



Unveiling Objects in Big Data

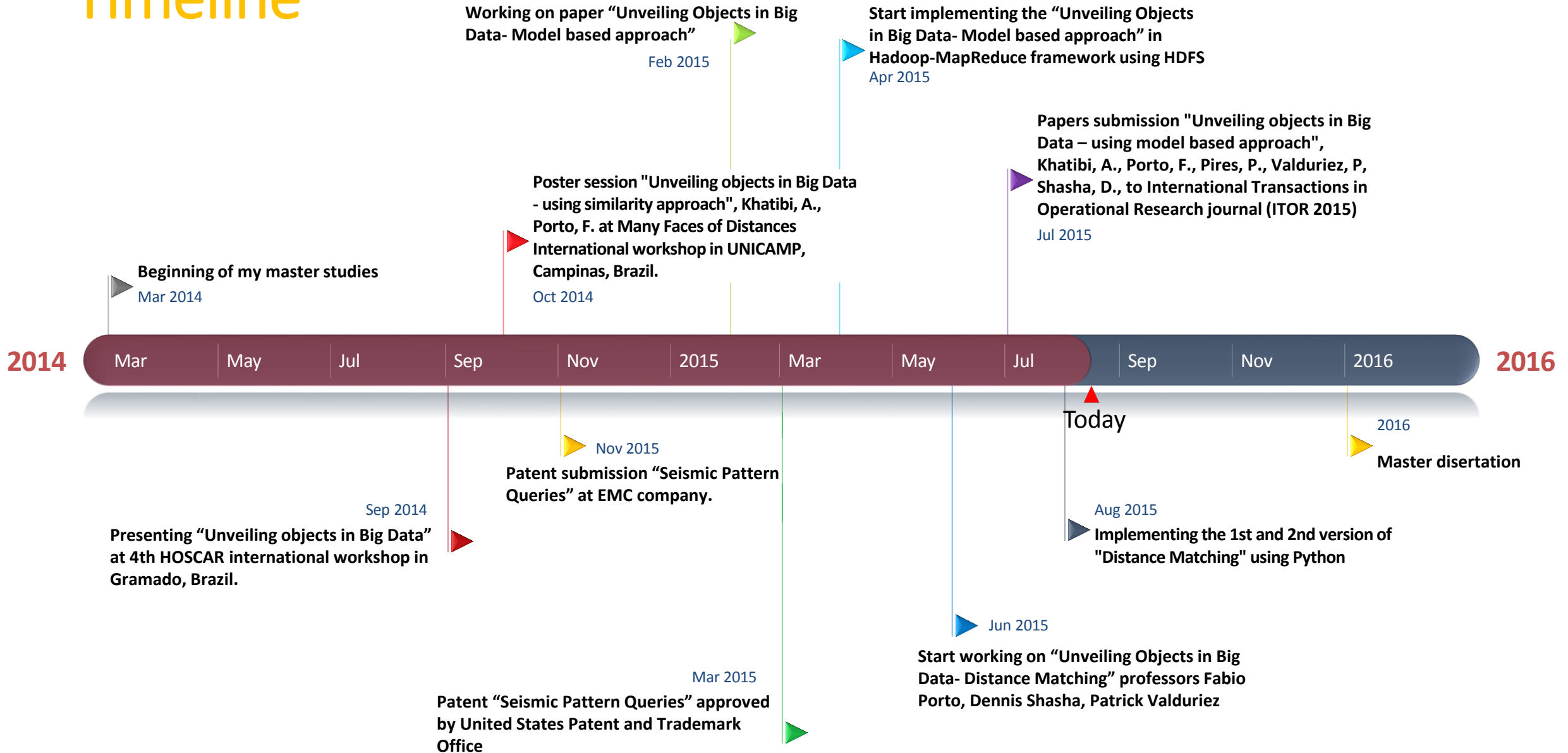
Distance Matching

Supervisor: Prof. Fabio Porto
Prof. Dennis Shasha, Prof. Patrick Valduriez

By Amir Khatibi

Aug 2015

Timeline



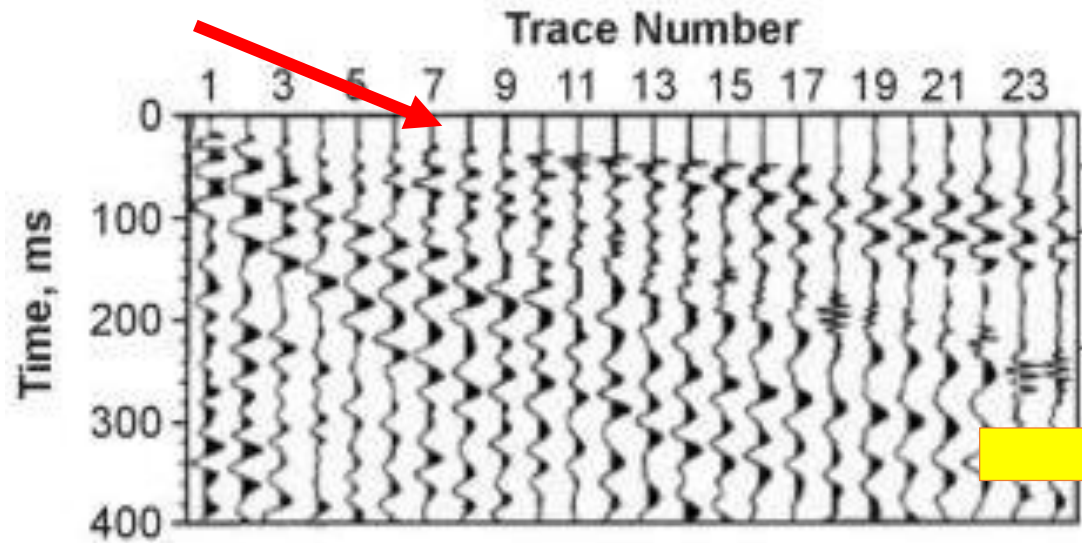
Unveiling objects

- Data deluge produced by the Big Data phenomenon, blurs **high-level objects** amongst billions of dataset elements.
- A **Sample Query** or a **Descriptive Language** specifies the elements that shall compose higher-level objects, and a composition model.

