# FocalPoint II

## A 300nS, 240 Gb/s switch/router

**FULCRUM** microsystems

---

**Agenda**

**Datacenter Interconnect Requirements**

**FocalPoint I Status Update**

**FocalPoint II (Bali) Overview**

**FULCRUM** microsystems

# Anatomy of the Multi-Fabric Data Center

*Inefficient islands add complexity and cost; limit scale-out*

**Assumption:**
Proprietary or single vendor fabrics are required to achieve latency and bandwidth needs

Front-End Servers (Clients)

Ethernet

*Bridge*

Comms Network

Cluster Network

Back-End Servers (Application Servers)

Compute Cluster

*Bridge*

Storage Network

**Assumption:**
Fibre Channel is required for lossless storage fabrics

Fibre Channel

Infiniband

Storage

**FULCRUM**
microsystems

---

# 10GE: Unifying Datacenter Interconnect

*Datacenter Ethernet enables full cross-sectional bandwidth and a single management domain over all three networks*

- Low latency
- 10G bandwidth
- Large-scale topologies
- QOS and flow control
- Rich ecosystem

Front-End Servers (Clients)

Comms Network

Cluster Network

Back-End Servers (Application Servers)

Compute Cluster

Storage Network

Legacy Storage

10G Ethernet

Clustered Storage

Fibre Channel

**FULCRUM**
microsystems

# Step 1: Solve Latency and Port Density

## *The world's most powerful 10G Ethernet switch chip*



- **Highest port density** (24 10GE ports)
- **Highest bandwidth** (240Gbps)
- **Lowest latency** (200ns)
- **Most scalable** (fat trees, 1,000s of ports)
- **Most integrated** (single chip)

### FocalPoint Evaluation Platforms
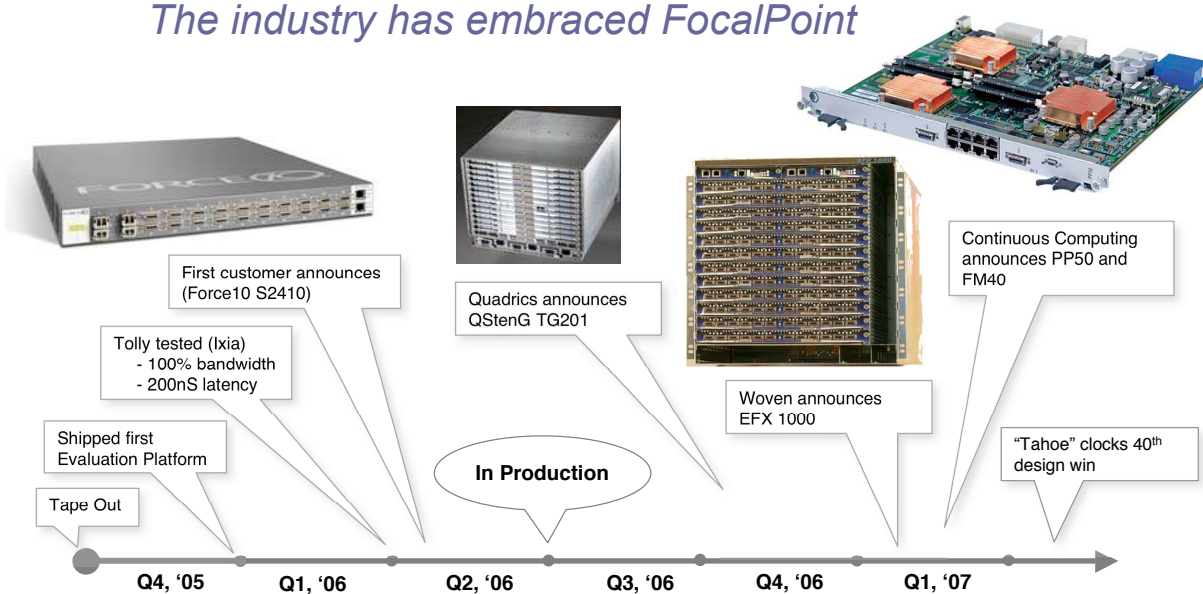(The world's most integrated 10G Ethernet systems)

| **Reno** | **San Marino** | **Heavenly** | **Vegas** |
|---|---|---|---|
| 24-port 10GE design | Highest 10GBase-T density | IBM BladeCenter-H fabric | Non-blocking 1GE platform |

