deal with new dialects and/or domain-specific lexicons, or other variations of existing languages.

The problem of adapting a model over time is addressed by the field of continual learning (CL), also referred to as lifelong

learning or incremental learning. It focuses on designing algorithms that continuously learn from a sequence of tasks.

A naive CL approach is to fine-tune the ASR model on data from new languages as they become available. Unfortunately, this often results in *catastrophic forgetting* [7], [8], a phenomenon that occurs when adjusting the model’s weights based on data from a different distribution. This distribution shift hampers the model’s performance on previously learned tasks. Various methods have been proposed to mitigate catastrophic forgetting in supervised learning, including rehearsal-based [9], [10], [11], [12], architecture-based [13], [14], [15], [16], and regularization-based [17], [18], [19] approaches. While these techniques have been extensively explored in vision and text domains for knowledge transfer across tasks, ASR has received limited attention, especially in multilingual settings.

One of the few works in this specific area is the one by Li et al. [20], who conduct an experimental study to analyze the impact of model’s capacity when incorporating additional languages in a massively multilingual ASR model. However, their investigation solely focuses on the application of incremental fine-tuning, thereby leaving room for exploring a wide range of CL techniques. Additionally, all their models are trained *from scratch*. Research by Ostapenko et al. [21] emphasizes the importance of employing large-scale pretrained models in the fields of computer vision and natural language processing (NLP) for CL. When training a model from scratch for CL, high-level features are often unstable. On the contrary, using large-scale pretrained models can yield to more robust and versatile hidden representations that could better generalize to newly introduced tasks. There are two types of pretrained models that are commonly used in ASR. The first consists of large-scale *supervised* models like Whisper [3], which encompass both an encoder and a decoder. The second includes *self-supervised* models such as wav2vec 2.0 [22], HuBERT [23], and WavLM [24], which are employed as encoder-only models for audio feature extraction. Given the potential differences in training dynamics when applied to CL, it is important to explore the use of both types of models.

To address the lack of research in exploring the potential of large-scale pretrained multilingual ASR models for CL, as well as the scarcity of diverse CL methods specifically designed for multilingual ASR, we introduce Continual Learning for Multilingual ASR (CL-MASR), a benchmark for continual learning applied to multilingual ASR. CL-MASR provides:

- A curated selection of challenging medium/low-resource languages from the Common Voice 13 [25] dataset to utilize for CL experiments in the context of multilingual ASR.

*Equal contribution.

L. Della Libera, P. Mousavi, and M. Ravanelli are with the Gina Cody School of Engineering and Computer Science, Concordia University, Montreal, Canada. (e-mail: luca.dellalibera@mail.concordia.ca; pooneh.mousavi@mail.concordia.ca; mirco.ravanelli@mail.concordia.ca).

S. Zaiem is with LTCI, Télécom Paris, Paris, France. (e-mail: salah.zaiem@telecom-paris.fr).

C. Subakan is with the Department of Computer Science and Software Engineering, Université Laval, Québec, Canada. (e-mail: cem.subakan@ift.ulaval.ca).

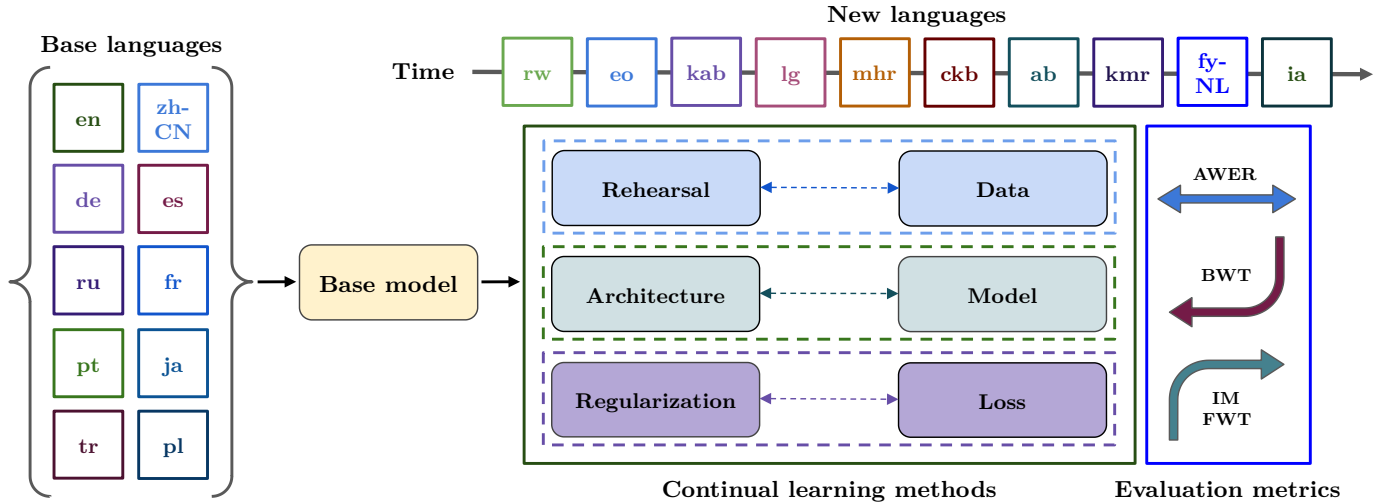


Fig. 1. The workflow of CL-MASR consists of initial pretraining of a base model using a set of base languages. New languages are added incrementally in the same order as depicted in the figure. Rehearsal-based methods operate on the data, architecture-based on the model, and regularization-based on the loss function. Average word error rate (AWER), backward transfer (BWT), intransigence measure (IM), and forward transfer (FWT) are used to assess the performance.

- A diverse set of well-known CL methods implemented on top of Whisper [3] and WavLM [24], along with standard evaluation metrics.
- A modular platform based on the popular SpeechBrain [26] toolkit that enables the easy extension to other CL strategies and/or pretrained models, and facilitates analysis and visualization of the experimental results. Our goal is to provide a friendly interface for CL researchers in order for them to easily test their novel methods on CL-MASR.

Based on the experiments conducted using CL-MASR, we conclude that experience replay [9] is on average one of the best performing CL approaches in the context of multilingual ASR.

II. RELATED WORK

Multilingual ASR. In recent years, we have witnessed impressive progress in the field of multilingual ASR, with the development of massively multilingual models capable of transcribing tens or even hundreds of languages. One such model is M-CTC-T [2], a transformer encoder architecture with a connectionist temporal classification (CTC) [27] head trained using a combination of supervised and semi-supervised techniques on Common Voice 6 [25] and VoxPopuli [28] datasets. Notably, it takes advantage of unlabeled data to improve its performance and supports 60 different languages. Another example is Whisper [3], a transformer-based encoder-decoder model trained in a *weakly* supervised manner using noisy data crawled from the web. It is designed to perform ASR as well as voice activity detection, language identification, speech translation, and timestamp prediction. Trained on 680,000 hours of annotated speech, it can successfully transcribe 99 different languages without requiring any fine-tuning.

Continual Learning in Vision and Text Domains. CL has been extensively explored in both vision and text [29],

[30], [31], [32], [33], [31], [34], [35], [32]. Several datasets and benchmarks exist for studying CL in the visual domain. For example, permuted MNIST [8], CORE50 [36], Split-MNIST [18], and Split-CIFAR [37], [12] employ synthetic data to evaluate various CL methods for image classification tasks. Visual Domain Decathlon [38] and CLEAR [39], on the contrary, propose CL benchmarks based on real-world images. In a more recent study by Bornschein et al. [40], a novel benchmark comprising more than 100 visual classification tasks is introduced. CL is also an active research area in NLP. CLIF [41] and Natural Language Decathlon [42] are benchmarks for CL in NLP tasks including entity typing, sentiment analysis, and natural language inference. Other works [43], [44] explore CL methods for multimodal tasks at the intersection of vision and language such as visual question answering. These studies demonstrate the benefits of CL approaches for learning new tasks while avoiding catastrophic forgetting in diverse domains. Their findings motivate us to investigate the effectiveness of CL applied to multilingual ASR.