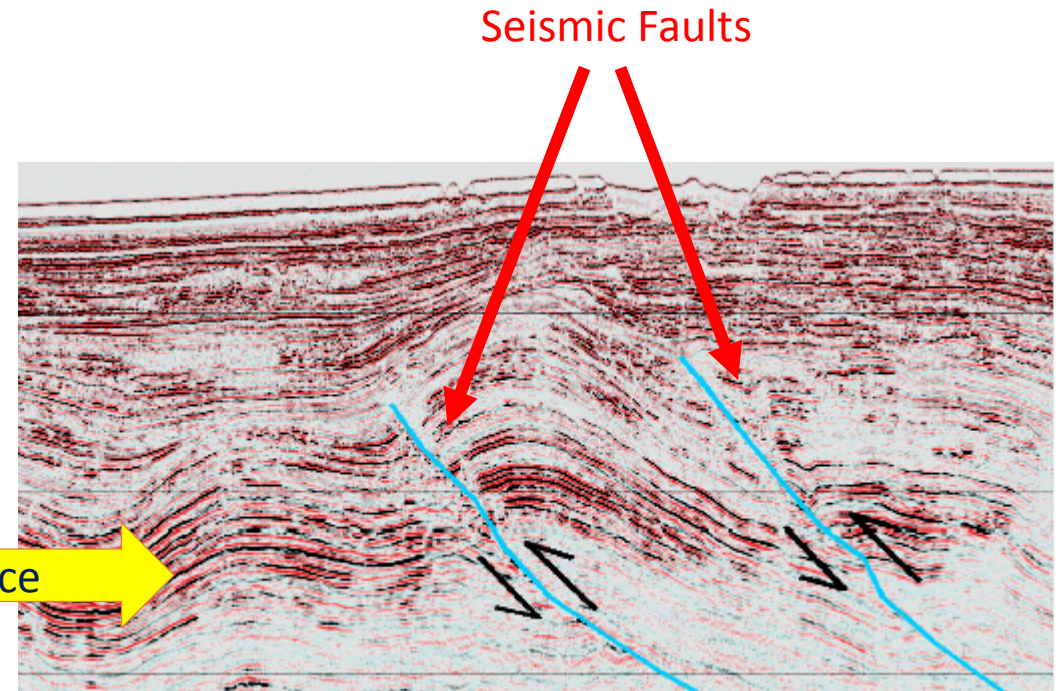
Unveiling objects – Seismic Data

- The **objects of interest** are components built from the elements held by the target dataset. Those features maybe obtained from the dataset through a **smart combination of low-level elements**.

