

PockEngine: Sparse and Efficient Fine-tuning in a Pocket

Ligeng Zhu

ligeng@mit.edu

MIT

Cambridge, MA, USA

Wei-Chen Wang

wweichen@mit.edu

MIT

Cambridge, MA, USA

Lanxiang Hu

hlxde2@gmail.com

UCSD

San Diego, CA, USA

Wei-Ming Chen

wmchen@mit.edu

MIT

Cambridge, MA, USA

Ji Lin

jilin@mit.edu

MIT

Cambridge, MA, USA

Chuang Gan

ganchuang@csail.mit.edu

MIT-IBM Watson AI Lab

Cambridge, MA, USA

Song Han

songhan@mit.edu

MIT, NVIDIA

Cambridge, MA, USA

ABSTRACT

On-device learning and efficient fine-tuning enable continuous and privacy-preserving customization (e.g., locally fine-tuning large language models on personalized data). However, existing training frameworks are designed for cloud servers with powerful accelerators (e.g., GPUs, TPUs) and lack the optimizations for learning on the edge, which faces challenges of resource limitations and edge hardware diversity. We introduce **PockEngine**: a tiny, sparse and efficient engine to enable fine-tuning on various edge devices. PockEngine supports **sparse backpropagation**: it prunes the backward graph and sparsely updates the model with measured memory saving and latency reduction while maintaining the model quality. Secondly, PockEngine is **compilation first**: the entire training graph (including forward, backward and optimization steps) is derived at compile-time, which reduces the runtime overhead and brings opportunities for graph transformations. PockEngine also integrates a rich set of **training graph optimizations**, thus can further accelerate the training cost, including operator reordering and backend switching. PockEngine supports **diverse applications, frontends and hardware backends**: it flexibly compiles and tunes models defined in PyTorch/TensorFlow/Jax and deploys binaries to mobile CPU/GPU/DSPs. We evaluated PockEngine on both vision models and large language models. PockEngine achieves up to $15 \times$ speedup over off-the-shelf TensorFlow (Raspberry Pi), $5.6 \times$ memory saving back-propagation (Jetson AGX Orin). Remarkably, PockEngine enables fine-tuning LLaMav2-7B on NVIDIA Jetson AGX Orin at 550 tokens/s, $7.9 \times$ faster than the PyTorch.

CCS CONCEPTS

- Computer systems organization → Neural networks.

Permission to make digital or hard copies of part or all 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).

MICRO '23, October 28–November 1, 2023, Toronto, ON, Canada

© 2023 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-0329-4/23/10.

<https://doi.org/10.1145/3613424.3614307>

KEYWORDS

neural network, sparse update, on-device training, efficient finetuning

ACM Reference Format:

Ligeng Zhu, Lanxiang Hu, Ji Lin, Wei-Chen Wang, Wei-Ming Chen, Chuang Gan, and Song Han. 2023. PockEngine: Sparse and Efficient Fine-tuning in a Pocket. In *56th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO '23), October 28–November 1, 2023, Toronto, ON, Canada*. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3613424.3614307>

1 INTRODUCTION

Edge devices are ubiquitous and produce an increasing amount of data in our daily lives. The need for intelligent, personalized, and private AI is rapidly growing, as a single model fails to fit different users' needs. However, while deep learning inferences are widely performed on edge devices, the training of deep neural networks is typically run on cloud GPU servers. Cloud-based training requires users to upload their personal data to the cloud, which not only incurs additional data transfer costs, but also brings privacy risks over sensitive data (e.g., healthcare data, keyboard input history, GPS location, etc.).

On-device training is a promising solution for model customization without sacrificing privacy (Figure 1). It allows a pre-trained model to continuously adapt to sensor data without sending it to the cloud. For example, the smart keyboard model can update itself to better predict the next word from users' typing history; the email assistant can learn from users' previous drafts and train personalized language models; vision models can automatically adapt to environments with domain shifts [54]). The near-sensor training paradigm also brings important benefits for energy and connectivity: it saves energy from data transmission (which is much more expensive than computation [36]); it also helps with applications like ocean sensing [26] and smart agriculture [57] that do not have physical access to the Internet.

Despite all the benefits, on-device training is difficult due to the following challenges:

(1) **Resource Limitations.** The capacity of edge devices is orders of magnitude smaller than cloud servers. People have been

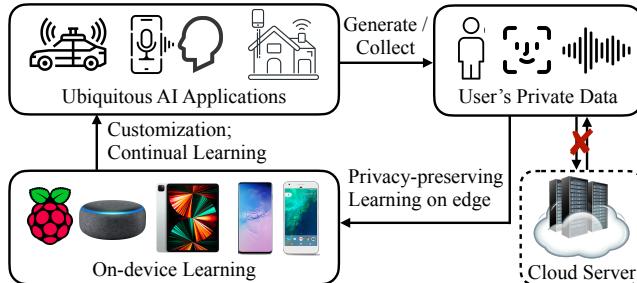


Figure 1. On-device learning and local fine-tuning enable customization, protect privacy, and form a virtuous cycle between user and devices.

trying hard to squeeze deep learning models just for edge *inference*, while model *training* and *fine-tuning* are more power-, computation-, and memory-expensive. We need extra memory to store all intermediate feature maps for backpropagation, and extra computation for the backward pass (roughly 3 \times compared to inference). Sometimes the training needs a larger batch size to ensure a stable convergence, making the process even more costly. For MobilenetV2 [51] the training memory is 14 \times and 7.3 \times larger than inference (batch size 8) and for BERT [19] the peak memory usage is 7.3 \times larger compared to inference. Furthermore, the optimizers also require extra memory (2x for Momentum and 3x for Adam [31]). With the current training framework, the training costs could soon exceed the resource limits of edge hardware.

(2) Hardware Diversity While the accelerators on cloud servers are dominated by GPUs, the hardware of edge platforms has a wide range of options on the market. The processor ranges from ARM microcontrollers to powerful Apple M1 chips, and the accelerator varies between Qualcomm Adreno GPUs, Hexagon DSPs, and edge TPUs. Each hardware comes with a different inference library. PockEngine can directly use these inference libraries for training by compiling the training graph into standard ONNX format. On the other hand, popular deep learning training frameworks like TensorFlow [5], PyTorch [47] and Jax [10] are developed for high-end cloud GPUs/TPUs. The performance is poor when directly applied to edge platforms*.

To address the above challenges, we introduce **PockEngine**, a tiny and efficient training engine designed for on-device training. We highlight the following properties:

- PockEngine provides system-level support for both **dense** and **sparse backpropagation**. Apart from updating the whole model, PockEngine supports flexible *sparse* update schemes by computing the gradients for only part of the weights, which proves to be a more efficient option for fine-tuning/transfer learning without harming the accuracy [11, 21, 24, 25, 38, 42, 43]. Existing training frameworks can only *simulate* the sparse backpropagation by computing the backward and mask out gradients, but cannot *realize* measured speed up and memory savings. PockEngine supports sparse

backpropagation via graph pruning and dead code elimination with the compilation nature, leading to smaller computation and memory usage.

- PockEngine is a **compilation-based** efficient training engine and enables many **inference-only framework to perform training**. Our compilation workflow helps to connect diverse model architectures and frontend options (e.g., vision/NLP models, PyTorch/TensorFlow/ONNX definitions) with various backend libraries (e.g., SNPE for Qualcomm, Metal for Apple Silicon, TVM), exposing a unified intermediate representation (IR). By sharing the same set of operators for both forward and backward operations, we not only enable inference frameworks to train neural networks, but also allow for various graph optimizations to improve efficiency (see Figure 4).
- PockEngine implements a rich set of **graph optimizations** to improve the efficiency on edge devices, including operator fusion, operator reordering, layout transforms, and backend switching that are conventionally used for *inference only*. We find that the training graphs actually have more optimization opportunities due to their complexity. By sharing the same operator set with inference graphs, PockEngine can well utilize the optimization techniques from inference engines (e.g., PockEngine utilizes previously inference-only winograd convolution to accelerate training).

We extensively evaluated PockEngine on six edge platforms and six deep learning tasks from vision to NLP. PockEngine achieves up to 11 \times speedup over TensorFlow for the same training workload. With sparse backpropagation, we can further improve the acceleration up to 21 \times without losing transfer learning accuracy on tiny microcontrollers. We hope our work can contribute to the thriving of on-device training by providing a *general-purpose, high-efficiency, user-friendly* training framework for edge devices.

2 RELATED WORK

2.1 Cloud Deep Learning Systems

The success of deep learning is built on top of popular training frameworks such as PyTorch [47], TensorFlow [6], MXNet [13], JAX [10], etc. These systems are designed for development flexibility and depend on a host language (e.g., Python) to execute. This brings significant memory overhead (>300MB) and makes the runtime especially slow on low-frequency CPU (e.g., ARM Cortex). Moreover, the operator kernels are optimized for high-end GPU devices and lack performance tuning for edge devices and some overheads such as extra gradient buffers for the optimizer step are not considered a bottleneck for powerful server hardware. PockEngine is a compilation-based framework thus the runtime does not rely on host languages as compared in Table 1. This moves most workloads from runtime to compile-time to minimize the runtime overhead and enables later optimizations to improve training throughput.

*The frameworks themselves cannot even be installed due to the tight resource constraints of low-end hardware like microcontrollers [42].

Table 1. Comparison between existing deep learning frameworks. “-” denotes the feature is not fully supported for training.

