# Exploring the Engineering Behind the Making of a Switch

BRKARC-3466

TOMORROW
starts here.

2

# Agenda

- Overview
- Concept
- System Design
- Mechanical / Physical Design
- Buffer Design
- Forwarding Design

- ASIC Engineering
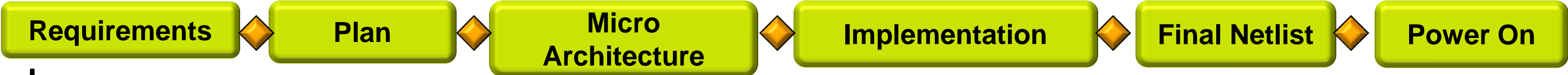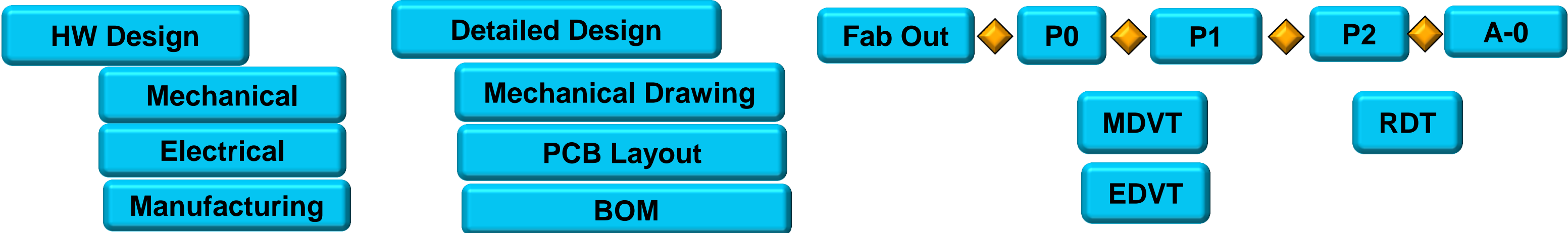- Hardware Engineering
- Software Engineering
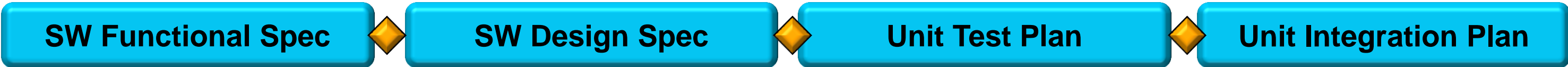
# Overview

# Timeline

**ASIC**

Product Requirements Document

| Requirements | Plan | Micro Architecture | Implementation | Final Netlist | Power On |

**Hardware**

| HW Design | | Detailed Design | | Fab Out | P0 | P1 | P2 | A-0 |

| Mechanical | Mechanical Drawing | | | MDVT | | RDT |

| Electrical | PCB Layout | | | EDVT |

| Manufacturing | BOM |

**Software**

| SW Functional Spec | SW Design Spec | Unit Test Plan | Unit Integration Plan |

**Software Test**

| Master Test Plans | Functional Test Plans | Automation | Regression | FCS |

Cisco*live!*

# Nexus 7000 and F2/F2E Modules



 Cisco Public

# Concept

# Concept
What customer problem will the product solve?

- Vision
- Market
- Cost
- Time to Market
- Differentiation
- Innovation

- Technology
- Life Cycle
- How Big?
- How many ports?
- Fixed vs Modular
- Backward Compatibility

# Nexus 7000 Vision (circa 2007)

**Cisco's End-to-End Data Centre Switching platform; providing solutions for 10G, 40G, and 100G for Access, Aggregation, and Core.**

**Consolidate IP, Storage, and IPC networks onto a single Ethernet fabric and deliver innovative features and services that provide value to our customers.**

# DC Evolutionary Innovation

**2013 Phase 3**

**2011 Phase 2**
½ Terabit Slot

**2009 Phase 1**
¼ Terabit Slot

**FCS**
**DC3**
**80G Slot**

10G Access
40 / 100G Aggregation
Unified Fabric

10GbE Access
10GbE Aggregation
Unified Fabric

10G Aggregation

Cisco Internal Slide CY2007
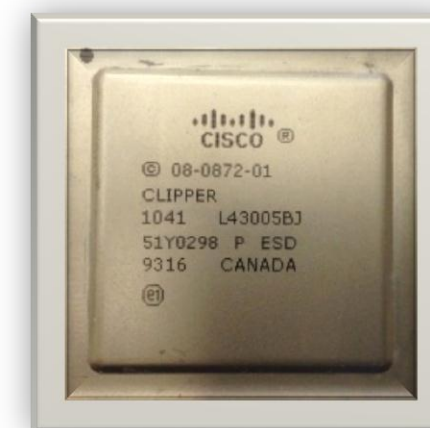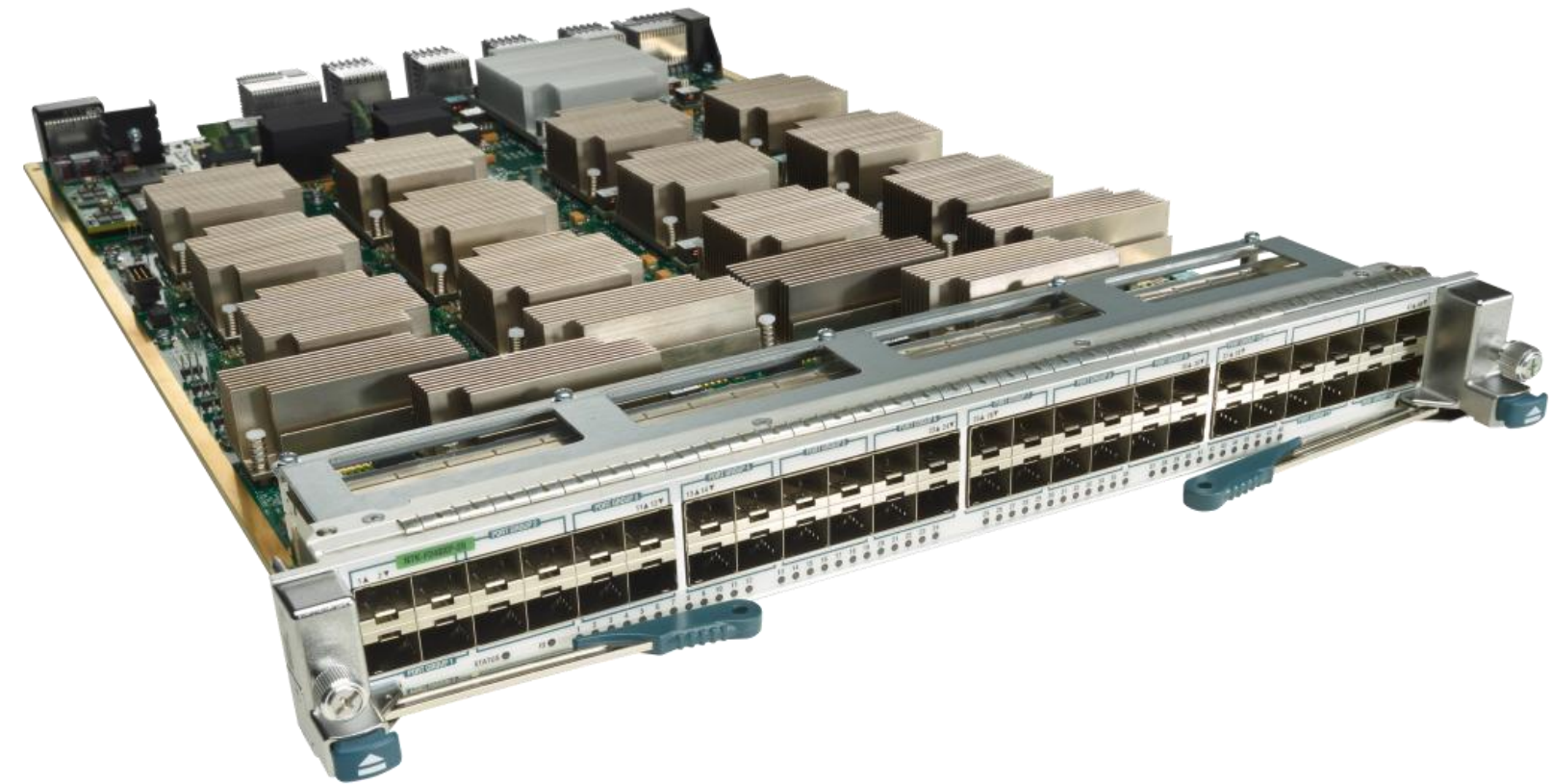
Cisco live!

# F2 Series
## High Level Goals

- 48 Ports 1/10G
  Line Rate 64 Bytes

- Low Latency

- L2MP, TRILL, FEX, FCoE, L3
  Forwarding

- Optimise for Data Centre

- IPv4 & IPv6
  Equal Performance

- Cost target

Cisco *live!*

# System Design

# Many Factors to Weigh
## Applicable to any Switch / Router Design

- Standards requirements
- Market requirements
- Designability
- Silicon technology
- Processor technology

- Manufacturability
- Time to market
- Flexibility
- Budget
- Modular / Fixed

# Many Factors to Weigh

Baseline Data Centre Switch Requirements

## Data Plane

- Buffering

- No packet drop

- Throughput

- Port count

- Modular

- No single point of failure

- In-order delivery

- Future protocol compatibility

## Control plane

- Modular

- Restartable (including active-active state handling)

- Non-disruptive code load & activation

- No single point of failure

- Scaleable

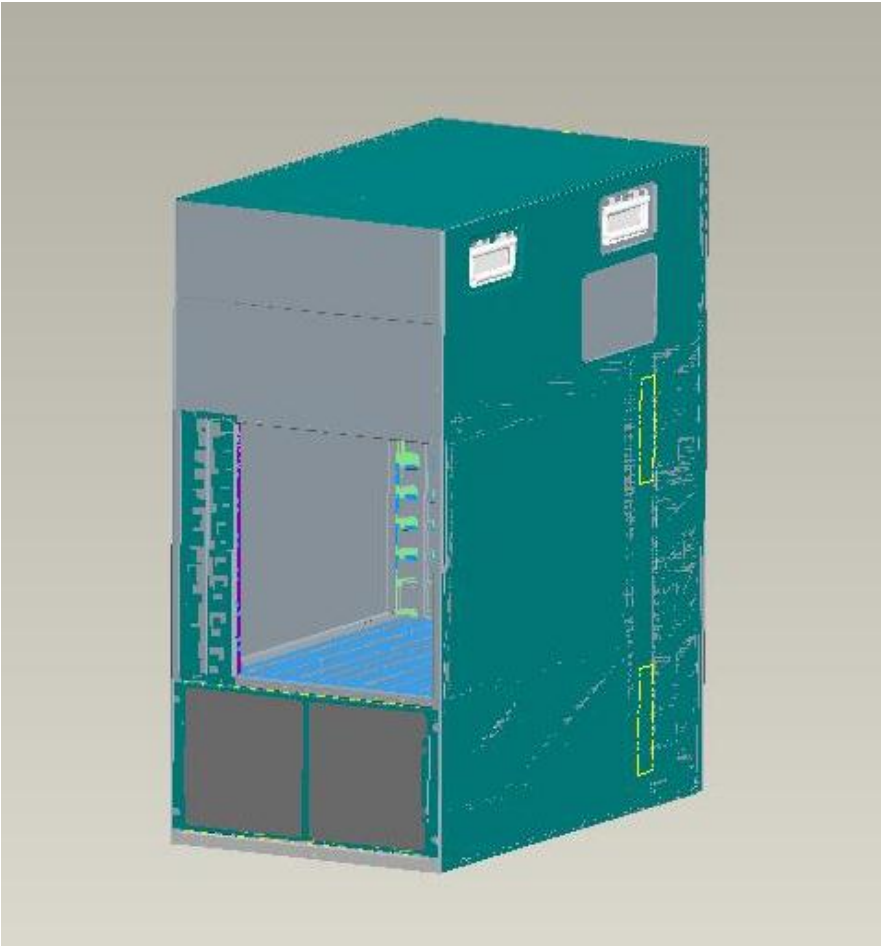- Unit Testable

- Future protocol compatibility

Cisco Public

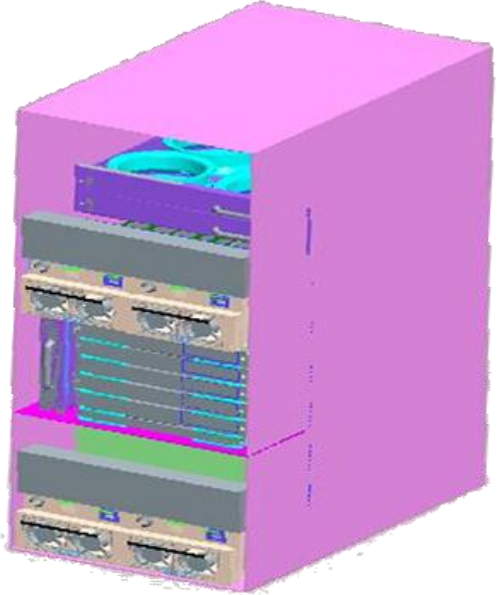Cisco live!
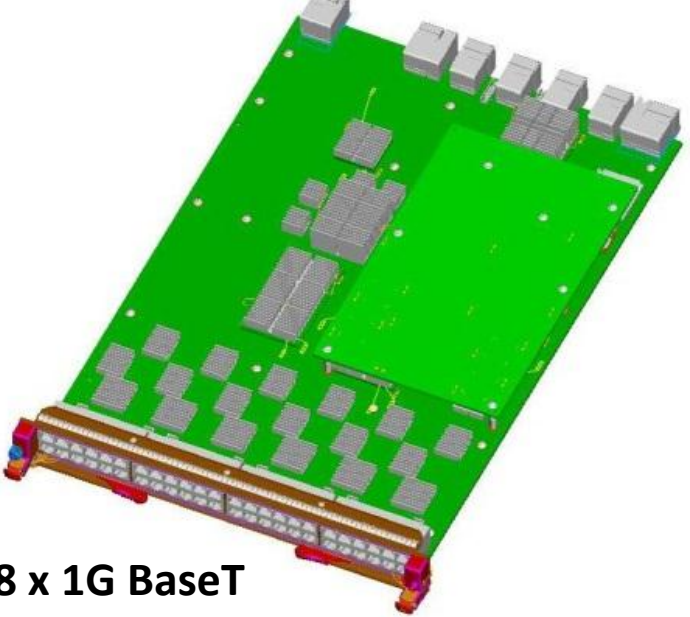
# Mechanical / Physical Design
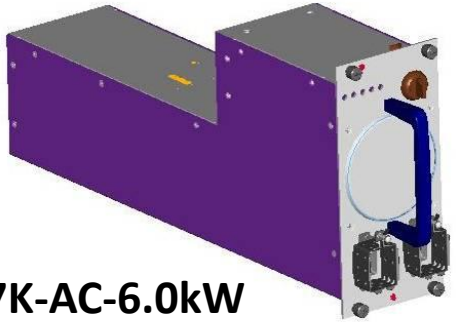
# Mechanical Design



