

# Application-Aware Resource Provisioning in a Heterogeneous Internet of Things

---

Eric Sturzinger\*, **Massimo Tornatore** † \*, and Biswanath Mukherjee\*

\*University of California, Davis

† Politecnico di Milano

16 May 2017

**UCDAVIS**

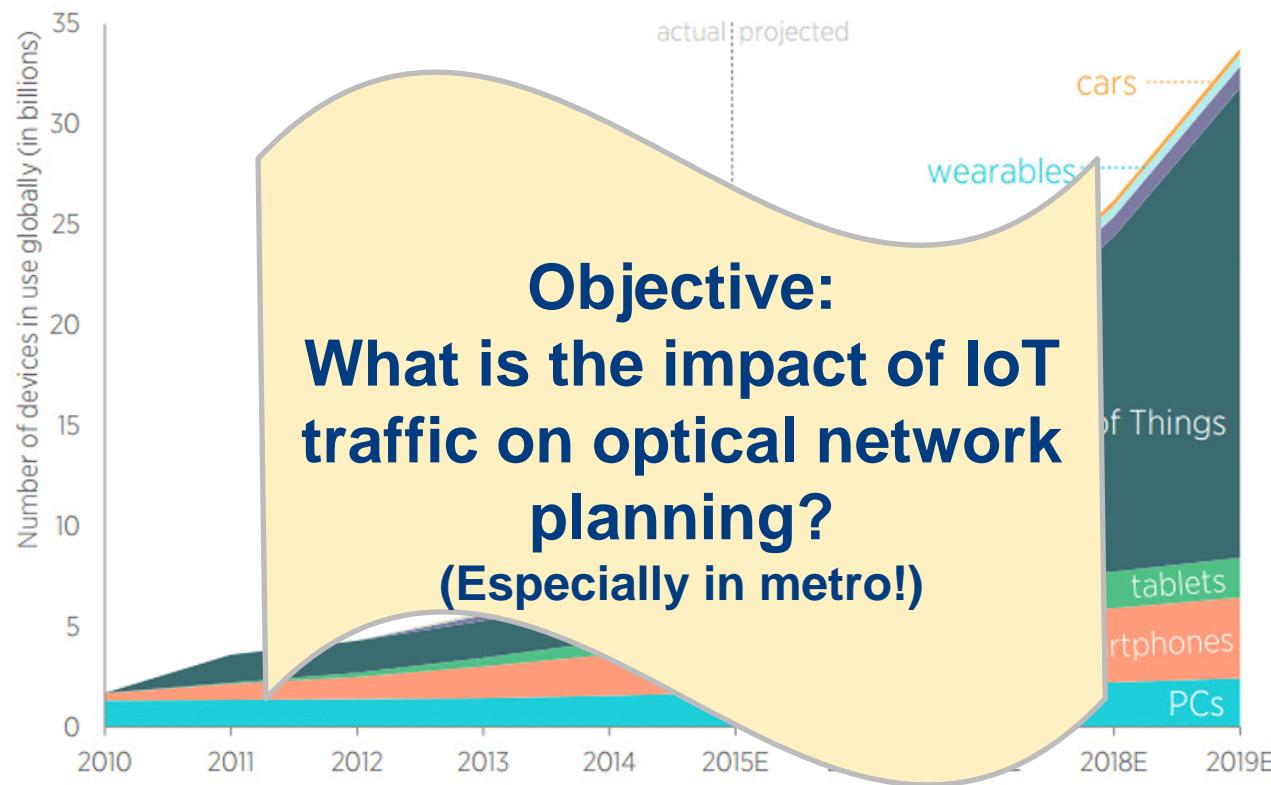


# Outline

- Motivation
- Application profiles
- Mathematical formulation
- Simulation Results
- Conclusion

## Motivation and objective

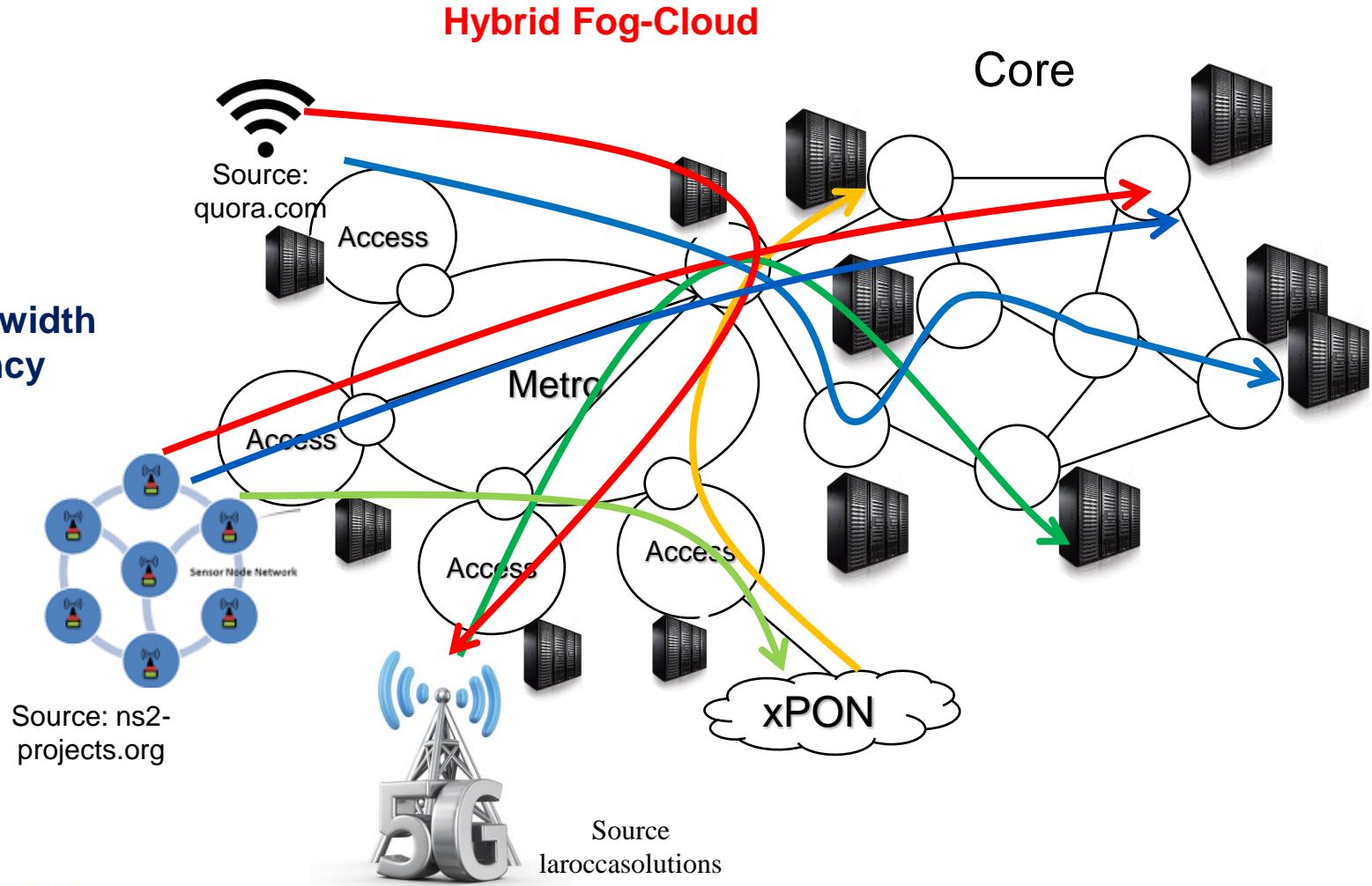
- Next big wave:  
**Internet of Things (IoT) and Machine to Machine (M2M) traffic**



Source: John Greenough, "The Internet of Everything 2015," *Business Insider Intelligence*. Produced by Adam Thierer and Andrea Castillo, Mercatus Center at George Mason University, 2015.

