# MNN-AECS: Energy Optimization for LLM Decoding on Mobile Devices via Adaptive Core Selection

### Zhengxiang Huang
Shanghai Jiao Tong University
Shanghai, China
huangzhengxiang@sjtu.edu.cn

### Chaoyue Niu[*]
Shanghai Jiao Tong University
Shanghai, China
rvince@sjtu.edu.cn

### Zhaode Wang
Alibaba Group
Hangzhou, China

### Jiarui Xue
Shanghai Jiao Tong University
Shanghai, China

### Hanming Zhang
Shanghai Jiao Tong University
Shanghai, China

### Yugang Wang
Shanghai Jiao Tong University
Shanghai, China

### Zewei Xin
Shanghai Jiao Tong University
Shanghai, China

### Xiaotang Jiang
Alibaba Group
Hangzhou, Zhejiang, China

### Chengfei Lv
Alibaba Group
Hangzhou, Zhejiang, China

### Fan Wu
Shanghai Jiao Tong University
Shanghai, China

### Guihai Chen
Shanghai Jiao Tong University
Shanghai, China

## ABSTRACT

As the demand for on-device Large Language Model (LLM) inference grows, energy efficiency has become a major concern, especially for battery-limited mobile devices. Our analysis shows that the memory-bound LLM decode phase dominates energy use, and yet most existing works focus on accelerating the prefill phase, neglecting energy concerns. We introduce Adaptive Energy-Centric Core Selection (AECS) and integrate it into MNN to create the energy-efficient version, MNN-AECS, the first engine-level system solution without requiring root access or OS modifications for energy-efficient LLM decoding. MNN-AECS is designed to reduce LLM decoding energy while keeping decode speed within an acceptable slowdown threshold by dynamically selecting low-power CPU cores. MNN-AECS is evaluated across 5 Android and 2 iOS devices on 5 popular LLMs of various sizes. Compared to original MNN, MNN-AECS cuts down energy use by 23% without slowdown averaged over all 7 devices and 4 datasets. Against other engines, including llama.cpp, executorch, mllm, and MediaPipe, MNN-AECS delivers 39%–78% energy saving and 12%–363% speedup on average.

## CCS CONCEPTS

• **Human-centered computing** → **Ubiquitous and mobile computing**; • **Software and its engineering** → **Power management**.
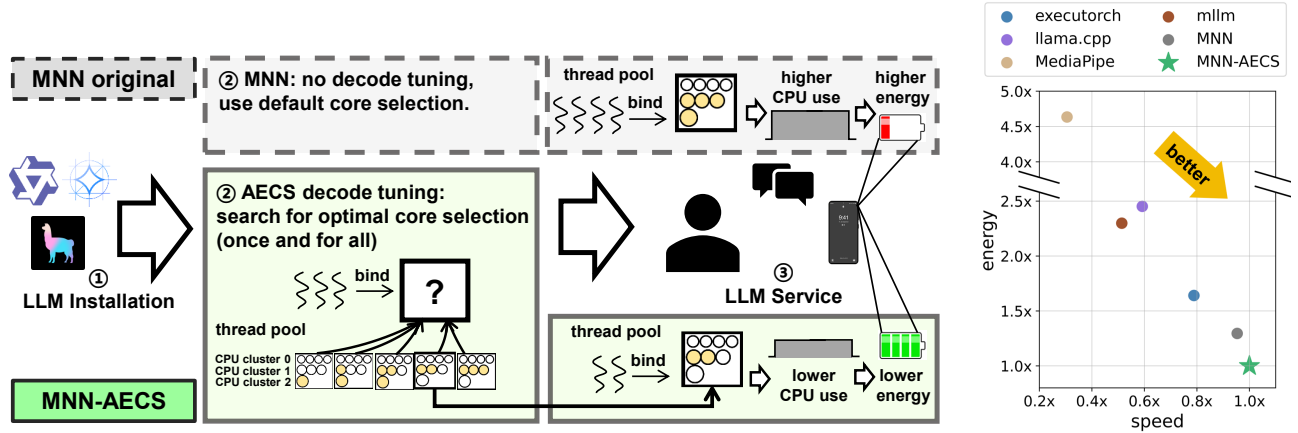
## KEYWORDS

On-Device LLM Inference, Energy Optimization

## 1 INTRODUCTION

### 1.1 Importance of On-Device LLM Inference

The daily demands for Large Language Model (LLM) service are rapidly growing, with most serving workload being processed on cloud servers currently [17, 55]. Yet, such on-cloud LLM service incurs high expenses to LLM service providers and raises privacy concerns to end users.

By contrast, on-device LLM inference as an alternative solution addresses these 2 issues: it significantly relieves workload on cloud servers and reduces serving expenses, and it's also much more privacy-preserving because no queries and generated texts are ever exposed to the servers. These advantages have garnered increasing attention from both LLM teams and mobile phone companies. Qwen2.5 series [42] provide distilled small models (0.5B, 1.5B, 3B) for on-device use cases in MNN [33] and GGUF [18, 19] format. Llama3.2 [7] collaborates with executorch [8] to launch on-device inference pipeline of its 1B/3B models. Google launches Gemini Nano (2B) with AI Edge SDK [20]. Apple Intelligence [22] adopts device-cloud collaborative serving framework [12], where a 3B model handles simple tasks on mobile devices to relieve cloud server workload.

---
[*]Chaoyue Niu is the corresponding author (rvince@sjtu.edu.cn).

(a) The 8 circles in the thread pool box represent 8 cores, and the yellow ones are selected. Bigger circles represent bigger cores in big-little core CPU architecture. MNN-AECS searches and selects less and smaller cores than original MNN in decode phase to reduce CPU utilization and frequency, resulting in lower energy.

(b) Energy and speed are presented in ratio relative to MNN-AECS.

**Figure 1: (a) MNN-AECS workflow. (b) MNN-AECS geometric mean performance over 5 baseline engines across 7 devices and 4 datasets. MNN-AECS (green star) consumes the lowest energy and is also the fastest.**

**Table 1: Energy consumption and CPU utilization of Qwen2.5-1.5B (4-bit quantized), running 20 data entries of multi-turn conversation sampled from ShareGPT dataset, with inference engine MNN.**

| device | total battery | consumption battery | energy | CPU use | power |
|--------|-------|---------|--------|---------|-------|
| Xiaomi 15 Pro | 6100 mAh | 386 mAh (6%) | 6031 J | 395% | 9.9 W |
| Mate 40 Pro | 4400 mAh | 717 mAh (16%) | 10438 J | 396% | 8.7 W |
| iPhone 12 | 2815 mAh | 708 mAh (25%) | 10379 J | 252% | 7.9 W |

## 1.2 On-Device LLM Inference Energy Issue

Energy is one of the major issues for on-device LLM inference. Heavy workload of LLM inference leads to high hardware utilization and power and prolongs execution time, resulting in intolerably high energy consumption (Table 1), draining mobile device battery and raising CPU temperature.

**Heavy workload of LLM inference.** Compared to traditional smaller discriminative Deep Learning (DL) models, LLMs induces much higher energy consumption due to their heavier workload caused by large model sizes and autoregressive decoding characteristic. 1) Model sizes of mobile LLMs typically scale up to 3B [7, 40, 42] currently, leading to 3-10 GFLOPS/token computation and up to 1-3 GB/token DRAM memory visits even after 4/8-bit quantization, resulting in heavy per-token workload both on processors and memory bus. Processor utilization is boosted to nearly maximum, raising the device power to as high as 8-10 W, as is shown in Table 1. 2) For each auto-regressive decoding loop, LLM inputs the last previously generated token and outputs a new token one at a time until reaching the end of sentence. Usually hundreds of tokens are generated in each round of

LLM decoding, thus prolonging execution time to minute-level, cumulating the heavy per-token energy consumption to an unprecedentedly high total amount.

**Limited battery capacity** of mobile devices is an important hardware constraint given the high energy consumption of LLM inference. Unlike cloud servers provisioned with continuous power supply, mobile phones, for unplugged and ubiquitous use, are equipped with limited incorporated battery. Table 1 shows the battery capacity of 3 common mobile devices: Xiaomi 15 Pro (high-end), Huawei Mate 40 Pro and iPhone 12 (middle-end), ranging from 2815 to 6100 mAh.

**Consequences of high energy consumption.** 1) The battery runs out too fast, which is the direct consequence of high energy consumption. For example, in Table 1, 20 conversations sampled from ShareGPT [39] consume up to $6\% \sim 25\%$ battery in less than 15 minutes. Hence, users may suffer from the inconvenience of frequent charging. 2) High temperature can be a malicious side-effect of high cumulative energy consumption. CPU temperature can climb up to $80°C$ during LLM inference. High temperature may degrade user experience due to CPU throttling [14] and gradually damage the hardware [27, 37].

## 1.3 Our Work

In this work, we target at the energy challenge of on-device LLM inference. Our energy analysis of the two LLM phases, prefill and decode, reveals that decode phase is the primary contributor to the total LLM inference energy consumption on common on-device text generation tasks. Thus, we focus our design on optimizing LLM decode energy.

We introduce a novel method: Adaptive Energy-centric Core Selection (AECS), to optimize the LLM decode energy

with speed guarantee (no more than $\epsilon$ decode slowdown). We also integrate it into on-device LLM engine MNN [33] and propose our energy-efficient version of it: MNN-AECS.

Figure 1(a) shows the workflow of MNN-AECS. Between installation and LLM service, we add a once-and-for-all AECS decode tuning, which searches for the optimal core selection that reduces decode energy with negligible speed loss by reducing CPU utilization, where a CPU core selection is a core binding plan on Android and a thread number on iOS. AECS adopts a heuristic search rather than a naive exhaustive one to guarantee the optimality of search results and also reduce search time by up to 90%. After tuning, the searched optimal core selection is input to the thread pool as decode configuration for all future services. MNN is modified so that prefill and decode can adopt different core selections, eliminating interference between prefill and decode.

To position our work, **MNN-AECS is a pure engine-level system solution** implemented outside OS without root access and also orthogonal to any algorithmic designs.

