



TestConX

中国

China

November 13, 2025
Shanghai

www.testconx.org

TestConX China 2025

MEDIATEK



Test Challenges and Direction in the Age of AI Everywhere



Harry H. Chen, IC Testing Scientist
Computing & AI Technology Group



— MEDIATEK OVERVIEW

MediaTek at a Glance

- 5th Largest fabless IC design company
- \$16.5B Revenue in 2024 (USD)
- \$4.1B R&D Investment in 2024 (USD)
- 2 Billion + Connected devices powered annually
- >21,000 MediaTek Group Employees

Growing



Automotive



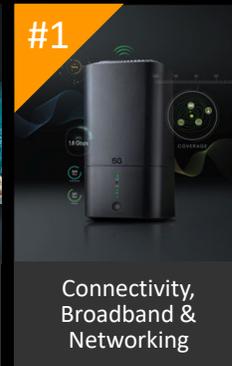
ASIC Solutions



Smartphone



Smart TV



Connectivity,
Broadband &
Networking



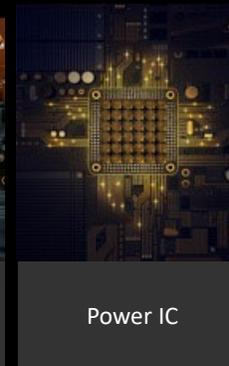
Arm Chromebook
& Android Tablets



Smart Home



IOT



Power IC

Source: IDC, Gartner, TechInsights, IHS and MediaTek company data (based on 2024 market share)

