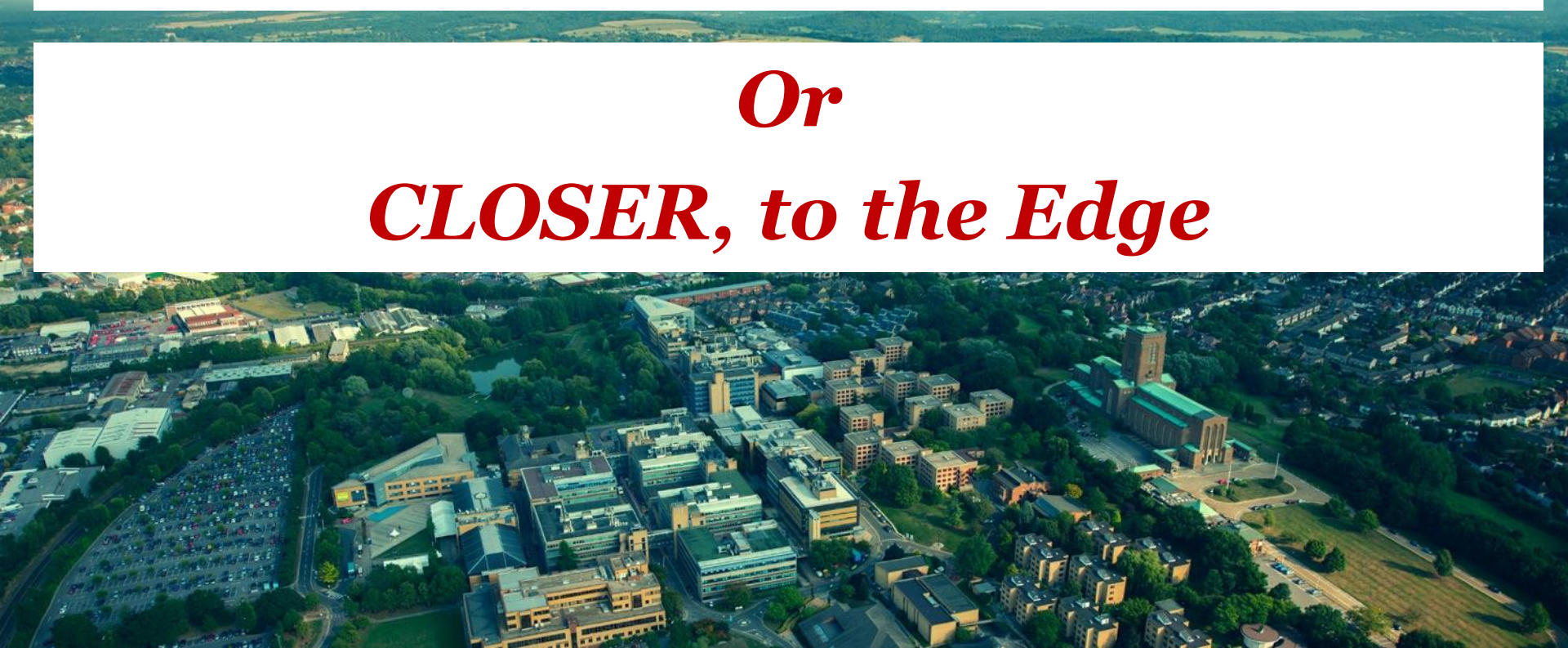


**CLOSER 2018 Keynote:
Lee Gillam, University of Surrey**

Will Cloud Gain an Edge?

Or

CLOSER, to the Edge



‘Traditional’ Cloud

- NIST SP800-145 (aka Mell and Grance): 3-4-5 & SPI
- Large (economically efficient, easily maintained – but still expensive) datacentres in relatively few, geographical locations (regions) to support large user numbers, *centralized* corporate entity
- A ‘Big Four’ in **Amazon**, Microsoft, IBM, Google

‘New’ Cloud

- Containers (*Docker, kernel-locked*), and Functions (*multiplicity of approaches*) – added “CaaS” and “FaaS”
- **Edge** (*multiplicity of approaches and concerns*)
- (Re-)distributed Computing, and ‘new’ problems (new ‘traditional’ problems)
- ‘Big Four’?

'Traditional' Cloud

(Big Four) Clouds are big

Cost and performance (=cost) variation

Performance variation and implications for energy efficiency

'New' cloud

'serverless' and performance

Multiplicity of Edges

'serverless' Edges

An application

Cloud Cars and exemplars

Summary and take home

‘Traditional’ Cloud

(Illusion of) Infinite capacity - consider one (big) provider:

AWS: 2014 based on 11 regions and 28 AZs: 2.8-5.6m servers (Morgan, 2014)
based on datacenter of up to 80k servers

12 Jun 2012, 1 trillion objects in S3.

13 April 2013, 2 trillion

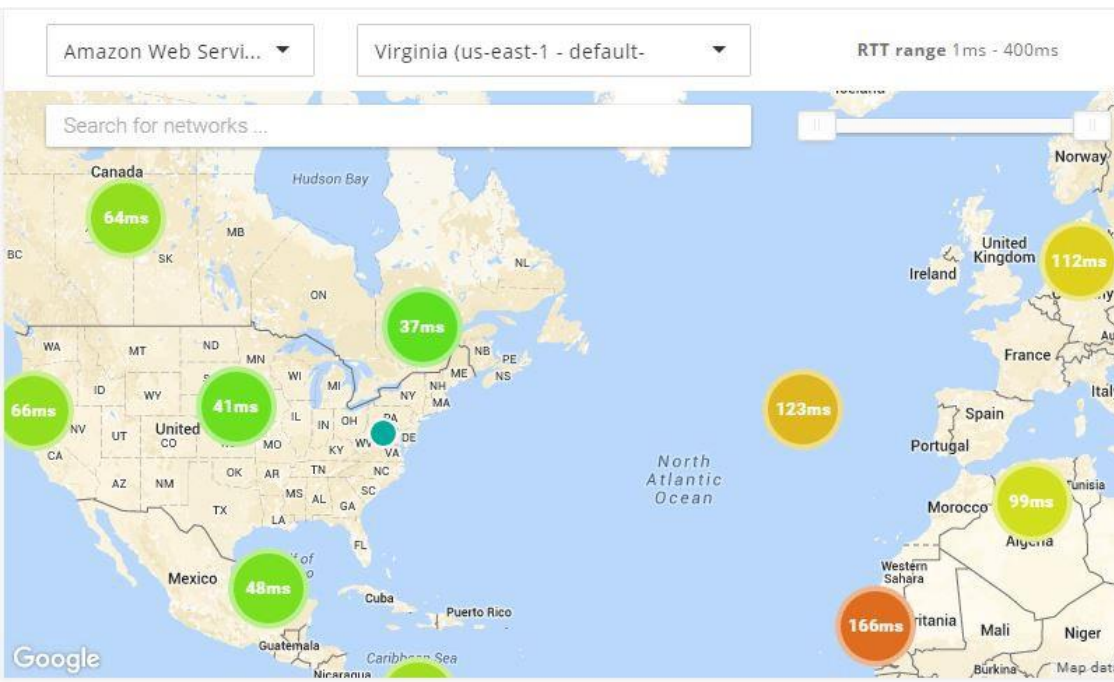
2018, **18** regions, **54** AZs, 5 more regions coming: ~**10m** servers?

T. P. Morgan, A rare Peek Intro The Massive Scale of AWS, 2014,
<https://www.enterprisetech.com/2014/11/14/rare-peek-massive-scale-aws/>

Clouds are 'big'

Regions: Contain big datacenters at distance but with big networking also

Edges: a shorter distance to something that does something useful for compute



Comparison of significant cloud data centres

Sq Footage

| | | |
|-----------|-----------------|---------|
| Google | Lenoir, NC | 476,000 |
| | Dallas, OR | 206,000 |
| Apple | Apple, NC | 500,000 |
| | Chicago, IL | 700,000 |
| Microsoft | San Antonio, TX | 470,000 |
| | Lockport, NY | 190,000 |
| YAHOO! | La Vista, NE | 350,000 |

Manchester United pitch ~80,000 sq ft

Figure sources: Datapath.io and Make IT Green: Cloud Computing and its Contribution to Climate Change, Greenpeace