Nexus 7010 Rear
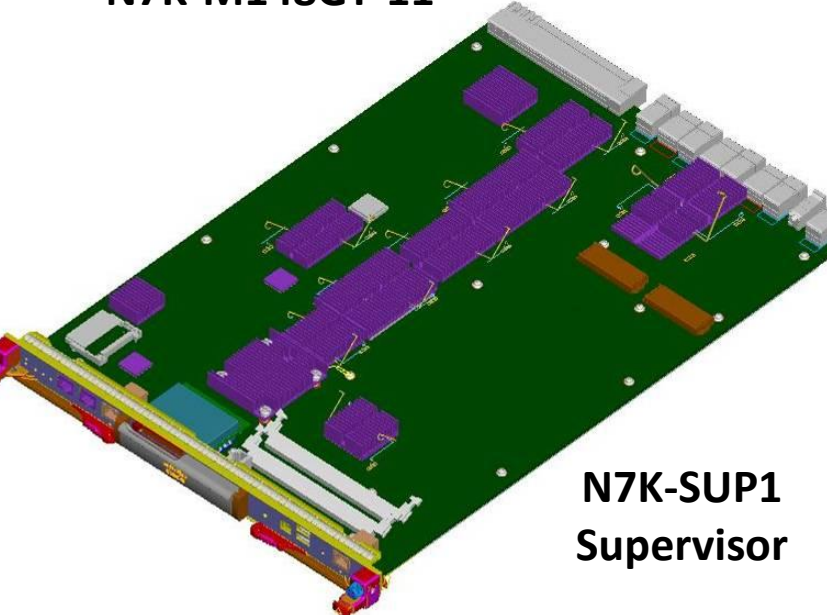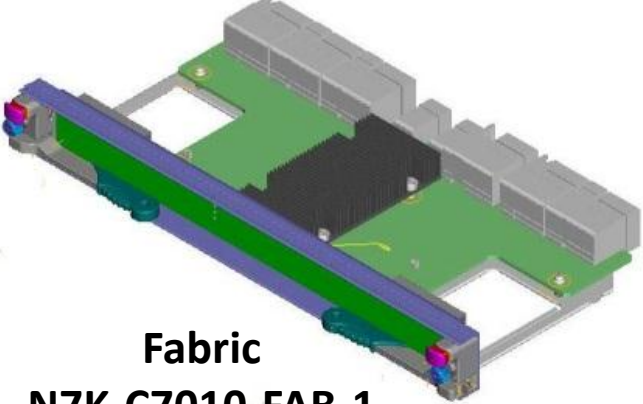


Nexus 7010 Front



**Nexus 7010 Rear**



**48 x 1G BaseT**
**N7K-M148GT-11**



**N7K-SUP1**
**Supervisor**



**N7K-AC-6.0kW**
**Power Supply**



**Fabric**
**N7K-C7010-FAB-1**



**32 x 10G SFP+**
**N7K-M132XP-12**

# Industrial Design / Usability



Ejectors

# Industrial Design / Usability

Cisco Public

# Buffer Design

# Memory Technology

DRAM:  Dynamic Random Access Memory
SDRAM: Synchronous Dynamic Random Access
eDRAM: embed Dynamic Random Access Memory

DDR3 Latency ~10ns
1 Transistor + 1 Capacitor
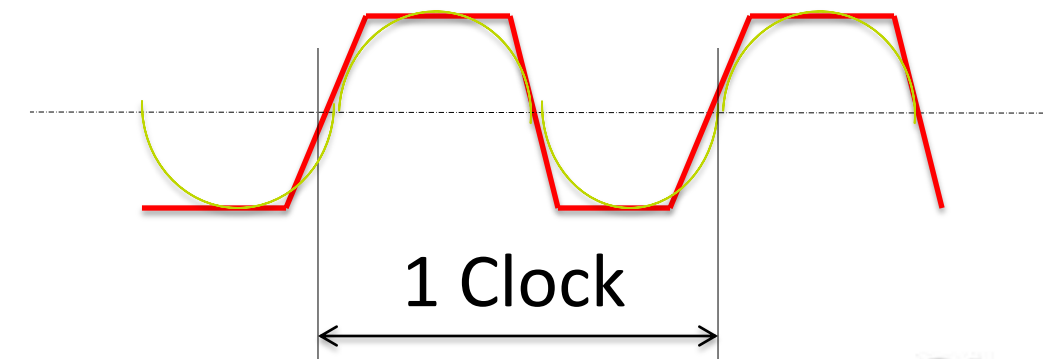Requires Refresh

SRAM:  Static Random Access Memory
SSRAM: Synchronous Static Random Access

SRAM Latency 1 cycle
6 Transistors

DDR: Double Data Rate – Transfer on Rising and Falling Edges of Clock
QDR: Quad Data Rate – Transfers on Rising and Falling and 2
                                    intermediate points between them

1 Clock

# Packet Buffer Design

Cell Pointers      Packet Ptr,Port, Queue

| Packet Ptr | Port | Queue |
|---|---|---|
|  | 1 | 2 |
|  | 4 | 2 |
|  |  |  |

xRAM Memory Banks

- **Fixed**
  - Memory segmented into a fixed cell size, like 128, 384, 512 bytes
  - If packet is smaller than cell size, then the left over space in page is unused
  - Easier to map to different memory banks

- **Packed**
  - Packets are placed head to tail into the packet memory
  - More efficient utilisation of the memory
  - More complex free space management and bank mapping

# One plus one does not equal Two
# Flow Balancing

Switch A ⟷ Switch B ≠ Switch A ⟷ Switch B

10 Links @ 1Gbps Each
Bandwidth = 10Gbps
Flow Bandwidth = 1Gbps
Serialisation Delay = 20uS

1 Link @ 10Gbps Each
Bandwidth = 10Gbps
Flow Bandwidth = 10Gbps
Serialisation Delay = 2uS

Cisco live!

# One Plus One Equal Two Word Spraying

1 Packet

X Bits

| N | ... | 5 | 4 | 3 | 2 | 1 |

Encode

FEC/ Scramble

**64/66**

5  1
6  2
7  3
8  4

FEC/DeScramble

Decode

Creates a link N x link speed, with small increase in Latency due to supporting difference in delays between physical lanes
Example 40/100G Ethernet, 180ns maximum skew between lanes

# Scrambing / Encoding

|  | Serdes (Gbps) | Encoding | Bandwdith after Encoding (Gbps) |
|---|---|---|---|
| PCI express v1 | 2.525 | 8b/10b | 2.02 |
| PCI express v2 | 5G | 8b/10b | 4 |
| PCI express v3 | 7.99 | 128b/130b | 7.867 |
| 10G Ethernet XAUI | 4 x 3.126G | 8b/10b | 10 |
| 10G Ethernet XFI | 10.3125 | 64b/66b | 10 |
| 40G Ethernet 4x XFI | 4x 10.3125 | 64b/66b | 40 |
| 100G Ethernet | 10 x 10.3125 | 64b/66b | 100 |
| 100G Ethernet 4x 25G | 4x 25.78125G | 64b/66b | 100 |
| 2G FC | 2.125 | 8b/10b | 1.7 |
| 4G FC | 4.250 | 8b/10b | 3.4 |
| 8G FC | 8.5 | 8b/10b | 6.8 |
| 16G FC | 14.025 | 64b/66b | 13.6 |

Cisco Public

# Single ASIC

- Scalability limited by memory bandwidth/size

- Typically optimised for fixed configuration

- Cost effective with small port counts

- Often used as building block

# Switch Architecture



Mesh

Crossbar

Clos / Fat Tree

Cisco Public

Cisco live!

# Complete System – Pull Fabric

# Forwarding Design

# High Level View of Forwarding

Parse Packet → Table Lookups → Forwarding Decision

- L2 Table
- L3 Table
- Classification

Cisco live!

# 10G Ethernet Forwarding Rate

## 1x10G Ethernet Forwarding Rate



10G Ethernet = 14.88Mpps @ 64 Bytes

67.2ns to receive a packet

100G Ethernet = 148.8Mpps @ 64 Bytes

6.72 ns to receive a packet

# Table Lookups
# CAMs, HASH Tables and *Tries

Input Key

TCAM

| 1 | 01001010 |
|---|----------|
| 2 | 010010XX |
| 3 | 01001XX0 |
| 4 | 01001XXX |

Hash Table

Trie

Result

Cisco live!

# CAMs

**Content Addressable Memory**

```
          01001110
             ↓
    1   01101010
    2   01101011
    3   01001110      ← Hit!
    4   01101100
             ↓
    3  ←  Result
```

**Ternary Content Addressable Memory**

```
    01001000    Lkup #1
    01001101    Lkup #2
    01001110    Lkup #3
             ↓
    1   01001010
    2   010010XX    ← Hit #1!
    3   01001XX0    ← Hit #3!
    4   01001XXX    ← Hit #2!
             ↓
    2  ←   Result #1
    4  ←   Result #2
    3  ←   Result #3
```

**Storing 1 bit in TCAM takes 10-12 transistors**

# Hash Tables

Input MAC Address

0000.c000.0001

Mathematical Functional produce value between 0 and Page Size

Compare if value in each page matches input value

Pages

Page Size

1 bit in SRAM takes 6 transistors

1 bit in DRAM takes 1 transistor

# Tries

- Many **different** *tries
  - Bitwise Trie
  - Balanced Trie
  - Patricia Trie
  - Fixed or Variable Stride Tries
- **Store** information in each leaf or pointer to table with information in it

# Algorithmic TCAMs

- From a software point of view looks like TCAM

  - May be all algorithmic or combination of TCAM and algorithmic

  - Software Driver takes TCAM representation and compiles the table to optimally utilise the underlying device

- Why? Algorithmic approaches allow the tables to scale with less than linear power increase

High Level View
Algorithmic TCAM

| 01001000 | Lkup #1 |
| 01001101 | Lkup #2 |
| 01001110 | Lkup #3 |

| 1 | 01001010 | |
| 2 | 010010XX | ← Hit #1! |
| 3 | 01001XX0 | ← Hit #3! |
| 4 | 01001XXX | ← Hit #2! |

2 ← Result #1

4 ← Result #2

3 ← Result #3

Implementation
Algorithmic TCAM

| 1 | 01001010 |
| 2 | 010010XX |
| 3 | 01001XX0 |
| 4 | 01001XXX |

# L3 Table: Design 1

| IPv4 Unicast FIB VRF / Prefix / Mask / Paths / Offset |
|---|
| 1 / 10.1.2.0 / 24 / 4 / 1 |
| 1 / 10.1.3.0 / 24 / 1 / 5 |
| 3 / 10.1.2.0 / 24 / 2 / 9 |
| 3 / 10.1.3.0 / 24 / 2 / 9 |
| |
| |

**H A S H**

| Rewrite Information |
|---|
| ADJ 1 - Rewrite SRC A+DST A MAC |
| ADJ 2 - Rewrite SRC A+DST B MAC |
| ADJ 3 - Rewrite SRC A+DST C MAC |
| ADJ 4 - Rewrite SRC A+DST D MAC |
| ADJ 5 - Rewrite SRC A+DST D MAC |
| ADJ 6 - Rewrite SRC A+DST F MAC |
| ADJ 7 - Rewrite SRC A+DST G MAC |
| ADJ 8 - Rewrite SRC A+DST H MAC |
| ADJ 9 - Rewrite SRC A+DST I MAC |
| ADJ 10 - Rewrite SRC A+DST J MAC |

Software View

# L3 Table: Design 2

| IPv4/v6 Unicast FIB VPN / Prefix / Mask / Paths / Offset |
|---|
| 1 / 10.1.2.0 / 24 / 4 / 1 |
| 1 / 10.1.3.0 / 24 / 1 / 5 |
| 3 / 10.1.2.0 / 24 / 2 / 6 |
| 3 / 10.1.3.0 / 24 / 2 / 6 |
| |
| |

**H A S H**

| Path Table |
|---|
| Path 1 |
| Path 2 |
| Path 3 |
| Path 4 |
| Path 1 |
| Path 1 |
| Path 2 |
| |
| |
| |

| Rewrite Information |
|---|
| ADJ 1 - Rewrite SRC A+DST A MAC |
| ADJ 2 - Rewrite SRC A+DST B MAC |
| ADJ 3 - Rewrite SRC A+DST C MAC |
| ADJ 4 - Rewrite SRC A+DST D MAC |
| ADJ 5 - Rewrite SRC A+DST E MAC |
| ADJ 6 - Rewrite SRC A+DST F MAC |
| ADJ 7 - Rewrite SRC A+DST G MAC |
| ADJ 8 - Rewrite SRC A+DST H MAC |
| ADJ 9 - Rewrite SRC A+DST I MAC |
| ADJ 10 - Rewrite SRC A+DST J MAC |

Software View

# L2 Table / Host Table / FIB
## Common Optimisation

- Hash tables take less space than TCAMs and Tries

- Instead of placing /32 or /128 entries for host entries into the FIB, place them into the hash table

- Common for the L2 table and the Host table to share the same memory

- Allows for the FIB Table to be smaller since it does not need to contain single path /32 and /128 entries

# Forwarding Design

## Design 1

```
Parse
Packet
```

```
L2 Table
```
```
L3 Table
```
```
Ingress
Security ACLs
```
```
Ingress QoS
ACL
```

```
Adjacency
Table
```

```
Egress
Security
ACLs
```
```
Egress QoS
ACL
```

```
Input / Output
Policing
```

```
Fwd Decision
```

```
Update
Statistics
```

## Design 2

```
Parse
Packet
```
```
L2 Table
```
```
VPN
CAM
```

```
L3 Table (x2)
```
```
Ingress Security
ACLs
```
```
Ingress QoS
ACL
```

```
Adjacency
Table
```

```
Egress
Security
ACLs
```
```
Egress QoS
ACL
```

```
Input / Output
Policing
```

```
Fwd
Decision
```

```
Update
Statistics
```

# Forwarding Design

## Design 2

| Parse Packet | | L2 Table | | VPN Table | | L3 Table (x2) | | Adjacency Table | | Egress Security ACLs | | Input / Output Policing | | Fwd Decision | | Update Statistics |

- L3 Table (x2)
- Ingress Security ACLs
- Ingress QoS ACL
- Egress Security ACLs
- Egress QoS ACL

Parse Packet → L2 Table → VPN Table → [Ingress Security ACLs / Ingress QoS ACL] → Adjacency Table → [Egress Security ACLs / Egress QoS ACL] → Input / Output Policing → Fwd Decision → Update Statistics