Penetration of AI into our daily lives

10 EVERYDAY EXAMPLES OF ARTIFICIAL INTELLIGENCE AI IN DAILY LIFE



Virtual Assistants



Smart Phone
Features



Personalized
Shopping



Email Spam
Filters



Music and Video
Streaming



Navigation and
Ride-Sharing Apps



Social Media



Banking



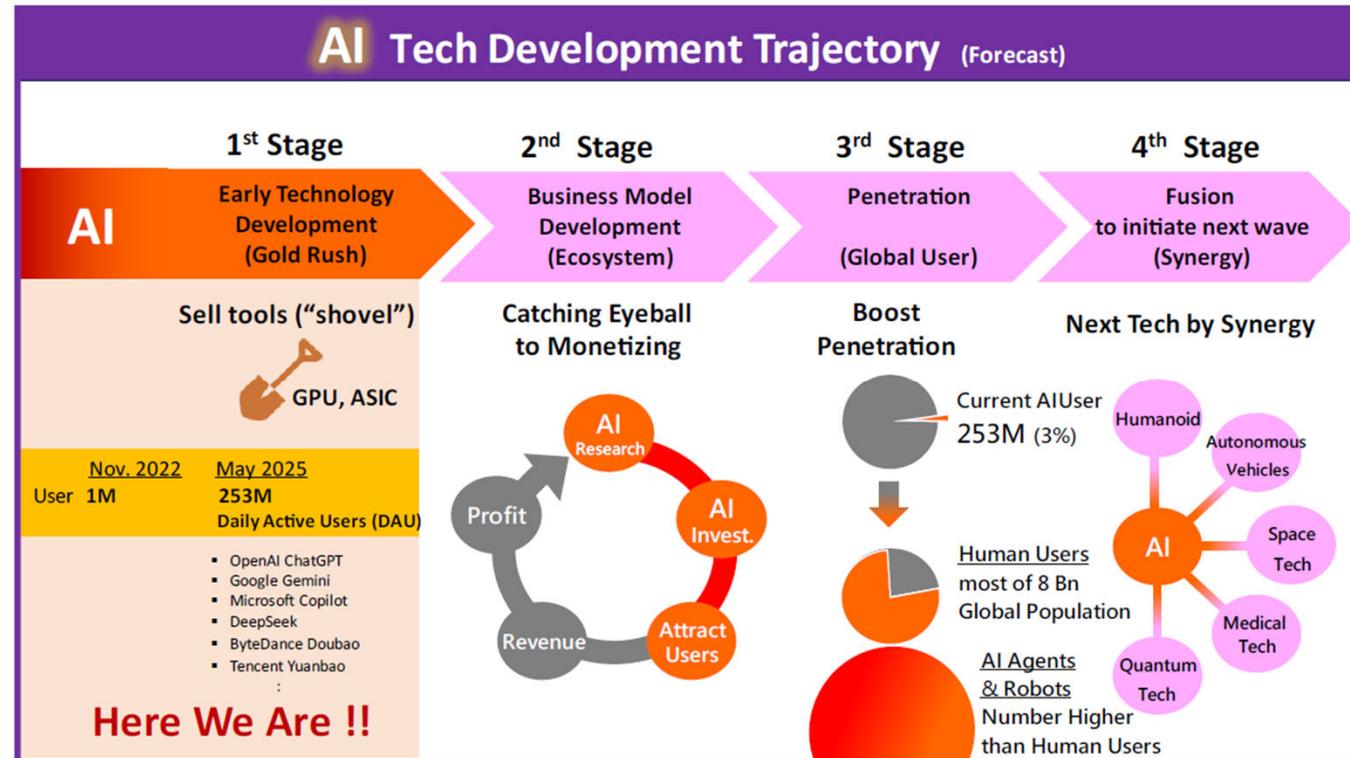
Smart Home
Devices



Chatbots and
Customer Support

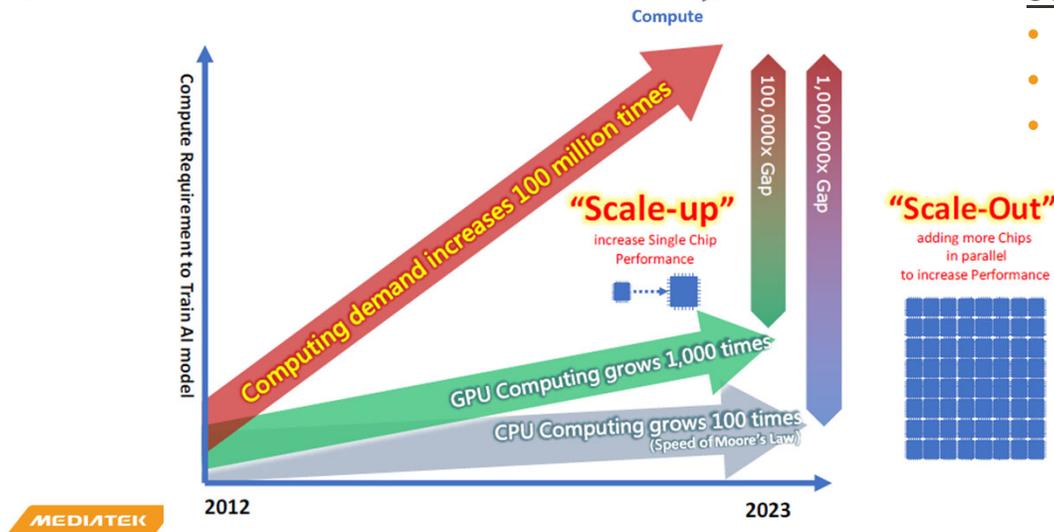
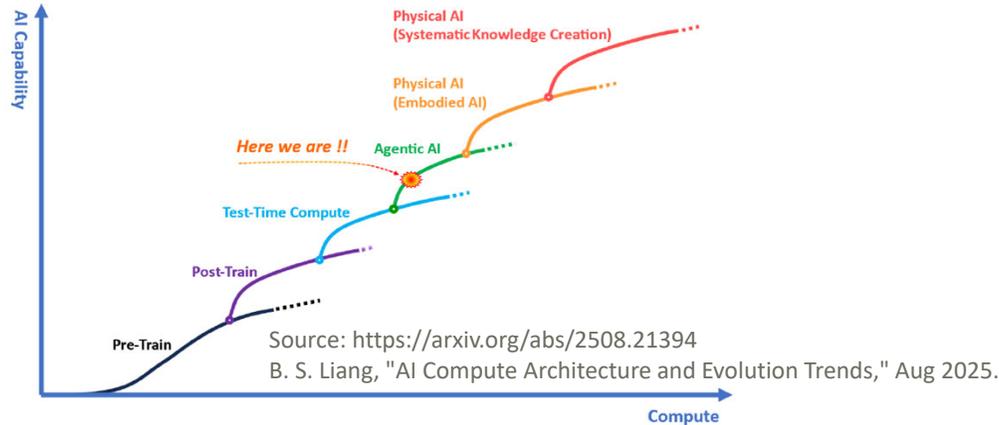
Source: <https://aiproinsight.com/10-everyday-examples-of-ai-ai-in-daily-life/>

Adoption of AI is still at early stage



Source: B. S. Liang, "AI Compute Architecture and Evolution Trends," 2025.
<https://arxiv.org/abs/2508.21394>

Driver of semiconductor growth for years to come



Scale-up

- Semiconductor process shrink to angstrom level
 - Moore's Law slowing down, reticle limit
- Multi-die integration via advanced packaging
 - HBM, chiplets, CoWoS, 3DIC, hybrid bonding
- Domain-specific architectures (DSA)
 - TPU, FP{32,16,8,4} numbering, HW/SW co-optimize
- Performance versus energy efficiency

Scale-Out

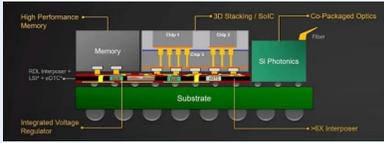
- Millions of cores deployed across data centers
- High-speed (optical) interconnect ↔ core clusters
- Multi-datacenter training & inference



Source: <https://cloud.google.com/blog/products/compute/how-cloud-tpu-v5e-accelerates-large-scale-ai-inference>

What about reliability, safety, and security?

DPPM = Defective Parts Per Million
SLT = System-Level Test

Integration Level	Trends	Challenges
<p style="text-align: center;">Die</p> 	<p>Smaller geometry ⇒ higher variation Higher density ⇒ thermal hotspots Complex fabrication steps ⇒ more defects Lower operating voltage ⇒ less margin Multi-core/domain ⇒ design complexity</p>	<p>Fault modeling gaps Marginal defect escapes Higher DPPM Structural test gaps More functional test & SLT</p>
<p style="text-align: center;">Package</p> 	<p>Denser & smaller connections Thermal, mechanical & material considerations High-speed SerDes, co-packaged optics (CPO) Backside power delivery Multi-die matching & interaction</p>	<p>New interconnect defect types Heat dissipation & thermal warpage Interconnect performance monitoring Test access, redundancy & repair SLT for assembled die stack</p>
<p style="text-align: center;">System</p> 	<p>Rising HW/SW system complexity Rising energy consumption Higher current flows Liquid cooling High field utilization stress-induced aging</p>	<p>High system DPPM Silent data corruption (SDC) Costly failure diagnosis Implementing shift-left SLT Cost-effective SLT methodology</p>

Limits of DPPM prediction and fault models

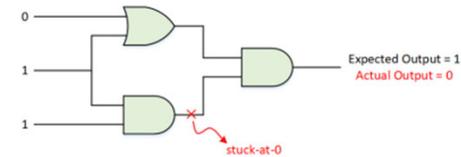
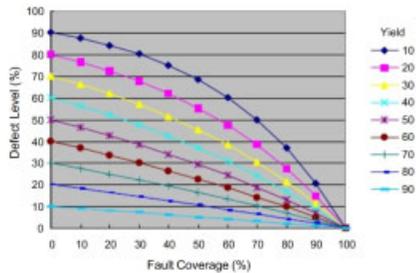
Equation used to predict DPPM

$$\frac{DPPM}{10^6} = 1 - Y^{(1-T)}$$

Y and T are fractional numbers in $[0, 1]$

Y = manufacturing yield (estimated by actual wafer test yield)

