

Designing High Availability MPLS Networks



Levente Laposi

IP Competence Center, Vienna

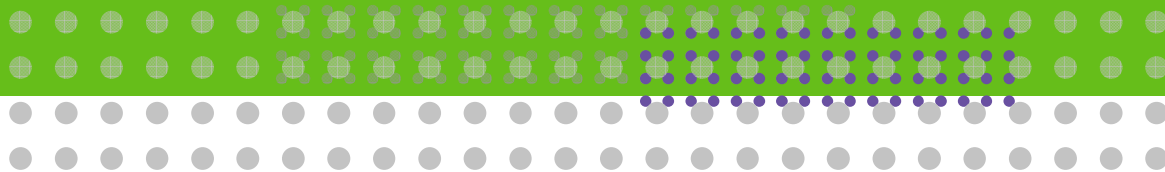
April 2009

Agenda

1. High Availability Challenges
2. Node Resiliency Features
3. Designing for Backbone Resilience
4. Fault Detection Methods
5. Designing for Service Availability: Customer Access Resiliency
6. Questions?

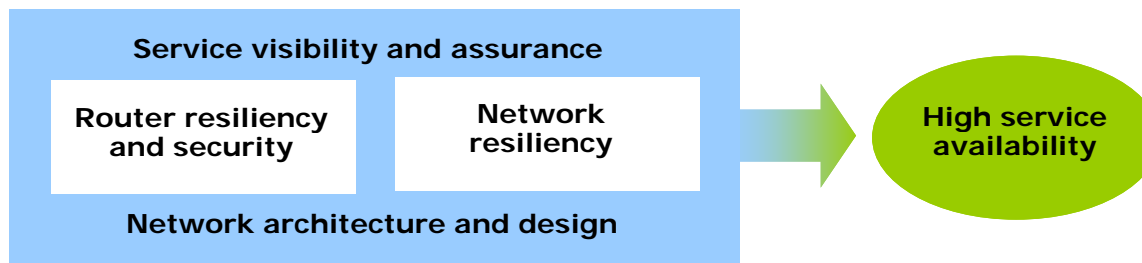
1

High Availability Challenges



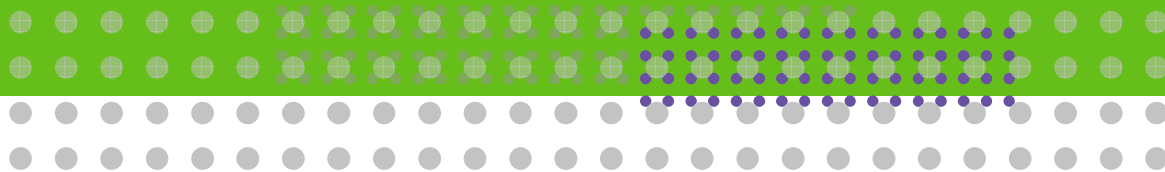
Model for achieving IP high availability

Good working definition of IP high availability is the continuous, uninterrupted delivery of all services with an actual *achievable* service availability of 99.999% or greater...



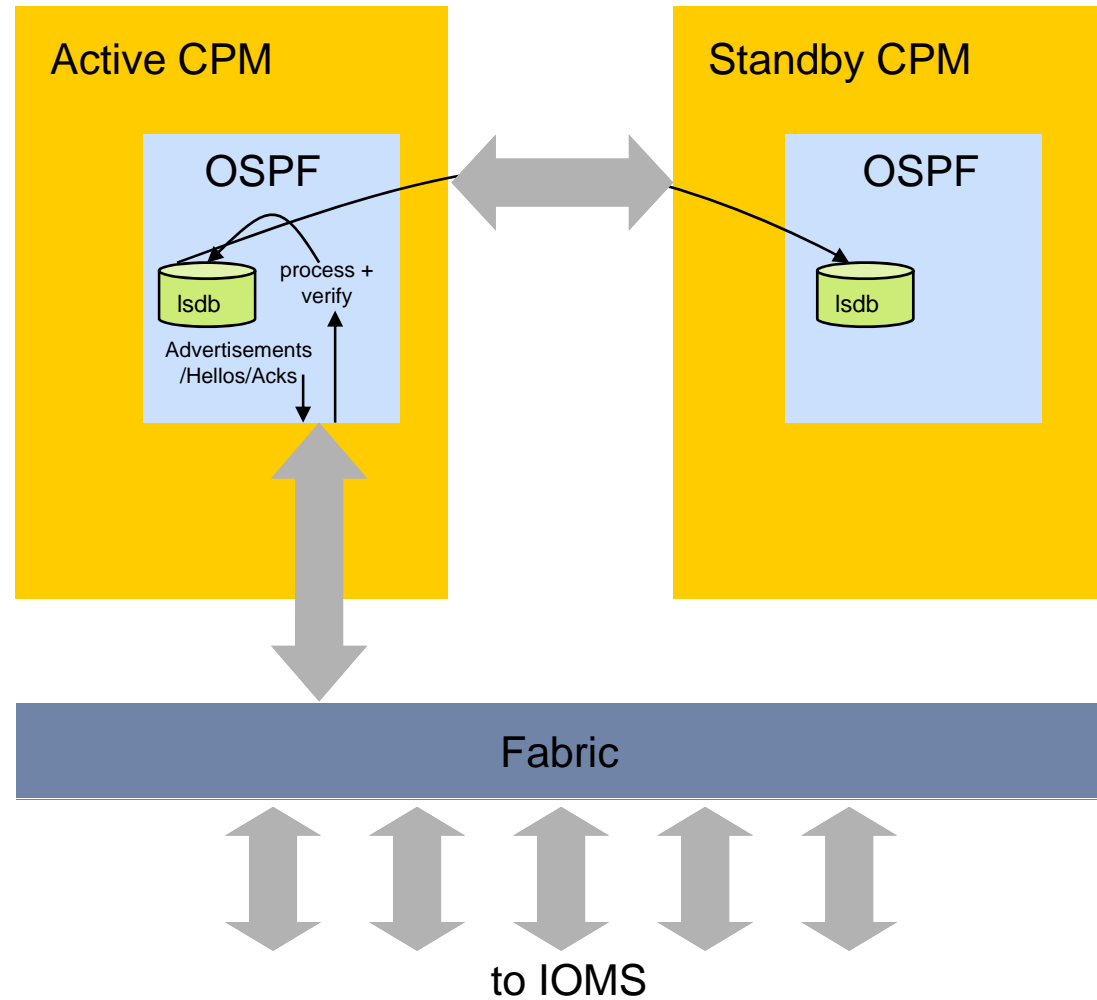
- Network architecture and design
 - Ring versus meshed architectures
 - Indirect versus direct PE uplink redundancy
- Node resiliency and security
 - Node requirements (NSR, ISSU)
- Network resiliency
 - Protocol reliability and interactions
- Service visibility and assurance
 - Continuous monitoring of all network layers, protocols, and services being delivered
 - Proactive detection and resolution of potential areas of outage before they occur
 - Rapid analysis, isolation, and recovery from a failure at any level to minimize network disruption and outage time

