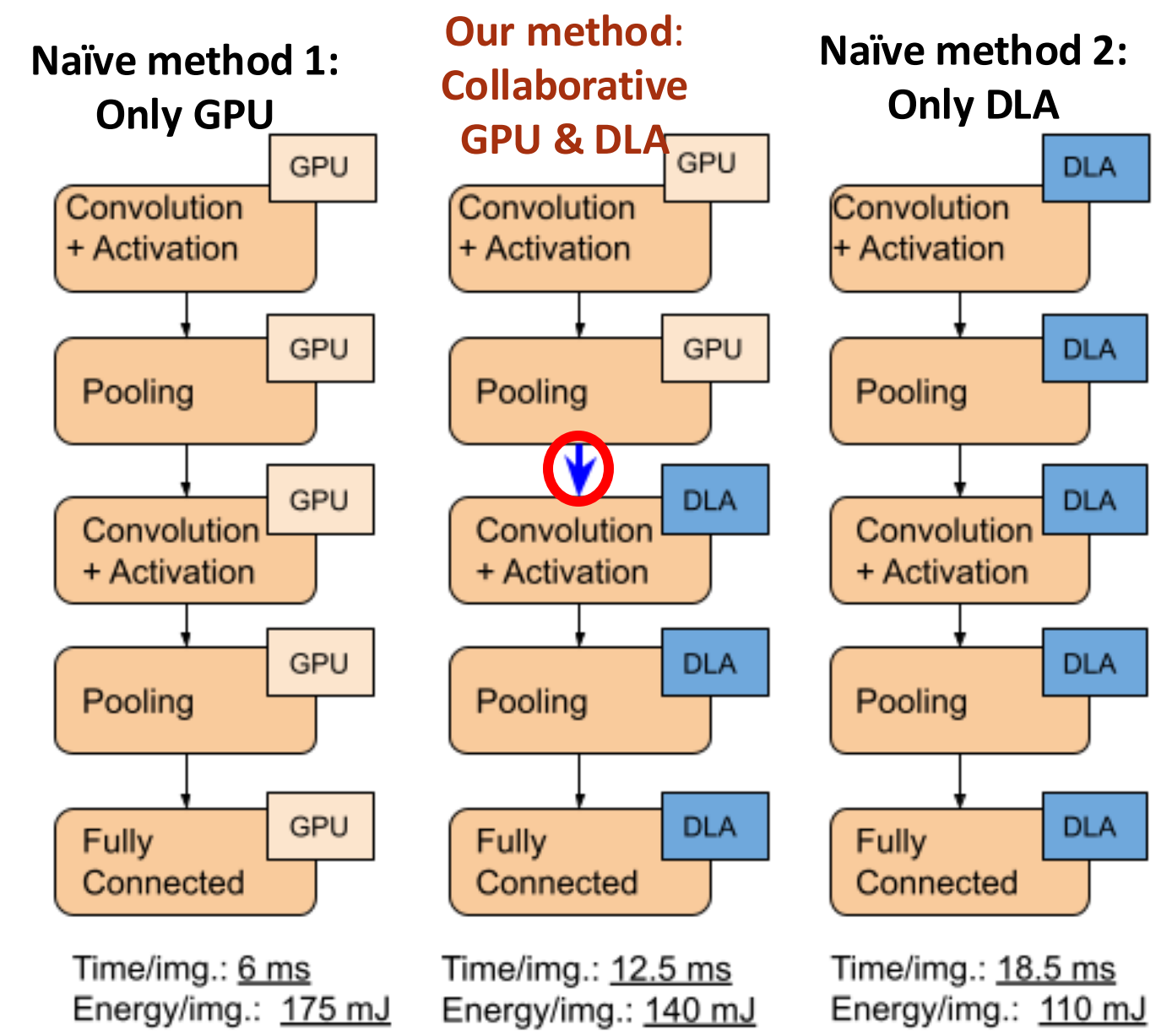




Energy-aware and Shared-memory-contention-aware Layer-wise Scheduling of DNN Inference