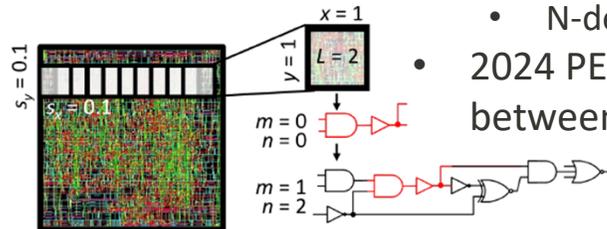
T = fault coverage (stuck-at fault model)



Problems

- Simplified statistics: uniform distribution of independent faults
- Non-digital logic areas? memory, analog, power network
- Stuck-at faults (60+ years) no longer reflect today's fab defects
- Defect-oriented models: delay, cell-aware, opens, bridges
- Generate tests to increase fortuitous defect detection
 - N-detect, pseudo-exhaustive
- 2024 PEPR testing study showed >95% discrepancy between fault model and extracted faulty function*

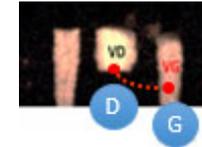
PEPR =
Pseudo-Exhaustive
Physically-Aware
Region



* C. Nigh et al., "Faulty Function Extraction for Defective Circuits," ETS 2024.

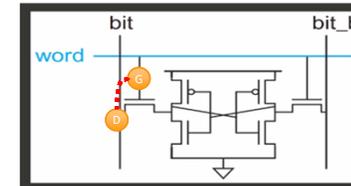
Example of 3nm defect escaping structural test

CMP scratch during MEOL step \Rightarrow
Soft resistive bridge connecting transistor Drain and Gate terminals



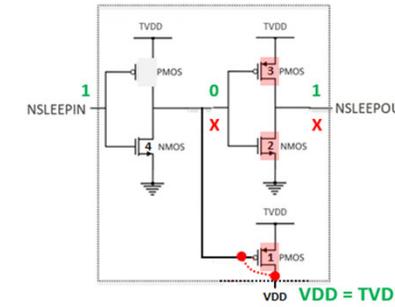
Impact depends on where the defect lands

At bit-cell in memory with MBIST \Rightarrow
Short between word and bit lines escape MBIST algorithms



Customer RMA

Power distribution: switch chain to control in-rush current \Rightarrow
Local voltage drop causes higher V_{min} running specific SLT workload

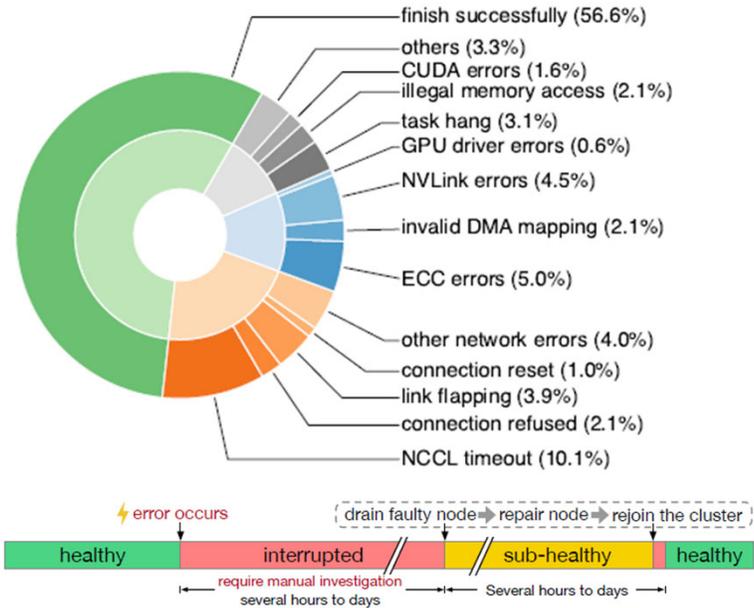


Customer RMA

$$VDD = TVDD \times \text{Resistive Divider}$$

Digital logic with scan-DFT \Rightarrow
Detectable by cell-aware structural ATPG test pattern?

System reliability is a serious concern at AI CSPs



Component	Category	Interruption Count	% of Interruptions
Faulty GPU	GPU	148	30.1%
GPU HBM3 Memory	GPU	72	17.2%
Software Bug	Dependency	54	12.9%
Network Switch/Cable	Network	35	8.4%
Host Maintenance	Unplanned Maintenance	32	7.6%
GPU SRAM Memory	GPU	19	4.5%
GPU System Processor	GPU	17	4.1%
NIC	Host	7	1.7%
NCCL Watchdog Timeouts	Unknown	7	1.7%
Silent Data Corruption	GPU	6	1.4%
GPU Thermal Interface + Sensor	GPU	6	1.4%
SSD	Host	3	0.7%
Power Supply	Host	3	0.7%
Server Chassis	Host	2	0.5%
IO Expansion Board	Host	2	0.5%
Dependency	Dependency	2	0.5%
CPU	Host	2	0.5%
System Memory	Host	2	0.5%

(148+17+6)/16K ≈ 10.7K GPU logic DPPM

Alibaba: 256 H800 GPUs

- GPT-3 training failure rate = 43.4% for top 5% of resource intensive tasks
- Recovery can take up to days

Source: "Unicon: Economizing Self-Healing LLM Training at Scale," 2023.
<https://arxiv.org/abs/2401.00134>

