# Jupiter rising: A decade of Clos topologies and centralized control in Google's datacenter networks

**Phil Sykes**
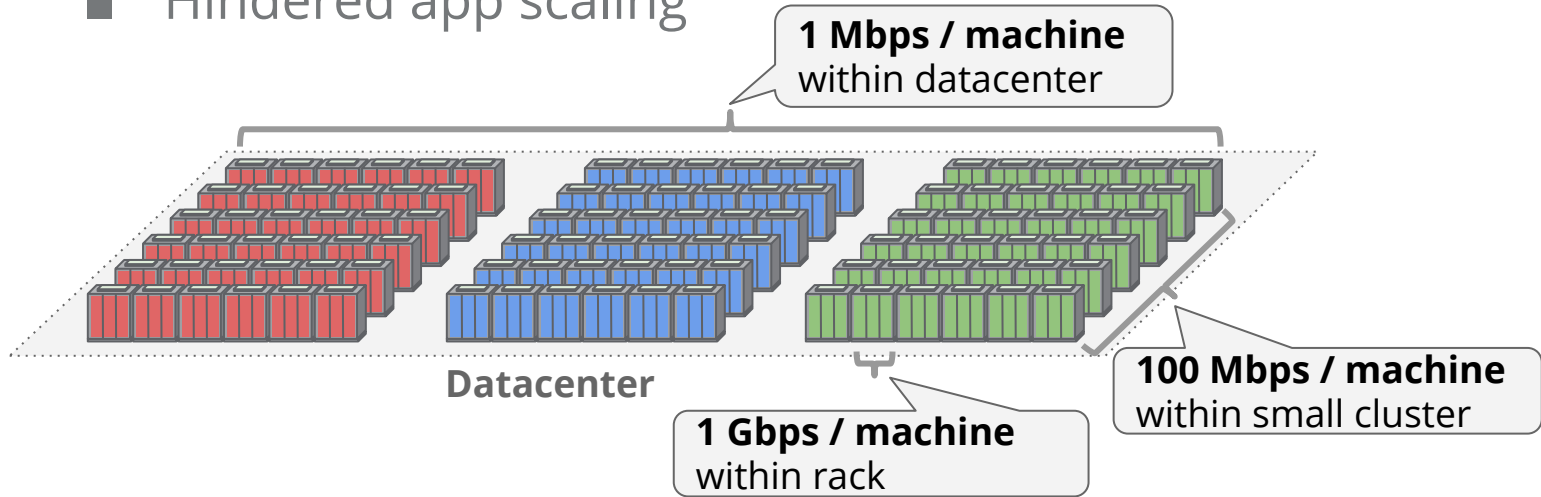**Google NetOps (Dublin)**
**philsykes@google.com**

Arjun Singh, Joon Ong, Amit Agarwal, Glen Anderson, Ashby Armistead, Roy Bannon, Seb Boving, Gaurav Desai, Bob Felderman, Paulie Germano, Anand Kanagala, Jeff Provost,  Jason Simmons, Eiichi Tanda, Jim Wanderer, Urs Hölzle, Stephen Stuart, and Amin Vahdat

On behalf of several teams in Google:

Platforms Networking Hardware and Software Development, Platforms SQA, Mechanical Engineering, Cluster Engineering, NetOps, Global Infrastructure Group (GIG), and SRE.
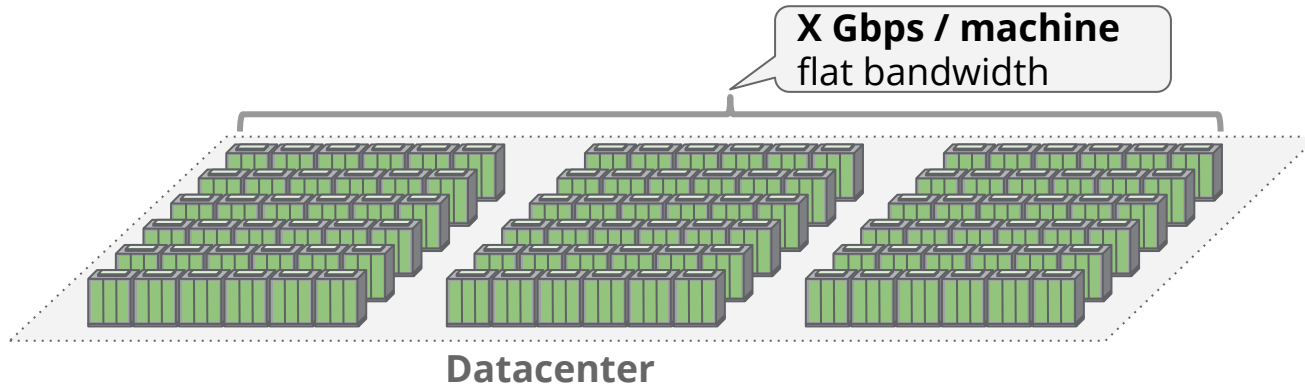
Google

