# FPGAs – EPIC Benefits

Philip Leong
Director, Computer Engineering Laboratory
http://phwl.org/talks

THE UNIVERSITY OF
SYDNEY

› Focuses on how to use parallelism to solve demanding problems

- Novel architectures, applications and design techniques using VLSI, FPGA and parallel computing technology

› Research

- Nanoscale interfaces

- Machine learning

- Reconfigurable computing

› Collaborations

- Consunet, DST Group

- Intel, Xilinx

› Ex-students

- Xilinx, Intel, Waymo

FPGAs

Applications

Our work
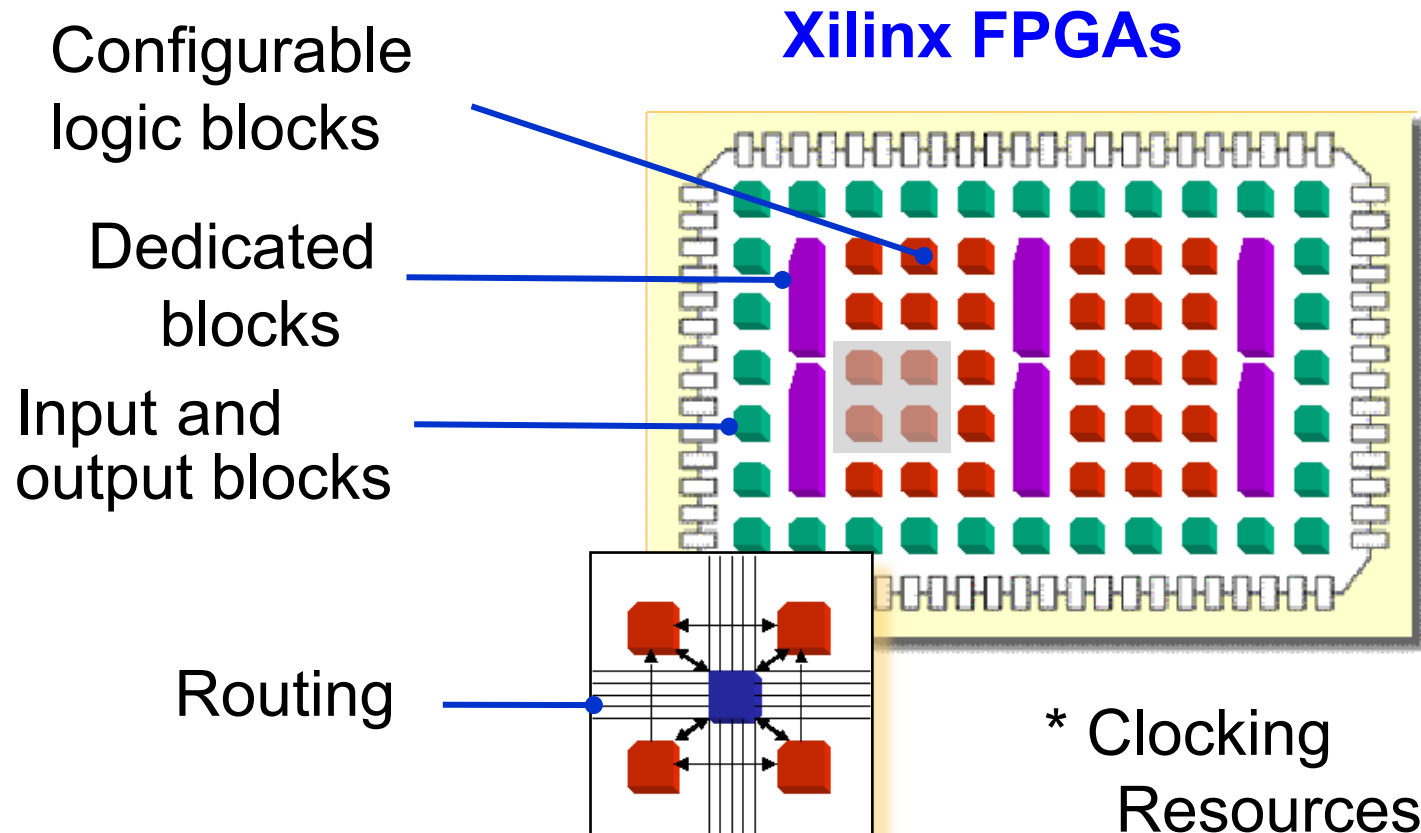
THE UNIVERSITY OF
SYDNEY

FPGAs

Applications

Our work

THE UNIVERSITY OF
SYDNEY
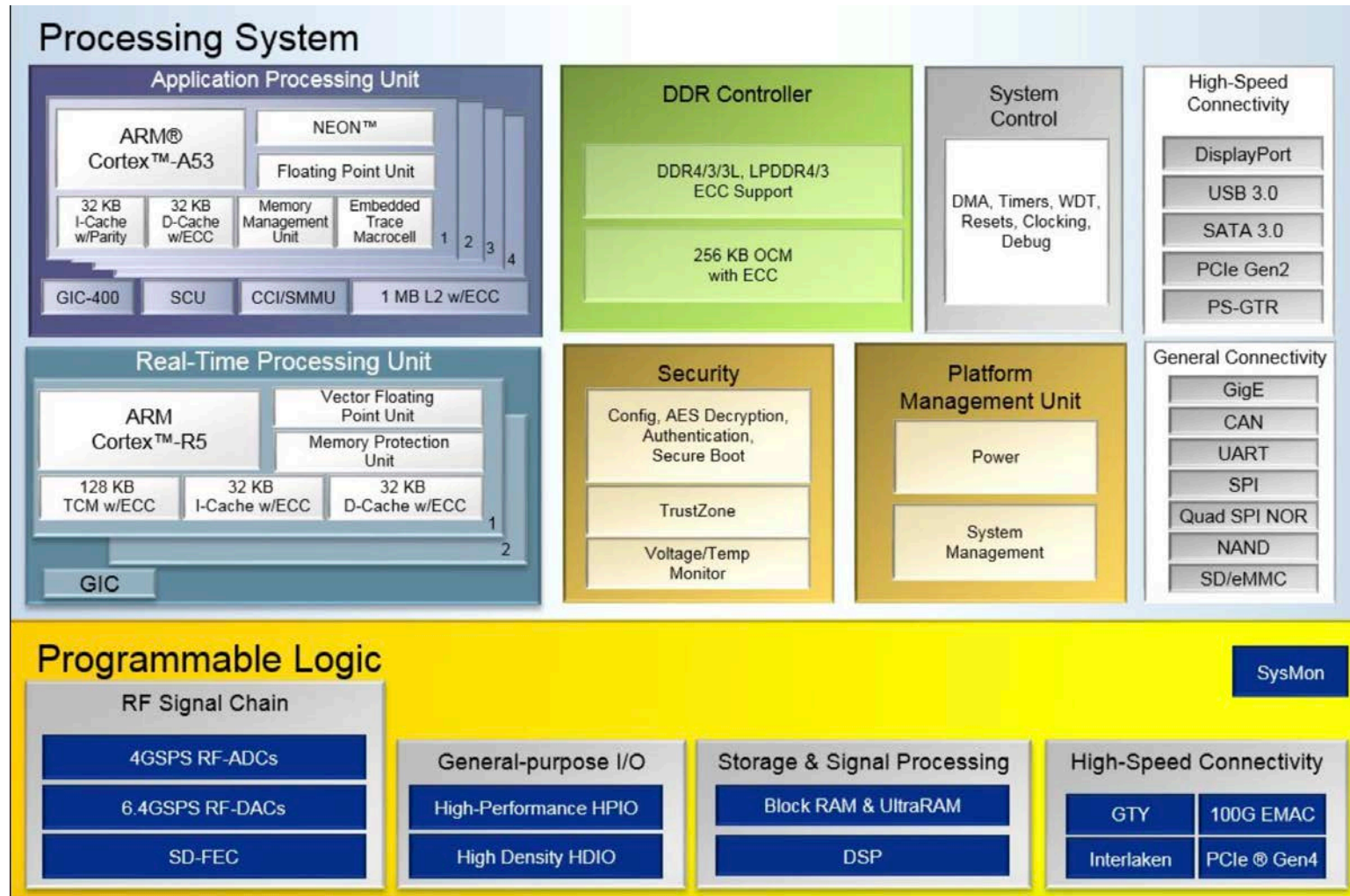
User-customisable integrated circuit

