

Innovating at the Edge Meetups

# Edge Big Bang

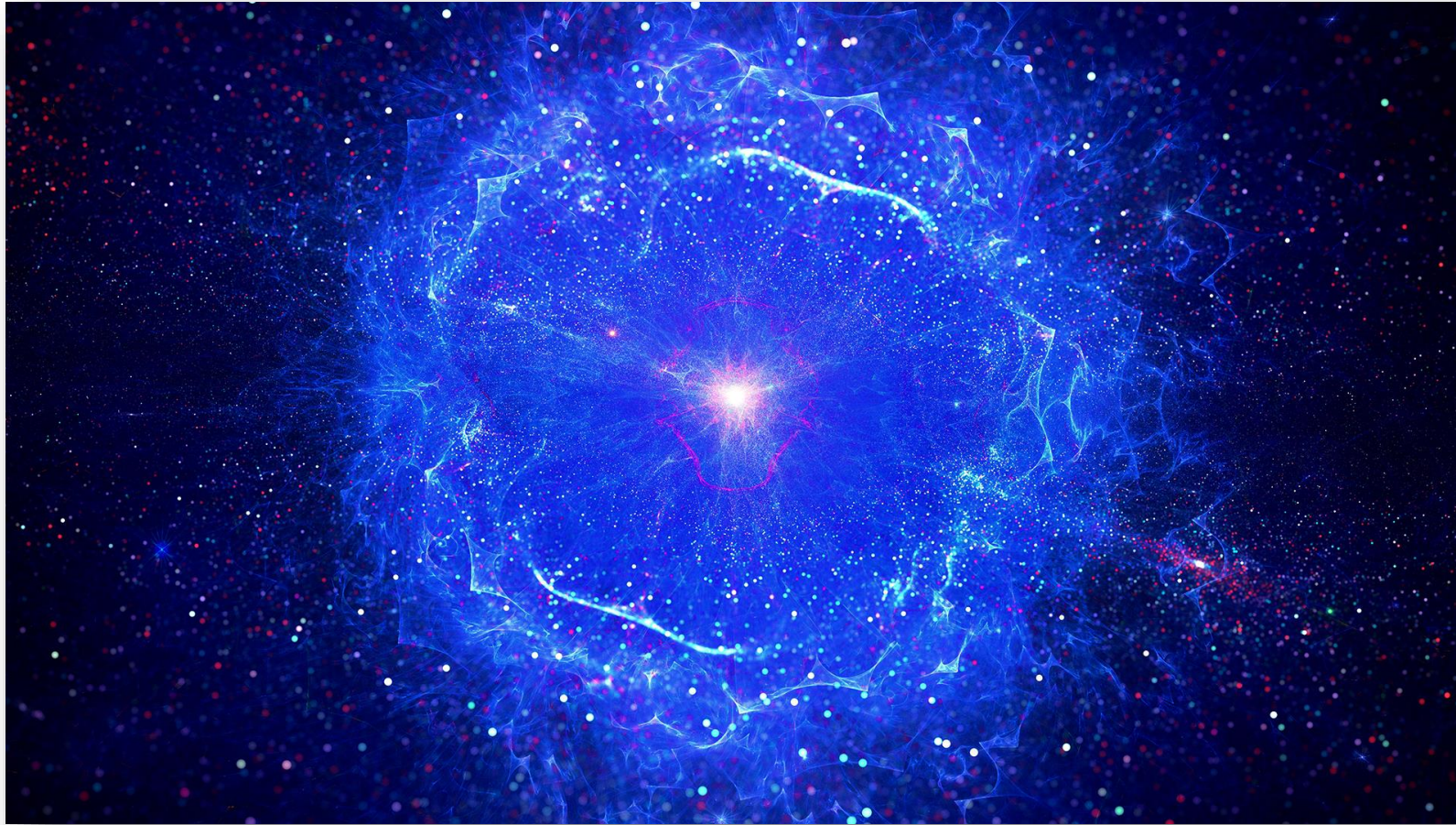
- Francesc Guim, Chief of Edge Architecture, Intel

The Intel logo is located in the bottom left corner of the slide. It consists of the word "intel" in a white, lowercase, sans-serif font, with a registered trademark symbol (®) to its upper right. The logo is positioned over a dark blue background that features a decorative graphic of several overlapping squares in various shades of blue and white, arranged in a stepped, staircase-like pattern.

intel<sup>®</sup>

# Notices and Disclaimers

- Intel technologies may require enabled hardware, software or service activation.
- No product or component can be absolutely secure.
- Your costs and results may vary.
- © Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.



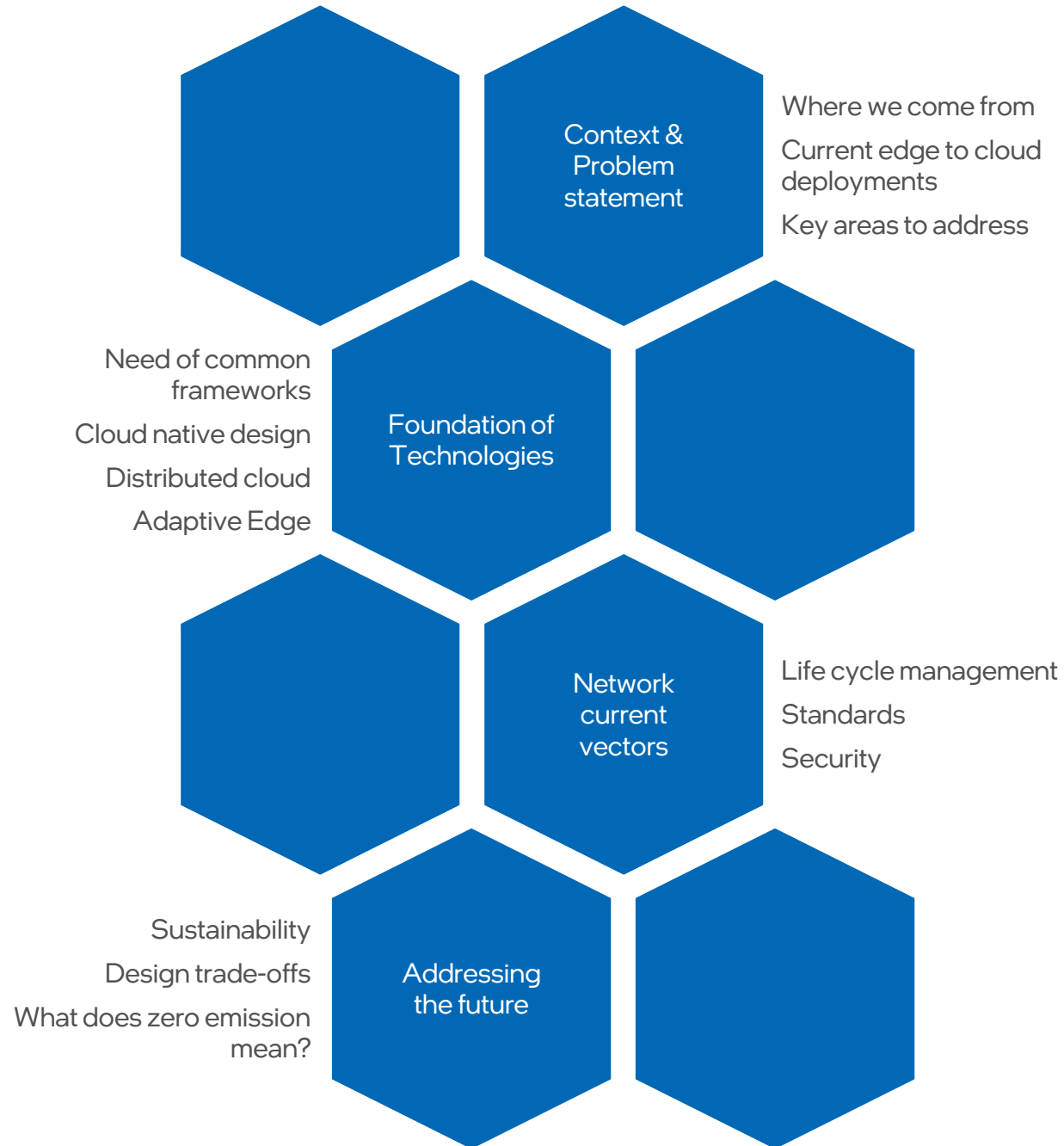
POCs  
On-Premise or Controlled env.  
10/100s Edges per deployment

Production deployments  
Far Edge, Near Edge, Data Center Edge, Cloud  
Up to 100K Edges per deployment  
Challenges: Temperature, Energy, Reliability and security @ scale

2016

2023

# Outline



# Outline



# Challenge 1: Vertical View of Edge

END USERS & DEVICE LOCATION +  
TRANSPORT TYPE

- STREET USERS



Telco Infrastructure

- VEHICLES
- STREET CAMERAS
- STREET SENSORS
- ...



IOT Connectivity & Services

- RETAIL SHOPS
- PUBLIC LOCATIONS (I.E: LIBRARIES)
- PRIVATE ENTERPRISE
- ...

