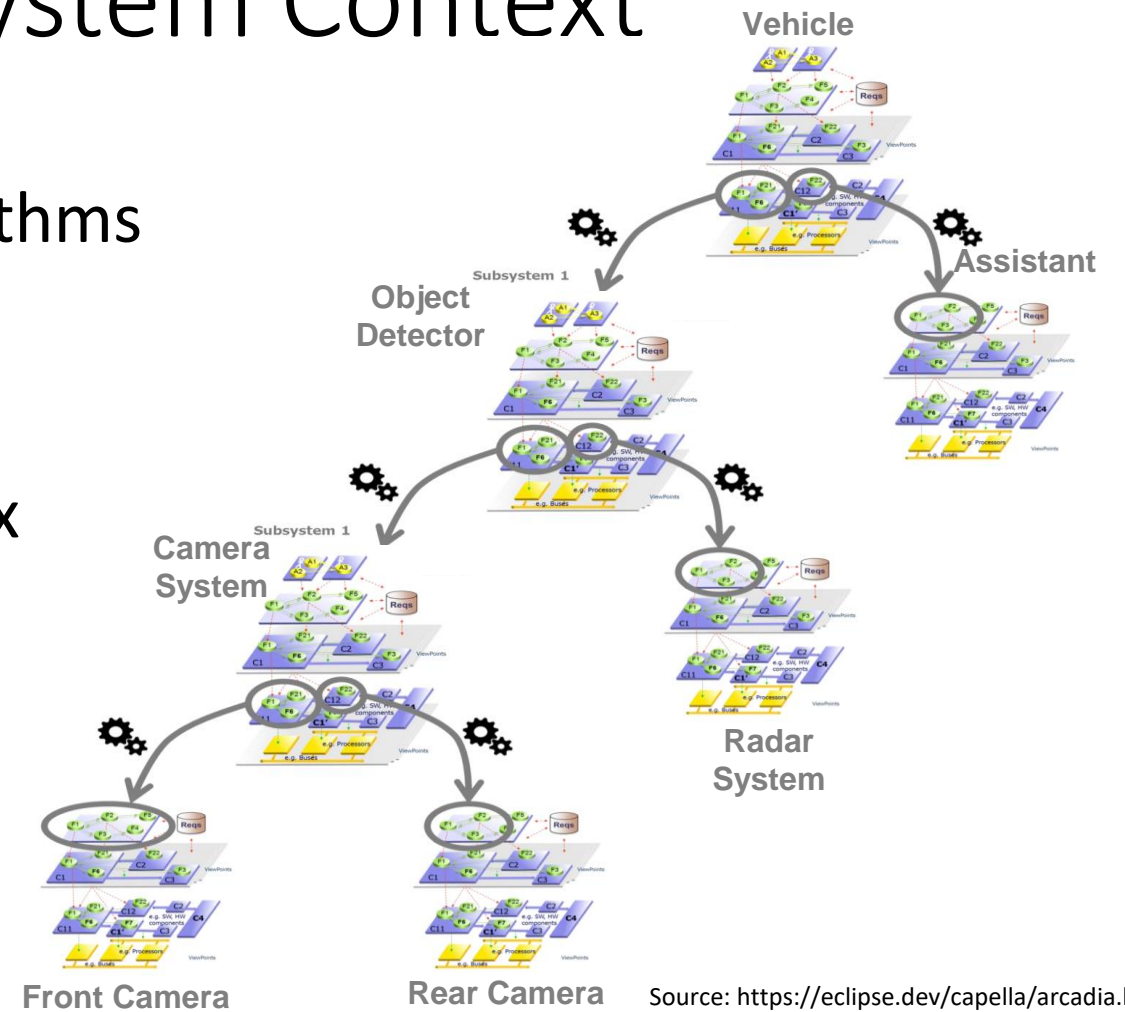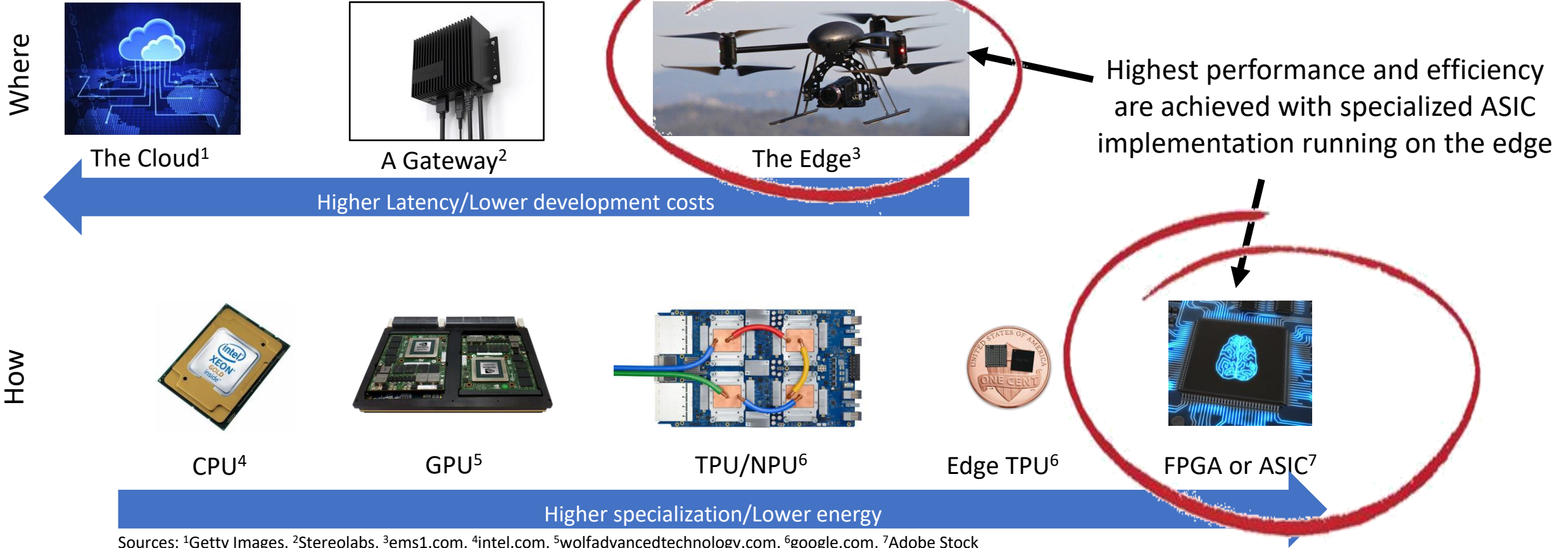# Artificial Intelligence in System Context

One system can have multiple AI algorithms

- Dedicated algorithms for different purposes
- Nested algorithms to provide complex functions, e.g.,
  - Filtering and FFT
  - Object recognition
  - Sensor fusion



Source: https://eclipse.dev/capella/arcadia.html

# Deploying Inferencing Systems, where and how

**Where**



The Cloud[1]



A Gateway[2]



The Edge[3]

Highest performance and efficiency are achieved with specialized ASIC implementation running on the edge

← Higher Latency/Lower development costs

**How**



CPU[4]



GPU[5]



TPU/NPU[6]



Edge TPU[6]



FPGA or ASIC[7]

Higher specialization/Lower energy →

Sources: [1]Getty Images, [2]Stereolabs, [3]ems1.com, [4]intel.com, [5]wolfadvancedtechnology.com, [6]google.com, [7]Adobe Stock

# System Architecture Considerations

**AI algorithms can be implemented in many different ways:**

- Pure software implementation
  - Very flexible and easy to update
  - Performance and timing issues in timing critical applications
- Software with generic hardware accelerator (GPU, NPU)
  - Relies on standard HW
  - Limited flexibility
  - Power consumption and timing issues
- Software with bespoke hardware accelerator
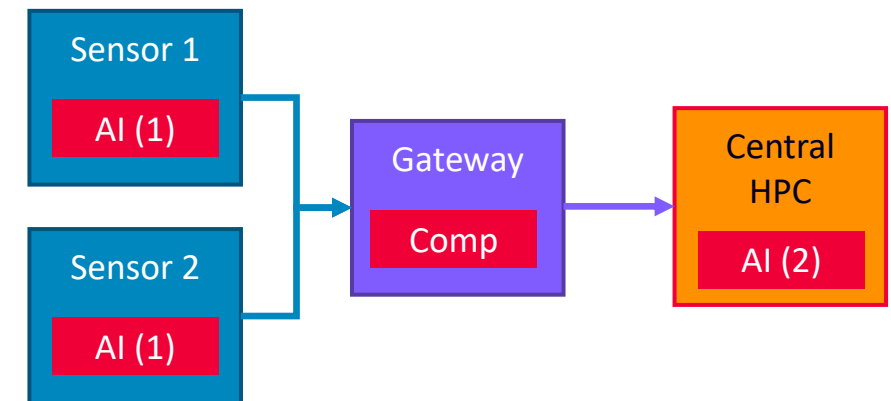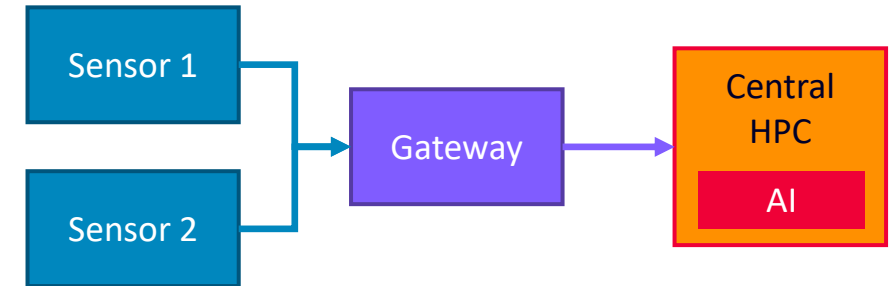  - Requires development of custom HW
  - Low power and predictable timing

# Centralized or Distributed Computation
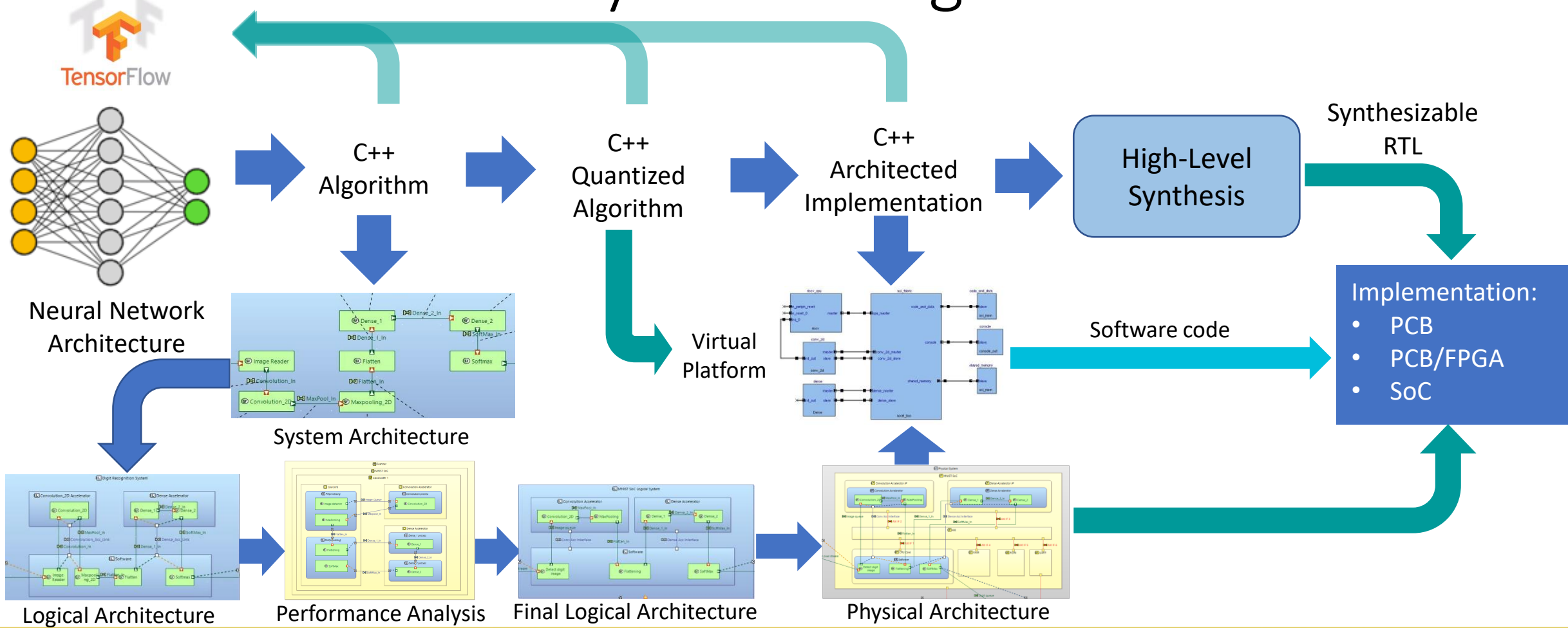
## Centralized computation

- Uncompressed data through network
- High computational load on HPC
- Flexible
- High power consumption

## Distributed computation

- Pre-processing data in its' origin
- Load shared across multiple components
- Data amount reduced by pre-processing and compression
- Low power consumption through dedicated HW

