# SOC Design of a Neural Network for Real-Time Semantic Segmentation of 2Kx1K@60fps Video

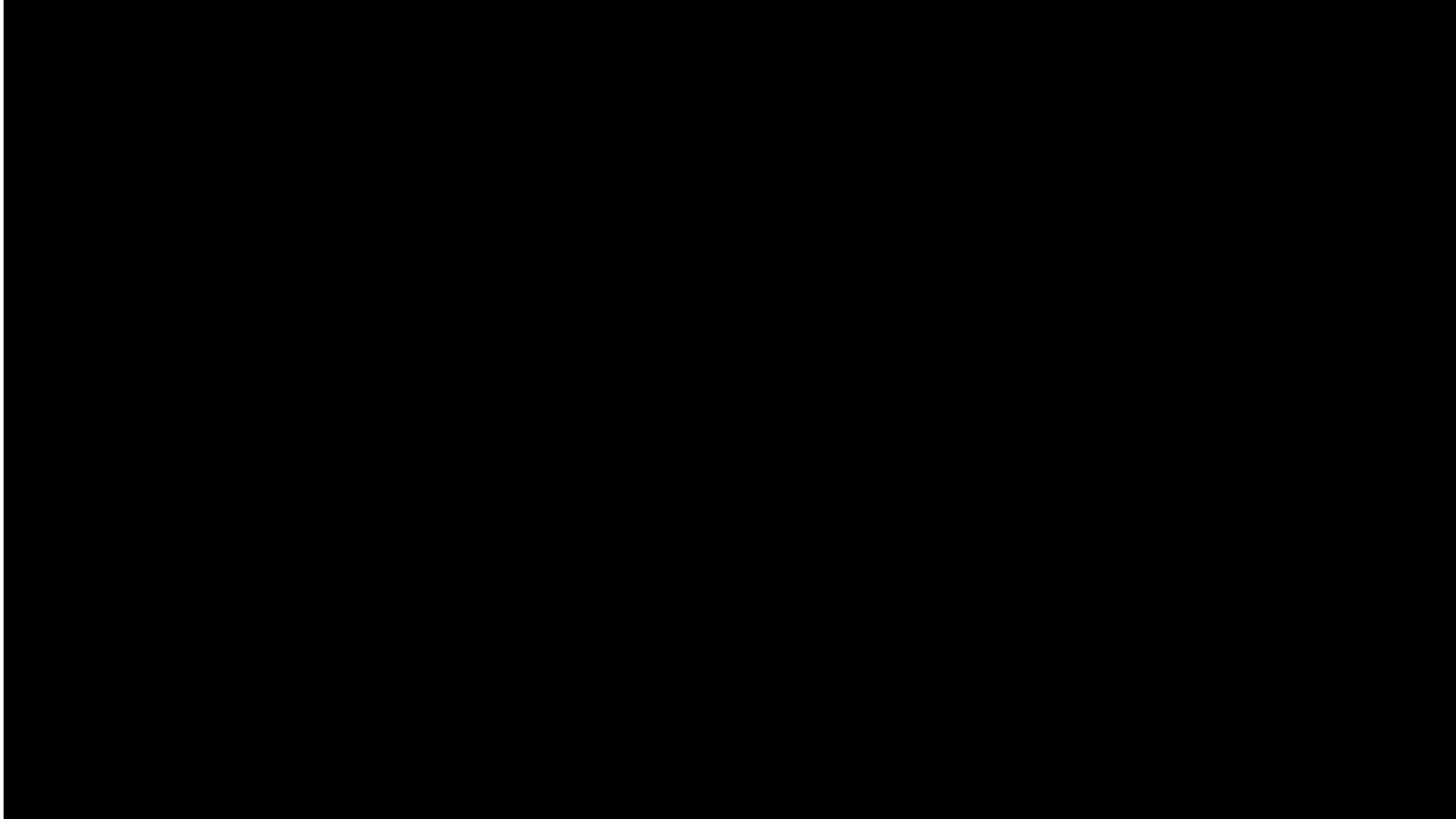Youn-Long Lin

Department of Computer Science

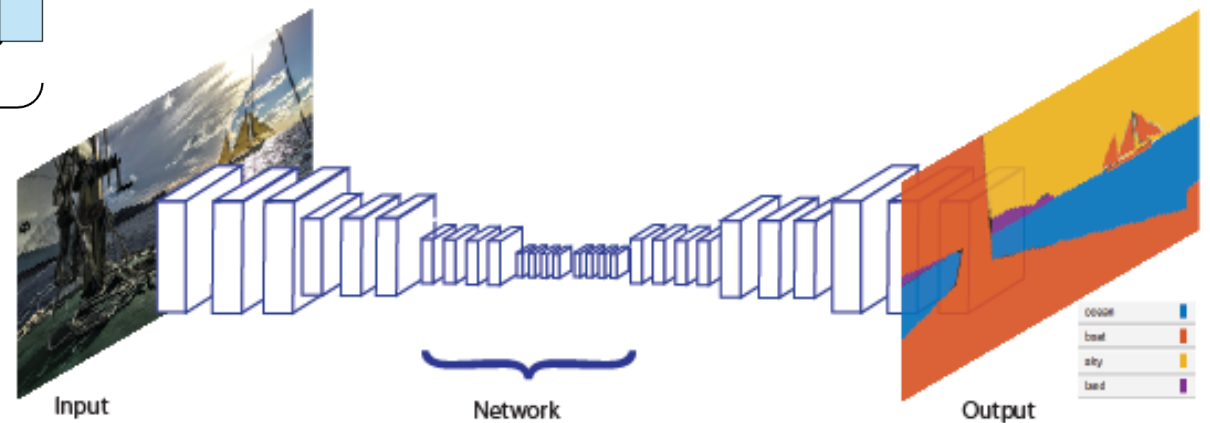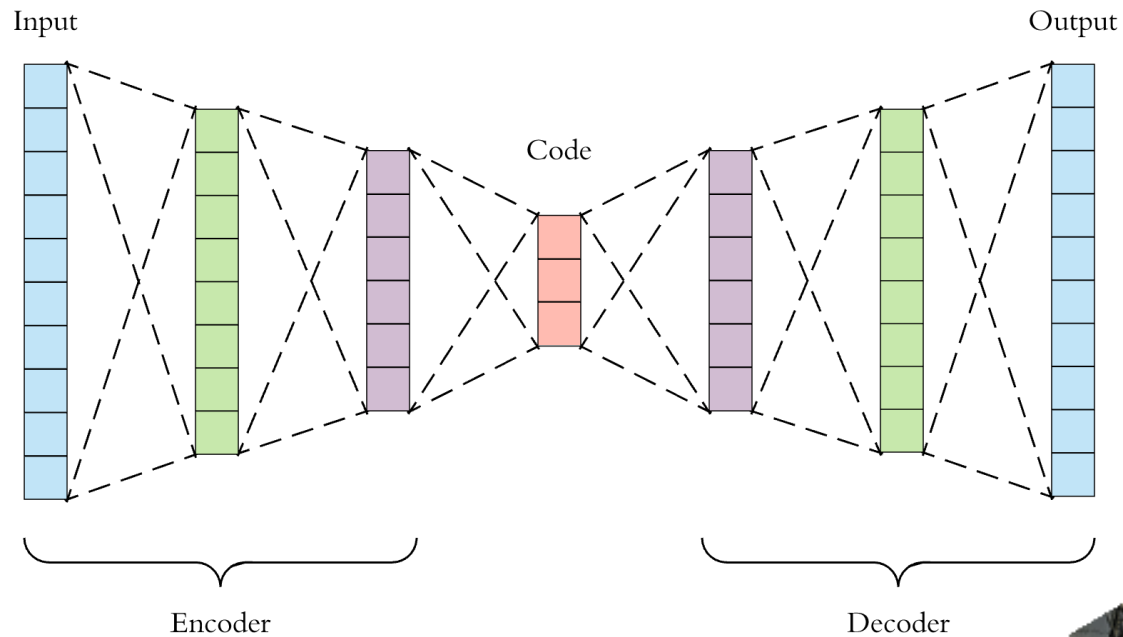National Tsing Hua University