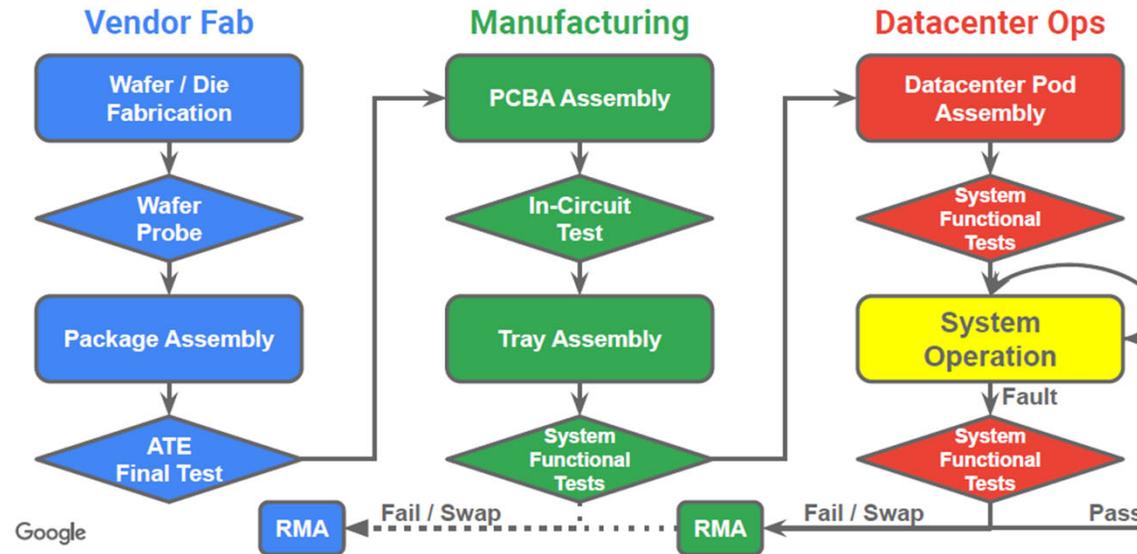
Meta: 16K H100 GPUs, 80GB HBM3

- Failures during 54-day period of Llama3 450B pre-training
- 58.7% failures GPU-related

Source: "The Llama 3 Herd of Models," 2024.
<https://arxiv.org/abs/2407.21783>



Google reports raise alarm & call for new test approach

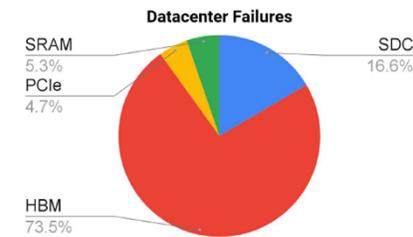
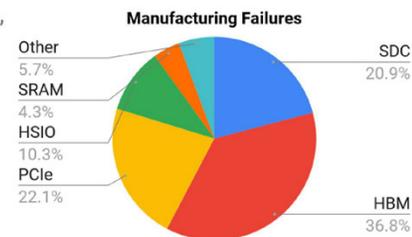


Our vendor-driven Fab testing follows high norms, built on high coverage by DFT structures.

However, NO mission-mode; Minimal functional testing @ ATE

- ATE Tests
 - >99% Stuck At Coverage
 - >94% Transition Delay Fault Coverage
 - Cell-Aware testing added to address SDC

With billions of transistors, high coverage can still leave millions of untested elements and paths!



Source: "Training in Turmoil: Silence Data Corruption in Systems at Scale," ITC 2021 Silicon Lifecycle Management Workshop.

10³ – 10⁴ DPPM far exceeds industry target of 100 DPPM

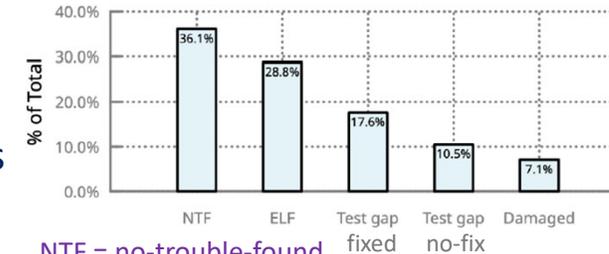
Google observations since 2016/2017 across multiple platform generations

1. Test escapes ~5000 DPPM over lifetime
2. Test escapes resulting in SDC ~1000 DPPM
3. SDC-causing chips produce errors frequently (~1M in 1B hours)
4. Limited diagnose and root-cause analysis of test escapes
5. Most test escapes are detected after deployment in data centers

Detection approach	% defective machines detected	Defective machines detected (per million)
Pre-deployment testing	12%	479
Post-deployment		
Online/Offline testing	29%	1099
System health and forensics	49%	1886
User-level	10%	393

Source:
“Silent Data Corruption by 10x Test Escapes Threatens Reliable Computing,” 2025.
<https://arxiv.org/abs/2508.01786v4>

4. Summary of vendor root-cause analysis



NTF = no-trouble-found

ELF = early-life failures

QED = Quick Error Detection

Call to action

A. Quick diagnosis of SLT fails

- * insights into causes of test escapes
- * need techniques such as QED & fault injection studies

B. In-field bad chip detection

- * functional ATPG & in-field scan test
- * vary test conditions {voltage, frequency, temperature}

C. New test experiments

- * fill current gaps and validate new test approaches
- * break silos, increase trust and collaboration

New focus on marginal defects and system health condition

- Traditional testing focused on “binary” pass/fail or good/bad
 - Mostly assumes gross defects that cause failure under any condition
- Marginal defects are more common in advanced nanometer nodes
 - Only cause failure under specific conditions and workloads
 - Usually escapes structural testing due to difference from system mode
 - May be detected by SLT but with high effort, but no guarantees
- System-level profiling and health condition check holds promise
 - CSPs use periodic health checks and checkpointing to enhance utilization and reliability
 - Meta CP-Bench*, a test suite to detect HW crash, performance degradation, and SDC

CP-Bench (C=configure, P=parameterize)

- ◆ 30+ PyTorch AI workloads
- ◆ In-house HW health checks
- ◆ Distributed and concurrent modes
- ◆ Chip vendor RMA criteria

Offline profiling of healthy HW

Perf: time to finish a single training step
SDC: model parameters checksum at every #N training steps

Online anomaly detection

Detect degradation if >3% performance drop
Detect SDC if different checksum values

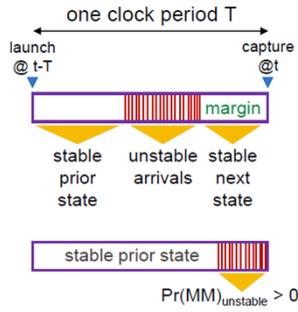
Run modes

Distributed: one model across host node GPUs (covers interconnect)
Concurrent: one model on each GPU (for GPU-level diagnosis)

* X. Jiao et al., “CP-Bench: A PyTorch Test Suite to Detect AI Hardware Failure, Performance Degradation, and Silent Data Corruption,” ITC 2025.

