# IP/MPLS High Availability Technologies

**Pranav Dharwadkar**

**CRBU Product Mgr (pranavd@cisco.com)**

**Originator/Patent Holder of Prefix Independent Convergence**
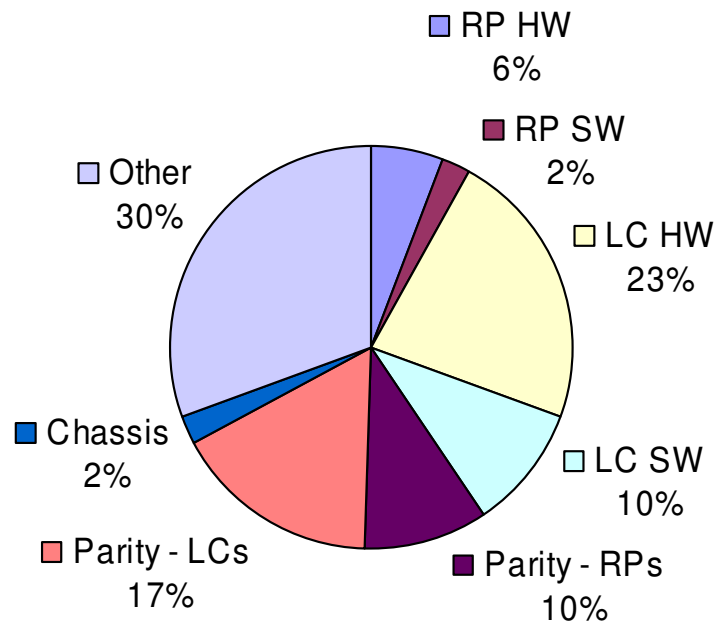
**December 4th 2008**

# Agenda

- High Availability – What is it ?

- What is Cisco's view on High Availability

- IP/MPLS High Availability Technologies

- High Availability Infrastructure Tools

# Quantification

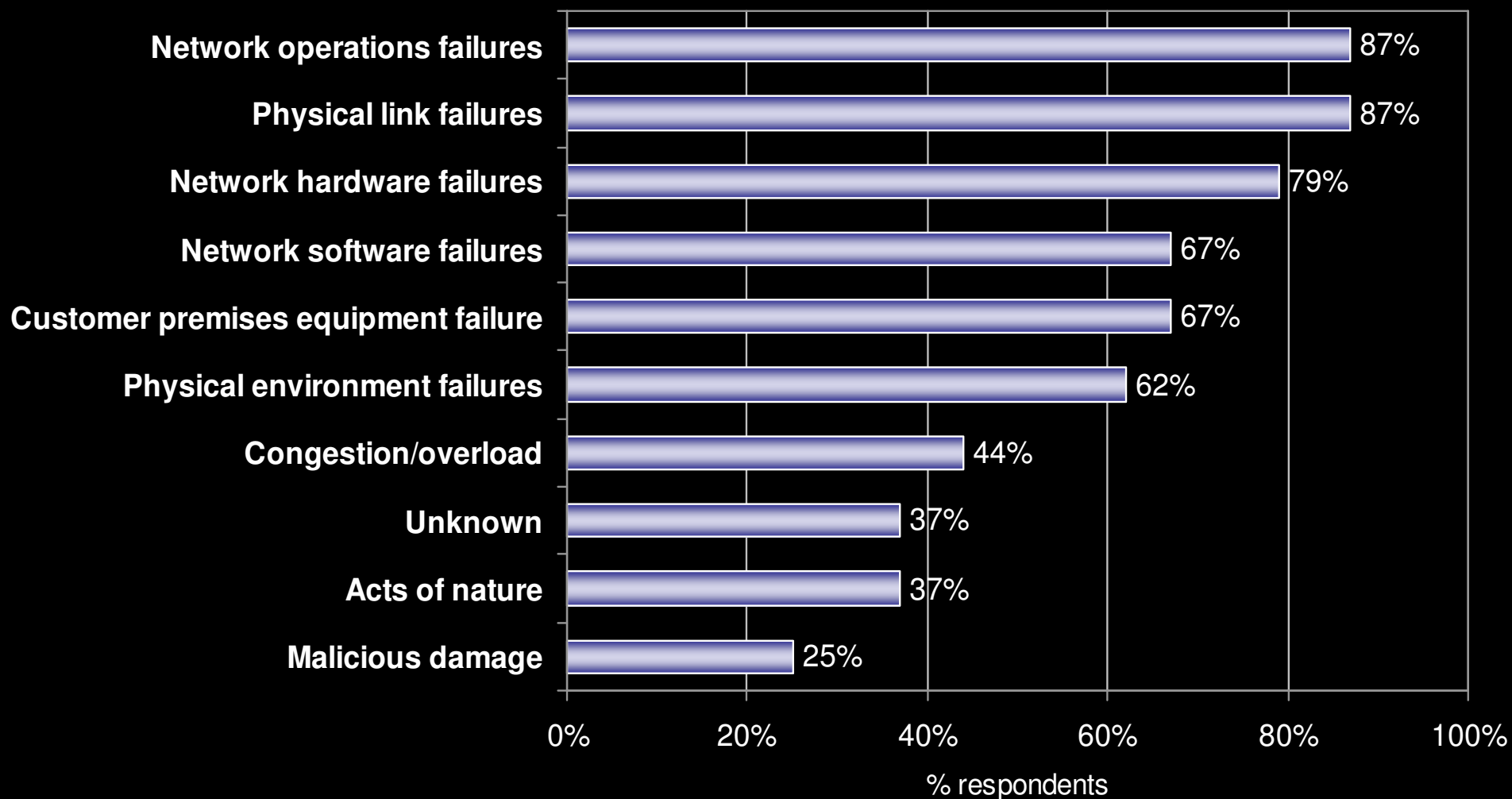| Percent Availability | N-Nines | Downtime Time Minutes/Year |
|---|---|---|
| 99% | 2-Nines | 5,000 Min/Yr |
| 99.9% | 3-Nines | 500 Min/Yr |
| 99.99% | 4-Nines | 50 Min/Yr |
| 99.999% | 5-Nines | 5 Min/Yr |
| 99.9999% | 6-Nines | .5 Min/Yr |

# Root Cause Analysis: Number of faults

**From 101 failures in a large IP Network (April – June 2002)**



**Other** category includes: Line problems, operator error, config errors, cables, etc.

# Causes of Unscheduled Downtime



SOURCE: Sage Research, *IP Service Provider Downtime Study: Analysis of Downtime Causes, Costs and Containment Strategies, August 17, 2001, Prepared for Cisco SPLOB*
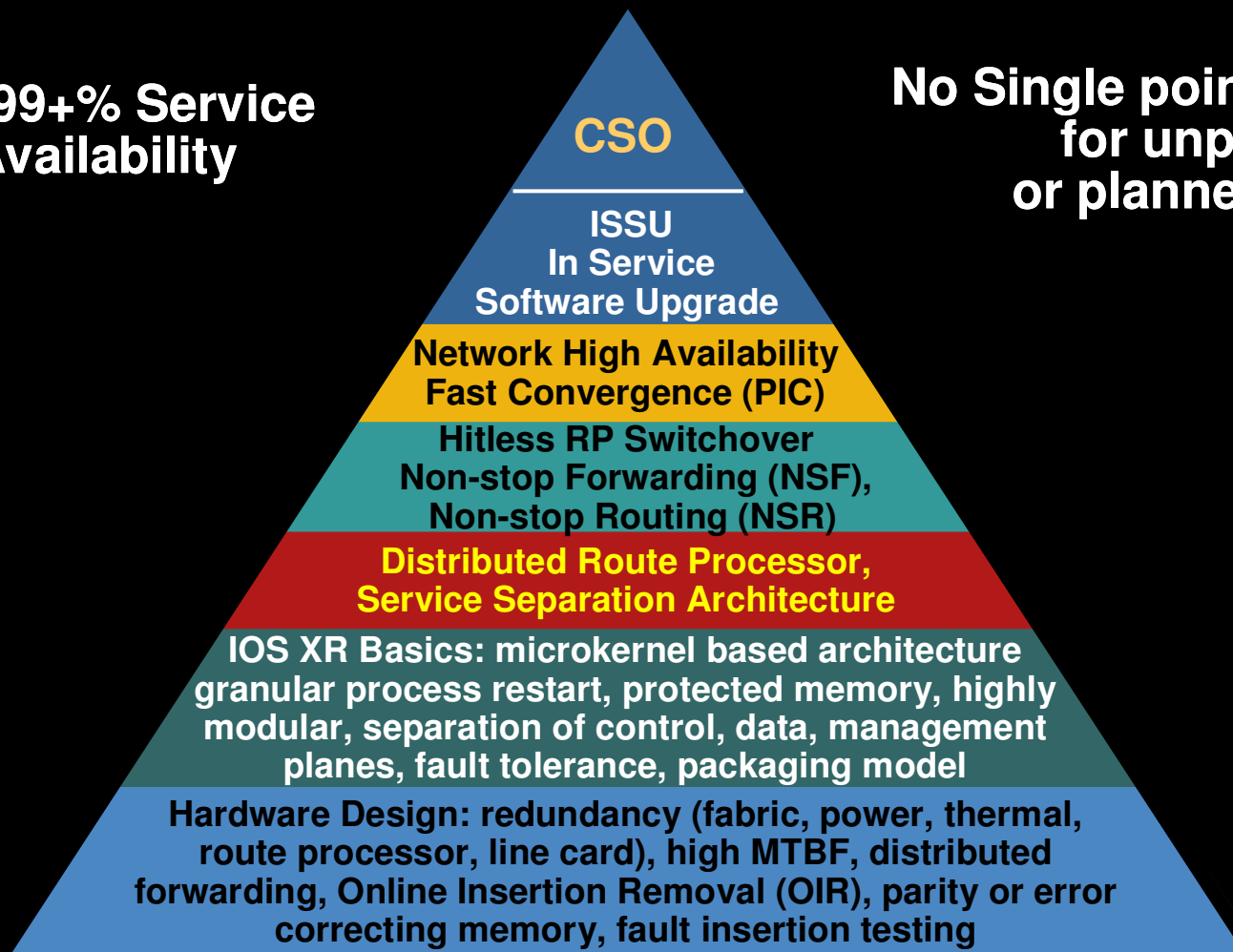
# HA – 4 Factors

- Network Factor
  - Link Failures – Fiber cuts
  - Link Failures – Forwarding logic failure
  - Node Failures – HW
  - Link Failures - SW
  - Congestion
  - Security attacks

- System Factor
  - HW Failure
  - SW Failure

- Operations Factor
  - Network Operations Failure
  - Out of resource conditions
  - Sparing & Support
  - Training

- Environment Factor
  - Physical Environment Failure
  - Malicious Damage

# Recipe for High Availability - Increase MTBF, Reduce MTTR

## Continuous Systems Operation

99.999+% Service Availability

No Single points of failure for unplanned or planned events

**CSO**

**ISSU In Service Software Upgrade**

**Network High Availability Fast Convergence (PIC)**

**Hitless RP Switchover Non-stop Forwarding (NSF), Non-stop Routing (NSR)**

**Distributed Route Processor, Service Separation Architecture**

**IOS XR Basics: microkernel based architecture granular process restart, protected memory, highly modular, separation of control, data, management planes, fault tolerance, packaging model**

**Hardware Design: redundancy (fabric, power, thermal, route processor, line card), high MTBF, distributed forwarding, Online Insertion Removal (OIR), parity or error correcting memory, fault insertion testing**
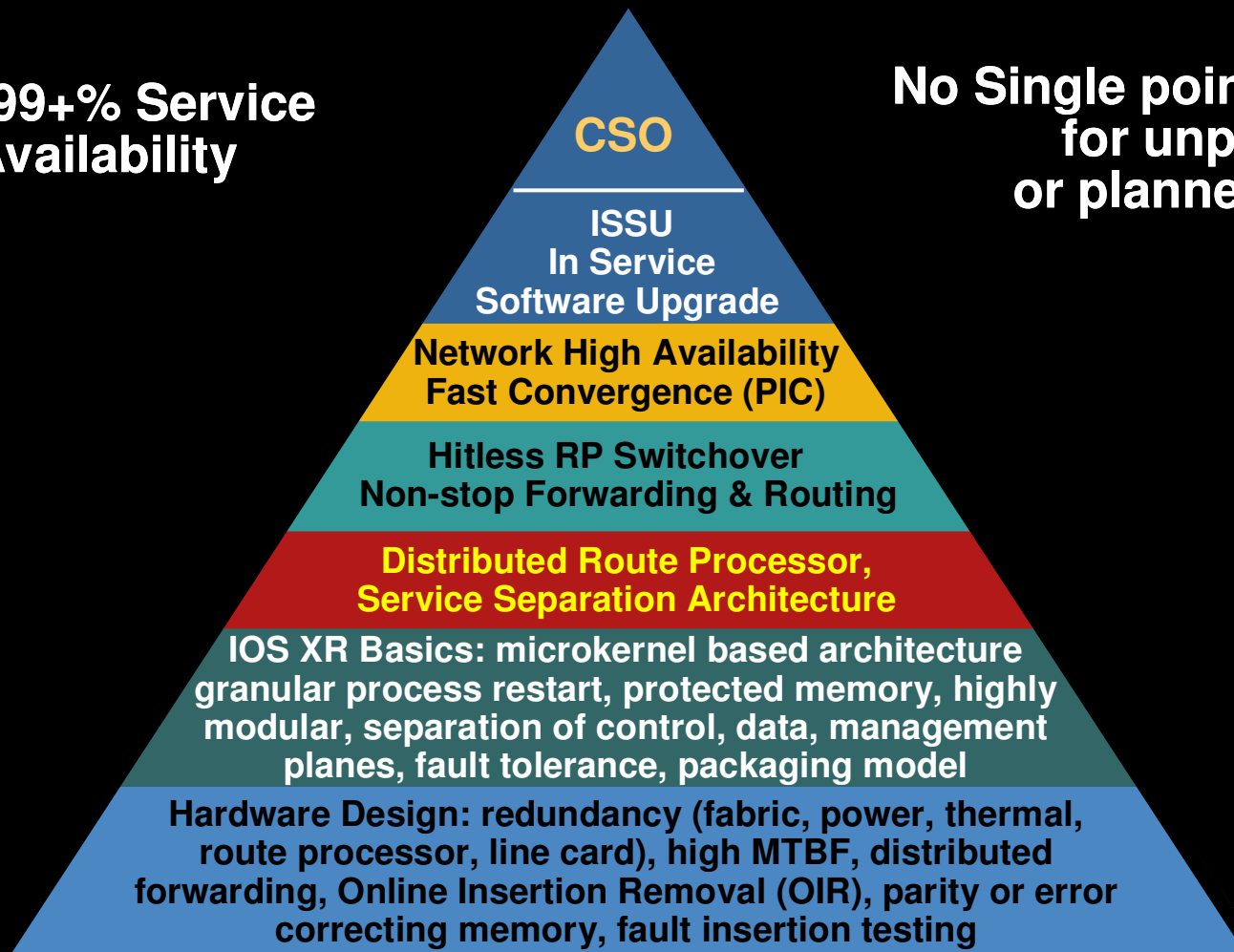
# Recipe for High Availability - Increase MTBF, Reduce MTTR

## Continuous Systems Operation

**99.999+% Service Availability**

**No Single points of failure for unplanned or planned events**

**CSO**

**ISSU In Service Software Upgrade**

**Network High Availability Fast Convergence (PIC)**

**Hitless RP Switchover Non-stop Forwarding & Routing**

**Distributed Route Processor, Service Separation Architecture**

**IOS XR Basics: microkernel based architecture granular process restart, protected memory, highly modular, separation of control, data, management planes, fault tolerance, packaging model**

**Hardware Design: redundancy (fabric, power, thermal, route processor, line card), high MTBF, distributed forwarding, Online Insertion Removal (OIR), parity or error correcting memory, fault insertion testing**

# HA: Hardware Design

## Redundancy Built into every piece of hardware

**No single Point of Failure**

**Failure of fabric, power, thermal, route processor results in immediate switchover to redundant hardware**

## Hardware Memory Error Detection & Correction

**Error Correction Seamless; only traffic hitting faulty memory affected during error correction**
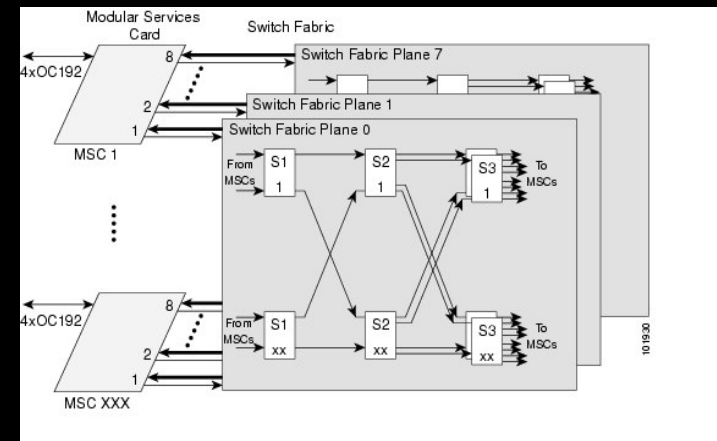
# Hitless Expansion:
## Multi-chassis
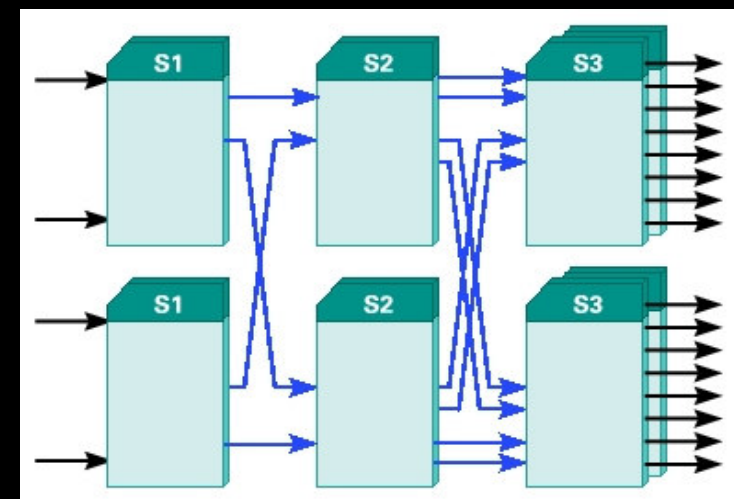
## Goal:

**Allow incremental capacity expansion with no service loss**

## Implementation:

Same switch fabric architecture for standalone and multi-chassis uses load-shared redundancy

Methodology: Take one plane out of service at a time in a standalone config, upgrade that plane to multi-chassis configuration and bring it back online



*CRS 8-plane switch fabric*



*Single Plane of 8-plane switch fabric*

# Hardware High Availability

## Full Hardware Redundancy

No Single Point of Failure

Online Insertion & Removal (OIR) for all cards
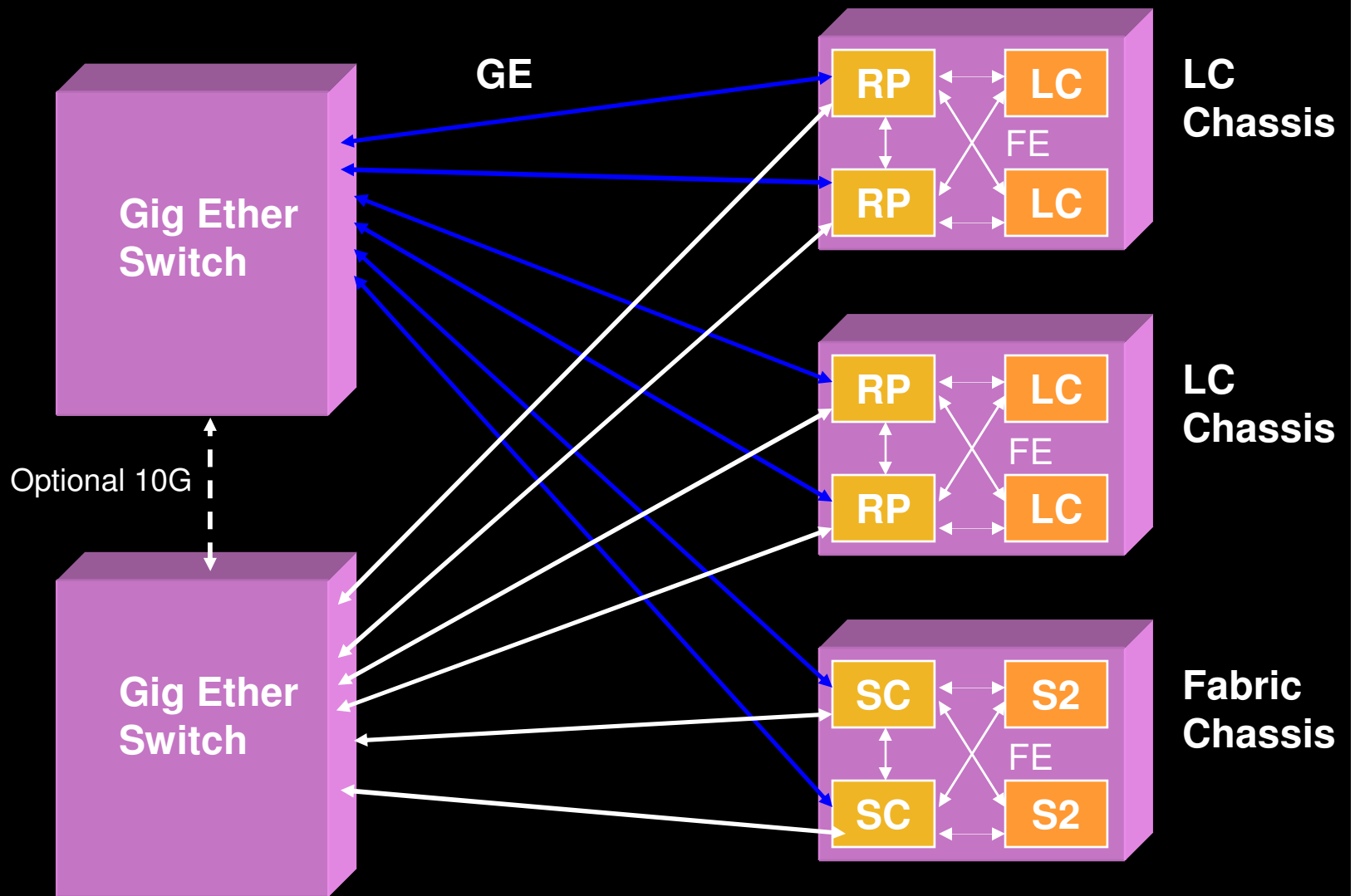
HW assist arbitration

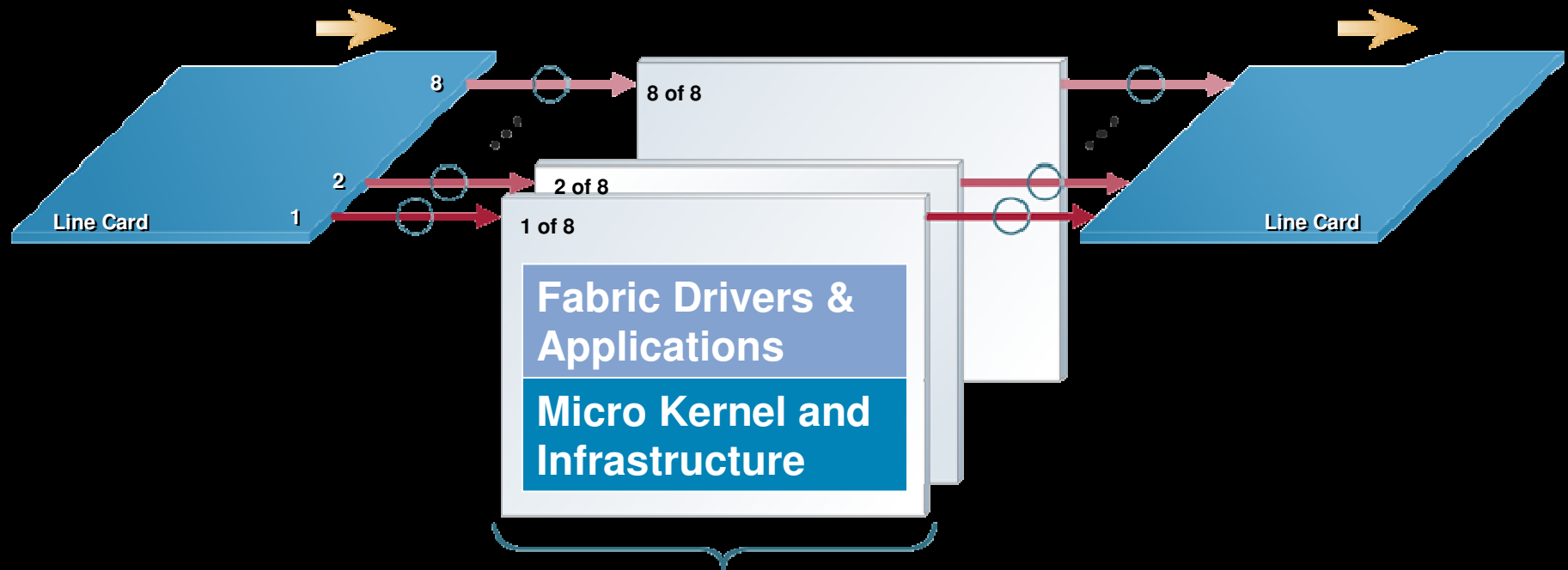No outage on Fabric Upgrade or Failure

ECC protected memory and R-S FEC for optical links

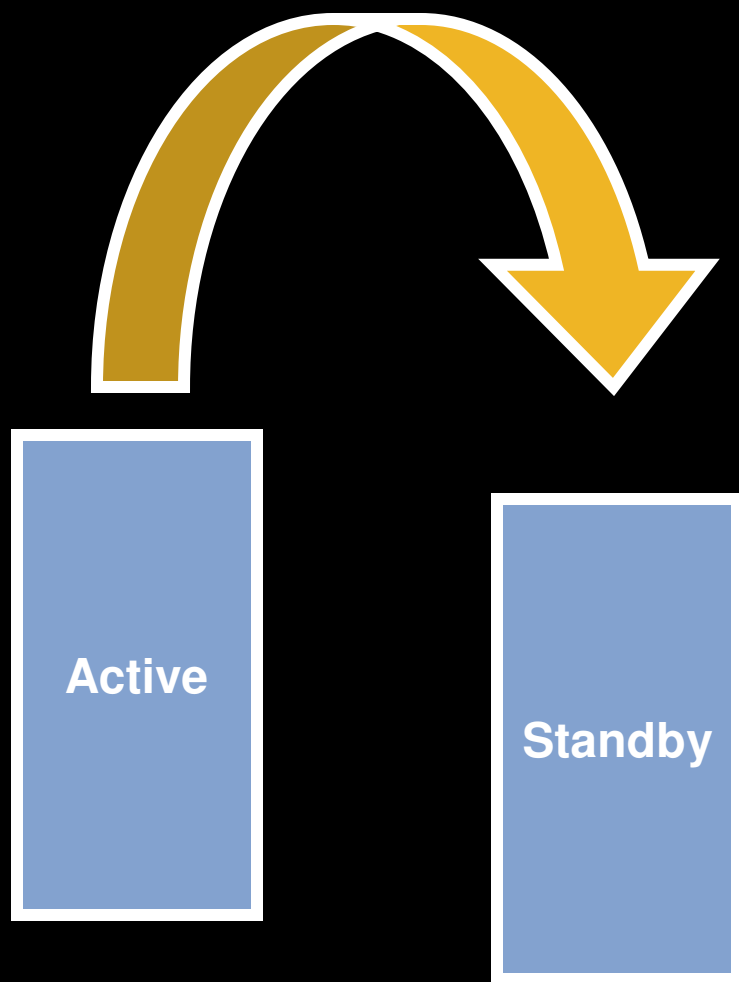Redundant Out-of-Band System Control Network