# Project Goals

1. Real-time Semantic Segmentation of HD Video (1Kx2K@60fps)

2. Hardware-friendly neural network architecture

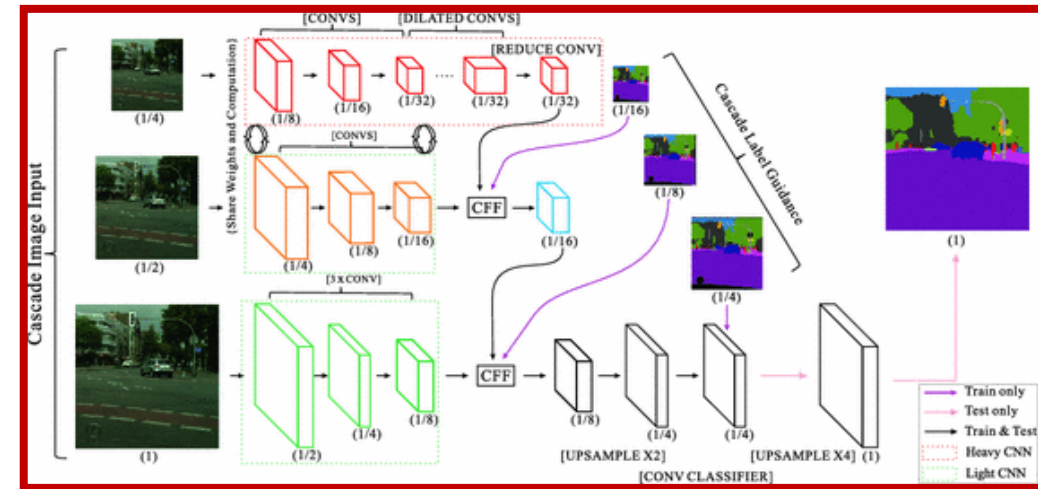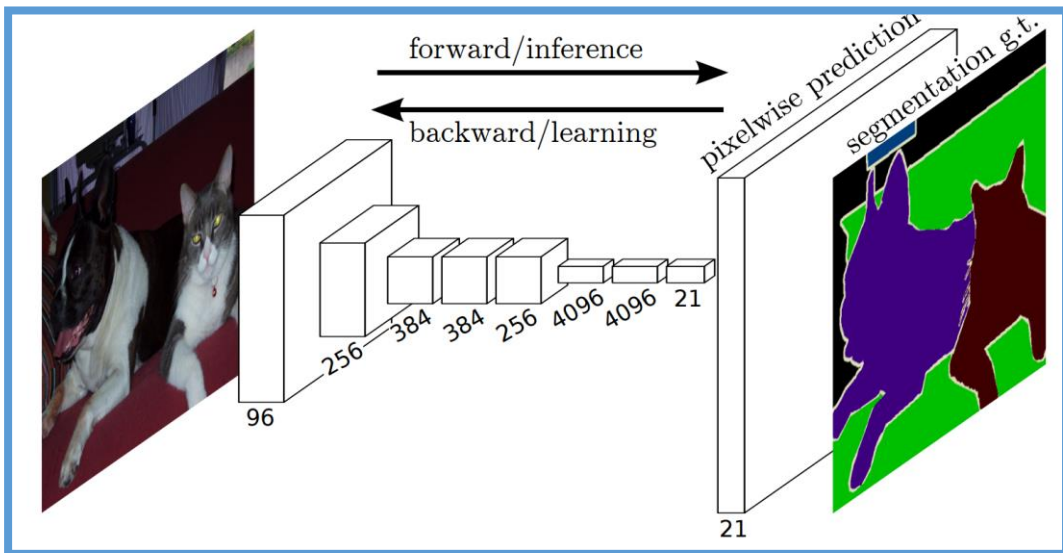3. Proof-of-Concept using GPU

