# Coevolution of Neural Network and Computer Architecture
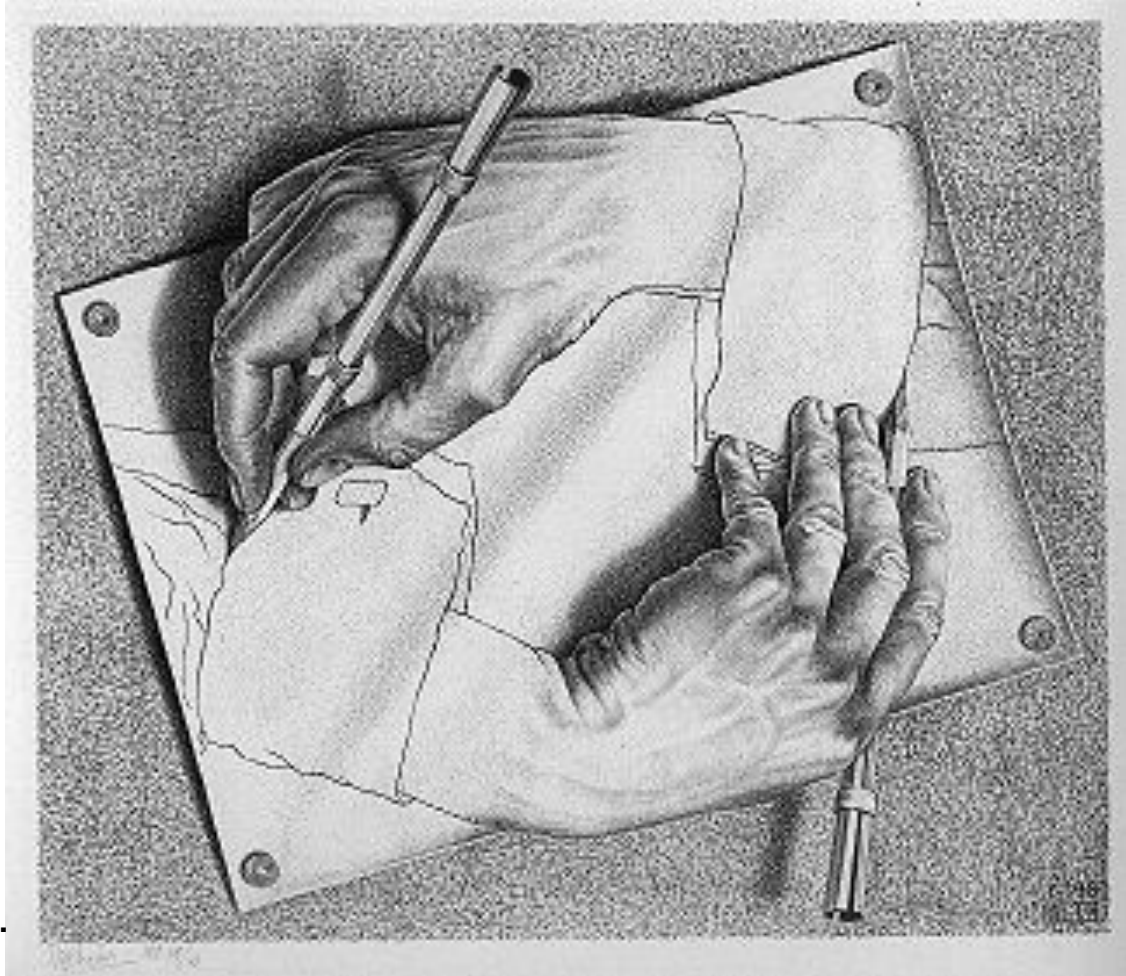
zsc@megvii.com

Aug. 2019

Face++ 旷视

Propose new kind of Neural Network for hardware.

