

Universidad de La Habana
Facultad de Matemática y Computación



Enfoques Zero-Shot para la Extracción de Conocimiento a partir de Lenguaje Natural

Autor: Rolando Sánchez Ramos

Tutor: Dr. Alejandro Piad Morffis

Trabajo de Diploma
presentado en opción al título de
Licenciado en Ciencias de la Computación



Noviembre de 2023

github.com/rolysr/nl2ql

Agradecimientos

Opinión del tutor

Dr. Alejandro Piad Morffis
Facultad de Matemática y Computación
Universidad de la Habana
Noviembre, 2023

Resumen

Esta tesis se centra en abordar la complejidad inherente a la consulta de bases de datos en forma de grafo, como Neo4J. Estas bases de datos a menudo requieren un conocimiento especializado en lenguajes de consulta, lo que limita su accesibilidad a un grupo reducido de usuarios con habilidades técnicas avanzadas. Para superar esta limitación, proponemos la aplicación del aprendizaje *Zero-Shot*, un enfoque innovador en el procesamiento del lenguaje natural. En esta investigación, se lleva a cabo un experimento basado en el modelo GPT-4 para traducir consultas de lenguaje natural a código *Cypher*. La evaluación se realiza utilizando el conjunto de datos de evaluación MetaQA, que abarca una amplia variedad de ejemplos de consultas. Los resultados obtenidos fueron del 76,53 %, 43,45 % y 31,03 % para los tres lotes de evaluación del *benchmark* utilizado, mejorando de esta forma el mejor resultado de modelos de lenguaje en la traducción de lenguaje natural a código *Cypher* sobre MetaQA mediante el aprendizaje *Zero-Shot*.

Abstract

This thesis focuses on addressing the inherent complexity of querying graph databases, such as Neo4J. These databases often require specialized knowledge in query languages, limiting their accessibility to a small group of users with advanced technical skills. To overcome this limitation, we propose the application of Zero-Shot learning, an innovative approach in natural language processing. In this research, an experiment is conducted based on the GPT-4 model to translate natural language queries into Cypher code. The evaluation is carried out using the MetaQA evaluation dataset, which covers a wide variety of query examples. The results obtained were 76,53 %, 43,45 %, and 31,03 % for the three evaluation lots of the benchmark used, thereby improving the best result of language models in translating natural language into Cypher code using Zero-Shot learning.

Índice general

Introducción	11
1. Estado del Arte	16
1.1. Preliminares	18
1.1.1. Bases de Datos <i>Neo4J</i>	18
1.1.2. Lenguaje <i>Cypher</i>	19
1.1.3. Grandes modelos de Lenguaje y Generación de Lenguajes de Consulta Formales	21
1.2. Extracción de conocimiento mediante enfoques neurosimbólicos basados en representación intermedia (IR, por sus siglas en inglés) de la consulta dada en lenguaje natural	22
1.2.1. Conjuntos de entranamiento y evaluación	22
1.2.2. Preprocesamiento	23
1.2.3. Postprocesamiento	23
1.2.4. Consulta	24
1.3. Extracción de conocimiento mediante enfoques basados en técnicas de <i>prompt engineering</i> mediante LLMs	24
1.3.1. Zero-Shot Learning	24
1.3.2. Few-Shot Learning	25
1.3.3. Chain-of-Thought Prompting	26
1.3.4. <i>Fine-Tuning</i>	27
1.4. Consideraciones generales	28
2. Propuesta de Solución	29
2.1. Definición del problema <i>Text-to-Cypher</i>	30
2.1.1. LLM para resolver <i>Text-to-Cypher</i>	30
2.2. Enfoques considerados	31
2.3. Propuesta de solución diseñada	32
2.3.1. Interacción con una base de datos <i>Neo4J</i>	33