Shift left to vendor ⇒ chip internal health profiling

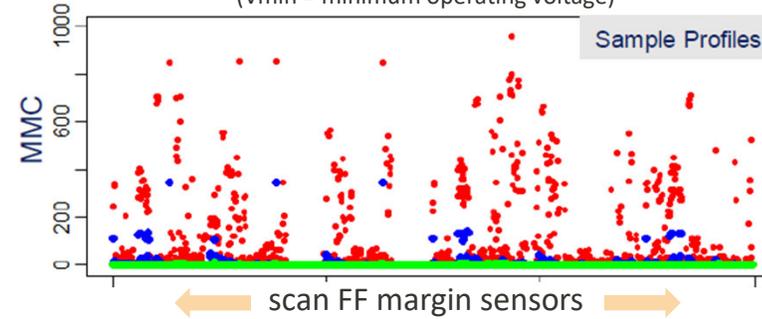
Fine-grain MMC profile of margin sensitivity at FF cones*



Scan FFs ⇒ margin sensors
 Apply tests under “stress”
 ↓ voltage, ↑ frequency
 Squeeze margin
 Force capture errors
 ⇒ mismatch count (MMC)
 US Patent 9465071 B2 Oct. 11, 2016.

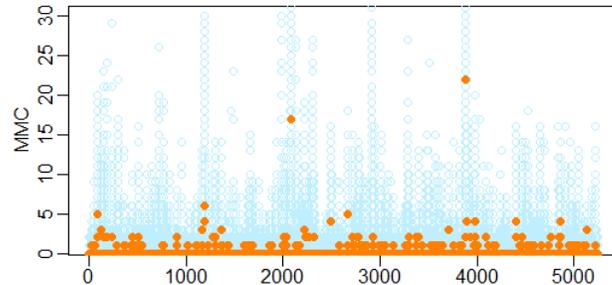
MMC profiles correlated to Vmin

3 samples: { low, middle, high } Vmin
 (Vmin = minimum operating voltage)

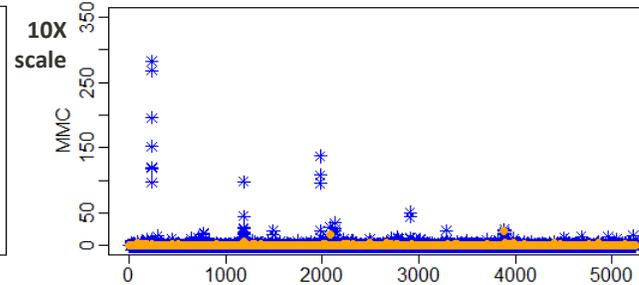


Q2s of samples MMC distribution

⇒ systematic design margin sensitivity



Outliers ⇒ marginal defectivity



Chip ↔ System Link

- ◆ Zero HW overhead for scan-DFT logic
- ◆ Also applies to in-field scan
- ◆ Not pass/fail ⇒ health binning
- ◆ Refine by correlation to system health
- ◆ System-level adaptive SLT
- ◆ System-level adaptive workload

* H. H. Chen, “Analysis of Vmin Variability in Complex Digital Logic via Post-Silicon Profiling,” VLSI-DAT 2023.

Implications for future development

Heterogeneous Integration Roadmap (HIR)

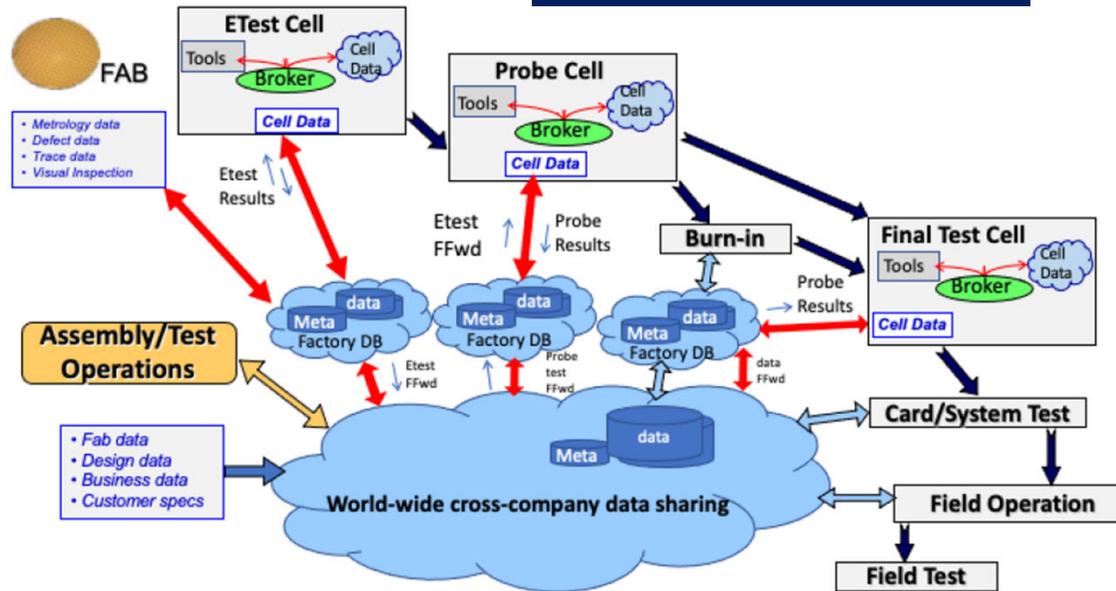
Test Chapter (2024 edition)

