



Innovations in Clouds,
Internet and Networks

19th
ICIN
CONFERENCE

PARIS
MARCH 1 - 3, 2016

Cloud RAN Challenges and Solutions

Rajeev Agrawal, Anand Bedekar, Troels
Kolding, Vishnu Ram

NOKIA

Tackling future network and service complexity

Operator key drivers for cloud-based radio

- (1) Reduce site rental, energy, and operational costs
- (2) Streamline maintenance operations
- (3) Faster rollout of new features, services, and capabilities
- (4) Lowering costs leveraging the economies of scale of the IT industry

Innovations drive data growth

Data growth drives cloud-based radio

Distributed RAN

Distributed & Centralized RAN

Centralized & Distributed Telco Cloud for RAN, Core, & OSS

Overview of Cloud RAN challenges and solutions

Enable Ethernet Fronthaul

Maximize use of server platforms

RAN software optimized for Cloud operation

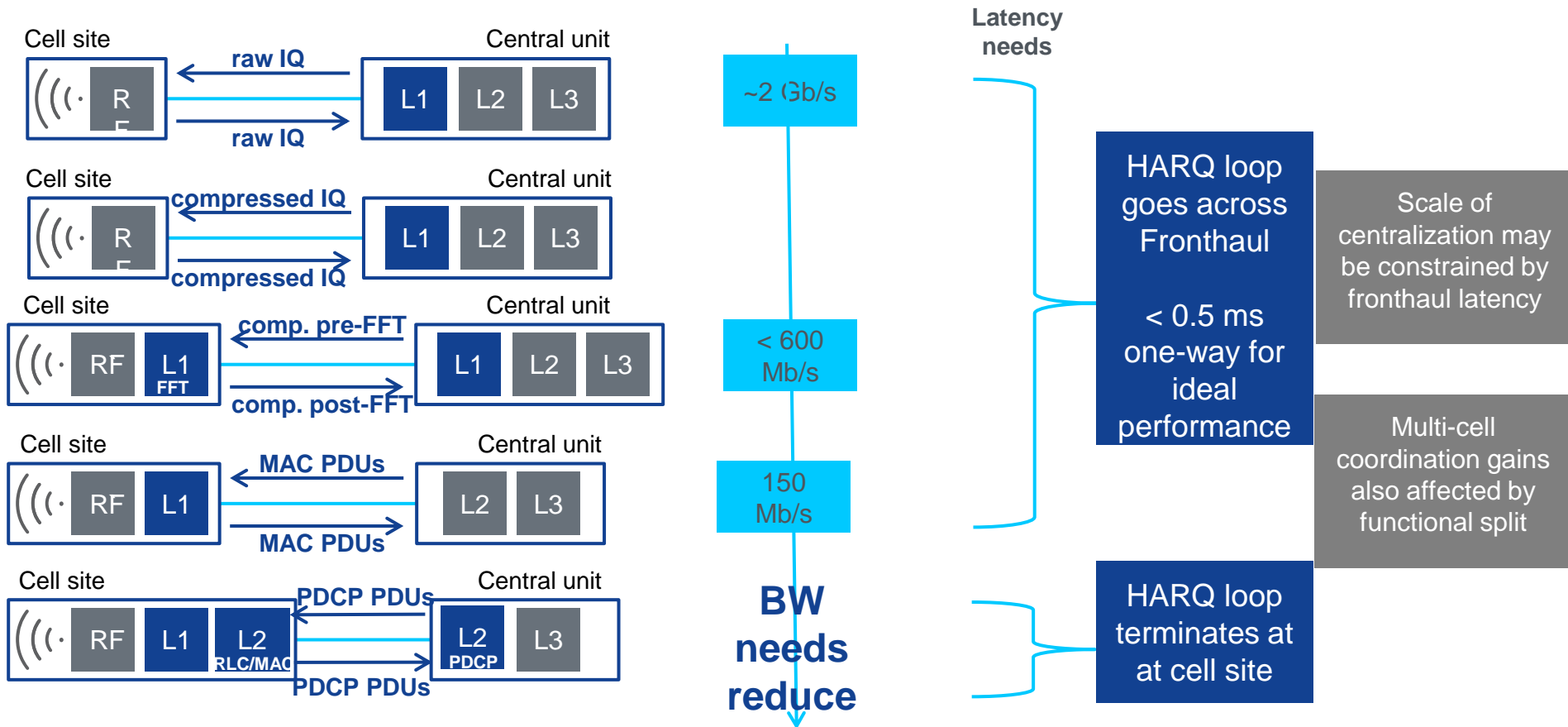
Cloud Management

Multi-cell coordination

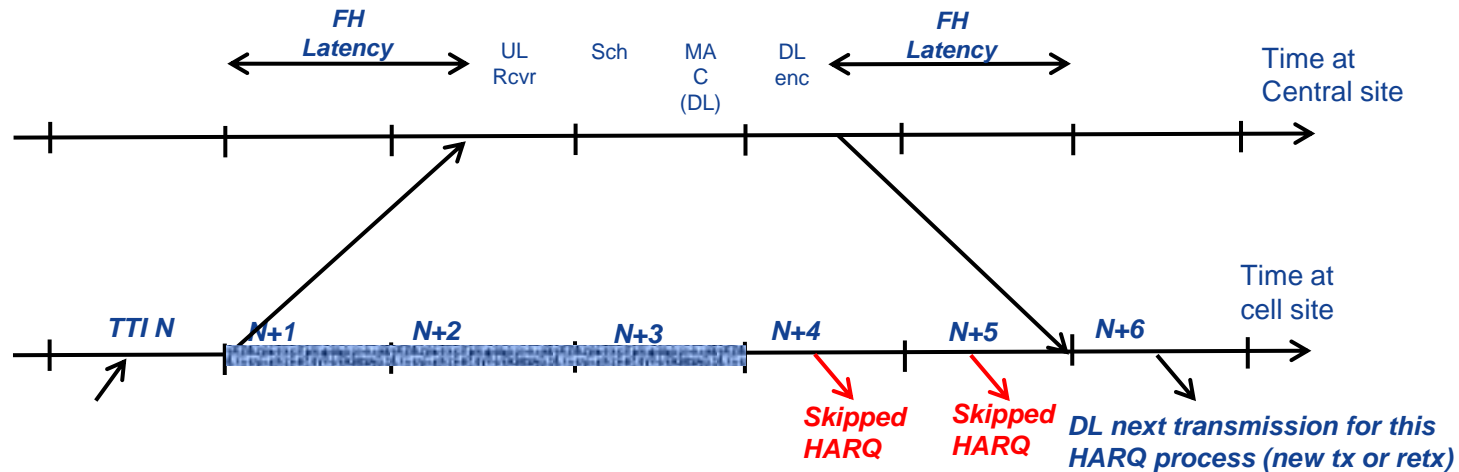
Optimization of latencies

- Different functional splits and algorithm enhancements matched to transport needs
- Improve transport networks
- Enable time-sensitive functions on GPP (x86, ARM)
- Some functions may still be on non-GPP hw
- Disaggregate SW architecture for maximizing pooling benefits
- Load-based elastic scaling
- Enhanced Orchestration and Virtualized Infrastructure Management for RAN
- Common algorithm architecture for distributed and centralized configurations
- Exploit Liquid Clusters
- Co-locating RAN across cells and with EPC enables latency optimizations

Fronthaul needs affected by Functional Splits



HARQ centralization constrains FH latency, Centralization scale



3GPP-standardized HARQ cycle:
UE tx in TTI $n \Rightarrow$ eNB acks, retx in $N+4$

LTE HARQ RTT limits time budget for Front-haul, Processing

Fiber prop delay: 100us per 20km
FH Jitter consumes time budget (dejittering)


FH Latency (and distance) limits scale of centralization of **real-time RAN**:
Likely to have many “**smallish clouds**”
(Non-real-time RAN allows larger scale)

