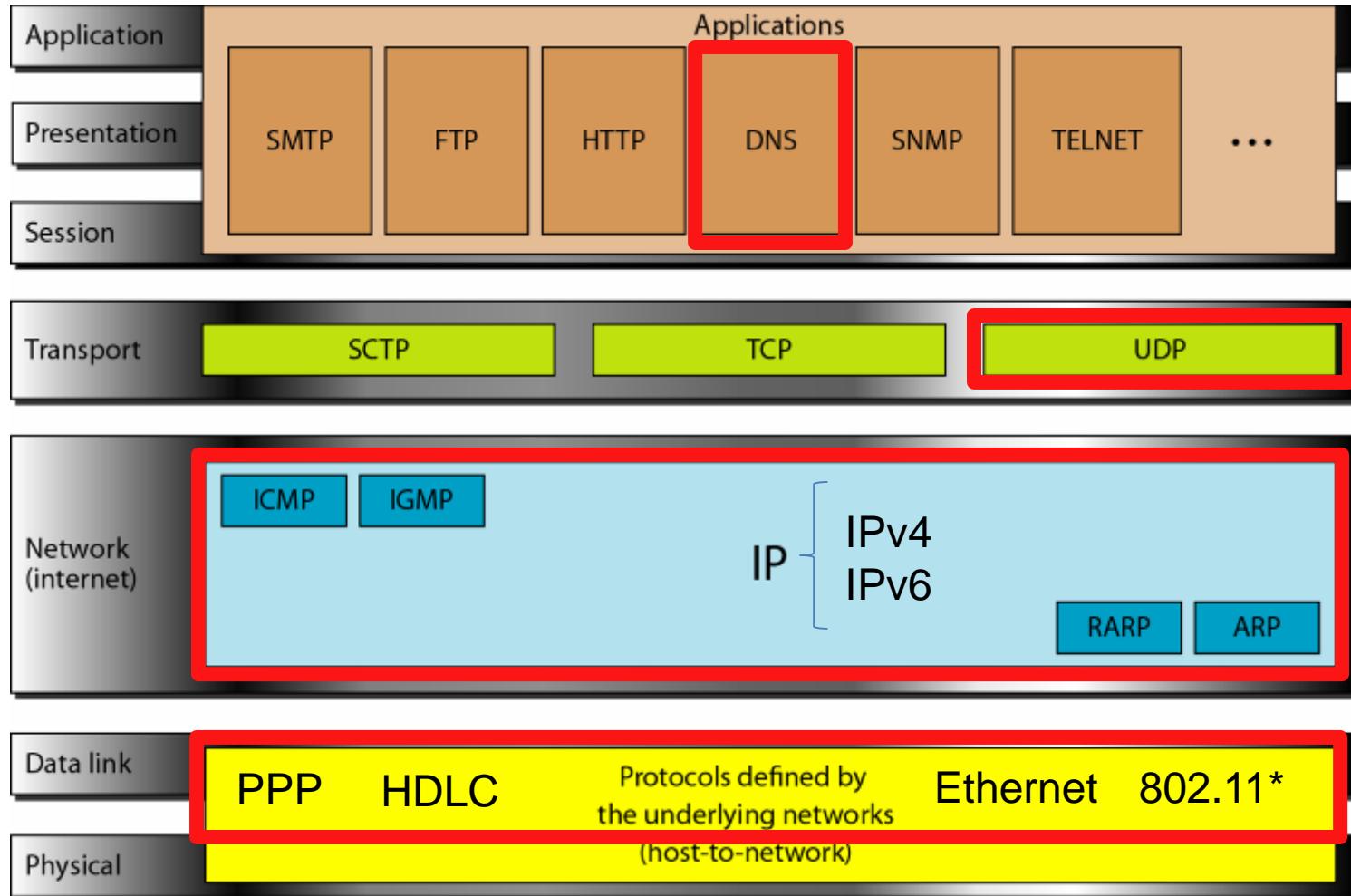


Overview

- Link Layer
- Network Layer
 - Addressing
 - Address Resolution (ARP)
 - Fragmentation
 - Intra-AS Routing
 - Distance Vector
 - Link State
 - Multicast Routing
 - IPv6
- Transport Layer
 - UDP
 - DNS



Protocols in the OSI Model



URLs to Names to Addresses

URL

DNS

IP Address

`http://www.wiki.com/index.html`

`www.wiki.com`

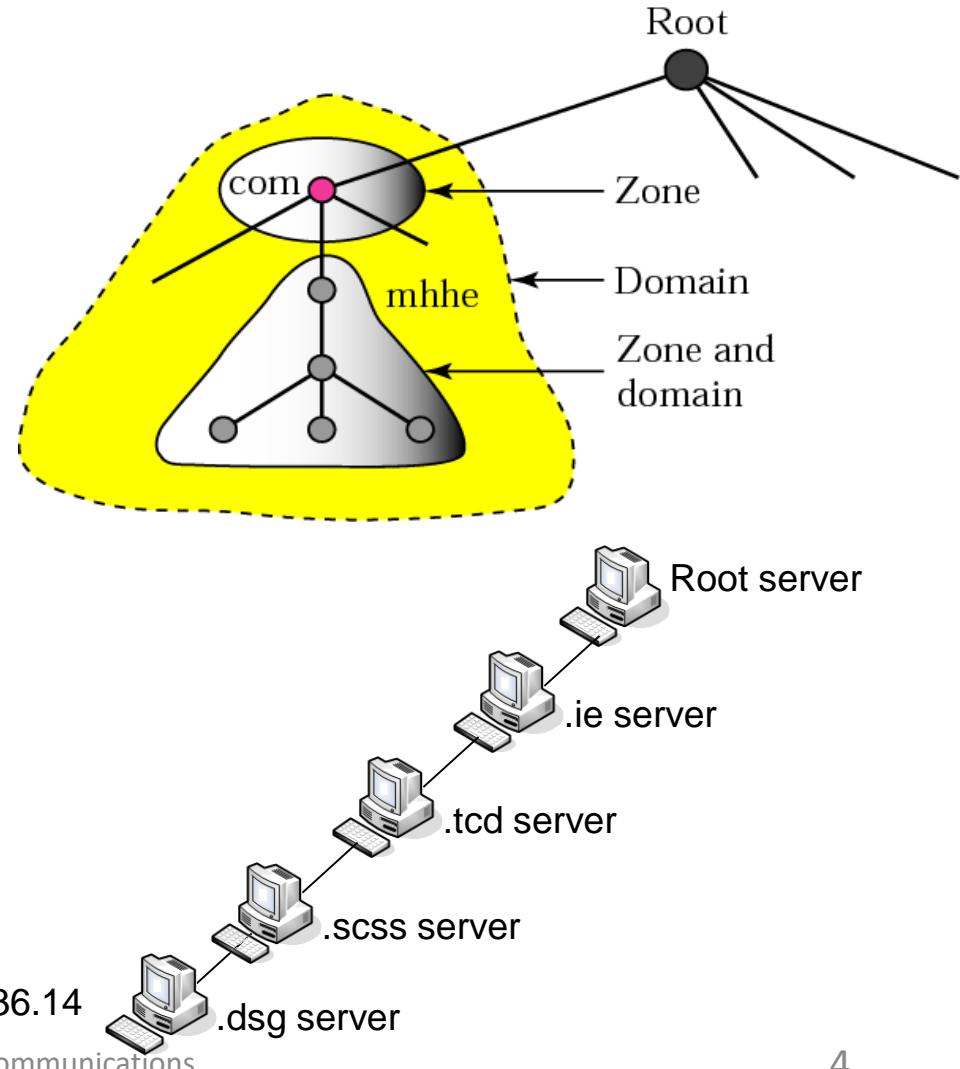
66.96.149.1

*URL = Uniform Resource Locator

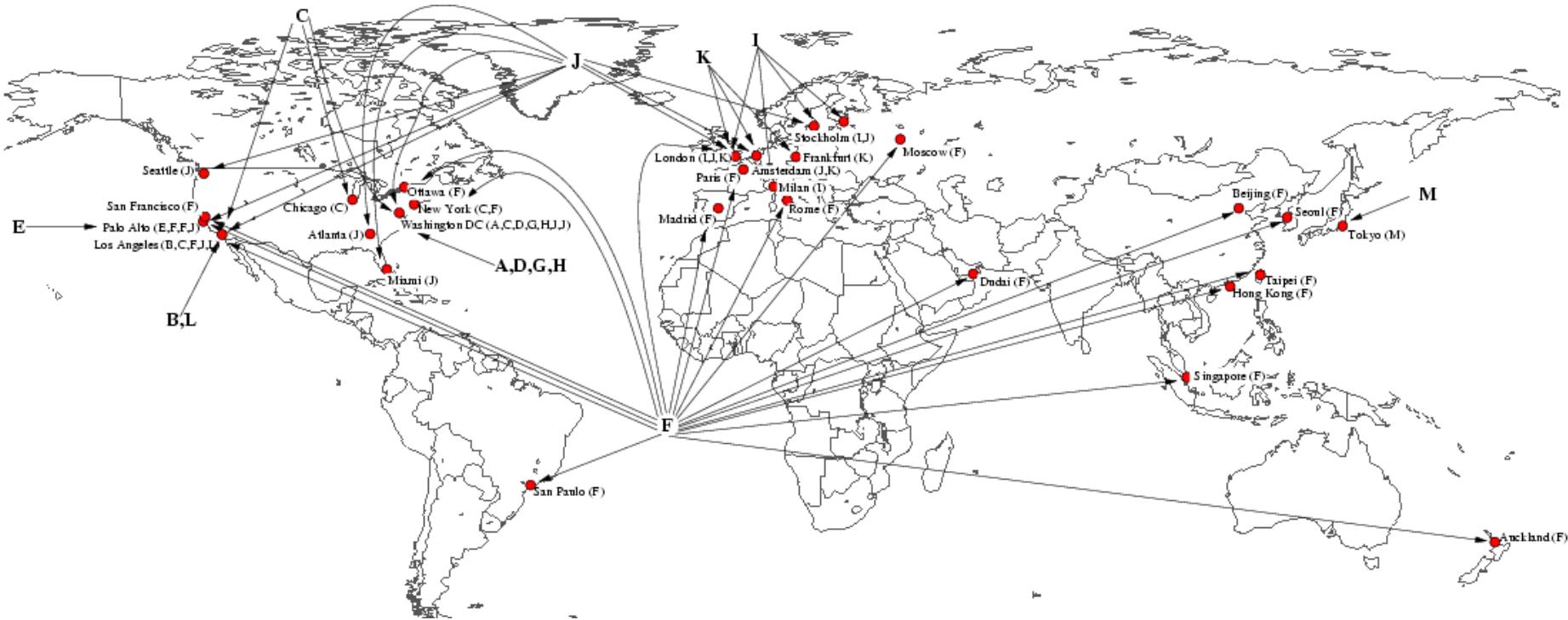


Hierarchy of Name Servers

- Every zone has a DNS server
- DNS server maintain lists of
 - Nodes in the zone
 - References to servers of zones underneath it

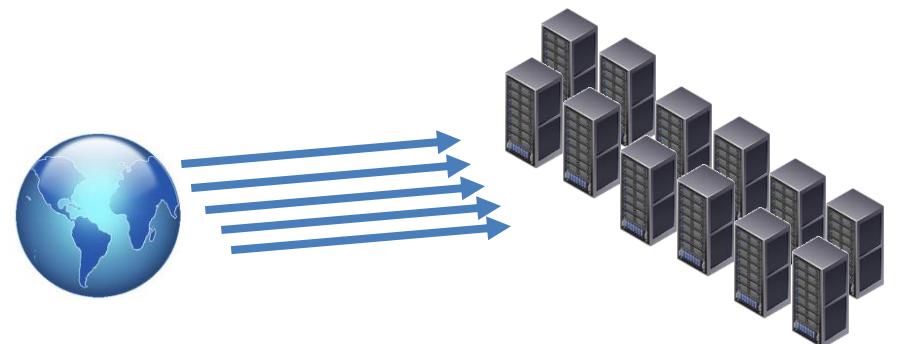
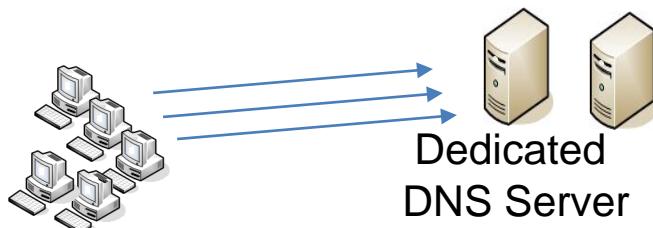


DNS Root Servers - Anycast



Mirai Bot & The Attack on Dyn

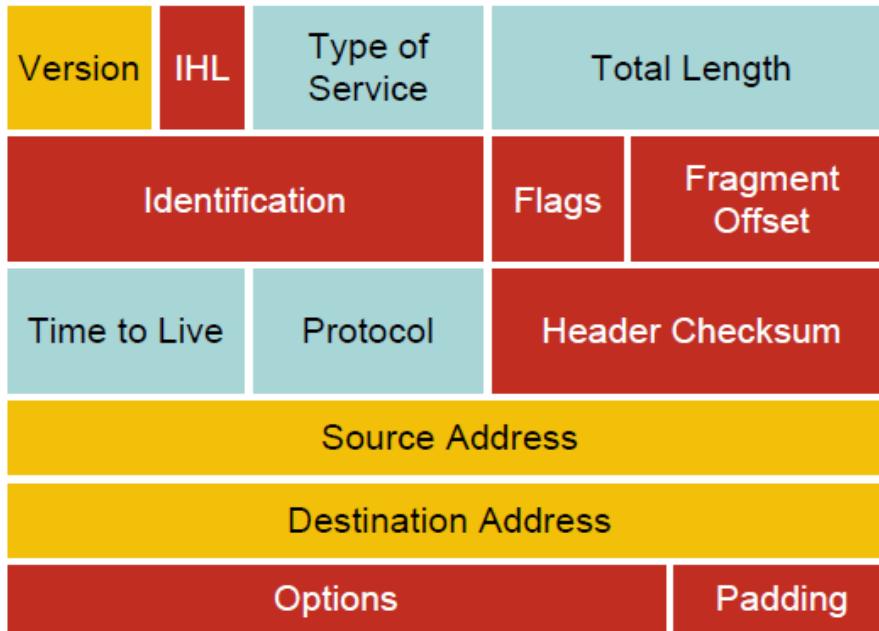
- Good example for scalable DNS services
- ... and the vulnerability of these services



Cloud/Datacentre-based
DNS Services

Header Comparison

IPv4 Header



IPv6 Header

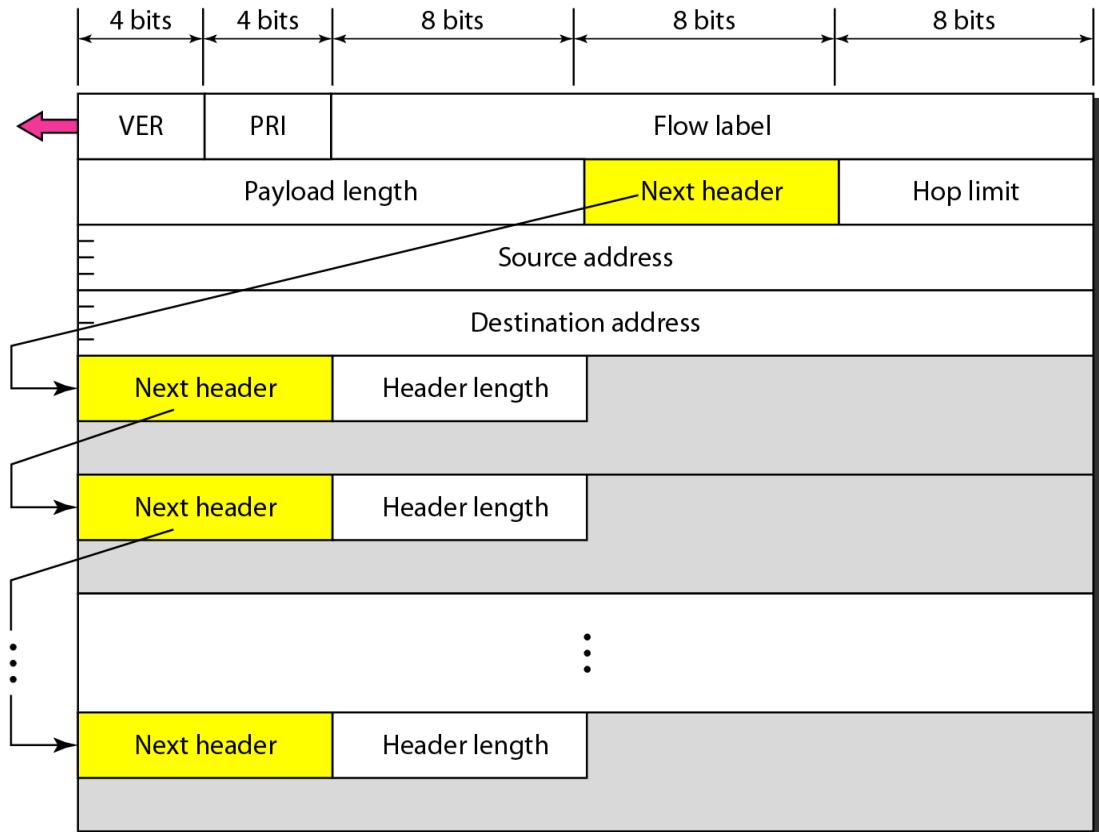


Legend

- Field's name kept from IPv4 to IPv6
- Fields not kept in IPv6
- Name and position changed in IPv6
- New field in IPv6

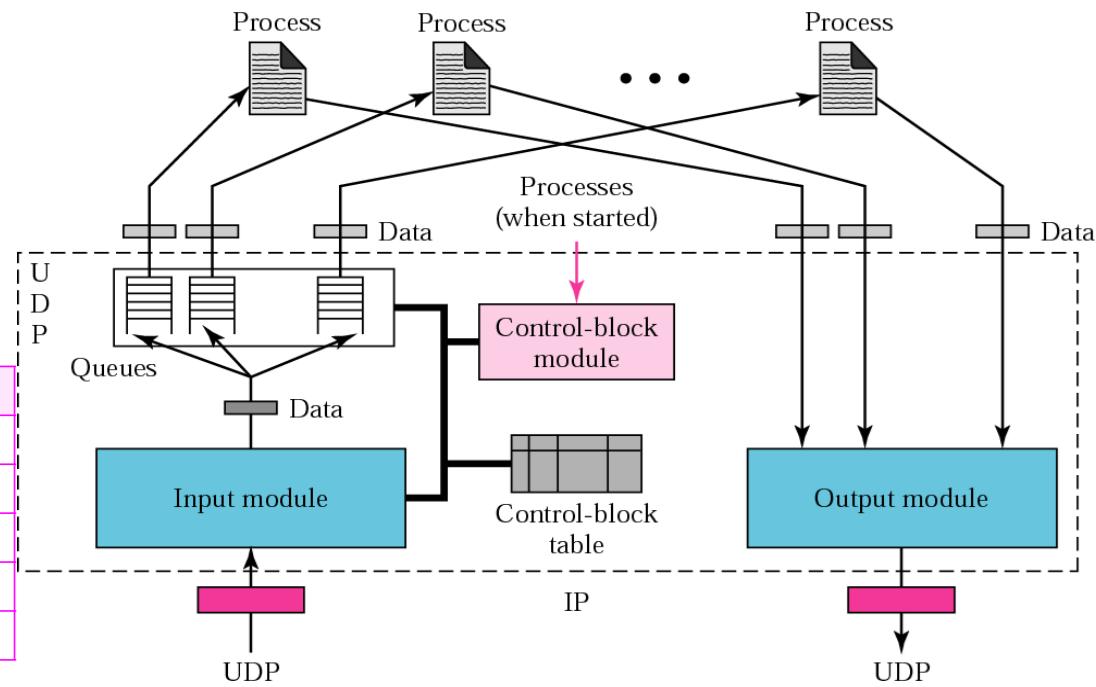
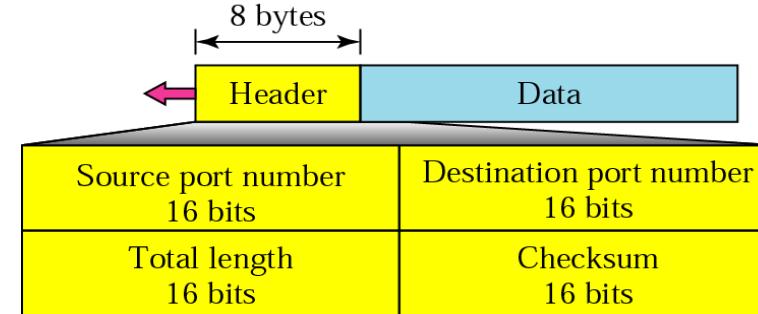


Extension Headers



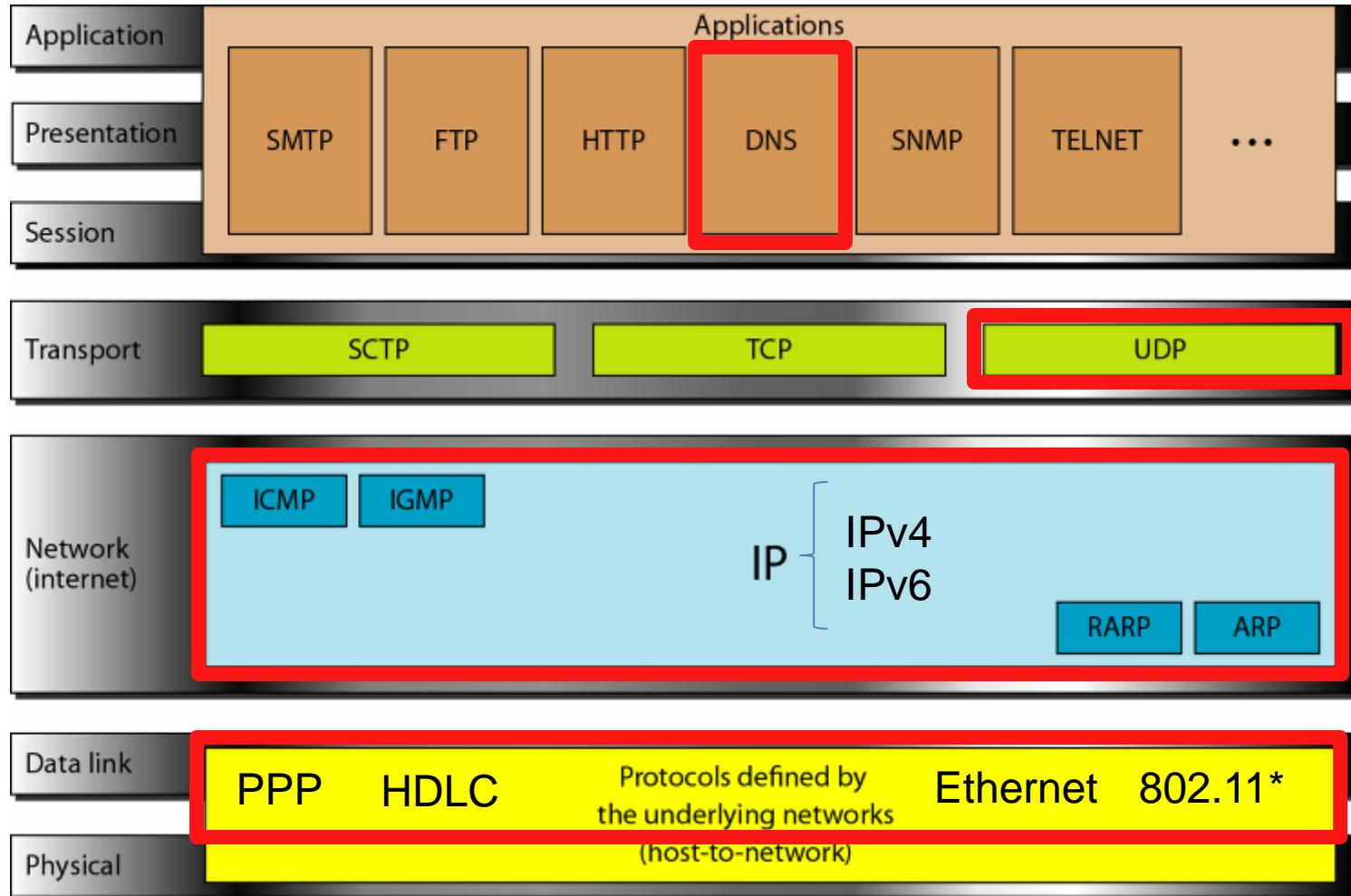
UDP – Portnumber id's Queues

- Connectionless, unreliable protocol
 - No flow and error control
 - Port numbers are used to multiplex data

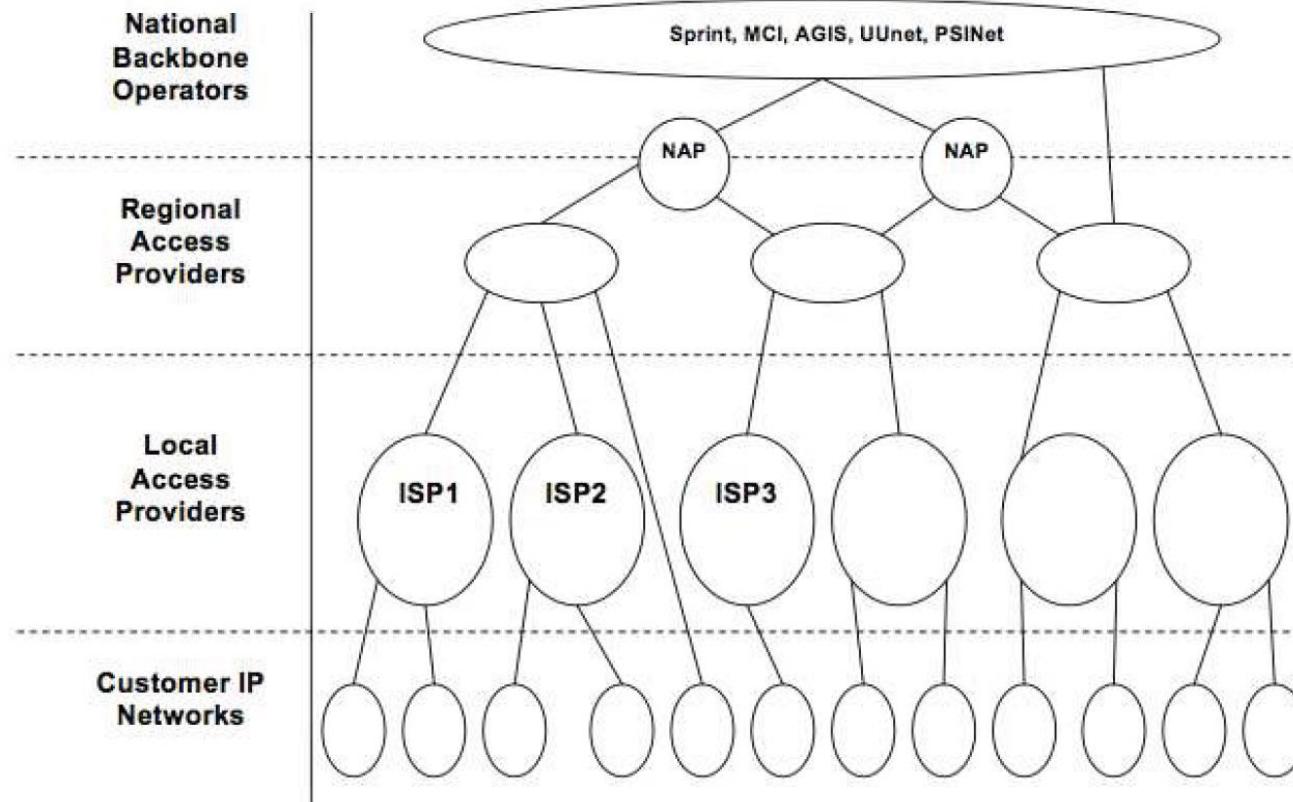


State	Process ID	Port Number	Queue Number
IN-USE	2,345	52,010	34
IN-USE	3,422	52,011	
FREE			
IN-USE	4,652	52,012	38
FREE			

Protocols in the OSI Model



Traditional Logical Internet Topology



CS Predictions

- "I think there is a world market for maybe five computers."