	Support Training	Support Sparse-BP	Run without Host Language	Kernel Optimized for Edge	Compile-Time AutoDiff	Graph Optimizations
PyTorch [47]	✓	✗	✗	✗	✗	✗
TensorFlow [5]	✓	✗	✗	✗	✗	-
Jax [10]	✓	✗	✗	✗	✗	✗
TVM [14]	✗	✗	✓	✓	-	✓
MNN [30]	✓	✗	✓	✓	✗	✗
PockEngine (ours)	✓	✓	✓	✓	✓	✓

2.2 Edge Deep Learning Systems

When deploying models on tiny edge devices, inference libraries like TVM [14], TF-Lite [3], NCNN [1], TensorRT [2], and OpenVINO [58] deliver optimized kernels for mobile platforms and provide a lightweight runtime without host language. However, they focus mostly on inference and do not support on-device training. MNN [30] has preliminary support for CNNs but the flexibility is rather limited and it does not optimize training memory usage. POET [48] applies rematerialization and paging to deal with restricted memory size, but it introduces extra computation, relies on large external Flash (e.g. 32GB SD Card) and does not support general model and workload definition. PockEngine provides complete training support for popular models at various scales including MCUNet [41], MobileNetV2 [51], ResNet [23], DistilBERT [52], and BERT [19]. PockEngine optimizes both computation and memory efficiency to make on-device training easy and realistic.

2.3 Efficient On-Device Learning Algorithms

Edge devices have limited computational capacity. Therefore, on-device training for edge devices often focuses on transfer learning [11, 34]. It first pre-trains the model on large-scale datasets to learn general and rich features, such as ImageNet [18] for ConvNets or BooksCorpus [65] for BERT. The model is then transferred to downstream tasks, such as Visual Wake Words [17] for vision or the GLUE benchmark [59] for language. After which, the model can be customized to a small amount of personal data (e.g., learning a user’s accent) to perform better at the *same* task.

Due to the smaller scale and diversity of the downstream data, people found that it is not always necessary to update the entire model to achieve a good performance. *Sparingly* updating part of the model proves to be a good solution that achieves similar or better performance at a smaller training cost [11, 21, 24, 25, 38, 42, 43]. The most straightforward method is to fine-tune only the classifier layer [12, 20, 22, 53], but the capacity is limited when the domain shift is large. For CNN models, people have investigated fine-tuning only biases [11, 62], batch normalization layers [21, 44], added parallel branches [11], etc. The sparse backpropagation scheme is even more popular for adapting pre-trained language models (e.g., BERT [19], GPT [50]) to various downstream tasks, which significantly reduce the trainable parameters [24, 25, 38]. However, sparse backpropagation lacks system support. Despite the great theoretical savings, existing training frameworks cannot realize

measured speedup or memory saving from sparse backpropagation. PockEngine provides system-level support for such flexible workloads to deliver a faster program and efficient runtime.

2.4 Computation Graph Transformation and Optimizations

There are plenty of graph transformations for inference scenarios. For example, one common transform used in edge deployment is data layout conversion, as the ‘NCHW’ preferred by GPU training is not efficient on the edge. Another common optimization technique is layer fusion. IO-intensive layers (e.g. ReLU) can usually be fused into preceding compute-intensive layers (e.g. CONV, LINEAR). In addition, MetaFlow [28] proposes functional-preserving graph transformations to optimize DNN architectures. TASO [27] further introduces automated generation of transformation rules using formal verification. These techniques have been proven effective in inference, but few studies have explored their performance on training, even though the training graph is much more complex. Standing on the shoulder of conventional wisdom, PockEngine is early exploration for apply these graph optimizations techniques to on-device training and discover more potential optimizations. PockEngine shows that these optimizations bring up to 1.2x speedup.

2.5 Compilation-Based Workflow

Existing training frameworks (e.g., PyTorch, TensorFlow) are based on runtime auto differentiation for flexibility. However, the design is not suitable for edge devices with limited memory and computation resources. Instead, PockEngine is based on a compilation-based workflow, sharing the following benefits:

Offload Workload from Runtime to Compile Time. With the compilation-centric design, we can offload part of the workload from runtime to compile time, like backward graph derivation with autodiff, memory scheduling, execution planning, etc. Modern neural network usually consists of thousands of operators, the overhead might be small for cloud servers but not negligible for edge devices (Figure. 7).

By offloading computation to the compiler, it is possible to perform more aggressive optimizations that would not be feasible or efficient to perform at runtime. For example, PockEngine performs graph pruning, fusions, and backend switching, which can lead to significant performance gains and memory saving.

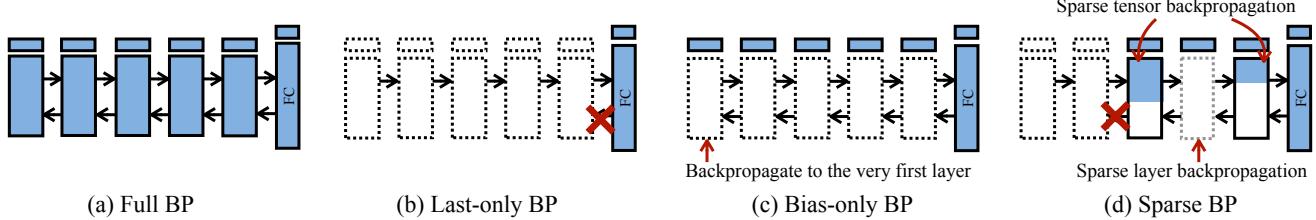


Figure 2. The computation graph of different backpropagation schemes on a five-layer model. We use blue to indicate the demanded intermediate activations during training. Sparse-BP delivers the best cost-quality tradeoff which we will show in Section 4.

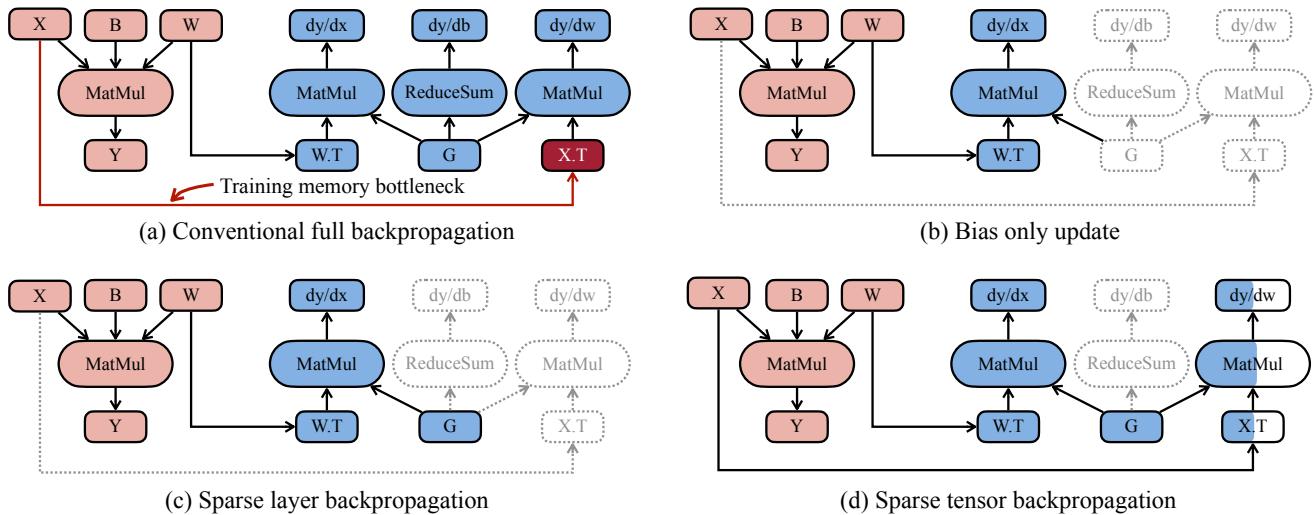


Figure 3. The computation graph of sparse backpropagation for a linear layer. Red and blue blocks indicate the forward and backward OPs respectively. The red line denotes the training memory bottleneck brought by storing activations, which can be avoided using bias only / sparse update as shown in (b) (c) (d).

Another advantage of compilation-based workflow is that it allows us to optimize the code across the entire program, rather than just focusing on optimizing individual operations at runtime. This not only allows us to compile used operators only to ship slim binaries, but also reveals the memory redundancy in the training loop (details in Section 3.2).

Support Diverse Frontends/Backends. Unlike the cloud, edge platforms are highly diverse, with different instruction sets, degrees of parallelism, etc. Our compilation-based workflow provides general support for various frontends/backends. It can effortlessly support *training* on hardware and vendor libraries that are designed specifically for *inference* (e.g., PockEngine can enable training on Qualcomm Hexagon DSPs with SNPE library).

The *PockEngine frontend* takes in a neural network represented in various representations (e.g., ONNX, torchscript, tf.graph) and analyzes the DAG structure. It will then perform automatic differentiation (autodiff) to derive the backward graph which computes the gradients w.r.t. the loss function (Figure 7). With the *static* forward and backward graph, PockEngine will convert it into a unified intermediate representation (IR), perform graph optimizations (will be introduced later), and generate the code for different backends. Only used operators will be compiled and PockEngine link these OPs

to build a light-weight executable binary. The *PockEngine* backend supports both vendor libraries (*e.g.*, SNPE for Snapdragon GPUs and DSPs, TensorRT for NVIDIA GPUs) and customized kernels (*e.g.*, TVM [14] tuning for ARM CPUs).

