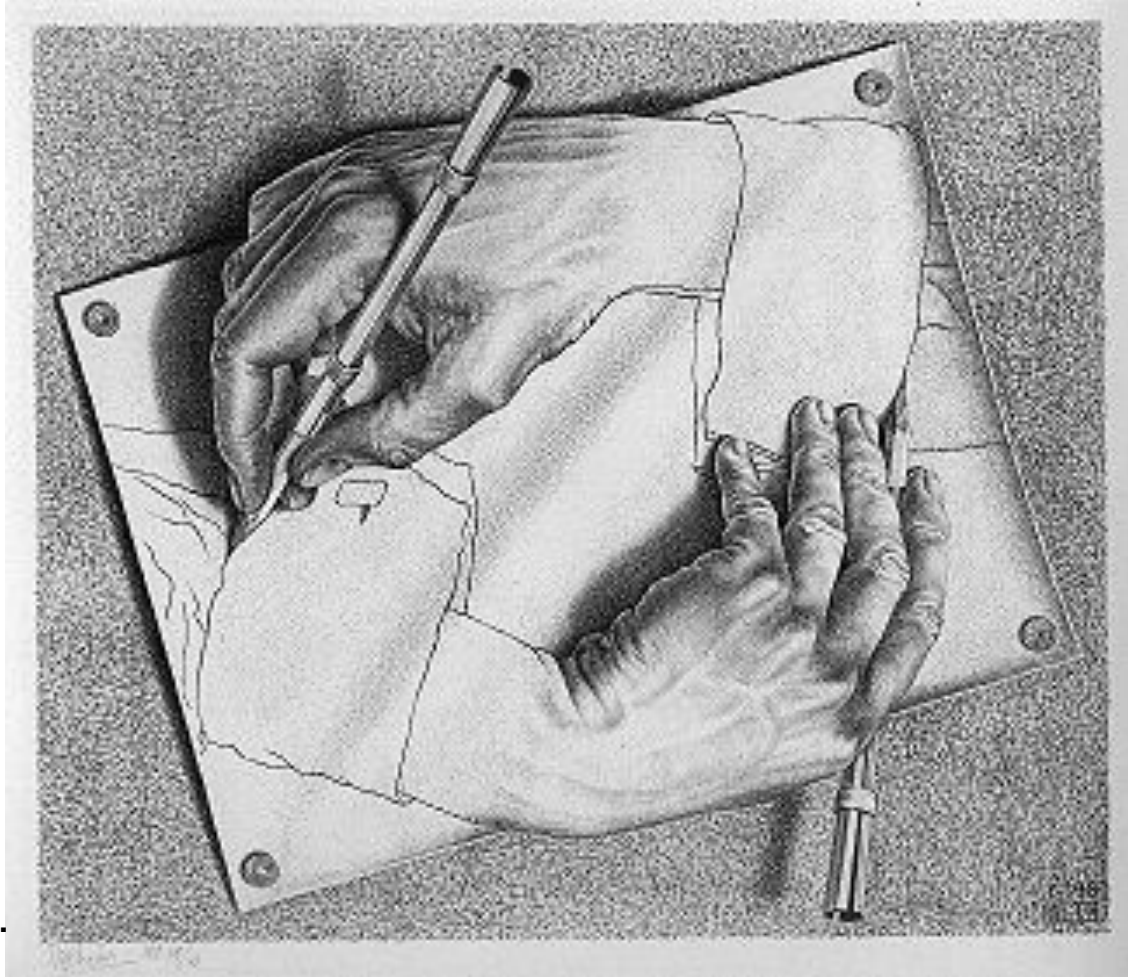


# Coevolution of Neural Network and Computer Architecture

[zsc@megvii.com](mailto:zsc@megvii.com)

Aug. 2019



Define new hardware for Neural Network.

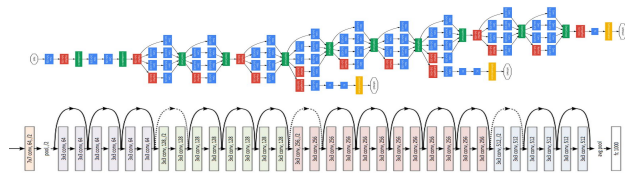
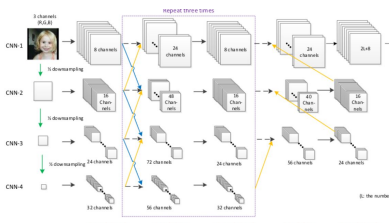
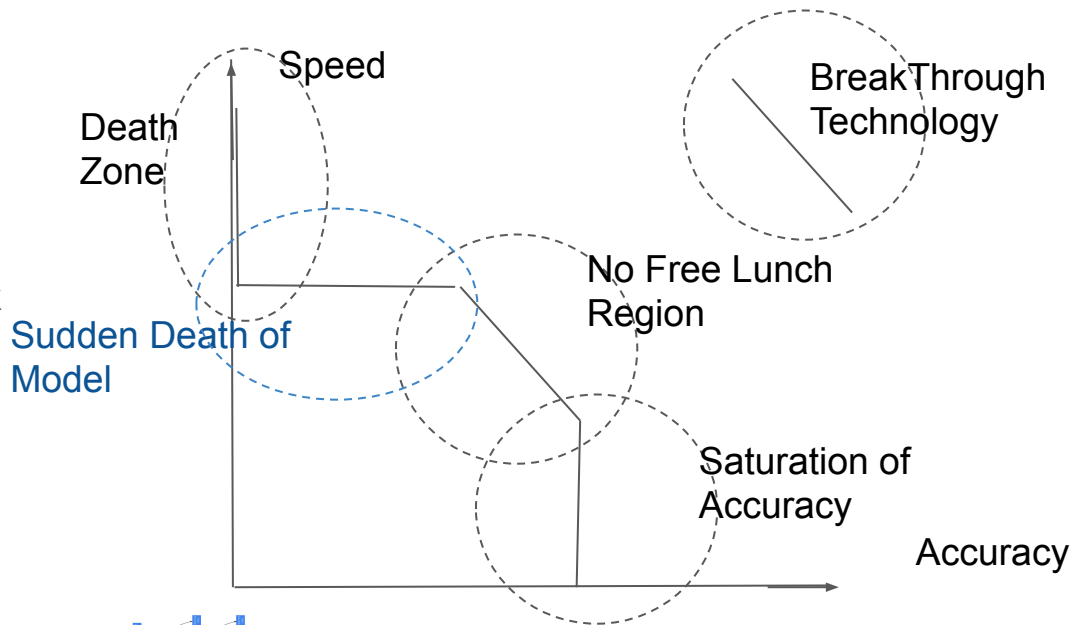
Propose new kind of Neural Network for hardware.