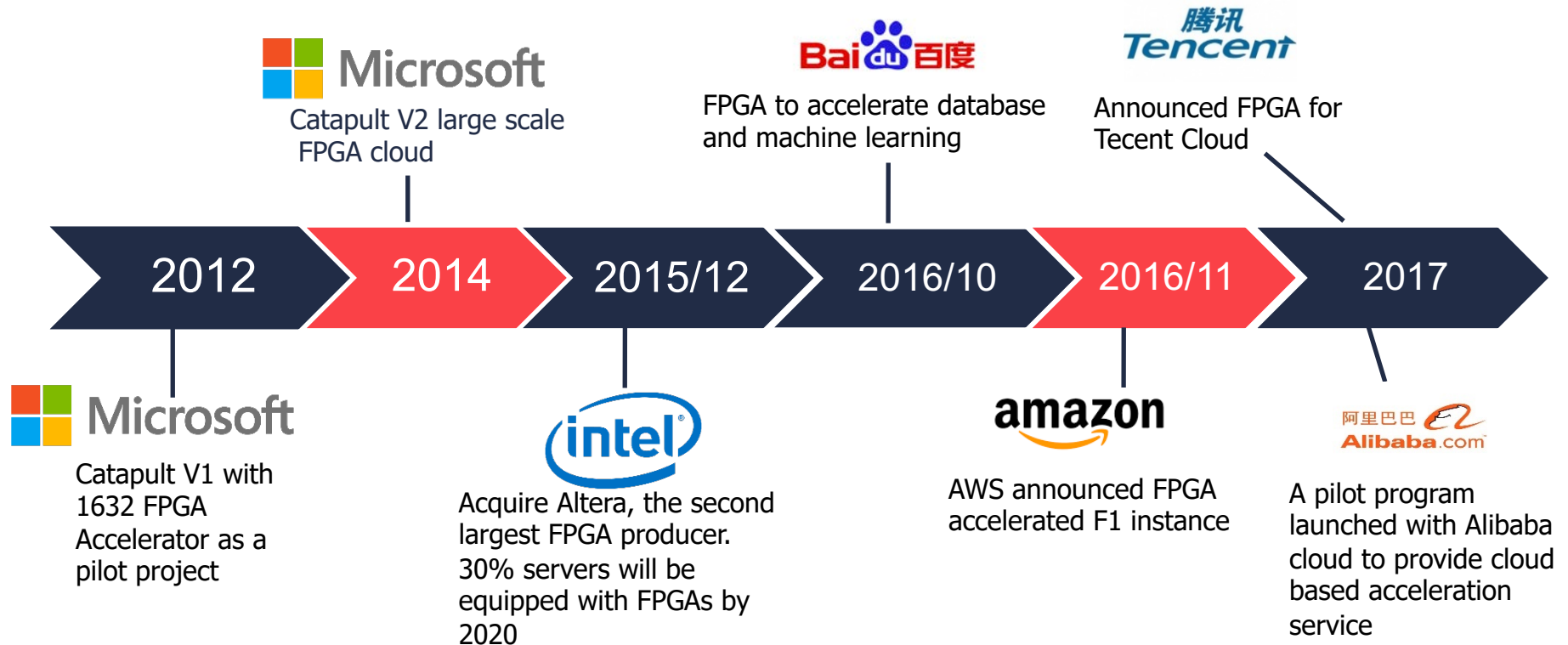› Dedicated blocks: memory, transceivers and MAC, PLLs, DSPs, ARM cores

