

AI-RANの新規技術開発発表 事前説明会

説明者



船吉 秀人

ソフトバンク株式会社

先端技術研究所 先端無線統括部

統括部長

ソフトバンク AI-RANの取り組み

ソフトバンクのAI-RANの取り組み



Press Releases 2023

NVIDIA Collaborates With SoftBank Corp. to Power SoftBank's Next-Gen Data Centers Using Grace Hopper Superchip for Generative AI and 5G/6G

Arm-Based Superchip and BlueField-3 DPU Power Revolutionary Architecture to Enable Generative AI-Driven Wireless Communications

May 29, 2023
NVIDIA
SoftBank Corp.

NVIDIA and SoftBank Corp. today announced they are collaborating on a pioneering platform for generative AI and 5G/6G applications that is based on the NVIDIA GH200 Grace Hopper™ Superchip and which SoftBank plans to roll out at new, distributed AI data centers across Japan.



Founding Members

SoftBank NVIDIA arm
NOKIA ERICSSON SAMSUNG
T Mobile Microsoft DEEPSIG
Northeastern University THE UNIVERSITY OF TOKYO

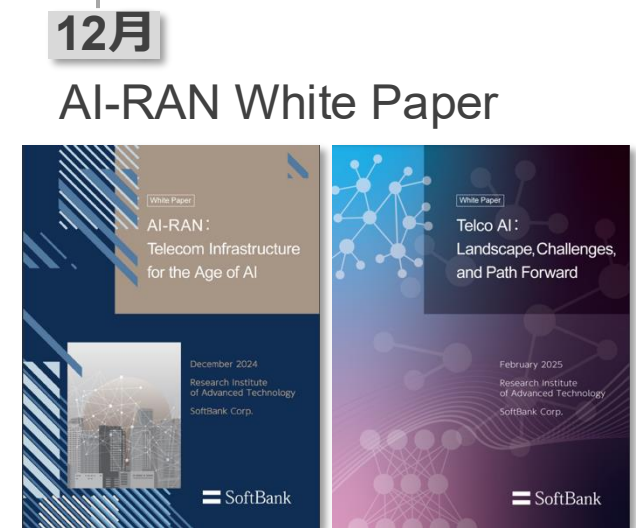
76 社

※2025/2/27時点



パートナーシップ

SoftBank NVIDIA
SoftBank Red Hat