Seismic Dataset



Inference

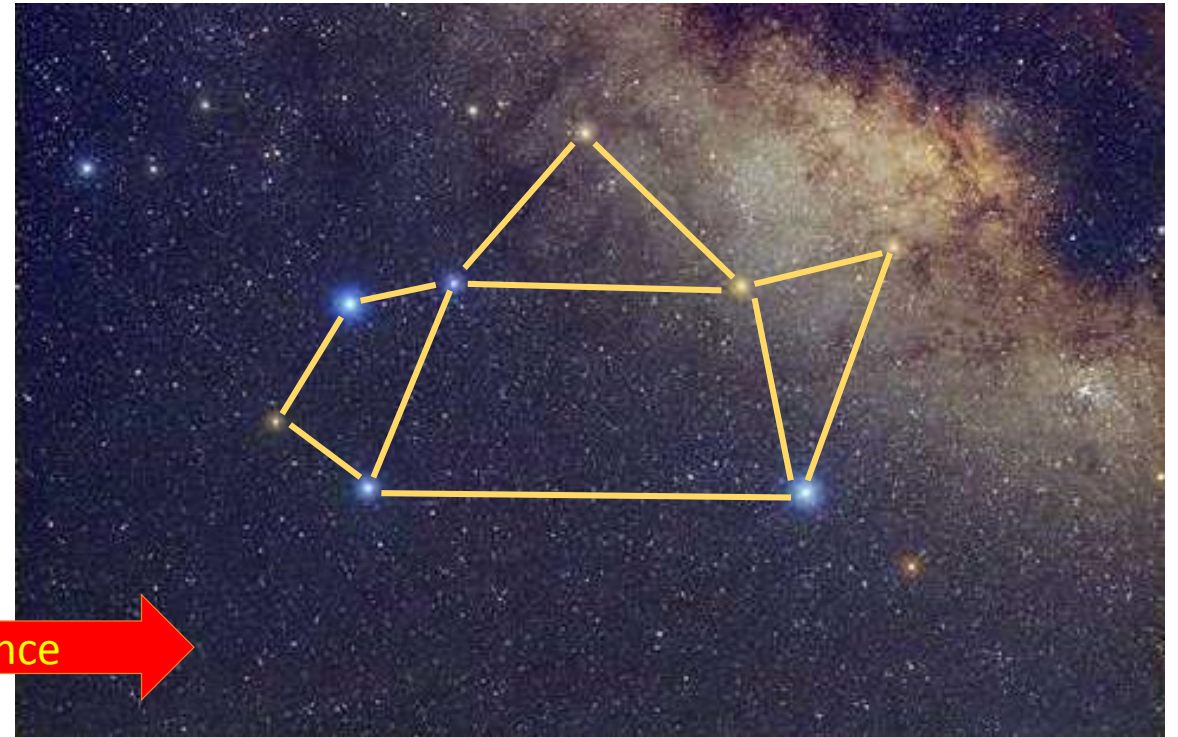


Unveiling objects – Astronomy

low level objects: stars and their attributes



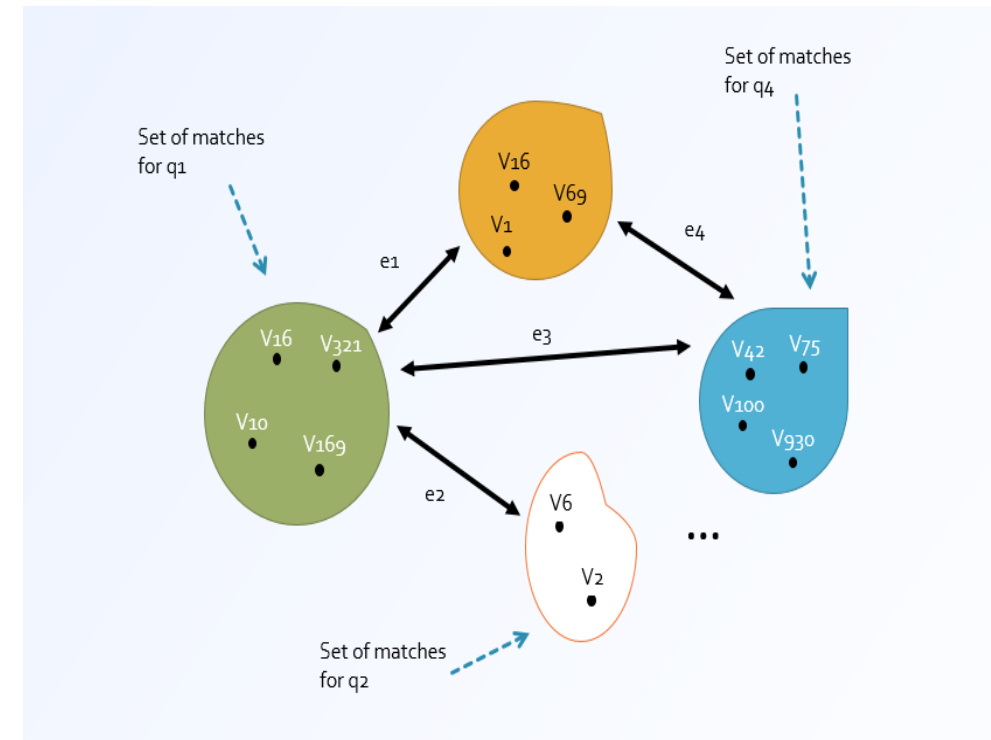
Higher level object: constellation



Inference

Unveiling objects

- Distance queries are important in many areas like:
 - Astronomy catalogues
 - Seismic Data
 - Environmental sensor data
 - ...
- We look for **efficient strategies** to search for such elements in huge datasets
