

Evaluación del patrón *Polyglot Persistence* en una aplicación de mediano porte

Grupo 5: Santiago Acevedo, Donato Aguirre *Instituto de Computación*
Facultad de Ingeniería, Universidad de la República
Montevideo, Uruguay

Resumen

Este documento describe el proceso de implementación del proyecto final de la materia bases de datos no relacionales, donde se aborda la descripción de la solución propuesta, así como los principales desafíos encontrados.

I. INTRODUCCIÓN

Desde la tercera revolución de las bases de datos que comenzó en el 2005, se reconoce que una sola arquitectura de datos no es suficiente para cubrir la variedad de desafíos que trajo consigo la era digital moderna, dando nacimiento a diversas alternativas al modelo relacional.

En los últimos años, muchas empresas aclamadas por sus productos tecnológicos han llevado el lema "*the best tool for the job*" ("la mejor herramienta para la tarea en cuestión") aún más allá, promocionando la construcción de aplicaciones políglotas, esto es, que hacen uso de diferentes tecnologías de bases de datos en simultáneo.

En este trabajo se busca descubrir la implementación óptima para un sistema de pequeña a mediana escala y un cierto conjunto de requerimientos. En particular, se quiere verificar si la construcción de una aplicación donde algunos requerimientos se resuelvan utilizando una base de datos y otros otra, resulta en mejor rendimiento, ponderando también la experiencia del desarrollador. Para encontrar la combinación óptima, primero se implementa el sistema en diferentes versiones, utilizando una base de datos relacional (Postgres), una documental (Mongodb) y finalmente una de grafos (Neo4j).

Tras una fase de experimentación donde se comparan los tiempos de respuesta bajo estrés, los tiempos de carga de datos y la experiencia del desarrollador, se concluye que Neo4j destaca por la legibilidad del código de las consultas y su rendimiento en específicamente en uno de los requerimientos. Mongo destaca por su rendimiento en otro de los requerimientos. Postgres es consistentemente satisfactorio en todos los aspectos. Sin embargo, se concluye que para la escala y requerimientos estudiados, si el sistema se implementa con Postgres, no obtiene ganancias substanciales de la introducción de las otras dos que necesariamente justifique la complejidad agregada.

El documento se organiza en seis secciones, comenzando con la sección corriente, que corresponde a la introducción I. En la sección II se presentan los trabajos relacionados. En la sección III, se introduce la temática y los requerimientos de la aplicación. En la sección IV, se expande en las decisiones y detalles que conciernen a la implementación. En la sección V, se describen la fase de análisis presentando los resultados obtenidos. Finalmente, en la sección VI se recopilan las conclusiones y se plantea posible trabajo futuro.

II. TRABAJOS RELACIONADOS

Este trabajo expande en el patrón "*Polyglot Persistence*" que, según Martin Fowler, fue originalmente planteado por Scott Leberknight en 2008 [9] [3]. Scott se inspiró en el término "*Polyglot Programming*" [2], que expresa que las aplicaciones deben hacer uso de los diferentes lenguajes de programación para beneficiarse del hecho que diferentes lenguajes son apropiados para atacar diferentes problemas. La idea de "*Polyglot Persistence*" es la misma pero aplicada a los datos: Utilizar el almacenamiento de datos correcto en las partes de la aplicación que corresponda. En lugar de utilizar un solo tipo de almacenamiento y un solo tipo de consistencia, el patrón busca flexibilidad, incorporando "*NoSQL data stores*" para satisfacer requerimientos específicos. Con la aparición de micro-servicios el término se popularizó aún más, remarcando que cada micro-servicio puede utilizar su propia tecnología de almacenamiento, de manera transparente para sus consumidores.

En su planteo original el patrón no define una única forma de implementación, al momento ya habiéndose sido aplicado de diferentes maneras por las grandes empresas. El artículo [6] es uno de los varios donde se presentan diferentes formas de implementación del patrón ya sea en sistemas distribuidos o dentro de la misma aplicación.

Sin embargo, su utilización no viene sin sus varias desventajas, que van desde la pérdida de la consistencia y atomicidad en algunos casos, hasta el aumento de la complejidad y costo del mantenimiento del sistema, por lo que la decisión de introducirse debe ser consciente.

En los últimos años, muchos libros y artículos de empresas tecnológicas reconocidas han promocionado sus experiencias implementando este patrón. Dado que la mayoría de los proyectos en la práctica no tienen esta escala, no es fácil determinar la línea de cuándo amerita la introducción de estas técnicas.

Se han llevado a cabo algunos experimentos y trabajos que comparan soluciones en sistemas políglotas, como es el ejemplo del artículo encontrado en [7]. En él se hace una comparación entre un sistema clásico implementado en ArangoDB versus un enfoque políglota construido en base a Neo4j y Mongo, concluyendo que el sistema políglota presenta mejor rendimiento para una cantidad de datos grande.

La motivación de este proyecto es profundizar en el caso de una aplicación estándar de mediano porte, considerado a partir del momento donde los datos almacenados ya no quepan confortablemente en memoria. Se pone énfasis en los atributos de calidad rendimiento, escalabilidad y mantenibilidad del código de las consultas, ignorando enteramente la consistencia e integridad.

A su vez, varias publicaciones comparan el rendimiento de ciertas consultas entre bases de datos. Por ejemplo, el artículo de la página de Neo4j [11] cita que la popular consulta 'amigos de amigos de amigos' es 180 veces más performante en Neo4j que en Postgres. En este estudio se incorpora un requerimiento de este tipo con la hipótesis de que la adición de Neo4j para resolverlo será conveniente.

III. DESCRIPCIÓN DE LOS REQUERIMIENTOS DE LA APLICACIÓN

La temática de la aplicación gira entorno al mundo empresarial, con la misión de brindar, tanto al emprendedor, como al inversor, herramientas de análisis para apoyar sus decisiones. Se relevaron tres requerimientos funcionales que la plataforma debería cumplir para prototipar una solución. Estos son:

- Visualización del perfil de una *startup*: Dado el nombre de una empresa, obtener sus datos básicos y la lista de sus rondas de inversión, con los datos de los inversores que participaron en cada una.
- Recomendar potenciales inversores para una *startup*: Dada una ciudad y un rubro de empresa, listar los veinte inversores que han invertido la mayor cantidad de veces en esa ciudad y rubro, ordenados por la cantidad de inversiones de forma descendente.
- Recomendar nuevas inversiones a un inversor: Dado el *id* de un inversor, listar veinte recomendaciones para su próxima inversión, esto es, las empresas con la mayor cantidad de inversiones hechas por sus co-inversores.

IV. DESCRIPCIÓN DE LA SOLUCIÓN

IV-A. Conjuntos de datos utilizados

La fuente más popular y fiable para obtener información empresarial sobre compañías privadas y públicas es Crunchbase [1]. Esta plataforma permite a sus usuarios *enterprise* exportar toda la información actualizada en formato CSV. Sin embargo, para minimizar costos, en el proyecto se utilizó un *dump* de Crunchbase del año 2013 que tiene el mismo formato pero fue extraído de la plataforma Kaggle [5] de manera gratuita. Dicho *dump* se encuentra disponible en [4].

El *dump* contiene once archivos en formato CSV, de los cuales solo tres se determinaron necesarios para resolver los requerimientos relevados. Una descripción más detallada del formato de los archivos se aprecia en el anexo en la sección VII-A. El tamaño aproximado de los archivos del *dataset* es el siguiente:

Archivo	Tamaño MB	Cantidad filas	Cantidad columnas
Objects	285	463k	40
Funding Rounds	13	53k	23
Investments	6	81k	6

Dado a que el *dump* contiene algunas filas corruptas, esto es, *foreign keys* a filas inexistentes en los otros archivos, fue necesario realizar una sanitización previa. Para esto, se realizó un pequeño *script* en python para eliminar estas filas. El tamaño del *dataset* resultante siguió siendo representativo dado que los datos corruptos eran menos de mil en más de 400 mil registros.