**FULCRUM** microsystems

---

# FocalPoint Status Report

## *The industry has embraced FocalPoint*

First customer announces
(Force10 S2410)

Tolly tested (Ixia)
- 100% bandwidth
- 200nS latency

Shipped first
Evaluation Platform

Tape Out

Quadrics announces
QStenG TG201

**In Production**

Woven announces
EFX 1000

Continuous Computing
announces PP50 and
FM40

"Tahoe" clocks 40[th]
design win

| Q4, '05 | Q1, '06 | Q2, '06 | Q3, '06 | Q4, '06 | Q1, '07 |

**FULCRUM** microsystems
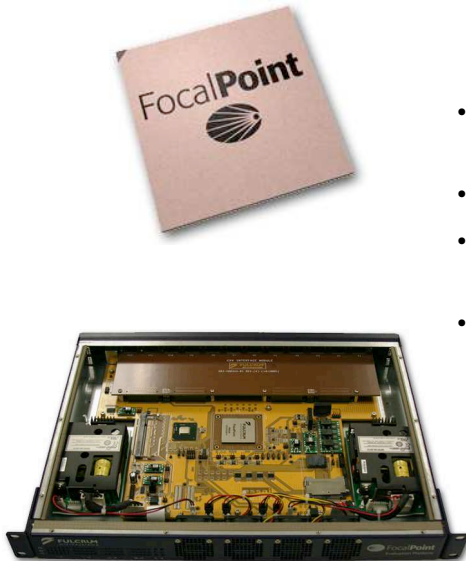
# Step 2: Routing & Network Performance

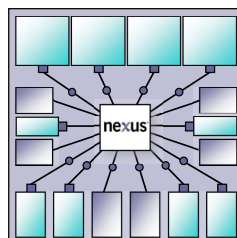## *FocalPoint II (Bali) project goals*

- **Maintains Gen I Performance**
    - **24 10GE ports, 200nS, 360 MPPS**
    - **But increase to 2MB memory**
- **IPv4 & IPv6 unicast & multicast routing**
    - **16k IP addresses**
- **L2-L4+ ACLs with deep inspection**
- **Chip cascades**
    - Virtual switch of Clos, rings & meshes
- **Clos Improvements**
    - As close to full bandwidth as possible

    **Converged Enhanced Ethernet (CEE)**
    - **Enable lossless Ethernet fabrics**

**FULCRUM** microsystems

---

# Architecture Enabling Circuits

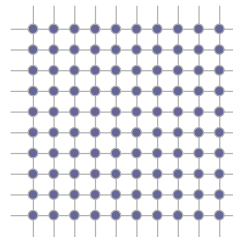## *Two key IP blocks demonstrate the virtues of the technology*

### Nexus[*]
### *(Terabit Crossbar)*

- Gigahertz performance
- Terabit capacity
- Nanosecond latency
- No power penalty

\* Licensed to **PMC** for SoC interconnect
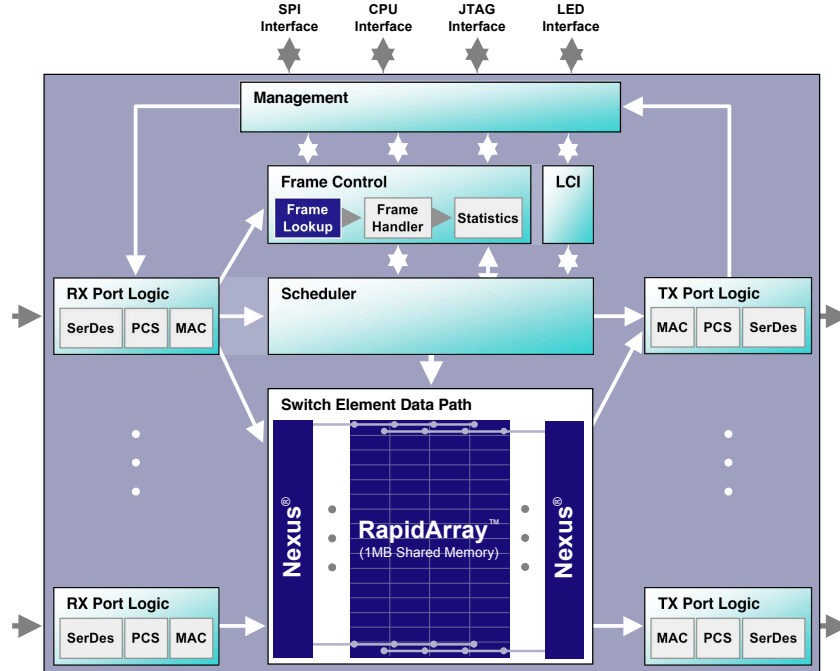**PMC-SIERRA**

