# AttestLLM: Efficient Attestation Framework for Billion-scale On-device LLMs

Ruisi Zhang*
ruz032@ucsd.edu
UC San Diego

Yifei Zhao*
yifei.zhao@ucf.edu
University of Central Florida

Neusha Javidnia
njavidnia@ucsd.edu
UC San Diego

Mengxin Zheng
mengxin.zheng@ucf.edu
University of Central Florida

Farinaz Koushanfar
farinaz@ucsd.edu
UC San Diego

## Abstract

As on-device LLMs (e.g., Apple on-device Intelligence) are widely adopted to reduce network dependency, improve privacy, and enhance responsiveness, verifying the legitimacy of models running on local devices becomes critical. Existing attestation techniques are not suitable for billion-parameter Large Language Models (LLMs), struggling to remain both time- and memory-efficient while addressing emerging threats in the LLM era.

In this paper, we present AttestLLM, the first-of-its-kind attestation framework to protect the hardware-level intellectual property (IP) of device vendors by ensuring that only authorized LLMs can execute on target platforms. AttestLLM leverages an algorithm/software/hardware co-design approach to embed robust watermarking signatures onto the activation distributions of LLM building blocks. It also optimizes the attestation protocol within the Trusted Execution Environment (TEE), providing efficient verification without compromising inference throughput. Extensive proof-of-concept evaluations on LLMs from Llama, Qwen, and Phi families for on-device use cases demonstrate AttestLLM's attestation reliability, fidelity, and efficiency. Furthermore, AttestLLM enforces model legitimacy and exhibits resilience against model replacement and forgery attacks.

**CCS Concepts:** • **Security and privacy → Trusted computing**.

**Keywords:** Attestation, Trusted Execution Environment, On-device LLMs

## 1 Introduction

In recent years, on-device large language models (LLMs) [36, 58, 60] have been widely deployed in consumer electronic devices, such as smartphones, augmented/virtual reality headsets, and smart home devices, to enable local model inference. Companies such as Apple and Qualcomm [5, 45] invest considerable effort in optimizing both model compression algorithms and hardware architecture to reduce latency, enhance user privacy, and decrease reliance on network connectivity. However, this shift toward local inference introduces a largely underexplored challenge: end users and untrusted third parties along the supply chain (e.g., resellers) have control over the device and its rich execution environment (REE) [7, 63]. This opens the door to potential unauthorized misuse, where the pre-installed, vendor-optimized LLMs can be replaced with illegal and unauthorized models. Such tampering can not only violate the hardware vendor's intellectual property (IP) by running undesired models on the platform, but also erode its reputation and customer trust if devices exhibit degraded performance, unsafe behaviors, or inconsistent user experiences.

Existing approaches to intellectual property (IP) protection for large language models (LLMs) primarily focus on the software level by watermarking. These algorithms [30, 61, 63] embed unique signatures in model weights or trigger identifiable model behavior patterns under specific input. They are, however, not suitable for protecting the device from hardware-level IP infringement, where vendors are keen to have strong enforcement over models running on hardware platforms.

Attestation [7] is introduced to provide the execution control of models with the help of trusted execution environment (TEEs). It securely verifies the legitimacy of the model by checking if its decoded watermark matches the device-specific signature. Only watermarked models that passed attestation can run inference in REE, whereas unauthorized model execution is aborted. However, such attestation is typically performed for deep neural networks with millions of parameters, and scaling the attestation to support billion-parameter LLMs remains challenging. The bottlenecks mainly come from three aspects: (i) **hardware constraints**: on-device TEE, such as Arm TrustZone [44], only has 10−32 MB secure memory capacity, which is insufficient to store LLMs' billion-scale parameters and intermediate activations; (ii) **efficiency**: the attestation process shall incur minimal runtime overhead to ensure smooth user experience. This requires a co-design of the watermark verification algorithm and the attestation pipeline to reduce the latency; (iii) **emerging threats**: the LLMs deployed on the edge device are typically fine-tuned from a limited set of foundation models, which may be inferred by the adversary. Thus, the on-device LLMs become susceptible to unauthorized model

---

Ruisi Zhang, Yifei Zhao, Neusha Javidnia, Mengxin Zheng, and Farinaz Koushanfar
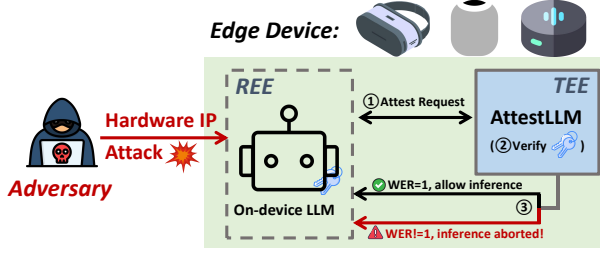


**Figure 1.** Overview of AttestLLM's attestation. On-device LLM is deployed in the REE, where the potential adversary may conduct hardware IP attacks. Our framework AttestLLM requests attestation periodically and only allows execution of authenticated models while blocking unauthorized ones.

replacement or forgery attacks. Designing watermarking algorithms that are resilient to such attacks becomes critical.

In this paper, we present **AttestLLM**, the first attestation framework to protect the hardware-level intellectual property of on-device LLMs. The framework consists of two key components: offline LLM watermarking and online attestation. **Offline LLM watermarking** encodes unique watermarking signatures onto the activation outputs of each transformer block [1]. This design enables every block to serve as a standalone verification unit, allowing attestation to selectively sample and verify a subset of blocks. It avoids the execution of full LLM inference, thereby reducing memory footprint and runtime overhead. To preserve model quality, AttestLLM employs a sensitivity analysis across transformer blocks and allocates signature lengths adaptively under a global watermarking budget, where performance-critical blocks receive shorter signatures, while less sensitive blocks are assigned longer ones. Then, these signatures are embedded onto the edge LLMs via an optimization-based watermark insertion algorithm, which encodes the watermark without hurting model performance. **Online attestation** overcomes the memory limitations of Arm TrustZone [44] by adopting a virtualization-based secure enclave, which isolates TEE resources from the REE to enable scalable attestation under constrained environments. To enhance efficiency, in every step, AttestLLM dynamically samples a subset of transformer blocks for verification, while ensuring sufficient attestation strength to resist evasion attempts by malicious users. Those blocks are executed in parallel, and AttestLLM decouples verification computations to overlap them with secure copy communication to further reduce attestation latency.

As detailed in Figure 1, AttestLLM enables attestation for on-device LLMs deployed on edge devices. The LLM, embedded with watermarking, is executed within the REE. Prior to inference, ① an attestation request is issued, and ② our

proposed AttestLLM running in the TEE verifies the embedded watermarking keys. ③ If the watermark extraction rate (WER) equals 1, the model is authorized to run; otherwise, inference is blocked, thereby preventing the execution of unauthorized or tampered models, especially in the presence of potential hardware IP attacks.

In summary, our contributions are summarized as follows:

■ We present AttestLLM, the first-of-its-kind attestation framework designed to protect hardware-level IP for hardware device vendors, allowing only authorized LLMs to run on the target platform.

■ We co-design algorithm/software/hardware to scale the attestation to billion-parameter LLMs with: (i) a **quality-preserving and robust watermarking algorithm** to effectively encode device-specific signatures onto the activation distribution of each transformer block; (ii) **efficient and scalable attestation under hardware resource constraints**, achieved through TEE-aware optimizations, including virtualization based secure enclave, sub-block sampling, and communication/computation overlapping, to overcome the memory limitations and reduce attestation overhead.

■ We validate AttestLLM's effectiveness on a variety of mainstream on-device LLMs, demonstrating: (i) AttestLLM keeps the model's fidelity while ensures 100% watermark extraction accuracy; (ii) AttestLLM can efficiently attest the LLM by introducing minimal additional overhead over untrusted REE inference, achieving at least 12.5× lower latency overhead and 9.5× lower energy overhead compared to TEE-shield LLM inference [34]. (iii) AttestLLM resists various model replacement and forgery attacks from adversaries and only allows authorized model to be run on the hardware platform.

## 2 Background and Related Work

In this section, we first introduce the related work for on-device large language model watermarking. Then, we present the background and literature on the trusted execution environment for LLM attestation.

