



Progetto Batch processing

SABD 2018-2019

Montesano, Perrone, Pusceddu



1 Architettura

Data Ingestion

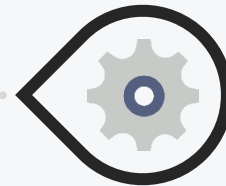
2



3 Query1 Core & SQL

Query2 Core & SQL

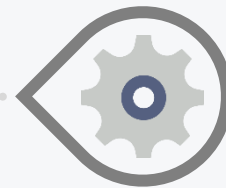
4



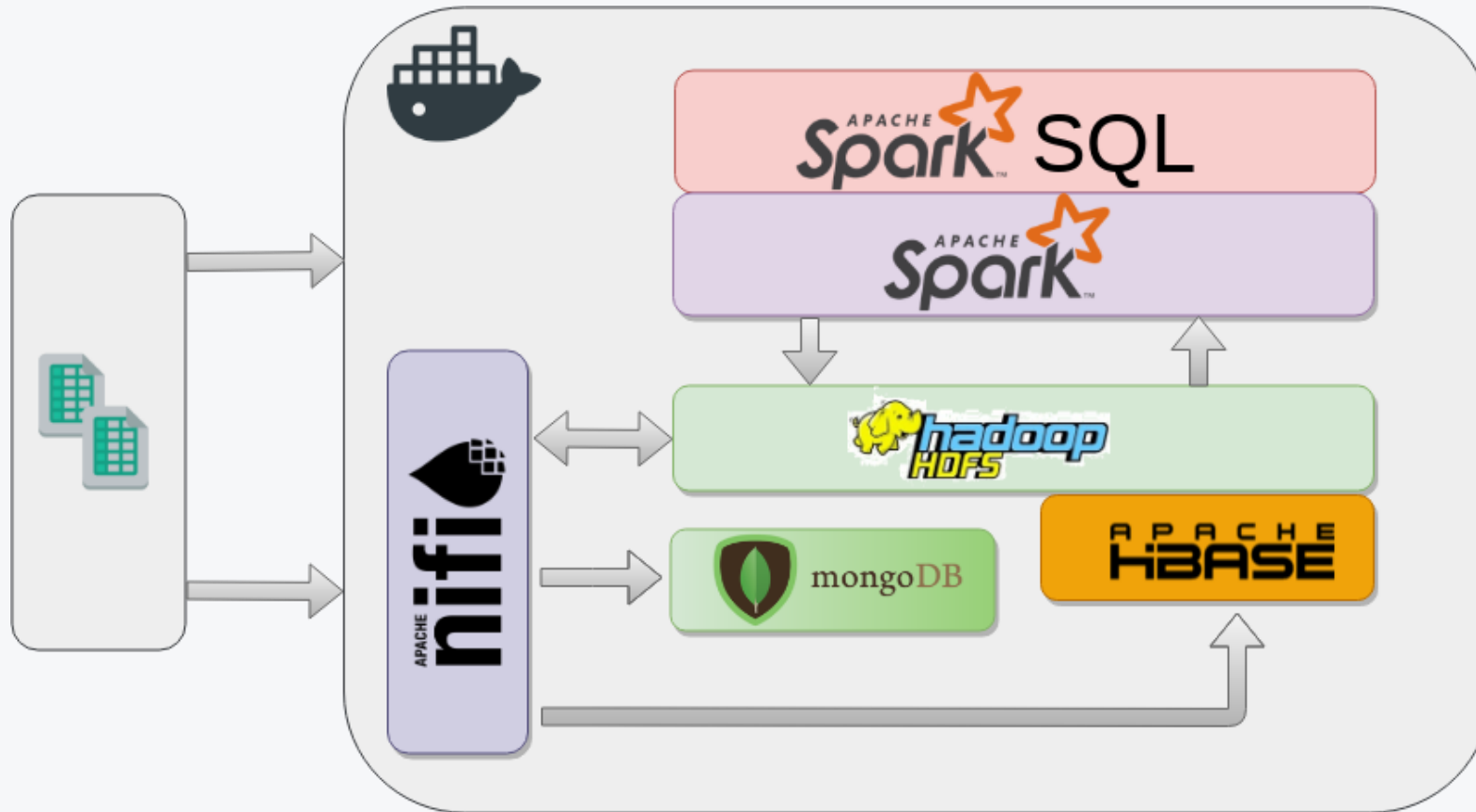


5 Query3 Core & SQL

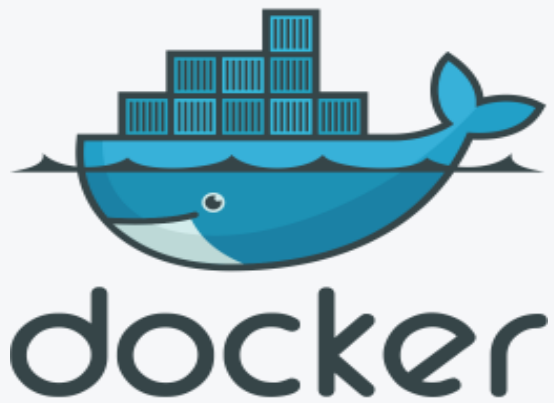
Discussione dei risultati 6



Architettura



- Apache Spark esegue con un master e due worker con cluster manager locale
- Ogni worker ha 1Gb di memoria
- Linguaggio utilizzato: Java
- Hdfs esegue con un master e 3 worker



Local Deployment

- Il deploy in locale è costruito utilizzando le Docker Image
- Vengono avviati i container di : Hdfs, Nifi, Hbase, MongoDB ,Spark
 - I container lavorano in modalità cluster
- I container vengono eseguiti sotto la stessa rete denominata: net



Data Ingestion

- L'ingestion dei dati è stato realizzato con il framework Apache Nifi
- Effettua l'iniezione dei dati nel file system distribuito HDFS in due diversi formati : csv e parquet
- Trasferisce i risultati delle query verso due database non relazionali Apache Hbase e MongoDB

Scrittura risultati

I risultati vengono salvati su due database in formato Json. Sono stati usati sia Hbase e MongoDB a scopo didattico.

- Per MongoDB è stata creata una collection diversa per ogni query per ogni modalità di spark: Core ed SQL



- Per Hbase sono state create 3 tabelle, una per query
- Le tabelle per la query 1 e query 3 sono state divise in due famiglie di colonne: Core ed Sql
- La tabella della query 2 è stata divisa in 6 famiglia di colonne, 3 per i file di Spark Core e 3 per i file di Spark SQL

QUERY



1

Per ogni anno del dataset individuare le città che hanno almeno 15 giorni al mese di tempo sereno nei mesi di marzo, aprile e maggio.

2

Individuare, per ogni nazione, la media, la deviazione standard, il minimo, il massimo della temperatura, della pressione e dell'umidità registrata in ogni mese di ogni anno.

