# Scalable Spatio-temporal Indexing and Querying over a Document-oriented NoSQL Store

**Nikolaos Koutroumanis**, Christos Doulkeridis

Department of Digital Systems
University of Piraeus (UPRC), Greece

# Contents

# Introduction/Motivation (1/2)

- In recent years the increasing size of the spatio-temporal data requires new approaches for their storage and retrieval

- Scalable querying of spatio-temporal data management is a challenging topic

- Uber reports Millions of trips per quarter year;

*Uber Technologies, Inc. Q4 2020 Earnings, Supplemental Data*



| Q4 2019 | Q1 2020 | Q2 2020 | Q3 2020 | Q4 2020 |
|---------|---------|---------|---------|---------|
| 1,907 | 1,661 | 737 | 1,184 | 1,443 |

# Introduction/Motivation (2/2)

- NoSQL stores are exploited by modern applications for data storage and querying, providing <span style="color:red">scalability</span> and <span style="color:red">availability</span>
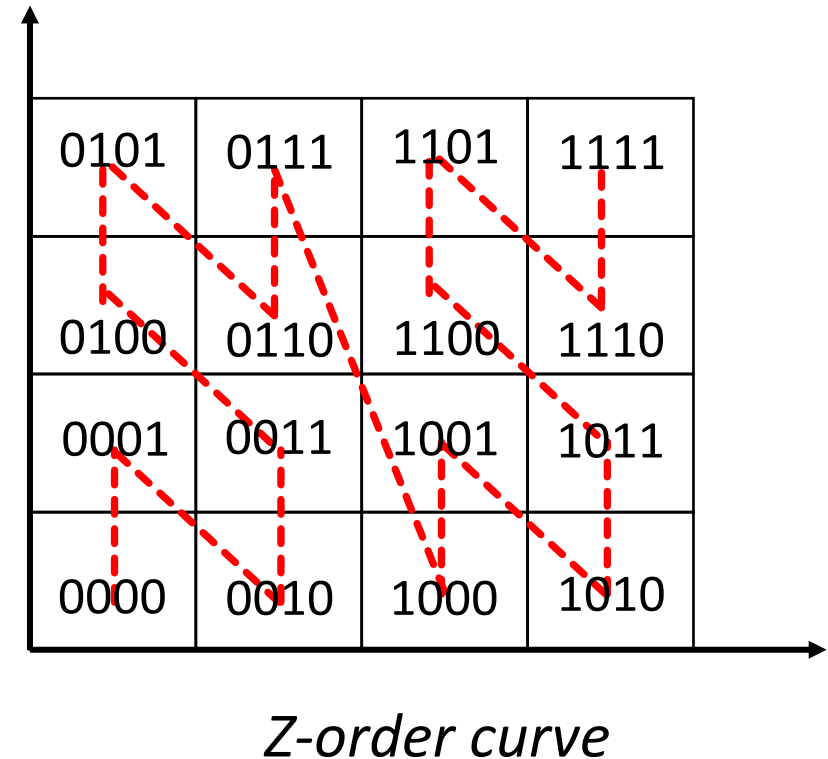


- Despite the popularity of NoSQL systems, they are not optimized for spatial data

- In this work we opt for MongoDB store to support efficient spatio-temporal querying, as it embeds geospatial features
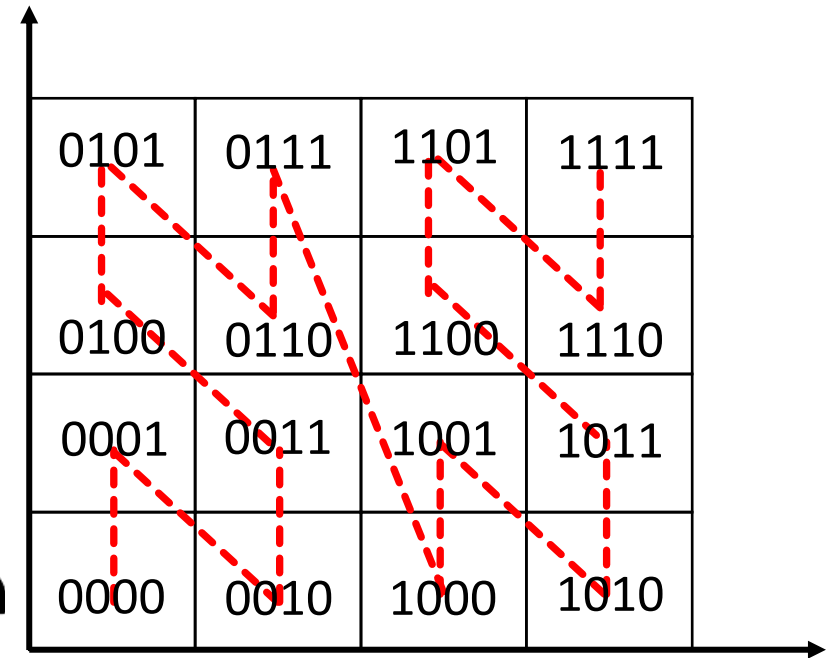
# Baseline (1/3)