# Grand challenge for datacenter networks

- Tens of thousands of servers interconnected in clusters
- *Islands of bandwidth* a key bottleneck for Google a decade ago
  - Engineers struggled to optimize for b/w locality
  - Stranded compute/memory resources
  - Hindered app scaling

**1 Mbps / machine**
within datacenter

**Datacenter**

**1 Gbps / machine**
within rack

**100 Mbps / machine**
within small cluster

3

Google

# Grand challenge for datacenter networks

- **Challenge: Flat b/w profile across all servers**
  - Simplify job scheduling (remove locality)
  - Save significant resources via better bin-packing
  - Allow application scaling

**X Gbps / machine**
flat bandwidth

**Datacenter**

Google

# Motivation

- **Traditional network architectures**

  - Cost prohibitive

  - Could not keep up with our bandwidth demands

  - Operational complexity of "box-centric" deployment

- **Opportunity: A datacenter is a single administrative domain**

  - One organization designs, deploys, controls, operates the n/w

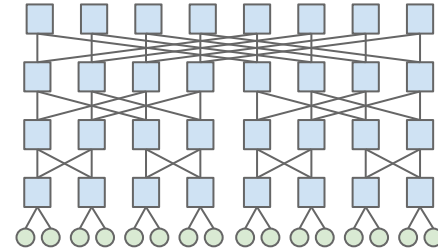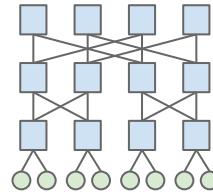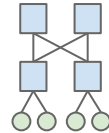  - ...And often also the servers

Google

# Three pillars that guided us

**Merchant silicon:** General purpose, commodity priced, off the shelf switching components

**Clos topologies:** Accommodate low radix switch chips to scale nearly arbitrarily by adding stages

**Centralized control / management**

# SDN: The early days

- Control options
  - Protocols: OSPF, ISIS, BGP, etc; *Box-centric* config/management
  - Build our own

- Reasons we chose to build our own central control/management:
  - Limited support for **multipath forwarding**
  - No **robust** open source stacks
  - **Broadcast** protocol scalability a concern at scale
  - Network **manageability** painful with individual switch configs