Variability in compute/execution environment eats budget – leaves less for FH

Tight control of jitter in processing time

Multi-cell Coordination with Liquid Clusters

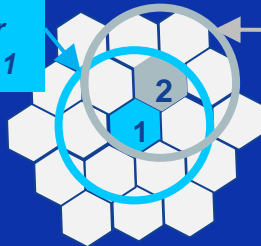
3-cell Intra-site cluster



Intra-site Cluster

We recommend harvesting intra-eNB coordination gains before inter-eNB coordinated scheduling

Liquid Cluster of Cell 1

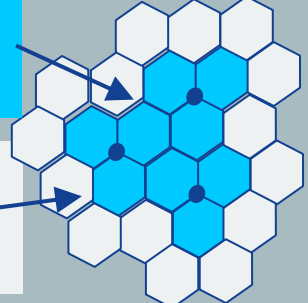


Liquid Cluster of Cell 2

Liquid Cluster

Coordination with right set of neighbors
Decentralized coordination is best suited for this

9-cell non-overlapping cluster



Non-Overlapping Cluster

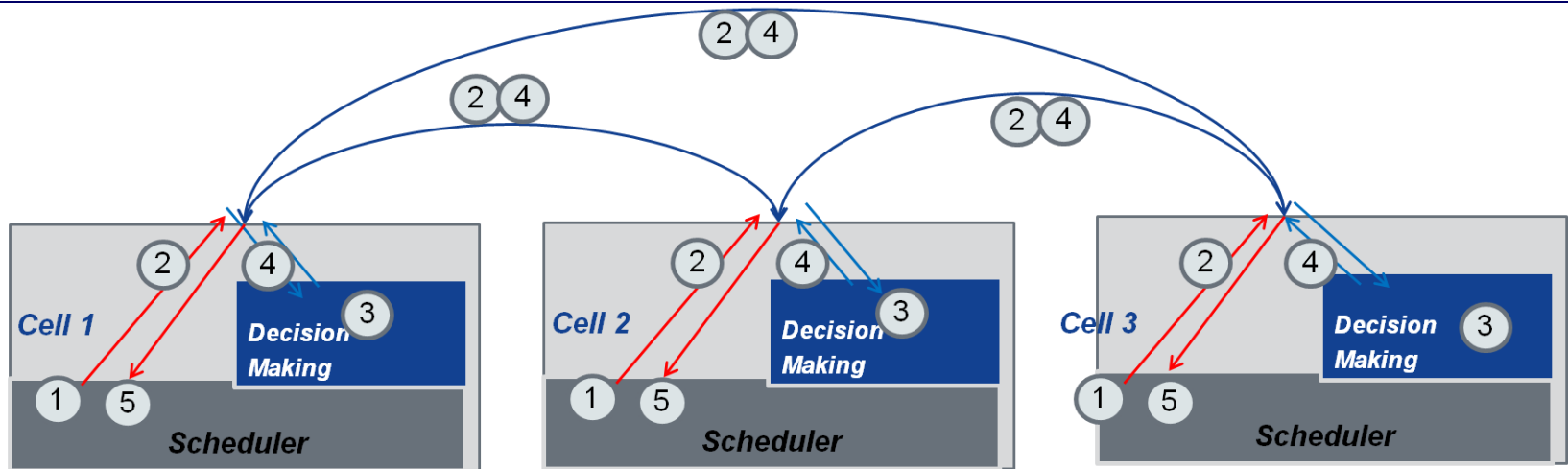
The 9 cells in the cluster cannot coordinate with these adjacent cells

No coordination benefits at cluster borders
Typically used by centralized coordination

Maximize cooperation gains – no hard cluster boundaries
Applies in a variety of coordination algorithms – Coordinated Scheduling, Dynamic Point Selection, UL Joint Reception CoMP, ...

