# Dynamic Resource Aware Task Scheduling for Mobile Edge Cloud Computing

Indrajeet Roy College of Engineering Northeastern University Boston, MA, United States roy.i@northeastern.edu Ajith Kalpathi Ganesh
College of Engineering
Northeastern University
Boston, MA, United States
kalpathiganesh.a@northeastern.edu

Abdulrazaq Surakat College of Engineering Northeastern University Boston, MA, United States surakat.a@northeastern.edu

Abstract—This paper presents an enhanced implementation of the Mobile Cloud Computing (MCC) task scheduling algorithm originally proposed by Lin et al.[1] in Energy and Performance-Aware Task Scheduling in a Mobile Cloud Computing Environment. Although the original algorithm effectively optimized energy consumption under strict completion time constraints, its reliance on static power models, fixed network conditions, and predetermined task characteristics limited its real-world applicability. Our implementation addresses these limitations by integrating dynamic power profiles that adjust to varying device states, modeling realistic network conditions with throughput fluctuations, and adapting task attributes to simulate computational complexity and data intensity. In addition, we extend Lin's binary device-cloud architecture into a three-tier framework that incorporates edge computing resources, providing additional offload options with various energy-performance trade-offs. To further enhance decisionmaking, we employ a Q-learning approach for task migration, which adaptively explores the solution space to identify optimal offloading patterns. In general, our approach bridges the gap between theoretical MCC scheduling algorithms and practical requirements, enabling efficient energy-efficient task offloading across a range of device power states, network conditions, and computational loads without sacrificing algorithmic efficiency.

#### I. Introduction

Mobile Cloud Computing (MCC) augments the limited computational, storage, and battery resources of handheld devices by offloading tasks to remote servers. As applications such as real-time object recognition and on-device language models become prevalent, task placement decisions must optimize across three competing objectives: energy consumption, execution latency, and deadline satisfaction. This optimization is complex, requiring consideration of task dependencies, computational demands, network conditions, and device power states while ensuring application deadlines are met.

Lin et al.[1] introduced a two-phase scheduler for energy and performance-aware task offloading in a binary devicecloud setting. Their approach first generates a minimal-delay schedule using a HEFT-inspired algorithm, then iteratively migrates tasks between device and cloud to minimize energy consumption while respecting deadline constraints. Although effective in controlled settings, their approach relies on simplifying assumptions such as static power models, constant bandwidth, and binary execution options that limit its practical applicability.

Edge computing provides additional options for the task offloading set. Edge devices, positioned closer to end users than traditional cloud datacenters, provide an intermediate execution tier that reduces latency compared to cloud only solutions while typically consuming less energy than local execution. This three-tier architecture introduces new trade-offs and significantly expands the scheduling solution space beyond Lin et al.[1]'s binary model. We present three key extensions to Lin's MCC task scheduling algorithm:

- Dynamic resource modeling: Static power models are extended with dynamic profiles that account for device operating states, core-frequency scaling, and workload-dependent power consumption. Network conditions are modeled using time-varying bandwidth and signal-strength-dependent energy costs.
- Three-Tier Model Integration: The binary device—cloud scheduler is extended to incorporate edge servers, permitting task placement on the mobile device, an edge node, or the cloud. The original two phase procedure delay optimal initial scheduling followed by energy focused migration is adapted to the expanded architecture, enabling offloading strategies that utilize the low-latency and moderate-energy characteristics of edge resources.
- Alternative Optimization Strategies: To efficiently explore the expanded solution space, two refinement approaches are implemented and compared for the energy-optimization phase:
  - A heuristic migration strategy extended to three tiers, which deterministically selects task migrations based on estimated energy/time improvements.
  - A reinforcement learning (Q-learning) migration strategy, which adaptively explores migration decisions across diverse task graphs and operating conditions to learn policies that may outperform fixed heuristics in complex, dynamic scenarios.

#### II. RELATED WORK

Prior work in mobile cloud computing (MCC) covers dynamic resource modeling, edge computing, and reinforcement-learning-based task offloading. Lin et al. [1] introduced a two-phase scheduling algorithm that employs static power models and a binary device—cloud architecture. Subsequent studies refined task scheduling under MCC—Guo et al. [5] explored workflow-based offloading, and Wu et al. [6] added context awareness—but none extended the model to a three-tier device—edge—cloud environment.

Power modeling has evolved from simplistic constant-power approaches to sophisticated state-dependent models. Carroll and Heiser [7] conducted empirical measurements of smartphone power consumption, while Huang et al. [8] characterized energy impacts of different network interfaces under variable conditions. Our dynamic resource models incorporate these insights with battery-level sensitivity and workload-dependent energy estimation. Edge computing has transformed mobile offloading from binary decisions to multi-tier architectures. Satyanarayanan et al. [9] introduced cloudlets as trusted, resource-rich computers available for nearby mobile devices, while Shi et al. [10] defined edge computing as computation occurring at the network edge. Our three-tier extension builds on these concepts while addressing algorithmic challenges of task scheduling across heterogeneous tiers.

For decision optimization, Tang and Wong [2] applied deep reinforcement learning to task offloading in edge computing systems, demonstrating adaptability to changing conditions. Huang et al. [11] developed distributed reinforcement learning for multi-user computation offloading. Our Q-learning approach specifically targets energy-time trade-offs in multi-tier environments, integrating reinforcement learning as an enhancement to a proven algorithmic framework.

Task characterization has evolved from simplified executiontime models to nuanced representations. Wang et al. [12] and Flores et al. [13] established taxonomies based on computational and data intensity, informing offloading policies. Our approach incorporates these insights while adding dynamic attributes that drive tier-specific offloading decisions for compute-intensive, data-intensive, and balanced tasks.

# III. FOUNDATIONAL ENERGY AND PERFORMANCE-AWARE TASK SCHEDULING ALGORITHM

Lin et al.[1]'s two-phase task-scheduling algorithm for Mobile Cloud Computing (MCC) serves as the baseline. The following section summarizes the original system model, formalizes the problem, and reviews the scheduling methodology, while identifying limitations that motivate the extensions introduced in this paper.

#### A. System Model and Problem Formulation

Lin et al.[1] modeled an application as a directed acyclic graph (DAG) G = (V, E), where each node  $v_i \in V$  represents a task and each directed edge  $(v_i, v_j) \in E$  imposes a precedence constraint requiring task  $v_i$  to complete before task  $v_j$  can start. The execution environment consists of a mobile device with K heterogeneous local cores and a remote cloud with virtually unlimited computational resources, enabling task offloading and increased parallelism.

The energy consumption model proposed by Lin et al.[1] splits the device's energy into three parts:

## Local-execution energy

$$E_i^{\text{device}} = P_k T_{i,k}^l,$$

where  $P_k$  is the constant power draw of core k and  $T_{i,k}^l$  is the time required to execute task  $v_i$  on that core.

# · Offloading energy

$$E_i^{\mathrm{RF}} = P_{\mathrm{d2c}} T_i^s, \qquad T_i^s = \frac{d_{i,\mathrm{d2c}}^{\mathrm{send}}}{r_{\mathrm{eff}}(\mathrm{d2c})},$$

where  $P_{\rm d2c}$  is the RF-transmitter power,  $d_{i,\rm d2c}^{\rm send}$  is the input-data size (bytes), and  $r_{\rm eff}(\rm d2c)$  is the effective uplink rate. (Energy for receiving the result is included in background overhead and is not modeled separately.)

 Cloud-computation energy Computation energy on the cloud is treated as zero from the mobile device's perspective.

The total energy consumed by the mobile device,  $E_{\rm total}$ , aggregates individual task energies based on their respective execution assignments (local core or cloud offloading):

$$E_{ exttt{total}} = \sum_{v_i \in V} egin{cases} E_i^{ exttt{device}}, & ext{if tier}(v_i) = ext{DEVICE} \ E_i^{ ext{RF}}, & ext{if tier}(v_i) = ext{CLOUD} \end{cases}$$

The scheduling problem formulated by Lin et al.[1] chooses, for each task  $v_i \in V$ , both an execution location and schedule to minimize the total mobile-device energy  $E_{\rm total}$ , subject to:

- 1) **Precedence Constraints:**  $\forall (v_i, v_j) \in E$ , task  $v_j$  cannot start until  $v_i$  has completed and its results are available.
- 2) Deadline Constraint: The makespan

$$T_{\text{total}} = \max_{v_i \in \text{exit\_tasks}} \{ \max(FT_i^l, FT_i^{wr}) \}$$

must satisfy  $T_{\text{total}} \leq T_{\text{max}}$ , where  $FT_i^l$  is the finish time if  $v_i$  runs locally, and  $FT_i^{wr}$  is the finish time when  $v_i$  is offloaded (including round-trip transmission).

Solving the problem entails three scheduling decisions:

 (i) Task placement—determine for each task whether it executes on the mobile device or is offloaded to the cloud.

- (ii) *Core assignment*—map locally executed tasks to the heterogeneous cores available on the device.
- (iii) *Start-time scheduling*—select the exact start time on the assigned resource, i.e., a local core or the device-to-cloud communication channel.
- 1) Phase 1: Initial Scheduling (Minimal Makespan):

The goal of this phase is to generate a valid task schedule with the smallest possible makespan  $T_{\rm total}$ , satisfying the deadline  $T_{\rm max}$ . In practice,  $T_{\rm max}$  is set equal to this minimal makespan. The procedure is based on the HEFT algorithm and consists of three steps:

1) Primary Assignment (Local vs. Cloud). Assign each task  $v_i$  to the cloud if its estimated remote execution time

$$T_{i,\text{cloud}}^{\text{est}} = T_i^s + T_i^c + T_i^r$$

(upload + compute + download) is less than its minimal local execution time

$$T_{i,\text{device}}^{\text{est}} = \min_{k} T_{i,k}^{l}.$$

2) **Task Prioritization (Rank).** Compute an upward rank for each task  $v_i$ :

$$rank(v_i) = \overline{w}_i + \max_{v_j \in succ(v_i)} rank(v_j),$$

where

$$\overline{w}_i = \begin{cases} \frac{1}{K} \sum_{k=1}^{K} T_{i,k}^l, & \text{if } v_i \text{ is local,} \\ T_{i,\text{cloud}}^{\text{est}}, & \text{if } v_i \text{ is remote.} \end{cases}$$

Higher rank values indicate tasks on the critical path.

3) **Execution-Unit Selection.** Tasks are scheduled sequentially in descending order of their computed priorities. For each available execution resource (local cores *k* and the cloud), the Earliest Start Time (EST) and Earliest Finish Time (EFT) are computed, ensuring adherence to task precedence constraints:

$$EST(v_i, k) = \max \left\{ RT_i^l, \, FAT_{\text{core}}(k), \, \max_{v_j \in \text{pred}(v_i)} FT_{\text{final}}(v_j) \right\}$$

$$EST(v_i, C) = \max \left\{ RT_i^{ws}, \ FAT_{\text{ws}}, \ \max_{v_j \in \texttt{pred}(v_i)} FT_{\texttt{final}}(v_j) \right\}$$

 $RT_i^l$  and  $RT_i^{ws}$  represents the earliest times task  $v_i$  can start locally or begin uploading data to the cloud.  $FAT_{\rm core}(k)$  and  $FAT_{\rm ws}$  represent the finish available times for the respective resources.  $FT_{\rm final}(v_j)$  represents the effective finish time of predecessor tasks, accounting for the location of their results. The resource providing the lowest EFT is selected for scheduling task  $v_i$ .

2) Phase 2: Energy Reduction (Task Migration): Phase 2 minimizes the total energy consumption  $E_{\rm total}$  of the Phase 1 schedule while enforcing the makespan constraint  $T_{\rm max}$ . Starting from a minimal-makespan schedule (i.e.,  $T_{\rm total} \leq T_{\rm max}$ ), tasks are iteratively migrated to alternative execution units to reduce  $E_{\rm total}$  without exceeding the deadline.

This process consists of four main steps:

- Migration Evaluation: The algorithm evaluates potential migrations of tasks initially assigned to local cores, considering two migration types: (i) from one local core to another (potentially more energy-efficient) local core, and (ii) from a local core to the cloud.
- 2) Selection Criteria: During each iteration, the algorithm selects the most beneficial migration based on energy savings and makespan implications. The decision follows a two-tier criterion:
  - (i) Migrations that yield the highest energy reduction
     (ΔE = E<sub>before</sub> E<sub>after</sub>) without increasing the current
     makespan (T<sub>total</sub>) are given priority.
- (ii) If no migrations satisfy the first criterion, the migration that maximizes the energy-to-time improvement ratio  $(\Delta E/(T_{\rm after}-T_{\rm before}))$  is selected, provided the updated schedule remains within the allowed deadline  $(T_{\rm total} \leq T_{\rm max})$ .
- 3) Linear-Time Rescheduling: After identifying the optimal migration candidate, a linear-time kernel routine (i.e., O(|V|), where |V| is the total number of tasks) updates the existing schedule by adjusting the sequences of tasks on the affected execution units and recalculating their start and finish times.
- 4) **Iteration and Convergence:** Steps 1 through 3 are repeated iteratively until no further beneficial migrations can be identified that adhere to the original makespan constraint  $T_{max}$ .
- B. Limitations of the Original Approach

Although Lin et al.'s[1] algorithm establishes a baseline for energy and performance aware task scheduling in MCC, its reliance on static assumptions limits effectiveness under dynamic mobile cloud computing operating conditions:

- Static Power Consumption Model: The original method assumes constant power parameters (P<sub>k</sub>, P<sub>d2c</sub>) for device cores and wireless transmission. In practice, device power consumption varies with battery state, CPU-frequency scaling, and workload characteristics,making the static model inaccurate and inflexible.
- Constant Network Condition Assumption: Lin et al.'s algorithm uses fixed network parameters—constant communication rates  $r_{\rm eff}({\rm d}2{\rm c})$ ,  $r_{\rm eff}({\rm c}2{\rm d})$ —and static transmission energy costs, which do not simulate real-world wire-

less conditions where throughput, latency, and power consumption vary with signal strength, congestion, mobility, and interference.

- Binary Architectural Constraint: The original framework considers only a two-tier (device-cloud) architecture and omits edge computing, which can offer lower latency than cloud and lower energy costs than local execution, providing effective intermediate offloading options.
- Simplified Task Characterization: The original model considers only execution times (T<sub>i,k</sub>, T<sub>i</sub><sup>c</sup>) and ignores key factors such as computational complexity (CPU-bound vs. I/O-bound), data intensity (transfer volume), and pertask communication overhead that critically influence offloading decisions.

#### IV. DYNAMIC RESOURCE MODELING

The proposed framework addresses the limitations of static resource models by integrating dynamic models that accurately simulate variability in real-world mobile computing scenarios.

#### A. Mobile Device Power Consumption Modeling

In comparison to Lin et al.'s[1] static power model, the framework utilizes dynamic, state-dependent, load-aware power models for mobile devices.

**Mobile Device Core Power Model:** The power consumption of a specific mobile device core k is modeled as

$$P_k = (b_k + c_k \cdot \text{load}) \, bf,$$

where  $b_k$  is the baseline idle power,  $c_k$  is the incremental active power coefficient,  $load \in [0, 1]$  is the normalized utilization, and bf is the battery-state scaling factor.