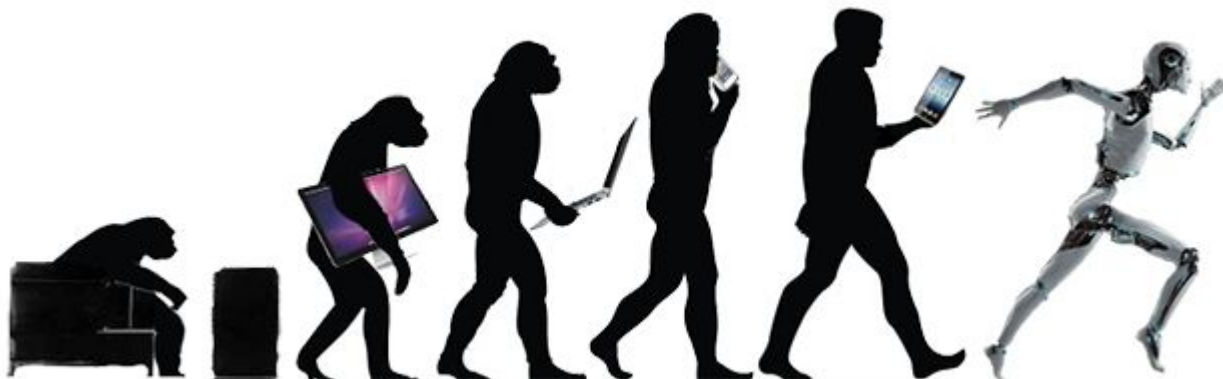Define new hardware for Neural Network.

Software-hardware co-evolution

# Deep Learning Challenge

- Make it start: Conceptual Breakthrough
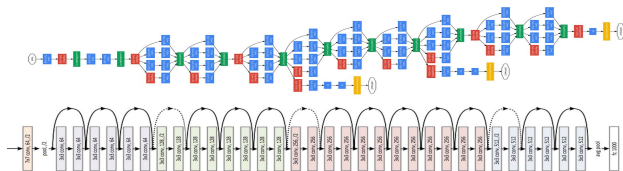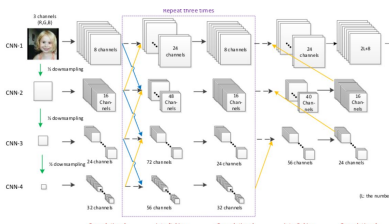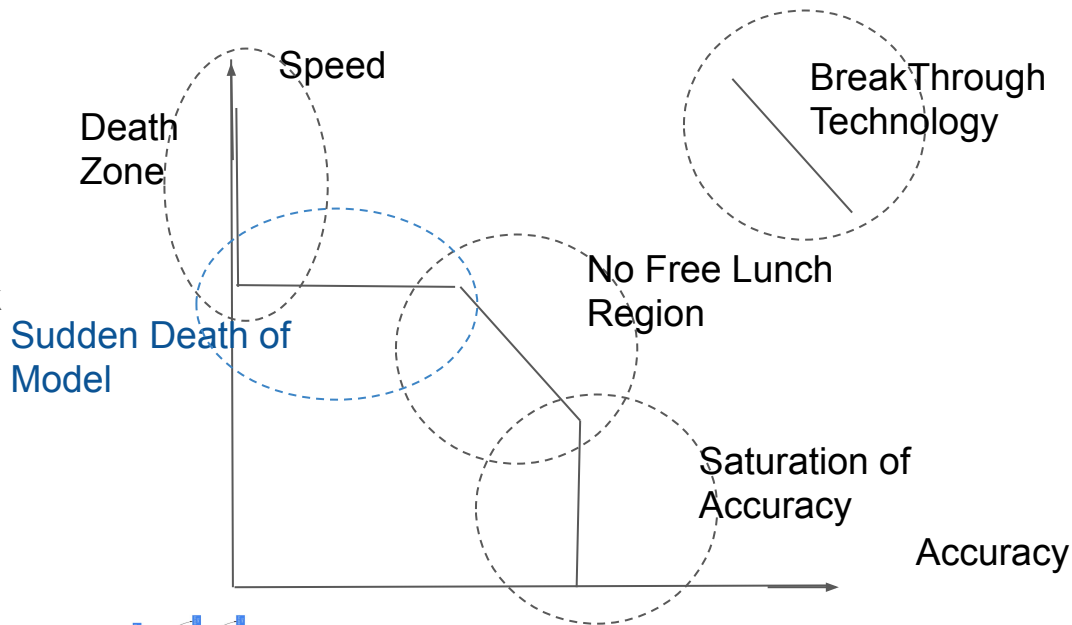- Make it work: Building product
- Make it cheap: Democratize



