

Next Generation Verification for the Era of AI/ML and 5G

Frank Schirrmeister, Pete Hardee, Larry Melling,
Amit Dua, Moshik Rubin

Cadence Design Systems, Inc.

 cadence®

Agenda

- The era of 5G and AI/ML
 - Challenges
 - What they mean for verification
 - Typical Designs
- Verification options to *enable* 5G/AI/ML
- AI/ML Inside/Outside for Verification
- Summary - Outlook

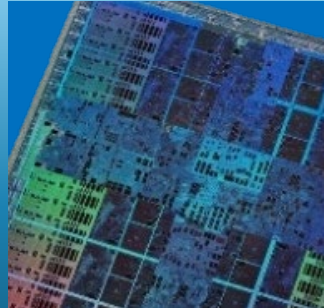
THE ERA OF 5G AND AI/ML

Electronics Innovation

2020



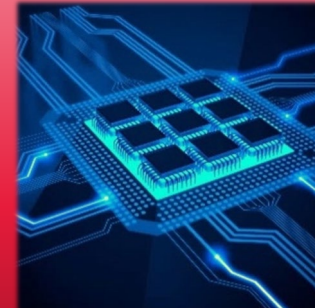
Edge-to-Cloud
and 5G



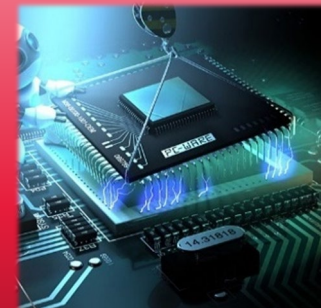
Specialized
SoCs



Artificial
Intelligence



Parallel and
Distributed

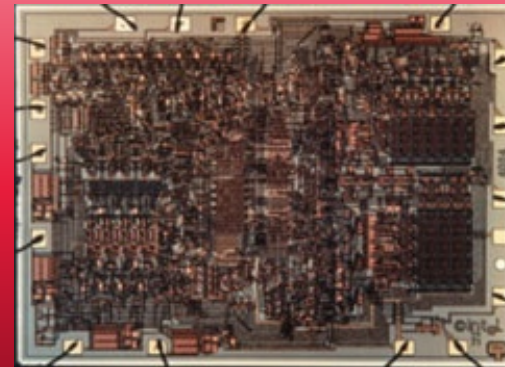
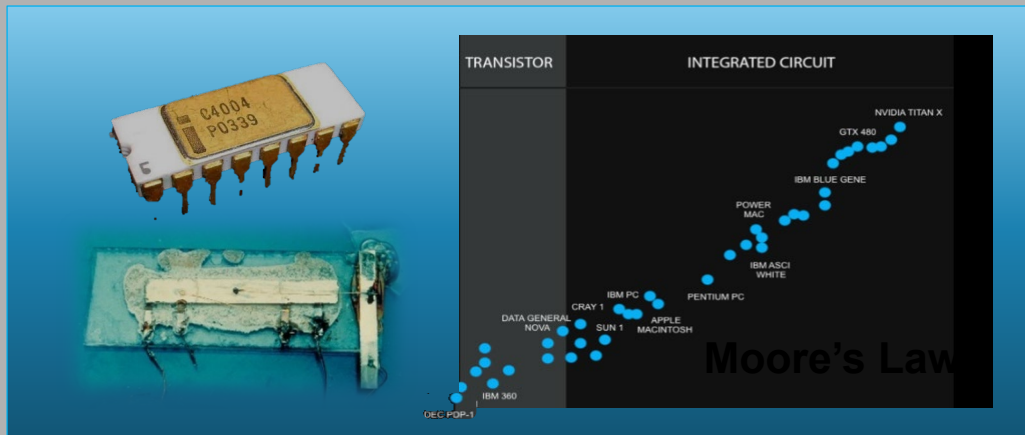


Full-flow
Solutions



Deep
Learning

1960

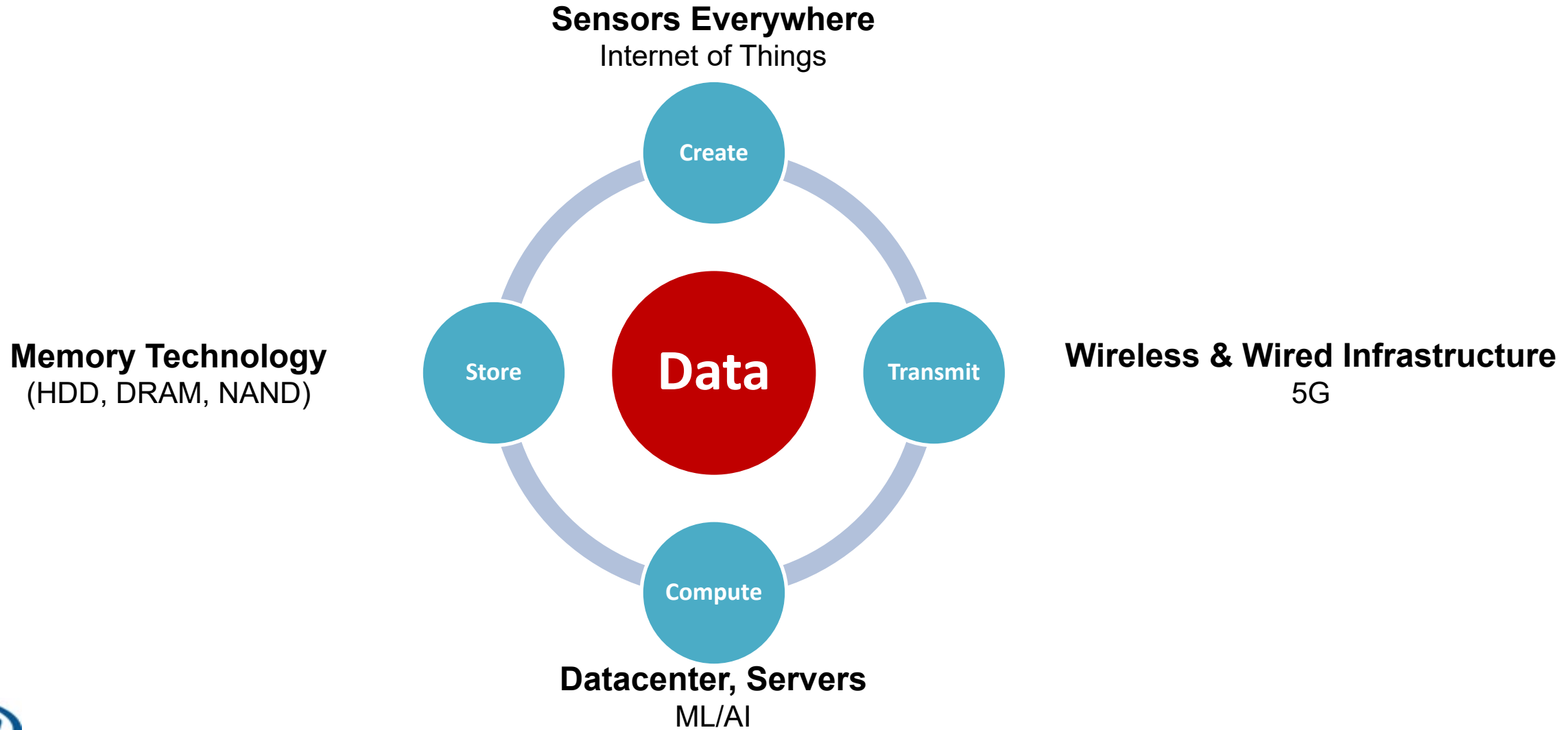


Semiconductors and Systems

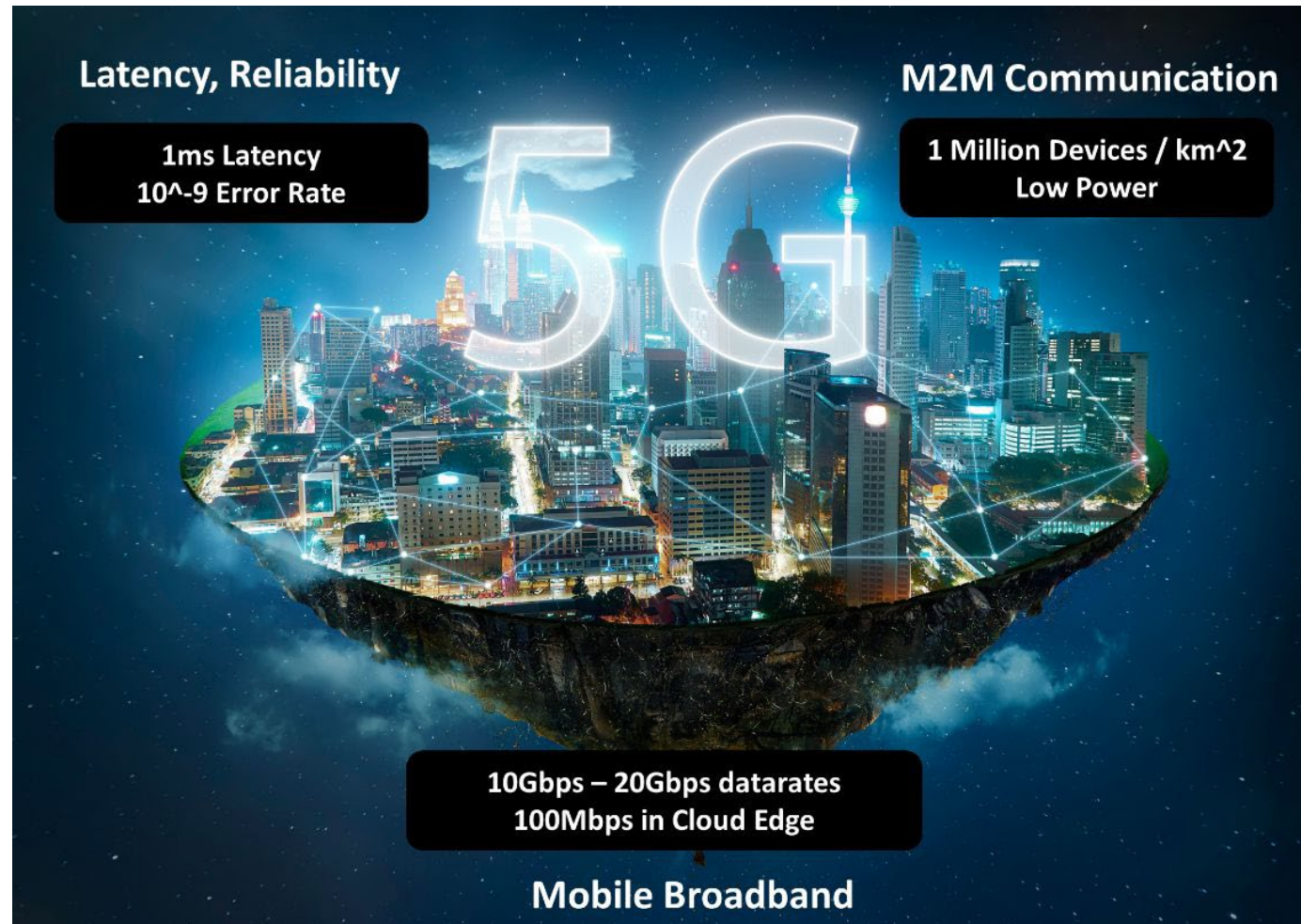
EDA system design enablement

© Cadence Design Systems, Inc - All Rights Reserved

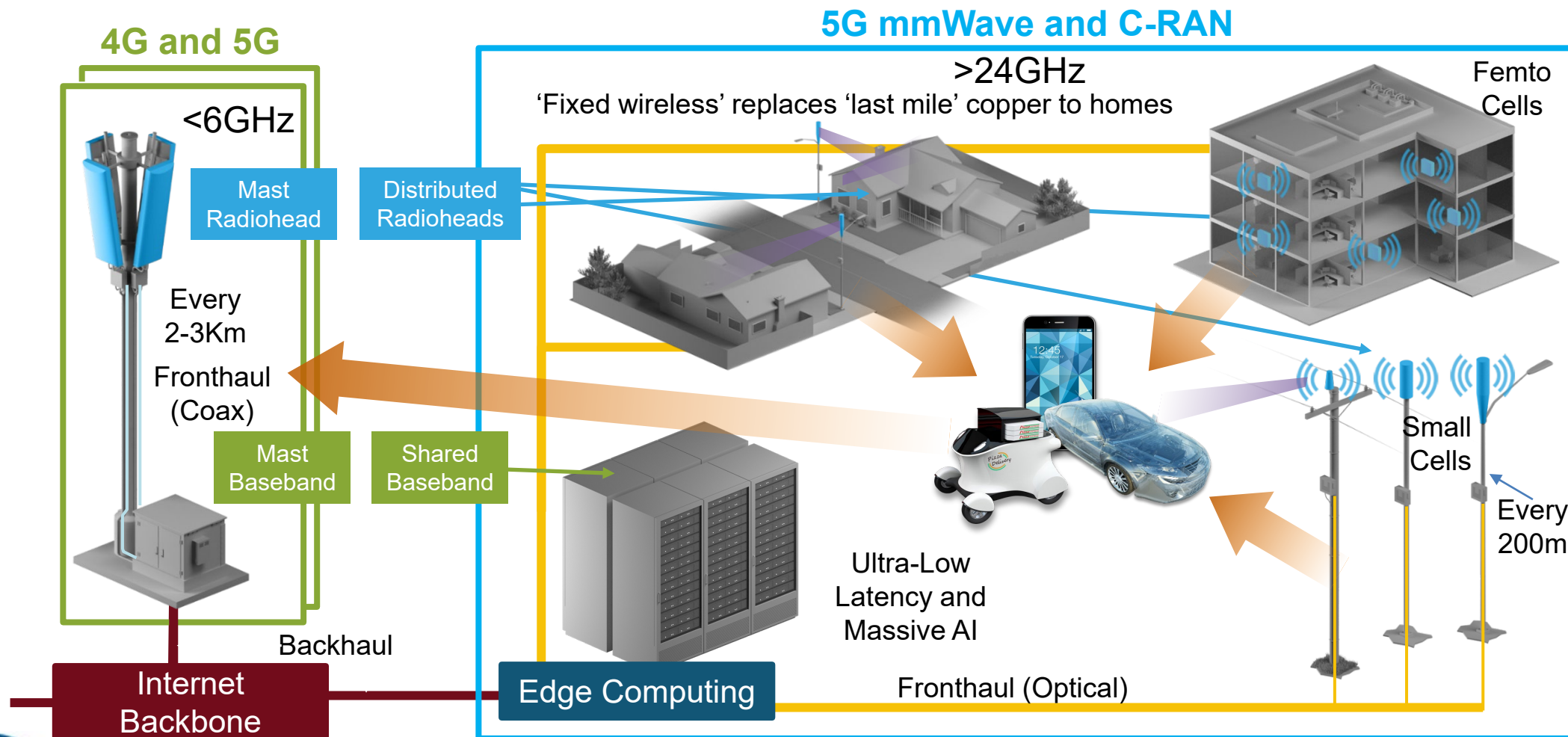
A Data-driven World



The Promise of 5G – Where to Start?



5G — New Bandwidth, New Subsystems



5G Requirements 1/2

Enhanced mobile broadband

- **Needs**
 - Faster speed, Lower latency, Greater capacity
 - On-the-go, ultra-high-definition video, virtual reality, and other advanced applications.
- **Design characteristics**
 - traditional large (>200MG), complex designs
 - requiring full-chip execution for SW - Emulation

Internet of Things

- **Needs**
 - Existing networks struggling
 - 5G unlocks (IoT) - more connections at once
 - Additional monthly revenues for carriers
 - IoT revenues smaller because of low usage
 - 5G competes against Wi-Fi and Zigbee.
- **Design characteristics**
 - much smaller (<32MG), very power sensitive
 - performance system dependent
 - multi-device simulation / emulation for QoS and performance validation

5G Requirements 2/2

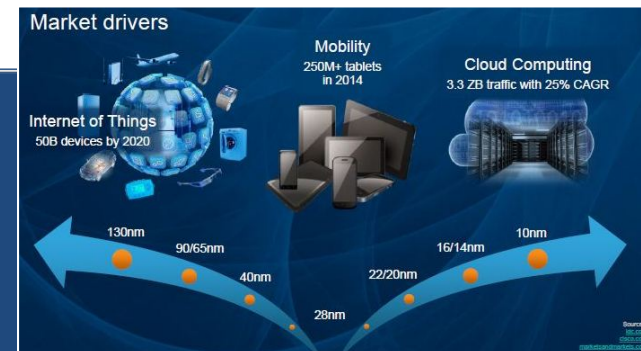
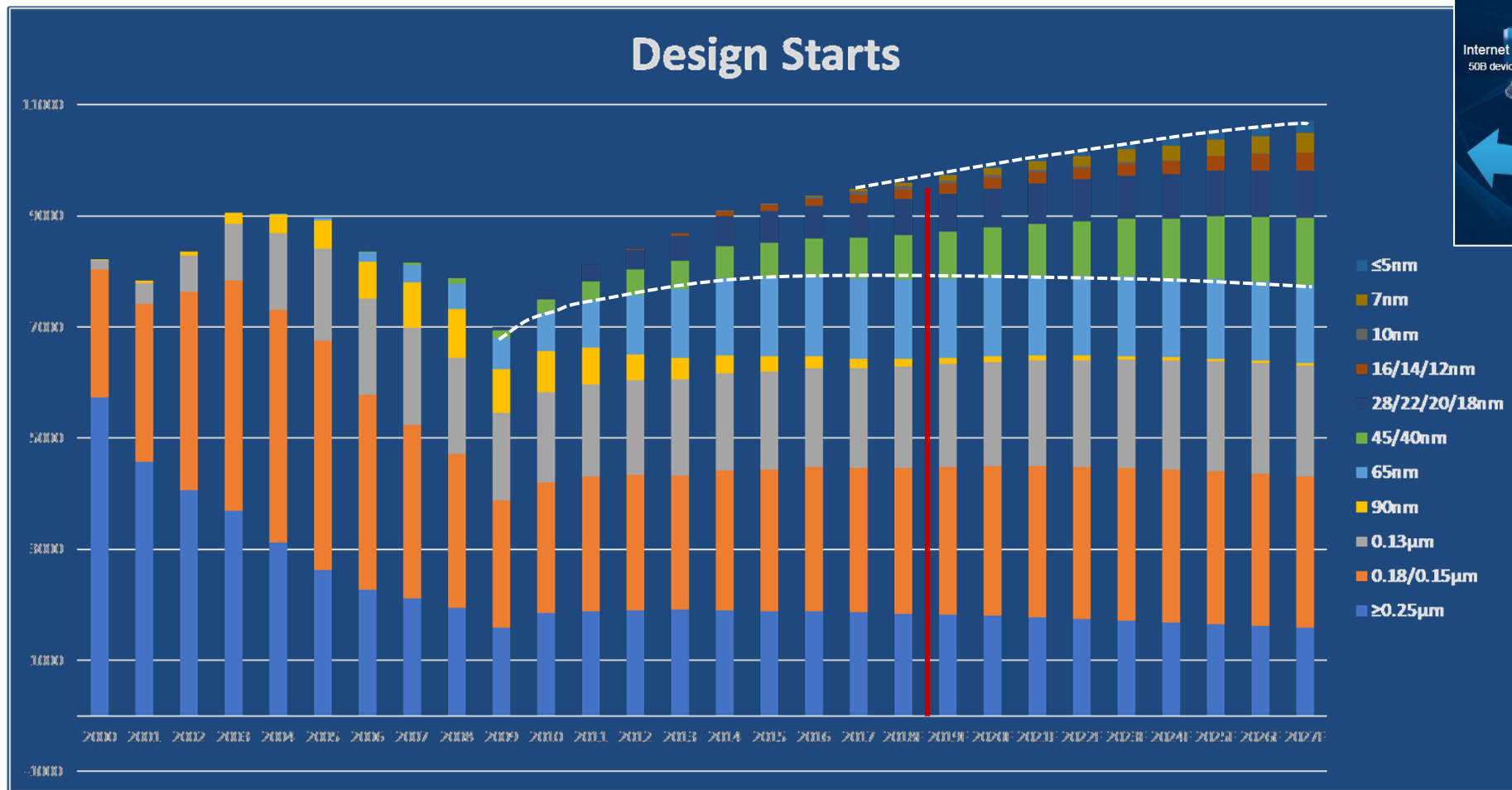
Mission-critical & control

- **Needs**
 - absolute reliability in medical, vehicle safety
 - Latency limiting factor
 - 5G delivering lower latency
 - New use cases in healthcare, utilities, traffic management, and other time-critical contexts
 - Operators expect only incremental revenue
- **Design characteristics**
 - Small-to-medium designs (<200MG)
 - Functional safety drives need for system emulation

Fixed wireless access

- **Needs**
 - 5G, millimeter wave spectrum, capable of delivering speeds of more than 100 Mbps to the home
 - Viable alternative to wired broadband
 - New revenue stream for wireless operators in areas with less fiber/cable access
- **Design characteristics**
 - Extension to traditional base-station developers
 - More complexity requiring system emulation
 - Design size expected to be large (>200MG)

Design Start Market: Bifurcation

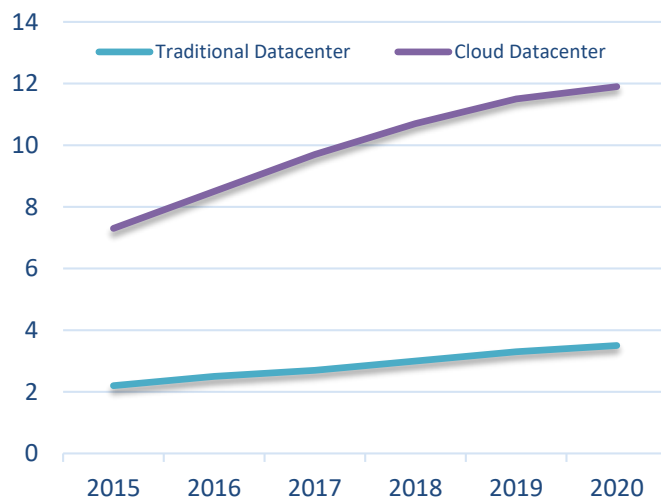


Source: IBS 2014 to 2018

Datacenter Opportunities

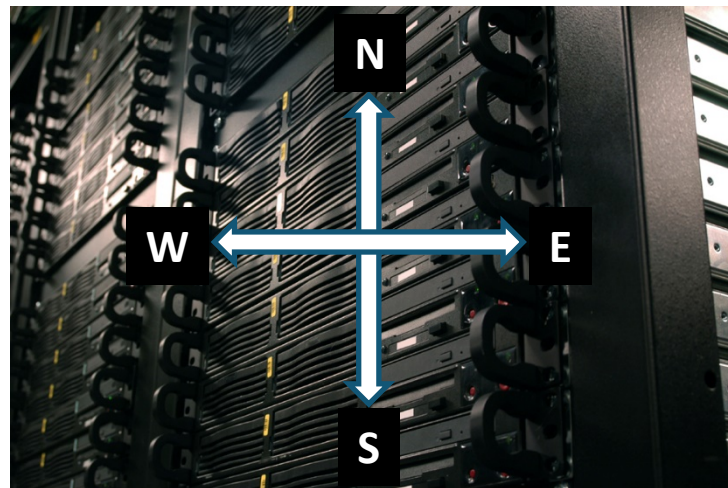
Workload-optimized, high-performance compute, connectivity, accelerators – AI/ML/DL

Typical Workloads per Server



Hyperscale Optimization CPU

- Workload optimized
- Machine learning
- Deep learning
- Accelerator offloads



Rack-Level Connectivity





- Leaf /spine
- Memory pool (HBM)
- Connectivity / SiP
- Reduced latency
- Mesh / 3D-torus / fabric



Scale Out Clusters

- DNN
- SSD / NVMe
- Coherency
- VM / containers
- Mesh/3D-torus

AI Chip Datacenter Technology

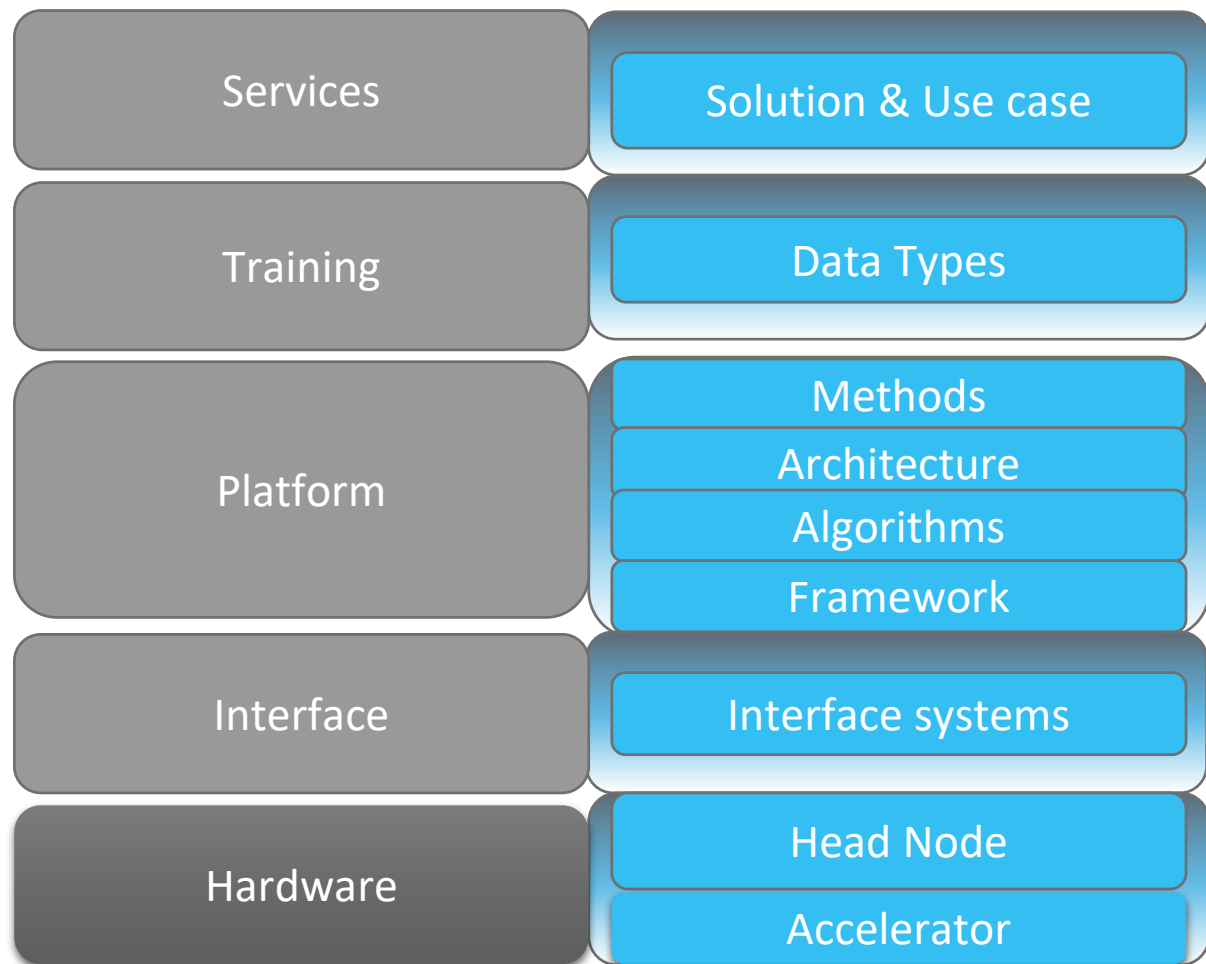
	Opportunities in existing market	Emerging opportunities
 <p>Compute</p>	<ul style="list-style-type: none"> Accelerators for parallel processing, such as GPUs¹ and FPGAs² 	<ul style="list-style-type: none"> Workload-specific AI accelerators Quantum computing, neuromorphic
 <p>Memory</p>	<ul style="list-style-type: none"> High-bandwidth memory On-chip memory (SRAM³) 	<ul style="list-style-type: none"> Quasi-volatile memory Non-volatile memory (NVM) PCM, MRAM
 <p>Storage</p>	<ul style="list-style-type: none"> Potential growth in demand for existing storage systems as more data is retained 	<ul style="list-style-type: none"> AI-optimized storage systems Emerging NVM (as storage device)
 <p>Networking</p>	<ul style="list-style-type: none"> Infrastructure for data centers 	<ul style="list-style-type: none"> Silicon photonics Programmable switches

¹Graphics-processing units. ²Field programmable gate arrays. ³Static random access memory.