# Enabler: hybrid fog-cloud computing

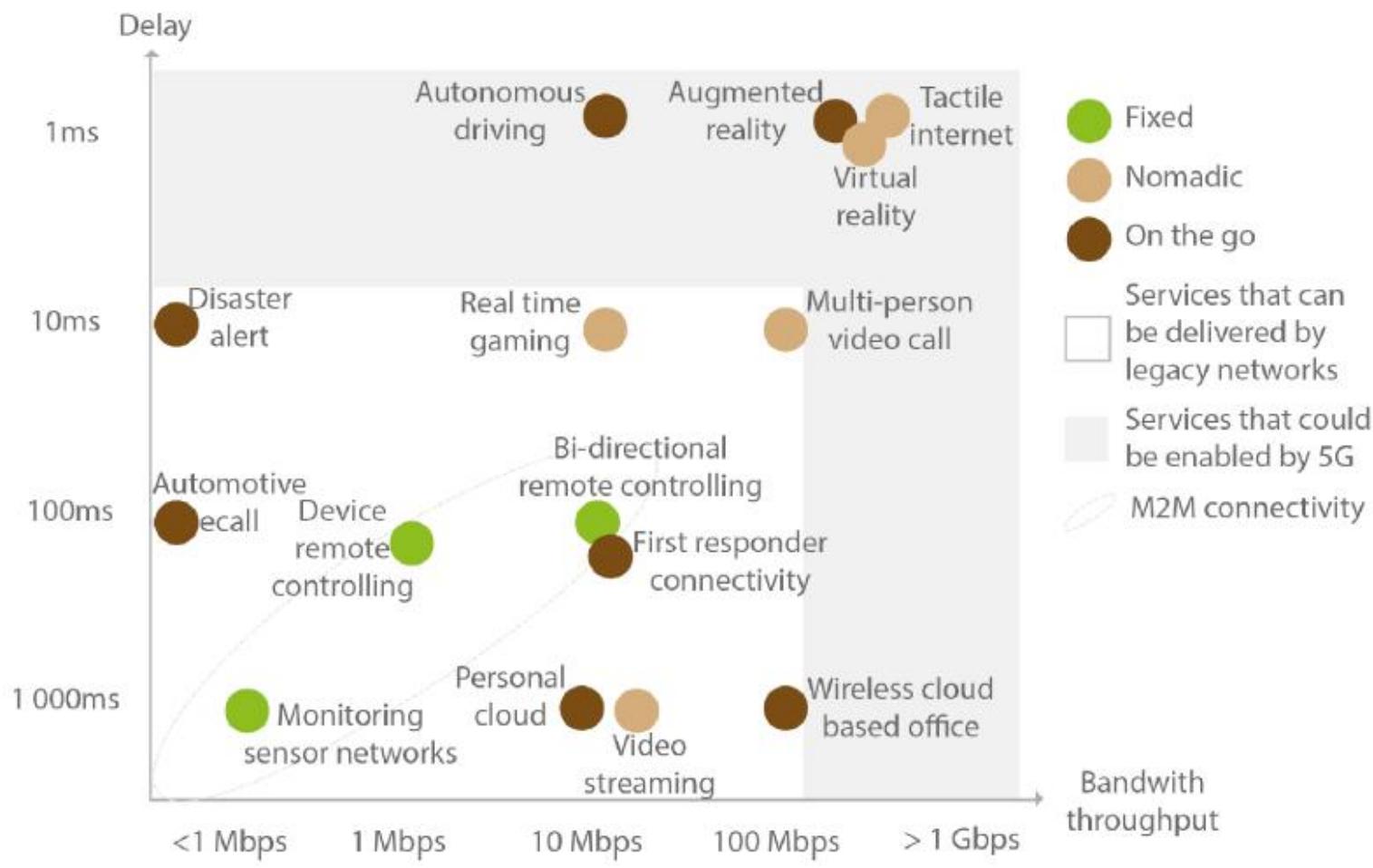
- Drivers
  - Bandwidth
  - Latency



# Challenges

1. Lack of quantitative characterization of IoT and M2M application
  - Identify **application profiles**
2. Network **cost** parameters
3. **A model to assign resources** (bandwidth, computing, storage):
  - tailored to application profile (*slicing*)
  - minimizing costs

# Application Profile (1)



Source: GSMA Intelligence, 2015

## Application Profile (2)

- Each application profile contains a unique combination of parameters
  - **$\Theta$ : Latency budget** (source to destination)
  - **$\kappa$ : Bandwidth**
  - **$\alpha$ : Computational complexity** (per unit of traffic)
  - **$\Lambda$ : Storage time**
  - **$\beta$ : Compression factor** (ratio of processed-to-raw data)

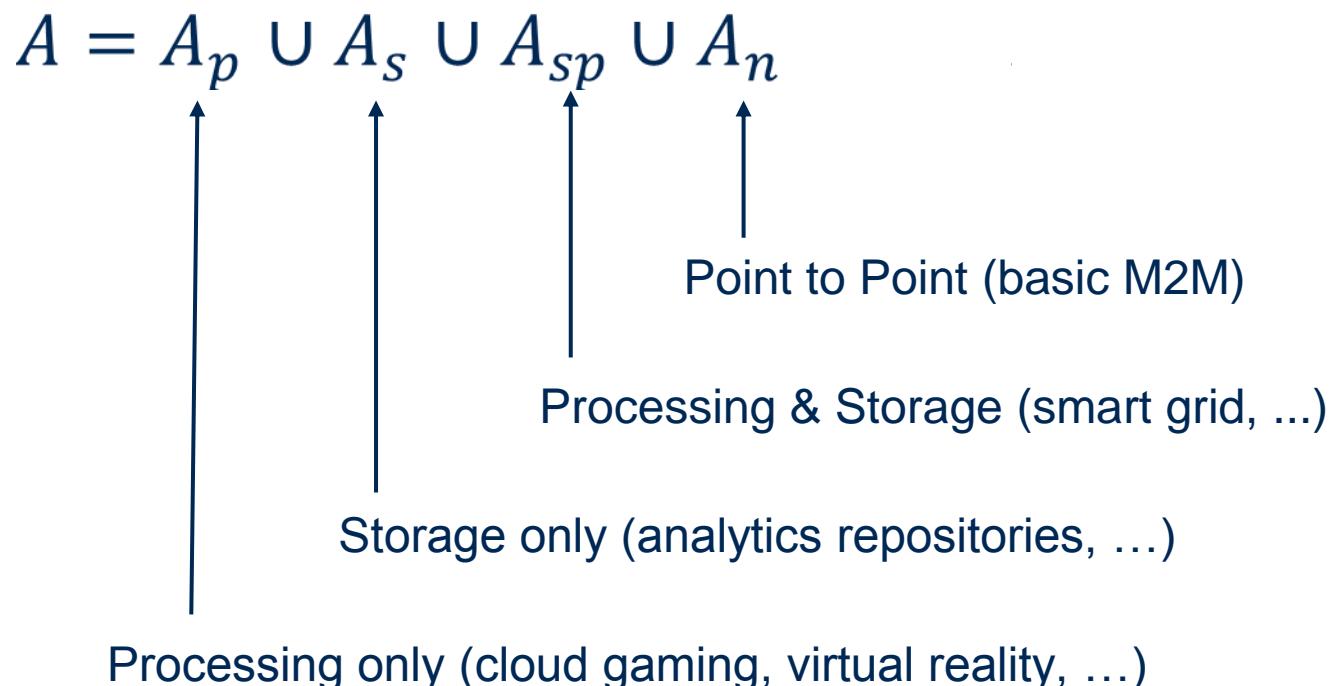
Examples	$\Theta$ (ms)	$\kappa$ (Mbps)	$\alpha$ (CPU/Mbps)	$\beta$	$\Lambda$ (hrs)
1- AR/VR	10	100	0.03	0.6	0
2 – Factory Automation	20	1	0.009	0.8	10
3 – Data Backup	1000	1	0	0	4
4 – Smart Grid	50	0.4	0.007	0.3	0
5 – Smart Home	60	.001	0	0	0
6 – Medical	40	2	0.02	0.2	0.1
7 – Environmental Mon.	1000	1	0.02	0.1	100
8 – Tactile Internet	1	200	.005	0.8	0

5G PPP, "5G Automotive Vision," white paper, 2015.

A. Frotzscher *et al.*, "Requirements and Current Solutions of Wireless Communication in Industrial Automation," *Proc.IEEE ICC Wksps.*, Sydney, Australia, 2014, pp. 67–72.

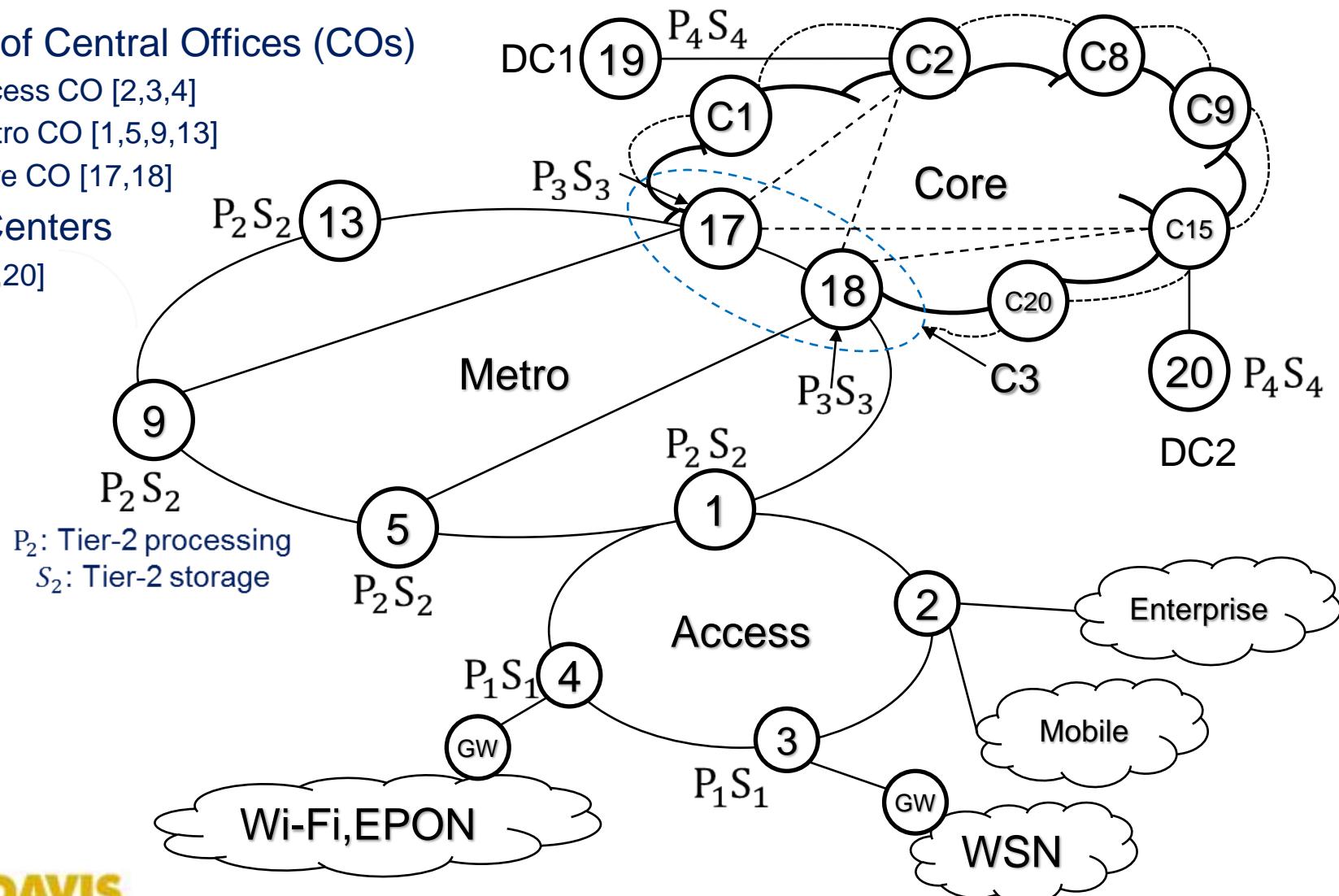
## Application Profile (3)

- Additional note: we classify applications in four categories



# Hybrid Fog-Cloud Network Architecture

- 3 tiers of Central Offices (COs)
  - Access CO [2,3,4]
  - Metro CO [1,5,9,13]
  - Core CO [17,18]
- Data Centers
  - [19,20]



