# ECE618
# Hardware Accelerators for Machine Learning (Spring 2022)

## Lecture 1: Course Information &
## Machine Learning and FPGA Accelerator Recap

Weiwen Jiang, Ph.D.

Electrical and Computer Engineering

George Mason University

wjiang8@gmu.edu

# Agenda

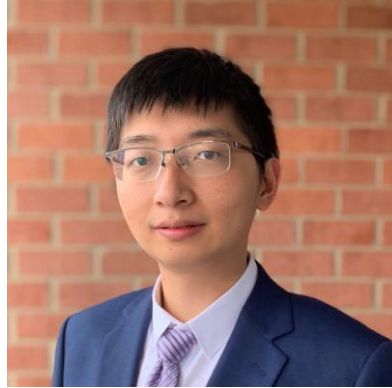**Course Information**

Course Resources

Course Policy

Tools for Lab

Motivation and Schedule

# Course Information

| | |
|---|---|
| **Instructor** | Dr. Weiwen Jiang |
| **E-Mail** | wjiang8@gmu.edu |
| **Phone** | (703)993-5083 |
| **Lecture Time** | **Monday 19:20 - 22:00** |
| **Location** | **Room 1002, Music/Theater Building** |
| **Office Hour** | Monday 16:30 - 17:30 |
| **Office** | Room 3247, Nguyen Engineering Building |
| **Zoom** | http://go.gmu.edu/zoom4weiwen |
| **Backup Course Zoom** | https://go.gmu.edu/ece618 **(Need Permission First)** |

# About Me.



Dr. Weiwen Jiang

- **Background**
  - Researcher at University of Pittsburgh (2017-2019)
  - Postdoc at University of Notre Dame (2019-2021)
  - George Mason University (2021 - present)
- **Research Interests**
  - HW/SW Co-Design
  - Quantum Machine Learning
- **Contacts:**
  - wjiang8@gmu.edu
  - Nguyen Engineering Building, Room3247
  - (703)993-5083
  - https://jqub.ece.gmu.edu/

# Teaching Assistant



Yi Sheng (Ph.D. Candidate)

ysheng2@gmu.edu

https://jqub.ece.gmu.edu/yi/

Office Hours: TBD

# Course Description

Covers the **hardware design** principles to **deploy** different machine learning algorithms. The emphasis is on understanding the fundamentals of **machine learning and hardware architectures** and determine plausible methods to **bridge them**.

Topics include precision scaling, in-memory computing, hyperdimensional computing, architectural modifications, GPUs and vector architectures, quantum computing as well as recent hardware programming tools such as **Xilinx AI Vitis, Xilinx HLS, and IBM Qiskit**.

# Recommend Prerequisite

- **ECE 554: Machine Learning for Embedded Systems**

- Good C programming
  - Especially required for FPGA-related project

- Familiar with Python and PyTorch

# Agenda

Course Information

**Course Resources**

Course Policy

Tools for Lab

Motivation and Schedule

# Course Resources

- **Blackboard:**
  - Assignments will be posted and submitted here!
  - Online discussion, shared documents, announcements.
    - Do NOT upload codes in discussion.
- **Course Website:**
  - https://jqub.ece.gmu.edu/2022/01/01/HA4ML/
  - Course information (TA time, location, zoom, etc.)
  - Slides, readings, and documents will be posted here!

# Agenda

Course Information

Course Resources

**Course Policy**

Tools for Lab

Motivation and Schedule

# Grading Policy

- Midterm Exam                                    10%
- Final Exam                                       20%
- Research Paper Presentation          20%
- Assignments and Labs                     20%
- Project                                             30%

# You Have Been Warned. Zero Tolerance!

- **No matter vaccinated or not, face mask is required in class**



NOTICE
FACE MASK
REQUIRED

- **Request to a Zoom access for a few classes if needed**

# You Have Been Warned. Zero Tolerance!

- **Lecture content and materials should NOT go online without explicit permission**



- **No plagiarism!**

  The most common sense of way interpreting no plagiarism:
  **You need to DO your work.**

# Agenda

Course Information

Course Resources

Course Policy

**Tools for Lab**

Motivation and Schedule

# Tools for lab

Google Colab

Xilinx High-Level Synthesis

IBM Qiskit

# Agenda

Course Information

Course Resources

Course Policy

Tools for Lab

**Motivation and Schedule**

# What Software to Be Accelerated? --- MLP/CNN

**Supervised Learning**

**Example:** Classification

**Training**

**Given:** <u>Labeled</u> data as training dataset

$(x_i, y_i)$: $x_i$ training data, $y_i$: label

$x_i =$ 　　$y_i = 3$

**Output:** A learned function $f$ from X to Y

$$f: x \mapsto y$$

**Inference/Execution**

**Given:** Unseen data test dataset
　　　　A learned function $f$

**Do:** $f($  $) = 3$

# What Software to Be Accelerated? --- MLP/CNN



INPUT 32x32 — C1: feature maps 6@28x28 — S2: f. maps 6@14x14 — C3: f. maps 16@10x10 — S4: f. maps 16@5x5 — C5: layer 120 — F6: layer 84 — OUTPUT 10

Convolutions — Subsampling — Convolutions — Subsampling — Full connection — Full connection — Gaussian connections

- Local receptive fields
- Shared weights
- Pooling (subsampling)

**Cat?**

**Dog?**

# What Software to Be Accelerated? --- RNN

**Supervised Learning**

**Example:** Classification

**Training**

**Given:** Labeled data as training dataset

$(x_i, y_i)$: $x_i$ training data, $y_i$: label

$$x_i = \text{(waveform)} \qquad y_i = \text{"can I"}$$

**Output:** A learned function $f$ from X to Y

$$f: x \mapsto y$$

**Inference/Execution**

**Given:** Unseen data test dataset
A learned function $f$

**Do:** $f(\text{(waveform)}) = \text{"brown fox"}$



VoxCeleb

*A large scale audio-visual dataset of human speech*

**SEP-28k Dataset**



Actual speech
Can I feed my d-dog uh uh [...] pea-ea-nut butter?
Intended speech
Can I feed my dog peanut butter?

# What Software to Be Accelerated? --- RNN



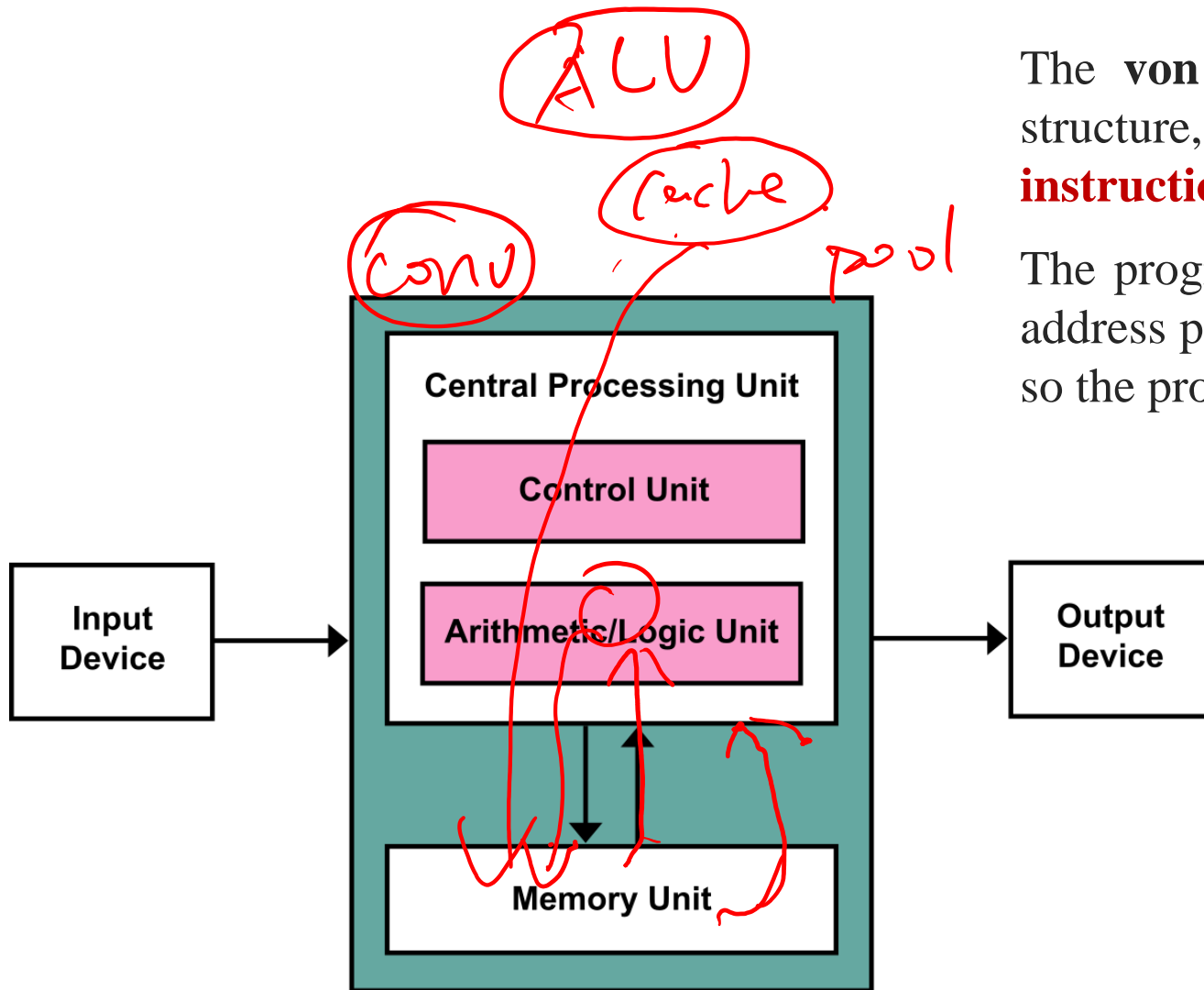| one to many | many to one | many to many | many to many |
|:---:|:---:|:---:|:---:|
| Image captioning (image → words) | Sentiment classification (words → sentiment) | Translation (words → words) | Video frame classification (frames → classes) |

# What Hardware Will Be Covered in This Course?

The **von Neumann structure**, also known as the Princeton structure, is a **memory structure** that merges **program instruction memory and data memory together**.

The program instruction memory address and the data memory address point to different physical locations in the same memory, so the program instruction and data are of the same width.