3

Individuare, per ogni nazione, le 3 città che hanno registrato nel 2017 la massima differenza di temperature medie nella fascia oraria locale 12:00-15:00 nei mesi di giugno, luglio, agosto e settembre rispetto ai mesi di gennaio, febbraio, marzo e aprile. Confrontare la posizione delle città nella classifica dell'anno precedente (2016).

Pre-Processing

Step by step



Fase di pre processamento comune a tutte le query per ripulire i dati spuri

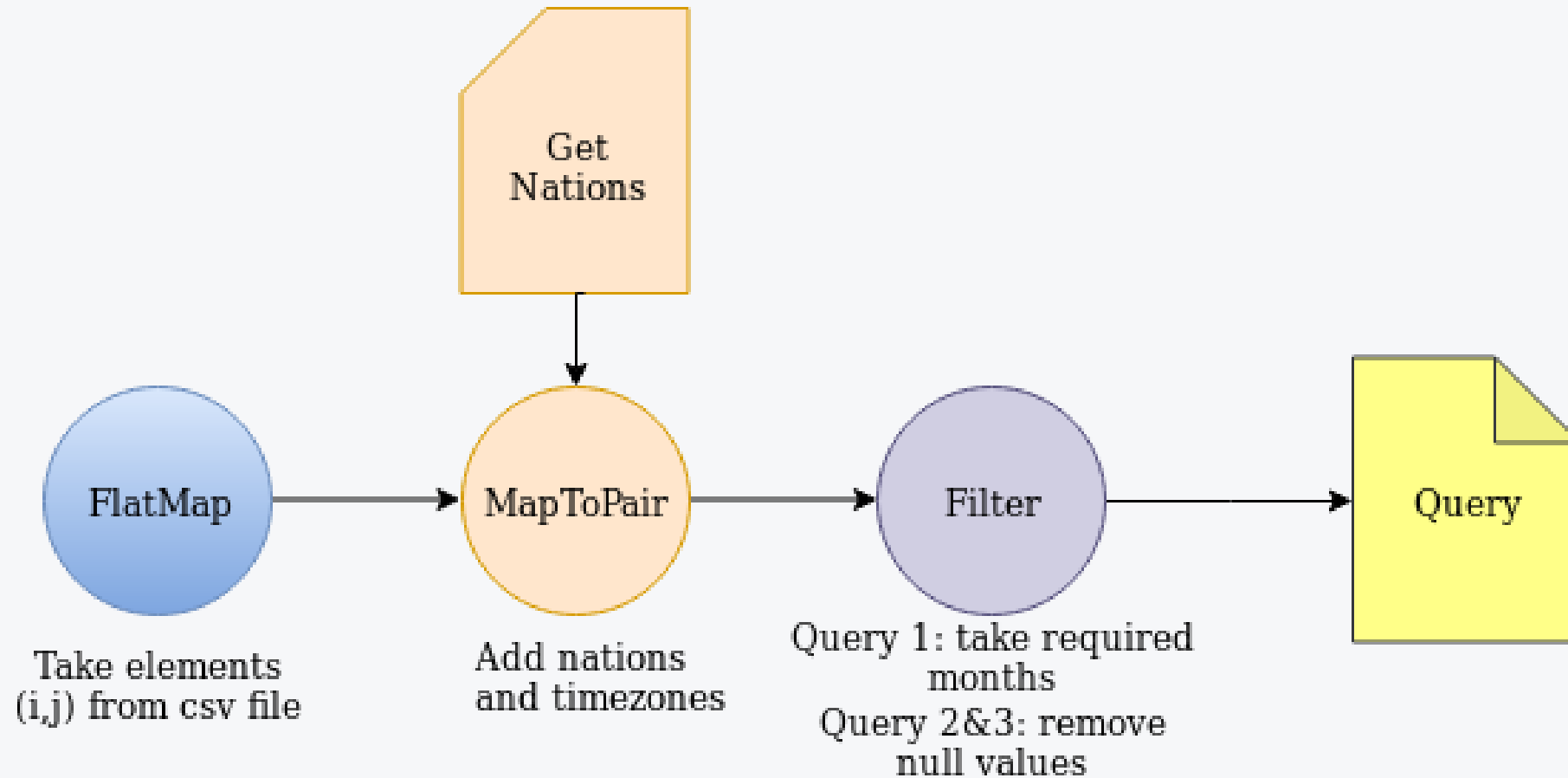


Chiamate REST al servizio Geonames per ottenere le nazioni tramite le coordinate e i fusi orari

