

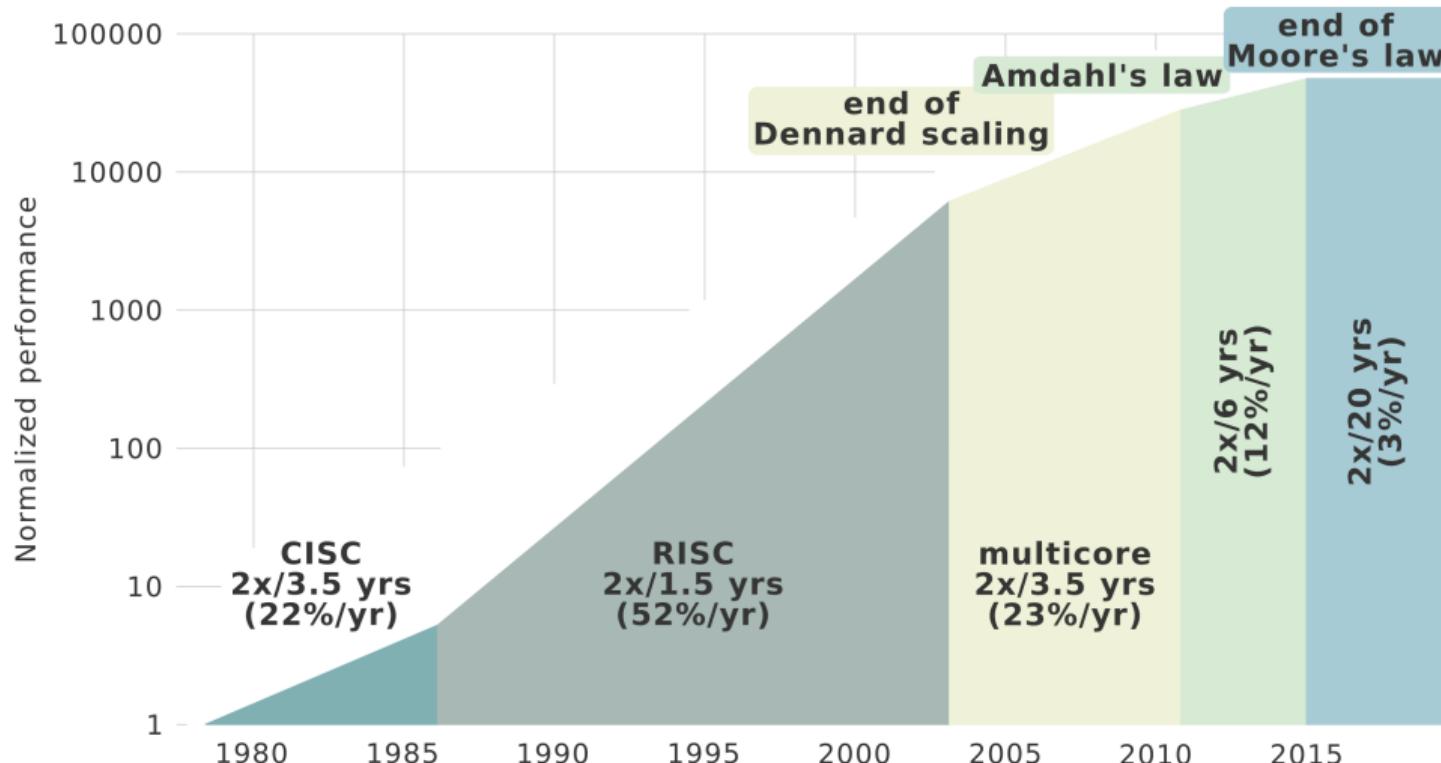
HW/SW Co-Design

SCF 2024

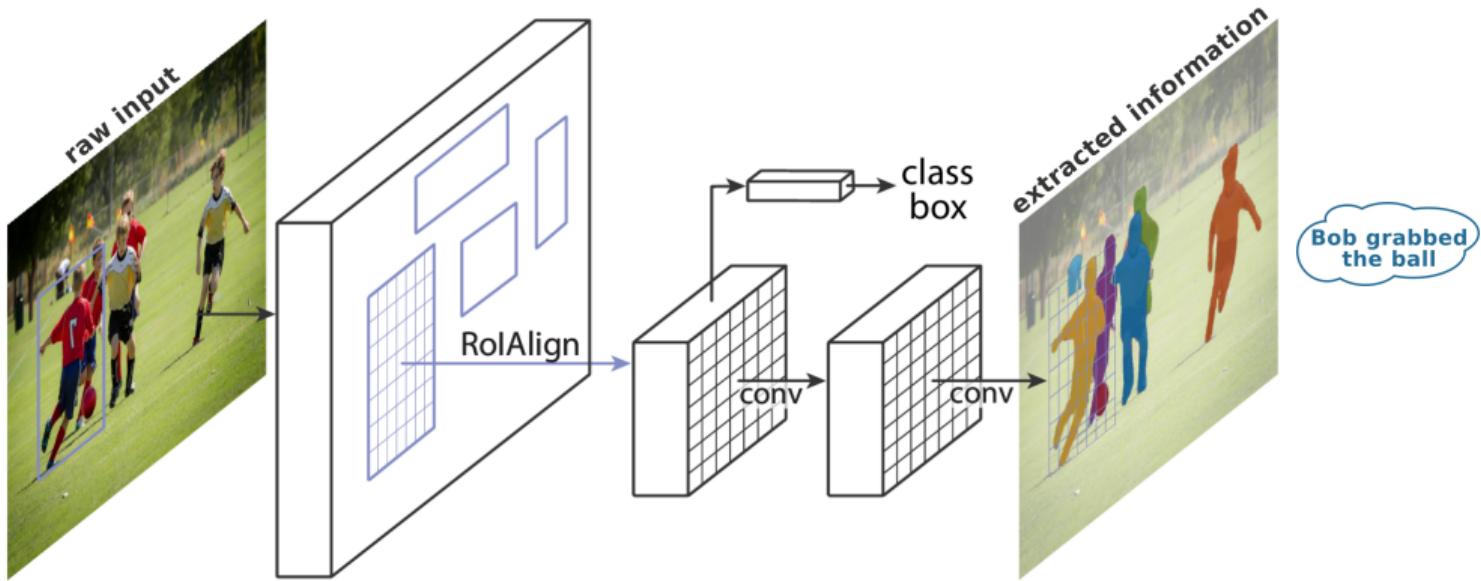
Alberto Dassatti



CPU PERFORMANCE IS STAGNATING



DATA AND PROCESSING



500 million
tweets



4 million
hours of content
on YouTube



4.3 billion
Facebook entries



3.6 billion
Instagram pictures

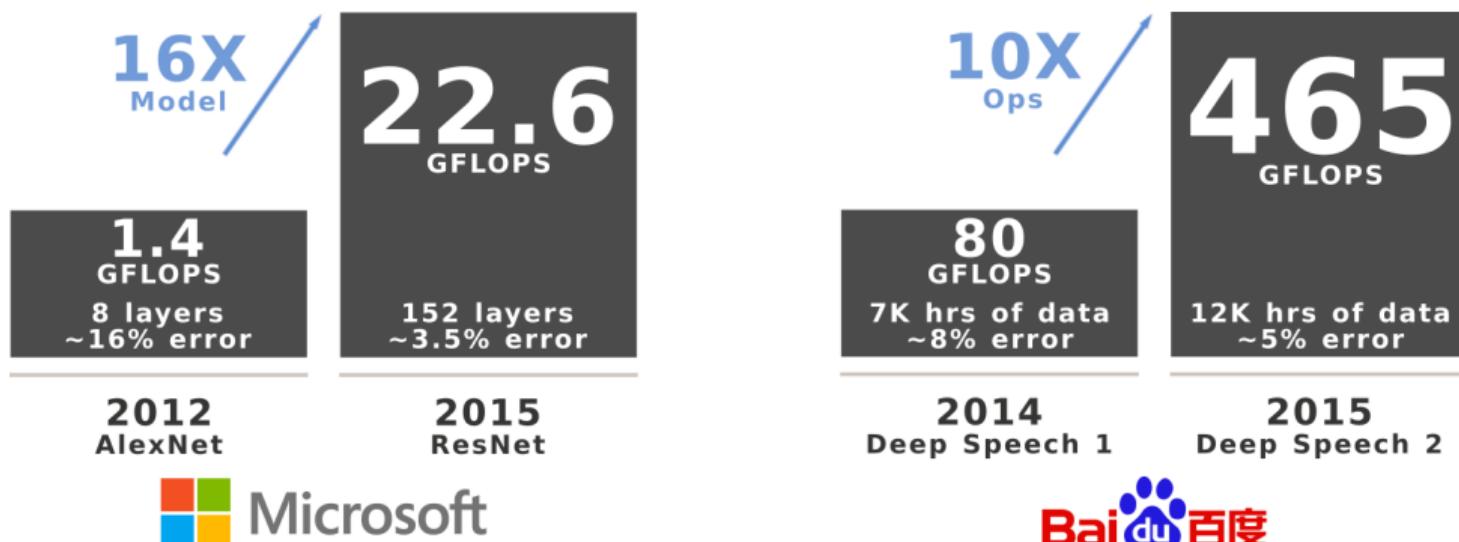


6 billion
Google searches



205 billion
e-mails

COMPUTATIONAL DEMAND IS INCREASING



PROCESSING IS POWER HUNGRY



PROCESSING IS POWER HUNGRY