# Challenges faced in building our own solution

- **Topology and deployment**
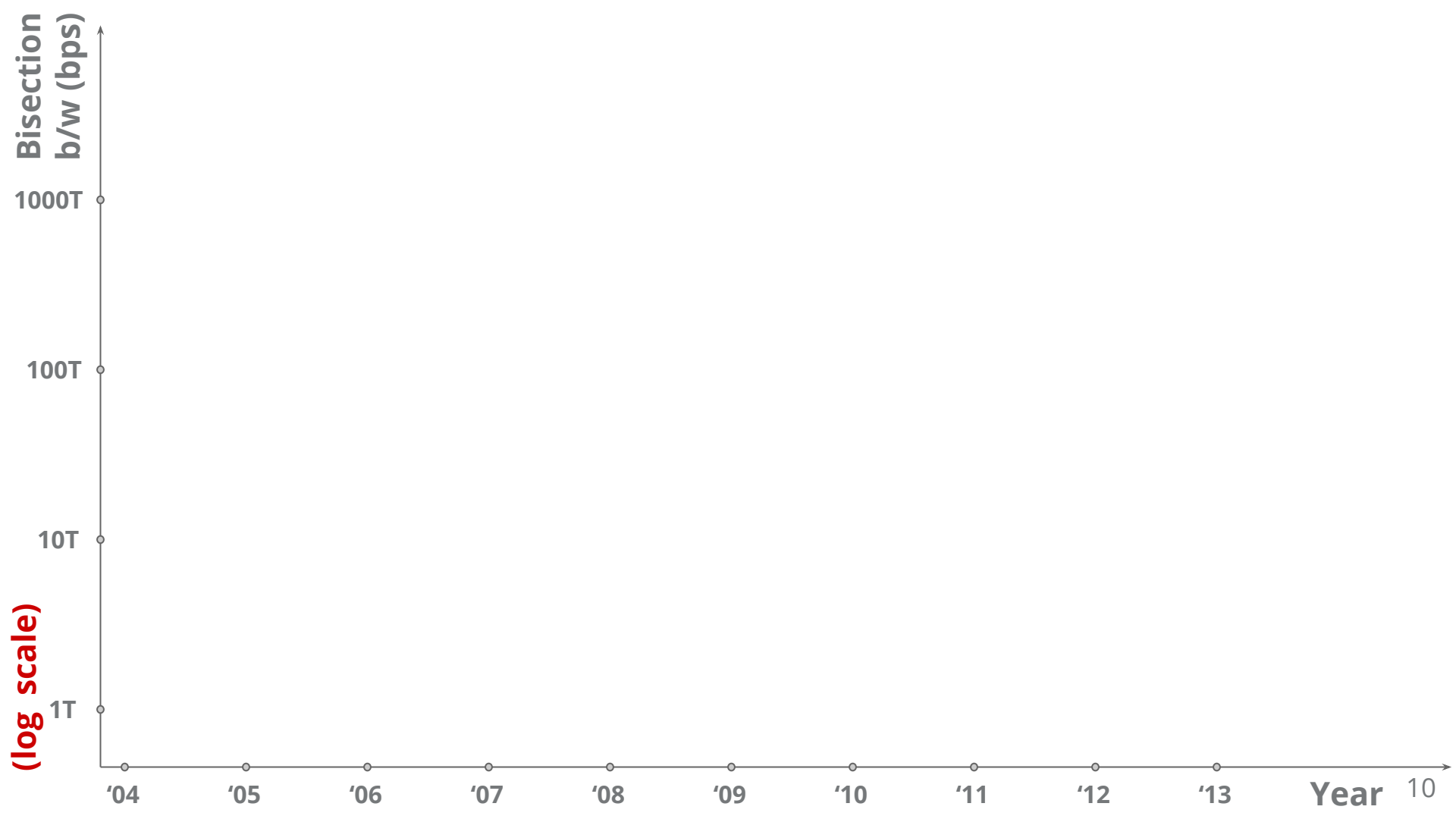  - Introducing our network to production
  - Unmanageably high number of cables/fiber
  - Cluster-external burst b/w demand
- **Control and management**
