

# Real-Time Image Recognition Using Collaborative IoT Devices

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

\*Georgia Institute of Technology

\*\*Indiana University; EgoVid Inc.

**Georgia  
Tech**



**compArch**



# Prevalence of IoT Devices

2

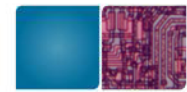
Internet of Things (IoT) devices are everywhere

- ▶ Smart Locks, Smart Sprinklers, Smart Plugs, Smart Baby Monitors, Smart Cookers, Smart Thermostats, Smart Mirrors, Smart Cleaners, and Smart Refrigerators



Many of which generate/capture abundance of real-time **raw** data such as images.

<https://www.pentasecurity.com/blog/10-smartest-iot-devices-2017/>



# How to Process IoT data?

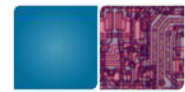
---

3

Advancements of deep neural networks (DNN) provides many high-accuracy solutions to previously impossible tasks:

- ▶ Image Recognition
- ▶ Face Recognition
- ▶ Video (Action Recognition)
- ▶ Voice Recognition

Performing these tasks in **real-time** requires high computational power.

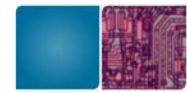


# Where to Process (I)

---

4

- ▶ **(Option A)** Use the individual IoT device
  - ▶ Limited energy (e.g., battery powered)
  - ▶ Limited compute power
  - ▶ So, unable to meet time constraints
  
- ▶ **(Option B)** Offload to Cloud
  - ▶ Such as Voice recognition service of Apple's Siri, Amazon's Echo, Microsoft's Cortana, and Google Home
  - ▶ Any problem?



# Where to Process (II)

5

## (Option B) Cloud processing is promising but:

- ▶ Not Scalable
  - ▶ More traffic, data, and storage
  - ▶ IoT devices outnumbered world population in 2017
- ▶ Privacy and Security
  - ▶ Voice recognition? Big Brother's spying devices in the novel 1984
  - ▶ Multiple layers: Network security, encryption, and etc.
- ▶ Quality of Service (QoS) and Reliability
  - ▶ We have a tight timing constraint for real-time recognition

F.Biscotti et al., "The Impact of the Internet of Things on Data Centers," Gartner Research, vol. 18, 2014.



# Where to Process (III)

6

**(Option C)** What if we could harvest the aggregated computational power of local IoT devices?

- ▶ At a given time, not all devices are fully utilized





# Collaborative IoT Devices

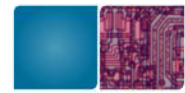
7

**(Option C)** We study such collaboration between IoT devices in our paper, *Musical Chair*.

- ▶ Our performance metric: Inferences per second
- ▶ We use same models, so we have same accuracy

In this work, we showcase the application of Musical Chair for Image recognition models on a farm of Raspberry Pis

Hadidi et al. "Musical Chair: Efficient Real-Time Recognition Using Collaborative IoT Devices." *arXiv preprint arXiv:1802.02138* (2018).



# Outline

---

8

Motivation

Musical Chair

- ▶ Data and Model Parallelism

Hardware and Software Overview

System Evaluations

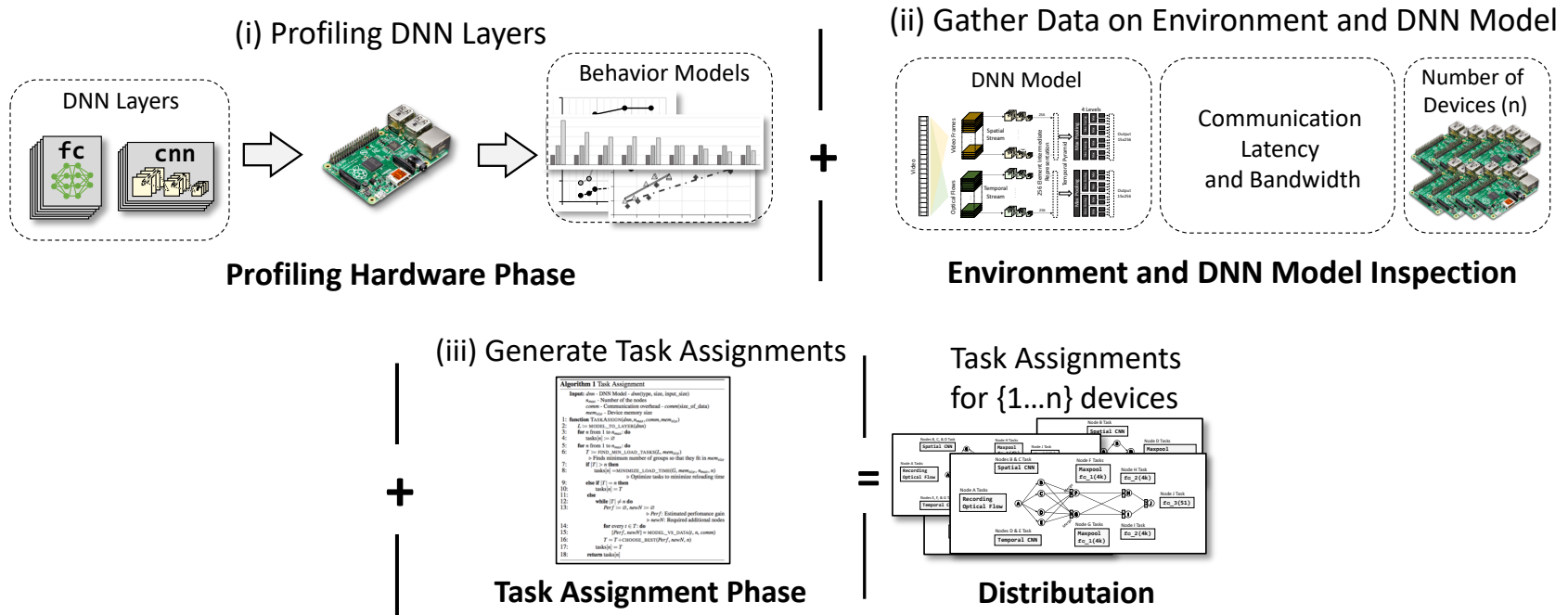
Conclusion



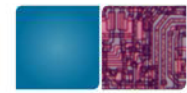


# Musical Chair

Musical Chair is a technique for distributing DNN computations over multiple IoT devices.