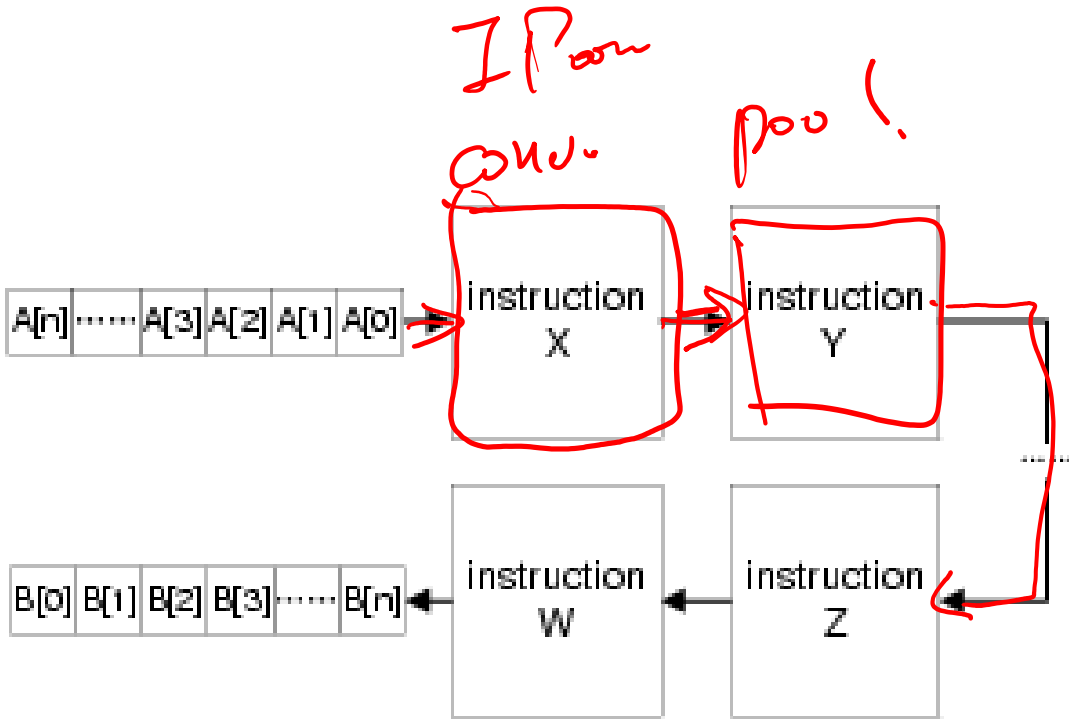Intel's 12th Gen "Alder Lake" 10nm Desktop CPU

NVIDIA RTX A6000 Workstation Graphics Card (in my lab)

ODROID-XU4 Single Board Computer with Quad Core 2GHz A15, 2GB RAM

NVIDIA Jetson Nano

# What Hardware Will Be Covered in This Course?

**Streaming architecture:** data items are pushed in and out as sequential streams, the **instructions are mapped into programmable circuit units** along the path from the input ports to output ports. Therefore, instead of fetching instructions and data back and forth from the memory, the computation gets performed as the **data streams flow** through the circuit units in one pass.
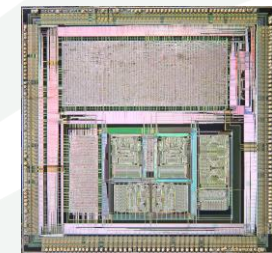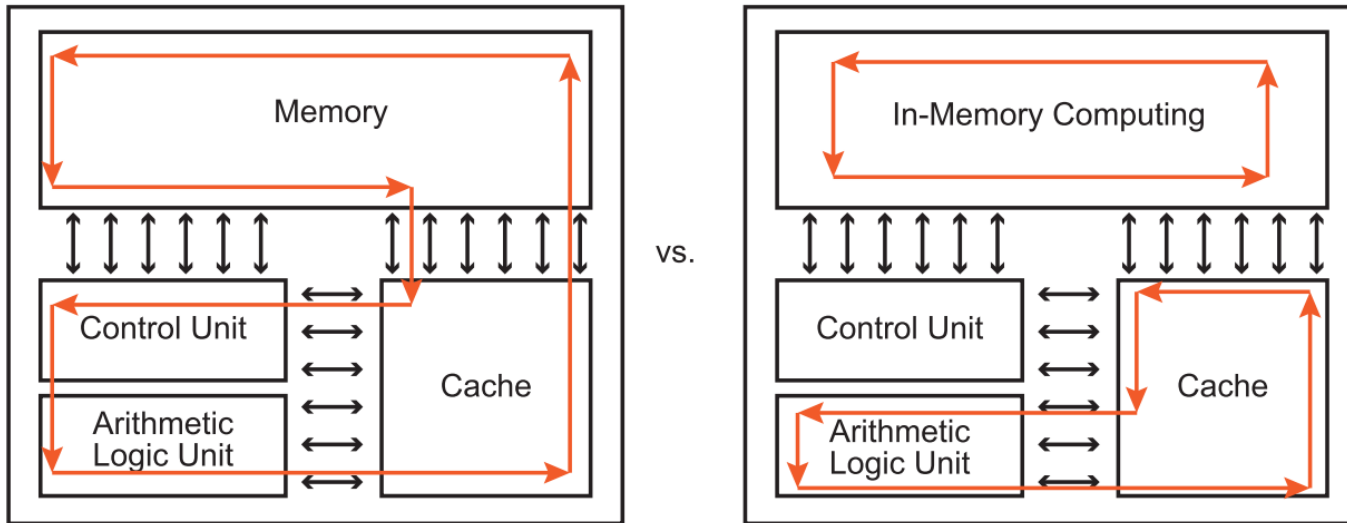


PYNQ

ZCU Series (102, 104, 106)

Xilinx Alveo U280 Data Center
Accelerator Card

ASIC

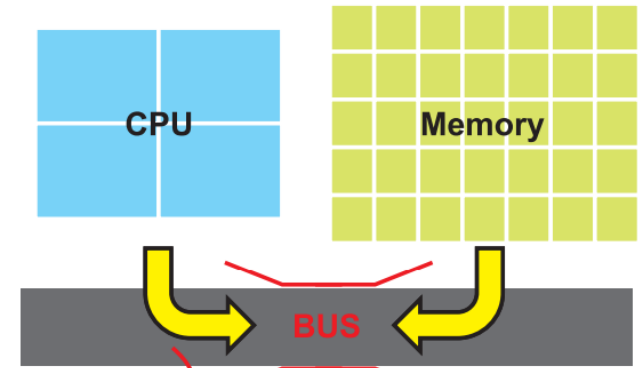# What Hardware Will Be Covered in This Course?

**In-memory computing** is the technique of **running computer calculations entirely in computer memory** (e.g., in RAM).
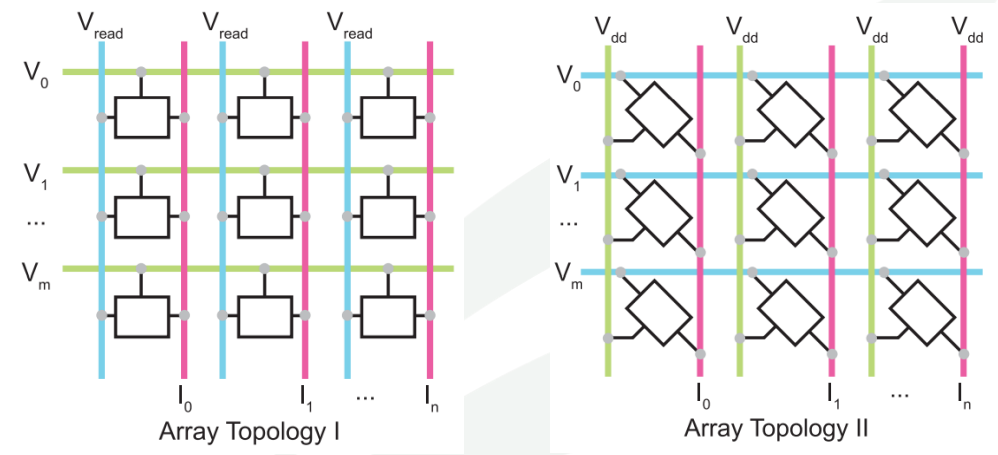


von Neumann Architecture vs. In-Memory Computing

immediate data flow — bus connection — WL (gate) — SL — BL

[**Yan,** Advanced Intelligent Systems 2019]



von Neumann bottleneck:
1. Memory data
2. Intermedia data
3. Computing results
4. Control Signal

Array Topology I   Array Topology II

# What Hardware Will Be Covered in This Course?

**Quantum computing** is a type of computation that harnesses the collective properties of **quantum states**, such as **superposition**, **interference**, and **entanglement**, to perform calculations.
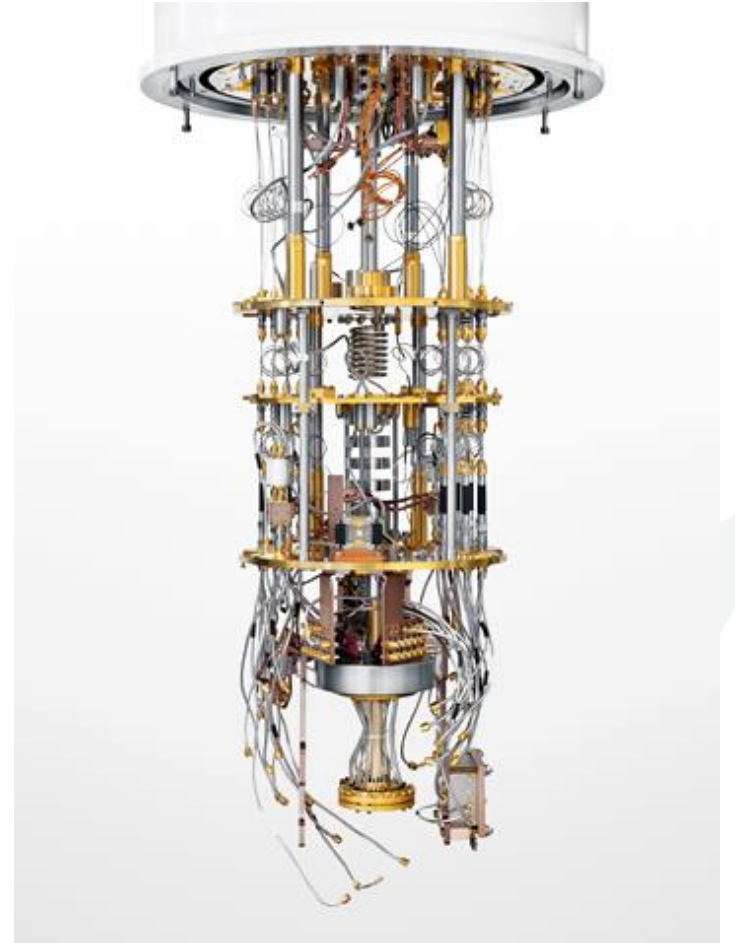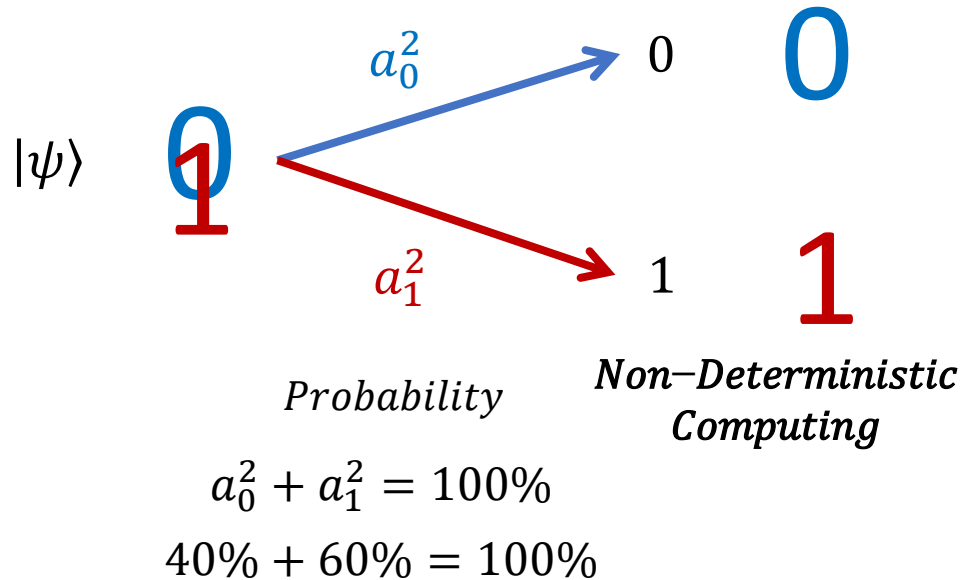
**Classical Bit**

