

Exercises Big Data Management

Database Technologies and Information Management (DTIM) group
Universitat Politècnica de Catalunya (BarcelonaTech), Barcelona
February 9, 2024

Contents

1	Big Data Design	5
	1.1 Theoretical questions	5
	1.2 Problems	5
2	Distributed Data	7
	2.1 Theoretical questions	7
	2.2 Problems	9
3	Distributed File System	14
	3.1 Theoretical questions	14
	3.2 Problems	14
4	Key-Value Stores	16
	4.1 Theoretical questions	16
	4.2 Problems	16
5	Document Stores	19
	5.1 Problems	19
6	New Relational Architecture	22
	6.1 Theoretical questions	22
	6.2 Problems	22
7	Distributed Processing Frameworks	24
	7.1 MapReduce in use	24
	7.1.1 Theoretical questions	24
	7.1.2 Problems	24
	7.2 MapReduce internals	25
	7.2.1 Problems	25
	7.3 Spark in use	30
	7.3.1 Problems	30
	7.4 Spark internals	34
	7.4.1 Theoretical questions	34
	7.4.2 Problems	34
8	Streams	36
	8.1 Problems	36
9	Architectures	39
	9.1 Theoretical questions	39
	9.2 Problems	39

1 Big Data Design

1.1 Theoretical questions

1. Briefly explain what “physical independence” is, which is the general position of NOSQL systems on this, and why they take this position.
2. In the framework of the RUM conjecture, name six data structures we saw in the course (not concrete tool implementations) and place them in the corresponding category:
 - Read optimized:
 - Write optimized:
 - Space optimized:
3. Compare a B-tree and an LSM-tree in the context of the **RUM conjecture** (i.e., as an answer to this question, three brief explanations of the form “From the perspective of X, Y-tree is better than Z-tree, because of this and that.” are expected).
 - (a) R
 - (b) U
 - (c) M

1.2 Problems

1. Consider:
 - a) ['BCN', POPULATION:{VALUE:'2,000,000'}, REGION:{VALUE:'CAT'}]
 - b) ['BCN', ALL:{VALUE:'2,000,000;CAT'}]
 - c) ['BCN', ALL:{POPULATION:'2,000,000';REGION:'CAT'}]

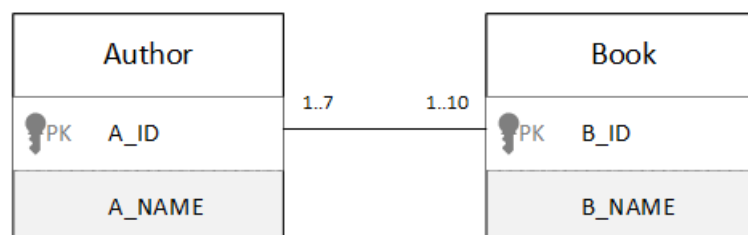
Place {a, b, c} in the table:

	Less variable schema	→	More variable schema
Less explicit schema			
↓			
More explicit schema			

Name two criteria you would use to choose among them:

- 1)
- 2)

2. Consider the following conceptual schema and propose several design alternatives to translate it to a logical representation (e.g., using JSON notation). Then, explain which is the best alternative and why.



2 Distributed Data

2.1 Theoretical questions

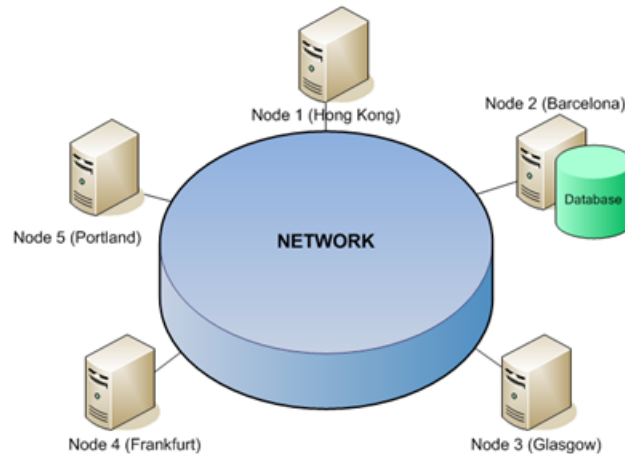
1. Give five types of data or software resource that can usefully be shared. Give examples of their sharing as it occurs in practice in distributed systems.
2. Consider the implementation strategies for massively multiplayer online. In particular, what advantages do you see in adopting a single server approach for representing the state of the multiplayer game? What problems can you identify and how might they be resolved?¹
3. The INFO service manages a potentially very large set of resources, each of which can be accessed by users throughout the Internet by means of a key (a string name). Discuss an approach to the design of the names of the resources that achieves the minimum loss of performance as the number of resources in the service increases. Suggest how the INFO service can be implemented so as to avoid performance bottlenecks when the number of users becomes very large.¹

4. A server program written in one language (for example, C++) provides the implementation of a BLOB object that is intended to be accessed by clients that may be written in a different language (for example, Java). The client and server computers may have different hardware, but all of them are attached to the Internet. Describe the problems due to heterogeneity of (a) programming languages, and (b) implementations by different developers that need to be solved to make it possible for a client object to invoke a method on the server object.¹

Give arguments for and against allowing the client requests to be executed concurrently by the server. In the case that they are executed concurrently, give an example of possible “interference” that can occur between the operations of different clients. Suggest how such interference may be prevented.

5. A service is implemented by several servers. Explain why resources might be transferred between them. Would it be satisfactory for clients to multicast all requests to the group of servers as a way of achieving mobility transparency for clients?¹
6. Explain what is (a) a distributed system, and (b) a parallel system. Compare both of them (i.e., what has one and not the other and vice-versa).
7. Which kind of database is this according to the distribution of data?

¹From G. Couloris et al. *Distributed Systems: Concepts and Design*, 5th Ed. Addison-Wesley, 2012



8. Which two kinds of schema information contains the Global Conceptual Schema that does not contain the Local Conceptual Schema in the Extended ANSI-SPARC Architecture for DDBMS?
 - 1.
 - 2.
9. In the context of distributed data management, name the **four big challenges** that need to be carefully considered in the presence of distribution from the tenant/user point of view.
 - I.
 - II.
 - III.
 - IV.
10. Name the three characteristics of fragmentation that make it correct, and for each of them give an example of fragmentation schema where it is **not fulfilled**.
 - (a)
 - (b)
 - (c)
11. Which is the main problem in having replicas, and which is the innovation introduced by some NOSQL tools to solve it.
 - Problem:
 - Innovation:
12. Given N replicas, let's call R the *ReadConcern* parameter of MongoDB and W the *WriteConcern* (which indicate respectively the number of copies it reads and writes, before confirming the operation to the user); give the equation involving those variables that corresponds to the **eventually consistent** configuration.
13. How might the clocks in two computers that are linked by a local network be synchronized without reference to an external time source? What factors limit the accuracy of the procedure you have described? How could the clocks in a large number of computers connected by the Internet be synchronized? Discuss the accuracy of that procedure.¹

14. What is the difference between query cost and query response time ...
 - a) in centralized systems (without parallelism)?
 - b) in distributed systems (inherently parallel)?
15. Name the two factors that make impossible having linear scalability according to the Universal Scalability Law.
 - (a)
 - (b)

2.2 Problems

1. Briefly explain (a) which fragmentation strategy has been applied for the tables below and whether this fragmentation strategy is (b) complete, (c) disjoint and (d) allows to reconstruct the global relations (if so, (e) indicate the operation).

Global Relations

```
Kids(kidId, name, address, age)
Toys(toyId, name, price)
Requests(kidId, toyId, willingness)
```

Note that `requests(kidId)` is a foreign key to `kids(kidId)` and similarly, `requests(toyId)` refers to `toys(toyId)`.

Fragments

```
K1= Kids[kidId, name]
K2= Kids[kidId, address, age]
```

```
T1= Toys(price ≥ 150)
T2= Toys(price < 150)
```

```
R1 = Requests ⋈ T1
R2 = Requests ⋈ T2
```

2. You are a customer using an e-commerce based on heavy replication (e.g., Amazon):
 - a) Show a database replication strategy (e.g., sketch it) where:
 - i. You buy an item, but this item does not appear in your shopping cart.
 - ii. You reload the page: the item appears.
 ⇒ What happened?
 - b) Show a database replication strategy (e.g., sketch it) where:
 - i. You delete an item from your command, and add another one: the shopping cart shows both items.
 ⇒ What happened? Will the situation change if you reload the page?