Optimizing Latency under Energy Constraint

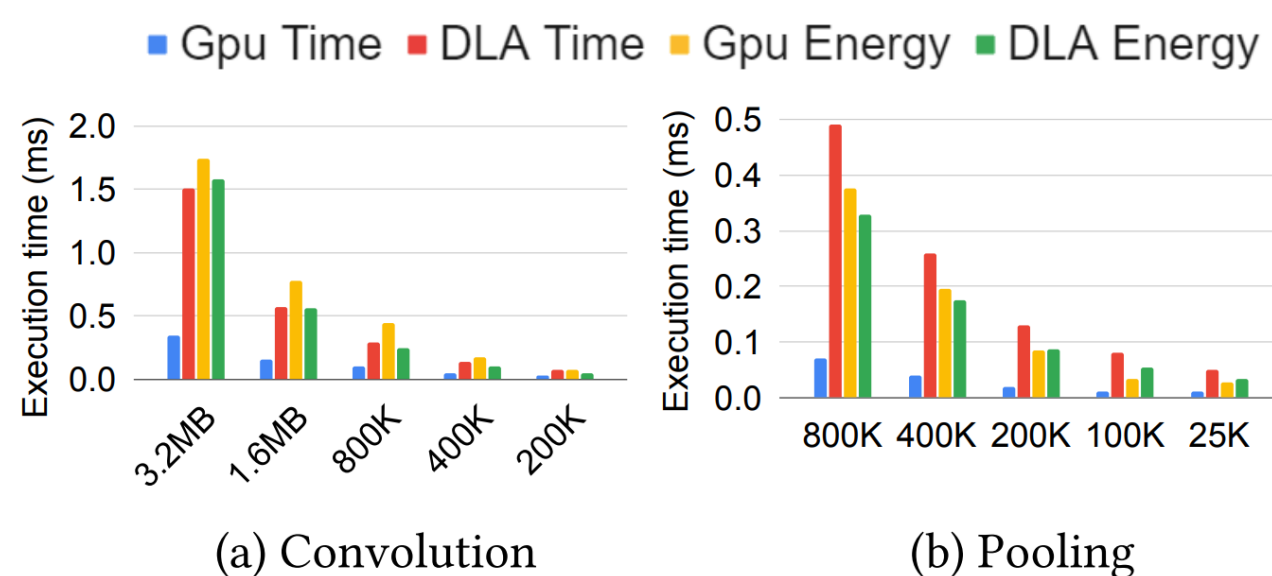


Target Hardware: Xavier AGX & Target Application: GoogleNet

Take-away: Collaborative multi-accelerator (GPU & DLA) execution of a single DNN inference in a single edge device can save significant time & energy.

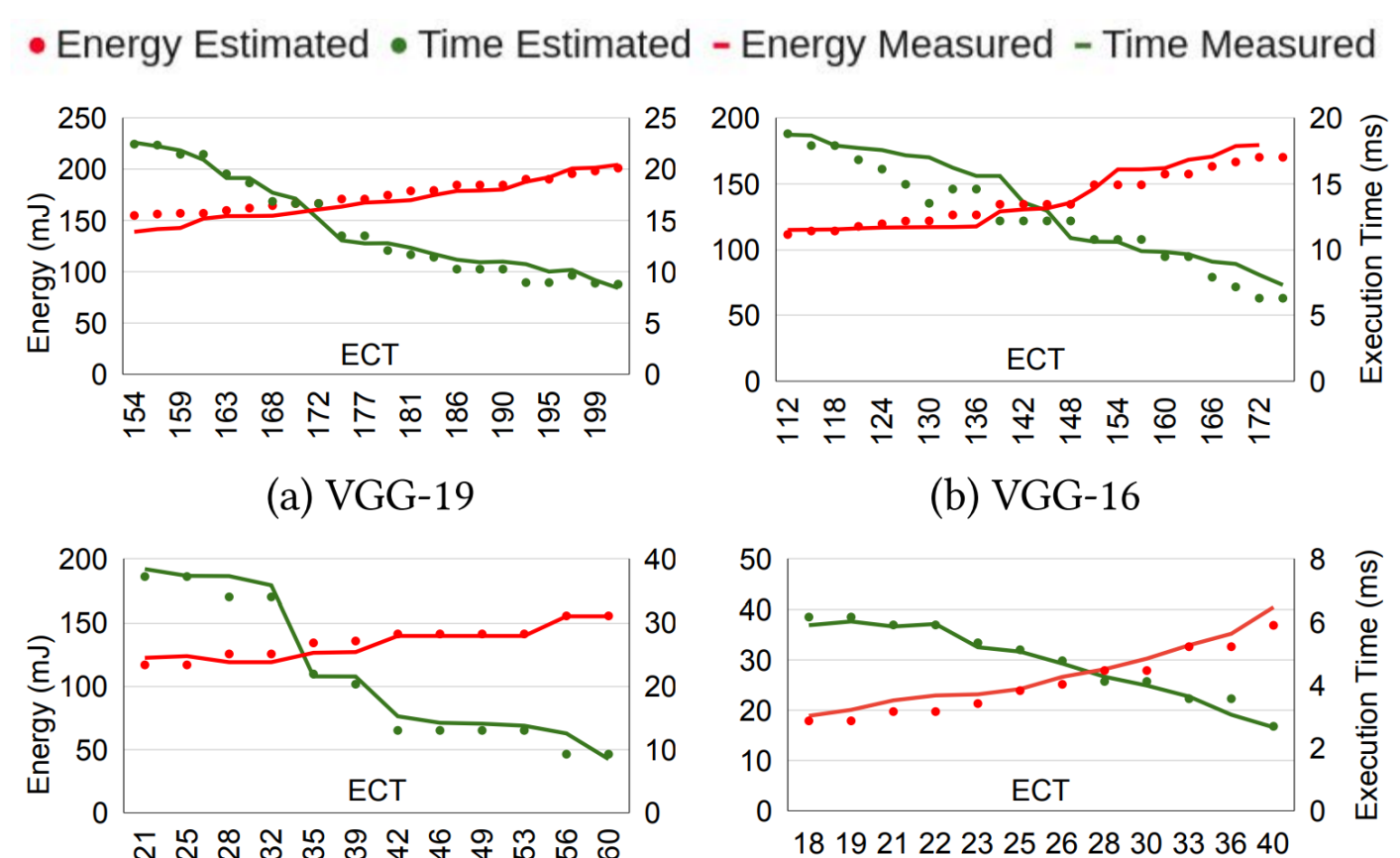
Observation-1: Getting benefits of Hardware Heterogeneity When Scheduling Tasks (layers)

