



Distributed Perception by Collaborative Robots

Interactive Talk

Ramyad Hadidi*, Jiashen Cao*, Matthew Woodward*,
Michael S. Ryoo**, and Hyesoon Kim*

*Georgia Institute of Technology

Click t ** EgoVid Inc.

**Georgia
Tech**



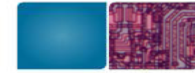
comparch



Robots and Their Environment

2

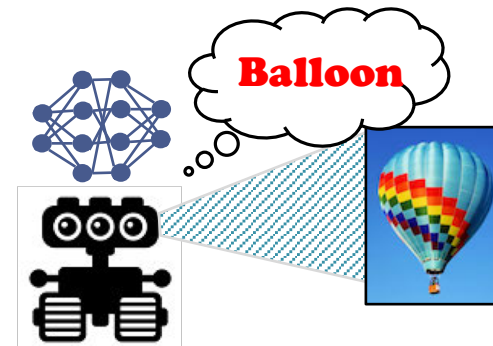
- ▶ Robots need to process lots of raw data in their environment.
 - ▶ Visual, Sounds, Temperature, ...
 - ▶ They need to understand it, to act upon it
 - ▶ For instance:
 - ▶ Drones that study an area after a disaster
 - ▶ Smart security system with lots of cameras
 - ▶ Swarm robots



Deep Learning (DL) and Robots

3

- ▶ How they should process complex raw data?
 - ▶ Use **deep learning**!
 - ▶ It can help in many tasks:
 - ▶ Object detection
 - ▶ Scene recognition
 - ▶ Action recognition
 - ▶ Speech recognition



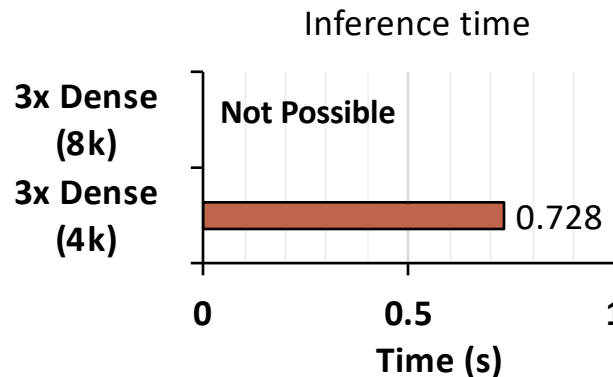
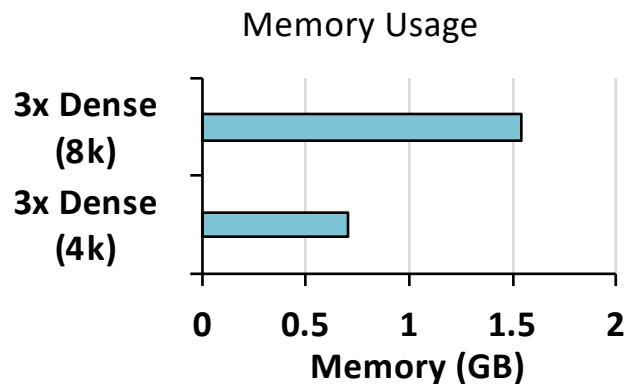


DL Computation is Heavy

4

- ▶ But DL models are computationally intensive and resource hungry specially for cheap robots.

An example of 3x dense layers on resource constrained device:





DL Computation is Heavy

- ▶ But DL models are computationally

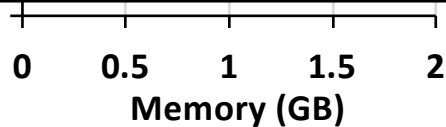
inter

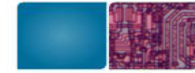
It is difficult to execute DL models on all kinds of robots.

An example

of robots because:

- 1) Usually good models have large memory footprint.
- 2) For a robot, latency for single inference is important.

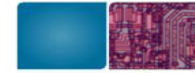




DL Computation is Heavy

6

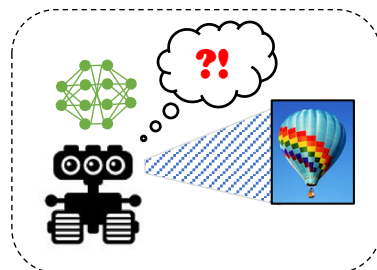
- ▶ So robots need the result fast and in **real time!**
- ▶ Then how resource-constrained robots can use DL to understand their surroundings?



Let's Collaborate

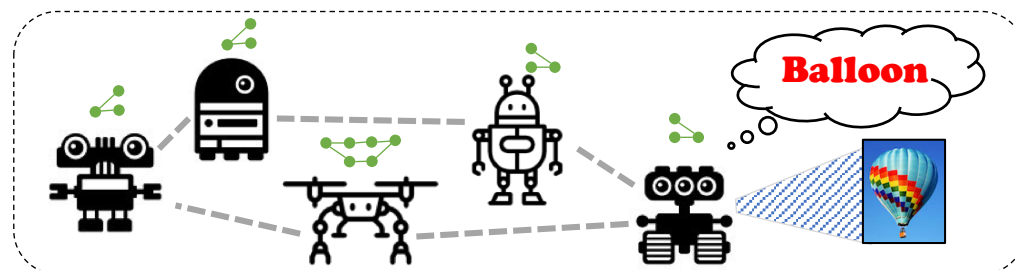
7

- ▶ Usually, many cheap robots share their environment.
- ▶ Not all robots need to perform computations at same time.
- ▶ So What if they share their knowledge and help each other?



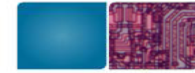
Computation Domain

(a) Single Robot



Computation Domain

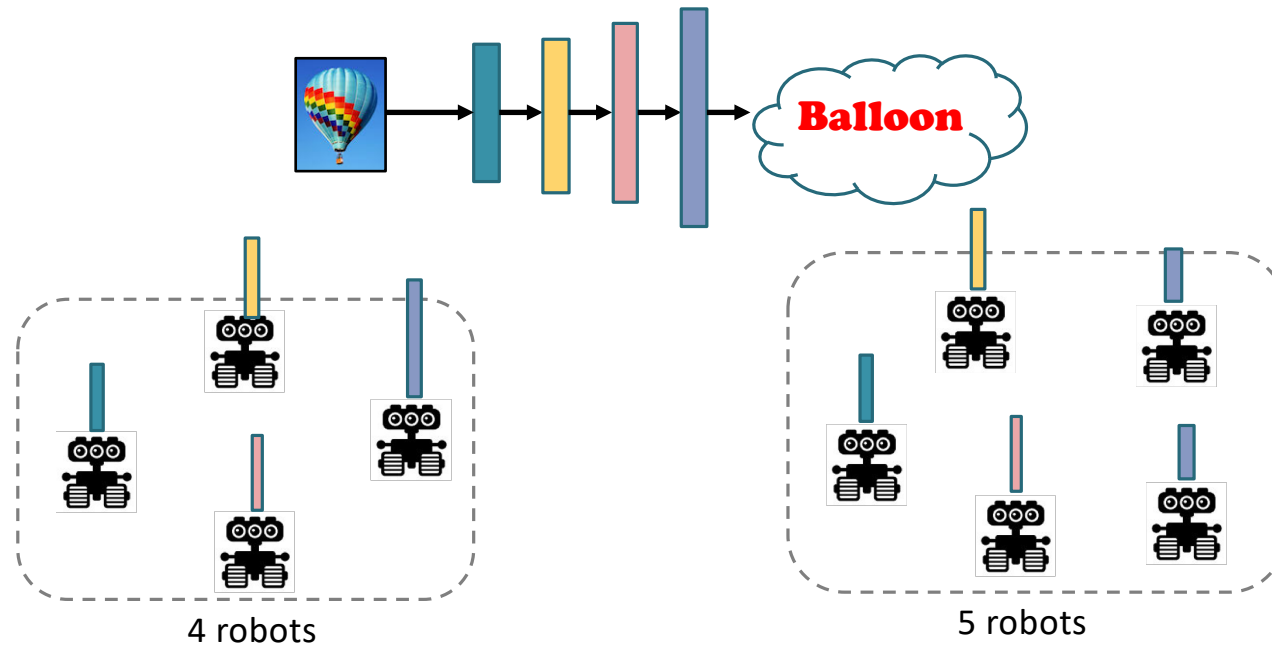
(b) Collaborative Robots



Our Work Overview

8

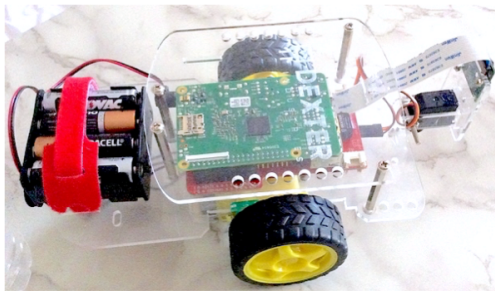
- ▶ We have proposed a technique to efficiently distribute DNN-based recognition.



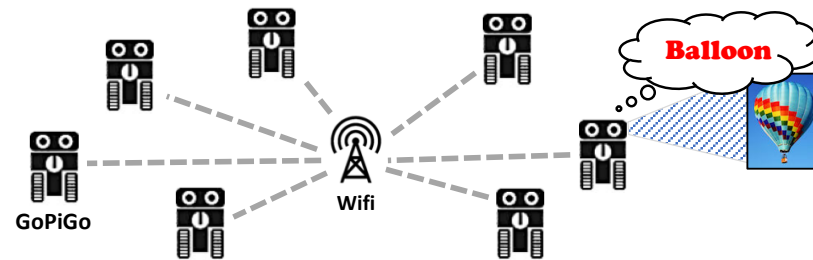


Our Work Overview

- ▶ We proposed an algorithm for deploying the distributed robot system only with **Raspberry Pis**.

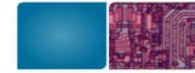


(a) GoPiGo Robot



(b) Our Distributed Robot System

- ▶ We used AlexNet, VGG16, and a video recognition model as example models.



