

Batch-MOT: Batch-Enabled Real-Time Scheduling for Multiobject Tracking Tasks

Donghwa Kang, Seunghoon Lee, Cheol-Ho Hong¹, *Member, IEEE*, Jinkyu Lee², *Senior Member, IEEE*, and Hyeongboo Baek³, *Member, IEEE*

Abstract—Targeting a multiobject tracking (MOT) system with multiple MOT tasks, this article develops Batch-MOT, the first system design that achieves both (G1) timing guarantee and (G2) accuracy maximization, by utilizing *batch execution* that allows multiple deep neural network (DNN) executions to perform simultaneously in a single DNN inference resulting in significantly decreased execution time without accuracy loss. To this end, we propose an adaptable scheduling framework that allows run-time execution behaviors deviated from our base scheduling algorithm (i.e., nonpreemptive fixed-priority scheduling) without compromising G1. Based on the adaptable framework, we then develop 1) a run-time batching mechanism that finds and executes a batch set of MOT tasks and 2) a run-time idling mechanism that waits for the future releases of MOT tasks for batch execution. Both run-time mechanisms can achieve G1 and G2 without incurring high run-time overhead, as they systematically exploit the run-time execution behaviors allowed by the adaptive framework. Our evaluation conducted with a real-world data set demonstrates the effectiveness of Batch-MOT in improving tracking accuracy while providing a timing guarantee compared to the state-of-the-art real-time MOT system for multiple MOT tasks.

Index Terms—Batch execution, multiobject tracking (MOT), real-time scheduling, timing guarantee.

I. INTRODUCTION

AS MODERN autonomous vehicles (AVs) are equipped with multiple cameras, they require performing multiple multiobject tracking (MOT) tasks under limited computing resources. Perception tasks, such as MOT are required to

complete their execution before specified deadlines because AVs' safety-related functions for path planning and vehicle control heavily rely on the timely perception, e.g., determining the time-to-collision with pedestrians and cars ahead as extensively discussed in the previous studies [1], [2], [3], [4]. Additionally, it is widely acknowledged that low accuracy in the perception tasks also compromises the safety of AVs [2], [3]. Therefore, an MOT system with multiple MOT tasks must achieve two goals simultaneously for every MOT task: 1) (G1) timing guarantee and 2) (G2) accuracy maximization.

The deep neural network (DNN)-based MOT approaches are increasingly deployed in modern AVs [5], [6], [7], [8], [9]. However, it is challenging to fully utilize them in order to achieve G1 and G2 for an MOT system with multiple MOT tasks due to the inherent tradeoff between G1 and G2. A recent study developed a scheduling framework that provides different execution options by efficiently utilizing the control knob of processing either a full-size or a down-scaled input image, which is the only existing study that addresses both G1 and G2 for multiple MOT tasks [3]. However, this control knob has its tradeoff; ensuring G1 might compromise G2 by necessitating downscaled image processing.

To overcome the tradeoff between G1 and G2, we utilize *batch execution* for multiple MOT tasks, which, supported by modern DNN models, concurrently processes multiple inputs in one DNN inference reducing execution time without accuracy loss by optimizing GPU resources [4], [10], [11]. As shown in Fig. 1(a) of the experiment results for a GPU of Tesla V100 (comparable to the NVIDIA Orin system-on-chip (SoC) providing similar GPU capability for Tensor core operations [12]) with the state-of-the-art DNN model (i.e., YOLOX [13]), 1.0 time unit taken for processing 12 *full-size* (size of 672×672) images *individually* (one by one) decreases to 0.46 time unit with *batch execution* without accuracy loss (maintaining 41%). Notably, this processing time is even smaller than 0.74 time unit taken for processing 12 *down-scaled* (size of 256×256) images *individually* (one by one) with accuracy drop to 17.7%. This is because it maintains nearly the same DNN inference time (see ii. in Fig. 1(b) for 19–21 ms) until the GPU reaches resource saturation, which occurs when processing more than ten input images for batch execution. On the other hand, the execution times of the other parts (to be detailed in Section II) linearly increase with the number of input images in a batch. We conducted the same experiments on a GPU-enabled embedded board (i.e., NVIDIA

Manuscript received 4 August 2024; accepted 4 August 2024. This work was supported in part by the National Research Foundation of Korea (NRF) Grant funded by the Korea Government (MSIT) under Grant 2022R1A4A3018824 and Grant RS-2024-00438248; in part by the National Research and Development Program through the National Research Foundation of Korea (NRF) funded by Ministry of Science and ICT under Grant 2021M3H2A1038042. The work of Hyeongboo Baek supported by the 2024 Research Fund of the University of Seoul. This article was presented at the International Conference on Embedded Software (EMSOFT) 2024 and the ESWEK-TCAD special issue. This article was recommended by Associate Editor S. Dailey. (*Corresponding authors: Jinkyu Lee; Hyeongboo Baek.*)

Donghwa Kang is with the Department of Computer Science and Engineering, Incheon National University, Incheon 22012, South Korea.

Seunghoon Lee and Jinkyu Lee are with the Department of Computer Science and Engineering, Sungkyunkwan University, Jongno 03063, South Korea (e-mail: jinkyu.yi.3@gmail.com).

Cheol-Ho Hong is with the Department of Intelligent Semiconductor Engineering, Chung-Ang University, Dongjak 06974, South Korea.

Hyeongboo Baek is with the Department of Artificial Intelligence, University of Seoul, Dongdaemun 02504, South Korea (e-mail: hbbaek@uos.ac.kr).

This article has supplementary downloadable material available at <https://doi.org/10.1109/TCAD.2024.3443002>, provided by the authors.

Digital Object Identifier 10.1109/TCAD.2024.3443002

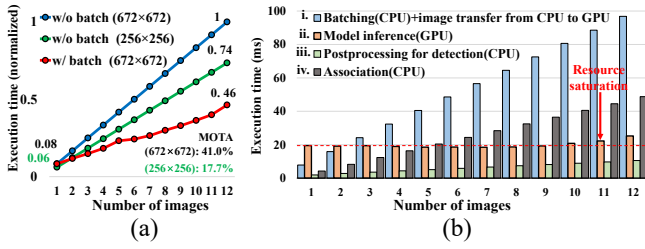


Fig. 1. Execution times for different number of input images of an MOT system with YOLOX on a Tesla V100 GPU. (a) Total execution. (b) Decomposition of batch execution.

76 Xavier SoC [14]) and observed that it can accommodate two
 77 input images (see Fig. 6(d) for 36–38 ms) for a batch execution
 78 before the resource saturation.

79 Motivated by the shorter execution time by batch execution
 80 without accuracy loss, we target a set of MOT tasks, in which
 81 G1 is compromised by processing the original DNN workloads
 82 (i.e., full-size image) *individually*¹ but is not compromised by
 83 processing the reduced DNN workloads (i.e., the down-scaled
 84 image that compromises G2) *individually*. Then, we consider
 85 the individual execution with the reduced DNN workloads in
 86 default and aim at performing batch execution with the original
 87 DNN workloads from multiple MOT tasks as frequently as
 88 possible at run-time, such that both G1 and G2 are achieved;
 89 this entails the following challenges.

90 C1: To determine the batch set to execute, we require an
 91 online mechanism to find feasible batch sets, along with
 92 run-time information on active MOT tasks, which results
 93 in significant run-time overhead; thus, a new scheduling
 94 framework with the *low run-time overhead* is essential.

95 C2: While batch execution can expedite overall processing,
 96 it may delay specific tasks due to factors like priority
 97 inversion (to be detailed in Fig. 4 in Section V), neces-
 98 sitating a runtime mechanism (based on the answer to
 99 C1) with a schedulability test to ensure the timely task
 100 execution.

101 C3: Since, any work-conserving scheduling cannot yield any
 102 batch execution under a situation where there is only one
 103 active task, we need a run-time idling mechanism (based
 104 on the answer to C1) that accelerates the batch execution
 105 for the situation while ensuring the timely execution of
 106 every task (based on the answer to C2).

107 In this article, we propose Batch-MOT, the first system
 108 design to achieve G1 and G2 by utilizing *batch execution* for
 109 multiple MOT tasks, which systematically tackles the chaining
 110 challenges C1–C3. Batch-MOT employs nonpreemptive fixed-
 111 priority scheduling (NPFP) as a base scheduling algorithm,
 112 in which each MOT task is executed nonpreemptively and the
 113 task priority is predefined.

114 To address C1, we analyse the offline schedulability test
 115 that provides timing guarantees under NPFP. Based on the
 116 analysis, we propose a new scheduling framework NP_{ADAPT}
 117 (the nonpreemptive ADAPTable scheduling framework) that

¹If not compromised, the computing resource is sufficient for achieving G1 and G2 without any advanced technique, which is not the scope of this article.

118 allows run-time execution behaviors to deviate from NPFP
 119 without compromising timing guarantees achieved by the
 120 offline schedulability test.

121 The strategy of NP_{ADAPT} to reduce runtime overhead
 122 involves the *offline* identification of the amount of allowable
 123 run-time execution behavior deviations (denoted by Δ_k defined
 124 in Section IV-B) from NPFP for each MOT task, providing
 125 an interface for developing a *runtime* batching mechanism that
 126 incurs low runtime overhead without compromising timing
 127 guarantees.

128 As to C2, we develop a run-time batching mechanism
 129 NPFP^B (NPFP with batch execution) based on NP_{ADAPT}.
 130 It finds a set of MOT tasks with a low run-time overhead,
 131 such that executing the set as a batch does not compro-
 132 mise the timely execution of any task. This is achieved by
 133 systematically exploiting the properties to be discussed in
 134 Section III and the run-time execution behaviors allowed by
 135 NP_{ADAPT}.

136 To address C3, we propose an advanced run-time batching
 137 mechanism with an idling scheme NPFP^{BI} (NPFP^B with
 138 idling), developed on top of NPFP^B. NPFP^{BI} further utilizes
 139 the run-time execution behaviors allowed by NP_{ADAPT} and
 140 determines the idling interval for each MOT task to wait for
 141 the future release(s) of the other MOT tasks to be executed as
 142 a batch, without incurring much run-time overhead. The rela-
 143 tionship among NP_{ADAPT}, NPFP^B, and NPFP^{BI} is described
 144 in Figure S.1 in the supplement [15].

145 We implemented Batch-MOT and evaluated it using an
 146 open MOT data set of the autonomous driving system.
 147 Our evaluation demonstrates that Batch-MOT exhibits higher
 148 tracking accuracy and lower run-time overhead without
 149 compromising timing guarantee, compared to the only exist-
 150 ing study addressing both G1 and G2 for multiple MOT
 151 tasks [3].

152 We clarify our novelty and contribution along with an
 153 explanation of related work as follows.

- 154 1) To the best of our knowledge, Batch-MOT is the first
 155 study providing a *strict* timing guarantee for batch DNN
 156 execution of multiple (camera) tasks. Even extending
 157 our interest to general DNN beyond MOT, the existing
 158 studies address different problems from ours. That is,
 159 [2], [3], [11], [16], [17] do not deal with a timing guar-
 160 tee for batch DNN execution; [4], [18] are designed for
 161 single-camera (i.e., single-task) systems; and [10], [19]
 162 aim at improving the overall FPS or minimizing the
 163 deadline miss ratio, therefore not addressing strict timing
 164 guarantees.
- 165 2) Batch-MOT enables the assurance of timing guaran-
 166 tees through a simple online test with low scheduling
 167 cost despite the complicated impact of batch DNN
 168 execution (both with/without idling) on the timing
 169 guarantee. This is achieved by a) the novel design
 170 of the scheduling frameworks NP_{ADAPT}, NPFP^B,
 171 and NPFP^{BI} (to be detailed in Sections IV–VI,
 172 respectively) and b) the mathematical foundation of
 173 their timely correctness, both of which are highly
 174 challenging.