Thomas Watson, President of IBM, 1943

- "There is no reason anyone would want a computer in their home."

Ken Olsen, Founder of Digital Equipment Corporation, 1977

- "I predict the Internet in 1996 will catastrophically collapse."

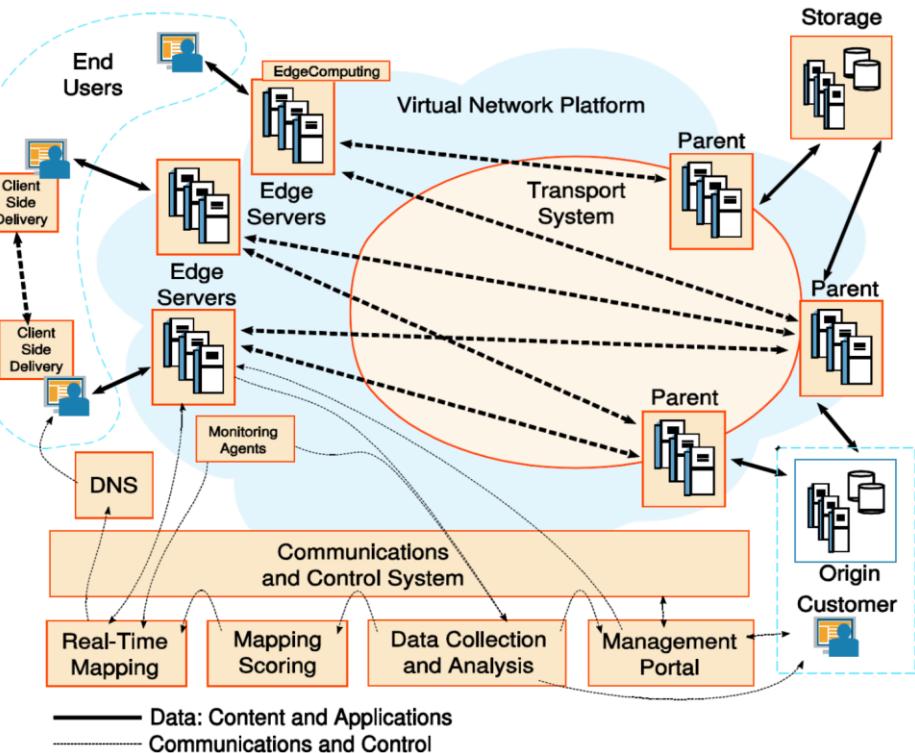
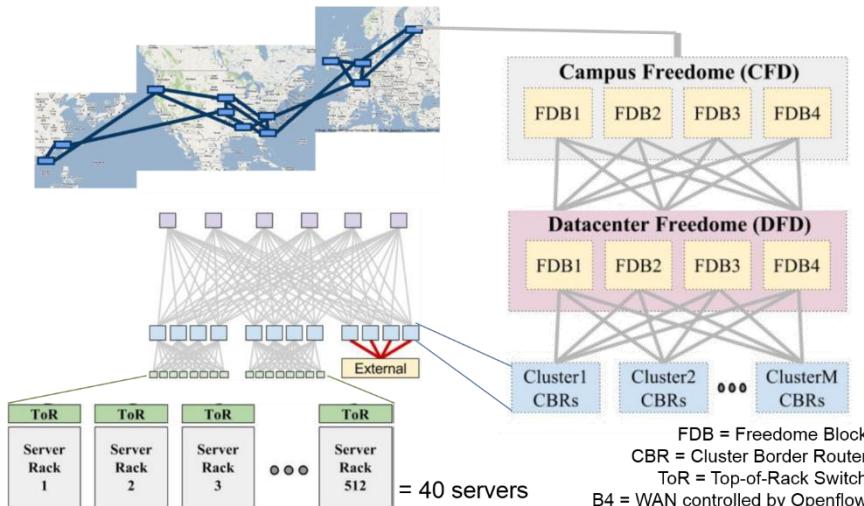
Robert Metcalfe, 1995

- "IPv6 is dead."

David Cheriton, 1999



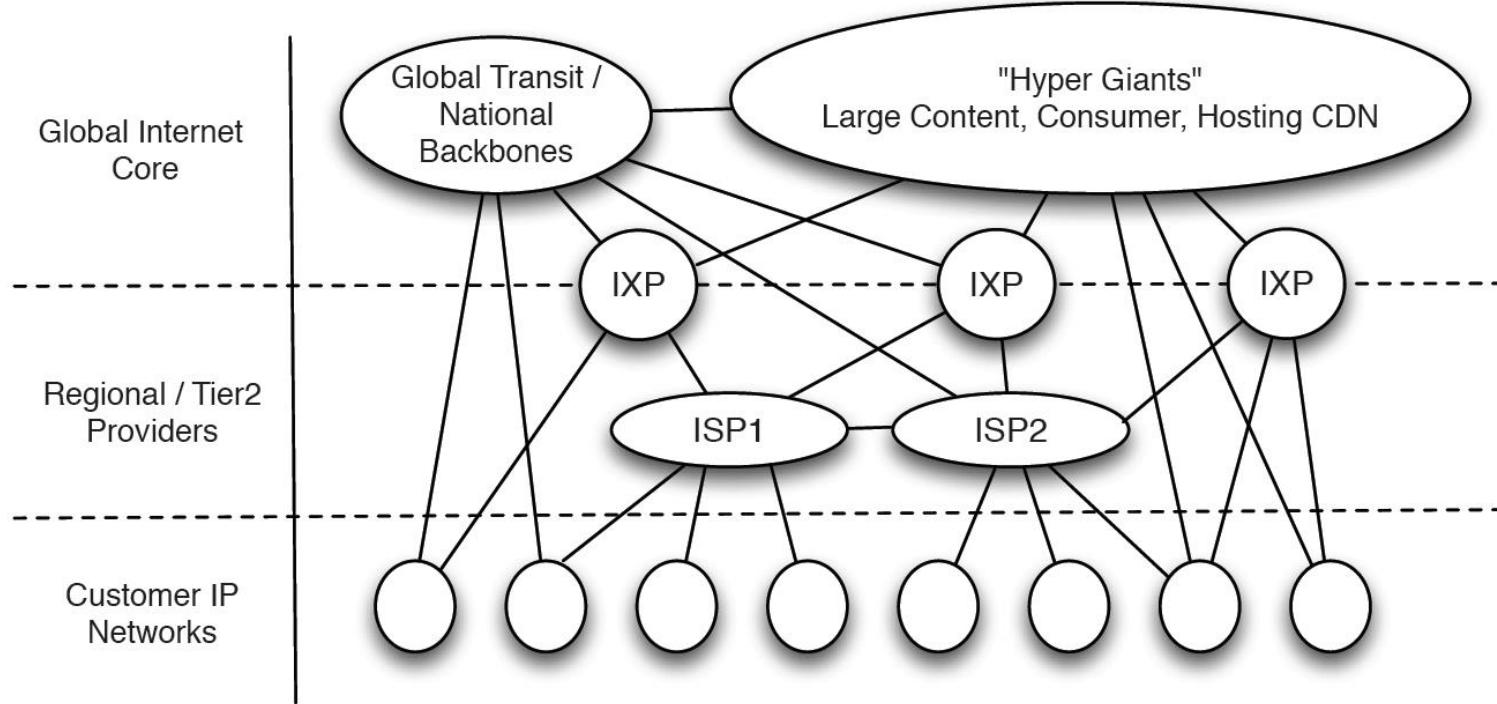
Hyper-Giants



- Google

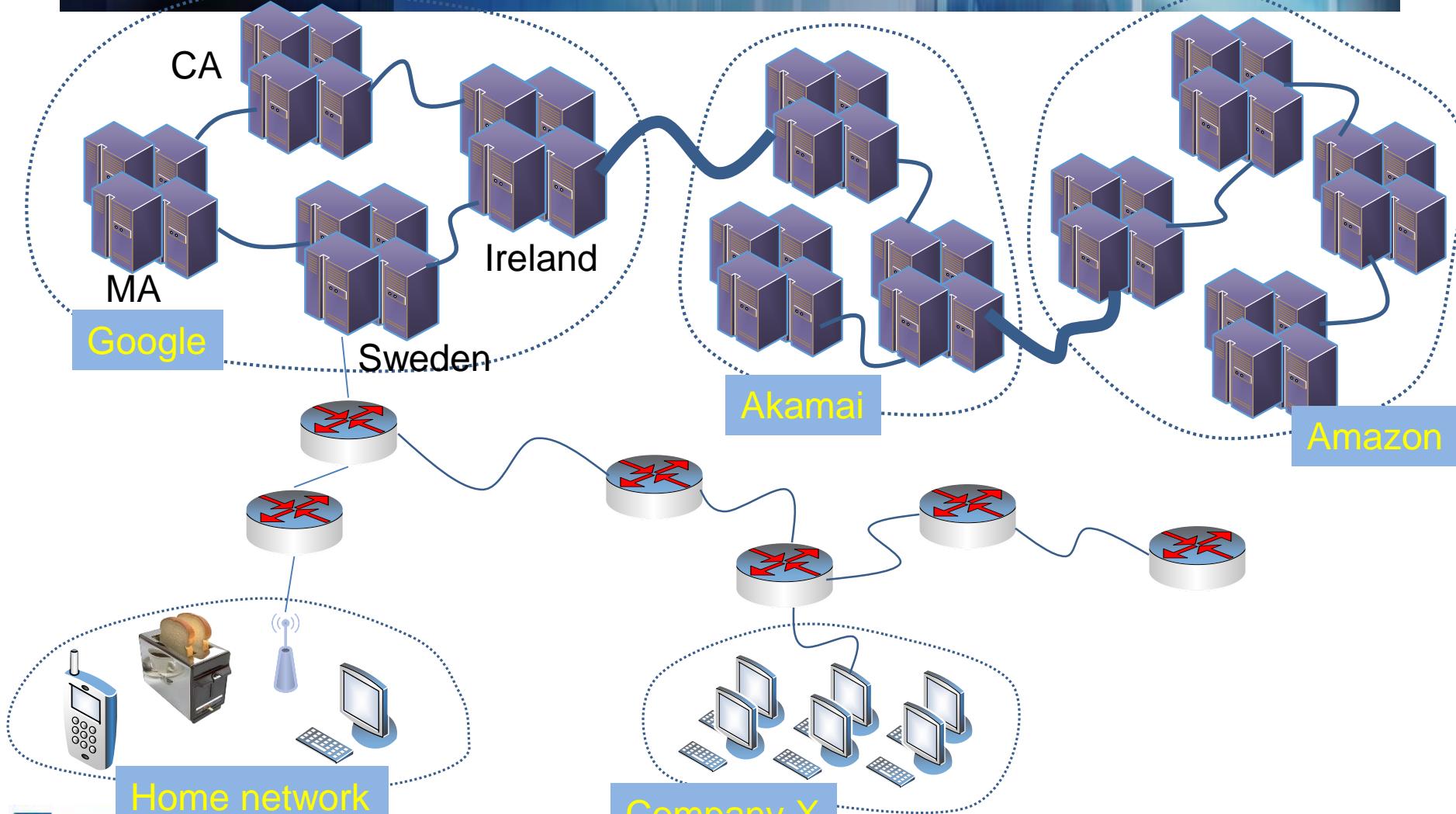
- Akamai

Emerging Logical Internet Topology



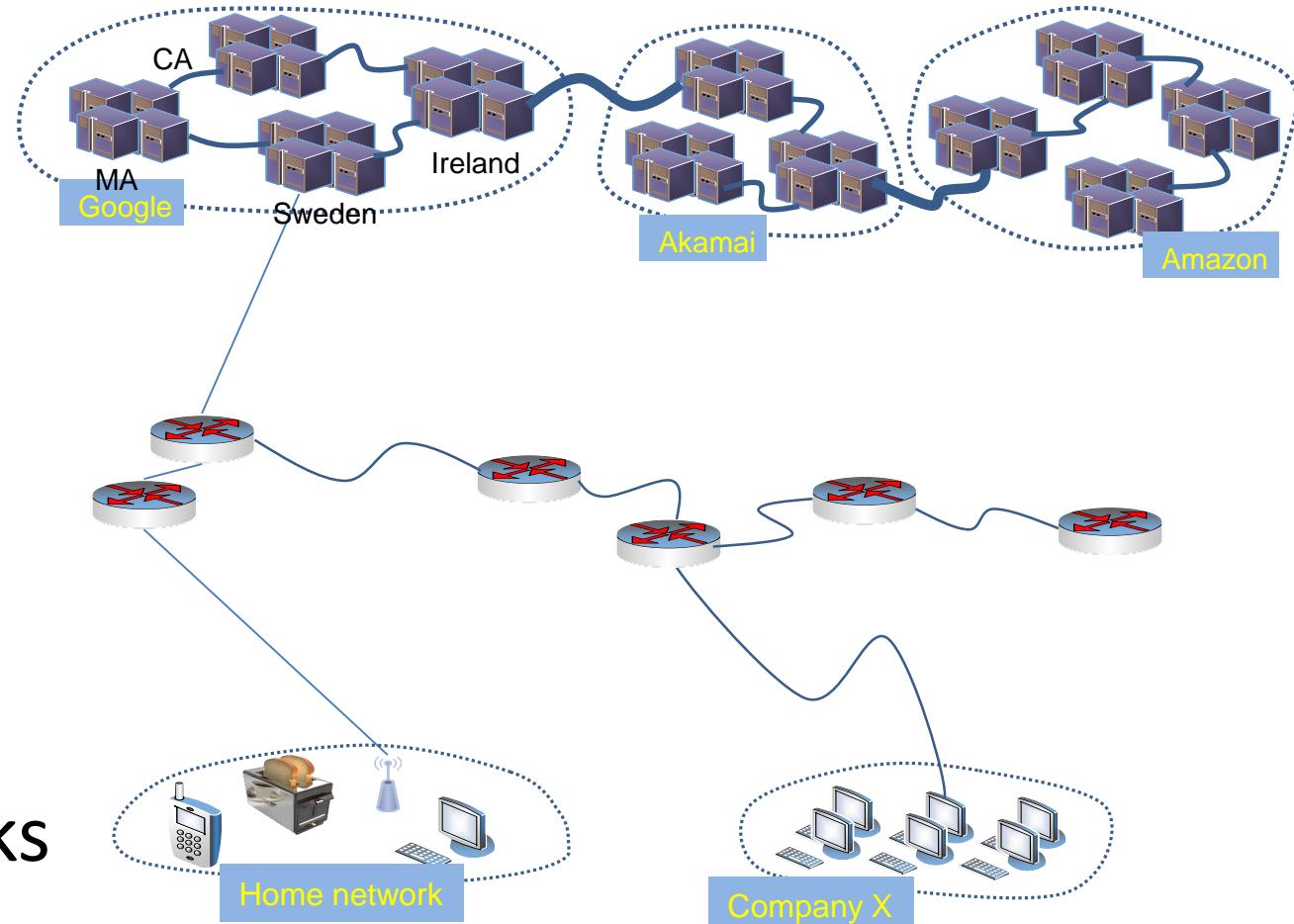
- According to a statement by Craig Labovitz in 2014:
 - 30 Entities created 50% of the traffic in the US at peak time

Datacentres at the Core?



My view of “future”* networking

- Sets of Datacentres



- Traditional Internet

- Edge networks

*or current?



Overview

- Link Layer
- Network Layer
 - Addressing
 - Address Resolution (ARP)
 - Fragmentation
 - Intra-AS Routing
 - Distance Vector
 - Link State
 - Multicast Routing
 - IPv6
- Transport Layer
 - UDP
 - DNS
- Software-Defined Networking / Openflow
- CLOS / Fat-tree
- ATM/MPLS



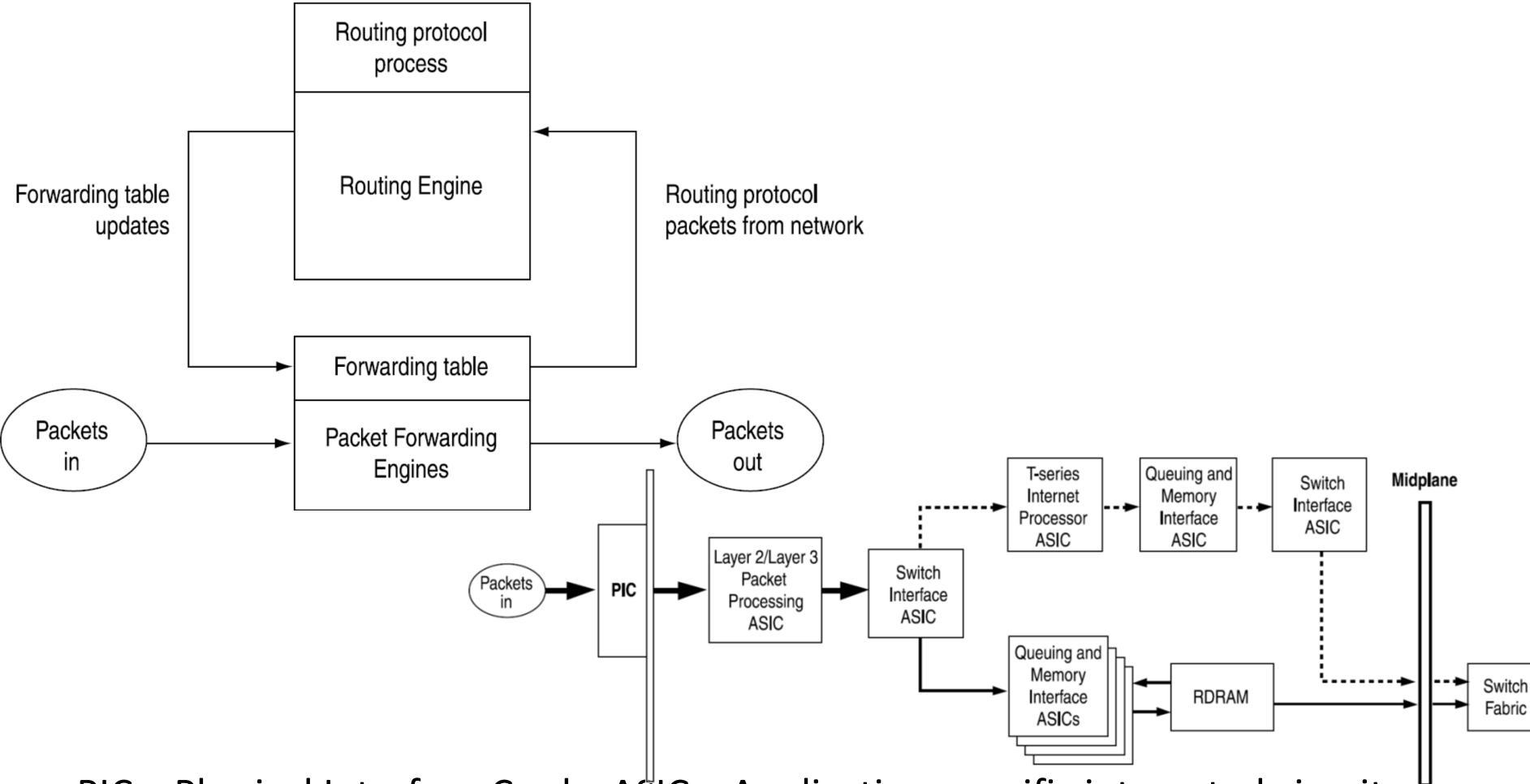


CS2031 Telecommunications II

SDN & Openflow



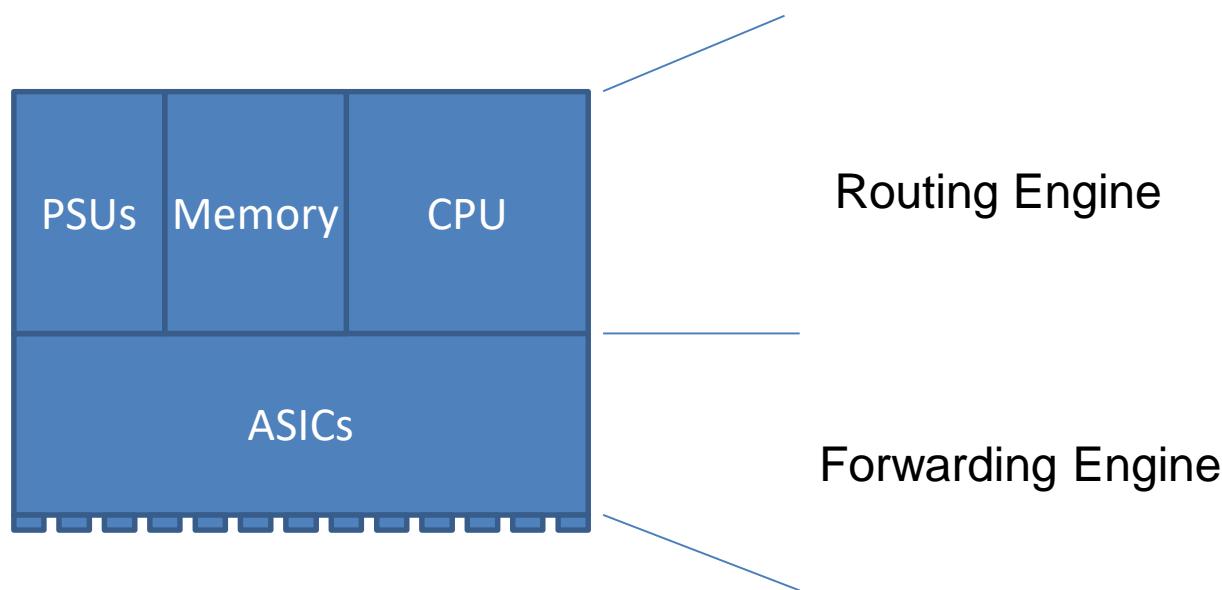
Switch/Router Internals



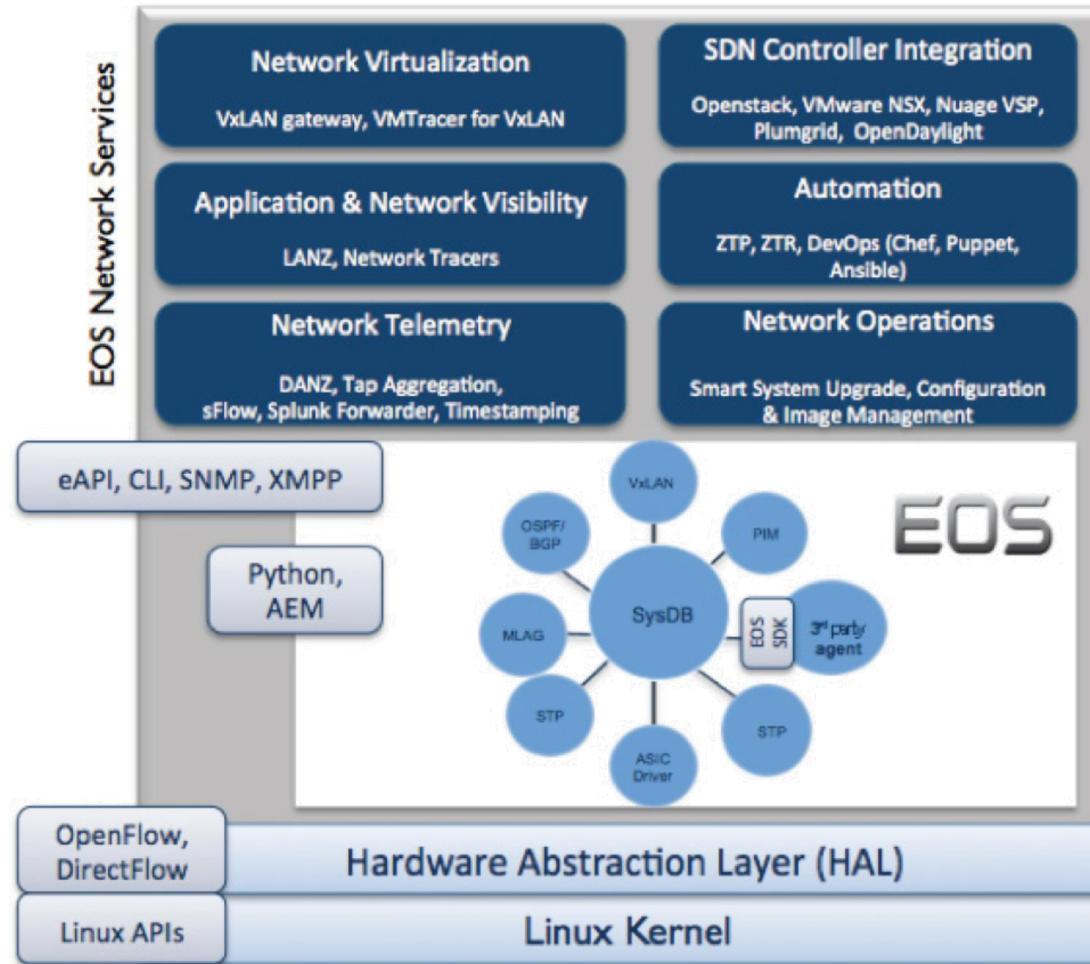
PIC = Physical Interface Card

ASIC = Application-specific integrated circuits

Switches/Routers



Example Switches/Router OS



Original Openflow Paper

- Openflow switches
- Controller w/ secure connection
- Configurable flow tables

OpenFlow: Enabling Innovation in Campus Networks

March 14, 2008

Nick McKeown
Stanford University

Guru Parulkar
Stanford University

Scott Shenker
University of California,
Berkeley

Tom Anderson
University of Washington

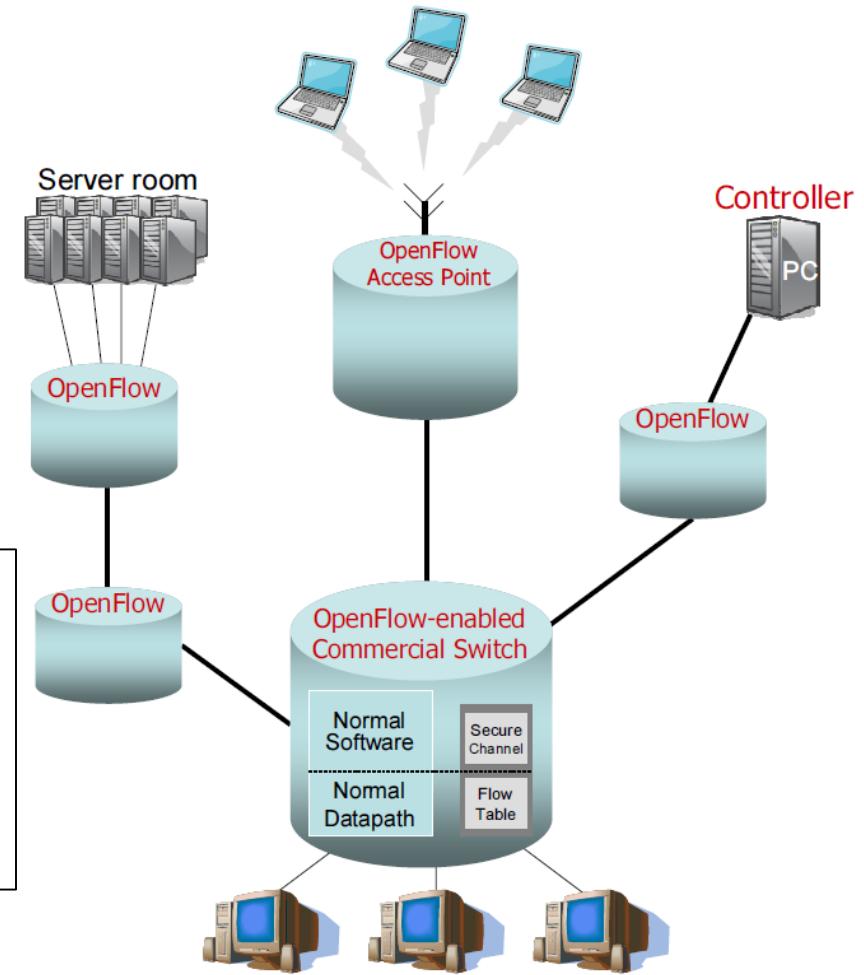
Larry Peterson
Princeton University

Jonathan Turner
Washington University in
St. Louis

Hari Balakrishnan
MIT

Jennifer Rexford
Princeton University

Header of 2008 Paper



Openflow Switch

Software Layer

OpenFlow Client

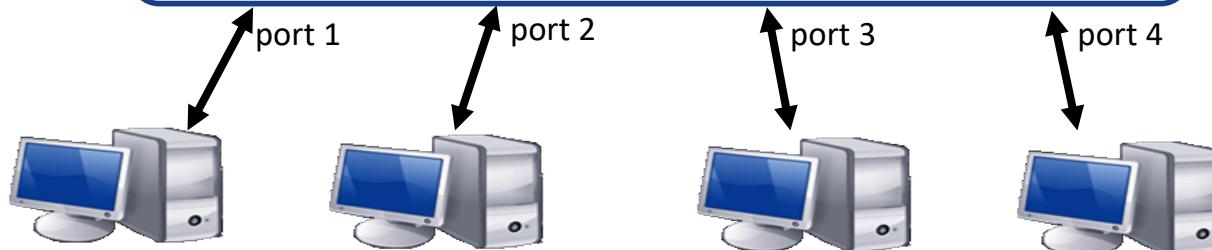
Hardware Layer

Flow Table						
MAC src	MAC dst	IP Src	IP Dst	TCP sport	TCP dport	Action
*	*	*	5.6.7.8	*	*	port 1