## Design 3

Parse Packet → L2 Table → VPN Table → [Ingress Security ACLs / Ingress QoS ACL] → L3 Table (x2) → Input Policing → Adjacency Table → [Egress Security ACLs / Egress QoS ACL] → Output Policing → Fwd Decision → Update Statistics
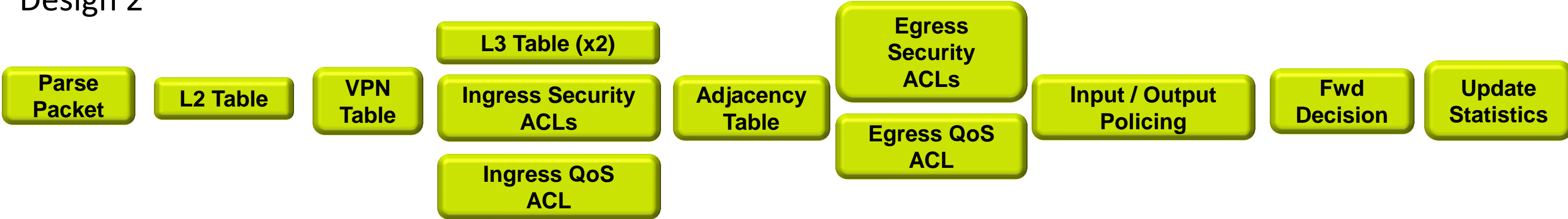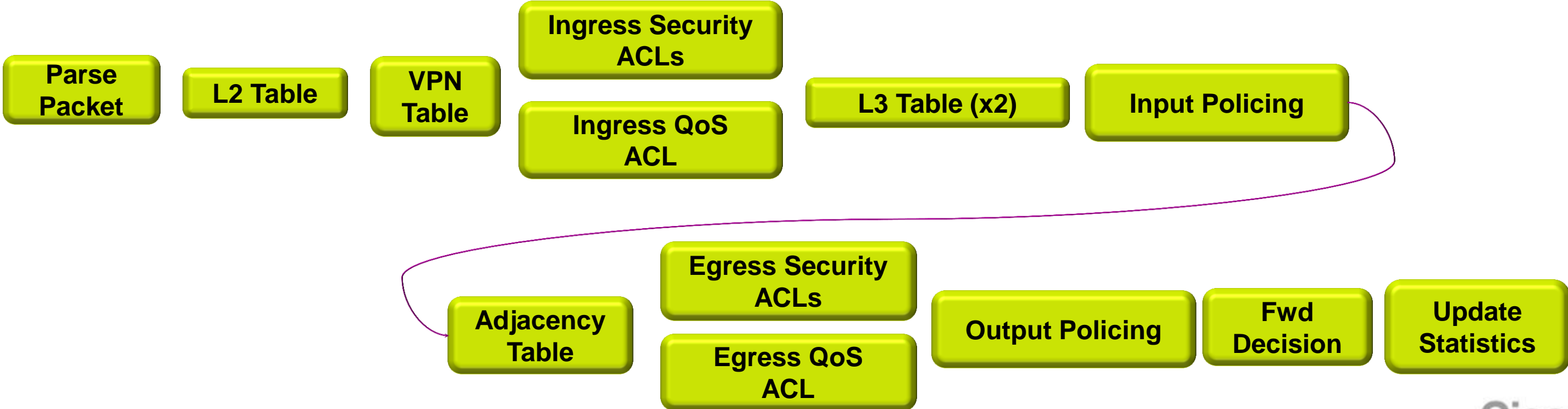
# References

- Network Algorithmics,: An Interdisciplinary Approach to Designing Fast Networked Devices George Varghese

- Art of Computer Programming Vol 1-4, Donald E. Knuth

- Introduction to Algorithms, Third Edition
  Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest and Clifford Stein

- IEEE SIGCOMM Papers

# ASIC Engineering

# ASICs vs FPGAs

ASIC - Application Specific Integrated Circuit

- A finished IC which is built to the exact specification & functionality of the customer
- Can make optimal use of the underlying silicon circuits
- Low part cost, High upfront investment
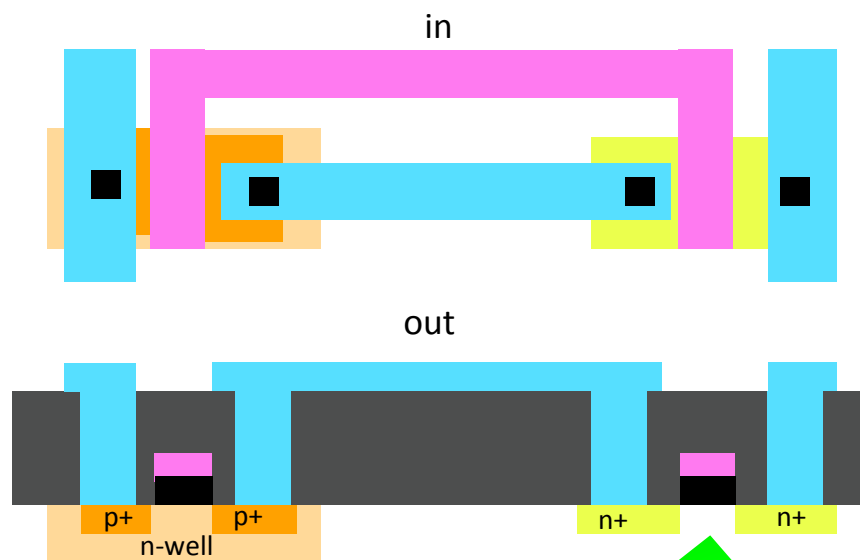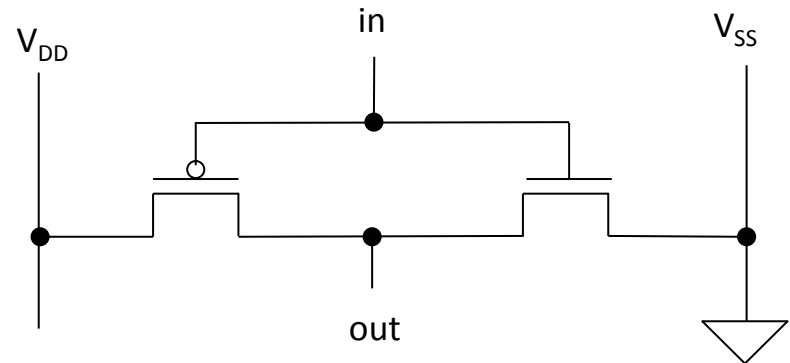- Significant development time

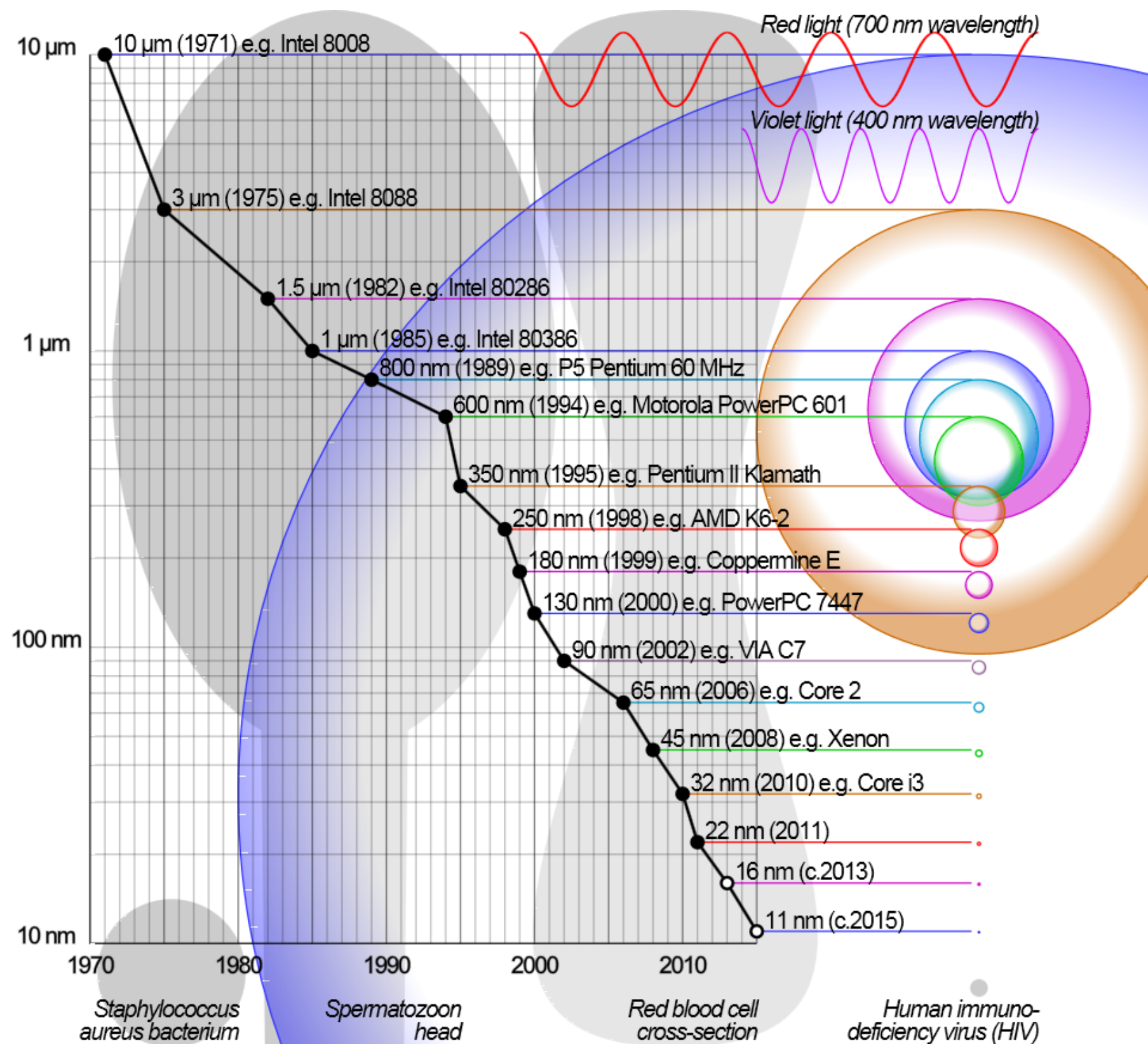FPGA (EPLD) Field Programmable Gate Array

- An IC that can be configured with the required functionality <u>after</u> it is installed into a target system
- Flexibility vs. sub-optimal use of underlying silicon circuits
- Higher part cost
- Shorter development time
- Main players: Xilinx, Altera

# CMOS



$V_{DD}$    in    $V_{SS}$

in

out

p+    p+         n+    n+
n-well

**"Feature size"**
This dimension is what Moore's Law is all about!



Red light (700 nm wavelength)

10 µm    10 µm (1971) e.g. Intel 8008

Violet light (400 nm wavelength)

3 µm (1975) e.g. Intel 8088

1.5 µm (1982) e.g. Intel 80286

1 µm    1 µm (1985) e.g. Intel 80386
800 nm (1989) e.g. P5 Pentium 60 MHz
600 nm (1994) e.g. Motorola PowerPC 601
350 nm (1995) e.g. Pentium II Klamath
250 nm (1998) e.g. AMD K6-2
180 nm (1999) e.g. Coppermine E
130 nm (2000) e.g. PowerPC 7447

100 nm    90 nm (2002) e.g. VIA C7
65 nm (2006) e.g. Core 2
45 nm (2008) e.g. Xenon
32 nm (2010) e.g. Core i3
22 nm (2011)
16 nm (c.2013)

10 nm    11 nm (c.2015)

1970    1980    1990    2000    2010

Staphylococcus    Spermatozoon    Red blood cell    Human immuno-
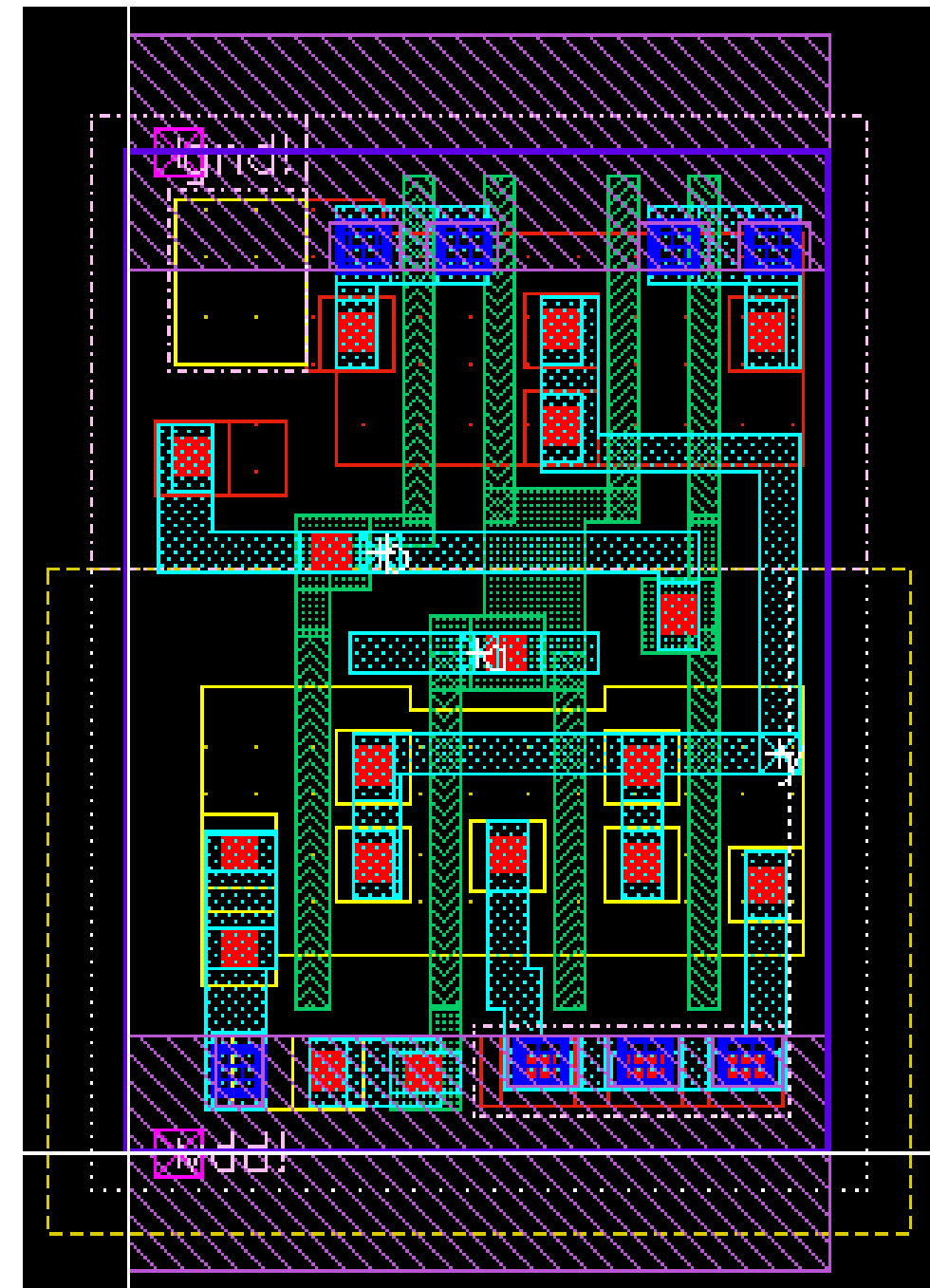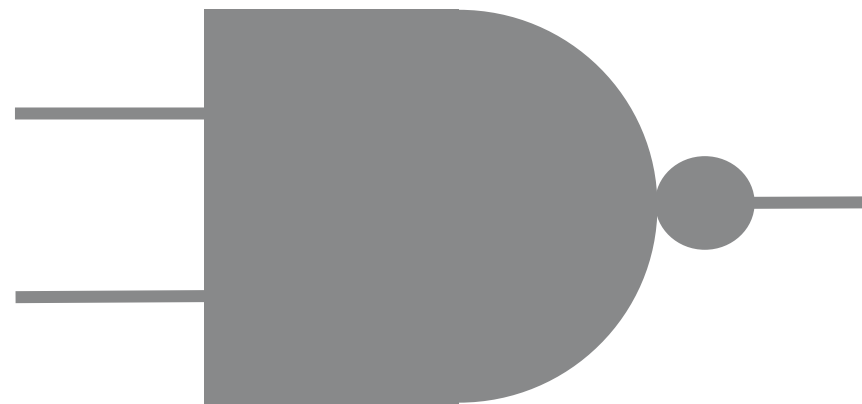aureus bacterium    head    cross-section    deficiency virus (HIV)