### 2.1 On-device Large Language Model Watermarking

To deploy large language model on edge devices, researchers typically leverage various compression techniques, such as knowledge distillation [28, 57] and quantization [14, 35], to obtain a compact and low-precision format and reduce the inference burden. One of the most significant steps is quantization, which maps the model from high-precision float32 [42] or bfloat16 [12] into INT8 or lower. These model optimizations, often times, are co-designed with target hardware platforms, optimizing kernel implementations [20, 24],

---

[1]Transformer block is the building unit of LLM, typically consists of multi-head self-attention, feed-forward neural networks, layer normalization, and some residual connections

memory layouts [62], and resource allocation [33] strategies to fully exploit hardware capabilities for faster inference.

AttestLLM enforces the usage control of the on-device LLMs running on the hardware platforms, and thus protects the hardware-level intellectual property (IP) of the device vendor [7]. In contrast, most existing watermarking research focuses on protecting the model's IP itself at the software level, thereby crafted to satisfy different criteria. Most of these software-level model watermarking goes from two directions: (i) **parameter-based watermarking** [10, 61, 63]: embedding unique identifiers directly into the model's parameters, enabling ownership verification when the model weights are accessible. It could be done during training [10] by adding an additional regularization term into the loss function, or post-hoc [61, 63] by weight sensitivity analysis to enable watermarking insertion without introducing quality degradations; and (ii) **backdoor-based watermarking** [27, 29, 30]: encoding secret triggers onto the LLM's token prediction probabilities, and produce distinctive behavior only when specific inputs are encountered. This is done typically by fine-tuning pre-trained language models with multi-task learning and specific triggers to enable robust watermark extraction even after downstream fine-tuning.

## 2.2 TEE for LLM Attestation

### 2.2.1 Trusted Execution Environments.
Trusted Execution Environment (TEE) is a secure area within the main processor that ensures the confidentiality and integrity of code and data through hardware-enforced isolation [48]. For x86 architectures, several TEE implementations have been introduced. Intel SGX [9] provides application-level enclaves with hardware-encrypted memory. AMD SEV [49], along with its extensions SEV-ES and SEV-SNP, offers the confidential virtual machine through memory and register encryption. On the edge side, most devices are based on the ARM architecture[46]. Arm TrustZone [44] is a hardware-based security extension that partitions the processor into two isolated execution environments: the Normal World (NW) and the Secure World (SW) with the SW designed to run critical applications [2].

### 2.2.2 Secure LLM Inference with TEEs.
Recent research has increasingly focused on enabling secure inference of LLMs using TEEs [32, 41, 52], which offer secure isolation to ensure the integrity and untampered execution. Most of these works, however, are performed on cloud environments, where secure enclaves have sufficient memory and compute capacity to run full-model inference entirely within the TEE. This setup ensures both confidentiality and integrity while still maintaining acceptable performance. Secure LLM inference in these systems is typically supported by platforms such as Intel processors with SGX [9], TDX [8], and AMD CPUs with SEV [56].

In contrast, edge devices face significant resource constraints with limited memory and compute power. Full model inference in these cases becomes infeasible. To address these limitations, prior work has explored several alternative strategies: (i) **TEE-shield Inference** [15, 34, 64]: slicing model weights into memory-fitting segments and executing each segment sequentially within the TEE. However, such designs involve frequent world switches and secure memory operations, leading to non-trivial communication and context-switch overhead that limits efficiency; (ii) **TEE-Based Attestation**: Instead of executing inference inside the TEE, this approach assigns attestation and access control to the secure world, while the normal world performs inference to maximize performance. The TEE verifies model integrity before and during execution to prevent both offline tampering (e.g., pre-execution model modification) and online tampering (e.g., runtime weight corruption). Existing frameworks, such as DeepAttest [7], rely on Intel SGX [9] to perform attestation, rather than the ARM-based architectures commonly used in edge devices today[46].

These designs were primarily developed for DNNs with relatively small parameter sizes. Scaling TEE-based frameworks to billion-parameter LLMs introduces additional challenges: (i) **strict hardware constraints**: ARM TrustZone [44], the dominant on-device TEE, provides only 10–32 MB of secure memory [43], has restricted instruction sets, and lacks efficient computation libraries, making it challenging to meet the memory and performance demands of LLM-scale attestation; (ii) **emerging security threats**: LLMs are typically fine-tuned from a limited set of foundation models, which makes it easier for adversaries to infer their architectures and perform model replacement attacks. Moreover, unlike DNNs, LLM inference proceeds sequentially, with each step conditioned on previously generated tokens. This sequential process increases susceptibility to memory attacks, including tampering with intermediate states or injecting malicious components such as LoRA adapters [25].

In this paper, AttestLLM adopts a hardware–software co-design approach to **co-optimize the watermarking algorithm with the underlying hardware platform verification**. This enables the first-of-its-kind efficient LLM attestation to protect the hardware-level intellectual property without degrading the device's performance or usability.

## 3 Threat Model

In this section, we introduce our threat model, including the scenario for applying AttestLLM, potential threats from adversaries, and criteria AttestLLM shall meet for usability and applicability.

---

[2] Throughout the paper, we use "normal world" and "rich execution environment "; "secure world" and "trusted execution environments" interchangeably.

## 3.1 Scenario

As in Figure 1, hardware providers [5, 36, 45] embed large language models (LLMs) directly onto edge devices to reduce network dependency, strengthen privacy, lower inference costs, and deliver more responsive performance. For instance, Apple integrates on-device intelligence across products such as the HomePod, iPhone, and iPad through Apple Silicon hardware [5]; and Qualcomm [45] enables on-device LLM inference on AR/VR platforms to support voice control, context-aware assistance, and seamless user experiences.

While these on-device deployments offer tangible benefits, they introduce a critical security challenge: end users and third parties along the supply chain (e.g., resellers) may be adversaries. With full access to the embedded LLM, they can exploit the device to run unauthorized models, which potentially repurpose the compute resources for other unapproved applications or deploy models that violate safety and ethical standards. This undermines both the intended functionality of the hardware and the provider's security assurances and reputation. With AttestLLM, hardware vendors can attest whether the LLMs running in REE are authenticated; if not, execution of the unauthorized LLMs is aborted.

## 3.2 Attackers' Capability and Potential Threats

**Attackers' Capability**: We consider the adversary possesses access to the system but cannot manipulate the secure Trusted Execution Environment (TEE) inside the device. They have prior knowledge of the LLM model architecture and the watermarking algorithm. However, they do not know the model weights since they are stored in a shuffled memory region and cannot obtain the watermark signatures in the secure memory. They aim to bypass the attestation to host another unauthorized model on the hardware device.

**Potential Threats**: To achieve this, the adversaries may execute the following attacks. We perform a detailed analysis and quantitative evaluations in Section 5.6.

- **Model replacement attack**: The adversary has prior knowledge of the model architecture, but do not have access to the parameters in the shuffled memory region. They thus host a malicious model with the same architecture but different parameters.
- **Watermark Forgery attack**: The adversary forges the signature stored in the secure memory inside the TEE. Then, they attempt to embed the signature onto the unauthorized LLM to fool the hardware platform to host an illegitimate model.

We also note that adversaries may leverage Low-Rank Adaptation (LoRA) [25] to fine-tune the model or partially replace the model submodules to bypass the attestation and host an illegitimate model. These attacks could be tackled by control flow attestation [4, 13] and memory-layout randomization [7, 38], which ensure that the model attested in the secure world is the same one executed in the normal world.

These are orthogonal to AttestLLM and can be combined to further boost security with protocols in Section 4.3.3.

## 3.3 Attestation Criteria

We design AttestLLM with the following criteria in mind to ensure its usability and practicality.

- **Efficiency**: The attestation process shall be lightweight, ensuring the usability of the LLM service in REE.
- **Fidelity**: The watermarked model's quality shall not be degraded after watermark insertion.
- **Robustness and Reliability**: The embedded watermark shall withstand various attacks from adversaries. The attestation framework shall only authorize the execution of watermarked LLMs, whereas unauthorized models' execution are blocked.

## 4 Method

### 4.1 Overview

As in Figure 2, AttestLLM consists of two main parts: (i) offline LLM watermarking, which encodes unique device-specific watermarks onto the model; and (ii) online attestation with TEE, which securely verifies the watermarks within the hardware's protected enclave, ensuring only the owner's authorized model can be inferred in the normal world. Details of each component are in Section 4.2 and Section 4.3, respectively.

### 4.2 Offline Large Language Model Watermarking

#### 4.2.1 Watermarking Workflow.
The watermarking in AttestLLM consists of three consecutive steps: (i) **sensitivity analysis**, (ii) **optimization-based watermarking insertion**, and (iii) **watermark verification**. As in Figure 2, at each attestation step, AttestLLM selects a subset of transformer blocks within the LLM to verify model legitimacy. Thus, our watermarking workflow is designed so that the watermark in each block can be independently extracted and verified as a standalone unit.

