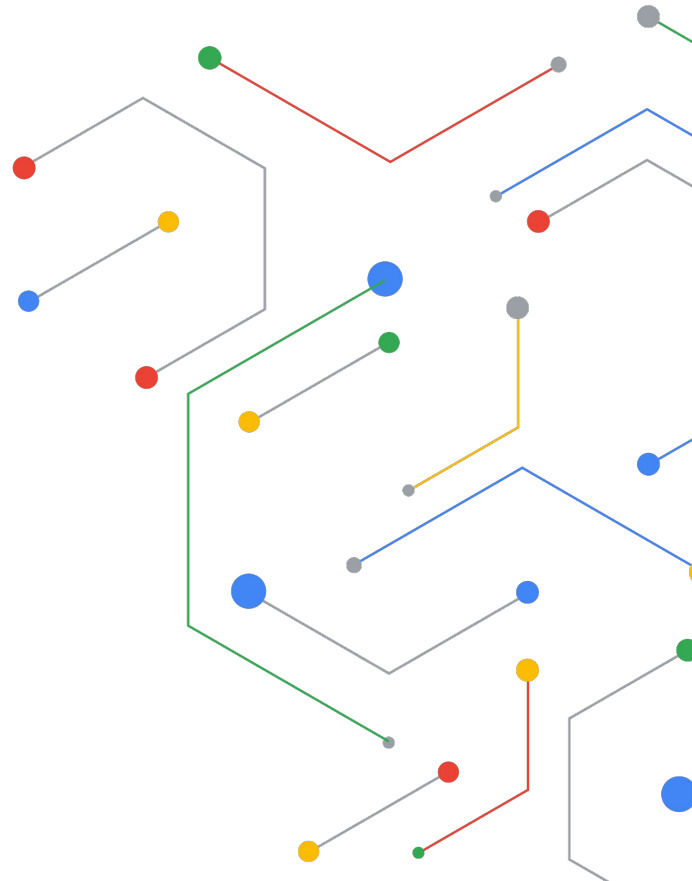*"Without networking, there is no cloud."*

# AI/ML Cloud Networking Innovation

Shaowen Ma, shaowen@google.com
Group Product Manager
Dec 2023

# 7DEOH RI &RQWHQWV

- **Cloud Network Evolution** — 01
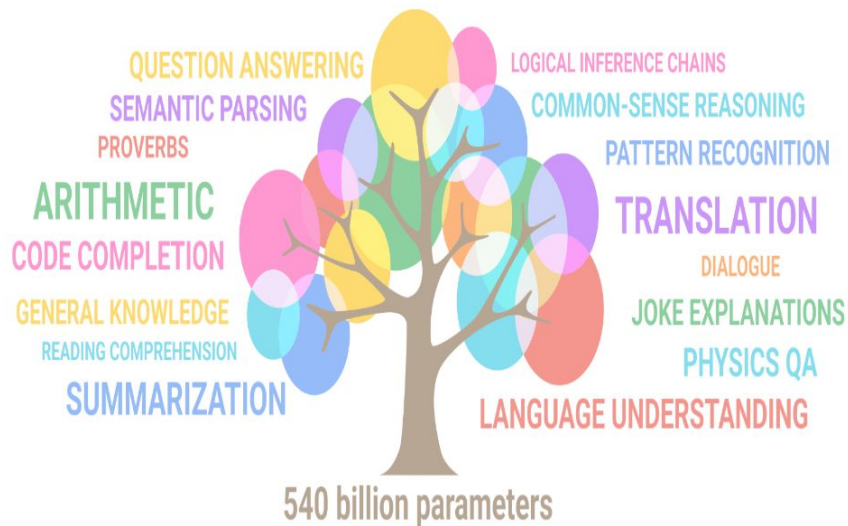- **AI Fabric Innovation** — 02
- **AI Smart Offload Transport Innovation** — 03
- **Summary** — 04

Google Cloud

# Cloud Network Evolution

# Exploding cost of powerful AI

**Model & Dataset Size**

**Capabilities**



QUESTION ANSWERING
SEMANTIC PARSING
PROVERBS
ARITHMETIC
CODE COMPLETION
GENERAL KNOWLEDGE
READING COMPREHENSION
SUMMARIZATION

LOGICAL INFERENCE CHAINS
COMMON-SENSE REASONING
PATTERN RECOGNITION
TRANSLATION
DIALOGUE
JOKE EXPLANATIONS
PHYSICS QA
LANGUAGE UNDERSTANDING

540 billion parameters

**Costs**



Training compute (petaFLOP)

- 1 billion
- 100 million
- 10 million
- 1 million
- 100,000
- 10,000

PaLM 2
PaLM (540B)
Chinchilla
AlphaCode
GPT-3 175B (davinci)
AlexaTM 20B
GPT-NeoX-20B
T5-11B
DALL-E
Whisper
T5-3B
NLLB
ALBERT-xxlarge
GPT-2
BERT-Large
GPT
Transformer

Jun 12, 2017          Jan 21, 2020          Jun 4, 2021          Oct 17, 2022
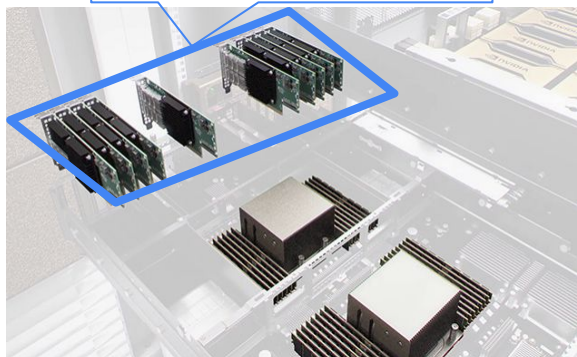
Publication date