GPU/DLA: Latency: 2x to 4.5x & Energy 1.1x to 2.4x

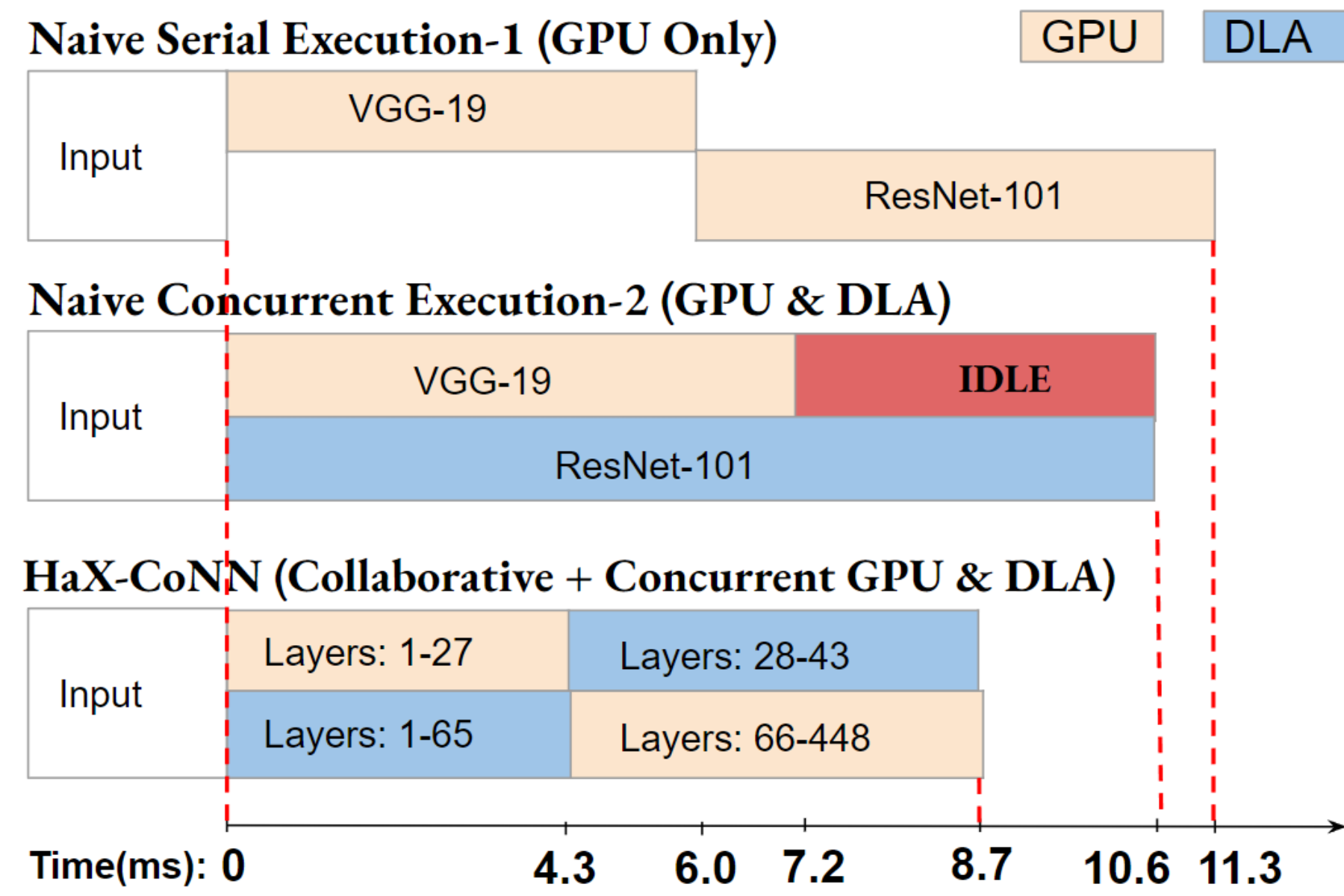


EVALUATION (Energy)

Comparison of measured results against estimation by our model



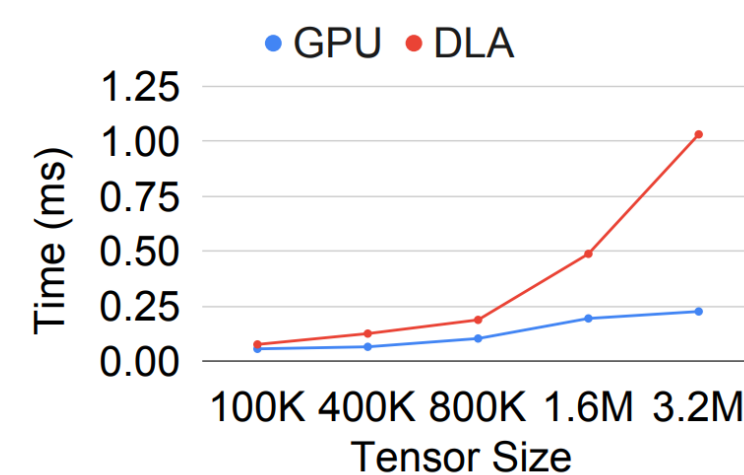
