

# MII: A Multifaceted Framework for Intermittence-aware Inference and Scheduling

Ziliang Zhang<sup>1</sup>, Cong Liu<sup>1</sup>, Hyoseung Kim<sup>1</sup>

<sup>1</sup>University of California, Riverside

{zzhan357, cong1, hyoseung}@ucr.edu

**Abstract**—The concurrent execution of deep neural networks (DNN) inference tasks on intermittently-powered batteryless devices (IPDs) has recently garnered much attention due to its potential in a broad range of smart sensing applications. While the checkpointing mechanisms (CMs) provided by the state-of-the-art make this possible, scheduling inference tasks on IPDs is still a complex problem due to significant performance variations across DNN layers and CM choices. This complexity is further accentuated by dynamic environmental conditions and inherent resource constraints of IPDs. To tackle these challenges, we present MII, a framework designed for intermittence-aware inference and scheduling on IPDs. MII formulates the shutdown and live time functions of an IPD from profiling data, which our offline intermittence-aware search scheme uses to find optimal layer-wise CMs for each task. At runtime, MII enhances job success rates by dynamically making scheduling decisions to mitigate workload losses from power interruptions and adjusting these CMs in response to actual energy patterns. Our evaluation demonstrates the superiority of MII over the state-of-the-art. In controlled environments, MII achieves an average increase of 21% and 39% in successful jobs under stable and dynamic energy patterns. In real-world settings, MII achieves 33% and 24% more successful jobs indoors and outdoors.

## I. INTRODUCTION

Intermittently-powered batteryless devices (IPDs) [1, 2] offer a promising pathway to zero carbon emissions and maintenance-free operations. Recent advances have enabled them to execute deep neural network (DNN) inference tasks [3–6], essential for smart sensing and IoT applications. These devices harvest ambient energy from the environment and store it in capacitors. Once sufficient energy accumulates, IPD executes tasks using this energy until depletion. IPDs are typically equipped with two types of memory: volatile memory (VM), which is fast but loses data upon shutdown, and non-volatile memory (NVM), which is slow but retains data after shutdown [7, 8]. Since an IPD turns on and off across power cycles, it must store intermediate computation results from VM to NVM before powering off [9–14].

Existing research on IPDs primarily centers around checkpointing mechanisms (CMs) that preserve execution progress across power failures. Broadly, these mechanisms fall into two types: just-in-time checkpointing (JIT), and static checkpointing (ST) using atomic blocks. JIT [10, 15–17] checkpoints the system state once at the end of each power cycle, achieving faster speeds but demanding a larger peak memory. On the other hand, ST [3, 5, 12, 14, 17–19] transforms the program code into smaller atomic blocks of various granularity (e.g., layers, filter, and tiles for DNNs), with checkpointing code at

the end of each block, offering a smaller peak memory but at the cost of speed (Sec. II-B).

Although existing studies have laid the groundwork for executing DNN inference tasks on IPDs, significant challenges persist for real-world deployment. First, layer-wise structural distinctions of DNNs lead to performance heterogeneity across layers, demanding an optimal CM for each layer (Sec. III-A). However, the limited VM size and intermittent power of IPDs make this particularly challenging due to the inevitable device shutdowns experienced by some layers (shutdown layers).

Second, the real-world environments present varying energy patterns, resulting in different shutdown layers for inference tasks at runtime. Consequently, CMs choices considered to be optimal for one environment may become the worst in another (Sec. III-B), necessitating a runtime adaption of CMs.

**Contributions.** To address the challenges above, we present MII: Multifaceted framework for Intermittence-aware Inference and scheduling. MII consists of two parts: offline and online. The offline phase addresses the first challenge which requires co-consideration of both shutdown layers and peak memory usage. Our offline intermittence-aware search method identifies the optimal CM for each layer under a given environment so that each task’s execution time is minimized and the memory constraint is met. The online phase addresses the second challenge, which requires a low-overhead algorithm that quickly captures the environment dynamics and makes adaptations accordingly. MII’s online phase makes scheduling decisions dynamically, aligns task execution with the power cycles, and adapts CMs according to the actual energy supply and usage patterns. MII also introduces a proactive shutdown feature to mitigate the wasted work problem in a mixed JIT and ST system. Compared to existing work, MII achieves an optimal execution of each inference task and adapts it to the runtime environment with its unique layer-wise CM design.

We implemented MII on an Apollo4 Blue Plus [1] Board and tested it against 8 DNNs trained from 6 datasets. We evaluate MII in both controlled and real-world environments and compare it with three state-of-the-art methods [4, 5, 15]. MII achieves an average increase of 21% and 39% in successful jobs than the other methods under stable energy patterns and dynamic energy patterns from the controlled environment. MII achieves 33% and 24% more successful jobs under indoor and outdoor real-world environments. These results demonstrate that MII achieves: (i) *Efficiency* in intermittent DNN inference execution via offline layer-wise CM selection and online algorithms; (ii) *Adaptability* to diverse environments, ranging

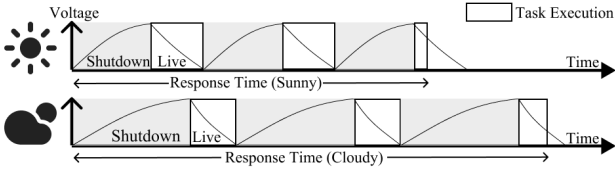


Fig. 1: Same task execution over shutdown and live times. Environment variation results in different response times.

from stable to dynamic energy patterns; and (iii) *Applicability* to real-world scenarios with significantly better performance over the state-of-the-art.

## II. BACKGROUND

### A. IPD and Intermittent Inference

An IPD harvests energy from ambient sources, e.g., solar, wind, radio waves, and vibration. Once sufficient energy is accumulated, the IPD turns on and begins program execution during *live time*. Although energy harvesting can continue during this phase, the device typically consumes energy at a faster rate than it accumulates. When energy is depleted, the IPD powers down, waiting for enough energy to be harvested (*shutdown time*) before restarting this process [9–14, 20]. Fig. 1 depicts an example of task execution over live and shutdown times. Several prior studies [4, 8, 17, 21–23] have enabled multiprogramming and priority-based scheduling of tasks on IPDs, with a timekeeping ability across power cycles using either the MCU’s deep sleep mode or external real-time clock (RTC). In this context, our focus is on inference tasks that should run periodically for practical use in areas such as smart sensing, which periodically samples readings and runs inference tasks for anomaly or object detection [3, 17, 24].

An intermittent inference task refers to a task executing the forward propagation of a DNN under intermittent power. Because of on-off power cycles, modern IPDs are equipped with both volatile memory (VM) and non-volatile memory (NVM). VM is typically SRAM which is fast but small (tens to hundreds of KB in most MCUs) and loses all data when powered off. NVM, such as MRAM [1] and FRAM [2], is larger and slower compared to VM, but it can maintain data when powered off. To fully utilize the speed of VM in DNN inference, existing work [5, 25] loads all necessary data to VM, including the input features map (IFM), weights (WEI), and output features map (OFM), and then performs computations using the VM data. Before shutdown, the calculated OFM from this power cycle is checkpointed to NVM, and once the device reboots, IPD resumes the remaining inference computations by fetching the checkpointed data to VM.

In contrast to other computational tasks, DNN inferences keep a large memory footprint during execution and have a magnitude more data than needed for checkpointing [3, 5, 25]. For example, a tiny 7-layer DNN performing a 32x32 pixel colored image classification needs to checkpoint 9216 output features to NVM for the largest layer, whereas non-inference tasks, such as thermometer sensing and alarm, only need to checkpoint less than 10 outputs [14, 17, 19, 21, 22]. Despite the large memory footprint, loading all corresponding data (including WEI) to VM during inference is necessary, as it

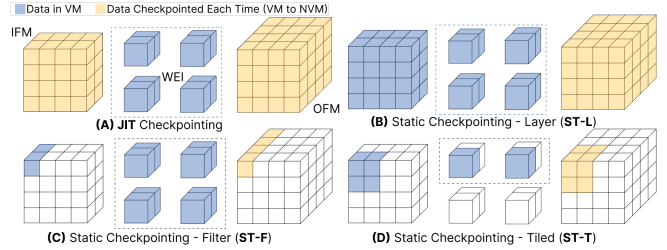


Fig. 2: 4 types of CMs applicable to DNNs: blue is the data read back to VM; yellow is the data checkpointed each time. significantly reduces NVM accesses and results in up to 51% less response time and 39% longer live time for the same 7-layer DNN compared to direct read and write in NVM [3, 5].<sup>1</sup>

### B. Checkpointing Mechanisms

State-of-the-art checkpointing mechanisms (CMs) fall into two categories: JIT checkpointing (JIT) and static checkpointing (ST). JIT [10, 15, 16] makes a checkpoint of the entire system’s states to NVM when shutdown is imminent. The device’s energy level, i.e., capacitor voltage, is constantly polled and compared with a predefined voltage threshold (JIT threshold) that guarantees a successful checkpointing [16, 17]. When the capacitor voltage falls below the JIT threshold, JIT checkpoints the system states to NVM so that the IPD can be safely shut down without losing its progress [15]. Although JIT enjoys fast execution speed by checkpointing only once per power cycle, it demands a substantial amount of memory. JIT needs to checkpoint both the IFM and OFM of the current layer since the previous checkpoint may not have saved the previous layer’s OFM (the current layer’s IFM). A detailed memory access pattern of JIT is shown in Fig. 2(A).

ST entails transforming the original task into atomic blocks and performing a checkpointing at the end of each atomic block [3, 5, 14, 17–19]. If a shutdown occurs in the middle of a block, the IPD resumes from the last checkpoint upon reboot and re-executes the block. Since any code with Write-After-Read (WAR) can disrupt idempotency, methods have been studied to construct atomic blocks to guarantee memory consistency and correct execution [12, 18, 26]. In the context of DNNs, an inference task can be divided into atomic blocks of various granularities shown in Fig. 2(B)(C)(D):<sup>2</sup>

- ST-L (layer): Due to its explicit IFM and OFM structures, each layer of a DNN can be naturally modeled as an atomic block, achieving ST at the layer-level granularity. ST-L is generally the fastest among all STs, but it needs to load all IFM, WEI, and OFM of that layer, resulting in the largest peak memory size among all STs.
- ST-F (filter): In convolution layers, a filter is convolved across the IFM to compute a feature vector output. By re-writing each filter convolution into a separate atomic block, ST is attained at the filter granularity [3, 4]. ST-F is

<sup>1</sup>This holds for IPDs that run CPU and SRAM at higher clock rates than MRAM or FRAM, such as our platform [1]. For others like MSP430 [2], loading WEI may be considered optional.

