

# Continual Error Correction on Low-Resource Devices

Kirill Paramonov\*  
Samsung R&D Institute UK  
United Kingdom

Mete Ozay  
Samsung R&D Institute UK  
United Kingdom

Aristeidis Mystakidis  
CERTH  
Greece

Nikolaos Tsalikidis  
CERTH  
Greece

Dimitrios Sotos  
CERTH  
Greece

Anastasios Drosou  
CERTH  
Greece

Dimitrios Tzovaras  
CERTH  
Greece

Hyunjun Kim  
Samsung Research  
South Korea

Kiseok Chang  
Samsung Research  
South Korea

Sangdok Mo  
Samsung Research  
South Korea

Namwoong Kim  
Samsung Research  
South Korea

Woojong Yoo  
Samsung Research  
South Korea

Jijoong Moon  
Samsung Research  
South Korea

Umberto Michieli\*<sup>†</sup>  
Samsung R&D Institute UK  
United Kingdom

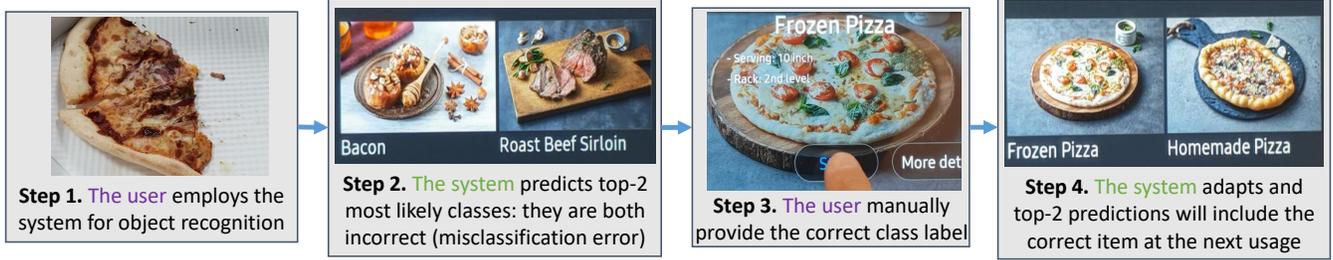


Figure 1: Proposed continual few-shot error correction system and user interactions flow in a food recognition use case.

## Abstract

The proliferation of AI models in everyday devices has highlighted a critical challenge: prediction errors that degrade user experience. While existing solutions focus on error detection, they rarely provide efficient correction mechanisms, especially for resource-constrained devices. We present a novel system enabling users to correct AI misclassifications through few-shot learning, requiring minimal computational resources and storage. Our approach combines server-side foundation model training with on-device prototype-based classification, enabling efficient error correction through prototype updates rather than model retraining. The system consists of two key components: (1) a server-side pipeline that

leverages knowledge distillation to transfer robust feature representations from foundation models to device-compatible architectures, and (2) a device-side mechanism that enables ultra-efficient error correction through prototype adaptation. We demonstrate our system’s effectiveness on both image classification and object detection tasks, achieving over 50% error correction in one-shot scenarios on Food-101 and Flowers-102 datasets while maintaining minimal forgetting (less than 0.02%) and negligible computational overhead. Our implementation, validated through an Android demonstration app, proves the system’s practicality in real-world scenarios.

\*Both authors contributed equally to this research.

<sup>†</sup>Corresponding author: u.michieli@samsung.com.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

MMSys '25, Stellenbosch, South Africa

© 2025 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-1467-2/2025/03

<https://doi.org/10.1145/3712676.3719269>

## CCS Concepts

• **Computing methodologies** → **Scene understanding; Object recognition; Object identification; Causal reasoning and diagnostics;** • **Human-centered computing** → **Ubiquitous and mobile computing systems and tools.**

## Keywords

Error Correction, Continual Learning, AI Mistakes

**ACM Reference Format:**

Kirill Paramonov, Mete Ozay, Aristeidis Mystakidis, Nikolaos Tsalikidis, Dimitrios Sotos, Anastasios Drosou, Dimitrios Tzovaras, Hyunjun Kim, Kiseok Chang, Sangdok Mo, Namwoong Kim, Woojong Yoo, Jijoong Moon, and Umberto Michieli. 2025. Continual Error Correction on Low-Resource Devices. In *ACM Multimedia Systems Conference 2025 (MMSys '25), March 31-April 4, 2025, Stellenbosch, South Africa*. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/3712676.3719269>

**1 Introduction**

The integration of AI functionalities into everyday digital devices has grown exponentially, particularly in discriminative applications like image classification and object detection [21]. From smart home appliances to mobile devices, AI models are becoming increasingly prevalent in enhancing user experience through automated recognition, classification, and detection tasks. However, these AI models frequently make prediction errors that frustrate users and diminish product reliability [31]. This is particularly problematic in consumer devices where users expect consistent, accurate performance and have no means to correct persistent errors.

