

# Applications of SDN

---

B4: THE GOOGLE USE CASE

# Applications of SDN

---

B4: THE GOOGLE USE CASE

# Initial SDN Adopters

---



Cloud providers, to enable efficient upscaling of data centers

Use cases:

- **Switching fabrics – SDN within data centers:**
  - From black box to white box switches, server interconnect controlled by software
- **Virtual networks – SDN to set up, manage, tear down:**
  - Easy-to-use VPNs and VLANs      VPN L3 / VLAN L2
- **Wide area networks – traffic engineering between sites:**
  - Dynamically managing end-to-end paths in links between data centers

handle traffic between nodes in DC

efficiently link in WAN

Google WAN

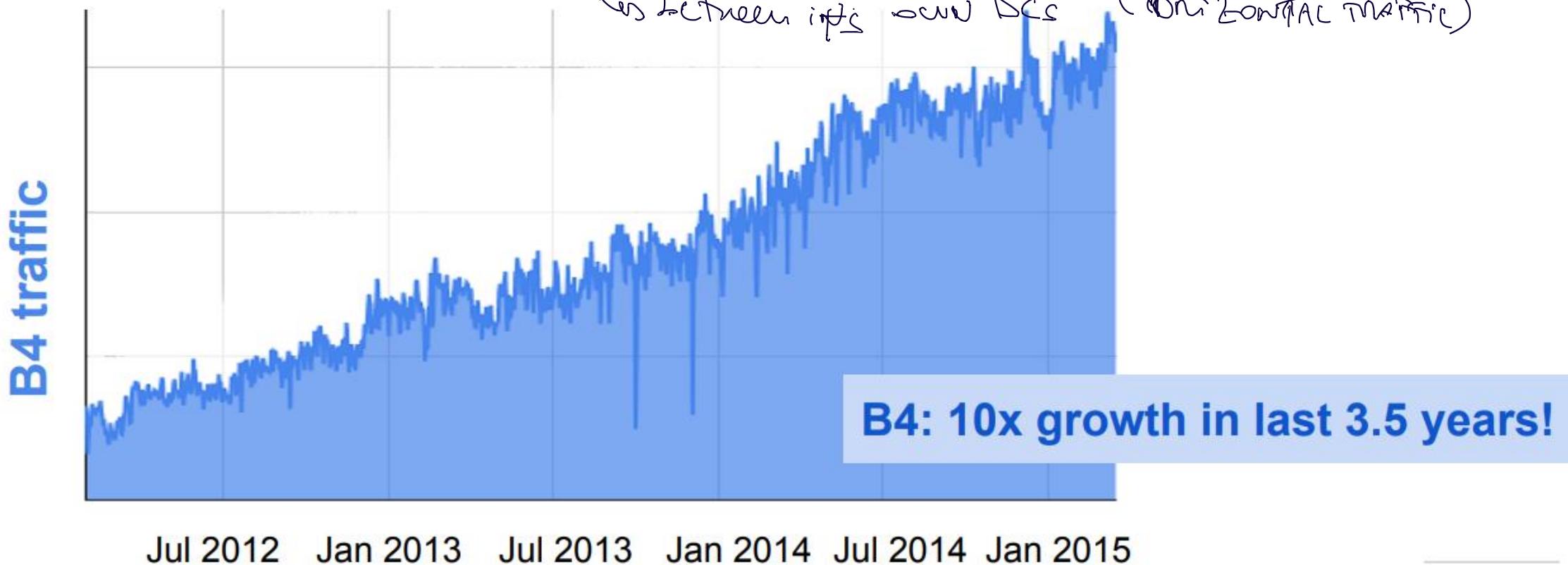
used to inter-DC traffic

# Two Backbones



Two separate backbones:

- 1) • B2: Carries Internet facing traffic → Growing faster than the Internet
- 2) • B4: Inter-datacenter traffic → More traffic than B2, growing faster than B2



# B4 – Google's software defined WAN

- Google's globally deployed WAN, datacenter backbone
  - private WAN connecting its data centers (NOT: Internet-facing backbone for user traffic)
- Modest number of sites (~ dozen)

*Between  
DCs*



# Motivation for Backend Backbone

Data centers deployed across the world

- Serve content with geographic locality (delay threshold)
- Replicate content for fault tolerance
  - Synchronized replication

