



UNIVERSITÀ DI PISA

Computer Engineering, Artificial Intelligence and Data
Engineering

Large-Scale and Multi-Structured Database

PokèMongo

Project Documentation

TEAM MEMBERS:

Edoardo Fazzari

Mirco Ramo

Olgerti Xhanej

Academic Year: 2020/2021

Contents

1	Introduction	3
1.1	Description	3
2	Analysis	5
2.1	Functional Requirements and Use Cases	5
2.1.1	Use Cases List	5
2.1.2	UML Use Cases Diagram	7
2.2	Non-Functional Requirements	9
2.3	Sources, Velocity properties and Volume of data	9
2.4	UML Entities Diagram	10
2.5	Main application queries	11
3	Project	13
3.1	Adopted Databases	13
3.2	Document Database	14
3.2.1	Queries Handled	14
3.2.2	Entities handled	14
3.2.3	Collections structure	15
3.2.4	Indexes	16
3.3	Graph Database	20
3.3.1	Queries Handled	20
3.3.2	Entities handled	20
3.3.3	Graph structure	20
3.3.4	Indexes	20
3.4	Redundancies and consistency management	20
3.5	Database Properties	20
3.5.1	Availability	20
3.5.2	Replicas	20
3.5.3	Eventual Consistency	20
3.5.4	Sharding	20
3.5.5	Pros and Drawbacks	20
3.6	Client, Server and Daemon Thread	20
3.7	Technologies and Frameworks	20
4	Implementation	21
4.1	Package structure and information hiding	21
4.1.1	Packaging strategy and information hiding	21
4.1.2	UML package diagram	21

4.2	APIs and SPIs	21
4.3	Main tools	21
4.3.1	GSON	21
4.3.2	Caching mechanism and multimedia management . . .	21
4.3.3	Password Encryptor	21
4.3.4	Logger	21
4.4	Analytics queries	21
4.4.1	User Rankings	21
4.4.2	Pokémon Rankings	21
4.4.3	Usage Statistics	21
4.4.4	Dynamic Catch Rate	21
4.5	Business logic	21
4.5.1	Points computing	21
4.5.2	Dynamic Catch Rate Computing	21
5	Test	22
5.1	Privacy and Security	22
5.2	Unit Test	22
5.3	Robustness	22
5.4	Performance	22

1 — Introduction

PokeMongo is a gaming application in which users compete each other to build up the best Team choosing between the set of Pokémon available.

1.1 Description

Every **User** can build up his own team. Every **Team** is composed by up to 6 distinct **Pokémon** and is assigned to a numerical value (points) based on features and properties of the chosen Pokémon, for ranking purposes.

A **User** can also follow other users in order to make new friends basing on common friends or common interests. Moreover users can express sentiments on **Pokémon**, choosing their favorite ones and posting or commenting on them.

Users can also navigate through the ranking in order to visualize the best teams (according to the values cited before) and the most used/caught **Pokémon**, both among their friends, grouped by country and among worldwide players.

User can browse for a specific **Pokémon** using the *Pokédex* tool, in which he/she can lookup for **Pokémon** according to search filters like *Pokémon name*, *Type* or *Points*.

Moreover, as a “real” Pokémon Trainer, the **User** is invited to *Catch ‘em‘ all*, i.e. to try to get a new **Pokémon** in order to create/update his/her own Team. Thus, it is provided to the **User** a prefix number of *daily Pokéball* to be used to try to capture them. At each **Pokémon** is associated a probability to catch it, the higher the Pokémon’s value, the lower the probability.

Furthermore, the **User** can exploit the social network structure of the application to make new **Friends** and discover new **Pokémon**. Indeed, he/she can search for new friends by *username* or choosing them among the provided recommended friends list. The **User** can choose his/her **favorite Pokémon**, obtaining in this way a shortcut to catch it faster, and can post or answer to **Posts** in order to express his/her opinion on that **Pokémon**.

In addition, to extend the dynamic behavior of the application, the *catch rate* (i.e. the probability to get a Pokémon using a Pokéball) changes in time depending on the number of **Users** who have that **Pokémon**: *the more it is popular, the harder will be to catch it*. Since the rankings’ points are computed based on the catch rate, the winning strategy could be on predicting which **Pokémon** will become popular in the near future and try to get it early! Every **User** has access to the visualization of the temporal drift of the

catch rate.

