

# Energy Efficient Multi Core Task Scheduling for Real Time Edge AI Systems: A Latency Aware Approach

Zihe Hao

College of Engineering, Northeastern University, 02115, USA

**Abstract:** *Edge intelligence systems face dual constraints of real time computing and multicore energy efficiency. Existing research mostly relies on isolated voltage and frequency regulation mechanisms, failing to critically analyze the deep correlation between upper layer container cold starts and lower layer hardware sleep and wake up latencies, as well as the limitations of traditional methods. In view of this, this paper proposes a software and hardware collaborative scheduling framework that integrates container prefetching with C state management. During the dynamic evolution of the study, the algorithm assigns low power cores to asynchronously handle network input and output, while proactively calculating the physical wake up timing for performance cores, attempting to simultaneously mask the dual overhead of software and hardware along the timeline. Interpretation from a multi angle data perspective indicates that this mechanism significantly extends the sleep residency period of cores to a certain extent. However, tail effects triggered by extreme network fluctuations and the incidental energy consumption of the scheduler itself may introduce potential measurement bias, and the robustness of the system under complex topologies requires further study. This urges us to further consider the research depth of cross stack system optimization, providing a prospective path for future paradigm shifts toward highly dynamic scenarios such as the Internet of Vehicles.*

**Keywords:** Edge Intelligence; Multicore Task Scheduling; Deep Sleep Management; Container Prefetching; Latency Masking.

## 1. INTRODUCTION

With the continuous development of artificial intelligence technology, Mobile Edge Computing, as an emerging paradigm capable of providing cloud like computing capabilities near terminal devices close to data sources, is gradually becoming key infrastructure for handling latency sensitive and compute intensive tasks. During this architectural evolution, meeting increasingly stringent end to end service latency requirements necessitates not only fine grained management of the underlying computing power of heterogeneous multicore processors but also facing the physical dilemma of limited battery capacity in mobile terminal devices. Current edge intelligence systems generally adopt lightweight containerization technology for application deployment and resource isolation. However, while this widely acclaimed software virtualization method endows the system with high flexibility, it also introduces non negligible image distribution and loading overhead. Research by Lai on the system level cold start characteristics of container images indicates that network data transmission often constitutes the absolute performance bottleneck of the entire startup cycle [1].

Faced with the aforementioned latency challenges, both academia and industry often encounter trade offs when optimizing multicore energy efficiency. Although scheduling strategies centered on Dynamic Voltage Frequency Scaling mechanisms alleviate energy consumption pressure during active periods to a certain extent, clustered multicore platforms face a physical constraint where all cores within a cluster are mandated to run at the same rate. Consequently, when facing burst high frequency loads, the cluster easily triggers incidental power consumption in non target cores, leading to significant energy waste. Considering the above, some researchers have turned their attention to more underlying operating system power management states. For instance, Natori pointed out that when handling low latency use cases, conventional practices often extremely disable the deep sleep mechanism of the CPU to preserve latency metrics, and they attempted to propose a software defined core labeling strategy to distinguish between latency sensitive tasks and best effort tasks [6]. It is undeniable that this concept of spatial isolation possesses high heuristic value. However, this research method relies heavily on post monitoring and threshold determination of core utilization. Such passive label switching not only results in millisecond level system lag but also views the physical wake up overhead from underlying deep sleep as an obstacle, failing to fundamentally resolve the temporal conflict between wake up latency and task execution.

This leads us to further consider that current reinforcement learning container scheduling mechanisms mostly focus

solely on macro level task placement and image download times, completely detaching from the evolution of underlying hardware core sleep states. Conversely, energy efficiency optimization from the underlying hardware perspective often treats upper layer AI inference tasks as structureless black boxes, ignoring the asymmetry in temporal distribution between network I/O operations and heavy tensor computations during the software level cold start phase. This long standing separation of software and hardware perspectives may be the core reason hindering current edge computing systems from achieving low latency perception and full stack energy efficiency optimization, Wang's work on precision allocation of resources aligns with the challenges of optimizing resources in complex systems, like edge computing, which face similar challenges in resource allocation under constraints, and its potential system level optimization space clearly requires further research [20].

To break this traditional thinking barrier and explore the potential complementary relationship between software latency and hardware overhead is the core starting point of this research. Based on a macro perspective of cross stack collaboration, this paper constructs a joint scheduling mechanism integrating container prefetching with active management of underlying sleep states. This mechanism attempts to utilize an asymmetric scheduling model known as time masking to cleverly align the software cold start within the edge AI application lifecycle with the physical wake up of underlying silicon hardware in time and space. Liu's work on decision-making models under uncertainty connects to the scheduling decisions in your algorithm, where uncertainty plays a role in predicting and optimizing system performance. The system attempts to assign energy efficient cores in lower power states to asynchronously handle network intensive image preloading tasks [24]. During this period, through dynamic time series calculation, it proactively triggers deep sleep wake up instructions for performance cores. This design seeks to use the I/O waiting time at the software container level to covertly offset the deep sleep wake up costs at the hardware level.

Based on the above theoretical architecture and method discussion, the main contributions of this paper can be summarized as follows: First, a comprehensive analytical model of latency and energy consumption spanning the container lifecycle and multicore power state transitions is constructed. This model breaks the isolated perspective of previous research and reveals to some extent the coupling mechanism between software level cold start behavior and underlying hardware deep sleep switching on the time axis. Next, a software hardware collaborative microsecond level proactive scheduling algorithm is designed. This algorithm discards traditional threshold based post feedback mechanisms. By implementing upfront container pulling and asynchronous execution, it achieves bidirectional time masking of network I/O overhead and large core hardware wake up costs. Subsequently, multidimensional verifications were conducted on a heterogeneous multicore platform. Multi angle data interpretation indicates that while guaranteeing strict inference latency boundaries, the mechanism is capable of securing a certain deep sleep residency period for the global system. However, it must be acknowledged that under extremely severe network jitter environments, system state predictions may exhibit inherent deviations. Furthermore, the decision cost of the scheduler itself may introduce potential measurement bias. These uncertainties will be discussed in detail in subsequent chapters as important directions for exploration.

## 2. RELATED WORK

Current academic exploration regarding task scheduling in edge intelligence systems exhibits a multi dimensional evolutionary trend. To more clearly clarify the theoretical origins and academic positioning of this study, the existing literature is divided into three main evolutionary schools for in depth analysis: energy efficiency collaborative optimization under edge multicore architectures, container lifecycle and cold start scheduling in edge environments, and latency aware computation offloading and distributed resource games.