*https://medium.com/global-silicon-valley/the-evolution-of-mobile-computing-d273f23eda61*
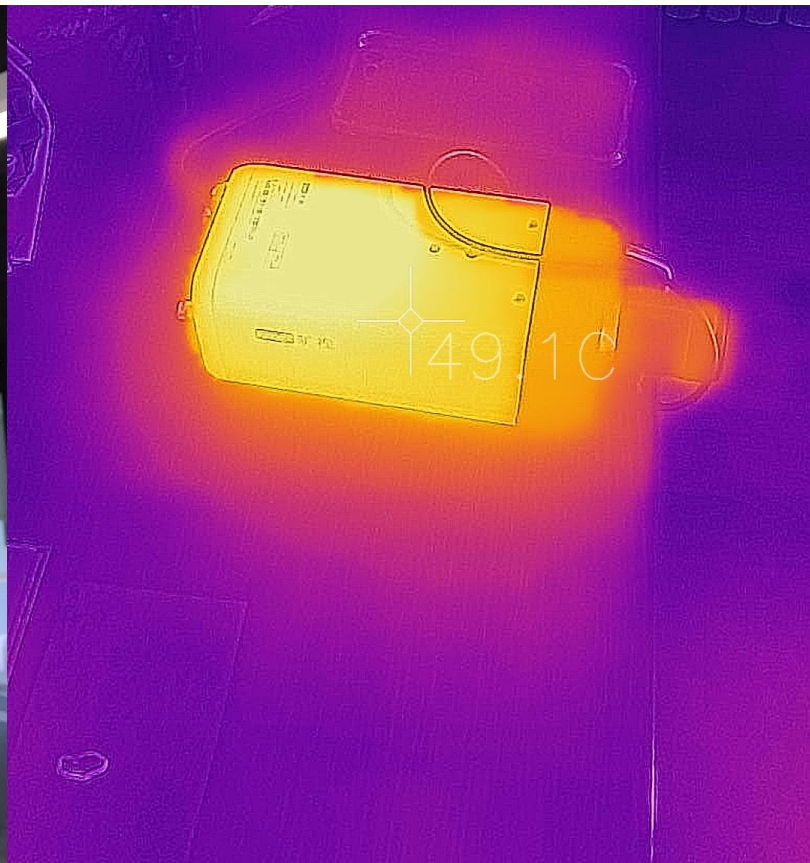
# Tradeoff between Accuracy and Speed

- Breakthroughs improve both accuracy and speed
  - Factorized Convolution (GoogleNet)
  - Skip connection (ResNet)
  - Fully Convolutional Network
  - Better Loss Function
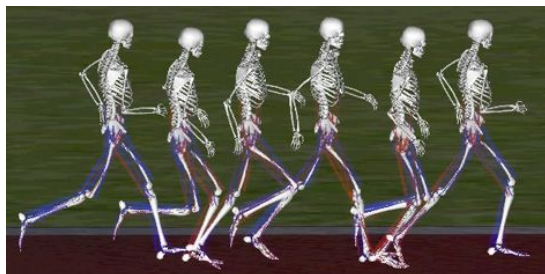  - Batch Normalization
  - Cyclic Learning Rate

