**Olha Horlova**
PhD candidate, 3rd year
Advisor: **Stefano Ceri**
Co-advisor: **Abdulrahman Kaitoua**
Thesis submission: May 2020
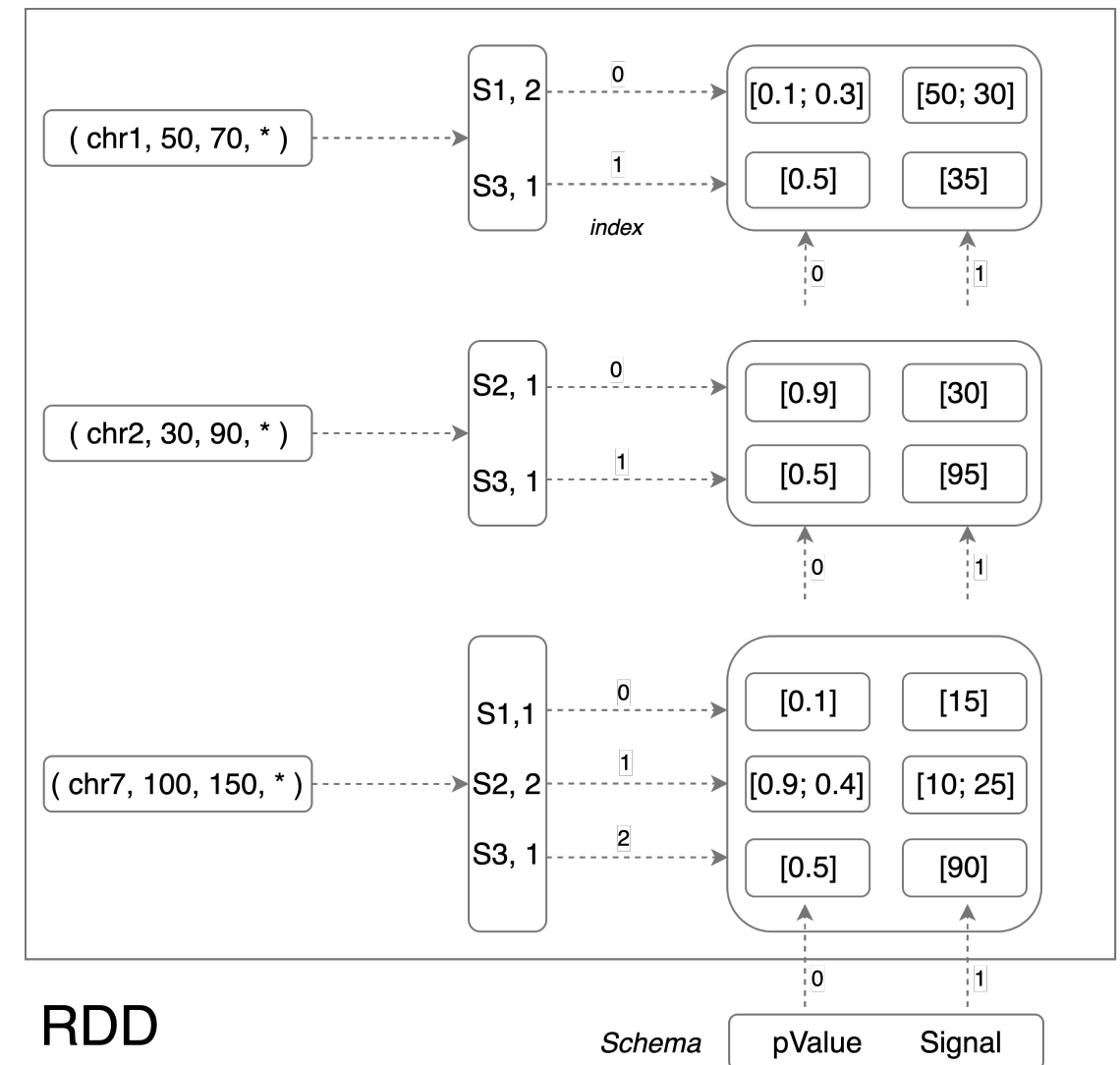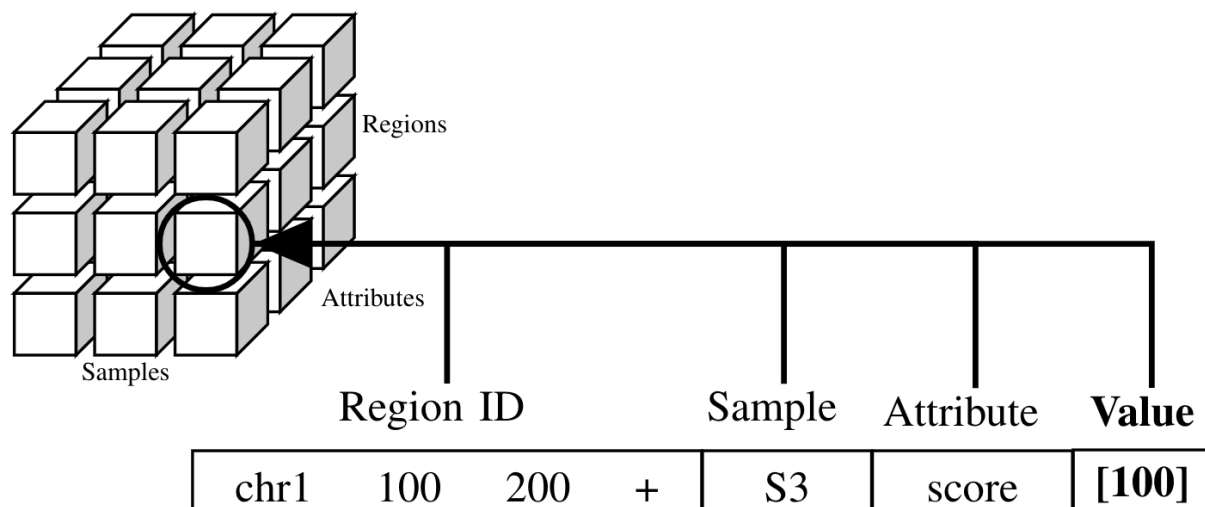
# THE EVOLUTION OF GENOMIC DATA MODEL FOR THE CLOUD