  - Operating at huge scale
  - Routing scalability / routing with massive multipath
  - Interop with external vendor gear
- **Performance and reliability**
  - Small on-chip buffers
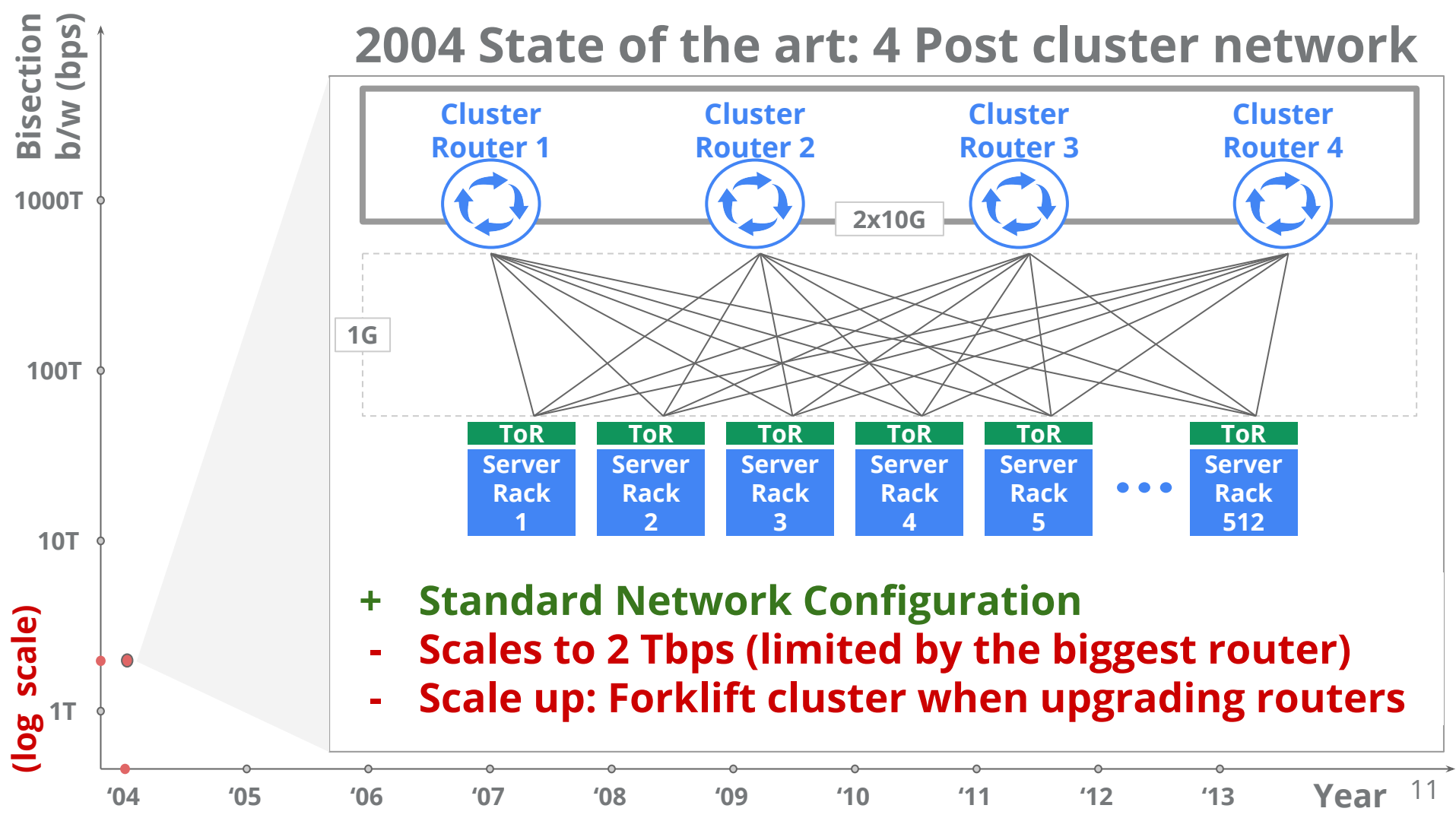  - High availability from cheap/less reliable components