SoftBank FUJITSU



AITRAS

AI-RANコンセプトにもとづいた
ソフトバンクオリジナルのプロダクト

AITRAS のシステム構成

: ソフトバンク開発部分

AITRAS

AITRAS オーケストレーター

管理・制御ソフト

Edge AIのサービス
メニューを提供



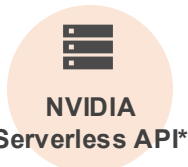
LLM
ロボット



マルチモーダル
LLM



高機能
RAG



NVIDIA
Serverless API*1

Edge AI Apps

NVIDIA AI Enterprise*1

RAN L2/L3 ソフトウェア

RAN L1 ソフトウェア

NVIDIA AI Aerial

無線信号の処理ソフトを
高度化

仮想化基盤を構築・実装

仮想化基盤

MIG/MPS



NVIDIA GH200*2

Arm Neoverse V2

Radio Unit



*1: Serverless API powered by NVIDIA AI Enterprise

*2: NVIDIA GH200 Grace Hopper Superchip

AITRAS のシステム構成

AITRAS

AITRAS オーケストレーター

Edge AI Apps

NVIDIA AI Enterprise

RAN L2/L3 ソフトウェア

RAN L1 ソフトウェア

NVIDIA AI Aerial

仮想化基盤

MIG/MPS

NVIDIA GH200

Arm Neoverse V2



NVIDIA GH200や
NVIDIA AI Enterpriseなど
さまざまなアセットを活用

FUJITSU

Red Hat

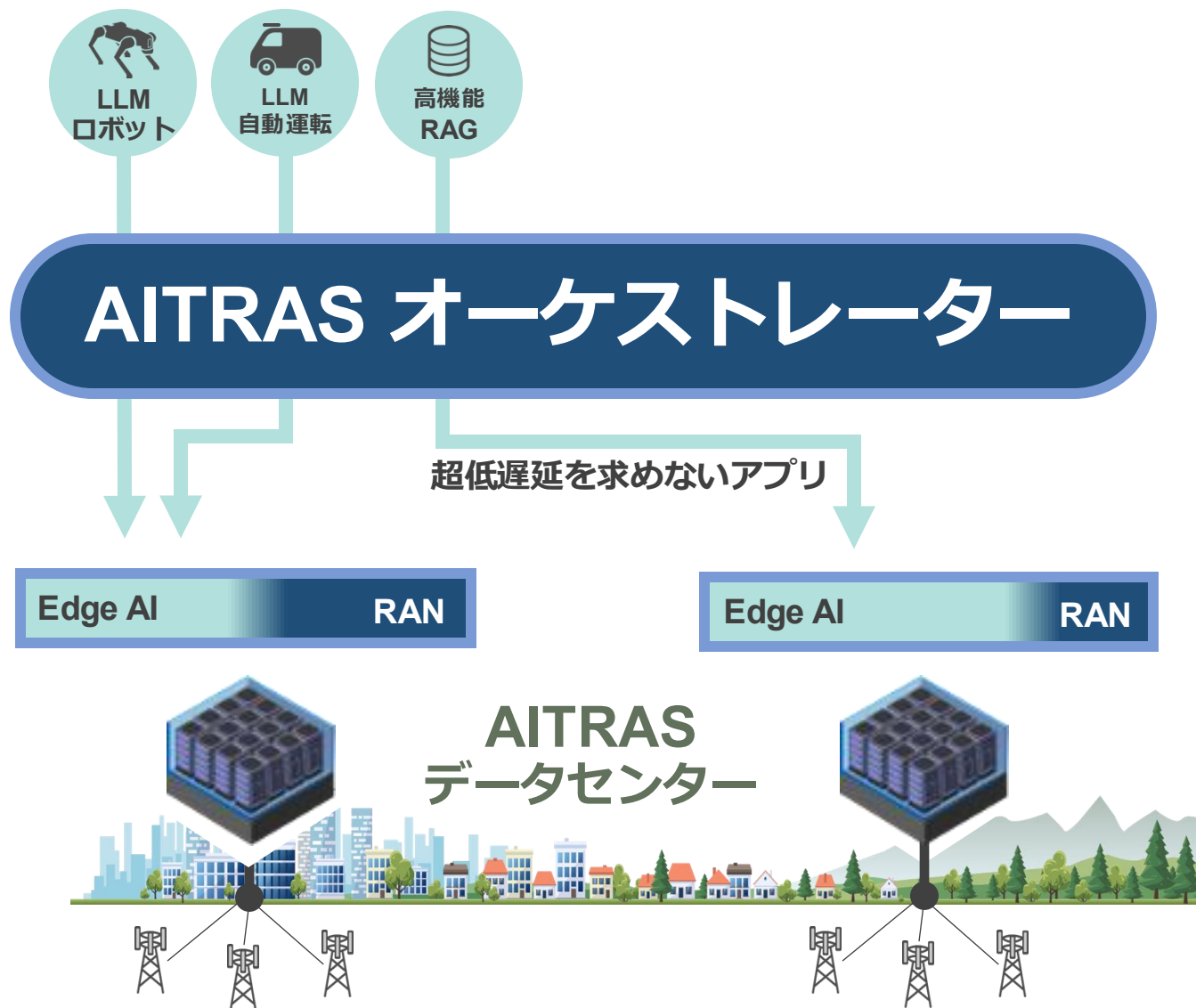
FUJITSU

Radio Unit



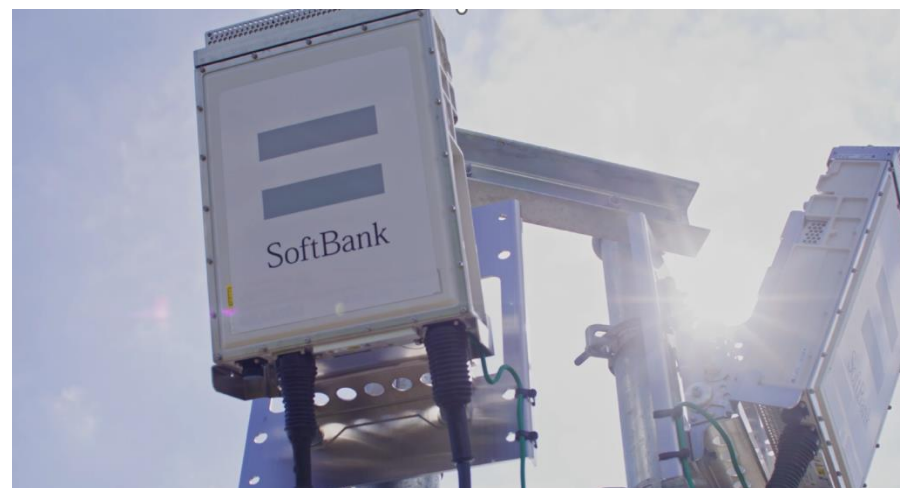
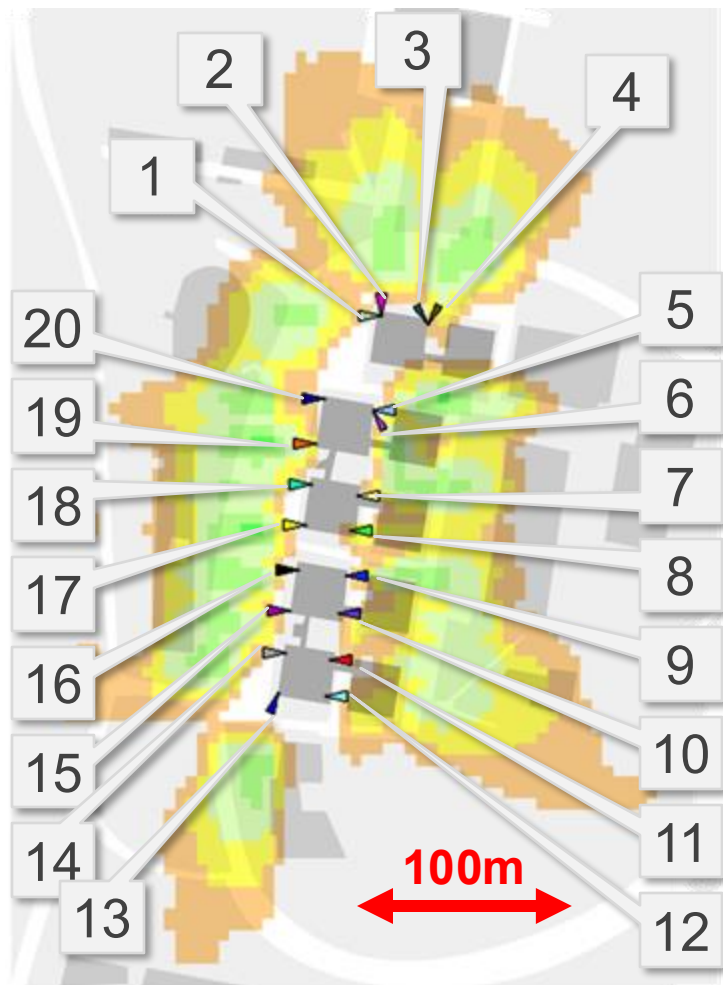
NVIDIA. arm
1サーバーで
20セルを実現