Liquid Clusters enable flexible configuration of cooperating cells

Decentralized Algorithm Architecture for Multi-Cell Coordination

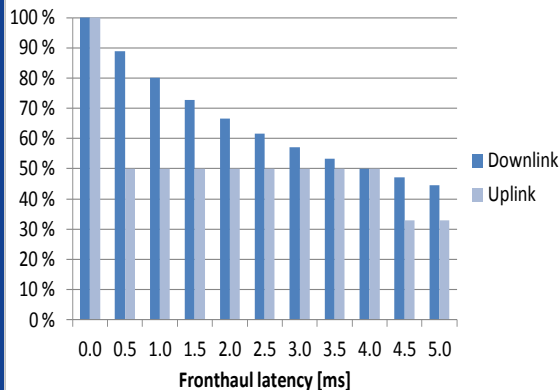


Decentralized algorithm architecture for Multi-cell coordination: can achieve equal or better performance than centralized algorithms

When scheduler (or L1) is centralized:
Within central site, coordination among distributed processors

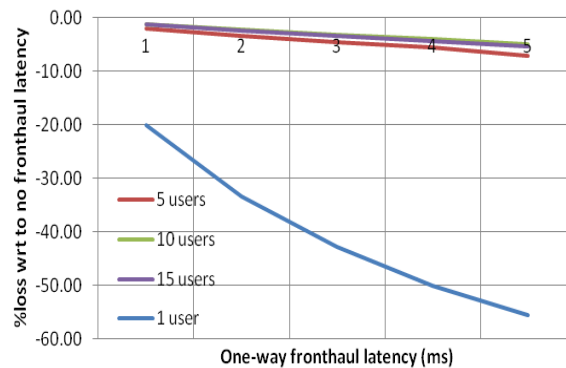
Liquid Clusters and Decentralized algorithm architecture enable flexible mapping of functionality within physical cloud architecture

Single user throughput

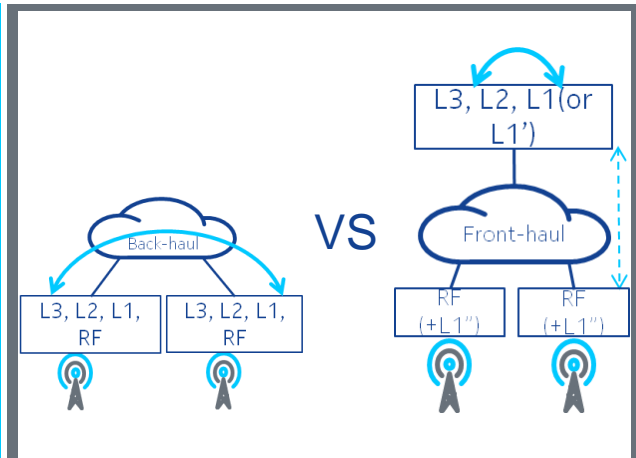


Single-user throughput:
Significant impact –
UL worse than DL

Degradation in Average Tput (MuCCS OFF)



DL Multi-user scenario:
At >5 UEs, recovers system
throughput
UL needs Sch enhancements



For a given latency (FH or BH), Distributed coordination will outperform coordination in central site

Virtualization of Real-time RAN Functions

L1 functions currently considered more efficient on DSP/SoC/FPGA hardware
For some L1 functions, special hw acceleration may be required (e.g. turbo decoder)

