

# Analysis and Modeling of Collaborative Execution Strategies for Heterogeneous CPU-FPGA Architectures

Sitao Huang<sup>1</sup>, Li-Wen Chang<sup>2</sup>, Izzat El Hajj<sup>3</sup>, Simon Garcia De Gonzalo<sup>1</sup>, Juan Gómez-Luna<sup>4</sup>, Sai Rahul Chalamalasetti<sup>5</sup>, Mohamed El Hadedy<sup>6</sup>, Dejan Milojicic<sup>5</sup>, Onur Mutlu<sup>4</sup>, Deming Chen<sup>1</sup>, Wen-mei Hwu<sup>1</sup>

<sup>1</sup>University of Illinois at Urbana-Champaign, USA

<sup>2</sup>Microsoft, USA

<sup>3</sup>American University of Beirut, Lebanon

<sup>4</sup>ETH Zürich, Switzerland

<sup>5</sup>Hewlett Packard Labs, USA

<sup>6</sup>California State Polytechnic University, Pomona, USA



**I** ILLINOIS

Electrical & Computer Engineering

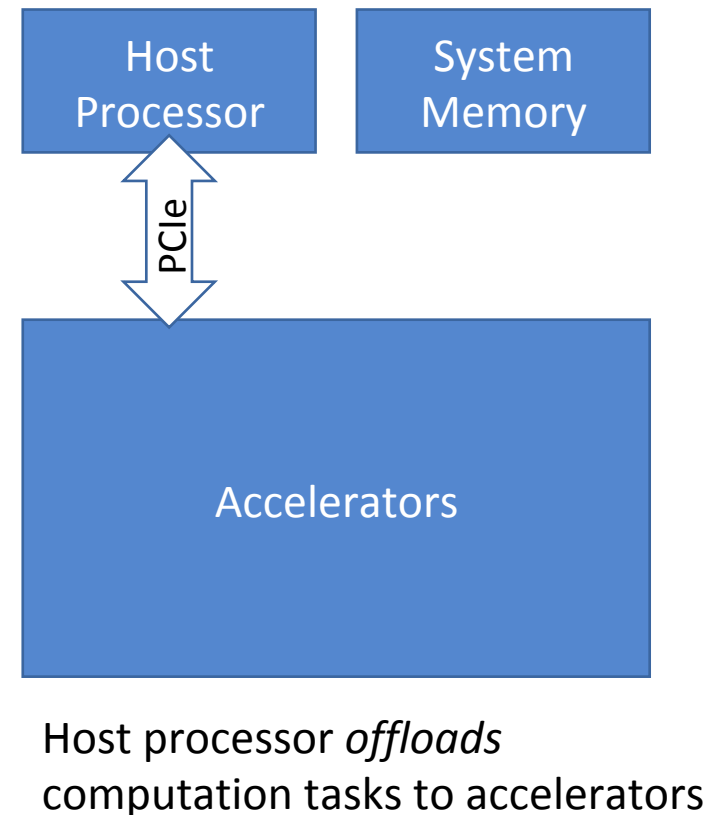
COLLEGE OF ENGINEERING

# Outline

- Introduction to Collaborative Computing
- Analytical Models of Collaborative Computing
- Evaluation
  - Evaluation platform
  - Data partitioning vs. task partitioning
  - Impact of kernel replication
- Key Insights
- **Chai** Benchmarks for CPU-FPGA Systems

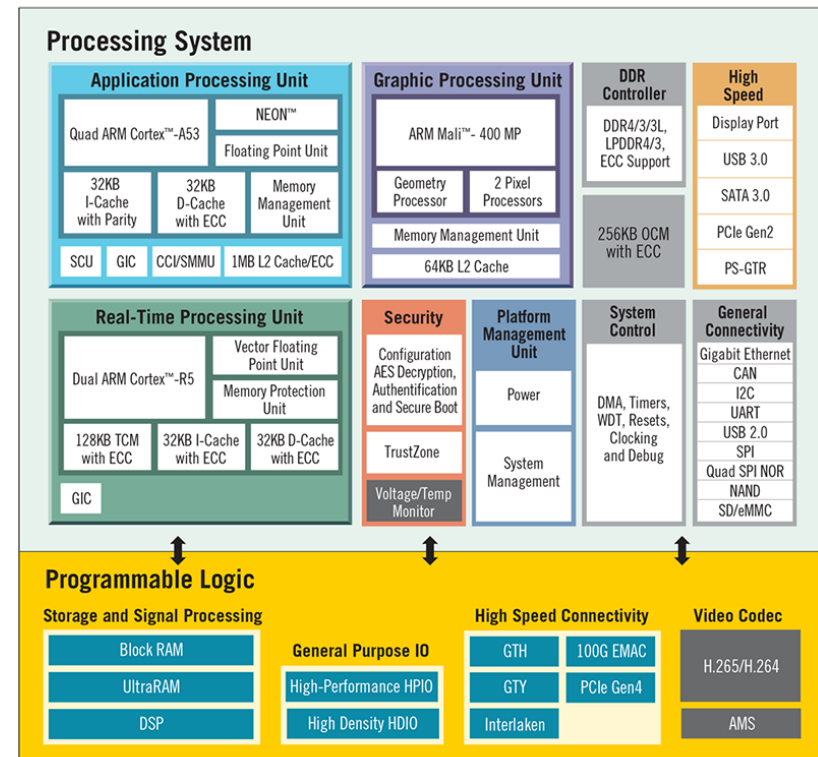
# Collaborative Computing

- Traditionally, accelerators (GPUs, FPGAs, etc.) have been used as *offload* engines



# Collaborative Computing

- Traditionally, accelerators (GPUs, FPGAs, etc.) have been used as *offload* engines
- Heterogeneous architectures moving towards tighter integration
  - Unified memory
  - System-wide atomics



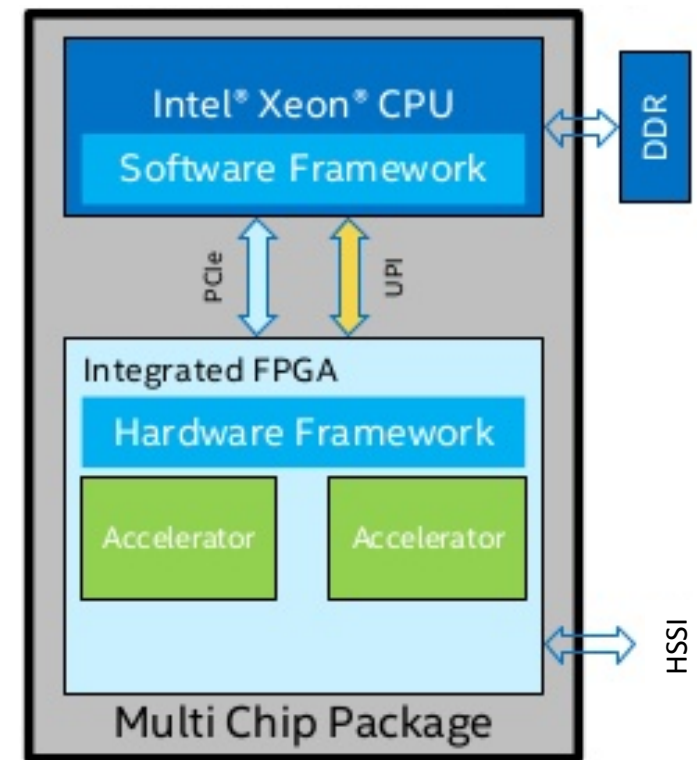
Note: Illustration not drawn to scale.

Xilinx Zynq UltraScale+ MPSoC

# Collaborative Computing

- Traditionally, accelerators (GPUs, FPGAs, etc.) have been used as *offload* engines
- Heterogeneous architectures moving towards tighter integration
  - Unified memory
  - System-wide atomics
- Tighter integration allows fine-grained collaboration

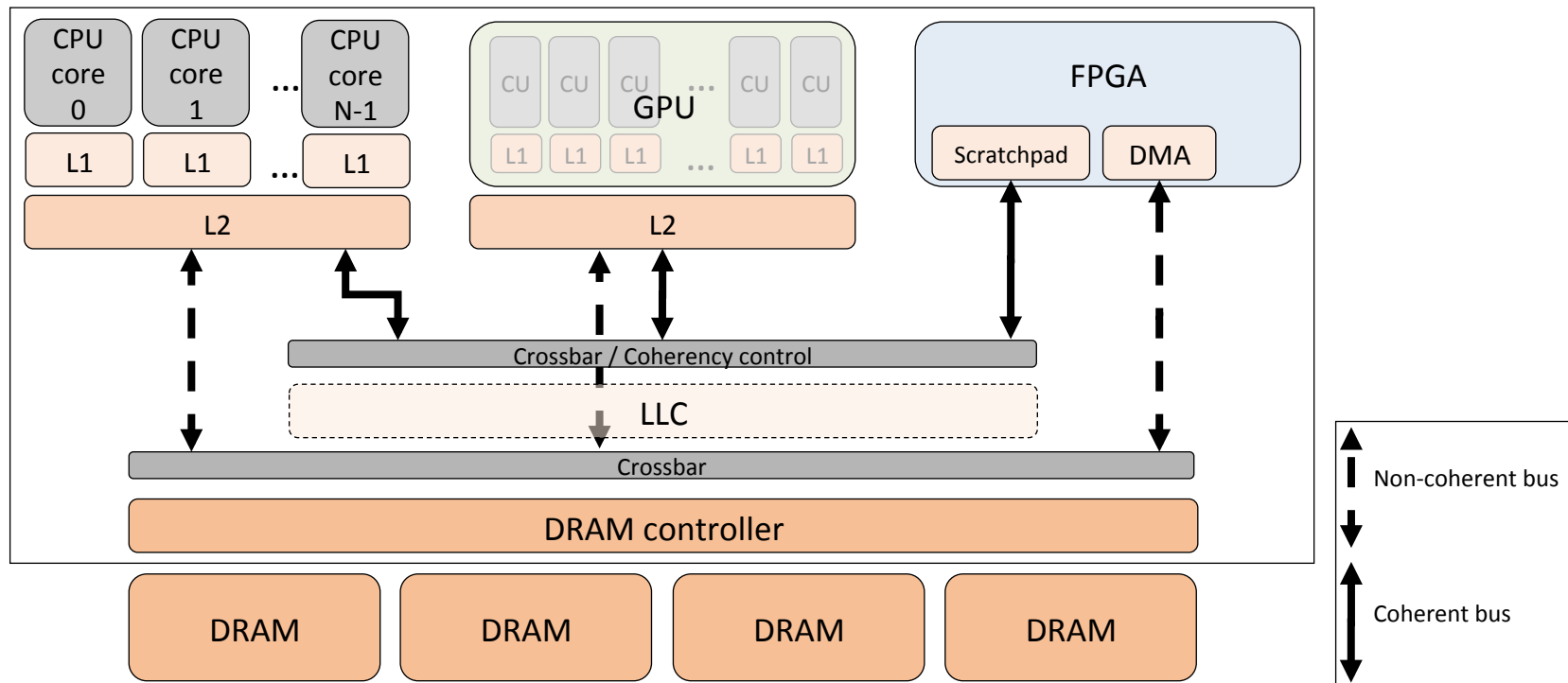
**Key challenge:** identify the best CPU-FPGA collaboration strategy



