



**FLAME**

FACILITY FOR LARGE-SCALE ADAPTIVE MEDIA EXPERIMENTATION

# Dynamic service delivery using network-aware graph analytics and endpoint controls

**Nikolay Stanchev**

*IT Innovation*

# Introduction



- *Estimating end-to-end delay measurements through graph-based analytics*
- *Building a network topology graph*
- *Building a temporal graph from time-series measurements of a media service*
- *Querying for round-trip time calculations*
- *Writing round-trip-time calculations as time-series measurements*
- *Automating the full graph monitoring pipeline*
- *Creating a state-change alert policy based on the new metric*

# FLAME CLMC – Cross Layer Management and Control

CLMC combines measurements from different layers of the FLAME platform and allows for the construction of new metrics which could give a better estimate of the performance of a given service.

A particular use case is the estimation of the end-to-end delay of a media service, that is the delay that a client would experience while using this service – a metric which could be broken down to two different factors:

- network-related measurements
- service-related measurements

## Measurement model of the end-to-end service delay

Overall, the assumption is that the total delay of a service function could be measured using the following formula:

$$\mathit{total\_delay} = \mathit{forward\_network\_delay} + \mathit{service\_delay} + \mathit{reverse\_network\_delay}$$

which could be extended to the following:

$$\begin{aligned} \mathit{total\_delay} = & \mathit{forward\_latency} \\ & + \mathit{forward\_data\_delay} \text{ (dependent on request size and bandwidth)} \\ & + \mathit{service\_delay} \\ & + \mathit{reverse\_data\_delay} \text{ (dependent on response size and bandwidth)} \\ & + \mathit{reverse\_latency} \end{aligned}$$

Full details of these calculations can be found at <https://gitlab.it-innovation.soton.ac.uk/FLAME/consortium/3rdparties/flame-clmc/blob/master/docs/total-service-request-delay.md> (contributed by Stephen Phillips).

Important assumption of this simplified model is that services can measure the processing time for requests which is **forward\_data\_delay + service\_delay + reverse\_data\_delay** (out-of-the-box support for nginx, tomcat, etc.)

# Building the network topology graph

CLMC utilises the northbound API of the SDN controller (Floodlight in current implementation) to build up the network topology graph and retrieve the switch-to-switch (a.k.a. service routers) network latencies.

The graph is then stored in Neo4j and is supposed to be managed by the platform provider.

This is implemented as a REST-like API with the following three endpoints:

- **POST** <http://platform/clmc/clmc-service/graph/network> – builds up the network topology graph and creates new nodes and links if needed
- **PUT** <http://platform/clmc/clmc-service/graph/network> – builds up the network topology graph, but also updates the latency measurement of already existing links
- **DELETE** <http://platform/clmc/clmc-service/graph/network> - completely deletes the network topology graph

The *Neo4j* browser could be used to view the graph and explore latency measurements between service routers – <http://platform/clmc/neo4j/browser>



## Building a temporal graph for a media service

A temporal graph is simply a graph representing the state of a media service in a given time window, e.g. from start of today until end of tomorrow.

CLMC builds such graphs to calculate the round-trip time from a specific user equipment to a specific service function endpoint.

In order to build this graph, three metrics must be measured for a given service function:

- **response\_time** (or service delay) – how much time it takes to process a request (seconds), that is from the moment the first byte of the request is read until the moment last byte of the response is sent.
- **request\_size** – the size of a request to the service function (bytes)
- **response\_size** – the size of a response from the service function (bytes)

# Building a temporal graph for a media service



An example with a *Tomcat*-based service:

- use the *Tomcat Telegraf* input plugin for monitoring – includes fields like ***bytes\_sent***, ***bytes\_received*** and ***processing\_time*** in measurement ***tomcat\_connector***.
- ***processing\_time*** is the total time spent processing incoming requests measured since the server has started.
- ***bytes\_sent*** and ***bytes\_received*** measured using the same approach.
- ***request\_count*** gives the number of requests since the server has started.