The challenge of AI model errors is compounded by several factors in real-world applications. First, the deployment environment often differs significantly from the training environment, leading to domain shift issues. Second, personal devices frequently encounter user-specific patterns and preferences that were not represented in the training data. Third, resource constraints on edge devices limit the applicability of traditional error correction approaches that require substantial computational power or storage.

Current approaches to handling AI errors primarily focus on detection rather than correction [9, 12, 23]. When correction methods exist, they typically require substantial resources: either collecting additional training data [29], finetuning models [19], or maintaining large storage capacities [7]. Such requirements make these solutions impractical for consumer devices with limited computational and storage resources. Moreover, these approaches often necessitate cloud connectivity for model updates, raising privacy concerns and limiting functionality in offline scenarios.

Furthermore, existing error correction methods often suffer from catastrophic forgetting, where correcting one type of error leads to degraded performance on previously well-handled cases. This creates a frustrating user experience where fixing one issue simply leads to another, ultimately reducing trust in the AI system. The challenge lies in developing a system that can adapt to user corrections while maintaining overall model performance and operating within device constraints.

We present a novel system that enables efficient, on-device error correction for AI models. Our solution allows users to correct model mistakes using just a few labelled examples without requiring significant computational resources or storage. For instance, as shown in Figure 1, users can correct misclassification errors with just one example (1-shot error correction), making the system highly practical for real-world applications.

Our technical approach combines three key innovations that address the challenges of on-device error correction:

- Server-side knowledge distillation from foundation models to create robust yet compact models suitable for resource-constrained devices, leveraging the power of large models while maintaining practical deployability;
- A lightweight, prototype-based classification pipeline provided with an efficient prototype update mechanism that allows continual adaptation to user corrections while preventing catastrophic forgetting;
- A practical demonstration via an Android application that showcases real-time error correction capabilities, validating our approach in practical scenarios.

Our system addresses several key requirements for actual deployment: (i) *resource efficiency*, as the system operates within the computational and storage constraints of typical consumer devices; (ii) *privacy preservation*, as all error corrections happen locally on the device, ensuring user data remains private; (iii) *ease of use*, as users can correct errors through simple interactions; (iv) *stability*, as the system maintains overall performance while adapting to user corrections; (v) *offline operation*, as the system does not require cloud connectivity.

Experimental results demonstrate the effectiveness of our system, achieving over 50% error correction accuracy in one-shot scenarios while maintaining minimal computational overhead. We validate these results through both benchmark evaluations and real-world testing via our Android demonstration app. Our experiments cover a range of scenarios, from standard image classification tasks to more complex object detection applications, showing the system's versatility across different use cases.

The rest of this paper is organized as follows: Section 2 discusses related work and existing approaches to AI error correction. Section 3 formulates the problem setup and introduces our evaluation metrics. Section 4 provides a detailed description of our system architecture, including both server-side and device-side components. Section 5 presents our experimental results and demonstration implementation. Finally, Section 6 concludes the paper.

**2 Related Work**

Existing work related to AI model errors can be categorized into three main approaches.

1) *Error Detection and Confidence Estimation*: most current solutions focus solely on analysing [27] or detecting misclassifications [7, 9, 12, 19, 23]. While these methods are effective at identifying errors, they don't provide correction mechanisms. In [17] misclassifications of trained models are reduced thanks to the addition of a regularization term in the optimization process. Related approaches [24, 25] estimate model confidence to implement quality of service policies, such as discarding low-confidence predictions. In contrast, our work assumes users can provide ground truth labels for misclassified samples and focuses on actually correcting these errors with no need for backpropagation.

2) *Domain Generalization and Adaptation*: several approaches attempt to improve model robustness through domain generalization at pre-training, aiming to handle multiple domains or adversarial samples [2, 6, 10, 34]. Similarly, domain adaptation techniques [3, 8, 30, 33] focus on adapting networks to domains different from

their pre-training distribution. However, these approaches typically require large amounts of labelled data and don't address the correction of individual model mistakes.

3) *Distribution Shift Detection*: a well-established research direction focuses on out-of-distribution detection, aiming to identify when input samples come from distributions unseen during pre-training [5, 15, 16, 32]. While these methods can identify potentially problematic inputs, they differ fundamentally from our approach as they do not provide mechanisms to adapt models and correct previous misclassifications.

Our work provides a practical and resource-efficient solution for error correction that operates directly on user devices. Unlike previous methods, our system enables immediate error correction with minimal user input while maintaining low computational and storage requirements. By focusing on prototype updates rather than model retraining, we achieve efficient personalization without the need for extensive data collection or computational resources.

### 3 Problem Formulation

We address the challenging problem of correcting misclassification errors in AI models deployed on resource-constrained devices. While this is a critical issue for practical AI applications, it has received limited attention in existing literature. Our focus is on discriminative AI models for computer vision applications.

#### 3.1 Problem Setup

Our system's deployment process consists of three main phases.