### RapidArray
### *(Asynchronous SRAM)*

- 720 MHz SRAM
- 1200 MHz interconnect
- 518 Gbps throughput
- Scalable for any use

## Key Benefits:

- **Easily integrates independent clock domains**
- **Provides 4x overspeed**
- **Reduces overall chip area**

- **2x the speed of vendor cores (same size, density, yield)**
- **Reduces power consumption (based on activity)**
- **CAM circuit is close relative**

**FULCRUM** microsystems
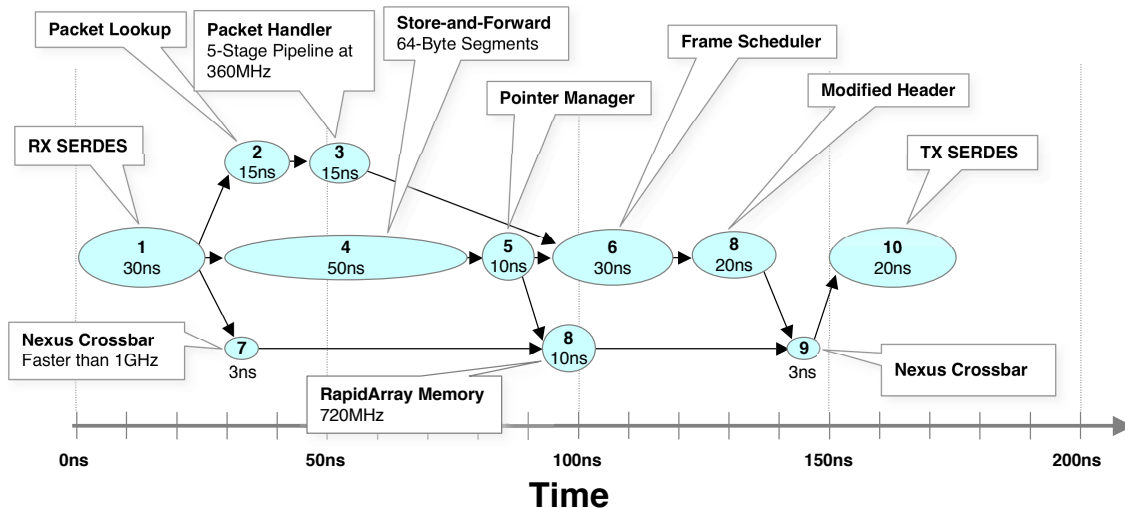
# FocalPoint I & II Architecture

*Modular architecture, centralized control*

SPI Interface   CPU Interface   JTAG Interface   LED Interface

**Management**

**Frame Control**
- Frame Lookup
- Frame Handler
- Statistics

**LCI**

**RX Port Logic**
- SerDes
- PCS
- MAC

**Scheduler**

**TX Port Logic**
- MAC
- PCS
- SerDes

**Switch Element Data Path**

Nexus®

**RapidArray™**
(1MB Shared Memory)

Nexus®

**RX Port Logic**
- SerDes
- PCS
- MAC

**TX Port Logic**
- MAC
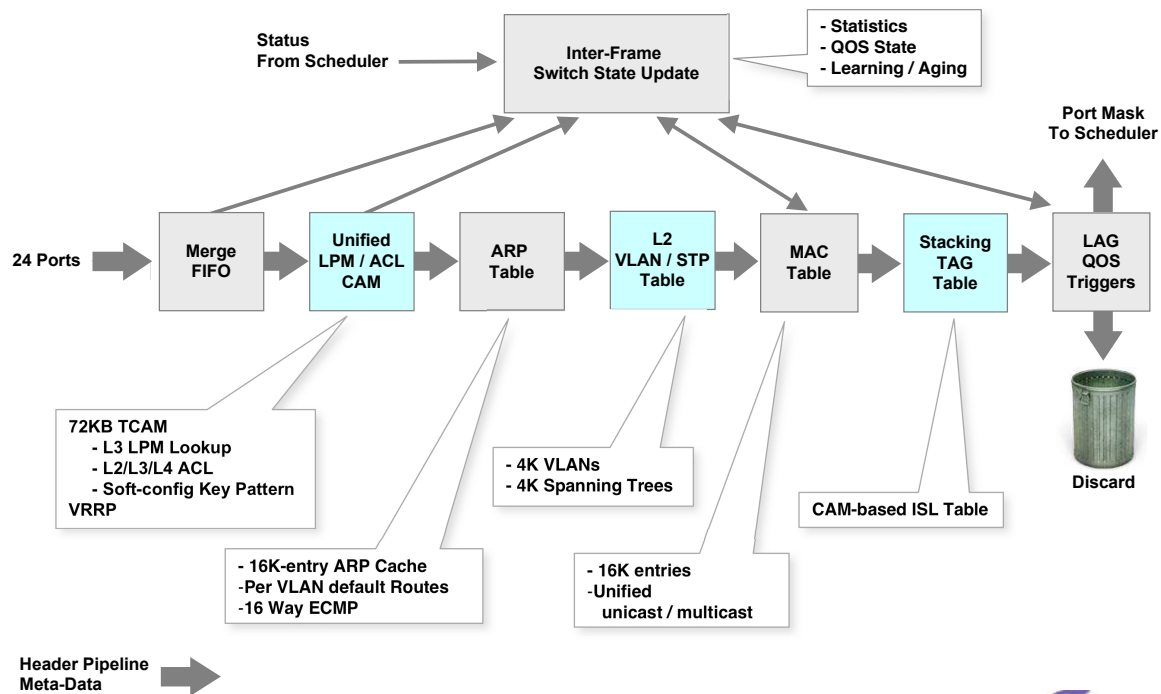- PCS
- SerDes

FULCRUM
microsystems

---

# FocalPoint I & II Latency Detail

*Requirement: Full L3 latency less than 300nS with 360 MPPS*

*Budget: 35 clock-cycle pipeline (2.6 nS / cycle) for all L3*

Packet Lookup

Packet Handler
5-Stage Pipeline at 360MHz

Store-and-Forward
64-Byte Segments

Frame Scheduler

Pointer Manager