Significant Slowdown when DNNs are assigned to GPU & DLA



Target Hardware: Orin AGX, Target Application: VGG19 & ResNet101

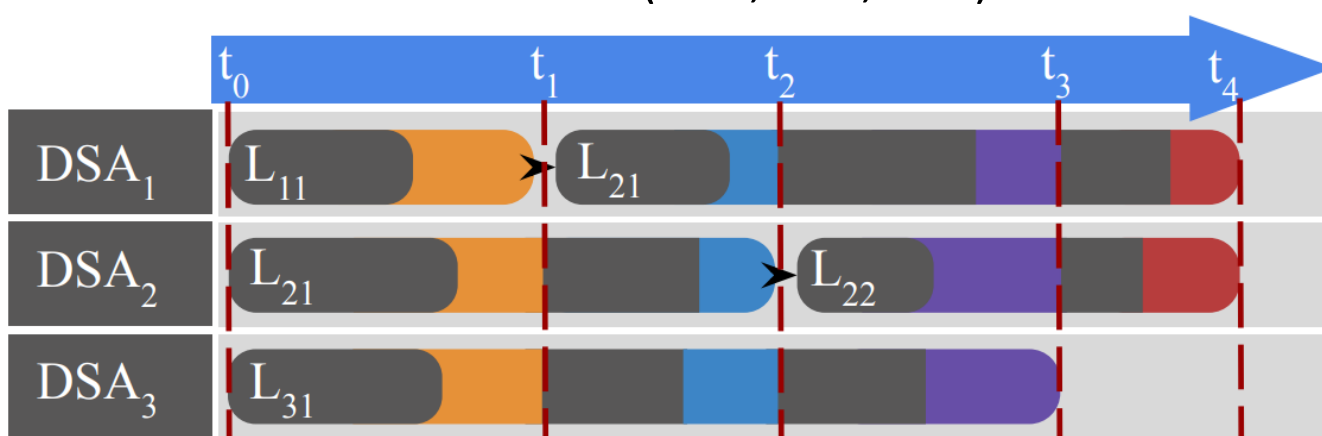
Take-away: Ensuring optimal concurrent execution requires careful consideration of shared memory contention.