**Xilinx FPGAs**

Configurable
logic blocks

Dedicated
blocks

Input and
output blocks

Routing

* Clocking
Resources

Source: Xilinx

# Recent Uptake in Reconfigurable Computing

**Microsoft**
Catapult V2 large scale FPGA cloud

**Bai du 百度**
FPGA to accelerate database and machine learning

**腾讯 Tencent**
Announced FPGA for Tecent Cloud



| 2012 | 2014 | 2015/12 | 2016/10 | 2016/11 | 2017 |

**Microsoft**
Catapult V1 with 1632 FPGA Accelerator as a pilot project

**intel**
Acquire Altera, the second largest FPGA producer. 30% servers will be equipped with FPGAs by 2020

**amazon**
AWS announced FPGA accelerated F1 instance

**阿里巴巴 Alibaba.com**
A pilot program launched with Alibaba cloud to provide cloud based acceleration service

Speed of Innovation Outpaces Silicon Cycles

Source: Xilinx

› FPGAs commercial off-the-shelf

› They offer an opportunity to implement complex algorithms with higher throughput, lower latency and lower power through

- **E**xploration– easily try different ideas to arrive at a good solution

- **P**arallelism – so we can arrive at an answer faster

- **I**ntegration – so interfaces are not a bottleneck

- **C**ustomisation – problem-specific designs to improve efficiency (power, speed, density)

Vitis: Unified Software Platform

Source: Xilinx

FPGAs

Applications

Our work

THE UNIVERSITY OF
SYDNEY

› Compact Muon Solenoid

- Few interesting events ~100 Higgs events/year

- 1.5Tb/s real-time DSP problem

- (2014) More than 500 Virtex and Spartan FPGAs used in real-time trigger

- (2019 doing FPGA-based DNN inference using Vivado HLS)





40 million collisions per second → PB/s → 100,000 selections per second → TB/s → 1,000 selections per second → GB/s →

This can generate up to a petabyte of data per second.
Filtering the data in real time, selecting potentially interesting events (trigger).

Source: Intel

› Uses FPGAs for DNNs, Bing search, and software defined networking (SDN) acceleration to reduce latency, while freeing CPUs for other tasks

- 2010: MSR study FPGAs to accelerate Web search
- 2012: Project Catapult's scale pilot of 1,632 FPGA servers deployed
- 2013: Bing decision-tree algorithms 40x faster than CPUs
- 2015: FPGAs deployed at scale in Bing and Azure datacenters (> 1M) - enabled 50% ↑ throughput, 25% ↓ latency.
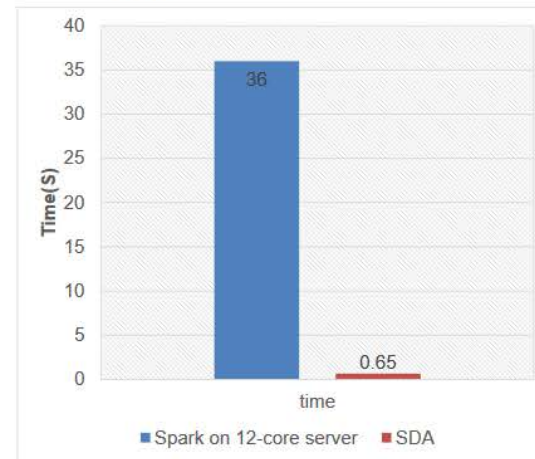
## World's fastest cloud network

Source: Microsoft

# Accelerator for SQL Queries (40% of their data analysis)



| Total data: | ~1EB |
| Processing data : | ~100PB/day |
| Total web pages: | ~1000 Billion |
| Web pages updated: | ~10Billion/day |
| Requests: | ~10Billion/day |
| Total logs : | ~100PB |
| Logs updated: | ~1PB/day |

Evaluation - real case query

- TPC-DS scale = 10 , query3
- Execution time
  - 55x

**Key Zynq UltraScale+ RFSoC Benefits:**

- Integrated Direct RF data converters for 4x4 TX/RX mobile backhaul architectures
- Multi-Level LDPC codec (SD-FEC) to meet 5G standards and support for custom codes
- Turbo Decode (SD-FEC) for 4G LTE-Advanced and 4G LTE Pro
- DSP48-rich fabric (6,620 GMACs) provides high-performance filtering and encoding/decoding
- 33 Gb/s transceivers for 12.2G CPRI and expansion into 16G & 25G CPRI