Outline

10

Introduction & Motivation

Data and Model Parallelism

- ▶ Fully Connected and Conv Layers

Distributing DNN

- ▶ Algorithm

System Evaluations

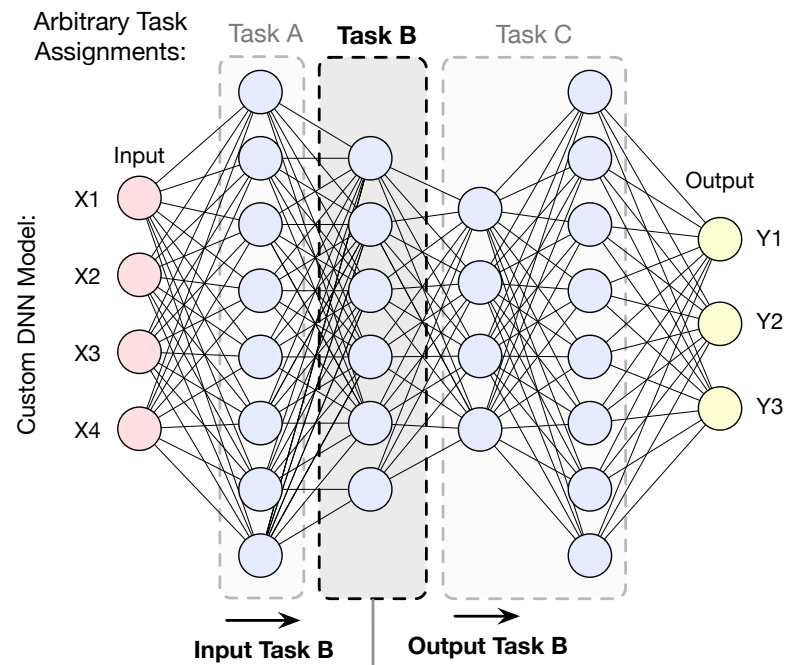
Conclusions



Model & Data Parallelism

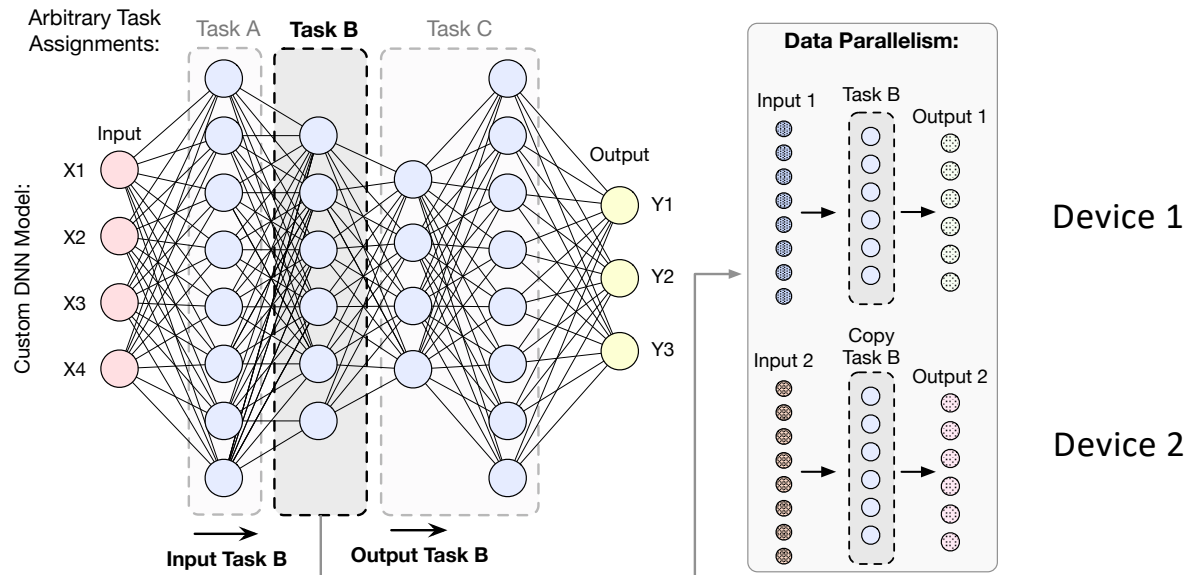
11

Assume a custom DNN model, divided layer by layer:



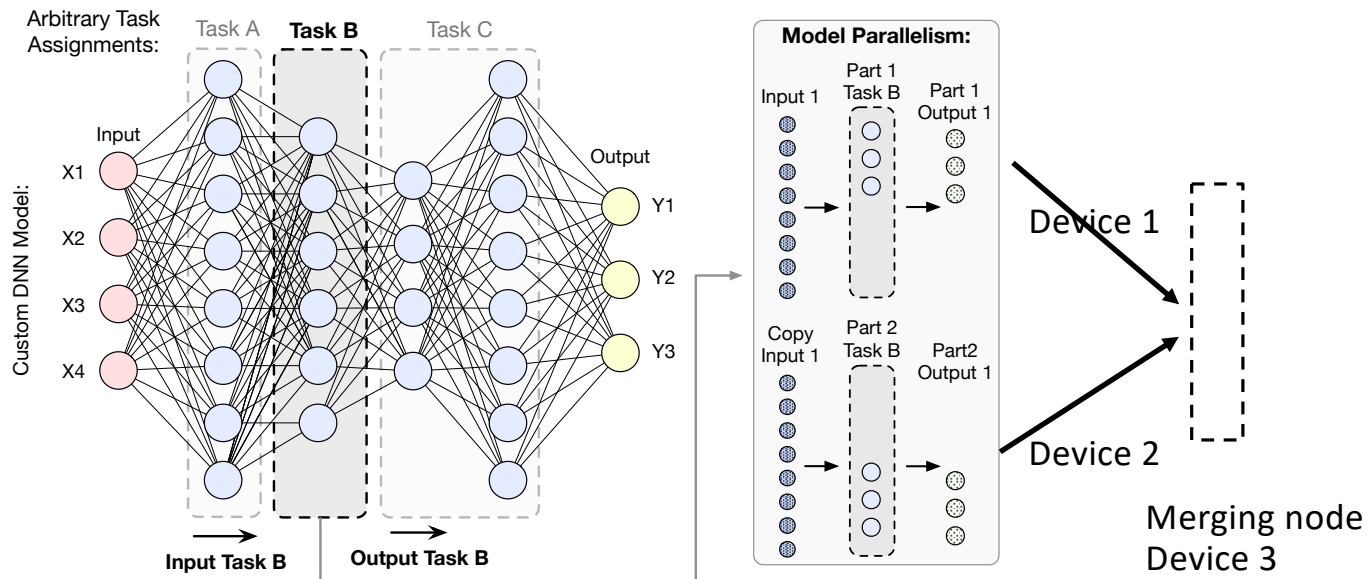


Data Parallelism



Data parallelism is providing the next input to multiple devices in a network.

Model Parallelism



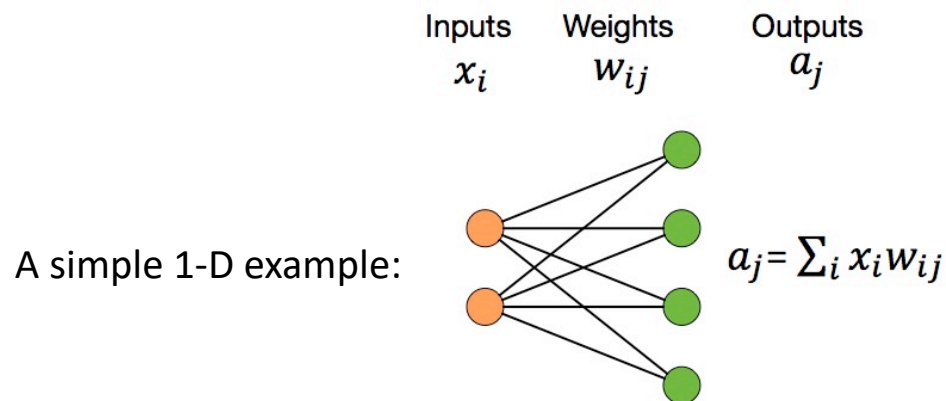
Model parallelism is splitting parts of a given layer or group of layers over multiple devices.

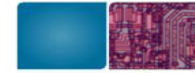


Fully Connected (FC) Layer

14

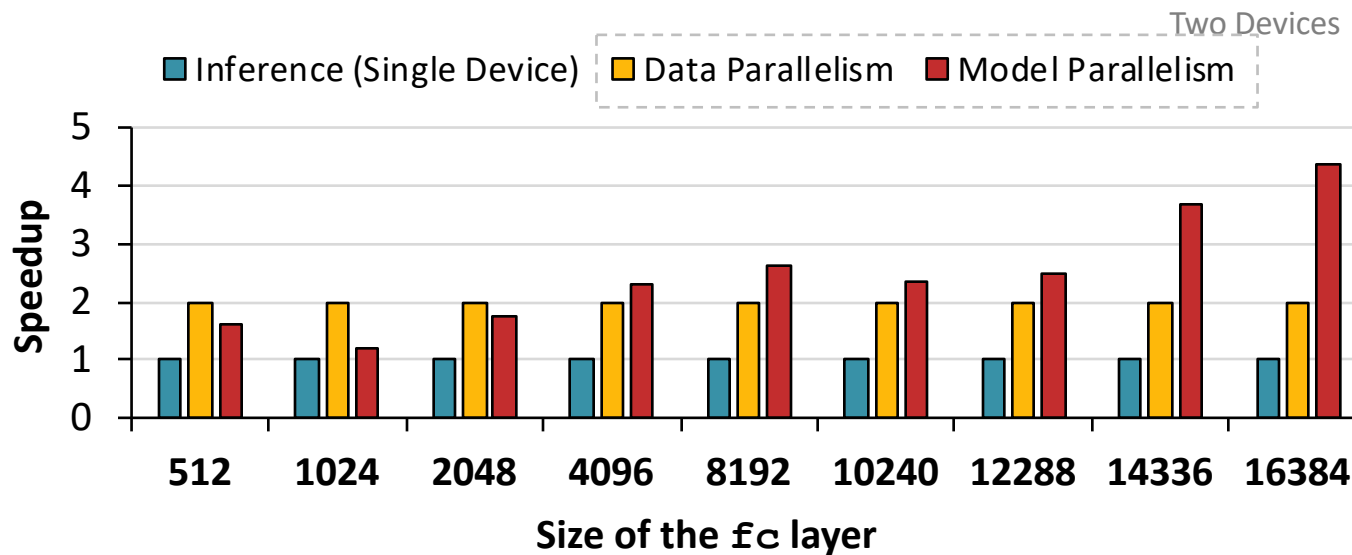
- ▶ Every output element is a summation of weighted inputs
- ▶ Each output element have its own set of weights
- ▶ A matrix multiplication