Observation-2: Transition Cost between GPU and DLA (with shared memory)



Observation-3: Slowdown Depends on Memory Demand by Layers and the Accelerator

Different layers observe varying slowdown on different accelerators (GPU, DLA, PVA)



FORMULATION

We integrate layer execution time, inter-accelerator transition time, and memory contention slowdown into a cost function formulate the scheduling problem.

$$\text{Objective: } \min T(L, P, S(L \rightarrow P))$$

$$\text{s.t. } E(L, P, S(L \rightarrow P)) < ECT$$

$$\text{Total Time: } T(L, P, S(L \rightarrow P), TR) = \sum_{i=1}^{\text{len}(L)} (t(L_i, s(L_i)) + (TR_i \times \tau(L_i, s(L_i), \text{OUT}[IN]) + (TR_i \times \text{pipeline}(L_i, S(L_i))))$$

$$\text{Trans. Time: } TR_i = \begin{cases} 1 & \text{if } S(i) \neq S(i+1) \\ 0 & \text{if } S(i) = S(i+1) \end{cases}$$

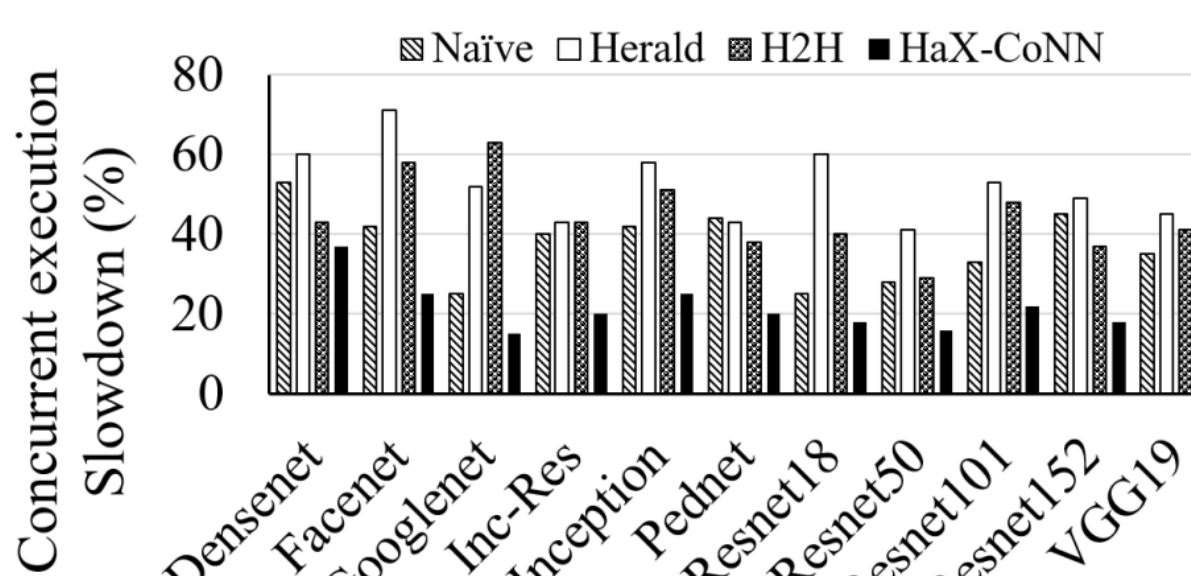
$$\text{Slowdown Amount: } C_{L_i, n, S(L_i), L} = \sum_{L_j \in L} \frac{t(L_i, n, S(L_i), L) \times \text{cont_model}(L_i, n, L_j)}{t(L_i, n, S(L_i), L) \times \text{len}(L_j)}$$

$$\text{where } 1 \leq j \leq \text{len}(DNN_n), 1 \leq n \leq \text{len}(DNN), L_j \in L_s$$

$$\text{Int}_k \cap [s_{L_i, n}, e_{L_i, n}] \neq \emptyset, \text{Int}_k \cap [s_{L_j, n}, e_{L_j, n}] \neq \emptyset$$