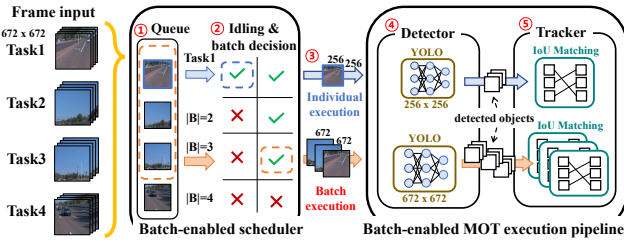


Fig. 2. Overview of Batch-MOT

II. OVERVIEW OF BATCH-MOT

As illustrated in Fig. 2, the core design features of Batch-MOT include the *batch-enabled MOT execution pipeline* and the *batch-enabled scheduler* to be detailed in this section. The MOT execution pipeline and scheduler are implemented as separate threads, and they communicate with each other by exchanging messages through the shared memory. The workflow of Batch-MOT is as follows. Each input frame of the MOT tasks is forwarded to a ready queue (① in Fig. 2). Once the batch-enabled scheduler determines the idling time and batch size of MOT tasks (②), some MOT tasks in the ready queue are combined into a batch or no batch is constructed according to the idling and batch decision (③). Then, detection (④) and association (⑤) are conducted sequentially for either batch or individual execution.

A. Batch-Enabled MOT Execution Pipeline

In the batch-enabled MOT execution pipeline, the *front-end* DNN-based detector identifies the position and size of each object’s bounding box in the input frame and sends the detection information to the *back-end* tracker. The tracker then conducts an association to match each detected object with one of the existing objects in the previous frame (called tracklet) based on the intersection over union-based (IoU-based) matching and updates the tracking information for the matched tracklet.

For the detector, Batch-MOT adopts any existing stand-alone DNN models (e.g., the YOLO series [13], [20]) that can accept variable input image sizes determining the tracking accuracy as shown in Fig. 1(a). The batch MOT execution pipeline supports two types of execution: 1) *individual* and 2) *batch* execution. For an input image of the size of 672×672 , the individual execution scale-downs the input image to the size of 256×256 . Then, it conducts the detection of the MOT tasks and shows decreased execution time at the expense of sacrificing the tracking accuracy. In a batch execution, multiple input images with the original size of 672×672 are combined in a batch for DNN inference, and it is transferred from the CPU to GPU memory [i. in Fig. 1(b)]. Then, the DNN inference is conducted on the GPU to detect candidate objects [ii] in Fig. 1(b), and the postprocessing, such as nonmaximum suppression (NMS) [20] is performed on the CPU to extract high-confident objects among detected candidates [iii] in Fig. 1(b). As the tracker uses IoU-based matching, it compares the position and size of tracklets with the objects detected in the current frame on a one-to-one

basis and matches two objects whose size of overlapping region is greater than a given threshold on the CPU. In the case of batch execution, after detection is performed for multiple MOT tasks, the associations for the batch are then performed sequentially on the CPU [iv] in Fig. 1(b); however, if the time cost for interprocess communication (IPC) on the platform is relatively low compared to the execution time of an association, the associations can be executed in parallel across multiple CPUs using multiprocessing, which decreases the overall execution time.

B. Batch-Enabled Scheduler

Batch-MOT supports a thread-level scheduler invoked when an MOT task is released or completed. The proposed batch-enabled scheduler operates as a background daemon and communicates with the MOT execution pipeline through the shared memory. To make Batch-MOT capable of addressing C1–C3, the batch-enabled scheduler is designed as follows.

To tackle C1–C3, Batch-MOT needs to implement a run-time batch decision mechanism that does not compromise timing guarantees while maintaining the low run-time overhead. This is challenging because batch execution affects the behavior of multiple tasks, including i) the target task; ii) the lower priority tasks; and iii) the higher priority tasks. In other words, the batch execution of given \mathcal{B} can change the execution of i), ii), and/or iii), and the impact of one is different from that of the others. Considering all the possible execution variations of batch executions and their influences on the other tasks at every scheduling decision may cause prohibitively high run-time overhead, which has not been addressed in the previous study [3].

To overcome the challenge, we first analyse the underlying principle of our base scheduling algorithm NPFP and its schedulability analysis that judges the timing guarantee of the given task τ_k by considering i), ii), and iii) to be detailed in Section IV-A. We then develop a new scheduling framework, NP_{ADAPT}, which enables the run-time execution behaviors of i), ii), and iii) deviated from NPFP while preserving the schedulability guaranteed under NPFP, by associating the run-time execution behaviors with the schedulability test to be detailed in Section IV-B. By using the run-time execution behaviors with the properties to be discussed in Section III, the scheduler finds and executes schedulable batch sets under work-conserving scheduling (addressing C1 and C2 to be detailed in Section V) and beyond work-conserving scheduling (addressing C1 and C3, to be detailed in Section VI), with the low run-time overhead.

III. SYSTEM MODEL

We consider an MOT system with multiple DNN-based MOT tasks $\tau = \{\tau_i\}_{i=1}^n$ [3] on a platform equipped with multiple CPUs and a single GPU. As an input video frame is provided periodically, an MOT task τ_i is considered a periodic real-time task with a timing constraint. That is, an MOT task τ_i invokes a series of jobs J_i , each separated by exactly T_i time units; once a job of τ_i is released at t , it should finish its execution no later than its deadline $t + T_i$. The period of

each task is not necessarily the same, which makes it possible to address the situation where the frame rate of each camera varies based on its intended use (e.g., side-facing cameras typically operate at lower rates while forward-facing cameras operate at higher rates [1]); of course, our task model also accommodates a set of tasks with the same period. A job is said to be *active* at t , if it has remaining execution at t . Let $\tau(t)$ denote a set of tasks, each of whose job is active at t . Let $r_i(t)$ denote the earliest job release time of τ_i after t . Since, each task is strictly periodic, it is possible to know $r_i(t)$ at t and indicate the earliest job deadline of τ_i after t . Let $\text{LP}(\tau_k)$ and $\text{HP}(\tau_k)$ denote a set of tasks whose priority is lower and higher than τ_k , respectively. Notations are summarized in Table S.1 in the supplement.

An MOT task set τ is said to be *schedulable* under a target scheduling framework if every job invoked by tasks in τ does not miss its deadline when the framework schedules τ . Since, there is at most one active job of $\tau_i \in \tau$ at any time, we use a task τ_i and a job of τ_i (denoted by J_i) interchangeably when no ambiguity arises. As presented, we aim to achieve G1 and G2 for every MOT task. Let C_i denote the worst-case execution time (WCET) of τ_i when each performs *individual* execution (as opposed to *batch* execution).

To schedule a set of MOT tasks, we decide to apply the following two policies. First, we enforce nonpreemptiveness between the detection and association of each job; that is, once a job of τ_i starts its execution, it sequentially performs the execution of its detection and association subjobs without any preemption. Second, we disallow individual MOT tasks to be executed in parallel (if not a part of a batch), while the associations of MOT tasks in a batch can be performed simultaneously on multiple CPUs depending on the time cost of IPC mentioned in Section II. The two policies not only reduce the run-time scheduling overhead but also significantly lower the complexity of considering various run-time scenarios that incur different interference/blocking; therefore, the policies make it possible to ensure offline timing guarantees through a simple online test with low scheduling cost to be developed in Sections IV–VI.

In this article, we process a down-scaled image (i.e., 256×256) as the *individual* execution of each job, while we do a full-size image (i.e., 672×672) as the *batch* execution of a group of jobs. Let \mathcal{B} denote a set of tasks whose jobs will be executed as a batch, and let $C_{\mathcal{B}}$ denote the WCET of \mathcal{B} , where $|\mathcal{B}| \geq 2$. We take the measurement-based approach to derive the WCET of MOT tasks, using the experiment setup in Section VII, with an in-depth discussion provided in Section VIII.

In this article, we use the following properties of batch execution.

- P1: $C_{\mathcal{B}} \geq \max_{\tau_i \in \mathcal{B}} C_i$;
- P2: $C_{\mathcal{B}} \leq \sum_{\tau_i \in \mathcal{B}} C_i$; and
- P3: $C_{\mathcal{B}} \leq C_{\mathcal{B}'}$, if $\mathcal{B} \subset \mathcal{B}'$.

Batch execution of multiple MOT tasks dramatically shortens total inference latency by reducing the number of GPU invocations. That is, *individual* execution requires as many invocations as the number of given MOT tasks, while *batch* execution performs inference with one invocation. Since, the

inference latency through a single GPU invocation increases monotonically according to the DNN workload of the invocation, P1 holds generally. Likewise, since DNN workload of \mathcal{B} will be increased if we add more job(s) to \mathcal{B} , P3 holds generally. Apart from P1 and P3, which typically holds, P2 holds under the following condition: the benefit of reducing the number of GPU invocations outweighs the increasing workload (from a smaller total DNN workload of multiple *individual* executions to a larger DNN workload of a single *batch* execution). To satisfy the condition, we need to deploy a detector that offers high optimization of batch execution. In this article, we target the state-of-the-art detectors optimized for batch execution that ensure P2 is met (e.g., YOLOX [13] illustrated in Fig. 1, YOLOv5 [20], Faster-RCNN [21], and others with varying sizes for both the downscaled and full-size input images).

IV. DEVELOPING ADAPTABLE SCHEDULING FRAMEWORK

A. Base Scheduling Algorithm NPFP

In this article, we employ NPFP [3], [22] as a base scheduling algorithm. As explained in Section III, once a job of τ_i starts its execution, it sequentially performs the execution of its detection and association subjobs without any preemption (by NP). Also, each task's priority is predefined, and each job inherits its invoking task's priority (by FP). Whenever there is at least one active job while the computing system is idle, NPFP selects the highest-priority active job and starts its execution. NPFP is work conserving, meaning that the computing system cannot be idle if there is at least one active job; also, the vanilla NPFP does not allow any batch execution.

The schedulability of a set of MOT tasks under NPFP is guaranteed by the next lemma [3], [23], which is a sufficient but not necessary schedulability test.

Lemma 1 (From [3], [23]): An MOT task set τ is schedulable by NPFP, if $R_k \leq T_k$ holds for every $\tau_k \in \tau$, where R_k (i.e., the response time of τ_k) can be calculated as follows. $R_k(x+1)$ is calculated by (1) sequentially with $x = 0, 1, 2, \dots$, starting from $R_k(0) = C_k + \sum_{\tau_h \in \text{HP}(\tau_k)} C_h + \max_{\tau_j \in \text{LP}(\tau_k)} C_j$, until $R_k(x+1) = R_k(x)$ (implying $R_k = R_k(x)$) or $R_k(x+1) > T_k$ (implying no bounded R_k)

$$R_k(x+1) = C_k + \sum_{\tau_h \in \text{HP}(\tau_k)} \left\lceil \frac{R_k(x)}{T_h} \right\rceil \cdot C_h + \max_{\tau_j \in \text{LP}(\tau_k)} C_j. \quad (1)$$

Proof: Here, we summarize the proof in [3] and [23]. Consider the following situation.