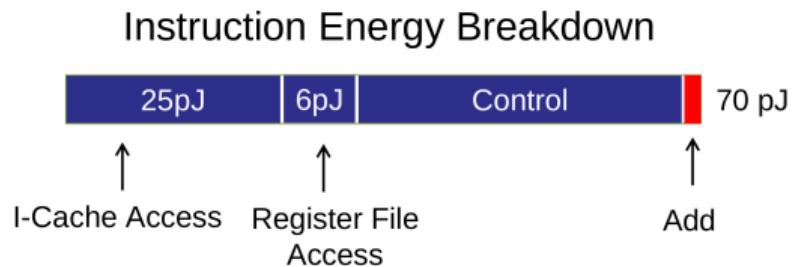


SOFTWARE: ENERGY COST OF COMPUTATION IS MINIMAL

ADD R3 R1 R2 // $R3 = R1 + R2$

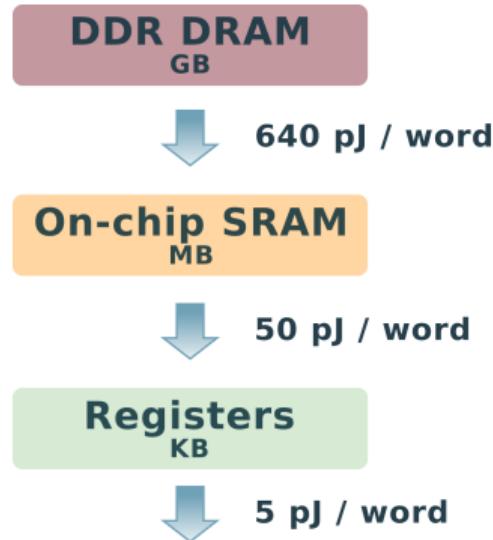
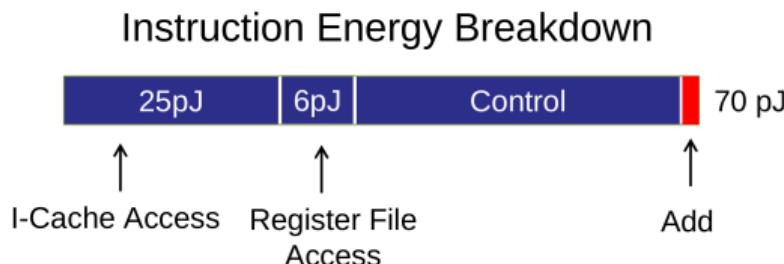
SOFTWARE: ENERGY COST OF COMPUTATION IS MINIMAL

ADD R3 R1 R2 // $R3 = R1 + R2$



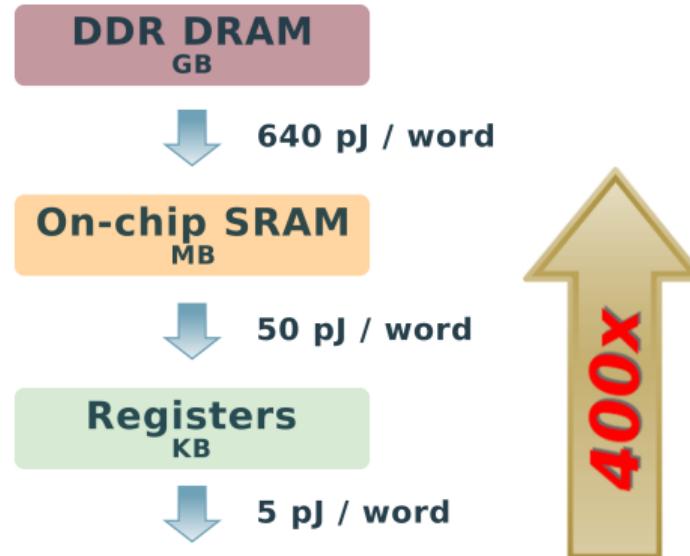
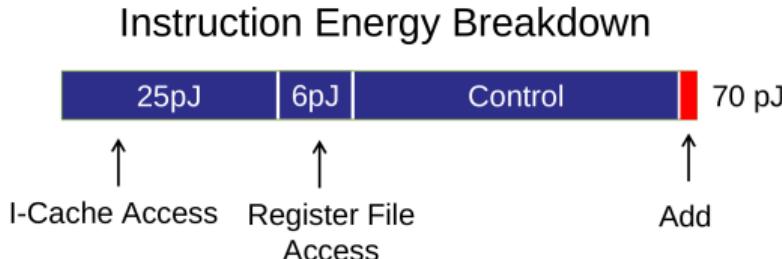
SOFTWARE: ENERGY COST OF COMPUTATION IS MINIMAL

```
LOAD R1 MEM_ADDR1  
LOAD R2 MEM_ADDR2  
ADD R3 R1 R2 // R3 = R1 + R2  
STORE R3 MEM_ADDR3
```

