

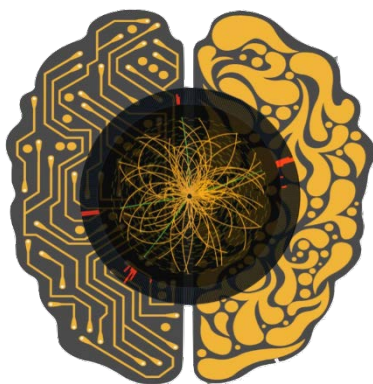


# Real-time AI Systems (Academia)

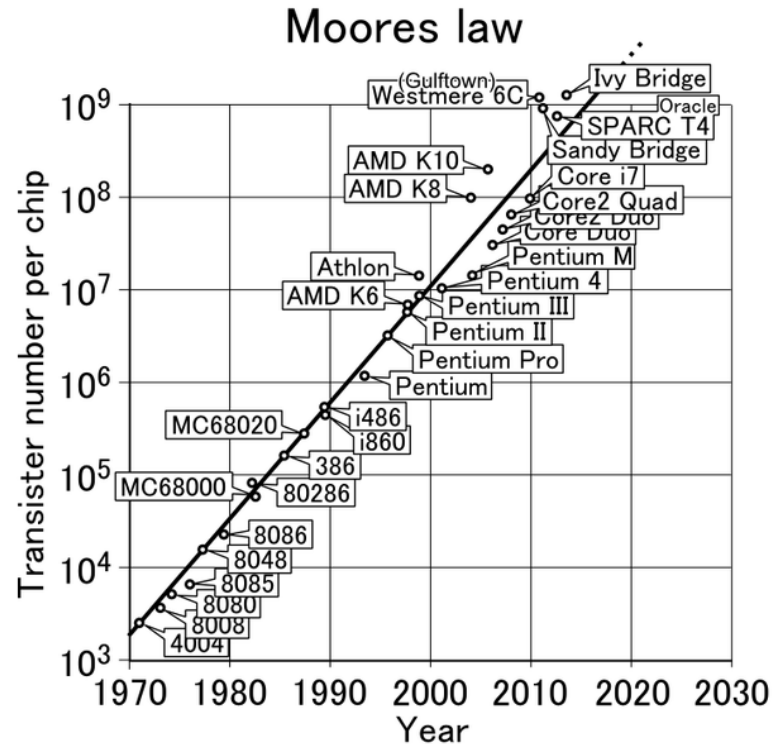
Giuseppe Di Guglielmo  
giuseppe@cs.columbia.edu

Columbia University

*Fast Machine Learning  
Fermi Lab  
September 10-13, 2019*



# Technology Trends

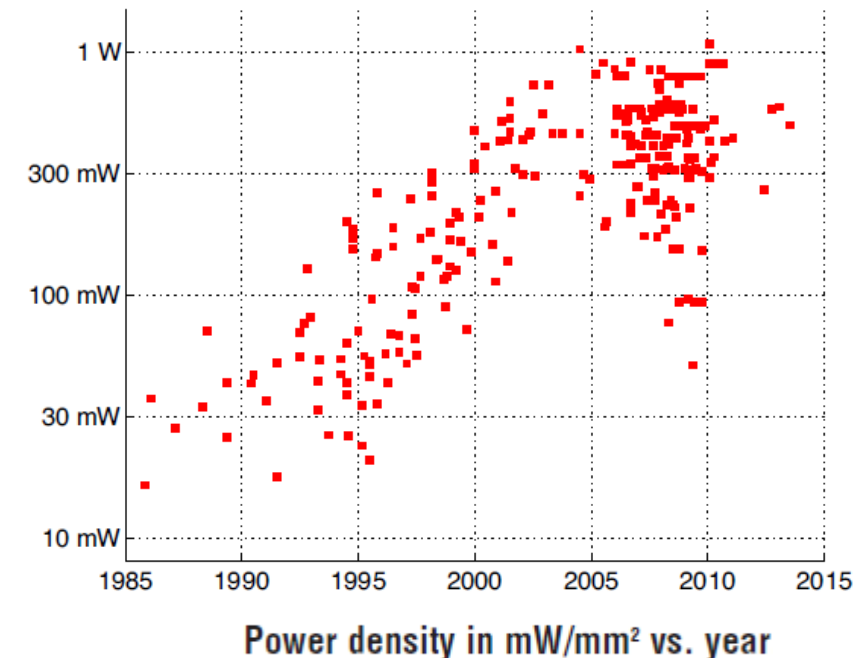


## • Moore's Law had many lives

- 2004: The end of Moore's Law?
- 2015: Beyond Moore's Law
- 2016: After Moore's Law
- 2017: A new way to extend Moore's Law

The Economist

- ... but Dennard's Scaling has stopped
  - On-chip power density cannot continue to grow



Source: Horowitz et al., 2014

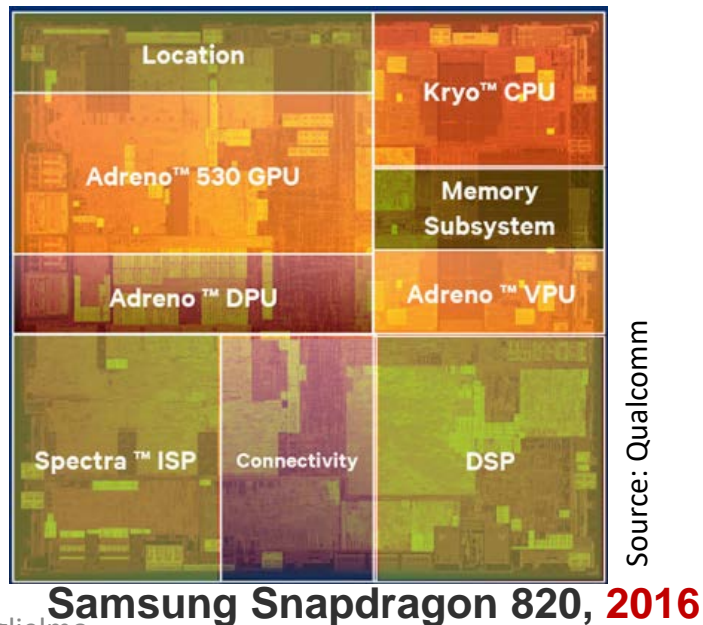
# Emerging Computing Platforms for AI

- **Heterogeneous Multi-Core Systems**

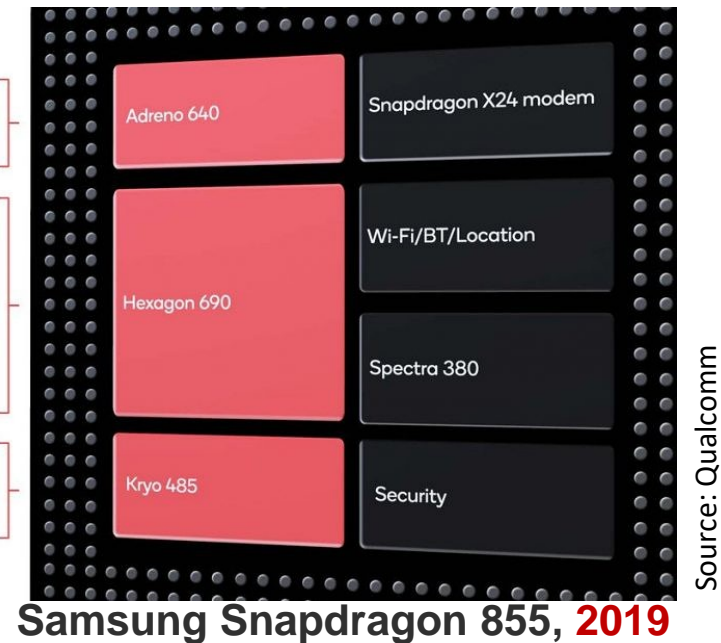
- Mix of processor cores and specialized hardware accelerators yields more energy-efficient computing