http://en.wikipedia.org/wiki/Semiconductor_device_fabrication

1.6 X Increase in usable gates between process nodes

# Gates

```
module nand2(a,b,c)
  input a,b;
  ouput c;
begin
    c <= !(a & b);
end
```
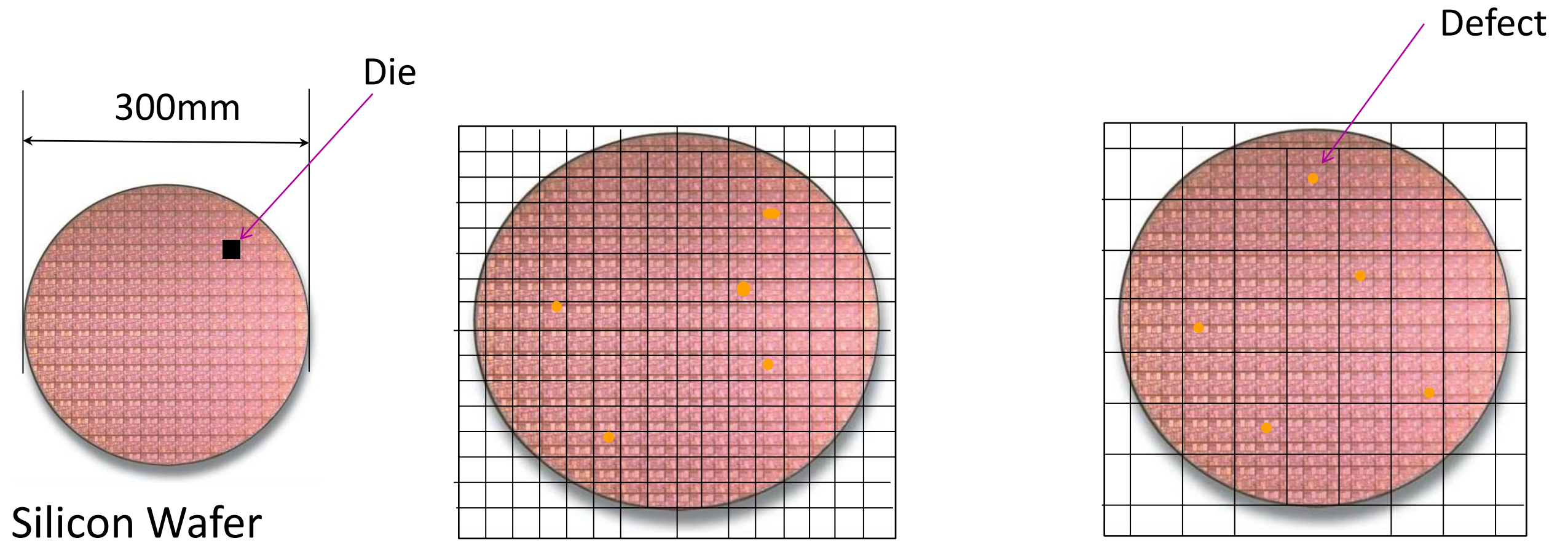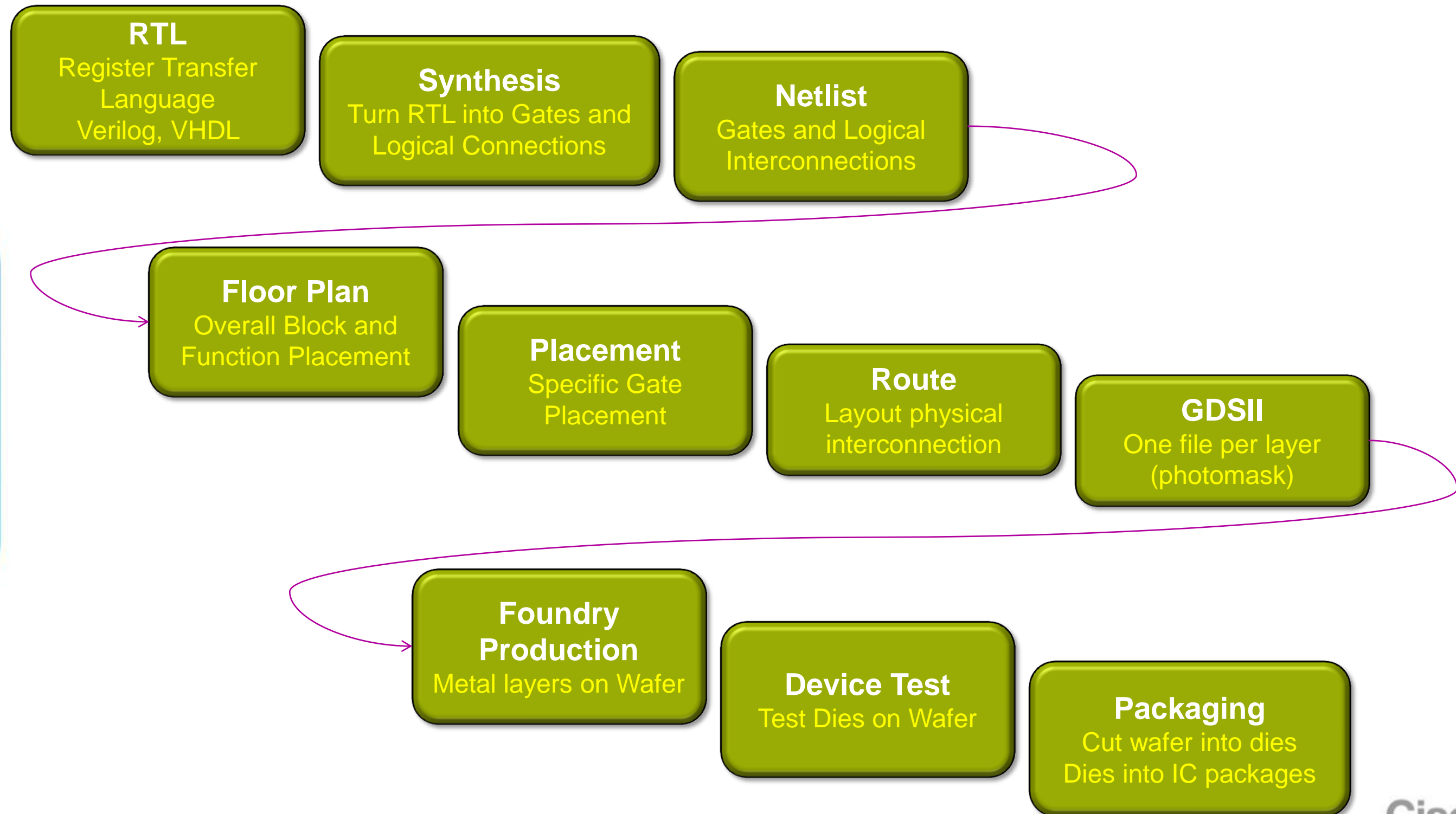
# Why is Die Size Important?



With same number of defects per wafer,
smaller die size results in higher yield per wafer

# Integrated Circuit Production

**RTL**
Register Transfer Language
Verilog, VHDL

**Synthesis**
Turn RTL into Gates and Logical Connections

**Netlist**
Gates and Logical Interconnections

**Floor Plan**
Overall Block and Function Placement

**Placement**
Specific Gate Placement

**Route**
Layout physical interconnection

**GDSII**
One file per layer (photomask)

**Foundry Production**
Metal layers on Wafer

**Device Test**
Test Dies on Wafer

**Packaging**
Cut wafer into dies
Dies into IC packages

Cisco *live!*

# Integrated Circuit Production

**RTL**
Register Transfer Language
Verilog, VDHL

**Synthesis**
Turn RTL into Gates and Logical Connections

**Netlist**
Gates and Logical Interconnections

ASIC Customer
- Cisco

**Floor Plan**
Overall Block and Function Placement

**Placement**
Specific Gate Placement

**Route**
Layout physical interconnection

**GDSII**
One file per layer (photomask)

ASIC Vendor
- Avago, IBM, TI, ST Micro
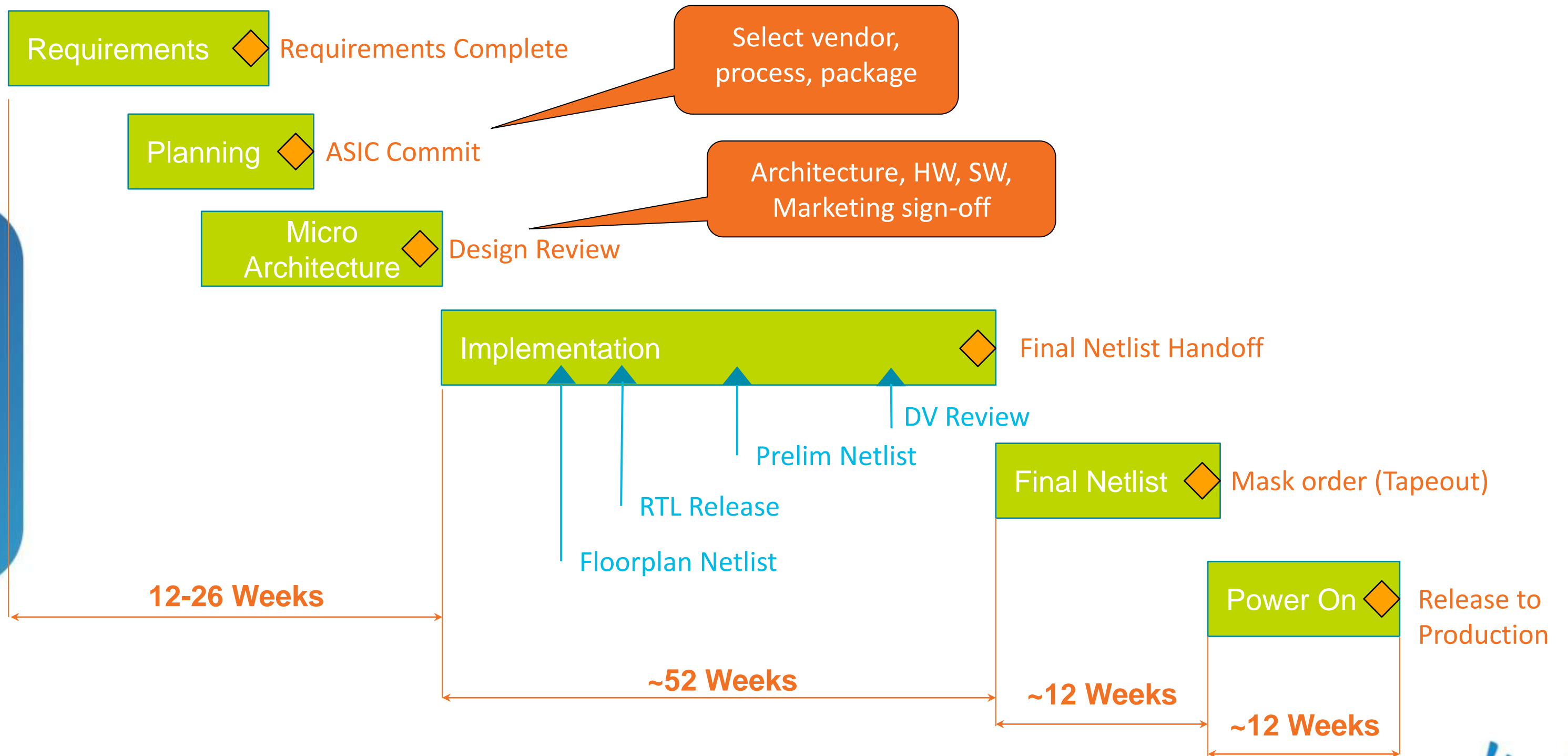- COT - Cisco

**Foundry Production**
Metal layers on Wafer

**Device Test**
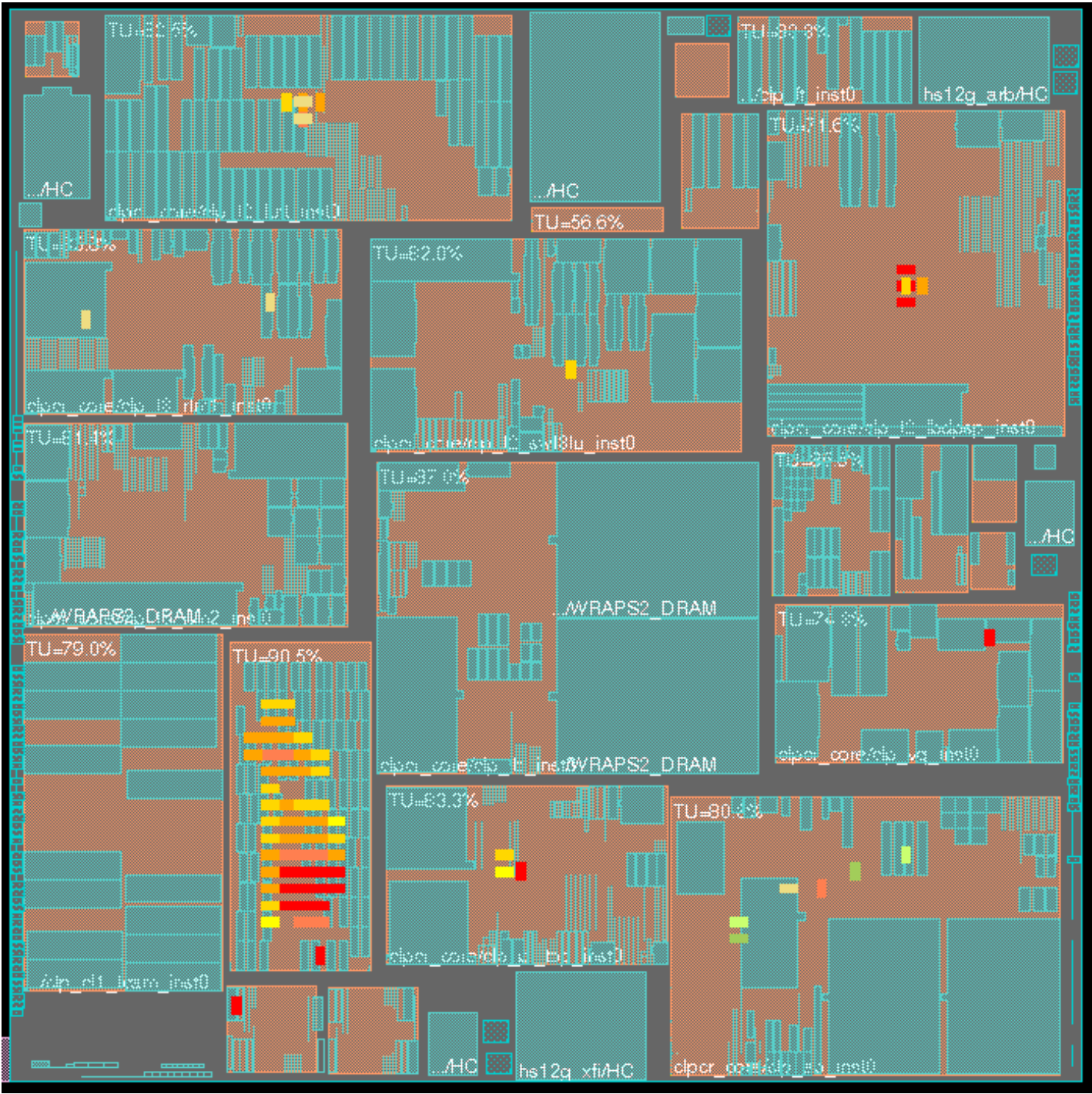Test Dies on Wafer

**Packaging**
Cut wafer into dies
Dies into IC packages

Silicon Foundry
- IBM, TSMC, Global Foundries

Cisco live!