*Device Requirements Analysis*. The process begins with analyzing the target device's capabilities and constraints. We consider:

- hardware constraints (memory, compute, power, storage);
- maximum allowable inference time;
- minimum required accuracy.

Based on these constraints, we select an appropriate model architecture  $M_I$  that offers the optimal accuracy-footprint tradeoff. Additional optimizations, such as quantization to 4-bit precision, may be applied if supported by the target hardware.

*Server-side Pre-training*. The selected model  $M_I$  undergoes pre-training on a training dataset  $\mathcal{D}_{train}$  to produce the trained model  $M_T$ . The training samples in the dataset are indexed by  $i = 1, \dots, n$  and are composed of input images  $X_i$  and corresponding ground truth label  $y_i \in C$ , namely  $\mathcal{D}_{train} = \{(X_i, y_i)\}_{i=1}^n$ . The pre-training techniques employed in our system are described in Section 4.1.

*On-device Error Correction*. After deployment, the model  $M_T$  processes a stream of data  $\mathcal{D}_{test} = (X_i^{test}, y_i^{test})_{i=1}^{n_{test}}$ , where labels  $y_i^{test}$  are only available when provided by users. The test dataset naturally partitions into: (i) correctly classified samples  $\mathcal{D}_C \subset \mathcal{D}_{test}$ , and (ii) misclassified ones  $\mathcal{D}_E \subset \mathcal{D}_{test}$ , where  $\mathcal{D}_E \cap \mathcal{D}_C = \emptyset$ .

Users can provide ground truth labels for  $s$  samples from  $\mathcal{D}_E$  for some classes, enabling few-shot error correction. The model  $M_T$  is then adapted using these annotations to produce  $M_A$ , with improved performance on previously misclassified samples. The adaptation is efficient and backpropagation-free (see Section 4.2).

### 3.2 Evaluation Metrics

We evaluate our system via three key metrics. We refer to *accuracy* as the percentage of correctly classified images out of the total number of images in the set. We computed:

- (1) **Base Recognition Accuracy** ( $\uparrow$ ): Accuracy of  $M_T$  on the test set  $\mathcal{D}_{test}$ , named  $Acc_{base}$
- (2) **Error Correction Accuracy** ( $\uparrow$ ): Accuracy of  $M_A$  on the misclassified set  $\mathcal{D}_E$ , named  $Acc_E$
- (3) **Forgetting Rate Percentage** ( $\downarrow$ ): we compute the accuracy of the adapted model  $M_A$  on the previously correctly-classified samples  $\mathcal{D}_C$ , named  $Acc_C$ . Then, the forgetting rate is computed as:  $For := 100 - Acc_C$ .

Therefore, a successful error correction system should achieve: (i) high error correction accuracy ( $Acc_E$ ), (ii) minimal user labelling requirement (small  $s$  for low-shot learning), (iii) low forgetting rate ( $For$ ), (iv) controlled footprints.

Our solution ensures deterministic control over inference time and memory usage increases, with built-in mechanisms to handle resource constraints (detailed in Section 4.2).

## 4 Our System

Our system consists of two sequential pipelines: (1) server-side pre-training followed by (2) on-device deployment and error correction. We first present the system for image classification and then extend it to object detection in Section 4.3.

### 4.1 Server-side Pipeline

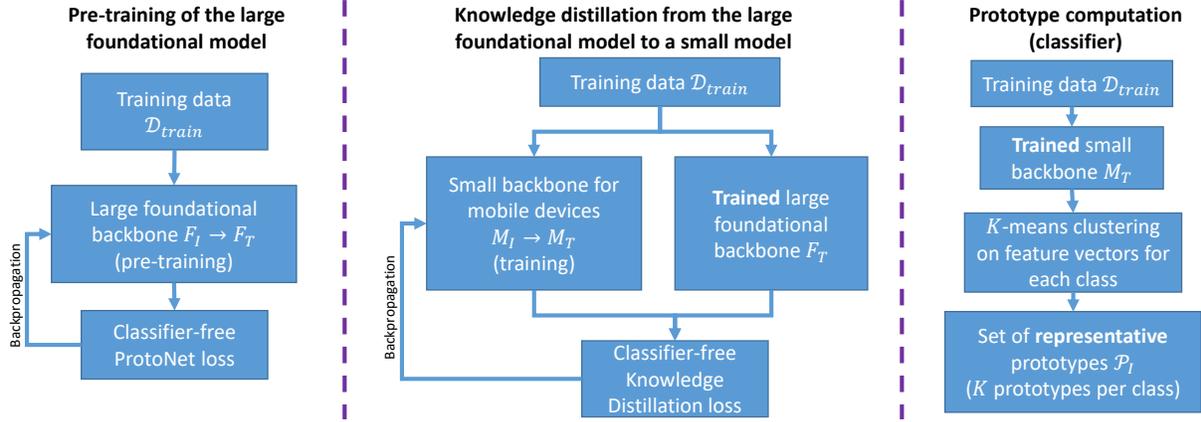
The server-side pipeline consists of three main stages, as illustrated in Figure 2.