- 1) The first job of τ_k and every first job of tasks whose priority is higher than τ_k are released at t_0 .
- 2) A job of a task whose C_i is the largest among tasks whose priority is lower than τ_k is released right before t_0 .
- 3) The following jobs of τ_k and those of tasks whose priority is higher than τ_k are released periodically.

It was proven that one of the jobs of τ_k released under the situation (but not necessarily the first job of τ_k released at t_0) yields the largest response time of τ_k [22].

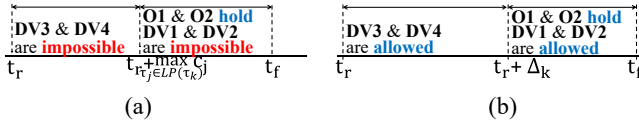


Fig. 3. Properties of NPFP and NP_{ADAPT}. (a) NPFP (b) NP_{ADAPT}.

386 Then, it is trivial that the response time of the first job of τ_k
 387 under the situation is upper-bounded by $R_k(x)$ that satisfies (1).
 388 The proof of [3, Lemma 1] proves that the response time of
 389 the $(x+1)^{th}$ job of τ_k cannot be larger than that of the x^{th} job
 390 (where $x \geq 1$), if $R_k \leq T_k$ holds with (1), meaning there is no
 391 self-pushing phenomenon issue [22] for τ_k if $R_k \leq T_k$ holds
 392 with (1). ■

393 B. Adaptable Scheduling Framework NP_{ADAPT}

394 Once Lemma 1 deems τ schedulable, we can guaran-
 395 tee timely execution of τ scheduled by NPFP. However,
 396 the accuracy of MOT tasks in τ scheduled by NPFP
 397 cannot be maximized, because every job under NPFP
 398 performs individual execution with a down-scaled image (as
 399 opposed to batch execution of multiple jobs with a full-size
 400 image). Therefore, targeting τ deemed schedulable under the
 401 vanilla NPFP (that does not employ batch execution), we
 402 want to maximize the MOT accuracy by performing batch
 403 execution as much as possible without compromising the
 404 schedulability.

405 However, a batch execution of a set of multiple jobs easily
 406 compromises each job's timely execution achieved by the
 407 individual execution of corresponding jobs. For example, if
 408 a higher-priority job of τ_i and a lower-priority job of τ_j are
 409 executed as a batch, the execution time of the batch could be
 410 larger than that of the job of τ_i solely (by P1), which may
 411 yield the deadline miss of the job of τ_i . Therefore, we need to
 412 identify which run-time execution behaviors (caused by batch
 413 execution) that deviate from NPFP do not compromise the
 414 schedulability.

415 To establish boundaries of run-time execution behavior devi-
 416 ation (e.g., execution time increment due to batch execution,
 417 intentional idling for batch execution) without compromising
 418 the schedulability, we target a task set that satisfies Lemma 1
 419 and analyse how Lemma 1 guarantees the schedulability under
 420 NPFP. To this end, we focus on a job of τ_k that is released
 421 and finished at t_r and t_f , respectively. First, by the property
 422 of the nonpreemptiveness of NPFP, a lower-priority job can
 423 block the execution of a higher-priority job *only when* the
 424 former starts its execution before the release of the latter. As
 425 addressed by the third term of the RHS of (1), the following
 426 holds in $[t_r + \max_{\tau_j \in LP(\tau_k)} C_j, t_f)$ under NPFP, illustrated in
 427 Fig. 3(a).

428 O1: Except the execution of the job of τ_k and jobs with
 429 higher priority than τ_k , any other run-time behav-
 430 ior (e.g., other jobs' execution, the system idling) is
 431 disallowed.

432 Second, after the lower-priority blocking, the only possible
 433 execution behaviors that affect the schedulability of the job
 434 of τ_k are the execution of the job of τ_k itself and jobs with

435 higher priority than τ_k . As addressed by the first two terms of
 436 the RHS of (1), the following holds in $[t_r + \max_{\tau_j \in LP(\tau_k)} C_j, t_f)$
 437 under NPFP if Lemma 1 holds. This is illustrated in
 438 Fig. 3(a).

439 O2: The amount of execution of the job of τ_k and jobs
 440 with higher priority than τ_k does not exceed $C_k +$
 441 $\sum_{\tau_h \in HP(\tau_k)} \lceil ([t_f - t_r]/T_h) \rceil \cdot C_h$.

442 Considering the two properties, we define a class of
 443 the nonpreemptive scheduling algorithms with the minimum
 444 requirements, which 1) allows run-time execution behavior
 445 deviated from NPFP to be potentially utilized for batch
 446 execution and 2) does not compromise the schedulability under
 447 NPFP guaranteed by Lemma 1.

448 *Definition 1:* We define NP_{ADAPT} associated with given
 449 $\{\Delta_k \geq 0\}_{\tau_k \in \tau}$ (the nonpreemptive ADAPTable scheduling
 450 framework), as any nonpreemptive scheduling algorithm in
 451 which every job of $\tau_k \in \tau$ (that is released and finished at t_r
 452 and t_f , respectively) satisfies the following features, which are
 453 illustrated in Fig. 3(b).

454 F1: In $[t_r + \Delta_k, t_f)$, O1 holds.

455 F2: In $[t_r + \Delta_k, t_f)$, O2 holds.

456 Since, F1 and F2 are the only requirements, NP_{ADAPT} with
 457 given $\{\Delta_k\}$ can accommodate the following possible run-time
 458 execution behaviors deviated from NPFP *as long as F1 and*
 459 *F2 hold*. Recall that C_i for each τ_i is the WCET when its job
 460 performs individual execution with a down-scaled image (as
 461 opposed to batch execution of multiple jobs with a full-size
 462 image).

- 463 1) In $[t_r + \Delta_k, t_f)$, (DV1) the job of τ_k may execute for
 464 more than C_k .
- 465 2) In $[t_r + \Delta_k, t_f)$, (DV2) a job of $\tau_h \in HP(\tau_k)$ may execute
 466 for more than C_h .
- 467 3) In $[t_r, t_r + \Delta_k)$, (DV3) a job of $\tau_j \in LP(\tau_k)$ executes for
 468 more than C_j .
- 469 4) In $[t_r, t_r + \Delta_k)$, (DV4) the computing system becomes
 470 idle even though there is an active job, meaning that
 471 NP_{ADAPT} is not work conserving.

472 The above run-time execution behaviors will be used
 473 for the run-time batching mechanism to be explained in
 474 Sections V and VI. We would like to emphasize that
 475 DV1–DV4 are run-time execution behaviors deviated from
 476 NPFP, as NPFP does not allow them in the corresponding
 477 intervals; in other words, NPFP disallows DV1 and DV2
 478 in $[t_r + \max_{\tau_j \in LP(\tau_k)} C_j, t_f)$, and DV3 and DV4 in $[t_r, t_r +$
 479 $\max_{\tau_j \in LP(\tau_k)} C_j)$.

480 Considering that NPFP satisfies O1 and O2 in $[t_r +$
 481 $\max_{\tau_j \in LP(\tau_k)} C_j, t_f)$ if Lemma 1 holds, the following prop-
 482 erty holds: for τ that is deemed schedulable by Lemma 1,
 483 NPFP belongs to NP_{ADAPT} associated with $\{\Delta_k =$
 484 $\max_{\tau_j \in LP(\tau_k)} C_j\}$. From F1 and F2 for NP_{ADAPT}, we can easily
 485 derive the schedulability analysis of NP_{ADAPT}, by replacing
 486 $\max_{\tau_j \in LP(\tau_k)} C_j$ with Δ_k in Lemma 1.

487 *Theorem 1:* An MOT task set τ is schedulable by NP_{ADAPT}
 488 associated with given $\{\Delta_k\}$, if $R_k \leq T_k$ holds for every $\tau_k \in \tau$,
 489 where R_k (i.e., the response time of τ_k) can be calculated as
 490 follows. $R_k(x+1)$ is calculated by (2) sequentially with $x =$
 491 $0, 1, 2, \dots$, starting from $R_k(0) = C_k + \sum_{\tau_h \in HP(\tau_k)} C_h + \Delta_k$,
 492 until $R_k(x+1) = R_k(x)$ (implying $R_k = R_k(x)$) or $R_k(x+1) >$

493 T_k (implying no bounded R_k)

$$494 \quad R_k(x+1) = C_k + \sum_{\tau_h \in \text{HP}(\tau_k)} \left\lceil \frac{R_k(x)}{T_h} \right\rceil \cdot C_h + \Delta_k. \quad (2)$$

495 *Proof:* Suppose that, a job of τ_k whose release time is
 496 t_r does not finish its execution until $t_r + R_k$, although R_k
 497 satisfies (2). First, we consider that the amount of execution
 498 of τ_k and its higher-priority tasks in $[t_r + \Delta_k, t_r + R_k]$ is not
 499 larger than $C_k + \sum_{\tau_h \in \text{HP}(\tau_k)} \lceil (R_k/T_h) \rceil \cdot C_h$. From (2), $R_k -$
 500 $\Delta_k = C_k + \sum_{\tau_h \in \text{HP}(\tau_k)} \lceil (R_k/T_h) \rceil \cdot C_h$ holds. Therefore, if the
 501 amount of execution of τ_k and its higher-priority tasks in $[t_r +$
 502 $\Delta_k, t_r + R_k]$ is not larger than $C_k + \sum_{\tau_h \in \text{HP}(\tau_k)} \lceil (R_k/T_h) \rceil \cdot C_h$,
 503 the supposition contradicts F1. Second, we consider that the
 504 amount of execution of τ_k and its higher-priority tasks in $[t_r +$
 505 $\Delta_k, t_r + R_k]$ is larger than $C_k + \sum_{\tau_h \in \text{HP}(\tau_k)} \lceil (R_k/T_h) \rceil \cdot C_h$, which
 506 immediately contradicts F2. Therefore, the supposition always
 507 contradicts. ■

508 Provided that $\Delta_k \geq \max_{\tau_j \in \text{LP}(\tau_k)} C_j$, no execution of τ_j
 509 can occur within the interval $[t_r + \Delta_k, t_f)$, where t_r and t_f
 510 represent the release and finishing times of τ_k , respectively.
 511 This condition ensures the sustainability property with respect
 512 to $\{C_j\}$.

513 In the next section, we will develop a run-time batching
 514 mechanism by utilizing the capability of NP_{ADAPT} in achiev-
 515 ing the schedulability even in the presence of the run-time
 516 execution behavior deviated from NPFP (as long as F1 and
 517 F2 are satisfied). To this end, we will deploy the largest Δ_k
 518 for NP_{ADAPT} to accommodate a longer blocking period due
 519 to batch execution or a longer idling for batch execution to
 520 be performed in the future. Let Δ_k^* denote the largest Δ_k that
 521 does not compromise the schedulability of τ_k in Theorem 1,
 522 and let R_k^* denote R_k with Δ_k^* . We can easily verify that if τ is
 523 deemed schedulable by Lemma 1, $\Delta_k^* \geq \max_{\tau_j \in \text{LP}(\tau_k)} C_j$ holds
 524 for every $\tau_k \in \tau$.