2.3.2.	Extracción de información de una base de datos <i>Neo4J</i>	34
2.3.3.	Almacenamiento de información en una base de datos <i>Neo4J</i>	36
2.3.4.	Selección del modelo	37
2.3.5.	Diseño de la información entrada al LLM	39
2.3.6.	Caso de estudio	41
3.	Detalles de Implementación	44
3.1.	Despliegue de una instancia de una base de datos <i>Neo4J</i> . . .	44
3.2.	GraphContractor	46
3.3.	KnowledgeBase	47
3.4.	DBSeeder	48
3.5.	SchemaMaker	49
3.6.	GPT-4	50
4.	Análisis Experimental	52
4.1.	Evaluación sobre el <i>benchmark MetaQA</i> [?]	52
4.1.1.	Resultados	54
4.2.	Discusiones	59
	Conclusiones	60
	Bibliografía	64

Índice de figuras

1.1. Ejemplo de representación de información en una base de datos <i>Neo4J</i> . Nótese la estructura de nodos para representar información de entidades y de aristas que corresponden a relaciones entre estas.	19
1.2. Ejemplo de código <i>Cypher</i>	20
1.3. Secuencia de flujo representativa del Estado del Arte de traducción de Lenguaje Natural a Lenguaje Formal utilizando el enfoque neurosimbólico basado en (IR).	22
1.4. Ejemplo del flujo de trabajo de experimentos basados en ZSL para la traducción de lenguaje natural a código en <i>SQL</i>	25
1.5. Ejemplo del flujo de trabajo de experimentos basados en FSL para la traducción de lenguaje natural a código en <i>SQL</i>	26
1.6. Ejemplo del flujo de trabajo de experimentos basados en CoT para la traducción de lenguaje natural a código en <i>SQL</i>	27
2.1. Flujo de funcionamiento del componente <i>GraphContractor</i>	34
2.2. Flujo de funcionamiento del componente <i>KnowledgeBase</i>	36
2.3. Flujo de funcionamiento del componente <i>DBSeeder</i>	37
2.4. Ejemplo básico sobre cómo utilizar GPT-4 para traducir lenguaje natural en lenguaje de consulta <i>Cypher</i>	39
2.5. Ejemplo del proceso que se realiza para obtener un formato verbalizado de una base de datos <i>Neo4J</i> con el uso del <i>SchemaMaker</i>	40
2.6. Flujo de entrada y salida en el proceso de traducción de lenguaje natural a código en <i>Cypher</i>	41
2.7. Arquitectura y funcionamiento de la propuesta de solución. . . .	42
3.1. Comando utilizado para desplegar una base de datos <i>Neo4J</i>	45
3.2. Implementación de la clase <i>GraphContractor</i>	46

3.3. Implementación de la clase KnowledgeBase.	47
3.4. Implementación de la clase DBSeeder.	48
3.5. Implementación de la clase Schema Maker.	49
3.6. Implementación para utilizar el modelo GPT-4.	51

Introducción

En la época actual, asistimos a un constante aumento en la producción de información en diversos formatos: visual, auditivo y textual, que abarca todos los ámbitos de la sociedad [9]. De manera particular, resulta sumamente intrigante la información generada a través del ingenio creativo y la investigación humana. Estos tipos de datos se almacenan debido a su relevancia y a la necesidad de acceder a ellos en el futuro, pudiendo optar por una organización estructurada o no. Sorprendentemente, solo alrededor del 20% de la información a nivel mundial se encuentra estructurada [1].

Las bases de conocimiento constituyen un tipo particular de bases de datos diseñadas para la administración del saber. Estas bases brindan los medios para recolectar, organizar y recuperar digitalmente un conjunto de conocimientos, ideas, conceptos o datos [2]. La ventaja fundamental de mantener la información de manera estructurada radica en su facilidad para ser consultada, ampliada y modificada. Debido a su utilidad y prevalencia, la recuperación de información a través de consultas en bases de conocimiento se ha convertido en una tarea esencial.

Es primordial que la información almacenada en bases de conocimiento adopte un formato adecuado para permitir búsquedas ágiles y precisas. Entre los formatos más comunes se encuentran los modelos de Entidad-Relación y el modelo Relacional. A pesar de ser enfoques más antiguos, el modelo Relacional (BDR) sigue siendo el más ampliamente utilizado en la actualidad [6]. No obstante, en ocasiones, las características específicas del problema demandan un formato más expresivo, y es en este punto donde las bases de datos orientadas a grafos (BDOG) [17] entran en juego.