# User cases: Deep Learning

# User cases: Reinforcement Learning

Characteristics: require fast & complex simulations

OpenSim
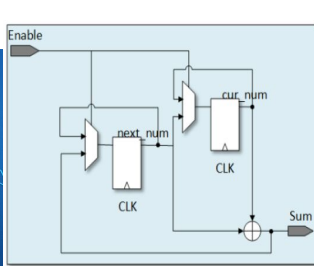


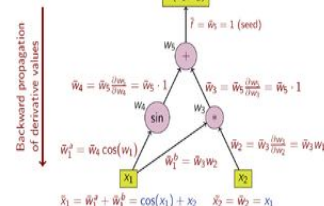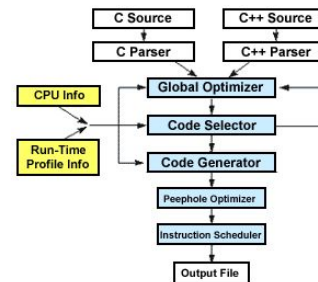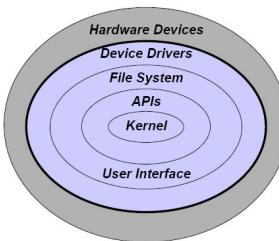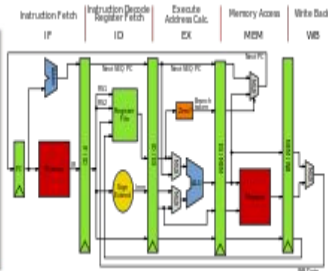A human skeleton model for locomotive task modeling.

GTA 5
AirSim



Simulation for self-driving car/ADAS and Drones.

# Computer Architecture answer to Deep Learning Challenge

- Make it start: Conceptual Breakthrough
  - GPU: flexible powerhouse
- Make it work: Building product
  - ISA & Programming models: Graph Compiler and Execution Engine
- Make it cheap: Democratize
  - ASIC, Edge Computing, Cloud computing: mass production of all-in-one chips

# Computation Stack



**Silicon**
- Partitioning & Planning
- Place & Route
- Timing Closure

**Verilog**
- Karnaugh map
- Finite State Machine

**Architecture**
- ISA
- Micro-code
- Resource allocation
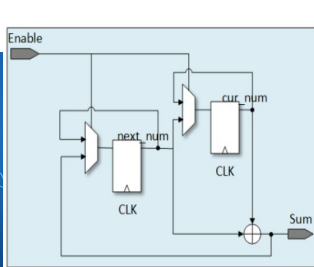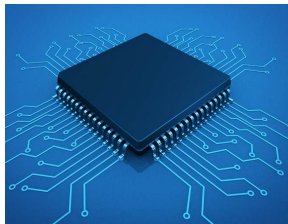
**Operating System**
- Page table
- File system
- Interrupts

**Compiler**
- Parallelism mining
- Memory latency hiding

**Computation Graph Engine**
- Kernels
- Execution Plan

# Computation Stack



**Silicon**
- Partitioning & Planning
- Place & Route
- Timing Closure

**Verilog**
- Karnaugh map
- Finite State Machine

**Architecture**
- ISA
- Micro-code
- Resource allocation

**Operating System**
- Page table
- File system
- Interrupts
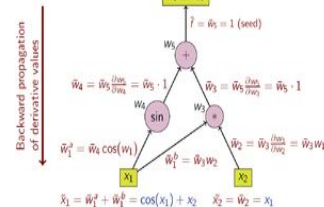
**Compiler**
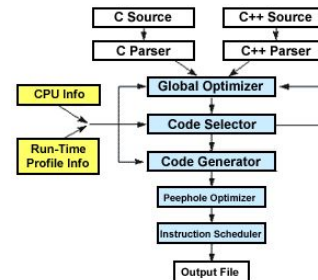- Parallelism mining
- Memory latency hiding

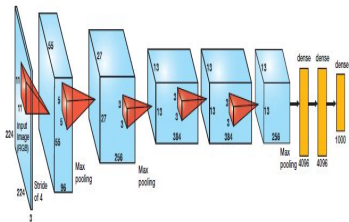**Computation Graph Engine**
- Kernels
- Execution Plan

