

Meta Silicon Infrastructure and Evolution with AI

Rupa Raghavan Farishta Mahzoz Salina Dbritto

Agenda

Silicon Validation Infrastructure - Rupa Raghavan

Englab Operations - Farishta Mahzoz

Evolution with AI - Salina Dbritto

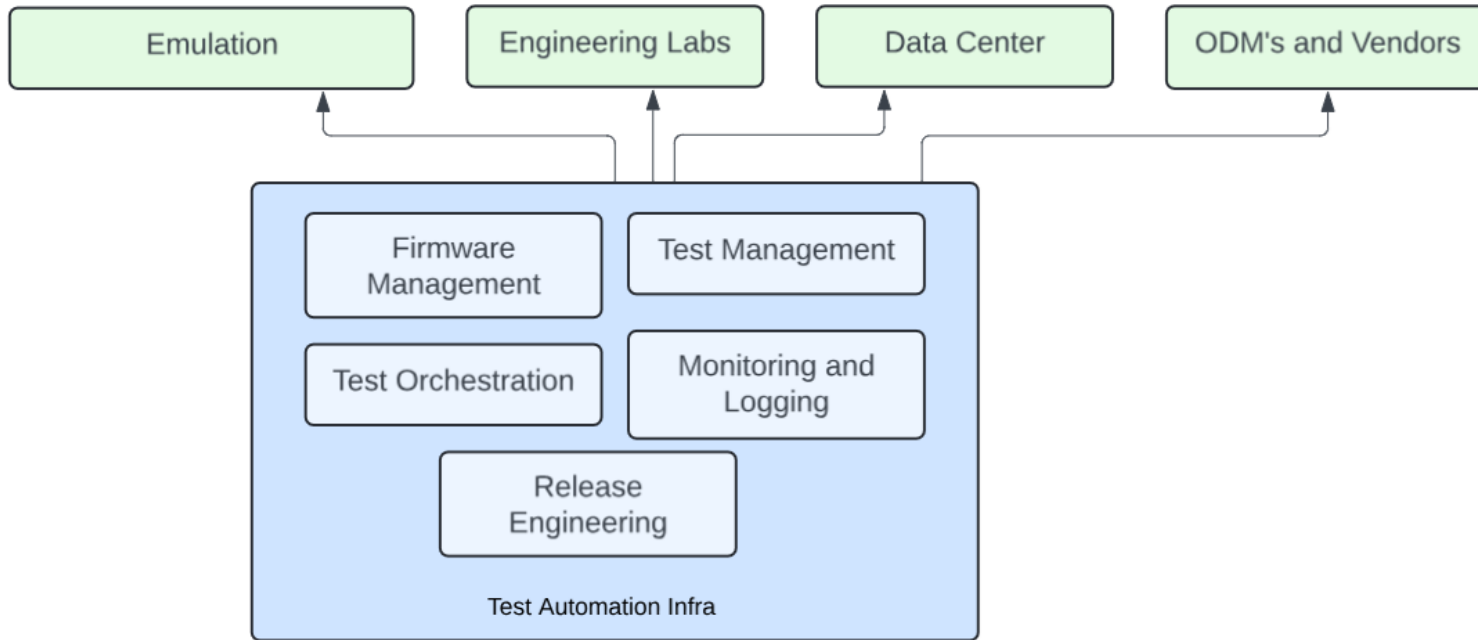
Introduction

- Meta is developing specialized silicon to support AI workloads
- Required a paradigm change in our approach to test and validate AI systems and platforms
- Transitioned from a component-centric to system-level approach
- Involves integrating compute, network, and liquid cooling strategies
- Goal is to enable interoperability and integrated racks from early NPI phases.