**Mobile Device Battery-Level Sensitivity:** Battery charge levels significantly impact mobile device power efficiency, especially at lower charge states. This sensitivity is modeled using a multiplicative battery scaling factor bf, which adjusts both idle and active power consumptions for all cores and RF components as follows:

$$bf = \begin{cases} 1.0, & \text{if } B > 30\\ 1.0 + 0.01 \times (30 - B), & \text{if } B \le 30 \end{cases}$$

where B represents the current battery level. This simulated up to a 30% increase in power consumption as the battery level drops below 30%, modeling real-world power inefficiencies at low battery states.

The proposed implementation models a big.LITTLE architecture with three core classes—high-performance, mid-range, and efficiency—each exhibiting distinct power—performance trade-offs. These differentiated profiles enable energy-aware scheduling by assigning tasks to cores based on their computational demands and the battery-scaling factor bf. Table I summarizes the power model parameters for each core class.

TABLE I: Core power models differentiated by performance class.

Core Class	Base idle (W)	<b>Total Power Expression</b>
High-performance	0.10	$(0.10 + 0.20 + 1.80 \log b)  bf$
Mid-range	0.05	$(0.05 + 0.10 + 1.40  \mathrm{load})  bf$
Efficiency	0.03	$(0.03 + 0.05 + 0.95  \mathrm{load})  bf$

#### B. Mobile Device RF Transmission Power Modeling

The framework defines RF power consumption models for device-to-edge and device-to-cloud transmissions, simulating realistic variations due to available bandwidth, and signal strength. r represents the data rate (Mbps), s the signal strength (dBm), bf the battery-scaling factor, and  $r_{\rm eff}$  the RF circuit efficiency. The models are given by

$$P_{\text{d2e}} = \frac{bf}{r_{\text{eff}}} \left[ 0.1 + 0.4 \frac{r}{10} (1 + 0.02 (70 - s)) \right]$$

$$P_{\rm d2c} = \frac{bf}{r_{\rm eff}} \left[ 0.15 + 0.6 \, \frac{r}{5} \left( 1 + 0.03 \left( 70 - s \right) \right) \right]$$

These models extend beyond the constant transmission power assumption in Lin et al.[1], enabling a more detailed simulation of mobile RF energy consumption. The coefficient values are derived from empirical measurements and vendor specifications to simulate realistic operating conditions:

- Baseline Offsets (0.1 for device-to-edge, 0.15 for device-to-cloud): Simulates the RF transceiver's idle power and control-plane overhead during connected standby (beacon signaling, channel maintenance, handshakes) even when little or no user data are sent. Essentially, the keep-alive cost of staying attached to the network.
- Data-Rate Scaling Factors (0.4 for device-to-edge, 0.6 for device-to-cloud): Simulates the additional RF power demanded by higher throughput. The coefficients convert requested data rate into the incremental power budget, with the larger cloud value reflecting the higher effective-isotropic-radiated-power (EIRP) and modulation complexity of wide-area links.
- Signal-Strength Adjustment Factors (0.02 for device-to-edge, 0.03 for device-to-cloud): Model automatic transmit-power-control (TPC) responses to path loss and fading. These factors represent the extra power required to maintain the target signal-to-noise ratio as received-signal strength degrades, with the larger cloud coefficient accounting for longer propagation distances and harsher transmission environments.

The RF power models explicitly differentiate between short-range edge offloading and long-range cloud transmissions, enabling the scheduler to generate more accurate and realistic energy-cost estimates for task migration decisions.

#### C. Edge Device and Cloud Server Power Modeling

Load-dependent power consumption models are utilized for edge devices and cloud servers to simulate hardware heterogeneity and dynamic utilization characteristics.

1) Edge Server Power Model: Power consumption for each edge server node  $m \in \{1, ..., M\}$  and each core  $c \in \{1, \dots, C_m\}$  is defined by assigning a core-specific efficiency factor, efficiency m.c, modeling hardware diversity:

efficiency<sub>m,c</sub> = 
$$\max (0.5, [1.0 - 0.05(m - 1) - 0.02(c - 1)] \times \delta$$

where  $\delta \sim \mathcal{U}(0.9, 1.1)$  introduces a small random variation, and the linear terms represent gradual efficiency degradation across edge devices and cores. A lower bound of 0.5 ensures realistic operational limits.

From this, a power scaling factor is defined as:

$$\kappa_{m,c} = \frac{1.0}{\text{efficiency}_{m,c}},$$

where higher values correspond to lower efficiency (higher power consumption). The power consumption of each edge core is computed as:

$$P_{\text{edge}}(m, c, \ell) = P_{\text{idle}}(m, c) + P_{\text{dynamic}}(m, c, \ell),$$

where:

$$\begin{split} P_{\text{idle}}(m,c) &= 5.0 \times \kappa_{m,c}, \\ P_{\text{dynamic}}(m,c,\ell) &= (3.0 + 12.0\,\ell) \times \kappa_{m,c}, \end{split}$$

with  $\ell \in [0,1]$  representing normalized core utilization.

The numerical coefficients used in the models simulate standard server benchmarks and typical industry data:

- Edge idle power baseline (5.0 W): Simulates low-power edge hardware commonly deployed in close-proximity infrastructures.
- Edge dynamic power scaling  $(3.0 + 12.0 \ell \text{ W})$ : Simulates realistic edge-server CPU power consumption ranges from idle (minimal load) to maximum load scenarios, consistent with low-to-medium performance edge CPUs.
- 2) Cloud Server Power Model: Cloud servers are modeled as high-capacity infrastructure, with increased idle power consumption and load-dependent scaling. The power model is:

$$P_{\texttt{cloud}}(\ell) = P_{\texttt{idle}}(\texttt{cloud}) + P_{\texttt{dynamic}}(\texttt{cloud}, \ell),$$

where:

$$P_{\rm idle}({\rm cloud}) = 50.0,$$
 
$$P_{\rm dynamic}({\rm cloud},\ell) = 20.0 + 180.0\,\ell,$$

and  $\ell \in [0,1]$  representing normalized utilization.

The numerical coefficients used in the models simulate standard server benchmarks and typical industry data:

- Cloud idle power (50.0 W): Simulates a typical idle power draw for standard dual-socket cloud servers according to SPECpower benchmarks.
- Cloud dynamic power scaling (20.0 + 180.0 ℓ W): Simulates realistic incremental power consumption for cloud servers, scaling with workload intensity and aligning with standard 2-socket high-performance CPU server configurations.

These parameter values enable accurate, realistic evaluation  $\mathsf{efficiency}_{m,c} = \max{(0.5,\ [1.0 - 0.05(m-1) - 0.02(c-1)] \times \delta)}$  of task-offloading strategies in the three-tier MCC framework, addressing the limitations imposed by the static and simplified modeling assumptions.

# D. Network Variability Modeling

Lin et al.'s[1] model assumes fixed communication rates, which fails to capture real-world bandwidth fluctuations and congestion. The proposed framework defines separate upload and download base rates  $r_{\text{base,up}}(\ell)$ ,  $r_{\text{base,down}}(\ell)$  for each link type  $\ell$  in the three-tier architecture. Table II lists example values.

TABLE II: Base upload/download bandwidths for link types

Link type $\ell$	$r_{\mathrm{base,up}}(\ell)$ [Mbps]	$r_{\mathrm{base,down}}(\ell)$ [Mbps]
$\overline{\text{Device} \to \text{Edge } (\text{d2e}_m)}$	10	12
Device $\rightarrow$ Cloud (d2c)	5	6
$Edge \leftrightarrow Edge$	30	30
$Edge \rightarrow Cloud$	50	60

**Device**  $\rightarrow$  **Edge links**: Represent short-range connectivity (e.g., Wi-Fi, local 5G microcell) between a mobile device and a nearby edge server. Uplink speeds are modeled slightly lower than downlink to reflect common access-point asymmetry.

**Device** → **Cloud links**: Model wide-area cellular or public internet connections from the mobile device to a remote cloud data center. These links exhibit lower uplink throughput (e.g., due to shared spectrum and ISP traffic shaping) compared to local offloading paths.

Edge ↔ Edge links: Simulates high-capacity, symmetric backhaul connections (e.g., fiber or metro Ethernet) interconnecting edge nodes supporting low-latency, peer-to-peer task migration.

Edge → Cloud links: Simulates network links between edge infrastructure and central cloud regions to accommodate high symmetric bandwidth bulk data exchange and inter-datacenter synchronization.

By distinguishing these link types and their characteristic rates, the scheduler selects the appropriate parameters for each offloading path, improving the accuracy of transfer time and energy cost estimates under heterogeneous, time varying network conditions.

Scenario Scaling: A global bandwidth scaling factor  $f_{\rm BW}$  is applied to each base rate to simulate different network-quality scenarios:

$$r_{\rm eff}(\ell) = f_{\rm BW} \ r_{\rm base}(\ell).$$

Values  $f_{\rm BW} > 1$  simulate favorable conditions (e.g., 5G or wired backhaul), while  $f_{\rm BW} < 1$  represent degraded environments. This approach enables simulation across a range of network states without modifying link-specific parameters. In the three-tier architecture, distinct base bandwidths are defined for device–edge (d2e<sub>m</sub>), edge–cloud, and device–cloud (d2c) links. Scaling these rates by  $f_{\rm BW}$  simulates symmetric and asymmetric variations, allowing the scheduler to opportunistically exploit higher-bandwidth paths. Such dynamic network modeling is critical for the design and evaluation of effective scheduling strategies in MCC environments.

#### E. Task Characterization

Lin et al.'s[1] original model characterizes each task only by its execution times  $(T_{i,k}^l, T_i^c)$ , ignoring important factors that affect offloading efficiency. The proposed framework augments task representations with:

- Computation Category: Classifies tasks as computebound, data-bound, or balanced, enabling the scheduler to prioritize appropriate resources for each task type. This differentiation allows compute-intensive tasks to be directed to high-performance cores or cloud resources, while keeping data-intensive tasks closer to their data sources to minimize transfer overhead.
- Data Payload Size: Quantifies the input/output data volumes for each task, allowing precise calculation of energy and time costs for data transfers. This enables the scheduler to make intelligent tradeoffs between local execution (avoiding transfer costs) and remote execution (gaining computational benefits) based on actual data requirements.
- Communication Overhead: Models path-specific transfer times based on data sizes and bandwidth rates for each network link, enabling the scheduler to account for heterogeneous network conditions when making migration decisions. This enhances energy efficiency by avoiding costly transfers over constrained or power-hungry channels.

These heterogeneous attributes collectively enable more accurate, context-aware scheduling that dynamically adapts to both the computational nature of tasks and the current state of the three-tier network environment.

Tasks are classified into three categories—Compute-Intensive, Data-Intensive, and Balanced. Categories are assigned probabilistically according to P(Compute) = 0.3, P(Data) = 0.3, P(Balanced) = 0.4, to simulate mixed application workflows. Examples include encryption and image processing for Compute-Intensive tasks, and database queries and large file transfers for Data-Intensive tasks.

Task Attributes: Complexity and Data Intensity: Each task  $v_i$  carries two stochastic attributes:

- Computational complexity complexity  $(v_i) \in [\gamma_{\min}, \gamma_{\max}]$
- Data intensity intensity $(v_i) \in [\delta_{\min}, \delta_{\max}]$

Task class	$\mathtt{complexity}(v_i) \; \mathtt{range}$	$\mathtt{intensity}(v_i)$ range
Compute-Intensive	$[0.7\gamma_{ m max},\gamma_{ m max}]$	$[\delta_{\min},2\delta_{\min}]$
Data-Intensive	$[\gamma_{ m min}, 2\gamma_{ m min}]$	$[0.7\delta_{ m max},\delta_{ m max}]$
Balanced	$[\gamma_{\min},\gamma_{\max}]$	$[\delta_{\min},\delta_{\max}]$

## • Task specialization classification:

For compute-intensive tasks, complexity  $(v_i) \sim U[0.7\,\gamma_{\rm max},\,\gamma_{\rm max}];$  for data-intensive tasks, intensity  $(v_i) \sim U[0.7\,\delta_{\rm max},\,\delta_{\rm max}].$  This ensures compute-heavy tasks typically demand more CPU than balanced or data tasks of equal size, and vice versa.

#### • Offloading incentive structure:

The non-dominant attribute is drawn from its lower range: intensity $(v_i) \sim U[\delta_{\min}, 2 \delta_{\min}]$  for compute tasks, and complexity $(v_i) \sim U[\gamma_{\min}, 2 \gamma_{\min}]$  for data tasks.

- CPU-heavy tasks incur low transmission cost but high local compute cost, promoting offloading when bandwidth permits.
- Data-heavy tasks incur low compute cost but high transmission cost, favoring local or edge execution under limited uplink.

#### • Energy-delay scaling alignment:

The 30%/70% splits keep attribute multipliers within the linear regions of the mobile-core power model  $P_k$  and the RF-energy models  $(P_{\rm d2e}, P_{\rm d2e})$ , ensuring proportional scaling of simulated energy and delay without saturating either model.

By sampling tasks as either strongly CPU-bound or strongly data-bound, the scheduler receives clear energy-versus-delay trade-off signals, simplifying evaluation of the three-tier of-floading logic and preventing mixed workloads from obscuring validation of the dynamic power and network models.

# V. THREE-TIER HEURISTIC EXTENSION

This section extends Lin et al.'s[1] two-phase scheduling algorithm to a three-tier (device-edge-cloud) architecture. The heuristic framework employs the dynamic resource models from Section IV and minimizes mobile-device energy consumption  $E_{\rm total}$  under the makespan constraint  $T_{\rm max}$ . The algorithm retains the original two phases:

- (i) *Makespan minimization:* A HEFT-inspired kernel algorithm produces a minimal-delay schedule.
- (ii) Energy optimization: A linear-time (O(|V|)) kernel algorithm iteratively migrates tasks to reduce  $E_{\rm total}$  while ensuring  $T_{\rm total} \leq T_{\rm max}$ .

## A. Phase 1: Initial Scheduling (Minimal Makespan)

This phase extends the HEFT algorithm to a three-tier (Device–Edge–Cloud) setting to produce a schedule minimizing the total completion time  $T_{\rm total}$ , subject to task precedence constraints. The resulting makespan is used to set  $T_{\rm max}$  for Phase 2.

- 1) Step 1.1: Preliminary Tier Estimation: For each task  $v_i \in V$ , we three completion times assuming no contention are computed for each tier, and the tier with the minimum estimate is selected:
- a) Device estimate:

$$T_{i,\text{dev}}^{\text{est}} = \min_{0 \le k < K} T_{i,k}^l$$

b) Cloud estimate:

$$\begin{split} T_{i,\text{cloud}}^{\text{est}} &= T_{\text{transfer}}(\text{d2c},\, d_{i,\text{d2c}}^{\text{send}}) + T_{i}^{c} \\ &+ T_{\text{transfer}}(\text{c2d},\, d_{i,\text{c2d}}^{\text{recv}}) \end{split}$$

c) Edge estimate:

$$\begin{split} \widehat{T}_{i,\mathrm{edge}}(m,c) &= T_{\mathrm{transfer}}(\mathtt{d2e}_m,\,d_{i,\mathtt{d2e}_m}^{\mathrm{send}}) + T_{i,m,c}^e \\ &\quad + T_{\mathrm{transfer}}(\mathtt{e2d}_m,\,d_{i,\mathtt{e2d}_m}^{\mathrm{recv}})\,, \\ T_{i,\mathrm{edge}}^{\mathrm{est}} &= \min_{m,c} \widehat{T}_{i,\mathrm{edge}}(m,c) \end{split}$$

d) Tier selection:

$$Tier_i^{pre} = \arg\min_{\tau \in \{\text{dev,edge,cloud}\}} T_{i,\tau}^{\text{est}}$$

The preliminary tier  $\operatorname{Tier}_i^{\operatorname{pre}}$  is used to compute the average  $\operatorname{cost}\overline{w}_i$  in the prioritization step (Step 1.2). This assignment is non-binding: the final execution unit selection (Step 1.3) may reassign  $v_i$  based on actual resource contention and dependencies to minimize its finish time.