The screenshot shows a dark-themed interface with a list of fields. At the top, there are three dropdown menus: "Group by: auto", "Compare: none", and "Fill: null". The list of fields includes:

Field Name	Action
bytes_received	1 Function
bytes_sent	1 Function
current_thread_count	
current_threads_busy	
error_count	
max_threads	
max_time	
processing_time	1 Function
request_count	1 Function

# Building a temporal graph for a media service



For CLMC to understand how to build a media service graph, the graph must be described in JSON format. The description must define the following:

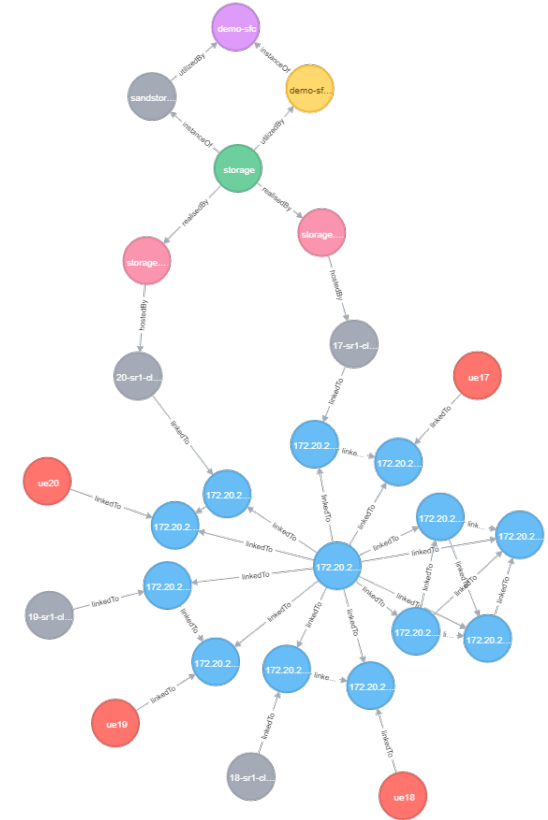
- the time window for the temporal part of the graph (service function endpoint nodes)
- the media service (a.k.a. service function chain) identifiers
- the service function packages the media service is using
- a partial influx query for obtaining the measurement values described in the previous slide (**request/response size** and **service delay**), basically an aggregation function with a field name, e.g. *mean(processing\_time)*
- the measurement name where the fields from these partial queries reside, e.g. *tomcat\_connector*



# Building a temporal graph for a media service

The JSON description is then sent to CLMC as the body of a POST request to `/clmc/clmc-service/graph/temporal`

```
{
  "from": 1549881060,
  "to": 1550151600,
  "service_function_chain": "demo-sfc",
  "service_function_chain_instance": "demo-sfc_1",
  "service_functions": {
    "sandstorage": {
      "response_time_field": "(last(processing_time) - first(processing_time)) / ((last(request_count) - first(request_count)) * 1000)",
      "request_size_field": "(max(bytes_received) - min(bytes_received)) / (last(request_count) - first(request_count))",
      "response_size_field": "(max(bytes_sent) - min(bytes_sent)) / (last(request_count) - first(request_count))",
      "measurement_name": "tomcat_connector"
    }
  }
}
```



## Building a temporal graph for a media service

The JSON description is fixing a time window (from **1549881060** to **1550151600**, UNIX timestamps) and defining how to query the average response time, request size and response size of all endpoints that use the *sandstorage* service function package.

Since Tomcat's measurement model is reporting continuously increasing metrics (the value since the server has started), the partial influx queries look a bit more complicated.