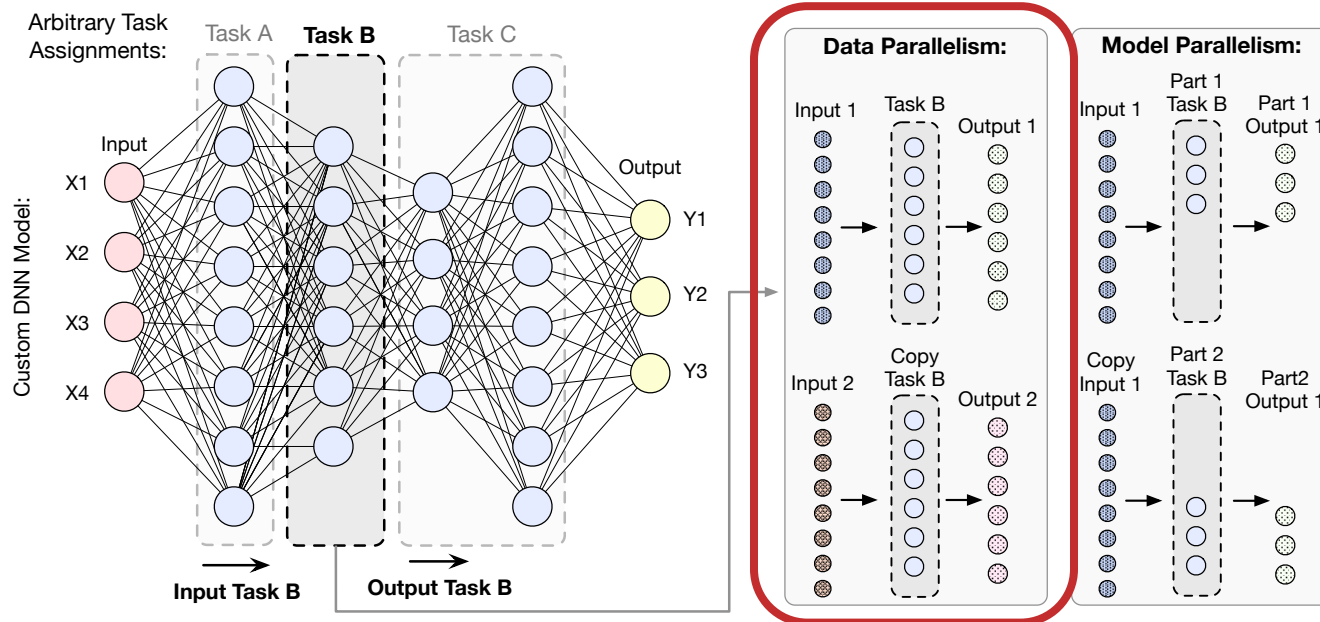
Hadidi et al. "Musical Chair: Efficient Real-Time Recognition Using Collaborative IoT Devices." *arXiv preprint arXiv:1802.02138* (2018).



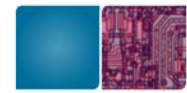
# Model & Data Parallelism

10

Two forms of distribution:

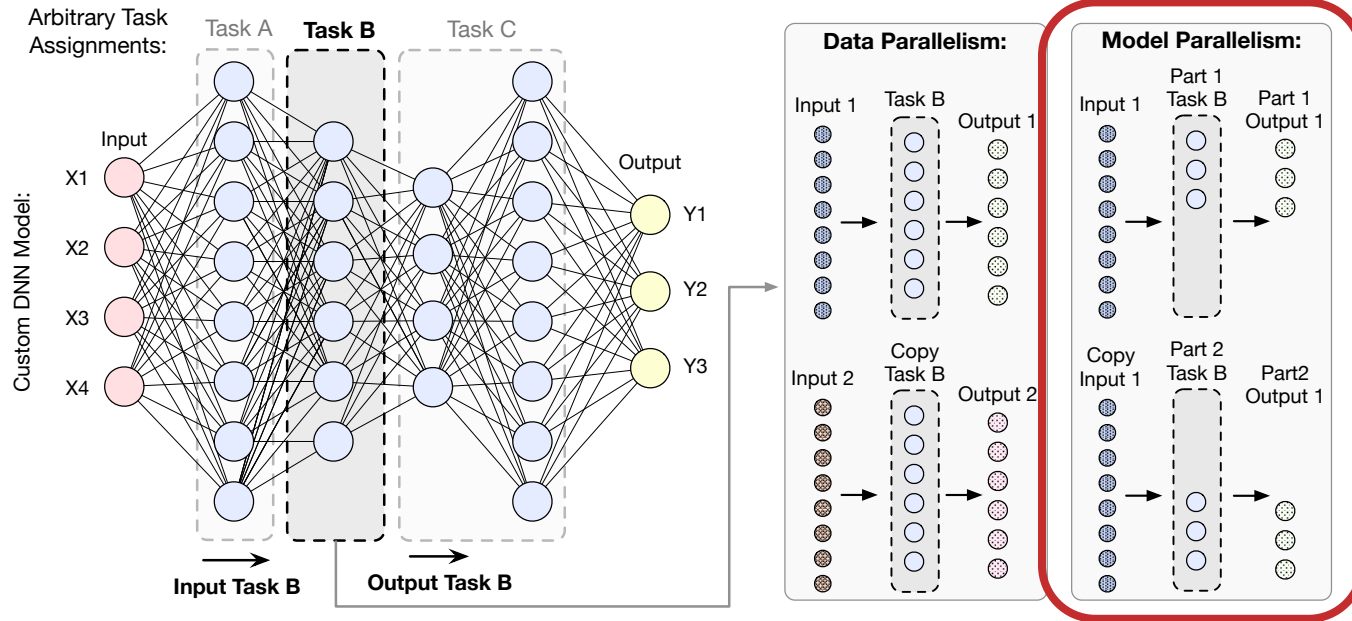


Data parallelism is providing the next input to multiple devices in a network.



# Model & Data Parallelism

Two forms of distribution:



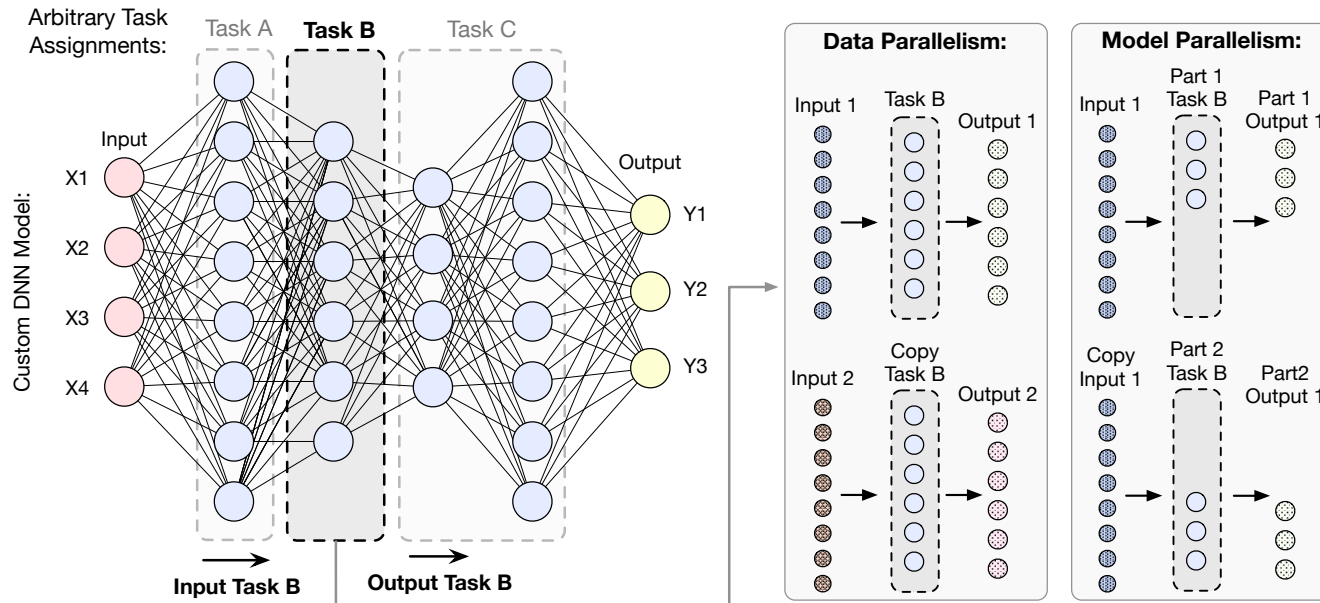
Model parallelism is splitting parts of a given layer or group of layers over multiple devices.



# Model & Data Parallelism

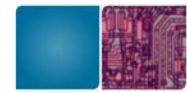
12

Two forms of distribution:



**Convolution Layers:** Mostly data parallelism

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

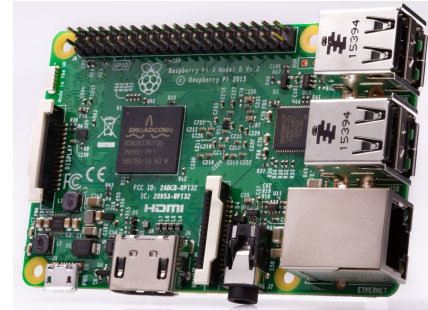


# Hardware Overview

13

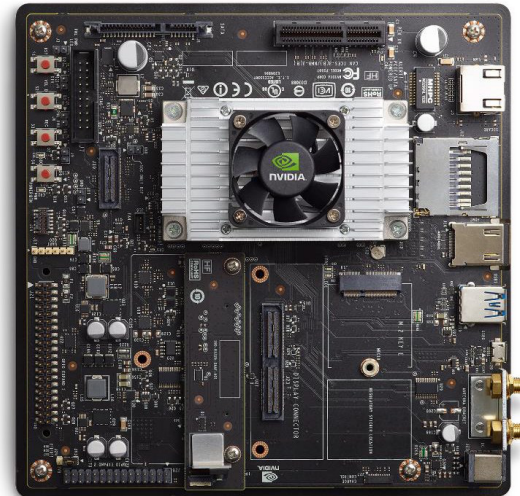
## Raspberry PI 3:

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



## Nvidia Jetson TX2:

- ▶ High-end embedded platform
- ▶ Has a GPU



Moreover, we measured whole system power with a power analyzer



# Software Overview

14

## Dependencies:

- ▶ Ubuntu 16.04
- ▶ Keras 2.1
  - ▶ With Tensorflow backend for Raspberry Pis
  - ▶ With Tensorflow-GPU backend for TX2
- ▶ Apache Avro for procedure call and data serialization



## Image Recognition Models:

- ▶ AlexNet
- ▶ VGG16



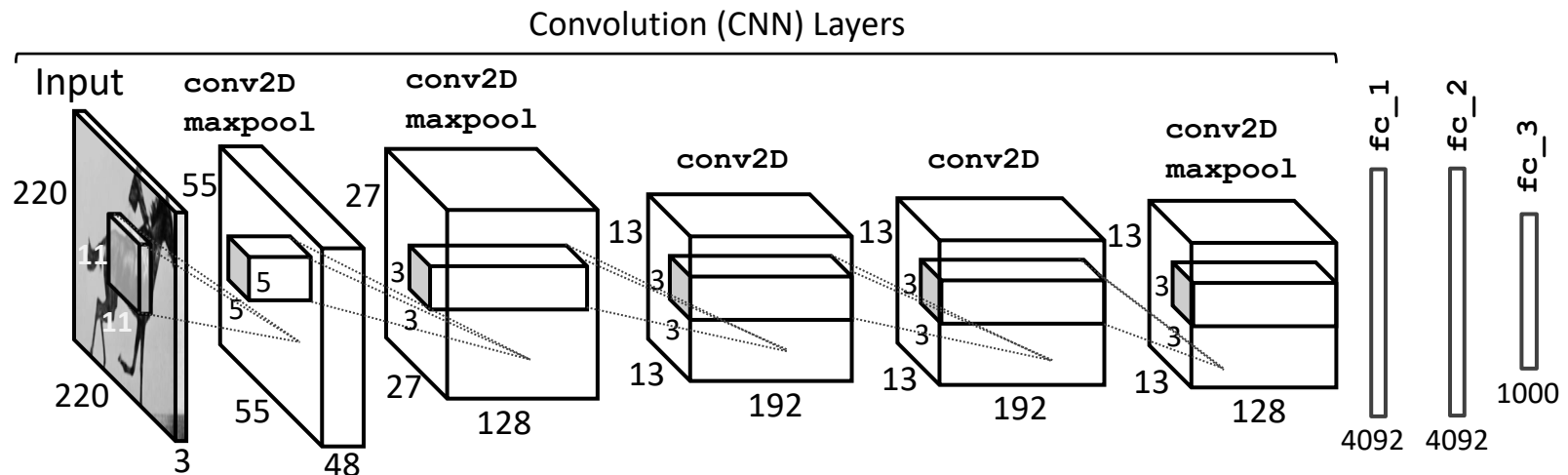
# AlexNet

15

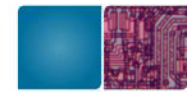
Input Size: 220x220x3

Five convolution layers

Three fully connected layers

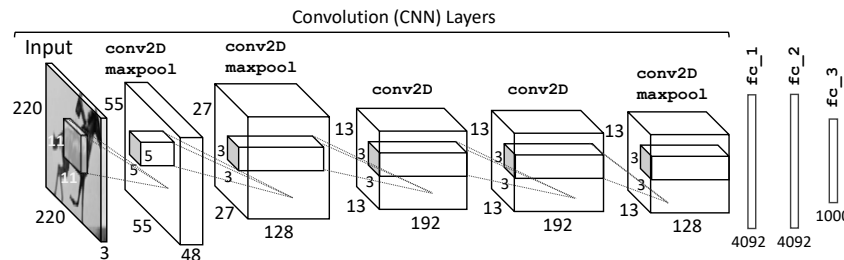
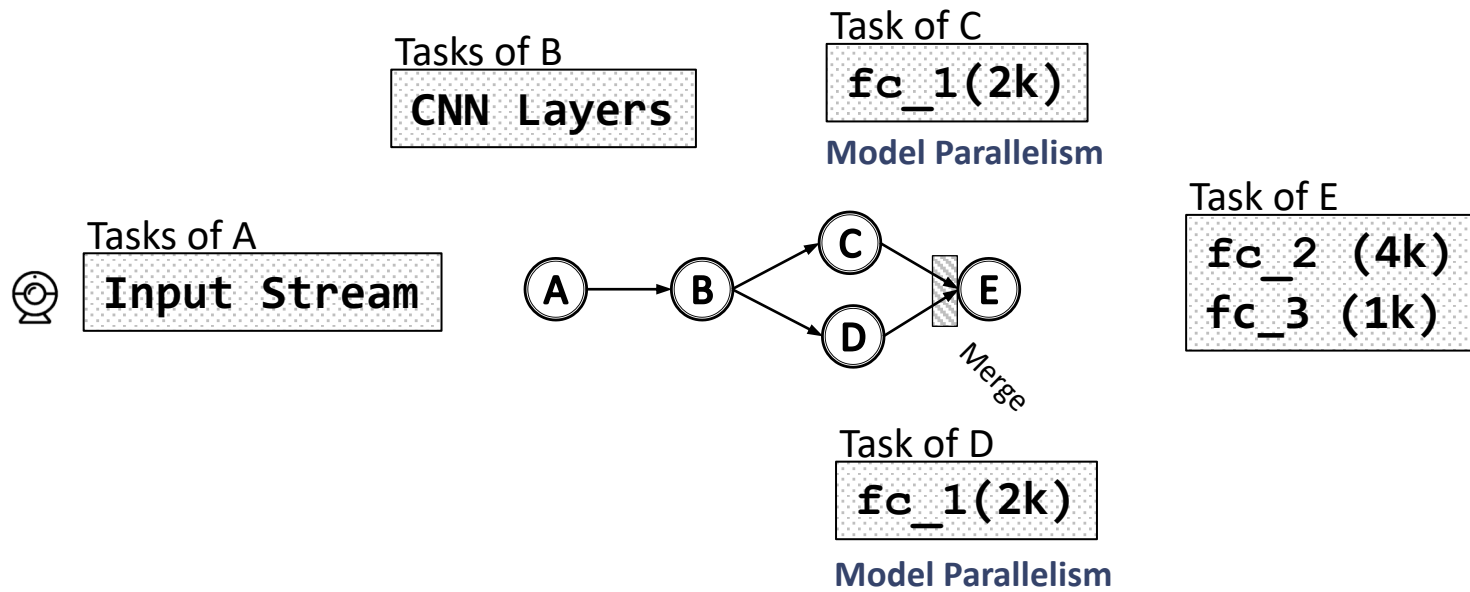