To evaluate MNN-AECS, we develop a cross-platform LLM energy evaluation testbed to measure energy and speed of on-device LLM engines on both Android and iOS devices. The energy profiling module of the testbed provides precise energy data probed from OS interfaces and passed to MNN-AECS and evaluation pipeline though JNI on Android and tunneld for iOS. MNN-AECS is compared against original MNN, llama.cpp, executorch, mllm, and MediaPipe across 5 Android and 2 iOS devices, using 5 LLMs: Qwen2.5-1.5B/3B, Llama3.2-1B/3B, and Gemma2-2B. Results in Figure 1(b) summarize MNN-AECS normalized performance across all devices and datasets, reducing energy by 23% over its base engine MNN without slowdown, and by 39%–78% compared to all other engines.

In summary, the contributions of this paper include:

- We analyze the energy challenge of on-device LLM inference and identify that LLM decoding dominates energy consumption in Section 2.
- We propose MNN-AECS, a novel energy-efficient version of MNN that optimizes LLM decoding energy while preserving speed in Sections 3 and 4.
- We develop a cross-platform LLM energy evaluation testbed, and evaluate MNN-AECS against state-of-the-art on-device LLM inference engines in Section 5.
- We explore potential future extensions to mobile XPUs in Section 7.

## 2 BACKGROUND

### 2.1 On-Device LLM Inference

LLM model architecture consists of multiple sequential transformer blocks, each containing an attention module [44] and an FFN (Feed Forward Network) [44] or Gated-MLP

[7, 41, 42]. LLM inference pipeline can be divided into 2 phases: prefill phase and decode phase, showing different operational intensity.

**Compute-bound prefill phase.** Prefill phase takes a user prompt as input, with length $p$ after tokenization. Denote the hidden dimension of LLM as $d$, and then all the intermediate activations $H$ are of the shape $p \times d$, being matrices. For attention modules, these matrices $H$ are input and linearly affined to

$$Q, K, V = H \cdot W^{Q,K,V}, \ H \in \mathbb{R}^{p \times d}. \tag{1}$$

Then, self-attention inputs Q, K, V, and computes as

$$\mathbf{attention}(H) = \mathbf{softmax}(\mathbf{mask}(\frac{Q \cdot K^T}{\sqrt{d}})) \cdot V. \tag{2}$$

Subsequently, FFN/GatedMLP are formulated as

$$\mathbf{FFN}(H) = F_{act}(H \cdot W_{up} + b_{up})W_{down} + b_{down},$$
$$\mathbf{GatedMLP}(H) = (F_{act}(H \cdot W_{gate}) * (H \cdot W_{up}))W_{down}, \tag{3}$$

where the input matrix $H$ is multiplied by weight matrices.

Both prefill attention and FFN/GatedMLP are dominated by Matrix-Matrix Multiplication (GEMM). Hence, LLM prefill phase is compute-bound.

**Memory-bound decode phase.** During decoding, $K, V$ tensors of previous tokens are cached in a structure called KV cache for reuse, not needed to compute again [25]. Therefore, only 1 token needs to be processed during each decoding, causing the shape of intermediate activations $h$ to be $1 \times d$, a vector. For each attention module, vectors $q, k, v$ are computed from $h$, and $k, v$ are concatenated to previous KV matrices $K_C$ and $V_C$, formulated as

$$q, k, v = h \cdot W^{Q,K,V}, \ h \in \mathbb{R}^{1 \times d},$$
$$K_c := \begin{bmatrix} K_c \\ k \end{bmatrix}, V_c := \begin{bmatrix} V_c \\ v \end{bmatrix}. \tag{4}$$

Then, Self-attention is performed between vector $q$ and matrices $K_C, V_C$ as follows:

$$\mathbf{attention}(h) = \mathbf{softmax}(\frac{q \cdot K_c^T}{\sqrt{d}}) \cdot V_c. \tag{5}$$

After that, the intermediate state vector $h$ is input to:

$$\mathbf{FFN}(h) = F_{act}(h \cdot W_{up} + b_{up})W_{down} + b_{down},$$
$$\mathbf{GatedMLP}(h) = (F_{act}(h \cdot W_{gate}) * (h \cdot W_{up}))W_{down}, \tag{6}$$

where vector $h$ is multiplied by weight matrices.

Compared with prefill operations in Equations 1, 2, 3, matrix $H$ is replaced with vector $h$, and $Q$ is also replaced with vector $q$, so that both decode attention and FFN/Gated-MLP are dominated by Matrix-Vector Multiplication (GEMV), and thus memory-bound.

---

[1]All the Android CPU start with A or X stands for ARM Cortex-A or ARM Cortex-X. Some of the iPhone information is inaccessible and left blank.

**Table 2: Mobile devices used in our experiments and their hardware and OS specification. [1]**

| | device | DRAM | battery | SoC | CPU | GPU | NPU | OS | freq scheduler |
|---|---|---|---|---|---|---|---|---|---|
| **Android** | Huawei Mate 40 Pro | 8GB | 4400 mAh | Hisilicon Kirin 9000 | 1*A77(3.13GHz) +3*A77(2.54GHz) +4*A55(2.05GHz) | Mali-G78 (24 cores) | Da Vinci (2+1) | HarmonyOS 4.2.0 | schedutil |
| | Honor V30 Pro | 8GB | 4100 mAh | Hisilicon Kirin 990 | 2*A76(2.86GHz) +2*A76(2.36GHz) +4*A55(1.95GHz) | Mali-G76 (16 cores) | Da Vinci (2+1) | HarmonyOS 4.2.0 | schedutil |
| | Samsung Galaxy A56 | 8GB | 5000 mAh | Exynos 1580 | 1*A720(2.9GHz) +3*A720(2.6GHz) +4*A520(1.95GHz) | Xclipse 540 (2 cores, 256 shaders) | 1 core | One UI 7.0 | schedutil |
| | Meizu 21 | 12GB | 4800 mAh | Snapdragon 8 Gen 3 | 1*X4(3.3GHz) +3*A720(3.15GHz) +2*A720(2.96GHz) +2*A520(2.27 GHz) | Adreno 750 | Hexagon | FlymeOS 11.2.0 | walt |
| | Xiaomi 15 Pro | 16GB | 6100 mAh | Snapdragon 8 Elite | 2*Oryon(4.32GHz) +6*Oryon(3.53GHz) | Adreno 830 | Hexagon | HyperOS 2.0.23 | walt |
| **Apple** | iPhone 12 | 4GB | 2815 mAh | Apple A14 | 2+4 | 4 cores | 16 cores | iOS 18.3.2 | |
| | iPhone 15 | 8GB | 3349 mAh | Apple A16 | 2+4 | 5 cores | | iOS 18.4.1 | |

Moreover, given that normally on-device LLM serves only 1 user and thus processes 1 query at a time (i.e., batch size = 1), decoding can't be batched [9, 54] to remove its memory-boundedness. Therefore, we conclude that on-device LLM decoding is memory-bound.

## 2.2 Decoding Dominates Energy

Energy consumption is the product of time and power. Among the 2 LLM phases, decode phase has a longer execution time and similar power compared to prefill, dominating the overall LLM inference energy use. Figure 2(d) shows that decode energy is 16× to 26× more than prefill.
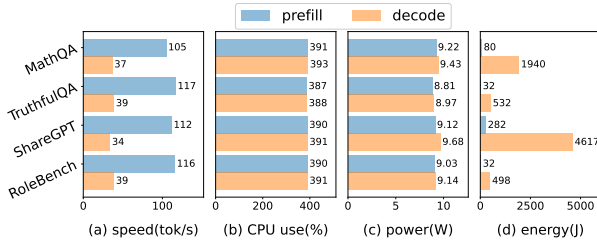


**Figure 2: Comparison of prefill and decode speed, CPU use, power, and energy of Qwen2.5-1.5B across 4 datasets on Xiaomi 15 Pro.**

**Decode phase takes much longer time,** in that decode speed is slower and decode length is longer. 1) Decode speed is 3× slower than prefill on Xiaomi 15 Pro CPU as an example in Figure 2(a). Due to memory-boundedness, all XPU cores share and wait for the same underlying memory bus on mobile devices under Unified Memory Architecture (UMA), making decoding hard to speed up at system level. By contrast, compute-bound prefill can be accelerated well by workload balance [45], layout transformation [36], and XPU [51, 52]. 2) Decode length is 3.5× longer than prefill length on conversational text generation such as ShareGPT [39] and RoleBench

[46], shown in Figure 3. This is because user queries (prefill) are usually shorter than generated responses (decode).
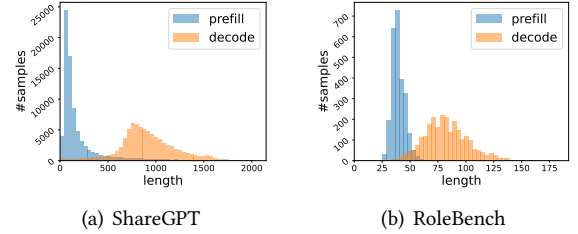


(a) ShareGPT  (b) RoleBench

**Figure 3: Prefill length and decode length distributional characteristic. (decode length ≈ 3.5× prefill length)**

**Decode phase power is comparable to prefill.** Figures 2(b) and 2(c) show that CPU utilization is similarly high during both prefill and decode phases, and so is the power, as power is positively correlated with CPU utilization. This phenomenon is due to existing engines' strategies of utilizing the same number of CPU cores in both phases.

Thus, we **focus on decode phase energy optimization**, especially reducing decode power, because decode time can hardly speed up due to memory-boundedness. Decode power can be reduced by decreasing CPU utilization during decode via leveraging mobile hardware and OS characteristics.

## 2.3 Mobile Hardware Characteristics

Multi-core XPU is a major hardware characteristic of smartphones from mainstream brands [21, 23, 34, 38, 50], which motivates our AECS design: adaptive selection across CPU multi-cores, introduced in Section 3. We also clarify the rationale of utilizing CPU for LLM decoding in our design.