$$X = 0 \ \textit{or} \ 1$$

**Quantum Bit (Qubit)**

$$|\psi\rangle = |0\rangle \ \textbf{\textit{and}} \ |1\rangle$$

$$|\psi\rangle = a_0|0\rangle + a_1|1\rangle$$

**Reading out Information from Qubit (Measurement)**

$|\psi\rangle$

$a_0^2 \rightarrow 0$   0

$a_1^2 \rightarrow 1$   1

*Non–Deterministic Computing*

*Probability*

$$a_0^2 + a_1^2 = 100\%$$

$$40\% + 60\% = 100\%$$

# Why Need Specialized Hardware Accelerators?

- Specialized High-Efficiency Computing!
- Why specialization?
  - **Power constraint of modern computers**
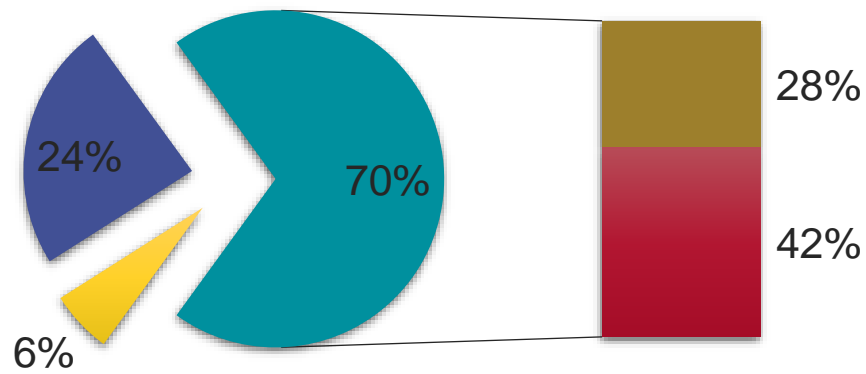


| << 1W | ~ 1W | ~ 15W | ~ 50W | ~ 100W | ~ 100W |

# Why Need Specialized Hardware Accelerators?

- Specialized High-Efficiency Computing!

- Why specialization?
  - Power constraint of modern computers
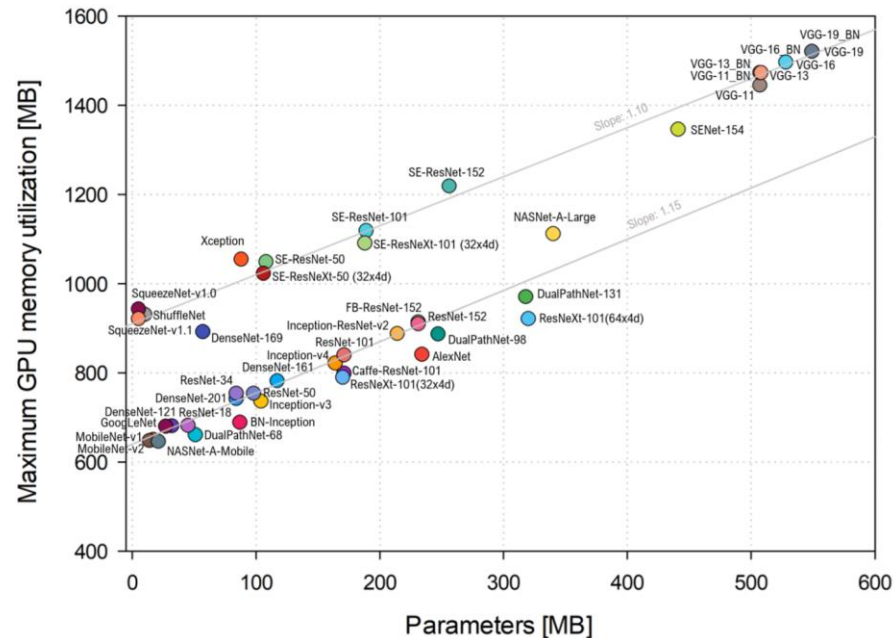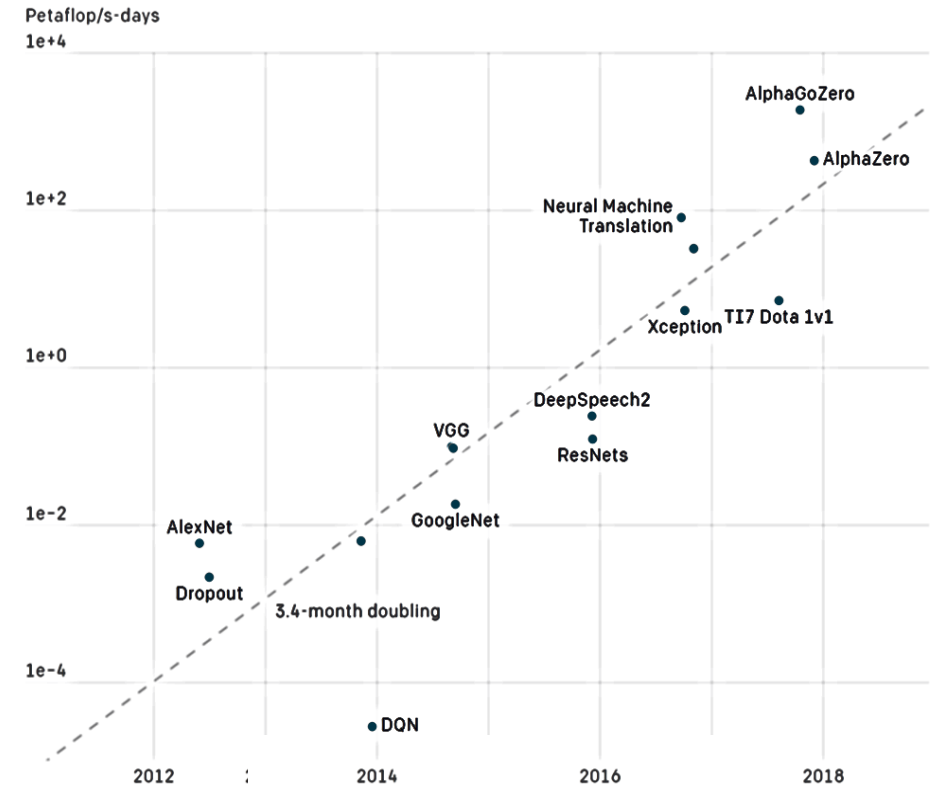  - **In-efficiency of general–purpose computing**

**Embedded Processor Energy Breakdown**

■ Arithmetic  ■ Clock and control  ■ Data supply  ■ Instruction supply

24%
6%
70%
28%
42%

Dr. Weiwen Jiang, ECE,  [Images credit]: Prof. Callie Hao @ GATech  University

# Why Need Specialized Hardware Accelerators?

- ## Specialized High-Efficiency Computing!

- ## Why specialization?
  - Power constraint of modern computers
  - In-efficiency of general–purpose computing
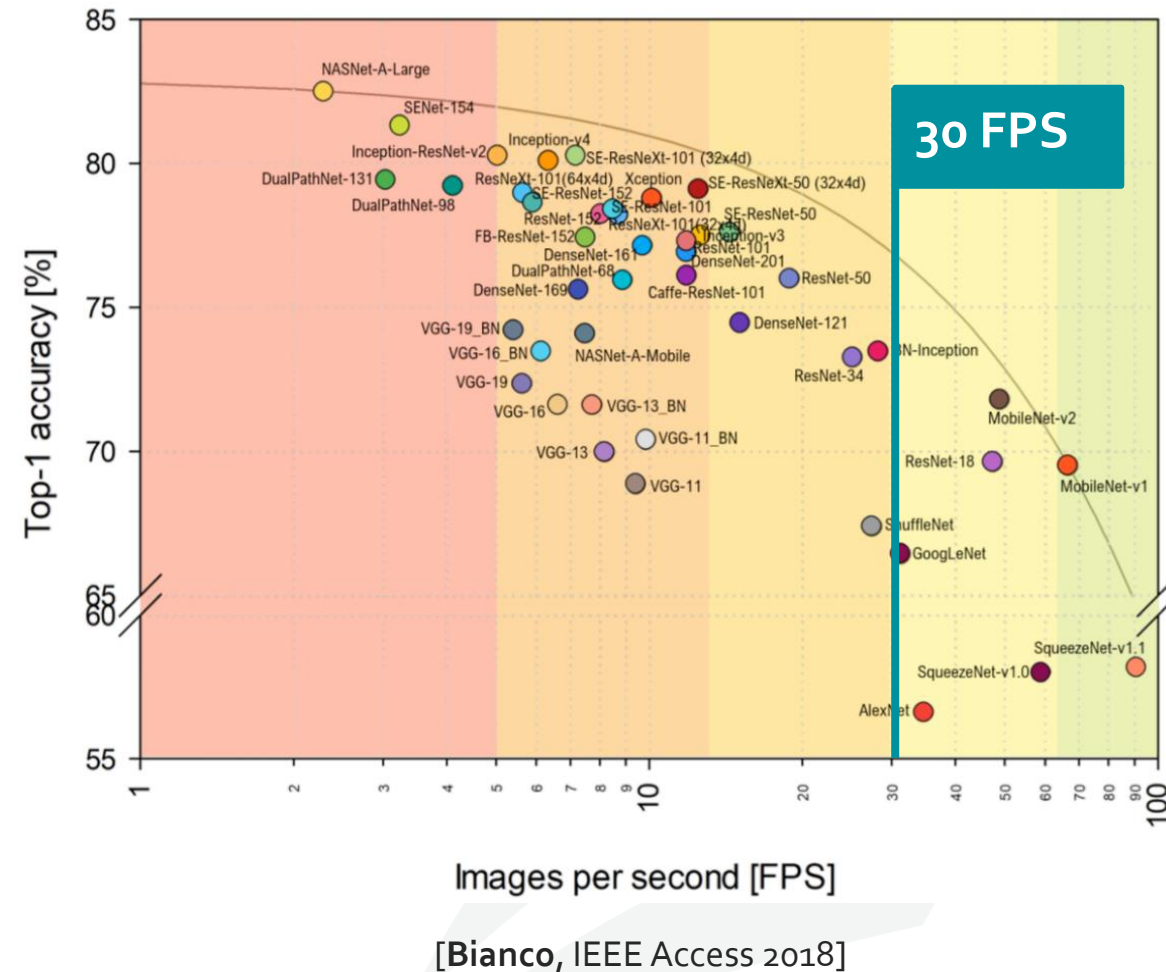  - **Data and computation explosion (big data, AI)**



[**Bianco,** IEEE Access 2018]



https://openai.com/blog/ai-and-compute/

[Images credit]: Prof. Callie Hao @ GATech

# Why Need Specialized Hardware Accelerators?

- Specialized High-Efficiency Computing!