AITRAS オーケストレーター



- AIによる自動リソース割当
- サーバーの役割変更
- NVIDIA Serverless API, NVIDIA AI Enterprise 対応

AITRAS 20 セルをオンエア



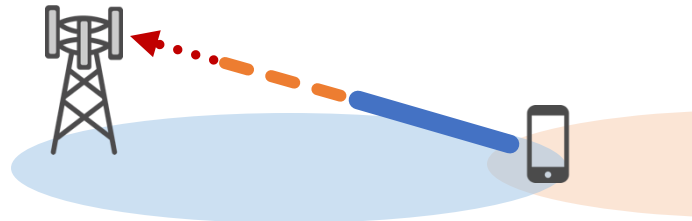
AI-RANに関する 新規開発の発表

①

AI技術によるRANの性能向上効果を実証

AI for RAN技術を開発

UL Signal Processing



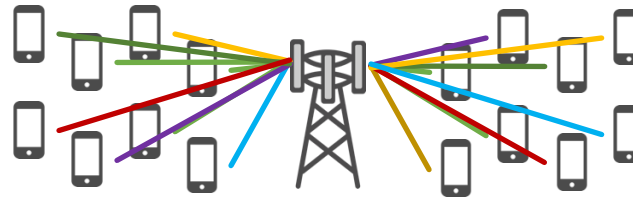
ノイズが多いエリアやセルエッジでは受信感度が低く
チャンネルエラーに影響

AITRAS



チャンネル補間

MU-MIMO



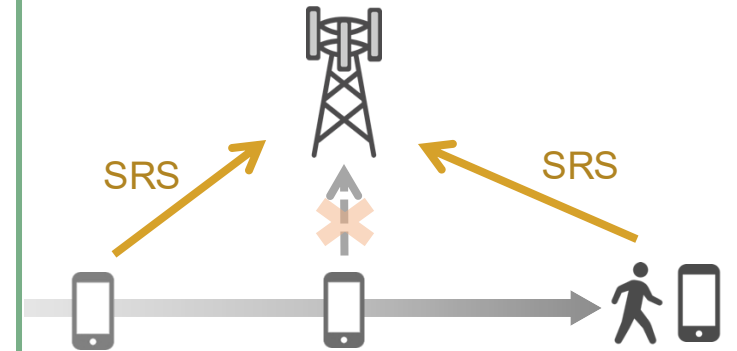
端末のペアリングに関する
膨大なマトリクス計算が発生

AITRAS



ユーザーペアリング
の最適化

Beamforming



接続端末が増えると
SRSの送信間隔が大きくなる

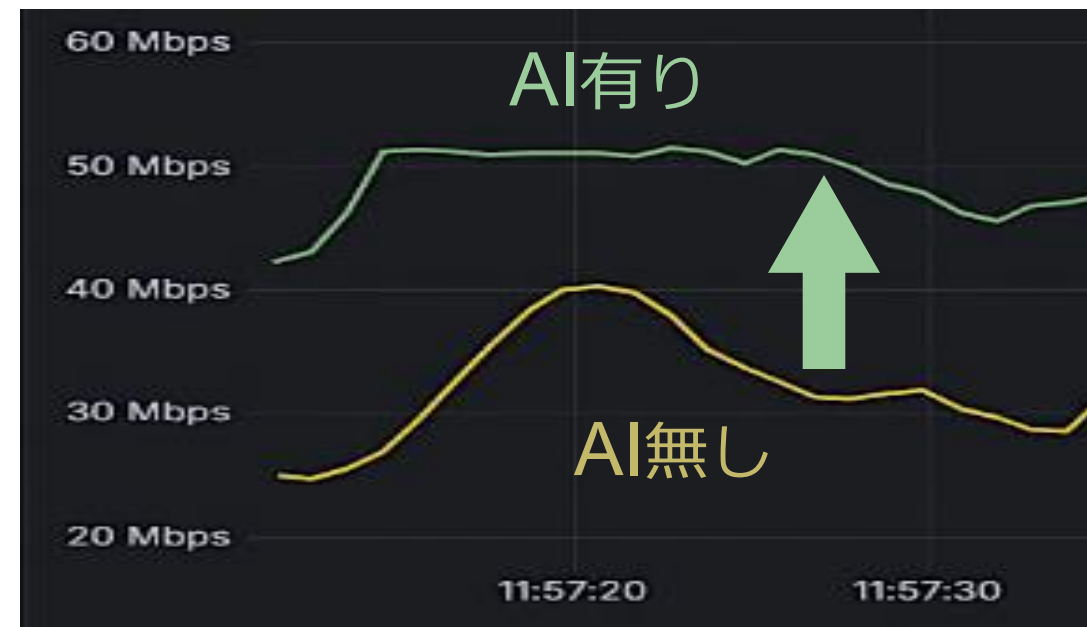
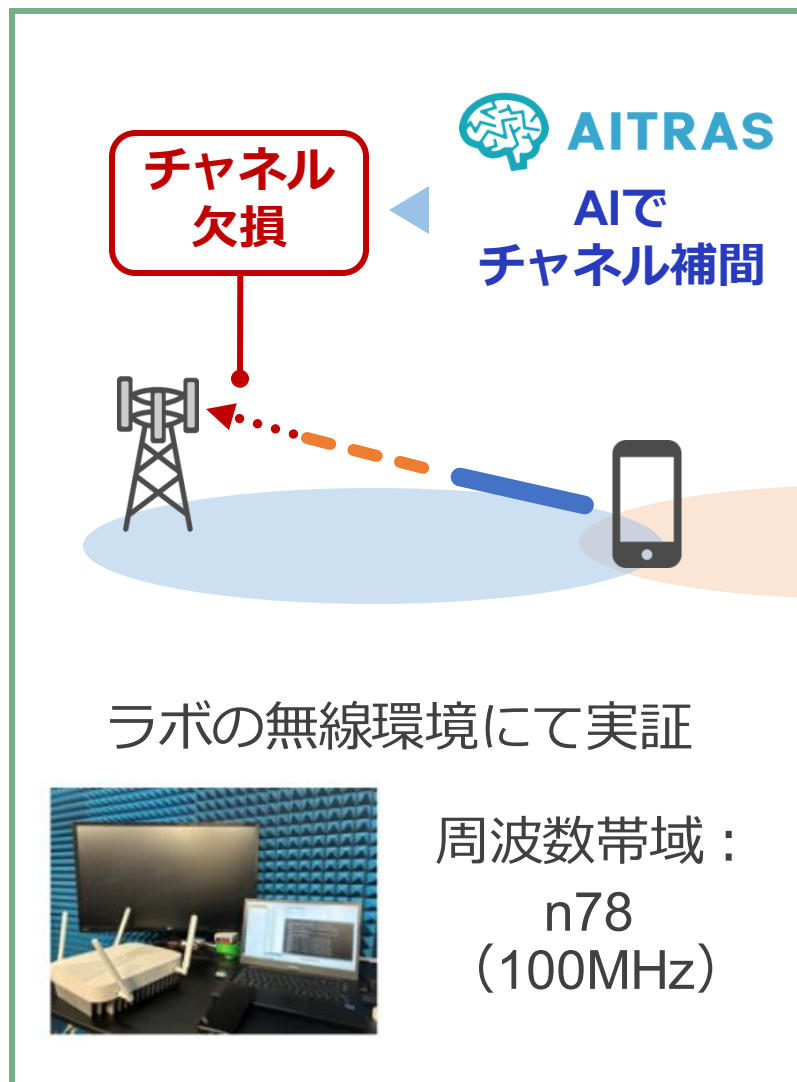
SRS : Sounding Reference Signal
サウンディング参照信号