Modified Header

RX SERDES

TX SERDES

**2** 15ns   **3** 15ns

**1** 30ns   **4** 50ns   **5** 10ns   **6** 30ns   **8** 20ns   **10** 20ns

Nexus Crossbar
Faster than 1GHz

**7** 3ns   **8** 10ns   **9** 3ns   Nexus Crossbar

RapidArray Memory
720MHz

0ns   50ns   100ns   150ns   200ns

**Time**

FULCRUM
microsystems

# FP II L2-L4 Packet Processing Pipeline



Status From Scheduler

Inter-Frame Switch State Update

- Statistics
- QOS State
- Learning / Aging

Port Mask To Scheduler

24 Ports

Merge FIFO

Unified LPM / ACL CAM

ARP Table

L2 VLAN / STP Table

MAC Table

Stacking TAG Table

LAG QOS Triggers

72KB TCAM
- L3 LPM Lookup
- L2/L3/L4 ACL
- Soft-config Key Pattern
VRRP

- 4K VLANs
- 4K Spanning Trees

CAM-based ISL Table

Discard

- 16K-entry ARP Cache
-Per VLAN default Routes
-16 Way ECMP

- 16K entries
-Unified
  unicast / multicast

Header Pipeline Meta-Data

FULCRUM
microsystems

---

# L2-L4 Packet Parsing and Manipulation



Frame Data Storage

## Header Fields Pipeline

Dest MAC
Source  MAC
VLAN 1
VLAN 2
VLAN Priority
Etype
F64
C.N. headers
SIP v4 or V6
DIP v4 or V6
TTL
TOS/DSCP
Protocol
L4 Source
L4 Dest
TCP Flags
Deep Inspection A..D
CRC & Checksum

New Dest MAC
New Souce MAC
New VLAN 1
New VLAN Priority
New Etype
New F64
New C.N. headers
New TOS/DSCP
New TTL
2nd CRC Check