Notably, instead of binding each operator with a backward implementation (*e.g.*, `matmul`, `matmul_backward`), PockEngine uses the same set of primitive operations as inference to construct the training graph, allowing us to utilize inference-only backends (*e.g.*, SNPE, TensorRT, TVM) for training, achieving high efficiency at minimal engineer effort.

2.6 Sparse Backpropagation and Computation Graph Pruning

Edge devices have a limited computation capacity compared to the cloud. Therefore, on-device training on edge usually targets a transfer learning/fine-tuning scenario. Due to the smaller scale and diversity of the downstream data, people found that updating the entire model may not always lead to the best performance due to over-fitting and feature distortion [11, 34]. Updating only a subset of the models is proven to be a good solution that achieves similar or better performance at a much smaller training cost, including

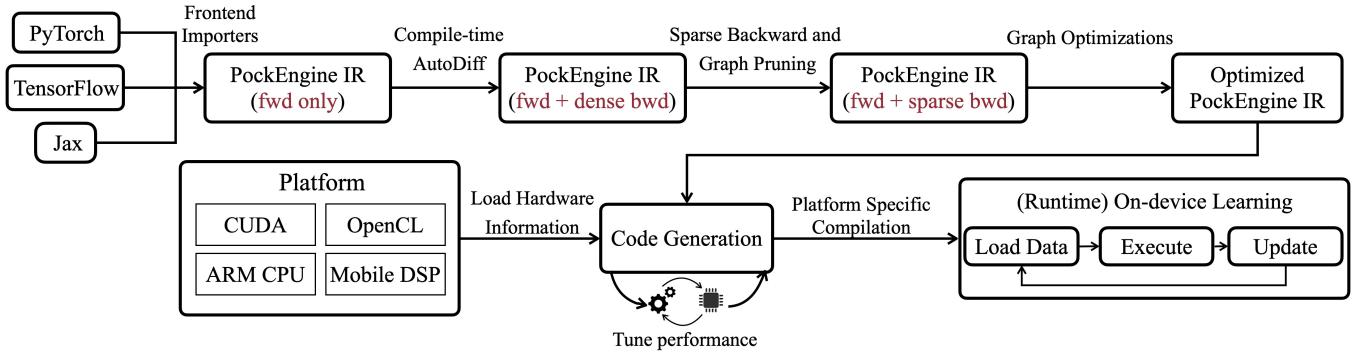


Figure 4. The workflow of PockEngine. PockEngine performs the auto-diff at compile-time, prunes the computation graph to support sparse backpropagation, and enables previously inference-only hardware platforms to perform backpropagation. PockEngine enables efficient fine-tuning on resource-constrained devices like NVIDIA Jetson and mobile devices.

updating bias terms [11] and the normalization layers [21] for vision models training the low-rank parts [25] and input prompts for language models [38], and sparsely update the important modules [42]. PockEngine aims to generally support on-device training for various workloads and we focus on the sparse update to reduce training costs.

During the compilation, PockEngine takes in a user-defined sparse backpropagation scheme and will *prune* the corresponding subgraphs of backpropagation calculation. PockEngine flexibly supports the following sparse backpropagation patterns:

Bias-only Update. Bias-only update does not require saving the intermediate activation [11], which significantly reduces memory usage (consider a linear layer $y = Wx$, $dW = f_1(dy, x)$, $db = f_2(dy)$, only the weight gradient requires saving the input). It also saves the computation by 1/3 by skipping dW computation.

Layer-wise Sparse Backpropagation. Not all the layers/weight tensors are equally important for transfer learning [42]. For transfer learning to a downstream task, we find that part of the layers can be kept frozen without affecting the transfer learning performance (we can find the layers to freeze by sensitivity analysis [42]; detailed in Section 4.1). Therefore, we can skip the computation of part of the layers to further improve the training throughput.

Sub-layer Sparse Backpropagation. For edge devices with limited capacity (e.g., microcontrollers), we further support sub-layer level sparse BP, where only part of the channels of a layer (convolutional layers and linear layers) are updated[†]. It further reduces the memory cost for storing intermediate activation (we do not need to store activation for the frozen channels) and the computation cost for gradient calculation.

3 POCKENGINE

Compared to conventional training frameworks, sparse backpropagation has the following unique advantages

- Expensive intermediate activations can be released immediately after forward When either learning the bias-only

[†]Following [42], we simply update the first k channels of a layer. k is the #channels to update.

(dy/db and dy/dx) or fully skipping the layer (only dy/dx to keep chain-rule). Thus sparse backpropagations greatly reduce the main memory bottleneck of training (the red connection line in Figure 3.a).

- Sparse back-propagation does not back-propagate the very first layers in DNN models since there is no need to compute gradients to the front layers if they do not require gradients (the red X mark in Figure 5).

None of the prior work can convert the theoretical savings into measured speed-up and memory savings. PockEngine provides systematic support for sparse BP and is able to actually reduce the on-device training cost and we expand as follows

3.1 Searching for Sparse Backpropagation Scheme

Not all the weights are equally important for transfer learning [21, 35, 42]. We aim to fine-tune only the **important weights** to reduce the training costs while preserving the model’s accuracy.

Cost Model and Search Criterion. In order to find the training scheme, we build cost models for model quality and training cost. Following [42], we first fine-tune only one linear (conv, fc) layer until convergence, and then repeat this process for all layers. This is an offline analysis and we use the accuracy improvement/degradation as the “contribution” of the weights of i^{th} layer (Δacc_{Wi}). Similarly, we obtain the results for bias terms of k^{th} layer (Δacc_{bk}) and then iteratively repeat the same operations to all weights and biases to estimate their performance.

For the training cost, we focus on the memory as edge devices usually have limited memory and will easily get OOM. Thus we profile the feature map size and record it as $Memory_{k,i,r}$. We then solve the following optimization:

$$\begin{aligned} k^*, i^*, r^* = \max & \left(\sum_{k,i,r} \Delta acc_{bk} + \sum_{i,i,r} \Delta acc_{W_{i,r}} \right) \\ \text{s.t. } & Memory(k, i, r) \leq \text{constraint}, \end{aligned} \quad (1)$$

where i is the layer index of weights, k is the layer index of biases and r is the ratio of learnable weights. Optimizing the objectives

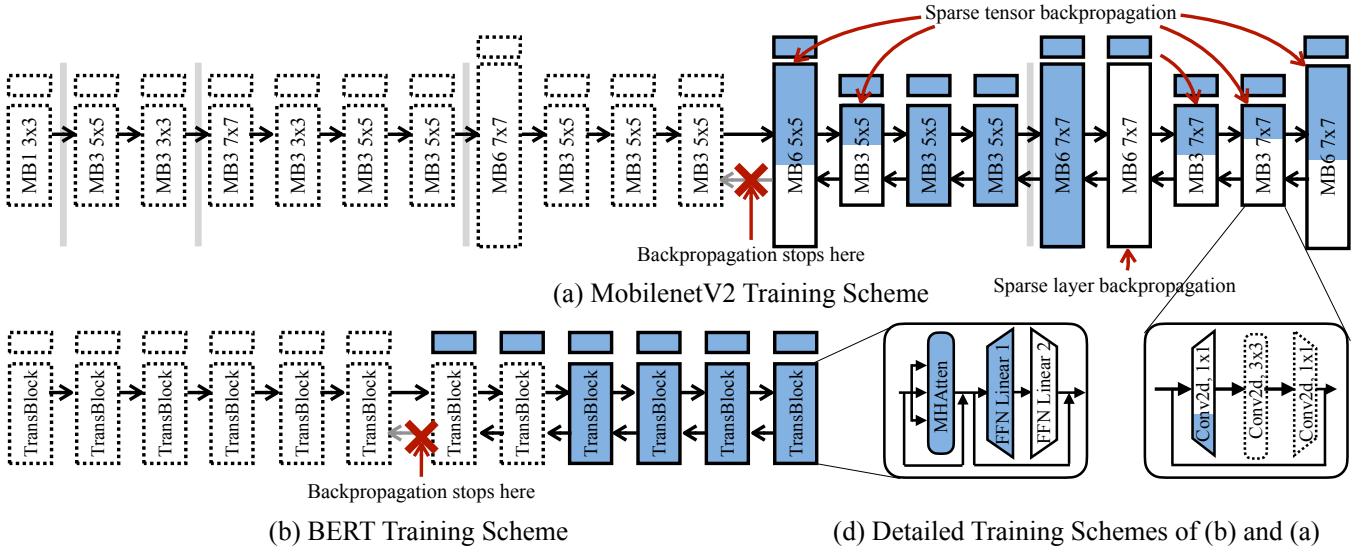


Figure 5. The computation graph of sparse backpropagation for ConvNet and Transformers.

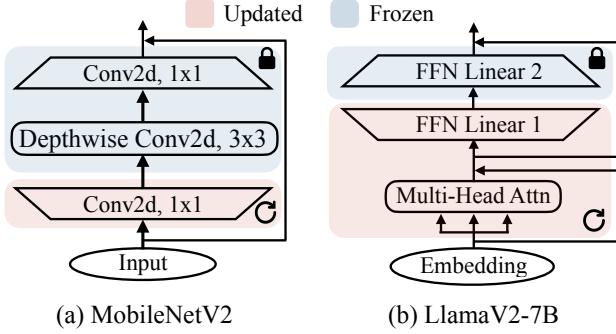


Figure 6. The sparse backpropagation schemes for MobileNetV2 and LlamaV2-7B building blocks. The first point-wise convolution plays an important role for ConvNet, while for Llama models, the attention module and first FNN layer are more important.

finds the optimal update config where total contributions are maximized and the memory footprint does not exceed the constraint. We assume that the accuracy contribution of each tensor (Δ_{acc}) can be summed up thus the problem can be efficiently solved with evolutionary search.