Software-hardware co-evolution

# Tradeoff between Accuracy and Speed

- Breakthroughs improve both accuracy and speed

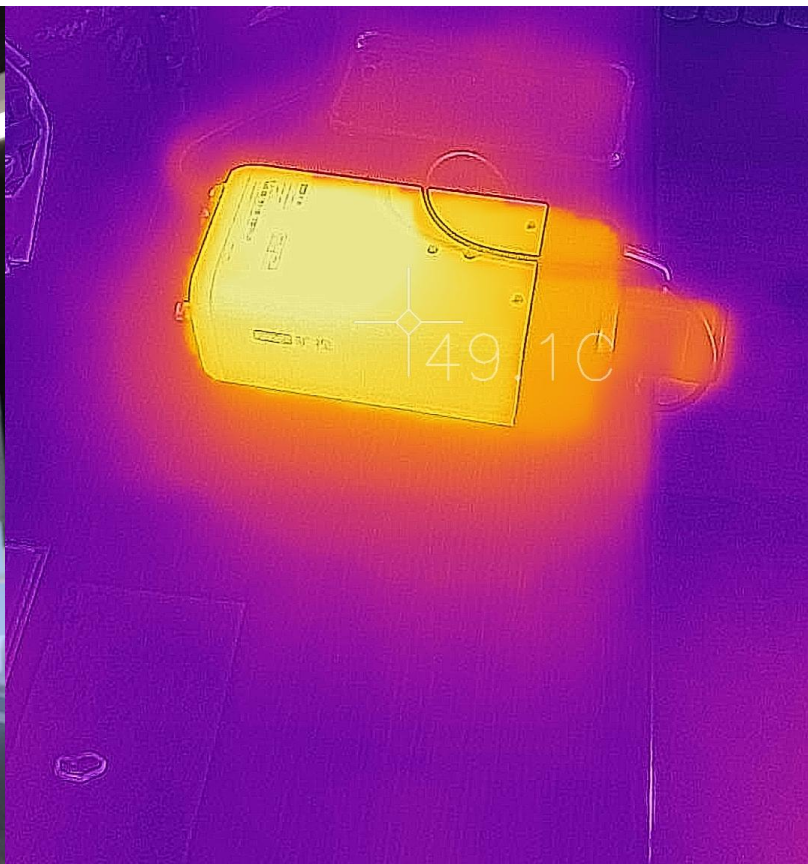
- Factorized Convolution (GoogLeNet)
- Skip connection (ResNet)
- Fully Convolutional Network
- Batch Normalization
- Cyclic Learning Rate
- NAS
- Transfer Learning in NLP





1, QA: 0.81 6, QA: 0.86 7, QA: 0.87 8, QA: 0.89 9, QA: 0.86 10, QA: 0.81 11, QA: 0.92 12, QA: 0.88 13, QA: 0.9 14, QA: 0.81 15, QA: 0.88 16, QA: 0.81 17, QA: 0.86 18, QA: 0.89 19, QA: 0.81 20, QA: 0.92 21, QA: 0.82 22, QA: 0.83 23, QA: 0.87 24, QA: 0.91 25, QA: 0.81 26, QA: 0.89 27, QA: 0.87 28, QA: 0.87 29, QA: 0.87 30, QA: 0.9 31, QA: 0.88 32, QA: 0.88 33, QA: 0.83 34, QA: 0.89 35, QA: 0.81 36, QA: 0.85 37, QA: 0.83 38, QA: 0.89 39, QA: 0.89 40, QA: 0.84 41, QA: 0.9 42, QA: 0.84 43, QA: 0.81 44, QA: 0.88 45, QA: 0.9 46, QA: 0.82 47, QA: 0.83 48, QA: 0.83 49, QA: 0.88 50, QA: 0.84 51, QA: 0.81 52, QA: 0.85 53, QA: 0.83 54, QA: 0.85 55, QA: 0.8 56, QA: 0.85 57, QA: 0.92 58, QA: 0.85 59, QA: 0.81 60, QA: 0.87 61, QA: 0.82 62, QA: 0.82 63, QA: 0.81 64, QA: 0.86 65, QA: 0.89 66, QA: 0.87 67, QA: 0.87 68, QA: 0.85 69, QA: 0.85 70, QA: 0.89 71, QA: 0.91 72, QA: 0.84 73, QA: 0.85 74, QA: 0.85 75, QA: 0.85 76, QA: 0.85 77, QA: 0.85 78, QA: 0.85 79, QA: 0.85 80, QA: 0.81 81, QA: 0.85 82, QA: 0.85 83, QA: 0.84 84, QA: 0.86 85, QA: 0.86 86, QA: 0.87 87, QA: 0.88 88, QA: 0.87 89, QA: 0.87 90, QA: 0.87 91, QA: 0.88 92, QA: 0.87 93, QA: 0.87 94, QA: 0.83 95, QA: 0.87 96, QA: 0.88 97, QA: 0.88 98, QA: 0.88 99, QA: 0.87 100, QA: 0.87

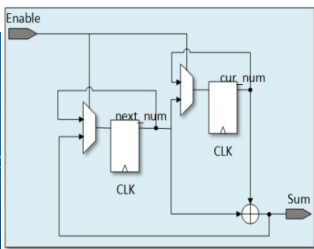
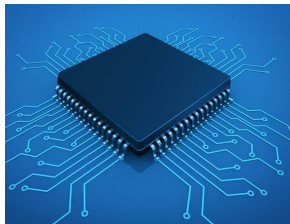
# User cases: Deep Learning