Source: Xilinx 2019

› Amadeus IT Group S.A adjusted profit €1.27B in 2019

› Accelerated inference of gradient boosted decision trees for search queries and quantified cost

https://github.com/Xilinx/Vitis_Libraries

Source: Xilinx

https://xilinx.github.io/Vitis_Libraries/data_analytics/2020.1/benchmark/result.html

Source: Xilinx

Overview
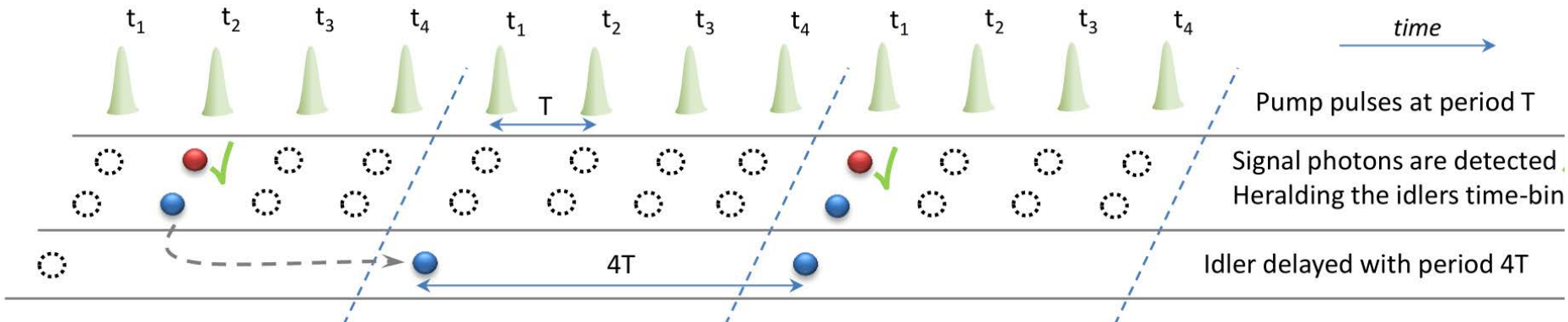
FPGAs

Applications
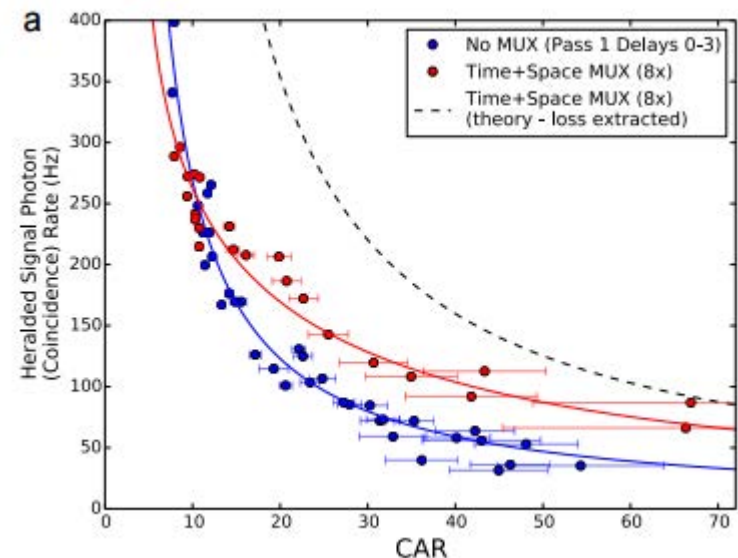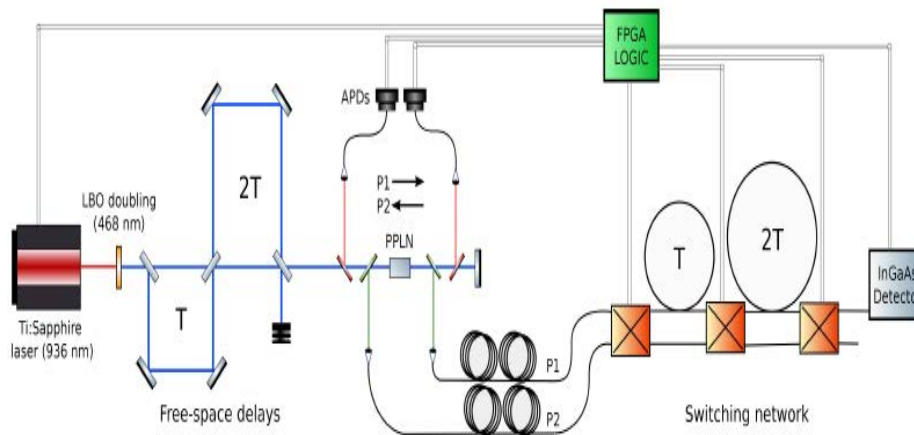
Our work