Generalization and Acceleration. It is worth noting that the sparse update scheme is general and universal across different datasets. We only perform ONE scheme search on CIFAR (for vision models) and CoLA (for language models) and sparse-BP demonstrates good generalization capability. The schemes achieve competitive training accuracy compared to full fine-tuning (Table 2 and Table 3). Specifically, we find that for CNNs: it is most effective to update the weights of the first convolution in each block, while for transformer blocks, the weights in the attention module and the first linear layer in the Feed-Forward Network (FFN) are more important (Figure 6). Such schemes are also memory-efficient: the

depthwise conv and second pointwise conv in the inverted bottleneck block (Figure 6.a) and the second linear layer in the FFN (Figure 6.b) have the largest input activation, while our update scheme does not require saving these large features.

After finding and specifying the gradients needed for the on-device training, PockEngine automatically traces dependency and analyzes the updated topology, then prunes the training graph using dead code elimination (DCE) to prune the computation graph and remove intermediate nodes and buffers that are no longer needed for the training. Because the pruning is performed on graph level at compile-time, it can deliver measured memory saving and throughput improvement.

3.2 Training Graph Optimization

After we get the static, pruned training graph, PockEngine applies various graph optimization techniques on the unified IR before translating to different backends, which further improves the training efficiency.

Operator Reordering and In-place Update. Different execution orders lead to different life cycles of tensors and the overall/peak memory footprint will be also affected even for the same computational graphs. This has been well-studied for inference [7, 39] but less discussed for training because the backward graph is usually derived during runtime and the compiler/scheduler does not have global information of the training process.

A concrete example is the optimizer, where the gradients are applied to update the model parameters. In conventional training, frameworks calculate all gradients and then apply the update. This is common among frameworks like PyTorch and TensorFlow as the optimizer and forward-backward are separate components in the system design. However, such a practice leads to significant memory waste for storing the gradients. In small batch training with sparse backpropagation, the cost of storing parameter gradients is close to peak memory usage in forward and backward as shown in Table 4:

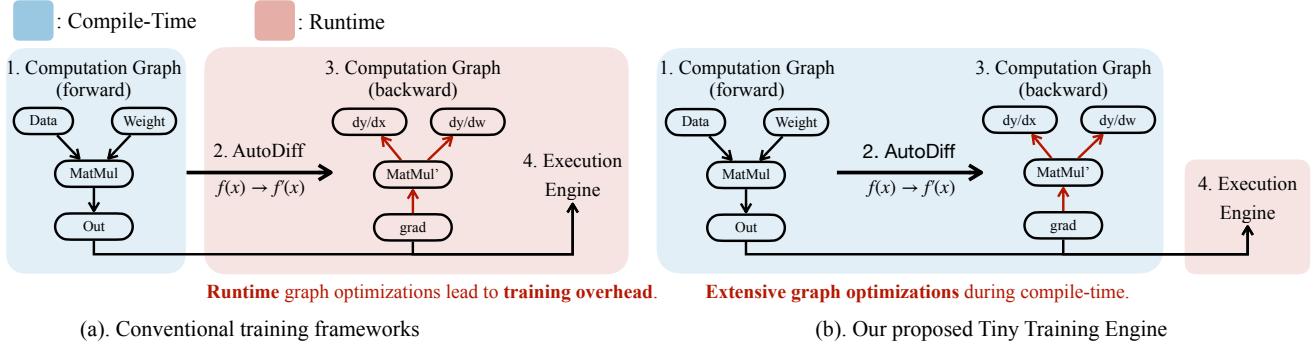


Figure 7. The comparison between runtime auto-differentiation and our compile-time differentiation. By offloading the differentiation to compile time, PockEngine not only simplifies the runtime, but also enables plenty of optimization opportunities, which will be detailed in Section 3.2.

To address the overhead, PockEngine obtains all tensor information and plans for a better execution schedule. By reordering operators, the gradients can be immediately applied to the corresponding parameters before back-propagating to earlier layers. We further trace the life-cycle of all tensors (weights, activations, gradients) and re-order the schedules to reduce memory usage, leading up to 21x savings on microcontrollers for MCUNet.

Operator Fusion. In most deep learning frameworks, a simple operation usually requires a number of fine-grained kernels to implement. For example, a single-layer normalization operation requires three kernel calls and two memory reads and writes for forward, and six kernel calls and five memory reads and writes for backward. Moreover, transformations such as fusing cheap operations into expensive ones (e.g. CONV-BN-ReLU), and parallel linear operations (e.g. batch matmul) have been shown effective in improving the inference. During compilation and codegen, PockEngine fuse these kernels into a single one and results in less memory IO and kernel calls.

Functional-Preserving Graph Transformation. Existing DNN frameworks optimize a computation graph by applying rules either designed by domain experts [2, 5] or automatically discovered by program [27, 29]. There are more optimization opportunities but previous research is unable to utilize them since the backward graph was derived at runtime in earlier frameworks. Extensive investigation of potential graph optimizations will lead to slow training and incur undesired runtime overhead.

Our engine integrates these optimization techniques and is an early trial to apply to the training graph. PockEngine transforms the data layout for different hardware. For vision tasks, NCHW is the most widely used layout. But this format is only efficient on accelerators like GPU. When training on mobile CPUs / DSPs, such format is no longer optimal and PockEngine will transform the layout at compile-time to facilitate runtime training efficiency.

Furthermore, PockEngine explores different implementations of kernels. For example, Winograd has been widely used in inference because of its faster computation. However, the savings are not free: it requires extra pre-processing of the weights. If the weights are not static, then the transformation needs to be applied every epoch and the total FLOPs can be even higher than normal convolution.

Hence it was utilized in inference and not incorporated into training frameworks. For on-device training scenarios, there are many frozen layers where the weights are not being changed during training [11, 62]. These layers in fact can utilize Winograd to accelerate but such opportunities are ignored in current frameworks even if the `requires_grad` attribute is set to `False`. PockEngine obtains the complete training graph during compile-time thus knowing the updating information of each parameter. Therefore, we can analyze the tensor and graph information, knowing whose weights are static and whose are dynamic. PockEngine can bind operation to the fastest implementation and enable the chance to utilize Winograd even in the training.

4 RESULTS

In this section, we comprehensively evaluate the performance of PockEngine. We first study the effectiveness of sparse backpropagation, then present the experimental results on different hardware and platforms, compared with other training frameworks. Finally, we discuss the graph optimization results.

4.1 Setups

Models. We evaluate PockEngine on popular vision and language models. For vision tasks, we choose MCUNet [41] (5FPS model), MobilenetV2 [51] (width multiplier 0.35 and 1.0), and ResNet-50 [23]. All normalization layers (e.g. BatchNorm) are fused into the linear operations (e.g. Conv, Linear). For masked language models, we choose the base-uncased version of BERT [19] and DistilBERT [52] to benchmark the performance.

Datasets. For vision models, we first pre-trained them on ImageNet [18] with resolution 224×224 (except 128×128 for MCUNet), and then fine-tuned on a set of downstream tasks to evaluate the transfer learning accuracy (including Cars [32], CIFAR-10 [33], CUB [61], Flowers [45], Foods [9], Pets [46], and VWW [17] conventional TinyML setting used in [11, 42]). The NLP models (BERT and DistilBERT) are pre-trained on Wikipedia and BookCorpus [66]. We evaluate their transfer learning performance on the GLUE [59] benchmark (including CoLA, MNLI, MRPC, QNLI, QQP, RTE, SST-2). For the chatbot models, we use LlamaV2 [56] then fine-tuned with instructions from Stanford Alpaca dataset [55]. We follow

Table 2. Sparse BP achieves comparable transfer learning performance (< 1% degradation on average) compared to the full update for vision models at various scales, while reducing the cost of on-device training.