# Network (Cost) Parameters

## •• Costs (in COs)

- $\mu$ : processing cost
- $v$ : storage cost
- $\Lambda$ : metro bandwidth cost
- $\epsilon_{up}, \epsilon_{down}$ : core bandwidth cost

Tier	$\mu$ (\$/CPU/Mo)	$v$ (\$/GB/Mo)	$\Lambda$ (\$/Mbps/Mo)	$\epsilon_{up}, \epsilon_{down}$ (\$/Mbps/Mo)
1 – Access CO	90	0.0042	~1	~1/1
2 – Metro CO	70	0.004	~1	~1/1
3 – Core CO	50	0.0035	~1	~1/1
4 - DC	25	0.0025	~1	~1/1

<https://cloud.google.com/compute/pricing>

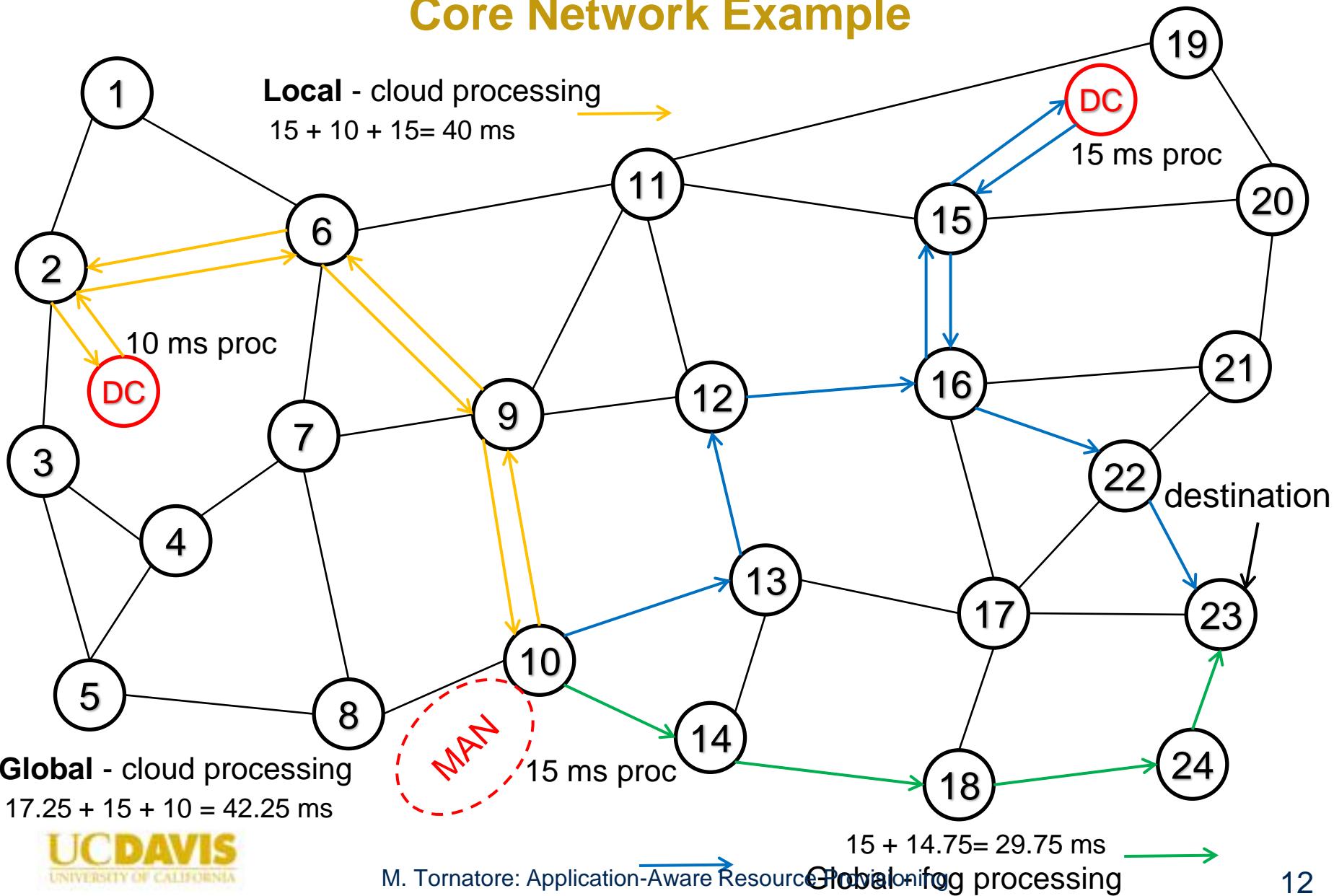
<https://cloud.google.com/storage/pricing#pricing-example-simple>

<http://drpeering.net/white-papers/Internet-Transit-Pricing-Historical-And-Projected.php#>

# Problem Statement

- **Inputs**
  - Offered traffic (per s-d pair, per application)
  - Application profiles:
    - $\Theta, \kappa, \alpha, \beta, \Delta$
  - Hybrid fog-cloud network topology  $G(N,L)$
- **Objective function**
  - Minimize total resource provisioning cost
- **Constraints**
  - Node (processing, storage) and link capacity
  - Latency
- **Outputs**
  - “Slice” (per s-d pair, per application) consisting of:
    - Path(s) (with required bandwidth)
    - Required processing and storage resources at each node

# Core Network Example



# Mathematical Formulation

Variables

## Objective Function:

$$\min(Cost_p^1 + Cost_s^2 + Cost_u^3 + Cost_d^4 + Cost_c^5)$$

$$x_{a,f}^{s,m} \in \{0, 1\}$$

$$r_{a,k,f}^{s,m} \in \{0, 1\}$$

$$r_{a,k,f}^{\prime,s,m} \in \{0, 1\}$$

$$r_{a,k,f}^{\prime\prime,s,m} \in \{0, 1\}$$

### 1 Processing

$$Cost_p = \sum_{m \in \mathcal{N}_p} \mu_m \sum_{a \in \mathcal{A}_p \cup \mathcal{A}_{sp}} \alpha_a \sum_{s \in \mathcal{N}_g} \sum_{f \in \mathcal{N}_g \cup \mathcal{F}_c} x_{a,f}^{s,m} v_{a,f}^{s,m}$$

### 2 Storage

$$Cost_s = \sum_{a \in \mathcal{A}_s} \Delta_a \sum_{f \in \mathcal{N}_s} \nu_f \sum_{s \in \mathcal{N}_g} x_{a,f}^{s,m} v_{a,f}^{s,m} + \sum_{a \in \mathcal{A}_{sp}} \beta_a \Delta_a \sum_{f \in \mathcal{N}_s} \nu_f \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_p} x_{a,f}^{s,m} v_{a,f}^{s,m}$$

### 3 Upstream core BW

$$Cost_u = \epsilon_{up} \left[ \sum_{a \in \mathcal{A}_p \cup \mathcal{A}_{sp}} \beta_a \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_{pl}} \sum_{f \in \mathcal{N}_{DC} \cup \mathcal{F}_c} x_{a,f}^{s,m} v_{a,f}^{s,m} + \sum_{a \in \mathcal{A}_p \cup \mathcal{A}_s \cup \mathcal{A}_{sp}} \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_{DC}} \sum_{f \in \mathcal{N}_g \cup \mathcal{N}_s \cup \mathcal{F}_c} x_{a,f}^{s,m} v_{a,f}^{s,m} + \right. \\ \left. \sum_{a \in \mathcal{A}_n} \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_c} \sum_{f \in \mathcal{F}_c} v_{a,f}^{s,m} \right]$$

### Downstream core BW

$$Cost_d = \epsilon_{down} \sum_{a \in \mathcal{A}_p \cup \mathcal{A}_{sp}} \beta_a \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_{DC}} \sum_{f \in \mathcal{N}_g \cup \mathcal{N}_{sl}} x_{a,f}^{s,m} v_{a,f}^{s,m}$$

4

$$Cap_{1,i,j} = \sum_{a \in \mathcal{A}_n \cup \mathcal{A}_s} \sum_{s \in \mathcal{N}_g} \left[ \sum_{f \in \mathcal{N}_g \cup \mathcal{N}_{sl}} \sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{s,m} v_{a,f}^{s,m} + \sum_{f \in \mathcal{F}_c} \sum_{m \in \mathcal{N}_c} \sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{s,m} v_{a,f}^{s,m} \right]$$

5

### Metro BW

$$Cap_{2,i,j} = \sum_{a \in \mathcal{A}_p \cup \mathcal{A}_{sp}} \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_p} \sum_{f \in \mathcal{N}_g \cup \mathcal{N}_s \cup \mathcal{F}_c} \sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{s,m} v_{a,f}^{s,m}$$

$$Cap_{3,i,j} = \sum_{a \in \mathcal{A}_p \cup \mathcal{A}_{sp}} \beta_a \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_p} \sum_{f \in \mathcal{N}_g \cup \mathcal{N}_s \cup \mathcal{F}_c} \left[ \sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{\prime s,m} + \sum_{d \in \mathcal{N}_c} \sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{\prime\prime s,m,d} \right] v_{a,f}^{s,m}$$

$$Cost_c = \Lambda \sum_{i,j \in \mathcal{L}_l} \sum_h Cap_{h,i,j}$$

# Mathematical Formulation (cont.)

## Constraints:

### Processing/Storage Assignments

$$\sum_{m \in \mathcal{N}_p} x_{a,f}^{s,m} = 1, \forall (a \in \mathcal{A}_p, s \in \mathcal{N}_g, f \in \mathcal{N}_g \cup \mathcal{F}_c)$$

$$\sum_{f \in \mathcal{N}_s} x_{a,f}^{s,m} = 1, \forall (a \in \mathcal{A}_s, s \in \mathcal{N}_g, m = f)$$

$$\sum_{m \in \mathcal{N}_p} \sum_{f \in \mathcal{N}_s} x_{a,f}^{s,m} = 1, \forall (a \in \mathcal{A}_{sp}, s \in \mathcal{N}_g)$$

$$\sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{s,m} D_k^{s,m} \leq \theta_{a,f}, \forall (a \in \mathcal{A}_n \cup \mathcal{A}_s, s \in \mathcal{N}_g, m \in \mathcal{N}_l \cup \mathcal{N}_s, f \in \mathcal{N}_g \cup \mathcal{N}_s \cup \mathcal{F}_c)$$