2 Node Resiliency Features



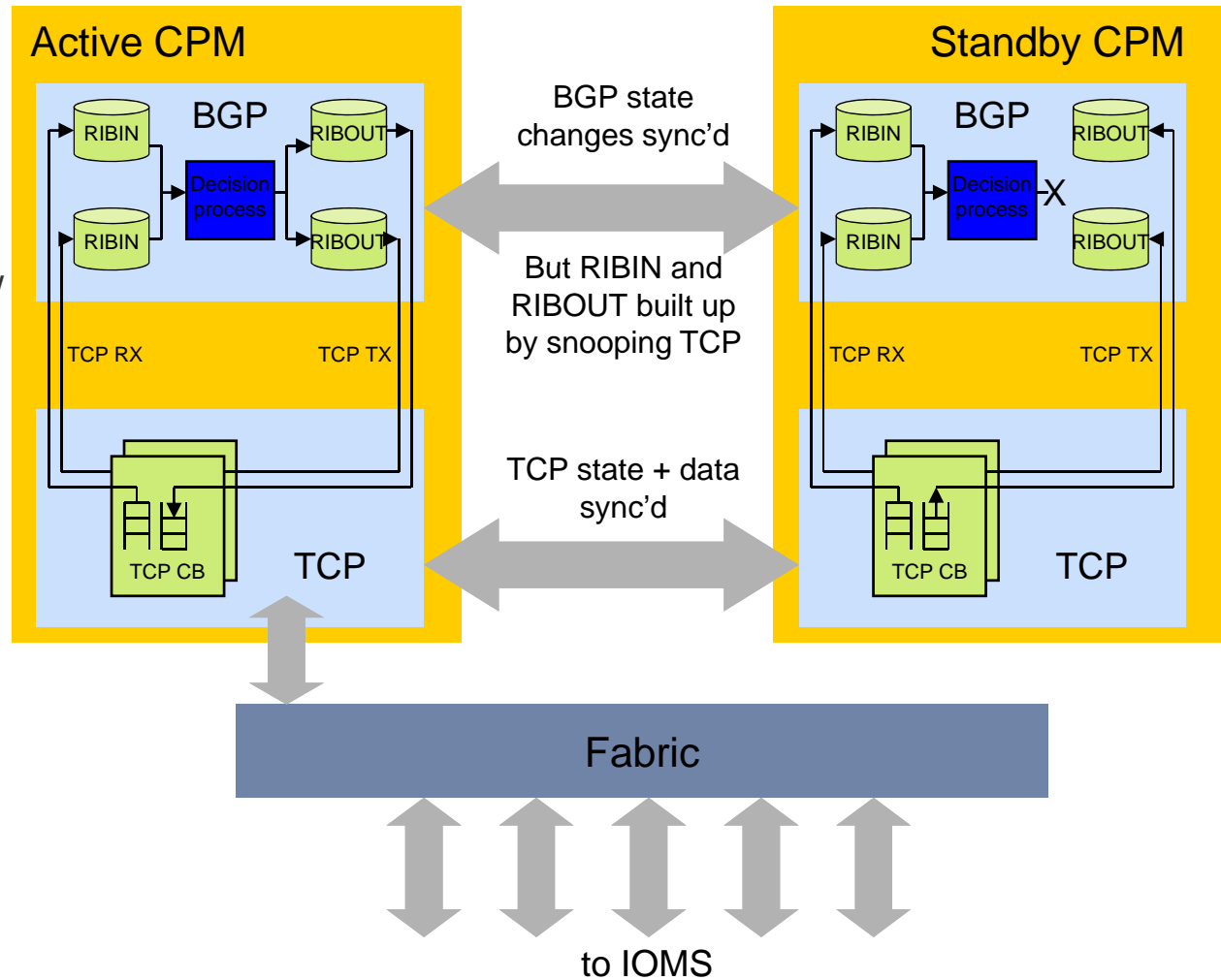
Hot Redundancy = High Availability

- How you do it is important!
 - Not parallel processing, successful operations go to the standby
- Active can never rely on anything from the standby
- Results are sent over (if possible)
- Your architecture must support it! Must be designed in from the start.
- Super fast, reliable link between active and standby
- Fast failover detection