Las BDOG han ganado progresivamente popularidad como una manera efectiva de almacenar información en los últimos años. Estas bases tienen la capacidad de modelar una diversidad de situaciones del mundo real al tiempo que mantienen un alto nivel de simplicidad y legibilidad para

los seres humanos. Las BDOG presentan numerosas ventajas en comparación con las bases de datos relacionales. Esto incluye un mejor rendimiento, permitiendo el manejo más rápido y eficaz de grandes volúmenes de datos relacionados; flexibilidad, ya que la teoría de grafos en la que se basan las BDOG permite abordar diversos problemas y encontrar soluciones óptimas; y escalabilidad, ya que las bases de datos orientadas a grafos permiten una escalabilidad eficaz al facilitar la incorporación de nuevos nodos y relaciones entre ellos. Ejemplo de un sistema de gestión de BDOG es *Neo4J* [13], a través del cual es posible construir instancias de este tipo de base de datos e interactuar con las mismas a través del lenguaje de programación *Cypher* [14], el cual posee una sintaxis declarativa similar a *SQL* [16].

Por otro lado, el avance en la comprensión del lenguaje natural se ha visto potenciado con el surgimiento de los grandes modelos de lenguajes (LLMs) [18] como GPT-4 [11] o LLaMA-2 [5], los cuales presentan una serie de habilidades emergentes como elaboración de resúmenes de textos, generación de código, razonamiento lógico, traducción lingüística entre otras [10]. Dichas herramientas constituyen modelos de *Machine Learning* entrenados con un gran volumen de datos, lo cual es posible gracias al número de parámetros con los que estos son configurados [18].

Usualmente, para el uso de los LLMs basta con ofrecerles como dato de entrada un texto [18], el cual describe la tarea que se espera que estos realicen. Además, son muchas las técnicas existentes para elaborar una entrada de calidad, esto con el objetivo de que la respuesta por parte de dicho modelo de lenguaje ofrezca resultados alentadores al respecto, lo cual se conoce como *prompt engineering* [7]. Una técnica bastante común es *Zero-Shot Learning* (ZSL) [12], la cual consiste en describirle a un LLM un procedimiento a realizar sin ofrecer de antemano ejemplos de cómo resolverlo, como por ejemplo, en tareas relacionadas con la generación de código, donde algunos de estos son capaces de generar algoritmos expresados en un lenguaje de programación formal a partir de una sentencia o consulta en lenguaje natural sin recibir como entrada del usuario algunos ejemplos de código, o especificaciones de cómo funciona el lenguaje objetivo a generar [8] [3].

En lo que respecta a la comprensión del lenguaje natural y su uso en consultas a bases de conocimiento, existen diversas vías llevadas a cabo y con resultados diversos, donde se hacen análisis sintácticos y semánticos sobre la consulta, muchas veces asistidos por diccionarios o mapas sobre la base de conocimiento en cuestión. Se usan modelos de paráfrasis como técnica de aumento de datos y finalmente *Transformers* o incluso LLMs para llevar de la consulta ya curada al lenguaje de consulta formal o a un

lenguaje intermedio capaz de expresar a esta a alto nivel [4] [15].

Por las razones anteriormente expuestas, resulta interesante la investigación sobre la tarea de generación de código de consulta formal a partir de una sentencia en lenguaje natural mediante el uso de LLMs, especialmente el diseño e implementación un experimento capaz de demostrar las capacidades reales de estos para dicho acometido, lo cual designará la importancia de continuar el estudio de dichas herramientas con el objetivo de mejorar los sistemas de extracción de conocimientos en BDOG.