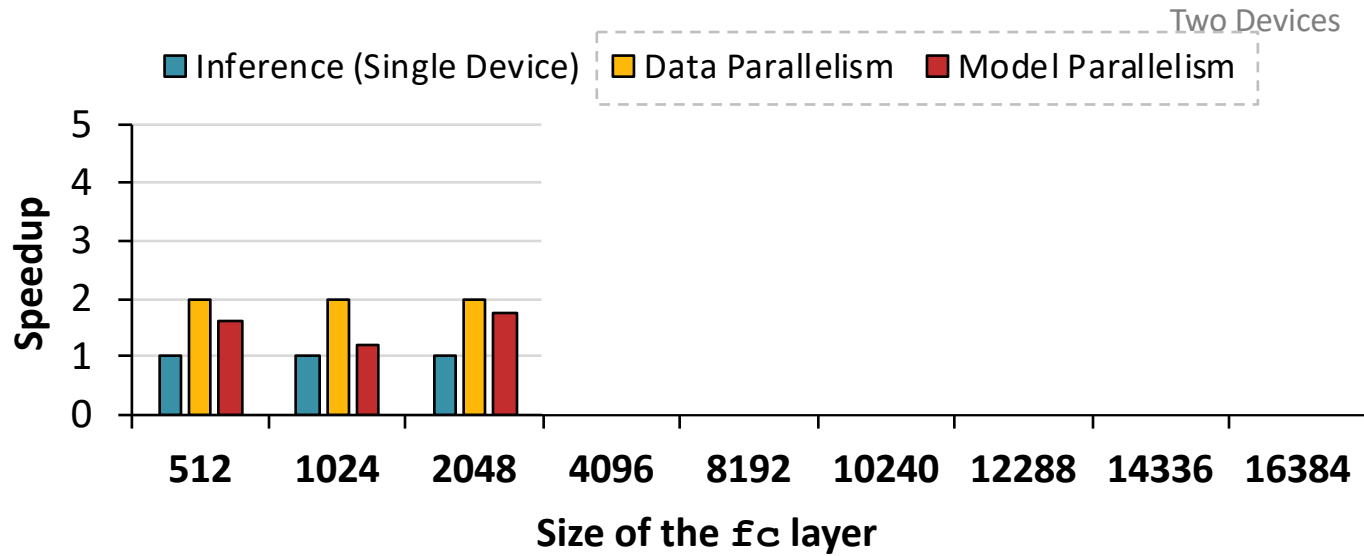


Model & Data Parallelism for FC

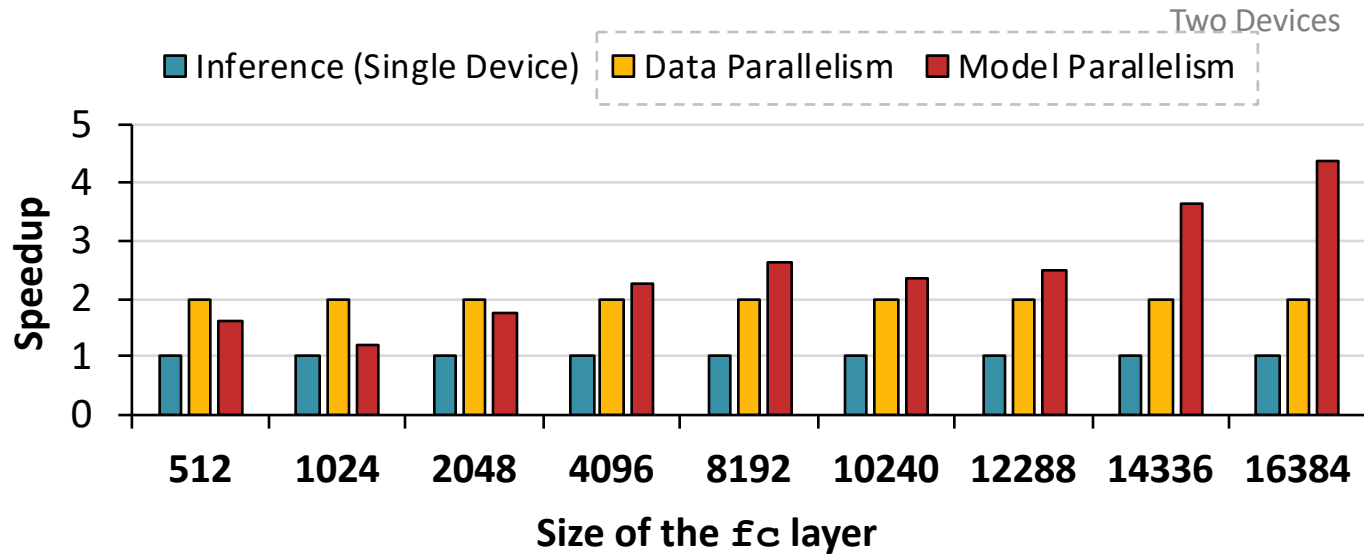
- ▶ Every output element is a summation of weighted inputs
- ▶ Each output element have its own set of weights
- ▶ A matrix multiplication



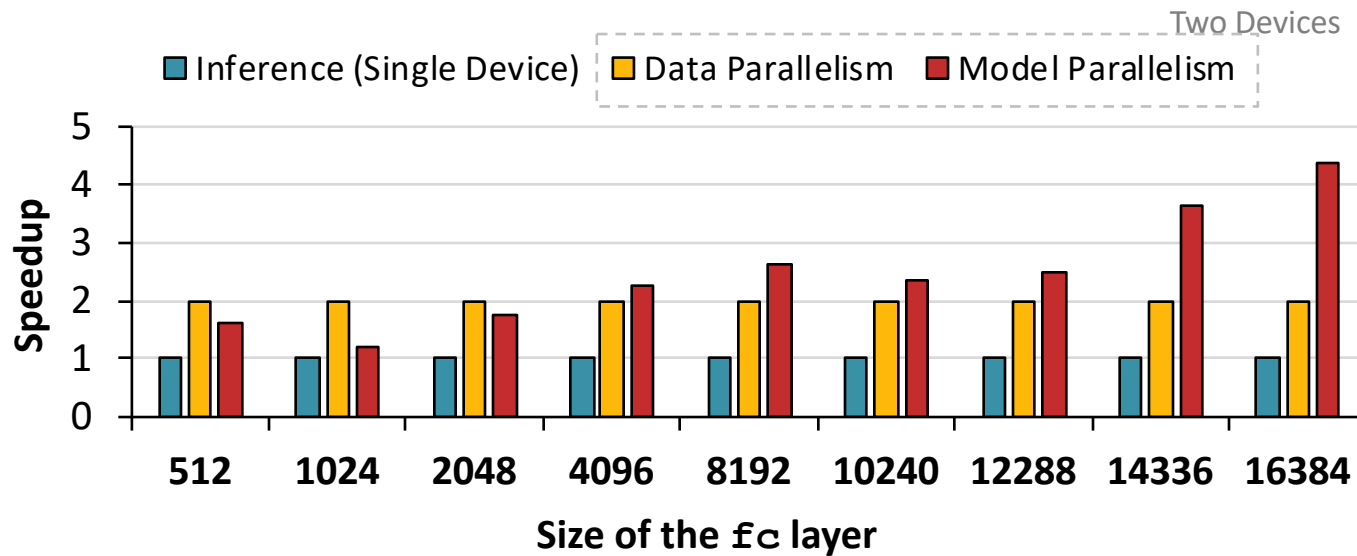
Model & Data Parallelism for FC



Model & Data Parallelism for FC

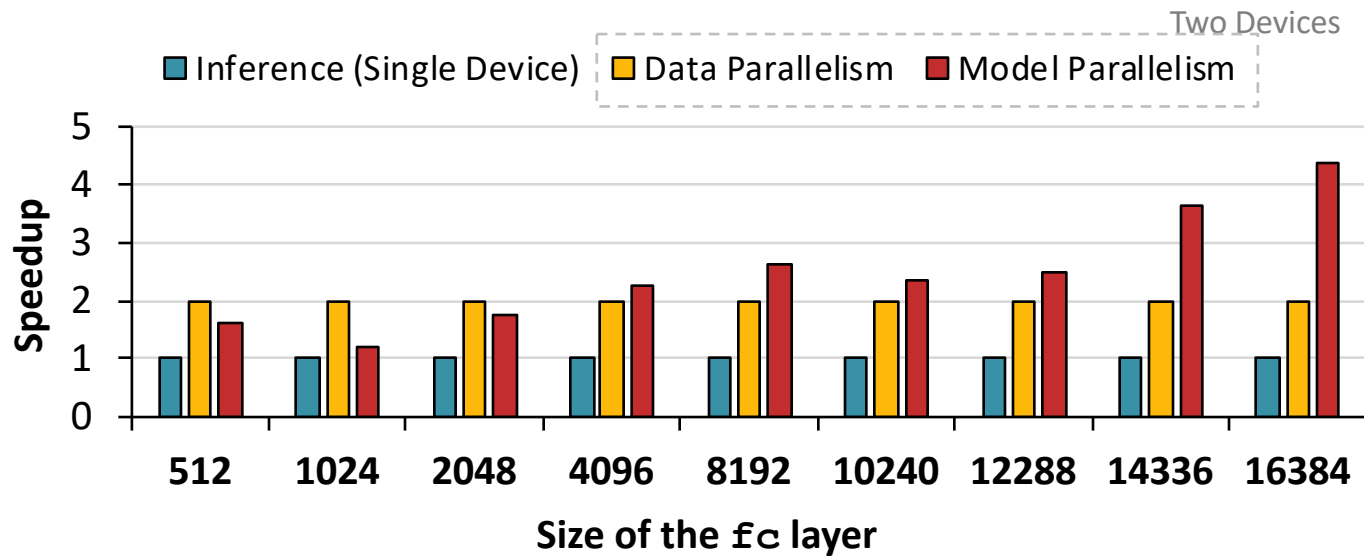


Model & Data Parallelism for FC



Model parallelism reduces memory footprints.
So, we avoid slow hard drive accesses