# Model-Based AI System Design

# Model-Based Architecture Exploration

Model-Based Systems Engineering (MBSE)
- Formalized application of modeling to support system design
- Covers design, analysis, verification and validation activities throughout development
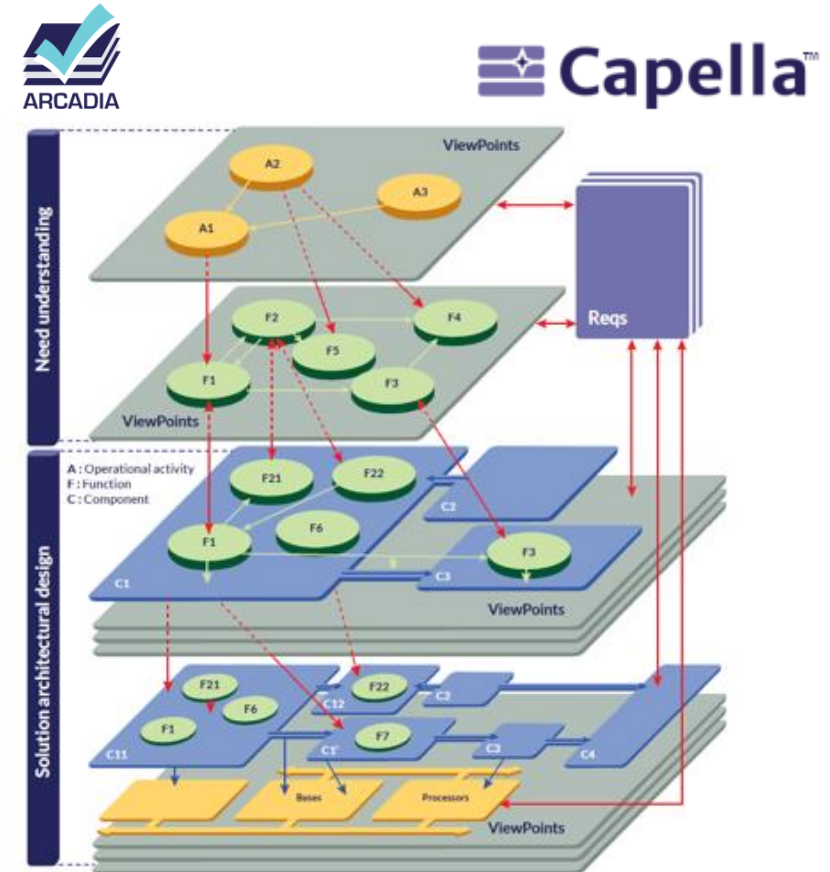
MBSE enables abstract level system architecture exploration
- Functional analysis
- Mapping functions to architectural elements
- Architecture performance analysis
- Subsystem transitioning of system components
- Several formalized approaches available
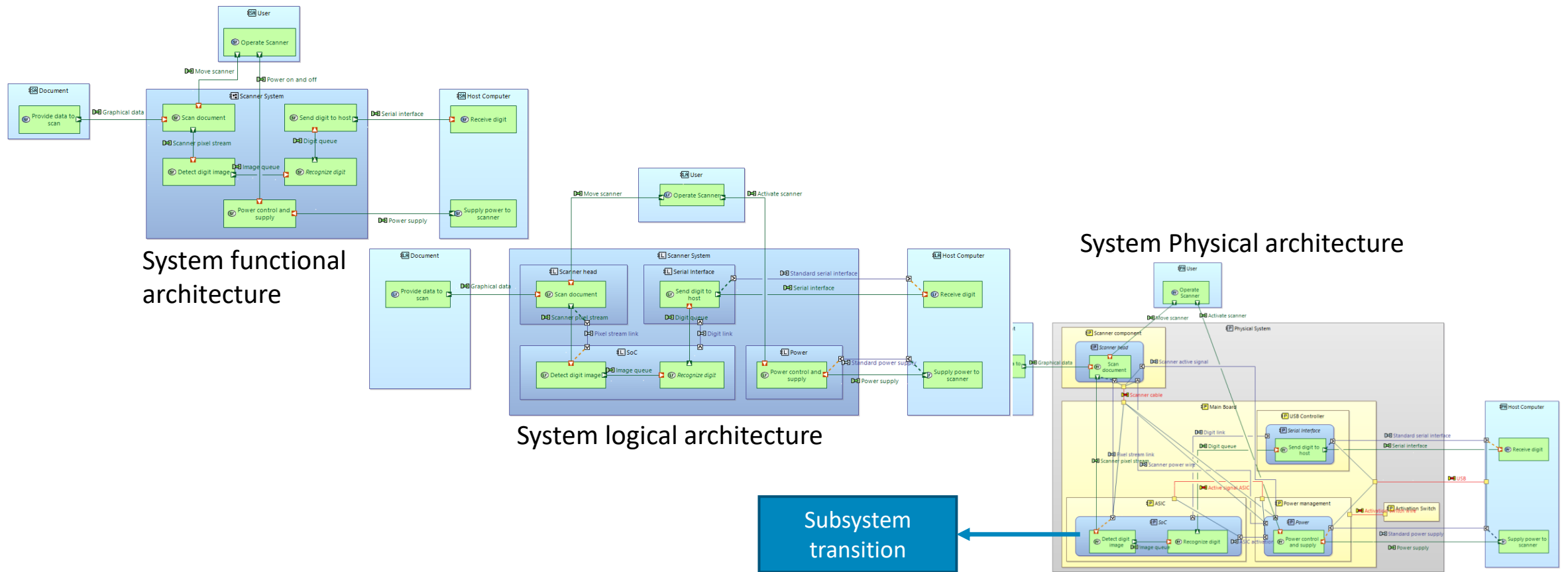
# ARCADIA Methodology

Tooled method to define, analyze, design &
verify system, SW and HW architectures

- Operational Analysis

  - What the users of the system need to accomplish

- System Analysis

  - What the system has to accomplish for the user

- Logical Architecture

  - How the system will work to fulfill expectations

- Physical Architecture
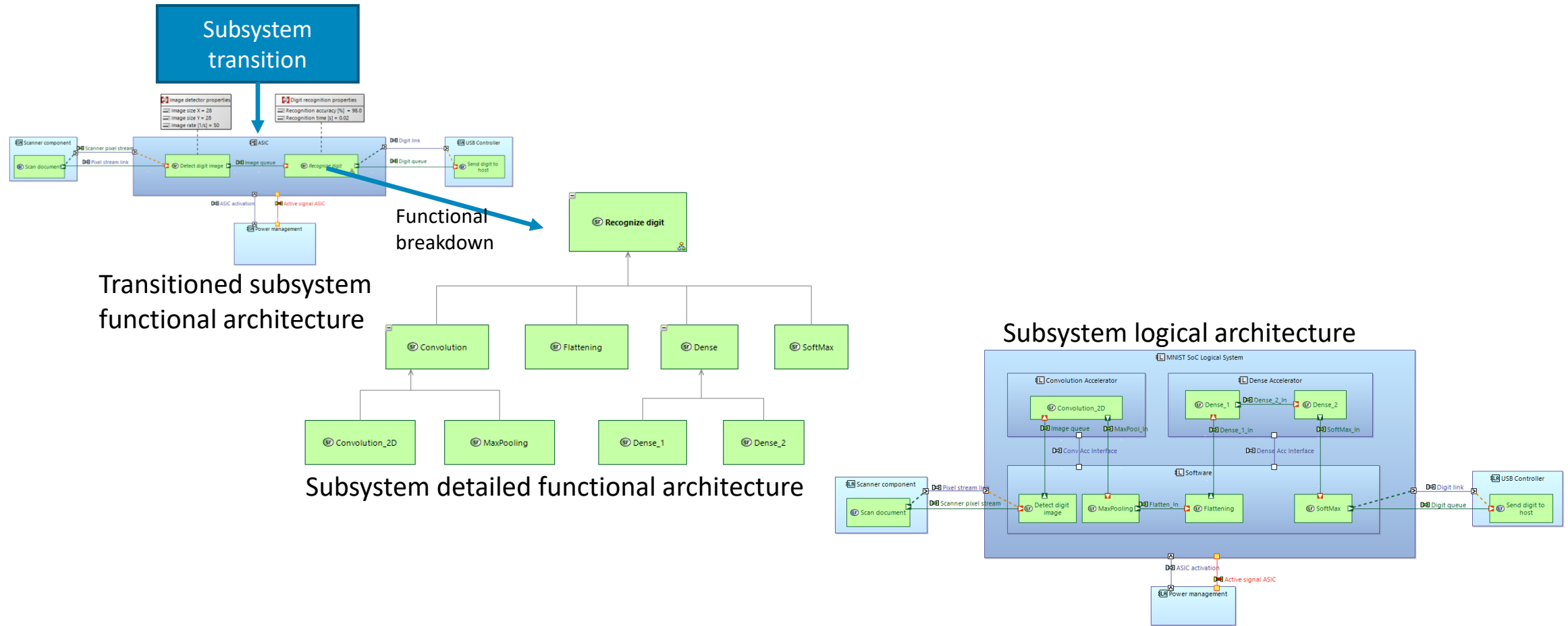
  - How the system will be developed and built
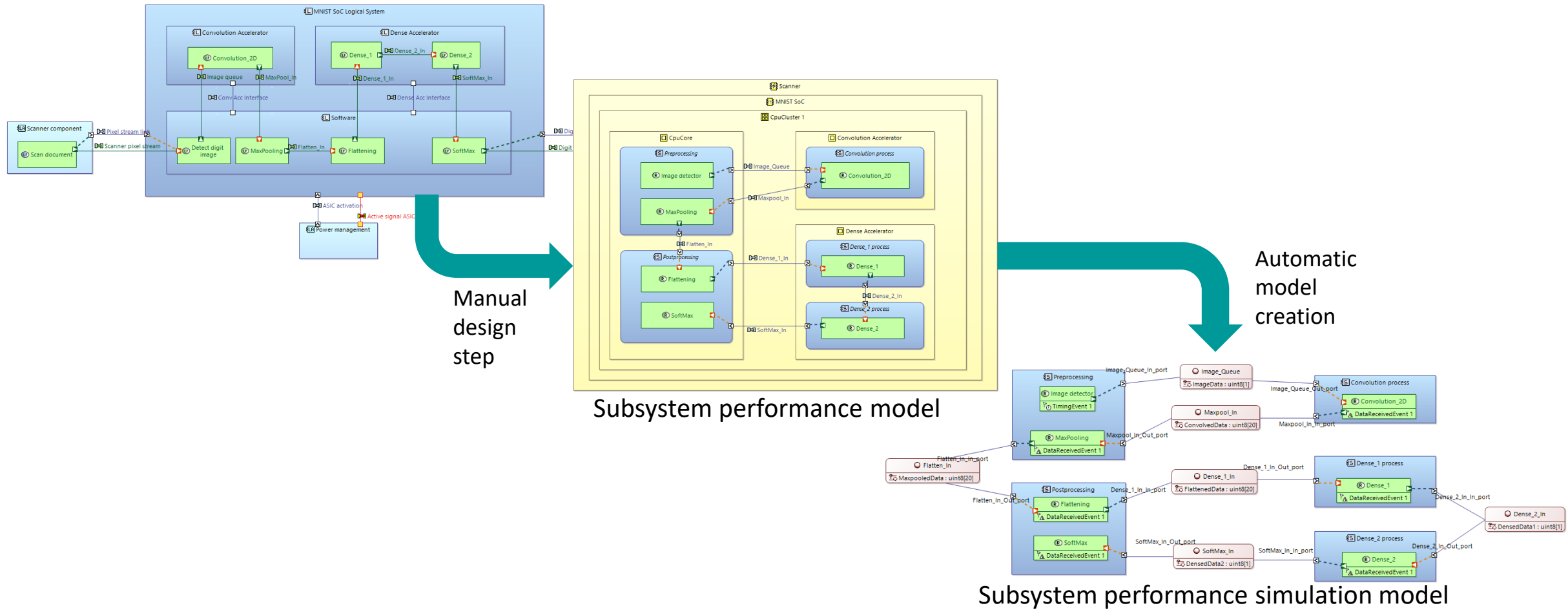


Source: https://eclipse.dev/capella/arcadia.html

# Top-Level System Exploration



System functional architecture

System logical architecture

System Physical architecture

Subsystem transition

# Subsystem Transition and Decomposition



Subsystem transition

Transitioned subsystem functional architecture

Functional breakdown

Subsystem detailed functional architecture

Subsystem logical architecture

# Subsystem Architecture Exploration



Manual design step

Subsystem performance model

Automatic model creation

Subsystem performance simulation model

# Analyzing Performance Simulation Results

Full SW implementation with 3 CPU cores

Accelerated implementation with 2 CPU cores

# Analyzing Performance Simulation Results contd
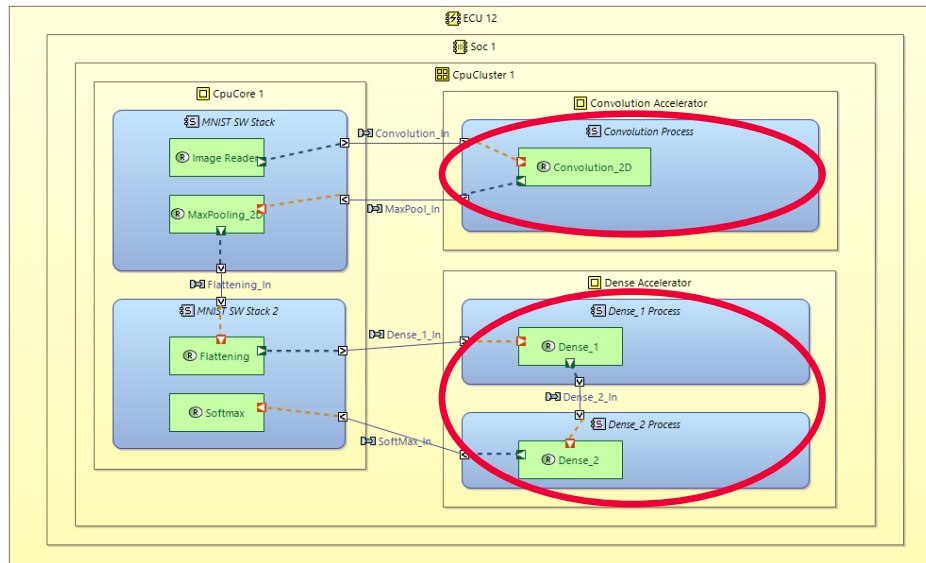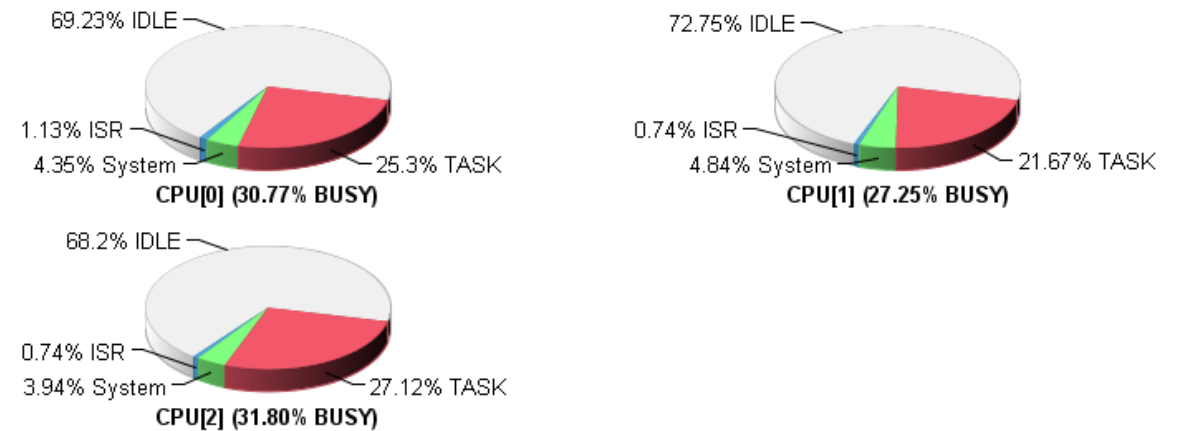
**Architecture with 2 HW accelerators and one CPU core**