Problemática

Para utilizar el lenguaje de consulta formal *Cypher* se requiere de conocimientos básicos de programación, lo cual consume cierto tiempo y esfuerzo. Esto tiene como consecuencia que, solo aquellas personas con experiencia en el uso de lenguajes de programación puedan hacer uso de la mayoría de los sistemas de almacenamiento de datos desarrollados con esta tecnología y teniendo en cuenta la necesidad de poseer un conocimiento del dominio sobre el cual está construida la base de datos a consultar. Por lo tanto, llevar a cabo una mejora en las herramientas orientadas a democratizar dicho proceso permitiría hacer más rápido y eficiente dicha consulta en cuanto a tiempo y recursos computacionales. Debido a dicha situación, se propone una experimentación basada en un LLM capaz de traducir una consulta en lenguaje natural a un código en *Cypher*, donde a su vez se verifique la efectividad de este a partir de enfoques basados en ZSL, los cuales intuitivamente pueden ofrecer como resultado una cota inferior para la efectividad de sistemas desarrollados en base a dichos algoritmos de aprendizaje. Además, actualmente la implementación de sistemas de generación de código de consulta formal está principalmente orientada al lenguaje *SQL*, mientras que para el lenguaje *Cypher*, no existen suficientes estudios recientes que avalen la calidad de tales herramientas para dicho caso de uso, incluso cuando las bases de datos orientadas al segundo lenguaje representan muchos de los sistemas de almacenamiento de conocimientos, como por ejemplo, las redes semánticas.

Objetivos

Dadas las ideas anteriores, los objetivos principales del trabajo consistirán en diseñar e implementar una estrategia experimental capaz de verificar la capacidad mínima de los LLMs para la consulta en lenguaje natural a

bases de conocimiento estructuradas con independencia del dominio, para lo cual se empleará un enfoque basado en ZSL.

Para lograr los objetivos generales se trazaron los siguientes objetivos específicos:

1. Estudiar el estado del arte de los modelos de Aprendizaje Automático capaces de hacer predicciones de tipo texto-a-texto.
2. Analizar el trabajo de tesis sobre este tema anteriormente desarrollado en la facultad.
3. Implementar un modelo de Aprendizaje Automático capaz de convertir una consulta en lenguaje natural humano a un lenguaje formal que permita obtener datos a partir de una base de conocimiento.
4. Explorar las capacidades de enfoques *Zero-Shot* para la traducción de lenguaje natural al lenguaje *Cypher* con el fin de desarrollar un modelo capaz de realizar dicha tarea sin necesidad de ser entrenados directamente para la misma.
5. Mejorar el sistema de evaluación de resultados permitiendo que el conjunto de datos de prueba y evaluación sea lo más realista posible y con una mayor complejidad.

Organización de la tesis

El presente documento está estructurado en 5 capítulos que engloban las etapas cubiertas en la investigación. En el capítulo 1, Estado del Arte, se reseña el estado actual de la teoría, herramientas y técnicas más usadas en los temas tratados. En el capítulo 2, Propuesta de Solución, se propone una sistema que responde a algunas de las limitaciones principales de los modelos desarrollados en el estado del arte. Para ello propone una vía de resolver el problema basada fundamentalmente en el aprendizaje *Zero-Shot*. En el capítulo 3, Detalles de Implementación, se describen por menores en la implementación del modelo propuesto como solución, se esclarecen decisiones de diseño, y se muestran porciones del código referente a los componentes principales desarrollados. En el capítulo 4, Análisis Experimental, se sugiere un marco experimental para analizar los resultados obtenidos durante la investigación y se comprueba experimentalmente una mejora con respecto a modelos utilizando enfoques similares en el estado del arte. Finalmente, se formulan las conclusiones, que recogen los resultados obtenidos en la investigación en función de los objetivos definidos, así

como las recomendaciones, donde se proponen siguientes líneas de trabajo a ser exploradas en continuación de la investigación actual. Para finalizar se indican las referencias bibliográficas consultadas, con el fin de complementar la información provista en el trabajo.