RX Port Logic

SerDes | PCS | MAC

TX Port Logic

MAC | PCS | SerDes

FULCRUM
microsystems
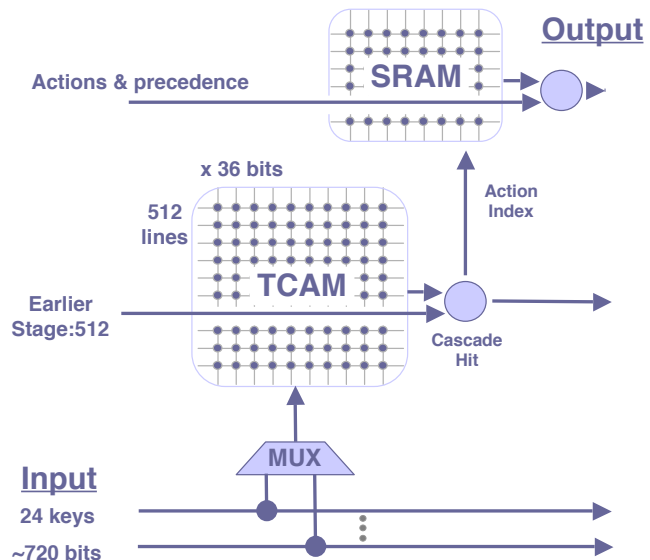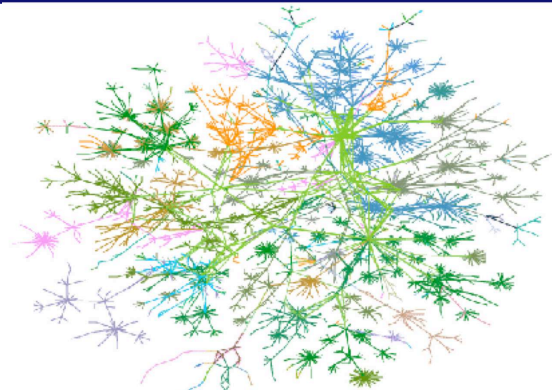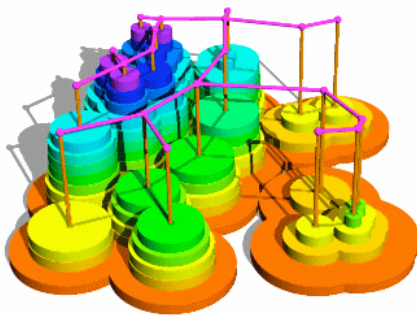
# Filter & Forwarding (FFU) Organization

- **Input header keys**
  - Contains input fields
  - Any key available to any and all banks
- **TCAM Organization**
  - 32 banks, 72 KB, 16k min rules
  - Combine up to 32 banks
  - SRAM encodes 1 or multiple actions
  - Precedence in action combine allows multiple levels of non-orthogonal rules
- **Performance**
  - 1 gate delay per stage
  - 6 stages per flop
  - 6 clocks overall: 15nS of Latency
  - 360 MHz

**Output**

**SRAM**

Actions & precedence

x 36 bits

512 lines

**TCAM**

Action Index

Earlier Stage:512

Cascade Hit

**MUX**

**Input**

24 keys

~720 bits

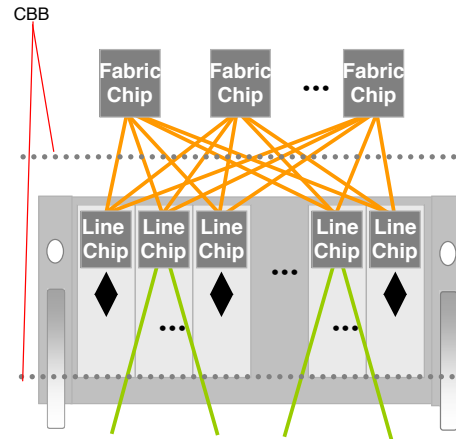FULCRUM
microsystems

---

# Topology Enhancements and Ethernet

- **An Ethernet switch uses a single spanning tree**
  - Networks should scale in a non-blocking fashion
  - Spanning tree hashing and resilience