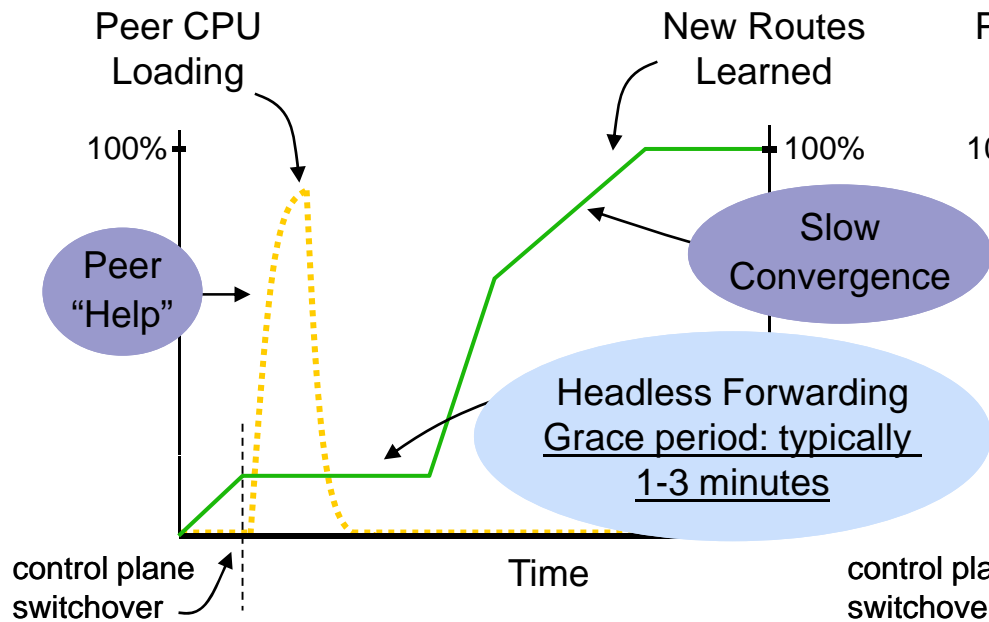
BGP Based High Availability

- Preventing bug propagation (especially important for upgrades)
 - By building RIB-OUT on the standby based on what the active transmitted (by snooping active's TCP transmits), we know exactly what has been advertised
 - Any inconsistencies can be corrected!
- Can data be missed?
 - Standby ack's before packets are transmitted (even if data hasn't yet been processed on the standby)
 - Doesn't this violate one of our rules?
- But what about a packet causing us to crash?
 - Standby can never get ahead of active
 - death packet handling