L1 may require non-GPP hardware

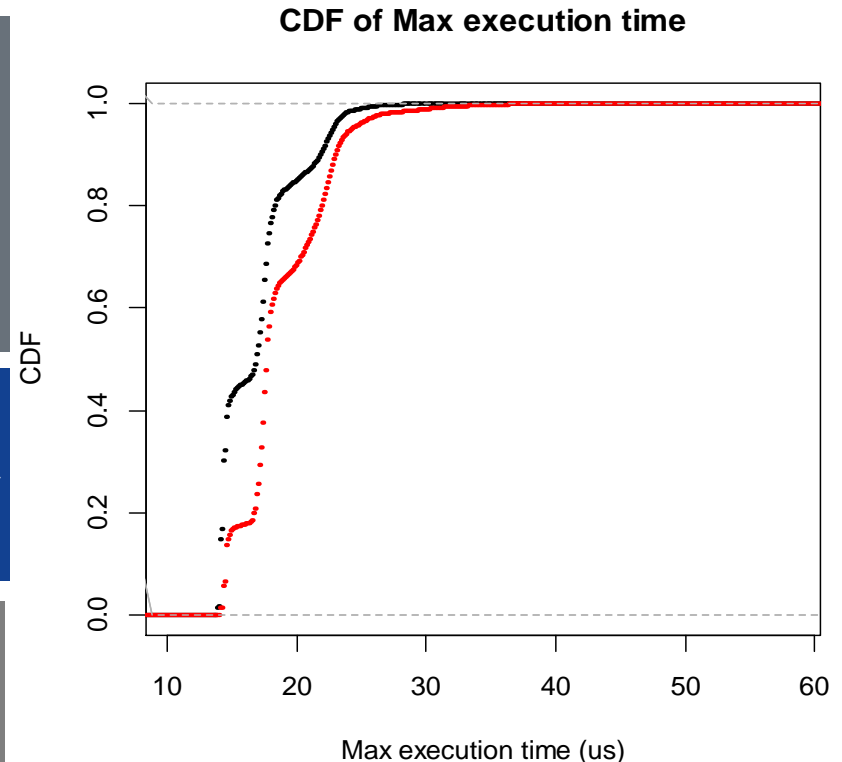
L2, Scheduler can be virtualized on GPP/Linux

“Real-time” functions L2, Scheduler can be virtualized

Methods of controlling execution jitter:

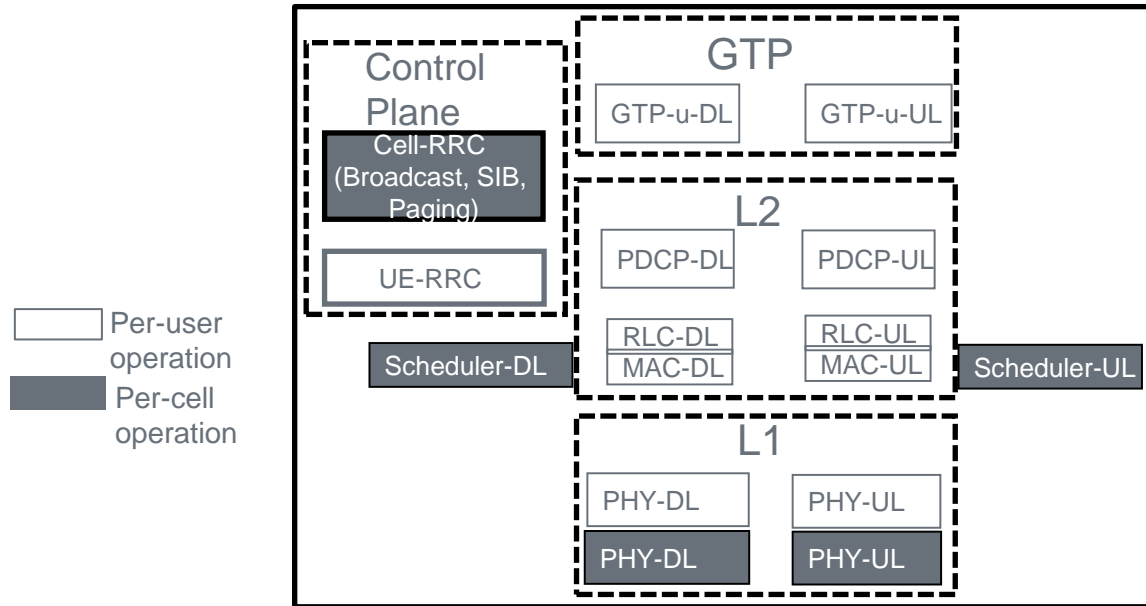
Poll-driven task dispatch, rather than interrupts
Careful mgmt. of host OS and BIOS interrupt handling
CPU core isolation and pinning, etc.
Both guest as well as host OS need careful tuning