For example, the average response time query:  $(last(processing\_time) - first(processing\_time)) / ((last(request\_count) - first(request\_count)) * 1000)$

Simple scenario:

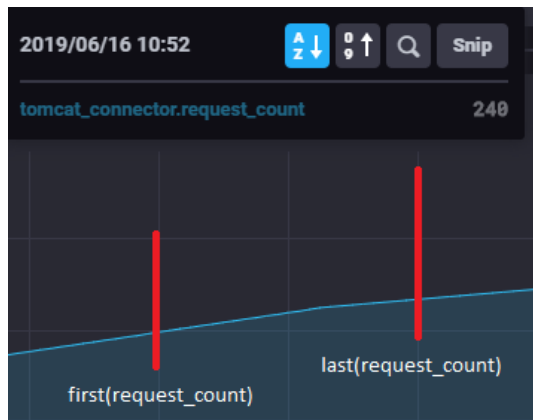
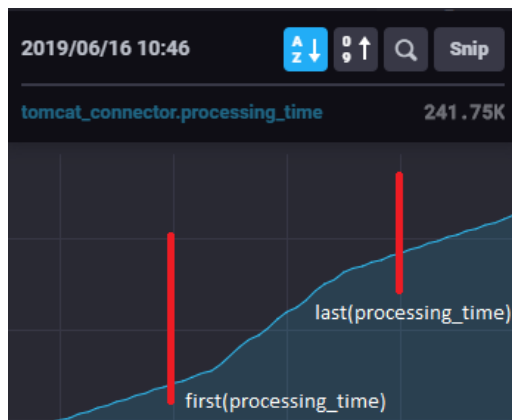
- processing\_time measurements (measured in milliseconds) received in the time period – 41629, 41641, 41793, 41839
- requests\_count measurements received in the time period – 102, 103, 108, 110
  
- $last(processing\_time) - first(processing\_time) = 210$  (milliseconds used to process all requests in this time period)
- $last(request\_count) - first(request\_count) = 8$  (total of 8 requests processed in this time period)
- $(last(processing\_time) - first(processing\_time)) / (last(request\_count) - first(request\_count)) = 26.25$

The query defined above will evaluate to 26.25ms = 0.02625s and will give us the the average delay of the service per request. The same reasoning is used for *bytes\_sent* and *bytes\_received* to calculate the average request size and average response size.

## Building a temporal graph for a media service

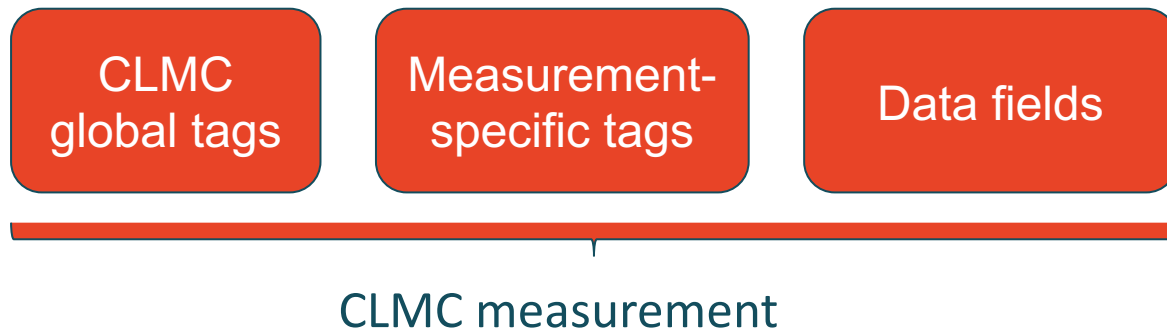
Visualisation is always more helpful than simply looking at the numbers:

- measurement time window for 'processing\_time'
- measurement time window for 'request\_count'



## From time-series to graph data

CLMC monitoring data model:



- Media service global tags – *flame\_sfc, flame\_sfp, flame\_sfe, flame\_location, etc.*
- Through the context given by these measurements CLMC can extract the graph nodes and relationships from the time-series data.