3. Consider the following architectures and answer the questions:²

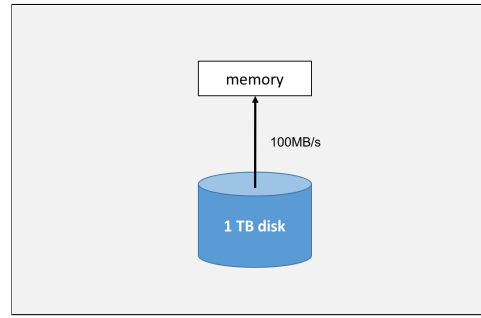


Figure 1: Centralized architecture

Type	Latency	Bandwidth
Disk	$\approx 5 \times 10^{-3}$ s (5 millisc.)	At best 100 MB/s

- How long would it take (i.e., response time) to read 1TB with sequential access (Figure 1)? (in secs)
- How long would a single random access (i.e., reading one tuple, of for example 100 bytes, through an index) take (i.e., response time), assuming we already have the physical address? (in secs)

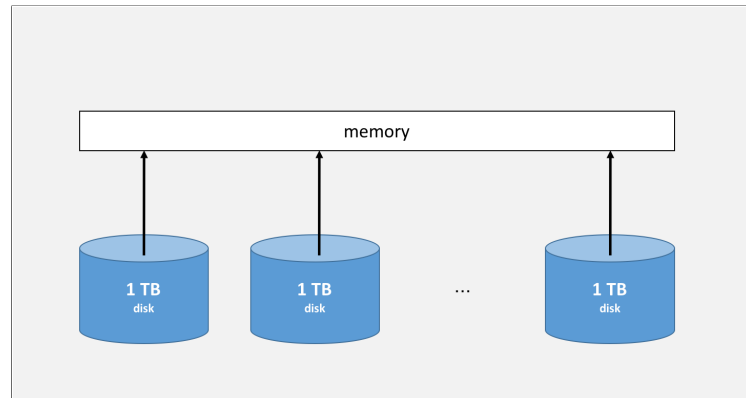


Figure 2: Shared-memory architecture

Type	Latency	Bandwidth
Disk	$\approx 5 \times 10^{-3}$ s (5 millisc.)	At best 100 MB/s

- How long would it take (i.e., response time) to read 1TB with parallel access (Figure 2)? Assume 100 disks (i.e., 100 replicas of the whole data) on the same machine with shared-memory and infinite CPU capacity.
- How long would a single random access (i.e., reading one tuple, of 100 bytes, through an index) take (i.e., response time), assuming we already have the physical address? (in secs)

²From S. Abiteboul et al. *Web Data Management*. Cambridge Press, 2011.

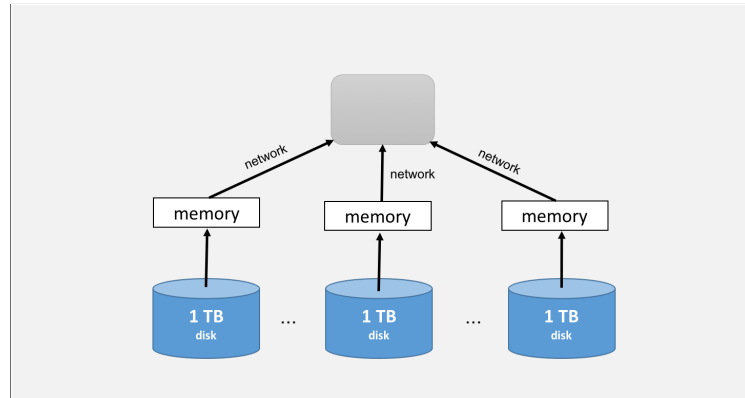


Figure 3: Shared-nothing architecture

Type	Latency	Bandwidth
Disk	$\approx 5 \times 10^{-3}$ s (5 millisc.)	At best 100 MB/s
LAN	$\approx 1 - 2 \times 10^{-3}$ s (1-2 millisc.)	≈ 1 GB/s (single rack) ≈ 10 MB/s (switched)
Internet	Highly variable (typically 10-100ms)	Highly variable (typically a few MB/s, e.g., 10MB/s)

Note 1: It is approx. one order of magnitude faster to exchange main memory data between 2 machines in a data center, that to read on the disk.

Note 2: Exchanging through the Internet is slow and unreliable with respect to LANs.

- e) How long would it take (i.e., response time) to read 1TB with distributed access (Figure 3)? Assume 100 shared-nothing machines (with all data replicated in each of them) in a star-shape LAN in a single rack where all data is sent to the center of the star in only one network hop.
- f) How long would it take (i.e., response time) to read 1TB with distributed access (Figure 3)? Assume 100 shared-nothing machines (with all data replicated in each of them) in a star-shape cluster of machines connected through the Internet where all data is sent to the center of the star in only one network hop.
- g) How long would a single random access (i.e., reading one tuple, of 100 bytes, through an index) take (i.e., response time) in the case of a LAN, assuming we already have the physical address? (in secs)
- h) How long would a single random access (i.e., reading one tuple, of 100 bytes, through an index) take (i.e., response time) in the case of an Internet connection, assuming we already have the physical address? (in secs)
4. What are the main differences between these two distributed access plans? Under which assumptions is one or the other better?

```
SELECT *
FROM employee e, assignedTo a
WHERE e.#emp = a.#emp AND a.responsability = 'manager'
```

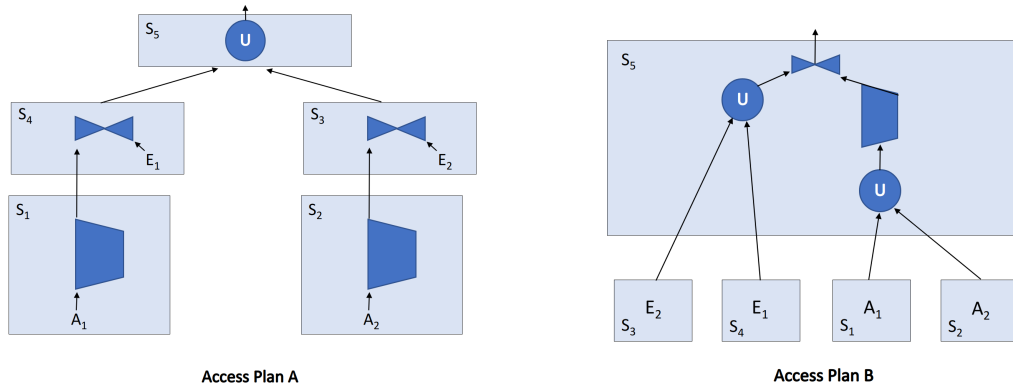


Figure 4: Distributed Access Plans

```

Employee(#emp, empName, degree)
S4: E1 = Employee(#emp ≤ 'E3')
S3: E2 = Employee(#emp > 'E3')

AssignedTo(#emp, #proj, responsibility, fullTime)
FK: #emp references Employee
S1: A1 = AssignedTo(#emp ≤ 'E3')
S2: A2 = AssignedTo(#emp > 'E3')

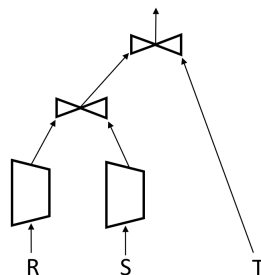
```

5. Compute the fragment query (data location stage) for the database setting and query below and find in how many ways we can assign the operations to the different sites.

⇒ **Database setting:**

- A distributed database with 5 sites (i.e., database nodes): S_1, S_2, S_3, S_4 and S_5 .
- 3 relations in the database R, S and T
- Each relation is horizontally fragmented in two fragments (we refer to them by the name of the relation and a subindex, for example: R_1, R_2). You can consider them to be correct (i.e., complete, disjoint and reconstructible).
- Each fragment is replicated at all sites.
- We have the following query: $Q_1 = \sigma(R) \bowtie \sigma(S) \bowtie T$

⇒ **Process tree of the query:**



6. Consider a left-deep process tree corresponding to a query, where each internal node is a join, and every leaf a data source (e.g., relational table). Knowing that the tree contains 9 nodes (including leaves), the system has as much parallelism capacity as needed to run all the joins in pipelining mode (no other kind of parallelism is available), which is the occupancy if the overall cost of the serial query is 4 seconds? Explicit any assumption you need to make.

3 Distributed File System

3.1 Theoretical questions

1. How does HDFS decide where to place the different chunks of a file?

3.2 Problems

1. Given a table with ten blocks ($B = 10$) and 17 tuples per block ($R = 17$) for a total of 170 tuples, compare the cost of reading the whole table sequentially, against the cost of accessing seven random (potentially repeated) tuples in an unknown order. Assume that seek time is $12ms$, average rotation time is $3ms$ and transferring one block is $0.03ms$. You should consider which is the probability of two tuples accessed consecutively are actually in the same block, but ignore the presence of any cache or buffer pool mechanism.
2. Consider a table stored in HDFS with the following characteristics:

size(T)	= 256 MB
T	= 64
size(row)	= 4 MB
cols(T)	= 4
size(cell)	= 1 MB
size(Header)	= 0
size(Footer)	= 0

- a) Given the configuration parameters, how much space would you need to store it in a horizontal layout (i.e., Avro)?