The safeguard and the improvement of the application is in charge of **Admin** users. They are able to *ban mischievous users, delete inappropriate posts or comments, add/remove Pokémon* to the collection, *consult geo-temporal usage statistics* which are useful to make new business plans.

2 — Analysis

2.1 Functional Requirements and Use Cases

2.1.1 Use Cases List

- An *unregistered user* can
 - Register
- A *registered user* can
 - Login
 - Consult Pokédex
 - * Search by Name
 - * Search by Type(s)
 - * Search by Pokédex ID
 - * Search by Catch Rate
 - * Search by Points
 - * Search by Pokemon characteristics like Height or Weight
 - Consult ranking:
 - * Most popular Pokémon among all Users
 - * Most popular Pokémon in each Country
 - * Best World Teams
 - * Best Teams among Friends
 - * Best Teams by Country
 - Find Users:
 - * See recommended users based on common friends
 - * See recommended users based on common Pokémon interests
 - * Find users by username
 - * Follow/Unfollow them
 - Interact with Pokémon network:
 - * Insert/Remove a Pokémon in his/her own favorite Pokémon list
 - * Create a post on a Pokémon to share opinions
 - * Add answers to posts

- * Follow/Unfollow them
 - * The post owner can also remove the post at his/her will
- Team handling:
 - * Remove Pokemon from the team
 - * View team
 - * Change name of the Team
 - * Save modified team
 - * View the value of the team
- Catching:
 - * Browse a Pokémon you want to catch searching it by name
 - * Select a Pokémon you want to catch from the list of favorites
 - * Try to catch a Pokemon to add to your Team
- Settings:
 - * Change Email
 - * Change Password
 - * Change Country
- Logout:
 - * Exit from the account
 - * Return to the sign in window
- At each time can:
 - * See the remaining daily Pokèballs
 - * Mute/Unmute Music
 - * By clicking on a Pokémon name, visualize all the information about it
- An *admin* can
 - Sign In
 - Add Pokèmon to the Pokédex
 - Remove Pokèmon from the Pokédex
 - See the number of registered Users in time
 - See the numbers of login per day
 - See the numbers of login per day in every Country
 - Remove a User from the system
 - Remove Posts/Answers from the system

- Consult Rankings
- Logout
- The *system* should
 - Daily update Pokeball number of each user
 - Periodically update Pokemon catch rates based on the number of users that own that pokemon
 - Update team points if the user has 6 Pokémon of different types
 - Periodically compute usage statistics to be consulted by the administrators