Vision Model	Method	Avg. Acc	On-Device Training Vision Datasets						
			Cars	CIFAR	CUB	Flowers	Foods	Pets	VWW
MCUNet-5FPS [41]	Full BP	74.1%	56.7 ± 1.1%	86.0 ± 0.7%	56.2 ± 0.5%	88.8 ± 0.2%	67.1 ± 0.3%	79.5 ± 0.4%	88.7 ± 0.3%
	Bias Only	72.7%	52.4 ± 1.4%	83.4 ± 0.5%	55.2 ± 0.6%	86.7 ± 0.4%	65.0 ± 0.4%	78.0 ± 0.3%	88.1 ± 0.3%
	Sparse BP	74.8%	55.2 ± 1.3%	86.9 ± 0.6%	57.8 ± 0.4%	89.1 ± 0.3%	64.4 ± 0.3%	80.9 ± 0.3%	89.3 ± 0.4%
MobileNetV2 [51]	Full BP	89.2%	87.1 ± 0.9%	96.0 ± 0.5%	76.6 ± 0.8%	95.4 ± 0.2%	83.9 ± 0.2%	90.7 ± 0.4%	94.5 ± 0.2%
	Bias Only	87.3%	85.8 ± 0.8%	94.0 ± 0.7%	74.5 ± 0.7%	95.1 ± 0.5%	82.0 ± 0.6%	87.6 ± 0.5%	92.4 ± 0.3%
	Sparse BP	88.5%	86.4 ± 1.0%	95.0 ± 0.9%	76.4 ± 1.0%	95.4 ± 0.3%	81.5 ± 0.5%	90.4 ± 0.3%	94.2 ± 0.3%
ResNet-50 [23]	Full BP	90.5%	88.2 ± 0.5%	96.8 ± 0.4%	79.9 ± 0.6%	94.2 ± 0.3%	85.2 ± 0.4%	93.6 ± 0.2%	95.3 ± 0.1%
	Bias Only	87.8%	84.3 ± 0.6%	93.7 ± 0.7%	75.0 ± 0.3%	92.5 ± 0.5%	83.7 ± 0.3%	91.8 ± 0.4%	93.8 ± 0.1%
	Sparse BP	90.3%	86.7 ± 0.7%	96.2 ± 0.6%	81.0 ± 0.7%	95.6 ± 0.3%	84.0 ± 0.3%	93.4 ± 0.5%	95.1 ± 0.1%

Table 3. For language models, sparse BP maintains the fine-tuning accuracy for at a reduced training cost. Results are reported with mean and standard deviation for 3 runs.

Language Model	Method	Avg.	On-Device Training Language Datasets						
			CoLA	MNLI	MRPC-acc	QNLI	QQP-acc	RTE	SST-2
Distill-BERT [52]	Full BP	76.9%	46.6 ± 1.2%	81.9 ± 0.2%	83.8 ± 1.9%	88.3 ± 0.1%	90.0 ± 0.2%	59.6 ± 1.9%	90.8 ± 0.8%
	Bias Only	72.8%	44.6 ± 0.9%	73.2 ± 0.7%	78.9 ± 2.3%	83.4 ± 1.4%	83.6 ± 0.5%	57.8 ± 2.5%	88.0 ± 1.3%
	Sparse BP	77.0%	47.9 ± 1.5%	81.1 ± 0.3%	84.2 ± 1.8%	87.8 ± 0.1%	88.5 ± 0.3%	58.0 ± 1.6%	90.6 ± 0.5%
BERT [19]	Full BP	81.8%	59.9 ± 1.5%	84.0 ± 0.1%	85.8 ± 1.9%	90.9 ± 0.2%	90.8 ± 0.3%	68.2 ± 2.0%	92.7 ± 0.7%
	Bias Only	78.1%	51.1 ± 0.5%	78.6 ± 0.8%	83.6 ± 2.6%	88.5 ± 1.0%	86.0 ± 0.1%	67.9 ± 3.3%	90.7 ± 1.3%
	Sparse BP	81.7%	58.6 ± 0.8%	84.4 ± 0.2%	86.2 ± 1.6%	90.8 ± 0.1%	90.3 ± 0.6%	69.4 ± 1.8%	91.8 ± 0.4%

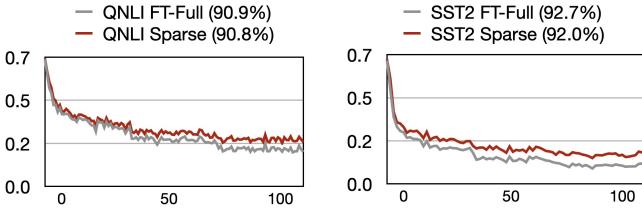


Figure 8. The training loss curves of FT-Full and our used sparse update on QNLI and SST-2 dataset using BERT. Sparse updates slightly slow down the training curve, but do not degrade the final accuracy

alpaca-eval [37] and MT-Bench [16, 63] to evaluate the response quality.

Runtime adaptations. PockEngine is a compilation-based framework, and the compilation workflow helps to handle various frontends as well as adapt to different backends and runtimes. PockEngine takes models defined in PyTorch/TensorFlow/Jax and translates them into IR where graph optimizations can be applied. The optimized training graph is then fed to deep learning libraries like TVM [14], SNPE [49] and TinyEngine [40] to generate platform-specific binaries.

Sparse-BP Schemes for Fine-tuning. We obtain sparse backpropagation schemes for CNN models (MCUNet [41], MobileNetV2 [51], ResNet-50 [23]) and transformer [58]-based NLP models

(BERT [19] and DistilBERT [52]) by sensitivity analysis [42] while considering computation/memory cost. Our final update schemes are:

- MCUNet: we further support sparse tensor backpropagation to handle the extreme memory constraints of IoT devices [41]. We update the biases of the last 7 blocks and update {100%, 100%, 50%, 100%} of the weights of the first convolutions for the intermediate 4 blocks.
- MobileNetV2: update the *biases* and the *weights* of the first 1x1 convolution for the last 7 blocks (out of 19).
- ResNet-50: update the *biases* and the *weights* of the first 1x1 convolution for the last 8 blocks (out of 16).
- BERT: update the *biases* of the last 6 blocks (out of 12) and the *weights* of the attention module and the first linear in FFN for the last 4 blocks.
- Distill-BERT: update the *biases* of the last 3 blocks (out of 6) and the *weights* of the attention module and the first linear in FFN for the last 2 blocks.
- LlamaV2-7B: update the *biases* of the last 5 blocks (out of 32) and the *weights* of the attention module and the first linear in FFN for the last 5 blocks.

4.2 Effectiveness of Sparse Backpropagation

Sparse Backpropagation Achieves Full Accuracy. On-device training aims to continually improve user experience by local data thus the number of training samples is usually small compared

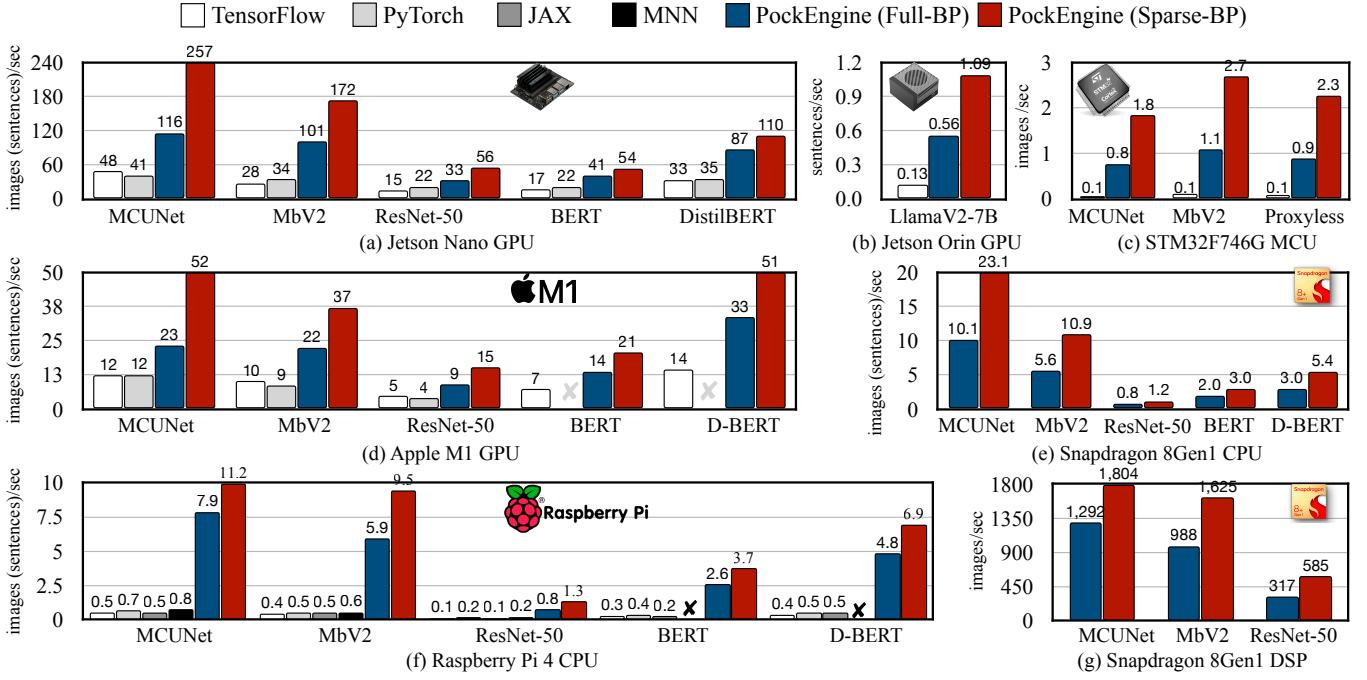


Figure 9. Training speed comparison between other frameworks and PockEngine of popular deep learning models on various hardware platforms. PockEngine consistently outperforms existing frameworks and sparse bp further speedups the training throughput.

to large-scale pre-training. Therefore, it is not necessary to perform full backpropagation to update the whole model as shown in Table 2. Remarkably, the downstream accuracy of sparse backpropagation can match the full-propagation baselines on both vision and language models (<1% performance drop). On some downstream datasets, the performance of sparse backpropagation is even higher surpassing the full baselines such as Flower in vision and mrpc-acc in language. The performance is far above the common requirements for TinyML [8] (80% accuracy on VWW), suggesting sparse propagation is a good strategy for on-device training.