AITRAS



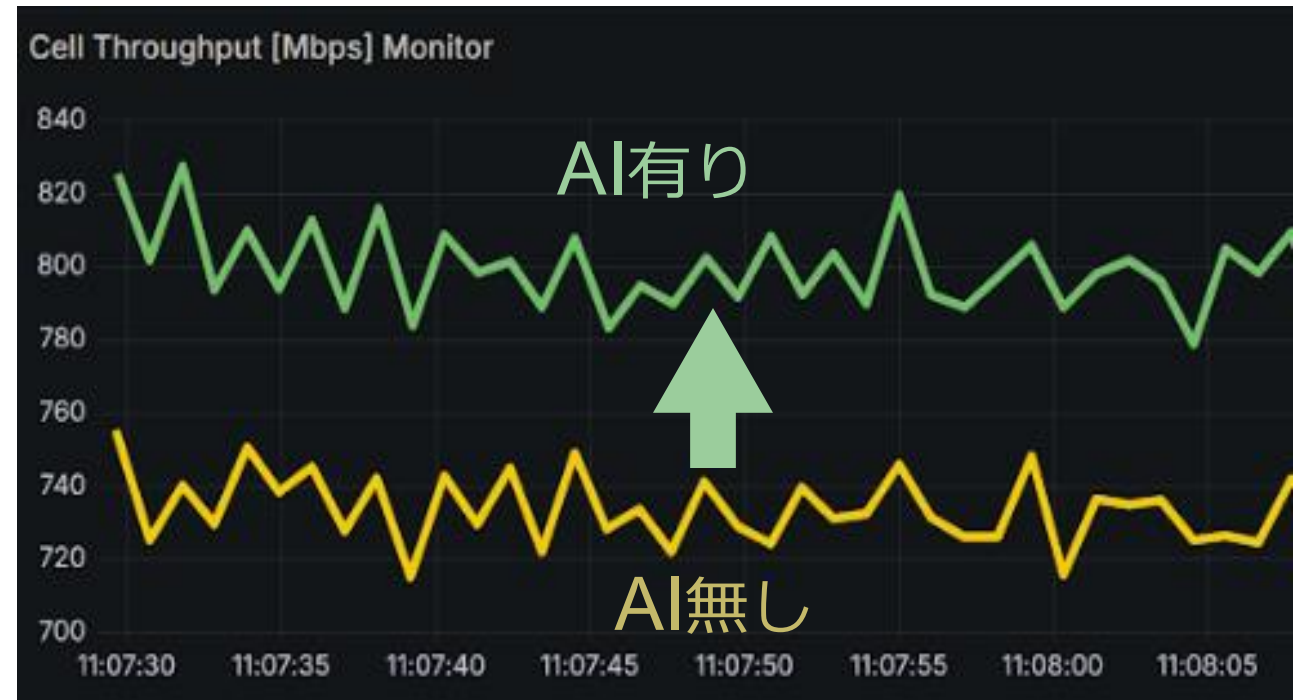
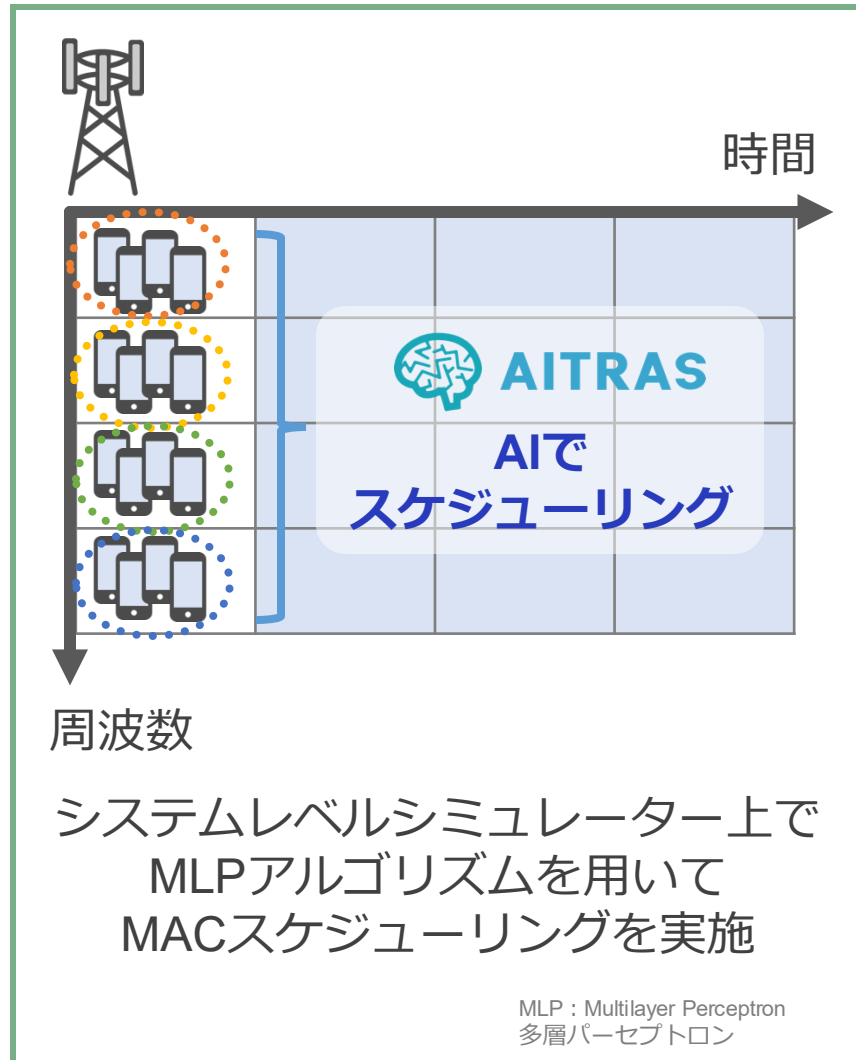
SRS 予測