- The **approach to AI** is getting **heterogeneous**

- DSP, GPU, CPU, FPGA, custom hardware...
  - ... whatever is the best piece of hardware to handle an AI task in the most power efficient way



- 50% More ALUs  
FP32 & FP16
- New Tensor Accelerator  
4x Vector eXtensions
- Optimized Scalar  
Voice Assistant  
INT16, INT8 & Mixed
- New dot product  
instructions  
FP32 & INT8



# Heterogeneous Multi-core System-on-Chip

**intel**

**MOBILEYE**  
Our Vision. Your Safety.™

**March 2017**

**Google**

Cloud TPU

**Coral**

**May 2016 (1<sup>st</sup> gen.)**  
**May 2017 (2<sup>nd</sup> gen.)**

**amazon web services**

**EC2 F1**

**intel**

Hardware Development Kit → Custom Logic → Amazon FPGA Image (AFI) → AWS Marketplace → Attach your AFI to an F1 Instance

**November 2016**

SKT 0

SKT 1 Module

QPI

JTAG

Intel Xeon® E5-2600 v2 Product Family

ATERA® Stratix V

DDR3

HSS0 v8

Top Side High Speed Connector

**June 2016** Ivy Town Xeon + FPGA: The HARP Program

**intel**

**Myriad X**

8.8 mm

16 programmable vector engines

8.1 mm

Dedicated Accelerator for Neural Networks

Intelligent Memory Fabric

Dedicated Imaging and Vision Accelerators

**August 2017**

**XILINX**

**Zynq UltraScale+ MPSoC**

- ARM Cortex A53 Application Processors 64-bit Dual-Core with Virtualization
- ARM Cortex R5 Real-Time Processors 32-bit Dual-Core Application Offload
- mali H.265 HEVC Graphics/Video ARM Mali-480MP H.265/264 CODECS
- UltraScale FPGA Logic UltraRAM, PCIe Gen4, 10Gb Ethernet, AMS
- Power Management Multiple Power Domains Power Gated Islands
- ISO Safety & Reliability IEC61508, ISO26262 System Isolation & Error Mitigation, Lockstep
- Security Information Assurance, Trust, Anti-Tamper, TrustZone Key and Vault Management
- High Speed Peripherals USB 3.0, PCIe Gen2, GbE SATA2.0, DisplayPort
- Runtime SW & Tools OS, RTOS, AMP Hypervisor Development, Heterogeneous Debug, Hardware/Software Profiling & Performance Analysis

**December 2013**

Sources: Intel, Google, Xilinx, Amazon  
Real-time AI Systems

**XILINX**

**Versal**

Programmable I/O

PCIE & CCIX

DDR

HBM

112Gbps

58Gbps

32Gbps

Multi-Gig Ethernet

600G Const

Mini LVDS

RF

Network-On-Chip

Platform Management Controller

Dual-core Arm Cortex-R5 Real Time Processor

Dual-core Arm Cortex-A72 Application Processor

Adaptable Hardware Engines

Intelligent Engines

AI Engines

DSP Engines

Scalar Engines

**October 2018**

# Accelerator

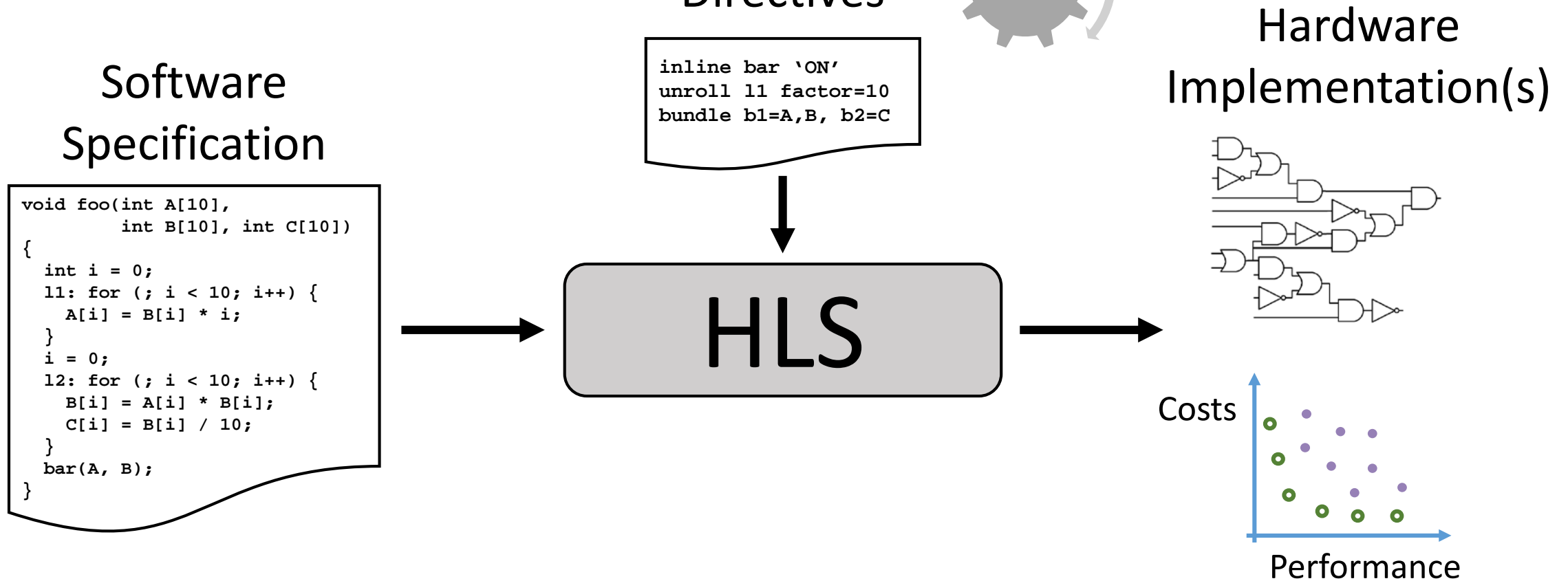
- A **special-purpose hardware** that is optimized to perform a **dedicated function(s)** as part of a **general-purpose computational system**
  - While being part of a larger system, it spans a scale from being closely integrated with a general purpose processor to being attached to a computing system
  - Increasingly accelerators are migrating into the chip, thus leading to the rise of heterogeneous multi-core SoC architectures”
- Implemented as
  - Application Specific Integrated Circuit (**ASIC**)
  - Field-Programmable Gate Array (**FPGA**)

Source: CSEE\_E4848, Columbia University

# High-Level Synthesis

**hls4ml**

- RF
- Precision
- Latency/Resources



# AI > ML > NN

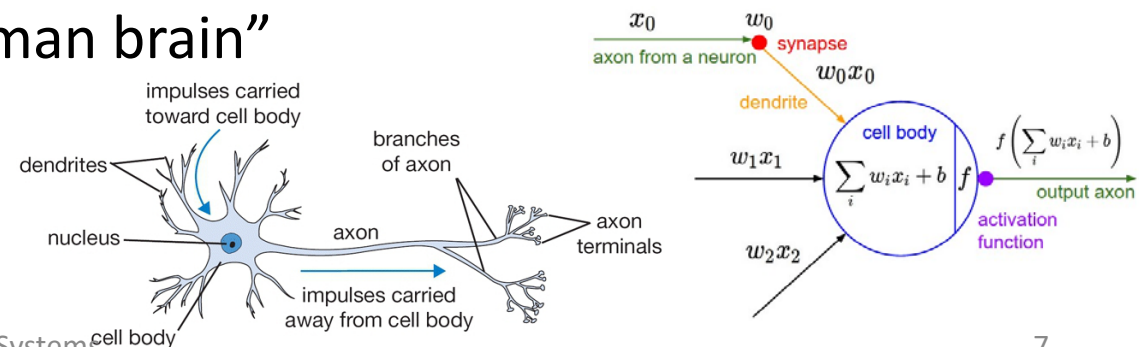
- **Artificial Intelligence (AI)**
  - Computer Vision
  - Pattern Recognition
  - ...
- **Machine Learning (ML)**
  - Linear Regression
  - K-Means Clustering
  - Decision Trees
  - ...
- **Neural Networks (NN)**
  - Convolutional Neural Networks
  - Binary Neural Networks
  - Recurrent Neural Networks
  - ...