**Latency – Pt to Pt, Storage**

$$\sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{s,m} D_k^{s,m} + \gamma_{a,m} + \sum_k r_{a,k,f}^{ts,m} D_k^{m,f} \leq \theta_{a,f,m}, \forall (a \in \mathcal{A}_p \cup \mathcal{A}_{sp}, s \in \mathcal{N}_g, m \in \mathcal{N}_p, f \in \mathcal{N}_g \cup \mathcal{N}_s)$$

**Latency – Local Proc, Proc/Storage**

### Solenoidality

$$\sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{s,m} = x_{a,f}^{s,m}, \forall (a \in \mathcal{A}_p \cup \mathcal{A}_{sp} \cup \mathcal{A}_s, s \in \mathcal{N}_g, m \in \mathcal{N}_p, f \in \mathcal{N}_g \cup \mathcal{N}_s \cup \mathcal{F}_c)$$

$$\sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{ts,m} = x_{a,f}^{s,m}, \forall (a \in \mathcal{A}_p \cup \mathcal{A}_{sp}, s \in \mathcal{N}_g, m \in \mathcal{N}_p, f \in \mathcal{N}_g \cup \mathcal{N}_s)$$

$$\sum_{m \in \mathcal{N}_c} \sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{ts,m,d} = x_{a,f}^{s,m}, \forall (a \in \mathcal{A}_p, s \in \mathcal{N}_g, f \in \mathcal{F}_c)$$

### Processing Delay

$$\gamma_{a,m} = \alpha_a \kappa_a \tau_m$$

$$\sum_{a \in \mathcal{A}_p \cup \mathcal{A}_{sp}} \alpha_a \sum_{s \in \mathcal{N}_g} \sum_{f \in \mathcal{N}_g \cup \mathcal{N}_s \cup \mathcal{F}_c} x_{a,f}^{s,m} v_{a,f}^{s,m} \leq C_m, \forall m \in \mathcal{N}_p$$

### Storage Capacity

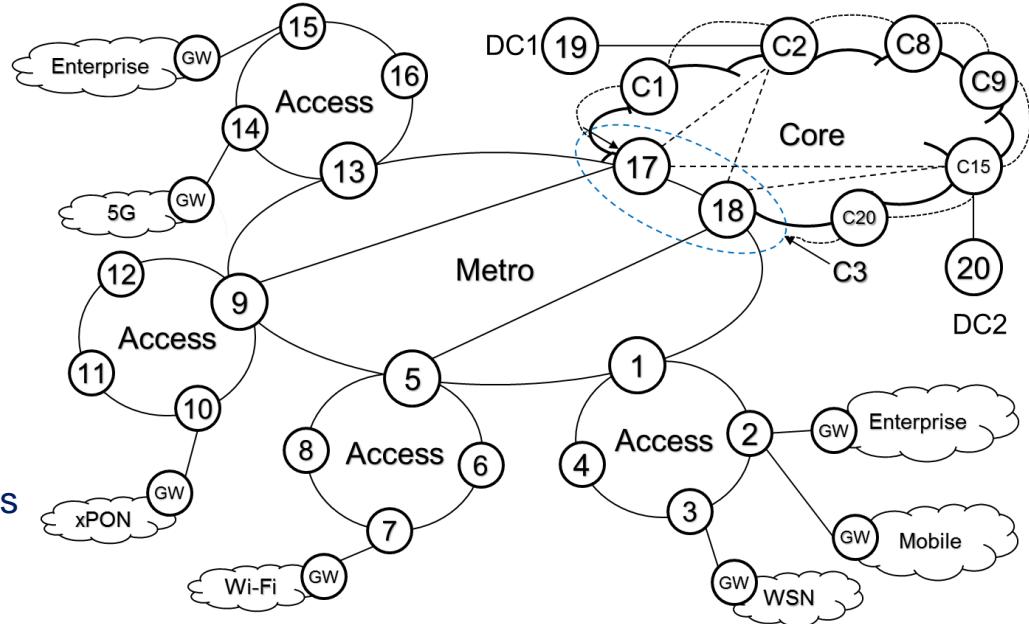
$$\sum_{a \in \mathcal{A}_s} \Delta_a \sum_{s \in \mathcal{N}_g} x_{a,f}^{s,m} v_{a,f}^{s,m} + \sum_{a \in \mathcal{A}_{sp}} \beta_a \Delta_a \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_p} x_{a,f}^{s,m} v_{a,f}^{s,m} \leq S_f, \forall f \in \mathcal{N}_s$$

$$\sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{s,m} D_k^{s,m} + \gamma_{a,m} + \sum_k r_{a,k,f}^{ts,m,d} D_k^{m,d} \leq \theta_{a,f}, \forall (a \in \mathcal{A}_p, s \in \mathcal{N}_g, m \in \mathcal{N}_p, d \in \mathcal{N}_c, f \in \mathcal{F}_c)$$

**Latency – Global Destination**

# Simulation Setup

- Traffic Volume – 5 Tbps
- Profiles
  - Computational complexity:
    - 0.005 - 0.03 CPU/Mbps
  - Compression factor: 0.1 – 1
  - Latency: 10 – 100 ms
    - Real-time ~ 10-50 ms
    - Near real-time ~ 50-100 ms
- Network resource costs
  - Processing cost: 25, 50, 70, 90 \$/CPU/Mo
    - Tier 4 (DC) costs [Google]
  - Storage cost: 2.50, 5, 7, 9 \$/TB/Mo
  - Metro/Core BW ~ 1\$ / Mbps

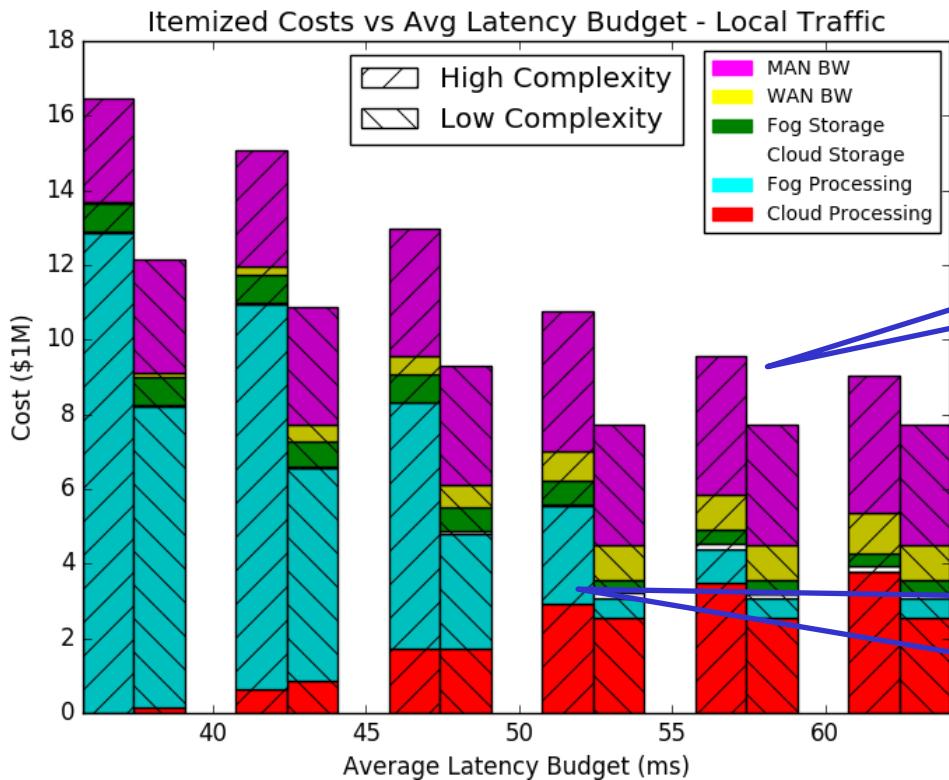


<https://cloud.google.com/compute/pricing>

<https://cloud.google.com/storage/pricing#pricing-example-simple>

<http://drpeering.net/white-papers/Internet-Transit-Pricing-Historical-And-Projected.php#>

## Results: latency effect

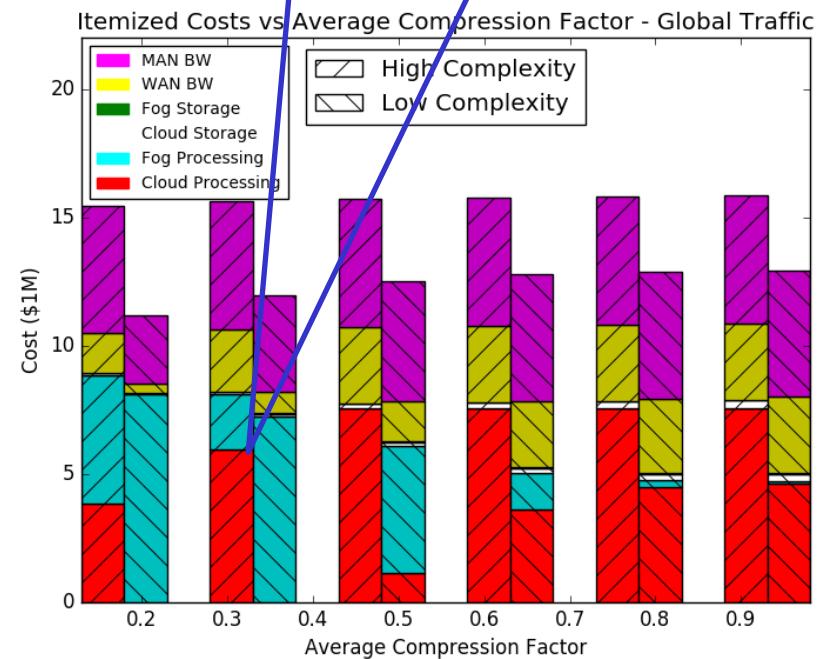
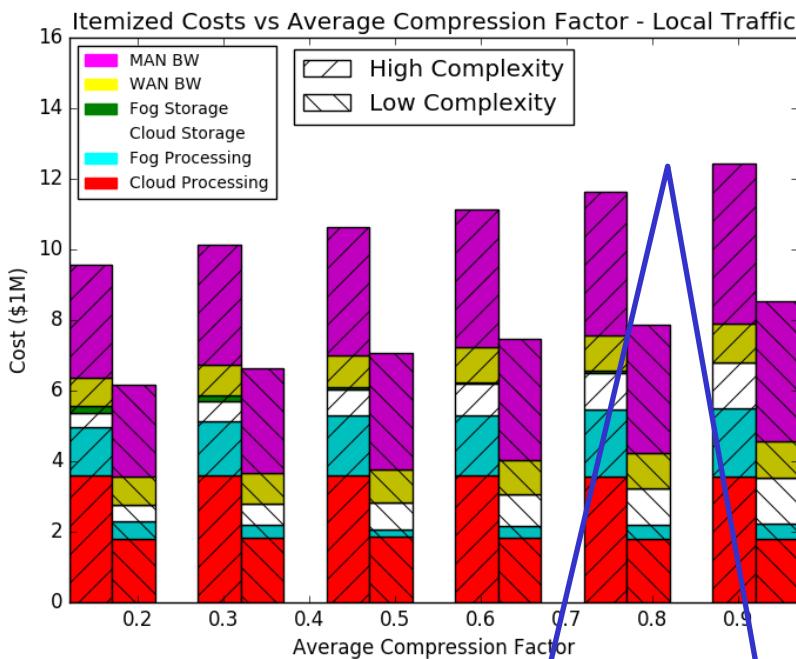


