

# HPCネットワークの多様化に挑む マルチベンダー×マルチOSで支えるHPCネットワーク運用の実際 (事前資料)

さくらインターネット株式会社

黒澤 潔裕



はじめに

アーキテクチャの試行

ホワイトボックスへの挑戦

自動化への取り組み

まとめ

今日の話は、  
こんな人に向けています

Closを「理想的な構成」だと思っている人

新たなベンダーへの挑戦はコスト削減策だと思っている人

自動化は“楽をするため”だと思っている人

さくらインターネットが挑戦した結果、皆さんがどう思うかを議論したい

## 構築済みHPCクラスタをマネージドサービスとして提供

- **DC基盤**
  - 電力（MW級）、冷却（空冷/液冷）、PUE最適化
- **計算リソース**
  - 計算用のGPU（B200 / H200 / H100）
- **ストレージ**
  - Lustre
- **ソフトウェア基盤**
  - Slurm / MPI / CUDA ・ ROCm / NCCL
- **ネットワーク**
  - Multi-Tenancy/Traffic/Lossless Ethernet/Cabling

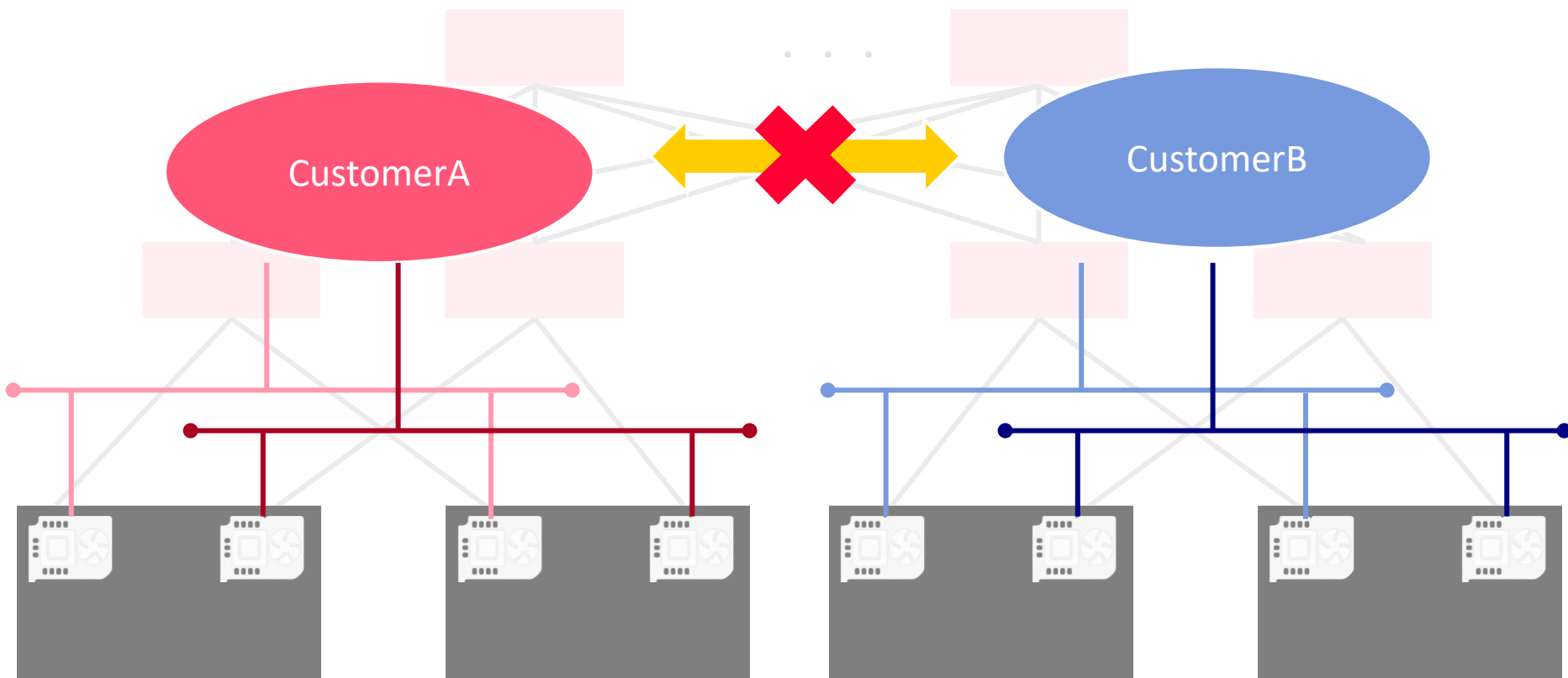


透明性の高い技術による実装を重視

# GPU クラスタにおける、高帯域かつlosslessの実現

## Multi-Tenancy

VRF  