# Algorithm 1 Primary tier assignment for each task $v_i \in V$

```
Require: Task set V; execution times T_{i,k}^l, T_{i,m,c}^e, T_i^c; data sizes
           d_{i,\ell}^{\text{send}}, d_{i,\ell}^{\text{recv}}; link rates r_{\text{eff}}(\ell)
     1 for all v_i \in V do
                     T_{i,\text{device}}^{\text{est}} \leftarrow \min_{k=0}^{K-1} T_{i,k}^{l}
                     t_{\text{up}} \leftarrow T_{\text{transfer}}(\text{d2c}, d_{i, \text{d2c}}^{\text{send}})
                     t_{\text{down}} \leftarrow T_{\text{transfer}}(\text{c2d}, d_{i,\text{c2d}}^{\text{recv}})

T_{i,\text{cloud}}^{\text{est}} \leftarrow t_{\text{up}} + T_i^c + t_{\text{down}}
                      T_{i, \text{edge}}^{\text{ést}} \leftarrow \infty \\  \text{for } m \leftarrow 1 \text{ to } M \text{ do} 
    7
                               for c \leftarrow 1 to C_m do
    8
                                         t_{\text{up}} \leftarrow T_{\text{transfer}}(\text{d2e}_m, d_{i, \text{d2e}_m}^{\text{send}})
t_{\text{down}} \leftarrow T_{\text{transfer}}(\text{e2d}_m, d_{i, \text{e2d}_m}^{\text{recv}})
    9
   10
                                          \widehat{T} \leftarrow t_{\text{up}} + T_{i,m,c}^e + t_{\text{down}}
   11
                                         T_{i,\text{edge}}^{\text{est}} \leftarrow \min(T_{i,\text{edge}}^{\text{est}}, \widehat{T})
   12
                     PreliminaryTier(v_i) \leftarrow \arg\min_{\tau \in \{\text{dev}, \text{edge}, \text{cloud}\}} T_{i,\tau}^{\text{est}}
   13
```

- 2) Step 1.2: Task Prioritization: After the primary assignment (Step 1.1) provides an initial estimate  $PreliminaryTier(v_i)$  for each task, this step calculates a priority value  $priority(v_i)$  for every task  $v_i \in V$ . This priority determines the order in which tasks are considered during the Execution Unit Selection phase (Step 1.3), giving precedence to tasks on the critical path of the application DAG. The calculation uses the upward rank method, common in list scheduling algorithms like HEFT.
- a) Compute Task Cost  $(\overline{w}_i)$ : An average or representative computation/execution cost  $\overline{w}_i$  is determined for each task  $v_i$ . This cost estimate is based on the task's characteristics and its preliminary tier assignment PreliminaryTier $(v_i)$  from Step 1.1:
  - If  $PreliminaryTier(v_i) = DEVICE$ : The cost is the average execution time across all K local device cores.

$$\overline{w}_i = \frac{1}{K} \sum_{k=0}^{K-1} T_{i,k}^l$$

• If PreliminaryTier $(v_i)$  = CLOUD: The cost is the total estimated path time  $T_{i,\text{cloud}}^{\text{est}}$  for cloud execution calculated in Step 1.1.

$$\overline{w}_i = T_{i, \text{cloud}}^{\text{est}}$$

• If PreliminaryTier $(v_i) = \text{EDGE}$ : The cost is the minimum estimated path time  $T_{i, \text{edge}}^{\text{est}}$  across all edge cores calculated in Step 1.1.

$$\overline{w}_i = T_{i, \text{edge}}^{\text{est}}$$

This cost  $\overline{w}_i$  represents the expected duration contribution of task  $v_i$  itself along the critical path, considering its likely execution environment.

- b) Compute Upward Rank Priority (priority( $v_i$ )): The priority of each task  $v_i$  is computed recursively, starting from the exit tasks of the DAG and moving upwards towards the entry tasks. The priority indicates the length of the critical path from task  $v_i$  to the end of the application, based on the estimated costs  $\overline{w}_i$ .
  - For an **exit task**  $v_{\text{exit}}$  (a task with no successors,  $\text{succ}(v_{\text{exit}}) = \emptyset$ ):

$$\mathtt{priority}(v_{\mathsf{exit}}) = \overline{w}_{\mathsf{exit}}$$

• For any other **non-exit task**  $v_i$ :

$$\texttt{priority}(v_i) = \overline{w}_i + \max_{v_j \in \text{succ}(v_i)} \{\texttt{priority}(v_j)\}$$

This recursive calculation ensures that a task's priority incorporates its own cost  $\overline{w}_i$  and also the maximum cost priority $(v_j)$  of the path following it through its successors  $\operatorname{succ}(v_i)$ . Tasks with higher priority values are deemed more critical.

c) Create Prioritized Task List  $(L_{prio})$ : All tasks  $v_i \in V$  are stored in a list  $L_{prio}$ , sorted in descending order based on their calculated priority $(v_i)$  values. Tasks with equal priority may be ordered arbitrarily or using secondary criteria (e.g., task ID). This list  $L_{prio}$  dictates the order in which tasks will be selected for scheduling in the next step.

Algorithm 2 Upward-Rank Prioritization (Step 1.2)

```
Require: Task set V, DAG G = (V, E), preliminary tiers, execu-
      tion times, no-contention estimates
Ensure: Priority list L_{\text{prio}} (descending)
      Stage 1 — Tier-aware cost
   1 for all v_i \in V do
            if PreliminaryTier(v_i) = 	ext{DEVICE} then \overline{w}_i \leftarrow \frac{1}{K} \sum_{k=0}^{K-1} T_{i,k}^l else if PreliminaryTier(v_i) = 	ext{CLOUD} then
  3
  4
                  \overline{w}_i \leftarrow T_{i,\text{cloud}}^{\text{est}}
  5
  6
            else
                  \overline{w}_i \leftarrow T_{i,\text{edge}}^{\text{est}}
      Stage 2 — Upward-rank recursion
  8 R \leftarrow \text{reverseTopo}(G)
  9 for all v_i \in R do
            if \operatorname{succ}(v_i) = \emptyset then
  10
                  priority(v_i) \leftarrow \overline{w}_i
 11
 12
                  \mathtt{priority}(v_i) \leftarrow \overline{w}_i + \max_{v_j \in \mathtt{succ}(v_i)} \mathtt{priority}(v_j)
 13
      Stage 3 — Sorted priority list
  14 L_{\text{prio}} \leftarrow \text{sort}(V, \text{priority}, \text{desc})
 15 return L_{\text{prio}}
```

3) Step 1.3: Execution Unit Selection (Earliest Finish Time Minimization): This step performs the core resource allocation for initial scheduling by assigning each task  $v_i$  to an execution unit u (device core k, edge core (m,c), or cloud server C) and computing its start time. The assignment minimizes the task's effective earliest finish time (EFT) while enforcing DAG precedence constraints and respecting resource availability across all three tiers. Tasks are scheduled sequentially in the order defined by the priority list  $L_{\rm prio}$  from Step 1.2.

For each task  $v_i$  selected from  $L_{prio}$ , the scheduler evaluates its potential execution on every available unit u. For each potential assignment  $(v_i, u)$ , the scheduler calculates the earliest time  $v_i$  could possibly finish its entire execution process, with results becoming available at the device. This involves calculating two key values:

Data Ready Time: The earliest time that all data dependencies for v<sub>i</sub> are met, meaning all its immediate predecessors v<sub>p</sub> ∈ pred(v<sub>i</sub>) have completed and their results are available at the device (FT<sub>final</sub>(v<sub>p</sub>)).

$$\mathsf{DataReadyTime}(v_i) = \max\left(0, \max_{v_p \in \mathsf{pred}(v_i)} \{FT_{\mathtt{final}}(v_p)\}\right)$$

This corresponds to  $RT_i^l$  but is foundational for all tiers.

• Earliest Start Time  $(EST(v_i, u))$ : The earliest moment task  $v_i$  can actually start its primary operation

(local computation or remote upload) on the first resource required by unit u. This depends on both data readiness and the availability of that specific resource.

$$EST(v_i, u_{initial}) = \max(DataReadyTime(v_i), FAT(u_{initial}))$$

where  $u_{\text{initial}}$  is the resource needed first (e.g., core k, channel d2c, channel d2e<sub>m</sub>), and  $FAT(u_{\text{initial}})$  is its Finish Available Time. Subsequent phases (compute, download) for remote tiers have their own ESTs determined by the preceding phase's finish time and the availability (FAT) of their respective resources  $(C \text{ or c2d for cloud; edge core } (m,c) \text{ or e2d}_m \text{ for edge}).$ 

The scheduler computes the final effective EFT for  $v_i$  for every possible assignment u. This EFT represents the time the task's results are available back at the mobile device  $(FT_{\text{final}}(v_i))$ . The calculation requires simulating the sequence of operations (upload, compute, download for remote tiers) and considering the availability (FAT) of each involved resource sequentially.

- a) Detailed Earliest Finish Time (EFT) Calculation and Resource Update Logic: The scheduler computes the effective earliest finish time (EFT), representing when task  $v_i$ 's results are available at the mobile device, for each candidate execution unit u by simulating the necessary steps and respecting resource availability:
- 1) **Device core**  $k \in \{0, \dots, K-1\}$ : For local execution on a device core k, the task can only start after all its predecessors' results are available at the device (DataReadyTime $(v_i)$ ) and the core itself is free  $(FAT_{core}(k))$ . The Earliest Start Time is the maximum of these two times:

$$EST(v_i, k) = \max(DataReadyTime(v_i), FAT_{core}(k)).$$

The final finish time on the device, which is also the effective EFT for this option, is simply the start time plus the local computation duration:

$$FT_i^l(k) = EST(v_i, k) + T_{i,k}^l.$$

2) **Cloud** C: Cloud execution involves three sequential phases. First, the device uploads data. This upload can start at  $EST_{ws}$ , which is the later of the data readiness time (DataReadyTime( $v_i$ )) and the upload channel availability ( $FAT_{ws}$ ). The upload finishes at  $FT_i^{ws}$ .

$$EST_{ws} = \max(\text{DataReadyTime}(v_i), FAT_{ws}),$$
  
 $FT_i^{ws} = EST_{ws} + T_{\text{transfer}}(\text{d2c}, d_{i,\text{d2c}}^{\text{send}}).$ 

Second, cloud computation starts at  $EST_c = RT_i^c$ , which requires the upload to be complete  $(FT_i^{ws})$  and any predecessors also processed on the cloud  $(v_p \in \operatorname{pred}_C(v_i))$ 

to have finished their computation  $(FT_p^c)$ . Computation finishes at  $FT_i^c$ .

$$RT_i^c = \max\left(FT_i^{ws}, \max_{v_p \in \text{pred}_C(v_i)} \{FT_p^c\}, 0\right),$$
  

$$EST_c = RT_i^c,$$
  

$$FT_i^c = EST_c + T_i^c.$$

Finally, the results are downloaded to the mobile device. This download can start at  $EST_{wr}$ , the later of the cloud computation finish time  $(FT_i^c)$  and the download channel availability  $(FAT_{wr})$ . The effective EFT for the cloud option is the final download finish time  $FT_i^{wr}$ .

$$EST_{wr} = \max(FT_i^c, FAT_{wr}),$$
  

$$FT_i^{wr} = EST_{wr} + T_{\text{transfer}}(\text{c2d}, d_{i,\text{c2d}}^{\text{recv}}).$$

3) **Edge core** (m, c) **for all** m, c: Edge execution is similar to the cloud path logic but uses edge-specific resources. The device first uploads data to edge node m, starting at  $EST_{des} = \max(\text{DataReadyTime}(v_i), FAT_{d2e}(m))$  and finishing at  $FT_{i,m}^{des}$ .

$$EST_{des} = \max(\text{DataReadyTime}(v_i), FAT_{dee}(m)),$$
  
 $FT_{i,m}^{des} = EST_{des} + T_{\text{transfer}}(\text{d2e}_m, d_{i,\text{d2e}_m}^{\text{send}}).$ 

Edge computation on core c can then potentially start at  $RT_{i,(m,c)}^e$ , which depends on the upload finishing  $(FT_{i,m}^{des})$  and any predecessors assigned to the same edge node m  $(v_p \in \operatorname{pred}_E(v_i, m))$  completing their computation  $(FT_{p,m,c'}^e)$ . The actual computation starts at  $EST_e$ , the later of this readiness time and the specific edge core's availability  $(FAT_{\text{edge}}(m,c))$ . Computation finishes at  $FT_{i,m,c}^e$ .

$$RT_{i,(m,c)}^{e} = \max\left(FT_{i,m}^{des}, \max_{v_p \in \operatorname{pred}_E(v_i,m)} \{FT_{p,m,c'}^{e}\}, 0\right),$$

$$EST_e = \max\left(RT_{i,(m,c)}^{e}, FAT_{\operatorname{edge}}(m,c)\right),$$

$$FT_{i,m,c}^{e} = EST_e + T_{i,m,c}^{e}.$$

Finally, results are downloaded back to the mobile device starting at  $EST_{er} = \max(FT_{i,m,c}^e, FAT_{\text{e2d}}(m))$  and completing at  $FT_{i,m}^{er}$ . This final time,  $FT_{i,m}^{er}$ , is the effective EFT for assigning task  $v_i$  to edge core (m,c).

$$EST_{er} = \max(FT_{i,m,c}^{e}, FAT_{\texttt{e2d}}(m)),$$
  
$$FT_{i,m}^{er} = EST_{er} + T_{\text{transfer}}(\texttt{e2d}_{m}, d_{i,\texttt{e2d}_{m}}^{\texttt{recv}}).$$

- b) Assignment and State Update: After evaluating the effective EFT for task  $v_i$  on all possible execution units u within the set  $\mathcal{U}$  (which includes all device cores k, the cloud C, and all edge cores (m,c)), the scheduler makes the final assignment and updates the system state:
- 1) **Select Best Unit:** Identify the execution unit  $u^*$  that resulted in the minimum effective EFT calculated across all possibilities. This represents the assignment that allows  $v_i$  to finish (with results available at the mobile device) earliest, given the current schedule and resource contention.

$$u^* = \arg\min_{u \in \mathcal{U}} \{\text{EffectiveEFT}(v_i, u)\}.$$

,where EffectiveEFT $(v_i, u)$  is  $FT_i^l(k)$ ,  $FT_i^{wr}$ , or  $FT_{i,m}^{er}$  depending on u

- 2) **Set Task Assignment:** Permanently assign task  $v_i$  to the selected unit  $u^*$ . This involves updating the task's internal state:
  - Set the numerical assignment index assignment(v<sub>i</sub>) to identify u\*.
  - Set the categorical execution tier tier(v<sub>i</sub>) to DEVICE,
     EDGE, or CLOUD based on u\*.
- 3) **Record Scheduled Finish Times:** Store the calculated finish times for all relevant phases associated with the chosen execution path  $(u^*)$  with the task  $v_i$ . This makes these times available for calculating the ready times of successor tasks. For example:
  - If assigned to device core k, record the final  $FT_i^l(k)$ .
  - If assigned to cloud C, record  $FT_i^{ws}$ ,  $FT_i^c$ , and the effective finish time  $FT_i^{wr}$ .
  - If assigned to edge core (m, c), record  $FT_{i,m}^{des}$ ,  $FT_{i,m,c}^{e}$ , and the effective finish time  $FT_{i,m}^{er}$ .