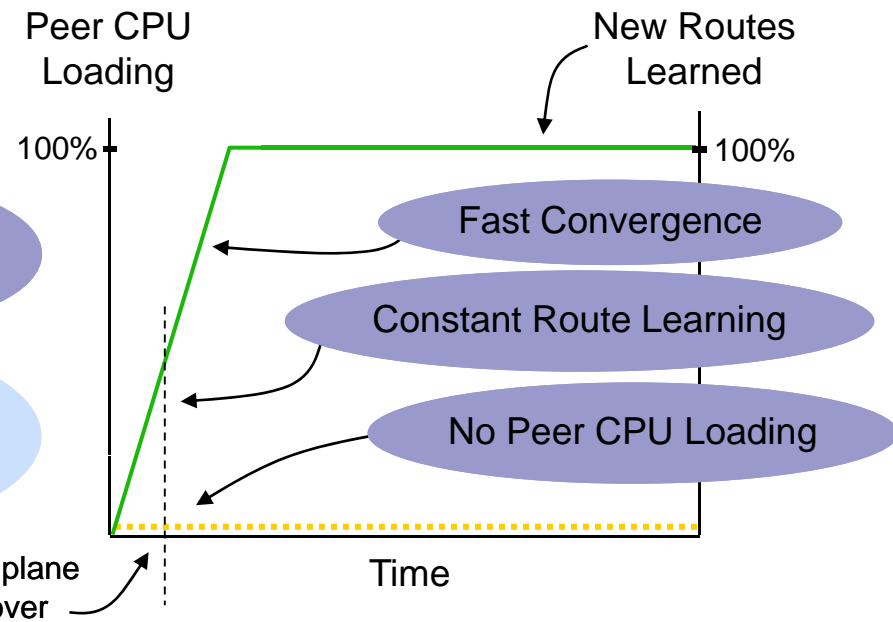


Recovery Comparison between GR and Non-Stop Forwarding – Expected Results

Graceful Restart



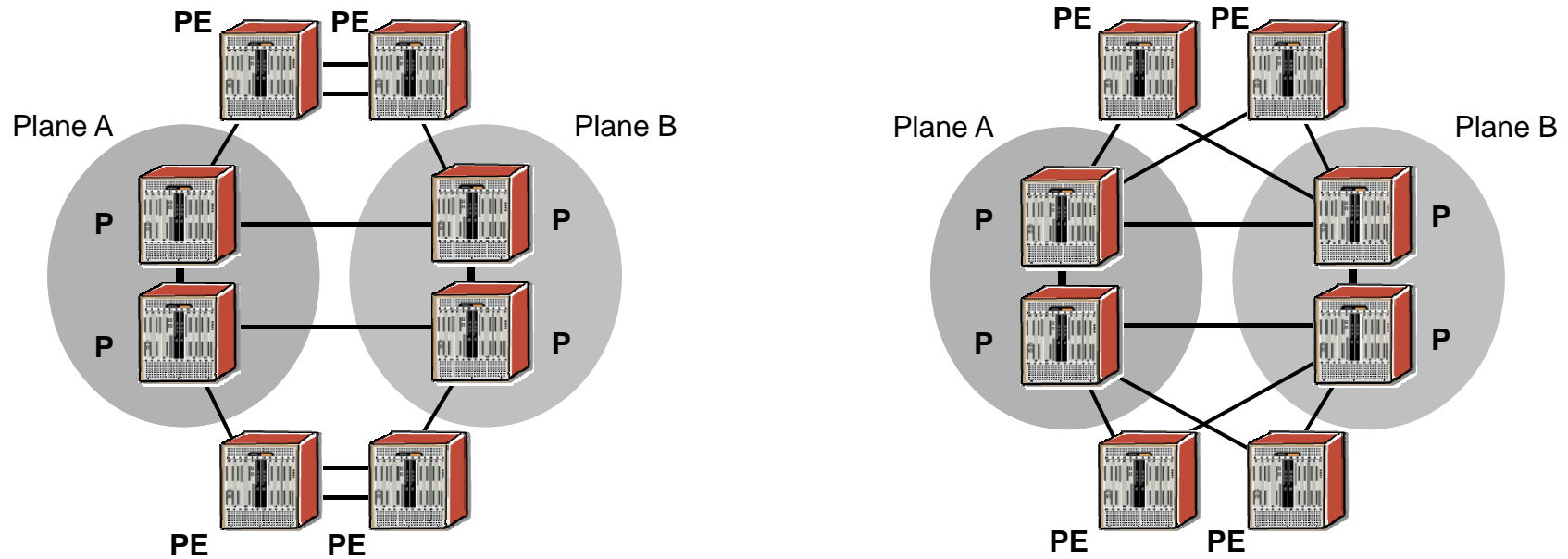
Non-Stop Routing



"Stop-and-Restart-Routing"
Network-Impact, Peers Help
Vendor X

"Non-Stop" Routing
Self-Contained & Transparent
Alcatel-Lucent

Network Recovery Comparison between GR and Non-Stop Forwarding - Expected Results



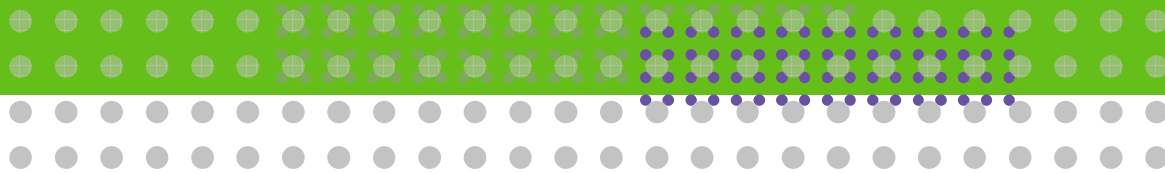
Graceful restart introduces UNCERTAINTY and does not guarantee SLA requirements below the grace period.

P/PE node reload gives better SLA guarantees

SOLUTION=HIGH AVAILABILITY

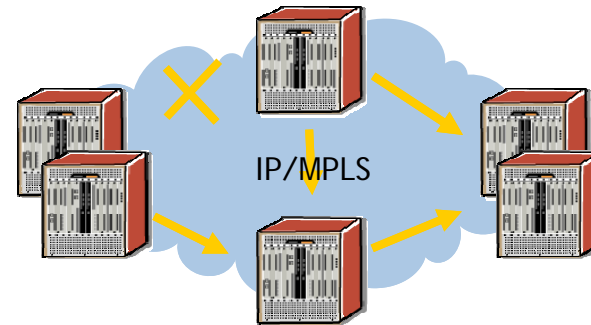
3

Designing for Backbone Resilience



Network layer convergence change propagation: IGP

- Upon detection of a link failure, the local node must generate a new IS-IS LSP/OSPF LSA to reflect the current state of its local interfaces
- The time for a network to fully converge following a link-state change is essentially derived from the following inputs:
 - Time taken for the source system to generate and flood the LSP/LSA to adjacent neighbors
 - Time taken for the LSP/LSA to propagate to adjacent neighbors
 - Time taken for the adjacent neighbors to re-flood the LSP/LSA and subsequently execute an SPF to re-compute the SPT topology. It is worthy of note that the LSP/LSA must be re-flooded BEFORE an SPF is executed.

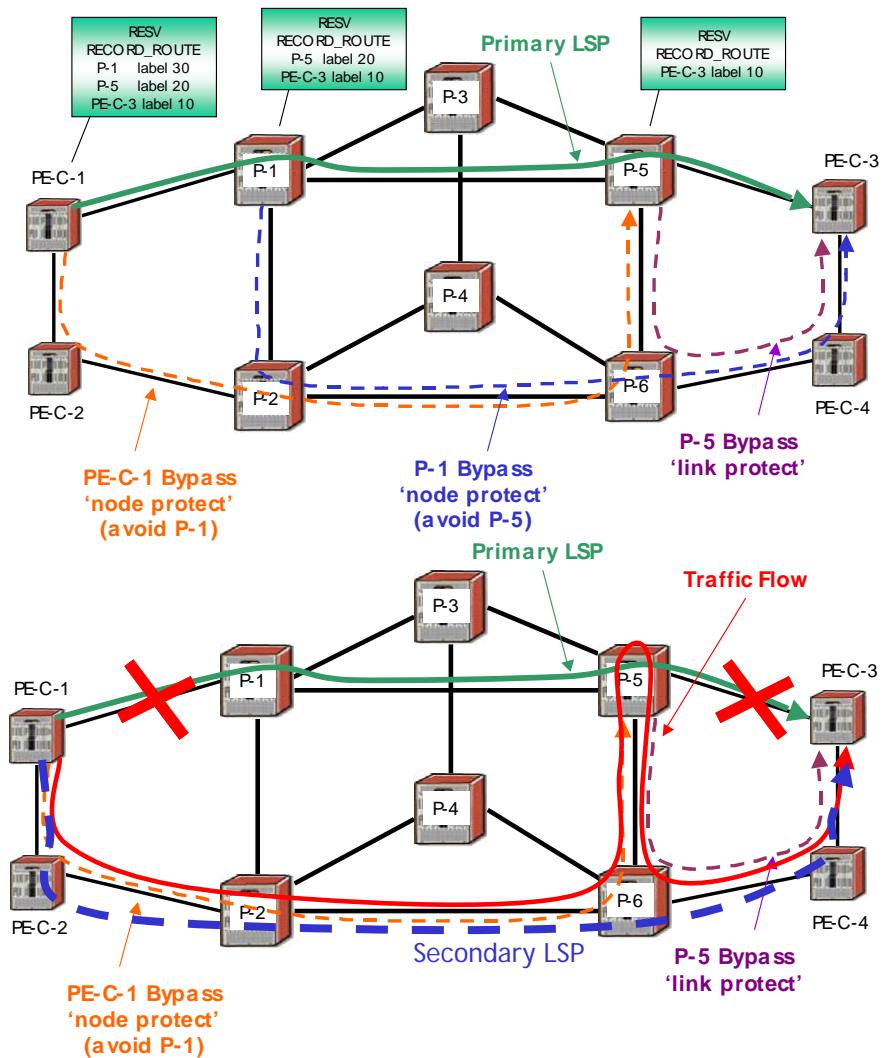


- OSPF
 - spf-wait <max-spf-wait (msec)> <spf-initial-wait (msec)> <spf-second-wait (msec)>
 - lsa-arrival <msec>
 - lsa-generate <msec>
- ISIS
 - spf-wait <spf-wait (s)> <initial-wait (msec)> <second-wait (msec)>
 - lsp-wait <lsp-wait (s)> <initial-wait (s)> <second-wait (s)>