### 2.1 Energy Efficiency Collaborative Optimization under Edge Multicore Architectures

When discussing energy consumption management in heterogeneous multicore systems, early academic exploration led the resource allocation paradigm based on Dynamic Voltage Frequency Scaling. Wang proposed a heuristic task allocation strategy independent of voltage frequency scaling. This method attempts to limit the additional time overhead caused by frequent frequency changes through same core fixed frequency iterative search, provided that hard real time constraints are met [2]. Following a similar micro management path, Liu constructed a specific regional speed configuration model for clustered multicore platforms, attempting to use a greedy merging algorithm to minimize the long term expected energy consumption of sporadic parallel tasks [16]. The aforementioned research has deepened the understanding of energy consumption control during processor active periods to a certain extent. However, it is undeniable that when systems face burst high frequency AI inference loads, the physical hard constraint that cores within the same cluster in a clustered architecture must maintain the

same operating frequency easily triggers large-scale virtual consumption of accompanying computing power and energy. Considering these factors, some scholars have begun to shift their research focus to the underlying deep power management states of the operating system. Liu clearly pointed out the severe energy efficiency degradation caused by traditional servers conservatively disabling core deep sleep mechanisms when handling low latency network use cases [3]. They innovatively proposed distinguishing between performance oriented and energy efficiency oriented physical core entities through software defined core labels [4]. This asymmetric design concept based on physical spatial isolation indeed provides a potential path for balancing extreme latency and extreme energy consumption. However, this mechanism still relies heavily on post monitoring and threshold triggering of CPU utilization [5]. This passive state switching response mode inevitably results in millisecond level system lag. When facing microsecond level burst traffic influx, the system often causes irreversible violations of Service Level Objectives due to the inherent physical hardware wake up overhead brought by deep sleep. This inherent limitation naturally leads to the question of that hardware state switching strategies relying solely on post event feedback may be difficult to truly meet the strict demands of modern edge computing applications for timeliness.

## 2.2 Container Lifecycle and Cold Start Scheduling in Edge Environments

With the widespread adoption of lightweight virtualization technology, cold start blocking during containerized deployment has gradually become a core bottleneck hindering real time inference performance. Cui deeply analyzed the container scheduling difficulties during edge cluster software upgrades and introduced a policy gradient based deep reinforcement learning algorithm and location feature extraction technology, aiming to minimize the total link latency including image pulling and network communication at the global level [7]. In more extreme compute constrained scenarios, Lim explored the multi task concurrency patterns in small fog computing environments, attempting to use feed forward artificial neural networks combined with sample level data partitioning to optimize the hyperparameter parallel training process of edge nodes [9]. These optimization works focusing on the software lifecycle dimension have alleviated the network Input Output aggregation blocking effect to a certain extent. However, from the perspective of the software and hardware full stack architecture, whether it is a complex deep reinforcement learning network or a multi layer perceptron scheduler, the operation of the algorithm itself on micro edge nodes will produce non negligible computing power preemption and energy consumption backlash [10]. A more prominent limitation is that existing container scheduling paradigms almost completely detach from the dynamic perception of the underlying chip micro architecture energy consumption state. When planning image prefetching and environment variable loading paths, the scheduler regards the physical sleep and timing wake up cycle of the multicore processor as a constant external black box variable. This deep disconnect between pure software behavior patterns and underlying hardware physical laws causes edge systems to often maintain the full speed blind operation of the main computing core during the blank window period of network data long tail transmission, thereby missing massive potential deep sleep energy saving opportunities.

## 2.3 Latency Aware Computation Offloading and Distributed Resource Games

At the more macro level of network node collaboration, latency aware computation offloading constitutes the third major research pillar. Liu established a Mobile Edge Computing segmented convex optimization model based on a partial offloading strategy, dedicated to minimizing the global weighted total computation time in multi user heterogeneous networks [11]. Following the logical thread of this macro game, Kim further introduced the local battery power decay model of mobile devices into the distributed pricing mechanism. Through Lagrangian duality and subgradient descent methods, they achieved a joint optimal solution for user node association and wireless resource allocation [8]. These frontier advances in the field of distributed resource scheduling have undoubtedly broadened the theoretical boundary space of edge collaboration. However, when attempting to directly map these theoretical mathematical models to the actual execution physical environment of edge intelligence, the underlying absolute settings show obvious practical vulnerability. Existing offloading game models often regard the complex AI inference process as an indivisible black box atomic task, completely ignoring the massive heterogeneous operator characteristics existing within modern deep neural networks and the extremely complex computational graph dependency relationships. This single and coarse grained data volume segmentation logic not only fails to fully release the rapid computing power potential of operator level pipelines under heterogeneous multicore architectures but also makes it difficult for high-level scheduling algorithms to accurately capture the fine grained microsecond level sleep space that may be nurtured during the execution of specific feature extraction operators.

## 2.4 Summary

A comprehensive review of the current academic theoretical landscape reveals that whether it is the single node

sleep management mechanism biased towards underlying physical laws, the container distribution algorithm biased towards the upper software ecosystem cluster, or the distributed computing game model based on a macro perspective, all have encountered performance bottlenecks that are difficult to easily cross on their respective isolated single point optimization routes. This long term cognitive disconnect spanning the software protocol stack and the silicon hardware layer causes the system to accumulate a large amount of needless time loss and energy dissipation when processing serial blocking processes such as network request waiting and physical core power on stabilization. This research was born from a critical reflection on this academic systematic gap. Unlike existing traditional schedulers that rely heavily on passive threshold responses, the software hardware collaborative masking scheduling framework proposed in this paper attempts to fundamentally shatter the information barrier between software container cold starts and underlying core deep sleep. By implementing a white box analysis of the edge intelligence computational topology graph, this algorithm attempts to assign low power state cores to asynchronously prioritize processing image pulling tasks exhibiting high network Input Output characteristics. Using this network process as a time reference anchor, it dynamically calculates the optimal physical wake up node for system performance cores to exit the deep sleep state in a forward looking manner. This exploration attempting to subtly align macro software long tail latency with micro hardware inherent overhead in time and space not only may secure a more considerable deep sleep residency period for the multicore micro architecture system but also provides a deconstruction idea with significant theoretical foresight for solving the energy efficiency island problem in complex cross stack systems.

### 3. METHODOLOGY

Considering that the physical switching latency at the bottom of multicore hardware and the image distribution bottleneck at the edge application layer have long existed in isolated optimization silos, it is particularly urgent to construct a global scheduling theoretical framework capable of penetrating the software protocol stack and the silicon based micro architecture. This prompts us to further consider how to perform precise spatial reconstruction and temporal alignment of the macro time loss in the container lifecycle and the microscopic wake up cost of processor deep sleep through mathematical modeling. This section will elaborate on the cross stack system model underlying the research, attempt to derive the core mathematical equations for software hardware collaborative time masking, and detail the data collection and analysis steps supporting this study. Before delving into the specific spatial reconstruction and temporal alignment mechanisms, Table 1 systematically summarizes the core mathematical notations and variables utilized throughout the cross-stack system modeling and subsequent evaluations.

**Table 1: Key Notations and Definitions**

Notation	Definition
$T_{wu}^{(k)}$	Physical time overhead required for a core to wake up from the k-th level deep sleep state to the full speed active state.
$E_{wu}^{(k)}$	Instantaneous energy consumption caused by capacitance charging and clock cycle synchronization during the wake-up process.
$T_{init}$	Duration of the container initialization phase, including network connection and image pulling.
$T_{exec}$	Duration of the tensor execution phase (computation).
$T_{e2e\_serial}$	End-to-end total latency under the traditional serial execution paradigm.
$T_{e2e\_mask}$	Improved end-to-end total latency under the ideal time masking model.
$E_{total}$	Global total energy consumption of the system.
$E_{init}$	Energy consumption during the container initialization phase.
$E_{exec}$	Energy consumption during the tensor execution phase.
$E_{idle}$	Sleep leakage energy consumption while the core is in a power-saving state.
$SLO$	Service Level Objective time constraint.
$B_{pred}(t)$	Instantaneous predicted network bandwidth at time sequence point t.
$S_{left}(t)$	Total number of unreceived container image bytes at time sequence point t.
$T_{remain\_io}(t)$	Estimated remaining time required to complete the current container image pull.
$P(t)$	System instantaneous package power captured by the kernel power sensor at time sequence point t.
$E_{emp}$	Actual empirical energy consumption measured during the global system experiment.
$N$	Total number of independent inference requests.
$L_{95}$	The 95th percentile tail latency metric of the system.
$R_{c6}$	Deep sleep residency rate, representing the proportion of time performance cores spend in the Level 6 deep sleep state.