**Take-Away 1:** Total cost decreases and stabilizes as application traffic migrates to cloud thanks to less restrictive latency budget

**Take-Away 2:** With increasing latency budgets, cloud processing for high complexity applications increases faster than for low complexity apps

## Simulation Results (cont.)

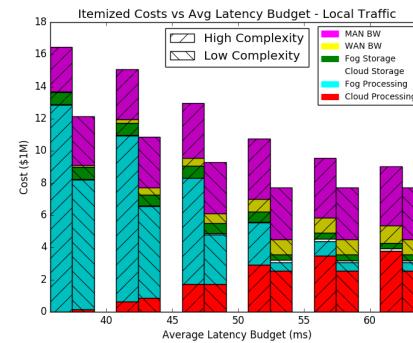
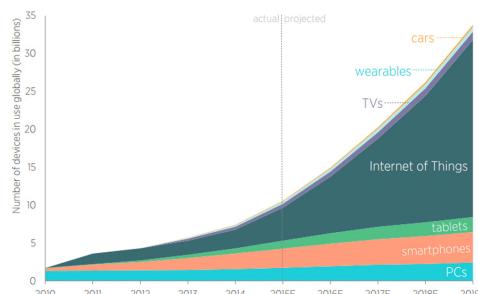
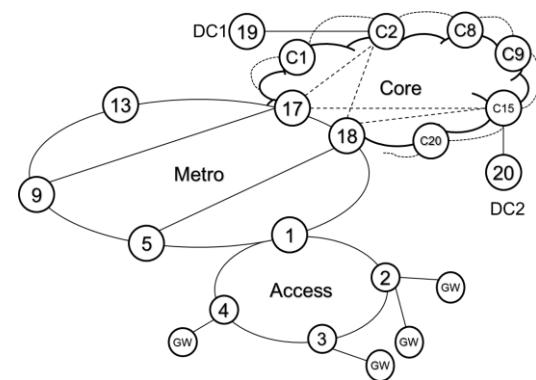
**Take-Away 4:** Total cost is stable, but processing location depends on compression factor!



**Take-Away 3:** For local traffic, cost increases with computation factor due to increase of WAN and MAN BW. **Processing location is unaffected**

# Conclusion

- Motivation
  - Lack of quantitative analysis of how specific application traffic affects resource provisioning
- Modeling work
  - a parameterized application profile:  $A = A_p \cup A_s \cup A_{sp} \cup A_n$ 
    - $\Theta, \kappa, \alpha, \beta, \Lambda$
  - Network costs for 4-tier hybrid fog-cloud architecture
- Developed a model for resource assignment
  - High flexibility: decoupling of storage and computing
- Simulation Results
  - Show impact not only of latency&BW, but also other aspects (compression factor, etc...)

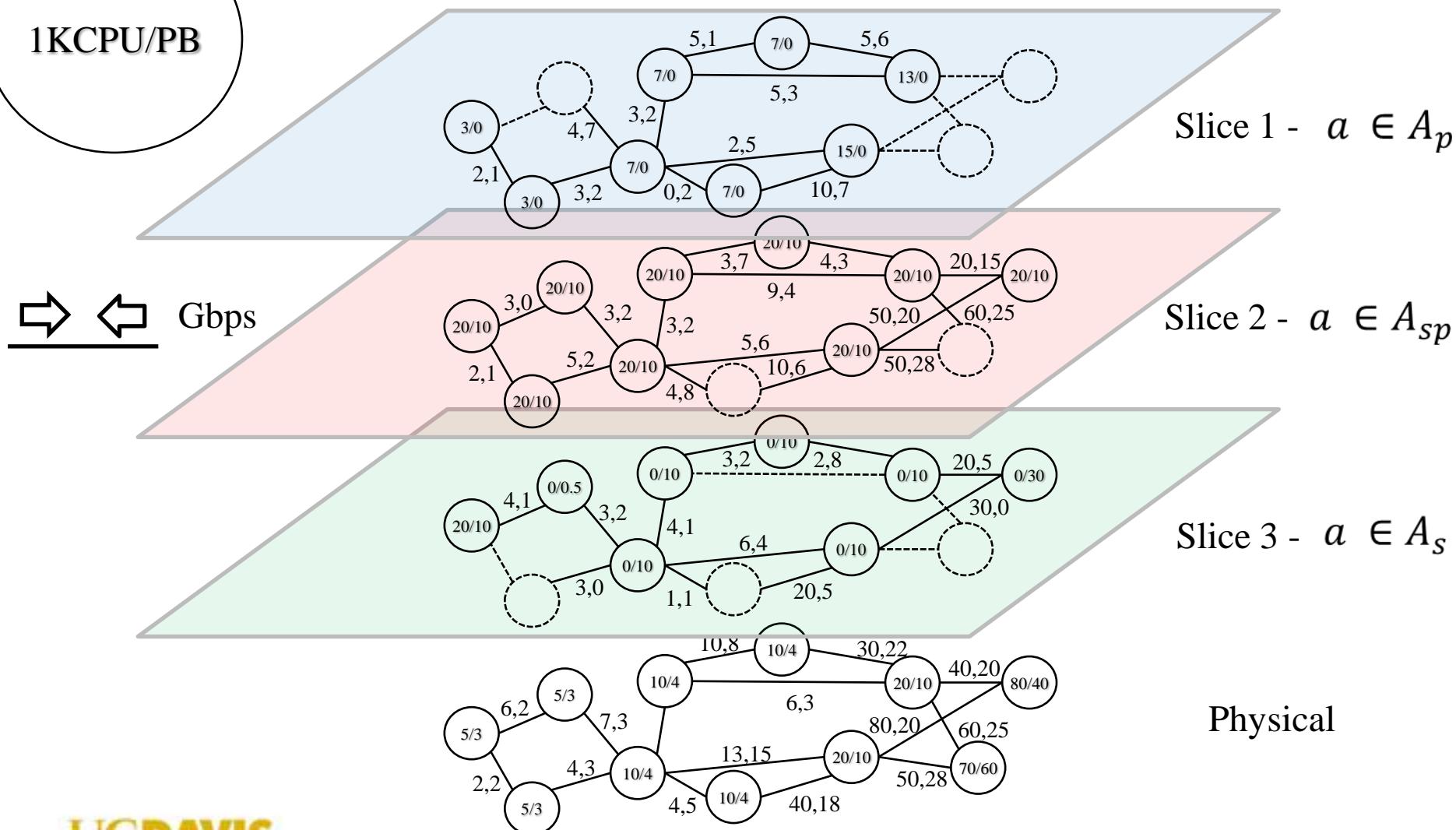




[massimo.tornatore@polimi.it](mailto:massimo.tornatore@polimi.it)