THE UNIVERSITY OF
SYDNEY

**Initially expectation** : Heralded single photon rate should enhance significantly without degrading coincidence to accidental ratio (CAR)
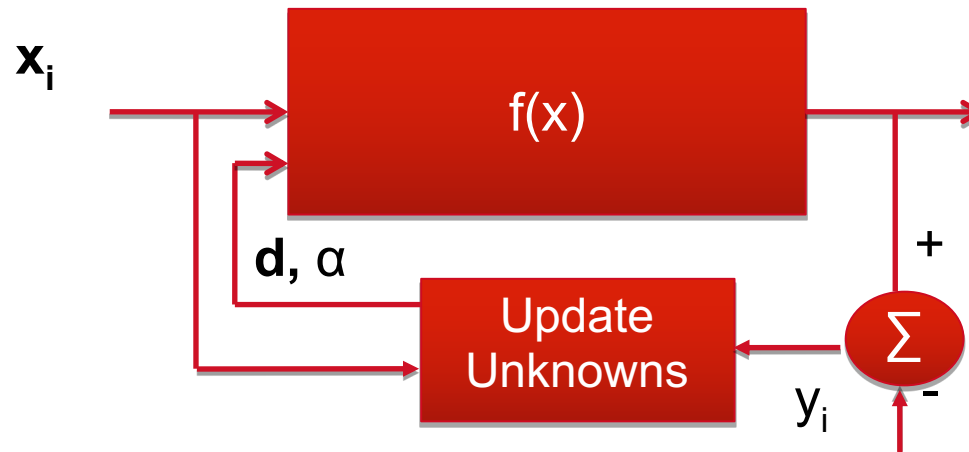

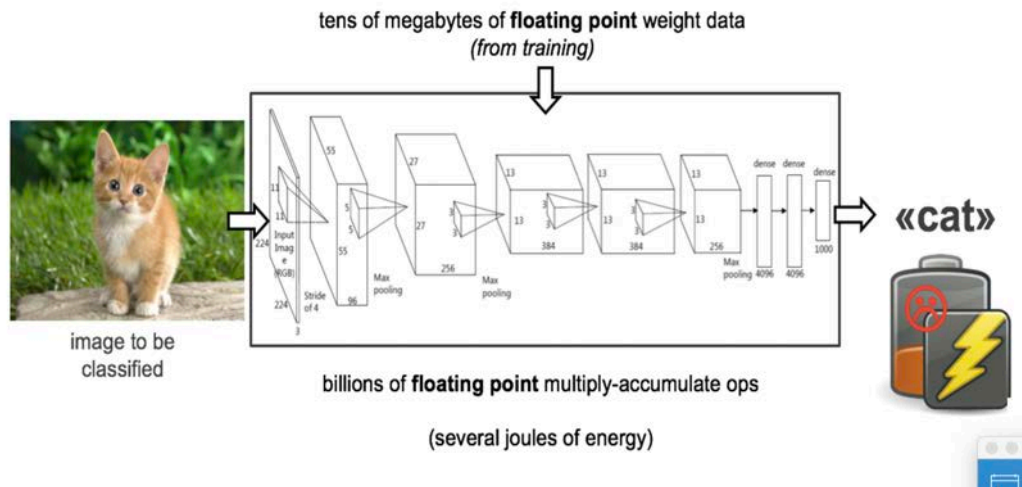
**Enhancement : 33%~59%**

*ARC Linkage with Exablaze*



› A family of kernel methods that can do simultaneous learning and inference

- Highest reported throughput 80 Gbps (TRETS'17)