ULチャネル補間



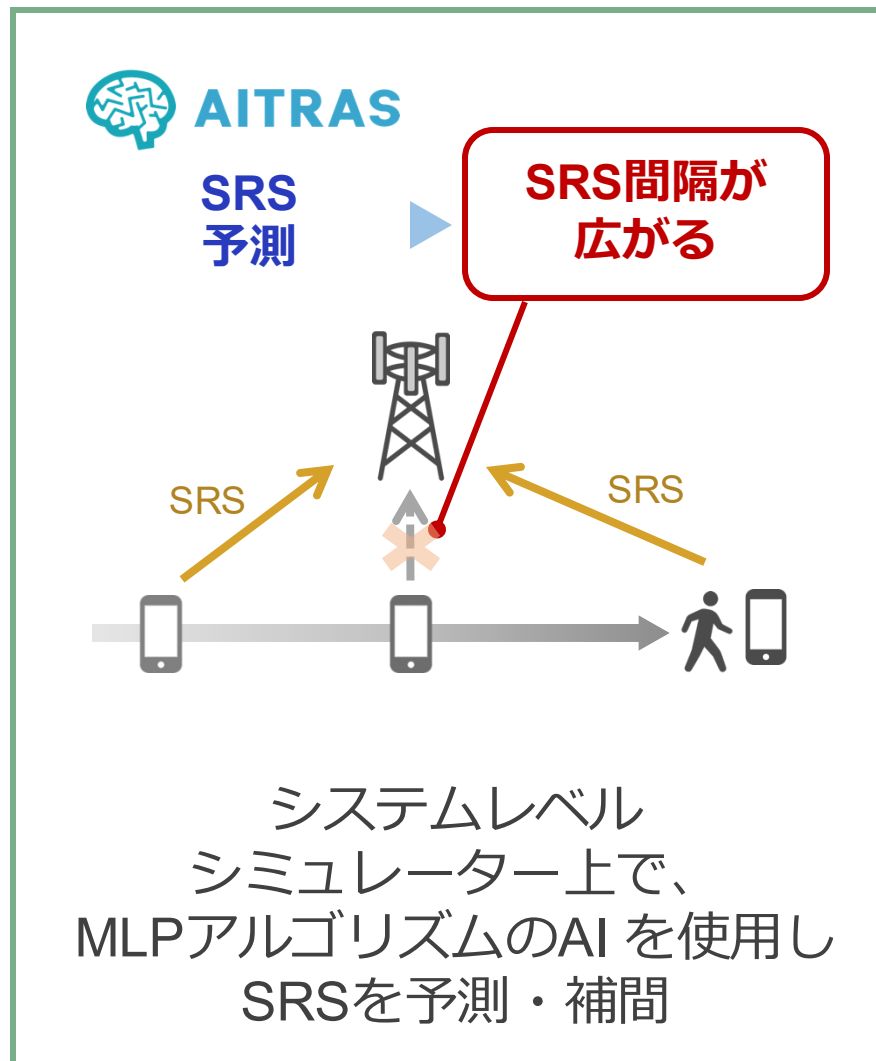
**スループットが
約20%向上**

MU-MIMO : ユーザーペアリングの最適化



セルスループットが
約9%向上

Beamforming : SRS 予測



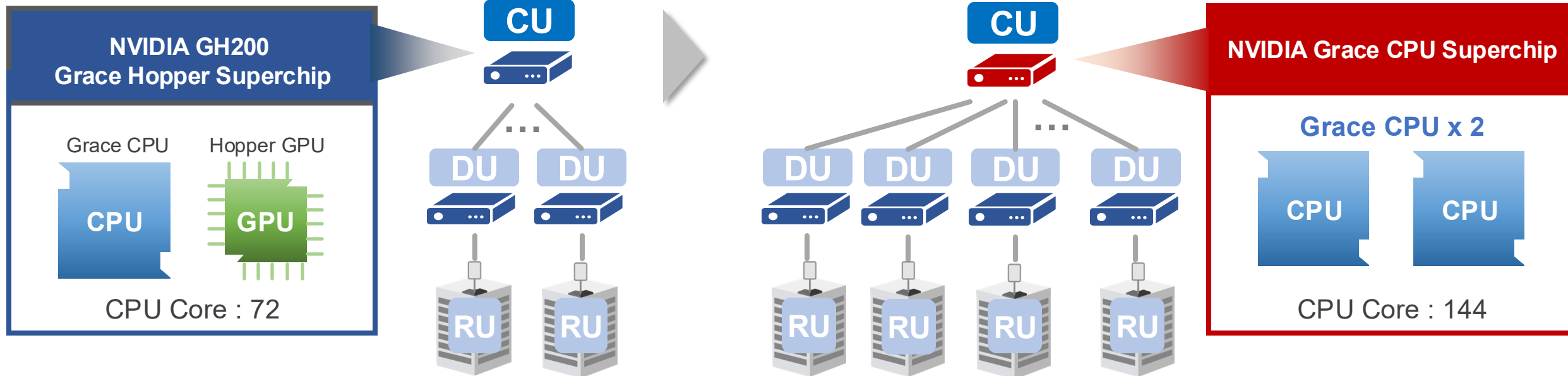
ユーザースループットが
約13%向上

②

「AITRAS」、
ArmベースのNVIDIAプラットフォームを活用した
C-RANとD-RANの
AI-RANアーキテクチャ実装の完了について

AITRASの進化①

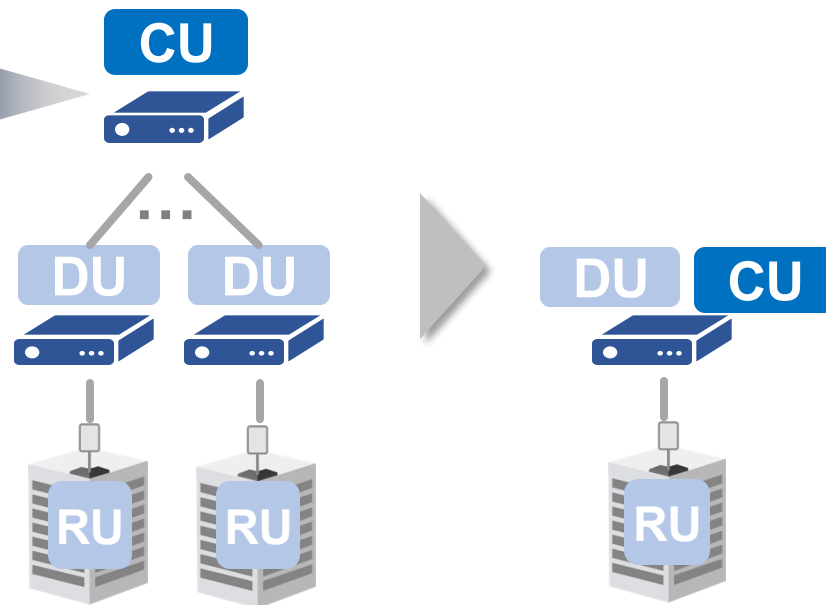
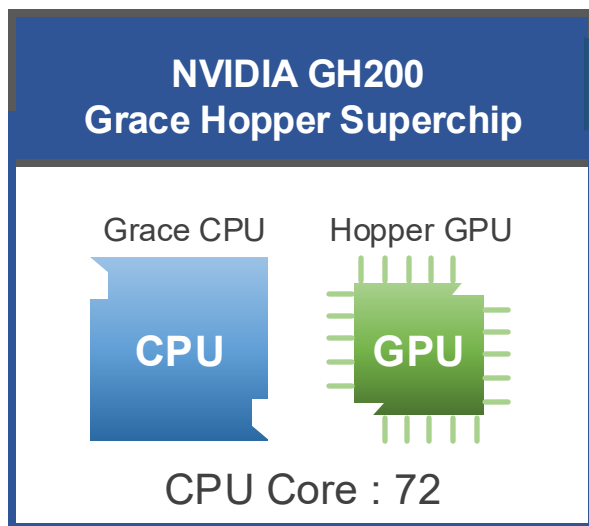
NVIDIA Grace CPU Superchip ServerへのCU実装



