

DeepTrust^{RT}: Confidential Deep Neural Inference Meets Real-Time!

Mohammad Fakhruddin Babar ✉ 

Electrical Engineering and Computer Science, Washington State University, Pullman, WA, USA

Monowar Hasan ✉ 

Electrical Engineering and Computer Science, Washington State University, Pullman, WA, USA

Abstract

Deep Neural Networks (DNNs) are becoming common in “learning-enabled” time-critical applications such as autonomous driving and robotics. One approach to protect DNN inference from adversarial actions and preserve model privacy/confidentiality is to execute them within trusted *enclaves* available in modern processors. However, running DNN inference inside limited-capacity enclaves while ensuring timing guarantees is challenging due to (a) large size of DNN workloads and (b) extra switching between “normal” and “trusted” execution modes. This paper introduces new time-aware scheduling schemes—DeepTrust^{RT}—to *securely* execute deep neural inferences for learning-enabled real-time systems. We first propose a variant of EDF (called DeepTrust^{RT}-LW) that *slices* each DNN layer and runs them sequentially in the enclave. However, due to extra context switch overheads of individual layer slices, we further introduce a novel *layer fusion* technique (named DeepTrust^{RT}-FUSION). Our proposed scheme provides hard real-time guarantees by *fusing* multiple layers of DNN workload from multiple tasks; thus allowing them to fit and run concurrently within the enclaves while maintaining real-time guarantees. We implemented and tested DeepTrust^{RT} ideas on the Raspberry Pi platform running OP-TEE+DarkNet-TZ DNN APIs and three DNN workloads (AlexNet-squeezed, Tiny Darknet, YOLOv3-tiny). Compared to the layer-wise partitioning approach (DeepTrust^{RT}-LW), DeepTrust^{RT}-FUSION can schedule up to 3x more tasksets and reduce context switches by up to 11.12x. We further demonstrate the efficacy of DeepTrust^{RT} using a flight controller (ArduPilot) case study and find that DeepTrust^{RT}-FUSION retains real-time guarantees where DeepTrust^{RT}-LW becomes unschedulable.

2012 ACM Subject Classification Computer systems organization → Real-time systems; Security and privacy → Systems security

Keywords and phrases DNN, TrustZone, Real-Time Systems

Digital Object Identifier 10.4230/LIPIcs.ECRTS.2024.10

Supplementary Material *Software (Source Code)*: <https://github.com/CPS2RL/DeepTrust-RT>

Funding This research is partly supported by the US National Science Foundation Award 2312006 and Washington State University Grant PG00021441. Any findings, opinions, recommendations, or conclusions expressed in the paper are those of the authors and do not necessarily reflect the sponsor’s views.

1 Introduction

The emergence of modern IoT-specific applications (such as autonomous vehicles, drones, and cognitive robots, among others) coupled with advances in computing power and hardware efficiency pushed artificial intelligence toward embedded devices. Engineers in modern safety-critical applications are progressively deploying more complex deep learning models to meet the need for on-device intelligence [1]. Many safety-critical learning-enabled systems also have stringent timing (viz., “real-time”) requirements. For example, an autonomous vehicle must periodically scan and recognize objects in its surroundings. This is often performed by executing a deep neural network (DNN) inference chain. Any delay in the object recognition



© Mohammad Fakhruddin Babar, Monowar Hasan;
licensed under Creative Commons License CC-BY 4.0
36th Euromicro Conference on Real-Time Systems (ECRTS 2024).

Editor: Rodolfo Pellizzoni; Article No. 10; pp. 10:1–10:24

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

10:2 DeepTrust^{RT}: Confidential Deep Neural Inference Meets Real-Time!

45 process may jeopardize decision-making, thus threatening the safety of the vehicle, passengers,
46 and others around it.

47 Executing DNN models on end-user devices introduces new security and confidentiality
48 challenges. Since most autonomous systems collect sensitive information (e.g., location
49 and driving dynamics, reconnaissance images, medical records), data leakage from the
50 learning task results in privacy concerns. Further, a compromised system could redistribute
51 proprietary models (e.g., parameters, intermediate results, final outputs), leaking the
52 intellectual property of the model provider. For instance, researchers have demonstrated
53 attacks such as membership inference [2, 3], fault injection [4, 5], and input reconstruction [6],
54 which can leak private model information and cause misclassification.

55 One way to tackle the DNN confidentiality problem is to execute the inference tasks
56 inside trusted enclaves such as Intel SGX [7] or ARM TrustZone [8] available in modern
57 processors. Enabling trusted execution for DNN workloads is challenging as most DNN
58 tasks are compute/memory-heavy and do not fit within the enclave memory. To put
59 this in context, VGG-16 [9], an image classification task, requires 528 MB of memory for
60 runtime computations. In contrast, OP-TEE [10], an open-source TrustZone stack, only
61 supports 16 MB of enclave memory. To address this issue, researchers proposed various
62 techniques to “slice” and execute DNN workload inside trusted enclaves [11]. However, they
63 are designed for general-purpose mobile/edge computing platforms (i.e., do not provide
64 real-time guarantees). Simply retrofitting existing frameworks without considering periodic,
65 deadline-based real-time tasks will not effectively ensure the dependability requirements of
66 learning-enabled hard real-time systems. For instance, although we can split the DNN model
67 into multiple partitions [12–15], as we shall see in this paper (Section 4.2), slicing and moving
68 inference tasks back and forth between the trusted and normal execution mode results in
69 extra delays due to high context switch overheads. This may cause some (or all!) critical
70 tasks to miss their deadlines [16].

71 Hence, our research aims to address the following problem.

The confidentiality of DNN models running on an untrusted device can be maintained by executing them within trusted enclaves.

Challenge: *How do we ensure compute-heavy real-time DNN tasks fit in limited capacity enclaves while retaining their timing guarantees (deadlines)?*

72
73 In response to the above problem, we develop scheduling models (called DeepTrust^{RT})
74 that ensure a set of learning-enabled real-time tasks can retain model confidentiality. Our
75 first attempt to make the DNN inference tasks *trusted* and *time-aware* relies on a slicing
76 mechanism that partitions DNN models *layer-by-layer* [13]. The idea is to sequentially send
77 one DNN layer at a time to the enclave, perform its computation, and get the results back.
78 However, a single layer may often not fit in the enclave due to its large size. Hence, we
79 use the Deep Compression [17] to reduce DNN model size considering the enclave capacity
80 (Section 4.1). We then use the compressed model and enable real-time scheduling capabilities
81 for the existing (non-real-time) layer-wise partitioning idea. We call our first approach
82 DeepTrust^{RT}-LW. We also derive related schedulability conditions (Section 4.2).

83 We find that despite real-time guarantees, DeepTrust^{RT}-LW results in poorer throughput
84 (i.e., fewer tasks are schedulable) due to high context switch overheads introduced by
85 the layer-by-layer partitioning technique. Hence, we further optimize DeepTrust^{RT}-LW
86 with a novel “fusion” approach that selectively *groups multiple layers from multiple tasks*,
87 considering enclave capacity and deadline constraints (Section 5). We name this technique
88 DeepTrust^{RT}-FUSION. Figure 1 illustrates the key intuition of DeepTrust^{RT}-FUSION for a

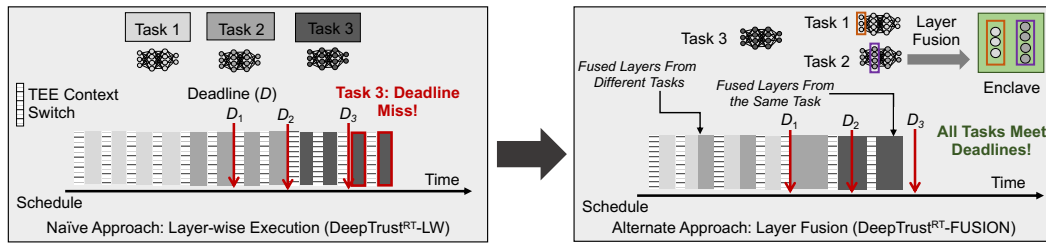


Figure 1 High-level schematic of the scheduling techniques used in the work. Due to the large size of a DNN model, often it is not feasible to fit within the enclave. If we slice the model *layer-by-layer* to fit in the enclave and send them sequentially (left rectangle, named DeepTrust^{RT}-LW), extra context switch overheads may violate real-time constraints. Hence, we also introduce a novel “layer fusion” technique (named DeepTrust^{RT}-FUSION, right rectangle) that *groups multiple layers from multiple tasks* together to reduce context switch costs and results in better schedulability.

89 three-task system. When DNN layers are sent sequentially to the enclave, extra context switch
 90 overheads cause longer response time, and one of the tasks misses the deadline. In contrast,
 91 fusing multiple layers saves context switch delays, thus resulting in faster response times.
 92 As a result, all tasks meet deadlines. Our proposed layer fusion approach (a) better utilizes
 93 enclave capacity, and (b) reduces context switches (normal-to-secure and secure-to-normal)
 94 incurred by the layer-by-layer partitioning technique (up to 11.12x for a set of 25 tasks at 50%
 95 utilization, see Section 6). As a result, DeepTrust^{RT}-FUSION achieves better schedulability
 96 compared to the DeepTrust^{RT}-LW approach.

97 In this work, we focus on ARM-based enclaves (viz., TrustZone) as ARM is the dominant
 98 architecture for embedded applications. However, our ideas are general and can be adapted
 99 to other enclaves (such as SGX) without loss of generality. We tested both DeepTrust^{RT}-LW
 100 and DeepTrust^{RT}-FUSION on three realistic workloads (e.g., AlexNet-squeezed [18], Tiny
 101 Darknet [18], YOLOv3-tiny [19]) running on Raspberry Pi and performed design-space
 102 exploration (Section 6). We also conducted a case study using the ArduPilot UAV autopilot
 103 system [20] and DNN-enabled workload (YOLOv3-tiny, Tiny Darknet) and found that
 104 DeepTrust^{RT}-FUSION can meet all deadlines for critical tasks while DeepTrust^{RT}-LW misses
 105 some (Section 6.2).

106 **Our Contributions.** This research is one of the first efforts to enable time-aware trusted
 107 DNN execution for learning-enabled real-time systems. Our contributions are as follows:

- 108 ■ Enabling *timing guarantees* for performing *confidential deep inference* in latency-critical
 109 learning-enabled systems.
- 110 ■ A new *scheduling framework and analytical model* (DeepTrust^{RT}-LW) to determine the
 111 feasibility of deploying a given real-time DNN workload on TrustZone enclaves (Section 4).
- 112 ■ A *novel task fusion approach* (DeepTrust^{RT}-FUSION) to further reduce TEE context
 113 switch overheads while retaining real-time guarantees (Section 5).

114 We performed thorough design-space exploration using three DNN architectures (AlexNet-
 115 squeezed, Tiny Darknet, YOLOv3-tiny). Our case study on a UAV system (ArduPilot)
 116 running on Raspberry Pi further demonstrates the efficacy and feasibility of the proposed
 117 techniques.

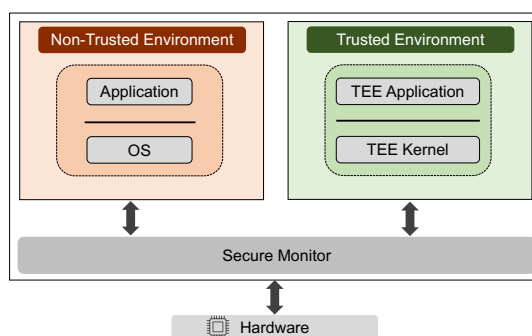
118 We now start with a background on trusted enclaves (TrustZone) and DNN architecture
 119 (Section 2) before introducing our model and related assumptions (Section 3).