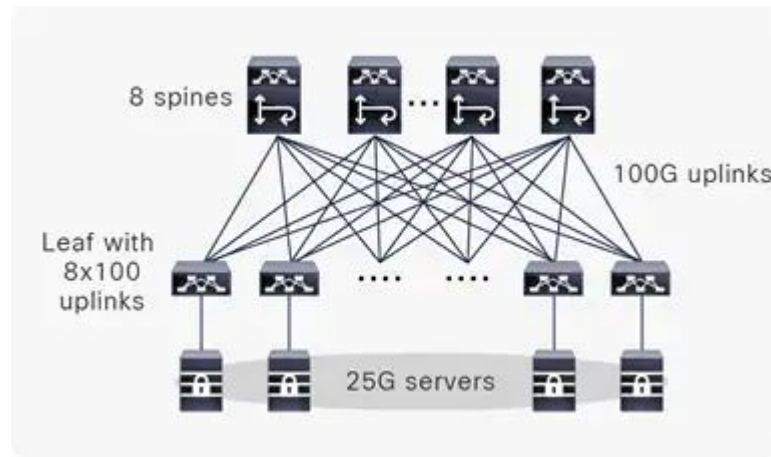)URP FOXVWHUV WR ZDUHKRXVH VFDOH FRPSXWHU

# $, 1HWZRUN 'HVLJQ 'LIIHUHQFH  6HUYHU
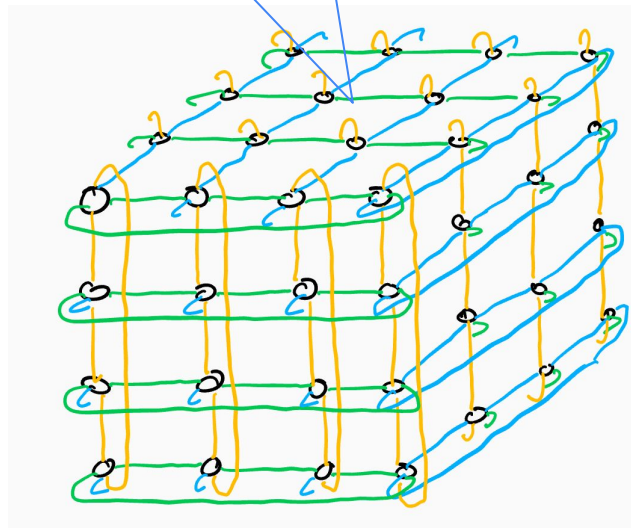
8x400GE= 3.2Tbps+

vs CPU(25/100Gbps)

32~128 times

>>

**Huge bandwidth GPU/TPU**
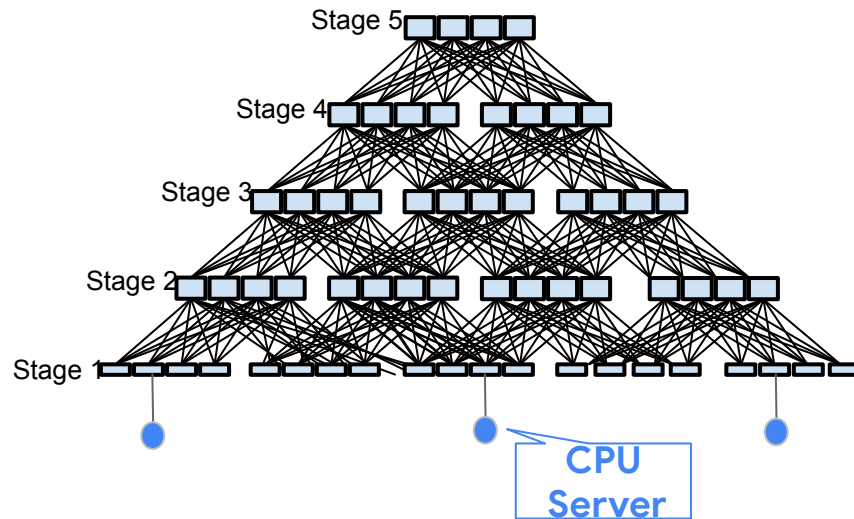1.6Tbps/3.2Tbps

**CPU Server**
still 25Gbps/100Gbps



8 spines

100G uplinks

Leaf with 8x100 uplinks

25G servers

# $, 1HWZRUN 'HVLJQ 'LIIHUHQFH  1HWZRUN IR

TPU Server

Google Jupiter Rising

Stage 5

Stage 4

Stage 3

Stage 2

Stage 1

CPU Server

**3D Torus**
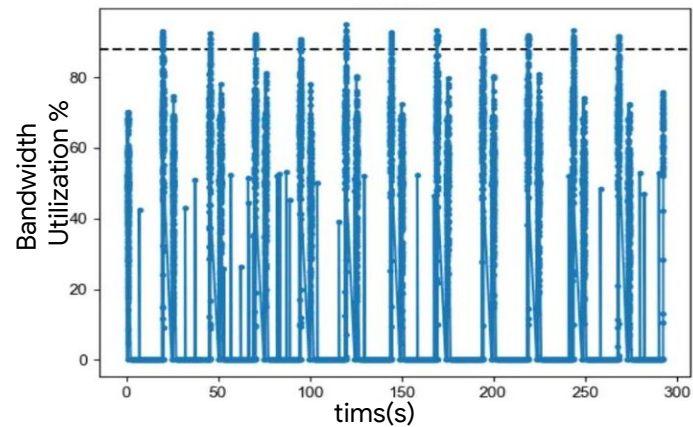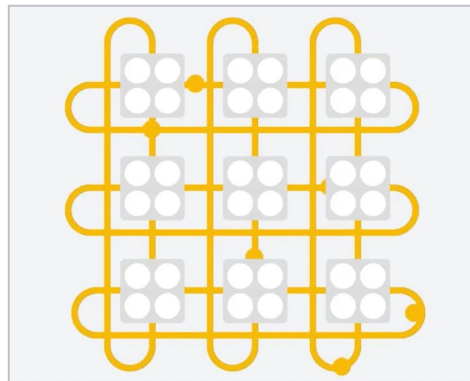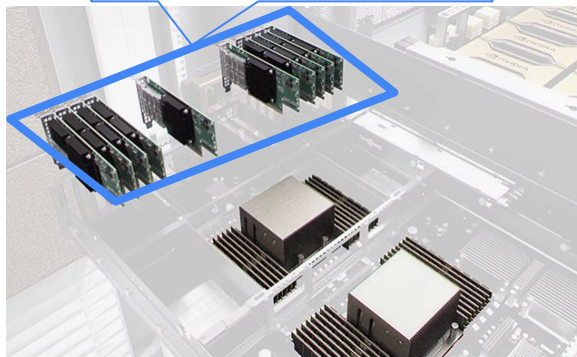**Each node Direct Connect**
**to 6 others**

