

Transparent Deployment of TFLite Workloads on Lightweight Many-Accelerator Architectures

Kuan-Lin Chiu, Guy Eichler, Chuan-Tung Lin
Giuseppe Di Guglielmo and Luca Carloni

June 29, 2024

Motivations & Goals

Optimizing Software for Accelerators in heterogeneous SoC architectures

- Modify software applications to offload specific computation kernels to accelerators
- Update the software algorithm when necessary
- Be able to run legacy software code

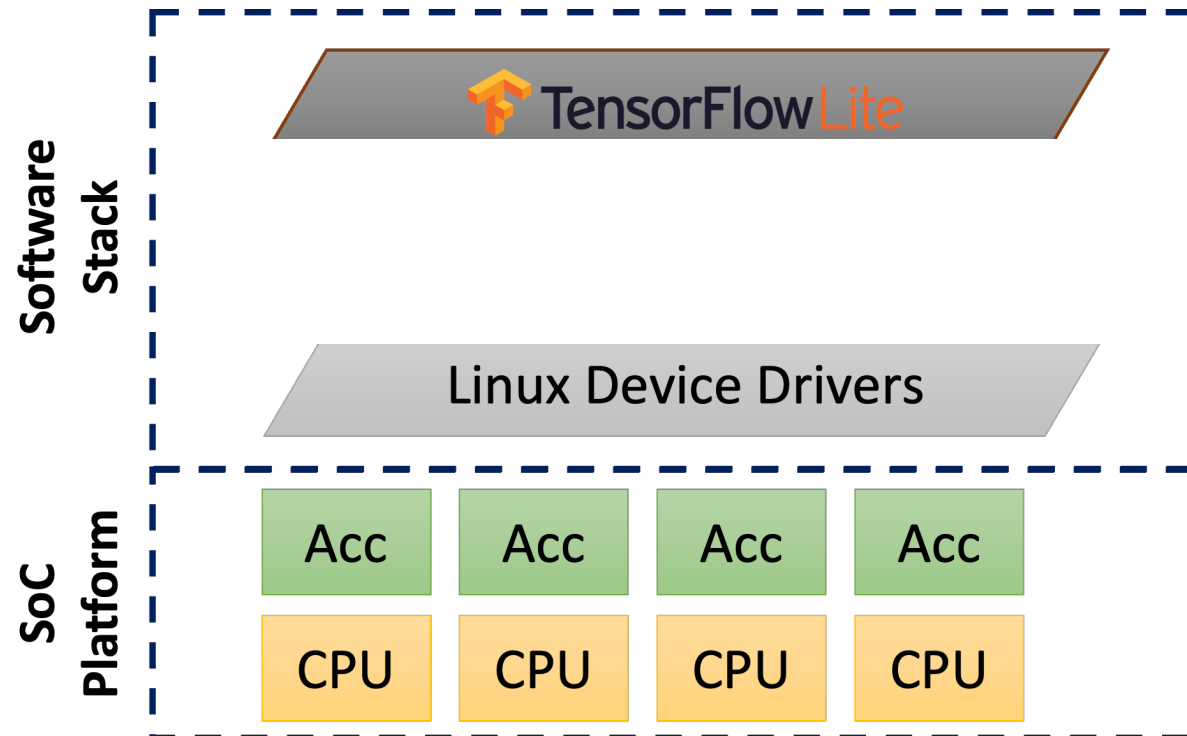
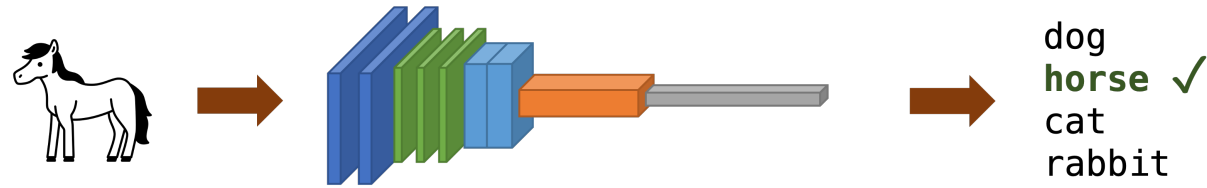
- Need: an **efficient** and **transparent** way to deploy software applications

Managing ML Applications Running in Parallel on the heterogeneous SoC

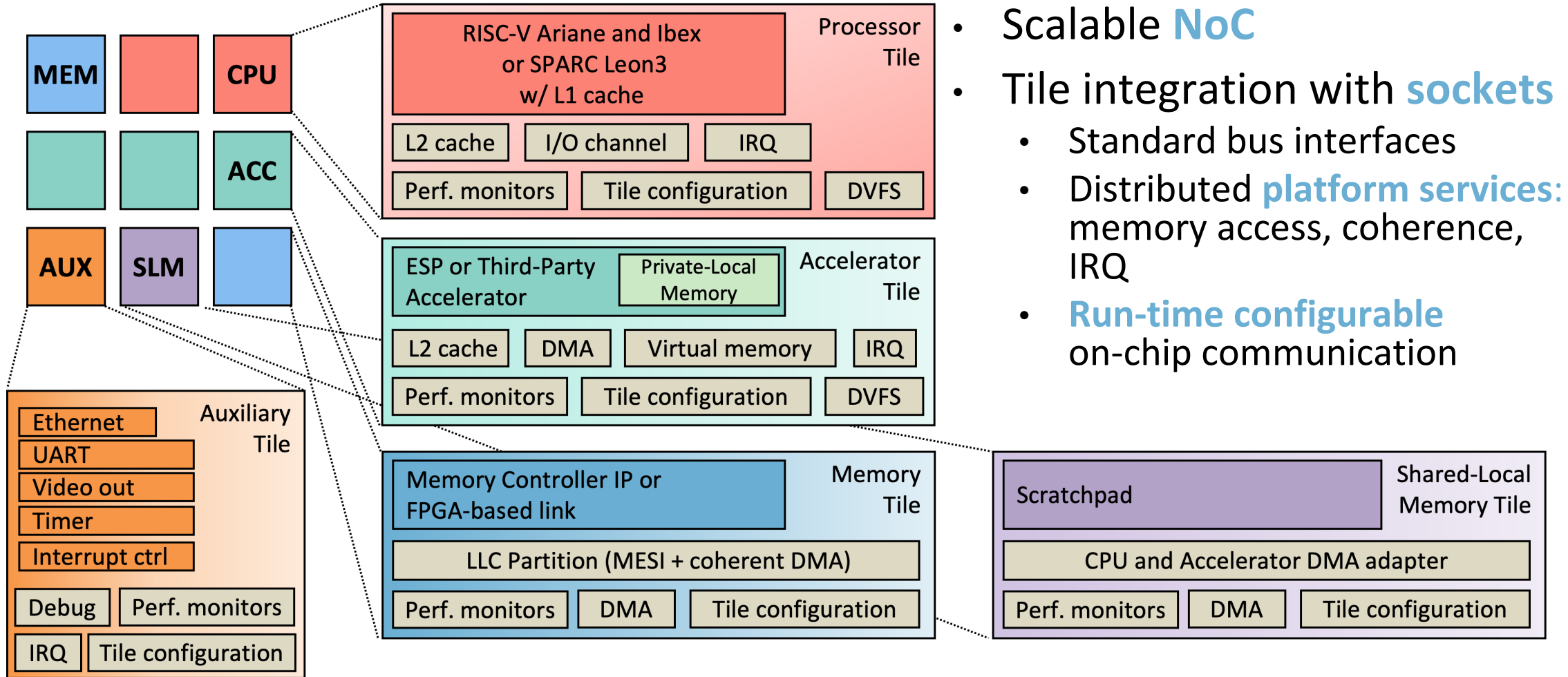
- Multi-tenant scenarios in real-world applications, ex: autonomous vehicles, drones, etc
- Inevitably, applications will compete in accessing accelerators