Google

# Outline

- Motivation
- **Network evolution**
- Centralized control / management
- Experience

# 2004 State of the art: 4 Post cluster network

**Cluster Router 1**  **Cluster Router 2**  **Cluster Router 3**  **Cluster Router 4**

2x10G

1G

ToR — Server Rack 1
ToR — Server Rack 2
ToR — Server Rack 3
ToR — Server Rack 4
ToR — Server Rack 5
• • •
ToR — Server Rack 512

**+ Standard Network Configuration**

**- Scales to 2 Tbps (limited by the biggest router)**

**- Scale up: Forklift cluster when upgrading routers**

Bisection b/w (bps) (log scale)

1000T
100T
10T
1T

'04  '05  '06  '07  '08  '09  '10  '11  '12  '13  **Year**

# DCN bandwidth growth demanded much more



50x — Traffic generated by servers in our datacenters

Aggregate traffic

1x

Time

Jul '08    Jun '09    May '10    Apr '11    Mar '12    Feb '13    Dec '13    Nov '14

Google

# Five generations of Clos for Google scale

**Firehose 1.1**

Bisection b/w (bps)

1000T

100T

**Firehose 1.0**

10T

**4 Post**

(log scale)

1T

**+ Chassis based solution (but no backplane)**
**- Bulky CX4 copper cables restrict scale**

'04  '05  '06  '07  '08  '09  '10  '11  '12  '13  **Year**

14

# Challenges faced in building our own solution

- **Topology and deployment**
  - **Introducing our network to production**
  - Unmanageably high number of cables/fiber
  - Cluster-external burst b/w demand
- **Control and management**
  - Operating at huge scale
  - Routing scalability / routing with massive multipath
  - Interop with external vendor gear
- **Performance and reliability**
  - Small on-chip buffers
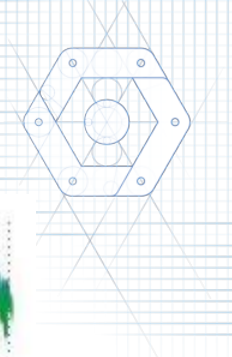  - High availability from cheap/less reliable components