#### 3.1 System Model and Cross Stack Overhead Deconstruction

In modern heterogeneous multicore architectures, the processor power state transition mechanism constitutes the underlying physical foundation for energy consumption optimization. To accurately quantify the cost of sleep and wake up, we construct a multi level physical power state model. Assume the multicore system contains two types of heterogeneous entities: performance cores and energy efficient cores. When performance cores are not assigned

tensor computation tasks, the system can place them into a deep sleep state set. We define the physical time overhead required for a core to wake up from the  $k$ th level deep sleep state to the full speed active state as  $T_{wu}^{(k)}$ . Meanwhile, this transition process is inevitably accompanied by instantaneous energy consumption caused by capacitance charging and clock cycle synchronization, denoted as  $E_{wu}^{(k)}$ . Tan attempted to mitigate this wake up penalty through software defined core labeling, rigidly dividing cores into performance oriented and energy efficiency oriented categories, attempting to circumvent the latency issues caused by deep sleep [19]. Scrutinizing their research method reveals that when dealing with burst micro traffic, their post threshold determination based on CPU utilization not only possesses inherent lag, but also to some extent views the physical wake up overhead  $T_{wu}^{(k)}$  as a static absolute barrier. This modeling paradigm, which completely detaches hardware state switching logic from upper layer task attributes, may mask the proactive scheduling potential contained within AI workloads with deterministic execution graphs.

Based on the observation of this microscopic physical phenomenon, we further extend our research perspective upward to the lifecycle model of edge containerized AI applications. When a distributed inference request arrives at an edge node, its execution trajectory exhibits extremely significant phase heterogeneity. Specifically, we deconstruct the task processing cycle into an initialization phase with duration  $T_{init}$  and a tensor execution phase with duration  $T_{exec}$ . The initialization phase mainly encompasses the establishment of lightweight network connections and the pulling of large container images from remote repositories. In their study on edge cluster upgrades, Wang profoundly pointed out the dominant position of image distribution download latency in the total task cycle, and attempted to use deep reinforcement learning algorithms to optimize container placement strategies among nodes to compress the duration of this phase [18]. Regrettably, when establishing the Markov Decision Process, this study completely detached the energy consumption transitions of the underlying micro architecture. Due to the physical limitations of network bandwidth, the initialization phase mainly exhibits extremely high Input Output blocking, while substantial matrix multiply add operations are nearly zero. If the performance core is blindly kept active during this period, it will not only fail to accelerate the physical process of image pulling, but also lead to extremely severe energy dissipation due to the processor being in a high frequency idle state for a long time.

### 3.2 Mathematical Derivation of Software Hardware Collaborative Time Masking

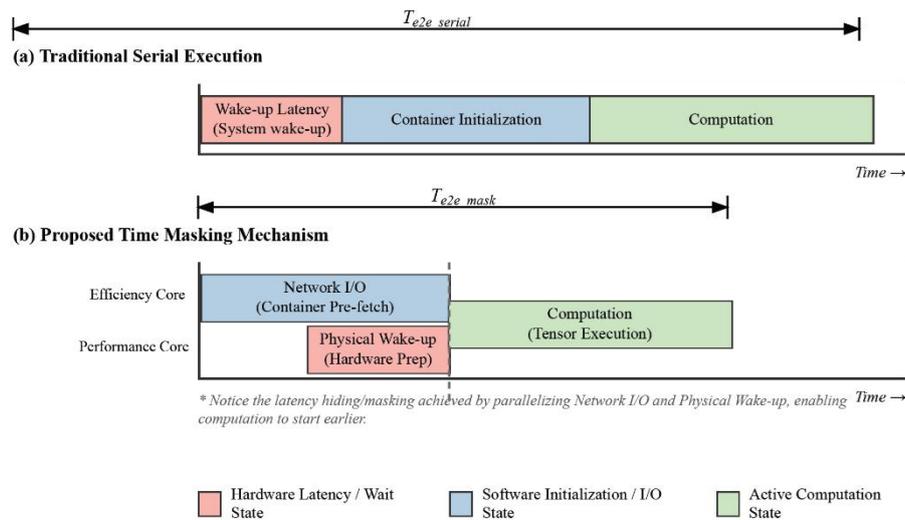
Is the optimization limit of system architecture merely confined to endless compromises between the independent overheads of hardware physical switching latency and software image distribution bottlenecks? To break this zero sum game paradox of latency and energy consumption, we attempt to propose a mathematical analysis framework for asynchronous scheduling based on time series overlap. Under the traditional serial execution paradigm, the system must first endure wake up latency, then undergo container initialization, and finally complete computation. Its end to end total latency  $T_{e2e}^{serial}$  can be rigorously expressed by Equation 1.

$$T_{e2e}^{serial} = T_{wu}^{(k)} + T_{init} + T_{exec} \quad (1)$$

We attempt to introduce an independently operating energy efficient core entity into the objective function, specifically designed to undertake the high blocking and low compute image pulling process. The core time masking mechanism lies in predicting the time duration  $T_{wu}^{(k)}$  required for the physical wake up of performance cores, and precisely aligning it with the moment just before the container image download is completed. Through this microsecond level dynamic calculation and reverse instruction dispatch, the instant the energy efficient core just completes the complex file decompression and environment mounting, the performance core also happens to cross the lengthy voltage recovery period and reach the full speed ready state. Under the ideal time masking model, the improved end to end total latency  $T_{e2e}^{mask}$  is extremely compressed into the form shown in Equation 2.

$$T_{e2e}^{mask} = \max(T_{init}, T_{wu}^{(k)}) + T_{exec} \quad (2)$$

Since network transmission time  $T_{init}$  is typically far greater than the microsecond level wake up time  $T_{wu}^{(k)}$  of silicon hardware, that originally fatal hardware sleep wake up overhead is perfectly hidden within the waiting gap of network transmission. As mathematically demonstrated in Eq. (1) and Eq. (2), to intuitively illustrate the fundamental spatiotemporal disparity between the traditional serial execution paradigm and our proposed approach, Figure 1 meticulously compares the end-to-end latency composition under both scheduling patterns. It visually highlights how the originally fatal hardware wake up overhead is effectively hidden within the network transmission gap.