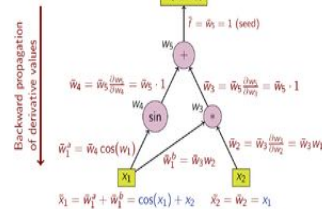
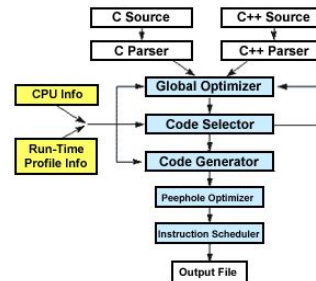
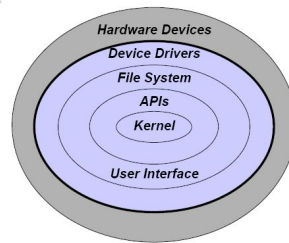
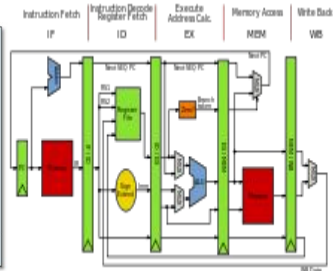
# Computer Architecture answer to Deep Learning Challenge

- 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



Circuit to Generate Fibonacci Series (Fig. 13)



## 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

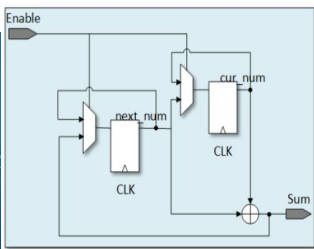
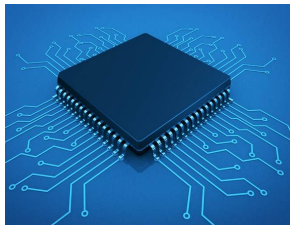
## Compiler

- Parallelism mining
- Memory latency hiding

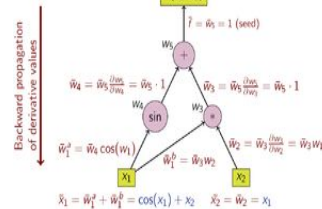
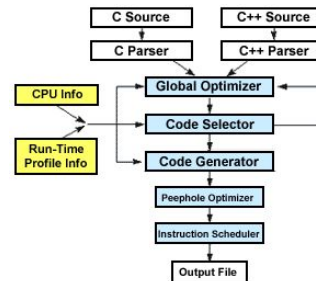
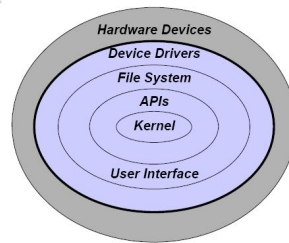
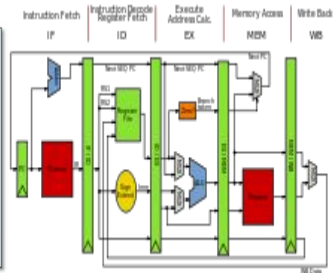
## Computation Graph Engine

- Kernels
- Execution Plan

# Computation Stack



Circuit to Generate Fibonacci Series (Fig. 13)



## 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

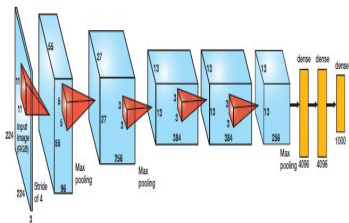
*How will this stack deal with changes?*



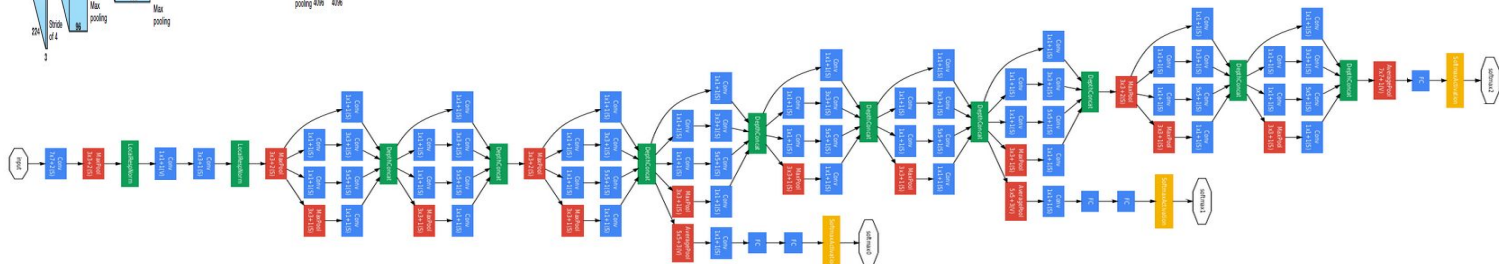
# Case study: Large Neural Networks

Characteristics: many channels + side-branches + many layers

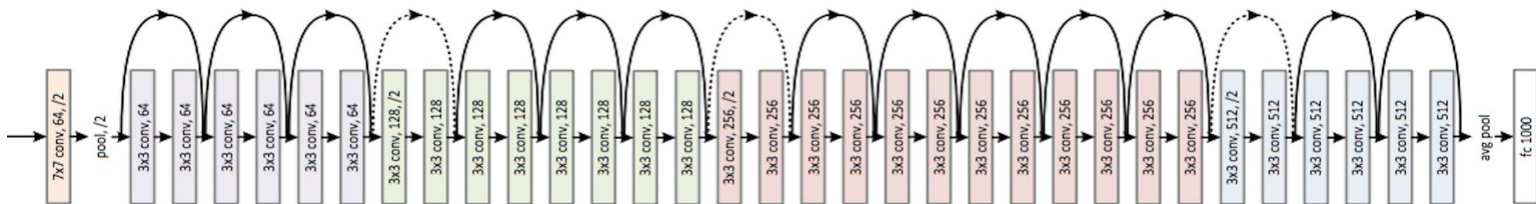
AlexNet



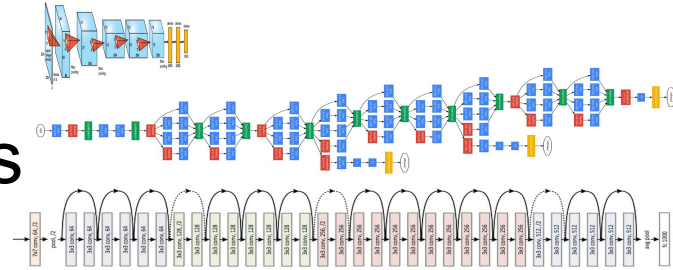
GoogLeNet



ResNet



# Case study: Large Neural Networks



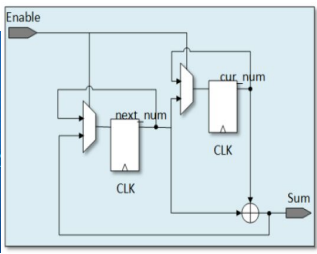
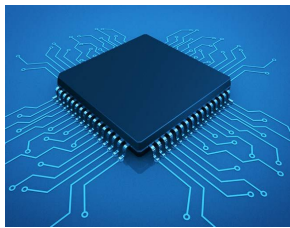
On-Chip-Memory for caching feature maps