# Firehose 1.1

**Four-post cluster routers**

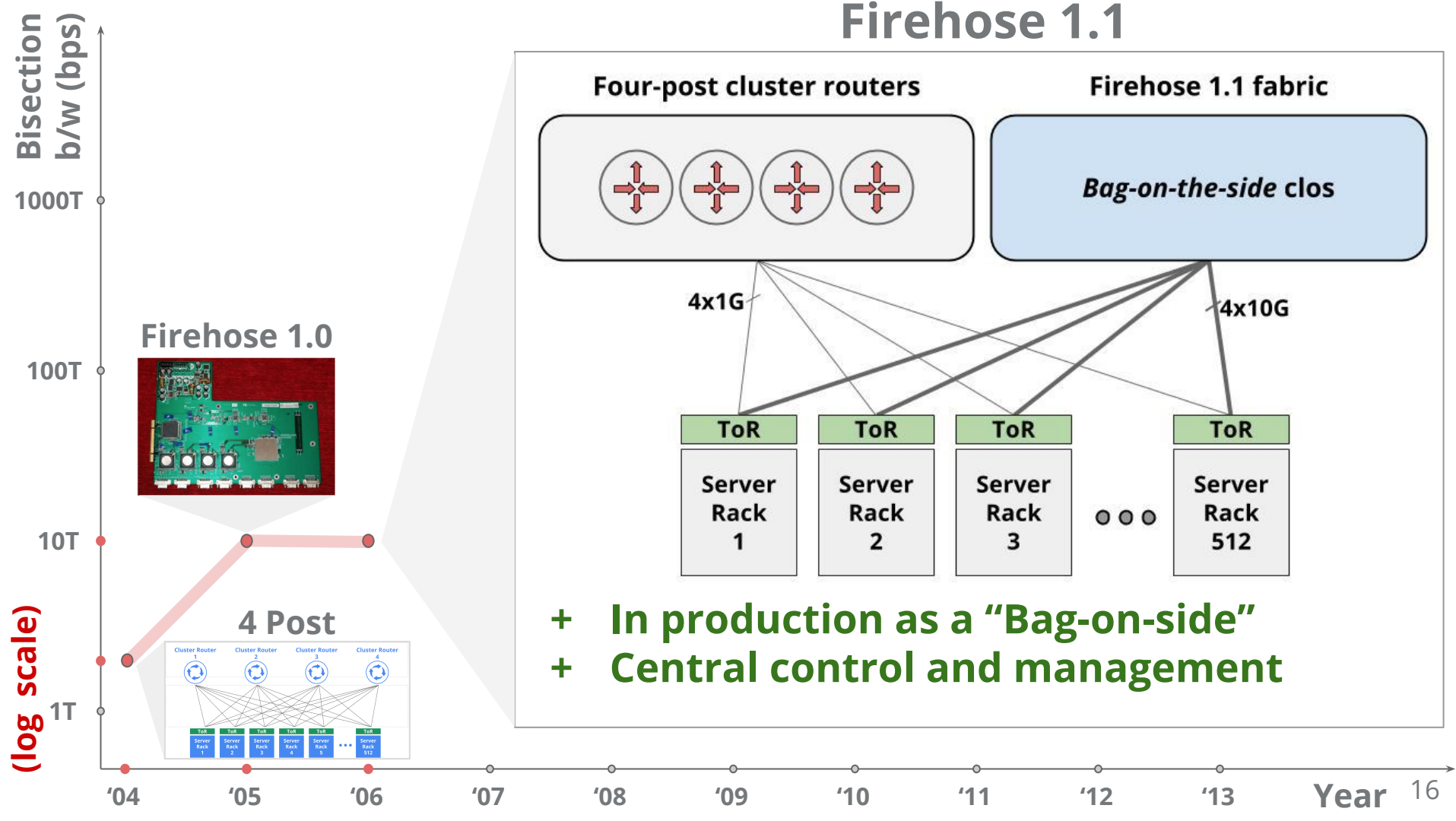**Firehose 1.1 fabric**

*Bag-on-the-side* clos

4x1G

4x10G

ToR · ToR · ToR · ToR

Server Rack 1 · Server Rack 2 · Server Rack 3 · · · Server Rack 512

**Firehose 1.0**

**4 Post**

Cluster Router 1 · Cluster Router 2 · Cluster Router 3 · Cluster Router 4

ToR · ToR · ToR · ToR · ToR · · · ToR

Server Rack 1 · Server Rack 2 · Server Rack 3 · Server Rack 4 · Server Rack 5 · Server Rack 512

+ **In production as a "Bag-on-side"**
+ **Central control and management**

Bisection b/w (bps)

1000T

100T

10T

1T

**(log scale)**

'04  '05  '06  '07  '08  '09  '10  '11  '12  '13   **Year**

16

Bisection b/w (bps)

(log scale)

1000T

100T

**Firehose 1.0**

10T

**Firehose 1.1**

**4 Post**

1T

'04  '05  '06  '07  '08  '09  '10  '11  '12  '13   **Year**

**Watchtower**

Bisection b/w (bps)

1000T

100T — Firehose 1.0

10T — 4 Post

1T

**(log scale)**

+ **Chassis with backplane**
+ **Fiber (10G) in all stages**
+ **Scale to 82 Tbps fabric**
+ **Global deployment**

Firehose 1.1

'04  '05  '06  '07  '08  '09  '10  '11  '12  '13  **Year**  18

**Bisection b/w (bps)**

1000T

100T

10T

**(log scale)**

1T

**Saturn**

**Watchtower**

**Firehose 1.0**

**Firehose 1.1**

**4 Post**

Cluster Router 1  Cluster Router 2  Cluster Router 3  Cluster Router 4

ToR  ToR  ToR  ToR  ToR  ToR
Server Rack 1  Server Rack 2  Server Rack 3  Server Rack 4  Server Rack 5  · · ·  Server Rack 512

+  **288x10G port chassis**
+  **Enables 10G to hosts**
+  **Scales to 207 Tbps fabric**
+  **Reuse in WAN (B4)**

'04   '05   '06   '07   '08   '09   '10   '11   '12   '13   **Year**

Bisection b/w (bps)

**(log scale)**

1000T

100T

10T

1T

**Watchtower**

**Jupiter**

**Saturn**

**Firehose 1.0**

**Firehose 1.1**

**4 Post**

'04  '05  '06  '07  '08  '09  '10  '11  '12  '13  **Year** 20

**+ Scales out building wide 1.3 Pbps**

# Jupiter racks



+ **Enables 40G to hosts**
+ **External control servers**
+ **OpenFlow**

Google

**Bisection b/w (bps)**

**(log scale)**

1000T

100T

10T

1T

**Watchtower**

**Jupiter (1.3P)**

**Firehose 1.0**

**Saturn**

**Firehose 1.1**

**4 Post**

'04   '05   '06   '07   '08   '09   '10   '11   '12   '13   **Year**

23

# Challenges faced in building our own solution

- **Topology and deployment**
  - Introducing our network to production
  - Unmanageably high number of cables/fiber
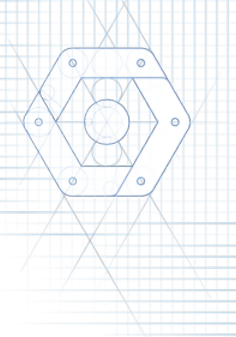  - Cluster-external burst b/w demand
- **Control and management**
  - **Operating at huge scale**
  - **Routing scalability / routing with massive multipath**
  - Interop with external vendor gear
- **Performance and reliability**
  - Small on-chip buffers
  - High availability from cheap/less reliable components