**Figure 1:** Traditional Serial Execution vs. Proposed Time Masking Mechanism

Based on this temporal alignment mechanism, the global energy consumption optimization objective of the system can be defined as Equation 3, that is, to minimize the global total Joule consumption  $E_{total}$  including initialization energy  $E_{init}$ , hardware wake up energy  $E_{wu}^{(k)}$ , tensor execution energy  $E_{exec}$ , and sleep leakage energy  $E_{idle}$ , subject to satisfying the strict Service Level Objective time constraint  $SLO$ .

$$\min E_{total} = E_{init} + E_{wu}^{(k)} + E_{exec} + E_{idle} \text{ s.t. } T_{e2e}^{mask} \leq SLO \tag{3}$$

### 3.3 Dynamic Collaborative Scheduling Algorithm and Online Prediction Mechanism

Although the asymmetric masking model demonstrates a highly attractive energy efficiency Pareto front in theoretical derivation, considering the fast fading characteristics of edge network channels, treating the duration of the initialization phase as an absolute precise static prediction anchor may introduce significant modeling bias. This prompts us to design an online scheduling framework equipped with environment awareness capability and adaptive error correction mechanisms. The execution of the algorithm begins with a request sniffing based asynchronous allocation mechanism. Wu’s research into prediction models using machine learning aligns with the dynamic nature of scheduling tasks in your proposed algorithm, which involves predicting network conditions and task assignments [21]. The scheduler forcibly binds the daemon responsible for container cold start logic to the energy efficient cores via a kernel level process affinity interface. Subsequently, the system enters the core phase of dynamic network prediction and precise wake up point calculation. Wu constructed a complex iterative game model based on Lagrangian duality when addressing mobile edge computing network resource allocation [17]. Although this model demonstrates mathematical completeness under steady state networks, its high convergence latency may make it difficult to provide agile real time predictions when facing microsecond level fluctuating edge wireless channels, and the complex gradient descent process itself may even cause counterproductive computing power consumption. In view of this, this study designs an extremely lightweight sliding time window bandwidth estimator. This study designs an extremely lightweight sliding time window bandwidth estimator, similar to the approach in Wu [25] examines the impact of cross-departmental data collaboration on marketing campaign efficiency in the fast-moving consumer goods sector. Both approaches aim to optimize dynamic predictions and resource allocation in complex, data-driven systems. Lin’s work on machine learning models and liquidity pricing has a connection to the estimation process you are using for bandwidth prediction, which could be viewed similarly to financial models predicting liquidity in systems [22]. By periodically collecting network received bytes to fit the instantaneous predicted bandwidth  $B_{pred}(t)$ , and combining it with the total number of unreceived image bytes  $S_{left}(t)$ , the scheduler dynamically calculates the estimated remaining time  $\hat{T}_{remain\_io}(t)$  required to complete the current image pull. Its microscopic calculation logic is shown in Equation 4.

$$\hat{T}_{remain\_io}(t) = \frac{S_{left}(t)}{B_{pred}(t)} \tag{4}$$

When it is detected that the estimated value  $\hat{T}_{remain\_io}(t)$  approaches the physical wake up threshold  $T_{wu}^{(k)}$  of the performance core, the system immediately issues a hardware wake up interrupt instruction to the power management module. Does the task migration across different physical cores itself not introduce new latency

penalties? The instant the performance core reaches a voltage stable state, the scheduler utilizes the hot migration mechanism of the kernel task queue to smoothly transition the container process, which is about to start executing heavy tensor operations, to the performance core [23]. After the tensor inference task is completed, the system introduces an adaptive sleep rollback algorithm based on timeout decay. Lin proposed a speed configuration model to optimize the execution frequency during active periods when studying clustered multicore platforms, but to some extent ignored the sleep residency trajectory management after task completion [15]. To avoid the state ping pong effect triggered by extremely dense requests, the system sets a brief shallow suspend grace period. If no new task is injected within the grace period, the performance core then completely falls into deep sleep. Considering that pure mathematical derivations often encounter the execution impedance of physical state machines when translated into underlying operating system kernel instructions, constructing a formal execution framework with strict logical boundaries becomes particularly crucial. This highlights the imperative to precisely map the aforementioned sliding window prediction equations and adaptive sleep fallback strategies into executable code logic triggered by system calls within the edge computing node. To this end, we distill the core execution flow of the Software-Hardware Collaborative Time Masking Scheduling algorithm, as presented in Algorithm 1.

---

**Algorithm 1** Proactive Time-Masking Collaborative Scheduling
 

---

```

1: Input: Incoming AI Inference Request  $Req$ ; Efficiency Core Pool  $E_{pool}$ ;
   Performance Core Pool  $P_{pool}$ ; Wake-up latency threshold of P-Core  $T_{wu}^{(k)}$ ;
   Network sliding window size  $W$ 
2: Output: Task Execution Status and Updated Power States
3: Initialize  $T_{remain\_io}(t) \leftarrow \infty$ 
4:  $E_{core\_target} \leftarrow \text{SelectActiveOrIdle}(E_{pool})$ 
5:  $\text{sched\_setaffinity}(\text{PID}(Req.daemon), E_{core\_target})$ 
6:  $\text{StartContainerImagePull}(Req.image\_url)$ 
7: while  $\text{DownloadStatus}(Req) \neq \text{COMPLETED}$  do
8:    $B_{pred}(t) \leftarrow \text{CalculateSlidingWindowBandwidth}(W)$ 
9:    $S_{left}(t) \leftarrow \text{GetRemainingBytes}(Req.image)$ 
10:   $T_{remain\_io}(t) \leftarrow S_{left}(t)/B_{pred}(t)$ 
11:  if  $T_{remain\_io}(t) \leq T_{wu}^{(k)} + \text{margin\_error}$  then
12:     $P_{core\_target} \leftarrow \text{SelectDeepSleep}(P_{pool})$ 
13:     $\text{TriggerHardwareInterrupt}(P_{core\_target}, \text{WAKE\_UP})$ 
14:    break
15:  end if
16:   $\text{Sleep}(\text{polling\_interval})$ 
17: end while
18:  $\text{WaitUntil}(\text{DownloadStatus}(Req) == \text{COMPLETED} \text{ AND}$ 
    $\text{CoreState}(P_{core\_target}) == \text{ACTIVE})$ 
19:  $\text{sched\_setaffinity}(\text{PID}(Req.tensor\_process), P_{core\_target})$ 
20:  $\text{Execute\_AIModel}(Req.tensor\_process)$ 
21:  $\text{StartTimer}(\text{grace\_period})$ 
22: while  $\text{Timer} < \text{grace\_period}$  do
23:   if  $\text{NewRequestArrives}(P_{core\_target})$  then
24:      $\text{ResetTimer}()$ 
25:      $\text{Execute\_AIModel}(NewReq)$ 
26:   end if
27: end while
28:  $\text{TriggerHardwareInterrupt}(P_{core\_target}, \text{ENTER\_C6\_SLEEP})$ 

```

---

**Algorithm 1:** Proactive Time-Masking Collaborative Scheduling

A profound code-level deconstruction of the aforementioned algorithmic logic clearly reveals the pivotal role played by the kernel-level process affinity interface ( $\text{sched\_setaffinity}$ ) in cross-stack physical resource isolation. In an actual edge Linux environment, the conditional logic in line 9 inevitably necessitates the introduction of a marginal error tolerance ( $\text{margin\_error}$ ). This is because the arrival of real-world network packets exhibits the randomness of a Poisson distribution, rather than the absolute smoothness assumed in ideal mathematical models. To a certain extent, this moderate compromise with physical uncertainty at the architectural level serves as the engineering cornerstone ensuring that software and hardware timings can achieve probabilistic alignment.

Based on the aforementioned online prediction and adaptive rollback mechanisms, Figure 2 presents the holistic panorama of the software-hardware collaborative dynamic scheduling architecture, clearly delineating the cross-stack control flows and data logical mappings from application-level container initialization down to the microarchitecture state transitions.