VLAN  
EVPN/VXLAN

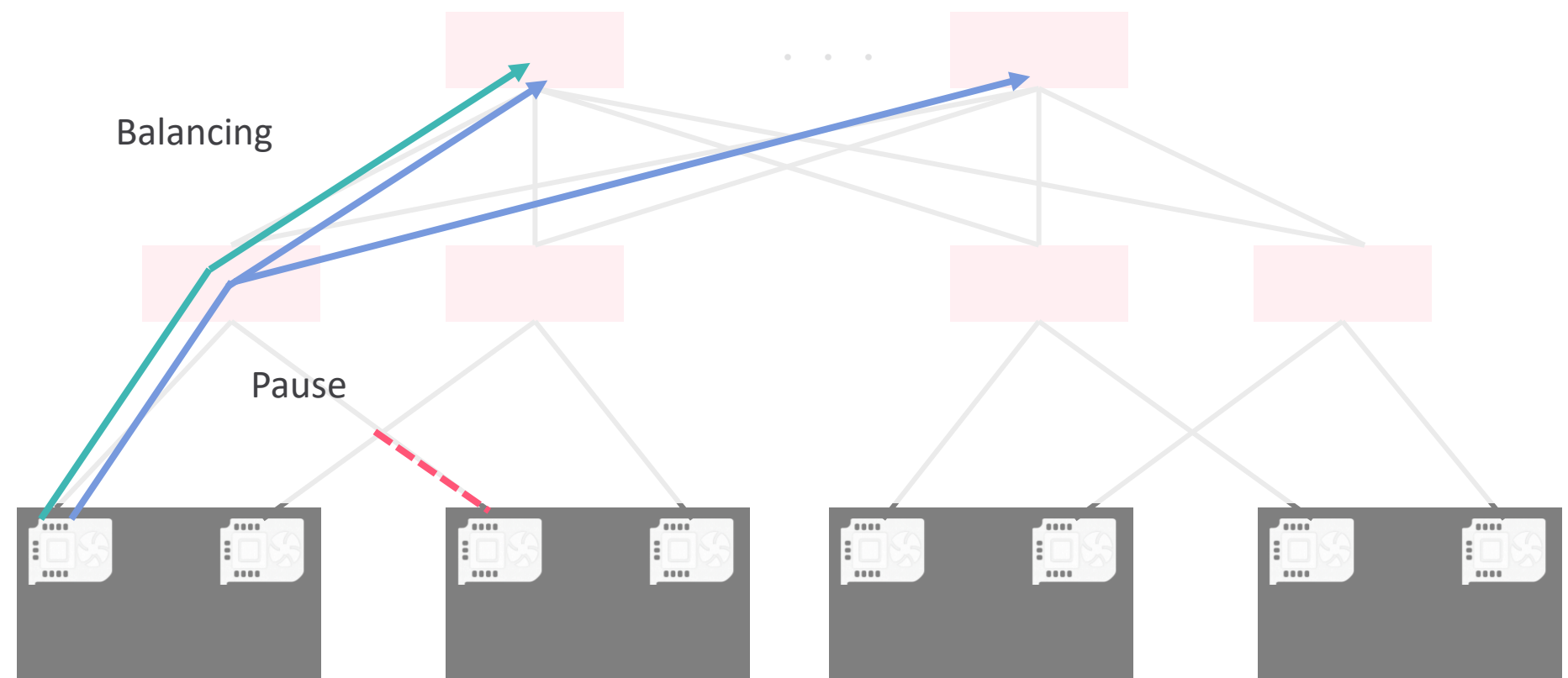


## High Traffic

Switch間接続 800Gbps  
GPU接続 400Gbps

## Lossless Ethernet

RoCEv2(ECN/PFC/CNP)  
Dynamic Load Balancing



## Cabling

高密度収容を実現するケーブル実装  
Switch: 800GBASE-SR8  
GPU: 400GBASE-DR4

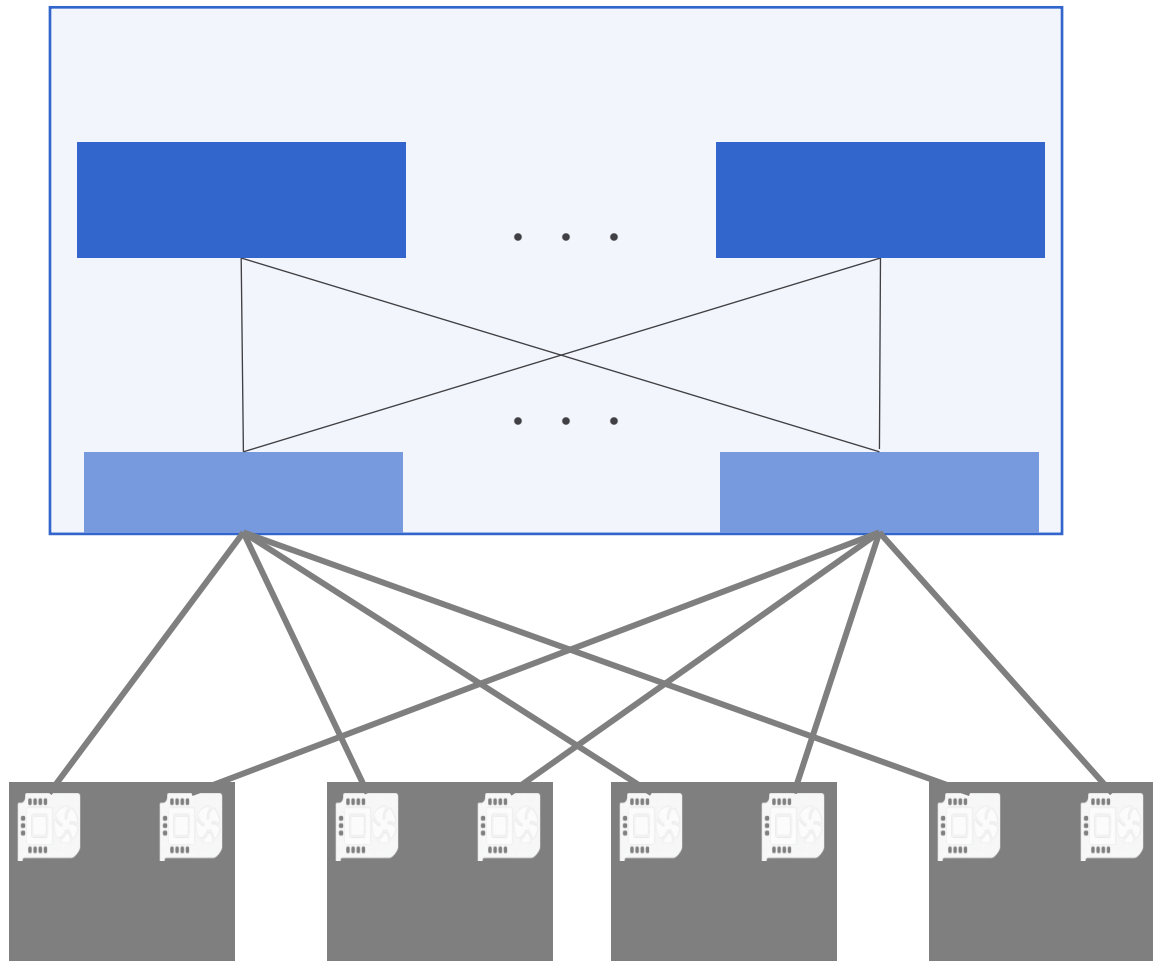