Continual Learning in ASR. Most of the research on CL in the ASR domain revolves around domain-incremental learning (DIL) [45], where tasks share the same data label space but come from different distributions. Sadhu et al. [46] propose a method in which the likelihood of a monolingual HMM-DNN ASR model is decomposed into sub-models specific for each data domain. Chang et al. [47] extend a monolingual hybrid CTC-transformer model to new data distributions. Li et al. [48] leverage the self-supervised pretrained BEST-RQ [49] and JUST hydra [50] joint training strategy for fine-tuning using both source and target domain data. These studies focus primarily on examining the effects of CL methods in DIL scenarios, leaving room for further exploration of CL techniques in settings such as multilingual ASR, where the data label space for each task may differ [45]. Relatedly, limited attention has been devoted to CL for multilingual

ASR. Li et al. [20] analyze how the capacity of a massively multilingual ASR model influences its performance when integrating additional languages. Kessler et al. [51] employ a pretrained wav2vec 2.0 [22] model to address the challenge of continually learning new language representations. However, there is a lack of research on the benefits of applying CL methods to expand the capabilities of large pretrained ASR models like Whisper [3]. In this work, our goal is to introduce the first CL benchmark for multilingual ASR that incorporates large pretrained models.

III. PROBLEM FORMULATION

CL focuses on the sequential learning of a model for multiple tasks while preserving the knowledge obtained from previous tasks. The primary assumption is that we have no or limited access to data from previous tasks. The goal is to learn new tasks without experiencing catastrophic forgetting. CL techniques can be categorized into different groups based on their intended use. One such category is task-incremental learning (TIL) [29], which is well-suited for scenarios where tasks have distinct data label spaces and the identity of the task is provided during both training and inference phases. This is the case, for example, of indoor vs outdoor scene classification. On the other hand, blurred boundary continual learning (BBCL) [29] is specifically designed for situations where task boundaries are unclear, and data label spaces are not entirely disjoint. For instance, in reinforcement learning, it is common for the agent to experience a gradual change of tasks as it takes actions. Multilingual ASR lies at the intersection of TIL and BBCL, as some labels (i.e. tokens) are shared across multiple languages, while others are specific to certain languages. This similarity among languages requires the utilization of approaches that can handle both distinct and overlapping data labels effectively.

In this work, we consider the learning of each language as a distinct task. Formally, let θ_0 denote the parameters of a **base model** that can transcribe speech from a set of L_{base} **base languages** and $\mathcal{D}_{\text{new}} = \cup_{i=1}^{L_{\text{new}}} \{(X_i, Y_i)\}$ a dataset of L_{new} **new languages** with X_i representing speech samples from the i -th new language and Y_i the corresponding transcriptions. We aim to incrementally train a sequence of models $\theta_1, \dots, \theta_{L_{\text{new}}}$ on $(X_1, Y_1), \dots, (X_{L_{\text{new}}}, Y_{L_{\text{new}}})$ where the i -th model can successfully transcribe all the languages up to i -th, including the base ones. The main challenge is to prevent each model from forgetting the previously learned languages.

IV. BENCHMARK DESIGN

We introduce CL-MASR, the first benchmark for CL in multilingual ASR. It includes the following four components: a dataset of speech-transcription pairs from multiple languages, a selection of base models to train incrementally, a variety of CL methods along with their implementations, and a standard set of evaluation metrics. The workflow is depicted in Fig. 1. The benchmark platform, based on the popular SpeechBrain [26] toolkit, is available at <https://github.com/speechbrain/benchmarks> and is licensed under Apache 2.0¹.

¹<https://www.apache.org/licenses/LICENSE-2.0>

A. Dataset

Our benchmark builds upon Common Voice 13² [25]. This public dataset, obtained through crowd-sourcing, consists of short audio recordings and their corresponding transcriptions for 108 languages, comprising a total of 17,690 validated hours divided into training, validation, and test splits. We select from it the following two sets of languages ($L_{\text{base}} = L_{\text{new}} = 10$):

- **base languages:** English (en), Chinese (zh-CN), German (de), Spanish (es), Russian (ru), French (fr), Portuguese (pt), Japanese (ja), Turkish (tr), and Polish (pl).
- **new languages:** Kinyarwanda (rw), Esperanto (eo), Kabyle (kab), Luganda (lg), Meadow Mari (mhr), Central Kurdish (ckb), Abkhaz (ab), Kurmanji Kurdish (kmr), Frisian (fy-NL), and Interlingua (ia).

The languages in the first group have a substantial amount of data available on the web and are typically well-supported by multilingual ASR systems [2], [3], [4], [5]. On the other hand, the second group consists of medium/low-resource languages that are often overlooked in ASR research. These languages have limited data available, making them more challenging to work with compared to the base languages. For example, Whisper [52] cannot transcribe any of them. For each language, we randomly extract up to 10 hours of data for training, 1 for validation, and 1 for testing, respectively. This approach not only reduces the computational burden when fine-tuning large multilingual ASR models, but also reflects a more realistic scenario where very limited data are available when incrementally learning new tasks. In order to improve the quality and consistency of the data, we apply some minimal preprocessing steps. First, we filter out utterances longer than 10 seconds, as they are likely to be recordings from open microphones or contain excessive background noise. Then, we discard utterances with transcriptions longer than 200 characters to limit memory usage. Additionally, we perform basic transcript normalization by removing punctuation marks and collapsing repeated spaces into a single one. For detailed information about the data distribution, refer to Table I.

B. Models

There are two main types of pretrained models for multilingual ASR: *supervised* and *self-supervised*. We explore the application of both categories to CL.

As a supervised pretrained ASR model, we employ the large-v2 version of Whisper³ [3]. Based on the encoder-decoder transformer [53] architecture, the encoder extracts audio features from input Mel spectrograms. Conditioned on the encoder’s hidden representations, the autoregressive decoder generates the corresponding transcriptions in a multi-task format that involves the use of special tokens as task specifiers. Since Whisper already supports all the base languages, there is no need for extra fine-tuning and can be directly used in a CL fashion to sequentially add new languages. Note however that CL with Whisper presents additional challenges compared to standard sequence-to-sequence models. First, the

²<https://commonvoice.mozilla.org/en>

³<https://huggingface.co/openai/whisper-large-v2>

TABLE I
LANGUAGES FROM COMMON VOICE 13 [25] USED IN OUR BENCHMARK.
REPORTED VALUES REFER TO THE RAW SUBSAMPLED DATA BEFORE
APPLYING ANY PREPROCESSING.

Language	ISO 639-1	Duration (minutes)		
		Training	Validation	Test
Base languages				
English	en	600	60	60
Chinese	zh-CN	600	60	60
German	de	600	60	60
Spanish	es	600	60	60
Russian	ru	600	60	60
French	fr	600	60	60
Portuguese	pt	600	60	60
Japanese	ja	586	60	60
Turkish	tr	600	60	60
Polish	pl	600	60	60
New languages				
Kinyarwanda	rw	600	60	60
Esperanto	eo	600	60	60
Kabyle	kab	600	60	60
Luganda	lg	600	60	60
Meadow Mari	mhr	600	60	60
Central Kurdish	ckb	484	60	60
Abkhaz	ab	600	60	60
Kurmanji Kurdish	kmr	296	60	60
Frisian	fy-NL	330	60	60
Interlingua	ia	313	60	60

multi-task training format requires predicting a language-specific token at the beginning of each transcript to allow for language identification. Therefore, a new token embedding must be trained for each new language. Second, since Whisper employs a universal byte-level BPE [52] tokenizer, its token space of size 51,865 is shared among all the languages. We observed that, because of this, Whisper is more susceptible to catastrophic forgetting. Finally, given the model size of $\sim 1,550$ M parameters, memory is a limiting factor when applying CL strategies.