Real-time virtualization requires control of execution jitter



Execution jitter of <30us shown to be achievable on x86/Linux
→ Good for LTE where TTI=1ms:
5G may need even tighter control due to shorter TTIs

SW Architecture – How to achieve Pooling for RAN functions



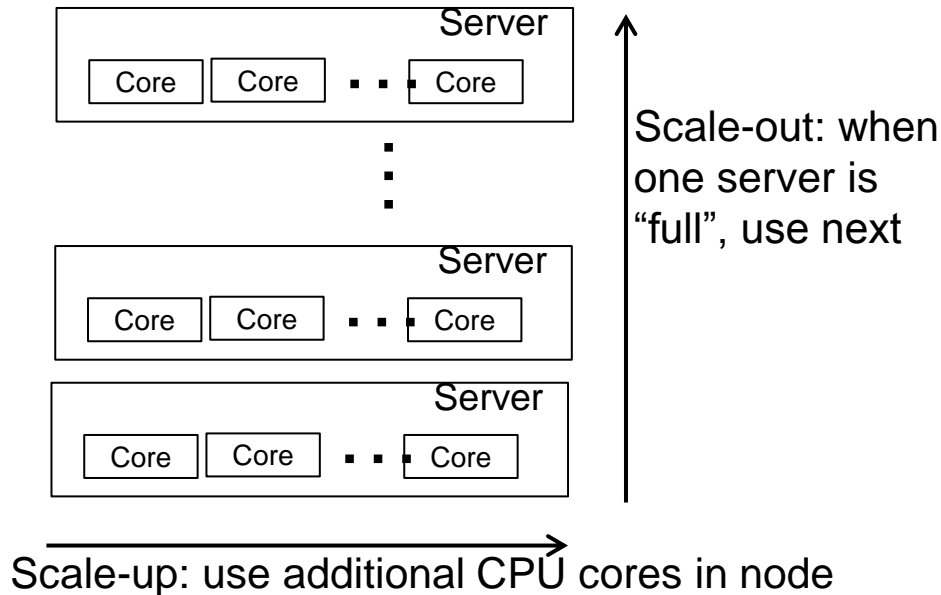
Certain functions have *per-user operation*, while others have *per-cell operation*

Per-user operation: Compute requirements scale with the number of users

Per-cell operation: Compute requirements scale with number of cells

Disaggregate and Reassemble RAN functions: provide the right form of instantiation for each function