- Why specialization?
  - Power constraint of modern computers
  - In-efficiency of general–purpose computing
  - Data and computation explosion (big data, AI)
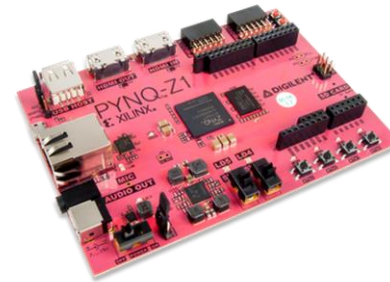  - **Real-time processing requirement**



[**Bianco,** IEEE Access 2018]

Dr. Weiwen Jiang, ECE, GMU  [Images credit]: Prof. Callie Hao @ GATech  ersity

# An Overview of Hardware Accelerators
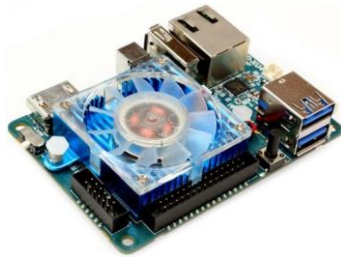
Intel's 12th Gen "Alder Lake" 10nm
Desktop CPU

NVIDIA RTX A6000 Workstation
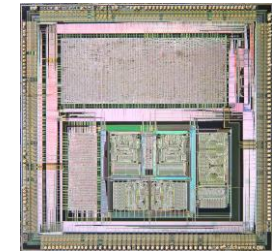Graphics Card (in my lab)

PYNQ

ZCU Series (102, 104, 106)

ODROID-XU4 Single Board Computer with
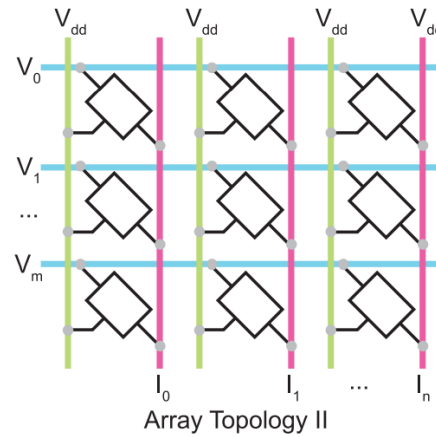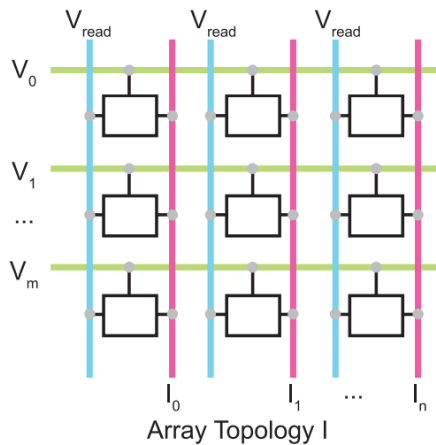Quad Core 2GHz A15, 2GB RAM

NVIDIA Jetson Nano

Xilinx Alveo U280 Data Center
Accelerator Card

ASIC

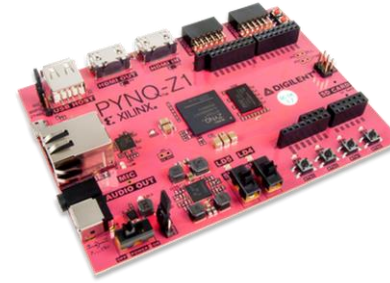Array Topology I

Array Topology II

# Schedule



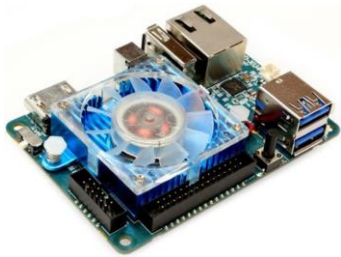Intel's 12th Gen "Alder Lake" 10nm Desktop CPU



NVIDIA RTX A6000 Workstation Graphics Card (in my lab)
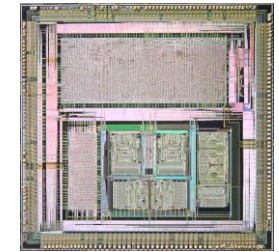


PYNQ



ZCU Series (102, 104, 106)



ODROID-XU4 Single Board Computer with Quad Core 2GHz A15, 2GB RAM



NVIDIA Jetson Nano



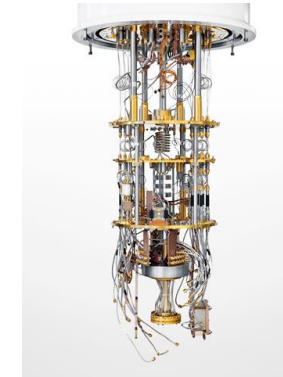Xilinx Alveo U280 Data Center Accelerator Card



ASIC

## Session I: Classical Computing Accelerators for Machine Learning

| Date | Topic |
|------|-------|
| Jan. 24 | Course Information & Machine Learning and FPGA Accelerator Recap |
| Jan. 31 | Vector Architectures, FPGAs and GPU Architectures |
| Feb. 7 | ASIC Accelerators |

# Schedule



Array Topology I

Array Topology II

## Session II: Novel Post-Moore Computing Accelerators for ML

| Date | Topic |
|---|---|
| Feb. 14 | In-Memory Computing Accelerator Design |
| Feb. 21 | Neuromorphic Accelerators |
| Feb. 28 | Hyperdimensional Computing Accelerators |
| Mar. 07 | Quantum Neural Network Accelerators |

# Schedule

## Session III: Other Accelerator Related Topics

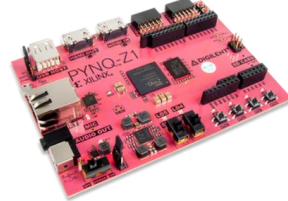| Date | Topic |
|------|-------|
| Mar. 28 | Project Proposal |
| Apr. 04 | Distributed Learning |
| Apr. 11 | Hands-on Accelerator Design (1) |
| Apr. 18 | Project Overview |
| Apr. 25 | Hands-on Accelerator Design (2) |
| May 02 | Project Presentations |
| May 11-18 | Final exam |

# Expectation & Final Project

- Implement ML on any hardware in a team with 1-3 students



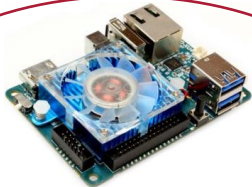Intel's 12th Gen "Alder Lake" 10nm Desktop CPU

NVIDIA RTX A6000 Workstation Graphics Card (in my lab)

PYNQ

ZCU Series (102, 104, 106)

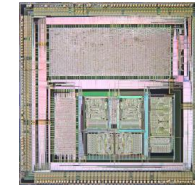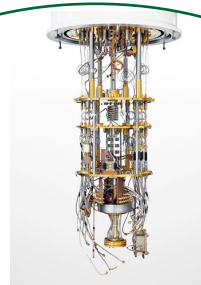ODROID-XU4 Single Board Computer with Quad Core 2GHz A15, 2GB RAM

NVIDIA Jetson Nano

Xilinx Alveo U280 Data Center Accelerator Card

ASIC

Array Topology I

Array Topology II

Order date: Dec 2, 2021    Dell Purchase ID: 2008410187802

Received By Dell    Confirmed    In Production    Shipped    Delivered
Dec 2, 2021    Dec 2, 2021        Revised ship date: Feb 28, 2022    Revised arrival: Mar 7, 2022

# What Did We Learn in ECE 554? (Recap)

## Application
- Computer Vision
- Natural Language
- Games
- ...

## ML/DL Algorithms
- AN
- MLP
- CNN
- RNN
- LSTM
- RL
- Transformer
- MLP-Mixer
- ...

## Optimization
- Network Inference
- Network Training
- Model Compression
- Network Design

## Embedded Systems

## Performance Model
- Resource Usage
- Power
- Latency

## Optimization
- Hardware Design
- Mapping/Scheduling
- Communication

# ECE 554 Course Recap

- **Machine Learning Basis:**
  - Different neural networks: **MLP**, **CNN**, **RNN**, **RL**
  - Training (**Gradient Descent**) and inferencing neural networks using Pytorch
  - Implement convolution using "**for loops**"

# Biological Neuron


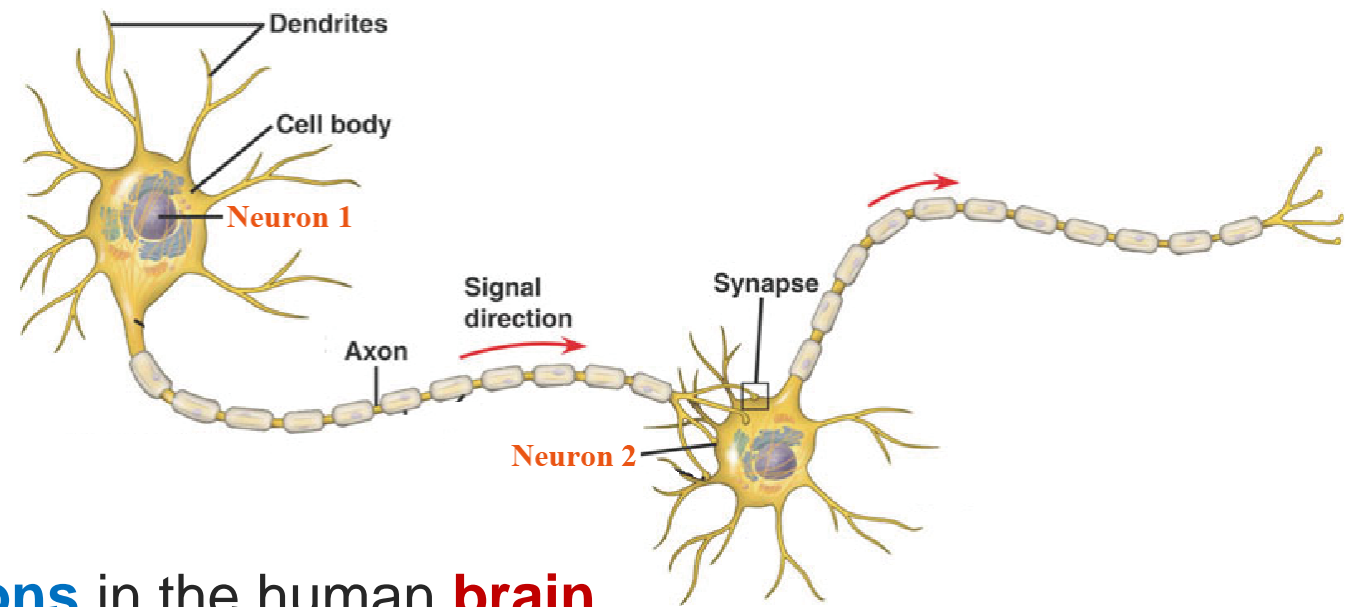Dendrites, Cell body, Neuron 1, Axon, Signal direction, Synapse, Neuron 2