- Lowest reported latency 80 ns (FPT'15)

- Higest capacity (FPGA'18)

*Collaboration with Xilinx*
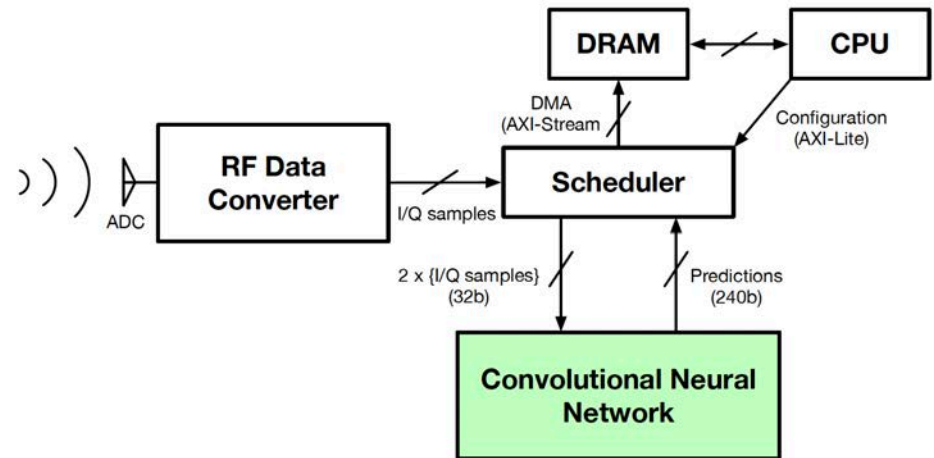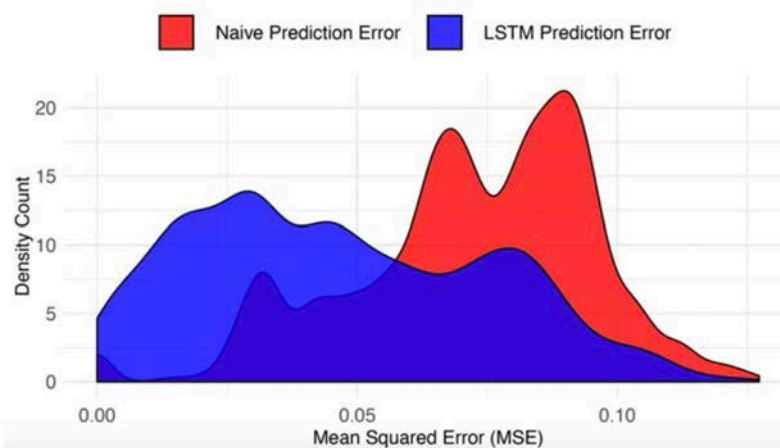


Ours is the most accurate and fastest reported FPGA-based CNN inference implementation CIFAR10: 90.9% acc, 122K fps (TRETS'19)

*Next Generation Technology Fund*

› Processing RF signals remains a challenge

  - FPGAs allow integration of radio, machine learning and signal processing



LSTM Spectral prediction: 4.3 µs latency on Ettus X310 XC7K410T (MILCOM'18)

Ternary Modulation classifier: 488K class/s, 8us latency, Xilinx ZCU111 RFSoC (FPT'19)

*Defence Innovation Hub*

› Implementation of a neuromorphic high dynamic range camera-based object detector on FPGAs