- Instructions for convolutions & non-linearity
- Systolic Array

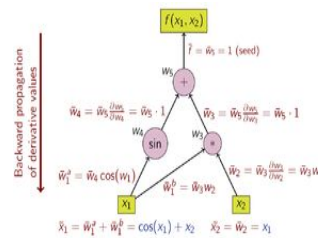
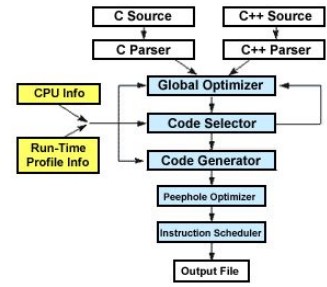
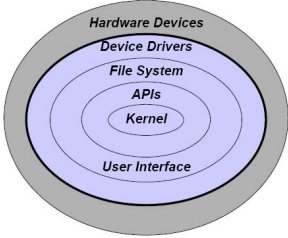
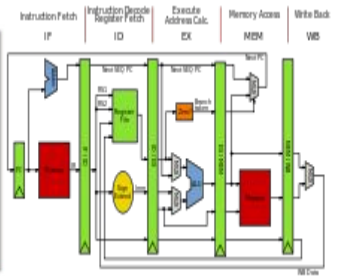
Large page-table

Auto-SIMD

Static analysis + dynamic profiling for kernel selection + execution plan



Circuit to Generate Fibonacci Series (Fig. 13)



Silicon

Verilog

Architecture

Operating System

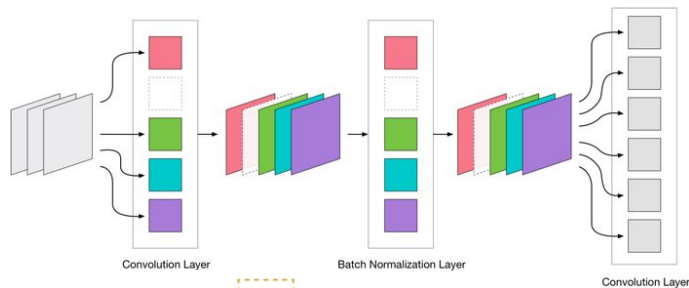
Compiler

Computation Graph Engine

# Case study: Small Neural Networks

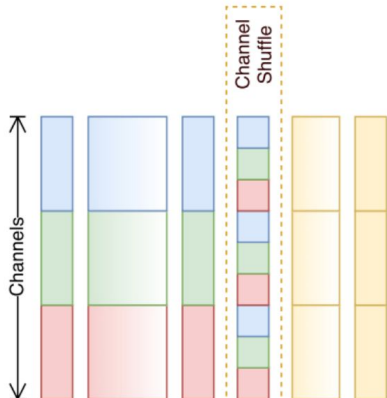
Characteristics: few channels + 1x1 convolutions

MobileNet



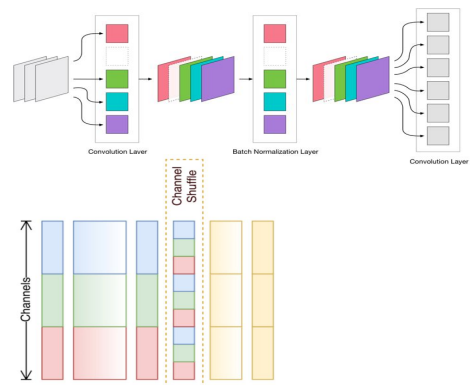
Lack of shortcut hurts its transfer learning ability.

ShuffleNet



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

# Case study: Small Neural Networks

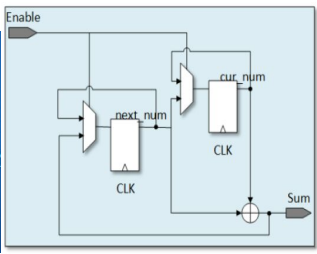
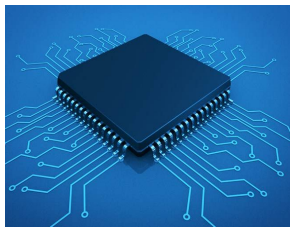


On-Chip-Memory may be more important.

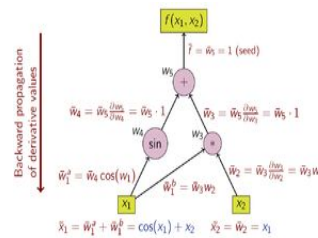
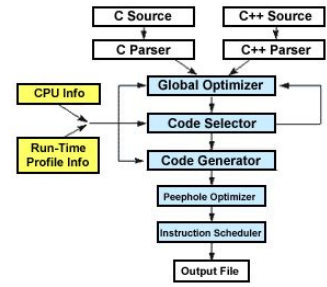
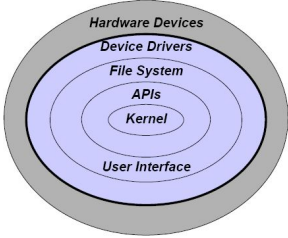
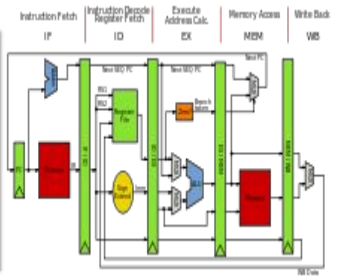
- Specialized support for few channel layers and 1x1 convolutions.
- Lower overhead
- Non-batch perf.
- Different batching
- Page coloring

Auto-SIMD

Fusion of layers + handcrafted kernels



Circuit to Generate Fibonacci Series (Fig. 13)