Model & Data Parallelism for FC



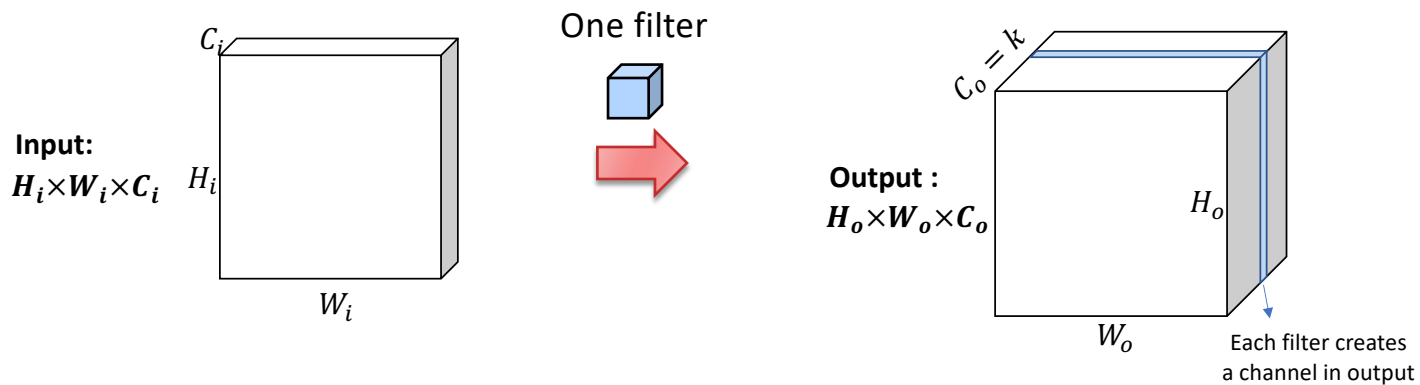
Fully Connected Layers: Either data or model parallelism depending on size of the layer, input, and memory



Convolution Layer

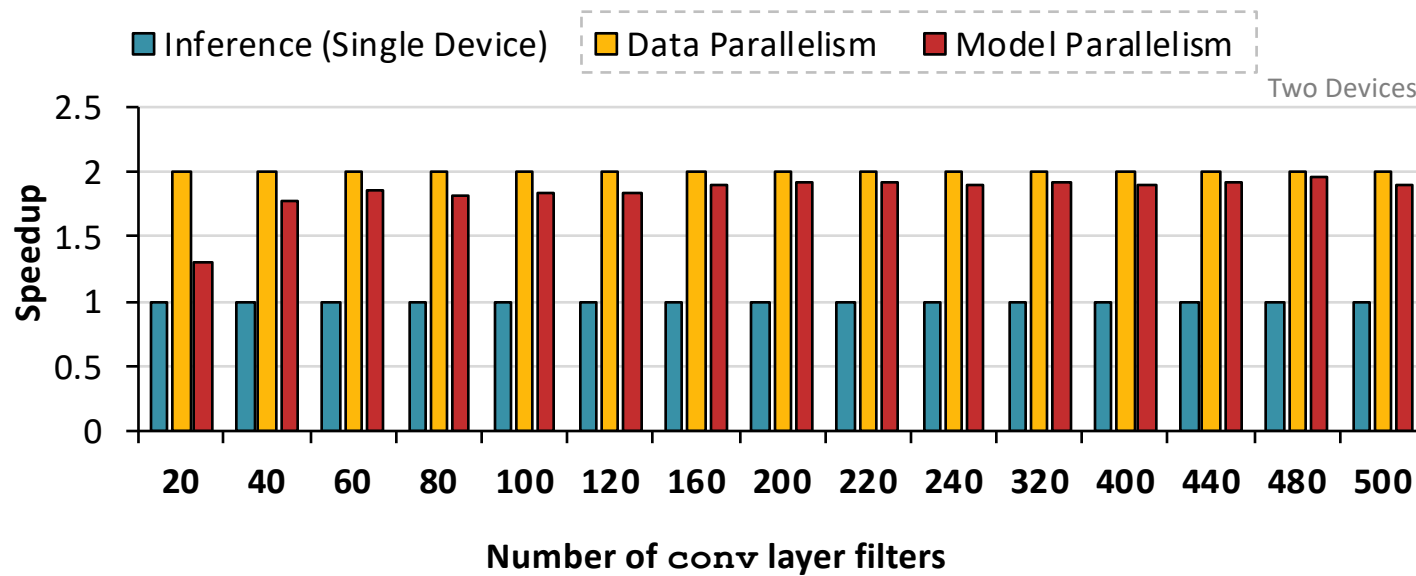
20

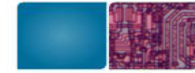
- ▶ Every channel in the output is derived from applying the **same** filter on the input
- ▶ Memory footprint size is smaller → fit into device memory size
- ▶ Model parallelism: split based on filters, and one more merging node at the end



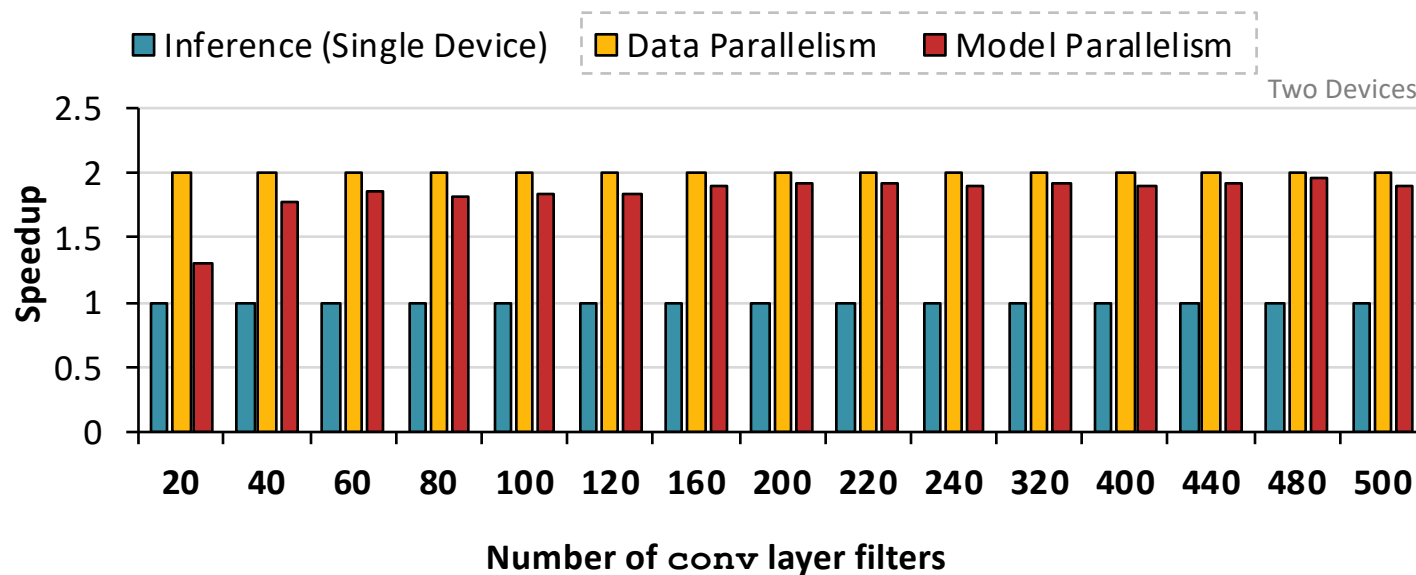


Model & Data Parallelism for Conv

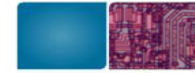




Model & Data Parallelism for Conv



Convolution Layers: Data parallelism is better



Outline

23

Motivation

Background

- ▶ ML Models Overview

Data and Model Parallelism

- ▶ Fully Connected and Conv Layers

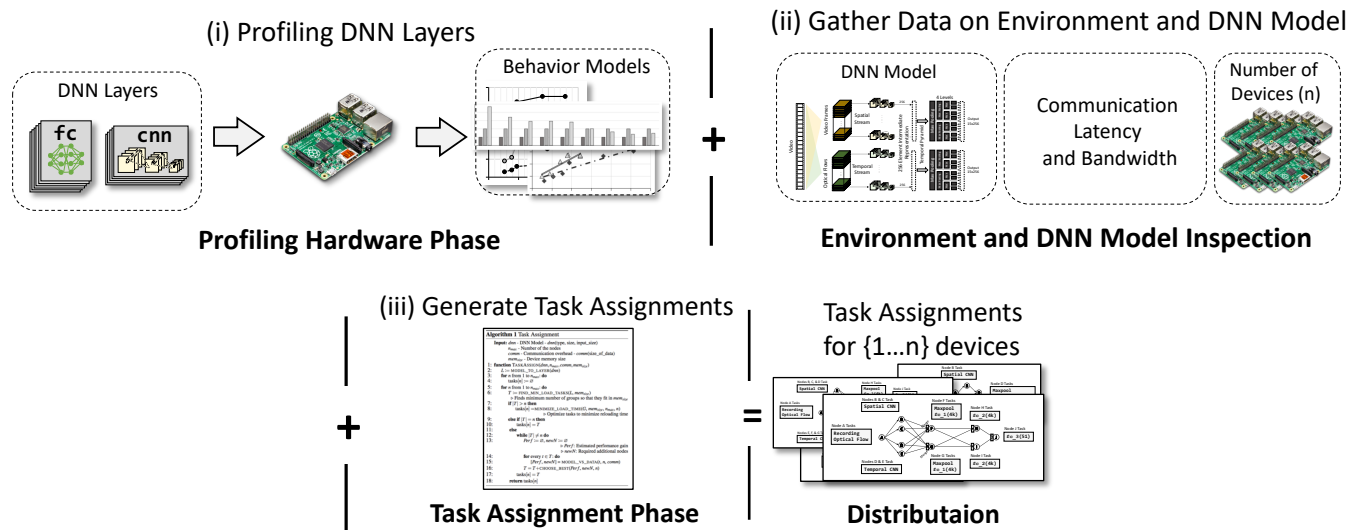
Distributing DNN

- ▶ **Algorithm overview**

System Evaluations

Conclusions

Distributing a DNN Model



Details are in the paper



Outline

25

Motivation

Background

- ▶ ML Models Overview

Data and Model Parallelism

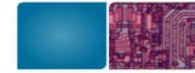
- ▶ Fully Connected and Conv Layers

Distributing DNN

- ▶ Algorithm overview

System Evaluations

Conclusions



Software & Hardware

26

Software:

- ▶ Keras 2.1
 - ▶ With Tensorflow backend for Raspberry Pis
 - ▶ With Tensorflow-GPU backend for TX2
- ▶ Apache Avro for procedure call (RPC) and data serialization



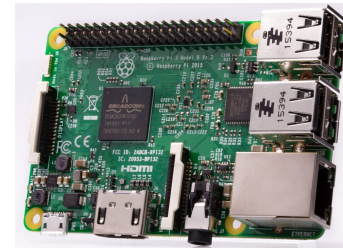


Hardware Overview

27

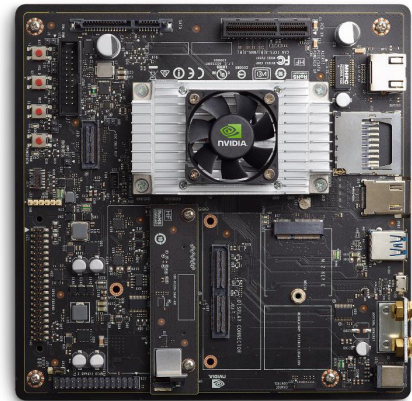
Raspberry PI 3:

- ▶ Cheap and accessible platform
- ▶ Connected via a Wifi router
- ▶ No GPU
- ▶ \$40



Nvidia Jetson TX2:

- ▶ High-end embedded platform
- ▶ Has a GPU
- ▶ \$600

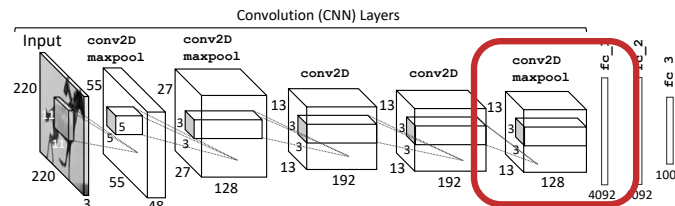
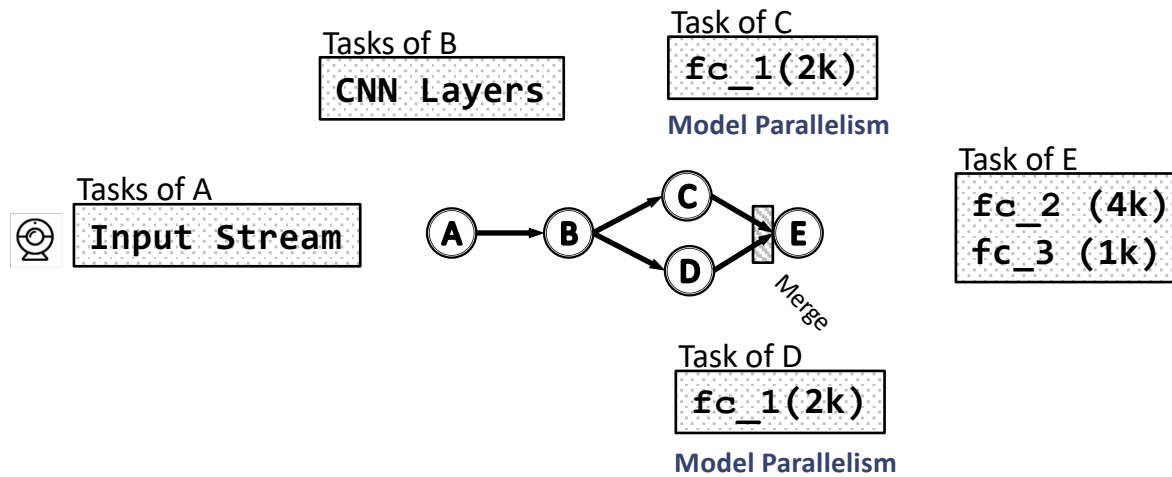


Moreover, we measured whole system power with a power analyzer

AlexNet Distribution I

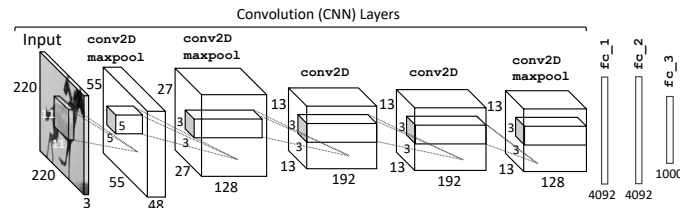
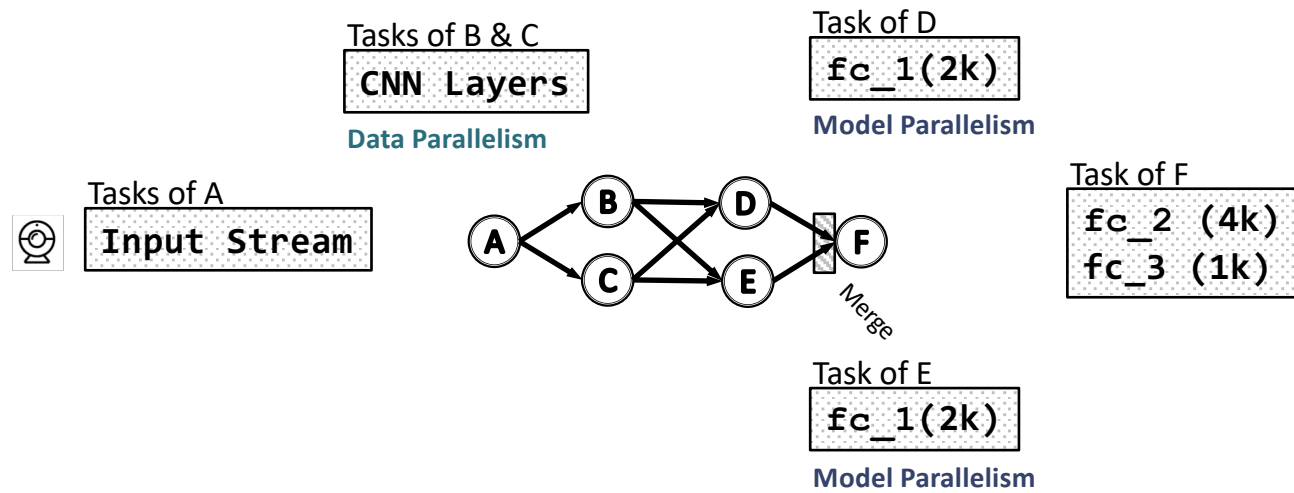


Five-device system:

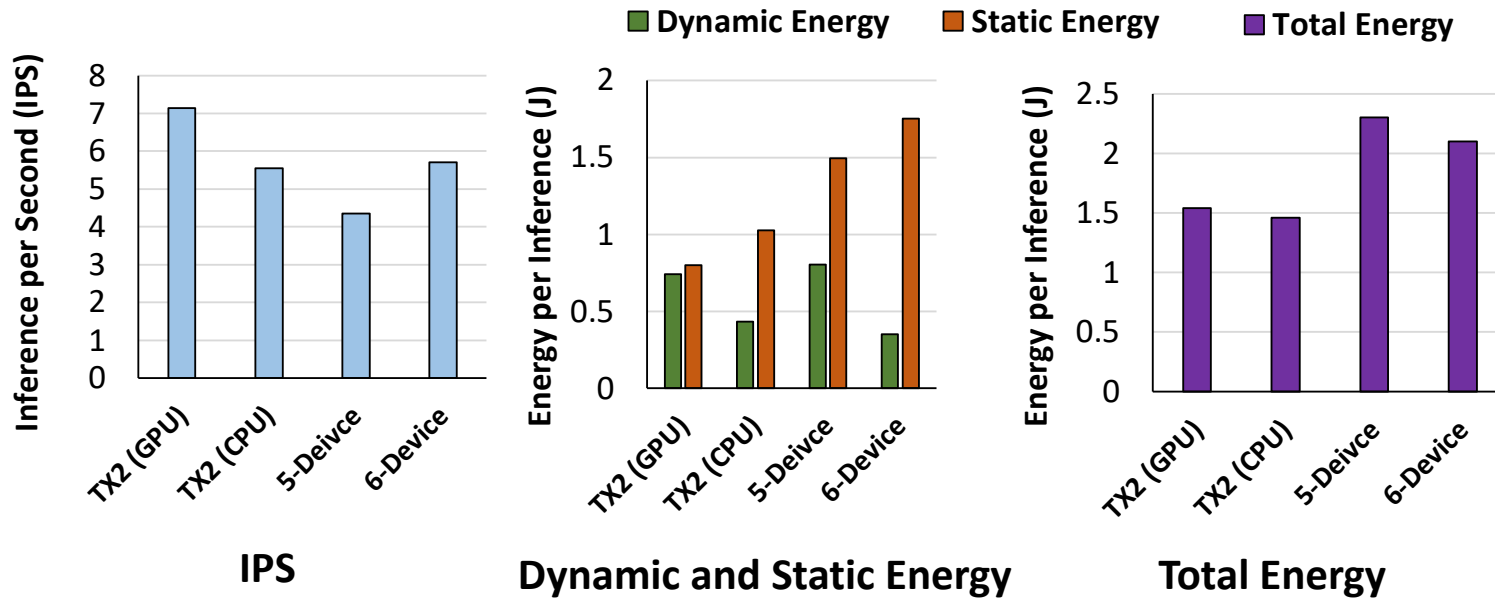


AlexNet Distribution II

Six-device system:



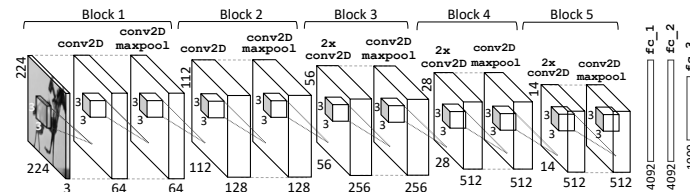
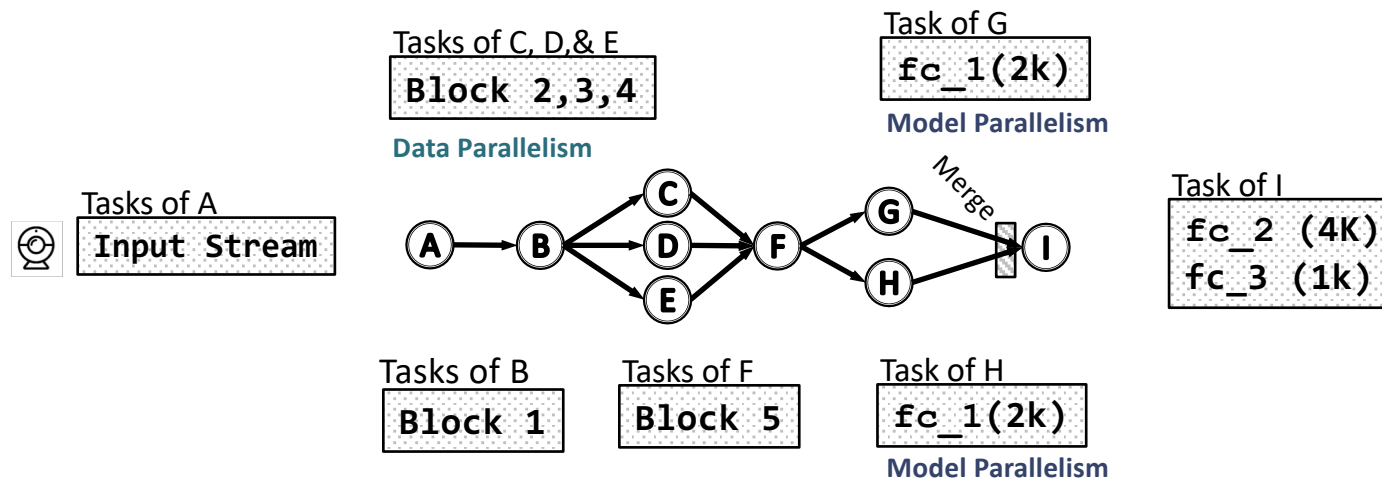
AlexNet Results



Comparable IPS with TX2 (-30%)
Lower dynamic energy consumption

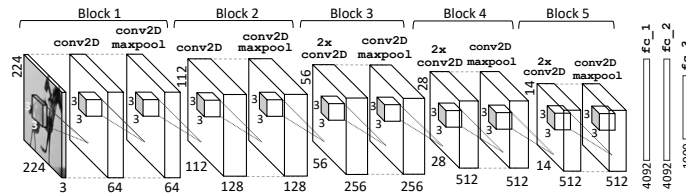
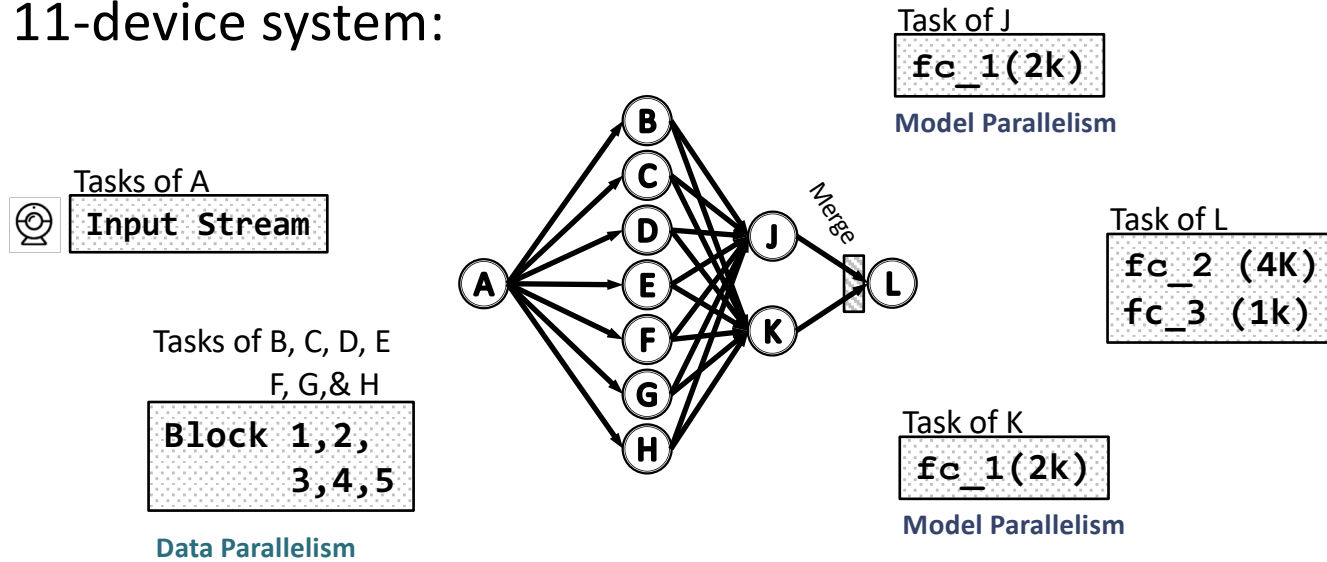
VGG16 Distribution I

Nine-device system:

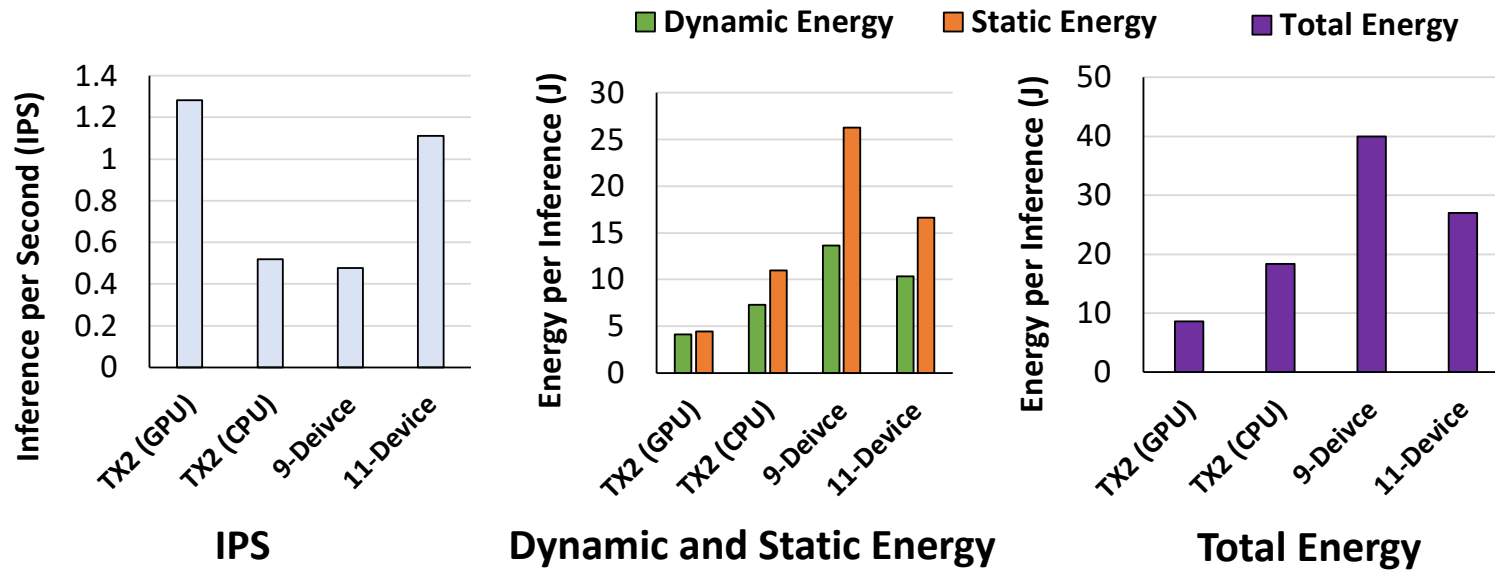


VGG16 Distribution II

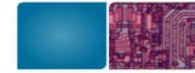
11-device system:



VGG16 Results



Comparable IPS with TX2 (-15%)
We achieve 2.3x speedup, by reassigning CNN blocks



Outline

34

Motivation

Background

- ▶ ML Models Overview

Data and Model Parallelism

- ▶ Fully Connected and Conv Layers

Distributing DNN

- ▶ Algorithm overview

System Evaluations

Conclusions



Conclusions

35

- ▶ We used a farm of Raspberry Pis for DNN processing
- ▶ Our technique achieves acceptable real-time performance
 - ▶ Compared to TX2 CPU, we achieve similar performance with 6 robots for AlexNet
 - ▶ 11 robots for VGG-16 compared to TX2 GPU

Future Work:

- ▶ Study the robustness of such systems
- ▶ Apply our technique to more DNN models
- ▶ Implement our model on distributed robot systems







Backup Slides



Layers of ML Models

39

- ▶ Convolution: Applies several filters to the input
 - ▶ Compute bound, more locality
- ▶ Activation: Introduces non-linearity
 - ▶ e.g., ReLU $f(x) = \max(0, x)$, not compute intensive
- ▶ Fully Connected (Dense)
 - ▶ i.e., matrix multiplication, bandwidth bound
- ▶ Pooling
 - ▶ Reduces dimensions, simple doing max, average, and ... on a subset of input

Image Recognition Models (I)



