



# Internet of Things Devices

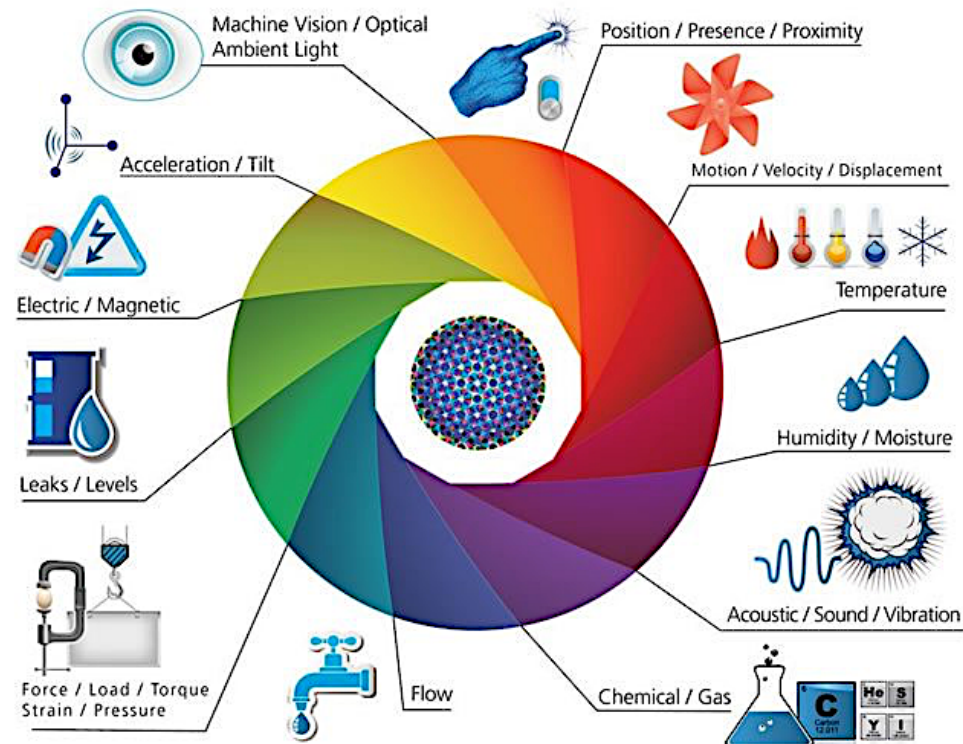
- ▶ Internet of Things (**IoT**) devices
  - ▶ Have access to an abundance of raw data
  - ▶ In home, work, or vehicle

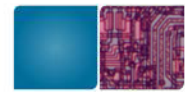




# IoT: Raw Data & Processing

- ▶ IoT is gaining ground with the widespread of
  - ▶ Embedded processors
  - ▶ Ubiquitous wireless networks
- ▶ Access to raw data
  - ▶ Understand it!
  - ▶ Real-time constraints
  - ▶ Limited resources
    - ▶ Power
    - ▶ Compute





# IoT: DNN-based Processing

---

- ▶ With deep neural networks (**DNNs**):
  - ▶ With DNNs IoTs can
    - ▶ Process several new data types and
    - ▶ Understand behaviors
  - ▶ Speech, vision, video, and text
- ▶ But, DNNs are resource hungry
  - ▶ Cannot meet real-time constraints on IoT devices
  - ▶ Several DNNs cannot be executed on IoTs