**Multi-core architecture of mobile XPU.** Table 2 presents the devices for evaluation in this work. All of them feature multi-core XPU. For CPU, they consist of several heterogeneous multi-core clusters, e.g., ARM big.LITTLE core CPU

architecture [32]. Such heterogeneity aims at assigning only less power-consuming clusters for small bursty workloads to save energy. GPUs such as Adreno and Mali also consist of a series of cores, each containing hundreds of shaders [24, 26]. NPUs are specialized in neural network operations, especially Matrix Multiplication. Huawei Da Vinci series also feature heterogeneous cores for different kinds of neural network operations.

**Table 3: Core Configurability of mobile XPUs.**

|  | CPU | GPU | NPU |
|---|---|---|---|
| core frequency | ✗ | ✗ | ✗ |
| core number | ✔ | ✗ | ✗ |
| core affinity | Android: ✔ iOS: ✗ | ✗ | ✗ |

**Rationale of using CPU for LLM decoding.** 1) CPU is the most configurable backend outside OS without root access, shown in Table 3. It offers us the largest design space and opportunity. Thread number and even affinity (on Android) can be designated by LLM engine, which is impossible on GPU/NPU. These interfaces can be utilized to search for energy-efficient core selection. 2) CPU decoding speed is comparable to mobile GPU/NPU recently due to memory-boundedness, so that most existing engines majorly utilize CPU decoding [19, 33, 51, 52]. CPU's configurability and fast decoding speed motivates us to design energy-efficient CPU decoding.

## 2.4 Mobile OS Characteristics

Mobile OS is capable of manipulating hardware frequency based on process elapsed time for general-purpose energy saving, but not suitable for LLM energy optimization. Outside OS without root access, only CPU affinity or thread number setting are possible, yet still being our design opportunity.

**Inside OS exists general-purpose energy saving via frequency scheduling, but incapable of reducing LLM decoding energy.** By controlling the operating frequency of mobile XPUs, OS DVFS (Dynamic Voltage Frequency Scaling), specifically CPUFreq [47] and devfreq [2] subsystems, adopt governors such as schedutil [4] and walt [6] to schedule frequency and save energy when system estimated workload is low [14]. Workload is estimated based on tasks elapsed time. However, such workload-based scheduling is incapable of reducing LLM decode phase energy, because it can't distinguish decode from prefill and leverage the memory-boundedness of LLM decoding to scale down frequency and consequently save energy without speed loss.

**Outside OS exist LLM decoding energy saving opportunities through direct frequency setting with root access and through core selection system calls without root access.** With root access, it is possible to directly set the core frequency to reduce decode energy. However, obtaining root access outside the OS is impractical for commercial use. Without root access, we leverage core selection interfaces instead, including CPU affinity and thread number setting. On Android, CPU affinity (binding threads to designated cores) is available. On iOS, though affinity is unavailable, thread number can be controlled. With proper core selection via core binding and thread number setting, we can leave cores idle during memory-bound decoding, which can lead to significant energy reduction without intolerable speed loss, as detailed in Section 3.2. This core selection approach represents a key design opportunity for our work.

## 2.5 Low-Power Consideration in Existing On-Device LLM Inference Engines

Given the practical challenge of reducing energy consumption of LLM inference, on-device LLM inference engine as a middle-ware is motivated to leverage hardware and OS characteristics to encounter it. However, current on-device LLM frameworks lack sufficient focus on low-power optimization. For example, MNN [33] offers only a *Power_Low* mode that forces single-thread execution, which cannot meet the speed demands of LLM inference. Other engines, such as llama.cpp[19], executorch [8], MediaPipe [16], mllm [51], and PowerInfer-2 [52], focus on accelerating the prefill phase to reduce prefill energy but fail to address the more energy-consuming decode phase. This is because their optimizations do not apply to decoding, which is memory-bound. Therefore, an engine-level design is in need to enable energy-efficient LLM decoding.

## 3 DESIGN

## 3.1 Optimization Objective

Our major objective is to reduce LLM decoding energy, the bottleneck. Meanwhile, decode speed also matters for user experience. Hence, **we optimize energy under speed constraint**. We formulate the optimization objective as follows:

$$\min_{I \in S} E(I) = P(I) \cdot t(I), \ t(I) \propto \frac{1}{speed(I)}$$
$$s.t. \frac{speed(I)}{\max_J speed(J)} \geq 1 - \epsilon. \tag{7}$$

Core selection $I$ serves as the decision variable and is chosen from the search space $S$. The objective $E(I)$ is the empirical energy consumption associated with $I$. $E(I)$ is computed as the product of power $P(I)$ and time $t(I)$. The constraint term ensures that the decode speed $speed(I)$ of any feasible $I$ shall be at most $\epsilon$ slower than the fastest decode speed $\max_J speed(J)$. We empirically select $\epsilon = 8\%$, which is a slowdown not noticeable to smartphone users.

(a) **llama.cpp**: all cores are selected.



(b) **MNN**: cluster 2 and 1 are selected.



(c) **MNN-AECS**: 2 cores from cluster 1 is selected.

**Figure 4: CPU frequency curve of 3 core selections of LLM decoding on Huawei Mate 40 Pro.**

In summary, the optimal core selection $I^*$ of Equation 7 is the most energy efficient core selection with an $\epsilon$-suboptimal decode speed guarantee.

## 3.2 Search Space Design

The search space $S$ in Equation 7 is all CPU core binding plans on Android and all possible thread numbers on iOS. For example, Mate 40 Pro has 3 clusters, from big (cluster 0) to small (cluster 2), containing 1, 3, and 4 cores respectively. The design space of Mate 40 Pro spans all possible combinations of these 8 cores. In contrast, iPhone 12 has 6 cores, and its design space is defined by the thread number ranging from 1 to 6. We first conduct preliminary experiments to verify such search space design, demonstrating that the optimized core selection can indeed reduce CPU utilization and frequency, and consequently the power in decode phase.

**Rationale of Android search space design.** Figure 4 shows the preliminary experiments on Mate 40 Pro. llama.cpp selects all 8 cores available including efficient cores, resulting in peaked frequencies across all clusters. MNN selects all 4 cores in cluster 2 and 1, and the frequency of these 2 selected clusters are boosted to peak. By contrast, the optimal solution (discovered by MNN-AECS) only selects 2 cores form cluster 1, successfully reducing the average frequency of idle clusters by 50%. The numerical results presented in Table 4 indicate that the frequency reduction contributes to 29% and 65% energy reduction over MNN and llama.cpp respectively.

Besides, our optimal solution only slows down 5% compared to MNN, because of memory-boundedness.

The results indicate that a proper CPU binding can leave CPU cores or even clusters idle, reducing frequency and power consumption of them.

**Table 4: Average CPU frequency, speed, power, and energy of 3 engines on Mate 40 Pro, Qwen2.5-1.5B.**

| engine | CPU frequency (GHz) | | | speed | power | energy |
|---|---|---|---|---|---|---|
| | cluster 2 | cluster 1 | cluster 0 | (tok/s) | (W) | (mJ/tok) |
| llama.cpp | 3.08 | 2.50 | 2.01 | 10.2 | 8.8 | 860 |
| MNN | 3.05 | 2.47 | 1.04 | **21.7** | 8.7 | 403 |
| MNN-AECS | **1.61** | **2.33** | **0.96** | 20.6 | **6.2** | **300** |

**Rationale of iOS search space design.** Table 5 shows the preliminary experiments on iPhone 12. Our optimal solution, selecting only 1 thread, is 42% and 48% more energy-saving and 14% and 100% faster than MNN and llama.cpp, by saving over 2× more CPU resources. Hence, a properly set thread number can also leaves cores or even clusters idle.

**Table 5: Average CPU frequency, speed, power, and energy of 3 engines on iPhone 12, Qwen2.5-1.5B.**

| engine | thread | CPU use | speed (tok/s) | relative power |
|---|---|---|---|---|
| llama.cpp | 2 | 197% | 15.3 | 955 |
| MNN | 4 | 230% | 27.6 | 871 |
| MNN-AECS | **1** | **97%** | **31.5** | **506** |

These preliminary experiments demonstrate how the optimal core selection saves considerable energy without losing much speed or even faster than existing works, as long as such solution can be found by our searching algorithm.

## 3.3 Adaptive Energy-Centric Core Selection (AECS)

We present our algorithm AECS: search for a core selection that optimizes Equation 7 based on CPU information and the once-and-for-all on-device searching.

Algorithm 1 shows the workflow of AECS. It has 2 search stages. **Stage 1** searches for the fastest core selection $\tilde{I}$ so that following searches can compare with speed of it $speed(\tilde{I})$ check the speed constraint. It also serves as the initializer in stage 2 search. **Stage 2** introduces a heuristic power estimation function $h(I)$ to guide the search. Based on stage 1 result and the heuristic function, stage 2 generates a heuristically pruned candidate tree $S_h(\tilde{I})$ as the candidate set. Then each candidate core selection on the tree is tested and profiled. Finally, the best candidate core selection $I^*$ minimizing the heuristically modified energy function $E_h(\cdot)$ and satisfying the speed constraint is returned.

**Stage 1: searching for the fastest core selection.** We first search for the core selection $\tilde{I}$ that achieves the fastest

---

**Algorithm 1:** AECS

---
   **Input**   :CPU information.
   **Output**:The optimized core selection $I^*$.
1  Read and analyze CPU info.
2  // **Stage 1:**
3  Search for the fastest core selection $\tilde{I}$.
4  // **Stage 2:**
5  Generate candidate set $T_h(\tilde{I})$ based on $\tilde{I}$.
6  **foreach** $I$ **in** $T_h(\tilde{I})$ **do**
7     |  run $I$ and measure $E(I)$ and $speed(I)$, record them.
8     |  **if** $speed(I) \geq speed(\tilde{I}) \cdot (1 - \epsilon)$ **then**
9     |    |  // $I$ violates speed constraint, pop it
10    |    |  $T_h(\tilde{I})$.pop(I)
11 $I^* \longleftarrow \arg\min_{I \in T_h(\tilde{I})} E_h(I)$
12 **return** $I^*$

---

decoding, so that we can check and verify the speed constraint in the following stage 2 energy optimization search.
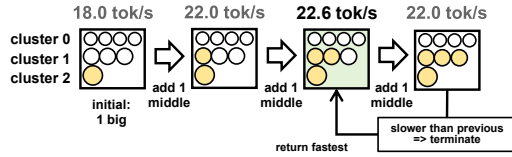


**Figure 5: Stage 1 process on Mate 40 Pro: Starting from 1 big core, middle cores are added sequentially. Stage 1 terminates with selecting 1 big and 2 middle cores.**

The input of stage 1 is the CPU information. Specifically, the max frequency of each CPU cluster is used to distinguish big cores from small ones. The output of this stage is the fastest core selection $\tilde{I}$. The search process starts from 1 prime core, and adds CPU cores greedily from big to small one at a time, meanwhile testing the speed for each core selection plan. Stage 1 terminates if adding 1 more core doesn't speed up any more, or no prime cores and performance cores are left. Efficient cores are not considered. The example search process of stage 1 on Mate 40 Pro is shown in Figure 5.