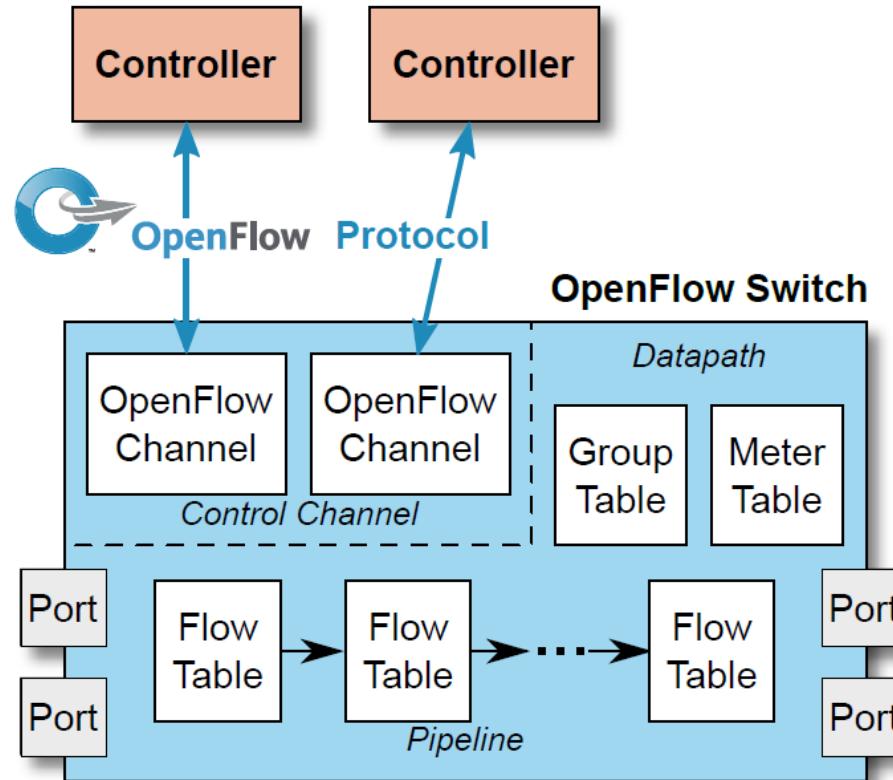


SSH conn.

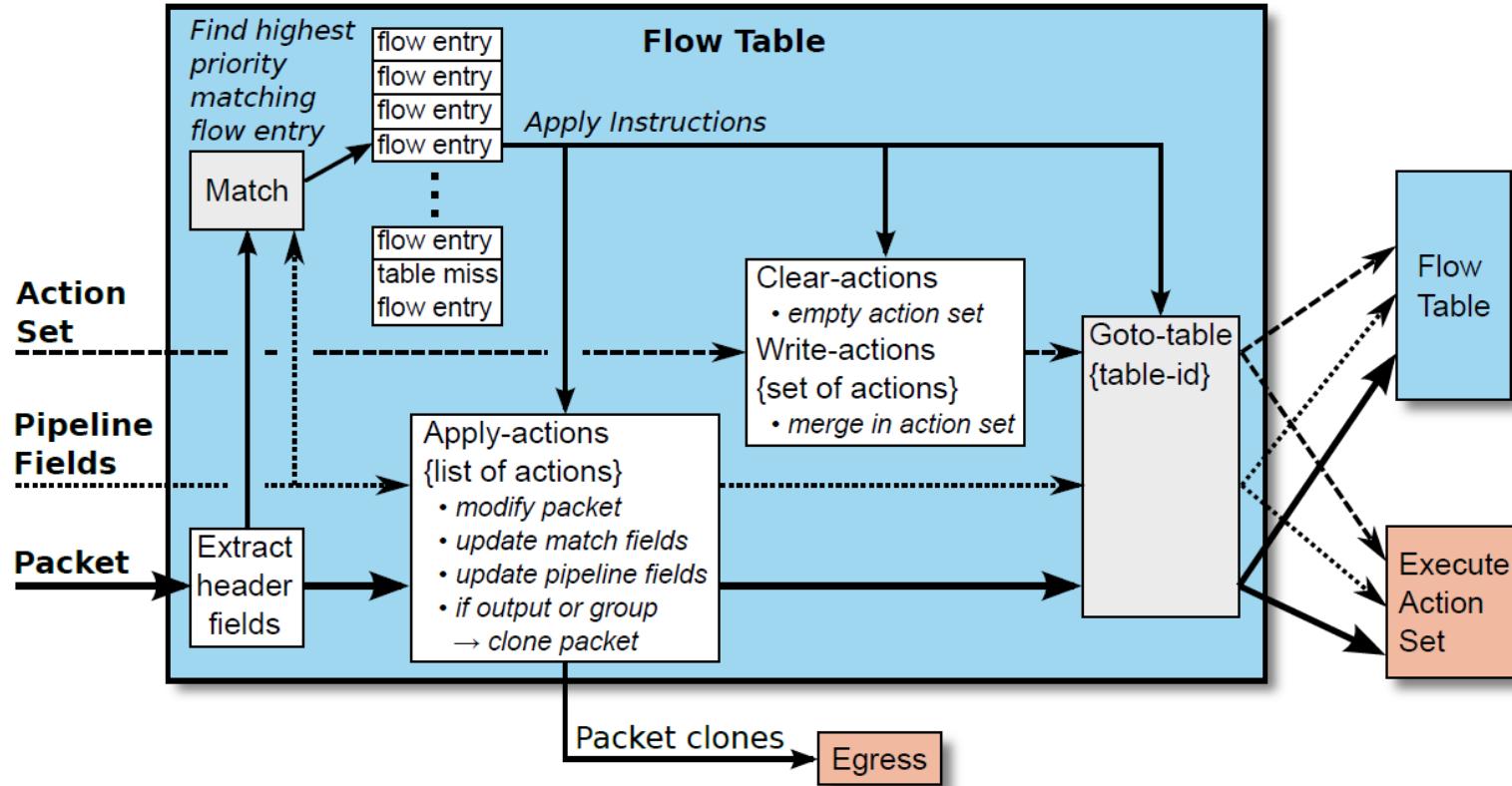
Controller



Components of Openflow Switches



Flowtable Processing



Counters in Flowtables

Counter	Bits	
Per Flow Table		
Reference Count (active entries)	32	<i>Required</i>
Packet Lookups	64	<i>Optional</i>
Packet Matches	64	<i>Optional</i>
Per Flow Entry		
Received Packets	64	<i>Optional</i>
Received Bytes	64	<i>Optional</i>
Duration (seconds)	32	<i>Required</i>
Duration (nanoseconds)	32	<i>Optional</i>
Per Port		
Received Packets	64	<i>Required</i>
Transmitted Packets	64	<i>Required</i>
Received Bytes	64	<i>Optional</i>
Transmitted Bytes	64	<i>Optional</i>
Receive Drops	64	<i>Optional</i>
Transmit Drops	64	<i>Optional</i>
Receive Errors	64	<i>Optional</i>
Transmit Errors	64	<i>Optional</i>
Receive Frame Alignment Errors	64	<i>Optional</i>
Receive Overrun Errors	64	<i>Optional</i>
Receive CRC Errors	64	<i>Optional</i>
Collisions	64	<i>Optional</i>
Duration (seconds)	32	<i>Required</i>
Duration (nanoseconds)	32	<i>Optional</i>
Per Queue		





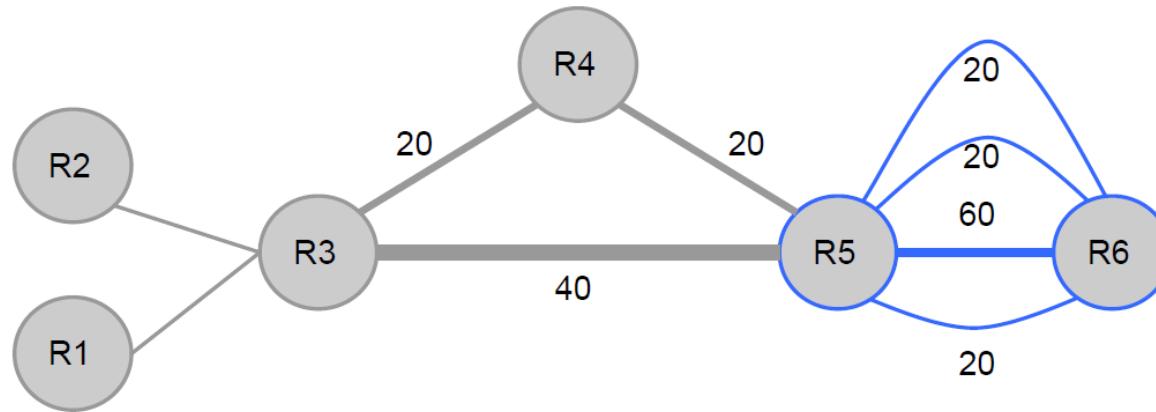
Why????

As in: Why Openflow? ☺



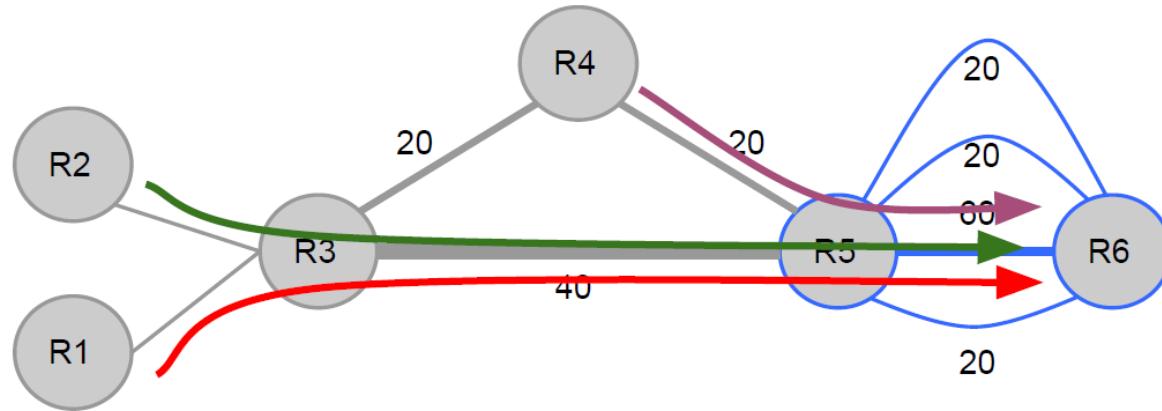
Traditional Net. Example

- Flows: R1->R6: 20; R2->R6: 20; R4->R6: 20



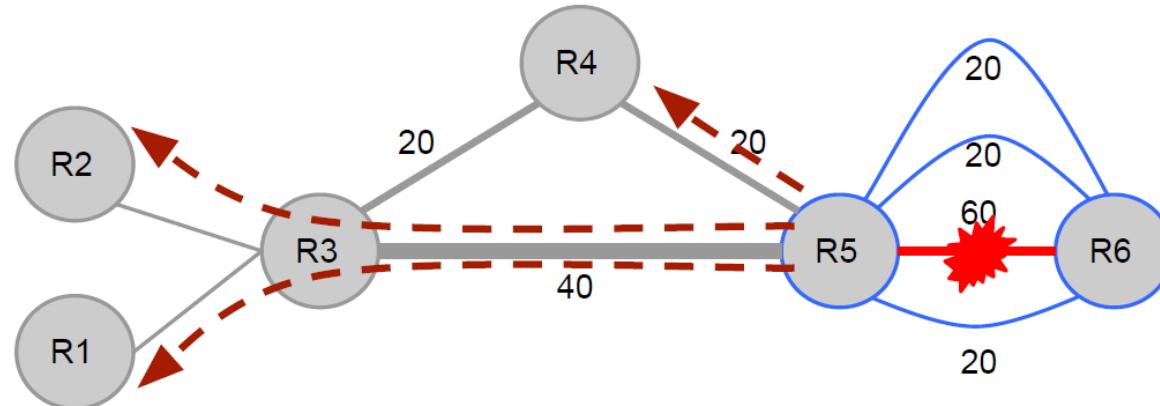
Traditional Net. Example

- Flows: R1->R6: 20; R2->R6: 20; R4->R6: 20



Traditional Net. Example

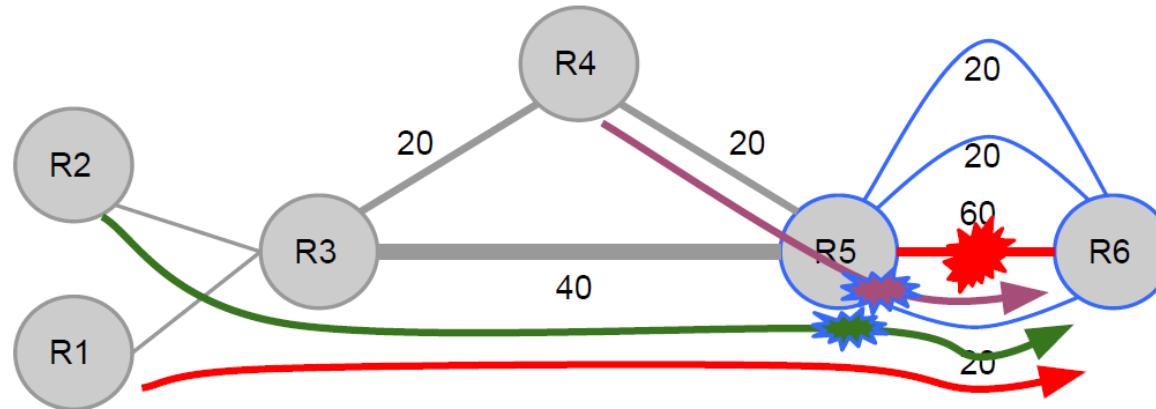
- Flows: R1->R6: 20; R2->R6: 20; R4->R6: 20



- R5-R6 link fails
 - R1, R2, R4 autonomously try for next best path

Traditional Net. Example

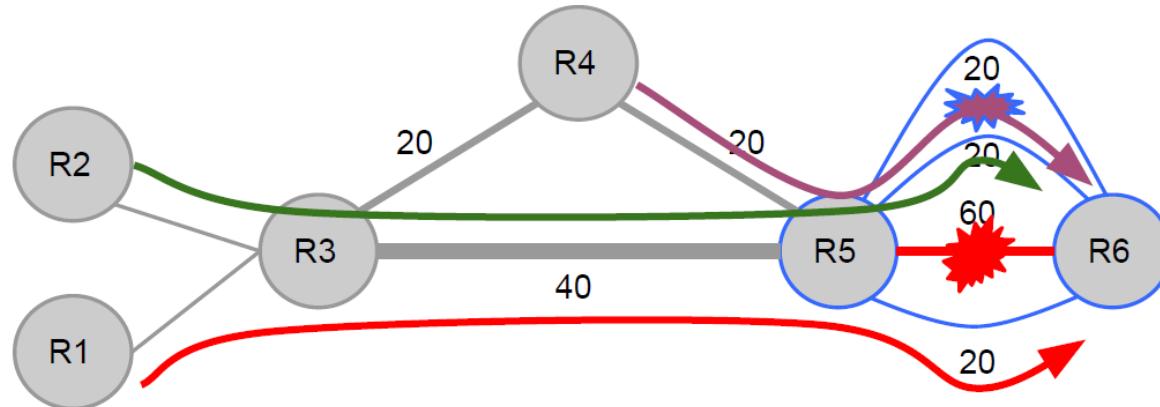
- Flows: R1->R6: 20; R2->R6: 20; R4->R6: 20



- R5-R6 link fails
 - R1, R2, R4 *autonomously* try for next best path
 - R1 wins, R2, R4 retry for next best path

Traditional Net. Example

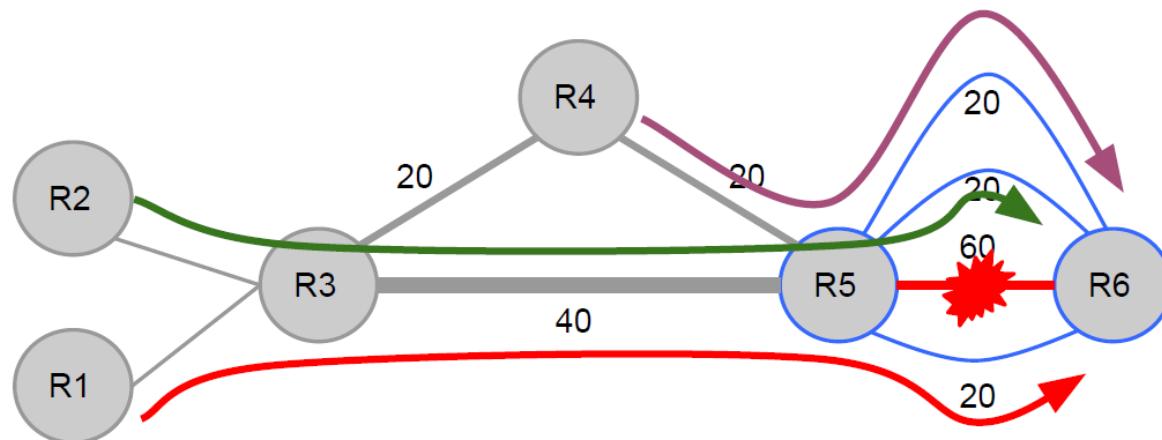
- Flows: R1->R6: 20; R2->R6: 20; R4->R6: 20



- R5-R6 link fails
 - R1, R2, R4 *autonomously* try for next best path
 - R1 wins, R2, R4 retry for next best path
 - R2 wins this round, R4 retries again

Traditional Net. Example

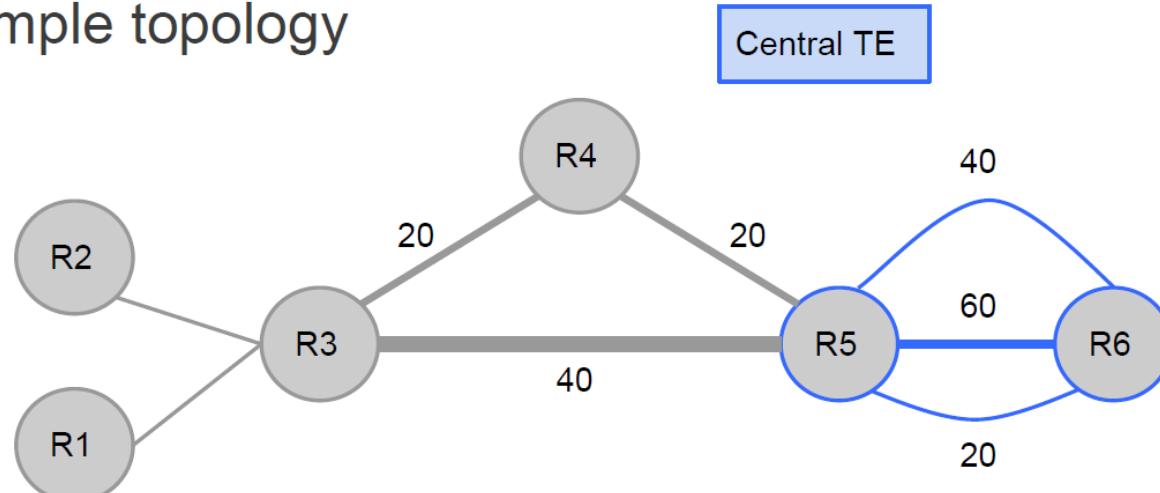
- Flows: R1->R6: 20; R2->R6: 20; R4->R6: 20



- R5-R6 link fails
 - R1, R2, R4 *autonomously* try for next best path
 - R1 wins, R2, R4 retry for next best path
 - R2 wins this round, R4 retries again
 - R4 finally gets third best path

Topology with Central TE

- Simple topology



- Flows:
 - R1->R6: 20; R2->R6: 20; R4->R6: 20

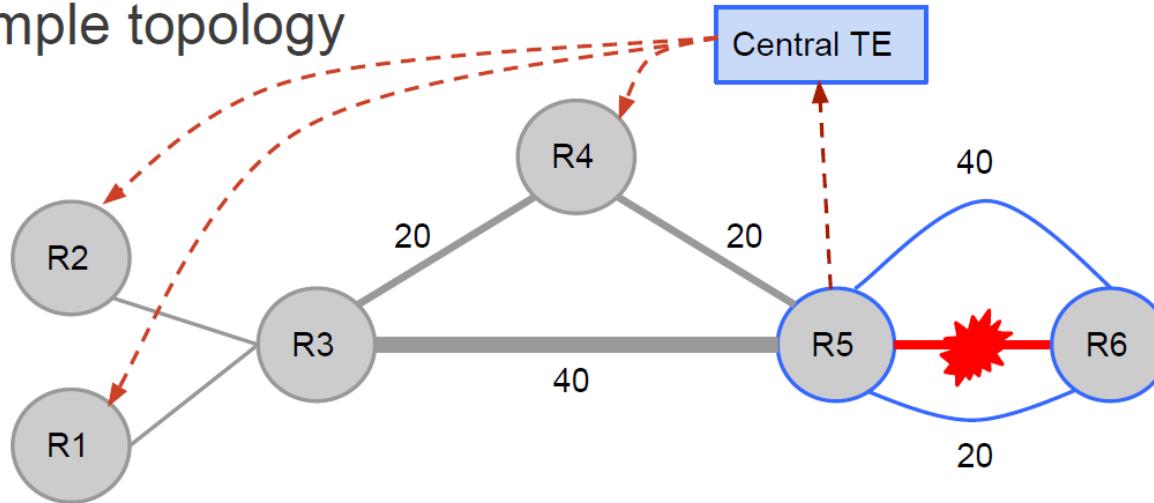
TE = Traffic Engineering

34

• Urs Hoelzle, Open Network Summit 2012

Topology with Central TE

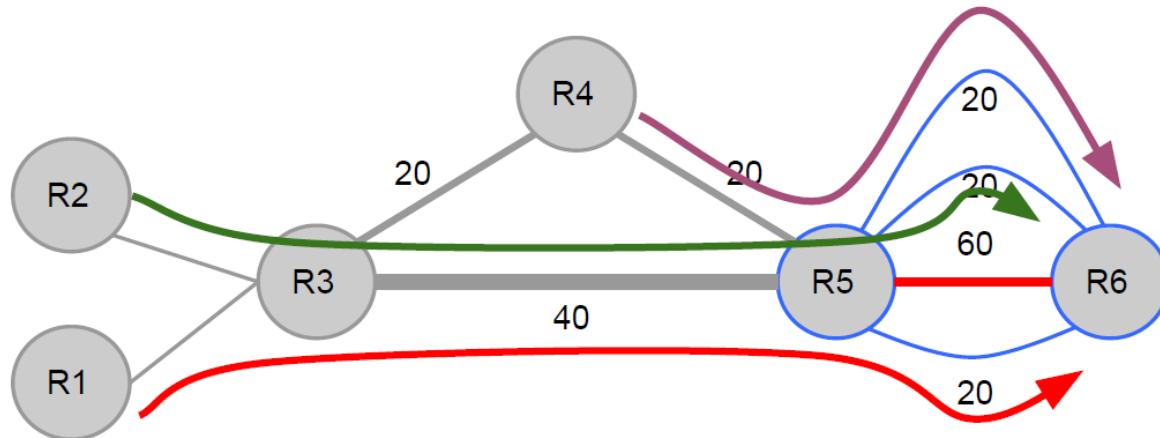
- Simple topology



- Flows:
 - R1->R6: 20; R2->R6: 20; R4->R6: 20
- R5-R6 fails
 - R5 informs TE, which programs routers in one shot