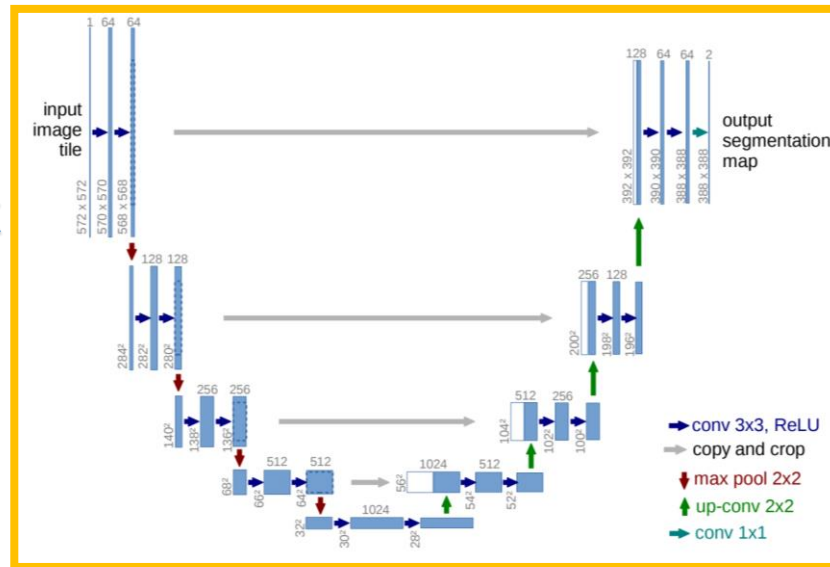4. ASIC

# Semantic Segmentation – CityScapes Dataset

# How to do it?
# Autoencoder, Naturally



Input

Output

Code

Encoder

Decoder

Input

Network

Output

# FCN

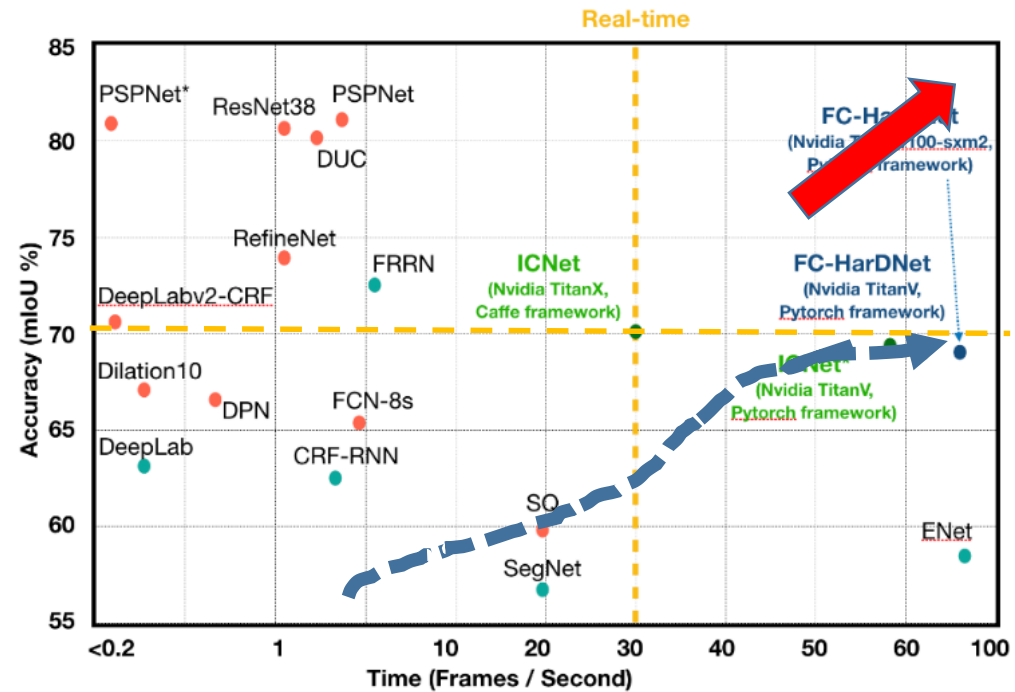## Fully Convolutional Network

# ICNet

# U-Net

- Long, J., Shelhamer, E., & Darrell, T. (2015). Fully convolutional networks for semantic segmentation. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 3431-3440).

- Ronneberger, O., Fischer, P., & Brox, T. (2015, October). U-net: Convolutional networks for biomedical image segmentation. In International Conference on Medical image computing and computer-assisted intervention (pp. 234-241). Springer, Cham.