**Stage 2: searching for the optimal core selection.** In stage 2, we search for the optimal core selection $I^*$ of Equation 7. To better carry out our search, we reformulate the original Equation 7 into

$$I^* = \arg\min_{I \in S_h(\tilde{I})} E_h(I),$$
$$s.t., \frac{speed(I)}{\max_J speed(J)} \geq 1 - \epsilon. \qquad (8)$$

Compared with Equation 7, 1) we change the objective from empirical energy $E(I)$ to heuristic objective $E_h(I)$ by introducing a power heuristic function $h$ to guide the search and

enhance searching robustness, and 2) the search space is scaled down from all core selections $S$ to a candidate subset $S_h(\tilde{I})$ filtered by heuristic function $h$. In what follows, we further introduce the detailed design of power heuristic function $h(I)$, modified objective function $E_h(I)$, and candidate set $S_h(I)$.

*Power heuristic function $h(I)$.* The power heuristic function is modeled based on 4 major hardware and OS characteristics, including power-frequency relationship, CPU type modeling, idleness modeling, and CPUFreq governor modeling.

- Power-frequency relationship is modeled quadratically along with a static power $Ps$: $h(I) \propto f^2 + Ps$, because recent work [14] suggests power increases super-linearly corresponding to frequency on modern cores.
- CPU type also determines CPU power. Therefore, each cluster of different CPU type is assigned a distinct scaling factor $a_i$, so that $h(I) \propto \sum_{i=0}^{n-1} a_i \cdot f_i^2 + Ps$, where $n$ is the number of CPU clusters.
- Idle cores inside a cluster is assigned a reduced factor $b < 1$ based on CPU idle state behaviors [1]. Denote the number of cores in cluster $i$ as $|C_i|$ and selected ones as $|I_i|$, the idleness reduced factor of cluster $i$ is summed up to $(|I_i| + (|C_i| - |I_i|)b)$. Up to now, $h(I)$ is formulated as $h(I) \propto \sum_{i=0}^{n-1} a_i(|I_i| + (|C_i| - |I_i|)b)f_i^2 + Ps$.
- CPUFreq governor behavior is modeled to estimate the operating frequency of cluster $i$ $f_i$. We check the source code[2] in latest Android kernel [5] and discover that the estimated workload is scaled by the capacity factor $s_I = \frac{selected\ biggest\ capacity}{biggest\ capacity}$, and thus the assigned frequency is usually the max frequency scaled by this factor, $f_i = f_{max,i} \cdot s_I$.

Combining these features, the final power heuristic function is modeled as

$$h(I) = \sum_{i=0}^{n-1} a_i(|I_i| + (|C_i| - |I_i|)b)(f_{max,i} \cdot s_I)^2 + Ps. \quad (9)$$

*Objective function:* $E_h(I) = (1 - \alpha)E(I) + \alpha h(I)t(I)$. The empirically measured energy $E(I)$ suffers from intolerable system fluctuation in energy and speed, as high as 5%, which can easily skew the search results. To address this issue, we modify the original objective to be a weighted average of the observed energy $E(I)$ and heuristic estimation $h(I) \cdot t(I)$ to enhance the robustness of the search for the optimal result. Ablation studies in Section 5.5 verifies this design.

*Generation of candidate set (heuristic tree):* $S_h(\tilde{I})$. To generate competitive candidates, we start from the stage 1 result $\tilde{I}$ as the root and grow a candidate tree. The candidate set $S_h(\tilde{I})$ is defined as all the nodes on the tree. Heuristic transformations are defined to generate succeeding nodes from

---

[2]function named *scale_load_to_cpu* in file kernel/sched/sched.h

parents, so that succeeding nodes are more energy-efficient and comparable in speed. On Android, 4 transformations are designed, including a) removing 1 smallest core, b) removing 2 smallest cores, c) changing 1 bigger core to a smaller one in another selected cluster, and d) changing a selected cluster of bigger cores to an unselected cluster of smaller cores. Besides, efficient cores such as Cortex-A5 series are not ignored from selection. The tree depth is limited less than 2, and transformations a) and b) are only allowed in level 1 to avoid plan duplication. The tree size is therefore capped small. In Figure 6, transformations a), b), c) are applicable to the root, which has 1 big core and 2 middle, generating 3 nodes at level 1. The first node at level 1 can further generate a node in level 2 with c), resulting in a tree with 5 candidates.
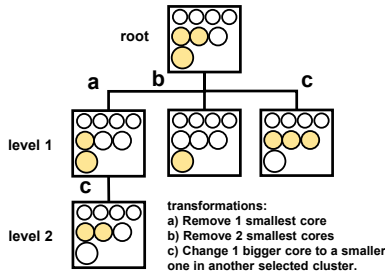


**Figure 6: Stage 2 heuristic tree on Mate 40 Pro: transformations a), b), c) are used to generate nodes.**

On iOS, the generation process is much easier. Only reducing 1 thread is valid to generate candidate node.

With the heuristic candidate tree, candidate set is 5-10× smaller while preserving the optimality of the search. A detailed analysis is presented in Section 5.5.

To effectively solve the optimization defined in Equation 8, stage 2 searching is introduced, presented in lines 5-11 of Algorithm 1. Stage 2 inputs CPU information to compute the power heuristic function, along with the stage 1 result, and outputs the energy-optimal core selection, which is also the final output of entire AECS searching. First, the candidate set $S_h(\tilde{I})$ is generated. Then, stage 2 searching proceeds over the candidate set $S_h(\tilde{I})$ and measures the speed and energy of each candidate. If a candidate violates the speed constraint, it's immediately removed from the candidate set. After all candidates are tested, the algorithm outputs the candidate with minimum value on the modified energy objective.

## 4 IMPLEMENTATION

We integrate AECS into the on-device inference engine MNN and implement an energy profiling module to form our version, MNN-AECS. Core functionality is implemented in approximately 6,000 lines of C++, Java, Swift, Objective C, and Python code.

**Table 6: Model compatibility of different engines on Android. ✔ denotes support, oom denotes out of memory on some devices, and ✗ denotes not supported.**

|  | Qwen2.5-1.5B | Qwen2.5-3B | Llama3.2-1B | Llama3.2-3B | Gemma2-2B |
|---|---|---|---|---|---|
| executorch | ✗ | ✗ | ✔ | ✔ | ✗ |
| llama.cpp | ✔ | ✔ | ✔ | ✔ | ✔ |
| MediaPipe | ✗ | ✗ | ✗ | ✗ | ✔ |
| mllm | ✔ | oom | ✔ | oom | ✗ |
| MNN | ✔ | ✔ | ✔ | ✔ | ✔ |

### 4.1 Engine Integration (MNN-AECS)

We choose MNN as our base engine because of its superior model compatibility (Table 6) and its fast CPU LLM decoding speed. MNN decodes 1.1× to 3× faster than nearly all other engines [8, 16, 19, 51]. This speed advantage is attributed to its contiguous KV cache and weight layout.

The integration involves outside-engine AECS search and inside-engine thread pool and memory pool modifications.

**Outside-engine AECS search.** AECS searching in Algorithm 1 is implemented outside MNN. Algorithm inputs CPU information including CPU capacity, max frequency and cluster information, and CPU type from */sys/devices/system/cpu/* and */proc/cpuinfo*. During searching, AECS calls the newly-added core selection interface of modified MNN and profiles energy through energy profiling module. After searching, the optimized core selection is recorded and used in following MNN LLM decoding.

**Inside-engine modifications.** The primary requirement for MNN modification is to enable rapid switching between different core selections. This functionality is crucial for AECS searching, which necessitates efficient testing of various core selections. Besides, our design employs distinct core selections for prefill and decode to prevent interference, where rapid switching is essential for inference speed. However, the original MNN doesn't support changing core selection. To address this challenge, we modify MNN thread pool and memory pool.

**Thread pool modifications.** We utilize Android system call *__NR_sched_setaffinity* to set and reset CPU bindings, while exposing an additional thread number reset interface on iOS. Thanks to the modified interfaces, core selection can be changed with a simple function call.

**Memory pool modifications.** Originally, MNN's KV cache memory buffer layout depends on thread number. We eliminate such dependency in our modified memory pool, so that KV cache can be shared by prefill and decode with different thread number.

### 4.2 Energy Profiling Module

**Requirements.** Energy profiling module needs to satisfy 3 key requirements: 1) *Precision*: Energy profiling is required to be highly precise. 2) *Interoperability with MNN-AECS*: Energy

data from OS interface, regardless of the interface language, shall be able to pass to our C++ engine and Python evaluation pipeline. 3) *Without developer mode*: Developer mode is not permitted for a commercial app. However, it can be used during evaluation, as evaluation is not part of the LLM serving app.

**Challenges.** For Android, energy profiling is only available from Java via *android.os.BatteryManager*. The OS interface updates every 250 ms, making tests shorter than 500 ms imprecise. Besides, the interface is in Java, but our engine is implemented in C++.