- In the case of Big data, processing the **distributed data store** in a **parallelized environment** is obvious.



Astronomy application – sample query

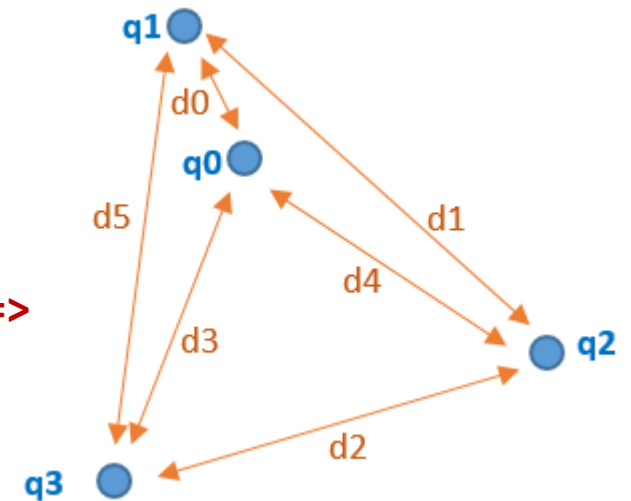
- Astronomy catalogue:

Cat (ra, dec, flux, photo-z, u, g, r, i, z,...)

magnitudes of light emitted by an sky object

position of sky objects specifies by the values of their coordinate in right-ascension (ra) and declination (dec).