# Network control and config

New conventional wisdom from engineering systems at scale

- **Logically centralized control** plane beats full decentralization

- **Centralized configuration and management** dramatically simplifies system aspects

- **Scale out** >> Scale up

# Challenges faced in building our own solution

- **Topology and deployment**
  - Introducing our network to production
  - Unmanageably high number of cables/fiber
  - Cluster-external burst b/w demand
- **Control and management**
  - Operating at huge scale
  - Routing scalability / routing with massive multipath
  - Interop with external vendor gear
- **Performance and reliability**
  - **Small on-chip buffers** ⟹ **Tune switches (eg ECN) and Hosts (DCTCP)**
  - High availability from cheap/less reliable components

# Challenges faced in building our own solution

- **Topology and deployment**
  - Introducing our network to production
  - Unmanageably high number of cables/fiber
  - Cluster-external burst b/w demand
- **Control and management**
  - Operating at huge scale
  - Routing scalability / routing with massive multipath
  - Interop with external vendor gear
- **Performance and reliability**
  - Small on-chip buffers
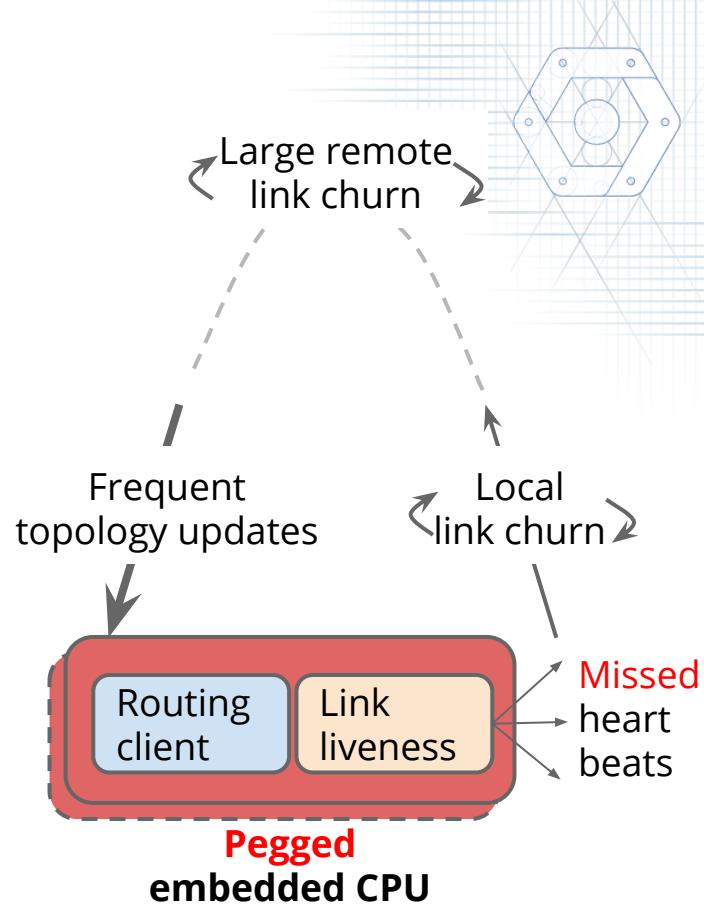  - **High availability from cheap/less reliable components**

    ↳ **Redundancy; diversity; implement only what was needed**

# Experience: Outages

Three broad categories of outages:

- **Control software failures at scale**
  - Cluster-wide reboot did not converge
    - Liveness protocol contended for cpu with routing process
  - Cannot test at scale in a hardware lab
    - Developed virtualized testbeds
- **Aging hardware exposes corner cases**
- **Component misconfigurations**



Large remote link churn

Frequent topology updates

Local link churn

Routing client

Link liveness

Missed heart beats

**Pegged embedded CPU**

# Grand challenge for datacenter networks

- **Challenge:  Flat b/w profile across all servers**
  - Simplify job scheduling (remove locality)
  - Save significant resources (better bin-packing)
  - Allow application scaling
- Scaled datacenter networks to Petabit scale in under a decade
- Bonus: reused solution in campus aggregation and [WAN](WAN)

**X Gbps / machine**
flat bandwidth

**Datacenter**

Google