# System Control Network
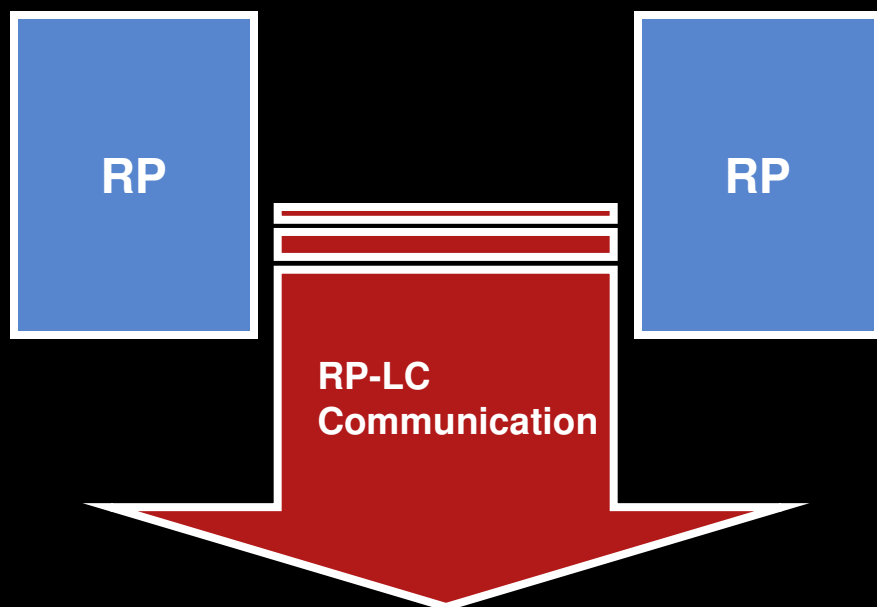
# No single points of Failure – Switch Fabric



8

8 of 8

2

2 of 8

Line Card    1

1 of 8

**Fabric Drivers & Applications**

**Micro Kernel and Infrastructure**

Line Card

- Supports 1:N redundancy
- Allows fabric upgrade one card at a time

# No Single points of Failure - RP

**Active**

**Standby**

- HW errors detected on active RP card

- Control plane lockup on active card

- Routing Protocol Crashes on Active RP

- Costly to recover on the same node

# Control Plane – Data Plane Separation

**RP**

**RP**

**RP-LC Communication**

**LC**

**LC**

- CRS and 12K has dedicated packet forwarding hardware (SPP / ISE)

- Packet forwarding on LC's can function autonomously during control plane outages

- Packet forwarding un-affected by:
  - ISIS, OSPF, BGP, MPLS mcast process restart
  - Infrastructure process restarts
  - RP failover

# Hitless Disk Replacement:
## Disk Mirroring

### Goal:

**Handle disk failures without causing a RP switchover**
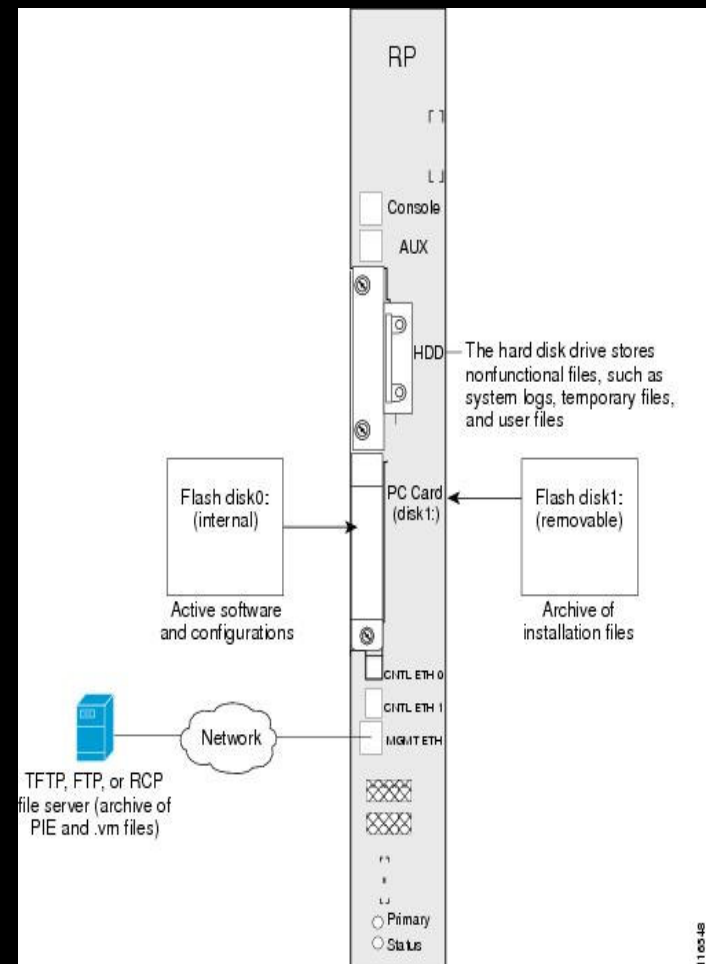
**Allow replacement of faulty disks**

### Implementation:

2 flashdisks on CRS RP extend persistent storage.

- Disks configured as redundant pairs

- All important data (image and config) are replicated on the 2 disks

Any disk outage (software or hardware issues) can be handled locally using the built-in redundancy



RP

Console

AUX

HDD — The hard disk drive stores nonfunctional files, such as system logs, temporary files, and user files

Flash disk0: (internal)

PC Card (disk1:) ← Flash disk1: (removable)

Active software and configurations

Archive of installation files

CNTL ETH 0
CNTL ETH 1
MGMT ETH

TFTP, FTP, or RCP file server (archive of PIE and .vm files)

Network

Primary
Status

*CRS Route Processor with dual flashdisks*

# HA – System Factor

Cisco Confidential

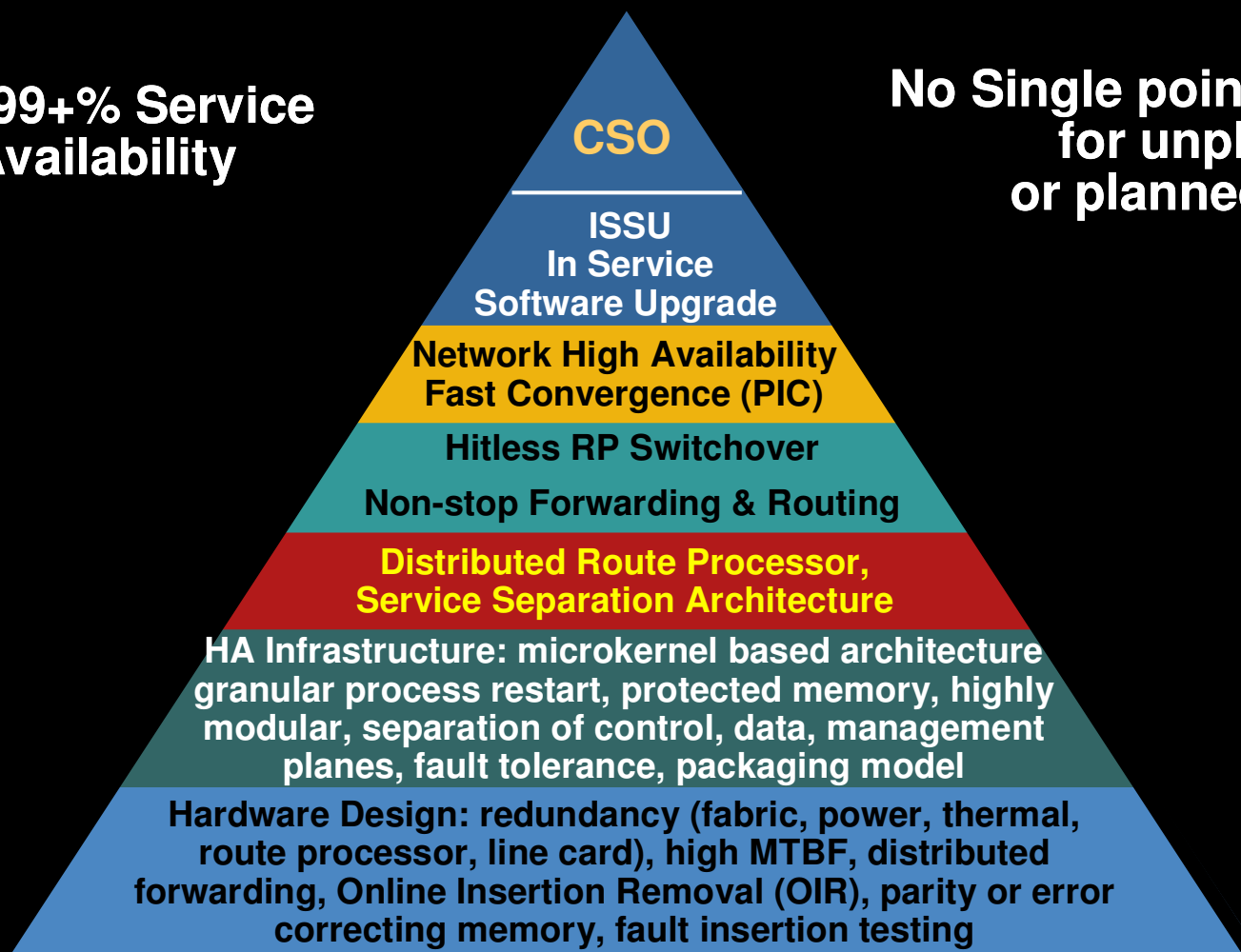# Avoid Single Points of Failure…

- Support Redundant RP and Minimize or Eliminate
    - Switchover or upgrade causes an outage
    - Time spent in troubleshooting the problem
    - The probability of no on-site spares or the spares don't work

- Support Link protection using links on different line cards – So LC is not a single point of failure
    - e.g., ECMP
    - e.g., VRFs hosted on multiple Line Cards, Multiple Interfaces to CE
    - e.g., Link Bundles

- Support Redundancy in Switching Fabric – So Switching Fabric is not a single point of failure
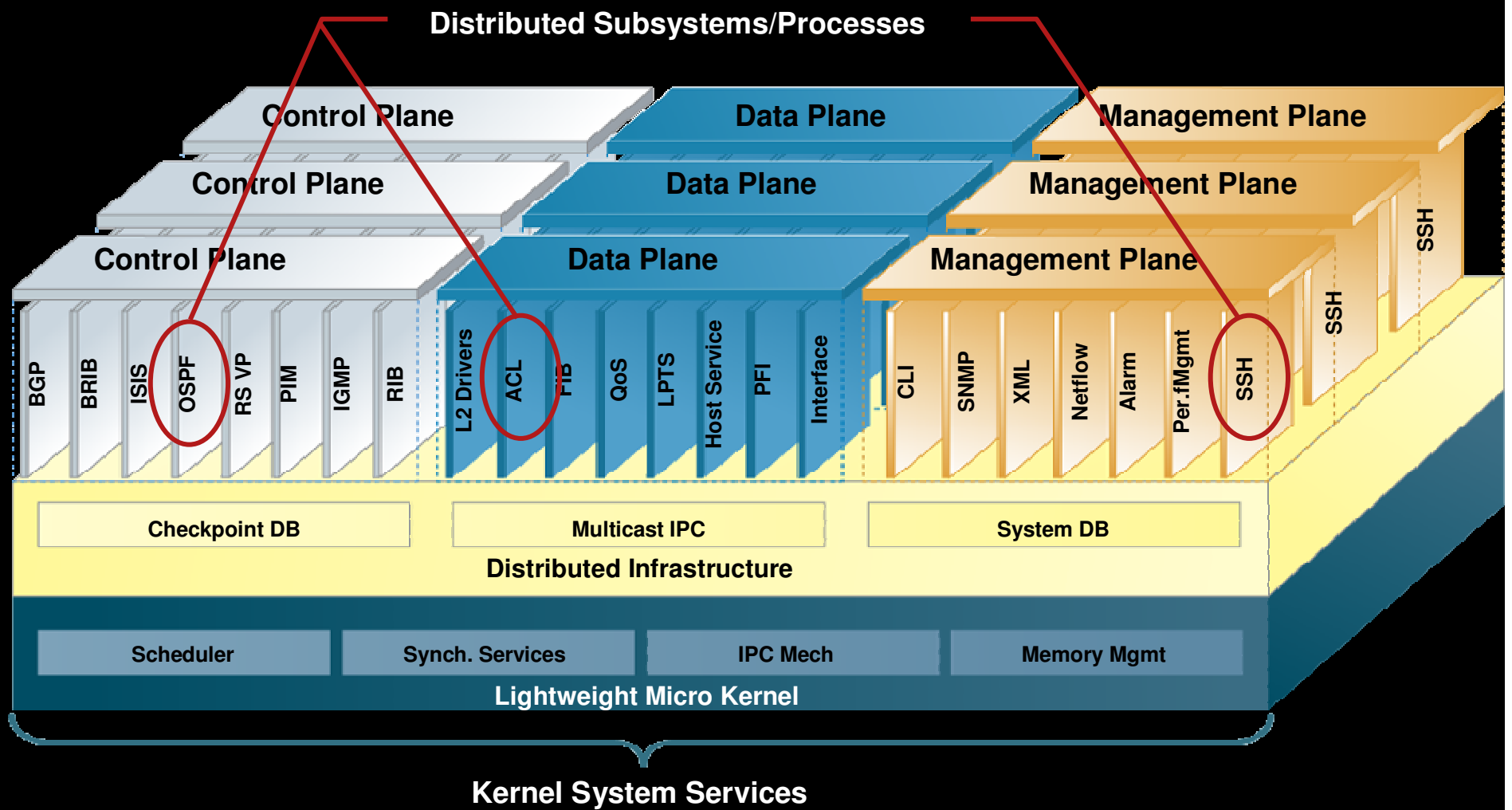
# Recipe for High Availability - Increase MTBF, Reduce MTTR

## Continuous Systems Operation

99.999+% Service Availability

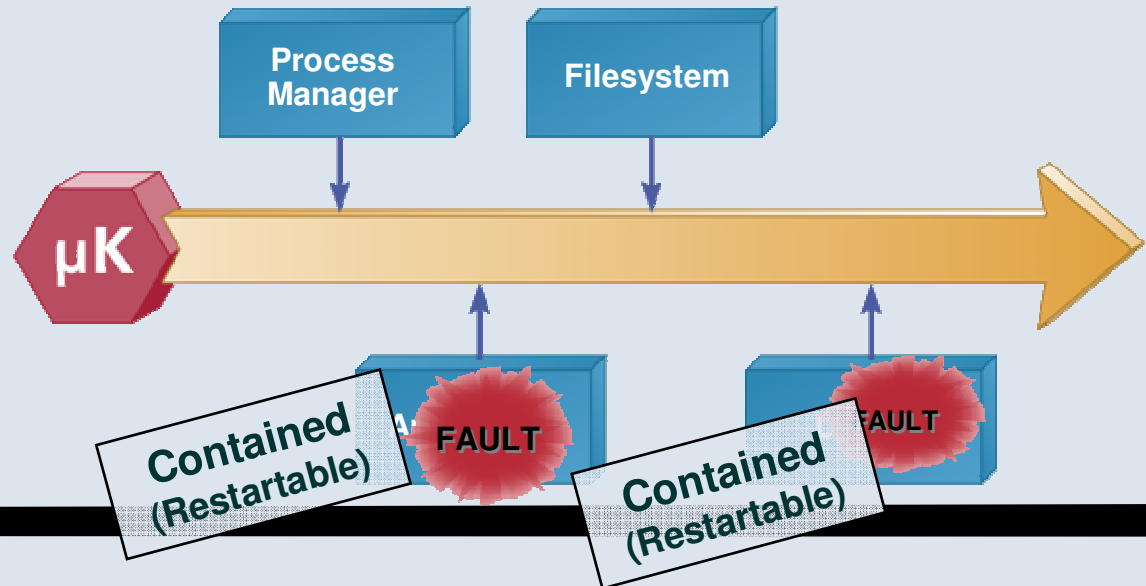No Single points of failure for unplanned or planned events

**CSO**

ISSU
In Service
Software Upgrade

Network High Availability
Fast Convergence (PIC)

Hitless RP Switchover

Non-stop Forwarding & Routing

Distributed Route Processor,
Service Separation Architecture

HA Infrastructure: microkernel based architecture granular process restart, protected memory, highly modular, separation of control, data, management planes, fault tolerance, packaging model

Hardware Design: redundancy (fabric, power, thermal, route processor, line card), high MTBF, distributed forwarding, Online Insertion Removal (OIR), parity or error correcting memory, fault insertion testing

# Distributed OS for Next Generation Networks



**Distributed Subsystems/Processes**

Control Plane
Control Plane
Control Plane

Data Plane
Data Plane
Data Plane

Management Plane
Management Plane
Management Plane

BGP | BRIB | ISIS | OSPF | RS VP | PIM | IGMP | RIB

L2 Drivers | ACL | FIB | QoS | LPTS | Host Service | PFI | Interface

CLI | SNMP | XML | Netflow | Alarm | Per.fMgmt | SSH | SSH | SSH

| Checkpoint DB | Multicast IPC | System DB |

**Distributed Infrastructure**

| Scheduler | Synch. Services | IPC Mech | Memory Mgmt |

**Lightweight Micro Kernel**

**Kernel System Services**

# IOS XR - Micro Kernel Architecture

**TRUE Microkernel
(Mach, QNX)**

MMU with full protection

Applications, drivers, and protocols are protected

Process Manager

Filesystem

µK

**Contained
(Restartable)**

FAULT

FAULT

**Contained
(Restartable)**

**Monolithic Kernel
(BSD/Linux, NT)**

MMU with partial protection

Applications are protected

**Contained
(Restartable)**

FAULT

Application

System Wide Corruption

Network Driver

FAULT

# High Availability Infrastructure

Contained

| BGP | IS-IS | RIB | | QoS | FIB | IP Stack | | CLI | XML | Alarm | File System |

**OS** → Distributed Middleware →

| OSPF | PIM | IGMP | | ACL | PFI | L2 Drivers | | Netflow | SNMP | SSH | Inter Process Communication |

Contained

- Granular process restart allows for fast recovery from failures
- Leverage hardware redundancy like link and RP redundancy
- Graceful restart mechanisms in routing protocol

Cisco Confidential and Proprietary

# Modular OS



- **Upgrade specific packages/Composites**
  - Across Entire system
    - Useful once a feature is qualified and you want to roll it without lot of commands
  - Targeted Install to specific cards
    - Useful while a feature is being qualified–reduces churn in the system to card boundary
- **Point Fix for software faults**

# Software Upgrade 101 – Bug Fixes
## What is available from 1st Release?



Installation Impact

**SMU ( Software Maintenance Unit, point fix )**

1. Hitless SMU : **64% of all bug fixes posted have NO traffic impact**

2. Traffic impact SMU: **11% of all bug fixes have traffic impact**

3. Reload SMU: **25% of bug fixes require a reload**

**Next step**

1. Complete the ISSU building blocks

# SMU Installation Impact



Pie chart legend: Hitless, Reload, Traffic Loss

- Hitless: 63%
- Reload: 26%
- Traffic Loss: 11%

# Recipe for High Availability - Increase MTBF, Reduce MTTR

## Continuous Systems Operation

**99.999+% Service Availability**

**No Single points of failure for unplanned or planned events**

**CSO**

**ISSU In Service Software Upgrade**

**Network High Availability Fast Convergence (PIC)**

**Hitless RP Switchover**

**Non-stop Forwarding & Routing**

**Distributed Route Processor, Service Separation Architecture**

**IOS XR Basics: microkernel based architecture granular process restart, protected memory, highly modular, separation of control, data, management planes, fault tolerance, packaging model**

**Hardware Design: redundancy (fabric, power, thermal, route processor, line card), high MTBF, distributed forwarding, Online Insertion Removal (OIR), parity or error correcting memory, fault insertion testing**

# Distributed Route Processor and IOS XR Service Separation Architecture

- Scaling control plane beyond basic RP with DRP

  Available since 2005

  Competitors had to react with announcement in 2008 !

- Service Domain Routers

  Independent/isolated physical routing instances

  Contain a subset of (d) RPs and LCs within a common (multi-) chassis.

  Solution available since 2005

  Lead customers: BT & Comcast



**SDR**

**SDR**

**SDR**

**Complete H/W Fault Isolation** ← **Service Separation Architecture** → **Complete Resource Sharing**

**Routing Instance A**

**Routing Instance B**

**Routing Instance C**

# Distributed Control Plane



- Routing protocols and signaling protocols can run in one or more (D)RP

- Each (D)RP can have redundancy support with standby (D)RP

- Out of resources handling for proactive planning

# Applications—BGP Multi-speakers



- Distributed BGP speakers to multiple RP and DRPs
- Single unified BGP RIB to external peers
- Achieve BGP peering scalability

# Recipe for High Availability - Increase MTBF, Reduce MTTR

## Continuous Systems Operation

**99.999+% Service Availability**

**No Single points of failure for unplanned or planned events**

**CSO**

**ISSU In Service Software Upgrade**

**Network High Availability Fast Convergence (PIC)**

**Hitless RP Switchover**

**Non-stop Forwarding & Routing**

**Distributed Route Processor, Service Separation Architecture**

**IOS XR Basics: microkernel based architecture granular process restart, protected memory, highly modular, separation of control, data, management planes, fault tolerance, packaging model**

**Hardware Design: redundancy (fabric, power, thermal, route processor, line card), high MTBF, distributed forwarding, Online Insertion Removal (OIR), parity or error correcting memory, fault insertion testing**

# High Availability Infrastructure



- **Distribution improves fault tolerance and recovery time by localizing the database and system management functionality to each node**

- **Granular process restart allows for fast recovery from failures**

- **IOS XR is designed to optimize the switch over between redundant hardware elements (RP, SC, PS, Fan C.)**

    **IOS XR is designed to route around fabric failure**

    **Line cards are protected by link bundling, APS, IPS, ECMP etc.**

# Continue to Route Despite Failure:
## Non Stop Routing (NSR)

**NSR Operation**

Cisco NSR-Enabled
Router Running BGP

**Goal:** **Maintain routing sessions during primary Route Processor failure**

**Implementation:**

- Currently, upon RP failover, routing sessions terminate. Protocols on standby RP reestablish sessions.

- With NSR, sessions are migrated from active to standby RP without notifying peers.

  No software upgrade on all routers

  No manual tuning of timers

  No additional load on peering routers

# What is NSR?

- **NSR is a self-contained solution to maintain the routing service (& hence the forwarding service) during:**

  RP/DRP fail-over

  Process restart

  In Service Software Upgrade

  Rack OIR in the case of Multi-chassis

- **No disruption to the routing protocol interaction with other routers.**

# What's the behavior today?

- **During RP failover:**

  Routing sessions terminate.

  Once standby RP becomes active, protocols reestablish the sessions, relearn the routes, and populate forwarding.

  NSF is designed to preserve the forwarding state while this is happening.

  But other routers in the network detect the failure and try to route around – huge network churn, could lead to forwarding loops and/or traffic loss.

- **Protocol Extensions – Graceful Restart (GR)**

  Neighboring routers detect the failure, but do not propagate the failure in the network.

  They also assist the router in coming back up.

# GR vs. NSR

- **Graceful Restart**

  Requires the software on all routers to be upgraded.

  Requires manual tuning of timers – If not correctly done, GR won't help.

  Adds load on the peering routers which could cause instability.

  Introduces a window in which forwarding loops and traffic loss can happen.

  Not all vendors have implemented GR

- **With Non-stop Routing:**

  Sessions don't terminate during failover.

  Routing interaction continues on the newly active RP without peers being aware.

# Evolution

# NSF, GR, NSR…

| | NSF | GR | NSR |
|---|---|---|---|
| **Forwarding plane kept intact** | Yes | Yes | Yes |
| **Session Failure** | Yes | Yes | No |
| **Failure propagation in network** | Yes | No | No |
| **Handling topology changes** | No | No | Yes |
| **Protocol extensions needed** | No | Yes | No |

# Complementarity

- NSF is the building block – needed for any HA solution

- GR is _not_ required for local control plane failure

    Sessions stay up with NSR

- GR required for remote control plane failure that don't support NSR

    Session will go down

    GR helper role has to get triggered

- GR also helps as a fallback option

# NSR & ISSU

- ## NSR is the building block for ISSU

  After standby RP upgrade to a newer version, it needs to get to NSR ready state before the RP failover step.

© 2007 Cisco Systems, Inc. All rights reserved.

# Design



- ❑ **Mirror canonical state to standby**
  - ❑ **Minimize chatter; Use efficient, asynchronous, ISSU-adhering communication**
- ❑ **Utilize H/W capabilities for receiving & sending protocol packets on both active & standby**
- ❑ **Make standby processes run "independent"**

# NSR Phases

- ## Synchronization

  Initial state mirroring between active and standby RPs to get standby processes up-to-date.

- ## NSR-Ready

  The active and standby stacks operate independently.

  Incoming packets are replicated to both RPs.

  Outgoing segments are sent out through standby RP.

  Infrequently synchronizing application state (optimize Post-FO behavior)

  Use asynchronous IPC between active and standby stacks.

- ## Switch-over

  Restore keepalive functionality to keep sessions up.

  Continue from where (previously) active RP left.

