

banks on drones or solar-charged batteries in a forest. There can be *power limits* to avoid overheating of poorly ventilated IP-67 enclosures in an industrial site [15]. There can also be *time constraints*, e.g., to meet a training deadline for an application or limit the billed cost. In federated learning, such time estimates may also guide device selection [16, 17]. So, it is important to *configure* the edge accelerator to adapt to these constraints optimally and offer *predictable power, energy and time estimates* for DNN training.

1.1. Challenges

A unique feature of Nvidia’s Jetson edge devices is their fine-grained control over CPU, GPU and memory frequencies and active CPU cores, enabled by dynamically setting their *power modes*. E.g., the Orin AGX offers 29 CPU frequencies up to 2.2GHz, 13 GPU frequencies up to 1.3GHz, 4 memory frequencies and 12 CPU core-count settings for a total of $\approx 18,000$ possible power modes (Table 2). These can help limit the power envelope in a constrained setting while throttling the performance. However, given this vast parameter space, selecting an *optimal power mode* to make such trade-offs for training workloads that arrive dynamically is non-trivial. Some frequencies span an order of magnitude, and our experiments show that they can have up to 36 \times impact on DNN training time and 4.3 \times impact on power usage. For instance, using a low power mode to train the ResNet model on ImageNet data (see § 2 for details) on Orin AGX takes 112 mins per epoch and consumes about 11.8 W of power, while increasing it to a higher MAXN power mode reduces the training time to 3.1 mins but increases the power load to 51.1 W.

These also vary a lot across *DNN workloads* that may arrive continuously over time, e.g., the BERT model on SQuAD dataset using the same MAXN power mode for Orin AGX takes a much longer 68.7 mins with a power of 57 W compared to ResNet which only took 3.1 mins above. As expected, these also vary substantially across *Jetson device generations*, which may co-exist in a deployment as new devices are deployed or old ones upgraded in a rolling manner. E.g., the AGX Xavier takes 8.47 mins per epoch and 36.4 W to train the ResNet model on ImageNet dataset, as against 3.1 mins and 51.1 W for AGX Orin, on the MAXN power mode on both devices.

Hence, a wrong power mode choice can cause high power and performance penalties, violating time and power load constraints or, in the worst case, destroying the device due to overheating. So, accurately estimating the power and time taken for training a given DNN workload on the edge accelerator using a specific power mode is critical ¹.

A naïve solution of benchmarking all power modes for the device for each new DNN model is intractable. E.g., it takes 16.3 h to profile even 25% of power modes for Orin AGX for the ResNet model on ImageNet data, and this will not scale for every workload, especially if it arrives dynamically and has time constraints. Random or even targeted sampling of some

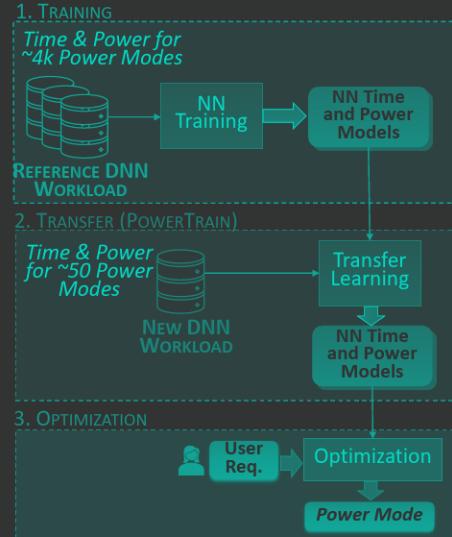


Figure 1: Time and Power Prediction Models using PowerTrain: (1) Initial training, (2) Transfer, and (3) Use for Optimization

power modes is inaccurate for optimizations, as we show later. Practical deployments may have multiple generations of accelerated devices, amplifying these costs. So, a *principled approach with low overheads and generalizability* is required to build prediction models to estimate training time and power for diverse DNN workloads on edge computing systems.

1.2. Gaps

There is substantial work on optimizing the energy consumption of DNN training on GPU servers [18] and predicting their runtime [19], and power and energy consumption [20]. But these do not translate directly to the edge due to architectural differences, such as the shared memory across CPU and GPU for edge accelerators, and the use of ARM-based rather than x86 CPUs. Similarly, several studies have examined energy and performance profiling for edge inferencing rather than training [21, 22, 23]. Some have also explored power modes and frequency scaling for micro-benchmarks on older Jetson architectures [3]. However, inferencing workloads have low compute demand compared to training, which can stress all resources of the edge device.

