



# IPv6 Microsegmentation

Ivan Pepelnjak ([ip@ipSpace.net](mailto:ip@ipSpace.net))  
Network Architect

ipSpace.net AG

# Who is Ivan Pepelnjak (@ioshints)

## Past

- Kernel programmer, network OS and web developer
- Sysadmin, database admin, network engineer, CCIE
- Trainer, course developer, curriculum architect
- Team lead, CTO, business owner



## Present

- Network architect, consultant, blogger, webinar and book author
- Teaching the art of Scalable Web Application Design

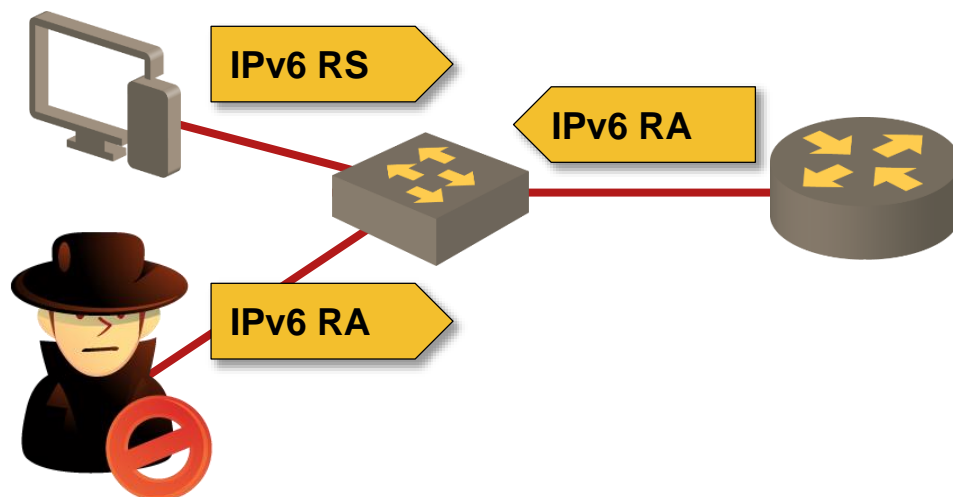
## Focus

- Large-scale data centers, clouds and network virtualization
- Scalable application design
- Core IP routing/MPLS, IPv6, VPN

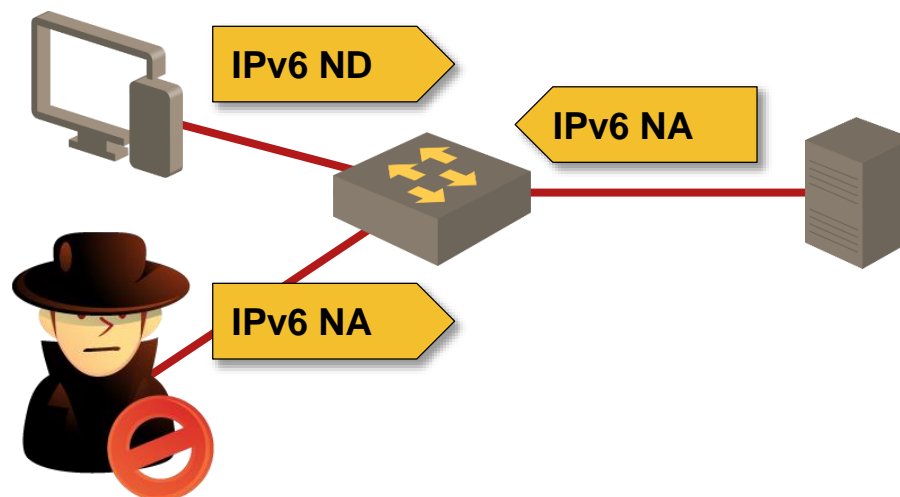


# IPv6 Layer-2 Security Challenges

# The Problem



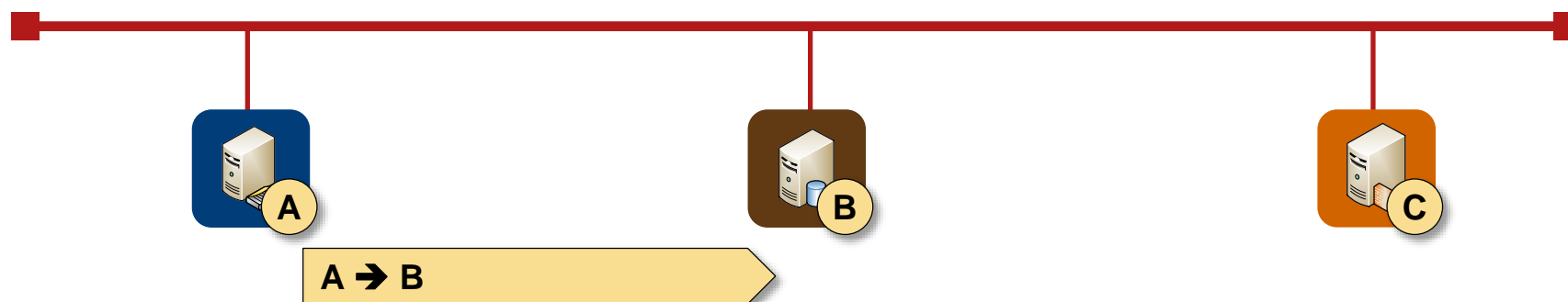
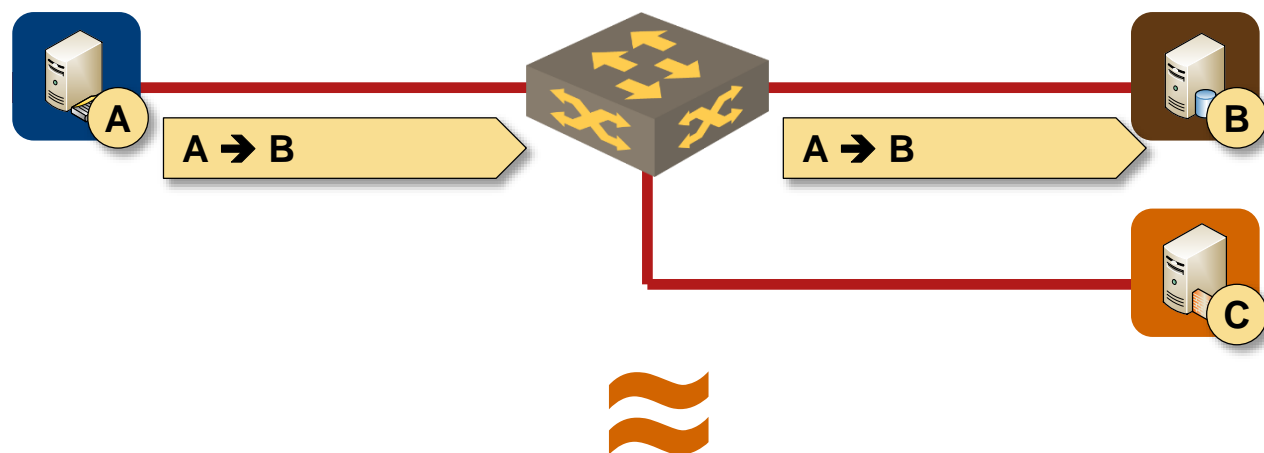
- **Assumption:** one subnet = one security zone
- **Corollary:** intra-subnet communication is not secured
- **Consequences:** multiple first-hop vulnerabilities



Sample vulnerabilities:

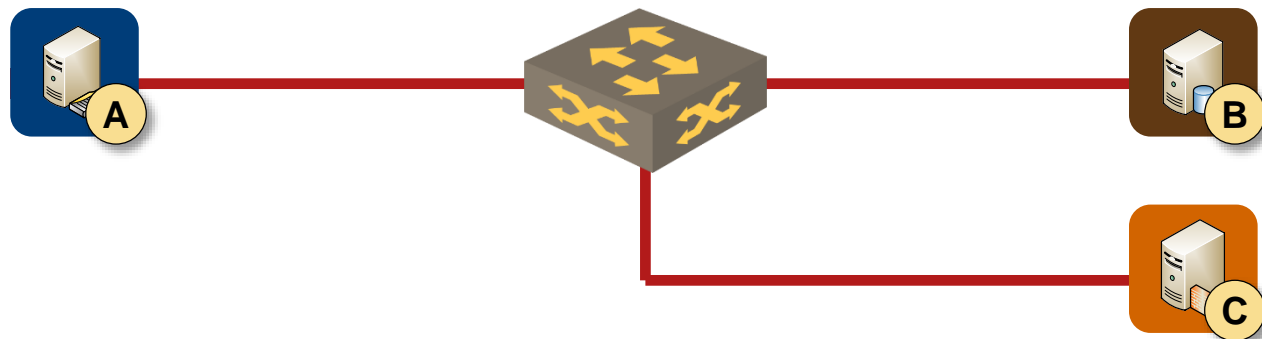
- RA spoofing
- NA spoofing
- DHCPv6 spoofing
- DAD DoS attack
- ND DoS attack

# Root Cause



**All LAN infrastructure we use today emulates 40 year old thick coax cable**

# The Traditional Fix: Add More Kludges



## Typical networking industry solution

- Retain existing forwarding paradigm
- Implement layer-2 security mechanisms

## Sample L2 security mechanisms

- RA guard
- DHCPv6 guard
- IPv6 ND inspection
- SAVI

## Benefits

- Non-disruptive deployment (clusters and Microsoft NLB still works)
- No need to educate customers

## Drawbacks

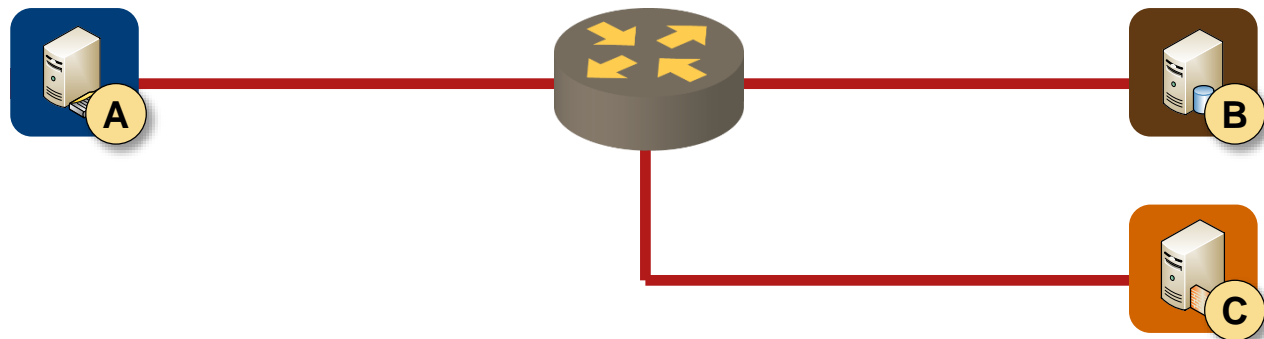
- Not available on all platforms
- Expensive to implement in hardware
- Exploitable by infinite IPv6 header + fragmentation creativity

Can we do any better than that?

# Layer-3-Only IPv6 Networks



# Goal: Remove Layer-2 from the Network



## Change the forwarding paradigm

- First-hop network device is a router (layer-3 switch in marketese)
- Fake router advertisements or ND/NA messages are not propagated to other hosts

## Simplistic implementation

- Every host is in a dedicated /64 subnet
- Default behavior on 3GPP and xDSL networks
- Somewhat harder to implement on Carrier Ethernet, hard on cable networks

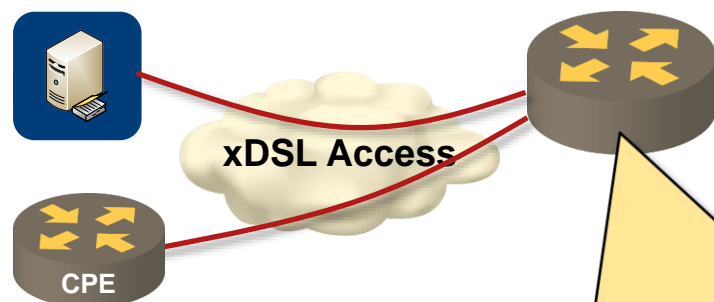


# IPv6 over 3GPP and PPPoX Networks



- Each device-to-network connection is a separate dial-up interface on BRAS/GGNS
- Customer device (phone, computer, CPE) interacts directly with the first-hop router
- A /64 subnet is allocated to each dial-up interface (usually from a local pool)
- Aggregate IPv6 prefix is advertised to the network core to minimize number of prefixes advertised in the core

# Sample IPv6 over PPPoX BRAS Configuration



```
interface Virtual-Template10
  mtu 1480
  peer default ipv6 pool PPP
  ipv6 enable
  ipv6 nd other-config-flag
  no ipv6 nd ra suppress
!
ipv6 local pool PPP 2001:DB8:CAFE::/48 64
ipv6 route 2001:DB8:CAFE::/48 Null0
!
router bgp 65000
  address-family ipv6
    network 2001:DB8:CAFE::/48
```

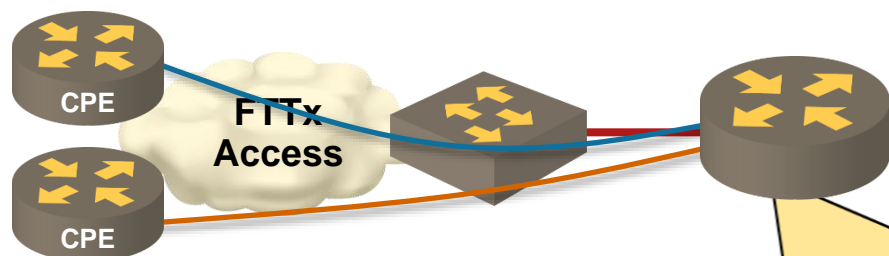
More details in *Building Large IPv6 Networks* webinar

# IPv6 Microsegmentation over Carrier Ethernet Networks



- Option#1: First-hop network device is a layer-3 switch (example: Cisco ME 3600)
- Option#2: Each customer resides in a dedicated VLAN (extensive service automation is highly recommended)

# Configuring VLAN-Based IPv6 Microsegmentation

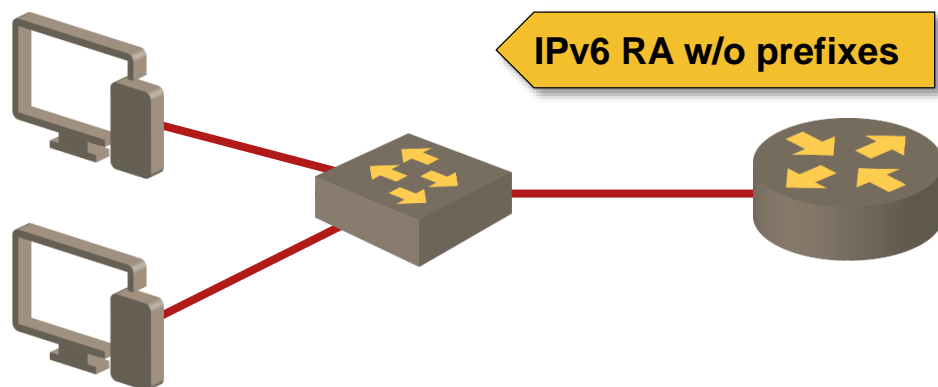


```
interface FastEthernet0/0.100
  encapsulation dot1Q 100
  ipv6 address 2001:DB8:ABBA:100::1/64
!
interface FastEthernet0/0.101
  encapsulation dot1Q 101
  ipv6 address 2001:DB8:ABBA:101::1/64
!
ipv6 route 2001:DB8:ABBA::/48 Null0
!
router bgp 65000
  address-family ipv6
    network 2001:DB8:ABBA::/48
```

More details in *Building Large IPv6 Networks* webinar

# Layer-3-Only Shared Subnets

# Tweaking On-net Determination



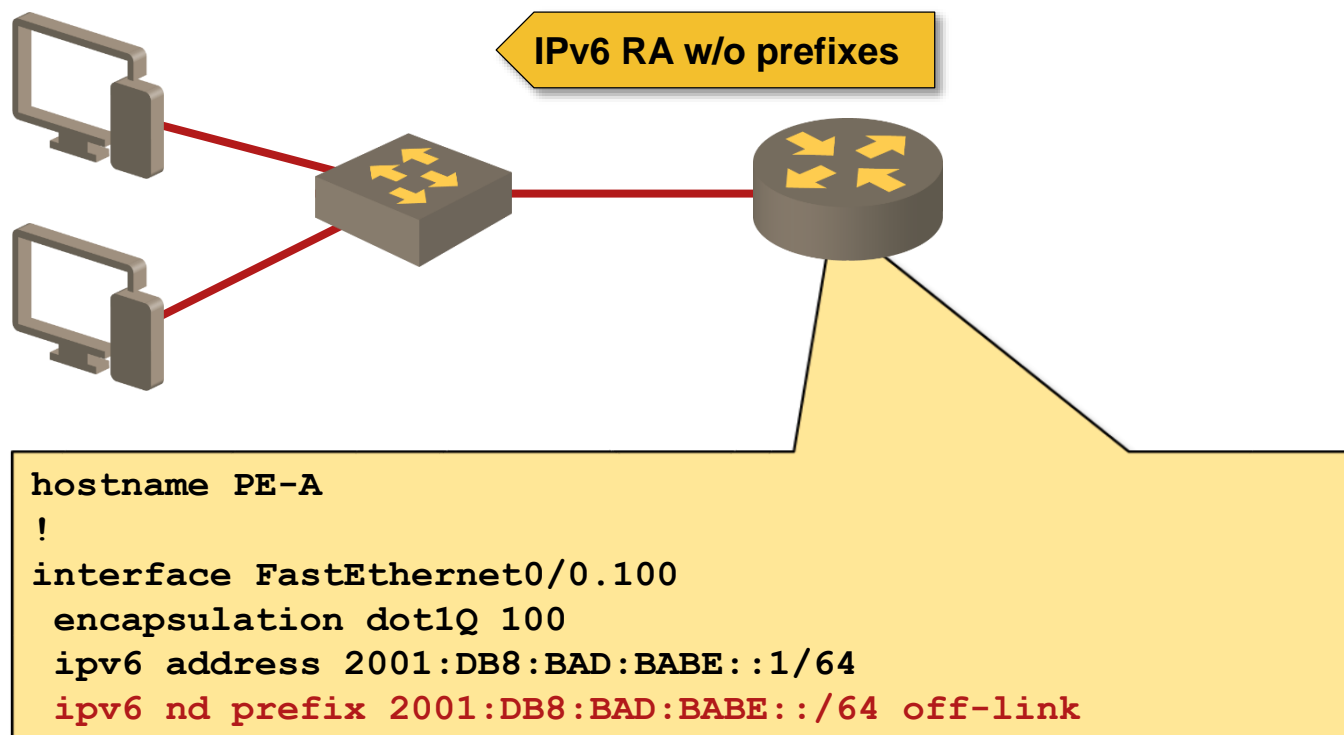
Local subnet is not advertised in RA messages

- IPv6 hosts cannot perform on-net check
- All intra-subnet traffic goes through the first-hop router
- Access lists on first-hop router enforce segmentation

## Drawbacks

- Relies on proper IPv6 host behavior
- RA and ND attacks are still possible without IPv6 first-hop security

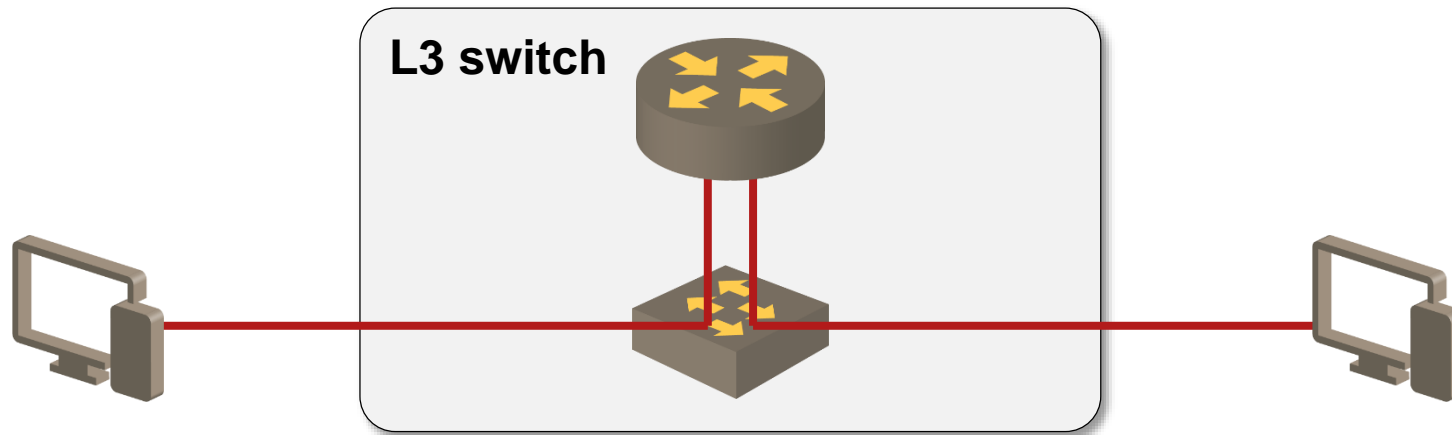
# Configuring Off-Link Local Prefix



- **Off-link** prefix enables SLAAC, but not host-to-host traffic
- **No-advertise** prefix disables SLAAC (combine with **managed-config-flag** to enforce DHCPv6)



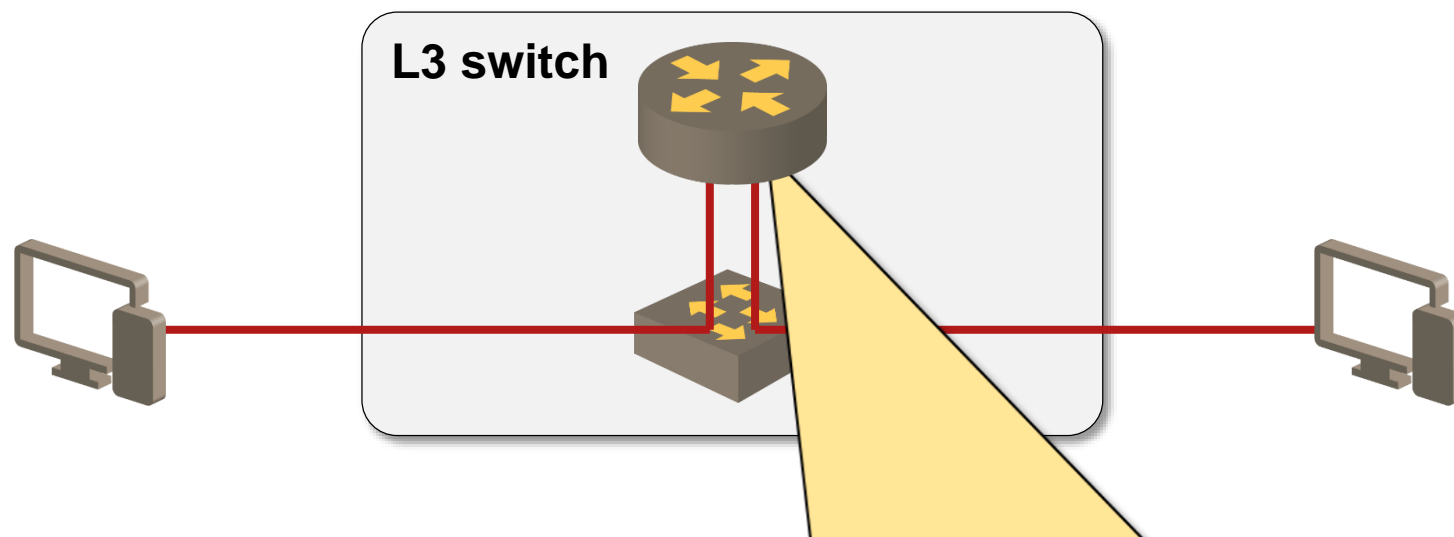
# Tweaking On-net Determination + PVLAN



**Private VLANs** can be used to enforce L3 lookup

- Force traffic to go through L3 device (router / L3 switch)
- Potential solution for campus environments with low-cost L2-only switches or virtualized environments
- L3 device **must not** perform mixed L2/L3 forwarding (hard to implement on a L2/L3 switch)
- This solution could break DAD process → use DAD proxy on the router

# Configuring Duplicate Address Detection Proxy



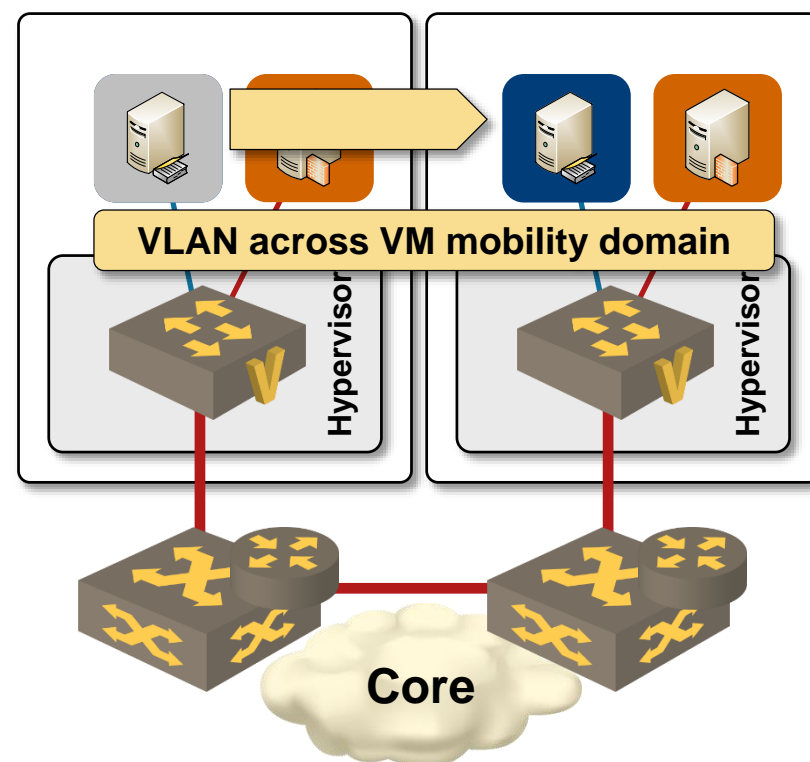
```
hostname PE-A
!  
interface FastEthernet0/0.100  
  encapsulation dot1Q 100  
  ipv6 address 2001:DB8:BAD:BABE::1/64  
  ipv6 nd prefix 2001:DB8:BAD:BABE::/64 off-link  
  ipv6 nd dad-proxy
```

# Data Center Considerations

# Implications of Live VM Mobility

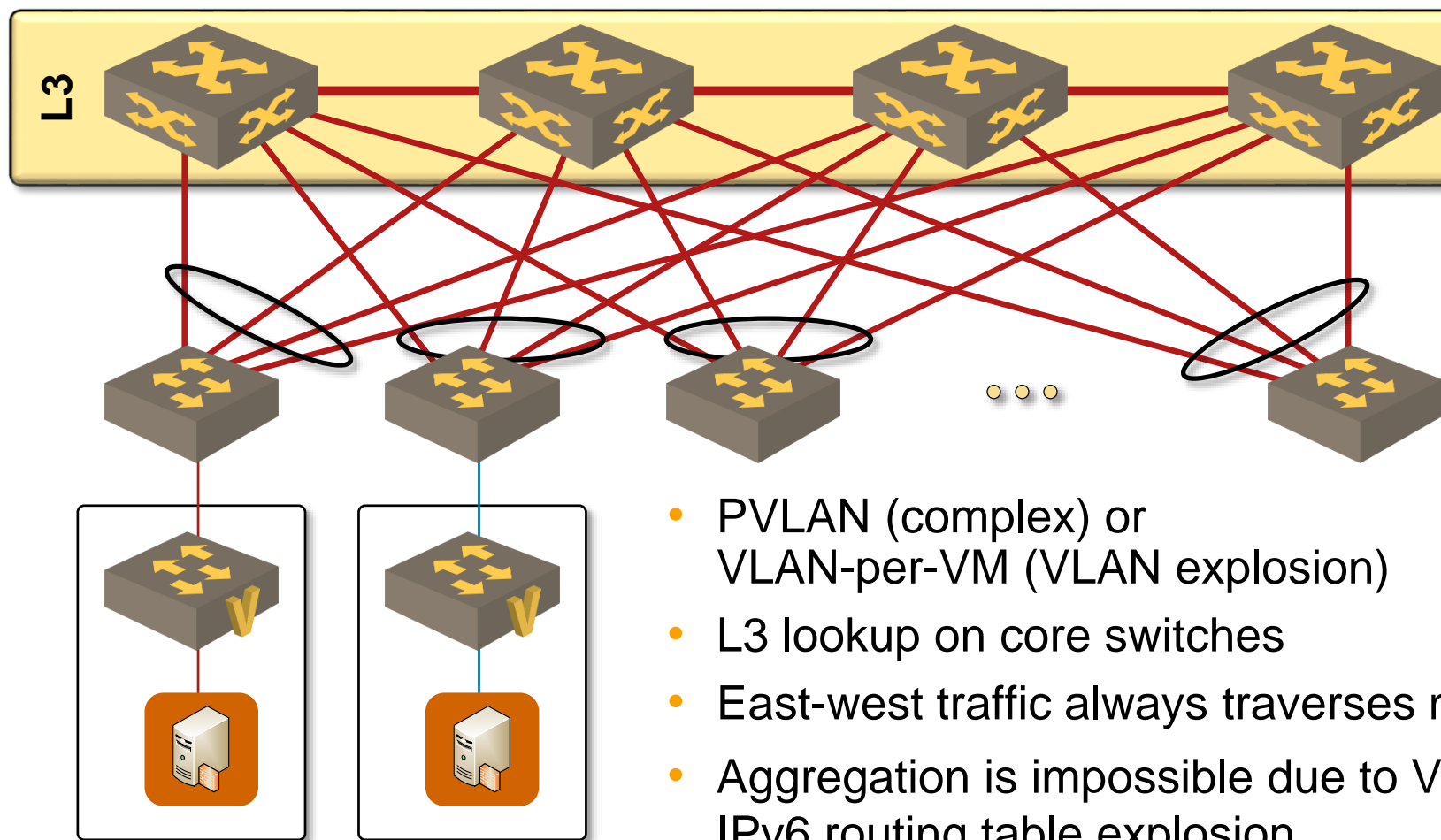
## Challenges

- VM moved to another server must retain its IPv6 address and all data sessions
- Existing L3 solutions are too slow for non-disruptive VM moves
- Live VM mobility usually relies on L2 connectivity between physical servers
- Large VLANs must span the whole VM mobility domain



More details in *VMware Networking* and *Cloud Networking* webinars


# Live VM Mobility with IPv6 Microsegmentation



We need something better in data centers


# Arista Spline Switches

Switch model	Ports	MAC	IPv4	ARP	IPMC	IPv6
7304	128 x 40GbE 512 x 10GbE 192 x 10GBASE-T	288K	16K	208K	104K	8K
7308	256 x 40GbE 1024 x 10GbE 384 x 10GBASE-T					
7316	512 x 40GbE 2048 x 10GbE 768 x 10GBASE-T					



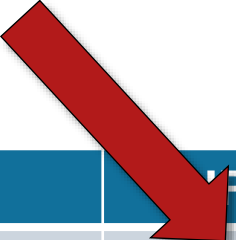
# Brocade VDX ToR Switches

## Port density

Switch model	GE ports	10GE ports	40GE ports	FC ports
VDX 6710	48	6	-	-
VDX 6720-24	24		-	-
VDX 6720-60	60		-	-
VDX 6730-32	24		-	8
VDX 6730-76	60		-	16
VDX 6740 	48		4	

## Table sizes





Switch	MAC	IPv4	ARP	IPv6
VDX 6740	160K	12K	32K	3K
VDX 67xx	32K	2K	12K	-






# Nexus 6000 and 9300 Series Overview

## Port density

Switch	1G	10GE	40GE
9396PX 	48 (SFP+)	48	12
9396TX 	48 (10GBASE-T)	48	12
9336PQ 			36
93128PX 	96 (10GBASE-T)	96	8
Nexus 6001 (48 x SFP+, 4 x QSFP)	48	64	4
Nexus 6004 (96 x QSFP)		384	96

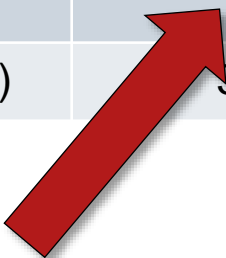
## Table sizes

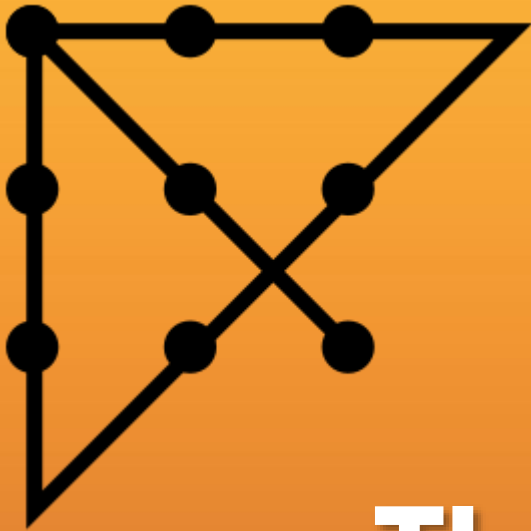
Switch	MAC	IPv4	ARP	IPv6	ND
Nexus 9300	96K	16K	88K	6K	20K
Nexus 6000	115K	24K	64K	8K	32K



# Fixed Data Center Switches – EX Series

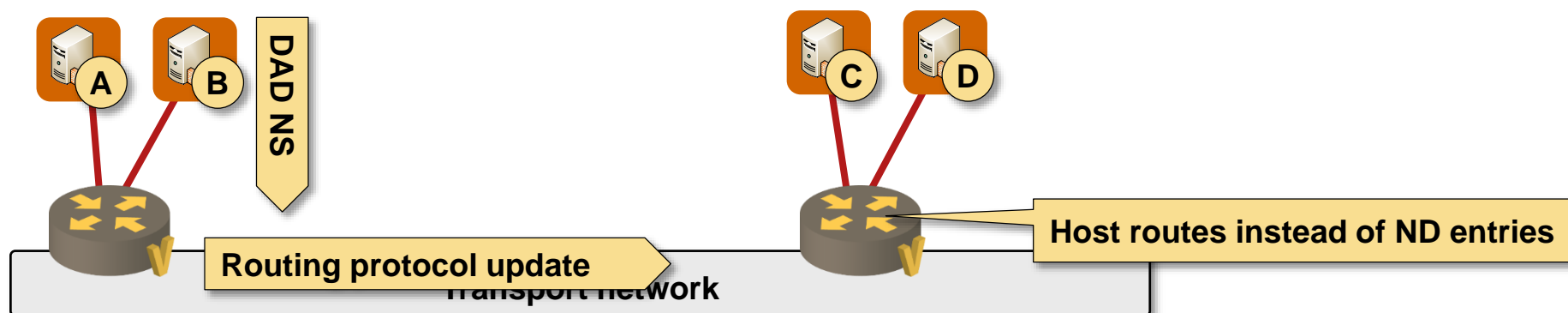
Model	EX4200	EX4300 <sup>New</sup>	EX4500	EX4550
Typical role	ToR	ToR	Tor/Core	ToR/Core
Max ports	48 x 1GE 2 x 10GE	24 / 48 GE 4 / 8 10GE	40 – 48 x 10GE	32 – 48 x 10GE 2 x 40GE
MAC table	32K	64K	32K	32K
IPv4 table	16K	4K	10K	10K
ARP	16K	64K	8K	8K
IPMC	8K	8K	4K	4K
IPv6 table	4K	1K	1K	1K
IPv6 ND	16K (shared)	32K	1K	1K





# Thinking Outside of the Box

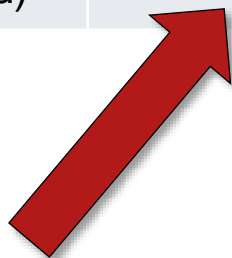
# Intra-Subnet (Host Route) Layer-3 Forwarding



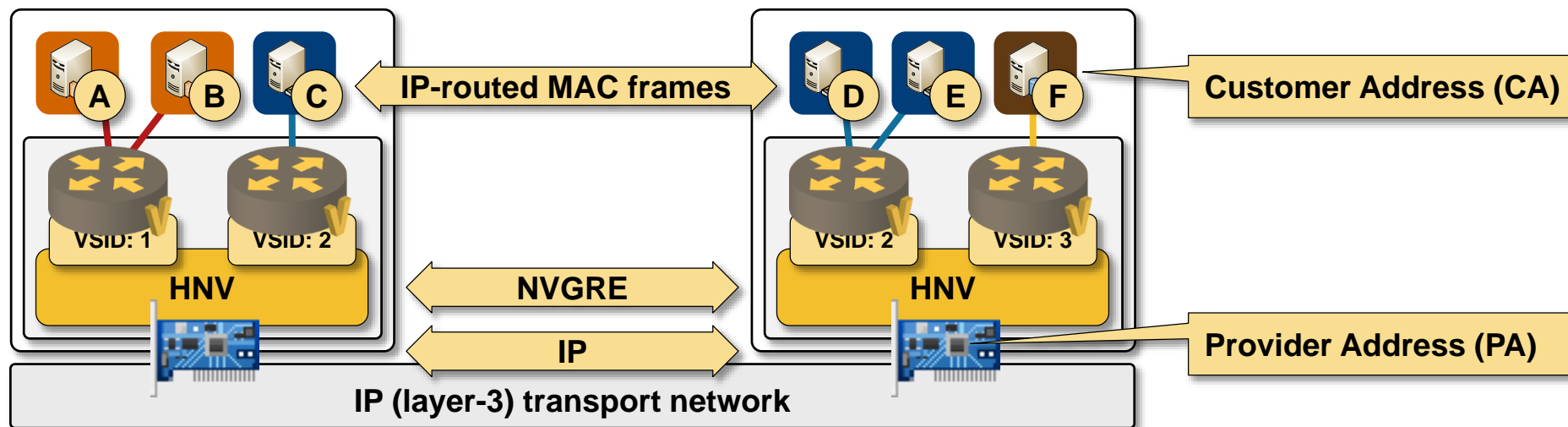
- Hosts are connected to layer-3 switches (routers)
- Numerous hosts share a /64 subnet  
→ a /64 subnet spans multiple routers
- First-hop router creates a host route on DAD, ND or DHCPv6 transaction
- IPv6 host routes are propagated throughout the local routing domain
- Host-side IPv6 addressing and subnet semantics are retained
- IPv6 ND entries are used instead of IPv6 routing table entries

# Fixed Data Center Switches – EX Series

Model	EX4200	EX4300 <sup>New</sup>	EX4500	EX4550
Typical role	ToR	ToR	Tor/Core	ToR/Core
Max ports	48 x 1GE 2 x 10GE	24 / 48 GE 4 / 8 10GE	40 – 48 x 10GE	32 – 48 x 10GE 2 x 40GE
MAC table	32K	64K	32K	32K
IPv4 table	16K	4K	10K	10K
ARP	16K	64K	8K	8K
IPMC	8K	8K	4K	4K
IPv6 table	4K	1K	1K	1K
IPv6 ND	16K (shared)	32K	1K	1K



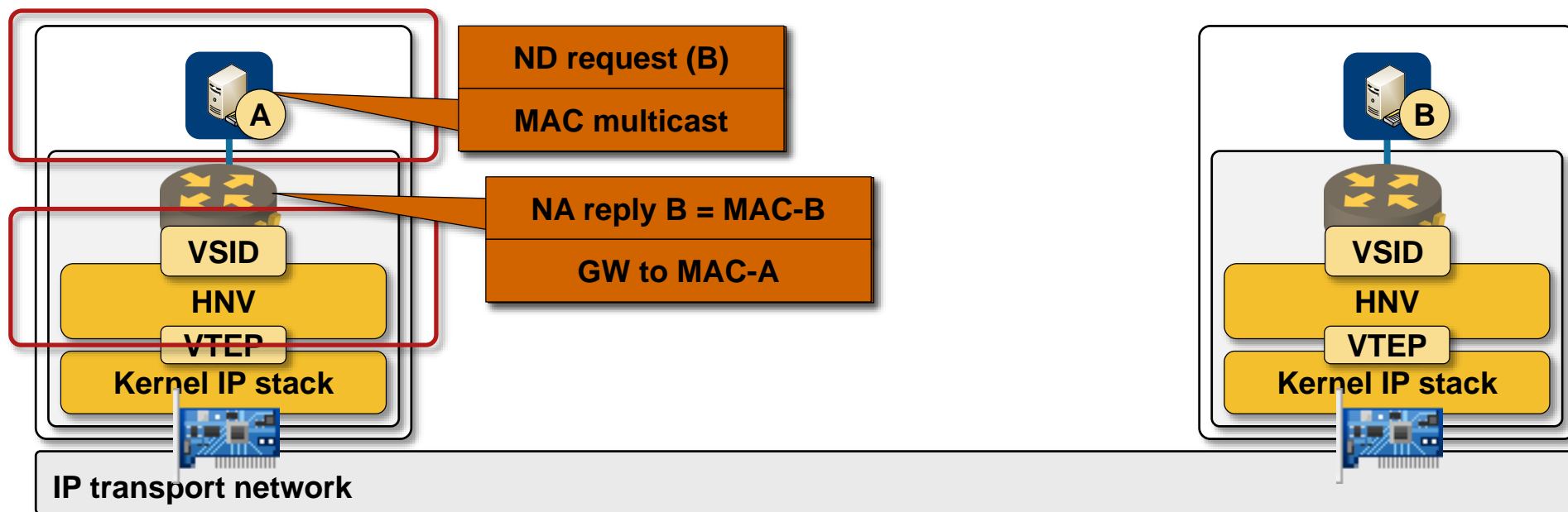
# Example: Hyper-V Network Virtualization



Full layer-3 switch in the hypervisor (distributed routing functionality)

- L3-only switching for intra-hypervisor and inter-hypervisor traffic
- IPv4 and IPv6 support in customer (virtual) and provider (transport) network
- ARP and ND proxies → no ARP or unknown unicast flooding
- Source node flooding or Customer → Provider IP multicast mapping

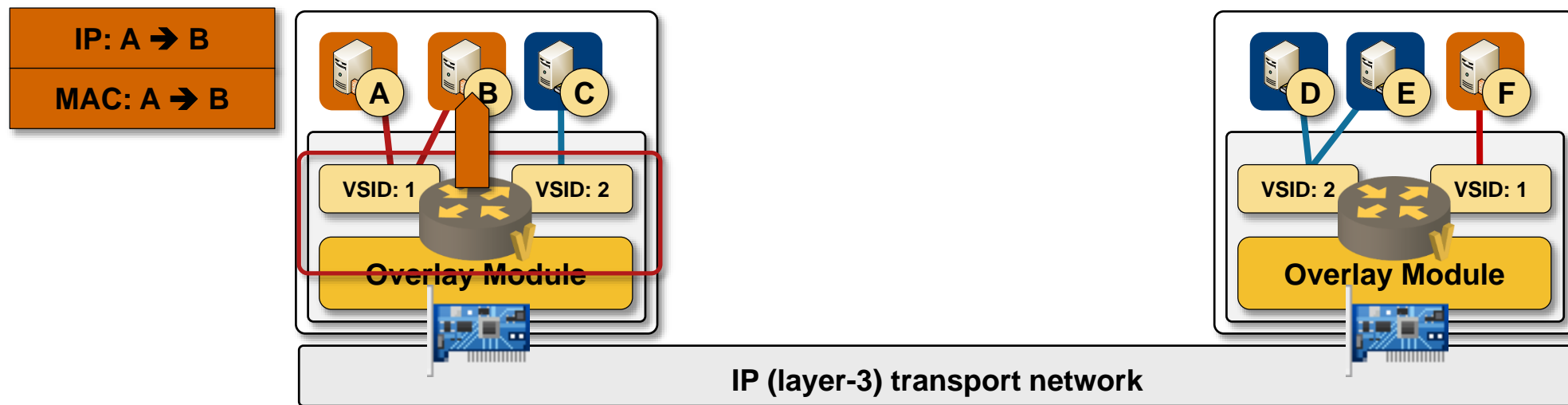
# Hyper-V Network Virtualization ND Proxy



- VM generates ND multicast
- L2 broadcast/multicast intercepted by Hyper-V kernel module
- Local Hyper-V replies to ND request with MAC address of remote VM
- Remote hypervisor is not involved
- Unicast ND requests are forwarded to target VM (NUD probes)



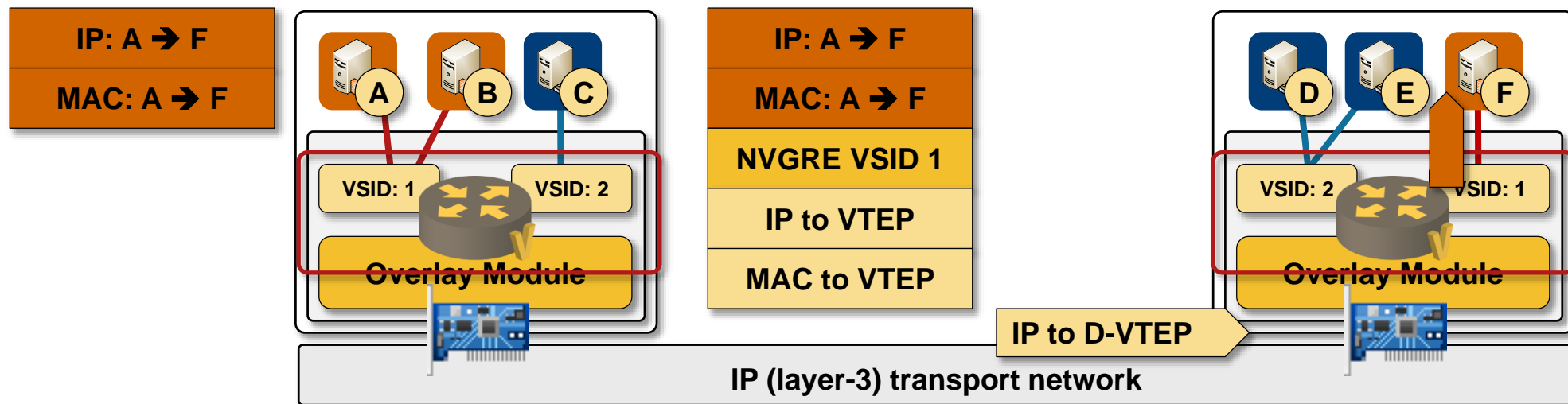
# HNV Local Switching



A → B

- On-link, sent directly to MAC-B
- L3 switched within the hypervisor (based on destination IPv6 address)
- IPv4, IPv6 and ARP packets are forwarded, all other traffic is dropped
- Ethernet frame delivered to target VM

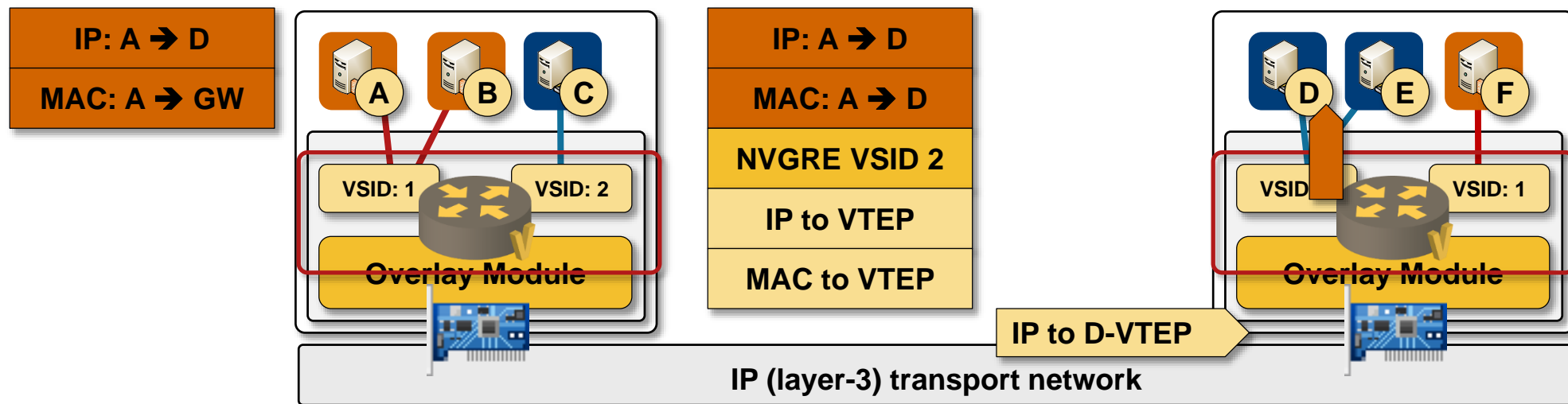
# HNV Remote Switching within a Subnet



A → F

- On-link, sent directly to MAC-F
- L3 switched within the hypervisor (based on destination IPv6 address)
- Destination VTEP is remote → build NVGRE envelope and send packet
- Packet received by remote hypervisor
- L3 switching within the routing domain (based on NVGRE VSID)
- Ethernet frame delivered to target VM

# HNV Remote Switching across Subnets

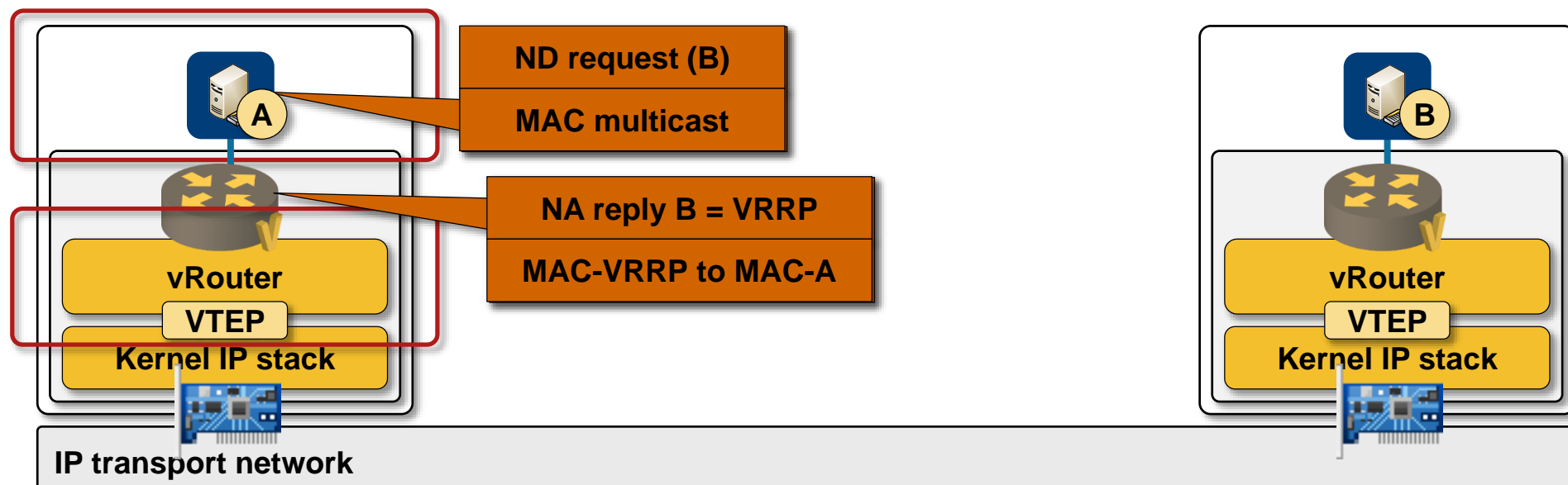


A → D

- Off-link, sent to GW MAC address
- L3 switched within the hypervisor (based on destination IPv6 address)
- Switching across subnets → MAC rewrite
- Destination VTEP is remote → build NVGRE envelope and send packet
- Packet received by remote hypervisor
- L3 switching within the routing domain (based on NVGRE VSID)
- Ethernet frame delivered to target VM

**HNV does not rewrite source MAC address or decrement TTL**

# Juniper Contrail ARP/ND Handling



- VM generates ARP broadcast or ND multicast
- ARP/ND requests (+ DNS and DHCP requests) are intercepted by local vRouter
- vRouter replies to all ARP/ND requests with VRRP MAC address
- Packet forwarding almost identical to Hyper-V case (no forwarding of unicast ND packets)

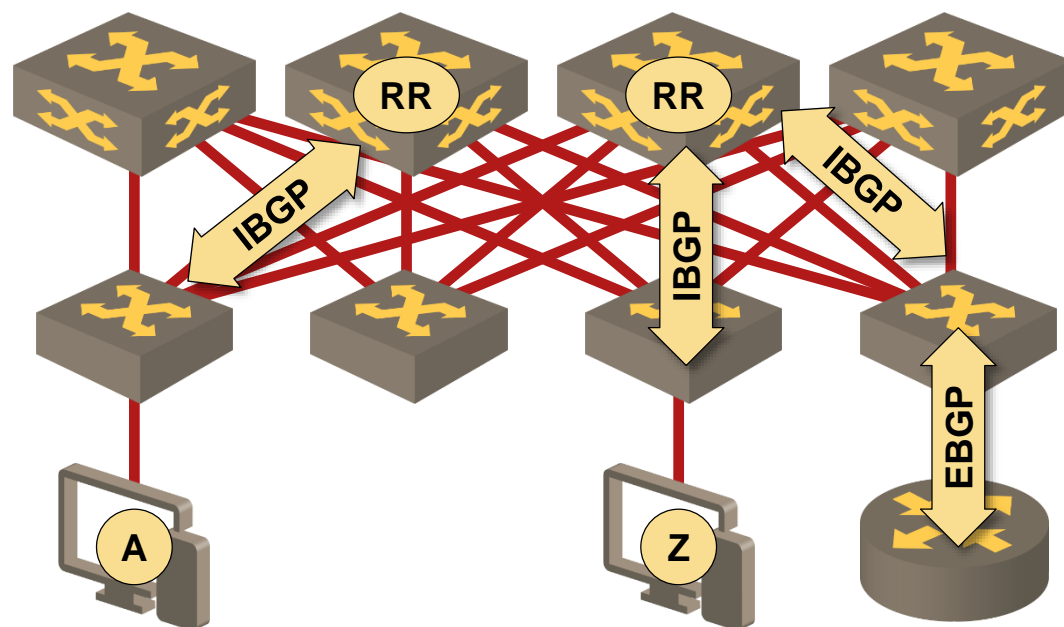
# IP Routing in Cisco Dynamic Fabric Automation (DFA)

## IP routing information distribution

- Host routes generated from ARP/ND/DHCP information or based on VDP messages (Nexus 1000v only)
- Subnet routes generated from configuration information
- External routes learned through routing protocols
- All IP routes inserted into MP-BGP and distributed across fabric

## Each fabric node knows

- All intra-fabric host routes
- All intra-fabric subnets
- All external routes



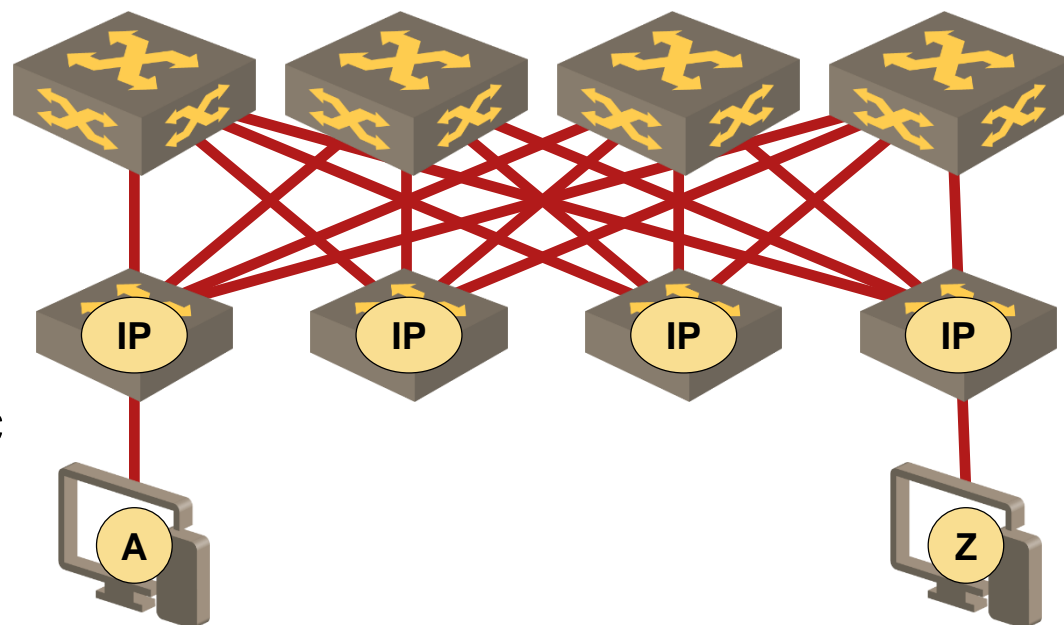
# Optimal Layer-3 Forwarding in Cisco DFA

All layer-3 leaf nodes share

- Default gateway IP address
- Default gateway MAC address
- All ARP/ND requests are answered with GW MAC address (proxy gateway mode)
- Integrates seamlessly with VM mobility

Typical packet forwarding

- Layer-3 lookup on ingress → egress next hop
- Layer-2 forwarding across fabric
- Layer-3 lookup on egress → delivered to destination