Silicon

Verilog

Architecture

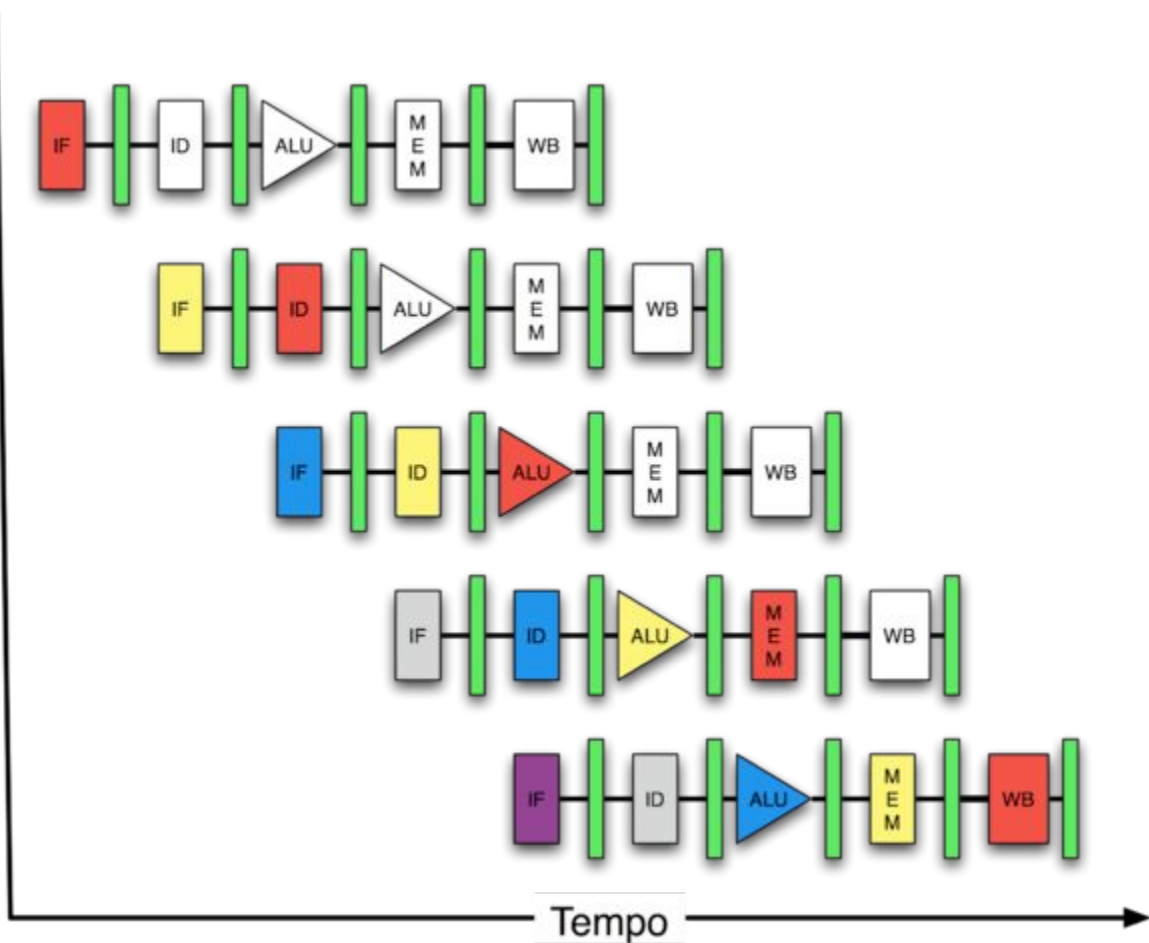
Operating System

Compiler

Computation Graph Engine

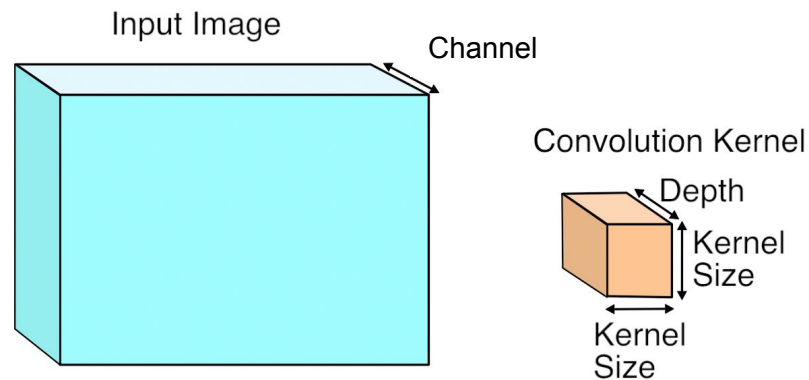
# Pipeline approach

- With DL, Need to deal with batch data to improve computation / memory ratio



# Instruction Size vs. Feature Size

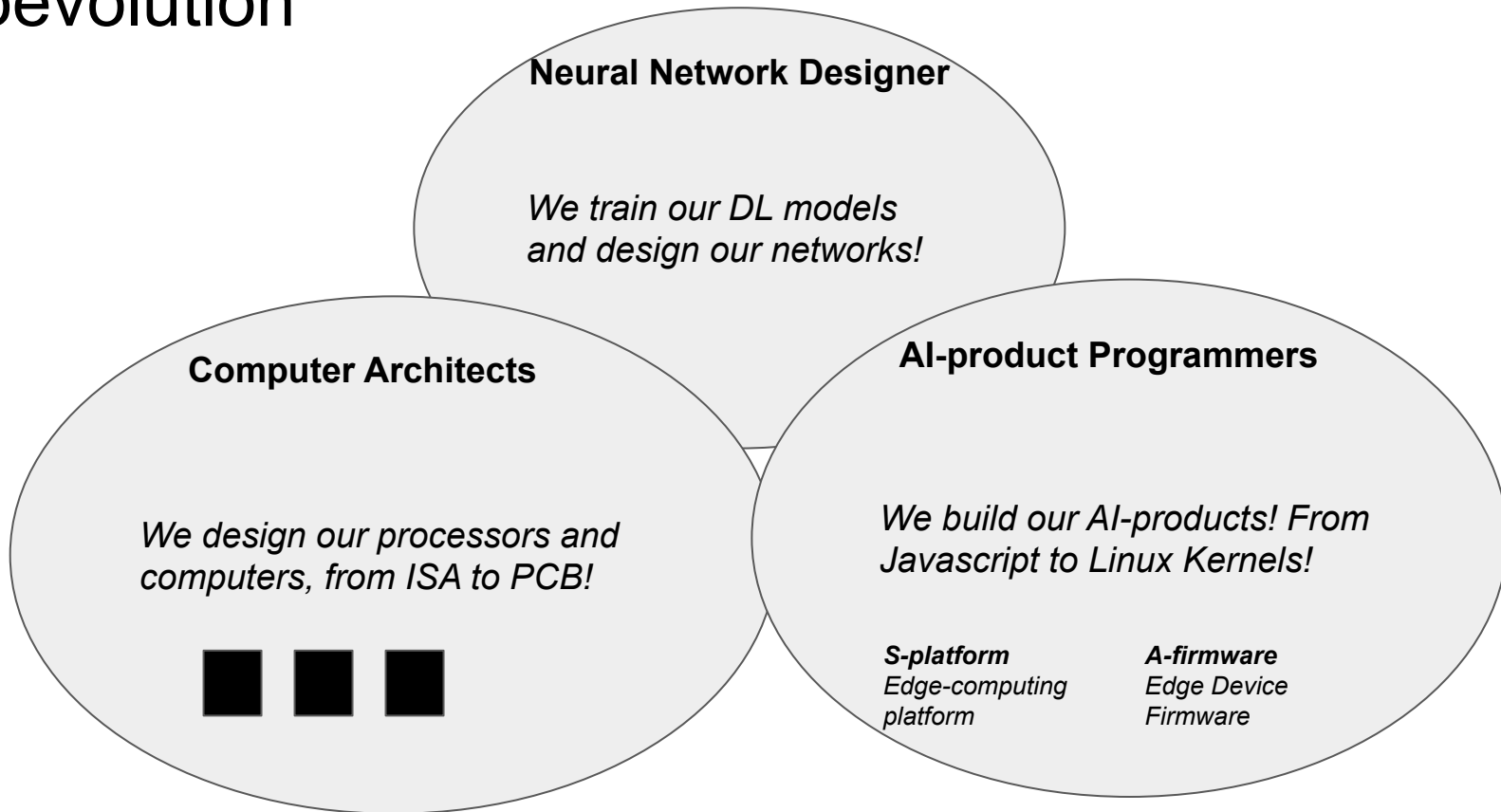
- Composition of Instruction
  - Type: conv, concat, unpool, stride
  - Weights
- Feature / Weight =  $N * H * W / (K * Kh * Kw)$ 
  - Feature =  $N * C * H * W$
  - Weight =  $K * C * Kh * Kw$
  - Getting smaller when later in Detection/Recognition NN
  - Relatively stable in image generation NN
- Efficiency suffers from low feature / weight ratio when small batch size



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

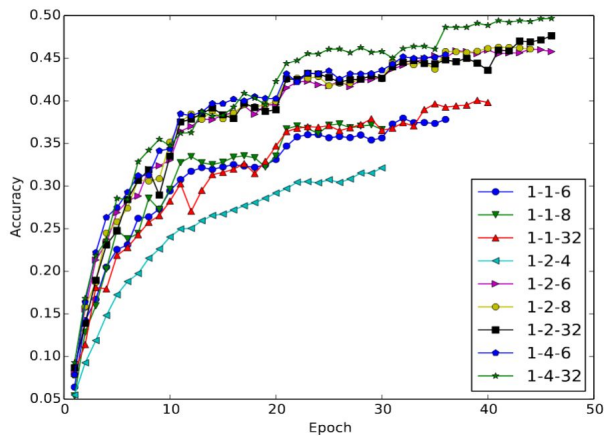
- Designer wants
  - A reliable performance model
    - Open architecture design and assembly/microcode level exposure
  - Better profilers for runtime diagnostics and analyzers
  - Support for sparse matrices, dynamic operations
- Architect wants
  - Batch operations with constant delays
  - Regular 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