# Supported HA events

- **Supported HA events**

    RP/DRP switchover

    Process restarts, SMU/pie installs

    - Supported by triggering RP switchover via configuration or fault manager scripts

    - SMU/pie installs need manual ordered activation (activate first on the standby)

# Other Protocols

- ISIS

    Checkpoints necessary state to the standby

    Able to keep adjacencies and retrieve all the necessary state with the current protocol mechanisms

- RSVP

    Checkpoints necessary state to the standby

    Retrieve the necessary state from neighbors (soft-state) with and without refresh reduction

- PIM

    Checkpoints necessary state to the standby

    Sends hellos with a different GenID to get the refreshed state from peers

- Summary

    All these protocols are able to recover state from network within the default timeout period

    There is still some churn noticeable in the network

# TCP NSR Send and Receive Paths

**ACTIVE RP**

BGP

**1**

TCP

**3**

**2**

**STANDBY RP**

BGP

TCP

**4**

Egress LC

**Send Path**

---

**ACTIVE RP**

route DB

BGP

**4**

Session DB

TCP

**3** ACK

**2**

**STANDBY RP**

route DB

BGP

**4**

Session DB

TCP

Ingress LC

**1**

**Receive Path**

# Recipe for High Availability - Increase MTBF, Reduce MTTR

## Continuous Systems Operation

**99.999+% Service Availability**

**No Single points of failure for unplanned or planned events**

**CSO**

**ISSU In Service Software Upgrade**

**Network High Availability Fast Convergence (PIC)**

**Hitless RP Switchover**

**Non-stop Forwarding & Routing**

**Distributed Route Processor, Service Separation Architecture**

**HA Basics: microkernel based architecture granular process restart, protected memory, highly modular, separation of control, data, management planes, fault tolerance, packaging model**

**Hardware Design: redundancy (fabric, power, thermal, route processor, line card), high MTBF, distributed forwarding, Online Insertion Removal (OIR), parity or error correcting memory, fault insertion testing**

# Principle of Simplicity

- "Simplicity is prerequisite for reliability"
  Edsger Dijkstra

- "Simplicity is the ultimate sophistication"
  Leonardo da Vinci

- Kiss: Keep
       It
       Simple
       Straighforward

Under Cisco NDA – Extremely Confidential. Do not distribute     46

# Kiss Principle



Gain

Academic optimum

KISS optimum

Complexity (cost)

# Routing Resilience

- Focuses on the <u>loss of connectivity</u> during a rerouting event

# Requirement

- < 1 or 2 second: human does not bother

- < 200msec: human does not notice

- < 50msec: human has the perception that
  he would be better off

# Requirement

- 99.999% Availability

  25920 LoC in msec / month

  518 core events @ 50msec per month

  Note: an event counts only if it impacts the reachability of a specific customer. Furthermore, the unreachability is most likely partial.

# Principle of Simplicity

- Principle of Simplicity asks

  what is the cost/complexity of fullfilling each requirement?

  what is the frequency of these events

- KISS: if very complex and infrequent, don't do it

- KISS: 200msec in any cases and 50msec in most cases might be the optimum

# Fast Convergence Roadmap

## Multi-second Convergence

- Fast External Failover
- Next Hop Tracking
- VRF scoping (IOS XR)
- BGP Local Convergence (IOS XR)

## Sub-second Convergence (typically sub-200 msec)

- BGP Prefix Independent Convergence (PIC) – Core
- BGP Prefix Independent Convergence (PIC) – Edge MultiPath
- BFD
- IGP & LDP Non Stop Routing
- ISIS & OSPF Prefix Prioritization
- BGP Non Stop Routing
- BGP Prefix Independent Convergence (PIC) – Edge Unipath Primary/Backup
- BGP/VPN Internet Path Diversity

## 50 msec Convergence

- MPLS TE Fast ReRoute
- ISIS IP Fast ReRoute
- Multicast only Fast ReRoute
- OSPF IP Fast ReRoute
- MPLS LDP Fast ReRoute
- Label Switched Multicast

# IGP convergence

# The "base"

# What is IGP convergence?



- An event within the Service Provider network that causes a RIB change for one or more IGP prefixes
  - What event: link (PE-P, P-P) or node (P,PE) failure
  - Indeed, edge router failure is a special case of IGP convergence
- This event may result in loss to BGP/VPN destinations
- BGP bestpath may get triggered upon IGP convergence

# Convergence time defined



- Assume a flow from Src A to Dest C

- T1: when L dies, the best path is impacted: loss of traffic

- T2: when the network converges, traffic reaches the destination again

- Loss of Connectivity for flow AC: T2 – T1, called "convergence" hereafter

# IGP convergence always matters

- BGP nhop availability

- PIM source availability

- MPLS TE topology and resource information

- Unplanned protection
  - Most MPLS FRR designs only protect link
  - Unknown SRLG

- Catastrophic event

# IGP Convergence - State of the Art

- Prefix Prioritization is THE key behavior

    CRITICAL: IPTV SSM sources

    HIGH: Most Important PE's

    MEDIUM: All other PE's

    LOW: All other prefixes

- Prefix prioritization customization is required for CRITICAL and HIGH

# IGP Convergence – State of the Art



36112-crs1-P-ISIS00-500Pr-2kT-i2i-Layer1down-050708-lnlb

# IGP Convergence – cannot be any simpler

- No tuning required

- Works for any design, for any failure

- Is anyway required

Pranav Dharwadkar Under Cisco NDA – Extremely Confidential. Do not distribute

# Flooding matters within a limited diameter



A                                    Z

**What matters is when this router converges.**

**What the study showed is that this "rerouting" router is never in a different continent and very rarely further away than 5 hops**

**Sure, this router at the end of the network will also converge.**

**Sure, there may be much more than 5 hops in the worst diameter case.**

**But this worst-case is not seen in practice, because a router closer to the failure will have rerouted earlier**

# IGP Convergence - Flooding





- ICI Sponsored Research – Olivier Bonaventure and Pierre Francois:

    a link failure within a continent very rarely requires a rerouting in a different continent. Propagation is thus bounded by 25msec (5000km of fiber)

    it is very rare for a failure to require rerouting further than 5 hops away from the failure. Flooding is thus bounded by 5*5msec.

- Intuitively, this rule is expected: designers build networks with resilience in mind

    most often (at the POP level)

# IGP Convergence - Flooding



Number of FIB corrections upon early FIB update after a node failure

X-Axis : pairs of (node failure, rerouting node for that failure). All rerouting nodes are plotted for each node failures

Y-Axis : The number of prefixes that got updated twice because the first one considering the link failure was wrong

Note: it doesn't account the FIB updates of the second run that were not performed in the first one. This is not "bad".

- ICI Sponsored Research, Pierre Francois:

- Aggressive Link-Oriented IGP convergence very rarely incur overhead in case of node failures. There are some rare cases where a node would have to update for half of (and one for almost all of) its IGP originated destinations.

- Intuitively, this is expected: very often, once you divert due to a link failure, you completely avoid the previous path and hence avoid the node

# IGP Convergence – SPT computation



iox-crs1-P-ISIS-███-dpIGP i2 i-dpBGP ipv4 i2 i-041406-r-nlb

(crs1e8-1,SPTend) - (crs1e8-1,SPTstart)

- 900-router ISIS network… without leveraging i-SPT

  with iSPT, most runs under 1msec

# LossLess Managed Operation

  Cisco Confidential

# LossLess Link In/Out service

- Managed event

- Link is brought out or in service

- LossLess methodology

  shut the IGP adjacency first

  wait a few seconds

  shut the interface or unplug the fiber



36112-crs1-P-ISIS00-500Pr-2kT-i2i-isisIntdown-050808-lnlb

Agilent measurements

```
RP/0/RP1/CPU0:crs1e8-1(config)#router isis 1

RP/0/RP1/CPU0:crs1e8-1(config-isis)#int GigabitEthernet0/6/0/4

RP/0/RP1/CPU0:crs1e8-1(config-isis-if)#shut

RP/0/RP1/CPU0:crs1e8-1(config-isis-if)#commit
```

# LossLess Link In/Out service

- Operation only required on one end of the link

  when the IGP adjacency is shut, the IGP sends a hello which brings the adj down on the other end

  when the IGP adjacency is shut, BFD sends a notification (down) to the other end

  when the IGP adjacency is shut, an LSP/LSA is resent. Two-Way-Connectivity-Check at the other end fails upon reception of this LSP/LSA

- Less error prone than using max-metric

  the operator might not reconfigure the correct metric

- Can be extended to SPA and Linecard IN/OUT of service

# Faster Convergence Post-Failure:
## Prefix Independent Convergence (PIC)

## Goal:

**Reduce convergence time to 100s of *msec* for all prefixes**

## Implementation:

Spans IGPs, BGP, MPLS (LDP), & FIB components

Applies to Core and/or Edge link and node failures

Assumptions:

Single failure

Alternate paths exist in the core

**PE1's RIB:**

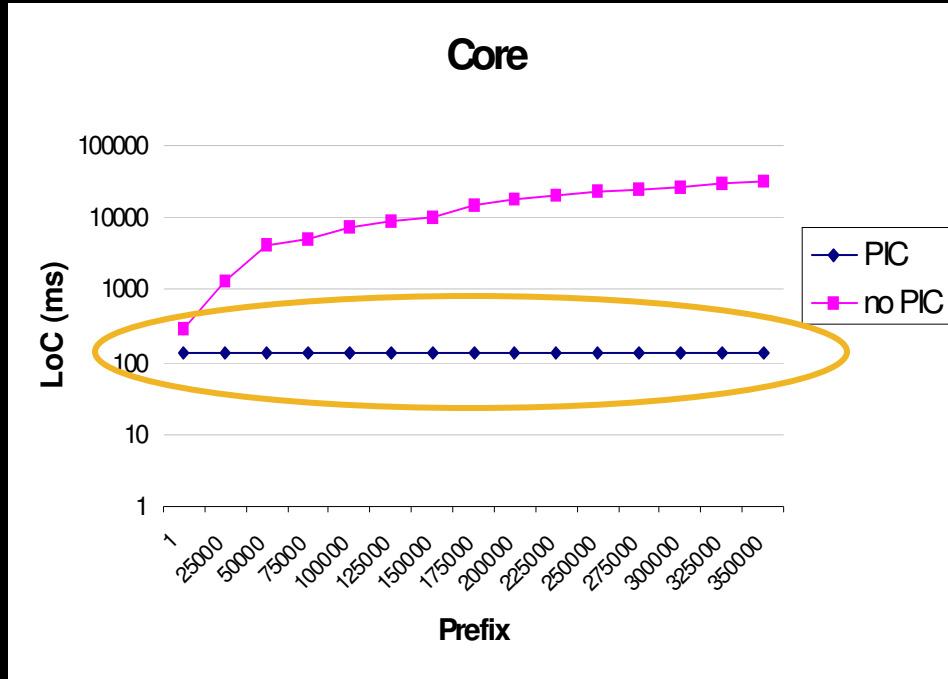10.1.1/24 via PE2
10.1.2/24 via PE2

When ?

....

99.99.99/24 via PE2

**PE2 is reachable via   PoS2**

PE2

PoS1

PE1

PoS 2

# PIC network scenarios



- Looking back to all IGP related failures, we distinguish between:
    - PIC Core: path to a NH changes (failure 1, 2 and 3)
    - PIC Edge: NH delete (failure 4, 5)

# BGP PIC Core
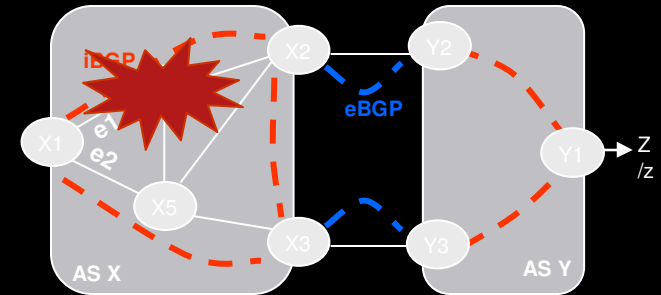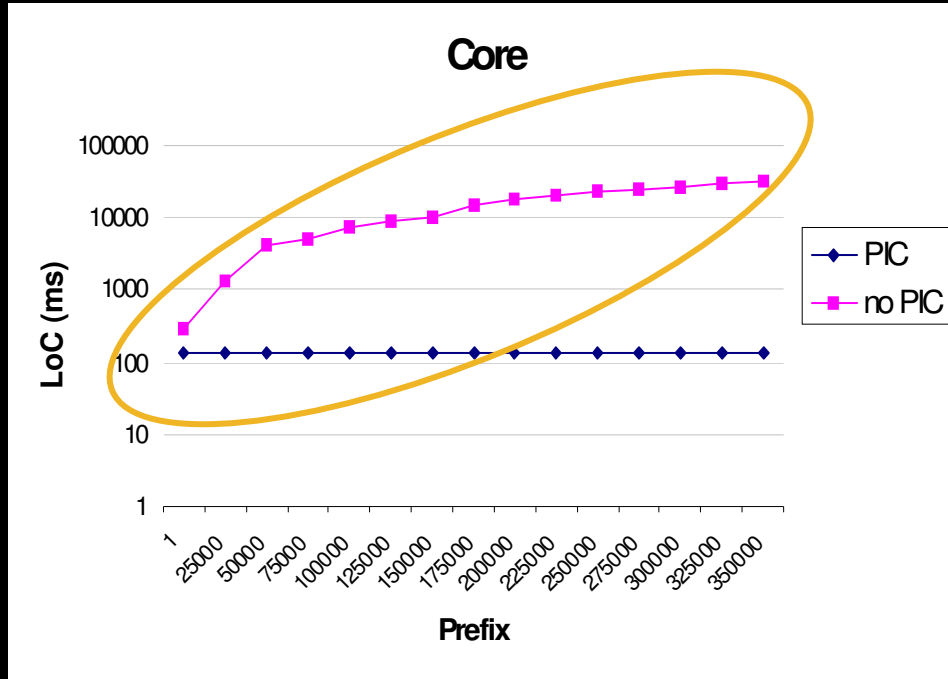
## Convergence Independent of VPN/BGP Route Scale

**P1**

**P2**

**P3**

**PE3**

**P4**

**CE1**

**PE1**

**VPN1
site1**

**CE2**

A/a
B/b
.......
N/n

**VPN1 Site2**

# Characterization
# BGP PIC Core Analysis



**Core**



- **BGP PIC Core:**

  Sub-second convergence upon PE uplink failure
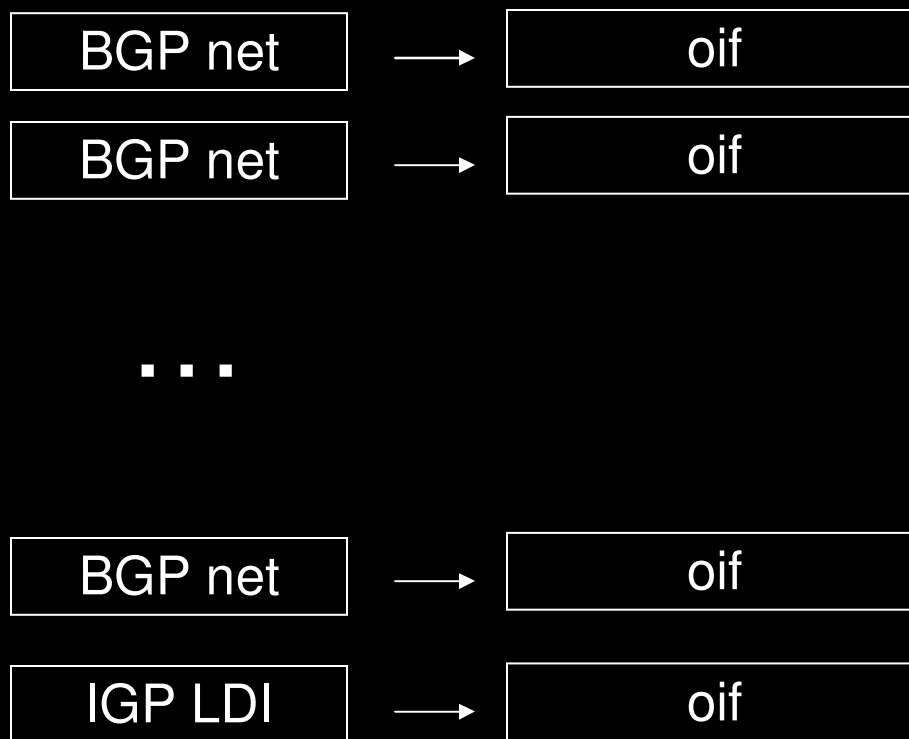
# Characterization Without BGP PIC Core



**Core**

- PIC
- no PIC

Axis labels: LoC (ms) vs Prefix

Prefix values: 1, 25000, 50000, 75000, 100000, 125000, 150000, 175000, 200000, 225000, 250000, 275000, 300000, 325000, 350000



X2

| X2 via X4 | X2 via X5 | Z1/z1 | Z2/z2 | ... |

Z349999/z349999

Z350k/z350k

Adj to X4 (E1 with DMAC …)

Adj to X5 (E2 with DMAC …)

- Without BGP PIC Core:

  Up to 10's of seconds of loss for PE uplink failures

# The right architecture: hierarchical FIB

BGP net

BGP net

...

BGP net

BGP LDI → IGP LDI → oif

BGP nexthop(s)　　IGP nexthop(s)　　Output Interface

- Pointer Indirection between BGP and IGP entries allow for immediate leveraging of the IGP convergence
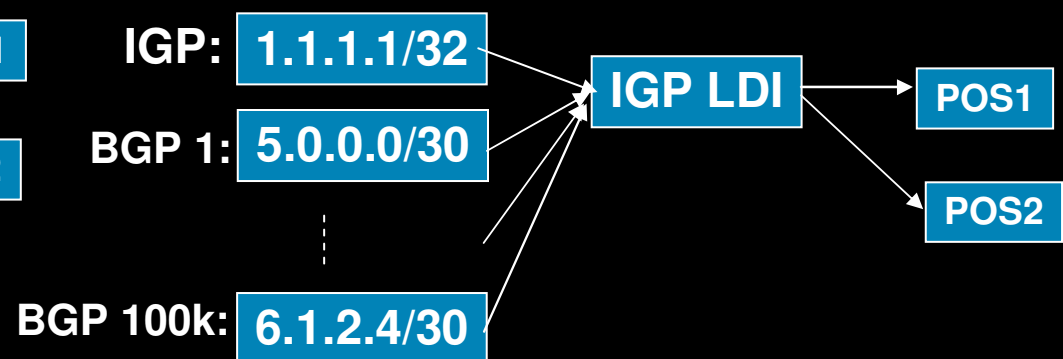
# The unoptimal way: flattened FIB

| BGP net | → | oif |
|---------|---|-----|

| BGP net | → | oif |
|---------|---|-----|

**. . .**

| BGP net | → | oif |
|---------|---|-----|

| IGP LDI | → | oif |
|---------|---|-----|

- Control Plane flattens the recursion such that any BGP FIB entry has its own local oif information

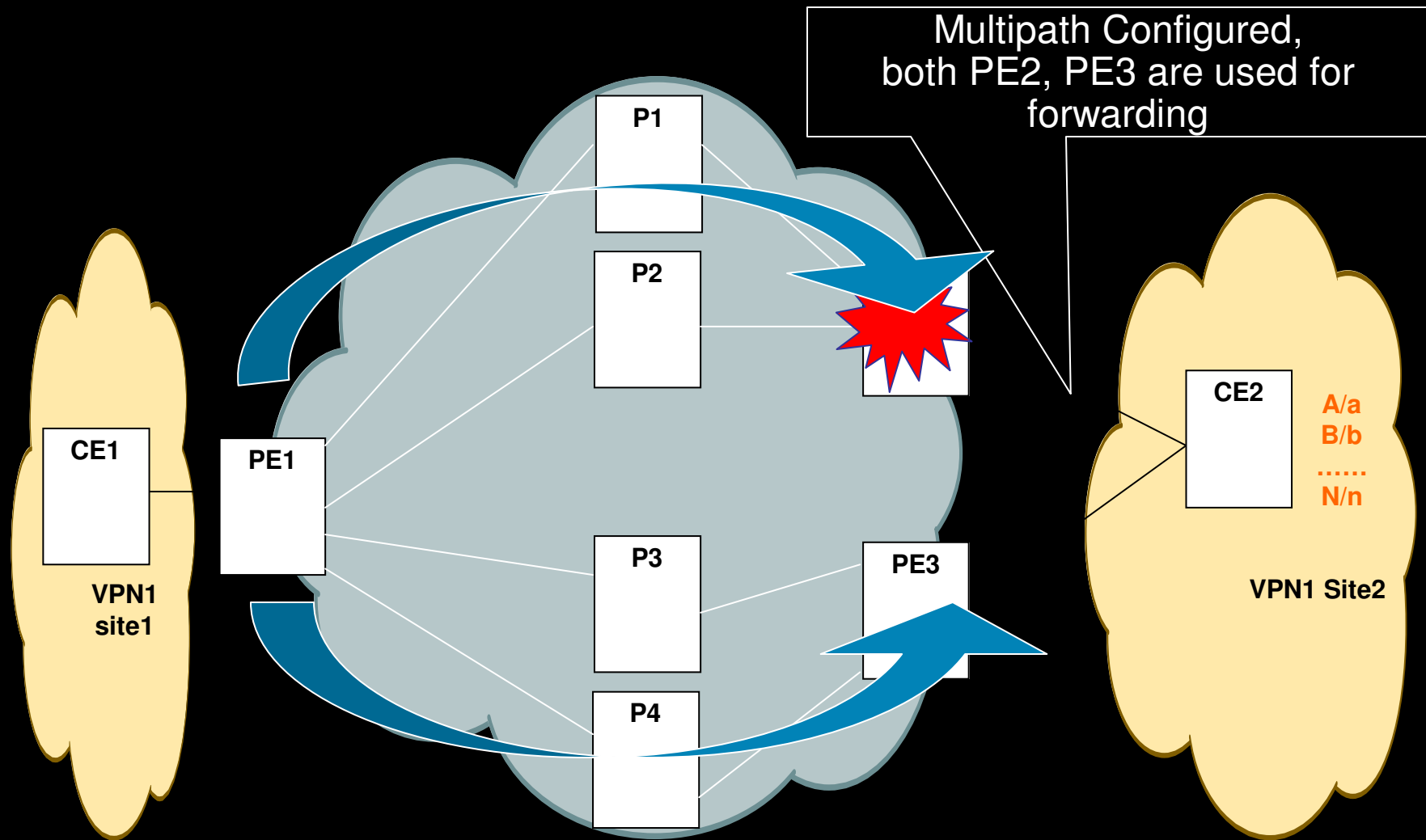Pranav Dharwadkar

# PIC vs no PIC example

**No PIC**                                          **PIC**

IGP: 1.1.1.1/32   POS1        IGP: 1.1.1.1/32   IGP LDI   POS1
BGP 1: 5.0.0.0/30   POS2      BGP 1: 5.0.0.0/30           POS2
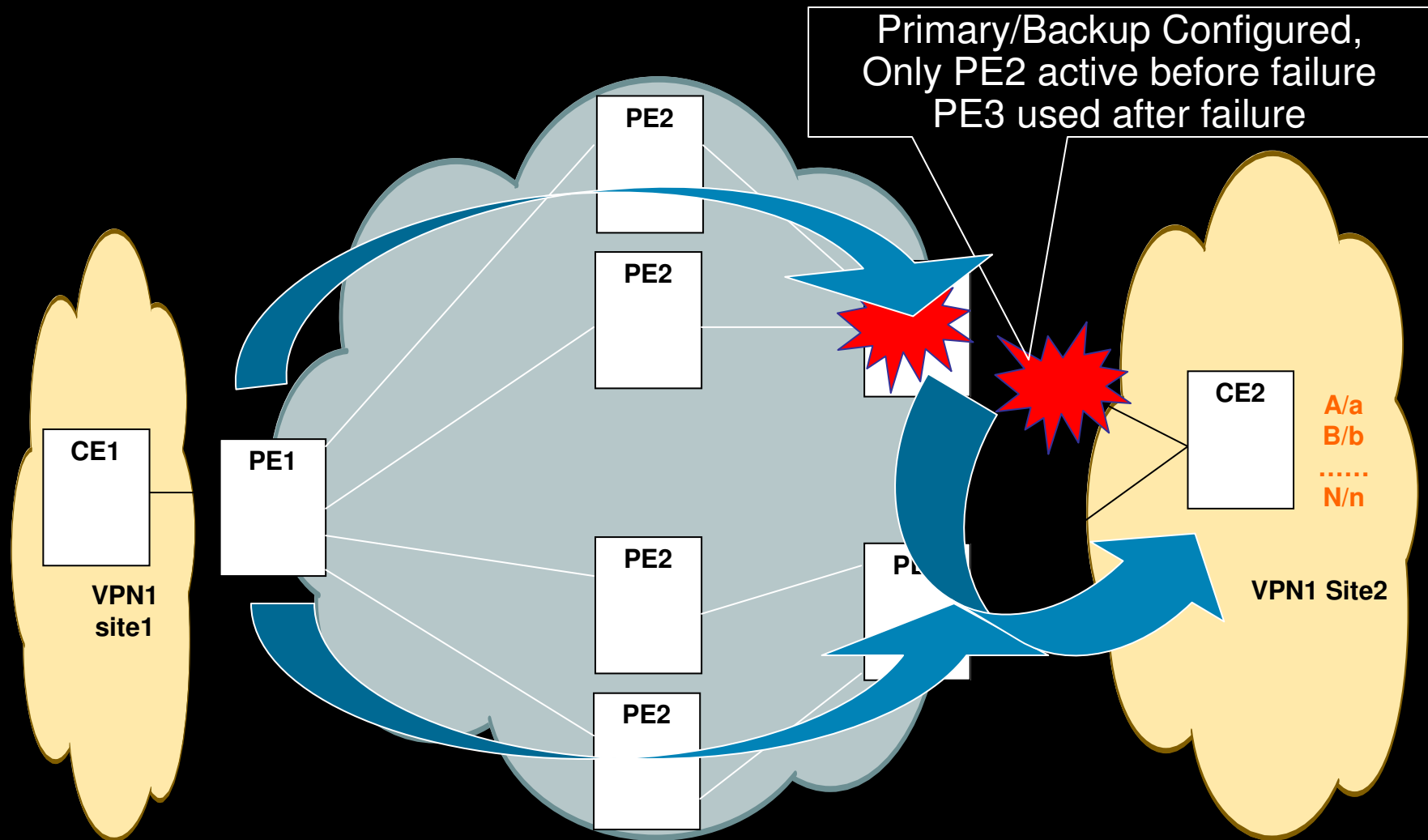BGP 100k: 6.1.2.4/30          BGP 100k: 6.1.2.4/30

- No PIC: extra processing needed to 'backwalk' all dependent BGP prefixes and re-resolve them.

- PIC: No extra processing needed thanks to the creation of IGP LoadInfo.