Need a network to connect these data centers to one another

- Not on the public Internet
- Cost effective network for high volume traffic
- Bursty/bulk traffic (not smooth/diurnal)

minimize cost, try to use links as much as possible

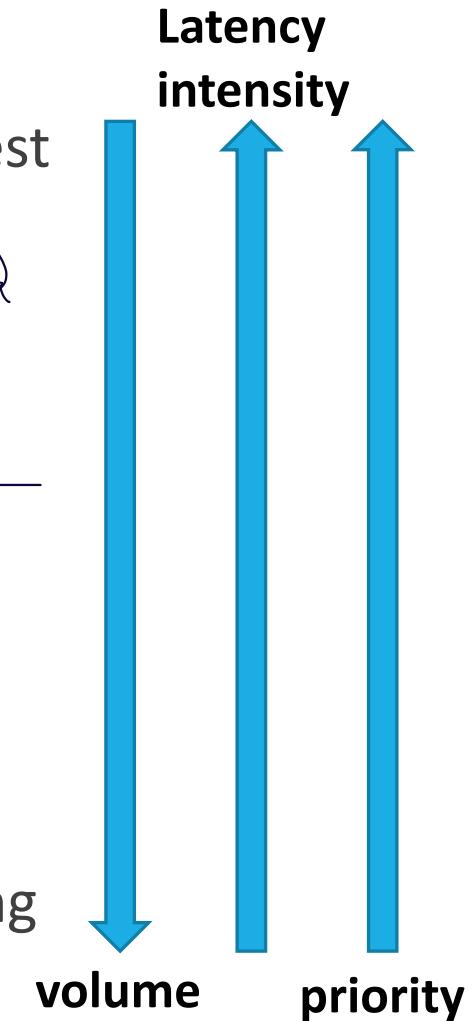
## WAN Intensive Apps

YouTube Web Search  
Google+ Maps AppEngine  
Photos and Hangouts  
Android/Chrome Updates

# B4 – Google's software defined WAN

## Main Traffic flow types

1. **User data copies** to remote data centers for availability/reliability (lowest volume, most latency intensive, highest priority)  
*DC → DC  
COMPUTATION TASK → MOVE DATA THAT ARE DISTRIBUTED*
2. **Remote storage access** for computation over distributed data sources  
*ACCESS TO DISTRIBUTED DATA TO BE PROCESSED*
3. Large-Scale data push **synchronizing state across multiple data centers**  
(highest volume, least latency intensive, lowest priority)



Cost is a major issue (high growth rate of WAN traffic)

- **WAN links are typically utilized 30-40% 2-3 x bandwidth overprovisioning**

# Traditional WAN routing

---

Treat all bits the same



30% ~ 40% average utilization



Cost of bandwidth, High-end routing gear

→ inter DC traffic

## B4 requirements

APPS HAVE  
BUT elastic

(1)

- **Elastic bandwidth demand**: no interactive apps, sync, update replicas...
- Applications can tolerate periodic failures or temporary bandwidth reductions

(2)

- **Moderate number of sites** to connect
- Few dozen of sites, no need for large routing tables

(3)

- **End application control** by Google → can schedule tasks for traffic
- Fine-grained application control (no need for link overprovisioning) (no all load at same time)

(4)

- **Cost sensitivity**
- Private intercontinental links cost, so it needs to be used at its full capacity (100%)

# Main design decisions

1)

**separate HW from SW** : simple / low cost HW and put intelligence in SW

- Customize routing and monitoring protocols to requirements
- Rapid custom protocol development
- Leverage powerful compute in Google servers

2)

**low cost switching hardware**

- Edge application control limits need for large buffers
- limited number of sites -> no need of large forwarding tables

3)

**centralize traffic engineering (TE)**

- Allow more optimal and faster TE than distributed control routing
- Share bandwidth among competing applications
- TE server is an application on top of an SDN distributed control layer

4)

**increase link utilization**

- Efficient use of expensive long haul transport
- Largest bandwidth consumers adapt to available bandwidth
- Try to share in more efficient way bandwidth between apps