- MongoDB offers spatial indexes (2d/2dsphere)

- The geospatial **2dsphere** index is based on geohashing

  **{location: "2dsphere"}**

- The geohash values are indexed by *B-trees*

```
{
 _id: 1,
 location: {"type": Point, coordinates: [37.983810, 23.727539]},
         ...
}
```



*Z-order curve*

# Baseline (2/3)

- MongoDB also offers compound indexes for indexing two or more fields in a single structure

- A spatio-temporal index in MongoDB is created in the two following ways:
  - `{ location: "2dsphere", date: 1 }`
  - `{ date: 1, location: "2dsphere" }`



```
{
 _id: 1,
 location: {"type": Point, coordinates: [37.983810, 23.727539]},
 date: ISODate("2018-09-12T12:15:17.777Z"),
      ...
}
```
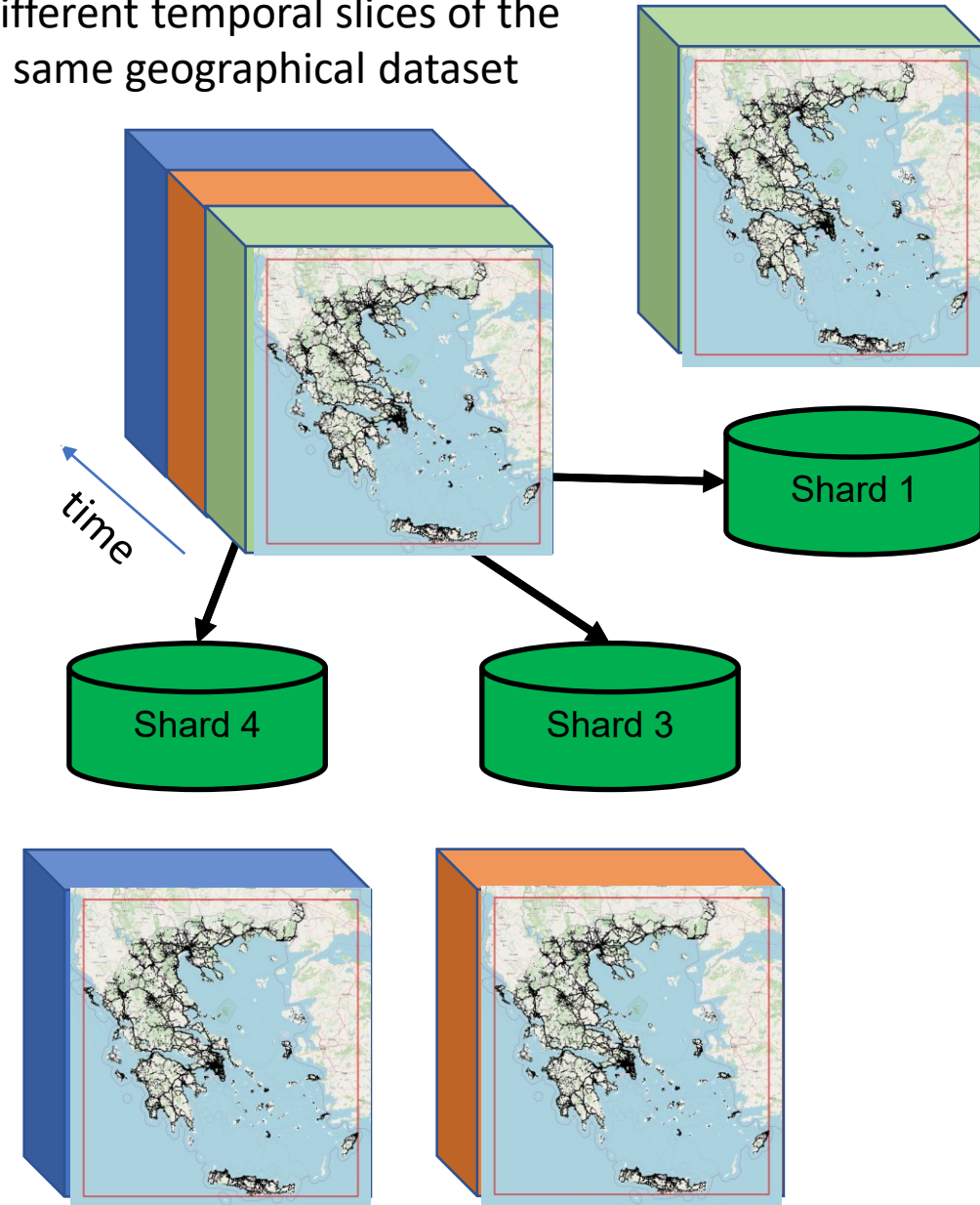
*Z-order curve*