- **Note: this slide only intends to demonstrate the extra operations needed when PIC is not supported and is not a representation of the real FIB structures at all!!**
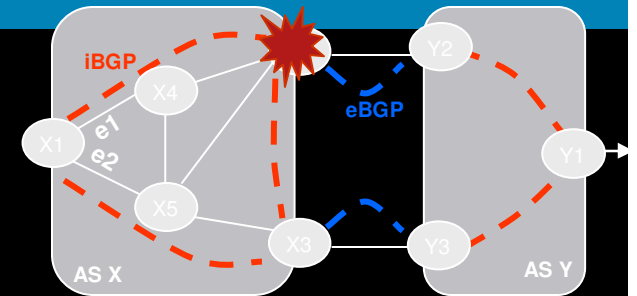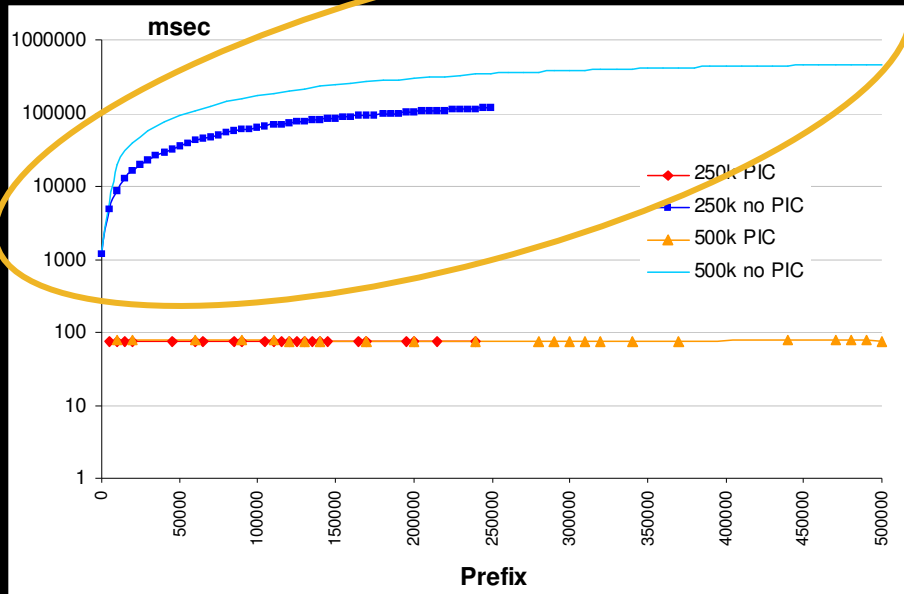
# BGP PIC Edge-Multipath

Multipath Configured,
both PE2, PE3 are used for
forwarding

**P1**

**P2**

**CE2**

A/a
B/b
......
N/n

**CE1**

**PE1**

**P3**

**PE3**

**VPN1 Site2**

**VPN1
site1**

**P4**

# BGP PIC Edge Primary/Backup

Primary/Backup Configured,
Only PE2 active before failure
PE3 used after failure

**PE2**

**PE2**

**PE2**

**PE2**

**CE1**

**PE1**

**PE2**

**CE2**

A/a
B/b
......
N/n
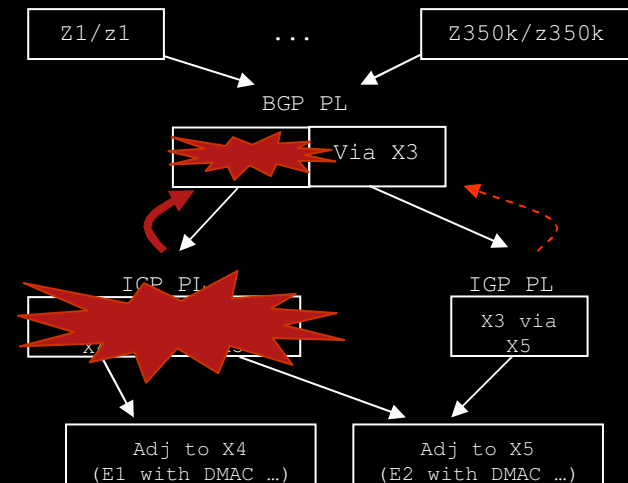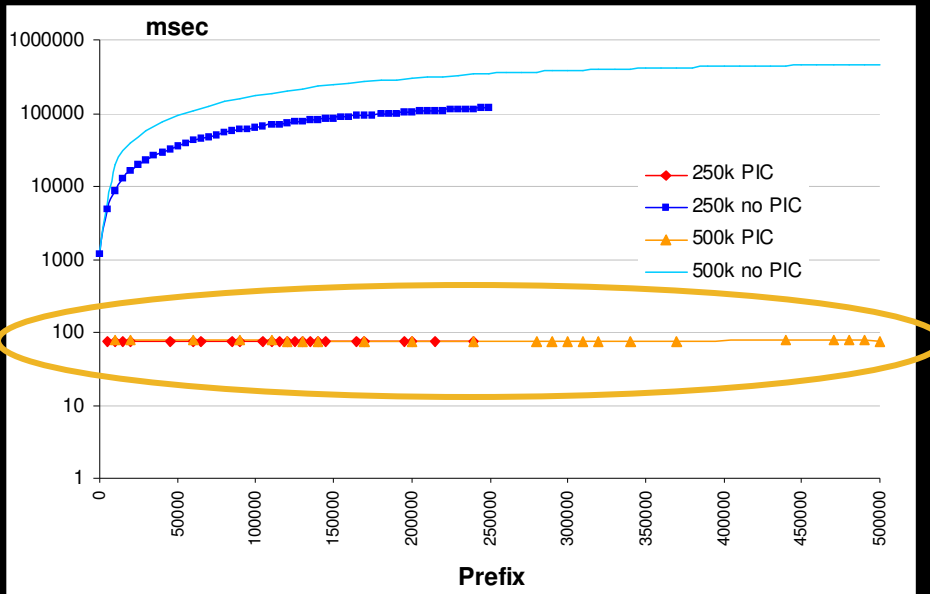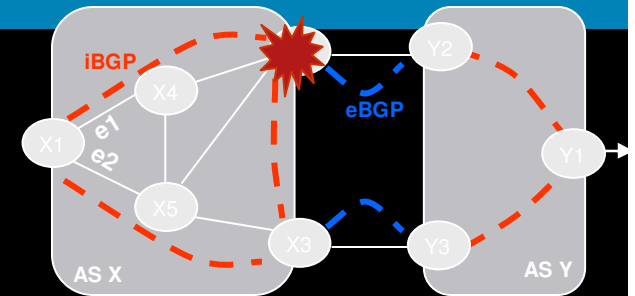
**VPN1 site1**

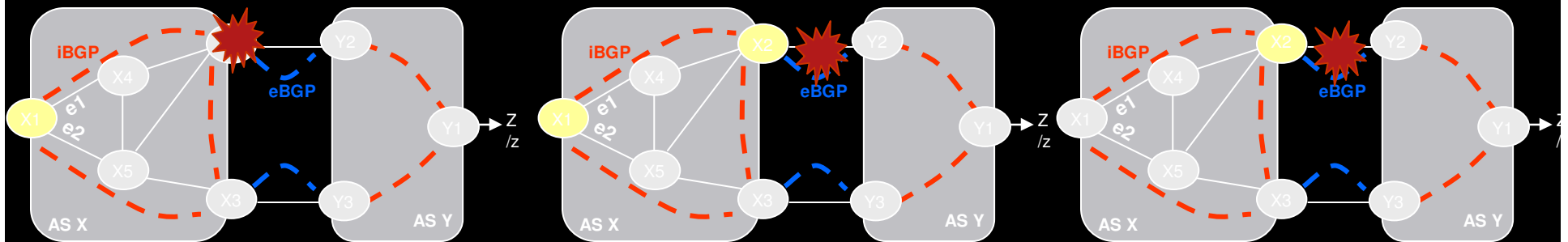**VPN1 Site2**

# Characterization
# Without BGP PIC Edge



- At IGP Convergence time, in a flattened dataplane FIB, all the BGP entries recursing via X2 point to an invalid path. No dataplane protection is possible.

- The control plane convergence is now required to move each BGP entry onto an alternate next-hop and then update the flattened dataplane FIB accordingly. This may take minutes.

# Characterization With BGP PIC Edge







- At IGP Convergence time, the complete IGP PL to X2 is deleted. SW FIB walks the linked list of parent BGP PL and in-place modify them to use alternate ECMP best nhops or enable alternate next-best nhops. This is quick because the BGP PL sharing is efficient.

- The control plane convergence still occurs in the background (blue curve) but its slowness does not impact dataplane connectivity and hence the T-SLA experience

# BGP PIC Edge application point



X2 does not set next-hop-self

X2 sets next-hop-self

X1 must be the reacting point as X2 is down.

X1's reaction is triggered by IGP convergence

X1 and X2 may be the reacting point.

X1 reaction is triggered by IGP convergence

X2 reaction is triggered by local interface failure

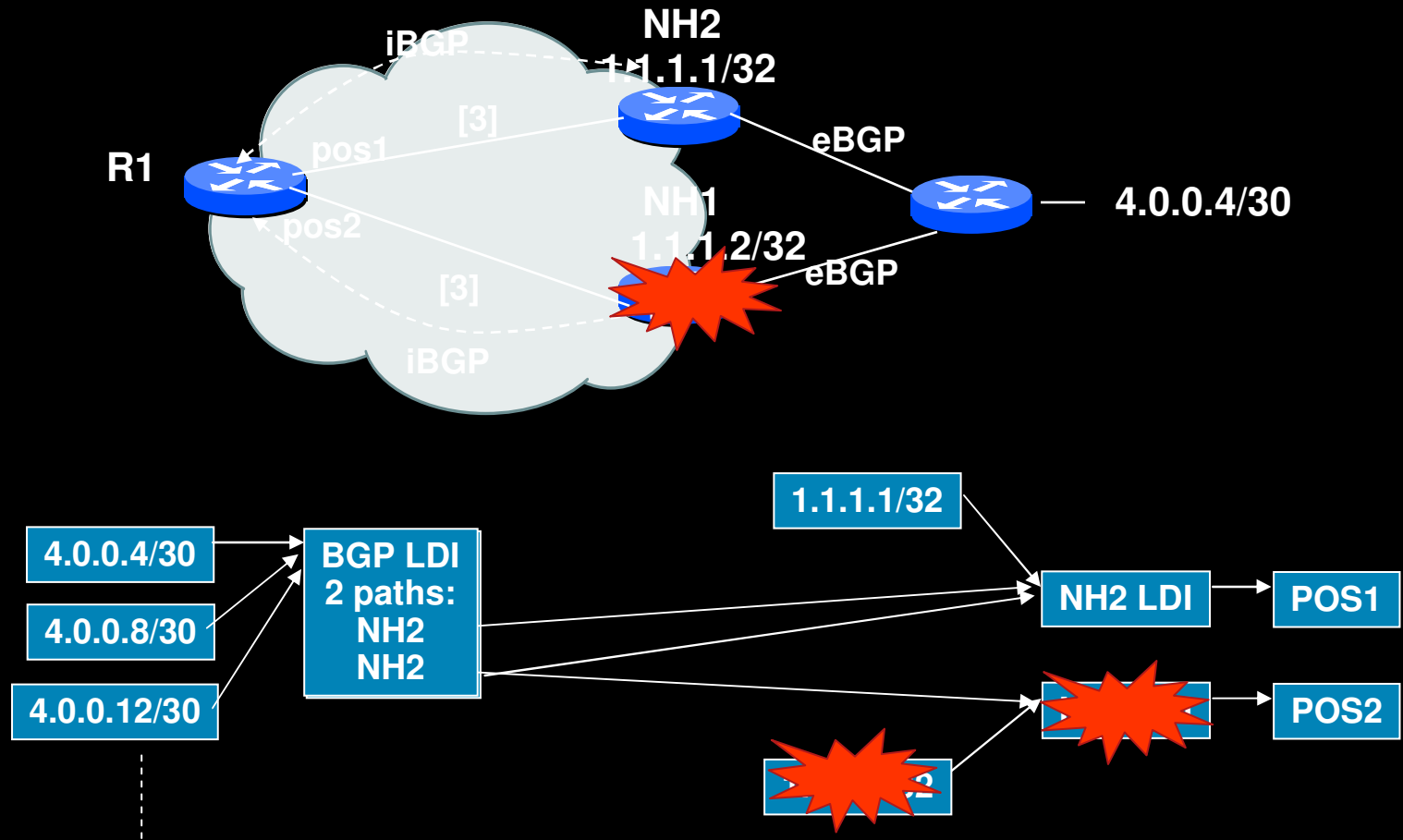X2 is the reacting point.

X2 reaction is triggered by local interface failure

Note: X1 is blind in this case as the next-hop is X2

# PIC Edge FIB perspective: example

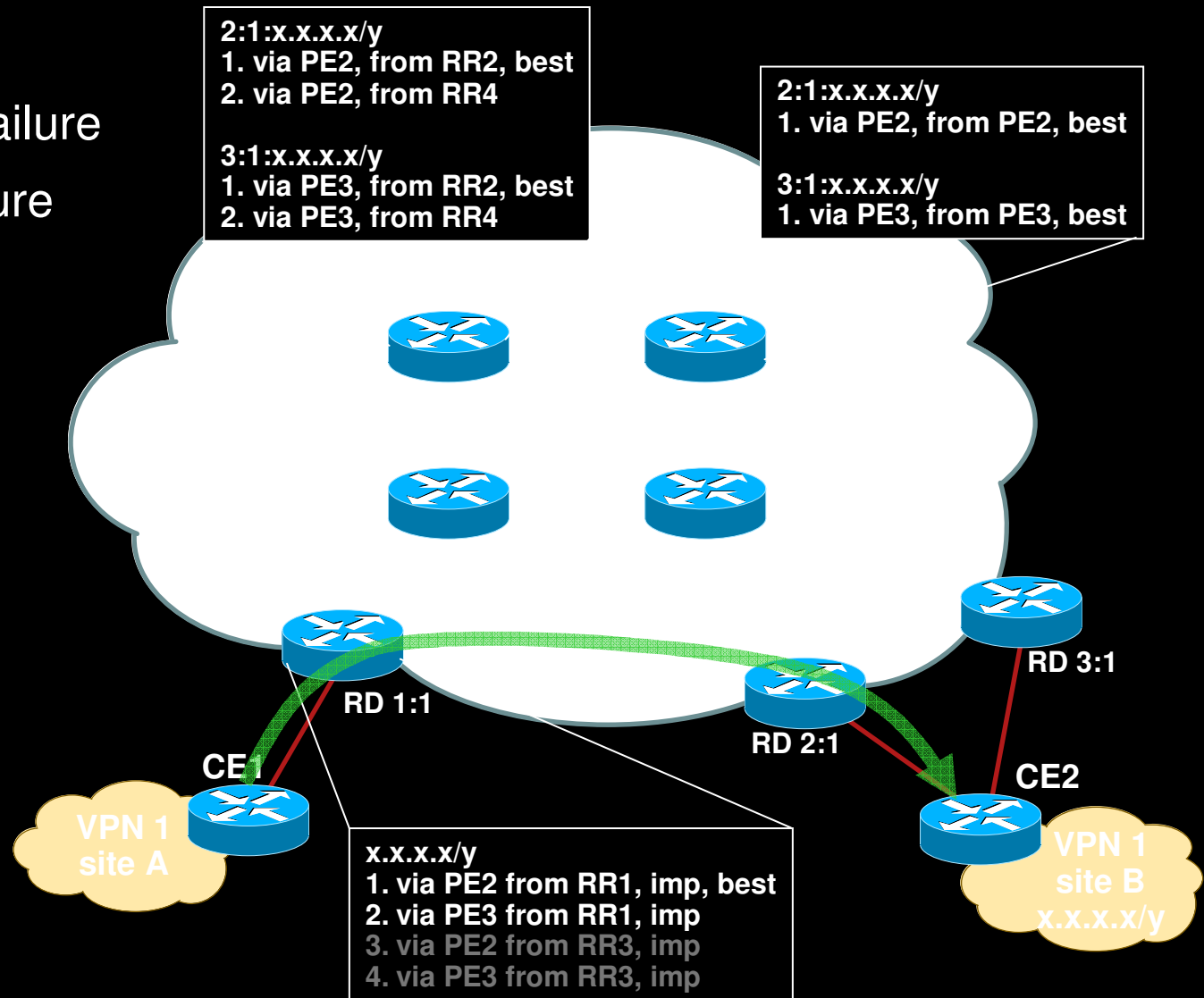# PIC Edge requires at least two paths

- For edge failures, our BGP PIC solution requires the availability of at least two paths to any BGP destination which require tight availability.

- Unique-RD-allocation guarantees this property for Classical L3VPN, IAS-A and is the key building block for the IAS-B, IAS-C and CsC solutions.

# PIC Edge requires at least two paths
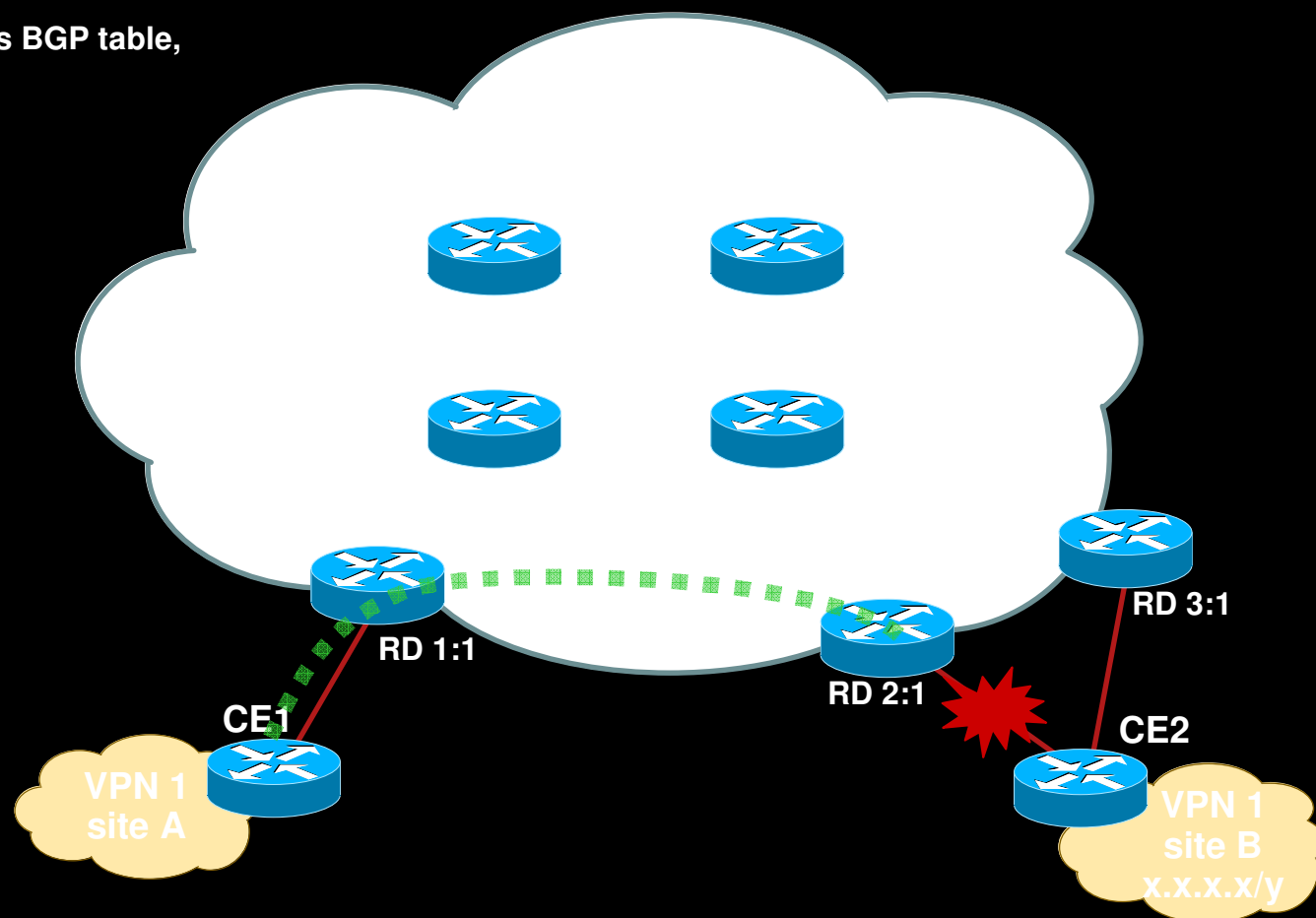
- Internet

- IAS-B

- IAS-C / CsC

# Edge Failures

- PE-CE link failure
- PE node failure

**2:1:x.x.x.x/y**
**1. via PE2, from RR2, best**
**2. via PE2, from RR4**

**3:1:x.x.x.x/y**
**1. via PE3, from RR2, best**
**2. via PE3, from RR4**

**2:1:x.x.x.x/y**
**1. via PE2, from PE2, best**

**3:1:x.x.x.x/y**
**1. via PE3, from PE3, best**

**RD 1:1**

**RD 3:1**

**RD 2:1**

**CE1**

**CE2**

**VPN 1 site A**

**VPN 1 site B x.x.x.x/y**

**x.x.x.x/y**
**1. via PE2 from RR1, imp, best**
**2. via PE3 from RR1, imp**
**3. via PE2 from RR3, imp**
**4. via PE3 from RR3, imp**

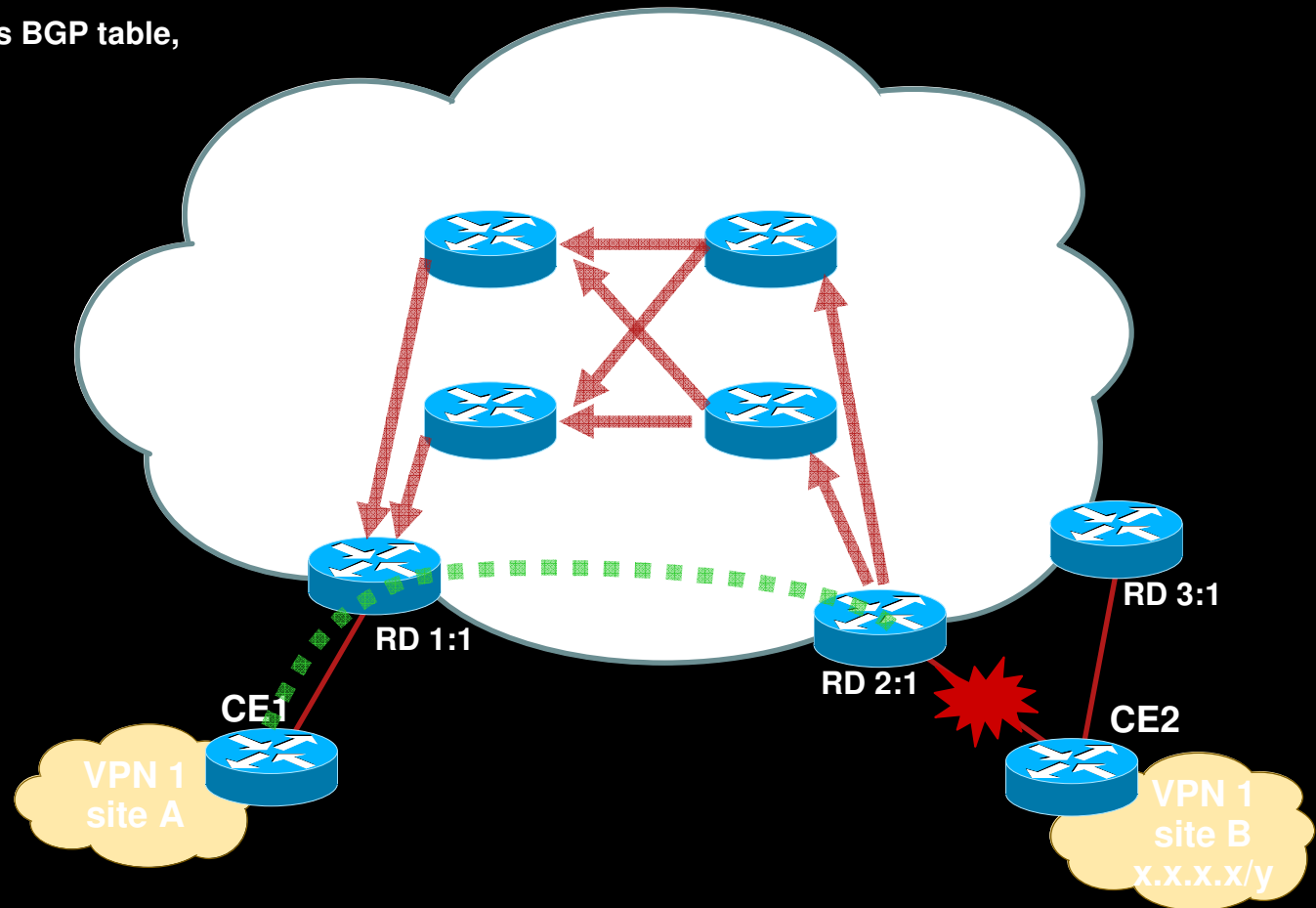# PE-CE link failure

1. link PE2-CE2 fails
2. Fast External Fallover scans BGP table, calculating new bestpaths



RD 3:1

RD 1:1

RD 2:1

CE1

CE2

VPN 1
site A

VPN 1
site B
x.x.x.x/y

# PE-CE link failure

1. **link PE2-CE2 fails**
2. **Fast External Fallover scans BGP table, calculating new bestpaths**
3. **PE2 withdraws paths**
4. **RR2 and RR4 propagate withdraws**
5. **RR1 and RR3 propagate withdraws**