Clouds are 'big' – with cost variation

| 27/4/16 | vCPU | ECU | Mem (GiB) | US-E (NV) | EU-W (Ire) | EU-W (Fra) | SA (SP) |
|-------------|------|----------|-----------|-----------|------------|------------|----------|
| t2.nano | 1 | Variable | 0.5 | \$0.0065 | \$0.007 | \$0.0075 | \$0.0135 |
| t2.micro | 1 | Variable | 1 | \$0.013 | \$0.014 | \$0.015 | \$0.027 |
| m4.xlarge | 4 | 13 | 16 | \$0.239 | \$0.264 | \$0.285 | N/A |
| m4.2xlarge | 8 | 26 | 32 | \$0.479 | \$0.528 | \$0.57 | N/A |
| m4.4xlarge | 16 | 53.5 | 64 | \$0.958 | \$1.056 | \$1.14 | N/A |
| m4.10xlarge | 40 | 124.5 | 160 | \$2.394 | \$2.641 | \$2.85 | N/A |
| m3.medium | 1 | 3 | 3.75 | \$0.067 | \$0.073 | \$0.079 | \$0.095 |
| m3.large | 2 | 6.5 | 7.5 | \$0.133 | \$0.146 | \$0.158 | \$0.19 |
| m3.xlarge | 4 | 13 | 15 | \$0.266 | \$0.293 | \$0.315 | \$0.381 |
| m3.2xlarge | 8 | 26 | 30 | \$0.532 | \$0.585 | \$0.632 | \$0.761 |
| c4.large | 2 | 8 | 3.75 | \$0.105 | \$0.119 | \$0.134 | N/A |
| c4.xlarge | 4 | 16 | 7.5 | \$0.209 | \$0.238 | \$0.267 | N/A |
| c4.2xlarge | 8 | 31 | 15 | \$0.419 | \$0.477 | \$0.534 | N/A |
| c4.4xlarge | 16 | 62 | 30 | \$0.838 | \$0.953 | \$1.069 | N/A |
| c4.8xlarge | 36 | 132 | 60 | \$1.675 | \$1.906 | \$2.138 | N/A |
| c3.large | 2 | 7 | 3.75 | \$0.105 | \$0.12 | \$0.129 | \$0.163 |
| c3.xlarge | 4 | 14 | 7.5 | \$0.21 | \$0.239 | \$0.258 | \$0.325 |
| c3.2xlarge | 8 | 28 | 15 | \$0.42 | \$0.478 | \$0.516 | \$0.65 |
| c3.4xlarge | 16 | 55 | 30 | \$0.84 | \$0.956 | \$1.032 | \$1.3 |
| c3.8xlarge | 32 | 108 | 60 | \$1.68 | \$1.912 | \$2.064 | \$2.6 |

Clouds are 'big' – with hardware variation

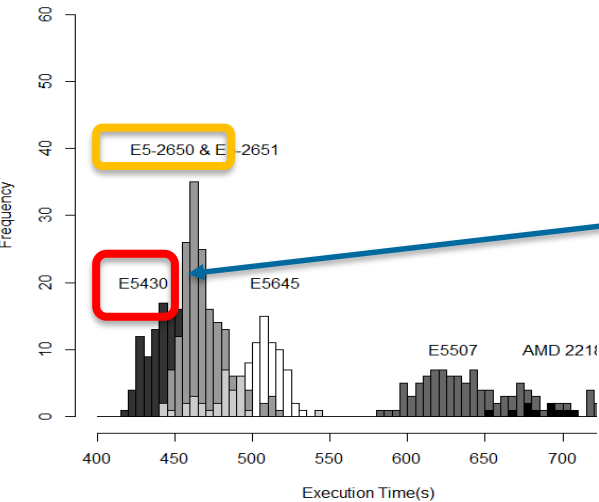
CPU model discovery - for ~700 EC2 FGS Instances for 1 user – 1 instance type

| Region | AZ | <i>E5430</i> | <i>E5-2650</i> | <i>E5645</i> | <i>E5507</i> |
|--|------------------------|---------------------|-----------------------|---------------------|---------------------|
| US East N. Virginia 2006 [year Region started] <i>Cheapest – but latencies</i> | <i>us-east-1a</i> | 31% | 0 | 25% | <u>44%</u> |
| | <i>us-east-1b</i> | 5% | <u>59%</u> | 29% | 7% |
| | <i>us-east-1c</i> | 0 | 47% | <u>52%</u> | 1% |
| | <i>us-east-1d</i> | 18% | 31% | <u>44%</u> | 7% |
| EU West Dublin 2007 | <i>eu-west-1a</i> | 4% | <u>75%</u> | 19% | 2% |
| | <i>eu-west-1b</i> | 28% | 0 | <u>44%</u> | 28% |
| | <i>eu-west-1c</i> | 4% | 0 | <u>63%</u> | 33% |
| US West N. California 2009 | <i>us-west-1b</i> | 0 | 0 | 13% | <u>87%</u> |
| | <i>us-west-1c</i> | 8% | 0 | 18% | <u>74%</u> |
| SA San Paulo 2011 | <i>sa-east-1a</i> | 0 | <u>81%</u> | 19% | 0 |
| | <i>sa-east-1b</i> | 0 | <u>86%</u> | 14% | 0 |
| US West Oregon 2011 | <i>us-west-2b</i> | 0 | <u>73%</u> | 27% | 0 |
| Asia Pacific Sydney 2012 | <i>ap-southeast-2a</i> | 0 | <u>64%</u> | 36% | 0 |
| | <i>ap-southeast-2b</i> | 0 | <u>75%</u> | 25% | 0 |

Cost, latency, computational capability (moving up stack)

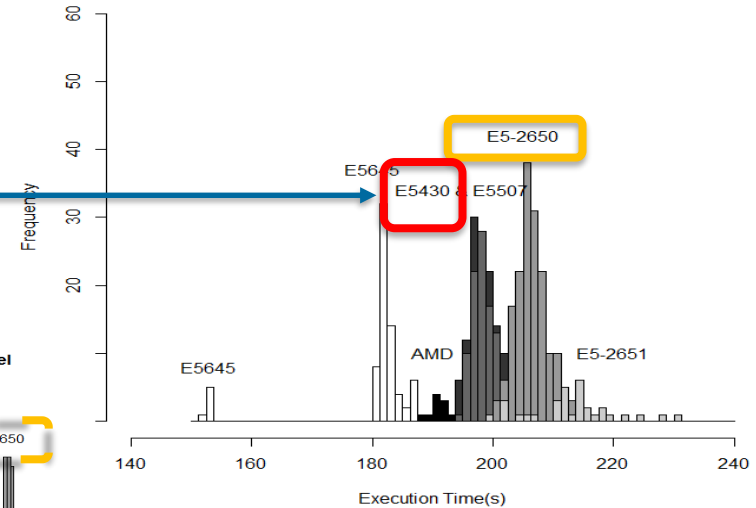
Performance varies by CPU

Histogram of bzip2 Results By CPU Model



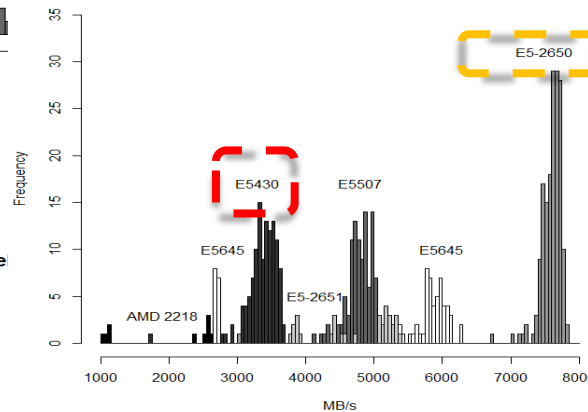
int – left better

Histogram of GNUGO Results By CPU Model

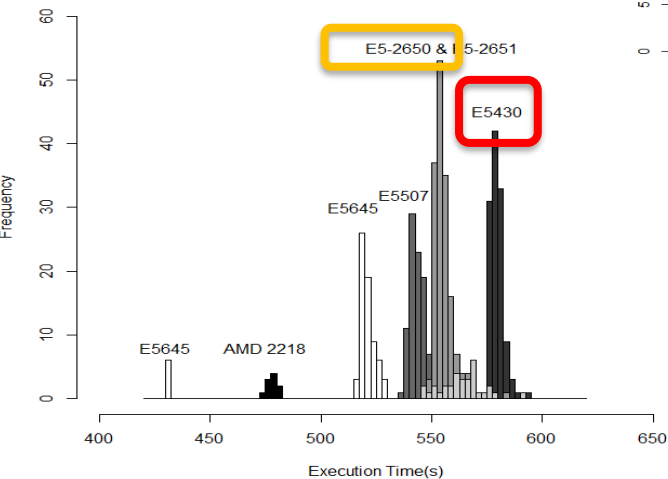


STREAM – right better

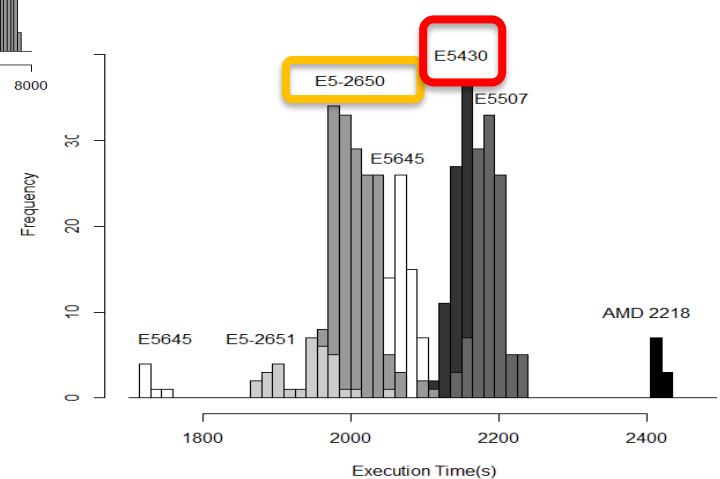
Histogram of STREAM Triad Results By CPU Model



Histogram of POV-Ray Results By CPU Model



Histogram of NAMD Results By CPU Model



fp – left better

Cost-efficient use

Heterogeneous hardware complicates costs

(at minimum) a user needs to:

- Identify suitable (cost-based?) instance offerings (**price determination**)
 - Rank 'best' by workload (**performance determination**)
 - Determine AZs (latency) offering those resources (**location selection**)
 - Attempt to obtain them (**instance lottery**)
 - For users with more than one account this may need to be done per account basis (**account selection**)
- **Costs** are incurred in (1) **performance determination** (2) **location selection** and (3) **instance lottery**, for every (4) **account selection** (*AND infrastructure composition changes over time*)

Spend involved with getting cost-efficiency (performance gaming /
deploy-and-ditch) - potentially expensive!