# ASIC Design Process

Requirements — ◇ Requirements Complete

Select vendor, process, package

Planning — ◇ ASIC Commit

Architecture, HW, SW, Marketing sign-off

Micro Architecture ◇ Design Review

Implementation ◇ Final Netlist Handoff

DV Review

Prelim Netlist

Final Netlist ◇ Mask order (Tapeout)

RTL Release

Floorplan Netlist

**12-26 Weeks**

**~52 Weeks**

Power On ◇ Release to Production

**~12 Weeks**

**~12 Weeks**

Cisco live!

# F2/F2E ASIC - Clipper



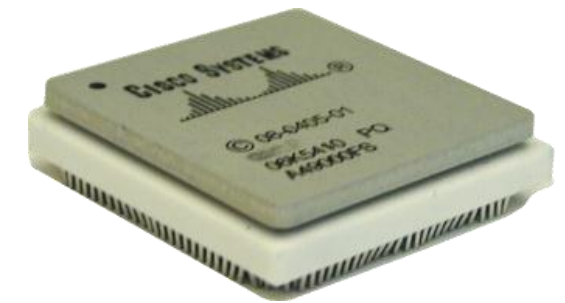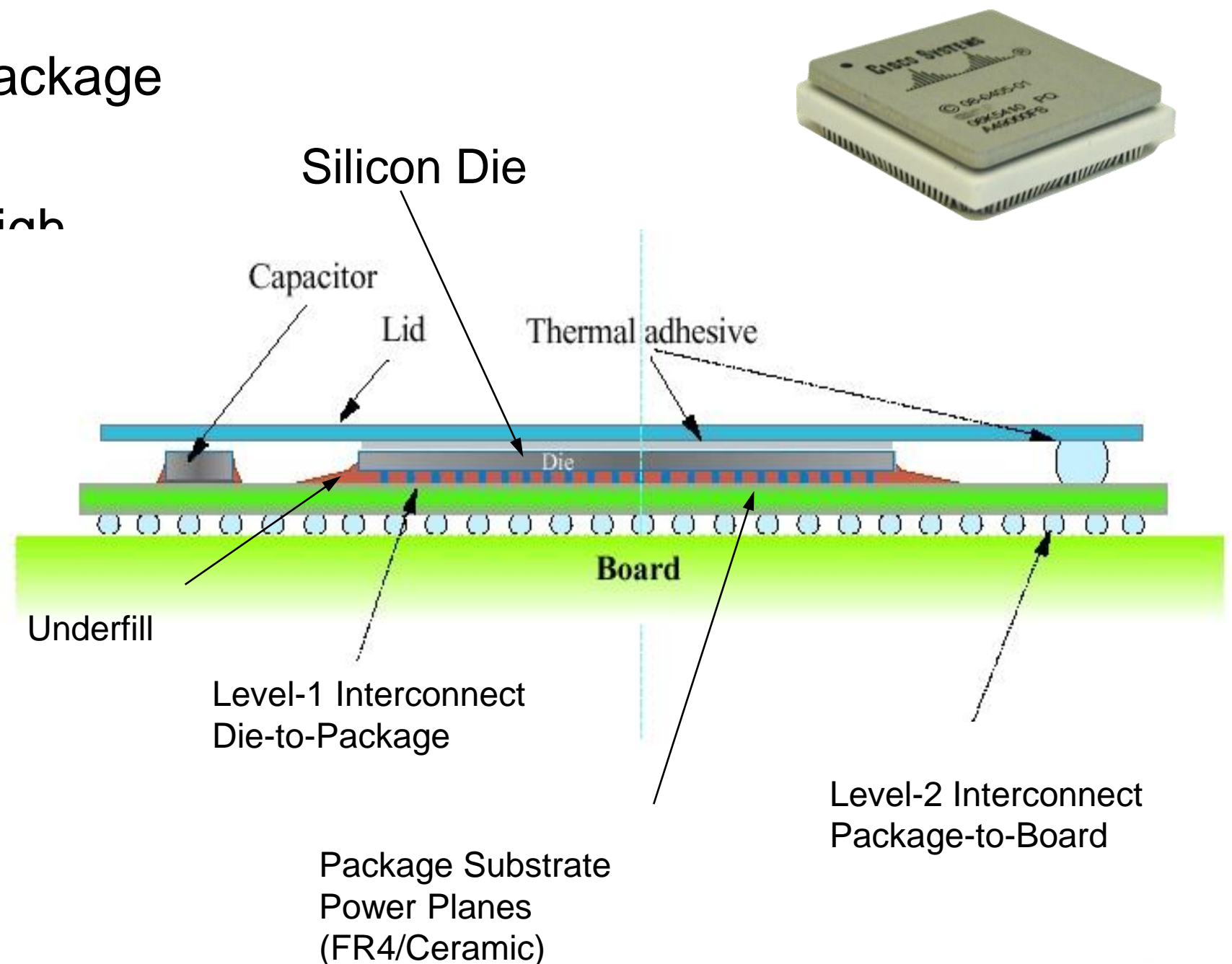| | F2 | F2E |
|---|---|---|
| Technology | IBM Cu-65 | IBM Cu-45 |
| Die Size | 18.0x18.3mm | 12.28x12.28mm |
| Total SRAM | 33.3Mb | 33.3Mb |
| Total eDRAM | 134Mb | 134Mb |
| Total TCAM | 2.94Mb | 2.94Mb |
| Register Array | 1.34Mb | 1.34Mb |
| Logic Gates | 45M | 45M |
| Signal Pin | 186 | 186 |
| Package IO | 840 | 840 |

Cisco Public

# Memory and Packet Corruption Protection

- No ECC or Parity – no way to determine if a software or hardware problem

- Parity – will detect single bit errors

- ECC – will detect 2 bit errors, and correct single bit

- Parity and ECC apply to a word (32 or 64 bits)

- CRC – Detect if a set of bytes (normally a packet) has been corrupted

# ASIC Packaging

- Electrical parasitics of the chip package are critical

- Impacts electrical properties of high speed signals

- Manufacturing tolerances constrain minimum ball pitch

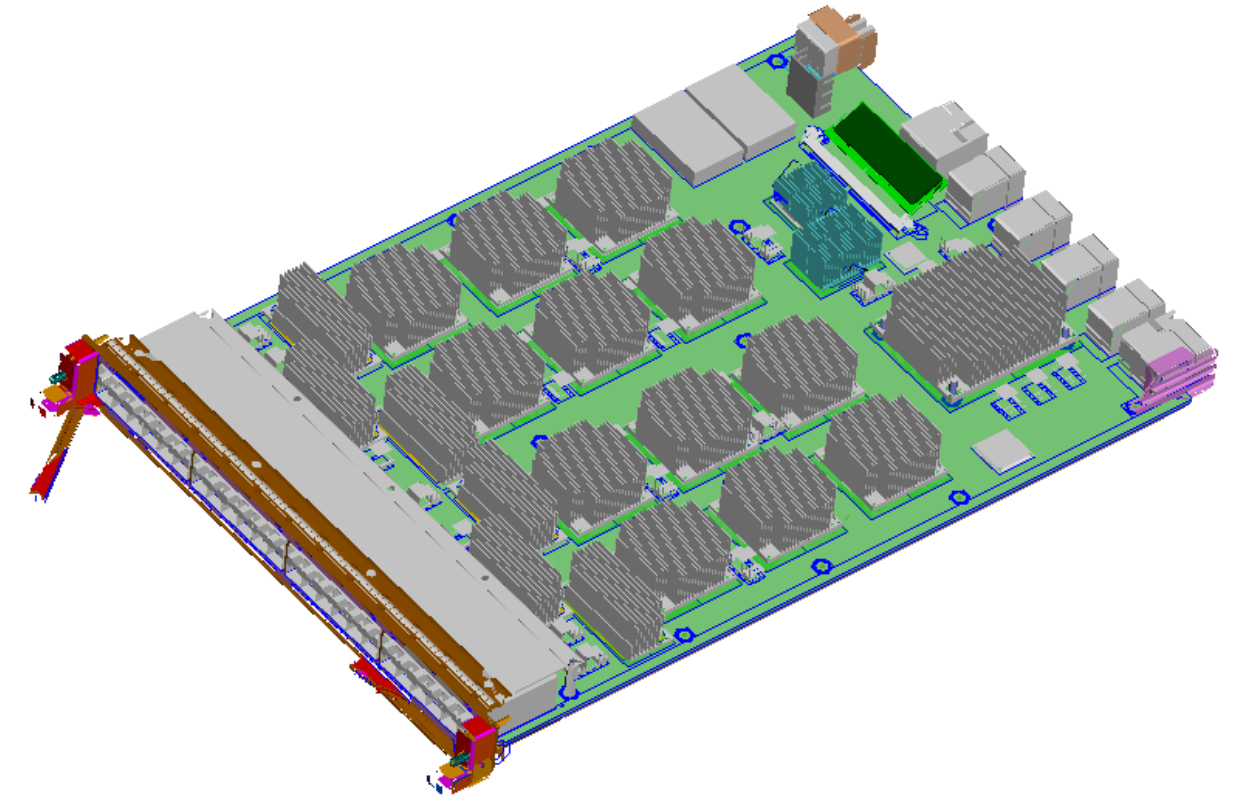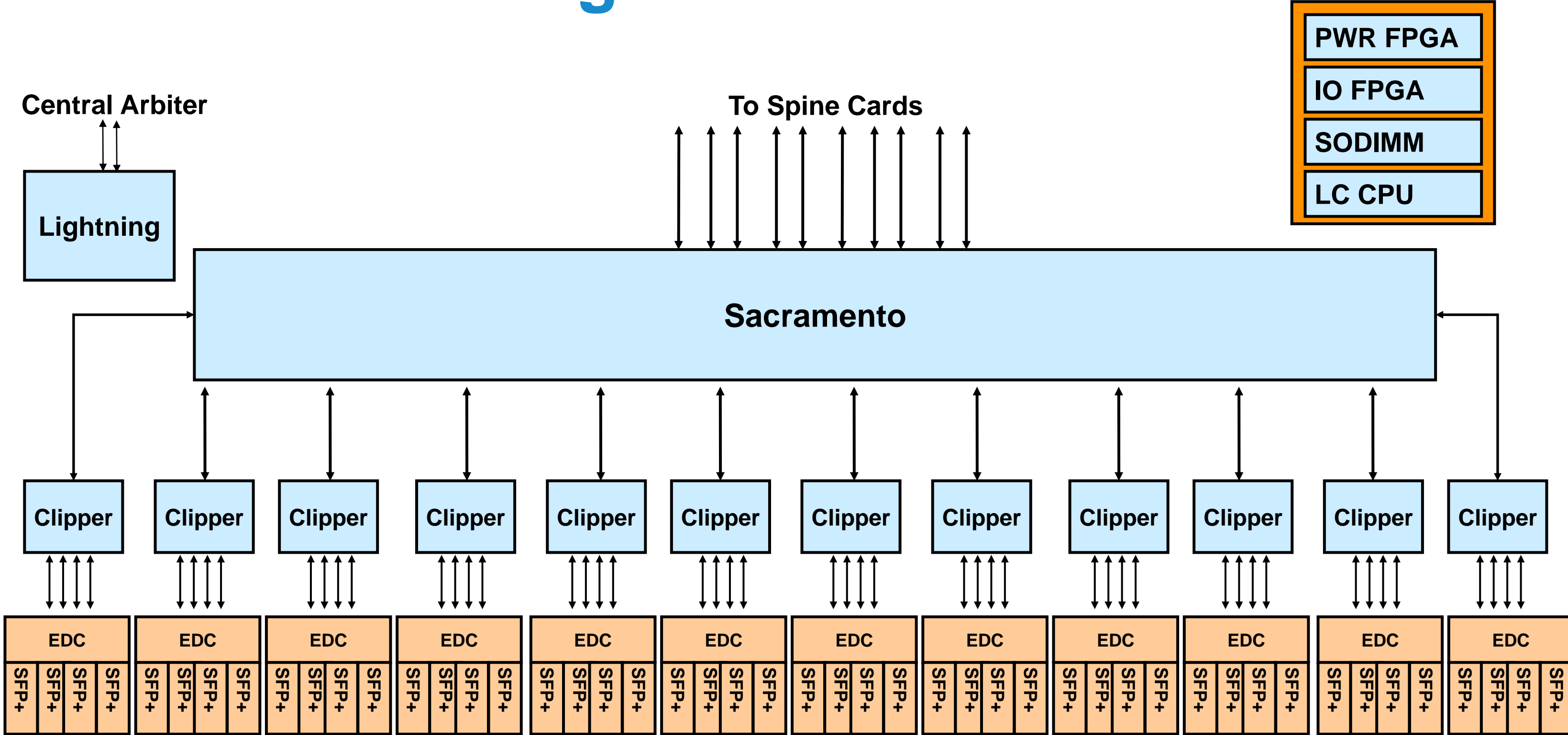- Limit to number of available signal I/O pins

Silicon Die

Capacitor

Lid

Thermal adhesive

Die

Board

Underfill

Level-1 Interconnect
Die-to-Package

Level-2 Interconnect
Package-to-Board

Package Substrate
Power Planes
(FR4/Ceramic)