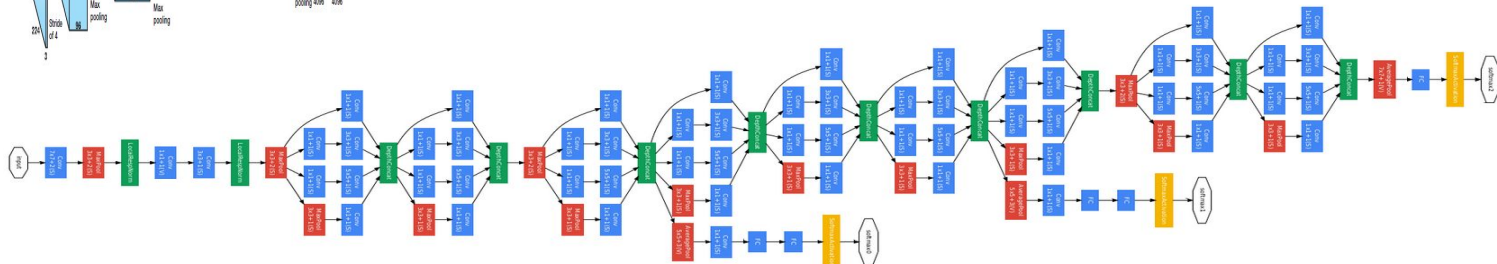*How will this stack deal with changes?*

# Case study: Large Neural Networks

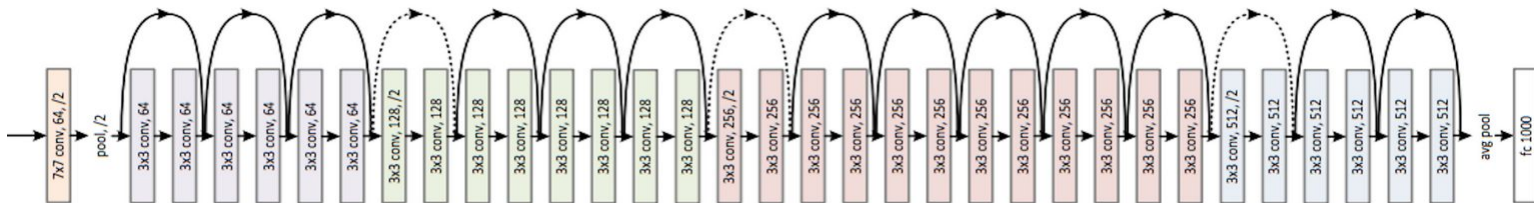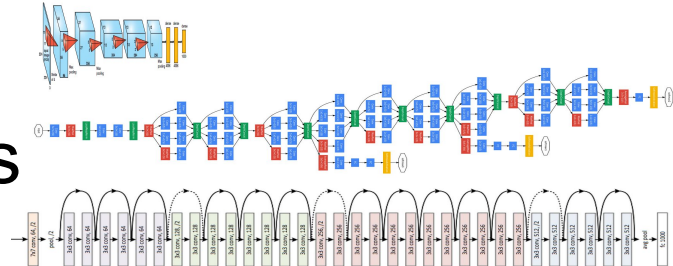Characteristics: many channels + side-branches + many layers

AlexNet



GoogLeNet


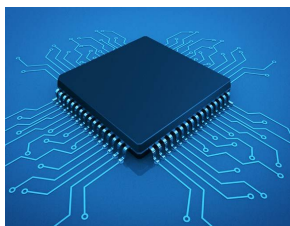
ResNet

# Case study: Large Neural Networks



On-Chip-Memory for caching feature maps

- Instructions for convolutions & non-linearity
- Systolic Array
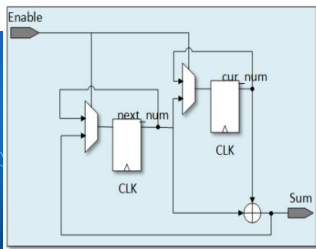
Large page-table

Auto-SIMD

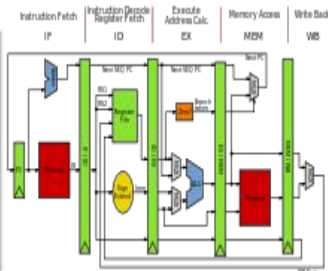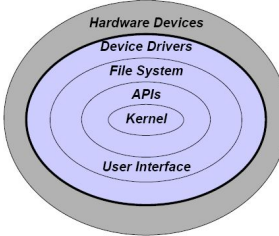Static analysis + dynamic profiling for kernel selection + execution plan