- Zhao, H., Qi, X., Shen, X., Shi, J., & Jia, J. (2018). Icnet for real-time semantic segmentation on high-resolution images. In Proceedings of the European Conference on Computer Vision (ECCV) (pp. 405-420).
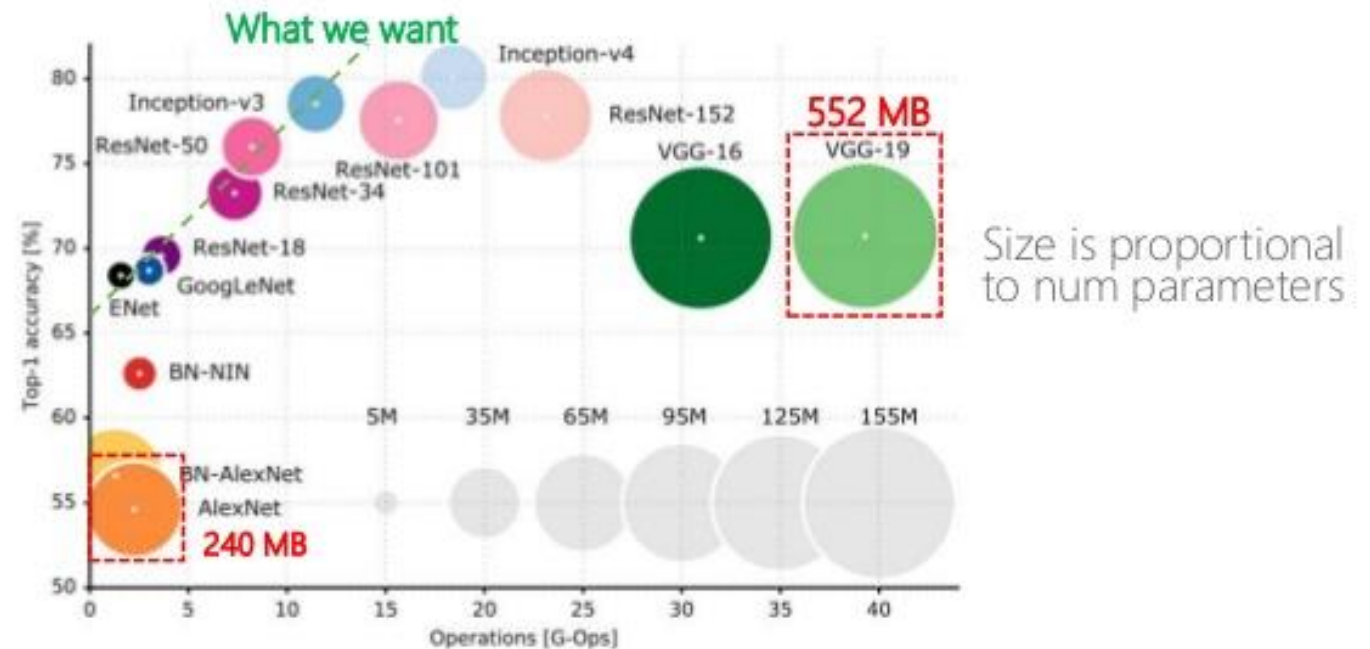
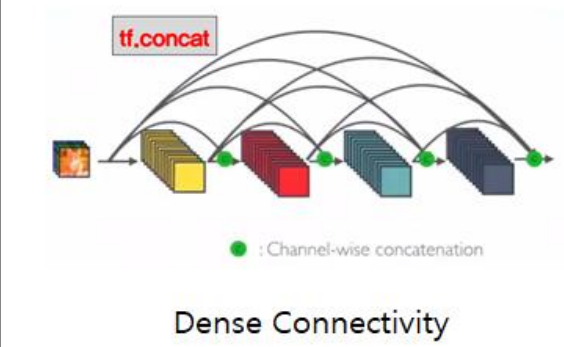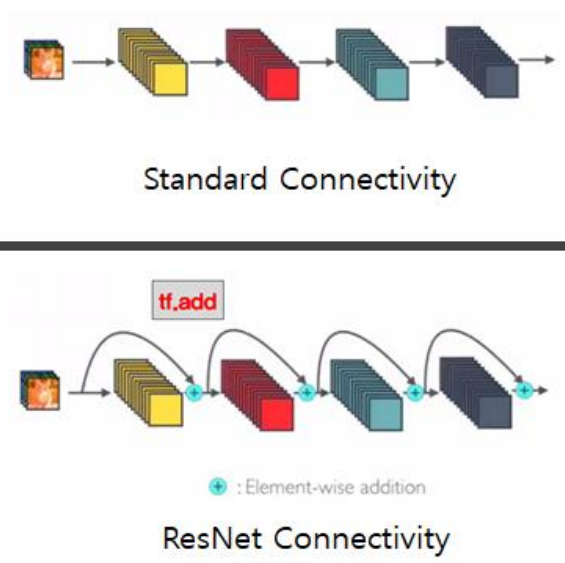# What we want for semantic segmentation?

## Accuracy vs frames per second

mean IoU
(Intersect over Union)

# Network design tradeoffs

- Performance (Accuracy)
- Cost
  - Hardware gate count
  - Network size (Parameters)
- Inferencing
  - Time (# Operations?)
  - Energy consumption