Our previous work focused on characterizing the performance of the previous generation Jetson Xavier AGX and Nano devices for DNN training, but offered only limited preliminary results on actual performance predictions or tuning of power modes [4]. Nvidia offers a limited set of pre-defined power mode choices (3 for Orin AGX apart from the default MAXN) with varying power budgets, but these are too few to allow targeted optimizations. As we show later, Nvidia’s own tool highly overestimates the power load for Orin AGX, causing excess training time when meeting a power budget. In summary, there is a clear need for modeling, prediction and optimization studies for DNN training workloads on accelerated edges, but none exist. *We address this gap, and offer the first principled modeling and prediction approach of its kind in this work.*

¹Since $\text{energy (mWh)} = \text{power (mW)} \times \text{time (h)}$, we focus on predicting power and time for training, and can derive energy estimates from these.

Table 1: Summary of scenarios and solution approaches

Scenario	Frequency of training	DNN Workload Changes?	DNN Training Time	Suitable Solution	Data Collection Overhead for Solution
Training once on a large dataset	One time	Never	Few days	Brute force; Profile all power modes and pick the best	1200–1800 min
Fine-tuning a Model	Occasional	Rare	Few hrs	Neural Network; Profile at least 100 power modes	20–50 min
Continuous learning	Periodic	Rare	< 1 hr	PowerTrain; profile 50 power modes and transfer	10–20 min
Federated learning on edge cloud	Often	Often	Unknown	PowerTrain; profile 50 power modes and transfer	10–20 min

the optimal for 5 out of 6 cases (except ResNet 15W) compared to Nvidia’s suggestions (NV), while falling within the power limits in most cases.

Lastly, in Figure 2b, we see that the PT-based optimization also outperforms simpler baselines like choosing the default MAXN power mode and doing random (RND) profiling of 50 power modes to build a partial Pareto and selecting the best, and our better baseline using a Neural Network (NN) trained on 50 power profile samples. For a set of optimization problems with power limits varying from 17–50W, PT has the lowest time penalty of 1% in excess of the optimal while also having the lowest percentage (26.5%) that exceeds the power limit by over 1W above the threshold.

These exemplars motivate the need for and illustrate the success of our robust PowerTrain prediction model to support power mode tuning for training DNN workloads on Jetson edge devices.

1.5. Discussion

It is useful to understand the scenarios under which these different approaches, including baselines and our proposed ones, are best suited for. In Table 1, we outline four potential application and deployment scenarios, their characteristics and the solution approach is most applicable. These estimates are based on our detailed experiments later reported in Section 4.

For instance, there may be a one-time training workload submitted to an edge service that runs new and large DNN model on a large corpus of user data. Since the training is time-consuming, running into days, having the perfect power mode is helpful even if the profiling time is costly. Here, a brute force approach may be recommended since the data collection time is approximately a day, which is well within the runtime of the workload. This kind of training workload is uncommon on the edge, and this approach is therefore impractical.

Fine-tuning a pre-trained DNN on a smaller dataset is a more common edge usecase (e.g., model personalization on a private edge device) where the DNN training time is typically a few hours. In this scenario, the DNN workload seldom changes and our NN approach using profiling over 100s of power models is feasible because of the longer training time.

Continuous learning is another common scenario on edge devices where the same DNN is regularly retrained with small batches of new incoming data to address data drift. Since the

training workload runs for a shorter time here, PowerTrain is the best approach since it has a smaller data collection time. Further, when the edge device is part of a federated learning setup or a private edge-cloud supporting multiple applications, the device could be assigned any DNN training workload, submitted often and with an unknown training duration. In such a case, it is best to use PT to minimize data collection time for the new workload and ensure that the optimization criteria are met before the workload changes.

1.6. Contributions

We make the following specific contributions in this article:

1. We develop a Neural Network (NN) based prediction approach to *estimate the time and the power* for DNN training workloads on the Nvidia Jetson Orin AGX for any given power mode (§ 3) using a large profiling corpus. This is extended into PowerTrain (PT), a *Transfer Learning* based approach that greatly reduces the online profiling overheads for any new DNN workload on Orin or another Jetson device.
2. We comprehensively validate NN and PT for accuracy and generalizability across 7 different DNN workloads and a different Jetson edge device. (§ 4).
3. Finally, we apply PT predictions to optimize the training time and power limit trade-offs for several DNN workloads (§ 5), and demonstrate that it out-performs other baselines, including NN.