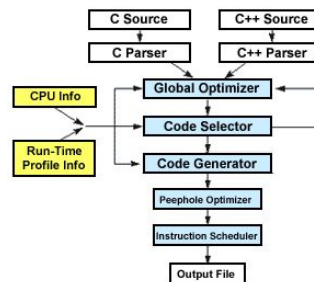Circuit to Generate Fibonacci Series *(Fig. 13)*
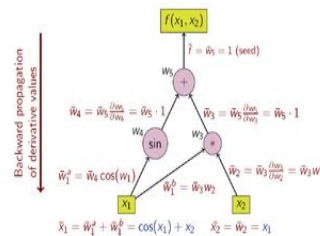
Silicon

Verilog

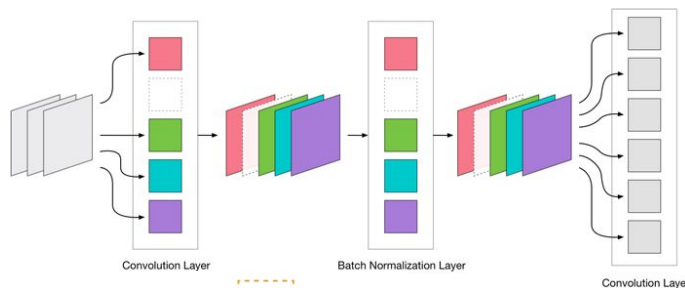Architecture

Operating System

Compiler

Computation Graph Engine

# Case study: Small Neural Networks

Characteristics: few channels + 1x1 convolutions

MobileNet



ShuffleNet



Lack of shortcut hurts its transfer learning ability.

The shuffle operation is an efficient way of information mixing, but its uniqueness slows its adoption.

# Case study: Small Neural Networks



On-Chip-Memory may be more important.

- Specialized support for few channel layers and 1x1 convolutions.
- Different batching

Lower overhead

Auto-SIMD

Fusion of layers + handcrafted kernels



Silicon

Verilog

Architecture

Operating System

Compiler

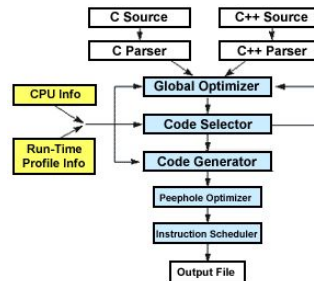Computation Graph Engine

# When a Neural Network Designers, a Computer Architect, a Compiler Expert and an OS Guru meet

- Designer wants
  - A reliable performance model
    - Open architecture design and assembly/microcode level exposure
  - Better profilers for runtime diagnostics and analyzers
  - Support for sparse matrices, dynamic operations
- Architect wants
  - Batch operations with constant delays
  - Regular memory access pattern subject to locality and many reuses
  - Streamlined memory/computation usage, no overwhelming peaks
  - Less number of operators
- Compiler Expert and OS Guru wants
  - To broker between the Designer and the Architect
    - Have a slow fallback for bizarre operators
    - Cutting peaks

# IC team @ Megvii


*GeneGAN*


*Learning2Paint*

## Neural Network Designer

*We train our DL models and design our networks!*

## Computer Achitects

*We build our processors and computers, from ISA to PCB!*
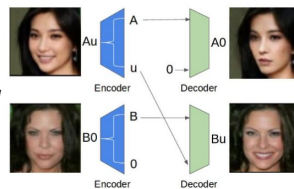
## AI-product Programmers