Size(MetaRow)	= 0.1MB
Size(MetaBody)	= 0.1MB

- b) Given the configuration parameters, how much space would you need to store it in a hybrid layout (i.e., Parquet)?

Size(MetaCol)	= 0.2MB
Size(MetaRowGroup)	= 0.5MB
Size(RowGroup)	= 8MB

- c) Given the configuration parameters, how much data would you retrieve in a hybrid layout (i.e., Parquet) to project two columns?

Size(MetaCol)	= 0.2MB
Size(MetaRowGroup)	= 0.5MB
Size(RowGroup)	= 8MB

- d) Given the configuration parameters, how much data would you need to retrieve it in a hybrid layout (i.e., Parquet) to select one row (e.g., given its key) if the table is not sorted?

Size(MetaCol)	= 0.2MB
Size(MetaRowGroup)	= 0.5MB
Size(RowGroup)	= 8MB

3. Given a file with C chunks and N nodes in the cluster, find the probability that HDFS uses all the nodes in case $C \geq N$, by using Stirling numbers of the second kind.³
4. Consider an HDFS cluster with 100 data nodes, without replication. If I upload a file with 10 chunks and 10 disk blocks each, answer the following questions and briefly justify your answer:
 - a) Which is the maximum number of machines containing data?
 - b) Which is the probability of the maximum number of machines actually contain data?
5. Given a file with 3.2GB of raw data stored in an HDFS cluster of 50 machines, and containing $16 \cdot 10^5$ rows in a **Parquet file**; consider you have a query over an attribute " $A = \text{constant}$ " and this attribute contains only 100 different and equiprobable values. Assuming any kind of compression has been disabled, explicit any assumption you need to make and give the amount of raw data (i.e., do not count metadata) it would need to fetch from disk.
 - Replication factor: 3 (default)
 - Chunk size: 128MB (default)
 - Rowgroup size: 32MB
6. Consider a cluster of ten worker machines and a single coordinator, which contains a sequence file of 128MB stored in HDFS with an OS block size of 32KB in all machines and a chunk size of 64MB. On the event of a client retrieving that file, give the **number of control messages** (do not consider data messages) that will travel the network in the case of (a) a client cache miss and in the case of (b) a client cache hit. Briefly justify both numbers.

³https://en.wikipedia.org/wiki/Stirling_numbers_of_the_second_kind

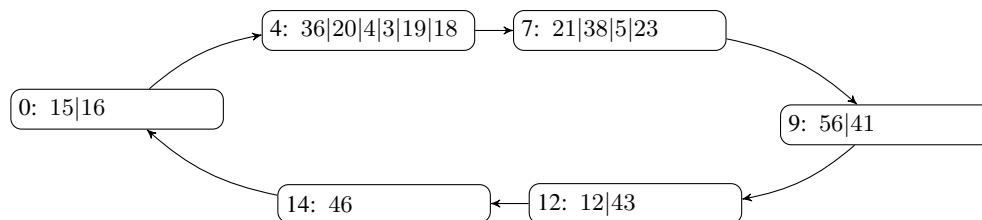
4 Key-Value Stores

4.1 Theoretical questions

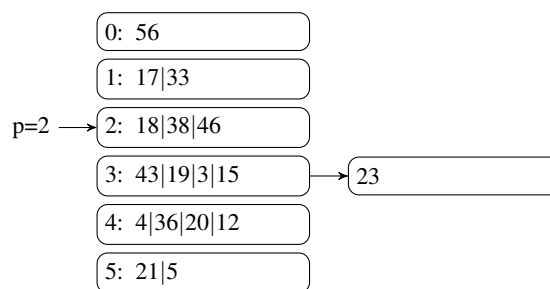
1. Which is the main difference between the hash functions used in the linear hash and consistent hash algorithms?
2. With respect to distributed systems, explain what is a distributed hash table (DHT), and provide a brief description of how consistent hashing guarantees balancing keys when adding new servers.

4.2 Problems

1. Let's assume we have a *Consistent hash* with $D = 16$, and the hash function is simply the module of the IP address or the key, and suppose the current state of the consistent hash is (position_in_the_ring:key|key|...):



- a) What happens when we insert objects 30 and 58? Draw the result.
 - b) What happens in the structure when we register a new server with IP address 37? Draw the result.
2. Let's suppose we have a *Linear Hash* and the hash function is simply the module of the key, the capacity of a bucket is only four entries, and current state of the linear hash is (bucketID: key|key|...):



- a) What happens in the structure when we insert keys 14, 27, 37, and 44? Draw the result.
3. Let's suppose that, we have an *LSM Tree* that reached the threshold to consider the MemStore is full, and it contains four entries with format $[key, value, timestamp]$ needing ten characters each. The content of the different structures is:

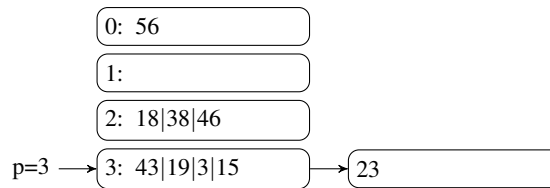
- MemStore: $[1, v, t50], [15, v, t49], [17, v, t47], [29, v, t48]$
- Commit Log: $[17, v, t47], [29, v, t48], [15, v, t49], [1, v, t50]$

- $SSTable_{Data}$: $[13, v, t23], [25, v, t17], [35, v, t40], [59, v, t38]$
- $SSTable_{Index}$: $[13, 0], [25, 30], [35, 60], [59, 90]$

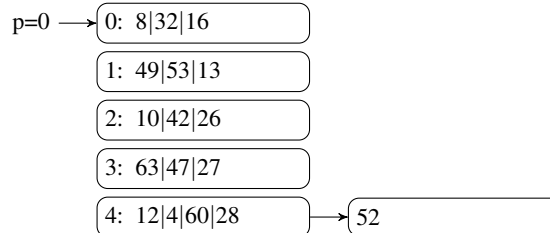
Assuming that the minimum size of an $SSTable$ is 120 characters and on having two $SSTables$ a minor compaction is automatically triggered, explicit the content of all structures once the compactation is done.

- MemStore:
- Commit Log:
- $SSTable_{Data}$:
- $SSTable_{Index}$:

4. Briefly explain what is wrong in this linear hash structure, or if you think it is right, explicitly say so.



5. Given the **linear hash** underneath with $f(x) = x$ (i.e., we directly apply the module to the keys), and a capacity of four keys per bucket, indicate if it corresponds to a state **valid or not**. If valid, give a possible insertion order leading to it. If not, clearly explain why.



6. Suppose you have a hash function whose range has size 100 (i.e., $D=100$), and a Consistent Hash structure with 5 machines ($M1...5$) whose identifiers map to values $h(M1) = 0$, $h(M2) = 20$, $h(M3) = 40$, $h(M4) = 60$, $h(M5) = 80$. What happens if you have an object mapped to value $h(O) = 90$?

7. Given an empty **Consistent Hash** with $h(x) = x\%32$ (i.e., we directly take module 32 to both the keys and the bucket IDs), and unlimited capacity in each bucket, consider you have a cluster of four machines with IDs 19, 22, 75, 92, and draw the result of inserting the following keys in the given order: 12, 4, 10, 49, 42, 60, 63, 53, 47, 27, 26, 28, 13, 52.

8. Suppose you implement a system to store images in hundreds of machines with thousands of users using HBase with a single column-family. These images taken at time VT belong to a person P who tags each with a single subject S (e.g., family, friends, etc.) and are concurrently uploaded into the system at time TT in personal batches containing multiple pictures of different subjects taken at different times. Each person can then retrieve all his/her pictures of one single subject that were taken after a given time. Precisely define the key you would use (which cannot be a hash) if you exclusively prioritize (i.e., do not consider any other criteria)...

- (a) Load balancing on ingestion:
 - ⇒ Assumptions made:
 - (b) Load balancing on querying:
 - ⇒ Assumptions made:
 - (c) I/O cost (i.e., minimum blocks flushed) on ingestion:
 - ⇒ Assumptions made:
 - (d) I/O cost (i.e., minimum blocks retrieved) on querying:
 - ⇒ Assumptions made:
9. Given an empty **linear hash** with $f(x) = x$ (i.e., we directly apply the module to the keys), and a capacity of four keys per bucket, draw the result of inserting the following keys in the given order: 12, 4, 10, 49, 42, 60, 63, 53, 47, 27, 26, 28, 13, 52.