Source: McKinsey analysis and Cadence

© Cadence Design Systems, Inc - All Rights Reserved

AI/ML Technology Stack

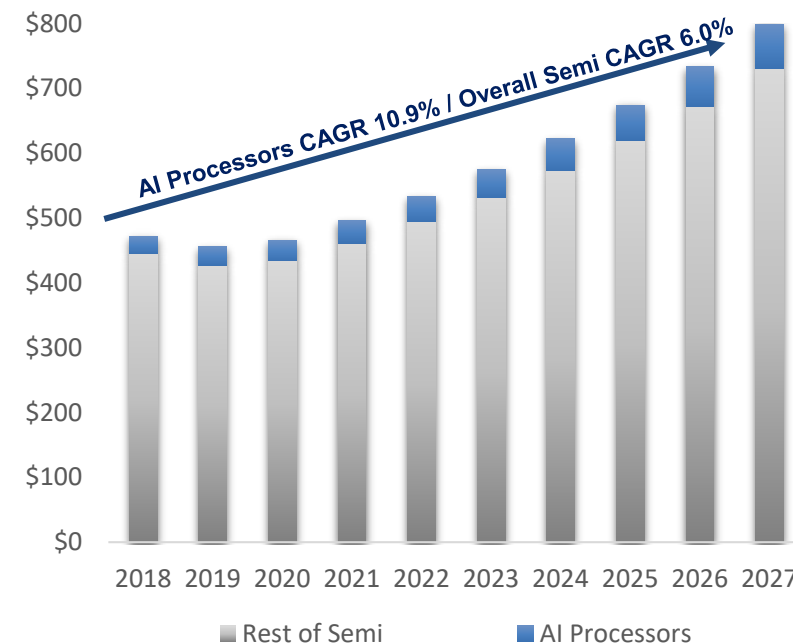
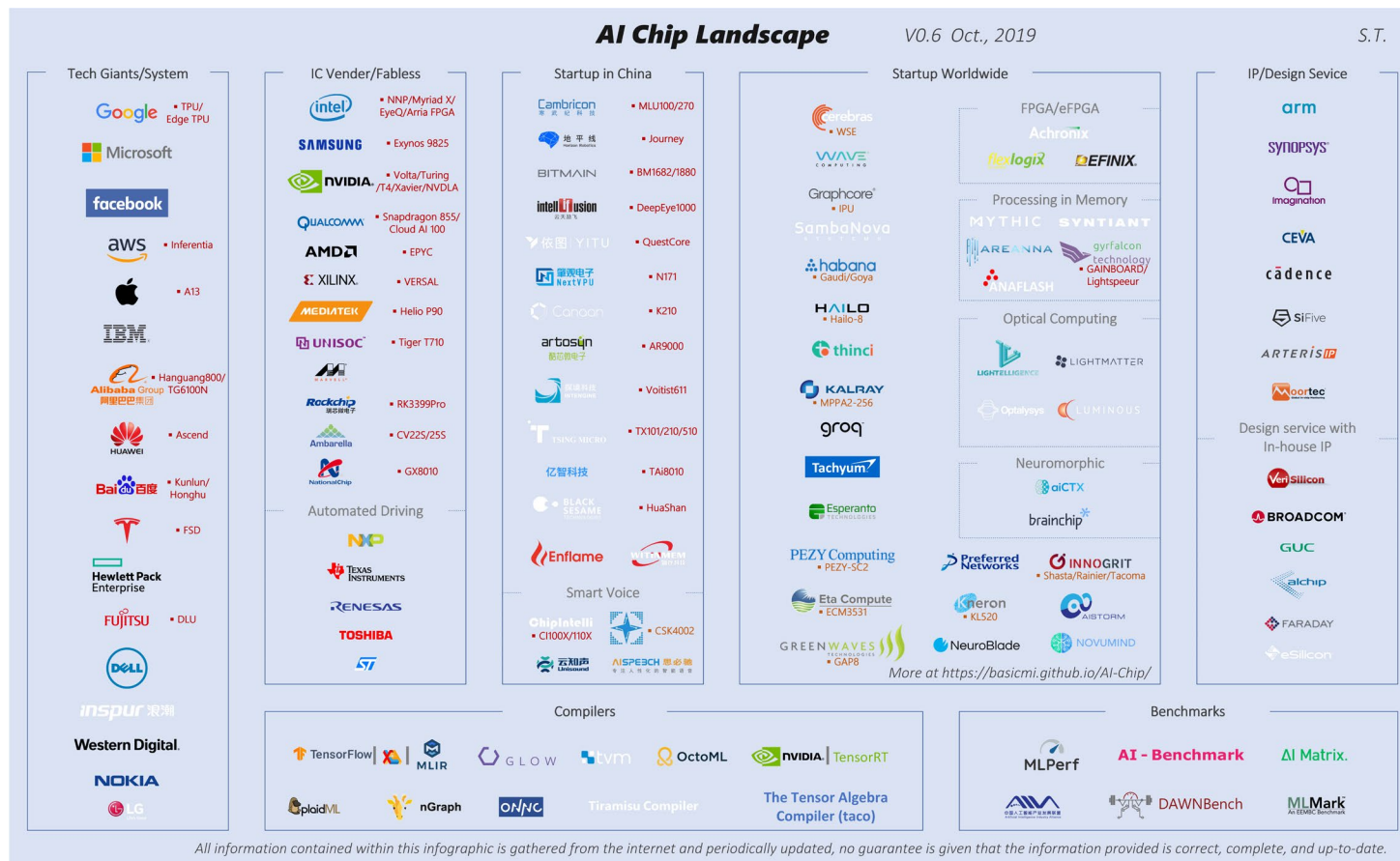


- Complex optimization across hardware, software, platforms and services
- Silicon designed to perform highly parallel operations required by AI enables simultaneous computations

“AI could allow semiconductor companies to capture 40 to 50 percent of total value from the technology stack, representing the best opportunity they’ve had in decades.”

Source: McKinsey analysis and Cadence

AI Chip Landscape, Key Ecosystems

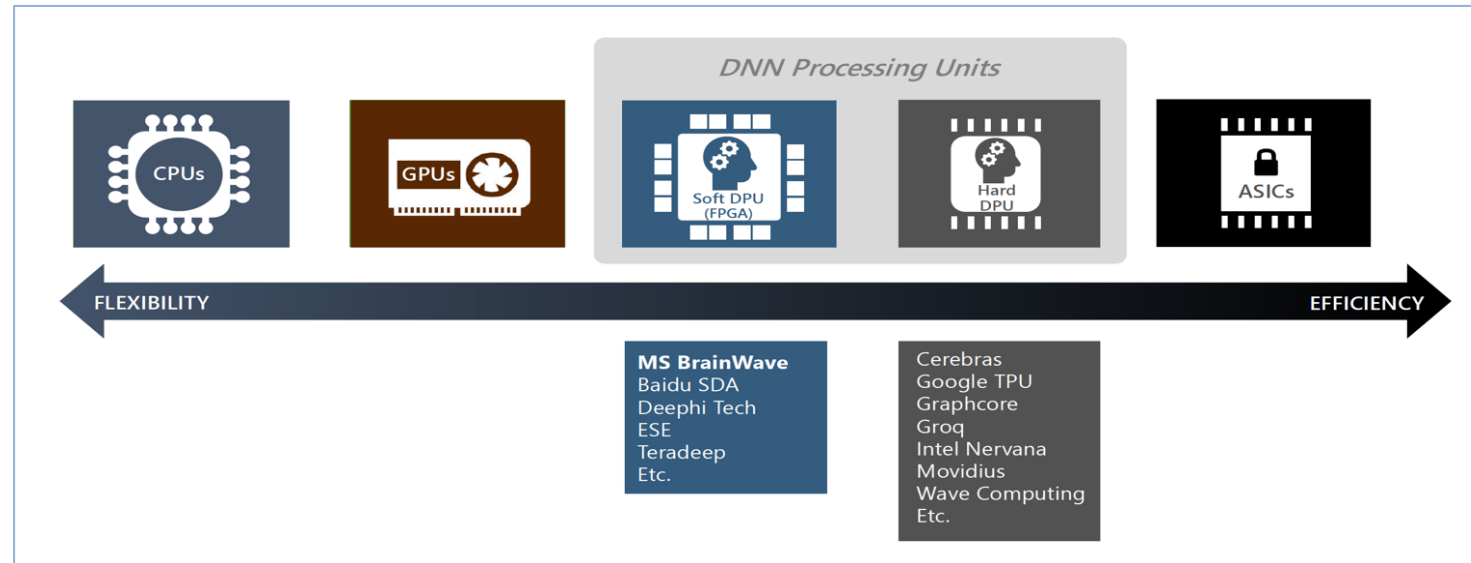


Source: International Business Strategies, Inc. (IBS)
Semiconductor Market Data Centers, April 2019.

AI processors expected growth:
3X faster than total semi market

<https://github.com/basicmi/AI-Chip/>

Many AI/ML Implementation Options



Source:
Microsoft. Hot Chips 2017

- Deep learning algorithms show a high degree of parallelism
 - General-purpose CPUs designed for sequential workloads
 - GPUs with massively parallel execution capability have popularized DNNs over the last few years
 - Proliferation of FPGAs and domain-specific accelerator ASICs for more efficiency

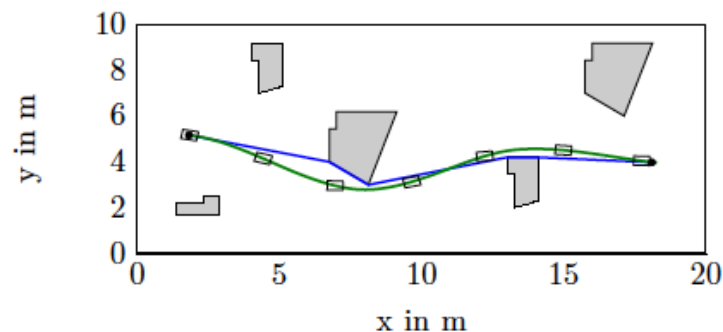
Decision Making is Shifting from MPC to AI

Example: Path Planning

Traditional:

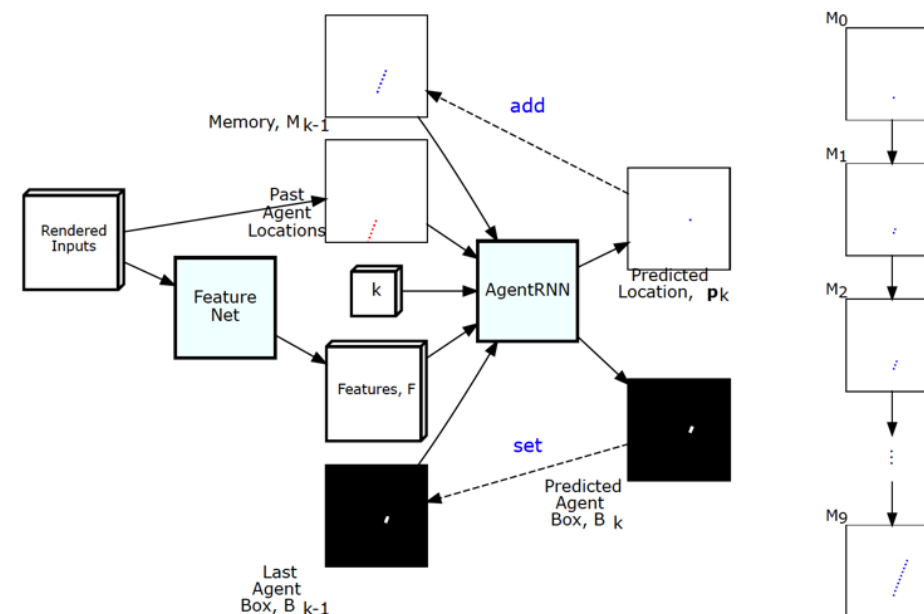
Model Predictive Control (MPC)

MPC frameworks running on CPU to generate safe trajectories

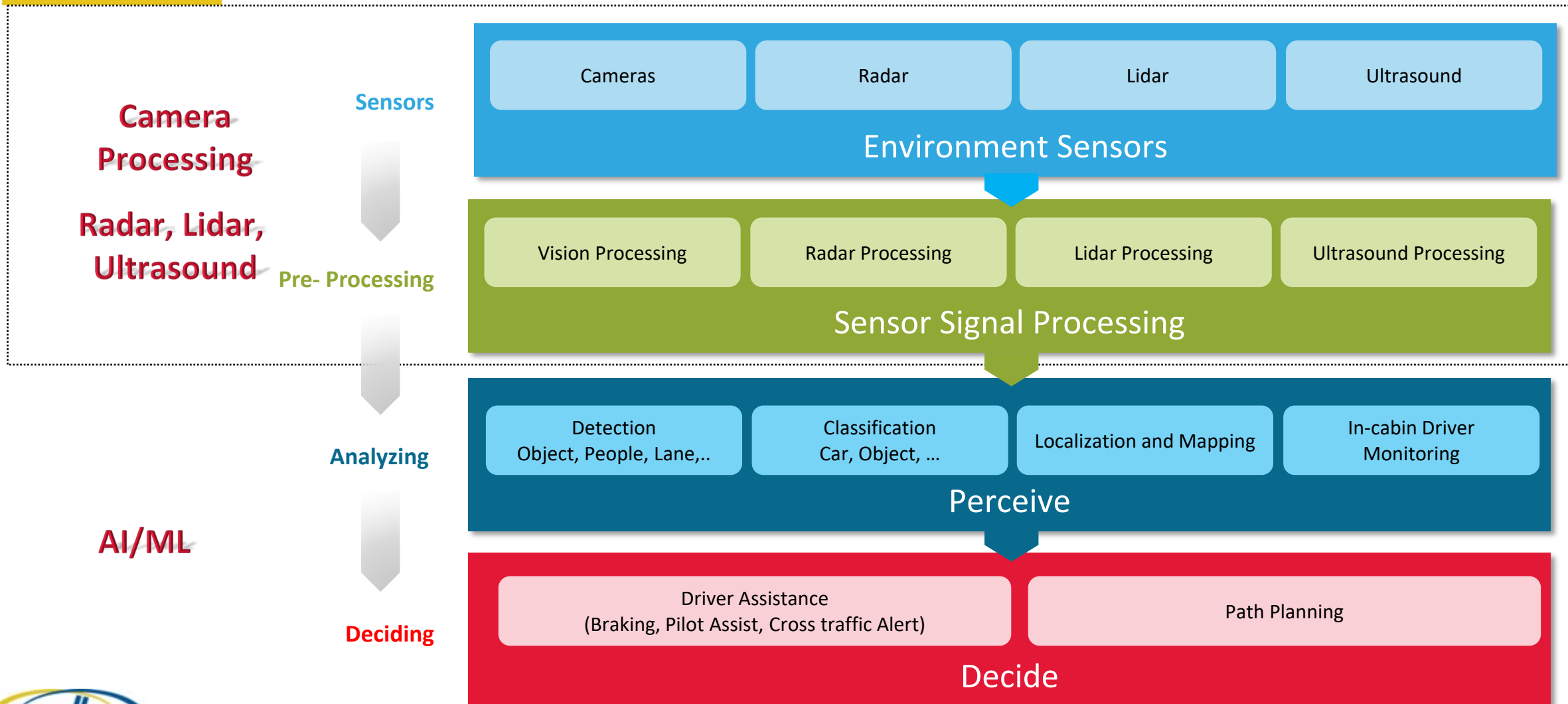


Present and Future:

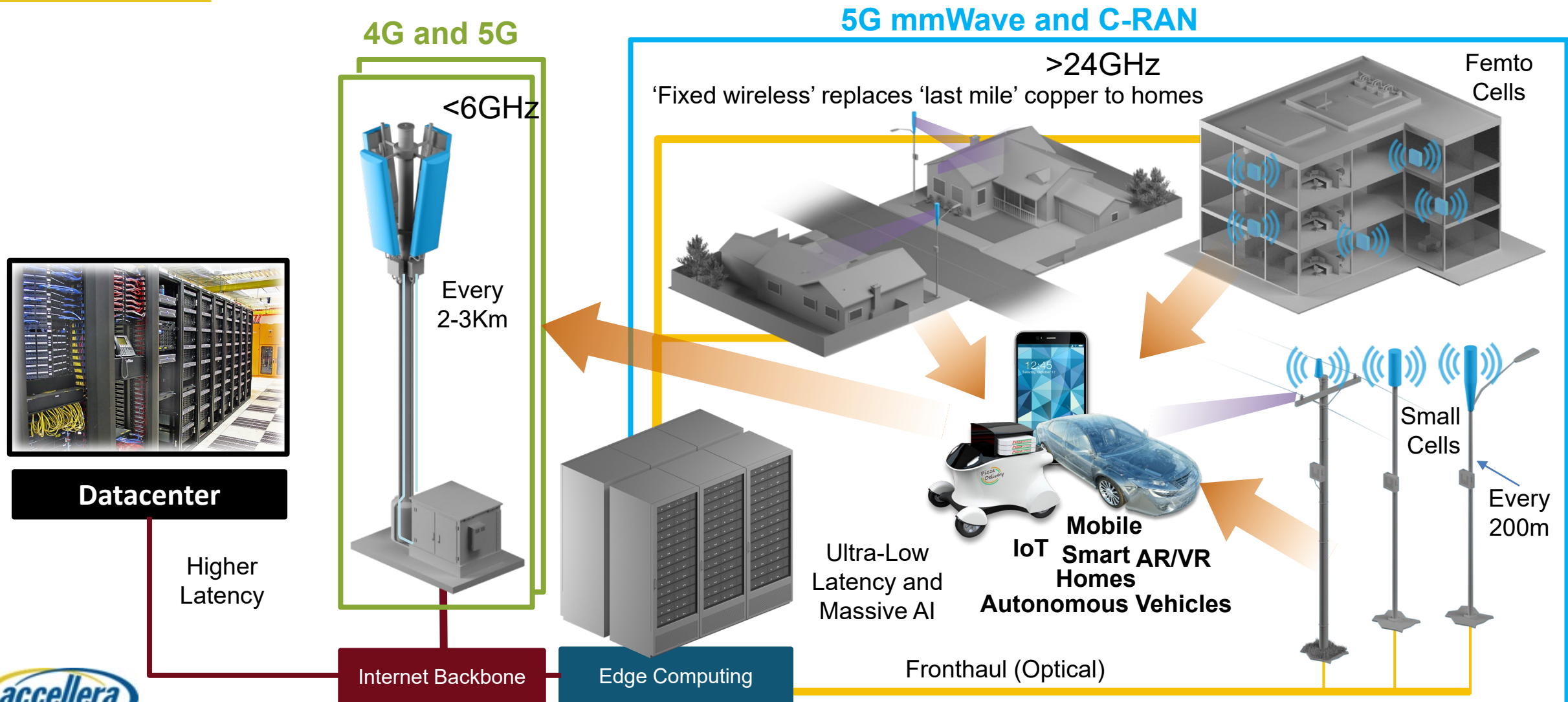
ChauffeurNet - Google Waymo
Convolution and recurrent neural networks



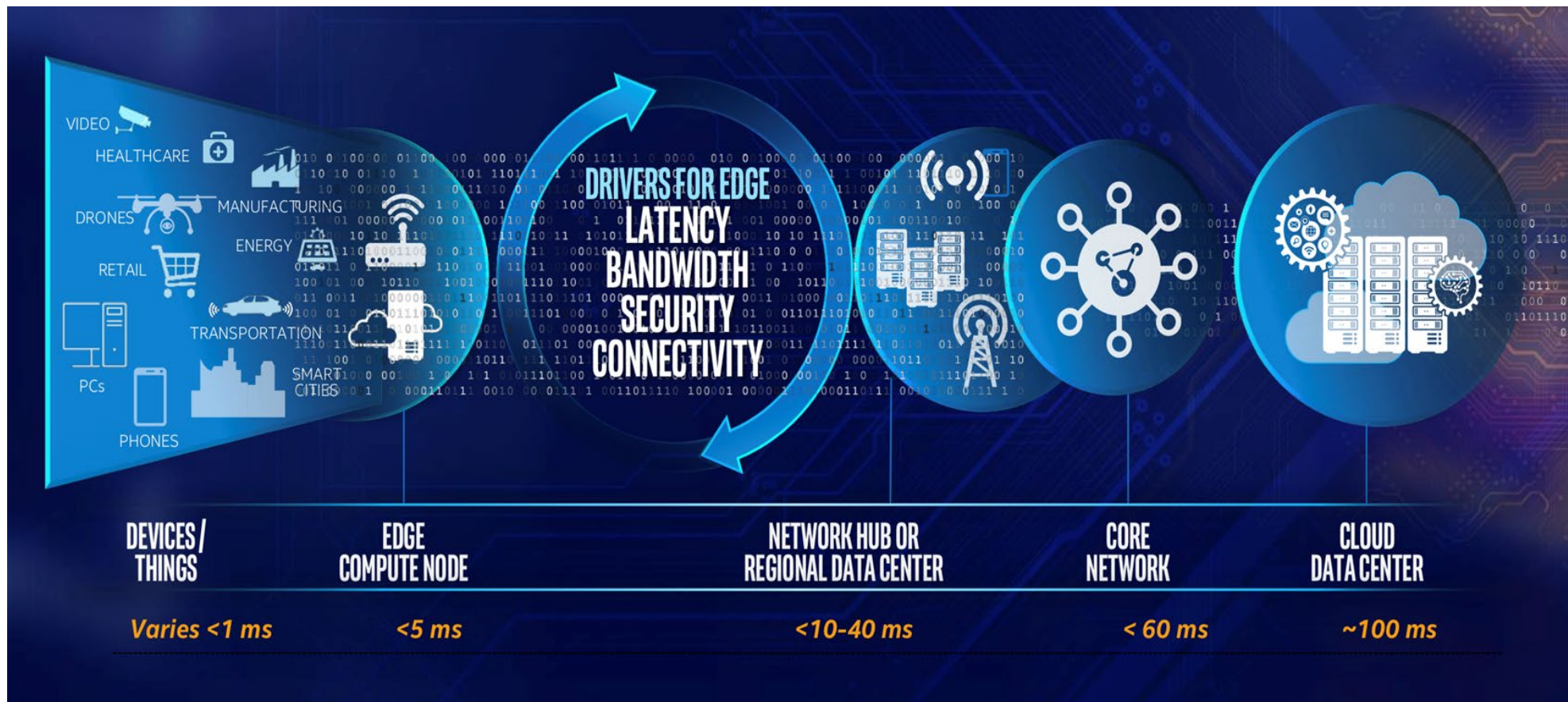
Accelerating AV and ADAS Experiences



AI, ML and 5G are connected



Data Drives Edge Computing



Source: Intel

Location of data is crucial!

Datacenter Training

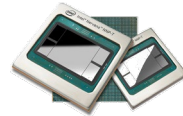
- High complexity NNs for ADAS
- Analysis to predict our behavior
- ...



Google Cloud TPUs



Nvidia GPUs



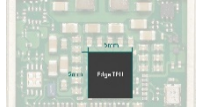
Intel Nervana NNP-T

Datacenter Inferencing

- More complex data - higher latency
- Train and deploy inferenced NN
- ...

Edge Training

- Face recognition to unlock phone
- Update with localized data
- ...



Google Coral



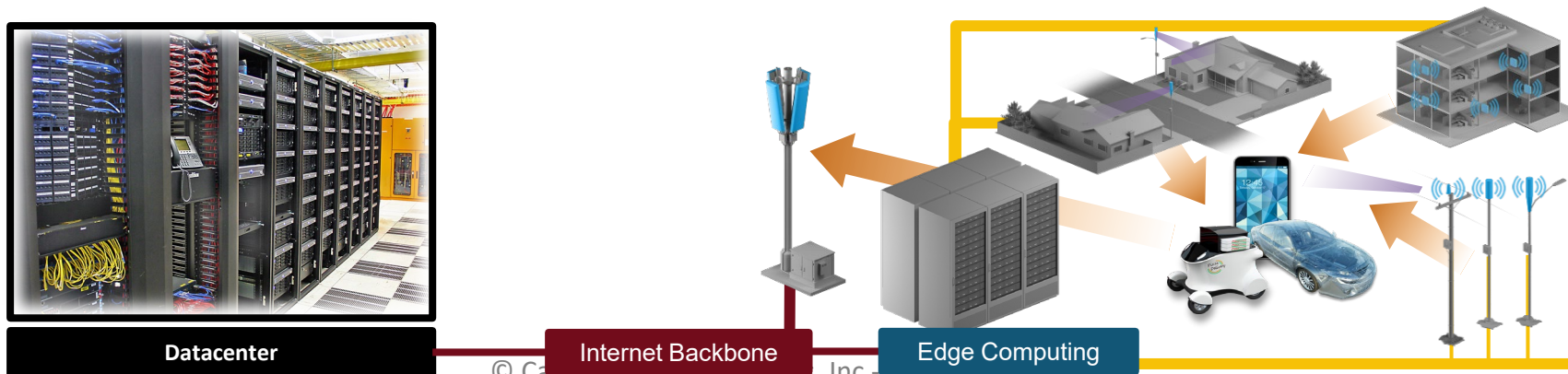
Nvidia Jetson Nano



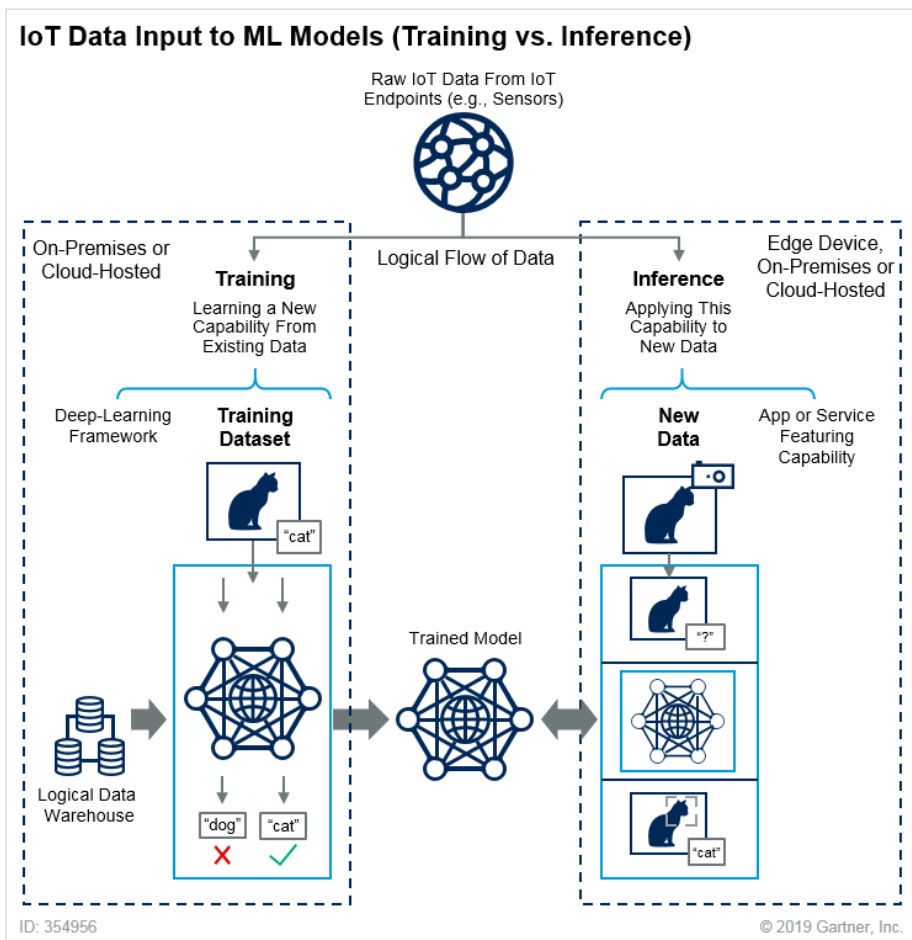
Intel Up Squared

Edge Inferencing

- Real time decisions (Vision)
- Phone unlock
- ...



Training & Inference at Datacenter & Edge



- **AI in Datacenter**
 - Scientists directly craft models for all situations
 - Model training is automated
 - Supervision is required in most cases to label data
 - Full fidelity data must be used for training
- **AI at the Edge**
 - Automatically adapt, deploy on available infrastructure
 - Performance needs to be achieved within the edge products' power and bandwidth constraints

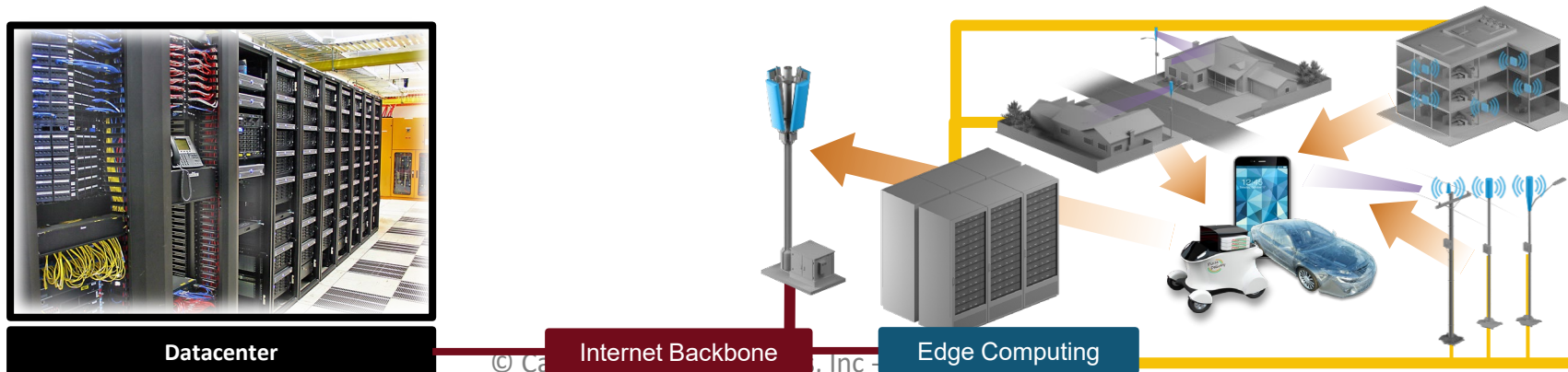
Design & Verification Requirements

Datacenter

- Big chip solutions for GPUs, NOCs, Workload specific AI accelerators
- High capacity verification, Advanced nodes
- Specific models for IP, verification
- Advanced flows Low power with high performance
- Workload optimization server flows (SBSA, ...)