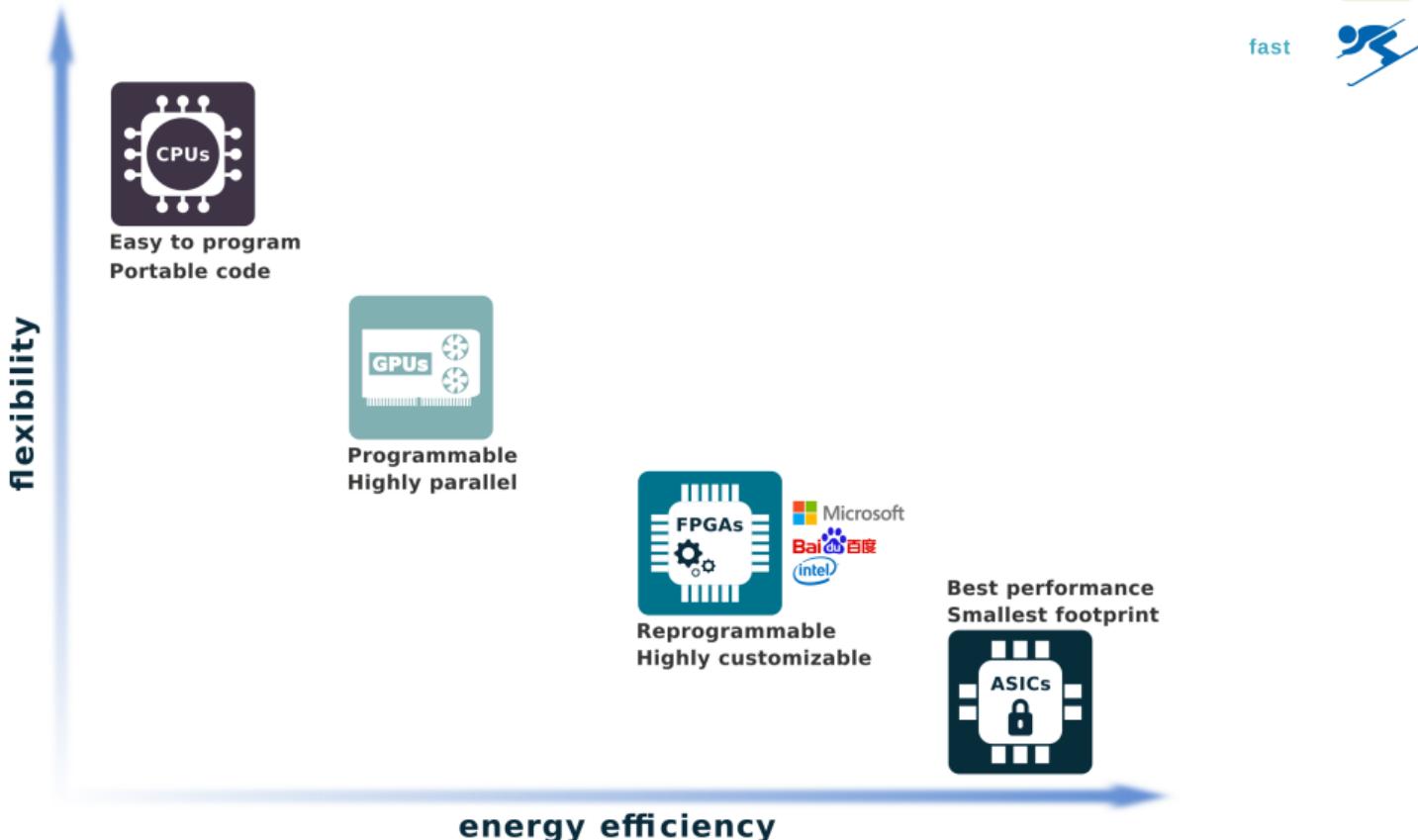


SOFTWARE: ENERGY COST OF COMPUTATION IS MINIMAL

```
LOAD R1 MEM_ADDR1  
LOAD R2 MEM_ADDR2  
ADD R3 R1 R2 // R3 = R1 + R2  
STORE R3 MEM_ADDR3
```



