# MicroT: Low-Energy and Adaptive Models for MCUs

Yushan Huang
Imperial College London
London, UK

Ranya Aloufi
Imperial College London
London, UK

Xavier Cadet
Imperial College London
London, UK

Yuchen Zhao
University of York
York, UK

Payam Barnaghi
Imperial College London
York, UK

Hamed Haddadi
Imperial College London
London, UK

## Abstract

We propose MicroT, a low-energy, multi-task adaptive model framework for resource-constrained MCUs. We divide the original model into a feature extractor and a classifier. The feature extractor is obtained through self-supervised knowledge distillation and further optimized into part and full models through model splitting and joint training. These models are then deployed on MCUs, with classifiers added and trained on local tasks, ultimately performing stage-decision for joint inference. In this process, the part model initially processes the sample, and if the confidence score falls below the set threshold, the full model will resume and continue the inference. We evaluate MicroT on two models, three datasets, and two MCU boards. Our experimental evaluation shows that MicroT effectively improves model performance and reduces energy consumption when dealing with multiple local tasks. Compared to the unoptimized feature extractor, MicroT can improve accuracy by up to 9.87%. On MCUs, compared to the standard full model inference, MicroT can save up to about 29.13% in energy consumption. MicroT also allows users to adaptively adjust the stage-decision ratio as needed, better balancing model performance and energy consumption. Under the standard stage-decision ratio configuration, MicroT can increase accuracy by 5.91% and save about 14.47% of energy consumption.

## CCS Concepts

• **Computer systems organization** → **Embedded systems**; • **Computing methodologies** → **Artificial intelligence**; **Machine learning**.

## Keywords

Resource Constraints, Machine Learning, Low-Energy, MCUs

## 1 Introduction

The deployment of Deep Neural Networks (DNNs) on resource-constrained devices, known as Tiny Machine Learning (TinyML), has garnered widespread attention across both academia and industry [1–3]. Unlike traditional Machine Learning (ML), TinyML necessitates the reduction in size and computational demand of DNNs to align with the stringent resource limitations of the device. These devices are typically deployed in specific environments to perform localized tasks. However, due to the constrained communication and privacy [4, 5], the cloud often lacks data specific to these localized tasks [6, 7]. Therefore, conventional cloud-trained, local-deployed, and task-specific models may not be effectively sufficient for the diverse requirements of multiple local tasks [8, 9]. This disparity presents a significant challenge to the application of TinyML.

Inspired by these challenges, several methods have been developed in the community. Popular approaches can be categorized into three types: (i) Federated Learning (FL) [10], (ii) Split Learning (SL) [11], and (iii) Multitask Learning (MTL) [12]. FL designs a collaborative approach between the cloud and local resources, training a global model in the cloud and fine-tuning the complete or partial model locally to adapt to specific local tasks. SL also involves cooperative training, with the model divided between the cloud and local devices. However, both FL and SL require regular and frequent communication between the cloud and local devices [13, 14], which can be challenging to apply in some real-world situations. For example, in deep-sea environments [15], limited communication may impede joint training and the transmitting of extensive training parameters [16]. In addition, transmitting model parameters incurs excessive communication and energy costs [17, 18]. MTL learns multiple specific local tasks in the cloud and shares portions of the model's structure, but still needs a small portion of local data [19]. Despite these advancements, none of the existing methods can fully address the multi-task challenges in Microcontroller Units (MCUs), thus limiting their broader application.

Recently, novel model compression methods have been proposed to further reduce the energy consumption of running DNNs on MCUs, such as model pruning [20], model quantization [21], and knowledge distillation (KD) [22]. Model pruning involves cutting off unimportant parameters in the model. However, in the context of TinyML, where models are inherently small, pruning might significantly decrease accuracy [23]. Model quantization converts the data type of model parameters from `FLOAT32` to more efficient formats like `INT8`, offering a balance between energy savings and acceptable accuracy loss. KD involves a smaller model

learning from a larger one, thereby achieving both model compression and energy reduction. However, these methods are usually executed in the cloud before deployment on MCUs [24–26]. This approach hinders their adaptability in locally and dynamically balancing energy consumption with model performance when processing multiple local tasks, limiting their application and effectiveness on MCUs.

**Contribution.**

In this study, we introduce MicroT, a low-energy and multi-task adaptive model framework [1] designed for MCUs. MicroT incorporates a powerful and tiny feature extractor and a classifier locally trained for multiple local tasks. The feature extractor, developed by self-supervised knowledge distillation (SSKD), inherits and learns general features from a teacher model, thereby improving model performance and reducing size. The classifier, trained locally on the MCU, leverages these general features, enabling MicroT to utilize simple classifiers to achieve low-energy training and sufficient model performance. To further reduce energy consumption, MicroT employs strategies including model segmentation, joint training, and joint inference (stage-decision). Additionally, MicroT offers user-configurable stage-decision ratios and thresholds, providing a way to adaptively adjust the balance between model performance and energy cost.

To demonstrate its feasibility, we implement and evaluate MicroT on two models, three datasets, and two MCU boards. These experiments involve evaluating the model performance for multiple local tasks, the effectiveness of model segmentation and stage-decision, and the system cost. The results show that MicroT effectively improves model performance for multiple local tasks on MCUs and can achieve low-energy classifier training and model inference. MicroT is practical as it can be easily extended to various models and MCUs. Compared to unoptimized feature extractors, MicroT can improve model performance by up to 9.87%. On MCUs, compared to standard full-model inference, MicroT can save up to about 29.13% of energy cost. With the standard stage-decision ratio of 0.5, MicroT can improve model performance by 5.91% and save about 14.47% of energy cost.

## 2 Background and Related Work

In this section, we provide the necessary background to comprehend MicroT.

**Self-Supervised Learning.** Self-supervised learning (SSL) automatically generates training signals by designing proxy tasks, such as predicting a part of an image or the next word in a text, allowing algorithms to learn useful data representations from unlabeled data [27]. A major advantage of SSL is its independence from labeled data, which also enhances the model's generalization capabilities. Recently, self-supervised

models (e.g., DinoV2 [28]) have reached, and sometimes surpassed, the performance of their supervised learning counterparts, with some studies currently utilizing the more generalized embedding features extracted from these models for specific tasks [29, 30]. We believe that such general features also have application potential for MCU multi-task problems.

**TinyML and Model Compression.** TinyML focuses on resource-constrained hardware. Model compression is essential for efficient TinyML to achieve tiny models that are compatible with devices with constrained computational capacity and minimal storage. Model compression utilizes several techniques such as pruning [31], quantization [21], and KD [22]. Typically, KD aims at learning the teacher model's logit output [32], class distribution [33], or embedding features [34]. When the teacher model's extracted features are both advanced and general, learning from these embedding features equips the student model to excel with new and unseen data [35]. This ability is particularly advantageous to address the multi-task challenge on MCUs.

**Joint Training and Inference.** Joint training involves the simultaneous training of multiple models or tasks within an integrated framework. This approach facilitates the sharing of information and learning of features across different models or tasks, enhancing the overall model performance [36, 37], and has been widely used in multi-task issues [38, 39]. During the inference stage, joint inference enables collaboration among various models, integrating their outputs to obtain more comprehensive and accurate outcomes. This approach is particularly important in applications that require rapid and accurate responses. Leveraging the advantages of joint training and inference, we optimize it for the multi-task challenge on MCUs. Rather than employing two separate, independent models, our strategy involves using a large and a small model with shared parameters, thereby minimizing memory demands for MCU deployment. Special emphasis is placed on the performance of the smaller model to decrease reliance on the larger one, which contributes to reduced energy consumption during joint inference.

**Transfer Learning and Multi-Tasking.** Transfer learning involves pre-training models on data-rich, large-scale tasks to acquire universal features, followed by transferring the knowledge to specific tasks [40]. This approach is particularly effective for scenarios demanding complex tasks processing with limited resources [41, 42]. In the context of MCUs, Wu et al. [43] propose EMO, an approach for emotion recognition across various task objects in practical applications. EMO trains a feature extractor in the cloud and fine-tunes a Kmeans-based classifier on the MCU. However, our evaluation reveals that EMO underperforms in processing complex images and tasks, achieving an average accuracy of only 24.5% (details in Section. 5.4). These results may stem
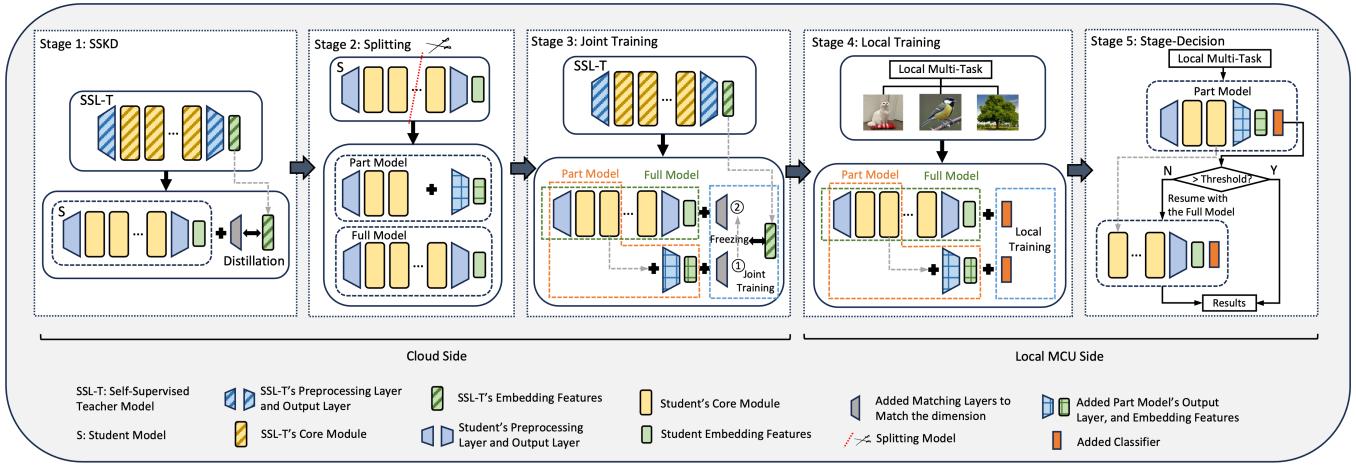
---

[1]We will release our code and design as open source upon acceptance.