<sup>2</sup>This approach of leveraging the DNN structure is motivated by iNAS [5], which offers benefits over general programming language-based [27] and compiler-assisted [28] methods by ensuring that blocks fit into device memory without requiring extensive manual effort and code changes.

generally slower than ST-L due to its finer-grain block size; however, ST-F requires less memory than ST-L by loading partial IFM and OFM. Note that ST-F still needs to load the entire WEI of the layer for its filter-wise computation.

- ST-T (tile): Inference can be further broken down by re-organizing into a tiled structure [5]. ST-T can be achieved by converting each tile’s execution into an atomic block. Hence, unlike ST-F, ST-T can do computation with only a portion of WEI, thereby further reducing peak memory with a potentially longer execution time.

### C. Environmental Effects

In real-world scenarios, the dynamics of the environment lead to significant changes in the energy harvesting rate, resulting in different response times for the same task on an IPD. Fig. 1 illustrates this phenomenon under solar energy. Compared to the case of sunny light conditions, the shutdown time under cloudy conditions is obviously longer due to the lower harvesting rate. The live time is shorter because the device still harvests energy while it is executing the inference task but the harvesting rate is lower.

The variation of the environment is therefore the key challenge in scheduling tasks on an IPD. Existing work addresses this in two categories: energy prediction and workload reduction. Energy prediction methods [21–23] assume a priori knowledge of future energy patterns or predict based on previous patterns. However, the prediction can never be perfect due to the sporadic nature of the environment, and the use of more complex models increases overhead. Conversely, workload reduction [4, 19, 29–31] reacts to environmental changes by reducing the workload (e.g., skipping some layers of DNNs, called ‘early termination’ [4]) as the harvesting rate reduces. Its limitations include the degradation in output quality, and the extra efforts and overhead to enable early termination. More importantly, in DNN inference, the entire layer OFM has to be loaded in VM for an early-exit classifier or model to begin execution [4, 32]. This makes it unable to keep the peak memory usage smaller than the layer OFM, potentially limiting the IPD from running multiple inference tasks.

### D. System Model

We consider an IPD equipped with a fixed-size capacitor and a solar panel for energy harvesting, and with an RTC for timekeeping. The system has  $m$  periodic DNN inference tasks. Each task  $\tau_i$  releases a job according to its period, and each job performs one inference of the task’s DNN with  $\eta_i$  layers. Due to the nature of intermittent computing, we do not aim to execute all jobs of tasks; instead, we focus on maximizing the number of jobs that successfully complete their execution. If a job cannot finish before the start of the next period, it continues executing in the next period and the next period’s job is skipped to prevent overloads. The system may have other non-inference tasks, e.g., sensor and peripheral tasks, but they are not the main focus of this paper and their CMs are statically determined as done in prior work [17, 22–24, 30].

The system has  $V_{ipd}$  KB of memory in VM for inference tasks. In practice, the DNN models for MCUs dynamically

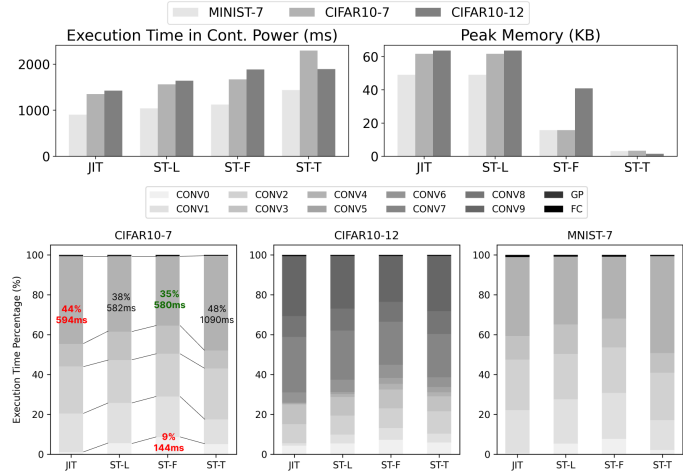


Fig. 3: Top: JIT vs. STs of various granularity tradeoff space. Bottom: layer performance heterogeneity under JIT and STs.

allocate and free memory from the heap space for each layer’s execution. Thus,  $V_{ipd}$  essentially indicates the heap area size, and for each task, the layer that uses the most heap memory determines the peak memory usage of that task.

## III. MOTIVATION

### A. Checkpointing Tradeoff on Intermittent Inference

To understand the effect of CMs on the execution time and memory usage of a DNN inference job, we set up an experiment evaluating both JIT and all granularities of ST (ST-L, ST-F and ST-T) with 3 DNNs on MNIST [33] and CIFAR10 [34] datasets. Each DNN name is given by ‘dataset name - # of layers’. We chose the DNNs and datasets following the prior work [3–5] as they are used in real-world applications like wildlife monitoring. We tested the DNNs on an Apollo4 Blue Plus [1] evaluation board due to its sufficient VM size to run these DNNs under all four CMs.

We first applied each CM to the entire DNN, as done in prior work [3–5] under continuous power from USB. As shown in the top of Fig. 3, JIT yields shorter execution times and consumes larger memory than ST, whereas ST uses a small memory size at the cost of longer execution time. This tradeoff occurs due to their inherent differences in checkpointing. JIT loads the entire layer’s IFM, OFM, and weights to the VM during execution. This results in the peak memory size matching that of the largest layer within the system. On the contrary, ST only fetches a portion of layer data as described in Fig. 2 so that it can maintain a small memory footprint. However, it needs to checkpoint the calculated results to NVM after each block. As the granularity of ST decreases from layer level to tile level, we observe an increase in inference latency and a decrease in peak memory usage.

We further break down the overall inference execution time into individual layers and characterize the layer-wise performance. Fig. 3 bottom shows the layer-wise relative execution time of DNN inference under different CMs for all 3 DNNs. As depicted in the figure, each layer experiences a different execution time depending on the CM used, which is not consistent with the general expectation that JIT is faster than ST. For instance, when executing CONV4, which has a

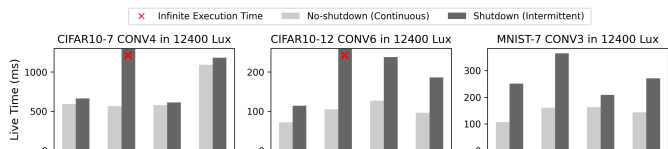


Fig. 4: Total live time of layers in the presence of shutdown.

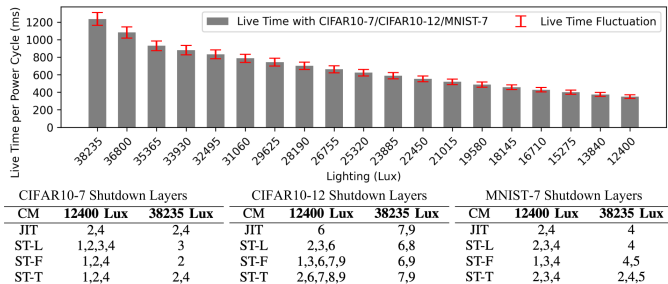


Fig. 5: Top: live time per power cycle when running CIFAR10-7 in various light conditions (CIFAR10-12 and MNIST-7 have the same pattern). Bottom: the shutdown layers of all three DNNs in two light conditions.

large IFM and small OFM, JIT takes longer execution time and yields worse performance compared to ST-F and ST-L (left-most red text) because JIT has to checkpoint the entire IFM of CONV4. ST-F achieves the best performance (top green text) due to the small OFM of CONV4, which reduces the checkpointing overhead of ST-F. However, if we choose ST-F as the CM across all layers of the job, it gives the worst performance for layer CONV0, which has a large OFM and small WEI (bottom red text). This is because the large OFM of CONV0 results in ST-F having a greater checkpointing overhead than the other CMs.

**Obs. 1.** For DNN inference tasks, relying on a single CM, as done in prior work, may result in sub-optimal performance, thereby requiring an optimal CM choice for each layer.

We therefore advocate a layer-wise adoption of different CMs. Although it may seem straightforward to choose the optimal CM for each layer, solving such a problem under an intermittent power condition is challenging. Fig. 4 compares the cumulative live time of each layer, i.e., the time that the device is live for inference and checkpointing, under intermittent power with a 2mF capacitor and 12400 Lux lighting condition. As depicted, the optimal CM choice under continuous power (‘No-shutdown’ in the legend) becomes often sub-optimal under intermittent power which causes the device to shut down at least once (‘Shutdown’). Also, the optimal CM chosen under continuous power can make a layer trapped in the endless loop of re-execution. For instance, ST-L is the optimal CM for the CONV4 layer of CIFAR10-7 under continuous power. However, ST-L makes it unable to checkpoint before shutdown under intermittent power, resulting in infinite execution time.

**Obs. 2.** Due to shutdown, the optimal CM choice for each layer of an interference task under intermittent power may diverge from the one made under continuous power.

### B. Intermittent Inference under Environment Variations

To find out how the change in energy harvesting patterns affects inference performance, we first focus on the live time

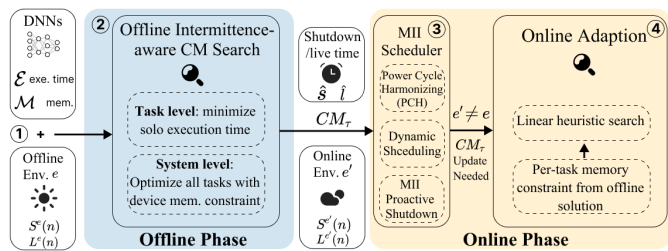


Fig. 6: Overview of the proposed MII framework

of an IPD, which directly affects the response time of an inference task, as discussed in Sec. II-C. We use the same three DNNs and run them each separately in various light conditions. The top part of Fig. 5 shows the average live time per power cycle when running the CIFAR10-7 model under different CMs and light conditions. From the results, we find that live time varies significantly with light conditions, and that under the same light condition, live time is only marginally affected by CMs and tasks. For example, under 38325 Lux lighting, the live time in each power cycle for all three DNNs is 1236ms. On the other hand, if we change the lighting to 29625 Lux, the live time per power cycle changes to 743ms. Some fluctuations may be observed depending on CMs or tasks, but they are within the 6-7% range of the live time.