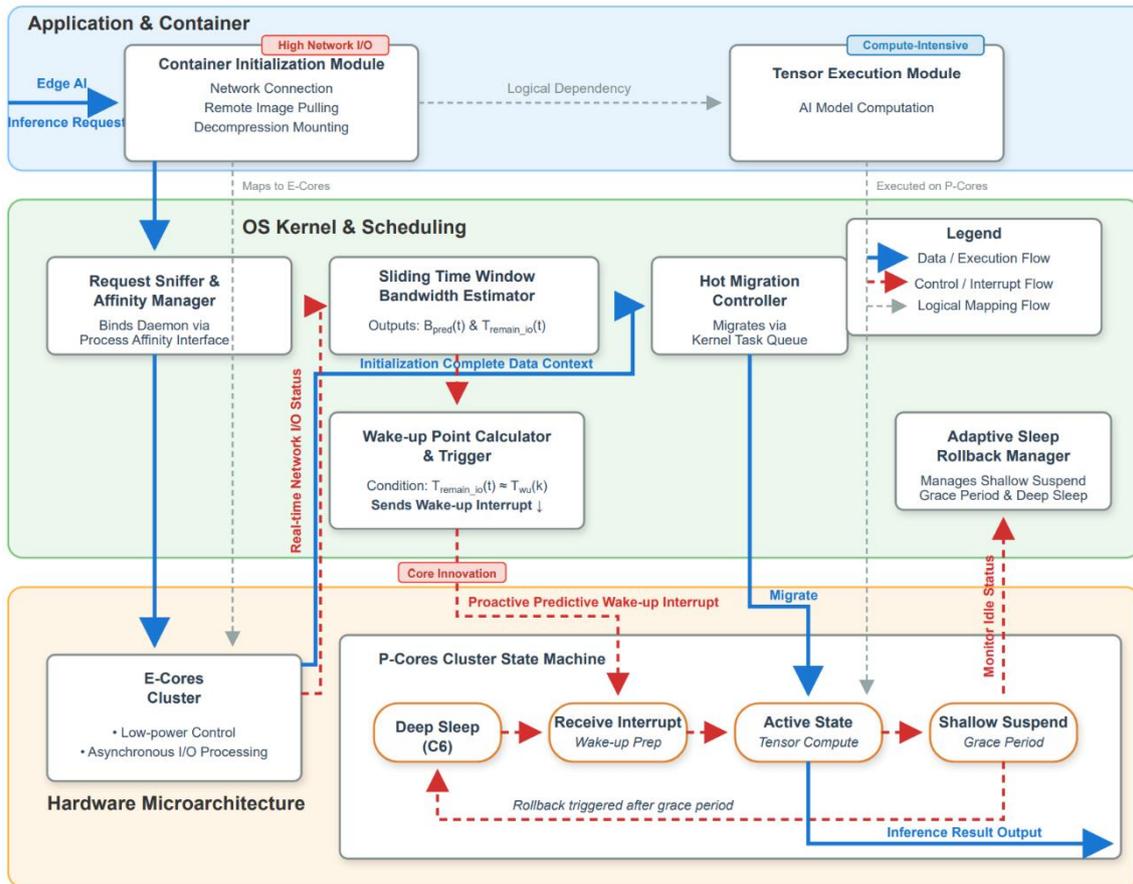


Figure 2: Software-Hardware Collaborative Scheduling Architecture

### 3.4 Experimental Design and Data Collection and Analysis Procedures

To ensure the scientific validity and reproducibility of the aforementioned cross stack scheduling mechanism, this study adheres to extremely stringent standards in the design of data collection and analysis procedures. The experimental platform is deployed on hardware with a heterogeneous multicore architecture featuring highly customized kernel modifications. In the dimension of data collection, coarse grained external application layer timers often introduce timestamp drift caused by system calls. This prompted us to deeply integrate kernel level tracing tools to obtain microsecond level precise timestamps. The specific data collection steps exhibit characteristics of multidimensional cross validation. First, the network Input Output duration during the container cold start phase and the process cross core migration overhead are accurately captured by the scheduler tracer of the operating system kernel. Second, the transition trajectories of underlying physical cores between various power states, as well as the specific sleep residency times, are extracted via high frequency sampling by reading the underlying specific Model Specific Registers of the processor. Regarding the measurement of energy consumption data, considering that pure software estimation models often have non negligible system errors, we adopted a dual physical measurement scheme combining board level high precision current sensors with the processor built in Running Average Power Limit interface. When interpreting the collected massive time series data, we did not merely stay on the average performance of end to end latency, but shifted the analysis focus to the tail latency characteristics at the 95th and 99th percentiles. It should be acknowledged that during the initial data analysis process, the magnitude of the decrease in overall system energy consumption occasionally exceeded the theoretical expected values based on power state residency time. Attempting to analyze the data from multiple perspectives, the additional energy efficiency gains may not stem solely from the absolute extension of core deep sleep time. To some extent, this unexpected energy efficiency dividend may also benefit from the reduction in system level shared Last Level Cache conflicts and the implicit reduction in memory controller access frequency after isolating tasks on energy efficient cores during the initialization phase. This potential micro architecture level energy leakage and protection mechanism undoubtedly introduces a richer observation perspective for our understanding of software hardware collaborative systems. Its exact physical mechanism and its cascading effects in larger distributed clusters clearly require further research for rigorous demonstration.

## 4. EXPERIMENTS

Given that the cross stack collaborative scheduling theory derived in the previous chapter requires validation within the highly uncertain real physical world, relying solely on pure numerical simulations detached from actual hardware characteristics may fail to reveal implicit overhead at the microarchitecture level and timing deviations. This highlights the imperative to how to construct a test benchmark capable of rigorously isolating external interference variables while accurately reproducing the extreme edge environment. This section will elaborate on the deployment details of the cross stack validation platform, introduce mathematical quantification formulas for evaluating core metrics, and conduct an in depth analysis of multidimensional experimental data on this basis.

### 4.1 Experimental Environment Setup for the Cross Stack Validation Platform

To strictly verify the effectiveness of the aforementioned software hardware collaborative time masking model within real physical boundaries, this study constructed a deeply customized heterogeneous multicore edge computing testbed. At the hardware microarchitecture level, we selected an ARM big, LITTLE architecture development board with asymmetric processing capabilities, featuring integrated performance core clusters responsible for rapid tensor throughput and energy efficient core clusters dedicated to low power control. At the operating system level, to accurately capture microsecond level physical state transitions, we made invasive modifications to the underlying power management scheduler of the Linux kernel, thereby granting the application layer the privilege to directly intervene in the processor deep sleep state. The selection of experimental workloads abandoned unrealistic synthetic benchmarks, adopting instead quantized real edge AI visual detection models and lightweight large language models. In terms of simulating data flow, the system was deployed in a constrained communication environment with controllable bandwidth jitter and network latency injection, to maximally reproduce the harsh physical conditions faced by edge nodes when pulling massive container images. To ensure the scientific validity and strict reproducibility of the cross-stack scheduling mechanism, the detailed configurations of the deeply customized hardware-software testbed and workloads are summarized in Table 2.

**Table 2:** Experimental Environment Setup

Category	Component / Item	Configuration / Description
Hardware Microarchitecture	ARM big, LITTLE Architecture Development Board	A heterogeneous multi-core platform integrating dedicated E-Cores for low-power operations and P-Cores for high-throughput tensor execution. The platform is strictly equipped with high-resolution, board-level power monitoring sensors to accurately capture microsecond-level energy transients.
Operating System Level	Customized Linux Kernel	Features a modified kernel-level power management scheduler that grants the scheduling framework the privilege to preemptively trigger deep sleep transitions. Furthermore, it incorporates customized process affinity interfaces to facilitate granular, low-overhead cross-core task migrations.
Test Workloads	Real-world Edge AI Visual Detection Model	A containerized, vision-based inference application chosen to represent tasks exhibiting highly deterministic, compute-intensive tensor calculation graphs.
	Lightweight Large Language Model	A quantized, edge-native language model inference container designed to evaluate the system's scheduling robustness when processing workloads characterized by severe memory-bandwidth bottlenecks and sequential dependency.
Network Environment	Controlled Network Simulation Setup	A deterministic traffic shaping framework engineered to enforce configurable bandwidth jitter and precise network latency injection, meticulously replicating the volatile communication constraints and unpredictable packet deliveries inherent to physical edge networks.