AttestLLM takes the original full-precision LLM $M$ and the signature sequence $B = \{b_1, b_2, ..., b_{|B|}\}$ as input. In $B$, each element $b_i \in \{0, 1\}$. At the $i$-th transformer block, the assigned signature is $B_i$ from the total signature sequence $B$, and the activation $\mathcal{A}_i$ from a trigger dataset $\mathcal{D}_{wm}$ is projected into a watermarking space using a projection matrix $WM_i$. The watermark is then embedded by adjusting and updating the model's weights $W_i$ so that the projected activations $WM_i(\mathcal{A}_i)$ closely align with the corresponding signature bits in $B_i$. During attestation, these projections are checked against the expected signature to confirm the model's authenticity.

#### 4.2.2 Sensitivity Analysis.
Since different transformer blocks contribute unevenly to model performance because
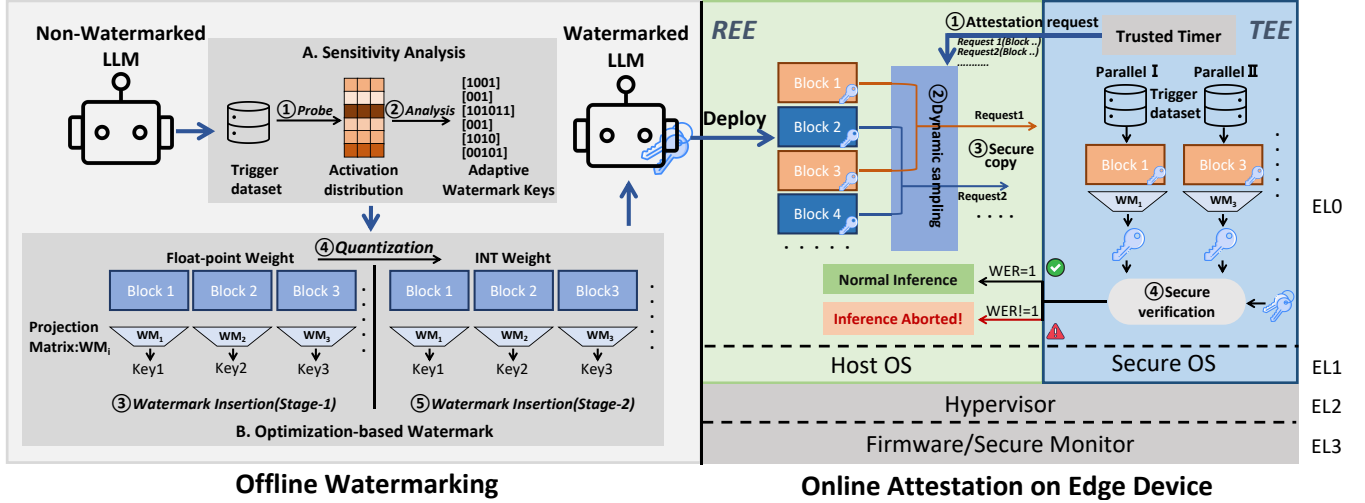
**Figure 2.** AttestLLM Framework. The left is Offline Watermarking. Device-specific watermarking signatures are encoded into Transformer blocks by (A) probing activations with the trigger dataset to allocate adaptive watermark keys capacity for each block and (B) optimization-based two-stage watermark insertion (Stage-1 before quantization and Stage-2 after quantization). The right is Online attestation on the edge device. The watermarked model is deployed on the REE side. TEE periodically requests attestation, dynamically samples blocks, securely copies them, and in parallel verifies the embedded watermarks with the trigger set. Successful verification permits normal inference, otherwise execution is aborted. The figure also illustrates isolation across privilege levels (EL0–EL3).

they handle different stages of the input feature embedding[23]. Some blocks can carry more watermark signature bits without harming utility, while others are more sensitive and should be assigned fewer bits. Weight parameter saliency correlates strongly with activation magnitude, where channels with larger activations process more information and are therefore more salient [31].

Inspired by this, say there are $|M|$ transformer blocks and a total watermark length of $|B|$. We allocate signature length to block $i$ as $B_i$. As in Equation 1, $B_i$ is inversely proportional to the block's peak activation magnitude $\max |\mathcal{A}_i|$, ensuring that more salient blocks receive fewer signature bits and less salient blocks absorb more.

$$B_i = |B| * \frac{1/\max |\mathcal{A}_i|}{\sum_1^{|M|} \max |1/\mathcal{A}_i|} \qquad (1)$$

To improve attestation efficiency, we reduce computational overhead by encoding the watermark on only a subset of activation channels, denoted as $|C|$, from each block. This lowers the cost of the projection performed by $WM_i$, as a smaller matrix dimension is required to perform projection. The subset is selected through multinomial sampling, where each channel's probability $p_{i,c}$ is defined in Equation 2. Channels with higher activation magnitudes are assigned lower sampling probabilities to preserve model fidelity, while those with lower activations are sampled more often to embed richer signature information. For ease of explanation, we use $\mathcal{A}_i$ to denote $\mathcal{A}_i[C]$ throughout the remainder of the paper.

$$p_{i,c} = \frac{1/|\mathcal{A}_{i,c}|}{\sum_{c=1}^{|\mathcal{A}_i|} \frac{1}{|\mathcal{A}_{i,c}|}}, \qquad (2)$$

**4.2.3 Optimization-based Watermark Insertion.** To deploy the model on an edge device, the most significant step of model compression from the perspective of watermarking is quantization, which substantially changes weight distributions and data types. Thus, we devise a quantization-aware optimization-based watermarking insertion algorithm in Algorithm 1. Given the activation from the last transformer block as $\mathcal{A}_{i-1}$, the activation $\mathcal{A}_i$ from the $i$-th transformer block $M_i$ is mapped into the watermarking space by the projection matrix $WM_i$. In each transformer block, the projection matrix $WM_i$ approximates the ground-truth signature $B_i$ with the optimization objective in Equation 3.

$$\min_{M_i, WM_i} \left| WM_i \big( M_i(\mathcal{A}i - 1) \big) - \mathcal{B}_i \right| \qquad (3)$$

To achieve faithful signature embedding while preserving model utility, AttestLLM introduces a two-stage watermark insertion before and after quantization. In the pre-quantization stage, the optimization is performed on the full-precision model. Thus, a gradient-based approach is employed. This facilitates rapid convergence and a good initial approximation of the watermark signature by adjusting both $M_i$ and $WM_i$ with backpropagation. To preserve model utility, we add a penalty $\alpha \|M_i' - M_i\|^2$ to the loss to ensure the

watermarked model parameters $M_i'$ do not diverge far from the original parameter $M_i$.

The quantized model weights are discrete, and gradient-based approaches are unable to optimize the model to approximate the objective in Equation 3. Hence, in the post-quantization stage, zeroth-order optimization [51] is applied. As in Algorithm 1, a subset of parameters $\Theta$ is perturbed along random directions, and the losses $L^+$ and $L^-$ are estimated through finite-difference evaluations [47]. These estimates guide the updates to fine-tune the quantized transformer block $M_{qi}$, reducing the discrepancy between the predicted watermark signature and the target signature $B_i$. This two-stage design ensures that watermark signatures remain well-aligned with $B_i$ after quantization, while limiting updates to a small parameter subset so as to preserve the overall model fidelity.