120 2 Background

121 2.1 Trusted Execution and ARM TrustZone

122 Trusted execution environments (TEEs) offer a secure and isolated runtime environment.
 123 TEEs ensure confidentiality and integrity for the code and data that can not be exploited even
 124 if the host (i.e., main) OS is compromised. Among various off-the-shelf TEE implementations,
 125 Intel SGX [7] and ARM TrustZone [8] are two widely used technologies in many IoT/mobile
 126 applications. SGX is usually used for general-purpose computers and servers. In contrast,
 127 TrustZone architecture is more suitable for embedded applications and hence is the focus of
 128 our work.

129 The runtime operations in TrustZone are divided into “normal” and “secure”
 130 worlds, each having its own kernel, user, and memory space (see Fig. 2). In the
 131 normal world, a conventional operating system (e.g., Linux/RTOS) provides the
 132 execution environment, whereas the secure world uses a minimal trusted kernel (e.g.,
 133 OP-TEE [10]). The state of the current processor is determined by a specialized
 134 bit called the non-secure (NS) bit. The NS bit has two states: $NS = 1$ for non-secure
 135 execution and $NS = 0$ for secure execution. TrustZone employs a customized mechanism
 136 known as secure monitor call (SMC) to switch between these two states. When an SMC
 137 instruction is invoked in the normal world, the processor cores perform a context switch
 138 from the normal world to the secure world and pause normal world operations. As a security
 139 precaution, the normal world cannot access secure memory, while the secure world can access
 140 normal world memory. TrustZone also ensures the isolation of external peripherals.
 141
 142
 143
 144
 145
 146

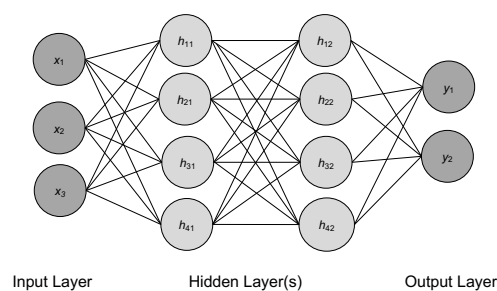


147 ■ **Figure 2** TrustZone architecture.

147 2.2 Deep Neural Inference

148 DNNs comprise an input layer, one or more
 149 hidden layers, and an output layer. Each
 150 layer consists of interconnected nodes,
 151 where the connections between nodes are
 152 represented as edges, each associated with
 153 a weight, and each node has an associated
 154 threshold value. When the output of a
 155 node exceeds its threshold, the node is
 156 activated, and the data is propagated to
 157 the next layer. Figure 3 illustrates a
 158 simplified DNN architecture.

159 DNN algorithms establish models
 160 based on training data, enabling them to generate predictions without explicit programming.
 161 During the inference process, input data is passed through the layers, and each layer performs
 162 matrix multiplications on the data. The final layer’s outputs can be numerical or classified
 163 outputs, depending on the specific application.



164 ■ **Figure 3** A simplified neural network structure. The input and output layers are connected through multiple hidden layers.

2.3 Confidential DNN

Many DNN-based applications (such as image processing, object detection, medical records, and financial transactions) deal with sensitive data and require protection against tampering or theft of intellectual property. When protecting model parameters is needed, an emerging approach is to execute critical DNN layers inside trusted enclaves. As enclaves have limited memory and DNN models are generally large, one common approach (used for general-purpose systems) is to execute the DNN workload layer-by-layer. This is known as “layer-based partitioning,” where each layer forms an independent partition. Each partition includes weights and biases and is stored in a separate encrypted file. The enclave stores the decryption key. The (encrypted) partition file is loaded into shared memory and decrypted by a trusted application on the secure side.

Note: *We do not intend to modify or improve the confidential DNN techniques—a large number of literature exists to show how DNN models can be protected using TrustZone or other TEEs [13–15, 21]. Our focus in this paper is to make existing (non-time-aware) confidential DNN ideas used for general-purpose systems adaptable for time-critical applications (viz., a set of periodic tasks with deadline constraints controlled by a real-time scheduling policy).*

3 Model and Assumptions

3.1 System Model

We consider a uniprocessor real-time system running on a TEE-enabled platform. The system consists of n real-time tasks $\Gamma = \{\tau_1, \dots, \tau_n\}$ performing DNN inference. Each task τ_i is characterized by $\tau_i = \{C_i^a, T_i, D_i, L_i, \mathcal{W}_i\}$, where C_i^a is the worst-case execution time (WCET) of the task inside the enclave, T_i is the period of task τ_i , D_i is the deadline, L_i is the number of layers of the DNN task, and \mathcal{W}_i is the set of sizes of each layer of the DNN task τ_i where $\mathcal{W}_i = \{w_{i1}, w_{i2}, \dots, w_{iL}\}$. Here, w_{ik} is the size of the weights associated with the edges between nodes (neurons), activation, and bias of nodes. In addition, W_i is the size of the DNN task τ_i where $W_i = \sum_{k=1}^{L_i} w_{ik}$. As mentioned earlier (Section 2.3), each layer partition, which includes weights and biases, is stored in an encrypted file. This encrypted file is loaded into shared memory and decrypted by a trusted application on the secure side. Let us denote $C_i^a = C_i^{dec} + C_i^{com}$ is the computation inside the enclave, where C_i^{dec} is the time required for decryption of the layers information and C_i^{com} is the computation time of task τ_i .

We assume the tasks follow the earliest deadline first (EDF) scheduling policy [22]. We consider an implicit deadline system (i.e., each task’s period is equal to its deadline, $D_i = T_i$). The taskset Γ is “schedulable” if the response time of each task (the time between arrival to completion) is less than its deadline. The trusted enclave has a finite capacity δ , i.e., it can execute $L_i \geq 1$ layers together as long as the total resource requirements of those layers are less than δ . We consider the size of each layer of a task less than δ . Invoking a TEE session involves a series of API calls. For instance, OP-TEE OS [10] requires 5 API calls for instantiating and terminating a TEE session (see Table 1). Each time the DNN layers enter the enclave, the data needs to be transferred into the enclave. Let $C_{s_i}^{st}$ be the SMC setup time and $C_{s_i}^d$ be the SMC cleanup time. Hence, $C_i^{cs} = C_{s_i}^{st} + C_{s_i}^d$ captures this data copy overhead. Note that the parameter C_i^{cs} is not part of the worst-case execution time (C_i^a). If a task requires n_i^{cs} many context switches (to-and-from normal to secure world), the total data copy overhead will be $n_i^{cs} \times C_i^{cs}$. In Section 5.2, we derive bounds on the number of context switches.

10:6 DeepTrust^{RT}: Confidential Deep Neural Inference Meets Real-Time!

```
1 Prepare session with the TA
2 Begin darknet
3 ...
4 ...
5 # darknetp:TEEC Invoke_Command(CONV) failed
```

■ **Listing 1** OP-TEE Log: Failed Invocation of a Large Model on Raspberry Pi.

209 Existing TEE implementations (for instance, OP-TEE) use non-preemptive enclave
210 execution. We incorporate this behavior, i.e., when a task performs DNN inference inside
211 the enclave, other higher-priority tasks will be “blocked” until the currently running task
212 releases the enclave.

213 3.2 Assumptions on Adversarial Capabilities

214 The pre-trained model (e.g., parameter and hyperparameter values such as structure or
215 properties of the particular DNN) is deployed to the real-time platform before the system is
216 operational. We consider an adversary seeks to gain access to sensitive information of the
217 model. Our focus is on protecting the inference operations of the DNN model. Attackers
218 may have access to the input data, but they will not obtain any information about the model
219 architecture or the final inference as long as they execute within the enclave. The attackers
220 may know all the task periods and their execution times. We further assume that the
221 adversaries cannot bypass the TEE protection mechanisms. Note that similar assumptions
222 are used by other researchers [21, 23].

223 Following the convention, we assume that the parameters of the pre-trained model are
224 kept encrypted in the local storage. The hyperparameters of the model are kept unencrypted
225 in the normal world as they are generally well-known to the public and do not disclose any
226 sensitive information about the input or training data [24]. During inference (i.e., invocation
227 of a job), the input data and model parameters are loaded into the enclave memory, and the
228 required inference operation is carried out after the decryption of the model parameters.

229 4 Time-Aware Confidential Deep Learning

230 In the vanilla case (i.e., when model confidentiality is not a concern), the weights and biases of
231 each neuron in a DNN architecture can be loaded into memory to calculate neuron activation.
232 However, a system with confidentiality requirements (execute models within an enclave)
233 presents challenges when it comes to preloading all the necessary values (e.g., weights, biases)
234 due to limited enclave size, which could be as low as 8 MB for some systems [25]. In contrast,
235 most DNN models need 100+ MB [9]. If a DNN model is too large, then the model may fail
236 to execute inside the enclave. To illustrate this, we conducted experiments on Raspberry Pi
237 running OP-TEE and Darknet [12]. For large models, Darknet was unable to load to model
238 (Listing 1). Hence, we used a model compression technique using *Deep Compression* [17]
239 to reduce the model size as we present below. Listing 2 shows the case after trimming that
240 allowed us to load and run the model successfully.

241 4.1 Resizing (Trimming) the Model

242 Deep Compression is a three-stage pipeline that reduces the storage requirement of neural
243 networks by 35x to 49x without compromising their accuracy. The pipeline consists of

```

1 Prepare session with the TA
2 Begin darknet
3 ...
4 ...
5 user CPU start: 0.029851; end: 0.029851
6 kernel CPU start:2.916167; end: 2.917143
7 Max: 2756 kilobytes
8 vmsize:545460850536; vmrss:365072222916;
9   vmdata:545460847464; vmstk:132;
10   vmexe:412; vmlib:545460848828
11
12 // Successful inference and job completion

```

■ **Listing 2** OP-TEE Log: Successful Inference.

■ **Algorithm 1** Model Compression

1: **Input:** w_{ij} , λ
2: **Output:** Compressed Size (w'_{ij})
3: Prune the network below a certain threshold λ following state-of-the-art techniques [26].
4: Retrain the network.
5: Quantize the weights of model: $r = \frac{mb}{m \log_2 k + kb}$ \triangleright *Plugging the value of r from approximation percentage (i.e., $(1 - \theta_{ij}\%$) to get the value of k*
6: Huffman coding to the quantized weights \triangleright *final compressed weight*
7: **return** w'_{ij}

244 pruning, trained quantization, and Huffman coding [25]. The first stage prunes the network
245 by learning only the essential connections, and the second stage quantizes the weights to
246 enforce weight sharing. Finally, the pipeline applies Huffman coding. The method reduces
247 the storage required by AlexNet-squeezed by 35x (from 240 MB to 6.9 MB), and VGG-16 by
248 49x (from 552 MB to 11.3 MB), without any significant loss of accuracy. This enables the
249 large model to fit inside TEE, tackling the memory constraints.