- **A Chip is not an Ethernet switch**
  - We may want to use one or many chips in a "switch"
  - We need to link data plane port state between chips
  - Multi-chip and multi-box link aggregation and management
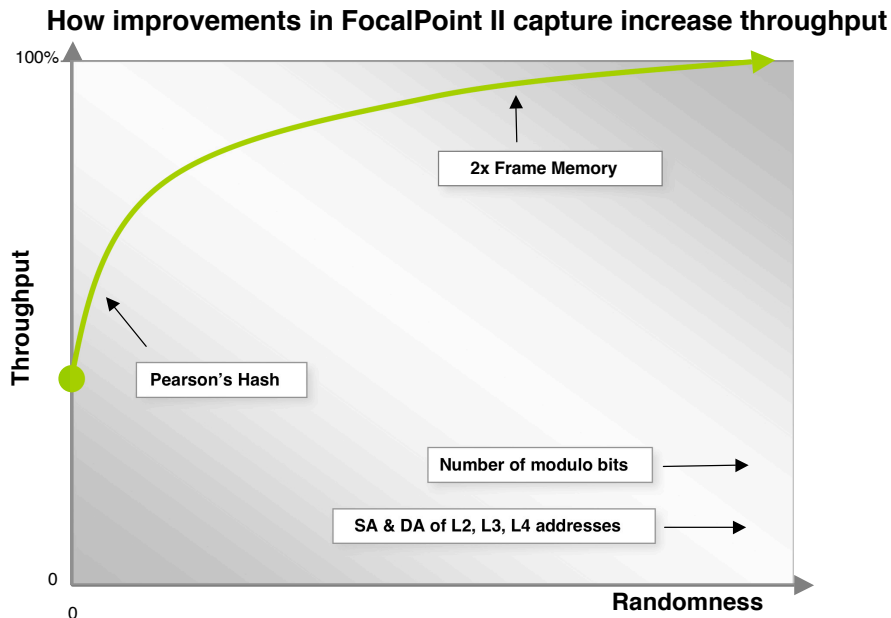
FULCRUM
microsystems

# Clos Architecture Theoretical Performance

- **Significant Economic advantages are driving the move to Clos**
    - 10x reduction in per port price
    - 10x reduction in latency
    - 2-8x increase in port density / BW

- **All modern Clos architectures are really statistical multi-path, multi-hop networks**
    - Examples are Infiniband (Mellanox), fiberchannel (Brocade), Ethernet (Fulcrum)

- **Clos architectures achieve ideal performance if**
    - The path selection (often hashing) is sufficiently stochastic
    - There is enough over-speed in the system to compensate for any non-ideal path selection
    - There is enough memory per switch to compensate for collisions so that flow control is infrequent

CBB

Fabric Chip    Fabric Chip    ···    Fabric Chip

Line Chip  Line Chip  Line Chip    Line Chip  Line Chip

FULCRUM
microsystems

---

# FocalPoint II achieves new performance levels

*All flow based systems face the same challenges*

**How improvements in FocalPoint II capture increase throughput**



100%

Throughput

2x Frame Memory

Pearson's Hash

Number of modulo bits

SA & DA of L2, L3, L4 addresses

0

0                                           Randomness

FULCRUM
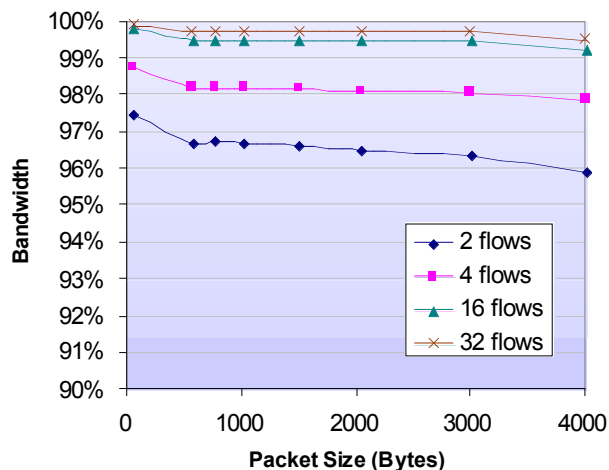microsystems

# Simulated Clos Performance

- **Standard Full Mesh**
  - Each port sends to every other port in the system every cycle
  - Every packet is randomly assigned to a flow on its port
  - 2-32 flows per port simulated
  - 576-9216 system flows (288P)
  - 288P system is made from 36 24P switch chips in a Clos configuration
- **Amount of over-speed from line card switching (24P line cards)**
  - 3% in 288P Clos
  - 23% in 48P Clos
- **Conclusion**
  - 16 flows per port is nearly perfect in a 288P system (four flow is very good)
  - 1 flow per port is nearly perfect in a 48P system given the 23% overspeed