RD 1:1

RD 2:1
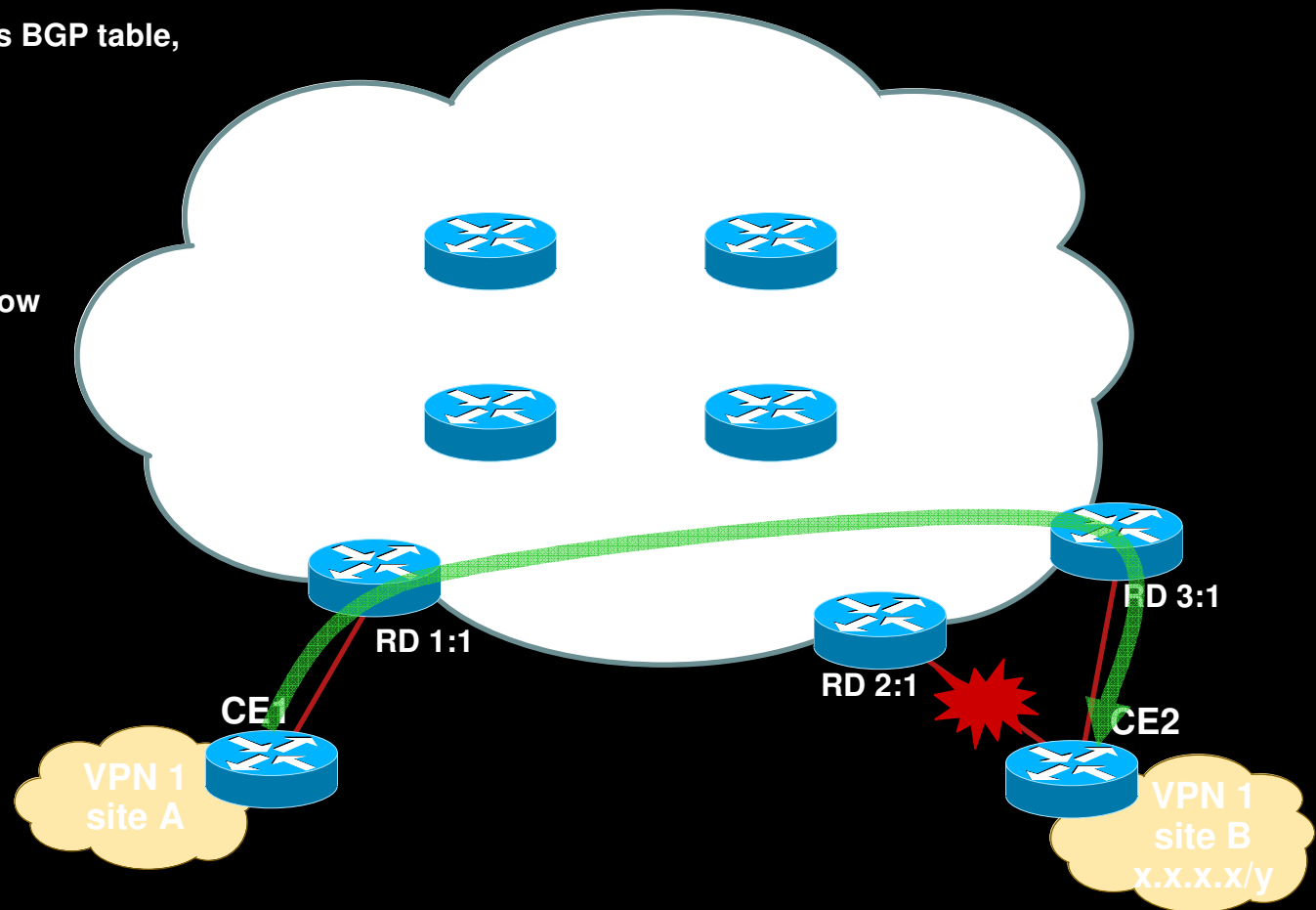
RD 3:1

CE1

CE2

VPN 1 site A

VPN 1 site B
x.x.x.x/y

# PE-CE link failure

1. **link PE2-CE2 fails**
2. **Fast External Fallover scans BGP table, calculating new bestpaths**
3. **PE2 withdraws paths**
4. **RR2 and RR4 propagate withdraws**
5. **RR1 and RR3 propagate withdraws**
6. **PE1 deletes path via PE2, now going via PE3**

**RD 3:1**

**RD 1:1**

**RD 2:1**

**CE1**

**CE2**

**VPN 1 site A**

**VPN 1 site B x.x.x.x/y**

# PE-CE link failure

- Convergence depends on
  - time to detect failure
  - time to scan BGP table (full versus scoped walk)
  - time to generate/propagate all withdraws
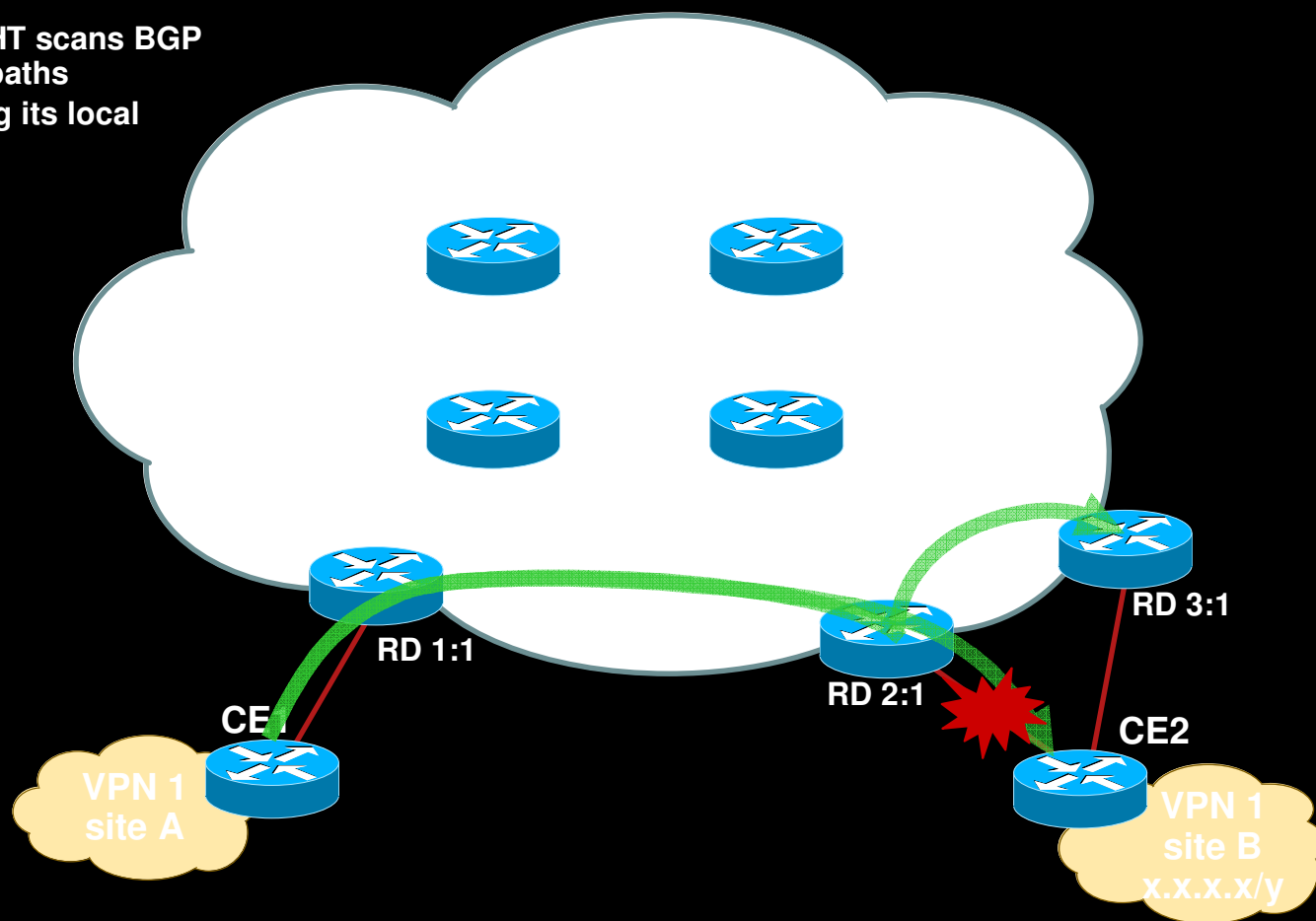  - time to process withdraws and update FIB

**Eliminate this with "local convergence" functionality, "local convergence" is prefix dependent convergence**

**Eliminate this with "PIC Edge" functionality, "PIC Edge" is Prefix Independent Convergence (PIC)**

# PE-CE link local convergence

1. **link PE2-CE2 fails**
2. **Fast External Fallover or NHT scans BGP table, calculating new bestpaths**
3. **PE2 updates FIB, preserving its local label for x.x.x.x/y**



**RD 3:1**

**RD 1:1**

**RD 2:1**

**CE1**

**CE2**

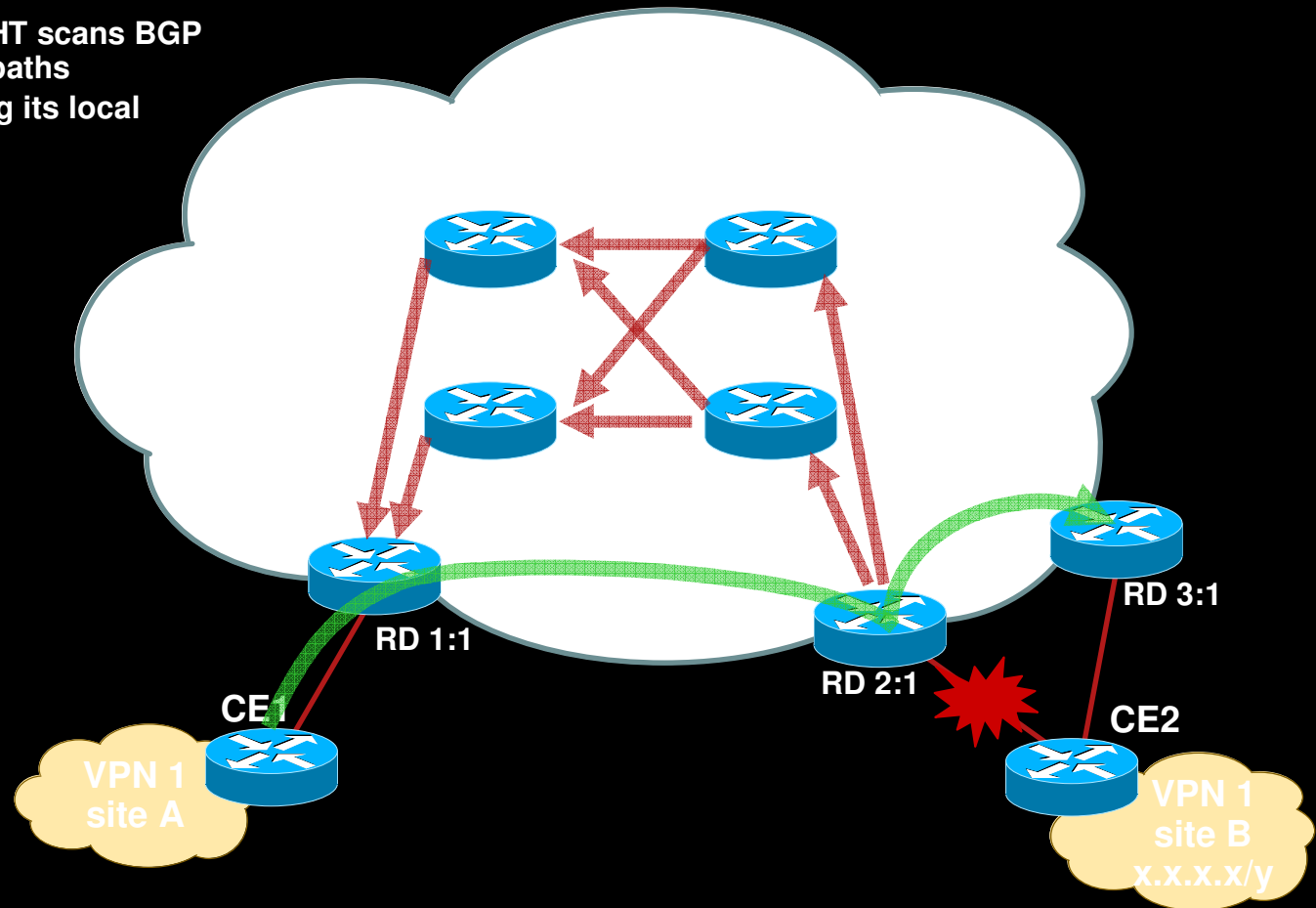**VPN 1 site A**

**VPN 1 site B x.x.x.x/y**

# PE-CE link local convergence

1. **link PE2-CE2 fails**
2. **Fast External Fallover or NHT scans BGP table, calculating new bestpaths**
3. **PE2 updates FIB, preserving its local label for x.x.x.x/y**
4. **PE2 withdraws paths**
5. **RR2 and RR4 propagate withdraws**
6. **RR1 and RR3 propagate withdraws**



RD 1:1

RD 2:1

RD 3:1

CE1

CE2

VPN 1 site A

VPN 1 site B x.x.x.x/y

# PE-CE link local convergence

1. **link PE2-CE2 fails**
2. **Fast External Fallover or NHT scans BGP table, calculating new bestpaths**
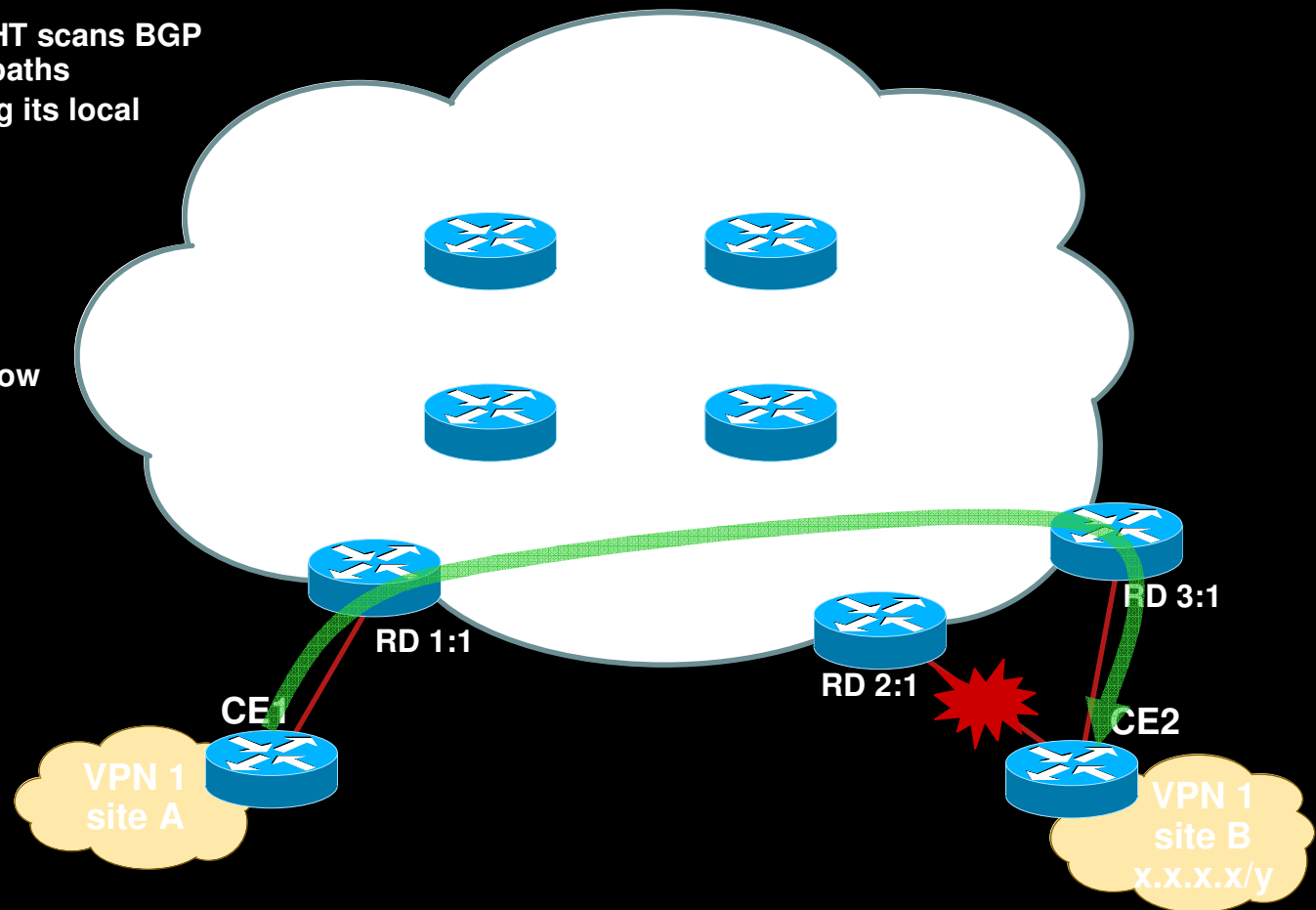3. **PE2 updates FIB, preserving its local label for x.x.x.x/y**
4. **PE2 withdraws paths**
5. **RR2 and RR4 propagate withdraws**
6. **RR1 and RR3 propagate withdraws**
7. **PE1 deletes path via PE2, now going via PE3**



RD 1:1

RD 3:1

RD 2:1

CE1

CE2

VPN 1 site A

VPN 1 site B x.x.x.x/y

# PE-CE link local convergence

- local label preservation

  IOX

  Connected prefixes use **aggregate label**

  Static routes, eBGP prefixes need **per-VRF label**

  IOS

  keep **"zombie" label** for 5 minutes (EDCS-500998)

# PE-CE link failure

- Convergence depends on

  D: time to detect failure

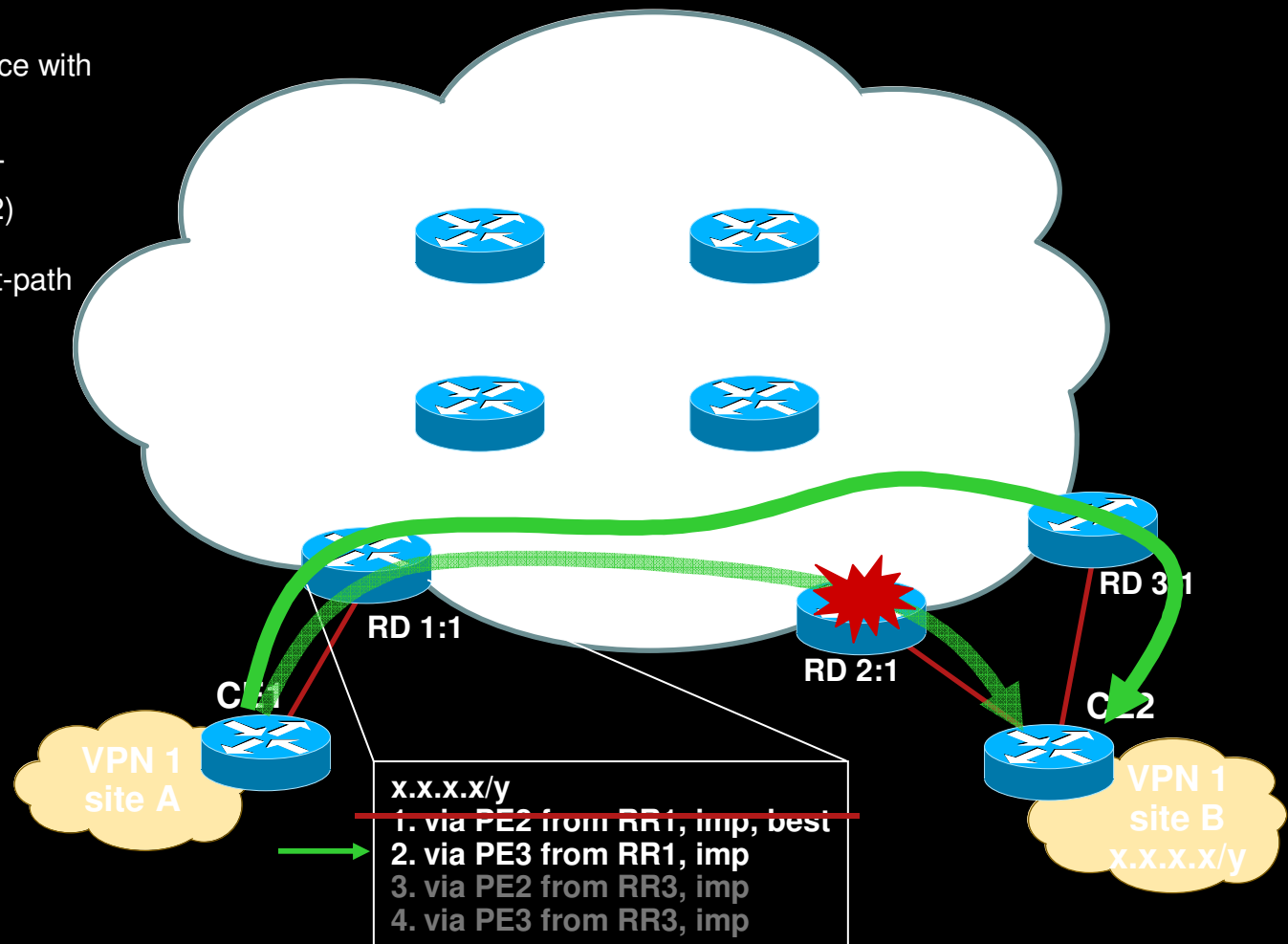|  | |
|---|---|
| **Eliminated with BGP PIC** | S(p): time to scan BGP table |
| |     Per-RD walk for VPNv4 and then IPv4 |
| | B(p): time to compute bestpath for impacted routes and upd FIB |
| | Wtx(p): time to generate/propagate all withdraws |
| **Eliminated With BGP Local Convergence and BGP PIC** | RR(p): time for the RR reflection |
| | Wrx(p): time to receive and process all withdraws |
| | B(p): time to compute bestpath for impacted routes and upd FIB |
| | Where X(p) means that this component scales with the table size |

# PE Node Failure

1. IGP neighbors detect the loss of PE2, re-originate and flood their LSP

2. PE1 completes IGP convergence with delete(PE2)

3. On PE1, RIB notifies BGP NHT of the loss of a BGP nhop (PE2)

4. On PE1, BGP recomputes best-path and selects an alternate nhop

RD 3:1

RD 1:1

RD 2:1

CE1

CE2

VPN 1 site A

VPN 1 site B
x.x.x.x/y

x.x.x.x/y
1. via PE2 from RR1, imp, best
2. via PE3 from RR1, imp
3. via PE2 from RR3, imp
4. via PE3 from RR3, imp

# PE Node Failure

- Convergence depends on

    D: time to detect failure

    IGP: time to complete a simple IGP convergence

    simple: leaf node deletion leads to maximum iSPFgain and very few prefixes deletions

**Eliminated with BGP PIC**

S(p): time to scan BGP table

Full VPNv4 walk(*) and then IPv4 walk

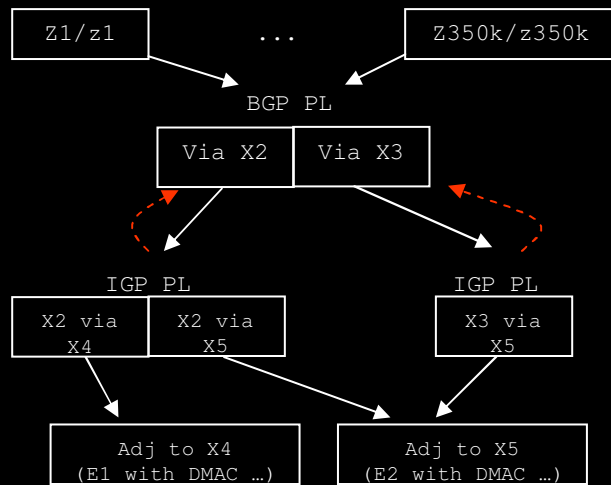B(p): time to compute bestpath for impacted routes and update FIB

(*) CSCsm80316: 0ms for default VPNv4 NHT timer
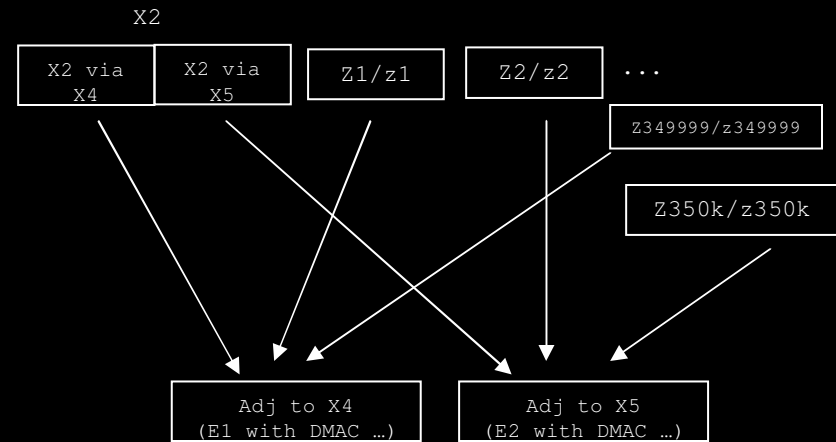
(*) Future: RD-based scoped walk upon NHT trigger for iBGP nhop invalidation

# Load-Balancing Efficiency
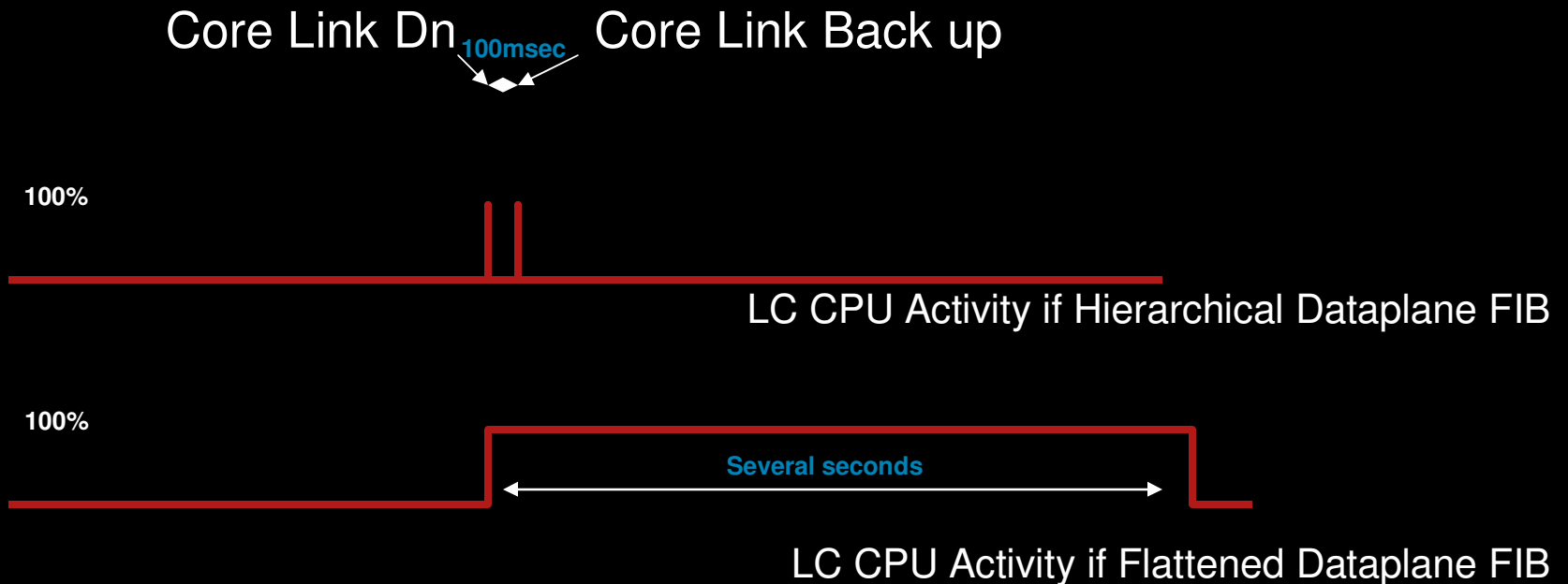
## Hierarchical Dataplane Fib

```
┌─────────┐              ┌──────────────┐
│ Z1/z1   │     ...      │ Z350k/z350k  │
└─────────┘              └──────────────┘
        ↘            ↙
          BGP PL
        ┌─────────┬─────────┐
        │ Via X2  │ Via X3  │
        └─────────┴─────────┘

          IGP PL                IGP PL
    ┌─────────┬─────────┐   ┌─────────┐
    │ X2 via  │ X2 via  │   │ X3 via  │
    │ X4      │ X5      │   │ X5      │
    └─────────┴─────────┘   └─────────┘

  ┌────────────────────┐   ┌────────────────────┐
  │    Adj to X4       │   │    Adj to X5       │
  │ (E1 with DMAC …)   │   │ (E2 with DMAC …)   │
  └────────────────────┘   └────────────────────┘
```