Furthermore, when evaluating language models, sparse backpropagation also maintains the finetuning accuracy at a reduced training cost. The average performance degradation is within 1%. This means that the use of sparse backpropagation can effectively reduce the time and cost required for training language models, without sacrificing accuracy. In fact, the results show that sparse backpropagation can even improve the model’s performance on certain sub-tasks (e.g. MRPC and RTE). By making training more efficient, sparse backpropagation could help to accelerate progress in these fields and enable the development of more advanced language models.

Sparse Backpropagation Reduces Training Time and Memory. Besides the comparable performance when transferring to downstream tasks, sparse backpropagation greatly reduces the training peak memory and improves the training speed.

Shown in Table 4, the training memory grows rapidly w.r.t the batch size and soon exceeds the limitation for edge devices (e.g. 1GB for Raspberry Pi), using swap or re-materialization techniques [48]

will introduce extra computation and energy cost. Sparse backpropagation cuts down peak memory usage (2.2× to 21.3×) and the saving is general across models and applications. Even when batch size grows, the required memory is still small and the memory cost of training MCUNet-5FPS sparse-BP with batch size 8 is still smaller than batch 1. Batched training helps improve device utilization as well as training stability.

When applying sparse backpropagation, operations and tensors related to frozen layers are automatically trimmed from the training graph via dead-code elimination, resulting in less computation and higher training throughput. Figure 9 shows that the sparse backpropagation can further accelerate the speed by 1.3x to 1.6x on Raspberry Pi 4. Previous efficient training algorithms only discuss the theoretical performance and PockEngine provides system-level support and translates into measured reduction.

4.3 PockEngine Speedups On-Device Training

We compare PockEngine with other training frameworks in Figure 9. PockEngine enables training on various hardware platforms, including Raspberry Pi 4, Snapdragon CPU and DSP, Apple M1, Jetson Nano, and microcontrollers. It also supports a wide range of models, such as MCUNet, MobilenetV2, ResNet-50, BERT, and DistilBERT. PockEngine effortlessly supports diverse models through its frontend, which converts neural networks represented in various formats to a unified intermediate representation.

Furthermore, the compilation-based workflow allows us to choose the best runtime backend for different training scenarios, including both vendor libraries (e.g. SNPE for Snapdragon GPUs and DSPs, TensorRT for NVIDIA GPUs) and customized kernels (e.g.,

Table 4. The training memory usage comparison of full backpropagation and sparse backpropagation. We report actual memory usage measured on Jetson AGX Orin. The saving ratios are more significant as batch sizes increase. “-” denotes that the experiments cannot fit into devices.

Platform	Model	#Params	Method	Training Memory		
				bs=1	bs=4	bs=16
MCU	MCUNet	0.6M	Full-BP	3.6MB	-	-
			Sparse-BP	173KB	-	-
Jetson Nano	MobilenetV2	3.4M	Full-BP	729MB	910MB	1.2GB
			Sparse-BP	435MB	501MB	819MB
Jetson Nano	ResNet50	26M	Full-BP	827MB	1.1GB	2.1GB
			Sparse-BP	663MB	723MB	885MB
Jetson AGX Orin	BERT	125M	Full-BP	1.7GB	3.6GB	5.7GB
			Sparse-BP	1.4GB	1.9GB	2.3GB
Jetson AGX Orin	LlamaV2	7B	Full-BP	43.1GB	-	-
			Sparse-BP	31.2GB	-	-

TVM-tuned kernels for ARM CPUs and Apple M1). We present a comparison of training workflows in Figure 9 and discuss it below:

Edge CPU. For platforms like the Raspberry Pi, PockEngine offers 13 to 21 × better performance compared to popular DNN training frameworks. This speedup is due to kernel tuning, which existing frameworks either overlook in favor of GPU kernel implementations (PyTorch, TensorFlow, Jax) or optimize only for the inference pipeline and operators (MNN). The corresponding ARM kernels do not provide ideal performance, let alone the overhead brought by frameworks.

Edge GPU. We benchmark edge GPU platforms using NVIDIA Jetson Nano and Jetson AGX Orin due to their widespread use in edge applications. GPUs have a much higher degree of parallelism and better training throughput than CPUs. The faster training speed of PockEngine (2.2x to 2.6× speedup) is mainly due to the compilation process: The host language Python is typically slow on low-frequency CPUs, while PockEngine’s compiled graph can run without host languages. While other frameworks like TensorRT [2] may also achieve this, they are limited to inference only and do not provide training support.

Apple M-Chip. The Apple M1 chip is a relatively new platform for training. While PyTorch and Tensorflow have preliminary GPU support, the compatibility is not ideal [‡]. Even with the latest build (commit ID: *c9913cf*), PyTorch throws errors when launching training for BERT and DistilBERT. On the other hand, PockEngine compiles the training graph to Metal, providing better compatibility and faster training speeds.

Mobile DSP. For Qualcomm DSP, we integrate SNPE [49] to deliver the final binaries. It is worth noting that SNPE is a conventionally inference-only library for integer models and our PockEngine easily extends it with training capability. As shown in Figure. 9 (g), the peak performance of DSP is impressive and even on par with edge GPUs.

Microcontrollers. For the microcontroller platform, we integrate TinyEngine [41] to perform the codegen and enable training under extremely limited memory constraints. Previous frameworks like TF-Lite-Micro [4] is *inference-only* and we report the projected latency. Show in Figure. 7 (c), the speed is much lower than Pock-Engine.

PockEngine enables efficient on-device training by compilation and adaptation to various runtimes. It further supports advanced backpropagation schemes and advanced graph optimization, which we will expand further in the following section.

5 FINE-TUNING CHATBOT WITH POCKENGINE

With the growing attention ChatGPT has received, the demand for fine-tuning one’s own Chatbot models has also been increasing. This allows users to tailor the model to their domain-specific needs (e.g., law, biomedical, health care) and ensures privacy (e.g. private emails, personal assistant) by not uploading information to the cloud. By fine-tuning our own language model, we can address these concerns and obtain a high-quality language model that meets our needs. In this section, we demonstrate how PockEngine can efficiently fine-tune a chatbot on edge platform (Jetson AGX Orin).

Models. We choose Meta’s LlamaV2 [56] and choose the 7B model as the backbone for our experiments. This decision was based on the trade-off of model quality and device resources. The detailed fine-tuning settings are discussed below.

Evaluation. For evaluation, we follow Alpaca-Eval [37] and MT-Bench [63] to use LLMs as the automated evaluator for benchmark generation and performance assessments. The quality of the answers is evaluated based on helpfulness, relevance, accuracy, and details from 805 questions[§] and 80 questions from Vicuna project[¶]. This is a pair-to-pair comparison and we choose *text-davinci-003* for Alpaca-Eval win rate (%) and *ChatGPT-3.5-Turbo* for MT-Bench Score.

[§]https://tatsu-lab.github.io/alpaca_eval

[¶]<https://github.com/lm-sys/FastChat>

[‡]<https://github.com/pytorch/pytorch/issues/77764>

Table 5. Instruction tuning comparisons between PyTorch and PockEngine. The pre-trained model is LlamaV2-7B [56] and we fine-tune the models following Stanford Alpaca’s setting [55]. We report the training loss and Alpaca-eval score [37] (reference model *text-davinci003*. PockEngine shows significant speedup over PyTorch on Jetson AGX Orin while fully matching the training quality. With the sparse update, PockEngine further improves the training throughput while maintaining the response quality.

Framework	Method	Iteration Latency (↓)	GPU Memory(↓)	Loss(↓)	Alpaca-Eval Winrate(↑)	MT-Bench score(↑)
PyTorch	FT-Full	7.7s	45.1GB	0.761	44.1%	6.1
PyTorch	LoRA (rank=8)	7.3s	30.9GB	0.801	43.1%	5.1
PockEngine	FT-Full	1.8s	43.1GB	0.768	43.7%	6.1
PockEngine	Sparse	0.9s	31.2GB	0.779	43.1%	5.7

Datasets. To align pre-trained language models with instructions, we follow the self-instruct [60] and adapt data from Stanford Alpaca [55]. The total training set has 52K examples containing diverse instructions and responses¹

Instruction: What is the meaning of the following idiom?
Input: It's raining cats and dogs.
Output: The idiom "it's raining cats and dogs" means that it is raining heavily.

Example Record from Alpaca Dataset.

Fine-tuning. We fine-tune the models for 3 epochs using a learning rate of 10^{-4} and no weight decay. The optimizer we use is memory-efficient Lion [15], and the maximum sentence length is limited to 512. The instruction tuning batch size is 1, and the gradient is accumulated over 16 steps. We sparsely update the biases of the last 5 blocks (out of 24) and the weights of the attention module and the first linear layer in the FFN for the last 5 blocks. We further freeze the layer-norm layers to reduce training costs and speed up training.

5.1 Quantitative Comparison.

PockEngine Accelerates Training. As shown in Table 5, PyTorch can train on Jetson AGX Orin, but one iteration takes more than 7 seconds for LlamaV2-7B. Fine-tuning on 1000 records would require 2 hours while PockEngine accelerates training by 4.4x and can finish in less than half an hour.

Sparse Backpropagation Accelerates Training. For popular parameter-efficient fine-tuning methods like LoRA [25], although they can effectively reduce the memory footprint (from 45.1GB to 30.9GB), the training cost is not significantly improved as they still need to backpropagate the first layer. In contrast, Sparse backpropagation reduces the backpropagation depth and significantly improves training speed (from 1768ms to 914ms, 1.9x faster).