# TopologyとNOSの二軸で新規設計を訴求

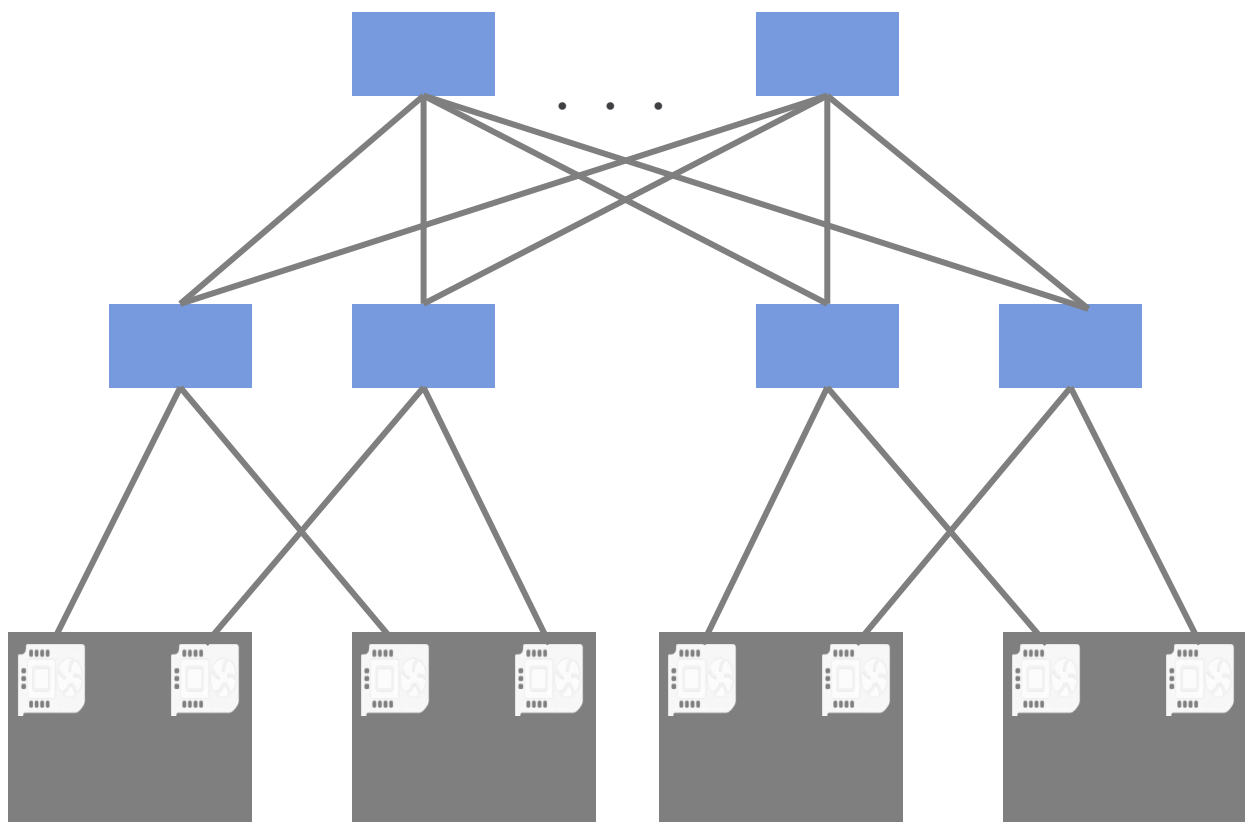
## Architecture

収容上限突破を目的としたClos Topology検討

Chassis Switch



Clos Topology



## Network OS

デリバリー速度と中立的な技術の採用



ARISTA



SONiC  
Edge-core  
NETWORKS

## 私たちの試行錯誤をベースに、広く下記のような議論を行いたい

### アーキテクチャの選定基準

皆さんはどのような基準でTopologyを選択しますか？