TECHNOLOGY VIEW



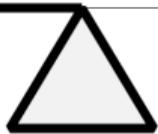


CLASSICAL APPROACH

Improve Performance
Improve Energy Efficiency
Reduce Power Density



**Implement
more in Hardware**

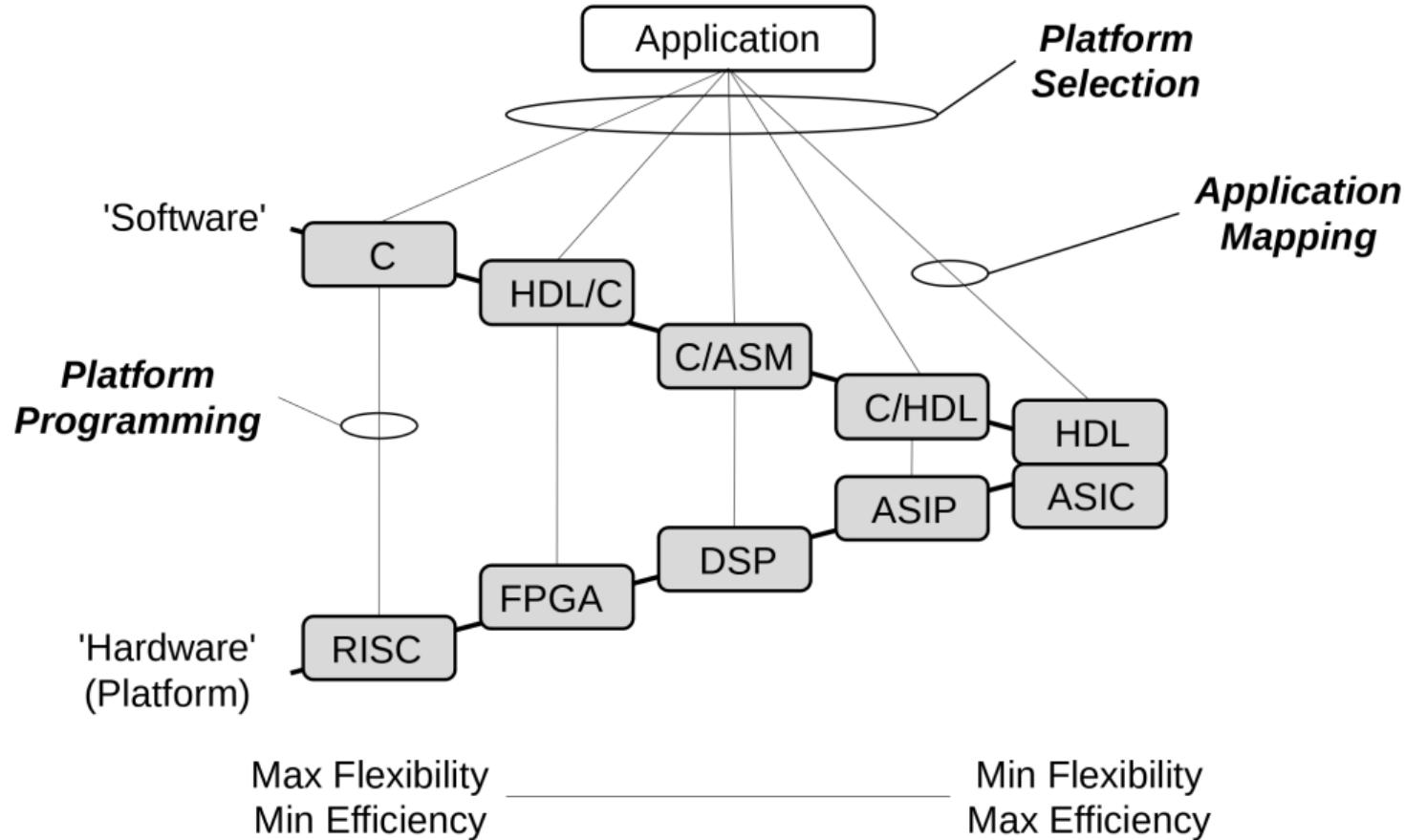


Manage Design Complexity
Reduce Design Cost
Stick to Design Schedule
Handle Deep Submicron