This modeling work leverages a prior in-depth study that we conducted on the effect of various factors, including power modes, on the Jetson devices for DNN training [4]. There, the focus was on the characterization and as a potential use-case of that study, we had proposed a simple linear regression technique (Section 5.7 in [4]) for predicting the DNN training time and energy on the Jetson devices, validated on 10 power models for one Jetson device, with mixed results. The present article focuses entirely on predicting the time and power for DNN workloads on Jetsongs, proposing a wholly different, novel NN-based prediction technique along with the PowerTrain transfer learning approach for fast retraining and robust estimates. We further present optimization results as a case study of using our prediction models. All these contributions are novel, with no more than 20% conceptual overlap with the prior work.

Table 2: Specifications of NVIDIA Jetson Orin AGX and Xavier AGX devkits used in evaluations

Feature	Orin AGX	Xavier AGX
CPU Architecture	ARM A78AE	ARM Carmel
# CPU Cores	12	8
GPU Architecture	Ampere	Volta
# CUDA/Tensor Cores	2048/64	512/64
RAM (GB)/Type	32/LPDDR5	32/LPDDR4
Peak Power (W)	60	65
Form factor (mm)	110 × 110 × 72	105 × 105 × 65
Price (USD)	\$1999	\$999

Features that vary across Power Modes	Orin AGX	Xavier AGX
CPU core counts	1..12	1..8
# CPU freqs.	29	29
Max. CPU Freq. (MHz)	2200	2265
# GPU freqs.	13	14
Max. GPU Freq. (MHz)	1300	1377
# Mem freqs	4	9
Max. Mem. Freq. (MHz)	3200	2133
# Power modes	18,096	29,232

Table 3: DNN workloads and Datasets used in Experiments. All models are trained with batch size of 16. Estimated epoch training time (mins) for MAXN fastest power mode on Orin AGX is given as reference.

Task	Model	# Layers	# Params	FLOPs [†]	Dataset	# Samples	Size	Epoch Time (min)
Image classification	MobileNetv3 [25]	20	5.5M	225.4M	GLD23k [26]	23k	2.8GB	2.3
Image classification	ResNet-18 [27]	18	11.7M	1.8G	ImageNet Val. [28]	50k	6.7GB	3
Object detection	YOLO-v8n [29]	53	3.2M	8.7G	COCO minitrain [30]	25k	3.9GB	4.9
Question answering	BERT base [31]	12	110M	11.5T	SQuAD [32]	70k	40.2MB	68.6
Next word prediction	LSTM [33]	2	8.6M	3.9G	Wikitext [34]	36k	17.8MB	0.4

[†] As per the typical practice, FLOPs reported correspond to a forward pass with minibatch size 1.

part of the data being fetched from memory and the rest from disk. This gives consistent minibatch training performance.

2.5. Data Collection from Profiling Power Modes

Our prediction models use as input the profiling information collected during DNN workload training for a specific power mode of the device. When profiling each power mode, we train the candidate DNN workload for ≈ 40 minibatches, and record the training time per minibatch and the sampled power usage during the training period. There is negligible variation in the time or power measured across the minibatches during DNN training. However, we note that the very first minibatch’s training takes a long time to train, possibly due to PyTorch performing an internal profiling to select the best possible kernel implementation in this phase. Therefore, we discard the first minibatch profiling entry. Additionally, we notice that the power measurements when training using a new power mode take 2–3s to stabilize. So, we use a sliding window logic to detect when the power stabilizes, and use the profiling entries subsequent to that while ensuring 40 “clean” minibatch training samples are collected per power mode.

For the three default workloads, ResNet, MobileNet and YOLO, we collect and assemble a large corpus of ground truth power mode profiling data to help us with prediction model training and validation of their performance. The Orin power modes span 29 CPU frequencies, 13 GPU frequencies, 4 memory frequencies and 12 CPU cores for a total of 18,096 possible power modes⁷. From these, we choose 4,368 power modes for profiling that are uniformly distributed through the solution

space; these include all combinations of GPU (13) and memory (4) frequencies, even number of CPU cores (6), and every alternate CPU frequency that is available, excluding the two slowest ones (14)⁸. The data collection time depends on the power mode chosen, e.g., a lower configuration will take longer to profile ≈ 40 minibatches of training, and we report these times as part of our results.

This large corpus of profiling data for different power modes for diverse DNN training workloads helps us rigorously evaluate our prediction models. This dataset collected is available at Anonymized and in itself forms a rich community asset.

3. ML-driven Modeling and Prediction

Here, we propose two data-driven approaches based on ML-based modeling to predict the power and training time per minibatch for a DNN training workload.

Our proposed methods were arrived at after exploring and eliminating other simpler prediction methods such as linear regression, random forest and decision trees, which did not perform well: DNN workload performance appeared inherently non-linear and linear regression was inaccurate; and random forest and decision trees suffered from overfitting. We also explored using a Multi-Layer Perceptron (MLP) to predict time and power. While this performed well with very few samples, our eventual solution based on Neural Networks (PT) outperforms it with a slightly larger number of samples but with the key added benefit of being able to transfer the learning to a new workload.