Section 8: System Level Test

Section 9: Data Analytics

Adaptive Test Architecture

- *dynamic data-driven*
- *cross-insertion data sharing*
- *feed forward & backward*



- ◆ Achieve full vision of adaptive test
- ◆ Massive test data volume to analyze
⇒ insights ⇒ optimized decisions
- ◆ System context at all test steps
- ◆ System-level fault modeling
→ failure analysis
→ test generation
→ fault tolerant design
- ◆ AI techniques broadly used to conquer complexity & achieve practicality

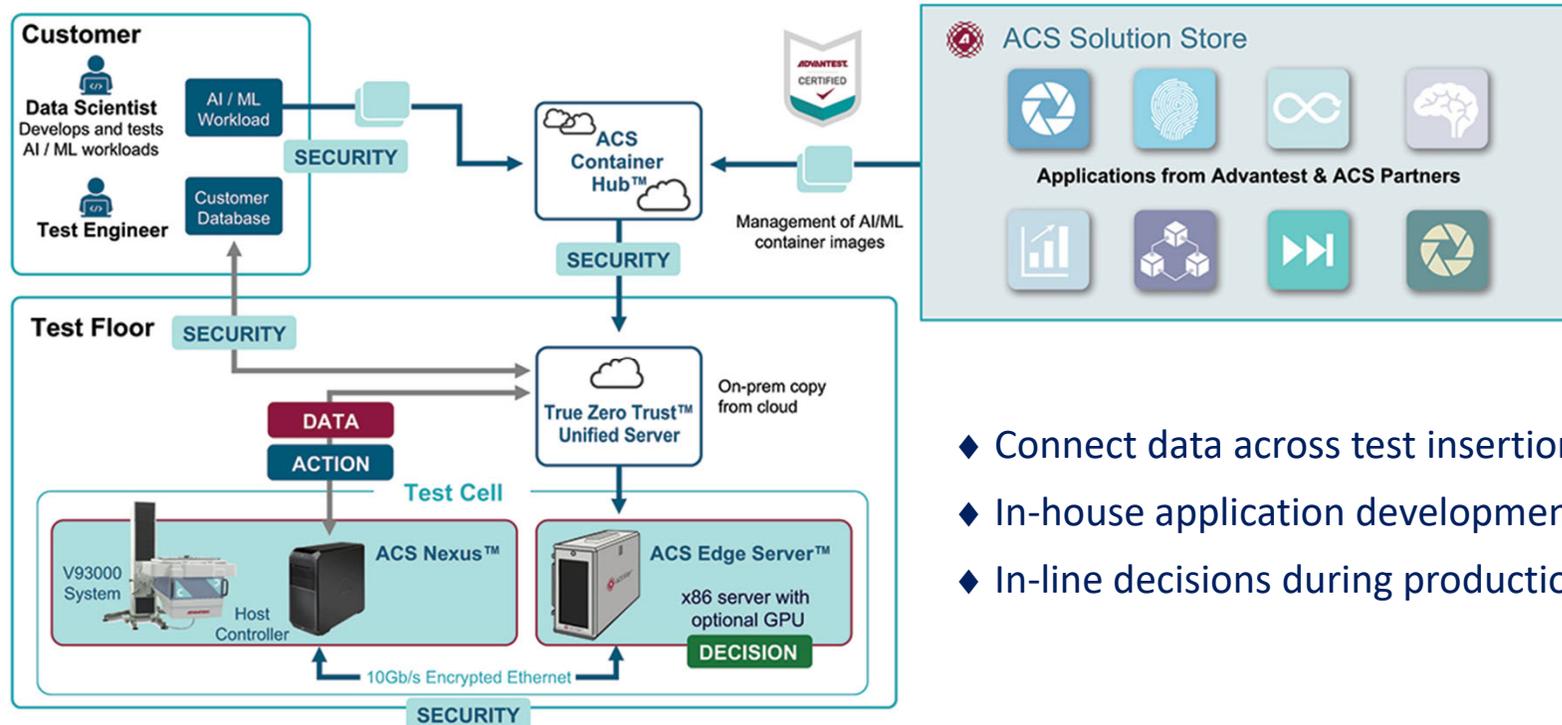


Source: https://eps.ieee.org/images/files/HIR_2024/HIR_2024_ch17_Test_Technology.pdf

Advantest ACS RTDI

ACS = Advantest Cloud Solutions
RTDI = Real-Time Data Infrastructure

Enables shift from traditional test workflows to adaptive AI-driven systems



- ◆ Connect data across test insertions
- ◆ In-house application development
- ◆ In-line decisions during production

Source: <https://www.advantest.com/en/products/acs/rtdi/>

Teradyne Titan HP for AI device SLT

AI device has larger physical size, draws high current, and runs hot
Reduces “form factor” mismatch to run more system functional content



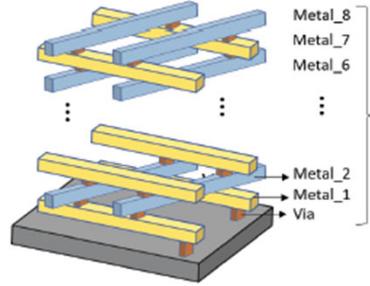
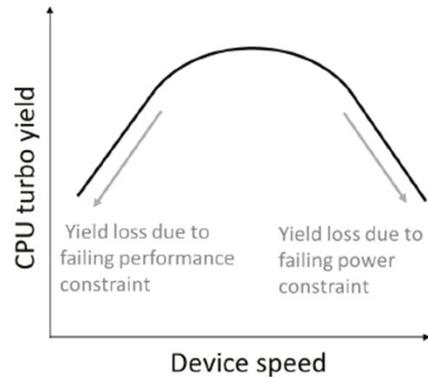
- ◆ High per device kW power delivery
- ◆ Active thermal control
- ◆ Shift tests between ATE & SLT