# From time-series to graph data



What happens in the background is that a time-series data query is executed and the results are used to generate the media service graph in the Neo4j database:

```
SELECT {0} AS mean_response_time, {1} AS mean_request_size, {2} AS mean_response_size FROM "{3}"."{4}".{5} WHERE "flame_sfc"='\{6}\' and "flame_sfci"='\{7}\' and "flame_sfp"='\{8}\' and time>={9} and time<{10} GROUP BY "flame_sfe", "flame_location", "flame_sf"
```

The placeholders are filled from the JSON configuration described in the previous slides. Depending on the results we can identify:

- Service function endpoint nodes (from the *flame\_sfe* tag)
- Service function nodes (from the *flame\_sf* tag)
- Service function package nodes (from the *flame\_sfp* tag)
- Service function chain nodes (from the *flame\_sfc* and *flame\_sfci* tag)
- Cluster nodes (from the *flame\_location* tag)

```
Endpoint <id>: 29 name: storage.demo-sfc.ict-flame.eu-172.90.12.52
ServiceFunction <id>: 7 name: storage
ServiceFunctionPackage <id>: 86 name: sandstorage
ServiceFunctionChain <id>: 24 name: demo-sfc
Cluster <id>: 71 name: 20-sr1-cluster1-cluster
```

Thus, the full media service graph is built on top of the network topology graph – cluster nodes being the intersection point between the two graphs.

## Query for round-trip-time estimation

Once the full graph has been built, CLMC offers an API endpoint to query for round-trip-time from a start point (Cluster, Switch or a User Equipment node) to an end point (Service Function Endpoint). The response is a breakdown of the round-trip-time measurement:

- **GET** <http://platform/clmc/clmc-service/graph/temporal/<uuid>/round-trip-time?startpoint=<cluster, switch or ue>&endpoint=<SF endpoint>>

The API endpoint is basically a query to the temporal graph for round-trip-time; the UUID parameter uniquely identifies the subgraph.

The UUID of the temporal graph can be retrieved from the response of the build request described in the previous slides.

```
1  {
2    "total_forward_latency": 0.008,
3    "local_tags": {
4      "traffic_source": "ue17"
5    },
6    "response_size": 1314.5,
7    "response_time": 0.003727272727272727,
8    "forward_latencies": [
9      0,
10     0.002,
11     0.006,
12     0
13   ],
14   "reverse_latencies": [
15     0,
16     0.006,
17     0.002,
18     0
19   ],
20   "global_tags": {
21     "flame_sfp": "sandstorage",
22     "flame_sfc": "demo-sfc",
23     "flame_server": "20-sr1-cluster1-cluster",
24     "flame_sfci": "demo-sfc_1",
25     "flame_sfe": "storage.demo-sfc.ict-flame.eu-172.90.12.52",
26     "flame_sf": "storage",
27     "flame_location": "20-sr1-cluster1-cluster"
28   },
29   "total_reverse_latency": 0.008,
30   "request_size": 0,
31   "round_trip_time": 0.019727272727272725
32 }
```

## Combining network and service measurements



In the background, a Neo4j Cypher query for shortest-path is executed (based on number of hops) to retrieve the network path between the start point and the end point defined in the URL of the round-trip time query request:

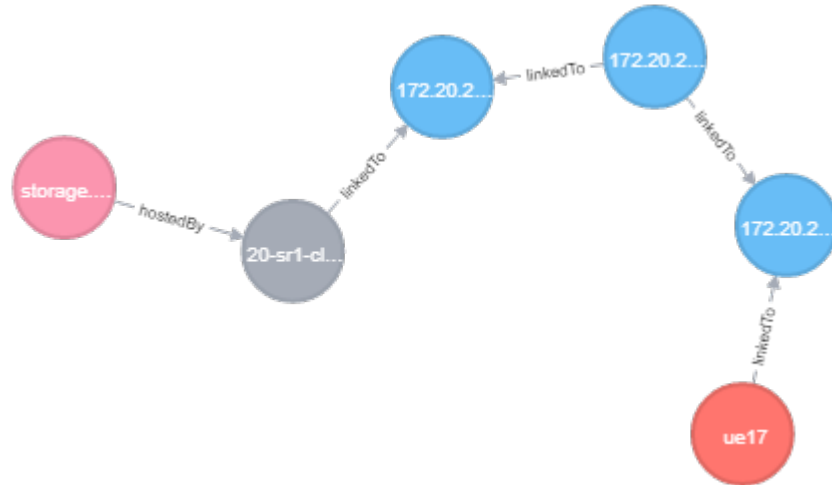
```
MATCH (startpoint:{0} {{ name: '{1}' } }),(endpoint:Endpoint {{ name: '{2}', uuid: '{3}' } }},  
path = shortestPath((startpoint)-[*]-(endpoint))  
WHERE ALL(r IN relationships(path) WHERE type(r)='linkedTo' or type(r)='hostedBy' )  
WITH extract(y in filter(x in relationships(path) WHERE type(x) = 'linkedTo') | y.latency) as latencies,  
endpoint.response_time as response_time, endpoint.request_size as request_size,  
endpoint.response_size as response_size  
RETURN latencies as forward_latencies, reverse(latencies) as reverse_latencies, response_time,  
request_size, response_size
```