The critical value stored is  $FT_{\text{final}}(v_i)$ , which equals the effective EFT of the chosen path.

- 4) **Update Resource Availability (FAT):** Update the Finish Available Time (FAT) for each resource that was utilized along the selected execution path. The FAT of a resource is set to the time when task  $v_i$  finished using it. This ensures that subsequent tasks consider the correct availability when being scheduled. For instance:
  - If  $v_i$  used device core k, update  $FAT_{core}(k) \leftarrow FT_i^l(k)$ .
  - If  $v_i$  used the cloud path, update  $FAT_{ws} \leftarrow FT_i^{ws}$  and  $FAT_{wr} \leftarrow FT_i^{wr}$ .
  - If  $v_i$  used edge core (m,c), update  $FAT_{d2e}(m) \leftarrow FT_{i,m}^{des}$ ,  $FAT_{\text{edge}}(m,c) \leftarrow FT_{i,m,c}^{e}$ , and  $FAT_{e2d}(m) \leftarrow FT_{i,m}^{er}$ .

This prevents resource over-subscription in future scheduling decisions.

This greedy assignment minimizes each task's finish time under current contention, producing a makespan-optimized initial schedule for the three-tier system.

### **Algorithm 3** Execution Unit Selection (Step 1.3)

```
Require: Prioritized list L_{\text{prio}}; DAG G = (V, E); current FT_{\text{final}}(\cdot); FAT(\cdot); sets \mathcal{K}, \mathcal{M}, \mathcal{C}_m; times T^l_{i,k}, T^c_i, T^e_{i,m,c}; transfer times T_{\text{transfer}}(\ell, d_{i,\ell}) Ensure: Updated assignment (v_i), tier (v_i), FTs, and FATs
   1 for all v_i \in L_{\mathrm{prio}} do
                  DRT \leftarrow \max(0, \max_{v_p \in \text{pred}(v_i)} FT_{\text{final}}(v_p))
EFT^* \leftarrow \infty, \ u^* \leftarrow \text{null}
    3
                                                                                                                                                    for all k \in \mathcal{K} do
                           EST_k \leftarrow \max(DRT, FAT_{core}(k))
    5
                         FT_l \leftarrow EST_k + T_{i,k}^l
if FT_l < EFT^* then
EFT^* \leftarrow FT_l, \ u^*
    6
    8
                                                                                                                                                        EST_{ws} \leftarrow \max(DRT, FAT_{ws})
                  FT_{ws} \leftarrow EST_{ws} + T_{transfer}(d2c, d_{i,d2c}^{send})
RT_c \leftarrow \max(FT_{ws}, \max_{p \in \text{pred}_C(v_i)} FT_p^c, 0)
 10
 11
 12
                  FT_c \leftarrow RT_c + T_i
                 \begin{array}{l} T_{L_c} \leftarrow H_{L_c} + I_i \\ FT_{wr} \leftarrow \max(FT_c, FAT_{wr}) + T_{\mathrm{transfer}}(\mathrm{c2d}, d_{i,\mathrm{c2d}}^{\mathrm{recv}}) \\ \text{if } FT_{wr} < EFT^* \text{ then} \\ EFT^* \leftarrow FT_{wr}, \ u^* \leftarrow C \end{array}
 13
 14
 15
                                                                                                                                                       16
                  for all m \in \mathcal{M} do
 17
                          for all c \in \mathcal{C}_m do
 18
                                   EST_{des} \leftarrow \max(DRT, FAT_{d2e}(m))
                                  \begin{split} & EST_{des} \leftarrow \max(DRT, FAT_{d2e}(m)) \\ & FT_{des} \leftarrow EST_{des} + T_{\text{transfer}}(\text{d2e}_m, d_{i,\text{d2e}_m}^{\text{send}}) \\ & RT_e \leftarrow \max(FT_{des}, \max_{v_p \in \text{pred}_E(v_i, m)} FT_{p,m,c'}^e, 0) \\ & FT_e \leftarrow \max(RT_e, FAT_{\text{edge}}(m, c)) + T_{i,m,c}^e \\ & FT_{er} \leftarrow \max(FT_e, FAT_{e2d}(m)) + T_{\text{transfer}}(\text{e2d}_m, d_{i,\text{e2d}_m}^{\text{recv}}) \\ & \text{if } FT_{er} < EFT^* \text{ then} \\ & EFT^* \leftarrow FT_{er}, \ u^* \leftarrow (m, c) \end{split}
 19
20
21
 22
 23
 24
                                                                                                                > Commit assignment and update state
                  assignment(v_i) \leftarrow u^*; tier(v_i) \leftarrow \text{tier}(u^*)
FT_{\text{final}}(v_i) \leftarrow EFT^*
25
 26
 27
                  Update all affected FAT entries to the finish times on the chosen path
```

#### B. Phase 2: Task Migration (Energy Optimization)

After Phase 1 produces a minimal-makespan schedule (minimizing  $T_{\rm total}$ ), Phase 2 iteratively refines this schedule to reduce the mobile device's energy consumption ( $E_{\rm total}$ ). This optimization is performed while strictly adhering to the application's deadline constraint ( $T_{\rm total} \leq T_{\rm max}$ ), which is set to the makespan achieved by the end of Phase 1. The core mechanism involves strategically migrating selected tasks to different execution units (other device cores, edge, or cloud) which will result in a task schedule that reduces mobile device energy consumption.

- 1) Step 2.1: Migration Candidate Generation: In each iteration of the optimization loop, the algorithm first identifies a set of potential task migrations worth evaluating. A migration is represented by a pair  $(v_{\rm tar}, u')$ , indicating that task  $v_{\rm tar}$ , currently assigned to unit  $u_{\rm curr} = {\tt assignment}(v_{\rm tar})$ , is being considered for moving to a different target unit u' ( $u' \neq u_{\rm curr}$ ). Since the goal is to minimize energy consumed by the mobile device, the algorithm primarily focuses on migrations that are most likely to reduce the sum of local-core computation energy ( $E_i^{\rm device}$ ) and RF-upload energy ( $E_i^{\rm RF}$ ).
- 1) Migrations originating from the device tier: The most promising energy-saving candidates are tasks currently running locally. For every task with  $tier(v_{tar}) = DEVICE$  (currently assigned to device core k), the following potential target units u' are considered:

- To another device core k' ( $k' \neq k$ ): This explores utilizing device heterogeneity. Moving  $v_{\rm tar}$  might save energy if the target core k' is inherently more power-efficient (i.e.,  $P_{k'} < P_k$ ). However, this potential energy gain ( $\Delta E^{\rm device}$ ) must be weighed against any potential increase in the overall makespan  $T_{\rm total}$ , which could occur if  $T^l_{{\rm tar},k'} > T^l_{{\rm tar},k}$  or if core k' is less available ( $FAT_{\rm core}(k')$  is later). The resulting schedule's makespan must still satisfy  $T_{\rm max}$ .
- To the cloud C: This considers offloading the computation entirely. Moving  $v_{\rm tar}$  eliminates its local computation energy ( $E_{\rm tar}^{\rm device}$ ) but incurs energy cost for the RF upload:

$$E_{\text{tar}}^{\text{RF}} = P_{\text{d2c}} \times T_{\text{transfer}} (\text{d2c}, d_{\text{tar,d2c}}^{\text{send}}).$$

This migration is considered a valid candidate only if the saved computation energy is greater than the incurred RF energy ( $E_{\rm tar}^{\rm device} > E_{\rm tar}^{\rm RF}$ ). The projected final finish time of the task via the cloud path ( $FT_{\rm tar}^{\it wr}$ ) after rescheduling must not cause the overall application makespan to exceed  $T_{\rm max}$ .

• To an edge core (m,c): Similar to cloud offloading, moving  $v_{\rm tar}$  to edge core (m,c) eliminates  $E_{\rm tar}^{\rm device}$  and introduces RF upload energy specific to the edge path:

$$E_{\text{tar}}^{\text{RF}} = P_{\text{d2e}} \times T_{\text{transfer}} (\text{d2e}_m, d_{\text{tar,d2e}_m}^{\text{send}}).$$

This migration is considered a valid candidate only if energy saved outweighs the energy spent on upload ( $E_{\rm tar}^{\rm device} > E_{\rm tar}^{\rm RF}$ ). Additionally, the resulting makespan, influenced by the edge path's final finish time ( $FT_{{\rm tar},m}^{er}$ ), must remain within the deadline  $T_{\rm max}$ . This evaluation is performed for all available edge cores (m,c).

- 2) Exclusion of other migration types: This heuristic excludes two classes of migrations, since neither reduces mobile-side energy  $E_{\text{total}}$ :
  - Remote-to-device moves: Any migration from edge or cloud back to the device reintroduces local computation energy  $E_{\rm tar}^{\rm device}$ , directly contradicting the objective of minimizing mobile energy.
  - Remote-to-remote moves: Migrations between edge and cloud (or between different edge nodes) do not affect the device's computation or RF-upload energy, and thus have no impact on the primary mobile device energy optimization target  $E_{\rm total}$ .
- 3) Construct Candidate Set M: Based on the above criteria, form the set

$$\mathcal{M} = \big\{ (v_{\text{tar}}, u') \mid \text{tier}(v_{\text{tar}}) = \text{DEVICE}, \\ u' \in (\mathcal{K} \setminus \{k\}) \cup \{C\} \cup \{(m, c)\} \big\}.$$

 ${\cal K}$  is the set of device cores, and (m,c) ranges over all edge cores. Early pruning by the condition

$$E_{\rm tar}^{
m device} > E_{
m tar}^{
m RF}$$

may be applied at this stage to reduce  $\mathcal{M}$  before detailed evaluation.

Each potentially beneficial candidate migration  $(v_{\text{tar}}, u') \in \mathcal{M}$  identified in this step is then passed to Step 2.2 for evaluation using the kernel rescheduling algorithm to determine its impact on the overall application execution schedule makespan  $T_{\text{total}}$  and the mobile device energy consumption  $E_{\text{total}}$ .

2) Step 2.2: Migration Evaluation: Once the candidate set  $\mathcal{M}$  of migrations  $(v_{\text{tar}}, u')$  is generated (Step 2.1), we must efficiently predict each move's effect on:

$$T'_{\text{total}} = \max_{v_i \in V} FT_{\text{final}}(v_i), \quad E'_{\text{total}} = \sum_{v_i \in V} \left( E_i^{\text{device}} + E_i^{\text{RF}} \right).$$

Re-running Phase 1 for every candidate is not efficient. Instead, the kernel rescheduling procedure of Lin et al[1] is adapted to the three-tier model. The kernel rescheduling algorithm updates only the affected task and its dependents in time  $\mathcal{O}(|V| + |E|)$  allowing efficient recomputation of all EST/FT values after a single assignment change.

The evaluation for a single candidate migration  $(v_{\text{tar}}, u')$  involves simulating the change on a temporary copy of the current schedule state (which includes all task attributes like finish times, assignments, and the execution sequences S) using two sub-steps:

- 1) **Sequence Construction:** The task execution sequences  $S_u$  are modified to simulate the proposed migration. This involves:
  - Removal: Task  $v_{\text{tar}}$  is located in and removed from its current execution sequence  $S_{u_{\text{curr}}}$ , where  $u_{\text{curr}} = \text{assignment}(v_{tar})$ .
  - Insertion: Task  $v_{\rm tar}$  is inserted into the sequence  $S_{u'}$  corresponding to the target unit u'. The insertion point is determined using binary search to maintain an order that respects task readiness, ensuring  $v_{\rm tar}$  is placed after tasks likely to start earlier. The key used for this binary search depends on the target unit u':
    - If u' is a device core k', the key is the task's earliest ready time for local execution,  $RT_{tar}^{l}$ .
    - If u' is the cloud C, the key is the earliest ready time for cloud upload,  $RT_{\text{tar}}^{ws}$ .
    - If u' is an edge core (m,c), the key is the estimated finish time of the device-to-edge upload,  $FT_{\mathrm{tar},m}^{\mathrm{des,est}} = RT_{\mathrm{tar},\mathrm{d2e}_m} + T_{\mathrm{transfer}}(\mathrm{d2e}_m,d_{tar,\mathrm{d2e}_m}^{\mathrm{send}}).$  This heuristic approximates the task's readiness at the edge for sequencing purposes before the full reschedule.
  - Temporary State Update: On the copied task object for  $v_{\text{tar}}$ , its assignment and tier attributes are updated to indicate the target unit u', and its scheduling status state( $v_{tar}$ ) is reset to UNSCHEDULED, indicating it needs to be rescheduled by the kernel.

- This sub-step produces a new set of task sequences S' simulating the hypothetical post-migration ordering.
- 2) **Kernel Scheduling:** The kernel algorithm is executed using the modified sequences S' and the temporarily updated task attributes. It recalculates all task timings (EST and FT values:  $FT_i^l, FT_i^{ws}, \ldots, FT_{i,m}^{er}, FT_{\text{final}}(v_i)$ ) for the entire DAG in linear time with respect to the number of tasks and edges  $(\mathcal{O}(|V| + |E|))$ . This process correctly propagates the timing changes resulting from the migration throughout the graph, respecting all dependencies and current resource availabilities (FAT).

After the kernel algorithm completes its run on the temporary schedule copy, the projected makespan  $T'_{\text{total}}$  is determined by finding the maximum  $FT_{\text{final}}(v_i)$  among all exit tasks. The projected mobile energy consumption  $E'_{\text{total}}$  is recalculated based on the new assignments and potentially altered execution/transfer durations seen in the updated FT values. These projected metrics  $(T'_{\text{total}}, E'_{\text{total}})$  quantify the expected outcome of the candidate migration and are passed to the selection logic in Step 2.3. This entire two-sub-step evaluation process is repeated for every candidate migration  $(v_{\text{tar}}, u')$  in the set  $\mathcal{M}$ .

3) Step 2.3: Heuristic Migration Selection: After evaluating all candidate migrations  $(v_{\text{tar}}, u') \in \mathcal{M}$  via the kernel rescheduling algorithm (Step 2.2) and obtaining their respective projected makespans  $(T'_{\text{total}})$  and mobile-energy consumptions  $(E'_{\text{total}})$ , this step implements the heuristic decision rule to select at most *one* migration to apply permanently in the current optimization iteration. The heuristic prioritizes maximal energy reduction, first considering only options that do not increase the makespan, and then considering energy-efficient time/energy trade-offs, always respecting the overall deadline  $T_{\text{max}}$ .