Can't eliminate resource uncertainty

Performance can be 'stable' over a long period for a given benchmark – past a good indicator of future - but may be subject to abrupt changes and severe degradation

One instance, 1379s for POV-Ray – ~13 standard deviations from the mean (639, 54)

Rarer: **'The requested Availability Zone is currently constrained and we are no longer accepting new customer requests for X/Y/Z instance types'**

- Go elsewhere, but other AZs may not be cost-efficient – AZ lock-in

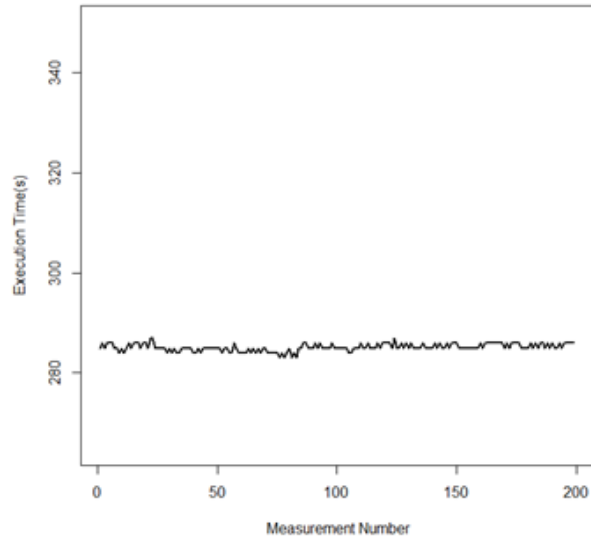
Unusually for a service: better can be cheaper

But much work around performance continues to assume homogeneity

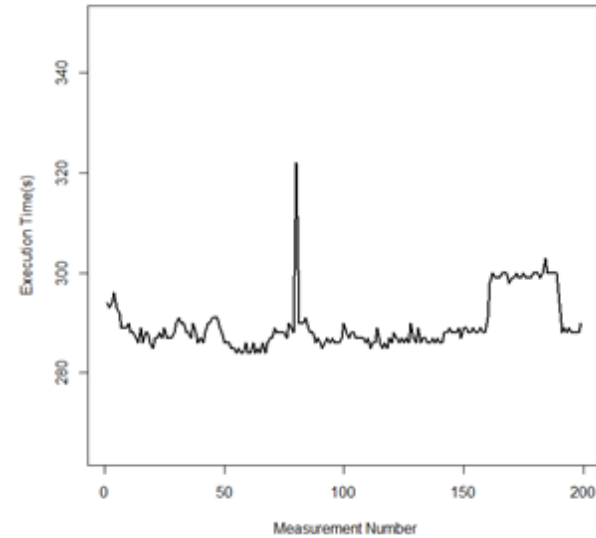
John O.Loughlin and Lee Gillam (2014) "Should Infrastructure Clouds be Priced Entirely on Performance? An EC2 Case Study". International Journal of Big Data Intelligence

Cost-efficient use

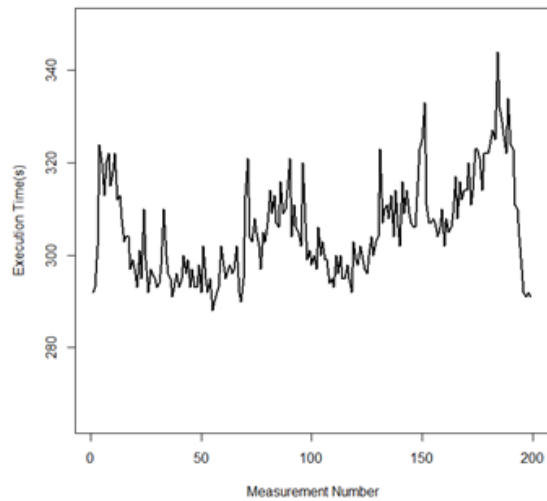
Stationary Time Series, mean = 285s, sdev = 0.7s



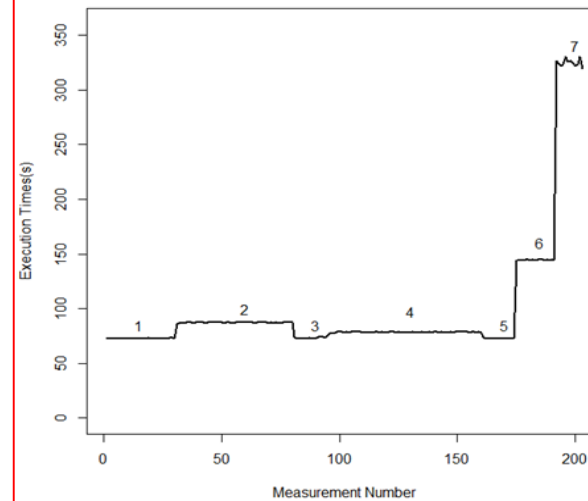
Non-Stationary Time Series



Non-Stationary Time Series



Timeplot of pbzip2



An aside – ‘brokers’ seem popular

Cloud Service Brokers (e.g. aggregators) might address performance issues

Instances with known performance characteristics

- A match-making service between user application performance needs and available resources

Re-price based on desired performance

- Make more suitable instances more expensive, and less suitable less so.

Extensive simulations suggested:

- Assuming clouds are opaque makes it difficult to avoid instance gaming.
- **Very difficult to make a profit**, even with careful pool management! - high vol.
- Opportunities in value of utility rather than price

Rare to find discussion of broker profit – Rogers & Cliff has been a notable exception

O'Loughin, J. (2018): A Workload-Specific Performance Brokerage for Infrastructure Clouds (unpub PhD thesis).
Rogers, O. & Cliff, D., 2012. A Financial Brokerage Model for Cloud Computing. Journal of Cloud Computing: Advances, Systems and Applications, 1(2)

Performance and energy trade-off for different kinds of workload

- *runtime variable with hardware (heterogeneity)*
- how much power needed to deliver runtime on given hardware
- best performance might not equate to most energy efficient
- performance \rightarrow runtime; user cost higher with longer runtime

Put workloads on best machines for it: consolidation (implies migration)

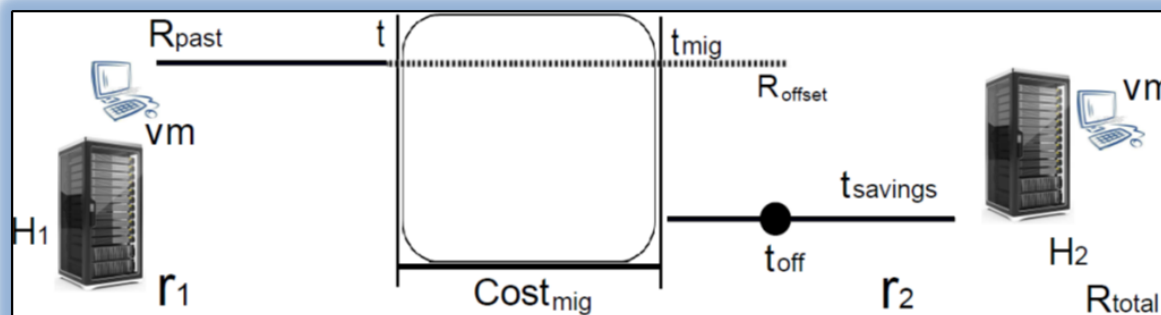
- **risk of being own noisy neighbour** \rightarrow longer runtimes
- additional energy use for period of migration: at least costs of equivalent resources plus network
- question of recoverability depends on continued use
- for providers, opportunity to switch off / maintain (may not be an incentive to)

Workload-related CPU model ranking

- $w1: A > B > C > D$; $w2: D > C > B > A$
- VM allocation: B allocated as available; A preferred
- VM consolidation: migration beneficial if workload can recoup cost of migration – implies performance maintained (contention)

Consolidation with Migration Cost Recovery (CMCR)

- Migrate to more efficient hosts
- VM terminated before $[t_{\text{off}}]$, the effort is wasted
- Recover migration overhead $[\text{Cost}_{\text{mig}}]$, efficient gain after $[t_{\text{off}}]$