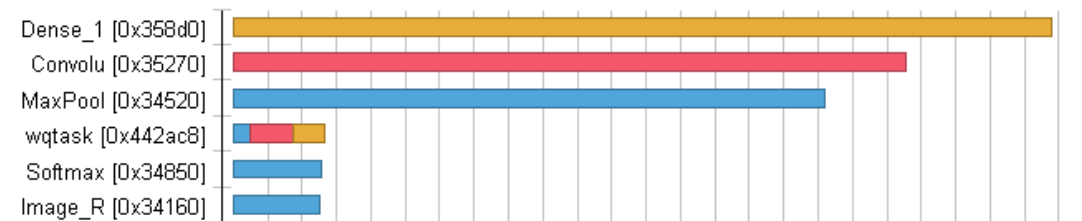- Clock frequency: 500 MHz
- Interval between inferences: 20ms

# Final Subsystem Architecture

Logical architecture is updated based on the simulation results and transitioned to Physical Architecture



Assisted Transition

# HW/SW Co-Architecting

**Iterative multi-abstraction level process**

1. Individual pieces of algorithm developed separately
2. Complete functional model in C/C++
3. Architecture exploration
4. Partitioned HW/SW model
5. Virtual Hardware model with custom accelerators
6. Partitioning optimization
7. High-Level Synthesis of custom accelerators

Continuous verification throughout the process

# Model-Based AI System Design

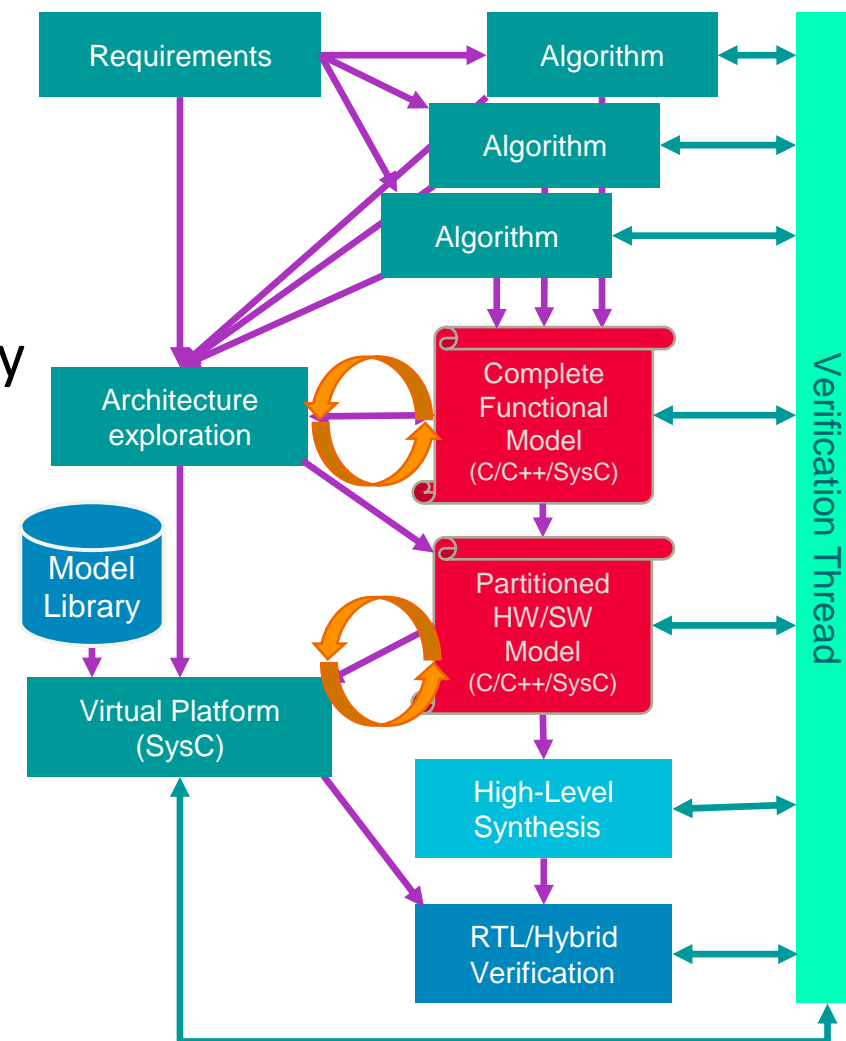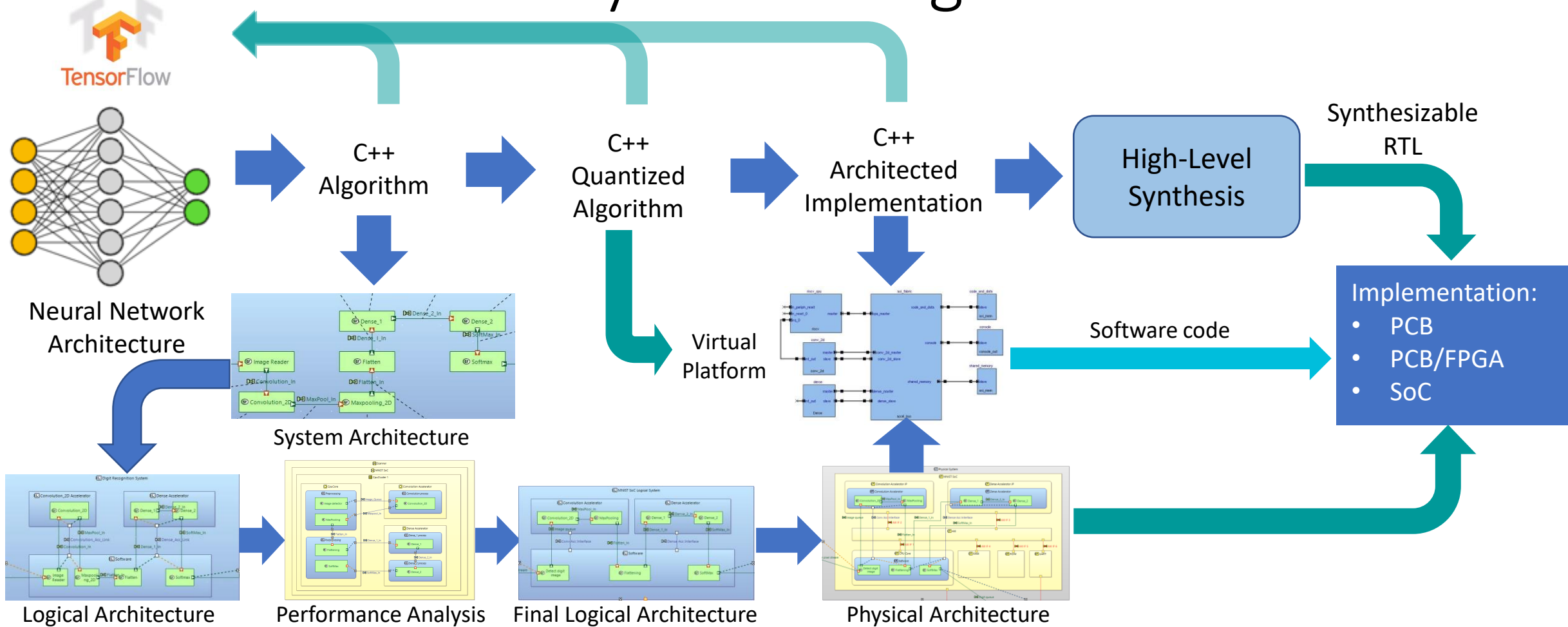# Optimize Neural Network



TensorFlow

Neural Network
Architecture

- In the datacenter, size and speed are secondary considerations

- On the edge we need to be fast and compact
  - Fewer layers
  - Fewer channels

**Original MNIST network**

```
MAC operations:          2357000
Number of parameters:    1966030
Minimum data transfer:   1971244 words
```

**Optimized MNIST network**

```
MAC operations:           235400
Number of parameters:      39690
Minimum data transfer:     42464 words
```

# Create an equivalent C++ model



C++
Algorithm

- To use co-architecting flow we need to convert the algorithm to C++
- Only include programmability in the C++ for parameters you will vary

# Create System (Functional) Architecture



Neural Network Architecture

C++ Algorithm

System Architecture

- Specifies different functions and their dependencies

- Levels of details depends on the system hierarchy level

- Relevant parameters attached to functions as properties

- Functional breakdown needed, when moving down in subsystem hierarchy

# Functional Breakdown and Parameterization



Subsystem detailed parameterized functional architecture

Subsystem transition

Transitioned subsystem functional architecture

Functional breakdown

Subsystem functional breakdown

# Create Initial Logical Architecture



TensorFlow

Neural Network Architecture

C++ Algorithm

System Architecture

Logical Architecture

- Initial system architecture
- Group system functions into logical components based on
  - Performance metrics (MAC operations, profiling data, etc.)
  - Data sizes of the function exchanges
- Allocate function exchanges to component exchanges

# Explore Different Architecture Options



Neural Network Architecture

C++ Algorithm

System Architecture

Logical Architecture

Performance Analysis

- Create initial performance analysis model based on the logical architecture

- Explore different allocations
  - Multi-core
  - Multi-cluster
  - Architectures with hardware accelerators
  - Multi-chip
  - Multiboard
  - …

# Runtime and Data Communication Analysis

## Analyzing performance simulation results
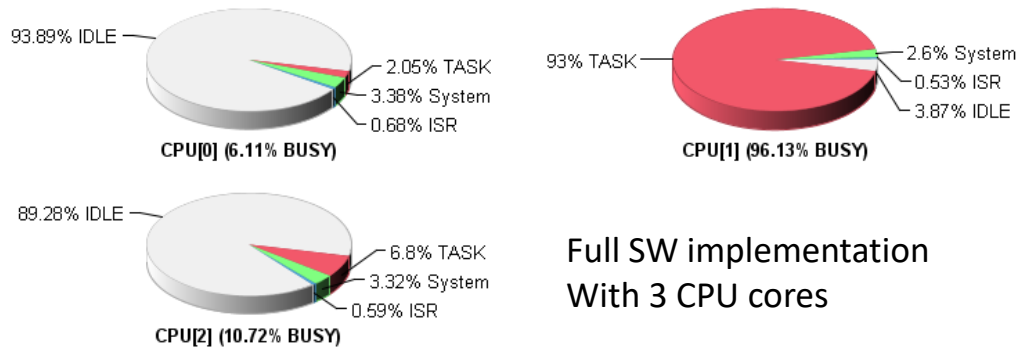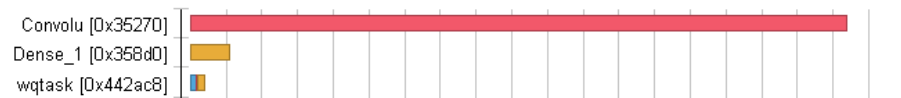


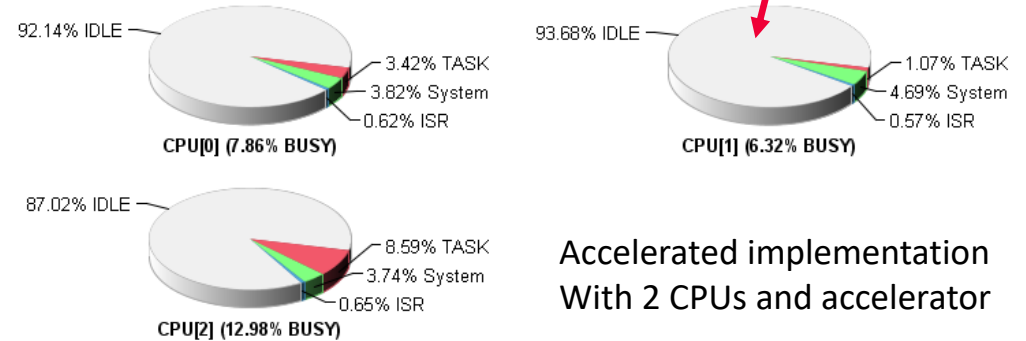**Kernel Object Run-Time Distribution per CPU**

**CPU Execution Run-Time Statistics**

93.89% IDLE — CPU[0] (6.11% BUSY): 2.05% TASK, 3.38% System, 0.68% ISR

93% TASK — CPU[1] (96.13% BUSY): 2.6% System, 0.53% ISR, 3.87% IDLE

89.28% IDLE — CPU[2] (10.72% BUSY): 6.8% TASK, 3.32% System, 0.59% ISR

Full SW implementation
With 3 CPU cores

**Task Execution Statistics**

Convolu [0x35270]
Dense_1 [0x358d0]
wqtask [0x442ac8]

---

**Kernel Object Run-Time Distribution per CPU**

**CPU Execution Run-Time Statistics**

92.14% IDLE — CPU[0] (7.86% BUSY): 3.42% TASK, 3.82% System, 0.62% ISR

93.68% IDLE — CPU[1] (6.32% BUSY): 1.07% TASK, 4.69% System, 0.57% ISR

87.02% IDLE — CPU[2] (12.98% BUSY): 8.59% TASK, 3.74% System, 0.65% ISR

Accelerated function

Accelerated implementation
With 2 CPUs and accelerator

**Task Execution Statistics**

Dense_1 [0x358d0]
wqtask [0x442ac8]
ISR [0x1d]
ISR [0x2]
Convolu [0x35270]
MaxPool [0x34520]

