





# Distributed inference in IoT clusters

Dr. Victor Cionca Munster Technological University



#### Bio

Lecturer in Computer Science @ MTU Cork

PhD in Wireless Sensor Networks from University of Limerick, 2012

Lead PI in MoNet project (MTU-TUTF funding). FI in SFI CONFIRM Centre. FI in SFI CONNECT Centre. NI in SFI ADVANCE CRT.

High speed wireless networks.

Network control, optimization, management.

Programmable networks.





#### Today's outline

The challenges of ML-based inference in IoT devices

Distributing inference between IoT devices

Communication challenges.



#### Machine Learning in IoT

Federated Learning

- Train on local data (very slow)
- Aggregate model centrally.

#### Inference

- Far more common
- Still quite slow
- 2017 flagship smartphone ~100s for image classification. (Huynh et al, "Deepmon...", ACM MobiSys 2017)











#### Delay constrained applications

- Ex: 10FPS processing
- <100ms for data processing on the IoT device</li>

Current practice: offload processing to the cloud.

MTU Victor.Cionca@mtu.ie















Examples:

- Waste classification for sorting plastics
- Urban monitoring

Runs if needed

Common feature: slow, computational classifiers are only executed *sometimes*. Per-device workload is *not homogeneous*.

MTU Victor.Cionca@mtu.ie



MTU Victor.Cionca@mtu.ie





#### **Parallelising DNNs**

DNNs are predominantly sequential tasks.

Typical offloading uses *model* partitioning.

P2 must wait until P1 completes.

No gains if P2 has same computational power.



AlexNet

Output: 1 of 1000 classes Image copyright Cmglee CC (src: wikimedia)



## **Parallelising DNNs**

DNNs are predominantly sequential tasks.

Typical offloading uses *model* partitioning.

P2 must wait until P1 completes.

No gains if P2 has same computational power.

#### Input partitioning

- Partitions can run in parallel
- Performance gains regardless of computational power



#### Distributed inference in IoT clusters



Work with Mr. Jamie Cotter (PhD student), Dr. Ignacio Castineiras, Dr. Donna O'Shea @ MTU Cork, CS Dept. Funded by SFI ADVANCE CRT.



## Computation vs **Communication**

## With computational offloading, **communication is an overhead**.

DNN inference time: ~100's ms

Communication delays (RTT):

- WiFi <10ms
- Mobile: LTE ~50ms, 5G promised <1ms
- To nearest datacentre +10-20ms.

This is the ideal case, without communication errors.

We must improve the communication.





MTU Victor.Cionca@mtu.ie

#### **Faster communication**

More bandwidth - example mmwave (2GHz in 60GHz band  $\rightarrow$  ~35Gbps).

More transmission streams (MIMO).

Higher modulation orders (more bits per radio symbol).

However...

Faster transmission also has higher probability of errors.





#### Reliability - avoiding transmission errors

Instantaneous signal power



Errors

Adding channel estimation data

Adding error checking bits

Adding guard times

Adding error correction bits

Adding complex error correction algorithms etc.

Some applications have very stringent reliability requirements.

Ex industrial control loops require 10<sup>-9</sup> bit error probability.





Results obtained by Mr. Yasantha Samarawickrama (PhD student) SFI CONFIRM Centre.



Confirm

**Smart Manufacturing** 

Reliability and latency are most times inversely proportional.

Some alternatives: spatial diversity, cooperation, NOMA.

• Early work with PhD student Ms. Tabinda Ashraf

#### URLLC: Ultra Reliable, Low Latency Communication

Defined by the application.

Related to Age of Information.

Several definitions

V1: Maximise the reliability, subject to delay constraintsV2: Minimise the latency, subject to maximum transmission error.V3: Achieve transmission rate within given time.

We need joint optimization of reliability and latency based on channel characteristics.



## Achieving URLLC



Work conducted with Ms. Tabinda Ashraf, Mr. Yasantha Samarawickrama, Dr. Álvaro de Medeiros. Funded by SFI CONFIRM Centre.



#### Challenge: when and how to configure the comms?

The meaning of optimal configuration

- System is configured to perform at max parameters for given channel state
- Any unexpected change in channel state leads to failures.

Basic: expected channel state.

SotA: expected distribution of channel states.

We observed that channel parameters can change drastically  $\rightarrow$  different channel.

Currently investigating optimization in such scenarios.



Results obtained by Dr. Alvaro de Medeiros based on channel measurements published by NIST TN 1951.



#### Conclusions

- ML-based inference is still an issue in IoT
- Offloading between IoT devices possible due to non-homogeneous workloads
- Leads to more efficient resource utilization

Still challenges

- Dealing with delay constraints
- Dealing with wireless communication errors.

Thank you for your time!