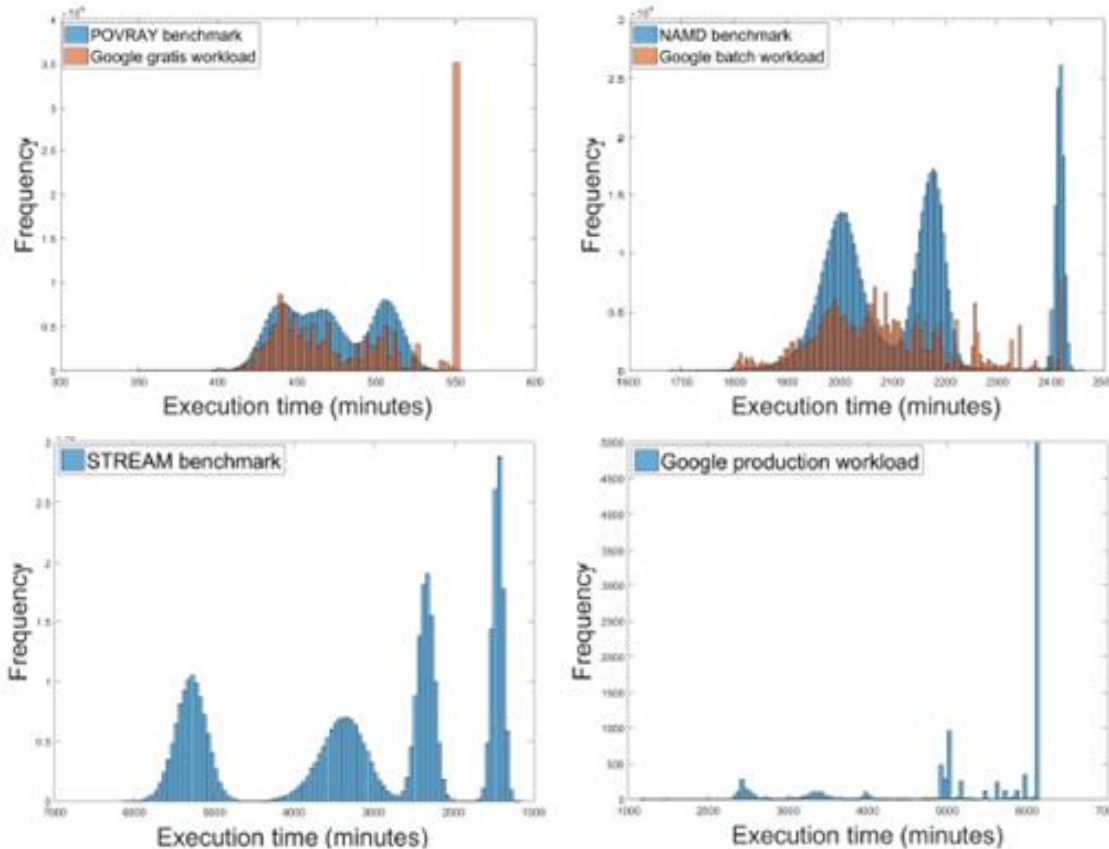


Broad characterisation: **9** scheduling approaches, several types of consolidation (+ none), CloudSim model using **12,583** heterogeneous hosts, **25m** VMs (tasks) in Google workload trace data, **5** minimum-runtime settings, migration rounds at 5 minute intervals (host utilization < 20%). **On-demand** VM allocation

CPU info not provided, so map Google priorities to benchmark results (range scaled) for preferences (Gratis (0) : POVRAY, Batch (2) : NAMD, Production (9) : STREAM). *Distributions typically skewed lognormal per CPU model.*

Then relate power ratings of these 'Cloud' CPUs – simulation results can be related to real VMs.

| Workload | Bench mark | CPU Model | Real benchmarks | | | | | Google data | | | | |
|------------|------------|-----------|-----------------|--------------|------|------|-------|-------------|--------------|------|------|-------|
| | | | (μ) | (σ) | Min | Max | CoV | (μ) | (σ) | Min | Max | CoV |
| Gratis | POVRAY | E5430 | 439 | 11 | 421 | 467 | 0.025 | 438.06 | 9.42 | 421 | 467 | 0.022 |
| | | E5-2650 | 468 | 12 | 451 | 500 | 0.026 | 473.87 | 11.93 | 451 | 500 | 0.025 |
| | | E5645 | 507 | 10 | 490 | 535 | 0.02 | 498.55 | 10.44 | 490 | 535 | 0.021 |
| Batch | NAMD | E5-2651 | 1994 | 41.9 | 1952 | 2036 | 0.021 | 1991 | 39.51 | 1800 | 2040 | 0.02 |
| | | E5-2650 | 2007 | 28.5 | 1978 | 2036 | 0.014 | 1963.4 | 28.41 | 1900 | 2015 | 0.015 |
| | | E5645 | 2043 | 96.4 | 1946 | 2140 | 0.047 | 1931.4 | 93.43 | 1800 | 2170 | 0.048 |
| | | E5430 | 2160 | 20.7 | 2135 | 2189 | 0.01 | 2103.6 | 22.1 | 2080 | 2150 | 0.011 |
| Production | STREAM | E5507 | 2187 | 18.1 | 2162 | 2217 | 0.008 | 2191.8 | 15.69 | 2150 | 2200 | 0.007 |
| | | E5430 | 1446 | 66 | 1328 | 1572 | 0.045 | 1404.4 | 44.33 | 1328 | 1572 | 0.032 |
| | | E5507 | 2348 | 104 | 2078 | 2448 | 0.044 | 2346.7 | 107.21 | 2078 | 2448 | 0.046 |
| | | E5645 | 3395 | 287 | 2995 | 4008 | 0.085 | 3388.7 | 238.22 | 2995 | 4008 | 0.07 |
| | | E5-2650 | 5294 | 191 | 4935 | 5860 | 0.036 | 5294.5 | 197.52 | 4935 | 5860 | 0.037 |



Findings confirmed: *sensible* allocation better (easier) than consolidation; migrate longer running VMs – but **assumes clouds are not opaque**. (i.e. provider has knowledge of workload)

Zakarya, M. (2017): A Workload-Specific Performance Brokerage for Infrastructure Clouds (unpub PhD thesis) .

‘New’ Cloud

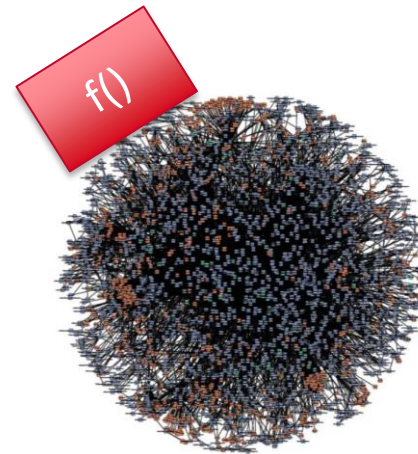
Relevance to 'new' Cloud

A look at so-called 'Function as a Service'

'serverless' computing (yet servers are **essential**)

- You're not supposed to "worry" about provisioning
- Billing per 100 milliseconds (AWS Lambda, Google Cloud Functions; Azure has at least 2 ways to pay incl. based on memory consumption)
- Functions may be time-constrained – 5 mins, Lambda/Azure; though HTTP timeout of 30s (e.g. AWS API Gateway) gives 29s runtime

Are performance/cost questions relevant?



Relevance to 'new' Cloud - performance

A look at so-called 'Function as a Service' – AWS Lambda

AWS Lambda runs a Function in a Container on a VM ('serverless')

IP address may over time change – 2 functions run gave e.g.: ip-10-23-17-3, ip-10-14-98-122.

Short runtimes good: a small test - 113.44 ms, then 114.34 ms, 102.65 ms, 113.57 ms – all rounded to nearest 100 (200ms) for billing; reasonable consistency.

Underlying process:

| USER | PID | %CPU | %MEM | VSZ | RSS | TTY | STAT | START | TIME | COMMAND |
|---|-----|------|------|--------|-------|-----|------|-------|------|---------|
| 490 | 1 | 1.3 | 0.3 | 212024 | 15372 | ? | Ss | 16:49 | 0:00 | |
| /usr/bin/python2.7 /var/runtime/awslambda/bootstrap.py | | | | | | | | | | |
| 490 | 7 | 0.0 | 0.0 | 117208 | 2476 | ? | R | 16:49 | 0:00 | ps auxw |

Through several uses, underlying process remains.

Relevance to 'new' Cloud - performance

A look at so-called 'Function as a Service' – AWS Lambda

Limitations exist, e.g. can't run 'ifconfig', no 'sudo' so can't install as root, and can't get at AWS metadata of VM. But can find out *CPU model* (/proc/cpuinfo [dual 'core' **c4**]) and underlying system (uname [**AWS Linux**]),

```
model name      : Intel(R) Xeon(R) CPU E5-2666 v3 @ 2.90GHz
cpu MHz         : 2900.066
cache size      : 25600 KB
```

```
Linux ip-10-23-17-3 4.4.35-33.55.amzn1.x86_64 #1 SMP ...
x86_64 x86_64 x86_64 GNU/Linux
```

Others have seen:

- **Intel(R) Xeon(R) CPU E5-2680 v2 @ 2.80GHz** - a **c3** instance:
<http://zqsmm.qiniucdn.com/data/20150416152509/index.html>.

Relevance to 'new' Cloud - performance

A look at so-called 'Function as a Service' – AWS Lambda

Can run a [small –time limit!] arbitrary linux application – e.g. a benchmark such as STREAM [**~2GB/s**], if:

- Precompiled elsewhere, downloaded into local filestore for Function (/tmp), made executable (chmod), executed and output returned
- *Variations per execution, with rounding; **performance/location/lottery/account remains important***

```
-----  
STREAM version $Revision: 5.10 $  
-----
```

```
...
```

```
-----  
Function      Best Rate MB/s  Avg time     Min time     Max time  
Copy:         1979.3   0.091534     0.080836     0.116964  
Scale:        1974.5   0.099291     0.081033     0.117087  
Add:          2370.4   0.122582     0.101247     0.157255  
Triad:        2362.9    0.126896     0.101569     0.157570
```

Multiplicity of Edges

Will Cloud gain an **Edge**?

Network **Edge** devices – devices that have sensors (a mobile phone?)