**CLOS**
**Each node up to 9 hops(switches)**
**to others**

# $, 1HWZRUN 'HVLJQ 'LIIHUHQFH

**8x400GE= 3.2Tbps+**
vs CPU(25/100Gbps)



**Huge bandwidth GPU/TPU**
**(compute still 25G/100G Server)**
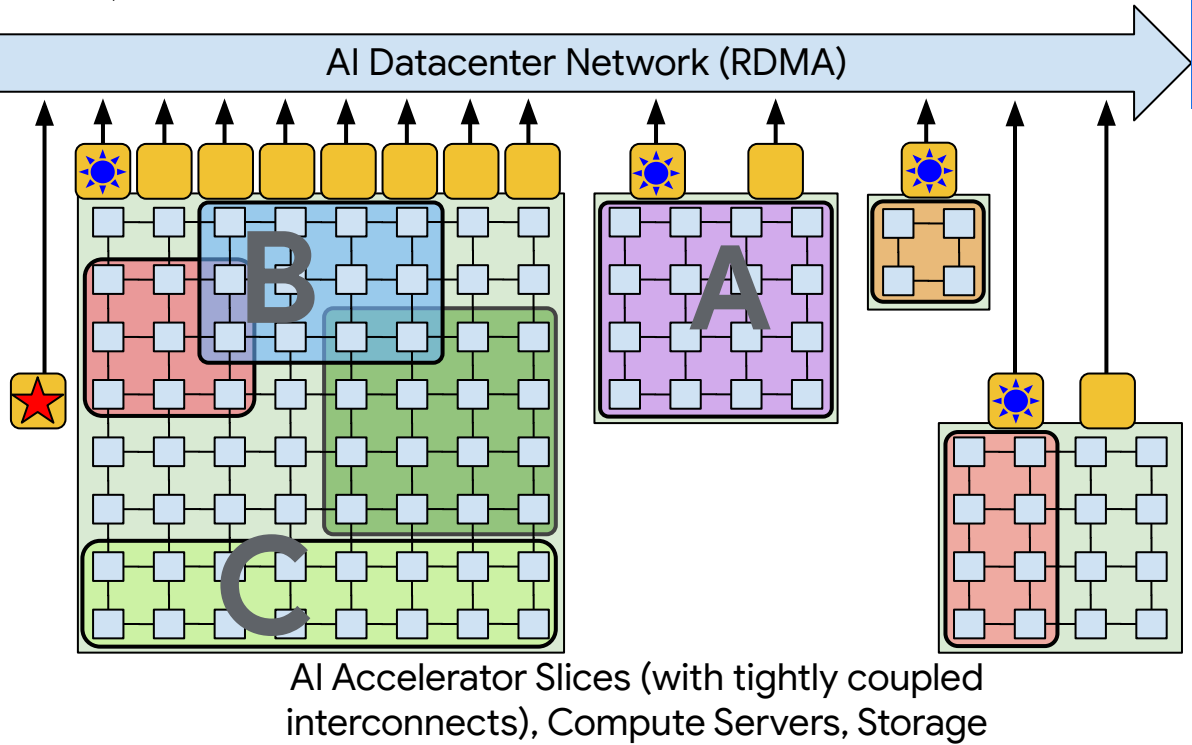


**AllReduce, All2All**
**(compute still CLOS)**



Bandwidth Utilization %

tims(s)

**Burst to 3.2Tbps/100%**

# AI Fabric Innovations

$, 1 H W Z R U N & O X V W H U V 7 R S R O R J \

AI Datacenter Network (RDMA)

B
A
C

AI Accelerator Slices (with tightly coupled interconnects), Compute Servers, Storage

Host (many per island)
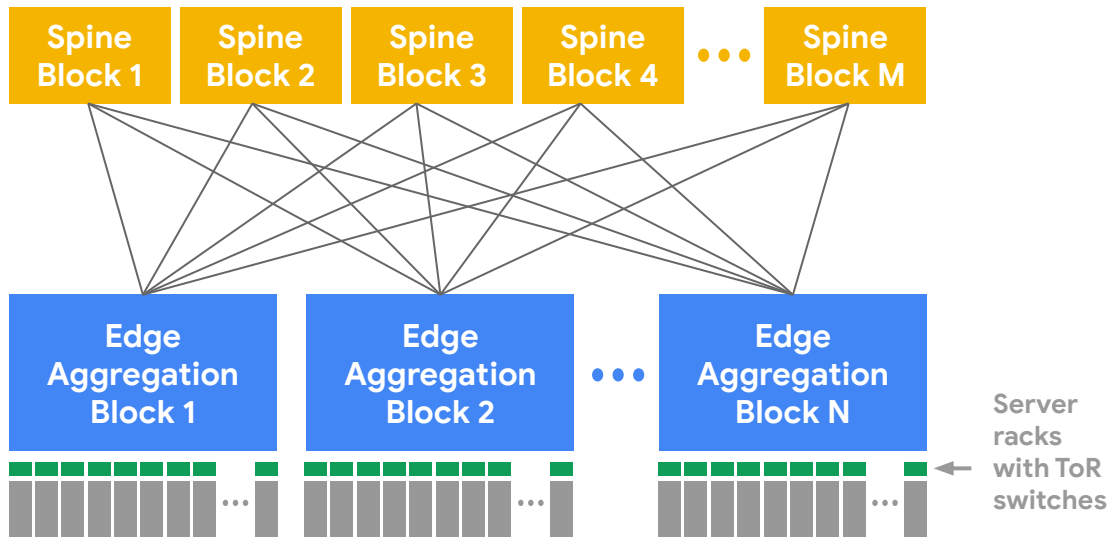Resource Manager (global)
Scheduler (per island)