⁷ Interestingly, the number and actual frequencies available change with the Board Support Package (BSP) version. We also noticed that the same Orin device, but with a different JetPack version, had a different set of frequencies available. Nvidia’s usually active forums do not have any comments related to this.

⁸ During our experiments, we found that the Jetson device only supports changing from higher to lower CPU and GPU frequencies. Other changes require reboots, and cause an error otherwise. We came up with an ordering of the power modes that satisfies this requirement.

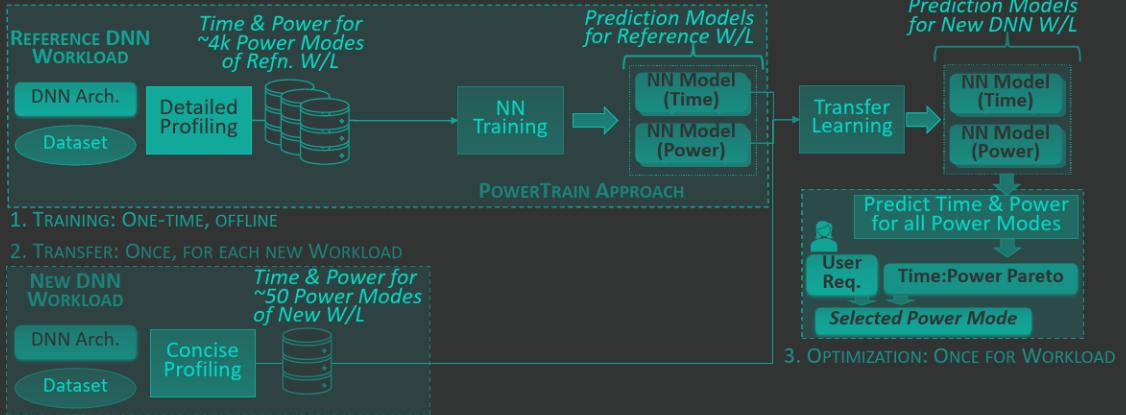


Figure 3: Time and Power Prediction Models: (1) Initial training, (2) Transfer, and (3) Use for Optimization

A key aspect that we consider here is the *overhead* for training the prediction models for a new DNN workload and, potentially, new hardware as well. This overhead is manifest primarily in the form of *profiling time* since the actual time to perform prediction model training itself is negligible, taking 15 minutes at most. NN models tend to be more accurate with more training samples, yet collecting profiling data for a growing number of power modes can take a linearly longer time. While the effort spent profiling can be reused as part of the actual training activity of the DNN workload (i.e., it is not wasted even from the user’s perspective), the benefits of constructing the prediction model and performing optimizations on it are subdued if this data collection time starts being a large part of the multi-epoch DNN workload training time.

3.1. Neural Network (NN) Prediction Model per Workload

Neural Networks (NN) [39] consist of multiple layers of interconnected neurons, which allow it to learn complex, non-linear relationships in the data. Neural networks usually contain one input layer, one or more hidden layers, and an output layer. We first develop simple NN architectures to predict the training time and the power for a given power mode on an edge accelerator, for a specific DNN workload.

Our *NN architecture* has 4 dense layers with 256, 128, 64 and 1 neurons, respectively. We use the ReLU activation function for the first 3 dense layers, and a linear activation function for the final layer. This was arrived at based on a grid search for hyper-parameters using *sklearn* library’s GridSearchCV. The input feature vector consists of the configuration for a power mode: CPU cores, CPU frequency, GPU frequency and memory frequency, while the output layer returns the predicted training time per minibatch or power. Each input feature is normalized to a value between 0.0 to 1.0 using the *sklearn* library’s StandardScaler to ensure better model generalization and faster convergence by providing consistent scales for all inputs. Adam is used as the optimizer, with a learning rate of 0.001 and the Mean Squared Error (MSE) serves as the loss function. We also introduce two dropout layers after the first and second dense layers to avoid overfitting. By default, we train the NN for 100 epochs and verify convergence. We use model checkpointing to

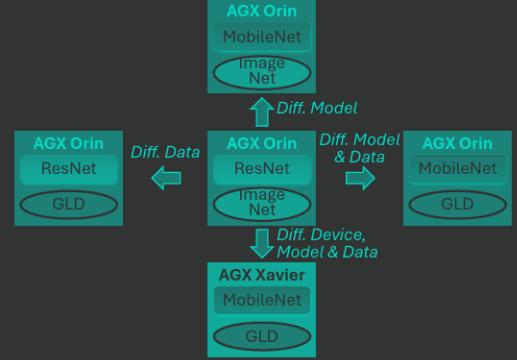


Figure 4: Generalizability of PowerTrain.

save the best weights (i.e., model with the least validation loss) seen during training and use it as the final model.