---

**Algorithm 1** Quantization-aware two-stage WM embedding onto transformer block $M_i$

---

**Require:** Base transformer blocks $M_i$, activation from $(i-1)$-th transformer block $A_{i-1}$, signature $B_i$, projection matrix $WM_i$

**Ensure:** Watermarked and quantized transformer block $M_{qi}$

    **for** epoch = 1 to Epoch$_{\text{pre}}$ **do**

        Compute activations: $A_i = M_i(A_{i-1})$

        Compute WM prediction: $B_i'(x) = WM_i \cdot A_i$

        Compute losses:

$$L = \sum_i \|B_i' - B_i\| + \alpha \|M_i' - M_i\|^2$$

        Update $M_i', WM_i$ via gradient descent with learning rate $\eta_{\text{pre}}$

    **end for**

    Quantize transformer block to obtain $M_{qi}$

    **for** t = 1 to Epoch$_{\text{post}}$ **do**

        Select a subset $\Theta$ of parameters for post-quant update

        Sample a unit perturbation $u$ in parameter space of $\Theta$

        Evaluate $L^+ = L(M_{qi} + \mu u, WM_i, B_i, A_{i-1})$

        Evaluate $L^- = L(M_{qi} - \mu u, WM_i, B_i, A_{i-1})$

        Estimate gradient: $\hat{g} \leftarrow \frac{L^+ - L^-}{2\mu} \cdot u$

        Update: $M_{qi} \leftarrow M_{qi} - \eta_{\text{post}} \hat{g}$

    **end for**

---

#### 4.2.4 Watermark Verification.

Given the trigger dataset $D_{wm}$, in the $i$-th transformer block $M_{qi}$, AttestLLM will obtain the activation $\mathcal{A}_i$ and map it to the watermarking space using the secret projection matrix $WM_i$. The signature is decoded as $B_i'$. Then, the watermark extraction accuracy is computed using Equation 4, where $|B_i|$ is the length of the inserted signature, and $|B_i|'$ is the number of matching signature bits.

$$\%WER = 100 \times \frac{|B_i|'}{|B_i|} \tag{4}$$

We note that this watermark verification can be performed in parallel by checkpointing the activation output from the $i-1$-th transformer block, denoted as $\mathcal{A}_{i-1}$. AttestLLM, then, can bypass recomputing the first $i-1$ transformer blocks with $D_{wm}$, and instead only execute the $i$-th block computation using the stored activation $\mathcal{A}_{i-1}$.

### 4.3 Online Attestation with TEE

#### 4.3.1 TEE Setup and Attestation Workflow.

To overcome the inherent limitations of TrustZone, such as limited secure memory and restricted instruction sets. We adopt a virtualization-based approach using protected Kernel Virtual Machines (pKVM) [11]. pKVM provides software-isolated enclaves enforced by a minimal, privileged hypervisor, which operates at a higher privilege level (EL2) and is solely responsible for enclave isolation and security [26, 40]. The hypervisor partitions memory into REE and TEE, each running in isolation. It enforces strict memory isolation by tracking exclusive page ownership and mediating all transitions through Stage-2 MMU translation to prevent unauthorized access. In our design, the REE side handles LLM inference, while the TEE enclave performs attestation. This setup supports Linux-based stacks and larger memory footprints, making it well-suited for secure and efficient LLM attestation on edge devices.

The attestation process begins with a secure platform boot anchored by a hardware root of trust, which verifies the privileged hypervisor and boot chain before the TEE is created. Once the enclave is initialized, it requests verification of the LLM model. To accelerate secure copy, the REE stages the attestation-selected transformer blocks in a read-only memory region exposed via virtio-pmem [53], allowing the enclave to copy them directly into private pages for verification. The enclave also retrieves the encrypted trigger dataset, decrypts it using the AES algorithm[3], and uses it to perform secure verification of the model blocks. If verification succeeds, inference continues in the REE; otherwise, the enclave raises an alert or may enforce a device reset to block untrusted execution. To maintain model integrity during runtime, attestation is periodically triggered using the same procedure to detect any unauthorized modification.

#### 4.3.2 Optimizations for Efficient Attestation.

*Dynamic attestation.* To ensure both security and efficiency in resource-constrained edge environments, we propose a *dynamic attestation mechanism* that activates adaptively based on system state. The design goal is to offer a *tunable trade-off* between security level and attestation overhead by adjusting runtime parameters.

We define an *attestation interval parameter $f$*, which controls how often the model is verified. Specifically, attestation is triggered once every $f$ generated tokens during inference. The value of $f$ is dynamically adjusted according to

the device's workload and operational context: for privacy-critical scenarios, a small $f$ (e.g., $f = 50$) ensures frequent verification, while in latency-sensitive or conversational applications, $f$ can be relaxed (e.g., $f = 200$–$500$) to reduce runtime overhead.

When attestation is triggered, the system randomly samples $k$ transformer blocks out of $L$ total blocks. The weights of the selected blocks are securely copied into the TEE, where the watermarks are extracted and compared against encoded signatures. The model is allowed to proceed only if the WER is one; otherwise, execution is aborted.

The probability that an attacker tampers with $t$ transformer blocks and evades detection in a single attestation is $P_{\text{miss}} = \frac{\binom{L-t}{k}}{\binom{L}{k}}$, where $L$ is the total number of transformer blocks and $k$ is the number of sampled blocks for attestation. Over $r = \left\lfloor \frac{m}{f} \right\rfloor$ attestation rounds during $m$ generated tokens (with attestation triggered every $f$ tokens), the overall evasion probability is

$$P_{\text{evasion}} = \left( \frac{\binom{L-t}{k}}{\binom{L}{k}} \right)^r . \tag{5}$$

***Attestation pipeline optimization.*** To further minimize attestation latency without exceeding the resource limits of the TEE, we design a pipelined and partially parallel verification mechanism. Specifically, we divide each attestation round into three stages: (i) secure memory copy of $k$ selected transformer blocks from the REE into the enclave, (ii) watermark extraction from each block, and (iii) signature verification. These stages are decoupled and overlapped whenever possible to enable a streamlined data flow. Within the enclave, we implement controlled parallelism for the extraction and verification stages. Rather than processing all $k$ blocks sequentially, we dynamically allocate a bounded number of worker threads based on available enclave resources, such as secure memory size and the number of usable computational cores. Each worker independently processes one transformer block, extracting the watermarks following Section 4.2.4 and decides if the LLM inference can be performed in REE.

Besides, we optimize performance by adopting an early-exit strategy: if a mismatch is detected in any block, the verification procedure aborts immediately without processing the remaining blocks. This mechanism not only accelerates the rejection of tampered models but also reduces unnecessary computation when integrity is already compromised.

### 4.3.3 Protocols for Robust Security.
Our attestation framework detects both offline and runtime tampering of model weights. However, in edge devices, adversaries with access to the REE may still launch memory attacks to bypass the attestation, including adding LoRA adaptors [25] that hook into intermediate states to fine-tune the model, or replacing part of the LLM submodules. These, however, could be defended by **memory-layout randomization**

(MLR) [7, 37, 50]. Such re-randomization dynamically relocates code and data sections and disrupts the persistent layouts, which largely reduces the possibilities of brute-force attempts to search stored model weights and watermark keys. In addition, we use control-flow attestation (CFA)[4, 13] to verify the execution path of the REE program, ensuring that it matches the program attested by the TEE. Together, these protocols ensure that AttestLLM's attestation in the TEE is faithfully propagated and that the hardware-level IP of device vendors remains protected.

### 4.4 AttestLLM API

The AttestLLM API can be integrated with real-world on-device LLM deployments. It includes two primary functions: the **offline embedding** API to insert device-specific watermarks into LLMs before distribution, and the **online attestation** API to verify those watermarks inside a secure enclave during inference. Developers can configure security parameters such as the trigger dataset, attestation interval $f$, and the number of transformer blocks $k$ for verification, to achieve a balance between security and overhead.

## 5 Experiments

We conduct experiments to seek answers to the following critical research questions:
❶ Whether the encoded signatures can be reliably attested?
❷ Whether watermarking large language models will degrade the model performance?
❸ Whether constant model attestation will incur significant time/memory overhead on edge devices?
❹ How robust are the encoded signatures to malicious attacks from adversaries?

We empirically answer questions ❶ and ❷ in Section 5.3 and Section 5.4. Question ❸ and ❹ are addressed in Section 5.5 and 5.6 respectively.

### 5.1 Experiment Setup

We perform evaluations on language models targeting on-device use cases, including Llama3 1B, 3B, and 8B [21], Qwen3 4B and 8B [59], and Phi-4 15B [2]. The details of models are summarized in Table 1.