“Machine mimics **cognitive functions** such as learning and problem solving”

“Gives computers the ability to **learn** without being explicitly programmed”

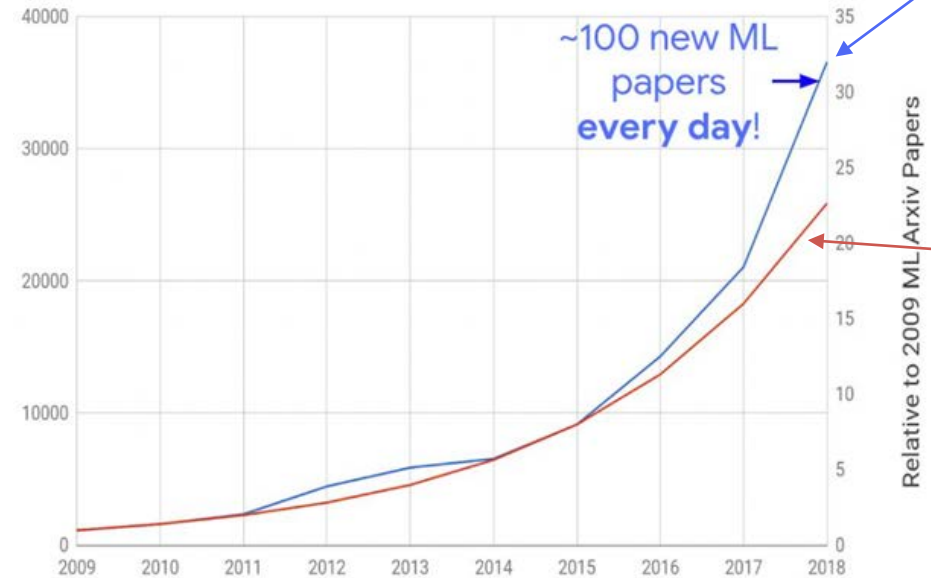
Training Phase  
Inference Phase

“Mimics the **physiology** of human brain”



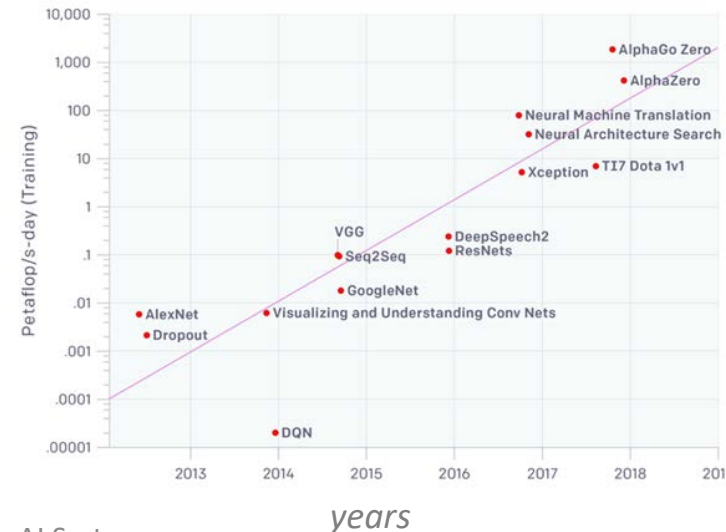
# Exponential Growth in Deep Learning

- There are several factors
  - Algorithms
  - Hardware
  - Data
  - Capital
  - Talent
  - Application



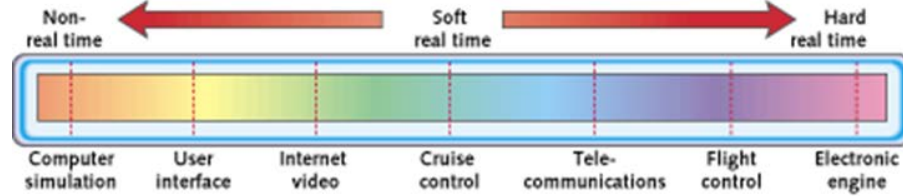
ML arXiv Papers

Moore's Law growth rate (~ 2x/2 year)



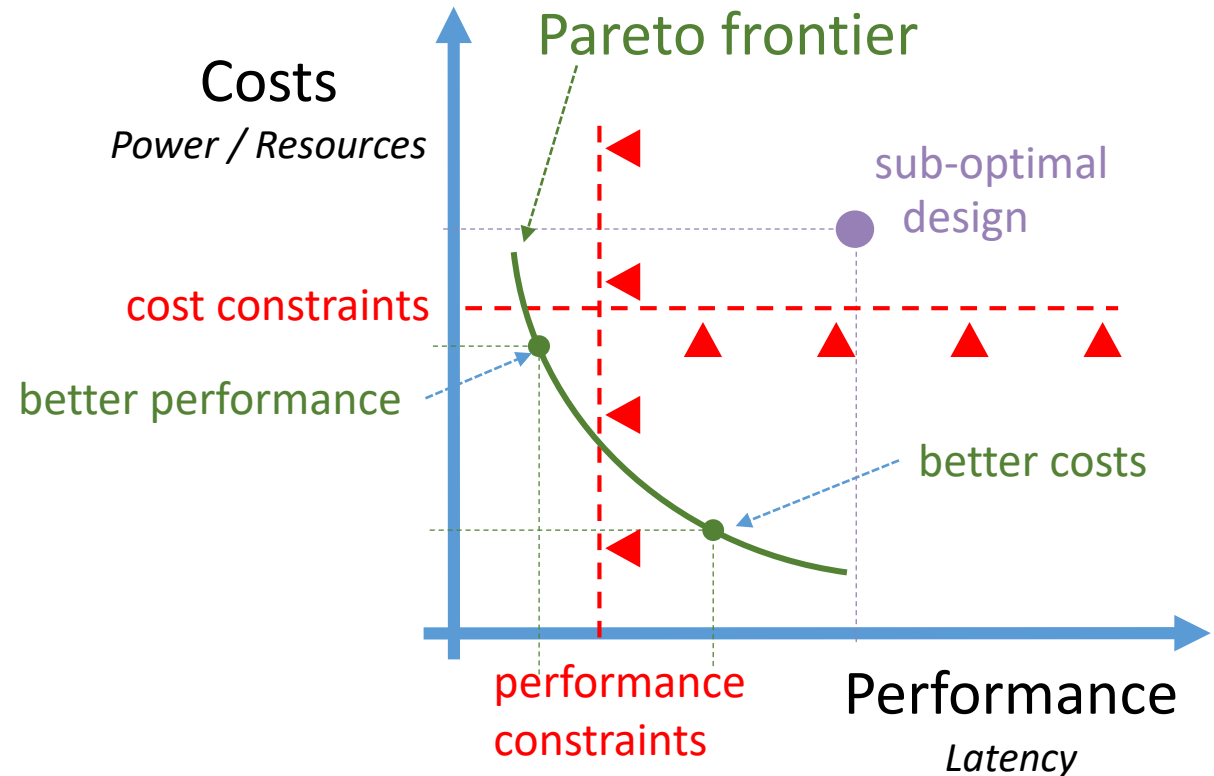


# (Hardware) Real-Time System



Source: eet.com

- Deadline-driven design
  - Constraints
    - Period/Frequency
    - Deadline
- Predicable system
  - “all tasks meet all deadlines”
- **Deterministic** system
  - A predictable system where the “timing behavior can be predetermined”



- HW design **optimization** space
  - Pareto optimality: “we cannot improve performance without paying costs, and vice versa”