### GPU基盤におけるマルチベンダーの検討

皆さんはどのように「マルチベンダー」を実現していますか？

### 増大する運用コストへのアプローチ

どのような解決を考え、どのようなスキルセットが必要でしたか？

当日は皆さんの考えと共に広く議論ができればと思います

# Appendix

# Aristaシャーシを用いたGPUクラスタ設計

JANOG54にて発表。シャーシスイッチによる設計当時の登壇資料

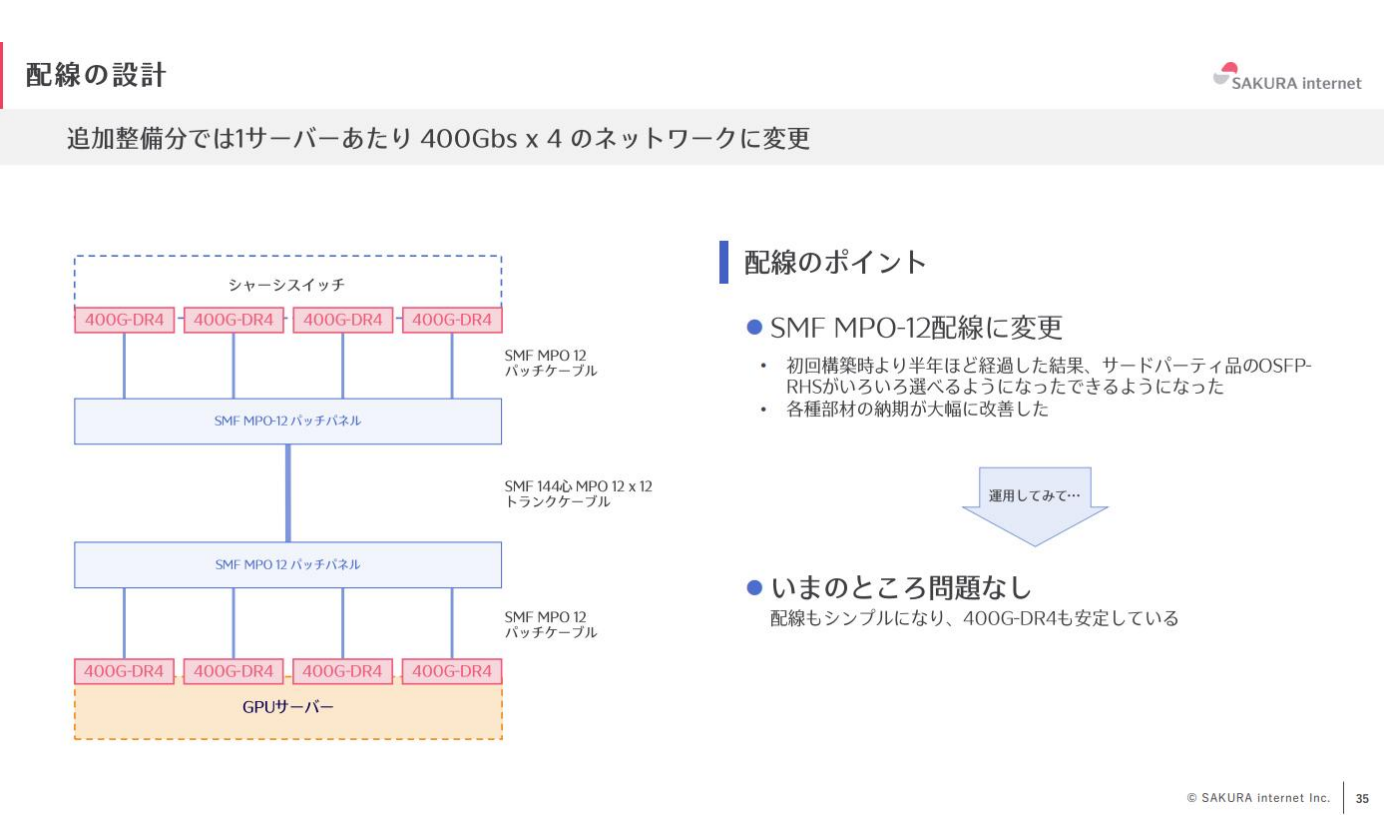
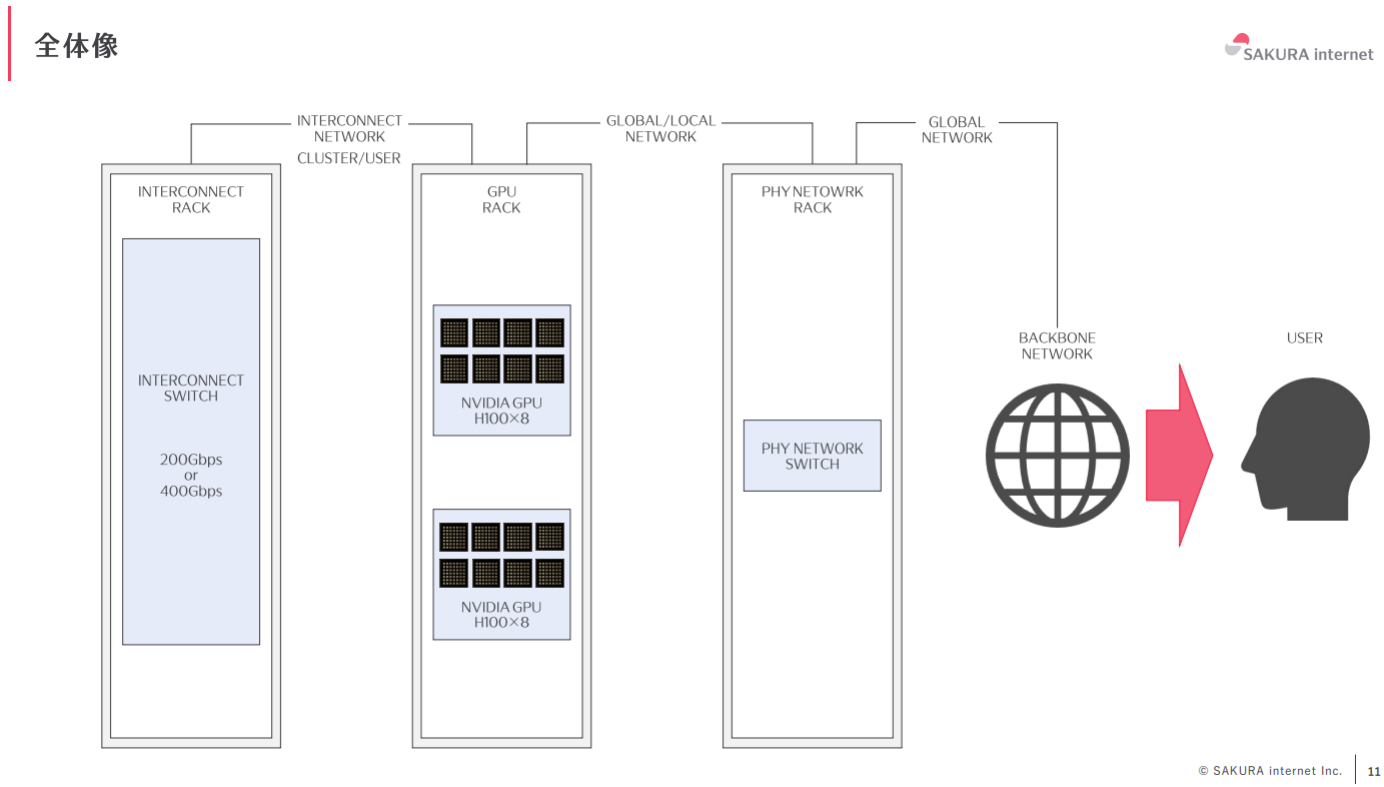
## JANOG54 meeting