A. Krizhevsky et al., "Imagenet Classification With Deep Convolutional Neural Networks," in NIPS 2012



# AlexNet Distribution I

## Five-device system:

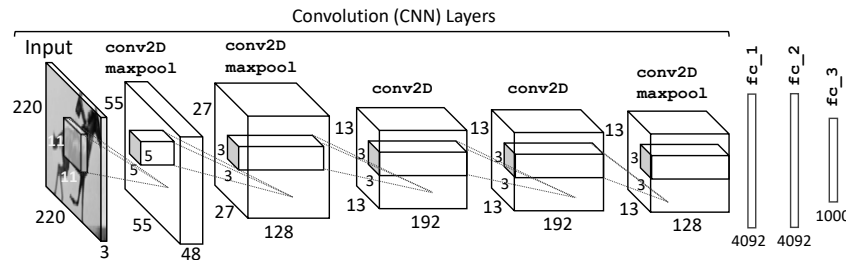
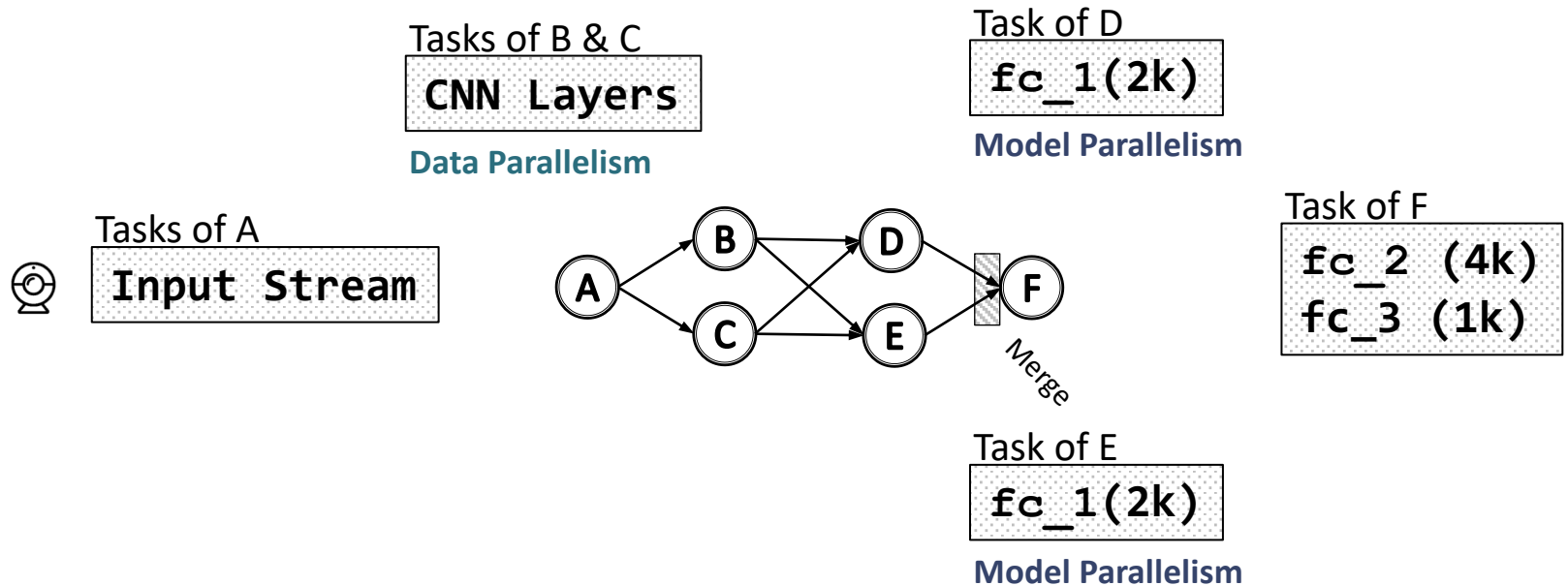