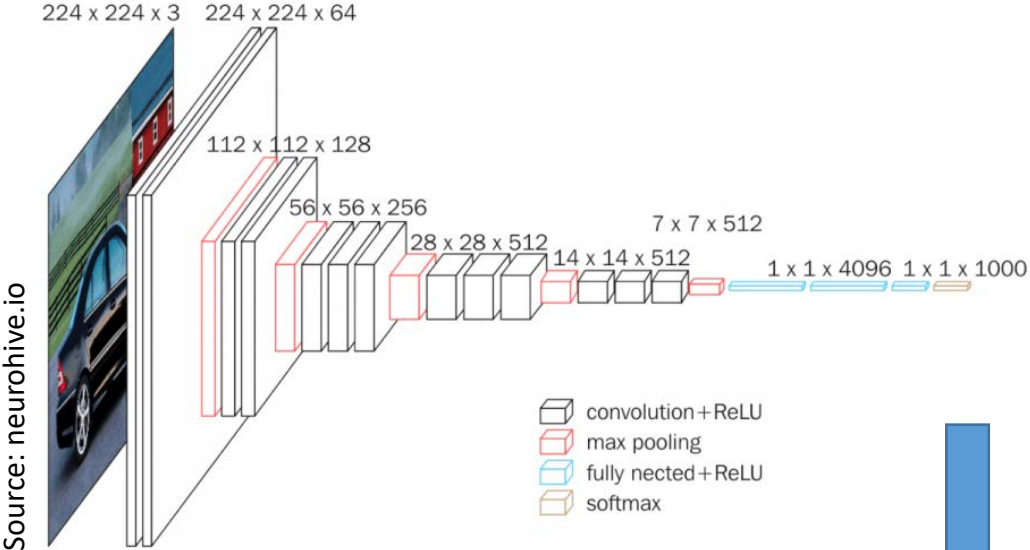
# A Too Short Blanket Problem



Source: <http://dogatesketchbook.blogspot.com/2008/02/illustration-friday-blanket.html>

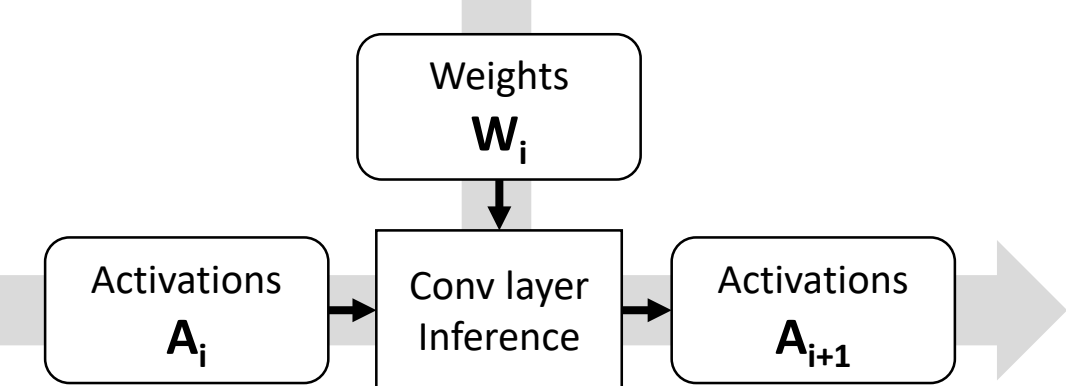
# NN = Deterministic Performance

- We can **statically** compute
  - Memory occupation
  - Number of operations
  - Data movement
  - ...



### Convolutional layer

- $W_i = IN\_Ch_i * OUT\_Ch_i * F_i^2 * bits$
- $A_i = BATCH * F_i^2 * IN\_Ch_i * bits$
- $A_{i+1} = BATCH * F_i^2 * OUT\_Ch_i * bits$

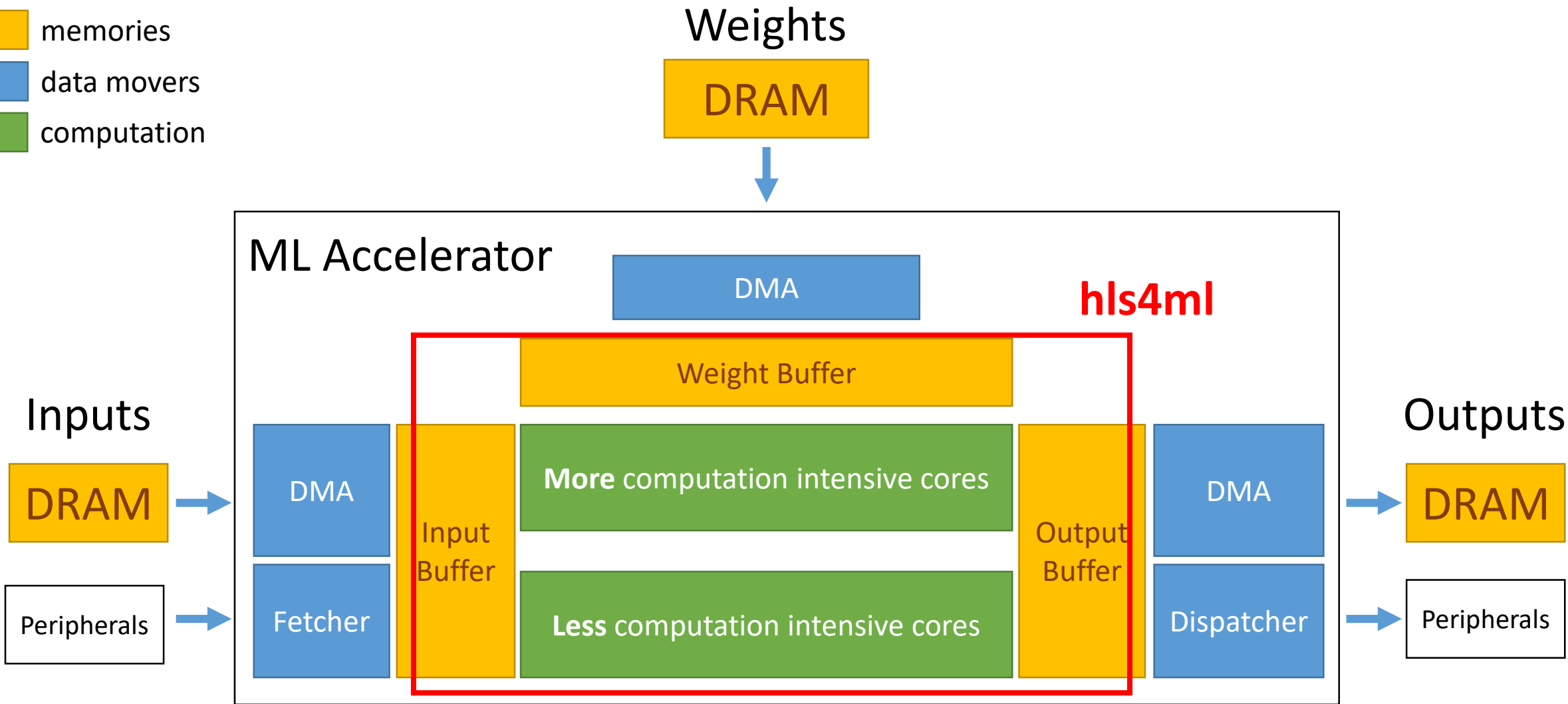


Feature Maps - Input						
	H	W	C	Size (#)	Size (B)	Size (MB)
1_1	224	224	3	150528	602112	0.5742188
2_1	224	224	64	3211264	12845056	12.25
2_2	112	112	64	802816	3211264	3.0625
3_1	112	112	128	1605632	6422528	6.125
3_2	56	56	128	401408	1605632	1.53125
3_3	56	56	256	802816	3211264	3.0625
3_4	56	56	256	802816	3211264	3.0625
4_1	28	28	256	200704	802816	0.765625
4_2	28	28	512	401408	1605632	1.53125
4_3	28	28	512	401408	1605632	1.53125
5_1	14	14	512	100352	401408	0.3828125
5_2	14	14	512	100352	401408	0.3828125
5_3	14	14	512	100352	401408	0.3828125
AVG	75	75	286			