**Figure 1: System Overview of MicroT**

from the feature extractor's lack of universality. Nonetheless, EMO still inspires the application of transfer learning to multi-task scenarios, as training only the classifier is efficient for the resource-constrained MCUs (details in Section. 5.6).

## 3 MicroT Design

**Assumptions.** Our study focuses on image data. We consider a realistic scenario: (i) On the cloud, there are public datasets available, but no local target task datasets. While the local MCU can access labeled data for various local target tasks. For example, on the cloud, the ImageNet and Sea Animals Image datasets [44] are available, but there are no datasets for local target tasks such as pet recognition [45]. (ii) The communication between the MCU and the cloud is unstable, as might be encountered in challenging environments like deep-sea [46]. This unstable communication makes learning approaches like FL challenging to be applied in these scenarios. (iii) The local MCUs are energy-cost-sensitive, such as systems associated with energy harvesting devices [47].

## 3.1 Overview

We propose a low-energy and multi-task adaptive model framework designed to enable MCUs to efficiently process multiple local tasks with low-energy training and inference. Fig. 1 shows the overview of this framework. MicroT's design provides the following functionalities:

**Powerful and Tiny Feature Extractor.** MicroT splits the original model into a feature extractor and a classifier. The feature extractor is adept at capturing general features, making it ideal for a range of local tasks. For its training, we employ Self-Supervised Knowledge Distillation (SSKD), a blend of SSL [28] and KD. Note that, the feature extractor is trained on the cloud, utilizing public datasets instead of specific local task data. In SSKD, we learn the embedding

features instead of the logit knowledge extracted by the teacher model (details in Section. 3.2).

**Low-Energy MCU Local Training.** MicroT facilitates low-energy classifier training on MCUs for multiple local tasks. This low energy consumption is achieved by leveraging the general and universal features extracted by the powerful and tiny feature extractor. Consequently, MicroT requires only simple classifiers to ensure adequate performance across multiple tasks. To optimize memory usage on MCUs, MicroT strategically offloads the memory burden of the feature extractor during the classifier training stage, maintaining just the essential extracted features. This approach effectively frees up additional memory space on the MCUs for classifier training (details in Section. 3.4).

**Low-Energy MCU inference.** MCUs employ a joint inference mechanism, called stage-decision, to reduce inference energy costs. The stage-decision means that samples are first processed by the part model, followed by a score calculation. If the score is below the set threshold, the full model will resume and continue the inference from the features before the part model's output layer. This mechanism enables efficient and low-energy local inference (details in Section. 3.3 and 3.5).

**Dynamic Parameter Adjustment.** On the MCUs, MicroT's stage-decision mechanism offers dynamic adjustment, enabling users to modify the stage-decision ratio according to the requirements of the balance between model performance and energy cost for various local tasks. This ratio is crucial as it dictates the proportion of samples processed only by the part model. MicroT sets a standard ratio as a benchmark. Users have the flexibility to decrease this ratio to improve model performance, or conversely, to raise it for enhanced energy efficiency (details in Section. 3.5).

## 3.2 Self-Supervised Knowledge Distillation

The design of the feature extractor has two considerations: (i) The necessity for the extracted features to possess universality for the multiple local tasks. (ii) The unavailability of local task data in the cloud for training. To address these considerations, SSKD is employed, utilizing KD to obtain a compact student feature extractor that demands less size and computational power. Concurrently, SSKD enables the student to learn general SSL-extracted features from the teacher.

Firstly, we apply SSL to the teacher. The rationale for not applying SSL directly to the student model includes: (i) SSL generally enhances greater model performance improvements in larger models, given the simpler structure of the student model, direct SSL application may only offer limited improvements [48]. (ii) Currently, there are several large models available that have already utilized SSL [28, 49, 50], these pre-trained models can be efficiently employed for specific task [29, 30].

Next, we use KD to obtain a compact student feature extractor. There are several challenges: (i) Without knowledge of the local task classes, we cannot directly learn the logits from the teacher model. (ii) The student model's output layer dimensions $P$ differs from the teacher model's $Q$. (iii) The absence of a local task dataset in the cloud and the difficulties of performing KD on the MCU directly. To address the (i) and (ii) challenges, we introduce a fully connected layer (Matching Layer) after the student model to align the output dimensions $(P, Q)$ [51]. For the (iii) challenge, we utilize a non-target public dataset in the cloud, which requires a variety of image categories and ample learning samples. This ensures that the student feature extractor can adequately inherit the teacher model's powerful representation capabilities across a broad range of image content. As shown in Fig. 1 Stage 1, we first obtain a teacher model with SSL capabilities, then use its extracted general features as learning targets for the student model. Upon integrating the matching layer into the student model, we perform KD from scratch, using the same dataset. The distillation loss, considering $p$ as the student's extracted features, $q$ as the teacher's extracted features, and $L_{\text{MSE}}$ as the Mean Squared Error, can be formulated as follows:

$$L_{\text{distill}} = L_{\text{MSE}}(p, q) \tag{1}$$

## 3.3 Model Segmentation and Joint Training

To enhance the efficiency of local inference on MCUs, MicroT employs model segmentation and joint training. It segments the feature extractor into a part model and a full model, facilitating stage-decision and joint inference on the MCU. To rationally split the feature extractor, we propose a model split fused score $F_{\text{score}}$, and adopt joint training to further optimize the performance of the part model and the joint

---

**Algorithm 1** Determine Optimal Model Segmentation Point

1: Initialize:
2: $PM \leftarrow$ Split full model into part model based on module
3: $A_{full} \leftarrow$ Compute accuracy of the full model
4: $M_{full} \leftarrow$ Compute MAC of the full model
5: $F_{max} \leftarrow 0$
6: $OptimalSplit \leftarrow 1$
7: **for** $i = 1$ to length of $PM$ **do**
8:     $A_i \leftarrow$ Compute accuracy of $PM[i]$
9:     $M_i \leftarrow$ Compute MAC of $PM[i]$
10:     **if** $i > 1$ **then**
11:         $\Delta A \leftarrow A_i - A_{i-1}$
12:         $\Delta M \leftarrow M_i - M_{i-1}$
13:         $G \leftarrow \text{Normalize}(\Delta A / \Delta M)$
14:     **end if**
15:     $R_A \leftarrow (A_{full} - A_i)/A_{full}$
16:     $R_A^{norm} \leftarrow \text{Normalize}(R_A)$
17:     $R_M \leftarrow (M_{full} - M_i)/M_{full}$
18:     $R_M^{norm} \leftarrow \text{Normalize}(R_M)$
19:     $F_{score} \leftarrow [3*(1-R_A^{norm})*G*R_M^{norm}]/[(1-R_A^{norm})+G+R_M^{norm}]$
20:     **if** $F_{score} > F_{max}$ **then**
21:         $F_{max} \leftarrow F_{score}$
22:         $OptimalSplit \leftarrow i$
23:     **end if**
24: **end for**
25: **return** $OptimalSplit$

---

model. Considering the absence of local tasks in the cloud, we utilize the public non-task dataset to calculate the optimal segmentation point. The definition of $F_{\text{score}}$ is:

$$F_{\text{score}} = \frac{3 \times (1 - R_A^{\text{norm}}) \times G \times R_M^{\text{norm}}}{(1 - R_A^{\text{norm}}) + G + R_M^{\text{norm}}} \tag{2}$$

where $F_{\text{score}}$ is the score used to evaluate the model segmentation point. $R_A^{\text{norm}}$ represents the normalized ratio of accuracy loss, which measures the decrease in part model's accuracy relative to the full model's accuracy. $R_M^{\text{norm}}$ is the normalized Multiply–Accumulate (MAC) reduction ratio, which quantifies the decrease in part model's MAC relative to the full model's MAC. $G$ is the normalized gain ratio, representing the trade-off between the increased accuracy and MAC from the part model split by the current point to the previous.