# Baseline (3/3)

- However, MongoDB cannot distribute documents on a shards based on their spatial information

- Thus, we are restricted to the integration of the temporal field only as a shard key for spatio-temporal querying:

$$\texttt{\{date: 1\}} \quad \textbf{(bsl)}$$

- On each shard we create one of the following compound indexes (co-existing with a local single index based on the date field)

- `{ location: "2dsphere", date: 1 }` **(BslST)**

- `{ date: 1, location: "2dsphere" }` **(BslTS)**

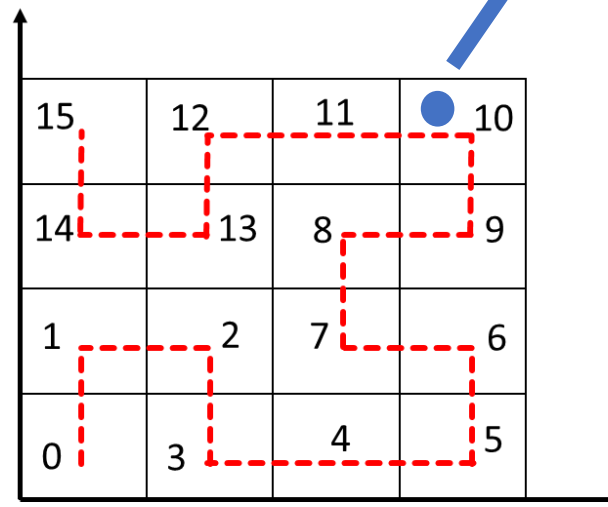Different temporal slices of the same geographical dataset

# Our approach (1/2)

- We exploit the Hilbert space-filling curve and integrate the 1D numeric value in a new field in each document

- We form the compound indexes based on the fields:

  **{hilbertIndex: 1, date: 1}**

```
{
 _id: 1,
 location: {"type": Point,
        coordinates: [37.983810, 23.727539]},
 date: ISODate("2018-09-12T12:15:17.777Z"),
 hilbertIndex: NumberLong(10),
 ...
}
```
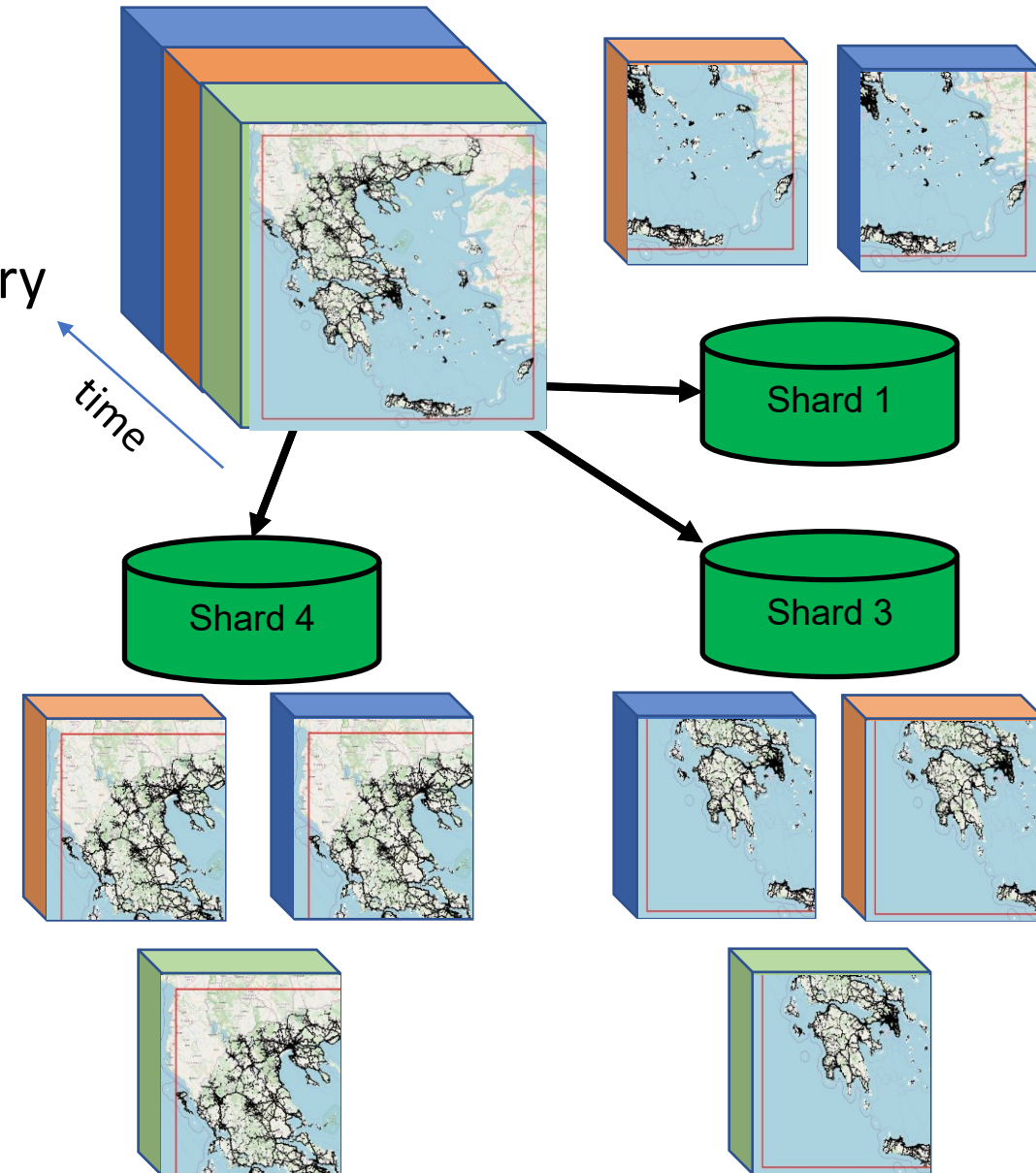


*Hilbert curve*

# Our approach (2/2)

- Since the spatial information is embedded in every document as a numeric field, we use the shard key:
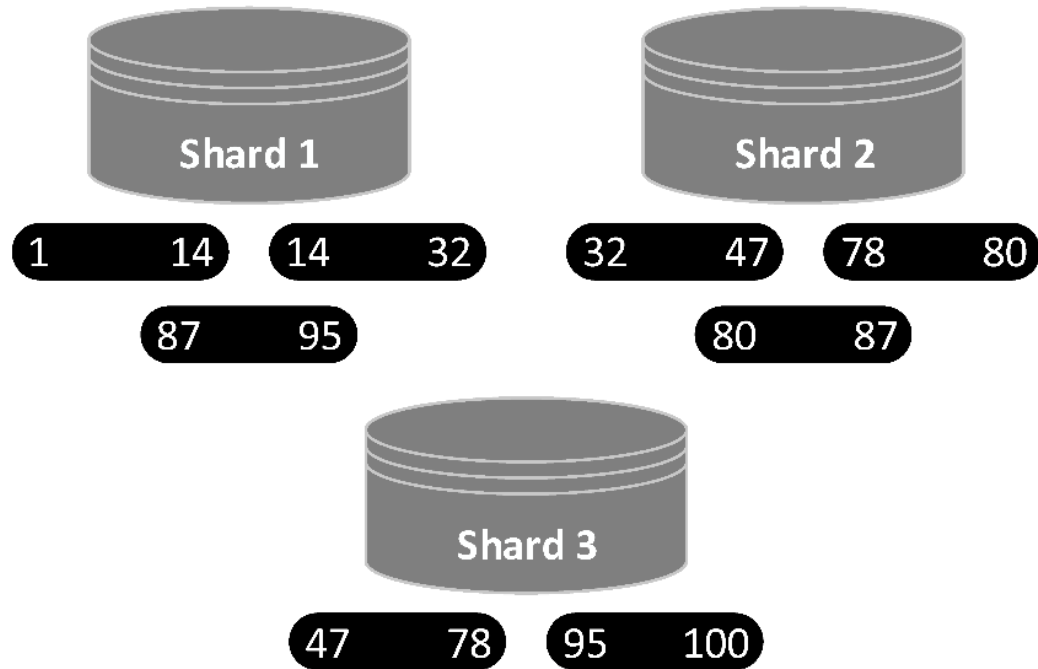
   **`{hilbertIndex: 1, date: 1}`** **(hil)**