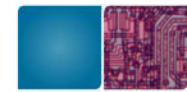




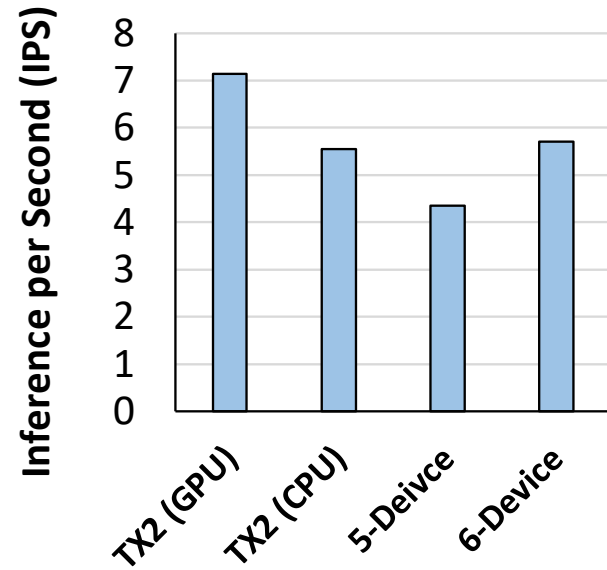
# AlexNet Distribution II

## Six-device system:

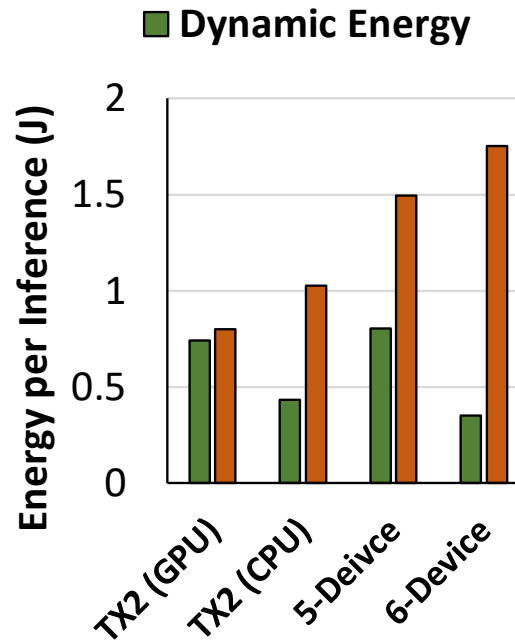




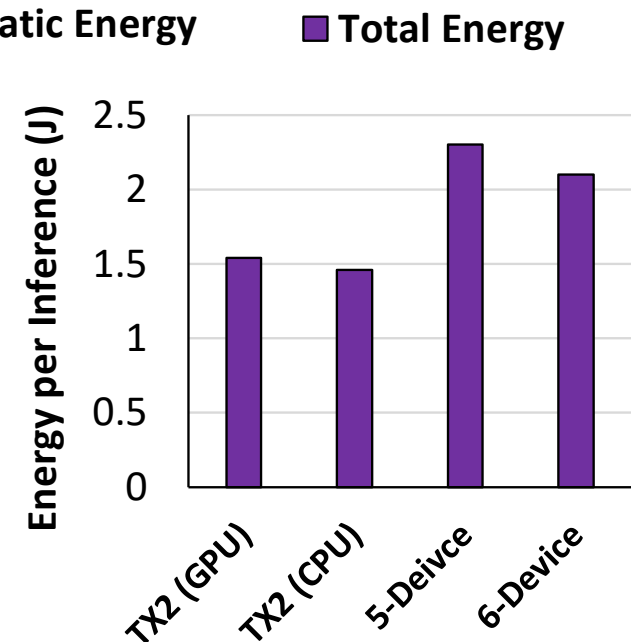
# AlexNet Results



IPS



Dynamic and Static Energy



Total Energy

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



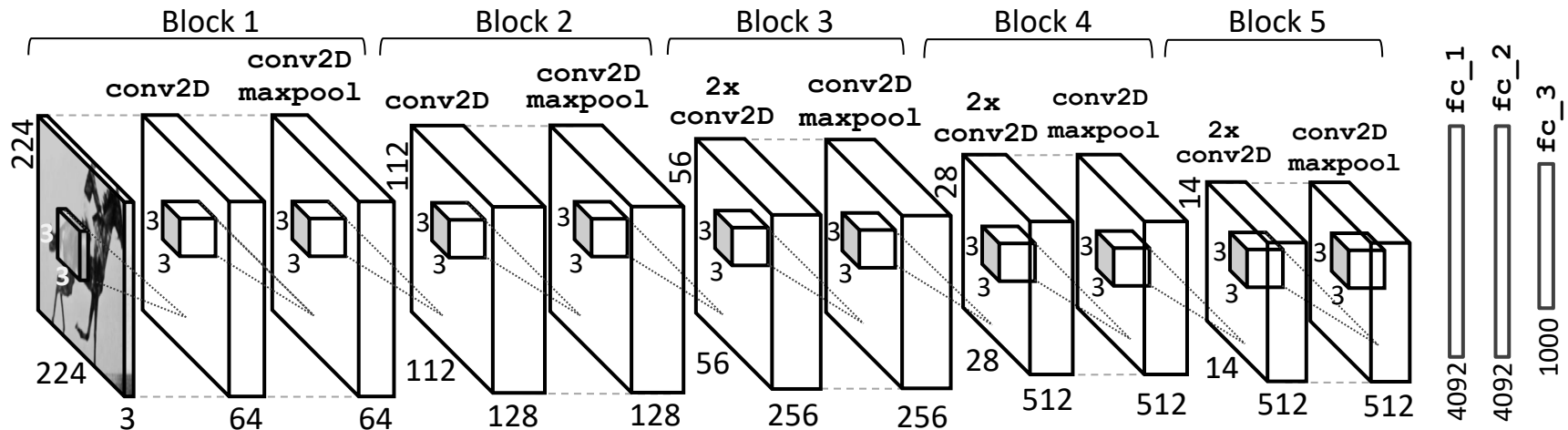
# VGG16

19