Big Excel file

# Architecture of a ML Accelerator

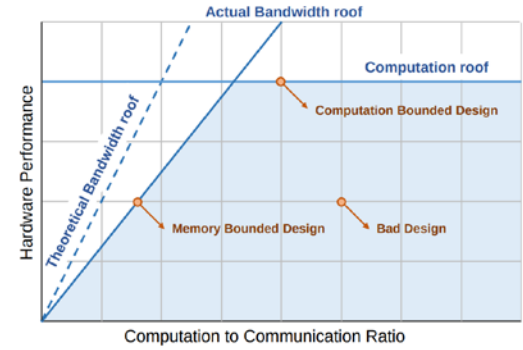
- memories
- data movers
- computation



# Pain Points\*

- memories
- data movers
- computation

Weights and Inputs fetching: bandwidth vs. performance



A Survey of FPGA Based Neural Network Accelerator, K. Guo et al., May 2018

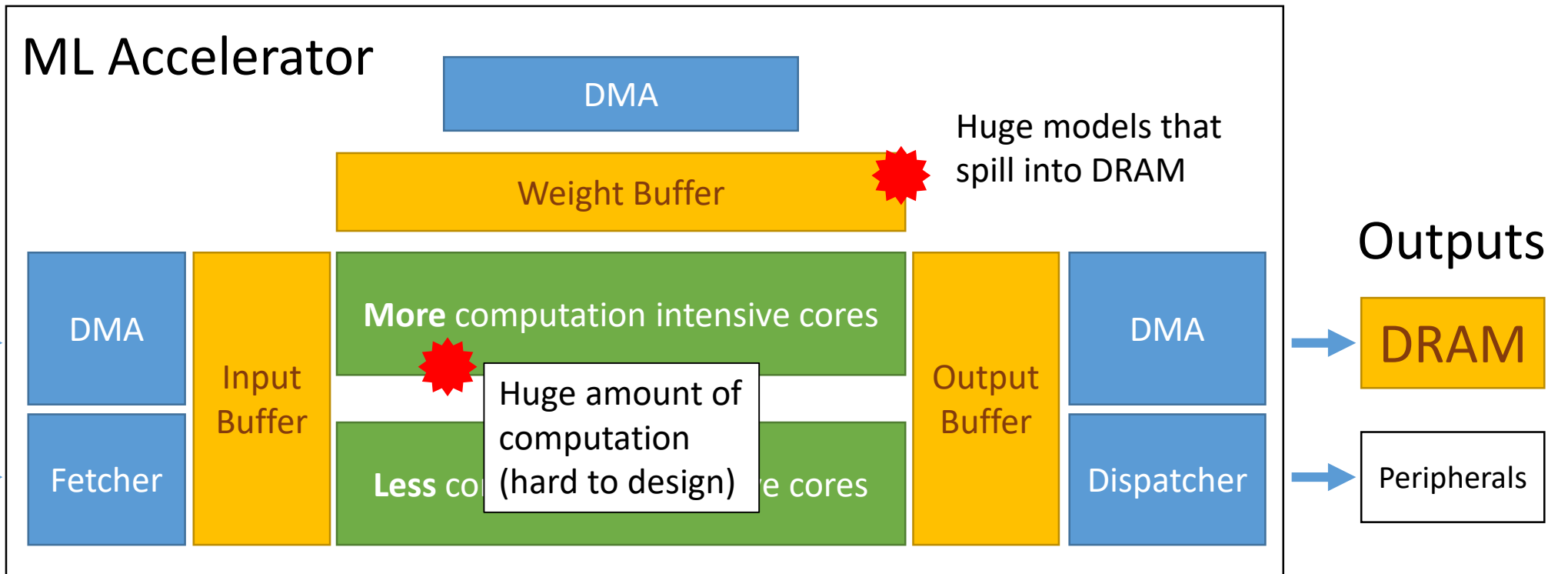
Caching and interconnects may introduce unpredictability

## Weights



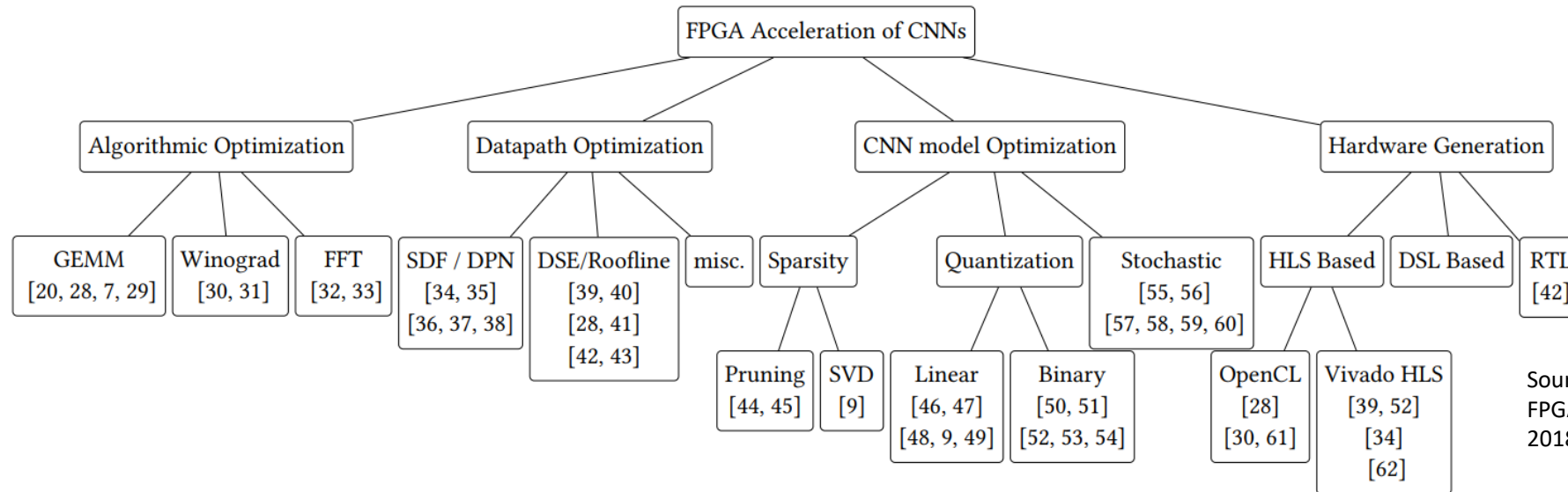
Power consumption (embedded domain)

Latency (real-time domain)



\* Michaela Blott, Principal Engineer in Xilinx, Architecture for Accelerating DNNs, Hot Chips 2018

# Cost & Performance Optimization Techniques



Source: Accelerating CNN Inference on FPGAs: A Survey, K. Abdelouahab et al., 2018

