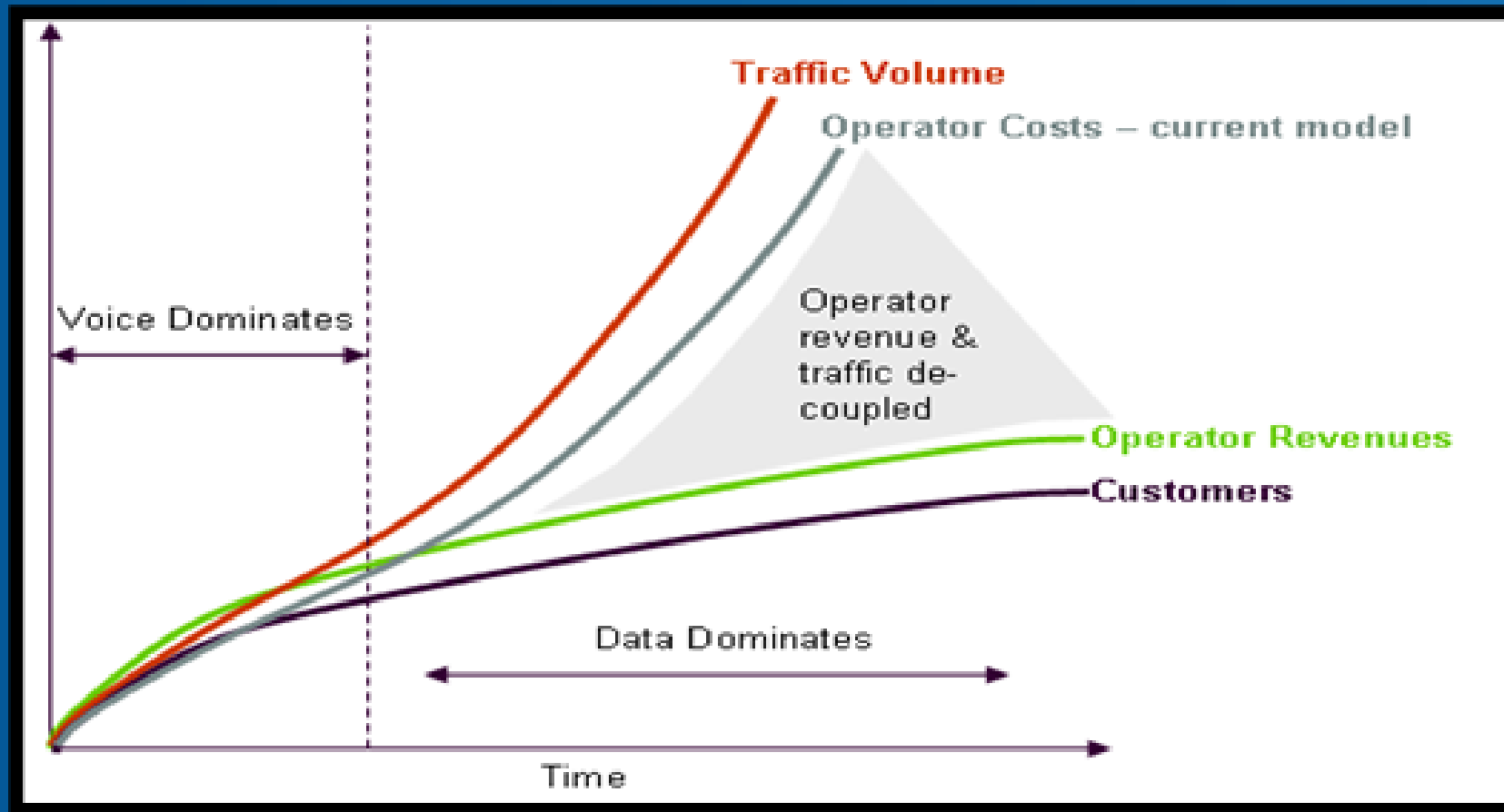# O-RAN
A L L I A N C E

# Confluence of ORAN, AI and Chip

*Alex Jinsung Choi*

Chair of O-RAN ALLIANCE

SVP Deutsche Telekom

1

# Discrepancies between traffic growth and revenue growth (Source: Accenture)



https://telecom.altanai.com/category/access-and-physical-layer/

# Applying Cost Criteria when prioritizing work items

**Cost Accelerators**

New feature addition