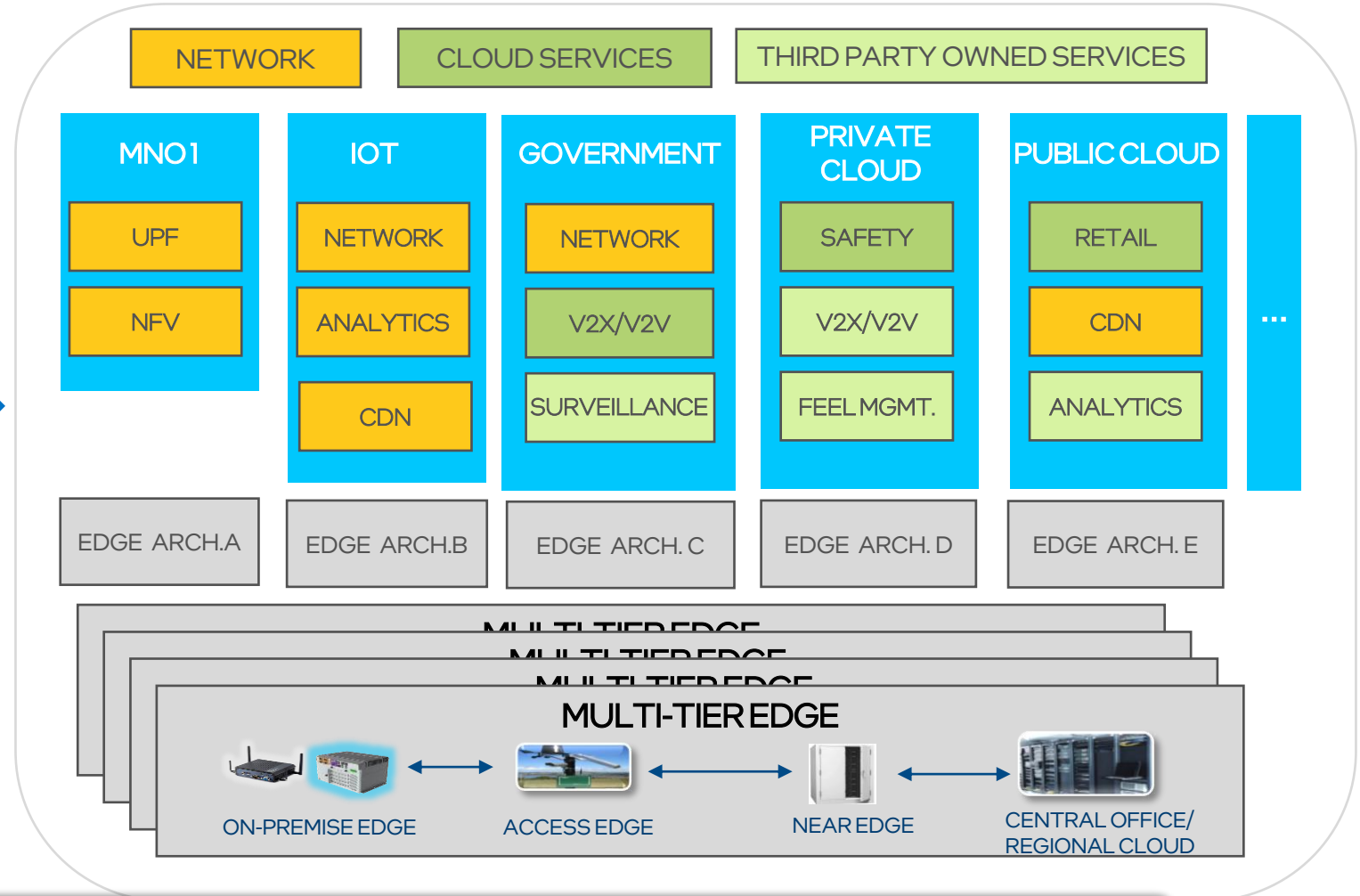


Enterprise Connectivity & Services

TELCO  
5G/LTE

PRIVATE 5G/LTE, WIFI,  
WIRED, TETRA, LORA ETC.

VERTICALIZED AND NON-CONSOLIDATED ARCHITECTURE



**Challenges:** CSPs/CoSPs have multiple edge verticals (IOT, Enterprise, Telco) with differing architectures. How do we drive a scalable & converged architecture?

# Challenge 2: Edge is Distributed, has Life and Requires Scales

DEPLOYMENT TIERS	<b>ON-PREMISE EDGE</b>		<b>NETWORK EDGE</b>	<b>DC EDGE</b>	<b>CORE DC</b>	<b>PUBLIC CLOUD</b>			
	<b>INTELLIGENT SENSOR/GW</b>	<b>INTELLIGENT EDGE</b>							
DEPLOYMENT LIMITATIONS	AVAILABLE POWER: [5,60W] THERMALS: -40 to +85C	AVAILABLE POWER: [100,3000W] THERMALS: -40 to +70C	AVAILABLE POWER: [400,3000W] THERMALS: -20 to +70C	AVAILABLE POWER: 9KW/Rack THERMALS: 20 to 40C	TYPICAL AVAILABLE POWER: GEO DEPENDENT 7KW – +20KW/RACK				
FORM FACTORS	 Ruggedized/ Custom FF	 Ruggedized/ Custom FF	 Server/FF optimized for extreme conditions	 Short depth server Chassis	 Standard 1U/2U Chassis	 CSP Custom FF			
PROCESSOR TYPE	 ENTRY ATOM	 ENTRY CLIENT	 XEON-D(SoC)	 ENTRY XEON	 ENTRY ATOM	 XEON-D(SoC)	 ENTRY XEON	 XEON FAMILY	
ACCELERATION	 VPU	 FPGA LOW POWER	 INTEL eASIC    INTEL STRATIX    INTEL FPGA    GPU (<75W)				 VPU	 FPGA	 GPU (<150W)
STORAGE/MEMORY		 HIGH PERF SSD		 HIGH PERF SSD		 PERSISTENT MEMORY	 HIGH PERF SSD		

# Challenge 3: Amount of Technology Building Blocks

	End-to-End Services								
	AI for Imaging	Content Distribution	VNFs	NLP	KVS	AR/VR			
Service Orchestration	K8s/K3 based Local/Multi-Edge Orchestrator		Hybrid-cloud / edge		Cloud Service Based Orchestration				
SW Cloud Service Provides Solutions	AWS (outpost, sage maker ...)		Azure (ARC, stack edge, HCI ..)		Google (functions, ..)				
VIM	KVM	OpenStack	ACRN	Containerization tech (VM or bare metal)					
Infrastructure	Infrastructure Telemetry			SD-WAN					
Network Infra	NIC	Smart NIC / IPU	Ethernet Switch	CXL switching	NZTP	iPXE			
Frameworks	EdgeX	OpenVINO™	Intel® Smart Edge	Edge Insights for Industrial	Edge Control Stack				
Managed Network Services	CUPS vEPC User Plane (PGW-U / SGW-U)		CUPS vEPC Control Plane (MME / HSS / SGW-C / PGW-C)		5G FWA – Fixed Wireless Access				
	Full vEPC (MME / HSS / PGW / SGW)		vBNG	FMC – AGF - FMIF	Ethernet / MPLS	vRAN			
Access	LTE / 5G NR (ORAN/FlexRan)	Fiber	WiFi / Industrial WiFi	LoRA	Ethernet	IIoT protocols			
Form Factors	Smart Camera (VPU)	No AC Cooling	AC/Liquid Cooling	Industrial Fanless	AI Appliance	Edge Servers	Edge DC Servers	Data center Servers	
SW Platform	Security (Isecl)		Quality of Service	Acceleration	Management	Telemetry			
xPUs	GPU	AI Accelerators	FPGA(s)	Low power CPU (Atom)	Client CPU (i3, i5, i7)	IPU	Xeon D	Xeon E	Xeon SP
Platform and xPU Features	Data protection	Workload Isolation	Edge Attestation	Edge Quality of	AMT	TSN, TCC			



# How Do We Architect Edge To Cloud?

- At scale
- With security
- With an optimized total cost of ownership
- With compute democratization
- Satisfying edge use cases requirements
- Being sustainable and scalable at the long run