Capítulo 1

Estado del Arte

Con el incremento constante de la cantidad de información generada en todo el mundo, la recuperación de información se ha convertido en un aspecto de creciente relevancia tanto en el ámbito industrial como en el académico. En consecuencia, la reducción del tiempo transcurrido entre el momento en que un usuario desea acceder a la información y el momento en que efectivamente puede hacerlo ha sido objeto de un número creciente de investigaciones científicas en los últimos años. Este capítulo se dedica a la evaluación de diversas estructuras de interfaces entre el ser humano y bases de conocimiento, las cuales tienen como objetivo abordar esta problemática.

Las Interfaces de Lenguaje Natural a Bases de Datos (NLIDB, por sus siglas en inglés) representan un campo de investigación dinámico centrado en facilitar las interacciones entre humanos y computadoras con bases de datos relacionales utilizando consultas en lenguaje natural. A lo largo de las últimas décadas, el desarrollo de NLIDB ha pasado por varias fases transformadoras, impulsadas por avances tecnológicos y metodológicos, así como por una creciente demanda de una mejor accesibilidad a las bases de datos.

Las etapas iniciales del desarrollo de NLIDB se caracterizaron por sistemas específicos de dominio. Estos sistemas fueron diseñados para trabajar dentro de áreas de conocimiento bien definidas, donde se utilizaba el procesamiento de lenguaje natural controlado para garantizar la comprensión de las consultas y la interacción con la base de datos. Por ejemplo, algunos trabajos pioneros [?] [?] demostraron NLIDBs que se adaptaban a dominios específicos, lo que los hacía altamente efectivos pero limitados en alcance. Del mismo modo, en años posteriores [?] [?] se continuó la exploración del

uso de interfaces de lenguaje natural controlado dentro de dominios de conocimiento particulares.

Otro enfoque durante esta fase implicó NLIDBs basados en reglas. Algunos sistemas propuestos al respecto [?] dependían de reglas predefinidas para traducir consultas en lenguaje natural en declaraciones *SQL* para la recuperación de datos en la base de datos. Si bien estos sistemas ofrecían ciertas ventajas, carecían de versatilidad para manejar una amplia gama de consultas de usuarios en diferentes dominios.

A medida que avanzaba la investigación en NLIDB, hubo un cambio hacia la independencia de dominio y la flexibilidad. Los sistemas recientes han buscado reducir la dependencia del conocimiento específico del dominio y las reglas. Algunos investigadores [?] [?] han desarrollado NLIDBs que utilizan técnicas de aprendizaje supervisado, lo que los hace más adaptables a varios dominios y entradas de usuario. Además, un avance significativo se ha producido con la integración de redes neuronales profundas en el desarrollo de NLIDBs, donde varias investigaciones [?] [?] han demostrado el potencial del aprendizaje profundo en NLIDB, aprovechando vastos repositorios de texto y código para el entrenamiento. Este enfoque ha mejorado significativamente el rendimiento de NLIDB, permitiendo un procesamiento de consultas más natural y contextual.

Para el caso específico de NLIDB con respecto a BDOGs inicialmente se desarrollaron trabajos enfocados en técnicas similares a las empleadas para *SQL* [?] [?], donde se realizaba un preprocesamiento en la consulta dada, se aprovechaba la información ofrecida por el esquema [?] de la base de datos a consultar y finalmente dicho conocimiento era utilizado por un modelo de aprendizaje profundo entrenado sobre un conjunto de pares de lenguaje natural y lenguaje de consulta formal como por ejemplo *Cypher*.

Recientemente, los trabajos orientados a esta área de estudio han estado enfocados en dos metodologías principales que serán tratadas en las secciones 1.2 y 1.3:

1. Enfoques neurosimbólicos basados en representación intermedia de la consulta dada en lenguaje natural.
2. Enfoques basados en técnicas de *prompt engineering* mediante LLMs.

1.1. Preliminares

1.1.1. Bases de Datos *Neo4J*

El sistema *Neo4J* emerge como una plataforma de base de datos de grafos preeminente, distinguiéndose por su capacidad para manejar datos interconectados con una eficiencia y flexibilidad notables [?]. Este sistema gestiona las relaciones entre los datos con una estructura nodal y de aristas, lo que permite una representación más natural de las interconexiones inherentes a muchos conjuntos de datos [?].