Silicon Validation Infrastructure

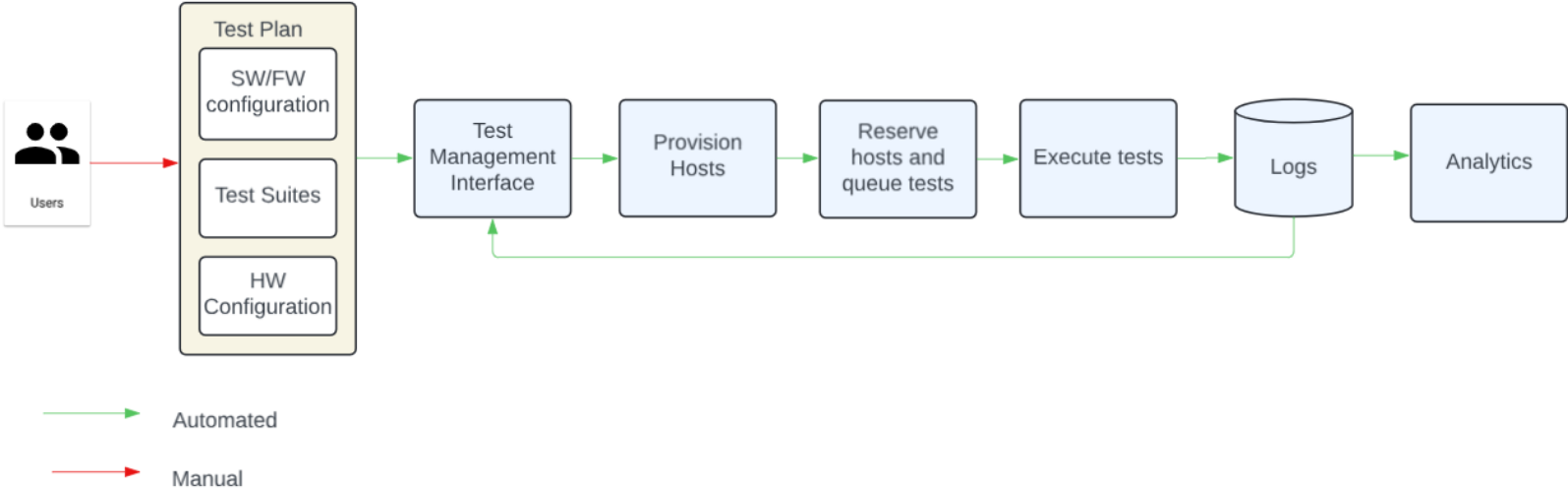
- Systems Validation at Meta happens in different environments
 - Emulation environment
 - Engineering Labs
 - Data Center
 - ODM/MFG
- Tests are developed to be environment agnostic
- Automation infra is scalable and portable across
 - Geographical regions
 - Validation environments
 - Platforms

Scalable Test Automation Infra



Silicon Automation Infrastructure

Silicon Automation Infra



Silicon Automation Infrastructure

- Primary focus is to left-shift validation and AI workload enablement to earlier NPI phases
- High signal and consistent tests that can be seamlessly ported
- Tests are automatically error categorized , results in faster debug/triage
- Continuous integration release engineering flow checks the SW/FW releases and maintains the overall health of the systems ensuring better utilization
- The automation infrastructure enables engineers with a seamless experience across all validation environments

ENGLAB INFRASTRUCTURE

Our Engineering Labs are located across US, Asia and Europe. These labs are foundational to infrastructure hardware delivery with following services:

- Test management and support
- Tooling
- Program Management
- Engineering Services
- Capacity
- Lab Ecosystem

ENGLAB INFRASTRUCTURE

- ***Test management and Tooling***
 - Before NPI validation, over 1 million tests were executed during the bring-up and post-silicon phases in our engineering labs, showcasing the rigorous and thorough validation process.
 - These tests were conducted using internal tools and processes, ensuring consistency, efficiency, and control over the testing environment.

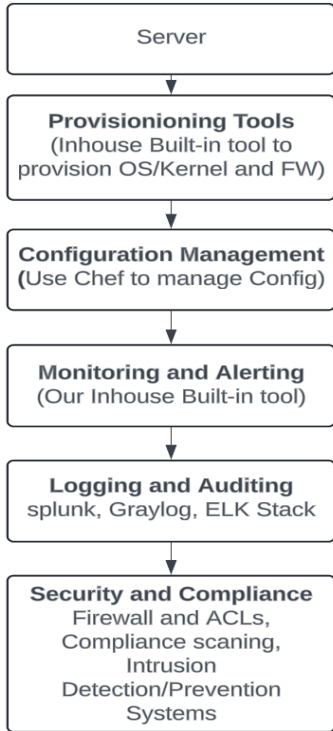
- ***Program Management and Engineering Services***
 - The increase in programs from 2020-2022 to 2022-2024 indicates a significant expansion of the silicon portfolio.
 - The increase in lab users suggests that more teams and individuals are utilizing the lab's resources, indicating increased collaboration and knowledge sharing across the organization.