250 Recall that, to fit the model in the TEE, the size of each layer must be less than the
251 enclave capacity δ . For a given DNN task τ_i , W_i is the size of the task, L_i is the total number
252 of layers, and then the set of size of the layers is $\mathcal{W}_i = \{w_{i1}, \dots, w_{iL_i}\}$, where $W_i = \sum_{j=1}^{L_i} w_{ij}$.
253 We check $\forall w_{ij}, w_{ij} < \delta$. If $w_{ij} > \delta$, we calculate the approximation $\theta_{ij} = w_{ij} - \delta$ required
254 for this layer. The approximate percentage is defined by $\theta_{ij}\% = \theta_{ij}/w_{ij}$. The first stage
255 of Deep Compression (see Algorithm 1) prunes the network by learning only the required
256 connections, and the second stage quantizes the weights to enforce weight sharing. In general,
257 for a network with m connections, each connection is represented by b bits, constraining the
258 connections to have only k shared weights will result in a compression rate of

$$259 \quad r = \frac{mb}{m \log_2 k + kb}. \quad (1)$$

260 Let us assume $(1 - \theta_{ij}\%)$ is the desired value for the compression rate r . Plugging the desired
261 compression rate $r = (1 - \theta_{ij}\%)$, we can find the cluster size k based on Eq. (1). After
262 checking and resizing all the layers, we will get the desired task ready that can fit within the
263 enclave (see Algorithm 2).

264 4.1.1 Formal Description of Model Trimming

265 Algorithm 1 and Algorithm 2 formally present our ideas for trimming a given DNN model.
266 The model compression process (Algorithm 1) initially prunes the network below a threshold
267 λ to remove less critical connections (Line 3). For this, we use the techniques Han et al. [26]

Algorithm 2 Resize all Layers

```

1: Input: Model size set ( $\mathcal{W}_i$ ), TEE Capacity  $\delta$ 
2: Output: Resized model size set ( $\mathcal{W}'_i$ )
3:  $\mathcal{W}'_i = []$   $\triangleright$  Initialize to an  $n$  empty array
4: for  $j = 1$  to  $L_i$  in  $\mathcal{W}_i$  do
5:   if  $w_{ij} > \delta$  then
6:     Optimized the layer using Algorithm 1
7:      $\mathcal{W}'_i \leftarrow w_{ij}'$ 
8:   else
9:      $\mathcal{W}'_i \leftarrow w_{ij}$ 
10:  end if
11:   $j \leftarrow j + 1$ 
12: end for
13: return Resized model ( $\mathcal{W}'_i$ )

```

268 described. We rerun the network to learn the final weights with pruned networks (Line 4).
 269 Then, the algorithm quantizes weights, determining the value of shared weights k plugging
 270 desired compression rate $r = (1 - \theta_{ij}\%)$ in Eq. (1) (Line 5). Finally, we apply Huffman
 271 coding [25] to the quantized weights (Line 6).

272 Following the steps in Algorithm 1 allows us to resize a single layer. We then use
 273 Algorithm 2 to resize *all* the layers of a task τ_i so that we can fit at least a single layer at
 274 a given time inside TEE. Algorithm 2 examines each layer of τ_i to determine if it exceeds
 275 the TEE capacity δ (Lines 4-12). For instance, if $w_{ij} > \delta$, the layer is optimized using
 276 Algorithm 1 (Line 6) and stores the resized layers information in \mathcal{W}'_i (Line 7). If $w_{ij} < \delta$,
 277 unchanged value of w_{ij} is stored in \mathcal{W}'_i (Line 9). This process is repeated for each layer of τ_i ,
 278 and resized layer information is stored in \mathcal{W}'_i .

279 We note that a compressed model may not fit into TEE due to limited enclave size (i.e.,
 280 $W_i = \sum w_{ij} > \delta$). In such cases, a known technique (used in general-purpose systems) is to
 281 split the DNN model into smaller parts [13, 27]. This partitioning method is beneficial as
 282 the only values needed at a given time are the activation of the previous layer, the weights,
 283 and biases for the current layer. To illustrate, for two fully connected layers, each with z
 284 neurons, it would require z activations, $z \times z$ weights, and z biases. This effectively reduces
 285 the instantaneous memory requirement to that of a single layer. The largest layer in the
 286 model determines the minimum amount of secure world memory needed for confidential
 287 DNN execution. However, this approach partitions each layer and transfers results back
 288 and forth from secure to the normal world. This extra context switch overhead could be
 289 a bottleneck for real-time applications. Thus, we need timing analysis and schedulability
 290 conditions to ensure all tasks retain real-time constraints. The following shows how we can
 291 adapt the EDF scheduler for a conventional layer-wise partitioning technique. We refer to
 292 this EDF variant as DeepTrust^{RT}-LW.

293 4.2 DeepTrust^{RT}-LW: Real-Time Layer-wise DNN Execution

294 Our first approach—DeepTrust^{RT}-LW—sequentially partitioned the layers. They are then
 295 transferred to the enclave and scheduled using EDF policy. This is feasible since there is
 296 no cross-dependency between any two layers, and each layer can be computed sequentially
 297 independently [27]. Traditional EDF schedulability conditions often involve checking many
 298 relative deadline points to assess the schedulability of a taskset up to the hyperperiod [28, 29].
 299 To speed up this process, Zhang et al. propose an improved algorithm (called QPA) that
 300 significantly reduces the computation required to check each relative deadline [28]. To
 301 determine the schedulability conditions for DeepTrust^{RT}-LW, we use the existing QPA-

Algorithm 3 DeepTrust^{RT}-LW Schedulability Checking

```

1: Input: Real-time taskset ( $\Gamma$ )
2: Output: Taskset schedulability
3:  $t \leftarrow \max\{T_1, T_2, \dots, T_n\}$ 
4:  $T_{min} \leftarrow \min\{T_1, T_2, \dots, T_n\}$ 
5: while  $t > T_{min}$  do
6:    $h(t) \leftarrow \sum_{i=1}^n \lfloor \frac{t}{T_i} \rfloor (C_i^a + L_i \times C^{cs})$   $\triangleright$  Calculate  $h(t)$  for the given  $t$ 
7:   if  $h(t) + b(t) > T_{min} \wedge h(t) + b(t) < t$  then
8:      $t \leftarrow h(t) + b(t)$ 
9:   else if  $h(t) + b(t) \leq T_{min}$  then
10:     Taskset is schedulable
11:     Break
12:   else
13:     Taskset is not schedulable
14:     Break
15:   end if
16: end while

```

302 based EDF timing analysis technique [28] and adapt it to our DNN-based workload. We
 303 choose QPA-based analysis instead of others [29] because (a) it provides us a modular
 304 model that can be extended to more general tasksets (arbitrary deadline systems) and (b)
 305 computational complexity of QPA is an offline (design-time) analysis which will not affect
 306 runtime performance of DeepTrust^{RT}.

307 Recall that the execution within the enclave is non-preemptive. Such behavior is modeled
 308 by incorporating a “blocking” delay in schedulability analysis. In EDF scheduling with
 309 blocking, a set of tasks is schedulable if $\forall t > 0, h(t) + b(t) \leq t$, where $h(t)$ is the *processor*
 310 *demand function* and $b(t)$ is the blocking delay [30, 31]. The function $h(t)$ calculates the
 311 maximum execution time required by the system for all tasks with both their arrival times
 312 and their deadlines in a contiguous interval of length t . The demand function $h(t)$ is given
 313 by: $h(t) = \sum_{i=1}^{i=n} \lfloor \frac{t}{T_i} \rfloor C_i$. In our context, the blocking delay is $b(t) = \max\{C_j^{cs} | D_j > t\}$.

314 For DeepTrust^{RT}-LW, the computing time is given by $C_i = C_i^a + n_i^{cs} \times C_i^{cs}$, where
 315 n_i^{cs} is the total SMC context switches. Hence, we can rewrite $h(t)$ as follows: $h(t) =$
 316 $\sum_{i=1}^{i=n} \lfloor \frac{t}{T_i} \rfloor (C_i^a + n_i^{cs} \times C_i^{cs})$, see Lemma 1 in Section 4.3 for a formal derivation. Note that,
 317 in DeepTrust^{RT}-LW, $n_i^{cs} = L_i$. The upper limit of t that needs to be checked is defined by
 318 $S = \max\{T_1, T_2, \dots, T_n\}$. The taskset is schedulable if $U < 1$ and $h(t) + b(t) \leq T_{min}$, where
 319 $T_{min} = \min\{T_1, T_2, \dots, T_n\}$.

320 4.3 Workflow and Analysis of DeepTrust^{RT}-LW

321 As mentioned before, in DeepTrust^{RT}-LW, layers are sequentially executed inside the enclave,
 322 and the tasks are scheduled using the EDF policy. Algorithm 3 presents steps for the
 323 schedulability checks. We start by finding the maximum task period in the taskset (Line 3).
 324 T_{min} stores the minimum value of the task period in the taskset (Line 4). The processor
 325 demand function $h(t)$ calculates the maximum execution time required by the system for
 326 given t (Line 6). If $h(t) + b(t) > T_{min}$ and $h(t) + b(t) < t$, we tighten the bound on processor
 327 demand to check if we can execute all ready tasks. This is done by changing the value of t
 328 to $h(t) + b(t)$ (Line 8). If $h(t) + b(t) \leq T_{min}$ at any time t , we can conclude that our system can
 329 execute all ready tasks without missing any deadlines. Therefore, the task set is schedulable
 330 (Line 9). If it finds $h(t) + b(t) > t$ at any time t , then we report that the taskset is not
 331 schedulable (Line 13).

10:10 DeepTrust^{RT}: Confidential Deep Neural Inference Meets Real-Time!

■ **Table 1** APIs Required for Invoking a TEE Call. The Overheads are Measured on Raspberry Pi.

API	Function	Overhead (μs)
TEEC_InitializeContext()	Initialize connection	240
TEEC_OpenSession()	Open a new TEE session	18000
TEEC_InvokeCommand()	Invokes a Command	280
TEEC_CloseSession()	Close the session	1180
TEEC_FinalizeContext()	Close connection	110

332 **Determining the Processor Demand Function.** The following lemma shows the
 333 expression for $h(t)$.

334 ► **Lemma 1.** *The maximum execution time required by the system contiguous interval of*
 335 *length t , is given by: $h(t) = \sum_{i=1}^{i=n} \lfloor \frac{t}{T_i} \rfloor C_i$.*

336 **Proof.** From traditional EDF timing analysis [28], $h(t) = \sum_{i=1}^{i=n} \max\{0, 1 + \lfloor \frac{t-D_i}{T_i} \rfloor\} \times C_i$.
 337 Replacing $D_i = T_i$ in the above equation (since we have an implicit deadline system) and
 338 after simplification, $h(t)$ can be rewritten as: $h(t) = \sum_{i=1}^{i=n} \max\{0, \lfloor \frac{t}{T_i} \rfloor\} \times C_i$. Note that,
 339 $\frac{t}{T_i}$ is a non-negative value. Hence, reduced form of $h(t)$ is $h(t) = \sum_{i=1}^{i=n} \lfloor \frac{t}{T_i} \rfloor C_i$. Replacing
 340 $C_i = C_i^a + n_i^{cs} \times C_i^{cs}$, $h(t)$ can be rewritten as: $h(t) = \sum_{i=1}^{i=n} \lfloor \frac{t}{T_i} \rfloor \times (C_i^a + n_i^{cs} \times C_i^{cs})$. ◀

341 **Overhead Analysis.** For a given task τ_i , the worst-case execution time of the model
 342 inside TEE is C_i^a , where $C_i^a = \sum_{j=1}^{j=L_i} C_{ij}^a$ and C_{ij}^a is the computation time for layer j . In
 343 DeepTrust^{RT}-LW, if a task τ_i has L_i number of layers, we need L_i number of context switches.
 344 The total execution time of task τ_i required in the layer-wise approach is $C_i = C_i^a + L_i \times C_i^{cs}$.
 345 We now explain the overhead of DeepTrust^{RT}-LW using a simple example.