In identifying the optimal model segmentation point, we consider both model performance and computational efficiency, acknowledging the impact of computational demand on energy cost [52, 53]. For example, we use accuracy as the performance indicator for the model and MAC as the indicator for computational cost. Specifically, as the Algorithm 1 shows, the algorithm initially splits the complete model into multiple part models based on modules, and then computes the accuracy ($A$) and MAC ($M$) for each part model. Subsequently, the algorithm calculates the accuracy improvement ($\Delta A$) and MAC change ($\Delta M$) for each part model relative to the previous part model, and from this calculates the normalized gain ratio ($G$). Additionally, the algorithm evaluates

---

**Algorithm 2** Stage-Decision

---

1: Initialize:
2: $S_{train} \leftarrow$ Set of N training samples to determine threshold
3: $S_{infer} \leftarrow$ Set of M new samples for inference
4: $PartModel \leftarrow$ Function to process samples using the part model
5: $FullModel \leftarrow$ Function to process samples using the full model
6: $C \leftarrow [\,]$ Array to store confidence scores
7: $R \leftarrow [\,]$ Array to store results after stage-decision
8: $AdjustFactor \leftarrow$ Default = 1, User-defined factor to adjust the threshold
9: Determine the median confidence threshold from training samples:
10: **for** $i$ = 1 to $N$ **do**
11:      $features_{ip}, c_{ip} \leftarrow PartModel(S_{train}[i])$
12:      Append $c_{ip}$ to $C$
13: **end for**
14: Sort $C$ in ascending order
15: **if** $N$ mod 2 = 1 **then**
16:      $C_{median} \leftarrow C[\frac{N+1}{2}]$
17: **else**
18:      $C_{median} \leftarrow \frac{C[\frac{N}{2}]+C[\frac{N}{2}+1]}{2}$
19: **end if**
20: $Threshold \leftarrow C_{median} \times AdjustFactor$
21: Perform stage-decision on inference samples:
22: **for** $i$ = 1 to $M$ **do**
23:      $result_{ip}, c_{ip} \leftarrow PartModel(S_{infer}[i])$
24:      **if** $c_{ip} < Threshold$ **then**
25:          $features_{ip} \leftarrow$ Features before output layer of the $PartModel$
26:          $result_{if} \leftarrow FullModel(features_{ip})$
27:          $R[i] \leftarrow result_{if}$
28:      **else**
29:          $R[i] \leftarrow result_{ip}$
30:      **end if**
31: **end for**
32: **return** $R$

---

the accuracy reduction rate ($R_A$) and MAC reduction rate ($R_M$) for each part model relative to the complete model and normalizes these ratios. By integrating these metrics, the algorithm calculates a composite score ($F_{score}$) for each part model. After completing the calculation of the composite scores for all part models, the algorithm selects the split point with the highest composite score as the optimal split point, and obtains the part model and full model, as shown in Fig. 1 Stage 2. This approach ensures that the selected optimal model segmentation point does not significantly affect the model performance while ensuring lower computational complexity, thereby reducing energy consumption.

When computing $F_{score}$, we utilize another available public non-task dataset in the cloud, which is different from the dataset used in SSKD. The main reasons for this are as follows: (i) The dataset used in SSKD requires a wide range of categories and ample samples, generally resulting in a large dataset. Using this dataset would reduce the efficiency of score computation; (ii) Using a dataset different from the one used in SSKD increases the credibility of the score, as the new public non-task dataset, like the MCU local task dataset, has not been learned before.

Following the model segmentation, we implement joint training. As shown in Fig. 1 Stage 3, this process involves augmenting the part model with an additional output layer and a matching layer to align the feature dimensions extracted by the teacher. Subsequently, the models undergo SSKD joint training, building upon the initial feature extractor. We first independently train the part model, and then delete the added layers and freeze its parameters to further train the full model. This joint training can enhance the part model's independence and performance, enabling the part model to process a larger proportion of samples more effectively, thereby reducing the dependency on the full model and saving energy. Moreover, joint training ensures cohesive integration of the part model within the full model. This integration facilitates parameter sharing between the two, thereby preventing the need for additional memory allocation when deploying these two models on MCUs. Due to model segmentation and the part model, this joint training approach may lead to some accuracy loss in the full model. However, the performance improvement from the part model can also transfer to the full model, thereby compensating for this accuracy loss. Our experiments also demonstrate that this loss in accuracy is acceptable (details in Section 5.3). The cooperation between the models, improved through joint training, thus optimizes both performance and energy cost, crucial for applications on resource-constrained devices.

## 3.4 Classifier Training

The MCU obtains the feature extractor (with INT8 quantization and only about 0.63% accuracy loss) from the cloud, and locally builds and trains classifiers to process multiple tasks, as shown in Fig. 1 Stage 4. The general features from the feature extractor enable the use of a basic classifier, which can achieve desirable model performance while mitigating memory constraints on the MCU. Classifier development and training on the MCU are executed in the $C$ programming language. For example, constructing a 2-layer Neural Network (NN) classifier entails several key steps: 1) Defining the classifier's architecture, including layer count, neurons per layer, input and desired output of the NN, layers (excluding the input layer), and the learning rate; 2) Structuring parameter matrices, encompassing weight matrix, bias array, output array, and error (the partial derivative of the total error with respect to the weighted sum); 3) Defining activation functions; 4) Implementing a function to load datasets and models; 5) Creating a model construction function for memory allocation and model establishment; 6) Developing a training function, which includes forward and backward propagation processes; 7) Instituting a function to save the model; 8) Developing a function to release memory. To alleviate memory load during classifier training, the MCU offloads

the memory used by the feature extractor prior to training, retaining only the extracted features. Given the MCU's limited memory capacity, the batch size is set to one. Since both the part model and full model are deployed on the MCUs and may process multiple tasks, there will be multiple classifiers, each of them trained separately.

## 3.5 Stage-Decision

Upon completing local training, MicroT implements stage-decision to further reduce energy cost during inference, as shown in Figure 1 Stage 5. Note that, in local training and stage-decision, data preprocessing is also necessary, including resizing and normalization of input images, which should match the preprocessing steps in the cloud. Stage-decision involves initially processing samples by the part model and calculating confidence scores. If the score falls below the set threshold, the full model will resume and continue inference from the stage preceding the part model's output layer. Setting an appropriate threshold is crucial, as it determines the stage-decision ratio, the ratio of samples processed only by the part model. The algorithm of stage-decision is detailed in Algorithm 2.

This threshold is established using confidence scores from the part model. The median of $C$, denoted as $C_{median}$, is selected as the threshold. Consequently, approximately 50% of the samples are processed only by the part model, establishing a stage-decision ratio of 0.5. Likewise, employing the 1/4 and 3/4 quantiles of $C$ as thresholds results in stage-decision ratios of 0.25 and 0.75, respectively. Our experiments (details in Section. 5.5) show that the threshold can be determined with just a few samples, avoiding excessive resource consumption on the MCU. Moreover, a ratio of 0.5, based on the median confidence score, balances model performance with energy efficiency. Therefore, this ratio is adopted as the standard stage-decision ratio (details in Section 5.4).

MicroT enables diverse balances between model performance and energy consumption by offering variable thresholds and ratios. It offers the *AdjustFactor* function, allowing users to dynamically modify the threshold according to their specific needs. In practical scenarios, users can customize their balance of model performance and energy efficiency. By raising the threshold, and consequently decreasing the stage-decision ratio, users can enhance model performance. Alternatively, lowering the threshold increases the ratio, leading to reduced energy consumption.

## 4 Implementation & Evaluation Setup

In this section, we first introduce the implementation of the MicroT system (Section 4.1), and then detail how we assess its performance on various DNN models and datasets (Section 4.2) using different metrics (Section 4.3).

### 4.1 MicroT Prototype

We deploy MicroT by STM32-X-CUBE-AI [54], and modify the original $C$ code to achieve the preprocessing, feature extractor implementation, classifier inference and training, dynamic threshold adjustment, and stage-decision. We run MicroT on STM32H7A3ZI and STM32L4R5ZI respectively. The STM32H7A3ZI is equipped with an ARM Cortex-M7 core, 2MB of Flash memory, and 1.18MB of RAM, with a maximum frequency of 280MHz. The STM32L4R5ZI features an ARM Cortex-M4 core, 2MB of Flash memory, and 640KB of RAM, with a maximum frequency of 120MHz. We utilize Monzoon High Voltage Power Monitor [55] to set the MCU power supply voltage to $1.9V$ [6]. We set the on-device inference and training with the MCU board frequency at 120MHz. Fig. 2 shows the overview of the MicroT prototype.

### 4.2 Models and Datasets

We select DinoV2 [28] as the teacher model due to its ability to learn general image features from a substantial collection of public image datasets and also benefits from the Vision Transformer model and SSL. For the student model, to balance accuracy, memory footprint, and energy consumption, we utilize ProxylessNAS_w0.3 [56] and MCUNet_int3 [57], which are the commonly used lightweight CNN and module-based models in the related literature [58, 59]. For MCUNet and ProxylessNAS, both architectures are fundamentally characterized by Mobile Inverted Bottleneck Convolutions (MBConv), as Fig. 3 shows. Note that, we remove the original classifier layers from ProxylessNAS and MCUNet, retaining only the remaining structures as feature extractors. Subsequently, we integrate either Logistic Regression (LR) or a 2-layer Neural Network (NN) as classifiers. The input dimension of the classifier aligns with the output dimension of the teacher model (DinoV2 = 384), suppose the number of output categories is $V$. The LR classifier follows a $(384, V)$ structure, while the 2-layer NN follows a $(384, 128, V)$ structure.

For SSKD on the initial feature extractor, we train the model from scratch with 0.01 learning rate, SGD optimizer, and 50 epochs. After splitting the initial feature extractor, in the joint training, we further jointly train the part model and full model with 0.005 learning rate, maintaining the same optimizer and epoch count. In addition, we consider the impact of the input image resolution, and utilize the 224 and 128 resolutions, which are practical for MCUs [60, 61]. To avoid complex retraining for 128 resolution, we fine-tune the model from its 224 resolution state, using 0.005 learning rate, SGD optimizer, and 50 epochs. On MCUs, classifier training and system cost evaluation are conducted using SGD optimizer, cross-entropy loss function, Softmax output layer activation function, RELU hidden layer activation function, 0.01 learning rate, and 200 epochs, without momentum [61] (details in

**Figure 2: MicroT Prototype Overview**



**Figure 3: Main Structure of ProxylessNAS and MCUNet**

Section. 5.2). Some examples of common abbreviations and their meanings in these experiments are as follows:

- MCUNet_r224: MCUNet without the original classifier, and image resolution is 224.
- MCUNet_r224_LR: MCUNet with LR classifier added after removing the original classifier, and 224 resolution.
- MCUNet_r224_part: The part model of MCUNet after model segmentation, 224 resolution, no original classifier.
- MCUNet_r224_part_NN: The part model of MCUNet with NN classifier added after removing the original classifier, and 224 resolution.
- MCUNet_r224_sd_0.5: The part model and full model with added classifier and execute stage-decision, the stage-decision ratio is 0.5, and 224 resolution.