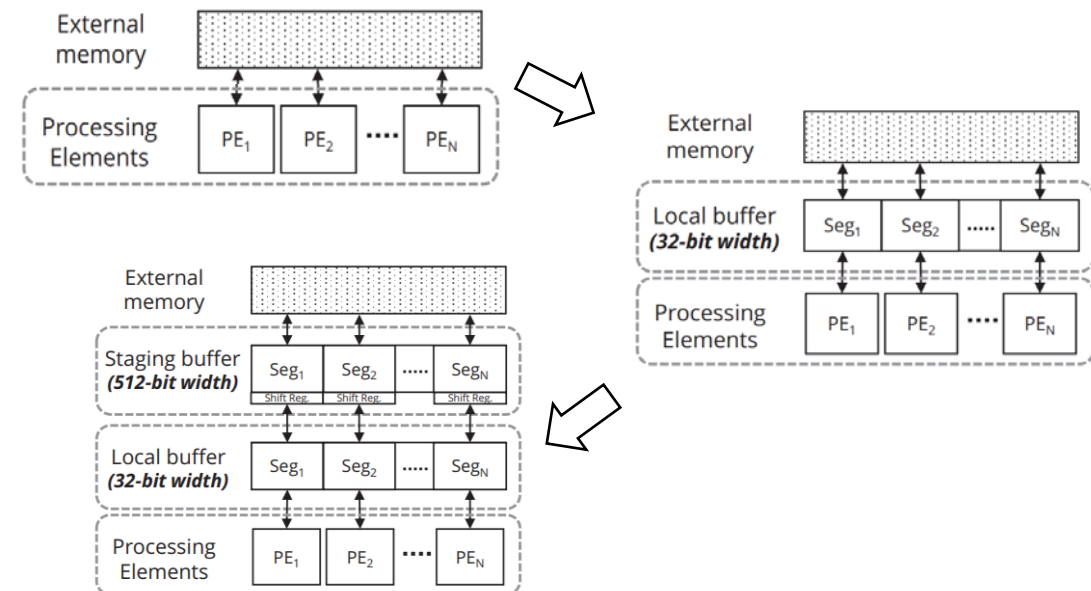
- Accelerating CNN Inference on FPGAs: A Survey, K. Abdelouahab et al., Jan 2018
  - <https://arxiv.org/pdf/1806.01683>
- Toolflows for Mapping Convolutional Neural Networks on FPGAs: A Survey and Future Directions, S. I. Venieris et al., Mar 2018
  - <https://arxiv.org/pdf/1803.05900>
- A Survey of FPGA Based Neural Network Accelerator, K. Guo et al., May 2018
  - <https://arxiv.org/pdf/1712.08934>
- Optimizing Memory Efficiency for Deep Convolutional Neural Networks on GPUs, C. Li et al, 2016
  - <https://arxiv.org/pdf/1610.03618>

# A Good Start is Half of the Work

- **Kernel** (Task)
  - Functionally and computationally important portion of an application
  - Well defined interface (inputs and outputs)
  - E.g. Conv2D
- **Algorithm**
  - An algorithm solves a particular kernel
  - E.g. Conv2D via GEMM, Winograd, FFT
- **Implementation**
  - Different implementations may exist of the same algorithm
  - HLS is very sensitive on the coding style

Source: MachSuite: Benchmarks for Accelerator Design and Customized Architectures, David Brooks, 2014

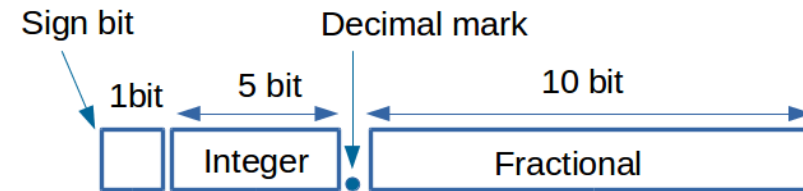
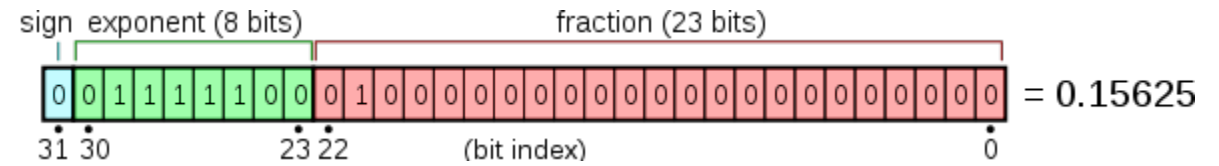
- Loop transformations to minimize memory access
  - Memory layout
- Create local memory buffers and keep data on-chip as much as



Source: Bandwidth Optimization Through On-Chip Memory Restructuring for HLS, Jason Cong et al., 2017

# Reducing Bit-Precision

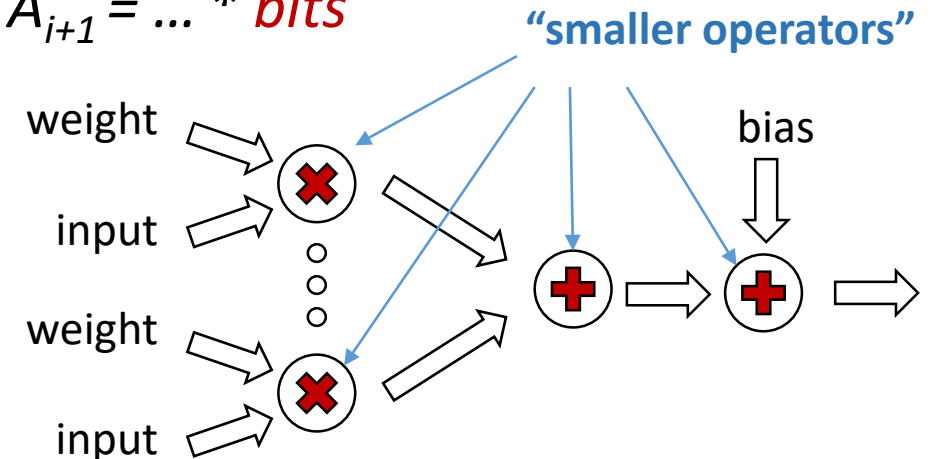
- Arithmetic
  - Floating Point
    - FP64, FP32, FP16, FP11 ...
  - Fixed Point (Linear Quantization)
    - ✓ A lot of freedom
    - ✓ “Essentially” integer arithmetic
    - ✗ Overflow and underflow problems
  - Binary quantization
- Linear reduction in **memory footprint**
  - Reduce the amount of data transfer
  - Model may fit local buffers (on-chip)
- Reduction of the **arithmetic logic**
  - Improve area, power, latency



$$W_i = \dots * \text{bits}$$

$$A_i = \dots * \text{bits}$$

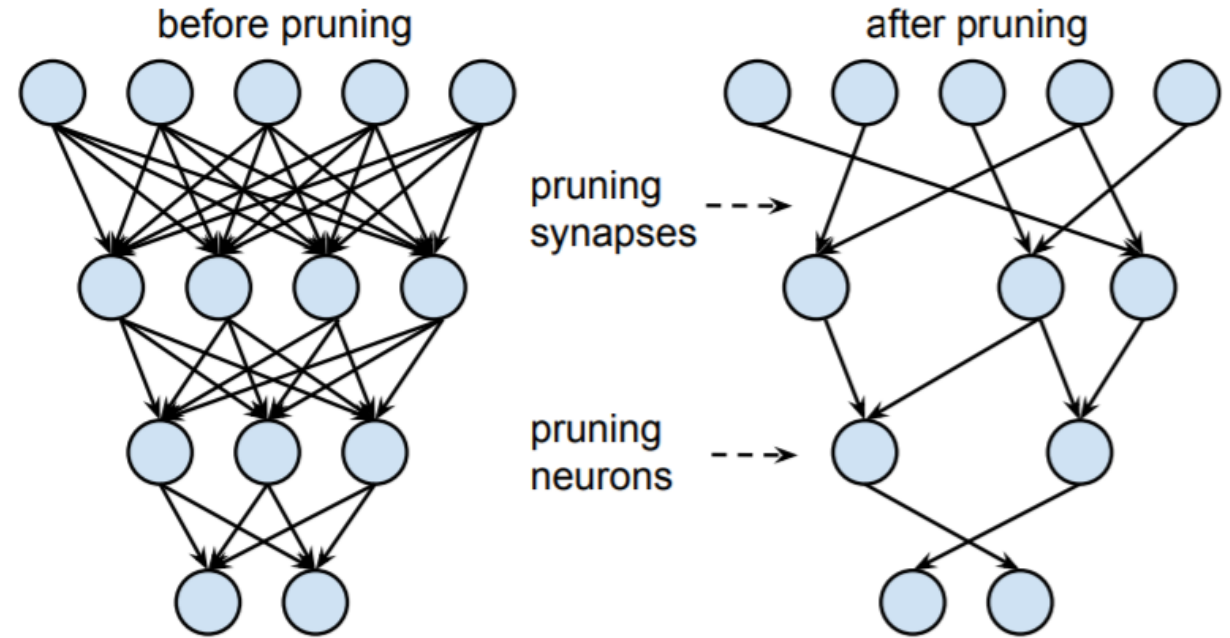
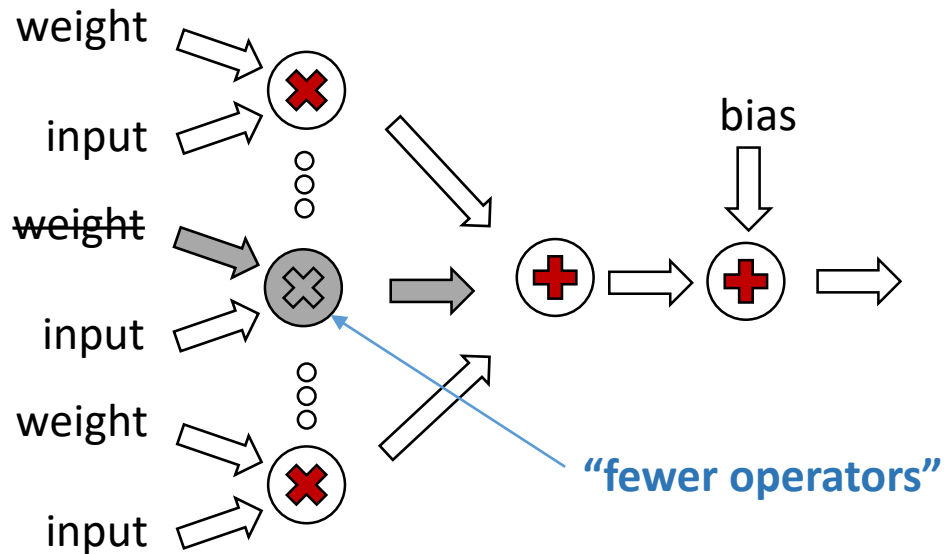
$$A_{i+1} = \dots * \text{bits}$$