- By default, each shard has compound indexes based on the declared fields in the shard key

- Shard overloading is unlikely to occur in case of spatial/spatio-temporal skewness
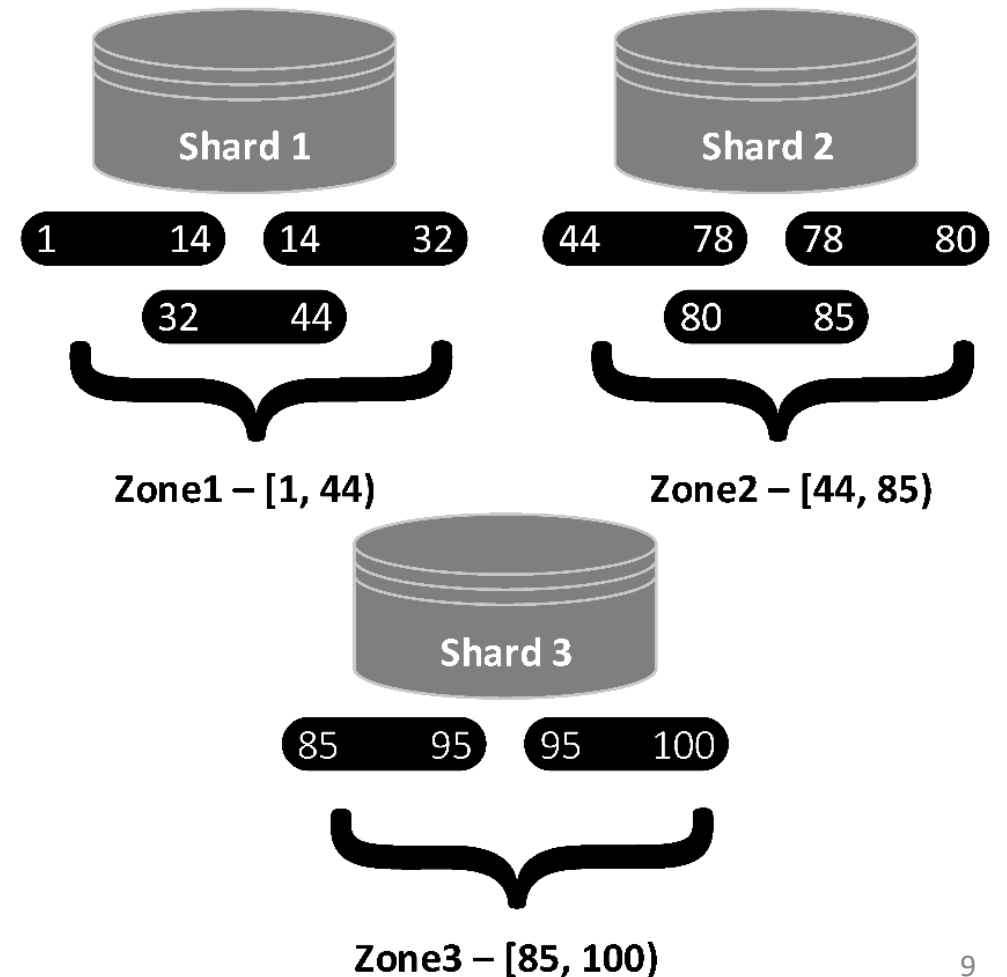
Different temporal slices of the same geographical dataset

# Usage of zones

**Default distribution**

**Zone usage**



Shard 1

| 1 | 14 | | 14 | 32 |

| 87 | 95 |

Shard 2

| 32 | 47 | | 78 | 80 |

| 80 | 87 |

Shard 3

| 47 | 78 | | 95 | 100 |

Continuous ranges of shard keys (*chunks*) are distributed evenly in the cluster

Shard 1

| 1 | 14 | | 14 | 32 |

| 32 | 44 |

Zone1 – [1, 44)