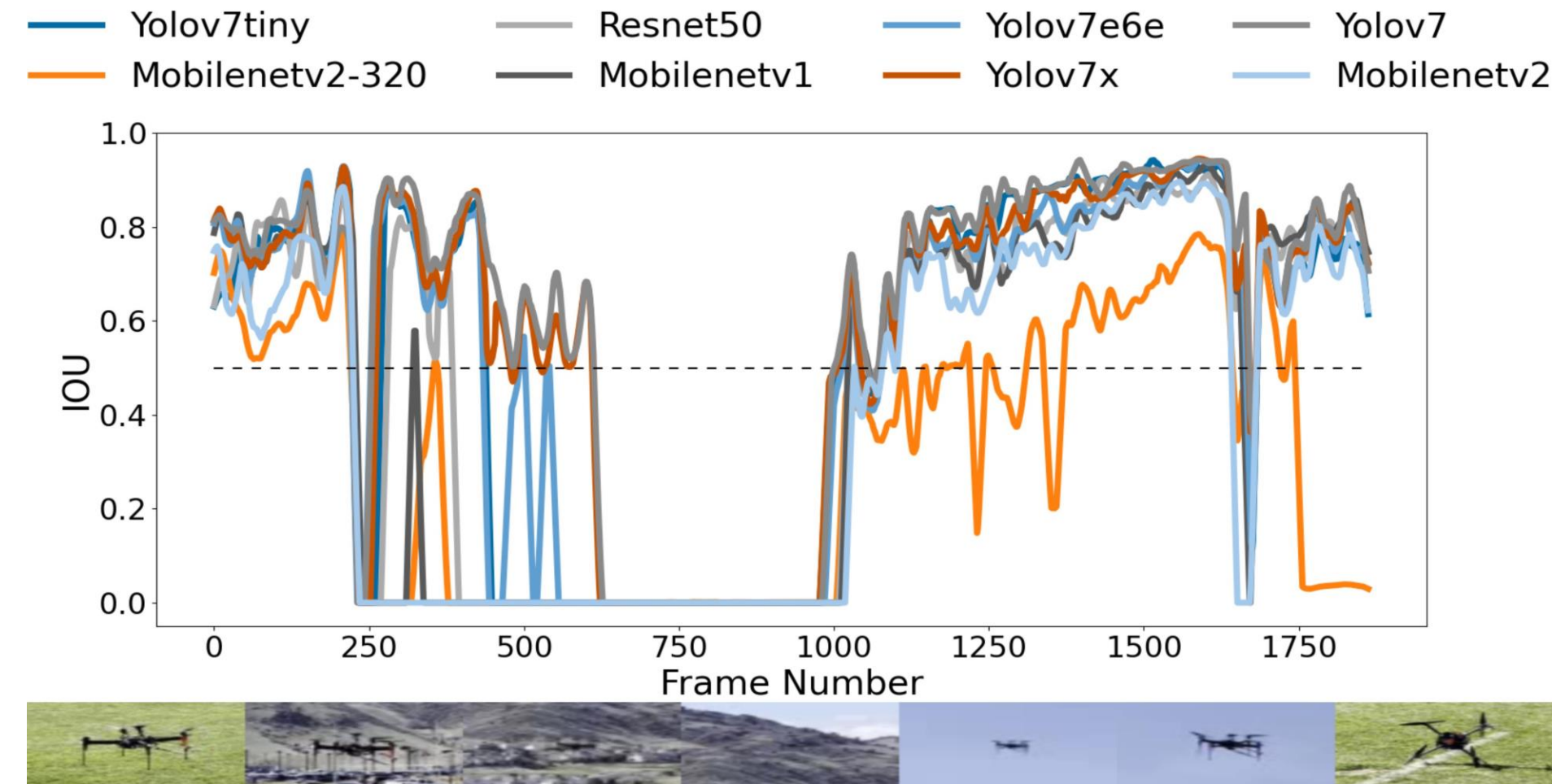
EVALUATION (Minimal Shared Memory Contention)

Slowdown amount per DNN when we run a DNN (x-axis) with other DNNs on the Orin AGX. HaX-CoNN (our work) minimizes the slowdown amount.