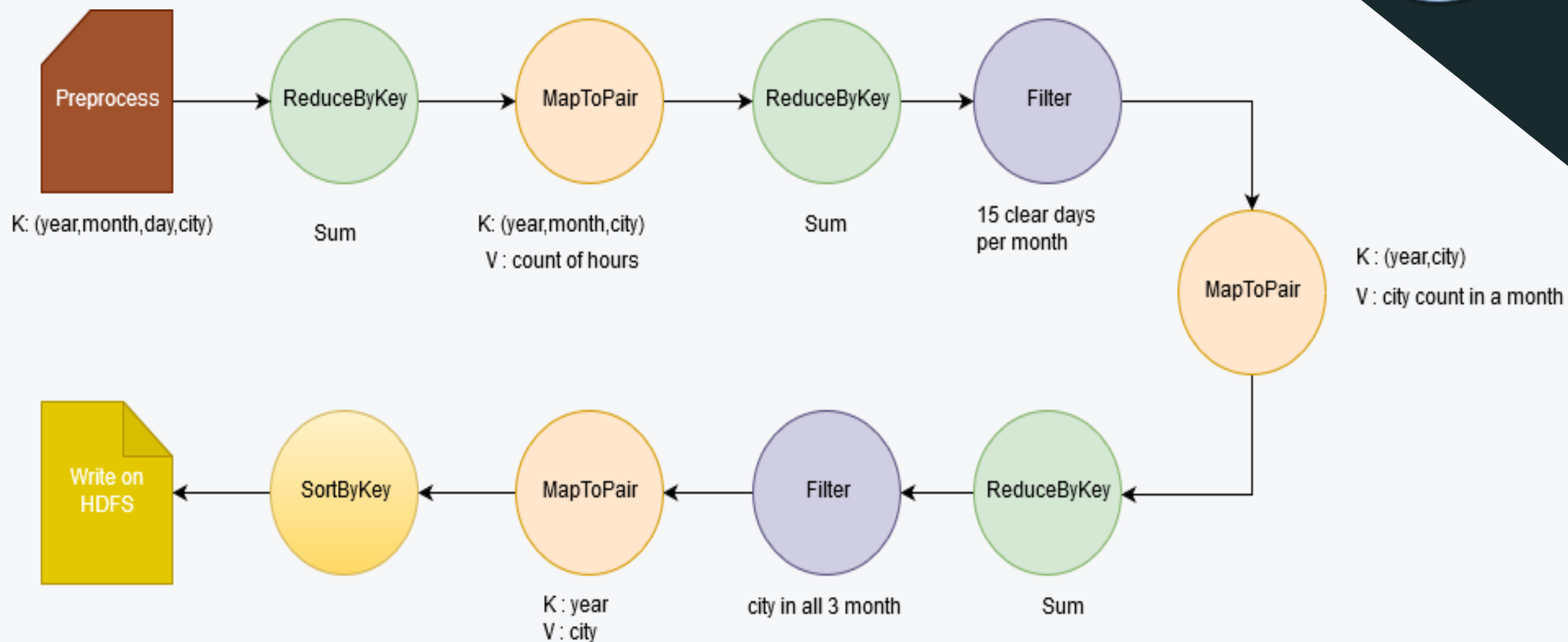


Pre processamento dedicato alla specifica query. Per la query 1 si filtrano i mesi di interesse, per le query 2 e 3 si filtrano i valori non conformi