5 Document Stores

5.1 Problems

1. Optimize the performance of the following document design:

```
{ "_id": 123,  
  "name": "Alberto"  
  "address": {  
    "street": "Jordi Girona",  
    "city": "Barcelona",  
    "zip_code": "08034"},  
  "telephone1": "93 4137889",  
  "telephone2": "93 4137840",  
  "telephone3": "123456789"  
}
```

2. Analyze (i.e., briefly give **pros and cons**) the following **JSON design** compared to other equivalent JSON designs from the three perspectives:

```
"BID-PRODUCT": {  
  "B_ID": int(4), "B_PRICE": int(4), "U_ID": int(4),  
  "PRODUCT": {  
    "P_ID": int(4), "P_INFO": varchar(100)  
  }  
}  
"PRODUCT-SELLER-REGION": {  
  "P_ID": int(4), "P_INFO": varchar(100),  
  "USER": {  
    "U_ID": int(4), "U_F_NAME": varchar(20),  
    "REGION": {  
      "R_ID": int(4), "R_NAME": varchar(10)  
    }  
  }  
}  
"PRODUCT-COMMENTS": {  
  "P_ID": int(4), "P_INFO": varchar(100),  
  "COMMENTS": [{  
    "C_ID": int(4), "C_TITLE": varchar(20), "U_ID": int(4)  
  }]  
}
```

(a) Read (a.k.a. Query)

- i. Pros:
- ii. Cons:

(b) Update

- i. Pros:
- ii. Cons:

(c) Memory (a.k.a. Space)

- i. Pros:
- ii. Cons:

3. Optimize the performance of the following pipe:

```

db.partsupp.aggregate([
  {"$sort": {"n_name": 1, "s_name": 1, "p_partkey": 1}},
  {"$match": {"part.p_size": 5}},
  {"$project": {
    "supplier.s_name": 1,
    "supplier.s_phone": 1,
    "supplier.s_comment": 1,
    "_id": 0
  }}
])

```

4. Assume a MongoDB collection of JSON documents following the structure shown below. Indicate a sequence of stage operators that can be used in a pipeline to compute the top five departments in terms of number of employees such that its location has at least one `performance_review` greater than 7. In the example below, since both locations have performance reviews greater than 7, manufacturing has a total of 90 employees, while billing and sales have 80 and 40 employees, respectively. You need only provide the names of the stage operators and the sequence in which they are applied. You should not use the `mapreduce` feature of MongoDB but rather the specific stage operators provided by MongoDB. You do not need to explicitly handle the case of ties (i.e., in the case of ties, you may return any of the tied documents to complete the top five).

```

{ "location": "Barcelona",
  "budget": 3.000.000,
  "departments": [ { "name": "sales", "employees": 40 }, { "name": "manufacturing",
    "employees": 75 } ],
  "performance_reviews": [ 7, 6, 3 ]
},
{ "location": "Brussels",
  "budget": 5.000.000,
  "departments": [ { "name": "manufacturing", "employees": 15 }, { "name": "billing",
    "employees": 80 } ],
  "performance_reviews": [ 6, 8, 7 ]
}

```

5. Imagine you want to design a system to maintain data about citizens and doubt whether to use HBase or MongoDB. Precisely, for each citizen, it will store: their personal data (with *pID*), their city data (with *cID*) and their employment data. We also know the workload (i.e., queries and frequency of execution) is as follows:

- Q_1 : Average salary in Barcelona (50% frequency) – information obtained from the set of city and employment data.
- Q_2 : Average weight for young people (less than 18 years old) (45% frequency) – information obtained from the set of personal data.
- Q_3 : Number of VIPs (5% frequency) – information obtained from the set of personal data.

⇒ Discuss your choice of technology and data model, without pre-computing the results of the queries. Clearly specify the structure of the data (i.e., tables/collections, keys, values, etc.), trade-offs and assumptions made.

6. Assume you have a MongoDB collection which occupies 6 chunks **evenly distributed** in 3 shards (i.e., 2 chunks per shard). Being the document *ID* also the shard key, the chunk of a document is determined **by**

means of a hash function. Assuming that accessing one document takes one time unit (existing indexes are used at no cost) and we have 6,000 documents in the collection, k of which have value “YYY” for attribute “other”, how many time units would take the following operations⁴:

- a) *FindOne*($\{_id : "XXX"\}$)
- b) *Find*($\{_id : \{\$in : [1, \dots, 3000]\}\}$), being $[1, \dots, 6000]$ the range of existing IDs.
- c) *Find*($\{other : "YYY"\}$), being the attribute indexed.
- d) *Find*($\{other : "YYY"\}$), being the attribute NOT indexed.

7. Assume you have a MongoDB collection which occupies 6 chunks **UNevenly distributed** in 3 shards (i.e., 1, 2 and 3 chunks per shard respectively). Being the document Id also the shard key, the chunk of a document is determined **by means of a hash function**. Assuming that accessing one document takes one time unit (existing indexes are used at no cost) and we have 6,000 documents in the collection, k of which have value “YYY” for attribute “other”, how many time units would take the following operations⁵:

- e) *FindOne*($\{_id : "XXX"\}$)
- f) *Find*($\{_id : \{\$in : [1, \dots, 3000]\}\}$), being $[1, \dots, 6000]$ the range of existing IDs.
- g) *Find*($\{other : "YYY"\}$), being the attribute indexed.
- h) *Find*($\{other : "YYY"\}$), being the attribute NOT indexed.

8. Assume you have a MongoDB collection which occupies 6 chunks and is **evenly distributed** in 3 shards (i.e., 2 chunks per shard) . Being the document Id also the shard key, the chunk of a document is determined **by range**. Assuming that accessing one document takes one time unit (existing indexes are used at no cost) and we have 6,000 documents in the collection, k of which have value “YYY” for attribute “other”, how many time units would take the following operations⁶:

- i) *FindOne*($\{_id : "XXX"\}$)
- j) *Find*($\{_id : \{\$in : [1, \dots, 3000]\}\}$), being $[1, \dots, 6000]$ the range of existing IDs.
- k) *Find*($\{other : "YYY"\}$), being the attribute indexed.
- l) *Find*($\{other : "YYY"\}$), being the attribute NOT indexed.

⁴As typically in RDBMS optimizers, assume uniform distribution of values and statistical independence between pairs of attributes.

⁵As typically in RDBMS optimizers, assume uniform distribution of values and statistical independence between pairs of attributes.

⁶As typically in RDBMS optimizers, assume uniform distribution of values and statistical independence between pairs of attributes.

6 New Relational Architecture

6.1 Theoretical questions

1. Briefly explain the concept of associativity in the context of memory usage. Identify the benefit of high (respectively low) associativity.

6.2 Problems

1. Assume you have a SanssouciDB table T stored **row-wise**, which occupies 300 blocks and contains tuples with three variable length attributes [A, B, C] (underlined attribute is declared to be the primary key of the table). Assume there is **no index** and give the **average cost** of each query assuming accessing one block is one second.

- SELECT A,B,C FROM T WHERE A=x; (being x a constant)
- SELECT SUM(B) FROM T;

2. Assume you have a SanssouciDB table T stored **column-wise**, which occupies 300 blocks and contains tuples with three attributes [A, B, C] (underlined attribute is declared to be the primary key of the table). Assume there is **no index**, storage of each attribute uses exactly the same space after compression (i.e., 100 blocks), and run length encoding has been applied for non-key attributes storing ending row position per run. Give the **average cost** of each query assuming accessing one block is one second and explicit any other assumption you make.

- SELECT A,B,C FROM T WHERE A=x; (being x a constant)
- SELECT SUM(B) FROM T;

3. Assume that there is a table T with attributes [A,B,C] (the underlined attribute is the primary key of the table), which occupies 300 disk blocks of 8KB each, with 100bytes per row. All three attributes require the same space (even in case of compression). Supposing that there is not any index, which would be the **minimum amount of memory required to process** (do not consider the memory required to store the final output of the query) the following query in either row storage or column storage (briefly **justify** your answer and explicit any assumption you make).

SELECT A FROM T WHERE B='x' AND C='y';

- Row storage
- Column storage

4. Without considering compression and assuming a disk block size of 8KB, is there any case where a query retrieving all tuples, but only half of the equally-sized attributes of a relational table performs better in **row storage** without any kind of vertical partitioning than in **columnar storage**? Numerically justify your answer.

5. Represent the given column with dictionary and run-length encoding storing end row position.

A
A
B
B
B
C
A
A

Dictionary

End Row

Code

Table

6. Given the two columns of a table represented with run-length encoding using dictionary, give the position of the row(s) that fulfill the predicate "Charlie AND Beta". Briefly explain how a column store would obtain them.

<table><tr><td>Bravo</td></tr><tr><td>Charlie</td></tr></table> <p>Dictionary</p> <table><tr><td>2</td><td>1</td></tr><tr><td>3</td><td>0</td></tr><tr><td>3</td><td>2</td></tr></table> <p># values Dictionary positions</p>	Bravo	Charlie	2	1	3	0	3	2	<table><tr><td>Alpha</td></tr><tr><td>Beta</td></tr></table> <p>Dictionary</p> <table><tr><td>1</td><td>1</td></tr><tr><td>5</td><td>0</td></tr><tr><td>2</td><td>2</td></tr></table> <p># values Dictionary positions</p>	Alpha	Beta	1	1	5	0	2	2
Bravo																	
Charlie																	
2	1																
3	0																
3	2																
Alpha																	
Beta																	
1	1																
5	0																
2	2																

7 Distributed Processing Frameworks

7.1 MapReduce in use

7.1.1 Theoretical questions

1. Classify a MapReduce system as either query-shipping (i.e., the evaluation of the query is delegated to the site where the data are stored), data-shipping (i.e., the data are moved from the stored site to the site executing the query) or hybrid (i.e., both query- and data-shipping). Clearly justify your answer.