To further explore the effect of CMs and environment conditions on intermittent inference, we characterize *shutdown layers*, which are the subset of layers of an inference job that experience shutdowns during their execution. The tables at the bottom of Fig. 5 depict the shutdown layers of each DNN job in two representative light conditions. The shutdown layers of each DNN vary significantly with the CM used and the given light condition. Recall our discussion in Sec. III-A and with Fig. 4. The optimal CM when the layer does not shut down becomes often sub-optimal or even the worst if a shutdown occurs for that layer. Vice versa, if we choose the optimal CM assuming that the layer always experiences a shutdown, such a CM will likely perform worse when there is no shutdown.

From the above two experiments on live time and shutdown layers, the following observation can be made.

**Obs. 3.** While both live time and shutdown layers play significant roles in determining optimal CMs, they can vary drastically with environmental conditions. Consequently, there is a strong need for runtime adaptation with low overhead.

## IV. MII DESIGN OVERVIEW

Fig. 6 presents the overview of our MII framework, designed to address the two key challenges elaborated with our observations in Sec. III. We illustrate each challenge in a separate paragraph and propose MII’s solutions.

Motivated by Obs. 1 and 2, MII introduces an *Offline Phase* to tackle the challenge of finding a layer-wise optimal CM solution under a given environment, while considering possible device shutdowns and device memory constraints. We opt for an offline solution because the profiling of DNN execution time and memory footprint on a per-layer basis is feasible only in an offline setting. In ①, the offline phase models the energy supply pattern of a given environment condition  $e$  into the cumulative shutdown time  $S^e(n)$  and cumulative live time  $L^e(n)$  functions where  $n$  is the number of power cycles. In ②,

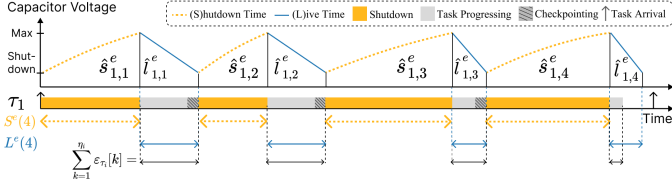


Fig. 7: Cumulative shutdown time  $S^e(4)$  and live time  $L^e(4)$  of an intermittent inference task  $\tau_1$ .

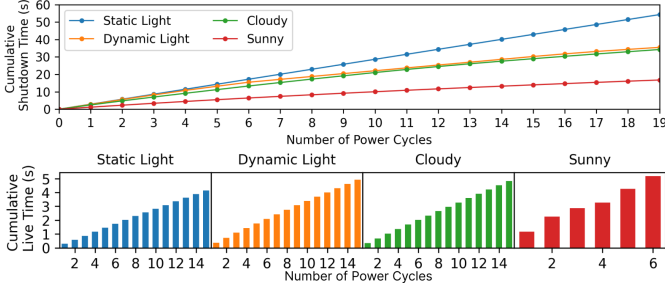


Fig. 8: Top: cumulative maximum shutdown time  $S^e(n)$  under different environment condition  $e$ . Bottom: cumulative minimum live time  $L^e(n)$  for CIFAR10-12 inference task  $\tau_{C12}$ .

the offline intermittence-aware CM search finds the optimal CMs for the layers of each task to minimize their execution time while adhering to the device memory constraint.

Stemming from Obs. 3, actual energy supply patterns can vary drastically at runtime, which requires a timely adaption of CMs to cope with changing environmental conditions. Therefore, MII presents an *Online Phase* to tackle this challenge. The MII scheduler in ③ takes into account the CMs found at the offline phase and updates the offline  $S^e(n)$  and  $L^e(n)$  to their online versions,  $S^{e'}(n)$  and  $L^{e'}(n)$ , based on the online shutdown and live times,  $\hat{s}$  and  $\hat{l}$ , collected using an on-board RTC. The scheduler harmonizes task execution with power cycles, makes scheduling decisions based on the online energy pattern, and employs proactive shutdown to mitigate the wasted work problem in a mixed JIT and ST system. Lastly, we introduce the online adaption method in ④, which adapts the CMs according to the latest  $S^{e'}(n)$  and  $L^{e'}(n)$  using a linear heuristic search for each scheduled inference task.

## V. OFFLINE PHASE

### A. Modeling Shutdown and Live Time Patterns

Before conducting the offline CM search, we model the energy supply pattern of a given environmental condition  $e$  based on profiling and formulate it into shutdown and live time functions,  $S^e(n)$  and  $L^e(n)$ .

Recall that the execution of a job of an inference task  $\tau_i$  can take multiple power cycles, as illustrated in Fig. 7. We denote the shutdown time of  $\tau_i$  during the  $j$ -th power cycle in an environment condition  $e$  as  $\hat{s}_{i,j}^e$ , and the live time as  $\hat{l}_{i,j}^e$ . During each  $\hat{s}_{i,j}^e$ , the IPD remains off, only harvesting energy. Conversely, during each  $\hat{l}_{i,j}^e$ , the IPD powers on and starts consuming the capacitor's energy to execute  $\tau_i$  while still harvesting energy. Given the IPD's fixed capacitor size, we make the assumption A1 about the shutdown time as below:

**A1.** The duration of shutdown in the  $j$ -th power cycle is solely affected by the environment condition  $e$ , not by any task

on the IPD. Hence, for ease of presentation, we use  $\hat{s}_j^e$  to denote the shutdown time for the  $j$ -th power cycle.

This is a valid assumption because at the beginning of the  $j$ -th power cycle, the voltage level of the IPD's capacitor is at the power-off voltage regardless of the type of tasks executed in the previous power cycle, and the IPD turns on only when it reaches the power-on voltage, the timing of which is affected by the energy harvesting rate. Hence, with A1 and the profiled  $\hat{s}_j^e$  data, we can represent the shutdown time as a function of the number of power cycles.

**Def. 1** (Cumulative Shutdown Time).  $S^e(n)$  gives the cumulative maximum shutdown time over  $n$  consecutive power cycles in an environment condition  $e$ . Given  $\mathcal{N} \gg n$  shutdown time profiles,  $S^e(n)$  can be obtained by:

$$S^e(n) = \max_{1 \leq k \leq \mathcal{N}-n+1} \sum_{j=k}^{n+k-1} \hat{s}_j^e \quad (1)$$

Fig. 8 top plot gives an example of  $S^e(n)$  in various real-world environment conditions (x-axis indicates  $n$ ). Static and dynamic light are collected under a controllable artificial light source, whereas sunny and cloudy are collected under natural sunlight under two different weather conditions. For each condition  $e$  in this figure,  $S^e(n)$  determines a conservative estimate of the total time required to charge the IPD for execution across  $n$  power cycles. Note that  $S^e(n)$  is non-linear, e.g.,  $S^e(n+1) \leq S^e(n) + S^e(1)$ .

Unlike the shutdown time, the live time of an inference task,  $\hat{l}_{i,j}^e$ , depends not only on the energy harvesting rate of the environment  $e$  but also on the energy consumption rate of the system. As shown in Sec. III-B with Fig. 5, the energy harvesting rate is the dominant factor in the live time per power cycle, while the variation due to the type of tasks or CMs is relatively small (less than 7% of the live time per power cycle). We therefore make the following assumption A2.

**A2.** While turned on, the energy consumption rate of the IPD is the same for all tasks.

Strictly speaking, this assumption is not necessarily true because tasks may have different memory and I/O access patterns. However, since our work focuses on DNN inference tasks that do not involve direct I/O access and the energy harvesting rate has a much higher impact on the device's live time per power cycle, we find A2 works well in practice. With A2, we use  $\hat{l}_j^e$  to denote the live time for the  $j$ -th power cycle and derive a live time function  $L^e(n)$ , similar to  $S^e(n)$ .

**Def. 2** (Cumulative Live Time).  $L^e(n)$  gives the cumulative minimum live time over  $n$  consecutive power cycles in an environment condition  $e$ . Given  $\mathcal{N} \gg n$  live time profiles,  $L^e(n)$  can be obtained by:

$$L^e(n) = \min_{1 \leq k \leq \mathcal{N}-n+1} \sum_{j=k}^{n+k-1} \hat{l}_j^e \quad (2)$$

Unlike  $S^e(n)$ ,  $L^e(n)$  captures the *minimum* cumulative time. This allows us to have a conservative estimate of the time available for task execution over  $n$  power cycles. We can use  $L^e(n)$  to find out how many power cycles are needed to execute a job of a task, assuming no other tasks are executing in the system. For instance, if a job has the execution time of

TABLE I: Execution time and memory usage profiles

Symbol	Meaning
$a_i(x, k)$	$k$ -th layer execution time under CM $x$ without shutdown
$d_i(x, k)$	$k$ -th layer execution time under CM $x$ with shutdown
$v_i(x, k)$	$k$ -th layer peak memory under CM $x$

$t$  units, finding  $n$  that satisfies  $L^e(n-1) < t \leq L^e(n)$  tells us the number of power cycles involved. Hence, we can derive an inverse function,  $\bar{L}^e(t)$ , which gives the number of power cycles for  $t$  units of execution.

Fig. 8 bottom illustrates  $L^e(n)$  during one job execution of the CIFAR10-12 inference task under the same four real-world environments. Obviously, it takes the least number of power cycles in the sunny condition. Both  $S^e(n)$  and  $L^e(n)$  are stored in the device’s NVM so that the scheduler can access them for online adaptation.