- Need: a mechanism to **automatically assign** available accelerators to applications

Wolt: SoC Deployment of Machine Learning Workloads

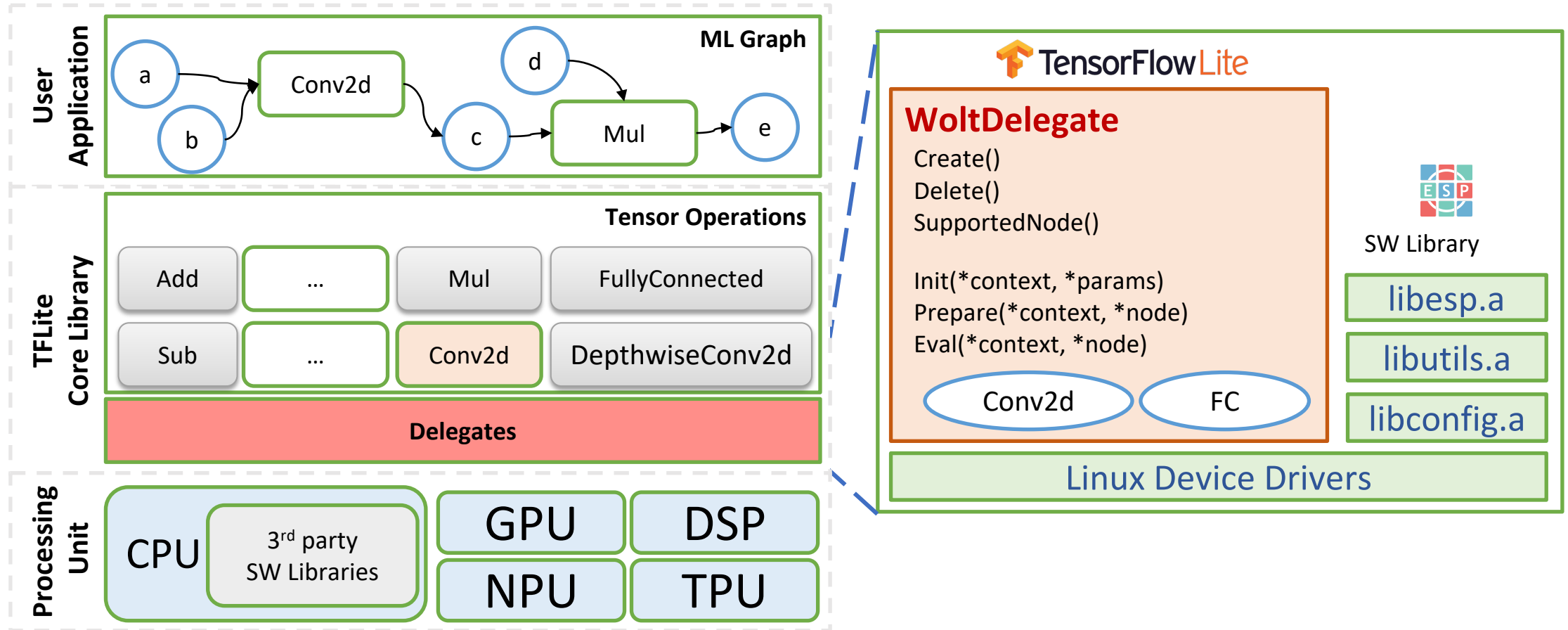


The ESP Tile-based SoC Architecture



- Scalable **NoC**
- Tile integration with **sockets**
 - Standard bus interfaces
 - Distributed **platform services**: memory access, coherence, IRQ
 - **Run-time configurable** on-chip communication

What is a Delegate? What is Wolt Delegate?



Init()

- Store indices of all delegated operations
- Initialize memories
- Initialize hardware buffer