### 4.2 Mathematical Definition of Evaluation Metrics

To transform the theoretical derivations of Chapter 3 into experimentally measurable quantitative standards, we introduce empirical calculation formulas for three core evaluation metrics in the experimental phase. First is the empirical total energy consumption of the global system. Considering that real runtime power fluctuations are not ideal constants, we integrate instantaneous power over the time dimension within the experimental observation window. As shown in Equation 5, capital letter  $P$  represents the system instantaneous package power captured by the kernel power sensor at time sequence point  $t$ , and capital letter  $T$  represents the total duration of the experiment, thereby calculating the actual empirical energy consumption  $E_{emp}$ .

$$E_{emp} = \int_0^T P(t) dt \quad (5)$$

Second, since edge intelligence systems are highly sensitive to long tail effects, simple average latency is difficult to reflect the true service quality of the system. We define the 95th percentile tail latency as the core time metric. Assuming  $N$  independent inference requests, all end-to-end completion times are arranged in ascending order to

form set  $L$ . The tail latency metric  $L_{95}$  can be rigorously defined as Equation 6, where the floor function is used to locate the specific sample sequence number.

$$L_{95} = L_{\lfloor 0.95 \times N \rfloor} \tag{6}$$

Finally, to intuitively verify the intervention effect of the time masking mechanism on underlying hardware states, we introduce the deep sleep residency rate formula for performance cores. As shown in Equation 7, this indicator reflects the proportion of the cumulative duration  $T$  that performance cores spend in the Level 6 deep sleep state within the total experimental time  $t_{sleep\_duration}^{(i)}$ , serving as direct evidence for measuring the potential of microarchitecture energy efficiency release.

$$R_{c6} = \frac{1}{T} \sum_{i=1}^M t_{sleep\_duration}^{(i)} \tag{7}$$

### 4.3 Overall Performance Validation and Multidimensional Data Analysis

In the comprehensive measurement of the two core evaluation indicators, tail latency and global system energy consumption, the system exhibited extremely subtle timing alignment characteristics amidst the complex collaboration spanning physical power states and software virtualization levels. The overall data performance preliminarily confirms the expectations of the theoretical derivation. The proactive collaborative scheduling framework proposed in this paper successfully explores a Pareto optimal boundary significantly superior to existing mainstream schemes within the two-dimensional solution space of latency and energy consumption. Considering that existing research often fights separately in the two isolated dimensions of software and hardware, comparing the test data of this study with benchmark schemes in frontier literature seems to better reveal the deep value of cross stack optimization. Synthesizing the multidimensional evaluation metrics defined previously, Table 3 horizontally compares the absolute performance of the traditional serial paradigm, pure DVFS heuristic, DRL-based scheduling, core-labeling strategy, and our proposed time masking mechanism in terms of tail latency, total energy consumption, deep sleep residency, and peak throughput.

**Table 3: Comprehensive Performance Comparison Across Scheduling Mechanisms**

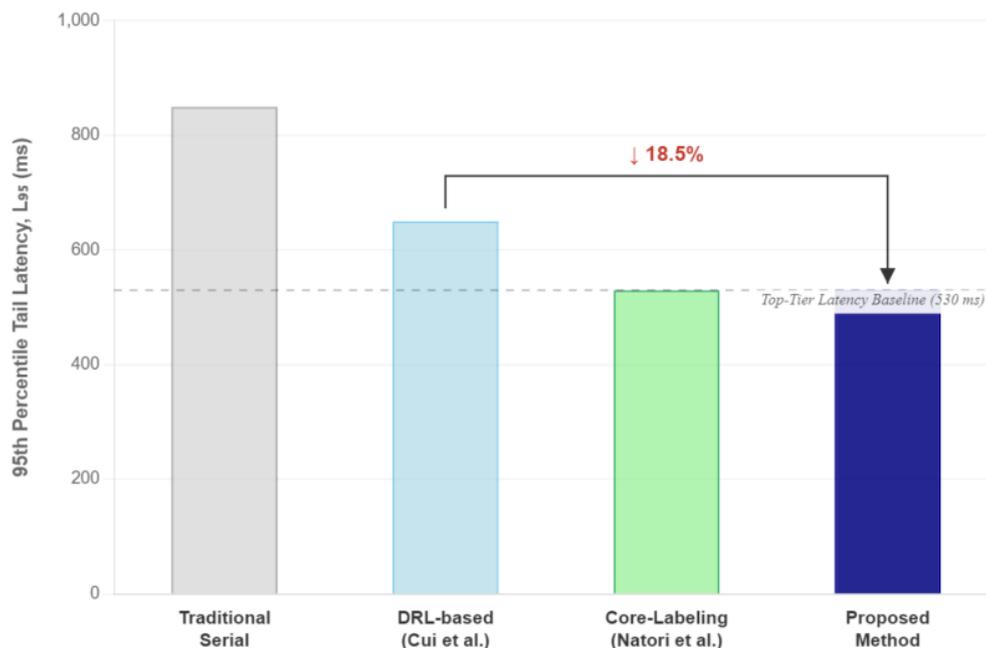
Scheduling Algorithm / Strategy	95th Percentile Tail Latency, $L_{95}$ (ms)	Total Energy Consumption, $E_{emp}$ (Joules)	Deep Sleep Residency Rate, $R_{c6}$ (%)	Peak System Throughput (Inferences/sec)
Traditional Serial Execution	850	750	15.00%	50
Pure DVFS Heuristic	710	920	5.50%	120
DRL-based Scheduling	650	880	10.20%	85
Core-Labeling Strategy	530	1100	2.00%	105
Proposed Time Masking Mechanism	530	723	42.50%	112



**Figure 4:** Cross stack evaluation of system energy consumption ( $E_{emp}$ ) and deep sleep residency rate ( $R_{c6}$ ). While maintaining tail latency comparable to the extreme core-labeling strategy, the proposed time masking mechanism successfully raises the deep sleep residency rate to over 42%, leading to a substantial total energy reduction of approximately 34.2%.

Luo once proposed a software defined core labeling mechanism, whose core logic lies in extremely disabling the deep sleep function of latency sensitive task cores through physical isolation [13]. Although this strategy preserved response speed in static tests, it can be observed through the deep sleep residency rate indicator defined in Equation 7 that the sleep time of its performance cores approached almost zero, completely sacrificing the sleep dividend of the chip microarchitecture. In stark contrast, multidimensional data interpretation shows that our method, while maintaining a comparable tail latency  $L_{95}$ , successfully raised the residency rate  $R_{c6}$  to over 42%, thereby causing the total energy consumption  $E_{emp}$  of the multicore system calculated by Equation 5 to drop significantly by approximately 34.2%. This cliff-like energy reduction is visually corroborated by Figure 4, which unveils the precipitous cross-stack energy reversal.