Placeholders are filled with the request URL query parameters.

# Combining network and service measurements

Through the result of the previous query, measurements from the network topology layer are combined with measurements from the application layer, i.e. the Service Function Endpoint level.

- Network measurements coming from the network path to the service function endpoint
- Application measurements coming from the temporal service function endpoint node





## Converting round-trip time results to time-series data

Once the graph has been built and a round-trip time query is executed, the results can be written back as a new measurement in InfluxDB.

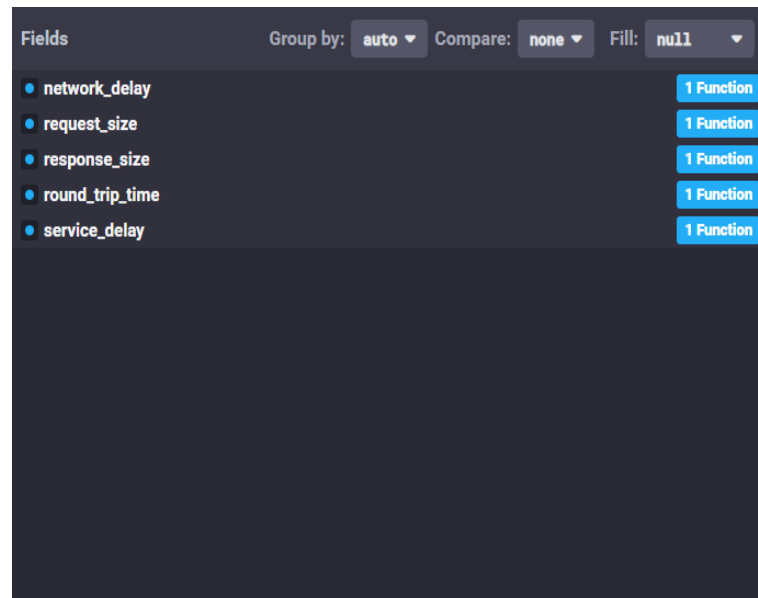
The following steps need to be followed to achieve this:

- convert the JSON response from the Graph API to Influx line protocol format

Example:

```
graph_measurements,flame_sfp=sandstorage,flame_sfc=demo-sfc,flame_server=17-sr1-cluster1-cluster,flame_sfci=demo-sfc_1,flame_sfe=storage.demo-sfc.ict-flame.eu-172.90.4.52,flame_sf=storage,flame_location=17-sr1-cluster1-cluster,traffic_source=ue18  
round_trip_time=0.002,service_delay=0.002,network_delay=0,request_size=0,response_size=1422.5  
1560675146000000000
```

- send a POST request to InfluxDB including the measurement line above



## Clean up the media service graph

Since the endpoints layer of the full graph is *temporal* and only valid for the defined time window, it needs to be deleted and created again if a more recent round-trip-time measurement is needed.