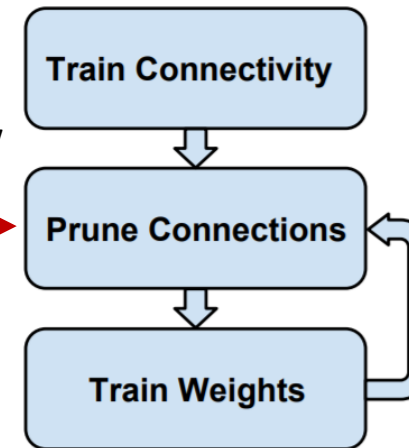


# Pruning

- Goal: Reduce storage, data transfer, and computation
- Reduction of the **model size** without loss of prediction accuracy
  - Alexnet 9x, VGG-16 13x



“We throw away those edges with weights below a certain threshold”



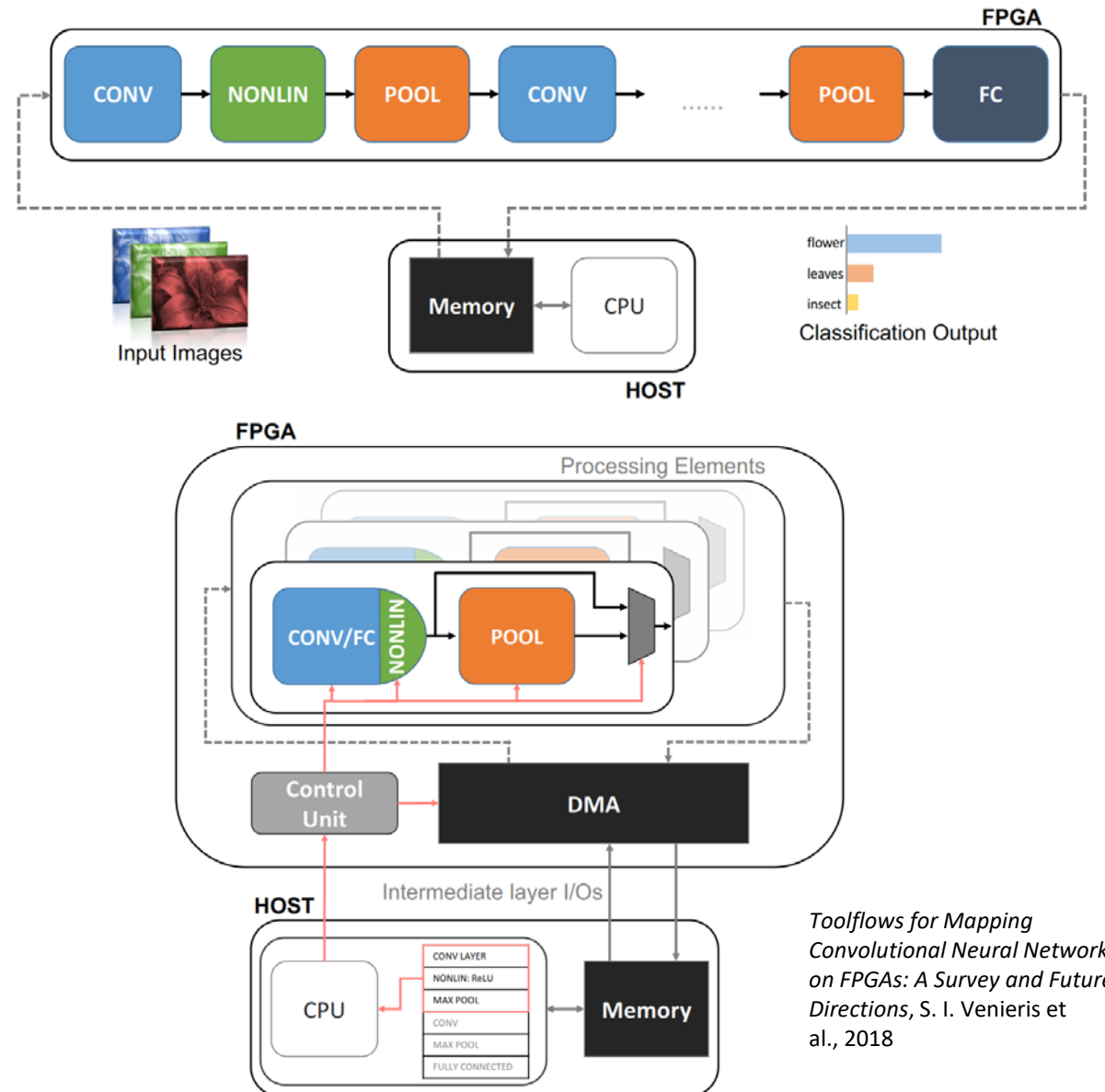
# Hardware Architecture

- Streaming

- One distinct hardware block for each CNN layer, where each block is optimized separately to exploit the inner parallelism of the layer
  - Weights: on-chip or off-chip
  - Controller: software or hardware

- Single Computation Engine

- Single computation engine that executes the CNN layers sequentially
  - Software controller of the scheduling of operations (may be inefficient)
  - Flexibility and reusability over customization
  - “One-size-fits-all” approach has higher performance on NN with a uniform structure



# Data Layout in Main Memory

- CNNs work on 4-dimensional matrices
  - Data can be stored in memory in 24 (=4!) different ways
- Data layout determines the memory-access patterns and has critical performance impact
- Algorithms and data layout should be designed to minimize the number of DMA transactions
  - Rule of thumb:
    - «Fewer and longer transactions»