# Quantize HW Functions & Update Logical Arch.



- Update logical architecture based on performance results
- Quantize functions that are allocated to hardware accelerators
- Functional verification of quantized functions

# Quantization of HW Functions

- Optimizing hardware word lengths to minimize HW area
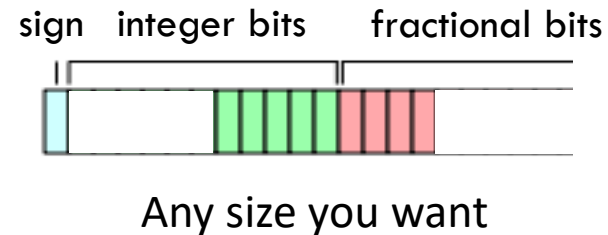  - Ideally every variable individually

- Fixed-point data types ideal
  - ac_fixed
  - sc_(u)fixed

sign   integer bits    fractional bits

Any size you want

- Several analysis methods available
  - Value Range Analysis -based (simulation based)
  - Static Analysis
  - Brute force

# Create Physical Architecture Model



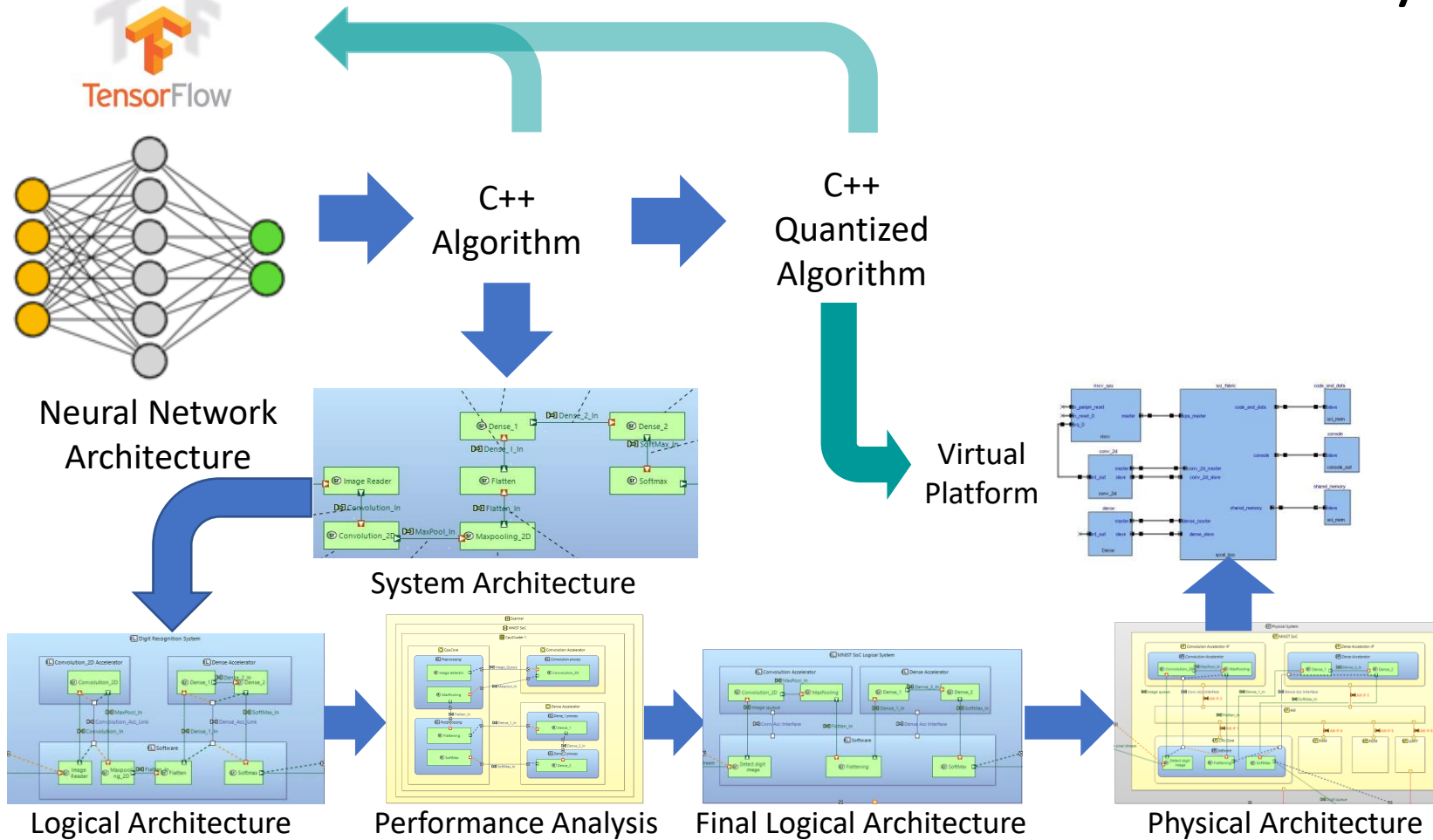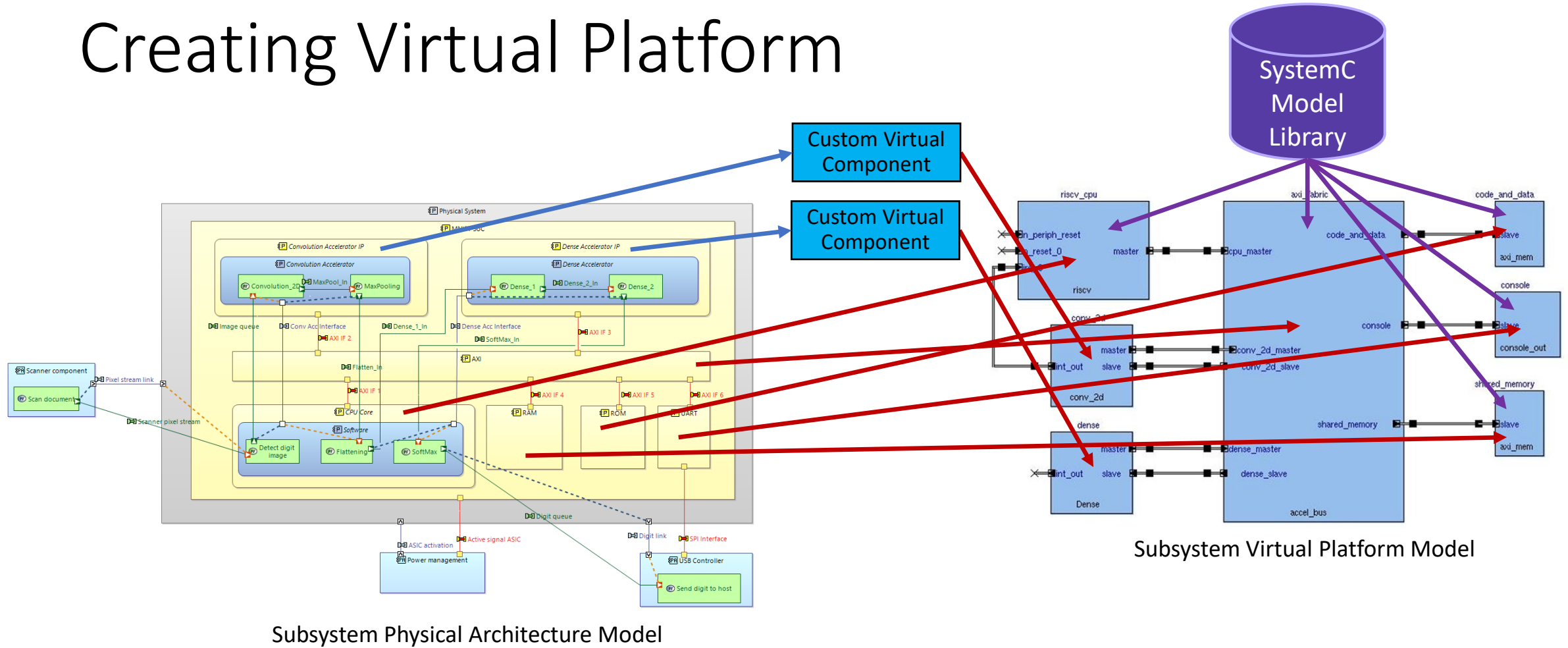- Assisted transition of model
- Non-functional components added manually
  - Memories
  - Interconnects
  - Peripherals

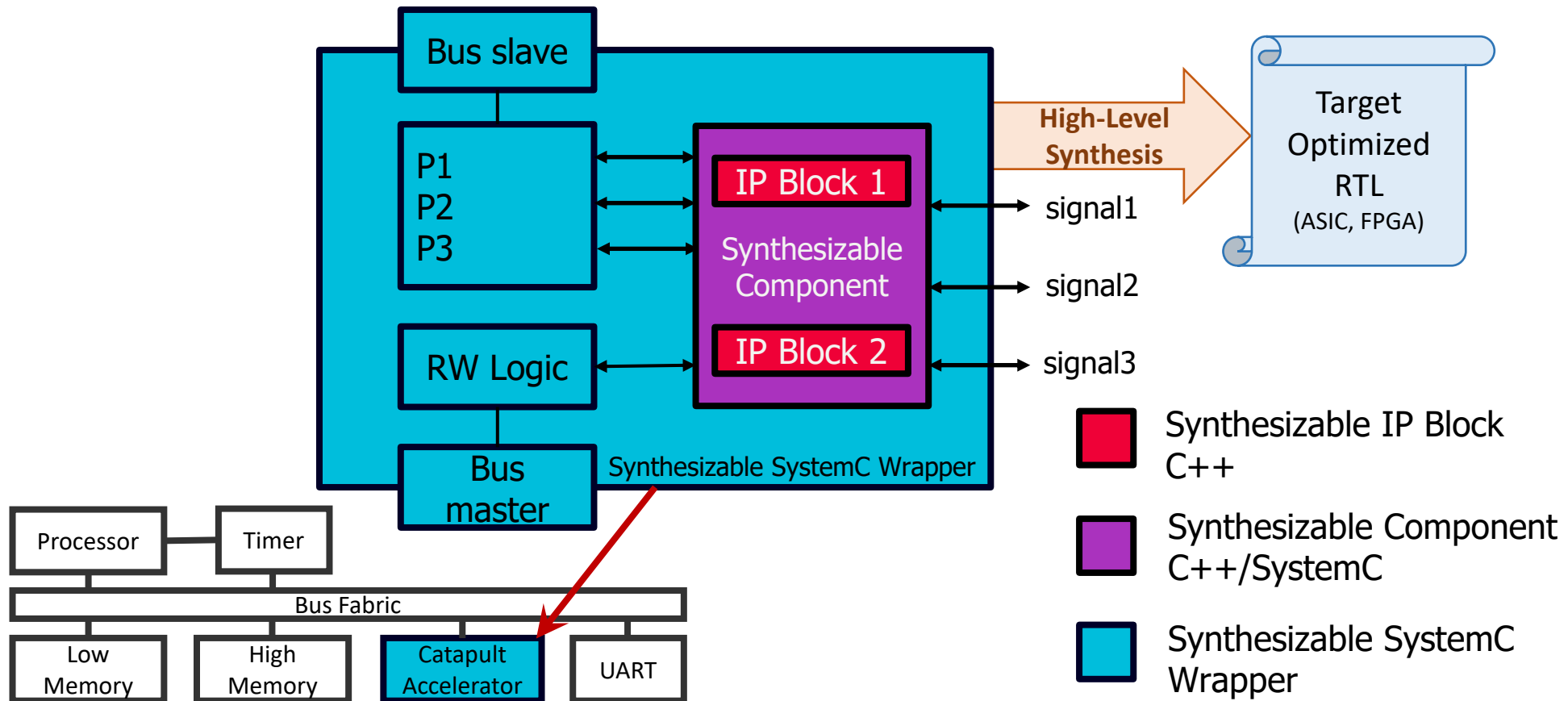# Create Virtual Platform from Physical Architecture



- Create SystemC model of hardware platform
- SW functions mapped to processors
- Initial drivers

# Creating Virtual Platform



Subsystem Physical Architecture Model

Subsystem Virtual Platform Model

# Creating Custom Virtual Component



Bus slave

P1
P2
P3

RW Logic

Bus master

IP Block 1

Synthesizable Component

IP Block 2

signal1

signal2

signal3

Synthesizable SystemC Wrapper

High-Level Synthesis

Target Optimized RTL
(ASIC, FPGA)