2.1.2 UML Use Cases Diagram

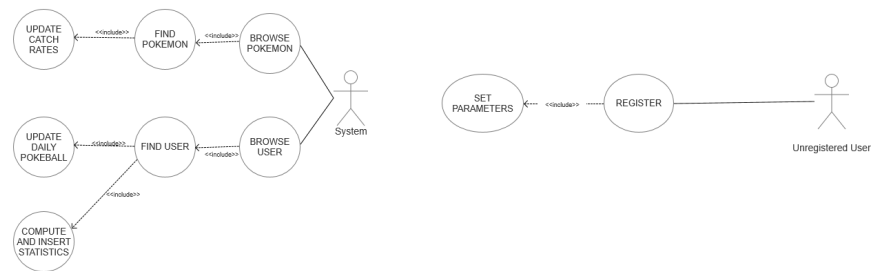


Figure 1: Use Case Diagram 1

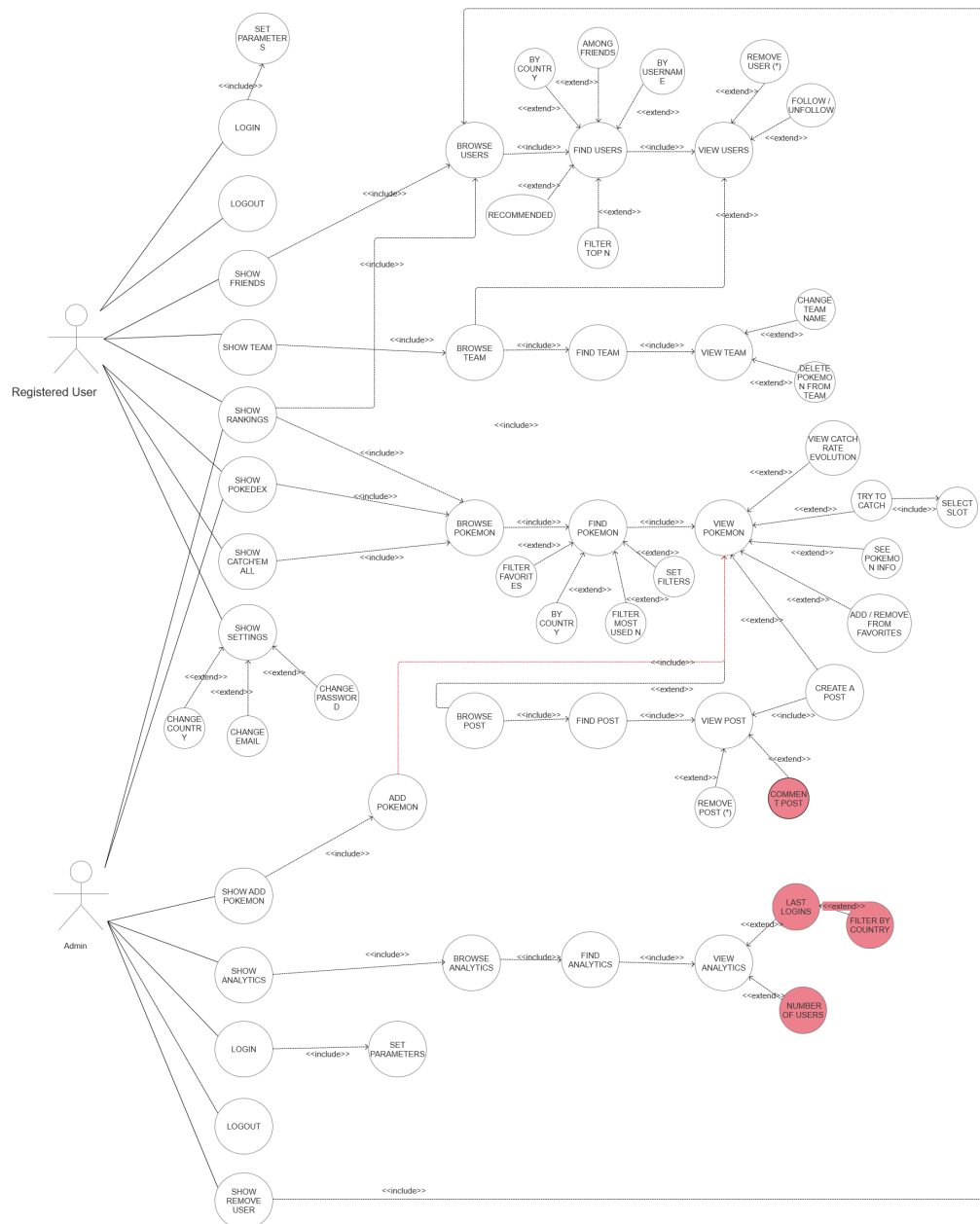


Figure 2: Use Case Diagram 2

(*) only for the User who created the Post and Admins, in **Red** Browse-find-view comments and browse-find-view answers had not been reported)

2.2 Non-Functional Requirements

- The application should guarantee a high availability. The application should guarantee a **high availability**
- It should be **easy to use**, especially for children and youngsters, and enjoyable
- It should have a **read-your-own-writes consistency** on each user's own team, so he/she can always be sure that Pokémon have been correctly caught/freed up
- The application should always provide to each user the most recent version of the rankings in order to permit him/her to immediately verify his/her progresses
- The statistics regarding usage and catch rate evolution are not needed to be real-time, they can be updated periodically and be eventually consistent
- Posts, comments and answers must follow a **causal-consistency**
- **Response time** is an important issue: redundancies and larger memory consumptions are preferred over high latencies
- **Passwords are crypted** for security reasons
- A graphical interface and the usage of multimedia are crucial for an involving game experience

2.3 Sources, Velocity properties and Volume of data

Data stored in the application backend has been downloaded and imported from the following sources:

1. **Pokémon Data** → <https://pokeapi.co>,
<https://bulbapedia.bulbagarden.net/wiki>
2. **Countries data** → <https://gist.github.com/kalincernev/486393efcca01623b18d>
3. **Data for the generation of realistic users** → <https://github.com/smashew/NameDatabases/blob/master/NamesDatabases/surnames/all.txt>

All the imported data has been modified, updated and preprocessed in order to satisfy the application needs. Users added have the only purpose of showing the application functionalities, **for privacy issues they are not real people**; anyway they have been created using *realistic criteria*.