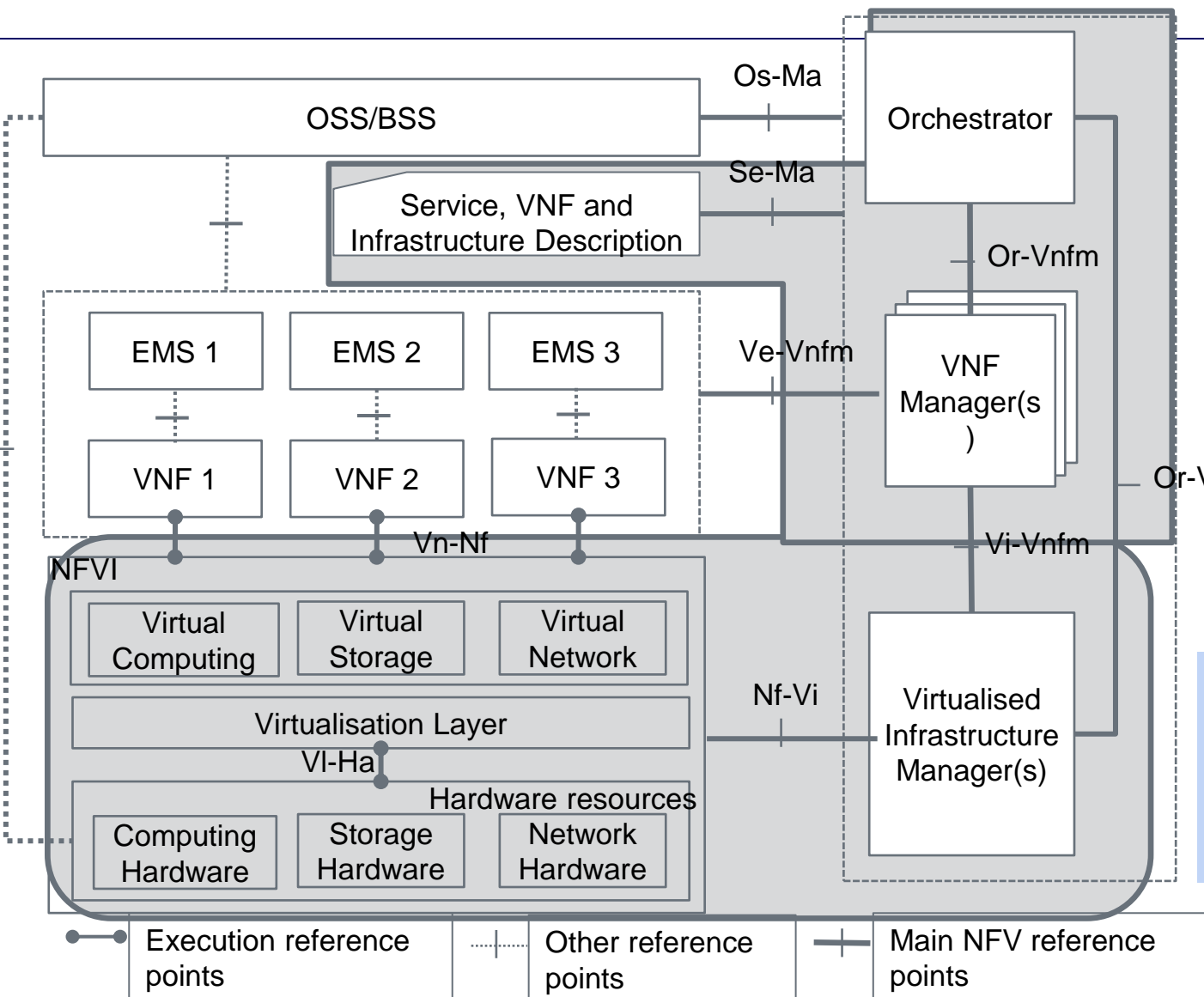
For Per-user operations, can do pooling with *horizontal aggregation* – any processing unit (e.g. VM) can handle any user from any cell



Elasticity: Dynamically change amount of processing resources consumed

Support **Scale-up** (cores within a server/VM) as well as **Scale-out** (number of servers/VMs)

Algorithm to automatically adapt number of CPUs based on load, transport latency/jitter, processor type, real-time deadlines, etc.



Consistent with NFV reference architecture

But RAN has special needs:

Functionality at guest, host OS for real-time operation
Mix of GPP, non-GPP HW
Elastic Scaling needs special cooperation b/w guest and host OS

Current tools e.g. Openstack need extensions to meet RAN needs

Maximize coordination gains,
flexible placement of coordinating
functions

Multi-cell coordination with Decentralized Algorithm Architecture and Liquid Clusters

Choice of functional split is a
complex Tradeoff BW,
latency/jitter, multi-cell
coordination, cost/complexity

Functional splits enable
fronthaul flexibility:
Multi-dimensional tradeoff

L1 may still need non-GPP
HW, accelerators
L2 and above can be
virtualized: but requires
control of execution jitter

Maximize use of GPP/Linux
Some functions need non-
GPP

Disassemble and Reassemble RAN
functions for optimal pooling
Elastic Scaling for automatic adaptation to
changing loads, real-time operation KPIs

Pooling: Exploit characteristics
of Per-user and Per-cell ops

Mix of GPP and non-GPP hardware in
compute infrastructure
Support operation and scaling of real-
time RAN ops on GPP

Cloud Management and
Orchestration may need special
support for RAN



Innovations in Clouds,
Internet and Networks

19th
ICIN
CONFERENCE

PARIS
MARCH 1 - 3, 2016

Thank you!

#ICIN2016

NOKIA