346 ► **Example 1.** Let us assume we have three tasks τ_1, τ_2, τ_3 each having 5 layers (i.e., $L_i = 5$)
 347 and $\delta = 7$. The size of τ_1 and τ_2 is 10, and the size of τ_3 is 5 units. We cannot execute all
 348 the layers of τ_1 inside the enclave as the size of $\tau_1 > \delta$. DeepTrust^{RT}-LW requires five SMC
 349 switches from the normal to secure world for five layers for each task τ_1, τ_2 , and τ_3 . Hence,
 350 we need $3 \times 5 = 15$ SMC switches to execute these three tasks inside TEE.

351 4.4 The Need for an Efficient Scheduler

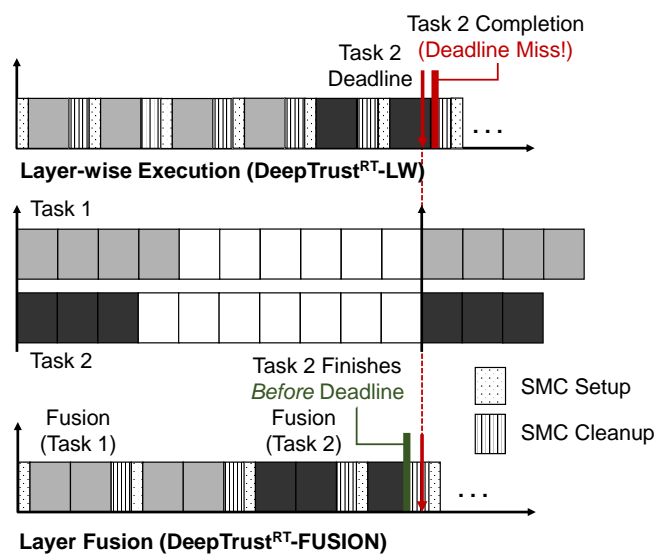
352 Despite DeepTrust^{RT}-LW ensures real-time guarantees (for schedulable tasksets), as we shall
 353 see below (and also from our evaluation in Section 6), it results in poorer schedulability. This
 354 is because each switching results in extra SMC invocation, which increases task response
 355 times. For example, OP-TEE (an open-source TrustZone port for Linux) [10] performs five
 356 API calls to initiate and teardown a TEE session. Each of these API calls takes a considerable
 357 amount of time. We carried out experiments to time each call on a Raspberry Pi platform.
 358 As the Table shows, initiating a TEE session, transferring data to/from the enclave, and
 359 cleanup steps take approximately 20 ms. In our context, each of the layer execution sessions
 360 will add up those TEE API overheads, thus potentially slowing down the inference task and
 361 may result in missed deadlines. We further illustrate this issue using a simple example.

362 ► **Example 2.** Let us consider the taskset listed in Table 2.

■ **Table 2** Example Taskset 1.

Task	L	C^{cs}/layer	C_i^a	C	T
τ_1	8	20	290	450	700
τ_2	6	20	270	390	1500
τ_3	8	20	290	450	3000

363 We now show how layer-wise execution in DeepTrust^{RT}-LW adds context switch overheads
 364 that can increase the overall execution time. There are three tasks τ_1, τ_2, τ_3 , where C_1^a, C_2^a, C_3^a
 365 are 450, 390, and 450 units respectively. The maximum blocking delay for task τ_2 is 450
 366 time units. The periods T_1, T_2, T_3 are 700, 1500, and 3000 time units, respectively. In this
 367 taskset, $\sum C_i^a/T_i = 0.69 < 1$. Let us assume the context switch overhead is 20 units per
 368 layer. Adding this context switch overhead leads execution times, C_1, C_2, C_3 to 450, 390, and
 369 450 units, respectively. As a result, the utilization is $\sum C_i/T_i = 1.05 > 1$. The taskset is
 370 not schedulable under DeepTrust^{RT}-LW since the system utilization is over 100%. In this
 371 example, we can see the summation of the actual execution time, $\sum C_i^a = 850$, and the
 372 summation of total execution time $\sum C = 1290$. This indicates an additional 34% context
 373 switching overhead in executing the taskset.



■ **Figure 4** Key intuition of model fusion: when the tasks are executed layer-wise (DeepTrust^{RT}-LW, top schedule), Task 2 misses deadline due to multiple context switch overheads. However, in a multi-layer execution approach (DeepTrust^{RT}-FUSION, bottom schedule), multiple layers are fused together, which reduces context switch overheads and all the tasks meet their deadlines.

374 To address this problem, we develop a simple yet compelling idea: instead of sending
 375 each layer sequentially, we propose to *group (fuse) multiple layers from multiple tasks*
 376 *(as long as they fit in the enclave) and send them together*. We refer to this technique
 377 DeepTrust^{RT}-FUSION. Figure 4 illustrates a high-level schematic for two tasks. In this
 378 case, DeepTrust^{RT}-LW misses deadlines for Task 2 due to multiple context switch overheads.
 379 However, when we fuse the layers in DeepTrust^{RT}-FUSION, we save context switch costs,
 380 thus allowing both tasks to meet deadlines.

10:12 DeepTrust^{RT}: Confidential Deep Neural Inference Meets Real-Time!

381 We note that task fusion has been used in literature for TEE-enabled conventional
 382 (i.e., non-learning enabled) fixed-priority real-time systems to reduce TEE context switch
 383 overheads [32]. We borrow a similar concept to group multiple layers of tasks and fit
 384 them within the enclave. For each decision instance, we group the tasks with shorter
 385 deadlines so that (a) enclave utilization (capacity usage) is maximized, i.e., fit as many
 386 layers as possible, and (b) tasks do not miss deadlines. Our selection process, as described in
 387 Section 5, is inspired by the bin-packing heuristics (such as best-fit) [33] used in partitioned
 388 multiprocessor scheduling. We now discuss DeepTrust^{RT}-FUSION scheduler in detail and
 389 derive schedulability conditions.

5 DeepTrust^{RT}-FUSION: Multi-layer Task Fusion

391 Fusing multiple layers from multiple tasks can significantly save context switch overheads
 392 compared to DeepTrust^{RT}-LW approach. For example, AlexNet-squeezed [34] has 16 layers.
 393 If the system follows DeepTrust^{RT}-LW, it needs 16 context switches (one for each layer).
 394 In contrast, assuming each layer is 1 MB in size and the enclave has 8 MB of memory, if
 395 we fuse layers using DeepTrust^{RT}-FUSION, we can finish the execution with two context
 396 switches. We now illustrate how DeepTrust^{RT}-FUSION improves schedulability. Fusion
 397 works independent of task period types (i.e., harmonic/non-harmonic), as we illustrate with
 398 an example.

399 ► **Example 3.** Let us consider the following taskset parameters (Table 3 and Table 4).

■ **Table 3** Example Taskset and Layer Size.

Task	L	Size of layers (MB)	Total Size (MB)
τ_1	8	{0.046, 0.186, 0.48, 0.39, 0.27, 5.84, 2.69, 1.50}	11.40
τ_2	6	{0.186, 0.48, 0.39, 5.84, 2.69, 1.50}	11.08
τ_3	8	{0.046, 0.186, 0.48, 0.39, 0.27, 5.84, 2.69, 1.50}	11.40

■ **Table 4** Example Taskset with Fusion Parameters.

Task	L	C^{cs}/layer	CS (fusion)	C_i^a	C	T (non-harmonic)	T (harmonic)
τ_1	8	20	2	290	330	700	700
τ_2	8	20	2	270	310	1500	1400
τ_3	8	20	2	290	330	3000	2800

400 We show for both harmonic and non-harmonic cases. Let us first consider the non-
 401 harmonic periods. In this case, $\sum C_i/T_i = 0.78 < 1$. We can calculate the schedulability
 402 conditions of as follows: (a) $t = 3000$, $h(t) = 2270$; (b) $t = 2270$, $h(t) = 1300$; and (c)
 403 $t = 1300$, $h(t) = 330$. We can see $h(t) < T_{min}$. Hence, the taskset is schedulable (recall: the
 404 same taskset is not schedulable using DeepTrust^{RT}-LW). Let us now consider the harmonic
 405 case. For harmonic periods, (a) $t = 2800$, $h(t) = 2270$; (b) $t = 2270$, $h(t) = 1300$; and (c)
 406 $t = 1390$, $h(t) = 330$. We can see $h(t) < T_{min}$. Hence, the taskset is schedulable.

5.1 Workflow of DeepTrust^{RT}-FUSION

408 DeepTrust^{RT}-FUSION aims to maximize the usage of TEE capacity. Hence, we send the
 409 maximum number of layers that TEE can support to reduce the SMC context switching

Algorithm 4 DeepTrust^{RT}-FUSION: Task Fusion and Scheduling

```

1: Input: Real-time taskset ( $\Gamma$ ), TEE-capacity  $\delta$ 
2: Output: Taskset schedulability decision
3: Compress the Model  $\triangleright$  See Algorithm 1
4: Resize Layers  $\triangleright$  See Algorithm 2
5:  $\Omega(t) = \{\mathcal{W}'_1, \mathcal{W}'_2, \dots, \mathcal{W}'_i\}$   $\triangleright$  Obtain the set of the weight of each task available at time  $t$ 
6:  $T_{hyp} = \text{LCM of } \{T_1, T_2, \dots, T_n\}$   $\triangleright T$  is the set of period of all DNN tasks
7: BEGIN  $\triangleright$  Find layers to send to TEE
8: while TRUE do
9:    $S = \text{FIND\_LAYERS\_TO\_SEND}\{\Omega(t)\}$   $\triangleright$  See Line 23 for definition
10:   Send  $S$  to TEE
11: end while
12: END

13: function FIND_TRANSITION_OF_LAYERS( $\Omega(t)$ )
14:    $i \leftarrow$  index of first task in  $\Omega(t)$ 
15:   while  $i \leq$  no of task available at time  $t$  do
16:     if  $\sum_{j=p}^{j=k} w_{ij} = \delta_1 < \delta$  and  $\sum_{j=p}^{j=k+1} w_{ij} > \delta$  then
17:        $i = i + 1$ 
18:       Remove  $w_{im}, \dots, w_{ik}$  from  $\Omega(t)$ 
19:     end if
20:   end while
21:   return  $w_{ip}, \dots, w_{ik}, w_{(i+1)p'}, \dots$ 
22: end function

23: function FIND_LAYERS_TO_SEND( $\Omega(t)$ )
24:   while  $\Omega(t) \neq \text{NULL}$  do
25:      $S = \text{FIND\_TRANSITION\_OF\_LAYERS}(\Omega(t))$   $\triangleright$  See Line 13 for definition
26:     Check schedulability using Lemma 3
27:     if Schedulable then
28:       Continue
29:     else
30:       break  $\triangleright$  Taskset is not schedulable
31:     end if
32:     if  $t \geq T_{hyp}$  then  $\triangleright T_{hyp}$  is the hyperperiod of  $T$ 
33:       break
34:     end if
35:   end while
36:   return  $S$ 
37: end function

```

410 overhead. For a given DNN task τ_i , the worst-case execution time of the model inside TEE
 411 is C_i^a , where $C_i^a = \sum_{j=1}^{j=L} C_{ij}^a$ and L_i is the number of layers. If there is L_i layers in task τ_i ,
 412 then the size of each layer will be $w_{i1}, w_{i2}, \dots, w_{iL_i}$, where $\sum_{j=1}^{j=L_i} w_{ij} = W_i$. We first check
 413 if the following condition holds: $(w_{i1} + w_{i2}) < \delta$. We find the maximum value of k where
 414 $\sum_{j=1}^{j=k} w_{ij} = \hat{\delta} < \delta$, $\sum_{j=1}^{j=k+1} w_{ij} > \delta$. If some extra capacity is left (i.e., $\delta - \hat{\delta}$), we check
 415 the subsequent task to fit within this extra space. We find the maximum value of k for the
 416 next task where $\sum_{j=1}^{j=k} w_{(i+1)j} < (\delta - \hat{\delta})$, $\sum_{j=1}^{j=k+1} w_{(i+1)j} > (\delta - \hat{\delta})$. We check all available
 417 candidate tasks at a given time t to check whether layers can fit inside the enclave. Once we
 418 obtain the schedule profile, we repeat the same steps for all subsequent task arrivals.