Capacity increase/QoS/Latency improving
Techs requiring additional resources and complex control mechanisms, spectrum-RAT dedication, High dimensional signal processing  etc.
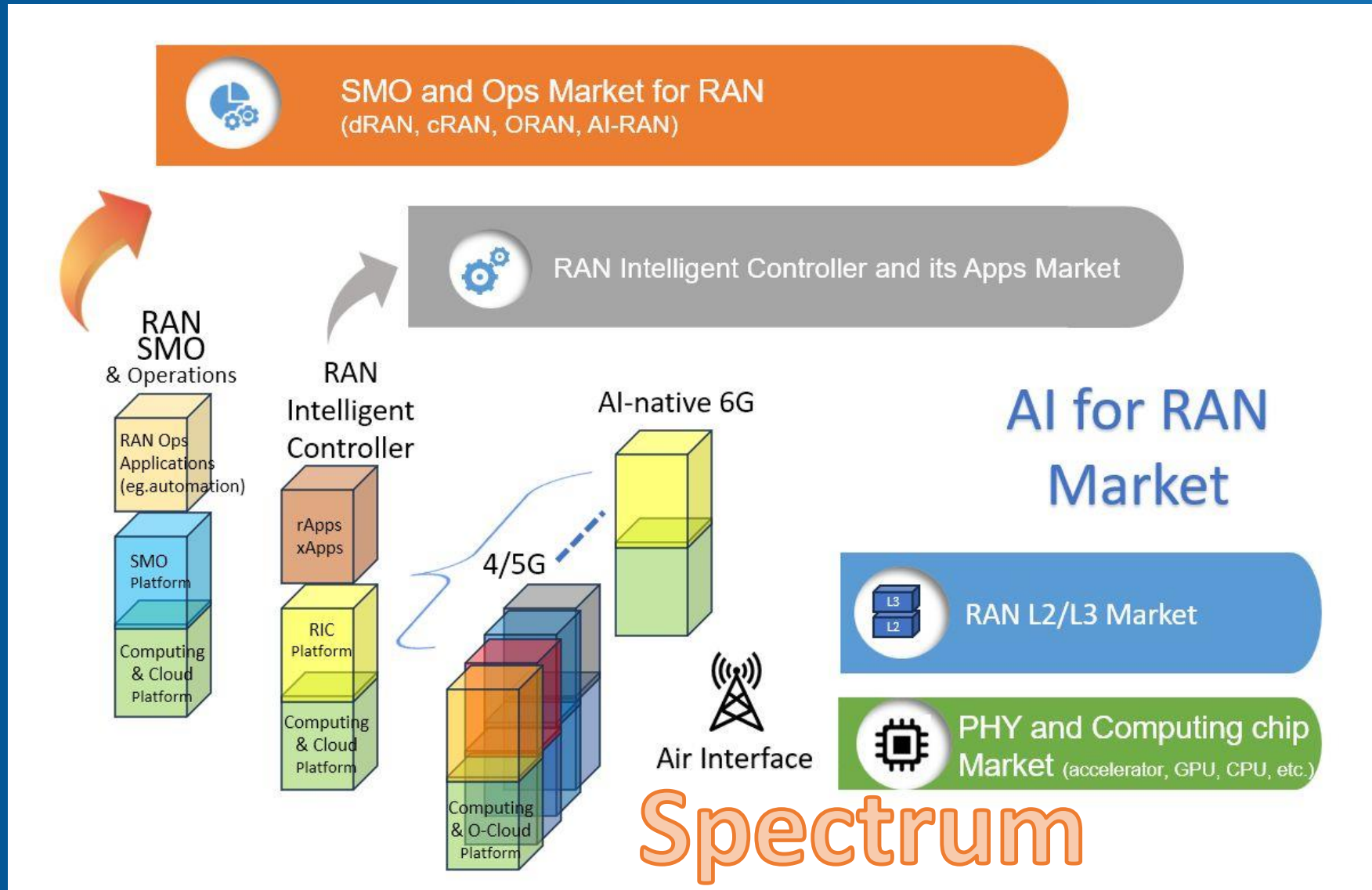
**Cost Optimizers**

Techs for Cost Optimization drivers: commoditization, NG OSS, merchant silicons, simplification, resource sharing, OPEX driven, automation
AIOps, DevOps, etc.