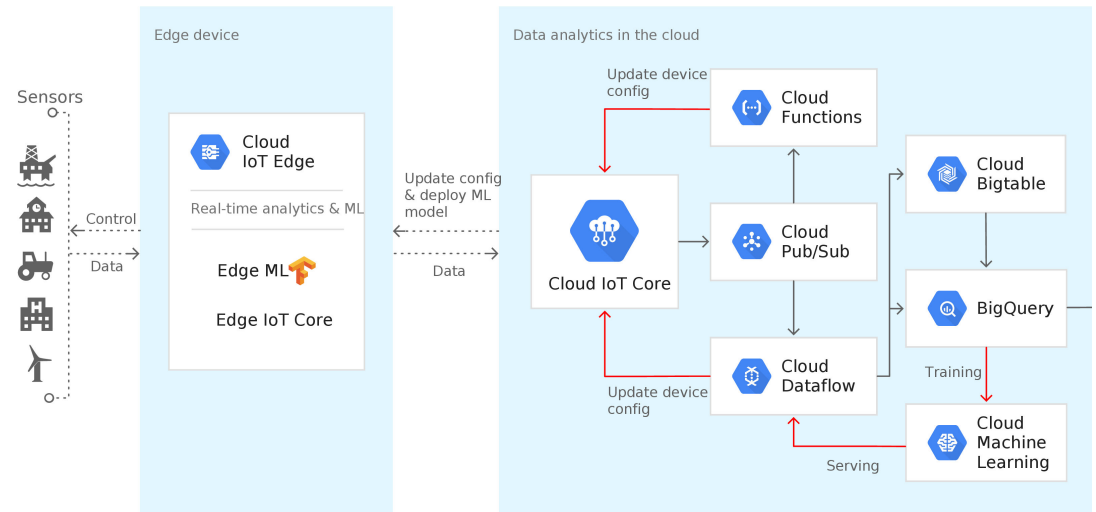
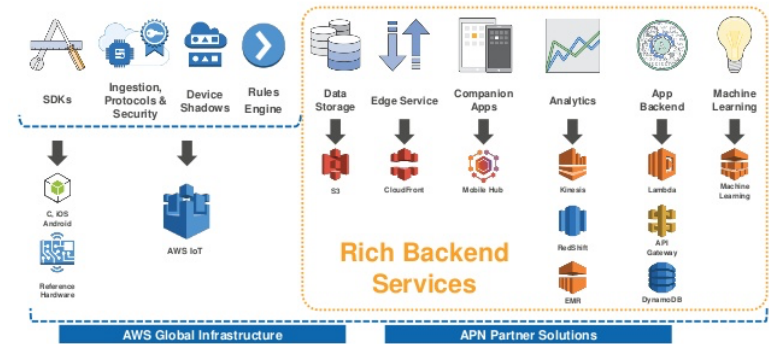