生成AI向けパブリッククラウドサービスをつくってみた話

### JANOG54 Meeting

### 生成AI向けパブリッククラウドサービス をつくってみた話

2024年7月4日  
さくらインターネット株式会社  
井上 喬規  
高峰 誠  
平田 大祐



# SONiCを用いたGPUクラスタの構成

SONiCにてClos Topologyを実装した際の記事（2025/6 ISC TOP500 49位にランクイン）

## Arxiv

SAKURAONE: Empowering Transparent and Open AI Platforms through Private-Sector HPC Investment in Japan  
Fumikazu Konishi


## The Linux Foudation User stories

Open Networking at Scale: How SAKURA internet Deployed a TOP500 GPU Supercomputer with SONiC

## SONiC Workshop Japan 2025

SONiCで構築・運用する生成AI向けパブリッククラウドネットワーク

SAKURAONE: EMPOWERING TRANSPARENT AND OPEN AI PLATFORMS THROUGH PRIVATE-SECTOR HPC INVESTMENT IN JAPAN

 **Fumikazu Konishi**  
Research Center  
SAKURA internet Inc.  
Japan

SAKURAONE

USER STORY

Open Networking at Scale: How SAKURA internet Deployed a TOP500 GPU Supercomputer with SONiC

SONiCで構築・運用する生成AI向けパブリッククラウドネットワーク

さくらインターネット株式会社  
黒澤 潔裕

  
© SAKURA internet Inc.