525 *Offline Time-Complexity:* We can find $\Delta_k^* \in [0, T_k - C_k]$
 526 using the binary search. Hence, the time complexity to find
 527 Δ_k^* and R_k^* for every $\tau_k \in \tau$ using Theorem 1 is $O(n^2 \cdot \log(n) \cdot$
 528 $\max(T_k))$, which is affordable as it is performed offline.

529 V. NPFP^B: ENABLING RUN-TIME BATCHING

530 As C1 and C2 in Section I indicate, utilizing batch execution
 531 necessitates a mechanism that efficiently finds a batch of jobs
 532 to be executed *at run-time* without compromising timing guar-
 533 antee. Therefore, this section develops a run-time mechanism
 534 that achieves the following goals to address C1 and C2 in
 535 Section I, respectively.

- 536 1) Perform batch execution as frequently as possible
 537 (for high accuracy) while minimizing the run-time
 538 complexity.
- 539 2) Do not compromise the schedulability of τ under
 540 NPFP, guaranteed by Lemma 1.

541 To this end, we develop NPFP^B (NPFP with batch execu-
 542 tion) associated with given $\{\Delta_k\}$. We address the second goal
 543 by making NPFP^B follow F1 and F2 (implying NPFP^B with
 544 $\{\Delta_k\}$ belongs to NP_{ADAPT} with $\{\Delta_k\}$). Then, the remaining
 545 step is how to design a run-time batching mechanism that
 546 addresses the first goal while satisfying F1 and F2.

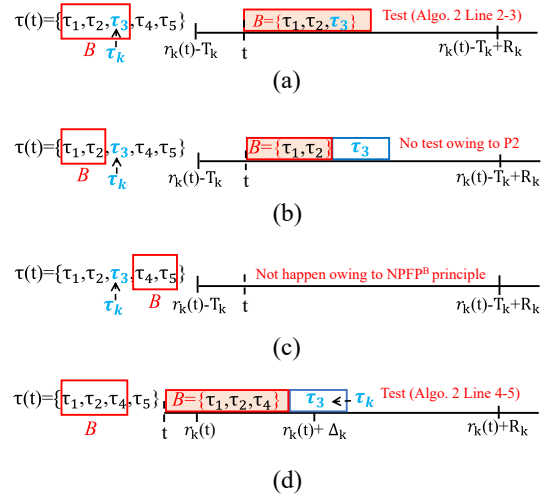


Fig. 4. Four batch execution cases for a set of tasks $\{\tau_1, \tau_2, \tau_3(=\tau_k), \tau_4, \tau_5\}$; the smaller index, the higher priority. (a) Case 1: $\tau_k \in \tau(t)$ and $\tau_k \in \mathcal{B}$. (b) Case 2: $\tau_k \in \tau(t)$, $\tau_k \notin \mathcal{B}$, and $\tau_h \in \text{HP}(\tau_k)$. (c) Case 3: $\tau_k \in \tau(t)$, $\tau_k \notin \mathcal{B}$, and $\tau_h \in \text{LP}(\tau_k)$. (d) Case 4: $\tau_k \notin \tau(t)$ (and therefore $\tau_k \notin \mathcal{B}$).

547 As a first step to develop NPFP^B associated with given
 548 $\{\Delta_k\}$, we investigate how each batch execution affects the
 549 schedulability under NPFP guaranteed by Lemma 1. From
 550 now on, we interpret a run-time execution behavior deviated
 551 from NPFP due to batch execution as a change of WCET
 552 (as well as the actual execution time) of the highest-priority
 553 task in the batch; therefore, the priority of a batch execution
 554 inherits the priority of the highest-priority task in the batch.
 555 For example, if a higher-priority job of τ_i and a lower-priority
 556 job of τ_j are executed as a batch, we regard this situation as
 557 an increase of WCET of the job of τ_i from C_i to $C_{\mathcal{B}}$ where
 558 $\mathcal{B} = \{\tau_i, \tau_j\}$.

559 Consider τ deemed schedulable by Lemma 1. Suppose
 560 we schedule τ by NPFP until t , but we are going to
 561 execute a set of jobs as a batch (denoted by \mathcal{B}) at t . We
 562 investigate how the schedulability of a job J_k of the task
 563 $\tau_k \in \tau$ is affected differently according to the following four
 564 cases, which are illustrated in Fig. 4 with a task set $\tau =$
 565 $\{\tau_1, \tau_2, \tau_3(=\tau_k), \tau_4, \tau_5\}$, in which a smaller task index implies
 566 a higher priority. Let τ_h denote the highest-priority task among
 567 tasks in \mathcal{B} ; recall that $\tau(t)$ is a set of tasks, each of whose job
 568 is active at t , and $r_k(t)$ is the earliest job release time of τ_k
 569 after t .

- 570 1) *Case 1 of $\tau_k \in \tau(t)$ and $\tau_k \in \mathcal{B}$:* The WCET of J_k is
 571 changed from C_k to $C_{\mathcal{B}}$, e.g., $\mathcal{B} = \{\tau_1, \tau_2, \tau_3(=\tau_k)\}$ in
 572 Fig. 4(a).
- 573 2) *Case 2 of $\tau_k \in \tau(t)$, $\tau_k \notin \mathcal{B}$, and $\tau_h \in \text{HP}(\tau_k)$:* The
 574 longest time for a job of τ_h to delay the execution of J_k is
 575 changed from C_h to $C_{\mathcal{B}}$, e.g., $\mathcal{B} = \{\tau_1, \tau_2\}$ in Fig. 4(b).
- 576 3) *Case 3 of $\tau_k \in \tau(t)$, $\tau_k \notin \mathcal{B}$, and $\tau_h \in \text{LP}(\tau_k)$:* J_k
 577 experiences an additional delay from a lower-priority job
 578 of τ_h for up to $C_{\mathcal{B}}$, which does not occur under NPFP,
 579 e.g., $\mathcal{B} = \{\tau_4, \tau_5\}$ in Fig. 4(c).
- 580 4) *Case 4 of $\tau_k \notin \tau(t)$ (therefore $\tau_k \notin \mathcal{B}$):* The longest time
 581 for a job of τ_h to delay the execution of J_k (to be released
 582 at $r_k(t) > t$) is changed from $\max(0, t + C_h - r_k(t))$ to
 583 $\max(0, t + C_{\mathcal{B}} - r_k(t))$, e.g., $\mathcal{B} = \{\tau_1, \tau_2, \tau_4\}$ in Fig. 4(d).

584 To make $\text{NPFPP}^{\mathcal{B}}$ associated with given $\{\Delta_k\}$ preserve
585 schedulability even in the presence of batch execution for
586 Cases 1–4, our basic design principle for the run-time mech-
587 anism of $\text{NPFPP}^{\mathcal{B}}$ is as follows.

588 We enforce the prioritization policy of FP on batch
589 execution. To this end, we disallow the batch execution of
590 \mathcal{B} at t , if there is an active job of $\tau_k \notin \mathcal{B}$ whose priority is
591 higher than the lowest-priority task in \mathcal{B} . The principle has
592 three distinct advantages as follows.

593 DP1: To make the schedule under $\text{NPFPP}^{\mathcal{B}}$ as similar as
594 possible to that under the base scheduling algorithm
595 NPFPP ,

596 DP2: To eliminate Case 3, which a) not only reduces a
597 burden to check the schedulability affected by given
598 batch execution b) but also helps to comply with F1
599 by preventing a lower-priority job from executing in
600 the interval of interest of F1 (i.e., $[t_r + \Delta_k, t_f]$).

601 DP3: To narrow down the number of possible choices of
602 the set of tasks to be executed as a batch, i.e., from
603 $2^{|\tau(t)|} - 1$ to $|\tau(t)| - 1$, which significantly reduces the
604 run-time overhead for checking different batch candi-
605 dates as well as the offline measurement/analysis
606 overhead for obtaining $C_{\mathcal{B}}$ for different \mathcal{B} .

607 Under the design principle, we handle Cases 1–4 for a given
608 batch execution. For each case, we either derive a condition
609 for the batch execution not to compromise the schedulability
610 (for Cases 1 and 4) or verify that the batch execution cannot
611 compromise the schedulability (for Cases 2 and 3), both of
612 which are achieved by satisfying F1 and F2.

- 613 1) *Case 1 of $\tau_k \in \tau(t)$ and $\tau_k \in \mathcal{B}$* : If $t + C_{\mathcal{B}} \leq r_k(t) - T_k + R_k$
614 holds, F2 is satisfied and the job of τ_k active at t does
615 not miss its deadline,² as exemplified in Fig. 4(a).
- 616 2) *Case 2 of $\tau_k \in \tau(t)$, $\tau_k \notin \mathcal{B}$, and $\tau_h \in \text{HP}(\tau_k)$* : The
617 longest time for jobs of tasks in \mathcal{B} (at most one job per
618 task) to be executed is decreased from $\sum_{\tau_h \in \mathcal{B}} C_h$ (in the
619 $\sum_{\tau_h \in \text{HP}(\tau_k)} \lceil (R_k(x)/T_h) \rceil \cdot C_h$ term in the RHS of (1)) to
620 $C_{\mathcal{B}}$ by P2, which cannot compromise the schedulability
621 of J_k that is active at t , as exemplified in Fig. 4(b).
- 622 3) *Case 3 of $\tau_k \in \tau(t)$, $\tau_k \notin \mathcal{B}$, and $\tau_h \in \text{LP}(\tau_k)$* : This
623 case does not occur under the design principle DP2, as
624 illustrated in Fig. 4(c).
- 625 4) *Case 4 of $\tau_k \notin \tau(t)$ (and Therefore $\tau_k \notin \mathcal{B}$)*: If $t +$
626 $C_{\mathcal{B}} \leq r_k(t) + \Delta_k$ holds, the batch execution of \mathcal{B} does
627 not violate F1 and does not affect F2, as exemplified in
628 Fig. 4(d).

629 The formal proof of the conditions/statements in Cases 1–4
630 will be in the proof of Theorem 2.

631 Using Cases 1–4, Algorithm 2 presents $\text{NPFPP}^{\mathcal{B}}$ associated
632 with $\{\Delta_k\}$, which is performed at t at which there is at least one
633 active job while the computing system is ready to work. After
634 defining $\mathcal{B}(n)$ as a set of the n highest-priority tasks among
635 $\tau(t)$ in line 1, we check whether there are multiple active jobs
636 (i.e., two or more tasks in $\tau(t)$) in line 2. We find the largest x ,
637 such that $\mathcal{B}(x)$ does not compromise the schedulability of all
638 the other jobs, by testing schedulability test for online batching

²Recall that R_k is the response time of τ_k calculated by Theorem 1 for given Δ_k .