**WHAT WOULD IT TAKE...**

... to achieve **Performance**, **Isolation** and **Efficiency** at scale for High Bandwidth, Low Latency Workloads on **today's Datacenter Ethernet** Networks?
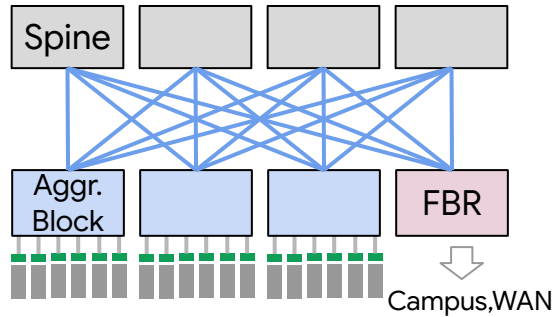
Google

# Five generations of clos topologies and software-defined networking



Spine Block 1 | Spine Block 2 | Spine Block 3 | Spine Block 4 | • • • | Spine Block M

Edge Aggregation Block 1 | Edge Aggregation Block 2 | • • • | Edge Aggregation Block N
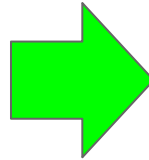
Server racks with ToR switches

A scalable, commodity data center network architecture, SIGCOMM 2008
Jupiter Rising: A Decade of Clos Topologies and Centralized Control in Google's Datacenter Network, SIGCOMM 2015
Orion: Google's Software-Defined Networking Control Plane, USENIX 2021

Google

**Electrical Packet Switch-based DCN**

**OCS-based LW Fabric**



[Jupiter Rising, Sigcomm'15]

[Jupiter Evolving, Sigcomm'22]

- Lightwave Fabric with OCS enables Datacenter Networks with
  - 30% reduction in CapEx and 40% reduction in power consumption
    - No Fiber Change for multiple generations. 100G to 200G to 400G+
    - No Electronic switching in LW Fabric, all Optical Switching.
  - Expansion, topology engineering, heterogeneous networking

Google Cloud

2 S W L F D O   F L U F X L W   V Z L W F K L Q J   L Q   W K H   G D W D



Dynamic functionality (bandwidth on demand)

Apollo OCS

WDM Module
multiple lanes - one fiber
(point-to-point)

Parallel Optics Module
multiple lanes - multiple fibers
(path diversity)

6

6

x8

x32

6

**30%** lower cost
**40%** lower power
**New capability:** application-specific topology!



Spine

100G    100G    100G    200G

100Gb/s links

200Gb/s links

Aggregation Block

100G    100G    100G    200G

ToR
Machine rack

MEMS mirror package die with an array of individually controllable micro mirrors

By actuating MEMS mirrors, the same input port (X) can be connected to a different output port (from Y to Z)



N Input Fibers

X

MEMS Mirror Array

optical core

MEMS Mirror Array

Y
Z

N Output Fibers

OCS switch (logical view)

Google

# 7 U D G L W L R Q D O   Y V   V R I W Z D U H   G H I L Q H G   Q H W Z

## Traditional network

## Software-defined network (SDN)

- O(log N) rounds to converge
- O(N log N) messages
- Global consistency and correctness near impossible

**2** (Hierarchical) Central SDN controller Compute best paths across network.

1000x compute power

SDN protocol

**3** Push forwarding rules to devices.

- O(1) rounds to converge
- O(N) messages
- Transition n/w from one globally consistent state to another

Each router exchanges routing info, independently builds reachability graph, populates switching table, and forwards traffic.

**1** Extract state, capacity, routing info.
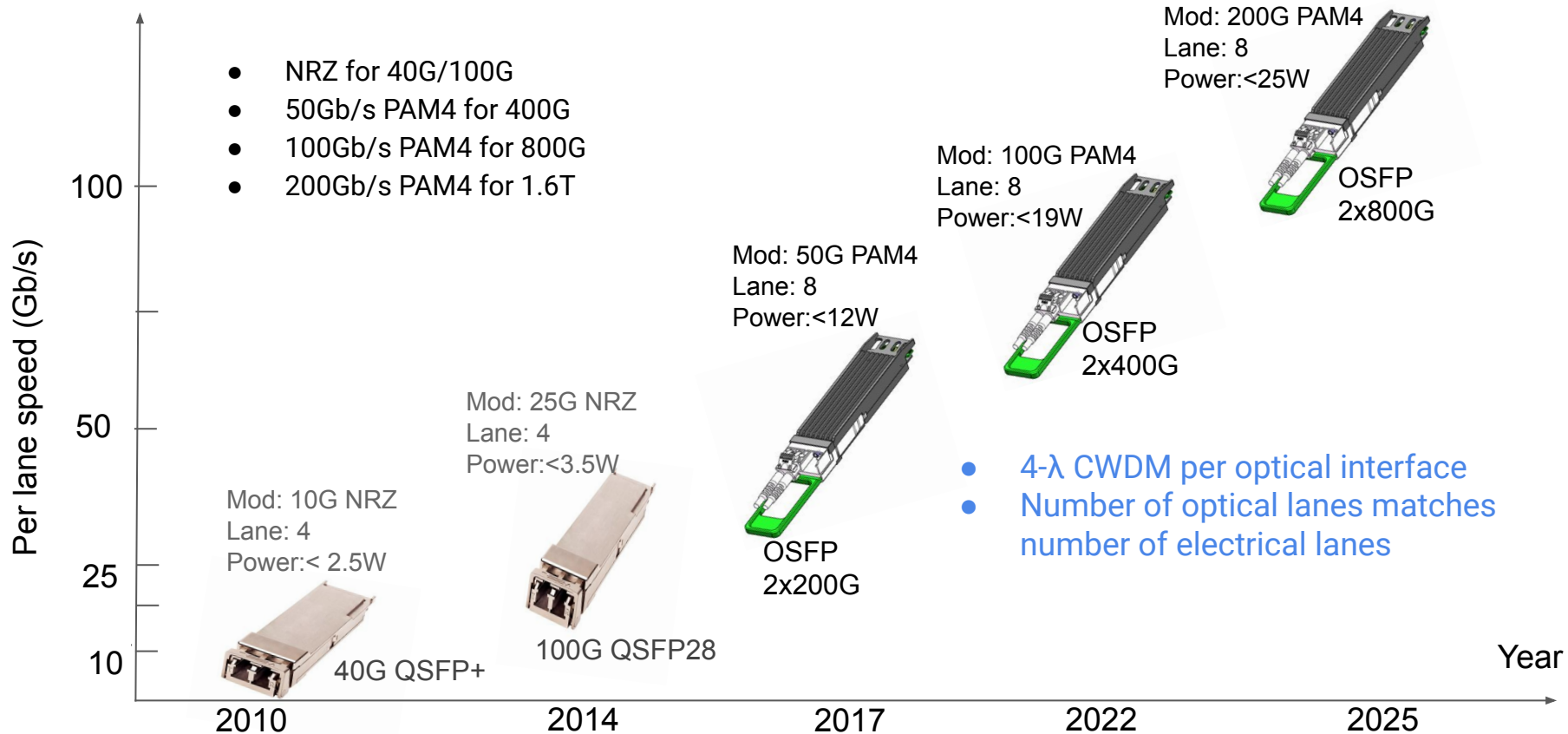
**4** Forward traffic along computed paths.