NVIDIA Grace CPUによって
CUの収容率はおおよそ2倍に

AITRASの進化②

1台のNVIDIA GH200 ServerへDU、CUの両方を実装



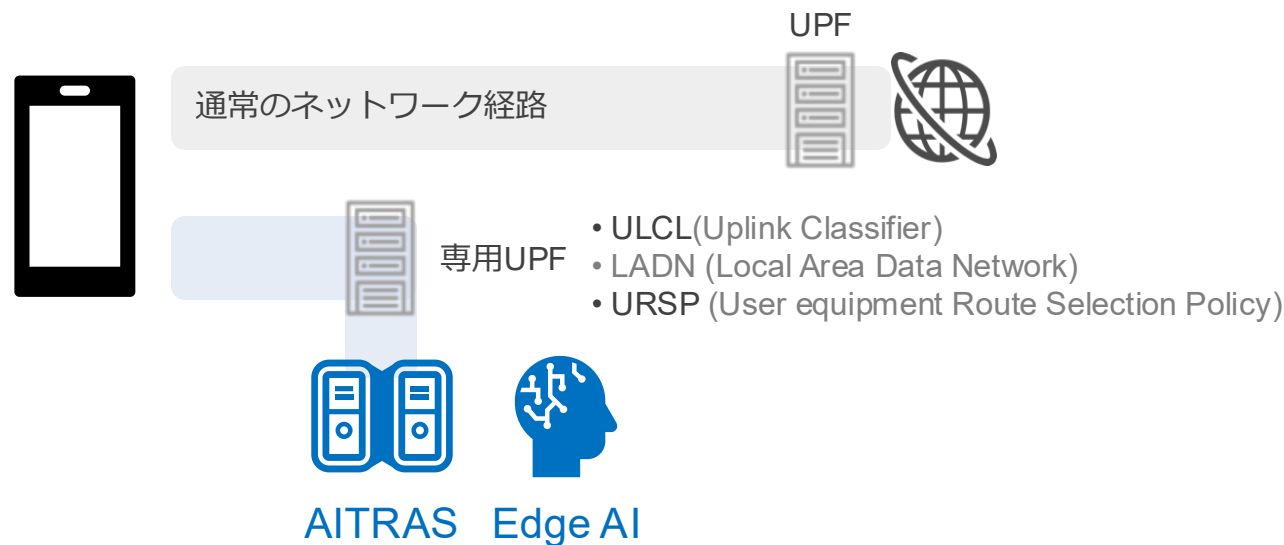
エンタープライズ施設
等で局所的にAI需要が
高まるエリアへ

D-RAN構成の
AITRASを展開可能に

③

ローカルブレイクアウト技術を活用し、
「AITRAS」上のエッジAIサーバーへセキュア
にアクセスする機能を開発

AI on RAN に向けた エッジルーティング技術の開発

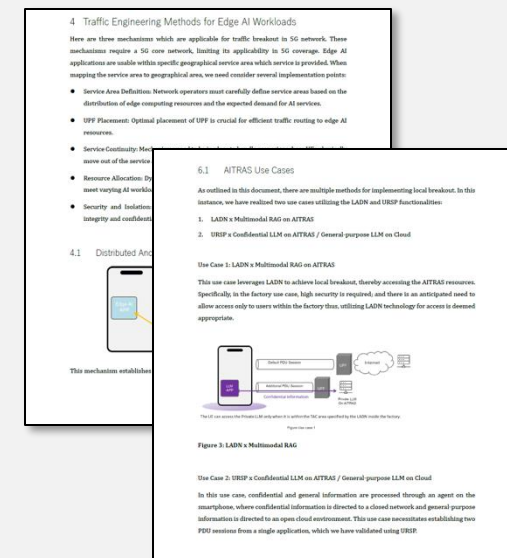
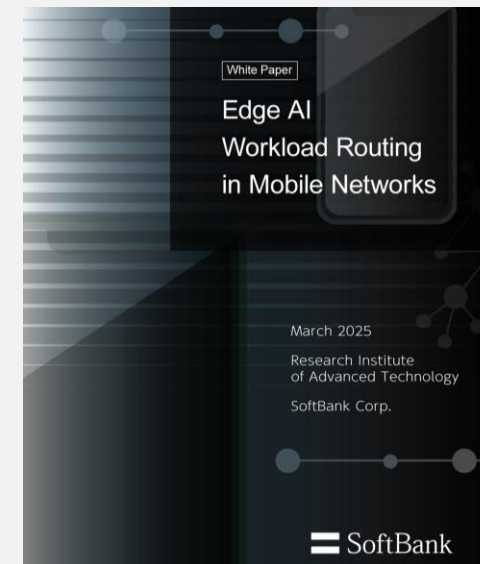


ネットワークエッジでの動的経路選択を可能にするルーティング技術

- Internetとローカルブレイクアウト経路間のスマートな切り替え
- AI on RANの展開オプションを拡大

White Paper 公開

“Edge AI Workload Routing in Mobile Networks”



④

ソフトバンクとレッドハット、
AI-RANのデータセンターにおける
消費電力を最適化するソリューションを開発

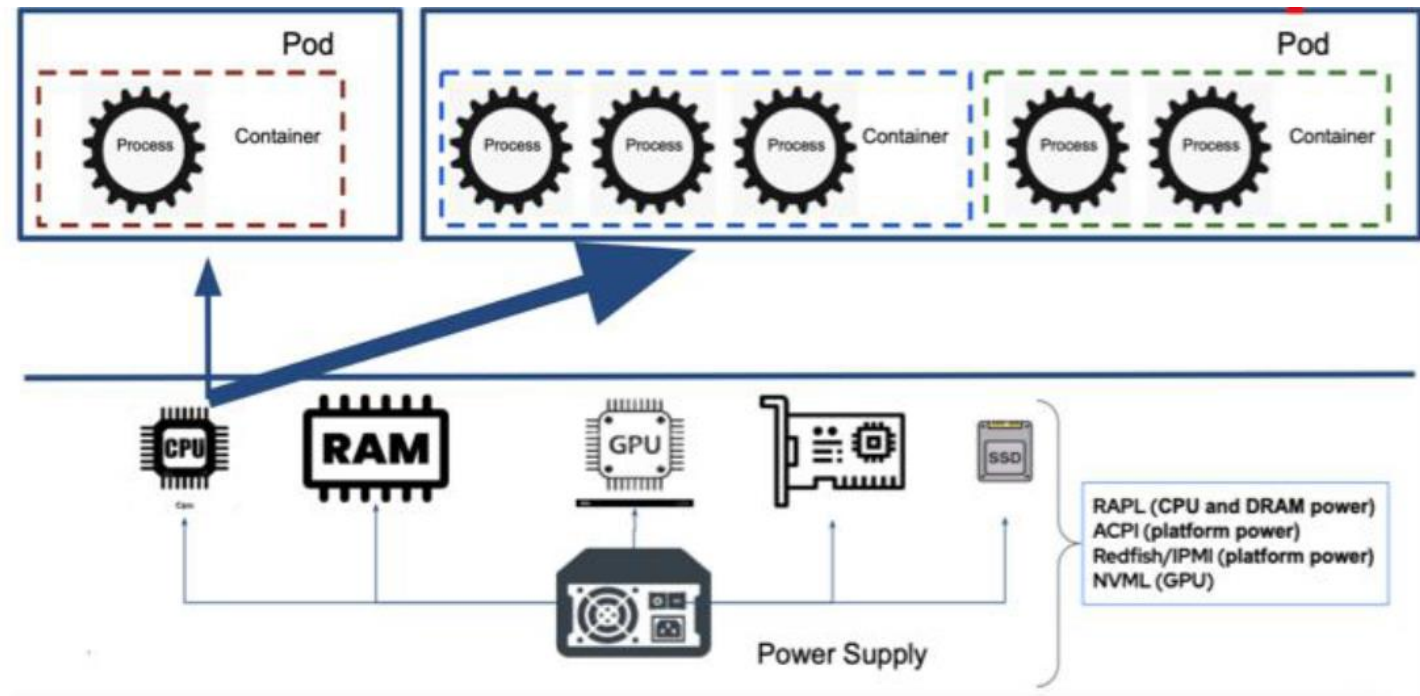
Red Hat Kepler プロジェクト



Red Hat



KEPLER

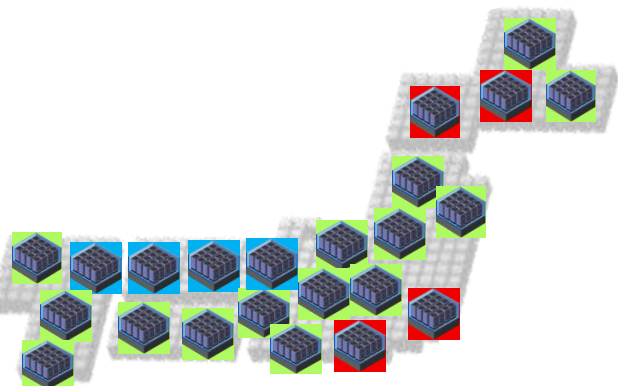


Red Hat Kepler プロジェクト :