*4.1.1 Foundation Model Pre-training*. Motivated by the recent success of large foundation vision models, we leverage the transformer-based DINO-v2 [22] model, finetuning it on domain-specific training data  $\mathcal{D}_{train}$  at the server side (e.g., food data for food recognition applications). More specifically, starting from pre-trained weights ( $F_T$ ), we train DINO-v2 without an output classifier (i.e., no output linear layer) using ProtoNet loss [28] to obtain  $F_T$ . Intuitively, the ProtoNet loss learns a metric space where classification can be performed by computing distances to class prototypes, providing a robust feature extractor while maintaining flexibility compared to a linear classifier head.

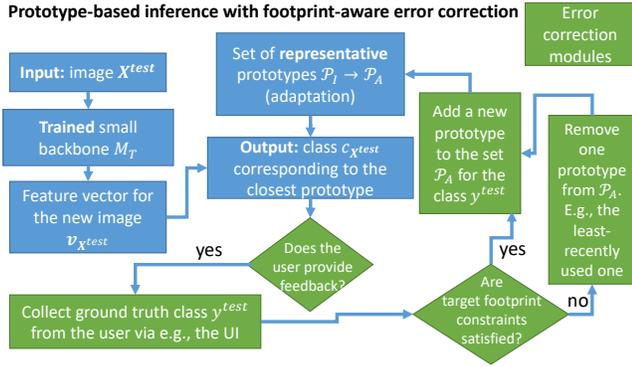
*4.1.2 Knowledge Distillation into Small Model*. Due to the limited-resource constraints on the device side, the finetuned large foundation model ( $F_T$ ) may not satisfy the target footprints (e.g., model size and inference time) on the device (see Section 3.1). If this is the case, we select a smaller model  $M_I$  (e.g., the convolutional-based MobileNet-V2 [26]). We load the pre-trained weights and continue training with knowledge distillation [11] from the frozen large (teacher) model  $F_T$  over the training data  $\mathcal{D}_{train}$ , obtaining  $M_T$ . As the distillation loss, we use an L1 distance between the class tokens of the teacher and student models, namely:

$$\mathcal{L}_{KD} = \frac{1}{n} \sum_{i=1}^n |F_T(X_i) - M(X_i)|, \quad (1)$$

where  $M$  is initialized as  $M_I$  and trained to obtain  $M_T$ .



**Figure 2: The server-side image classification pipeline is composed of three stages: (i) training of a large foundation model (Section 4.1.1), (ii) knowledge distillation from the trained foundation model to a smaller model (Section 4.1.2), and (iii) computation of prototypes as the classifier (Section 4.1.3).**



**Figure 3: The device-side image classification pipeline with error correction blocks in green.**

**4.1.3 Prototype Computation.** So far, we considered models without any classifier head. Therefore, the final stage at the server side is to compute a set of representative prototypes for each class  $c$  in  $\mathcal{D}_{train}$  using the trained model  $M_T$ .

Specifically, we compute the feature vectors of every sample,  $\mathbf{v}_{X_i} = M_T(X_i)$ ,  $\forall i = 1, \dots, n$ , and compute a  $K$ -means clustering for each class  $c$ . Namely, defining as  $\mathcal{V}$  the set of features associated to ground truth label  $c$ , i.e.,  $\mathcal{V} = \{\mathbf{v}_{X_i} \mid \forall i, s.t. y_i = c\}$ , the objective of the  $K$ -means for class  $c$  is to find:

$$\arg \min_{\mathcal{P}} \sum_{j=1}^K \sum_{\mathbf{v} \in \mathcal{V}} \|\mathbf{v} - \mathbf{p}_{c,j}\|^2 \quad (2)$$

where the  $K$  cluster centroids are denoted as  $\mathcal{P}_c = \{\mathbf{p}_{c,j}\}_{j=1}^K$  and are the representative prototypes for class  $c$ . In other words,  $\mathcal{P}_c$  is a set containing the  $K$  representative prototypes for class  $c$ .

The set of all representative prototypes across all classes together with their corresponding class labels is denoted as  $\mathcal{P}_I = \{(\mathcal{P}_c, c)\}$ ,  $\forall c \in \mathcal{C}$ .

## 4.2 Device-side Pipeline

The device-side pipeline is shown in Figure 3.