La potencia de *Neo4J* reside en su lenguaje de consulta, *Cypher*, que permite expresar consultas sobre grafos de manera declarativa. *Cypher* es específicamente diseñado para ser intuitivo y potente, proporcionando comandos que facilitan la realización de patrones de búsqueda complejos y análisis de relaciones en una forma compacta y legible [?].

En términos de rendimiento, *Neo4J* está diseñado para maximizar la velocidad y eficiencia en la recuperación y manejo de datos relacionales. Utiliza índices basados en árboles B+ y algoritmos de ruta como el de Dijkstra y A* para búsquedas optimizadas en el grafo [?]. Además, la integridad de los datos se garantiza mediante propiedades transaccionales ACID [?], que son fundamentales en aplicaciones críticas de negocio [?].

La escalabilidad de *Neo4J* es una de sus características más destacables, con soporte para configuraciones en clúster que permiten la replicación de datos y la tolerancia a fallos, asegurando la disponibilidad y la escalabilidad horizontal [?]. Esto es esencial para aplicaciones que requieren un rendimiento consistente bajo cargas de trabajo de lectura y escritura intensivas.

La versatilidad de integración de *Neo4j* también merece ser subrayada. Su compatibilidad con API REST, diversos lenguajes de programación y frameworks de desarrollo facilita su adopción en arquitecturas de software existentes [?]. Esto se complementa con una comunidad activa y recursos extensos para desarrolladores, lo que promueve una continua innovación y adopción en la industria [?].

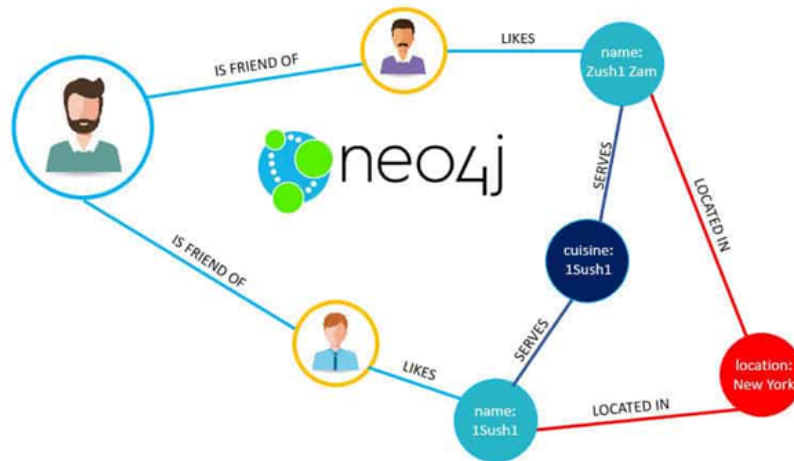


Figura 1.1: Ejemplo de representación de información en una base de datos *Neo4J*. Nótese la estructura de nodos para representar información de entidades y de aristas que corresponden a relaciones entre estas.

1.1.2. Lenguaje *Cypher*

El lenguaje de consulta utilizado en *Neo4J* es *Cypher*, un lenguaje declarativo centrado en el "qué" en lugar del "cómo" cuando se trata de recuperar datos. Esta naturaleza permite a los usuarios especificar patrones en los grafos sin requerir que ellos describan los algoritmos o pasos lógicos para encontrar esos patrones [?]. A continuación se mencionan algunos aspectos relevantes de este lenguaje de consulta formal y cómo facilitan la interacción con bases de datos de grafos como *Neo4J*.

1. **Patrones de coincidencia (Pattern Matching):** *Cypher* utiliza una sintaxis que se asemeja a los diagramas de entidad-relación ASCII para patrones de coincidencia, lo que facilita la representación visual de la estructura del grafo en el propio código [?]. Su sintaxis intuitiva hace que la comprensión y escritura de consultas sea más accesible, especialmente para usuarios nuevos en el manejo de bases de datos de grafos. Por ejemplo, para encontrar a un usuario y sus amigos en *Neo4J*, podríamos escribir una consulta como:
2. **Filtrado y condiciones:** *Cypher* permite incorporar condiciones dentro de los patrones de coincidencia o en cláusulas *WHERE* para filtrar

```
1 MATCH (user:Person)-[:FRIEND]->(friend) RETURN user, friend
```

Figura 1.2: Ejemplo de código *Cypher*.

resultados. Esto se asemeja al uso de `WHERE` en *SQL* pero está optimizado para trabajar con las conexiones entre nodos [?].