# Coevolution

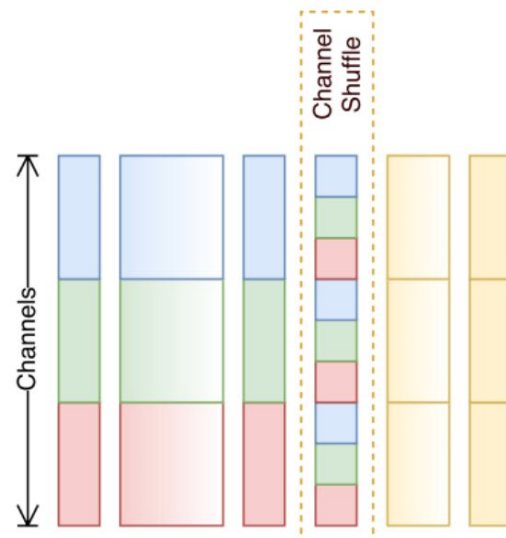




# Neural Network Designer



DoReFa-net



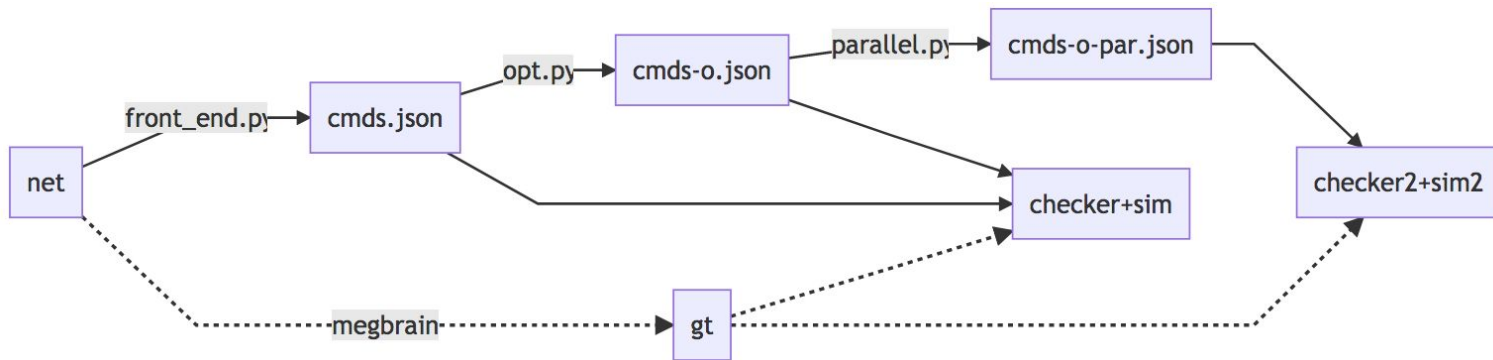
ShuffleNet

# Example: self-adjustable global channel quota

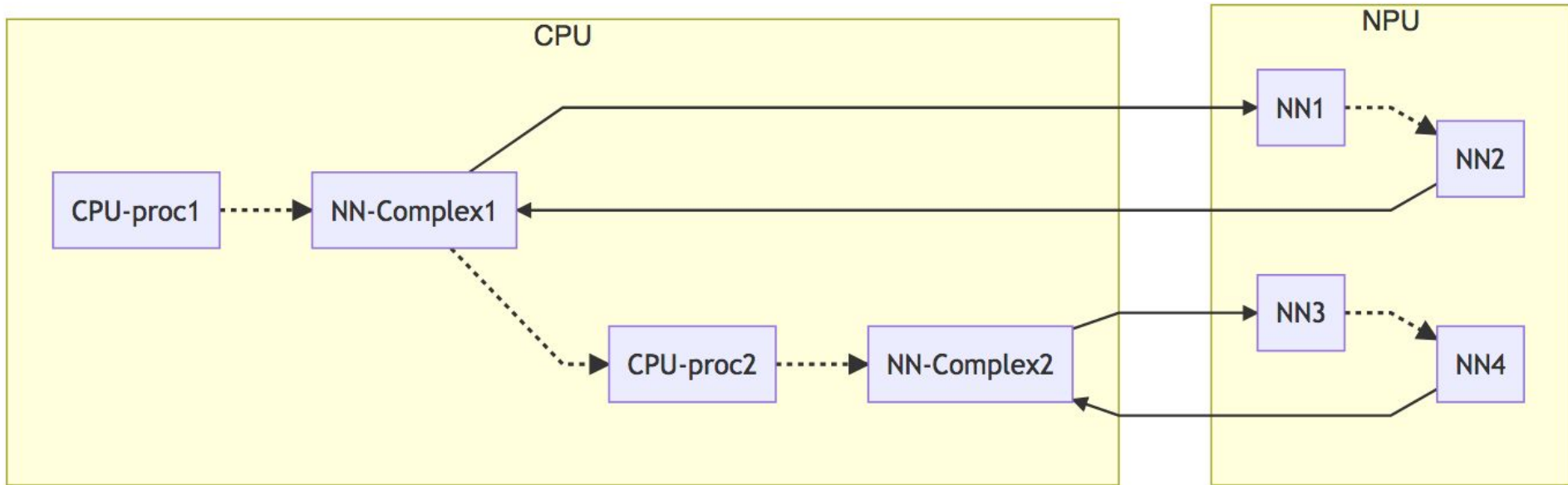
- k samples from n without replacement
  - follows multinomial hypergeometric distribution
  - satisfies  $EX_i = p_i$
- SVHN-cropped:
  - channel distribution
    - base: 3 96, 128, 128, 256 | 512, 10; train 19/sec, misclassify 0.027
    - self-adjusted: 3, 17, 26, 100 | 409, 10; train 62/sec, misclassify 0.031
      - 50% less #channel
    - self-adjusted: 4, 23, 51, 191 | 512, 10; train 41/sec, misclassify 0.028
      - 30% less #channel

# Computer Architect

"X Compiler": Optimizing & Autopar Compiler



# Pipeline Scheduling

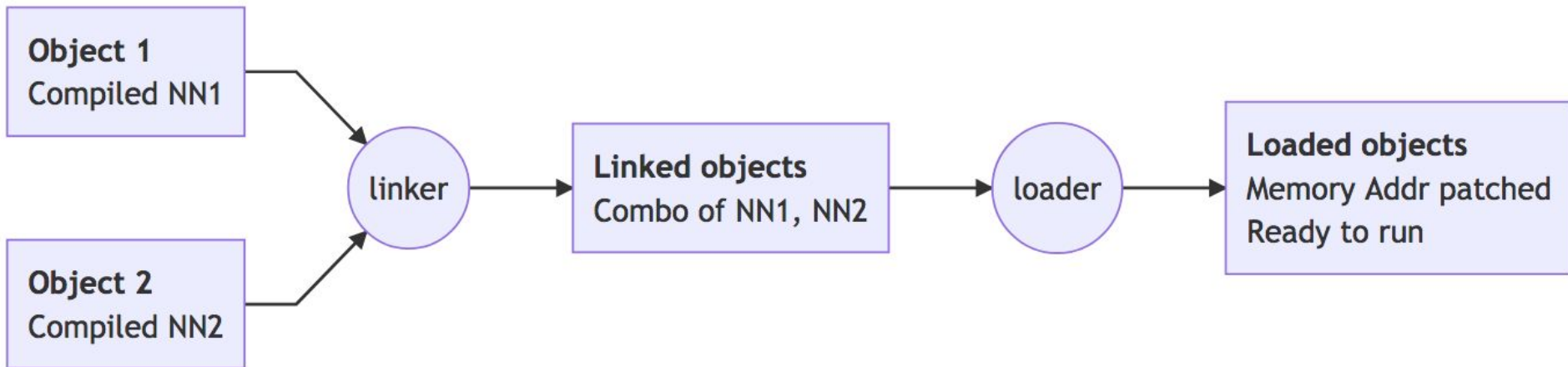


# Single Neural Network: Semi-Static Scheduling

- 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
  - Cache
- Dynamic Scheduling inevitable?

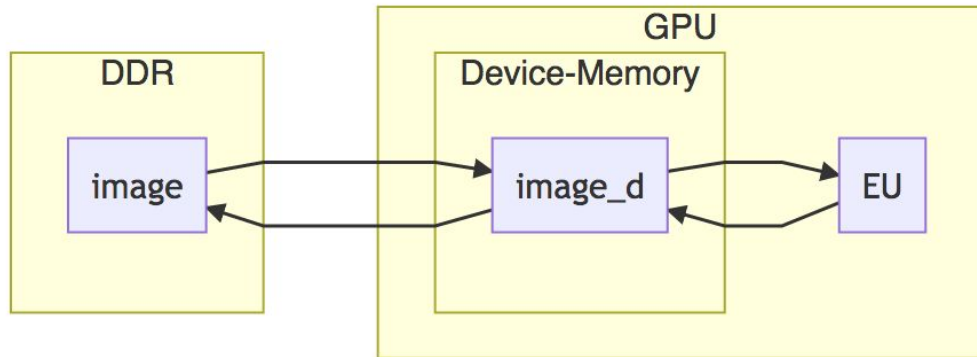
# Multiple Neural Network: Semi-Static Scheduling