For iOS, energy can only be monitored on PC through tunneld over WLAN. Furthermore, iOS energy profiling is impossible without developer mode. Though iOS Sysdiagnose [43] and Xcode energy gauge[3] can both measure energy performance of an iOS app, Sysdiagnose updates energy every 20s and suffers from $\geq 10\%$ fluctuation, failing to meet our precision requirement. In contrast, Xcode energy gauge is more precise but necessitates developer mode.
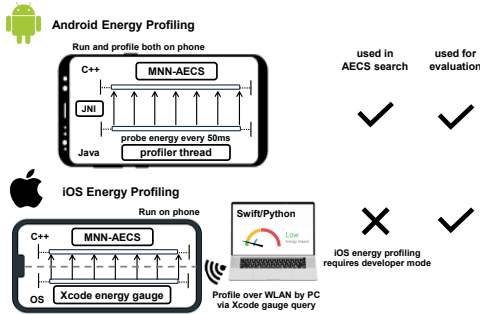


**Figure 7: MNN-AECS energy profiling module.**

**Implementations.** For Android, we pool current and voltage in a separate profiling thread to calculate energy, which is then passed to MNN-AECS over JNI (Java-Native Interface), illustrated in Figure 7. To make the energy profiling more precise, we probe *BatteryManager* every 50 ms so that OS energy updates can be monitored in time. For AECS search, we decode 50 tokens, so that the decode time exceeds the OS interface update interval, ensuring the precision.

For iOS, we opt for Xcode energy gauge to get more precise power data, despite it providing only relative power without exact units. Illustrated in Figure 7, energy of MNN-AECS is monitored by Xcode energy gauge on PC over WLAN though tunneld implemented in *pymobiledevice3* [3]. AECS search aborts energy profiling on iOS, using only the heuristics for optimization, due to the developer mode requirement of Xcode energy gauge. The iOS energy profiling module is utilized exclusively in evaluation experiments, standalone from MNN-AECS.

---

[3]*com.apple.xcode.debug-gauge-data-providers.Energy*

## 5 EVALUATION

### 5.1 Experimental Setup

We implement a novel cross-platform testbed to support our extensive experiments on both Android and iOS devices, comparing MNN-AECS energy and speed performance with 5 state-of-the-art on-device LLM engines across 5 popular LLMs over 4 real-world datasets.

**Devices.** To ensure extensive coverage of common-seen mobile devices, 5 different Android phones and 2 iOS phones are used in experiments, listed in Table 2.

**Baseline engines and LLMs.** On Android, experiments are conducted between our MNN-AECS and its base engine MNN, along with other 4 engines: llama.cpp [19], mllm [51], MediaPipe [16], and executorch [8], over Llama3.2-1B/3B [7], Qwen2.5-1.5B/3B [42], and Gemma2-2B [40]. If an engine doesn't support a model, shown in Table 6, we skip the model test on that engine. Qwen and Llama models are 4-bit quantized and Gemma is 8-bit quantized for fair comparison across engines. On iOS, we compare across ours, MNN, and llama.cpp, over Llama3.2-1B and Qwen2.5-1.5B.

**Testing benchmarks.** 2 types of tests are adopted in our analysis. 1) **Fixed length experiments:** In order to evaluate MNN-AECS performance under different prefill prompt lengths and decode generation lengths, we control LLM inputs and outputs, so that prefill lengths are in $\{64, 256, 1024\}$, and decode lengths in $\{128, 256, 512\}$. 2) **Dataset experiments:** To measure the performance of models on practical daily tasks, we conduct dataset experiments on 4 different mobile LLM use cases: math problem solving (MathQA [11]), open domain QA (TruthfulQA [30]), multi-turn conversation (ShareGPT [39]), and role play (RoleBench [46]). A representative subset randomly sampled from each dataset is used in test due to computation and battery constraints on phone.

**Metrics.** 4 major metrics are compared: energy (mJ/token), battery ($\mu$Ah/token), decode speed (token/s), and CPU use.

**Testing conditions.** To simulate real-world use condition, we conduct all tests **unplugged**. Besides, alike previous works [14], we cool down the phones to under $40°C$ before each test to prevent high temperature throttling. For those phones where temperature information isn't available, they sleep for about a minute to cool down. The battery charge is also kept above 50% to prevent low-battery throttling.

### 5.2 Tuned Results

The tuned core selections of MNN-AECS on 7 devices are presented in Table 7. These once-and-for-all tuned results are subsequently used across all our experiments.

**Lower CPU utilization.** Table 8 illustrates the significant CPU utilization reduction of MNN-AECS over the baselines. Other engines utilize $4 \sim 8$ cores. By contrast, MNN-AECS uses only $\leq 2$ cores across all 7 devices, reducing $50 - 75\%$

CPU use on almost all devices remarkably. Only llama.cpp on iPhone 15 has comparably low CPU utilization.

**Table 7: Tuned core selections of AECS on 7 devices. (Refer to Table 2 for complete device hardware info.)**

| | device | tuned core selections |
|---|---|---|
| **Android** | Mate 40 Pro | 2*A77(2.54GHz) |
| | Honor V30 Pro | 2*A76(2.36GHz) |
| | Galaxy A56 | 2*A720(2.6GHz) |
| | Meizu 21 | 1*X4(3.3GHz) +1*A720(2.96GHz) |
| | Xiaomi 15 Pro | 2*Oryon(4.32GHz) |
| **Apple** | iPhone 12 | 1 thread |
| | iPhone 15 | 2 thread |

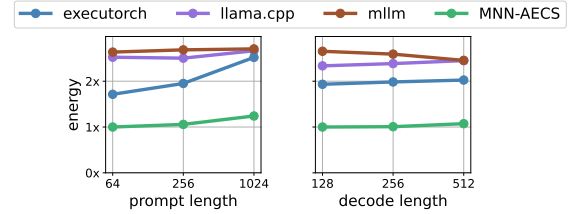**Table 8: CPU core utilization in decode phase. MNN-AECS reduces core utilization by** $50 - 75\%$**.**

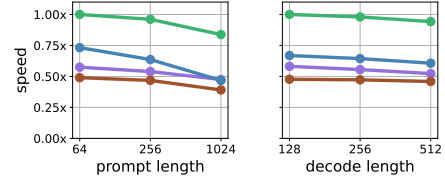| engine | selected core number | | | | | thread number | |
|---|---|---|---|---|---|---|---|
| | Mate 40 Pro | V30 Pro | Galaxy A56 | Meizu 21 | Xiaomi 15 Pro | iPhone 12 | iPhone 15 |
| executorch | 8 | 8 | 8 | 8 | 8 | - | - |
| llama.cpp | 8 | 8 | 8 | 8 | 8 | 2 | **2** |
| MediaPipe | 4 | 4 | 4 | - | 4 | - | - |
| mllm | 4 | 4 | 4 | 4 | 4 | - | - |
| MNN | 4 | 4 | 4 | 4 | 4 | 4 | 4 |
| MNN-AECS | **2** | **2** | **2** | **2** | **2** | **1** | **2** |

## 5.3 Fixed Length Experiments

To analyze energy and decode speed of MNN-AECS under different prompt and decode generation lengths, we normalize the geometric mean across all devices and models and compare MNN-AECS with executorch, llama.cpp, and mllm in Figure 8. Also, more detailed analyses on each device are provided between MNN-AECS and original MNN in Figures 9 and 10.

Illustrated in Figure 8, for energy, MNN-AECS achieves 36%-62%, 49%-60%, 52%-54% energy reduction under prompt length = 64, 256, 1024, and 49%-61%, 50%-61%, 50%-56% energy reduction under decode length = 128, 256, 512. The results indicate that MNN-AECS energy reduction is consistent under different LLM input and output. For decode speed, MNN-AECS achieves 33%-104%, 58%-104%, 69%-113% speedup under prompt length = 64, 256, 1024, and 43%-110%, 45%-107%, 45%-104% speedup under decode length = 128, 256, 512. The results indicate that MNN-AECS not only reduces the energy but also excels in decode speed.

Figures 9 and 10 compare MNN-AECS with its base engine MNN on 7 devices. Energy reduction of MNN-AECS is more significant for shorter prompts, where decode energy makes up for a higher proportion of total energy. Decode length has less impact on energy by contrast. MNN-AECS reduces 18% to 42% energy except for Meizu 21 (10%), whose OS does not scale down the CPU cluster frequency though idle making our strategy less effective.
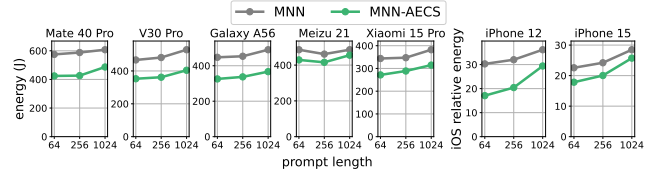


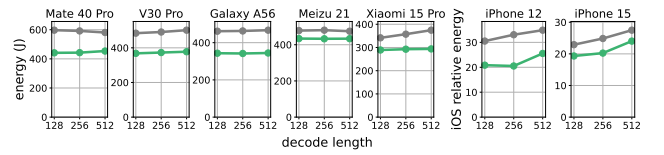(a) Impact of different prompt and decode lengths on energy.



(b) Impact of different prompt and decode lengths on decode speed.

**Figure 8: Normalized energy and decode speed comparison between MNN-AECS and baselines across 7 devices and 5 models.**



(a) Impact of prompt length on energy.



(b) Impact of decode length on energy.

**Figure 9: MNN-AECS energy reduction over MNN under different (a) prompt and (b) decode length.**

For decode speed, MNN-AECS only slows down for 6% on V30 Pro and speeds up 0 to 20% on all the rest, thanks to MNN-AECS's core affinity and less cores which contributes to less congestion and conflicts on bus. Prompt length also impacts more on decode speed, because longer prompts makes up large KV matrices for decode.

## 5.4 Dataset Experiments

Tables 9 and 10 show the dataset experiments results on 5 Android and 2 iOS devices, averaged over 4 real-world datasets specified in Section 5.1. We consumes the lowest energy and battery across all devices and models compared to our baselines. MNN-AECS is also the fastest solution in most cases, where slowdown is less than 7% in the rest.