Source: <https://www.teradyne.com/products/titan-hp/>

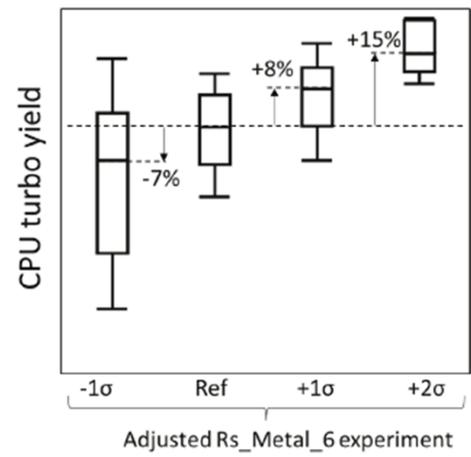
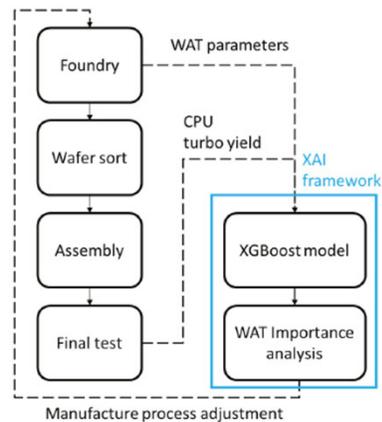
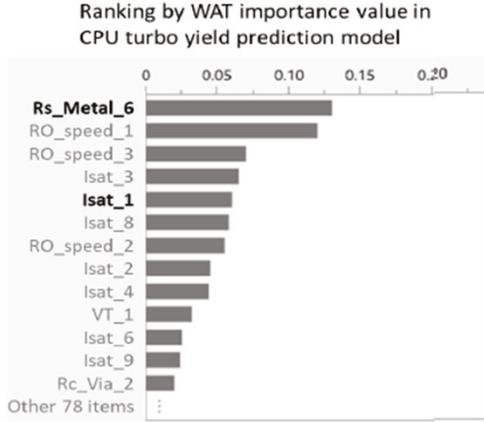
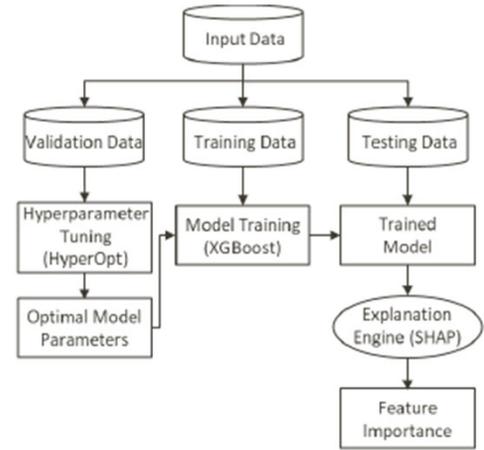
Use of AI to boost CPU turbo yield*

Turbo Yield

Higher operating frequency beyond sign-off without power increase



Related WAT: Rs_Metal and Rc_Via



* C. W. Lin et al., "Boost CPU Turbo Yield Utilizing Explainable Artificial Intelligence," ITC 2024.

Research in system-level fault modeling & analysis

Need to understand how marginal defects affect system failing behavior

Fault injection experiments at RTL or higher level of abstraction [1]

Use to validate fault tolerance features, e.g., protect more vulnerable design elements

Fault model: single-cycle register bit-flip (SC-RBF), transient fault not like stuck-at

Traditional fault simulation is too slow due to tremendous size of SC-RBF fault space

Developed ACE-Pro [2,3] to drastically reduce SC-RBF fault space

Obtain register R/W cycles from [single RTL good-machine simulation](#)

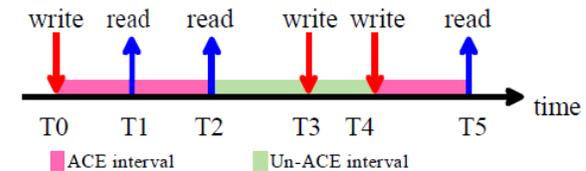
Use architecturally-correct execution (ACE) analysis to find un-R/over-W cycles ([Un-Ace interval](#))

For every register bit's def/use interval, [simulate SC-RBF at last-R cycle for one cycle](#)

Derive multi-cycle chains of equivalent faults for all time (E-chains)

Only need to fault simulate one SC-RBF from [non-masked E-chains](#)

Can achieve [more than 99% reduction in fault space size](#)



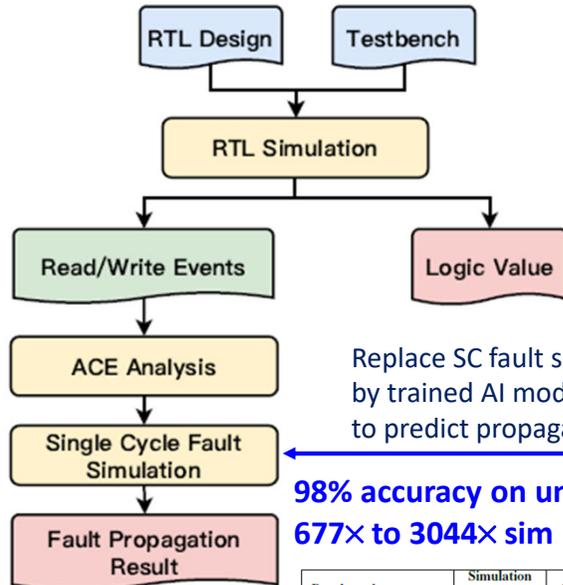
But simulating subset of ACE faults for even one cycle still takes too long!