The datasets are categorized into cloud datasets and local datasets. For the cloud datasets, we utilize ImageNet to train feature extractors, as it provides sufficiently rich samples for the student to learn from the teacher. We also utilize the Sea Animals Image (Sea) dataset [44] to identify the best model segmentation point, which is similar to the local dataset as it is unlearned by the student. For the local datasets, we utilize several datasets to represent multiple local tasks, including The Oxford-IIIT Pet (Pet) [45], CUB-200-2011 Bird (Bird) [62], and PlantCLEF 2017 (Plant) [63].

The Pet dataset [45] contains 3,680 images and 37 distinct pet species. The bird dataset [62] contains a collection of 11,788 images in 200 distinct categories. For the Plant dataset [62], we select a subset of it based on image quantities per category, containing 20 plant categories with a total of 11,660 training images. We randomly split the training and test datasets by the ratio of 0.7 and 0.3.

### 4.3 Performance Metrics

Our evaluation of MicroT focuses on two main aspects: (i) ML model performance on multiple local tasks. (ii) The system cost on the MCUs. All the results are the average values of ten repeat experiments.

**Model Performance.** Accuracy of MicroT models with different configurations on the multiple local datasets.

**Table 1: Overall performance of MicroT with varying stage-decision ratios. The red color indicates a decrease in accuracy, whereas green signifies improvements in accuracy and energy efficiency. MicroT's standard ratio is set at 0.5.**

| Model | Acc. | Acc. Improvement | Energy Cost | Energy Saving |
|---|---|---|---|---|
| Baseline | 51.86% | - | 5.39mJ | - |
| MicroT_0 | 61.73% | 9.87% | 5.39mJ | - |
| MicroT_0.25 | 60.17% | 8.31% | 5.00mJ | 7.21% |
| **MicroT_0.5** | **57.77%** | **5.91%** | **4.61mJ** | **14.47%** |
| MicroT_0.75 | 51.26% | -0.60% | 4.21mJ | 21.89% |
| MicroT_1 | 41.14% | -10.72% | 3.82mJ | 29.13% |

**Table 2: Results of training classifiers on MCU and GPU with varying batch sizes and partial samples from the Pet dataset.**

| Batch Size | Feature Extractor | Pet | | |
|---|---|---|---|---|
| | | LR | NN | Avg. |
| 128 (GPU) | MCUNet_r128_full | 71.47 | 76.87 | 74.17 |
| | Proxy_r128_full | 64.26 | 67.86 | 66.06 |
| 1 (MicroT) | MCUNet_r128_full | 69.74 | 74.48 | 72.11 |
| | Proxy_r128_full | 62.95 | 65.78 | 64.37 |

**System Cost.** We monitor the efficiency of local training and inference on MCUs, and measure the following costs: (i) Runtime ($s$): The time taken by the MCU to process operations, reported and calculated by reading the current time from the MCU's internal clock. (ii) Memory Usage ($MB$ or $KB$): The maximum Flash memory and RAM usage, obtained by STM32-CUBE-IDE [64]. (iii) Energy Consumption ($mJ$): The energy cost to perform operations, measured by Monsoon High Voltage Power Monitor [55] with $50Hz$ sampling rate. We utilize Monsoon [55] to set the input voltage ($U$) to 1.9V [6] and measure the time ($t$) and average current ($I$). Then we calculate the average power ($P = UI$) and average energy consumption ($E = Pt$). To obtain stable values, we only utilize the current recorded after running for 2 $min$.

## 5 Evaluation Results

In this section, we present the experimental evaluation of MicroT aiming to answer a set of key questions.

### 5.1 Overall Performance

Table. 1 shows an overall performance comparison between MicroT with varying stage-decision ratios and the baseline, including accuracy and energy cost. The model accuracy here represents average values taken from three datasets, two models, and two resolutions. The energy cost is averaged over two MCU boards, and the Plant [63] and Bird [62] datasets, as they have the smallest and largest number of classes. The baseline refers to the model employing an unoptimized feature extractor and standard full model inference. Analysis reveals that, across different ratios, MicroT's accuracy and energy consumption vary, achieving a maximum increase of 9.87% in accuracy and 29.13% in energy saving

compared to the baseline. The ratio is set at 0.5 by default, achieving a balance between accuracy and energy savings, leading to a 5.91% improvement in accuracy and a 14.47% reduction in energy usage. Table. 1 summarizes the results from all experiments, with details in the following section.

## 5.2 Does MicroT's Local Training Match Cloud-Level Accuracy?

For training the classifier on MCUs, due to memory limitations, the batch size is set to one. However, this reduces the efficiency of the experiments, as it does not utilize hardware parallelism. Therefore, at the start of the experiments, we investigate the performance gap between classifier training on MCUs with a single-batch and on the GPU with a batch size of 128, to see if the latter can be used as an approximation to improve the efficiency of the experiments. The training settings on the MCU and cloud are kept consistent, as illustrated in Section. 4.2. The reason for not using momentum is that for a single-batch setup, momentum does not benefit model performance and can increase memory usage [61]. The classifier training on the MCUs is based on $C$ language, while on the GPU is based on *Python*. Due to the low efficiency of single-batch training, we randomly select 20 images with 128 resolution from each class in the Pet dataset (a total of 740 images), the results are shown in Table. 2.

We find that, keeping the classifier training configuration consistent, local training on the MCU with single-batch and training on the GPU with a 128 batch are in similar accuracy. This allows us to use the results of batch training on the GPU for evaluation. Therefore, we default to reporting the GPU classifier training results with based on the aforementioned classifier training configuration, unless otherwise stated.

## 5.3 How Much does MicroT Improve Performance in Segmented Models?

MicroT divides the original model into a a part model and a full model to make stage-decision, thereby reducing energy consumption. A superior performance of the two models can potentially improve the performance of stage-decision. Therefore, we evaluate the performance improvement of the segmented models.

The first step is to determine the optimal model segmentation point by the model segmentation fused score (as illustrated in Section. 3.3), considering four aspects: the comparison of accuracy and MAC between the part model up to the current module and the previous module, and the comparison of accuracy and MAC between the part model the full model. We first train the ProxylessNAS and MCUNet feature extractors on the ImageNet dataset using SSKD. Subsequently, to find the optimal segmentation point, we utilize

the Sea dataset [44] which is also public on the cloud, different from the SSKD dataset (ImageNet), and has never been learned by the feature extractors (as illustrated in Section. 3.3). Note that, we do not want the part model to be too small to have a poor performance, also do not want it to be too close to the full model, as this would result in minimal optimization for the system cost. Therefore, for MCUNet, we only analyze modules 4 to 14, and for ProxylessNAS, we analyze modules 6 to 17. Since the output dimensions of different part models are different, we also add a matching layer. The model segmentation fused scores of MCUNet and ProxylessNAS are shown in Fig. 4.

We can observe that for 128 and 224 resolutions, the optimal segmentation point for the MCUNet is at the 9th module, while the ProxylessNAS is at the 14th module. In these four models, the MAC generally increases steadily with the depth of the network. However, at the optimal model segmentation points, there is a significant improvement in accuracy, indicating that the model has acquired essential information and extracted key features at this depth and structure. We believe this demonstrates that the current model depth offers the best cost-performance ratio when considering model performance and computational complexity. Furthermore, regardless of changes in resolution, the optimal model segmentation points are at the same depth for the same model, suggesting that for a given model, the optimal segmentation point may have a certain stability, and the main factor influencing this point is the model structure.

Subsequently, we split the MCUNet at its 9th module and the ProxylessNAS at its 14th module, resulting in part models and full models. These models are then further trained on ImageNet by SSKD and joint training. After adding the additional output layer and matching layer to the part model, we first train it, and then delete the added layers and freeze its parameters to further train the full model. We compare the performance of models with and without MicroT to evaluate the performance improvement that MicroT brings to the segmented models. The results are shown in Fig. 5.