## Accuracy vs Operations Per Image Inference



Alfredo Canziani, Adam Paszke, Eugenio Culurciello, "An Analysis of Deep Neural Network Models for Practical Applications" 2016
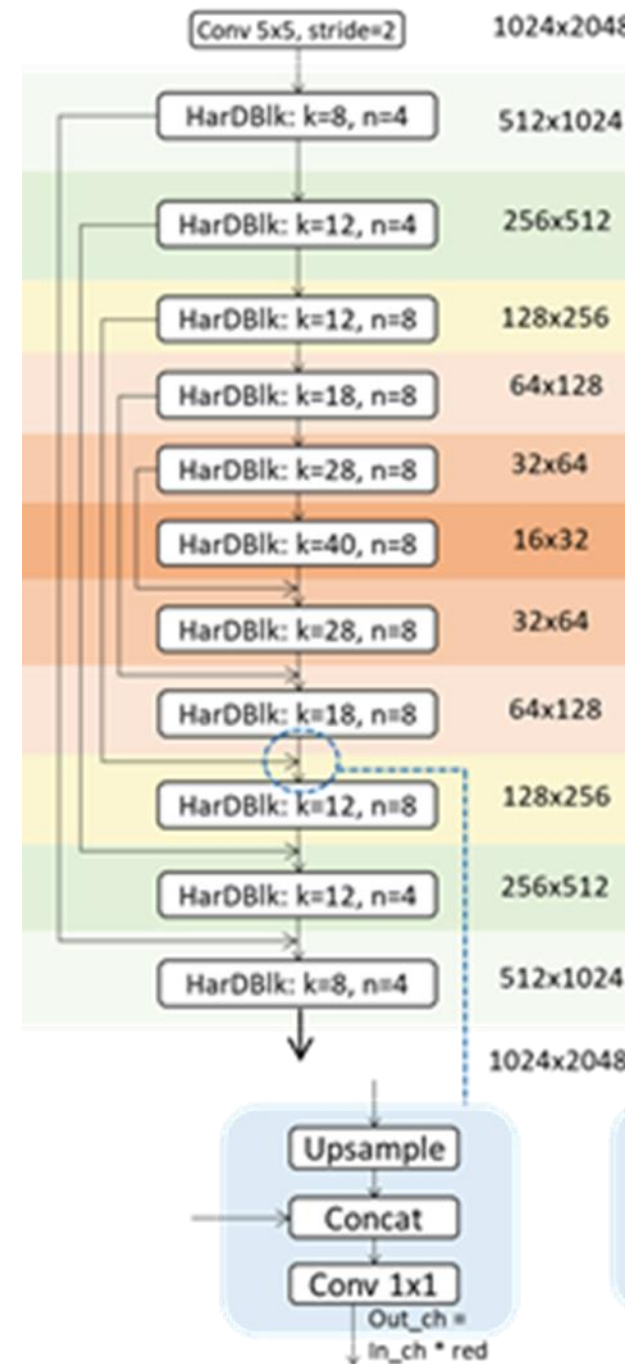
# U-Net & DenseNet-Inspired Network Design

# Proposed Network
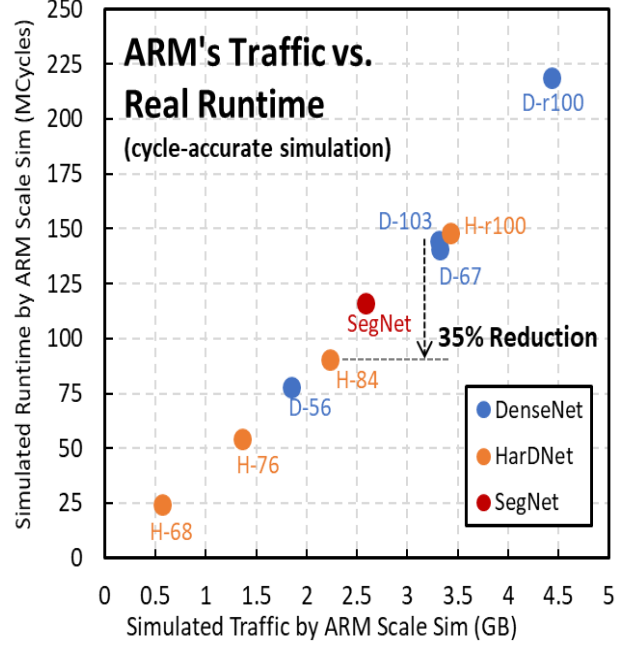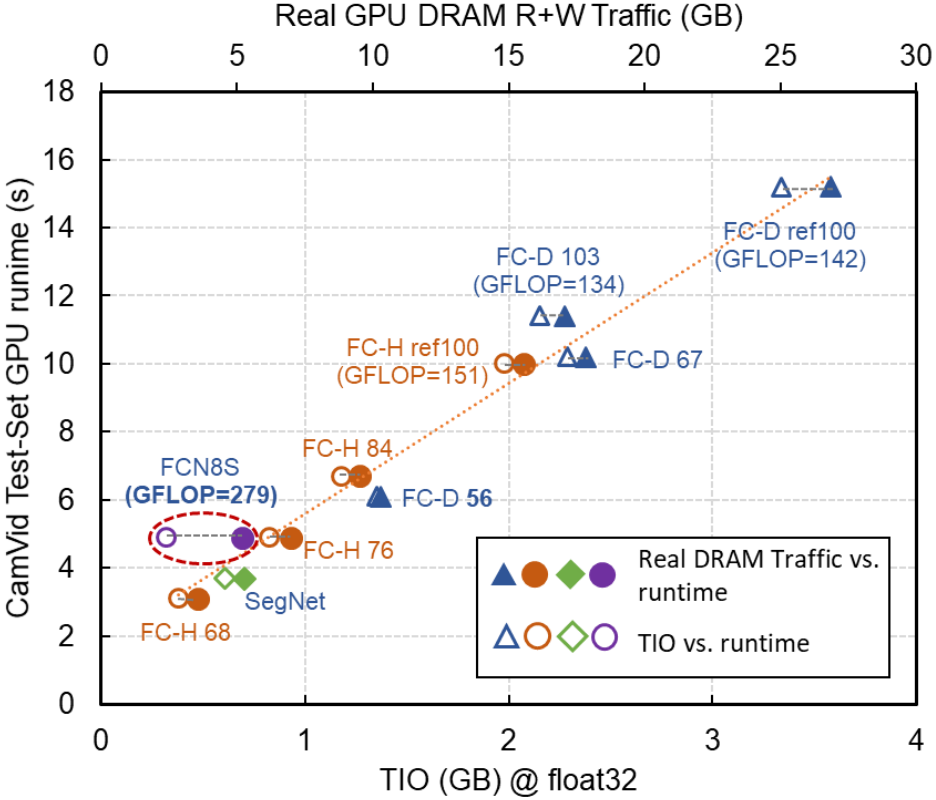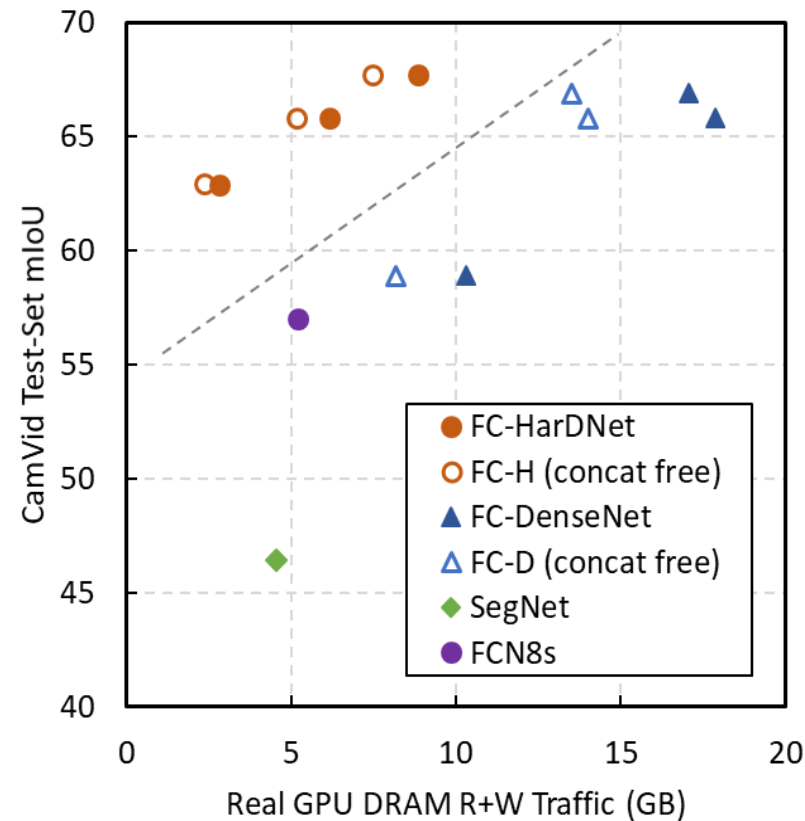
# DRAM Traffic vs Run-Time