Network layer convergence

Next best path calculation and path restoration: MPLS/IGP

- MPLS Fast Reroute (RFC4090) provides a standard restoration mechanism across the network
 - Facility: 1:n protection or Detour: 1:1 protection
 - Link and node protection
 - Restoration time is sub-50 ms
- MPLS primary/secondary
 - Fault propagated to the LER for primary/secondary decision
 - Head-end decision on primary/secondary
 - More control over the routing of the protected path
- LDP
 - Simple/scalable protocol to operate (multipoint to point approach)
 - Relies on the IGP convergence (200-300 ms depending on the network span)



Network layer convergence

LSP failover and recovery

- MPLS Fast Reroute usually refers to several aspects related to RSVP-TE failure recovery
 - Failover sequence of events (link/node)
 1. Failure detected by adjacent LSR (for example, link failure)
 2. Adjacent LSR moves traffic onto pre-established bypass LSP (sub-50 ms)
 3. LSR propagates error message back to LER
 4. LER moves traffic onto secondary LSP (order of magnitude is 100s of ms depending on configuration)
 5. LER attempts to re-establish new primary path (global revertive)
 - When primary LSP recovers, traffic re-establishes on the primary LSP
 - Make-before-break to re-established primary path
 - Typically, for critical applications or strict constraints on path routing, the recommendation is to build a primary LSP with a secondary backup LSP
 - Detour/Bypass for sub-50ms link/node protection (may be sub-optimal)
 - A secondary backup LSP can be pre-engineered to provide a 2nd best path during the primary failure
-

Network layer convergence

Next best path calculation and path restoration: LDP or RSVP

	LDP	RSVP
Protocol	Multipoint to point protocol	Point to point protocol
Scalability	Scales more simply, One LSP per PE	Mesh issues in large networks Complexity to mitigate
Simplicity	Simple operation	More control, more complexity
Convergence	200 msec - 2 sec depending on the network topology	+/- 50 msec depending on protection design
Traffic engineering	Not available	Available (intra-area) LDPoRSVP (inter-area)
Protection mechanisms	LDP ECMP	FRR, secondary standby LSPs

Alcatel-Lucent recommends RSVP. It is more future proof for traffic engineering.

Network layer convergence

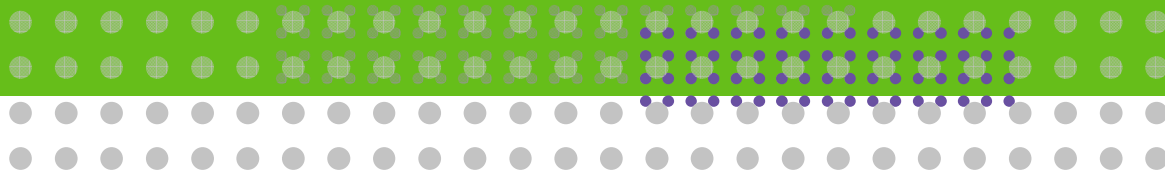
Next best path calculation and path restoration: RSVP options

	Strict	Loose	
Control	High - depends on the operator	Low - depends on the IGP	
Operation	Complex - full mesh to be setup	Easy - CSPF makes this happen	
Protection	Primary/Secondary	One-to-one (detour)	Facility
Control	High when used with strict paths	Low - depends on the IGP	Low - depends on the IGP
Operation	Complex - full mesh to be setup	Easy - CSPF makes this happen	Easy - CSPF makes this happen
Scalability	Lower - Per-LSP end-end	Low - Per-LSP across network	Best - RSVP scalable protection, shared across network

Alcatel-Lucent recommends RSVP with loose LSP(s) using facility protection scheme and secondary LSP(s)

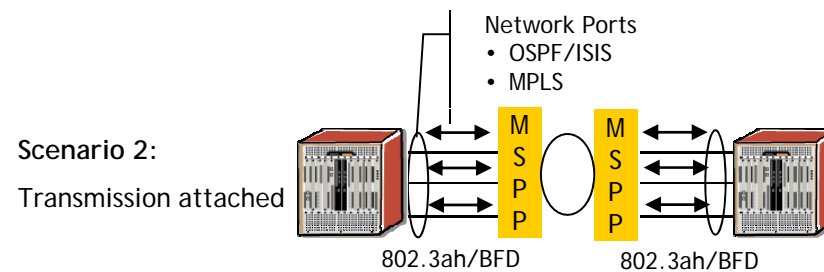
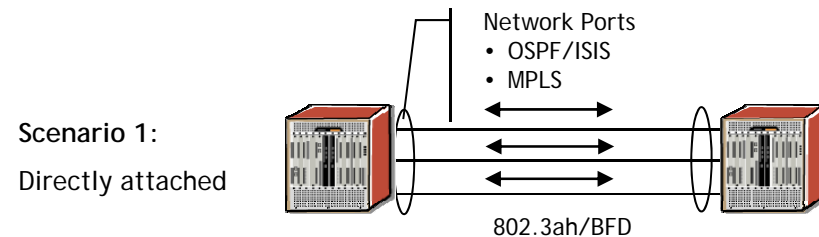
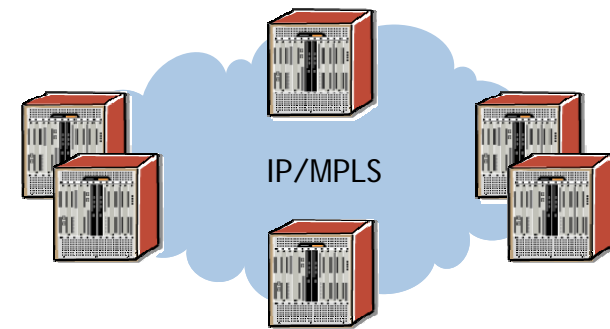
4