We can observe that MicroT can improve varying degrees of model performance on the part model and full model. Specifically, MicroT leads to an average increase of 9.87% in full models' accuracy and 12.56% in part models' accuracy. We believe that the reason for the improvement in the accuracy of part models is due to the SSKD and joint training of MicroT. SSKD helps the part model and full model learn more general embedding features from the teacher, which enhances the universality of the extracted features, thereby enabling stable performance across different local tasks. Additionally, in joint training, we pay extra attention to the part model training to improve its independence. However, we note that the improvement in accuracy of the full model is slightly less than that of the part model. We believe this is the
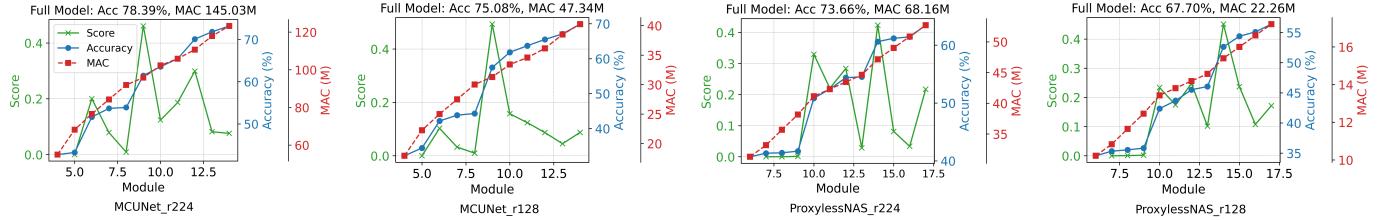
**Figure 4: The model segmentation fused score of MCUNet and ProxylessNAS. The X-axis represents the part model up to which module, and Y-axis represents the fused score (green), accuracy (blue), and MAC (red).**
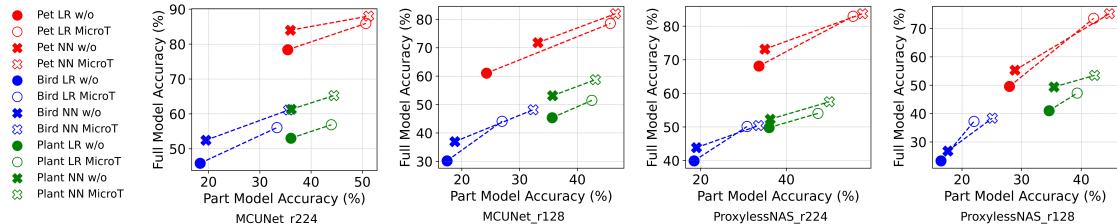


**Figure 5: The performance improvement by MicroT to the segmented models. The dashed connecting lines show the changes in model performance produced by the part model and full model with and without MicroT.**

result of a combination of several factors: (i) In joint training, the addition of the model segmentation point and the part model may lead to some degree of information loss. (ii) The accuracy improvement of the part model is also beneficial for the full model, as it can feed back and propagate to the full model. Nevertheless, the information loss in the full model is still acceptable considering the performance improvements of both the full model and part model.

## 5.4 How does the ML Performance of MicroT Compare to State-of-the-Art?

In these experiments, we compare the ML performance of MicroT with state-of-the-art (SOTA). For the selection of SOTA, Lin et al. [61] propose a method for on-device training, but this method is not designed for multi-task issues. Moreover, training the entire model on MCUs for local tasks leads to higher latency and longer model convergence time, which is not efficient for practical deployment and application. Therefore, we do not consider it as a comparative method in our research. Wu et al. [43] propose EMO, which involves training a CNN-based feature extractor and a Kmeans-based classifier in the cloud, and then fine-tuning the classifier on the MCUs based on local tasks. The classifier in EMO is an unsupervised classifier, which mainly adjusts the cluster point centers and the range of class distances on the MCUs. However, we find that EMO does not perform well on complex datasets like Pet [45], Bird [62], and Plant [63], with an average accuracy of only 24.5%. This may be caused by the simple CNN feature

extractor and unsupervised classifier. Therefore, we maintain the key method of EMO, modify the feature extractor to MCUNet and ProxylessNAS, change the classifier to LR or NN, and regard this as the baseline. MicroT with different stage-decision ratios is then compared with the baseline, the results are shown in Fig. 6.

We can observe that when the stage-decision ratios are 0, 0.25, and 0.5, the model performance consistently surpasses the baseline, with average improvements of 9.87%, 8.31%, and 5.91%. This indicates that at these ratios, the stage-decision between the part model and the full model effectively enhances ML performance. Lower ratios mean that the full model processes more low-confidence samples from the part model, thereby increasing accuracy but costing more energy. Conversely, at a ratio of 0.75, only some MicroT performance exceeds the baseline. This is because, at this threshold, more samples are processed only by the part model, speeding up processing while leading to a decline in ML performance. The different ML performances under various ratios emphasize the importance of threshold configuration in balancing model performance and computational efficiency. Appropriate threshold settings can ensure sufficient accuracy while reducing computational resource consumption, which is especially critical in resource-limited environments. We discuss experiments and analysis regarding computational efficiency and energy consumption in detail in Section. 5.6.

We also find that the decrease in model performance with increasing stage-decision ratio is not linear. Specifically, as
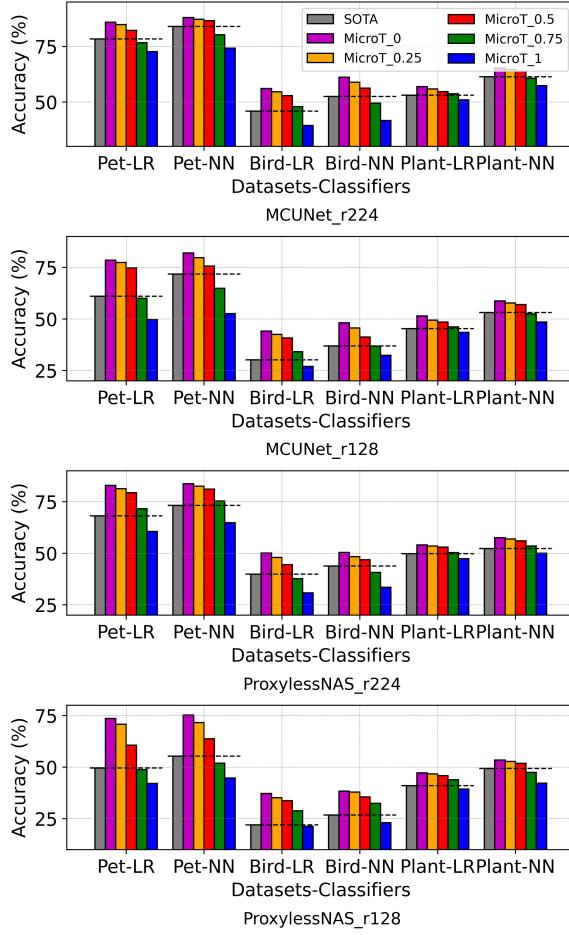
Figure 6: The ML performance of MicroT with different stage-decision ratios. The dashed lines show the accuracy of the baseline. The numbers in the legend represent the ratio.
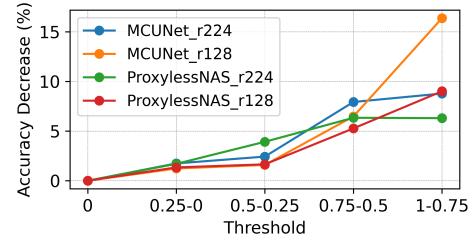


Figure 7: Model average performance decreases from the current stage-decision ratio to the previous ratio, with a fixed ratio increment of 0.25.
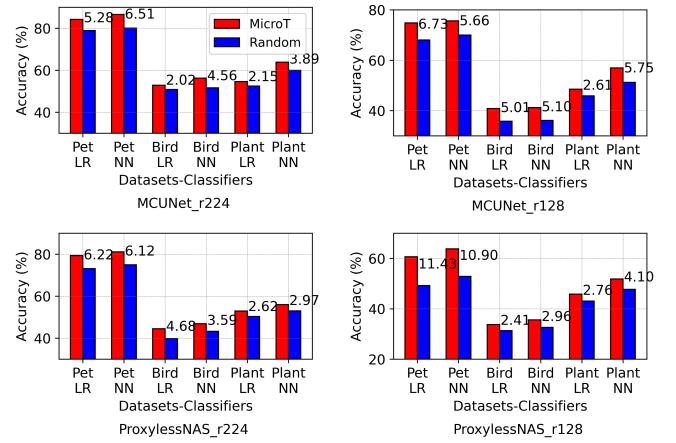


Figure 8: Achieve a stage-decision ratio of 0.5 by random selection and using the confidence median (MicroT).

shown in Fig. 7, with a fixed ratio increment of 0.25, we calculate the performance decrease with the previous ratio (for example, ratio 0.25 compared to threshold 0, ratio 0.5 compared to ratio 0.25). We notice that the decline in model performance exhibits a non-linear trend. This non-linearity may stem from several reasons: (i) The model exhibits noticeable differences in processing samples of varying difficulty, performance drops are particularly pronounced when the less capable part model processes more samples. (ii) The distribution of sample difficulties in the dataset might be non-uniform, leading to more pronounced performance decreases at certain ratio points. (iii) Although we determine the ratio by confidence scores, there may still be other factors that affect the model's handling of samples with varying difficulty. It's possible that the challenging samples are not always accurately identified and passed to the full model. At

a ratio of 0.5, we can observe a slower rate of performance decrease, indicating a reduced dependency on the full model while maintaining performance, thus lowering energy consumption. This setting not only considers the model performance in processing samples of different difficulties, but also the energy efficiency requirements on resource-constrained MCUs, making it a safe standard ratio setting in balancing performance and energy consumption.