1KCPU/PB

## Slice Per Application



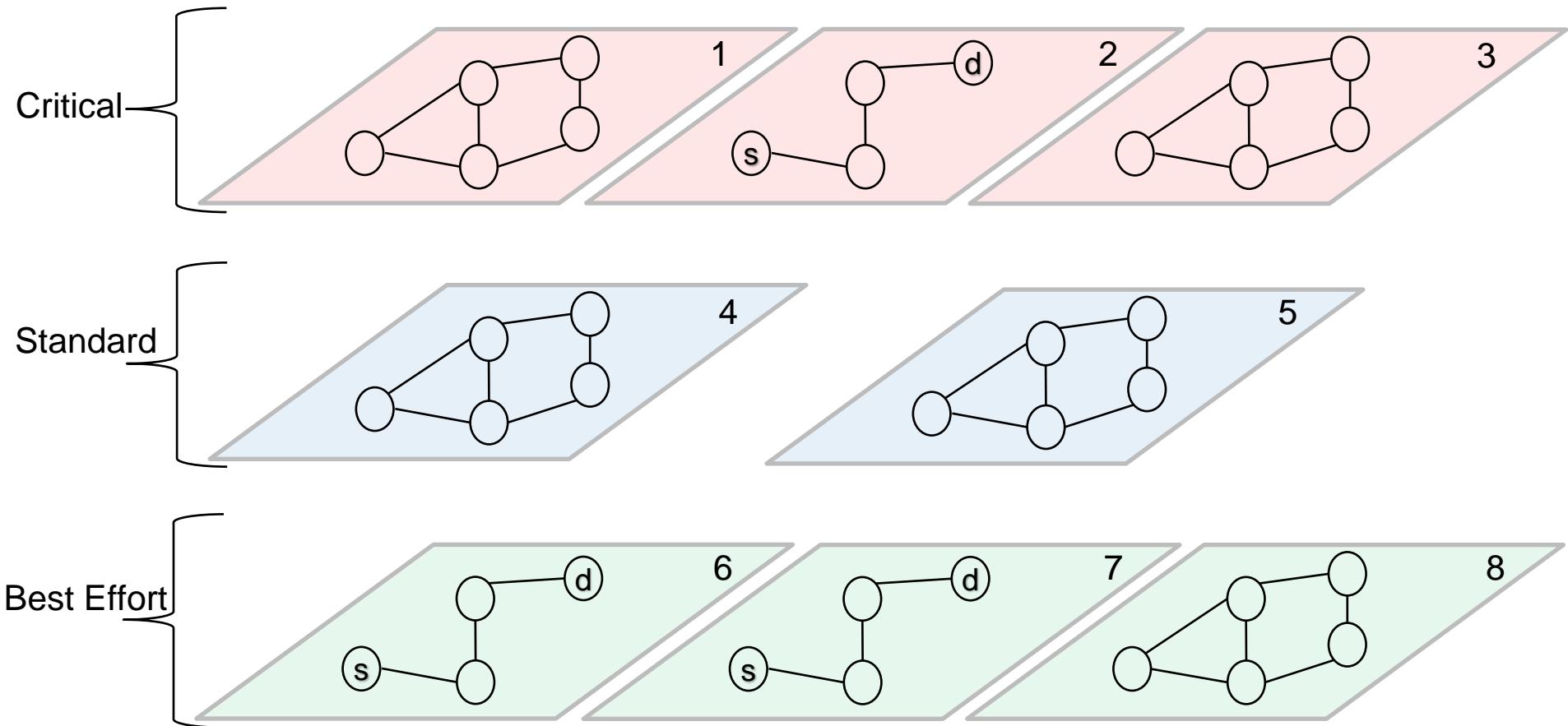
# Slice Priority

- Previous works categorize IoT/M2M slices/usage scenarios as:
  - Ultra-reliable and low latency communications (URLLC): autonomous driving, emergency services, automated manufacturing, remote medical surgery
  - Enhanced Mobile Broadband (eMBB): streaming video, high capacity multimedia, AR/VR
  - Massive Machine Type Communication (mMTC): (low power) sensor networks, smart metering, city, home (huge number of devices), less latency constrained
- Specific applications with parameterized profiles are assigned a slice of resources, which is then prioritized in a certain class
- Critical – Emerg. services, life/health/safety, remote surgery, auto. driving, factory automation/actuation
- Standard – AR/VR, gaming, Pokemon, smart grid/metering
- Best Effort – sensor data with no real-time actuation

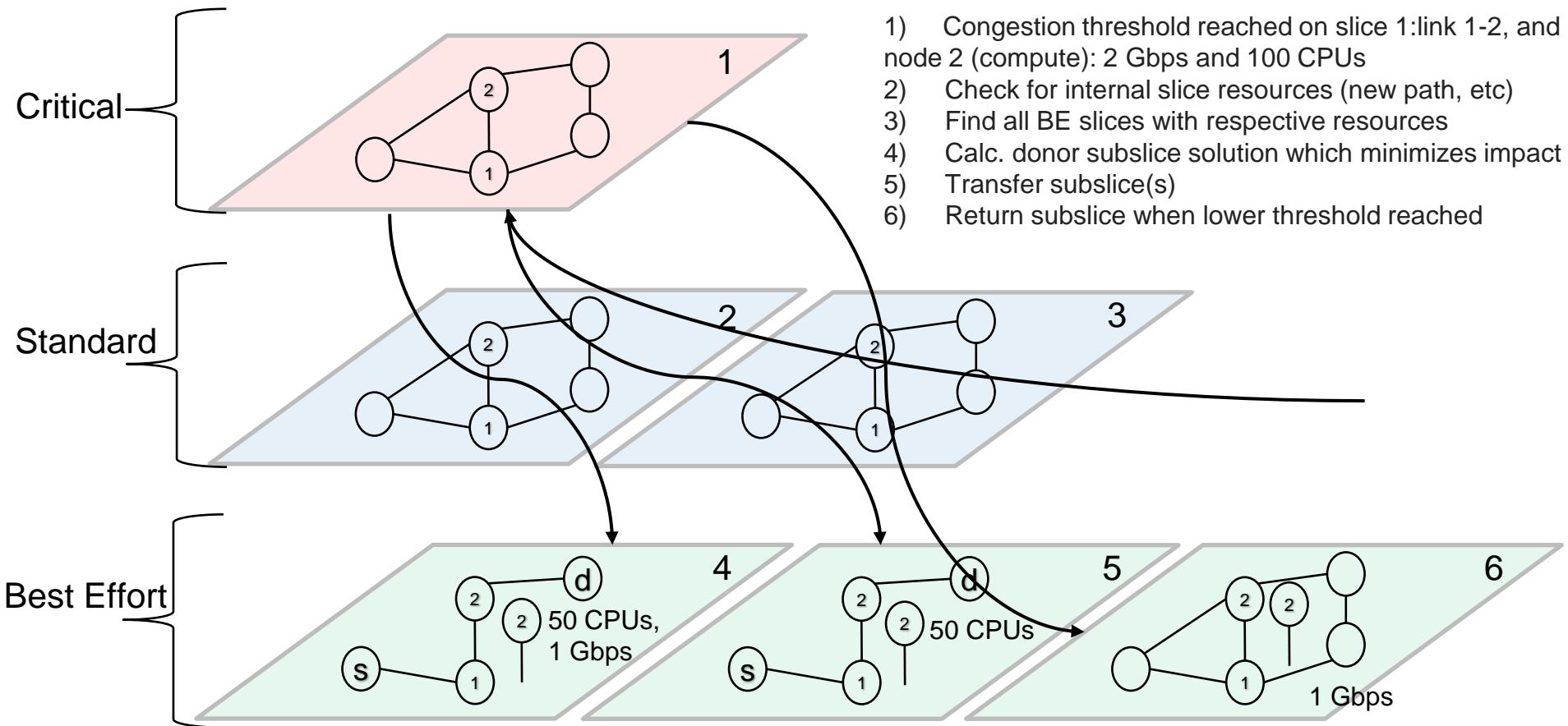
Nakao, A., Du, P., Kiriha, Y., Granelli, F., Gebremariam, A.A., Taleb, T. and Bagaa, M. End-to-End Network Slicing for 5G Mobile Networks. *Journal of Information Processing*, 25, pp.153-163, 2017.

The Fifth Generation Mobile Communication Forum (5GMF) White Paper. “5G Mobile Communications for 2020 and Beyond.” July, 2016.

## Slice Priority (cont.)



# Reslicing





# Questions

# IoT/M2M Application Popularity

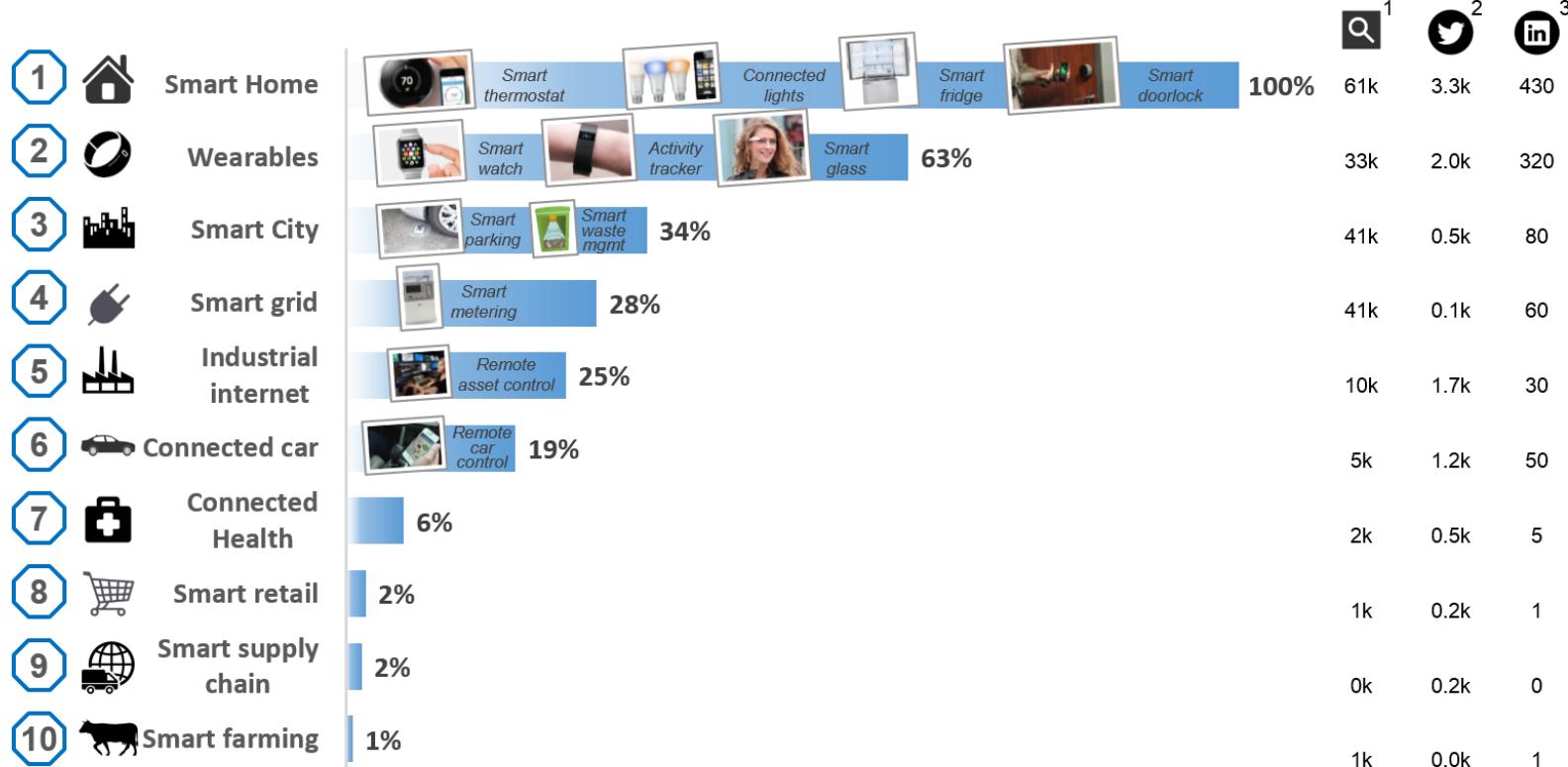


IoT Analytics – Quantifying the connected world

## Applications

## Overall popularity (and selected examples)

## Scores



1. Monthly worldwide Google searches for the application 2. Monthly Tweets containing the application name and #IOT 3. Monthly LinkedIn Posts that include the application name. All metrics valid for Q4/2014.  
Sources: Google, Twitter, LinkedIn, IoT Analytics