Edge

- Specialized IP, extendable processing
- Advanced IP models (verification, implementation)
- Comprehensive low power tool set for Digital and Analog IP



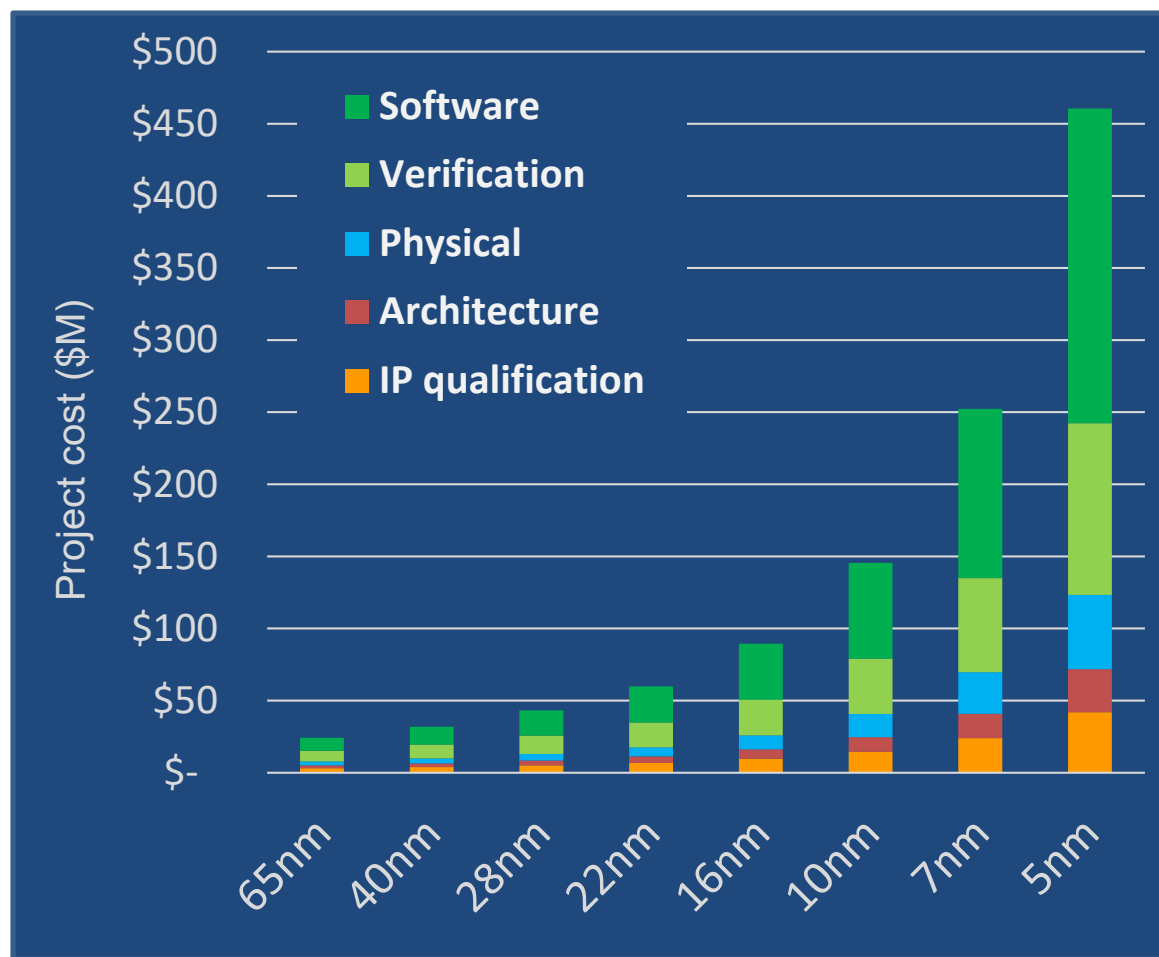
Design and Verification Challenges

- Diverse requirements
 - Training – High throughput, big designs
 - Inference – Flexibility, Lowest power
- Verification Challenges
 - Significant software content
 - Big Designs – Emulation throughput, debug
 - Physical and virtual interfaces , Virtualization
- Physical Design Challenges
 - Complex SystemVerilog Descriptions
 - 1st Floorplan uncertainty
 - Advanced node foundry limits and closure
- System Design Challenges
 - Mix of older and advanced nodes
 - New system design innovations with 2.5D/3D



Example: AI Processor on FPGA Prototype

Unifying Challenge: Verification and Software



Expanding
Software
Challenge

Expanding
Verification
Challenge

**VERIFICATION
&
SOFTWARE**

2^N

Source: IBS 2018

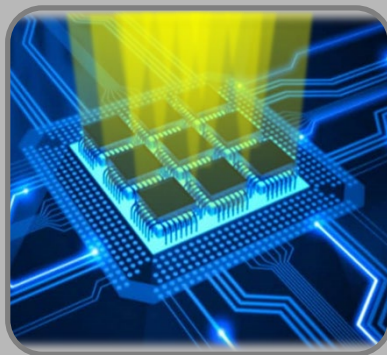
VERIFICATION OPTIONS TO *ENABLE* 5G/AI/ML

AI/ML Verification Requirements



Formal

- SoC **Scalability**
- **Smart Proof** technology
- Optimized **regressions**



Simulation

- **Fast simulation** for high-activity designs
- UVM **randomization**
- Fast elaboration for **replicated structures**
- Coverage **metrics**



Emulation

- **Billion-gate designs**
- Parallel Partition Compiler
- INT8 to INT64
- **Power / Performance**
- Memory **models**: HBMx
- **Senor model**: MIPI CSI

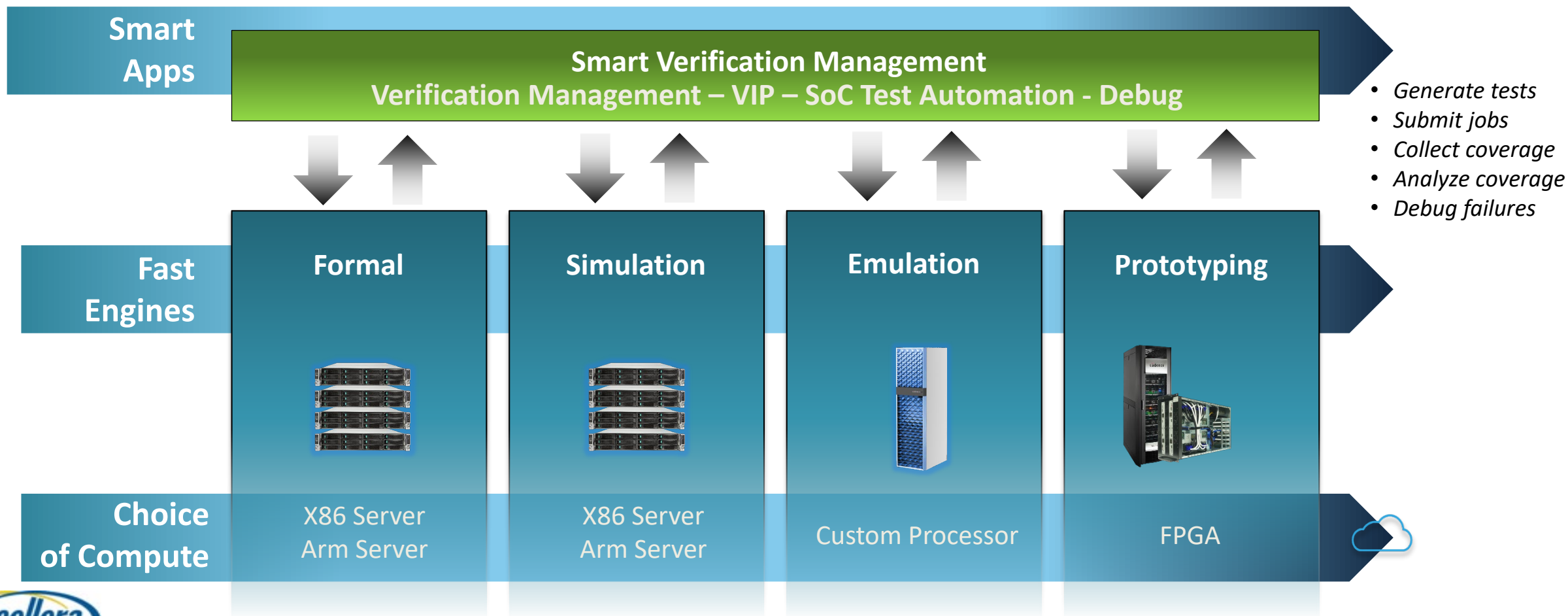


Prototyping

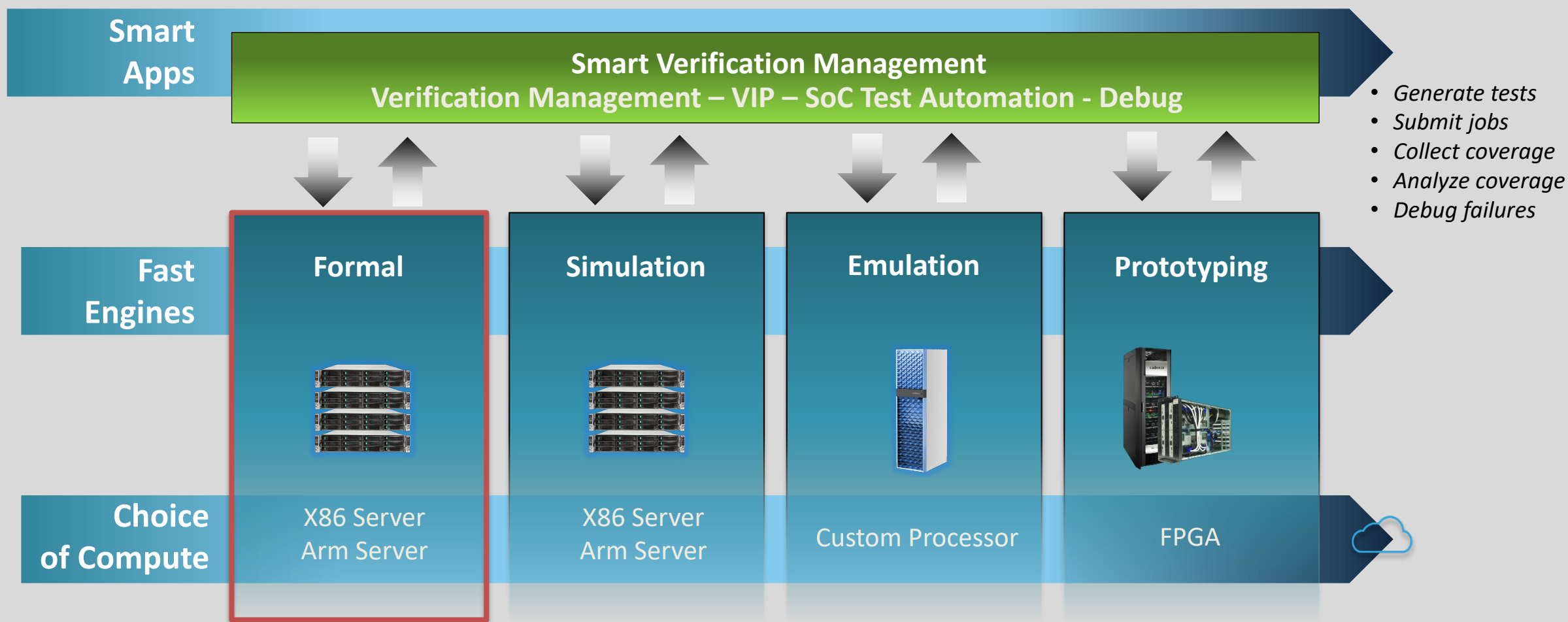
- **Scaling** to large designs
- Unified frontend
- ICE test suite: **Faster, data-extended (DL) regression**
- **SW driven validation** – deep learning data training refinement

Verification Throughput!

Find and fix the most bugs per \$ invested in bare metal compute



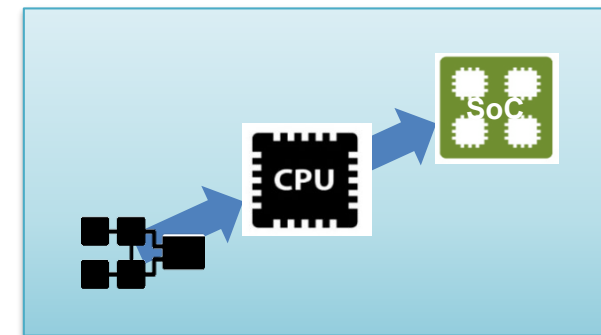
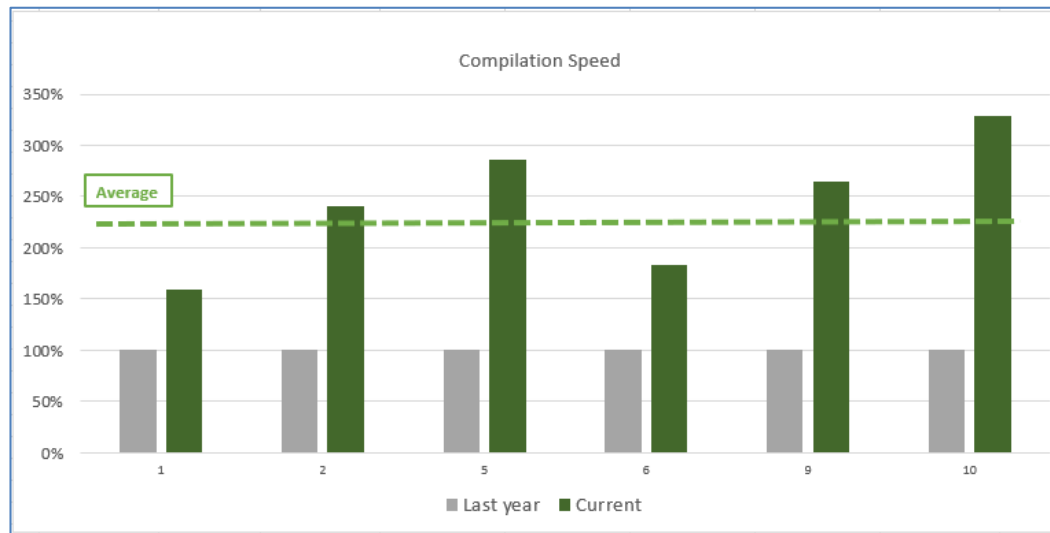
Formal Verification



SoC Design Scalability

Compilation Performance

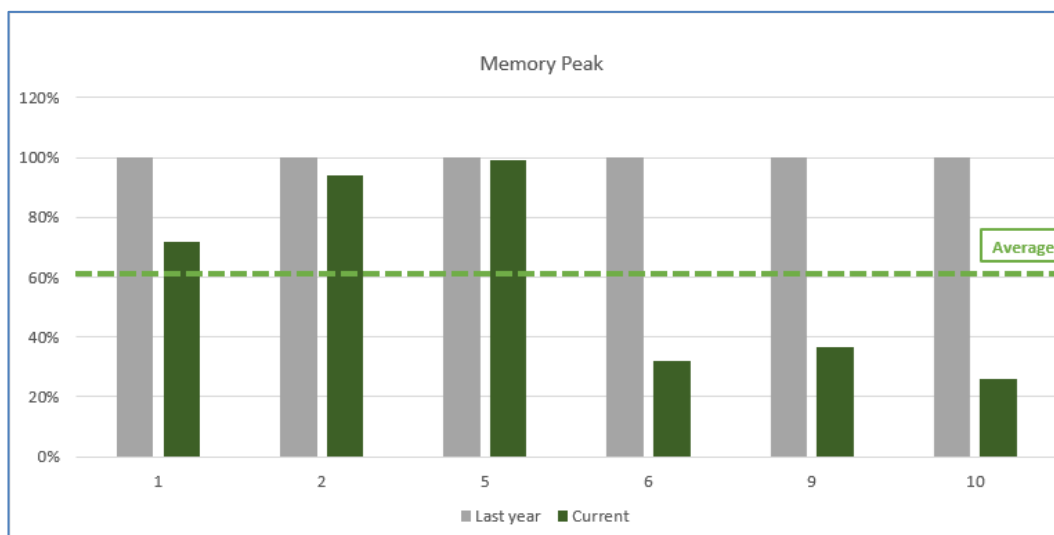
✓ **2X faster in 2019**



- 6 big designs from 4 different customers
- Min. compilation speed: 100M gates/hour
- Peak memory: less than 0.5 Kbyte/gate

Compilation Memory

✓ **40% smaller in 2019**



© Cadence Design Systems, Inc. All rights reserved.

Smart Proof Automation Framework



Component & Data Management

Proof Profiling Data

- Regular read/write

Proof Caching

- Cache storage in single file
- Automatic cleanup of old cache data

Multi Advisor Proof Orchestration

- Forced on
- No overwrite on engine mode

Engine Algorithm Selection

- Automatic training and inference

Learning

Machine Learning

Optimizes subsequent runs/regressions

Optimizes out-of-the-box proofs



Find more bugs



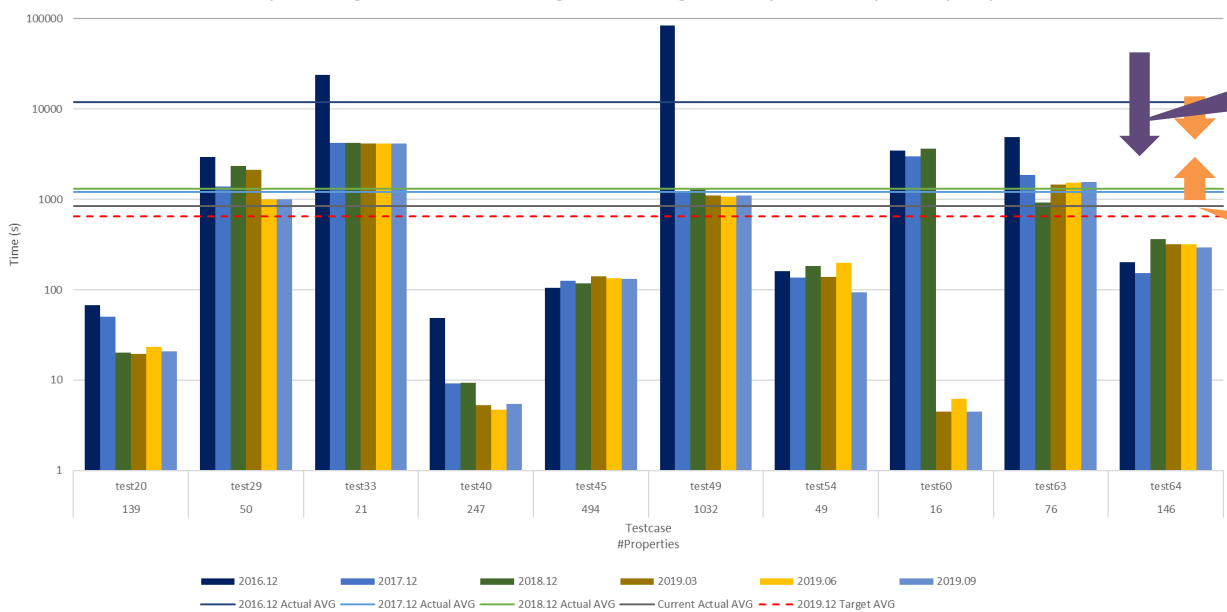
Better convergence



Faster proofs

Property Performance and Convergence

Fully Converged test cases: Average of Convergence Elapsed Time per Property



Avg. 14X
speedup in
last 3 years

Avg. 1.6X
speedup in
2019

Core Engine Performance

- ✓ 10 fully-converging designs
- ✓ 14X speedup last 3 years
- ✓ 1.6X speedup in 2019

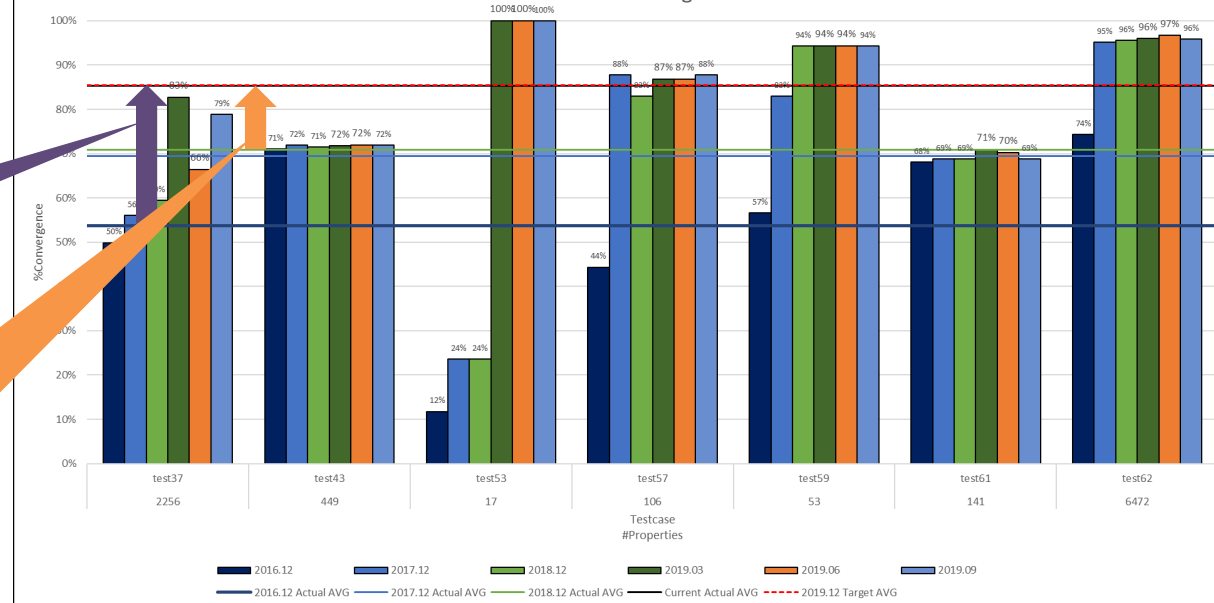
Core Engine Proof Success

- ✓ 7 “hard” designs: significant # undetermined properties
- ✓ 2X reduction in undetermined properties in 2019

Avg. 3.2X non-
converged
reduction in last
3 years

Avg. 2X non-
converged reduction
in 2019

Hard test cases: Convergence Rate



Proof profiling data (PPD) and proof caching

- Challenge
 - Running engines to reproduce previous tool results can often be extremely costly in terms of resources, especially in cases where the design/environment has minimally changed or not changed at all
 - Need ability to learn from past runs of a design, to optimize subsequent proofs
- Approach
 - Keep record of previous proof strategies that worked, so they can be tried again
 - `prove -save_ppd -with_ppd`
 - Be able to harvest previous proof results, when confirmed that they still apply
 - `set_prove_cache on`
- Benefits
 - Better convergence: faster proofs on properties determined before frees up resources to work on other properties
 - Smarter use of machine time
 - Helps with reproducibility of proof results

Proof profiling data (PPD)

- Saves knowledge collected during a prove command to be used as recommendations in subsequent runs

JasperGold®

✓	Type	Name	Engine	Bound	Time
✓	Cover	state_23	Q3	17 - 23	10.9
✓	Cover	state_24	U	22 - 24	80.8
?	Assert	property:0	B	25 -	0.3
?	Cover	all_positions	Bm	23 -	128.6
✓	Cover	all_but_central	L	23 - 24	97.0

prove -save_ppd

ygproject/sessionLogs/session_0/sessionPPD.ppd

- Benefits
 - Increase the chances of converging faster on properties determined before, and that did not have their results

JasperGold®

prove -with_ppd

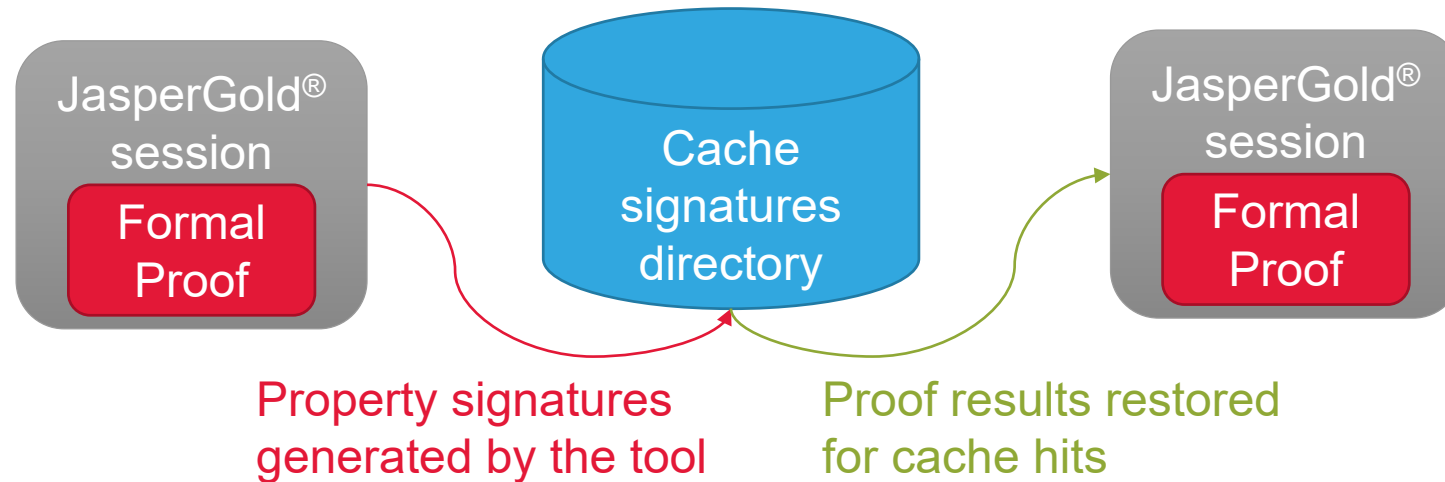
INFO (IPF149): Starting PPD exploitation on 36 matching properties

“Best engines” orchestration

Regular orchestration kicks when results in PPD fail to be reproduced

Proof caching

- Restores proof results based on properties' signatures
- Current version sensitive to small changes to design and environment



- Benefits
 - Save engine time when processing unchanged properties already determined before
 - Focus resources on properties that changed or that were never determined, improving convergence

A Continuum of Dynamic Engines

Verification and software platforms need to interoperate



SDK OS Simulation

Highest speed
Earliest in flow
Ignores HW
Easy replication
Cross-compile



Virtual Platform

Almost @ speed
Pre-RTL
Less accurate
TLM HW Debug
Great SW debug
Easy replication
Less HW detail
Slower with detail



HDL Simulation

KHz Range
Early RTL
Golden Reference
Best HW debug
Limited SW Debug
Easy replication
Mixed-abstractions
Slow SW execution



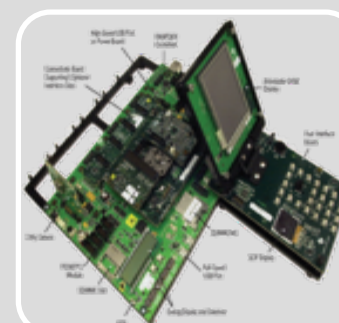
Acceleration Emulation

MHz Range
Early RTL
Min RTL mods
Detailed HW debug
Great SW Debug
Harder to replicate
Datacenter access
Contested Resource



FPGA Prototype

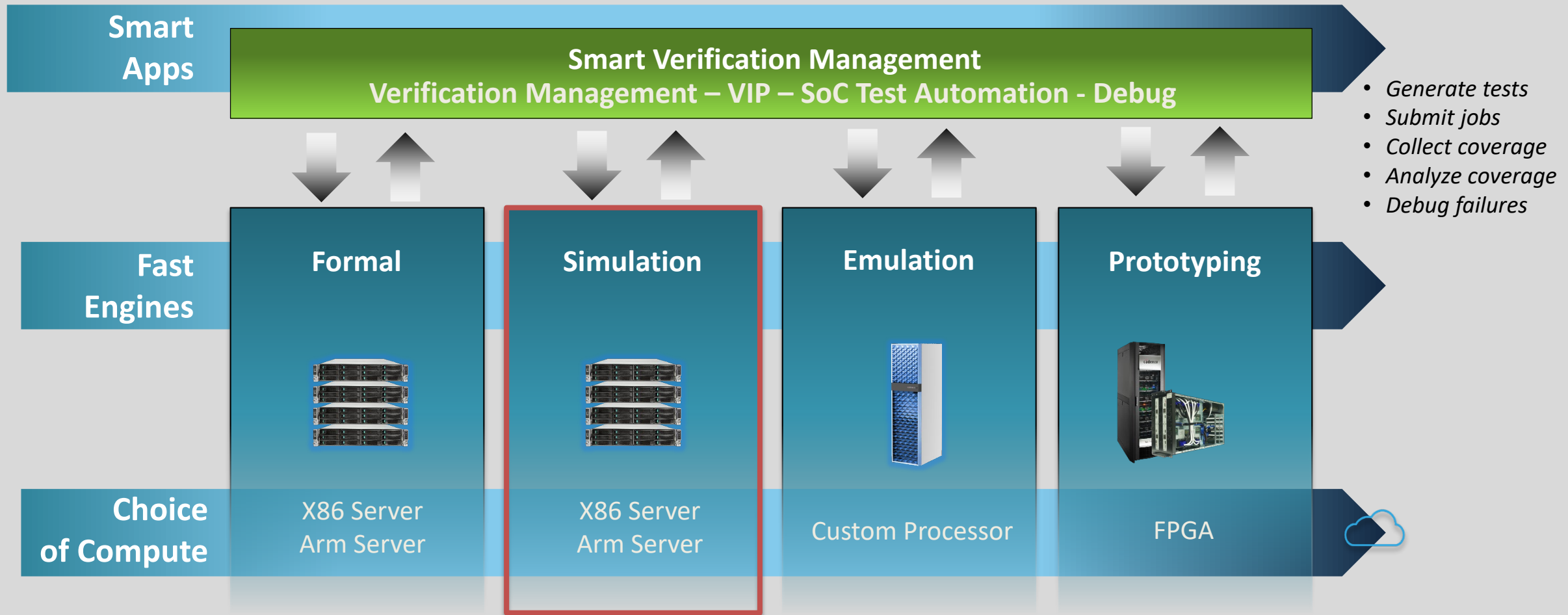
10's of MHz
Later RTL
Some RTL mods
Some HW debug
Great SW Debug
OK to replicate
Harder Bring-up