*We build our AI-products! From Javascript to Linux Kernels!*

*S-platform*
*Edge-computing platform*
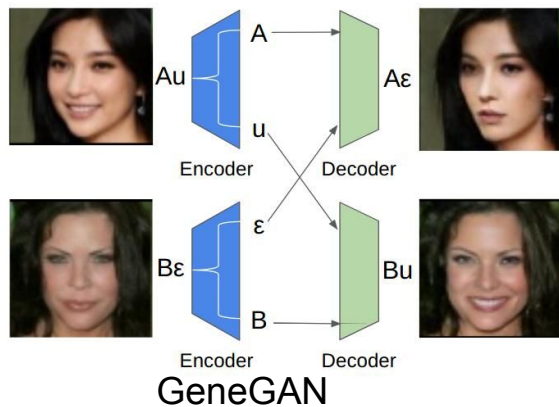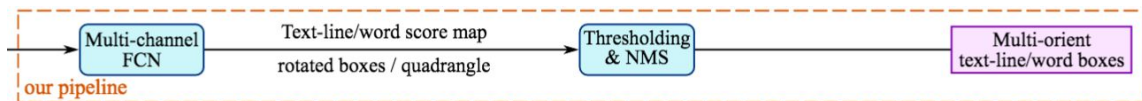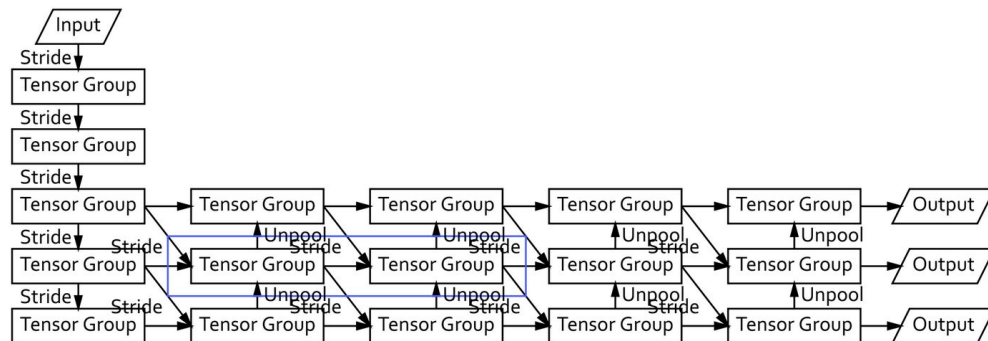
*A-firmware*
*Edge Device Firmware*