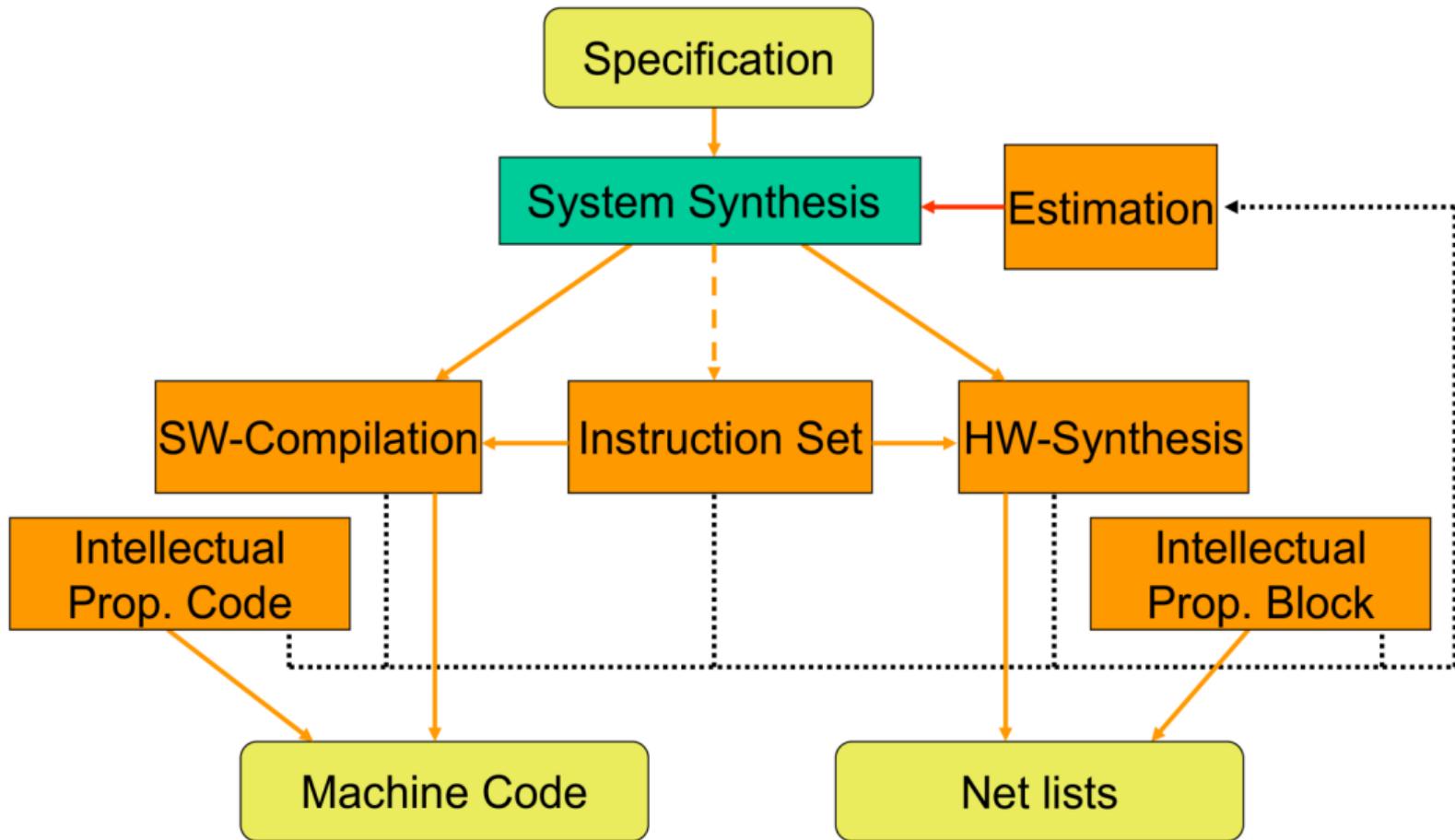
**Implement
more in Software**

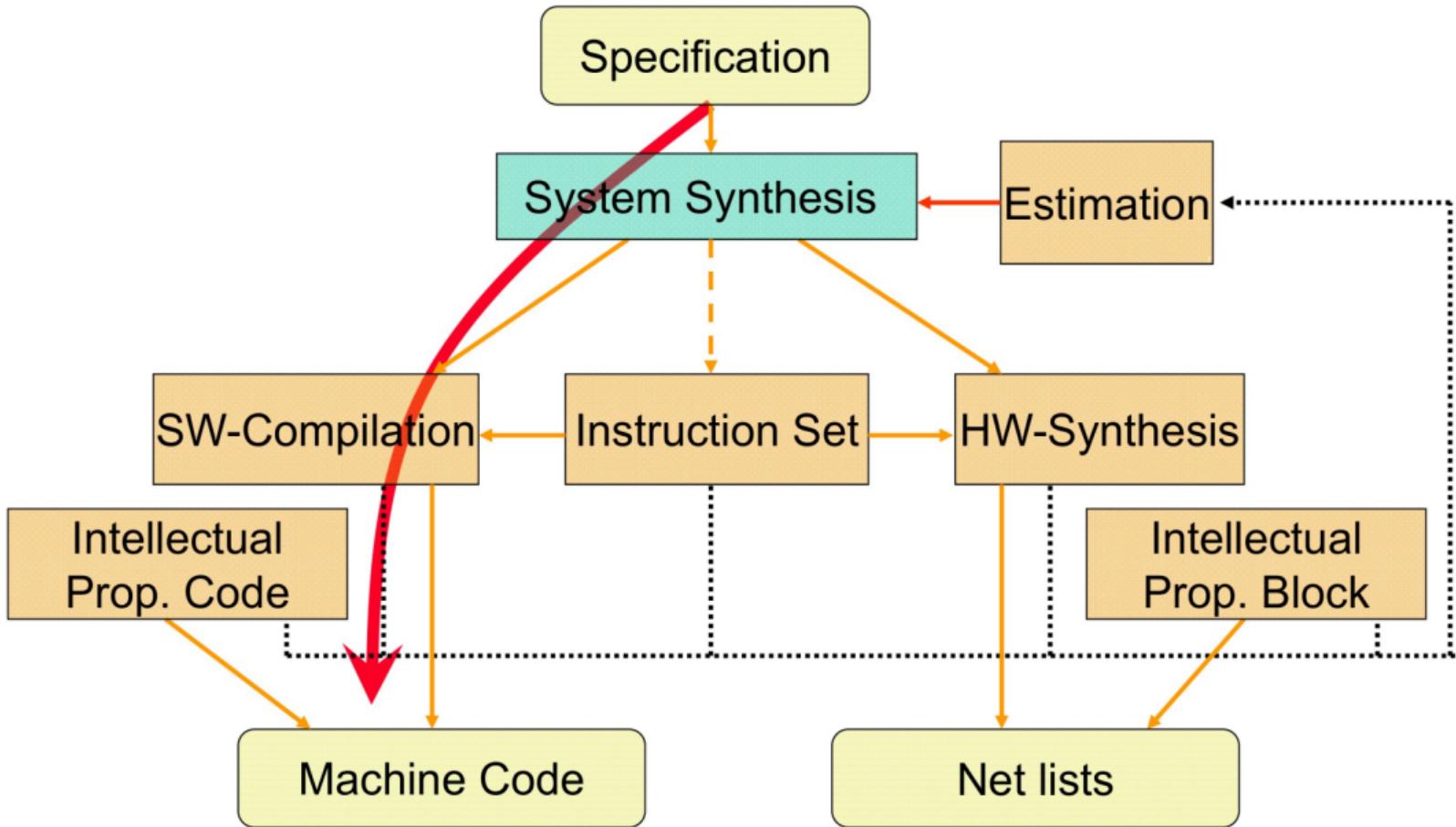
CLASSICAL APPROACH

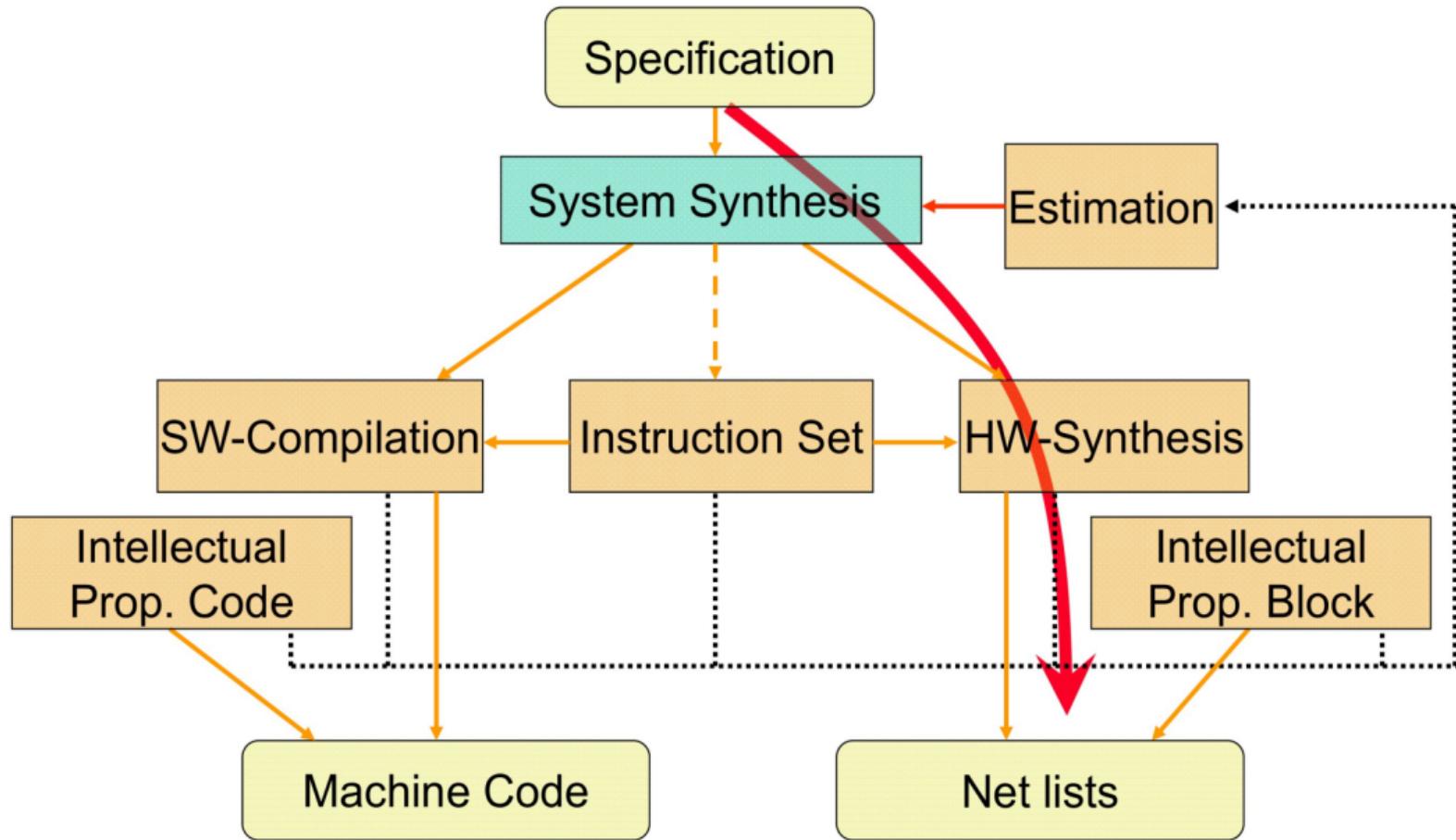


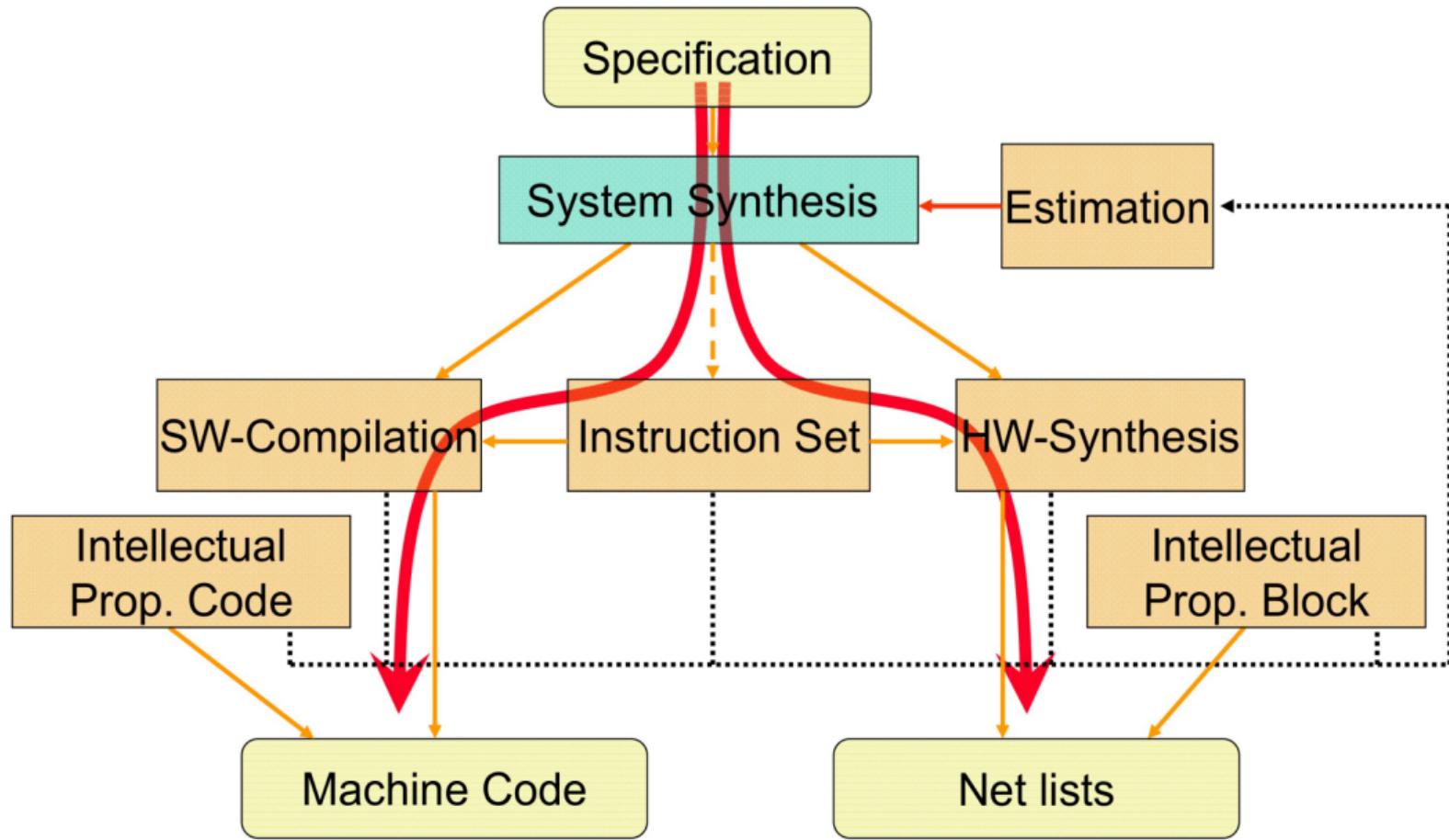










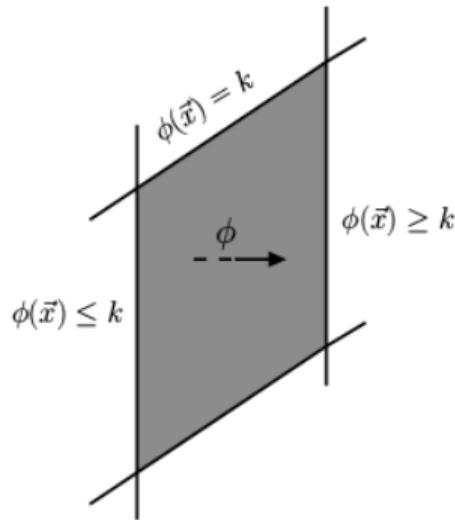
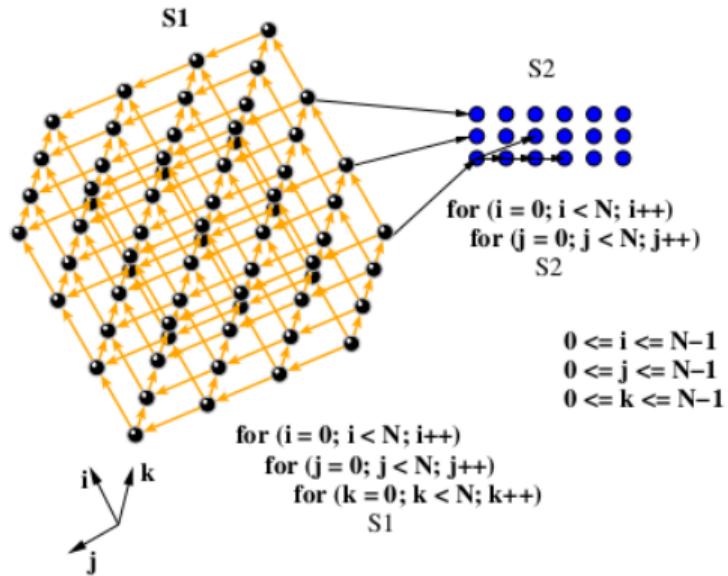


KEY TECHNOLOGIES

- ▶ Domain Specific Languages - DSL
- ▶ High Level Syntesis - HLS
- ▶ Formal methods
- ▶ Data oriented Design - DOD
- ▶ In-memory computing

DSL

- ▶ Capture the Expert's ideas - productivity and abstractions
- ▶ Decoupling computations and scheduling



HLS

- ▶ HDLs are too low level
- ▶ Use high level languages (C/C++; SYCL, OpenCL) to capture the computations
- ▶ Not fully automated, yet.
- ▶ There is a need for a run-time coordination
- ▶ Model the CPU-Accelerator paradigm, not much more

- ▶ Middle-ground: CGRAs (TFA)

FORMAL METHODS

- ▶ Math is super useful as abstraction (3k years proven!)
- ▶ Remove the need of simulation
- ▶ Remove debugging and tooling!
- ▶ They can be used for both design and verification
- ▶ We need more of them!