[1] Y. He et al., "Understanding and Mitigating Hardware Failures in Deep Learning Training Accelerator Systems," ISCA 2023.

[2] D. A. Yang et al., "ACE-Pro: Reduction of Functional Errors with ACE Propagation Graph," ITC 2021.

[3] D. A. Yang et al., "Transient Fault Pruning for Effective Candidate Reduction in Functional Debugging," ITC 2022.

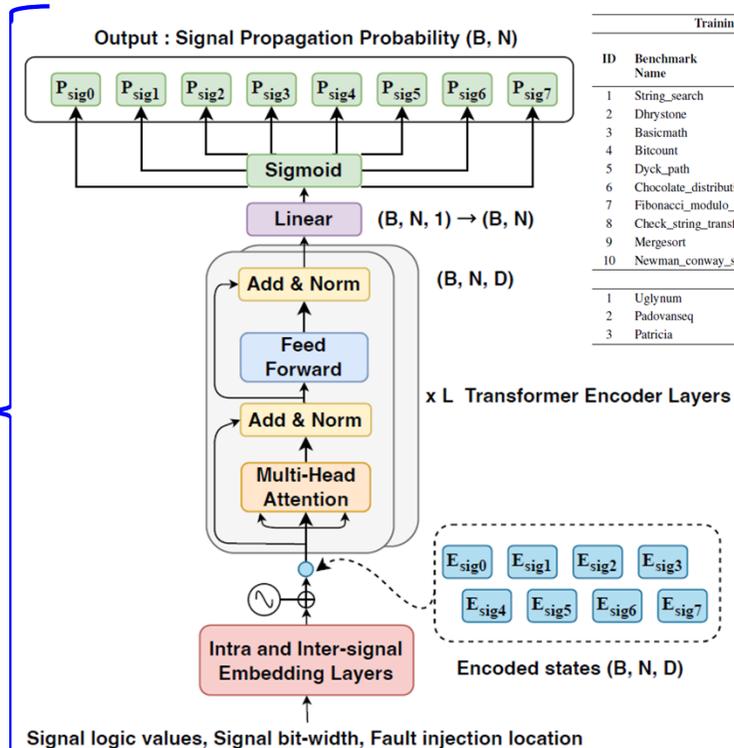
Predict fault propagation behavior with AI transformer*



Replace SC fault sim by trained AI model to predict propagation

98% accuracy on unseen benchmarks
677x to 3044x sim : inference speed-up

Benchmark	Simulation Sample Count	Fault Simulation Time (hrs)	Model Inference Time (sec)	Speedup (X)
Uglynum	16324	26.8	32.4	2974
Padovansq	17180	27.2	33.0	2970
Patricia	15909	26.0	31.7	2949
String_search	22434	34.1	43.83	2801
Dhrystone	20916	29.4	44.1	2403
Basicmath	21621	9.9	52.6	677
Bitcount	22172	19	46.0	1488
Dyck_path	15106	23.1	31.5	2641
Chocolate_distribution	16269	27	35.6	2727
Fibonacci_modulo_p	14802	23.2	29.7	2808
Check_string_transform	15771	26.2	31.3	3013
Mergesort	16671	26.5	31.3	3044
Newman_conway_seq	16420	26.8	37.1	2603



PicoRV32 RISC-V core

Training Group with Data Collection				
ID	Benchmark Name	Total Faults	Remained Fault Count After ACE (count/%)	Data Collection (Fault Sim) Time (hrs)
1	String_search	598968	121167 20.2%	27.3
2	Dhrystone	598968	112965 18.9%	24.5
3	Basicmath	598968	116637 19.5%	8.4
4	Bitcount	598968	110859 18.5%	10.9
5	Dyck_path	598968	84987 14.2%	19.3
6	Chocolate_distribution	598968	90001 15.0%	20.3
7	Fibonacci_modulo_p	598968	83643 14.0%	18.2
8	Check_string_transform	598968	87333 14.6%	20.9
9	Mergesort	598968	91692 15.3%	22.0
10	Newman_conway_seq	598968	90075 15.0%	21.1
Test Group (Unseen)				
1	Uglynum	598968	90111 15.0%	20.8
2	Padovansq	598968	94034 15.7%	20.8
3	Patricia	598968	88256 14.7%	18.2

Conclusion and take-aways

- ❑ Rise of AI driving more complexity and creating tough test challenges
- ❑ Imperative to move rapidly towards full data-driven adaptive test {methods, flows, practices}
- ❑ Deep collaboration throughout the supply chain becomes a necessity
- ❑ AI will play a key role in speeding the development of future solutions

From my TestConX China 2021 Keynote Talk



Test engineers → IC “medical” doctors

- DFX to make healthy IC & systems
- Run tests to probe health condition
- Diagnose IC illness to find correct fix

Requires multidisciplinary knowledge

- ✓ Testing and diagnosis methods
- ✓ Design flow from system to devices
- ✓ Data analytics, use of ML/AI as tools
- ✓ Reliability, Safety, Security

Thank You



谢谢

COPYRIGHT NOTICE

The presentation(s) / poster(s) in this publication comprise the Proceedings of the TestConX China 2025 workshop. The content reflects the opinion of the authors and their respective companies. They are reproduced here as they were presented at the TestConX China 2025 workshop. This version of the presentation or poster may differ from the version that was distributed at or prior to the TestConX China 2025 workshop.

The inclusion of the presentations/posters in this publication does not constitute an endorsement by TestConX or the workshop's sponsors. There is NO copyright protection claimed on the presentation/poster content by TestConX. However, each presentation / poster is the work of the authors and their respective companies: as such, it is strongly encouraged that any use reflect proper acknowledgement to the appropriate source. Any questions regarding the use of any materials presented should be directed to the author(s) or their companies.

“TestConX”, the TestConX logo, the TestConX China logo, and the TestConX Korea logo are trademarks of TestConX. All rights reserved.