Fault Detection Methods

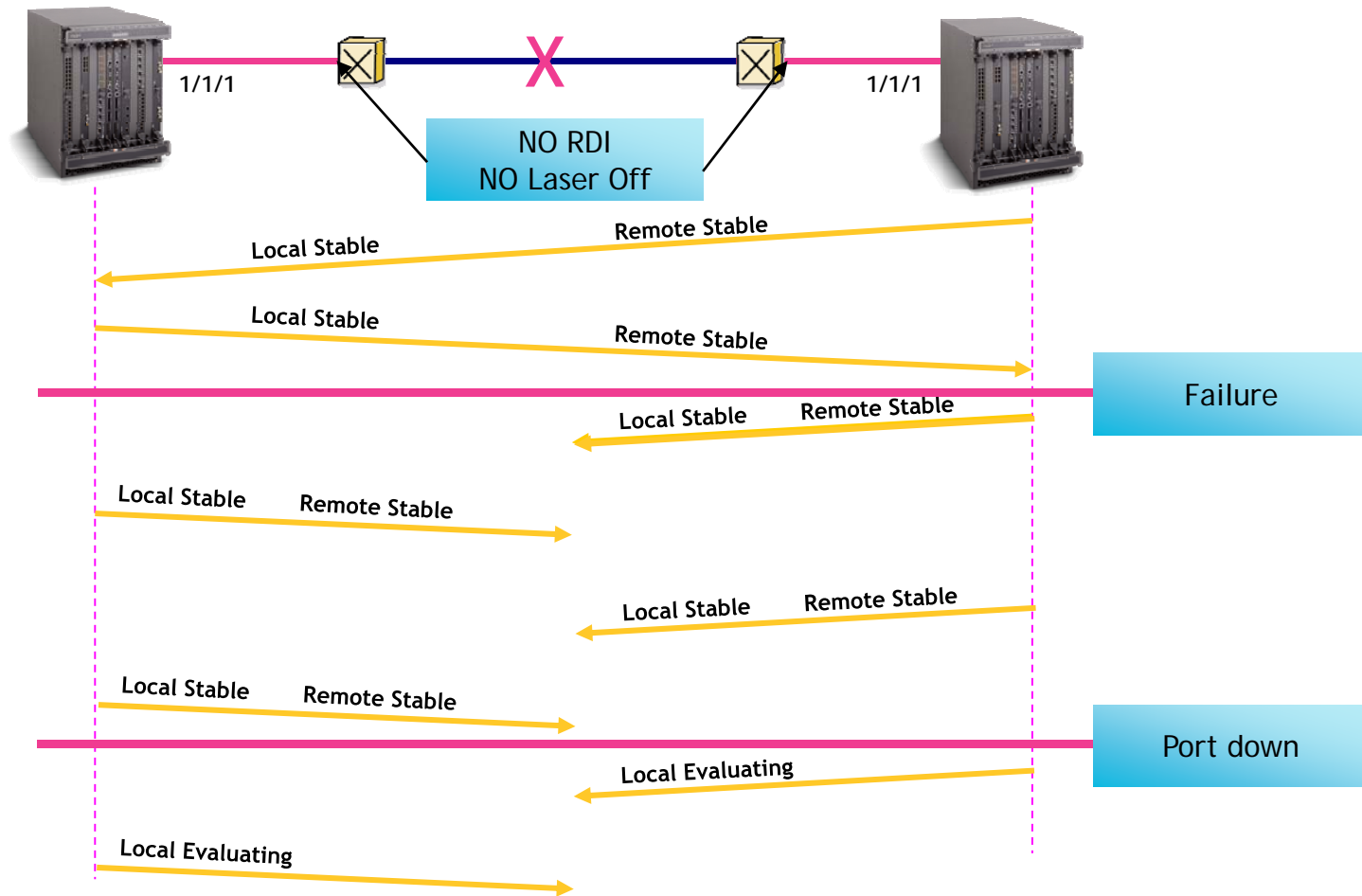


Network layer fault recovery: Fault detection

- Fault recovery time: detection + recovery
- A trigger is required to activate the bypass tunnel
 - OAM or LOS for Sonet/SDH
 - LOS or RDI for Ethernet
- Direct connection or indirect connection may change the available triggers
 - Local link failure results in loss-of-light (LOS) with rapid local detection
 - Failures within a transmission network require propagation of failure or a higher-layer trigger
- When no other trigger is available consider BFD or 802.3ah
 - Can also act as a last resort trigger even with other mechanisms available
 - 802.3ah for Fast Reroute trigger
 - BFD over LAG
 - BFD for RSVP (BFD triggers equivalent actions as interface down)