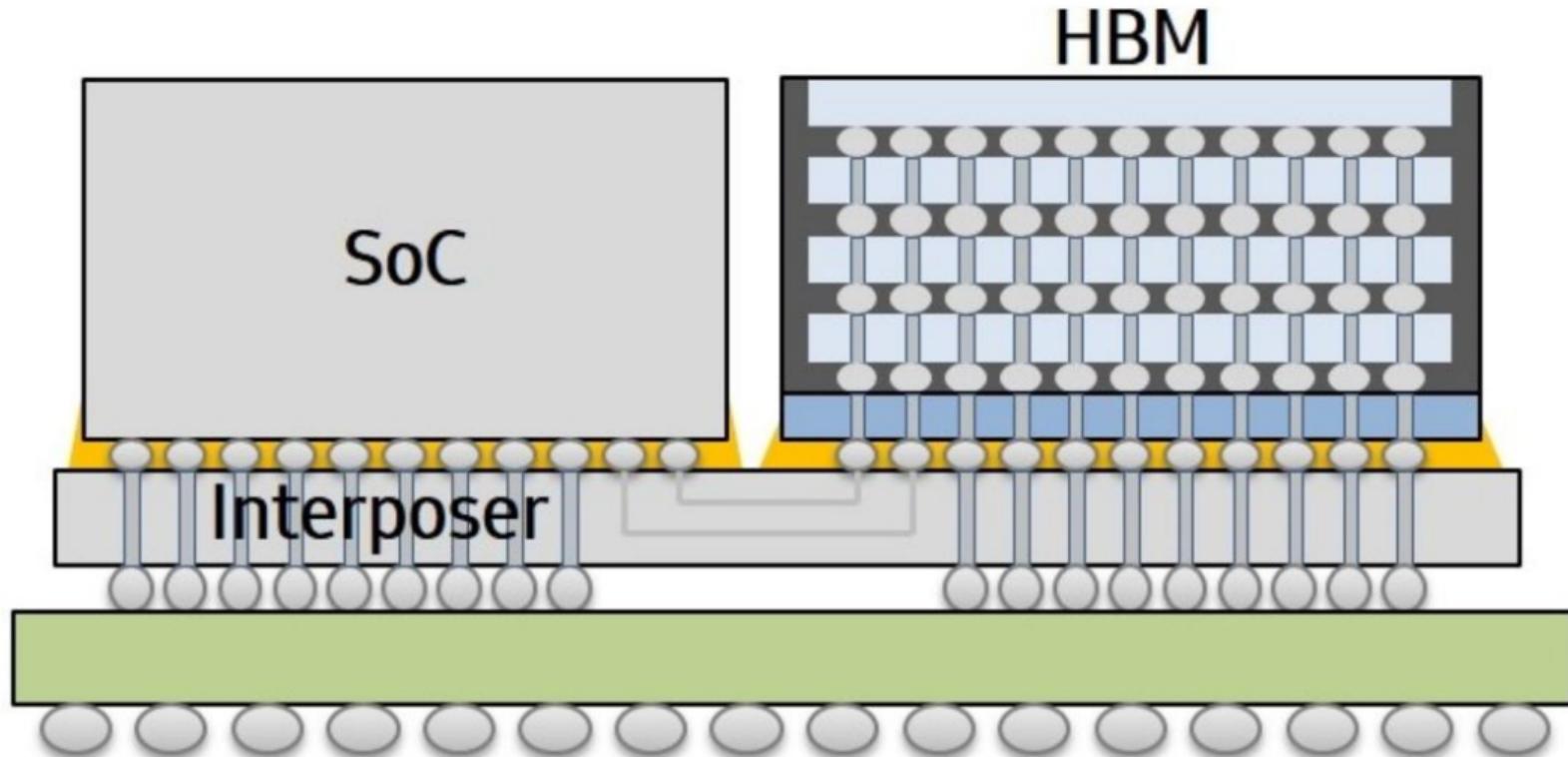
DoD

- ▶ Most code is written in OO;
- ▶ Programmers are trained to think in OO, but the execution overhead is huge
- ▶ DoD pushes to think about data and transformations
- ▶ Very popular in Game design
- ▶ Great performance
- ▶ Adoption requires a mental shift, new tools, and methodologies