Prototyping Board

Real time speed
Fully accurate
Actual Silicon
Difficult HW debug
OK SW Debug
Easy to replicate
HW changes hard

Simulation



Simulation Performance

COMPILE

Incremental Build Parallel Build Hierarchical build w/ Cloning

Up to 10X speed-up

FULL REGRESSION THROUGHPUT

Relentless focus on performance

*Continuous Core
Performance Enh's*

*Save/Restart w/
dynamic test reload*

LONG TEST LATENCY

Rocketick multi-core technology

*MC-Lite
1.2-1.8X*

*MC-Rocketick
3-5X*

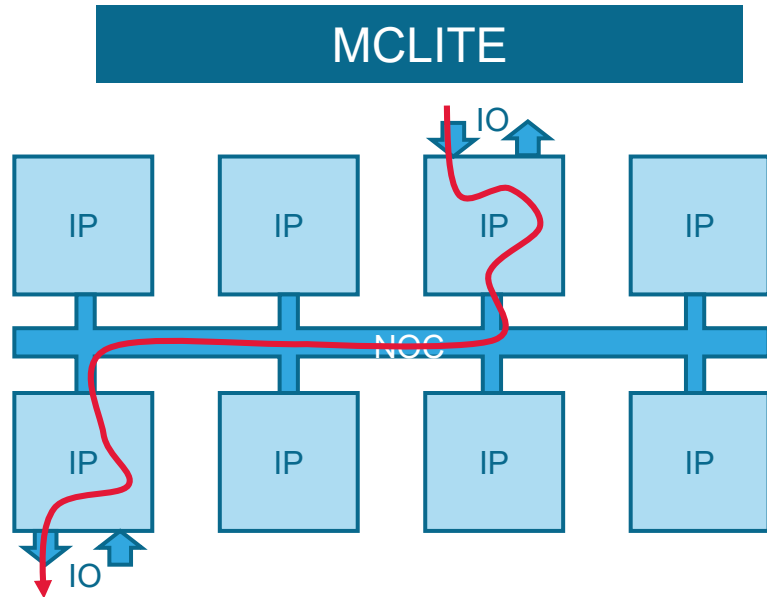
COMPUTE

X86

Arm

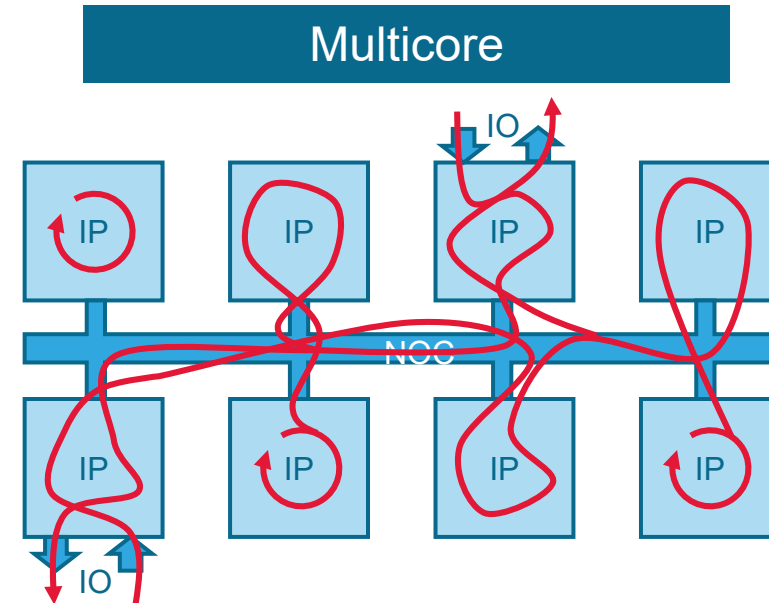
Cloud

Multi-Core Performance



Can use on any long running test

- Uses the single-core build
 - Decision to use mclite deferred to runtime
- Same as single-core scheduling
 - Guaranteed congruent results with SC
- **1.2 – 1.8X gains**

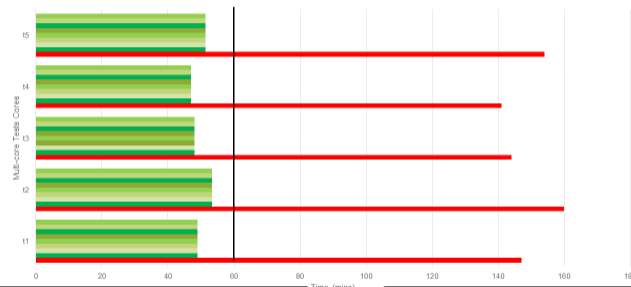


*Requires “high activity” DUT-heavy test
e.g. Gate Level ATPG*

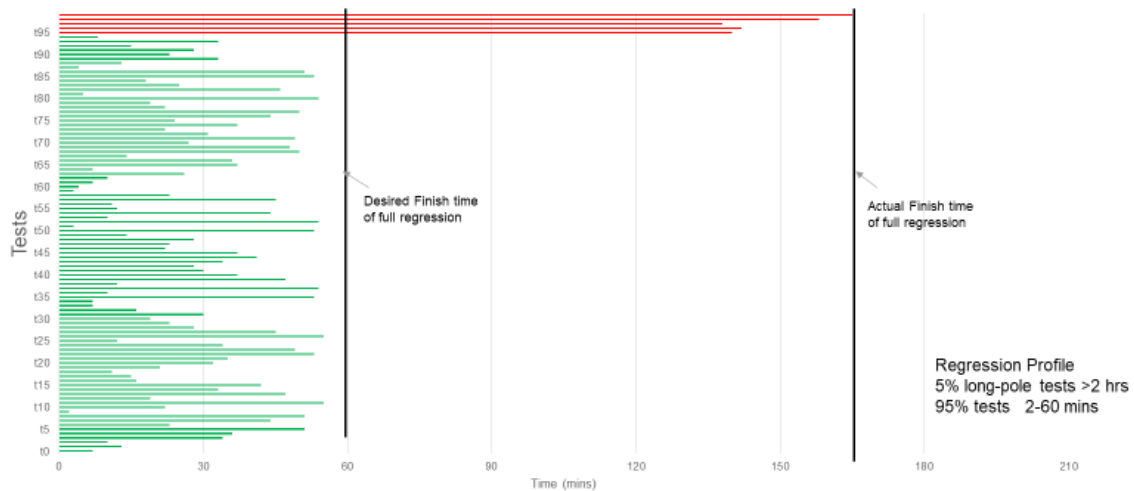
- Highly scalable parallel solution
 - Needs a special build
- Specialized (Rocketick) engine
 - Non-accelerable parts (testbench) on SC
- **3-5X gains**

Multi-Core and Simulation Regressions

Multi-Core Simulation: 3X faster using 8-cores



Simulation Regression: Xcelium Single-Core



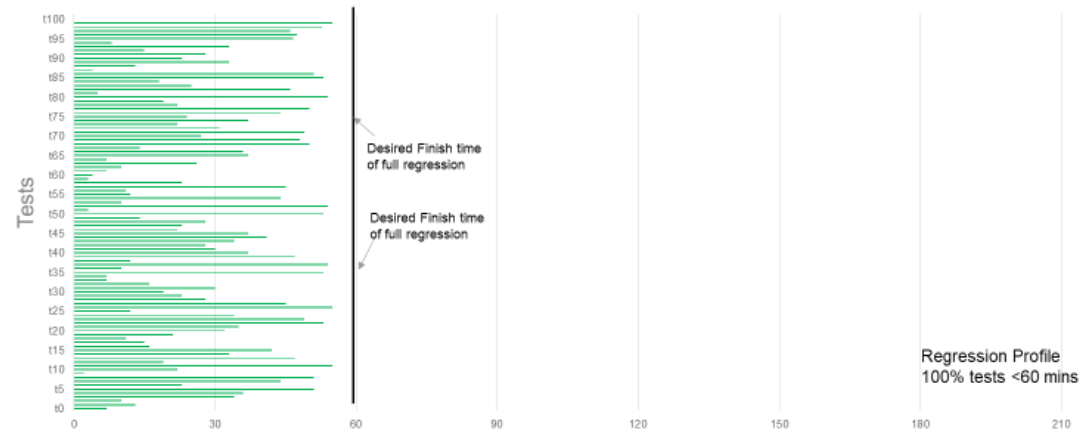
20

© 2019 Cadence Design Systems, Inc. All rights reserved.

cadence®

Simulation Regression: Multi-Core 3X on 5% Long-pole tests.

Actual finish time of regression = Desired finish time of regression

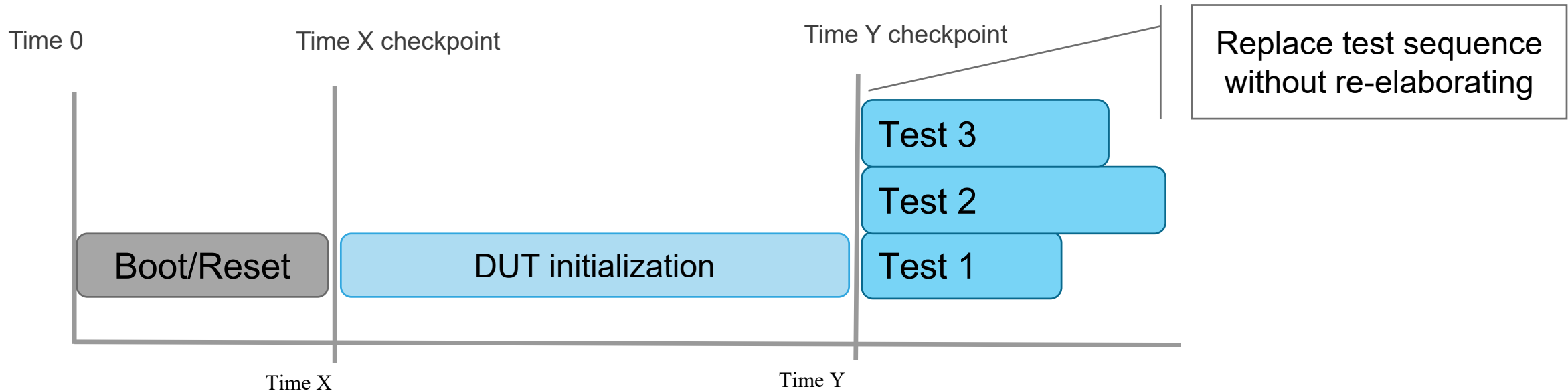


22

© 2019 Cadence Design Systems, Inc. All rights reserved.

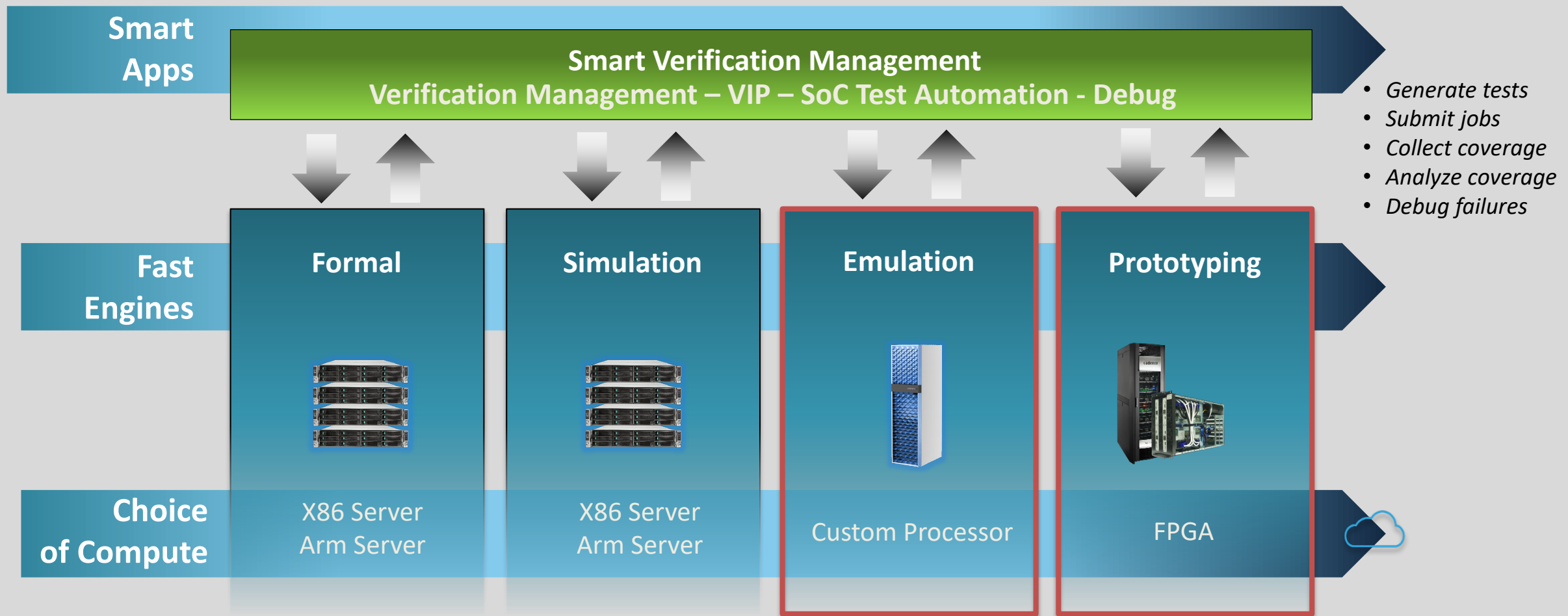
cadence®

Dynamic Test Reload for SystemVerilog



- SystemVerilog/UVM Dynamic Test Load
 - Load new SystemVerilog package into saved snapshot
 - Call testbench functions when the snapshot is reloaded
- Dynamic Base Snapshot: Time zero snapshot or saved snapshot
- Dynamic Test Snapshot: Contains the incremental new SystemVerilog package

Hardware Assisted Development



Emulation And Prototyping

Emulation “Debug your design”

Key Care Abouts

- **Predictable fast build:** Rebuild new
- **SOC level capacity:** IP level sim enough already
- **Fast and complete debug:** re

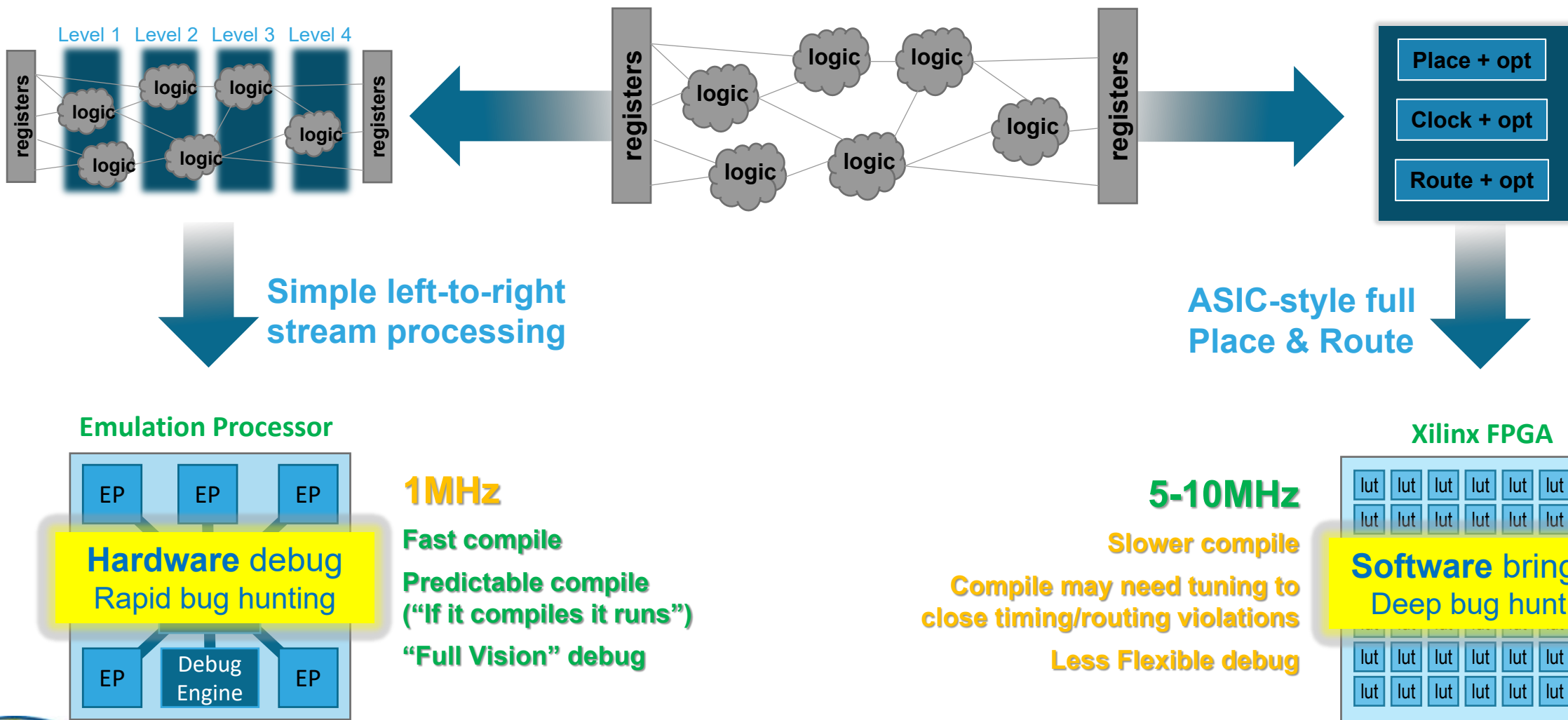
Fast RTL
verification &
debug

Prototyping “Debug your Software”

- **Build time and debug less important:** design 1-2 weeks
- **Highest performance:** software debug simulation runs
- **Lowest cost:** replicate one build for developers
- **SOC level capacity:** more and more software at SOC level

Early SW
development
& HW
regressions

Processor Based Emulation and FPGA Based Prototyping



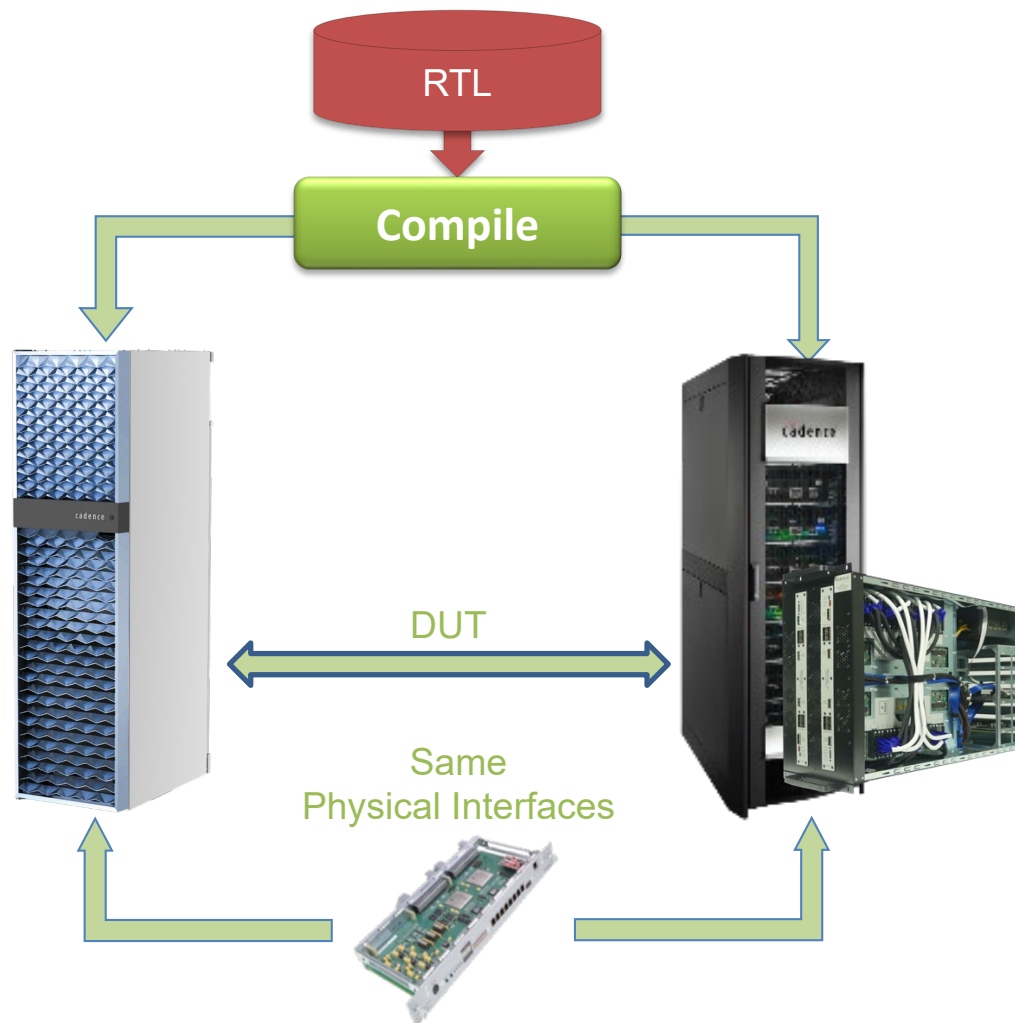
Dynamic Duo: Emulation & Prototyping

Emulation

- Optimized HW/SW debug
- SoC acceleration, HW/SW
- Power & Performance Analysis
- Advanced Use Models

Prototyping

- Automated Bring-Up
- Scalable performance
- SW development
- HW/SW regressions



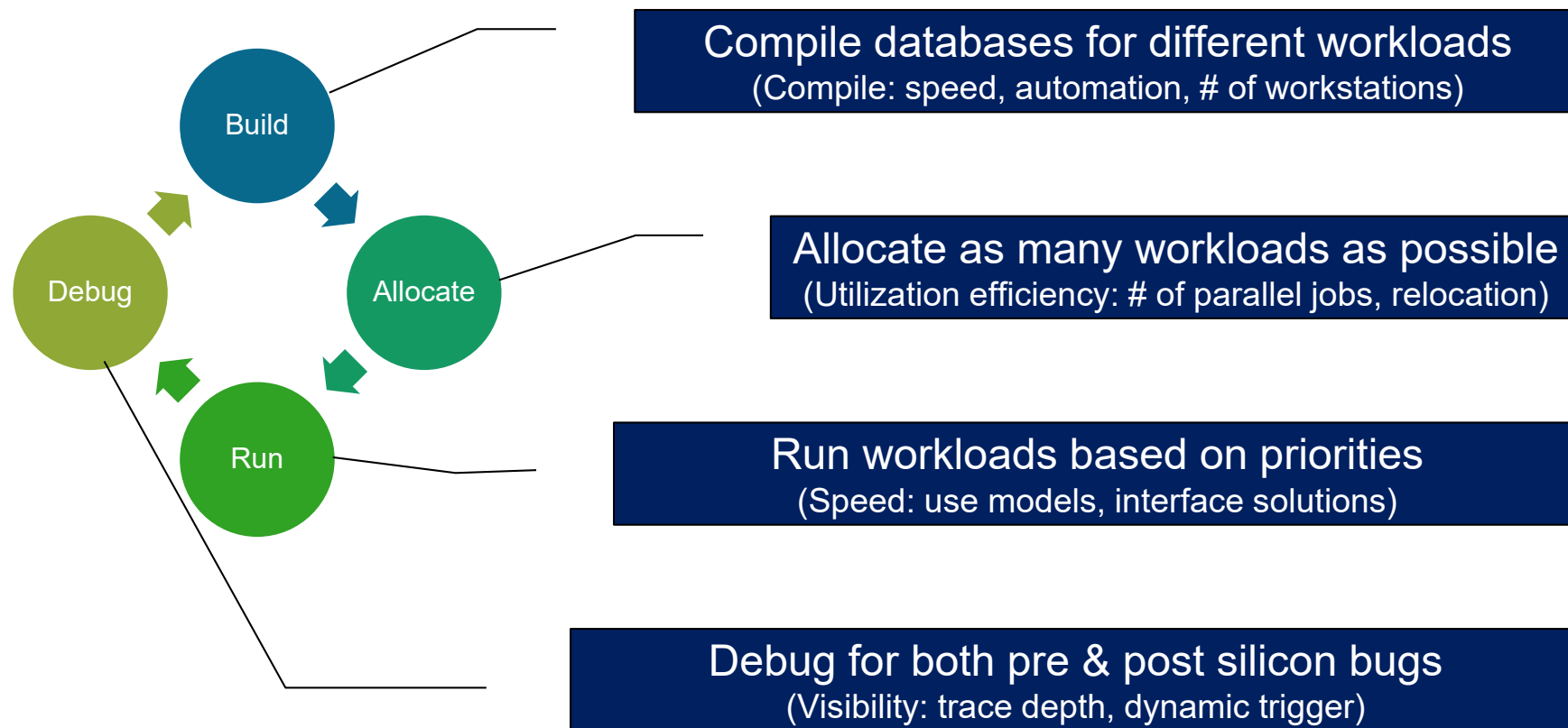
Congruency and common environment

Emulation Capabilities

- Palladium® Z1 **enterprise** emulation platform
 - Up to **5X** greater emulation **throughput**
- **Scalability** from IP blocks to full systems on chip
 - Capacity of up to **9.2 billion gates** with **2304 users**
- Best in class **total cost of ownership** (TCO)
 - **22 use models**
- New era of **datacenter-class** emulation
 - **Proven reliability**



HW-assisted verification productivity loop



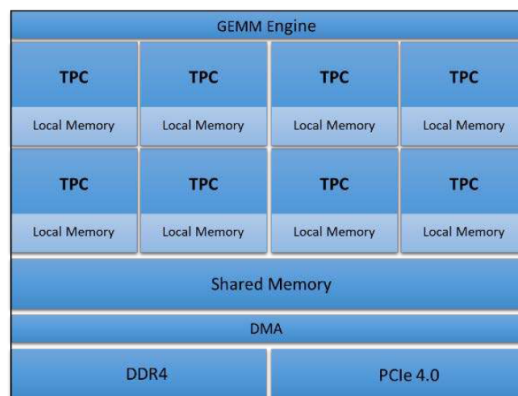
Billion Gate Design Examples

Habana Labs

Billion-gate class inference processor: Goya

Goya Processor Architecture

- Heterogenous compute architecture
 - 3 Engines: TPC, GEMM and DMA
 - Work concurrently using a shared SRAM
- Tensor Processor Core (TPC™)
 - VLIW SIMD vector core
 - C-programmable
- GEMM operations engine
- Tensor addressing
- Robust to any address stride
- Latency hiding capabilities
- PCIe Gen4.0 x16
- 2 DDR4 channels @ 2.667 GT/s, 40GB/s BW, 16GB capacity
- Dedicated HW and TPC ISA for special functions acceleration (e.g. Sigmoid/GeLU, Tanh)
- Mixed-precision data types: FP32, INT32, INT16, INT8, UINT32, UINT16, UINT8



TSMC – 16nm

Source: Hotchip Conference

Fujitsu

Billion-gate class HPC/AI Processor: A64FX

A64FX Chip Overview

Architecture Features

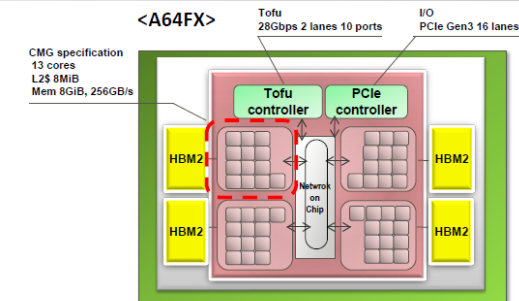
- Armv8.2-A (AArch64 only)
- SVE 512-bit wide SIMD
- 48 computing cores + 4 assistant cores*
*All the cores are identical
- HBM2 32GiB
- Tofu 6D Mesh/Torus 28Gbps x 2 lanes x 10 ports
- PCIe Gen3 16 lanes

7nm FinFET

- 8,786M transistors
- 594 package signal pins

Peak Performance (Efficiency)

- >2.7TFLOPS (>90%@DGEMM)
- Memory B/W 1024GB/s (>80%@Stream Triad)



	A64FX (Post-K)	SPARC64 Xifx (PRIMEHPC FX100)
ISA (Base)	Armv8.2-A	SPARC-V9
ISA (Extension)	SVE	HPC-ACE2
Process Node	7nm	20nm
Peak Performance	>2.7TFLOPS	1.1TFLOPS
SIMD	512-bit	256-bit
# of Cores	48+4	32+2
Memory	HBM2	HMC
Memory Peak B/W	1024GB/s	240GB/s x2 (in/out)

All Rights Reserved. Copyright © FUJITSU LIMITED 2018

Source: Hotchip Conference

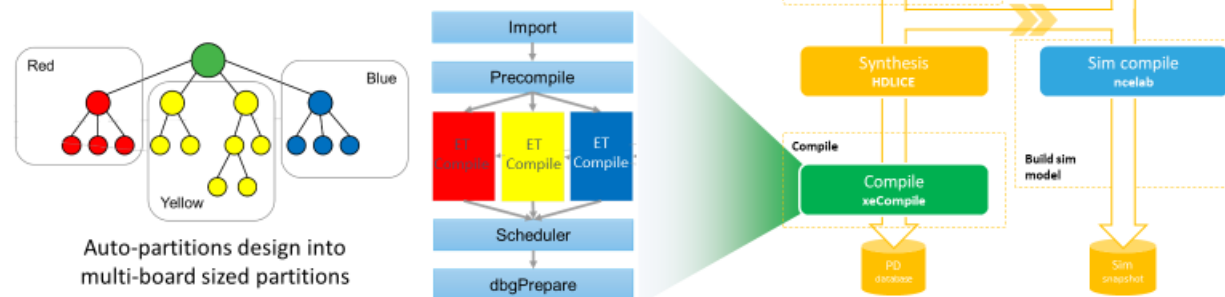
Emulation Model Scaling: Enabling AI/ML/5G



- Scaling from IP verification to multi-core/die emulation
 - 4 million gates to 7 billion gates
- HW/SW co-verification without performance impact
 - FullVision, waveform streaming, dynamic RTL, physical and virtual JTAG debuggers
 - Supports 3rd party 5G testers
- Emulation in the cloud eases variability of workload sizes and use models

Advanced Parallel Partition Compiler (PPC)

- 2nd generation automatic partitioning PPC
 - Improve compile efficiency by 60%
 - Generates better runtime and better utilization vs traditional
 - Partitions are auto-created from instrumented gate count



- Parallel Partition Compiler enables one billion gate emulation model to be compiled in about 4 hours
- 2nd generation PPC-based emulation models are fully automated – no need for manual partitioning
- PPC enables practical scaling of emulating billion-gate class AI/ML and 5G designs

INT64 to INT8 Computational Efficiency

Fujitsu

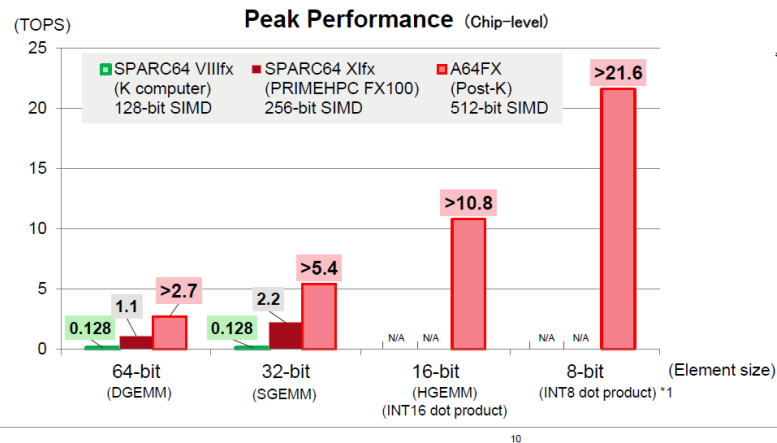
Billion-gate class HPC/AI Processor: A64FX

Execution Unit

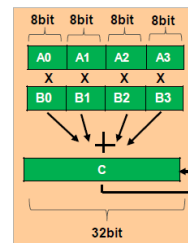
FUJITSU

Extremely high throughput

- 512-bit wide SIMD x 2 Pipelines x 48 Cores
- >90% execution efficiency in (D|S|H)GEMM and INT16/8 dot product



*1 INT8 dot product
 $C = \sum (A_i \times B_i) + C$



- For certain ML applications, accuracy may be traded off for faster and more power efficient implementation methods
- Depending on the target applications, design may scale down (INT8) or scale-across (INT 8 to INT 64) architecturally to ensure computational efficiency

Source: Hotchip Conference

Power / Performance Trade-offs

Fujitsu

Billion-gate class HPC/AI Processor: A64FX

Power Management (Cont.)