we just have place f(x)s and small enough destination

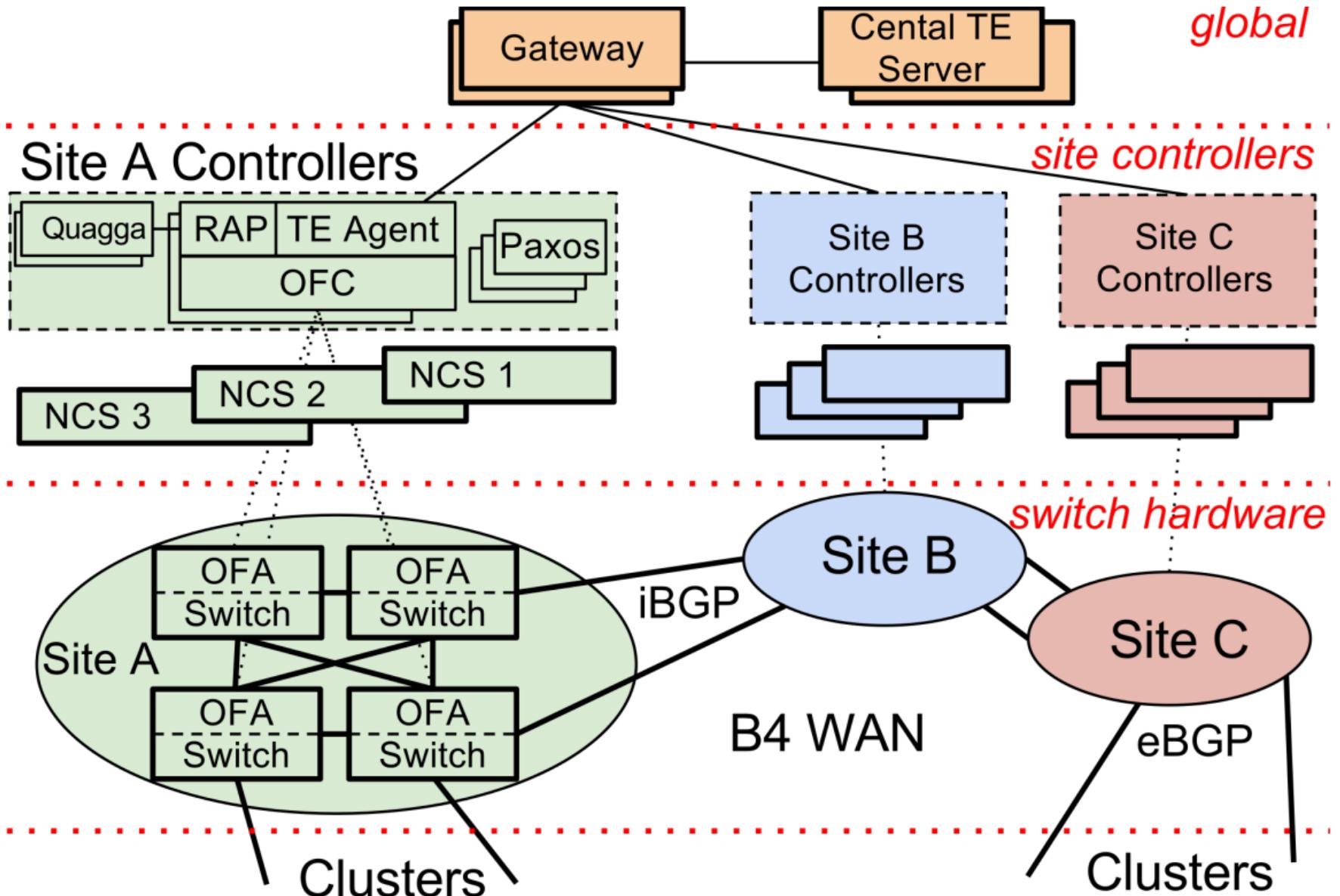
(SITES)  
=> slow buffer  
(NO LAG)  
FORWARDING TABLES

edge control in a net

**Traffic Engineering**: techniques to optimize and control traffic flows in a telecommunication network in order to ensure a maximum throughput and a sufficient QoS level for your flows