## Hybrid Fog-Cloud Architecture – Hierarchical

$P_t$ : Tier t processing

$S_t$ : Tier t storage

$P_4S_4$

$P_3S_3$

$P_2S_2$

$P_1S_1$

Capacity (+)

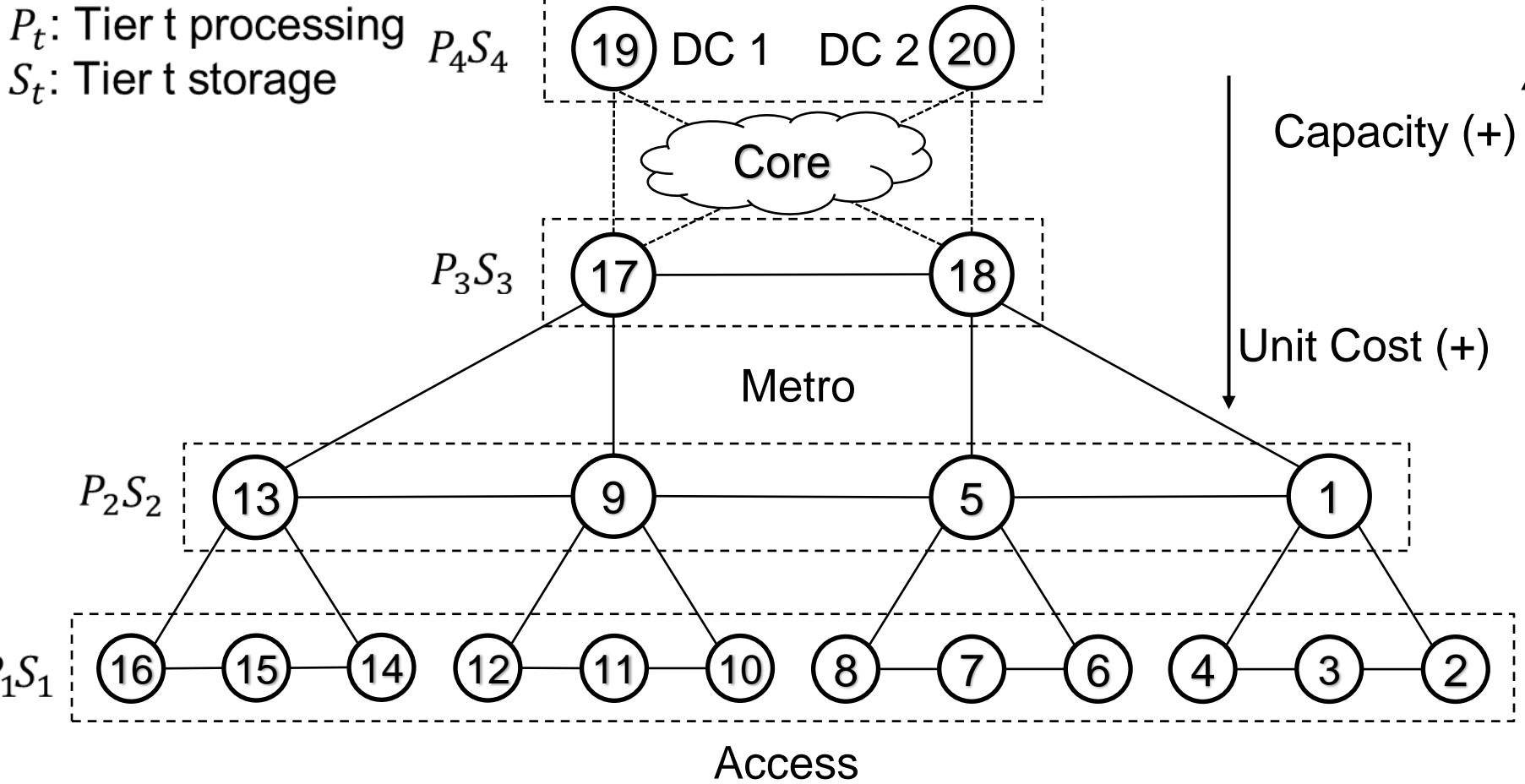
Unit Cost (+)