## Flattened Dataplane Fib

```
  X2
┌─────────┬─────────┐   ┌─────────┐   ┌─────────┐  ...
│ X2 via  │ X2 via  │   │ Z1/z1   │   │ Z2/z2   │
│ X4      │ X5      │   └─────────┘   └─────────┘
└─────────┴─────────┘
                                  ┌──────────────────┐
                                  │ Z349999/z349999  │
                                  └──────────────────┘

                                  ┌──────────────┐
                                  │ Z350k/z350k  │
                                  └──────────────┘

  ┌────────────────────┐   ┌────────────────────┐
  │    Adj to X4       │   │    Adj to X5       │
  │ (E1 with DMAC …)   │   │ (E2 with DMAC …)   │
  └────────────────────┘   └────────────────────┘
```

---

### Optimum

Any flow can be load-balanced on
any BGP path of the BGP PL and
any path of the IGP PL

# Robustness

Core Link Dn **100msec** Core Link Back up

100%

LC CPU Activity if Hierarchical Dataplane FIB

100%

**Several seconds**
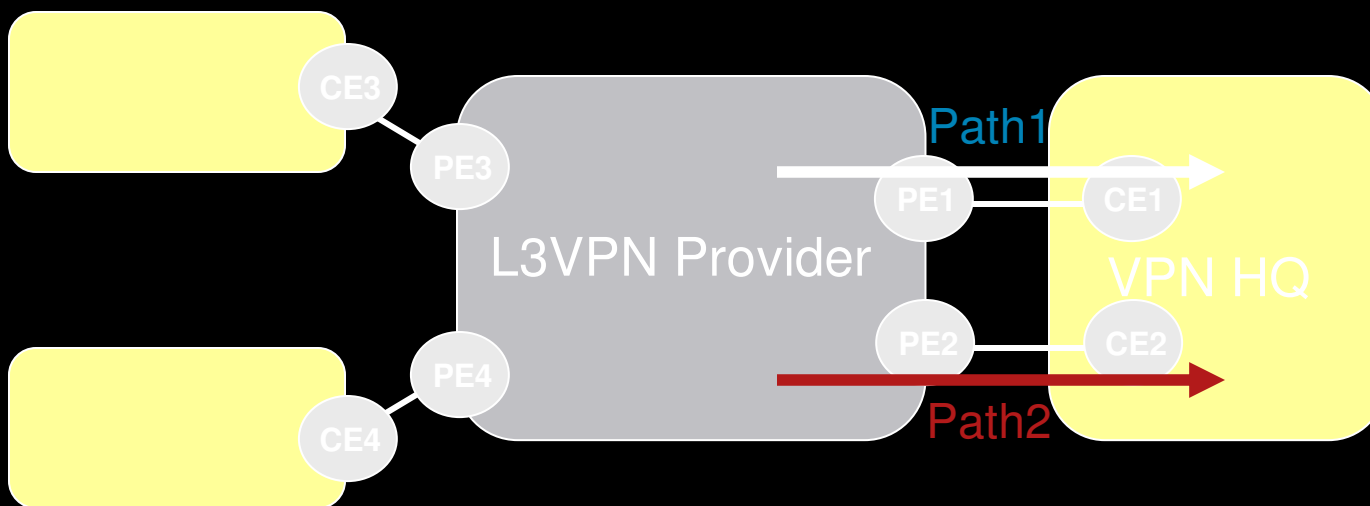
LC CPU Activity if Flattened Dataplane FIB

- At 10usec per update of a flattened entry, a 250k Internet table takes 2.5 seconds

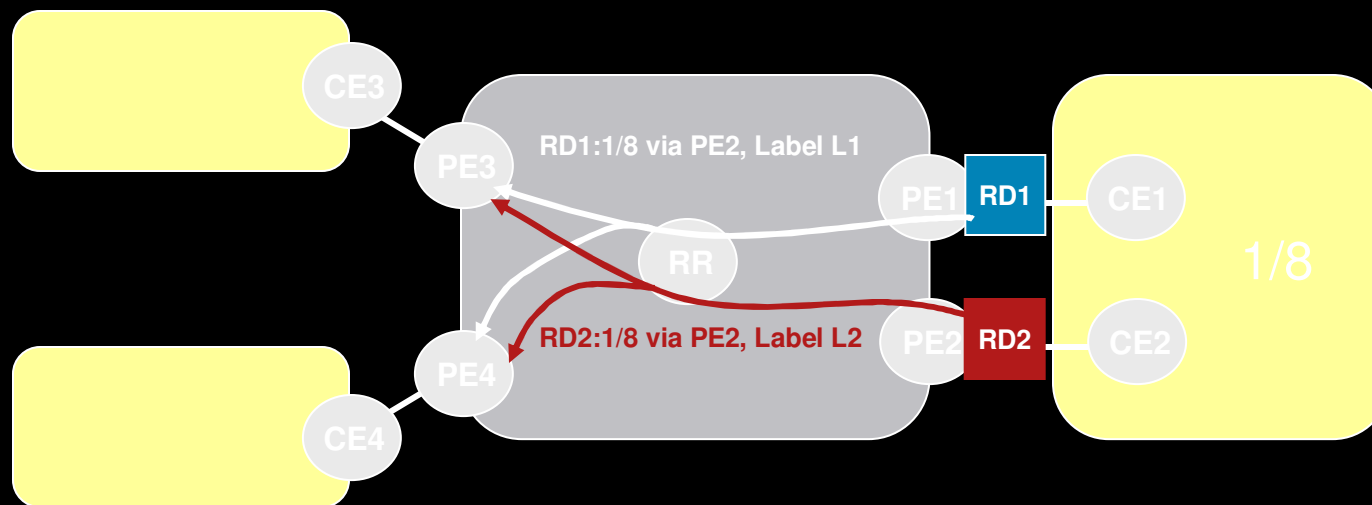# Providing BGP Path Diveristy

# L3VPN

- Important sites are dual homed



Spoke sites are numerous, less important and hence economically, a decision may be taken to single home them

HQ sites are rare, very important and hence are always multi-homed

# L3VPN

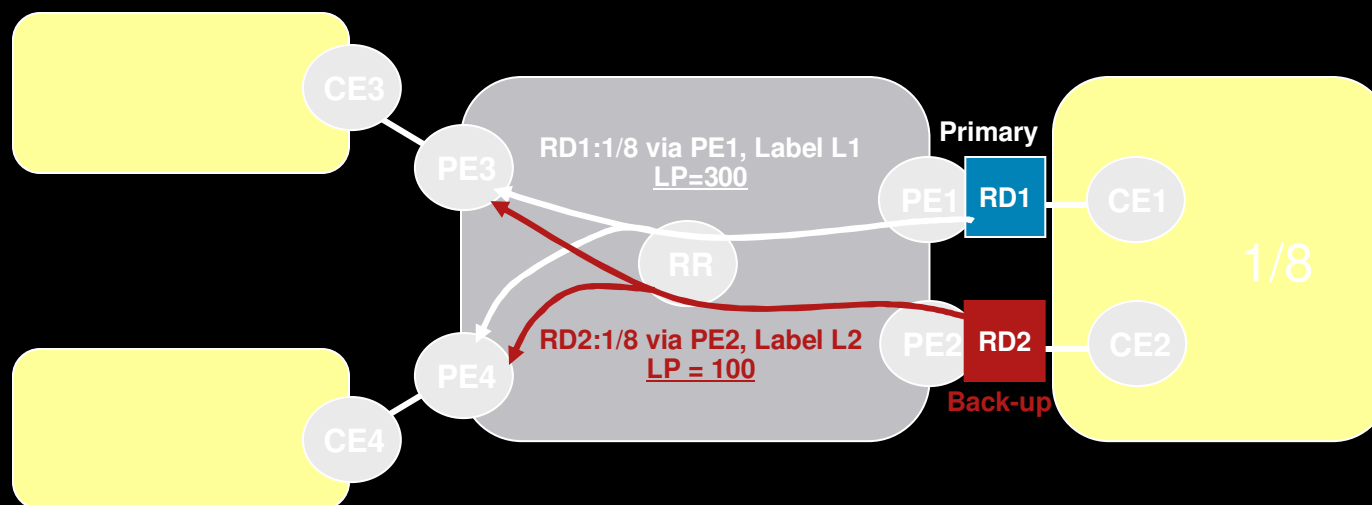- Unique RD allocation ensures both paths are learned, even through route reflectors

# L3VPN with Primary/Backup Policy
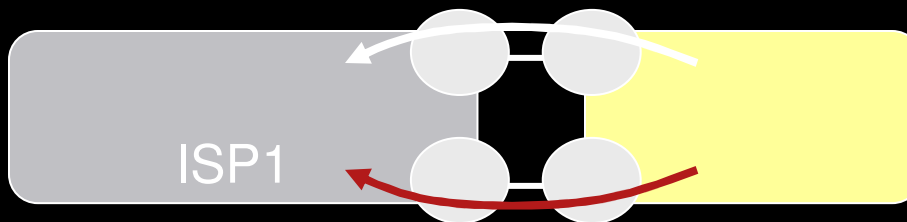
- Best-external Advertisement

   a PE always announces its best external path even if he himself selects an iBGP path

   the label bound to the advertised best-external path is installed in PE2's TFIB to enable remote PE's (PE1, PE3, PE4) to send traffic via that path without requiring PE2 to enable that path as overall best



RD1:1/8 via PE1, Label L1
LP=300

RD2:1/8 via PE2, Label L2
LP = 100

Primary

Back-up

CE3
PE3
PE4
CE4
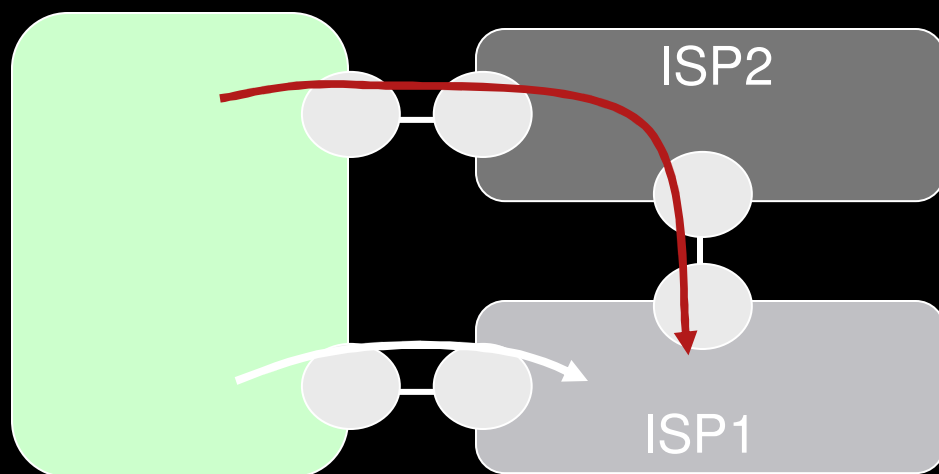RR
PE1
RD1
CE1
PE2
RD2
CE2
1/8

# Internet

- Disjoint paths to routes caring about Tight-SLA do exist at the AS boundary

- Case 1: an entity caring about its Internet connectivity dual-homes to the same provider
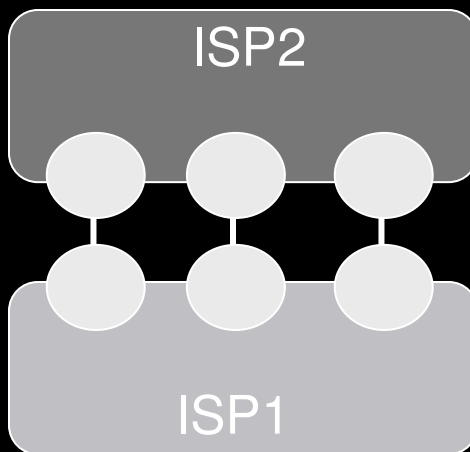
ISP1

# Internet

- Disjoint paths to routes caring about Tight-SLA do exist at the AS boundary

- Case 2: an entity caring about its Internet connectivity single-homes to two ISP's. These two ISP's vie for the same market segment in the same geography. These two SP's are thus guaranteed to peer

# Internet

- Disjoint paths to routes caring about Tight-SLA do exist at the AS boundary

- Case 3: aside for pure resilience requirements, the current load and the future growth leads to multiple links between Peers and Transit Suppliers
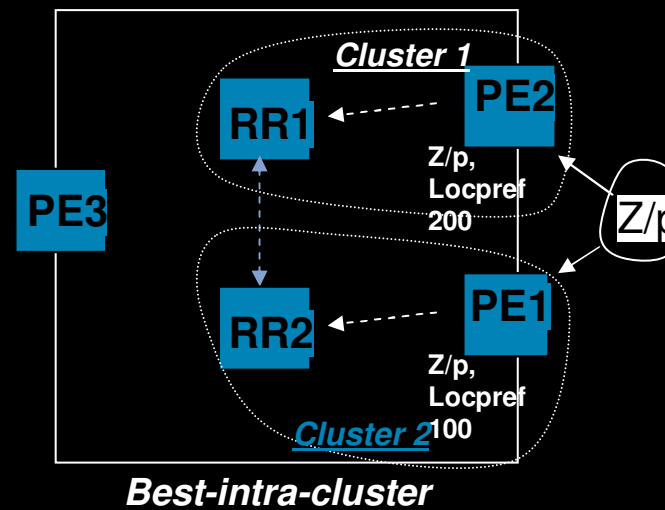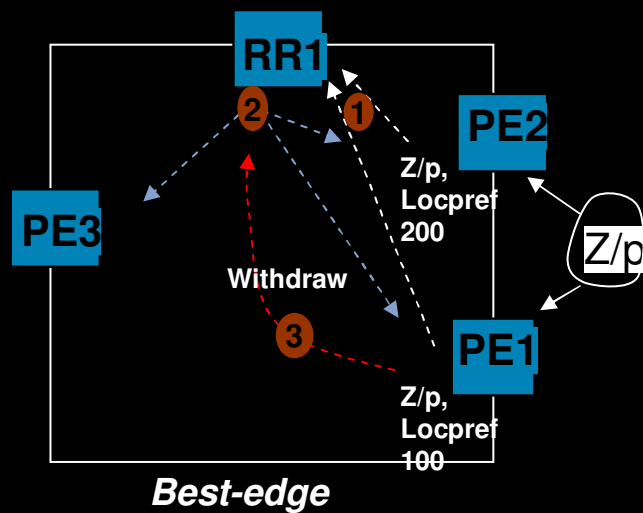
# Problem: Data hiding

- Path reduction at two places:
    - Less preferred border (AS or confed) routers don't announce their paths to iBGP
    - RRs (or confed-ebgp peers) hide all but the best path

- Thus ingress routers most often know about one exit point only

- When that exit point fails, traffic loss scales with control plane convergence
    - PIC can't get triggered

- Not knowing about more exit points also means the ingress routers can't do load balancing

- Not having path diversity has other issues as well:
    - Route oscillation: a protocol bug
    - Unneeded Churn (route hunting)

# Goal

- To improve path diversity in BGP topologies

    Assumption: multiple paths to the same prefix are generally available at the edge of the network

- Application

    Fast convergence

    Load balancing

    Less churn

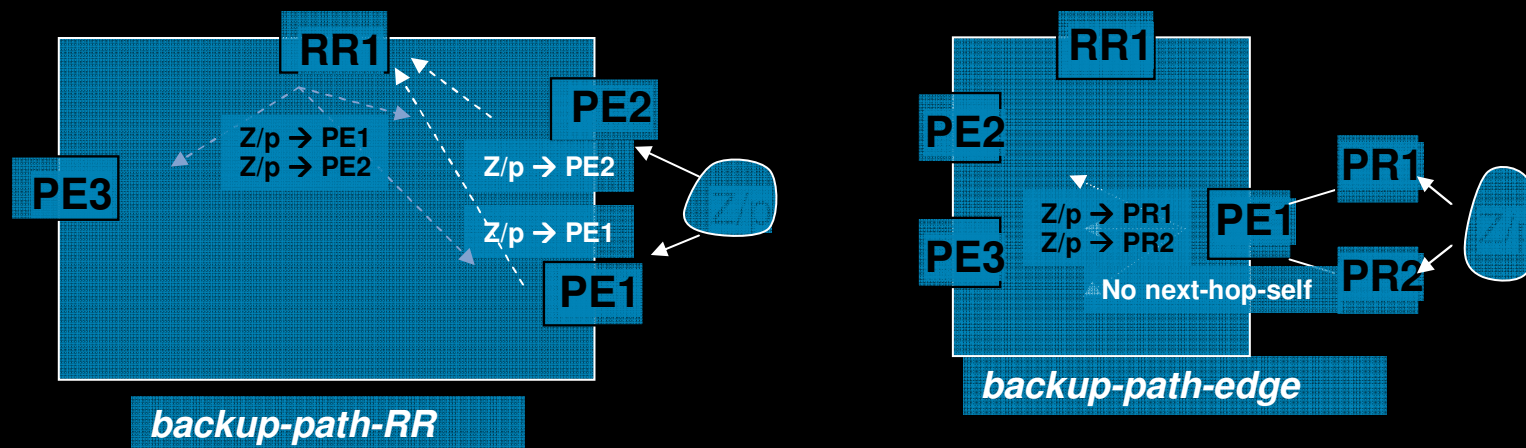    Eliminating route oscillation

    Route server

# Best-External

Less preferred border routers should announce their "own" path instead of withdrawing it



**Best-edge**

**Best-intra-cluster**

## Best-external

# Add-Path

Aggregators (RRs, confed border routers) should advertise backup paths



RR1
PE2
PE3
Z/p → PE1
Z/p → PE2
Z/p → PE2
Z/p → PE1
PE1

**backup-path-RR**

RR1
PE2
PE3
Z/p → PR1
Z/p → PR2
PE1
PR1
PR2
**No next-hop-self**

**backup-path-edge**

## Additional-path

# Application

- Best-External

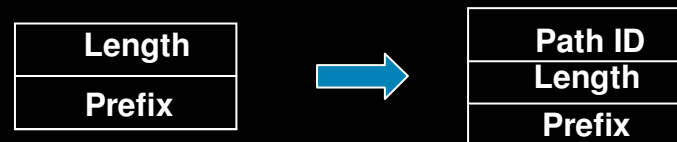  Internet

  VPN (Primary-Backup policy)

- Add-Path

  Internet

  IAS-B and IAS-C-bgp

  Note: not for classic VPN, as they are typically designed with unique RD.

# Add-path draft overview

- Extend NLRI format to include path-ID (so that multiple paths for the same prefix can be advertised).

| Length |
|--------|
| Prefix |

→

| Path ID |
|---------|
| Length  |
| Prefix  |

- Path-ID is application specific, but mostly an opaque ID that is pair-wise

    $id_1:z/p \neq id_2:z/p$

- Capability negotiation for add-path support per *[AFI, SAFI]*

# Implementation: what does it change?

- What paths to advertise? (when we don't want to advertise all)

    Selecting backup paths / second-best

- Update generation: per-path vs. per-prefix

    Adj-RIB-Out is per-prefix today since only best path is sent

    Maintain "send state" per "path to be announced"

- Update reception

    Control plane: process multiple instances of prefix, select second-best

    Data plane: to be able to install current best and second-best in forwarding for PIC

Pranav Dharwadkar

# Add-path: selecting second-best

## Simple rule

**1** Select best

**2** Remove all paths whose originator ID == best's (including best)

**3** Run bestpath selection again on the remaining paths to select backup

# Cost

- Memory overhead

  By how much?

- CPU cycle increase for update processing

  Update reception at edge routers increases proportional to #additional paths

  Update generation at aggregators also increases proportional to #additional paths

- CPU cycle increase for other internal processing as well

  E.g. Next-hop trigger

  Garbage collection

# Design parameters

| | |
|---|---|
| Memory | Shouldn't be prohibitively expensive |
| Convergence | Control-plane convergence shouldn't be degraded greatly; data-plane convergence: no impact |
| Extensibility | Extensible to advertise any #of paths (one to all) |
| Control plane churn | Avoid unnecessary control plane churn (e.g. path ordering changes) |

# Status

- Minor modifications to add-path draft

- Accompanying application drafts

  Fast convergence and Load balancing
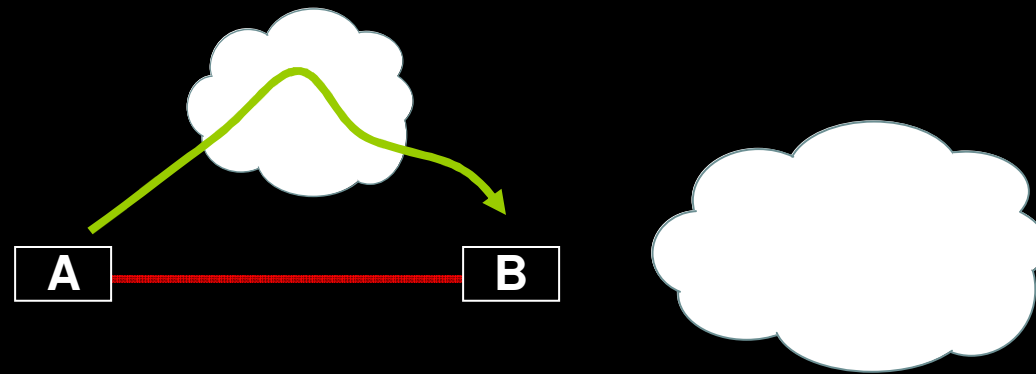
  Route oscillation

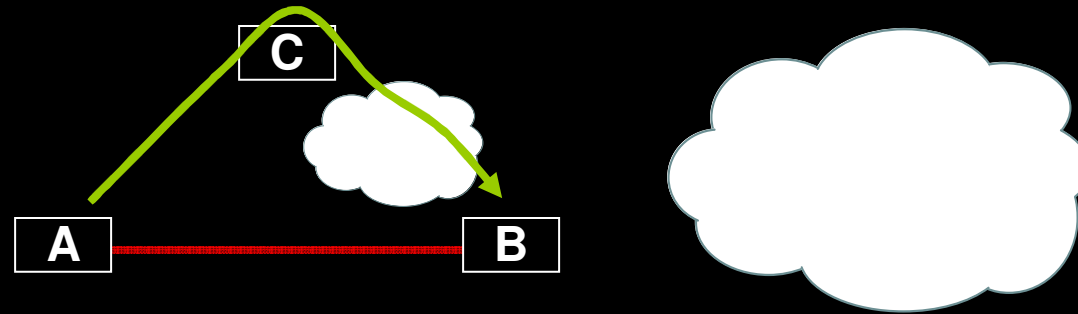- Submit the drafts at the next IETF

# IPFRR

# IPFRR - Protection

- Per-Link LFA (loop-free alternate)

- Per-Link PQ

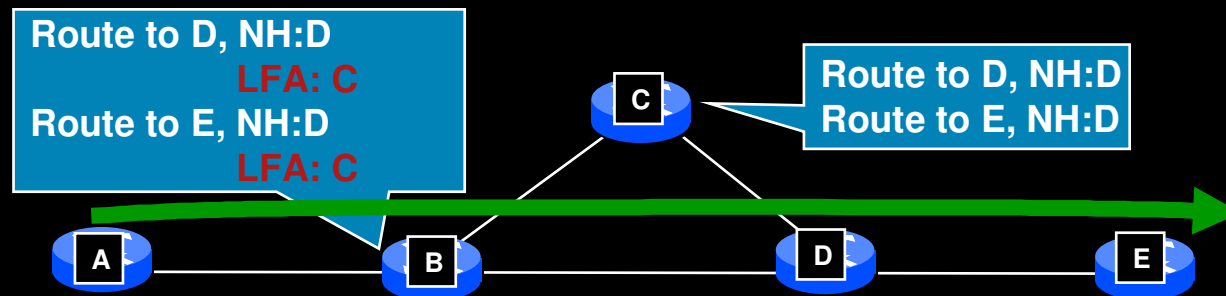- Per-Link Not-Via

# IPFRR – Protection
# Per-Link Property



- If A finds an alternate path to B

- Then this alternate path is valid for any destinations that A normally routes via B
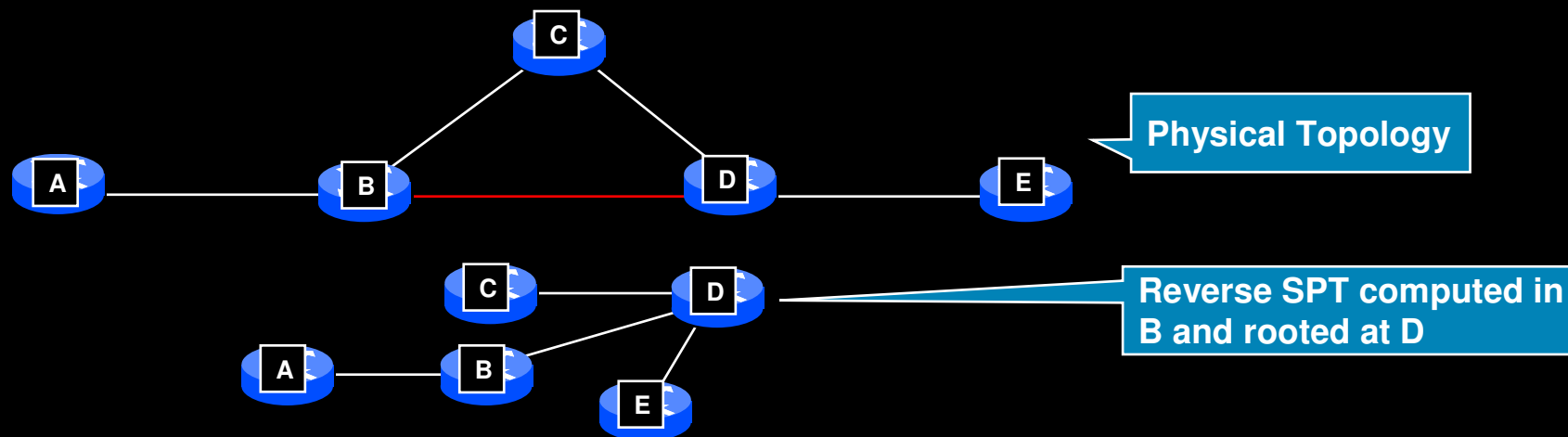
# IPFRR – Protection
## Per-Link LFA