IV-B. Implementación de la aplicación

Para alcanzar los objetivos del proyecto no fue necesaria la construcción de una interfaz de usuario. Se puso énfasis en la implementación del lado *back-end* del sistema. Se construyó una API en NodeJS, utilizando los *drivers* para Javascript de cada una de las tecnologías de bases de datos, sin la ayuda de un *ORM*.

IV-C. Modelado de las bases de datos

En esta sección se describe el modelado de la solución en cada paradigma junto con su justificación.

En los tres modelos se determinó necesario soportar que un inversor, ya sea persona, organización financiera o empresa, invierta en una ronda de inversión que puede pertenecer a una empresa, producto de una empresa o entidad financiera.

De las empresas, productos y entidades financieras interesa principalmente su nombre, categoría, ciudad, país, dirección, página web, descripción, logo, entre otros. De las personas se conoce su nombre, ciudad, página web, etc.

De las rondas de inversión se conoce su fecha, para qué empresa, producto o entidad financiera recaudan, el tipo de ronda (Serie A, B, C, etc.), y el total recaudado. Se conoce qué inversores, ya sea individuo, entidad financiera o empresa, participaron en una ronda. Lamentablemente Crunchbase no especifica la cantidad que invirtió cada inversor específico.

IV-C1. Modelado relacional: En el modelado relacional se representaron tres entidades que se traducen a tres tablas de Postgres. Estas son `Objects`, `FundingRounds` e `Investments`.

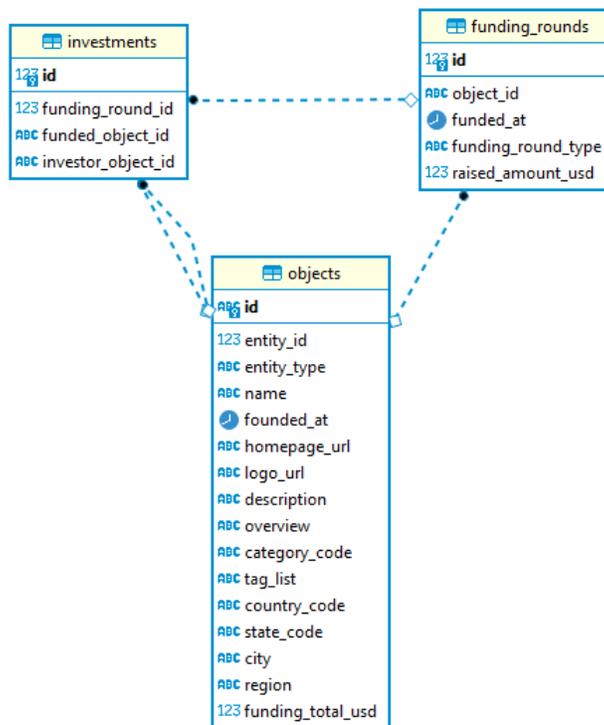


Figura 1: Modelo Postgres

La tabla `Objects` incluye una columna `entity_type` que puede tomar el valor de "Person", "Company", "Product" o "FinancialOrg". Esto es para soportar organismos que invierten y reciben inversiones al mismo tiempo. También permite distinguir inversiones en una empresa o en un producto específico de una empresa. Dado a que los diferentes tipos de `Objects` comparten la mayoría de los atributos, la cantidad de `null` en la tabla no se considera un problema.

`Investments` conecta una `FundingRound` con un `Object` que representa el inversor, conformando una relación N-N.

Un aspecto a notar es la de-normalización en la tabla `Investments`, donde se guarda el `investor_object_id` para poder accederlo directamente, evitando tener que hacer `join` con la tabla de `FundingRounds` para obtenerlo cada vez. La

aplicación de esta técnica tuvo un impacto positivo en el rendimiento de las tres consultas.

Con el mismo fin de optimizar las consultas, se indexaron todas las claves primarias y foráneas, así como otras columnas por las que se busca como *objects.name*. La columna *objects.name* es de tipo *CITEXT* para permitir la búsqueda ignorando el *case* por defecto, pero pudiendo hacer el uso del índice. También se colocó un índice compuesto por *objects.category_code* y *objects.city* ya que la segunda consulta busca por las dos al mismo tiempo.

IV-C2. Modelado documental: Al momento del diseño de la base Mongo se priorizó la primer consulta, esta es, la visualización del perfil de una empresa. Para eso se decidió modelar una única colección, *objects*, cuyos documentos se correspondan con el *JSON* final a retornar al buscar una empresa. Se hizo uso del patrón *Embedded Document* donde se definieron las *funding rounds* y los datos de sus *investors* dentro de cada *object*. Se siguió la recomendación en la documentación oficial de mongo, que se puede encontrar en [10].

Un ejemplo de un documento de la colección figura en el listado 1.

Listado 1 Ejemplo de documento Mongo

```

1  {
2  "id": "c:1",
3  "name": "Wetpaint",
4  "foundedAt": "2005-10-16T02:00:00.000Z",
5  "homepageURL": "http://wetpaint-inc.com",
6  "logoURL": "http://s3.amazonaws.com/crunchbase_prod_assets/assets/images/resized/0000/3601/3601v1-max
   -250x250.png",
7  "description": "Technology Platform Company",
8  "overview": "Wetpaint is a technology platform company that uses its proprietary state-of-the-art...",
9  "categoryCode": "web",
10 "tagList": "wiki, seattle, elowitz, media-industry, media-platform, social-distribution-system",
11 "countryCode": "USA",
12 "stateCode": "WA",
13 "city": "Seattle",
14 "region": "Seattle",
15 "fundingTotalUsd": "39750000",
16 "fundingRounds": [
17   {
18     "id": 888,
19     "type": "series-a",
20     "raisedAmountUSD": "5250000",
21     "fundedAt": "2005-09-30T03:00:00.000Z",
22     "investors": [
23       {
24         "id": "f:430",
25         "name": "Frazier Technology Ventures",
26         "type": "FinancialOrg"
27       }
28     ]
29   }
30 ]
31 }
```

En este modelo también se colocaron índices por las claves buscadas, como *objects.id* y *objects.name*. Por el mismo motivo que en el caso de Postgres, se colocó un índice compuesto por *objects.category_code* y *objects.city*. También fue necesario indexar por *fundingrounds.investors.id* para permitir que los *lookups* performen bajo un tiempo aceptable.

IV-C3. Modelado de grafos: Para la base Neo4j, se decidió modelar *city* y *category* como nodos separados para poder sacarle provecho a la segunda consulta, donde se sugieren inversores para un emprendimiento considerando su historial de inversión en la ciudad y categoría. Sin embargo, como se describirá en la sección de experimentación, esta decisión pudo haber tenido un pequeño impacto negativo en el rendimiento a causa de perder la oportunidad de utilizar un índice compuesto para encontrar objetos dada una ciudad y categoría.