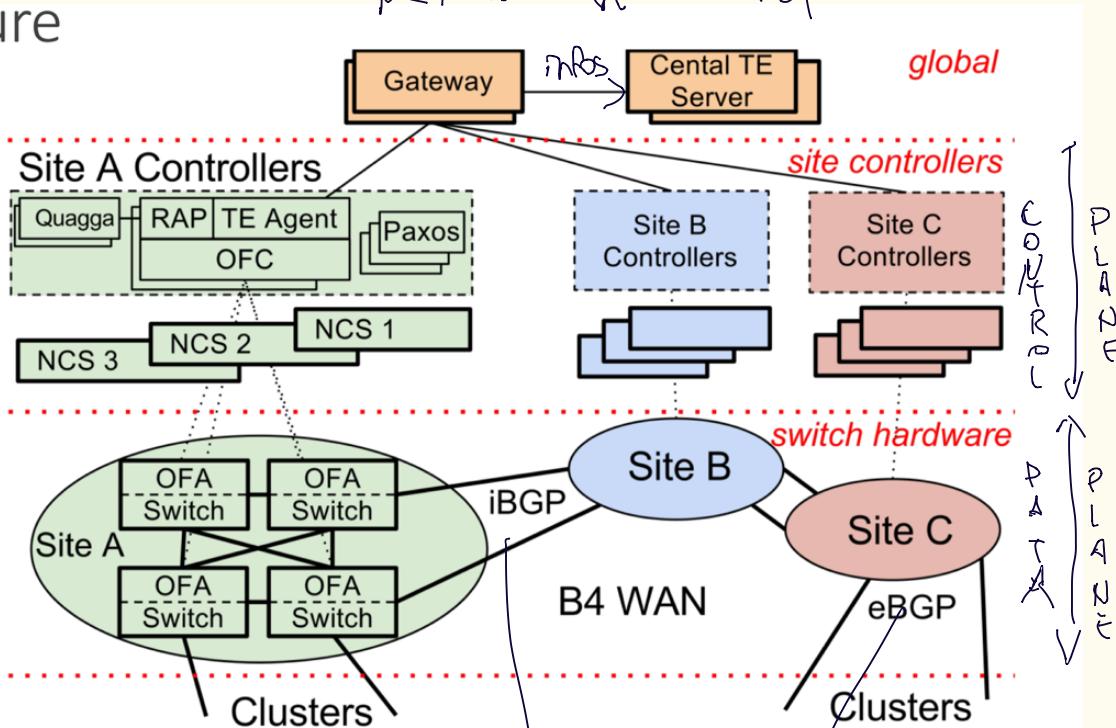
can be + flow

# B4 architecture overview



ture

Network-app. control



NCS: network control servers, host controllers

internal / external BGP session

=> DIFFERENT Flow: routes provided by controller & by BGP (local priority)

• SITE A, B, C = Datacenters

• SET of OpenFlow switch in each site

• OFA: open flow agent that communicates with the controller

• each SITE has a SITE CONTROLLER

and at site controller:

multiples SDN switches  
(replication, performance,...)  
reliability...

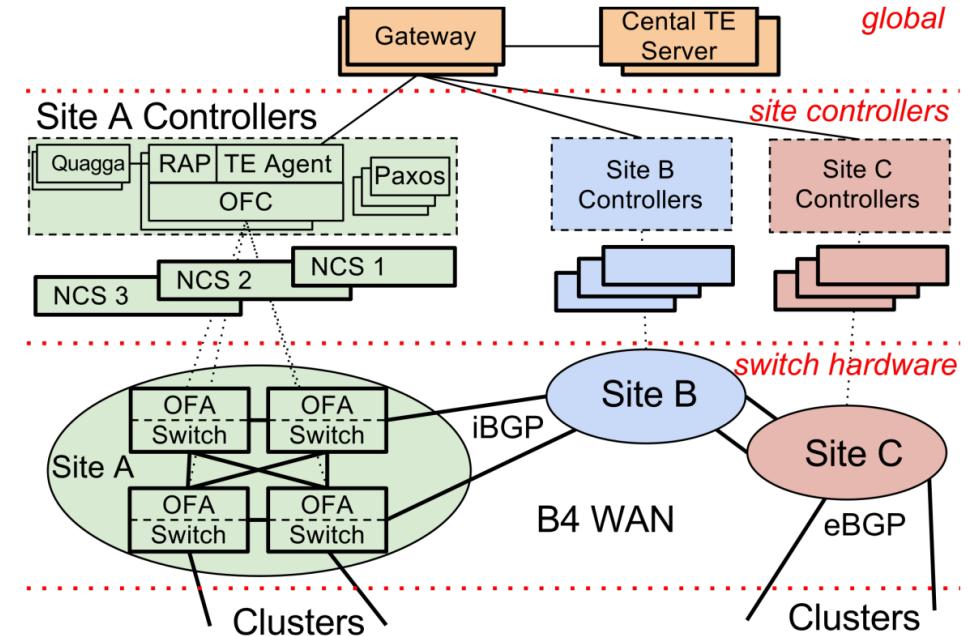
# B4 Architecture overview

## Switch hardware

- Forwards traffic (skewer by controller)
- No complex control software (low cost switch hw)
- Google custom designed with commodity silicon

## Control layer

- Network Control Servers (NCS) hosting both OpenFlow controllers (OFC) and Network Control Applications (NCAs).
- OFCs maintain network state based on NCA directives and switch events
- OFCs instruct switches to set forwarding table entries based on this changing network state