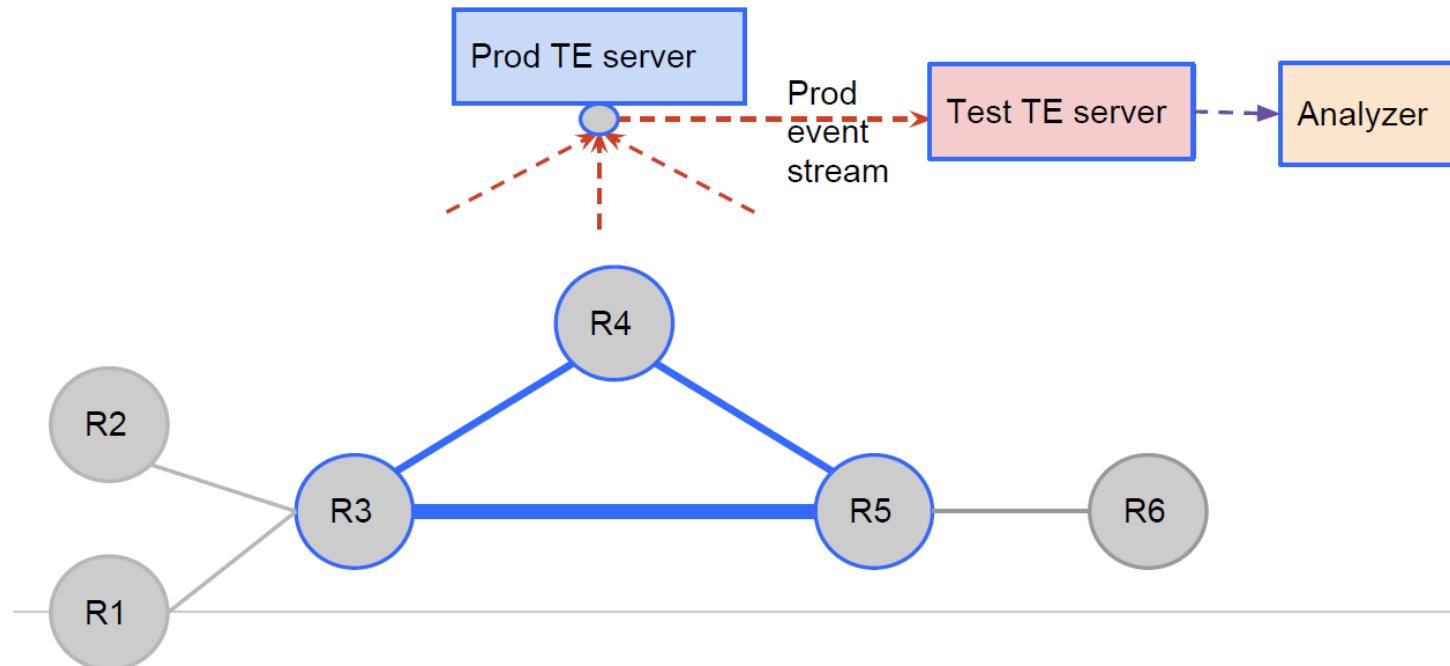
Topology with Central TE

- Simple topology

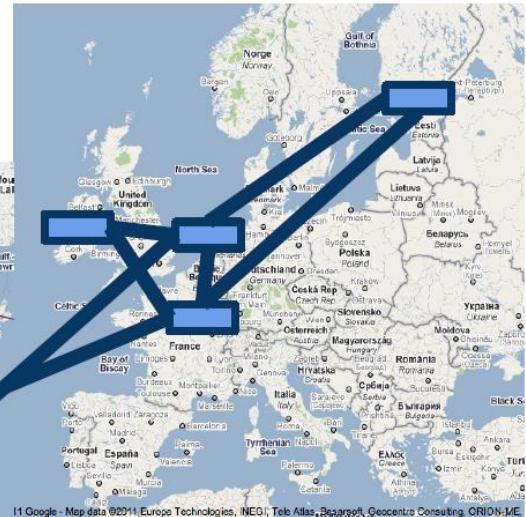
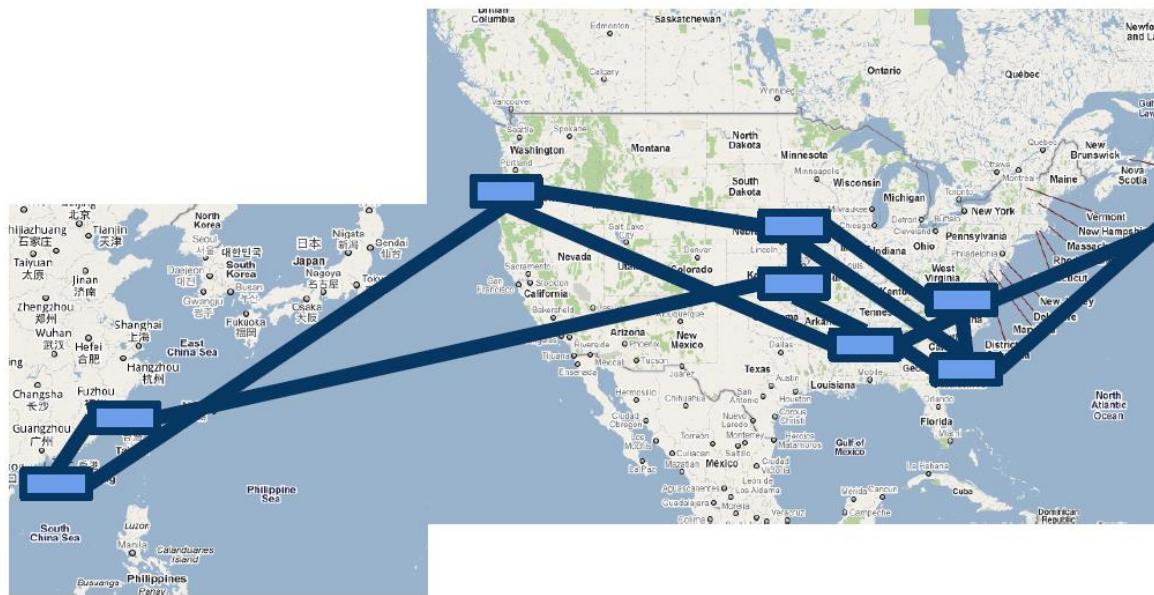


- Flows:
 - R1->R6: 20; R2->R6: 20; R4->R6: 20
- R5-R6 link fails
 - R5 informs TE, which programs routers in one shot
 - Leads to faster realization of target optimum

Traffic Engineering with Analyzer

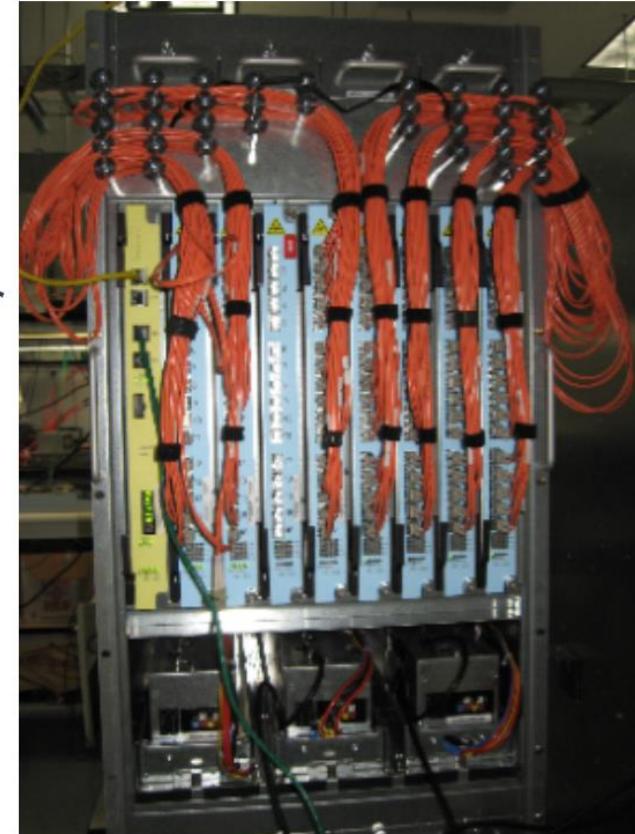


Google's WAN

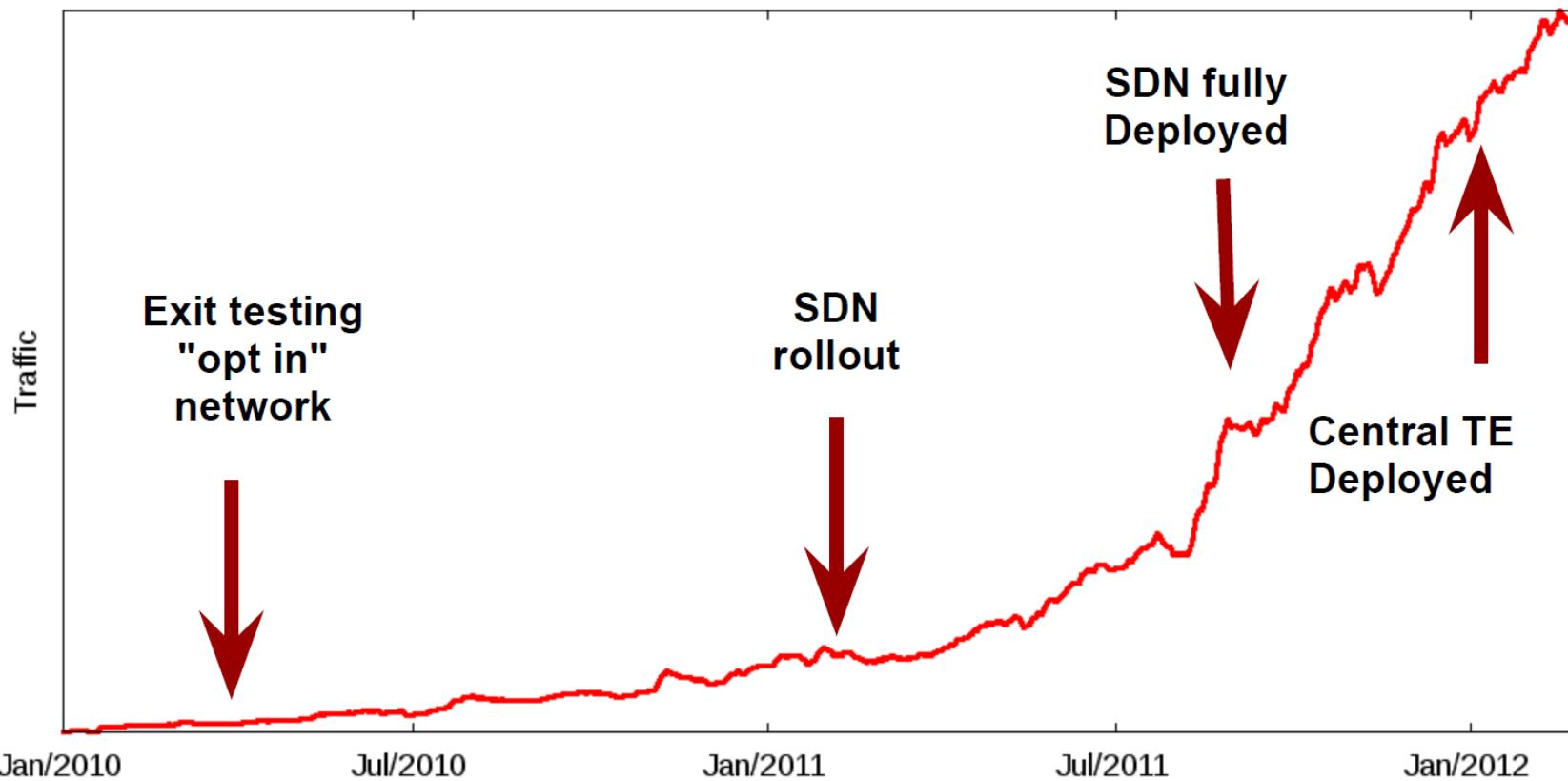


Google's Hardware

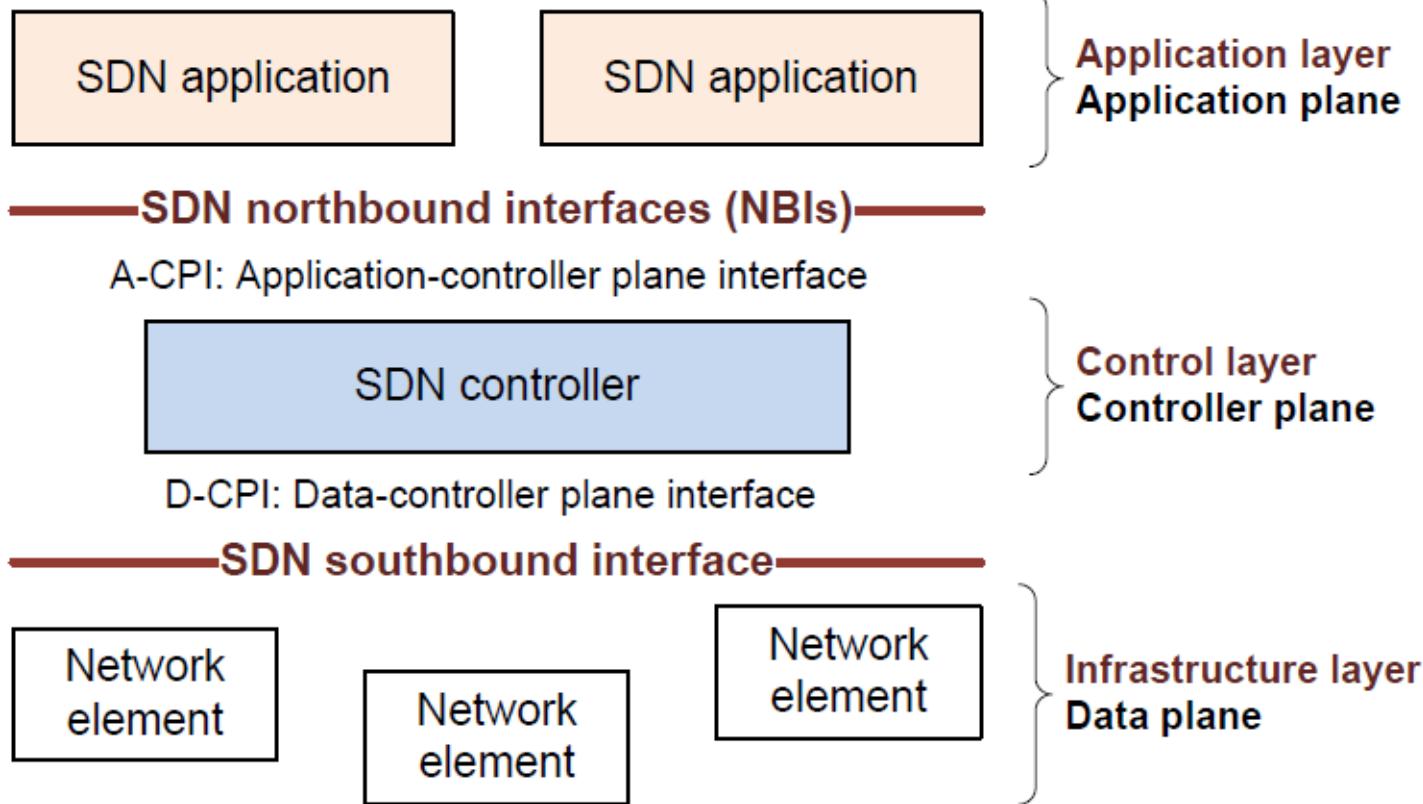
- Built from merchant silicon
 - 100s of ports of nonblocking 10GE
- OpenFlow support
- Open source routing stacks for BGP, ISIS
- Does not have all features
 - No support for AppleTalk...
- Multiple chassis per site
 - Fault tolerance
 - Scale to multiple Tbps



Timeline of Deployment



SDN Architecture



SDN Overview

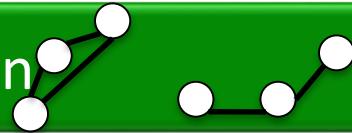
Specifies behavior

Control Program

Abstract Network Model

Compiles to topology

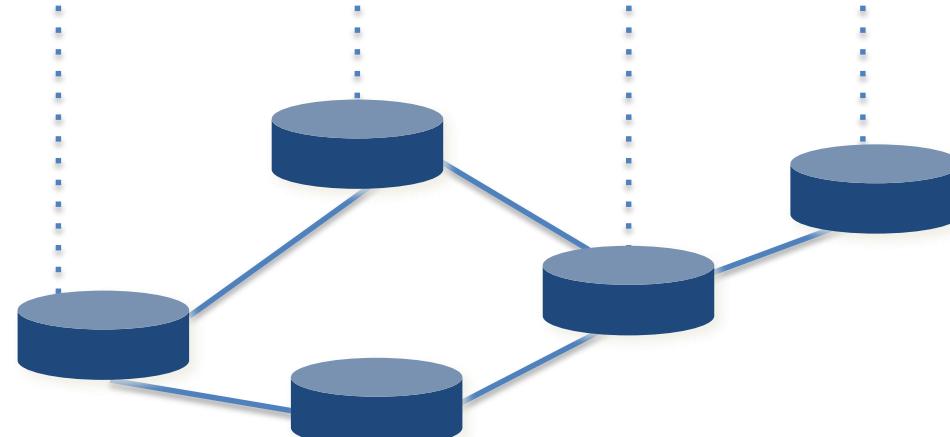
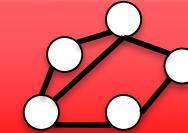
Network Virtualization

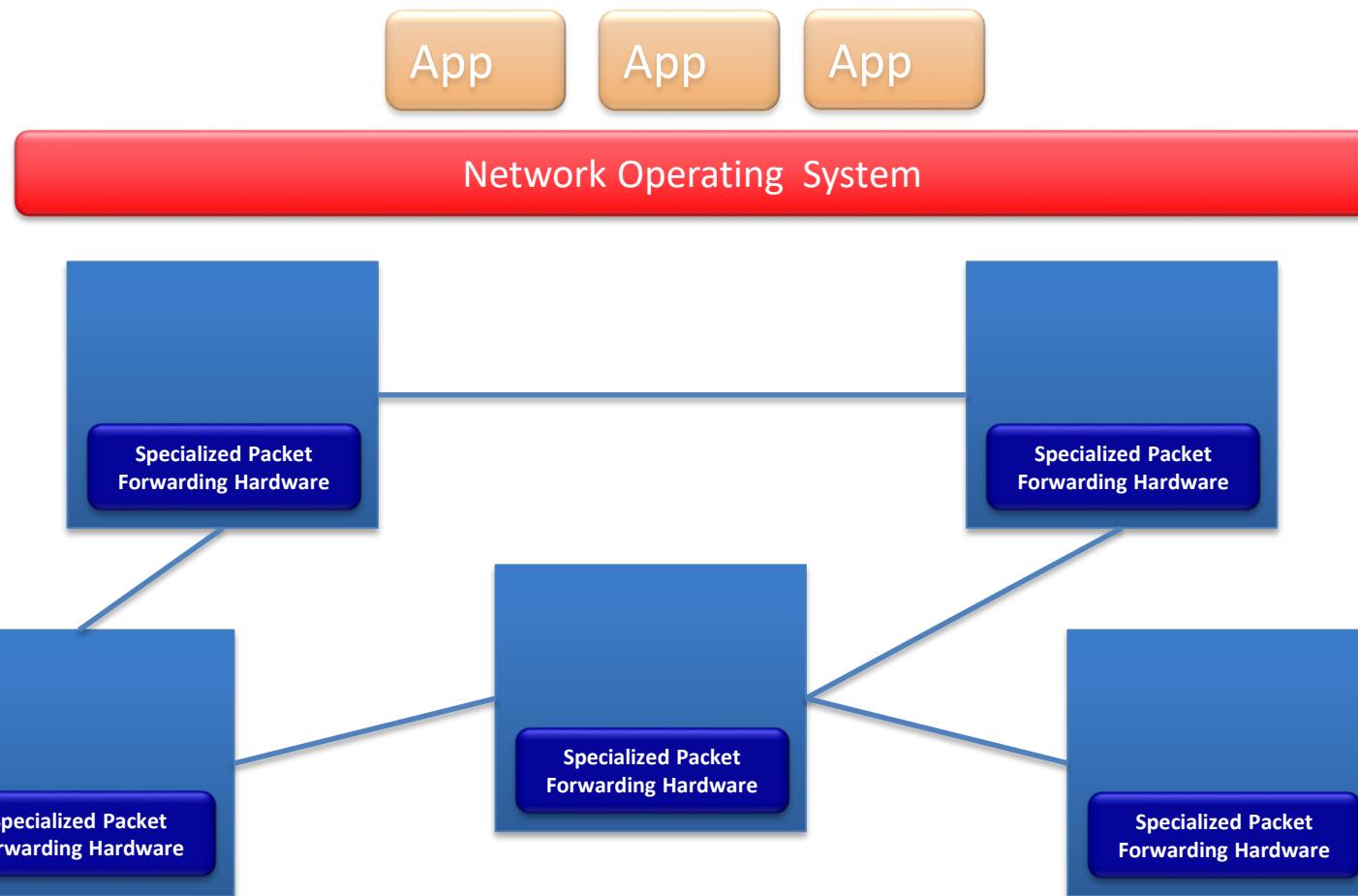


Global Network View

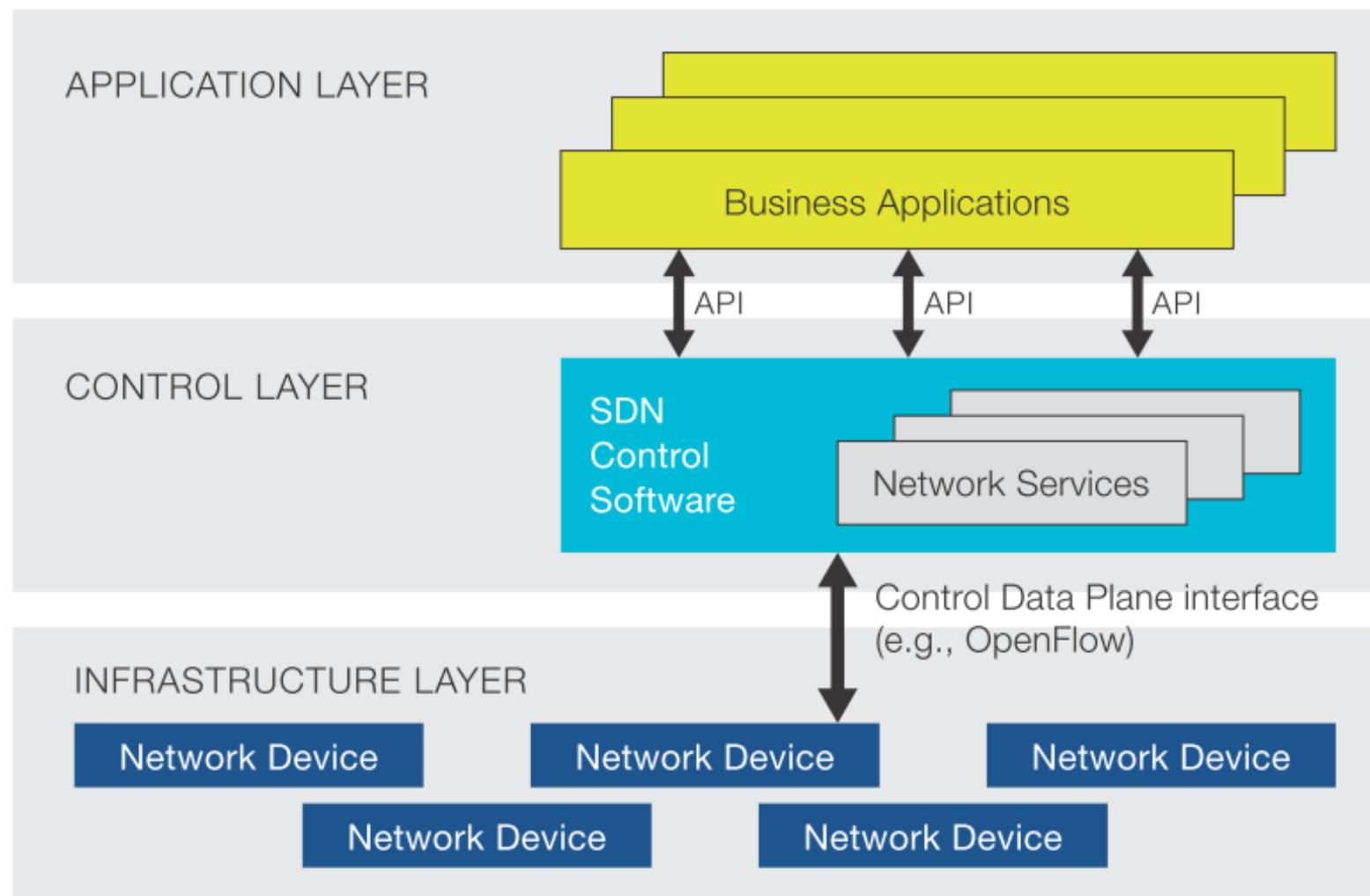
Transmits to switches

Network OS

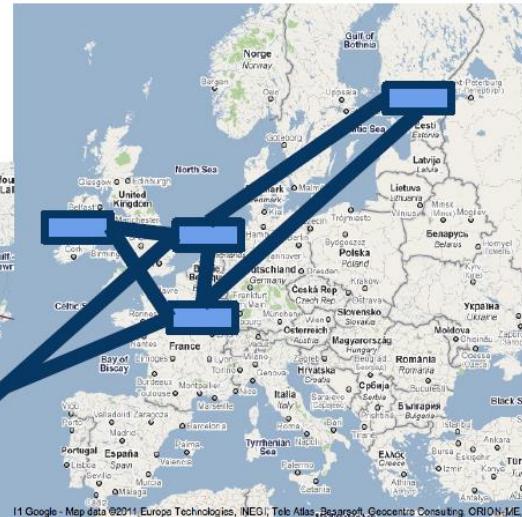
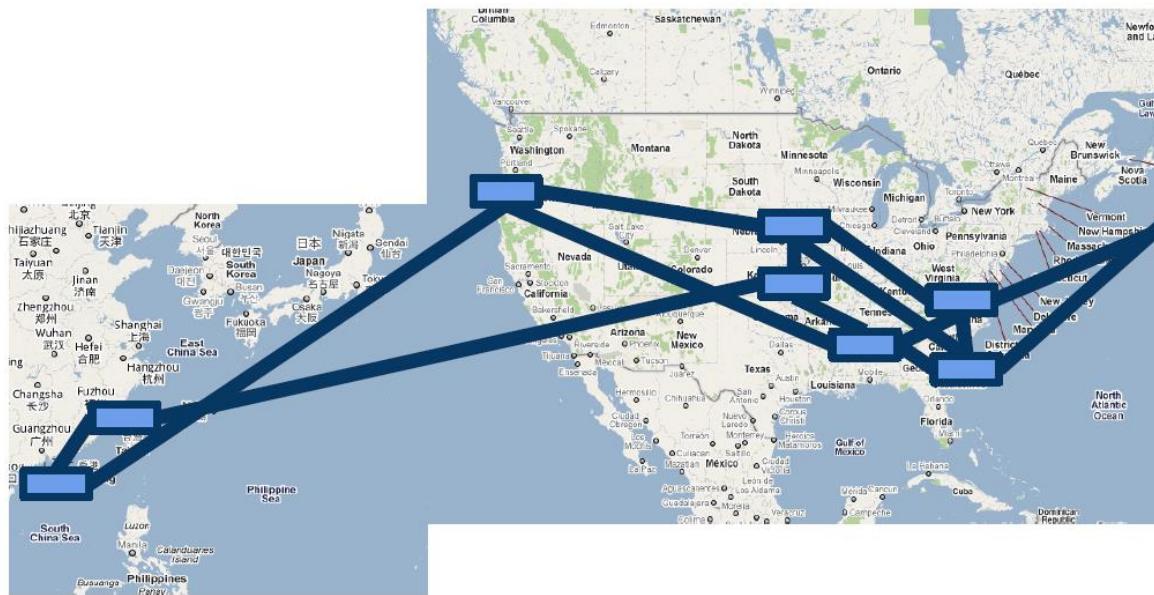