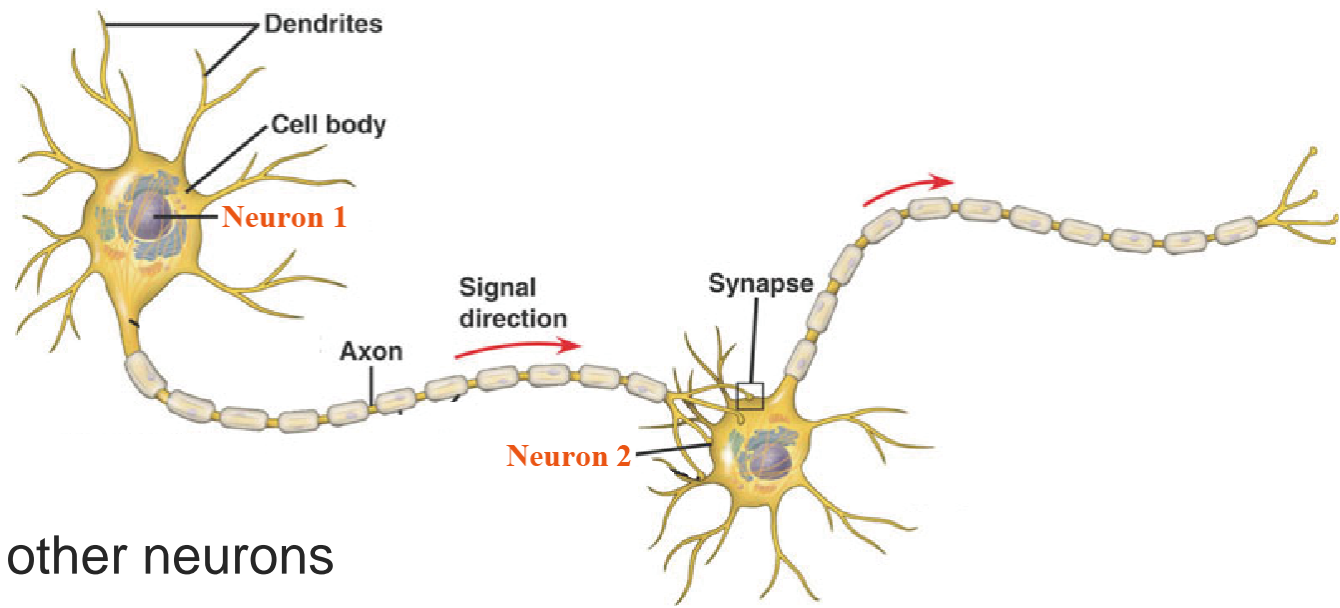**Human intelligence reside**

**in the brain:**

- Approximately **86 billion neurons** in the human **brain**
- The brain is a **network** of **neurons**, connected with nearly $10^{14} - 10^{15}$ **synapses**

**How to equip intelligence in the machine?**

- **To understand how the brain network is constructed**
- **To mimic the brain**

# Biological Neuron



**Neurons work together:**

- **Cell body** process the information
- **Dendrites** receive messages from other neurons
- **Axon** transmit the output to many smaller branches
- **Synapses** are the contact points between axon (Neuron 1) and dendrites (Neuron 2) for message passing

**Cell body** receives input signal from **dendrites** and produce output signal along **axon**, which interact with the next neurons via **synaptic weights**

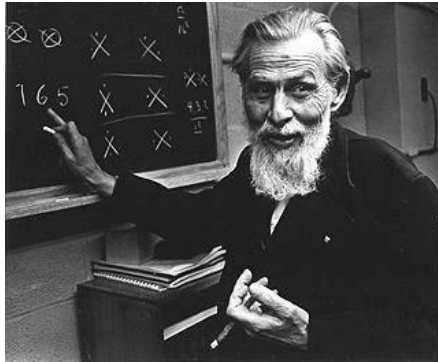**Synaptic weights** are **learnable** to perform useful computations

(e.g., Recognizing objects, understanding language, making plans, controlling the body.)
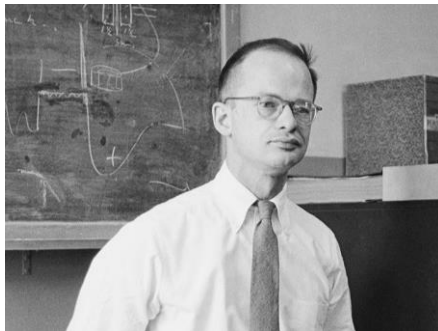
# Artificial Neuron Design

- **Idealized neuron models**

  - Idealization **removes complicated details** that are not essential for understanding the main principles.

  - It allows us to apply **mathematics** and to make **analogies**.

# McCulloch-Pitts (MP) Neuron
## The first computational model of a biological neuron @ 1943



Warren McCulloch

Walter Pitts

**Synapse**

$x_0$

+1

**Dendrite**

**Synapse**

$x_1$

+1

**Dendrite**

**Synapse**

$x_2$

+1

**Dendrite**

**Signal Direction**

**Cell Body**

$g|f$

**Axon**

$y = f(g(x))$

**Assumptions:**

- Binary devices (i.e., $x_i \in \{0,1\}$ and $y \in \{0,1\}$)

- Identical synaptic weights (i.e., +1)

- Activation function $f$ has a fixed threshold $\theta$

# Artificial Neuron Design

- **Idealized neuron models**

  - Idealization removes complicated details that are not essential for understanding the main principles.

  - It allows us to apply mathematics and to make analogies.

- **Break the limitations on MP Neuron**
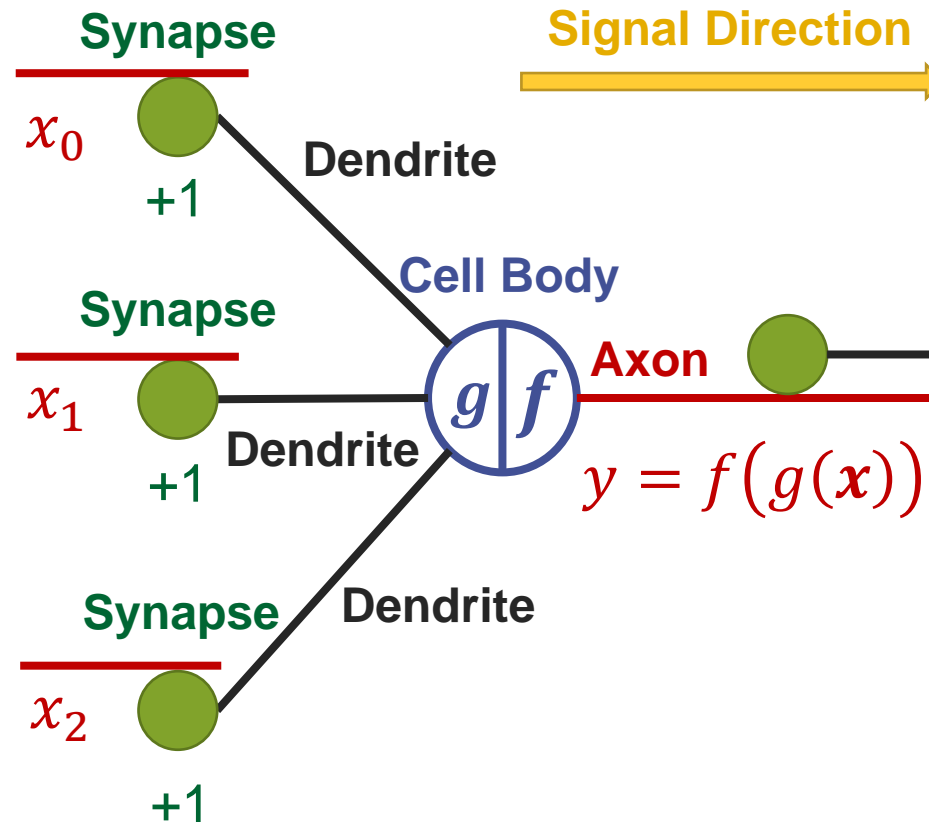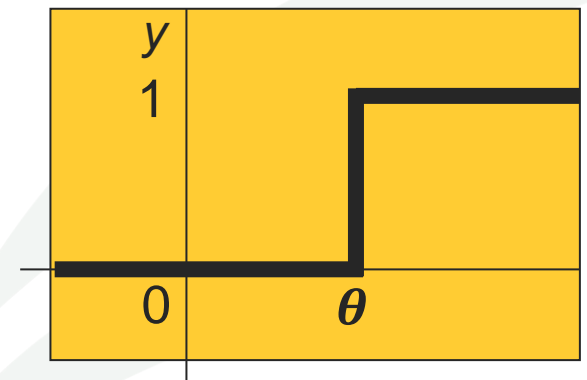
  - What about non-boolean inputs (say, real number)?

  - What if we want to assign more weight (importance) to some inputs?

  - What about functions which are not linearly separable ?

  - Do we always need to hand code the threshold?

# Multi-Layer Perceptron (MLP) – Lecture 2

- Input layer, output layer and hidden layers





A selection of commonly used activation functions for artificial neurons.

# Deep Convolutional Neural Networks (CNN) – <u>Lecture 3</u>

- One of the most widely used types of deep network

- Fully-connected nets treat far apart input pixels same as those close by
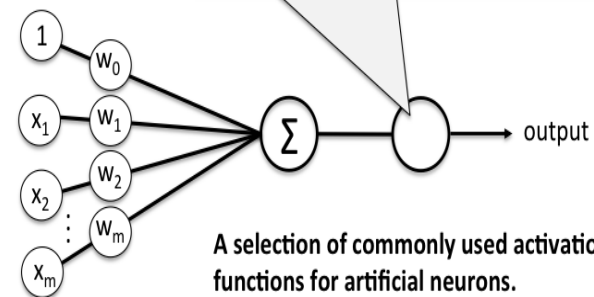    - Hence spatial information must be inferred from the training data

- In contrast, CNN proposes an architecture that inherently tries to take advantage of the spatial structure
    - Such an architecture makes convolutional networks fast to train
    - This, in turn, helps us train even deeper, many-layer networks

- Today, deep convolutional networks or some close variants are used in solving many interesting problems that go beyond image classification

- We will use image classification as a driving use case to explain the main concepts behind CNN

# Convolution – Lecture 3



**1.** Add p number of zeros around the image

**2.** The kernel jumps s pixels when being slid across the image

$S = d$

Input    Kernel    Output

**Parameters:**

- N: input channels
- M: output channels
- K: kernel size
- P: padding size
- S: stride
- *D: dilation*
- *R: rows*
- *C: columns*

CLASS

torch.nn.Conv2d(*in_channels*, *out_channels*, *kernel_size*, *stride=1*, *padding=0*, *dilation=1*, *groups=1*, *bias=True*, *padding_mode='zeros'*, *device=None*, *dtype=None*)