構成要素の検討



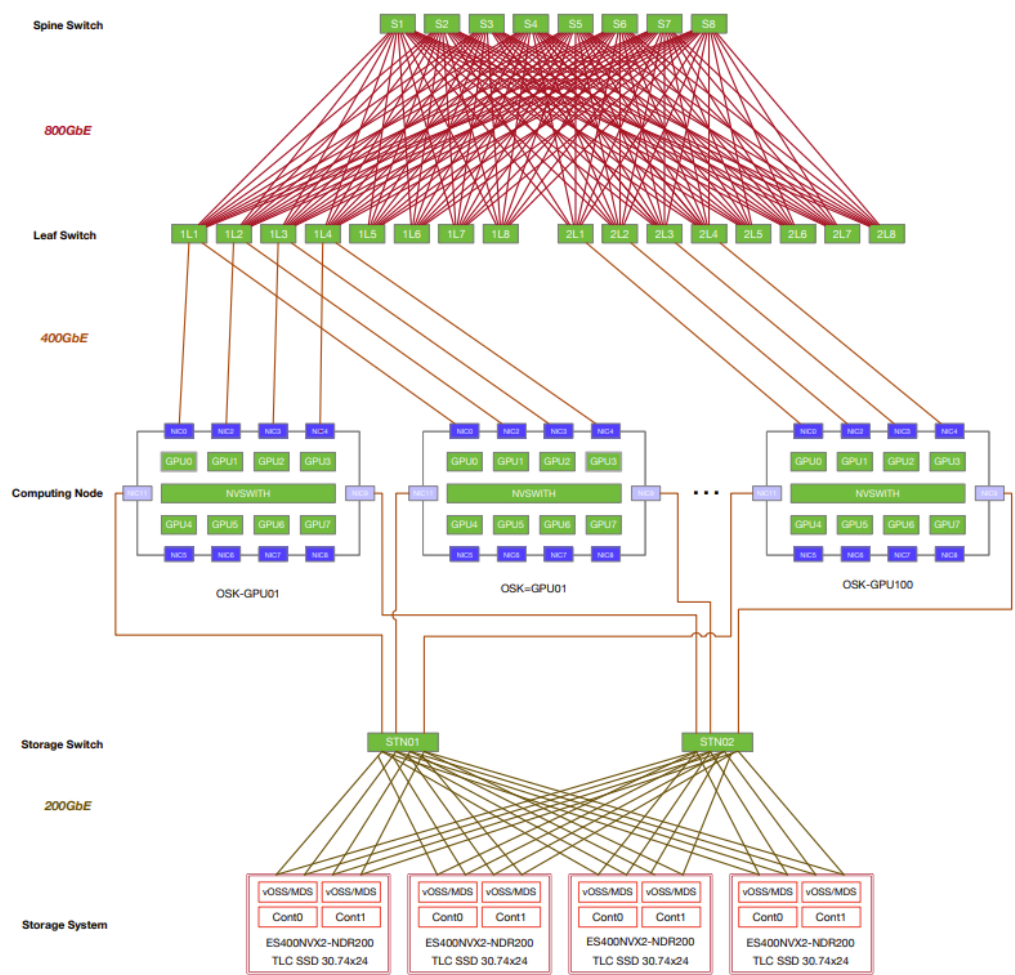


Figure 2: SAKURAONE Network Overview

## Organization

SAKURA internet Inc. is an internet company founded in 1996. Under the corporate philosophy of “Turning ‘what you want to do’ into ‘what you can do,’” we develop a variety of services to meet customer needs and propose DX solutions that cater to various industries.

Since our founding, which began with the provision of shared server services, we have expanded to offer services such as “Koukaryoku” to support generative AI, and “SAKURA Cloud,” which has been conditionally certified for use in government cloud systems. A key feature of our company is that we handle everything from development to operations in-house.

## Overview

SAKURA internet is responding to the growing demand for computational infrastructure driven by the rapid adoption of generative AI by continuously procuring next-generation GPUs and strengthening reliable operational systems in our own data centers. As a digital infrastructure company contributing to the sustainable development of the digital society, our mission is to provide cloud services for generative AI.