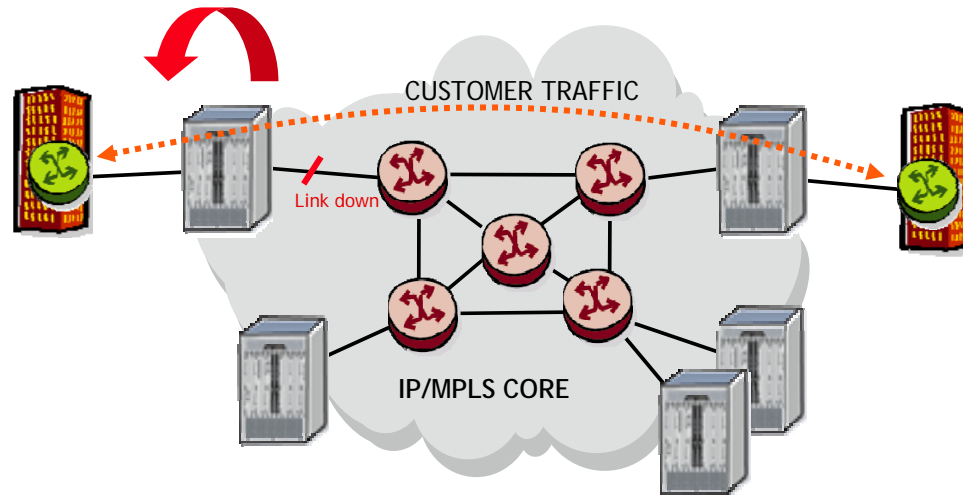


802.3ah Example - Using OAM bidirectional link failure indications



If MPLS FRR is enabled and possible, transmit-interval is set to 100 ms and multiplier is set to 2:
 Detection time is between 100 and 250 ms

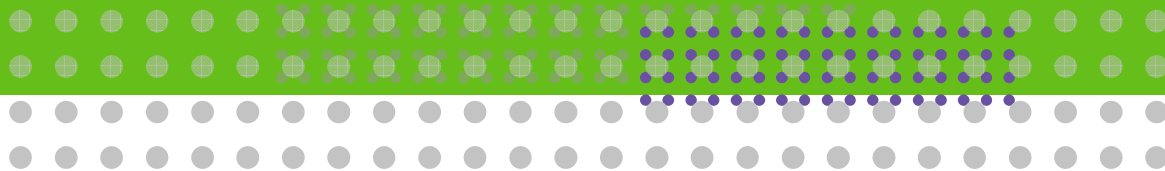
OAM enhancements: Link-level Loss Forwarding



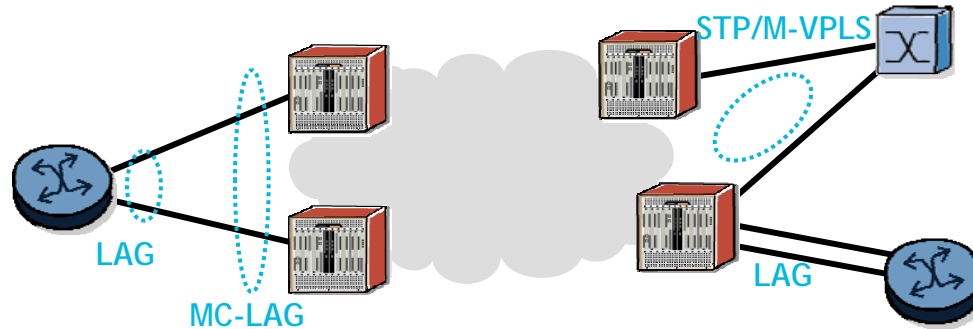
- Need e2e OAM fault notification for Ethernet VLL service
 - End user wants to activate a backup connectivity over another service provider
- This feature shuts down the laser on the interface to the CE under one of the following conditions:
 - Local fault on the PW or service
 - Remote fault on the SAP or PW (signaled with label withdrawal or T-LDP status bits)
- This feature currently applies to a NULL Ethernet SAP