Google

- 136x136 input/output ports
- Camera-based mirror control scheme simplifies design/manufacturing

Google Cloud

# :'0  2SWLFDO  7UDQVFHLYHU  IRU  '&1



- NRZ for 40G/100G
- 50Gb/s PAM4 for 400G
- 100Gb/s PAM4 for 800G
- 200Gb/s PAM4 for 1.6T

Per lane speed (Gb/s)

100

50

25

10

Mod: 10G NRZ
Lane: 4
Power:< 2.5W

40G QSFP+

Mod: 25G NRZ
Lane: 4
Power:<3.5W

100G QSFP28

Mod: 50G PAM4
Lane: 8
Power:<12W

OSFP
2x200G

Mod: 100G PAM4
Lane: 8
Power:<19W

OSFP
2x400G

Mod: 200G PAM4
Lane: 8
Power:<25W

OSFP
2x800G

- 4-λ CWDM per optical interface
- Number of optical lanes matches number of electrical lanes

Year

2010        2014        2017        2022        2025

Google Cloud

# 3 D O R P D U  2 S W L F D O  & L U F X L W  6 Z L W F K



[TPUv4, ISCA'23]

- 4x4x4 multi-TPU cubes tied together by LW Fabric
- LW Fabric enables reconfigurable interconnection between elemental cubes
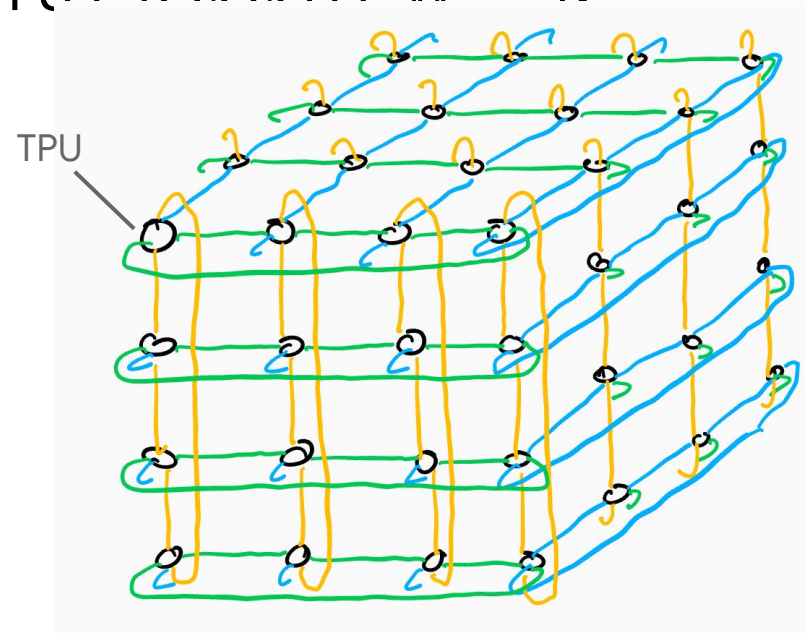  - ML Systems with improved scale, availability, utilization, modularity, deployment, security, power, and performance
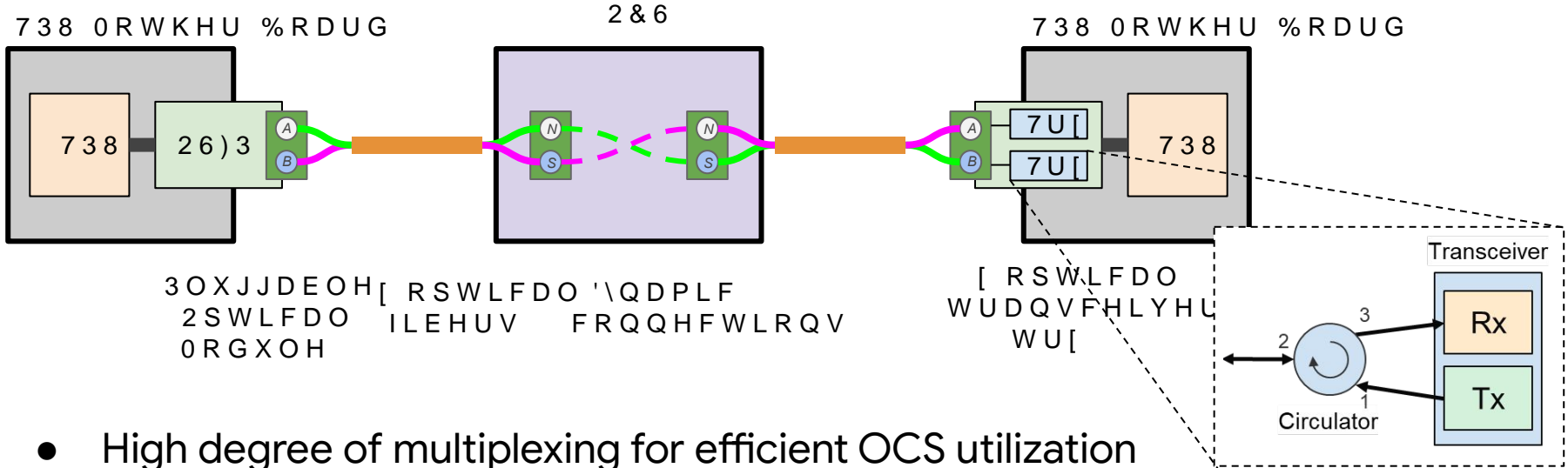
- TPUs are interconnected using a torus network
  - Each dimension is a ring, has three dimensions
- Well matched to requirements
  - High neighbor bandwidth: "allreduce"
  - Low radix
    - Workload is latency tolerant
    - Simpler router, integrated with TPU
  - Low cost ($-per-BW)
    - Mostly in-rack, passive electrical links
- Torus have a long, but niche history
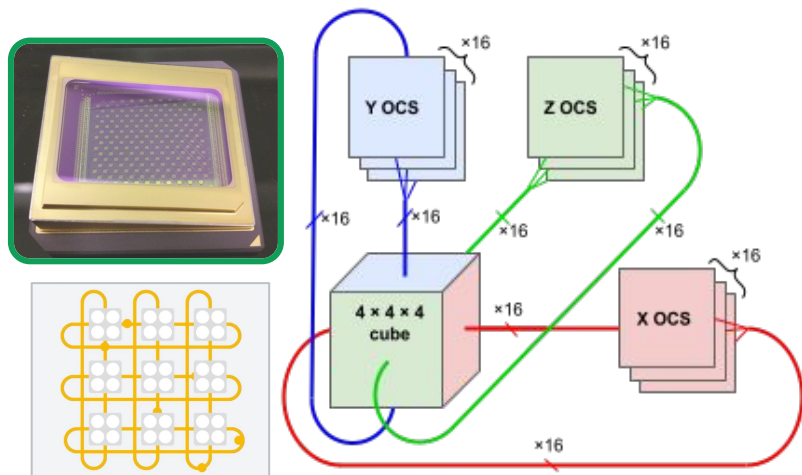  - Notable examples: Cray T3x, IBM BlueGene, Fujitsu K computer



