# DataFog: Towards a Holistic Data Management Platform for IoT Age at the Network Edge

——

Harshit Gupta, Zhuangdi Xu, Umakishore Ramachandran

Embedded Pervasive Lab (EPL), Georgia Institute of Technology

# Motivation

- **Situation awareness applications on edge**
  - **->** low-latency between sensing and actuation

- **Cloud-based data management**
  - **->** inevitable high latency

- **Bandwidth intensive IoT platforms**
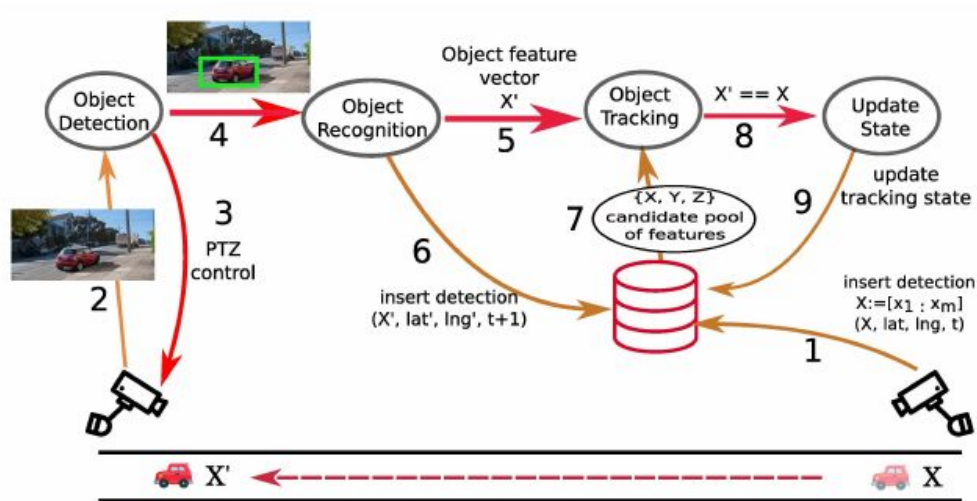  - **->** pressure on backhual bandwidth to transport data to the Cloud

The need to build a datastore at the edge of the network

# Challenges

- **Wide geo-distribution and heterogeneous of the edge infrastructure**

  - **->**  data-partitioning and replication policies

- **Scarcity of resources at the Edge**

  - -> interplay b/w the Edge (for low-latency) and the Cloud (for abundance of resources)

- **Resources at the Edge are more prone to failures**

  - Susceptible to geographically correlated failures

# Use case: Suspicious vehicle tracking

- Spatio-temporal range queries such as select all vehicle detections within 5km and 10 minutes to be efficient
- The distribution of workload is dependent on the distribution of vehicles in space, leading to hotspots
- For continuous operation, continuous streams of vehicle detections have to be saved in a datastore

# Key Characteristics

1. Spatio-temporal locality in range queries
2. Data-model: type, location, timestamp and value
3. Continuous generation of data
4. High availability requirements

# DataFog

A system that performs data partitioning between the Edge and the Cloud based on contextual relevance of data-items in space and time.

# Locality-aware distributed indexing

- Data-items are indexed based on their spatio-temporal attributes (e.g. Geohash)

- Consistent hashing for the location, timestamp and item-type attributes is used for partitioning data across nodes