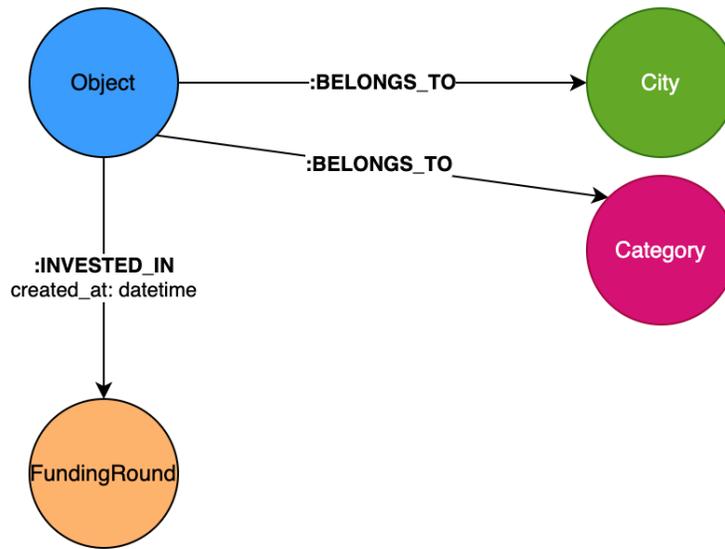


Figura 2: Modelo Neo4j

Dado a que el *dataset* cuenta con filas de *Objects* para los que no se conoce su ciudad y categoría, se decidió modelar un nodo de ciudad y un nodo de categoría que tuviesen el nombre *Unknown*. De esta manera, todos los *Objects* tienen una relación con una ciudad y una categoría y se permite que si se desea, se localicen rápidamente los *Objects* para los cuales no se conoce su ciudad y categoría.

Los índices creados fueron por *object.id*, *object.name*, *city.name*, *category.name* y *fundinground.id*.

IV-D. Carga de Datos

En esta sección se describe el procedimiento seguido para la carga de datos. En el caso de Postgres y Mongo, se implementaron y compararon dos enfoques distintos para cada una, en la búsqueda de un mecanismo performante y escalable. Todos los procedimientos se encuentran automatizados y versionados dentro del repositorio.

IV-D1. Base Postgres: Como se mencionó anteriormente, para la carga de la base de datos Postgres se probaron dos enfoques distintos.

El primero es a fuerza bruta, realizando una operación `INSERT` para cada fila de cada CSV, el cual presentó tiempos de ejecución de más de 15 minutos.

El segundo enfoque es aquel utilizado en la práctica y recomendado por la comunidad, que implica el uso de Streams para la lectura e inserción. Esta técnica permite consumir el archivo CSV de a tandas de filas, como un flujo continuo, evitando cargar todo el contenido del archivo a memoria de una sola vez. A medida que se lee las tandas, se las procesa e inserta a la base en *bulk*, resultando mucho más performante y escalable. Con este enfoque, la carga de las tres tablas finaliza en menos de un minuto.

IV-D2. Base Mongo: En el caso de Mongo, también se experimentó con dos enfoques diferentes. Ambos utilizan la base Postgres como la fuente original, como se realiza normalmente en la práctica.

El primer enfoque itera sobre cada fila de la tabla *Objects*. Para cada fila se realiza una consulta para formar la estructura final del documento a insertar en Mongo, esto es, un JSON con los datos de la empresa con sus rondas e inversores. Luego se realiza una consulta a la base Mongo para insertar el documento, y se procede a repetir para las demás filas. Este algoritmo presentó tiempos de ejecución de hasta 5 minutos.

El segundo enfoque se basa en primero exportar desde Postgres un JSON conformado por un array con los documentos a insertar. Esto se hizo mediante la consulta 22 que figura en el anexo. Luego se utiliza la función `COPY` de Postgres para guardar ese JSON en un archivo temporal, como se aprecia en el siguiente código:

Listado 2 Export Postgres a Mongo

```

1 COPY (
2     SELECT array_to_json(array_agg(results))
3     FROM (
4         ${sqlQuery}
5     ) results
6 ) TO '${path}' WITH (FORMAT text, HEADER FALSE)`

```

Luego, se utiliza el comando `mongoimport` para insertar cada elemento del array del archivo en la nueva colección, ejemplificado a continuación. Este algoritmo finalizó en 33 segundos, donde 11 segundos corresponden a la exportación de los datos de Postgres al archivo JSON.

Listado 3 Export Postgres a Mongo

```

1 mongoimport --uri mongodb://localhost:27017/bdnr --collection objects --jsonArray data/objects.json

```

IV-D3. Base Neo4j: El *script* de carga de datos de Neo4j asume que se corrió previamente el *script* para Mongo, ya que utiliza el mismo archivo JSON que se exportó desde Postgres. La carga de los datos del *array* se hizo mediante la consulta encontrada en el Anexo en 22 la cual utiliza la función `apoc.load.json`. Este algoritmo finalizó en 41 segundos.

IV-E. Descripción de los endpoints

En esta sección se describe la especificación de los *endpoints* implementados.

El primer *endpoint* se utiliza para obtener el perfil de una compañía y refiere al primer requerimiento descrito en la sección III. Este *endpoint* se utiliza mediante un *GET*, y requiere el comienzo del nombre de una empresa por la que se busca en los parámetros de la *request*. La búsqueda es indiferente a mayúsculas y minúsculas.

La firma de la consulta es la siguiente: `GET /company?name=companyName`.

Un ejemplo de petición y respuesta es:

Listado 4 Request Comapany

```

1 req: {
2     "method": "GET",
3     "url": "/v2/company?name=Fitbit",
4     "query": {
5         "name": "Fitbit"
6     },
7 }

```

Listado 5 Response Comapany

```

1 {
2     "company": {
3         "_id": "62bcbebe756234fc61bbf926",
4         "id": "c:10015",
5         "name": "Fitbit",
6         "founded_at": "2007-10-01",
7         "homepage_url": "http://www.fitbit.com",
8         "logo_url": "http://s3.amazonaws.com/crunchbase_prod_assets/assets/images/resized/0002/4083/24083v4-max-250x250.png",
9         "description": null,
10        "overview": "Fitbit inspires people to exercise more, eat better and live healthier lifestyles..",
11        "category_code": "health",
12        "country_code": "USA",
13        "state_code": "CA",
14        "city": "San Francisco",
15        "region": "SF Bay",
16        "funding_total_usd": 68069200,
17        "entity_type": "Company",
18        "funding_rounds": [
19            {
20                "id": 3619,

```

```

21     "raised_amount_usd": 2000000,
22     "funded_at": "2008-10-10",
23     "type": "series-a",
24     "investors": [
25         {
26             "id": "f:74",
27             "name": "True Ventures",
28             "entity_type": "FinancialOrg"
29         },
30         {
31             "id": "f:737",
32             "name": "SoftTech VC",
33             "entity_type": "FinancialOrg"
34         }
35     ]
36 },
37 ...
38
39 ]
40 }
41 }

```

El segundo *endpoint* corresponde al segundo requerimiento y se utiliza para recomendar potenciales inversores para una *startup* dado una categoría y una ciudad. Se debe introducir el nombre completo de la categoría y la ciudad por la que se filtra.

La firma es la siguiente: GET /investors?category=categoryName&city=cityName.

Un ejemplo de petición y respuesta es:

Listado 6 Request Investors

```

1 req: {
2     "method": "GET",
3     "url": "/v2/investors?category=software&city=London",
4     "query": {
5         "category": "software",
6         "city": "London"
7     }
8 }

```

Listado 7 Response Investors

```

1 {
2     "investors": [
3         {
4             "id": "f:2469",
5             "name": "Seraphim Capital",
6             "entity_type": "FinancialOrg",
7             "investments_count": 5
8         },
9         {
10            "id": "f:2485",
11            "name": "YFM Venture Finance",
12            "entity_type": "FinancialOrg",
13            "investments_count": 5
14        },
15        {
16            "id": "f:998",
17            "name": "DFJ Esprit",
18            "entity_type": "FinancialOrg",
19            "investments_count": 5
20        },
21        ...
22    ]
23 }

```

El tercer *endpoint* corresponde al tercer requerimiento y se utiliza para recomendar nuevas inversiones a un inversor, utilizando el identificador del inversor.

La firma de la consulta es la siguiente: GET /investmentRecommendations?investorId=investorId.