# References

- Indistinguishable From Magic: Manufacturing Modern Computer Chips
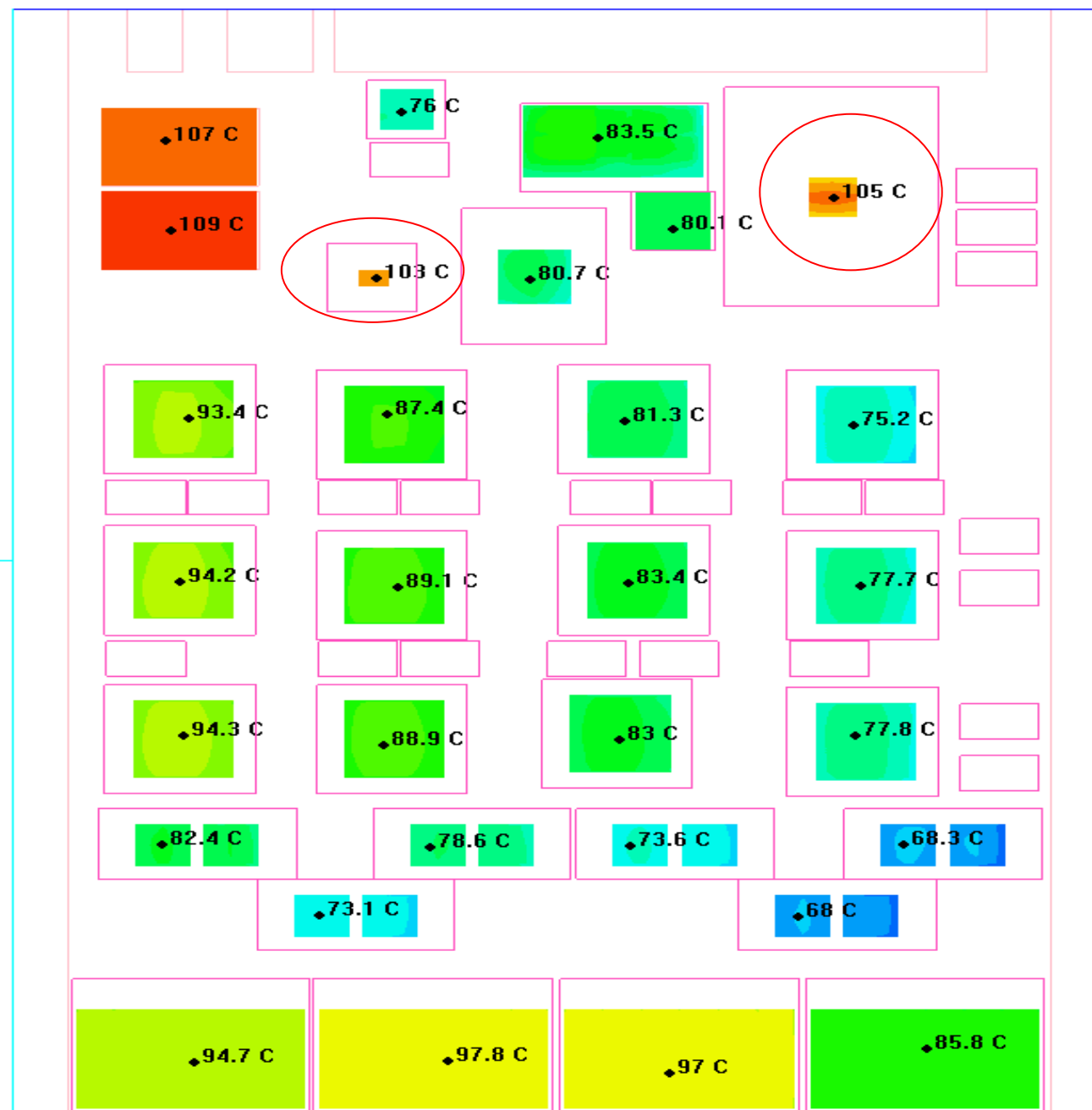  http://www.youtube.com/watch?v=NGFhc8R_uO4

 Cisco Public

# Hardware Engineering

# F2 Block Diagram

**Central Arbiter**

**To Spine Cards**

| PWR FPGA |
|---|
| IO FPGA |
| SODIMM |
| LC CPU |

**Lightning**

**Sacramento**

| Clipper | Clipper | Clipper | Clipper | Clipper | Clipper | Clipper | Clipper | Clipper | Clipper | Clipper | Clipper | Clipper |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

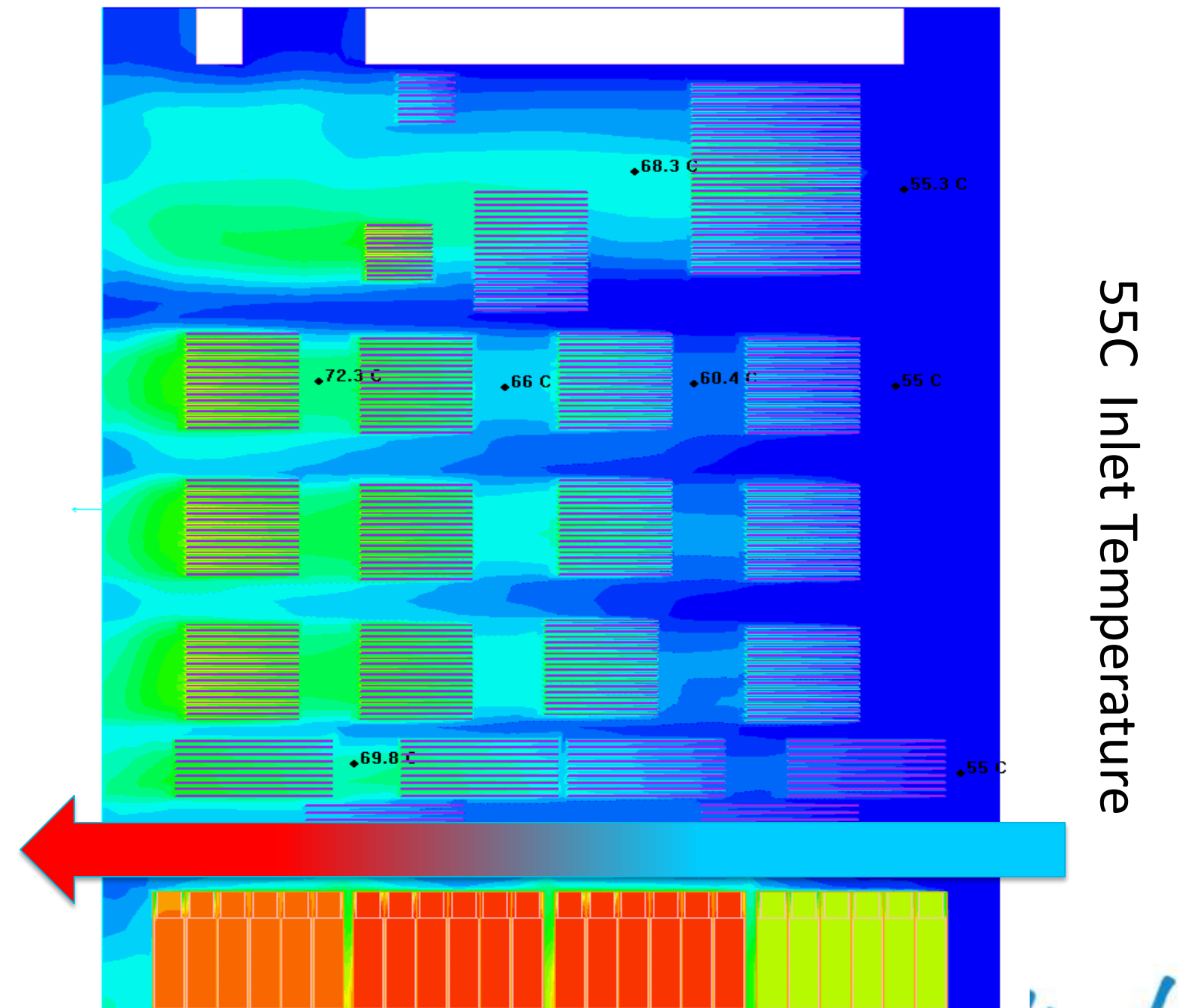| EDC | EDC | EDC | EDC | EDC | EDC | EDC | EDC | EDC | EDC | EDC | EDC | EDC |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

SFP+ SFP+ SFP+ SFP+ (repeated per EDC block)

Cisco *live!*

# Thermal Modeling

Component Case Temperatures

Temperature Contours

55C Inlet Temperature

Cisco Public
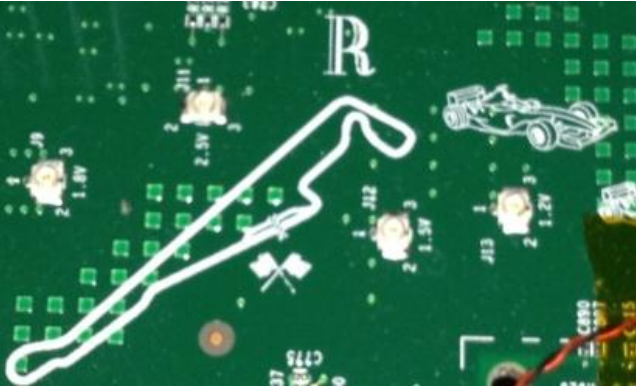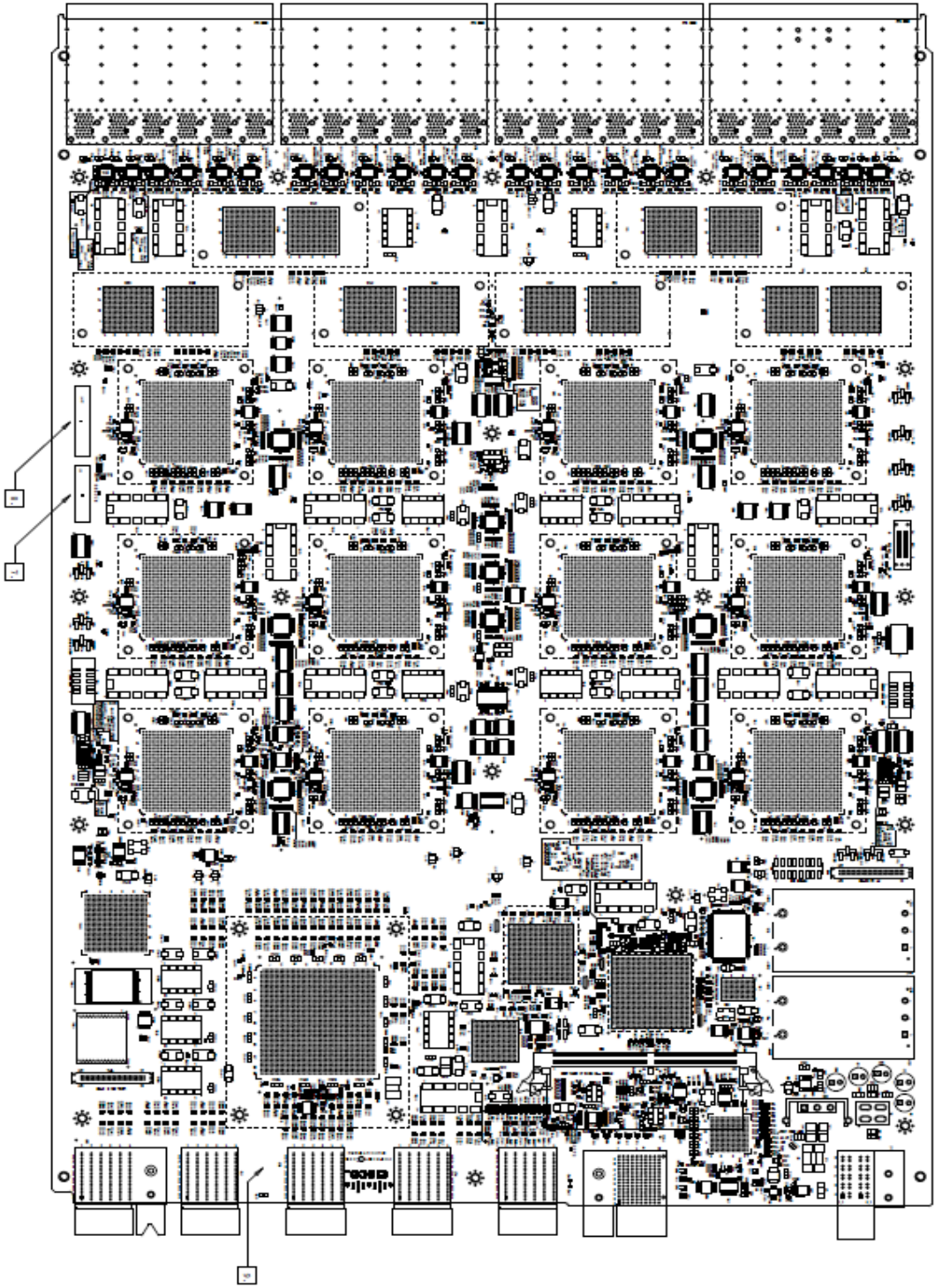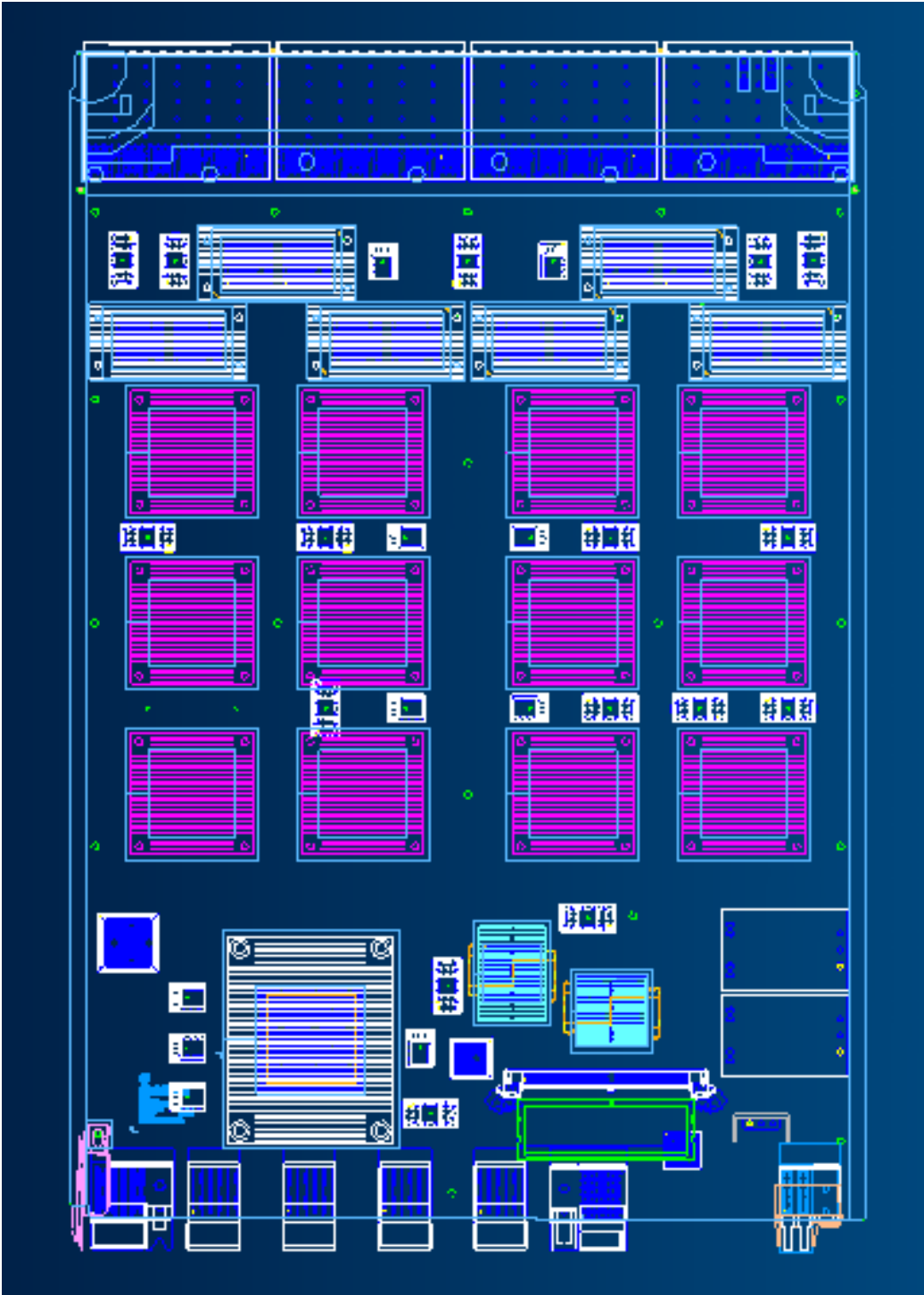
# Electrical / Mechanical Layout