To continuously meet the increasing demand, we recognized the necessity of resolving the following challenges:

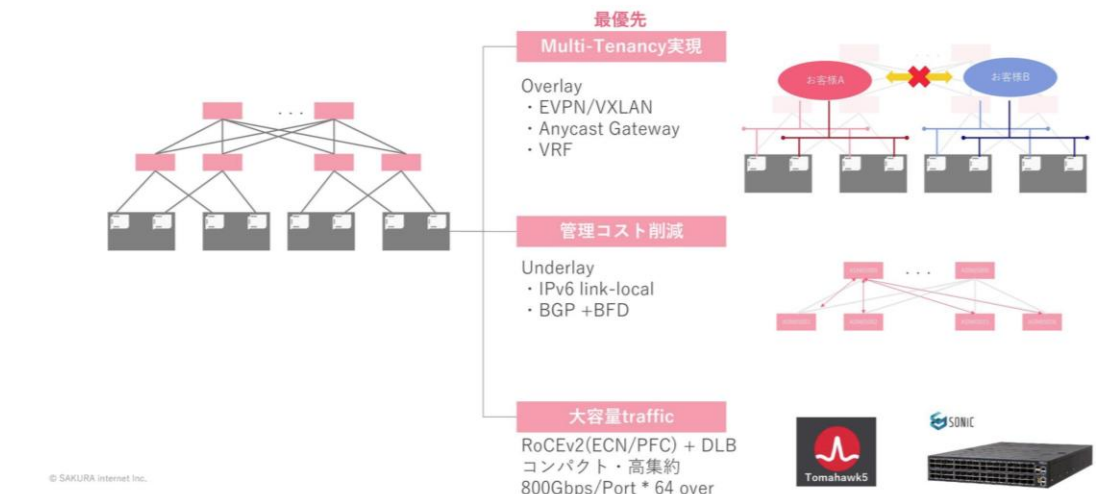
- Ensuring vendor-neutral supply to mitigate risks
- Adopting technologies with high neutrality
- Accelerating delivery speed

## Why SONiC?

We selected SONiC for our new 800-GPU cluster because it directly addressed these needs. SONiC provided:

- High transparency, as it is implemented on a Debian/ Linux platform
- Active development of new features supported by the global community
- The ability to streamline operations by leveraging the same technologies used for Linux servers

In addition, SAKURA internet has a corporate culture of leveraging OSS and bottom-up initiatives. This culture supported our adoption of SONiC and enabled us to launch a GPU cloud infrastructure in a very short period of time.



## 運用コマンドの拡充

未成熟機能はLinux/Pythonを用いて自社スクリプト開発

随接機器の確認

```
admin@gpmr-008:~$ show ip neighbor
Capability codes: R Router, B Bridge, O Other
LocalPort RemoteDevice RemotePort Capability RemotePortDesc
Ethernet01 de-csa3-gleaf-201 En24(Port24) BR To de-csa3-gpmr-008
Ethernet08 de-csa3-gleaf-201 En24(Port24) BR To de-csa3-gpmr-008
Ethernet16 de-csa3-gleaf-203 En24(Port24) BR To de-csa3-gpmr-008
Ethernet24 de-csa3-gleaf-204 En24(Port24) BR To de-csa3-gpmr-008
```

Port? Interface? = 混乱  
Descriptionは要らない

configの差分確認

```
admin@gpmr-001:~$ sudo diff startup-config.log running-config.log
(Sort違いにより大量の差分)
```

LLDP/BGPの情報を整理

```
admin@gpmr-008:~$ show neighbor
LocalPort RemotePort RemoteDevice RemotePort RemoteDevice RemotePort RemoteDevice RemotePort RemoteDevice RemotePort RemoteDevice RemotePort RemoteDevice RemotePort RemoteDevice
Ethernet01 Port1 de-csa3-gleaf-201 Ethernet184 Port24 Established Up
Ethernet08 Port2 de-csa3-gleaf-202 Ethernet184 Port24 Established Up
Ethernet16 Port3 de-csa3-gleaf-203 Ethernet184 Port24 Established Up
Ethernet24 Port4 de-csa3-gleaf-204 Ethernet184 Port24 Established Up
```

SONiCの設定を機序で差分比較

```
admin@gpmr-008:~$ show config compare
SUCCESS: SONiC Running Configuration の取得が完了しました。
SUCCESS: SONiC Startup Configuration の読み込みが完了しました。
ALERT: SONiC Config の差分が検出されました！

--- SONiC Config Running Config
+++ SONiC Config Startup Config
@@ -3888,7 +3888,8 @@
"PLEX_COUNTER_STATUS": "enable"
}
--- WRED, ECN, QUEUE:
- "PLEX_COUNTER_DELAY_STATUS": "false",
+ "PLEX_COUNTER_STATUS": "enable"
}
```

本当の差分はどこにある？