Shard 2

| 44 | 78 | | 78 | 80 |

| 80 | 85 |

Zone2 – [44, 85)

Shard 3

| 85 | 95 | | 95 | 100 |

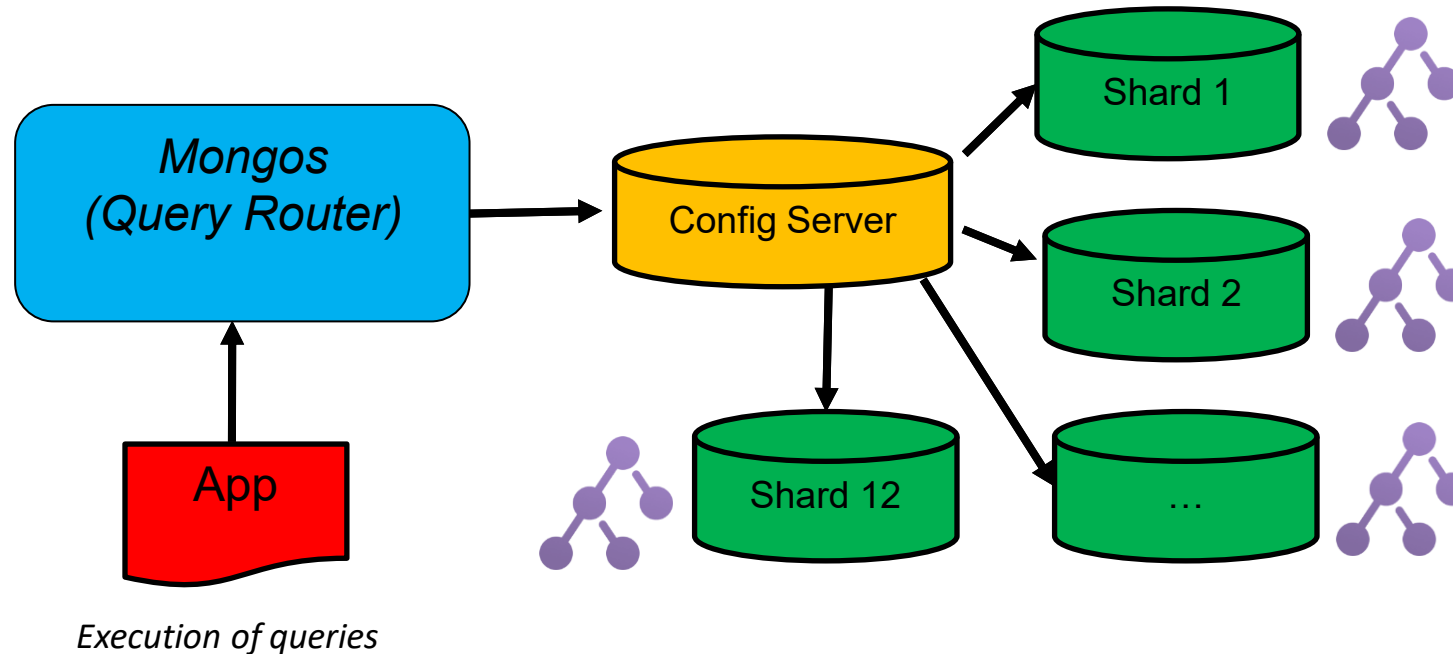Zone3 – [85, 100)

# Experiments (1/5) – Setup

- Assess the efficiency of the proposed Hilbert-based approach against the baseline on a MongoDB Cluster

- We consider two types of spatio-temporal queries; the "spatially small" ( $Q_x^s$ ) and the "spatially big" ( $Q_x^b$ )

| Real data set | |
| --- | --- |
| *Size in MongoDB* | *#points* |
| **~40GB** | **15.2M** |

| Time periods |
| --- |
| $Q_1^x$ – One hour |
| $Q_2^x$ – One day |
| $Q_3^x$ – One week |
| $Q_4^x$ – One month |