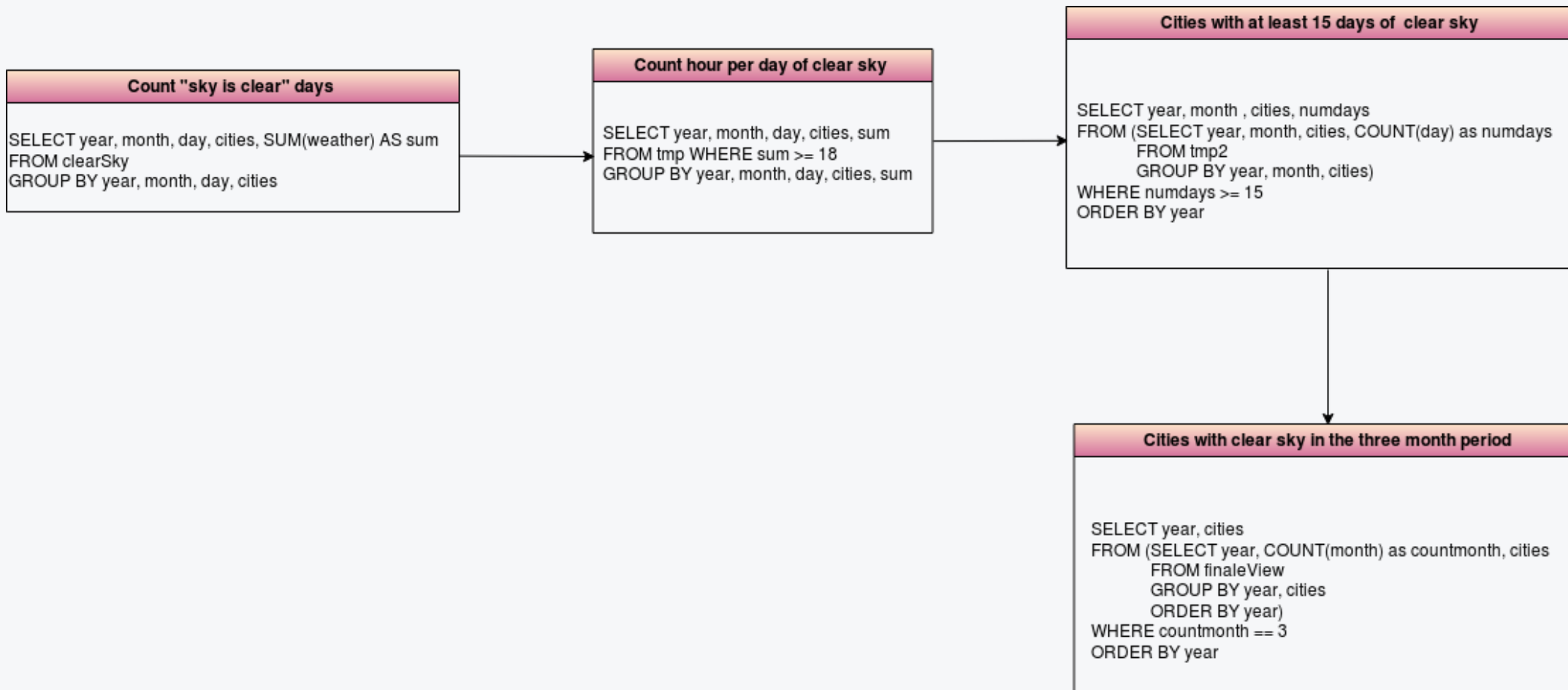
Pre-Processing



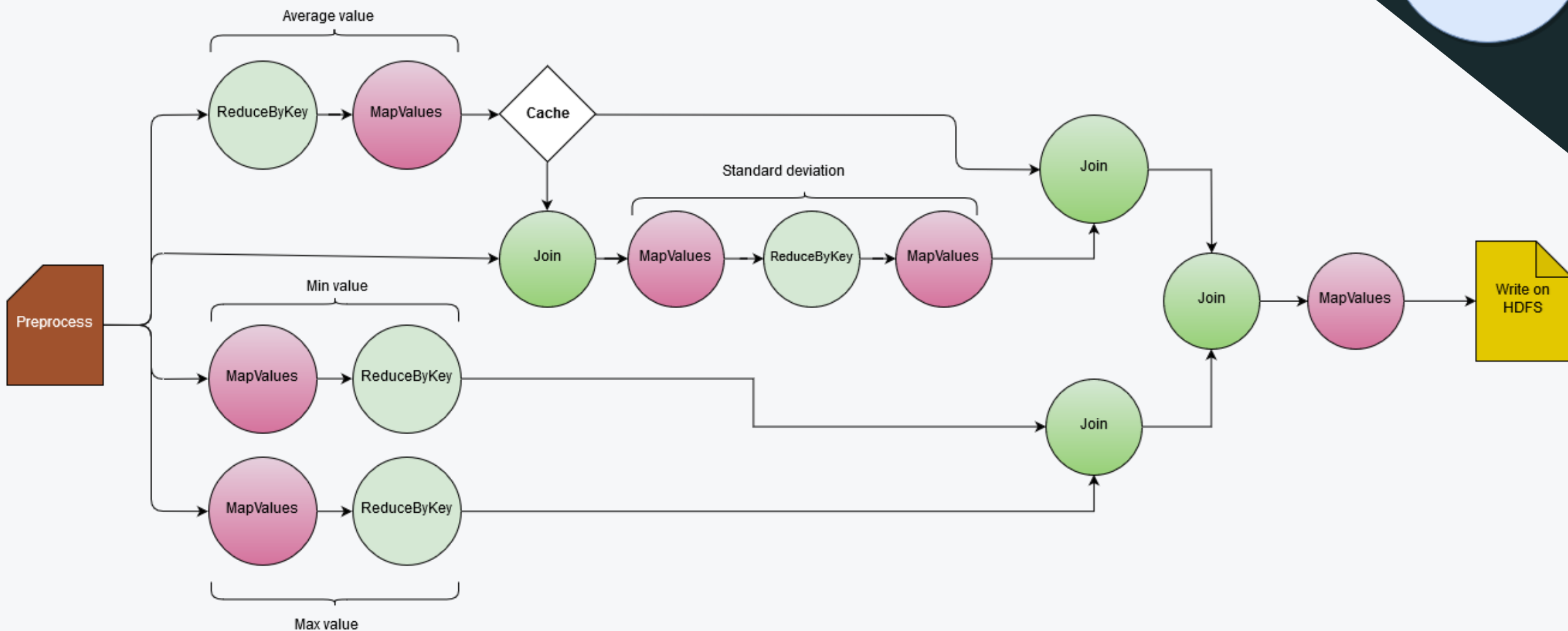
Query 1



Query 1 SQL



Query 2

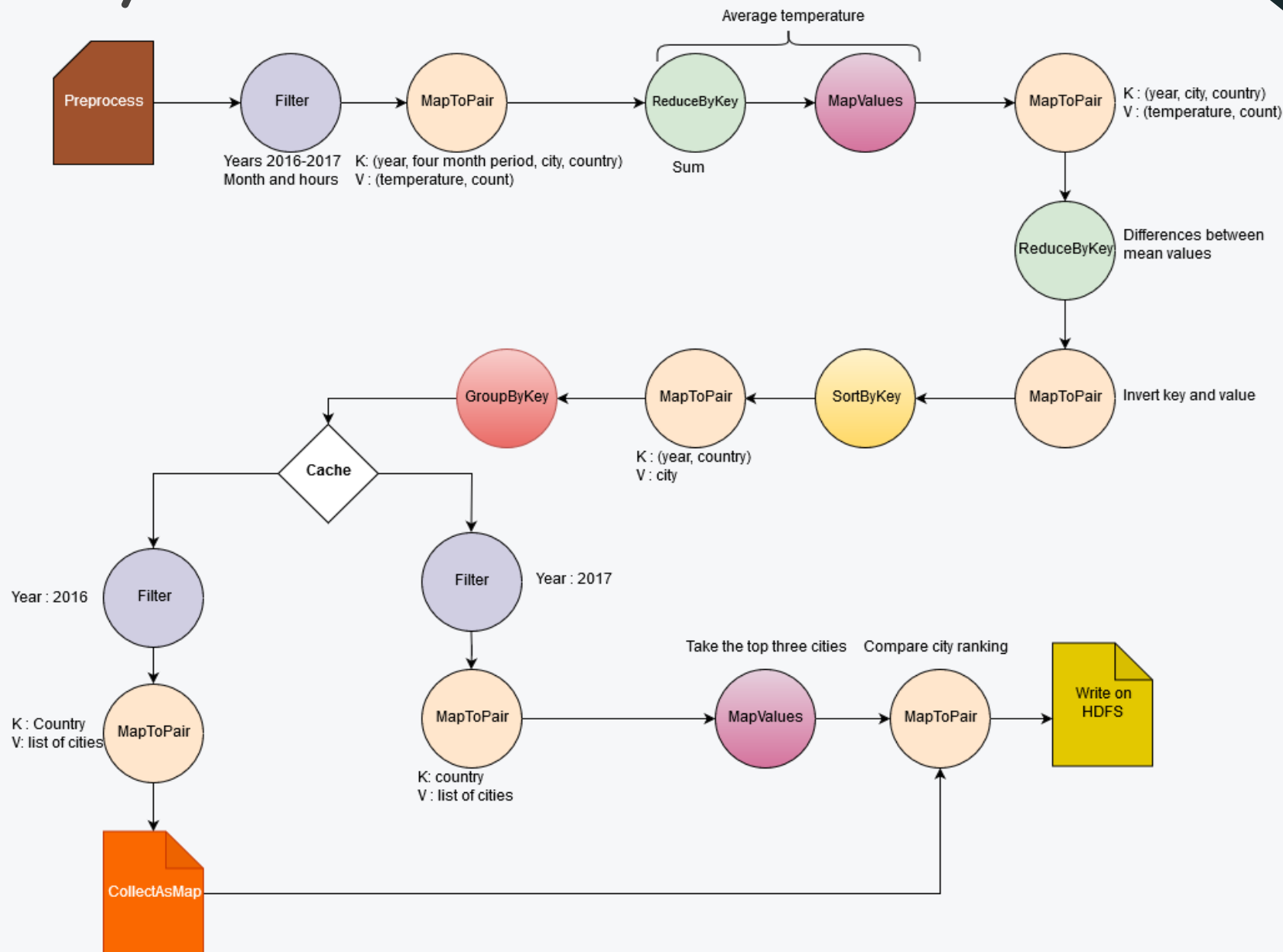


Query 2 SQL

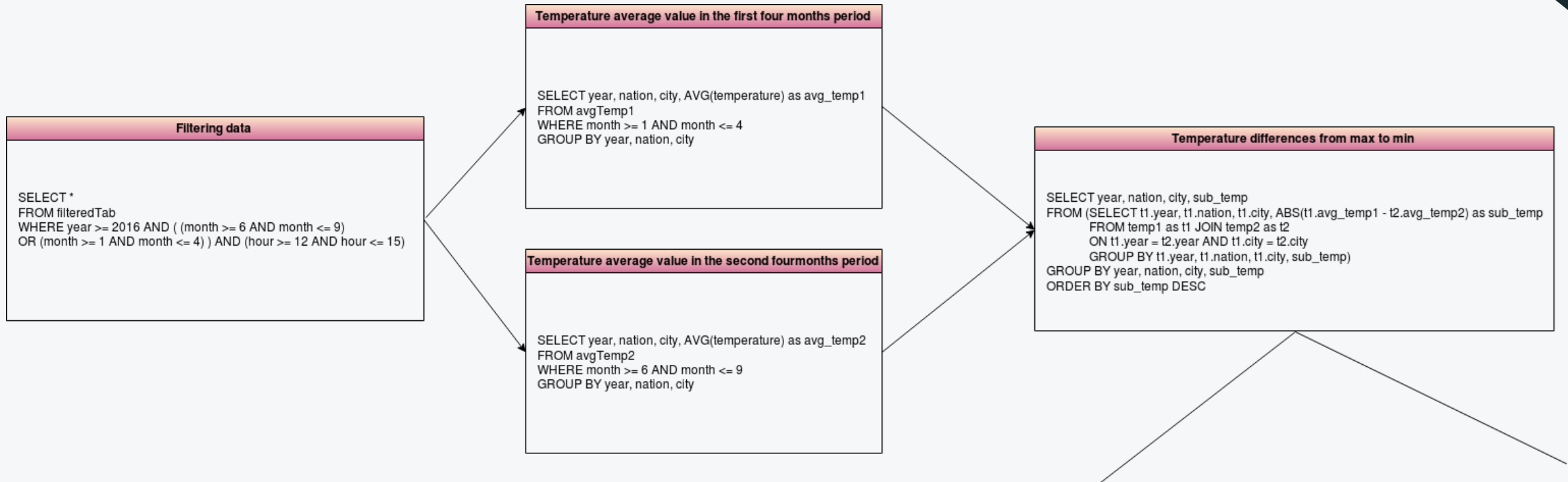
Statistics on humidity, temperature and pressure values

```
SELECT country,year, month, MEAN(value) AS mean ,MIN(value) as min, MAX(value) as max, STDDEV_SAMP(value) as stddev  
FROM statistics  
GROUP BY country,year,month
```