The execution time and memory usage of each layer of a task  $\tau_i$  are affected by the CM choice and whether the layer experiences shutdown during execution (Sec. III-B). Therefore, for each layer  $k$  of  $\tau_i$ , we record two lists of execution times with JIT, ST-L, ST-F, and ST-T:  $(a_{i,k}^{\text{JIT}}, a_{i,k}^{\text{ST-L}}, a_{i,k}^{\text{ST-F}}, a_{i,k}^{\text{ST-T}})$  represent times without shutdown (‘alive’), and  $(d_{i,k}^{\text{JIT}}, d_{i,k}^{\text{ST-L}}, d_{i,k}^{\text{ST-F}}, d_{i,k}^{\text{ST-T}})$  denote times with shutdown (‘dead’). We also record the maximum memory usage of the  $k$ -th layer of  $\tau_i$  under four CMs:  $(v_{i,k}^{\text{JIT}}, v_{i,k}^{\text{ST-L}}, v_{i,k}^{\text{ST-F}}, v_{i,k}^{\text{ST-T}})$ . For ease of reference, we introduce functions  $a_i(x, k)$  and  $d_i(x, k)$  to obtain the execution times of  $\tau_i$ ’s layer  $k$  under a given CM  $x$  when the layer is ‘alive’ or ‘dead’ respectively. We also introduce a function  $v_i(x, k)$  for the memory usage. Table I summarizes these profiled data.

### B. Offline Intermittence-aware CM Search

Our goal is to determine layer-wise CMs to minimize the solo execution time of a task  $\tau_i$  across power cycles, i.e., when  $\tau_i$  runs without temporal interference from other tasks, while ensuring that the collective peak memory usage of all tasks stays within the IPD’s memory constraint,  $V_{ipd}$ . We solve this problem through a two-level dynamic programming approach.

At first, we minimize each task  $\tau_i$ ’s execution time under a memory constraint  $V$ . Let us define  $\varepsilon_{\tau_i}[k][V]$  as follows:

$$\varepsilon_{\tau_i}[k][V] = \tau_i\text{'s collective execution time from layers 1 to } k, \text{ while not exceeding the memory constraint } V. \quad (3)$$

$$cm_{\tau_i}[k][V] = \text{CMs achieving } \varepsilon_{\tau_i}[k][V].$$

As each layer uses non-zero memory,  $\varepsilon_{\tau_i}[k][0] = \infty$  and  $cm_{\tau_i}[k][0] = \emptyset$  for all  $k \leq \eta_i$ . The execution time of layers 1 to  $k$  can be found by considering the occurrence of shutdowns. The peak memory usage of a task  $\tau_i$  is determined by the layer that uses the most memory among all layers. Hence, we can compute  $\varepsilon_{\tau_i}[k][V]$  and  $cm_{\tau_i}[k][V]$  using Alg. 1.

Alg. 1 iterates over the memory size  $V$  from 1 to  $V_{ipd}$ , and for each  $V$ , iterates over layers from 1 to  $\eta_i$ . For each layer  $k$ , it considers four CMs (line 5). Recall that the peak memory usage of  $\tau_i$  is determined by the layer with the maximum usage, not by the summation of all layers. Hence, if the use of a CM  $x$  for the  $k$ -th layer violates the memory constraint  $V$ , that CM  $x$  should be ignored (line 6). If the  $k$ -th layer has

Algorithm 1: Minimize task execution time

```

1  $\varepsilon_{\tau_i}[0][1\dots V_{ipd}] \leftarrow 0; cm_{\tau_i}[0][1\dots V_{ipd}] \leftarrow \emptyset;$ 
2 for  $V \in \{1\dots V_{ipd}\}$  do
3   for  $k \in \{1\dots \eta_i\}$  do
4      $min_\varepsilon \leftarrow \infty; min_{cm} \leftarrow \emptyset;$ 
5     for  $x \in \{\text{JIT, ST-L, ST-F, ST-T}\}$  do
6       if  $v_i(x, k) > V$  then
7         continue; /* Ignore CM violating mem limit */;
8       end
9        $val \leftarrow \varepsilon_{\tau_i}[k-1][V] + a_i(x, k)$  /* No shutdown */;
10      if  $\bar{L}^e(val) \neq \bar{L}^e(\varepsilon_{\tau_i}[k-1][V])$  then
11         $val \leftarrow \varepsilon_{\tau_i}[k-1][V] + d_i(x, k);$  /* Shutdown */;
12      end
13      if  $min_\varepsilon > val$  then
14         $min_\varepsilon \leftarrow val; min_{cm} \leftarrow x;$ 
15      end
16    end
17     $\varepsilon_{\tau_i}[k][V] \leftarrow min_\varepsilon;$ 
18     $cm_{\tau_i}[k][V] \leftarrow cm_{\tau_i}[k-1][V] \cup min_{cm};$ 
19  end
20 end

```

no shutdown, the cumulative execution time up to layer  $k$  is calculated by summing  $\varepsilon_{\tau_i}[k-1][V]$  and  $a_i(x, k)$  (line 9). If a shutdown occurs during layer  $k$ , the total number of power cycles up to layer  $k-1$  will be different from the number up to layer  $k$  (obtainable using the pseudo-inverse function  $\bar{L}^e(t)$ ), and  $d_i(x, k)$  needs to be used instead (line 11). Once the algorithm finishes,  $\varepsilon_{\tau_i}[\eta_i][V_{ipd}]$  gives the minimum execution time of  $\tau_i$  under the memory constraint.

The next step is to minimize the collective sum of solo execution times of all  $m$  tasks under the device’s memory constraint. The memory usage of the system  $\Gamma$  is determined by adding up the peak memory usage of each task  $\tau_i \in \Gamma$ . Let us define  $E_\Gamma[i][V]$  as the minimum sum of the solo execution times of  $i$  tasks (from  $\tau_1$  to  $\tau_i$ ) in  $\Gamma$  with the memory constraint of  $V$ , and  $CM_\Gamma[i][V]$  as the corresponding CM information. This can be solved by dynamic programming with the following recurrence relation:

$$E_\Gamma[i][V] = \min_{1 \leq j \leq V-1} E_\Gamma[i-1][j] + \varepsilon_{\tau_i}[\eta_i][V-j] \quad (4)$$

and with  $j$  found for  $E_\Gamma[i][V]$ ,

$$CM_\Gamma[i][V] = CM_\Gamma[i-1][j] \cup \{cm_{\tau_i}[\eta_i][V-j]\} \quad (5)$$

The initial conditions are:

$$E_\Gamma[0][1\dots V_{ipd}] = 0, E_\Gamma[1][1\dots V_{ipd}] = \varepsilon_{\tau_1}[\eta_1][1\dots V_{ipd}], \text{ and } CM_\Gamma[0][1\dots V_{ipd}] = \emptyset, CM_\Gamma[1][1\dots V_{ipd}] = cm_{\tau_1}[\eta_1][1\dots V_{ipd}].$$

For all  $m$  tasks in  $\Gamma$ , the solution is given by  $E_\Gamma[m][V_{ipd}]$  and  $CM_\Gamma[m][V_{ipd}]$ . We use  $CM_\Gamma[m][V_{ipd}]$  as the initial CM settings of tasks when the system is deployed. Also, we store each task’s peak memory usage corresponding to the solution as  $M_{\tau_i}$  in the device’s NVM since it will be used as guidance by the online adaptation.

## VI. ONLINE PHASE

In light of Obs. 3 from Sec. III-B, real-world energy pattern fluctuations directly influence shutdown layers. Therefore, we address this issue in the online phase and use  $e'$  to indicate the environment condition at runtime that is different from the profiled condition  $e$ . We present our online scheduler in Sec. VI-A and the details of online CM adaption in Sec. VI-B.

---

**Algorithm 2: MII Online Scheduler**


---

```

1 Input:  $S^{e'}(n)$ ,  $L^{e'}(n)$ ,  $\hat{l}$  and  $t_{prev}$  in NVM;
2  $t_{start} \leftarrow \text{RTC\_now}()$ ;  $\hat{s} \leftarrow t_{start} - t_{prev}$ ;
3 Power Cycle Harmonizing (PCH) */;
4  $Q_{tasks} \leftarrow$  Tasks arrived by  $t_{start}$  but not finished their jobs;
5 for  $\forall \tau_i \in Q_{tasks}$  do
6   | Assign_Priority( $\tau_i$ ); */ LST */;
7 end
8 CMs Adaptation */;
9 if  $S^{e'}(1) < \hat{s} \vee L^{e'}(1) > \hat{l}$  then
10  | Update  $S^{e'}(n)$  and  $L^{e'}(n)$ ; */ Stored in NVM */;
11  | for  $\forall \tau_i \in Q_{tasks}$  do
12    |  $cm_{\tau_i} \leftarrow \text{Online\_Adaptation}(\tau_i)$ ;
13  | end
14 end
15 Task Scheduling */;
16 while  $Q_{tasks} \neq \emptyset$  do
17  |  $\tau_i \leftarrow \text{Pick\_Highest\_Priority}(Q_{tasks})$ ;
18  | if  $cm_{\tau_i} \neq \text{JIT}$  then
19    | Check_Proactive_Shutdown();
20  | end
21  | Run( $\tau_i$ );
22  | if  $\tau_i$  completed its job then
23    |  $Q_{tasks} \leftarrow Q_{tasks} \setminus \tau_i$ ;
24  | end
25 end
26 */  $Q_{tasks} = \emptyset$  or JIT threshold or Proactive Shutdown triggered */;
27 Save the states of JIT tasks to NVM;
28  $t_{prev} \leftarrow \text{RTC\_now}()$ ; */ Store in NVM */;
29  $\hat{l} \leftarrow t_{prev} - t_{start}$ ; IPD Shutdown;

```

---

### A. MII Online Scheduler

The main goal of our scheduler is to make task scheduling decisions based on  $e'$ . The pseudocode of the scheduler is given in Alg. 2, which begins upon each reboot. Upon start, the scheduler uses the on-board RTC to compute the *online shutdown time* of the current power cycle, denoted as  $\hat{s}$ , by taking the difference between the timestamp recorded at the previous shutdown,  $t_{prev}$ , and the current timestamp at the start,  $t_{start}$  (line 2). It also estimates the *online live time*  $\hat{l}$ , calculated at the end of the last power cycle (line 29). These  $\hat{s}$  and  $\hat{l}$  are used to determine the condition to trigger online CM adaptation, presented in the next subsection. The other components of the scheduler are explained below.