**Classification Process.** The system receives an input image  $X^{test}$ . The image is passed through the small backbone  $M_T$  (pre-trained on the server) to obtain the feature vector  $\mathbf{v}_{X^{test}} = M_T(X^{test})$ . We initialize a new set of adapted prototypes  $\mathcal{P}_A = \mathcal{P}_I$ . Then, we compute the cosine distance between the feature vector  $\mathbf{v}_{X^{test}}$  and all prototypes  $\mathbf{p}_{c,j} \in \mathcal{P}_c, \forall c \in \mathcal{C}, \forall j = 1, \dots, K$  from the set of representative prototypes  $\mathcal{P}_A$  (constructed on the server). Finally, we select the class corresponding to the minimum cosine distance as the predicted inferred class  $c_{X^{test}}$ . In other words:

$$c_{X^{test}} := \arg \min_c (\cos \text{dist}(\mathbf{v}_{X^{test}}, \mathcal{P}_c)). \quad (3)$$

We note that cosine distance function  $\cos \text{dist}()$  computes the cosine distance between  $\mathbf{v}_{X^{test}}$  and all elements of  $\mathcal{P}_c$ .

**Error Correction.** If the inferred class  $c_{X^{test}}$  is incorrect (i.e., the prediction is different from the ground truth  $y^{test}$ ), the user could optionally decide to manually annotate the image  $X^{test}$  with the correct label  $y^{test}$ , via e.g., the user interface (UI). Then, our error correction system updates the array of representative prototypes:

$$\mathcal{P}_A := \{(\mathcal{P}_c, c) \in \mathcal{P}_A, c \neq y^{test}\} \cup \quad (4)$$

$$\{(\text{add}(\mathcal{P}_{y^{test}}, \mathbf{v}_{X^{test}}), y^{test}), (\text{remove}(\mathcal{P}_{y^{test}}, y^{test}), y^{test}) \in \mathcal{P}_A\}. \quad (5)$$

The first term ensures that representative prototypes of classes other than the current one are not modified. The second term adds the current feature vector  $\mathbf{v}_{X^{test}}$  to the current list of representative prototypes for class  $y^{test}$ , i.e.,  $\mathcal{P}_{y^{test}}$ .

To maintain resource constraints, the system drops one prototype from the set  $\mathcal{P}_A$  whenever the maximum memory or footprint target limits are hit. For example, the least-recently-used prototype can be removed. The system continually updates  $\mathcal{P}_A$  based on user feedback by adding incoming user-annotated feature vectors and uses  $\mathcal{P}_A$  to predict the class of any follow-up images.

### 4.3 Extension to Object Detection

Extending the system outlined so far to object detection can be done easily by incorporating an object detection network and image slicing mechanism based on the predicted bounding boxes.

*Server-side Modifications.* On the server side, we incorporate a state-of-the-art efficient object detection network (e.g., YoloV8 [13]) and finetune it to predict bounding boxes of objects of interest. In the food domain example, the object detector could be trained to identify bounding boxes corresponding to food items in the scene (i.e., one-class classification).

*Device-side Implementation.* On the device side, we use the fine-tuned model to extract the bounding boxes of the test image. Then, for each bounding box, we crop the image region and process it through the image classification network.

*Implementation Trade-offs.* Our current implementation uses separate networks for detection and classification, which offers maximum accuracy but higher computational cost.

To save complexity, one could use detection backbone features directly for classification. This optimization is demonstrated in [1]. The single backbone solution would reduce computational overhead at a cost of less than 5% relative accuracy. While further optimizations are possible by building prototypical inference directly on detection features, we leave this exploration for future work to maintain focus on our core contribution of error correction.

## 5 Experimental Results

### 5.1 Experimental Setup

*5.1.1 Implementation Details.* The server-side pipeline is implemented in Python using PyTorch, and training is conducted on a single NVIDIA GeForce RTX 3090 Ti GPU. Model quantization and initial server-side testing of the quantized models were implemented on an AMD Ryzen 9 5950X CPU. For on-device demonstrations, the model  $M_T$  is exported to ONNX and then converted to TFLite, followed by uniform 8-bit quantization. The quantized TFLite model is used with the TFLite C++ interpreter for inference. The device-side pipeline is integrated into NNTrainer [18] for deployment on Android devices. Demonstrations are conducted on a Samsung Galaxy S24 Ultra smartphone.

The complete codebase for both server-side and device-side pipelines will be open-sourced upon acceptance.

*5.1.2 Models.* For the Image Classification task, the transformer-based foundation model DINO-V2-small [22] is distilled into an efficient convolutional MobileNet-V2 [26]. For the Object Detection task, the YoloV8-nano model [13] is utilized.