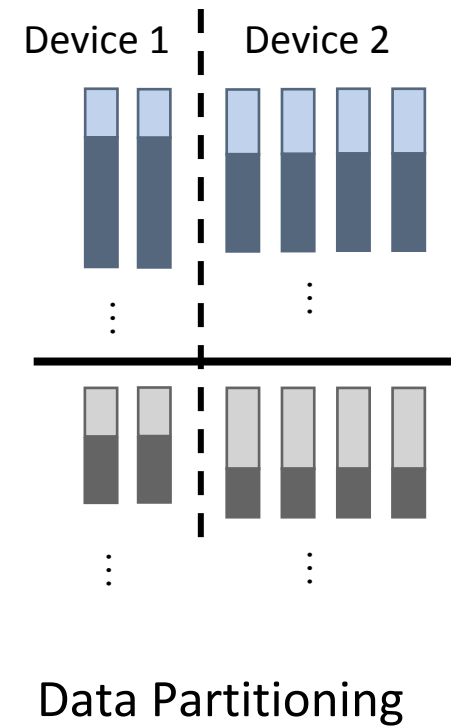
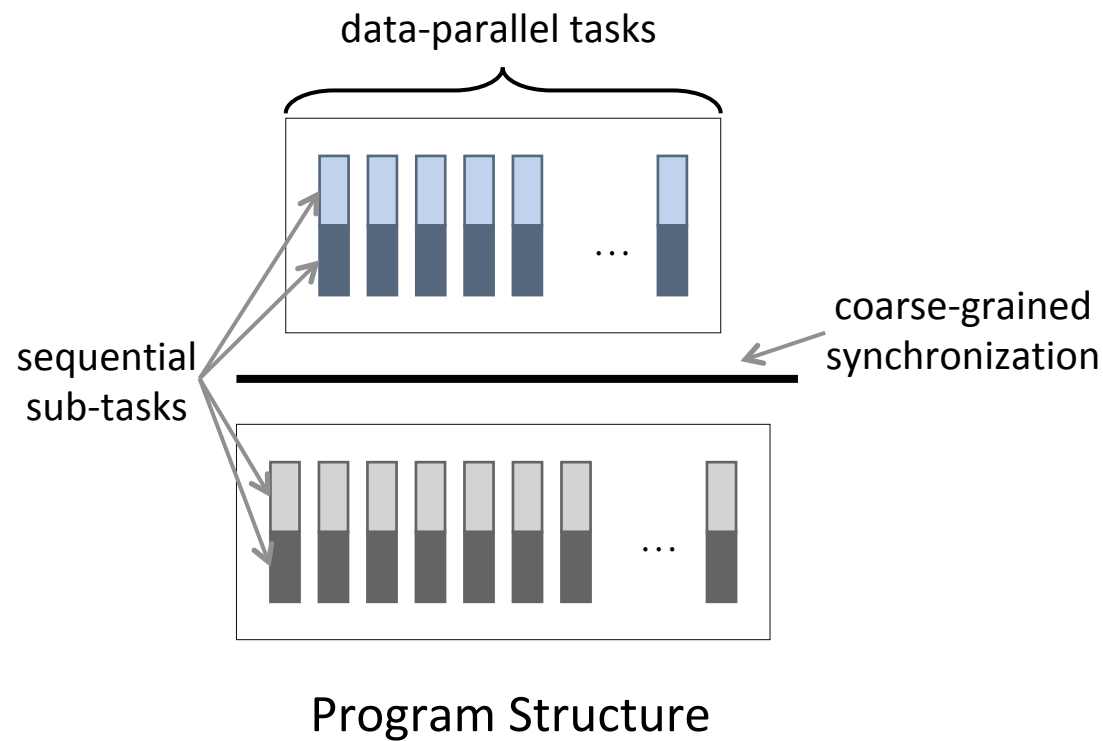
Intel Xeon + FPGA Integrated Platform (MCP)

1000000

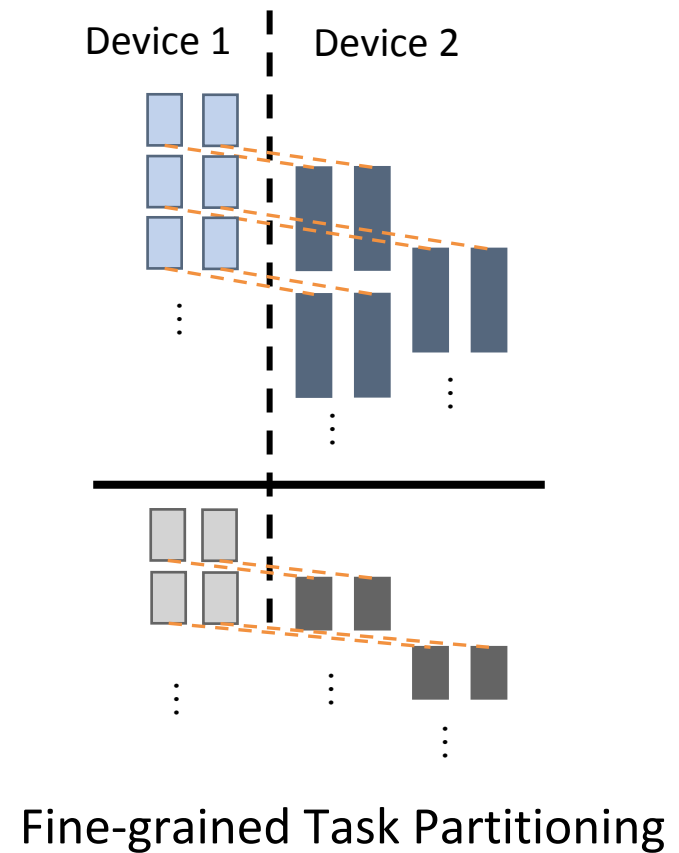
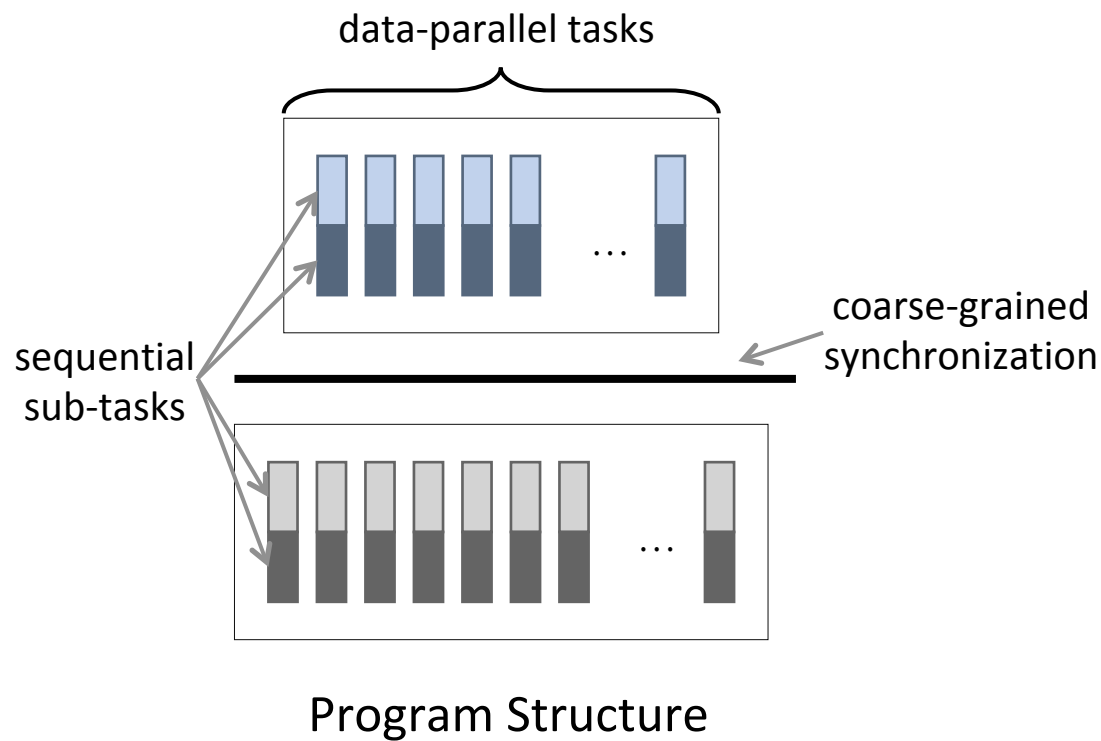
Our vision of an integrated heterogeneous system:



# Collaborative Patterns

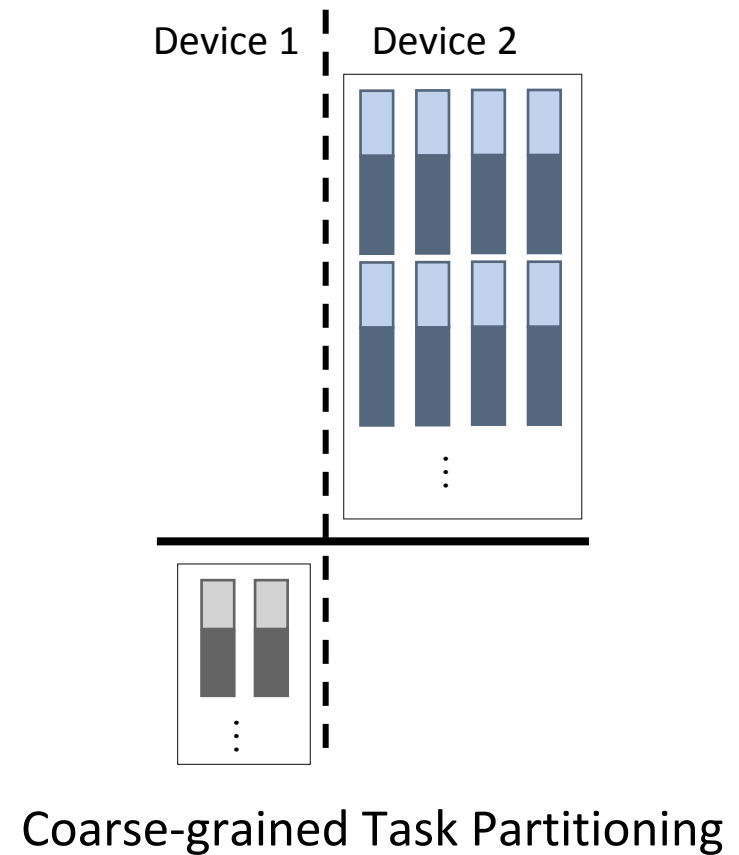
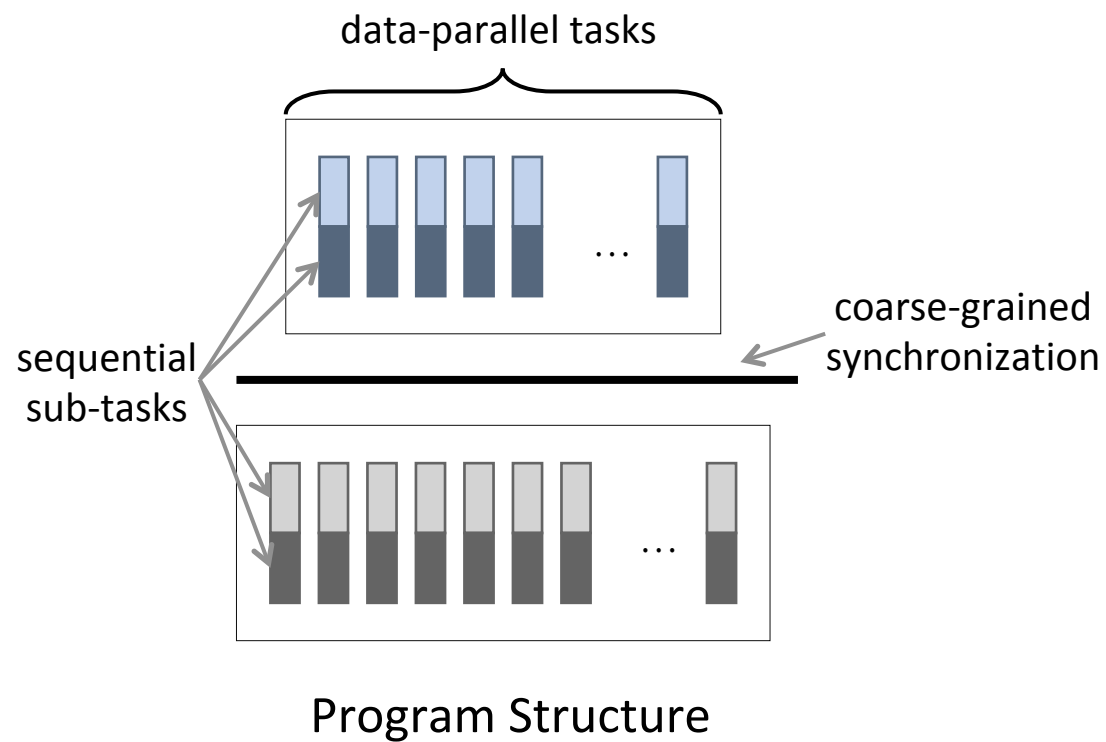