# Neural Network Designer


GeneGAN


EAST


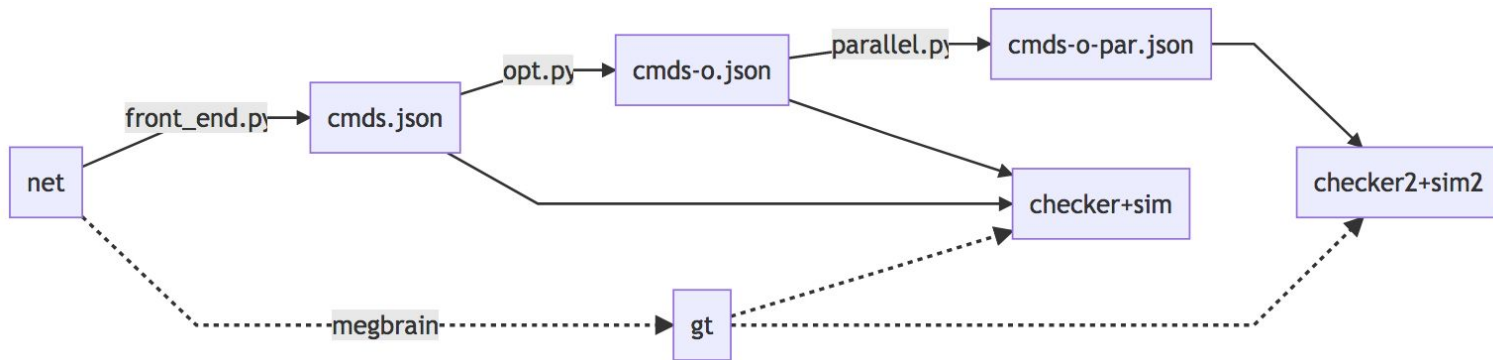Learning2Paint


Ъ-net

# Computer Architect

***"X Compiler"***: Optimizing & Autopar Compiler

# Static Scheduling
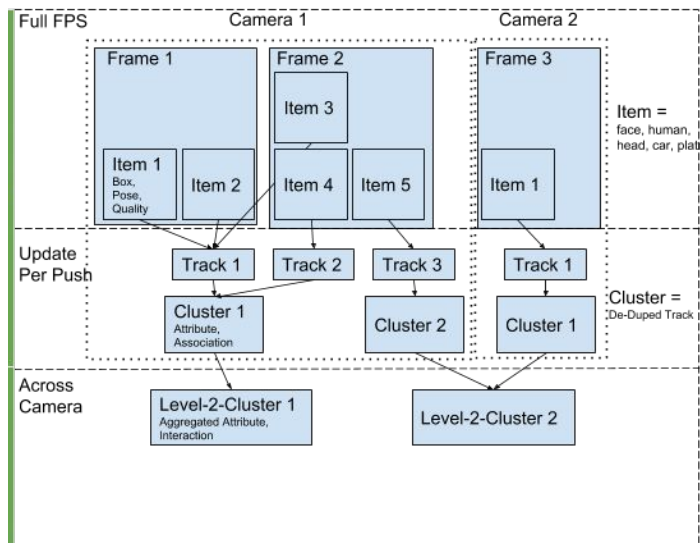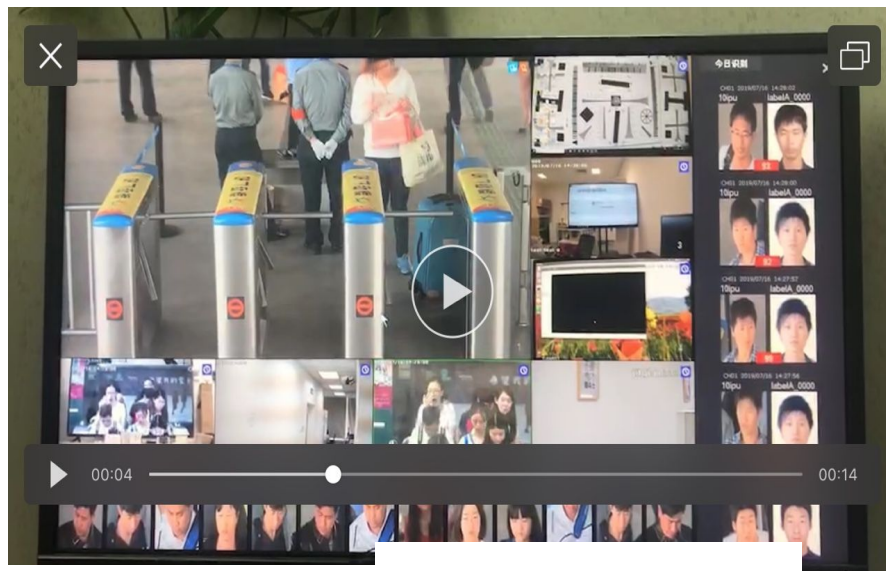
- Neural Networks are almost static
  - No branching
  - (almost always) Fixed length data: fixed input/output/intermediate size
  - Regular computation
- But there are "clouds"
  - DDR latency / bandwidth
  - Compression
  - Cache
  - Interrupts
- Dynamic Scheduling inevitable?

# AI-product Programmers

**S-platform**: Edge-computing platform



**A-firmware**: Edge Device Firmware

# About me

| | Natural Language Question & Answer<br>● Indoor Navigation with INS<br>● Group Orbit Optimization | |
|---|---|---|
| ● Source-to-source transformation<br>● Cache simulation | ● Natural Language Question & Answer<br>● Indoor Navigation with INS<br>● Group Orbit Optimization | ● OCR<br>● Quantized Neural Network<br>● Smart Camera<br>● Reinforcement Learning |
| Compiler Optimization | Machine Learning | Neural Network |

| 2007 | 2008 | 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 | 2018 | 2019 |
|------|------|------|------|------|------|------|------|------|------|------|------|------|