► Single-stream AlexNet

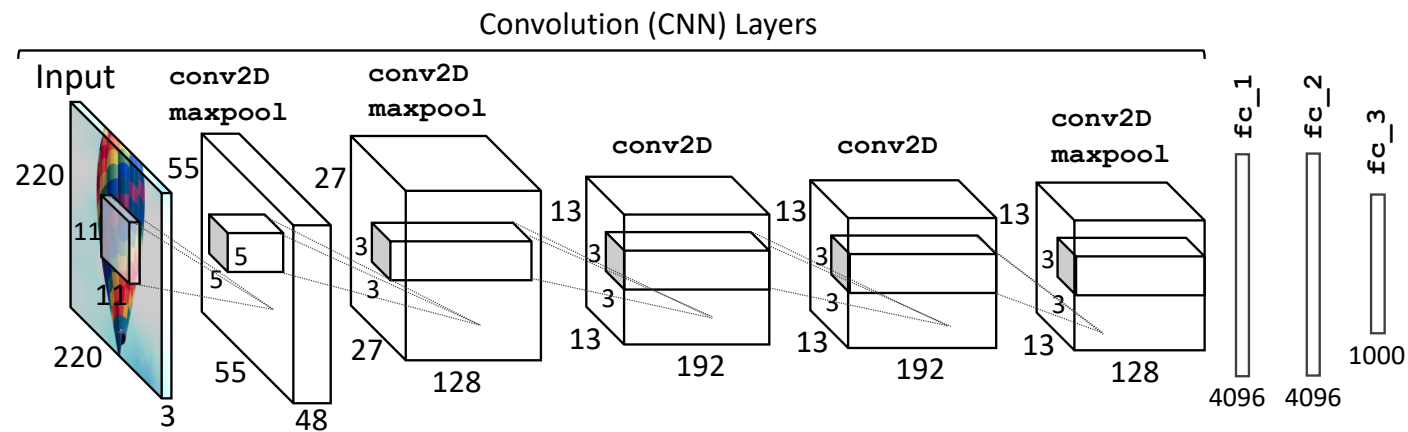
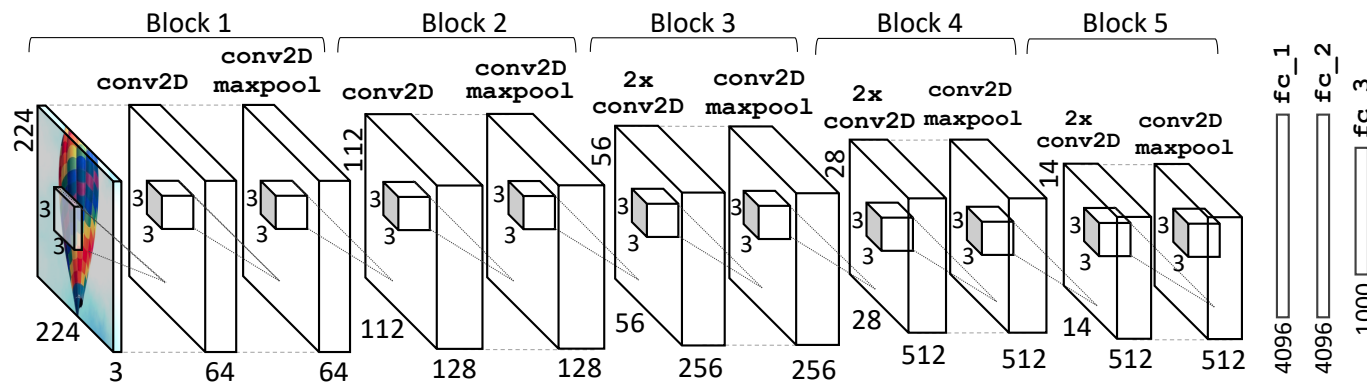


Image Recognition Models (II)

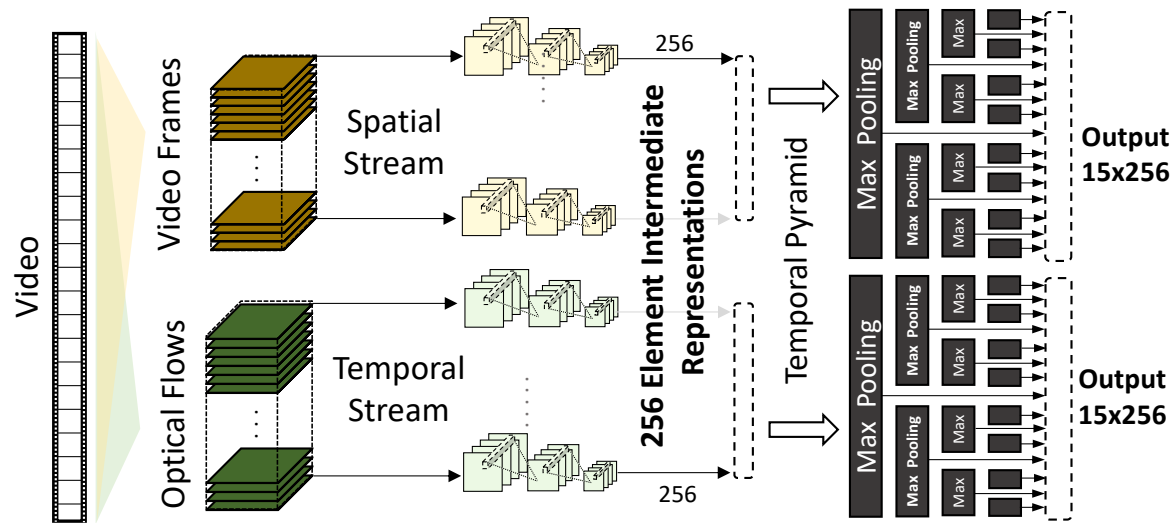


► VGG16



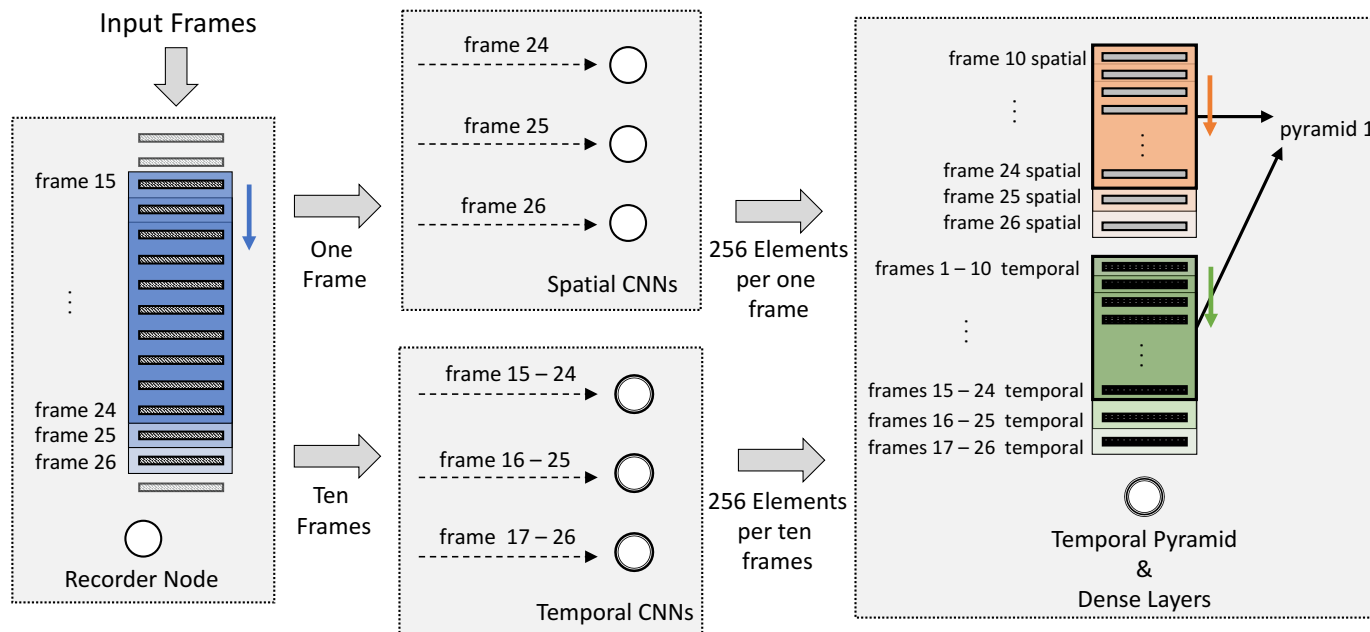
Vide Recognition Model

- ▶ i.e., Action recognition model
- ▶ Based on the two-stream model by Ryoo et al.^[1]



M. S. Ryoo, K. Kim, and H. J. Yang, "Extreme Low Resolution Activity Recognition with Multi-Siamese Embedding Learning," in *AAAI'18*. IEEE, Feb. 2018.

Sliding Window





Algorithm

Algorithm 1 Task Assignment Algorithm.

```
1: function TASKASSIGNMENT( $dnn, n_{max}, comm, mem_{size}$ )
   Inputs:  $dnn$ : DNN model in form of layers[type, size,  $input_{size}, \beta$ ]
             $n_{max}$ : Maximum number of the devices
             $comm$ : Communication overhead model ( $comm(size_{data})$ )
             $mem_{size}$ : Device memory size
2:    $L := \text{EXTRACT\_MODEL\_TO\_LAYERS}(dnn)$ 
3:   for  $n$  from 1 to  $n_{max}$ : do
4:      $tasks_{final}[n] := \emptyset$ 
5:   for  $n$  from 1 to  $n_{max}$ : do
6:      $TG, noFit := \text{FIND\_INITIAL\_TASKGROUP}(L, mem_{size})$ 
7:     if  $sizeof(TG) > n$  then
8:        $tasks[n] = \text{COMBINE\_TASKS}(TG, mem_{size}, n_{max}, n)$ 
9:     if  $sizeof(TG) = n$  then
10:       $tasks[n] = TG$ 
11:    if  $sizeof(TG) < n$  or  $noFit == True$  then
12:      while  $sizeof(TG) \neq n$  do
13:         $task_{variant} := \emptyset$ 
14:        for every  $t \in TG$ : do
15:           $[task_{variant}] += \text{PROFILED\_VARIANTS}(t, comm)$ 
16:           $task_{replaced}, task_{new} = \text{SELECT\_LOWEST}([task_{variant}])$ 
17:           $TG = TG - task_{replaced} + task_{new}$ 
18:           $tasks_{final}[n] = TG$ 
19:    return  $tasks_{final}$ 
```



Hardware Overview

Raspberry PI 3:

- ▶ Cheap and accessible platform
- ▶ Connected via a Wifi router
- ▶ No GPU
- ▶ \$40

Table 1: Raspberry PI 3 specification

CPU	1.2 GHz Quad Core ARM Cortex-A53	
Memory	900 MHz 1 GB RAM LPDDR2	
GPU	No GPGPU Capability	
Price	\$35 (Board) + \$5 (SD Card)	
Power Consumption	Idle (No Power Gating)	1.3 W
	%100 Utilization	6.5 W
	Averaged Observed	3 W

Nvidia Jetson TX2:

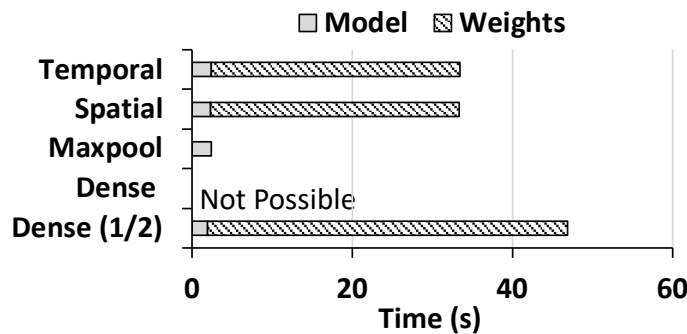
- ▶ High-end embedded platform
- ▶ Has a GPU
- ▶ \$600