[ref] Aqeel Anwar, What is Transposed Convolutional Layer? https://towardsdatascience.com/what-is-transposed-convolutional-layer-40e5e6e31c11

# From Static Image to Sequences of Data



Speech Signal

Deep Neural Network

Emotions

time=3

w1   w3   w4   w2

time=2

w1   w3   w4   w2

time=1

w1   w3   w4   w2

time=0

*Courtesy to Geff. Hinton*

- Assume each connection has 1 unit delay

- RNN can be unrolled into feedforward networks
  - Each layer keeps on reusing the same weights

$$
\begin{aligned}
\boldsymbol{a}^{(t)} &= \boldsymbol{b} + \boldsymbol{W}\boldsymbol{h}^{(t-1)} + \boldsymbol{U}\boldsymbol{x}^{(t)}, \\
\boldsymbol{h}^{(t)} &= \tanh(\boldsymbol{a}^{(t)}), \\
\boldsymbol{o}^{(t)} &= \boldsymbol{c} + \boldsymbol{V}\boldsymbol{h}^{(t)}, \\
\hat{\boldsymbol{y}}^{(t)} &= \mathrm{softmax}(\boldsymbol{o}^{(t)}),
\end{aligned}
$$

$$
h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)
$$

$$
y_t = W_{hy}h_t
$$

From GoodFellow et al.
Deep Learning

# ECE 554 Course Recap

- **Machine Learning Basis:**
  - Different neural networks: **MLP**, **CNN**, **RNN**, **RL**
  - Training (**Gradient Descent**) and inferencing neural networks using Pytorch
  - Implement convolution using "**for loops**"

- **Put Machine Learning onto Embedded Systems:**
  - Introduction to HLS (Lec 8-9)
    - Using **MLP** as example in class
    - Using **CNN** as example in Labs, which is based on the "**for loop**" implementation
  - Model compression on FPGA: pruning and quantization (Lec 10-11)
  - Neural architecture search (Lec 12)
    - Using **RNN-based RL** as controller/optimizer
    - Using **Gradient Descent** approach for optimization
  - Data movement in HLS-based FPGA implementation (Lec 13)
  - Co-explore neural architectures and FPGA design (Lec 14)

# High-Level Synthesis: HLS –

- High-Level Synthesis
  - Creates an RTL implementation from C, C++, System C, OpenCL API C kernel code
  - Extracts control and dataflow from the source code
  - Implements the design based on defaults and user applied **directives**

- Many implementation are possible from the same source description
  - Smaller designs, faster designs, optimal designs
  - Enables design exploration



**Accelerates Algorithmic C to RTL IP integration**

C, C++, SystemC, OpenCL API C

Vivado™ HLS

VHDL or Verilog

System IP Integration

- Algorithmic Specification
- Micro Architecture Exploration
- RTL Implementation
- Comprehensive Integration with the Xilinx Design Environment

# C Validation and RTL Verification – Lecture 8

- There are two steps to verifying the design
  - Pre-synthesis: C *Validation*
    - Validate the algorithm is correct
  - Post-synthesis: RTL *Verification*
    - Verify the RTL is correct
- C validation
  - A **HUGE** reason users want to use HLS
    - Fast, free verification
  - Validate the algorithm is correct **before** synthesis
    - Follow the test bench tips given over
- RTL Verification
  - Vivado HLS can co-simulate the RTL with the original test bench

# AXI_Stream – Lecture 13

Tr
Tc
K
Tn
Tm
FPGA_DATA
DMA_DATA
read_data
  input_dma_I
  % HLS INTERFACE axis register both port=input_dma_I
  output
  [] IFM
  ifm_input_dma
    data
    last
  I0
    I1
      I2
  O0
    O1
      O2

read_data_0



+ ap_ctrl
− input_dma_I_V
▶ input_dma_I_V_TVALID
◀ input_dma_I_V_TREADY
▶ input_dma_I_V_TDATA[63:0]
ap_clk
ap_rst_n

Vitis™ HLS

output_r_ap_vld
output_r[31:0]

Read_data (Pre-Production)

```cpp
1  #include <ap_fixed.h>
2  #include <hls_stream.h>
3
4  const int Tr=4;
5  const int Tc=4;
6  const int K=3;
7  const int Tn=3;
8  const int Tm=6;
9
10 //typedef ap_fixed<16,8,AP_TRN_ZERO, AP_SAT> FPGA_DATA;
11 typedef float FPGA_DATA;
12 struct DMA_DATA{
13     FPGA_DATA data;
14     bool last;
15 };
16
17 void read_data(hls::stream<DMA_DATA> &input_dma_I, FPGA_DATA *output){
18
19     static FPGA_DATA IFM[Tn][Tr+K-1][Tc+K-1];
20
21     DMA_DATA ifm_input_dma;
22     I0:for(int i=0;i<Tn;i++){
23         I1:for(int j=0;j<Tr+K-1;j++){
24             I2:for(int m=0;m<Tc+K-1;m++){
25                 ifm_input_dma=input_dma_I.read();
26                 IFM[i][j][m]=ifm_input_dma.data;
27             }
28         }
29     }
30
31     O0:for(int i=0;i<Tn;i++){
32         O1:for(int j=0;j<Tr+K-1;j++){
33             O2:for(int m=0;m<Tc+K-1;m++){
34                 output[i*(Tr+K-1)*(Tc+K-1)+j*(Tc+K-1)+m] = IFM[i][j][m]+2;
35             }
36         }
37     }
38
39 }
```

# Test Bench – Lecture 13

```
17
18  FPGA_DATA input[Tn*(Tr+K-1)*(Tc+K-1)] = { 0.47902363538742065,0.5932260751724243,0.59
19
20  void read_data(hls::stream<DMA_DATA> &input_dma_W, FPGA_DATA *output);
21
22  int main(){
23      hls::stream<DMA_DATA> input_dma_I("input_dma_I");
24
25      FPGA_DATA y[Tn*(Tr+K-1)*(Tc+K-1)]={0};
26
27      DMA_DATA ifm;
28      for(int i=0;i<Tn;i++){
29          for(int j=0;j<Tr+K-1;j++){
30              for(int m=0;m<Tc+K-1;m++){
31                  ifm.data = input[i*(Tr+K-1)*(Tc+K-1)+j*(Tc+K-1)+m];
32                  if(i==Tn-1 && j==Tr+K-1-1 && m==Tc+K-1-1)
33                      ifm.last = true;
34                  else
35                      ifm.last = false;
36                  input_dma_I.write(ifm);
37              }
38          }
39      }
40      read_data(input_dma_I,y);
41
42      for(int i=0;i<Tn;i++){
43          for(int j=0;j<Tr+K-1;j++){
44              for(int m=0;m<Tc+K-1;m++){
45                  printf("gap: %f\n", input[i*(Tr+K-1)*(Tc+K-1)+j*(Tc+K-1)+m]-y[i*(Tr+K
46              }
47          }
48      }
49      printf("Done!\n");
50      return 0;
51  }
52
```

Golden results generation code for Lab 4

```python
import torch
import torch.nn as nn

Tr=16
Tc=16
K=3
Tn=3
Tm=16

input = torch.rand(1,Tn,Tr+K-1,Tc+K-1)
weight = torch.rand(Tm, Tn, K, K)

conv = nn.Conv2d(3,2,3,bias=False)
conv.weight = torch.nn.Parameter(weight)

output = conv(input)

torch.set_printoptions(precision = 16)
torch.set_printoptions(profile="full")

print("FPGA_DATA input[Tn*(Tr+K-1)*(Tc+K-1)] = {", end=" ")
for elem in list(input.flatten()):
    print(float(elem), end=",")
print("};")

print("FPGA_DATA weight[Tm*Tn*K*K] = {", end=" ")
for elem in list(weight.flatten()):
    print(float(elem), end=",")
print("};")

print("FPGA_DATA output[Tm*Tr*Tc] = {", end=" ")
for elem in list(output.flatten()):
    print(float(elem), end=",")
print("};")

print()
```
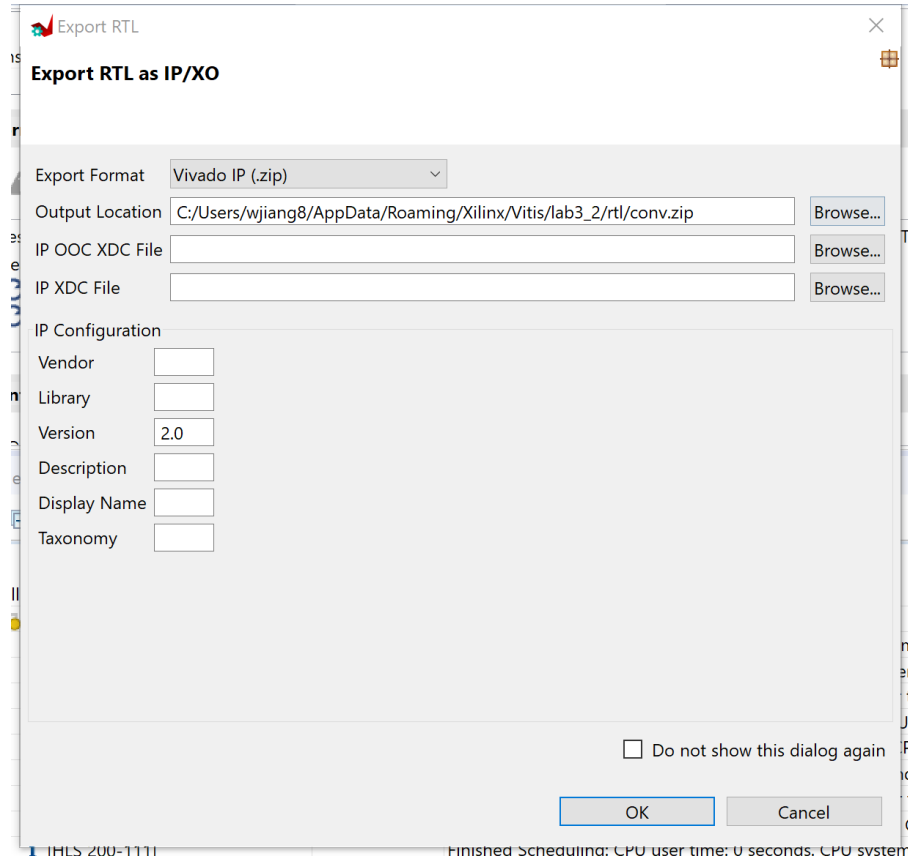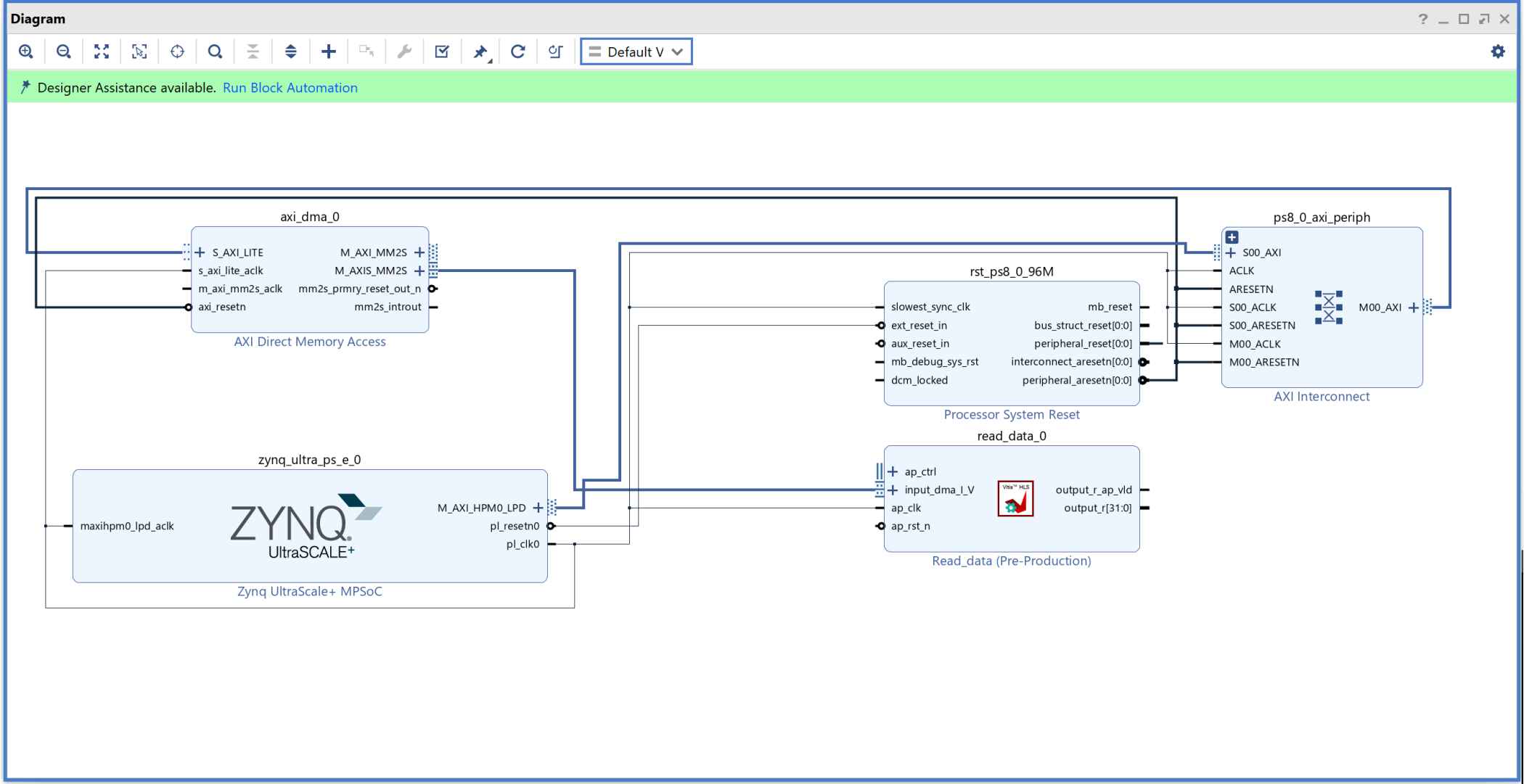
# Export RTL as IP Core – Lecture 13