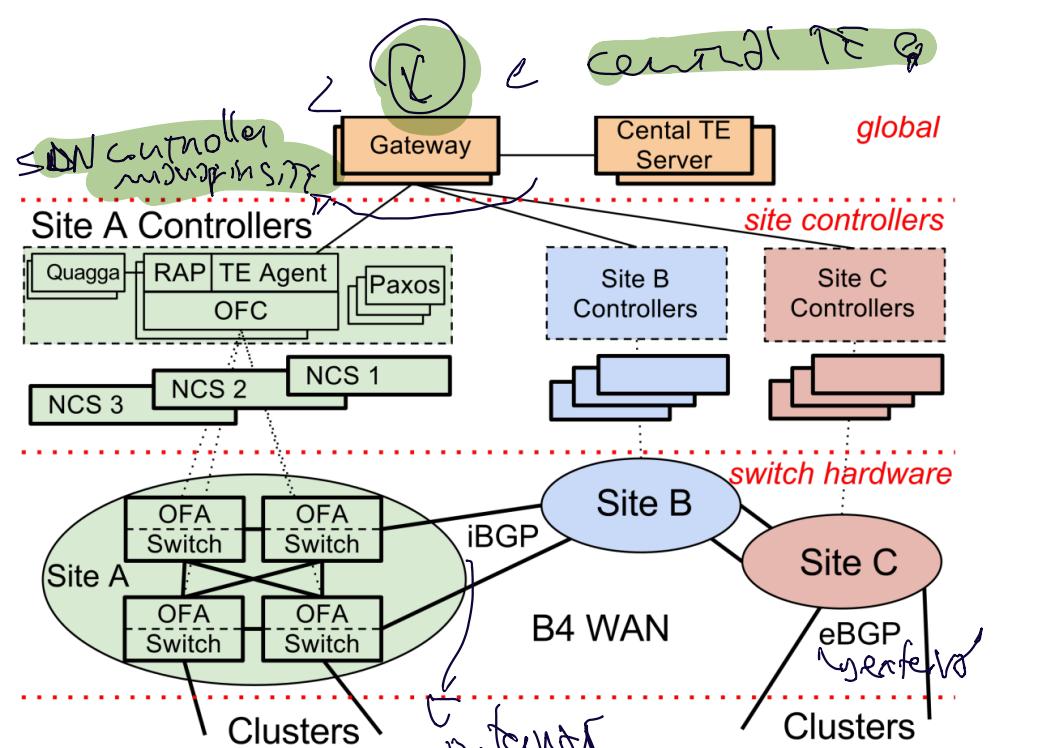
control TC  
Send  
directive  
To controller  
Translating entry  
in table + site

# B4 Architecture overview

## Global layer

- Logically centralized applications: SDN Gateway and a central TE server)
- SDN Gateway abstracts details of OpenFlow and switch hardware from the central TE server.
- Each B4 site consists of multiple switches with potentially hundreds of individual ports

• GATEWAY is an intermediator



## Traditional WAN integrated with SDN

Each cluster contains a set of BGP routers that peer with B4 switches at each WAN sites