**Table 1.** Configurations of models used for benchmarking.

| Model | Param. | Blocks | Hidden Size | Vocabulary Size |
|---|---|---|---|---|
| Llama 3 | 1B | 16 | 2048 | 128256 |
| Llama 3 | 3B | 28 | 3072 | 128256 |
| Llama 3 | 8B | 32 | 4096 | 128256 |
| Qwen 3 | 4B | 36 | 2560 | 151936 |
| Qwen 3 | 8B | 36 | 4096 | 151936 |
| Phi-4 | 15B | 40 | 5120 | 100352 |

The models are quantized to 8-bit using TorchAO [54] and 4-bit using AWQ [31]. The watermark insertion in Section 4.2 is performed on a CPU/GPU server with 4 NVIDIA RTX 6000

| Method | Metrics | INT4 Quantization | | | | | | | INT8 Quantization | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Model | Llama 3 | | | Qwen 3 | | Phi-4 | $\bar{\Delta}$ | Llama 3 | | | Qwen 3 | | Phi-4 | $\bar{\Delta}$ |
| | | 1B | 3B | 8B | 4B | 8B | 15B | | 1B | 3B | 8B | 4B | 8B | 15B | |
| non-WM | PPL $\downarrow$ | 10.84 | 8.22 | 6.64 | 14.88 | 9.98 | 6.78 | - | 9.88 | 7.88 | 6.36 | 13.82 | 9.63 | 6.75 | - |
| | Zero-shot Acc $\uparrow$ | 59.59% | 67.08% | 71.38% | 62.09% | 66.16% | 71.90% | - | 60.85% | 67.45% | 72.11% | 62.59% | 67.27% | 72.89% | - |
| AttestLLM | PPL $\downarrow$ | 10.92 | 8.24 | 6.66 | 15.07 | 9.96 | 6.85 | 0.62% | 9.92 | 7.89 | 6.37 | 13.72 | 9.62 | 6.76 | 0.07% |
| | Zero-shot Acc $\uparrow$ | 59.49% | 67.10% | 71.18% | 61.97% | 66.18% | 71.89% | 0.09% | 61.31% | 67.59% | 72.04% | 62.50% | 67.66% | 72.70% | 0.15% |
| | WM Extract. Acc | 100% | 100% | 100% | 100% | 100% | 100% | - | 100% | 100% | 100% | 100% | 100% | 100% | - |

**Table 2.** AttestLLM's performance in watermarking on-device large language models. Our evaluation includes Perplexity (PPL) to evaluate generated text fluency and Zero-shot Accuracy (Zero-shot Acc) for next-token prediction accuracy. $\bar{\Delta}$ is the average performance degradation compared with non-watermarked models.

GPUs, 256 GB RAM memory, and Intel Xeon CPUs connected via a 10Gbps Ethernet interface. In the first stage, gradient-based optimization, we use the Adam optimizer and set the learning rate $\eta_{pre}$ to 6e-6. The regularization parameter $\alpha$ is set to 1e-3. In the second stage, zeroth gradient optimization, the perturbation strength $\mu$ is set to 2 in INT4 quantization and 20 in INT8 quantization. The subset $\Theta$ size is set to 100, and the gradient $\eta_{post}$ is 0.5 in INT4 quantization and 0.1 in INT8 quantization. The signature length $|B|$ is 20-bit. The pre-quantization optimization is optimized for $Epoch_{pre} = 20$ epochs, and the post-quantization optimization is optimized for $Epoch_{post} = 40$ epochs. We use SimCSE dataset [17] for trigger dataset as a proof-of-concept. Then channel size $|C|$ is set to be 40% of the total hidden dimension size.

To evaluate our watermarking method and attestation protocol in realistic edge environments in Section 4.3, we conduct experiments on a commercial off-the-shelf device equipped with an 8-core ARM processor (4×Cortex-A76 + 4×Cortex-A55), 16GB LPDDR4X memory, and Android 13 as the host operating system. The LLM inference runs in REE, while the attestation framework is deployed inside a pKVM-protected enclave. The secure enclave is allocated 512MB of memory and runs a minimal Microdroid kernel[19] with a stripped-down root filesystem to reduce the trusted computing base (TCB). The pKVM-based setup allows us to instantiate a lightweight TEE on commodity ARM hardware without modifying the secure world or boot chain. For online attestation, the numbers of transformer blocks $L$ and sampled blocks $k$ $(L, k)$ are set to Llama3-1B (16, 2), Llama3-3B (28, 2), Llama3-8B (32, 4), Qwen3-4B/8B (36, 4), and Phi4-15B (40, 6). By default, we set to attestation interval to be $f = 100$ tokens.

### 5.2 Evaluation Criteria

To answer the aforementioned questions, we evaluate the system performance from the following aspects:
▶ Watermark Extraction Accuracy: The percentage of signatures successfully extracted for attestation.
▶ LLM Performance: We use *Perplexity (PPL)* to evaluate generated text fluency on WikiText Dataset [39], and, *Zero-shot Accuracy (Zero-shot Acc)* for next-token prediction evaluation on the mean of LAMBADA, HellaSwag, PIQA, and WinoGrande Datasets [16].

▶ Attestation Overhead: Attestation overhead is measured using two primary metrics: *Latency overhead (%)*, defined as the additional time required to generate one token compared to pure REE inference, and *Energy overhead (%)*, calculated by the additional energy consumed to generate one token compared to pure REE inference.

### 5.3 Watermark Extraction Accuracy

LLMs watermarked with AttestLLM can reliably pass the attestation with high watermark extraction accuracy, whereas the execution of non-watermarked ones will be aborted due to their low accuracy. As shown in Table 2, we encode watermarks onto large language models from the Llama, Qwen, and Phi families, following the methodologies detailed in Section 4.2. As shown, for LLMs quantized into INT4 and INT8, AttestLLM achieves a high extraction accuracy of 100%, ensuring the attestation framework can reliably identify authorized LLMs. For non-watermarked models, their watermark extraction rates are lower, and their execution attempts will be aborted.

### 5.4 Watermarked LLM Performance

Encoding watermarks onto LLMs will not degrade their quality and performance. We evaluate AttestLLM's performance on the generated text fluency using perplexity and predictive capability using zero-shot accuracy. As shown in Table 2, for INT4 quantization, compared to non-watermarked models, AttestLLM averagely degrades the Perplexity (PPL) by 0.62% and Zero-shot Accuracy by 0.09%. As for INT8 quantization, compared to non-watermarked models, AttestLLM degrades the PPL and Zero-shot Accuracy by 0.07% and 0.15%, respectively. These degradations are minimal and do not produce noticeable differences when users interact with the watermarked models [31].

### 5.5 Attestation Overhead

We evaluate the attestation overhead of LLMs in Table 1. The reported overhead is measured by the proportional latency/energy increase over pure REE execution. The attestation interval in Section 4.3.2 can be specified by user. Without loss of generality, we set the attestation interval to be $f = 100$ unless specified. This means the attestation is invoked for

every 100 generated tokens. We leverage the TEE-shield inference [15, 34, 64] as our baseline, in which the LLM is split according to secure memory capacity, and each segment is executed inside the TEE. It incurs frequent secure/normal world switches and data transfers, leading to substantial latency and energy overhead.

As shown in Figure 3 and Figure 4, we report the attestation latency and energy overhead of AttestLLM and TEE-shield inference. TEE-shield inference suffers from substantial overhead, with latency rising by 300–700% (averaging 484.3% for INT4 and 632.7% for INT8) and energy by 180–270% (averaging 192.5% for INT4 and 248.9% for INT8). In contrast, AttestLLM imposes only marginal additional costs over REE execution, with average overheads of 19.6% (INT4) and 16.8% (INT8) in terms of latency, and 14.5% (INT4) and 13.7% (INT8) for energy.
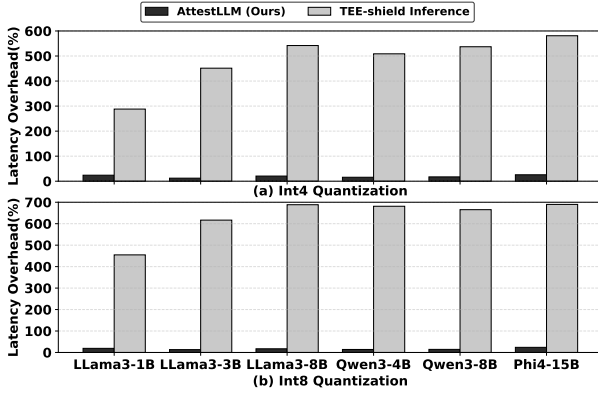


**Figure 3.** Attestation latency overhead (%) of models from Llama, Qwen, and Phi families when quantized into INT4 and INT8. Lower is better.
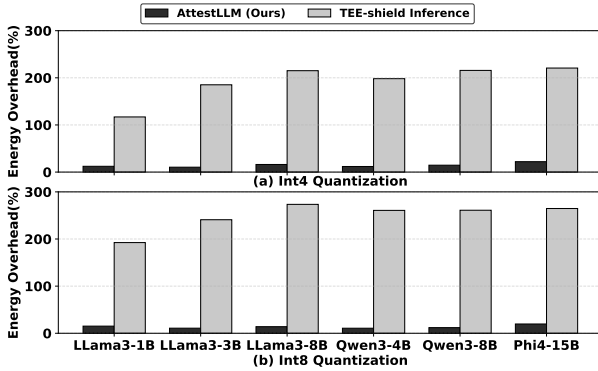


**Figure 4.** Attestation energy overhead (%) of models from Llama, Qwen, and Phi families when quantized into INT4 and INT8. Lower is better.