As a self-supervised pretrained ASR model, we employ the large version of WavLM⁴ [24] to extract audio features followed by two bidirectional LSTM (BiLSTM) [54] layers with a 1,024-dimensional hidden state and a linear projection to the token space. This architecture is widely recognized and has demonstrated reliable performance in the SUPERB [55] benchmark. First of all, we fit a character-level Sentence-Piece [56] tokenizer on the available transcriptions from both the base and the new languages, resulting in a vocabulary of 4,887 tokens. Then, we jointly fine-tune the model on all the base languages via CTC [27] loss. We train the model for 20 epochs using the AdamW [57] optimizer with an initial learning rate of 0.0001 and a batch size of 8. After each epoch, the learning rate is reduced by 20% if no validation performance improvement is observed. We clip the gradient L_2 norm to 5 to enhance stability. We also freeze the convolutional layers in WavLM’s encoder and enable automatic mixed-precision to reduce memory consumption and speed up training. Greedy decoding is used for inference. The fine-tuned model obtained from this process has a total of ~ 367 M parameters and

serves as the starting point for the CL experiments with self-supervised multilingual ASR.

C. Continual Learning Methods

CL algorithms can be categorized into three main groups: *rehearsal-based*, *architecture-based*, and *regularization-based* approaches. The first involve storing data from previous tasks or utilizing generative models. The second consist in progressively expanding the model with task-specific sub-networks. The third employ regularization techniques to induce knowledge sharing between tasks and prevent forgetting. In addition to naive fine-tuning, which provides a lower bound for the overall performance, we experiment with methods from all these categories. In particular, we implement:

- **Fine-tuning (FT)**: we sequentially fine-tune Whisper and WavLM on the new languages via cross-entropy and CTC [27] loss, respectively. We train the models for 2 epochs per language using the AdamW [57] optimizer with an initial learning rate of 0.0001 and the maximum batch size allowed by our hardware. After each epoch, the learning rate is reduced by 20% if no validation performance improvement is observed. We clip the gradient L_2 norm to 5 to enhance stability. We also freeze Whisper’s encoder and the convolutional layers in WavLM’s encoder and enable automatic mixed-precision to reduce memory consumption and speed up training. Greedy decoding is used for inference. When training Whisper, for each language a new randomly initialized embedding corresponding to the special token for language identification is appended to the decoder’s embedding layer and the entire layer is fine-tuned. Note that the embedding matrix is shared with the final linear projection. When testing Whisper on a given language, following [3], we manually force the correct special token for language identification instead of using the one predicted by the model. Unless otherwise specified, the same training setup is used for all the following methods.
- **Experience replay (ER)** [9]: this rehearsal-based method employs a buffer that retains a portion of data from previous tasks. These stored experiences are then replayed when learning about the current task in order to prevent catastrophic forgetting and promote knowledge transfer. ER can be implemented at either the batch level or the dataset level. In our experiments, we found the latter approach to be more effective. Before training, we mix data from the current task with randomly sampled data from each previous task with a replay ratio of 10%.
- **Averaged gradient episodic memory (A-GEM)** [10]: similarly to ER, this rehearsal-based method utilizes a replay buffer. However, it differs from ER by treating the losses on previous experiences as inequality constraints, avoiding their increase while allowing their decrease. This approach tries to actively prevent catastrophic forgetting while potentially improving the performance on previous tasks. Specifically, A-GEM calculates the average gradient of the replayed samples. If the mean gradient and the current gradient point in the same direction, the

⁴<https://huggingface.co/microsoft/wavlm-large>

current gradient is used to update the model’s parameters. Otherwise, the orthogonal projection to the averaged gradient is used. As in ER, we retain 10% of the data from each previous task.

- **Dark experience replay (DER)** [58]: this rehearsal-based method is an extension of vanilla ER that combines the replay of past experiences with knowledge distillation. In particular, it encourages the current task’s model (i.e. the *student*) to mimic the responses of a previous task’s model (i.e. the *teacher*) on data from the previous task. This is achieved by minimizing the loss on the current task plus the squared L_2 distance between the student’s output logits and the teacher’s output logits, computed on the corresponding previous task’s experiences sampled from a replay buffer. Note that, although we apply this technique in a task-incremental fashion, when paired with reservoir sampling [59], it can effectively deal with the challenging scenario where tasks boundaries are blurred and new data flow in continuously. The trade-off between mimicking the teacher and accommodating new knowledge from the current task is controlled by the hyperparameter α , which determines the relative strength of the regularization term. In our experiments, we set $\alpha = 1$ and, as in ER and A-GEM, we retain 10% of the data from each previous task.
- **Progressive neural networks (PNN)** [13]: this architecture-based method introduces identical sub-networks for each task and allows knowledge transfer among them via lateral adaptor connections. Such a progressive expansion allows the model to learn new tasks while retaining knowledge from previous ones without suffering from catastrophic forgetting. However, PNN requires the task identity to be provided at inference time. In our experiments, we extend Whisper by adding a final task-specific transformer [53] decoder layer and a corresponding embedding layer for each language. Similarly, for WavLM, we use two BiLSTM [54] layers with a 1,024-dimensional hidden state and a corresponding linear projection. All the other parameters stay frozen during training.
- **Piggyback (PB)** [14]: this architecture-based method involves selectively masking the frozen parameters of a base model. This is achieved by maintaining a set of learnable real-valued weights that undergo a deterministic thresholding function, resulting in binary masks. These masks are then applied to the existing parameters of the model. By updating the real-valued weights via gradient descent, the goal is to learn task-specific masks that are well-suited for the given task. Not only is this approach immune to catastrophic forgetting, but it is also agnostic to task ordering and memory-efficient as the masks incur a low overhead of only 1 bit per parameter. However, PB requires the task identity to be provided at inference time. In our experiments, due to memory constraints, we mask only the last two layers of each model. Furthermore, task-specific embedding layer and linear projection are trained for Whisper and WavLM, respectively. The real-valued weights are initialized to 0.01 and the masking threshold

is set to 0.005.

- **Learning to prompt (L2P)** [60]: this method falls under the umbrella of prompt tuning [61] techniques, which have recently emerged with the proliferation of large language models. However, it can be broadly classified as an architecture-based approach, since a prompt can be seen as a simple task-specific adapter network. L2P maintains a pool of *prompts*, i.e. learnable vectors, that are selectively injected into a frozen pretrained model to modify its behavior and adapt it to a new task. Prompts can be chosen either on a per-instance or per-task basis. Furthermore, they can be incorporated at different stages of the architecture such as the encoder’s or decoder’s input layer, the encoder’s last hidden layer, etc. via concatenation, Hadamard product, or other fusion operations. In our experiments, we use the simpler per-task variant. In particular, we maintain a pool of 10 prompts, one for each new language, of shape $d_{\text{model}} \times d_{\text{model}}$, where d_{model} is the dimensionality of the model’s hidden representations, and we post-multiply the encoder’s output by the prompt of the corresponding language. Similarly to PNN and PB, the task identity needs to be provided at inference time.
- **Elastic weight consolidation (EWC)** [17]: this regularization-based method estimates the importance of each parameter based on its contribution to the performance on previous tasks. It does so by computing the Fisher information matrix (FIM), which measures the sensitivity of the loss function with respect to the model’s parameters. Based on this sensitivity, EWC calculates a quadratic penalty term that is added to the loss function while training on the current task. This term encourages the model to preserve the important parameters for previous tasks while adapting to the new task. Such a trade-off is explicitly controlled by the hyperparameter λ , which determines the relative strength of the regularization term. Due to memory constraints, we resort to an online version of EWC and include the penalty term only for the last task. To mitigate the effect of this approximation, we introduce an hyperparameter, α , to control the influence of previous tasks in updating the parameter sensitivity. In our experiments, we set $\lambda = 5$ and $\alpha = 0.5$. Also note that the original training data for Whisper are not publicly available. Hence, we estimate the initial FIM using the training data from the base languages.
- **Learning without forgetting (LwF)** [19]: this regularization-based method involves using a frozen copy of the model trained up to the previous task as a *teacher* and the current model as a *student*. The goal is to enable the student to learn not only from the current task but also from the teacher. To achieve this, LwF utilizes knowledge distillation [62]. First, the output probabilities of the teacher are computed for each sample in the current task. Then, the student is trained to minimize the loss on the current task plus the cross-entropy between teacher’s and student’s output probabilities, smoothed by a temperature parameter T that increases the weight for smaller probabilities. The trade-off between mimicking the teacher and accommodating new knowledge from the

current task is controlled by the hyperparameter λ , which determines the relative strength of the regularization term. In our experiments, we set $T = 2$ and $\lambda = 10$.