We take two alternatives for providing the training data to the NN: (1) *All* gives it the full profiling data for a DNN workload using a 90:10 split of training and test, which is about 4.4k samples for our Orin AGX for the default DNN workloads. (2) The second provides a much smaller number of training data, from 10 to 100 uniformly sampled from the 4.4k, again using a 90:10 split. Clearly, the first approach is time-consuming to collect the profiling data, e.g., taking over 14h for ResNet, but can give better prediction accuracy. In contrast, a single epoch of training for ResNet using the MAXN faster power mode takes 3 mins to train, with a typical training lasting for 120 epochs, or 6 h. The latter sampling approach takes a much smaller time, running into 10s of minutes, but can suffer from lower accuracy. We explore these trade-offs in the evaluation.

3.2. PowerTrain: Generalizable Transfer Learning Model

While the NN approach is simple, it needs to be retrained for every new DNN workload we wish to run on the edge device. This requires fresh data collection overheads for different power modes of the new DNN workload. To overcome this limitation, we propose *PowerTrain*, which leverages Transfer Learning (TL) from one NN to another to offer a lower overhead generalizable approach.

<i>Transferred To</i>	MobileNet		ResNet		YOLO	
<i>From Base</i>	Time Val %	Power Val %	Time Val %	Power Val %	Time Val %	Power Val %
MobileNet	8.12	3.62	15.03	7.98	11.77	4.98
ResNet	14.53	5.62	9.34	4.06	11.50	4.95
YOLO	17.03	9.71	19.76	12.88	9.72	4.81

Figure 5: Effectiveness of PowerTrain using different *reference DNN workloads* to transfer from. MAPE for time and power predictions validated against ground truth is reported.

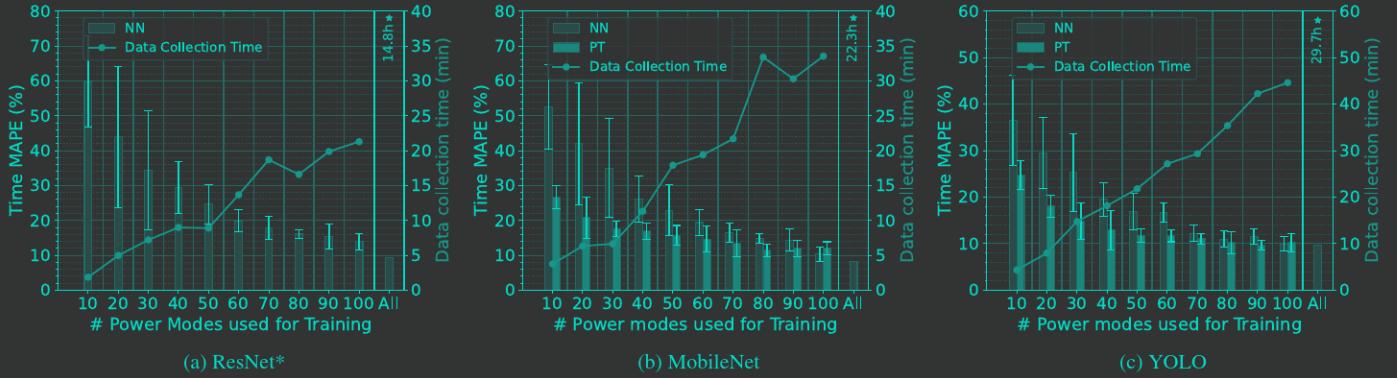


Figure 6: Prediction error and profiling overheads for *time prediction* for different DNN workloads

size of the training data used to train the model (NN), or update the transferred model (PT) from the base ResNet NN mode. *All* indicates NN training on all 4.4k samples. The median MAPE is shown as bars on the left Y axis with whiskers indicating $Q1-Q3$ error ranges, while the profiling time taken for these many training is shown on the right Y axis (line plot) as overheads. In both cases, smaller is better. Since ResNet is the reference model, we do not report PT results when ResNet is the target DNN workload.

PowerTrain performs better than NN on time predictions for a new DNN with fewer number of power modes profiled.