7.1.2 Problems

1. Consider the following implementation of the Cartesian Product with MapReduce (\oplus stands for concatenation) and answer the questions accordingly.

$$T \times S \implies \begin{cases} \text{map}(\text{key } k, \text{value } v) \mapsto \\ \begin{cases} [(h_T(k) \bmod D, k \oplus v)] & \text{if } \text{input}(k \oplus v) = T \\ [(0, k \oplus v), \dots, (D-1, k \oplus v)] & \text{if } \text{input}(k \oplus v) = S \end{cases} \\ \text{reduce}(\text{key } ik, \text{vset } ivs) \mapsto \\ \begin{cases} [\text{crossproduct}(T_{ik}, S)] \\ T_{ik} = \{iv | iv \in ivs \wedge \text{input}(iv) = T\}, \\ S = \{iv | iv \in ivs \wedge \text{input}(iv) = S\} \end{cases} \end{cases}$$

- (a) Which is the relationship of D with parallelism and scalability, if any?
 - (b) Which is the optimal value for D ?
2. Provide the MapReduce pseudo-code implementation of the following Relational operators. You can use “ \oplus ” symbol for concatenation, “ $\text{prj}_{a_{i_1} \dots a_{i_n}}(t)$ ” to get attributes “ $a_{i_1} \dots a_{i_n}$ ” in t , “ $\text{crossproduct}(S_1, S_2)$ ” to perform the cross product of two sets of rows, “ $\text{input}(t)$ ” to know if t comes from either T or S , and assume the “key” parameter contains the PK of the table and the “value” one all the others.
 - Aggregation ($\gamma_{A, \text{sum}(B)}(T)$)
 - Selection ($\sigma_{A=x'}(T)$)
 - Join ($T \bowtie_{T.A=S.B} S$)
 - Union ($T \cup S$)
 - Difference ($T \setminus S$)
 - Intersection ($T \cap S$)
 3. In relational algebra, the antijoin operator (\bowtie) is defined as the complement of the semijoin on the primary keys (PKs). Formally, assuming A and B are the PKs of R and S , respectively, then $R \bowtie S = R \setminus R \bowtie_{A=B} S$. Provide the MapReduce pseudo-code implementation of the antijoin operator. Assume the existence of the operator \oplus to concatenate strings, $\text{prj}_{att}(s)$ to project attribute att from the tuple s , and $\text{input}(s)$ to decide the origin (i.e., R or S) from a tuple s .

4. **Modify** the following **MapReduce code** in Python so that given a file with key-value pairs representing the edges of a graph⁷ (i.e., the key is the label of the first node and the value is the label of the second), it identifies all pairs of nodes connected by a path of length two (if two nodes are connected multiple times, they will be repeated in the output). Do not need to optimize it in any way, just make it work as expected.

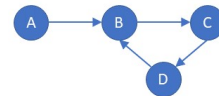
```
def map(k,v):
    return [(k,(k,v))]

def reduce(k, lv):
    outgoing = []
    incoming = []
    for v in lv:
        if k == v[0]:
            outgoing.append(v[1])
        elif k == v[1]:
            incoming.append(v[0])
    retValue = []
    for o in outgoing:
        for i in incoming:
            retValue.append(i+"-2->" +o)
    return retValue
```

Example of Input

Key Value

A B
B C
C D
D B



Expected output

A -2-> C
D -2-> C
B -2-> D
C -2-> B

7.2 MapReduce internals

7.2.1 Problems

1. Assume the following MapReduce program:

```
public void map(LongWritable key, Text value) {
    String line = value.toString();
    StringTokenizer tokenizer = new StringTokenizer(line);
    while (tokenizer.hasMoreTokens()) {
        write(new Text(tokenizer.nextToken()), new IntWritable(1));
    }
}

public void reduce(Text key, Iterable<IntWritable> values) {
    int sum = 0;
    for (IntWritable val : values) {
        sum += val.get();
    }
    write(key, new IntWritable(sum));
}
```

Consider the following data set:

- Block0: "a b b a c | c d c e a"
- Block1: "a b d d a | b b c c f"

Simulate the execution of the MapReduce code given the following configuration:

- The **map** and **reduce functions** are those of the wordcount
 - The **combine** function shares the implementation of the reduce
- One **Split** is one block

⁷Without reflexive edges like (A,A).

- The ”|” divides the **records** inside each block
 - We have two records per block
- Hadoop is configured with the parameter *dfs.replication=1*
- We can keep four pairs [key,value] per **spill**
- We have two **mappers** and two **reducers**
 - *Machine0*, contains block0, runs mapper0 and reducer0
 - *Machine1*, contains block1, runs mapper1 and reducer1
- The hash function used to shuffle data to the reducers uses the correspondence:
 - {b,d,f} → 0
 - {a,c,e} → 1

Fill the gaps in each step (numbers correspond to the phase in the MapReduce algorithm):

- 1) *Machine0* contains ... block(s).
Machine1 contains ... block(s).
- 2) We keep ... replica(s) (including the master copy) per block.
- 3) We have ... split(s) per machine.
- 4) Mapper0 reads ... records.
Mapper1 reads ... records.
- 5) Memory content during each spill in *Machine0*:

Spill1	Spill2	Spill3	Spill4
[,] [,]	[,] [,]	[,] [,]	[,] [,]
[,] [,]	[,] [,]	[,] [,]	[,] [,]

Spills in *Machine1*:

Spill1	Spill2	Spill3	Spill4
[,] [,]	[,] [,]	[,] [,]	[,] [,]
[,] [,]	[,] [,]	[,] [,]	[,] [,]

- 6) a) Partitions in *Machine0*:

Spill1		Spill2		Spill3	
Partition0	Partition1	Partition0	Partition1	Partition0	Partition1
[,] [,]	[,] [,]	[,] [,]	[,] [,]	[,] [,]	[,] [,]
[,] [,]	[,] [,]	[,] [,]	[,] [,]	[,] [,]	[,] [,]

Partitions in *Machine1*:

Spill1		Spill2		Spill3	
Partition0	Partition1	Partition0	Partition1	Partition0	Partition1
[,] [,]	[,] [,]	[,] [,]	[,] [,]	[,] [,]	[,] [,]
[,] [,]	[,] [,]	[,] [,]	[,] [,]	[,] [,]	[,] [,]

- b) Partitions in *Machine0*:

Spill1		Spill2		Spill3	
Partition0	Partition1	Partition0	Partition1	Partition0	Partition1
[,] [,]	[,] [,]	[,] [,]	[,] [,]	[,] [,]	[,] [,]
[,] [,]	[,] [,]	[,] [,]	[,] [,]	[,] [,]	[,] [,]

Partitions in *Machine1*:

Spill1		Spill2		Spill3	
Partition0	Partition1	Partition0	Partition1	Partition0	Partition1
[,][,]	[,][,]	[,][,]	[,][,]	[,][,]	[,][,]
[,][,]	[,][,]	[,][,]	[,][,]	[,][,]	[,][,]

c) Files in *Machine0*:

File0.0	File0.1	File0.2
[,][,][,][,]	[,][,][,][,]	[,][,][,][,]
File0.3	File0.4	File0.5
[,][,][,][,]	[,][,][,][,]	[,][,][,][,]

Files in *Machine1*:

File1.0	File1.1	File1.2
[,][,][,][,]	[,][,][,][,]	[,][,][,][,]
File1.3	File1.4	File1.5
[,][,][,][,]	[,][,][,][,]	[,][,][,][,]

7) Merges in *Machine0*:

Merge0	Merge1
[,][,][,][,]	[,][,][,][,]

Merges in *Machine1*:

Merge0	Merge1
[,][,][,][,]	[,][,][,][,]

8) Files in *Machine0*:

File0.0	File0.1
[,][,][,][,]	[,][,][,][,]

Files in *Machine1*:

File1.0	File1.1
[,][,][,][,]	[,][,][,][,]

9) Reducer0 reads ... from *Machine0* and ... from *Machine1*.
 Reducer1 reads ... from *Machine0* and ... from *Machine1*.

10) Merges in *Machine0*:

Merge0	Merge1
[,{ }][,{ }][,{ }][,{ }]	[,{ }][,{ }][,{ }][,{ }]

Merges in *Machine1*:

Merge0	Merge1
[,{ }][,{ }][,{ }][,{ }]	[,{ }][,{ }][,{ }][,{ }]

- 11) Reduce function is executed ... times in *Machine0*.
 Reduce function is executed ... times in *Machine1*.