TPU

64 TPUs arranged as a 4×4×4 torus network. Each TPU addressed by X,Y,Z coordinates and connected to six neighbors along X+, X−, Y+, Y−, Z+ and Z−.

Google Cloud

738 0RWKHU %RDUG

2&6

738 0RWKHU %RDUG

738    26)3    A B

N S    N S

A B    7U[ 7U[    738

3OXJJDEOH 2SWLFDO 0RGXOH

[RSWLFDO ILEHUV

2&6 '\QDPLF FRQQHFWLRQ

[RSWLFDO WUDQVFHLYHU WU[
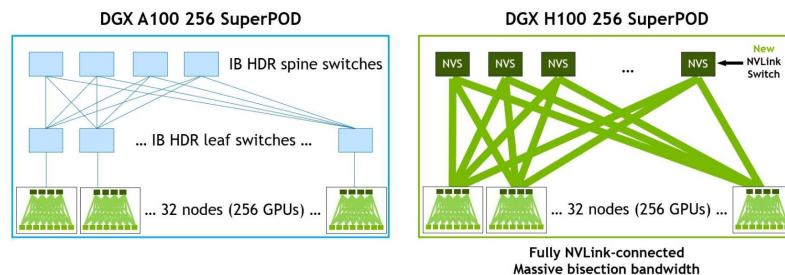
**Transceiver**

3 → Rx

2 ←→

Circulator

1 → Tx

- High degree of multiplexing for efficient OCS utilization
  - WDM transceivers
  - Circulator-based bidirectional links
- Benefits: OCS/fiber is data rate agnostic for extensibility, cost amortization
- Challenges: Higher losses, Multi-Path Interference (MPI) effects

Google Cloud

# OCS vs IB/NVLink and Ethernet



**DGX A100 256 SuperPOD** — IB HDR spine switches, ... IB HDR leaf switches ..., ... 32 nodes (256 GPUs) ...

**DGX H100 256 SuperPOD** — NVS NVS NVS ... NVS, New NVLink Switch, ... 32 nodes (256 GPUs) ...

Fully NVLink-connected
Massive bisection bandwidth

| | A100 SuperPod | | | H100 SuperPod | | | Speedup | |
|---|---|---|---|---|---|---|---|---|
| | Dense PFLOP/s | Bisection [GB/s] | Reduce [GB/s] | Dense PFLOP/s | Bisection [GB/s] | Reduce [GB/s] | Bisection | Reduce |
| 1 DGX / 8 GPUs | 2.5 | 2,400 | 150 | 16 | 3,600 | 450 | 1.5x | 3x |
| 32 DGXs / 256 GPUs | 80 | 6,400 | 100 | 512 | 57,600 | 450 | 9x | 4.5x |

**GCP Innovation with OCS**
Great for upgrade

**Industry with**
**Infiniband/ Nvlink Switch**

Google Cloud

| | Optical Circuit switching | Ethernet switching |
|---|---|---|
| Speeds | 100G/200G/400G/800G | 100G/200G/400G/800G |
| Upgrade | No Change | New Box |
| Upgrade cable | No Change | New Cables |
| Fanout | 136 | 128 |
| Latency | Low | High |
| Flow Control | Per Port/Channel | Per MAC/IP flow |
| Software Define/Error Isolation | Yes | yes |
| Power | 40% | 100% |
| Cost | 30% | 100% |