In Figure 6, the MAPE for PT with 10 power modes as training input for MobileNet is 26.7% while the error for NN using 10 power modes is 52.6%. With just 30 power modes, PT improves to a MAPE of < 20% for MobileNet while NN has a much higher 35% error. Similarly, for YOLO, at 30 power modes, PT's MAPE is < 15%, while at 10 samples it is 24.7%. Additionally, the $Q1-Q3$ whisker is much tighter for PT than NN, showing less variability on the accuracy even with different random subsets of power modes picked for training. This shows the potential of PowerTrain to generalize from a model trained fully on one reference DNN to another with few samples. The data collection time (right Y axis, line) for 30 power modes is under 15 mins for all 3 DNNs, and under 25 mins for 50 power modes. A typical training to convergence runs for 100s of epochs and many hours. Hence, the data collection overhead is a smaller fraction.

PowerTrain for power predictions achieves a higher accuracy for a new workload with fewer profiling samples than NN.

In Figure 7b, 20 power mode samples for MobileNet allows PowerTrain to have a MAPE of just 8.5%, as compared to 12%

for NN trained on a similar number of samples. Similarly, Figure 7c for YOLO shows PT achieve a power MAPE of 6.8% with 10 samples compared to 21% for NN. At larger number of samples of 50 or beyond, they have similar accuracies for YOLO though PT maintains a consistent 5% improvement over NN for MobileNet.

PT power prediction has lower MAPEs (5–10%) as compared to time predictions (10–20%) across all DNNs.

From Figures 7b and 6b, at 10 power mode samples for MobileNet, we observe a $\times 2$ lower MAPE for power as compared to time. This holds across DNN workloads, and for both PT and NN. At 50 power modes, which will be our default configuration for PT going forward, the power MAPEs for MobileNet and YOLO are 5.2% and 4.9% in contrast to time MAPEs of 15.7% and 11.7% for MobileNet and YOLO, respectively. This can be explained by the fact that we see little variation in power values for a given power mode, enabling the ML-based techniques to predict more accurately.

By 100 power modes, PowerTrain reaches close to the optimum accuracy for time and power prediction, and is comparable to NN training on all 4.4k power modes profiled.

In Figure 6b, at 100 power modes, PT has a MAPE of just 11.9% while NN trained on all samples has a MAPE of 8.1%. For power predictions on MobileNet, PT at 100 samples has a MAPE of 4.3% while the NN on all samples has a MAPE of 3.6%. For YOLO, the MAPE for time predictions is 10.18% for PT at 100 samples while the NN on all is a close 9.7%, and a similar trend is seen for power prediction as well. This again demonstrates the ability of PT to give higher accuracy with fewer samples.

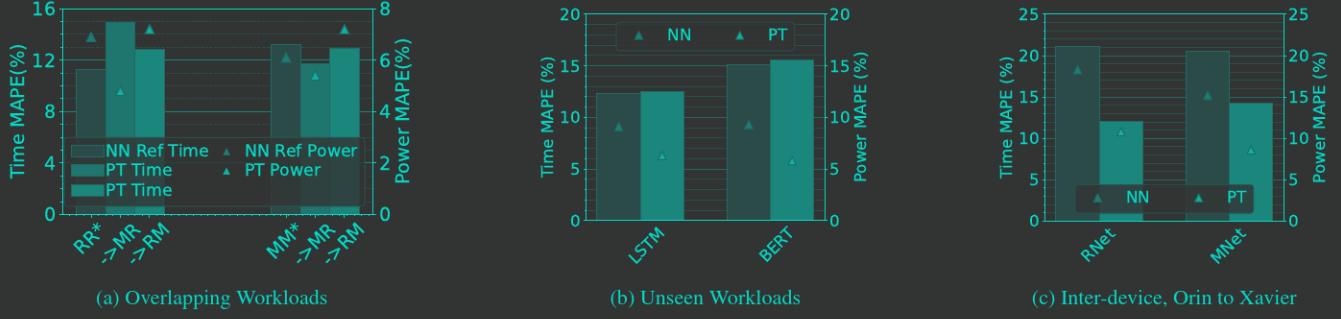


Figure 8: PowerTrain generalization when reference model is transferred to new workloads and device.

PowerTrain generalizes well to diverse DNN workloads with few samples, and outperforms NN on power predictions.

As seen from Figure 8b, for LSTM, the time MAPE is 12.5% for PT and 12.3% for NN, which are comparable. Similarly, for BERT, the time MAPEs of 15.6% and 15.1% are close. The benefits of PT is visible for power predictions, where PT has a MAPE of 6.3% while NN is 9.1% for LSTM, and likewise, PT has a 3.5% better MAPE for BERT.