 $T_{\text{total}}$  and  $E_{\text{total}}$  represent the makespan and mobile energy of the schedule before applying any migration in this iteration. For each evaluated candidate (represented by (v, u', T', E') containing the task, target unit, projected time, and projected energy), the change in energy  $\Delta E = E_{\text{total}} - E'$  and the change in makespan  $\Delta T = T' - T_{\text{total}}$  is calcualated. A positive  $\Delta E$  indicates energy savings, while a positive  $\Delta T$  indicates a makespan increase.

The selection follows a two-criteria approach:

- 1) Filter Valid Candidates: The set of evaluated candidates is filtered to include only candidates that are potentially beneficial and feasible. A candidate (v, u', T', E') is considered valid if it strictly reduces mobile energy (ΔE > 0) and its resulting makespan does not violate the deadline (T' ≤ T<sub>max</sub>). M<sub>valid</sub> represents this filtered set. If M<sub>valid</sub> is empty, no migration is possible in this iteration (proceed to step 4).
- Criteria 1 (Maximize Energy Reduction with No Time Increase): Prioritize migrations that offer energy savings without negatively impacting performance. Iden-

tify the subset  $\mathcal{M}_a = \{(v, u', T', E') \in \mathcal{M}_{\text{valid}} \mid \Delta T \leq 0\}$  containing valid migrations that do not increase the makespan. If  $\mathcal{M}_a$  is not empty, the migration  $(v_{\text{best}}, u'_{\text{best}})$  is selected from  $\mathcal{M}_a$  that provides the largest energy reduction (maximum  $\Delta E$ ):

$$(v_{\text{best}}, u'_{\text{best}}, T'_{\text{best}}, E'_{\text{best}}) = \underset{(v, u', T', E') \in \mathcal{M}_a}{\arg\max} \{\Delta E\}.$$

If a best migration is found under this criterion, it is chosen, and the algorithm proceeds to Step 2.4 to apply it.

3) Criteria 2 (Maximize Energy-per-Time Trade-off Efficiency): If no migration satisfies Criteria 1 (i.e.,  $\mathcal{M}_a$  is empty), the algorithm considers valid migrations that do increase the makespan but still meet the deadline. A subset  $\mathcal{M}_b = \mathcal{M}_{\text{valid}} \setminus \mathcal{M}_a$  is generated which contains valid migrations where  $\Delta T > 0$ . If  $\mathcal{M}_b$  is not empty,the migration  $(v_{\text{best}}, u'_{\text{best}})$  from  $\mathcal{M}_b$  that offers the best trade-off is selected, specifically the migration candidate that maximizes the ratio of energy saved per unit of makespan increase  $(\Delta E/\Delta T)$ :

$$(v_{\text{best}}, u_{\text{best}}', T_{\text{best}}', E_{\text{best}}') = \mathop{\arg\max}_{(v, u', T', E') \in \mathcal{M}_b} \left\{ \frac{\Delta E}{\Delta T} \right\}.$$

If a best migration is found under this criteria, it is chosen, and the algorithm proceeds to Step 2.4.

4) No Selection / Termination: If, after checking both criteria, no suitable migration is selected (i.e.,  $\mathcal{M}_{valid}$  was initially empty, or both  $\mathcal{M}_a$  and  $\mathcal{M}_b$  are empty), it indicates that no further energy reduction is possible under the heuristic rules within the deadline constraint. No migration is applied in this iteration, and the optimization process terminates, as described in Step 2.4.

This deterministic selection process ensures that the algorithm greedily pursues the most impactful energy savings available at each step, while performance (effect on makespan) is only traded off when necessary and efficient, always respecting the hard deadline  $T_{\max}$ .

#### Algorithm 4 SELECTMIGRATION( $\mathcal{L}, T_{\text{total}}, E_{\text{total}}, T_{\text{max}}$ )

Require: Evaluated list  $\mathcal{L} = \{(v, u', T', E')\}$  from Step 2.2; current makespan  $T_{\text{total}}$ ; current mobile energy  $E_{\text{total}}$ ; deadline  $T_{\text{max}}$  Ensure: best migration tuple or null 1 for all  $(v, u', T', E') \in \mathcal{L}$  do 2  $\Delta E \leftarrow E_{\text{total}} - E', \Delta T \leftarrow T' - T_{\text{total}}$  3 store  $(v, u', \Delta E, \Delta T, T', E')$  in  $\mathcal{L}'$  4  $\mathcal{M}_{\text{valid}} \leftarrow \{x \in \mathcal{L}' \mid \Delta E > 0 \land T' \leq T_{\text{max}}\}$  5  $\mathcal{M}_a \leftarrow \{x \in \mathcal{M}_{\text{valid}} \mid \Delta T \leq 0\}$  6 if  $\mathcal{M}_a \neq \emptyset$  then 7 return  $\arg\max_{x \in \mathcal{M}_a} \{\Delta E\}$  8  $\mathcal{M}_b \leftarrow \mathcal{M}_{\text{valid}} \setminus \mathcal{M}_a$  9 if  $\mathcal{M}_b \neq \emptyset$  then 10 return  $\arg\max_{x \in \mathcal{M}_b} \{\frac{\Delta E}{\Delta T}\}$ 

- 4) Step 2.4: Iteration and Convergence: This step concludes one iteration of Phase 2 and determines if the energy optimization loop continues.
- 1) **Apply Selected Migration:** If a migration  $(v_{best}, u'_{best})$  was chosen in Step 2.3:
  - The permanent assignment attributes  $assignment(v_{best})$  and  $tier(v_{best})$  are updated to indicate the new unit  $u'_{best}$ .
  - The modified execution sequences S' that were generated during the Sequence Construction sub-step (Step 2.2) for the evaluation of the chosen migration (v<sub>best</sub>, u'<sub>best</sub>) are retained.
  - All EST/FT values  $(FT_i^l, FT_i^{ws}, \dots, FT_{\text{final}})$  for all tasks  $v_i$  are updated to match the results obtained from the Kernel Scheduling sub-step (Step 2.2) for the selected migration.
  - The global task scheduling performance metrics  $T_{\rm total}$  and  $E_{\rm total}$  are updated based on the newly applied and rescheduled state.

The algorithm then loops back to Step 2.1 (candidate generation) for the next iteration, starting from this updated schedule state.

2) Check Convergence: If Step 2.3 resulted in no migration candidates being selected (either M<sub>valid</sub> was empty or neither criteria identified a suitable migration candidate), it indicates that no further improvements satisfying the heuristic criteria can be found. The optimization loop terminates.

The schedule  $\mathcal{S}^{\star}$  represents when the loop terminates is the final output of the heuristic optimization. It consists of the final task assignments (assignment( $v_i$ ), tier( $v_i$ )), the corresponding execution sequences  $S_u$ , all calculated EST/FT values, and the final makespan  $T_{\text{total}}$  and mobile energy  $E_{\text{total}}$ . By design, this final schedule satisfies  $T_{\text{total}} \leq T_{\text{max}}$  and represents a local optimum for mobile energy consumption under the defined heuristic migration strategy.

#### **Algorithm 5** Heuristic Energy-Optimization Loop (Phase 2)

```
Require: Initial schedule S_0; deadline T_{\text{max}}; task set V; DAG G = (V, E);
        resource sets \mathcal{K}, \mathcal{M}_{\text{nodes}}, \mathcal{C}_m; execution times T(\cdot); data sizes d(\cdot); rates r(\cdot);
power models P_{\mathrm{core}}(\cdot), P_{\mathrm{RF}}(\cdot)

Ensure: Schedule \mathcal{S}^{\star} with T_{\mathrm{total}} \leq T_{\mathrm{max}} and locally minimal E_{\mathrm{total}}
    1 \mathcal{S} \leftarrow \mathcal{S}_0
   2 T_{\text{total}} \leftarrow \text{CalculateMakespan}(\mathcal{S})
   3 E_{\text{total}} \leftarrow \text{CalculateMobileEnergy}(S)
               did\_migrate \leftarrow false
       T_{	ext{curr}} \leftarrow T_{	ext{total}}; E_{	ext{curr}} \leftarrow E_{	ext{total}}
Step 2.1 — Generate Migration Candidates
                \mathcal{M}_{cand} \leftarrow \emptyset
               for all v_i \in V with unit(v_i) = core_k \in \mathcal{K} do
                      E_{i,k}^l \leftarrow P_{\text{core}}(k) T_{i,k}^l
 10
                      for all k' \in \mathcal{K} \setminus \{k\} do
                             \mathcal{M}_{cand} \cup = \{(v_i, core_{k'})\}
 11
                      \begin{array}{l} \mathcal{M}_{\text{cand}} \cup = \{(v_i, col_{R'})\} \\ E_i^{\text{d2c}} \leftarrow P_{\text{RF}}(\text{d2c}) \ T_{\text{transfer}}(\text{d2c}, d_{i, \text{d2c}}) \\ \text{if } E_{i,k}^l > E_i^{\text{d2c}} \ \text{then} \\ \mathcal{M}_{\text{cand}} \cup = \{(v_i, C)\} \end{array} 
 13
 14
 15
                      for all m \in \mathcal{M}_{\text{nodes}} do
                            \begin{array}{l} \text{Transfer} \in \mathcal{S}_{\text{node}} \text{ sector} \\ E^{\text{d2e}}_{i,m} \in \mathcal{P}_{\text{RF}}(\text{d2e}_m) \ T_{\text{transfer}}(\text{d2e}_m, d_{i,\text{d2e}_m}) \\ \text{if } E^l_{i,k} > E^{\text{d2e}}_{i,m} \text{ then} \\ \text{for all } c \in \mathcal{C}_m \text{ do} \end{array}
 16
 17
 18
 19
                                           \mathcal{M}_{\text{cand}} \cup = \{(v_i, (m, c))\}
       Step 2.2 — Evaluate Candidates via Kernel Rescheduling
               \mathcal{R}_{eval} \leftarrow \emptyset
 20
               for all (v, u') \in \mathcal{M}_{cand} do
21
                      \begin{array}{l} (T^{\prime},E^{\prime}) \leftarrow \texttt{KernelEvaluate}(\mathcal{S},(v,u^{\prime})) \\ \mathcal{R}_{\text{eval}} \cup = \{(v,u^{\prime},T^{\prime},E^{\prime})\} \end{array} 
 22
 23
       Step 2.3 — Select Best Valid Migration
               26
              if \mathcal{M}_a \neq \emptyset then
                     (v^{\star}, u^{\star}, T^{\star}, E^{\star}) \leftarrow \arg\max_{x \in \mathcal{M}_a} (E_{\text{curr}} - E')
 28
                      \mathcal{M}_b \leftarrow \mathcal{M}_{\text{valid}} \setminus \mathcal{M}_a
                     if \mathcal{M}_b \neq \emptyset then (v^\star, u^\star, T^\star, E^\star) \leftarrow \arg\max_{x \in \mathcal{M}_b} \left( (E_{\mathrm{curr}} - E')/(T' - T_{\mathrm{curr}}) \right)
 30
 31
 32
                      else
 33
                             (v^{\star}, u^{\star}) \leftarrow \text{null}
       Step 2.4 — Apply Migration and Check Convergence
 34
              if (v^*, u^*) \neq \text{null then}
                      \mathcal{S} \leftarrow \text{APPLYMIGRATIONANDRESCHEDULE}(\mathcal{S}, (v^{\star}, u^{\star}))
 35
                      T_{\text{total}} \leftarrow \text{CalculateMakespan}(S)
                      E_{\text{total}} \leftarrow \text{CalculateMobileEnergy}(S)
                      did\_migrate \leftarrow true
 39 until not did_migrate
 40 return S
```

#### VI. Q-LEARNING ENHANCED ENERGY OPTIMIZATION

After establishing an initial three-tier schedule that minimizes delay (Phase 1), an alternative approach for Phase 2 energy optimization using Q-learning, a model-free reinforcement learning algorithm is tested. Unlike the deterministic heuristic migration strategy (Section V-B), Q-learning adaptively explores the vast solution space of possible task migrations across the device, edge, and cloud tiers. It learns a policy, encoded in a Q-table, to select migrations that optimize long-term rewards, potentially discovering more effective energy-saving strategies, especially in complex scenarios with multiple dependencies and resource interactions.

#### A. Q-Learning Framework Components

Q-learning models the task migration problem as a Markov Decision Process (MDP), where the agent learns to navigate the scheduling state space by taking migration actions to maximize cumulative rewards related to energy efficiency and maintain makespan deadline constraints.

- 1) State Representation  $(s \in S)$ : A state s provides a snapshot of the current scheduling configuration, representing task assignments and scheduling performance metrics and resource characteristics crucial for informed decision-making. A state s incorporates:
- Task Assignments: The assigned execution unit assignment $(v_i)$  for each task  $v_i$ .
- Current Performance: The current schedule makespan  $T_{\rm total}$  and mobile device energy  $E_{\rm total}$ .
- **Tier Distribution:** The number of tasks assigned to each tier (DEVICE, EDGE, CLOUD).
- **Tier Characteristics:** Average computational complexity (complexity) and data intensity (intensity) of tasks within each tier.

To manage the state space, continuous values like  $T_{\rm total}$  and  $E_{\rm total}$  are discretized and rounded before being used as keys for Q-table lookups.

- 2) Action Space  $(a \in \mathcal{A})$ : An action a corresponds to a potential task migration, represented by the pair  $(v_{tar}, u')$ , where  $v_{tar}$  is the task to be migrated and u' is the target execution unit (a specific device core k', the cloud C, or a specific edge core (m,c)), such that  $u' \neq \mathtt{assignment}(v_{tar})$ . The set of valid actions  $\mathcal{A}(s)$  from a given state s includes possible migrations  $(v_{tar}, u')$  not involving the task's current location.
- 3)  $Q ext{-Value}(Q(s,a))$ : The core of the learning process is the Q-table, which stores the Q-value, Q(s,a). This value represents the expected cumulative discounted future reward obtainable by taking action a in state s and following the optimal policy. Higher Q-values indicate more promising state-action pairs.
- 4) Learning Update Rule: The Q-values are learned iteratively using the Bellman equation update rule after observing a transition (s, a, r, s'):

$$Q(s, a) \leftarrow Q(s, a) + \alpha \Big[ r + \gamma \max_{a'} Q(s', a') - Q(s, a) \Big]$$

where

- $\alpha \in (0,1]$  is the *learning rate*, controlling how much new information overrides old estimates, often decayed over time.
- $\gamma \in [0,1)$  is the *discount factor*, weighting the importance of future rewards.
- r is the immediate reward received after taking action a in state s.
- s' is the successor state following action a.
- $\max_{a'} Q(s', a')$  is the maximum estimated future value from state s'.

This update minimizes the Temporal Difference (TD) error.

5) Reward Function (r): The reward function r guides the learning process by quantifying the immediate effect of a transition from state s to s' resulting from action a.

r =EnergyReward - TimePenalty

where:

- Energy Reward: Primarily driven by the energy reduction  $\Delta E = E_{\text{total}}(s) E_{\text{total}}(s')$ . Penalties apply if  $\Delta E < 0$ .
- Time Penalty: Applied if the deadline  $T_{\text{max}}$  is violated in state s', scaling quadratically with the magnitude of the violation  $\max(0, T_{\text{total}}(s') T_{\text{max}})$ .

This structure encourages energy savings while strongly discouraging deadline violations.

## B. Q-Learning Algorithm Execution

The optimization process runs over multiple episodes, with each episode consisting of several iterations (state transitions or steps).