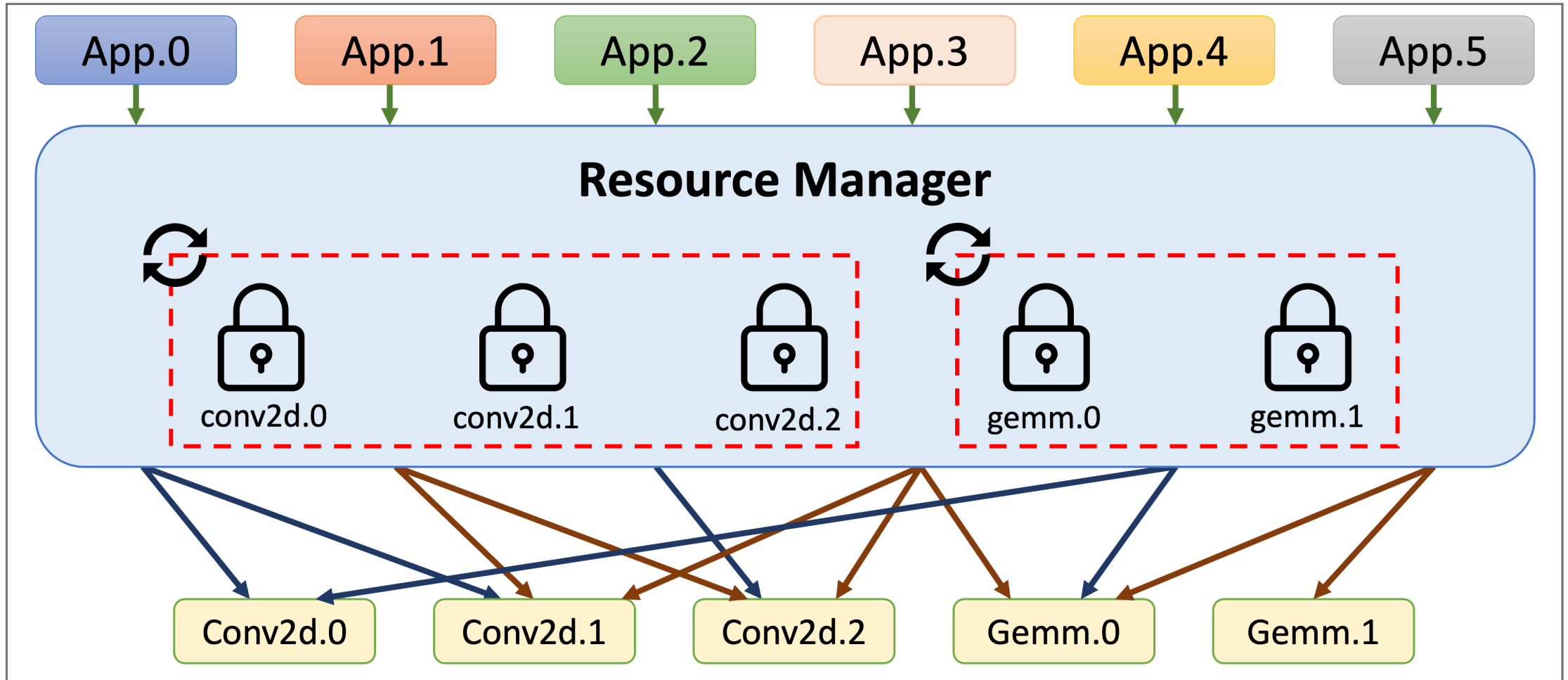
Prepare()

- Set parameters for each node (ex: dimensions of feature map and filter)

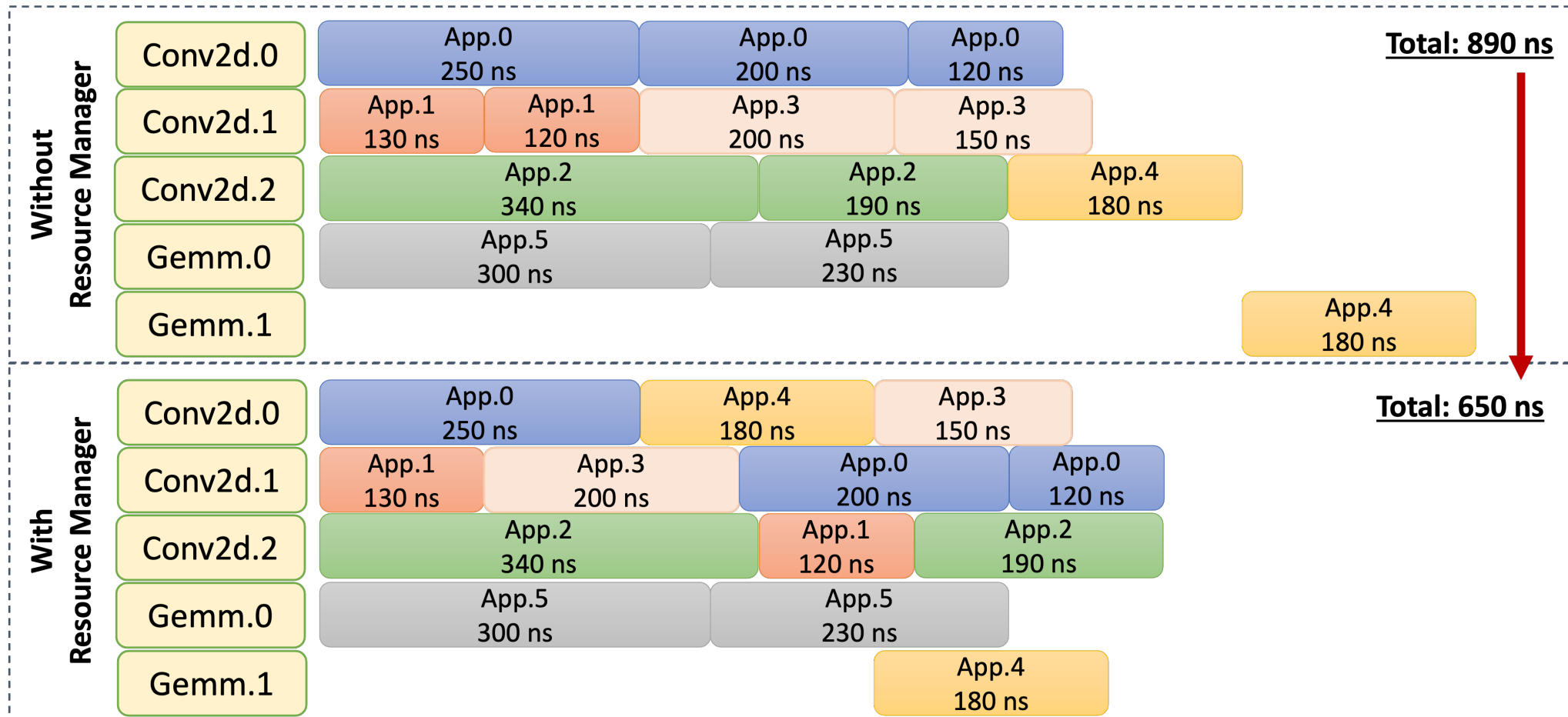
Eval()