# Experiments (2/5) - Setup



**Mongos (Query Router)**

**App**

*Execution of queries*

**Config Server**

**Shard 1**

**Shard 2**

**Shard 12**

**...**
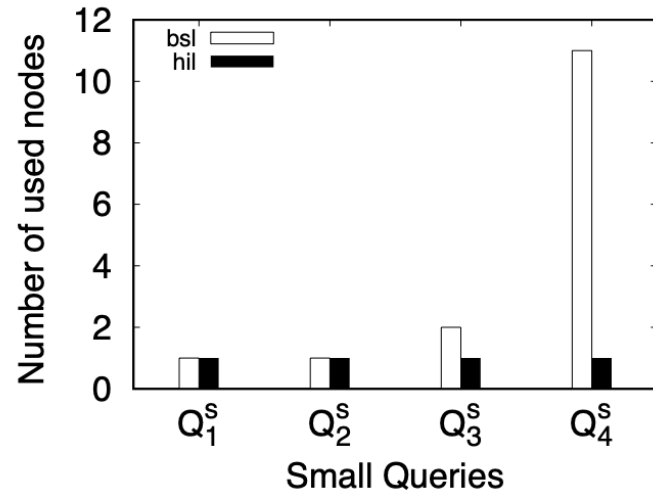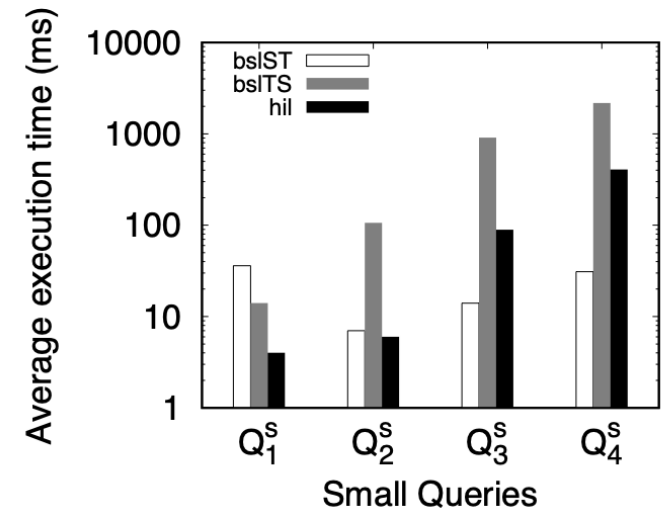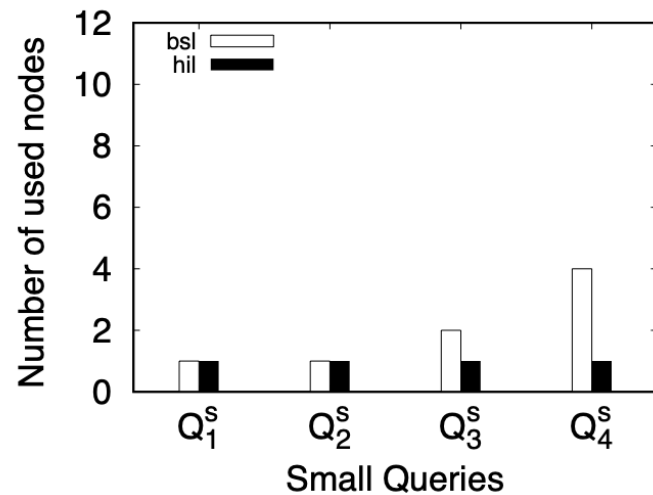
17 nodes, 8GB RAM, x4 CPU cores,
12 nodes are used as shards with 100GB disk