**Power Cycle Harmonizing (PCH).** An arbitrary arrival of periodic tasks is one of the reasons causing runtime deviations from the power cycle used by our offline search. To address this issue, our scheduler introduces PCH, which harmonizes task execution with the power cycle. PCH forces the scheduler to take into account only the tasks that have either arrived by  $t_{start}$  or those that have not finished their job execution within their periods (line 4,  $Q_{tasks}$ ). Therefore, the execution of any task that arrives during the live time of the current power cycle is deferred to the next power cycle, allowing the CM choices of tasks not to be disrupted by newly arriving tasks. Fig. 9 gives an example of the scheduling behavior with PCH. For  $n$ -th power cycle, the start of execution for all three tasks is harmonized to the end of  $\hat{s}_n$ , ensuring that the  $n + 1$  power cycle begins when the IPD turns off.

**Task Scheduling.** For the tasks found by PCH,  $Q_{tasks}$ , our scheduler uses a variant of the Least Slack Time (LST) scheduling policy to dynamically change task priorities, with each task's period as its deadline (lines 5-6). Our LST variant checks slack time at the boundary of each layer execution to

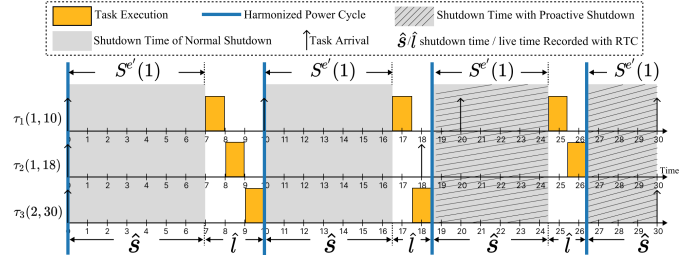


Fig. 9: MII online scheduler with 3 tasks. Each task  $\tau_i$  is characterized by (execution time, period). PCH delays the execution of  $\tau_2$ 's 2nd job. Proactive shutdown is shown by hatched grey boxes.

mitigate the overhead of the standard LST. The reason behind the use of LST is the following: Since DNN inference jobs are relatively long, their response times can be easily greater than their periods and the unscheduled jobs are skipped from execution. If we use other policies like EDF which is a job-level fixed-priority policy, a long-running job may dominate the live time over multiple power cycles as its priority does not change until completion, leading to starvation to other jobs and resulting in a disproportionate number of successful jobs per task compared to their periods. In other words, the use of LST can help achieve fairness in skipped jobs across tasks, as we will show in our evaluation. After updating CMs with the online adaption (lines 9-12), the scheduler proceeds to execute the tasks in  $Q_{tasks}$  in priority order and remove it from  $Q_{tasks}$  upon successful job completion (lines 16-25).

**Proactive Shutdown.** With ST, an IPD may shut down during the execution of an atomic block and re-execute it in the next power cycle. The re-executed portion is called *wasted work* [3, 9, 17, 18]. Although existing work mitigates this by voluntarily shutting down the IPD after a fixed number of calculations [5], it no longer works in a multi-task system with a mixture of JIT and ST because the JIT threshold can be triggered at any time of atomic block execution. To address this problem, we propose a proactive shutdown method. Proactive shutdown can be triggered by the scheduler before running any task  $\tau_i$  that uses ST as its CM (line 19). The scheduler proceeds with  $\tau_i$ 's execution only if the stored energy is enough to execute at least a single block of  $\tau_i$ . Otherwise, it makes the device shut down (line 26). This can also be triggered when no other task to execute ( $Q_{tasks} = \emptyset$ ) or the JIT threshold is met. The overhead of the proactive shutdown is negligible since it only requires an ADC reading at the end of each ST's block. Furthermore, for any JIT-enabled IPDs, ADC reading is already a prerequisite [10, 15, 16, 40].

One might have a concern that proactive shutdown turns the IPD off earlier and shortens the shutdown time  $\hat{s}$  (Fig. 9 hatched grey boxes), which might affect the condition to trigger online adaption. However, since the online CM adaption is triggered only when  $\hat{s}$  reaches a larger value than expected (explained in the next subsection), unnecessary online adaptations are effectively prevented.

### B. Online CM Adaption

To enable online CM adaptations, MII maintains  $S^{e'}(n)$  in NVM, which is an online version of the cumulative shutdown

TABLE II: DNNs Configuration

Test Case	DNN Name	Dataset	Layers Configuration	Params	Inputs	Top1 Acc.	Largest Layer
TC1	C7 [5]	CIFAR10 [34]	5xCONV2D+GP+FC	5,152	32x32x3	63%	63.9KB
TC1	M7 [5]	MNIST [33]	5xCONV2D+GP+FC	5,136	28x28x2	98%	49.0KB
TC1	C12 [5]	CIFAR10 [34]	10xCONV2D+GP+FC	15,722	32x32x3	78%	63.6KB
TC1	H5 [5]	HAR [35]	3xCONV1D+GP+FC	920	128x9	61%	3.5KB
TC2	FC4 [36]	KWS [37]	4xFC	263,436	49x10x1	87%	503.8KB
TC2	AutoEncoder [36]	ToyADMOS [38]	10xFC	269,992	1x640	85%	328.7KB
TC3	MBV1 [36]	VWW [39]	14xCONV2D+13xDCONV2D+AP+FC	221,794	96x96x3	80%	312.5KB
TC3	DSCNN [36]	KWS [37]	5xCONV2D+4xDCONV2D+AP+FC	24,908	49x10x1	90%	80.6KB

\* CONV2D/1D: Convolution 2D/1D, GP: Global Pooling, FC: Fully Connected, DCONV2D: Depthwise convolution 2D, AP: Average Pooling

time  $S^e(n)$  and initialized as  $S^{e'}(n) = S^e(n)$ . Also,  $L^{e'}(n)$ , an online version of the live time  $L^e(n)$ , is also maintained in NVM and initialized as  $L^{e'}(n) = L^e(n)$ . Both  $S^{e'}(n)$  and  $L^{e'}(n)$  are given as input to the scheduler since they are essential to quantify the deviation between the online environment condition  $e'$  and the profiled environment condition  $e$ .

Recall Defs. 1 and 2. If the IPD follows  $S^{e'}(n)$  and  $L^{e'}(n)$ , both  $\hat{s} \leq S^{e'}(1)$  and  $\hat{l} \geq L^{e'}(1)$  hold for one power cycle ( $n = 1$ ), indicating no shift of environment. Otherwise, either  $\hat{s} > S^{e'}(n = 1)$  or  $\hat{l} < L^{e'}(1)$  (line 9), meaning a potential shift of environment; hence, the scheduler first updates  $S^{e'}(n)$  and  $L^{e'}(n)$  with  $\hat{s}$  and  $\hat{l}$  (line 10), and then triggers the online adaptation (line 12). Since the direct use of the offline search algorithm (Sec. V-B) for online adaption introduces a huge overhead, we take a heuristic approach presented below to find a near-optimal solution by using the offline solution’s per-task memory usage,  $M_{\tau_i}$ , as a constraint for  $\tau_i$ .

The online adaptation is done individually for each task  $\tau_i \in Q_{tasks}$  to update task’s CM list  $cm_{\tau_i}[0 \dots \eta_i]$ . It initializes  $cm_{\tau_i}[1 \dots \eta_i]$  as follows: for each layer  $k$  of  $\tau_i$ ,  $cm_{\tau_i}[k] = \arg \min_x a_i(x, k)$ , where  $x \in \{\text{JIT, ST-L, ST-F, ST-T}\} \wedge v_i(x, k) \leq M_{\tau_i}$ . This discards any CM that causes the memory usage to exceed  $M_{\tau_i}$ , and chooses the CM giving the minimum execution time without considering shutdown during execution. It also uses a vector variable  $\varepsilon_{\tau_i}[1 \dots \eta_i]$  to keep track of each layer’s execution time corresponding to  $cm_{\tau_i}[1 \dots \eta_i]$ . Then, it takes the following steps to take into account the effect of shutdown for each  $\tau_i$ :

1. Start with the first power cycle  $n = 1$ .
2. Find a layer  $K$  that experiences a shutdown in the power cycle  $p$ . Such a layer  $K$  satisfies:  $\sum_{k=1}^{K-1} \varepsilon_{\tau_i}[k] \leq L^{e'}(n)$  and  $\sum_{k=1}^K \varepsilon_{\tau_i}[k] > L^{e'}(n)$ .
3. If  $K$  is found, this means the layer  $K$  is a shutdown layer. Hence, update  $\varepsilon_{\tau_i}[K] = \min_x d_i(x, K)$  and  $cm_{\tau_i} = \arg \min_x d_i(x, K)$ .
4. Repeat steps 2-3 until  $K$  reaches  $\eta_i$ .

This method takes a linear search approach, much faster than the offline algorithm. However, its optimality may be compromised due to its reliance on  $M_{\tau_i}$  determined offline. Nonetheless, by incorporating  $L^{e'}(n)$  which is continuously updated for the current environment  $e'$ , this approach offers substantial benefits as we will demonstrate in the evaluation.

## VII. MII IMPLEMENTATION

**Energy Source.** For the continuous power source, we used an X-NUCLEO-LPM01A [41], which provides a stable voltage of 1.8V to the IPD. The current consumption and

live voltage can be observed via a connected PC that directly controls LPM01A. For the intermittent energy source, we harvested energy using an LTC3588 energy harvester [42] and solar cells of 1.5W peak power [43]. We regulated the input voltage of IPD to 1.8V and stored the harvested energy in a set of capacitors with the size of 1mF. During operations, the capacitor’s voltage range is 2.87v to 4.03v [42].

**Evaluation Hardware.** We chose the Ambiq Apollo4 Blue Plus evaluation board [1] that has an ARM Cortex-M4 MCU, 2MB MRAM as NVM, an on-board RTC unit for timekeeping,<sup>3</sup> and an on-board ADC for capacitor voltage monitoring. In the version of FreeRTOS [45] provided by the board manufacturer, the heap area size is configured to 512KB. Therefore, we use 512KB as the device memory constraint for inference tasks in the system.