Un ejemplo de petición y respuesta es:

Listado 8 Request Investments Recommendations

```

1 req: {
2   "method": "GET",
3   "url": "/v3/investmentRecommendations?investorId=f%3A1294",
4   "query": {
5     "investorId": "f:1294"
6   },
7 }

```

Listado 9 Response Investments Recommendations

```

1 {
2   "recommendations": [
3     {
4       "id": "c:135515",
5       "name": "Romotive",
6       "matches_count": 4
7     },
8     {
9       "id": "c:29052",
10      "name": "Graphicly",
11      "matches_count": 4
12     },
13     {
14      "id": "c:4060",
15      "name": "Univa UD",
16      "matches_count": 4
17     },
18     ...
19   ]
20 }

```

IV-F. Implementación de las consultas

Para la implementación de las consultas se iteró para encontrar soluciones lo más optimizadas posibles, experimentando también en cuanto a los índices que se colocaron. Los resultados se validaron a mano para algunas entradas sencillas y se compararon entre las tecnologías para validar que retornaran lo mismo. A continuación, se describe brevemente la implementación de cada consulta a alto nivel. En los listados 13, 16 y 19 se encuentra el código completo de las consultas de Postgres, Mongo y Neo4j respectivamente.

IV-F1. Consulta 1: Para resolver el primer requerimiento, la API requiere retornar un JSON con los datos la empresa buscada por el comienzo del nombre, y debe incluir los datos de sus asociaciones de manera anidada.

En Mongo es trivial, debido a que el documento ya contiene la estructura final del *JSON* a retornar, alcanzando con un `findOne` con un filtro sobre el atributo `name` utilizando una expresión regular.

En el caso de Postgres y Neo4j resultó más fácil armar el *JSON* directamente desde la consulta que implementar código de la aplicación para construirlo. Si no fuese por esto las consultas hubiesen sido más directas, utilizando `LEFT JOIN` en Postgres, y pudiendo evitar usar `WITH` en Neo4j.

En Postgres se implementó utilizando sub-consultas para las asociaciones, y la función `array_agg` para convertir el resultado de las sub-consultas a un array para que el resultado final tenga el formato deseado.

En Neo4j se realiza primero un `MATCH` buscando nodos de tipo `Object` por el nombre dado, seguido de múltiples `OPTIONAL MATCH` para encontrar los nodos correspondientes a la ciudad, categoría, rondas de inversión e inversores. Es necesario `OPTIONAL MATCH` en el caso de las rondas de inversión e inversores, para asegurarse que empresas que no tengan estas relaciones sean retornadas de todas formas. Luego la consulta se delimita en partes separadas con `WITH`, para ir formando la estructura final deseada desde adentro hacia afuera. Esto es, primero definir mediante alias (`AS`) la estructura de cada `investor` de una `funding round` y luego, la estructura de la `funding round` de las empresas seleccionadas.

IV-F2. Consulta 2: Para la segunda consulta, que recomienda posibles inversores, se busca retornar los 20 inversores con la mayor cantidad de inversiones que pertenecen a una cierta categoría y ciudad.

En Postgres se resuelve fácilmente haciendo uso de `GROUP BY` por el ID del inversor, y `COUNT`, luego de haber localizado las inversiones en las empresas relevantes.

En Mongo, se implementa con *Aggregation Pipeline*. Primero se utiliza un `match` para localizar los documentos de empresas por categoría y nombre. Al igual que en Postgres, se verificó que internamente se estuviese haciendo uso del índice compuesto por estas claves. Luego es necesario hacer `unwind` de las `funding_rounds` y una vez más de sus respectivos `investors`. Hasta este punto se cuenta con un objeto por inversión, lo que permite que sean agrupados por ID del inversor utilizando `group` y contarlos utilizando `sum: 1` a medida que se agrupan.

En Neo4j, la consulta se resuelve con un solo gran `MATCH`, que busca los nodos de la categoría y ciudad dados, y navega hasta localizar todos los inversores que han invertido en empresas que pertenezcan. Luego se utiliza `WITH` para obtener un campo computado con `COUNT`. Cabe acotar que la razón por la que se puede hacer `MATCH`, en lugar de `OPTIONAL MATCH`, es que todos los objetos tienen una relación con un nodo de categoría y ciudad. Los `nulls` del *dataset* en categoría y ciudad se representaron como dos nodos con nombre `Unknown`.

Para las tres bases, los pasos finales de las consultas son muy similares, incluyendo ordenar el conjunto resultante y limitar la cantidad a retornar.

IV-F3. Consulta 3: En la tercer consulta, que recomienda nuevas inversiones, se busca retornar las 20 empresas o productos con mayor cantidad de inversiones hechas por co-inversores del inversor dado.

En Postgres, se resuelve haciendo una sub-consulta para encontrar los co-inversores, para luego obtener sus inversiones, agruparlas por empresa y contarlas. Finalmente se descartan las empresas en las que el inversor dado ya ha invertido, se ordena el conjunto resultante y se limita la cantidad a retornar.

En Mongo, se comenzó haciendo un `match` por `funding_rounds.investors.id` del ID dado, para buscar los `Objects` donde el inversor ha invertido. A continuación, se hace `unwind` de `funding_rounds` y una vez más de `funding_rounds.investors`. Hasta este punto se cuenta con un elemento por cada inversión hecha por los co-inversores del inversor dado, pero solamente en empresas donde el inversor ya ha invertido. Es por eso que luego se realiza una proyección para quedarse solamente con los IDs de los co-inversores. Luego se agrupa para obtener solo IDs únicos para hacer el siguiente paso más barato. A continuación, se realiza un `lookup` para encontrar todas las inversiones hechas por los co-inversoras. Más adelante se realiza un `group` por ID de la empresa, sumando uno en `matches_count` por cada elemento de una empresa que se agrupa, y utilizando `$addToSet` para mantener el array de IDs de inversores que participaron cada empresa. Este array luego se utiliza para remover las empresas donde ya se encuentre el inversor dado. Finalmente se ordena, se limita y se limpian los campos retornados.

En Neo4j, esta consulta requirió varias iteraciones y esfuerzo para lograr escribirla de una forma que fuese performante. La primera versión del código de la consulta era extremadamente sencilla y más legible que la actual, pero su rendimiento era diez veces peor. La versión actual realiza primero un `MATCH` para encontrar todas las empresas en las que el inversor con el ID dado invirtió. A continuación se busca la unicidad de los datos resultantes, por lo que hace uso de `COLLECT(distinct mycompanies)`, y precisando un `UNWIND` a continuación ya que el atributo computado es de tipo array. Este paso fue una de las adiciones de la nueva versión que tuvo significativas mejoras en el rendimiento. En un segundo `MATCH` se encuentran todos los inversores que invirtieron en estas empresas (co-inversores). Nuevamente se aplica `COLLECT(distinct coinvestors)` y `UNWIND` para obtener los co-inversores una sola vez cada uno y optimizar los siguientes pasos. En un tercer `MATCH` se buscan todas las empresas en las que los co-inversores invirtieron. Luego se quitan las empresas en las que el inversor dado ya invirtió. Luego se utiliza `WITH` para poder agrupar por ID del inversor y así contar cuántos co-inversores han invertido en cada una. Finalmente, se ordena el conjunto resultante y se limita la cantidad a retornar.

V. EXPERIMENTACIÓN

En esta sección se comparan las versiones en tres aspectos diferentes: Tiempo de carga de datos, rendimiento de las consultas y experiencia del desarrollador. Todas las pruebas fueron realizadas en el mismo hardware, bajo las mismas situaciones y con los mismos recursos asignados.

V-A. Hardware y versiones utilizadas

A continuación se especifica el hardware y versiones del software utilizados durante la experimentación.