{ "metric" : "ACV2351",
  "location" : {
    "latitude" : "33.42553",
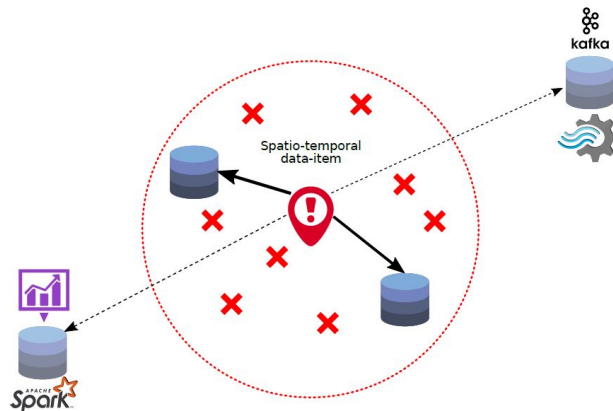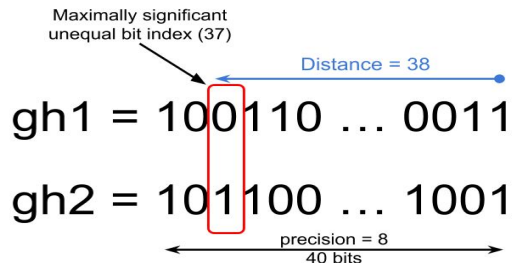    "longitude" : "-84.74456"
  }
  "timestamp" : "1520123197"
}

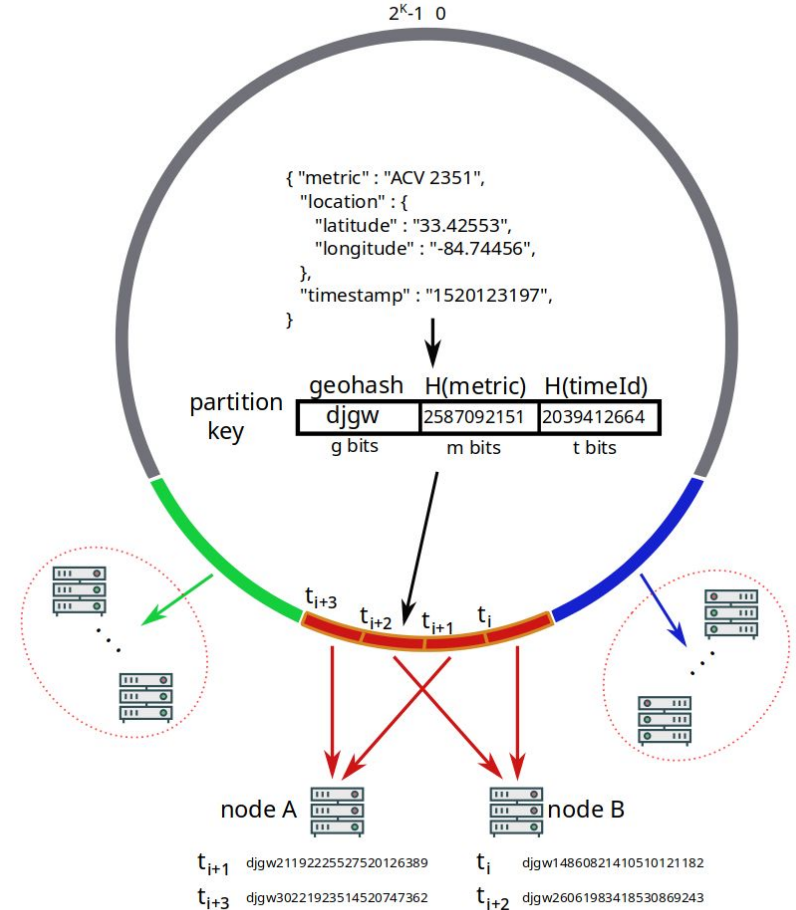| Geohash | H(metric) | H(timeId) |
|---|---|---|
| djgw | 258709251 | 2039412664 |

# Replication Policy

- Load-balancing and fault-tolerance

- Multiple replicas on Edge nodes for low latency

- Multiple replicas on remote datacenter nodes for tolerance from geographically correlated failures



Maximally significant unequal bit index (37)