ENGLAB INFRASTRUCTURE

- **Lab Ecosystem and Capacity**

- There is increase in lab locations implies that the infrastructure is expanding to accommodate the growing needs of the organization, providing more opportunities for teams to work on new AI projects and initiatives.
- Furthermore, the significant increase in devices brought up compared to 2020-2022 demonstrates the team's ability to adapt to new AI technologies and scale up operations to meet the demands of the rapidly evolving industry.

ENGLAB INFRASTRUCTURE

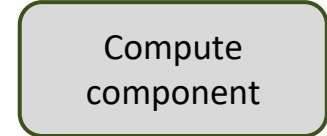
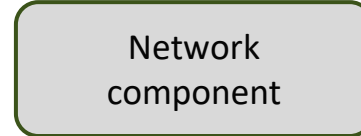


- The flow here shows the various tools and systems used in the bring up and management of servers at Meta.
- The provisioning tools are used to initially configure and set up the servers, while the configuration management tools are used to manage and maintain the servers' configurations over time.
- The monitoring and alerting tools are used to track the health and performance of the servers, while the logging and auditing tools are used to collect and analyze log data.
- Finally, the security and compliance tools are used to ensure the security and compliance of the servers with relevant regulations and standards.

Evolution with AI - Validation

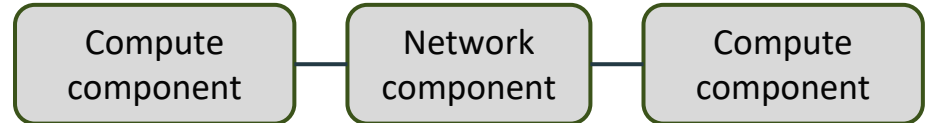
Component

- Component level validation (Power, performance, firmware) -> Englab environment (DiagOS)
- Prod testing using NOS -> Prod/DC environment
- Component not enough -> interop/link issues.



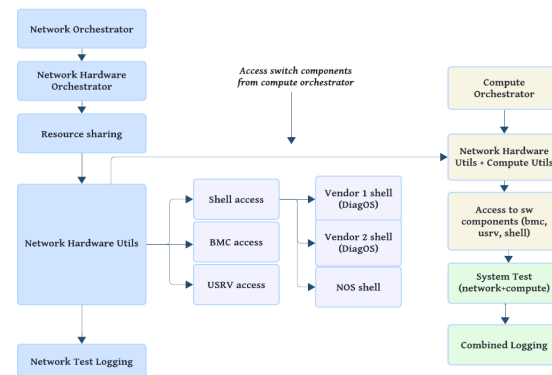
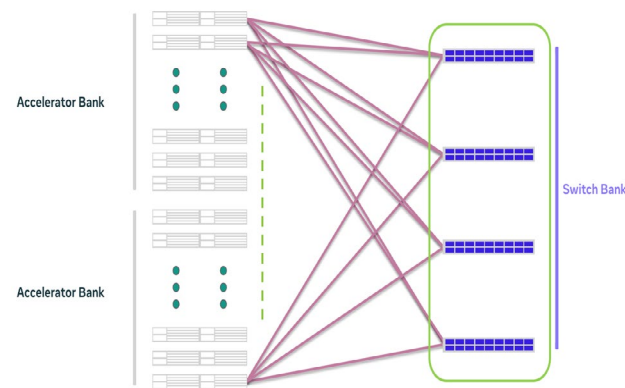
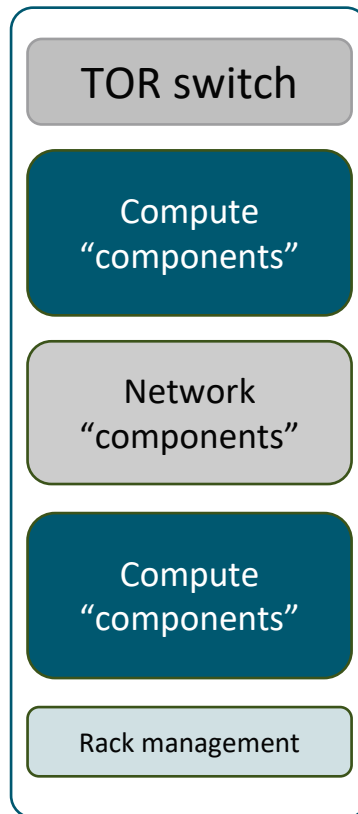
Component → Interop

- Setups to verify SerDes
- Capture any early link issues
- Compare and coordinate test results (Engineering lab vs At-scale datacenter)