419 Algorithm 4 formally presents DeepTrust^{RT}-FUSION approach. The fusion decision will
 420 be made when a task (a) arrives, (b) completes, or (c) returns from the enclave. Since the
 421 scheduler keeps track of the ready-queue and SMC returns (for example, OP-TEE tear-down
 422 APIs `TEEC_CloseSession()` and `TEEC_FinalizeContext()`), we know when to perform
 423 fusion decisions. For each scheduling decision event, DeepTrust^{RT}-FUSION scheduler picks
 424 the fuse candidates (for instance, the loops in Line 8-Line 11, Algorithm 4). Let $\Omega(t)$ be the
 425 set of all tasks scheduled by using the vanilla EDF (i.e., without any TEEs) algorithm at
 426 any given time t . We first calculate the hyperperiod of the taskset (Line 6). From $\Omega(t)$, we

10:14 DeepTrust^{RT}: Confidential Deep Neural Inference Meets Real-Time!

427 find the set of layers S to send to TEE (Line 9). We find the transition point k for each task
 428 and remove layers p to k from $\Omega(t)$, where p is an integer initialized to 0 (Line 18). Then,
 429 we calculate the corresponding candidate by following the condition (Line 16). We repeat
 430 this for all subsequent tasks available at that time using the `while` loop (Line 15-20). We
 431 return all the layers $w_{ip}, \dots, w_{ik}, w_{(i+1)p'}, \dots$ (Line 21) to S that is finding the set of layers
 432 to send to TEE (Line 9). We then check the schedulability condition (see Lemma 3 for a
 433 formal derivation). If the task is schedulable, we continue to find the next candidate to send
 434 to TEE and repeat this process till hyperperiod. In the following example, we demonstrate
 435 our proposed idea.

436 **► Example 4.** Let us assume we have three tasks τ_1, τ_2, τ_3 each having 5 layers and $\delta = 7$.
 437 The size of τ_1 and τ_2 is 10, and the size of τ_3 is 5 units. We consider the size of each layer to
 438 be the same for simplicity. We cannot execute all the layers of τ_1 inside the enclave as the
 439 size of $\tau_1 > \delta$. If we execute layer-by-layer, we need five SMC switching from the normal
 440 world for five layers for each task τ_1, τ_2 , and τ_3 . If we send multiple layers of τ_1 that can
 441 be supported by TEE, it still requires two SMC switching i.e., $\{w_{11}, w_{12}, w_{13}\}, \{w_{14}, w_{15}\}$.
 442 For task τ_2 , we also need two SMC switching $\{w_{21}, w_{22}, w_{23}\}, \{w_{24}, w_{25}\}$. For task τ_3 , we
 443 need one SMC context switching $\{w_{31}, w_{32}, w_{33}, w_{34}, w_{35}\}$. Hence, we need fifteen SMC
 444 switches for layer-by-layer operations to execute these three tasks. In contrast, it is possible
 445 to perform the same objective using only five SMC switches if we can send it by multiple
 446 layers. If we send multiple layers of τ_1 , we still have some extra capacity left ($\delta - \delta_1 = 1$).
 447 In this case, we check whether it is feasible to use that space capacity. In this example,
 448 $w_{11} + w_{12} + w_{13} + w_{31} = 7 \leq \delta$. Hence, we can fuse the first three layers from τ_1 and the first
 449 layer from τ_3 , and then send them together to the enclave. If we repeat the same operations for
 450 the rest of the layers we get the following pattern: $\{w_{11}, w_{12}, w_{13}, w_{31}\}, \{w_{14}, w_{15}, w_{21}, w_{32}\}$,
 451 $\{w_{22}, w_{23}, w_{24}, w_{33}\}, \{w_{25}, w_{34}, w_{35}\}$ i.e., we only need four SMC switches.

452 5.2 Schedulability Conditions and Overhead Analysis

453 Recall that, a taskset is schedulable if $\forall t, U(t) < 1$ and $h(t) \leq t$. We now derive the
 454 expressions for $U(t)$ and $h(t)$ for DeepTrust^{RT}-FUSION.

455 **► Lemma 2.** Let $n_i^s(t)$ is the number of context switches by applying fusion for a window of
 456 duration t . System utilization $U(t)$ for a given taskset at any given time t is given by

$$457 \quad U(t) = \sum_{i=1}^{i=n} \left(\frac{\lfloor \frac{t}{T_i} \rfloor \times C_i}{t} - \frac{n_i^s(t) \times C_i^{cs}}{t} \right). \quad (2)$$

458 **Proof.** To determine the system utilization for a given taskset, we assume that each task
 459 arrives at time $t = 0$. We then calculate the number of occurrences of each task within time
 460 t using the expression $\lfloor \frac{t}{T_i} \rfloor$, where T_i represents the period of task τ_i . The overhead of each
 461 task is then given by $\lfloor \frac{t}{T_i} \rfloor \times C_i$, where C_i represents the computation time required for task
 462 τ_i . At any given time t , system utilization is $\sum_{i=1}^{i=n} \frac{\lfloor \frac{t}{T_i} \rfloor \times C_i}{t}$. However, by applying layer
 463 fusion, we can reduce context switching overhead as $\frac{n_i^s(t) \times C_i^{cs}}{t}$, where $n_i^s(t)$ is the number of
 464 context switches by applying fusion for a window of duration t . Hence, we can calculate the
 465 system utilization at any given time t as follows: $U(t) = \sum_{i=1}^{i=n} \left(\frac{\lfloor \frac{t}{T_i} \rfloor \times C_i}{t} - \frac{n_i^s(t) \times C_i^{cs}}{t} \right)$ ◀

466 In the following, we derive the processor demand function $h(t)$ for DeepTrust^{RT}-FUSION.

467 ► **Lemma 3.** *The task set Γ is schedulable under DeepTrust^{RT}-FUSION if $\forall t > 0, t < T_{max};$
 468 $h(t) + b(t) \leq t, U(t) < 1,$ where*

$$469 \quad h(t) = \sum_{i=1}^{i=n} \left(\lfloor \frac{t}{T_i} \rfloor C_i - \frac{n_i^s(t) \times C_i^{cs}}{t} \right). \quad (3)$$

470 **Proof.** The demand function $h(t)$ calculates the maximum execution time required by all
 471 tasks that have both their arrival times and their deadlines in a contiguous interval of length
 472 t . Recall that, $h(t)$ is given by $h(t) = \sum_{i=1}^{i=n} \lfloor \frac{t}{T_i} \rfloor C_i$. DeepTrust^{RT}-FUSION can reduce up
 473 to $n_i^s(t)$ context switches for each task τ_i for a window of size t . Considering this reduction,
 474 we now rewrite $h(t)$ as: $h(t) = \sum_{i=1}^{i=n} \left(\lfloor \frac{t}{T_i} \rfloor C_i - \frac{n_i^s(t) \times C_i^{cs}}{t} \right)$. ◀

475 We now calculate the performance benefits of DeepTrust^{RT}-FUSION, i.e., the reduction
 476 in SMC context switch counts when we use layer fusion.

477 ► **Lemma 4.** *If we have z fused tasks in Γ , then the total context switch reduction within
 478 the hyperperiod is $\sum_{j=1}^{j=z} (k_j - 1)C_j^{cs}$, where k_j is the number of fused layers in j^{th} fused task,
 479 C_j^{cs} is the context switch overhead.*

480 **Proof.** If we can fuse k layers from different tasks that are available at time t , then j^{th} fused
 481 task τ_j^{fused} is defined as $(C_j^{\text{fused}}, n_j^{\text{fused}})$, where C_j^{fused} is the execution time of fused task
 482 and n_j^{fused} is the SMC context switching reduction due to j^{th} fused task. If we can fuse k_j
 483 layers, then C_j^{fused} can be measured using the following equation: $C_j^{\text{fused}} = C_j^{cs} + \sum_{i=1}^{i=k_j} C_{j_i}^a$,
 484 where $C_{j_i}^a$ is the computation time at i^{th} layer. If we can fuse k layers in j^{th} fused task, we
 485 can reduce $n_j^{\text{fused}} = (k_j - 1) \times C_j^{cs}$ context switches. If we have z number of fused tasks
 486 within the hyperperiod, we can define the total context switching overhead reduction as:

$$487 \quad n^s = \sum_{j=1}^{j=z} n_j^s = \sum_{j=1}^{j=z} (k_j - 1)C_j^{cs}. \quad (4)$$

488 ◀

489 6 Evaluation

490 We evaluate our techniques on two fronts: (a) design-space exploration with various DNN
 491 workloads (Section 6.1) and (b) case study with a UAV autopilot system (Section 6.2). Our
 492 implementation is available on GitHub: <https://github.com/CPS2RL/DeepTrust-RT>.

493 6.1 Design-Space Exploration with Deep Learning Workloads

494 6.1.1 Simulation Setup

495 We evaluate the performance of our scheme using synthetically generated workloads, with
 496 parameters similar to that used in prior work [35]. We vary the system utilization from 0% to
 497 100%. For each system utilization u in the range $[0, 10, \dots, 100]\%$, we generate 200 tasksets,
 498 each taskset containing 5 to 15 tasks. Task periods are randomly selected from 50 to 100 ms.
 499 For deep learning workload, we used three popular DNN architectures: AlexNet-squeezed [18],
 500 YOLOv3-tiny [19], and Tiny Darknet [18]. We also tested with a “random workload” where
 501 we randomly generated the number of layers, task period, size of layers, and computation
 502 time. We tested with two enclave capacities (δ): 8 MB for AlexNet-squeezed and Tiny
 503 Darknet and 16 MB for YOLOv3-tiny. We note that similar sizes of enclaves are used by

10:16 DeepTrust^{RT}: Confidential Deep Neural Inference Meets Real-Time!

■ **Table 5** Systems & Workloads.

Parameters	Description
Hardware	4x ARM Cortex A53, 1 GB RAM (Raspberry Pi 3 B)
Rich OS	Linux 6.2.0
Trusted OS	OP-TEE 3.19.0
Workloads	<ul style="list-style-type: none">• AlexNet-squeezed (Image Processing)• Tiny Darknet and YOLOv3-tiny (Object Detection)• Random: weights and run times are generated randomly

■ **Table 6** Simulation Parameters.