Algorithm 1 $\text{STOB}(t, \tau(t), \mathcal{B})$

```

1: for  $\tau_k \in \tau$  do
2:   if  $\tau_k \in \tau(t)$  and  $\tau_k \in \mathcal{B}$  then
3:     if  $t + C_{\mathcal{B}} > r_k(t) - T_k + R_k$  then return unschedulable
4:   else if  $\tau_k \notin \tau(t)$  then
5:     if  $t + C_{\mathcal{B}} > r_k(t) + \Delta_k$  then return unschedulable
6:   end if
7: end for
8: return schedulable

```

Algorithm 2 $\text{NPFPP}^{\mathcal{B}}$ Scheduling Algorithm

At t , at which a job is finished while there is at least one active job,
or at which at least one job is released while the system is idle,

```

1: Let  $\mathcal{B}(n)$  denote  $\{\tau_i(t)\}_{i=1}^n$ , where  $\tau_n(t)$  denotes the  $n^{\text{th}}$   
highest-priority task in the set of active tasks at  $t$  (i.e.,  $\tau(t)$ ) for  
 $1 \leq n \leq |\tau(t)|$ .
2: if  $|\tau(t)| \geq 2$  then
3:   Find the largest batch set  $\mathcal{B}(x)$  for  $2 \leq x \leq |\tau(t)|$  such that  
 $\text{STOB}(t, \tau(t), \mathcal{B}(x))$  in Algorithm 1 returns schedulable,  
using binary search; if such  $\mathcal{B}(x)$  exists, execute a set of  
active jobs invoked by tasks in  $\mathcal{B}(x)$  as a batch, and return.
4: end if
5: Execute the highest-priority active job, and return.

```

(STOB) ($t, \tau(t), \mathcal{B}(x)$) in Algorithm 1 (in Line 3); we will
639 detail Algorithm 1, including why it is possible to apply the
640 binary search. If such $\mathcal{B}(x)$ exists, execute a set of active jobs
641 invoked by tasks in $\mathcal{B}(x)$ as a batch (in line 3). Otherwise (or
642 there is only one active job at t), execute the highest-priority
643 job (in line 5), which is the same as NPFPP .
644

645 For a given \mathcal{B} , STOB in Algorithm 1 checks whether every
646 task τ_k satisfies the conditions in Cases 1 and 4. Lines 2 and
647 3 correspond Case 1, while lines 4 and 5 correspond Case 4.
648 Note that, by the statements of Cases 2 and 3, we do not need
649 to check the cases for schedulability. Now, we present why it
650 is possible to apply the binary search to find the largest $\mathcal{B}(x)$
651 in line 3 of Algorithm 2.

652 *Lemma 2*: Recall $\mathcal{B}(x)$ in line 1 of Algorithm 2. If
653 $\text{STOB}(t, \tau(t), \mathcal{B}(x+1))$ in Algorithm 1 returns schedulable,
654 then $\text{STOB}(t, \tau(t), \mathcal{B}(x))$ also returns schedulable.

655 *Proof*: Since $\mathcal{B}(x) \subset \mathcal{B}(x+1)$ holds, $C_{\mathcal{B}(x)} \leq C_{\mathcal{B}(x+1)}$ holds
656 by P3 in Section III. This implies that the opposite conditions
657 in lines 3 and 5 of Algorithm 1, respectively, satisfy $r_k(t) -$
658 $T_k + R_k \geq t + C_{\mathcal{B}(x+1)} \geq t + C_{\mathcal{B}(x)}$ and $r_k(t) + \Delta_k \geq t +$
659 $C_{\mathcal{B}(x+1)} \geq t + C_{\mathcal{B}(x)}$. Therefore, the lemma holds. ■

660 As we designed, NP_{ADAPT} subsumes $\text{NPFPP}^{\mathcal{B}}$ as follows.

661 *Theorem 2*: For τ with $\{\Delta_k \geq \max_{\tau_j \in \text{LP}(\tau_k)} C_j\}$ that is
662 deemed schedulable by Theorem 1, $\text{NPFPP}^{\mathcal{B}}$ associated with
663 $\{\Delta_k\}$ belongs to NP_{ADAPT} associated with $\{\Delta_k\}$.³

664 *Proof*: Suppose that a job of τ_k (denoted by J_k) scheduled
665 by $\text{NPFPP}^{\mathcal{B}}$ associated with $\{\Delta_k\}$ is released and finished at t_r
666 and t_f , respectively. We prove that J_k always satisfies F1 and
667 F2 of Definition 1, which proves the theorem.

668 First, we check whether F1 is satisfied. Since, $\text{NPFPP}^{\mathcal{B}}$ is
669 work conserving, it suffices to check (F1') whether there is no

³Since Theorem 1 with $\Delta_k = \max_{\tau_j \in \text{LP}(\tau_k)} C_j$ is equivalent to Lemma 1, τ
deemed schedulable by Theorem 1 with $\Delta_k \geq \max_{\tau_j \in \text{LP}(\tau_k)} C_j$ is also deemed
schedulable by Lemma 1 (i.e., NPFPP -schedulable).

670 execution of lower-priority jobs in $[t_r + \Delta_k, t_f)$. We consider
671 four cases.

672 (*Case F1a*): If a lower-priority job of τ_j starts its execution
673 at $t (< t_r)$ as individual execution, it will finish its execution
674 no later than $t + C_j (< t_r + C_j)$. Therefore, as long as $\Delta_k \geq$
675 $\max_{\tau_j \in \text{LP}(\tau_k)} C_j$ holds, F1' holds.

676 (*Case F1b*): If a batch (denoted by \mathcal{B}), including a lower-
677 priority job of τ_j starts its execution at $t (< t_r)$, it will finish
678 its execution no later than $t + C_{\mathcal{B}} (< t_r + C_{\mathcal{B}})$. By line 5 of
679 Algorithm 1, $t + C_{\mathcal{B}} \leq t_r + \Delta_k$ holds, meaning that F1' holds.

680 (*Case F1c*): If a lower-priority job of τ_j starts its execution
681 at $t (\geq t_r)$ although J_k is not finished until t , it violates line 5
682 of Algorithm 2.

683 (*Case F1d*): If a batch (denoted by \mathcal{B}), including a lower-
684 priority job of τ_j starts its execution at $t (\geq t_r)$, we consider
685 two subcases: (i) \mathcal{B} 's priority is lower than τ_k , and ii)
686 otherwise. Recall that \mathcal{B} has the priority of the highest-priority
687 task in \mathcal{B} ; so, the entire execution of \mathcal{B} has equal or higher
688 priority than τ_k . Therefore, i) violates line 5 of Algorithm 2,
689 and ii) does not violate F1 since it is regarded as an execution
690 whose priority is not lower than τ_k .

691 Second, we check whether F2 is satisfied with four cases.

692 (*Case F2a*): We consider there is no batch execution in
693 $[t_r, t_f)$. Similar to the proof of Lemma 1 for NPF \mathcal{B} , the
694 amount of execution of τ_k and its higher-priority tasks in $[t_r +$
695 $\Delta_k, t_f)$ is maximized when all the jobs of τ_k and its higher-
696 priority tasks are released at t_r . In this worst-case situation,
697 F2 trivially holds.

698 (*Case F2b*): We consider a batch (denoted by \mathcal{B}) starts its
699 execution at $t (< t_r)$ (and therefore \mathcal{B} cannot include τ_k). Then,
700 the batch execution will finish its execution no later than $t + C_{\mathcal{B}}$
701 ($< t_r + C_{\mathcal{B}}$). By line 5 of Algorithm 1, $t + C_{\mathcal{B}} \leq t_r + \Delta_k$
702 holds, meaning that the batch execution cannot contribute to
703 higher-priority execution in $[t_r + \Delta_k, t_f)$.

704 (*Case F2c*): We consider a batch (denoted by \mathcal{B}) that does
705 not include τ_k starts its execution at $t (\geq t_r)$. Then, the
706 WCET of \mathcal{B} is no larger than the sum of the corresponding
707 individual WCET (i.e., $C_{\mathcal{B}} \leq \sum_{\tau_h \in \mathcal{B}} C_h$) by P2. This is
708 equivalent to reducing the execution time of some tasks
709 $\{\tau_h \in \mathcal{B}\}$, such that $C_{\mathcal{B}} = \sum_{\tau_h \in \mathcal{B}} C'_h$, where C'_h denotes the
710 reduced execution time of τ_h . Therefore, the batch execution
711 does not compromise the bounded higher-priority execution
712 of Case F2a. Note that, the existence of τ_j belonging to both
713 \mathcal{B} and $\text{LP}(\tau_k)$ may compromise the bounded higher-priority
714 execution due to additional interference contribution by a
715 lower-priority task τ_j ; however, NPF \mathcal{B} disallows to execute
716 a batch that belongs to such τ_j .

717 (*Case F2d*): We consider a batch (denoted by \mathcal{B}) that
718 includes τ_k starts its execution at $t (\geq t_r)$. Different from
719 Case F2c, it is possible for a task whose priority is lower
720 than τ_k to be included in \mathcal{B} due to $\tau_k \in \mathcal{B}$. Recall that \mathcal{B} has
721 the priority of the highest-priority task in \mathcal{B} ; so, the entire
722 execution of \mathcal{B} has equal or higher priority than τ_k . Since
723 NPF \mathcal{B} is work-conserving nonpreemptive scheduling, the
724 execution of \mathcal{B} finishes no later than $t + C_{\mathcal{B}}$, which is no later
725 than $r_k(t) - T_k + R_k = t_r + R_k$ by line 3 of Algorithm 1. We
726 consider two cases: 1) $t_f = t_r + R_k$ and 2) $t_f < t_r + R_k$. In the

727 first case, if we apply F1, the amount of execution of the job
728 of τ_k and jobs of its higher-priority task (by either individual
729 execution or batch execution \mathcal{B}) in $[t_r + \Delta_k, t_f)$ is upper-
730 bounded by $R_k - \Delta_k$. Therefore, violation of F2 contradicts (2)
731 in Theorem 1. In the second case, the amount should be strictly
732 less than $C_k + \sum_{\tau_h \in \text{HP}(\tau_k)} \lceil (t_f - t_r) / T_h \rceil \cdot C_h$; otherwise, (2)
733 should hold for a value (denoted by $R'_k = t_f - t_r$) that is smaller
734 than R_k .

735 One may wonder whether the proof for Case F2d correctly
736 considers multiple jobs of $\tau_h \in \mathcal{B} \setminus \{\tau_k\}$ released after t . Since
737 \mathcal{B} starts at t and finishes at t_f , the job of a higher-priority
738 task τ_h released in (t, t_f) will start at t_f or later, which does
739 not belong to $[t_r + \Delta_k, t_f)$, the interval of interest of F2.
740 Instead, the schedulability of the job of a higher-priority task
741 τ_h released in (t, t_f) will be checked by Algorithm 1 when
742 $\tau_k = \tau_h$; if deemed unschedulable, the corresponding \mathcal{B} cannot
743 be scheduled. Therefore, the proof is correct. ■

744 As shown in line 5 of Algorithm 1, a larger Δ_k implies
745 a higher chance for the algorithm to allow the execution of
746 given \mathcal{B} . Therefore, we will use the largest $\{\Delta_k^*\}$ associated
747 with Theorem 1; we already explained how to calculate $\{\Delta_k^*\}$
748 in Section IV-B. Finally, we present the schedulability analysis
749 of NPF \mathcal{B} in the following theorem.