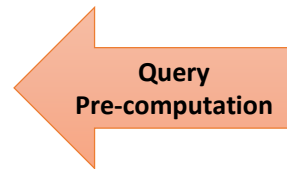
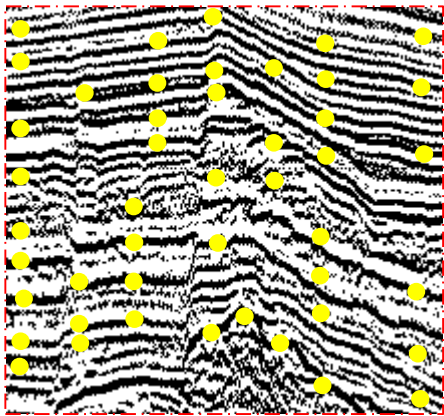
A sample query =>



Pre-computing the Query

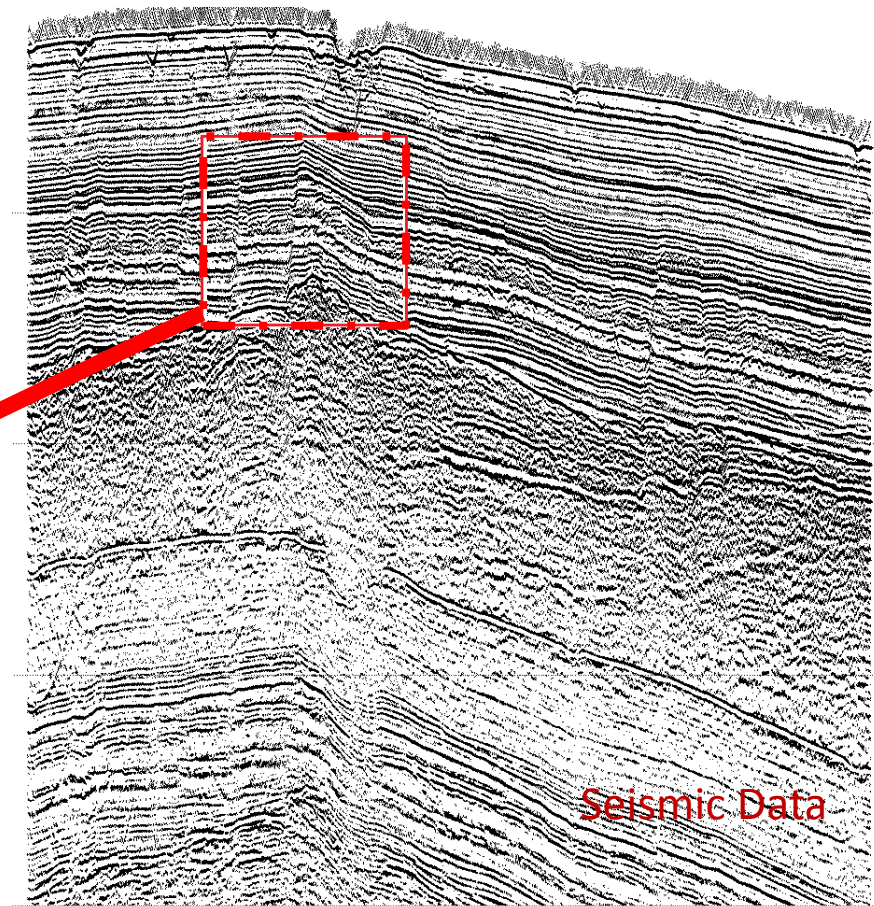
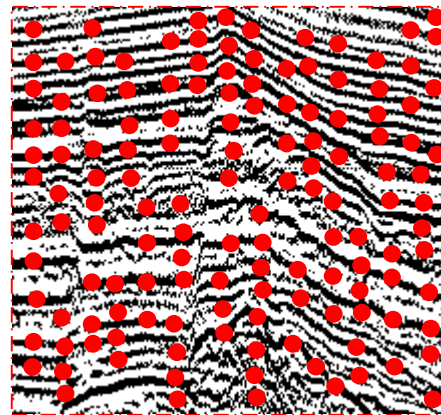
- Find the k most representative points in the query
 - Relieve us from many computations.