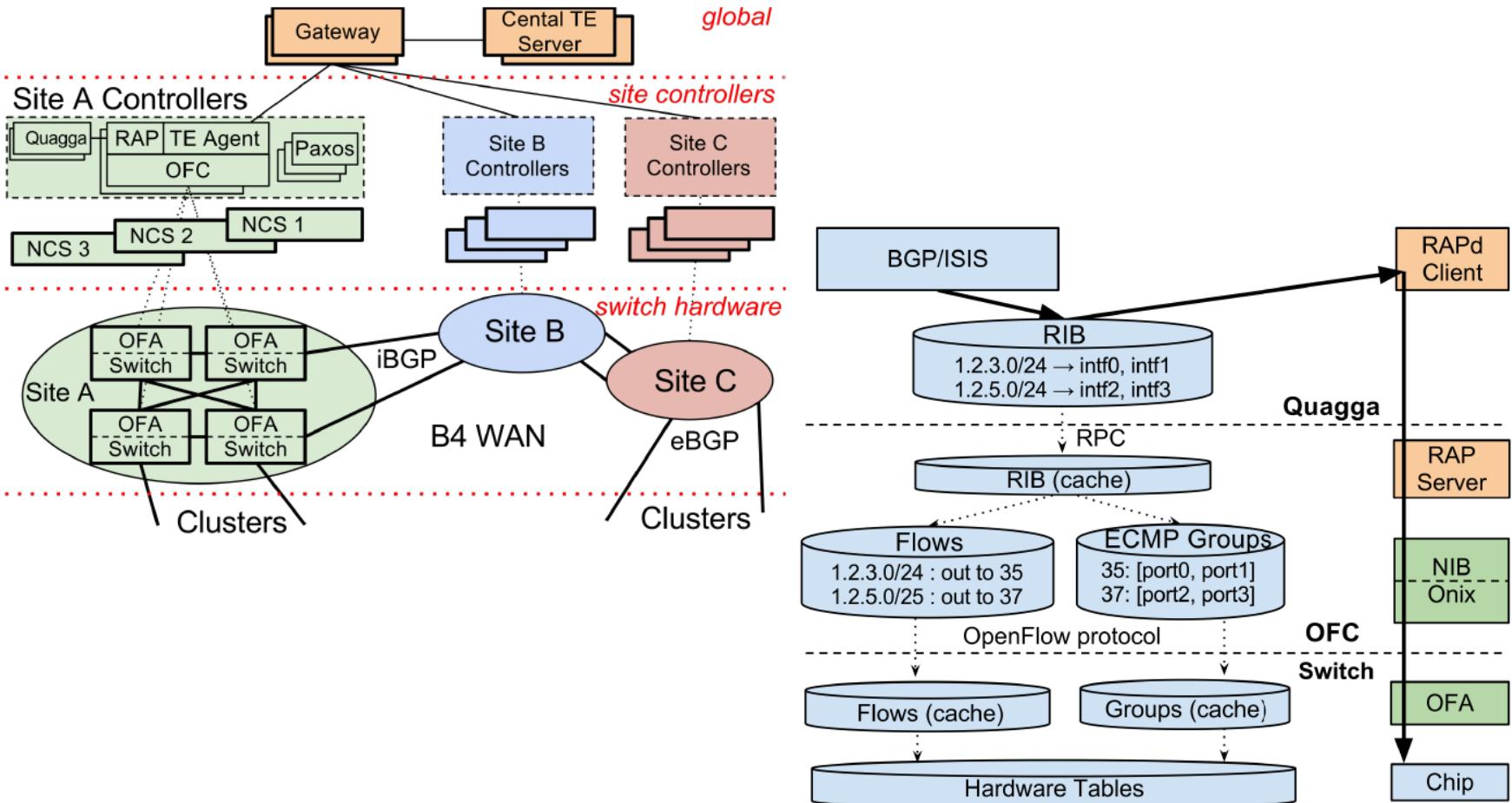
Openflow SDN Architecture



Google's WAN



Openflow in B4





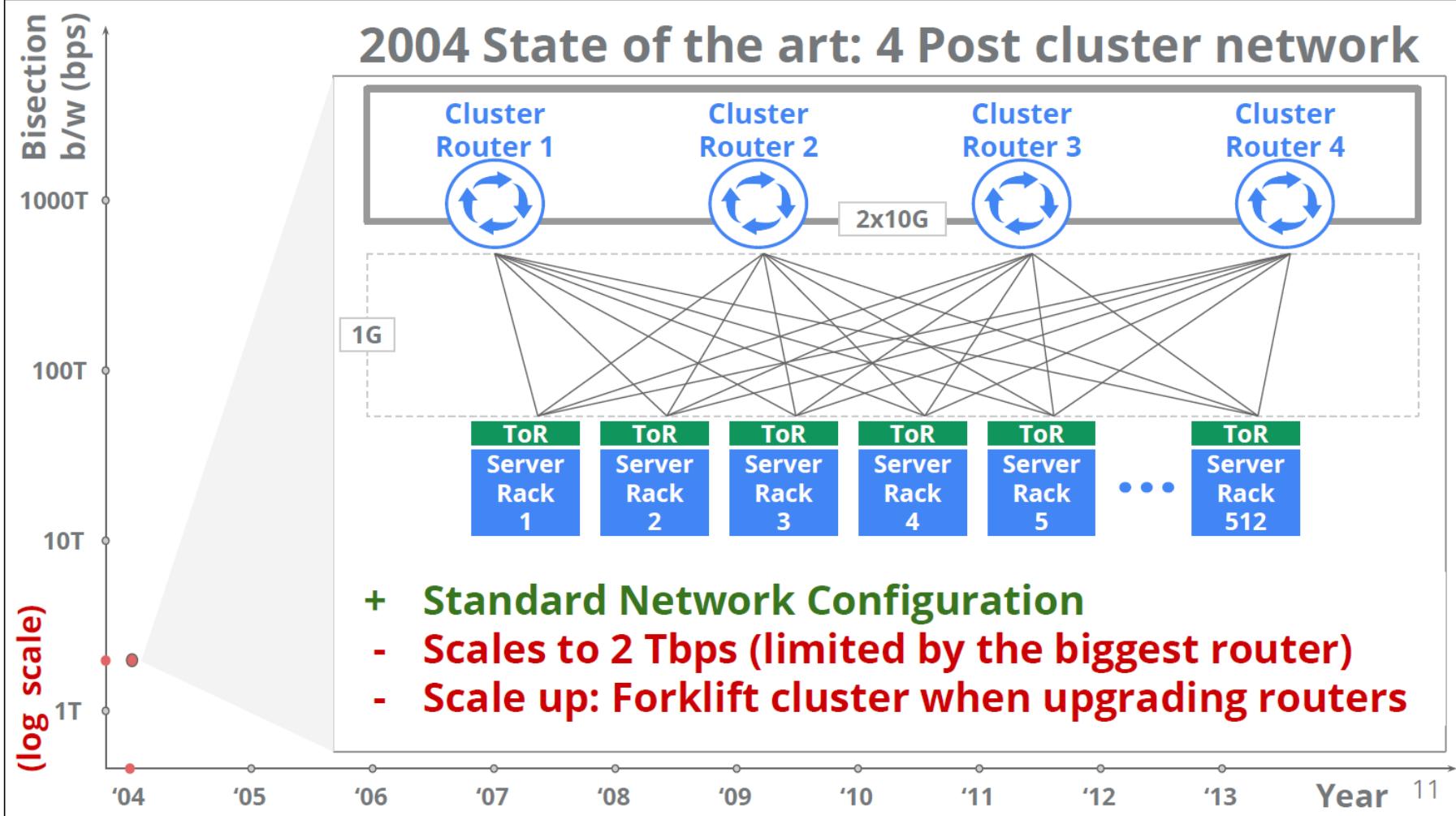
CS2031

Telecommunications II

Clos / Google's Infrastructure



Google's Infrastructure

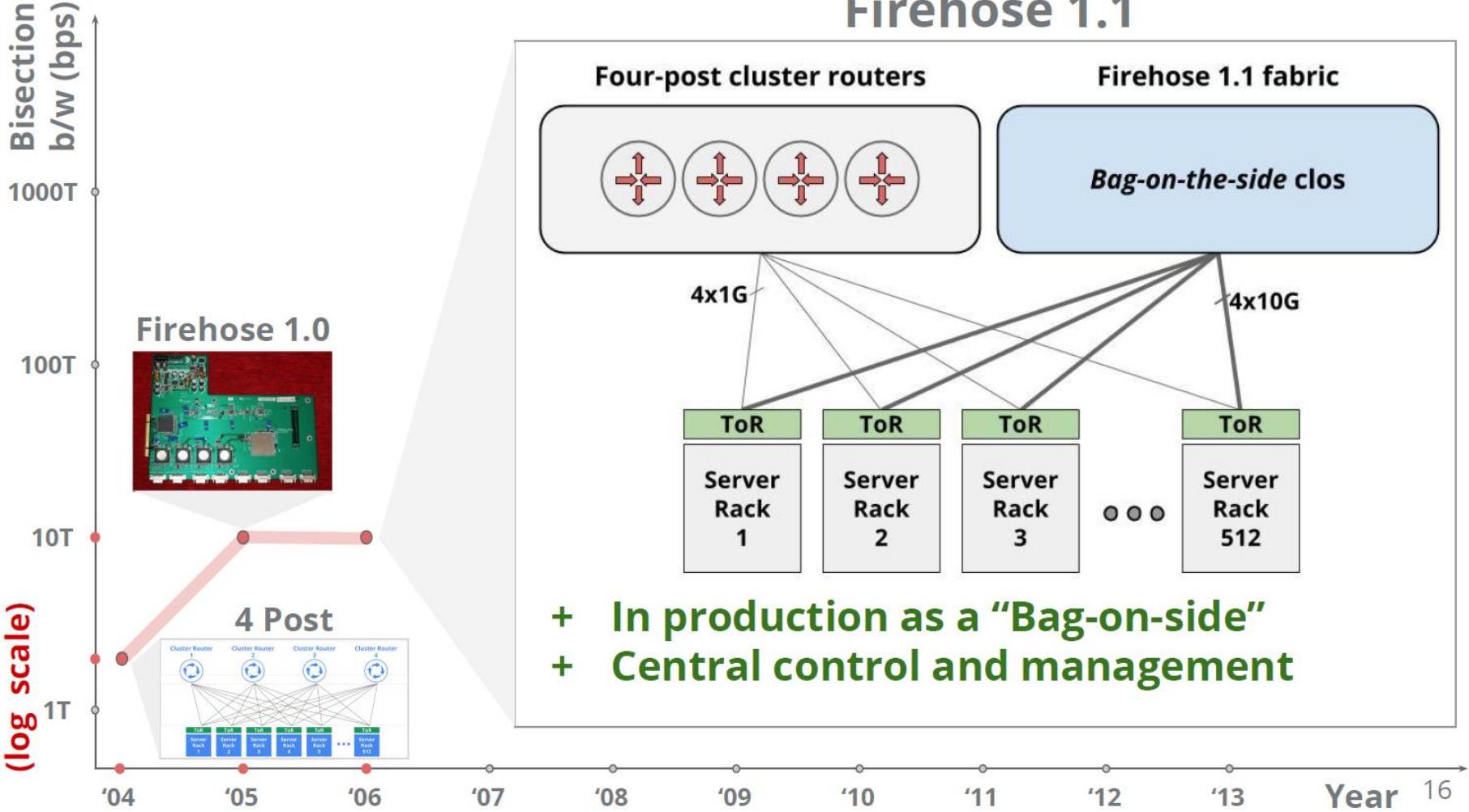


The Obligatory Datacentre Pic

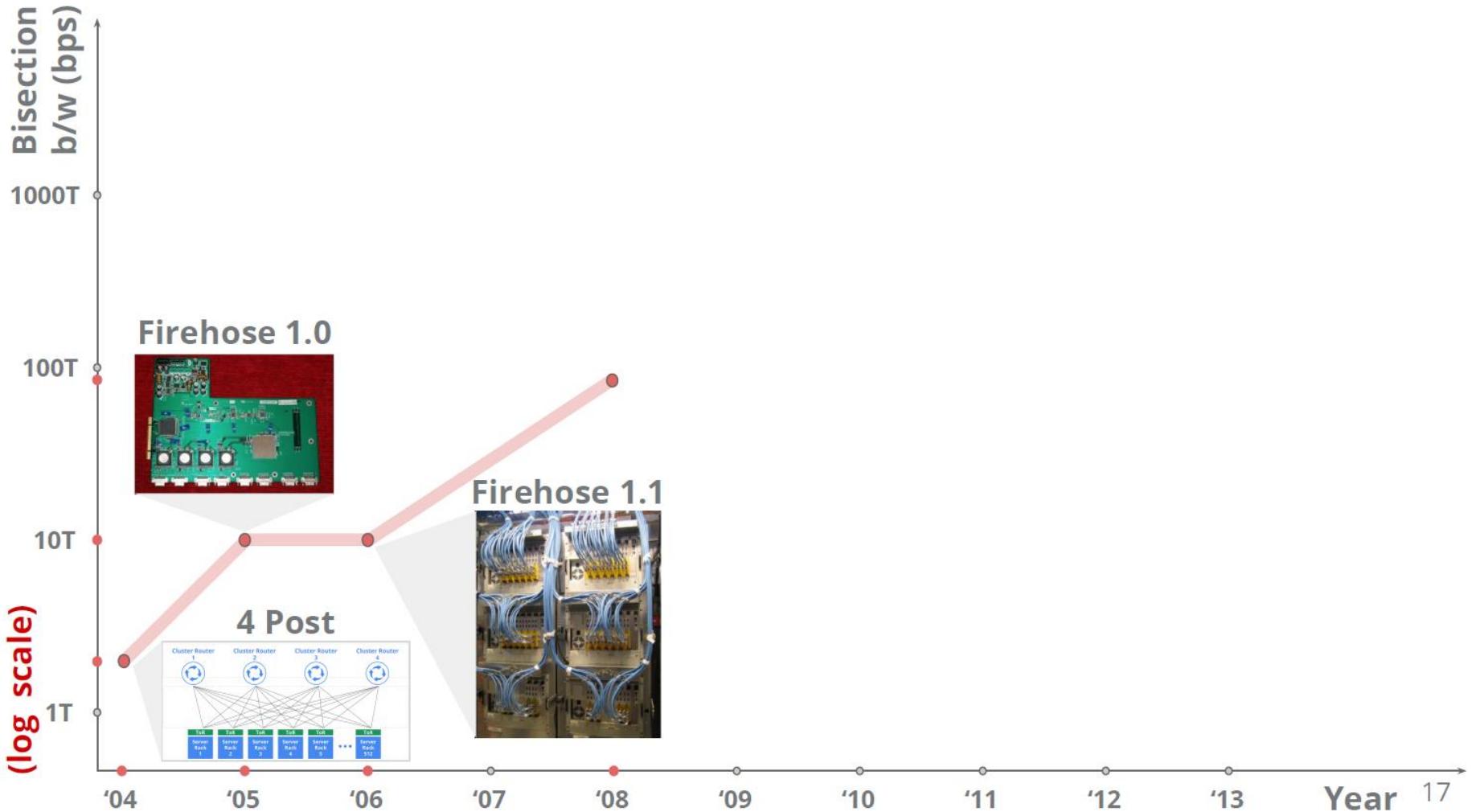


Google's Infrastructure

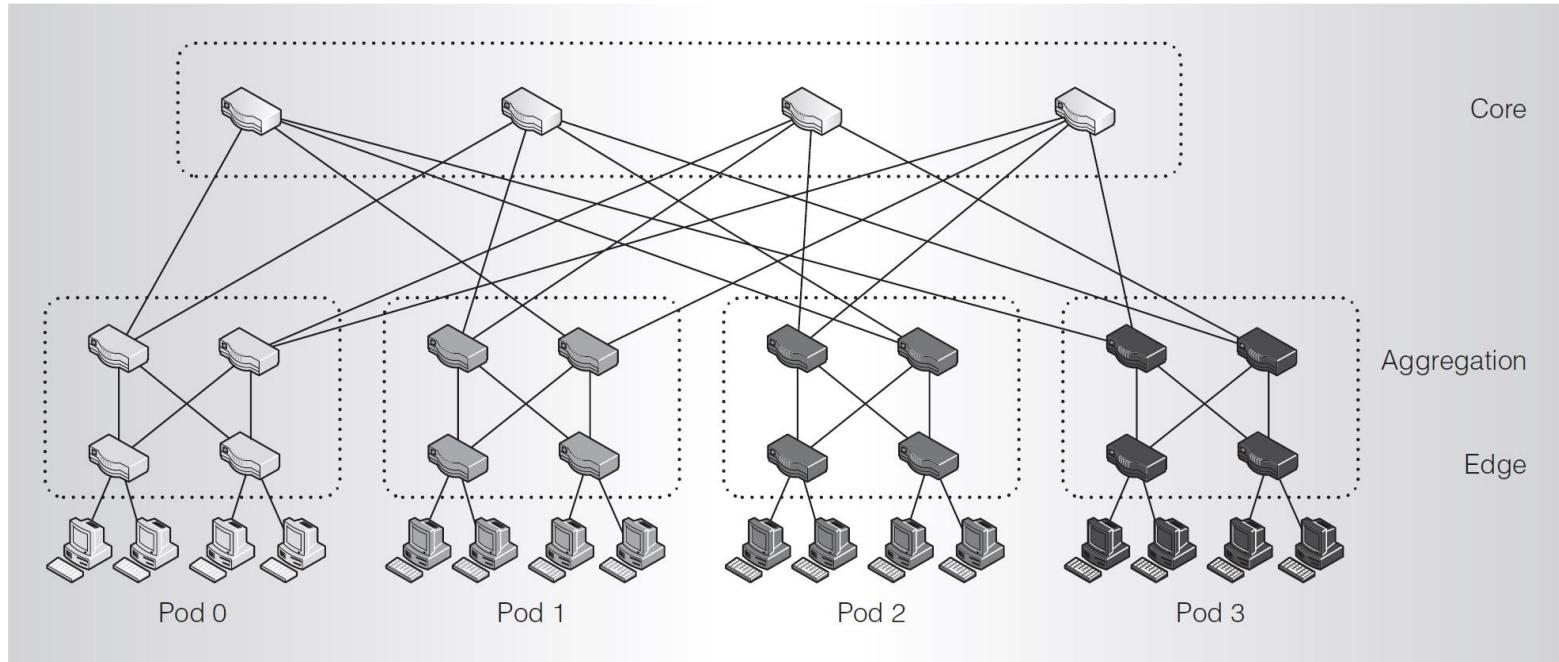
Firehose 1.1



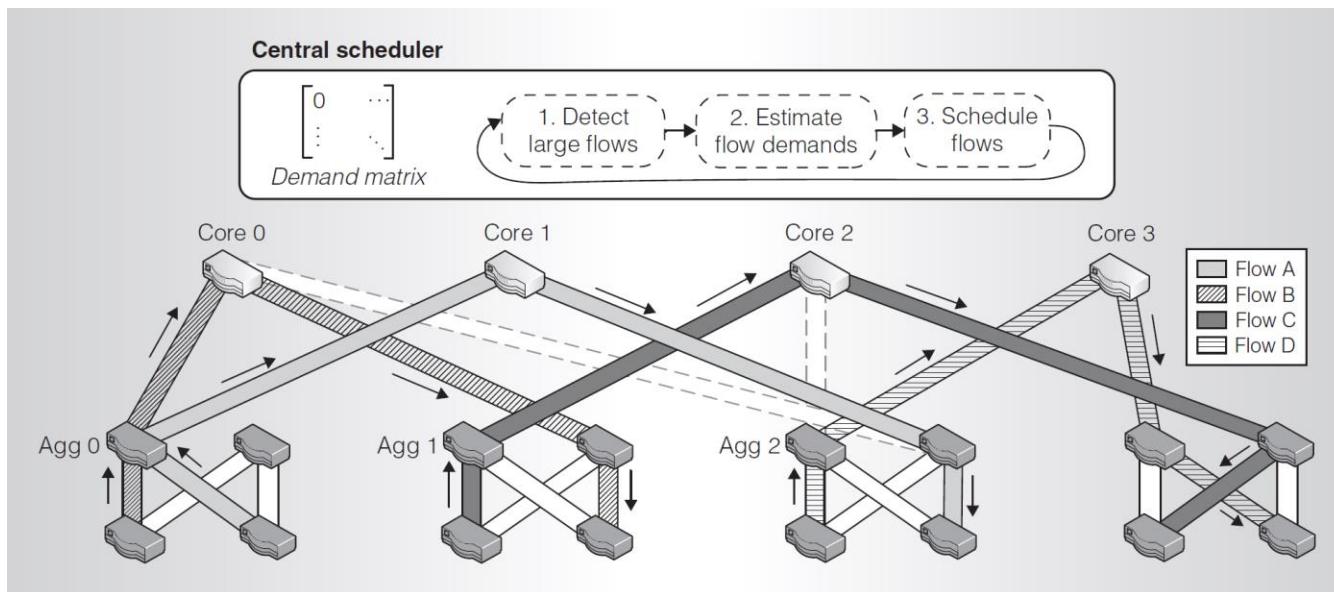
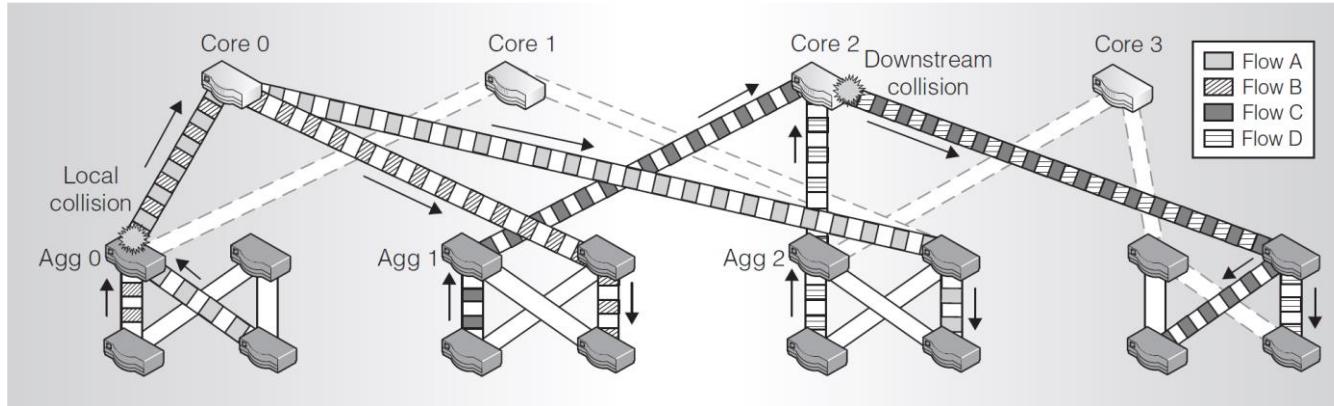
Google's Infrastructure



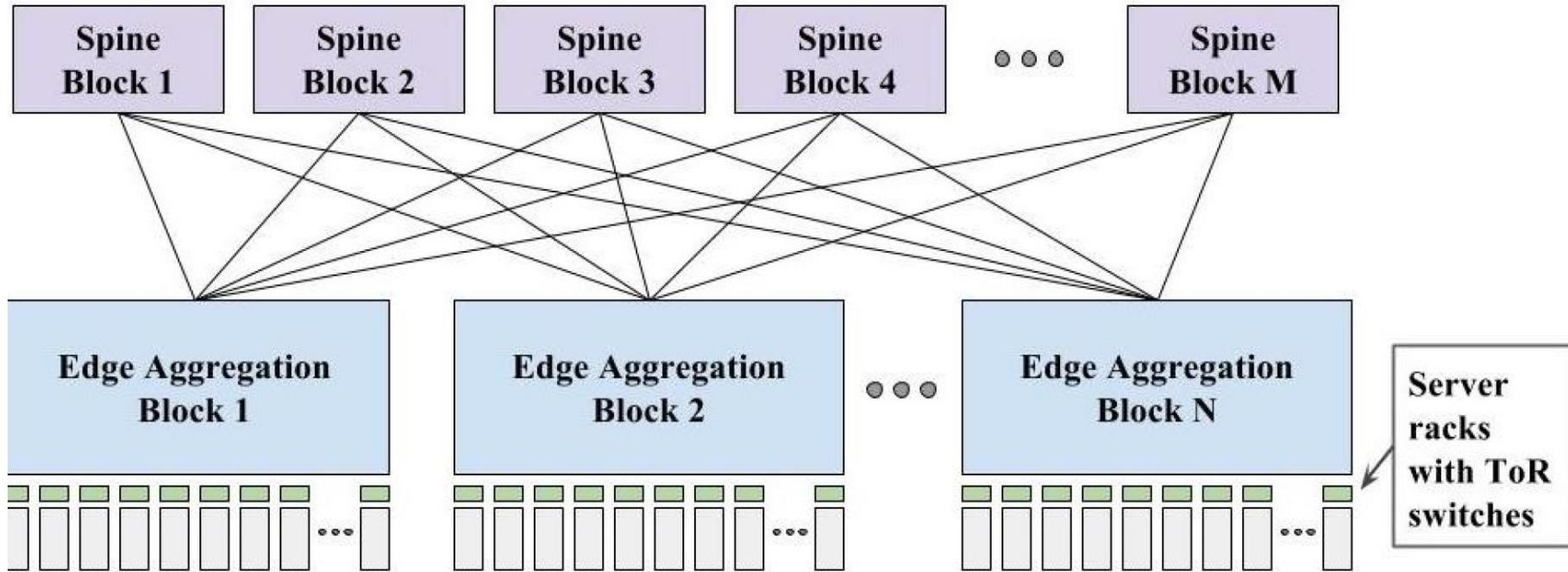
Fat-tree Structure



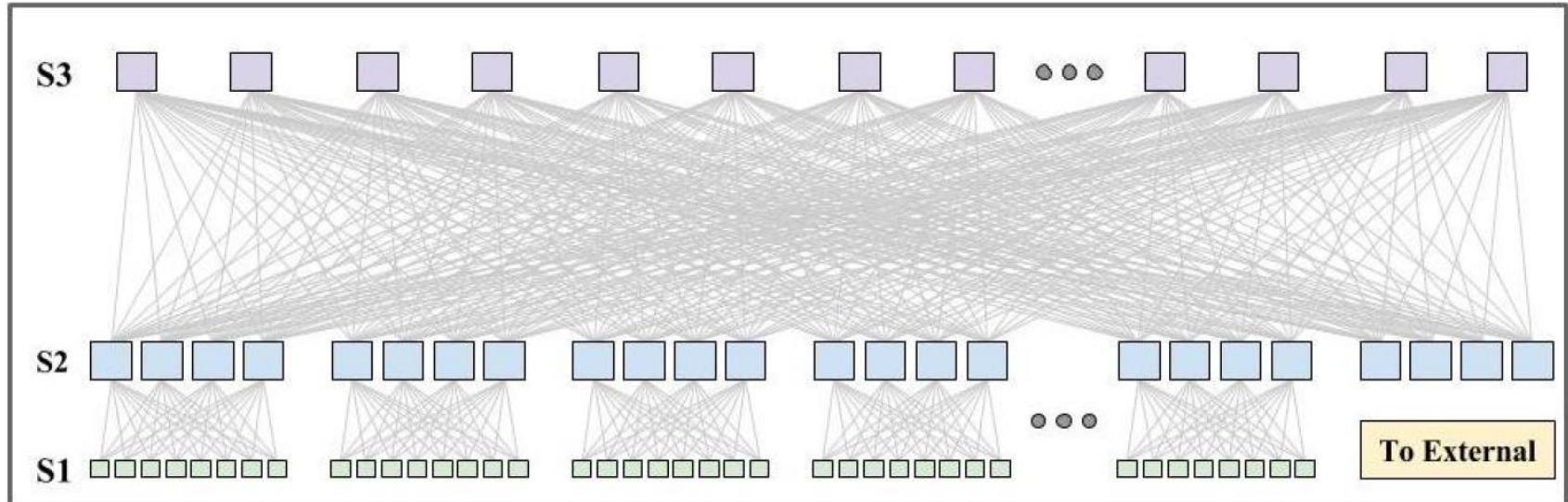
Equal-cost multipath (ECMP) forwarding collisions