Velocity is guaranteed by the dynamic catch rate mechanism: the popularity of a Pokémon influences both its catch rate and the amount of points that it will provide. As a consequence, Users are continuously stimulated by catching new Pokémon, in order to try to raise their amount of points: in this way old teams' data becomes quickly out-of-date.

Volume of data, considering 250K users, almost 1K Pokémon and about 500K posts is no lower than 100Mb.

2.4 UML Entities Diagram

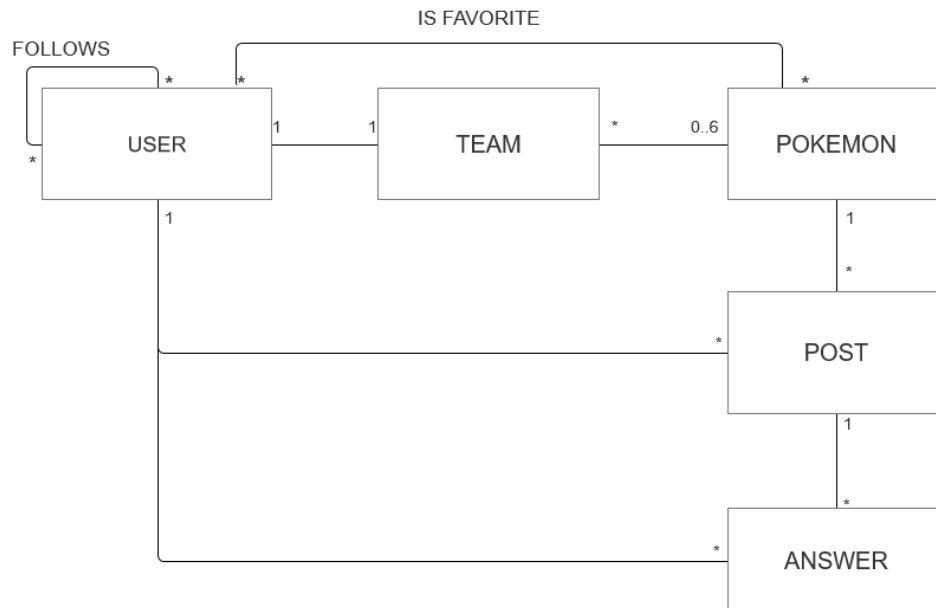


Figure 3: UML Entity Diagram

1. A **User** can build up only one **Team**: of course, each **Team** has just one owner.
2. A **Team** is composed of a maximum of six **Pokémon**, every **Pokémon** can be caught by anyone, so can belong to many **Teams**.

3. A **User** can follow many **Users**, in the meanwhile he/she can have many followers.
4. A **User** can have many favorites **Pokémon**. A **Pokémon** can be favorite of many **Users**.
5. A **Post** is created just by one **User** on one **Pokémon**. A **User** can create many posts and a **Pokémon** can have many **Posts** talking about it.
6. An **Answer** is written by one **User** and it refers to one **Post**. **Users** can submit many Answers and there can be many **Answers** behind a **Post**.

2.5 Main application queries

- Insert a **User** into the system at registration time
- Create a new **Pokémon** (admin only)
- Insert a **Pokémon** into a **Team**
- Create a new **Post**
- Create a new **Answer**
- Create a follow relationship
- Add a **Pokémon** to the favorites
- Retrieve **User** information at login time
- Retrieve a **User** by username when looking for a new friend
- Retrieve **Team** information based on user
- Retrieve **Pokémon** information using several filters
- Retrieve recommended **Users**
- Retrieve list of a **User**'s friends
- Retrieve a **Pokémon** by name when trying to catch it
- Retrieve all the **Posts** relative to a **Pokémon**
- Retrieve all the **Answers** to a **Post**