12) Files in *Machine0*:

Output0
[,][,][,][,]

Files in *Machine1*:

Output1
[,][,][,][,]

2. Let's suppose that we have a cluster of 100 machines and a MapReduce job with 1.000.000 key-value pairs in the input that generate 100.000 pairs in the output. Assume that both map and reduce functions generate one pair in the output per call. Assuming the reduce function is commutative and associative, is it worth to use the combine function? Briefly justify your answer.
3. Let's consider that the CPU time to process a key-value pair is always zero (both in the map and the reduce functions), as well as the time to retrieve a disk block, and we can send control messages between machines without any cost. However, our network bandwidth is 10MB/sec and we only have one processor per machine. We have ten machines, which have a direct connection between them (i.e., sending data from machine A to machine B is always one hop away, for any pair of machines, as depicted in the figure).

We are launching a MapReduce job with a mapper and a reducer in each of these machines and an input file of key-value pairs of an overall size of 1GB⁸ (whose chunks are already uniformly distributed in the ten machines, without any replication and not necessary of default size). The map function generates exactly one key-value pair per call, with exactly the same size as the key-value pair in the input. Assume that the hash function (without any cost) used in the partition phase perfectly maps the values in a uniform way; and also that when two machines are transferring data from one to another, they cannot perform any other data transfer activity in parallel (i.e., while A is transferring data to B, neither A nor B can send or receive anything from anyone else until this transfer is over, A can only send to B and not receive anything from anyone, B can only receive from A and not send anything to anyone, not even A). Supposing that scheduling the data transfer in the cluster of ten machines is perfect and there is not any idle period for any of the machines, which is the **duration of the job**? Explicit any **calculation and assumption** you need to make.

⁸Assume 1GB=1000MB.

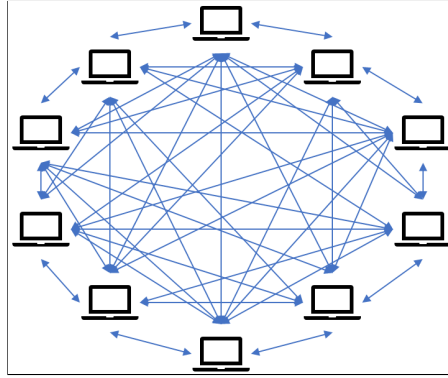


Figure 5: Fully connected architecture

Concept	Cost
Process any key-value pair	0s
Retrieve a disk block	0s
Send control messages between machines	0s
Network bandwidth	10MB/s

7.3 Spark in use

7.3.1 Problems

1. Consider a file (*wines.txt*) containing the following data:

```
wines.txt
type_1,2.064254
type_3,2.925376
type_2,2.683955
type_1,2.991452
type_2,2.861996
type_1,2.727688
```

Provide the ordered list of Spark operations (no need to follow the exact syntax, just the kind of operation and main parameters) you'd need to retrieve the minimum value per type. Do not use SQL and minimize the use of other Python libraries or code. Save the results in *output.txt*.

2. Consider two files containing the following data:

```
Employees.txt
EMP1;CARME;400000;MATARO;DPT1
EMP2;EUGENIA;350000;TOLEDO;DPT2
EMP3;JOSEP;250000;SITGES;DPT3
EMP4;RICARDO;250000;MADRID;DPT4
EMP5;EULALIA;150000;BARCELONA;DPT5
EMP6;MIQUEL;125000;BADALONA;DPT5
EMP7;MARIA;175000;MADRID;DPT6
EMP8;ESTEBAN;150000;MADRID;DPT6
```

```
Departments.txt
DPT1;DIRECCIO;10;PAU CLARIS;BARCELONA
DPT2;DIRECCIO;8;RIOS ROSAS;MADRID
DPT3;MARKETING;1;PAU CLARIS;BARCELONA
DPT4;MARKETING;3;RIOS ROSAS;MADRID
DPT5;VENDES;1;MUNTANER;BARCELONA
DPT6;VENDES;1;CASTELLANA;MADRID
```

Provide the ordered list of Spark operations (no need to follow the exact syntax, just the kind of operation and main parameters) you'd need to retrieve for each employee his/her department information. Do not use SQL and minimize the use of other Python libraries or code. Save the results in *output.txt*.

3. Consider an error log file (*log.txt*) like the one bellow:

```
log.txt
20150323;0833;ERROR;Oracle
20150323;0835;WARNING;MySQL
20150323;0839;WARNING;MySQL
20150323;0900;WARNING;Oracle
20150323;0905;ERROR;MySQL
20150323;1013;OK;Oracle
20150323;1014;OK;MySQL
20150323;1055;ERROR;Oracle
```

Provide the ordered list of Spark operations (no need to follow the exact syntax, just the kind of operation and main parameters) you'd need to retrieve the lines corresponding to both *errors* and *warnings*, but adding *Important:* at the beginning of the identifier of those of errors (i.e., only *errors*). Do not use SQL and minimize the use of other Python libraries or code. Save the results in *output.txt*.

4. Assume that “spark” variable is a **Spark** session and the `dataset.csv` contains the two columns “Color” and “Radius”. Clearly **identify the problems** you find in the following Spark code and propose some fix to obtain the expected result.

```
df0= spark.read.csv("dataset.csv", \
                    header='true', \
                    inferSchema='true', \
                    sep=',')
df1 = df0.select("Color", "Radius")
df2 = df1.withColumn("Pi", 3.141592)
df3 = df2.withColumn("Area", df2.Radius * df2.Radius * df2.Pi)
df4 = df3.groupBy("Color").max("Area")
df5 = df4.sort("Color")
```

5. **Given two files containing** the following kinds of data:

Employees.txt with fields: EmployeeID; EmployeeName; YearlySalary; CityOfResidence; SiteOfWork
 EMP4;RICARDO;250000;MADRID;DPT4
 EMP5;EULALIA;150000;BARCELONA;DPT5
 EMP6;MIQUEL;125000;BADALONA;DPT5
 EMP7;MARIA;175000;MADRID;DPT6
 EMP8;ESTEBAN;150000;MADRID;DPT6
 ...

Departments.txt with fields: SiteID; DepartmentName; StreetNumber; StreetName; City
 DPT1;DIRECCIO;10;PAU CLARIS;BARCELONA
 DPT2;DIRECCIO;8;RIOS ROSAS;MADRID
 DPT3;MARKETING;1;PAU CLARIS;BARCELONA
 DPT4;MARKETING;3;RIOS ROSAS;MADRID
 ...

Consider the following **PySpark code** and answer the questions bellow.

```
source1 = spark.read.format("csv").load("employees.txt", header='false', inferSchema='true', sep=";")
source2 = spark.read.format("csv").load("departments.txt", header='false', inferSchema='true', sep=";")
A = source1.toDF("eID","eName","eSalary","eCity","eDpt")
B = source2.toDF("dID","dArea","dNumber","dStreet","dCity")
C = A.select(A.eCity.alias("city"))
D = B.select("dArea")
E = D.crossJoin(C)
F = B.select("dArea",B.dCity.alias("city"))
G = E.subtract(F)
H = G.select("dArea")
result = D.subtract(H)
```

- (a) State in natural language the corresponding query it would answer?
- (b) Clearly indicate any mistake or improvement you can fix/make in the code? For each of them give (1) the line number, (2) pseudo-code to implement the fix, and (3) brief rationale.
6. Given two files containing the following kinds of data:

Employees.txt with fields: EmployeeID; EmployeeName; YearlySalary; CityOfResidence; SiteOfWork
 EMP1;RICARDO;250000;MADRID;DPT1
 EMP2;EULALIA;150000;BARCELONA;DPT2
 EMP3;MIQUEL;125000;BADALONA;DPT3
 EMP4;MARIA;175000;MADRID;DPT4
 EMP5;ESTEBAN;150000;MADRID;DPT3
 ...

Departments.txt with fields: SiteID; DepartmentName; StreetNumber; StreetName; City
 DPT1;DIRECCIO;10;PAU CLARIS;BARCELONA
 DPT2;DIRECCIO;8;RIOS ROSAS;MADRID
 DPT3;MARKETING;1;PAU CLARIS;BARCELONA
 DPT4;MARKETING;3;RIOS ROSAS;MADRID
 ...

Give a sequence of Spark operations in pseudo-code (resembling PySpark) to obtain for each city where employees that work in a site of a department in Barcelona live, the sum of the salaries of those employees. The result for the exemplary data would be:

MADRID;400000
 BADALONA;125000
 ...