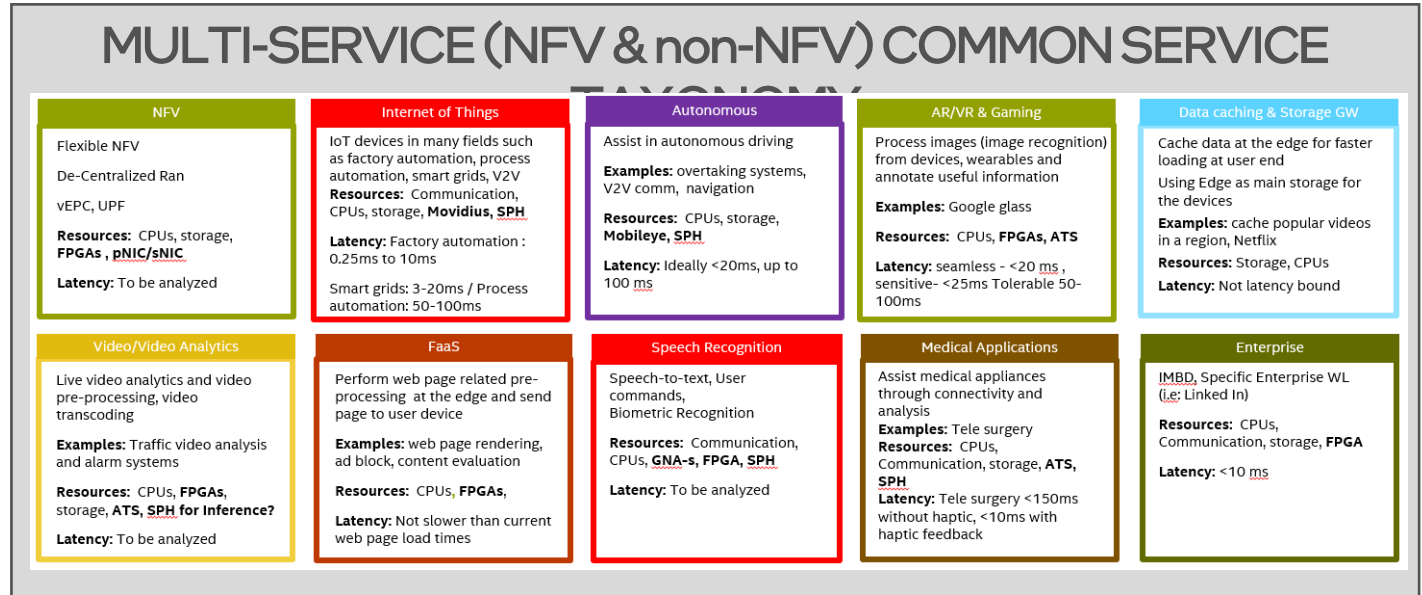
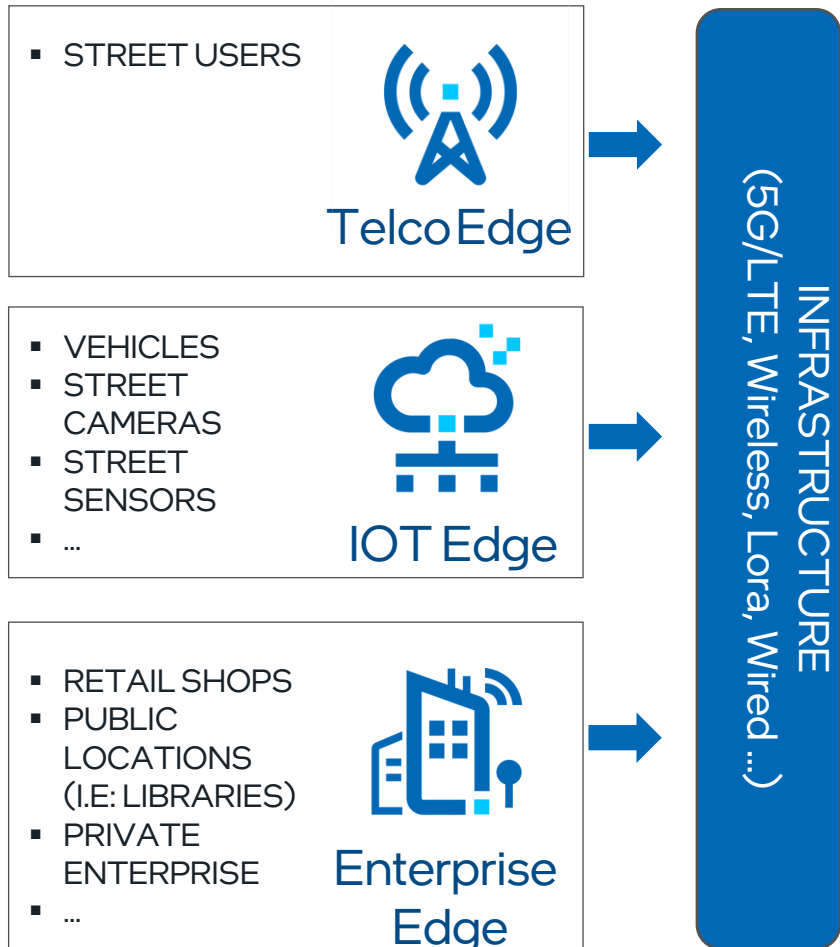
# Outline



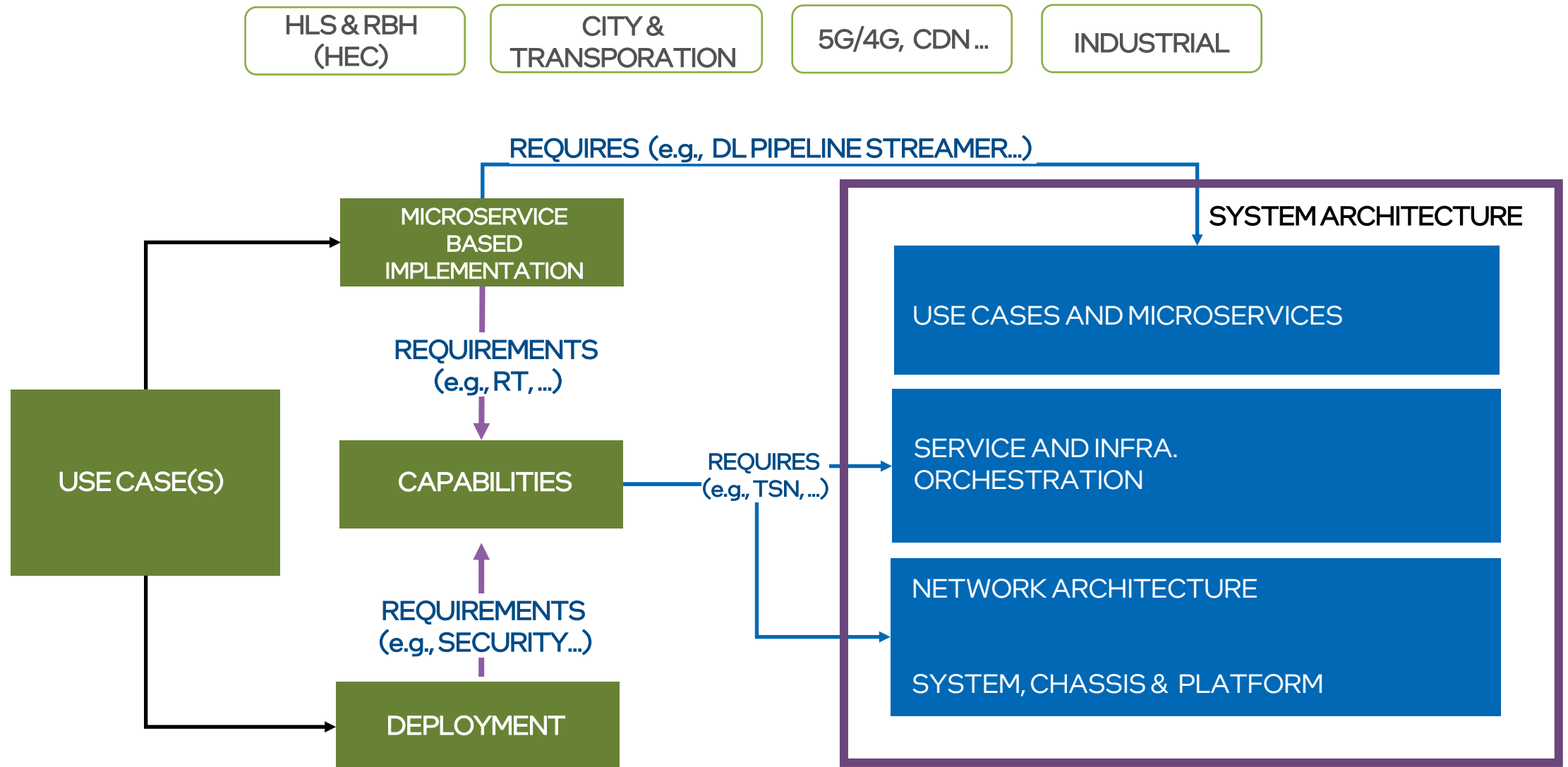
# Common Framework to Address Different Verticals