**DNN Setup.** To cover a wide range of modern DNN models runnable on IPDs, we selected 8 DNN models and categorized them into 3 test cases (TCs) based on their dominant layers. TC1 is obtained from existing inferences on IPDs researches [3–5], whereas TC2-3 are obtained from the MLPerf Tiny benchmark [36] covering anomaly detection, visual wake words, and keyword spotting. Details of these DNNs are found in Table II. **TC1** comprises 4 DNNs dominated by convolution 1D/2D (CONV1D/2D) layers. **TC2** includes 2 DNNs that only use fully connected (FC) layers. **TC3** includes MobileNetV1 (MBV1) [36] and Depthwise convolution neural network (DSCNN) [36], both of which predominantly use depth-wise 2D convolution (DW-CONV2D) layers.

**Method Implementation.** For the MII offline phase, we developed a PC-based profiler that enables profiling the execution time of tasks on the IPD in both continuous and intermittent power settings. It also includes our search algorithm described in Sec. V-B. The offline-searched CMs were written to the INFO0 section of NVM before running inference tasks. For the MII online phase, we implemented the MII scheduler in FreeRTOS [45] for the Apollo4 board. The online adaption overwrites offline CMs in NVM by atomically changing the CM list for each task when needed. Our current implementation did not use hardware acceleration for DNN execution. However, since MII performs all computations in the VM for both JIT and ST, it can be safely adapted to hardware acceleration features that usually require direct access to data in the VM, e.g., TI’s Low Energy Accelerator (LEA). We leave such extensions as part of future work.

<sup>3</sup>The RTC circuit ran with its own dedicated capacitor and it did not deplete during the experiment. For more reliable timekeeping, techniques like “persistent clocks” [44] can be considered.



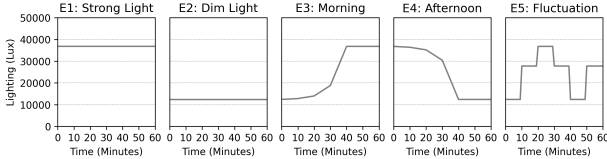


Fig. 10: Five Controlled Energy Patterns represented as the lighting in a unit of illumination (Lux). E1 and E2 are static lighting, whereas E3, E4, and E5 are dynamic lighting.

## VIII. CONTROLLED ENVIRONMENT EVALUATION

### A. Evaluation Setup

**Controlled Energy Patterns.** We conducted a 3 months study of lighting condition changes inside a greenhouse environment and identified 6 distinct energy patterns (E0-E5) that can capture more than 97% of the lighting conditions during daytime. Apart from E0, the continuous power, E1 to E5 are all intermittent power, and their energy patterns in one hour are shown in Fig. 10. Specifically, E1 and E2 represent static energy patterns, indicating unchanged light conditions throughout the duration. E3, E4, and E5 represent dynamic energy patterns that are either changing gradually or abruptly during the one-hour period.

**Baseline Configuration.** We compare MII against three established and state-of-the-art methods: QuickRecall [15] which is a JIT checkpointing-only system, iNAS [5] which uses ST-T with a tile size determined offline, and Zygarde [4] which uses ST-F and an early-exit model comprising mandatory and optional layers. For a fair comparison, we adjusted several settings of each method: In QuickRecall, we determined the JIT threshold to ensure the successful checkpointing of all JIT tasks in the system. In Zygarde, we set the first layer of each model as a mandatory layer since at least one layer OFM is needed for early exit. Also, we marked an inference result from a mandatory job as successful if it matched the result of a complete job. It is worth noting that the evaluation of Zygarde subsumes SONIC [3], which also focuses on intermittent DNN inference, because Zygarde extends SONIC’s APIs and has shown to outperform SONIC. In iNAS, we derived appropriate tile parameters by balancing peak memory usage and execution time. Although a larger tile size can shorten execution time by improving data reuse, it increases peak memory size, leading to out-of-memory if multiple tasks are executed concurrently.

**Test Cases and Job Generation.** The evaluation encompasses 8 DNNs in 3 test cases (TCs) given in Table II. To evaluate the performance of individual and combined TCs that represent different task sizes, we consider four scenarios: (A) **All TC1**: all four DNNs in TC1, (B) **TC1+TC2**: two DNNs (C7, M7) selected from TC1 and one model (FC4) from TC2, (C) **TC2+TC3**: one model (AutoEncoder) from TC2 and one model (DSCNN) from TC3, and (D) **TC1+TC3**: two DNNs (C7, C12) from TC1 and one model (MBV1) from TC3. Each model is executed by a distinct task on FreeRTOS. We schedule the tasks with specific periods, as will be shown by the number of generated jobs in later sections. We scale task periods w.r.t. the energy pattern because less lighting causes the sensor to sample fewer readings. For example, in All TC1 with E0, we set the period of C7 as 3.5s, C12 as 4.5s, M7 as

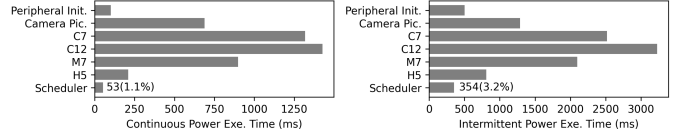


Fig. 11: System runtime breakdown when running 4 DNNs of TC1 under continuous (left) and intermittent (right) power.

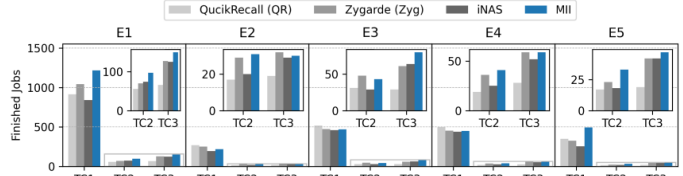


Fig. 12: Finished Jobs by running each TC separately in E1-E5. MII only uses offline phase to search CMs under E1.

3.6s, and H5 as 2.8s, whereas in All TC1 with E1, we multiply these periods by a factor of 1.6 due to longer harvesting time.

**System Overhead.** To evaluate the overhead of MII’s runtime scheduler and other non-inference tasks, we profiled the execution time of each software component when running all TC1 models on one downsampled image. The IPD first took a picture from the camera, stored a downsampled version in NVM, and then processed the image by running the inference of each DNN sequentially. Fig. 11 shows the runtime breakdown of the described workload under both continuous power provided by LPM01A [41] and an intermittent power source following E1 energy pattern. Recorded by a Discovery 3 digital oscilloscope [46], the IPD experiences 9 power failures under the intermittent power supply. The overhead of MII scheduler is composed of only 1.1% under continuous power and 3.2% under intermittent power of the entire system runtime.

### B. Offline Effectiveness

**CM Search Algorithm.** To evaluate the effectiveness of our offline search across various energy patterns, we use the offline searched optimal CMs found from the energy pattern E1 and apply these CMs to E2-E5. During a one-hour period, we measure for each inference task how many jobs execute successfully (denoted as *finished jobs*). For a level comparison, Zygarde’s finished jobs are captured by executing the complete DNN inference without relying on the early-exit feature. This feature will be assessed in the online evaluation in Sec. VIII-C. The results in Fig. 12 showcase that, for E1, MII outperforms QuickRecall, Zygarde, and iNAS by completing 41%, 18% and 40% more jobs, respectively. This is somewhat expected since the CMs are searched using E1. Although QuickRecall finished 7% more jobs on average than MII for smaller DNNs from TC1 in E2-4, it suffers from out-of-memory when running large DNNs and finishes the least TC2-3 inference jobs across all baselines. Interestingly, even when the IPD runs inference tasks using the E1-optimized CMs for E2-5, MII still achieves 10% and 28% more jobs than Zygarde and iNAS, respectively. These results show that the CMs searched by MII in a single energy pattern demonstrate efficacy across all examined energy patterns compared to the other methods. This underscores the robust generality of our search algorithm.

**Peak Memory.** Fig. 14 shows the peak memory usage when applying the CMs obtained from MII’s offline methods

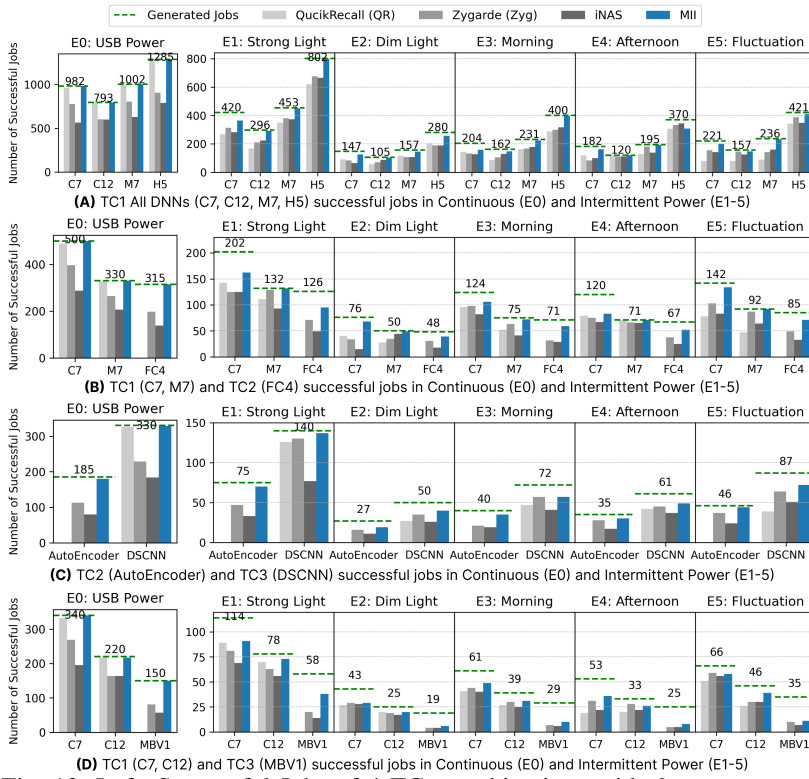


Fig. 13: Left: Successful Jobs of 4 TCs combinations with 6 energy patterns. Task periods are scaled with the energy pattern and are indicated by the numbers above green dashed lines. Right: Response Time of TCs under one stable energy pattern.