Refined query



Query
Pre-computation

Geologist's query

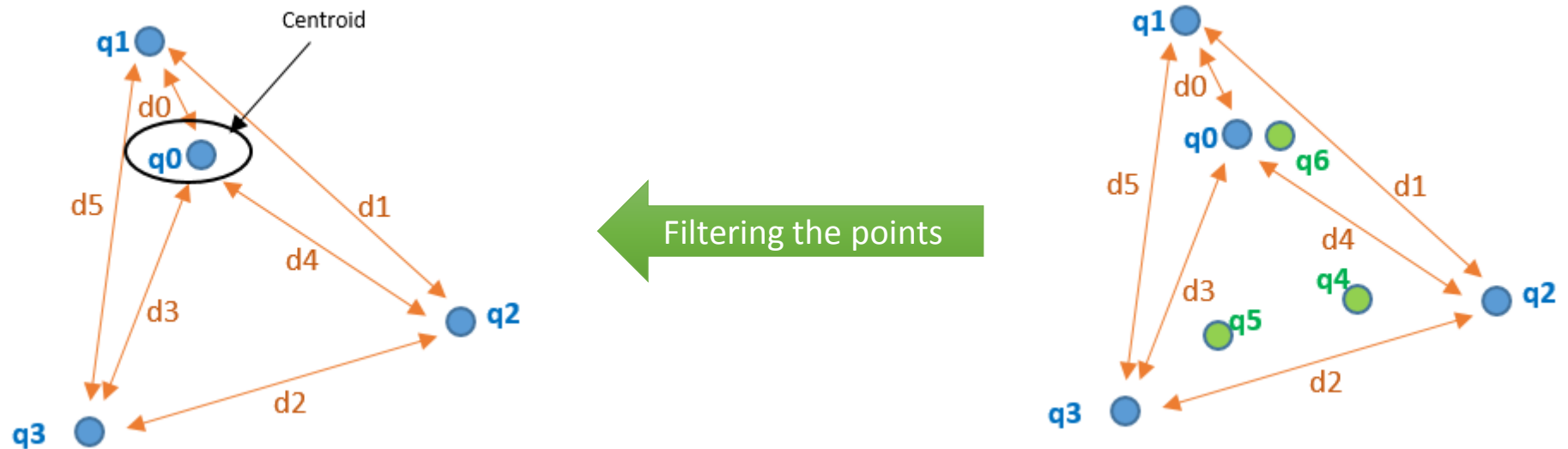


Seismic Data

ex: In a query with 50% of points very similar to each other, we can avoid from 50% of the computations in whole.

Pre-computing the Query

- Filter some small number of points that are far apart from each other.
 - Do the distance match for these points, firstly.
- Find the centroid point in the query
 - The nearest point to the Query's center of mass



Centroid-based approach

In this approach, each star is a potential centroid of a shape. Initial matching can be applied between the query centroid and the potential centroid of a shape:

- Query_centroid_distances (centroid_id, <<d1,q_1 >, ..., <dk,q_k >>)

Centroid's distances from other query elements

=> Pre-computing the dataset

(centroid_id,<<d1,star-id_1,ra_1,dec_1>,... ,<dn,star-id_n,ra_n,dec_n>>)

Pre-Computing the Dataset

- Pre-compute the distances of every sky object from other objects in a specific boundary **d-radius** of it and store an ordered list of distances and points
- Using a spatial index like **PH-Tree** to quicken the execution time.