EDGE END TO END ARCHITECTURAL FOUNDATION: INTER-OPERABLE, OPEN & SELF MANAGED

## END USERS & DEVICE LOCATION + TRANSPORT TYPE



# Use Case and End Customer Driven Architecture



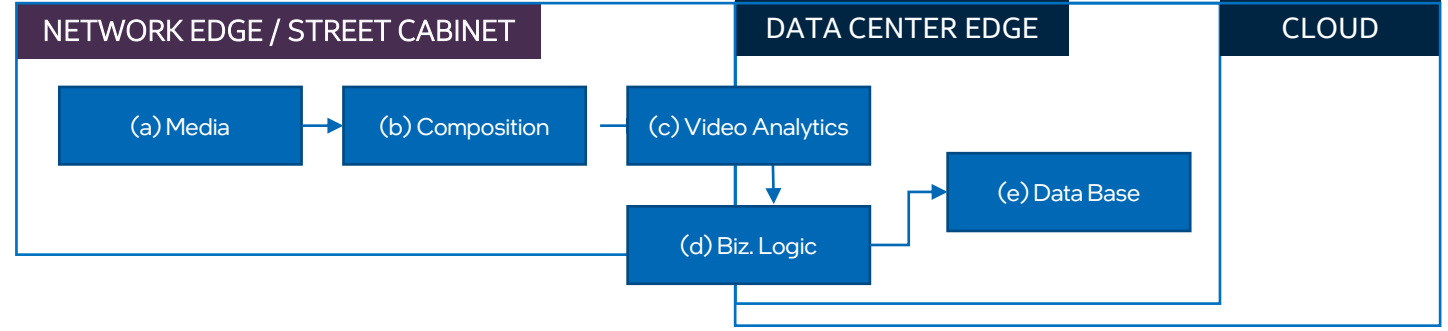
Tiers	On-Premise Edge				Network Edge	DC Edge	Public Cloud			
	Intelligent Sensor/GW		Intelligent Edge							
Network Latencies (Wire Round trip)	$< 1\text{ms}$		$< 1\text{ms}$		1-5ms	5ms + 1-2 ms (every 100kms)	5ms + 5 ms (every 400kms)			
Deployment Requirements	Compute Available Power: $< 50\text{ W}$ Form Factor: Small Box Thermals: NEBs Mgmt.: Remote		Compute Available Power: $\sim 10\text{KW}$ Form Factor: Rack(s) Thermals: NEBS or Standard DC Mgmt.: Remote		Compute Available Power: $< 600\text{ W}$ Form Factor: Pizza box Thermals: NEBS Mgmt.: Remote	Compute Available Power: $9\text{KW}/\text{rack} / 1\text{KW sqm}$ Form Factor: Rack(s) Thermals: NEBS or Standard DC Mgmt.: Remote	Standard Data Center (DC)			
Where, What & Why	Use Case	KPI	Use Case	KPI	Use Case	KPI	Use Case	KPI	Use Case	KPI
	Intelligent Transportation	Data Privacy, Backhaul Traffic Savings, Reliability	AR/VR	Latency, Backhaul Traffic Savings, Network scalability	Intelligent Transportation	Data Privacy, Backhaul Traffic Savings, Reliability throughput Latency	Intelligent Transportation	Data Privacy	CDN	Backhaul traffic savings, Throughput
	V2V	Latency					Drone/IoT	Same as Int. Transp.		
	Retail	Same as Int. Transp.	Retail	Data Privacy, Backhaul Traffic Savings, Reliability	V2V	Latency	Healthcare	Access to services		
	Video Analytics	Same as Int. Transp.			Video Analytics	Same as Int. Transp.	CDN & Storage	Backhaul Traffic Savings Throughput		
			RT Streaming Healthcare	Same as AR/VR Access to services	Drone/IoT	Same as Int. Transp.	Storage GW	Same as CDN		
					Rural	Access to services			FaaS	Latency
									AR/VR/MR	Latency

# Requirements to Distributed System Architecture

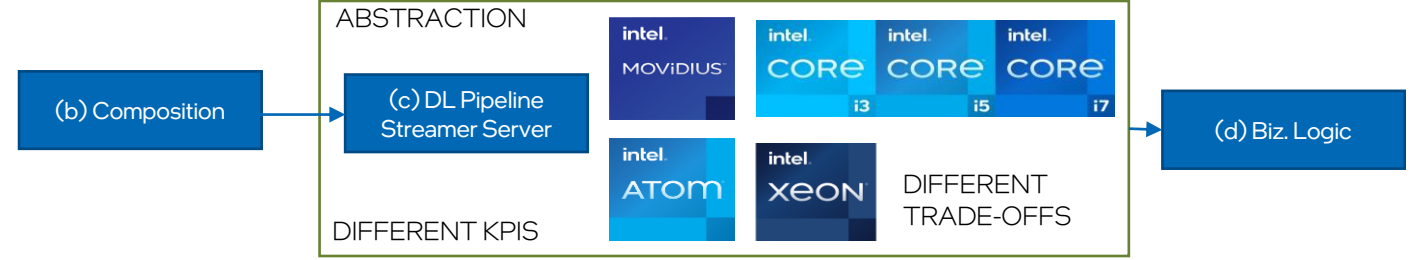


- Use case decomposition
- Intel HW Architecture Positioning w.r.t (1)
- Platform Architecture
- Edge Appliances (Rack, PODs, Platform etc..)