- Combo-NN
  - Multiple NN's may be triggered by the same chunk of input data
  - Though logically separate, can be "linked" together
  - Ad-hoc on-the-fly combo by JIT

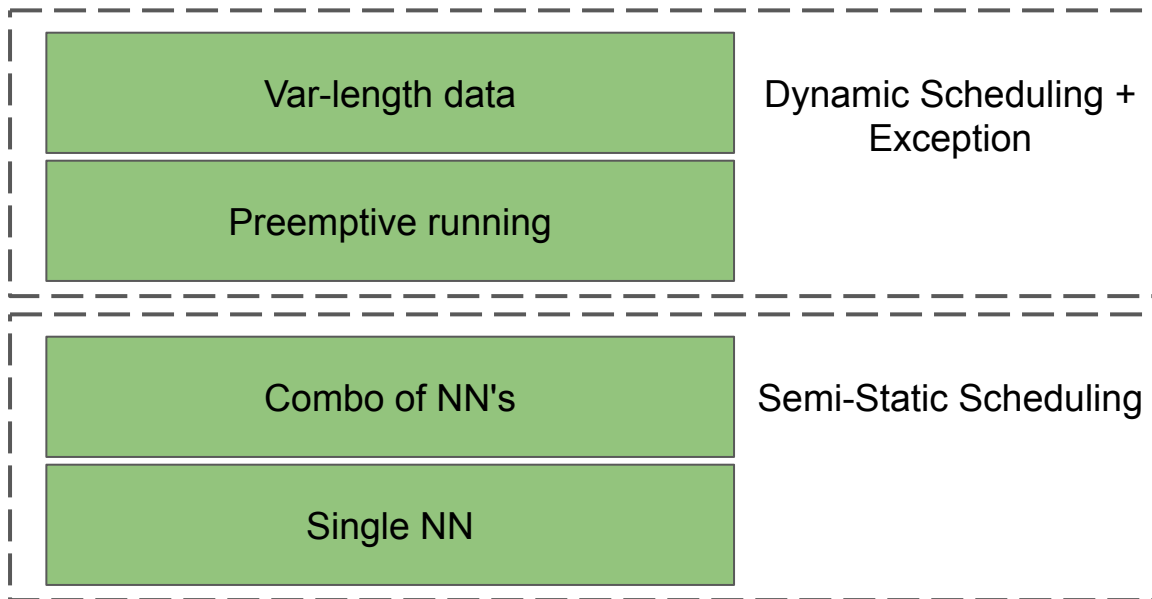


# Dyanmic Scheduling

- Complex memory state
  - Can Scheduling ensure OOM-free?
    - Interval analysis
  - Ensure proper recycling of resources when preemption
  - Exception mechanism
    - Spilling data to DDR when below watermark
    - May still not be safe
- Complex running time
  - Interrupts to CPU
  - Unbounded running time: Disk-level access



# Scheduling Overview

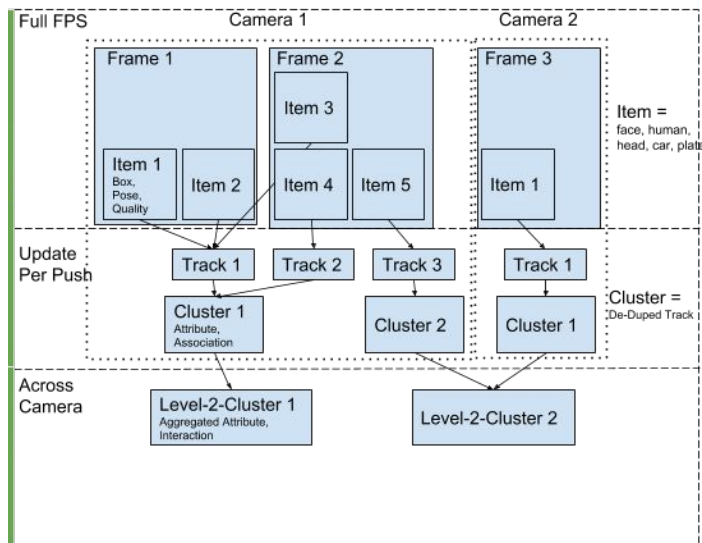




# AI-product Programmers

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

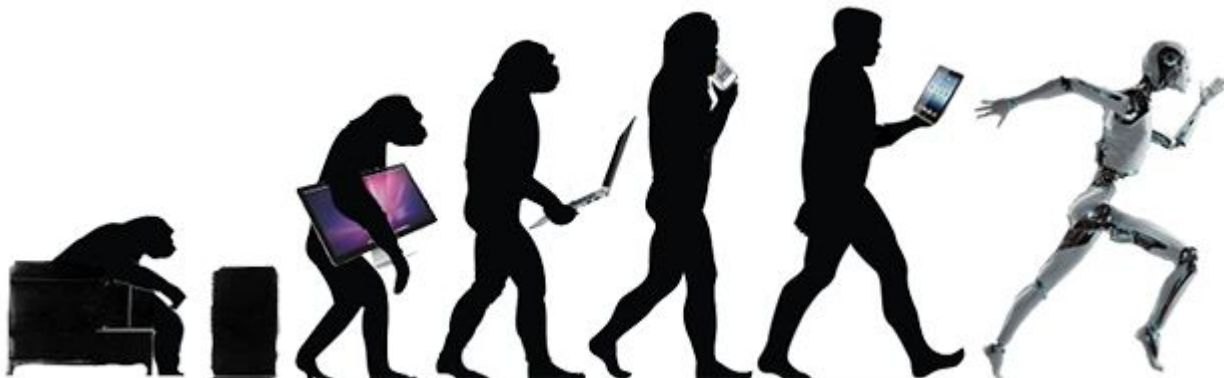
**A-firmware:** Edge Device Firmware



Backup after this slide

# Deep Learning Challenge

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

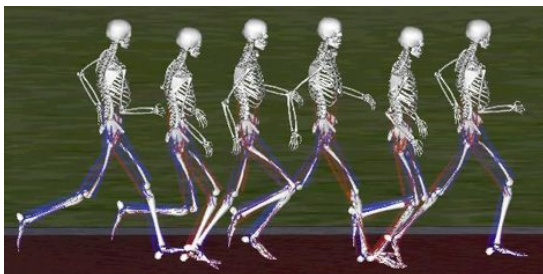


<https://medium.com/global-silicon-valley/the-evolution-of-mobile-computing-d273f23eda61>

# 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.