Customer **Edge** (/Edge router) – router on premises

Provider **Edge** – a provider's router

Edge Datacenter (/Cloudlet* /Content Delivery Network) – datacenter

Multi-access (/mobile) **Edge** – datacenter + RAN (e.g. 5G)

Also, **FOG**

Edge is also Microsoft's web browser, and a Ford vehicle

** Should not be confused with Cloudlet in CloudSim, which is a Task*

Fog “**complements and extends** the Cloud to the **edge** and endpoints” , Bonomi et al

- Fog is **additional to and complementing Cloud**,
- Example distributed applications: “A has one **Cloud component**, and two **Fog components** [...] B has one **cloud component**, one **component in the Core**, and a **Fog component**”.

OpenFog Consortium: “a system-level horizontal architecture that distributes resources and services of computing, storage, control and networking anywhere along **the continuum from Cloud to Things**”.

- OpenFog documents represent Fog diagrammatically as:
 - **between** Cloud and Things;
 - **including** Cloud;
 - **‘in’** Cloud Computing.
- OpenFog Consortium: Fog is “often erroneously called edge computing, but there are key differences. Fog **works with** the cloud, whereas **edge is defined by the exclusion of cloud**”.

Fog has not proven an entirely helpful notion

Multi-access Edge Computing (until recently, Mobile Edge Computing) – ETSI specification

- provide **capabilities of Cloud Computing close to the Radio Access Networks** in 4G and 5G telecommunications and converge with other radio access technologies (e.g. WiFi or Satellite).
- “can be seen as **a cloud server running at the edge of a mobile network**”.

ETSI MEC server supports VMs into which “MEC applications from vendors, service providers and third-parties are deployed and executed”.

VMs → Containers → Functions
[MEC authors: PaaS for “future releases”].

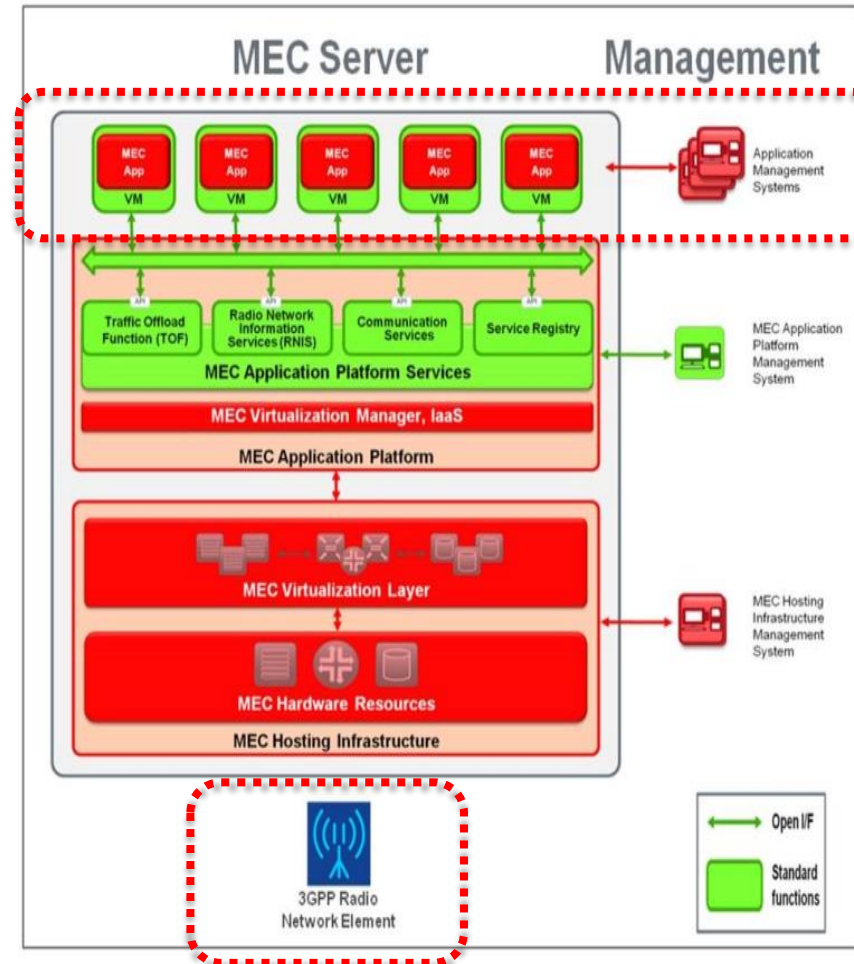


Figure 9: MEC server platform overview

Not just one MEC server per Edge?

Cloudlet (Carnegie Mellon University (CMU) notion) - middle tier of a 3-tier hierarchy: “mobile device - **cloudlet** – cloud”.

A “**datacentre in a box**”, implying multiple servers, local to the user, similar to MEC (and some characterisations of Fog).

*CMU formed Open Edge Computing initiative in 2015, together with Intel, Huawei, and Vodafone, intending to synchronize work with ETSI MEC, leading so far to OpenStack++ allowing for **migration between OpenStack clusters**.*

Integration with telecommunications per MEC does not appear yet to be paralleled.

Functions at the Edge?

Intended benefits of Edge: reduced end-to-end latency; smaller data volumes travelling shorter distances – computational capability and storage is nearer the user.

Intended benefits of Functions: small, fast-executing, provider-scaled capability.

Faster, smaller implies more suitable for cloud-assisted, or cloud-driven, control services.

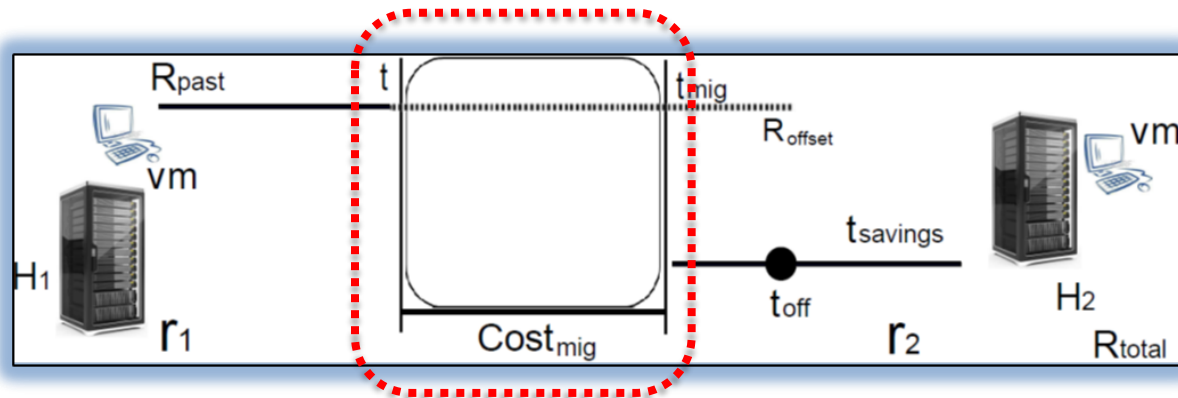
But: susceptible to hardware variations (performance), including due to contention, as well as provider-driven energy management.

*See, also, AWS Lambda at **Edge** (CDN-related FaaS).*

Edges running VMs / Containers / Functions, using various data

Migration when user moves to another edge

Execution would vary with hardware (slowdown/speedup may imply needing more/less resource for equivalence at the target – and need to know this)



Is a user highly likely to keep moving quickly between edges?

An Application



Socio-economic forces

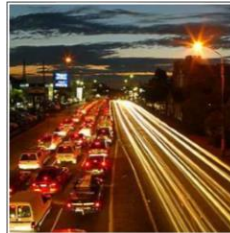


Safety

- 8m accidents, 1.3m fatalities, 7m injuries*

Economy

- 90b hours in traffic jams*



Environment

- 20-35% of global CO2 emission



Passenger Comfort



Market forces



(How) Can cars become fully independent?

What technology advances are required?



Antilocking Brake System (ABS)



Electronic Stability Program (ESP)



Adaptive Cruise Control (ACC)



Lane Keeping Assist



Parking Assist



Autonomous Driving in Traffic Jam and Highway



Full Autonomy

How will humans and vehicles interact with each other and their environment?