Google

# AI Smart Offload Transport Innovation

03

# CPU or SmartNic?

## Cloud Offload from X86



Only
CPU
Smart

GPU
SSD
FPGA
Andromeda
TOR

SSD    CPU    GPU    FPGA
PCI EXPRESS

Infrastructure
Processing
Unit
(IPU)

Andromeda
TOR

Google

# AI Transport Smart Offload From Google

**Falcon RDMA**



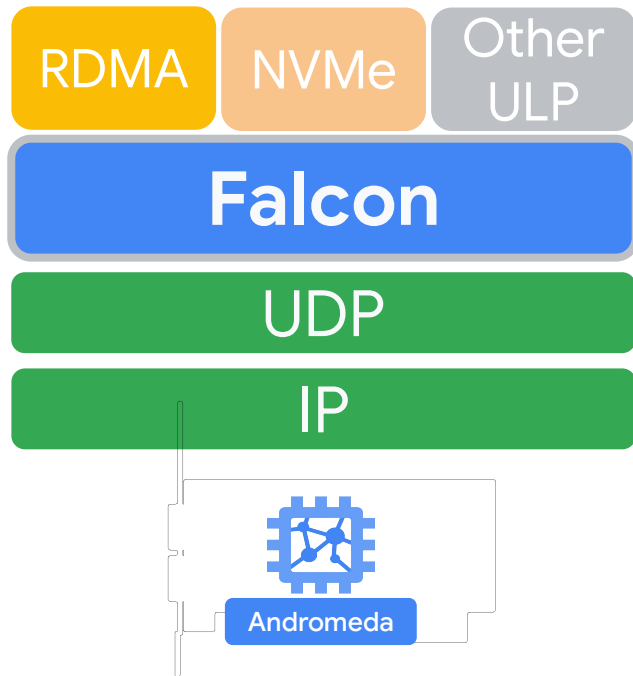**Predictable Efficiency performance @ warehouse-scale:**
hardware acceleration, offloads CPU from data movement,  Low-latency with OS-bypass, massive application bandwidth, mitigating congestion and efficient network utilization.

**Google Falcon**
introduces usability and scalability improvements via relaxed ordering and robust error handling.

**Need of the day:**
meets requirements of critical workloads, HPC and AI; also good for offloading Storage and RPC.

Google Cloud

Diagram showing protocol stack on the left: RDMA, NVMe, Other ULP layered above Falcon, which sits above UDP and IP, with an Andromeda hardware card below.

**Tail Latency** in Ethernet networks
→ HW assisted **delay-based Congestion Control**
→ **Selective ACKs** for fast loss recovery
→ **Multipath** capable connections

**Isolation and Visibility** at scale
→ μs-granularity per-flow **Traffic Shaping**
→ Fine-grained Stats for **Debuggability, Software Defined Network control**

**Efficiency and Security**
→ **Implemented in HW** for Low Latency, High Op Rate using Industry-standard Interfaces, and PSP encryption
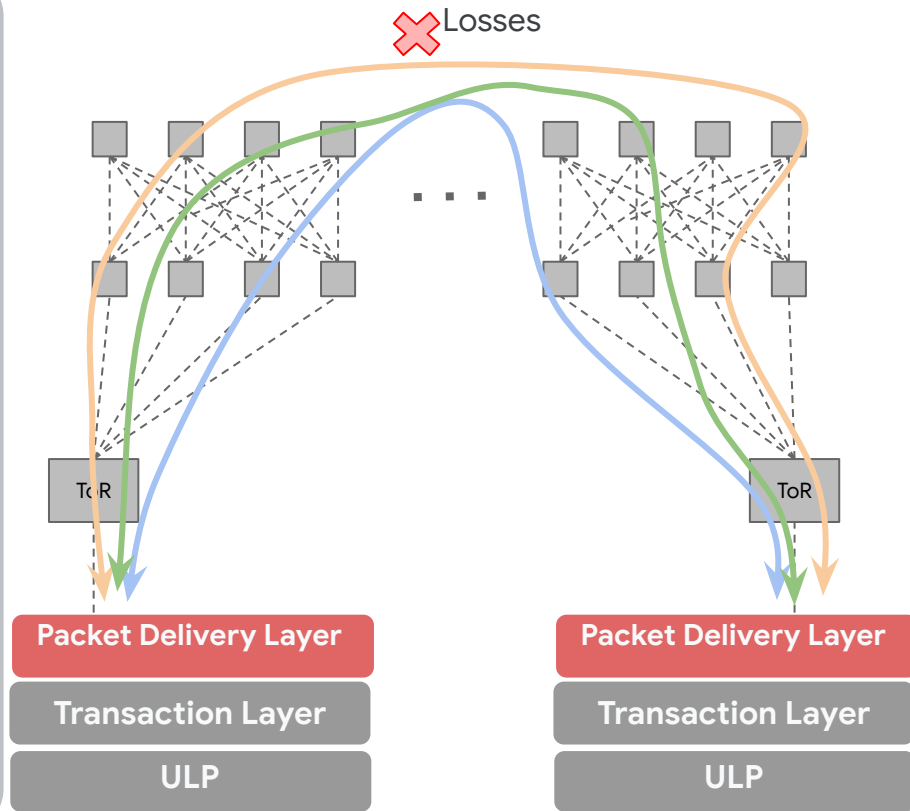
# Falcon Packet Delivery Layer