- Retrieve **User**'s favorite **Pokémon**
- Modify **User** settings (email, password, country)
- Update **Team**'s name
- Update **Team**'s points
- Update **Pokémon**'s catch rates Analytics: find % of **Users** that own that **Pokémon**
- Remove a **User** (admin only)
- Remove a **Pokémon** (admin only)
- Remove a **Post** (only admin and post's owner)
- Remove a follow relationship
- Remove a **Pokémon** from the favorite ones
- Analytics: ranking of most popular **Pokémon** in world/each country
- Analytics: ranking of best **Teams** in the world/each country/among friends
- Analytics: evolution on time of a **Pokémon** catch rate
- Analytics: evolution on time of number of logins per day/total **Users**/logins per day by country (admin only)

3 — Project

3.1 Adopted Databases

According to concept presented in the previous chapter we can make the following considerations:

1. Because of the performance constraint, a fast backend is required. Moreover, since the aim is to spread the application worldwide, the database infrastructure should be easy to distribute.
2. **Pokémon** must store heterogeneous data like URLs, different kinds of bios, float arrays and so on.
3. **Users** are divided into normal users and admins. Although the second ones are few, a denormalized approach could be better to handle the fact that these two categories have very different attributes.
4. Rankings are real-time OLAP queries: they need fast aggregation strategies.
5. Favorite Pokémon, Friends, Posts and Answers together form a real Social Network.
6. A **Team**, in a normalized relational model, could be seen as a relationship table between Users and Pokémon. Anyway, a huge table with a lot of duplicated PokémonID is not scalable due to the requirements of this application. There is a need to find the best way to perform quickly both the retrieving of a **User's** team and the ranking of the most used **Pokémon**, optimizing if possible memory consumption.

The points 1 to 4 guided the choice of a **Document Database** for handling User and Pokémon data. The flexibility, denormalization and performance of this kind of the database make it the most appropriate one.

The point 5 is best handled by a **Graph DB**, optimized for networks and different kinds of relationships. Moreover, we realized that the best way to handle a team is to decompose it in a set of Graph Relationships (USER – OWNS → POKEMON). Indeed, in this way queries mentioned at point 6 are very fast (just counting incoming/outcoming edges, see paragraph 3.3.1), and there are no useless, waste-memory repetitions of User IDs/Pokémon IDs.

Since each user can have only a team, team name and points are stored in the user collection.

3.2 Document Database

3.2.1 Queries Handled

- Insert a **User** into the system at registration time
- Create a new **Pokémon** (admin only)
- Retrieve **User** information at login time
- Retrieve a **User** by username
- Retrieve **Pokémon** information using several filters
- Retrieve a **Pokémon** by name when trying to catch it
- Modify **User** settings (email, password, country)
- Update **Team**'s name
- Update **Team**'s points
- Update **Pokémon**'s catch rates
- Remove a **User** (admin only)
- Remove a **Pokémon** (admin only)
- Analytics: ranking of best **Teams** in the world/each country/among friends
- Analytics: evolution on time of a **Pokémon** catch rate
- Analytics: evolution on time of number of logins per day/total **Users**/logins per day by country (admin only)

3.2.2 Entities handled

Document Database stores information about **Users** and **Pokèmons**.

In particular it remembers **User**'s anagraphics and login data, last login, remaining Pokéballs, team name and earned points; a Boolean field distinguish admin from normal users. Admins have no points nor team or Pokéballs.

In a separate collection are stored data about **Pokémon**: PokédexId (source: PokeAPI), characteristics, one or more types, bio, images URLs, current capture_rate and its last 30 catch_rates stored into an array of Floats.

The details of the collections are reported in the following paragraph.

3.2.3 Collections structure

<pre> 1 _id: ObjectId("5fd10b92b95ca407d0c0d726") 2 admin: false 3 username: "Caspar_Kolibius" 4 email: "Caspar.Kolibius@lsmdb.unipi.it" 5 password: "fd3b89c6e0a5a757d1a3ef9f12d99d53da55f72792881c4bd1fd16d91557089d" 6 surname: "Kolibus" 7 name: "Caspar" 8 birthDay: "2001-03-17T04:36:02.739Z" 9 country: "Israel" 10 teamName: "Team name" 11 lastLogin: "2020-12-09T18:38:24.970Z" 12 dailyPokeball: 10 13 points: 0 </pre>	<pre> ObjectId Boolean String String String String String String String String String Int32 Double </pre>
--	---

