# 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

FPGA Technology

Applications

Our work

FPGA Technology

Applications

Our work
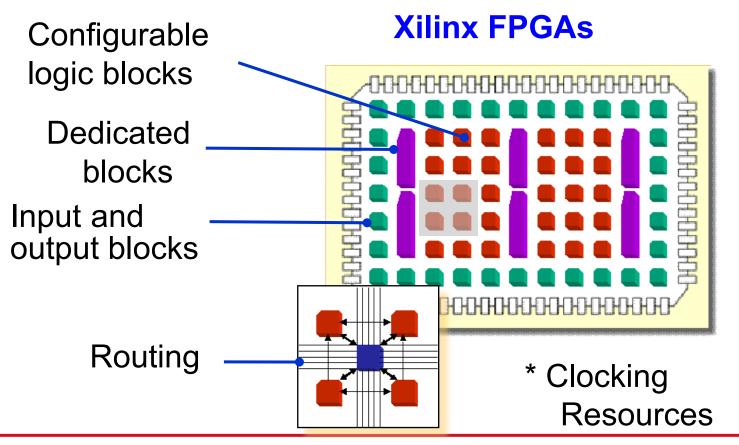
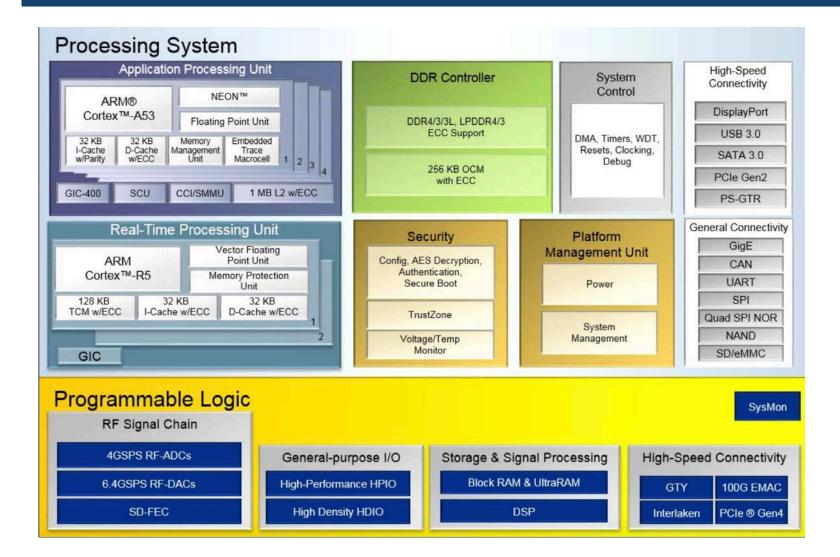User-customisable integrated circuit

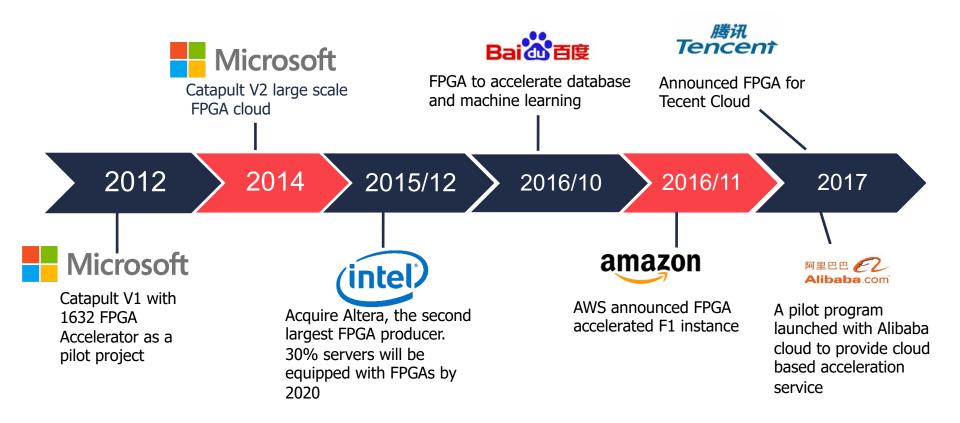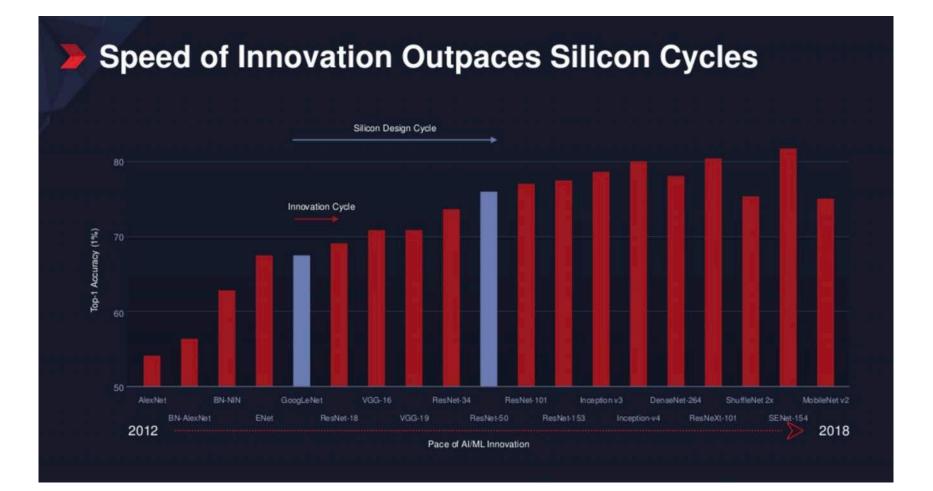› Dedicated blocks: memory, transceivers and MAC, PLLs, DSPs, ARM cores