5

Designing for Service Availability Customer Access Resiliency

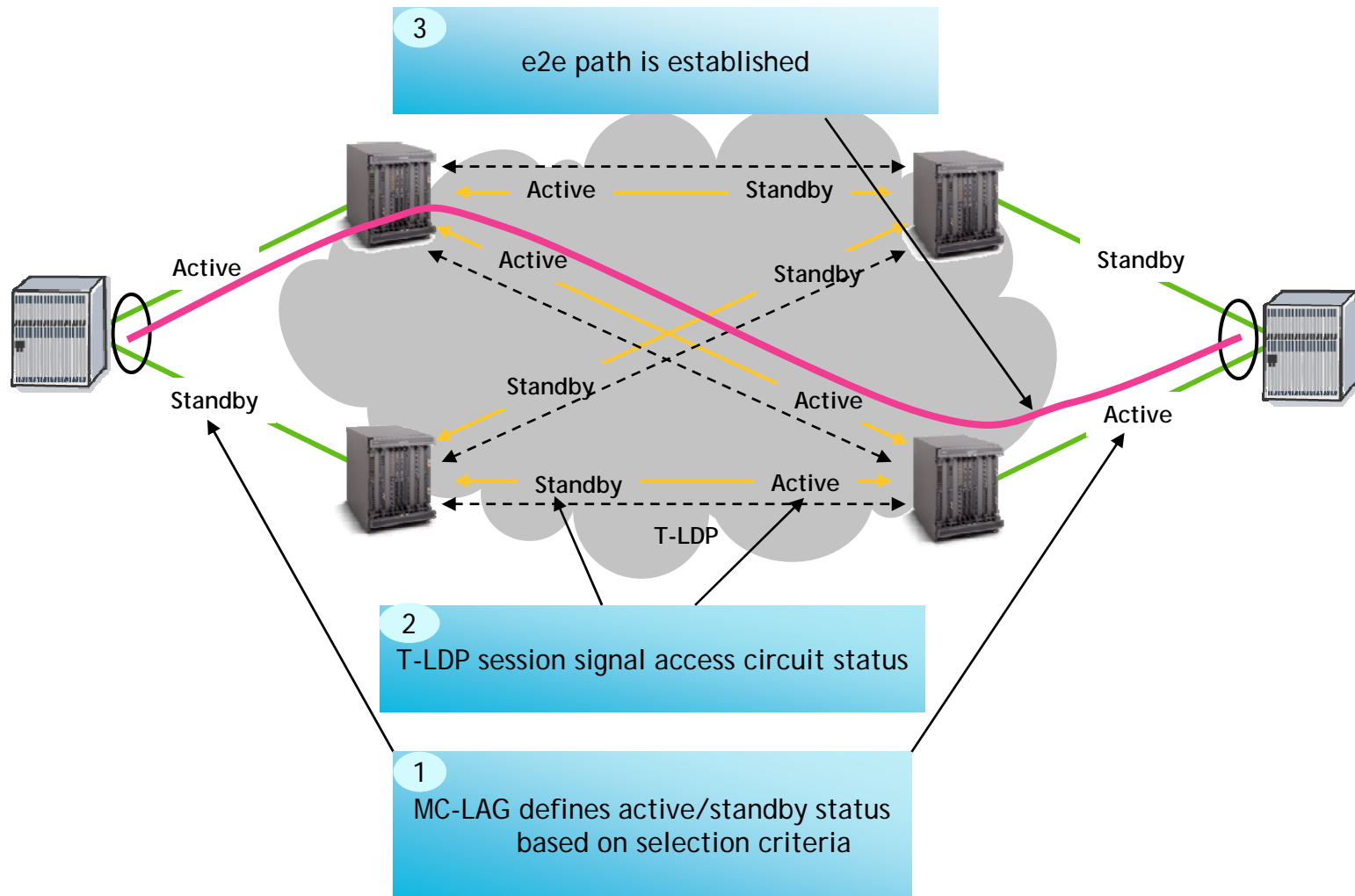


VPN Services: Access redundancy options



- **Link redundancy for access to single PE**
 - VLL service → LAG
 - VPLS service → LAG or STP
 - L3 services → LAG
- **Dual homing to two PEs**
 - VLL service → Multi-chassis LAG (MC-LAG) + active/standby PW
 - VPLS service → STP, M-VPLS or Multi-chassis LAG (MC-LAG)
 - L3 services → VRRP, SRRP, Routing/BFD

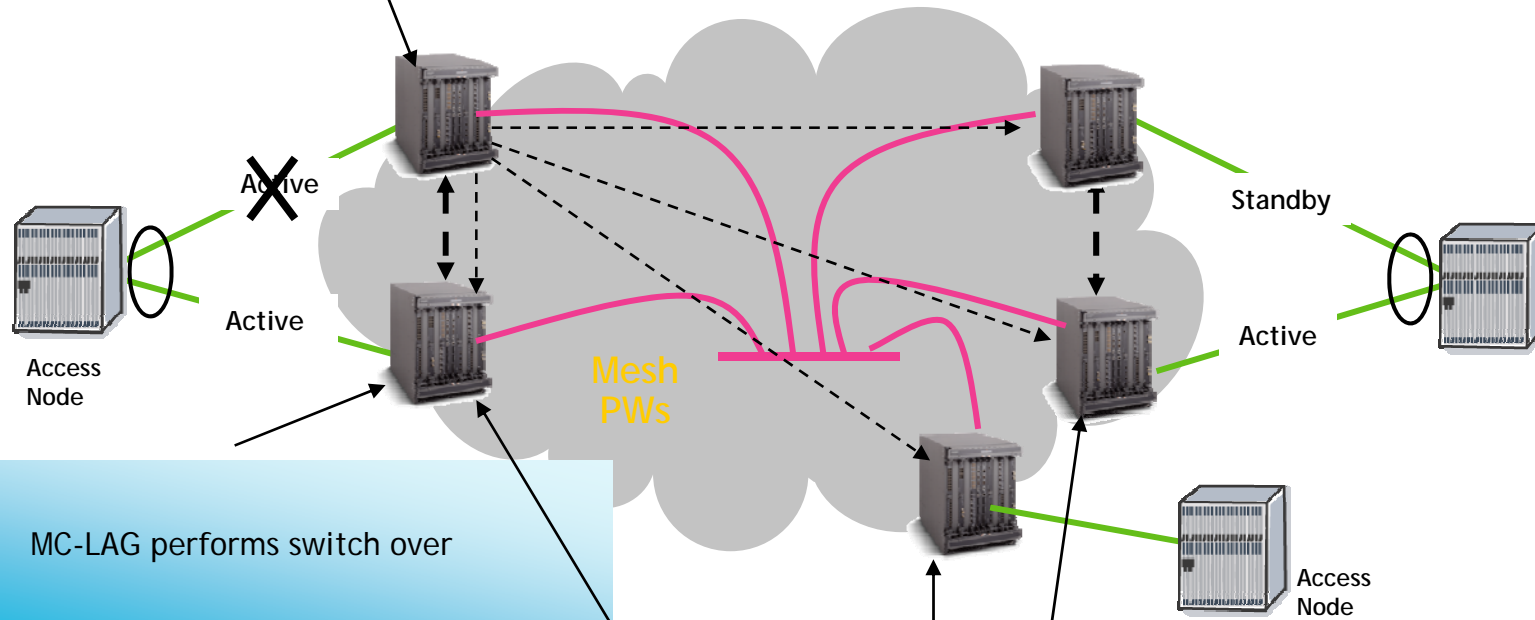
Dual homed point-to-point L2-VPNs (VLLs) Pseudowire redundancy operation



Dual-homing for multipoint L2 VPNs (VPLS)

MC-LAG and VPLS: LAG link failure

2 Sends LDP address withdraw message to all peers (referred to as forget me MAC-FLUSH)

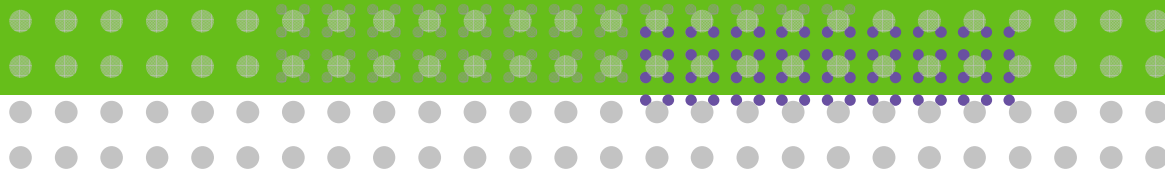


1 MC-LAG performs switch over

3 Relearning process starts again

6

Questions?



www.alcatel-lucent.com