4.3.3. Unseen Device from a Different Generation

In this experiment, we use PowerTrain to transfer a reference workload trained on Orin AGX to workloads that run on the Xavier AGX developer kit [36], which is the comparable previous generation top-end Jetson device. As shown in Table 2, compared to the Orin, the Xavier edge accelerator has an 8 core (vs. 12 for Orin) custom Carmel ARM CPU, a previous generation CUDA architecture (Volta rather than Ampere), and with $\frac{1}{4}$ th the number of CUDA cores (512 vs. 2048), and a previous generation RAM (LPDDR4 vs. LPDDR5). We randomly profile 1000 power modes out of the available 29,000 and collect power and time data for ResNet and MobileNet workloads in a similar manner as before. We use the ResNet workload trained on Orin AGX as our reference model for time and power predictions. Like before, PowerTrain then uses 50 random power modes for ResNet (or MobileNet) DNN workload profiled on Xavier to transfer-learn onto and validates the re-trained time and power prediction models using the remaining 950 power mode samples for that workload.

PT generalizes well to a device from a different generation and outperforms NN-based training.

Figure 8c, we see that when PT transfers from the reference DNN trained using ResNet workload on Orin to the same workload on Xavier (device changes), we see a time prediction MAPE of 12% and power prediction MAPE of 11%. This is much better than training an NN on just 50 power mode samples of the DNN workload on Xavier, which reports a MAPE of 21% for time and 18% on power. Similar benefits hold when both the device and the workloads are different for PT, where transferring the prediction models to MobileNet workload on Xavier has time MAPE of 14% and power MAPE of 9%, which are once again much better than NN.

This is a powerful result that suggests that there is sufficient similarity between the Jetson generations and DNN workload behaviors that can be learned during reference workload training, and adapted to a much different execution setting. This opens up opportunities to try even more diverse Jetsons like NX and Nano series, or even different edge models such as Raspberry Pi. This is left to future work.

5. Optimizing DNN Workloads using our Prediction Models

One of the applications of our time and power prediction models is to optimize the power mode configuration of DNN workloads. We formulate an optimization problem as follows. Given a DNN workload m_{tr} , choose a power mode pm from a set \mathbb{PM} that minimizes training time per epoch t_{tr} while ensuring that the power load P_{tr} falls within a user-specified budget P_b , i.e.,

$$\begin{aligned} & \text{Given } m_{tr}, \\ & \text{select } pm \in \mathbb{PM} \\ & \min t_{tr} \\ & \text{s. t. } P_{tr} \leq P_b \end{aligned}$$

As illustrated in the workflow in Figure 3, we use our prediction models to solve the optimization problem. For PowerTrain, we use the reference model trained using the ResNet workload on 4368 power modes, and transfer to other target workloads using 50 random power mode profiling samples. Since the time predictions are per minibatch, we scale them to per epoch times. As a ground truth optimization solution, we draw the observed Pareto front (*Obs Pareto* in Figure 9) using the profiling information for 4.4k power modes to minimize time and power. These set of points offer the least time value for a given power, and vice versa. For various power limits specified by the user, we can trivially search Pareto points to find the *optimal power mode* with a power value that is closest to but less than the power budget and report the training time for this.

Using the prediction models, we follow a similar approach, with the key difference that the models are used to estimate the training time and power for all possible (4.4k) power modes. We build a similar Pareto from PT predictions and this predicted Pareto (*PT Pred Pareto* in Figure 10) is used for the optimiza-

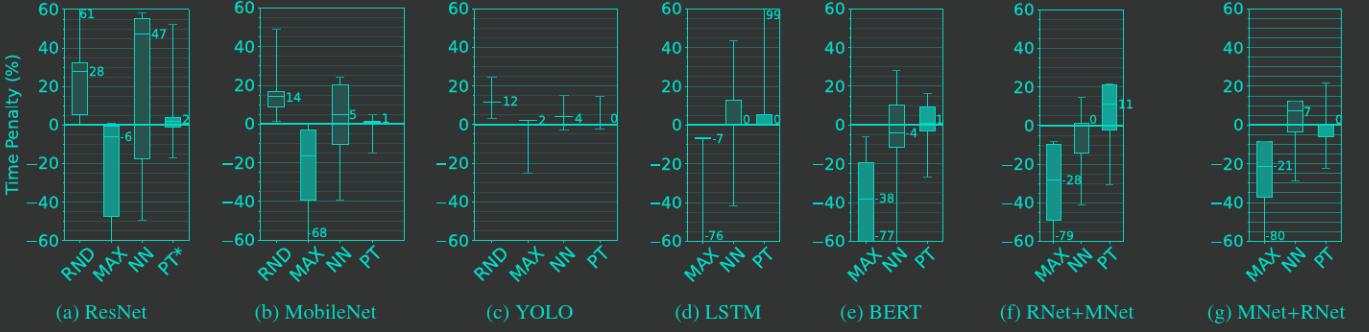


Figure 11: Time Penalty % for model-based optimization relative to ideal Pareto time. Lower is better. *PT for ResNet indicates training of base model on full data.