**Synthesis**

**Export RTL**

**Export RTL as IP/XO**

Export Format: Vivado IP (.zip)

Output Location: C:/Users/wjiang8/AppData/Roaming/Xilinx/Vitis/lab3_2/rtl/conv.zip    Browse...

IP OOC XDC File:    Browse...

IP XDC File:    Browse...

IP Configuration

Vendor:

Library:

Version: 2.0

Description:

Display Name:

Taxonomy:

☐ Do not show this dialog again

OK    Cancel

read_data_0

+ ap_ctrl
+ input_dma_I_V
ap_clk    output_r_ap_vld
ap_rst_n    output_r[31:0]

Read_data (Pre-Production)

# Import the IP into Block Design – Lecture 13

## AI Democratization —— Two Levels



**Level 1: Automation of Neural Network Design**

NAS

NASNet

MNasNet

ProxylessNAS

FBNet

FNAS

......

**Level 2: Automation of AI System Design**

(1) Cloud/Server

(2) Mobile Platform

(3) FPGA Accelerator

(4) ASIC Accelerator

(5) CiM Accelerator

(6) Quantum Comp.

◆ **Cloud / Server**

- Unlimited Resource
- Maximizing Accuracy
- AlexNet, VGGNet, ResNet, …

◆ **Mobile Phones**

- Fixed Hardware
- Accuracy v.s. Latency
- MnasNet, ProxylessNAS, …

◆ **Hardware Accelerators (e.g., FPGA)**

- Hardware Design Flexibility
- Accuracy, Latency, and Energy
- FNAS, SkyNet, EDDNet…

**Datasets / Applications**



**Hardware Platforms**



FPGA

**Neural Networks**

## 1 year for only 1 application

| Name | Year | Acc.(T5) |
|------|------|----------|
| AlexNet | 2012 | 83.4% |
| ZFNet | 2013 | 88.3% |
| VGGNet | 2014 | 92.7% |
| RestNet | 2015 | 96.4% |
| GoogleNet | 2016 | 96.9% |

**Problem**

- **Domain knowledge and excessive labor**

- **It takes too long to devise new architectures**

Manual AI Design

# AutoML: Neural Architecture Search (NAS) – Lecture 14

Sample architecture NN
with probability p



**Controller (RNN)**

**Train from Scratch To Obtain Accuracy A**

Compute gradient of p and
scale it by A to update the
controller

**Reinforcement Learning Based NAS**
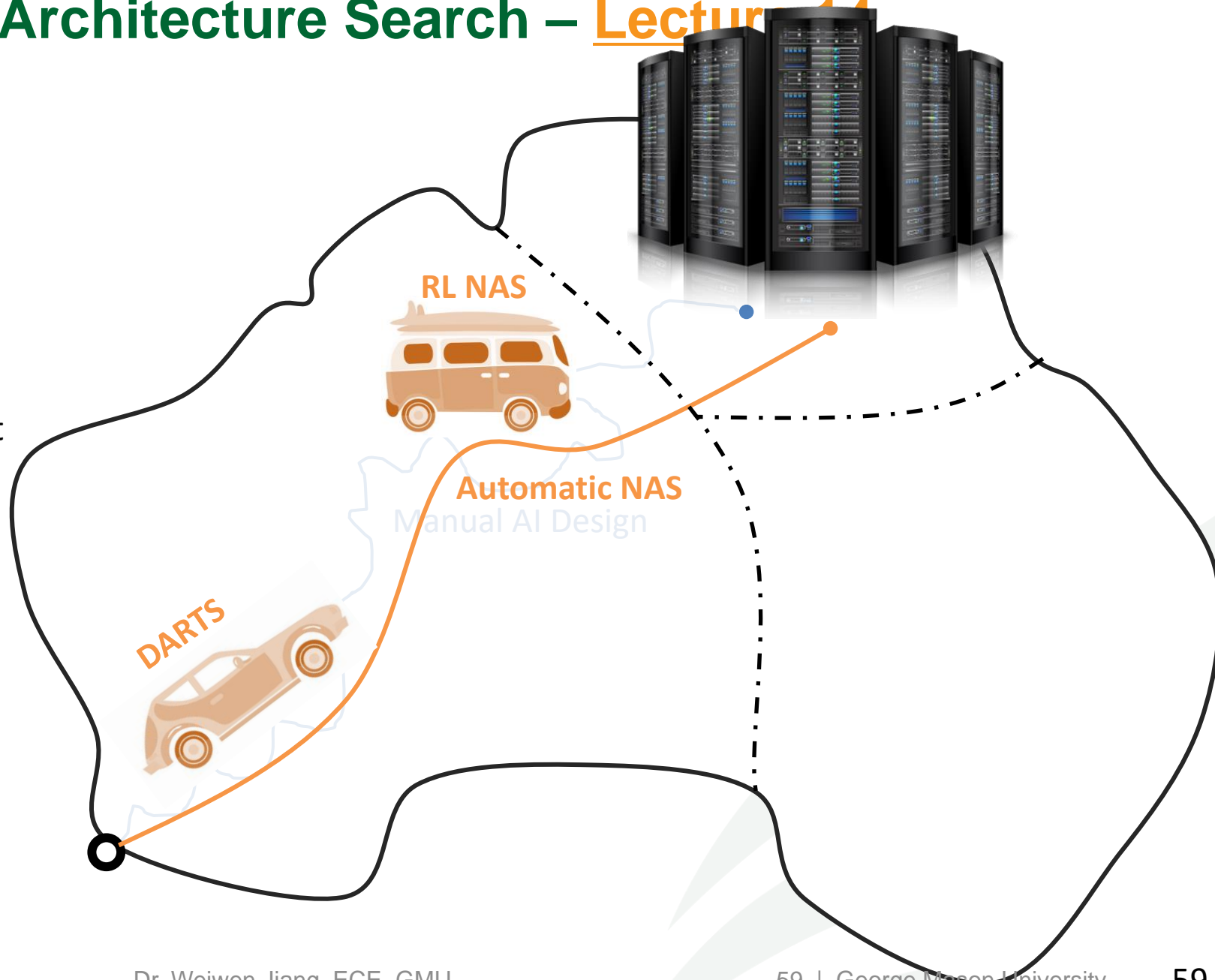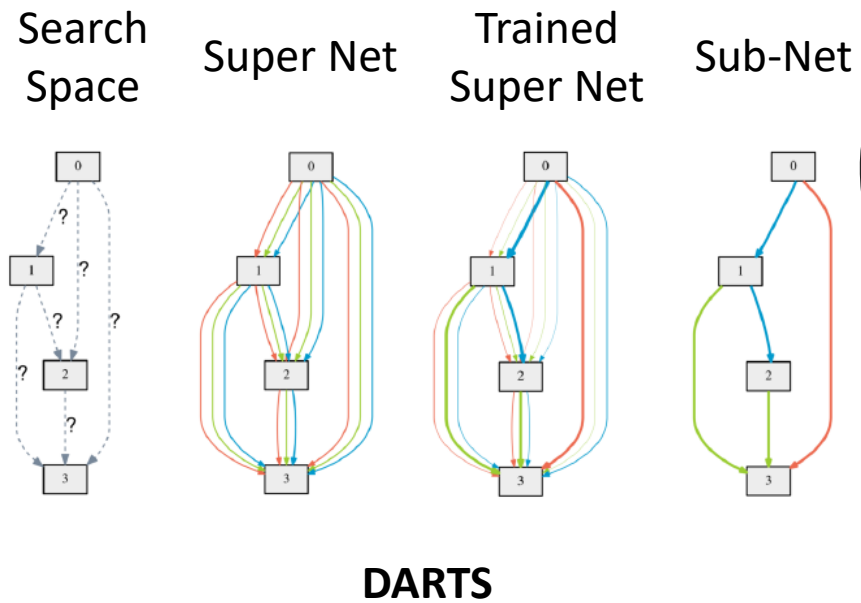
**RL NAS**