(1) DECOMPOSE USE CASE IN MICROSERVICES AND UNDERSTAND DEPLOYMENT REQUIREMENTS/OPTIONS



(2) UNDERSTAND HOW MICROSERVICES CAN ALLOW ADOPTION AND STICKINES

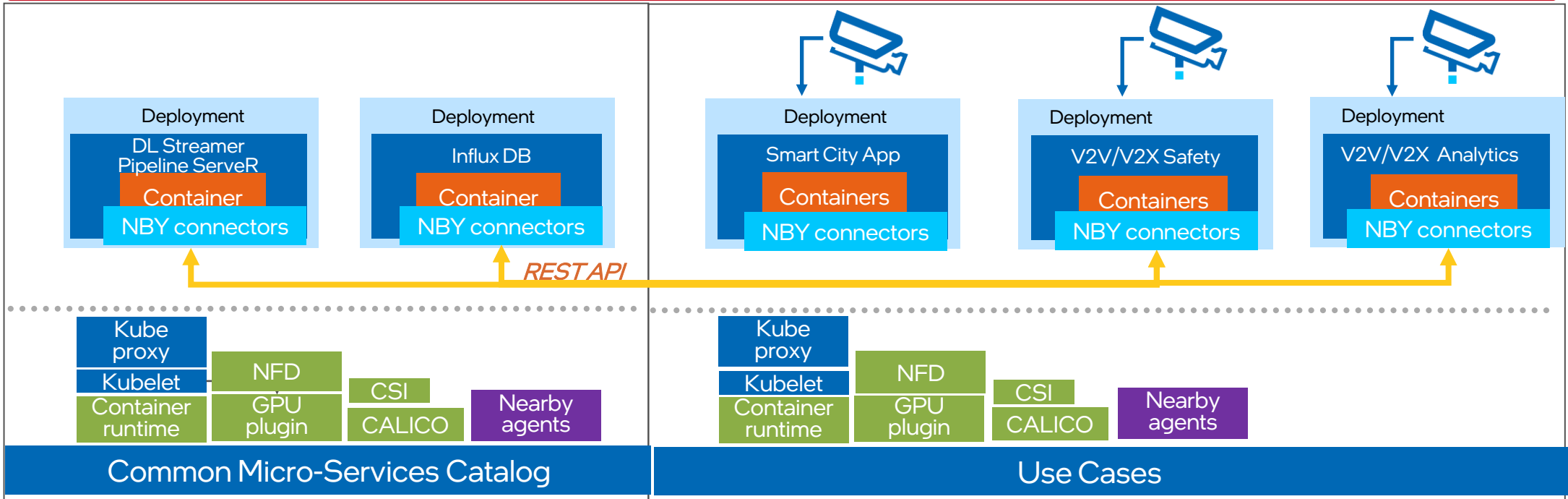


(3) TRANSLATE (1) AND (2) INTO SW/HW SYSTEM ARCHITECTURE



# End-to-End Edge Orchestration – Nearby One

Control Plane  
Workloads  
Worker Nodes  
Edge HW Architecture



### Distributed Edge Systems

Cellnex Zero Emission Site      Micropod Data Center Edge

### Edge Platforms

**Lenovo**

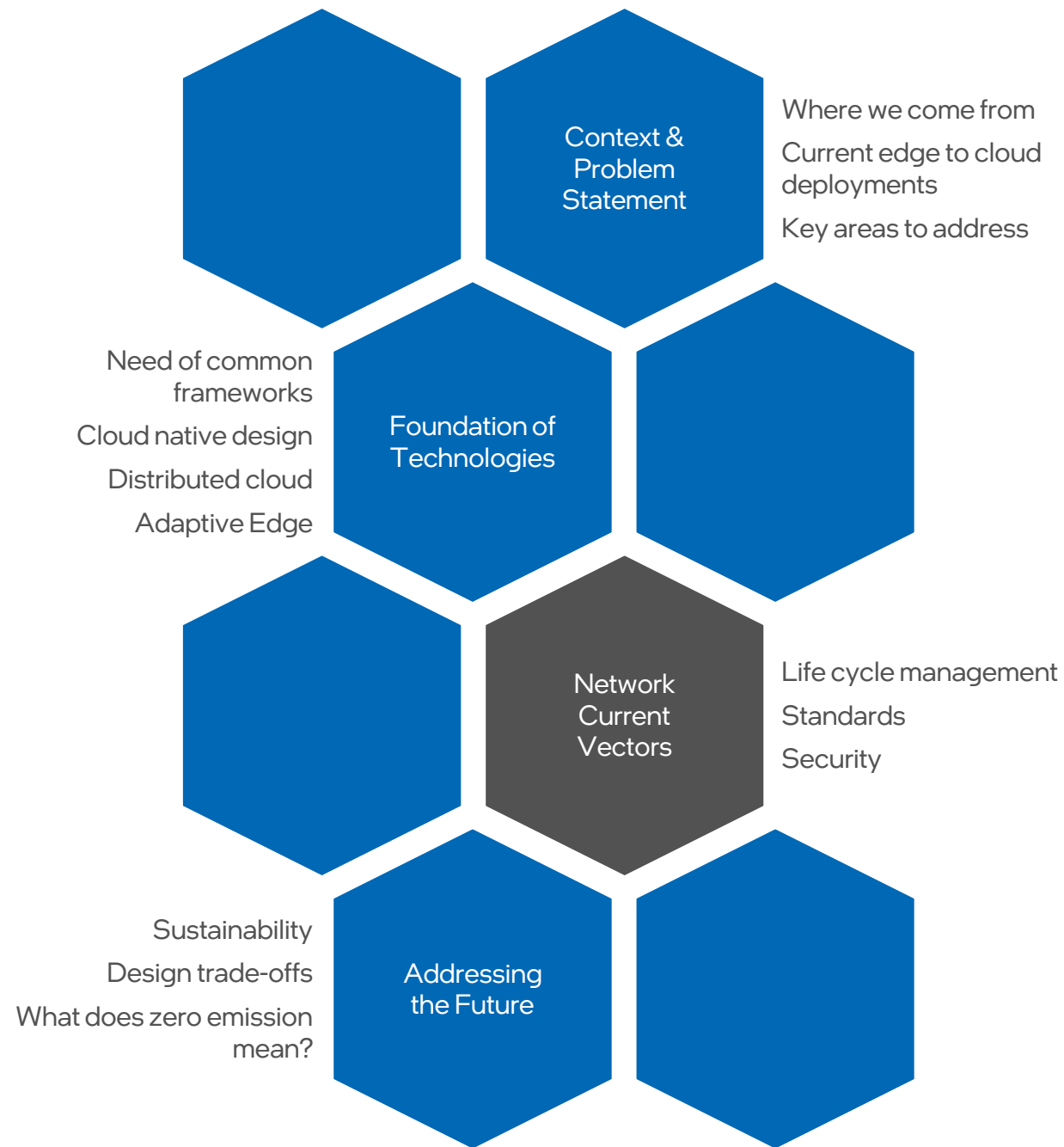
SE350      SE450

### Edge Compute Technologies

intel MOVIDIUS™      intel XEON™

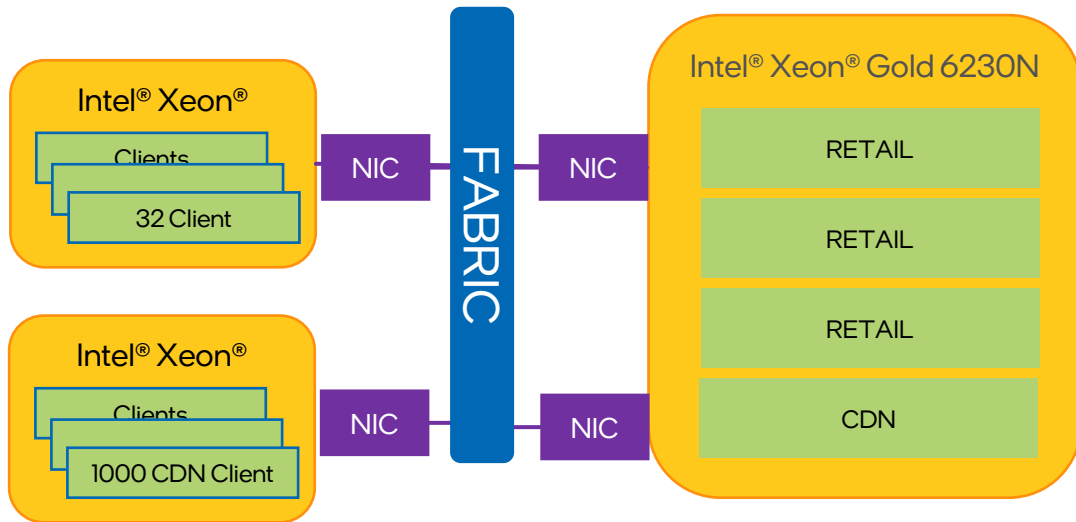
Intel® Xeon® Scalable Processors  
Intel® Gen 3 Movidius™ VPU-based Vision Accelerator Card

# Outline





# Network Matters for E2E QoS.. And Matters a Lot

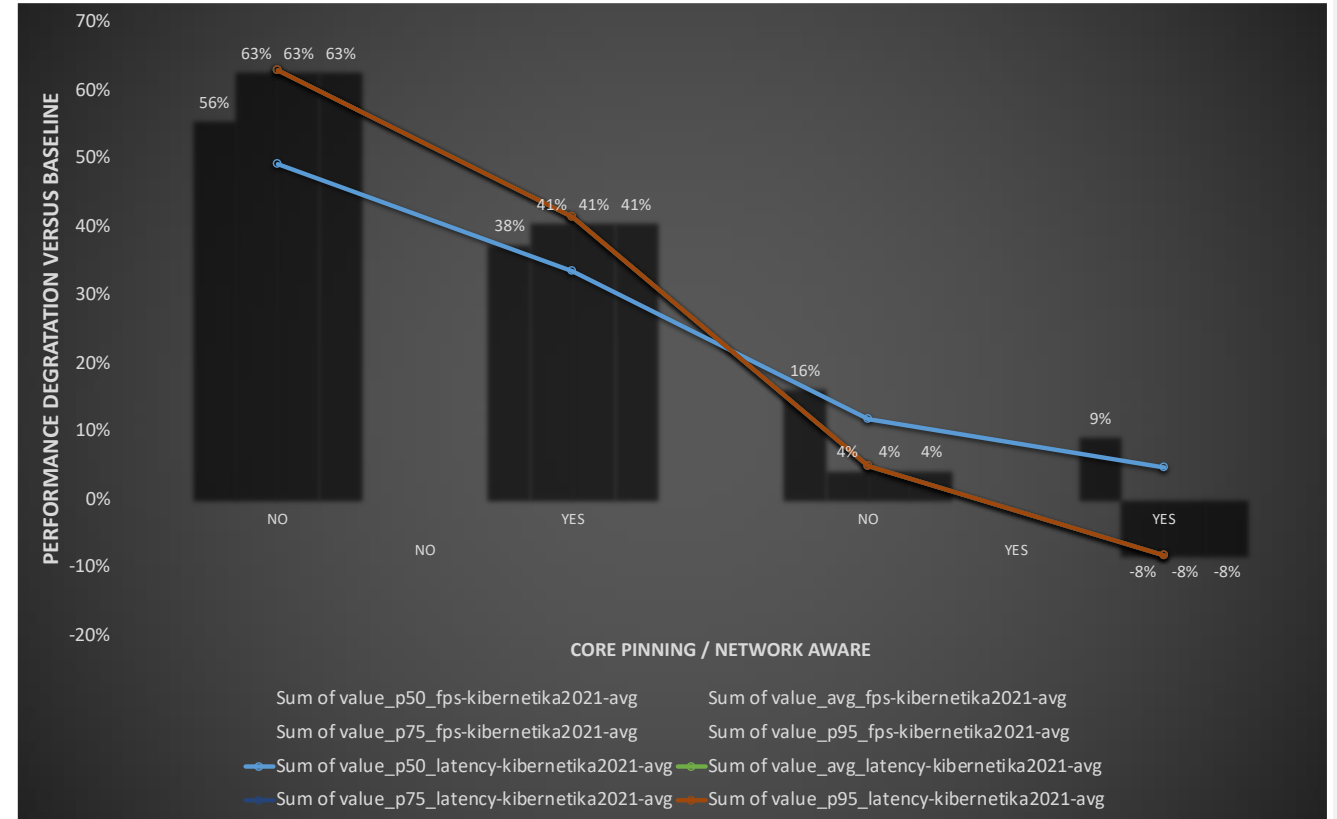


## CONDUCTED FOUR TYPE OF EXPERIMENTS:

- COMPUTE RESOURCE SELECTION (CORE PINNING)
  - a) YES = Workloads pinned into specific cores
  - b) NO = No core pinning
- NETWORK AWARE (INTERFACE SELECTION)
  - a) YES = Different interface for CDN and KIBERNETICA
  - b) NO = Same interface for both