7. Consider three files containing the following kinds of data:

Employees.txt
 EMP1,CARME,400000,MATARO,DEPT1,PROJ1
 EMP2,EULALIA,150000,BARCELONA,DEPT2,PROJ1
 EMP3,MIQUEL,125000,BADALONA,DEPT1,PROJ3

Projects.txt
 PROJ1,IBDTEL,TV,1000000
 PROJ2,IBDVID,VIDEO,500000
 PROJ3,IBDTEF,TELEPHONE,200000
 PROJ4,IBDCOM,COMMUNICATIONS,2000000

Departments.txt
 DEPT1,MANAGEMENT,10,PAU CLARIS,BARCELONA
 DEPT2,MANAGEMENT,8,RIOS ROSAS,MADRID
 DEPT4,MARKETING,3,RIOS ROSAS,MADRID

Provide the ordered list of Spark operations (no need to follow the exact syntax, but just the kind of operation and main parameters) you would need to obtain the departments with all employees assigned to the same project. The result must include department number. Save the results in `output.txt`. In the previous example, the result should be *DEPT2* and *DEPT4*.

8. Consider two files containing the following kinds of data:

Employees.txt
 EMP4;RICARDO;250000;MADRID;DPT4
 EMP5;EULALIA;150000;BARCELONA;DPT5
 EMP6;MIQUEL;125000;BADALONA;DPT5
 EMP7;MARIA;175000;MADRID;DPT6
 EMP8;ESTEBAN;150000;MADRID;DPT6

Departments.txt
 DPT1;DIRECCIO;10;PAU CLARIS;BARCELONA
 DPT2;DIRECCIO;8;RIOS ROSAS;MADRID
 DPT3;MARKETING;1;PAU CLARIS;BARCELONA
 DPT4;MARKETING;3;RIOS ROSAS;MADRID

Provide the ordered list of Spark operations (no need to follow the exact syntax, but just the kind of operation and main parameters) you'd need to retrieve the list of department IDs for those departments with workers from all cities where there are employees. Save the results in `output.txt`.

9. Consider two files containing the following kinds of data:

Employees.txt
 EMP1;RICARDO;250000€;MADRID;SITE2
 EMP2;EULALIA;150000€;BARCELONA;SITE1
 EMP3;MIQUEL;125000€;BADALONA;SITE3
 EMP4;MARIA;175000€;MADRID;SITE2
 EMP5;ESTEBAN;150000€;MADRID;SITE4

Departments.txt
SITE1;DPT.MANAGEMENT;FLOOR10;ST.PAU CLARIS;BARCELONA
SITE2;DPT.MANAGEMENT;FLOOR8;ST.RIOS ROSAS;MADRID
SITE3;DPT.MARKETING;FLOOR1;ST.PAU CLARIS;BARCELONA
SITE4;DPT.MARKETING;FLOOR1;ST.RIOS ROSAS;MADRID
SITE5;DPT.MARKETING;FLOOR5;ST.MARTI PUJOL;BADALONA

Provide the ordered list of Spark operations (no need to follow the exact syntax, but just the kind of operation and main parameters) you'd need to retrieve the list of department IDs for those departments with sites in all cities where employees live (these employees can be even from other departments). Save the results in `output.txt`. In the previous example, the result should be *DPT.MANAGEMENT*, because it has sites in all the three cities where there are employees (i.e., MADRID, BARCELONA and BADALONA). However, *DPT.MANAGEMENT* should not be in the result, because it does not have any site in BADALONA, where EMP3 lives.

10. Consider three files relating to a bibliographic database: `author.csv` relates authors with papers (you may assume that author names are unique, that authors have one or more papers, and that papers have one or more authors); `title.csv` gives the title of a paper (you may assume a paper has one title, but one title may be shared by many papers); and `citation.csv` indicates which papers cite which other papers (you may assume that each paper cites at least one other paper, that a paper may be cited zero or more times, and that a paper cannot cite itself).

author.csv		title.csv	
AUTHOR	PAPERID	PAPERID	TITLE
...	...	GP2014	Semantics of SPARQL
C. Gutierrez	GP2014	AGP2013	Deduction for RDF
C. Gutierrez	AGP2013	GZ2011	Graph databases
C. Gutierrez	GZ2011
...	...	citation.csv	
J. Perez	GP2014	CITES	CITED
J. Perez	AGP2013
J. Perez	P2017	GP2014	AGP2013
...	...	AGP2013	GZ2011
R. Angles	AGP2013	P2017	AKK2016
R. Angles	AKK2016
...

The headers are shown for illustration here. They do not need to be considered.

The count of self-citations for an author A , denoted $\text{self}(A)$, is defined as the number of citation pairs (P_1, P_2) where A is an author of both. The count of citations given by an author A , denoted $\text{give}(A)$, is the count of citation pairs (P_1, P_2) such that A is an author of P_1 . The count of citations received by A , denoted $\text{receive}(A)$, is the count of citation pairs (P_1, P_2) where A is an author of P_2 . The ratio of self-citations to all citations given and received are then defined, respectively, as $\frac{\text{self}(A)}{\text{give}(A)}$ and $\frac{\text{self}(A)}{\text{receive}(A)}$. In case that $\text{receive}(A) = 0$, you should omit the author A from the results (note that $\text{give}(A)$ cannot be 0 as an author must have at least one paper and a paper must have at least one citation). We provide an example output for the input data:

AUTHOR	SELF GIVE RATIO	SELF RECEIVE RATIO
C. Gutierrez	1.000	1.000
J. Perez	0.333	1.000
R. Angles	0.000	0.000
...

Headers are only shown for illustration. We will use Apache Spark to perform the analysis and compute the output. You should not assume any ordering of the input files. You do not need to order the output file in any particular way.

Given this input and desired output, design a Spark process to complete the required processing. In particular, you should draw the high-level DAG of operations that the Spark process will perform, detailing the sequence of transformations and actions. You should briefly describe what each step does, clearly indicating which steps are transformations and which are actions. You should also indicate which RDDs are virtual and which will be materialized. You should use caching if appropriate. You should provide details on any functions passed as arguments to the transformations/actions you use.

7.4 Spark internals

7.4.1 Theoretical questions

1. What indicates to Spark query optimizer the end of a stage and the beginning of the next one?

7.4.2 Problems

1. Considering the file and result example, briefly indicate the problems you find in the Spark code below (if any), and **modify the code** to fix them (if needed).

Employees.txt	
EMP1;CARME;40000;MATARO;DPT1;PROJ1	
EMP2;EUGENIA;35000;TOLEDO;DPT2;PROJ1	
EMP3;JOSEP;25000;SITGES;DPT3;PROJ2	
EMP4;RICARDO;25000;MADRID;DPT4;PROJ2	
EMP5;EULALIA;15000;BARCELONA;DPT5;PROJ2	
EMP6;MIQUEL;12500;BADALONA;DPT5;PROJ3	
EMP7;MARIA;17500;MADRID;DPT6;PROJ3	
EMP8;ESTEBAN;15000;MADRID;DPT6;PROJ3	

	Expected result
	[PROJ1,[EUGENIA, 35000, 37500.0]]
	[PROJ2,[EULALIA, 15000, 21666.6]]
	[PROJ3,[MIQUEL, 12500, 15000.0]]

```
rawEmps = sc.textFile(r"Employees.txt")
emps = rawEmps.map(lambda l: tuple(l.split(";"))).cache()
RDD1 = emps.map(lambda t: (t[5], (int(t[2]), 1)))
RDD2 = RDD1.reduceByKey(lambda t1, t2: (t1[0] + t2[0])/(t1[1] + t2[1]))
RDD3 = emps.map(lambda t: (t[5], t))
RDD4 = RDD3.join(RDD2)
RDD5 = RDD4.filter(lambda t: int(t[1][0][2])<t[1][1])
RDD6 = RDD5.map(lambda t: (t[0], (t[1][0][1],t[1][0][2],t[1][1])))
Result = RDD6.sortByKey()
```

2. Consider the following pipeline. This pipeline runs in a 4-machine cluster with HDFS to store the files and Spark to execute it. File1 and File2 are distributed in the cluster in 6 chunks each.

```
RddF1:= sc.textFile("... file1.txt")
RddF2:= sc.textFile("... file2.txt")
Rdd2:= RddF1.mapToPair(s.split(";")[0], s.split(";")[1-2])
Rdd3:= Rdd2.GroupByKey()
Rdd4:= Rdd3.MapValues(f1)
Rdd5:= RddF2.mapToPair(s.split(";")[0], s.split(";")[1-2])
Rdd6:= Rdd5.GroupByKey()
Rdd7:= Rdd6.MapValues(f2)
Rdd8:= Rdd4.join(Rdd7)
Rdd8.save("... file3.txt")
```

- a) How many stages will the scheduler generate for this pipeline? (Justify your answer)
- b) How many tasks will be generated within each stage? (Justify your answer)

3. Consider the legacy code written in MapReduce that specifies a `map()`, `combine()` and `reduce()` operations. This job reads from a text file f_1 . The combine and reduce operations coincide and you can assume all functions are correct. Write a Spark pipeline **equivalent to the MapReduce** job. Use `fmap` and `freduce` to refer to the code executed inside the map and combine / reduce operations. You can parametrise the `fmap` and `freduce` functions but resulting in minimal code adaptation.
4. Given a file of $3.2GB$ stored in an HDFS cluster of 50 machines, and containing $16 \cdot 10^5$ key-value pairs in a `SequenceFile`; estimate the execution time of a Spark job containing a single map transformation and an action storing the results in a file. Explicit any assumption you make and consider also the following parameters:
- Chunk size: $128MB$ (default)
 - Replication factor: 3 (default)
 - Map function (i.e., the parameter of the transformation) execution time: 10^{-3} sec/call (this is **the only cost** you have to consider)
 - Save action execution time: 0sec (do not consider its cost at all)

8 Streams

8.1 Problems

1. Let's suppose an arrival rate $\lambda = 50\text{messages/sec}$ and a service rate $\mu = 50\text{messages/sec}$. Given the message arrivals in the table below and a size of 100bytes per message, which is the minimum buffer size needed to guarantee that no message is lost?