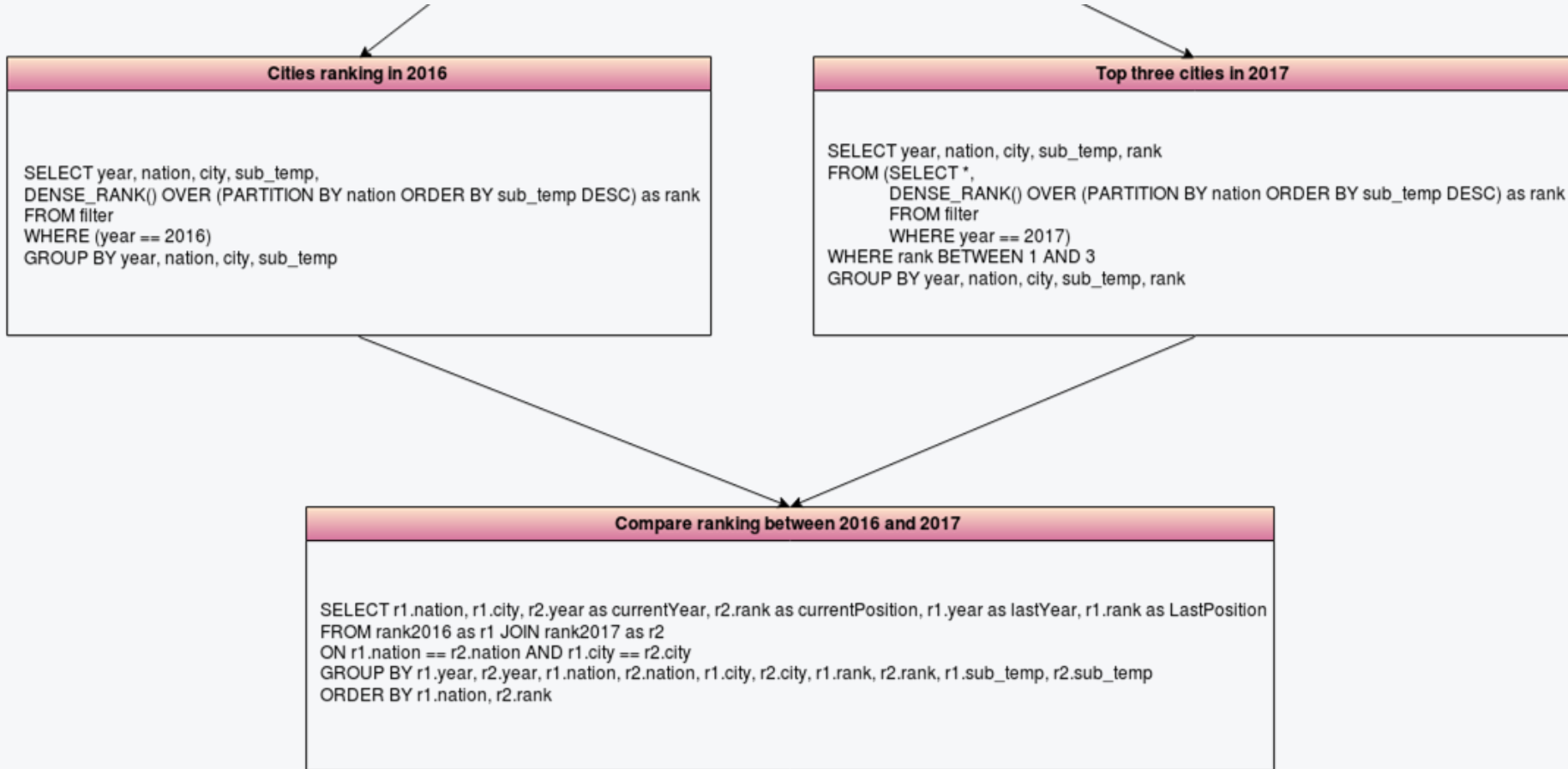
Query 3



Query 3 SQL - (1)



Query 3 SQL - (2)



Risultati query



Query
1

year	cities
2013	Las Vegas
2014	Las Vegas
2015	Phoenix
2015	Las Vegas
2016	Eilat
2016	Phoenix
2016	Las Vegas
2017	Eilat

Query
2

country	year	month	mean	min	max	std_dev
US	2017	4	289.0873198900438	268.85	309.43	6.920265998513651
IL	2013	11	290.0220528935104	280.368	305.67	5.526896959507662
IL	2015	4	290.0203029840284	276.984	307.246	4.467095681613358
US	2015	12	281.9348435066526	256.130423071	305.15	7.7474421142896
IL	2013	8	299.6482675851259	286.209	315.15	5.098120892435437
US	2014	9	293.56699400423815	271.170333333	311.889	5.999285077143429
US	2017	8	297.0385777080903	283.04	315.7	5.651514519969669
US	2014	2	277.7136012790829	247.13	302.91	10.59106048707934