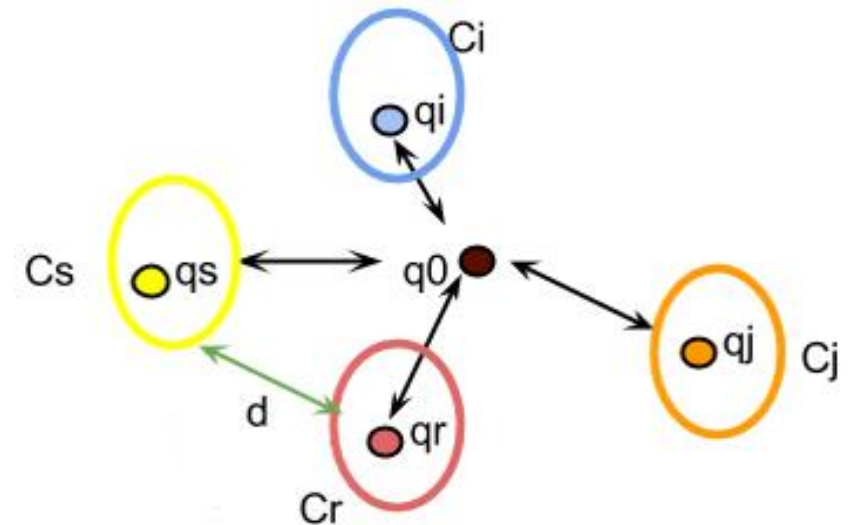


Initial matching

- We look for candidate matches with respect to the query centroid and the result would be stored in the following relation. Every centroid could have a set of candidates for p-th element in the query.

C_p (Bucket p)

(centroid_star_id, partner_star_id, Ra_partner, Dec_partner)



Pairwise distances

The second step is checking the pairwise distances in filtered centroids of a shape with respect to the query pairwise distances:

Element's distances from other query elements

- Query_pairwise_distances ($\langle\langle d_{ij}, q_i, q_j \rangle\rangle$)