Distance = 38

gh1 = 100110 … 0011

gh2 = 101100 … 1001

precision = 8
40 bits



kafka

Spatio-temporal data-item

Spark

# Handling workload skews

- Load-balancing region
- Partition key -> virtual node -> physical node
- Mechanisms for adapting to hotspots
  - Long-lived: launch and attach new datastore nodes to the running cluster
  - Short-lived: offload heavily loaded nodes's data items to lightly loaded nodes
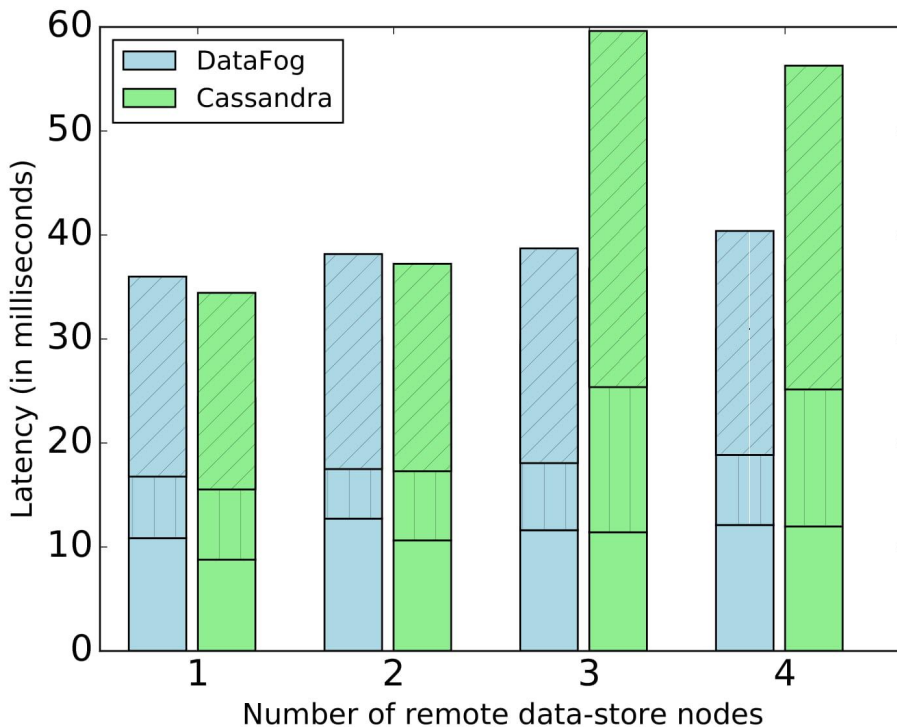
# Handling scarce resources at the edge

- TTL-based data eviction
    - Real-time analytics on temporal data
    - Batch-processing requires data spanning over a large period of time

- Data aggregation and compression
    - Omit redundant metadata to increase efficiency of storage utilization
    - Isomorphism of time series data

# Preliminary Evaluation

- Locality-aware distributed indexing

- SUMO to simulate vehicles movement in Georgia Tech campus equipped with 35 smart cameras

- MaxiNet to simulate 4 edge nodes within the campus and 4 remote nodes (CA, WA, IL and FL) on Microsoft Azure

- Range query: select all vehicle detections within 5km and 10 minutes

# Cont'

- Compared to a location-agnostic indexing done by off-the-shelf Cassandra
- Replication factor: 3
- When the number of remote nodes is 3 or greater, some data items end up having replicas only on remote nodes making the higher percentiles of latencies becoming higher

**Conclusion:**

1.  Present the case for a holistic management platform for IoT data at the network edge
2.  Identify the challenges and come up with algorithmic insights for addressing them
3.  Potential of such a platform as the improvement in performance by a replica placement approach based on spatial locality

**Open issues:** Interaction between datastore platforms owned by different stakeholders leads to a need of communication protocols and business models for sharing data across multiple edge administrative domains.

Q & A

# Conclusion

- Present the case for a holistic management platform for IoT data at the network edge
- Identify the challenges and come up with algorithmic insights for addressing them
- Potential of such a platform as the improvement in performance by a replica replacement approach based on spatial locality
- Future work: quantitative evaluation of the design decisions in comparison to state-of-the-art Cloud-based datastores

# Full evaluation

- Overhead of context-aware partitioning and replication
- Ability of load-balancing solutions to manage workloads with inherent skews
- The benefit of eviction-based strategy on utilization of storage resources at the edge
- Reduction of storage consumption
- Impact of parameters including replication distance, spatial encoding precision and etc