Processor    Timer

Bus Fabric

Low Memory    High Memory    Catapult Accelerator    UART

Synthesizable IP Block C++

Synthesizable Component C++/SystemC

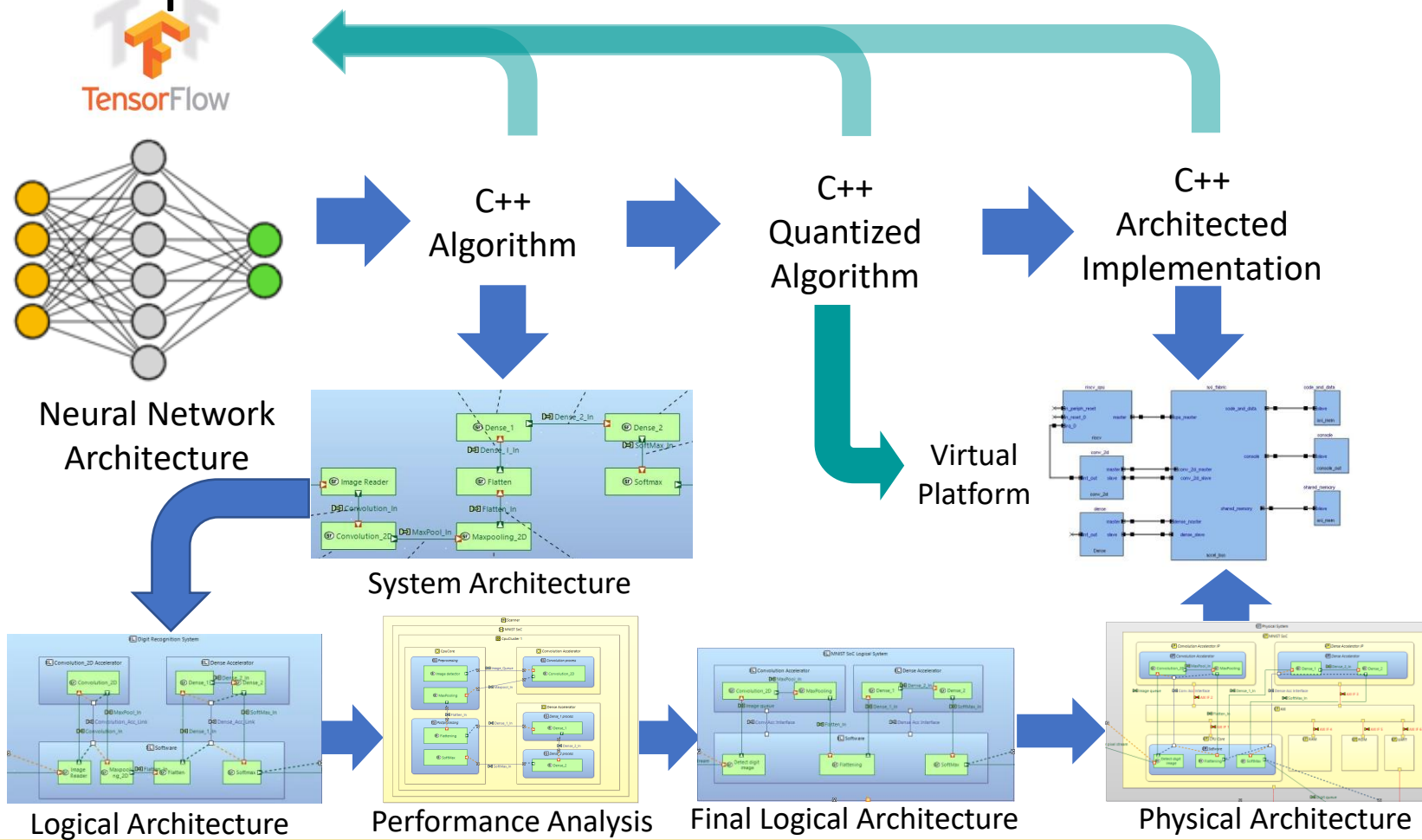Synthesizable SystemC Wrapper

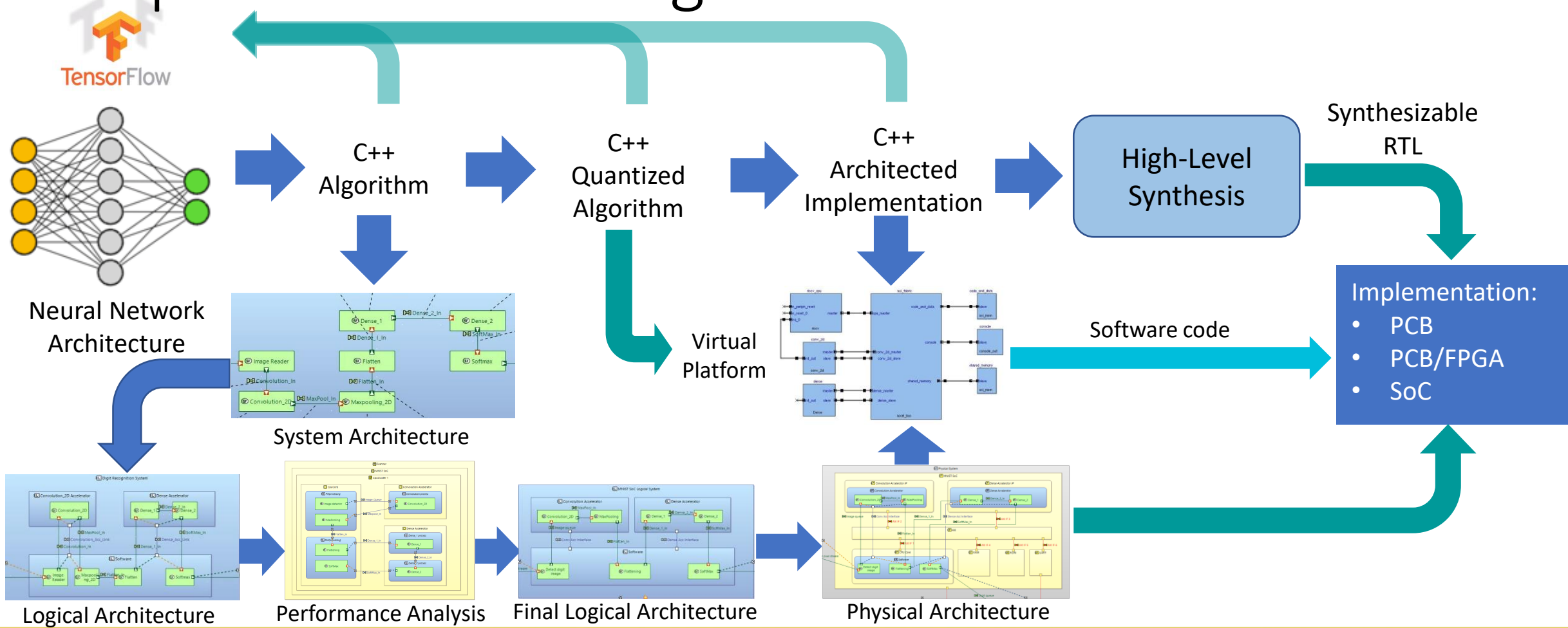# Optimize accelerator model for HLS

- Some code restructuring may be needed for optimal synthesis results
  - Block-level architecture
  - Loop order
  - Internal storages
  - Reusable functions

- Code modifications may influence algorithm accuracy
  - Continuous verification required

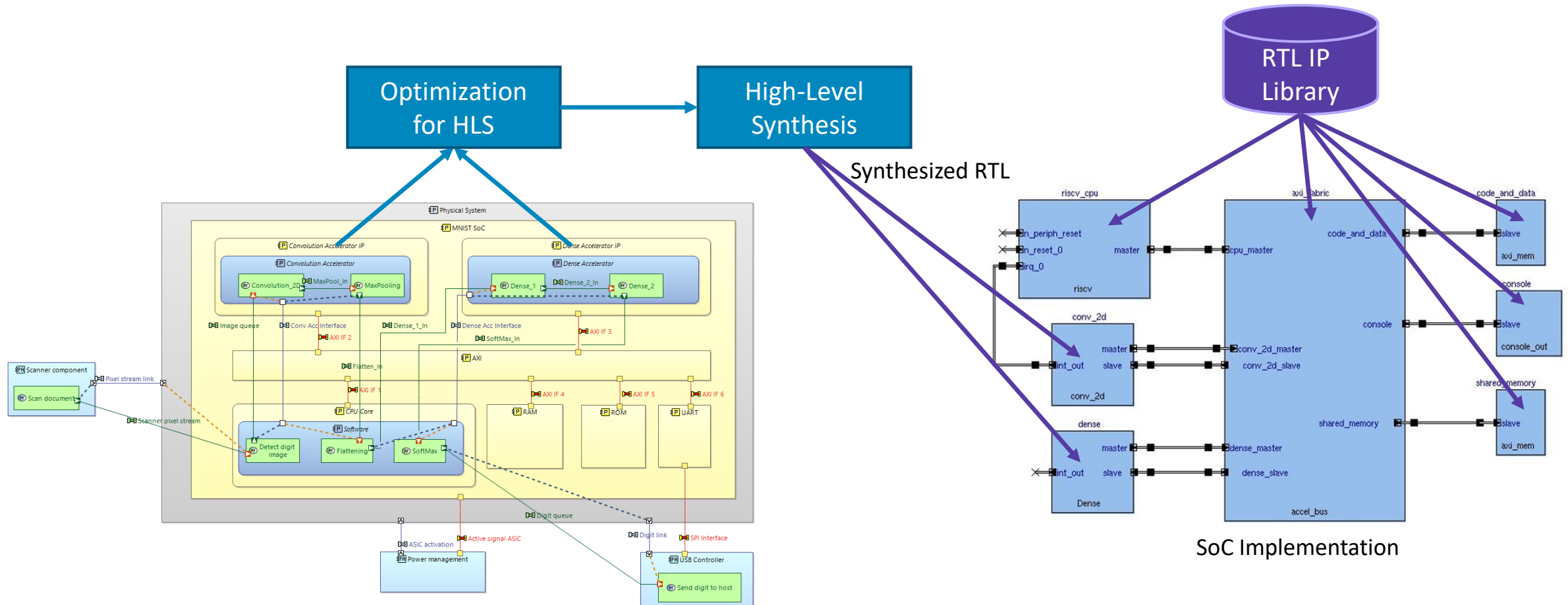# Update Virtual Platform with detailed Accelerator



- No need to change platform architecture
- Bit-accurate function of accelerator
- Final platform for SW development

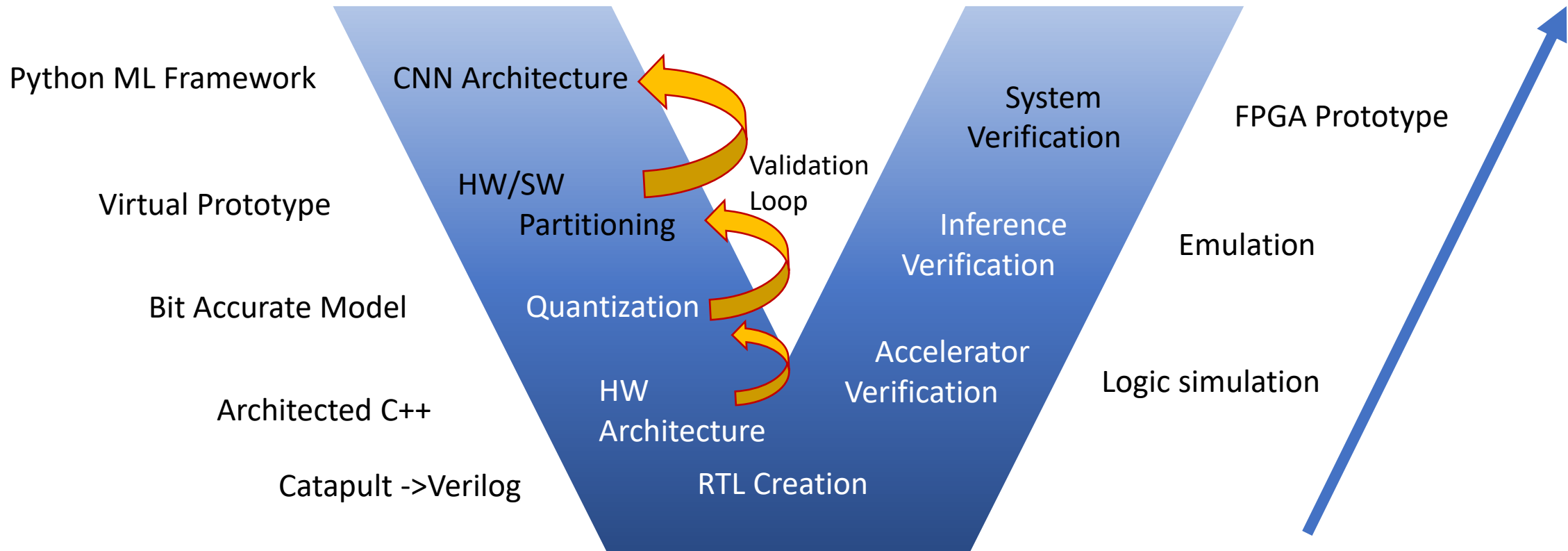# Implement and Integrate HW and SW Modules

# Hardware Implementation



Optimization for HLS

High-Level Synthesis

Synthesized RTL

RTL IP Library

SoC Implementation
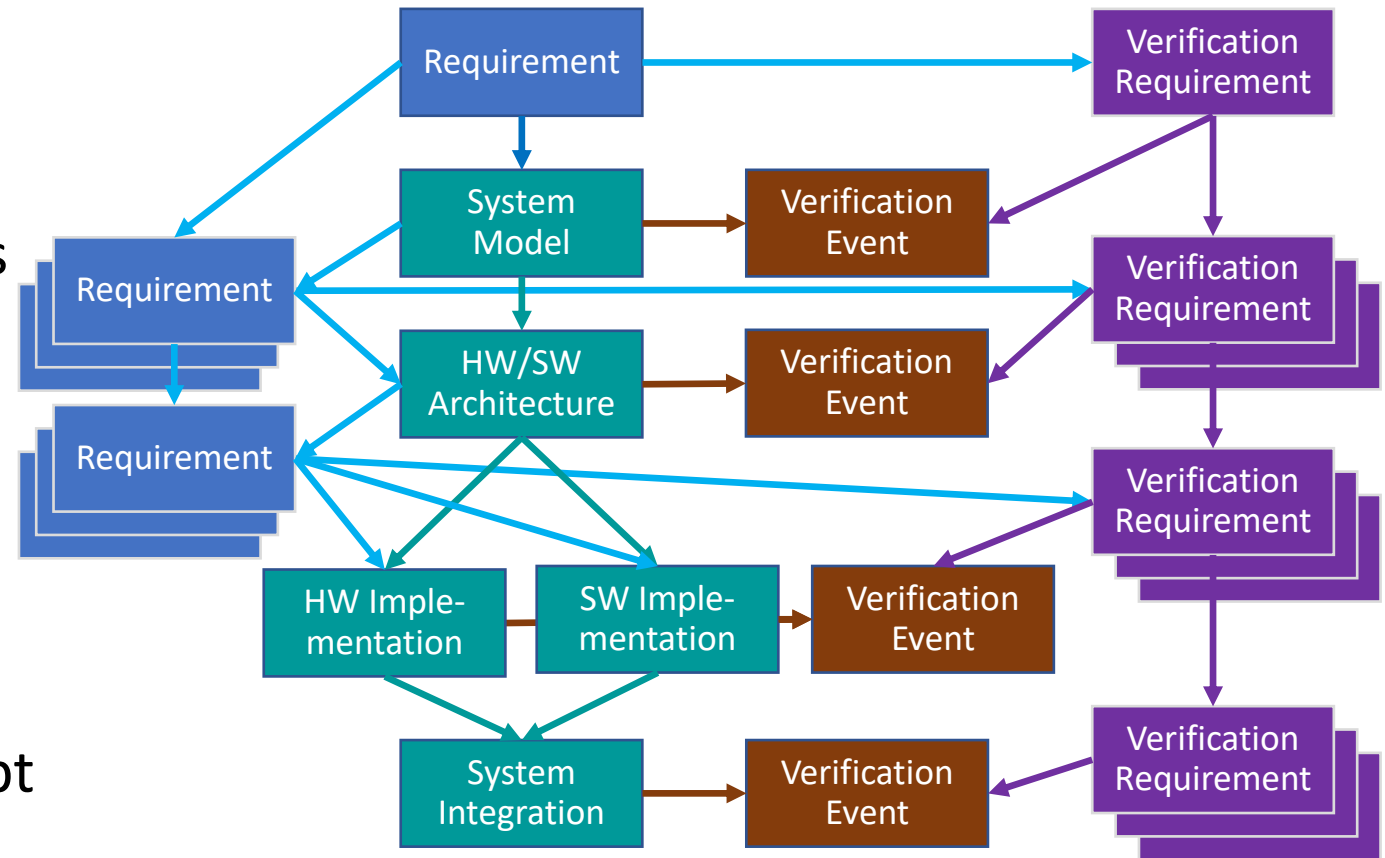
# Validation vs. Verification