Hardware y sistema operativo:

- Sistema Operativo: MacOS Monterrey
- Procesador: 3.1GHz Intel Core i7 4 núcleos
- RAM: 16GB 2133 Mhz LPDDR3

Versiones de herramientas:

- PostgreSQL: 14.3
- Mongo: v5.0.7
- Neo4j: 4.4.8
- Node: 16.14.2

V-B. Comparación de tiempos de carga de datos

En primer lugar, se compararon los distintos tiempos de carga de la base de datos descritos en la sección IV-D. Para esta etapa solo se compararon las soluciones finales ya optimizadas de los *scripts* de carga de datos de cada base. En otras palabras, no se consideraron los algoritmos a fuerza bruta que se implementaron inicialmente.

Los tiempos de ejecución comprenden únicamente los pasos relacionados a la carga de datos una vez que la base ya se encuentra vacía. En otras palabras, se descontó el tiempo de eliminado de los datos existentes, y, en el caso de Mongo y Neo4j, la generación del archivo JSON a partir de Postgres. Un mayor detalle del desglose de los tiempos de cada *script* se pueden observar en el Anexo en la sección VII-B.

Base de Datos	Tiempo de carga de datos
Postgres	36s
Mongo	21s
Neo4j	41s

Cuadro I: Comparación de tiempo de ejecución de carga de datos

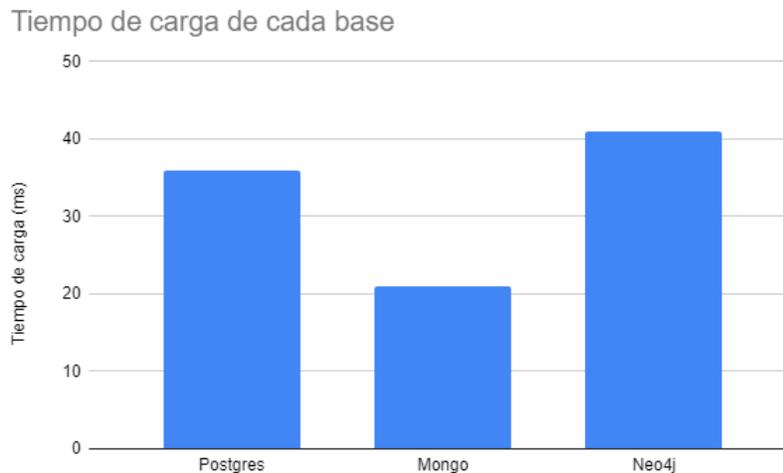


Figura 3: Comparación tiempos de carga

Como conclusión, vemos que fue posible construir un *script* de carga de datos escalable para cada una de las bases de datos, donde para el *dataset* de 300MB finalice por debajo de los 40 segundos. Se observó que recrear la base de Mongo a partir de un archivo *JSON* resultó un 50% más rápido que en Neo4J con el mismo archivo de entrada. Sin embargo, esto seguramente se debe al hecho de que la inserción en Mongo es directa (un elemento del *array* corresponde a un documento). En Neo4j, cada elemento del *array* implica varios pasos, donde se debe buscar o crear los nodos correspondientes a la ciudad, categoría,

rondas de inversión e inversores, resultando en una gran cantidad de nodos y relaciones, y consultas para buscarlos.

V-C. Comparación de tiempos de respuesta bajo carga

A continuación, se comparó la *performance* de las tres versiones de los tres *endpoints* de la aplicación. Las pruebas de *performance* se realizaron con la herramienta *Open Source K6* [8]. Los *scripts* se encuentran en el repositorio en la sección [12].

Para la selección de los parámetros de carga, se determinó arbitrariamente una duración de 30 segundos. Se lanzaron 10 hilos concurrentes, cada uno realizando peticiones sin descanso entre ellas, esperando a que una complete para lanzar la siguiente, todas las veces que sea posible hasta agotar el tiempo de la prueba. Las pruebas se corrieron por separado para que no se interfirieran entre ellas, de manera local en el *hardware* especificado anteriormente. Las tablas que resumen los resultados de los tiempos de respuesta de cada *endpoint* son las siguientes:

	Mediana	Throughput
Postgres	8,14ms	1167,6req/s
Mongo	5,8ms	1122,1req/s
Neo4j	9,9ms	793,0req/s

Cuadro II: Comparación de rendimiento del endpoint 1

	Mediana	Throughput
Postgres	6,1ms	1000,0req/s
Mongo	5,6ms	996,6req/s
Neo4j	11,2ms	648,7req/s

Cuadro III: Comparación de rendimiento del endpoint 2

	Mediana	Throughput
Postgres	80,0ms	98,8req/s
Mongo	196,9ms	41,1req/s
Neo4j	62,82ms	157,2req/s

Cuadro IV: Comparación de rendimiento del endpoint 3

De los resultados se confirmaron algunas hipótesis y se invalidaron otras. En lo que concierne al primer *endpoint*, todos los resultados se encontraron en el mismo orden de magnitud, con Mongo resultando un 40 % y 47 % más performante que Neo4j y Postgres respectivamente. Considerando que el modelado de Mongo fue optimizado para este propósito, se esperaba ver una diferencia incluso más notoria. Es por eso que se repitió la prueba con hasta 500 hilos concurrentes para descartar una carga poco desafiante. Sin embargo, los resultados se mantuvieron proporcionales. Se concluye que, o la carga es aún insuficiente, el *dataset* no tiene el tamaño suficiente, o la consulta de Postgres no es lo suficientemente costosa como para ameritar la introducción de Mongo al sistema únicamente por este requerimiento.

Con respecto al segundo *endpoint*, las tres versiones se consideran satisfactorias, con un promedio similar a los tiempos vistos en la primera consulta. Cabe destacar que Mongo obtuvo resultados muy similares a Postgres incluso cuando el código de la consulta de Mongo es, a primera vista, mucho más compleja. Neo4j performó casi un 40 % peor que las otras dos, aunque debido a la naturaleza de la consulta se esperaba lo opuesto. Sin embargo, se sospecha que la decisión de modelar la categoría y ciudad como nodos independientes en Neo4j puede haber tenido un impacto negativo, dado que se pierde la posibilidad de tener un índice compuesto como se hizo en Postgres y Mongo.

Finalmente en cuanto el tercer *endpoint*, se observaron tiempos mucho mayores a las primeras dos consultas. Esto es de esperar, dado que la consulta requiere la manipulación de muchísimos datos (todas las inversiones de todos los co-inversores de todas las rondas de inversión de un inversor dado). En este caso, Neo4j destacó, triplicando el *throughput* de Mongo y un 50 % mejor que Postgres.

V-D. Comparación de experiencia del desarrollador

Por último, se compara la facilidad de desarrollo para las tres consultas. Cabe mencionar que si bien la mantenibilidad del código no es el foco principal del estudio, en la práctica este atributo tiene un alto impacto en el costo del sistema a largo plazo, por que lo que es igual, o en algunos casos hasta más, importante que el rendimiento.

Para evaluar este atributo se le adjuntó una nota, del 1 al 5, a cada consulta realizada. Este puntaje busca denotar la complejidad cognitiva y el tiempo dedicado a construir y para entender el código de la consulta, donde 5 significa Muy Accesible y un 1 significa Muy Complejo.

A continuación se presenta una tabla con los puntajes asignados y una breve explicación de cada uno de estos.

	Consulta 1	Consulta 2	Consulta 3
Postgres	3	3	3
Mongo	5	2	1
Neo4j	3	5	4

Cuadro V: Comparación de facilidad de implementación de las consultas

La primera consulta en Mongo es trivial, debido a que el documento ya contiene la estructura final del *JSON* a retornar, alcanzando con un *findOne*. En cambio, la complejidad de la consulta en Postgres y Neo4j se vio afectada negativamente por el hecho de que la API requiere retornar un *JSON* con los datos de las asociaciones anidadas, y resultó más fácil obtenerlo directamente desde la consulta que implementar código de la aplicación para construirlo. Si no fuese por esto las consultas hubiesen sido más directas.