- 1) **Initialization:** Start with the schedule S from Phase 1. Initialize the Q-table and replay buffer. Set initial hyperparameters  $(\alpha, \gamma, \epsilon)$ .
- 2) **Episode Loop:** Repeat for a maximum number of episodes:
  - a) Reset episode state to the current best-known schedule.
  - b) **Iteration Loop:** Repeat for a maximum number of steps per episode or until convergence:
    - Get current state s.
    - Action Selection (Epsilon-Greedy): Select action  $a \in \mathcal{A}(s)$  either randomly (exploration, probability  $\epsilon$ ) or greedily ( $a = \arg\max_{a'} Q(s, a')$ , exploitation).
    - Action Evaluation (Simulation): Predict the outcome metrics (T', E') of taking action a by performing the sequence construction and kernel rescheduling described in Step 2.2 on a copy of the current state.
    - Reward Calculation: Calculate the immediate reward r based on the transition from the current  $(T_{\mathtt{total}}, E_{\mathtt{total}})$  to the predicted (T', E').
    - Action Execution (State Update): Permanently apply the chosen migration action a to the actual schedule  $\mathcal{S}$  by modifying the task assignment and sequences, then running the kernel rescheduling algorithm to update all timings (FTs) and resource availabilities (FATs). This transitions the system to the next state s'. Update the current  $T_{\text{total}}$  and  $E_{\text{total}}$ .
    - **Q-Table Update:** Update Q(s, a) using the observed reward r, the resulting state s', and the Bellman equation.

- Experience Replay: Store the transition (s, a, r, s') in the replay buffer. Sample experiences from the buffer and perform additional Q-value updates.
- Update  $s \leftarrow s'$ .
- Check termination conditions.
- c) **Decay Parameters:** Decrease exploration rate  $\epsilon$  and potentially learning rate  $\alpha$ .
- 3) **Output:** Return the best schedule (lowest valid energy respecting  $T_{\rm max}$ ) encountered.

This adaptive learning process allows the Q-agent to potentially find superior migration policies compared to the fixed rules of the heuristic approach.

# **Algorithm 6** Q-Learning-Enhanced Energy Optimization (Phase 2 – RL variant)

```
Require: Initial schedule S_0; deadline T_{\text{max}}; task set V; DAG G = (V, E);
       Q-learning hyper-parameters (\alpha, \gamma, \epsilon_{\text{start}}, \epsilon_{\text{end}}, \epsilon_{\text{decay}}); max episodes N_{\text{eps}}; max
steps N_{\text{steps}}; replay buffer B_{\text{replay}}; system parameters \mathcal{P} Ensure: Schedule \mathcal{S}^{\star} with T_{\text{total}}\!\leq\!T_{\max} and minimal learned E_{\text{total}}
   1 Initialise Q(s, a) \leftarrow 0 for all state-action pairs
  2 \ \mathcal{S}^{\star} \leftarrow \mathcal{S}_0;
  3 E^* \leftarrow \text{CALCULATEMOBILEENERGY}(S^*);
  4 T^* \leftarrow \text{CalculateMakespan}(S^*)
  6 for episode = 1 to N_{\rm eps} do
            \begin{array}{l} \mathcal{S} \leftarrow \operatorname{copy}(\mathcal{S}^{\star}) \\ T_{\text{total}} \leftarrow T^{\star}; E_{\text{total}} \leftarrow E^{\star} \end{array}
            for step = 1 to N_{\rm steps} do
      Step 2.1 — Candidate Generation
                    s \leftarrow \text{GetCurrentState}(\mathcal{S}, T_{\text{total}}, E_{\text{total}})
10
11
                    \mathcal{A}(s) \leftarrow \text{GETVALIDACTIONS}(s)
                   if A(s) = \emptyset then break
12
      Step 2.2 — Candidate Evaluation
13
                    if rand() < \epsilon then
                          a \leftarrow \text{random choice from } \mathcal{A}(s)
14
15
16
                          a \leftarrow \arg\max_{a' \in \mathcal{A}(s)} Q(s, a')
      (T',E') \leftarrow \texttt{EVALUATEMIGRATION}(\texttt{copy}(\mathcal{S}),a,\mathcal{P}) 
 Step 2.3 — Migration "Selection"
17
                    r \leftarrow \text{CALCULATEREWARD}(E_{\text{total}}, T_{\text{total}}, E', T', T_{\text{max}})
18
       Step 2.4 — Apply Migration & Learn
10
                    \mathcal{S} \leftarrow \overrightarrow{\mathsf{APPLYMIGRATIONANDRESCHEDULE}}(\mathcal{S}, a, \mathcal{P})
20
                    T_{\text{total}} \leftarrow \text{CalculateMakespan}(S)
                    E_{\text{total}} \leftarrow \text{CALCULATEMOBILEENERGY}(S)
21
22
                          \leftarrow GETCURRENTSTATE(S, T_{\text{total}}, E_{\text{total}})
23
                    \mathsf{bestQ} \leftarrow \max_{a'} Q(s', a')
24
                    Q(s,a) \leftarrow Q(s,a) + \alpha \left[r + \gamma \cdot \text{bestQ} - Q(s,a)\right]
Add (s,a,r,s') to B_{\text{replay}}
25
                    Sample mini-batch B_{\text{sample}} \subset B_{\text{replay}} for all (s_j, a_j, r_j, s'_j) \in B_{\text{sample}} do
26
27
                          \mathsf{bestQ}_j \leftarrow \mathsf{max}_{a'} \, Q(s'_j, a')
28
29
                          Q(s_j, a_j) \leftarrow Q(s_j, a_j) + \alpha \left[ \left. r_j + \gamma \cdot \mathrm{best} \mathbf{Q}_j - Q(s_j, a_j) \right. \right]
30
                    if T_{\mathrm{total}} \leq T_{\mathrm{max}} \wedge E_{\mathrm{total}} < E^{\star} then
31
                                \leftarrow \operatorname{copy}(\mathcal{S}); E^{\star} \leftarrow E_{\operatorname{total}}; T^{\star} \leftarrow T_{\operatorname{total}}
32
              \epsilon \leftarrow \max(\epsilon_{\text{end}}, \epsilon \times \epsilon_{\text{decay}})
             (optional) decay \alpha
34 return S^*
```

# VII. EVALUATION

This section evaluates the performance and energy efficiency of the proposed three-tier scheduling frameworks compared to the baseline two-tier heuristic.

1) Quantify the energy savings  $(E_{total})$  and makespan  $(T_{total})$  impact of incorporating the edge tier and dynamic resource models.

- 2) Compare the effectiveness of the extended three-tier heuristic migration (3T-H) against the adaptive Q-learning based migration (3T-QL) for energy optimization under deadline constraints ( $T_{max}$ ).
- Analyze the sensitivity of the algorithms to variations in key system parameters, network conditions, device states, and application characteristics.

## A. Simulation Setup

Experiments are carried out in a Python based simulator that incorporates the dynamic resource models of Section IV and supports the full three-tier architecture. Application workloads are expressed as DAGs generated with parameters for task count |V|, graph density, and task-type distribution (compute-intensive, data-intensive, balanced). Task attributes complexity( $v_i$ ) and intensity( $v_i$ ) and the per-link data sizes  $d_{i,\ell}^{\mathrm{send/receive}}$  are assigned stochastically according to task type.

Evaluated scheduling frameworks:

- 2T-H (baseline). Lin et al. [1] two-phase heuristic adapted to the proposed dynamic models in a two-tier (device-cloud) setting.
- **3T-H.** Proposed two-phase heuristic operating in the three-tier (device–edge–cloud) environment (Section V).
- **3T-QL.** Three-tier framework that shares the Phase-1 schedule of 3T-H and applies Q-learning for Phase-2 energy optimization (Section VI) with default hyper-parameters.

Parameter sweeps: Each test suite varies a single parameter while other parameters are held at nominal values:

- Battery level B: 20, 50, 80, 100 %.
- Bandwidth scaling factor  $f_{\text{BW}}$ : 0.5, 0.8, 1.1, 1.4, 1.7, 2.0.
- Number of edge nodes M: 0, 1, 2, 3.
- Number of device cores K: 2, 3, 4, 6, 8.
- Task count |V|: 20, 30, 50, 75.

Each configuration is executed three times with independent random seeds; results are reported as mean  $\pm$  standard deviation. The deadline is defined as  $T_{\rm max}=1.5\,T_{\rm min}$ , where  $T_{\rm min}$  is the minimal makespan produced by Phase 1 of the 2 step scheduling algorithm.

## B. Results and Analysis

This section presents the comparative results of the 2T-H, 3T-H, and 3T-QL frameworks, focusing on the insights gained from parameter sweeps. All plots show the mean results over three repetitions, with shaded areas representing  $\pm 1$  standard deviation.

1) Impact of Edge Tier and Dynamic Models (3T-H vs. 2T-H): Comparing the 3T-H framework against the 2T-H baseline consistently demonstrated the benefits of the edge tier across the tested parameter ranges (Figs. 1–5).

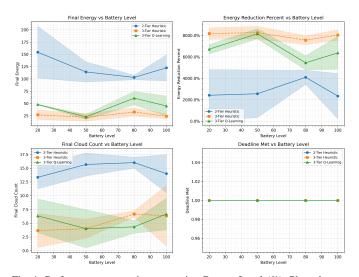


Fig. 1: Performance comparison sweeping Battery Level (%). Plots show Final Energy, Energy Reduction Percentage, Final Cloud Task Count, and Deadline Met Rate as battery level varies from 20% to 100%.

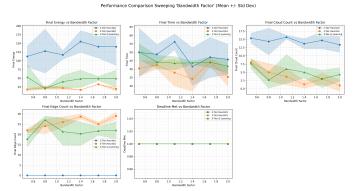


Fig. 2: Performance comparison sweeping Network Bandwidth Factor  $(f_{BW})$ . Plots show Final Energy, Final Time, Final Cloud Task Count, Final Edge Task Count, and Deadline Met Rate as bandwidth factor varies from 0.5 to 2.0.

- Energy Consumption: 3T-H consistently achieved significantly lower final mobile energy  $(E_{total})$  compared to 2T-H. As seen across all sweeps, the orange line (3T-H) lies substantially below the blue line (2T-H). The advantage was particularly highlighted at low battery levels  $(B \le 50\%, \text{ Fig. 1})$ , where 2T-H energy increased sharply while 3T-H remained relatively low. The presence of even one edge node (M = 1 vs. M = 0 in Fig. 5) caused a significant decrease in energy for 3T-H.
- Makespan and Deadline Constraint Adherence: 3T-H often finished faster (lower final  $T_{\rm total}$ ) than 2T-H after energy optimization, particularly noticeable when varying the number of mobile device cores (Fig. 3) or edge nodes (Fig. 5). This indicates the edge tier can improve performance even when the primary goal is energy saving. All frameworks displayed perfect deadline constraint adherence (Deadline Met = 1.0) across all sweeps.

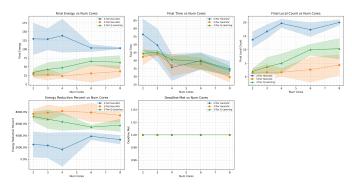


Fig. 3: Performance comparison sweeping Number of Device Cores (K). Plots show Final Energy, Final Time, Final Local Task Count, Energy Reduction Percentage, and Deadline Met Rate as core count varies from 2 to 8.

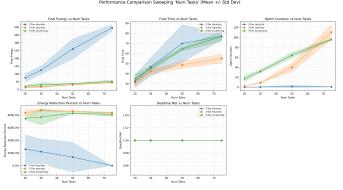


Fig. 4: Performance comparison sweeping Number of Tasks (|V|). Plots show Final Energy, Final Time, Optimization Duration, Energy Reduction Percentage, and Deadline Met Rate as task count varies from 20 to 75.

- Task Distribution: The results highlight the strategic use of the edge by 3T-H. The final Edge Count or the final count of number of tasks scheduled for execution on the edge tier in Fig. 5 increases sharply from 0 (for 2T-H at M=0) to over 25 tasks by M=2. Conversely, Fig. 2 shows edge tier utilization (Final Edge Count) peaking at lower bandwidth factors ( $f_{BW}\approx 0.8$ ) and then decreasing as bandwidth improves, while the Final Cloud Count generally increases, indicating a shift towards cloud when network conditions allow. This adaptability is absent in 2T-H.
- 2) Comparison of Optimization Strategies (3T-QL vs. 3T-H): Comparing the Q-learning optimization (3T-QL, green triangles) against the heuristic migration (3T-H, orange squares) for the energy optimization phase highlighted that under the tested conditions and default hyper-parameters, the heuristic consistently outperformed the Q-learning approach (Figs. 1–4):
- Energy Savings: Across all parameter sweeps, the 3T-H heuristic achieved lower final energy (E<sub>total</sub>) than 3T-QL.
   The orange line consistently sits below the green line in the Final Energy plots. While 3T-QL reduced energy

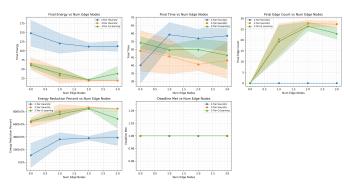


Fig. 5: Performance comparison sweeping Number of Edge Nodes (M). Plots show Final Energy, Final Time, Final Edge Task Count, Energy Reduction Percentage, and Deadline Met Rate as edge node count varies from 0 to 3.

- significantly compared to 2T-H, it did not match the effectiveness of the 3T-H heuristic in these experiments. The Energy Reduction Percent plots also show 3T-H generally achieving higher reduction percentages than 3T-QL.
- Makespan and Deadline Adherence: Both 3T-H and 3T-QL maintained perfect deadline constraint adherence (Deadline Met = 1.0). 3T-H consistently resulted in a lower final makespan ( $T_{\text{total}}$ ) compared to 3T-QL across all sweeps. The difference was noticeable when varying core counts, edge nodes, and bandwidth factors.
- Computational Overhead: The Optimization Duration vs Number of Tasks plot (Fig. 4) clearly shows that the optimization phase for 3T-QL required substantially more time than for 3T-H, and this overhead grew rapidly with the number of tasks. The heuristic's optimization duration remained very low and scaled minimally.
- Task Distribution Patterns: The final task distributions were similar, with both 3T frameworks utilizing the edge heavily when available (Fig. 5) and shifting load between edge and cloud based on bandwidth (Fig. 2).
- Potential Q-Learning Challenges: The consistent underperformance of 3T-QL relative to 3T-H in these tests suggests challenges in its current configuration for this specific problem. Potential reasons include sub-optimal default hyper-parameters, an inadequate state representation that does not fully capture system dynamics impacting the heuristic's effectiveness, or a reward function that does not perfectly align learning with the dual objectives under dynamic conditions.

Based on these results, the extended three-tier heuristic (3T-H) emerges as the superior algorithm, delivering better energy efficiency and faster execution times with significantly lower computational overhead compared to the tested Q-learning implementation (3T-QL).