## SHOWED RESULTS:

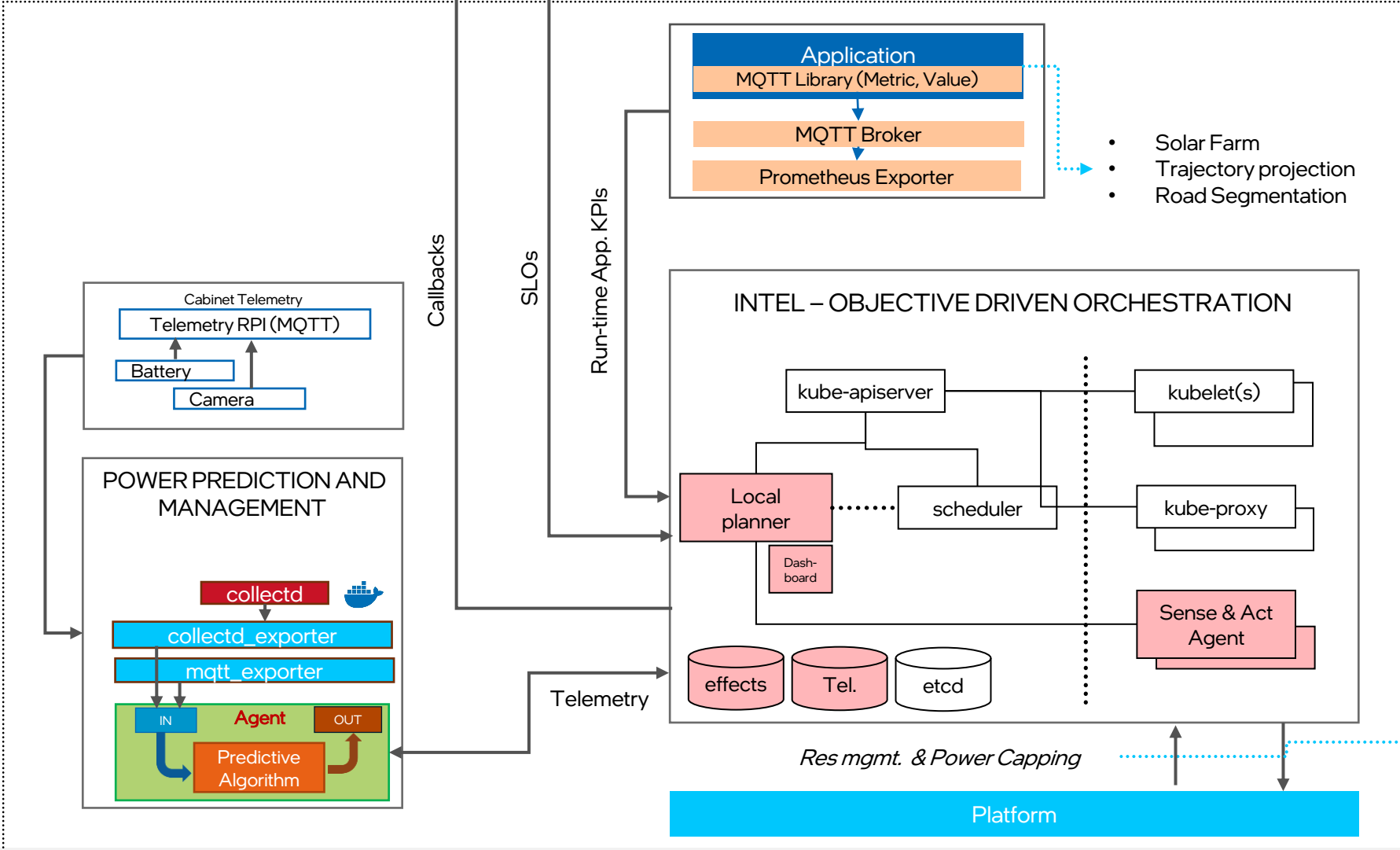
- Degradation of for each deployment with respect to baseline (Kibernetica alone)
- Application KPIS={fps, latency} – Statistics={avg, 50/75/95 percentile}



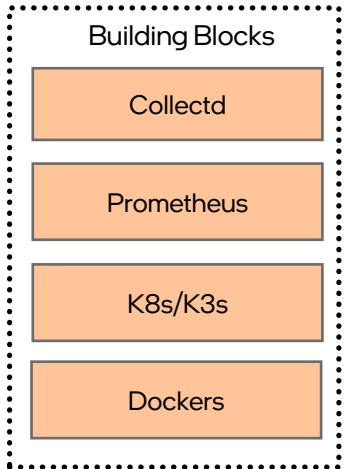
## NOTES:

- COMPUTE AWARE REDUCE DEGRADATION
- NETWORK AWARE REDUCE DEGRADATION
- **DISCLAIMER: we need to keep looking at the results**
- **DISCLAIMER 2: the fact that network + compute aware allocation has better performance than baseline can be something that we need more runs to normalize results**