- **DELETE** <http://platform/clmc/clmc-service/graph/temporal/<uuid>>

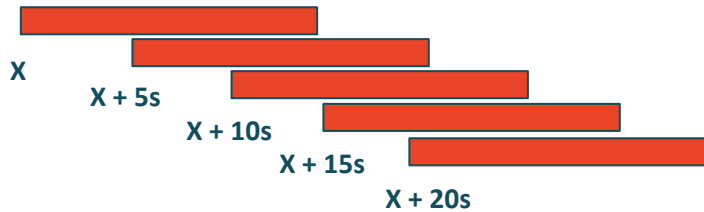
This finishes the lifecycle of a graph monitoring activity – build temporal graph, query it , delete it.

Alternatively, if the full graph of a service function chain must be deleted there is a separate API endpoint to use:

- **DELETE** <http://platform/clmc/clmc-service/graph/static/<SFC identifier>>

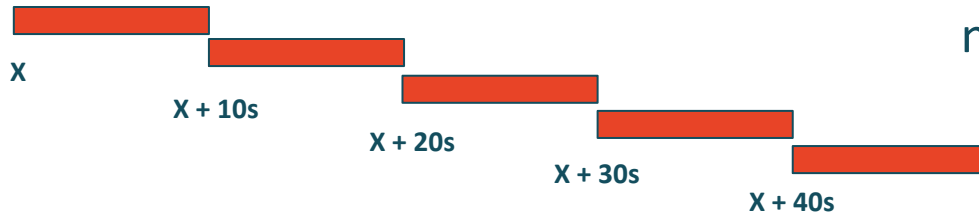
Ultimately, the media service provider can choose how to manage it. Two main strategies:

Sliding window:



overlapping time periods

Tumbling window:



Out-of-the-box support from CLMC



non-overlapping time periods

## Automating the end-to-end delay monitoring process

CLMC offers the full graph-based pipeline as a service. An API endpoint allows the activation of a graph monitoring process, running in the background on CLMC constantly executing the pipeline described in the previous slides.

A JSON configuration, similar to the one used in the build request for a temporal graph, is sent to CLMC to start a graph monitoring process

- **POST** <http://platform/clmc/clmc-service/graph/monitor>

The difference is that instead of defining a time window, we define a query period (e.g. 30 seconds, that is how often the pipeline script will execute) and the name of the measurement where results will be written in.

# Automating the end-to-end delay monitoring process

## Example JSON description:

```
{
  "query_period": 30,
  "results_measurement_name": "graph_measurements",
  "service_function_chain": "demo-sfc",
  "service_function_chain_instance": "demo-sfc_1",
  "service_functions": {
    "sandstorage": {
      "response_time_field": "(last(processing_time) - first(processing_time)) / ((last(request_count) - first(request_count)) * 1000)",
      "request_size_field": "(max(bytes_received) - min(bytes_received)) / (last(request_count) - first(request_count))",
      "response_size_field": "(max(bytes_sent) - min(bytes_sent)) / (last(request_count) - first(request_count))",
      "measurement_name": "tomcat_connector"
    }
  }
}
```

- the pipeline executes every 30 seconds building a temporal graph for the time window between *now()* – 30s and *now()*
- end-to-end delay metrics will be written into measurement named *graph\_measurements*