each switch  
will have  
its own OpenFlow

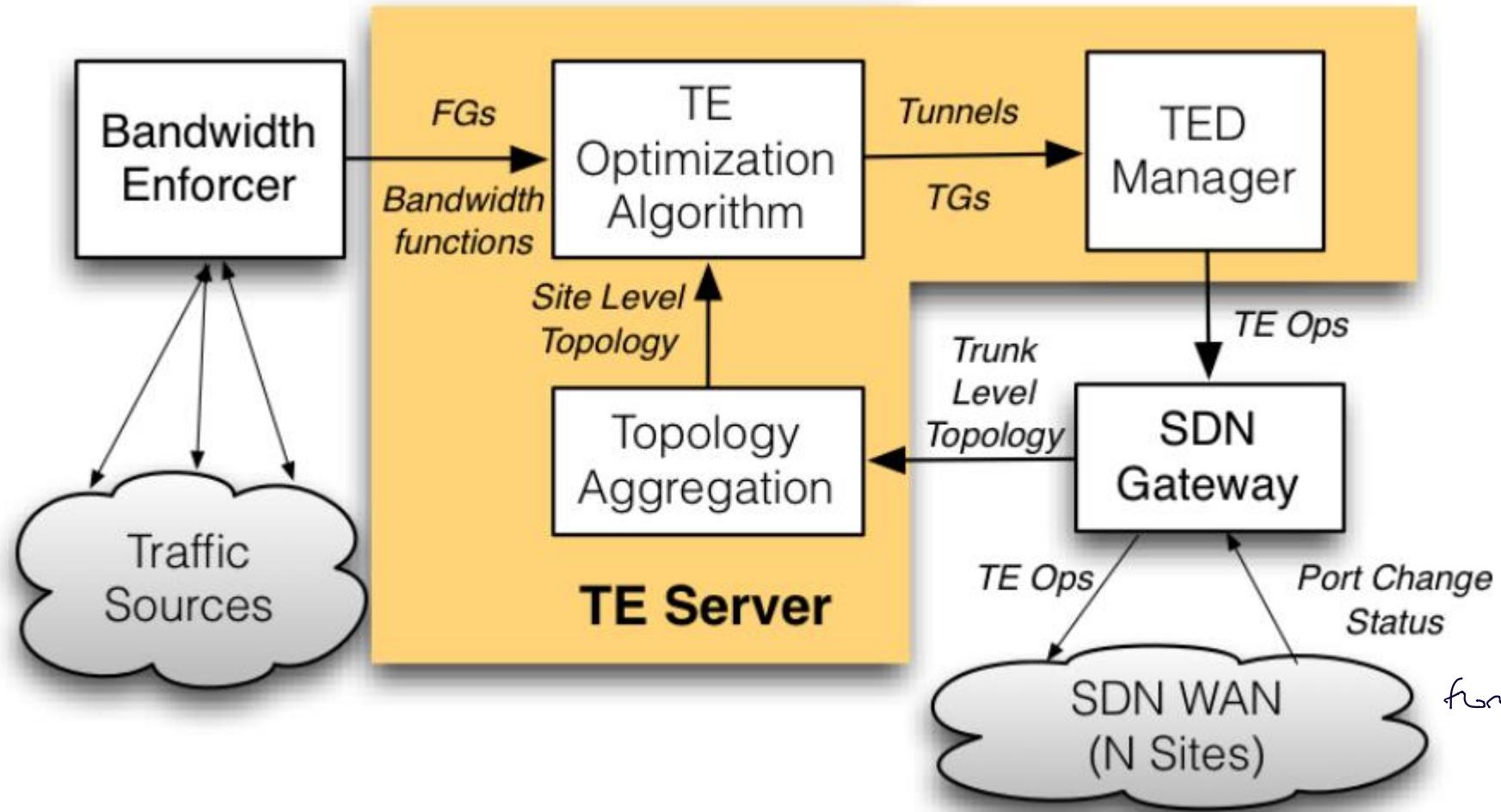
## Why BGP?

- operator familiarity with the protocol (BGP) → didn't disrupt network when introducing SDN
- used before introducing SDN
- enabling a gradual rollout

was  
WT site has routers running BGP and exchange BGP session  
MAO via BGP protocol