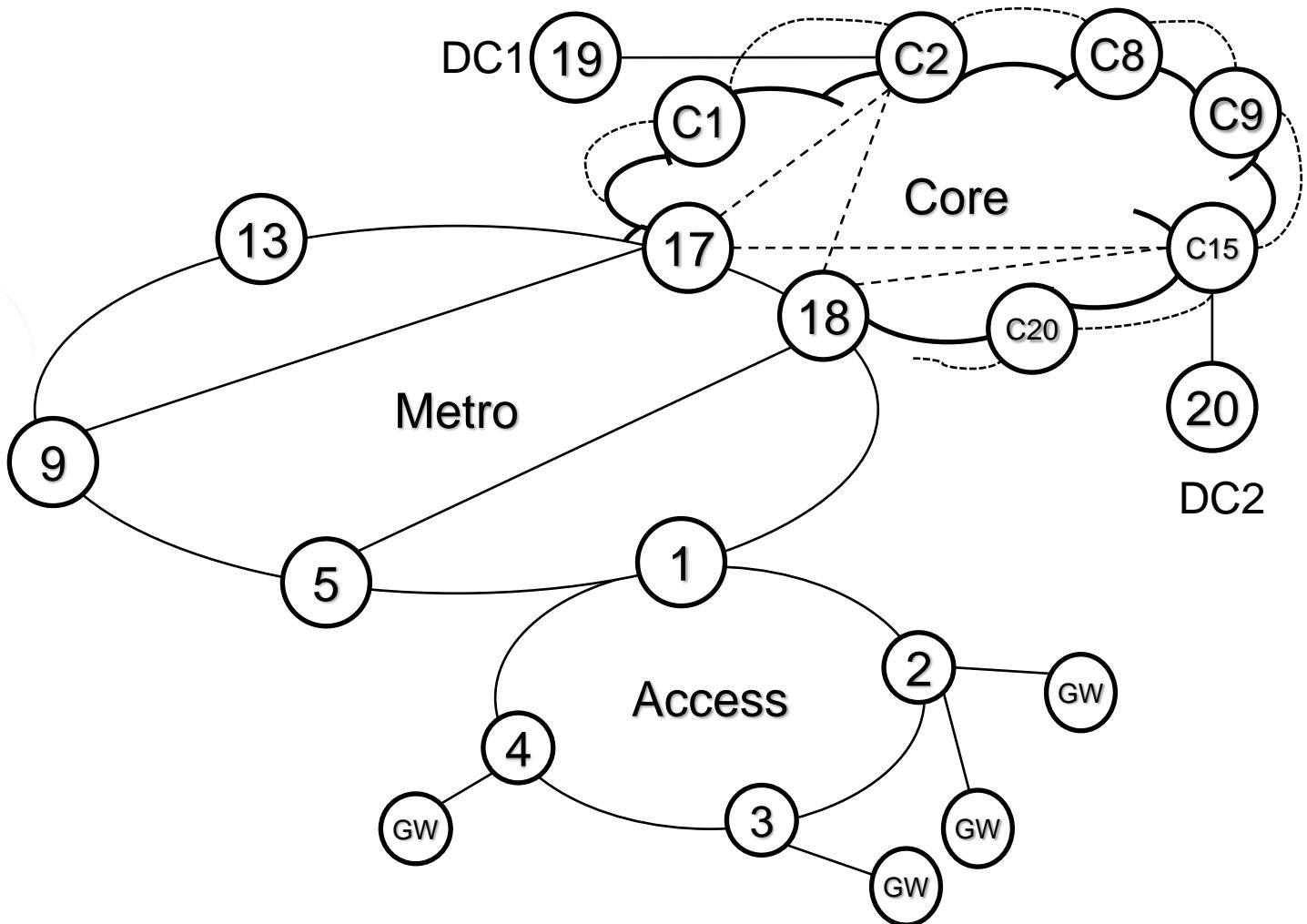
Access

Metro

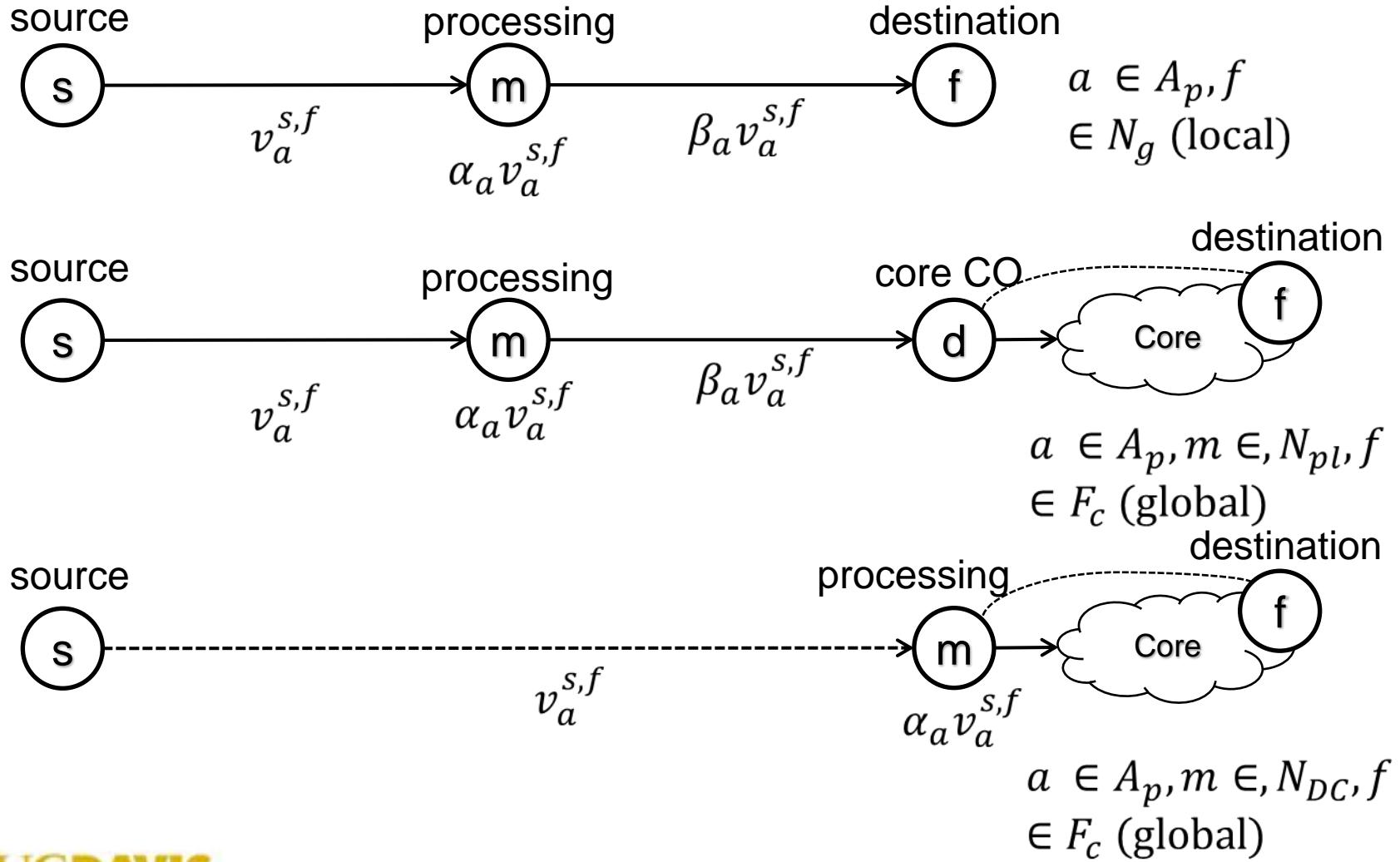
DC 1 DC 2

Core

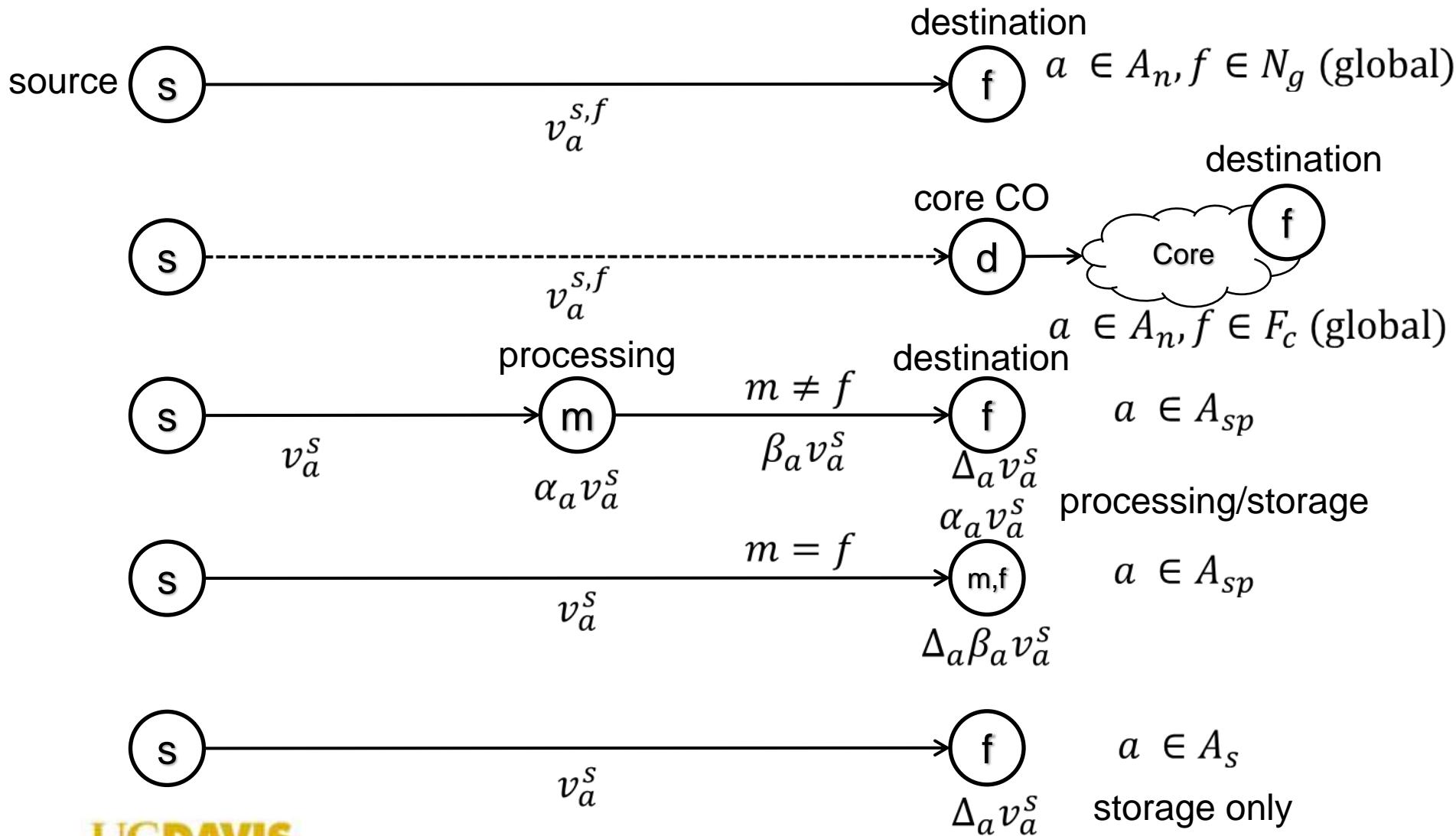




## Functional Scenarios

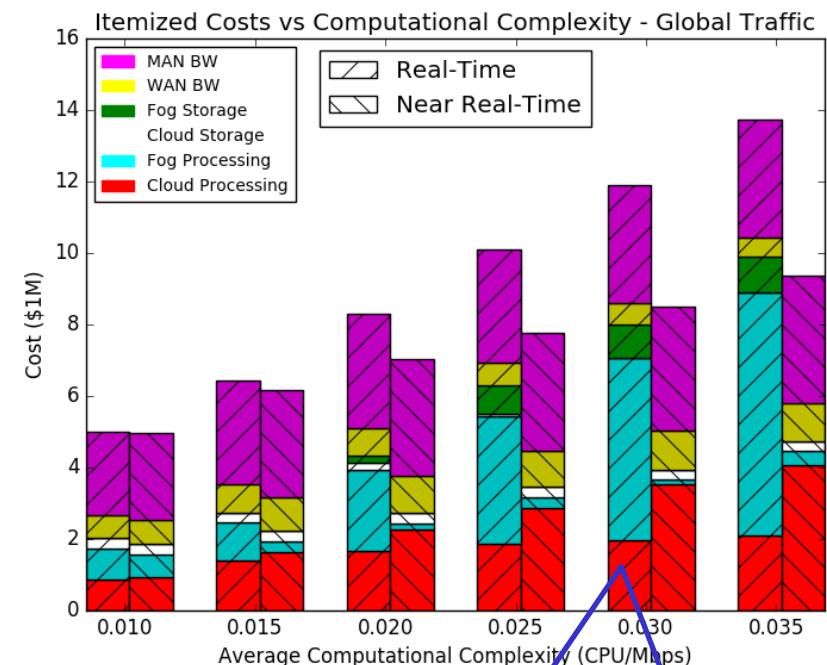
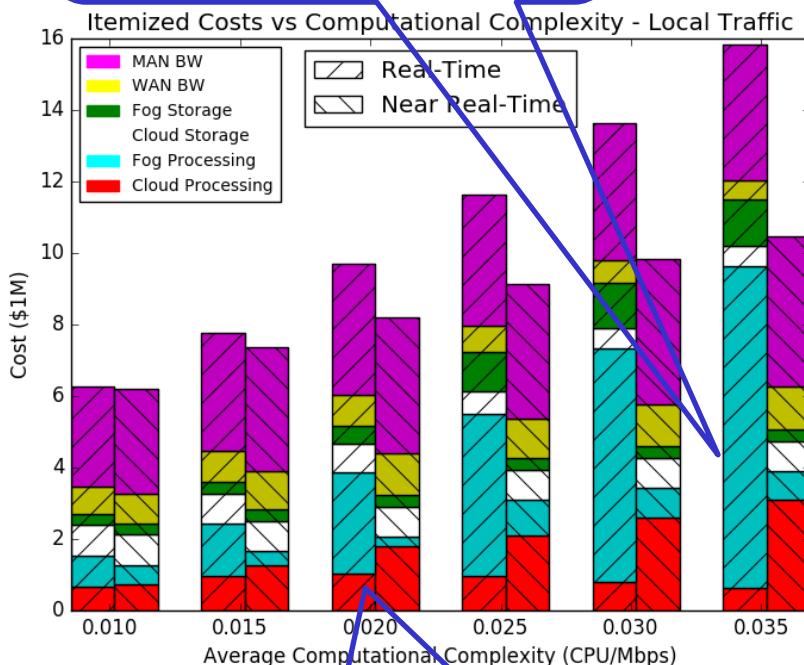


## Functional Scenarios (cont.)



# Take-Away Results: Effect of Computation complexity

for increasing computational complexity and real time services, fog processing costs steadily increases (no other choice!)



**Take-Away 2:** Cloud processing costs of real-time traffic start to decrease at .02 while near real-time cloud processing costs continue to increase with complexity

**Take-Away 1:** Cloud processing costs increase at much slower rate with increasing complexity for real-time traffic as DC compute locations restrict more RT traffic to fog processing