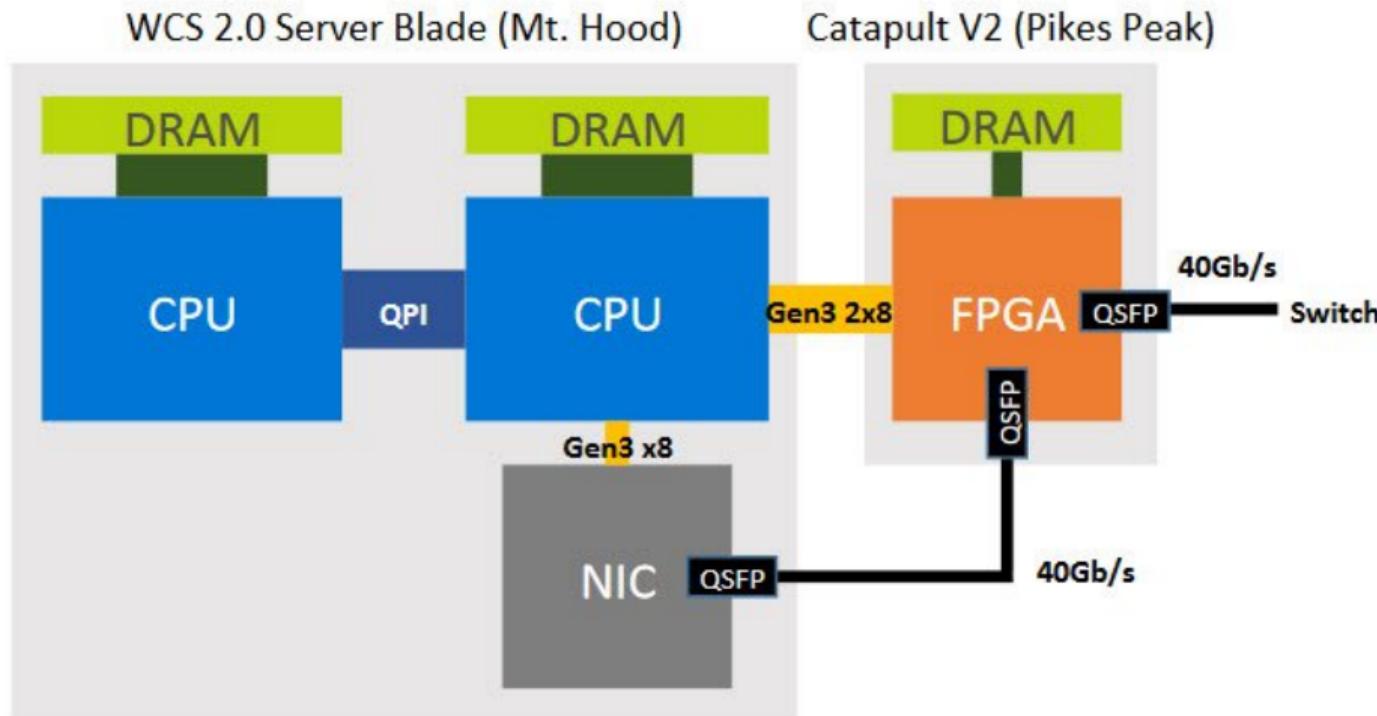
IN-MEMORY COMPUTING

- ▶ Instead of moving data, move processing (City center analogy)
- ▶ Requires new devices and abstractions (Smart-Drive)
- ▶ Requires new technological processes
- ▶ The energy gain can be huge
- ▶ Poor-man solutions: on-the-wire processing

REAL WORLD EXAMPLES: MEMORY INTEGRATION

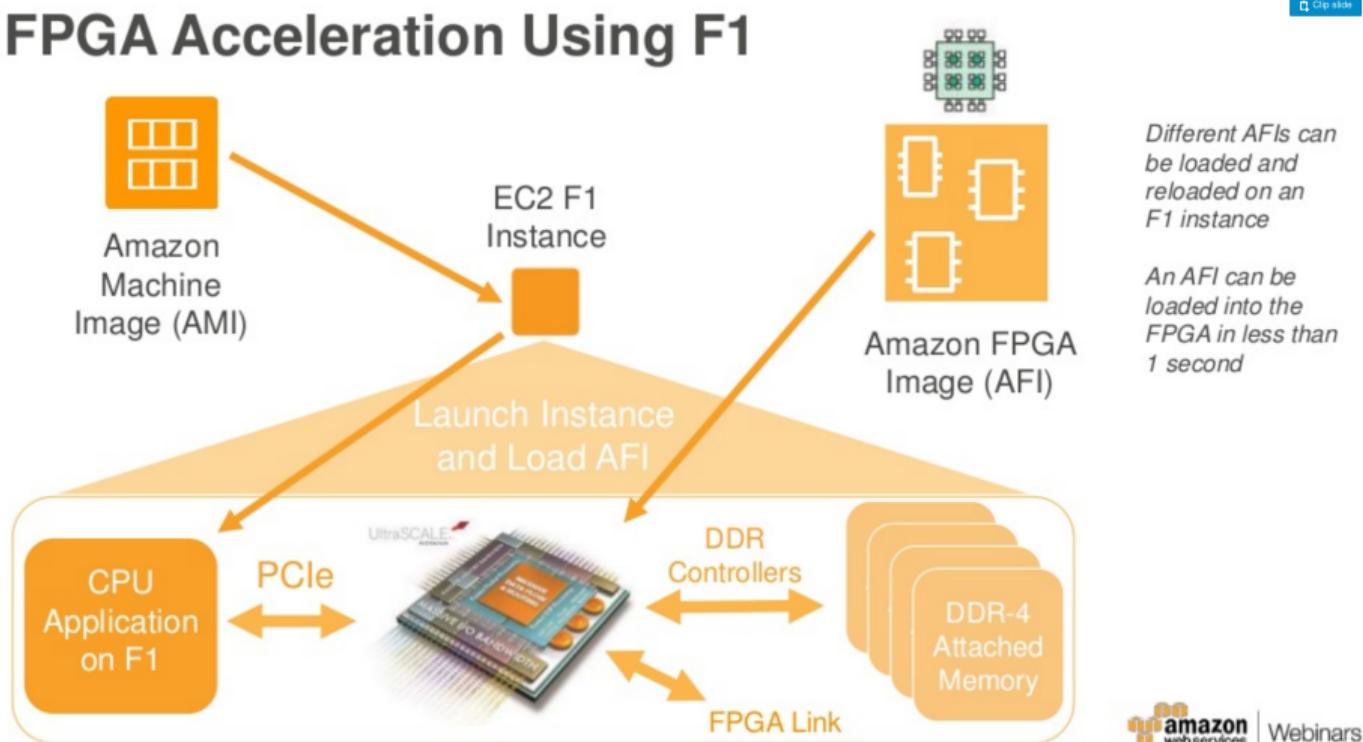


REAL WORLD EXAMPLES: MS CATAPULT

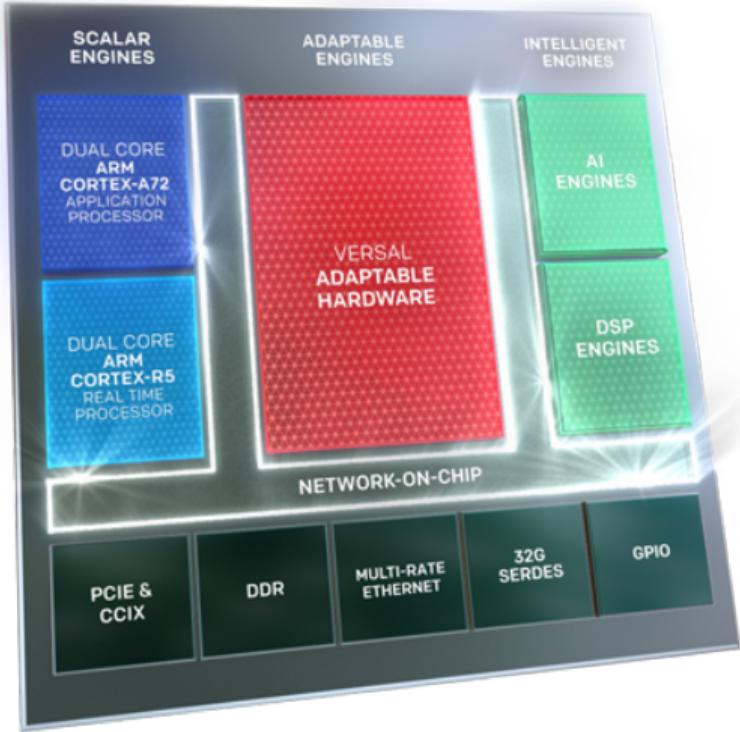


REAL WORLD EXAMPLES: AMAZON F1

FPGA Acceleration Using F1



REAL WORLD EXAMPLES: XILINX VERSAL



REAL WORLD EXAMPLES: XILINX VITIS

Caffe

TensorFlow

FFMPEG



XILINX
VITIS™

FPGA

SoC

ACAP

VALUE PROPOSITION

  Innovation friendly

 Energy efficient

 Versatile

 Fast