## 5.5 How Many Samples Do We Need to Determine the Confidence Thresholds?

The confidence threshold dictates the stage-decision ratio, determining the proportion of samples processed only by the part model. The threshold setting is based on the confidence score sequence from the part model. Specifically, after the part and full models have been trained on MCUs, the part model processes $N$ samples and generates a confidence score sequence. We then select the median of this sequence to achieve the standard stage-decision ratio of 0.5. Ideally, selecting the median of the entire test samples' confidence

**Table 3: Median value, stage-decision ratio, and model accuracy with different number of samples. The *N* represents the number of samples used to calculate the median value. The *Med.* represents the median value, the *Ratio* represents the stage-decision ratio, the *Acc.* represents the model accuracy (%).**

| N | MCUNet_r224_LR | | | | | | | | | ProxylessNAS_r128_NN | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Pet | | | Bird | | | Plant | | | Pet | | | Bird | | | Plant | | |
| | Med. | Ratio | Acc. | Med. | Ratio | Acc. | Med. | Ratio | Acc. | Med. | Ratio | Acc. | Med. | Ratio | Acc. | Med. | Ratio | Acc. |
| 5 | 0.688 | 0.486 | 82.35 | 0.408 | 0.456 | 54.35 | 0.467 | 0.447 | 54.72 | 0.183 | 0.444 | 64.89 | 0.183 | 0.423 | 35.86 | 0.298 | 0.456 | 51.92 |
| 10 | 0.674 | 0.504 | 82.12 | 0.400 | 0.465 | 52.45 | 0.439 | 0.484 | 54.15 | 0.168 | 0.519 | 64.43 | 0.154 | 0.516 | 35.80 | 0.273 | 0.530 | 51.79 |
| 20 | 0.676 | 0.501 | 82.17 | 0.384 | 0.491 | 52.76 | 0.439 | 0.484 | 54.15 | 0.167 | 0.525 | 64.39 | 0.153 | 0.519 | 35.74 | 0.280 | 0.512 | 51.85 |
| 40 | 0.679 | 0.497 | 82.26 | 0.386 | 0.488 | 52.74 | 0.436 | 0.491 | 54.49 | 0.170 | 0.507 | 64.16 | 0.161 | 0.492 | 35.57 | 0.282 | 0.504 | 51.85 |
| All | 0.677 | 0.500 | 82.21 | 0.378 | 0.500 | 52.85 | 0.430 | 0.500 | 54.61 | 0.172 | 0.500 | 63.76 | 0.159 | 0.500 | 35.57 | 0.283 | 0.500 | 51.82 |

**Table 4: The system costs during local training. *t* represents time (*s*), *P* represents average power (*mW*), and *E* represents energy consumption (*mJ*).**

| Object | STM32H7A3ZI | | | STM32L4R5ZI | | |
|---|---|---|---|---|---|---|
| | t | P | E | t | P | E |
| MCUNet_r224_full_NN | 2.86 | 3.90 | 11.16 | 5.19 | 2.79 | 14.48 |
| MCUNet_r224_part_NN | 2.18 | 3.91 | 8.52 | 3.85 | 2.70 | 10.40 |
| MCUNet_r128_full_NN | 1.25 | 3.81 | 4.76 | 2.15 | 2.72 | 5.85 |
| MCUNet_r128_part_NN | 1.02 | 3.76 | 3.83 | 1.75 | 2.51 | 4.39 |
| ProxylessNAS_r224_full_NN | 1.58 | 3.87 | 6.11 | 3.16 | 2.58 | 8.15 |
| ProxylessNAS_r224_part_NN | 1.32 | 3.82 | 5.04 | 2.56 | 2.53 | 6.48 |
| ProxylessNAS_r128_full_NN | 0.83 | 3.54 | 2.94 | 1.54 | 2.41 | 3.71 |
| ProxylessNAS_r128_part_NN | 0.75 | 3.47 | 2.60 | 1.31 | 2.38 | 3.12 |

**Table 5: The energy costs of MicroT under different stage-decision ratios. *E* represents energy consumption (*mJ*), *ES* represents the energy saving rate compared with the baseline (%), *Avg.* represents the average energy saving rate (%).**

| Model | STM32H7A3ZI | | STM32H7A3ZI | | Avg. |
|---|---|---|---|---|---|
| | E | ES | E | ES | |
| MCUNet_r224_sd_0 | 9.52 | - | 12.42 | - | |
| MCUNet_r128_sd_0 | 3.12 | - | 4.03 | - | |
| ProxylessNAS_r224_sd_0 | 4.64 | - | 6.19 | - | - |
| ProxylessNAS_r224_sd_0 | 1.45 | - | 1.78 | - | |
| MCUNet_r224_sd_0.25 | 8.87 | 6.83 | 11.48 | 7.57 | |
| MCUNet_r128_sd_0.25 | 2.89 | 7.37 | 3.70 | 8.19 | |
| ProxylessNAS_r224_sd_0.25 | 4.33 | 6.68 | 5.75 | 7.11 | 7.24 |
| ProxylessNAS_r128_sd_0.25 | 1.34 | 7.59 | 1.65 | 7.30 | |
| MCUNet_r224_sd_0.5 | 8.22 | 13.66 | 10.55 | 15.06 | |
| MCUNet_r128_sd_0.5 | 2.66 | 14.74 | 3.36 | 16.63 | |
| ProxylessNAS_r224_sd_0.5 | 4.02 | 13.36 | 5.31 | 14.22 | 14.47 |
| ProxylessNAS_r128_sd_0.5 | 1.23 | 15.17 | 1.51 | 15.17 | |
| MCUNet_r224_sd_0.75 | 7.57 | 20.48 | 9.61 | 22.62 | |
| MCUNet_r128_sd_0.75 | 2.42 | 22.44 | 3.03 | 24.81 | |
| ProxylessNAS_r224_sd_0.75 | 3.71 | 20.04 | 4.87 | 21.32 | 21.89 |
| ProxylessNAS_r128_sd_0.75 | 1.12 | 22.76 | 1.38 | 22.74 | |
| MCUNet_r224_sd_1 | 6.92 | 27.31 | 8.67 | 30.19 | |
| MCUNet_r128_sd_1 | 2.19 | 29.81 | 2.69 | 33.25 | |
| ProxylessNAS_r224_sd_1 | 3.40 | 26.72 | 4.43 | 28.43 | 29.13 |
| ProxylessNAS_r128_sd_1 | 1.01 | 30.34 | 1.24 | 30.34 | |

confidence scores for all test samples is not feasible. Therefore, to determine if it is efficient to determine the threshold locally, we utilize various quantities of samples to calculate their medians, and analyze these medians' corresponding stage-decision ratios and model performance. We test the MCUNet at 128 resolution with an LR classifier and ProxylessNAS at 224 resolution with an NN classifier, as these two models contain all types of experimental feature extractors, classifiers, and resolutions. The results are shown in Table. 3.

As Table. 3 shows, we find that as the number of samples increases, the median, stage-decision ratio, and model performance increasingly approximate those obtained by all test samples. However, even when using only five samples, these differences remain minor and acceptable, with an average difference of 0.048 in median and 0.38% in accuracy. This means that on MCUs, processing just five samples is sufficient to quickly and accurately determine the threshold.

Random selecting 50% samples is also a way to achieve the 0.5 stage-decision ratio. Therefore, regarding this random selection method as the baseline, we conduct a comparative experiment. For the baseline, we randomly select half of the samples to be processed only by the part model, the other half are further processed by the full model. For MicroT, the decision on which samples should be forwarded to the full model for further processing is based on the median of the confidence scores. As Fig. 8 shows, compared to a strategy relying on random selection, MicroT can improve the average accuracy by 4.83%. This finding suggests using the median of confidence scores as the sample selection strategy can enhance the model performance.

## 5.6 What is the System Cost of MicroT?

MicroT has two stages on the MCU: classifier training and stage-decision. In these experiments, we measure the maximum memory usage, time, average power, and energy consumption under different model and resolution conditions. We select the Plant [63] and Bird dataset [62] to measure the system cost, as these two datasets have the smallest and largest number of classes. To efficiently analyze the system cost, we use the average of these two datasets.

sequence can set the stage-decision ratio to 0.5, meaning 50% of the samples are processed only by the part model. However, in real-world applications, storing and computing the

Firstly, we measure the memory usage of MicroT by STM32-CUBE-IDE [64]. To simplify the process and show that MicroT can meet the memory requirements of the experimental MCU boards: (i) We only analyze the local training stage, since the memory usage of the local training stage is higher than the state-decision stage. (ii) We only analyze the full model, since the memory usage of the full model is higher than the part model. (iii) We only analyze the 2-layer NN classifier, since the memory usage of the NN classifier is higher than the LR classifier. For maximum RAM usage, MCUNet is $614.40KB$ (resolution 224) and $221.34KB$ (resolution 128), ProxylessNAS is $624.64KB$ (resolution 224) and $225.28KB$ (resolution 128). For flash memory usage, MCUNet uses about $0.91MB$ to store the model (feature extractor $0.67MB$, classifier $0.24MB$), and ProxylessNAS uses about $0.70MB$ (feature extractor $0.46MB$, classifier $0.24$ $MB$). These experimental results indicate that, from the perspective of memory usage, MicroT can meet the memory requirements of MCUs and can store several classifiers for multiple local tasks.

We utilize Monsoon [55] to measure and calculate the time, average power, and energy consumption (as illustrated in Section. 4.3). For efficient analysis: (i) We only analyze the local training stage, since the system costs of it are higher than the stage-decision stage. (ii) We only analyze the NN classifier, since the system costs of it are higher than the LR classifier. The results are shown in Table. 4.

We find that MicroT has fast processing speeds and low energy consumption during the local training stage. This is because only the classifier requires training, and it has a simple structure. This simple structure benefits from the general extracted features. We can also find that although the power of STM32H7A3ZI is higher, its processing speed is faster, resulting in overall lower energy consumption compared to STM32L4R5ZI. For example, with the MCUNet_r224_full_NN, the processing time on STM32H7A3-ZI is $2.86s$ with a power of $3.90mW$ and energy consumption of $11.16mJ$, whereas on STM32L4R5ZI, the processing time is $5.19s$, power is $2.79mW$, and energy consumption is $14.48mJ$. This could be attributed to the more efficient performance of the ARM Cortex-M7 core used in STM32H7A3ZI compared to the ARM Cortex-M4 core in STM32L4R5ZI. In terms of models, ProxylessNAS shows better energy and time efficiency than MCUNet. For example, on STM32H7A3ZI, the ProxylessNAS_r224_full_NN takes $1.58s$ and consumes $6.11mJ$, while the MCUNet_r224_full_NN takes $2.86s$ and consumes $11.16mJ$. Compared to the full model, the part model reduces training time and energy consumption. For example, the MCUNet_r224_full_NN and MCUNet_r224_part_NN on STM-32H7A3ZI require $2.86s/11.16mJ$ and $2.18s/8.52mJ$, respectively. This is because the local training includes feature extractor inference and classifier training, and the part model makes the feature extractor inference more efficient

and energy-saving. This energy-saving becomes more evident during the stage-decision inference stage, as the energy consumption for classifier inference is lower than during training, which can improve the proportion of the feature extractor's energy cost in the overall energy cost. These results demonstrate the effectiveness of the part model and the feasibility of MicroT's low-energy local training.