## Publications:

▸ **Olha Horlova**, Abdulrahman Kaitoua, Volker Markl, Stefano Ceri. Multi-Dimensional Genomic Data Management for Region-Preserving Operations. In Proceedings of the 35th IEEE International Conference on Data Engineering (ICDE 2019). Macau SAR, China. 8-11 April, 2019. DOI: https://doi.org/10.1109/ICDE.2019.00107

▸ **Olha Horlova**, Abdulrahman Kaitoua, Stefano Ceri. Array-based Data Management for Genomics. Accepted for the 36th IEEE International Conference on Data Engineering (ICDE 2020). Dallas, TX, USA. 20-24 April, 2020.

▸ Masseroli M, Canakoglu A, Pinoli P, Kaitoua A, Gulino A, Horlova O, Nanni L, Bernasconi A, Perna S, Stamoulakatou E, Ceri S. Processing of big heterogeneous genomic datasets for tertiary analysis of Next Generation Sequencing data. Bioinformatics (Oxford, England), 2018. DOI: https://doi.org/10.1093/bioinformatics/bty688

POLITECNICO
MILANO 1863

# Genomic Data Model: Row to Array

# Spark RDDs organization to support array model



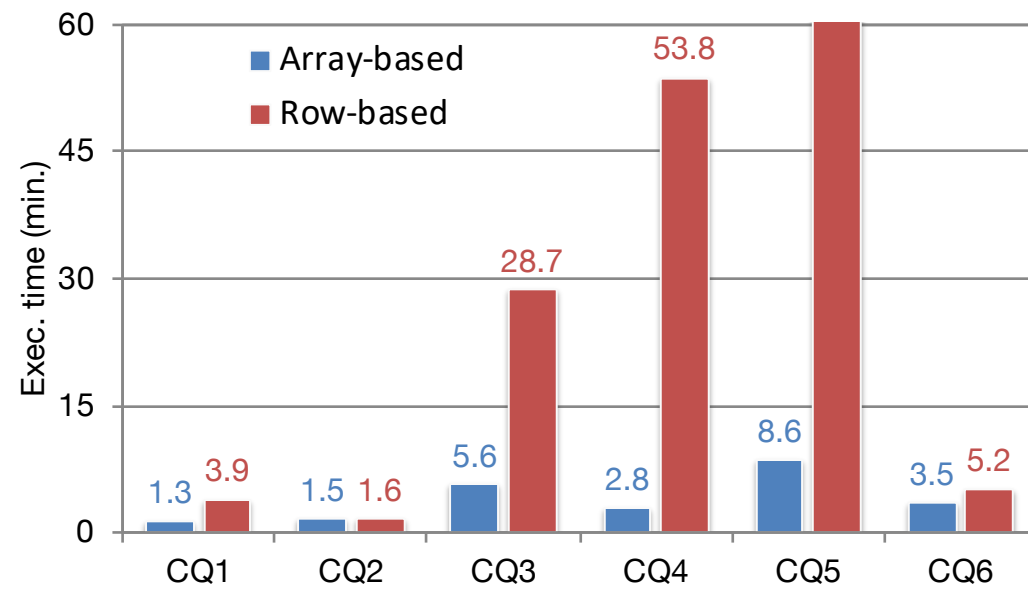| Region  ID | | | | Sample | Attribute | **Value** |
|---|---|---|---|---|---|---|
| chr1 | 100 | 200 | + | S3 | score | **[100]** |