Figure 3: Runtime vs. DRAM traffic measured by the simulation of ARM Scale.

# Low DRAM Traffic (Run-Time) and High Accuracy
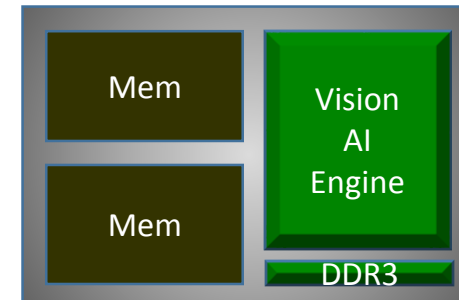
# To make an ASIC

USD8000/300W

USD100/10W

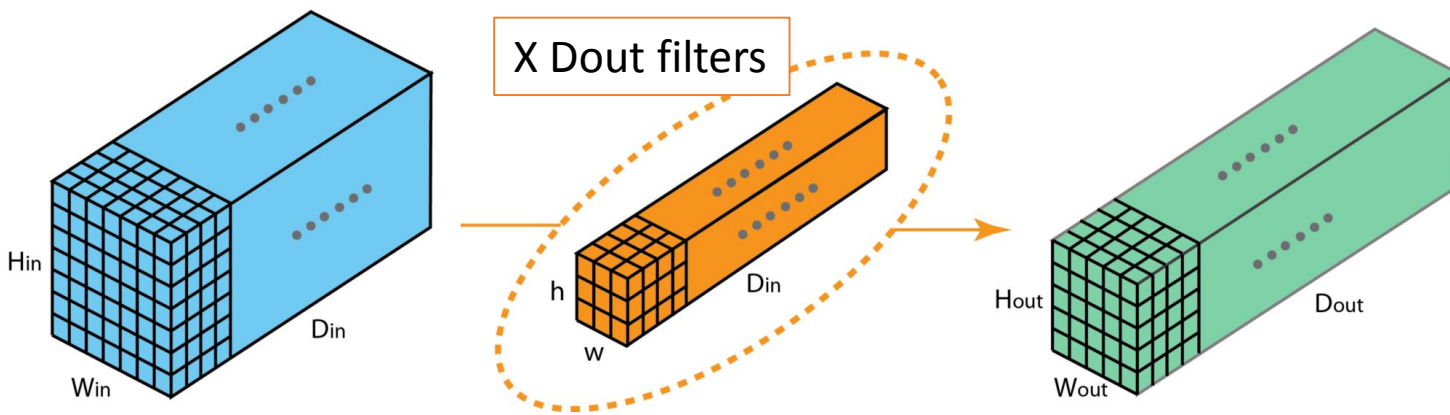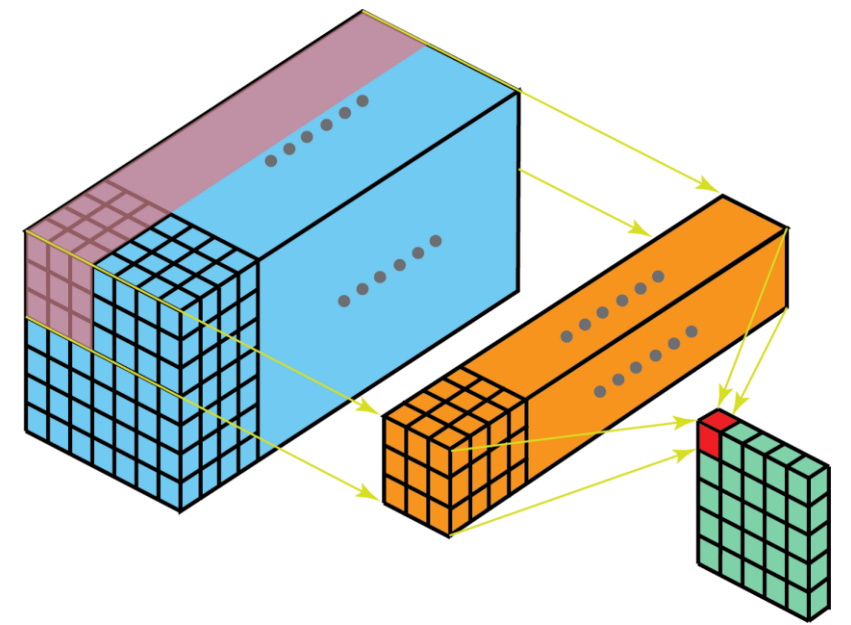ASIC