- **Memory Aware Synapses (MAS)** [63]: similarly to EWC, this regularization-based method mitigates catastrophic forgetting by penalizing large updates to parameters that contribute the most to the performance on previous tasks. Differently from EWC, it estimates parameter relevance as the average magnitude of the gradients of the squared L_2 norm of the learned function. Not only is this definition of importance agnostic to the loss function, but it can also be accumulated over tasks in an online manner. The trade-off between forgetting and learning is explicitly controlled by the hyperparameters λ and α , which determine the relative strength of the regularization term and the influence of previous tasks in updating the parameter importance, respectively. In our experiments, we set $\lambda = 1$ and $\alpha = 0.5$. Furthermore, as for EWC, we estimate the initial parameter importance using the training data from the base languages.

D. Evaluation Metrics

When evaluating CL methods, a primary consideration is to assess the *overall performance* across all learned tasks. To do so, we employ a variation of the average accuracy [64], [65], [19], the **average word error rate (AWER)**⁵, calculated incrementally after each newly introduced task:

$$\text{AWER}_t = \frac{1}{t} \sum_{i=1}^t \text{WER}_{t,i} \quad t = 1, \dots, T, \quad (1)$$

where T denotes the total number of tasks and $\text{WER}_{t,i}$ denotes the word error rate on the i -th task after the model finishes the learning on the t -th task. In particular, in our benchmark, we consider the 10 base languages as a single joint task. Hence, since we have 10 new languages, we set $T = 1 + 10 = 11$. Here, $\text{WER}_{t,1}$ represents the word error rate averaged over the 10 base languages after learning about the t -th task. Note that $\text{AWER}_t \in [0, \infty)$, with smaller values indicating better overall performance. It provides a comprehensive understanding of the model’s ability to retain and utilize knowledge from previously learned tasks while accommodating new information.

Another key aspect is memory *stability*, or robustness against *forgetting*, i.e. the impact of learning new tasks on the performance of previously learned tasks. More specifically, we aim to understand whether learning new tasks has any detrimental effect on the model’s performance on previously learned ones. To this end, we measure the **backward transfer (BWT)** [65], defined as

$$\text{BWT}_t = \frac{1}{t-1} \sum_{i=1}^{t-1} \text{WER}_{i,i} - \text{WER}_{t,i} \quad t = 2, \dots, T, \quad (2)$$

i.e. the average performance gain on previously learned tasks. Note that $\text{BWT}_t \in (-\infty, \infty)$ with positive values indicating

improvement on the previously learned tasks and negative ones indicating forgetting. However, $\text{BWT}_t \leq 0$ in most situations.

Learning *plasticity* is also a crucial factor to consider. It refers to the model’s capacity to effectively acquire new knowledge. Emphasizing too much on plasticity may lead to catastrophic forgetting, compromising the model’s ability to retain previously learned tasks. On the other hand, excessively prioritizing stability may hinder the model’s adaptability to new tasks. This is often referred to in the literature as the *stability-plasticity dilemma* [66]. We quantify plasticity via the **intransience measure (IM)** [64], defined as

$$\text{IM}_t = \text{WER}_{t,t} - \text{WER}_t^{\text{joint}} \quad t = 2, \dots, T, \quad (3)$$

where the reference value $\text{WER}_t^{\text{joint}}$ denotes the word error rate on the t -th task of the model jointly trained on all tasks (i.e. base + new languages) at the same time. Note that $\text{IM}_t \in (-\infty, \infty)$ with larger values corresponding to the inability of the model to effectively learn new tasks. Also note that the choice of the reference value is arbitrary, however joint training is the most natural option as it provides an upper bound for the overall performance.

Lastly, we are interested in examining the *influence* of previous tasks on learning the current task. While this concept is partly related to plasticity, it is worth noting that a model with large plasticity does not necessarily utilize knowledge from previous tasks to enhance performance on the current task. To capture this aspect, we propose a variation of **forward transfer (FWT)** [65], that we define as

$$\text{FWT}_t = \text{WER}_t^{\text{fine-tuned}} - \text{WER}_{t,t} \quad t = 2, \dots, T, \quad (4)$$

where $\text{WER}_t^{\text{fine-tuned}}$ denotes the word error rate on the t -th task of the model fine-tuned solely on that specific task. Note that $\text{FWT}_t \in (-\infty, \infty)$ with larger values corresponding to a stronger ability of the model to exploit knowledge from previous tasks. Also note that our definition significantly deviates from the one of Lopez-Paz et al. [65], as they interpret forward transfer as the zero-shot improvement in accuracy on future tasks compared to random guessing. However, in our setting, measuring performance on unknown tasks is not meaningful, as zero-shot transfer is difficult for multilingual ASR.

V. EXPERIMENTS

To show the utility of CL-MASR, we conduct a comparative study to determine the most effective combinations of CL methods and base models for multilingual ASR. Additionally, we analyze the impact of language ordering. All the experiments were run on 5 CentOS Linux machines with an Intel(R) Xeon(R) Silver 4216 Cascade Lake CPU with 32 cores @ 2.10 GHz, 64 GB RAM and an NVIDIA Tesla V100 SXM2 @ 32 GB with CUDA Toolkit 11.4. With the specified hardware configuration, approximately 10 days are necessary to complete all the experiments. For detailed information about the hyperparameters used in each experiment, refer to the official code repository.

⁵For Chinese and Japanese the concept of word is not well-defined, hence, following [3], we consider the character error rate instead of the word error rate.

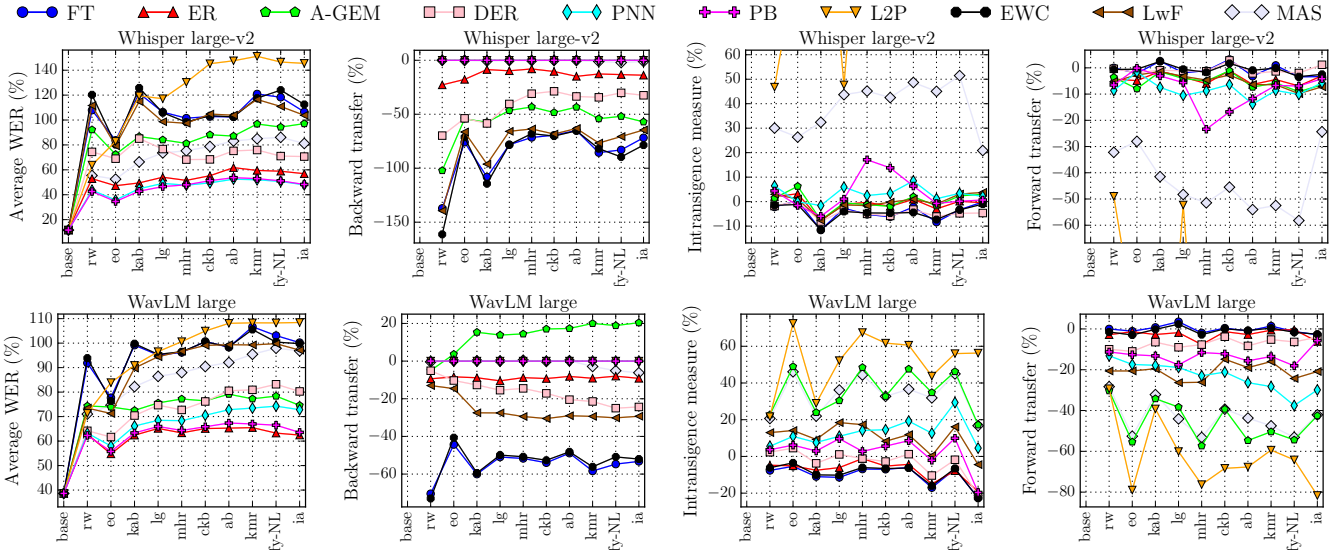


Fig. 2. Comparison of fine-tuning (FT), experience replay (ER), averaged gradient episodic memory (A-GEM), dark experience replay (DER), progressive neural networks (PNN), Piggyback (PB), learning to prompt (L2P), elastic weight consolidation (EWC), learning without forgetting (LwF), and memory aware synapses (MAS) applied to Whisper large-v2 and WavLM large on the base and new languages on Common Voice dataset. Note that we crop the top right figures as L2P is off-scale. ER is among the best performing methods together with PNN and PB, with the additional advantage of being task-agnostic.