**Automatic NAS**

Manual AI Design

**Problem**

- **Low Efficiency, hundreds even thousands of GPU hours**
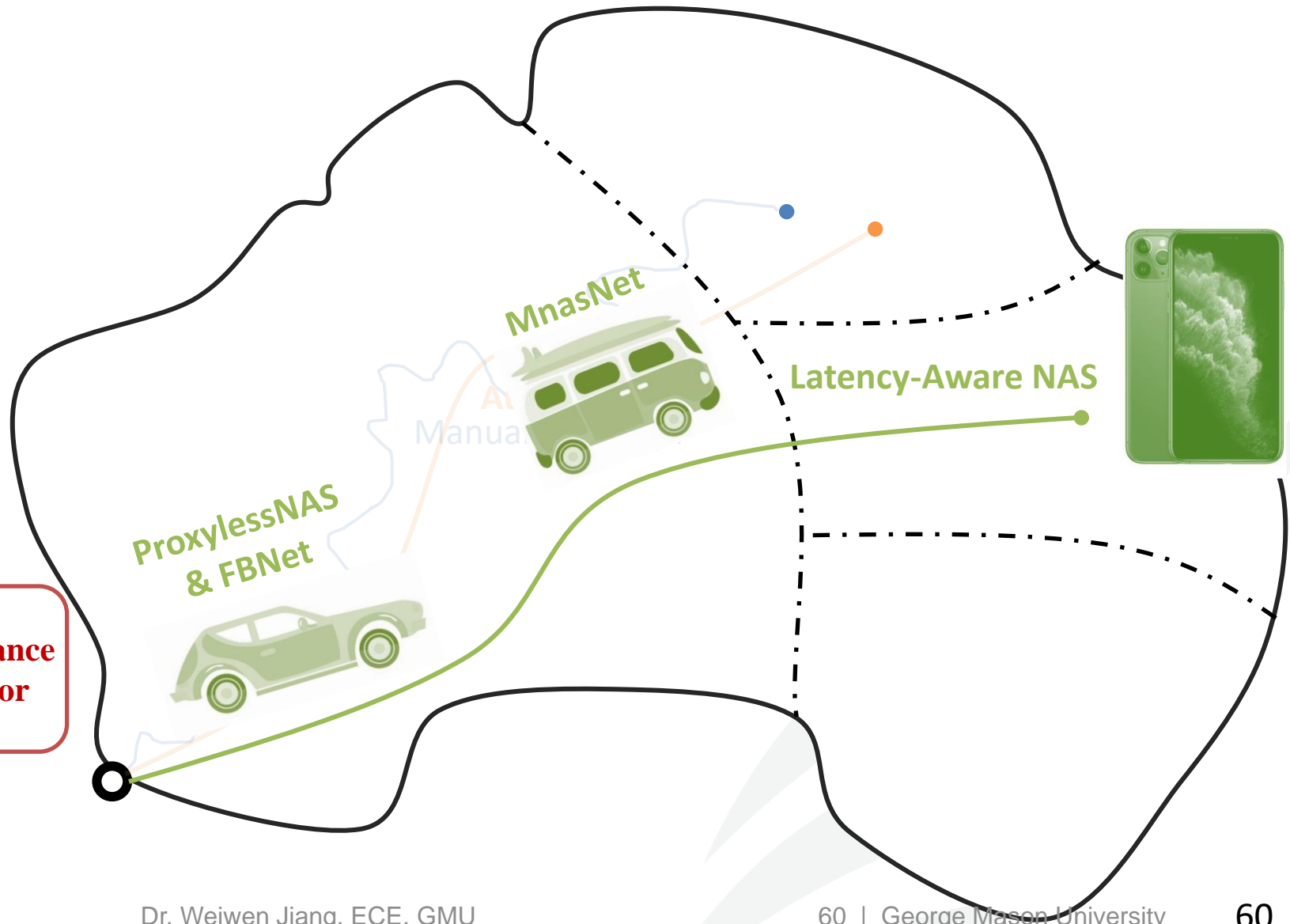- **Mono-Objective: Accuracy, leading network too complicated**

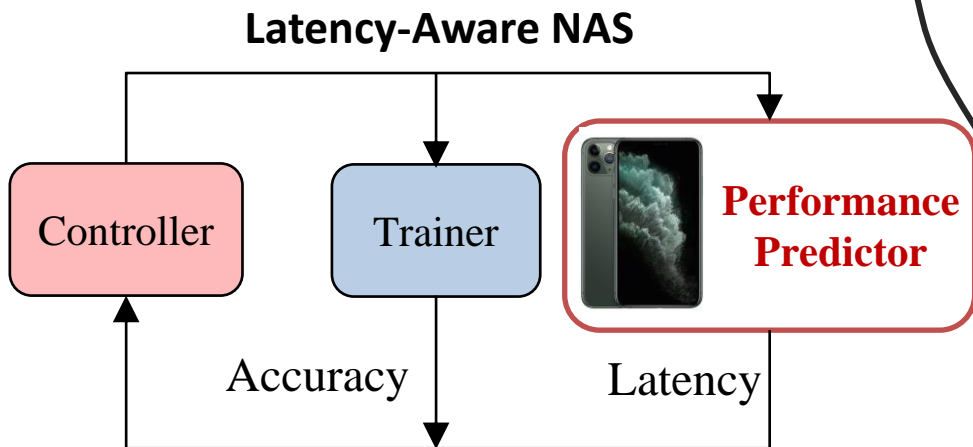| Name | Time |
|------|------|
| DARTS | Jun. 2018 |



Search Space    Super Net    Trained Super Net    Sub-Net

**DARTS**

RL NAS

Automatic NAS

Manual AI Design

DARTS

# AutoML: Hardware-Aware NAS – Lecture 14

| Name | Time |
|------|------|
| MnasNet | Jul. 2018 |
| ProxylessNAS | Dec. 2018 |
| FBNet | Dec. 2018 |

**Latency-Aware NAS**



MnasNet

Latency-Aware NAS

ProxylessNAS & FBNet

Controller — Trainer — **Performance Predictor**

Accuracy — Latency

# AutoML: Network-FPGA Co-Design Using NAS – Lecture 14

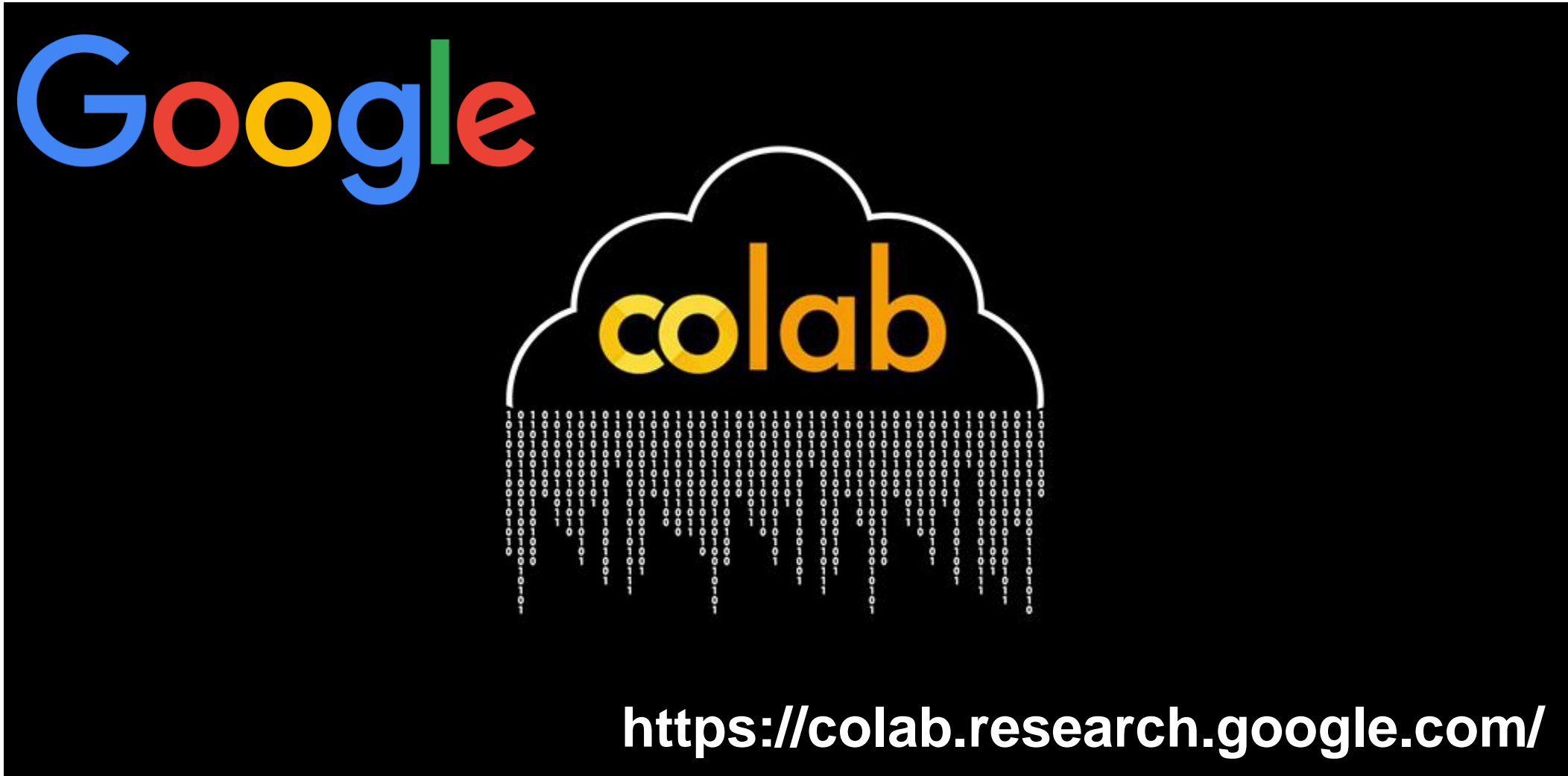| Name | Time |
|------|------|
| **FNAS (ours)** | Jan. 2019 |
| DNN/FPGA | Apr. 2019 |
| SkyNet | Sep. 2019 |
| EDDNet | May. 2020 |

**FNAS - DAC'19** (Best Paper Nomination)
**TCAD'20** (Best Paper Award)

# How to Conduct Neural Architecture Search – <u>Lecture 14</u>

- **Selection of the Backbone Architecture**
  - **VGG** (NAS with RL, FNAS), **GoogLeNet** (NASNet), **MobileNet** (FBNet, ProxylessNAS), etc.

- **Determination of the Search Space**
  - **Software:** Number of Channels, Kernel Size, Convolution Type, etc.
  - **Hardware**: Loop Titling Parameters, Loop Order, Schedule, etc.

- **Optimization Approaches**
  - **Deep Reinforcement Learning:** RNN based controller
  - **Gradient Descent:** DARTS
  - **Metaheuristics:** Swarm

- **Optimization Objective(s):**
  - **Software:** Accuracy, Robustness, Fairness, etc.
  - **Hardware:** Latency, Chip Area, Energy Efficiency, etc.

# Programming Platform



**https://colab.research.google.com/**

**GMU.EDU**

**George Mason University**

4400 University Drive
Fairfax, Virginia 22030

Tel: (703)993-1000