Sparse-BP Achieves Comparable Accuracy. Besides training throughput improvement, sparse backpropagation also maintains fine-tuning accuracy. When compared to full-BP, sparse-BP demonstrates similar performance, achieving an impressive Alpaca-Eval score of 43.7. This score closely matches the performance of full-BP, which has an Alpaca-Eval score of 44.1. Sparse-BP also performs favorably when compared to LoRA (Alpaca-Eval score 43.1).

¹<https://huggingface.co/datasets/tatsu-lab/alpaca>

5.2 Qualitative Comparison between Full-BP and Sparse-BP

In addition to quantitative analysis on accuracy and throughput, we also conduct qualitative analysis to compare the quality of sparse-BP and full-BP when performing instruction tuning on LlamaV2-7B. We make the assumption that a model’s knowledge and capabilities are learned primarily during pretraining, while alignment teaches LLMs how to interact with users [64]. Therefore, we examine the response quality from multiple perspectives:

Factual Recall. : To thoroughly evaluate the performance between Full- and Sparse-backpropagation, we conduct a detailed analysis of the corresponding models’ proficiency to answer questions related to general knowledge. We first examine the capability to correctly identify the capital city of a given country. All models consistently deliver accurate responses while the tuned models not only provide accurate answers but also offer more comprehensive information, enriching the overall user experience.

Input: What is the capital of United Status?
LlamaV2-7B-original: Washington, D.C.
LlamaV2-7B-tuned (Full-BP): The capital of the United States is Washington D.C. (District of Columbia).
LlamaV2-7B-tuned (Sparse-BP): The capital of the United States is Washington, D.C.

Example 1.The Capital of U.S.A.

Similarly, when we inquire about the dates of notable events, all models exhibit a remarkable level of confidence and accuracy in their responses. This suggests the reliability and correctness of both Full- and Sparse-backpropagation during fine-tuning.

Input: In what year was the Beijing Olympics held?
LlamaV2-7B-original: The Beijing Olympics were held in 2008..
LlamaV2-7B-tuned (Full-BP): The Beijing Olympics, also known as the Games of the XXIX Olympiad, were held in 2008. The opening ceremony took place on August 8, 2008, and the closing ceremony was held on August 24, 2008.
LlamaV2-7B-tuned (Sparse-BP): The Beijing Olympics, also known as the 2008 Summer Olympics, were held in Beijing, China in 2008.

Example 2. The Date of Beijing Olympics.

However, when the question becomes obscure, the original Llama model no longer provides accurate answers. For example, the raw Llama model (without fine-tuning) incorrectly states that Michael Jackson died in 2020 (the actual year is 2009), and it also provides

incorrect album names (the proper title is "Invincible"). Meanwhile, the Full-BP model and the Sparse-BP model both properly interpret the questions and provide accurate responses.

Input: what was Michael Jackson's last album?

LlamaV2-7B-original: Michael Jackson passed away on June 25, 2020, so he did not release an album after that. However, before his death, he released the album "Michael Jackson: The Ultimate Collection," which was a compilation of his previous studio albums and bonus tracks.

LlamaV2-7B-tuned (Full-BP): Michael Jackson's last album was 'Invincible', released in 2001.

LlamaV2-7B-tuned (Sparse-BP): Michael Jackson's last album, released during his lifetime, was "Invincible" which was released in 2001.

Example 3. The Title of Michael Jackson's Final Album

Operation and Logic Inference. Next, we evaluate a simple but interesting task: reverse the order of words. This task requires LLMs to understand the concept of reversing and properly segmenting the sentences. The original model misunderstood the task and reversed the words' meanings, leading to incorrect answers. Meanwhile, the Full-BP and sparse-BP models, which were fine-tuned, both correctly understood the term "reverse" and executed the given task, though the sparse-BP model made a tiny mistake and still misplaced the word "the".

Input: Please reverse the words in the sentence "I love the Micro conference"

LlamaV2-7B-original: I hate the conference Micro.

LlamaV2-7B-tuned (Full-BP): The reversed result is "Conference Micro the love I"

LlamaV2-7B-tuned (Sparse-BP): The reversed sentence would be "Conference the Micro love I".

Example 4. Reversing a sentence.

After comparing the response quality of executing given tasks, we then evaluated the logical inference capability of the fine-tuned models. We first briefly introduced the hosting location of the Micro conference for the past three years and asked where the conference was held last year. Without fine-tuning, the original Llama model failed to understand the question and only simply repeated the first sentence. However, for the fine-tuned models, both Full-BP and Sparse-BP were able to correctly name the place, with the response from Full-BP being more confident.

Input: Micro was running virtually online in 2021, and held physically at Chicago in 2022 and will soon be hosted in Toronto in 2023. Now is the year of 2023, where was Micro hosted last year?

LlamaV2-7B-original: Micro was hosted virtually in 2021

LlamaV2-7B-tuned (Full-BP): Micro was hosted physically in Chicago in 2022.

LlamaV2-7B-tuned (Sparse-BP): Based on the information provided, it seems that Micro was hosted physically at Chicago in 2022.

Example 5. Inferring Information Based on Context.

Maths. We then compared the mathematical abilities of different models. Mathematics is often considered a significant challenge for many language models, and obviously, the original LlamaV2 model struggled with math-related questions, while the Full-BP and

Sparse-BP models demonstrated certain capabilities for reasoning and properly calculating the results.

Input: Lee mows one lawn and charges \$33. Last week he mowed 16 lawns and three customers each gave him a \$10 tip. How many dollars did Lee earn mowing lawns last week?

LlamaV2-7B-original: Lee earned \$660 from mowing lawns last week.

LlamaV2-7B-tuned (Full-BP): Lee earned $\$33 \times 16 = \528 from mowing lawns. His total earnings from the 16 customers who gave him a \$10 tip is $\$10 \times 3 = \30 . Therefore, Lee earned a total of $\$528 + \$30 = \$558$ last week mowing lawns.

LlamaV2-7B-tuned (Sparse-BP): To determine how much Lee earned mowing lawns last week, we need to calculate the total amount of money he earned from mowing lawns and tips. Lee mowed one lawn and charged \$33, so he earned \$33 from that lawn. Last week, Lee mowed 16 lawns, and three customers each gave him a \$10 tip, so he earned $\$10 \times 3 = \30 in tips from those lawns and earned \$33 from mowing one lawn and total 16 lawns = \$528. Therefore, Lee earned \$558 last week from mowing lawns and receiving tips.

Example 6. Math Problem Solving

Note that this is a concise qualitative study comparing original, Full-BP fine-tuned, and Sparse-BP fine-tuned LLMs. We carefully selected representative samples for this study, although it is important to note that it is not comprehensive given the extensive range of responses the model can provide. The objective of this analysis is to present compelling evidence in support of two findings: (1) fine-tuning is an essential process for personalizing your own Chabot, and (2) Sparse-BP is capable of fine-tuning models with comparable quality with much reduced cost.

6 CONCLUSION

We present PockEngine, an efficient training framework for learning on edge. PockEngine has general support for various frontends/backends to deal with hardware heterogeneity on edge. It improves the efficiency of on-device training via (1) compilation-based auto-differentiation to offload overheads from runtime to compile time; (2) supporting sparse backpropagation with backward graph pruning; (3) training graph optimization including operator reordering/fusion and various function-preserving transforms.

Experiments on different edge devices show PockEngine can significantly speedup on-device training: 11.2 \times on ARM CPUs, 2 \times on Apple M1, and 2.7 \times on NVIDIA edge GPU, and 9.6 \times on microcontroller compared to TensorFlow. PockEngine supports sparse backpropagation, which further speeds up by 1.5 - 3.5 \times while matching the accuracy of full backpropagation. Further, PockEngine enables fine-tuning LLamaV2-7B language model on a Jetson AGX Orin at 914ms, 7.9 \times faster than the PyTorch baseline. We hope our engine design can facilitate AI applications with personalization and life-long learning capacity by democratizing learning on the edge.

ACKNOWLEDGMENTS

This work was supported by MIT-IBM Watson AI Lab, MIT AI Hardware Program, MIT-Amazon Science Hub, and NSF. Ligeng Zhu and Ji Lin were partially supported by the Qualcomm Innovation Fellowship.