# Collaborative Patterns





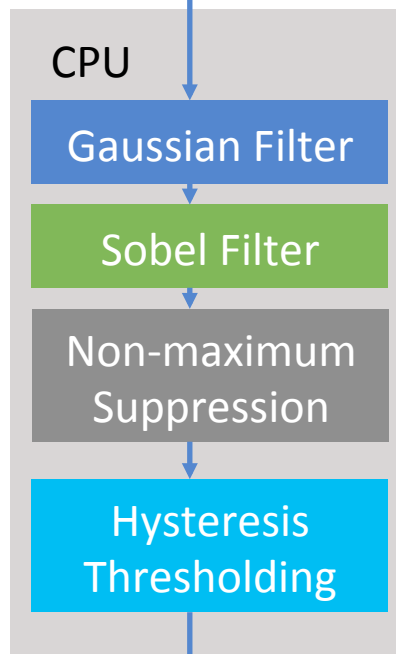
# Collaborative Patterns



# Data Partitioning

Using Canny Edge Detection (CED) as an example

Input images



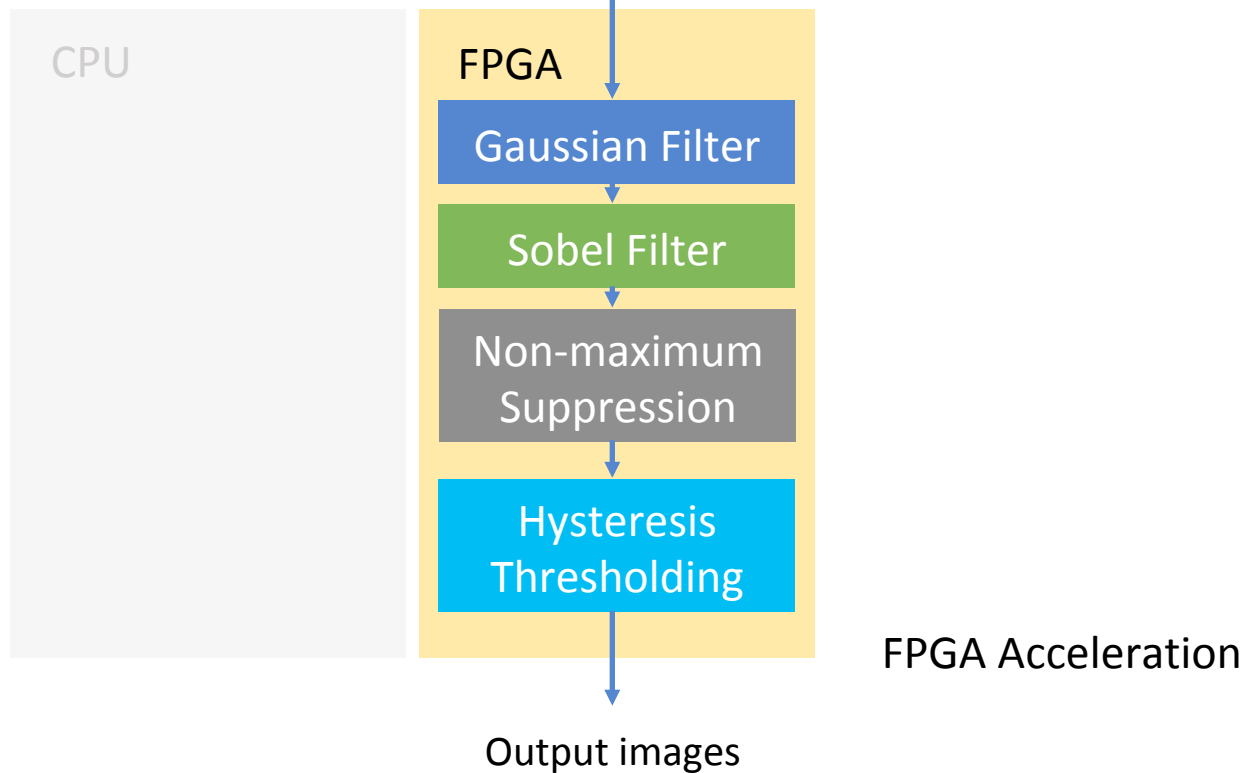
Output images

CPU Implementation

# Data Partitioning

Using Canny Edge Detection (CED) as an example

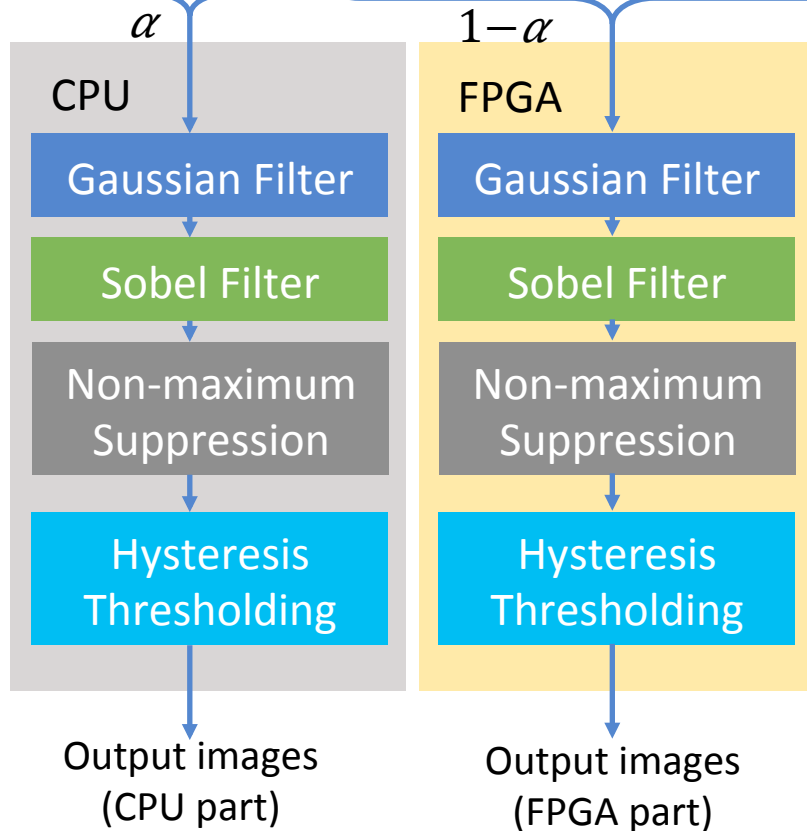
Input images



# Data Partitioning

Using Canny Edge Detection (CED) as an example

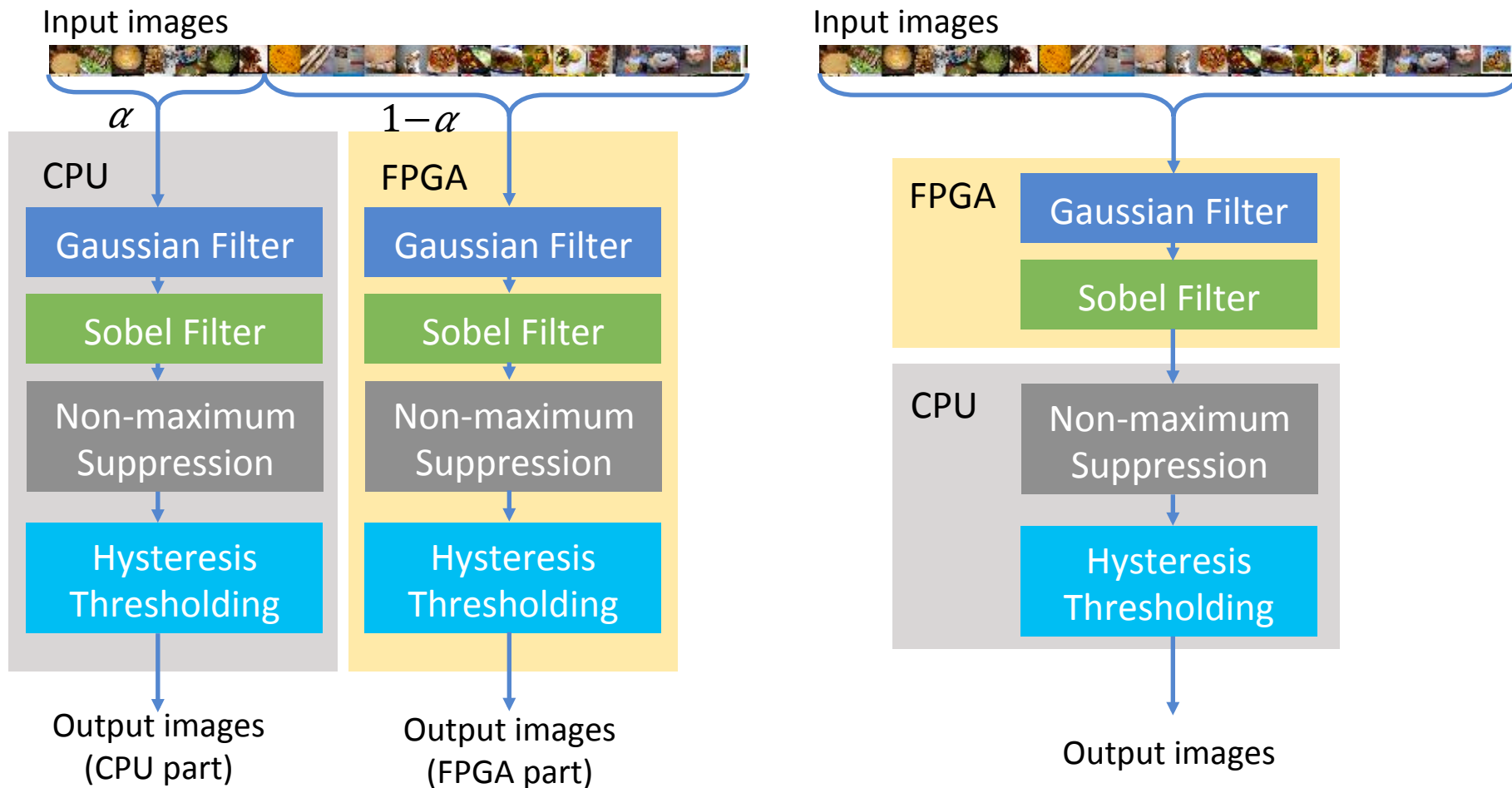
Input images



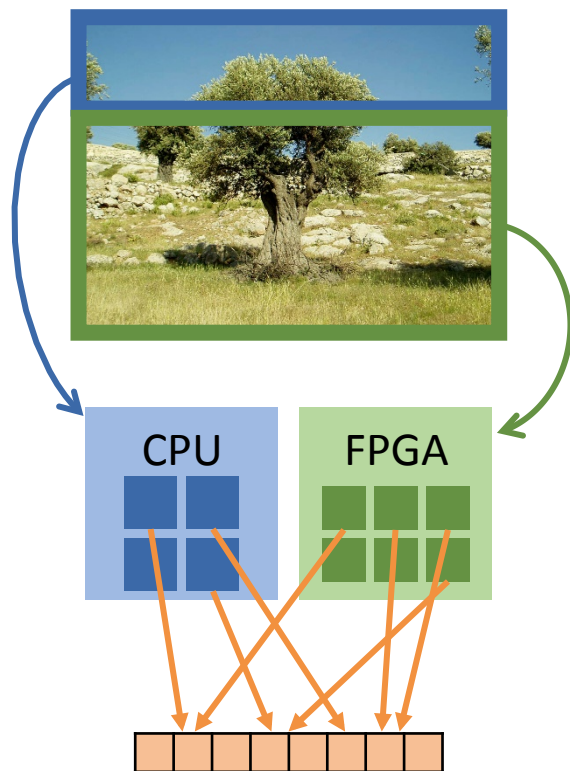
**CPU-FPGA Collaboration**

# Data Partitioning vs. Task Partitioning

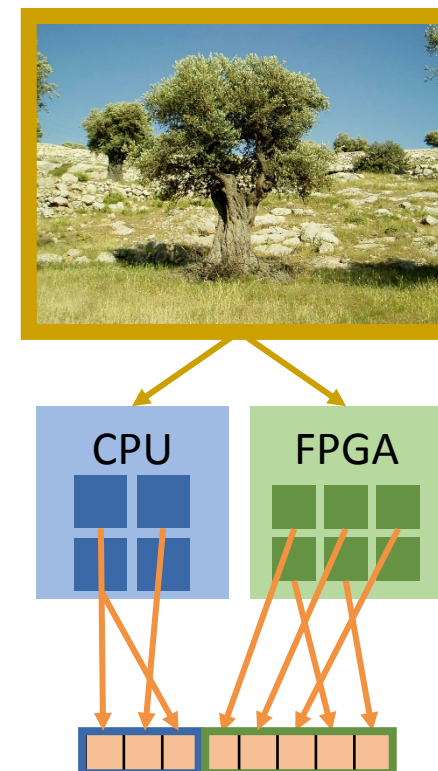
Using Canny Edge Detection (CED) as an example



## Another Data Partitioning Example: Image Histogram



**Input pixels** distributed across devices



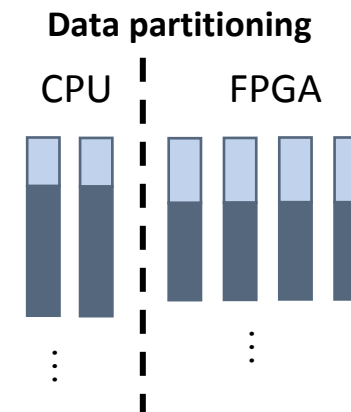
**Output bins** distributed across devices

# Analytical Models

- $N$ : Number of data parallel tasks in the application
- $t_{\downarrow i, C}$ : Execution time of sub-task  $i$  by a CPU worker
- $t_{\downarrow i, F}$ : Execution time of sub-task  $i$  by an FPGA worker
- $w_{\downarrow C}$ : Number of available CPU workers
- $w_{\downarrow F}$ : Number of available FPGA workers
- $\beta$ : Distribution and aggregation overhead factor
- $\alpha$ : Fraction of data parallel tasks assigned to CPU

# Analytical Models

- $N$ : Number of data parallel tasks in the application
- $t_{i,C}$ : Execution time of sub-task  $i$  by a CPU worker
- $t_{i,F}$ : Execution time of sub-task  $i$  by an FPGA worker
- $w_C$ : Number of available CPU workers



- $w_F$ : Number of available FPGA workers

- $\beta$ : Distribution and aggregation overhead factor
- The total execution time is

$\alpha$ : Fraction of data parallel tasks assigned to CPU

$$t_{data, total} = \beta \cdot \max \left( \frac{\alpha N \sum_{i=1}^N t_{i,C}}{w_C}, (1 - \alpha) \frac{N \sum_{i=1}^N t_{i,F}}{w_F} \right)$$

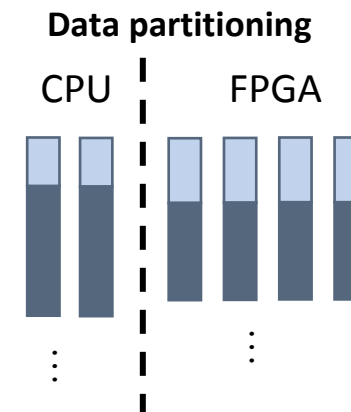
Total CPU execution time (sequential execution)

Total FPGA execution time (sequential execution)



# Analytical Models

- $N$ : Number of data parallel tasks in the application
- $t_{i,C}$ : Execution time of sub-task  $i$  by a CPU worker
- $t_{i,F}$ : Execution time of sub-task  $i$  by an FPGA worker
- $w_C$ : Number of available CPU workers
- $w_F$ : Number of available FPGA workers



- $\beta$ : Distribution and aggregation overhead factor
- The total execution time is

•  $\alpha$ : Fraction of data parallel tasks assigned to CPU

$$t_{data, total} = \beta \cdot \max \left( \alpha N \sum_{i=1}^N t_{i,C} / w_C, (1 - \alpha) N \sum_{i=1}^N t_{i,F} / w_F \right)$$

Fixing all the variables except  $\alpha$ , the optimal  $\alpha$  (global minimum point) is

$$\alpha^* = \frac{\sum_{i=1}^N t_{i,F} / w_F}{\sum_{i=1}^N t_{i,C} / w_C + \sum_{i=1}^N t_{i,F} / w_F}$$