# B4 – Traffic Engineering

Switch from BGP to SDN  
Migr. " to SDN

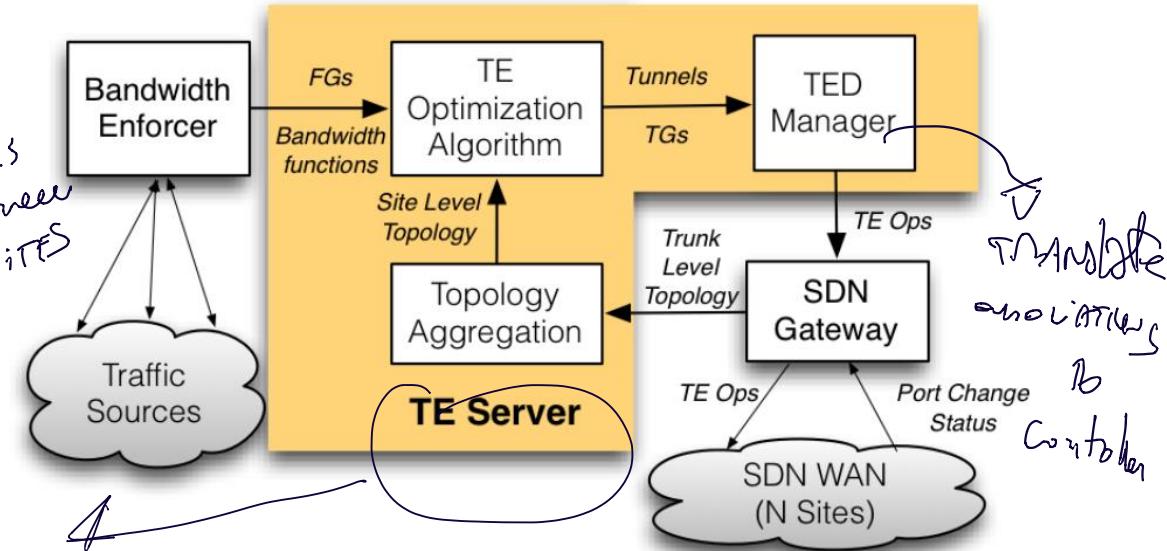


# B4 – Traffic Engineering

## TE Server

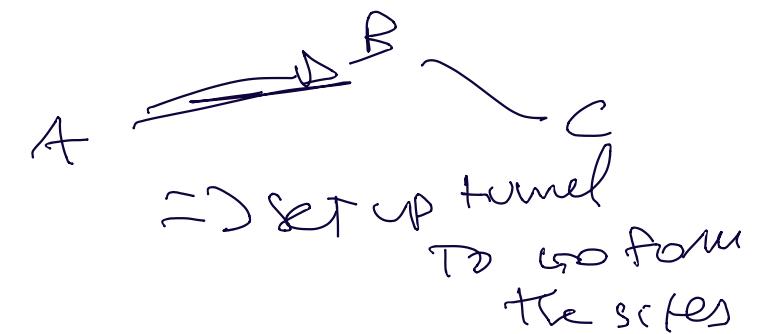
- aggregates trunks to compute site-site edges.
- This abstraction reduces the size of the topology graph in input to the TE Optimization Algorithm
- Applications are aggregated as a flow group (FG) defined as {source site, dest site, QoS} tuple
- A Tunnel (T) represents a site-level path in the network, e.g., a sequence of sites
- Tunnel Group (TG) maps FGs to a set of tunnels

IN POC JUST edge of links between sites  
SITES  
~~2 (no each port of sv. mes)~~ unites between sites



TE Boxes

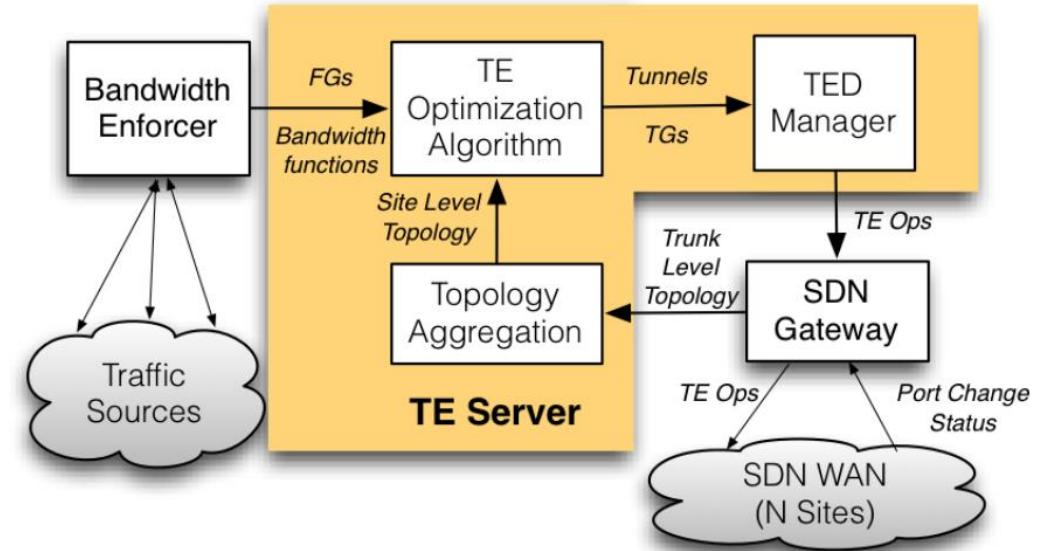
The std an EDGE



# B4 – Traffic Engineering

## Bandwidth Enforcer

- Specified the bandwidth allocation to a specific application (Bandwidth functions)
- These functions are derived from administrator-specified static weights



TE Server outputs the Tunnel Groups and, by reference, Tunnels and Flow Groups to the SDN Gateway.

The Gateway forwards these Tunnels and Flow Groups to OFCs that in turn install them in switches using OpenFlow