3) Evaluation across Parameter Ranges: Analysis under varying conditions of the parameter sweeps provide insights into how the dynamic resource modeling (Section IV), real-

istic variations in the system and application characteristics influence scheduling decisions and outcomes. The results from Figs. 1–5, align well with expected real-world behaviors:

- Battery Level (B, Fig. 1): Decreasing battery levels below 50% increases final mobile energy ( $E_{total}$ ) for all frameworks. This is a direct consequence of the battery factor (bf) in the power model (Section IV), which realistically simulates increased device power draw at low charge states. The 2T-H framework, reliant only on cloud offloading (which involves energy-intensive RF), is most affected. The 3T frameworks demonstrate greater resilience by shifting tasks towards the edge as indicated by the stable or slightly increasing edge/cloud counts relative to local increases in energy cost, mitigating the impact of low battery. This aligns with the realistic behavior that systems would avoid power-hungry operations or seek lower-energy offload options when battery is critical. The optimal energy point observed between 50–80% indicates a balance where the device is efficient, but offloading still offers benefits.
- Network Bandwidth  $(f_{BW}, \text{Fig. 2})$ : Network quality effects the offloading strategy, as expected in realistic system behavior. Low bandwidth  $(f_{BW} \leq 0.8)$  makes long-distance cloud transfers slow and energy-intensive due to RF activity captured by  $E_i^{\text{RF}}$  depending on  $T_{\text{transfer}}$ . Consequently, the 3T frameworks heavily utilize the closer edge tier (the Final Edge Count peaks), minimizing transfer overhead. As bandwidth improves  $(f_{BW} \geq 1.4)$ , the time and energy cost of cloud transfers decrease, making the powerful cloud resources more attractive; both 3T-H and 3T-QL adapt by increasing cloud usage (Final Cloud Count increases) and reducing edge usage. This dynamic shift between edge and cloud based on network conditions is a key realistic behavior enabled by the three-tier model and bandwidth scaling.
- Edge Nodes (M, Fig. 5): Increasing the number of edge devices highlights the realistic benefit of edge device availability in a cloud computing environment. Introducing even one edge node (M = 1) provides substantial energy and makespan reductions for 3T frameworks compared to the no-edge baseline (M = 0, equivalent to 2T-H), by offering a low-latency, moderate-energy offload target. Adding a second node (M = 2) yields further, though smaller, improvements as tasks can be better distributed. The observed diminishing returns at M = 3, where energy/time benefits plateau and edge task count saturates, realistically simulating that factors like task dependencies, communication bottlenecks to/from the device, or the suitability of some tasks for device/cloud limit or decrease the utility of adding more edge resources.
- Device Cores (K, Fig. 3): Increasing local processing power (K=2 to 8) reduces the application makespan ( $T_{\rm total}$ ) due to increased parallelism potential, aligning

- with realistic system behavior expectations for multicore devices. The impact on final energy  $(E_{total})$ : 2T-H sees moderate energy reduction as more tasks can run efficiently locally, reducing costly cloud offloads. For 3T frameworks, energy remains relatively flat or slightly increases; while individual local tasks might be cheaper on more/better cores (modeled via  $P_k$ ), the overall incentive to offload tasks to edge/cloud for significant energy savings (compared to any local execution) persists, especially as more cores potentially increase total device idle power. The steady increase in Final Local Count confirms increased local execution with higher K.
- Number of Tasks (|V|, Fig. 4): Increasing application complexity (|V|=20 to 75) leads to near-linear increases in both final makespan ( $T_{\rm total}$ ) and final energy ( $E_{\rm total}$ ) for the heuristic schedulers. The increasing energy gap between 2T-H and the 3T frameworks highlights the importance of efficient tier selection (including edge) for larger, more complex applications, a realistic scenario. The super-linear growth in optimization duration for 3T-QL highlights the scalability challenge inherent in applying reinforcement learning to larger state spaces, contrasting with the efficiency of the heuristics.

The analysis based on the dynamic resource models, produces results that are consistent with mobile cloud computing system behavior expectations and also highlights real-world system behaviors and trade-offs involving mobile device battery life, network quality, available computing resources (local, edge, cloud), and application scale.

#### VIII. CONCLUSION

This paper presented a comprehensive enhancement to Lin et al.'s[1] mobile cloud computing task scheduling algorithm by addressing three key limitations: static resource modeling, binary offloading architecture, and simplified task characterization. Our dynamic resource models incorporate battery-level sensitivity, workload-dependent power consumption, and variable network conditions, enabling more realistic simulation of mobile environments. By extending the architecture to include edge computing as an intermediate tier between devices and cloud resources, we provided additional offloading options with diverse energy-performance tradeoffs.

Experimental results across different configurations confirm that our three-tier framework consistently outperforms the binary device-cloud architecture, particularly for data-intensive tasks that benefit from edge execution. Furthermore, our Q-learning approach for task migration demonstrated superior energy savings while maintaining completion times within acceptable constraints. The integration of reinforcement learning enables the algorithm to discover non-intuitive optimizations that heuristic approaches might miss, evidenced by oscillatory migration patterns that ultimately converge to more efficient schedules.

#### IX. FUTURE WORK

Several promising research directions emerge from our enhanced three-tier scheduling framework. First, the dynamic power models could be further refined through real-world deployment and measurement studies on diverse mobile devices and network conditions. This would enable more accurate energy estimation and potentially reveal additional optimization opportunities specific to different hardware configurations.

Our Q-learning approach could be extended with more sophisticated reinforcement learning techniques such as deep Q-networks or actor-critic methods to better capture complex relationships between task characteristics, resource states, and optimal offloading decisions. Incorporating transfer learning capabilities would allow the system to leverage knowledge gained from previous applications to accelerate adaptation to new task graphs.

The current framework could be expanded to address multidevice scenarios where multiple mobile devices compete for limited edge and cloud resources. This introduces new dimensions of resource contention, fairness considerations, and potential for collaborative optimization across devices. Game-theoretic approaches might complement our reinforcement learning techniques in such distributed settings.

Incorporating predictive models for network conditions, user mobility, and workload characteristics could enable proactive scheduling decisions rather than purely reactive ones. By anticipating changes in resource availability or network quality, the scheduler could preemptively migrate tasks before conditions deteriorate, further improving both energy efficiency and performance reliability.

# APPENDIX A NOTATION REFERENCE

TABLE III: Comprehensive notation and definitions for the scheduling framework

Symbol	Meaning/Usage	Mathematical Formula/Definition
1. Fundamental C	Concepts	
G = (V, E)	Application represented as a Directed Acyclic Graph	V=task set, $E$ =dependency set
$v, v_i$ pred $(v)$	An individual task in the application graph Set of immediate predecessor tasks whose execu-	$v_i \in V$ $\{v_j \in V \mid (v_j, v) \in E\}$
$\operatorname{succ}(v)$	tion must complete before $v$ 's execution can start Set of immediate successor tasks whose execution requires $v$ 's execution to complete first	$\{v_j \in V \mid (v, v_j) \in E\}$
k	Index for a local core on the mobile device (Device Tier)	$k \in \{0, \dots, K-1\}$
(m,c)	Index for core $c$ on edge node $m$ (Edge Tier) Symbol representing the Cloud Tier resource	$m \in \{1, \dots, M\}, c \in \{1, \dots, C_m\}$
u	Identifier for any potential task execution location (unit)	$u \in \{k, (m, c), C\}$
$S_k$	Ordered sequence of tasks assigned to device core	
$S_{ m cloud}$	Ordered sequence of tasks assigned to the cloud tier	$S_{ ext{cloud}} = \{v_{i_1}, v_{i_2}, \ldots\}$
$S_{(m,c)}$	Ordered sequence of tasks assigned to edge core $(m,c)$	$S_{(m,c)} = \{v_{i_1}, v_{i_2}, \dots\}$
$\frac{S_u}{2}$ Core Invest Pos	General term for the execution sequence on unit <i>u</i>	$S_u \in \{S_k, S_{\text{cloud}}, S_{(m,c)}\}$
	ion Times & Characteristics	
$T_{i,k}^l$	Duration of $v_i$ computation executed locally on device core $k$	$T_{i,k}^l > 0, \ \forall i,k$
$T_{i,m,c}^e$	Duration of $v_i$ computation executed on edge server $m$ , core $c$ (Result of estimation model below)	$T_{i,m,c}^e > 0, \ \forall i,m,c$
$T_i^c$	Duration of $v_i$ computation performed on cloud resources	$T_i^c > 0, \ \forall i$
$\mathtt{complexity}(v_i)$ $\mathtt{intensity}(v_i)$	Measure of task $v_i$ 's computational requirement Measure of task $v_i$ 's data handling requirement	$ exttt{complexity} \in [\gamma_{\min}, \gamma_{\max}] \  exttt{intensity} \in [\delta_{\min}, \delta_{\max}]$
2.2 Communication	n Channels & Data Transfers	
$\overline{\ell}$	Identifier for a specific communication channel type	$\ell \in \{  extsf{d2c},  extsf{c2d},  extsf{d2e}_m,  extsf{e2d}_m, \ldots \}$
d2c	Channel for device uploading data to cloud (uplink)	_
c2d	Channel for device downloading data from cloud (downlink)	_
$d2e_m$	Channel for device uploading data to edge node $m$ (uplink)	
$\mathtt{e2d}_m$	Channel for device downloading data from edge node $m$ (downlink)	$m \in \{1, \dots, M\}$
$r_{ exttt{base}}(\ell)$	Base bandwidth provisioned for channel type $\ell$ (Mbps)	$r_{\mathrm{base}}(\ell) > 0$
$f_{BW}$	Global bandwidth scaling factor simulating overall network conditions	$f_{BW} > 0$
	<u>.</u>	Continued on next na

Symbol	Meaning/Usage	Mathematical Formula/Definition
$r_{ ext{eff}}(\ell) \ d_{i,\ell}^{ ext{send}}$	Effective data transfer rate achieved on channel $\ell$ Data volume (MB) sent for $v_i$ input/initialization over channel $\ell$	$\begin{split} r_{\rm eff}(\ell) &= r_{\rm base}(\ell) \times f_{BW} \\ d_{i,\ell}^{\rm send} &> 0 \end{split}$
$d_{i,\ell}^{\mathtt{receive}}$	Data volume (MB) of $v_i$ results received over channel $\ell$	$d_{i,\ell}^{\mathtt{receive}} > 0$
$T_{ exttt{transfer}}(\ell, d)$	Duration (seconds) required to transfer $d$ MB of data over channel $\ell$	$T_{\mathrm{transfer}}(\ell,d) = d/r_{\mathrm{eff}}(\ell)$
2.3 Mobile Device I	Power & Energy Models	
$\overline{B}$	Current battery level of the mobile device (%)	$B \in [0, 100]$
bf	Battery factor scaling device power based on current battery level $B$	$bf = \begin{cases} 1.0 & \text{if } B > 30\\ 1.0 + 0.01(30 - B) & \text{if } B \le 30 \end{cases}$
$b_k, c_k$	Base idle power and dynamic power coefficient for device core <i>k</i>	$b_k > 0, c_k > 0$
load	Normalized core utilization factor (for device core power calculation)	$\mathtt{load} \in [0,1]$
$P_k$	Instantaneous power consumption of device core $k$ given load and $bf$	$P_k = (b_k + c_k \cdot \mathtt{load}) \cdot bf$
$rf_{ t eff}$	Efficiency factor of the mobile device's radio frequency components	$rf_{ t eff} \in (0,1]$
s	Received signal strength indicator (RSSI) in dBm, influencing RF power	s < 0
r	Data transmission rate (Mbps) used in RF power calculation	$r=r_{\tt eff}(\ell)$ for the relevant channel $\ell$
$P_{ t d2e}$	Mobile device RF power draw during data upload to an edge node at rate $r$	$P_{\text{d2e}} = \frac{bf}{rf_{\text{eff}}} [0.1 + 0.4 \frac{r}{10} (1 + 0.02(70 - s))]$
$P_{ t d2c}$	Mobile device RF power draw during data upload to the cloud at rate $r$	$P_{\text{d2c}} = \frac{bf}{rf_{\text{eff}}} [0.15 + 0.6\frac{r}{5}(1 + 0.03(70 - s))]$
$E_i^{\tt device}$	Energy consumed by mobile device to compute $v_i$ locally on core $k$	$E_i^{\text{device}} = P_k \cdot T_{i,k}^l$
$E_i^{ t RF}$	Energy consumed by mobile device RF component to upload $v_i$ 's data over channel $\ell$	$E_i^{\mathtt{RF}} = P_\ell \cdot T_{\mathtt{transfer}}(\ell, d_{i,\ell}^{\mathtt{send}})$
2.4 Edge & Cloud S	Server Power Models (Non-Mobile Energy)	
$\overline{\texttt{efficiency}_{m,c}}$	Performance factor simulating heterogeneity of edge core $(m, c)$	E.g., $1.0 - 0.05(m-1) - 0.02(c-1)$
$\kappa_{m,c}$	Power scaling factor derived from edge efficiency	$\kappa_{m,c}=1.0/{ t efficiency}_{m,c}$
$P_{\mathtt{edge}}(m,c,\mathtt{load})$	Total power consumption of edge core $(m, c)$ at a	$(5.0+3.0+12.0  imes  exttt{load})  imes \kappa_{m,c}$
$P_{\mathtt{cloud}}(\mathtt{load})$	given load  Total power consumption of cloud resource at a given load	$50.0 + 20.0 + 180.0  \times \texttt{load}$
2.5 Edge Execution	Time Estimation Model (Deriving $T_{i,m,c}^e$ )	
$T_{\min}^{l}$	Minimum local execution time for $v_i$ across all device cores	$T_{\min}^l = \min_k T_{i,k}^l$
$T_i^s, T_i^r$	Simplified cloud send/receive times	Used only in $T_{\text{total}}^c$
$T_{ t total}^c$	Simplified total path time for cloud (input to edge estimate)	$T_{\text{total}}^c = T_i^s + T_i^c + T_i^r$
	Interpolation weight balancing device vs. cloud	$\alpha \in [0,1]$
$\alpha$	influence	
$lpha$ $T_{ ext{base}}^{ ext{edge}}$	influence Base edge computation time estimate derived from $T_{\min}^l, T_{\mathrm{total}}^c$	$T_{ exttt{base}}^{ ext{edge}} = lpha T_{ ext{min}}^l + (1 - lpha) T_{ ext{total}}^c$