# Approach 1: Offload to Cloud

► Send the request to cloud services

- AWS
- Google Cloud
- Microsoft

## AWS Platform For IoT Solutions





# Why Cloud is not Always a Solution

---

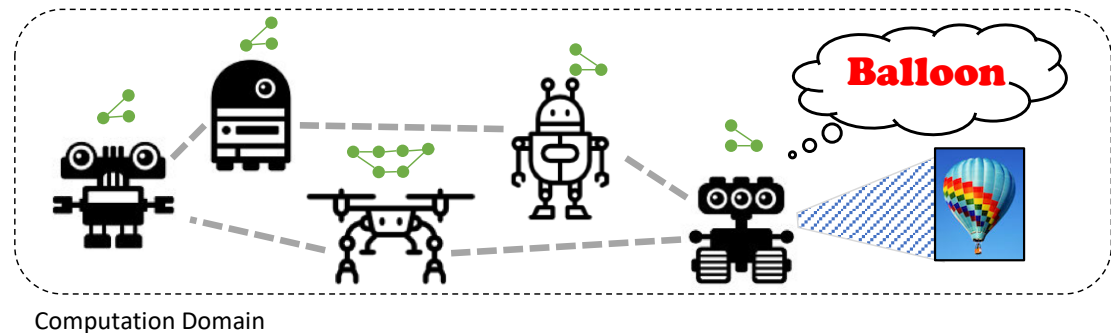
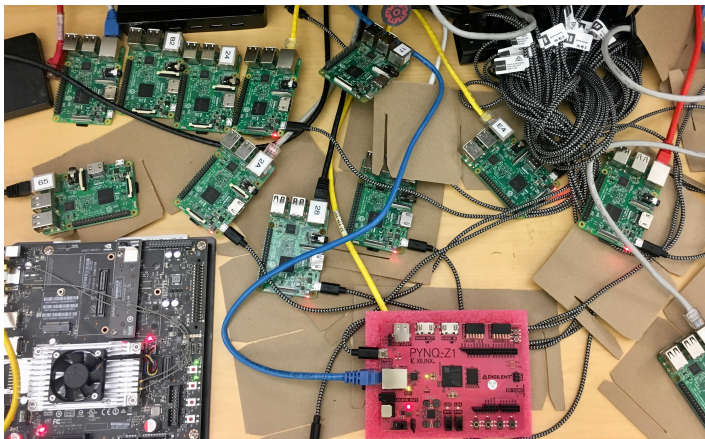
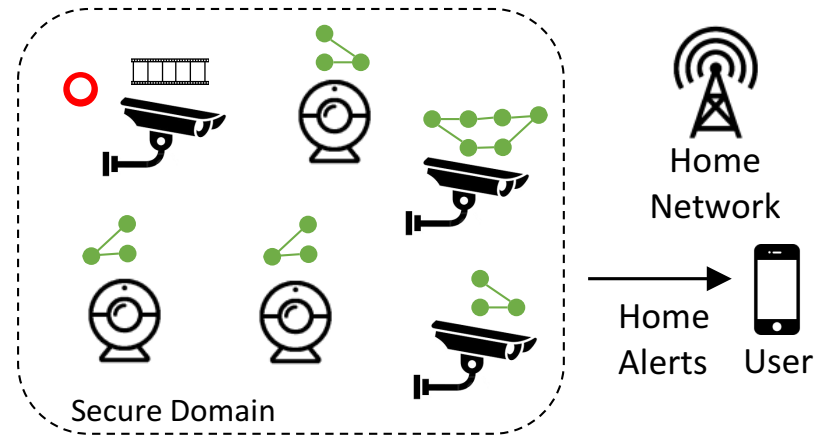
- ▶ Unreliable connections to the cloud
  - ▶ Plus low bandwidth and high latency
- ▶ Disconnected Devices
- ▶ **Privacy**
  - ▶ Privacy preserving learning (e.g., differential privacy)
  - ▶ Privacy preserving inference (e.g. homomorphic encryption)
- ▶ Personalization
- ▶ Federated learning



# Approach 2: IoT Collaboration

## ▶ Distribute computations with collaboration

- ▶ To meet demands of DNNs
- ▶ On top of common DNN techniques for constrained devices (e.g., pruning)





# IoT Collaboration Pros & Cons

---

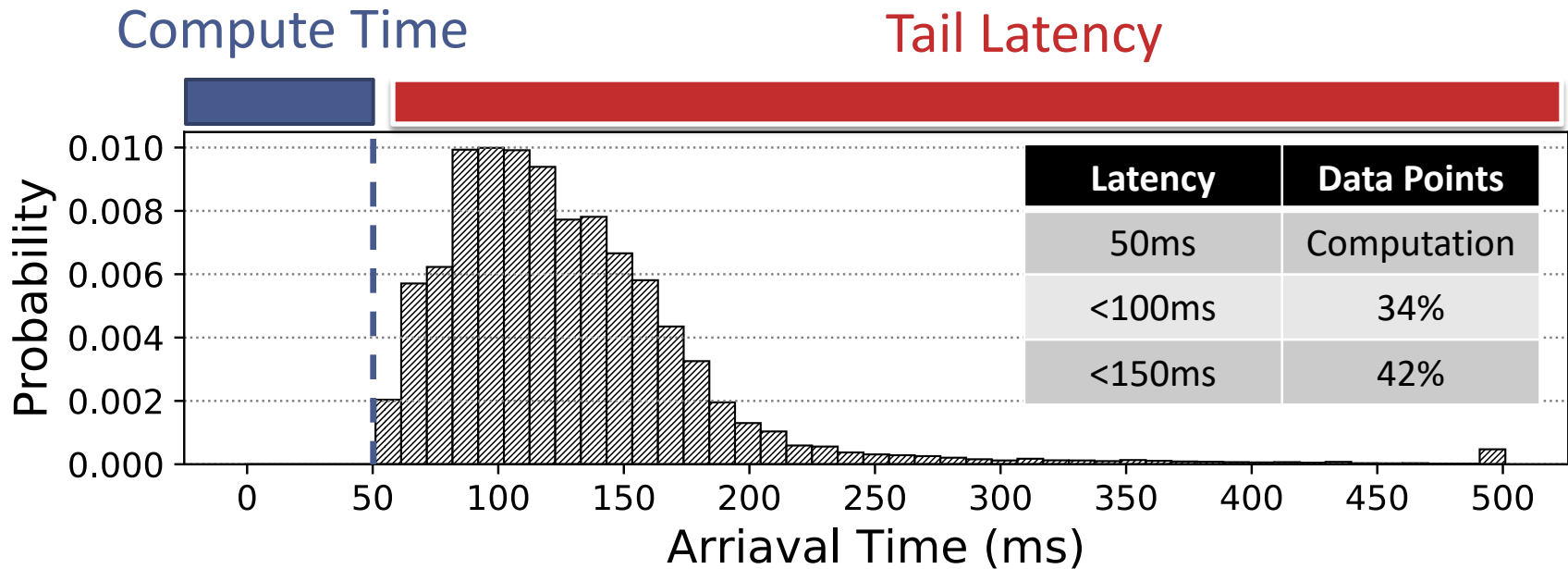
- ▶ Assuming DNN performance barrier is solved with collaboration among IoT devices