Parameters	Value
Enclave capacity, δ	8 MB
Utilization, U	0%-100%
Period T	[50, 100] ms
Number of layers, L	[5, 24]
Weight, W	[0.01, 7]
Execution time inside TEE per layer, c_{ij}^a	[0.1, 8] ms
SMC overhead, $c_s^{st} + c_s^d$	20 ms
Number of tasks, n	[5, 25]
Number of taskset for each utilization, N_u	200

504 OP-TEE. Unless otherwise specified, we consider SMC context switch overhead ($c_s^{st} + c_s^d$)
505 to be 20 ms. Table 5 summarizes platform and workload, and Table 6 lists key simulation
506 parameters.

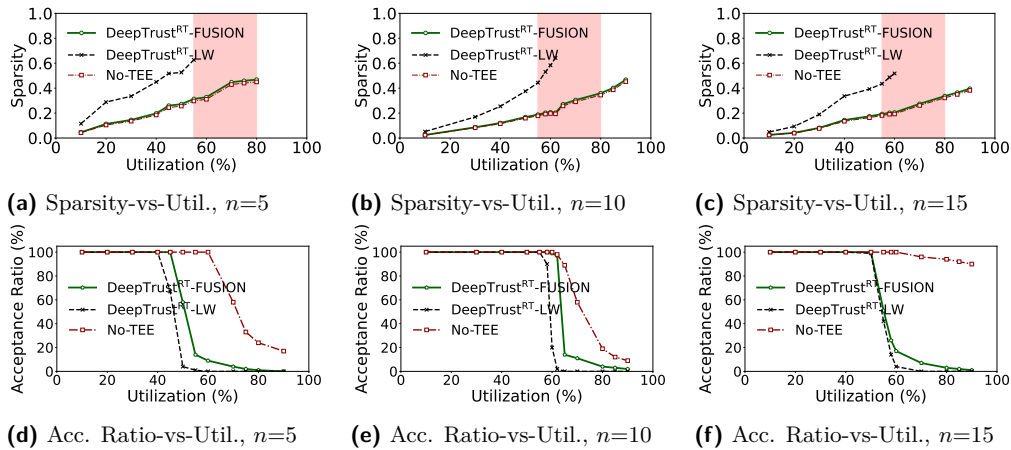
507 6.1.2 Schemes and Metrics

508 We compare DeepTrust^{RT}-FUSION with layer-wise execution technique (DeepTrust^{RT}-LW).
509 For completeness, we also study a “non-secure” variant that does not consider any enclave.
510 The schemes used in our evaluation are listed below.

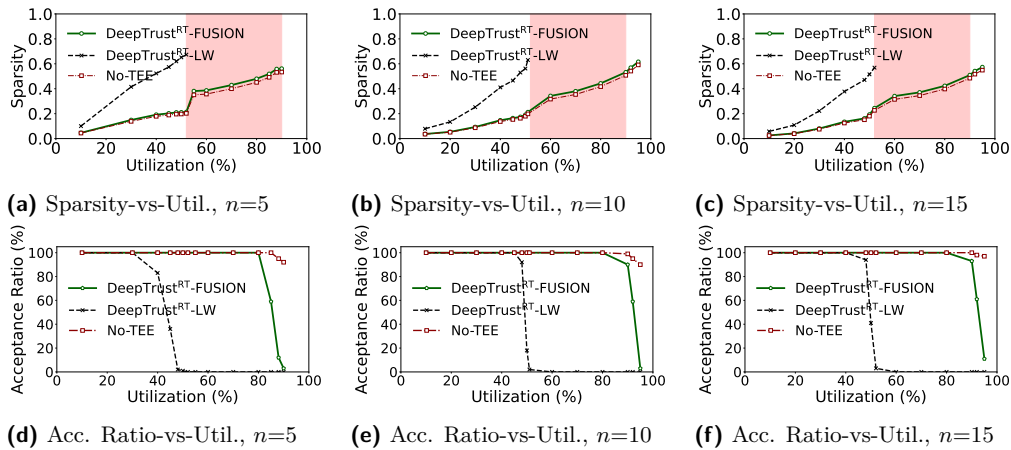
- 511 ■ **DeepTrust^{RT}-LW**: Sends the layers *sequentially* (layer-wise) to the enclave (Section 4.2).
- 512 ■ **DeepTrust^{RT}-FUSION**: Our proposed scheme that *fuses* multiple layers from multiple
513 tasks (Section 5).
- 514 ■ **No-TEE**: DNN task execution *without* any enclave. The tasks follow EDF scheduling
515 policy. In this case, model confidentiality is not enforced.

516 We tested the above schemes with the following two metrics.

- 517 ■ **Sparsity**: Our newly introduced metric that shows the “spread” of the task (viz., the
518 ratio between response time and period). Higher sparsity means tasks are completed late
519 and that may result in poorer QoS in terms of the DNN inference process. A Sparsity
520 value > 1 implies the task misses the deadline.
- 521 ■ **Acceptance Ratio**: A commonly used metric by the real-time community that represents
522 the fraction of tasks that meet deadlines over the total generated ones.



■ **Figure 5** Sparsity and Acceptance Ratio with varying system utilization for $\{5, 10, 15\}$ tasks using AlexNet-squeezed [18] architecture. The red shaded regions show cases where DeepTrust^{RT}-LW cannot find schedulable tasksets while other schemes can. DeepTrust^{RT}-FUSION result in better schedulability compared to DeepTrust^{RT}-LW as the utilization increases with performance penalty (i.e., both Sparsity and Acceptance Ratio are close to the No-TEE case).

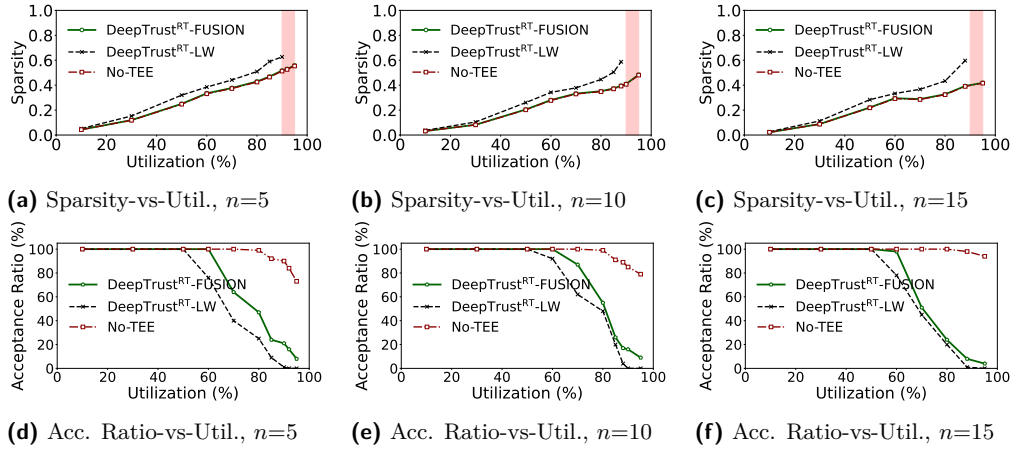


■ **Figure 6** Sparsity and Acceptance Ratio using Tiny Darknet [18] architecture using a setup identical to that of Fig. 5. The findings are similar.

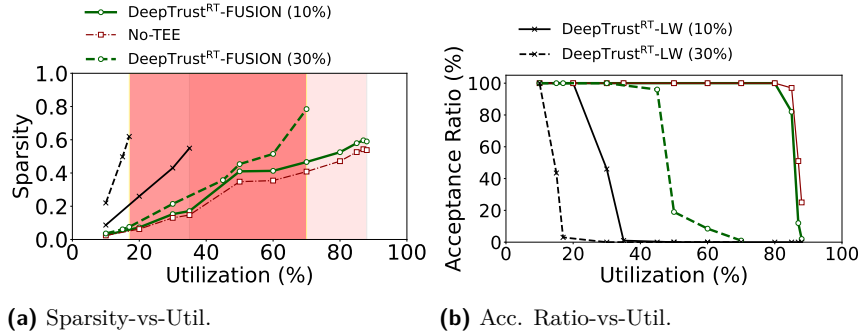
523 6.1.3 Results

524 We first show the Sparsity and Acceptance Ratio for varying numbers of tasks ($n = 5$, $n = 10$,
 525 and $n = 15$) for the DNN workloads listed on Table 5. The x-axis of Fig. 5 shows the
 526 various taskset utilization for randomly generated taskset. The y-axis of Fig. 5a and Fig. 5d
 527 shows Sparsity and Acceptance Ratio, respectively. We show the Sparsity and Acceptance
 528 Ratio for DeepTrust^{RT}-FUSION (Green), DeepTrust^{RT}-LW (Black), and No-TEE (Black)
 529 schemes. The red shaded regions in the figure represent the cases where DeepTrust^{RT}-LW is
 530 unable to find any schedulable candidate while DeepTrust^{RT}-FUSION finds some. For lower
 531 utilization, all schemes show similar behavior. However, DeepTrust^{RT}-FUSION outperforms
 532 DeepTrust^{RT}-LW up to 3x as the utilization increases (i.e., DeepTrust^{RT}-LW is unable to
 533 find schedulable tasksets as the utilization reaches 60%). This is expected because layer-wise

10:18 DeepTrust^{RT}: Confidential Deep Neural Inference Meets Real-Time!



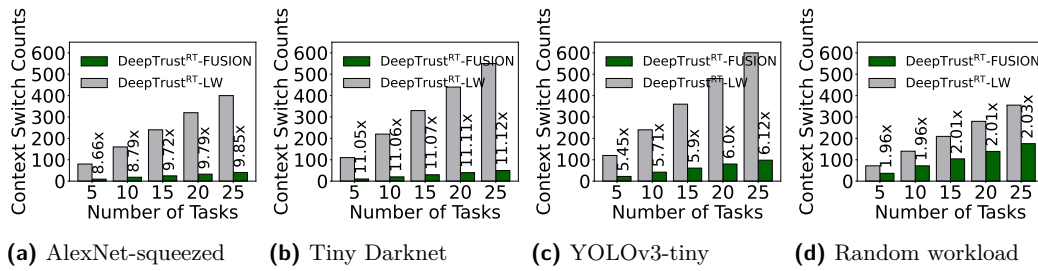
■ **Figure 7** Sparsity and Acceptance Ratio using YOLOv3-tiny [19] architecture using a setup identical to that of Fig. 5. Our findings are similar to Fig. 5 and Fig. 6.



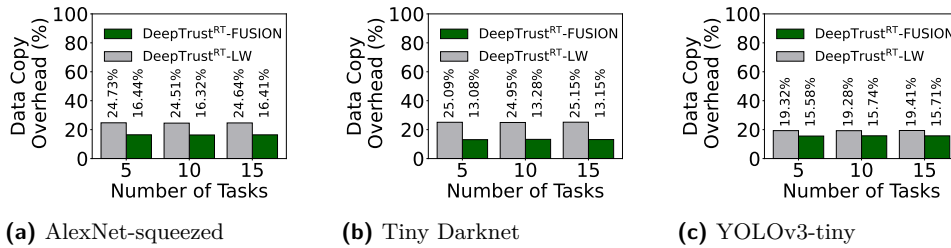
■ **Figure 8** Sparsity and Acceptance Ratio for a randomly generated workload with two different context switching overheads: (a) 10% of max(WCET) values (solid lines) and (b) 30% of max(WCET) values (dotted lines). Larger context switch delays result in higher Sparsity for DeepTrust^{RT}-LW when compared to DeepTrust^{RT}-FUSION, which in turn, reduces the percentage of schedulable tasksets. DeepTrust^{RT}-FUSION performs identically to No-TEE for lower context switch delays.