**Table 9: Dataset experiments results on 5 Android devices averaged over 4 datasets.**

| device | engine | Qwen2.5-1.5B | | | Qwen2.5-3B | | | Llama3.2-1B | | | Llama3.2-3B | | | Gemma2-2B | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | speed | battery | energy | speed | battery | energy | speed | battery | energy | speed | battery | energy | speed | battery | energy |
| Mate 40 Pro | executorch | - | - | - | - | - | - | 14.4 | 34.0 | 505.6 | 5.6 | 87.6 | 1246.6 | - | - | - |
| | llama.cpp | 10.1 | 53.0 | 772.2 | 6.3 | 77.7 | 1143.0 | 15.8 | 35.7 | 509.6 | 5.8 | 85.0 | 1223.7 | - | - | - |
| | MediaPipe | - | - | - | - | - | - | - | - | - | - | - | - | 3.1 | 372.5 | 5437.9 |
| | mllm | - | - | - | - | - | - | 7.3 | 73.0 | 1080.2 | - | - | - | - | - | - |
| | MNN | **20.7** | 26.8 | 389.4 | 10.5 | 50.0 | 737.7 | **26.2** | 21.2 | 314.6 | **9.7** | 54.1 | 794.9 | **12.3** | 44.1 | 662.2 |
| | MNN-AECS | 20.3 | **20.2** | **298.2** | **10.5** | **39.3** | **583.6** | 25.8 | **16.4** | **236.7** | 9.4 | **45.5** | **650.3** | 12.2 | **36.4** | **523.1** |
| V30 Pro | executorch | - | - | - | - | - | - | 16.9 | 28.6 | 410.8 | 7.2 | 65.6 | 955.1 | - | - | - |
| | llama.cpp | 9.9 | 50.3 | 728.1 | 6.3 | 73.9 | 1088.6 | 15.7 | 33.3 | 467.0 | 6.0 | 82.0 | 1166.0 | - | - | - |
| | MediaPipe | - | - | - | - | - | - | - | - | - | - | - | - | 4.7 | 96.7 | 1418.7 |
| | mllm | - | - | - | - | - | - | 7.7 | 64.3 | 926.8 | - | - | - | - | - | - |
| | MNN | **23.0** | 21.0 | 313.0 | **11.7** | 43.0 | 619.7 | **29.0** | 17.3 | 245.7 | **10.9** | 48.3 | 675.7 | **14.0** | 38.1 | 541.7 |
| | MNN-AECS | 20.9 | **17.2** | **250.3** | 11.1 | **32.8** | **479.9** | 26.4 | **13.4** | **199.8** | 10.3 | **36.7** | **525.2** | 13.3 | **27.7** | **400.2** |
| Galaxy A56 | executorch | - | - | - | - | - | - | 17.6 | 24.2 | 356.8 | 7.3 | 58.0 | 846.9 | - | - | - |
| | llama.cpp | 10.8 | 44.8 | 661.4 | 7.2 | 62.1 | 920.5 | 17.5 | 29.5 | 422.4 | **9.1** | 60.6 | 871.7 | - | - | - |
| | MediaPipe | - | - | - | - | - | - | - | - | - | - | - | - | 4.3 | 59.4 | 916.4 |
| | mllm | - | - | - | - | - | - | 6.9 | 45.7 | 669.0 | - | - | - | - | - | - |
| | MNN | **18.3** | 18.9 | 290.8 | 9.0 | 33.2 | 501.1 | 21.8 | 16.1 | 235.8 | 8.4 | 36.6 | 560.9 | 10.8 | 32.7 | 492.4 |
| | MNN-AECS | 18.0 | **16.1** | **237.2** | **9.2** | **28.3** | **429.5** | **22.2** | **12.0** | **185.0** | 8.4 | **31.5** | **469.5** | **11.0** | **25.3** | **375.0** |
| Meizu 21 | executorch | - | - | - | - | - | - | 22.4 | 25.4 | 368.8 | 9.5 | 57.4 | 841.6 | - | - | - |
| | llama.cpp | 11.3 | 56.9 | 853.4 | 7.8 | 81.3 | 1179.0 | 18.2 | 36.6 | 525.3 | 7.3 | 83.2 | 1226.8 | - | - | - |
| | mllm | - | - | - | - | - | - | 9.6 | 44.1 | 674.4 | - | - | - | - | - | - |
| | MNN | 21.8 | 19.2 | 283.6 | 11.6 | 35.9 | 548.6 | 26.8 | 14.8 | 228.6 | 10.3 | 39.9 | 585.2 | 11.1 | 37.2 | 552.0 |
| | MNN-AECS | **24.4** | **18.1** | **262.3** | **12.9** | **31.4** | **472.6** | **30.1** | **13.9** | **213.0** | **11.4** | **35.9** | **513.2** | **12.6** | **35.0** | **509.0** |
| Xiaomi 15 Pro | executorch | - | - | - | - | - | - | 39.8 | 13.8 | 195.3 | 21.1 | 33.1 | 480.3 | - | - | - |
| | llama.cpp | 33.5 | 21.2 | 318.7 | 14.1 | 50.2 | 756.8 | 41.5 | 18.1 | 267.8 | 15.1 | 50.1 | 742.2 | - | - | - |
| | MediaPipe | - | - | - | - | - | - | - | - | - | - | - | - | 12.6 | 46.1 | 692.5 |
| | mllm | - | - | - | - | - | - | 15.2 | 38.1 | 558.6 | - | - | - | - | - | - |
| | MNN | 37.1 | 14.2 | 213.9 | 20.1 | 28.8 | 416.4 | 47.2 | 9.8 | 143.5 | 19.0 | 29.7 | 443.0 | 23.4 | 25.4 | 359.4 |
| | MNN-AECS | **45.6** | **9.9** | **150.4** | **23.3** | **22.1** | **322.3** | **59.0** | **7.9** | **114.3** | **21.3** | **22.9** | **355.6** | **28.7** | **20.8** | **292.6** |



(a) Impact of prompt length on decode speed.



(b) Impact of decode length on decode speed.

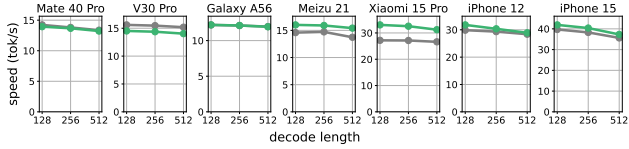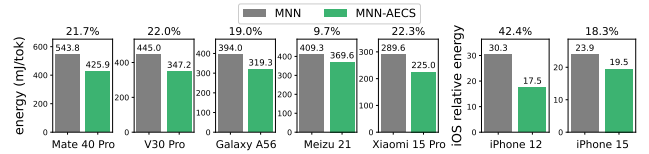**Figure 10: MNN-AECS decode speed compared to MNN under different (a) prompt and (b) decode length.**

Figure 12 compares MNN-AECS with executorch, llama.cpp, MediaPipe, and mllm on each devices, where energy and decode speed are geometric mean over all models and datasets tested on the device. MNN-AECS consistently reduces energy by 26% to 92% across all devices, and also speeds up by 12% to 363%.
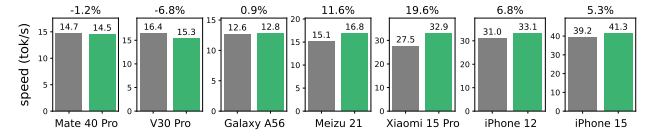
To boil down our energy reduction and speedup, we again compare MNN-AECS to original MNN. In Figure 11, we achieve 42% energy reduction on iPhone 12, about 20% reduction on 5 other phones, and 10% reduction on Meizu 21. Only 1% and 7% slowdown is witnessed on 2 Huawei phones, and we even speed up 1% to 20% on the rest 5. On the other

**Table 10: Dataset experiments results on 2 iPhones averaged over 4 datasets.**

| device | engine | Qwen2.5-1.5B | | Llama3.2-1B | |
|---|---|---|---|---|---|
| | | speed | energy | speed | energy |
| iPhone 12 | llama.cpp | 13.4 | 74.3 | 16.0 | 62.4 |
| | MNN | 27.9 | 34.0 | 34.5 | 27.1 |
| | ours | **29.3** | **19.6** | **37.3** | **15.5** |
| iPhone 15 | llama.cpp | 17.7 | 55.1 | 20.5 | 48.8 |
| | MNN | 33.4 | 28.6 | 46.1 | 19.9 |
| | ours | **35.5** | **23.4** | **48.1** | **16.3** |



(a) MNN-AECS energy reduction over MNN is presented atop subfigures (10%-42%).



(b) MNN-AECS speedup over MNN is presented atop subfigures (-7%-20%).

**Figure 11: (a) Energy and (b) decode speed comparison between MNN-AECS and MNN across 7 devices in dataset experiments.**

hand, Figure 13 shows that MNN-AECS energy reduction and speedup over MNN are consistent across 5 models.
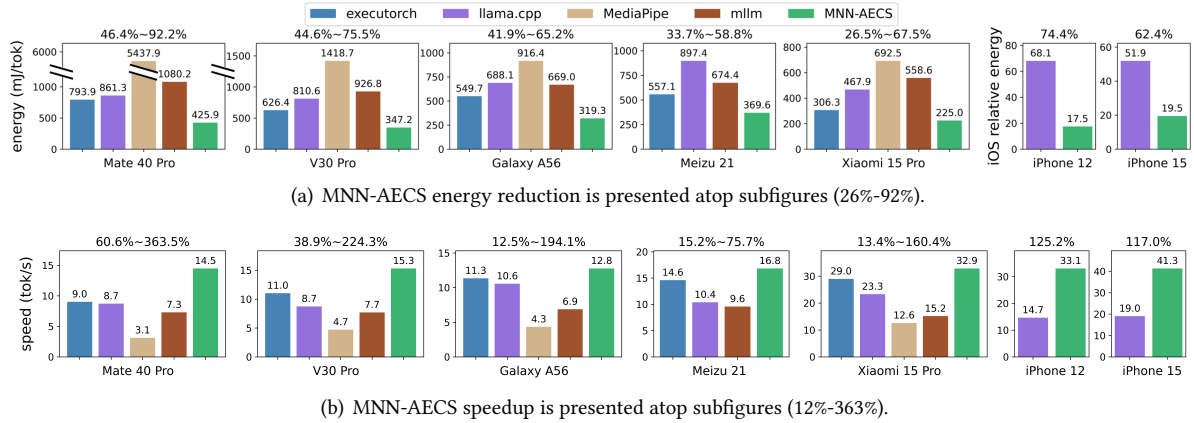
(a) MNN-AECS energy reduction is presented atop subfigures (26%-92%).