FUJITSU

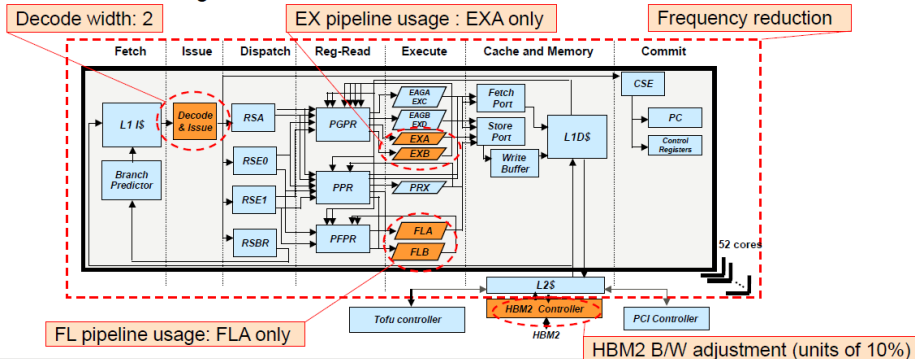
■ “Power knob” for power optimization

■ A64FX provides power management function called “Power Knob”

- Applications can change hardware configurations for power optimization

- Power knobs and Energy monitor/analyzer will help users to optimize power consumption of their applications

<A64FX Power Knob Diagram>

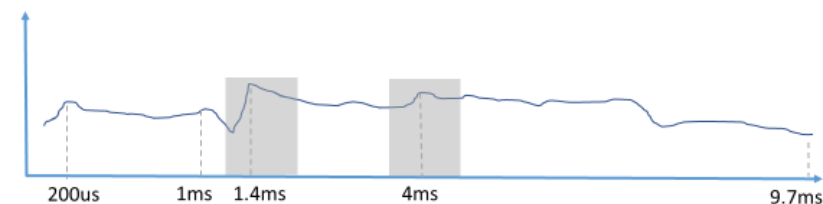


Source: Hotchip Conference

Customer Example

HW Toggle Count Case Study

Customer Example 1

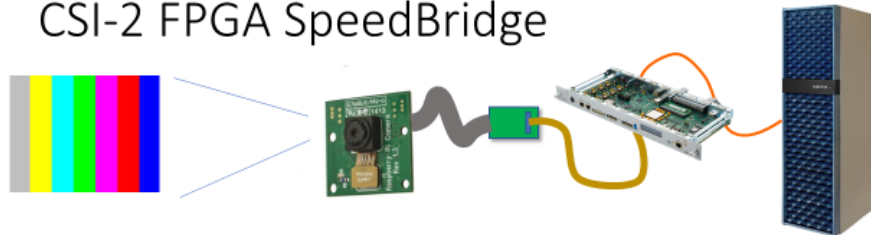


- Pre-silicon peak detection is critical in analyzing and correlating power and performance
- HW-WTC is 304x faster compared to SW-WTC (34.7 minutes vs 7 days) – enabling user to detect peaks faster and perform system analysis and design optimization

Source: Cadence

Senor Models – MIPI CSI-2

CSI-2 FPGA SpeedBridge

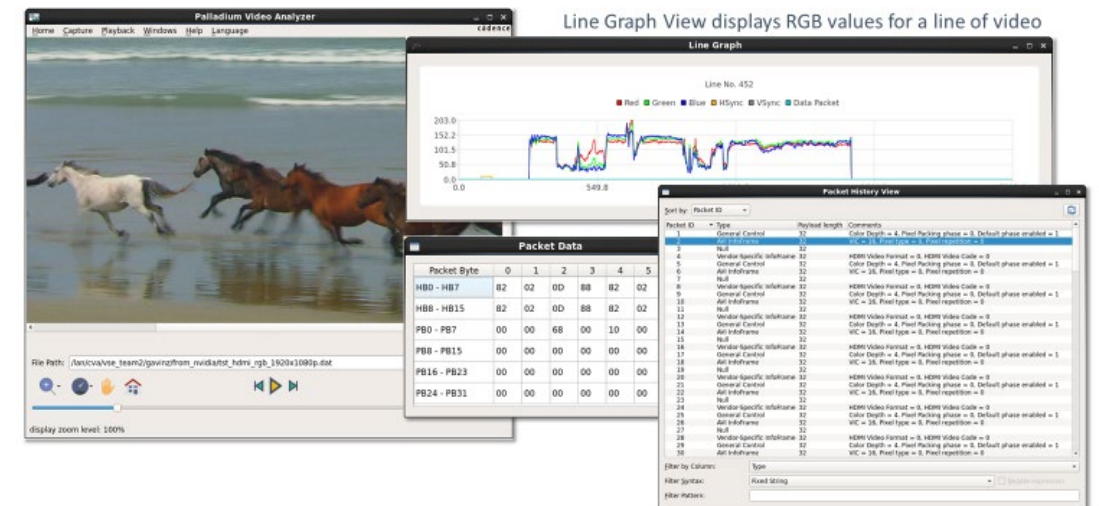


- Connects to real CSI-2 sensors for system-level design emulation
- Sensor video is scaled to emulation speed
- Users can verify their DUT with the actual sensor
 - Vendor specific register spaces, image processing functions, etc
 - Vendor specific embedded data in image packet stream
 - Real sensor behavior
- Daughtercards provide connectivity to different types of camera modules

SpeedBridge Capabilities:

- One sensor per SpeedBridge board
- CSI-2 RX only
- Sensor side up to 4 CSI-2 lanes
- Support max speed of 1.5Gbps
- Emulation side PPI interface up to 4 lanes
- Emulation side CCI interface
- Compatible with Palladium Video Analyzer for CSI stream viewing

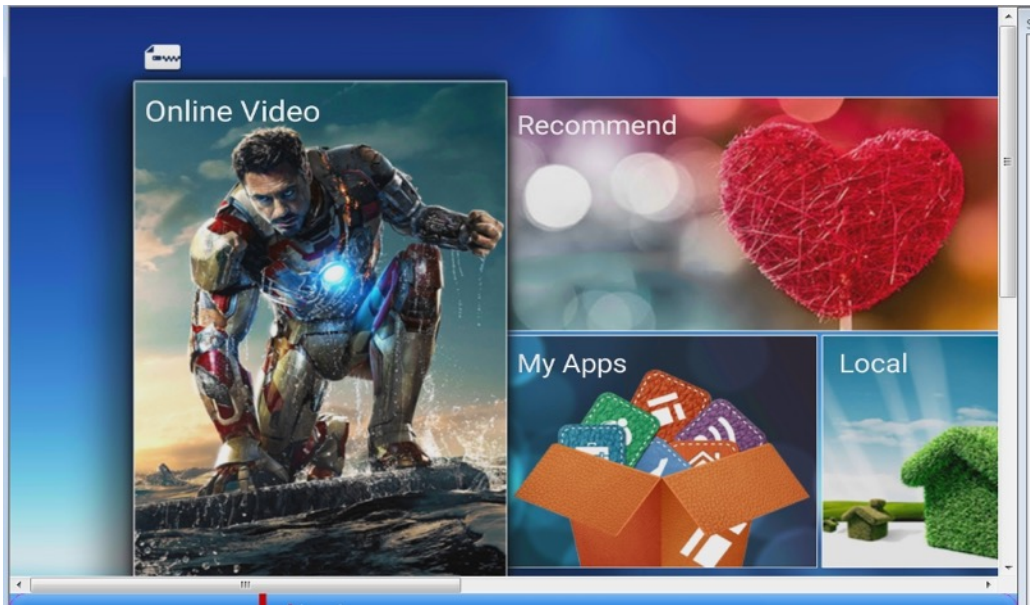
Palladium Video Analyzer Views



- Emulation with real CSI-2 sensor provides realistic live image capture and processing
- Allows user to conduct visual inspection (e.g. fast forward and replay)
- Video frames can be captured for detailed analysis (e.g. data packet, line graphs, etc.)

FPGA-based prototyping accelerates time to revenue

- Early, embedded software & firmware development
- Initial systems and/or proof of concept
- Pre-silicon chip (ASIC) verification

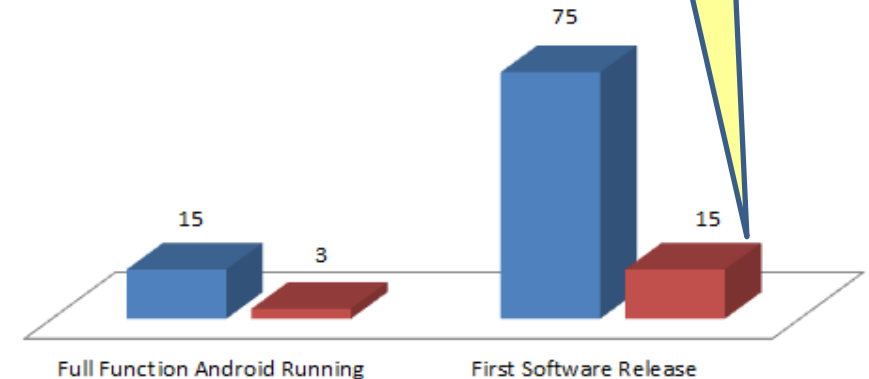


Software is ready when 1st silicon comes in

- Full Android boot 30 min after silicon is back
- Full system demo to customer in 3 days

Software Development Effort (# of days)

■ No Pre-Silicon System ■ With Pre-Silicon System



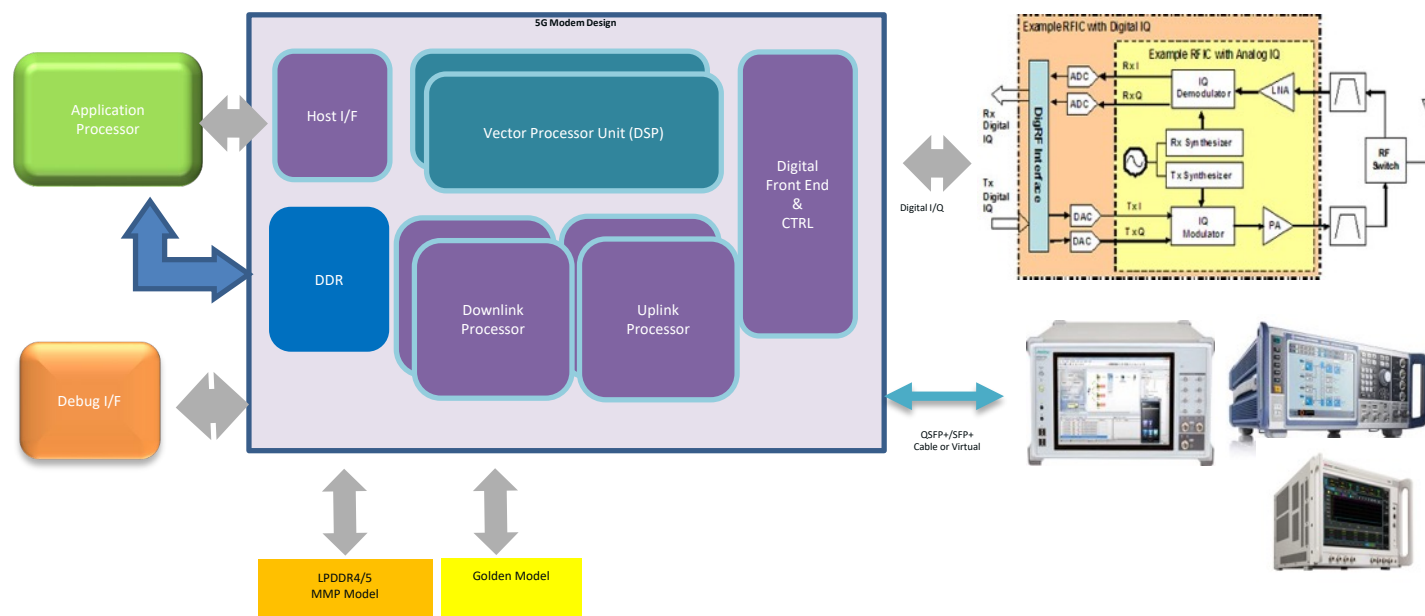
2 months faster
time to market!

Source: Amlogic

But ... your Mom's and Dad's prototyping
doesn't cut it any more



- Too many gates
- Too much memory
- Too many peripherals
- Too much software
- ...
- And not enough time



Protium X1 Enterprise Prototyping System



- **Performance**
 - Enabling early firmware and software development, automated bring-up
 - Up to 100MHz for single FPGA; up to 5MHz on billion gate designs
- **Capacity**
 - Advanced blade architecture scales to billions of gates
 - Ideal for AI, ML, 5G, mobile, and graphics applications
- **Fast Bring-up**
 - Unified Palladium® Z1 / Protium™ X1 compile ensures DUT congruency
 - Enables transition from emulation to prototyping in days
- **Multi-user**
 - Single-FPGA granularity assures high utilization and efficiency
 - Ideal for storage, automotive, image, consumer and medical applications

Scalable Capacity



- Blade architecture: scalability and flexibility
 - Each blade (of up to 150M gates) self-contained
 - Can be used as individual desktop system
 - Up to 8 blades mounted into standard 19" rack (1.2BG per rack)
- Racks connected for multi-billion gate prototyping
 - AI, 5G, graphic and mobile designs
- GUI-based, interactive configuration assistant
 - Create optimal configuration for user needs

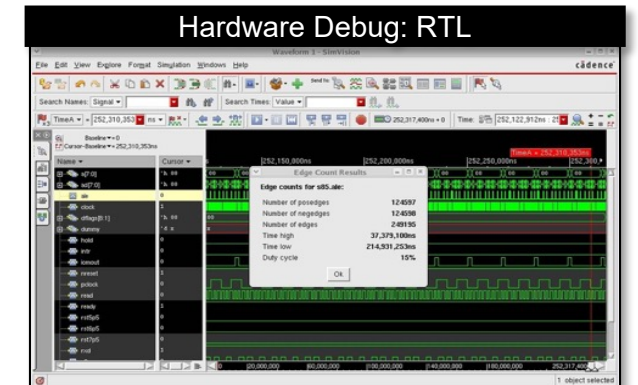
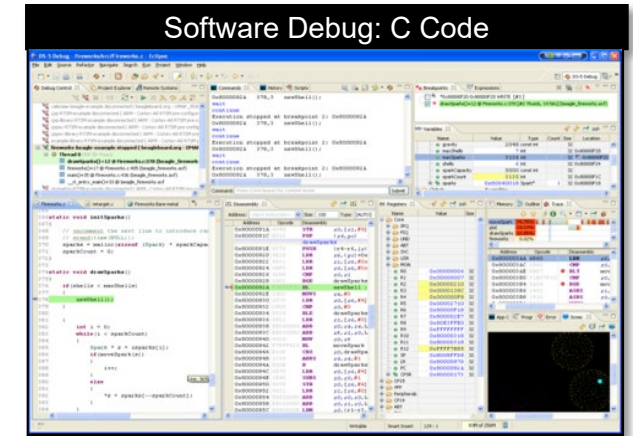
Scalable Performance

- New fully-automatic partitioning, technology mapping algorithms for best possible performance regardless of design size
 - New Pathfinder multi-FPGA partitioner
 - New multi-strategy, high-speed pin-multiplexing
- Manual optimization capabilities for higher performance up to 100MHz+
 - Black-boxing: native, high-speed interfaces
 - Manual partition guidance: sub-systems
 - Inter-FPGA hardware optimization to customize connectivity



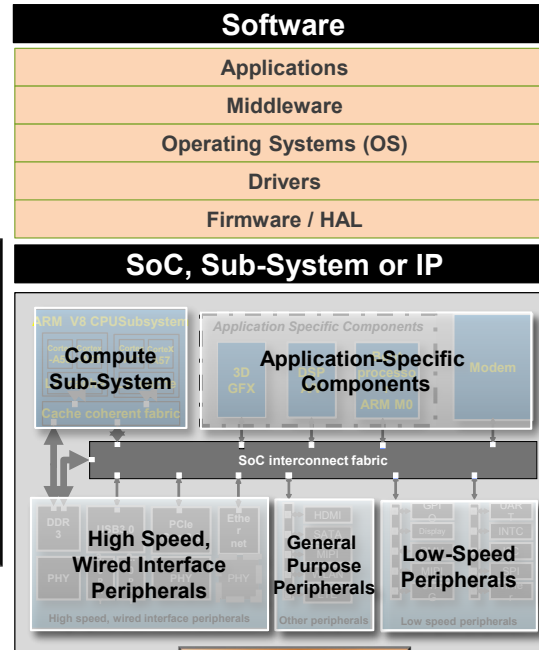
Advanced Debug

- **Software debug:** early firmware and software development
 - Memory (backdoor) upload and download
 - Clock control to stop and resume the hardware at any time
- Standard interfaces to industry-leading debuggers and software environments - use familiar tools
 - Joint Test Action Group (JTAG) and Universal Asynchronous Receiver/Transmitter (UART) interfaces
- Transaction interface
 - Directly connect to software models and virtual environments
- **Hardware debug:** bring-up design quickly, validate functionality
 - Force and release for internal notes
 - Prototyping full visibility
- Data capture card (DCC): 1000's of signals, millions of cycles
- Assertion checkers: automation of remote test execution

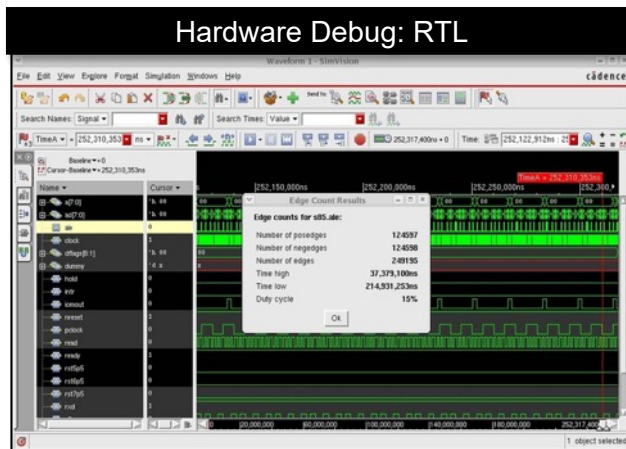


Advanced Prototyping Debug

- Waveforms **across partitions** - Design-centric view
- **Force/Release** - Predefined signals to "0" or "1" **during runtime**
- **Real-time monitoring** of predefined (at compile time) signals
- **Data Capture** - Thousands of signals for millions of cycles
- **Prototyping Full Visibility** Probe without recompile
- **Assertion Checkers** - Non-intrusive hardware monitors



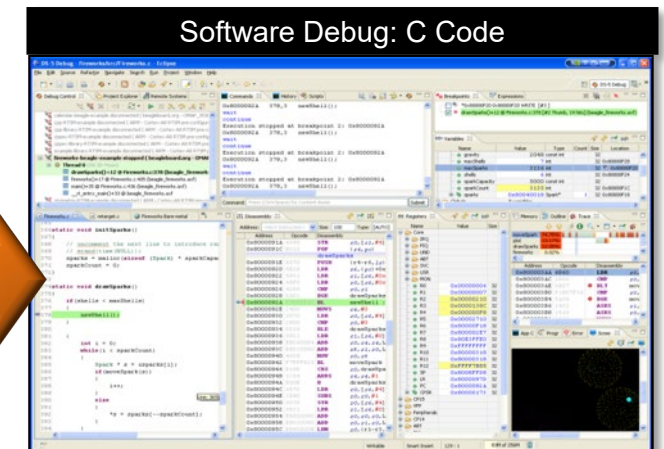
- **Backdoor memory** access - change boot code, software, etc.
- Clock **control** - **Start/Stop** the clock on demand
- Fully **scriptable** runtime environment
- **Remote access** - Network resource anytime from anywhere
- High-performance link to **software model**



Probes



JTAG



Daughter Cards & Peripherals

DeepChip – “Best of 2019”

The users have spoken!

- The Protium users also gushed a lot about the fact that they could *go back to Palladium* for fast debug and waves if needed.
 - “We run our design on Protium then go back to Palladium for debug. With Palladium, we can capture waves up and down the hierarchy of every net in our chip. It’s a really big advantage of Protium.”
 - “With Protium, we get the speed of an FPGA-based system, with the fast ramp-up, debug, and signal traces of Palladium. It only takes seconds for us to see all the waveforms.”
- Protium took 1.2 to 2.1 days to recompile vs. Synopsys HAPS taking 3.9 to 6.1 days to recompile.



And that's why Protium (actually the crazy fast incremental Protium compiles with FPGA 8.3 Mhz simulation speeds) wins the #1 Best of EDA in 2019 award from the end users this year.

<div> <div>What the EDA users REALLY think</div> <div>DEEPCHIP</div> <div>cadence Spectre[®] X Simulator</div> </div>		
(DAC'19 Item 1a) ----- [12/19/19]		
Subject: CDNS Protium crazy fast "Palladium-compiles" #1a for Best of 2019		
FAST COMPILES ROCK!: My quick-and-dirty summary of the emulator/prototyper world. Say you have two designs to simulate. One design is 200 million gates, the other design is 1 Billion gates.		
	Initial Ramp Up Time / Incremental Compile Time	Operating Speed
Palladium 200 M gates 1.0 B gates	initial ramp 2-4 weeks 1.0 hour 5.0 hours	1.2 Mhz 800 Khz
Zebu Server 4 200 M gates 1.0 B gates	initial ramp 4-6 weeks 25.8 hours (1.1 days) 41.2 hours (1.7 days)	2.0 Mhz 750 Khz
HAPS-80 200 M gates 1.0 B gates	initial ramp 2-3 months 93.6 hours (3.9 days) 146.4 hours (6.1 days)	20.0 Mhz 5.0 Mhz
Veloce Strato 200 M gates 1.0 B gates	initial ramp 3-5 weeks 5.1 hours 12.5 hours	1.6 Mhz 750 Khz
Protium 200 M gates 1.0 B gates	ramp 24 hours w/Palladium ramp 4-6 weeks w/o Palladium 28.8 hours (1.2 days) 50.4 hours (2.1 days)	8.3 Mhz 4.5 Mhz

Test Solutions for Emulation / Prototyping

- Next Gen 5G testers: serial high speed link
 - Due to higher 5G bandwidth requirements
 - Better physical constraints, standards based
 - 25G/40G Ethernet through QSFP/SFP+ fiber



- Off the shelf solution needed: rate adapt, convert to IQ data
 - Testers no longer slowdown parallel IQ data
 - Need to packetize/de-packetize Ethernet traffic
 - Need triggering method to rate adapt traffic flow

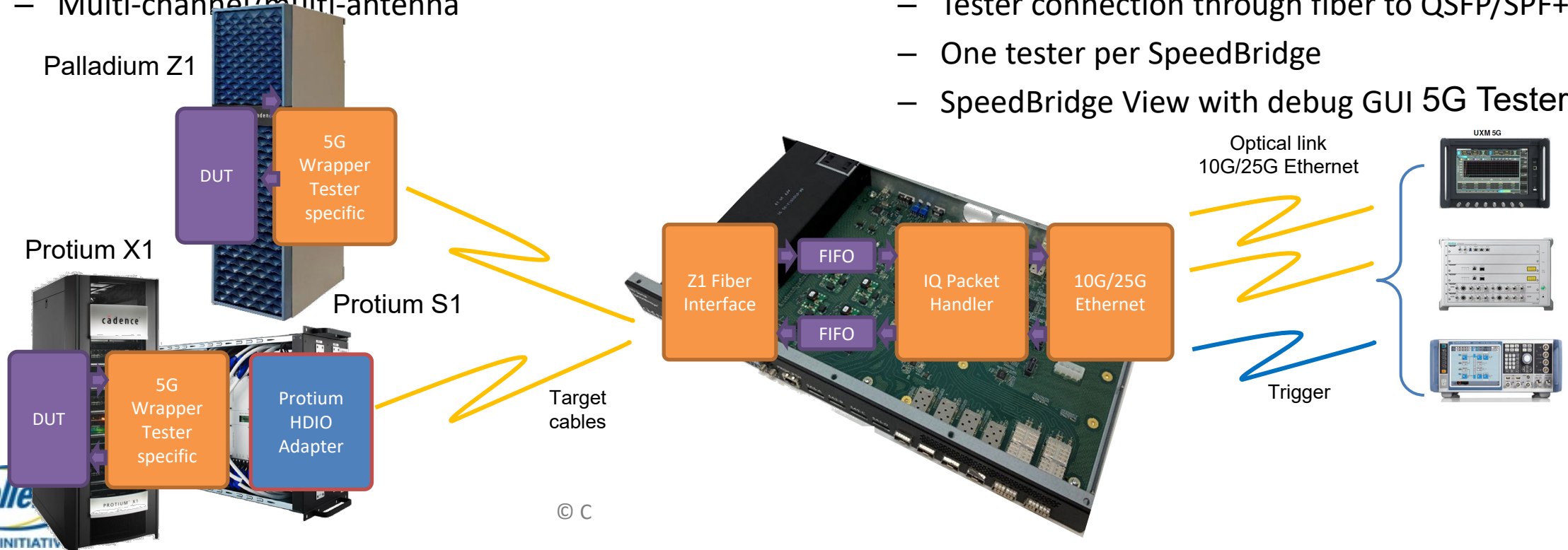
HD 5G Rate Adapter

• Capabilities

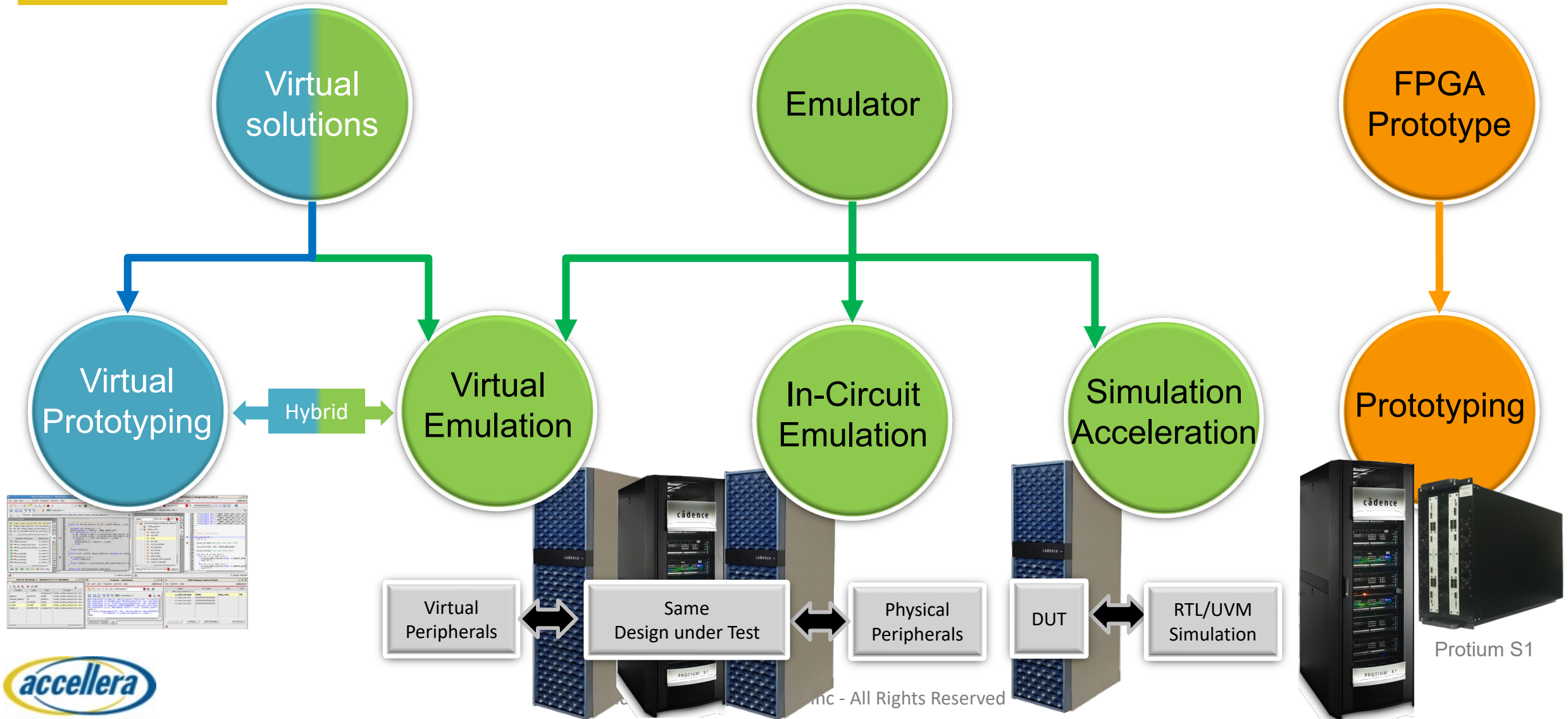
- R&S and Keysight
- Palladium Z1, Protium S1/X1
- Slow 5G IQ data to 5G Handset DUT
- Multi-channel/multi-antenna