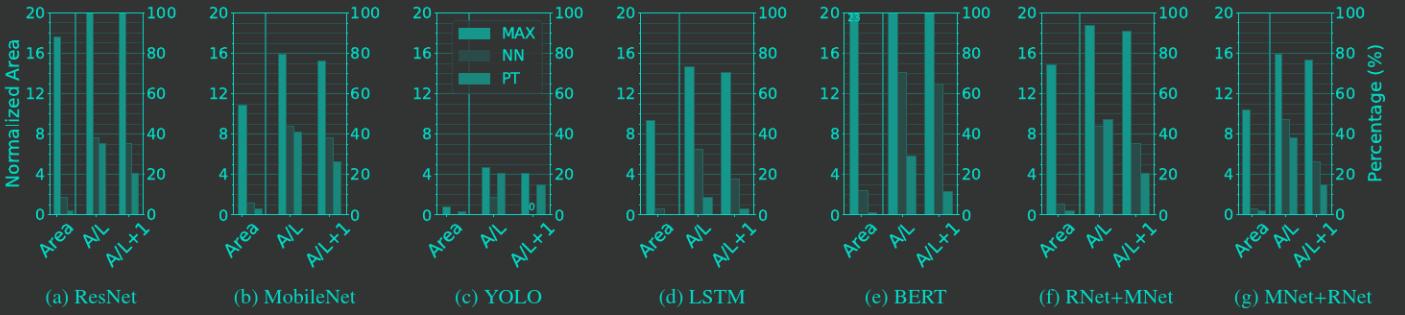


Figure 12: Pareto Power Errors for different DNN workloads

PT optimization results in a lower time penalty as compared to NN in most cases.

As seen from Figure 11, for MobileNet, the median time penalty for PT is just 0.7% over the optimal as compared to 5% for NN. Similarly, for YOLO, PT is as good as optimal with 0% penalty while NN has a median time penalty of 4%. For LSTM and BERT, PT has the same median time penalty as NN, however the distribution is tighter for PT. For MobileNet using ResNet training dataset (Figure 11f), PT has a lower median time penalty. PT does worse than NN only for ResNet with MobileNet data (Figure 11g).

PT optimization has the lowest normalized excess power AUC for most models and exceeds the power budget within a 1W threshold under 25% of the time.

In Figure 12, we report three metrics: Area, A/L and A/L + 1. Area stands for the *normalized area under the curve (AUC)* of power in excess of the given budget for the given solution by a strategy (W per solution); A/L stands for the *% of solutions that exceeds the power limit* for a strategy; and A/L + 1 stands for the *% of solutions that exceed the limit by more than a 1-Watt threshold*.

As seen from Figure 12, the *Area* is the lowest for PT across all DNN workloads except YOLO (6 out of 7 cases). Additionally, this budget + 1W buffer (A/L + 1) is breached in < 20% of the solutions for 6 out of 7 DNNs and is 25% for MobileNet. So PT's predictions and solutions help stay within the power budget most of the time.

MAXN almost always exceeds the power budget while offering the best time, while Random sampling has higher time penalties with no power violations

As seen from Figure 11, MAX has negative time penalties for ResNet, MobileNet and YOLO, i.e., is much faster than the optimal solution. But this is because of setting all frequencies to the maximum value, thus violating the power limit constraints often (Figure 12). In contrast, a Pareto from Random sampling is 12–28% slower than the optimal. So these simpler baselines are sub-optimal and potentially unusable.

6. Related Work

6.1. Energy, Power and Performance Modeling on Edge devices

Earlier works [3] use roofline modeling on a much older Jetson TK1 and TX1 to understand and characterize the performance of CPU and GPU micro-benchmarks for matrix multiplication. MIRAGE [42] uses gradient tree boosting to predict runtime and power consumption for DNN inference on a Jetson Xavier AGX. Others [21, 43] have investigated the impact of frequencies and cores with the latency, power and energy for inferencing on the Jetson Nano. Abdelhafez and Ripenau [22] examine the effect of frequencies on power consumption for stream processing workloads. Some [44] have developed methods to measure fine-grained layer-wise energy for inference workloads on the Jetson TX1. All of these have either looked at inferencing or micro-benchmarks, and not DNN training.