Optimizing Continuous Object Detection on Multi-Accelerator SoCs

Object Detection Model Accuracy



NO NEED FOR SOTA MODELS IN MOST SCENARIOS

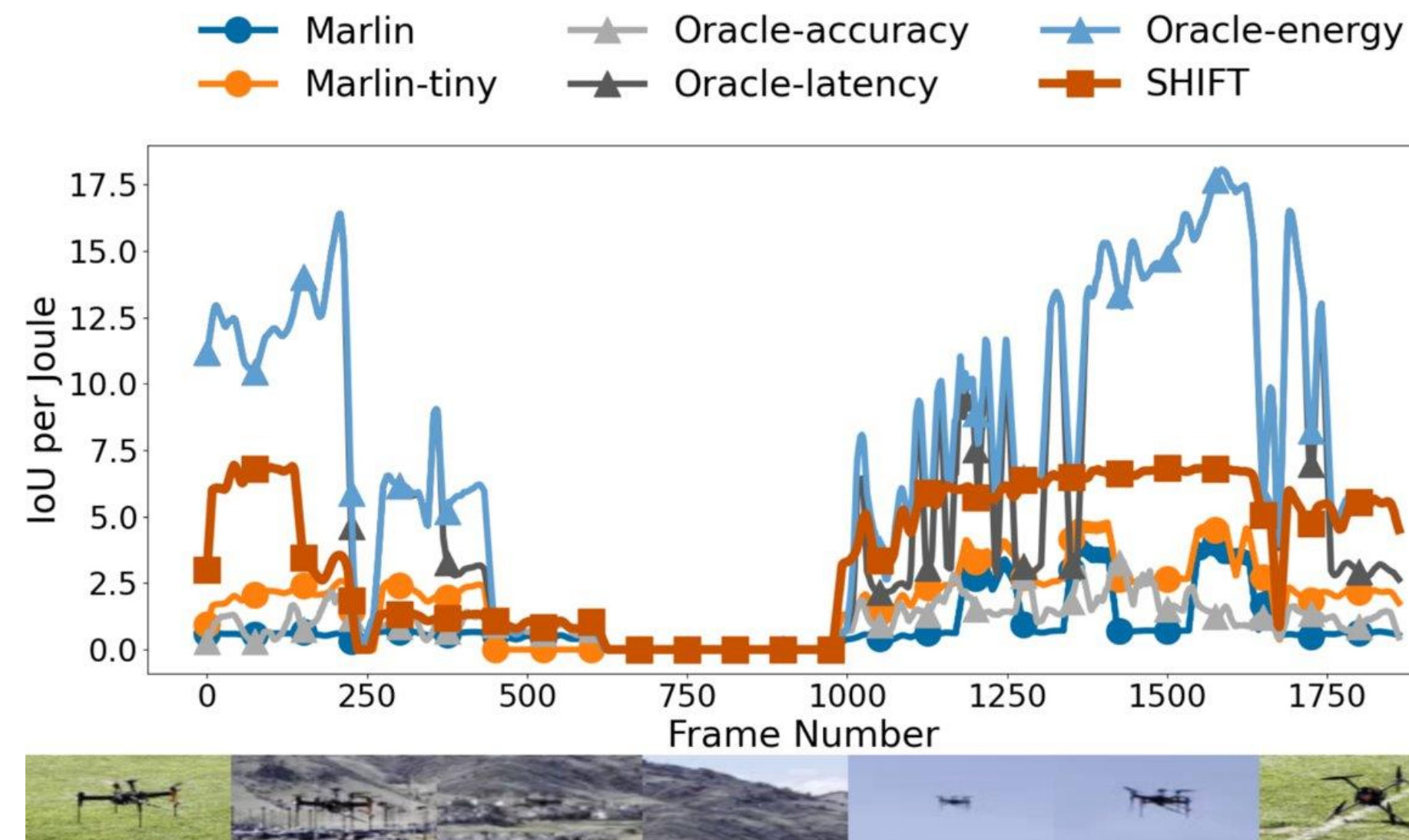
Take-away: Context-aware switching to non-SOTA object detection model architectures can improve energy/latency without significant accuracy loss



EXPLOIT NON-MONOTONIC RELATIONSHIPS BETWEEN MODEL ARCHITECTURES + HARDWARE

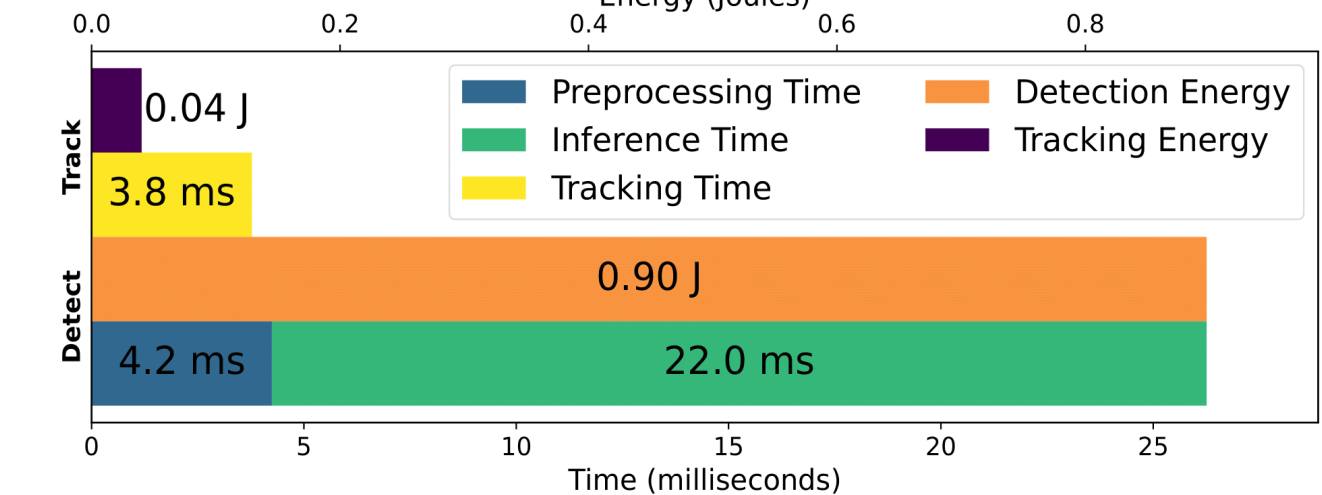
SOLUTION: Swap to models with lower computational cost if their predicted accuracy meets a threshold. Utilize a heuristic scheduler handling latency, energy, and accuracy.