=> Nested loops or Join Matrices

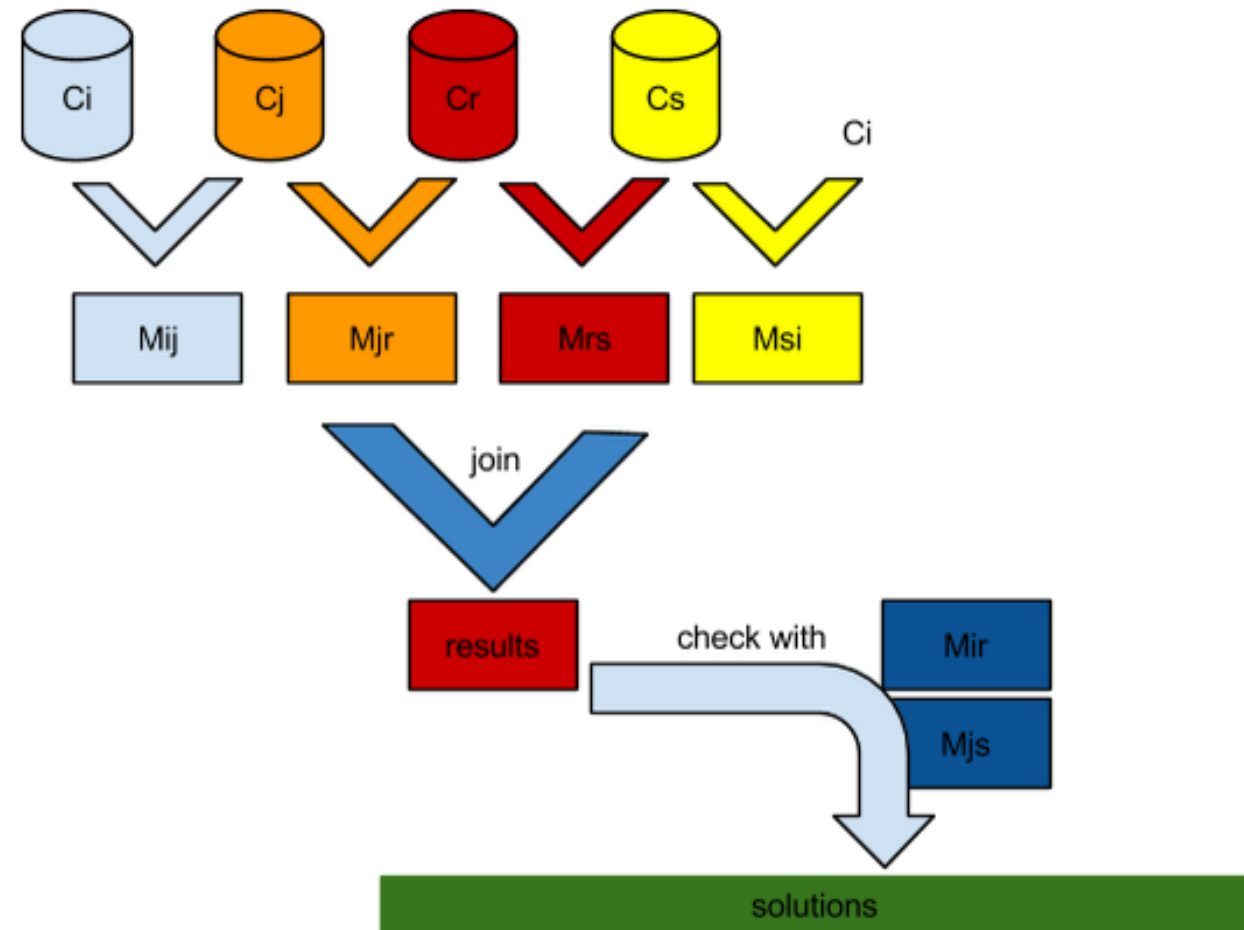
Matrices

- The idea of matrices aims to check the pairwise distances, **partially**.

- We create a set of matrices **Mij**

$M_{ij}(a,b) = 1$ iff:

dist (star a in C_i , star b in C_j) similar to dist (q_i, q_j)



Form the shapes

- According the density of intermediate results and the total number of solutions, we can do it differently:
 - **Nested Loops**
 - Using in-memory variables and for loops
 - **Nested Joins in SQLite**
 - Using Matrices and in-memory SQLite to join
 - **Matrix Multiplication**
 - Using binary matrix multiplication to filter the candidate centroids

Some Results

- Data-set size: **1000 sky objects**
- Pre-computing the dataset: **41 seconds**
- Average number of **neighbours** for each centroid: **170** (between 90 to 240)

(All the following results are only for the first 50 candidate centroids using **Python PL**)

Method		avg(Cp)= 10 Dense matrices (lots of solutions)	avg(Cp)= 10 sparse matrices	avg(Cp)= 102 Dense matrices (lots of solutions)	avg(Cp)= 102 sparse matrices
		#solutions: 72,187	#solutions: 39	#solutions: 29,790,096	#solutions: 60
Nested loops	Without matrices	1.35 s	65 ms	550 s	2.4 s
MM Filtering	sparse matrix multiplication	130 ms	99 ms	5.3 s	3.9 s
Matrices and Sqlite	with indices	1.11 s	111 ms	?more than 20 min	3.9 s
	without indices	297 ms	105 ms	51 s	3.7 s
Prediction		MM>NL	MM<NL	MM>NL	MM<NL
Conclusion		Totally True			

Thank you for your attention

ANY QUESTIONS?