- Execute delegated graph
- Invoke accelerators

The Wolt Resource Manager

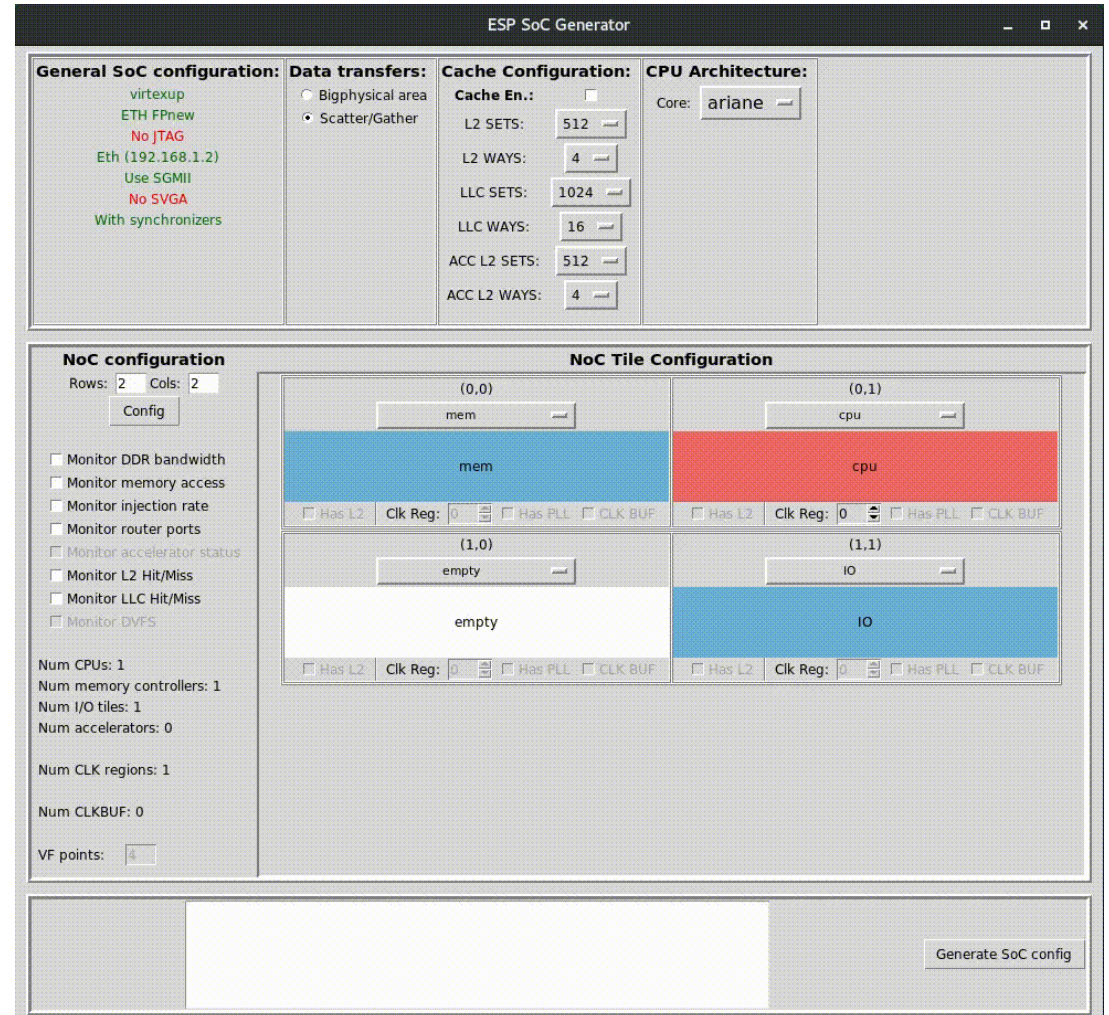
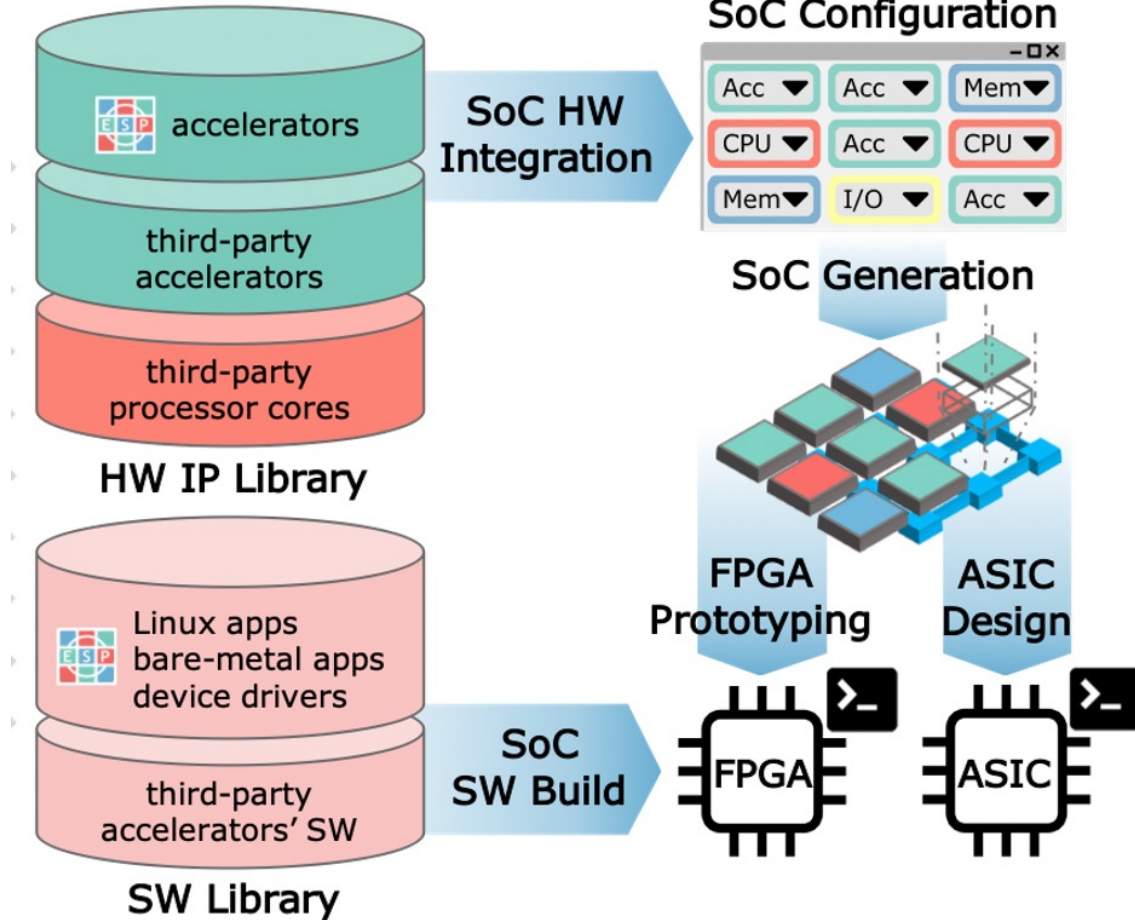