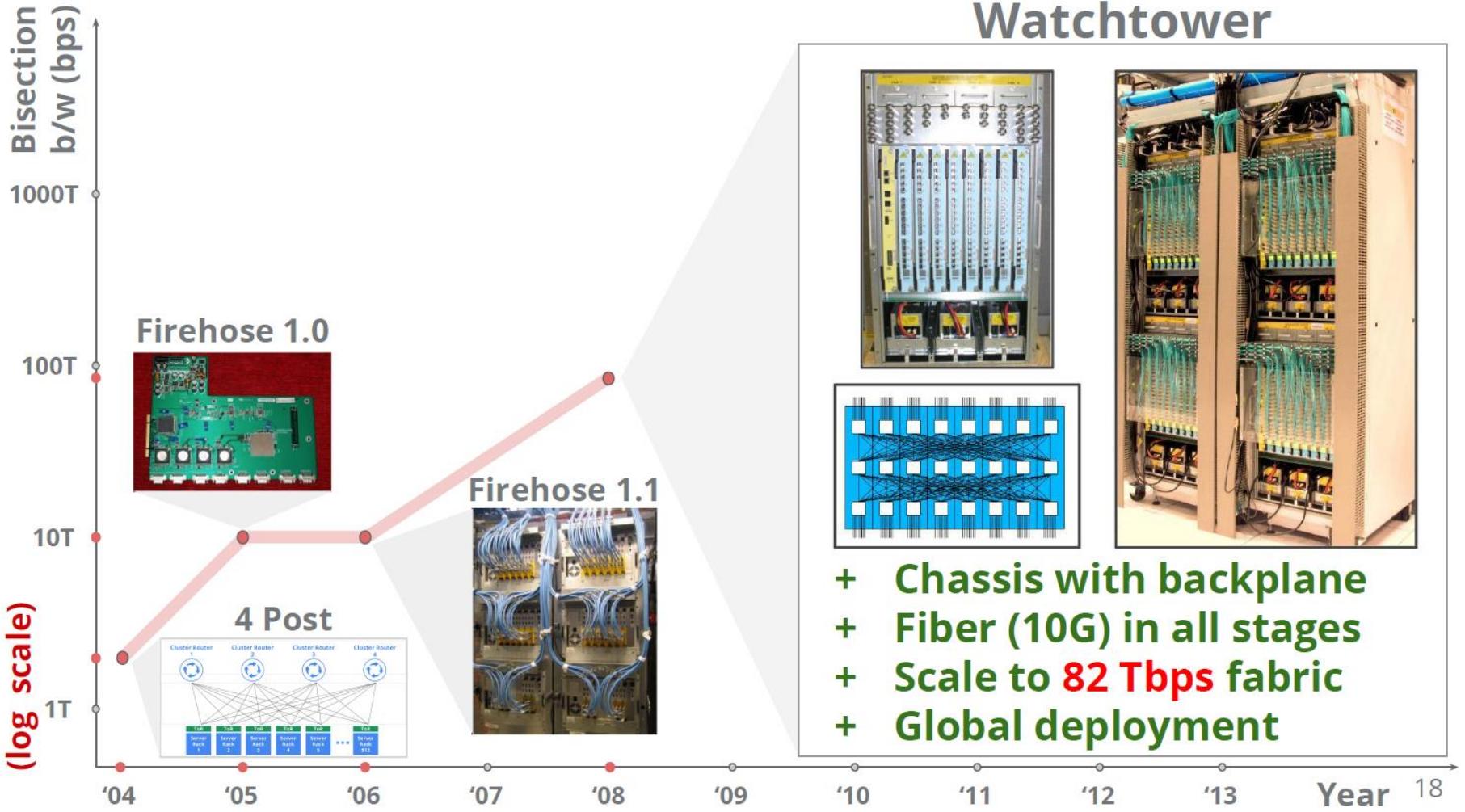
CLOS Architectures



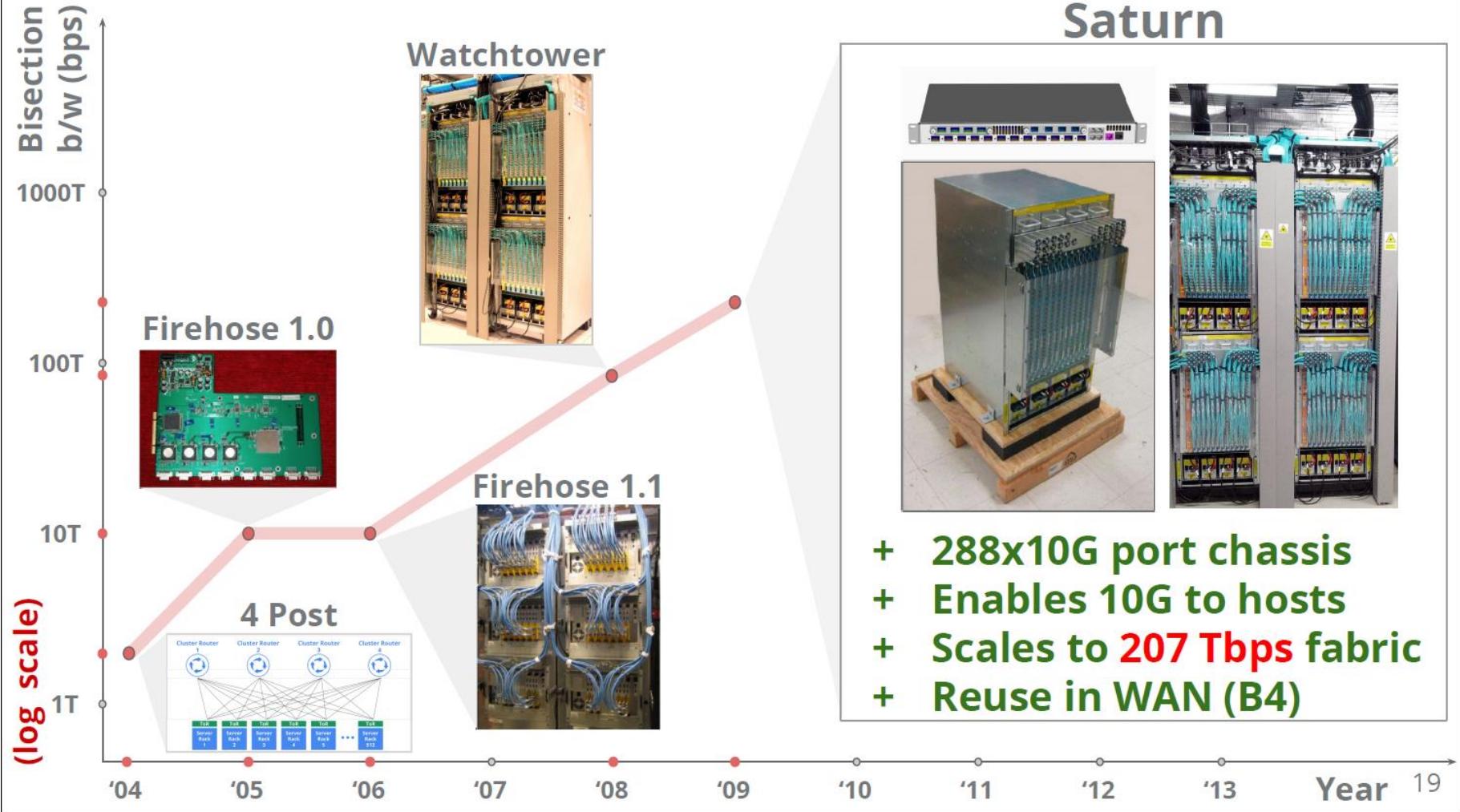
CLOS Architectures



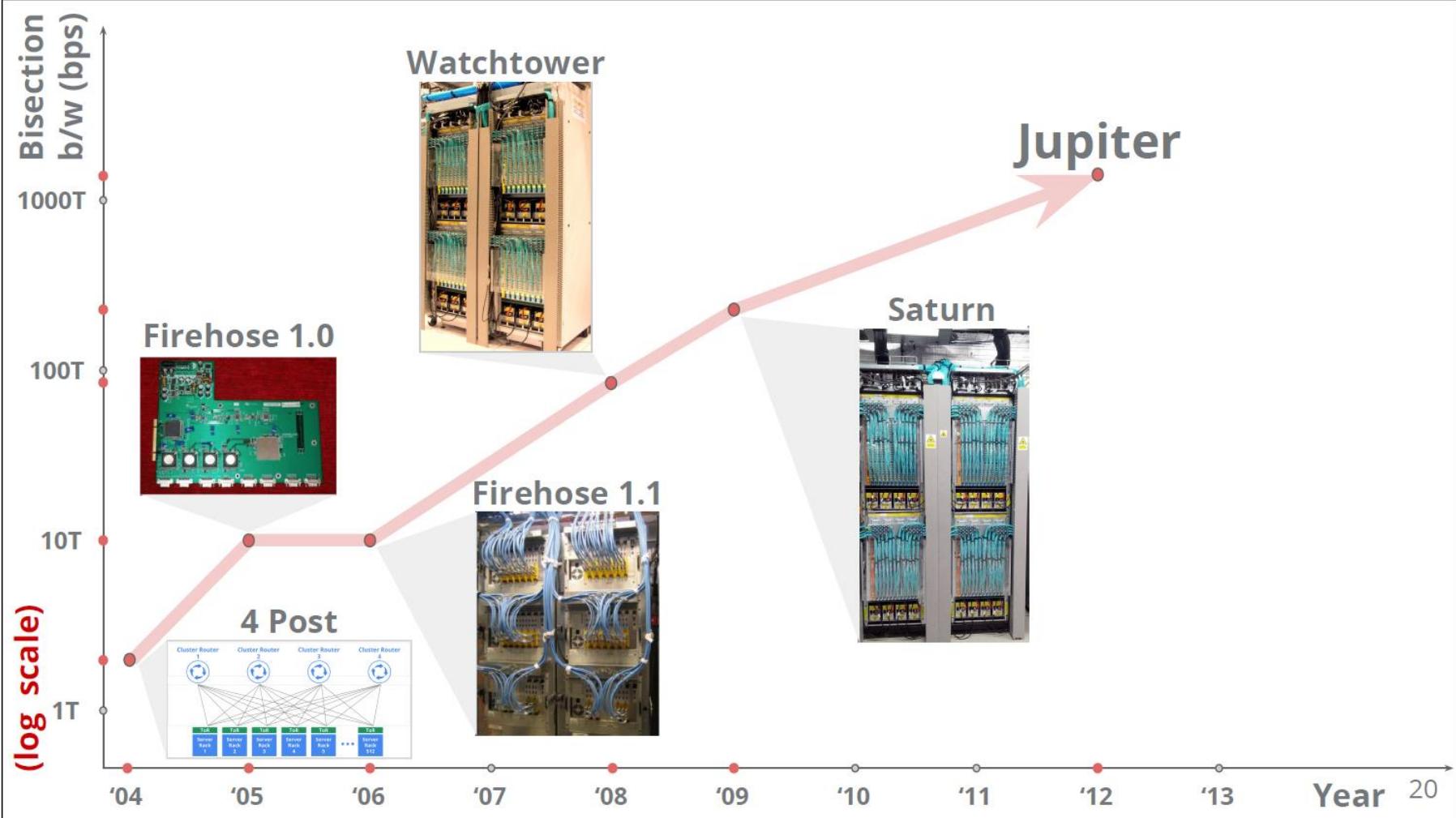
Google's Infrastructure



Google's Infrastructure

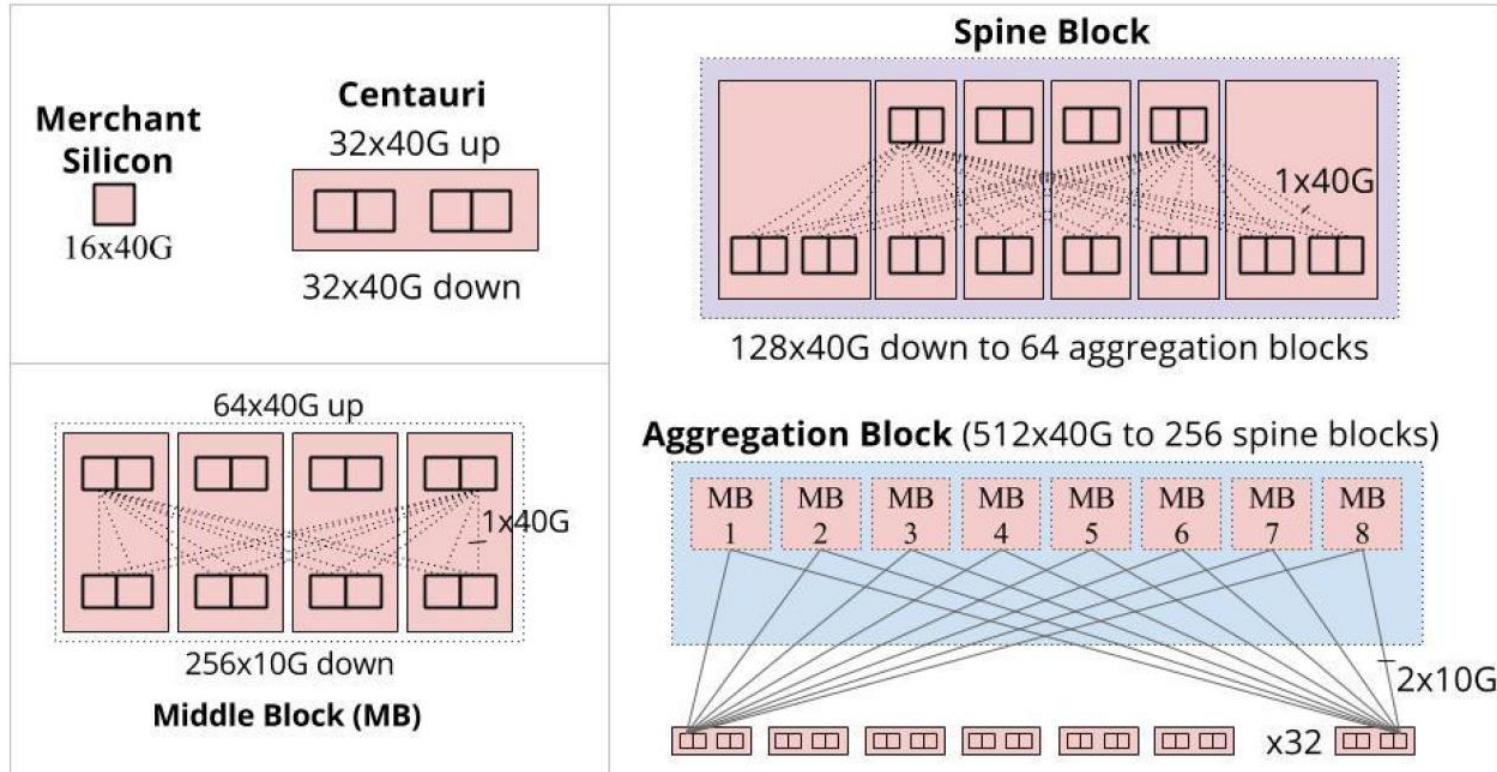


Google's Infrastructure



Google's Infrastructure

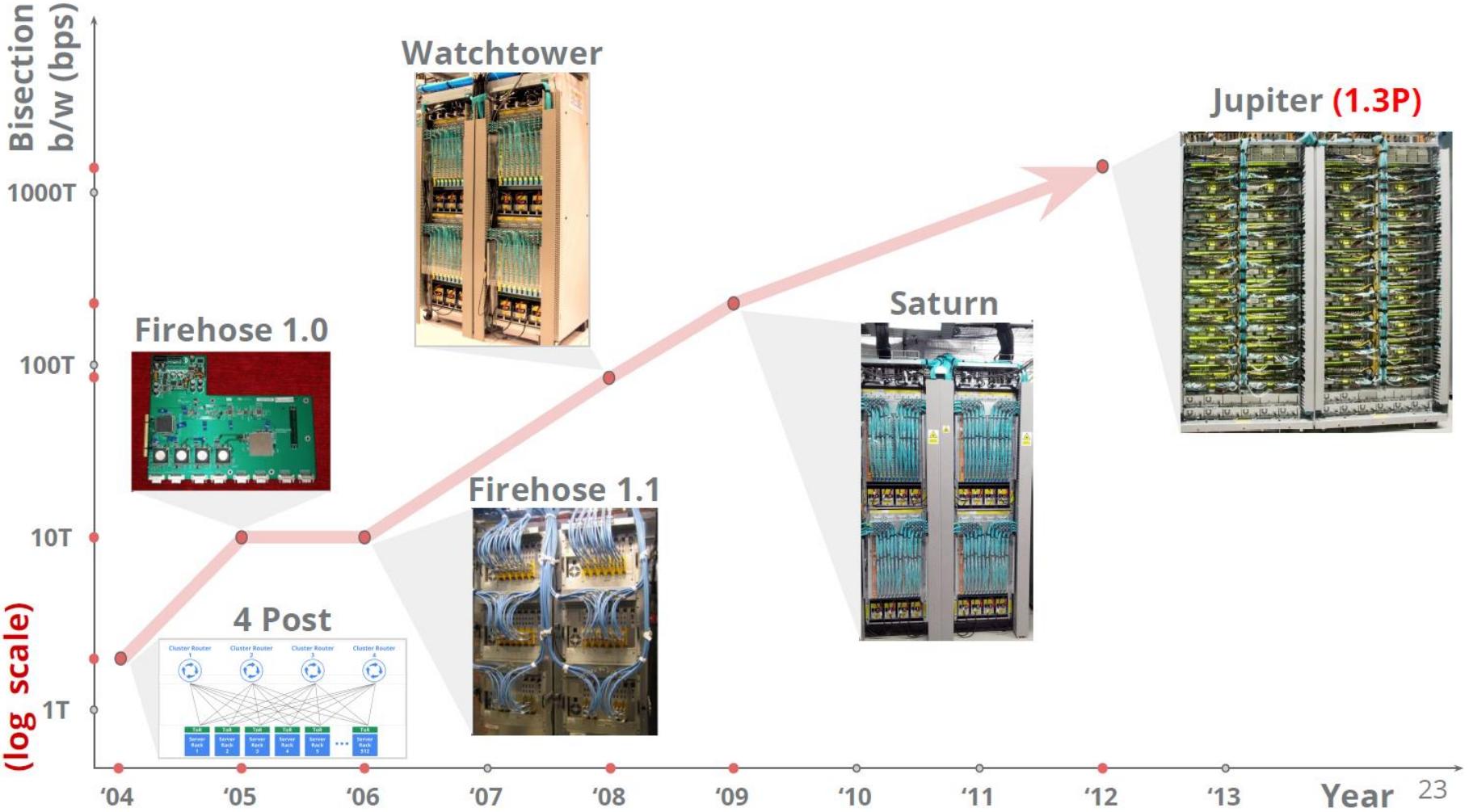
Jupiter topology



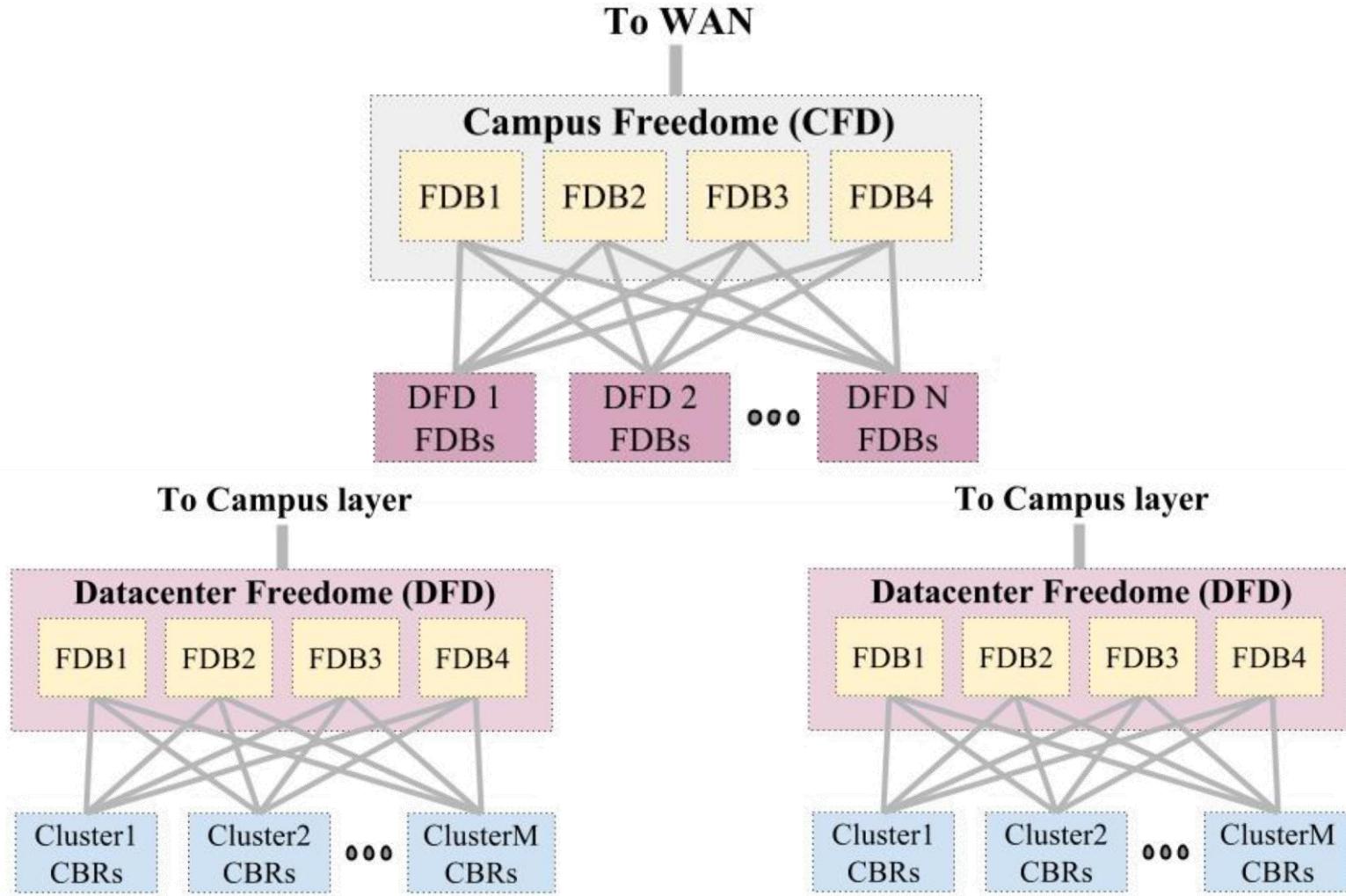
+ Scales out building wide 1.3 Pbps



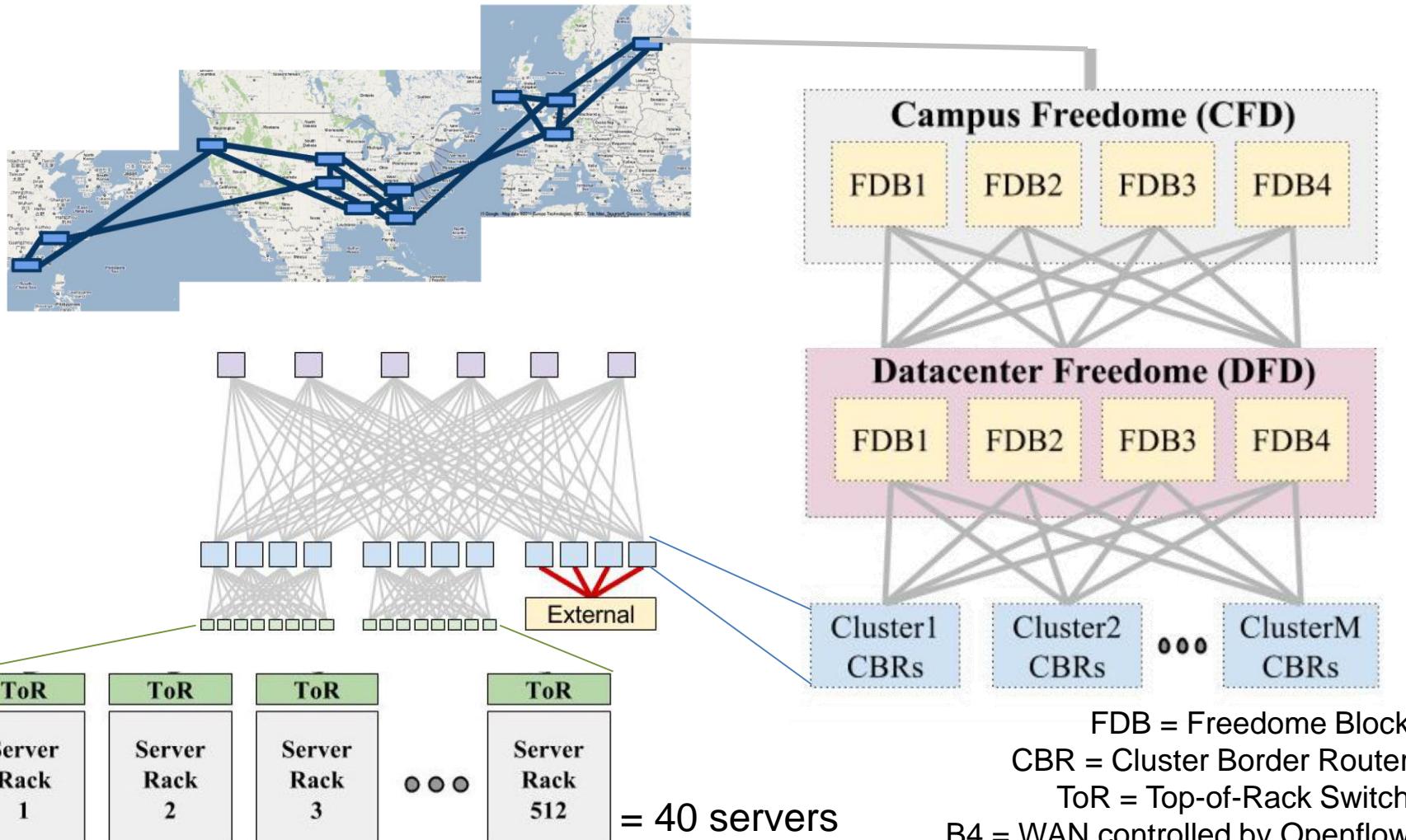
Google's Infrastructure



Google's Infrastructure

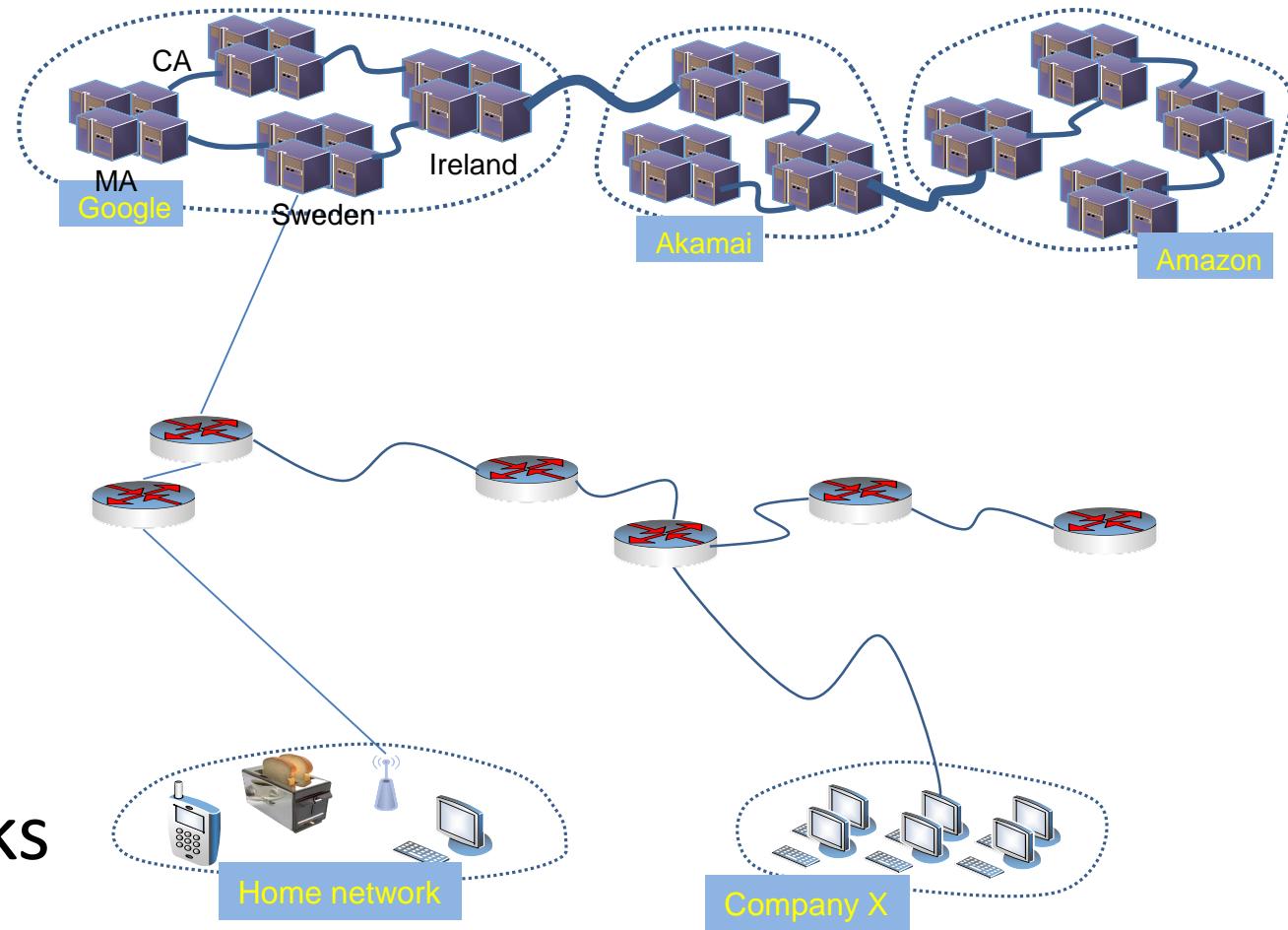


Google's B4 to Jupiter



My view of “future”* networking

- Sets of Datacentres
- Traditional Internet
- Edge networks



*or current?



9.12. & 16.12.
No Lectures & Tutorials



Akamai & IPv6

- Akamai will offer IPv6 services to its entire customer base in April - a long-awaited move that will be a major boon to the adoption rate of the next-generation internet protocol.

Carrying between 20% and 30% of the internet's web traffic on any given day, Akamai is the world's largest content delivery network (CDN). Akamai's engineering team has been working for two years to upgrade its 95,000 servers in 71 countries connected by 1,900 networks to support IPv6.

- Mar 27, 2012
<http://www.techworld.com/news/apps/ipv6-akamai-launch-new-internet-protocol-service-in-april-3347188/>

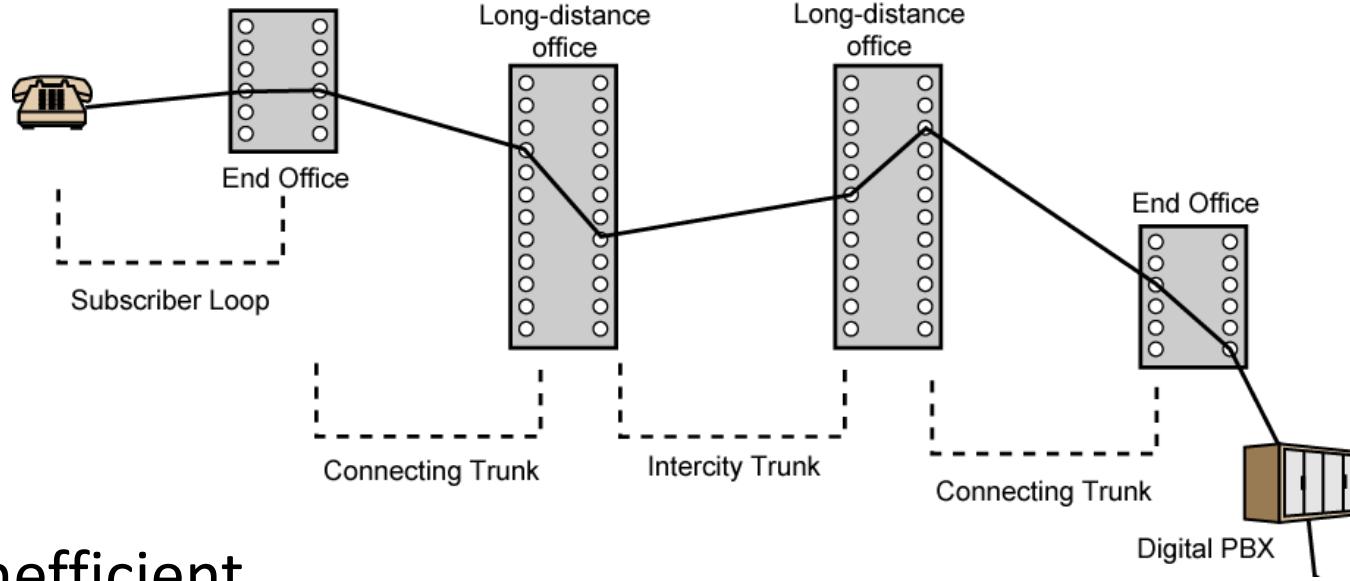


CS2031 Telecommunications II

Circuit Switching & Packet Switching