› Significantly improved accuracy in high contrast situations

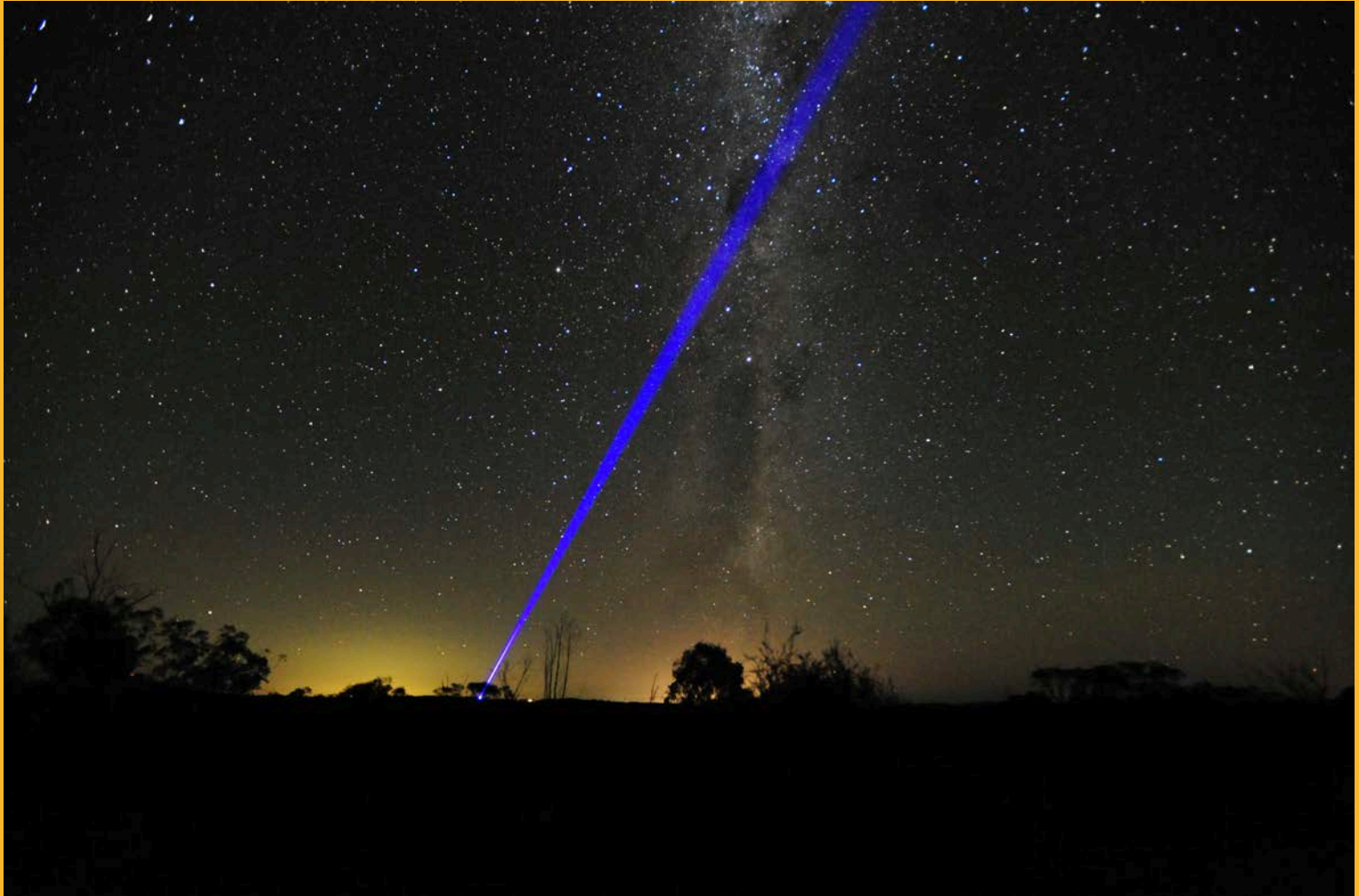Overview

FPGAs

Applications

Our work

› Industry Trends

- Cloud/edge unification

- **More** Sensors (video and hyperspectral); **more** nodes (edge devices/servers) generating data; **more** computation (DNNs, Monte Carlo methods); **more** bandwidth

- Real-time AI and data science applied at all levels

› FPGAs has advantages for these types of problems



Figure: Microsoft