The Wolt Resource Manager: Example



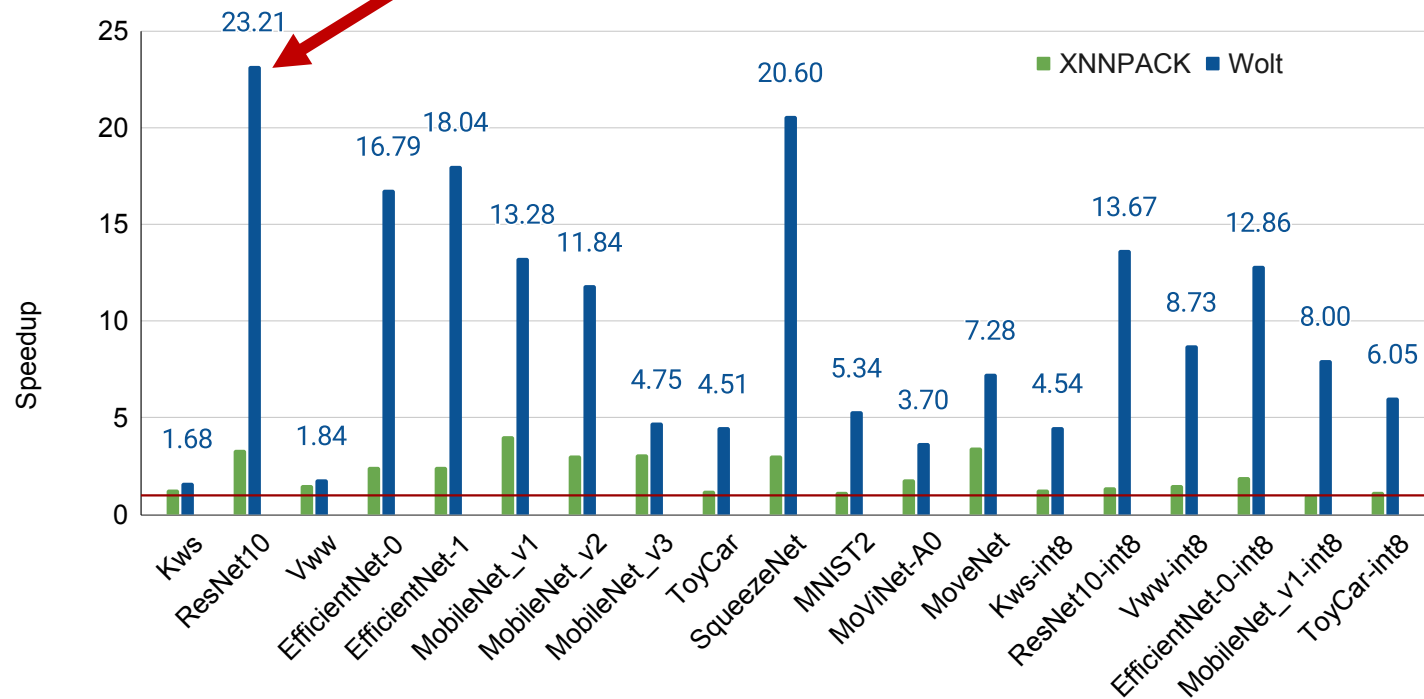
- A better accelerator assignment improves performance

ESP SoC Flow

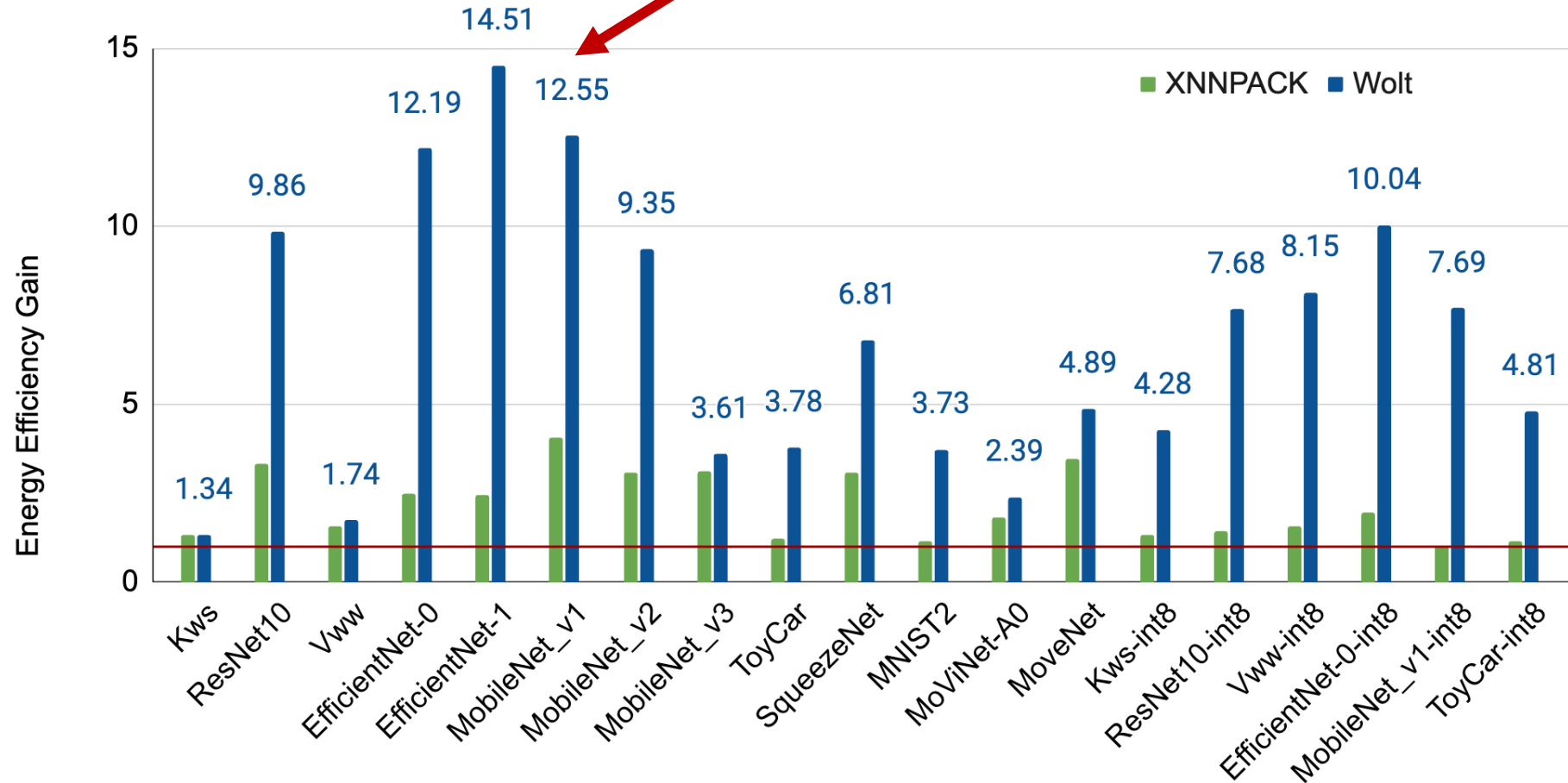


Evaluation: Performance

- Performance evaluation of Wolt with several different ML models
 - XNNPACK is the state-of-the-art optimized 3rd party software library for NN inference
 - The baseline is the software execution on a single RISC-V core