En lo que concierne a la segunda y tercera consulta, llegar a una primera versión que funcione en Mongo fue donde se requirió más tiempo. Los *aggregation pipelines* resultantes son extensos, sobretudo en la tercera consulta, requiriendo la adición de comentarios en cada paso para hacer acuerdo de la estrategia tomada y de qué representan los resultados intermedios.

En general, las consultas en Neo4j son más legibles y cortas, debido a la naturaleza de la sintaxis de Neo4j y las decisiones de modelado. Sin embargo, cabe destacar que al implementarlas inicialmente se cometieron errores que solo se identificaron al comparar los resultados con las otras versiones. Por ejemplo, el optimizar la tercer consulta para llegar al *throughput* descrito no fue tarea fácil y afectó negativamente a la legibilidad de la consulta, como se desarrolló en la sección IV-F.

VI. CONCLUSIONES Y TRABAJO FUTURO

El principal objetivo que se planteó inicialmente es el de descubrir la implementación más conveniente para un sistema de mediana escala y con cierto conjunto de requerimientos relevados. Para ello el sistema se implementó utilizando tres paradigmas de bases de datos diferentes: Relacional con Postgres, documental con Mongo, y de grafos con Neo4j. Para cada uno se experimentó el modelar la solución, automatizar el proceso de carga de datos utilizando diferentes estrategias, e implementar y optimizar consultas con diferentes complejidades. También se destacó la experiencia de desarrollar pruebas de carga, así como de implementar una API utilizando los *drivers* de las tres tecnologías, sin la ayuda de un *ORM*.

En lo que respecta a las conclusiones del análisis comparativo, Neo4j presentó mejores resultados en cuanto a facilidad de codificación de las consultas, pero significativamente peor rendimiento en todos las consultas menos la última, y Mongo exactamente lo contrario. Postgres presentó resultados consistentes y satisfactorios en todos los puntos evaluados. Si la prioridad se pone puramente en el rendimiento, se debe concluir que la solución final más performante consiste de un sistema políglota donde el primer y segundo *endpoint* son servidos desde Mongo, y el tercero desde Neo4j. Sin embargo, en la práctica no necesariamente se justifica la complejidad de mantener una aplicación políglota, considerando que Postgres resultó casi tan bueno como la mejor opción en todos los aspectos. Es por eso que se concluye que, en la práctica, este sistema se podría implementar utilizando solamente Postgres. En un futuro, si la escala de datos los requiere, se podría introducir Mongo y/o Neo4j.

Se logró evidenciar la popularidad de la elección de utilizar Postgres en los sistemas para *startups*. Postgres demostró ser una tecnología fácil de configurar, consultar y consistentemente performante para escalas como la de este estudio.

Trabajos a futuros pueden incluir el experimentar con escenarios más agresivos en las pruebas de carga, con un orden mayor de datos almacenados, o con otro conjunto de requerimientos para ver en qué punto Postgres deja de ser una buena opción.

VII. ANEXO

VII-A. Columnas de los archivos CSV

```

1 id
2 entity_type
3 entity_id
4 parent_id
5 name
6 normalized_name
7 permalink
8 category_code
9 status
10 founded_at
11 closed_at
12 domain
13 homepage_url
14 twitter_username
15 logo_url
16 logo_width
17 logo_height
18 short_description
19 description
20 overview
21 tag_list
22 country_code
23 state_code
24 city
25 region
26 first_investment_at
27 last_investment_at
28 investment_rounds
29 invested_companies
30 first_funding_at
31 last_funding_at
32 funding_rounds
33 funding_total_usd
34 first_milestone_at
35 last_milestone_at
36 milestones
37 relationships
38 created_by
39 created_at
40 updated_at

```

Listado 10 Estructura archivo Objects

```

1 id
2 funding_round_id
3 object_id
4 funded_at
5 funding_round_type
6 funding_round_code
7 raised_amount_usd
8 raised_amount
9 raised_currency_code
10 pre_money_valuation_usd
11 pre_money_valuation
12 pre_money_currency_code
13 post_money_valuation_usd
14 post_money_valuation
15 post_money_currency_code
16 participants
17 is_first_round
18 is_last_round
19 source_url
20 source_description
21 created_by
22 created_at
23 updated_at

```

Listado 11 Columnas del archivo Funding Rounds

```

1 id
2 funding_round_id
3 funded_object_id
4 investor_object_id
5 created_at
6 updated_at

```

Listado 12 Columnas del archivo Investments

VII-B. Resultado de ejecución de scripts de carga de datos

```

└─$ yarn setup:postgres
yarn run v1.22.19
$ node bin/setup/postgres.js
[18:34:39.867] INFO: Connection succeeded
[18:34:39.920] INFO: Enabled extensions
[18:34:40.071] INFO: Dropped table objects
[18:34:40.121] INFO: Dropped table funding_rounds
[18:34:40.126] INFO: Dropped table investments
[18:34:40.211] INFO: Created table objects
[18:34:40.225] INFO: Created table funding_rounds
[18:34:40.230] INFO: Created table investments
[18:34:40.230] INFO: Start processing file: objects.csv
[18:35:08.821] INFO: Parsed 462650 rows
[18:35:08.829] INFO: Finish processing file: objects.csv
[18:35:08.831] INFO: Start processing file: funding_rounds.csv
[18:35:11.081] INFO: Parsed 52626 rows
[18:35:11.756] INFO: Finish processing file: funding_rounds.csv
[18:35:11.757] INFO: Start processing file: investments.csv
[18:35:13.229] INFO: Parsed 80245 rows
[18:35:16.316] INFO: Finish processing file: investments.csv
✨ Done in 38.17s.

```

Figura 4: Carga base de datos Postgres

```

└─$ yarn setup:mongo
yarn run v1.22.19
$ node bin/setup/mongo.js
[18:51:41.974] INFO: Finished exporting PG objects to json.
  File path: /Users/santiago/Documentos Mac /bdnr/data/objects.json
  Elapsed time: 7 seconds
[18:51:41.974] INFO: Starting to sanitize file...
[18:51:45.226] INFO: Finished sanitizing file. Elapsed time: 11 seconds
[18:51:45.339] INFO: Finished dropping objects collection
2022-06-25T18:51:46.061-0300    connected to: mongodb://localhost:27017/bdnr
2022-06-25T18:51:49.061-0300    [###.....] bdnr.objects 42.4MB/301MB (14.1%)
2022-06-25T18:51:52.061-0300    [#####.....] bdnr.objects 87.5MB/301MB (29.1%)
2022-06-25T18:51:55.062-0300    [#####.....] bdnr.objects 134MB/301MB (44.6%)
2022-06-25T18:51:58.061-0300    [#####.....] bdnr.objects 181MB/301MB (60.3%)
2022-06-25T18:52:01.061-0300    [#####.....] bdnr.objects 224MB/301MB (74.6%)
2022-06-25T18:52:04.062-0300    [#####.....] bdnr.objects 266MB/301MB (88.6%)
2022-06-25T18:52:06.622-0300    [#####.....] bdnr.objects 301MB/301MB (100.0%)
2022-06-25T18:52:06.622-0300    462650 document(s) imported successfully. 0 document(s) failed to import.
[18:51:45.744] INFO: Finished indexing objects collection
[18:51:45.744] INFO: Importing into mongodb...
[18:52:06.634] INFO: Finished importing into mongodb. Elapsed time: 32 seconds
✨ Done in 33.27s.