Subsequently, we analyze the energy consumption during the stage-decision with different ratios. When the ratio is set to 0, meaning all samples are processed by the full model, we regard this condition as the baseline to assess how much energy MicroT can save. Since Section. 5.5 shows that the threshold calculated by five samples is not much different compared to the exact threshold, we use the exact threshold here to determine the ratio. The results are shown in Table. 5. We find that as the stage-decision ratio increases, MicroT achieves greater energy savings, with a minimum of 7.24% (ratio 0.25) and a maximum of 29.13% (ratio 1). For the standard ratio 0.5, MicroT saves 14.47% of energy. This indicates that MicroT can adaptively and effectively save energy while enhancing the performance of multi-task models on MCUs.

## 6 Discussion and Future Work

**Key Findings.** Our evaluation results showed that:

• SSKD can improve the model performance on multiple local tasks, with up to a 9.87% increase in accuracy compared with the SOTA.

• The model segmentation fused score and joint training can effectively optimize the model performance of the part and full model, laying the foundation for stage-decision.

• MicroT can achieve low energy local training on MCUs, with the lowest energy consumption at $2.60mJ$ for the part model, and $2.94mJ$ for the full model on the STM32H7A3ZI.

• The stage-decision ratio determines how many samples are processed only by the part model. And the stage-decision can effectively save energy consumption, up to 29.13% (ratio of 1). Under the standard ratio of 0.5, MicroT can save about 14.47% energy and improve 5.91% model performance.

• Determining the stage-decision ratio by the confidence threshold is efficient and accurate. MicroT only needs five local samples on the MCUs to obtain the confidence threshold. At the standard stage-decision ratio of 0.5, compared to randomly selecting 50% of the samples, MicroT can improve model performance by 4.83%.

**Model Architectures.** MicroT is adaptable to various models. In addition to the models used in the experiments, SqueezeNet's Fire Module [65], MobileNet's Inverted Residual Structure [66], models from Neural Architecture Search, and LSTM (regard each LSTM-layer as one module) can facilitate model segmentation by the fused score, then further apply the following steps of MicroT.

**Accelerating MCU Inference and Training.** The efficiency of MicroT's local inference and training on MCUs can be further improved. Integrating Sparse Updates [61], which minimize memory usage, could accelerate local training and inference within MicroT. This involves assessing parameter significance in both the part model and full model's shared sections. In the part model, only highly important parameters would be operational, while the full model would utilize all parameters. Such an approach could substantially decrease runtime, power, and energy usage, thereby improving MicroT's overall energy efficiency.

**Unsupervised Classifier.** MicroT relies on labeled data for classifier training on MCUs, but exploring unsupervised classifiers presents an interesting direction. While our results show that EMO's [43] unsupervised classifier struggles with complex datasets, its potential merits further investigation. The unsupervised classifier's poor performance on complex datasets may be caused by the vast number of categories and the dense feature space. Nonetheless, refining the differentiation of categories in this feature space, perhaps through methods like contrastive learning [50], could enhance the classifier's performance. Such advancements would broaden MicroT's applicability.

**Multimodal Feature.** MicroT focuses only on image data. However, the potential to handle diverse signal sources, like environmental sensor signals [67, 68], is essential for MCUs. Recent studies have investigated data and feature fusion techniques for multimodal data analysis on edge devices [69, 70]. MicroT could integrate its feature extractor with these fusion methods to map multimodal data into a unified feature space. This expansion not only widens MicroT's utility but also enhances its model robustness.

**Multiple Model Segmentation**. MicroT splits the initial feature extractor into a part model and a full model to enable low-energy training and inference on MCUs. In fact, MicroT can further improve the granularity of model segmentation by dividing the initial feature extractor into multiple part models, thus achieving higher adaptability.

## 7  Conclusion

In this paper, we introduced MicroT, a practical, low-energy, and open-source framework to address the multi-task challenge of MCUs. MicroT leverages a powerful and tiny feature extractor by SSKD, along with classifiers that can be specifically trained for local multiple tasks. To further reduce energy consumption, MicroT utilizes model segmentation, joint training, and stage-decision. Our design supports user-defined stage-decision ratio and threshold to adaptively balance model performance and energy consumption. Our experimental results show that MicroT exhibits an enhancement in model performance by up to 9.87%, and energy consumption by up to about 29.13%. Employing standard

settings for the stage-decision ratio and threshold, MicroT can achieve a 5.91% improvement in model performance and an energy saving of about 14.47%. MicroT's adaptability, energy efficiency, and sufficient model performance make it a feasible solution for the multi-task challenge of MCUs.

## References

[1] Lachit Dutta and Swapna Bharali. Tinyml meets iot: A comprehensive survey. *Internet of Things*, 16:100461, 2021.

[2] Youssef Abadade, Anas Temouden, Hatim Bamoumen, Nabil Benamar, Yousra Chtouki, and Abdelhakim Senhaji Hafid. A comprehensive survey on tinyml. *IEEE Access*, 2023.

[3] Swapnil Sayan Saha, Sandeep Singh Sandha, and Mani Srivastava. Machine learning for microcontroller-class hardware-a review. *IEEE Sensors Journal*, 2022.

[4] Deblina Bhattacharjee, Tong Zhang, Sabine Süsstrunk, and Mathieu Salzmann. Mult: An end-to-end multitask learning transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12031–12041, 2022.

[5] Yang Li, Amirmohammad Kazemeini, Yash Mehta, and Erik Cambria. Multitask learning for emotion and personality traits detection. *Neurocomputing*, 493:340–350, 2022.

[6] Tuochao Chen, Justin Chan, and Shyamnath Gollakota. Underwater messaging using mobile devices. In *Proceedings of the ACM SIGCOMM 2022 Conference*, pages 545–559, 2022.

[7] Yushan Huang and Hamed Haddadi. Poster: Towards battery-free machine learning inference and model personalization on mcus. In *Proceedings of the 21st Annual International Conference on Mobile Systems, Applications and Services*, pages 571–572, 2023.

[8] Jed Mills, Jia Hu, and Geyong Min. Multi-task federated learning for personalised deep neural networks in edge computing. *IEEE Transactions on Parallel and Distributed Systems*, 33(3):630–641, 2021.

[9] Qiong Chen, Zimu Zheng, Chuang Hu, Dan Wang, and Fangming Liu. On-edge multi-task transfer learning: Model and practice with data-driven task allocation. *IEEE Transactions on Parallel and Distributed Systems*, 31(6):1357–1371, 2019.

[10] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017.

[11] Praneeth Vepakomma, Otkrist Gupta, Tristan Swedish, and Ramesh Raskar. Split learning for health: Distributed deep learning without sharing raw patient data. *arXiv preprint arXiv:1812.00564*, 2018.

[12] Yu Zhang and Qiang Yang. An overview of multi-task learning. *National Science Review*, 5(1):30–43, 2018.

[13] Jiangshan Hao, Yanchao Zhao, and Jiale Zhang. Time efficient federated learning with semi-asynchronous communication. In *2020 IEEE 26th International Conference on Parallel and Distributed Systems (ICPADS)*, pages 156–163. IEEE, 2020.

[14] Zheng Chai, Yujing Chen, Ali Anwar, Liang Zhao, Yue Cheng, and Huzefa Rangwala. Fedat: A high-performance and communication-efficient federated learning system with asynchronous tiers. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–16, 2021.

[15] Sayed Saad Afzal, Waleed Akbar, Osvy Rodriguez, Mario Doumet, Unsoo Ha, Reza Ghaffarivardavagh, and Fadel Adib. Battery-free wireless imaging of underwater environments. *Nature communications*, 13(1):5546, 2022.

[16] Ahmed Allam, Waleed Akbar, and Fadel Adib. An analytical framework for low-power underwater backscatter communications. *The*

*Journal of the Acoustical Society of America*, 153(3_supplement):A376–A376, 2023.

[17] Jenny Hamer, Mehryar Mohri, and Ananda Theertha Suresh. Fedboost: A communication-efficient algorithm for federated learning. In *International Conference on Machine Learning*, pages 3973–3983. PMLR, 2020.

[18] Jed Mills, Jia Hu, and Geyong Min. Communication-efficient federated learning for wireless edge intelligence in iot. *IEEE Internet of Things Journal*, 7(7):5986–5994, 2019.

[19] Ronghang Hu and Amanpreet Singh. Unit: Multimodal multitask learning with a unified transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1439–1449, 2021.

[20] Davis Blalock, Jose Javier Gonzalez Ortiz, Jonathan Frankle, and John Guttag. What is the state of neural network pruning? *Proceedings of machine learning and systems*, 2:129–146, 2020.

[21] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference. In *Low-Power Computer Vision*, pages 291–326. Chapman and Hall/CRC, 2022.

[22] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129:1789–1819, 2021.

[23] Michael Zhu and Suyog Gupta. To prune, or not to prune: exploring the efficacy of pruning for model compression. *arXiv preprint arXiv:1710.01878*, 2017.