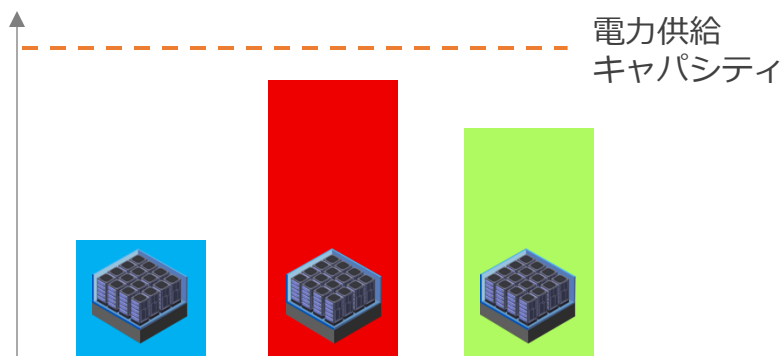
**Red Hat OpenShift上のPodのエネルギー消費を追跡し、
測定および最適化を行うオープンソースの取り組み**

AITRAS オーケストレーター × Red Hat Kepler

電力の需要と供給の
地域差



電力使用率



AITRAS
オーケストレーター

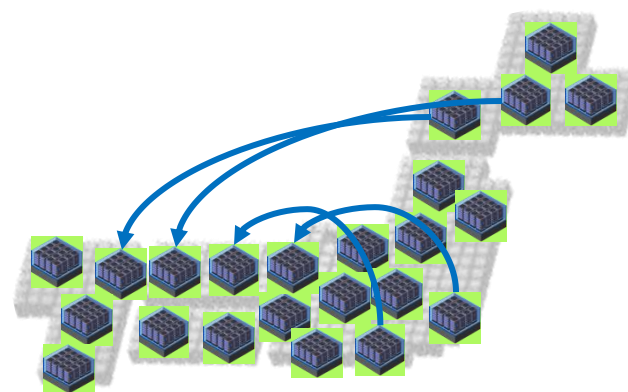
以下を考慮し、最適なDCに
OpenShift上のPodをデプロイ

- ✓ 場所
- ✓ 遅延の要求
- ✓ リソース状況
- ✓ + 電力消費

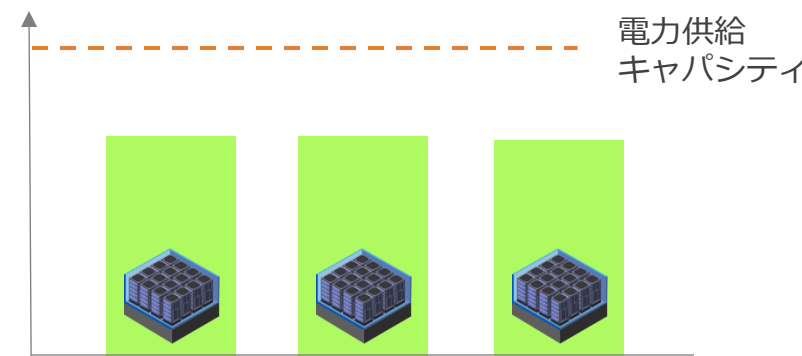
Red Hat
Kepler

Podの電力使用を予測

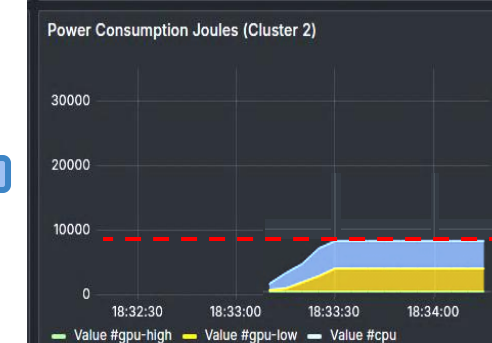
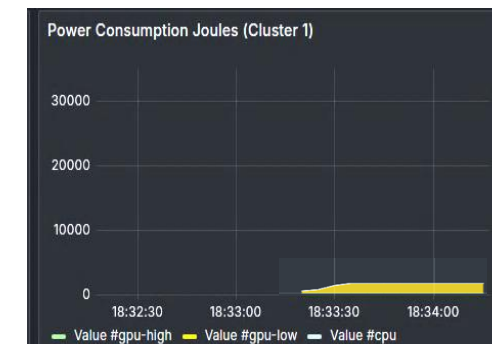
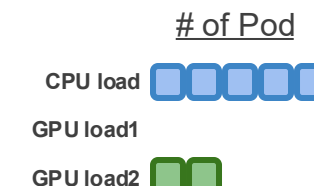
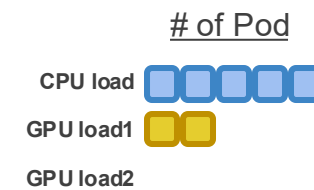
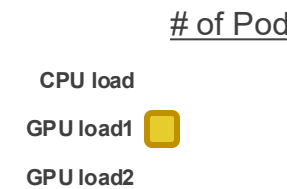
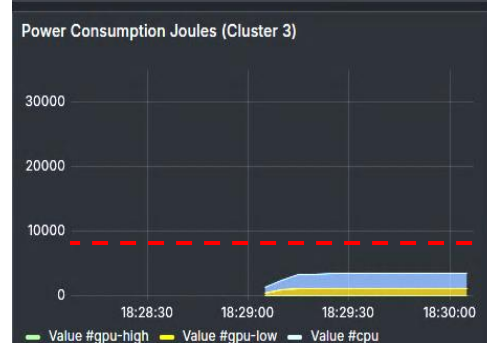
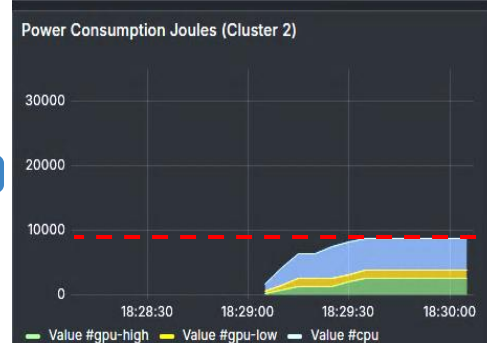
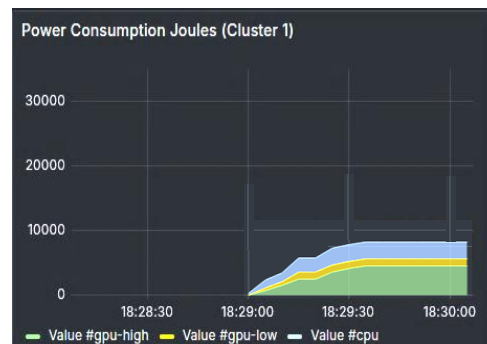
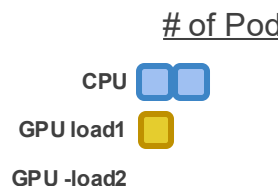
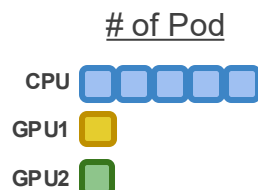
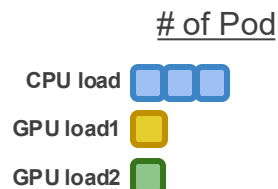
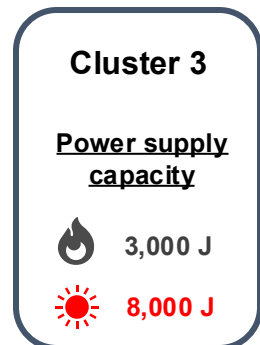
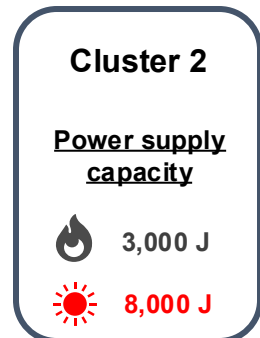
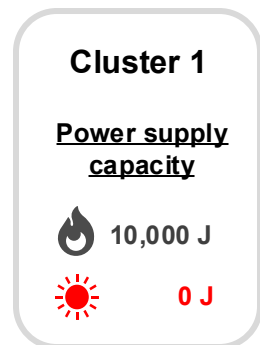
Podの電力使用料を予測し、
配置を最適化