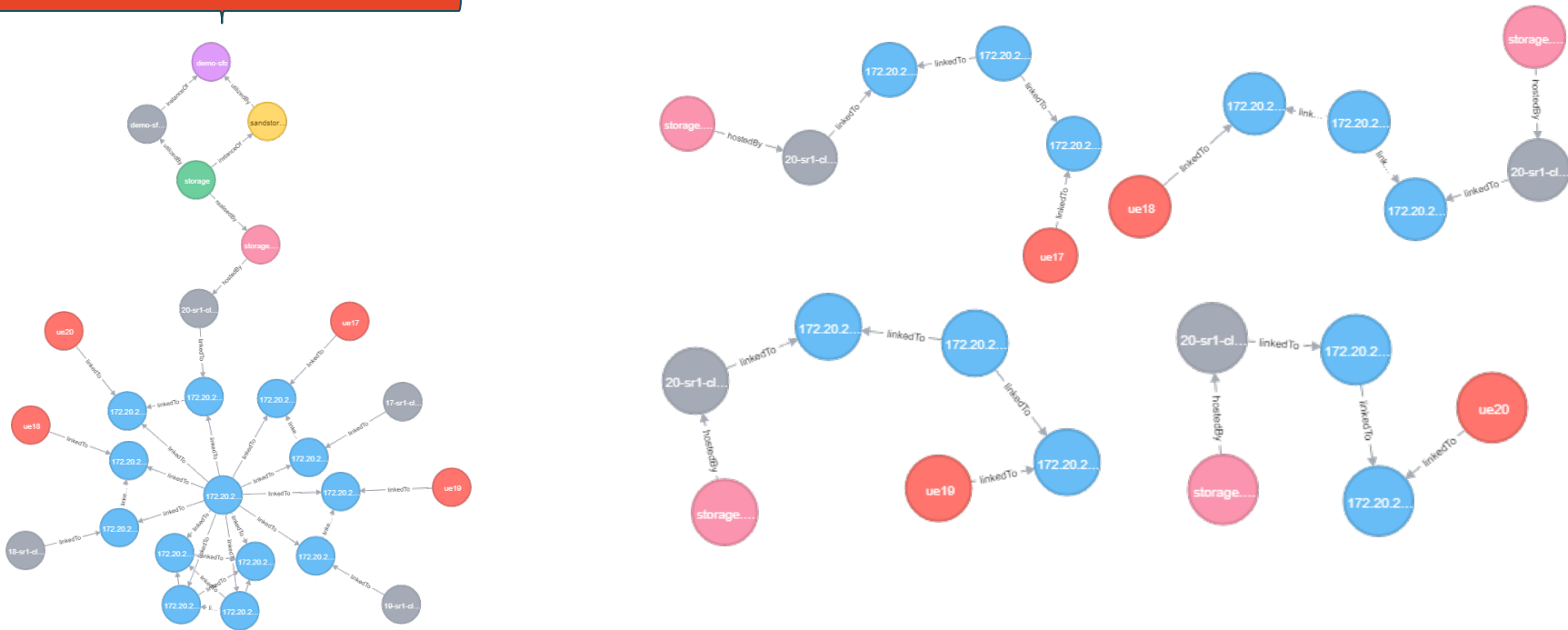


# Automating the end-to-end delay monitoring process

Runtime execution of the graph monitoring process – querying for round-trip time from all UE nodes:

Timestamp X

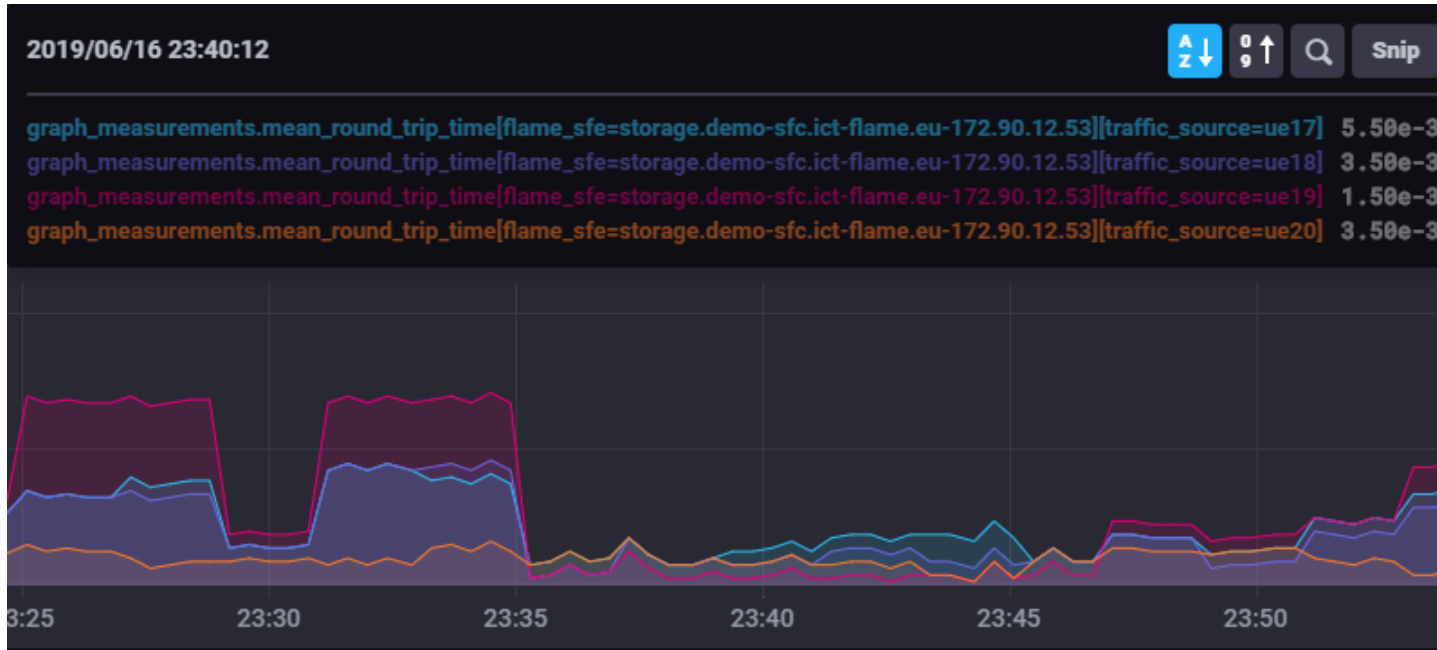
X + 30s



# Automating the end-to-end delay monitoring process



Runtime execution of the graph monitoring process – writing back the results from all round-trip time queries into InfluxDB and generating a new measurement with contextualised data:





## Managing the graph monitoring process

As with any other monitoring agents or measurement plugins, we need to be able to manage this graph pipeline.

A graph monitoring process can be stopped:

- **DELETE** `http://platform/clmc/clmc-service/graph/monitor/<uuid>`

or we can check its status:

- **GET** `http://platform/clmc/clmc-service/graph/monitor/<uuid>`

The UUID identifying a graph monitoring process is retrieved from the CLMC response for the request that started it.

## Creating an alert policy for the round-trip time metric



Now that we have a new metric, we can create a simple threshold alert policy which will boot a second service function endpoint once the average round-trip time performance exceeds the given value:

```
scale_out:
  event_type: threshold
  metric: graph_measurements.round_trip_time
  condition:
    threshold: 2
    granularity: 30
    comparison_operator: gte
    resource_type:
      flame_sf: storage
      flame_location: 20-sr1-cluster1-cluster
  action:
    implementation:
      - flame_sfemc
```

Usually, it is the *flame\_sfe* tag used to identifier an endpoint, runtime generation though – therefore, we use the combination of the two tags *flame\_sf* and *flame\_location* to identifies the first service function endpoint.

## Creating an alert policy for the round-trip time metric

In addition, we also create a trigger to stop the second service function endpoint when it is not doing any work:

```
scale_in:
  event_type: threshold
  metric: graph_measurements.round_trip_time
  condition:
    threshold: 0.5
    granularity: 65
    comparison_operator: lt
    resource_type:
      flame_sf: storage
      flame_location: 17-sr1-cluster1-cluster
  action:
    implementation:
      - flame_sfemc
```

Again, the combination of the two tags *flame\_sf* and *flame\_location* is used to identify the second service function endpoint.



**FLAME**



This project received funding from the European Union's Horizon2020 research and innovation programme under grant agreement No 731677

**THANK YOU FOR YOUR ATTENTION**



[ICT-FLAME.EU](http://ICT-FLAME.EU)



[@ICT\\_FLAME](https://twitter.com/ICT_FLAME)