750 *Theorem 3*: τ is schedulable by NPF \mathcal{B} associated with
751 $\{\Delta_k^*\}$ (i.e., the largest $\{\Delta_k\}$ that makes τ schedulable by
752 Theorem 1), if $\Delta_k^* \geq \max_{\tau_j \in \text{LP}(\tau_k)} C_j$ holds for every $\tau_k \in \tau$.

753 *Proof*: The theorem holds by Theorems 1 and 2. ■

754 *Run-Time Complexity*: At each t , which Algorithm 2
755 focuses on, we test the $O(\log(|\tau(t)|))$ batch sets by STOB,
756 each of which requires $O(|\tau|)$ time-complexity. Therefore, the
757 total run-time complexity is $O(|\tau| \cdot \log(|\tau(t)|))$, which is much
758 lower than $O(|\tau|^2 \cdot |\tau(t)|)$, the complexity of the existing study
759 for scheduling multiple MOT tasks in [3].

760 VI. NPF $\mathcal{B}^{\mathcal{I}}$: EXPLOITING IDLING FOR BATCHING

761 Although NPF \mathcal{B} efficiently finds and executes a set of
762 active jobs as a batch, it inherently cannot address the situation
763 where there is only one active job at t . Since, the situation
764 cannot be addressed by *any work-conserving scheduling*, we
765 need to develop a run-time idling mechanism that accelerates
766 batch execution to address C1 and C3 in Section I. To this
767 end, we develop NPF $\mathcal{B}^{\mathcal{I}}$ (NPF \mathcal{B} with idling) with given
768 $\{\Delta_k\}$. Based on NPF \mathcal{B} , NPF $\mathcal{B}^{\mathcal{I}}$ achieves the same goals:
769 1) maximizing the batch execution with minimum run-time
770 overhead, while 2) preserving the schedulability guaranteed by
771 Lemma 1.

772 Our design principles of the run-time idling mechanism of
773 NPF $\mathcal{B}^{\mathcal{I}}$ are as follows.

- 774 1) To prevent idling from compromising F2, any idling
775 cannot be overlapped with any $[t_r + \Delta_i, t_f)$ for any job
776 of τ_i released and finished at t_r and t_f , respectively.
- 777 2) Between the execution of a batch set at t and that of
778 the same batch set at $t' (> t)$, the latter cannot help
779 any job's timely execution. Therefore, we restrict time
780 instant candidates at which a batch set starts its execution
781 after idling, to the time instants at which any job (to be
782 executed as a batch) is released.

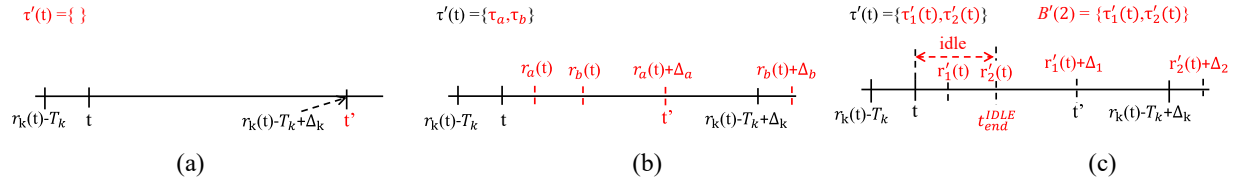


Fig. 5. Example scenario of the idling mechanism of NPFPP^{BI} in Algorithm 3. (a) Initialization (line 1). (b) Finding candidate tasks (lines 2–5). (c) Finding the largest schedulable batch set (lines 6 and 7).

Algorithm 3 Idling Mechanism of NPFPP^{BI}

- 1: Set $t' \leftarrow r_k(t) - T_k + \Delta_k$, and $\tau'(t) \leftarrow \emptyset$
 - 2: **for** $\tau_i \in \tau \setminus \{\tau_k\}$ sorted by $r_i(t)$ **do**
 - 3: **if** $r_i(t) \leq t'$ **then** set $\tau'(t) \leftarrow \tau'(t) \cup \{\tau_i\}$, and
 $t' \leftarrow \min(t', r_i(t) + \Delta_i)$
 - 4: **else** Exit the loop.
 - 5: **end for**
 - 6: Let $\tau'_n(t)$ denote the task with the n^{th} earliest next job release time after t among tasks in $\tau'(t)$; $r'_n(t)$ denote the next job release time of $\tau'_n(t)$ after t ; and $\mathcal{B}'(n)$ denote $\{\tau'_i(t)\}_{i=1}^n$, where $1 \leq n \leq |\tau'(t)|$.
 - 7: Find the largest batch set $\mathcal{B}'(x) \cup \{\tau_k\}$ for $1 \leq x \leq |\tau'(t)|$ such that $\text{STOB}(r'_x(t), \mathcal{B}'(x) \cup \{\tau_k\}, \mathcal{B}'(x) \cup \{\tau_k\})$ in Algorithm 1 returns schedulable, using the *binary search*; if such $\mathcal{B}'(x)$ exists, set $t_{\text{end}}^{\text{IDLE}} \leftarrow r'_x(t)$ where $t_{\text{end}}^{\text{IDLE}}$ denotes the latest idling time instant, and **return**.
 - 8: Execute the active job of τ_k at t , and **return**.
-

783 3) Once we determine to perform batch execution at t after
 784 idling, we include all the active jobs to the batch set to
 785 be executed, which eases the satisfaction of F1 and F2
 786 in the presence of idling.

787 Algorithm 3 presents the run-time idling mechanism of
 788 NPFPP^{BI} at t , at which there is only one active job of τ_k
 789 while the computing system is ready to work. Lines 1–5 find
 790 $\tau'(t) \subset \tau$, a set of candidate tasks to be executed with τ_k as
 791 a batch after idling. We aim at calculating $t'(> t)$, which is,
 792 the latest time instant at which all of the next released jobs in
 793 $\tau'(t)$ and the active job of τ_k can idle without violating F1. t' is
 794 determined by the earlier time instant between $r_k(t) - T_k + \Delta_k$
 795 and the earliest one among $r_i(t) + \Delta_i$ for every $\tau_i \in \tau'(t)$.
 796 Then, $r_k(t) - T_k + \Delta_k \leq t'$ and $r_i(t) + \Delta_i \leq t'$ hold for τ_k
 797 and every $\tau_i \in \tau'(t)$, respectively, which helps to achieve the
 798 first design principle by disallowing any job in the batch set
 799 to start its execution after the interval of interest of F2 for
 800 the job. In line 6, we define $\tau'_n(t)$ as the task with the n^{th}
 801 earliest next job release time after t among tasks in $\tau'(t)$,
 802 $r'_n(t)$ as the next job release time of $\tau'_n(t)$ after t , and $\mathcal{B}'(n)$
 803 as $\{\tau'_i(t)\}_{i=1}^n$.⁴ In line 7, we find the largest x such that the
 804 execution of a batch set of $\mathcal{B}'(x) \cup \{\tau_k\}$ does not compromise
 805 the schedulability of all the other jobs. To this end, we test
 806 $\text{STOB}(r'_x(t), \mathcal{B}'(x) \cup \{\tau_k\}, \mathcal{B}'(x) \cup \{\tau_k\})$ in Algorithm 1, meaning
 807 that we check whether a batch set of $\mathcal{B}'(x) \cup \{\tau_k\}$ can start its
 808 batch execution at $r'_x(t)$ at which jobs of $\mathcal{B}'(x) \cup \{\tau_k\}$ are the
 809 only active jobs; we will explain why it is possible to apply
 810 the binary search. If such $\mathcal{B}'(x) \cup \{\tau_k\}$ exists, we reserve that

⁴For tasks with the same next job release time, a higher priority implies an earlier next job release time.

a set of jobs invoked by tasks in $\mathcal{B}'(x) \cup \{\tau_k\}$ will be executed 811
 at $r'_x(t)$ by setting $t_{\text{end}}^{\text{IDLE}} \leftarrow r'_x(t)$. Otherwise, we execute the 812
 single active job at t immediately (in line 8), which is the same as 813
 as NPFPP . 814

Fig. 5 presents an example scenario at the current time 815
 instant t , at which an idling decision can be made with τ_k under 816
 Algorithm 3. As an initialization, t' and $\tau'(t)$ are set to $r_k(t) - 817$
 $T_k + \Delta_k$ and \emptyset , respectively, [line 1 of Algorithm 3, illustrated 818
 in Fig. 5(a)]. Then, two tasks τ_a and τ_b will be released at $t < 819$
 $r_a(t)$ and $t < r_b(t)$, and t' and $\tau'(t)$ are updated to $r_a(t) + \Delta_a$ 820
 and $\tau'(t) = \{\tau_a, \tau_b\}$, respectively, [lines 2–5 of Algorithm 3, 821
 illustrated in Fig. 5(b)]. Finally, $\mathcal{B}'(2) \cup \{\tau_k\}$ is determined 822
 as a schedulable batch set according to $\text{STOB}(r'_2(t), \mathcal{B}'(2) \cup 823$
 $\{\tau_k\}, \mathcal{B}'(2) \cup \{\tau_k\})$, and then $t_{\text{end}}^{\text{IDLE}}$ is set to $r'_2(t)$ [lines 6–7 of 824
 Algorithm 3 illustrated in Fig. 5(c)]. 825

We present why it is possible to apply binary search in line 7 826
 of Algorithm 3. 827

Lemma 3: Recall $r'_n(t)$ and $\mathcal{B}'(n)$ defined in line 6 of 828
 Algorithm 3. If $\text{STOB}(r'_{n+1}(t), \mathcal{B}'(n+1) \cup \{\tau_k\}, \mathcal{B}'(n+1) \cup 829$
 $\{\tau_k\})$ in Algorithm 1 returns schedulable, $\text{STOB}(r'_n(t), \mathcal{B}'(n) \cup 830$
 $\{\tau_k\}, \mathcal{B}'(n) \cup \{\tau_k\})$ also returns schedulable. 831

Proof: If we apply $C_{\mathcal{B}(n)} \leq C_{\mathcal{B}(n+1)}$ (from $\mathcal{B}(n) \subset \mathcal{B}(n+1)$ 832
 and P3 in Section III) and $r'_n(t) \leq r'_{n+1}(t)$ to the opposite 833
 conditions in lines 3 and 5 of Algorithm 1, the proof is similar 834
 to that of Lemma 2. ■ 835

Including the idling mechanism in Algorithm 3, we present 836
 the entire NPFPP^{BI} scheduling algorithm in Algorithm 4; note 837
 that, $t_{\text{end}}^{\text{IDLE}}$ is set to $-\infty$ when the system starts. In the case of 838
 $t < t_{\text{end}}^{\text{IDLE}}$ (lines 1 and 2), there should not be any execution, 839
 since the idling mechanism determines that all active jobs 840
 will be executed at $t_{\text{end}}^{\text{IDLE}}$, not at the current time instant t . 841
 In the case of $t = t_{\text{end}}^{\text{IDLE}}$ (lines 3 and 4), we start to execute 842
 all the active jobs as batch execution immediately, which is 843
 determined by the idling mechanism. In the case of $t > t_{\text{end}}^{\text{IDLE}}$ 844
 (lines 5 and 6), we consider two cases. First, if there is only 845
 one active job when the computing system is ready to work, 846
 we perform Algorithm 3. Second, if there is more than one 847
 active job when the computing system is ready to work, we 848
 perform Algorithm 2. 849