Table 2: Nvidia Jetson TX2 specifications

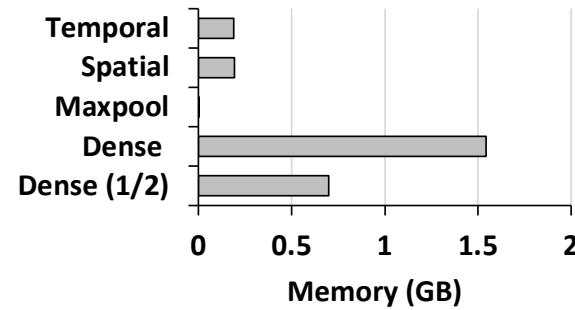
CPU	2.00 GHz Dual Denver 2 + 2.00 GHz Quad Core ARM Cortex-A57	
Memory	1600 MHz 8 GB RAM LPDDR4	
GPU	Pascal Architecture - 256 CUDA Core	
Total Price	\$600	
Power Consumption	Idle (Power Gated)	5 W
	%100 Utilization	15 W
	Averaged Observed	9.5 W

Moreover, we measured whole system power with a power analyzer

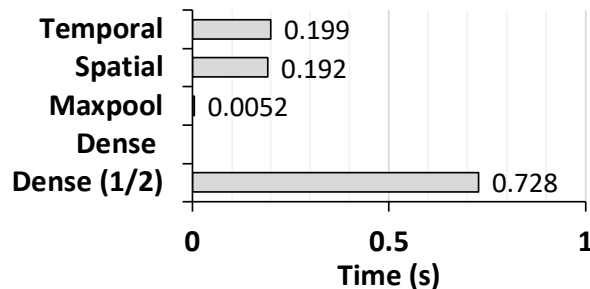
Video Recognition on Single PI



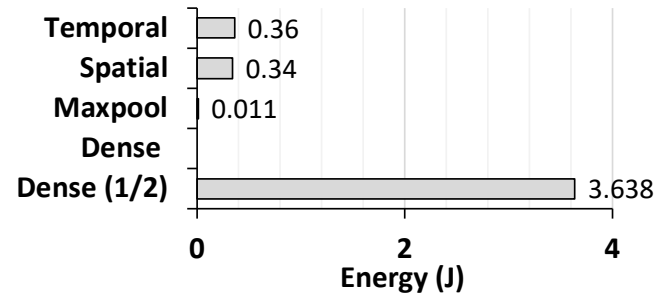
(a) Loading Time



(b) Memory Usage

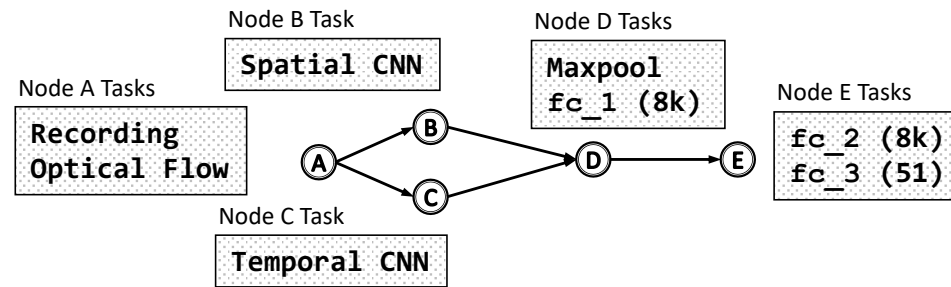


(c) Inference Time

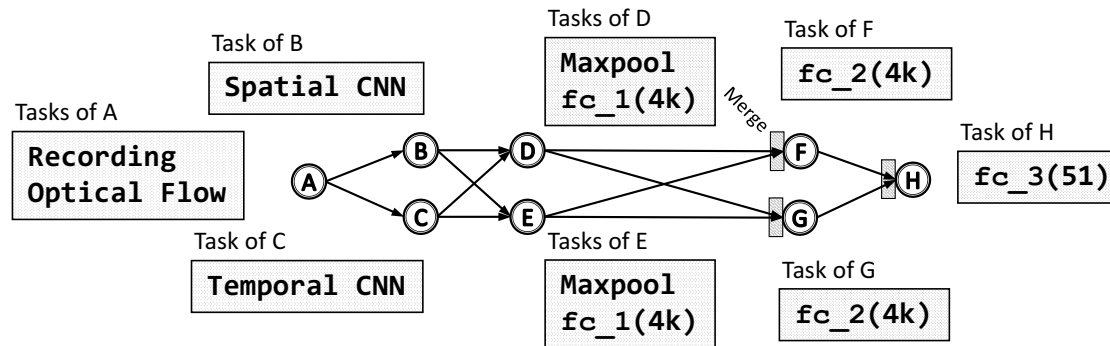


(d) Energy Per Inference

Video Recognition Distributions (I)

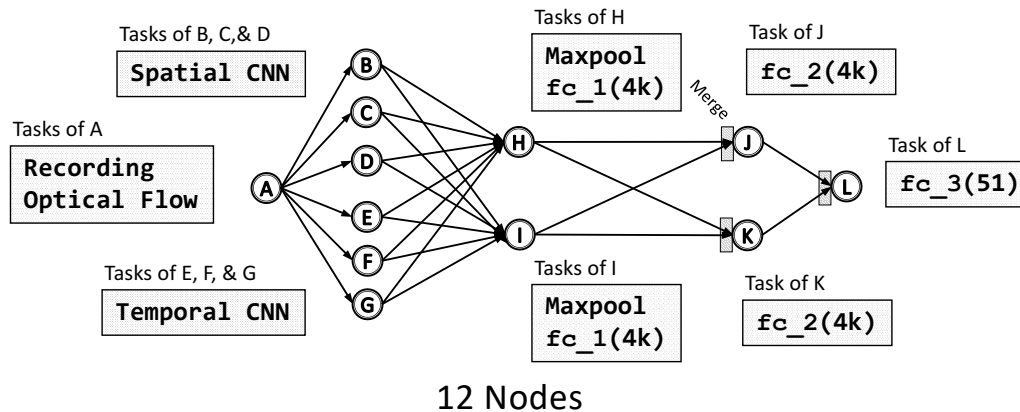
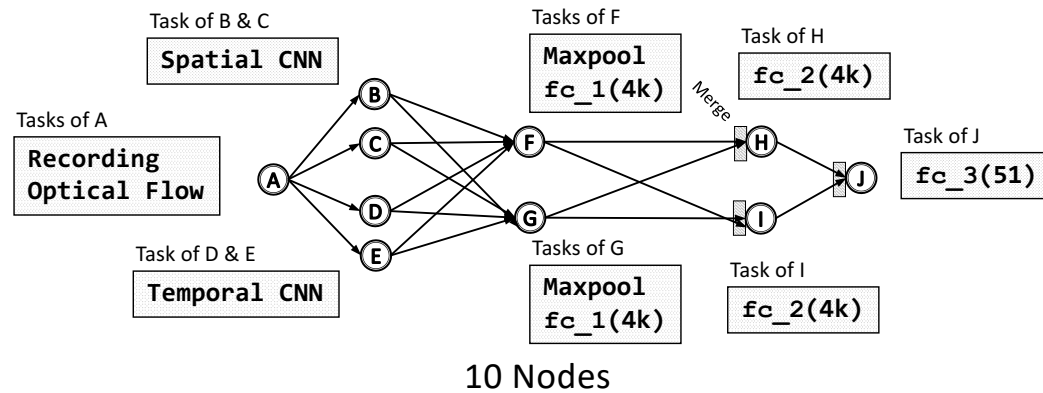


5 Nodes



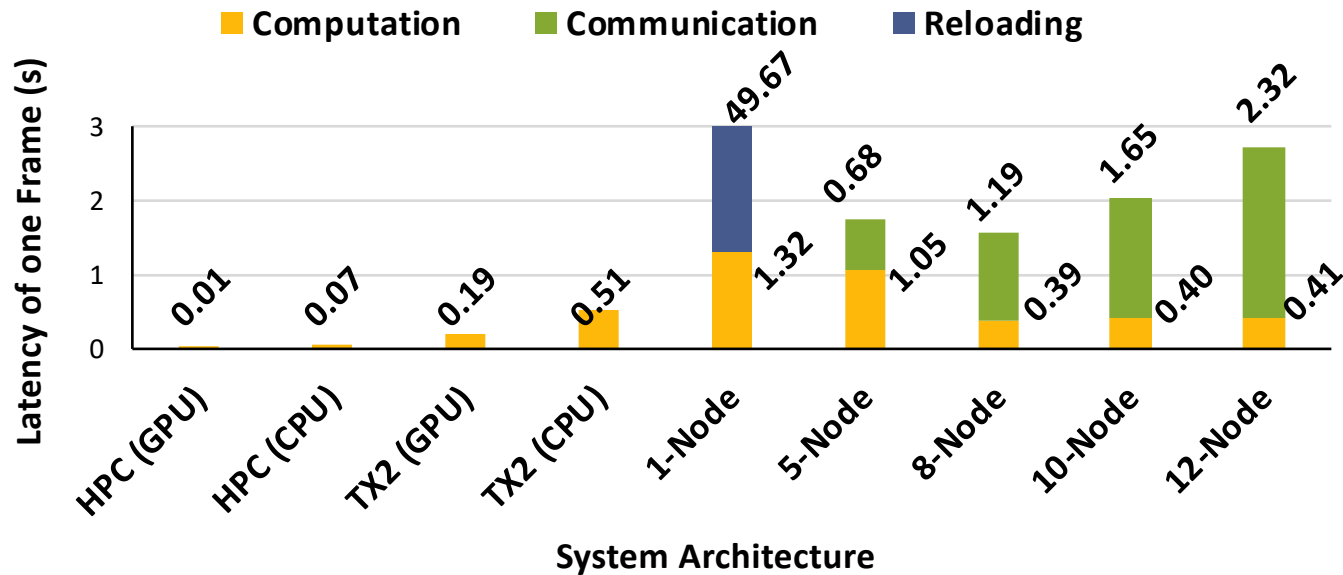
8 Nodes

Video Recognition Distribution (II)



Video Recognition Results (2)

Latency of one Frame (Seconds)



Video Recognition Results (3)



Energy:

