

# Data Center Networks

---

Enzo Mingozzi

Professor @ University of Pisa

[enzo.mingozzi@unipi.it](mailto:enzo.mingozzi@unipi.it)

# Data Center Networks

I SERVIZI TIACI CHE  
VENGONO FORNITI SONO  
SERVITI WEB-



NON E' UNA COESA  
SCONTATA

Edge  
Forwarders

CO SWITCINGA DENTRO  
LA FABRICA E' FATTOA  
L'USCITA

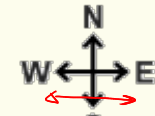
OUTSIDE  
(INTERNET, VPN)

WAN / DCI  
MPLS/IP Core

FWD  
A

FWD  
B

GATEWAYS



TRAFFICO  
WEST-EAST

RISORSE COMPUTAZIONALI  
SONO USATE DALLI CLIENT  
TRAMITE VIRTUALIZZAZIONE,  
MA SONO ALLEGATE AL  
DATA CENTER

Spine  
IP Switches

NETWORK

IP Fabric  
UNDERLAY

S0 ... SN

FORNISCIE UN  
SERVIZIO AD UN  
AGENZIA

NEI SERVER ABBIAMO  
APPLICAZIONI CHE  
FORNISCONO SERVIZI  
ATTIVO CUSTOMER  
NORTH-SOUTH  
TRAFFIC

Leaf  
IP Switches

L0 L1 L2 ..... LP LQ LR

VIRTUAL SWITCHES

Edge  
Forwarders

FWD  
0

FWD  
1

FWD  
2

FWD  
X

FWD  
Y

FWD  
Z

VM VM

CT CT

BM

CT CT

VM VM

BM

RISORSE VIRTUALIZZATE E  
FORNITE SOTTO FORMA  
DI CONTAINER O  
VIRTUAL MACHINE-

Legacy Hypervisor  
VM VM

BARE  
METAL

QUESTO HARDWARE E'  
DEDICATO O ADDEBITATO  
POSSESSO 'RACCOMANDA  
NON REVENIENT

Legacy Hypervisor  
VM VM

FORNISCIE UN  
SERVIZIO AD UN  
AGENZIA  
WEST-EAST  
TRAFFIC

SERVICES

SISTEMI LEGACY  
SENZA VIRTUAL  
SWITCHES

# Modern DCN requirements

- **Increased server-to-server communication**

- modern data center applications involve a lot of server-to-server communication (east-west)

- **Scale**

↳ MAGGIORE RISPETTO AL TRAFFICO NORTH-SOUTH (INPUT - OUTPUT)

MICROSERVICES, SERVERLESS -- AUMENTO DELLA CAPACITÀ DI ALLOCAZIONE DELLE RISORSE

↳ QUESTO HA FATTO CANTO ALLA DECISIONE DI PASSARE DA C2 A C3 FABRIC

- modern data centers range from a few hundred to a hundred thousand servers in a single physical location

↳ BISOGNA TENERE COSTO A DESIGN TIME

↳ CI SONO BISOGNO DI CENTINAIA DI MIGLIAIA DI PORTE (AREE FISICHE) PER COLLEGARE I SERVER

- **Resilience**

↳ QUANDO USIAMO CONTAINER OUR CLOUDS CRESCONO ASSAI

- The primary aim is to limit the effect of a failure to as small a footprint as possible

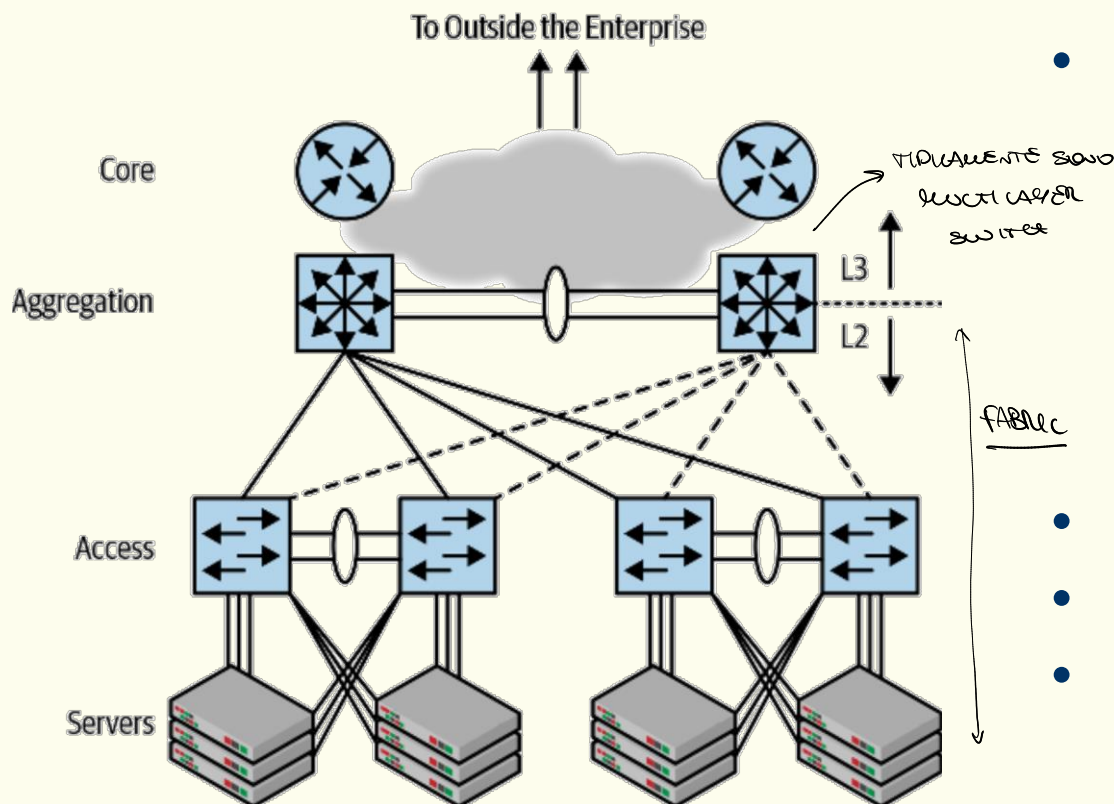
↳ DEFINE TOOL CHE PERMETTANO DI CRESTARE GLI ASSET A RONTINE

↳ CRESTARE E CONFINAMENTO GLI ASSET

# Legacy DCN topologies

## Access/Aggregation/Core network design

COSTURA CORE UNA RETE ARCSIMILE



- **Unscalability**

- Flooding
- VLAN limitations 12 BIT PER VLAN ID
- Burden of ARP
- Limitations of switches and STP → HAHA DECA POTENZA

- **Complexity**

- **Failure domain**

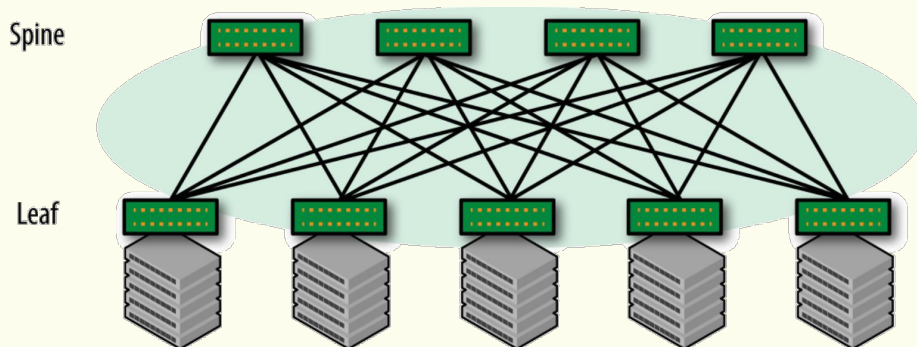
- **Unpredictability**

COMPUTAZIONE CHE NON  
PUOL USARE PER CORE  
E DENSITA' A RETE  
A LAYER 2.

# Modern DCNs

- The flexibility promised by bridging to run multiple upper-layer protocols is no longer needed. The only network-layer protocol that need be supported is IP!
- Modern DCNs are IP-based (**IP fabric**) with a **Leaf-and-Spine** topology

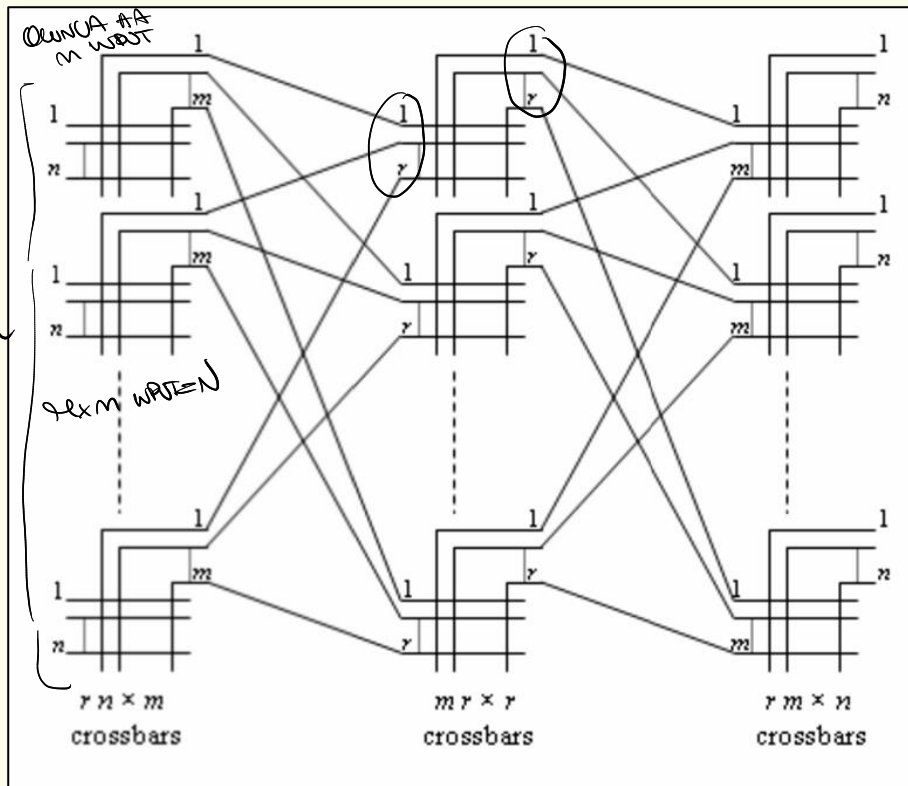
Quale è il problema?  
 Per essere usati hardware specializzati, con  
 software completamente proprietario.



# Clos networks



Originally invented by Charles Clos in the 1950s for old, circuit-switched, telephone networks



Charles Clos. "A Study of Non-blocking Switching Networks". Bell System Technical Journal, March 1953

**Scaling a switching matrix by decomposition:**

An  $N \times N$  matrix is realized by a multistage network made of smaller switching matrices organized into multiple layers or stages

RE E UNIDIREZIONALE → ABBIAMO WRT E COSTI SEPARATI E RENDIBILI CON UNA SOLA DIRECTIONE.

Let us consider a **3-stage network**

- Original matrix:  $N \times N$  → UN NUMERO DI RETTANGOLI DI SWITCHING SUL PRIMO PIANO
- Let us decompose  $N = r \times n$
- Input stage:  $r$  ( $n \times m$ ) switches
- Output stage:  $r$  ( $m \times n$ ) switches
- Intermediate stage:  $m$  ( $r \times r$ ) switches

NON HO UNO SWITCH CAPACE  $N \times N$ .

If the number of intermediate switches  $m$  is not sufficiently high, a blocking condition may occur

- **Non-blocking condition** (re-routing may be necessary):  $m \geq n$  LEVO GLI ALTRI CILINDRI E POSSO CAMBIARE IL CILINDRO.
- **Non-blocking condition without re-routing:**  $m \geq 2n - 1$  (Clos theorem)

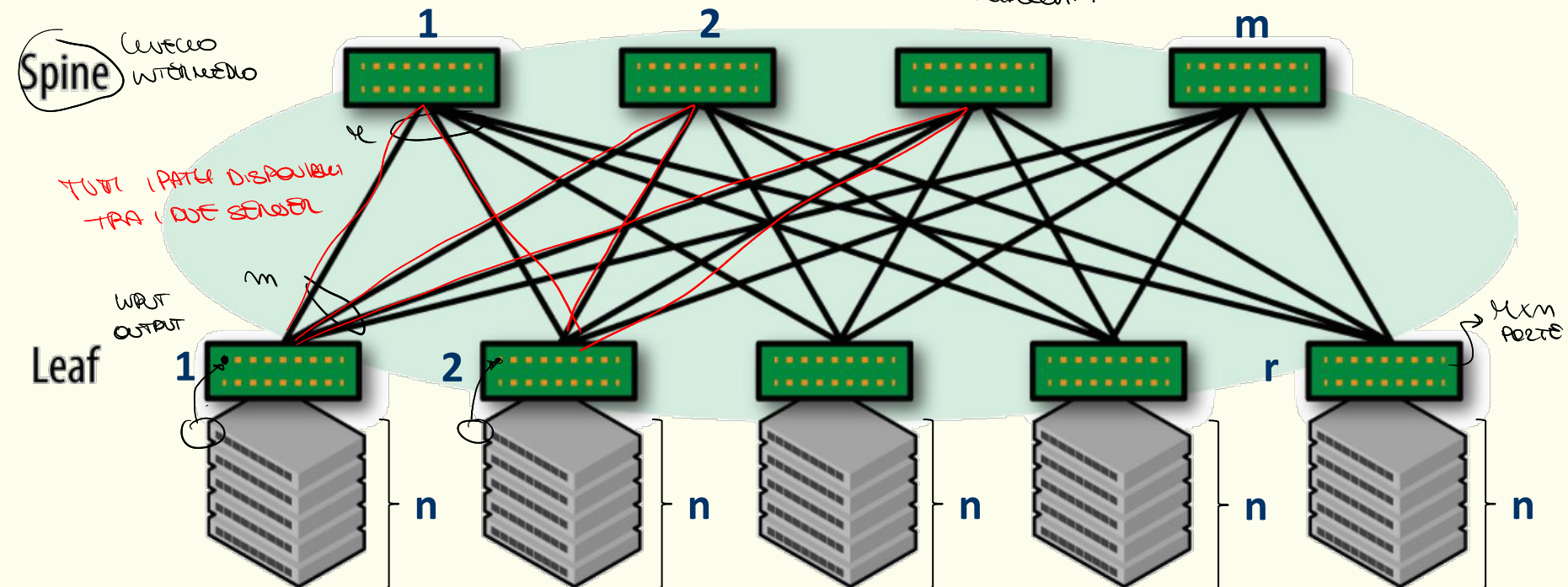


# Modern DCNs



- A simple Leaf-and-Spine topology is a folded (three-stage) Clos network

VEDADO 2 STADIUM ABBIATO UNA CLOS NETWORK A 3 STADI  
RUBEN A



$N = r \times n$  servers connected

# Modern DCNs

## Advantages of Leaf-and-Spine topologies

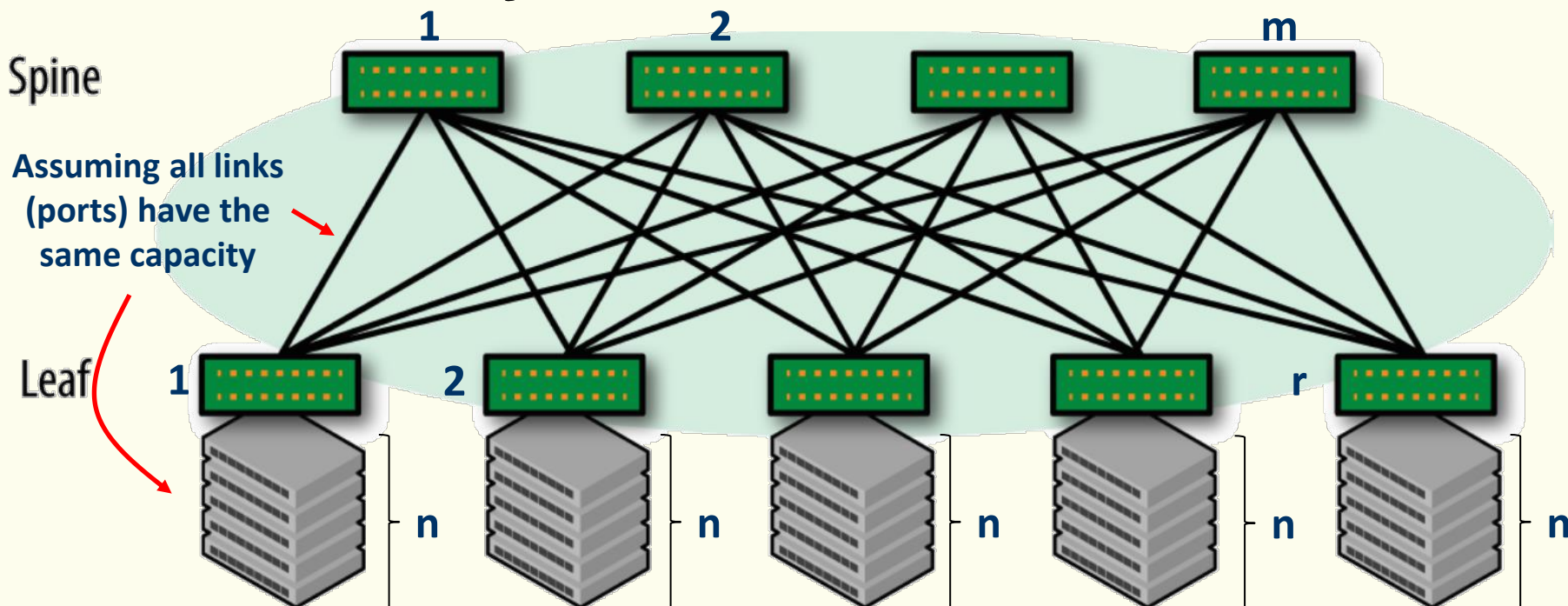
- **Simple and scalable**
  - Scale-out: More access ports → more leafs (if spines have sufficient downlink ports)
  - Scale-up: More capacity → more spines (and bandwidth on links)
- **Better load distribution:**  $m$  equal paths between each two leaves → ECMP for east-west traffic
- **Lower CAPEX and OPEX**
  - Economy of scale: switches with fixed configuration
  - Ease of configuration



# Capacity of the DCN

- How many access ports can be available (non-blocking)?  $N = r \times n$
- $R$  ports on Spine  $\rightarrow r = R \rightarrow$  numero massimo di server su rete
- $K$  ports on Leaves ( $n + m = K$ )  $\rightarrow n = m = K/2$    
 server SPWF

$$N_{\max} = R \times K / 2$$



$N = r \times n$  servers connected

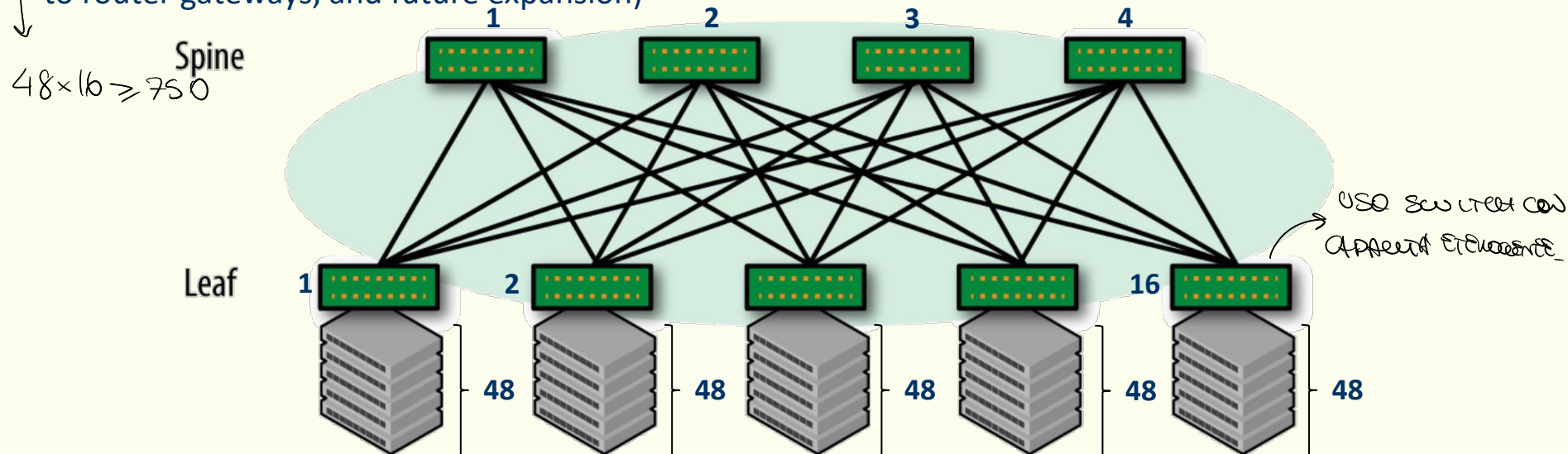
**Oversubscription ratio:** the ratio between the overall bandwidth server-side (access) and spine-side on a Leaf switch

$$48 \times 10 = 480$$
$$4 \times 40 = 160$$

PSN Core Case COS Cell SPWE\_ VAGNE  
Hides

- valone  
tires

**750 access ports required** → 4 Spine switches (using 16 ports per switch, 8 are left for interconnection to router gateways, and future expansion)



**N = 768=48×16 access ports**

# Capacity of the DCN

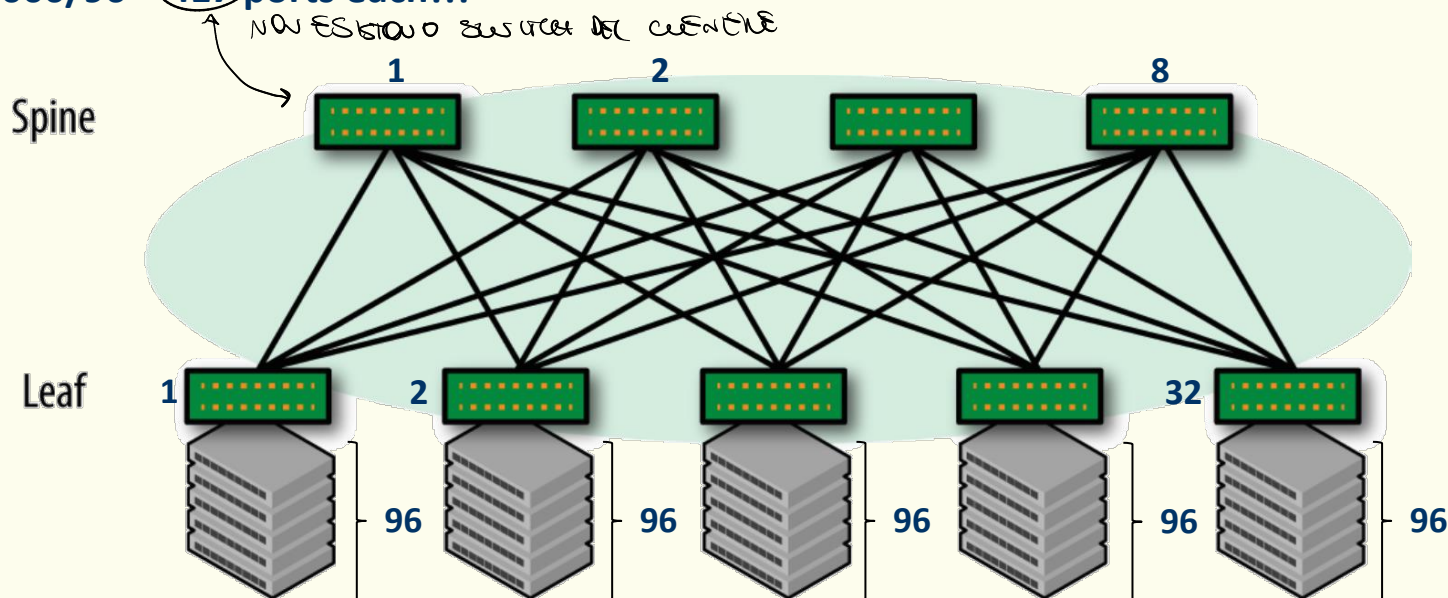
Another example (realistic figures considering switches available on the market):

- Leaf switches with **96 10Gb/s** access ports + **8 40Gb/s** uplink ports → oversubscription ratio **3:1**
- Spine switches with **32 40 Gb/s** ports

[Today, a 25Gb/s access link coupled with a 100Gb/s uplink is becoming the trend]

How many access ports are available at most?  $N = 96 \times 32 = 3.072$  *2 = 536 servers*

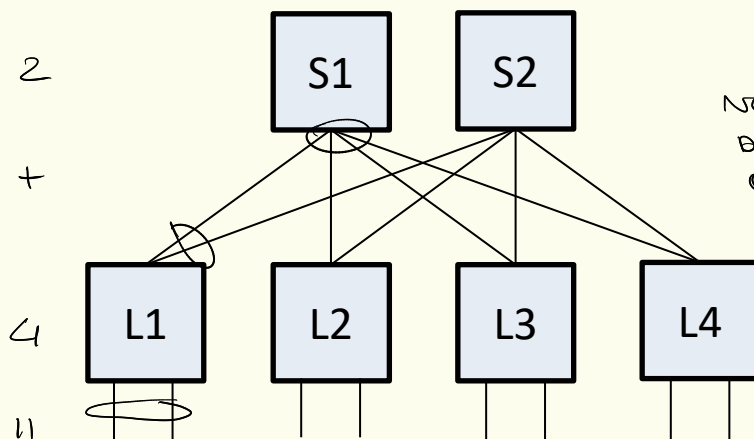
What if I need to accommodate **40.000** servers with 10 Gb/s access ports? Spine switches should have  $40.000/96 = 417$  ports each!!!



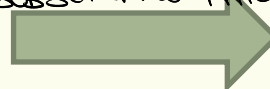
# Scaling Clos networks

Example: two-tier (three-stage folded) Clos network

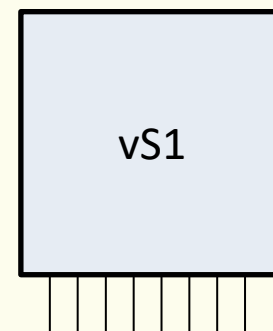
- 6 four-port switches
- 1:1 oversubscription



Now introduce bottleneck  
by removing the  
oversubscription



8-port vSpine switch



16 Switch

In general

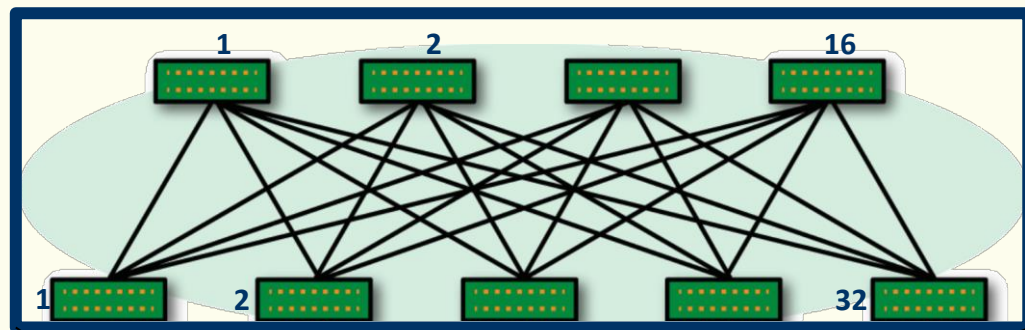
- $(N+N/2)$  N-port switches
- 1:1 oversubscription



$(N^2/2)$ -port vSpine  
switch



# Scaling Clos networks

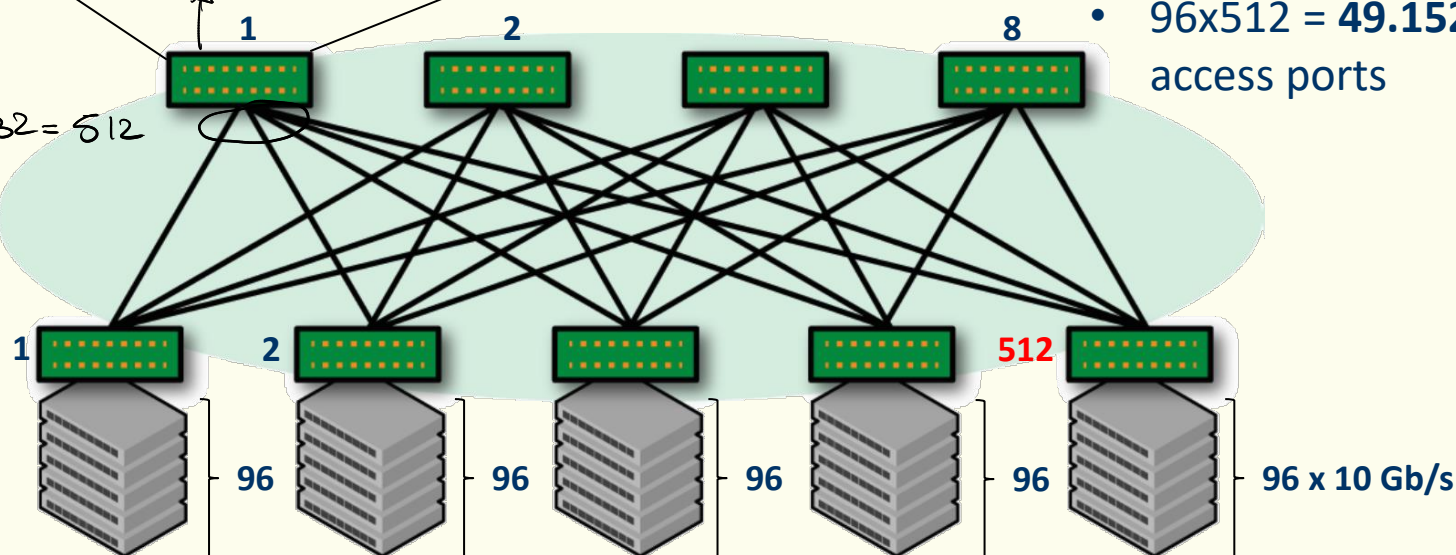


with a single required  
cost

Spine

$$16 \times 32 = 512$$

Leaf



## 3-tier (5-stage folded) Clos network

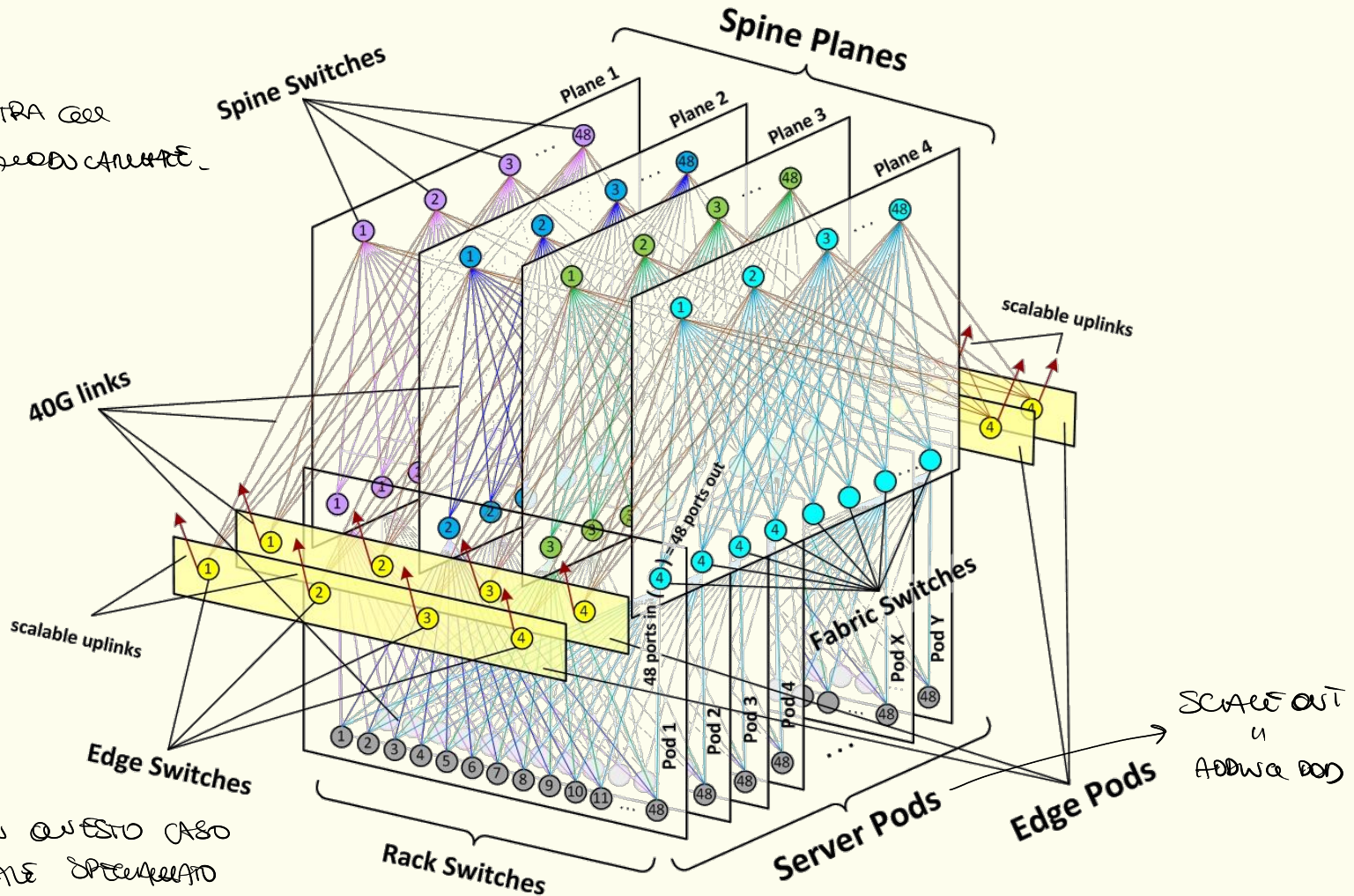
Example

- Leaf: **96 10Gb/s** access + **8 40Gb/s** uplink
- **3:1** oversubscription
- Spine: **32 40Gb/s** ports
- $96 \times 512 = \mathbf{49.152}$  **10Gb/s** access ports

# Data center networks



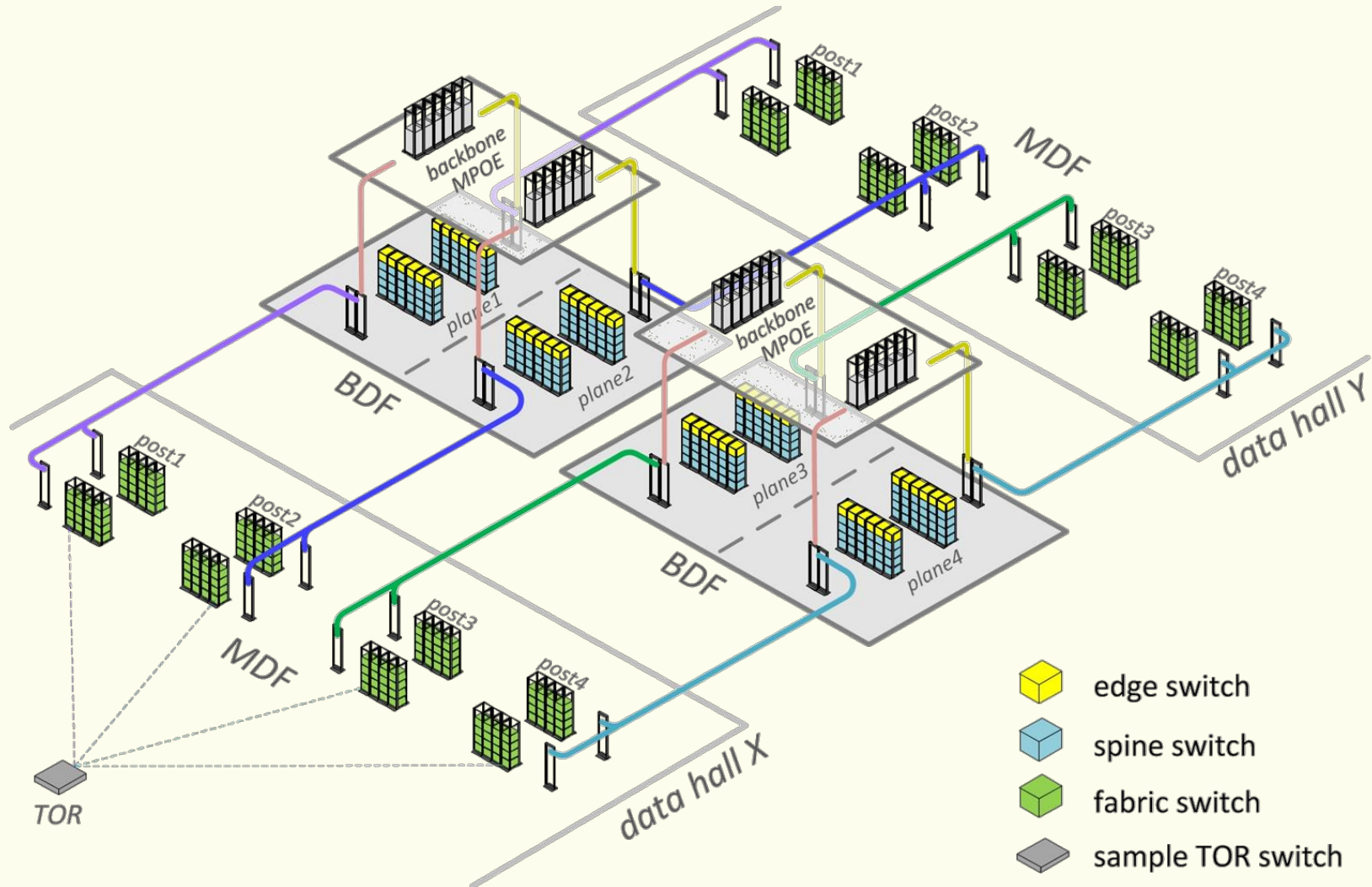
LE CONNESSIONI TRA GLI  
SWS NON SONO PUNTO A PUNTO



GLI SWITCHES IN QUESTO CASO  
SONO HARDWARE SPECIFICI



# Data center networks

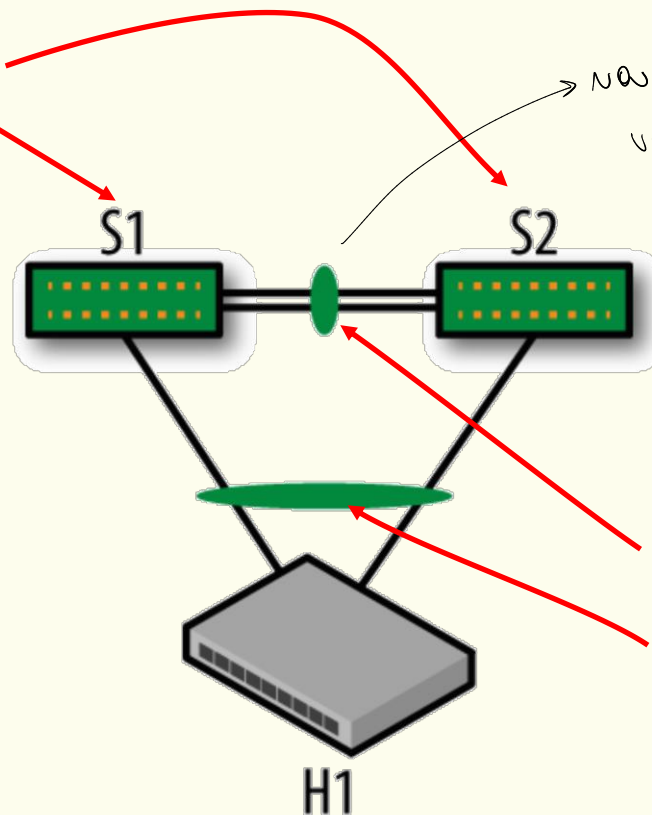


[Introducing data center fabric, the next-generation Facebook data center network - Facebook Engineering \(fb.com\)](#)

# Server attach models

- Single-attach vs. dual-attach server

To simplify cabling and facilitate rack mobility, usually both ToRs are in the same rack



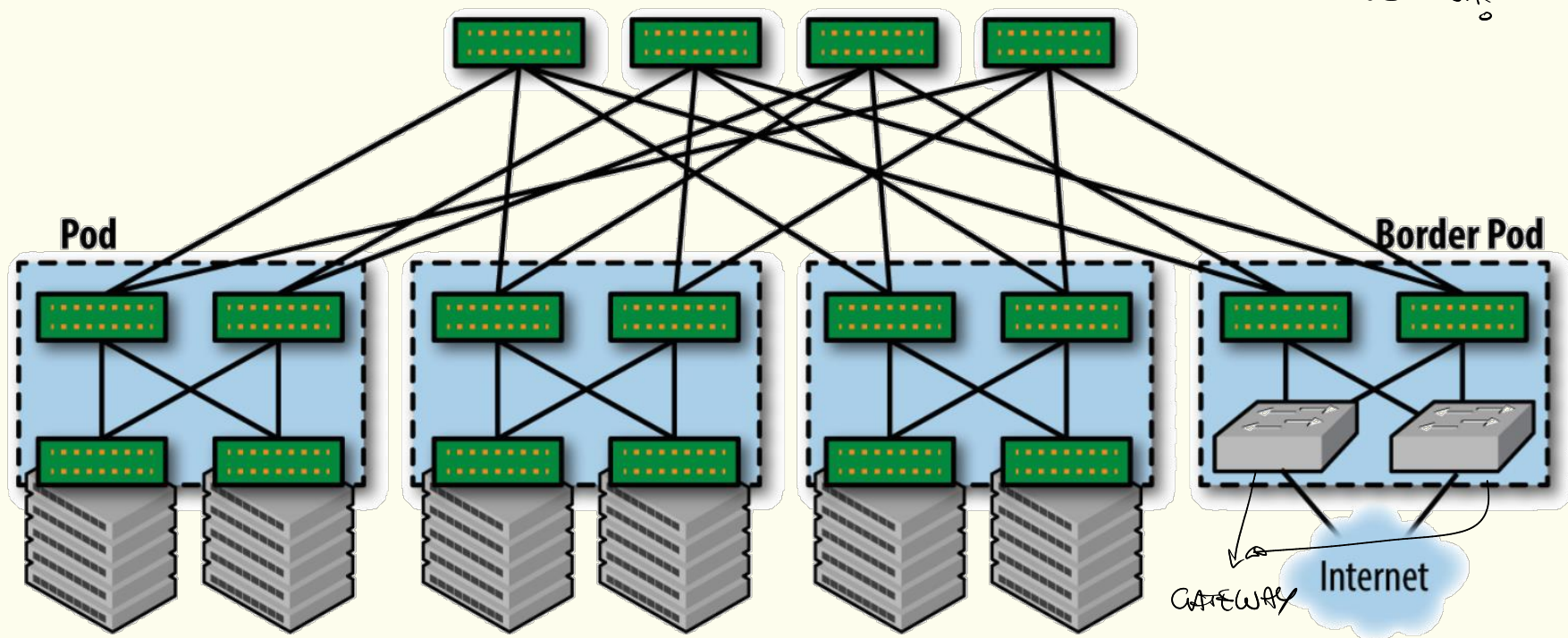
Now we have a single logical link between S1 and S2. This requires [vendor-proprietary protocols + LACP (Link Aggregation Control Protocol)]

dual links are aggregated into a **single logical link**  
This requires  
[vendor-proprietary protocols +  
**LACP** (Link Aggregation Control Protocol)

# External connectivity

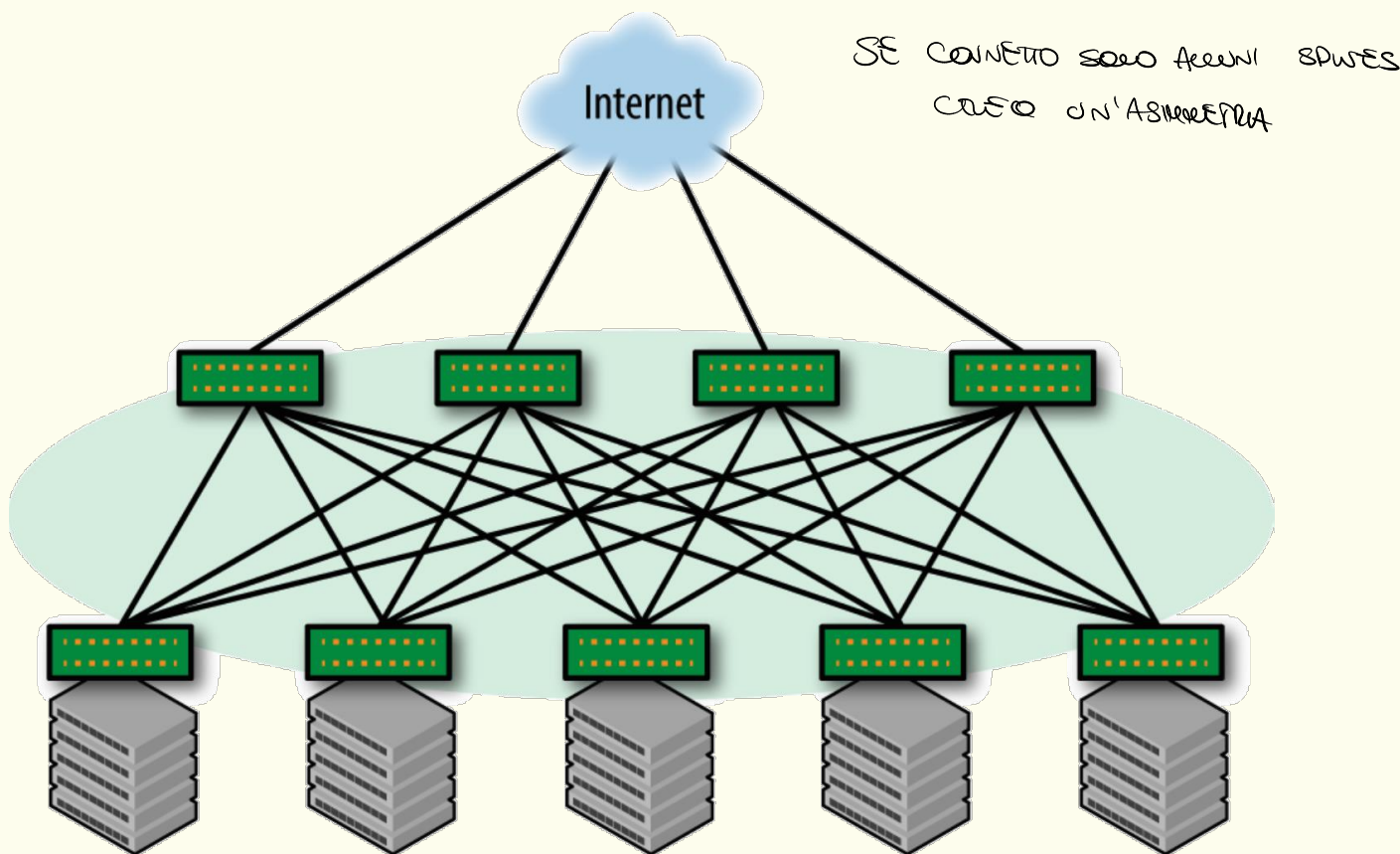
- Via *border* ToRs or pods

COME CONNETTAMO, SENZA  
UNA RETE ESTERNA?



# External connectivity

- Via Spine switches



# DCN Underlay – IP fabric

- Data plane: **IP**
  - In some cases, **MPLS** may be available
- Control plane
  - Distributed
    - **IGP**: OSPF or IS-IS
    - IGP-free: **eBGP**
  - Centralized/Hybrid
    - **SDN**

Come è implementato e come dipende dalla stack!

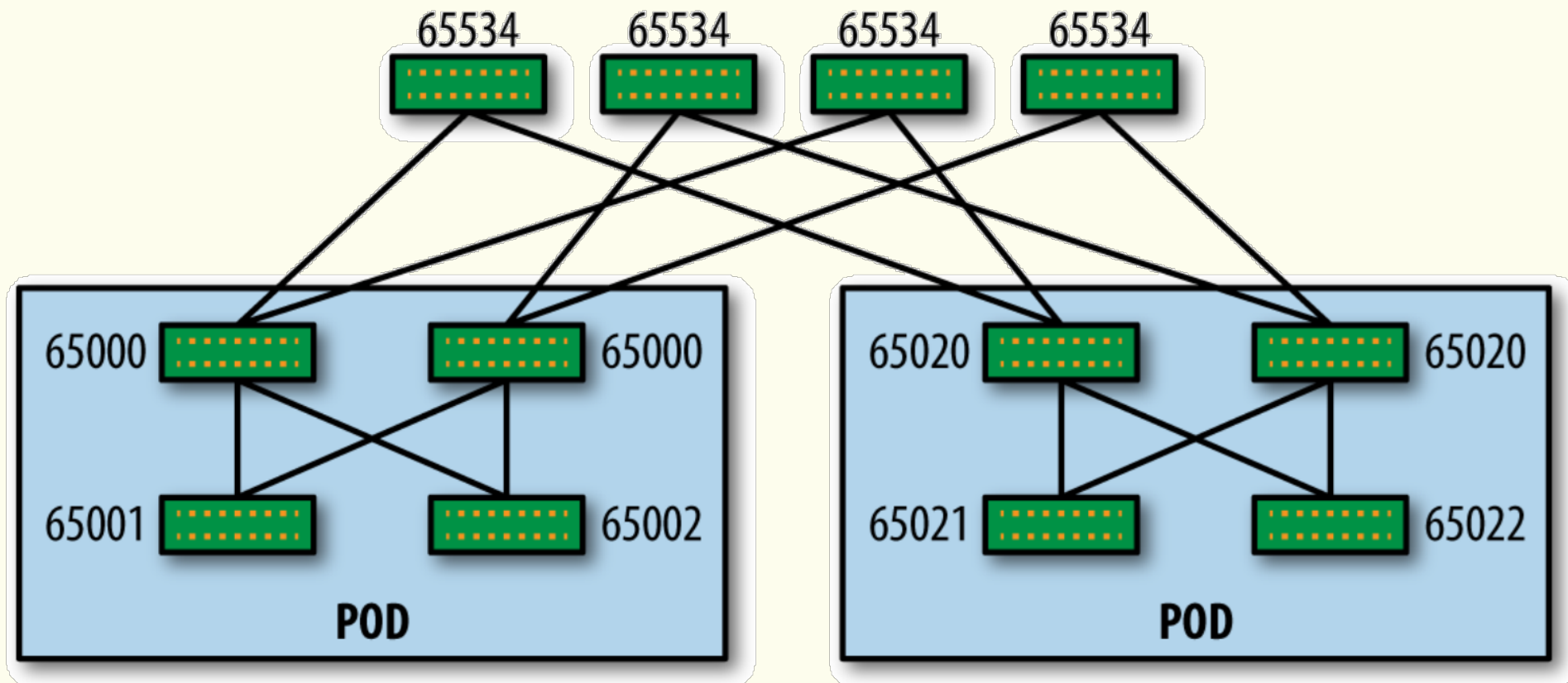
EXTERNAL BGP, così l'esterno è un altro network system -  
 Quando i router si parlano si trattano come AS diversi -

↓  
 Posso configurare tutto  
 centralmente



# DCN Underlay – IP fabric

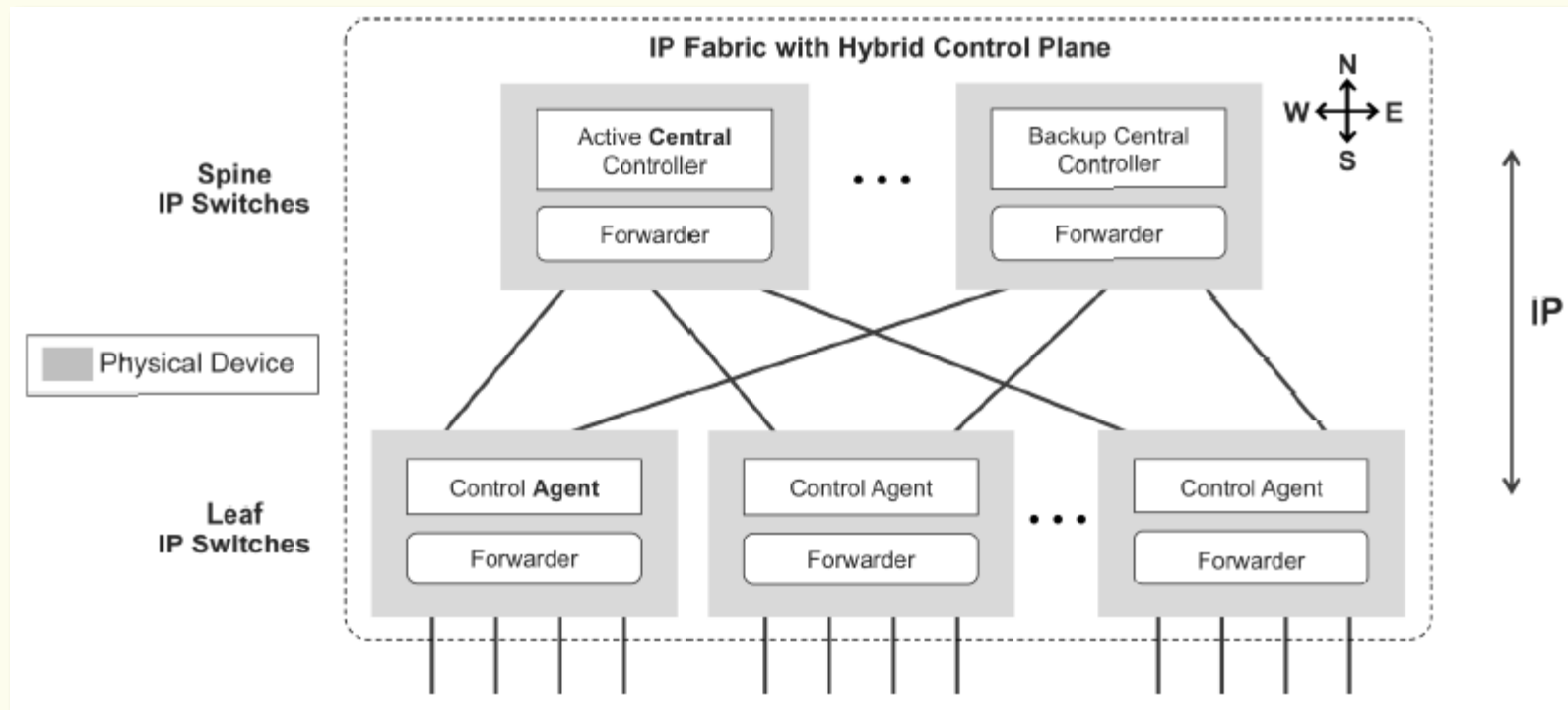
- Distributed control plane: eBGP





# DCN Underlay – IP fabric

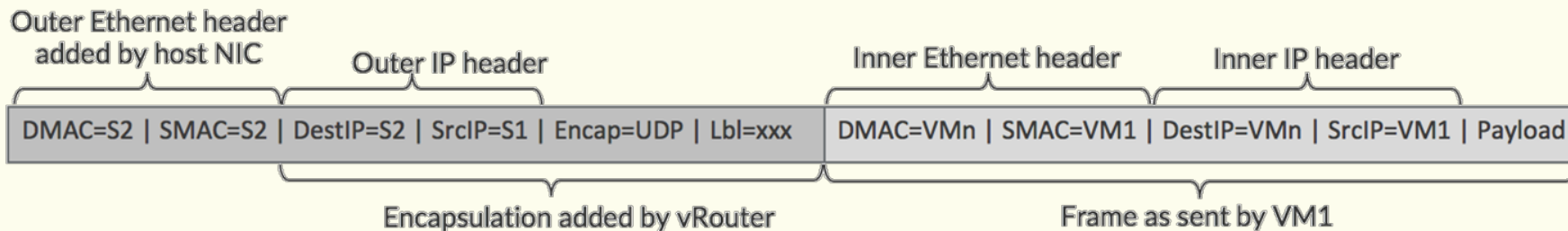
- Centralized/Hybrid control plane: either proprietary or use non-standard protocol extensions



# Network Virtualization Overlay



- Data plane: **L2**
  - Ethernet frames tunneled over the IP fabric (**VXLAN**, MPLSoUDP, MPLSoGRE, NVGRE, STT, ...)



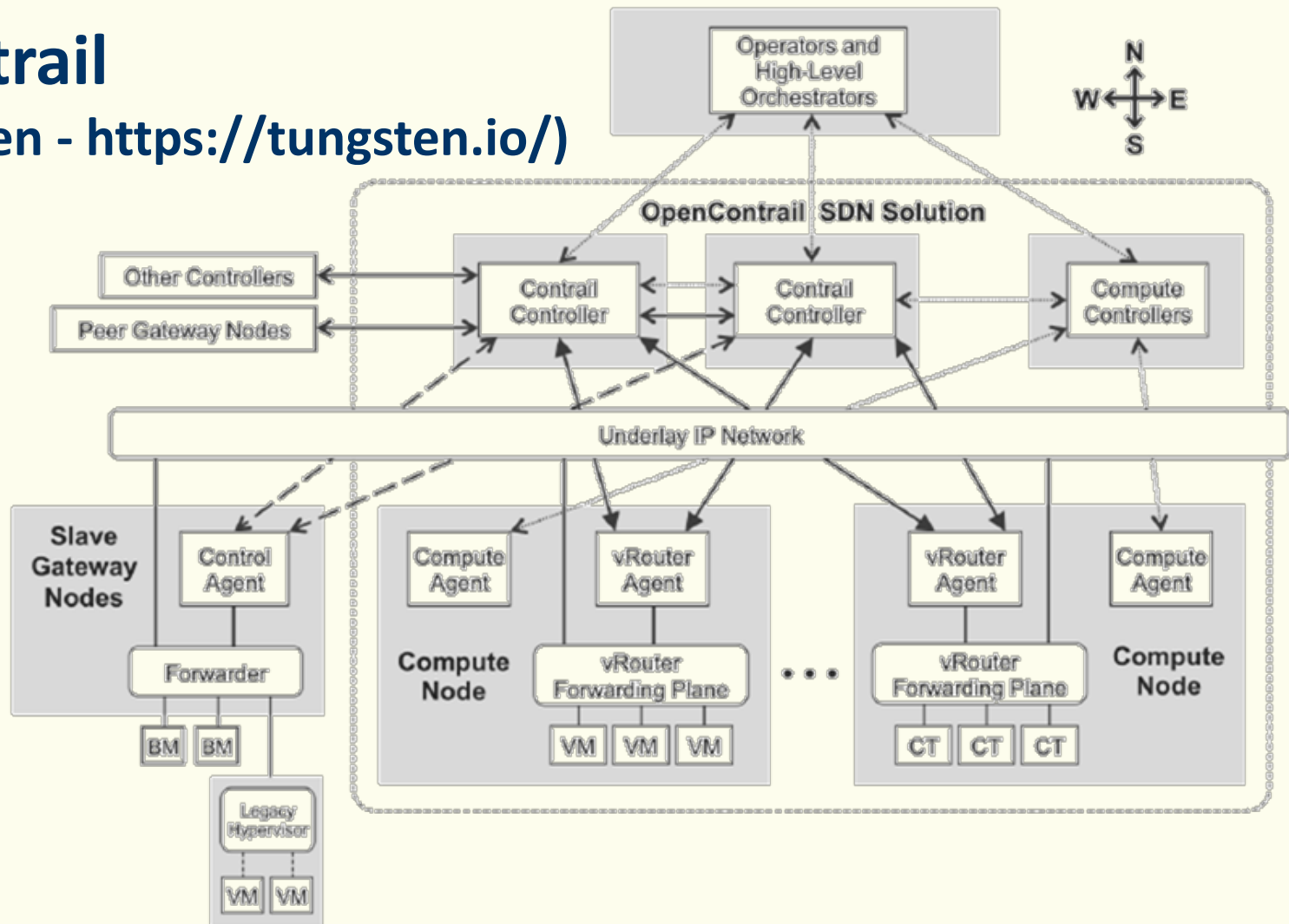
- Control plane: **controller**
  - Centralized: **SDN-based**
  - Protocol-based: **VXLAN + EVPN**

# Network Virtualization Overlay



## OpenContrail

(now Tungsten - <https://tungsten.io/>)



# References

- Dinesh G. Dutt, **Cloud Native Data Center Networking: Architecture, Protocols, and Tools**  
1st ed., O'Reilly, Dec. 2019
- RFC 7938 - **Use of BGP for Routing in Large-Scale Data Centers**
- RFC 8365 - **A Network Virtualization Overlay Solution Using Ethernet VPN (EVPN)**