Symbol	Meaning/Usage	Mathematical Formula/Definition
adj <sub>data</sub>	Edge time multiplier for data-intensive tasks	$=1.0-0.1(\mathtt{intensity}(v_i)/2)$
adj <sub>type</sub>	Adjustment factor based on $v_i$ 's task type	$\in \{ \text{adj}_{\text{compute}}, \text{adj}_{\text{data}}, 1.0 \}$
$\mathtt{efficiency}_{m,c}$	Performance factor simulating heterogeneity of edge core $(m, c)$	1.0 - 0.05(m - 1) - 0.02(c - 1)
$T_{i,m,c}^{ ext{final}}$		$T_{i,m,c}^{\texttt{final}} = T_{\texttt{base}}^{\texttt{edge}} \cdot \texttt{adj}_{\texttt{type}} \cdot \texttt{efficiency}_{m,c} \cdot \mathcal{U}(0.95, 1.05)$
-,,-	scheduling $v_i$	type the type
3. Phase 1 (Initial S	cheduling) Variables	
3.1 Preliminary Tier	Ranking Estimates (No Contention) (Step 1.1)	
$T_{i, {\tt device}}^{\tt est}$	Estimated time if $v_i$ runs locally (no channel and	$T_{i, \texttt{device}}^{\texttt{est}} = \min_k T_{i, k}^l$
Test	core contention)  Estimated time if a runs on cloud (no channel and	Test _ T (40a Jsend)   TC
$T_{i,\mathtt{cloud}}^{\mathtt{est}}$	Estimated time if $v_i$ runs on cloud (no channel and core contention)	$ \begin{array}{lll} T_{i, {\rm cloud}}^{\rm est} &=& T_{\rm transfer}({\rm d2c}, d_{i, {\rm d2c}}^{\rm send}) + T_i^c + \\ T_{\rm transfer}({\rm c2d}, d_{i, {\rm c2d}}^{\rm receive}) \end{array} $
$T_{i, \mathrm{edge}}^{\mathrm{est}}$	Estimated time if $v_i$ runs on best edge unit (no	$T_{i,\text{edge}}^{\text{est}} = \min_{m,c} (T_{\text{transfer}}(\text{d2e}_m, d_{i,\text{d2e}_m}^{\text{send}}) + T_{i,m,c}^e +$
	channel and core contention)	$T_{\mathtt{transfer}}(\mathtt{e2d}_m, d_{i,\mathtt{e2d}_m}^{\mathtt{receive}}))$
$\frac{\texttt{PreliminaryTier}(v_i)}{}$	Initial tier preference based on minimum $T_{i,\tau}^{\text{est}}$	$rg\min_{ au \in \{ ext{DEVICE,EDGE,CLOUD}\}} T_{i, au}^{ ext{est}}$
3.2 Prioritization Var		
		$\overline{w}_i = \begin{cases} T_{i, \text{cloud}}^{\text{est}} & \text{if } \text{PT}(v_i) = \text{CLOUD} \\ T_{i, \text{edge}}^{\text{est}} & \text{if } \text{PT}(v_i) = \text{EDGE} \\ \frac{1}{K} \sum_k T_{i,k}^l & \text{if } \text{PT}(v_i) = \text{DEVICE} \end{cases}$
$\overline{w}_i$	Cost weight for $v_i$ based on PreliminaryTier $(v_i)$	$\overline{w}_i = \left\{ T_{i, \text{edge}}^{\text{est}}  \text{if } PT(v_i) = \text{EDGE} \right.$
		$\frac{1}{K} \sum_{k} T_{i,k}^{l}$ if $PT(v_i) = DEVICE$
$\mathtt{priority}(v_i)$	Upward rank (critical path estimate) for scheduling	$\mathtt{priority}(v_i) = \overline{w}_i + \max_{v_j \in \mathtt{succ}(v_i)} \mathtt{priority}(v_j)$
-	order	
$\frac{L_{prio}}{}$	List of tasks sorted by $priority(v_i)$ descending	
4. Core Scheduling	Variables (Phase 1.3 & Kernel)	
4.1 Ready Times (Wh	nen a task can start)	
$RT_i^l$		$RT_i^l = \max(0, \max_{v_j \in \text{pred}(v_i)} FT_{\text{final}}(v_j))$
$DT^ws$	pendency met) Earliest device can start sending $v_i$ to cloud (data +	$DT^{ws} = max(DT^l E AT)$
$RT_i^{ws}$	channel ready)	$RI_i = \max(RI_i, FAI_{WS})$
$RT_i^c$		$RT_i^c = \max(FT_i^{ws}, \max_{v_p \in \text{pred}_C(v_i)} FT_p^c)$
ı	done + cloud preds done)	$t = (t - t) \text{ opcprox}(t_i) \text{ p}$
$RT_i^{wr}$	Earliest cloud results for $v_i$ are ready to send back	
$RT_{i, \mathtt{d2e}_m}$	Earliest device can start sending $v_i$ to edge $m$ (data	$RT_{i,d2e_m} = \max(RT_i^l, FAT_{d2e}(m))$
$DT^e$	+ channel ready)	D/T/e
$RT_{i,(m,c)}^e$	Earliest edge core $(m, c)$ can start computing $v_i$ (upload, preds, core ready)	$\begin{array}{ll} RT^e_{i,(m,c)} &= \\ \max(FT^{des}_{i,m}, \max_{v_p \in pred_E(v_i,m)} FT^e_{p,m,c'}, FAT_{edge}(m,c)) \end{array}$
42 C-1 - 1-1 - 1 C4	**	$\frac{\max(\Gamma T_{i,m}, \max_{v_p \in \text{pred}_E(v_i,m)} \Gamma T_{p,m,c'}, \Gamma \Lambda T \text{edge}(m,c))}{\Gamma T_{i,m}}$
	& Finish Times (Result of scheduling decision)	
$EST(v_i, u)$	Earliest possible start time for $v_i$ on unit $u$ considering data readings and resource qualitability	
	ering data readiness and resource availability	RT <sub>for_unit</sub> is:
		$\begin{cases} RT_i^l & \text{if } u = k \text{ (Device Core)} \\ RT_i^{loc} & \text{if } u = k \text{ (Climber 1.5)} \end{cases}$
		$RT_i^{ws}$ if $u = C$ (Cloud Upload Start)
		$\begin{cases} RT_{i,d2e_m} & \text{if } u = (m,c) \text{ (Edge Upload Start)} \end{cases}$
		$RT_i^c$ if $u = C$ (Cloud Compute Start)
		$RT_{i,(m,c)}^e$ if $u=(m,c)$ (Edge Compute Start)
		(etc. for download phases)
$SST(v_i, u)$	Actual scheduled start time of $v_i$ on its assigned	$SST(v_i, u) = EST(v_i, u)$
$ET^l$	unit u	$ET^l = CCT(s, k) + T^l$
$\frac{FT_i^l}{}$	Finish time of $v_i$ if scheduled on a device core $k$	$FT_i^l = SST(v_i, k) + T_{i,k}^l$

Symbol	Meaning/Usage	Mathematical Formula/Definition
$\overline{FT_i^{ws}}$	Finish time of sending $v_i$ 's data from device to	$FT_i^{ws} = SST_{ws} + T_{ exttt{transfer}}( ext{d2c}, d_{i, ext{d2c}}^{ ext{send}})$
E/T/C	cloud	ETC CCT LTC
$FT_i^c \\ FT_i^{wr}$	Finish time of $v_i$ 's computation within the cloud Finish time of receiving $v_i$ 's results from cloud at	$FT_i^c = SST_c + T_i^c$ $FT_i^{wr} = SST_{wr} + T_{renefer}(\text{c2d.} d_i^{\text{receive}})$
1 11	the device	$T_{i} = SSTwr + Transfer(SZG, w_{i,c2d})$
$FT_{i,m}^{\tt des}$	Finish time of sending $v_i$ 's data from device to	$FT_{i,m}^{\text{des}} = SST_{des} + T_{\text{transfer}}(\text{d2e}_m, d_{i,\text{d2e}_m}^{\text{send}})$
	edge node $m$	,
$FT_{i,m,c}^e$	Finish time of $v_i$ 's computation on edge core $(m, c)$	$FT_{i,m,c}^e = SST_e + T_{i,m,c}^e$
$FT_{i,m}^{ ext{er}}$	Finish time of receiving $v_i$ 's results from edge	$FT_{i,m}^{\text{er}} = SST_{er} + T_{\text{transfer}}(\text{e2d}_m, d_{i, \text{e2d}_m}^{\text{receive}})$
	node $m$ at the device	
	Effective finish time: when $v_i$ 's results are available at the device	$FT_i$ If $tier(v_i) = DEVICE$
$FT_{\mathtt{final}}(v_i)$	Effective finish time: when $v_i$ 's results are available	$FT_{\texttt{final}}(v_i) = \left\langle FT_i^{w_i}  \text{if } \texttt{tier}(v_i) = \texttt{CLOUD} \right\rangle$
	at the device	$ig(FT^{er}_{i,m}   ext{if tier}(v_i) =  ext{EDGE}$
4.3 Resource Avail	ability (Updated after each task scheduling)	
$\overline{FAT(u)}$	Finish Available Time: Timestamp when resource u	$FAT(u) = \max(0, \max_{v_j \in S_u} FT_{\text{on\_unit}}(v_j, u)) \text{ where}$
	completes its last assigned task	$FT_{\text{on\_unit}}$ is the relevant finish time on unit $u$ (e.g.,
		$FT_j^l$ if $u = k$ , $FT_j^{ws}$ if $u = ws$ , $FT_{j,m,c}^e$ if $u =$
(1)		(m,c) etc.)
$FAT_{\text{core}}(k)$	Availability time of device core k	$FAT_{\text{core}}(k) = \max(0, \max_{v_j \in S_k} FT_j^l)$
$FAT_{ m ws} \ FAT_{ m wr}$	Availability time of device-to-cloud send channel Availability time of cloud-to-device receive channel	$FAT_{ws} = \max(0, \max_{v_j \in S_{\text{cloud}}} FT_j^{ws})$ $FAT_{ws} = \max(0, \max_{v_j \in S_{\text{cloud}}} FT_{wr}^{wr})$
$FAT_{ m edge}(m,c)$	Availability time of edge core $(m, c)$	$FAI_{\text{wr}} = \max(0, \max_{v_j \in S_{\text{cloud}}} FI_j)$ $FAT_{\text{edge}}(m, c) = \max(0, \max_{v_j \in S_{(m,c)}} FT_{j,m,c}^e)$
$FAT_{d2e}(m)$	Availability time of device-to-edge $m$ send channel	$FAT_{\text{d2e}}(m) = \max(0, \max_{v_i \in S_{\text{new}}} FT_{i m}^{\text{de}})$
$FAT_{e2d}(m)$	Availability time of edge m-to-device receive chan-	$FAT_{\text{e2d}}(m) = \max(0, \max_{v_i \in S_{\text{edge}, m}} FT_{i, m}^{\text{er}})$
	nel	
5. Phase 2 (Energ	y Optimization) Variables	
$\overline{\mathcal{M}}$	Set of candidate migrations $(v_{tar}, u')$ generated	-
$E'_{ total}, T'_{ total}$	Projected Energy/Makespan after a candidate migration	Calculated by rescheduling all tasks using Kernel Rescheduling Algorithm after modifying sequences
$\Delta E$	Mobile energy reduction from a candidate migra-	for the migration $\Delta E = E$
<b>4</b>	tion	$\Delta E = E_{ exttt{total}} - E'_{ exttt{total}}$
$\Delta T$	Makespan change from a candidate migration	$\Delta T = T'_{ exttt{total}} - T_{ exttt{total}}$
Criterion 1	Select migration with max $\Delta E$ , given $\Delta T \leq 0$	$\arg\max_{\mathtt{migration} \in \mathcal{M}_a} \{\Delta E\}$
Criterion 2	Select migration with max $\Delta E/\Delta T$ , given $\Delta T > 0$	$\arg\max_{\mathtt{migration} \in \mathcal{M}_b} \{\Delta E / \Delta T\}$
	and $T'_{\text{total}} \leq T_{\text{max}}$	
6. Kernel Algoritl	hm Internal State	
$\mathtt{ready}_1(v_i)$	Count of unscheduled DAG predecessors of $v_i$	Integer $\geq 0$
$\mathtt{ready}_2(v_i)$	Sequence readiness of $v_i$ (0 if ready, 1 if waiting)	$\in \{0,1\}$
Q	Queue/Stack holding tasks ready for kernel	$Q \leftarrow \{v_i \mid \mathtt{ready}_1(v_i) = 0 \land \mathtt{ready}_2(v_i) = 0\}$
	scheduling	
	mance Metrics & Task State	
$T_{ total}$	Final application makespan (results available at device)	$T_{\text{total}} = \max_{v_i \in \text{exit tasks}} FT_{\text{final}}(v_i)$
$E_{\mathtt{total}}$	Total mobile device energy consumed	$E_{ exttt{total}} = \sum_{v_i \in V} (E_i^{ exttt{device}} + E_i^{ exttt{RF}})$
$T_{\max}$	Application deadline constraint	$T_{ t total} \leq T_{ ext{max}}$
$\mathtt{assignment}(v_i)$	Assigned execution unit index for task $v_i$	$\in \{0K - 1, K, K + 1K + M \cdot C_m - 1\}$
$tier(v_i)$	Assigned execution tier for task $v_i$	€ {DEVICE, EDGE, CLOUD}
$\mathtt{state}(v_i)$	Current scheduling phase status of task $v_i$	€ {UNSCHEDULED, SCHEDULED, KERNEL_SCHEDUI

#### REFERENCES

- [1] Energy and Performance-Aware task scheduling in a mobile cloud computing environment. (2014, June 1). IEEE Conference Publication IEEE Xplore. http://ieeexplore.ieee.org/document/6973741
- [2] Tang, M., & Wong, V. W. S. (2020, April 10). Deep reinforcement learning for task offloading in mobile edge computing systems. arXiv.org. https://arxiv.org/abs/2005.02459
- [3] Cong, P., Zhou, J., Li, L., Cao, K., Wei, T., & Li, K. (2020). A survey of Hierarchical Energy Optimization for Mobile Edge Computing. ACM Computing Surveys, 53(2), 1–44. https://doi.org/10.1145/3378935
- [4] Ali, F. A., Simoens, P., Verbelen, T., Demeester, P., & Dhoedt, B. (2015). Mobile device power models for energy efficient dynamic offloading at runtime. Journal of Systems and Software, 113, 173–187. https://doi.org/10.1016/j.jss.2015.11.042
- [5] S. Guo, B. Xiao, Y. Yang, and Y. Yang, "Energy-efficient dynamic offloading and resource scheduling in mobile cloud computing," in Proc. IEEE INFOCOM, 2016, pp. 1-9.
- [6] H. Wu, W. Knottenbelt, and K. Wolter, "Analysis of the energy-response time tradeoff for mobile cloud offloading using combined metrics," in Proc. 27th Int. Teletraffic Congress (ITC), 2015, pp. 134-142.
- [7] A. Carroll and G. Heiser, "An analysis of power consumption in a smartphone," in Proc. USENIX Annual Technical Conference, 2010, pp. 1-14.
- [8] J. Huang, F. Qian, A. Gerber, Z. M. Mao, S. Sen, and O. Spatscheck, "A close examination of performance and power characteristics of 4G LTE networks," in Proc. 10th Int. Conf. Mobile Systems, Applications, and Services, 2012, pp. 225-238.
- [9] M. Satyanarayanan, P. Bahl, R. Caceres, and N. Davies, "The case for VM-based cloudlets in mobile computing," IEEE Pervasive Comput., vol. 8, no. 4, pp. 14-23, Oct.-Dec. 2009.
- [10] W. Shi, J. Cao, Q. Zhang, Y. Li, and L. Xu, "Edge computing: Vision and challenges," IEEE Internet Things J., vol. 3, no. 5, pp. 637-646, Oct. 2016.
- [11] L. Huang, X. Feng, A. Feng, Y. Huang, and L. P. Qian, "Distributed deep learning-based offloading for mobile edge computing networks," Mobile Netw. Appl., vol. 23, pp. 1-8, Nov. 2018.
- [12] K. Wang, M. Shen, J. Cho, A. Banerjee, J. Van der Merwe, and K. Webb, "MobiScud: A fast mobile cloud gaming platform with dynamic resource scheduling," in Proc. 7th Int. Conf. Multimedia Systems, 2016, pp. 1-13.
- [13] H. Flores, P. Hui, S. Tarkoma, Y. Li, S. Srirama, and R. Buyya, "Mobile code offloading: From concept to practice and beyond," IEEE Commun. Mag., vol. 53, no. 3, pp. 80-88, Mar. 2015.