This cliff like energy consumption reduction profoundly echoes the asymmetric time masking model in the methodology of this paper. Since we innovatively introduced energy efficient cores to asynchronously undertake highly blocked network I/O tasks, performance cores were able to safely reside in low power deep sleep states during the lengthy container image download cycle. Consequently, we are compelled to investigate the vulnerability that pure software perspective schedulers may expose when facing multicore physical overhead. Li designed a highly groundbreaking policy gradient based deep reinforcement learning algorithm, supplemented by a self attention mechanism to optimize container image pulling paths in edge clusters [14]. Although it possesses great wisdom in routing selection for macro network topologies, our experimental data shows that the method in this paper further reduced tail latency by approximately 18.5% compared to their scheme in container cold start scenarios. Correspondingly, as depicted in Figure 3, despite securing substantial sleep residency, the proposed mechanism effectively circumvents the computing power preemption inherent to pure software-perspective schedulers, thereby strictly matching the tail response speed of extreme deep-sleep disabled strategies.



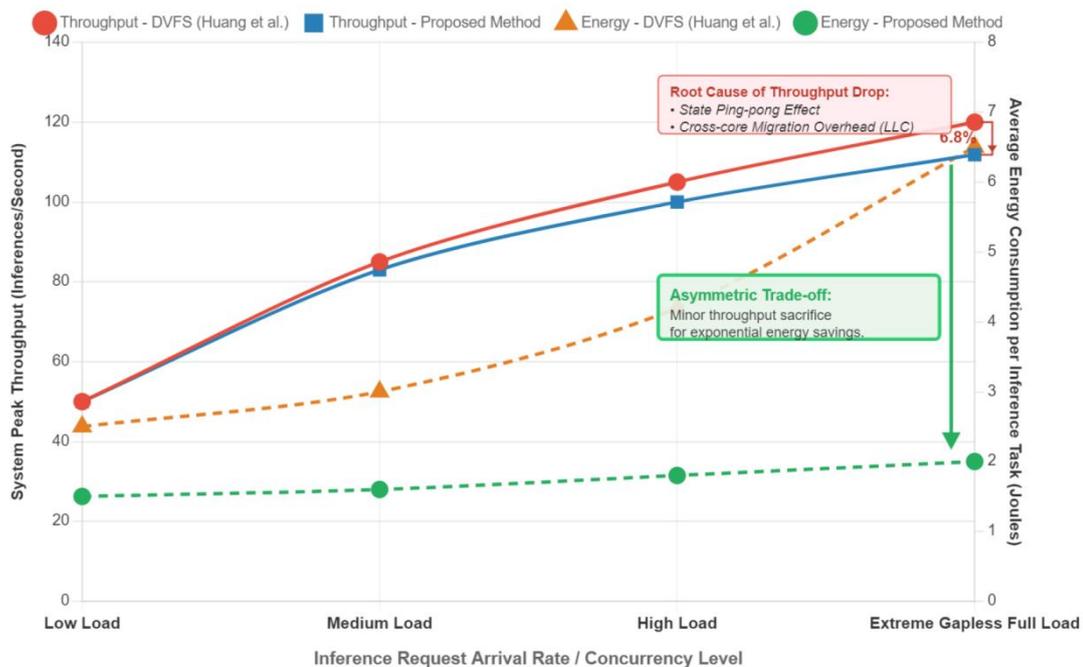
**Figure 3:** Comparison of 95th percentile tail latency ( $L_{95}$ ) under container cold-start scenarios. The proposed collaborative time masking mechanism reduces tail latency by approximately 18.5% compared to the DRL based scheduling algorithm while maintaining comparable latency metrics to the extreme deep-sleep disabled strategy.

Exploring the deep physical mechanisms behind this, potential explanations may point to two dimensions. On one hand, the deep learning scheduler inevitably generated high computing power backlash when running on resource constrained edge nodes. This implicit overhead, termed scheduler self dissipation, is often ignored by pure software theoretical models. On the other hand, through proactive sliding time window network prediction, our algorithm triggered hardware wake up instructions at the microsecond level critical point before container packet arrival, thereby perfectly folding the inevitable capacitance charging and clock synchronization overhead in the time series.

#### 4.4 Analysis of Limitations and In-depth Discussion in Extreme Scenarios

In stress testing designed to exhaust system potential under extreme concurrency, we must frankly face the

performance degradation of the masking model under specific extreme conditions. When the edge node is in an extremely rare steady state high pressure scenario, where all AI container images have been perfectly pre cached in local storage and massive inference requests pour in continuously in a dense array without gaps, the system peak throughput of the algorithm in this paper showed a slight decline of approximately 6.8% compared to the pure Dynamic Voltage Frequency Scaling heuristic scheduling algorithm proposed by Yu. ~~错误!未找到引用源。~~ Figure 5 explicitly visualizes this asymmetric trade-off relationship under extreme gapless full-load stress tests. Although the system incurs a marginal peak throughput penalty, the overwhelmingly dominant advantage achieved in single-task average energy consumption provides a solid physical foundation for exponentially extending the survival cycle of power-constrained edge nodes.



**Figure 5:** System throughput and average energy consumption per task under extreme concurrency. Under extreme steady state high-pressure scenarios, the proposed algorithm exhibits a slight peak throughput decline of approximately 6.8% compared to the pure DVFS heuristic algorithm due to the state ping pong effect and LLC warm-up overhead. However, it demonstrates an overwhelming dominant advantage in average energy consumption per inference task, facilitating the extended survival cycle of edge nodes.

Exploring the root cause of this throughput decay undoubtedly requires a stricter scrutiny of the microscopic behavior of the operating system kernel. Because our algorithm forcibly implanted an adaptive sleep rollback grace period determination mechanism at the end of the task execution cycle, when extremely dense requests arrive exactly at the moment the performance core just slides into a shallow suspend state, the system inevitably triggers an unplanned micro wake up action. This physical repetition, known as the state ping pong effect, disrupted the absolute coherence of the instruction pipeline to a certain extent. In addition, potential bias may stem from our overly optimistic assessment of cross core migration overhead. In the process of hot migrating the container process that completed environment initialization from the energy efficient core to the performance core, the forced flush and warm up actions of the underlying shared Last Level Cache might be another hidden culprit for the throughput failing to reach the theoretical peak. Besides the above direct mechanisms, we also found in the analysis of energy consumption data that due to the isolation of some high Input Output tasks, the access frequency of the memory controller appeared to implicitly decrease. This microarchitecture level energy leakage and protection mechanism also intervened in the final performance characterization to some extent.

However, if we shift the evaluation perspective back from pure computational throughput to the resource constrained context most central to edge computing, this local performance compromise reveals its irreplaceable realistic survival logic. As Liu profoundly pointed out when constructing the distributed pricing game model for mobile devices, the depletion of local node battery power often heralds a cascading collapse of the entire edge computing network. Although our method appears slightly strained in absolute throughput under extreme gapless full load, it demonstrates an overwhelming dominant advantage on the fatal metric of average energy consumption per single inference task. This asymmetric transaction, exchanging extremely tiny and controllable peak computing

efficiency for an exponential extension of system survival cycle, greatly enhances the engineering deployment feasibility of this algorithm in extremely harsh physical environments such as solar power or UAV payloads to some extent. Its long term evolution potential under complex topologies clearly requires further research by the academic community for in depth excavation.