- Validation == Are we doing the right thing?
  - Performed during design phase of the project
  - Comparing the model to higher abstraction level or requirement
  - Usually simulation between consecutive design steps
- Verification == Did we implement it correctly?
  - Performed during implementation/integration phase of the project
  - Implementation vs. requirements at the same level
  - Implementation vs. model at the same level
  - Test coverage, corner cases, etc.
  - Simulations, formal analysis, emulation, prototyping

# Validation and Verification

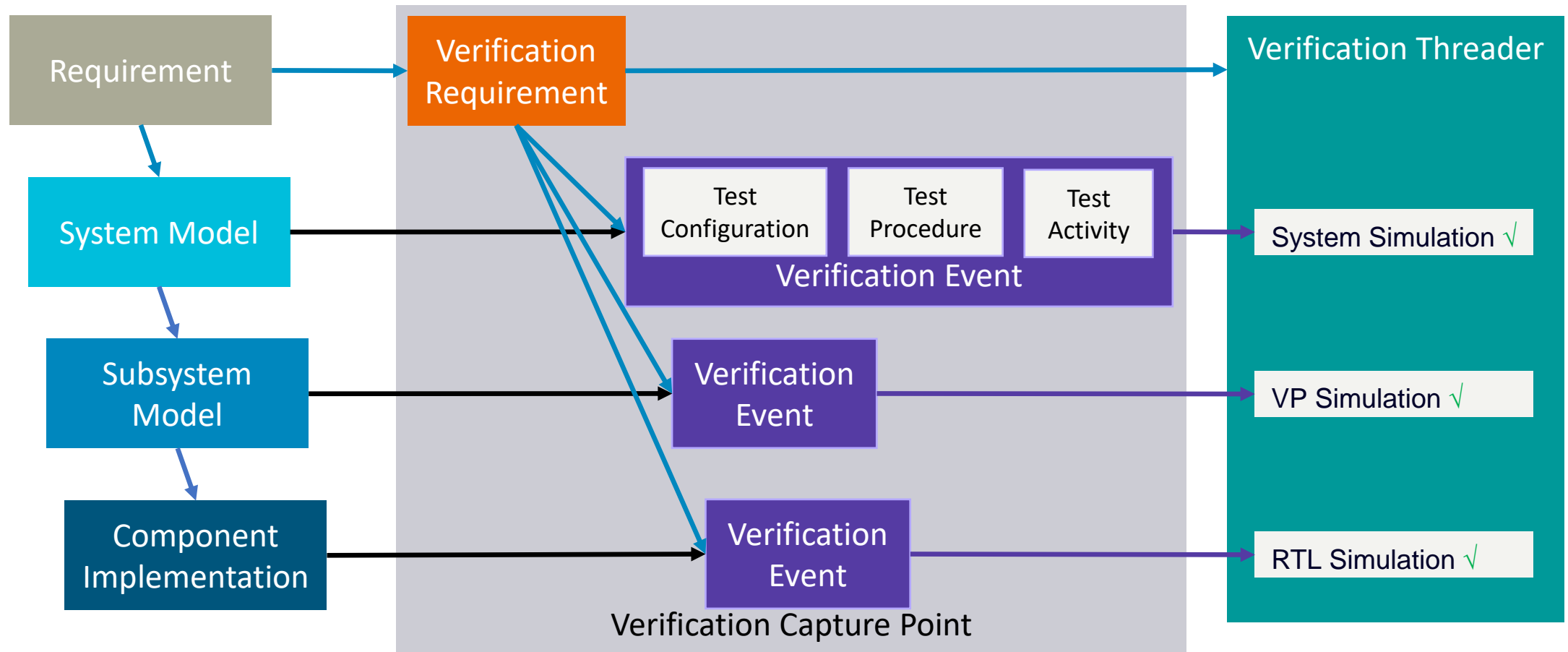# Multi-Abstraction-Level Verification Challenges

- Requirements driven verification with
  - Parameterizable requirements
  - Verification requirements
  - Requirement refinement
  - Hierarchical requirements

- Continuous verification and requirements tracing
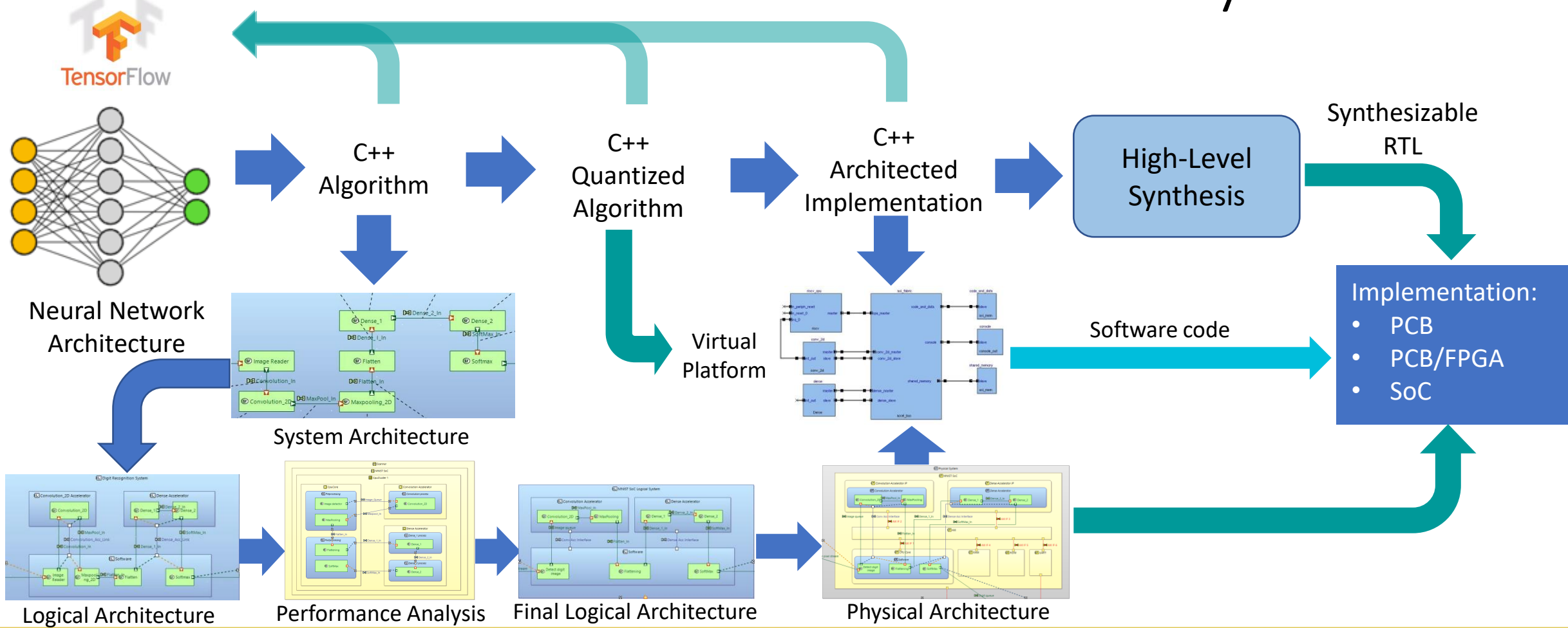  - Multi-layer verification concept

# Verification Process in Model-Based Design

- Requirements driven process:
  - Verification requirement defines test event
    - Test procedure
    - Test activity
    - Test configuration
  - Refined and hierarchical requirements need their own test events
- Verification Capture Point (VCP)
  - Bundles all test events related to one test requirement together
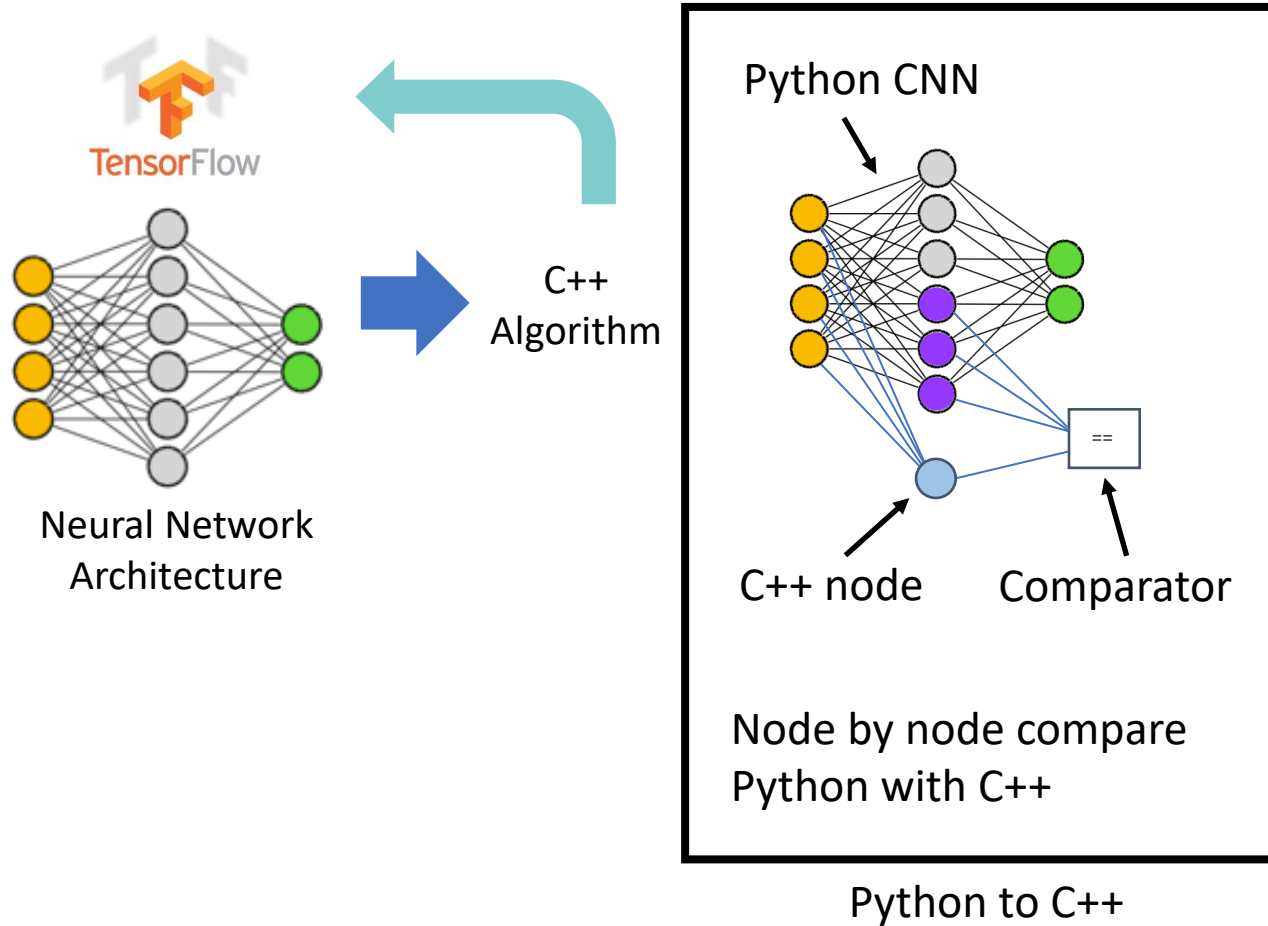  - Contains test events in different design phases and abstraction levels

# Verification Capture Point

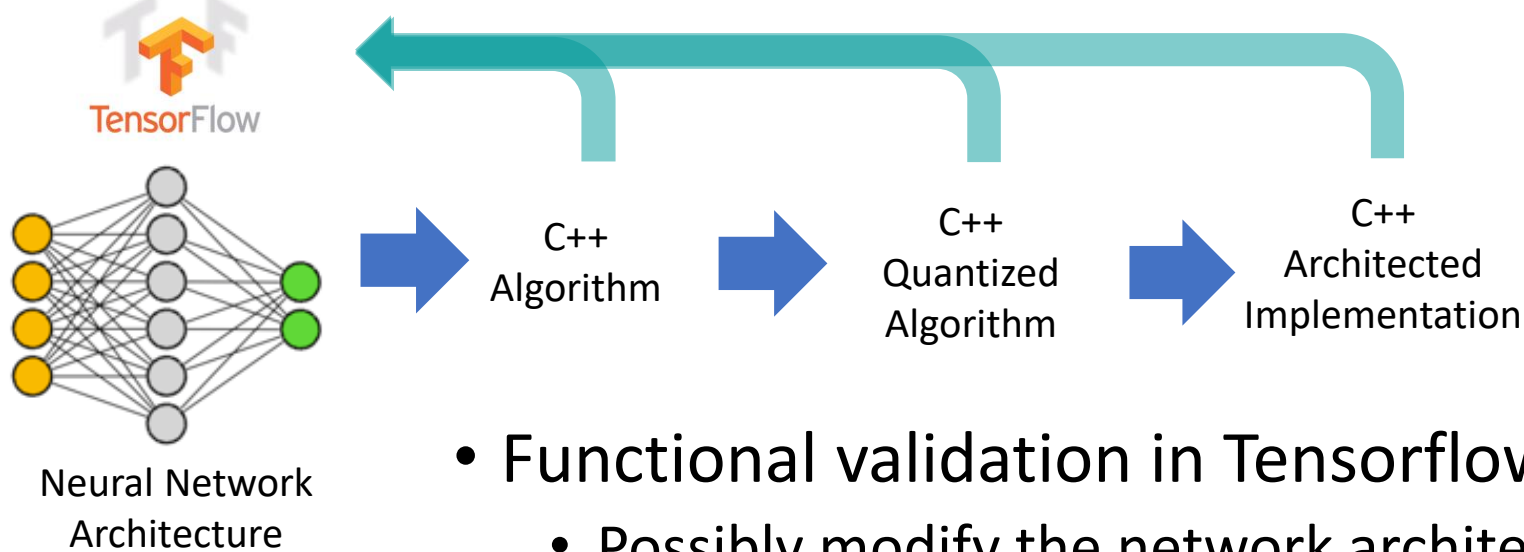# Validation and Verification of MNIST System