[24] Rick Pandey, Sebastian Uziel, Tino Hutschenreuther, and Silvia Krug. Towards deploying dnn models on edge for predictive maintenance applications. *Electronics*, 12(3):639, 2023.

[25] Zhepeng Wang, Yawen Wu, Zhenge Jia, Yiyu Shi, and Jingtong Hu. Lightweight run-time working memory compression for deployment of deep neural networks on resource-constrained mcus. In *Proceedings of the 26th Asia and South Pacific Design Automation Conference*, pages 607–614, 2021.

[26] Maxim Zemlyanikin, Alexander Smorkalov, Tatiana Khanova, Anna Petrovicheva, and Grigory Serebryakov. 512kib ram is enough! live camera face recognition dnn on mcu. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019.

[27] Xiao Liu, Fanjin Zhang, Zhenyu Hou, Li Mian, Zhaoyu Wang, Jing Zhang, and Jie Tang. Self-supervised learning: Generative or contrastive. *IEEE transactions on knowledge and data engineering*, 35(1):857–876, 2021.

[28] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.

[29] Van Nguyen Nguyen, Thibault Groueix, Georgy Ponimatkin, Vincent Lepetit, and Tomas Hodan. Cnos: A strong baseline for cad-based novel object segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2134–2140, 2023.

[30] Haonan Sun, Kai Zhang, Wei Lan, Qiufeng Gu, Guangxiang Jiang, Xue Yang, Wanli Qin, and Dongran Han. An ai dietitian for type 2 diabetes mellitus management based on large language and image recognition models: Preclinical concept validation study. *Journal of Medical Internet Research*, 25:e51300, 2023.

[31] Jiayi Liu, Samarth Tripathi, Unmesh Kurup, and Mohak Shah. Pruning algorithms to accelerate convolutional neural networks for edge applications: A survey. *arXiv preprint arXiv:2005.04275*, 2020.

[32] Ying Jin, Jiaqi Wang, and Dahua Lin. Multi-level logit distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24276–24285, 2023.

[33] Yunlong Yu, Bin Li, Zhong Ji, Jungong Han, and Zhongfei Zhang. Knowledge distillation classifier generation network for zero-shot learning. *IEEE Transactions on Neural Networks and Learning Systems*, 2021.

[34] Haoran Zhao, Xin Sun, Junyu Dong, Zihe Dong, and Qiong Li. Knowledge distillation via instance-level sequence learning. *Knowledge-Based Systems*, 233:107519, 2021.

[35] Guodong Xu, Ziwei Liu, Xiaoxiao Li, and Chen Change Loy. Knowledge distillation meets self-supervision. In *European Conference on Computer Vision*, pages 588–604. Springer, 2020.

[36] Thanh V Nguyen, Raymond KW Wong, and Chinmay Hegde. Benefits of jointly training autoencoders: An improved neural tangent kernel analysis. *IEEE Transactions on Information Theory*, 67(7):4669–4692, 2021.

[37] Iasonas Kokkinos. Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6129–6138, 2017.

[38] Zhiyuan Tang, Lantian Li, Dong Wang, and Ravichander Vipperla. Collaborative joint training with multitask recurrent model for speech and speaker recognition. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 25(3):493–504, 2016.

[39] Trevor Standley, Amir Zamir, Dawn Chen, Leonidas Guibas, Jitendra Malik, and Silvio Savarese. Which tasks should be learned together in multi-task learning? In *International Conference on Machine Learning*, pages 9120–9132. PMLR, 2020.

[40] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1):43–76, 2020.

[41] Kavya Kopparapu, Eric Lin, John G Breslin, and Bharath Sudharsan. Tinyfedtl: Federated transfer learning on ubiquitous tiny iot devices. In *2022 IEEE International Conference on Pervasive Computing and Communications Workshops and other Affiliated Events (PerCom Workshops)*, pages 79–81. IEEE, 2022.

[42] Qisong Hu, Xiaochen Tang, and Wei Tang. A real-time patient-specific sleeping posture recognition system using pressure sensitive conductive sheet and transfer learning. *IEEE Sensors Journal*, 21(5):6869–6879, 2020.

[43] Hao Wu, Jinghao Feng, Xuejin Tian, Edward Sun, Yunxin Liu, Bo Dong, Fengyuan Xu, and Sheng Zhong. Emo: Real-time emotion recognition from single-eye images for resource-constrained eyewear devices. In *Proceedings of the 18th International Conference on Mobile Systems, Applications, and Services*, pages 448–461, 2020.

[44] Abdallah Wagih Ibrahim AashiDutt. Sea animals image dataset. *https://www.kaggle.com/datasets/vencerlanz09/sea-animals-image-dataste*, 2022.

[45] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. The oxford-iiit pet dataset.

[46] Yuchen Zhao, Sayed Saad Afzal, Waleed Akbar, Osvy Rodriguez, Fan Mo, David Boyle, Fadel Adib, and Hamed Haddadi. Towards battery-free machine learning and inference in underwater environments. In *Proceedings of the 23rd Annual International Workshop on Mobile Computing Systems and Applications*, pages 29–34, 2022.

[47] Seunghyeok Jeon, Yonghun Choi, Yeonwoo Cho, and Hojung Cha. Harvnet: Resource-optimized operation of multi-exit deep neural networks on energy harvesting devices. In *Proceedings of the 21st Annual International Conference on Mobile Systems, Applications and Services*, pages 42–55, 2023.

[48] Shekoofeh Azizi, Basil Mustafa, Fiona Ryan, Zachary Beaver, Jan Freyberg, Jonathan Deaton, Aaron Loh, Alan Karthikesalingam, Simon Kornblith, Ting Chen, et al. Big self-supervised models advance medical image classification. In *Proceedings of the IEEE/CVF international*

*conference on computer vision*, pages 3478–3488, 2021.

[49] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR, 2021.

[50] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.

[51] Jacob Peplinski, Joel Shor, Sachin Joglekar, Jake Garrison, and Shwetak Patel. Frill: A non-semantic speech embedding for mobile devices. *arXiv preprint arXiv:2011.04609*, 2020.

[52] Yihui He, Ji Lin, Zhijian Liu, Hanrui Wang, Li-Jia Li, and Song Han. Amc: Automl for model compression and acceleration on mobile devices. In *Proceedings of the European conference on computer vision (ECCV)*, pages 784–800, 2018.

[53] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V Le. Mnasnet: Platform-aware neural architecture search for mobile. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2820–2828, 2019.

[54] STMicroelectronics. X-cube-ai. 2023.

[55] Monsoon Solutions Inc. Monsoon high voltage power monitor. *https://www.msoon.com/*.

[56] Han Cai, Ligeng Zhu, and Song Han. ProxylessNAS: Direct neural architecture search on target task and hardware. In *International Conference on Learning Representations*, 2019.

[57] Ji Lin, Wei-Ming Chen, Yujun Lin, Chuang Gan, and Song Han. Mcunet: Tiny deep learning on iot devices. *Advances in Neural Information Processing Systems*, 33, 2020.

[58] Sulaiman Sadiq, Jonathon Hare, Simon Craske, Partha Maji, and Geoff Merrett. Enabling imagenet-scale deep learning on mcus for accurate and efficient inference. *IEEE Internet of Things Journal*, 2023.

[59] Zhenhong Sun, Ce Ge, Junyan Wang, Ming Lin, Hesen Chen, Hao Li, and Xiuyu Sun. Entropy-driven mixed-precision quantization for deep network design. *Advances in Neural Information Processing Systems*, 35:21508–21520, 2022.

[60] Young D Kwon, Rui Li, Stylianos I Venieris, Jagmohan Chauhan, Nicholas D Lane, and Cecilia Mascolo. Tinytrain: Deep neural network training at the extreme edge. *arXiv preprint arXiv:2307.09988*, 2023.

[61] Ji Lin, Ligeng Zhu, Wei-Ming Chen, Wei-Chen Wang, Chuang Gan, and Song Han. On-device training under 256kb memory. *Advances in Neural Information Processing Systems*, 35:22941–22954, 2022.

[62] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011 dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.

[63] Herve Goeau, Pierre Bonnet, and Alexis Joly. Plant identification based on noisy web data: the amazing performance of deep learning (lifeclef 2017). CEUR Workshop Proceedings, 2017.

[64] STMicroelectronics. Stm32cubeide. 2023.

[65] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and< 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.

[66] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.

[67] Francesco Daghero, Alessio Burrello, Chen Xie, Marco Castellano, Luca Gandolfi, Andrea Calimera, Enrico Macii, Massimo Poncino, and Daniele Jahier Pagliari. Human activity recognition on microcontrollers with quantized and adaptive deep neural networks. *ACM*

*Transactions on Embedded Computing Systems (TECS)*, 21(4):1–28, 2022.

[68] Nattapol Kaewmard and Saiyan Saiyod. Sensor data collection and irrigation control on vegetable crop using smart phone and wireless sensor networks for smart farm. In *2014 IEEE Conference on Wireless Sensors (ICWiSE)*, pages 106–112. IEEE, 2014.

[69] Jinping Wang, Jun Li, Yanli Shi, Jianhuang Lai, and Xiaojun Tan. Am$^3$net: Adaptive mutual-learning-based multimodal data fusion network. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(8):5411–5426, 2022.

[70] Xiaomin Ouyang, Xian Shuai, Jiayu Zhou, Ivy Wang Shi, Zhiyuan Xie, Guoliang Xing, and Jianwei Huang. Cosmo: contrastive fusion learning with small data for multimodal human activity recognition. In *Proceedings of the 28th Annual International Conference on Mobile Computing And Networking*, pages 324–337, 2022.