電力使用率



AITRASオーケストレーター x Red Hat Kepler



**再生可能
エネルギーの
使用を最大化**

⑤

ソフトバンクとノキア、
1台のサーバー上で
AIとvRANの共存と最適なリソース割り当ての
自動化を実現

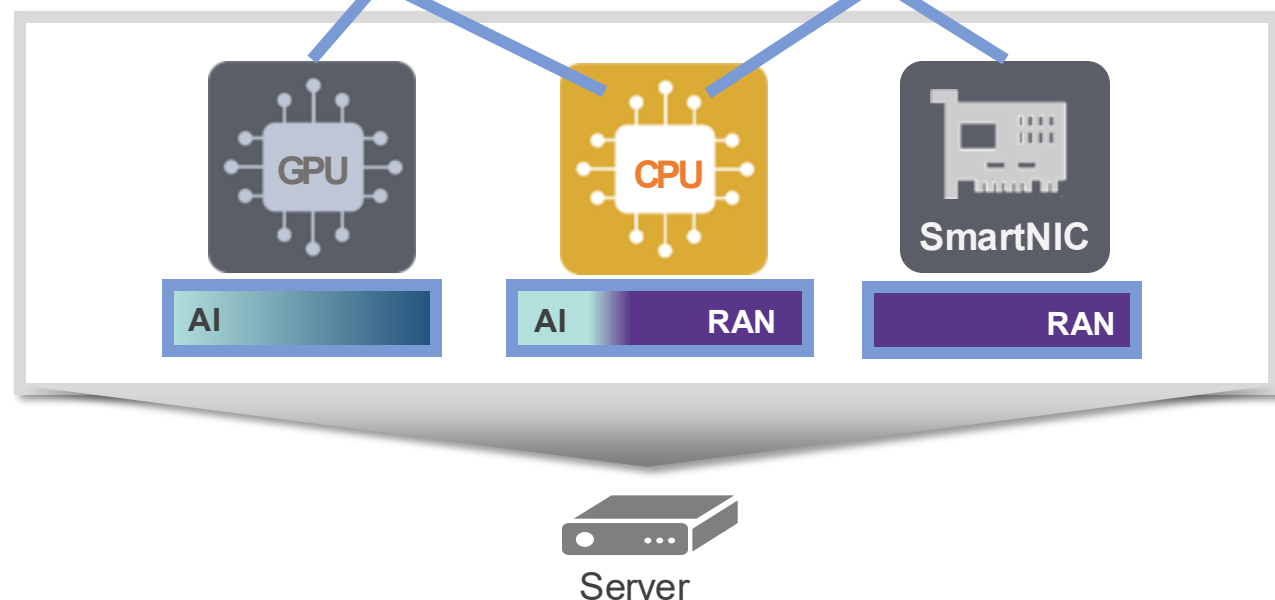
AIとvRANが1台のサーバー上で共存

SoftBank AITRAS オーケストレーター

Nokia
MantaRay NM

SoftBank AI

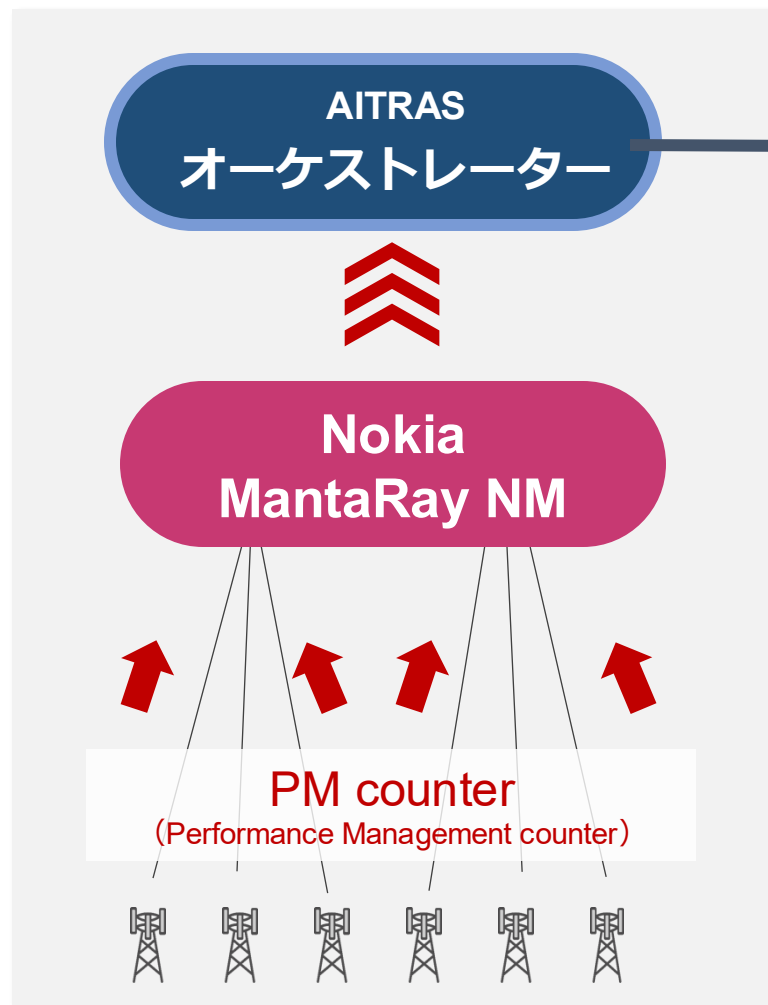
Nokia vRAN



- vDU処理の一部を担当する SmartNICがGPUサーバーをサポートするよう拡張
- AITRASオーケストレーターがリソース配分を最適化

NokiaのvRANソフトウェアを利用してAI and RANを実現

AITRAS オーケストレーター × Nokia MantaRay NM



- vRANトラフィックを予測
- AIとvRANのリソース配分を最適化