# Validating Python to C++ Translation Consistency



Neural Network Architecture

C++ Algorithm

Python CNN

C++ node          Comparator
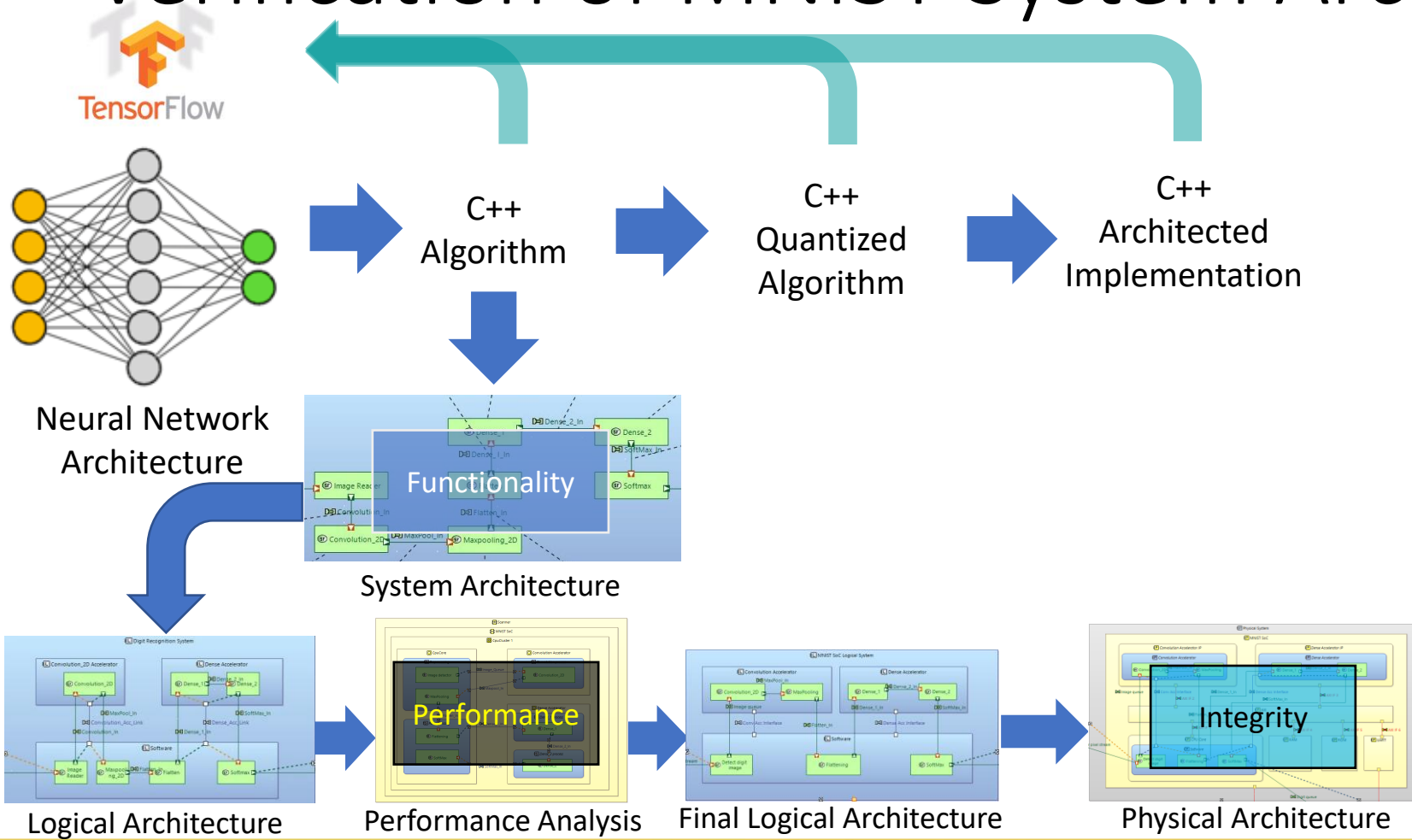
Node by node compare Python with C++

Python to C++

- Import a C++ node into Python and compare the outputs
- Start with one node
- Then one layer
- Then the whole network
- Should match Python results except for bottom 2 or 3 bits
- Variance results from order of operations
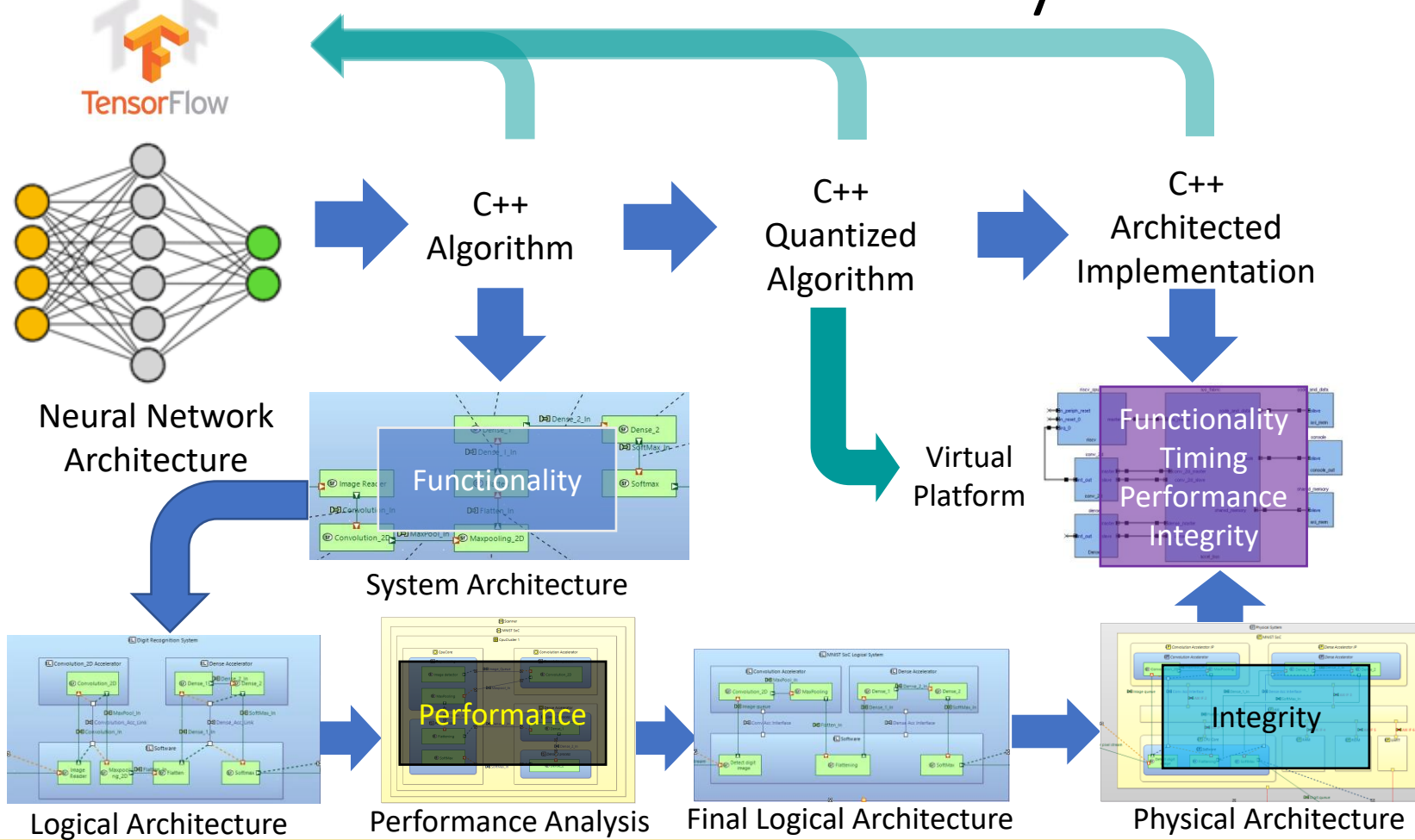
# Validation and Verification of C++ Code



Neural Network Architecture

C++ Algorithm → C++ Quantized Algorithm → C++ Architected Implementation

- Functional validation in Tensorflow
  - Possibly modify the network architecture for better PPA
- Verification of quantized and architected C++ code
  - Ensure that implementation meets requirements
  - In C++ domain with translated testbench
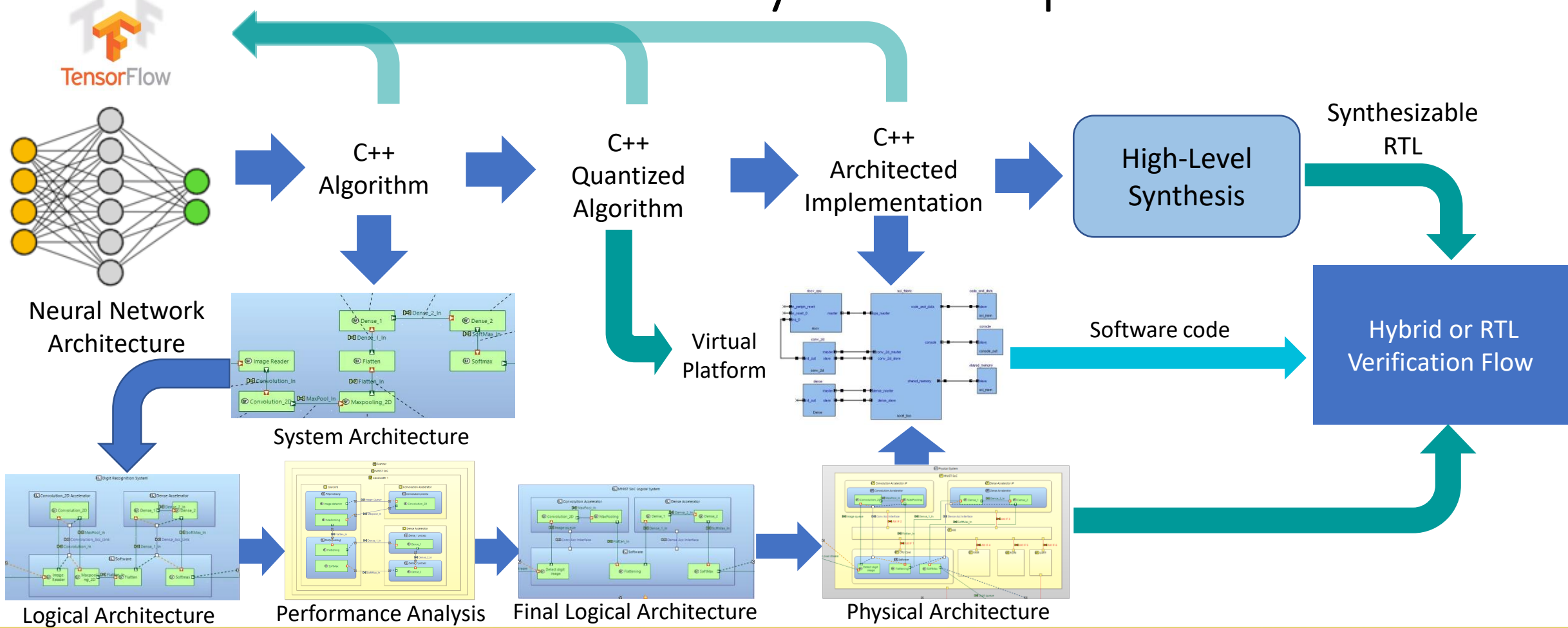  - Coverage analysis needed to ensure test quality

# Verification of MNIST System Architecture

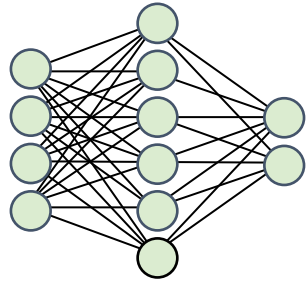# Verification of MNIST System Architecture

# Verification of MNIST System Implementation

# Verification – Before HLS



C++ Architected CNN
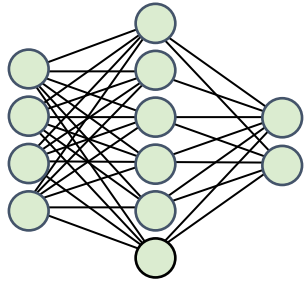
## Static Design Checks

Static code analysis and synthesis checks. Find coding errors and problem constructs

## Coverage Analysis

Determine completeness of test cases. Statement, branch and expression coverage as well as covergroups, coverpoints, bins and crosses