*5.1.3 Hyperparameters.* For the server-side pre-training, we train DINO-v2-small for 100 epochs with a batch size of 256 and a learning rate of  $1e-5$ . Knowledge distillation onto MobileNet-V2 is done for 100 epochs with a batch size of 256 and a learning rate of 0.1. Finally, we used YoloV8-nano [13] as the object detection model and we finetuned it on the OpenImages-V7 [14] dataset reaching 0.347 mAP50. We trained it for 200 epochs with early stopping enabled with patience 20 epochs, using a batch size of 16 and a learning rate

**Table 1: Error correction accuracy and forgetting rate of our system on two image classification datasets.**

| # shots $s$ | Food-101         |                  | Flowers-102      |                  |
|-------------|------------------|------------------|------------------|------------------|
|             | $Acc_E \uparrow$ | $For \downarrow$ | $Acc_E \uparrow$ | $For \downarrow$ |
| 1           | 51.1%            | 0.018%           | 54.3%            | 0.011%           |
| 2           | 67.3%            | 0.053%           | 68.9%            | 0.038%           |
| 3           | 75.0%            | 0.118%           | 78.4%            | 0.094%           |
| 4           | 79.9%            | 0.157%           | 81.1%            | 0.114%           |
| 5           | 83.2%            | 0.179%           | 85.7%            | 0.162%           |
| 7           | 87.3%            | 0.202%           | 89.0%            | 0.179%           |
| 10          | 91.0%            | 0.237%           | 92.6%            | 0.195%           |
| 20          | 93.2%            | 0.447%           | 96.7%            | 0.364%           |
| 50          | 96.4%            | 0.870%           | 98.8%            | 0.577%           |

of 0.01. For the server-side prototypes computation, we set  $K = 3$ , i.e., we store the three most representative prototypes per class.

*5.1.4 Datasets.* We used two datasets: the Food-101 [4] and the Oxford Flowers-102 [20], containing 101 and 102 classes, respectively. Both datasets represent domain-specific applications (e.g., food or flower recognition), where the user may want to provide the correct label to improve the AI model. As these datasets are out-of-domain for DINO-V2, domain finetuning is required. Each dataset is split into 70% training, 15% validation, and 15% testing per class; i.e., we preserve 70% of the images of each class in the training set. Few-shot error correction experiments are conducted with  $s \in \{1, 2, 3, 4, 5, 7, 10, 20, 50\}$  annotated samples per class.

### 5.2 Image Classification Results

Quantitative results of our system are reported in Table 1 in terms of error correction accuracy  $Acc_E$  and forgetting rate  $For$  for the image classification task. The base recognition accuracy  $Acc_{base}$  of the  $M_T$  model with initial representation prototypes  $P_I$  is 90.6% and 94.3% on the Food-101 and the Flowers-102 datasets, respectively.

Our system allows user to correct most of the misclassification errors with limited feedback. For example, correcting 51.1% of classification errors on Food-101 requires the user to annotate just one sample per class. At the same time, only 0.018% of previously-corrected samples are now misclassified by the adapted model.

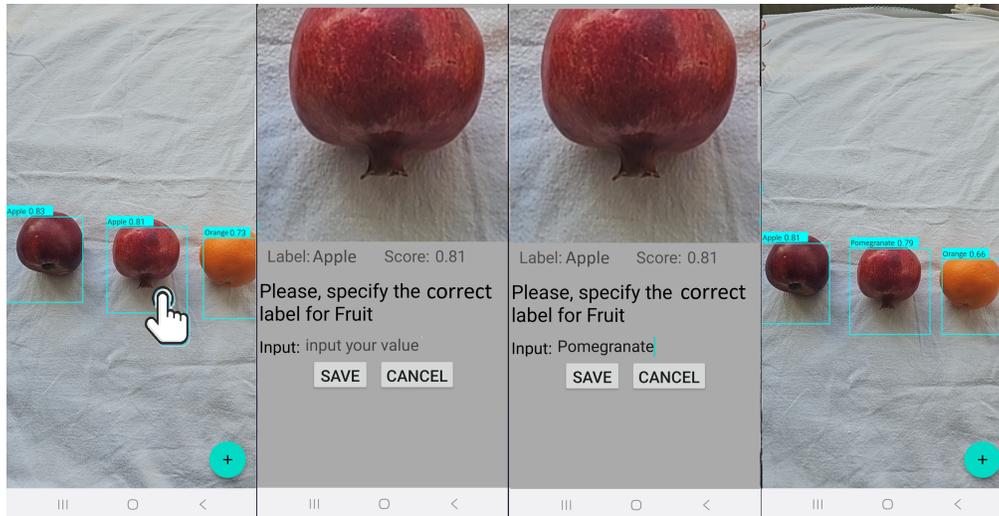
#### 5.2.1 Ablation Studies.

*Prototypes vs. Linear Classifier Head.* First, we checked that using a fully connected linear classifier head does not produce higher base classification accuracy (91.1% on Food-101 and 93.6% on Flowers-102) than prototypical inference (90.6% and 94.3%, respectively). Furthermore, fully-connected layers are less adaptable, and mistakes are harder to correct with efficient updates.

*Server-side Pipeline.* Table 2 highlights the performance of different server-side pre-training strategies on Food-101. Using only the small model finetuned on domain-specific data, i.e.,  $M_T$ , results in low  $Acc_E$  due to limited feature richness. Using the pre-trained foundation model  $F_I$  gives a good balance between error correction accuracy and forgetting, nonetheless the performance on the base classes is quite poor as the model is pre-trained on a large

**Table 2: Ablation studies on the server-side pipeline on the Food-101 dataset.  $Acc_{base}$  for all approaches is (from left to right): 82.6, 76.4, 92.4, and 90.6%.**

| # shots $s$ | $M_T$ only       |                  | $F_I$            |                  | $F_T$            |                  | Ours             |                  |
|-------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|
|             | $Acc_E \uparrow$ | $For \downarrow$ |
| 1           | 0.6%             | 0.1%             | 48.7%            | 0.1%             | 15.1%            | 0.1%             | 51.1%            | 0.0%             |
| 3           | 7.0%             | 0.1%             | 72.3%            | 0.3%             | 48.3%            | 0.1%             | 75.0%            | 0.1%             |
| 5           | 13.2%            | 0.3%             | 82.9%            | 0.4%             | 57.3%            | 0.1%             | 83.2%            | 0.2%             |
| 10          | 21.7%            | 0.6%             | 91.2%            | 0.6%             | 65.4%            | 0.2%             | 91.0%            | 0.2%             |

**Figure 4: Application workflow of our demo: (1) the user identifies misclassified objects and clicks anywhere within their bounding box, (2-3) the user provides the correct label, and (4) the user observes improved classification in subsequent iterations.**

generic dataset and not domain-finetuned. Additionally, the foundation model is typically an order of magnitude larger than the selected small model. The finetuned foundation model  $F_T$  overfits to the domain-specific data, achieving high  $Acc_{base}$  but poor adaptability to correct mistakes efficiently. Finally, our approach strikes an optimal balance, sacrificing only a small amount of error recognition accuracy (2.2%) to achieve significantly improved error correction capability and minimal forgetting, all while being more computationally efficient than foundation model-based methods.

### 5.3 On-device Demonstration

In our demo, we focus on the object detection task due to its practical value where users point a camera to target objects. Both MobileNetV2 and YoloV8-nano are quantized to 8 bits, achieving a  $\sim 3\times$  inference speedup with negligible accuracy degradation ( $<1\%$ ) compared to the 32-bit models. For instance, on a Samsung Galaxy S24, YoloV8-nano achieves 40 ms inference time (compared to 110 ms for 32-bit) while maintaining a mAP@50 of 0.347 (from 0.350). The 8-bit model has a compact size of 3.2 MB, enabling real-time inference on consumer devices.

Overall, the system runs in real-time on consumer smartphones. Figure 4 illustrates the application workflow. In the first screenshot, the user can benefit from the detection model to identify some food

items. In case of a misclassification error, the user may decide to provide the correct ground truth: in our app, the user should click on the bounding box and enter the correct label (second and third screenshots). Finally, the system updates its internal prototypical representations and will correctly identify the item at subsequent iterations (third screenshot).

## 6 Conclusion

In this work, we introduced a simple yet effective system for few-shot continual error correction in discriminative AI models, such as those used for object recognition. Our approach empowers end users to optionally and seamlessly correct misclassification errors in AI models, enhancing their overall utility and reliability. The proposed system comprises a server-side pre-training stage to develop a robust feature extractor while adhering to the resource constraints of target devices, such as memory, storage, and inference time. Additionally, we designed a novel on-device mechanism for few-shot continual error correction, which efficiently updates the stored prototypical representations. Our system represents a step forward in enabling pervasive AI applications, setting new benchmarks for accuracy, adaptability, and practicality in consumer products.