Input Size: 224x224x3

13 convolution layers

Three fully connected layers

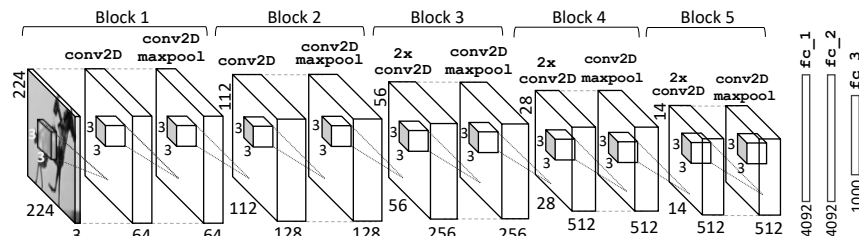
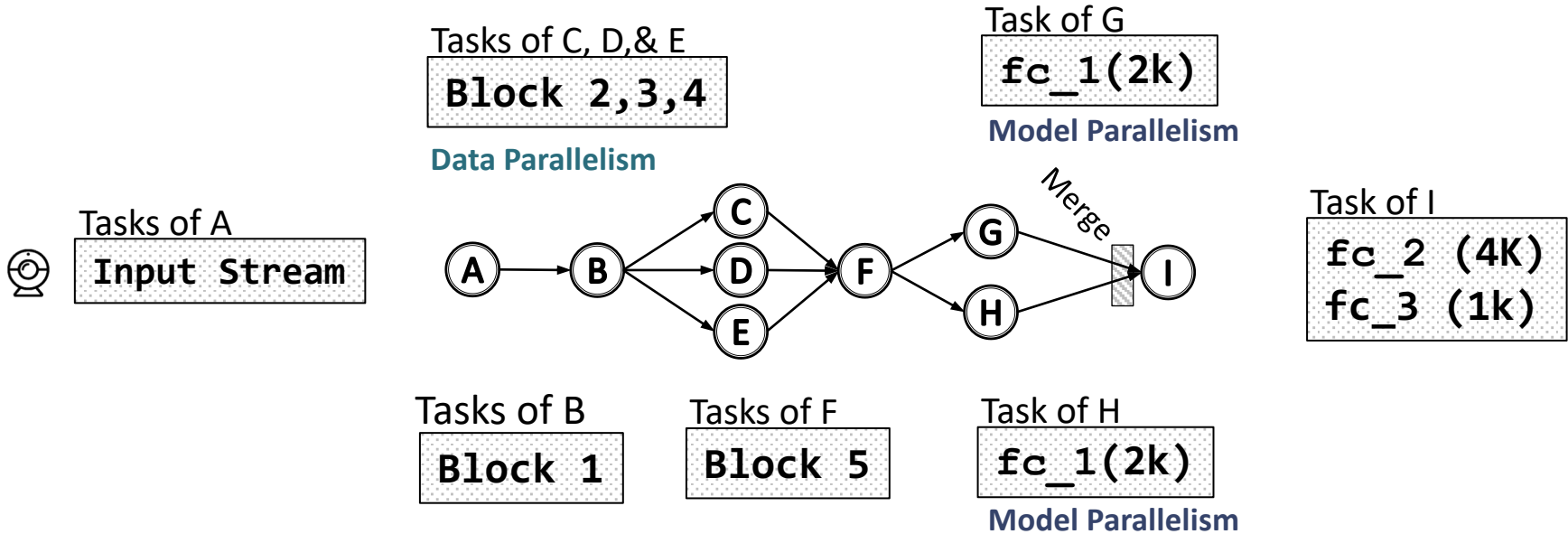


K. Simonyan et al., "Very Deep Convolutional Networks for Large-Scale Image Recognition," in ICLR, 2015.



# VGG16 Distribution I

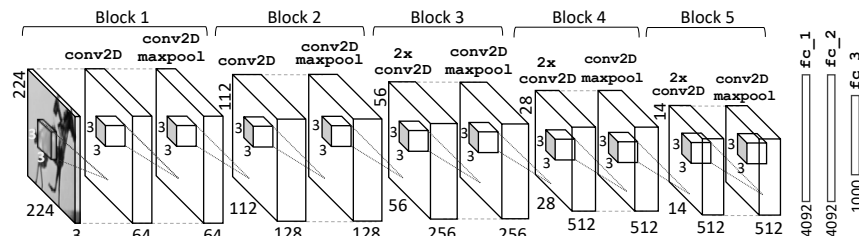
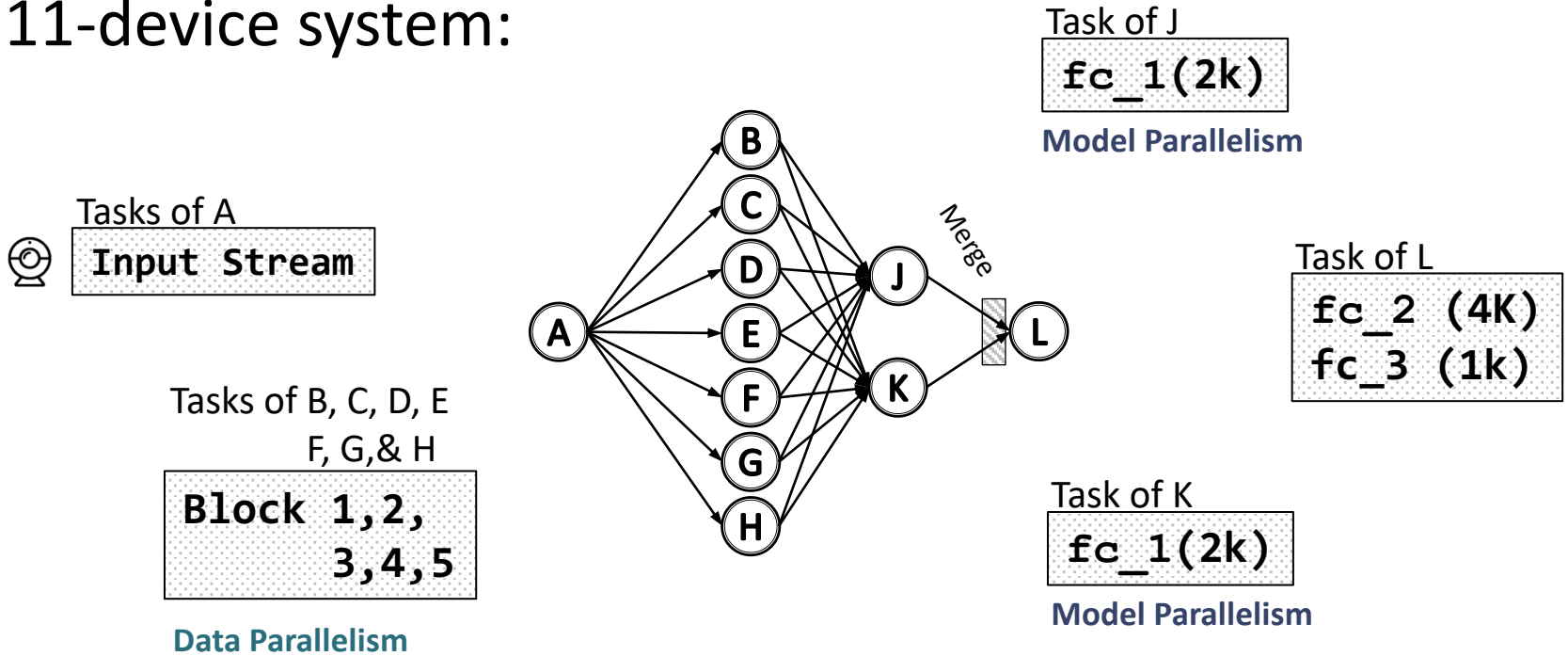
## Nine-device system:

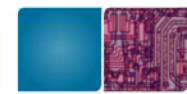




# VGG16 Distribution II

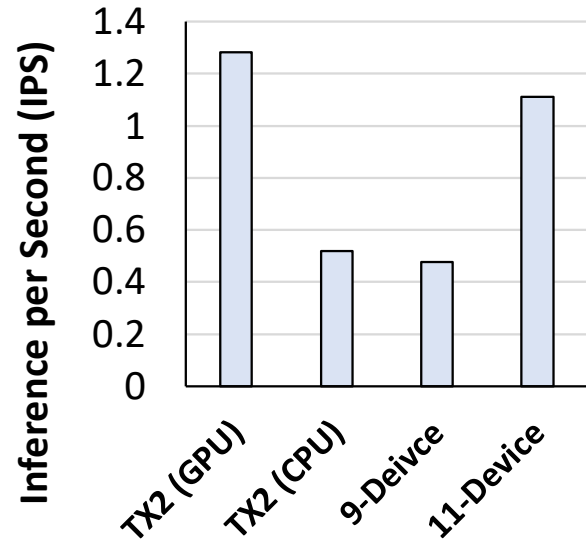
11-device system:



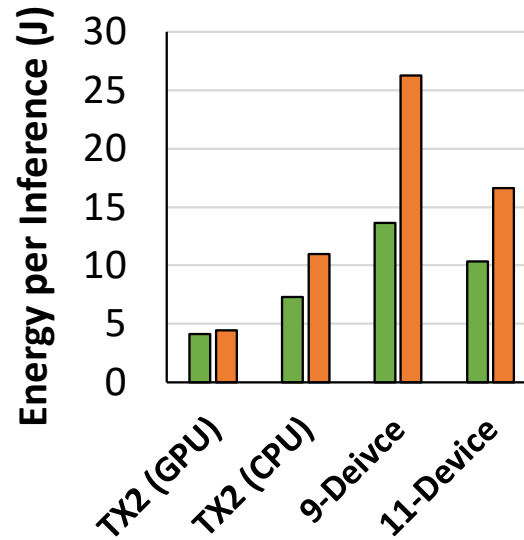


# VGG16 Results

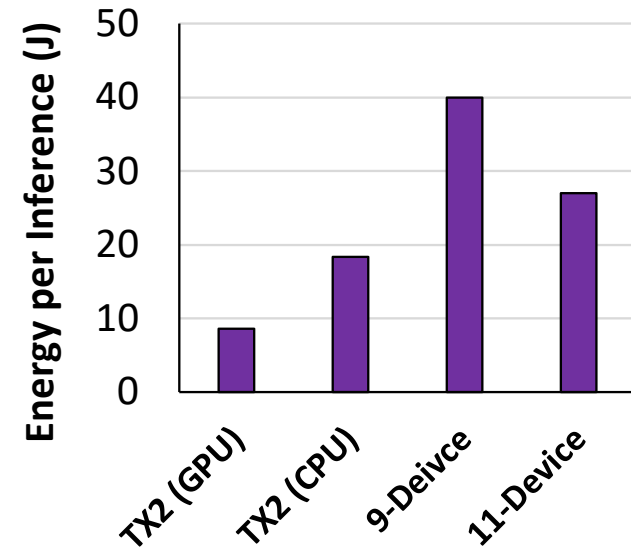
■ Dynamic Energy   ■ Static Energy   ■ Total Energy



IPS



Dynamic and Static Energy



Total Energy

Comparable IPS with TX2 (-15%)

We achieve 2.3x speedup, by reassigning CNN blocks



# Conclusions

---

23

- ▶ We used a farm of Raspberry Pis for DNN processing
- ▶ We are able to process IoT data locally by distribution
- ▶ Our technique achieves acceptable real-time performance

## Future Work:

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