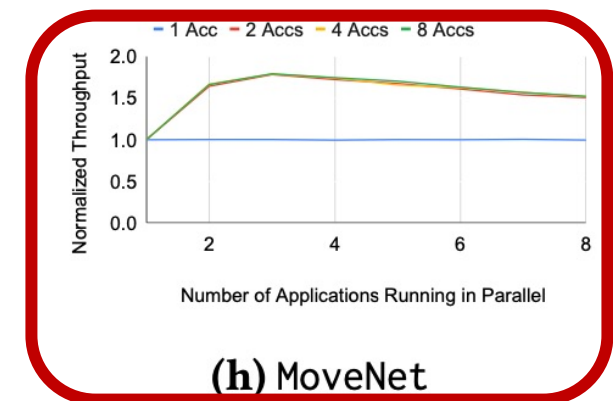
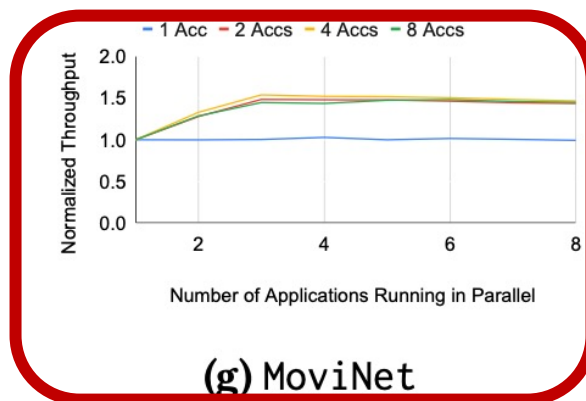
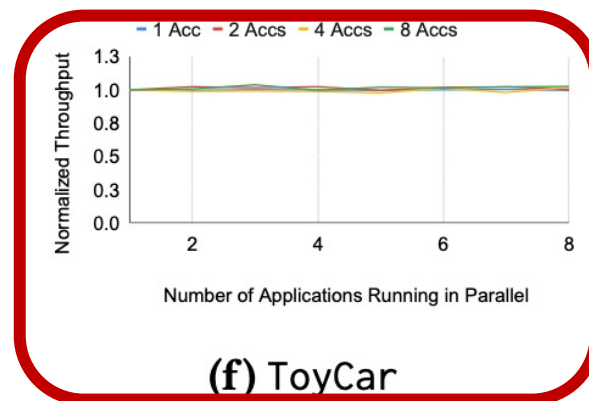
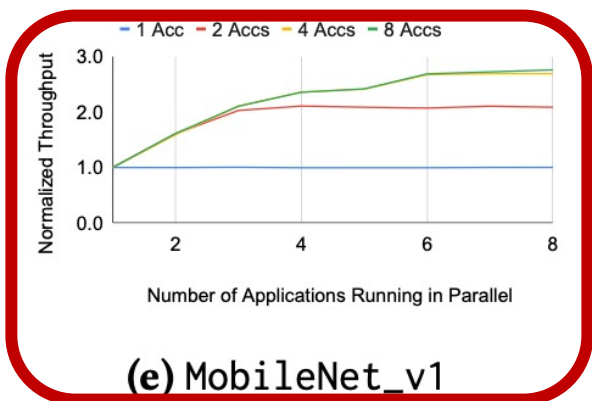
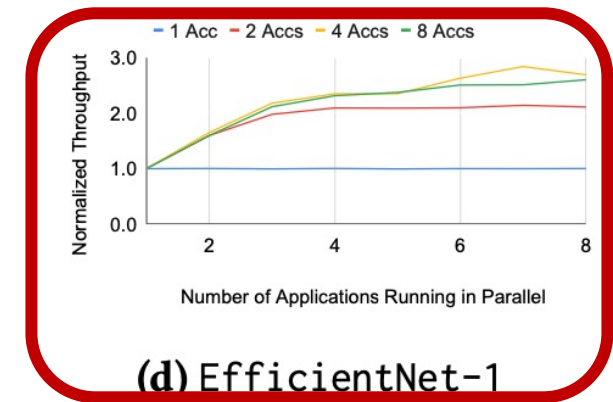
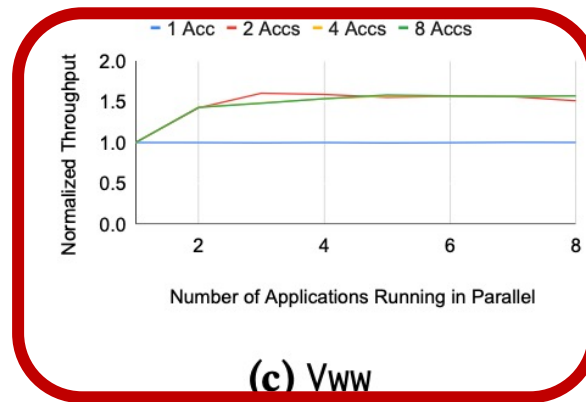
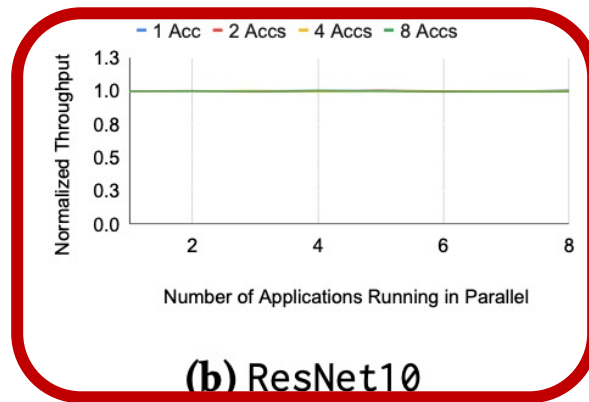
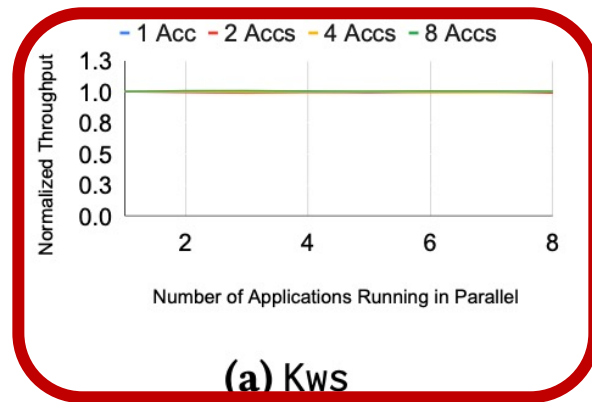


Evaluation: Energy Efficiency Gain



Experimental Evaluation: Wolt Resource Manager

- Evaluation of the resource manager with 8 different ML models



Time duration to run on the CPU > Time duration to run on the accelerators

Applications in Parallel with 2 CPUs

Amdahl's Law:

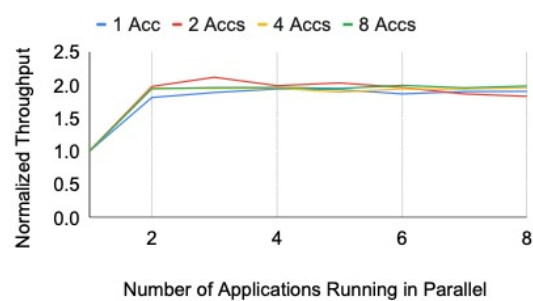
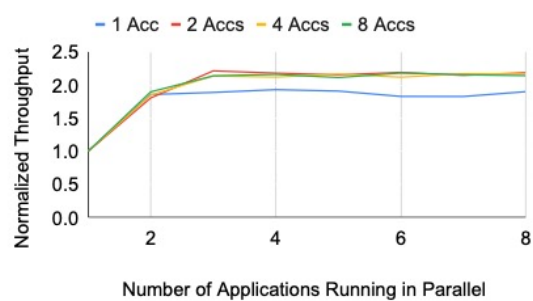
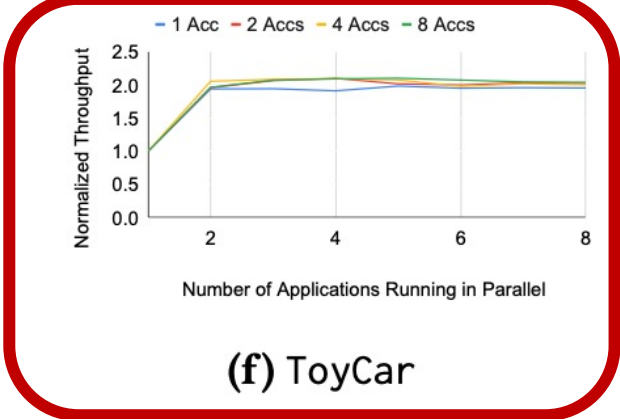
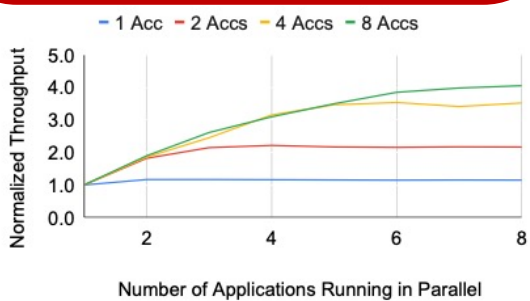
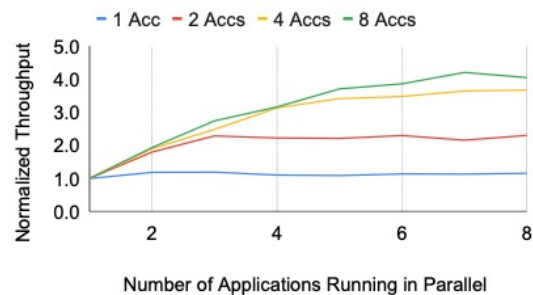
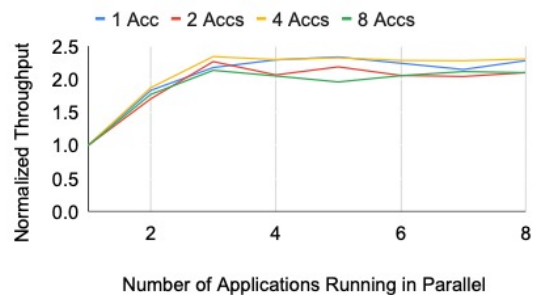
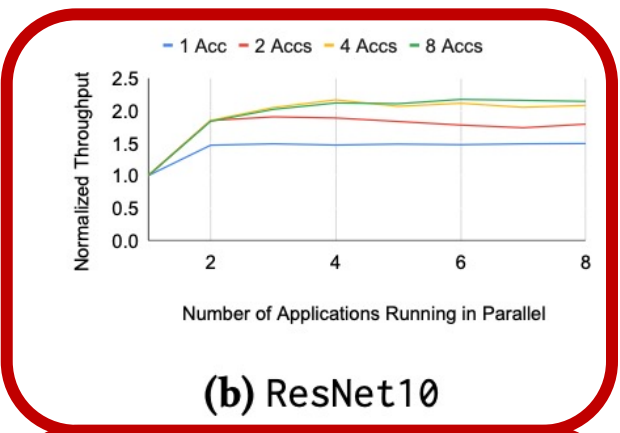
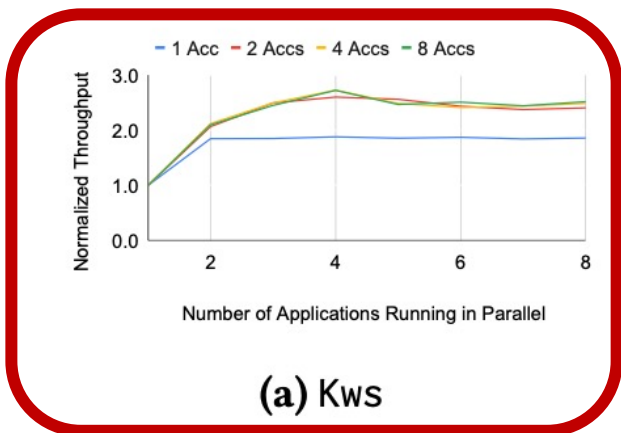
$$speedup = \frac{S + P}{S + \frac{P}{N}}$$