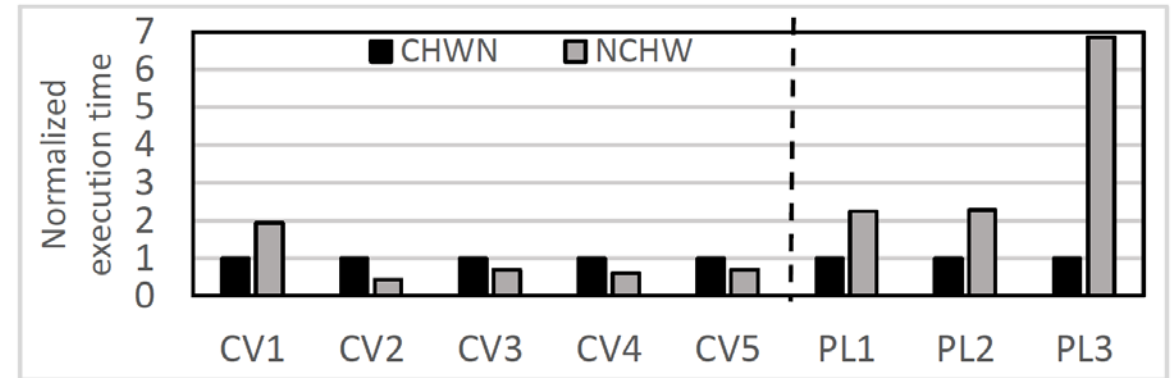
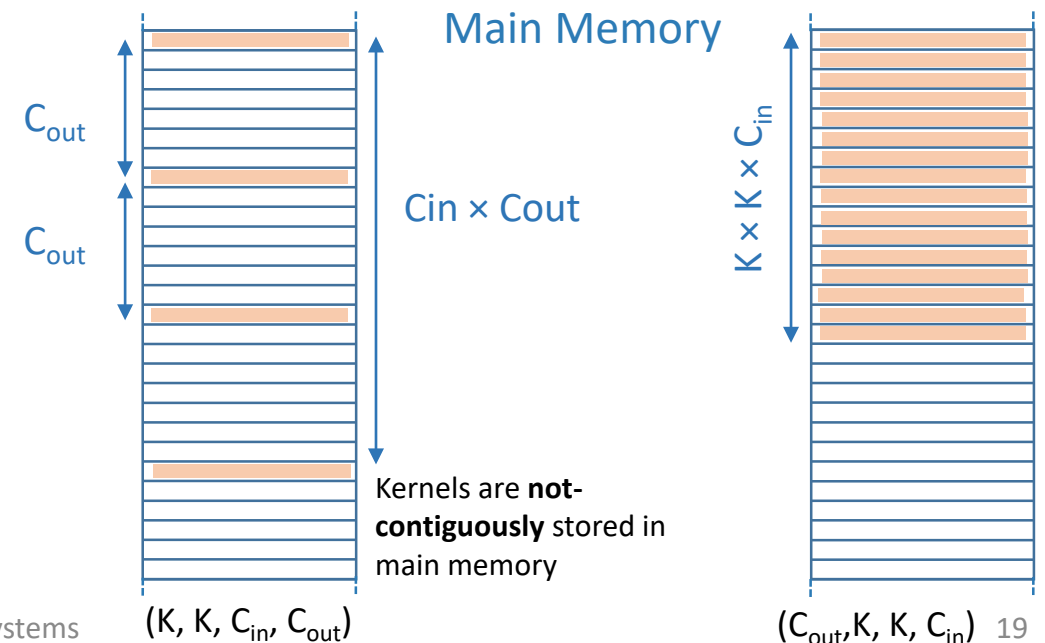


Fig. 1. Performance comparison between the CHWN layout (cuda-convnet2) and NCHW layout (cuDNNv4) on convolutional and pooling layers in AlexNet [12]

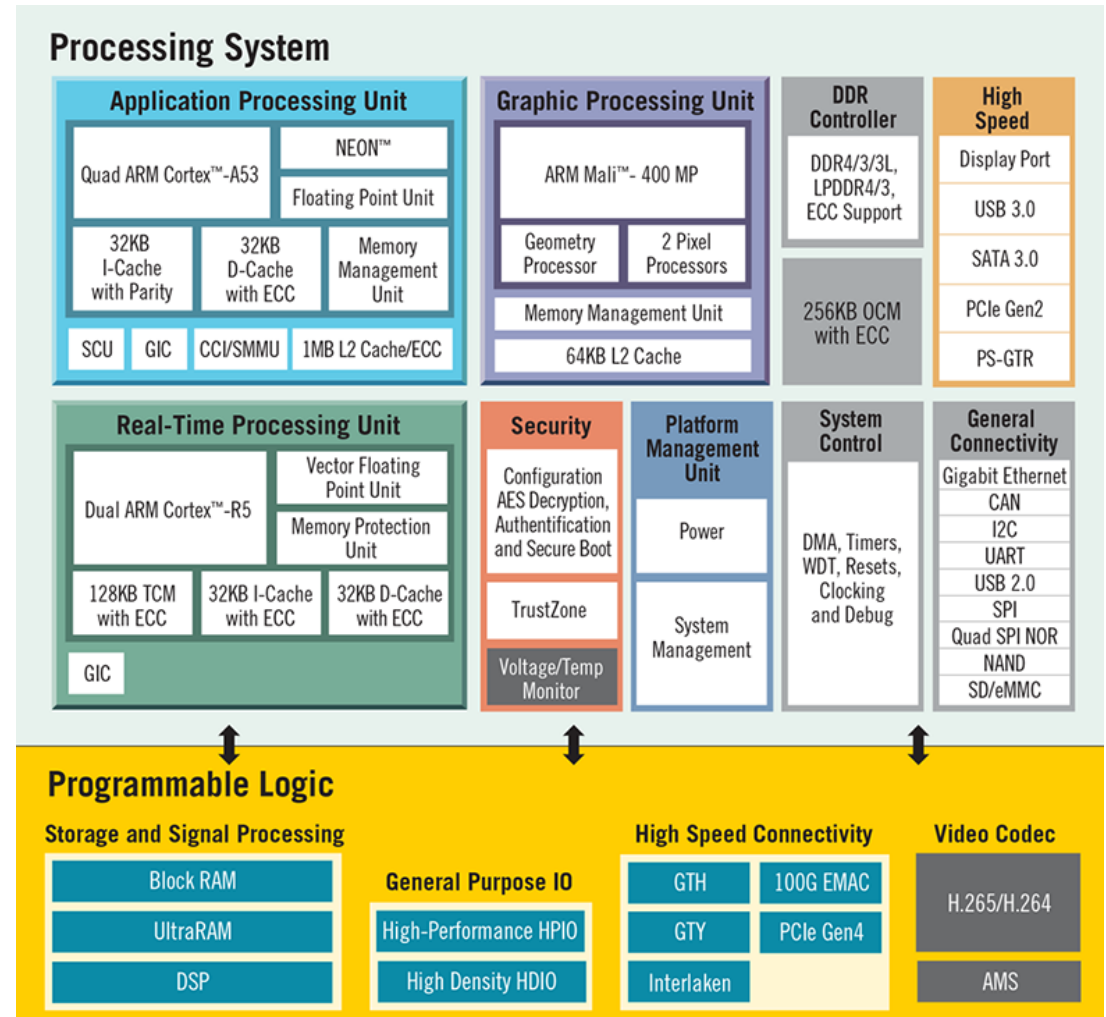
*Optimizing Memory Efficiency for Deep Convolutional Neural Networks on GPUs, C. Li et al, 2016*



# OS, Memory Hierarchy, RT-CPU

Zynq UltraScale+ MPSoC

- Real-time operating system
  - Predictable / Deterministic
  - Concept of running tasks
    - Read in data over an interface
    - Perform an operation on the data
    - ...
- “Caches are evil”
  - Disable caches
  - Cache partitioning



# Conclusions

- The design of a real-time AI system is more complicated than just meeting all of the deadlines
- Trade off between your timing constraints and area and power costs
- Exciting area to work on because of the constantly increasing importance of AI and variety of application domains



Source: <http://dogatesketchbook.blogspot.com/2008/02/illustration-friday-blanket.html>

Q/A

Thank you!