20 Layers

# EDVT
## Electronic Design Validation Test

- All tests performed using offline diagnostics and again with NXOS

- On-board power supplies have voltages margined to +5% & -5%

- Temperature testing occurs while

- Soaking for 12 hours at $55^o$ C and -$5^o$ C

- Ramping between extremes at $1^o$ C per minute

- Power cycle testing occurs during 12-hour soak
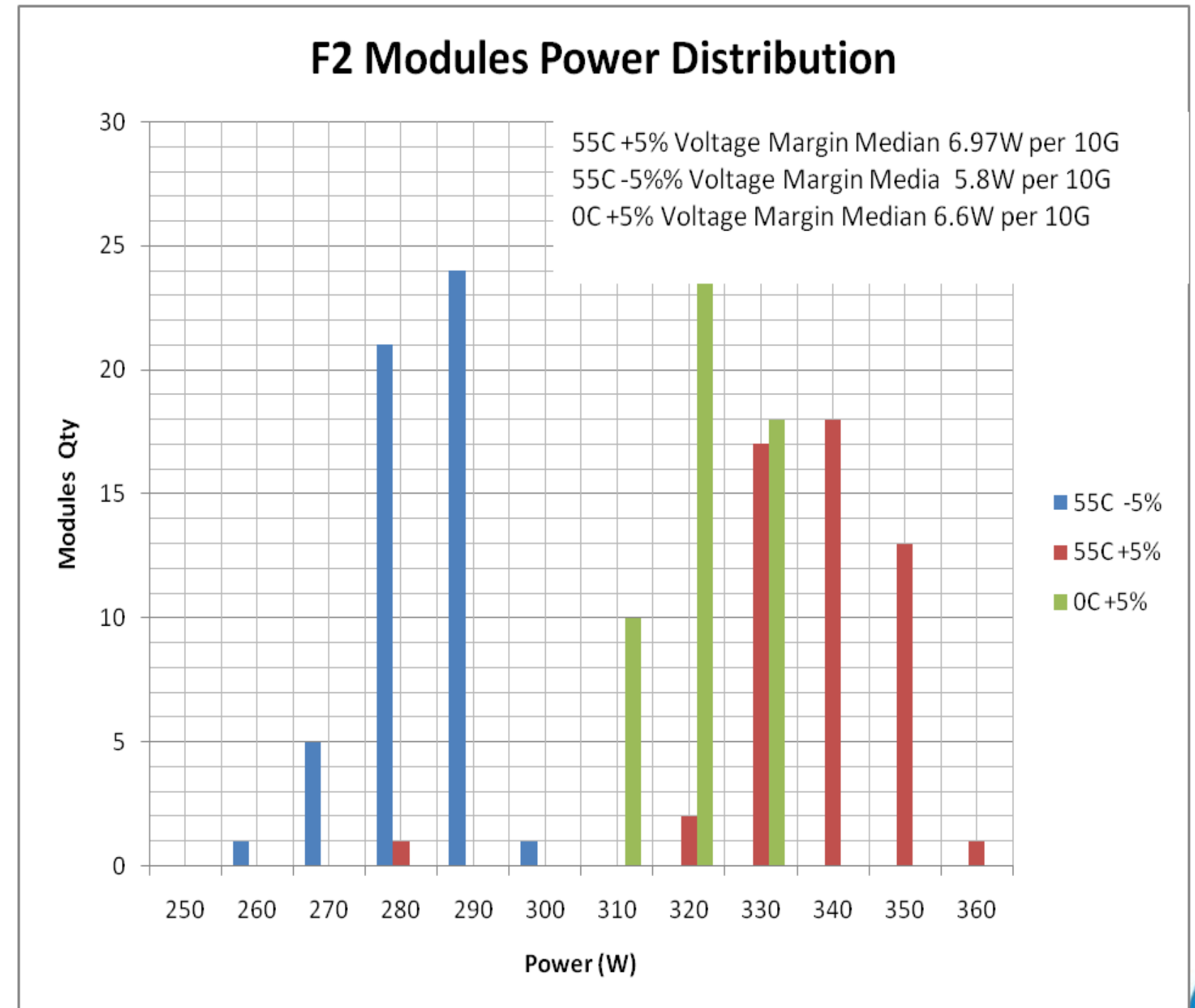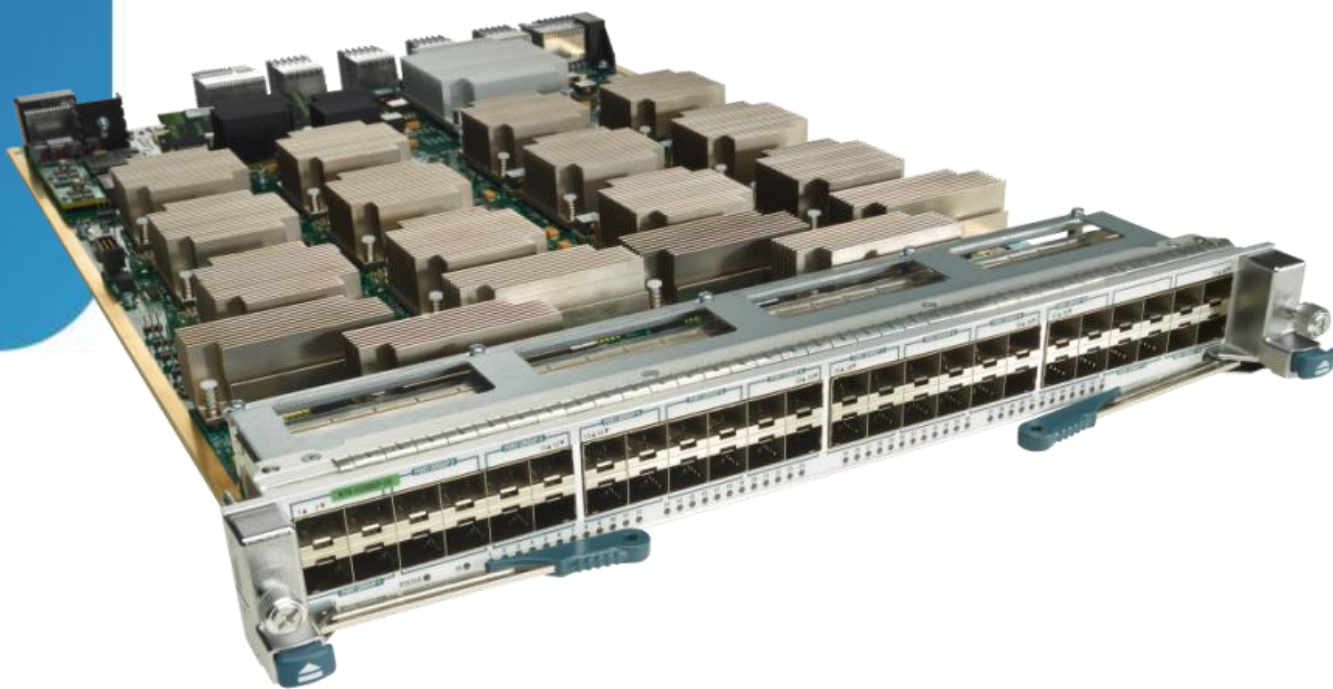


 Cisco Public

# RDT

## Reliability Demonstration Test

- The Reliability Demonstration Test (RDT) is Cisco's approach to verifying the stated reliability of a product prior to production release.

- The reliability to be demonstrated is the product's MTBF (Mean Time Between Failure).

- RDT replicates the end user operating environment and application through accelerated test time. It is expected that all hardware features are exercised in RDT.

- All new products including systems and boards are subject to RDT.

Cisco live!

# Power Consumption

- Skew Parts
- Data Sheet
  - Typical 340W
  - Maximum 450W



F2 Modules Power Distribution

55C +5% Voltage Margin Median 6.97W per 10G
55C -5%% Voltage Margin Media  5.8W per 10G
0C +5% Voltage Margin Median 6.6W per 10G

Cisco Public

# Generic Online Diagnostics

Generic Online Diagnostics provide a diagnostic framework for detecting hardware faults and verifying the health of hardware components throughout the chassis.

Diagnostics run during system Boot-Up, after OIR, On-Demand using the CLI, or as Health Checks in the background.
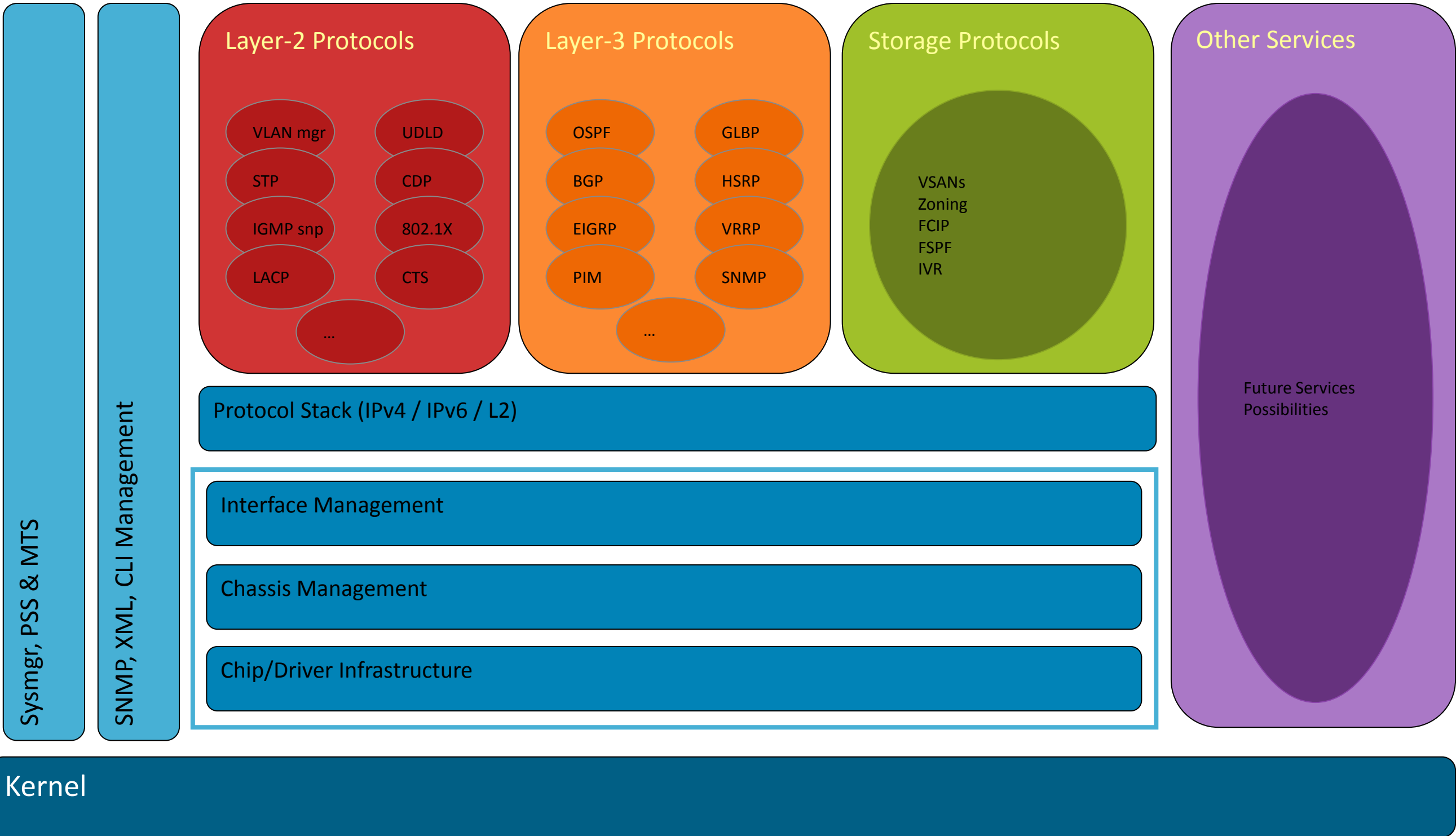
**Problem Areas:**

- Hardware Components (ASICs)

- Interfaces (Ethernet, SFP+, etc…)

- Connecters (loose connectors, bent pins, etc…)

- Memory Failure (Failure over time)

- Solder Joints

# Software Engineering

# NXOS Architecture

**Sysmgr, PSS & MTS**

**SNMP, XML, CLI Management**

## Layer-2 Protocols

- VLAN mgr
- UDLD
- STP
- CDP
- IGMP snp
- 802.1X
- LACP
- CTS
- ...

## Layer-3 Protocols

- OSPF
- GLBP
- BGP
- HSRP
- EIGRP
- VRRP
- PIM
- SNMP
- ...

## Storage Protocols

VSANs
Zoning
FCIP
FSPF
IVR

## Other Services

Future Services
Possibilities

**Protocol Stack (IPv4 / IPv6 / L2)**

**Interface Management**

**Chassis Management**

**Chip/Driver Infrastructure**

**Kernel**

Cisco *live!*

## Multi-threaded

Scalability with SMP and multi-core CPUs

Faster Route Re-convergence

Lower mean-time-to-recovery

## Real-Time

Real-Time preemptive scheduling

System operational when CPU is 100%

## Modularity

Most of the features are conditional

Can be enabled/disabled independently

Maximises efficiency

Minimises resources utilisation

## Separation Control Plane and Data Plane

No "software forwarding feature"