Query
3

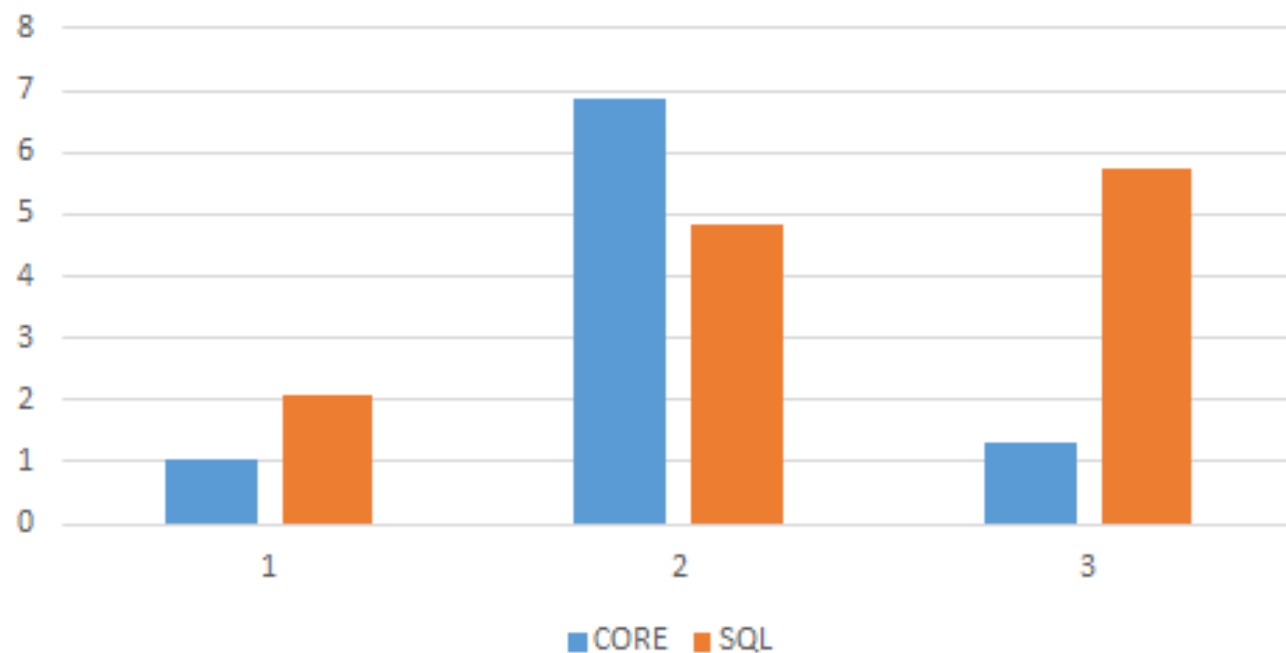
country	city	currentYear	currentPosition	lastYear	LastPosition
IL	Beersheba	2017	1	2016	1
IL	Eilat	2017	2	2016	5
IL	Haifa	2017	3	2016	2
US	Minneapolis	2017	1	2016	2
US	Chicago	2017	2	2016	3
US	Detroit	2017	3	2016	1

Media e varianza - Core vs SQL

Hardware: Intel-Core I5 9600K 6-core, SSD con interfaccia PCIe su slot M.2 e protocollo NVMe.