## References

- [1] Francesco Barbato, Umberto Michieli, Jijoong Moon, Pietro Zanuttigh, and Mete Ozay. 2024. Cross-architecture auxiliary feature space translation for efficient few-shot personalized object detection. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 8587–8594.
- [2] Francesco Barbato, Umberto Michieli, Mehmet Kerim Yucel, Pietro Zanuttigh, and Mete Ozay. 2024. A Modular System for Enhanced Robustness of Multimedia Understanding Networks via Deep Parametric Estimation. In *Proceedings of the 15th ACM Multimedia Systems Conference (Bari, Italy) (MMSys '24)*. Association for Computing Machinery, New York, NY, USA, 190–201.
- [3] Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. 2006. Analysis of representations for domain adaptation. *Advances in neural information processing systems* 19 (2006).
- [4] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. 2014. Food-101—mining discriminative components with random forests. In *Computer vision—ECCV 2014: 13th European conference, Zurich, Switzerland, September 6–12, 2014, proceedings, part VI 13*. Springer, 446–461.
- [5] Elena Camuffo, Umberto Michieli, Simone Milani, Jijoong Moon, and Mete Ozay. 2024. Enhanced Model Robustness to Input Corruptions by Per-corruption Adaptation of Normalization Statistics. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 11558–11565.
- [6] Guangyao Chen, Peixi Peng, Li Ma, Jia Li, Lin Du, and Yonghong Tian. 2021. Amplitude-phase recombination: Rethinking robustness of convolutional neural networks in frequency domain. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 458–467.
- [7] Charles Corbière, Nicolas Thome, Avner Bar-Hen, Matthieu Cord, and Patrick Pérez. 2019. Addressing failure prediction by learning model confidence. *Advances in Neural Information Processing Systems* 32 (2019).
- [8] Yaroslav Ganin and Victor Lempitsky. 2015. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*. PMLR, 1180–1189.
- [9] Dan Hendrycks and Kevin Gimpel. 2016. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136* (2016).
- [10] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. 2019. Augmix: A simple data processing method to improve robustness and uncertainty. *arXiv preprint arXiv:1912.02781* (2019).
- [11] Geoffrey Hinton. 2015. Distilling the Knowledge in a Neural Network. *arXiv preprint arXiv:1503.02531* (2015).
- [12] Heinrich Jiang, Been Kim, Melody Guan, and Maya Gupta. 2018. To trust or not to trust a classifier. *Advances in neural information processing systems* 31 (2018).
- [13] Glenn Jocher, Ayush Chaurasia, and Jing Qiu. 2023. Ultralytics YOLOv8. <https://github.com/ultralytics/ultralytics>
- [14] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, et al. 2020. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *International journal of computer vision* 128, 7 (2020), 1956–1981.
- [15] Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. 2017. Enhancing the reliability of out-of-distribution image detection in neural networks. *arXiv preprint arXiv:1706.02690* (2017).
- [16] Weitang Liu, Xiaoyn Wang, John Owens, and Yixuan Li. 2020. Energy-based out-of-distribution detection. *Advances in neural information processing systems* 33 (2020), 21464–21475.
- [17] Elyes Manai, Mohamed Mejri, and Jaouhar Fattahi. 2024. Minimizing Model Misclassification Using Regularized Loss Interpretability. In *2024 16th IIAI International Congress on Advanced Applied Informatics (IIAI-AAI)*. IEEE, 231–236.
- [18] Jijoong Moon, Hyun Suk Lee, Jiho Chu, Donghak Park, Seungbaek Hong, Hyungjun Seo, Donghyeon Jeong, Sungsik Kong, and MyungJoo Ham. 2024. A New Frontier of AI: On-Device AI Training and Personalization. In *2024 IEEE/ACM 46th International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*. IEEE Computer Society, 323–333.
- [19] Makaela Nartker, Zhenglong Zhou, and Chaz Firestone. 2023. When will AI misclassify? Intuiting failures on natural images. *Journal of Vision* 23, 4 (2023), 4–4.
- [20] Maria-Elena Nilsback and Andrew Zisserman. 2008. Automated flower classification over a large number of classes. In *2008 Sixth Indian conference on computer vision, graphics & image processing*. IEEE, 722–729.
- [21] Numalis. 2023. How AI is revolutionizing the home appliance industry. *Communication team* (2023).
- [22] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. 2023. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193* (2023).
- [23] Xin Qiu and Risto Miikkilainen. 2022. Detecting misclassification errors in neural networks with a gaussian process model. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 36. 8017–8027.
- [24] Haoxuan Qu, Lin Geng Foo, Yanchao Li, and Jun Liu. 2023. Towards more reliable confidence estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023).
- [25] Haoxuan Qu, Yanchao Li, Lin Geng Foo, Jason Kuen, Jiuxiang Gu, and Jun Liu. 2022. Improving the reliability for confidence estimation. In *European Conference on Computer Vision*. Springer, 391–408.
- [26] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. 2018. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4510–4520.
- [27] Daniel Sikar, Artur Garcez, Robin Bloomfield, Tillman Weyde, Kaleem Peeroo, Naman Singh, Maeva Hutchinson, Dany Laksono, and Mirela Reljan-Delaney. 2024. The Misclassification Likelihood Matrix: Some Classes Are More Likely To Be Misclassified Than Others. *arXiv preprint arXiv:2407.07818* (2024).
- [28] Jake Snell, Kevin Swersky, and Richard Zemel. 2017. Prototypical networks for few-shot learning. *Advances in neural information processing systems* 30 (2017).
- [29] Nathan TeBlunthuis, Valerie Hase, and Chung-Hong Chan. 2024. Misclassification in automated content analysis causes bias in regression. Can we fix it? Yes we can! *Communication Methods and Measures* (2024), 1–22.
- [30] Marco Toldo, Andrea Maracani, Umberto Michieli, and Pietro Zanuttigh. 2020. Unsupervised domain adaptation in semantic segmentation: a review. *Technologies* 8, 2 (2020), 35.
- [31] Kemal Toprak Uçar. 2023. A Comprehensive Guide to Error Analysis in Machine Learning. *Medium* (2023).
- [32] Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. 2024. Generalized out-of-distribution detection: A survey. *International Journal of Computer Vision* 132, 12 (2024), 5635–5662.
- [33] Kaichao You, Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I Jordan. 2019. Universal domain adaptation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2720–2729.
- [34] Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. 2022. Domain generalization: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45, 4 (2022), 4396–4415.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009