

Selection of Services and relation to Compute Metrics

Sabine Randriamasy (*Nokia Bell Labs*)

Side-meeting

Exposure of Network and Compute information to Support Edge Computing Applications

IETF 119, Brisbane, March 19 2024

Context

- Revision of
 - [Joint Exposure of Network and Compute Information for Infrastructure-Aware Service Deployment \(ietf.org\)](https://ietf.org)
- Abstract
 - The presentation will highlight the need to consider different levels of granularity & abstraction that can be used to select the decision metrics. It will also propose a modular approach to combine and adapt metrics to jointly select egress routers and edge servers (or their LB) depending on context.

Metric selection and handling

Dimension	Definition	Examples	Impact on metrics
Decision	what are the metrics used for	monitoring, benchmarking, service selection and placement	<ul style="list-style-type: none">• Metrics are explicit or service and/or infrastructure KPI-driven• Metrics should impact KPIs• Requires different metric aggregation levels• Sensitive to metrics dynamicity, in/off path?• Needs synchronized information
Driving KPI	what is assessed with the metrics	speed, scalability, cost, stability	
Decision scope	different granularities	infrastructure node/cluster, compute service, end-to-end application	
Receiving entity	receiving metrics	router, centralized controller, application management	
Deciding entity	computing decisions	router, centralized controller, application management	

Table 4: Dimensions to consider when identifying compute metrics.

Information consumer scenarios & abstraction level

- Consumer is an ISP that has no access to full compute information
 - Compute metrics will likely be estimated
- Consumer is an application that has no access to information while the ISP does and is willing to provide guidance
 - Application can get with abstracted net & comp information
- Consumer has access to full network and compute information and wants to use it for fine-grained decision making e.g. at the node/cluster level
- Consumer has access to full information but essentially needs guidance with abstracted information
- Consumer has access to information that is abstracted or detailed depending on the metrics

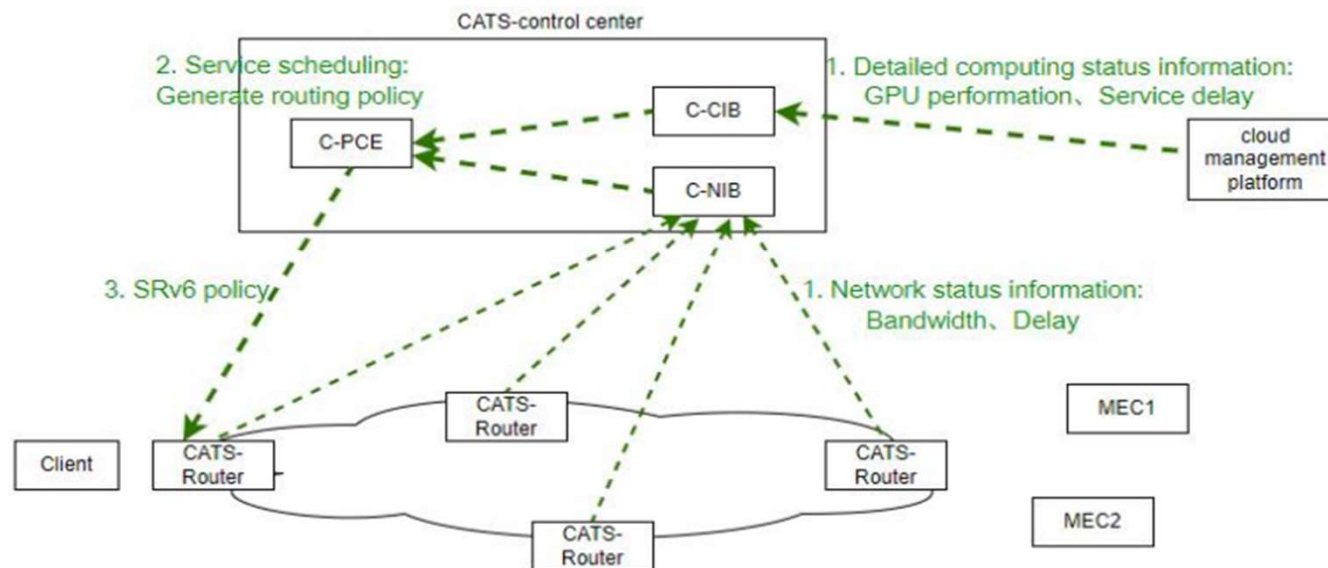
Flexible sets of decision metrics

- Metrics may be application KPIs or lower-level metrics impacting these KPIs
- A common metric such as network and compute delay may be convenient to handle if relevant and sufficient
- Application quality may be impacted by several metrics
- A vector of one or more metrics allows a more accurate and numerically consistent evaluation
- A metric may have a varying importance depending on context
 - E.g. bandwidth has more importance in peak hours or on lower capacity paths
- Metrics may have varying relative impact
 - fast computing requires both high capabilities in both network and computing
 - scheduled batch computing requires high computation capabilities when computation is done and high network capabilities when data is circulated for the computation
- Simultaneous selection of egress router and server (access point) is preferable