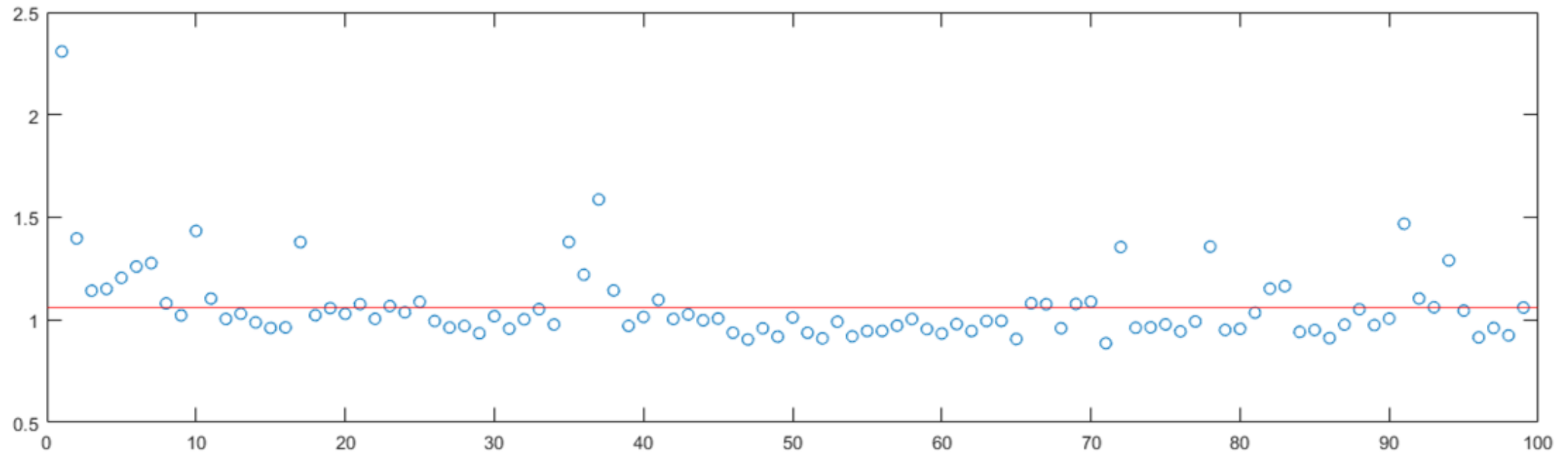


Mean Values Comparison

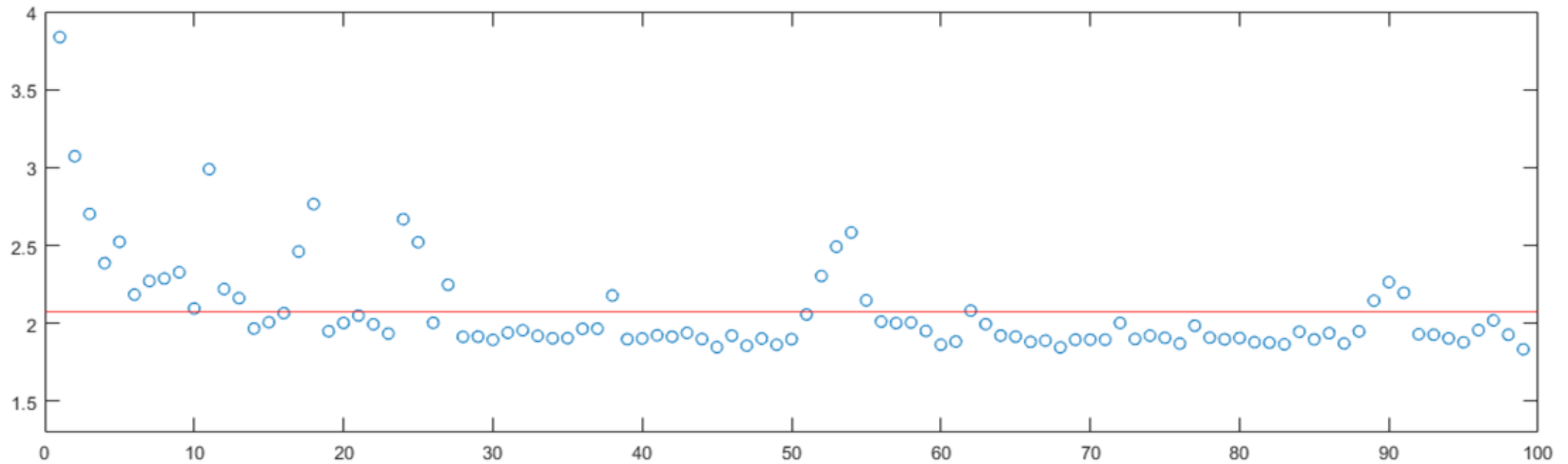


Varianza	1	2	3
Core	0.0353	0.471	0.0132
SQL	0.0951	0.2872	0.3044

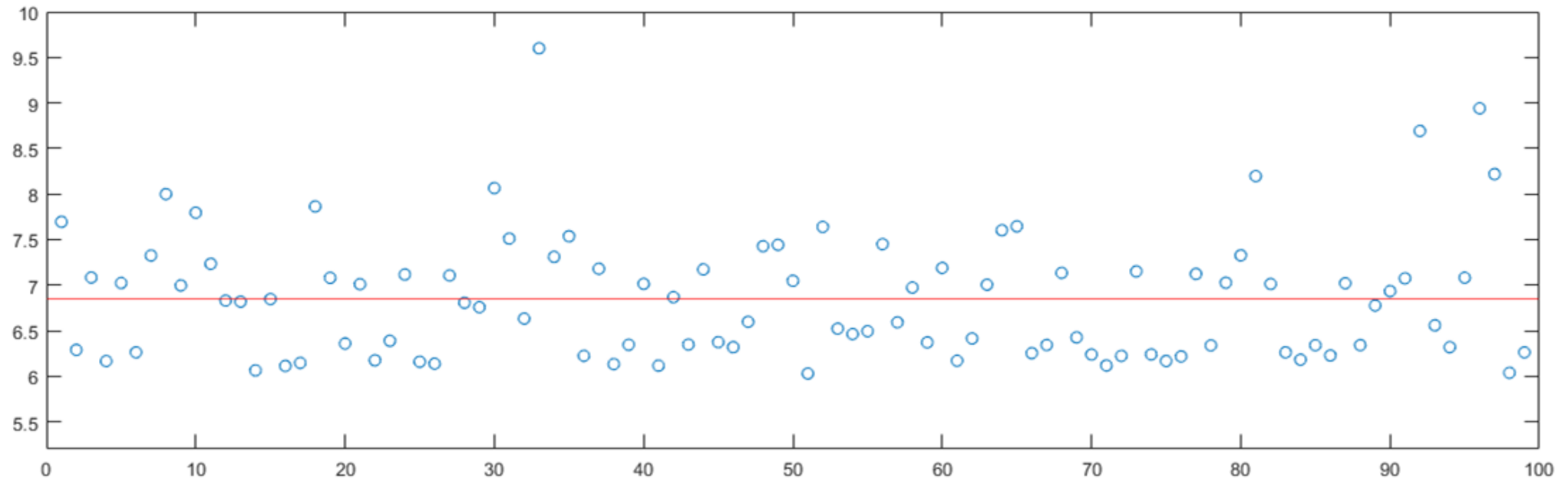
Query1 Core & SQL



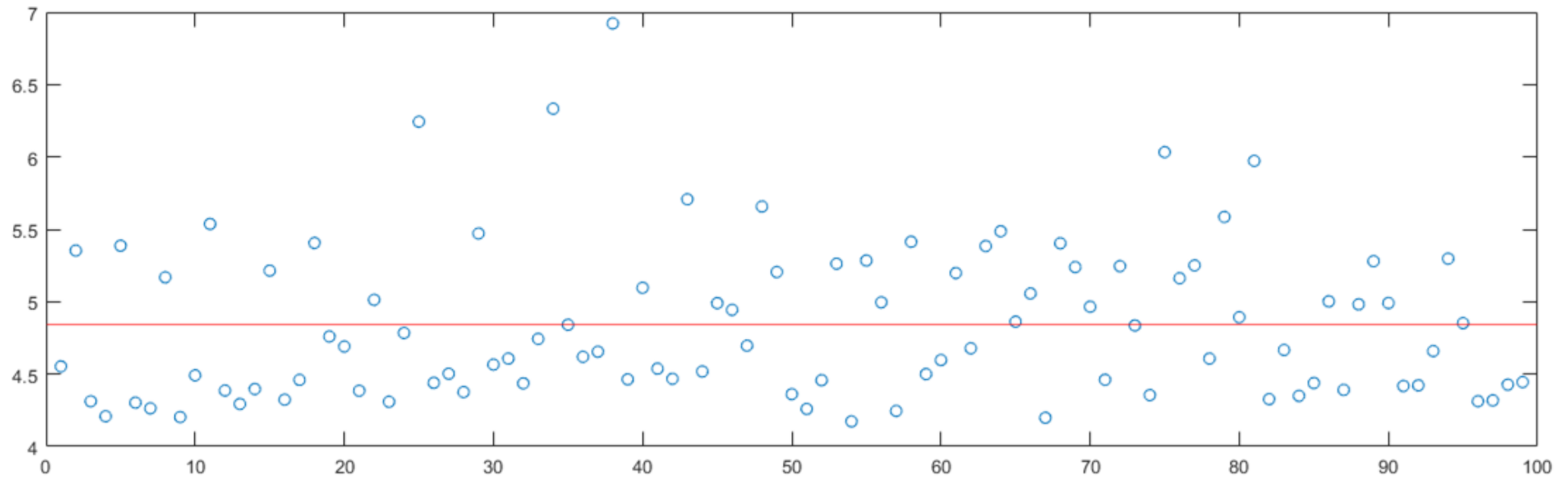
Query1 Core & SQL



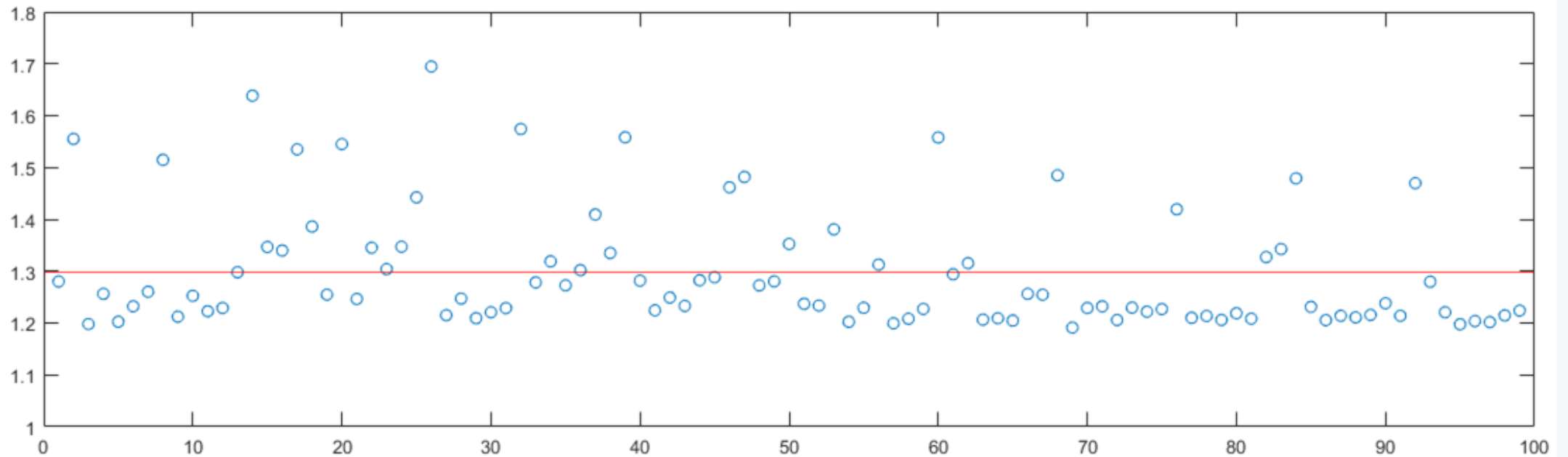
Query2 Core & SQL



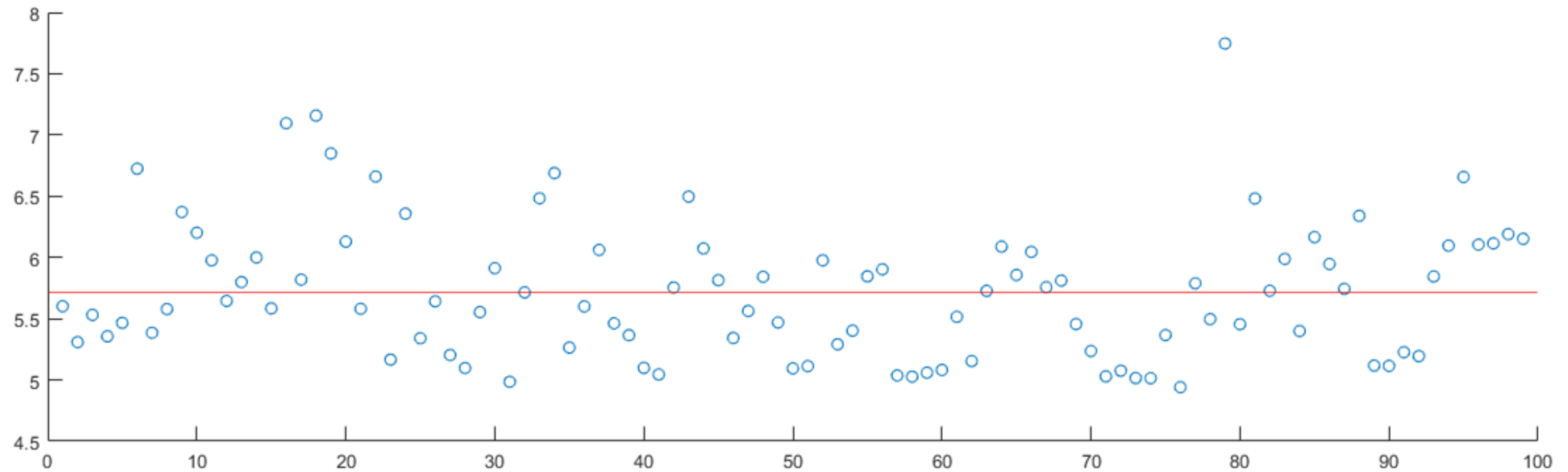
Query2 Core & SQL



Query3 Core & SQL



Query3 Core & SQL



Tempi Ni-Fi



Lettura da locale fino a scrittura su HDFS ~31 s

	05/29/2019 08:20:17.736 UTC	CREATE	de84144d-8515-4e99-b5a4-5...	0 bytes	ListFile	ListFile	
	05/29/2019 08:20:48.886 UTC	DROP	de84144d-8515-4e99-b5a4-5...	19.09 MB	PutParquet	PutParquet	

Scrittura su DB per Query Core

		PutHBaseJSONquery1 PutHBaseJSON 1.9.2 org.apache.nifi - nifi-hbase-nar
In	8 (244 bytes)	5 min
Read/Write	244 bytes / 0 bytes	5 min
Out	0 (0 bytes)	5 min
Tasks/Time	1 / 00:00:02.888	5 min

		PutMongoQuery1 PutMongo 1.9.2 org.apache.nifi - nifi-mongodb-nar
In	8 (244 bytes)	5 min
Read/Write	244 bytes / 0 bytes	5 min
Out	0 (0 bytes)	5 min
Tasks/Time	8 / 00:00:00.749	5 min

		PutHBaseJSONquery2Humidity PutHBaseJSON 1.9.2 org.apache.nifi - nifi-hbase-nar
In	123 (13.57 KB)	5 min