- To protect AB, A may reroute packets via C

    if, what ever link AB status, the IGP route from C to B avoids AB

# IP FRR: Loop Free Alternate routes

**Route to D, NH:D**
         **LFA: C**
**Route to E, NH:D**
         **LFA: C**

**Route to D, NH:D**
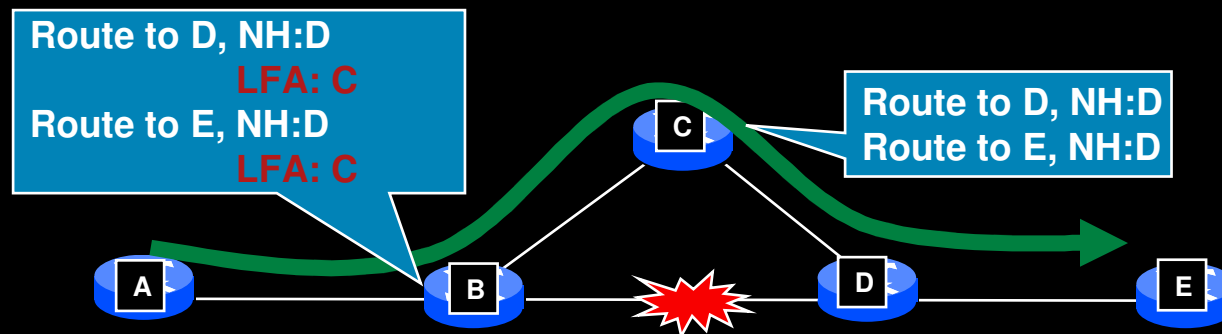**Route to E, NH:D**

C

A        B        D        E

- Used when another neighbor can be safely used as an alternate next-hop for protected traffic

- Upon BD link failure, B can safely reroute to C traffic it used to send to D

    No loop will be formed

    C will forward to D and not back to B

- Pre-computation without any new topology information

    B just leverages its link-state database

# IP FRR: LFA route computation



**Physical Topology**
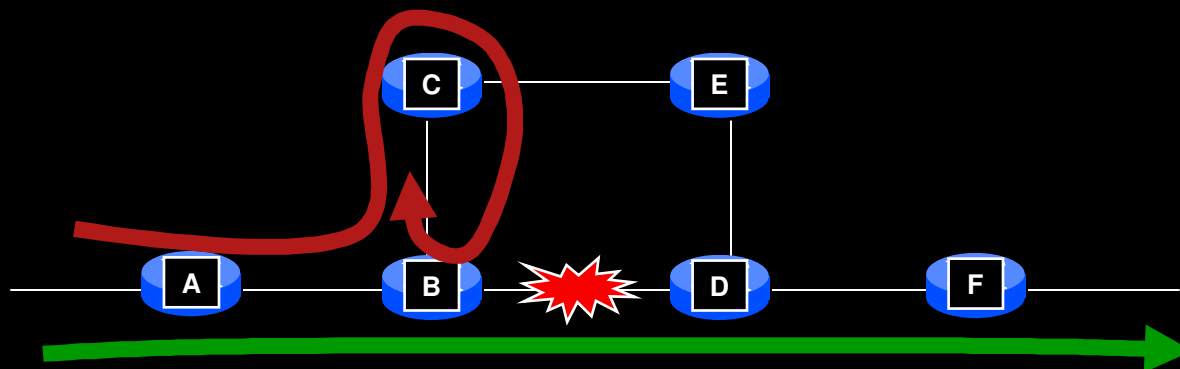
**Reverse SPT computed in B and rooted at D**

- B computes a reverse SPF rooted at D

    Neighbor at the other side of the protected link

- From computing router perspective, a valid LFA is a neighbor that does not belong to the same Sub-Tree (branch)

Pranav Dharwadkar

# IP FRR: Switchover on link failure

**Route to D, NH:D**
           **LFA: C**
**Route to E, NH:D**
           **LFA: C**

**C**

**Route to D, NH:D**
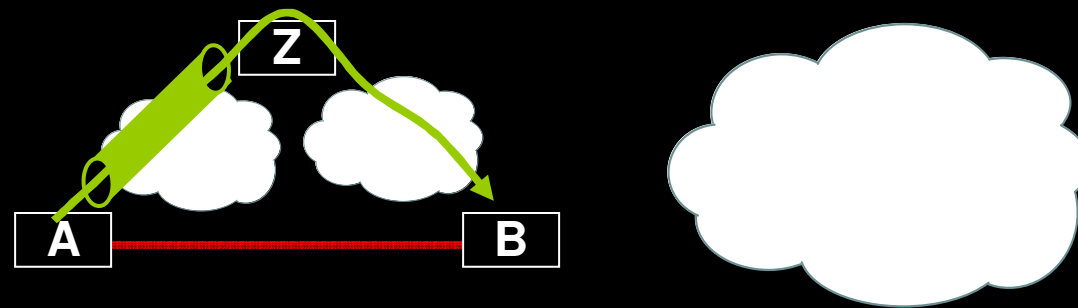**Route to E, NH:D**

**A**   **B**   **D**   **E**

- When link failure is detected, traffic is forwarded according to LFA backup entry

- Local decision in the rerouting node

   No need to signal anything

   No need for any kind of interoperability

- Traffic is rerouted and meanwhile the IGP converges

# Per-Link LFA is topology dependent



- B has no Link-LFA for B-D

    if B would push the traffic to F towards C,
    then C may loop it back to B

- Mitigation

    study across several SP's show 75 to 80% coverage !

    Cariden supports an IPFRR coverage module

# IPFRR – Protection
# Per-Link PQ



- A may <u>encapsulate</u> packets to Z to protect link AB
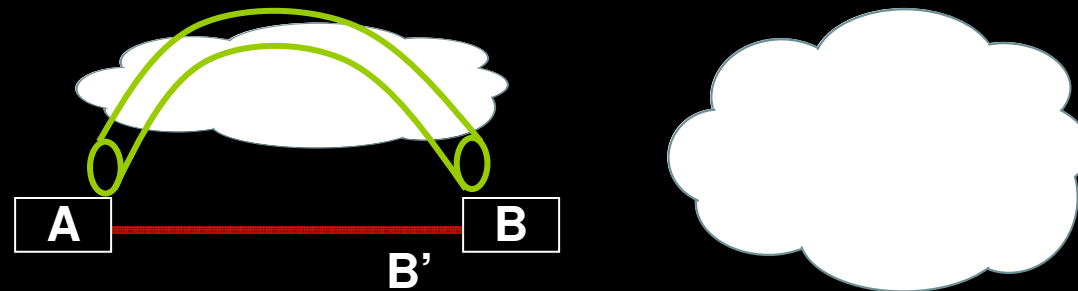
    If, what ever link AB status, the IGP route from A to **Z** avoids AB

    If, what ever link AB status, the IGP route from Z to B avoids AB

Note: only the leg from A to the PQ node is encapsulated

Note: the first condition is called the P condition, the second condition is called the Q condition… hence the name PQ node for a node which meets both conditions.

# IPFRR – Protection
# Per-Link NotVia

- B advertises an additional subnet (B') for adjacency AB

- <u>All nodes in the IGP topology</u> compute a route to B' based on a unique topology for link AB (the normal topology minus link AB)

- A may <u>encapsulate</u> packets to B' to protect link AB

Pranav Dharwadkar

# IPFRR – Protection Benefits

**% of protected links**



- Rule of thumb:
  per-link LFA: ~ 75% coverage
  PQ: ~ 100% coverage

# IPFRR – Protection Complexity (Cost)

| | LFA | PQ | Not-Via |
|---|---|---|---|
| SPF complexity | negligible | negligible | important |
| LSDB increase | no | no | yes |
| Encapsulation | no | yes | yes |
| IETF protocol change | no | no | yes |
| Network-wide migration | no | no | yes |
| Incremental deployment | yes | yes | no |

# IPFRR – Protection
# Kiss Optimum

- **IGP FC** as the base behavior

    anyway required in any FRR scheme

    ~ 200msec for 25% of the link failures

- Enforce **LossLess Maintenance**

- Leverage **Per-Link LFA**

    automated

    incremental deployment

    no protocol change or inter-operability validation

    ~ 50msec for ~75% of the link failures

# Do not re-invent the wheel

- If you really want 50msec upon any link and node failure for any topology: MPLS FRR

    link and node protection

    100% coverage

    deployed

    available since 1999

- Leverage automatic MPLS FRR link/node protection if per-link LFA not available

Under Cisco NDA – Extremely Confidential. Do not distribute

# IPFRR Per-Link LFA and Dual-homed PE's



**Backbone**

- **PE's are usually dual-homed to two interconnected DR's**

- **IPFRR Per-Link LFA is a perfect solution on such PE's**

# IPFRR to scale MPLS FRR



- IPFRR in the PoP

- MPLS FRR between PoP's

- Simpler full-mesh of TE tunnels (scale, inter-area)

Pranav Dharwadkar

# IPFRR and Node Protection

- Core link failures are much more frequent than core node failures. Hence optimize for links.

- If really a 100% coverage is required, just use MPLS FRR which is a mature and deployed technology

# Per-Link LFA support for LDP LSP's



LFA for BD:
**Push 42 and forward to C**

C

Label: 42
Label: 12
Dest: P1

Label: 12
Dest: P1

A

B

D

E

Label: XX
Dest: P1

Dest: P1

Prefix P1

- Upon failure, packets are encapsulated to LFA neighbor
- Using neighbor label as advertised by LFA neighbor
- Per platform label space
- Similar to MPLS-FRR (but without signaling)

# LDP support for Per-Prefix LFA

- Benefits

  improved capacity planning

  packets leverage shortest-paths beyond the first-hop LFA

  good likelihood to protect for node failure

  no label push

- Drawback

  Per-Prefix rewrite

  Mitigation: eng review confirmed commonality of infrastructure with BGP PIC Edge Primary Backup. Commonality of roadmap for the two projects (4.0/4.1).

# IPFRR and KISS

# IPFRR and the principle of Simplicity

- KISS optimizes requirement vs cost/complexity

- KISS Optimum

    IGP Convergence Enhancement

    LossLess Maintenance

    Per-Link LFA

# KISS Optimum

**Gain** (measured in reduction of msecs)

**Complexity**

(always impact OPEX, may impact CAPEX)

**IGP FC**

**Per-Link LFA when possible**

**LossLess Link IN/OUT of service**

**MPLS TE Mesh and MPLS FRR**

**full coverage, explicit engineering, mature, deployed, can be scaled through FC/IPFRR in the PoP's to avoid too large full-mesh**

# IOS-XR IPoDWDM / Routing Integration

# IOX IPoDWDM/Routing Integration

**Working path** | **Switchover lost data** | **Protected path**

**Working path** | **Protected**

## POS on router

**BER**

LOF

## Transponder

**Corrected bits**

FEC limit

**Optical impairments**

## 10GE WDM + G.709 on router

**BER**

Hitless Switch

**Corrected bits**

FEC limit

Protection trigger

**Optical impairments**

**Visibility of router into transmission layer performance allows for superior protection compared to transponder based networks**

# Which Failure Modes are Covered?

- Slow signal degradation due to:

  Aging of fiber plant

  Pinched patch cord

  High PMD – while PMD can be compensated for, this scheme allows to reduce the level of compensation (and sometimes to eliminate it), by covering rare cases where PMD exceeds the systems ability to compensate

- Even fast events may be covered:

  Cable cuts by backhoe take 100's of ms – this has not been tested but cutting a patch cord w scissors seems to take 100ms

  Human error: pulling a connector takes 10's ms – may or may not apply

  EDFA failure modes?

# IOX IPoDWDM/Routing Integration

- Generic implementation supporting

    MPLS FRR

    IP FRR

    FC

# IPoDWDM FRR based on pre-FEC errors

**(d) While "old" packets continue to flow along the original path, new packets are redirected to a protect path hitlessly**

**R3**

**(a) Increasing BER in one direction – no problem in the opposite direction**

**R1**  **R2**  **R4**  **R6**

**(c) Head end triggers a "forced switch" event**

**(b) Pre-FEC error backward defect indicator (BDI) flag raised by tail-end node (in G.709 frame)**

**R5**

Forward Error Correction (FEC)
Bit Error Rate (BER)

# HA in IPTV World

Under Cisco NDA – Extremely Confidential. Do not distribute

# Technical Challenges Introduced by IPTV

- **Network Transport**

    need sufficient capacity – 10GE$\rightarrow$100GE, Multi-TB Routers

    Native IP or MPLS – Different approaches exist

- **Dealing with (Video) Packet Loss**

    Minimize due to topology changes

    Potential for Zero Packet Loss exists

- **Monitoring and Troubleshooting (Video SLA Conformance)**

    measure and ensure IP network can support video flows

    determine where in network problem exists

# Video and Packet Loss

- Losing packets from video flow is not good and user will likely see impairment

- 4 Primary causes for packet loss:

  Excess delay – use QoS to minimize

  Congestion – use QoS and CAC to minimize

  L1/L2 errors – more likely is access, use FEC or retransmission to recover

  Network Reconvergence – impact reduced via FC/FRR but packets will still be lost

- What about 50ms recovery?

  Artifact of SONET/SDH, not really a video requirement

  Due to compression and placement of I-frame in packet, even one packet loss will impact user quality

# MPEG-2 Encoding and Loss Impact

| Sequence Header | GOP Hdr | $I_1$ | B | B | $P_1$ | B | B | $P_2$ | B | B | $P_3$ | B | B | $P_4$ | B | B | GOP Hdr | $I_2$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

*An MPEG-2 ES with 15-Frame Group-of-Pictures (GOP)*

- **Video sequence consists of I Frames separated by B and P frames**
    - **I Frames are anchors**
    - **P frames built from previous I or P frames**
    - **B frames built from preceeding and succeeding I or P frames**
    - **Update rate is ~30 Frames Per Second (FPS)**

- **Duration of impairment depends upon which frames are impacted**
    - **e.g. if $I_1$ lost then need to wait ~500ms for next $I_2$**
    - **e.g. If $P_1$ lost then need to wait until next I-frame**

- **Larger GOP sizes offer better compression but longer visual impairments if I or P are lost early in the sequence**

- **Evans, J., Greengrass, J., Begen, A., "Not All Packets Are Equal: The Impact of Network Packet Loss on Video Transport", to be published in 2008, IEEE Publication**

# MPEG Frame Impact from Packet Loss
## GOP Size: 500 ms (I:P:B = 7:3:1)

**I-Frame Loss Probability (%)** (y-axis: 0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100)

**Outage Duration (ms)** (x-axis: 0, 50, 100, 150, 200, 250, 300, 350, 400, 450, 500, 550, 600)

**50ms Requirement**

- Legacy telco voice (not video) requirement

- Artifact of SONET/SDH

**FACTS**

- At 50 ms outage, an I-frame will be lost with a chance of 34%

# Towards Lossless IPTV Transport:
## *Deployment Scenarios*

**Increasing Cost and Complexity**

Can compare different approaches in terms of
- Lossless or not
- Bandwidth usage in network working and failure cases
- Delay impact
- Cost and complexity of network design and deployment and application infrastructure

**Network Re-engineering**

**No network Re-engineering**

**MPLS TE FRR**

**Increasing Loss**

**Fast Convergence**

**Lossless**

**1 GOP Impacted**

source: John Evans

# Fast Convergence or Fast Reroute



- Network reconverges / reroutes on core network failure (link or node); loss of connectivity is experienced before the video stream connectivity is restored

- Fast Convergence or Fast Reroute

  - ✓ Lowest bandwidth requirements in working and failure case

  - ✓ Lowest solution cost and complexity

  - ! Requires fast converging network to minimize visible impact of loss

  - ✗ Is not hitless – will result in a visible artifact to the end users

source: John Evans

# Real topology – Before failure

# Real topology – Failure 1

# Summary across all available CRS1 data



iox_34-ci_mcast_failure1A_mroutes400_u50k_m40k_isis2500_acl_20061026-185358

Agilent measurements

Legend:
- 0%
- 50%
- 90%
- 100%

x-axis: prefix nr / channel nr

- **MC impact on UC is negligeable**
  First 500 ISIS prefixes < 350ms, lsp-gen 50msec, spf iw =50msec

Pranav Dharwadkar

# Summary across all available CRS1 data



**SSM Convergence as a function of the number of IPTV channels**

Legend:
- max of max
- median of median

Categories: 4000 IPTV channels, 800 IPTV channels, 400 IPTV channels

X-axis (ms): 0, 100, 200, 300, 400, 500, 600, 700, 800, 900, 1000

- And most important… CRS1 SSM Convergence is rather quick!

# MoFRR

- Automated Disjoint Routing: deliver two disjoint branches of the same IPTV PIM SSM tree to the same PE

- <50msec: the PE locally switches to the backup branch upon detecting a failure on the primary branch

    IPTV Inter-packet Gap is 0(1msec). Upon not receiving any packet from the primary branch for 50msec, switch-over to the backup feed

- Hitless: the PE uses the two branches to repair losses and present lossless data to its IGMP neighbors

    Leverage RTP sequences to repair losses

# Multicast only FRR

not wasted bandwidth

alt data path

data path

S

A

R
B

alt
path

join
path

C

7) If upstream of D there are
receivers, bandwidth is only
wasted from that point to D

8) When C fails or DC link fails, D makes
local decision to accept packets from B
9) Eventually unicast routing says B is new
RPF path

D

wasted bandwidth

R

data path

rpf path (RPF join)

alt join (sent on non-rpf)

○ interface in oif-list

✚ link down or RPF-failed packet drop

1) D has ECMP path {BA, CA} to S
2) D sends join on RPF path through C
3) D can send alternate-join on BA path
4) A has 2 oifs leading to a single receiver
5) When RPF path is up, duplicates come to D
6) But D RPF fails on packets from B

Pranav Dharwadkar

# Two-Plane Network Design

- ■ Many SP networks apply the Two-Plane Design

  two symetric backbone planes (blue and red)

  interconnected by grey links with large metrics to ensure that a flow entering the red plane goes all the way to its exit via the red plane

  pop's are dual-homed to each plane

  important content (IPTV source) is dual-homed to both planes



IPTV source

Pop1

Pop2

PopN

# Two-Plane Network Design

- An IPTV SSM Tree for a premium channel is densely covering the two-plane design

- From a capacity planning viewpoint, all Blue and Red routers in a PoP are or must be assumed to be connected to the tree

IPTV source

Pop1

Pop2

PopN

# MoFRR PIM Enhancement

- Send an additional join to an ECMP neighbor to the source

- Simple

    requires no protocol modification

    requires no new inter-operability testing

- Incremental Deployment

# MoFRR Applied to Two-Plane Network Design

- MoFRR only needs to be deployed on PE's (!)

- Does not create any additional capacity demand (!)

- Disjointness does not need to be created by explicit routing techniques. This is a native property of the design (!)

# MoFRR based on IGP trigger

- Upon failure along the primary path, IGP converges and the best path to the source is modified.

- This triggers the use of the already-established MoFRR backup branch

  gain over FC: no time incurred due to the building of the new branch

- Current Behaviour

- Target: sub200msec

# MoFRR 50msec Switch-Over

- IPTV Inter-packet Gap is 0(1msec).

- Monitor SSM (S, G) counter and if no packet received within 50msec switch onto the backup branch

- Feasibility, Scaling and Performance (WIP: Please let us know your interest)

# MoFRR Zero-Loss

- IPTV flows to use RTP

- MoFRR PE device to repair any loss thanks to RTP sequence match on the disjoint branch

- Feasibility, Scaling and Performance (WIP: Please let us know your interest)

# Towards Lossless IPTV Transport:
## *Deployment Scenarios*

**Increasing Cost and Complexity**

**Network Re-engineering**

**No network Re-engineering**

MPLS TE FRR + FEC or TR

MPLS TE FRR

Fast Convergence + FEC or TR

Fast Convergence

**Increasing Loss**

**Lossless**

**1 GOP Impacted**

source: John Evans

# Forward Error Correction (FEC)



- FEC adds redundancy to the transmitted data to allow the receiver to detect and correct errors (within some bound) without the need to resend any data

- Forward Error Correction

  - ✓ Supports hitless recovery from loss due to core network failures if loss can be constrained

  - ✓ No requirement for network path diversity – works for all topologies

  - ! Requires fast converging network to minimize FEC overhead

  - ✗ Higher overall bandwidth consumed in failure case compared to live / live

  - ✗ Incurs delay – longer outages require larger overhead or larger block sizes (more delay)

source: John Evans

# Temporal Redundancy



- With temporal redundancy the transmitted stream is broken into blocks, each block is then sent twice, separated in time

- If block separation period is greater than the loss of connectivity, at least one packet should be received and video stream play-out will be uninterrupted.

- Temporal Redundancy

  ✓ Supports hitless recovery from loss due to core network failures if loss can be constrained

  ✓ No requirement for network path diversity – works for all topologies

  ! Requires fast converging network to minimize block separation period

  ✗ Incurs 100% overhead

  ✗ Incurs delay – longer outages require larger block separation period

source: John Evans

# Towards Lossless IPTV Transport:
## *Deployment Scenarios*

**Increasing Cost
and Complexity**

**Network Re-engineering**

**No network Re-engineering**

TE +
Live / Live

MTR
+ Live / Live

MPLS TE FRR
+ FEC or TR

MPLS TE FRR

MoFRR +
Live / Live

Fast
Convergence +
FEC or TR

MoFRR

Fast
Convergence

**Increasing
Loss**

**Lossless**

**1 GOP Impacted**

source: John Evans

Pranav Dharwadkar

# MPLS Multicast (1)

- ## Drivers for MPLS Multicast

  High-rate single-sources (e.g., video/TV distribution)

  Network optimization (not all traffic on shortest path)

  QoS guarantees, Fast restoration

- ## Two Solutions:

  P2MP TE – one most discussed

  mLDP – extends "PIM-like" receiver-driven model into MPLS

- ## We note that MPLS multicast is in its infancy. IP Multicast is mature, deployed and carrying lots of broadcast video

# MPLS Multicast (2)

**P2MP TE**

head-end

S — R1 — R4 — R2 — D1
           R3 — D2

- Source-driven
  - R1 signals S2L LSP to R2
  - R1 signals S2L LSP to R3
  - R4 merges
- Ideal for single-source multicasts with few leafs
- Applications
  - IPTV distribution/staging
  - Wholesale multipoint transport

**mLDP**

S — R1 — R4 — R2 — D1
           R3 — D2

- Receiver-driven (like PIM)
  - R2 and R3 send mLDP "join" messages towards R1
  - runs over existing LDP/TCP sessions
- Ideal for dynamic, receiver-driven multicasts with many leafs
- Applications
  - Dynamic IP Multicast, mVPN, SDV

# RSVP TE based P2MP
## (draft-ietf-mpls-rsvp-te-p2mp.txt)

**Receivers**

**Receivers**

**Receivers**

**Key Features:**
1. **Bandwidth Reservation**
2. **Bandwidth Optimization**
3. **Explicit Routing**
4. **FRR**
5. **Broadcast TV**
6. **Sender Driven**

**Receivers**

**Video Source**

**Receivers**

**Receivers**

**Receivers**

# MLDP

# MoFRR and MPLS Transport

- MoFRR is as applicable to MLDP as to PIM

Under Cisco NDA – Extremely Confidential. Do not distribute

# Recipe for High Availability - Increase MTBF, Reduce MTTR

## Continuous Systems Operation

99.999+% Service Availability

No Single points of failure for unplanned or planned events

**CSO**

**ISSU In Service Software Upgrade**

**Network High Availability Fast Convergence (PIC)**

**Hitless RP Switchover**

**Non-stop Forwarding & Routing**

**Distributed Route Processor, Service Separation Architecture**

**IOS XR Basics: microkernel based architecture granular process restart, protected memory, highly modular, separation of control, data, management planes, fault tolerance, packaging model**

**Hardware Design: redundancy (fabric, power, thermal, route processor, line card), high MTBF, distributed forwarding, Online Insertion Removal (OIR), parity or error correcting memory, fault insertion testing**

# In Service Software Upgrade (ISSU)

## Goals

- **Reload software while node in production network**
- **Minimize outages to**
  - Control Plane sessions – near-hitless upgrade
  - Data Plane traffic: minimal traffic loss expected

## Implementation

- (D)RP SO + Parallel LC MDR
- Versioning Infrastructure
- ISSU State Machine
- MDR of Fabric Cards
- (Near) Hitless Upgrades
- Rollback Support
  - Active RP and LCs run same version;
  - Standby RP runs older image



Network Level Resiliency to Minimize Impact

Upgrades Performed In-Service to Minimize Impact

Removed from service for upgrade

In-service upgrade

Pranav Dharwadkar

# Reload SW without HW Restart
## Minimal Disruptive Restart (MDR)

### Goals

- **Reload software while hardware continues to function**
- **Minimize outages to Control Plane traffic (I.e. no loss of L2/L3 Sessions) and Data Plane traffic**

### Implementation

- Requires HW Capability to perform MDR
- Avoid reset of ASICs
- Preserve forwarding tables
- Negotiate protocol timeouts with peers prior to MDR (e.g. BFD).
- Micro code download causes only ~50 ms traffic hit on CRS LCs.
- Coordinated handling of HW Data Structure Changes - Hardware Dependent (Future)

**IOX Availability:** Starting in 3.3., Work in progress

**Controlled update**

**CPU**

**SW**

**HW**

**PSE**

**PLU** | **TLU** | **TCAM** | **Queuing**

1-2 Million prefixes | 128K adjacencies | 256K ACL entries | 8K queues

# ISSU – Overall Process Overview
## From RP to Line Cards…



**Legend**

- Old Active
- Old Standby
- New Active
- New Standby

This command forces the switchover of RP and MDR on other Nodes

This command initiates the Rollback Process

# Specific ISSU Technologies

- **ISSU is built upon HA Infrastructure**

  - <span style="color:orange">**It utilizes process restart where possible (example patches / SMU)**</span>

  - **Graceful / Stateful Switchover for Redundant Elements e.g. RP, DRP**

  - **Minimum Disruptive Restart (MDR) for non-redundant elements and Switching Fabric**

- **ISSU Requires some way to protect Forwarding Plane while Control Plane is being upgraded**

  - **NSF using Graceful Restart Extensions where possible**

  - **NSR requiring no Protocol Extensions for all cases**

- **ISSU Requires some way to protect Control Plane while Forwarding plane is being upgraded, options include**

  - **Line card with no Forwarding Impacting Changes**

  - **Line card MDR with Forwarding Impacting Changes**

  - **Line card reload**

# High Availability Tools

© 2007 Cisco Systems, Inc. All rights reserved.