Then, we prove that NP_{ADAPT} subsumes NPFPP^{BI} as 850
 follows. 851

Theorem 4: For τ with $\{\Delta_k \geq \max_{\tau_j \in \text{LP}(\tau_k)} C_j\}$ that is 852
 deemed schedulable by Theorem 1, NPFPP^{BI} associated with 853
 $\{\Delta_k\}$ belongs to NP_{ADAPT} associated with $\{\Delta_k\}$. 854

Proof: Suppose that, a job of τ_k scheduled by NPFPP^{BI} 855
 associated with $\{\Delta_k\}$ (denoted by J_k) is released and finished 856
 at t_r and t_f , respectively. We prove that J_k always satisfies F1 857
 and F2 of Definition 1, which proves the theorem. 858

Algorithm 4 NPFP^{BT} Scheduling Algorithm

At t , at which a job is finished while there is at least one active job, or at which at least one job is released while the system is idle,

- 1: **if** $t < t_{end}^{IDLE}$ **then**
- 2: Any job cannot start its execution.
- 3: **else if** $t = t_{end}^{IDLE}$ **then**
- 4: Execute all the active jobs at t as a batch.
- 5: **else if** $t > t_{end}^{IDLE}$ **then**
- 6: **if** there is only one active job at t **then** Perform Algorithm 3.
- 7: **else** Perform Algorithm 2.
- 8: **end if**

859 First, we check whether F1 is satisfied. Since, we add the
860 run-time idling mechanism to NPFP^B, we focus on proving
861 that the idling mechanism does not compromise F1. By the
862 selection of $\tau'(t)$ (in lines 1–5 of Algorithm 3) and the
863 selection of the time instant at which a batch execution starts
864 after idling (in line 7 of Algorithm 3), any batch execution
865 after idling cannot start its execution in $[t_r + \Delta_k, t_f]$. Also, a
866 batch execution after idling can be performed only if line 5
867 of Algorithm 1 guarantees that the execution finishes before
868 $t_r + \Delta_k$. This proves the satisfaction of F1.

869 Second, we check whether F2 is satisfied. Since, the idling
870 mechanism complies with F1, we can check F2 when a
871 batch execution starts. This is achieved by calling Algorithm 1
872 by line 7 of Algorithm 3, which corresponds to line 3 of
873 Algorithm 2 for NPFP^B. Therefore, the remaining proof is
874 similar to that for NPFP^B in Theorem 2. ■

875 Finally, we present the schedulability analysis of NPFP^{BT}
876 in the following theorem.

877 *Theorem 5:* τ is schedulable by NPFP^{BT} associated with
878 $\{\Delta_k^*\}$ (i.e., the largest $\{\Delta_k\}$ that makes τ schedulable by
879 Theorem 1), if $\Delta_k^* \geq \max_{\tau_k \in \text{LP}(\tau_k)} C_j$ holds for every $\tau_k \in \tau$.

880 *Proof:* The theorem holds by Theorems 1 and 4. ■

881 *Run Time-Complexity:* At each t , which Algorithm 3
882 focuses on lines 1–5 check at most $|\tau|$ tasks, and lines 6–8
883 test the $O(\log(|\tau(t)|))$ batch sets by STOB, each of which
884 requires $O(|\tau|)$ time complexity; hence, the total run-time
885 complexity of Algorithm 3 is $O(|\tau| \cdot \log(|\tau(t)|))$. By the run-
886 time complexity of Algorithms 2 and 3, that of Algorithm 4 is
887 also $O(|\tau| \cdot \log(|\tau(t)|))$, which is much lower than $O(|\tau|^2 \cdot |\tau(t)|)$
888 for the existing study [3].

VII. EVALUATION

A. Experiment Setup

891 We consider four different computing systems. The first one
892 is equipped with Intel Xeon Silver 4215R CPUs @ 3.20 GHz,
893 251.5 GB RAM, and a Tesla V100 GPU. We also consider
894 three GPU-enabled embedded boards: 1) NVIDIA Jetson TX2;
895 2) Xavier; and 3) Orin. The MOT execution pipeline and
896 scheduler run on Python and Pytorch, and the model precision
897 is set to FP16; on the same experiment setting in Fig. 1, we
898 observed nearly the same tracking accuracy with FP32 owing
899 to the mixed precision training. As the object detector, we
900 consider the YOLO series [13], [20] trained with the COCO
901 Dataset [24]. We use SORT [6] as the object tracker of the
902 two-stage methods. The performance was evaluated using the

Waymo Open Dataset [25], the autonomous driving data set
collected by autonomous driving cars. For the evaluation, we
use the measured WCETs in Fig. 8.

B. Experiment Results

906 In this section, we demonstrate the effectiveness of the
907 proposed run-time batching and idling mechanisms in improv-
908 ing the tracking accuracy by comparing the following three
909 approaches that operate on the architecture of Batch-MOT.
910 Note that, we apply the rate monotonic (RM) [26] to FP, for
911 all the approaches in this section. 912

- 913 1) NPFP in Section IV-A, in which all the MOT tasks
914 perform individual execution with down-scaled images. 914
- 915 2) NPFP^B in Section V, in which down-scaled and full-
916 size images are processed for the individual and batch
917 execution, respectively. 917
- 918 3) NPFP^{BT} in Section VI, in which down-scaled and full-
919 size images are processed for the individual and batch
920 execution, respectively. 920

921 We also compare our approaches with the only existing
922 study that addresses G1 and G2 for multiple MOT tasks as
923 follows. 923

- 924 1) RT-MOT, a flexible MOT execution scheduling frame-
925 work on the architecture proposed in [3] with the YOLO
926 series and DeepSORT [5] as its detector and tracker. 926

927 For all the target approaches, we provide a run-time option
928 called the individual full-size execution policy (IFP). Under
929 IFP, each target approach performs a full-size execution at t
930 only if: i) an MOT task τ_i is active alone at t ; ii) the full-size
931 execution (even with its WCET) of the only active task at t
932 can be completed by t_{next} (the earliest future release of any
933 task later than t), i.e., $t + C_i^{full} \leq t_{next}$, where C_i^{full} is the
934 WCET for the full-size execution of the task; and iii) if the
935 target approach is NPFP^{BT}, the system is not idle at t under
936 NPFP^{BT}, i.e., $t < t_{end}^{IDLE}$ in Algorithm 3. Conditions i) and ii)
937 ensure IFP improves accuracy through the full-size execution
938 without compromising timing guarantees of the other tasks.
939 Condition iii) ensures IFP can be incorporated into NPFP^{BT}
940 without conflicting with its idling mechanism. Note that, RT-
941 MOT inherently employs IFP. 941

942 While the existing DNN-based MOD techniques
943 (e.g., [2], [10], and [11]) could be considered for comparison,
944 adapting them for the MOT systems would require new
945 contributions. For example, extending DNN-SAM [2] would
946 need alignment in tracking algorithm, ROI identification, and
947 other features of RT-MOT for fairness. Although DNN-SAM
948 is optimized for accuracy within the ROI, it requires two
949 separate DNN inferences: one for the ROI and another for
950 outside areas, hindering high overall accuracy. Executing two
951 DNN inferences not only doubles the computational tasks of
952 pre/postprocessing [as shown in i) and iii) of Fig. 1], but also
953 limits the benefits of increased GPU utilization from batch
954 execution. 954

955 Since, the approaches share the same offline schedulability
956 test in Lemma 1, we compare their tracking accuracy of the
957 task sets whose schedulability is guaranteed by the test. We
958 use multiple object tracking accuracy (MOTA) [27], a primary
959 metric to evaluate the tracking accuracy; tracking accuracy 959

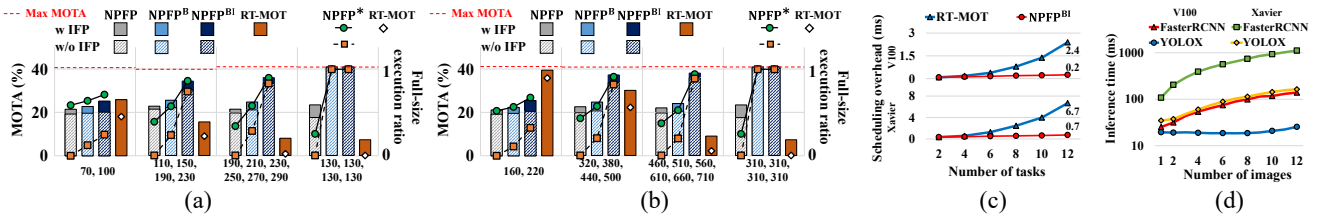


Fig. 6. Comparison of different approaches on YOLOX (a), (b), and (c), and different DNN models (d).

under MOTA is derived by counting miss detection, false detection, and miss tracking; we obtained similar experimental results using IDF-1 [27], another widely recognized metric for the tracking accuracy. In addition, to evaluate the effectiveness of resource utilization of each approach, we measure the ratio of the number of full-size image executions to the total number of the MOT executions (referred to as the full-size execution ratio).

Fig. 6(a) and (b) compare the tracking accuracy and full-size execution ratio of the four approaches using YOLOX on the Tesla V100 and Jetson Xavier. Similar results were observed for the Jetson TX2 and Orin (also with YOLOv5 [20]). We consider four sets of MOT tasks [periods shown on the x-axis in Fig. 6(a) and (b)] that pass the test in Lemma 1, but schedulability is not guaranteed when all the tasks use full-size input images. Note that, the two computing systems provide different WCETs, resulting in different task periods in each set. The bar and line in each graph represent the average MOTA score and full-size execution ratio, respectively. The red dotted line indicates the maximum achievable tracking accuracy.

As shown in Fig. 6(a) and (b), a higher full-size execution ratio leads to higher accuracy. The accuracy of RT-MOT shows a significant decrease when the full-size execution ratio is low as observed in the third and fourth task sets of Fig. 6(a) and (b). This decline is attributed to the unique approach of RT-MOT, where it detects objects only in a partial region (i.e., the region of interest) of the input image when the full-size execution cannot be performed. In contrast, NPFP^B and NPFP^{BI} consider the down-scaled entire region, making them more resilient to a low full-size execution ratio. As the number of tasks increases, the accuracy of RT-MOT dramatically decreases for both the computing systems. However, NPFP^B and NPFP^{BI} maintain high accuracy by securing the chance of batch execution. In the case of a set of tasks with equal periods (e.g., the fourth task set), NPFP^B and NPFP^{BI} achieve maximum accuracy owing to the high chance of batch execution. The IFP approach contributes significantly to improving accuracy in both the computing systems.