Figure 4: User Collection

Relevant Attributes:

- *Admin*: **true** → admin, **false** → normal user
- *Username*: unique mnemonic ID of the user
- *Email*: must respect typical e-mail format
- *Password*: encrypted version of the user-chosen password
- *Last Login*: timestamp of the last time the user logged into the application
- *dailyPokeball*: number of daily Pokéballs left. They are up to 10 per day
- *points*: worth of his/her team


```

1  _id: ObjectId("5fc257e9ae36f63454f88a4b")
2  id: 1
3  name: "bulbasaur"
4  weight: 69
5  height: 7
6  capture_rate: 45
7  biology: "A strange seed was planted on its back at birth. The plant sprouts and grows with this POKÉMON."
8  types: Array
9    0: "grass"
10   1: "poison"
11  portrait: "https://raw.githubusercontent.com/PokeAPI/sprites/master/sprites/pokemon/other/official-artwork/1.png"
12  sprite: "https://raw.githubusercontent.com/PokeAPI/sprites/master/sprites/pokemon/1.png"
13  capture_rates: Array
14    0: 44.9
15    1: 44.6

```

ObjectId
 Int32
 String
 Int32
 Int32
 Int32
 String
 Array
 String
 String
 String
 String
 Array
 Double
 Double

Figure 5: Pokèmon Collection

Relevant Attributes:

- *Id*: Pokédex ID (unique)
- *Name*: unique mnemonic ID of the Pokémon
- *Capture_Rate*: current index of probability to catch the Pokémon
- *Portrait/Sprite*: URLs of the graphical representations of this Pokémon
- *Capture_Rates*: array of the last 30 values of the capture_rate, one for each of the last 30 days.

3.2.4 Indexes

Username The first field in which we study the possibility of indexing is the *username* one in the **User** collection. A *username* is a REQUIRED and UNIQUE field of each **User**, and it is his/her mnemonic id inside the application. The field *username* is involved in the following queries:

Type	Query
W1	Insert a new username at registration time of an arbitrary user
W2	Remove a username when an admin delete's a user from the system
R1	Check uniqueness of a username at registration time
R2	Check user's credential at login time
R3	Find a user by username when a new follow request is submitted

Assuming that a registered user will play the game for about 100 days before “getting bored”, we can state that the number of logins-per-day will be 100 times the number of registrations-per-day: this means that the queries R1+R2 are submitted 101 times more than query W1. Moreover, we can assert that query W2 will be very rare, while R3 is a

popular query among the network structure of the application, say 30 times the number of registered users: we find out that read operations on this field are about 130 times the number of write operations. Now consider MongoDB performances with and without using an index on the username field, in a Database populated by 250k users.

```
1 > db.user.find({username:"eee"},{username:1}).explain("
    executionStats")
```

After submitting the previous command the following results are obtained.

```
"executionStats" : {
  "executionSuccess" : true,
  "nReturned" : 1,
  "executionTimeMillis" : 181,
  "totalKeysExamined" : 0,
  "totalDocsExamined" : 250464,
```

(a) Results without index

```
"executionStats" : {
  "executionSuccess" : true,
  "nReturned" : 1,
  "executionTimeMillis" : 2,
  "totalKeysExamined" : 1,
  "totalDocsExamined" : 1,
```

(b) Results with index

In the picture on the left is reported the output of the query when we do not use an index. Execution time is huge due to the very high number of docs examined. On the contrary, with an index, the same query need an execution time almost 100 times lower, and of course thanks to the index, DBMS only need to examine one document. Moreover the unique property permits to eliminate the need of submitting query R1 at each registration. Considering the very high speed-up ratio of the indexing and the high frequency of this kind of queries w.r.t. the write operations (as explained before), a UNIQUE INDEX on username has been created.

Country As seen before, starting from the application queries we demonstrate the benefits of an index in the field *country*.

Type	Query
W1	Insert the country data at registration timer
W2	Remove all the user's data if a user is banned by an admin
W3	Changing of settings after a user changes residence's country
R1	Rank all users by country
R2	Rank countries with the highest logins-per-day ratios

Let x be the number of registrations-per-day (W1), w.r.t this number W2 and W3 are very rare operations. Indeed, even though we can expect mischievous behaviors from some user, the number of country changes will never be comparable with x .