Workloads of CPU and FPGA workers are balanced

# Analytical Models

- $N$ : Number of data parallel tasks in the application
- $t_{i,C}$ : Execution time of sub-task  $i$  by a CPU worker
- $t_{i,F}$ : Execution time of sub-task  $i$  by an FPGA worker
- $w_C$ : Number of available CPU workers
- $w_F$ : Number of available FPGA workers

## Fine-grained task partitioning

$\beta$ : Distribution and aggregation overhead factor  
The total execution time is

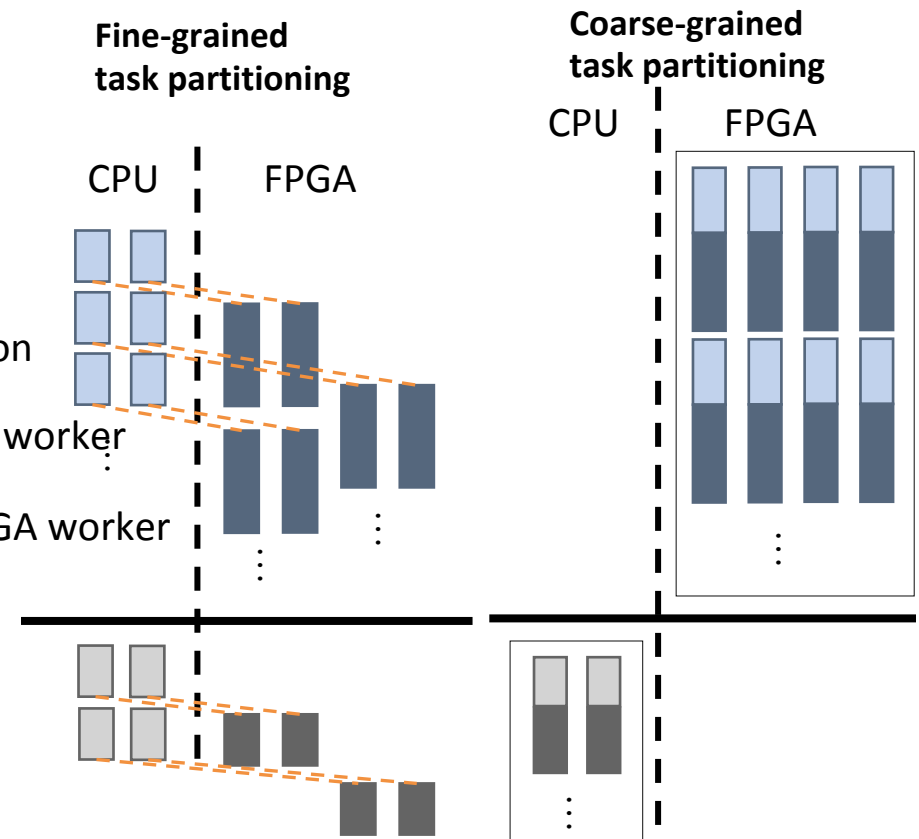
$$t_{task, total} = \beta \cdot N \cdot \max \left( \sum_{i \in S_C} t_{i,C} / w_C, \sum_{i \in S_F} t_{i,F} / w_F \right)$$

(Assume sub-tasks are very fine-grained)

## Coarse-grained task partitioning

The total execution time is

$$t_{task, total} = \beta \cdot N \cdot \left( \sum_{i \in S_C} t_{i,C} / w_C + \sum_{i \in S_F} t_{i,F} / w_F \right)$$



# Chai Benchmark Suite

**Chai:** Collaborative Heterogeneous Applications for Integrated-architectures

- Chai benchmark suite:

[chai-benchmarks.github.io](https://chai-benchmarks.github.io)

- 14 benchmarks covers data partitioning, fine-grain task partitioning, and coarse-grain task partitioning patterns
- OpenCL, C++ AMP, and CUDA versions
- Unified memory and system-wide atomic versions and traditional discrete architecture versions



# Evaluated Chai Benchmarks

Benchmark	Description	Strategy
<b>CED-D</b>	Canny Edge Detection	Data Partitioning
<b>CED-T</b>	Canny Edge Detection	Task Partitioning
<b>RSC-D</b>	Random Sample Consensus	Data Partitioning
<b>RSC-T</b>	Random Sample Consensus	Task Partitioning
<b>BS</b>	Bézier Surface	Data Partitioning
<b>HSTO</b>	Image Histogram	Data Partitioning
<b>SSSP</b>	Single-Source Shortest Path	Task Partitioning
<b>TQ</b>	Task Queue System (Synthetic)	Task Partitioning
<b>TQH</b>	Task Queue System (Histogram)	Task Partitioning

OpenCL-D (OpenCL discrete architecture) versions of these benchmarks are used.

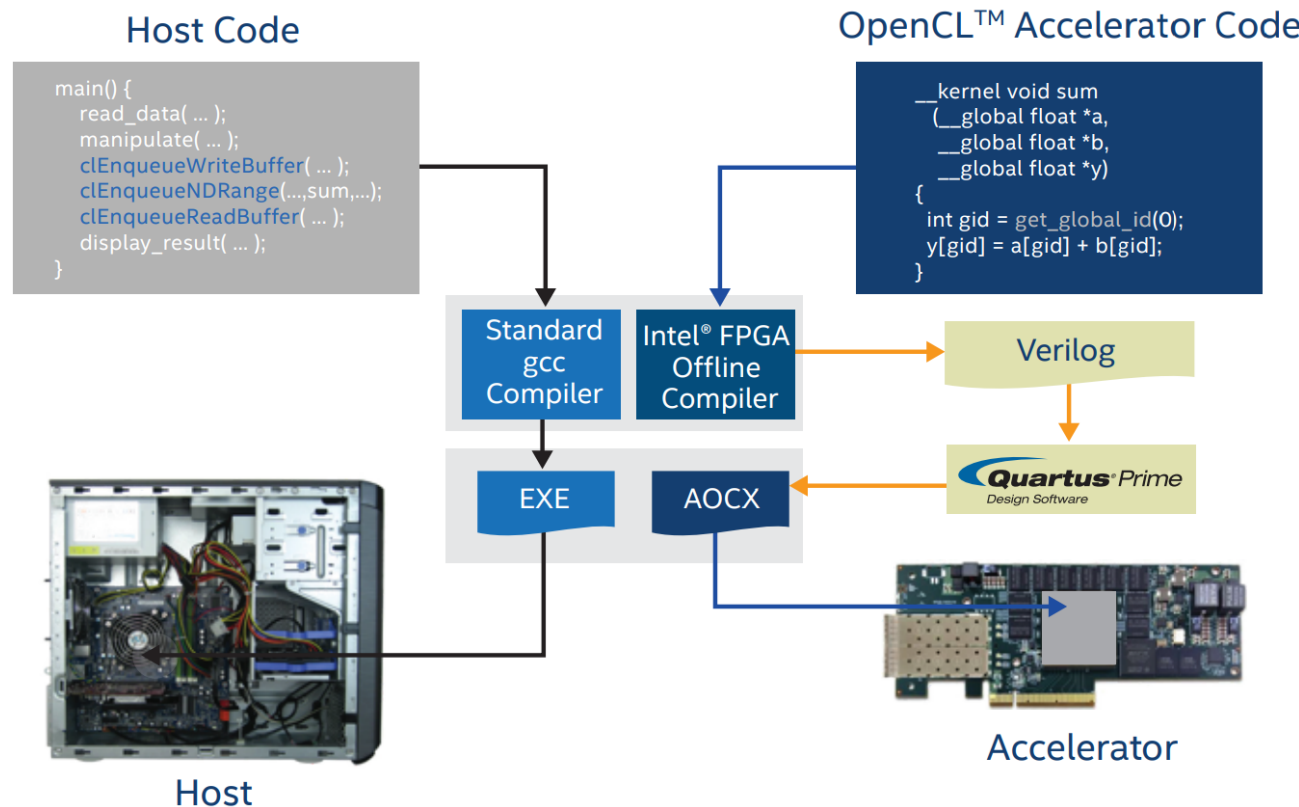
# Evaluation Platforms



	Platform A	Platform B
<b>FPGA Board</b>	Terasic DE5-Net	Nallatech 510T
<b>FPGA Chip</b>	<b>Intel Stratix V GX</b>	<b>Intel Arria 10 GX</b>
<b>On-Board Memory</b>	4 GB (DDR3)	8 GB (DDR4)
<b>Host CPU</b>	Intel Xeon E3-1240 v3	Intel Xeon E5-2650 v3
<b>Host Memory</b>	8 GB (DDR3)	96 GB (DDR4)
<b>Interface</b>	PCIe gen3.0 x8	PCIe gen3.0 x8

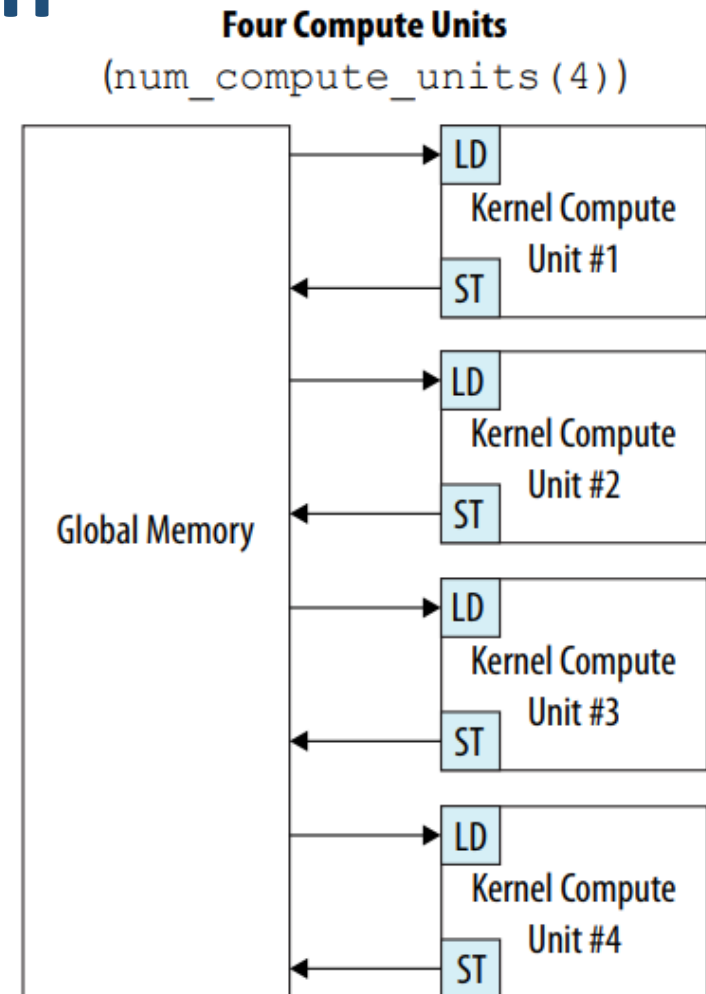
# Intel OpenCL SDK for FPGA

- Intel OpenCL SDK for FPGA is used to compile and synthesize host executable and FPGA design