## 5.6 Robustness against Attacks

As in Section 3.2, the adversary is an end-user or an untrusted third party along the supply chain. They have full access to the hardware device and attempt to run unauthorized models on the platform that violate hardware IP. However, they cannot control the TEE environment used for secure verification. Here, we use Llama3-8B [21] model as a proof-of-concept demonstrating the robustness of AttestLLM. The adversary has prior knowledge of the model architecture and the watermarking algorithm. However, he/she does not know model weights, and it is hard to recover the weights because they are shuffled and randomized in Section 4.3.3. They also do not have access to encrypted watermark keys and watermark seed.

***Model replacement attack.*** We consider the model running on the hardware platform that uses the Llama3-8B architecture, but the provider fine-tunes the model to meet their specialized requirements. The adversary has knowledge of the architecture and thus uses a similar model to substitute the original LLM. Here, we consider the original LLM to be a Llama3-8B model quantized into 8-bit, whereas the unauthorized model is distilled from DeepSeek-R1 output but has a Llama3-8B architecture [22]. As in Table 3, AttestLLM can effectively tell if the model has been watermarked or not and yield low watermark extraction accuracy for unauthorized models. Thus, AttestLLM is resilient toward model replacement attacks.

**Table 3.** Model replacement attack performance.

| Model | Status | WM Extract. Acc |
|---|---|---|
| Llama-8B | watermarked | 100% |
| DeepSeek-R1 Distilled | unauthorized | 78.47% |

***Watermark forgery attack.*** The adversary forges another set of signatures and embeds it onto the unauthorized model. As such, the adversary attempts to bypass the attestation with the forged signature. However, it can be hard to forge the watermark due to the following reasons: (i) the project matrix $WM$ and watermark signatures $B$ are kept confidential and encrypted in secure memory. It can be hard for the adversary to reconstruct the watermark signatures; (ii) The marked LLMs are stored in shuffled and randomized memory. It can be prohibitively expensive to brute-force search this information.

***Security Analysis.*** Due to memory-layout randomization in Section 4.3.3, we consider a stealthy adversary tried to brute-force search and tampered with $t=2$ transformer blocks in the on-device LLM. AttestLLM triggers the attestation every $f=100$ generated tokens. For a conversational session producing $m=1000$ tokens, this results in $r=10$ attestation rounds. Using Equation 5, the cumulative evasion probability within a single conversation ranges from $2.2\times10^{-1}$ in the

most vulnerable case (Llama3-3B, $L$=28, $k$=2) to $3.7\times10^{-2}$ in the most secure case (Phi4-15B, $L$=40, $k$=6). Over multiple conversations, the probability decreases rapidly: for Llama3-3B, it drops to $2.77\times10^{-7}$ after 10 sessions, enabling AttestLLM to reliably detect and abort compromised LLMs.

We do highlight that this performance is achieved *without* introducing large latency or energy overhead as shown in Fig. 3 and Fig. 4. Moreover, our dynamic attestation strategy is tunable: by slightly increasing $k$ or reducing $f$, the evasion probability can be further decreased, enabling flexible security–performance trade-offs.

### 5.7 Ablation Studies and Analysis
In this section, we perform additional ablation studies and analysis on AttestLLM under different settings. The performance is evaluated on the Llama3-1B [21] model quantized into INT4 unless specified.

#### 5.7.1 Watermarking before/after quantization.
We compare AttestLLM's two-stage watermark insertion performance with only inserting watermarks before/after quantization in Table 4. As seen, inserting watermarks solely before/after quantization would degrade the performance. Post-quantization insertion introduces more degradation because zeroth-order optimization must approximate gradients from a discrete weight distribution, causing larger weight updates than gradient-based optimization to embed the same signature. Conversely, pre-quantization insertion causes less quality degradation but suffers from reduced extraction accuracy, since quantization perturbs the activation outputs and weakens the embedded signal. By combining both stages, AttestLLM achieves a better balance, preserving model quality while maintaining high watermark extraction accuracy.

**Table 4.** Performance of inserting watermark before/after quantization.

|  | PPL | Zero-shot Acc | WM Extract. Acc |
|---|---|---|---|
| AttestLLM | 9.91 | 60.90% | 100% |
| Before Quantization | 10.26 | 59.58% | 97.19% |
| After Quantization | 10.84 | 58.51% | 100% |

#### 5.7.2 Attestation overhead under different hyperparameters.
AttestLLM allows configuration of the attestation interval $f$ and the number of attestation blocks per request $k$, enabling a trade-off between security strength and attestation efficiency. As discussed in Section 4.3.2, smaller $f$ and larger $k$ provide stronger protection for on-device LLMs but incur higher latency and energy overhead. We benchmark this trade-off in Fig. 5 using Qwen3-4B. It shows that overhead decreases linearly with larger $f$ and increases with larger $k$. Under our default setting ($f$=100, $k$=4), generating $m$=1000 tokens ($r$=10 rounds) bounds the cumulative evasion probability to ~0.10, achieving a practical balance between

security and efficiency without requiring higher $k$ for more blocks to verify watermark during attestation.
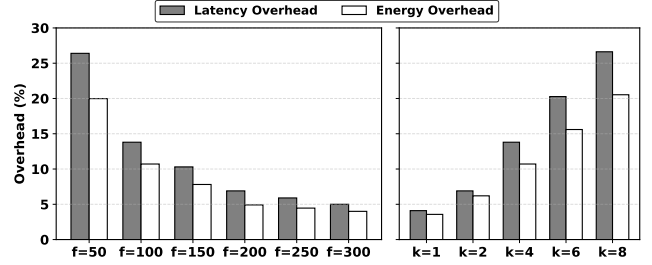


**Figure 5.** Left: As the attestation interval $f$ increases (less frequent checks), both latency and energy overhead decrease. Right: As the number of blocks attested per request $k$ grows, latency overhead and energy overhead increase.

#### 5.7.3 Attestation overhead with/without efficiency optimization.
Without attestation pipeline optimization, secure copy of sampled blocks and subsequent secure verification are executed in sequence. As detailed in Section 4.3.2, we optimize the process by combining a *pipelined verification* and *pre-allocated* buffers that eliminate allocator churn and TLB pressure. Within the secure enclave, multiple blocks are verified in parallel. The results are in Table 5, under the INT8 quantization of Qwen3-4B, latency overhead drops from 17.8% to 13.9%, and energy overhead from 14.3% to 10.7%. Thus, including those efficiency optimizations can reduce the overhead in attestation.

**Table 5.** Attestation overhead with/without efficiency optimizations.

| Attestation Overhead | w/o optimization | w optimization |
|---|---|---|
| Latency | 17.8% | 13.9% |
| Energy | 14.3% | 10.7% |

#### 5.7.4 Overhead breakdown.
Fig. 6 shows the breakdown of attestation overhead, including secure copy, decryption (which in our breakdown also includes trigger dataset retrieval), and secure verification. As seen, secure verification dominates the cost since it performs large-dimensional matrix multiplication operations to decode the watermark, which is computation- and memory-intensive. Secure copy through the virtio-pmem interface and decryption take ~10% of the total attestation time, benefiting from hardware acceleration on ARMv8-A cores with Cryptographic Extensions [55].

#### 5.7.5 Watermarking signature distribution among blocks.
We show the distribution of signature length among the blocks in Figure 7. As shown, the signatures are almost evenly distributed among the transformer blocks, except
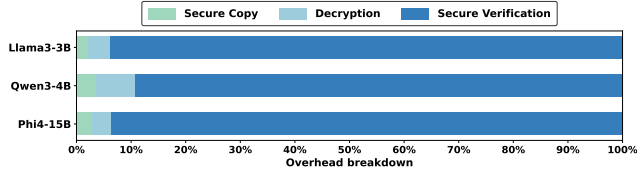
**Figure 6.** Attestation overhead, which is decomposed into three components: secure copy, decryption, and secure verification.

for some dynamics that depend on the absolute activation of each transformer block. The latter blocks are assigned fewer signatures, as they are critical for predicting the next generated token.
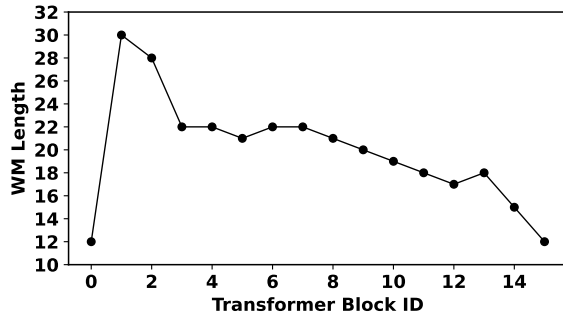


**Figure 7.** The distribution of watermarking signature length among blocks.

## 6 Discussion

### 6.1 Porting across SoCs

AttestLLM is designed for broad deployment across heterogeneous edge platforms. The watermarking mechanism is model-agnostic and compatible with diverse LLM architectures and quantization settings, ensuring high fidelity and robustness across models.