Pros	Cons
Not Dependent on Cloud	<b>Unreliable Latencies</b>
Privacy Preserving	<b>Accuracy Drop due to Data Loss &amp; Device Failure</b>
Enables Personalized Insight	



# Challenges Impact: Unreliable Latencies

- ▶ Histogram of arrival times in 4-node system performing AlexNet (model parallelism).



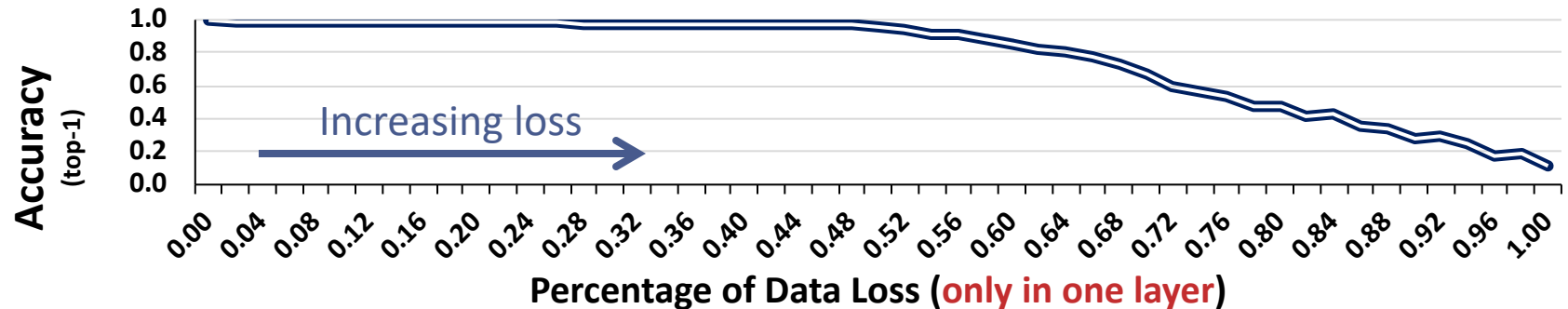
▶ Long Tail and Max Latency -> Straggler Problem