P: fraction can be parallelized
S: fraction that is serial
N: number of CPUs



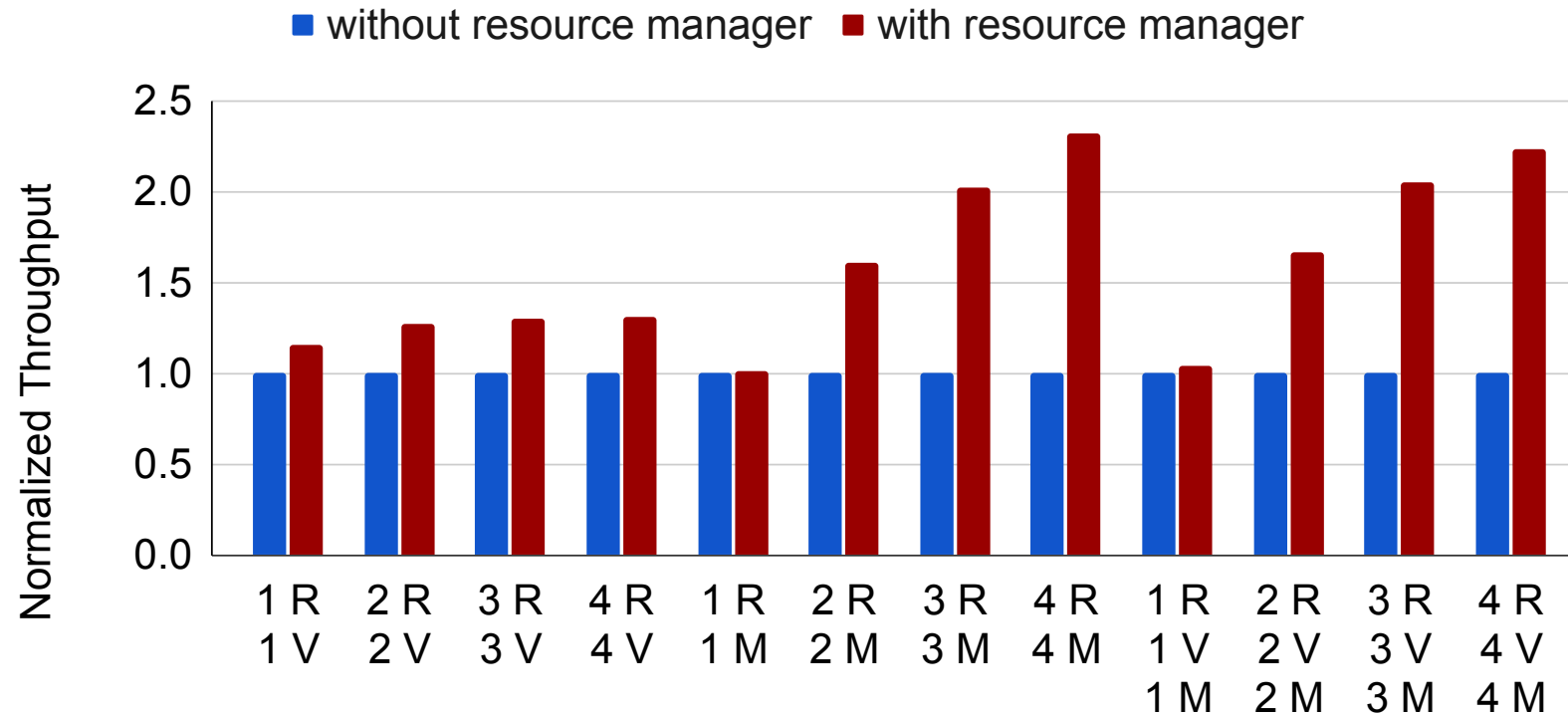
$$speedup = \frac{S + P}{\frac{S}{M} + \frac{P}{N}}$$

P: time duration on CPU
S: time duration on ACC
N: number of CPUs
M: number of ACCs



Experimental Evaluation: Multi-Tenant Scenario

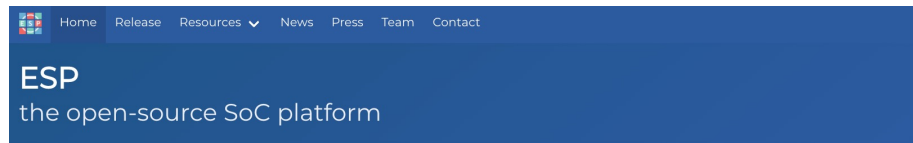
- Running a combination of ResNet10 (R), Vww (V), and MobileNet_v1 (M)
- SoC configuration: one CPU core and four Conv2d accelerators



Conclusion and Future Work

- With Wolt, we developed
 - a **transparent software layer** for TensorFlow Lite ML applications that decouples the development of software and hardware
 - a **resource manager** for running multi-tenant ML applications that enables the auto-assignment of accelerators and improves the performance for parallel executions
- We plan to
 - extend the Wolt's support to other ML frameworks (e.g., PyTorch and TVM)
 - improve the resource manager by adding more sophisticated scheduling mechanisms and supporting dynamic delegation of a given functionality to different accelerator instances that offer different design trade-offs

ESP : An Open-Source Platform for SoC Design



Tutorials

News: Upcoming tutorial on ESP on April 12 at ASPLOS 2021.

The ESP Vision

ESP is an open-source research platform for heterogeneous system-on-chip design that combines a scalable tile-based architecture and a flexible system-level design methodology.

ESP provides three accelerator flows: RTL, high-level synthesis (HLS), machine learning frameworks. All three design flows converge to the ESP automated SoC integration flow that generates the necessary hardware and software interfaces to rapidly enable full-system prototyping on FPGA.

Overview

Latest Posts

Upcoming tutorial at ASPLOS 2021

On April 12 we will present a tutorial on ESP at ASPLOS 2021.

[Read more](#)

Published: Apr 2, 2021

Release 2021.1.0

The GitHub Release 2021.1.0 of ESP is now available.

[Read more](#)

Published: Jan 26, 2021

ISCA 2024 tutorial The ESP Approach to Agile Chip Design

Latest update: 2024-06-20

- [Logistics](#)
- [Tutorial Overview](#)
- [Preliminary setup](#)
- [Program](#)
- [Video recording](#)
- [Team](#)
 - [Speakers](#)

- 9:00 am Jun 30, 2024
- Location: Quebracho A

www.esp.cs.columbia.edu



Thank you from the **ESP** team!



sld.cs.columbia.edu



[ColumbiaSld](https://twitter.com/ColumbiaSld)



esp.cs.columbia.edu



[sld-columbia/esp](https://github.com/sld-columbia/esp)



[c/ESP-platform](https://www.youtube.com/c/ESP-platform)

Transparent Deployment of TFLite Workloads on Lightweight Many-Accelerator Architectures

Kuan-Lin Chiu, Guy Eichler, Chuan-Tung Lin

Giuseppe Di Guglielmo and Luca Carloni