# Experiments (3/5) - Performance of "spatially small" ST queries
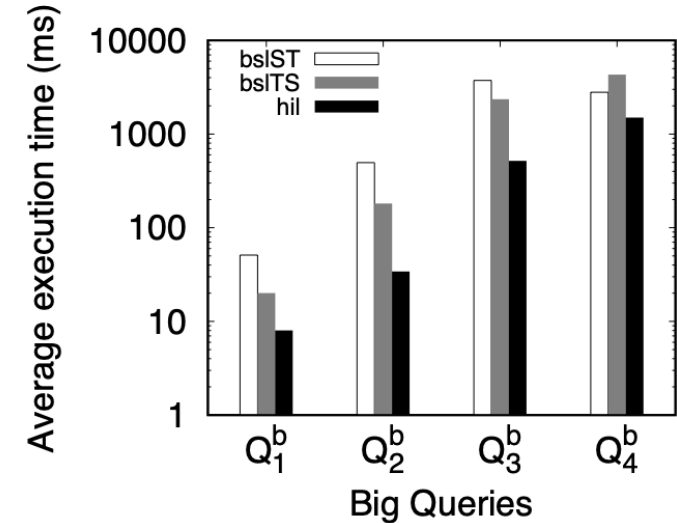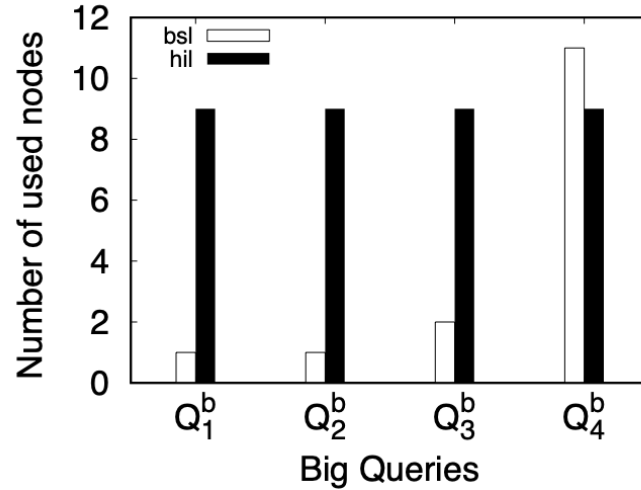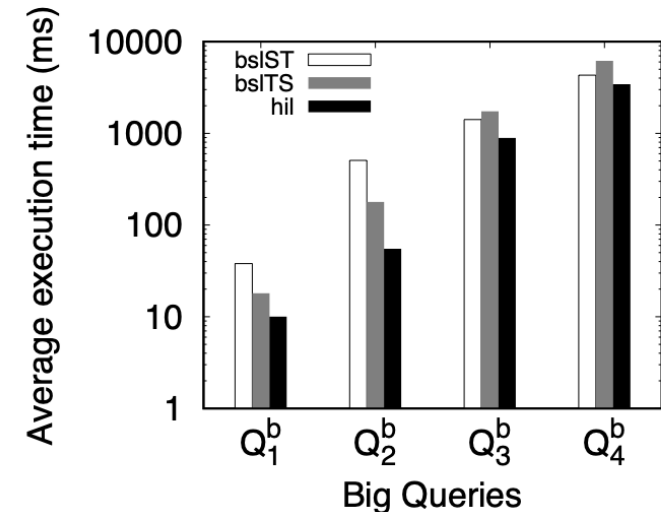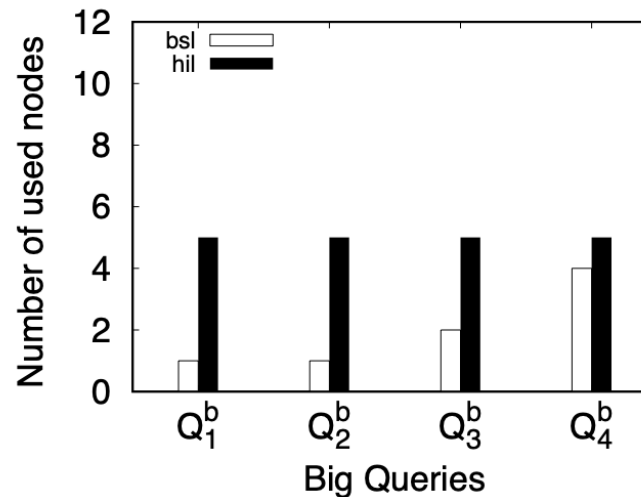
Default distribution

Zone usage

# Experiments (4/5) - Performance of "spatially big" ST queries

Default distribution

Zone usage

# Experiments (5/5) - Scalability study

- We record the performance of $Q_2^b$

| | Real data set | | |
|---|---|---|---|
| Size Factor | Size (GB) in MongoDB | #points (M) |
| x1 ($R_1$) | ~40 | 15.2 |
| x2 ($R_2$) | ~83.87 | 31.4 |
| x3 ($R_3$) | ~127.21 | 47.7 |
| x4 ($R_4$) | ~171.39 | 63.9 |

# Conclusions

- We proposed a spatio-temporal approach for indexing and sharding on MongoDB

- The referred approaches were collated and evaluated on a MongoDB cluster

- Demonstrated the advantages of the proposed approach on specific queries

- Can be easily adopted as a solution on top of MongoDB

# Thank you for your attention

**_More info_**:

our group: http://www.datastories.org/

project Track & Know: https://trackandknowproject.eu/

project Chorologos: https://www.ds.unipi.gr/chorologos/

e-mail: koutroumanis@unipi.gr