AttestLLM leverages pKVM [1, 19] to instantiate a lightweight, virtualized secure enclave without requiring modifications to the secure world (EL3) or the Android boot chain. This architecture is fully compatible with ARM-based SoCs that support the standard virtualization extension at EL2, and those implementing non-VHE (nVHE) mode. Such support is already present in a wide range of commercial SoCs, including Qualcomm Snapdragon XR2, Snapdragon 8 Gen series, and other Cortex-A55/A76+ platforms. AttestLLM also supports attestation with accelerators (e.g., NPUs), where enclaves are securely assigned accelerators by configuring Stage-2 translations along with per-device IOMMU contexts. AttestLLM can be readily ported across platforms and integrated into consumer-facing products such as smart home assistants, AR/VR headsets, and other privacy-sensitive edge

applications, providing a scalable and practical solution for secure on-device LLM deployment. Furthermore, our design is compatible with ARMv9 Confidential Computing Architecture (CCA) [6]. AttestLLM can directly leverage the Realm Management Extension (RME) to provide hardware-enforced isolation, further reducing the reliance on hypervisors and simplifying deployment on next-generation SoCs.

### 6.2 Attest ultra-large LLMs on cloud

AttestLLM is designed to attest LLMs for on-device use cases, so we validate its effectiveness on models up to 15 billion parameters, which is a size commonly used in on-device deployments. However, the same co-design attestation principles can extend to ultra-large LLMs running in cloud environments [18, 52]. In such contexts, clients may want to verify they are using the specific LLM version they paid for, or model providers can ensure that only authenticated third-party hardware or infrastructure runs their models. These use cases have become viable due to NVIDIA's confidential computing support in its Hopper and Blackwell architecture GPUs [65], which integrate trusted execution environments. While attesting ultra-large LLMs in cloud are orthogonal to our work, AttestLLM could potentially benefit future research in such security scenarios.

## 7 Conclusion

In this paper, we present AttestLLM, the first-of-its-kind attestation framework to avoid unauthorized LLMs running on the hardware platform and thus protecting the hardware IP of device vendors. We leverage a hardware-software co-design approach to scale models up to billion-scale while maintaining efficiency and robustness. AttestLLM leverages an optimization-based watermark insertion to encode signatures onto the LLM activation distributions, while maintaining the model quality and security. In online attestation, AttestLLM leverages a virtualization-based secure enclave, dynamic sub-block sampling, and verification computation overlapping with communication to reduce the attestation time and memory overhead. Extensive evaluations of on-device LLMs demonstrate the effectiveness of AttestLLM to authorize only legitimate models can be run on the target hardware platform.

## References

[1] [n. d.]. Protected Virtual Machine Firmware (pvmfw). https://android.googlesource.com/platform/packages/modules/Virtualization/+/HEAD/guest/pvmfw.

[2] Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, et al. 2024. Phi-4 technical report. *arXiv preprint arXiv:2412.08905* (2024).

[3] Ako Muhamad Abdullah et al. 2017. Advanced encryption standard (AES) algorithm to encrypt and decrypt data. *Cryptography and Network Security* 16, 1 (2017), 11.

[4] Tigist Abera, N Asokan, Lucas Davi, Jan-Erik Ekberg, Thomas Nyman, Andrew Paverd, Ahmad-Reza Sadeghi, and Gene Tsudik. 2016. C-FLAT: control-flow attestation for embedded systems software. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security.* 743–754.

[5] Apple Machine Learning Research. 2024. Introducing Apple's On-Device and Server Foundation Models. [Online]. Available: https://machinelearning.apple.com/research/introducing-apple-foundation-models.

[6] Arm Ltd. 2021. Arm Confidential Compute Architecture. https://www.arm.com/architecture/security-features/arm-confidential-compute-architecture.

[7] Huili Chen, Cheng Fu, Bita Darvish Rouhani, Jishen Zhao, and Farinaz Koushanfar. 2019. DeepAttest: An end-to-end attestation framework for deep neural networks. In *Proceedings of the 46th International Symposium on Computer Architecture.* 487–498.

[8] Pau-Chen Cheng, Wojciech Ozga, Enriquillo Valdez, Salman Ahmed, Zhongshu Gu, Hani Jamjoom, Hubertus Franke, and James Bottomley. 2024. Intel tdx demystified: A top-down approach. *Comput. Surveys* 56, 9 (2024), 1–33.

[9] Victor Costan and Srinivas Devadas. 2016. Intel SGX explained. *Cryptology ePrint Archive* (2016).

[10] Bita Darvish Rouhani, Huili Chen, and Farinaz Koushanfar. 2019. Deepsigns: An end-to-end watermarking framework for ownership protection of deep neural networks. In *Proceedings of the twenty-fourth international conference on architectural support for programming languages and operating systems.* 485–497.

[11] W. Deacon. 2020. Virtualization for the Masses: Exposing KVM on Android. In *KVM Forum.*

[12] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Marc'aurelio Ranzato, Andrew Senior, Paul Tucker, Ke Yang, et al. 2012. Large scale distributed deep networks. *Advances in neural information processing systems* 25 (2012).

[13] Ghada Dessouky, Shaza Zeitouni, Thomas Nyman, Andrew Paverd, Lucas Davi, Patrick Koeberl, N Asokan, and Ahmad-Reza Sadeghi. 2017. Lo-fat: Low-overhead control flow attestation in hardware. In *Proceedings of the 54th Annual Design Automation Conference 2017.* 1–6.

[14] Kazuki Egashira, Mark Vero, Robin Staab, Jingxuan He, and Martin Vechev. 2024. Exploiting llm quantization. *Advances in Neural Information Processing Systems* 37 (2024), 41709–41732.

[15] Akshay Gangal, Mengmei Ye, and Sheng Wei. 2020. HybridTEE: Secure mobile DNN execution using hybrid trusted execution environment. In *2020 Asian Hardware Oriented Security and Trust Symposium (AsianHOST).* IEEE, 1–6.

[16] Leo Gao et al. 2021. *A framework for few-shot language model evaluation.* doi:10.5281/zenodo.5371628

[17] Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple Contrastive Learning of Sentence Embeddings. In *Empirical Methods in Natural Language Processing (EMNLP).*

[18] In Gim, Caihua Li, and Lin Zhong. 2024. Confidential prompting: Protecting user prompts from cloud llm providers. *arXiv preprint arXiv:2409.19134* (2024).

[19] Google. 2024. *Android Virtualization Framework (AVF) Overview.* https://source.android.com/docs/core/virtualization

[20] Dibakar Gope, David Mansell, Danny Loh, and Ian Bratt. 2024. Highly optimized kernels and fine-grained codebooks for llm inference on arm cpus. *arXiv preprint arXiv:2501.00032* (2024).

[21] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).

[22] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025.

[23] Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025).

[23] Shwai He, Guoheng Sun, Zheyu Shen, and Ang Li. 2024. What matters in transformers? not all attention is needed. *arXiv preprint arXiv:2406.15786* (2024).

[24] Pin-Lun Hsu, Yun Dai, Vignesh Kothapalli, Qingquan Song, Shao Tang, Siyu Zhu, Steven Shimizu, Shivam Sahni, Haowen Ning, and Yanning Chen. 2024. Liger kernel: Efficient triton kernels for llm training. *arXiv preprint arXiv:2410.10989* (2024).

[25] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. *ICLR* 1, 2 (2022), 3.

[26] Dingji Li, Zeyu Mi, Yubin Xia, Binyu Zang, Haibo Chen, and Haibing Guan. 2021. Twinvisor: Hardware-isolated confidential virtual machines for arm. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles.* 638–654.

[27] Linyang Li, Botian Jiang, Pengyu Wang, Ke Ren, Hang Yan, and Xipeng Qiu. 2023. Watermarking llms with weight quantization. *arXiv preprint arXiv:2310.11237* (2023).

[28] Lei Li, Yongfeng Zhang, and Li Chen. 2023. Prompt distillation for efficient llm-based recommendation. In *Proceedings of the 32nd ACM international conference on information and knowledge management.* 1348–1357.