A. Comparison of Continual Learning Methods

Fig. 2 shows the main results of our comparative study. First of all, we observe that, when applying naive FT, the AWER is $\sim 100\%$ after learning each new language, which means that the model can hardly retain any knowledge about previous tasks if no CL intervention is made. This highlights the challenging nature of our benchmark, which makes it well-suited for the development of robust CL techniques. Among the evaluated CL methods, two architecture-based approaches, namely PNN and PB, obtain the smallest AWER. Notably, they perform well despite the addition of *small* adapter layers. This emphasizes the benefits of using large-scale pretrained models in CL. However, L2P, which is also architecture-based, performs poorly. We hypothesize this is due to the complex nature of the sequence-to-sequence mapping task, which is significantly more challenging than the original task L2P was designed for, i.e. image classification. While ER performs comparably to PNN and PB, we would like to emphasize that it has the advantage of being task-agnostic at inference time, making it a more robust choice for multilingual ASR. DER is also competitive to some degree, but the trade-off between forgetting and learning is too biased towards the latter, leading to an overall performance that is inferior to ER. Finally, the regularization-based methods, namely EWC, LwF, and MAS, exhibit larger AWER, in most cases similar to naive FT. Interestingly, MAS excels at mitigating forgetting but struggles with learning, possibly due to an overly strong regularization effect introduced by the parameter importance penalty.

Regarding the stability-plasticity trade-off, most methods are characterized by $BWT \leq 0$ and $FWT < 0$. In particular, architecture-based strategies are immune to forgetting by design. However, this hinders their ability to learn new tasks, resulting in poor IM and FWT. On the other hand,

FT demonstrates superior adaptability to new tasks but falls short in mitigating forgetting. These observations generally hold true for both base models. The only exception is A-GEM, which yields to $BWT > 0$ when applied to WavLM. Although not surprising (A-GEM explicitly aims to improve the model’s performance on previous tasks), this effect is not observed for Whisper. A possible explanation is that Whisper’s tokens are shared among languages, potentially introducing conflicts when optimizing A-GEM’s loss function. For additional experiments and more extensive results, refer to Appendix B and Appendix D, respectively.

B. Comparison of Base Models

Based on Fig. 2, we also draw a comparison between Whisper and WavLM. We observe that Whisper tends to outperform WavLM in terms of AWER, especially on the base languages. This is in line with expectations, as it was pretrained on a vast and diverse dataset encompassing 99 languages, unlike WavLM, which was pretrained on the 10 base languages only. However, it is important to note that different stability-plasticity trade-offs are achieved by the two models. Whisper exhibits better capacity in learning new languages, as indicated by smaller IM and larger FWT across all methods. On the other hand, WavLM demonstrates superior performance in terms of BWT for both regularization-based and rehearsal-based methods, especially in mitigating forgetting for the base languages. This difference can be attributed to several factors, including the pretraining strategies (supervised versus self-supervised) and, as previously noted, the token space (byte-level BPEs versus characters). In particular, we suspect that the type of tokens used in each model has an impact on the balance between learning and forgetting. To gain a deeper understanding and analyze these effects more comprehensively, further investigation is necessary. Another interesting

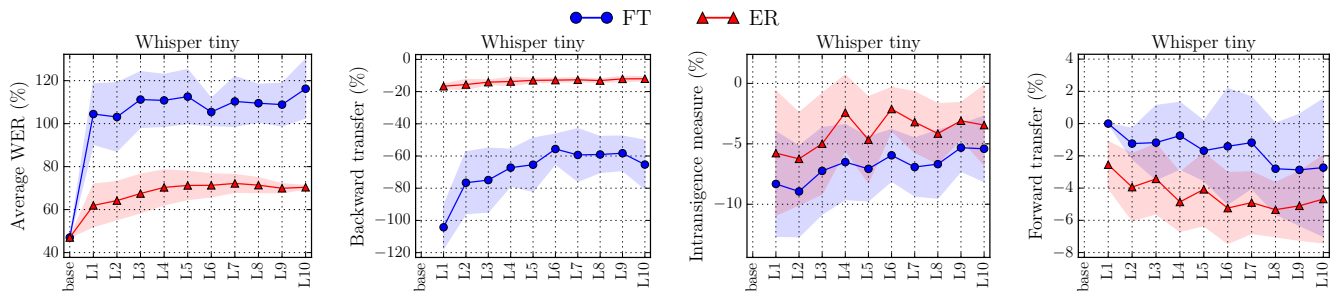


Fig. 3. Comparison of fine-tuning (FT) and experience replay (ER) applied to Whisper tiny. Curves represent the mean ± 1 standard deviation over 10 random language (L) orderings. ER is effective in mitigating forgetting and reducing ordering sensitivity.

aspect is the steep increase in AWER and forgetting for both Whisper and WavLM after learning the first new language. Our analysis suggests that this phenomenon is primarily associated with imbalance in terms of amount of data and/or number of training steps between the initial supervised joint pretraining and the subsequent incremental training. For more details, refer to Appendix C.

C. Impact of Language Ordering

Performance variability due to task ordering is a primary concern in CL studies. In Fig. 3, we examine the impact of language ordering. Due to resource constraints, we limit our analysis to the tiny version of Whisper⁶ [3] and we consider only FT and ER methods. We observe that FT exhibits a significant fluctuation of 10-20% in both AWER and BWT, depending on the language sequence. On the other hand, ER is effective in mitigating forgetting and reducing ordering sensitivity: as more languages are added, the variance of AWER decreases, since the ordering becomes less important when a fraction of data from all previous tasks is accessible. However, the ordering of languages still introduces substantial variance with respect to FWT and IM. This highlights the importance of task ordering when transferring knowledge from previously learned languages to a new one.

VI. CONCLUSION

We introduce CL-MASR, the first benchmark for continual learning in multilingual ASR. CL-MASR offers a curated selection of medium/low-resource languages, a modular and flexible platform for executing and evaluating various CL methods on top of existing large-scale pretrained multilingual ASR models, and a standardized set of evaluation metrics. Through a comparative study, we show that experience replay is one of the most effective strategies for combating catastrophic forgetting in multilingual ASR. By releasing the code, we hope that CL-MASR will promote further research in the field and serve as a valuable real-world testbed for novel CL algorithms.

As a future work, we are planning to further explore other categories of CL such as continual self-supervised representation learning [67], [68], [69]. Additionally, deploying large-scale pretrained models poses challenges in practical

applications due to the high computational burden. Using knowledge distillation [70], [71] to transfer knowledge from larger teacher models to smaller and more memory-efficient student models as an alternative for base models could open up interesting directions for future work.

ACKNOWLEDGEMENTS

We thank Reza Davari for valuable discussions. We acknowledge the support of the Natural Sciences and Engineering Research Council of Canada (NSERC) and the Digital Research Alliance of Canada (alliancecan.ca).