EVALUATION: On Same Scenario



Target Hardware: Xavier NX + Luxonis OAK-D Lite

Latency of Object Detection vs. Tracking in MOT



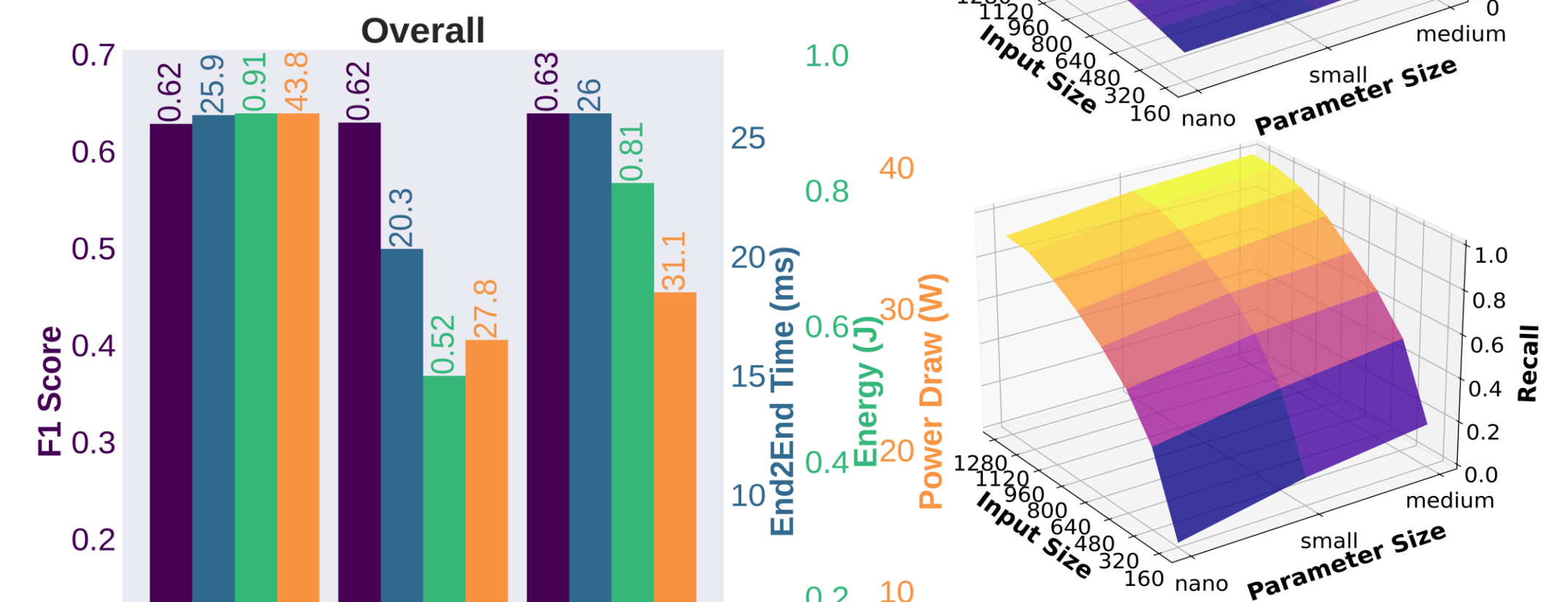
DETECTION IS THE BOTTLENECK FOR LATENCY AND ENERGY

Take-away: Splitting detection can reduce computational bottlenecks on edge SoCs



SOLUTION: Split frames at runtime and schedule models to process regions across hardware to minimize latency. Regions with low-priority objects are aggregated into a single region for efficient processing.

EVALUATION: Methods on MOT17 Dataset



INCREASING SINGLE MODEL ACCURACY HAS QUICKLY DIMINISHING RETURNS, MORE COMPLEX SOLUTION REQUIRED

RESULTS: SHIFT & PMOTHS

- SHIFT** has up to **7.5x** reduction in energy usage and **2.5x** improvement in latency
- PMOTHS** demonstrates up to **3.0x** reduction in latency, **6.5x** in energy, and **2.0x** in power draw on MOT17 dataset.