Instant (sec)	1	2	3	4	5	6	7	8	9	10
#Messages	2	10	100	50	20	75	150	100	40	10

2. Let's suppose a stationary situation in the processing of a data stream, where the whole memory available is occupied by messages. Assume also the existence of a separated buffer big enough to allow us the complete processing of these messages already in memory without being concerned with new arrivals. Assume we have $M = 10$ memory pages with $R_S = 10$ messages per page. If the processing of these messages consists only of a lookup and in-memory processing/comparison of a message and a tuple, whose execution time is thousand times smaller than a disk access, give the cost in these two situations considering that every disk block of the table contains also $R_T = 10$ tuples:
 - a) Going through an index for each message with cost $h = 3$ for the index and one more disk access to the lookup table per message.
 - b) Bringing each block of the lookup table (whose size is $B = 100$) into memory one by one and checking all messages in memory against all its tuples.
3. Let's suppose that setting the execution environment for the processing of messages in a stream is $S = 100\text{ms}$ (e.g., placing&retrieving all the information to/from the stack), independently of their number. Both packing and unpacking the messages in a batch (i.e., a list of elements) have the same cost of $Pk = 1\text{ms}$ per message, and processing each message is $Ps = 10\text{ms}$.
 - a) Which is the cost of processing truly streaming 10 messages one at a time?
 - b) Which is the cost of processing one packed micro-batch of 10 messages at once?
4. Let's suppose we have a log file recording the events coming from different machines. Thus, for each event we have the following information:

$$(\log ID, \text{traceID}, \text{eventID}, \text{duration})$$

The *logID* corresponds to the IP of the machine; the *traceID* identifies the transaction inside the machine (i.e., two traceIDs can coincide in different machines); the *eventID* identifies the kind of action performed by the machine; finally, the *duration* is the number of milliseconds taken to implement the action. Assuming that we cannot keep all log entries in memory, and we decide to randomly sample them, **give the attributes (up to three) you would use as parameters of the hash function implementing such sampling**, so that each of the following queries gives a result as accurate as possible.

- a) For each machine, estimate the average number of actions per transaction.
...
...
...

b) Estimate the average sum of the duration per transaction.

...
...
...

c) Estimate the fraction of transactions where the same kind of action appears more than once.

...
...
...

d) Estimate the number of transactions with more than two actions taking more than 100ms.

...
...
...

5. Let's suppose we have a black-list of IP addresses (whose packages we do not want to cross our firewall), which is too long to be kept in memory (10^7 elements). Thus, we decide to implement a **Bloom filter** (to avoid further processing of black-listed addresses), with only 10^8 bits.⁹

a) How many hash functions would you use?

b) What's the probability of a false positive in that case?

c) Briefly explain what is the consequence of a false positive in this concrete case.

6. Let's suppose we have a log file recording the events coming from different machines. Thus, for each event we have the following information:

$(logID, traceID, eventID, duration)$

The *logID* corresponds to the IP of the machine; the *traceID* identifies the transaction inside the machine; the *eventID* identifies the kind of action performed by the machine; finally, the *duration* is the number of milliseconds taken to implement the action. Consider the following table and assume that each machine generates the same number of events and at the same pace, and use an **exponentially decaying window model** with a constant $c = 0.5$.

Time	logID	traceID	eventID	duration
1	1	1	A	1
2	2	1	B	100
3	1	1	C	10
4	2	1	D	10
5	1	1	A	100
6	2	1	C	1

a) Give the milliseconds (and details of calculation) used per machine (i.e., for Machine 1 and Machine 2) when the last event arrives, provided that the current milliseconds are more valuable than the oldest ones.

⁹Note: $\ln 2 \approx 0.693$

7. Assume we ingest a stream with an event every time a ticket is sold at a theater. Precisely, the stream has the following structure:

$(movieID, theaterID, timestamp, price)$

Next, we ingest the following ordered set of events:

- $(m1, t1, 12h, 10\text{€})$
- $(m2, t1, 14h, 12\text{€})$
- $(m2, t2, 15h, 18\text{€})$
- $(m1, t3, 18h, 6\text{€})$
- $(m3, t2, 19 : 15h, 13\text{€})$
- $(m1, t3, 19 : 30h, 10\text{€})$
- $(m2, t3, 19 : 45h, 25\text{€})$
- $(m3, t1, 20h, 17\text{€})$
- $(m1, t3, 20 : 30h, 10\text{€})$
- $(m2, t2, 21h, 8\text{€})$

Provide a detailed answer (i.e., describe the process) for the following questions:

- (a) Which theaters would be considered heavy hitters by the algorithm, given a required frequency of 50%?
 - (b) Under the exponentially-decaying window model, using an aging constant of $c = 0.1$ and a purging threshold of 0.6 (i.e., we'll remove movies as soon as their popularity falls strictly below 0.6). What would be the most popular movies at 21h, considering an aging tick every 15 minutes? Give the value corresponding for each of the movies in that moment.
8. If you have a stream of messages that arrive at a rate of 10,000 per second accounting for 512MB per hour, how many RDDs will process your Spark Streaming process in one hour? Explicit all the assumptions you make.
9. Assume we ingest a stream with an event every time a ticket is sold at a theater. Precisely, the stream has the structure $(movieID, theaterID, timestamp, price)$. Next, we ingest the following ordered set of events:

- Event 1: $(m_3, t_4, 12h, 10\$)$
- Event 2: $(m_1, t_2, 13h, 17\$)$
- Event 3: $(m_2, t_4, 14h, 11\$)$
- Event 4: $(m_4, t_1, 15h, 8\$)$
- Event 5: $(m_1, t_3, 16h, 9\$)$
- Event 6: $(m_3, t_4, 17h, 5\$)$
- Event 7: $(m_6, t_1, 18h, 15\$)$
- Event 8: $(m_5, t_2, 19h, 12\$)$
- Event 9: $(m_7, t_5, 20h, 17\$)$
- Event 10: $(m_1, t_1, 21h, 11\$)$

Which theaters would be considered heavy hitters (using the approximate method) considering a required frequency of 33%? Provide a detailed answer (i.e., describe the content of the auxiliary structure after processing every message).

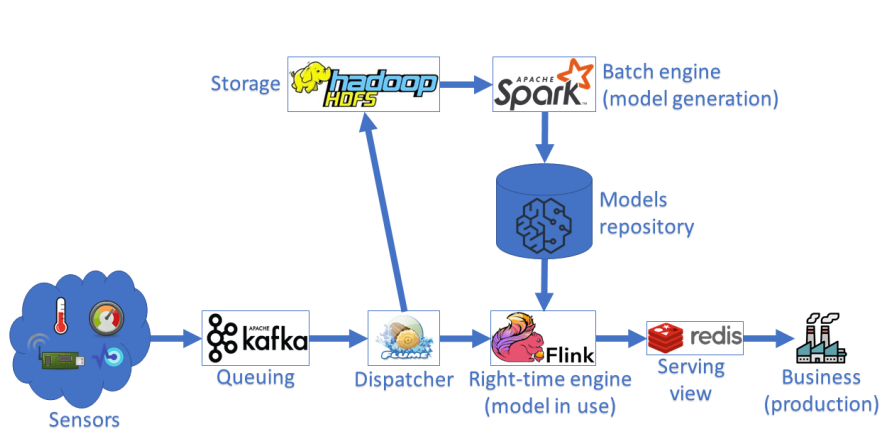
9 Architectures

9.1 Theoretical questions

1. What problem do Data Lakes solve?

9.2 Problems

1. Consider the software architecture below and briefly explain the main **management (maintenance)** risk it has.



2. Assume we have a **MongoDB** collection in a distributed cluster, which contains prices of apartments without any secondary index. Such collection is big enough not to completely fit in memory. We want to **use Spark** to compute the standard deviation per neighbourhood. Clearly identify the most efficient option and briefly justify the choice (it is **not necessary to provide the Spark code**).
 - (a) Use Spark only to push the query to MongoDB aggregation framework and simply get the result.
 - (b) Push only some of the operations to MongoDB aggregation framework and the run the rest in Spark.
 - (c) Load the whole collection to an RDD and perform all computations in Spark.