REFERENCES

- [1] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei, "Scaling laws for neural language models," *arXiv preprint arXiv:2001.08361*, 2020.
- [2] L. Lugosch, T. Likhomanenko, G. Synnaeve, and R. Collobert, "Pseudo-labeling for massively multilingual speech recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 7687–7691.
- [3] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," *arXiv preprint arXiv:2212.04356*, 2022.
- [4] Y. Zhang, W. Han, J. Qin, Y. Wang, A. Bapna, Z. Chen, N. Chen, B. Li, V. Axelrod, G. Wang, Z. Meng, K. Hu, A. Rosenberg, R. Prabhavalkar, D. S. Park, P. Haghani, J. Riesa, G. Perng, H. Soltau, T. Strohmaier, B. Ramabhadran, T. Sainath, P. Moreno, C.-C. Chiu, J. Schalkwyk, F. Beaufays, and Y. Wu, "Google USM: Scaling automatic speech recognition beyond 100 languages," *arXiv preprint arXiv:2303.01037*, 2023.
- [5] V. Pratap, A. Tjandra, B. Shi, P. Tomasello, A. Babu, S. Kundu, A. Elkahky, Z. Ni, A. Vyas, M. Fazel-Zarandi, A. Baevski, Y. Adi, X. Zhang, W.-N. Hsu, A. Conneau, and M. Auli, "Scaling speech technology to 1,000+ languages," *arXiv preprint arXiv:2305.13516*, 2023.
- [6] M. P. Lewis, G. F. Simon, and C. D. Fennig, "Ethnologue: Languages of the world, nineteenth edition," Online version: <https://www.ethnologue.com>, 2016.
- [7] R. M. French, "Catastrophic forgetting in connectionist networks," *Trends in Cognitive Sciences*, vol. 3, no. 4, pp. 128–135, 1999.
- [8] I. J. Goodfellow, M. Mirza, D. Xiao, A. Courville, and Y. Bengio, "An empirical investigation of catastrophic forgetting in gradient-based neural networks," in *International Conference on Learning Representations (ICLR)*, 2014.
- [9] D. Rolnick, A. Ahuja, J. Schwarz, T. P. Lillicrap, and G. Wayne, "Experience replay for continual learning," *International Conference on Neural Information Processing Systems (NeurIPS)*, vol. 32, pp. 350–360, 2019.
- [10] A. Chaudhry, M. Ranzato, M. Rohrbach, and M. Elhoseiny, "Efficient lifelong learning with A-GEM," in *International Conference on Learning Representations (ICLR)*, 2019.
- [11] H. Shin, J. K. Lee, J. Kim, and J. Kim, "Continual learning with deep generative replay," in *International Conference on Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017, pp. 2994–3003.

⁶<https://huggingface.co/openai/whisper-tiny>

- [12] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, “iCaRL: Incremental classifier and representation learning,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5533–5542.
- [13] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu, R. Pascanu, and R. Hadsell, “Progressive neural networks,” *arXiv preprint arXiv:1606.04671*, 2016.
- [14] A. Mallya, D. Davis, and S. Lazebnik, “Piggyback: Adapting a single network to multiple tasks by learning to mask weights,” in *European Conference on Computer Vision (ECCV)*, 2018, pp. 72–88.
- [15] R. Aljundi, P. Chakravarty, and T. Tuytelaars, “Expert gate: Lifelong learning with a network of experts,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 7120–7129.
- [16] A. Mallya and S. Lazebnik, “PackNet: Adding multiple tasks to a single network by iterative pruning,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7765–7773.
- [17] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell, “Overcoming catastrophic forgetting in neural networks,” *Proceedings of the National Academy of Sciences*, vol. 114, no. 13, pp. 3521–3526, 2017.
- [18] F. Zenke, B. Poole, and S. Ganguli, “Continual learning through synaptic intelligence,” in *International Conference on Machine Learning (ICML)*, 2017, pp. 3987–3995.
- [19] S. Hou, X. Pan, C. C. Loy, Z. Wang, and D. Lin, “Learning a unified classifier incrementally via rebalancing,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 831–839.
- [20] B. Li, R. Pang, Y. Zhang, T. N. Sainath, T. Strohmaier, P. Haghani, Y. Zhu, B. Farris, N. Gaur, and M. Prasad, “Massively multilingual ASR: A lifelong learning solution,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 6397–6401.
- [21] O. Ostapenko, T. Lesort, P. Rodriguez, M. R. Arefin, A. Douillard, I. Rish, and L. Charlin, “Continual learning with foundation models: An empirical study of latent replay,” in *Conference on Lifelong Learning Agents*, vol. 199, 2022, pp. 60–91.
- [22] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” in *International Conference on Neural Information Processing Systems (NeurIPS)*, vol. 33, 2020, pp. 12 449–12 460.
- [23] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhota, R. Salakhutdinov, and A. Mohamed, “HuBERT: Self-supervised speech representation learning by masked prediction of hidden units,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.
- [24] S. Chen, C. Wang, Z. Chen, Y. Wu, S. Liu, Z. Chen, J. Li, N. Kanda, T. Yoshioka, X. Xiao, J. Wu, L. Zhou, S. Ren, Y. Qian, Y. Qian, J. Wu, M. Zeng, X. Yu, and F. Wei, “WavLM: Large-scale self-supervised pre-training for full stack speech processing,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 6, pp. 1505–1518, 2022.
- [25] R. Ardila, M. Branson, K. Davis, M. Kohler, J. Meyer, M. Henretty, R. Morais, L. Saunders, F. Tyers, and G. Weber, “Common Voice: A massively-multilingual speech corpus,” in *Twelfth Language Resources and Evaluation Conference*, 2020, pp. 4218–4222.
- [26] M. Ravanelli, T. Parcollet, P. Plantinga, A. Rouhe, S. Cornell, L. Lugosch, C. Subakan, N. Dawalatabad, A. Heba, J. Zhong, J.-C. Chou, S.-L. Yeh, S.-W. Fu, C.-F. Liao, E. Rastorgueva, F. Grondin, W. Aris, H. Na, Y. Gao, R. D. Mori, and Y. Bengio, “SpeechBrain: A general-purpose speech toolkit,” *arXiv preprint arXiv:2106.04624*, 2021.
- [27] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *International Conference on Machine Learning (ICML)*, 2006, pp. 369–376.
- [28] C. Wang, M. Riviere, A. Lee, A. Wu, C. Talnikar, D. Haziza, M. Williamson, J. Pino, and E. Dupoux, “VoxPopuli: A large-scale multilingual speech corpus for representation learning, semi-supervised learning and interpretation,” in *Annual Meeting of the Association for Computational Linguistics and the International Joint Conference on Natural Language Processing*, 2021, pp. 993–1003.
- [29] L. Wang, X. Zhang, H. Su, and J. Zhu, “A comprehensive survey of continual learning: Theory, method and application,” *arXiv preprint arXiv:2302.00487*, 2023.
- [30] M. De Lange, R. Aljundi, M. Masana, S. Parisot, X. Jia, A. Leonardis, G. Slabaugh, and T. Tuytelaars, “A continual learning survey: Defying forgetting in classification tasks,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 7, pp. 3366–3385, 2021.
- [31] Z. Ke and B. Liu, “Continual learning of natural language processing tasks: A survey,” *arXiv preprint arXiv:2211.12701*, 2022.
- [32] Z. Mai, R. Li, J. Jeong, D. Quispe, H. Kim, and S. Sanner, “Online continual learning in image classification: An empirical survey,” *Neuro-computing*, vol. 469, pp. 28–51, 2022.
- [33] A. Antoniou, M. Patacchiola, M. Ochal, and A. Storkey, “Defining benchmarks for continual few-shot learning,” *arXiv preprint arXiv:2004.11967*, 2020.
- [34] V. Lomonaco, L. Pellegrini, A. Cossu, A. Carta, G. Graffieti, T. L. Hayes, M. De Lange, M. Masana, J. Pomponi, G. M. van de Ven, M. Mundt, Q. She, K. Cooper, J. Forest, E. Belouadah, S. Calderara, G. I. Parisi, F. Cuzzolin, A. S. Toliás, S. Scardapane, L. Antiga, S. Ahmad, A. Popescu, C. Kanan, J. van de Weijer, T. Tuytelaars, D. Bacciu, and D. Maltoni, “Avalanche: an end-to-end library for continual learning,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2021, pp. 3595–3605.
- [35] M. Masana, X. Liu, B. Twardowski, M. Menta, A. D. Bagdanov, and J. van de Weijer, “Class-incremental learning: survey and performance evaluation on image classification,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [36] V. Lomonaco and D. Maltoni, “CORe50: a new dataset and benchmark for continuous object recognition,” in *Conference on Robot Learning*, 2017, pp. 17–26.
- [37] A. Krizhevsky and G. Hinton, “Learning multiple layers of features from tiny images,” Technical report, 2009.
- [38] S.-A. Rebuffi, H. Bilen, and A. Vedaldi, “Learning multiple visual domains with residual adapters,” *International Conference on Neural Information Processing Systems (NeurIPS)*, vol. 30, p. 506–516, 2017.
- [39] Z. Lin, J. Shi, D. Pathak, and D. Ramanan, “The CLEAR benchmark: Continual learning on real-world imagery,” in *International Conference on Neural Information Processing Systems (NeurIPS): Datasets and Benchmarks Track*, 2021.
- [40] J. Bornschein, A. Galashov, R. Hemsley, A. Rannen-Triki, Y. Chen, A. Chaudhry, X. O. He, A. Douillard, M. Caccia, Q. Feng, J. Shen, S.-A. Rebuffi, K. Stacpoole, D. de las Casas, W. Hawkins, A. Lazaridou, Y. W. Teh, A. A. Rusu, R. Pascanu, and M. Ranzato, “NEVIS’22: A stream of 100 tasks sampled from 30 years of computer vision research,” *arXiv preprint arXiv:2211.11747*, 2022.
- [41] X. Jin, B. Y. Lin, M. Rostami, and X. Ren, “Learn continually, generalize rapidly: lifelong knowledge accumulation for few-shot learning,” in *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2021, pp. 714–729.
- [42] B. McCann, N. S. Keskar, C. Xiong, and R. Socher, “The Natural Language Decathlon: Multitask learning as question answering,” *arXiv preprint arXiv:1806.08730*, 2018.
- [43] C. Greco, B. Plank, R. Fernández, and R. Bernardi, “Measuring catastrophic forgetting in visual question answering,” in *Increasing Naturalness and Flexibility in Spoken Dialogue Interaction: 10th International Workshop on Spoken Dialogue Systems*. Springer, 2021, pp. 381–387.
- [44] T. Srinivasan, T.-Y. Chang, L. Pinto Alva, G. Chochlakakis, M. Rostami, and J. Thomason, “CLiMB: A continual learning benchmark for vision-and-language tasks,” in *International Conference on Neural Information Processing Systems (NeurIPS)*, vol. 35, 2022, pp. 29 440–29 453.
- [45] G. M. van de Ven, T. Tuytelaars, and A. S. Toliás, “Three types of incremental learning,” *Nature Machine Intelligence*, pp. 1–13, 2022.
- [46] S. Sadhu and H. Hermansky, “Continual learning in automatic speech recognition,” in *Interspeech*, 2020, pp. 1246–1250.
- [47] H.-J. Chang, H.-y. Lee, and L.-s. Lee, “Towards lifelong learning of end-to-end ASR,” in *Interspeech*, 2021, pp. 2551–2555.
- [48] B. Li, D. Hwang, Z. Huo, J. Bai, G. Prakash, T. N. Sainath, K. Chai Sim, Y. Zhang, W. Han, T. Strohmaier, and F. Beaufays, “Efficient domain adaptation for speech foundation models,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.
- [49] C.-C. Chiu, J. Qin, Y. Zhang, J. Yu, and Y. Wu, “Self-supervised learning with random-projection quantizer for speech recognition,” in *International Conference on Machine Learning (ICML)*, 2022, pp. 3915–3924.
- [50] D. Hwang, K. C. Sim, Z. Huo, and T. Strohmaier, “Pseudo label is better than human label,” in *Interspeech*, 2022, pp. 1421–1425.
- [51] S. Kessler, B. Thomas, and S. Karout, “An adapter based pre-training for efficient and scalable self-supervised speech representation learning,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 3179–3183.
- [52] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, “Language models are unsupervised multitask learners,” Technical report, 2019.
- [53] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *International*