Fig. 6(c) presents the run-time overhead of the scheduling algorithm for NPFP^{BI} and RT-MOT. As discussed in Section VI, the run-time complexity of NPFP^{BI} is $O(|\tau| \cdot \log(|\tau(t)|))$, which is much lower than $O(|\tau|^2 \cdot |\tau(t)|)$, the complexity of RT-MOT. As the number of MOT tasks increases, the difference between the run-time scheduling overhead of NPFP^{BI} and RT-MOT becomes larger. For example, the run-time scheduling overhead of RT-MOT is about 12 times (2.4/0.2) and 9.6 times (6.7/0.7) larger than that of NPFP^{BI}

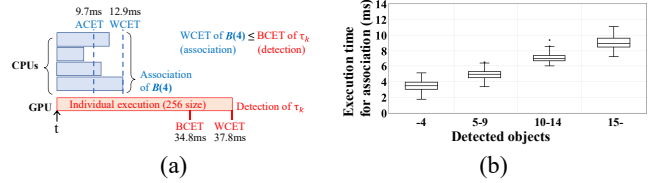


Fig. 7. CPU/GPU parallel execution example. (a) Varying WCET with the different number of objects on Jetson Orin. (b) Varying association WCET.

for a set of 12 MOT tasks on the two considered systems, respectively.

Fig. 6(d) shows the variation in average DNN inference time (excluding pre/post processing) based on the number of input images for batch execution, using different DNN models and computing systems. YOLOX demonstrates GPU resource saturation with more than the ten input images on Tesla V100 and two on Jetson Xavier, highlighting the effectiveness of batch execution in reducing inference time; similar trends were observed for YOLOv5. In contrast, Faster-RCNN [21] saturates with a single input image due to its two-stage design, which splits computation into region proposal and classification, resulting in higher serialization during classification as noted in [4].

VIII. DISCUSSION

CPU/GPU Parallelism: Contrary to Batch-MOT's assumption, consider CPU/GPU parallel execution where, at time t , four associations of $\mathcal{B}(4)$ are executed in parallel on multiple CPUs, while the another task τ_k is performed individually on the GPU as shown in Fig. 7(a). If the WCET of $\mathcal{B}(4)$'s association (e.g., 12.9 ms) is less than the best case execution time (BCET) of τ_k (e.g., 34.8 ms), then τ_k 's association can be executed nonpreemptively, consistent with Batch-MOT's assumption. Our experiments confirmed that this condition always holds. Then, we conducted the accuracy evaluation, including CPU/GPU parallel execution on Jetson Orin without modifying the Batch-MOT's offline tests. The experiments demonstrated that the CPU/GPU parallel execution did not incur any deadline misses and resulted in a marginal accuracy improvement (up to 0.81%) compared to the case without CPU/GPU parallel execution (see Figure S.4(a) in supplement). The reason for the marginal improvement is that, as shown in Figure S.4(a), the average case execution time (ACET) of the association (e.g., 9.7 ms) is much smaller than that of detection, so the reduction in response time achieved by CPU/GPU parallel execution is minimal.

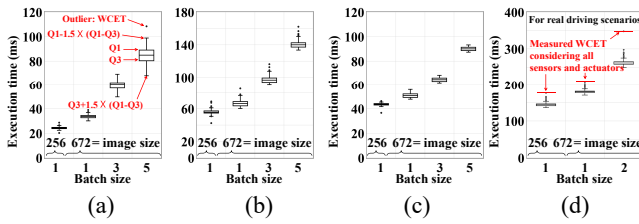


Fig. 8. Execution time measurements on Waymo Dataset (a)–(c) and our real-world driving scenarios (d).

1043 *Flexible WCET*: Association compares objects detected in
 1044 the t -th frame with those in the $(t-1)$ th frame, associating the
 1045 most similar pairs. Therefore, the WCET of association in the
 1046 t -th frame depends strictly on the number of detected objects
 1047 as shown in Fig. 7(b). If Batch-MOT splits an MOT task
 1048 into detection and association subtasks and allows scheduling
 1049 decisions between them, the WCET of the association subtask
 1050 (measured offline based on the number of objects) can be
 1051 dynamically determined by the number of objects detected in
 1052 the detection subtask. This approach would require additional
 1053 alters to the Batch-MOT’s current schedulability tests.

1054 *WCET Measurement*: Fig. 8 shows the execution time
 1055 for the down-scaled (256×256) images and batch execution
 1056 of full-size (672×672) images for up to five MOT tasks
 1057 using YOLOX on the four computing systems evaluated in
 1058 Section VII. Measurements are with 1000 iterations to obtain
 1059 WCET [e.g., Outlier: WCET in Fig. 8(a)]; more details are
 1060 in Figure S.2 and Table S.2 in the supplement. The Waymo
 1061 Dataset was used for Tesla V100, Jetson Xavier, and Orin,
 1062 while real driving scenario videos were used for Jetson
 1063 TX2. Fig. 8(d) considers all the communication overheads,
 1064 including sensors (e.g., camera and LiDAR) and actuators.
 1065 Note that, Batch-MOT does not predict execution time at run-
 1066 time but uses offline WCET. Recent studies on offline WCET
 1067 of DNN execution [2], [3], [4], [11], [18] are widely accepted.
 1068 It is generally reasonable to assume an upper bound with
 1069 high confidence through intensive measurement plus a safety
 1070 margin.

IX. CONCLUSION

1072 In this article, we proposed a novel system design, Batch-
 1073 MOT, that enables batch execution of multiple MOT tasks to
 1074 maximize the tracking accuracy while providing timing guar-
 1075 antees. Using a new scheduling framework, NP_{ADAPT}, which
 1076 allows run-time execution deviations with timing guarantees,
 1077 we developed a run-time batching mechanism, NPPF^B, and
 1078 a run-time idling mechanism, NPPF^{BI}. These mechanisms
 1079 efficiently find and execute MOT tasks as a batch without
 1080 compromising timely execution. Experiments demonstrated
 1081 that Batch-MOT improves the tracking accuracy over the
 1082 state-of-the-art real-time MOT systems while ensuring timing
 1083 guarantees.

REFERENCES

1085 [1] M. Yang et al., “Re-thinking CNN frameworks for time-sensitive
 1086 autonomous-driving applications: Addressing an industrial challenge,”
 1087 in *Proc. IEEE Real-Time Embed. Technol. Appl. Symp. (RTAS)*, 2019,
 1088 pp. 305–317.

[2] W. Kang et al., “DNN-SAM: Split-and-merge DNN execution for real-time object detection,” in *Proc. 28th IEEE Real-Time Embed. Technol. Appl. Symp. (RTAS)*, 2022, pp. 160–172.

[3] D. Kang et al., “RT-MOT: Confidence-aware real-time scheduling framework for multi-object tracking tasks,” in *Proc. IEEE Real-Time Syst. Symp. (RTSS)*, 2022, pp. 318–330.

[4] S. Liu et al., “Self-cueing real-time attention scheduling in criticality-aware visual machine perception,” in *Proc. IEEE Real Time Technol. Appl. Symp. (RTAS)*, 2022, pp. 173–186.

[5] N. Wojke, A. Bewley, and D. Paulus, “Simple online and realtime tracking with a deep association metric,” in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, 2017, pp. 3645–3649.

[6] A. Bewley, Z. Ge, L. Ott, F. Ramos, and B. Uppcroft, “Simple online and realtime tracking,” in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, 2016, pp. 3464–3468.

[7] Y. Zhang, C. Wang, X. Wang, W. Zeng, and W. Liu, “FairMOT: On the fairness of detection and re-identification in multiple object tracking,” *Int. J. Comput. Vis.*, vol. 129, no. 11, pp. 3069–3087, 2021.

[8] P. Chu, J. Wang, Q. You, H. Ling, and Z. Liu, “TransMOT: Spatial-temporal graph transformer for multiple object tracking,” 2021, *arXiv:2104.00194*.

[9] J. Pang et al., “Quasi-dense similarity learning for multiple object tracking,” in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021, pp. 164–173.

[10] H. Zhou, S. Bateni, and C. Liu, “S³DNN: Supervised streaming and scheduling for GPU-accelerated real-time DNN workloads,” in *Proc. IEEE Real-Time Embed. Technol. Appl. Symp. (RTAS)*, 2018, pp. 190–201.

[11] Y. Xiang and H. Kim, “Pipelined data-parallel CPU/GPU scheduling for multi-DNN real-time inference,” in *Proc. IEEE Real-Time Syst. Symp. (RTSS)*, 2019, pp. 392–405.

[12] “NVIDIA Orin developer kit.” Accessed: Mar. 27, 2023. [Online]. Available: <https://www.nvidia.com/ko-kr/autonomous-machines/embedded-systems/jetson-orin/>

[13] Z. Ge, S. Liu, F. Wang, Z. Li, and J. Sun, “YOLOX: Exceeding YOLO series in 2021,” 2021, *arXiv:2107.08430*.

[14] “NVIDIA xavier developer kit.” Accessed: Jul. 9, 2018. [Online]. Available: <https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/jetson-agx-xavier>

[15] “Supplement.” Accessed: Feb. 8, 2024. [Online]. Available: <https://www.bit.ly/24EMSOFT-Batch-MOT-supplement>

[16] A. Soyyigit, S. Yao, and H. Yun, “Anytime-Lidar: Deadline-aware 3D object detection,” in *Proc. IEEE Int. Conf. Embed. Real-Time Comput. Syst. Appl. (RTCSA)*, 2022, pp. 31–40.

[17] S. Heo, S. Jeong, and H. Kim, “RTScale: Sensitivity-aware adaptive image scaling for real-time object detection,” in *Proc. Leibniz Int. Proc. Inform. (LIPIcs)*, vol. 231, 2022, pp. 1–22.

[18] S. Lee and S. Nirjon, “SubFlow: A dynamic induced-subgraph strategy toward real-time DNN inference and training,” in *Proc. IEEE Real-Time Embed. Technol. Appl. Symp. (RTAS)*, 2020, pp. 15–29.

[19] S. Liu et al., “On removing algorithmic priority inversion from mission-critical machine inference pipelines,” in *Proc. IEEE Real-Time Syst. Symp. (RTSS)*, 2020, pp. 319–332.

[20] “YOLOv5.” Accessed: Nov. 23, 2022. [Online]. Available: [Online]. Available: <https://github.com/ultralytics/yolov5>

[21] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards real-time object detection with region proposal networks,” in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 28, 2015, pp. 1–9.

[22] L. George, N. Rivierre, and M. Spuri, “Preemptive and non-preemptive real-time uniprocessor scheduling,” INRIA, Paris, France, Rep. RR-2966, 1996. [Online]. Available: <https://who.rocq.inria.fr/Laurent.George/#Publication>

[23] G. Yao, G. Buttazzo, and M. Bertogna, “Feasibility analysis under fixed priority scheduling with fixed preemption points,” in *Proc. IEEE Int. Conf. Embed. Real-Time Comput. Syst. Appl. (RTCSA)*, 2010, pp. 71–80.

[24] T.-Y. Lin et al., “Microsoft COCO: Common objects in context,” in *Proc. 13th Eur. Conf. Comput. Vis. (ECCV)*, 2014, pp. 740–755.

[25] P. Sun et al., “Scalability in perception for autonomous driving: Waymo open dataset,” in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2020, pp. 2446–2454.

[26] J. Lehoczky, L. Sha, and Y. Ding, “The rate monotonic scheduling algorithm: Exact characterization and average case behavior,” in *Proc. IEEE Real-Time Syst. Symp. (RTSS)*, 1989, pp. 166–171.

[27] E. Ristani, F. Solera, R. Zou, R. Cucchiara, and C. Tomasi, “Performance measures and a data set for multi-target, multi-camera tracking,” in *Proc. Eur. Conf. Comput. Vis. Workshops (ECCVW)*, 2016, pp. 17–35.