```

Figura 5: Carga base de datos Mongo

```

└─$ yarn setup:neo4j
yarn run v1.22.19
$ node bin/setup/neo4j.js
/Users/santiago/DocumentosMac/bdnr/data/objects.json -> /usr/local/Cellar/neo4j/4.4.8/libexec/import/objects.json
[20:15:46.657] INFO: STEP 1/4 Finished copying file into imports directory: /usr/local/Cellar/neo4j/4.4.8/libexec/import/objects.json. Timestamp: 0 seconds
[20:15:46.658] INFO: STEP 2/4 Starting to destroy existent nodes. Timestamp: 0 seconds
[20:15:46.741] INFO: STEP 2/4 Deleted all Neo4j nodes. Timestamp: 1 seconds
[20:15:46.741] INFO: STEP 3/4 Starting to create indexes. Timestamp: 1 seconds
[20:15:46.762] INFO: STEP 3/4 Created indexes successfully. Timestamp: 1 seconds
[20:15:46.763] INFO: STEP 4/4 Starting import. Timestamp: 1 seconds
52626
[20:16:27.088] INFO: STEP 4/4 Imported objects successfully. Timestamp: 41 seconds
✨ Done in 42.76s.

```

Figura 6: Carga base de datos Mongo

VII-C. Consultas

Listado 13 Consulta 1 PostgreSQL

```

1  SELECT
2  c.id,
3  c.name,
4  c.founded_at,
5  c.homepage_url,
6  c.logo_url,
7  c.description,
8  c.overview,
9  c.category_code,
10 c.country_code,
11 c.state_code,
12 c.city,
13 c.region,
14 c.funding_total_usd,
15 (
16  SELECT array_agg(f)
17  FROM (
18    SELECT
19      f.id,
20      f.raised_amount_usd,
21      f.funded_at,
22      f.funding_round_type as type,
23      (
24        SELECT array_agg(x)
25        FROM (
26          SELECT
27            i.id,
28            i.name,
29            i.entity_type
30          FROM investments x
31          LEFT JOIN objects i
32            ON i.id = x.investor_object_id
33          WHERE f.id = x.funding_round_id
34        ) x
35      ) AS investors
36    FROM funding_rounds f
37    WHERE c.id = f.object_id
38    ORDER BY f.funded_at
39  ) f
40 ) AS funding_rounds
41 FROM objects c
42 WHERE lower(c.name) LIKE lower($1)
43 ORDER BY c.name ASC
44 LIMIT 1

```

Listado 14 Consulta 2 PostgreSQL

```

1  SELECT
2  i.id as id,
3  i.name as name,
4  i.entity_type,
5  investments_count
6  FROM (
7    SELECT
8      i.id,
9      i.name,
10     i.entity_type,
11     COUNT(*) AS investments_count
12    FROM objects c
13    INNER JOIN funding_rounds f
14      ON c.id = f.object_id
15    INNER JOIN investments x
16      ON f.id = x.funding_round_id
17    INNER JOIN objects i
18      ON i.id = x.investor_object_id
19    WHERE c.category_code = $1
20    AND c.city = $2
21    GROUP BY i.id, i.name, i.entity_type
22  ) i

```

```

23 ORDER BY
24   investments_count DESC,
25   id ASC
26 LIMIT 20

```

Listado 15 Consulta 3 PostgreSQL

```

1 SELECT
2   o.id,
3   o.name,
4   COUNT(DISTINCT x.investor_object_id) AS matches_count
5 FROM objects o
6 INNER JOIN investments x
7   ON o.id = x.funded_object_id
8 WHERE x.investor_object_id IN (
9   -- my co-investors ids:
10  SELECT DISTINCT
11    others_investments.investor_object_id AS investor_id
12  FROM investments my_investments
13  INNER JOIN investments others_investments
14    ON others_investments.funded_object_id = my_investments.funded_object_id
15  AND others_investments.investor_object_id <> $1
16  WHERE my_investments.investor_object_id = $1
17 )
18 AND o.id NOT IN (
19   -- ids of companies I already invested in:
20  SELECT DISTINCT my_investments.funded_object_id
21  FROM investments my_investments
22  WHERE my_investments.investor_object_id = $1
23 )
24 GROUP BY o.id
25 ORDER BY COUNT(DISTINCT x.investor_object_id) DESC, o.id ASC
26 LIMIT 20

```

Listado 16 Consulta 1 Mongo

```

1
2 const query = { name: {

```

Listado 17 Consulta 2 Mongo

```

1
2 const pipeline = [
3   {
4     $match: {
5       category_code: category,
6       city
7     }
8   },
9   {
10    $project: {
11      funding_rounds: 1
12    }
13  },
14  {
15    $unwind: '$funding_rounds'
16  },
17  {
18    $unwind: '$funding_rounds.investors'
19  },
20  {
21    $group: {
22      _id: '$funding_rounds.investors.id',
23      id: {
24        $first: '$funding_rounds.investors.id'
25      },
26      name: {
27        $first: '$funding_rounds.investors.name'
28      },
29      entity_type: {
30        $first: '$funding_rounds.investors.entity_type'

```

```

31     },
32     investments_count: {
33       $sum: 1
34     }
35   },
36 },
37 {
38   $sort: {
39     investments_count: -1,
40     _id: 1
41   }
42 },
43 {
44   $limit: 20
45 },
46 {
47   $unset: '_id'
48 }
49 ]
50
51 collection.aggregate(pipeline)

```

Listado 18 Consulta 3 Mongo

```

1
2 const pipeline = [
3   {
4     $match: {
5       'funding_rounds.investors.id': investorId // select the objects the investor invested in. Uses
6         index
7     },
8   },
9   {
10    $unwind: '$funding_rounds'
11  },
12  {
13    $unwind: '$funding_rounds.investors' // get one record per co-investor investment
14  },
15  {
16    $project: {
17      investor_id: '$funding_rounds.investors.id' // we only care about co-investor's ids
18    }
19  },
20  {
21    $group: {
22      _id: '$investor_id' // get distinct records only
23    }
24  },
25  {
26    $set: {
27      investor_id: '$_id'
28    }
29  },
30  {
31    $unset: '_id'
32  }, // up until here we have the list of distinct co-investors ids
33  {
34    $lookup: { // get each coinvestor's investments
35      from: 'objects',
36      localField: 'investor_id',
37      foreignField: 'funding_rounds.investors.id',
38      as: 'investments_tmp',
39      pipeline: [
40        {
41          $project: {
42            _id: 0,
43            id: '$id',
44            name: 1
45          }
46        }
47      ]
48    }
49  },
50  {
51    $unwind: '$investments_tmp'

```

```

51   },
52   {
53     $project: {
54       investor_id: 1,
55       id: '$investments_tmp.id',
56       name: '$investments_tmp.name'
57     }
58   }, // up until here we have an array with all coinvestor's investments. Eg: [{ investor_id: 'c:70937',
59     // id: 'c:189377', name: 'Polar' }]
60   {
61     $group: { // count number of occurrences by object id (how many investments each company got from
62       // the co-investors)
63       _id: '$id',
64       name: {
65         $first: '$name'
66       },
67       matches_count: {
68         $sum: 1
69       },
70       investor_ids: { $addToSet: '$investor_id' }
71     },
72     {
73       $match: {
74         $expr: { $not: { $in: [investorId, '$investor_ids'] } } // filter out objects the current investor
75           // already invested in
76       }
77     },
78     {
79       $sort: {
80         matches_count: -1,
81         _id: 1
82       }
83     },
84     {
85       $limit: 20
86     },
87     {
88       $project: {
89         _id: 0,
90         id: '$_id',
91         name: '$name',
92         matches_count: '$matches_count'
93       }
94     }
95   ]
96 }
97 collection.aggregate(pipeline)