534 execution in DeepTrust^{RT}-LW increases delay due to additional context switches. At higher
 535 utilization, that causes more tasks to miss deadlines and results in lower acceptance. We also
 536 note that the performance of our scheme (both in terms of Sparsity and Acceptance Ratio)
 537 is close to No-TEE case (recall: No-TEE does not provide model confidentiality). Hence,
 538 DeepTrust^{RT}-FUSION can improve the security posture of the DNN tasks without significant
 539 overhead (close to the vanilla execution that does not have TEE support). In Figs. 6-7, we
 540 repeat the experiments with ImageNet1k datasets and obtain similar results. As the number
 541 of tasks increases (i.e., $n = 15$), the impact of context switching becomes more apparent,
 542 and hence, Acceptance Ratio in No-TEE case significantly outperforms DeepTrust^{RT}-LW
 543 and DeepTrust^{RT}-FUSION in highly utilized systems.

544 To further analyze the effect of context switches on Sparsity and Acceptance Ratio, we
 545 vary the SMC overheads as a percentage of WCET. Let $\max(WCET)$ denote the maximum
 546 WCET value observed in our experiments. The solid lines in Fig. 8 show the context switch
 547 cost as 10% of $\max(WCET)$ values of all tasks while dotted lines are generated with SMC
 548 overheads with 30% of $\max(WCET)$. As the figures show, the effect of larger context switch



■ **Figure 9** Context switches overhead comparison for three architectures (e.g., Tiny Darknet, AlexNet-squeezed, YOLOv3-tiny) and one for a random taskset. DeepTrust^{RT}-FUSION reduces context switch overhead 1.96x-11.12x compared to DeepTrust^{RT}-LW.



■ **Figure 10** DeepTrust^{RT} data copy overheads for various DNN workloads. Inference confidentiality increases response times due to additional data transfers and SMC calls. However, this overhead remains constant with the increasing number of tasks. The overheads of DeepTrust^{RT}-FUSION are less compared to DeepTrust^{RT}-LW due to fewer number of context switches.

549 costs causes DeepTrust^{RT}-LW to perform poorly as delays accumulating by higher context
 550 switch duration lead to longer response times (higher Sparsity values), that in turn, cause
 551 more tasks to miss their deadlines (result in lower Acceptance Ratio).

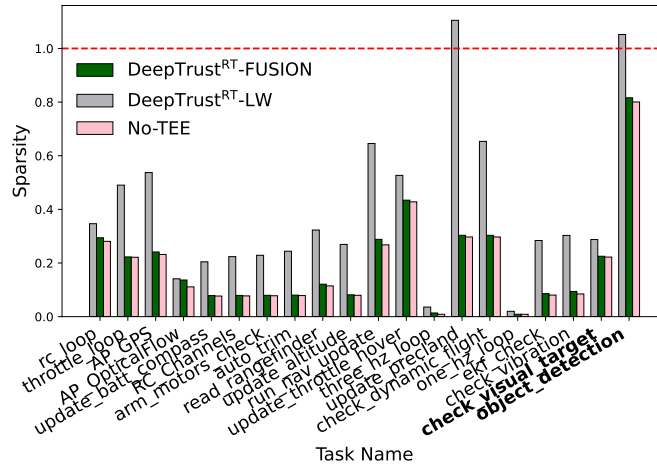
552 *DeepTrust^{RT}-FUSION outperforms DeepTrust^{RT}-LW, especially for high utilization scenarios. Further, the overhead of DeepTrust^{RT}-FUSION is negligible as its performance is close to the No-TEE case. Systems with longer TEE context switch delay can be significantly benefited by layer fusion compared to layer-wise partitioning.*

553 In the next set of experiments (Fig. 9), we measure the number of SMC context switches
 554 for DeepTrust^{RT}-FUSION and DeepTrust^{RT}-LW. For this experiment, we set the system
 555 utilization 50%. Note that, as No-TEE does not have any enclave, there are no SMC calls
 556 (context switches). Hence, our plots exclude No-TEE in this case. As the figures show,
 557 DeepTrust^{RT}-FUSION enables us to achieve a significant reduction in context switch counts
 558 compared to DeepTrust^{RT}-LW (5.45x-11.1x) for all three architectures. This is because
 559 DeepTrust^{RT}-FUSION groups multiple layers, hence reducing overall SMC calls.

560 *DeepTrust^{RT}-FUSION can significantly reduce the number context switches (1.9x-11.1x) (see Fig. 9). This reduction of context switches also contribute to higher acceptance rate (see Figs. 5-7).*

561 Recall from Section 3.1 that each time a context switch is performed, normal world
 562 (encrypted) data needs to be transferred to the secure world. We now analyze this data copy
 563 overhead. The experiments in Fig. 10 show the overheads for the various DNN workloads

10:20 DeepTrust^{RT}: Confidential Deep Neural Inference Meets Real-Time!



■ **Figure 11** Sparsity for ArduPilot controller tasks. Bold tasks are DNN inference workload. The red horizontal line denotes the deadline. The increasing number of context switches in DeepTrust^{RT}-LW caused a larger spread of tasks (higher Sparsity), and as a result, two tasks missed deadlines. Under DeepTrust^{RT}-FUSION, all tasks were able to meet deadlines.

564 (AlexNet-squeezed, Tiny Darknet, and YOLOv3-tiny) and a varying number of tasks ($n = 5$,
565 $n = 10$, and $n = 15$) running on Raspberry Pi and OP-TEE. To calculate the end-to-end data
566 copy overheads, we first measured the response times for No-TEE case and then subtracted
567 these values from the response times of each of the DeepTrust^{RT} schemes. Finally, we
568 normalized them with the task periods (i.e., calculated Sparsity) and obtained the overhead
569 percentage. We only considered schedulable tasksets. For each data point, we generated 100
570 samples and took the 90th percentile value. As the figure shows, enabling confidential inference
571 comes with a cost, i.e., increase in response times. This data copy overhead is system (i.e.,
572 underlying SMC implementations) and workload (i.e., DNN layers/architecture) dependent.
573 For instance, we find that the additional delay in response times due to transferring context
574 for DeepTrust^{RT}-LW and DeepTrust^{RT}-FUSION are (a) 2.39 s and 1.34 s (AlexNet-squeezed),
575 (b) 1.54 s and 2.95 s (Tiny Darknet), and (c) 5.62 s and 6.96 s (YOLOv3-tiny), respectively
576 on Raspberry Pi+OP-TEE setup (recall: each SMC overheads could be as high as 20 ms, see
577 Table 1). As the figure shows, the data copy overhead scales well with the increasing number
578 of tasks (remains constant). Further, DeepTrust^{RT}-FUSION incurs lower overheads due to a
579 reduced number of context switches, as we also observed in prior experiments (Fig. 9).

580 *Confidential deep inference comes with a cost: it increases response times due to additional data transfer between normal and secure worlds. However, this data transfer overhead does not increase significantly with the increasing number of tasks.*

581 6.2 Case Study with a UAV Controller System

582 In the final set of experiments (Fig. 11), we evaluate DeepTrust^{RT}-FUSION, DeepTrust^{RT}-LW,
583 and No-TEE with a UAV autopilot system (ArduPilot [20]) running on Raspberry Pi 3 [36].
584 The ArduPilot controller has 18 real-time tasks (defined in `/ArduCopter/Copter.cpp`).
585 Since the vanilla controller does not have any DNN workload, we included two additional
586 inference tasks (i.e., `check_visual_target()` and `object_detection()`) that use Tiny
587 Darknet and YOLOv3-tiny models, respectively to perform object detection. The periods of

our DNN tasks were 5. The total system utilization (including two of the included DNN tasks) was 0.75. Each of the bars in Fig. 11 shows the various tasks and their Sparsity for each of the three schemes. The figure shows that due to high context switches, DeepTrust^{RT}-LW misses deadlines for two real-time tasks (i.e., Sparsity > 1). In contrast, both DeepTrust^{RT}-FUSION and No-TEE were able to meet all deadlines.

High SMC context switch overheads cause DeepTrust^{RT}-LW to miss deadlines for two real-time tasks. DeepTrust^{RT}-FUSION, in contrast, was able to meet all deadlines.

7 Discussion

In this work, we assume all layers execute inside TEEs. There exist use cases where not all layers have confidentiality requirements. For example, in image/voice recognition applications where the user may not want to reveal input and processed data, running initial input layers and final output layers inside TEE should be sufficient. Our future work will explore the variable number of TEE executions and analyze the performance trade-offs in a real-time context. Our research focuses on scheduling within a single enclave as existing TrustZone implementations support a single enclave. We will further investigate the feasibility and performance benefits of DeepTrust^{RT} running on multiple enclaves.

DeepTrust^{RT}-FUSION currently selects a whole slice of a layer and fuses it with another task. For example, consider τ_i has four layers $\{l_i^{11}, \dots, l_i^{14}\}$ and τ_j has three layers $\{l_j^{11}, \dots, l_j^{13}\}$. If feasible (i.e., enclave has capacity), DeepTrust^{RT}-FUSION fuses all seven layers. It could also be possible to obtain a “partial” slice of a layer in case a complete slice is not fit in the enclave (or the enclave has a little extra capacity). For instance, in the example above, $\{l_i^{11}, l_i^{12}, l_j^{12}, l_j^{13}\}$ could form a fusion group in case all seven layers do not fit to further improve schedulability. Sub-layer-based partitioning ideas will be explored in our future work.

We assume only inference tasks use TEEs. In practice, other (non-DNN) tasks could also use TEEs, thus potentially limiting enclave availability. DeepTrust^{RT}-FUSION can be extended for such scenarios considering extra *slack* reclaimed from other non-inference tasks. The overall security of DeepTrust^{RT} relies on the underlying TEE architecture. However, TEEs could also be vulnerable, especially exposed to schedule-based attacks [37] for real-time context. One approach to limit such observability is to introduce “noise” in the scheduler [38]. For instance, instead of fusing the same set of tasks, DeepTrust^{RT}-FUSION can be extended to select fusion candidates from different groups, thus limiting the predictability and hence reducing the chances of information leakage.

8 Related work

Research in confidential deep inference for real-time context is still in the early stages. In our preliminary (workshop) paper [39], we propose a layer-grouping idea for fixed priority systems. Unlike DeepTrust^{RT}, our prior work does not provide formal timing guarantees. AegisDNN [23] proposes to execute only a few layers that will be executed inside SGX-based TEEs. However, AegisDNN [23] is primarily designed for soft real-time systems and allows deadline misses. In contrast, DeepTrust^{RT} is designed for hard real-time systems.

There exists other work for general-purpose systems. DarknetZ [21] proposes to execute only a few layers that will be executed inside TEE, which is not suitable for applications that require executing all layers within TEE. Layers that execute outside of the secure

10:22 DeepTrust^{RT}: Confidential Deep Neural Inference Meets Real-Time!

630 world expose information to the untrusted normal world, raising data privacy concerns. A
631 similar line of work exists (e.g., HybridTEE [40], Confidential DL [13], Occlumency [41]),
632 for executing machine learning workloads inside TEEs. However, none of them consider
633 real-time constraints.