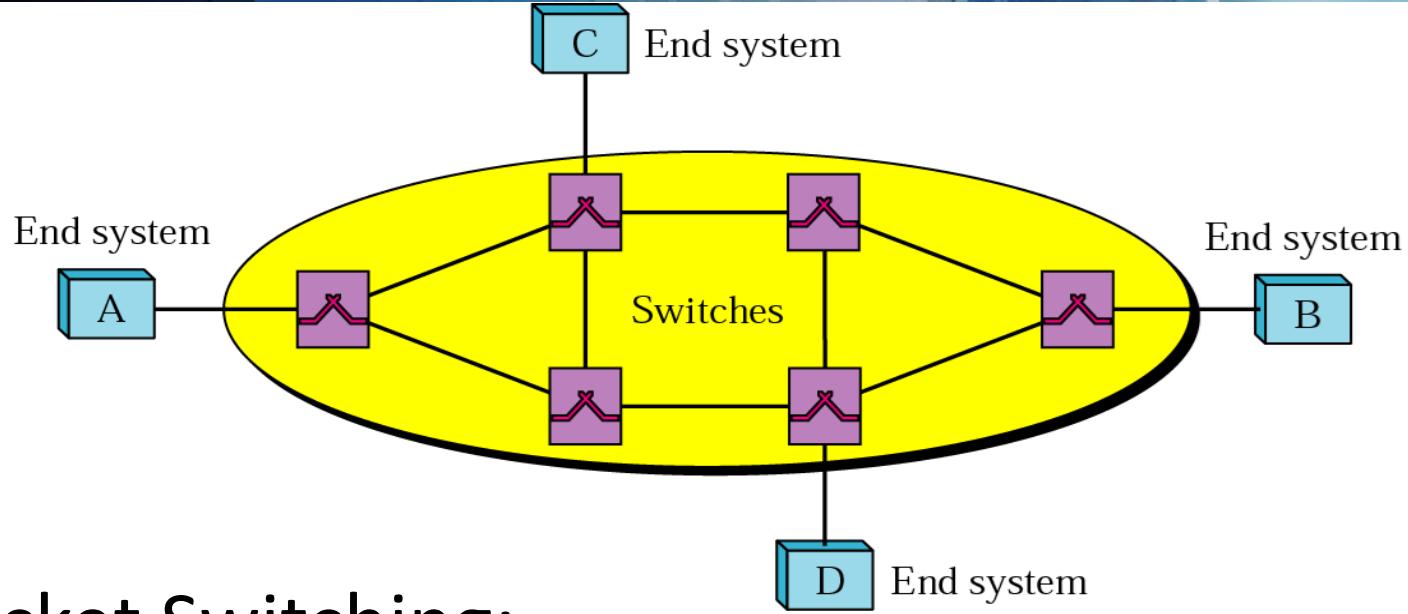


Public Circuit Switched Network



- Inefficient
 - Channel capacity dedicated for duration of connection
 - If no data, capacity wasted
- Set up of connection takes time
- Once connected, transfer is transparent

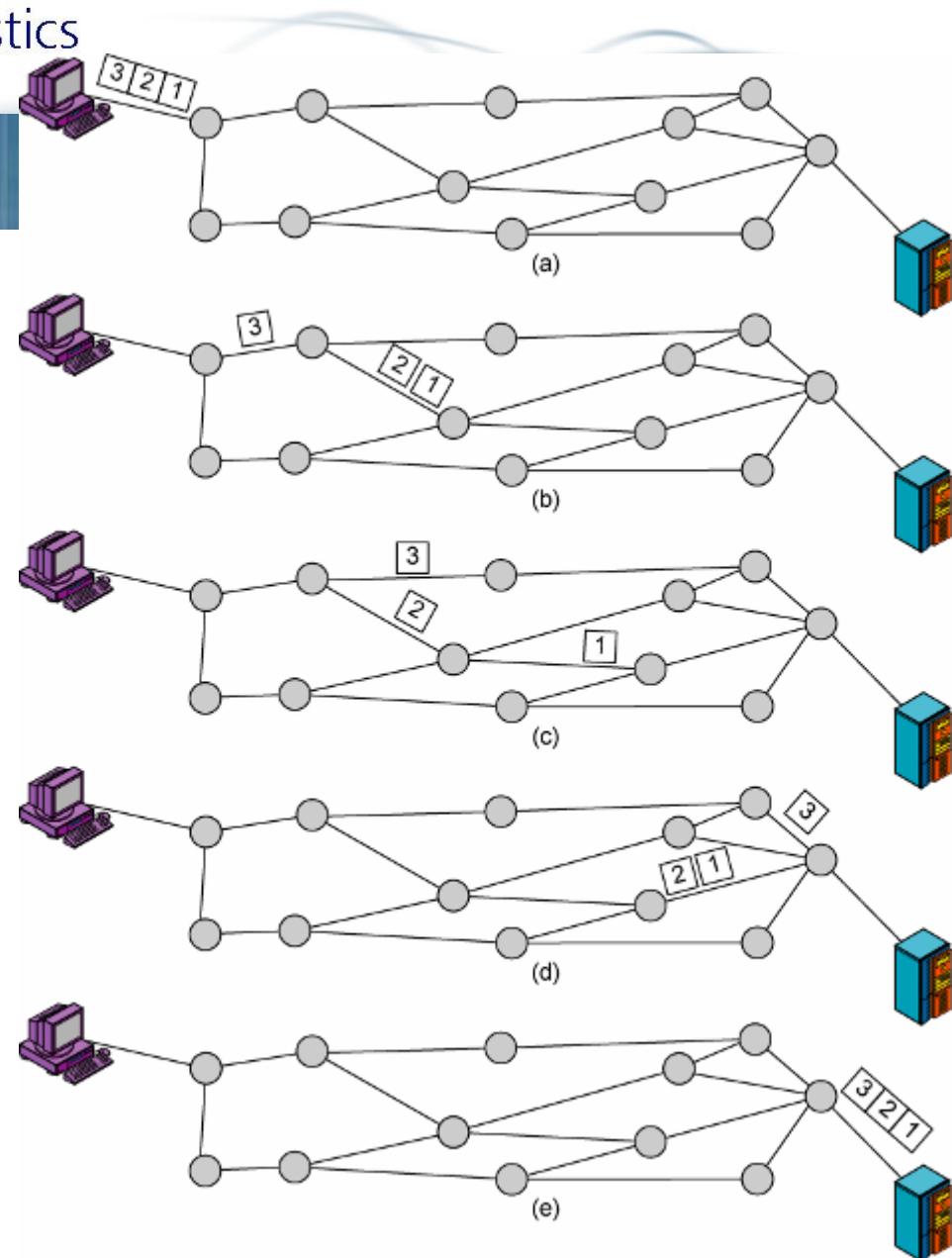
Switched Networks



- **Packet Switching:**
 - Switching decisions are made on individual packets
- **Virtual Circuit Switching:**
 - A circuit is setup explicitly for individual connections

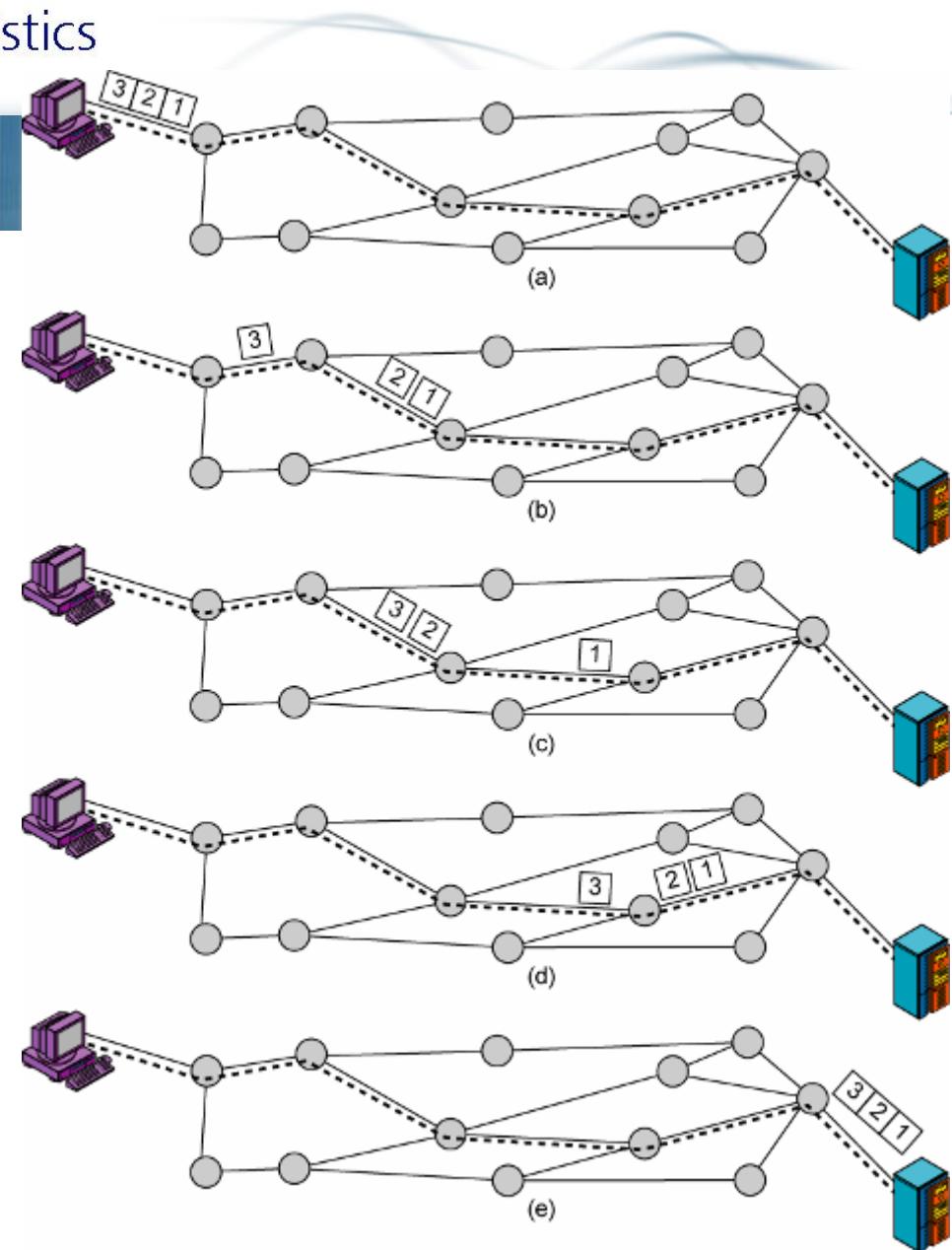
Packet switching

- Frames can be transferred over different paths in the network
- Reliability is generally delegated to higher layers
- Order is not necessarily maintained



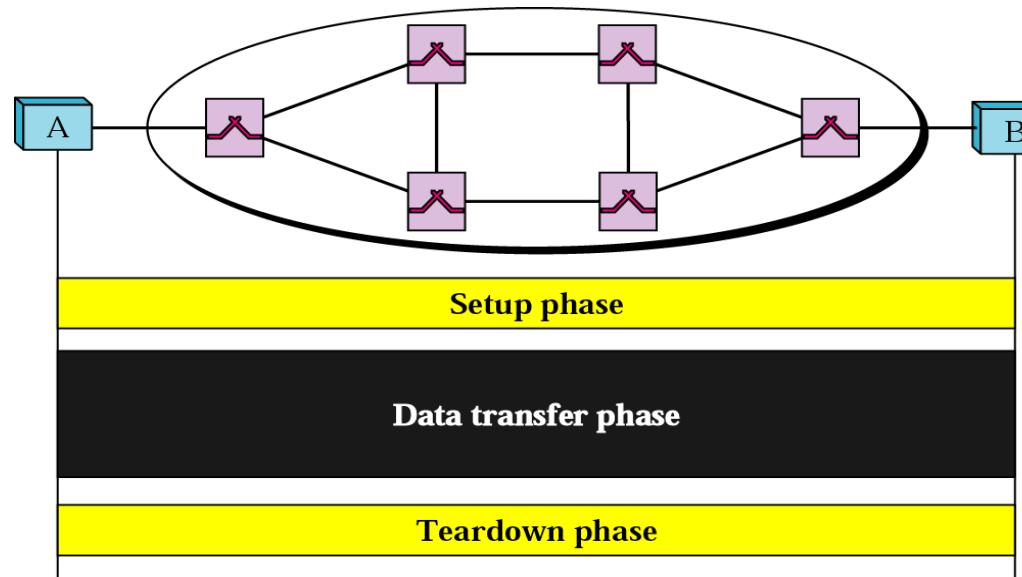
Virtual Circuits

- Connection-oriented communication
- Connection is established before communication
- The network maintains order

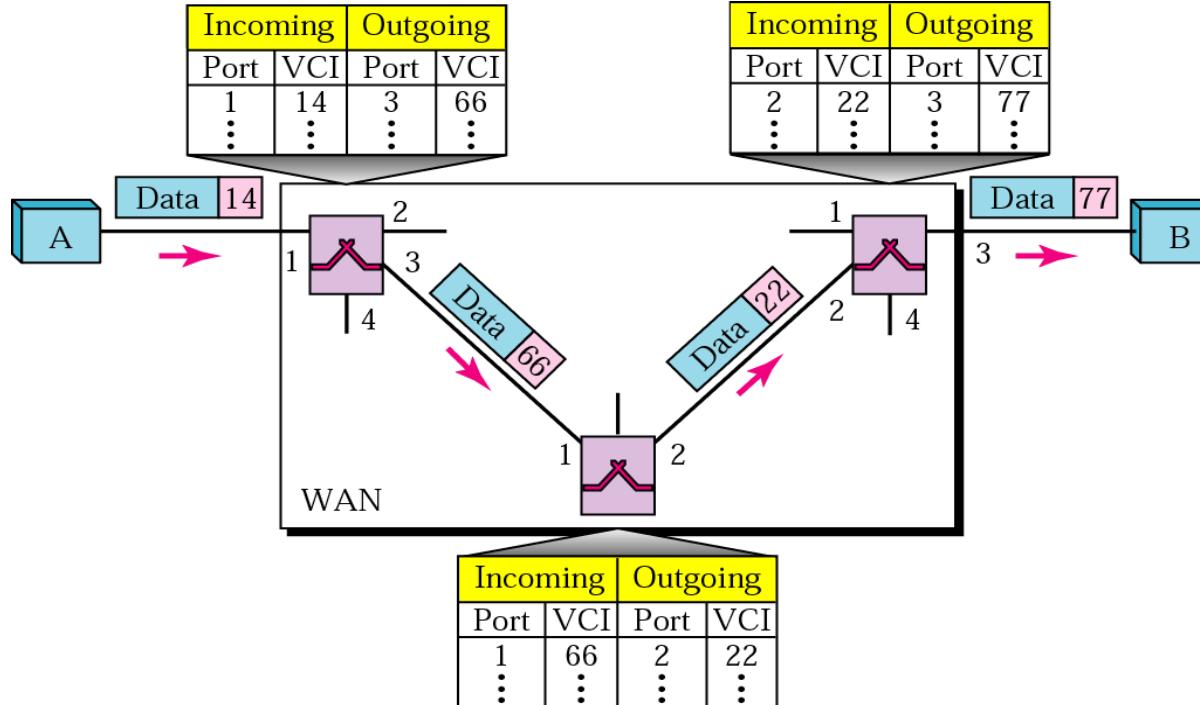


Phases for Virtual Circuit Communication

- Three phases
 - Connection setup
 - Data Transfer
 - Connection termination



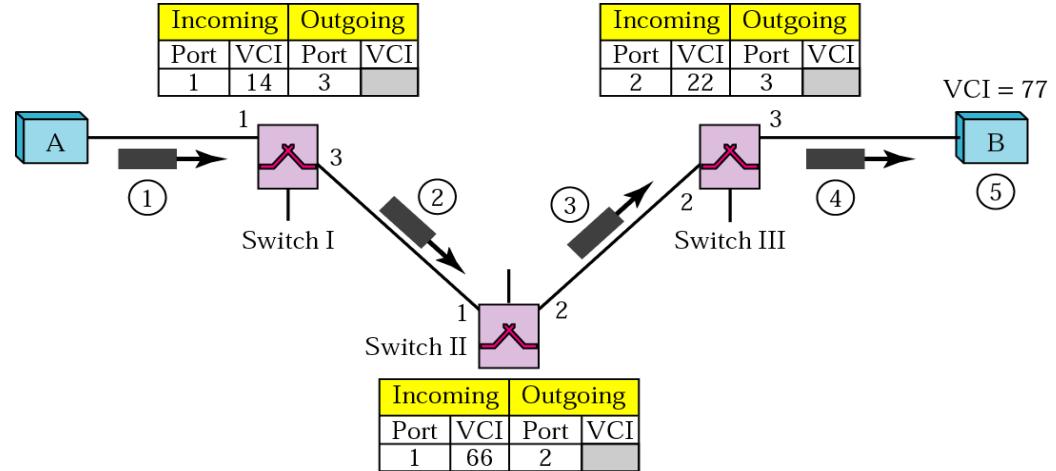
Virtual Circuit Switching



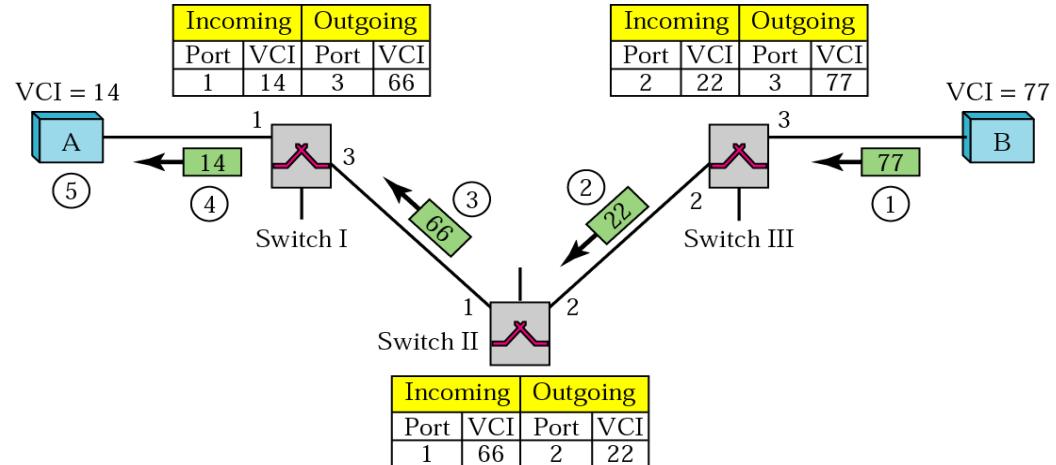
- Every switch maintains a table
 - For duration of communication one entry for incoming and outgoing line
 - Incoming and outgoing line are identified by port number and virtual circuit identifier

Setup Phase

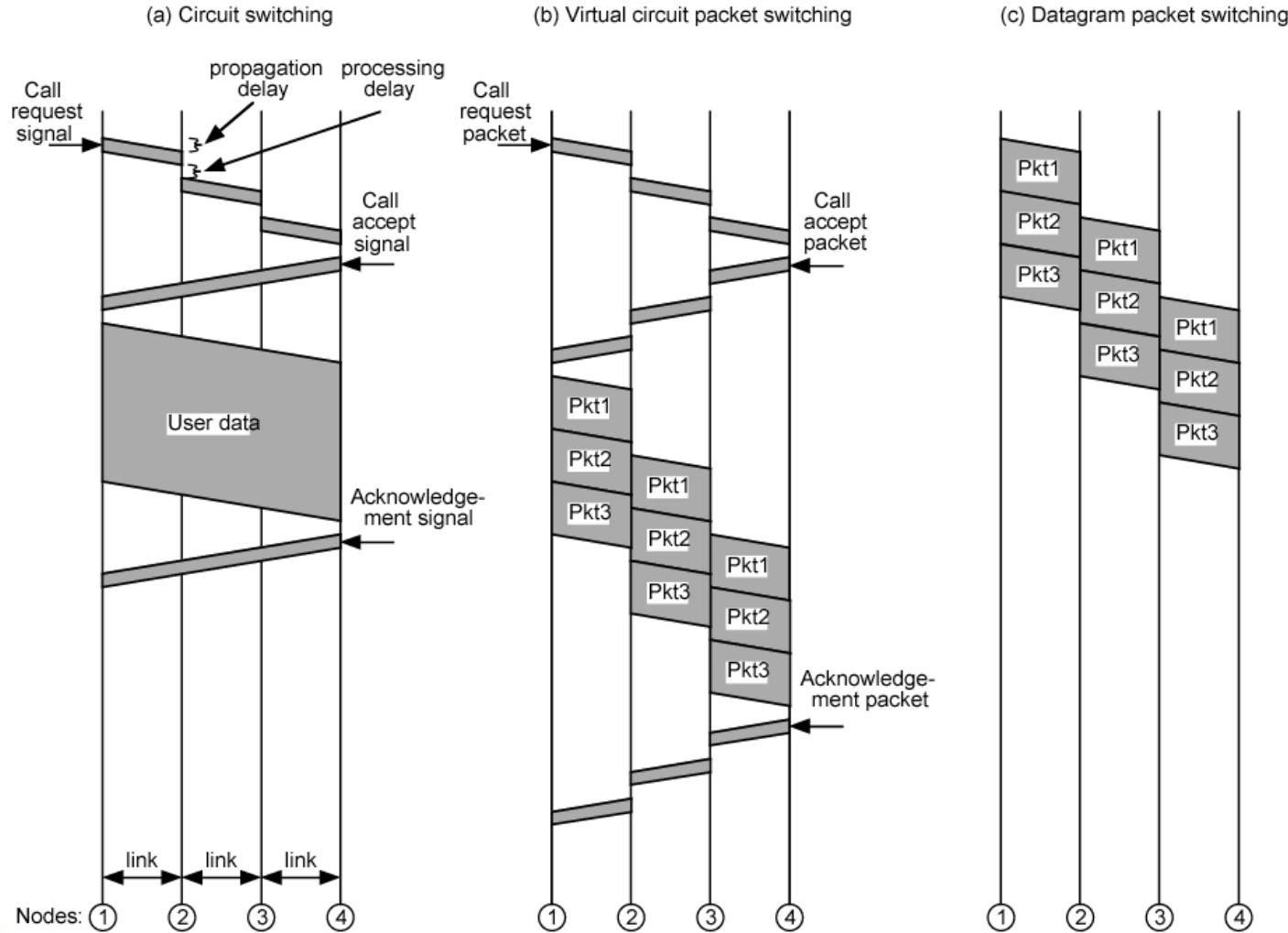
- Setup request



- Setup acknowl.