# C++ to RTL consistency



C++ Architected CNN

**Formal**

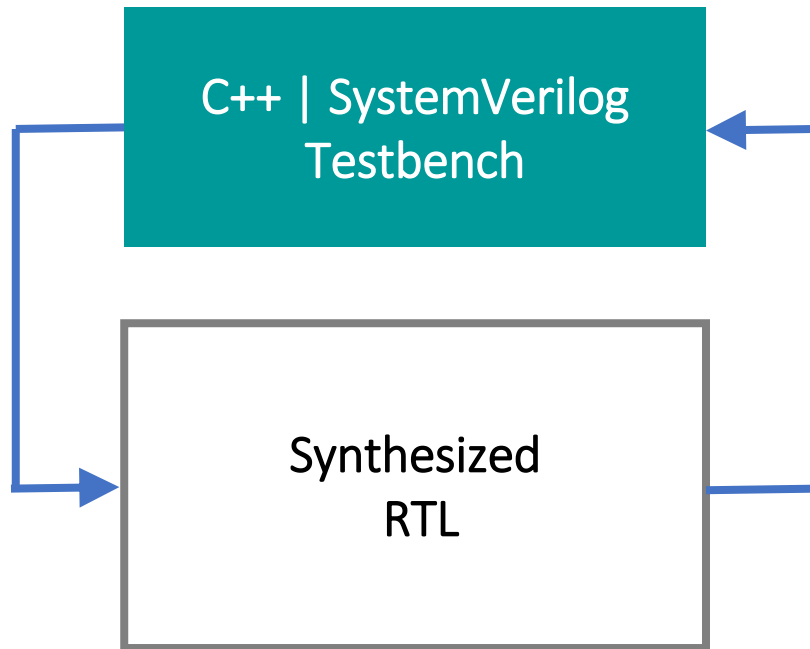Using formal techniques, prove as much equivalency as possible

**UVM**

Architected C++ is used as a predictor for RTL verification

**RTL Coverage**

Determine remaining verification effectiveness through RTL coverage metrics
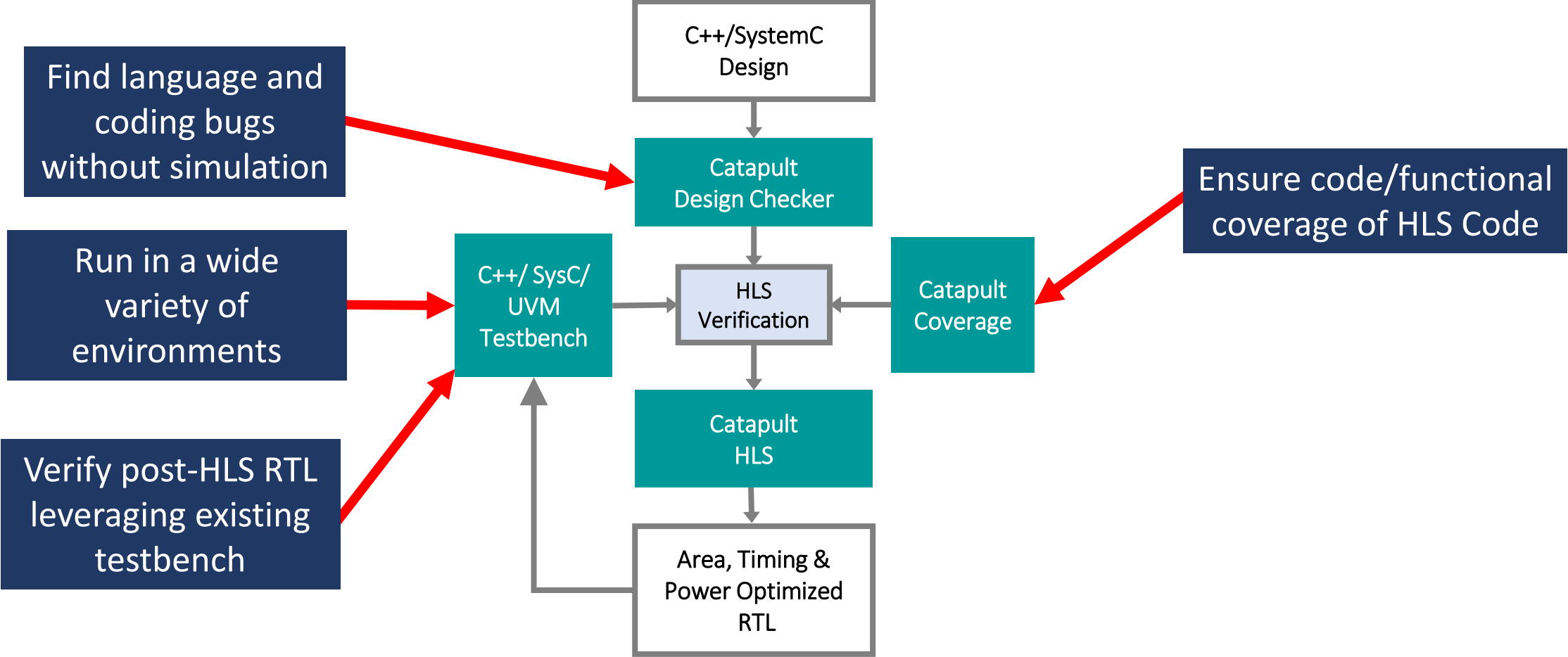
```
44
45 always @(posedge clk_i or negedge rstn_i) begin
46 if(!rstn_i)  begin
47    ring_cnt <= 2'b00 ;
48    grnt_o <= 4'b0;
49 end
50 else begin
51    if(req_i == 4'b000)
52       grnt_o <= 4'b0000;
53
54    if(timer_start)  begin
55
56       if(req_i[3:0] == 4'b0001)
57          grnt_o <= 4'b0001;
58       else if (req_i == 4'b0010)
59          grnt_o <= 4'b0010;
60       else if (req_i == 4'b0100)
61          grnt_o <= 4'b0100;
62       else if (req_i == 4'b1000)
63          grnt_o <= 4'b1000;
64       else begin
65
66          if (timer_exp)
67             ring_cnt <= ring_cnt +1 ;
68
69
70          if(ring_cnt == 2'b00) begin
71             if(req_i[0])
72                grnt_o <= 4'b0001;
73             else
74                ring_cnt <= ring_cnt +1 ;
75          end
76
77          if(ring_cnt == 2'b01) begin
78             if(req_i[1])
79                grnt_o <= 4'b0010;
80             else
81                ring_cnt <= ring_cnt +1 ;
82          end
83
84          if(ring_cnt == 2'b10) begin
85             if(req_i[2])
86                grnt_o <= 4'b0100;
87             else
88                ring_cnt <= ring_cnt +1 ;
89          end
90
```
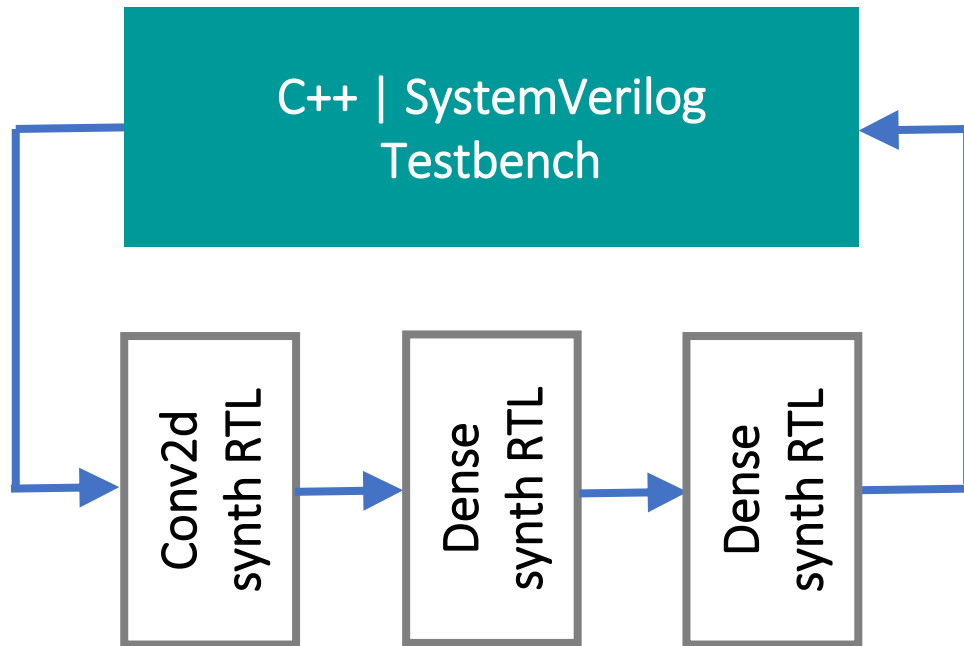
# Block Level Verification



- Re-use C++ or employ new System Verilog testbench
- Prove correctness
  - Cover corner and exception conditions
- Repeat for each accelerator

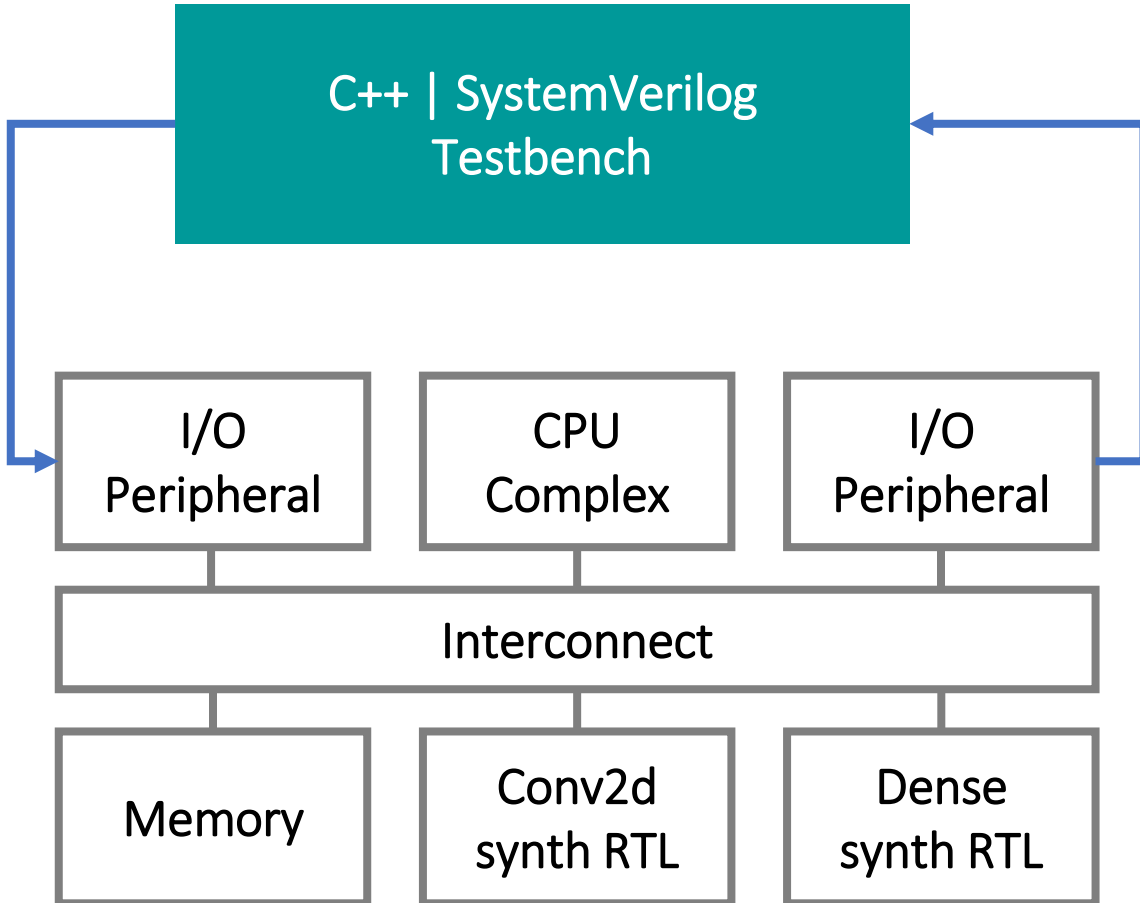# Block Level Verification - HLS Flow Tools
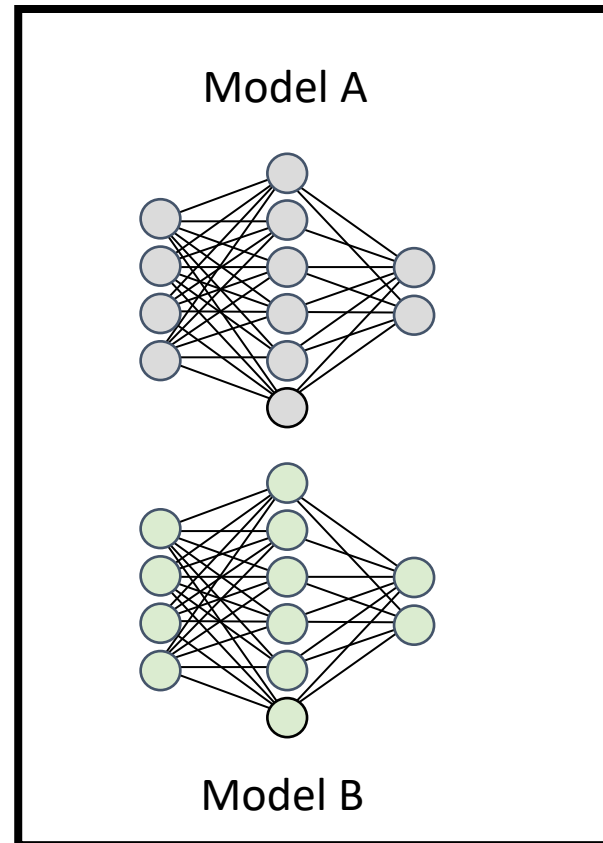
# Sub-System Verification



- Prove the correctness of a collection of accelerator

- Exercise larger functions
  - In this case inference
  - Low level coverage is not important here

- Run times could be impractical for logic simulation
  - May require acceleration (emulation | FPGA prototype)

# System Verification



- Includes processors, software, interconnect
  - Exercises HW and SW interfaces
- Execute typical and exceptional use cases
- Testbench drives I/O, clock, and reset
  - Software and processor orchestrate operation (as in final system)
- Likely to require FPGA prototype

# Debugging – When Things Go Wrong


Model A

Model B

- Log all intermediate values to memory or log file
  - This includes output from each layer

- Have scripts that can compare intermediate values from different model representations
  - This identifies the first point of divergence between models
  - Immediately find layer and node where problem resides

- Intermediate values from the Python can be recorded to a file for comparison

# Questions?

# Model-Based Approach for Developing Optimal HW/SW Architectures for AI systems

Petri Solanti, petri.solanti@siemens.com

Russell Klein, russell.klein@siemens.com