- Conference on Neural Information Processing Systems (NeurIPS)*, 2017, pp. 6000–6010.
- [54] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [55] S. wen Yang, P.-H. Chi, Y.-S. Chuang, C.-I. J. Lai, K. Lakhota, Y. Y. Lin, A. T. Liu, J. Shi, X. Chang, G.-T. Lin, T.-H. Huang, W.-C. Tseng, K. tik Lee, D.-R. Liu, Z. Huang, S. Dong, S.-W. Li, S. Watanabe, A. Mohamed, and H. yi Lee, “SUPERB: Speech processing Universal PERFORMANCE Benchmark,” in *Interspeech*, 2021, pp. 1194–1198.
- [56] T. Kudo and J. Richardson, “SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing,” in *Conference on Empirical Methods in Natural Language Processing (EMNLP): System Demonstrations*, 2018.
- [57] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” in *International Conference on Learning Representations (ICLR)*, 2019.
- [58] P. Buzzega, M. Boschini, A. Porrello, D. Abati, and S. Calderara, “Dark experience for general continual learning: a strong, simple baseline,” in *International Conference on Neural Information Processing Systems (NeurIPS)*, vol. 33, 2020, pp. 15920–15930.
- [59] J. S. Vitter, “Random sampling with a reservoir,” *ACM Trans. Math. Softw.*, vol. 11, no. 1, pp. 37–57, 1985.
- [60] Z. Wang, Z. Zhang, C.-Y. Lee, H. Zhang, R. Sun, X. Ren, G. Su, V. Perot, J. Dy, and T. Pfister, “Learning to prompt for continual learning,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 139–149.
- [61] B. Lester, R. Al-Rfou, and N. Constant, “The power of scale for parameter-efficient prompt tuning,” in *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2021, pp. 3045–3059.
- [62] G. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” in *NeurIPS Deep Learning and Representation Learning Workshop*, 2015.
- [63] R. Aljundi, F. Babiloni, M. Elhoseiny, M. Rohrbach, and T. Tuytelaars, “Memory aware synapses: Learning what (not) to forget,” in *European Conference on Computer Vision (ECCV)*, 2018, pp. 144–161.
- [64] A. Chaudhry, P. K. Dokania, T. Ajanthan, and P. H. Torr, “Riemannian walk for incremental learning: Understanding forgetting and intransigence,” in *European Conference on Computer Vision (ECCV)*, 2018, pp. 532–547.
- [65] D. Lopez-Paz and M. Ranzato, “Gradient episodic memory for continual learning,” in *International Conference on Neural Information Processing Systems (NeurIPS)*, 2017, pp. 6470–6479.
- [66] M. Mermillod, A. Bugajska, and P. Bonin, “The stability-plasticity dilemma: Investigating the continuum from catastrophic forgetting to age-limited learning effects,” *Frontiers in psychology*, vol. 4, p. 504, 2013.
- [67] J. Gallardo, T. L. Hayes, and C. Kanan, “Self-supervised training enhances online continual learning,” in *British Machine Vision Conference*, 2021.
- [68] D. Madaan, J. Yoon, Y. Li, Y. Liu, and S. J. Hwang, “Representational continuity for unsupervised continual learning,” in *International Conference on Learning Representations (ICLR)*, 2022.
- [69] D. Rao, F. Visin, A. Rusu, R. Pascanu, Y. W. Teh, and R. Hadsell, “Continual unsupervised representation learning,” in *International Conference on Neural Information Processing Systems (NeurIPS)*, vol. 32, 2019.
- [70] C.-Y. Hsieh, C.-L. Li, C.-K. Yeh, H. Nakhost, Y. Fujii, A. Ratner, R. Krishna, C.-Y. Lee, and T. Pfister, “Distilling step-by-step! Outperforming larger language models with less training data and smaller model sizes,” *arXiv preprint arXiv:2305.02301*, 2023.
- [71] K. J. Liang, W. Hao, D. Shen, Y. Zhou, W. Chen, C. Chen, and L. Carin, “MixKD: Towards efficient distillation of large-scale language models,” in *International Conference on Learning Representations (ICLR)*, 2021.
- [72] A. Conneau, M. Ma, S. Khanuja, Y. Zhang, V. Axelrod, S. Dalmia, J. Riesa, C. Rivera, and A. Bapna, “FLEURS: Few-shot learning evaluation of universal representations of speech,” in *IEEE Spoken Language Technology Workshop (SLT)*, 2023, pp. 798–805.
- [73] N. Goyal, C. Gao, V. Chaudhary, P.-J. Chen, G. Wenzek, D. Ju, S. Krishnan, M. Ranzato, F. Guzmán, and A. Fan, “The FLoRes-101 evaluation benchmark for low-resource and multilingual machine translation,” *Transactions of the Association for Computational Linguistics*, vol. 10, pp. 522–538, 2022.