• Features

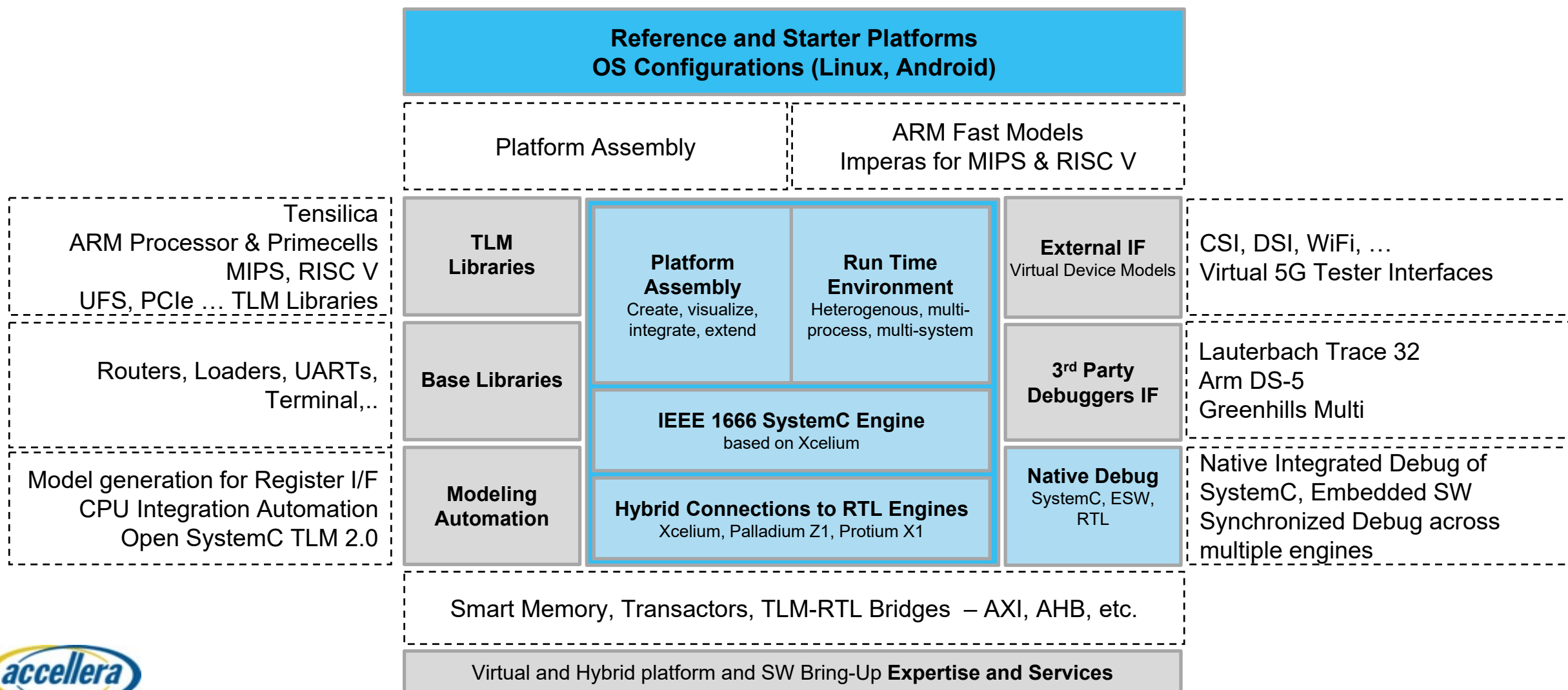
- Direct fiber connection (Z1)
- HDIO adapter (Protium S1/X1)
- 1U 19" Rack mountable
- Tester connection through fiber to QSFP/SPF+
- One tester per SpeedBridge
- SpeedBridge View with debug GUI 5G Tester



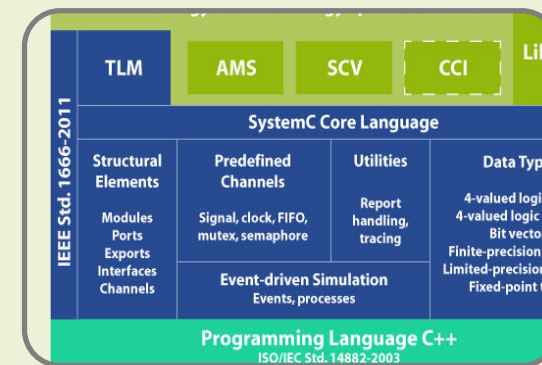
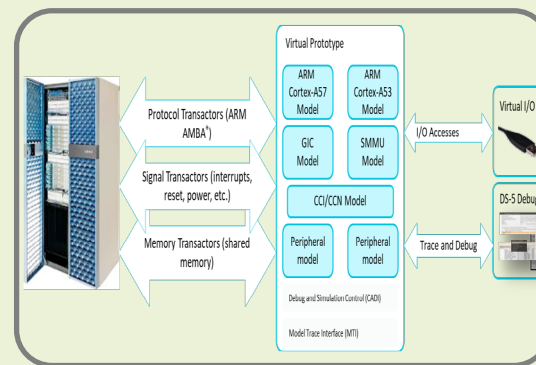
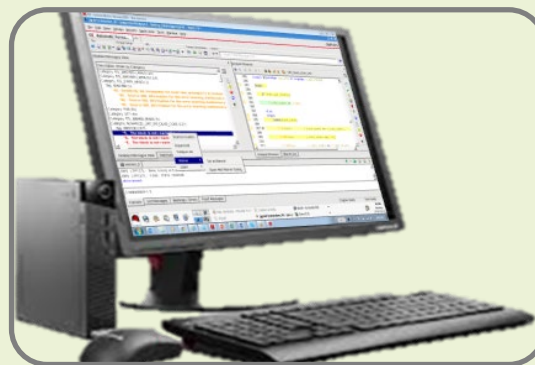
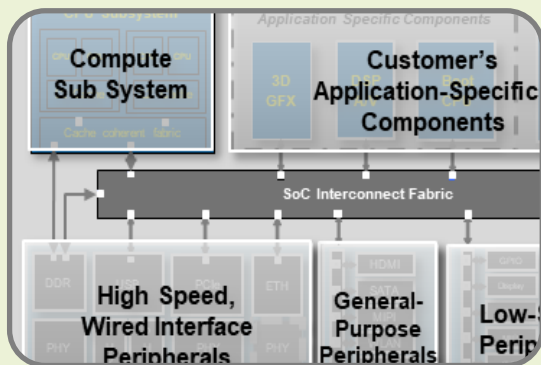
Virtual, Emulation & Prototyping



Virtual Prototyping and Hybrid Environment



Use Models



Architecture Analysis

Accuracy, Accuracy
Best achieved "RTL-up",
auto generation for
interconnect

AT modeling requires
careful effort-accuracy-
performance
considerations

Software Development

Speed, Speed, Speed
Loosely timed
"just enough" accuracy is
sufficient

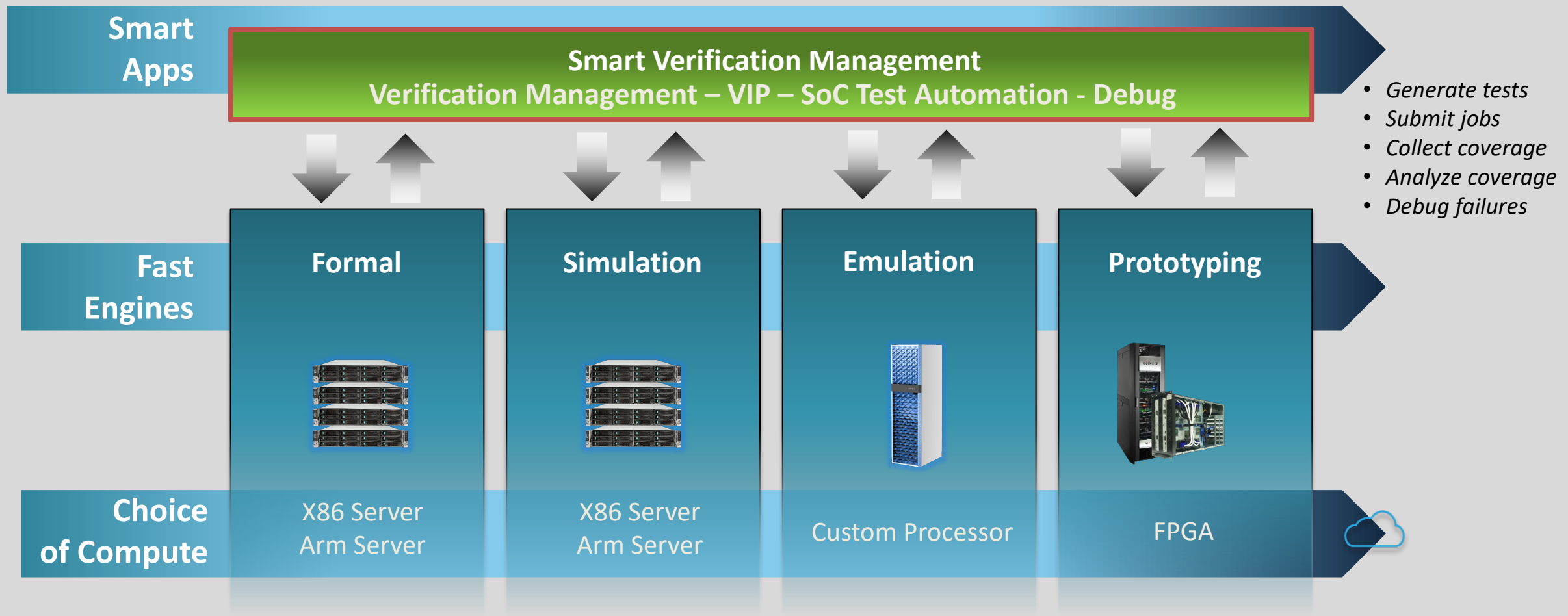
Mixed Fidelity Hybrids

Speed & Fidelity
Details in hardware
Keep processor sub-
system and other
peripherals virtual
Smart synchronization
between TLM and RTL

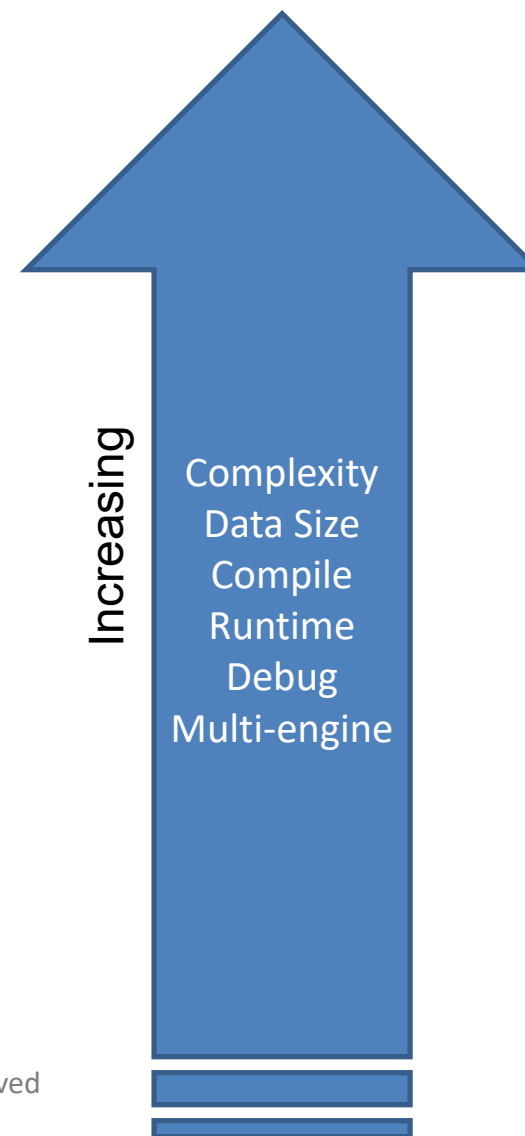
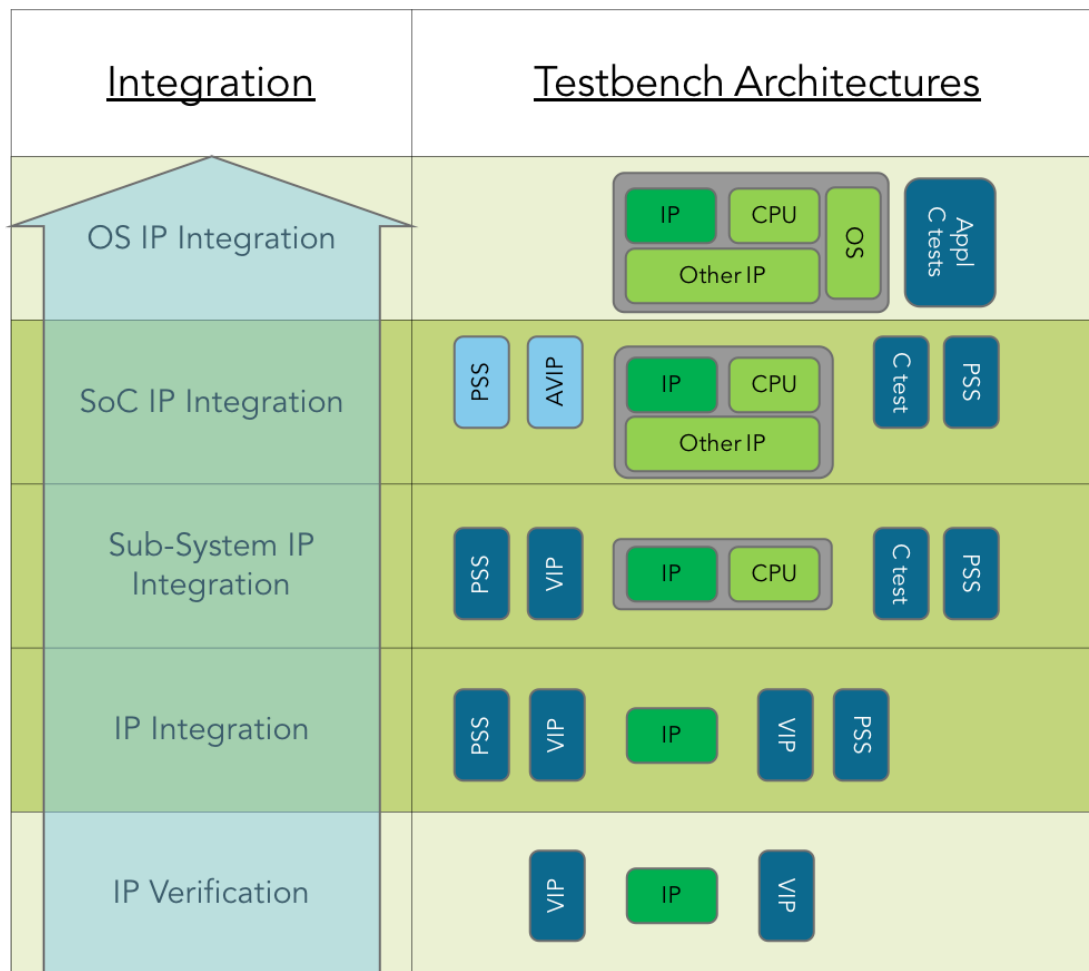
Hardware Verification

SystemC to drive
hardware DUT

Smart Verification Management



The Problem



Formal Simulation Emulation Prototyping

© Cadence Design Systems, Inc - All Rights Reserved

Verification Management: Data-Driven

Use-case-based

- Define legal operations
- Workload matters: must represent real operation

Data Collection

- Non-intrusive data collection
- Use the right execution platform

Analysis

- Correlate, filter, learn, predict
- Anomaly Detection

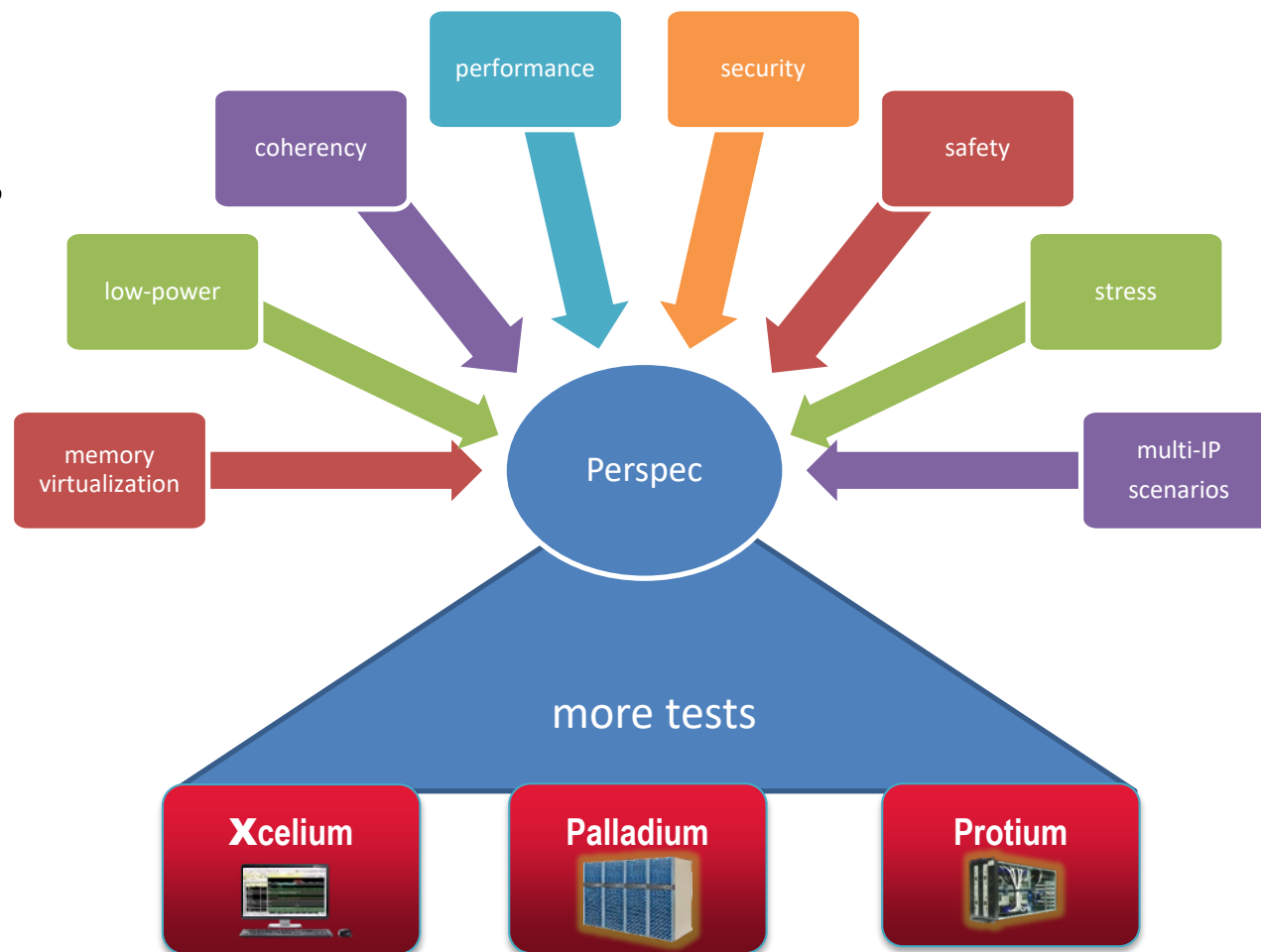
Goal-based

- Verification throughput
- Smarter bug hunting

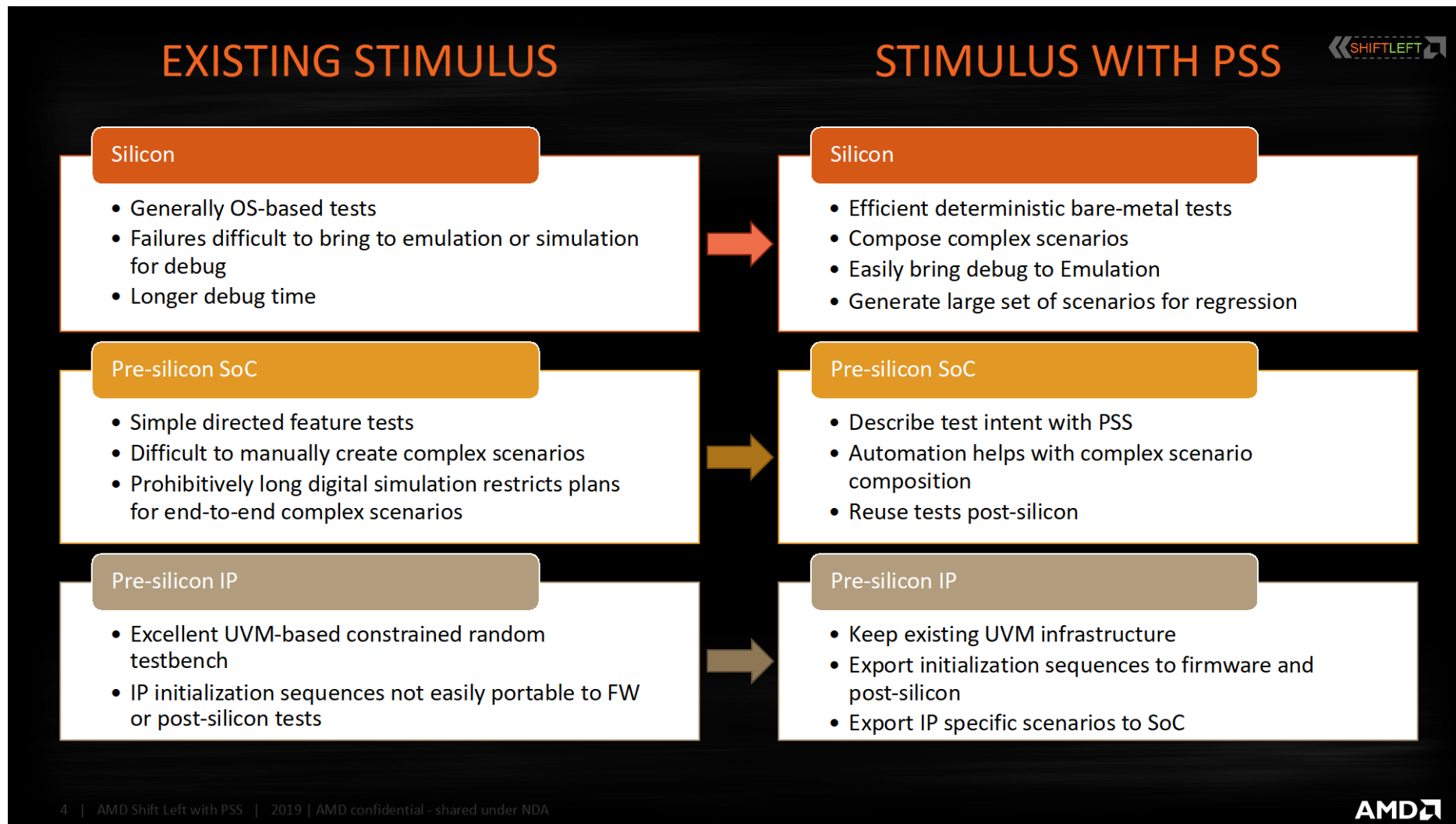
Use-case Based Test Generation

Accellera Portable Stimulus Standard

- Describe Test Intent and Design Behaviors
 - Use-cases
 - Legal scenario space
- Deliver Test Portability
 - Vertical reuse: From IP to SoC
 - Horizontal reuse: from Simulation to Emulation to Post Silicon
- Across Users
 - Abstract modeling
 - Actions, inputs, outputs, resources



PSS Impact on Stimulus



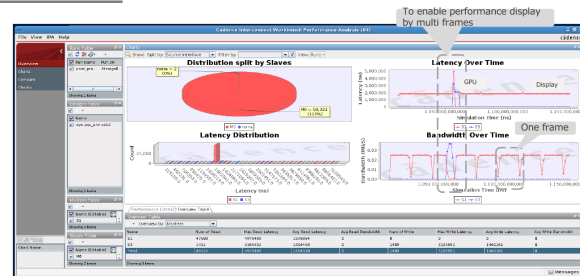
Renesas Performance Verification with Perspec Generated Use-cases

Application Example of IWB/Palladium for Performance Verification

2016.07.15

第一要素技術事業部 デザインメソッド部
機能検証技術課
Mr Makoto Matsumoto
Renesas System Design Corporation

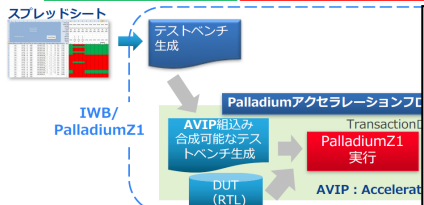
IWB/PalladiumZ1 Analysis (1)



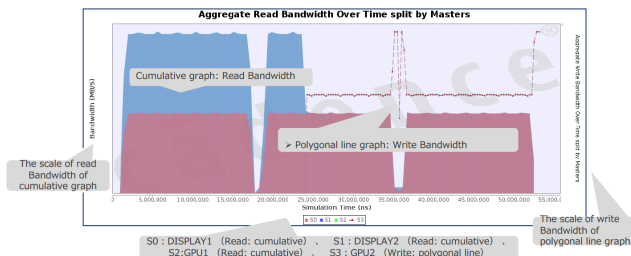
IWB/PalladiumZ1適用フロー

- IWBから生成されるテストベンチを用いPalladiumZ1での実行
- PalladiumZ1組み込みのAVIPからTransactionログを出力。Performanceログへ変換
- PerformanceログからIWBにて性能表示