Fully distributed hardware forwarding

## Line Card Offloading

Offload to line card CPUs

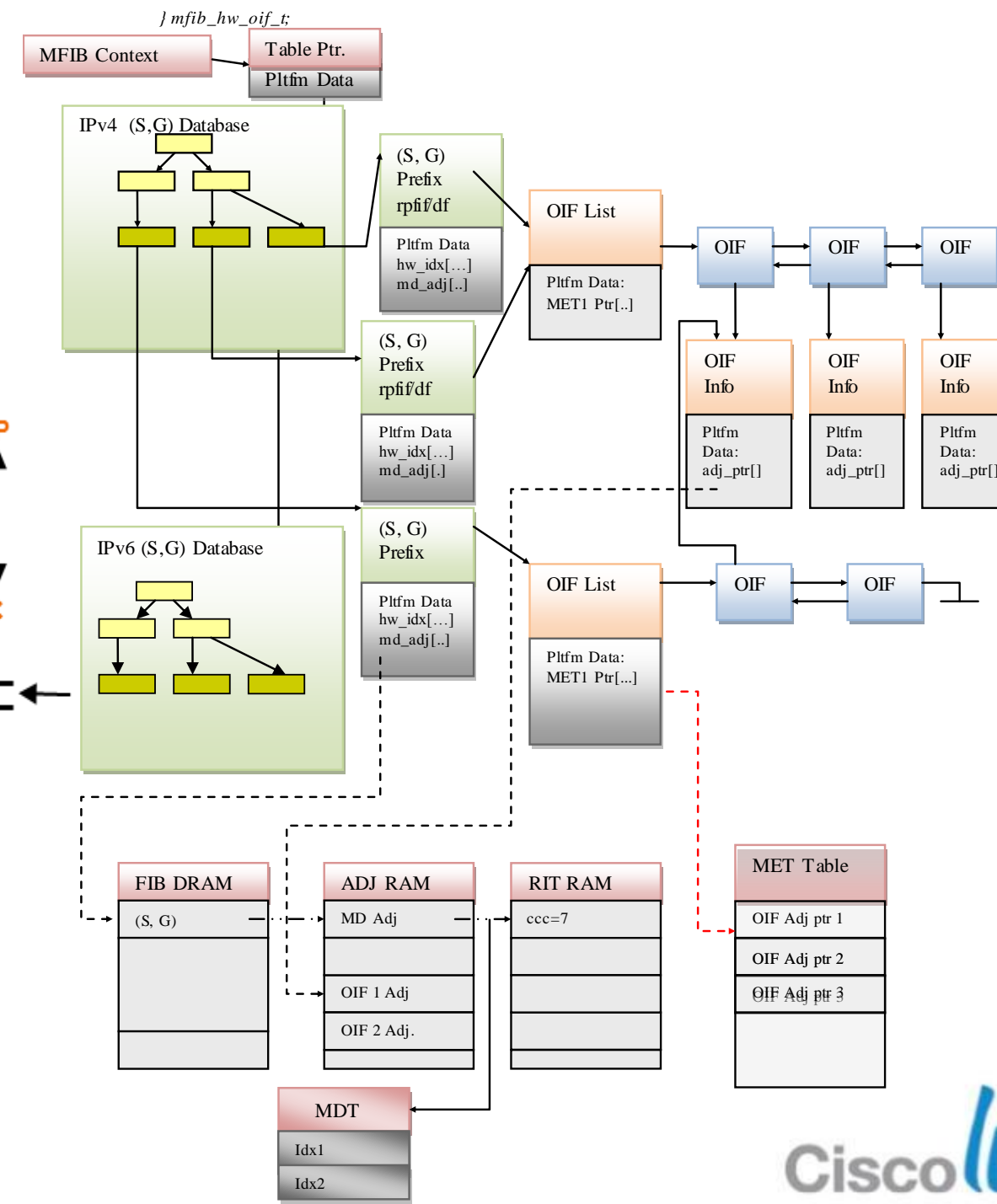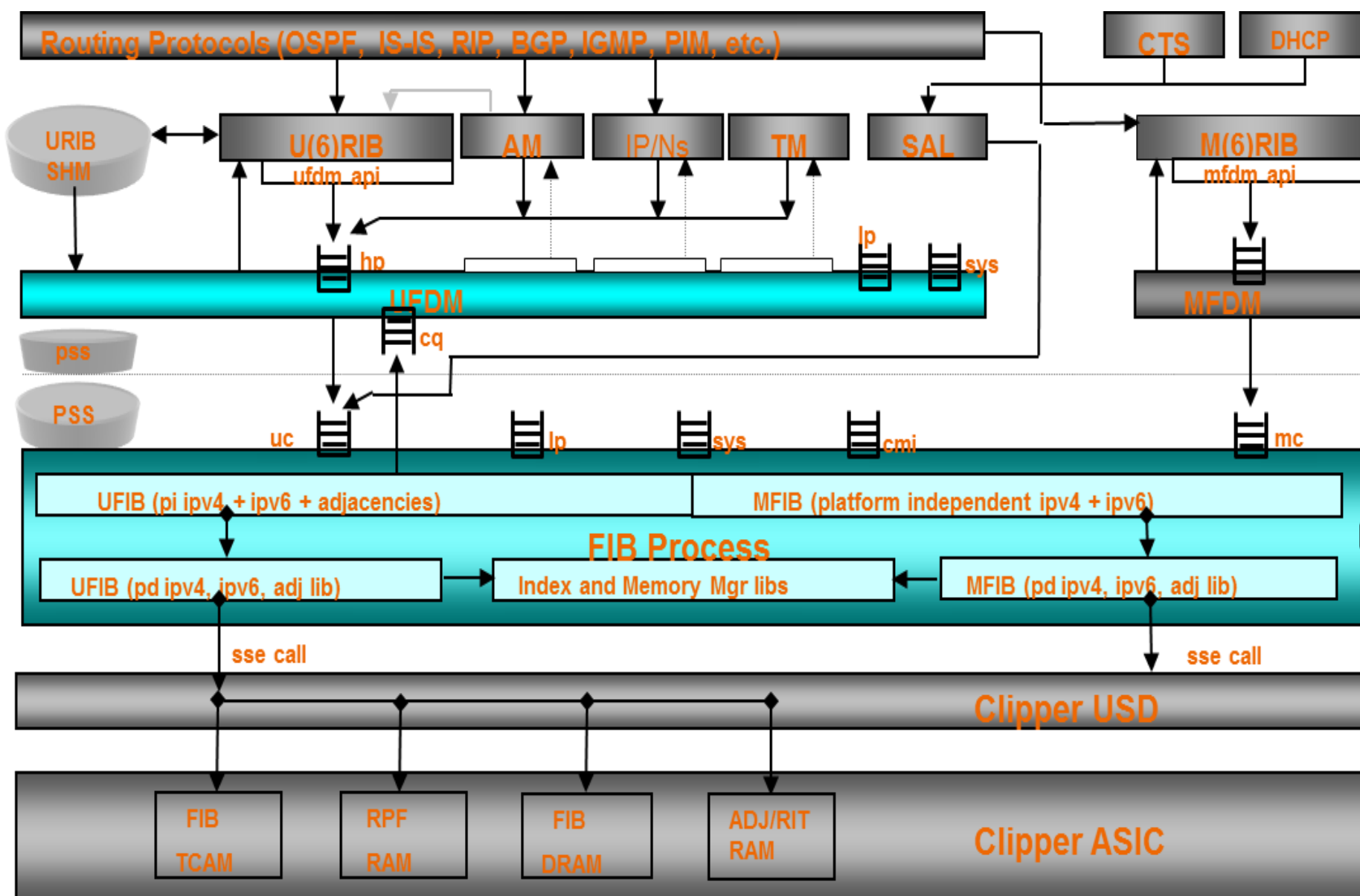Scales with # of line cards

Optimal hardware programming

# Software Engineering

# Design Review
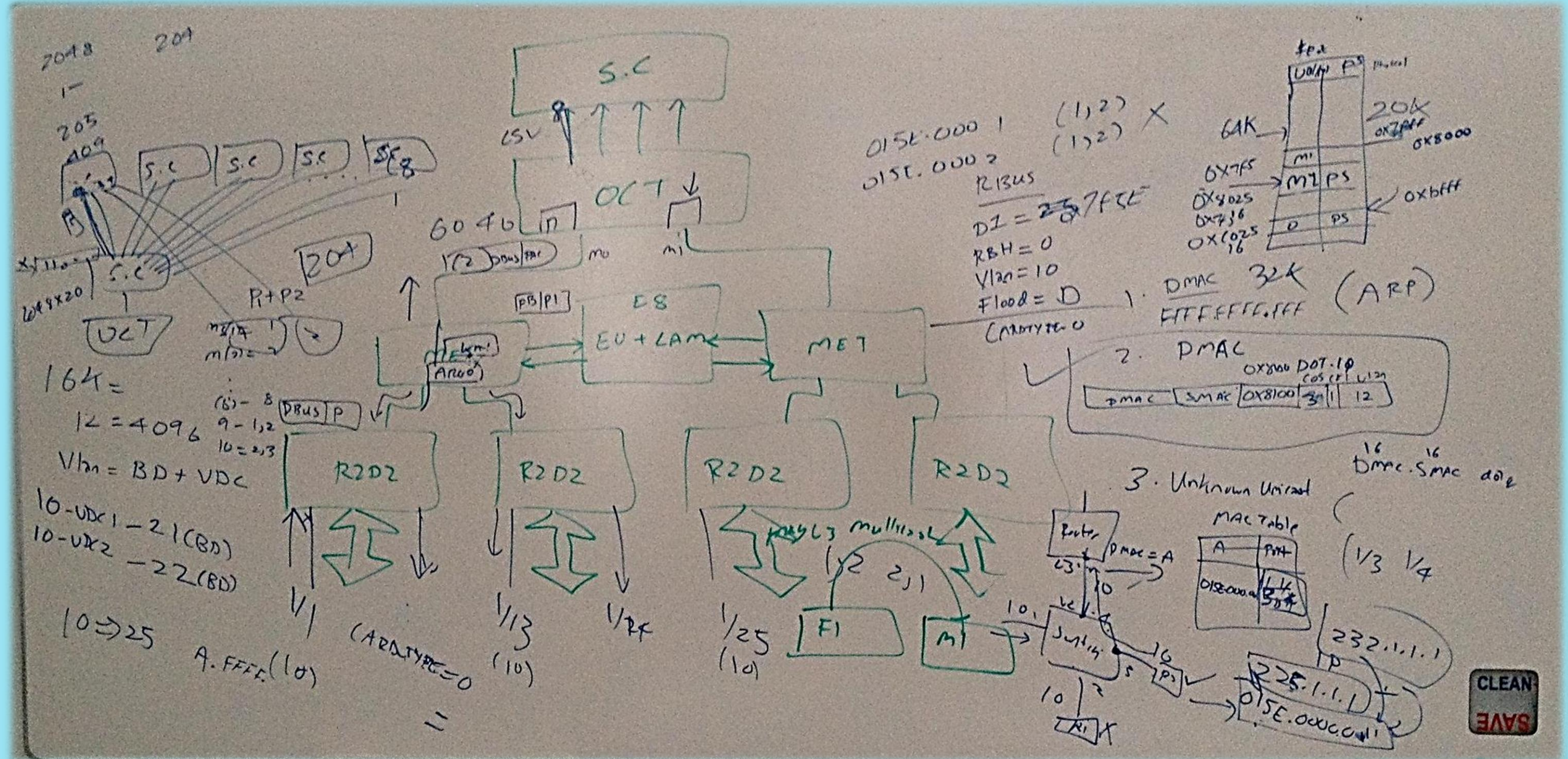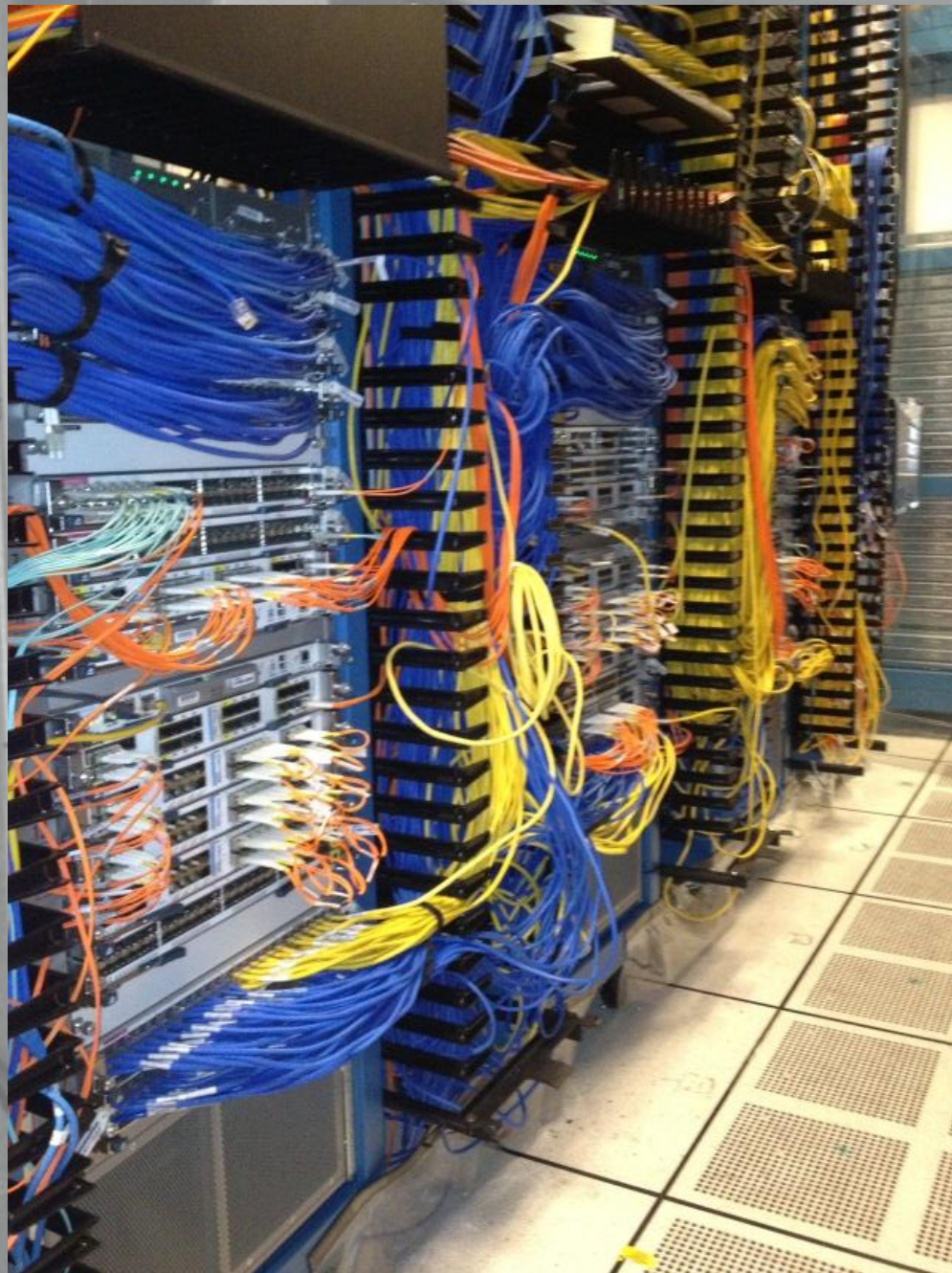
# Development Test

| Master Test Plans | Functional Test Plans | Automation | Regression | FCS |

- Testing of completed integrated feature
- Test for interactions with other features and functions
- Test for interoperability with Cisco and 3rd party devices
- Build scripts to automate testing so is repeatable on future releases

# First Customer Ship

**Product Requirements Document**

**ASIC**

| Requirements | PI... | ...lementation | Final Netlist | Power On |

**Hardware**

| HW Des... | P1 | P2 | A-0 |

| Mec... |

| Elect... |

| Manufactu... |

**Software**

| SW Functional Spec | ...st Plan | Unit Integration Plan |

**Software Test**

| Master Test Plans | ...nal Test Plans | Automation | Regression | FCS |

Cisco live!

# Q & A

# Complete Your Online Session Evaluation

## Give us your feedback and receive a Cisco Live 2013 Polo Shirt!

Complete your Overall Event Survey and 5 Session Evaluations.

- Directly from your mobile device on the Cisco Live Mobile App
- By visiting the Cisco Live Mobile Site www.ciscoliveaustralia.com/mobile
- Visit any Cisco Live Internet Station located throughout the venue

Polo Shirts can be collected in the World of Solutions on Friday 8 March 12:00pm-2:00pm

Don't forget to activate your Cisco Live 365 account for access to all session material, communities, and on-demand and live activities throughout the year.  Log into your Cisco Live portal and click the "Enter Cisco Live 365" button.

www.ciscoliveaustralia.com/portal/login.ww

Cisco Public