Read/Write	13.57 KB / 0 bytes	5 min
Out	0 (0 bytes)	5 min
Tasks/Time	3 / 00:00:02.310	5 min

		PutMongoQuery2Humidity PutMongo 1.9.2 org.apache.nifi - nifi-mongodb-nar
In	123 (13.57 KB)	5 min
Read/Write	13.57 KB / 0 bytes	5 min
Out	0 (0 bytes)	5 min
Tasks/Time	123 / 00:00:01.780	5 min

		PutHBaseJSONquery3 PutHBaseJSON 1.9.2 org.apache.nifi - nifi-hbase-nar
In	0 (0 bytes)	5 min
Read/Write	0 bytes / 0 bytes	5 min
Out	0 (0 bytes)	5 min
Tasks/Time	1 / 00:00:02.868	5 min

		PutMongoQuery3 PutMongo 1.9.2 org.apache.nifi - nifi-mongodb-nar
In	0 (0 bytes)	5 min
Read/Write	0 bytes / 0 bytes	5 min
Out	0 (0 bytes)	5 min
Tasks/Time	0 / 00:00:00.700	5 min

Considerazioni su Parquet

Utilizzando Parquet le dimensioni dei file cambiano nel modo seguente:

city_attributes.csv	1.01KB	city_attributes.parquet	1.89KB
humidity.csv	7.97MB	humidity.parquet	1.39MB
pressure.csv	10.68MB	pressure.csv	1.14MB
temperature.csv	12.08MB	temperature.parquet	5.88MB
weather_description.csv	19.09MB	weather_description.parquet	806.65KB

I tempi in lettura sono dimezzati, ma siccome erano già dell'ordine dei 300ms per brevità si è deciso di non riportarli.

Grazie per l'attenzione

Progetto Batch processing– SABD 2018-2019
Montesano,Perrone,Pusceddu