(b) MNN-AECS speedup is presented atop subfigures (12%-363%).

**Figure 12: Energy and decode speed in dataset experiments, geometrically averaged across 5 LLMs on 7 devices.**



(a) MNN-AECS relative energy reduction over MNN on 5 LLMs.



(b) MNN-AECS relative speedup over MNN on 5 LLMs.

**Figure 13: (a) Energy and (b) decode speed comparison between MNN-AECS and MNN across 5 models in dataset experiments.**

## 5.5 Ablation Studies

We compare AECS to exhaustive traversal of all possible core selections to demonstrate optimality and efficiency. Besides, we also analyze how our heuristic-averaged energy objective enhances robustness. Results are shown in Table 11.

**Optimality and Efficiency: compare to exhaustive traversal.** We conduct exhaustive traversal on all test devices and observe that the searched results are exactly the same as AECS results, which means the optimality rate of AECS is 100%. The search space of exhaustive traversal scales up to 20-71, taking up 10-20 minutes of foreground execution, which is intolerable for users. AECS successfully removes uncompetitive candidates so that the search time is reduced to 1-2 minutes. The results indicate that AECS has the same optimality rate as exhaustive search while being 10× faster.

**Robustness: compare to AECS without power heuristic.** In AECS design, the energy objective is a weighted average of measured energy and the power heuristic. Table 11

**Table 11: Ablation studies: compare AECS to exhaustive traversal and AECS without heuristic.**

|  | exhaustive (optimal) | AECS w/o heuristic | AECS |
|---|---|---|---|
| search space | 20-71 | 4-9 | 4-9 |
| serach time (min) | 10-20 | 1-2 | 1-2 |
| optimality rate | 100% | 60%-90% | 100% |

shows that removing the power heuristic leads to fluctuation in searching and consequently lower optimality rate.

## 6 RELATED WORKS

### 6.1 On-device LLM Inference Engine

A series of well-known mobile LLM inference engines are available these days, including MNN [33], llama.cpp [19], executorch [8], mllm [51], MediaPipe [16] (based on LiteRT [15]), and MLC-LLM [35] (based on TVM [13]). Among these, MNN and llama.cpp support the widest range of models. MNN achieves leading inference speed on both mobile CPUs and GPUs with its optimized KV cache and weights layout that ensures contiguous memory reads. llama.cpp is highly extensible and thus also the most popular. executorch integrates SpinQuant in its engine to improve accuracy, though only supporting Llama series. mllm leverages NPU GEMM kernels to accelerate prefill. MediaPipe and MLC-LLM achieve leading GPU prefill performance through their optimized GPU GEMM kernels. Additionally, PowerInfer-2 [52] , though not accessible yet, also leverages NPU and activation sparsity to accelerate and accommodate very large LLMs on mobile devices.

These engines primarily focus on prefill acceleration via optimized GEMM kernels, accuracy improvements through quantization, and memory reduction with sparsity. However, none of them address the LLM decoding energy challenge.

## 6.2 Mobile OS DVFS for DL Tasks

Several mobile DVFS algorithms are designed for DL tasks, such as Asymo [45] and GearDVFS [28]. These methods directly manipulate CPU frequency, which requires root access or customized Android kernel. While they are valuable for OS-level design, they are impractical for engine-level design in our context. Moreover, they specifically target at traditional CNN-based DL tasks, which have significantly different computational characteristics compared to LLM.

## 6.3 LLM Quantization and Sparsity

Quantization and sparsity are 2 major algorithmic methods for LLM speedup and memory reduction. In practice, quantization algorithm, such as k-quant [19], SmoothQuant [48], SpinQuant [31], AWQ [29], are commonly adopted to quantize models to 4/8-bit with tolerable accuracy loss to fit the model into the limited memory on mobile devices. On the other hand, sparsity in activation [10, 52], KV cache [49], and MoE structure [53] are also promising directions for reducing memory and computation of LLM inference. Both LLM quantization and sparsity are orthogonal to our core selection system design of core selection.

## 7 DISCUSSION ON EXTENSIONS TO MOBILE XPU

During memory-bound LLM decoding, XPUs frequently wait for the bus, resulting in limited performance gains from increasing parallelism due to the shared memory bus bottleneck. MNN-AECS exploits such characteristics on CPU to save energy under speed guarantee. Such abstraction and system design can potentially be extended to mobile GPU/NPU.

Latest Snapdragon Adreno GPU automatically reduces shader processors (SP) frequency and mainly leveraging texture processors (TP) during image read [24], which is an idea similar to ours and can be utilized to save energy during GPU LLM decoding. Our preliminary trials to utilize such features in MNN successfully reduce GPU main frequency on Xiaomi 15 Pro (Snapdragon 8 Elite). Though memory bus peaks all the time, ALUs in SP are mostly idle and thus GPU utilization and frequency keep low. On the other hand, existing mobile NPUs primarily aim at GEMM acceleration. LLM decoding can be more energy-efficient, if a hardware API is available to designate a smaller core with less ALUs and registers to process memory-bound Vec-Mat Multiplication.

## 8 CONCLUSION

We have presented MNN-AECS, a novel engine-level system design for energy-efficient LLM decoding on mobile devices. By adaptively selecting CPU cores, the most energy-efficient core selection is identified for LLM decoding without compromising speed. MNN-AECS has achieved 23% energy reduction without slowdown compared to original MNN and delivered 39%–78% energy reduction and 12%-363% decode speedup over 4 other popular on-device LLM engines.

## REFERENCES

[1] [n. d.]. ARM idle states binding description — ARM Linux Kernel documentation. https://www.kernel.org/doc/Documentation/devicetree/bindings/arm/idle-states.txt. [Accessed 03-05-2025].

[2] [n. d.]. Device Frequency Scaling — The Linux Kernel documentation — docs.kernel.org. https://docs.kernel.org/driver-api/devfreq.html. [Accessed 03-05-2025].

[3] 2013-2025. GitHub - doronz88/pymobiledevice3: Pure python3 implementation for working with iDevices (iPhone, etc...). — github.com. https://github.com/doronz88/pymobiledevice3. [Accessed 09-05-2025].

[4] 2022. kernel/sched/cpufreq_schedutil.c. https://android.googlesource.com/kernel/msm/+/refs/tags/android-12.1.0_r0.32/kernel/sched/cpufreq_schedutil.c. [Accessed 03-05-2025].

[5] 2022. kernel/sched/sched.h. https://android.googlesource.com/kernel/msm/+/refs/tags/android-12.1.0_r0.32/kernel/sched/sched.h. [Accessed 03-05-2025].

[6] 2022. kernel/sched/walt.c. https://android.googlesource.com/kernel/msm/+/refs/tags/android-12.1.0_r0.32/kernel/sched/walt.c. [Accessed 03-05-2025].

[7] 2024. Llama 3.2: Revolutionizing edge AI and vision with open, customizable models — ai.meta.com. https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/. [Accessed 19-04-2025].

[8] 2025. GitHub - pytorch/executorch: On-device AI across mobile, embedded and edge for PyTorch — github.com. https://github.com/pytorch/executorch. [Accessed 21-04-2025].

[9] Amey Agrawal, Ashish Panwar, Jayashree Mohan, Nipun Kwatra, Bhargav S. Gulavani, and Ramachandran Ramjee. 2023. SARATHI: Efficient LLM Inference by Piggybacking Decodes with Chunked Prefills. arXiv:2308.16369 [cs.LG] https://arxiv.org/abs/2308.16369

[10] Keivan Alizadeh, Iman Mirzadeh, Dmitry Belenko, S. Karen Khatamifard, Minsik Cho, Carlo C Del Mundo, Mohammad Rastegari, and Mehrdad Farajtabar. 2024. LLM in a Flash: Efficient Large Language Model Inference with Limited Memory. In ACL. https://arxiv.org/pdf/2312.11514

[11] Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. MathQA: Towards Interpretable Math Word Problem Solving with Operation-Based Formalisms. arXiv:1905.13319 [cs.CL] https://arxiv.org/abs/1905.13319

[12] Apple. 2024. Introducing Apple's On-Device and Server Foundation Models. https://machinelearning.apple.com/research/introducing-apple-foundation-models. [Accessed 11-05-2025].

[13] Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Meghan Cowan, Haichen Shen, Leyuan Wang, Yuwei Hu, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. 2018. TVM: an automated end-to-end optimizing compiler for deep learning. In Proceedings of the 13th USENIX Conference on Operating Systems Design and Implementation (Carlsbad, CA, USA) (OSDI'18). USENIX Association, USA, 579–594.

[14] Xinglei Dou, Lei Liu, and Limin Xiao. 2025. An Intelligent Scheduling Approach on Mobile OS for Optimizing UI Smoothness and Power. ACM Trans. Archit. Code Optim. 22, 1, Article 12 (March 2025), 27 pages. https://doi.org/10.1145/3674910

[15] Google AI Edge. 2025. LiteRT overview — ai.google.dev. https://ai.google.dev/edge/litert. [Accessed 19-04-2025].

[16] Google AI Edge. 2025. MediaPipe Solutions guide — ai.google.dev. https://ai.google.dev/edge/mediapipe/solutions/guide. [Accessed 19-04-2025].

[17] Shiwei Gao, Youmin Chen, and Jiwu Shu. 2025. Fast State Restoration in LLM Serving with HCache. In *Proceedings of the Twentieth European Conference on Computer Systems* (Rotterdam, Netherlands) *(EuroSys '25)*. Association for Computing Machinery, New York, NY, USA, 128–143. https://doi.org/10.1145/3689031.3696072

[18] Georgi Gerganov. 2024. GGUF. https://github.com/ggml-org/ggml/blob/master/docs/gguf.md. [Accessed 11-05-2025].