# Application Life Cycle Management Matters a Lot

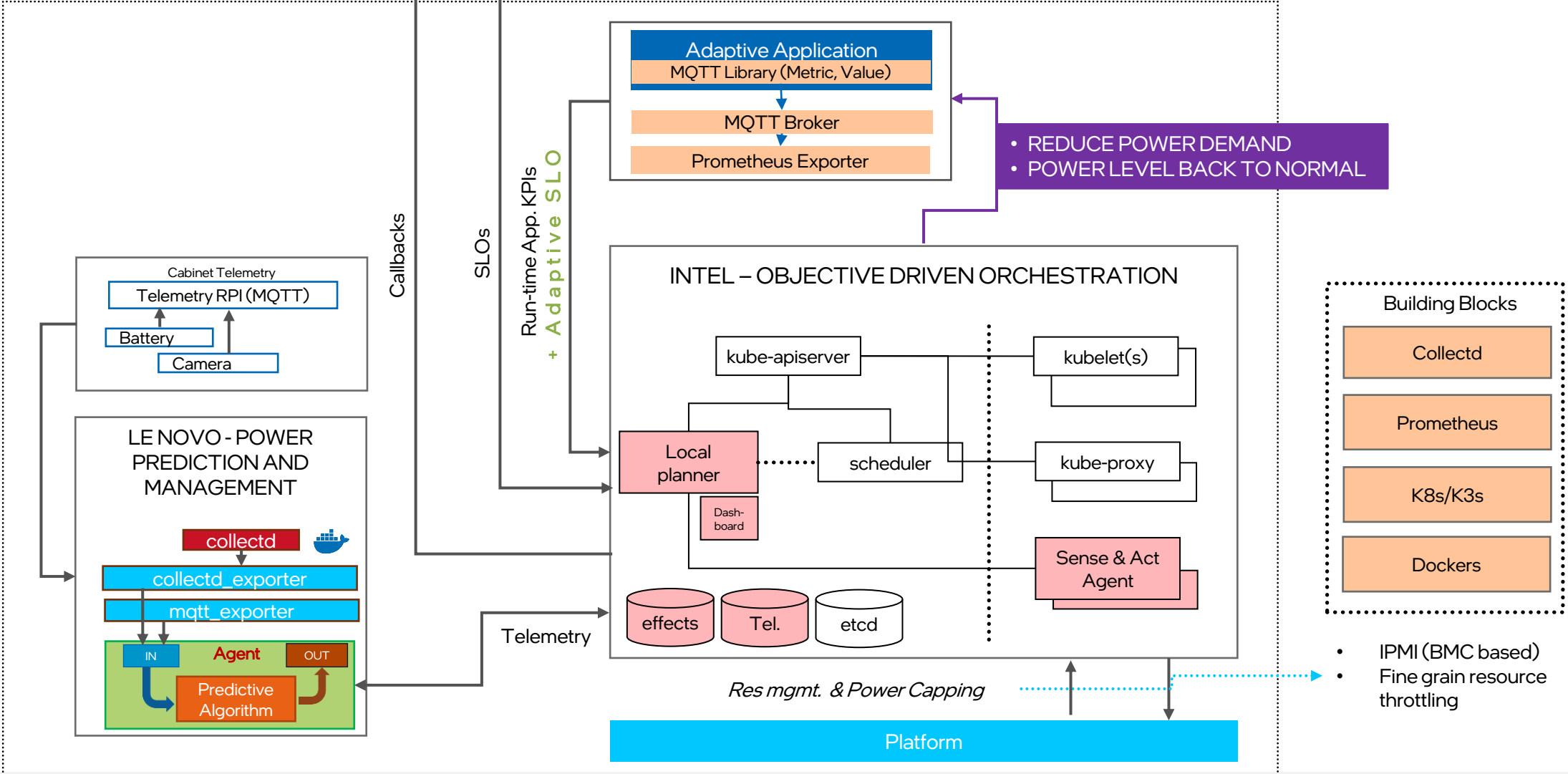


- Solar Farm
- Trajectory projection
- Road Segmentation



- IPMI (BMC based)
- Fine grain resource throttling

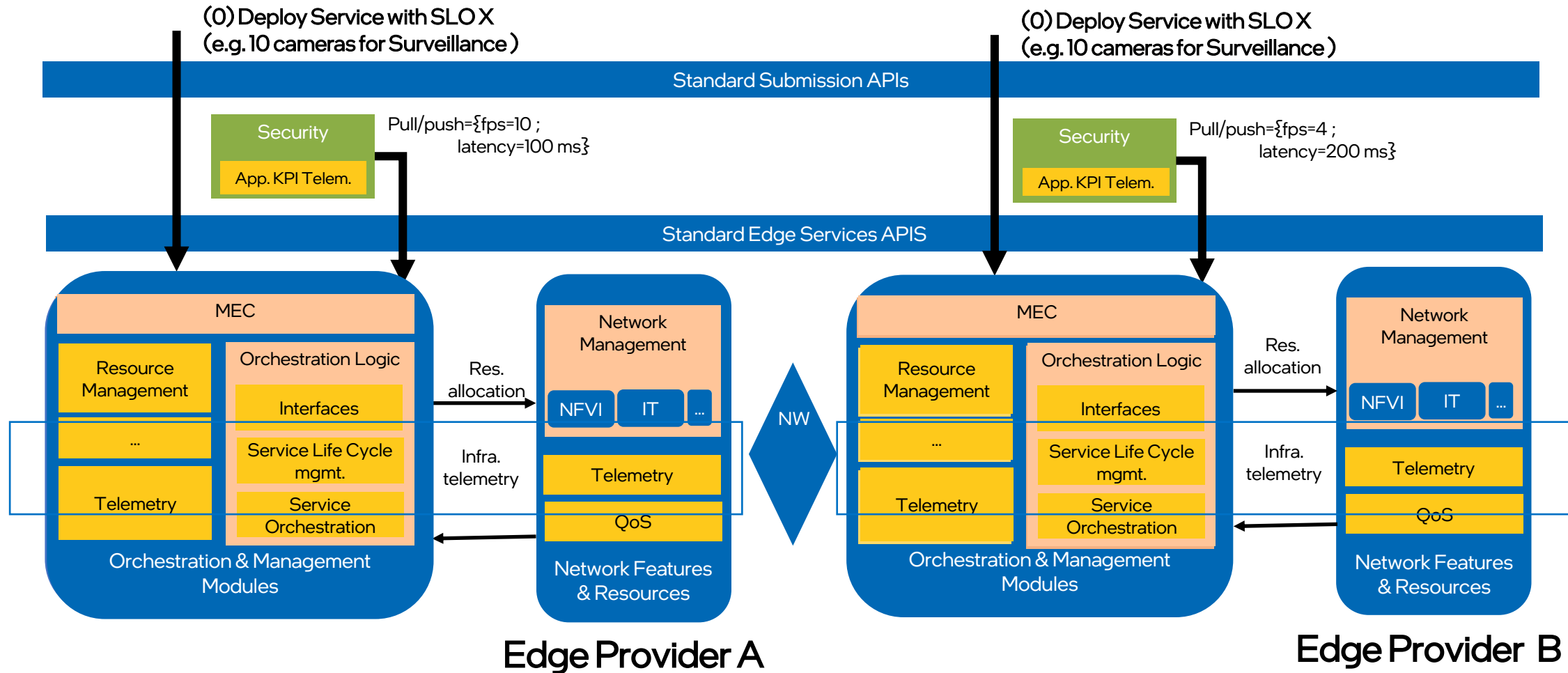
# Application and Service Can Adapt Sometimes



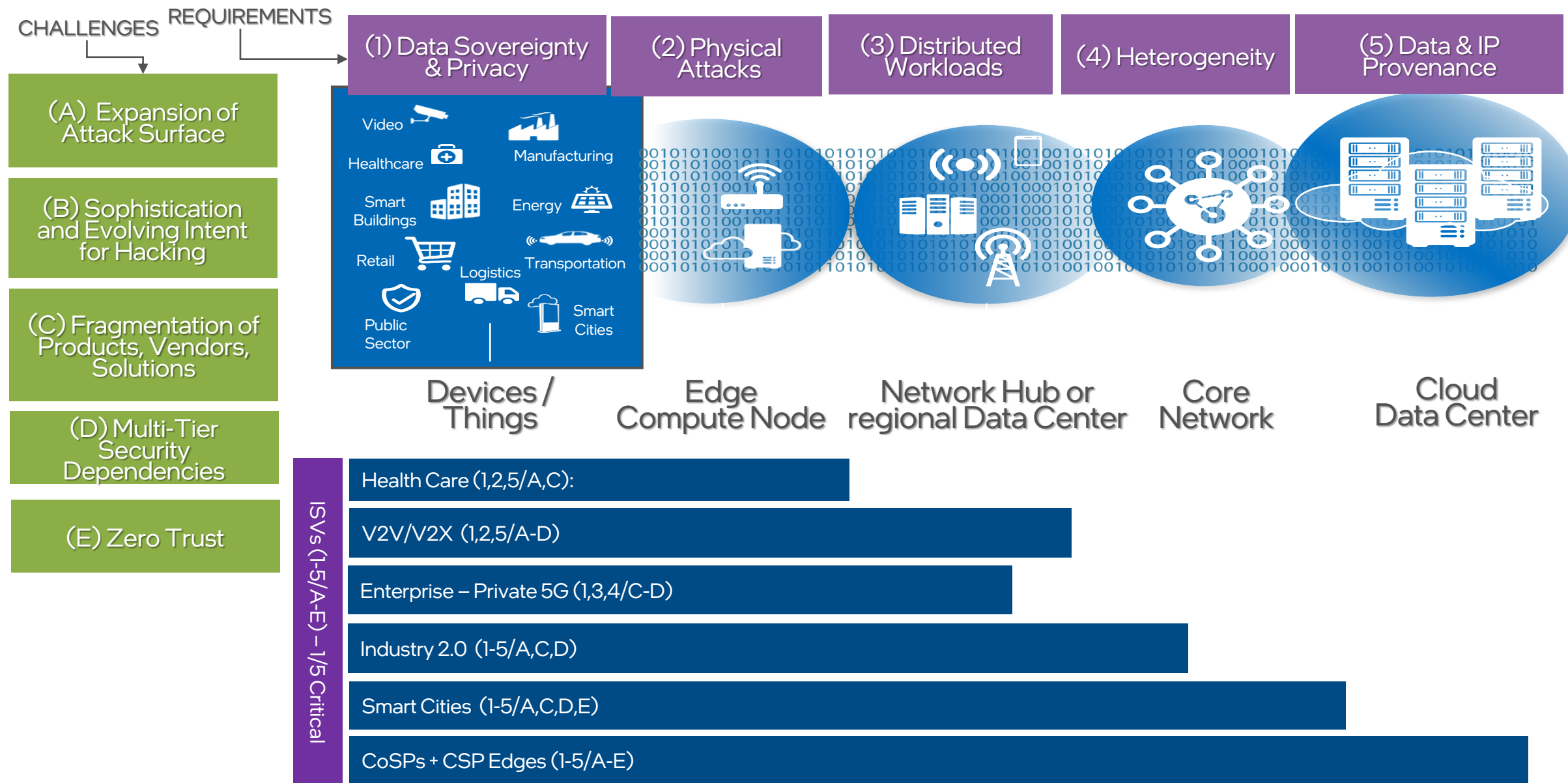
# But ... We Still Have Some Important Challenges

- Network end to end observability
- Network monitoring at service or application level
- Establish end-to-end QoS hooks between Edges and UE
- Application metrics vs system metrics in multitenant deployments
- And ....