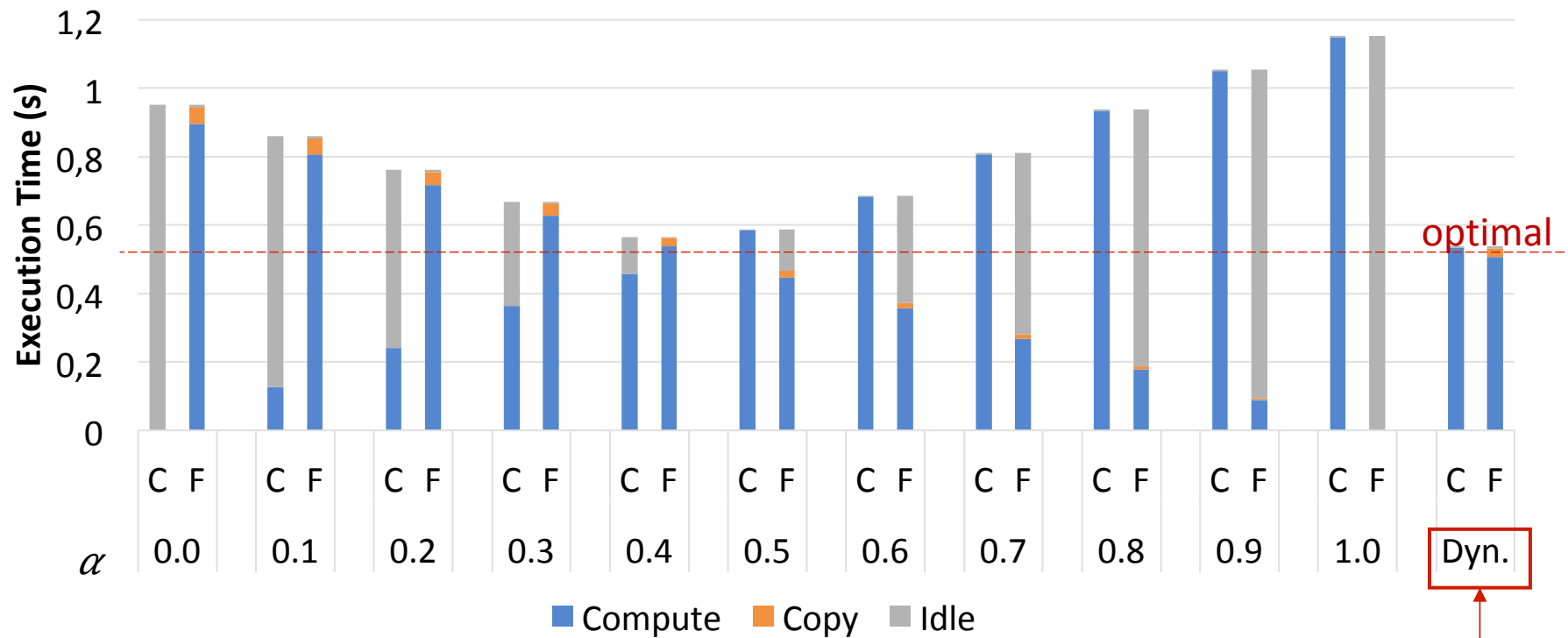
# Compute Unit Replication

- OpenCL kernels are synthesized to compute units on FPGA
- The compute units on FPGA can be replicated by adding `num_compute_units` attribute in the OpenCL kernel code
- `num_compute_units` attribute modifies the **number of compute units to which work-groups can be scheduled**, which also modifies the **number of times a kernel accesses global memory**
- We evaluate the impact of compute unit replication