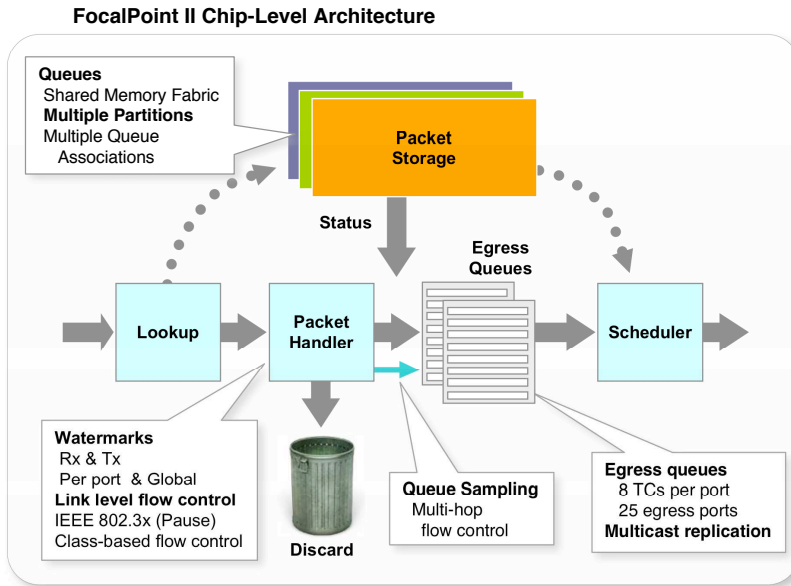### 288P Full Mesh Performance

FULCRUM
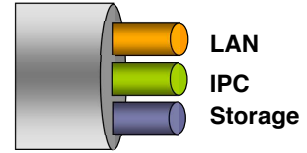microsystems

---

# Congestion Control

- **There are three drivers for congestion control enhancements in the enterprise datacenter**
  - Storage & IPC traffic is highly loss and jitter sensitive
  - Multiple traffic types, like storage and LAN, require different best in class congestion control practices
  - Cost oriented fully integrated full bandwidth switch chips are required to use memory very efficiently

- **As a result there is a race to produce lossless, non-HOL blocking, low latency fabrics with optimal bandwidth**
  - Congestion control is being standardized by the IEEE in 802.1au and potential future working groups

FULCRUM
microsystems

## Traffic separation enables virtual switching

**FocalPoint II Chip-Level Architecture**

**Queues**
Shared Memory Fabric
**Multiple Partitions**
Multiple Queue
Associations

**Packet Storage**

**Status**

**Egress Queues**

**Lookup**

**Packet Handler**

**Scheduler**

**Watermarks**
Rx & Tx
Per port & Global
**Link level flow control**
IEEE 802.3x (Pause)
Class-based flow control

**Discard**

**Queue Sampling**
Multi-hop
flow control

**Egress queues**
8 TCs per port
25 egress ports
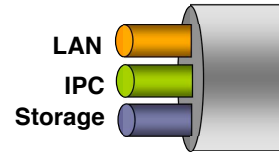**Multicast replication**

**Egress Features**

**LAN**
**IPC**
**Storage**

•Multi-level scheduling
   Bandwidth groups with priority
•Traffic distribution
   Deficit weighted round robin

**Ingress Features**

**LAN**
**IPC**
**Storage**

•Flow Control
   Link & Per class pause
•Static & Dynamic rate-limiting
   Pause pacing & policing

**FULCRUM** microsystems

---

# Bali Chip Plot

## Fabricated in TSMC 0.13um
## 250 Million Transistors



**TCAM**

**RapidArray Memory 2 MB**

**MAC Table**
- 16K addresses

**Ethernet Port Logic**
- Phy (SerDes)
- PCS
- MAC

**Frame Control**
- Frame handler
- Lookup
- Statistics

**Nexus Crossbar**
- Terabit capacity
- 3ns latency

**Scheduler**
- Highly optimized
- High event rate

**FULCRUM** microsystems

# Thank You!

**Uri Cummings**
*Founder, CTO*
uri@**fulcrum**micro.com

**FULCRUM**
microsystems

**818.871.8100**
www.**fulcrum**micro.com

**26630 Agoura Road**
**Calabasas, CA 91302**

**FULCRUM**
microsystems