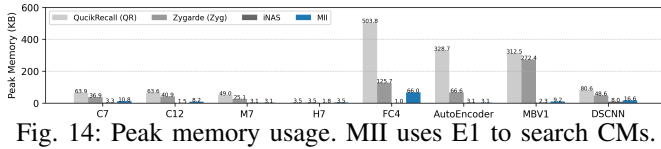
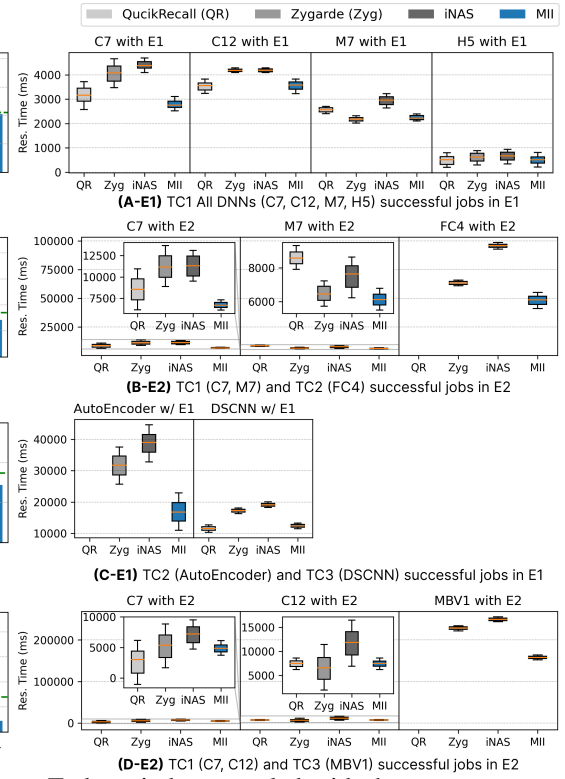


Fig. 14: Peak memory usage. MII uses E1 to search CMs.

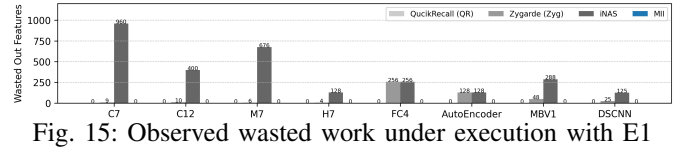


Fig. 15: Observed wasted work under execution with E1

under the E1 energy pattern. From these results, MII can effectively reduce the peak memory compared to QuickRecall and Zygarde. QuickRecall consistently accesses the layer’s IFM, OFM, and weights, resulting in the highest peak memory. Although Zygarde often needs to load the entire layer’s OFM for an early exit, its memory usage is still smaller than QuickRecall. iNAS has the smallest because it only loads partial IFM, weights, and OFM in each power cycle; however, it suffers from increased execution time. MII’s reduction in peak memory, compared to QuickRecall and Zygarde, is primarily due to that MII can choose a non-JIT CM or finer-grained ST for a large layer by imposing the memory constraint.

**MII Offline Phase:** The derived optimal layer-wise CMs effectively tackle the challenges highlighted in Obs. 1 and Obs. 2 and are applicable to various energy patterns.

### C. Online Scheduling and Adaption Effectiveness

**Successful Jobs.** The left part of Fig. 13 illustrates the number of successful jobs in four combinations of TC1-4, with the total number of generated jobs indicated by green dashed lines. If an optimal solution exists for each online energy pattern (E1-E5), its number of successful jobs is upper-bounded by the green dashed lines. Hence, the gap in each subplot between the green line and each bar represents the deviation between each baseline and the optimal solution. Although it does not reach the optimal performance because

of the heuristic approach used in online adaption, MII still achieves on average a 21% increase of successful jobs in stable energy conditions (E1-E2) and on average 39% successful jobs increase under dynamic energy conditions (E3-E5) compared to other methods. When comparing these results with the offline evaluation, MII’s online phase significantly improves job success rates over the other three methods. Specifically, for (A) TC1 All in E1, there is a 56% increase in successful jobs, thanks to the online scheduler and adaptation of MII.

In addition to the MII’s significant improvement in total successful jobs, the following two observations can be made: (i) The LST variant of MII mitigates starvation of short-period jobs and helps achieve fairness across tasks. For example, in Fig. 13(B)-E5, when using MII, the short-period task C7 can complete up to 94% of its generated jobs, and the long-period tasks M7 and FC4 complete 99% and 84% of their generated jobs. On the other hand, when using Zygarde, which uses an EDF variant, one of the long-period tasks M7 completes 95% while the other two tasks complete only 73% and 58% of their jobs, showing a much higher discrepancy in successful jobs per task than MII. (ii) If we ignore the correctness of inference results of jobs, Zygarde has executed more jobs than MII in most cases due to its early-exit feature, which executes only a mandatory portion of the job. However, if a task uses a large DNN model, this feature makes Zygarde suffer from low accuracy, resulting in a lower number of successful jobs.

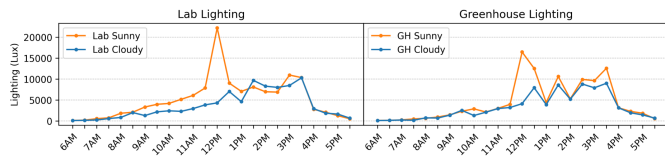


Fig. 16: Real-world solar energy pattern

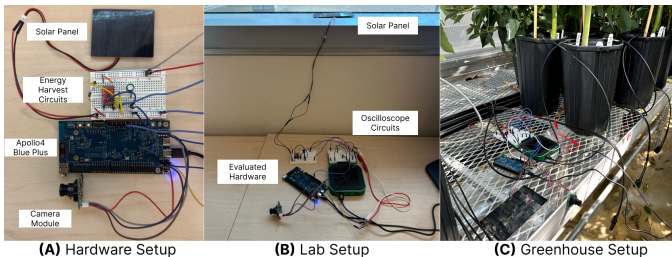


Fig. 17: Experiment setup in Lab and Greenhouse.

Therefore, in Fig. 13(D)-E5, although Zygard achieves 2% more successful jobs than MII for a small DNN task C7, it achieves 30% and 10% less successful jobs than MII for large DNNs tasks C12 and MBV1, respectively.

**Response Time.** The right side of Fig. 13 showcases the response time of DNN inference tasks under one stable energy pattern (E1 or E2). This supplements the successful job results. Although the response time is not the major optimization objective of MII (successful job is), MII tends to yield more preferable results overall, with a shorter average response time and a smaller variation. QuickRecall suffers from the out-of-memory issue when executing a large model such as AutoEncoder and MBV1. Hence, although it gives a shorter response time for small models (A-E1 in the figure), it completes 35%, 40%, and 22% less successful jobs on average compared to MII for TC1+TC2, TC2+TC3, and TC1+TC3, respectively. Zygard often gives a shorter response time due to its early-exit feature, but it comes at the cost of low accuracy.

**Wasted Work.** We characterize the amount of wasted work by profiling the number of output features discarded during shutdown. Fig. 15 depicts the wasted work of each DNN when running the four combinations of TCs under the E1 energy pattern. Both Zygard and iNAS show some wasted work due to the issue discussed in Sec. VI-A. For QuickRecall, since it uses JIT, it obviously has zero wasted work. However, to ensure that JIT successfully checkpoints all tasks' largest layers, a higher JIT threshold is required for QuickRecall compared to MII, i.e., 3.7V vs. 3.3v. MII has zero wasted work while keeping a smaller JIT threshold because of our proactive shutdown technique (Sec. VI-A) as well as the ability to avoid using JIT for large layers.

## IX. REAL WORLD EVALUATION

**Experiment Setting.** The IPD was deployed in two settings: lab window side (Lab) and Greenhouse (GH). We evaluated the system under Sunny and Cloudy conditions, both occurring within a single day. The lighting changes for each setting are depicted in Fig. 16. In Lab, the IPD was positioned under direct sunlight with minimal interference. However, in GH, IPD faced consistent interference from leaf shades and plant shadows. To demonstrate effectiveness in common smart sensing applications, all TC1 DNNs listed in Table II were

TABLE III: Successful jobs of MII compared to Zygard

DNN	Lab Sunny	Lab Cloudy	GH Sunny	GH Cloudy
C7	+17%	+51%	+20%	+61%
C12	+38%	+47%	+38%	+9%
M7	+18%	+34%	+33%	+6%
H5	+19%	+35%	+33%	-7%

selected for the experiment and ran continuously from 6AM to 6PM on a day (Fig. 16). A comparison setup, running the same set of DNNs with Zygard [4], was placed adjacent to the experimental setup. Zygard was chosen as it is state-of-the-art and capable of adapting to environmental changes. Fig. 17 illustrates our setup in both Lab and GH.

**Successful Jobs.** Table III presents the percentage difference in the number of successful jobs executed by MII compared to Zygard. In the GH Cloudy setting, MII outperformed Zygard for all three DNNs, executing 61%, 9%, and 6% more jobs than Zygard. However, for H5 in the same setting, MII completed 7% fewer jobs. It was mainly due to that, although Zygard executed only mandatory layers under the cloudy condition, it still managed to produce around 95% of inference results that matched the results of H5's complete inference. However, this favorable outcome for Zygard is largely confined to small models with inherently low accuracy, which are not substantially affected by the accuracy degradation of early exits. Overall, MII outperformed Zygard by an average of 33% more successful jobs in the Lab and 24% in the GH.

**II Online Phase:** With the runtime scheduler and adaptation methods, MII outperforms the state-of-the-art across diverse energy patterns and in real-world environments, addressing the challenge from Obs. 3.

## X. RELATED WORK

### A. Intermittent Computing

In the context of IPD software systems, prior work has focused on issues such as memory inconsistency [7–9, 12, 13, 15, 29, 47], task idempotency [10, 12, 16, 18, 26], and sensing/computation task coordination [14, 17, 19, 31]. There are numerous studies on IPD hardware improvement [11, 40, 48, 49] and energy harvesting circuits [27, 50].