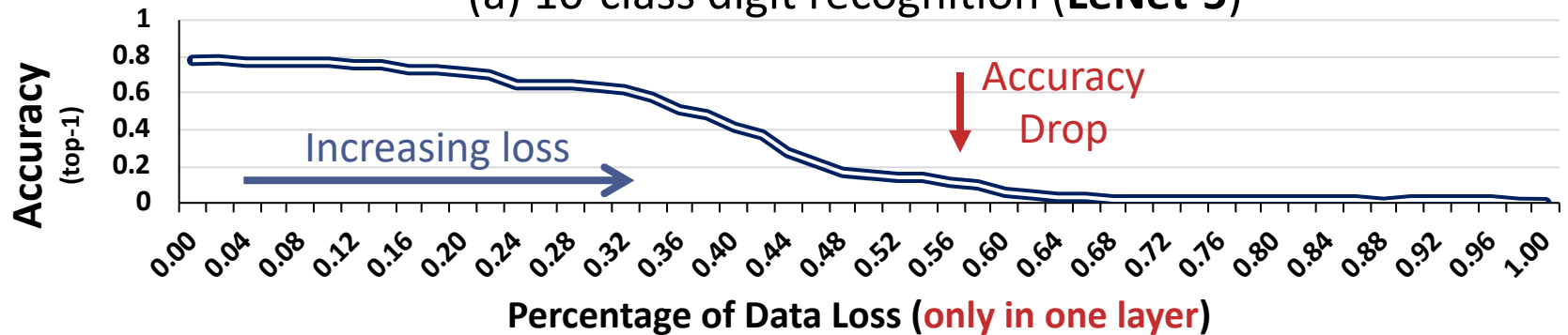


# Challenges Impact: Accuracy Drop

► Common to loose data parts due to



(a) 10-class digit recognition (LeNet-5)



(b) 1000-class image recognition (Inception v3)

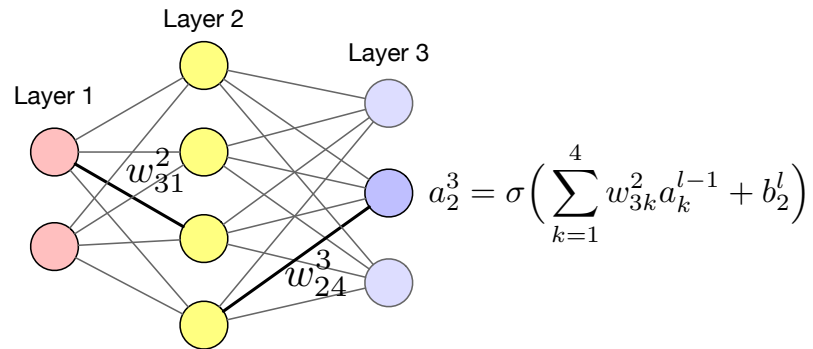
► **High Accuracy Drop**



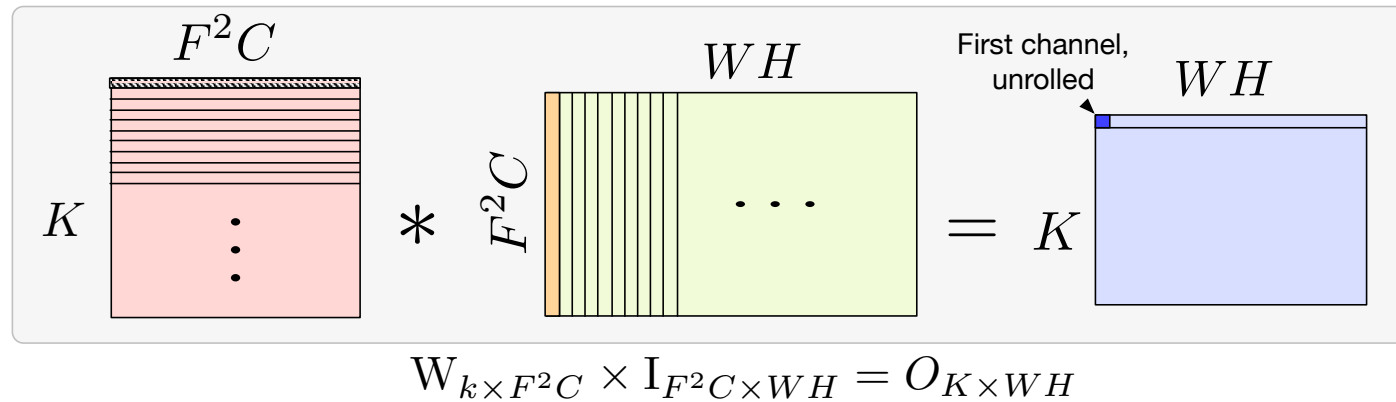
# Computation of DNNs

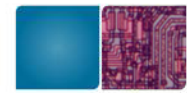
- ▶ Each layer's computations can be represented as matrix-matrix multiplication (GEMM kernels).