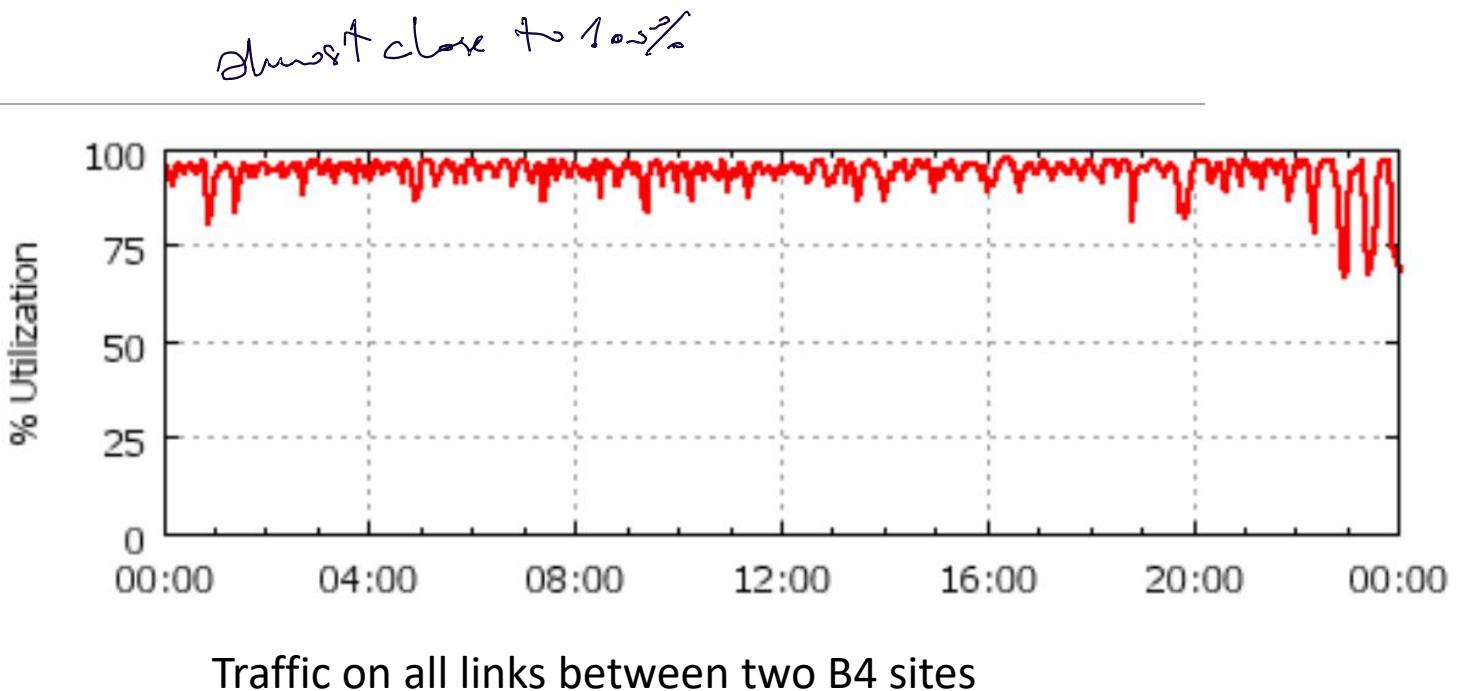
# At the end...

## OpenFlow has helped Google improving B4 backbone<sup>a</sup>

- Utilization near 100% for elastic loads (Sudarshan)
- Mirror production events for testing on virtualized switches
- Reduced complexity
- Reduced costs

## Follow-up Improvements<sup>b</sup>

- hierarchical topology
- decentralized TE algorithms



### References

- Sushant Jain et al. "B4: experience with a globally-deployed software defined wan". In Proc. of the ACM SIGCOMM 2013 ACM, New York, NY, USA, 3-14
- Hong CYY et al.,.. B4 and after: managing hierarchy, partitioning, and asymmetry for availability and scale in google's software-defined WAN. In Proc. Proc. of the ACM SIGCOMM 2018 Aug 7 (pp. 74-87).

# Additional use cases

---

Transport network and enterprise network services

- <https://www.infinera.com/control-automation/>
- <https://www.telstra.com.au/business-enterprise/products/networks/sdn/telstra-programmable-network>
- <https://www.blueplanet.com/resources/analysys-mason-telefonica-germany-is-partnering-with-blue-planet-to-execute-its-ifusion-transport-sdn-strategy.html>
- <https://www.timenterprise.it/connectivity/rete-fissa/tim-sdwan>

# References

---

- a. Sushant Jain et al. “B4: experience with a globally-deployed software defined wan”. In Proc. of the ACM SIGCOMM 2013 ACM, New York, NY, USA, 3-14
- b. Hong CYY et al.,. B4 and after: managing hierarchy, partitioning, and asymmetry for availability and scale in google's software-defined WAN. In Proc. Proc. of the ACM SIGCOMM 2018 Aug 7 (pp. 74-87).
- c. Lessons Learned from B4, Google's SDN WAN, Subhasree Mandal, July 9, 2015

# Switches with OpenFlow support - examples

---

Cisco Catalyst 9300 series

<https://www.cisco.com/c/en/us/support/docs/switches/catalyst-9300-series-switches/217210-understand-openflow-on-catalyst-9000-ser.html>

Noviflow

<https://noviflow.com/noviswitch/>

Juniper

<https://www.juniper.net/documentation/us/en/software/junos/sdn-openflow/topics/concept/junos-sdn-openflow-supported-platforms.html>