REFERENCES

- [1] [n. d.]. NCNN : A high-performance neural network inference computing framework optimized for mobile platforms. <https://github.com/Tencent/ncnn>.
- [2] [n. d.]. NVIDIA TensorRT, an SDK for high-performance deep learning inference. <https://developer.nvidia.com/tensorrt>.
- [3] [n. d.]. TensorFlow Lite. <https://www.tensorflow.org/lite>.
- [4] [n. d.]. TensorFlow Lite Micro. <https://github.com/tensorflow/tflite-micro>.
- [5] Martin Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. <https://www.tensorflow.org/>. Software available from tensorflow.org.
- [6] Martin Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. 2016. Tensorflow: A system for large-scale machine learning. In *OSDI*.
- [7] Byung Hoon Ahn, Jinwon Lee, Jamie Menjy Lin, Hsin-Pai Cheng, Jilei Hou, and Hadi Esmaeilzadeh. 2020. Ordering chaos: Memory-aware scheduling of irregularly wired neural networks for edge devices. *arXiv preprint arXiv:2003.02369* (2020).
- [8] Colby R Banbury, Vijay Janapa Reddi, Max Lam, William Fu, Amin Fazel, Jeremy Hollerman, Xinyuan Huang, Robert Hurtado, David Kanter, Anton Lokhmotov, et al. 2020. Benchmarking TinyML systems: Challenges and direction. *arXiv preprint arXiv:2003.04821* (2020).
- [9] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. 2014. Food-101—mining discriminative components with random forests. In *European conference on computer vision*. Springer, 446–461.
- [10] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. 2018. JAX: composable transformations of Python+NumPy programs. <http://github.com/google/jax>
- [11] Han Cai, Chuang Gan, Ligeng Zhu, and Song Han. 2020. TinyTL: Reduce Activations, Not Trainable Parameters for Efficient On-Device Learning. *arXiv preprint arXiv:2007.11622* (2020).
- [12] Ken Chatfield, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2014. Return of the devil in the details: Delving deep into convolutional nets. In *BMVC*.
- [13] Tianqi Chen, Mu Li, Yutian Li, Min Lin, Naiyan Wang, Minjie Wang, Tianjun Xiao, Bing Xu, Chiyuan Zhang, and Zheng Zhang. 2015. Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv preprint arXiv:1512.01274* (2015).
- [14] Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Meghan Cowan, Leyuan Wang, Yuwei Hu, Luis Ceze, et al. 2018. {TVM}: An automated end-to-end optimizing compiler for deep learning. In *OSDI*.
- [15] Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Yao Liu, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, and Quoc V. Le. 2023. Symbolic Discovery of Optimization Algorithms. *arXiv:2302.06675* [cs.LG]
- [16] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An Open-Source Chatbot Impressing GPT-4 with 90%* ChatGPT Quality. <https://vicuna.lmsys.org>
- [17] Aakanksha Chowdhery, Pete Warden, Jonathon Shlens, Andrew Howard, and Rocky Rhodes. 2019. Visual wake words dataset. *arXiv preprint arXiv:1906.05721* (2019).
- [18] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR*.
- [19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [20] Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Darrell. 2014. Decaf: A deep convolutional activation feature for generic visual recognition. In *ICML*.
- [21] Jonathan Frankle, David J Schwab, and Ari S Morcos. 2020. Training BatchNorm and Only BatchNorm: On the Expressive Power of Random Features in CNNs. *arXiv preprint arXiv:2003.00152* (2020).
- [22] Chuang Gan, Naiyan Wang, Yi Yang, Dit-Yan Yeung, and Alex G Hauptmann. 2015. Devnet: A deep event network for multimedia event detection and evidence recounting. In *CVPR*, 2568–2577.
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *CVPR*.
- [24] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for NLP. In *International Conference on Machine Learning*. PMLR, 2790–2799.
- [25] Edward Hu, Yelong Shen, Phil Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Lu Wang, and Weizhu Chen. 2021. LoRA: Low-Rank Adaptation of Large Language Models. *arXiv:2106.09685* [cs.CL]
- [26] Junsu Jang and Fadel Adib. 2019. Underwater backscatter networking. In *Proceedings of the ACM Special Interest Group on Data Communication*, 187–199.
- [27] Zhihao Jia, Oded Padon, James Thomas, Todd Warszawski, Matei Zaharia, and Alex Aiken. 2019. TASO: optimizing deep learning computation with automatic generation of graph substitutions. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*, 47–62.
- [28] Zhihao Jia, James Thomas, Todd Warszawski, Mingyu Gao, Matei Zaharia, and Alex Aiken. 2019. Optimizing DNN computation with relaxed graph substitutions. *Proceedings of Machine Learning and Systems* 1 (2019), 27–39.
- [29] Zhihao Jia, James Thomas, Todd Warszawski, Mingyu Gao, Matei Zaharia, and Alex Aiken. 2019. Optimizing DNN computation with relaxed graph substitutions. *Proceedings of Machine Learning and Systems* 1 (2019), 27–39.
- [30] Xiaotang Jiang, Huan Wang, Yiliu Chen, Ziqi Wu, Lichuan Wang, Bin Zou, Yafeng Yang, Zongyang Cui, Yu Cai, Tianhang Yu, et al. 2020. MNNT: A universal and efficient inference engine. *arXiv preprint arXiv:2002.12418* (2020).
- [31] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [32] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 2013. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, 554–561.
- [33] Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learning multiple layers of features from tiny images. (2009).
- [34] Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. 2022. Fine-tuning can distort pretrained features and underperform out-of-distribution. *arXiv preprint arXiv:2202.10054* (2022).
- [35] Yoonho Lee, Annie S Chen, Fahim Tajwar, Ananya Kumar, Huaxiu Yao, Percy Liang, and Chelsea Finn. 2022. Surgical Fine-Tuning Improves Adaptation to Distribution Shifts. *arXiv preprint arXiv:2210.11466* (2022).
- [36] Philip Lewis, Neil Patel, David Culler, and Scott Shenker. 2004. Trickle: A self-regulating algorithm for code propagation and maintenance in wireless sensor networks. In *Proc. of the 1st USENIX/ACM Symp. on Networked Systems Design and Implementation*, Vol. 25, 37–52.
- [37] Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. AlpacaEval: An Automatic Evaluator of Instruction-following Models. https://github.com/tatsu-lab/alpaca_eval
- [38] Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190* (2021).
- [39] Edgar Liberis and Nicholas D Lane. 2019. Neural networks on microcontrollers: saving memory at inference via operator reordering. *arXiv preprint arXiv:1910.05110* (2019).
- [40] Ji Lin, Wei-Ming Chen, Han Cai, Chuang Gan, and Song Han. 2021. Mcunetv2: Memory-efficient patch-based inference for tiny deep learning. *arXiv preprint arXiv:2110.15352* (2021).
- [41] Ji Lin, Wei-Ming Chen, Yujun Lin, John Cohn, Chuang Gan, and Song Han. 2020. Meunet: Tiny deep learning on iot devices. In *NeurIPS*.
- [42] Ji Lin, Ligeng Zhu, Wei-Ming Chen, Wei-Chen Wang, Chuang Gan, and Song Han. 2022. On-Device Training Under 256KB Memory. In *NeurIPS*.
- [43] Pramod Kaushik Mudrakarta, Mark Sandler, Andrey Zhmoginov, and Andrew Howard. 2018. K for the price of 1: Parameter-efficient multi-task and transfer learning. *arXiv preprint arXiv:1810.10703* (2018).
- [44] Pramod Kaushik Mudrakarta, Mark Sandler, Andrey Zhmoginov, and Andrew Howard. 2019. K for the Price of 1: Parameter-efficient Multi-task and Transfer Learning. In *ICLR*.
- [45] 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.
- [46] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. 2012. Cats and dogs. In *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 3498–3505.
- [47] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019).
- [48] Shishir G Patil, Paras Jain, Prabal Dutta, Ion Stoica, and Joseph Gonzalez. 2022. POET: Training Neural Networks on Tiny Devices with Integrated Rematerialization and Paging. In *International Conference on Machine Learning*. PMLR, 17573–17583.
- [49] Qualcomm. [n. d.]. Snapdragon Neural Processing Engine SDK. <https://developer.qualcomm.com/sites/default/files/docs/snpe/overview.html>
- [50] Alex Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. (2018).
- [51] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. 2018. MobileNetV2: Inverted Residuals and Linear Bottlenecks. In *CVPR*.

- [52] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108* (2019).
- [53] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. 2014. CNN features off-the-shelf: an astounding baseline for recognition. In *CVPR Workshops*.
- [54] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. 2020. Test-time training with self-supervision for generalization under distribution shifts. In *International conference on machine learning*. PMLR, 9229–9248.
- [55] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford Alpaca: An Instruction-following LLaMA model. https://github.com/tatsu-lab/stanford_alpaca.
- [56] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* (2023).
- [57] Deepak Vasishth, Zerina Kapetanovic, Jongho Won, Xinxin Jin, Ranveer Chandra, Sudipta Sinha, Ashish Kapoor, Madhusudhan Sudarshan, and Sean Stratman. 2017. {FarmBeats}: An {IoT} Platform for {Data-Driven} Agriculture. In *14th USENIX Symposium on Networked Systems Design and Implementation (NSDI 17)*. 515–529.
- [58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *arXiv preprint arXiv:1706.03762* (2017).
- [59] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461* (2018).
- [60] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. Self-Instruct: Aligning Language Model with Self Generated Instructions. *arXiv:2212.10560* [cs.CL]
- [61] Peter Welinder, Steve Branson, Takeshi Mita, Catherine Wah, Florian Schroff, Serge Belongie, and Pietro Perona. 2010. *Caltech-UCSD Birds 200*. Technical Report CNS-TR-201. Caltech. /se3/wp-content/uploads/2014/09/WelinderEtal10_CUB-200.pdf. <http://www.vision.caltech.edu/visipedia/CUB-200.html>
- [62] Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. 2021. BitFit: Simple Parameter-efficient Fine-tuning for Transformer-based Masked Language-models. *CoRR abs/2106.10199* (2021). arXiv:2106.10199 <https://arxiv.org/abs/2106.10199>
- [63] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging LLM-as-a-judge with MT-Bench and Chatbot Arena. *arXiv:2306.05685* [cs.CL]
- [64] Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivas Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023. LIMA: Less Is More for Alignment. *arXiv:2305.11206* [cs.CL]
- [65] Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*. 19–27.
- [66] Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning Books and Movies: Towards Story-Like Visual Explanations by Watching Movies and Reading Books. In *The IEEE International Conference on Computer Vision (ICCV)*.