Fully-connected layer:



Conv. layer:





# Computation Distribution of DNNs

## ► Methods distributing computation of a model\*

Fully-connected Layers

Output splitting:

$$\begin{bmatrix} w_{11} & w_{12} & \dots & w_{1k} \\ w_{21} & w_{22} & \dots & w_{2k} \\ w_{31} & w_{32} & \dots & w_{3k} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m1} & w_{m2} & \dots & w_{mk} \end{bmatrix}_{m \times k} \times \begin{bmatrix} a'_1 \\ a'_2 \\ a'_3 \\ \vdots \\ a'_k \end{bmatrix}_{k \times 1} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_m \end{bmatrix}_{m \times 1}$$

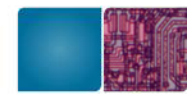
Weights (divided among nodes)      Inputs (every node needs a copy)      Outputs (each node independently)

Input splitting:

$$\begin{bmatrix} w_{11} & w_{12} & \dots & w_{1k} \\ w_{21} & w_{22} & \dots & w_{2k} \\ w_{31} & w_{32} & \dots & w_{3k} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m1} & w_{m2} & \dots & w_{mk} \end{bmatrix}_{m \times k} \times \begin{bmatrix} a'_1 \\ a'_2 \\ a'_3 \\ \vdots \\ a'_k \end{bmatrix}_{k \times 1} = \begin{bmatrix} \delta a_1 \\ \vdots \\ \delta a_m \end{bmatrix}_{m \times 1} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_m \end{bmatrix}_{m \times 1}$$

Weights (divided among nodes)      Inputs (divided among nodes)      Outputs (each node calculates partial sums)

- Same can be applied on conv. layers\*
  - Channel , spatial , and filter splitting

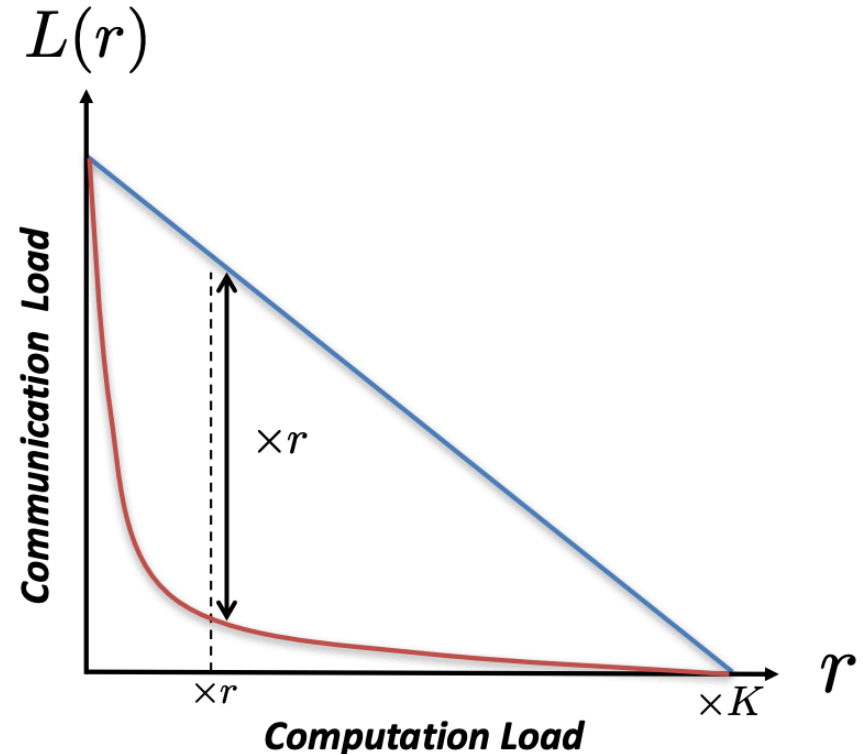


# Coded Distributed Computing (CDC)

- ▶ Designed for MapReduce workloads (2018)\*
- ▶ Performing redundant or coded computer per node to reduce communication.