**Xilinx FPGAs**

Configurable
logic blocks

Dedicated
blocks

Input and
output blocks

Routing

* Clocking
Resources

# Recent Uptake in Reconfigurable Computing

**Microsoft**
Catapult V2 large scale FPGA cloud

**Baidu 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

› 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

https://github.com/Xilinx/Vitis_Libraries

Source: Xilinx

Random Forest Classification Training

Dataset:

1 - HEPMASS (https://archive.ics.uci.edu/ml/datasets/HEPMASS)

2 - HIGGS (https://archive.ics.uci.edu/ml/datasets/HIGGS)

| Dataset | Sample Num | Tree Depth | Tree Num | End-to-End (s) | Speedup | Thread num | Spark (s) |
|---------|-----------|-----------|----------|----------------|---------|-----------|-----------|
| 1 | 7000000 | 5 | 512 | 61.20 | 10.2 | 28 | 622.30 |
| 1 | 7000000 | 5 | 1024 | 121.20 | 15.3 | 16 | 1849.724 |
| 2 | 8000000 | 5 | 512 | 70.30 | 13.3 | 28 | 933.83 |
| 2 | 8000000 | 5 | 1024 | 138.84 | 15.5 | 16 | 2154 |

K-Means Clustering Training

Dataset:

1 - NIPS Conference Papers (http://archive.ics.uci.edu/ml/datasets/NIPS+Conference+Papers+1987-2015)

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

Source: Xilinx

FPGA Technology

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: Geoff Hall, Imperial College

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.





Source: Microsoft

## World's fastest cloud network

Source: Microsoft https://docs.microsoft.com/en-us/azure/networking/azure-network-latency

› 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



Spark on 12-core server ■ SDA

**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

Source: Xilinx

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

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

FPGA Technology

Applications

Our work

**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)

- Highest capacity (FPGA'18)

http://phwl.org/assets/papers/dknlms_trets20.pdf

*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)

http://phwl.org/assets/papers/amc_raw20.pdf

*Defence Innovation Hub*

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

› Significantly improved accuracy in high contrast situations

See paper for details

### Algorithm 1: DNN Training

Define: layer $l$; time $t$; 8-bit weights $\bar{W}_l^t$; input activations $x_l^t$;
deltas $\nabla x_l^t$; weight updates $\nabla W_l^t$; quantisation
functions $Q_w$, $Q_a$, $Q_e$; quantisation scaling coefficients
$qw$, $qa$, $qe$; gemm inputs $A$, $B$; gemm output $C$;
batch size $K$;

1. Forward:
   Software:
   1  $\bar{x}_l^t, qa = Q_a(x_l^t)$; $B = im2col(\bar{x}_l^t)$;
   Hardware:
   2  $A = (\bar{W}_l^t)^T$;
   3  $C = tofloat(gemm(A, B), qw, qa)$;
   Software:
   4  $x_{l+1}^t = C$;
2. Backward:
   Software:
   5  $\nabla \bar{x}_{l+1}^t, qe = Q_e(\nabla x_l^t)$; $tmp = im2col(\bar{x}_l^t{}^T)$;
   6  for $i = 1, 2, \ldots, K$ do
   Hardware:
   7  $A = \nabla \bar{x}_l^t(i)^T$; $B = tmp(i)$;
   8  $C = tofloat(gemm(A, B), qe, qa)$;
   Software:
   9  $\nabla W_l^t \mathrel{+}= C$;
   10  end
   Hardware:
   11  $A = \bar{W}_l^t$; $B = \nabla \bar{x}_l^t$;
   12  $C = tofloat(gemm(\ldots$
   Software:
   13  $\nabla x_l^t = col2im(C)$;

Float (A53 only) — 680, 172, 5.6
8-bit (A53+FPGA) — 38.6, 9.89, 1.58

Time per iteration (sec) vs Batch Size (1, 32, 128)

- **17x speed-up over ARM**

## FPGA

- Low-Precision (8-bit)
  - All matrix multiplications
  - >95% of DNN operations

## ARM

- High-Precision
  - Everything else!
  - Of particular importance is the **weight update and gradient accumulator**

- Suits a Zynq platform
  - Fast DDR, shared between PL and floating-point

http://phwl.org/assets/papers/lptrain_fpt19.pdf

FPGA Technology

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