3. **Agregación de datos:** *Cypher* proporciona funciones de agregación, como `COUNT`, `SUM`, `AVG`, `MAX` y `textttMIN`, que permiten realizar cálculos sobre grupos de datos. Estas funciones son esenciales para resumir información sobre los datos conectados [?].
4. **Modificación de grafos:** Además de recuperar datos, *Cypher* puede ser utilizado para crear, actualizar y eliminar nodos y relaciones. Esto incluye la capacidad de manejar transacciones y asegurar la integridad de los datos [?].
5. **Optimización de consultas:** *Cypher* está diseñado para optimizar las consultas de grafos de forma automática. El planificador de consultas de *Neo4J* reorganiza y optimiza las operaciones de consulta para una ejecución eficiente, lo que abstrae una capa de complejidad para los desarrolladores [?].
6. **Extensibilidad:** *Cypher* es extensible, lo que significa que se pueden crear funciones definidas por el usuario y procedimientos almacenados que se pueden invocar dentro de las consultas. Esto permite personalizar y ampliar la funcionalidad del lenguaje para satisfacer necesidades específicas [?].
7. **Interoperabilidad:** A través de su protocolo *Bolt* y la API REST, *Cypher* puede ser utilizado desde una variedad de lenguajes de programación y entornos, permitiendo que sistemas externos interactúen con *Neo4J* [?].

1.1.3. Grandes modelos de Lenguaje y Generación de Lenguajes de Consulta Formales

Los Grandes Modelos de Lenguaje (LLMs) son una clase de modelos de procesamiento de lenguaje natural que han revolucionado la forma en que las computadoras comprenden y generan texto. Estos modelos se destacan por su capacidad para generar texto coherente y contextualmente relevante en función del contexto proporcionado. Esta habilidad es esencial cuando se trata de generar código de consulta en lenguajes formales como *Cypher* o *textSQL* [?].

Una de las ventajas clave de los LLMs es su capacidad para traducir preguntas en lenguaje natural en consultas en lenguaje formal de manera automática. Esto significa que pueden tomar preguntas como «¿Cuál es la población de Nueva York?» y convertirlas en consultas *SQL* precisas, como `SELECT population FROM cities WHERE name = 'Nueva York'`. Esta traducción automática simplifica significativamente la interacción entre humanos y sistemas de bases de datos, ya que elimina la necesidad de que los usuarios aprendan el lenguaje formal [?].

Otra característica destacada es la versatilidad de los LLMs en términos de lenguajes de consulta. No están limitados a un lenguaje específico; pueden generar consultas en varios lenguajes formales. Esto significa que pueden interactuar con diferentes sistemas de bases de datos que utilizan diferentes lenguajes de consulta, desde *SQL* hasta *Cypher* y *SPARQL*, entre otros [?].

Los LLMs también demuestran un fuerte entendimiento de contextos complejos en sus respuestas. Pueden manejar preguntas que involucran múltiples condiciones, cláusulas *JOIN*, filtros y agregaciones en las consultas, lo que los hace aptos para abordar consultas complejas en bases de datos [?].

Por último, estos modelos tienen la capacidad de inferir relaciones y estructuras de datos a partir del contexto proporcionado. Esto les permite generar consultas que exploran conexiones y patrones en los datos de una base de datos de manera inteligente, lo que es especialmente valioso en aplicaciones de análisis de datos y minería de información [?].

1.2. Extracción de conocimiento mediante enfoques neurosimbólicos basados en representación intermedia (IR, por sus siglas en inglés) de la consulta dada en lenguaje natural

Con respecto a la recuperación de información de una base de conocimiento con este enfoque, mediante consultas hechas en lenguaje natural se han hecho varias investigaciones científicas que se pueden agrupar bajo el patrón de la Figura 1.2.