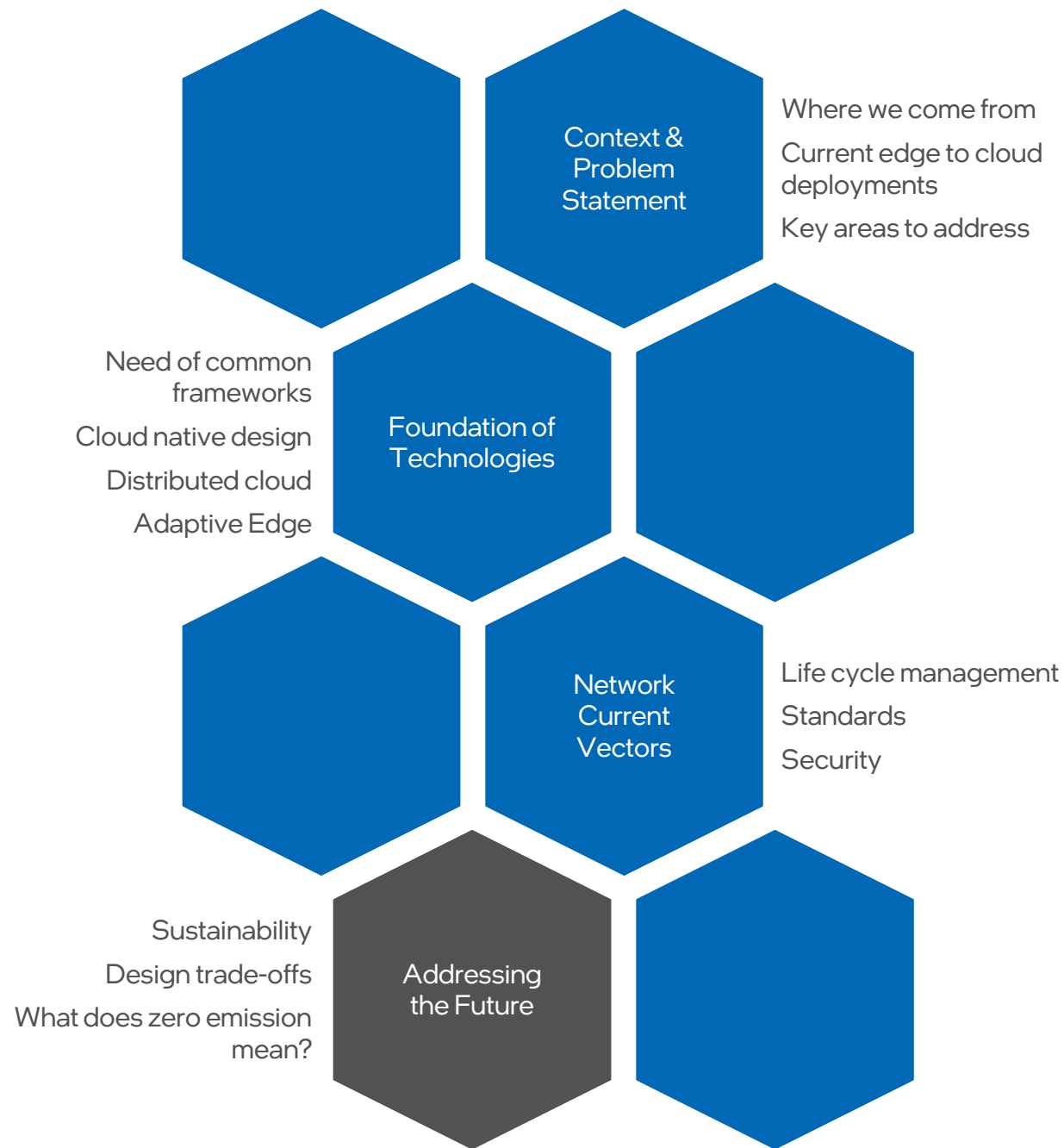
We need standards and common APIs  
 We need network to as part of the service life cycle mgmt.  
 Re CAMARA PROJECT, OPG ...



# ... and we need it secure



# Outline



# What Sustainable Distributed Edge May Require ...

SYSTEM SUSTAINABILITY KPIS

E2E INFRA & SERVICE ORCH.

DATA TO COMPUTE VS COMPUTE TO DATA

NEW ORCH. MECHANISMS

NETWORK

WATTS/BYTE  
(constant versus dynamic depending on load or technologies)

PHYSICAL INFRA.

HEAT RE-USE  
(constant versus dynamic depending on load or technologies)

## EDGE APPLIANCE

ENERGY IN ENERGY OUT

RENEWAL  
(solar, wind..)

ENERGY RE-USE (HEAT RE-USE)

SYS. ARCH INGREDIENTS

IPU ROLE

RESOURCE POOLING

NEW CHASIS AND MATERIALS

NEW COOLING SCHEMES

SILICON INGREDIENTS

POWER OPTIMIZATIONS

SILICON PARTITIONING

SUSTAINABLE-SST



# We Need to Adapt Distributed Edge System Architectures

Energy mgmt.,  
Observability,  
Networking,  
Connectivity ..



Chassis architecture,  
Cabinet design,  
Telemetry, Security ..



## System Trade-offs

- CORE VS XEON
- HEAT REUSE
- CORE VS XEON VS VPU VS GPU
- PERF VS POWER
- POWER VS RELIABILITY
- IMMERSION LIQUID VERSUS FAN
- PERF / \$\$
- CAPEX VERSUS OPEX
- SALABILITY
- RESOURCE DISAGGREGATION
- ACCELERATION VS GP
- ...

Lenovo



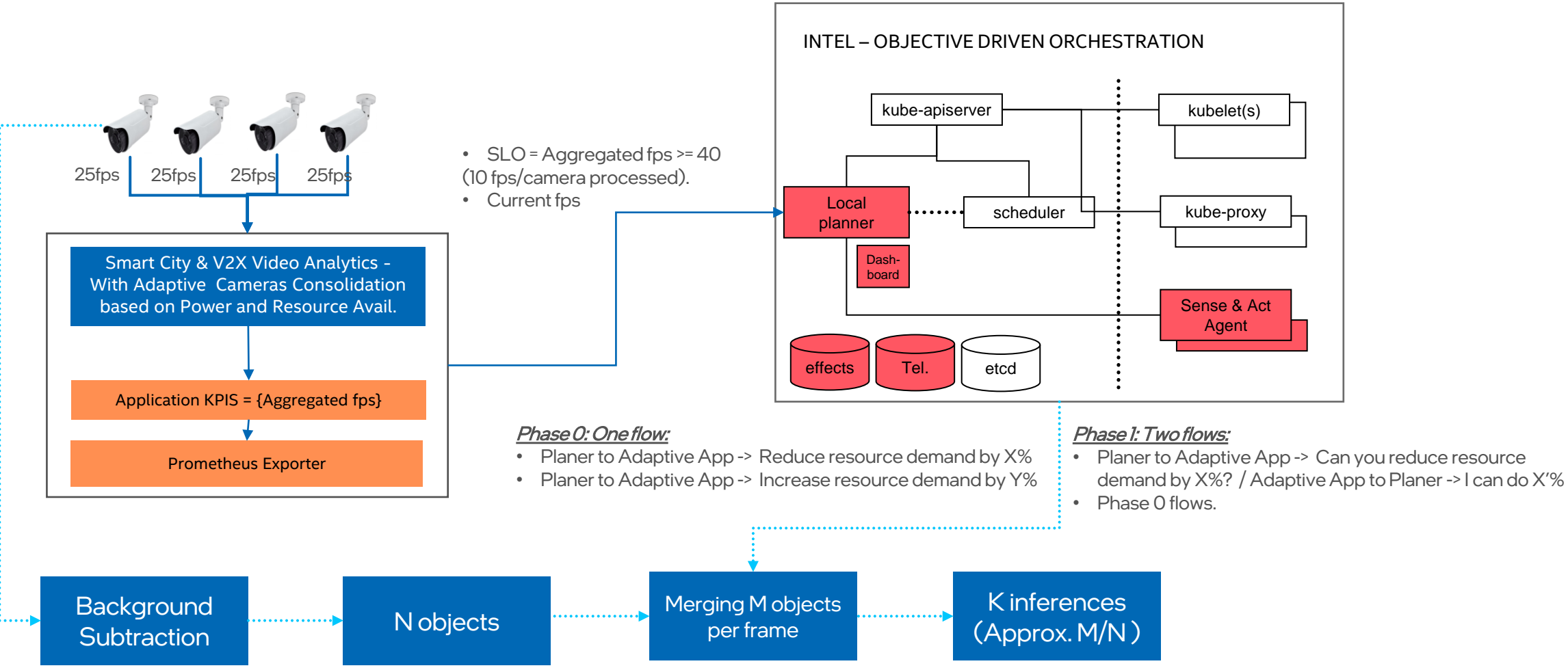
cellnex



Edge Emersion Cooling Edge  
Appliance (FF, mgmt., security,  
connectivity, power,..)..  
Platform architecture



# We Need Application Awareness



M = 1 -> Highest accuracy Highest Power Consumption  
 M > 1 -> The bigger M the lower accuracy and the Lower Power Consumption

# Call For Action

# Call For Action

1. Common services framework for Edge to Cloud for application development:

“Change the paradigm from “local/enterprise/private” or “specific cloud” to anywhere (Edge to Cloud)”

2. Event driven functions and Open APIs for interoperability

3. End-End Sustainability for Edge-Cloud

- Awareness at infrastructure (HW/SW) and application level
- Reporting, Telemetry, and feedback loop to achieve sustainability goals
- Using sustainability as a metric for deployment of distributed apps from Edge to Cloud

# References

# References

## Conference and White Papers

- Composable Architectures for a Sustainable Edge
  - <https://networkbuilders.intel.com/solutionslibrary/composable-architectures-for-a-sustainable-edge>
- Autonomous Lifecycle Management for Resource-Efficient Workload Orchestration for Green Edge Computing
  - <https://ieeexplore.ieee.org/document/9612603>
- Convergence of Edge Services & Edge Infrastructure
  - <https://ieeexplore.ieee.org/document/9665021>
  - <https://dblp.org/pid/08/2452.html>

# Questions?

Francesc Guim, Solutions Architect  
[francesc.guim@intel.com](mailto:francesc.guim@intel.com)

The Intel logo is centered on a solid blue background. It features the word "intel" in a white, lowercase, sans-serif font. A small blue square is positioned above the letter 'i'. To the right of the word "intel" is a registered trademark symbol (®).

intel®