On the other hand, in order to guarantee a read-your-own-write eventual consistency on ranking R1, this query is recomputed every time a user asks to see the ranking itself. Thus, since the gameplay is highly based on rankings, we can estimate that R1 frequency will be about $400x$.

Furthermore we have to consider R2. Despite the fact that this query is executed just once per day (so $frequency(R2) \ll x$), it is an asynchronous procedure sensitive to execution time since it needs to lock the entire collection, make it unavailable to users for a while.

As seen before, let us compare DBMS performances with and without a country index.

```
1 > db.user.find({country:"Italy"}).explain("executionStats")
```

```
"executionStats" : {
  "executionSuccess" : true,
  "nReturned" : 989,
  "executionTimeMillis" : 291,
  "totalKeysExamined" : 0,
  "totalDocsExamined" : 250464,
```

(a) Results without index

```
"executionStats" : {
  "executionSuccess" : true,
  "nReturned" : 989,
  "executionTimeMillis" : 5,
  "totalKeysExamined" : 989,
  "totalDocsExamined" : 989,
```

(b) Results with index

Considering again about 250k users, without an index we need to scan the whole database, which means a medium-high execution time for each request.

On the contrary, we have a very high increase of performances introducing and index on country: execution time is about 58 times lower and the only documents examined are the ones that must be returned.

To summarize, considering the difference in frequency between reads and writes and the high decrease of execution time, an index on country has been introduced.

Pokemon Name Queries on Pokémon's name:

Type	Query
W1	Insert a new Pokémon into the Database
W2	Delete a Pokémon from the Database
R1	Search a Pokémon by name in the Pokédex
R2	Browse a Pokémon by name in Catch'Em'All in order to try to catch it
R3	Check name's uniqueness of each Pokémon when added to the database

Again, W1 and W2 are rare and admin-related operations: this means that this queries will not require a frequent update of the index. On the

contrary R1 and especially R2 are very frequent gameplay queries inside the application: we can estimate that R1+R2 frequency will be several orders of magnitude higher than W1+W2 one.

R3 instead is a query always required before W1, but it can be managed by DBMS adding a unique property to the index, thus reducing computational cost of the operation itself.

In terms of execution time, the final report is the following:

```
1 > db.user.pokemon({name: "pikachu"}).explain("
  executionStats")
```

```
"executionStats" : {
  "executionSuccess" : true,
  "nReturned" : 1,
  "executionTimeMillis" : 1,
  "totalKeysExamined" : 0,
  "totalDocsExamined" : 893,
```

(a) Results without index

```
"executionStats" : {
  "executionSuccess" : true,
  "nReturned" : 1,
  "executionTimeMillis" : 0,
  "totalKeysExamined" : 1,
  "totalDocsExamined" : 1,
```

(b) Results with index

Even if we have little changes on execution time due to the limited number of Pokémon, we can see how the index permits to decrease very much the number of examined documents.

For the reasons explained before and because of the very high ratio between reads and writes, we consider this little improvement enough relevant for the application purposes.

3.3 Graph Database

3.3.1 Queries Handled

3.3.2 Entities handled

3.3.3 Graph structure

3.3.4 Indexes

3.4 Redundancies and consistency management

3.5 Database Properties

3.5.1 Availability

3.5.2 Replicas

3.5.3 Eventual Consistency

3.5.4 Sharding

3.5.5 Pros and Drawbacks

3.6 Client, Server and Daemon Thread

3.7 Technologies and Frameworks

4 — Implementation

4.1 Package structure and information hiding

4.1.1 Packaging strategy and information hiding

4.1.2 UML package diagram

4.2 APIs and SPIs

4.3 Main tools

4.3.1 GSON

4.3.2 Caching mechanism and multimedia management

4.3.3 Password Encryptor

4.3.4 Logger

4.4 Analytics queries

4.4.1 User Rankings

4.4.2 Pokémon Rankings

4.4.3 Usage Statistics

4.4.4 Dynamic Catch Rate

4.5 Business logic

4.5.1 Points computing

4.5.2 Dynamic Catch Rate Computing

5 — Test

5.1 Privacy and Security

5.2 Unit Test

5.3 Robustness

5.4 Performance