[19] Georgi Gerganov. 2025. GitHub - ggml-org/llama.cpp: LLM inference in C/C++ — github.com. https://github.com/ggml-org/llama.cpp. [Accessed 19-04-2025].

[20] Google. 2025. Gemini Nano with the Google AI Edge SDK | Android Developers — developer.android.com. https://developer.android.com/ai/gemini-nano. [Accessed 11-05-2025].

[21] Huawei. 2025. Product List | HUAWEI Support Global — consumer.huawei.com. https://consumer.huawei.com/en/support/product/. [Accessed 03-05-2025].

[22] Apple Inc. 2024. Apple Intelligence. https://www.apple.com/apple-intelligence/. [Accessed 03-05-2025].

[23] Apple Inc. 2025. iPhone — apple.com. https://www.apple.com/iphone/. [Accessed 03-05-2025].

[24] Qualcomm Technologies Inc. 2023. Qualcomm Snapdragon Mobile Platform OpenCL General Programming and Optimization. https://docs.qualcomm.com/bundle/publicresource/80-NB295-11_REV_C_Qualcomm_Snapdragon_Mobile_Platform_Opencl_General_Programming_and_Optimization.pdf. [Accessed 03-05-2025].

[25] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In *Proceedings of the 29th Symposium on Operating Systems Principles* (Koblenz, Germany) *(SOSP '23)*. Association for Computing Machinery, New York, NY, USA, 611–626. https://doi.org/10.1145/3600006.3613165

[26] Arm Limited. 2020. The Bifrost Shader Core Version 1.0 — developer.arm.com. https://developer.arm.com/documentation/102546/0100/?lang=en. [Accessed 03-05-2025].

[27] Arm Limited. 2020. Principles of High Performance guide Version 1.0 — developer.arm.com. https://developer.arm.com/documentation/102544/0100/?lang=en. [Accessed 03-05-2025].

[28] Chengdong Lin, Kun Wang, Zhenjiang Li, and Yu Pu. 2023. A Workload-Aware DVFS Robust to Concurrent Tasks for Mobile Devices. In *Proceedings of the 29th Annual International Conference on Mobile Computing and Networking* (Madrid, Spain) *(ACM MobiCom '23)*. Association for Computing Machinery, New York, NY, USA, Article 19, 16 pages. https://doi.org/10.1145/3570361.3592524

[29] Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. 2024. AWQ: Activation-aware Weight Quantization for On-Device LLM Compression and Acceleration. In *Proceedings of Machine Learning and Systems*, P. Gibbons, G. Pekhimenko, and C. De Sa (Eds.), Vol. 6. 87–100. https://proceedings.mlsys.org/paper_files/paper/2024/file/42a452cbafa9dd64e9ba4aa95cc1ef21-Paper-Conference.pdf

[30] Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring How Models Mimic Human Falsehoods. arXiv:2109.07958 [cs.CL] https://arxiv.org/abs/2109.07958

[31] Zechun Liu, Changsheng Zhao, Igor Fedorov, Bilge Soran, Dhruv Choudhary, Raghuraman Krishnamoorthi, Vikas Chandra, Yuandong Tian, and Tijmen Blankevoort. 2025. SpinQuant: LLM Quantization with Learned Rotations. In *The Thirteenth International Conference on Learning Representations*. https://openreview.net/forum?id=ogO6DGE6FZ

[32] Arm Ltd. 2025. big.LITTLE: Balancing Power Efficiency and Performance — arm.com. https://www.arm.com/technologies/big-little. [Accessed 03-05-2025].

[33] Chengfei Lv, Chaoyue Niu, Renjie Gu, Xiaotang Jiang, Zhaode Wang, Bin Liu, Ziqi Wu, Qiulin Yao, Congyu Huang, Panos Huang, Tao Huang, Hui Shu, Jinde Song, Bin Zou, Peng Lan, Guohuan Xu, Fei Wu, Shaojie Tang, Fan Wu, and Guihai Chen. 2022. Walle: An End-to-End, General-Purpose, and Large-Scale Production System for Device-Cloud Collaborative Machine Learning. In *Proceedings of USENIX Symposium on Operating Systems Design and Implementation (OSDI)*. USENIX, Carlsbad, CA, USA, 249–265.

[34] Meizu. 2025. AI Phones - Meizu Global — meizu.com. https://www.meizu.com/global/product-list/phones. [Accessed 03-05-2025].

[35] MLC team. 2023-2025. *MLC-LLM*. https://github.com/mlc-ai/mlc-llm

[36] Wei Niu, Md Musfiqur Rahman Sanim, Zhihao Shu, Jiexiong Guan, Xipeng Shen, Miao Yin, Gagan Agrawal, and Bin Ren. 2024. SmartMem: Layout Transformation Elimination and Adaptation for Efficient DNN Execution on Mobile. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3* (La Jolla, CA, USA) *(ASPLOS '24)*. Association for Computing Machinery, New York, NY, USA, 916–931. https://doi.org/10.1145/3620666.3651384

[37] Tom Olson. 2013. How low can you go? Building low-power, low-bandwidth ARM Mali GPUs — community.arm.com. https://community.arm.com/arm-community-blogs/b/mobile-graphics-and-gaming-blog/posts/how-low-can-you-go-building-low-power-low-bandwidth-arm-mali-gpus. [Accessed 22-04-2025].

[38] Samsung. 2025. All New Samsung Mobiles Prices & models | Samsung Jordan — samsung.com. https://www.samsung.com/levant/smartphones/all-smartphones/. [Accessed 03-05-2025].

[39] shareAI. 2023. ShareGPT-Chinese-English-90k Bilingual Human-Machine QA Dataset. https://huggingface.co/datasets/shareAI/ShareGPT-Chinese-English-90k.

[40] Gemma Team. 2024. Gemma 2: Improving Open Language Models at a Practical Size. arXiv:2408.00118 [cs.CL] https://arxiv.org/abs/2408.00118

[41] Llama Team. 2024. The Llama 3 Herd of Models. arXiv:2407.21783 [cs.AI] https://arxiv.org/abs/2407.21783

[42] Qwen Team. 2025. Qwen2.5 Technical Report. arXiv:2412.15115 [cs.CL] https://arxiv.org/abs/2412.15115

[43] Apple Device Support Tutorial. 2024. Using Sysdiagnose to Troubleshoot iOS or iPadOS. https://it-training.apple.com/tutorials/support/sup075/. [Accessed 09-05-2025].

[44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf

[45] Manni Wang, Shaohua Ding, Ting Cao, Yunxin Liu, and Fengyuan Xu. 2021. AsyMo: scalable and efficient deep-learning inference on asymmetric mobile CPUs. In *Proceedings of the 27th Annual International Conference on Mobile Computing and Networking* (New Orleans, Louisiana) *(MobiCom '21)*. Association for Computing Machinery, New York, NY, USA, 215–228. https://doi.org/10.1145/3447993.3448625

[46] Zekun Moore Wang, Zhongyuan Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao

Ni, Man Zhang, Zhaoxiang Zhang, Wanli Ouyang, Ke Xu, Wenhu Chen, Jie Fu, and Junran Peng. 2023. RoleLLM: Benchmarking, Eliciting, and Enhancing Role-Playing Abilities of Large Language Models. *arXiv preprint arXiv: 2310.00746* (2023).

[47] Rafael J. Wysocki. [n. d.]. CPU Performance Scaling — The Linux Kernel documentation — docs.kernel.org. https://docs.kernel.org/admin-guide/pm/cpufreq.html. [Accessed 03-05-2025].

[48] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. 2023. SmoothQuant: accurate and efficient post-training quantization for large language models. In *Proceedings of the 40th International Conference on Machine Learning* (Honolulu, Hawaii, USA) *(ICML'23)*. JMLR.org, Article 1585, 13 pages.

[49] Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024. Efficient Streaming Language Models with Attention Sinks. *ICLR* (2024).

[50] Xiaomi. 2025. Xiaomi Global Home — mi.com. https://www.mi.com/global/product-list/phone/. [Accessed 03-05-2025].

[51] Daliang Xu, Hao Zhang, Liming Yang, Ruiqi Liu, Gang Huang, Mengwei Xu, and Xuanzhe Liu. 2025. Fast On-device LLM Inference with NPUs. In *Proceedings of the 30th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1* (Rotterdam, Netherlands) *(ASPLOS '25)*. Association for Computing Machinery, New York, NY, USA, 445–462.

https://doi.org/10.1145/3669940.3707239

[52] Zhenliang Xue, Yixin Song, Zeyu Mi, Xinrui Zheng, Yubin Xia, and Haibo Chen. 2024. PowerInfer-2: Fast Large Language Model Inference on a Smartphone. arXiv:2406.06282 [cs.LG] https://arxiv.org/abs/2406.06282

[53] Rongjie Yi, Liwei Guo, Shiyun Wei, Ao Zhou, Shangguang Wang, and Mengwei Xu. 5555. EdgeMoE: Empowering Sparse Large Language Models on Mobile Devices . *IEEE Transactions on Mobile Computing* 01 (Feb. 5555), 1–16. https://doi.org/10.1109/TMC.2025.3546466

[54] Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, Soojeong Kim, and Byung-Gon Chun. 2022. Orca: A Distributed Serving System for Transformer-Based Generative Models. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*. USENIX Association, Carlsbad, CA, 521–538. https://www.usenix.org/conference/osdi22/presentation/yu

[55] Yinmin Zhong, Shengyu Liu, Junda Chen, Jianbo Hu, Yibo Zhu, Xuanzhe Liu, Xin Jin, and Hao Zhang. 2024. DistServe: disaggregating prefill and decoding for goodput-optimized large language model serving. In *Proceedings of the 18th USENIX Conference on Operating Systems Design and Implementation* (Santa Clara, CA, USA) *(OSDI'24)*. USENIX Association, USA, Article 11, 18 pages.