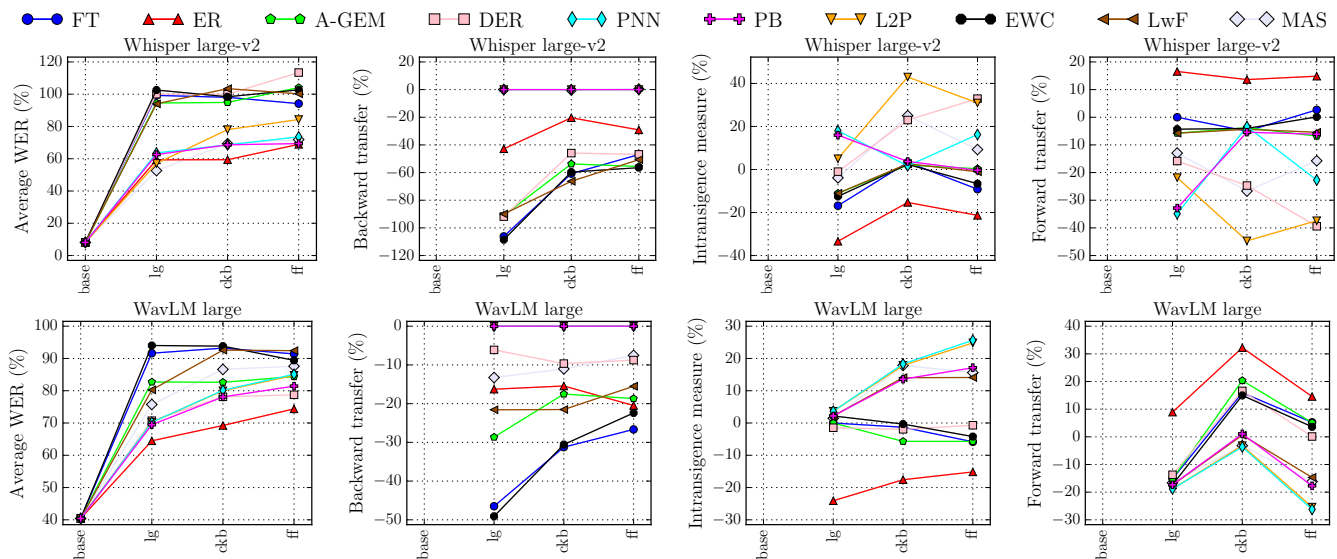


Fig. 4. Comparison of fine-tuning (FT), experience replay (ER), averaged gradient episodic memory (A-GEM), dark experience replay (DER), progressive neural networks (PNN), Piggyback (PB), learning to prompt (L2P), elastic weight consolidation (EWC), learning without forgetting (LwF), and memory aware synapses (MAS) applied to Whisper large-v2 and WavLM large on the base and new languages. ER achieves the smallest AWER across both Whisper and WavLM.

APPENDIX

A. General Information

1) *Dataset Documentation:* The Common Voice 13 [25] is an openly accessible speech dataset that collects the voices of volunteers from all over the world. For extensive documentation refer to the official website (<https://commonvoice.mozilla.org/en/about>).

2) *Intended Uses:* CL-MASR is intended for researchers in continual learning and related fields in order for them to easily test their novel methods on a challenging real-word task.

3) *Hosting and Maintenance Plan:* CL-MASR platform is hosted and version-tracked via GitHub. It is available at <https://github.com/speechbrain/benchmarks>. The download link for the Common Voice 13 [25] dataset can be found on the official website (<https://commonvoice.mozilla.org/en/datasets>).

CL-MASR is a community-driven and open-source initiative. We plan to extend it by running additional experiments and including new continual learning methods and base models. We welcome external contributors.

4) *Licensing:* Our work is licensed under Apache 2.0 (<https://www.apache.org/licenses/LICENSE-2.0>). The Common Voice 13 [25] dataset is licensed under CC0 public domain Creative Commons (<https://creativecommons.org/share-your-work/public-domain/cc0>).

5) *Author Statement:* We, the authors, will bear all responsibility in case of violation of rights.

B. FLEURS Dataset

To showcase the flexibility of our benchmark platform and to further validate the results of our analysis, we extend our comparative study to the **Few-shot Learning Evaluation of Universal Representations of Speech (FLEURS)**⁷ [72] dataset. Derived from the machine translation FLoRes-101 [73] benchmark, it consists of n -way parallel natural human speech and text in 102 languages, with approximately 12 hours of annotated data per language divided into training, validation, and test splits. Following the same procedure as Common Voice 13, we select from it two sets of languages ($L_{\text{base}} = L_{\text{new}} = 3$):

- **base languages:** English (en), German (de), and Arabic (ar).
- **new languages:** Luganda (lg), Central Kurdish (ckb), and Fula (ff).

For each language, we randomly extract up to 10 hours of data for training, 1 for validation, and 1 for testing respectively. For the preprocessing and experimental evaluation, we use the same hyperparameters as Common Voice 13.

Fig. 4 shows the outcomes of our experiments on FLEURS. The main findings align with those from Common Voice 13: rehearsal-based and architecture-based methods are generally effective for CL in ASR, while regularization-based approaches tend to perform poorly. Remarkably, ER achieves the smallest AWER across both Whisper and WavLM. However, we observe some noteworthy exceptions. When applied to Whisper, DER struggles to both learn new tasks and mitigate forgetting. In contrast, MAS exhibits exceptional performance on Whisper, on par with ER. Lastly, when applied to WavLM, PNN shows scarce plasticity (WER $\sim 100\%$ for each new

⁷https://huggingface.co/datasets/google/xtreme_s

⁸Licensed under CC-BY public domain Creative Commons (<https://creativecommons.org/licenses/by/4.0/>).

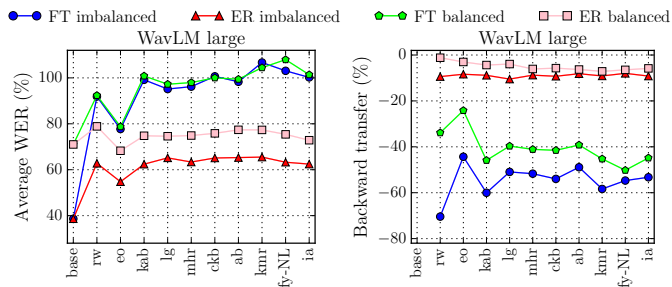


Fig. 5. Comparison of the imbalanced and balanced variants of fine-tuning (FT) and experience replay (ER) applied to WavLM large on the base and new languages (10 hours per language). The steep increase in AWER and forgetting after learning the first new language is observed only for the imbalanced variants. ER effectively mitigates forgetting in both the balanced and imbalanced scenarios.

language), performing similarly to L2P. This highlights how the performance of most CL methods is highly sensitive to both the dataset and the selection of hyperparameters.

C. Effect of Imbalance between Base and New Languages

We hypothesize that the substantial drop in accuracy observed after learning the first new language (see Fig. 2) is associated with imbalance in terms of amount of data and/or number of training steps between the initial supervised joint pretraining and the subsequent incremental training. To confirm this, we conduct an experiment involving imbalanced and balanced variants of fine-tuning (FT) and experience replay (ER) applied to WavLM large on the base and new languages (10 hours per language). Imbalanced means initial supervised joint pretraining on the base languages for 20 epochs, followed by incremental training for 2 epochs per language (as in the main results). Balanced means initial supervised joint pretraining on the base languages for 2 epochs, followed by incremental training for 2 epochs per language. Fig. 5 shows that in the balanced scenario, the increase in forgetting is significantly smaller. This validates our hypothesis that imbalance is the primary cause of the observed phenomenon. Even though the imbalanced case is characterized by this spike in WER, overall, AWER is smaller compared to the balanced case. This is expected, as the imbalanced case is trained with more data and/or epochs, and is likely the practical use case scenario. Remarkably, ER effectively mitigates forgetting in both balanced and imbalanced scenarios.

D. Additional Results

Tables II, III, V, and IV present a more comprehensive results from our experiments on the Common Voice dataset.