Emergent research about running machine learning workloads on IPDs is still in its very early stage. Existing work focuses on learning [25] as well as inference tasks [3–6, 18, 26, 31] on IPDs. However, none of these has studied the tradeoff between JIT [10, 15, 16] and ST [3, 5, 12, 14, 17–19] checkpointing mechanisms as well as the different granularity of atomic blocks in ST for DNN inference tasks. Although there have been some attempts to co-use JIT and ST in the same system [17, 31], they use ST exclusively for peripheral access and JIT for computational tasks, meaning that when applied to inference tasks, all will be governed by JIT.

Recent work [40] has proposed to put the device into sleep mode and trigger JIT checkpointing only when no more energy is available. While this has the potential to improve the standard JIT method used in MII, depending on the leakage current in sleep mode, it can increase the recharge time of the capacitor. For example, in our evaluation platform [1], using its default deep sleep mode results in at least 58% more time to fully recharge under the E2 energy pattern in Sec. VIII-A.

## B. Real-time Scheduling on IPD

Supporting task Scheduling on IPDs has been considered a big challenge because environment fluctuations lead to varying task response times [14, 19]. Existing work has approached this issue in two ways: energy prediction [21–23], which analyzes the supply energy pattern and makes scheduling decisions accordingly; and workload reduction [4, 19, 29–31], which reactively changes the amount of workload with respect to environment variations. Specifically, Celebi [21] leverages the energy prediction approach and presents both offline and online schedulers to meet task deadlines on IPDs. Zygarde [4] trains a DNN to be early exit-able at every layer and splits the layers into mandatory and optional parts. It then schedules either only mandatory layers or both types of layers based on ambient conditions. However, none of the existing work has studied the effect of execution patterns given different CMS. MII addresses these limitations.

## XI. CONCLUSION

This paper presents MII: a Multifaceted framework for Intermittent Inference and scheduling. The design of MII originated from the three key observations in Sec. III and divided into offline and online phases. MII is compared with three representative state-of-the-art methods [4, 5, 15] through a controlled environment evaluation as well as a real-world field study. The results show that MII is able to achieve performance efficiency in intermittent DNN inference, adaptability to environment changes, and applicability to real-world scenarios. Future work includes the consideration of different learning tasks, such as deep reinforcement learning, on IPDs.

## ACKNOWLEDGEMENT

This work was supported by the U.S. NSF grants CNS 1943265, CPS 2230969, CNS 2300525, CNS 2343653, CNS 2312397 and the USDA/NIFA SCRI grant 2020-51181-32198. The work was also supported by the IITP grant (No.2021-0-00360) funded by the Korean government (MSIT).

## REFERENCES

- [1] Ambiq, “Apollo4 Blue Plus KBR SoC Eval Board,” 2022.
- [2] Texas Instruments, “MSP430FR5994,” 2022.
- [3] G. Gobieski *et al.*, “Intelligence Beyond the Edge: Inference on Intermittent Embedded Systems,” in *ASPLOS*, 2019.
- [4] B. Islam *et al.*, “Zygarde: Time-sensitive on-device deep inference and adaptation on intermittently-powered systems,” *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, vol. 4, no. 3, sep 2020.
- [5] H. R. Mendis *et al.*, “Intermittent-aware neural architecture search,” *ACM Trans. Embed. Comput. Syst.*, vol. 20, no. 5s, sep 2021.
- [6] C.-K. Kang *et al.*, “Everything leaves footprints: Hardware accelerated intermittent deep inference,” *IEEE TCAD*, vol. 39, no. 11, 2020.
- [7] H. Volos *et al.*, “Mnemosyne: Lightweight persistent memory,” *SIGARCH Comput. Archit. News*, vol. 39, no. 1, p. 91–104, mar 2011.
- [8] M. Buettner *et al.*, “Dewdrop: An Energy-Aware runtime for computational RFID,” in *NSDI*, 2011.
- [9] B. Lucia *et al.*, “A simpler, safer programming and execution model for intermittent systems,” in *PLDI*, 2015.
- [10] D. Balsamo *et al.*, “Hibernus: Sustaining computation during intermittent supply for energy-harvesting systems,” *IEEE Embedded Systems Letters*, vol. 7, no. 1, pp. 15–18, 2015.
- [11] J. V. D. Woude *et al.*, “Intermittent computation without hardware support or programmer intervention,” in *OSDI*, Nov. 2016.
- [12] B. Ransford *et al.*, “Mementos: System support for long-running computation on rfid-scale devices,” in *ASPLOS*, 2011.
- [13] A. Mirhoseini *et al.*, “Idetic: A high-level synthesis approach for enabling long computations on transiently-powered asics,” in *PerCom*, 2013.
- [14] J. Hester *et al.*, “Timely execution on intermittently powered batteryless sensors,” in *SenSys*, 2017.
- [15] H. Jayakumar *et al.*, “Quickrecall: A hw/sw approach for computing across power cycles in transiently powered computers,” *J. Emerg. Technol. Comput. Syst.*, vol. 12, no. 1, aug 2015.
- [16] D. Balsamo *et al.*, “Hibernus++: A Self-Calibrating and Adaptive System for Transiently-Powered Embedded Devices,” *IEEE TCAD*, vol. 35, no. 12, pp. 1968–1980, 2016.
- [17] K. Maeng *et al.*, “Supporting peripherals in intermittent systems with just-in-time checkpoints,” in *PLDI*, 2019.
- [18] —, “Alpaca: Intermittent Execution without Checkpoints,” *Proc. ACM Program. Lang.*, vol. 1, no. OOPSLA, oct 2017.
- [19] —, “Adaptive dynamic checkpointing for safe efficient intermittent computing,” in *OSDI*, 2018.
- [20] M. Chetto *et al.*, “Feasibility analysis of periodic real-time systems with energy harvesting capabilities,” *Sustainable Computing: Informatics and Systems*, vol. 22, pp. 84–95, 2019.
- [21] B. Islam *et al.*, “Scheduling computational and energy harvesting tasks in deadline-aware intermittent systems,” in *RTAS*, 2020, pp. 95–109.
- [22] M. Karimi *et al.*, “Real-time task scheduling on intermittently powered batteryless devices,” *IEEE Internet of Things Journal*, vol. 8, no. 17, pp. 13 328–13 342, 2021.
- [23] —, “Energy-adaptive real-time sensing for batteryless devices,” in *RTCSA*, aug 2022.
- [24] A. Bukhari *et al.*, “Opense: An open-world sensing framework for incremental learning and dynamic sensor scheduling on embedded edge devices,” *IEEE Internet of Things Journal*, 2024.
- [25] S. Lee *et al.*, “Intermittent learning: On-device machine learning on intermittently powered system,” *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, vol. 3, no. 4, sep 2020.
- [26] A. Colin *et al.*, “Chain: Tasks and channels for reliable intermittent programs,” *SIGPLAN Not.*, vol. 51, no. 10, p. 514–530, oct 2016.
- [27] A. Colin, E. Ruppel, and B. Lucia, “A reconfigurable energy storage architecture for energy-harvesting devices,” in *ASPLOS*, 2018.
- [28] J. Choi *et al.*, “Compiler-directed high-performance intermittent computation with power failure immunity,” in *RTAS*, 2022.
- [29] K. S. Yildirim *et al.*, “Ink: Reactive kernel for tiny batteryless sensors,” in *SenSys*, 2018.
- [30] K. Maeng *et al.*, “Adaptive low-overhead scheduling for periodic and reactive intermittent execution,” in *PLDI*, 2020.
- [31] M. Surbatovich *et al.*, “A type system for safe intermittent computing,” *Proc. ACM Program. Lang.*, vol. 7, no. PLDI, jun 2023.
- [32] H. Hu *et al.*, “Learning anytime predictions in neural networks via adaptive loss balancing,” in *AAAI*, 2017.
- [33] Y. LeCun *et al.*, “The mnist database of handwritten digits,” 2005.
- [34] A. Krizhevsky *et al.*, “Learning multiple layers of features from tiny images,” University of Toronto, Tech. Rep. 0, 2009.
- [35] J. Reyes-Ortiz *et al.*, “Human Activity Recognition Using Smartphones,” 2012.
- [36] C. Banbury *et al.*, “Mlperf tiny benchmark,” *arXiv preprint*, 2021.
- [37] P. Warden, “Speech commands: A dataset for limited-vocabulary speech recognition,” *CoRR*, vol. abs/1804.03209, 2018.
- [38] Y. Koizumi *et al.*, “Toyadmos: A dataset of miniature-machine operating sounds for anomalous sound detection,” 2019.
- [39] A. Chowdhery *et al.*, “Visual wake words dataset,” 2019.
- [40] K. Akhunov, E. Yildiz, and K. S. Yildirim, “Enabling efficient intermittent computing on brand new microcontrollers via tracking programmable voltage thresholds,” in *ENSYS*, 2023.
- [41] STMicroelectronics, “X-NUCLEO-LPM01A,” 2018.
- [42] SparkFun, “Sparkfun energy harvester breakout - LTC3588,” 2021.
- [43] L. Seed Technology Co., “Monocrystaln solar cell 1.5w 8.2v,” 2021.
- [44] J. Hester *et al.*, “Persistent clocks for batteryless sensing devices,” *ACM TECS*, vol. 15, no. 4, pp. 1–28, 2016.
- [45] “FreeRTOS™,” 2022, <https://www.freertos.org/index.html>.
- [46] Digilent, “Analog Discovery 3,” 2023.
- [47] A. Colin *et al.*, “Termination checking and task decomposition for task-based intermittent programs,” in *CC*, 2018.
- [48] M. Surbatovich *et al.*, “Towards a formal foundation of intermittent computing,” *Proc. ACM Program. Lang.*, vol. 4, no. OOPSLA, nov 2020.
- [49] A. Colin *et al.*, “An energy-interference-free hardware-software debugger for intermittent energy-harvesting systems,” in *ASPLOS*, 2016.
- [50] P. Nintanavongsa *et al.*, “Design optimization and implementation for rf energy harvesting circuits,” *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 2, no. 1, pp. 24–33, 2012.