5 year UK research programme - Cloud-Assisted Real-time Methods for Autonomy (CARMA) project

[Skip Nav](#) | [Accessibility](#) | [Disability tools](#) | [Media Enquiries](#)

[Communities login](#) 

EPSRC

Engineering and Physical Sciences
Research Council

 search

|  | FUNDING | RESEARCH | INNOVATION | SKILLS | NEWS, EVENTS AND PUBLICATIONS | ABOUT US |
|---|---------|----------|------------|--------|-------------------------------|----------|
|---|---------|----------|------------|--------|-------------------------------|----------|

News, events and publications >

News >

[Home](#) / [News, events and publications](#) / [News](#) / [Jaguar Land Rover and EPSRC announce £11 million autonomous vehicle research programme](#)

Jaguar Land Rover and EPSRC announce
£11 million autonomous vehicle
research programme

See also

[TASCC: Towards Autonomy -
Smart and Connected Control](#)

Related links

[RAS 2020 Robotics and
Autonomous Systems \(PDF
1.71MB\) \[connect -
InnovateUK\]](#)

[Jaguar Land Rover](#)

Issue date: 09 October 2015

Tag: Press Release

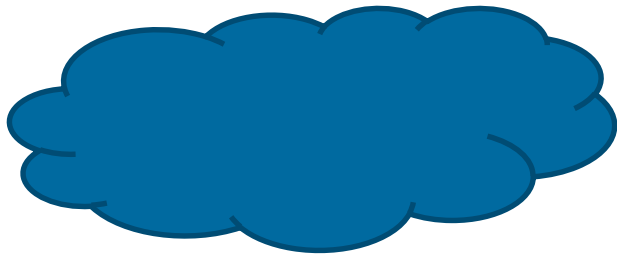
Related themes: [Engineering](#) [ICT](#) [Manufacturing the future](#)

TOWARDS AUTONOMY
SMART AND CONNECTED CONTROL



THALES





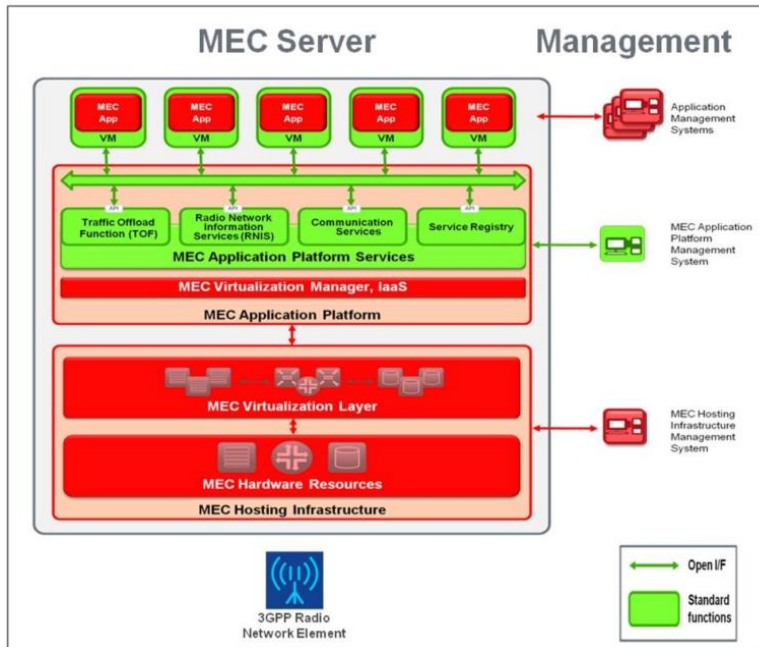
Internet of Things (Tractors, Kettles, Fridges, Cars)

Waymo generates 1GB data **per second** (2 PB/car/year)

Estimate of 2 billion total vehicles by 2020

Some of this 1GB/s may be usefully processed close to the vehicles, with vehicles connecting through the RAN to access or provide e.g. local road information

Some data may aggregated over MEC Servers, or where capability is not needed immediately, or if not available locally



CARMA's vision: design and validate a novel, secure framework to *enable* implementation of safe and robust semi-autonomous and fully autonomous functions

Main objectives

Address key technical research challenges
Validate through proof-of-concept demonstrators
Evaluate scalability

Multi-access/mobile edge, 5G , Cloud.

Security and effects on performance and latency

Safety remains paramount

CARMA Core (Cloud):

- Based on commercially available **public cloud** resources
- Services where **higher latency** is tolerable, information is coming from a **wider geography**, longer term storage needs, and so on.

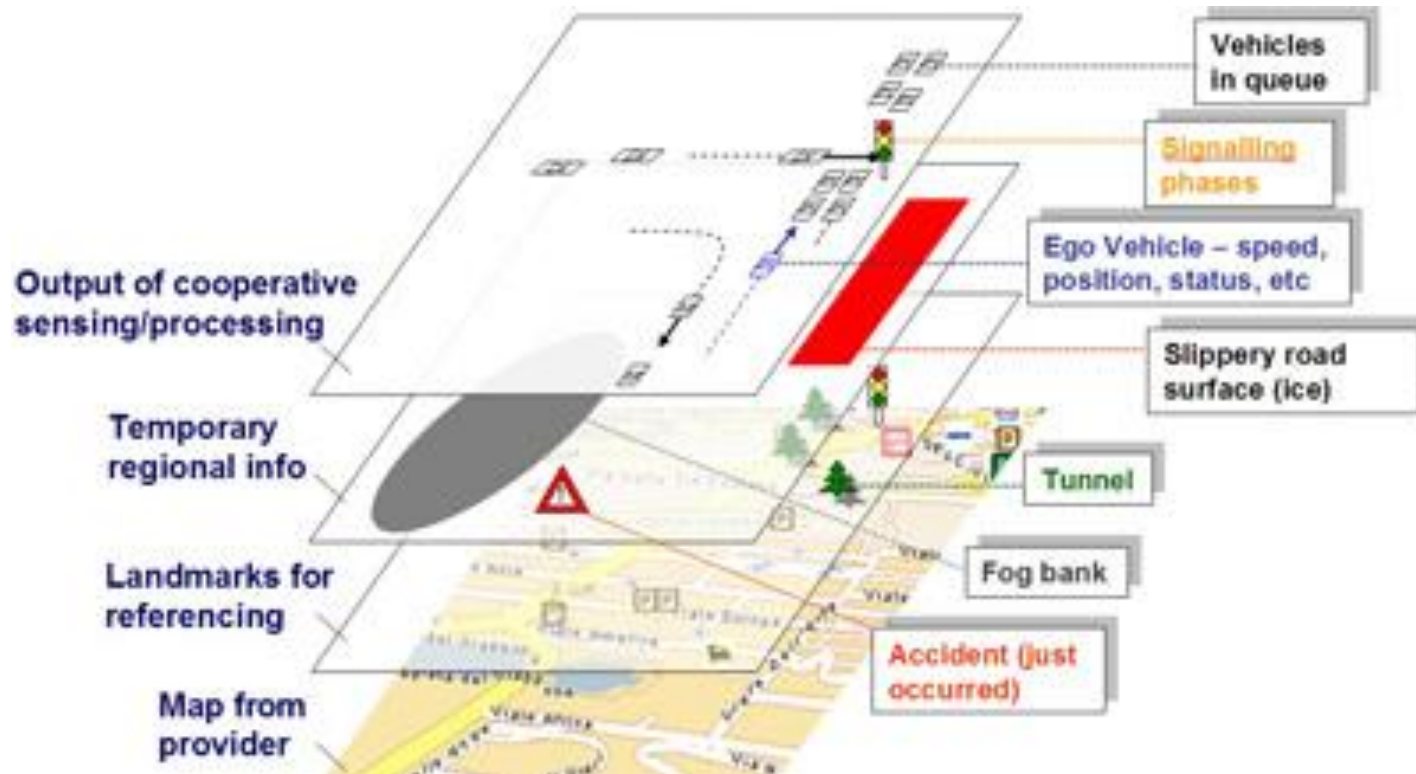
CARMA Edge (MEC, 5G):

- Host beneficially **off-board** (low latency) processes
- Information collected from *around* vehicles to support **cooperativity** and computation beyond capabilities, including sensor ranges, of a given vehicle.

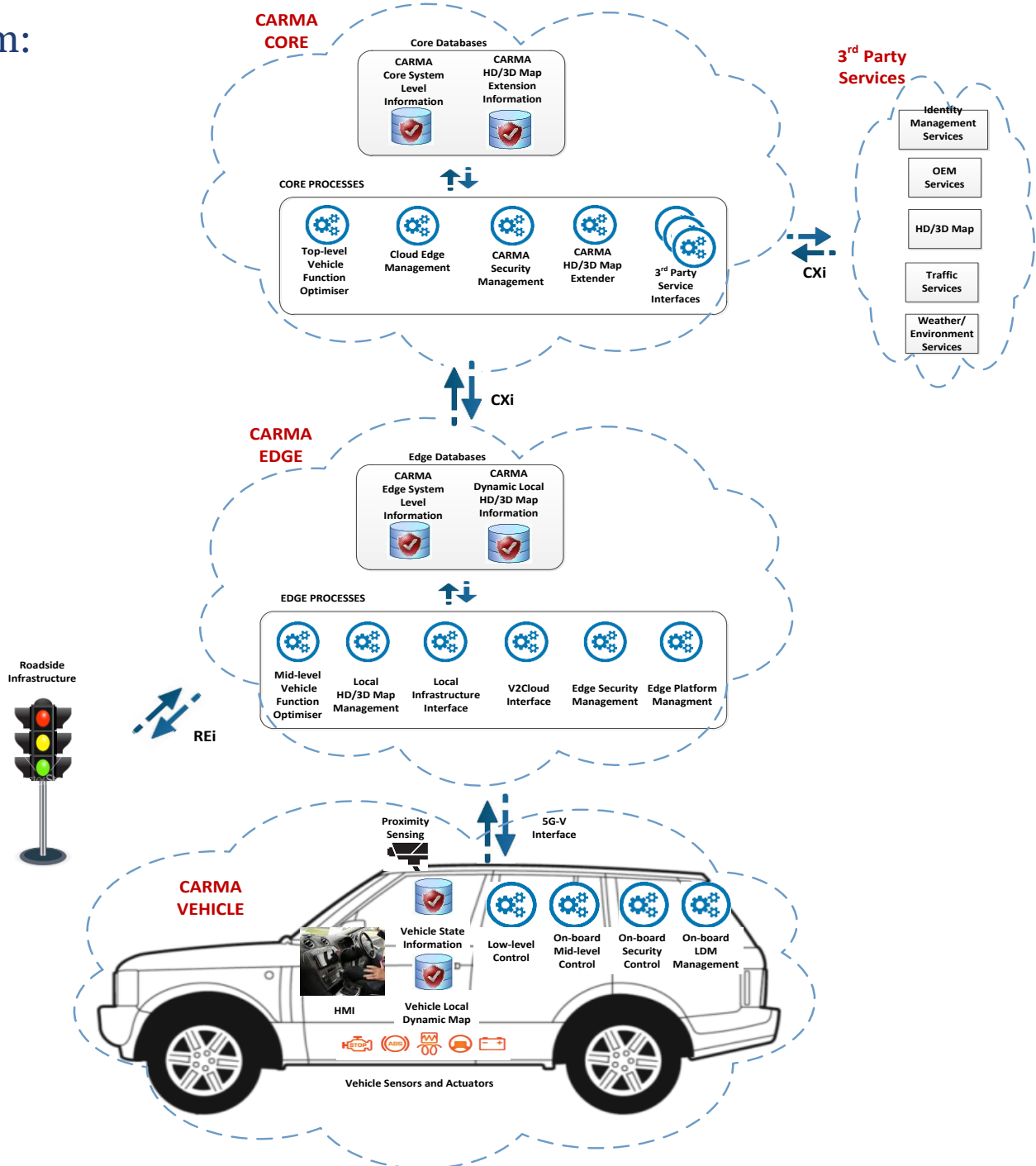
CARMA Vehicle:

- **On-board vehicle network** across all sensors, infotainment, actuators etc.
- *Significant increases in on-board computational capability to be expected.*

Cooperativity beneficial for maps where information has varying transience – information around vehicles and beyond sensor ranges



CARMA Platform: Logical Design



Public Cloud Exemplars

IBM Cloud-connected Vehicles

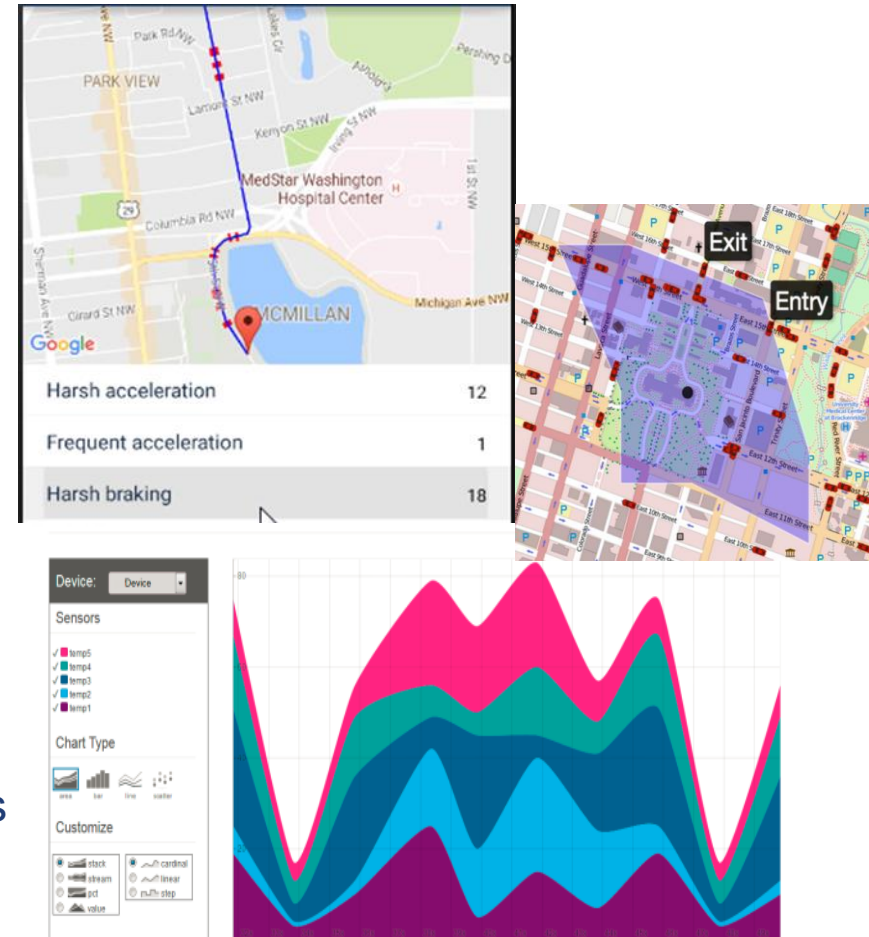
Largely Vehicle to Cloud

- IBM Watson IoT Driver Behavior Service
- IBM Connected-car IoT app with Geospatial Analytics.
- Microsoft Connected Vehicle Platform
- Google Connected Vehicle Platform
-

A hint of Edge, in -

- AWS Connected Vehicle Solution
- IBM Edge (Apache Edgent) analytics
- ...

None of these address Autonomous



Public Cloud Exemplars

AWS Connected Vehicle Solution

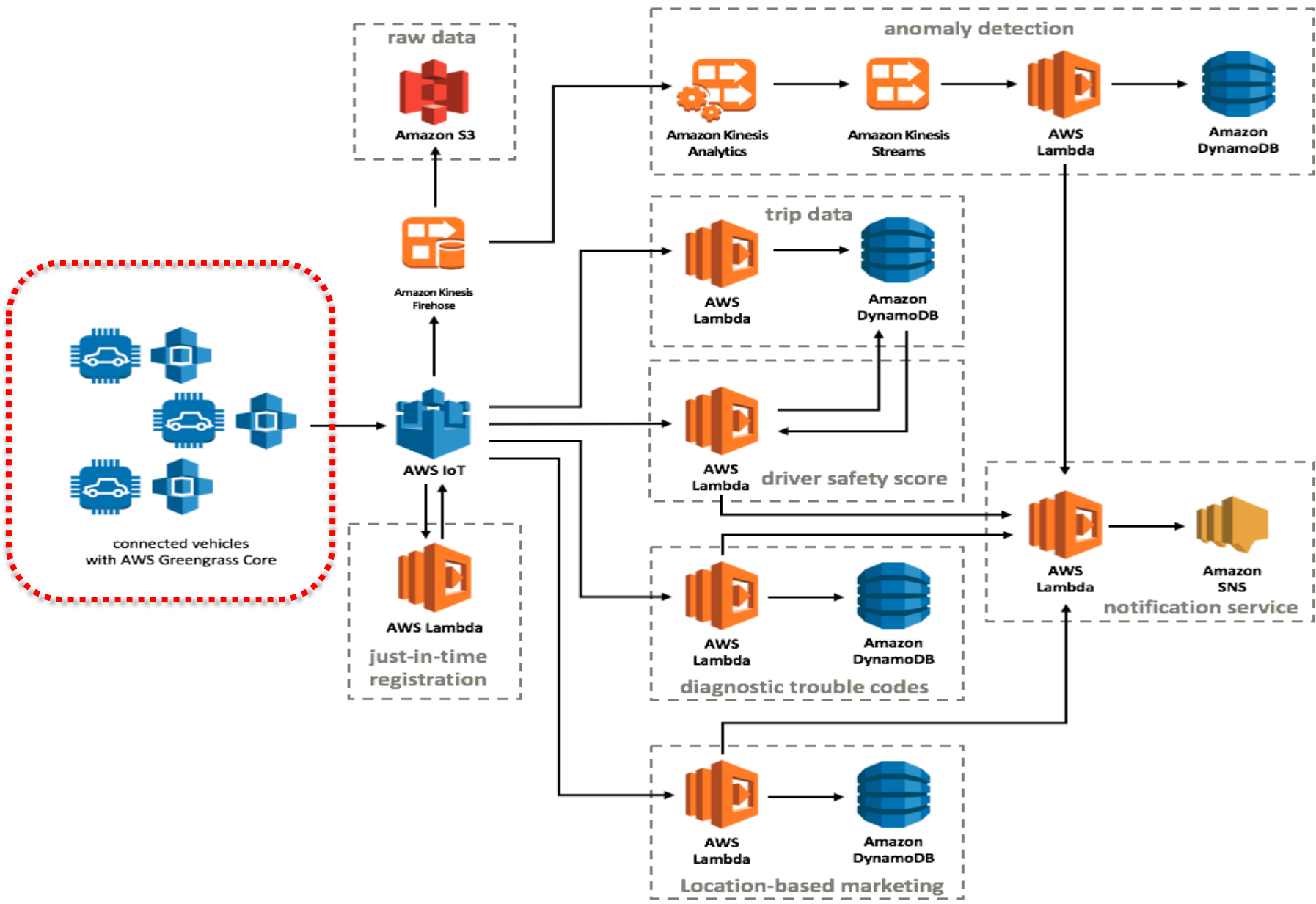
AWS suggests an application stack launchable from a template for multiple services

- Long term data storage
- Treatment of telematics anomalies
- Capture of trip data
- Calculation of a driver safety score
- Diagnostics and reporting

Messages through IoT gateway (Message Queue Telemetry Transport – MQTT, pub/sub): hierarchy of topics, and data payloads; rules on messages trigger execution of Lambda Functions for applications.