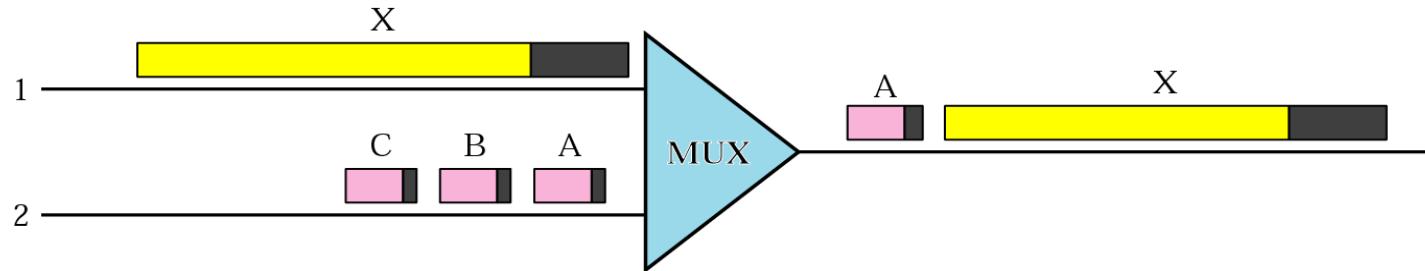
Event Timing



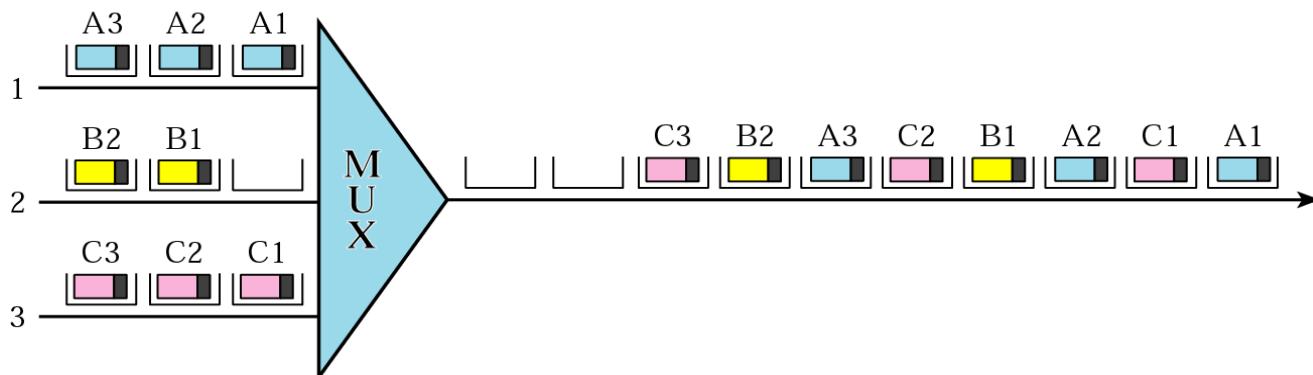
Asynchronous Transfer Mode (ATM)

- Example of virtual circuit switching
 - Cell-Switching
- Similarities between ATM and packet switching
 - Transfer of data in discrete chunks
 - Multiple logical connections over single physical interface
- In ATM flow on each logical connection is in fixed sized packets called cells
- Minimal error and flow control
 - Reduced overhead

Motivation for ATM

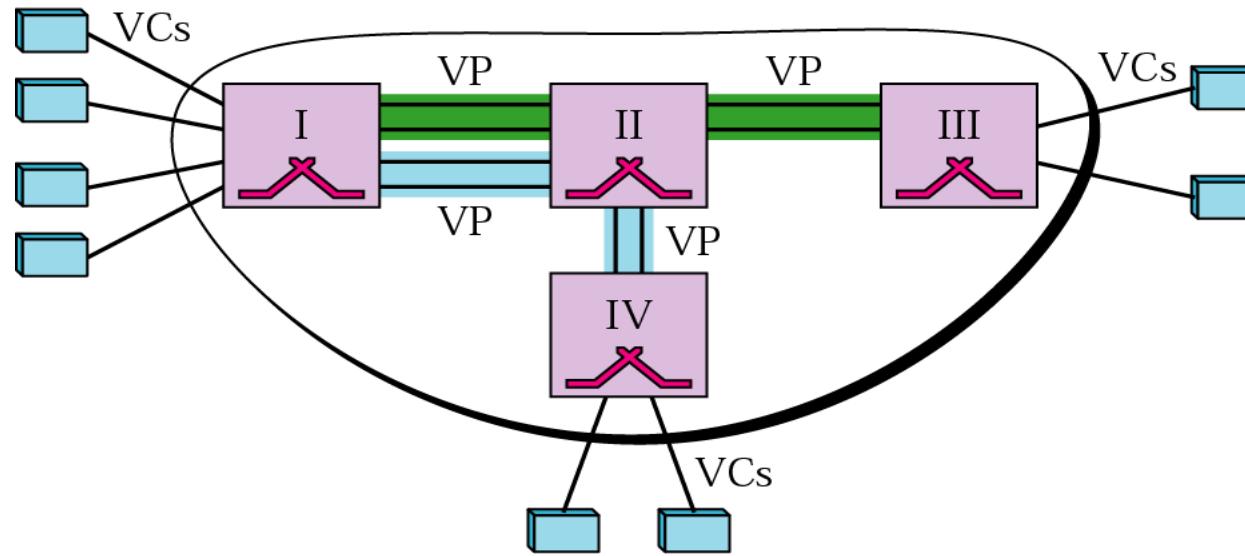


- Frames at a switch may be handled in any order and occupy switch for underspecified time



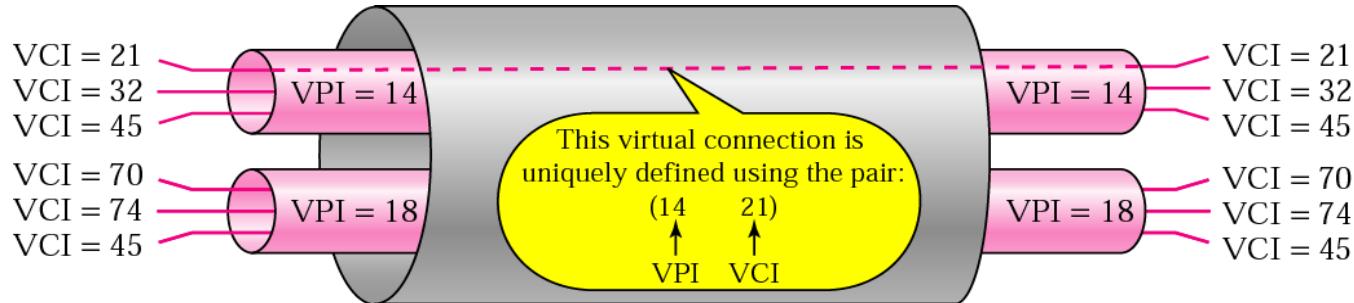
- Small, fixed-size frames allow simple, fast switches

Virtual Circuits / Virtual Paths

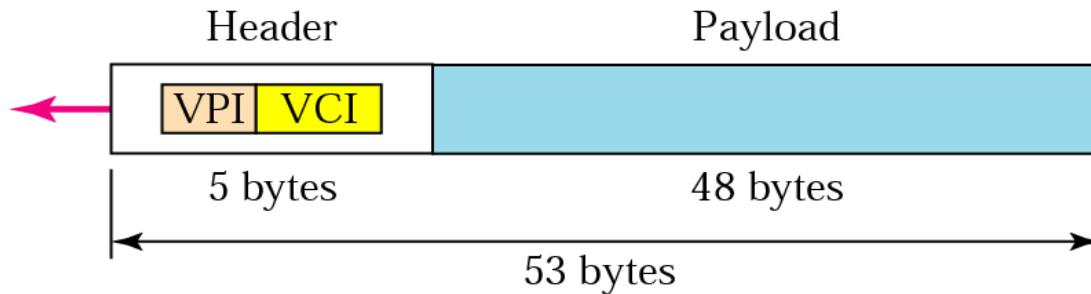


- Virtual circuits are collected into virtual paths

ATM Packet

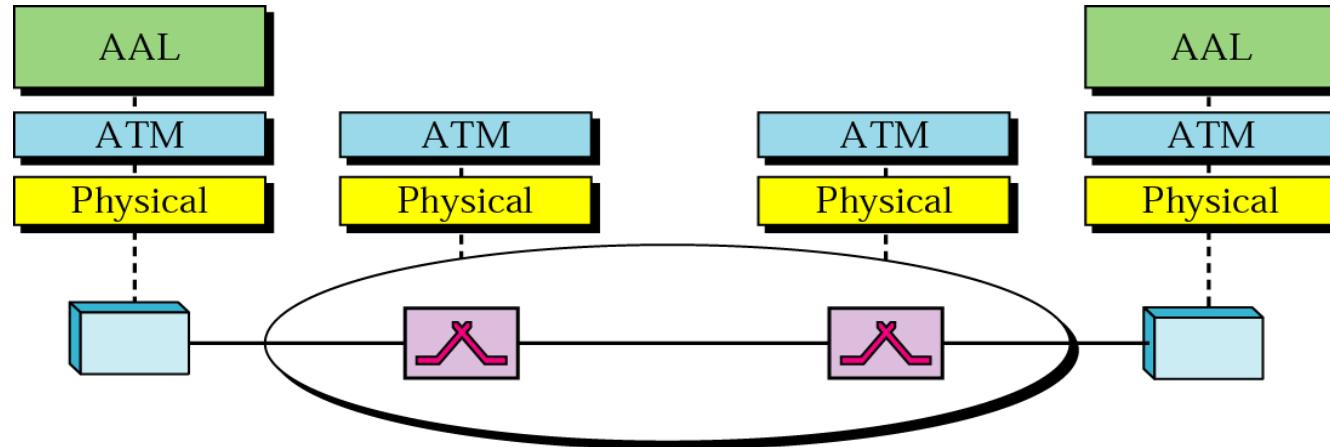


- Connection is specified by combination of Virtual Path ID and Virtual Circuit ID



- Every frame is exactly 53 bytes

Application Adaptation Layer (AAL)



- ATM defined a number of AALs for various purposes (each has its own header format):
 - AAL1: Constant bit rate e.g. multimedia
 - AAL2: Variable-data-rate
 - AAL3/4: Connection-oriented data services
 - Sequencing and Error Control
 - AAL5: Simple and efficient adaptation layer (SEAL)

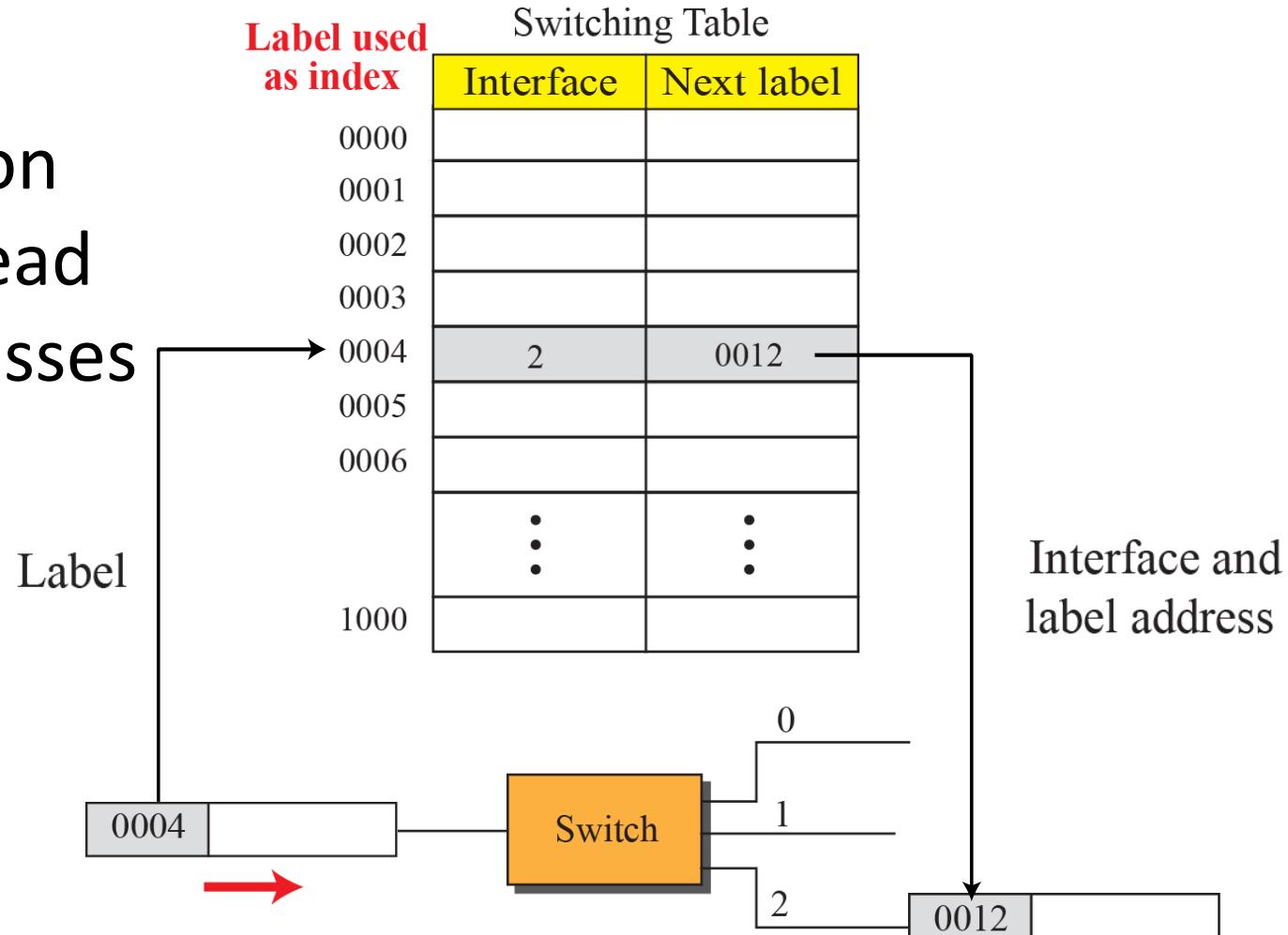
ATM – It Didn't Happen

- From Tanenbaum:

“ATM was going to solve all the world’s networking and telecommunications problems by merging voice, data, cable television, telex, telegraph, carrier pigeons, ...”
- It didn’t happen:
 - Bad Timing
 - Technology
 - Implementation
 - Politics

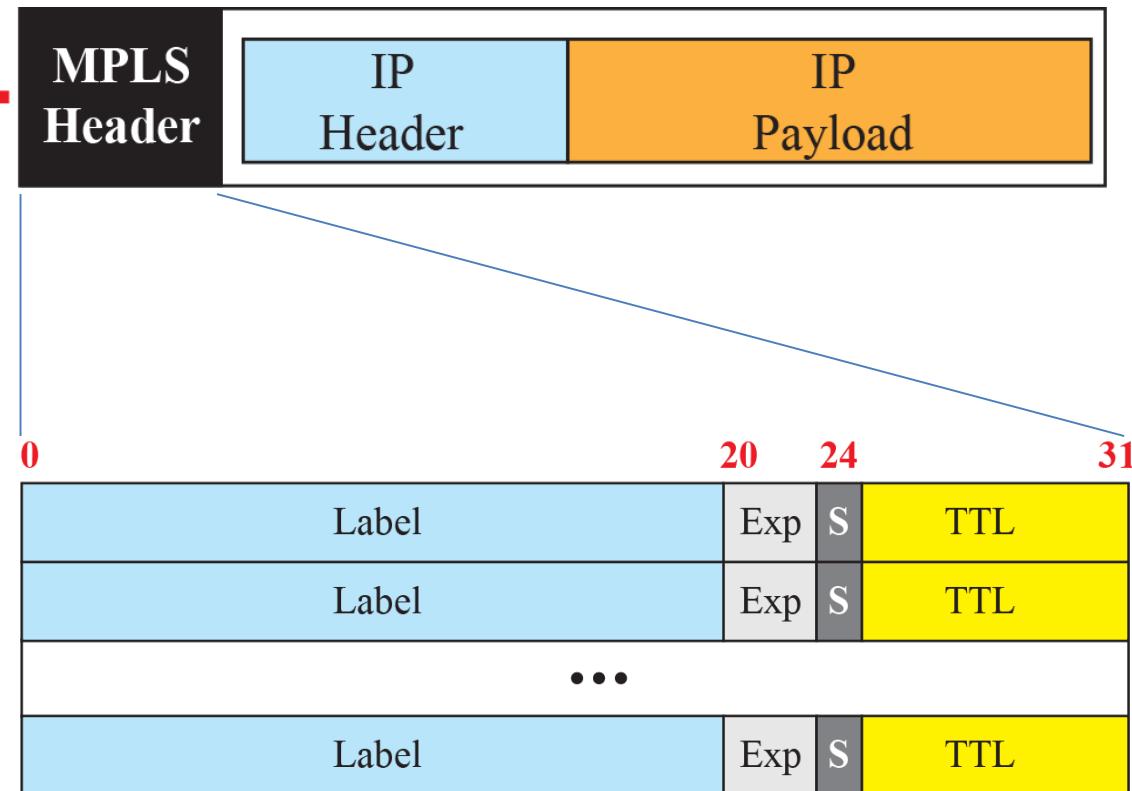
Multiprotocol Label Switching (MPLS)

- Enables switching on labels instead of IP addresses

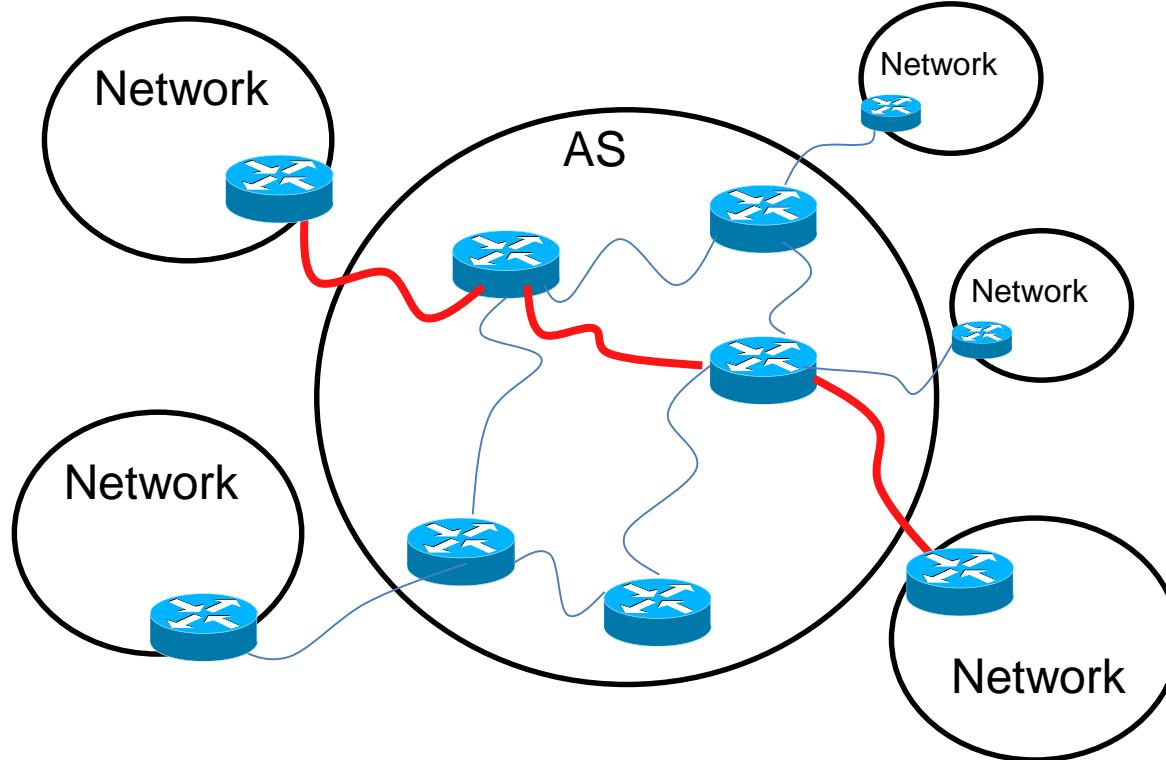


MPLS Header

- MPLS header
as stack of
labels

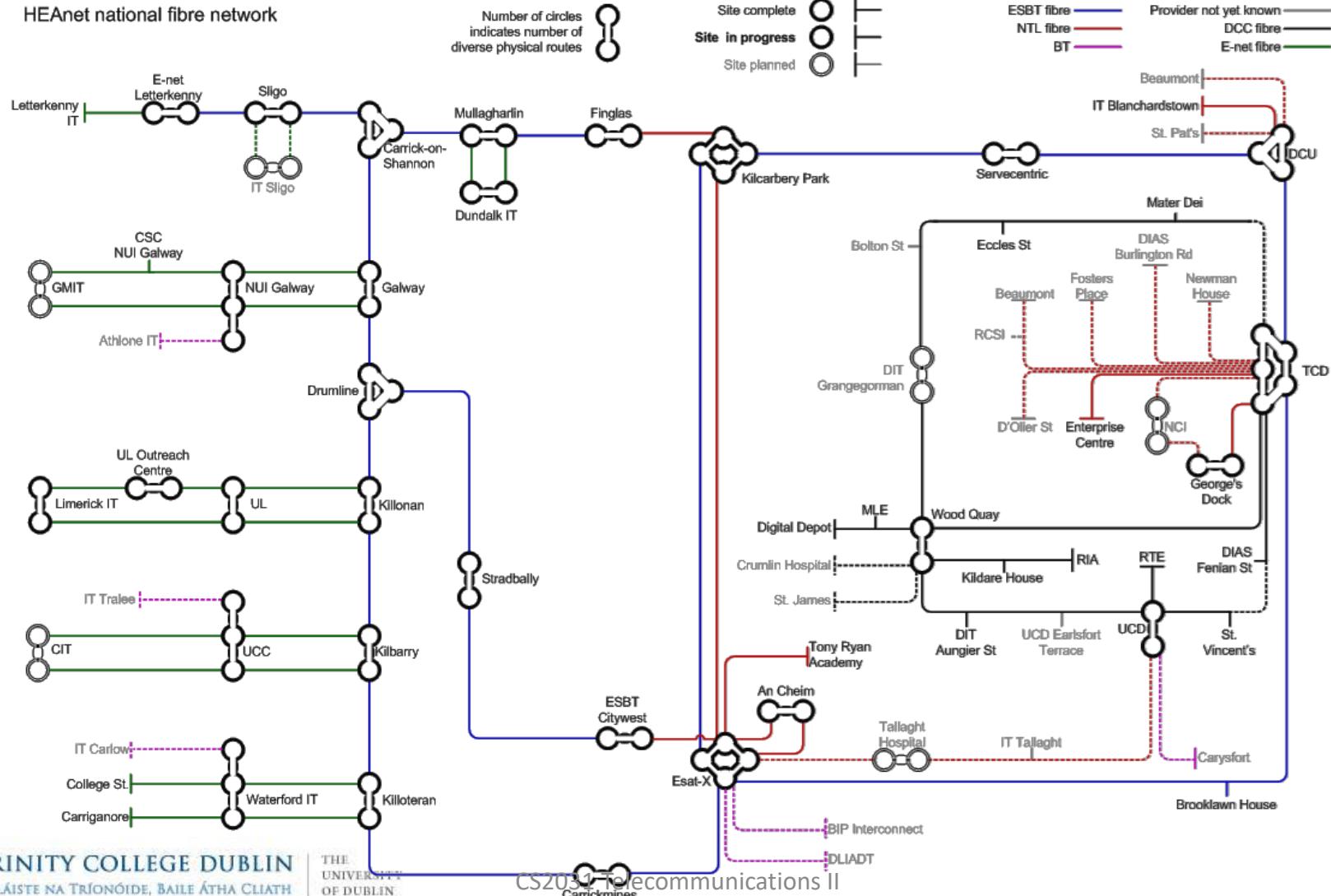


MPLS Use Case



- Creating a virtual network

HEAnet Fibre Network



Summary: Virtual Circuit Switching – ATM

- Virtual Circuit Switching
 - Preplanned route established before any frames sent
 - Call request and call accept frames establish connection (handshake)
 - Each frame contains a virtual circuit identifier instead of destination address
 - No routing decisions required for each frame
 - Clear request to drop circuit
 - Not a dedicated path
- Asynchronous Transfer Mode (ATM)
 - Example for virtual circuit switching
 - Cells consist of 5-byte header and 48-byte payload
 - Circuits identified by virtual circuit ID and virtual path ID
 - Application adaptation layer (AAL) for specific application areas