環境構築 → 性能測定 → 性能解析



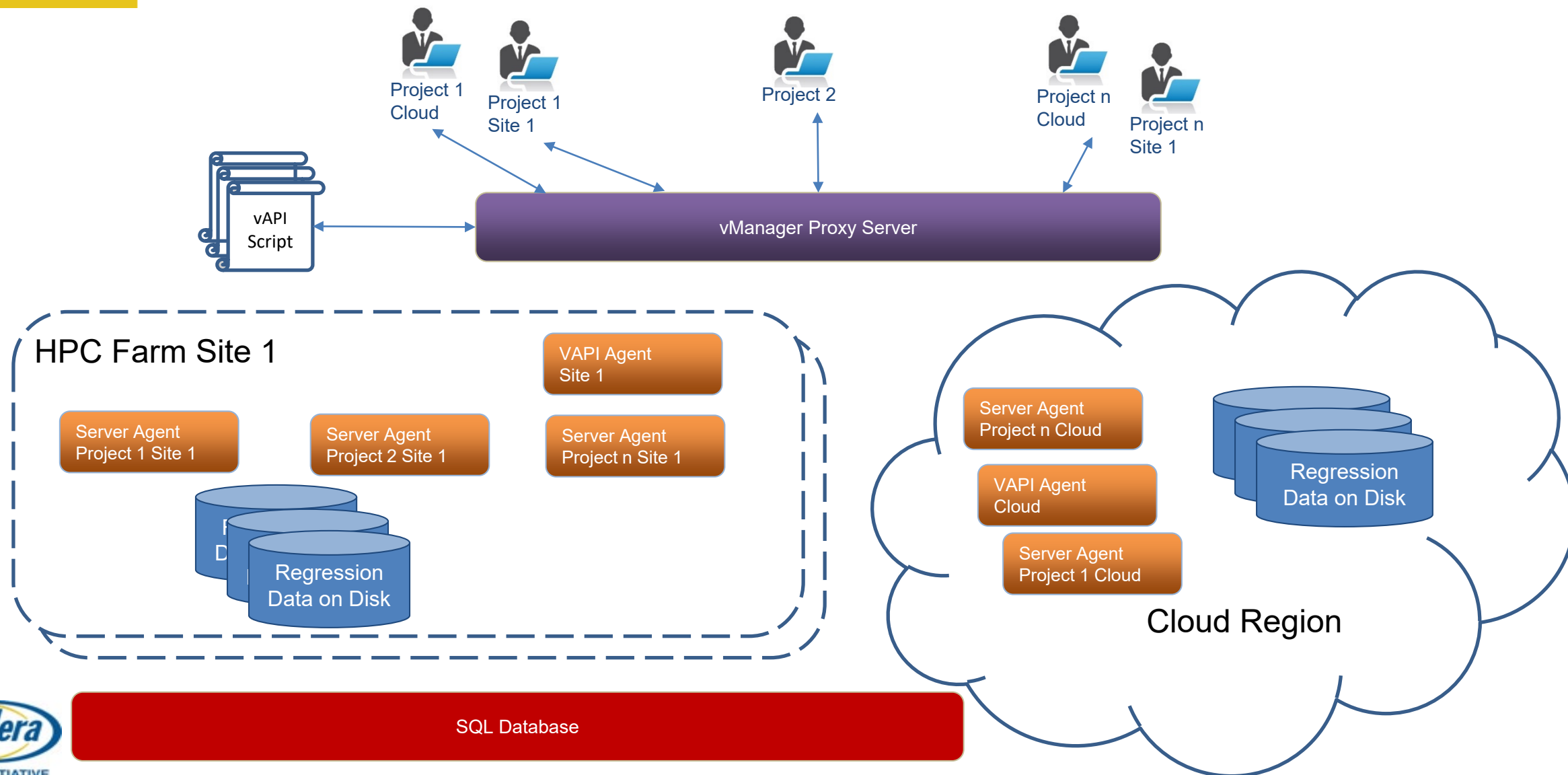
IWB/PalladiumZ1 Analysis (2)



- Leading industrial and automotive MCUs
 - Number of integrated IPs is increasing
 - Switched interconnect
 - Configuration has big impact on performance
- Interconnect Workbench performance analysis
 - Early performance characterization
 - Interconnect tuning to optimize performance
 - Use case performance validation
- Palladium Z1 with Perspec use cases
 - Bring-up the entire design and software
 - Perspec generating use case tests
 - Reduce from 50 hour simulations to 12 minutes

-
- The diagram illustrates the Accellera Perspec architecture for SOC-level tests. At the top, the **accellera** logo is shown. Below it, a green document icon represents the **Portable Stimulus Specification (PSS)** for Scenario-based **SOC-level tests**. A blue arrow points down to a blue box labeled **Perspec**. From **Perspec**, four blue arrows point down to four red boxes representing different simulation engines: **JasperGold**, **Xcelium**, **Palladium**, and **Protium**. Each box contains a small image of the respective engine's interface or hardware. To the left of these boxes, the word **Design** is written in red, with two document icons and a blue arrow pointing towards the simulation engines. Below the simulation engines, a blue arrow points down to a blue box labeled **vManager**. To the right of the simulation engines, a blue arrow points down to a blue box labeled **Silicon**, which contains a small image of a silicon chip. Below **Silicon**, the text **SOC-level Coverage vs. Test Plan** is written in black. A blue arrow points from **Silicon** to a screenshot of a software interface showing a table of test results. The table has columns for Name, Overall Average Grade, Status Grade, Assertion Passed, and Assertion Failed. The table lists various test cases and their corresponding grades and status.
- accellera
- Portable Stimulus Specification (PSS)
for Scenario-based **SOC-level** tests
- Perspec
- Design
- JasperGold
- Xcelium
- Palladium
- Protium
- Silicon
- SOC-level**
Coverage vs. Test Plan
- vManager
- WIP
- | Name | Overall Average Grade | Status Grade | Assertion Passed | Assertion Failed |
|--|-----------------------|--------------|------------------|------------------|
| 2.1 SWC | 56.47% | OK | 0 | 0 |
| 2.1.1 HW_flash_SIW integration (VIRTUAL) | 58.09% | OK | 0 | 0 |
| 2.1.2 SWC Case (SIW DRIVEN) | 56.47% | OK | 0 | 0 |
| 2.1 Subsystem | 51.72% | OK | 0 | 0 |
| 2.1.1 ESW (SIW ACCEL) | 69.62% | OK | 0 | 0 |
| 2.2 SWC Subsystem | 39.8% | OK | 0 | 0 |
| 2.2.1 Interconnect Verification (SIW) | 44.44% | OK | 0 | 0 |
| 2.2.2 IP Connectivity (FORMAL) | 8% | OK | 0 | 0 |
| 2.2.3 Power Intact (SIMULATION) | 54.95% | OK | 0 | 0 |
| 2.3 IP | 59.32% | OK | 0 | 0 |
| 3.1 UART (SIMULATION) | 73.98% | OK | 0 | 0 |
| 3.2 SWC FORMAL | 60.18% | OK | 0 | 0 |
| 3.3 PLL DRIVEN SIGNAL (SIW) | 55.27% | OK | 0 | 0 |
| 3.4 GPIO (SIMULATION) | 55.25% | OK | 0 | 0 |
| 3.5 SPI (SIMULATION) | 30.73% | OK | 0 | 0 |

Distributed Data: Centralized Management



vManager: Multi-engine Coverage

Combined Metrics

Xcellium

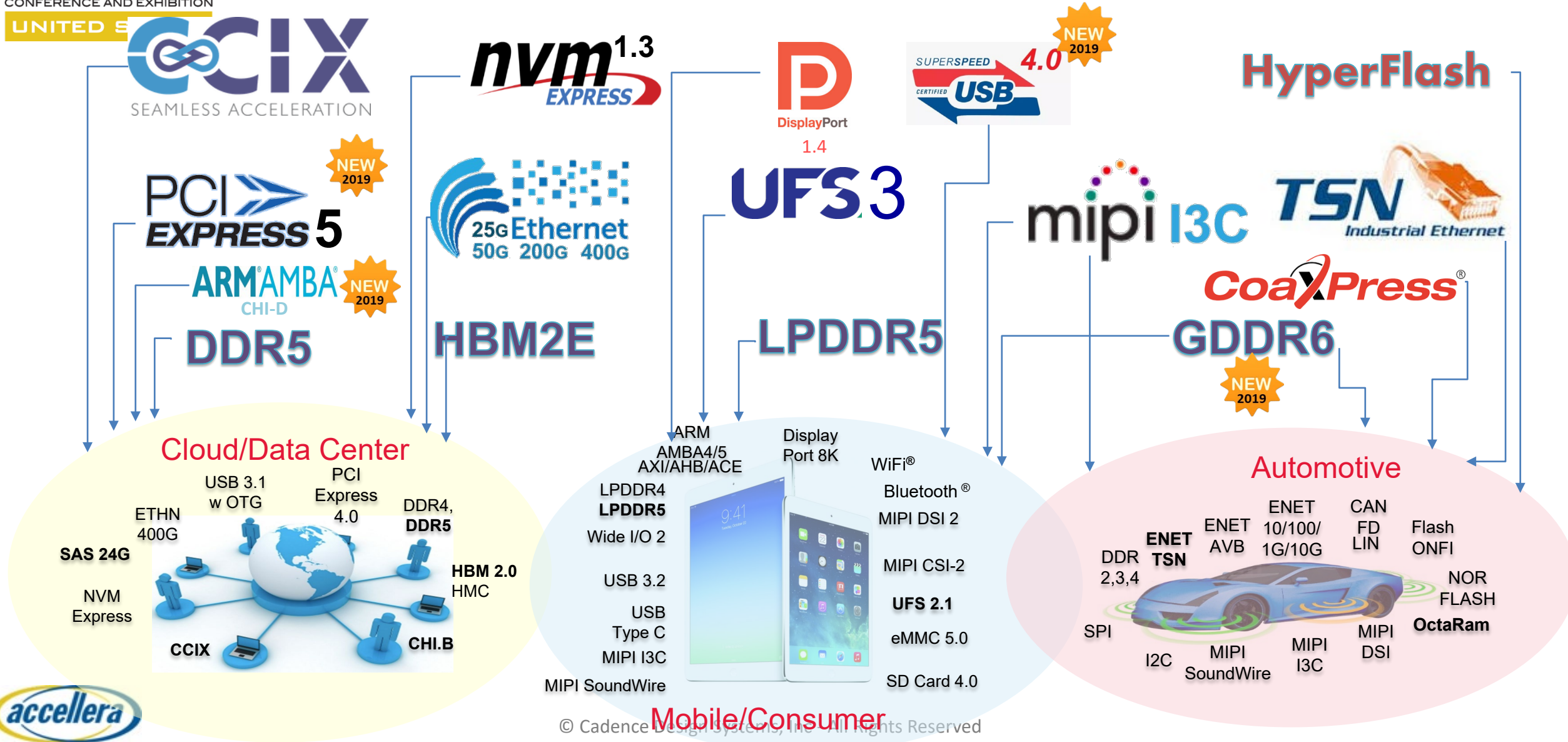
JasperGold

Verification Hierarchy									
default									
Ex	UNR	Name	Combined Average Grade	Combined Status Grade	Overall Average Grade	Assertion Status Grade	Formal Average Grade	Formal Status Grade	Valid Metrics
		(no filter)	(no filter)	(no filter)	(no filter)	(no filter)	(no filter)	(no filter)	[b,s,a]
		Verification Metrics	<div><div></div></div> 86.84%	<div><div></div></div> 44.34%	<div><div></div></div> 59.33%	<div><div></div></div> 45.31%	<div><div></div></div> 64.66%	<div><div></div></div> 19.38%	- - - - -
		Types	<div><div></div></div> 67.03%	<div><div></div></div> 45.31%	<div><div></div></div> 64.63%	<div><div></div></div> 45.31%	<div><div></div></div> 0%	<div><div></div></div> 0%	- - - - -
		Instances	<div><div></div></div> 95.77%	<div><div></div></div> 43.36%	<div><div></div></div> 54.02%	<div><div></div></div> 45.31%	<div><div></div></div> 90.96%	<div><div></div></div> 38.76%	- - - - -
		top_tgen_pvmem	<div><div></div></div> 95.77%	<div><div></div></div> 43.36%	<div><div></div></div> 54.02%	<div><div></div></div> 45.31%	<div><div></div></div> 90.96%	<div><div></div></div> 38.76%	- - t - -
		u_tgen_dut	<div><div></div></div> 97.86%	<div><div></div></div> 43.36%	<div><div></div></div> 63.76%	<div><div></div></div> 45.31%	<div><div></div></div> 95.6%	<div><div></div></div> 38.76%	b - t s - -
		u_apb_biu	<div><div></div></div> 100%	n/a	<div><div></div></div> 73.05%	n/a	<div><div></div></div> 100%	n/a	b - t s - -
		u_core_blk	<div><div></div></div> 99.25%	n/a	<div><div></div></div> 58.41%	n/a	<div><div></div></div> 99.25%	n/a	b - t s - -
		u_axi_biu	<div><div></div></div> 98.29%	n/a	<div><div></div></div> 66.34%	n/a	<div><div></div></div> 97.07%	n/a	b - t s - -
		u_chk_tgen	<div><div></div></div> 96.67%	<div><div></div></div> 36.92%	<div><div></div></div> 42.27%	<div><div></div></div> 44.62%	<div><div></div></div> 16.47%	<div><div></div></div> 22.22%	b - t - - a c
		u_axi_checker_tgen	<div><div></div></div> 83.84%	<div><div></div></div> 45.55%	<div><div></div></div> 82.96%	<div><div></div></div> 45.55%	<div><div></div></div> 0%	<div><div></div></div> 41.44%	b - t - - a -
		u_pvmem	<div><div></div></div> 90.03%	n/a	<div><div></div></div> 59.5%	n/a	<div><div></div></div> 79.52%	n/a	- - t - - -

vPlan Hierarchy					
Ext	UNR	Name	Overall Average Grade	User Entered Grade	User Calc Grade
		[-] V APB_UART	<div><div></div></div> 74.51%	n/a	<div><div></div></div> 72.22%
		[-] 1 APB_UART	<div><div></div></div> 67.89%	n/a	<div><div></div></div> 63.3%
		[-] 1.1 Interfaces	<div><div></div></div> 69.79%	n/a	<div><div></div></div> 69.79%
		[-] 1.2 Functional Features - B	n/a	n/a	<div><div></div></div> 49.5%
		[-] 1.2.1 blah	n/a	n/a	<div><div></div></div> 60.5%
		[-] 1.2.1.1 grade 1	n/a	<div><div></div></div> 77%	<div><div></div></div> 77%
		[-] 1.2.1.2 grade 2	n/a	<div><div></div></div> 44%	<div><div></div></div> 44%
		[-] 1.2.2 blah blah	n/a	n/a	<div><div></div></div> 38.5%
		[-] 1.2.2.1 grade 3	n/a	<div><div></div></div> 53%	<div><div></div></div> 53%
		[-] 1.2.2.2 grade 4	n/a	<div><div></div></div> 24%	<div><div></div></div> 24%
		[-] 1.2.2.3 grade 5	n/a	n/a	n/a

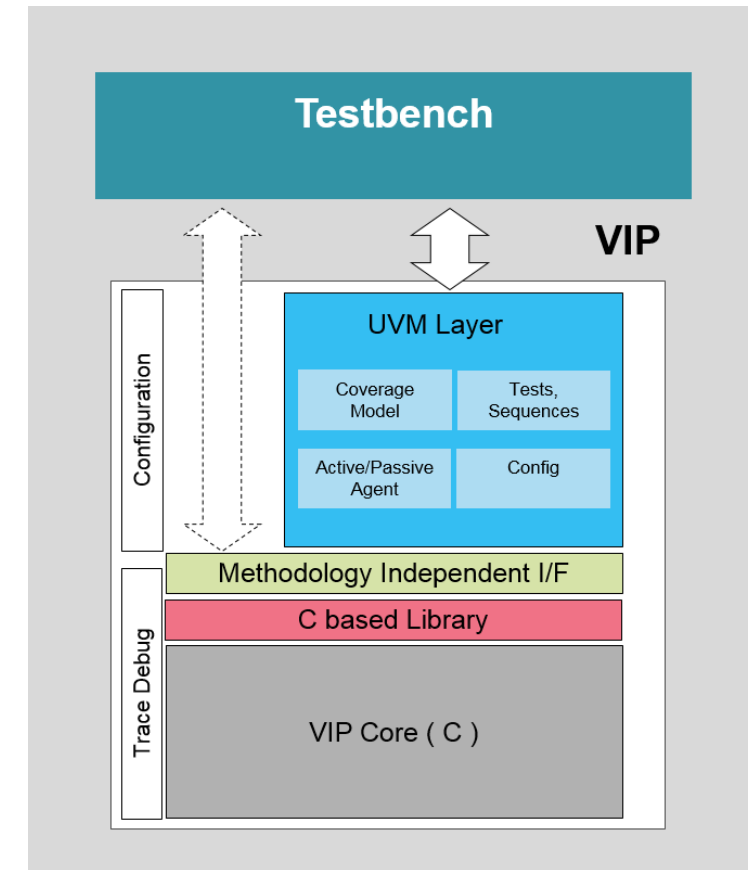
- Single click merge across regressions
- Multi-engine coverage
- User-defined grade calculation

VIP Catalog – Protocol Support



Cadence VIP Architecture

- Fast VIP
 - All VIPs implemented in C for the most optimized performance!
 - Uses dynamic memory allocation for all internal memories/registers/data structures
- Portable and Scalable architecture
 - Seamless IP=> SoC and Project => Project transition
- Support all languages, simulators, methodologies
 - Native SV, Verilog, OVM, UVM, VMM, C, SystemC, etc.
- Consistent user experience across multiple protocols
 - Common UVM library for all VIPs and Memory models



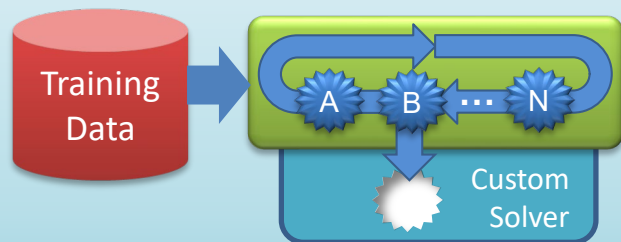
AI/ML INSIDE/OUTSIDE FOR VERIFICATION

Glossary

- Machine learning
 - The study of algorithms and statistical models that computer systems use to perform a specific task effectively without using explicit instructions
- Machine learning model
 - An approximative model for a function, automatically generated from training data, that allows inference over new data
- Reinforcement learning
 - A machine learning task that allows computer systems to automatically and dynamically determine the ideal behavior within a specific context to maximize its performance, based on observation, reward and action
- Supervised learning
 - A machine learning task in which the algorithm builds a model of an unknown function from a subset of its inputs and the desired outputs
 - In this type of learning, one needs to supply labels to the output data

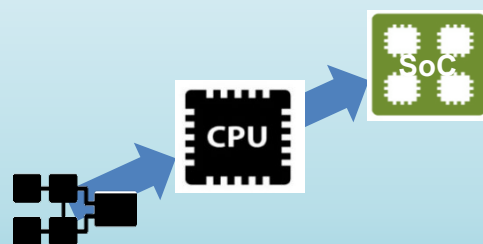
ML-enabled Formal Verification

Smart Proof Technology



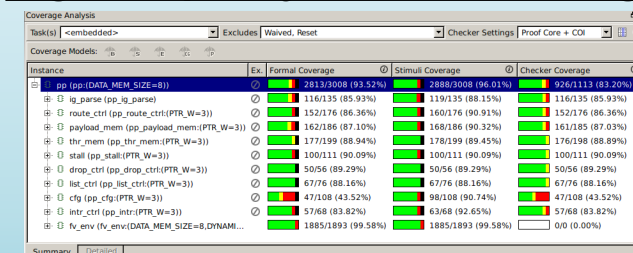
ML for solver inference and multi-advisor orchestration

Advanced Design Scalability



2x design capacity increase and 50% memory footprint reduction

Signoff-quality Formal Coverage



Signoff-accurate formal coverage with new intuitive analysis GUI



Third-Generation JasperGold® Formal Verification Platform



*“We measured **2x faster** proofs out-of-the-box, **5x faster regressions** and non-converged properties **reduced by 50%**”*

-Mirella Negro Marcigaglia, digital design verification manager, STMicroelectronics

Smart Proof automation framework



Component & Data Management

Proof Profiling Data

- Regular read/write

Proof Caching

- Cache storage in single file
- Automatic cleanup of old cache data

Multi Advisor Proof
Orchestration

- Forced on
- No overwrite on engine mode

Engine Algorithm
Selection

- Automatic training and inference

Learning

Machine Learning

Optimizes subsequent runs/regressions

Optimizes out-of-the-box proofs



Find more bugs

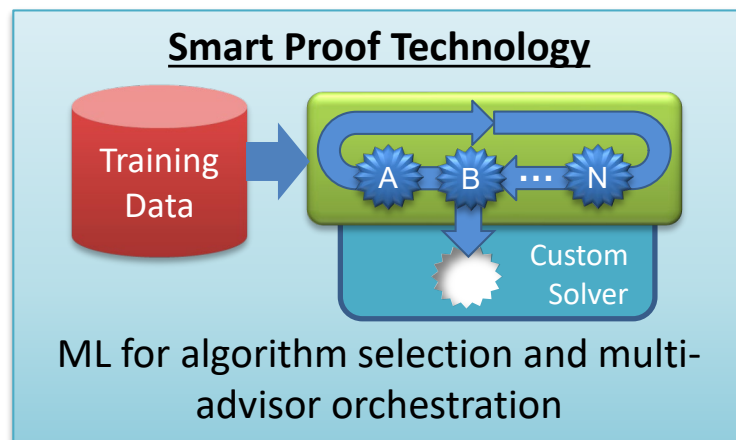


Better convergence



Faster proofs

Smart Proof technology: ML-enabled optimizations



Algorithm Selection

Supervised learning: ML uses training data from 500+ customer designs

Supervised inference: Best-fit core engine selected and configured to create custom solver

=> **Up to 4X (2X avg.) better out-of-the-box proofs**

Multi-Advisor Orchestration

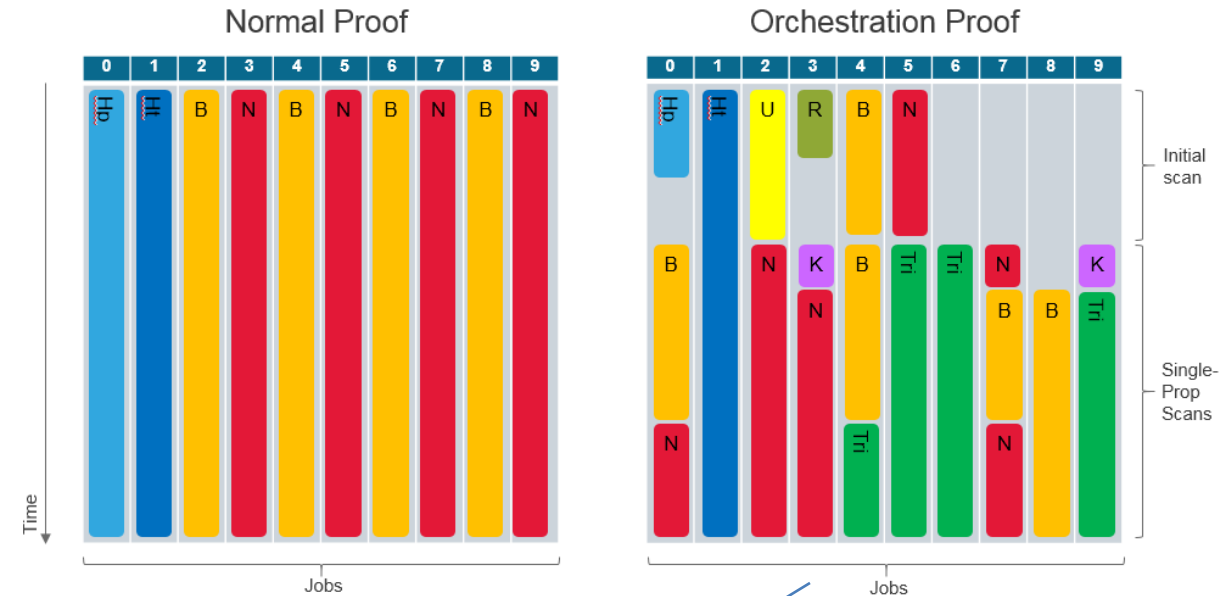
Multiple proof advisors use reinforcement learning to improve proof efficiency

Uses training data for **better out-of-the-box proofs**

Adjusts training data using proof profiling for **up to 6X (5X avg.) better subsequent proofs and regressions**

What is multi advisor proof orchestration?

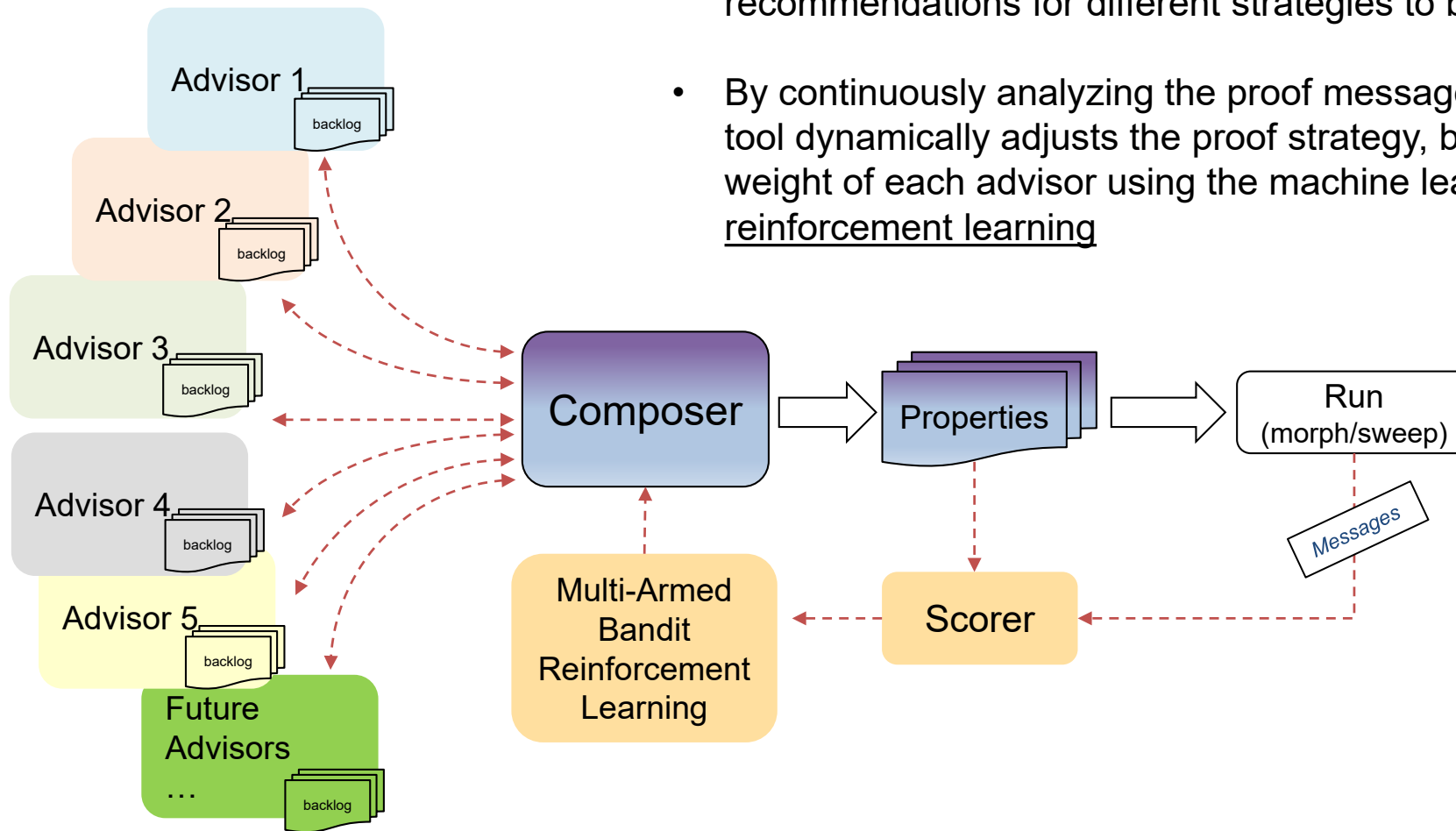
- Algorithm that dynamically adjusts engine selection, time limits, license usage, etc, during a proof, enabled by reinforcement learning
 - Hides complexity from the user: hand technology control over to the tool
 - Focuses on resources (max licenses, max jobs, global time limit), which are still respected in orchestration
 - Uses time slicing to switch engines during run, allowing engine mode diversity
 - Opens the door to future automatic optimizations (partitioning, bug hunting, etc)
 - No need to deploy every new strategy
- Continuously expanded to include new technology



+ Custom engines dynamically created during proof, to further explore non-default engine settings

Multi advisor proof orchestration overview

- A group of advisors with different weight provide recommendations for different strategies to be run during a proof
- By continuously analyzing the proof messages and results, the tool dynamically adjusts the proof strategy, by updating the weight of each advisor using the machine learning strategy of reinforcement learning



How to enable multi advisor orchestration?

- Orchestration is on by default whenever the user doesn't change the engine mode for a proof, but can be forced to be turned on with configuration commands

```
[<embedded>] % prove -task <embedded>
INFO (IPF036): Starting proof on task: "<embedded>"
INFO (IPF031): Settings used for proof thread 1:
orchestration          = on (auto)
time_limit             = 86400s
per_property_time_limit = 1s * 10 ^ scan
engine_mode            = auto
proofgrid_per_engine_max_jobs = 1
max engine jobs        = auto
proofgrid_mode         = local
proofgrid_restarts     = 10
```

```
[<embedded>] % prove -task <embedded> -engine mode {B R D}
INFO (IPF036): Starting proof on task: "<embedded>", 53 properties
INFO (IPF031): Settings used for proof thread 3:
orchestration          = off (auto)
time_limit             = 86400s
per_property_time_limit = 1s * 10 ^ scan
engine_mode            = B R D
proofgrid_per_engine_max_jobs = 1
max engine jobs        = B R D, total 3
proofgrid_mode         = local
proofgrid_restarts     = 10
```

```
1558735989: INFO (IPF152): Orchestration is evaluating 3 properties using the following settings:
per_property_time_limit = 2s
engine_mode             = M AD AMcustom2 *
```

...

```
1558736021: INFO (IPF152): Orchestration is evaluating 3 properties using the following settings:
per_property_time_limit = 8s
engine_mode             = D I Ncustom3 *
```

...