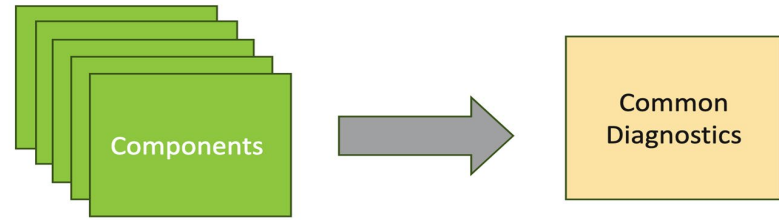
System Integration

- **Component** → **Interop** → **System**
- AI systems: **Integrated network and compute**
- **Component and Interop** : Not enough
- **System Integration Strategies** (frameworks, tooling, provisioning, firmware, data logging etc.)
- **Test categories**
 - Link stability tests, System health, Traffic tests, leak detection etc.
 - Pre-Silicon Post-Silicon
 - Pre-Provisioning Post-Provisioning
- **Rack Observability** : Link stats, Counters, SerDes

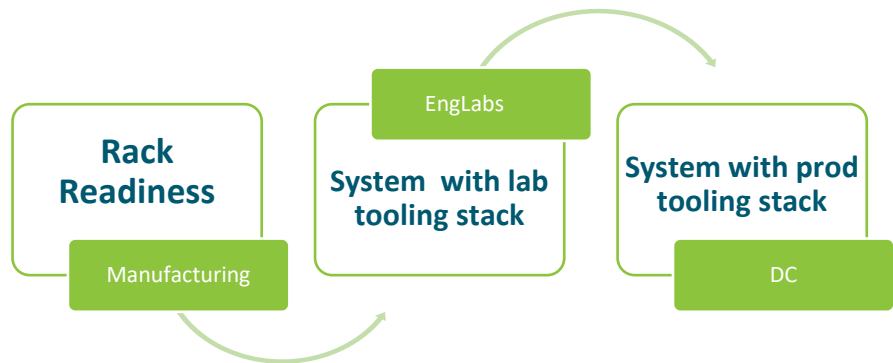


Network & Compute

- **Hardware management and bringup**
 - Physical lab space
 - Asset mgmt and system health checking
 - Provisioning readiness (Component focused)
- **Validation**
 - Test launchers, CI flows
 - Output formats
 - DUT definition
 - Testing ideology
 - Test data analysis and metrics
- **Processes**
 - Alignment. (Network and Compute programs)
 - Firmware managed
- **Overlap (tools and utilities). Collaboration opportunities**

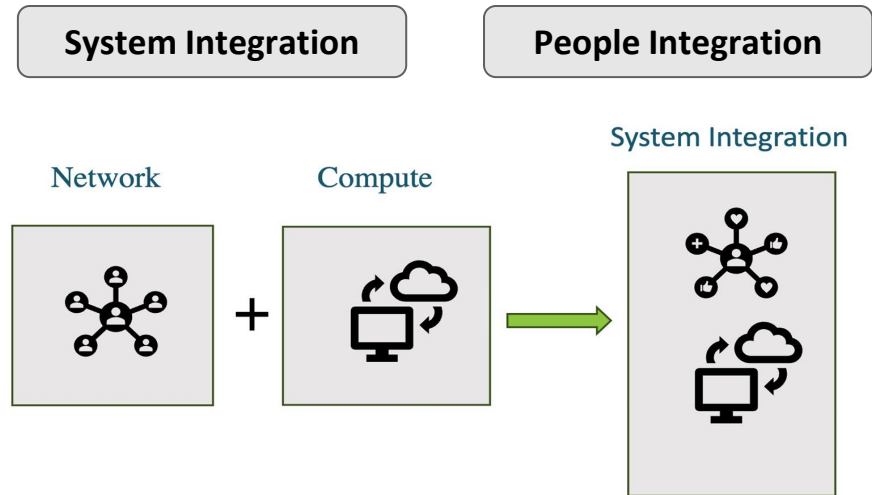


Observability info from components within AI Systems



Cohesive ecosystem → Cultural shift

- Breakdown Silos
- Improve collaboration, communication
- Increase efficiency
- Better support AI and data driven initiatives
- Flexible Infra
- Devops adoption
- Shared goals, performance metrics
- Innovation mindset , Openness
- Paradigm Shift



We have to work across teams and partners to deploy these complex systems in a timely manner

Thank you!!

Questions?