RDD

```
1  ArrayModel(key: RegionKey, value: RegionData)
2
3  RegionKey(chrom: String, start: Long, stop: Long, strand: Char)
4
5  RegionData(Replication: Array[(Long, Int)], Attribute: Array[Array[Array[GValue]]])
```
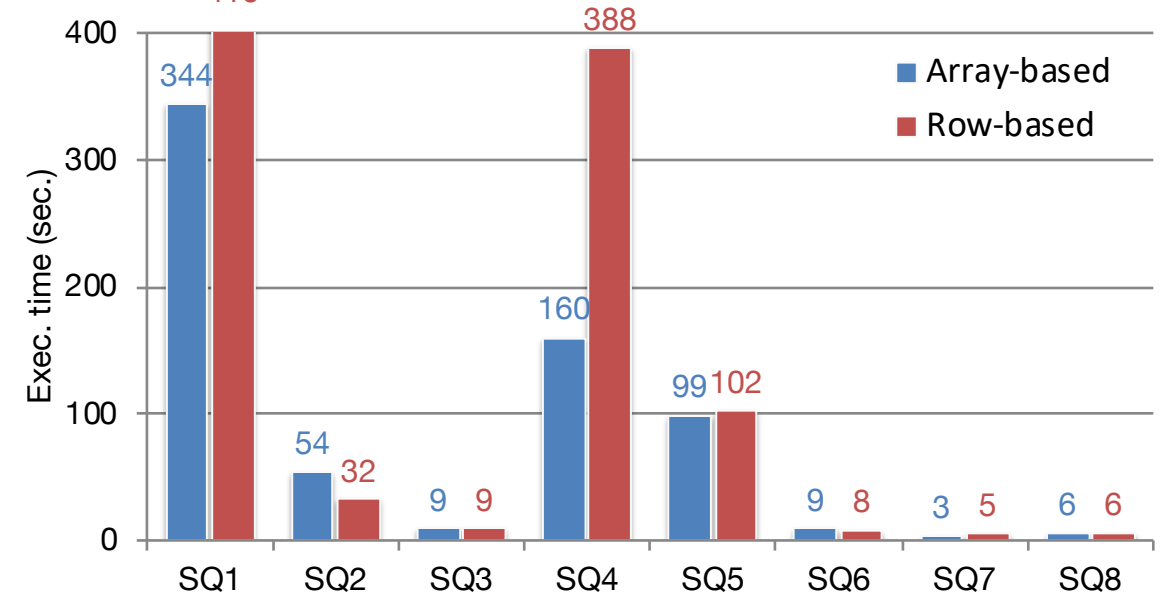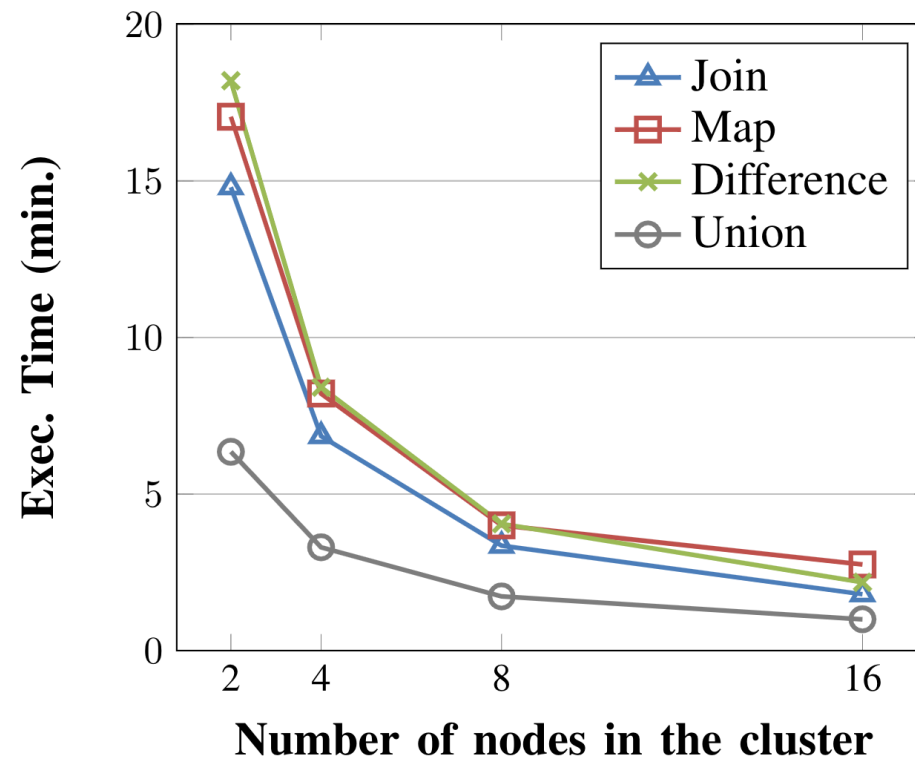
# Model Evaluation



Complex queries — Exec. time (min.): Array-based, Row-based
CQ1: 1.3, 3.9; CQ2: 1.5, 1.6; CQ3: 5.6, 28.7; CQ4: 2.8, 53.8; CQ5: 8.6, 300; CQ6: 3.5, 5.2

Simple queries — Exec. time (sec.): Array-based, Row-based
SQ1: 344, 410; SQ2: 54, 32; SQ3: 9, 9; SQ4: 160, 388; SQ5: 99, 102; SQ6: 9, 8; SQ7: 3, 5; SQ8: 6, 6

Scaling — Exec. Time (min.) vs Number of nodes in the cluster: Join, Map, Difference, Union

# Towards spatial and temporal applications

We can map genomic coordinates to:

▸ longitude and latitude of locations in spatial data

▸ time intervals of temporal data

Example #1: Find minimum distance offices of public or private organizations closest to given locations (e.g. for all banks, the closest bank office from home)

Example #2: Find the closest time events in different countries when a certain climate event occurred (e.g., for each nation/state/region, the event closest in time to Xmas 2019 when temperature was higher than 40 degrees Celsius)

# THANK YOU!