More in *Data Center Fabrics* webinar

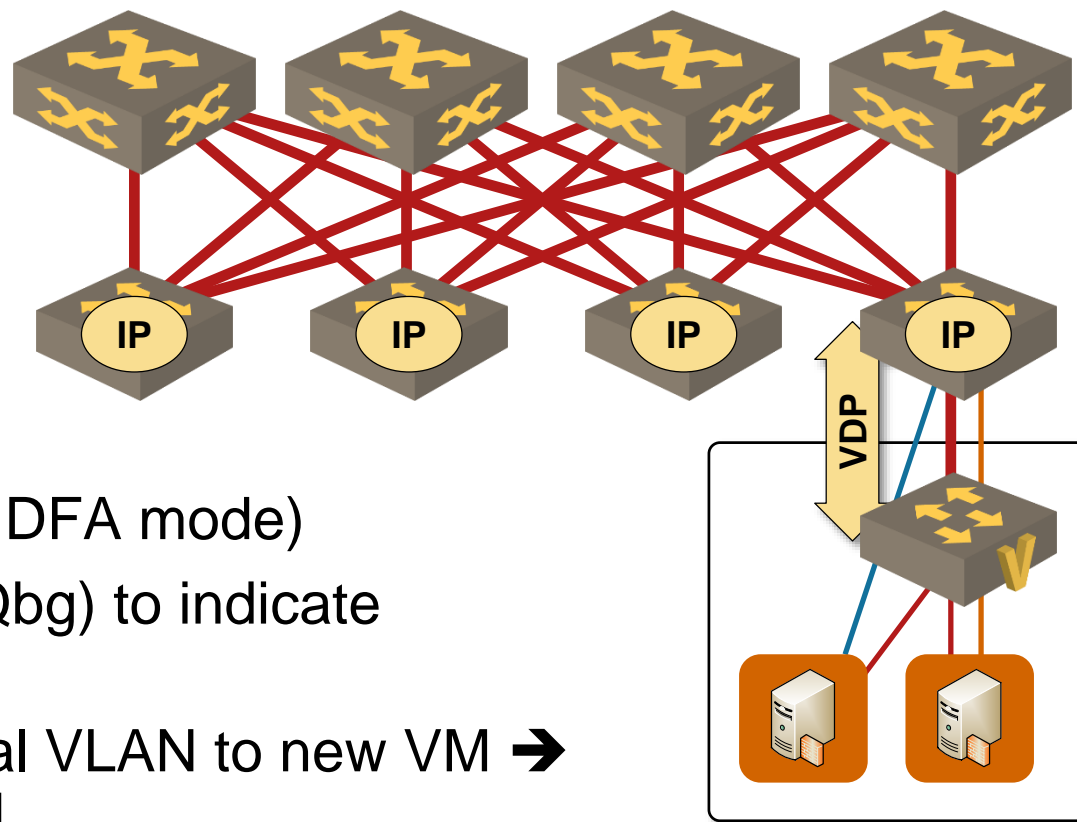
# Virtual Machine Microsegmentation with Cisco DFA

## Problem:

- Cisco DFA integrates with existing L2-only hypervisors
- Microsegmentation between virtual machines running on the same hypervisor is impossible

## Solution (requires Nexus 1000v in DFA mode)

- Nexus 1000v uses VDP (802.1Qbg) to indicate VM connectivity requirements
- DFA leaf assigns a dynamic local VLAN to new VM → each VM is in a dedicated VLAN
- L3 traffic is terminated at DFA leaf





# Summary

# IPv6 Microsegmentation Solutions

Why?

- Removes first-hop (L2) IPv6 security challenges

How?

- Dedicated dynamic interface per host (mobile, PPPoX)
- Dedicated VLAN per host (Carrier Ethernet, campus, data center)
- Host routing

# Implementations of Host Route-Based Forwarding

## IPv6 and IPv4

- Hyper-V Network Virtualization
- Juniper Contrail
- Cisco Dynamic Fabric Automation (DFA)

## IPv4 only

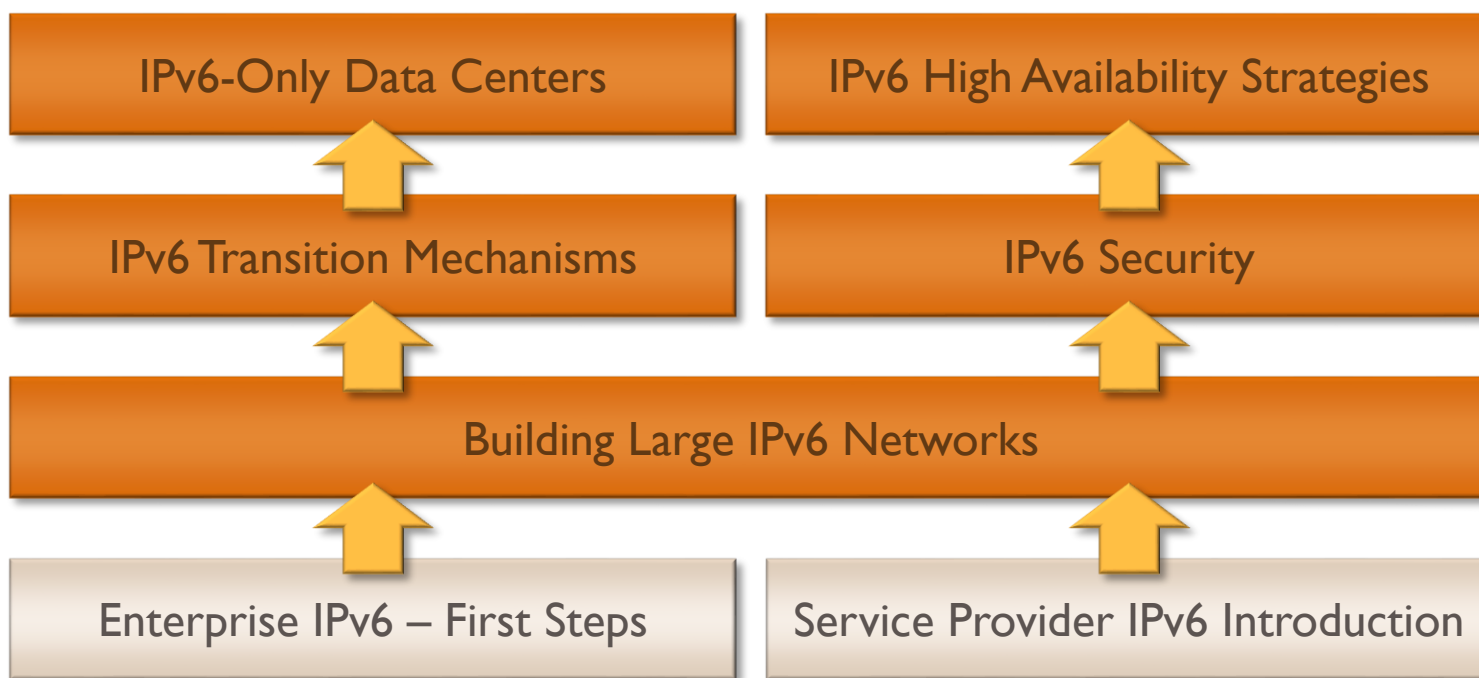
- Nuage Virtual Services Platform (VSP)
- Cisco Application Centric Infrastructure (ACI)

## Unrelated honorable mention

- IPv6 RA guard and ND inspection implemented on VMware NSX

**Hint: vote with your wallet!**

## More Information



### Availability

- Live sessions
- Recordings of individual webinars
- **Yearly subscription**

### Other options

- Customized webinars
- ExpertExpress
- On-site workshops



# Questions?

## Paperwork issues

- Follow-up email
- Please fill in the evaluation form
- Recording available within 24 hours
- PDF materials always available for download
- Discount for future webinars – register through [my.ipspace.net](http://my.ipspace.net)
- Upgrade to yearly subscription
- Please spread the word!

Send them to [ip@ipspace.net](mailto:ip@ipspace.net) or [@ioshints](https://twitter.com/ioshints)

## Stay in Touch

Web: [ipSpace.net](http://ipSpace.net)  
Blog: [blog.ipSpace.net](http://blog.ipSpace.net)  
Email: [ip@ipSpace.net](mailto:ip@ipSpace.net)  
Twitter: [@ioshints](https://twitter.com/ioshints)



SDN: [ipSpace.net/SDN](http://ipSpace.net/SDN)  
Webinars: [ipSpace.net/Webinars](http://ipSpace.net/Webinars)  
Consulting: [ipSpace.net/Consulting](http://ipSpace.net/Consulting)