# > 99% Computes are 2D Convolution (Multi-Channel, Multi-Filter)



X Dout filters



A Comprehensive Introduction to Different Types of Convolutions in Deep Learning -- Towards intuitive understanding of convolutions through visualizations
Kunlun Bai

```
for (m = 0; m < numOutputLayers; m++)        //Loop 1
    for (n = 0; n < numInputLayers; n++)     //Loop 2
        for (h = 0; h < outputHeight; h++)   //Loop 3
            for (w = 0; w < outputWidth; w++)    //Loop 4
                for (i = 0; i < kernelHeight; i++)   //Loop 5
                    for (j = 0; j < kernelWidth; j++) //Loop 6
                        out[m][h][w] +=
                            in[n][h * S + i][w * S + j] *
                            kernel[m][n][i][j];
```
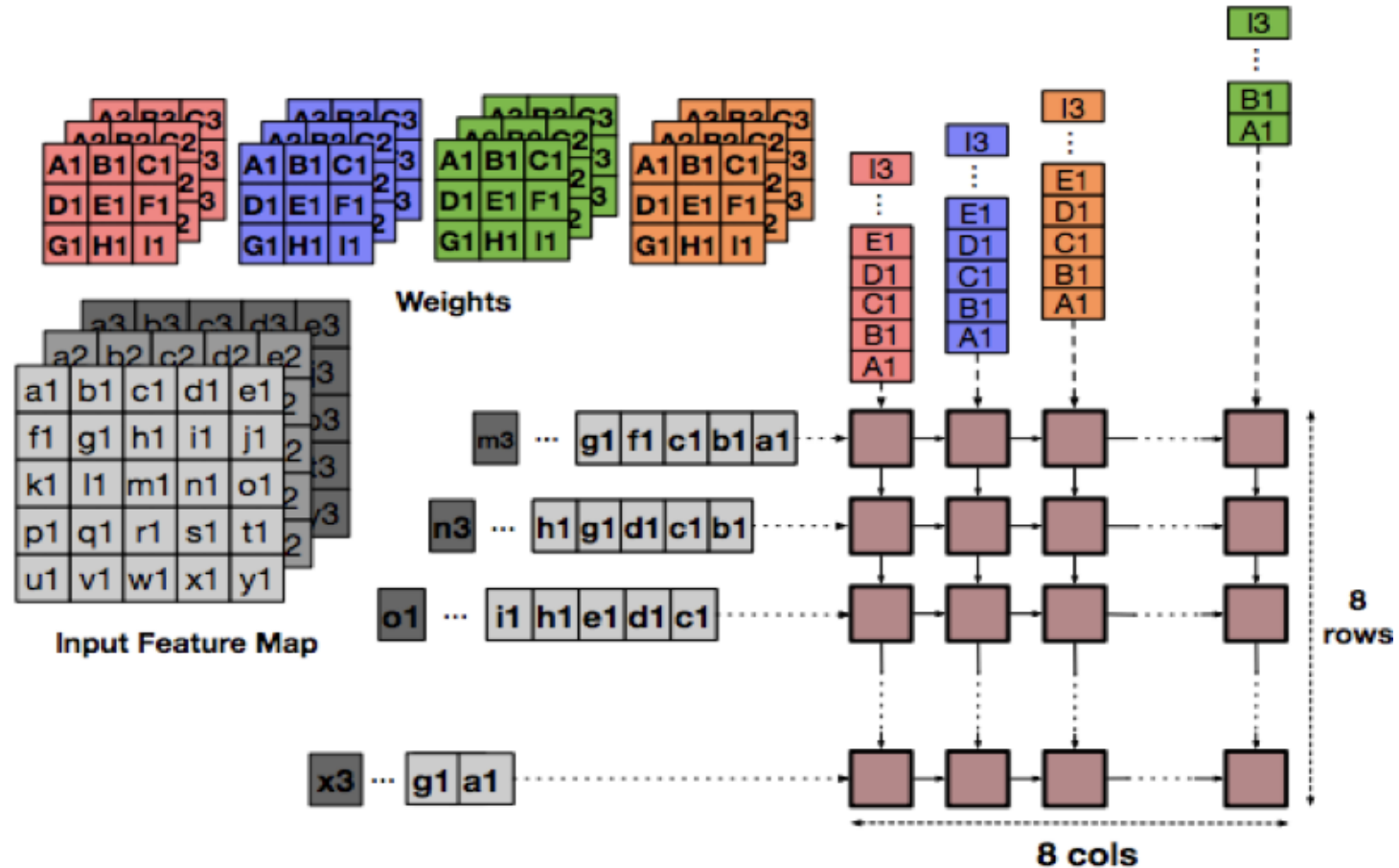
# 2D Convolution with Systolic Array