[29] Shuai Li, Kejiang Chen, Kunsheng Tang, Jie Zhang, Weiming Zhang, Nenghai Yu, and Kai Zeng. 2023. Turning Your Strength into Watermark: Watermarking Large Language Model via Knowledge Injection. *arXiv preprint arXiv:2311.09535* (2023).

[30] Shen Li, Liuyi Yao, Jinyang Gao, Lan Zhang, and Yaliang Li. 2024. Double-i watermark: Protecting model copyright for llm fine-tuning. *arXiv preprint arXiv:2402.14883* (2024).

[31] 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. *Proceedings of machine learning and systems* 6 (2024), 87–100.

[32] Zechao Lin, Sisi Zhang, Xingbin Wang, Yulan Su, Yan Wang, Rui Hou, and Dan Meng. 2025. LoRATEE: A Secure and Efficient Inference Framework for Multi-Tenant LoRA LLMs Based on TEE. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP).* IEEE, 1–5.

[33] Chang Liu and Jun Zhao. 2024. Resource allocation for stable llm training in mobile edge computing. In *Proceedings of the Twenty-fifth International Symposium on Theory, Algorithmic Foundations, and Protocol Design for Mobile Networks and Mobile Computing.* 81–90.

[34] Zhuang Liu, Ye Lu, Xueshuo Xie, Yaozheng Fang, Zhaolong Jian, and Tao Li. 2021. Trusted-DNN: A TrustZone-based adaptive isolation strategy for deep neural networks. In *Proceedings of the ACM Turing Award Celebration Conference-China.* 67–71.

[35] Zechun Liu, Changsheng Zhao, Igor Fedorov, Bilge Soran, Dhruv Choudhary, Raghuraman Krishnamoorthi, Vikas Chandra, Yuandong Tian, and Tijmen Blankevoort. 2024. Spinquant: Llm quantization with learned rotations. *arXiv preprint arXiv:2405.16406* (2024).

[36] Zechun Liu, Changsheng Zhao, Forrest Iandola, Chen Lai, Yuandong Tian, Igor Fedorov, Yunyang Xiong, Ernie Chang, Yangyang Shi, Raghuraman Krishnamoorthi, et al. 2024. Mobilellm: Optimizing sub-billion parameter language models for on-device use cases. In *Forty-first International Conference on Machine Learning.*

[37] Kangjie Lu, Wenke Lee, Stefan Nürnberger, and Michael Backes. 2016. How to Make ASLR Win the Clone Wars: Runtime Re-Randomization.. In *NDSS.*

[38] Hector Marco-Gisbert and Ismael Ripoll Ripoll. 2019. Address space layout randomization next generation. *Applied Sciences* 9, 14 (2019), 2928.

[39] Stephen Merity et al. 2016. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843* (2016).

[40] Myungsuk Moon, Minhee Kim, Joonkyo Jung, and Dokyung Song. 2025. ASGARD: Protecting On-Device Deep Neural Networks with Virtualization-Based Trusted Execution Environments. In *Proceedings 2025 Network and Distributed System Security Symposium*.

[41] Avanika Narayan, Dan Biderman, and Christopher Re. [n. d.]. Proof-of-Concept for Private Local-to-Cloud LLM Chat via Trusted Execution Environments. In *ES-FoMo III: 3rd Workshop on Efficient Systems for Foundation Models*.

[42] Institute of Electrical, Electronics Engineers. Computer Society. Standards Committee, and David Stevenson. 1985. *IEEE standard for binary floating-point arithmetic*. IEEE.

[43] OP-TEE core maintainers. 2024. OP-TEE trusted OS. [Online]. Available: https://github.com/OP-TEE/optee_os.

[44] Sandro Pinto and Nuno Santos. 2019. Demystifying arm trustzone: A comprehensive survey. *ACM computing surveys (CSUR)* 51, 6 (2019), 1–36.

[45] Qualcomm. 2022. Whitepaper: The future of AI is hybrid. *Qualcomm blog* (2022).

[46] Tahmid Noor Rahman, Nusaiba Khan, and Zarif Ishmam Zaman. 2024. Redefining Computing: Rise of ARM from consumer to Cloud for energy efficiency. *arXiv preprint arXiv:2402.02527* (2024).

[47] Marco Rando, Cesare Molinari, Lorenzo Rosasco, and Silvia Villa. 2023. An optimal structured zeroth-order algorithm for non-smooth optimization. *Advances in Neural Information Processing Systems* 36 (2023), 36738–36767.

[48] Mohamed Sabt, Mohammed Achemlal, and Abdelmadjid Bouabdallah. 2015. Trusted execution environment: What it is, and what it is not. In *2015 IEEE Trustcom/BigDataSE/Ispa*, Vol. 1. IEEE, 57–64.

[49] AMD Sev-Snp. 2020. Strengthening VM isolation with integrity protection and more. *White Paper, January* 53, 2020 (2020), 1450–1465.

[50] Zhuojia Shen, Komail Dharsee, and John Criswell. 2022. Randezvous: Making randomization effective on MCUs. In *Proceedings of the 38th Annual Computer Security Applications Conference*. 28–41.

[51] James C Spall. 2000. Adaptive stochastic approximation by the simultaneous perturbation method. *IEEE transactions on automatic control* 45, 10 (2000), 1839–1853.

[52] Jianchang Su and Wei Zhang. [n. d.]. Runtime Attestation for Secure LLM Serving in Cloud-Native Trusted Execution Environments. In *Machine Learning for Computer Architecture and Systems 2025*.

[53] The QEMU Project Developers. 2024. virtio-pmem. https://www.qemu.org/docs/master/system/devices/virtio-pmem.html.

[54] torchao. 2024. *TorchAO: PyTorch-Native Training-to-Serving Model Optimization*. https://github.com/pytorch/torchao

[55] KC Wang. 2023. ARMv8 Architecture and Programming. In *Embedded and Real-Time Operating Systems*. Springer, 505–792.

[56] Luca Wilke, Jan Wichelmann, Anja Rabich, and Thomas Eisenbarth. 2023. Sev-step: A single-stepping framework for amd-sev. *arXiv preprint arXiv:2307.14757* (2023).

[57] Zhuofeng Wu, He Bai, Aonan Zhang, Jiatao Gu, VG Vydiswaran, Navdeep Jaitly, and Yizhe Zhang. 2024. Divide-or-Conquer? Which Part Should You Distill Your LLM? *arXiv preprint arXiv:2402.15000* (2024).

[58] Jiajun Xu, Zhiyuan Li, Wei Chen, Qun Wang, Xin Gao, Qi Cai, and Ziyuan Ling. 2024. On-device language models: A comprehensive review. *arXiv preprint arXiv:2409.00088* (2024).

[59] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388* (2025).

[60] Rongjie Yi, Liwei Guo, Shiyun Wei, Ao Zhou, Shangguang Wang, and Mengwei Xu. 2023. Edgemoe: Fast on-device inference of moe-based large language models. *arXiv preprint arXiv:2308.14352* (2023).

[61] Shuguang Yuan, Xingyu Su, Peizhuo Lv, Weiji Xue, Jing Yu, Xiaojie Zhu, and Chi Chen. 2025. An Efficient White-box LLM Watermarking for IP Protection on Online Market Platforms. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*. 3680–3691.

[62] Chen Zhang, Kuntai Du, Shu Liu, Woosuk Kwon, Xiangxi Mo, Yufeng Wang, Xiaoxuan Liu, Kaichao You, Zhuohan Li, Mingsheng Long, et al. 2025. Jenga: Effective Memory Management for Serving LLM with Heterogeneity. *arXiv preprint arXiv:2503.18292* (2025).

[63] Ruisi Zhang and Farinaz Koushanfar. 2024. EmMark: Robust watermarks for IP protection of embedded quantized large language models. In *Proceedings of the 61st ACM/IEEE Design Automation Conference*. 1–6.

[64] Ziqi Zhang, Chen Gong, Yifeng Cai, Yuanyuan Yuan, Bingyan Liu, Ding Li, Yao Guo, and Xiangqun Chen. 2024. No privacy left outside: On the (in-) security of tee-shielded dnn partition for on-device ml. In *2024 IEEE Symposium on Security and Privacy (SP)*. IEEE, 3327–3345.

[65] Jianwei Zhu, Hang Yin, Peng Deng, Aline Almeida, and Shunfan Zhou. 2024. Confidential Computing on NVIDIA Hopper GPUs: A Performance Benchmark Study. *arXiv preprint arXiv:2409.03992* (2024).