# Highly Scalable Detection of Failures:
## Distributed Bidirectional Forwarding Detection (BFD)

## Goal

**Increase number of concurrently active BFD sessions to improve fault detection efficiency**

## Implementation

Applies to both Physical and Logical interfaces

**Scale per LC:**

100 sessions at 15ms timer, 2 retries

Max 1000 Sessions at 150ms timer, 2 retries

**For Bundled Interfaces:**

1000 sessions at 250ms

**Maximum Scale per system:**

\# LC * # Sessions,

CPU Utilization <=13% (Target)

## IOX Availability

**Introduced in 3.3**

## IOS Availability

**Roadmap item in IOS**

Routing Protocols

BFD Server

RP

LC

BFD Agent

Fast Protect

FRR

Regular failure notification

Fast failure notification

# Drill to root cause without Line Card outage
## LC Root Cause Analysis

### Goal

**LC troubleshooting options based on type of failure**

### Implementation

**Graceful re-routing of Data Traffic upon LC failure**

- BFD process on LC notifies peer "LC is going down"
- Shelfmgr notifies routing protocols

**Maintenance mode for debugging**

- Special mode to allow live debugging on LC
- Minimal set of processes run on LC
- CLI available to enable/disable Maintenance mode

### IOX Availability: 3.4

# Need for Video Monitoring (Vidmon)

- IP/MPLS network:

    delivers broadcast TV streams to many users
    delivers VoD streams to individual paying users

- We know that video is extremely sensitive to packet loss

    chucked packets will have deleterious effect on QoE (aka meltdown)

- Need techniques to isolate where and what video packets are being dropped in the network

"I expect my routers and switches to work, what I need Cisco to tell me is how I am going to troubleshoot video problems in this new network and avoid the need to constantly rollout trucks."

"Big IPTV/VoD Provider"

# Monitoring a Video Transport Service



- Transport service provider must ensure proper service from source boundary to receiver or access service boundaries

source: Clarence Filsfils

# Active Monitoring: IPTV SLA



- IPSLA probes are sent from the boundary source routers to some/all boundary edges. Probes are sent with the Video DSCP.

  End-to-end loss measurement

  End-to-end latency/jitter measurement

  Statistical measurement

source: Clarence Filsfils

# Passive per flow video transport monitoring: Router Vidmon

- Embedded router technique (MDI) to isolate where and what video packets are being dropped in the network on a per-flow basis

- Discriminates between problems at the source boundary, at the edge boundary, within the network

- Complements IPSLA functionality

- Focuses on loss monitoring

- Scales to 100's of flows

- Leverages MDI, a well-known industry metric: RFC4445

source: Clarence Filsfils

# Router Vidmon: Detecting Core Router Problem



Video Source

MDI(CTV)
0:0

MDI(CTV)
0:0

MDI(CTV)
0:0

MDI(CTV)
0:0

MDI(CTV)
0:0

MDI(CTV)
-50:0.1

Access Service and Receivers

Access Service and Receivers

Access Service and Receivers

**MDI (CiscoTV, CRS1) = DF: MLR**

- Delay Factor (DF): a measurement at router CRS1 of the accumulated jitter for flow "CiscoTV"
- Media Loss Rate (MLR): a measurement at router CRS1 of the accumulated loss for flow "CiscoTV"

- Passive per flow monitoring complements IPSLA by detecting an individual per-flow issues at identified core routers

- Pervasive router support and deployment allows for troubleshooting the root cause location

- Significantly reduced CAPEX and OPEX compared to external probes

# Router Vidmon: Detecting Source Problem



**Video Source**

**MDI**
**MDI(CTV)**
**-70:3**

**MDI(CTV)**
**-70:3**

**MDI**
**MDI(CTV)**
**-70:3**

**MDI(CTV)**
**-70:3**

**MDI**
**MDI(CTV)**
**-70:3**

**MDI**
**MDI(CTV)**
**-70:3**

**MDI**
**MDI(CTV)**
**-70:3**

Access Service and Receivers

Access Service and Receivers

Access Service and Receivers

**MDI (CiscoTV, CRS1) = DF: MLR**

- Delay Factor (DF): a measurement at router CRS1 of the accumulated jitter for flow "CiscoTV"
- Media Loss Rate (MLR): a measurement at router CRS1 of the accumulated loss for flow "CiscoTV"
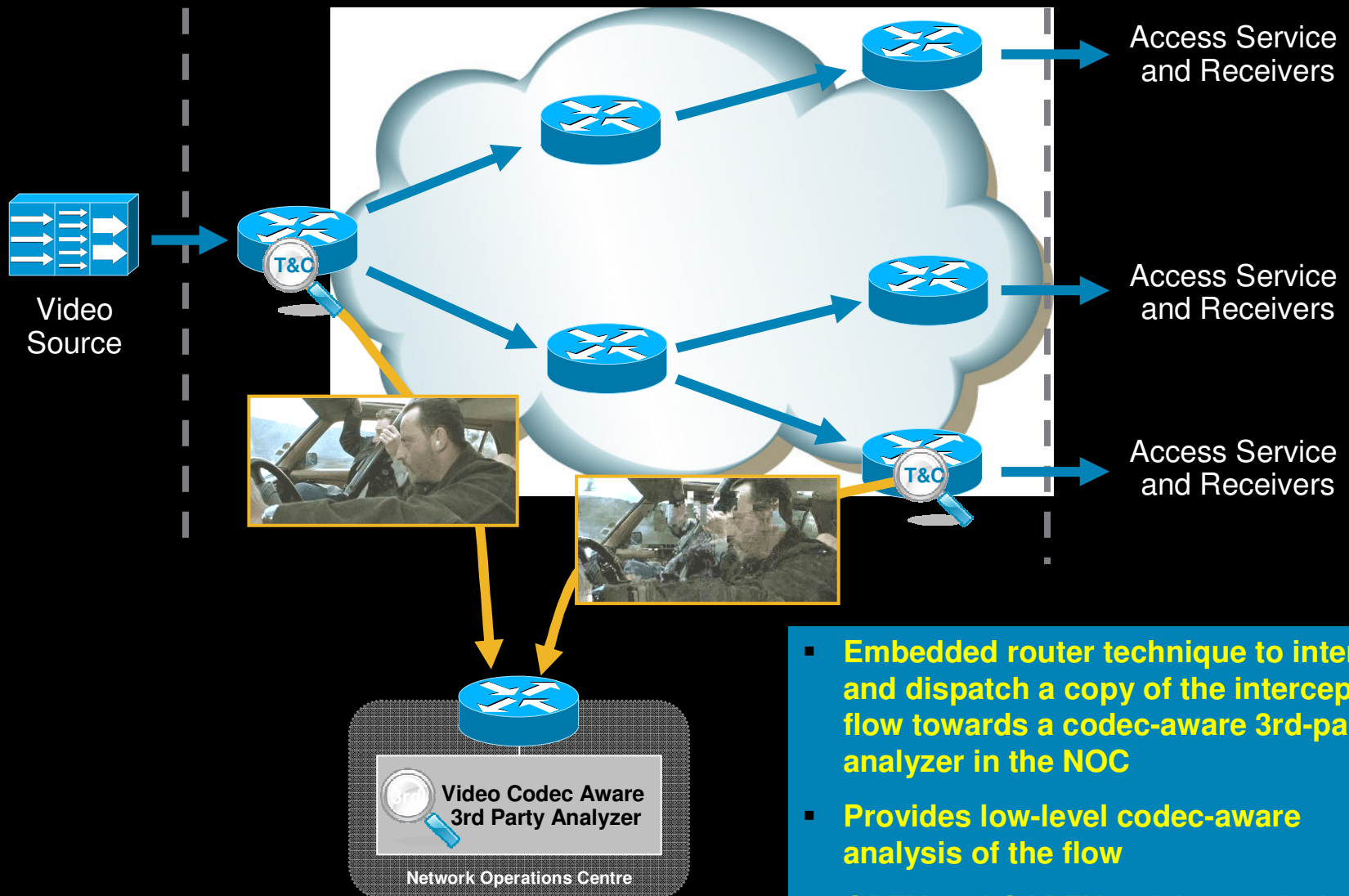
- Complements IPSLA by detecting a problem with the source

source: Clarence Filsfils

# Router Vidmon: Trap and clone



Video Source

Access Service and Receivers

Access Service and Receivers

Access Service and Receivers

T&C

T&C

Video Codec Aware
3rd Party Analyzer

**Network Operations Centre**

- **Embedded router technique to intercept and dispatch a copy of the intercepted flow towards a codec-aware 3rd-party analyzer in the NOC**

- **Provides low-level codec-aware analysis of the flow**

- **OPEX and CAPEX efficient technique for fine-grained video-aware analysis**

# Easier, Faster Troubleshooting
## Device D&I: Current work areas

**Debug-ability enhancements:**

Syslogd: to minimize dependencies.

Dumping (full memory, faster dump).

Process tracing.

Implement On Board Failure Logging (OBFL).

Implement Generic On Line Diagnostics (GOLD).

Additional monitoring of resource usage, IPC latencies, CPU utilization

**Fault Management:**

Develop Service Availability Infrastructure and Fault Correlator with visibility to all HW and SW error input sources.

**Fault Monitoring:**

Add kernel shell access to standby RP, enhanced CLI.

Develop Device Health Monitor

# Enhancements to dumping, logging, tracing
## Generic On Line Diagnostics

*Cisco GOLD Framework*

**Goal**

- **Provide early warning of issues with key HW infrastructure**

**Implementation**

- Detect a discrete set of issues within Switch-fabric and Control-Ethernet
    - Does not look at data forwarding or operations
    - Provides notification of failure

- No impact to router operation – devices remain in-service

- Recommended all customers enable diag .pie and configure online tests as 'best practice' from 3.3 onwards

**IOS XR Availability:** Introduced in 3.3, enhanced in 3.4, 3.5 & 3.6

# Enhancements to dumping, logging, tracing
## On Board Failure Logging

**Goal**

**Save crash info onto NVRAM for troubleshooting**

**Implementation**

Information saved includes

- Boot time, run time
- Boot temperatures and voltages
- On Board memory errors and ASIC errors
- Field diagnostic results
- Crash logs and dumps
- Syslog Events

Think of this as crashinfo for XR, but stored on the failed device itself. Provides a 'failure history' for the board

**IOX Availability:** Introduced in 3.4

# Enhancements to dumping, logging, tracing
## Fault Manager

## Goal

**Device-based Intelligent and programmable Event Collection and Correlation**
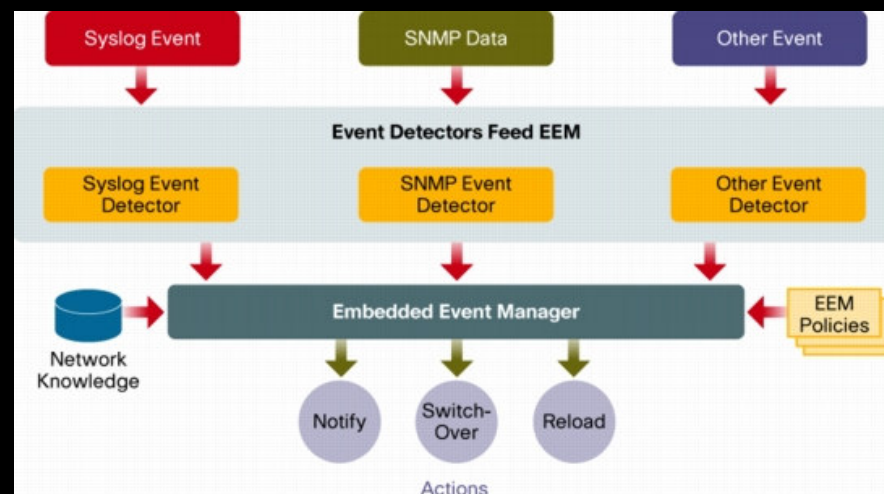
## Implementation

Fault Manager (FM) is a TCL engine inside IOS XR

Runs TCL scripts provided by network operators

interact with router via CLI

Events triggered by a variety of methods

Often used for Periodic Health Checks & Debugging

**IOX Availability:** Day 1
Design with enhancements in 3.6. and future

# Enhancements to dumping, logging, tracing
## Service Availability Infrastructure (Foundation of Health Monitor)

## Goal:

**Monitor health of internal system services for better troubleshooting**

## Implementation:

Standard way to manage outages of GSP, QNET, SYSDB and related services

Will support policy definition

## IOS XR Availability:
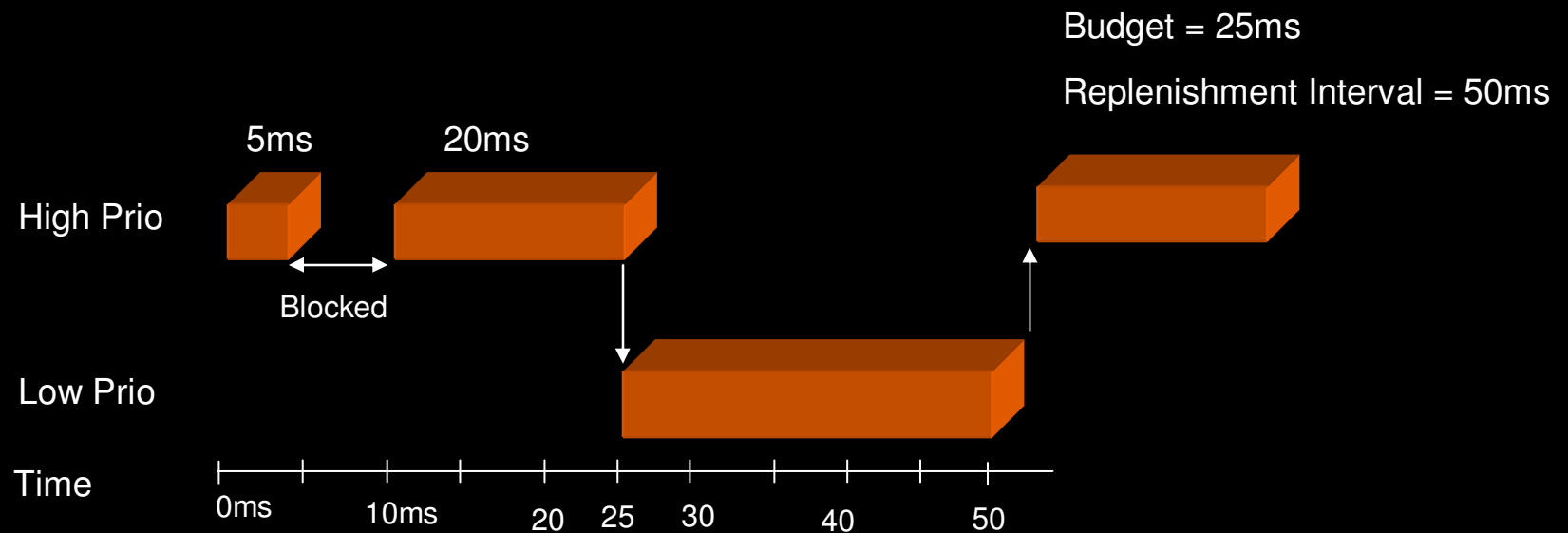
Infrastructure in 3.6.

Producers (major infrastructure services) and Consumers (applications): future

Long term, will integrate with fault correlation and Embedded Fault Manager for better fault handling

# Sporadic Scheduling

- Sporadic Scheduled processes are allowed to run at a high priority for a certain interval ("Budget")

- Once budget is exhausted, they will be dropped to a lower priority for a certain interval ("Replenishment Interval")

- Once budget is replenished, they will be bumped back up to the higher priority

- With the right Budget and Replenishment period, a busy process gives up enough CPU to let others run occasionally

Budget = 25ms

Replenishment Interval = 50ms

5ms          20ms

High Prio

Blocked

Low Prio

Time    0ms      10ms      20  25  30      40      50

# HA – Network & Operations Factor

# Typical POP Architecture

**Edge**　　　　**Distribution**　　　　**Core**　　　　**Peering**



High Speed Provider Edge Routers

Low Speed Edge Routers

Dial-up / ADSL / MetroE

Stm64

GE/FE

10GE

10GE

FE / GE

Other POPS

CRS-1 CORE

Other POPS

Class 4 and Class 5 VoIP services

Voice PE

Stm64

SBC Peering

Internet Peering

Internet Provider VPN

# Potential Problems and possible Solutions

| Potential Failures | Possible Solutions |
|---|---|
| Network Operations Failure | 1. Improved FCAPS Procedures<br>2. Improved Monitoring and Troubleshooting tools |
| Edge Router Hardware or Software Failure Potential Single Point of Failure | 1. Use NSR or Graceful Restart to protect against control plane outages - helps increase MTBF and improves MTTR when redundant path does not exist<br>2. Dual homing Access Layer into different Edge Devices - redundant path in network<br>3. Reduce chances of router failure (See system Factor) |
| Single Link Failure between Core and Edge Router | Recovered using link protection:<br>1. ECMP Path through another core router<br>2. Using Bundled links between core and edge router<br>3. Using TE/FRR |
| Core Router Failure | Recovered using Node protection:<br>1. Using NSR or Graceful Restart to avoid network wide reconvergence or recovery - Improves MTTR and increases MTBF<br>2. Pick a path through the other redundant core router in case of ECMP<br>3. Use IGP Fast Convergence to recover<br>4. Reduce chances of Router Failure (See System Factor) |

# Fault Detection and Recovery

| Potential Problems | Possible Solutions |
|---|---|
| Process Failures | 1. Run-time Modularity that minimizes collateral damage<br>2. Memory Protection between processes<br>3. Fault tolerance and recovery<br>4. Perform Fault insertion Testing to find problems earlier |
| Kernel Failures | 1. Reduce probability of Kernel Failures<br>2. Utilize hardware redundancy<br>3. Perform Fault Insertion Testing to find problems earlier |
| Process Hangs | 1. Detect and recover from process Hangs<br>2. Utilize hardware redundancy if necessary |
| Out or Resource condition | 1. Provide tools to proactively detect out of resource conditions<br>2. Detect and recover from out of resource condition |

Pranav Dharwadkar

# Network Level High Availability

| Failure Type | XR Technologies to use |
|---|---|
| Link Failure | LOS based IGP and FRR trigger<br>BGP Next Hop Tracking – helps BGP convergence due to IGP changes<br>BFD based trigger for IGP, BGP and FRR<br>Traffic Engineering & Fast Reroute<br>IGP and BGP Fast Convergence<br>LDP Session Protection<br>Link Bundling |
| Link Recovery | Interface Dampening<br>Link Verification<br>Fast Convergence<br>LDP-IGP Synchronization |
| Neighbor Node Failure | BFD based trigger for IGP, BGP and FRR<br>FRR Node Protection<br>Reroute using ECMP<br>Fast convergence |
| Network Congestion | Support for robust QOS capability<br>Priority Queue Support<br>Minimum Bandwidth Support<br>Policing (Ingress and Egress)<br>Shaping (Ingress and Egress)<br>Marking (Ingress and Egress)<br>WRED |

# Summary

- Continuous Service Operation remains our highest priority goal

- Challenges are significant given the scale and complexity of networks we are building today

- Substantial progress has been made on the CRS1 and IOX platforms. We are continuing with dedicated focus on this.

- We continue to seek your advice for improving platform and system availability.

  Please continue to tell us what you think!

# Questions?

# HA Summary

| Feature | IOS XR | Competition (TTM Solution) |
|---|---|---|
| Application Modularity | Day 1 Design | Not Modular enough |
| Infrastructure Modularity | Day 1 Design | No |
| Platform SW Modularity (Device Drivers, Microcode) | Day 1 Design | No |
| Online Diagnostics and Troubleshooting | R 3.3. | ? |
| Kernel Enhancements (Sporadic Scheduler support) | R3.4. | No |
| NSR | Today: ISIS, MPLS-TE, Multicast<br><br>CY2007: OSPFv2, LDP<br><br>CY2008: BGP | Juniper Just Announced<br><br>Alcatel Claims Support<br><br>Scale and performance impact not known<br><br>Any known deployment experience? |
| MDR | Incremental Delivery starting in R3.3.<br><br>Some future uncommitted enhancements needed for full MDR | No |
| ISSU | Patch / SMU Support: Most components<br><br>Maintenance & Feature ISSU - Future | Partial Support; Have not demonstrated support for forwarding plane or infrastructure changes |

# CRS-1 ISSU Phases - Roadmap

**Not CC/ECed**   \* Admin Plane must be backward Compatible

**Phase 1, 2 for up to 2+1; Phase 3 for all ; All phases have per-SDR support\***

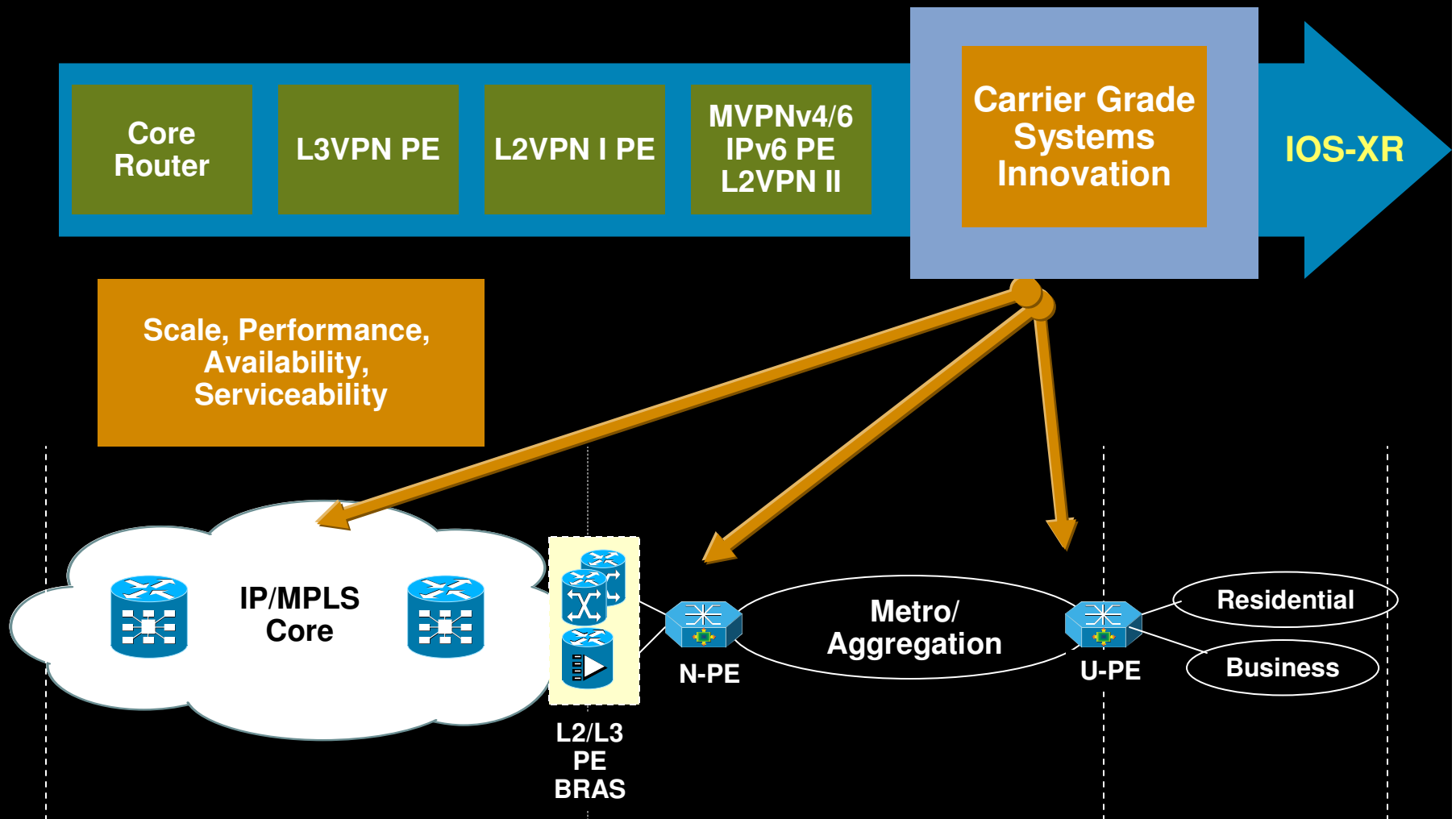| ISSU Phase | Deliverable | Technique | Forwarding Impact | Target Delivery Dates |
|---|---|---|---|---|
| 1 | Maintenance Release | Parallel Upgrade <br> RP SO + LC Reload | < 30 min | 1H CY 2009 |
| 1a | Maintenance Release <br> No forwarding change | Parallel Upgrade <br> RP SO + LC MDR (No ASIC Reset) | < 1 second | 1H CY 2009 |
| 1b | Maintenance Release <br> With forwarding change | Parallel Upgrade <br> RP SO + LC MDR (RSM/ASIC reset) | < 60 seconds | 1H CY 2009 |
| 2 | Feature Release <br> With forwarding change | Parallel Upgrade <br> RP SO + LC MDR (RSM/ASIC reset) | < 60 seconds | CY 2010 |
| 3 | Feature Release <br> With forwarding change | Granular Upgrade <br> Card-by-Card Upgrades | < 60 seconds | CY2011 |

# ISSU Next steps

1. **CC ISSU framework**
   **Target CC:  Jun 2008**

2. **EC entire Phase 1x**
   **Target EC:  Aug 2008**
   **Target FCS: By 1HCY2009**

3. **Target Phase 2 in 2010**

4. **Target Phase 3 in 2011**

**Start**

# Anchor Point: Carrier Grade Innovation

| Core Router | L3VPN PE | L2VPN I PE | MVPNv4/6 IPv6 PE L2VPN II | Carrier Grade Systems Innovation | IOS-XR |

**Scale, Performance, Availability, Serviceability**

IP/MPLS Core

L2/L3 PE BRAS

N-PE

Metro/ Aggregation

U-PE

Residential

Business

# Cisco's IP NGN Framework
## IOS XR playing a key role in Cisco NGN Vision