Citation: Ananda Samajdar, Yuhao Zhu, Paul Whatmough, Matthew Mattina, and Tushar Krishna. Scale-sim: Systolic cnn accelerator. arXiv preprint arXiv:1811.02883, 2018.
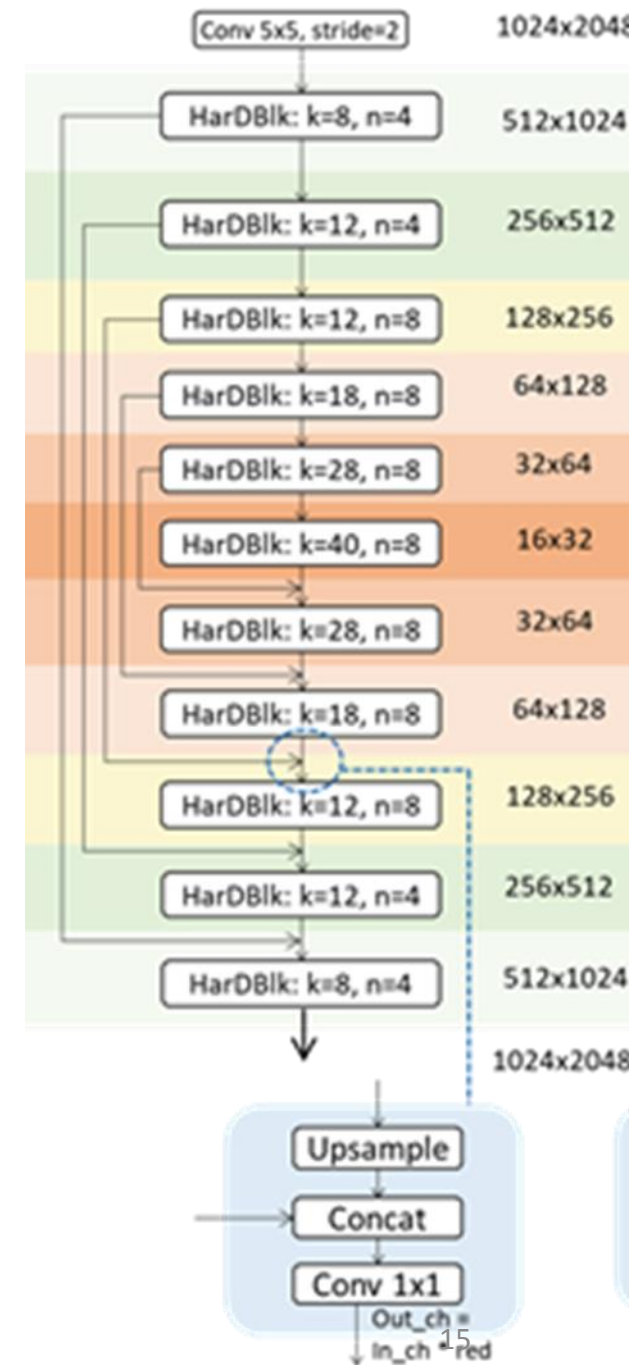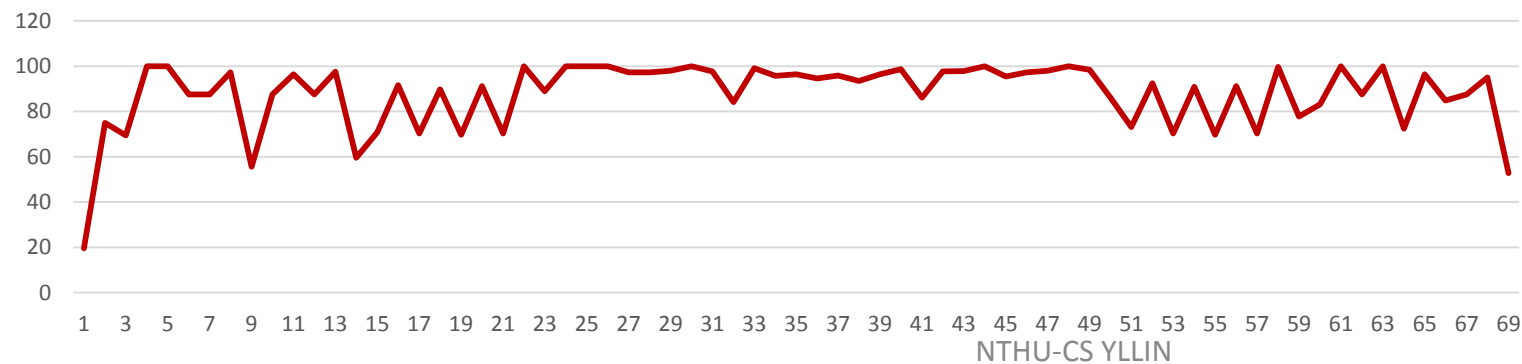
# Memory bound vs Compute Bound

### Per Layer # Paras vs #Ops



### Per Layer Utilization of 9216 MACs

# Results

- Network
  - 69-Layer Convolutional Neural Network (Other versions: 84, …)
  - 3.8M parameters
  - 59.685G Operations per 1Kx2K frame inference
- PyTorch on GPU Implementation
  - 80fps on a TWCC nVidia Tesla V100 32GB GPU (300Watt, USD8,000)
  - 59.658Gops * 80 / 120Tops        ~ 4% utilization
- Preliminary ASIC Design
  - 9216 MACs (Mutli-Add) = 18,432 PEs
  - Peak performance 9.216 Tera ops @ 500 MHz
  - 3.846M Clock Cycles to inference a 1Kx2K frame
  - 59.685G / (3.846M * 18432) ~ 85% utilization

# Thank you!!