CATS WG - Slides presented at IETF 117

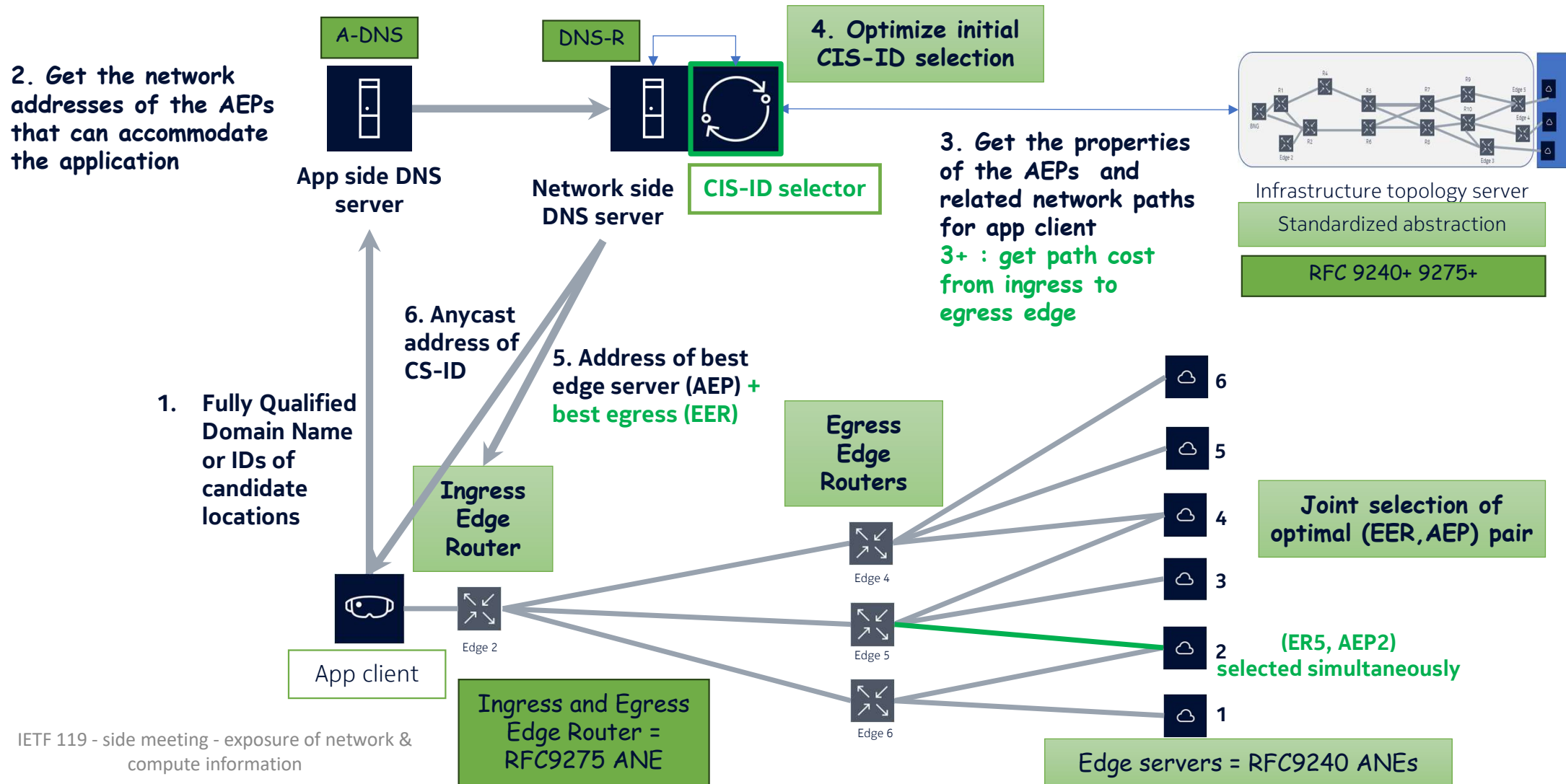
Workflow – A centralized model

- Computing information: aware by the CATS-control center by restful interface from cloud management platform.
- Network information: aware by the BGP-LS or telemetry interface from routers.
- CATS-control center performs service scheduling according to the detailed computing information and network information, then generates routing policy and **sends to CATS ingress router.**



Application to centralized CATS deployment model with ALTO-supported CATS Control Center