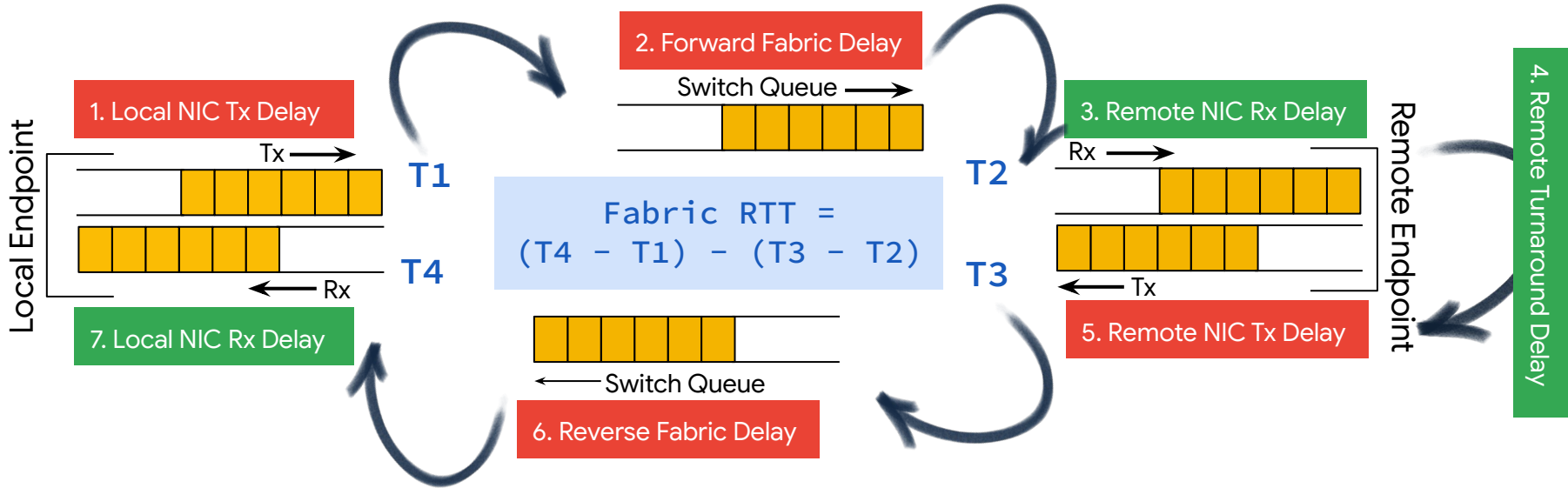**Delay-based Congestion Control** for low latency and high utilization.

**Leverages multiple paths** in the network fabric transparently to applications.

**End-to-end reliable delivery**
- Timely retransmission of lost packets.
- Hardware based retransmission.
- Ack coalescing/piggybacking for high Op rate.



Losses

ToR

ToR

Packet Delivery Layer

Packet Delivery Layer

Transaction Layer

Transaction Layer

ULP

ULP

Google

# Swift Congestion Control as Baseline



**2. Forward Fabric Delay**

Switch Queue →

**1. Local NIC Tx Delay**

Tx →

T1

T4

Rx ←

Local Endpoint

**7. Local NIC Rx Delay**

$$\text{Fabric RTT} = (T4 - T1) - (T3 - T2)$$

**3. Remote NIC Rx Delay**

Rx →

T2

T3

← Tx

Remote Endpoint

**4. Remote Turnaround Delay**

**5. Remote NIC Tx Delay**

← Switch Queue

**6. Reverse Fabric Delay**

Swift* is a delay based  congestion-control for Datacenters that achieves low-latency, high-utilization, near-zero loss implemented completely at end-hosts and NICs supporting diverse workloads like large-scale incast across latency-sensitive, bursty and IOPS-intensive applications working seamlessly in heterogeneous datacenters.

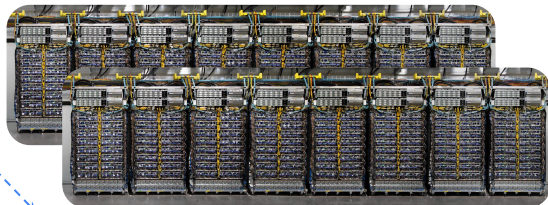*Swift: Delay is simple and effective for congestion control in the Datacenter, SIGCOMM 2020.

# Summary

# Rapid Innovation with Cloud TPUs



**Cloud TPU v2**
- Domain-specific AI supercomputing
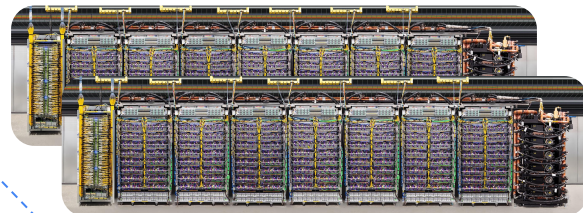- 256 chips distributed shared memory

8x

**Cloud TPU v4**
- Optically reconfigurable 3D Torus
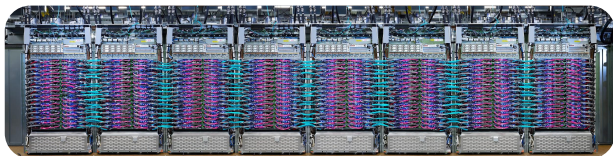- 4k chips with distributed shared memory

20x

**Cloud TPU v5p**
- Programmable Sparsecores for embeddings
- 9k chips with distributed shared memory

2018

2020

2022

2023

2024

**Cloud TPU v3**
- Liquid cooling
- 1k chips distributed shared memory

**Cloud TPU v5e**
- Efficient and scalable training and serving
- 256 chips, horizontally scalable to 10s of k

Google Cloud

# 6 X P P D U \

- **AI Fabric Innovation: What are the Benefits of Lightwave Fabrics for DCN and ML?**
    - Provides direct optical connections (circuits) between network endpoints
    - Comprised of an optical circuit switch (OCS), WDM optical transceivers (trx) co-designed with OCS, circulators, and the hardware/software control plane
    - Reconfigurable, extensible fabric for both datacenter networks (DCN) and ML
    - Enables performant, cost & energy efficient DCNs and & ML supercomputers
        - DCN: 30%/40% reduction of CapEx/OpEx
        - ML: Ability to run large/multi-k node systems; Up to 3.3x speed up in model training

- **AI Software Transport Innovation**
    - Design new Protocol for Remote DMA for AI Network(GPU/CPU and TPU)

Google

# Thank you