634 SuperTEE [32] aims to reduce TEE task switching overhead. However, SuperTEE [32]
635 is not designed for learning-enabled real-time systems and can not be directly adapted for
636 DNN workloads. Researchers also propose techniques (e.g., Subflow [35], AppNet [42]) to
637 make deep learning “time-aware,” but they do not consider trusted execution aspects. The
638 proposed research is one of the fundamental works that investigates time-aware confidential
639 deep learning techniques for hard autonomous systems.

9 Conclusion

641 This research introduces techniques to enable real-time guarantees for confidential deep
642 learning using trusted enclaves. We show how to slice a partition in a large deep-learning
643 model to schedule using real-time schedulers such as EDF. We further propose an optimization
644 using a novel idea of “layer fusion” that selectively groups multiple layers from various tasks
645 to minimize TEE context switch overheads. By using the approach presented in this work,
646 engineers of future autonomous systems will be able to design/schedule systems efficiently
647 and measure overheads of deep neural inference workloads in a “quantifiable” way.

References

- 649 1 J. Lee, N. Chirkov, E. Ignasheva, Y. Pisarchyk, M. Shieh, F. Riccardi, R. Sarokin, A. Kulik,
650 and M. Grundmann, “On-device neural net inference with mobile GPUs,” *arXiv preprint*
651 *arXiv:1907.01989*, 2019.
- 652 2 R. Shokri, M. Stronati, C. Song, and V. Shmatikov, “Membership inference attacks against
653 machine learning models,” in *2017 IEEE symposium on security and privacy (SP)*, pp. 3–18,
654 IEEE, 2017.
- 655 3 N. Li, W. Qardaji, D. Su, Y. Wu, and W. Yang, “Membership privacy: A unifying framework
656 for privacy definitions,” in *Proceedings of the 2013 ACM SIGSAC conference on Computer &*
657 *communications security*, pp. 889–900, 2013.
- 658 4 Z. Chen, G. Li, K. Pattabiraman, and N. DeBardeleben, “BinFI: An efficient fault injector for
659 safety-critical machine learning systems,” in *Proceedings of the International Conference for*
660 *High Performance Computing, Networking, Storage and Analysis*, pp. 1–23, 2019.
- 661 5 G. Li, K. Pattabiraman, and N. DeBardeleben, “TensorFI: A configurable fault injector for
662 TensorFlow applications,” in *IEEE International symposium on software reliability engineering*
663 *workshops (ISSREW)*, pp. 313–320, IEEE, 2018.
- 664 6 H. Benkraouda and K. Nahrstedt, “Image reconstruction attacks on distributed machine
665 learning models,” in *Proceedings of the 2nd ACM International Workshop on Distributed*
666 *Machine Learning*, pp. 29–35, 2021.
- 667 7 “Intel software guard extensions: Intel SGX SDK for Linux OS. <http://intel.com>. accessed:
668 2020-06-30.”
- 669 8 T. Alves, “TrustZone: Integrated hardware and software security,” *Information Quarterly*,
670 vol. 3, pp. 18–24, 2004.
- 671 9 K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image
672 recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- 673 10 Linaro, “Open portable trusted execution environment.” <https://www.op-tee.org>, Accessed
674 on 2021.
- 675 11 M. F. Babar and M. Hasan, “Trusted deep neural execution—a survey,” *IEEE Access*, 2023.

- 676 12 J. Redmon, “Darknet: Open source neural networks in C.” <http://pjreddie.com/darknet/>,
677 2013–2016.
- 678 13 P. M. VanNostrand, I. Kyriazis, M. Cheng, T. Guo, and R. J. Walls, “Confidential deep learning:
679 Executing proprietary models on untrusted devices,” *arXiv preprint arXiv:1908.10730*, 2019.
- 680 14 R. Liu, L. Garcia, Z. Liu, B. Ou, and M. Srivastava, “SecDeep: Secure and performant
681 on-device deep learning inference framework for mobile and IoT devices,” in *Proceedings of the
682 International Conference on Internet-of-Things Design and Implementation*, pp. 67–79, 2021.
- 683 15 T. Elgamal and K. Nahrstedt, “Serdab: An IoT framework for partitioning neural networks
684 computation across multiple enclaves,” in *20th IEEE/ACM International Symposium on
685 Cluster, Cloud and Internet Computing (CCGRID)*, pp. 519–528, IEEE, 2020.
- 686 16 M. F. Babar and M. Hasan, “Real-time scheduling of TrustZone-enabled DNN workloads,” in
687 *Proceedings of the 4th Workshop on CPS & IoT Security and Privacy*, pp. 63–69, 2022.
- 688 17 S. Han, H. Mao, and W. J. Dally, “Deep compression: Compressing deep neural networks with
689 pruning, trained quantization and Huffman coding,” *arXiv preprint arXiv:1510.00149*, 2015.
- 690 18 F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, “SqueezeNet:
691 AlexNet-level accuracy with 50x fewer parameters and < 0.5 mb model size,” *arXiv preprint
692 arXiv:1602.07360*, 2016.
- 693 19 P. Adarsh, P. Rathi, and M. Kumar, “Yolo v3-tiny: Object detection and recognition using
694 one stage improved model,” in *2020 6th international conference on advanced computing and
695 communication systems (ICACCS)*, pp. 687–694, IEEE, 2020.
- 696 20 “<https://github.com/ardupilot/ardupilot>.”
- 697 21 F. Mo, A. S. Shamsabadi, K. Katevas, S. Demetriou, I. Leontiadis, A. Cavallaro, and
698 H. Haddadi, “DarkneTZ: towards model privacy at the edge using trusted execution
699 environments,” in *Proceedings of the 18th International Conference on Mobile Systems,
700 Applications, and Services*, pp. 161–174, 2020.
- 701 22 J. A. Stankovic, M. Spuri, K. Ramamritham, and G. Buttazzo, *Deadline scheduling for
702 real-time systems: EDF and related algorithms*, vol. 460. Springer Science & Business Media,
703 1998.
- 704 23 Y. Xiang, Y. Wang, H. Choi, M. Karimi, and H. Kim, “AegisDNN: Dependable and timely
705 execution of DNN tasks with SGX,” in *IEEE Real-Time Systems Symposium (RTSS)*, pp. 68–81,
706 IEEE, 2021.
- 707 24 K. Kim, C. H. Kim, J. J. Rhee, X. Yu, H. Chen, D. Tian, and B. Lee, “Vessels: Efficient
708 and scalable deep learning prediction on trusted processors,” in *Proceedings of the 11th ACM
709 Symposium on Cloud Computing*, pp. 462–476, 2020.
- 710 25 A. Moffat, “Huffman coding,” *ACM Computing Surveys (CSUR)*, vol. 52, no. 4, pp. 1–35,
711 2019.
- 712 26 S. Han, J. Pool, J. Tran, and W. Dally, “Learning both weights and connections for efficient
713 neural network,” *Advances in neural information processing systems*, vol. 28, 2015.
- 714 27 M. S. Islam, M. Zamani, C. H. Kim, L. Khan, and K. W. Hamlen, “Confidential execution of
715 deep learning inference at the untrusted edge with ARM TrustZone,” in *Proceedings of the
716 Thirteenth ACM Conference on Data and Application Security and Privacy, CODASPY’23,
717 (New York, NY, USA)*, p. 153–164, Association for Computing Machinery, 2023.
- 718 28 F. Zhang and A. Burns, “Schedulability analysis for real-time systems with EDF scheduling,”
719 *IEEE Transactions on Computers*, vol. 58, no. 9, pp. 1250–1258, 2009.
- 720 29 C. L. Liu and J. W. Layland, “Scheduling algorithms for multiprogramming in a hard-real-time
721 environment,” *Journal of the ACM (JACM)*, vol. 20, no. 1, pp. 46–61, 1973.
- 722 30 A. Singh and S. Baruah, “Global EDF-based scheduling of multiple independent synchronous
723 dataflow graphs,” in *2017 IEEE Real-Time Systems Symposium (RTSS)*, pp. 307–318, 2017.
- 724 31 R. I. Davis, “A review of fixed priority and EDF scheduling for hard real-time uniprocessor
725 systems,” *ACM SIGBED Review*, vol. 11, no. 1, pp. 8–19, 2014.

10:24 DeepTrust^{RT}: Confidential Deep Neural Inference Meets Real-Time!

- 726 **32** A. Mukherjee, T. Mishra, T. Chantem, N. Fisher, and R. Gerdes, “Optimized trusted execution
727 for hard real-time applications on COTS processors,” in *Proceedings of the 27th International
728 Conference on Real-Time Networks and Systems*, pp. 50–60, 2019.
- 729 **33** J.-J. Chen, “Partitioned multiprocessor fixed-priority scheduling of sporadic real-time tasks,”
730 in *2016 28th Euromicro Conference on Real-Time Systems (ECRTS)*, pp. 251–261, IEEE,
731 2016.
- 732 **34** A. Gholami, K. Kwon, B. Wu, Z. Tai, X. Yue, P. Jin, S. Zhao, and K. Keutzer, “SqueezeNext:
733 Hardware-aware neural network design,” in *Proceedings of the IEEE conference on computer
734 vision and pattern recognition workshops*, pp. 1638–1647, 2018.
- 735 **35** S. Lee and S. Nirjon, “SubFlow: A dynamic induced-subgraph strategy toward real-time dnn
736 inference and training,” in *2020 IEEE Real-Time and Embedded Technology and Applications
737 Symposium (RTAS)*, pp. 15–29, IEEE, 2020.
- 738 **36** M. Richardson and S. Wallace, *Getting started with Raspberry Pi*. O’Reilly Media, Inc., 2012.
- 739 **37** M. A. Aguida and M. Hasan, “Work in progress: Exploring schedule-based side-channels in
740 TrustZone-enabled real-time systems,” in *2022 IEEE 28th Real-Time and Embedded Technology
741 and Applications Symposium (RTAS)*, pp. 301–304, IEEE, 2022.
- 742 **38** C.-Y. Chen, D. Sanyal, and S. Mohan, “Indistinguishability prevents scheduler side channels
743 in real-time systems,” in *Proceedings of the 2021 ACM SIGSAC Conference on Computer and
744 Communications Security*, pp. 666–684, 2021.
- 745 **39** M. F. Babar and M. Hasan, “Real-time scheduling of Trustzone-enabled DNN workloads,” in
746 *Proceedings of the 4th Workshop on CPS & IoT Security and Privacy*, pp. 63–69, 2022.
- 747 **40** A. Gangal, M. Ye, and S. Wei, “HybridTEE: Secure mobile DNN execution using hybrid
748 trusted execution environment,” in *Asian Hardware Oriented Security and Trust Symposium
749 (AsianHOST)*, pp. 1–6, IEEE, 2020.
- 750 **41** T. Lee, Z. Lin, S. Pushp, C. Li, Y. Liu, Y. Lee, F. Xu, C. Xu, L. Zhang, and J. Song,
751 “Oclumency: Privacy-preserving remote deep-learning inference using SGX,” in *The 25th
752 Annual International Conference on Mobile Computing and Networking*, pp. 1–17, 2019.
- 753 **42** S. Bateni and C. Liu, “ApNet: Approximation-aware real-time neural network,” in *2018 IEEE
754 Real-Time Systems Symposium (RTSS)*, pp. 67–79, IEEE, 2018.