```

Listado 19 Consulta 1 Neo4j

```

1  MATCH (c:Object)
2  WHERE c.name STARTS WITH $name
3  OPTIONAL MATCH (c)-[:BELONGS_TO]->(city:City)
4  OPTIONAL MATCH (c)-[:BELONGS_TO]->(category:Category)
5  OPTIONAL MATCH (f:FundingRound)-[:BELONGS_TO]->(c)
6  OPTIONAL MATCH (i:Object)-[:INVESTED_IN]->(f)
7  WITH c, i, f, city, category
8  ORDER BY c.name ASC
9  LIMIT 1
10 WITH c,f,city,category,
11 {
12   id: i.id,
13   name: i.name,
14   type: i.entity_type
15 } as investors
16 WITH c,city,category,
17 {
18   id: f.id,
19   type: f.type,
20   funded_at: f.funded_at,
21   raised_amount_usd: f.raised_amount_usd,
22   investors: collect(investors)
23 } as funding_rounds
24 WITH {

```

```

25     category_code: category.name,
26     city: city.name,
27     id: c.id,
28     name: c.name,
29     founded_at: c.founded_at,
30     homepage_url: c.homepage_url,
31     logo_url: c.logo_url,
32     description: c.description,
33     overview: c.overview,
34     country_code: c.country_code,
35     state_code: c.state_code,
36     region: c.region,
37     funding_total_usd: c.funding_total_usd,
38     funding_rounds: collect(funding_rounds)
39   } as companies
40   RETURN companies
41   LIMIT 1

```

Listado 20 Consulta 2 Neo4j

```

1  MATCH (i:Object)-[:INVESTED_IN]->(f:FundingRound)-[:BELONGS_TO]->(c:Object),
2     (c)-[:BELONGS_TO]->(city:City { name: $city }),
3     (c)-[:BELONGS_TO]->(category:Category { name: $category })
4  WITH i, COUNT(f) AS investments_count
5  RETURN i.id AS id, i.name AS name, i.entity_type AS entity_type, investments_count
6  ORDER BY investments_count DESC, i.id ASC
7  LIMIT 20

```

Listado 21 Consulta 3 Neo4j

```

1  MATCH
2     (me:Object {id: $investorId})-[:INVESTED_IN]->(f:FundingRound)-[:BELONGS_TO]->(mycompany:Object)
3  WITH me, collect(distinct mycompany) AS mycompanies
4  UNWIND mycompanies AS mycompany
5  MATCH
6     (coinvestor:Object)-[:INVESTED_IN]->(f:FundingRound)-[:BELONGS_TO]->(mycompany)
7  WITH mycompanies, collect(distinct coinvestor) AS coinvestors
8  UNWIND coinvestors AS coinvestor
9  MATCH
10     (coinvestor)-[:INVESTED_IN]->(f:FundingRound)-[:BELONGS_TO]->(othercompany:Object)
11 WHERE NOT othercompany IN mycompanies
12 WITH othercompany, COUNT(distinct coinvestor) AS matches_count
13 ORDER BY matches_count DESC, othercompany.id ASC
14 LIMIT 20
15 RETURN othercompany.id AS id, othercompany.name AS name, matches_count

```

Listado 22 Consulta PostgreSQL para exportar datos para Mongo y Neo4j

```

1  SELECT
2     c.id,
3     c.name,
4     c.founded_at,
5     c.homepage_url,
6     c.logo_url,
7     c.description,
8     c.overview,
9     c.category_code,
10    c.country_code,
11    c.state_code,
12    c.city,
13    c.region,
14    c.funding_total_usd,
15    c.entity_type AS entity_type,
16  (
17    SELECT array_to_json(array_agg(f))
18  FROM (
19    SELECT
20     f.id,
21     f.raised_amount_usd,
22     f.funded_at,
23     f.funding_round_type as type,

```

```

24     (
25     SELECT array_to_json(array_agg(x))
26     FROM (
27     SELECT
28     i.id,
29     i.name,
30     i.entity_type as entity_type
31     FROM investments x
32     LEFT JOIN objects i
33     ON i.id = x.investor_object_id
34     WHERE f.id = x.funding_round_id
35     ) x
36     ) AS investors
37     FROM funding_rounds f
38     WHERE c.id = f.object_id
39     ORDER BY f.funded_at
40     ) f
41 ) AS funding_rounds
42 FROM objects c
43 ORDER BY c.name ASC

```

Listado 23 Consulta carga datos Neo4j

```

1 CALL apoc.load.json('file://\${jsonFileName}')
2 YIELD value AS value
3 MERGE (city:City {
4   name: coalesce(value.city, "Unknown")
5 })
6 MERGE (category:Category {
7   name: coalesce(value.category_code, "Unknown")
8 })
9 CREATE (o:Object {
10  id: value.id,
11  name: value.name,
12  logo_url: value.logo_url,
13  homepage_url: value.homepage_url,
14  entity_type: value.entity_type,
15  founded_at: value.founded_at,
16  homepage_url: value.homepage_url,
17  description: value.description,
18  overview: value.overview,
19  funding_total_usd: value.funding_total_usd,
20  country_code: value.country_code,
21  state_code: value.state_code,
22  region: value.region
23 })
24 CREATE (o)-[:BELONGS_TO]->(city)
25 CREATE (o)-[:BELONGS_TO]->(category)
26 WITH o, value
27 UNWIND value.funding_rounds AS funding_round
28 MERGE (f:FundingRound {
29   id: toInteger(funding_round.id),
30   type: funding_round.type,
31   raised_amount_usd: funding_round.raised_amount_usd,
32   funded_at: coalesce(funding_round.funded_at, "")
33 })
34 CREATE (f)-[:BELONGS_TO]->(o)
35 FOREACH (investor IN funding_round.investors |
36   MERGE (i:Object {
37     id: investor.id
38   })
39   MERGE (i)-[:INVESTED_IN]->(f)
40 )
41 RETURN count(*)

```

REFERENCIAS

- [1] Crunchbase. Crunchbase official page. <https://www.crunchbase.com/>. Accessed: June 2022.
- [2] Neal Ford. Polyglot Programming - Neal Ford. <http://memeagora.blogspot.com/2006/12/polyglot-programming.html>. Accessed: June 2022.
- [3] Martin Fowler. PolyglotPersistence - Martin Fowler. <https://martinfowler.com/bliki/PolyglotPersistence.html>. Accessed: June 2022.
- [4] Kaggle. Dump de datos sobre Crunchbase 2013. <https://www.kaggle.com/datasets/justinas/startup-investments?select=objects.csv>. Accessed: June 2022.
- [5] Kaggle. Kaggle Datasets. <https://www.kaggle.com/datasets>. Accessed: June 2022.
- [6] Pwint Khine and Zhaoshun Wang. A review of polyglot persistence in the big data world. *Information*, 10:141, 04 2019.

- [7] Neha Singh Rohit Chandra Sunil B. Shambharkar J. J. Kulkarni. Review of nosql databases and performance comparison between multimodel and polyglot persistence approach. *INTERNATIONAL JOURNAL OF SCIENTIFIC TECHNOLOGY RESEARCH*, 9:5, 01 2020.
- [8] Grafana Labs. K6 Load testing tool. <https://k6.io>. Accessed: June 2022.
- [9] Scott Leberknight. Polyglot Persistence - Scott Leberknight's Weblog. http://www.sleberknight.com/blog/sleberkn/entry/polyglot_persistence. Accessed: June 2022.
- [10] MongoDB. Model One-to-Many Relationships with Embedded Documents. <https://www.mongodb.com/docs/manual/tutorial/model-embedded-one-to-many-relationships-between-documents/#embedded-document-pattern>. Accessed: June 2022.
- [11] Neo4j. How much faster is a graph database, really? <https://neo4j.com/news/how-much-faster-is-a-graph-database-really/>. Accessed: July 2022.
- [12] Donato Aguirre Santiago Acevedo. Repositorio BDNR Grupo 5. <https://gitlab.fing.edu.uy/santiago.acevedo/bdnr-grupo5>. Accessed: June 2022.