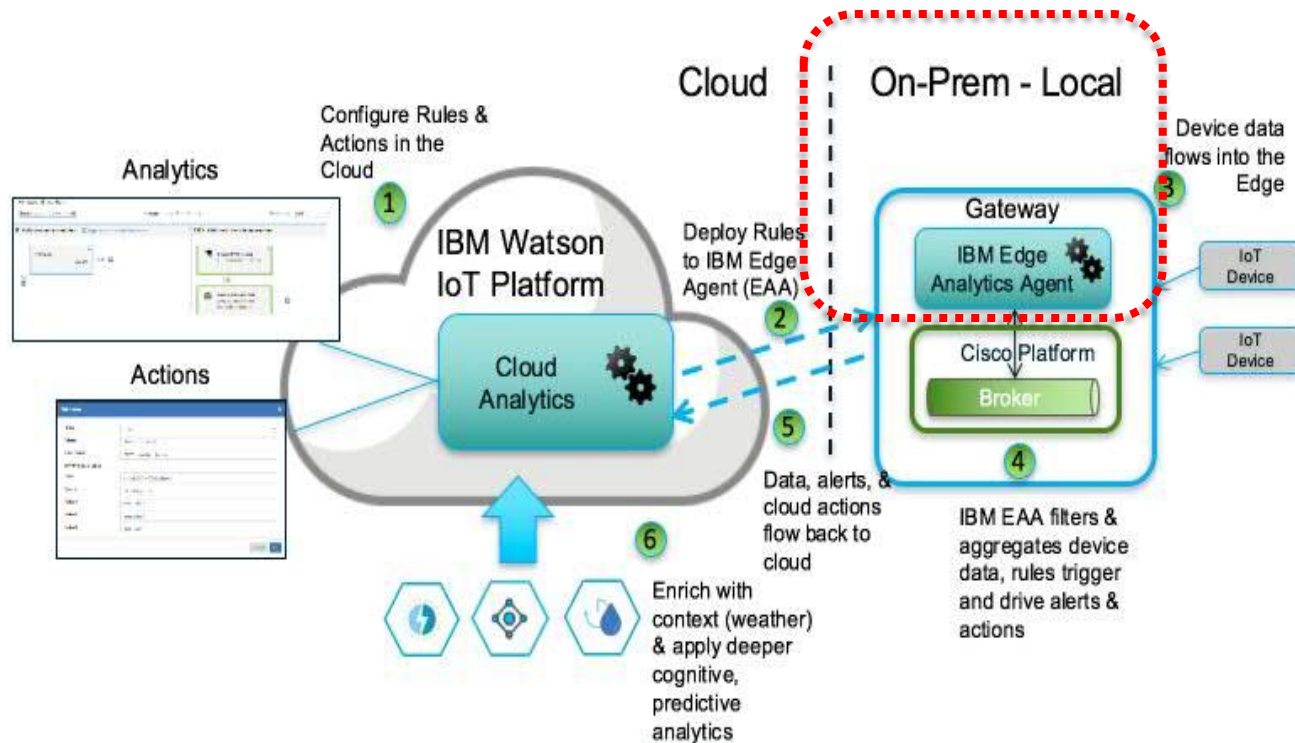
Includes AWS Greengrass for the Edge, but most of the work in Cloud

Public Cloud Exemplars



IBM Edge Analytics: Edge Agent on a Raspberry Pi and the DGLux tool.

How Edge Analytics works



Edge offerings largely not yet replicating Clouds; actions, per provider, required in order to achieve.

- **AWS** IoT rules not yet deployable to (Greengrass) Edge; local versions of other services largely also unavailable.
- **IBM** “Edge Analytic Rule(s) are pushed to the edge device” but no other services
- *Microsoft IoT Edge requires a device running Windows 10 or Windows server with Docker*
- *Google IoT seems not to relate to Edge capability as yet.*
- Vendors also promoting their **own** flavour of Function as a Service.

MQTT (Advanced Message Queuing Protocol), tending to be supported – AMQP and others less so.

This is an improving situation.

Summary and take home

Summary

Will Cloud gain an Edge?

Which Edge? For us, MEC, but market forces...

'Traditional' Cloud

(Big Four) Clouds are big

Cost and performance (=cost) variation

Performance variation and implications for energy efficiency

'New' cloud

'serverless' and performance

Multiplicity of Edges

'serverless' Edges

An application

Cloud Cars and exemplars

Summary and take home

Take home messages 1/2

Will Cloud gain an Edge?

Provider challenges include:

- What to provide in hardware
 - light lifting (rPis)
 - heavy lifting (servers, stacks)
 - telecommunications connectivity
 - “Moore’s law”
- What to provide on hardware
 - Bare metal
 - VMs/Containers/Functions
- Where to locate
- How to maintain
- Support for **Migration**
- How to secure
- To what quality of service
- At what price



Application/user challenges, as well as services to prioritise, include:

- Which Clouds/Edges to use? Vendor(s) lock-in? (hostage to functionality)
- What runs where in V/E/C, when?
 - Fixed, or dynamic, accounting for limited, heterogeneous, resources under contention? (dynamic reconfiguration)
 - What **Performance** – guarantees?
- **Migration** between V/E/C and/or across Edges? (c/w elastic / scalable)
 - Live migration – uninterruptible?
- How to secure? (What to secure? When to secure?)
- How to price across V/E/C? Who pays whom for what?
 - Fixed / dynamic
 - Based on services and/or demand?



An opportunity for an interested researcher

Cloud Assisted Real-time Methods for Autonomy (CARMA)

Research Fellow in Mobile/Multi-access Edge Cloud
Computing

Deadline: 8 April

<https://jobs.surrey.ac.uk/Vacancy.aspx?id=4701>

- A. Stevens, M. Dianati, K. Katsaros, C. Han, S. Fallah, C. Maple, F. McCullough, and A. Mouzakitis. Cooperative automation through the cloud: The CARMA project. In Proceedings of 12th ITS European Congress. 2017.
- 5G Automotive Vision. White Paper, 5G-PPP, 2015.
- P. Mell, T. Grance, NIST Definition of Cloud Computing, NIST, 2011
- T. P. Morgan, A rare Peek Into The Massive Scale of AWS, 2014, <https://www.enterprisetech.com/2014/11/14/rare-peek-massive-scale-aws/>
- F. Bonomi, R. Milito, J. Zhu, and S. Addepalli, Fog computing and its role in the internet of things. Proc. First Edition of the MCC Workshop on Mobile Cloud Computing, MCC '12, 13–16, New York, USA, 2012.
- ETSI. Mobile Edge Computing (MEC). Introductory Technical White Paper, ETSI, 2014.
- S. Senior, C. Rec, H. Nishar, T. Horton, “AWS Connected Vehicle Solution: AWS Implementation Guide”, 2017.
<https://s3.amazonaws.com/solutions-reference/connected-vehicle-cloud/latest/connected-vehicle-solution.pdf>
- Microsoft Corporation, “Empowering automotive innovation”, 2017.
http://download.microsoft.com/download/6/9/D/69D92EB1-F1EE-4893-ABE1-C005D7F9FF57/Microsoft_Connected_Vehicle_Platform_Whitepaper_EN_US.pdf
- Google, “Designing a Connected Vehicle Platform on Cloud IoT Core” 2017.
<https://cloud.google.com/solutions/designing-connected-vehicle-platform>

- Lee Gillam, Konstantinos Katsaros, Mehrdad Dianati and Alex Mouzakitis (2018) "Exploring Edges for Connected and Autonomous Driving". IEEE INFOCOM Workshops: CCSNA 2018.
- Muhammad Zakarya and Lee Gillam (2017) "An Energy Aware Cost Recovery Approach for Virtual Machine Migration". In Bañares, José Ángel, Tserpes, Konstantinos, Altmann, Jörn (Eds.), Proc. GECON 2016.
- Nick Antonopoulos and Lee Gillam (2017) "Cloud Computing: Principles, Systems and Applications". *Second Edition*. Springer-Verlag.
- John O'Loughlin and Lee Gillam (2017) "A Performance Brokerage for Heterogeneous Clouds". FGCS.
Muhammad Zakarya and Lee Gillam (2017) "Energy Efficient Computing, Clusters, Grids and Clouds: A Taxonomy and Survey". Journal of Sustainable Computing, Informatics and Systems.
- John O'Loughlin and Lee Gillam (2016) "Sibling Virtual Machine Co-location Confirmation and Avoidance Tactics for Public Infrastructure Clouds". Journal of Supercomputing 72(3):, 961-984
- John O'Loughlin and Lee Gillam (2015) "Addressing Issues of Cloud Resilience, Security and Performance through Simple Detection of Co-locating Sibling Virtual Machine Instances". 5th International Conference on Cloud Computing and Services Science (CLOSER 2015).
- John O.Loughlin and Lee Gillam (2014) "Should Infrastructure Clouds be Priced Entirely on Performance? An EC2 Case Study". International Journal of Big Data Intelligence
- John O.Loughlin and Lee Gillam (2014) "Performance Evaluation for Cost-Efficient Public Infrastructure Cloud Use". In J. Altmann et al. (Eds.): GECON 2014, LNCS 8914, pp. 133–145, 2014.

Thank You

l.gillam@surrey.ac.uk

Further information:

Publications: <https://sites.google.com/site/drleegillam/publications>

Twitter: @leegillam

Journal: <http://www.journalofcloudcomputing.com/> - fully open access to in-depth Cloud research

<https://jobs.surrey.ac.uk/Vacancy.aspx?id=4701>

Acknowledgements:

Drs. O'Loughlin & Zakarya; CARMA team, JLR, EPSRC
And, where needed, Google Images

Journal of Cloud Computing

Advances, Systems and Applications

- Cloud Architectures
- Cloud Technologies
- Cloud Services

Springer Open