Simultaneous selection of servers and egress routers



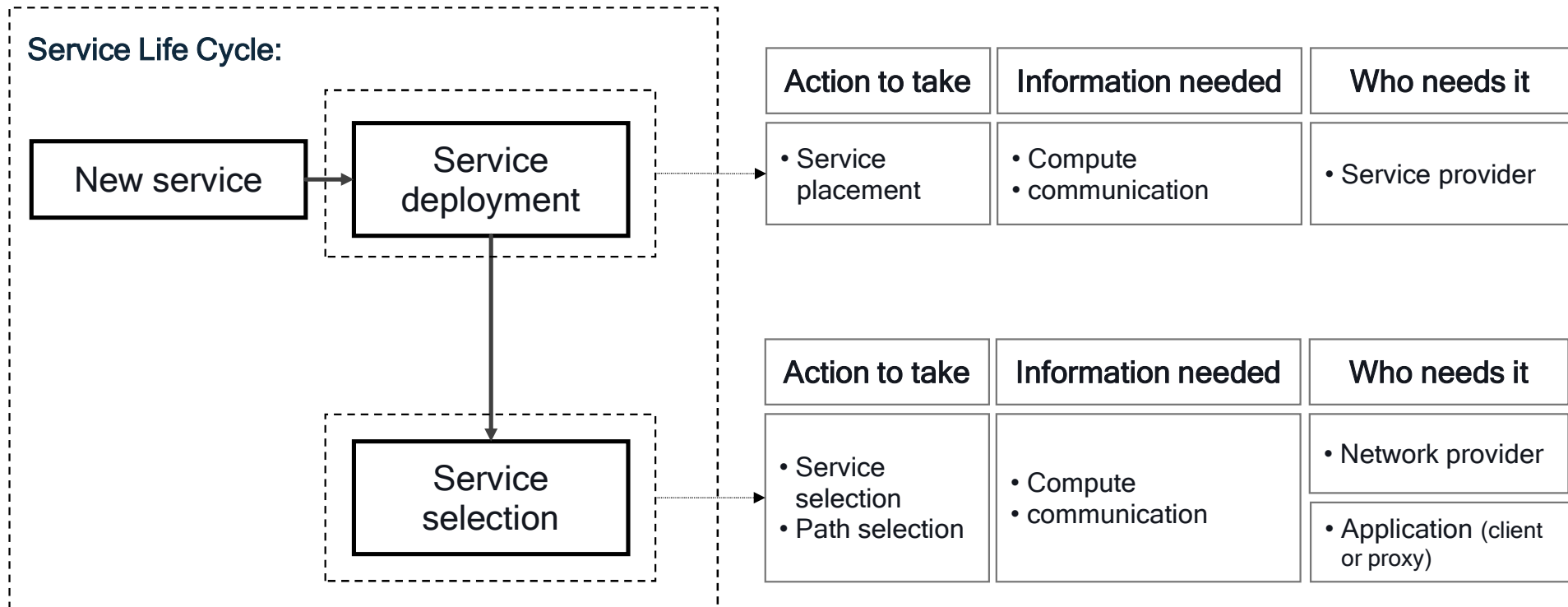
Thank you – Questions?

Back-up

Motivation

- Standardization of network information is quite mature but is in progress for compute information.
- There is a need to define a set of compute metrics to support various use cases being served in the IETF.
- Some ad hoc work exists in the IETF:
 - CATS (e.g., draft-du-cats-computing-modeling-description)
 - ALTO (e.g., draft-contreras-alto-service-edge)
 - OPSAWF (e.g., RFC 7666 MIB)
- Metrics are also defined in other bodies such as the Linux Foundation, DMTF, ETSI NFV, etc:
 - Raw compute infrastructure metrics (e.g., processing, memory, storage)
 - Compute virtualization resources and service quality metrics (e.g., VNF resources in VMs)
 - Service metrics including compute-related information (e.g., service delay, availability)

Problem space



History and updates from IETF 118

- -01 version presented at IETF 118, collecting initial feedback
- Updates from -01 (now in -03 version)
 - Use cases documented for better illustrating the problem space
 - New section on “Production and Consumption Scenarios of Compute-related Information”
 - Reference to raw resources and allocated resources
 - New section on “Metrics Selection and Exposure”
 - Reference to how the metrics are exposed and (2) which kind of metrics need to be exposed
 - Discussion on dimensions to consider when identifying compute metrics
 - Discussion on abstraction levels and information access
 - Reference to distribution and exposure mechanisms
- Added Roland as co-author

Example of simultaneous selection of servers and egress routers

Assumptions for the example

- all vector components are criteria to be maximized
- « raw weights » $w_n = w_c = 1$

Normalized weights

- $W_{un} = w_{un}/(w_{un}+w_{uc})$
- $W_{uc} = w_{uc}/(w_{un}+w_{uc})$

Global utility value UG

$$UG = W_{un} * U_n + W_{uc} * U_c$$

Sequential selection: (S3, R3) UG = 0.875
 (if server first then egress)

Simultaneous selection: (S1, R1) UG = 0.8888

Non-normalized selection hard to decide