Intel® FPGA SDK for OpenCL™ Best Practices Guide

# Evaluation: Canny Edge Detection (Data Partitioning)

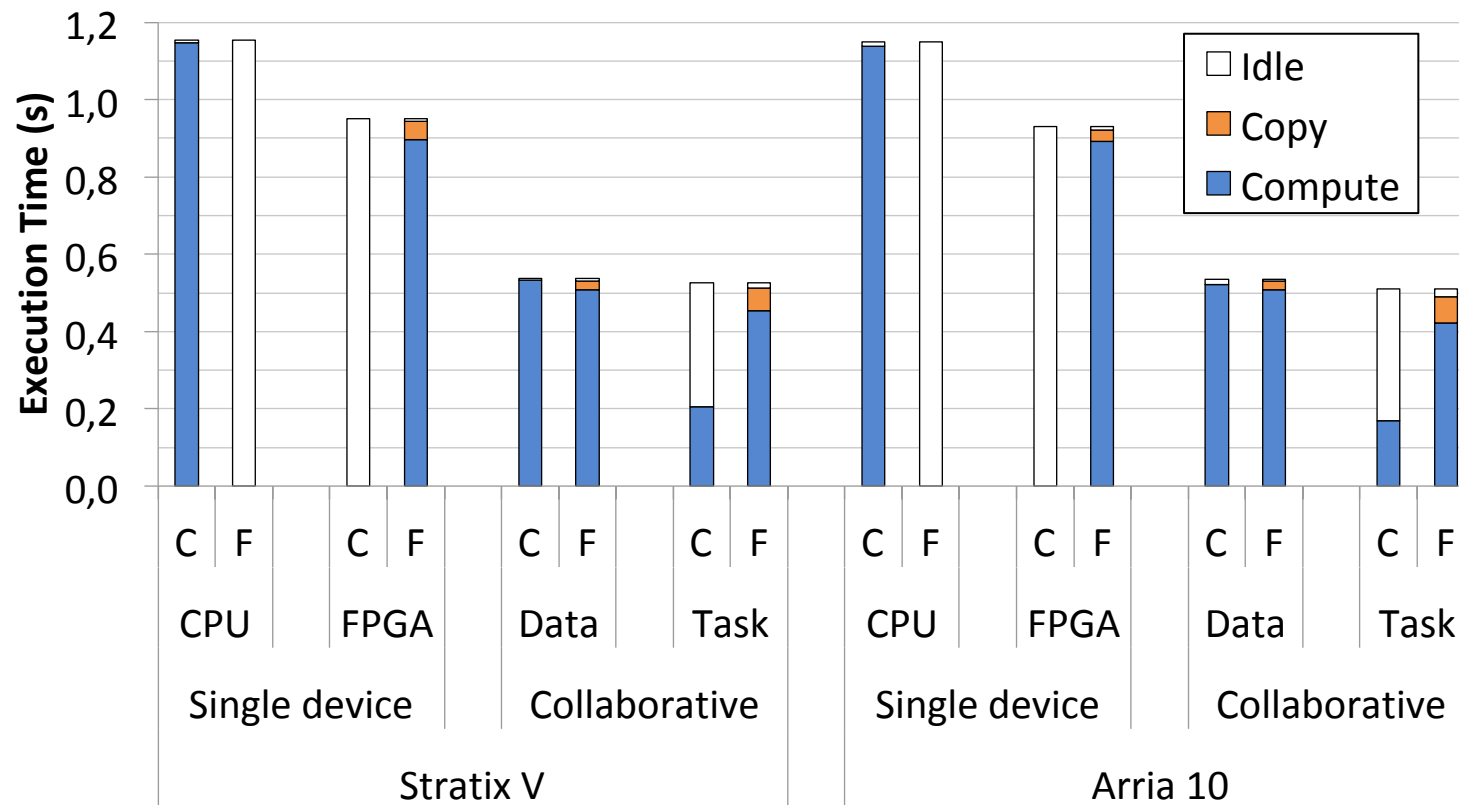


- C: CPU; F: FPGA
- $\alpha$ : Fraction of data parallel tasks assigned to CPU

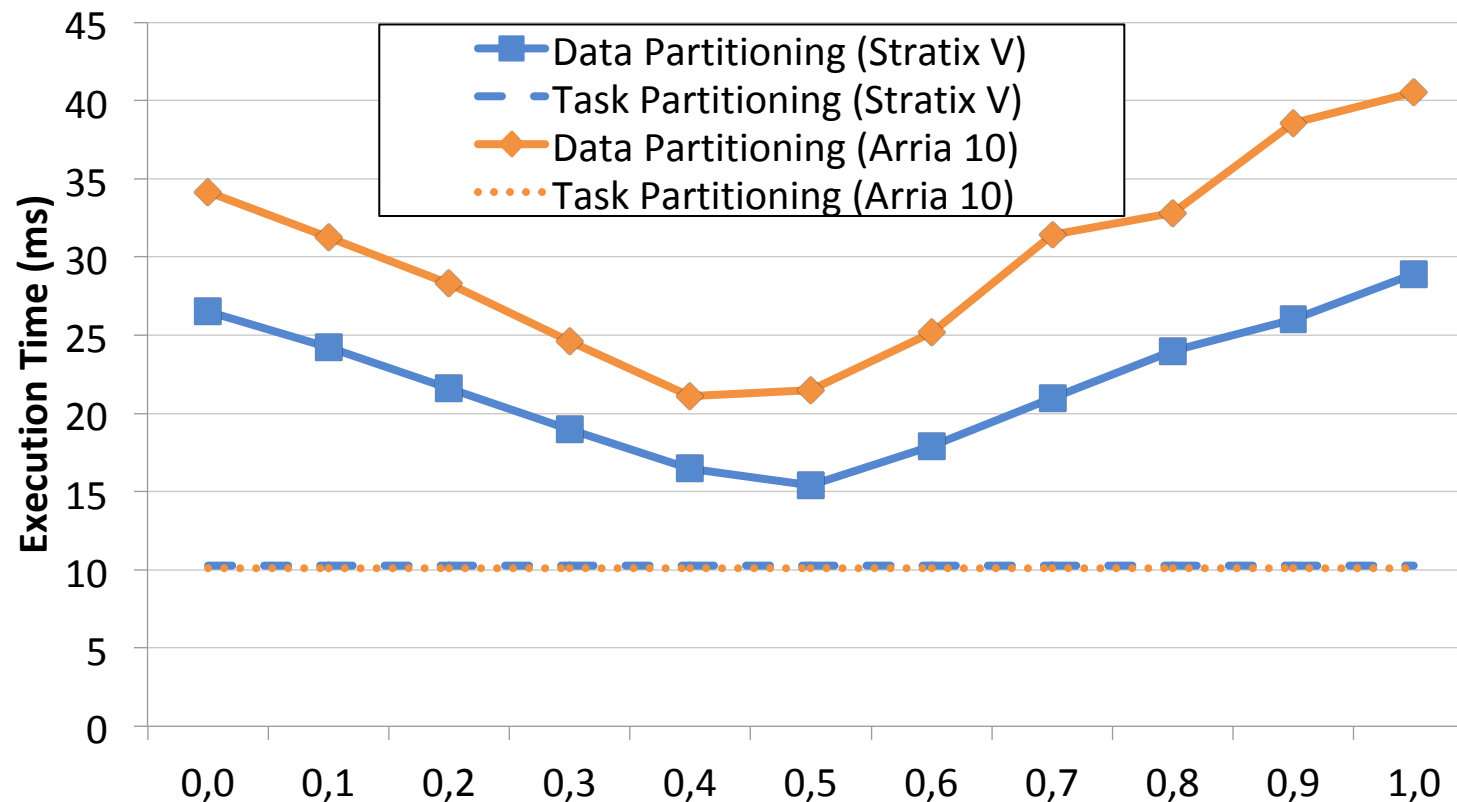
**Dynamic data partitioning:**  
assigning batch of data parallel  
tasks to idle devices, achieving  
dynamic workload balance



# Evaluation: Canny Edge Detection (Data Partitioning and Task Partitioning)



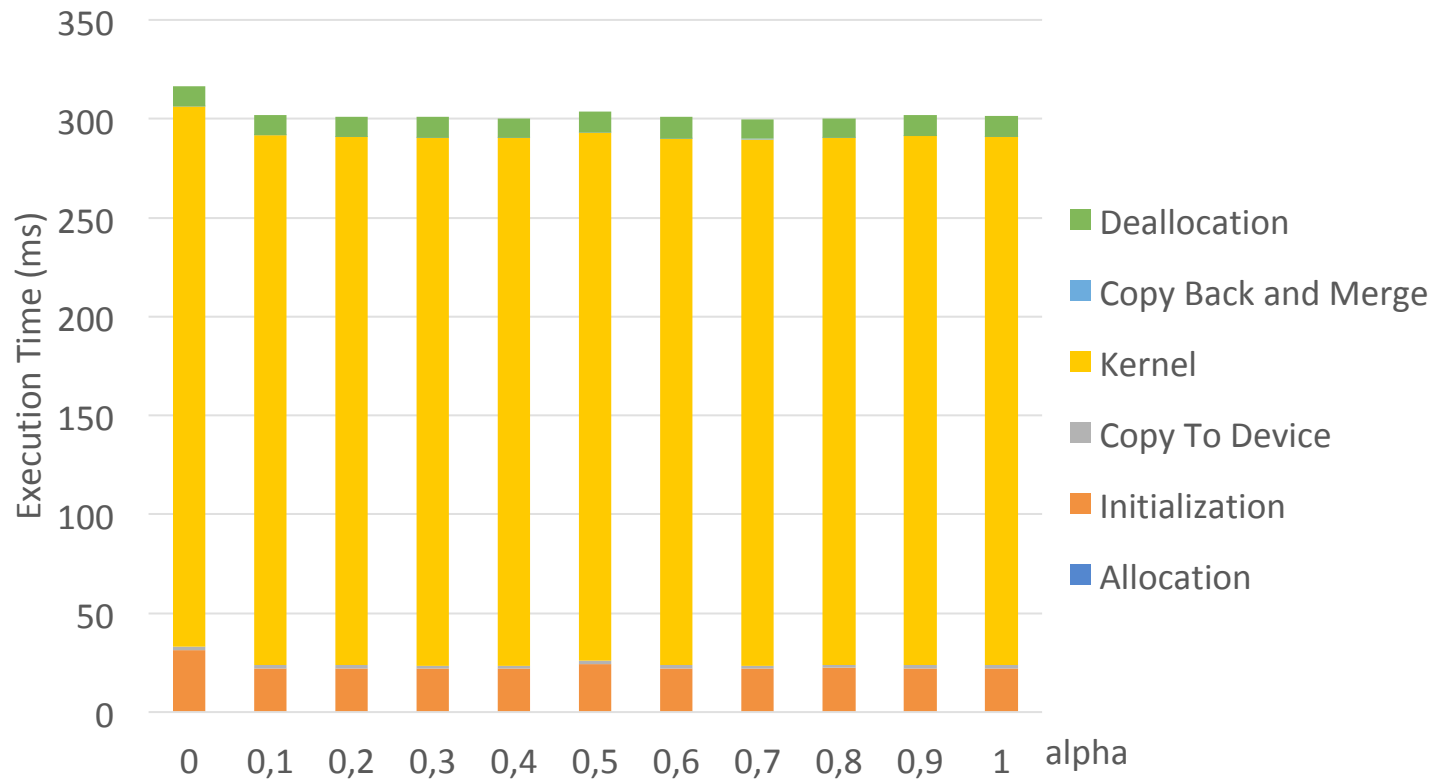
# Evaluation: Random Sample Consensus (Data Partitioning and Task Partitioning)



# Bézier Surface (BS, Data Partitioning)

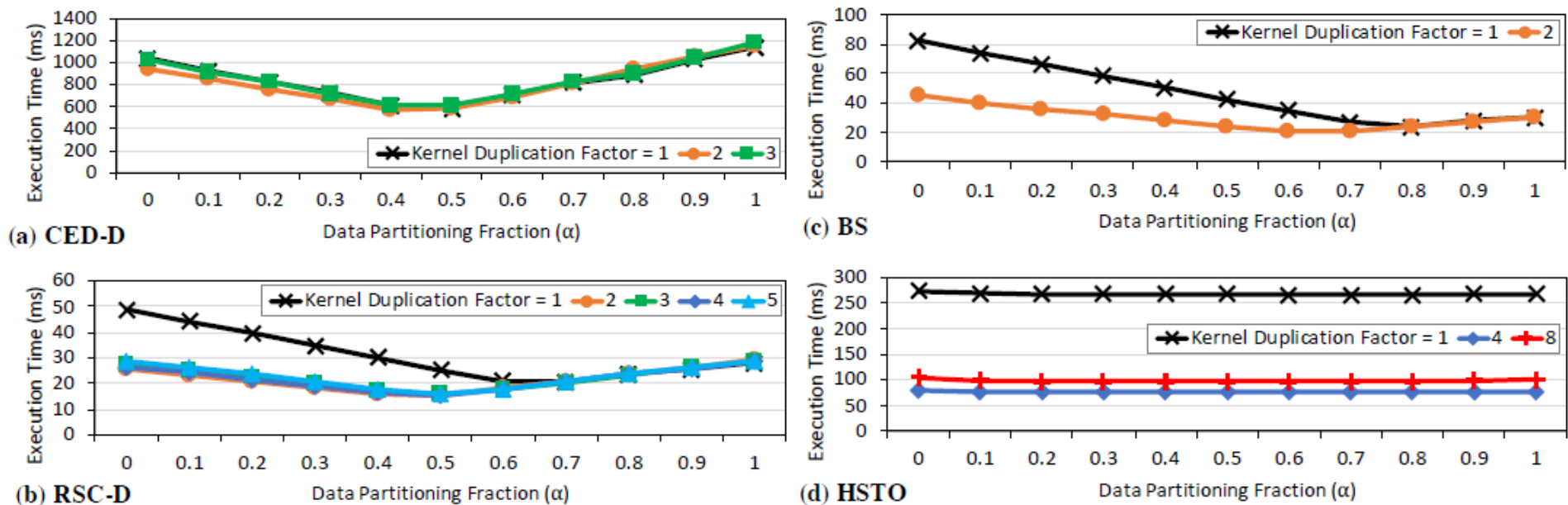


# Histogram (HSTO, Output Data Partitioning)



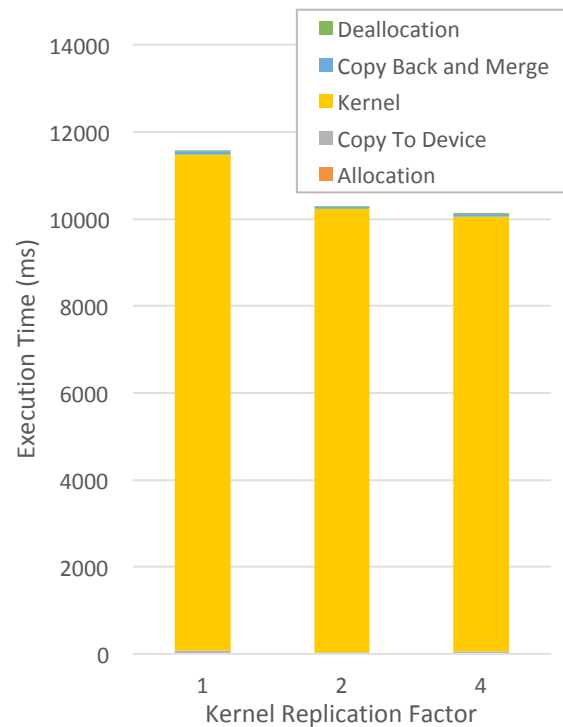
# Kernel Replication – Data Partitioning

We evaluated the performance under different kernel replication factors.