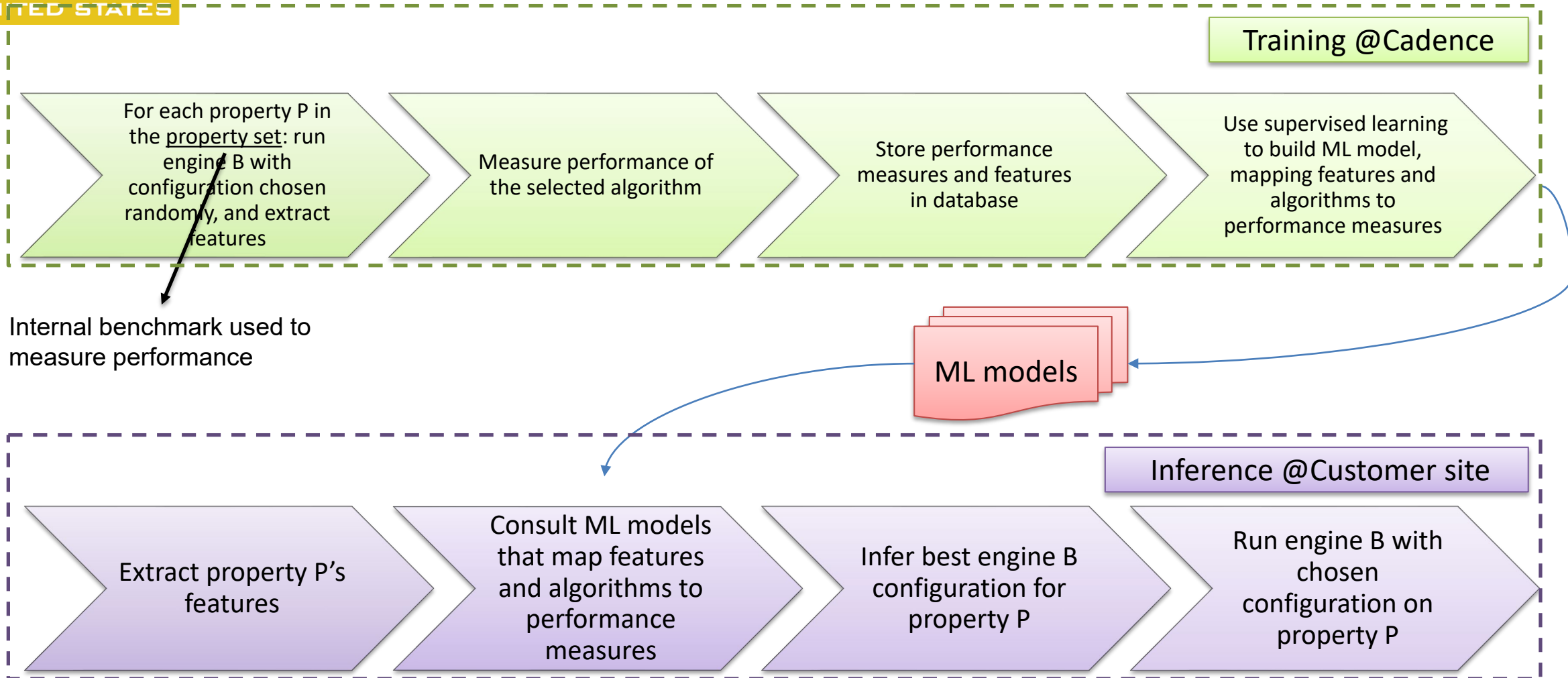
```
1558739909: INFO (IPF152): Orchestration is evaluating 3 properties using the following settings:
per_property_time_limit = 512s
engine_mode             = Tri AMcustom2 Ncustom4 *
```

*Custom engines created dynamically in the proof to further explore non default settings

Algorithm selection concept: example engine B

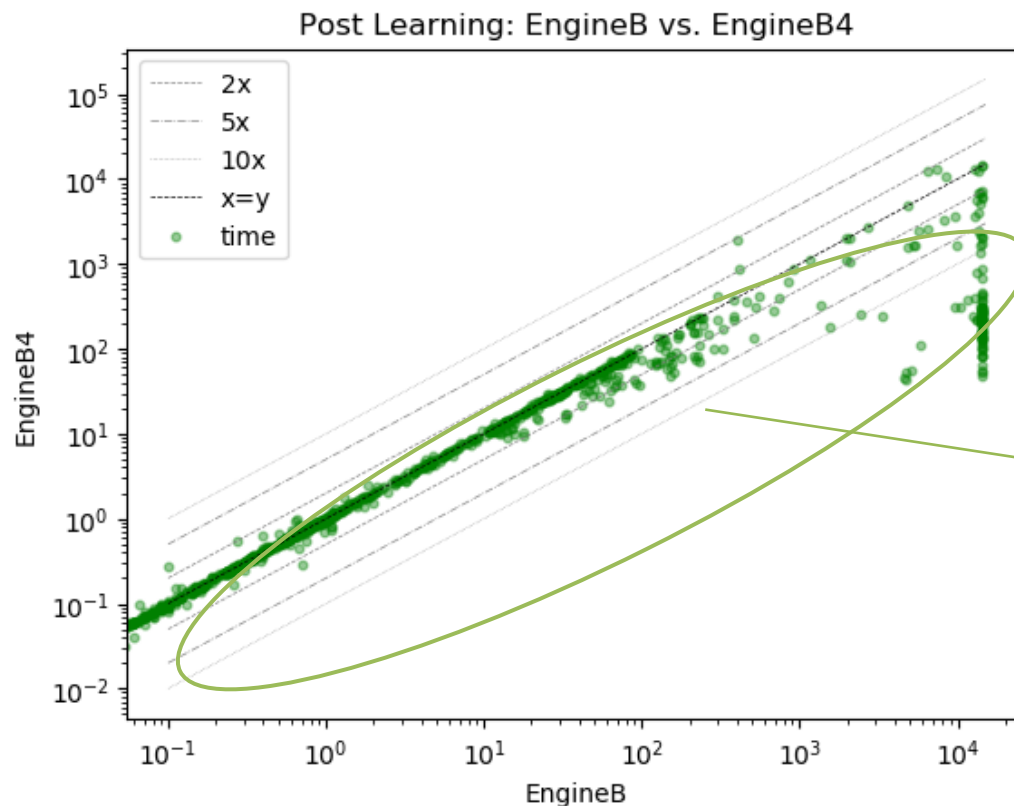
- Challenge
 - JasperGold has different technologies which are used internally by engine B
 - The challenge is to choose the best engine B parameter configuration to run on each given property
- Approach
 - Use supervised learning to create a machine learning model trained on our internal benchmarks across multiple customer designs, to try to infer on the fly better-than-default engine B configuration on each given property
 - Wrap solution into a new engine called B4, which can be added to the engine mode using `set_engine_mode` command like any other engine

Algorithm selection for engine B overview



Algorithm selection for engine B results

- Inference evaluation against engine B showed overall performance boost when using engine B4



Dots below the diagonal line correspond to properties which were proven faster with engine B4, when compared to engine B

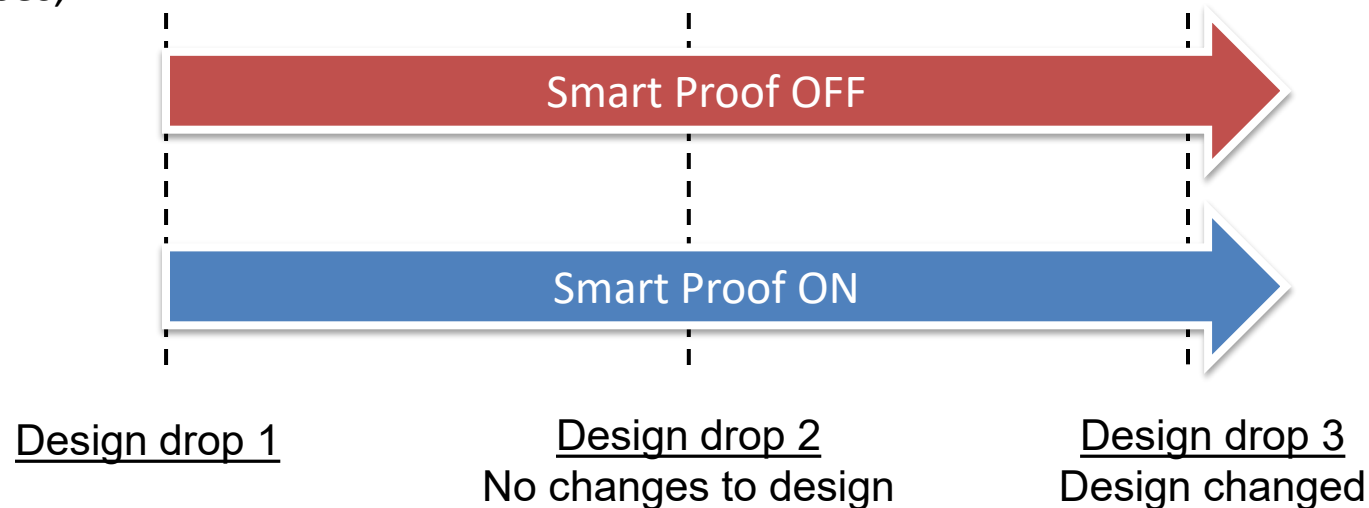
Smart Proof case study

A customer scenario

- Three versions of a customer design (processor core)
- Goal: compare performance of proof using Smart Proof vs. proof without it over time
- Resources
 - Dedicated machine
 - 20 jobs, 20 licenses, 24h

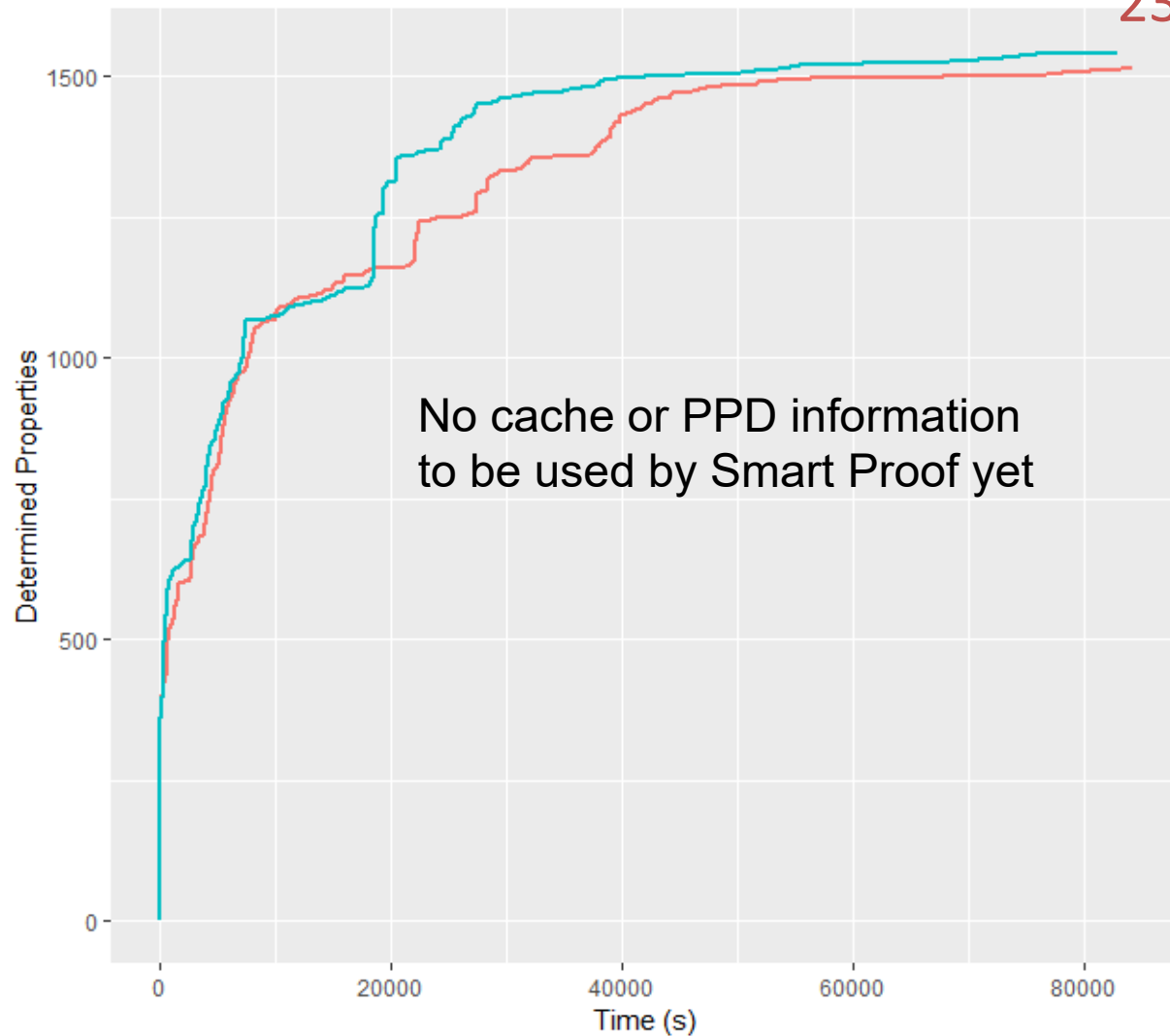
Characteristics for design drop 1

# Properties	2340
# Gates	290k
Convergence after 24h with 20 jobs (Smart Proof off)	65%



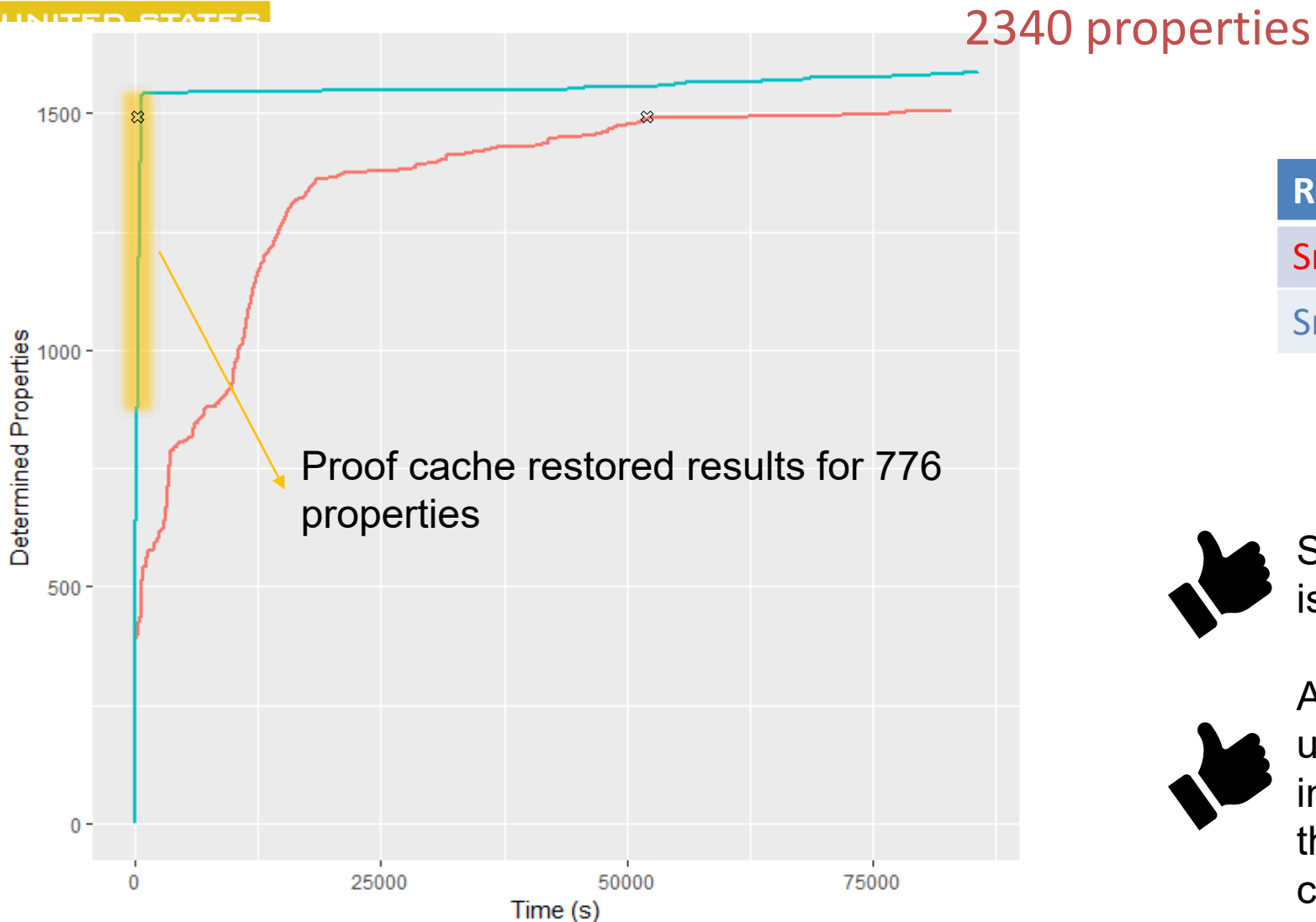
Design drop 1

2340 properties



Run	# Determined
Smart Proof OFF	1516
Smart Proof ON	1543

Design drop 2: no changes to design



Run	# Determined
Smart Proof OFF	1506
Smart Proof ON	1590

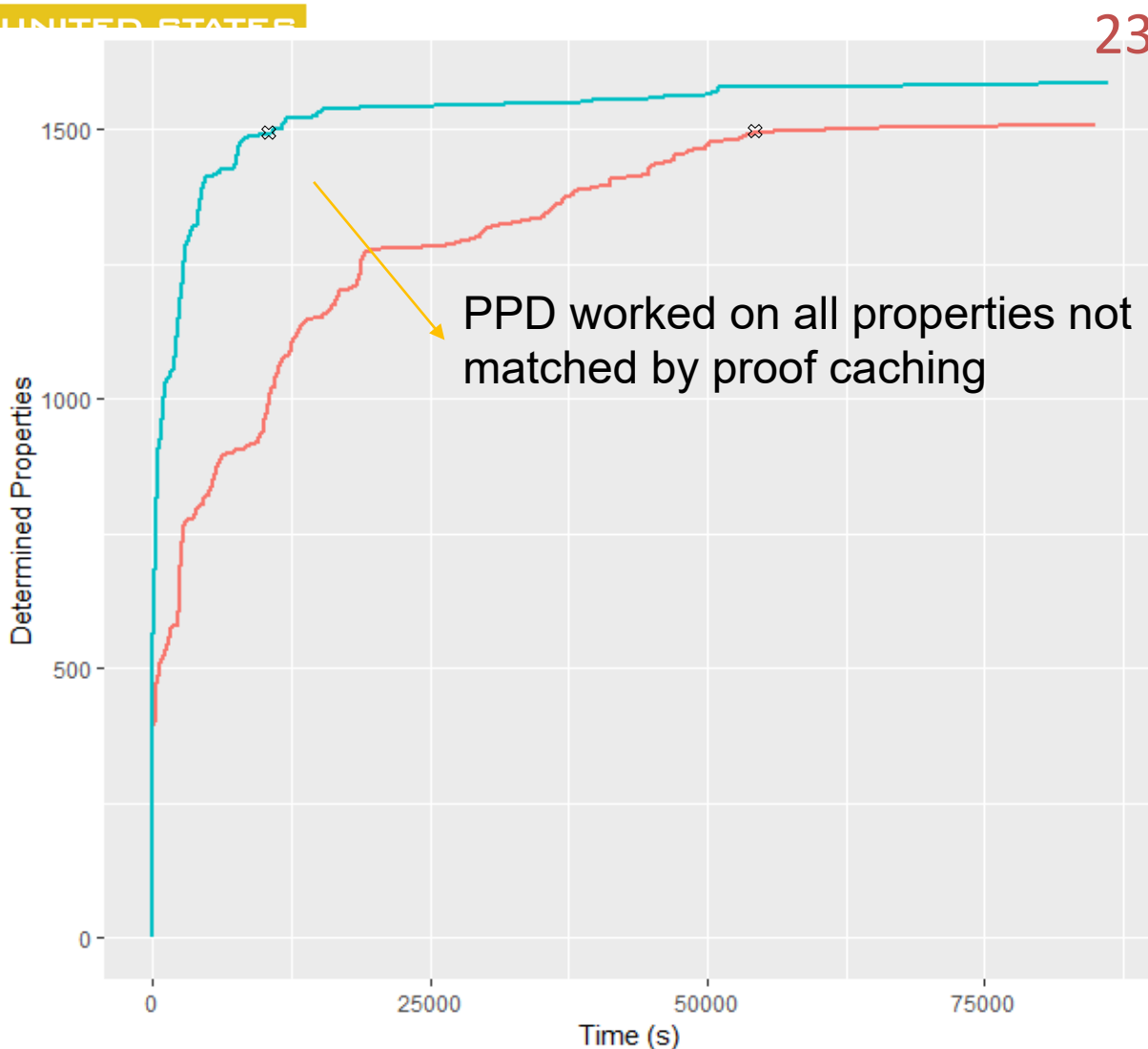


Speedup: run with Smart Proof ON is **32x faster** than proof without it



All properties determined in drop 1 using Smart Proof were reproduced in drop 2 (the same is not true for the runs with Smart Proof off in this case study)

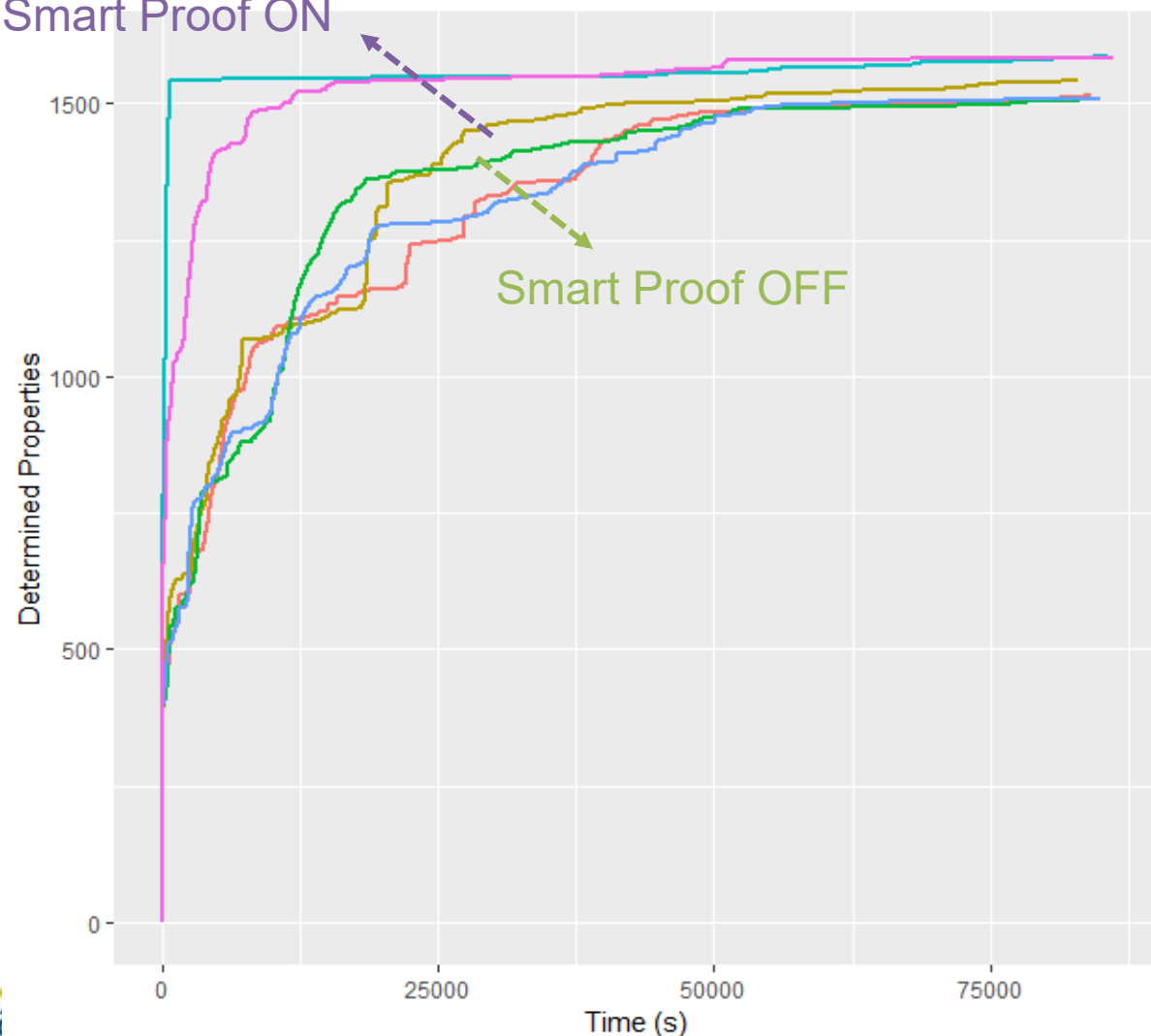
Design drop 3: design changed



Speedup: run with Smart Proof ON is **5x faster** than proof without it

Case study take away

Smart Proof ON



- By learning from previous runs, Smart Proof frees up resources, which are better utilized to solve hard properties
- When the design is stable, Smart Proof is a great tool to help reproduce previous proof results

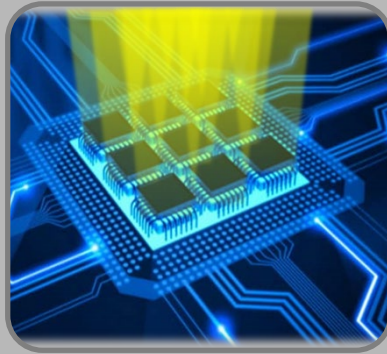
SUMMARY - OUTLOOK

Enabling Verification of AI/ML Designs



Formal

- High **Capacity**
- Regression improvements
- SAT Solver Inferencing



Simulation

- **Fast simulation** for high-activity designs
- UVM **randomization**
- Fast elaboration for **replicated structures**
- Coverage **metrics**



Emulation

- **Billion-gate designs**
- Parallel Partition Compiler
- INT8 to INT64
- **Power / Performance**
- Memory **models**: HBMx
- **Senor model**: MIPI CSI

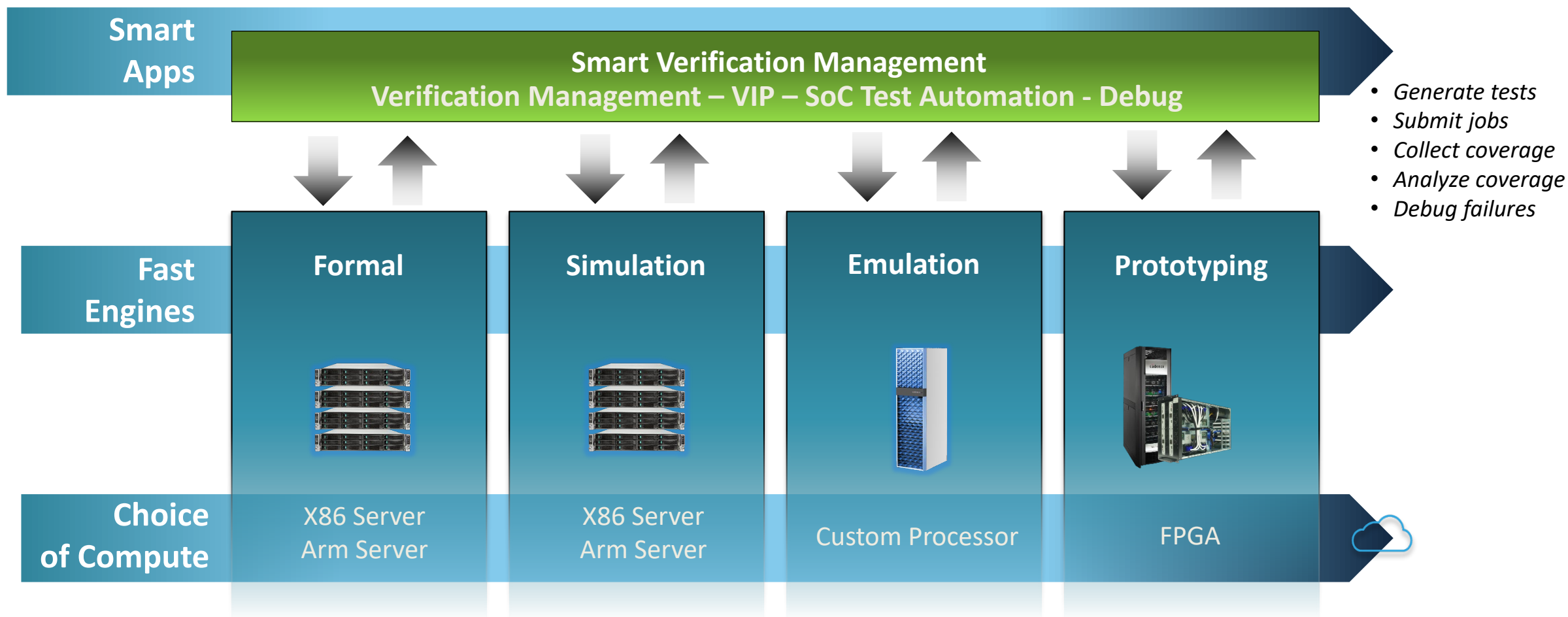


Prototyping

- **Scaling** to large designs
- Unified frontend
- ICE test suite: **Faster, data-extended (DL) regression**
- **SW driven validation** – deep learning data training refinement

Verification Throughput!

Find and fix the most bugs per \$ invested in bare metal compute



Some Examples

- Datacenter / Edge limited by compute power
- Accelerators
 - Custom to the ML application. No one size fits all.
 - HW/SW Co-validation at architectural level - performance using up to date user models
- Training/Inference in datacenter:
 - Massive # of systems - every milliwatt counts
 - HW/SW Co-validation with power analysis, user models
- Inference at Edge:
 - Very limited power budget
 - Validating power spec required before fabrication



Palladium Z1 “instrumental for Gaudi”
Source: Cadence Earnings Call Q1’19

Palladium Z1
Source: Cadence Earnings Call Q1’19

HW Portfolio
Source: Cadence Earnings Call Q2’19

Full Verification Suite
Source: Press Release

HW Portfolio – Capacity, Debug
Source: Video @ DVCON Keynote