## Cost
Spectrum Cost
RAN Cost
xHAUL Cost
Cloud Infra Cost
SW Licensing Cost
Site Cost
Energy Cost
Operational Cost
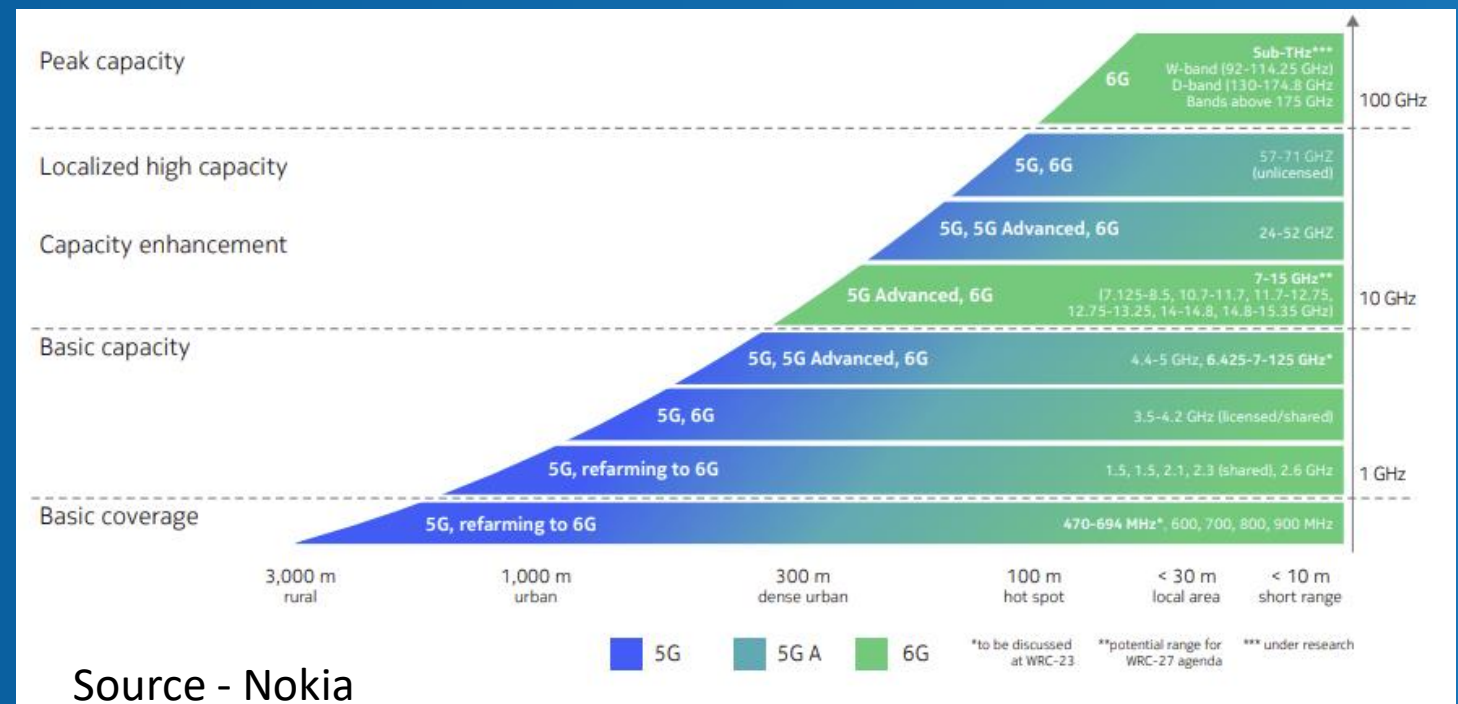Integration Cost
AI/ML infra cost
etc.

# O-RAN Landscape

# RAN - Spectrum Engineering

Spectrum engineering is defined as the discipline within telecommunications focused on the strategic management and technical optimization of radio frequency (RF) spectrum. It involves the analysis, planning, and regulation of the spectrum to ensure efficient, fair, and effective use of this finite resource. Key activities in spectrum engineering include:

- Spectrum Modelling and Analysis

- Interference Management

- Spectrum Sharing and Re-farming

- Coordination and Compliance

- Technical Analysis and Simulation



Source - Nokia

- Support for Emerging Technologies: As new technologies such as Dynamic Spectrum Sharing, MRSS, Advanced Cognitive Radio/Software Defined Radio, Sensing in 5G/5G-Advanced/6G

# O-RAN Accelerator Adaption Layer and Silicon
## why relevant?

# AI for RAN





Note - **Illustration created by adding O-RAN AI/ML to the picture in Jessica Chuang's great LinkedIn post with the title of 'AI in 5G Use Case - Summary of 3GPP's Work & Study on AI/ML for 5G System'**

# Three Types of AI/ML Models for RAN

# AI for RAN use cases including RAN-Ops

# Advanced RIC Apps - Real Time RAN Control

# How to measure success of O-RAN Alliance

# of Big Innovations

Adoption & Scale by Operators (Brownfield & Greenfield)

# of Certificates & Badges

Adoption by SDOs, 6G, ITU

Org Growth and Mix

O-RAN Open RAN Market Share

# of PlugFest Participants

Performance Competence (incl. Cost- Energy Efficiency)

# of Startups

TCO Competence (CAPEX, OPEX)