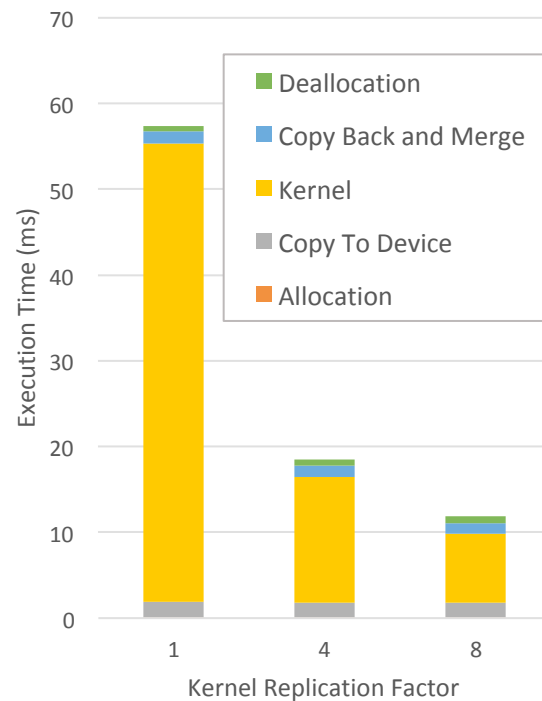


# Kernel Replication – Task Partitioning

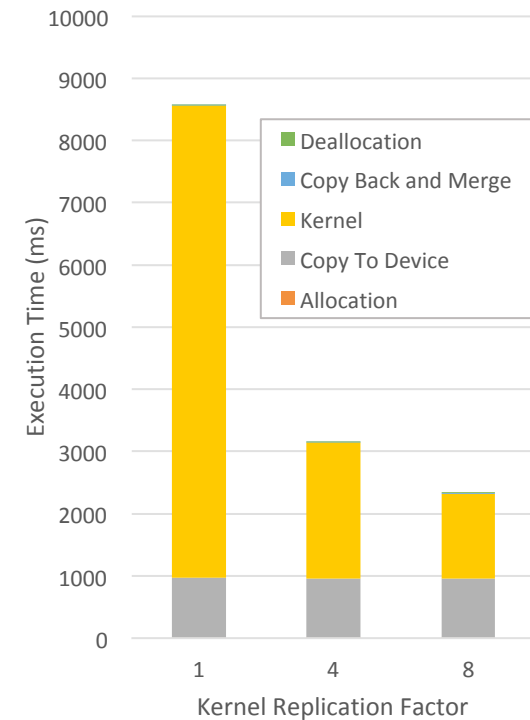
We evaluated the performance under different kernel replication factors.



SSSP



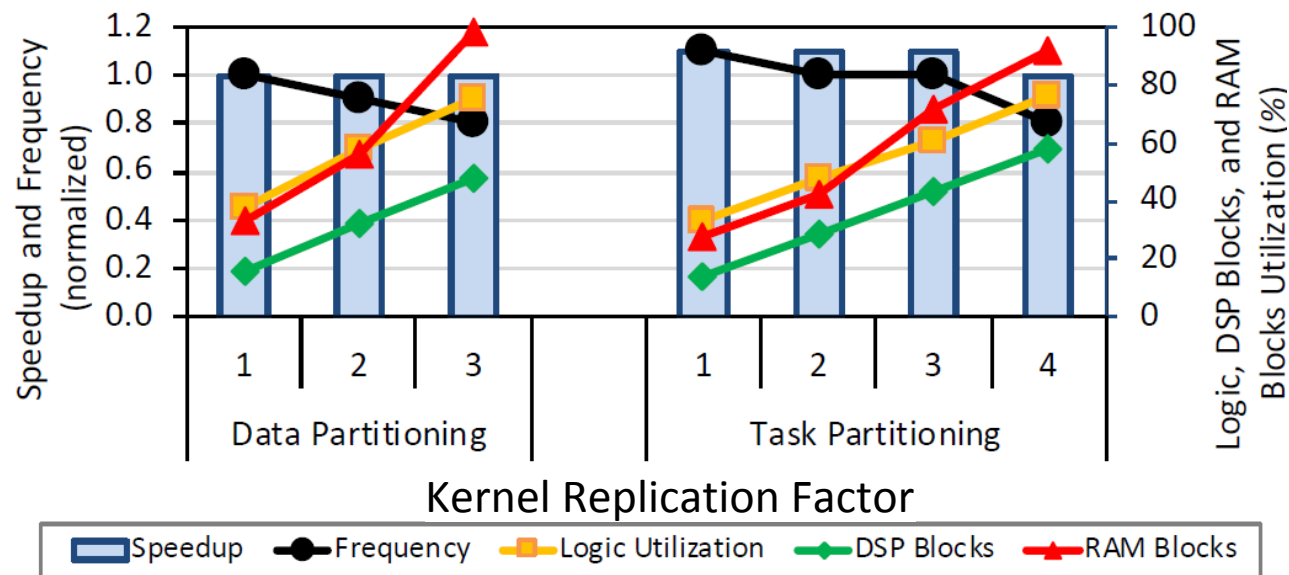
TQ



TQH

# Impact of Replication

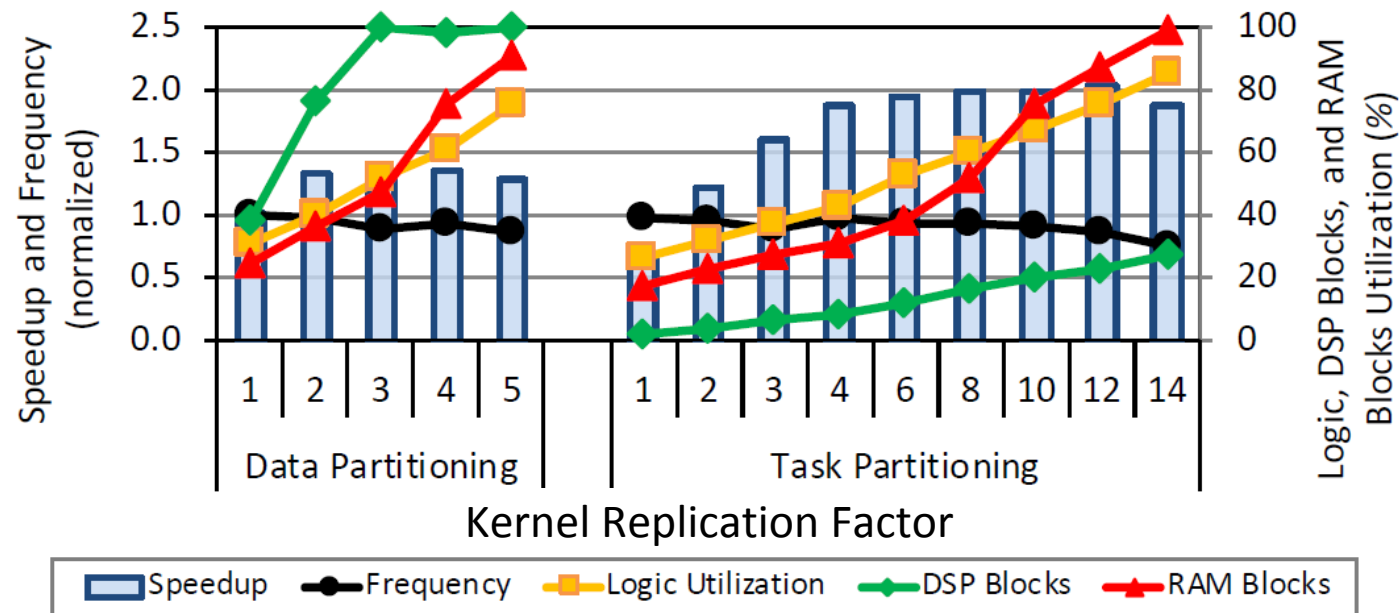
## Canny Edge Detection



- Replication factor for this application has little impact on performance
- Further profiling reveals the reason of performance saturation is the saturation of the memory bandwidth
- Task partitioning can afford a larger replication factor

# Impact of Replication

## Random Sample Consensus



- Replication improves performance of this application
- Bounding resource: DSP blocks
- Task partitioning releases the pressure on DSP block and thus can afford a larger replication factor



# Key Insights

- Collaborative execution is beneficial
- Data partitioning requires careful choice of partitions to provide the highest performance
- Task partitioning generally enables more kernel replication on the FPGA than data partitioning
- Data partitioning inflicts less burden on programmers and has less communication overhead than task partitioning
- OpenCL stack provides a convenient programming model while there is still room for better programmability and higher performance

# Chai Project

- Papers:
  - **Analysis and Modeling of Collaborative Execution Strategies for Heterogeneous CPU-FPGA Architectures.** *ICPE'19*. (this work)
  - **Collaborative Computing for Heterogeneous Integrated Systems.** *ICPE'17 Vision Track*.
  - **Chai: Collaborative Heterogeneous Applications for Integrated-architectures.** *ISPASS'17*.
- **Chai Benchmark Suite:**
  - Website: [chai-benchmarks.github.io](http://chai-benchmarks.github.io)
  - Code: [github.com/chai-benchmarks/chai](https://github.com/chai-benchmarks/chai)
  - Online Forum: [groups.google.com/d/forum/chai-dev](https://groups.google.com/d/forum/chai-dev)



# Analysis and Modeling of Collaborative Execution Strategies for Heterogeneous CPU-FPGA Architectures

Sitao Huang<sup>1</sup>, Li-Wen Chang<sup>2</sup>, Izzat El Hajj<sup>3</sup>, Simon Garcia De Gonzalo<sup>1</sup>, Juan Gómez-Luna<sup>4</sup>, Sai Rahul Chalamalasetti<sup>5</sup>, Mohamed El Hadedy<sup>6</sup>, Dejan Milojicic<sup>5</sup>, Onur Mutlu<sup>4</sup>, Deming Chen<sup>1</sup>, Wen-mei Hwu<sup>1</sup>

<sup>1</sup>University of Illinois at Urbana-Champaign, USA

<sup>2</sup>Microsoft, USA

<sup>3</sup>American University of Beirut, Lebanon

<sup>4</sup>ETH Zürich, Switzerland

<sup>5</sup>Hewlett Packard Labs, USA

<sup>6</sup>California State Polytechnic University, Pomona, USA

## Thanks!



**I ILLINOIS**

Electrical & Computer Engineering

COLLEGE OF ENGINEERING