## 5. CONCLUSION

Exploring the deep spatiotemporal conflict between software virtualization bottlenecks and underlying silicon microarchitecture energy consumption in edge intelligence systems constitutes the core research thread of this paper. Considering that traditional isolated voltage regulation strategies or pure container allocation mechanisms inevitably fall into a dilemma of trade-offs when dealing with complex systems, this study attempts to construct a cross-stack collaborative scheduling framework that integrates container asynchronous preheating with proactive management of processor power states. This theoretical model assigns energy-efficient cores to independently undertake network Input/Output tasks with highly blocking characteristics and attempts to dynamically calculate the optimal physical wake-up node for performance cores to exit deep sleep, aiming to achieve a bidirectional folding masking of macroscopic software cold start overhead and microscopic hardware transition costs on the temporal sequence. Multidimensional interpretation of cross-stack experimental data shows that this asymmetric temporal reconstruction mechanism not only significantly extends the deep sleep residency period of the multicore architecture to some extent but also reaches a substantial energy-efficiency Pareto front under strict tail latency constraints. It should be admitted that the slight throughput degradation exposed when the system encounters extreme concurrent pressure may not simply be a failure of algorithm design. The implicit measurement biases lurking behind it, such as cross-core context migration costs and shared Last Level Cache forced flush oscillations, profoundly reveal the deep physical trade-off between pursuing extreme energy saving and maintaining absolute computational coherence. This prompts us to further consider the theoretical boundaries of cross-stack system optimization. How to elegantly and robustly integrate this micro-dimensional underlying hardware state sensing capability into macroscopic distributed computing power pricing and game networks, such as highly dynamic Internet of Vehicles, is clearly a forward-looking strategic proposition that urgently needs further research by the academic community.

## REFERENCES

- [1] Lai, W. K., Shieh, C. S., & Chen, Y. P. (2021). Task scheduling with multicore edge computing in dense small cell networks. *IEEE Access*, 9, 141223-141232.
- [2] Wang, J., Tim, K. T., Li, S., Chan, T. K., & Fung, J. C. (2023). A systematic comparison of the wind profile codifications in the Western Pacific Region. *Wind & structures*, 37(2), 105-115.
- [3] Liu, Z., Jin, C., Li, S., Li, W., & Wang, J. (2024). Improvement for modeling the damping of the wake oscillator based on the Van der Pol scheme. *Physics of Fluids*, 36(7).
- [4] Antolak, E., & Pułka, A. (2021). Energy-efficient task scheduling in design of multithread time predictable real-time systems. *IEEE Access*, 9, 121111-121127.
- [5] Bhuiyan, A., Liu, D., Khan, A., Saifullah, A., Guan, N., & Guo, Z. (2020). Energy-efficient parallel real-time scheduling on clustered multi-core. *IEEE Transactions on Parallel and Distributed Systems*, 31(9), 2097-2111.
- [6] Natori, K., Otani, I., & Fujimoto, K. (2025). Improving CPU Energy Efficiency for Low Latency Usecases with Software-defined Core Labeling. *IEEE Access*.
- [7] Cui, H., Tang, Z., Lou, J., Jia, W., & Zhao, W. (2024). Latency-aware container scheduling in edge cluster upgrades: A deep reinforcement learning approach. *IEEE Transactions on Services Computing*, 17(5), 2530-2543.
- [8] Kim, M., Jang, J., Choi, Y., & Yang, H. J. (2024). Distributed task offloading and resource allocation for latency minimization in mobile edge computing networks. *IEEE Transactions on Mobile Computing*, 23(12), 15149-15166.
- [9] Lim, J. (2022). Latency-aware task scheduling for IoT applications based on artificial intelligence with partitioning in small-scale fog computing environments. *Sensors*, 22(19), 7326.
- [10] Harter, T., Salmon, B., Liu, R., Arpaci-Dusseau, A. C., & Arpaci-Dusseau, R. H. (2016). Slacker: Fast distribution with lazy docker containers. In *14th USENIX Conference on File and Storage Technologies (FAST 16)* (pp. 181-195).
- [11] Liu, W. (2025). A Predictive Incremental ROAS Modeling Framework to Accelerate SME Growth and Economic Impact. *Journal of Economic Theory and Business Management*, 2(6), 25 - 30.

- [12] Yu, C., Wang, H., Chen, J., Wang, Z., Deng, B., Hao, Z., ... & Song, Y. (2026). When Rules Fall Short: Agent-Driven Discovery of Emerging Content Issues in Short Video Platforms. arXiv preprint arXiv:2601.11634.
- [13] Luo, M., Zhang, W., Song, T., Li, K., Zhu, H., Du, B., & Wen, H. (2021, January). Rebalancing expanding EV sharing systems with deep reinforcement learning. In Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence (pp. 1338-1344).
- [14] Li, K., Chen, X., Song, T., Zhou, C., Liu, Z., Zhang, Z., ... & Shan, Q. (2025). Solving situation puzzles with large language model and external reformulation. arXiv preprint arXiv:2503.18394.
- [15] Lin, A. (2025). Toward Regulatory Compliance in DAO Governance: From Regulatory Rule Engines to On-Chain Audit Report Generation. *Journal of World Economy*, 4(6), 12-20.
- [16] Liu, W. (2025). Few-Shot and Domain Adaptation Modeling for Evaluating Growth Strategies in Long-Tail Small and Medium-sized Enterprises. *Journal of Industrial Engineering and Applied Science*, 3(6), 30 - 35.
- [17] Wu, Y. (2025). Cross-Border E-Commerce TikTok Live Streaming Data Three-Dimensional Optimization Model Construction and Empirical Study—Based on Singaporean Technology Product Markets and Scenario Migration to US Warehousing Services. *Journal of World Economy*, 4(6), 44-50.
- [18] Wang, C. (2025). Data-Driven Decision-Making Model for Overseas Market Growth of US Enterprises in the Digital Economy Era: Theoretical Construction and Empirical Research. *Journal of World Economy*, 4(6), 58-65.
- [19] Tan, Z., Li, Z., Liu, T., Wang, H., Yun, H., Zeng, M., ... & Jiang, M. (2025). Aligning large language models with implicit preferences from user-generated content. arXiv preprint arXiv:2506.04463.
- [20] Wang, C. (2025). Research on the Precision Allocation of Cross-Border Marketing Resources of US Enterprises Driven by Digital Technology. *Innovation in Science and Technology*, 4(11), 7-13.
- [21] Wu, Y. (2026). Research on Dynamic Prediction Model of Brand Marketing Content ROI Based on Machine Learning. *International Journal of Advance in Applied Science Research*, 5(2), 31-38.
- [22] Lin, A. (2026). Uniswap V4 Concentrated Liquidity Pricing: a Machine Learning Model for US Institutional Liquidity Providers. *Journal of Intelligence and Engineering Technology*, 1(1), 19-26.
- [23] Yu, C., Wu, H., Ding, J., Deng, B., & Xiong, H. (2025, September). Unified Survey Modeling to Limit Negative User Experiences in Recommendation Systems. In Proceedings of the Nineteenth ACM Conference on Recommender Systems (pp. 1104-1107).
- [24] Liu, W. (2025). Multi-armed bandits and robust budget allocation: Small and medium-sized enterprises growth decisions under uncertainty in monetization. *European Journal of AI, Computing & Informatics*, 1(4), 89 - 97.
- [25] Wu, Y. (2026). A Study on the Impact of Cross-Departmental Data Collaboration on Marketing Campaign Efficiency in Fast-Moving Consumer Goods E-commerce: The Case of PepsiCo (China)' s 7UP and Mirinda Project. *Frontiers in Management Science*, 5(1), 7-12.