This work: **DNNs on IoT**

**More Compute / Node  
=  
More Reliability**



\* Li, Songze, et al. "A fundamental tradeoff between computation and communication in distributed computing." *IEEE Transactions on Information Theory* 64.1 (2018): 109-128.



# Using CDC for Robustness

---

- ▶ Add column-wise summation of the weights:

$$\begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{11} + w_{21} & w_{12} + w_{22} \end{bmatrix} \times \begin{bmatrix} a'_1 \\ a'_2 \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ a_1 + a_2 \end{bmatrix}$$

- ▶ The new weights are constant, so done in offline

$$\begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{:1}^{cdc} & w_{:2}^{cdc} \end{bmatrix} \times \begin{bmatrix} a'_1 \\ a'_2 \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ a^{cdc} \end{bmatrix}$$

- ▶ Distribute outputs among nodes
    - ▶ Thus, applicable only to output-splitting methods
-



# How to Distribute CDC and Recover?

## ▶ Add column-wise summation of the weights:

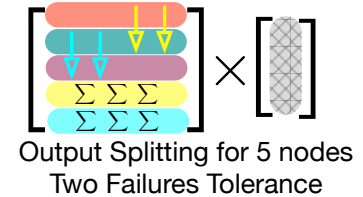
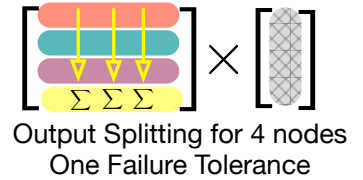
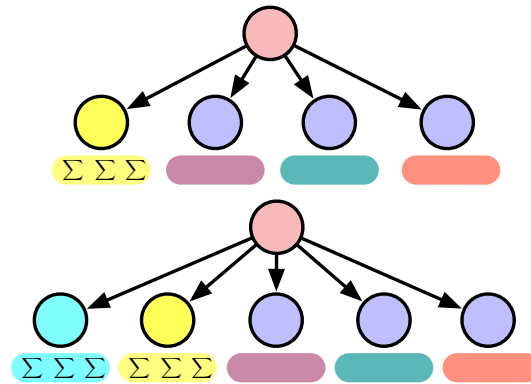
▶ Simple example  
(one output/device)

▶ Recovery

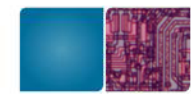
▶ **Subtraction vs. Multiplication**

▶ You also needs the weights, that you would not have in the final node

▶ Multiple out/device: Just create a new weight matrix



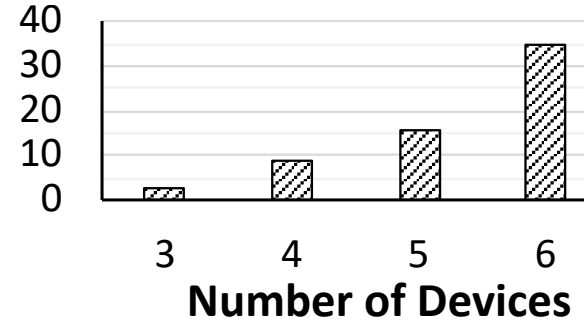
$$\begin{bmatrix} w_{11} + w_{(\frac{m}{2}+1)1} & w_{12} + w_{(\frac{m}{2}+1)2} & \dots & w_{1k} + w_{(\frac{m}{2}+1)k} \\ w_{21} + w_{(\frac{m}{2}+2)1} & w_{22} + w_{(\frac{m}{2}+2)2} & \dots & w_{2k} + w_{(\frac{m}{2}+2)k} \\ \vdots & \vdots & \ddots & \vdots \\ w_{\frac{m}{2}1} + w_{m1} & w_{\frac{m}{2}2} + w_{m2} & \dots & w_{\frac{m}{2}k} + w_{mk} \end{bmatrix}_{\frac{m}{2} \times k}$$



# Straggler Mitigation & Failure Coverage

Do not need to wait for all devices to send data:  
(AlexNet)

Performance Improvement (%)



Better Coverage versus with 2-modular redundancy (2MR):