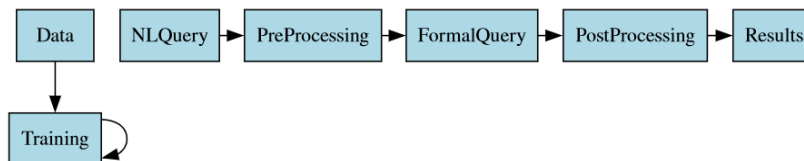


Figura 1.3: Secuencia de flujo representativa del Estado del Arte de traducción de Lenguaje Natural a Lenguaje Formal utilizando el enfoque neurosimbólico basado en (IR).

1.2.1. Conjuntos de entranamiento y evaluación

En cuanto a los datos para el entrenamiento existen dos opciones. Básicamente se puede buscar un conjunto de datos de referencia (*Benchmark*) en los que entrenar y probar el sistema, o se crea uno. La mayoría de las investigaciones existentes elijen la primera opción. Algunos de los *Benchmarks* más populares son:

- **WikiSQL:** WikiSQL [?] es el banco de datos más grande y más utilizado, contiene 26531 tablas y 80654 pares de consultas en lenguaje natural y lenguaje *SQL*. Las tablas se extraen de tablas *HTML* de Wikipedia. Luego, cada consulta en *SQL* se genera automáticamente para una determinada tabla bajo la restricción de que la consulta produce un conjunto de resultados no vacío.
- **Spider:** Este *Benchmark* [?] es un punto de referencia multidominio a gran escala con 200 bases de datos de 138 dominios diferentes y 10.181 pares de consultas.

- **MetaQA:** El conjunto de datos *MetaQA* [?] contiene más de 400000 pares de preguntas y respuestas de múltiples pasos obtenidos de la base de conocimiento WikiMovies [?]. Mientras que investigaciones previas se han centrado principalmente en la anotación SPARQL [?], nuestra innovación implica reconfigurar METAQA en *Cypher*, estableciéndolo como un valioso punto de referencia para el aprendizaje de pocos ejemplos.

En el caso alternativo, se suelen usar las propias bases de conocimiento objeto de estudio para crear un conjunto de entrenamiento. Una de las técnicas empleadas para esto es Random Walk [?], en la que se hace un recorrido aleatorio sobre un subconjunto de las entidades y relaciones de la base de conocimiento, y se elaboran consultas artificiales que respondan a dichas entidades y relaciones [?].

1.2.2. Preprocesamiento

En general las investigaciones en el área realizan algún tipo de preprocesamiento a la consulta. Algunos realizan parafraseo para llevar la consulta a representaciones canónicas [?], en su lugar otros realizan un análisis morfológico léxico, con toquenización, lematización, eliminación de stop-words, tagueo de partes de la oración (*POS-tagging*), etc. [?]. También se encuentran las investigaciones que usan en esta etapa traducciones basadas en diccionarios especializados en el dominio, o de propósito general, al igual que ontologías, para "suavizar" vocablos difíciles de entender por el resto del sistema [?]. Además se usan técnicas como vectorización de palabras (*word2vector*), y transformación de la consulta a una representación en grafos [?].

1.2.3. Postprocesamiento

En la fase de Post-Procesamiento, se observan diversas enfoques en las investigaciones. La mayoría de los investigadores realizan un análisis semántico que implica la clasificación según tipos de datos, el uso de ontologías y bases de conocimiento externas para realizar mapeos [?] [?]. Algunos optan por convertir la consulta en una representación canónica, mientras que otros la transforman en un lenguaje intermedio antes de transpilarla al lenguaje objetivo [?]. También hay quienes la codifican directamente en forma de grafo, y algunos la convierten en embeddings [?].

1.2.4. Consulta

En la fase de construcción de la consulta, existen tres enfoques principales. El primero es un enfoque manual que implica la búsqueda y formateo de palabras clave como *"where"* y *"select"*, además del mapeo de atributos a tablas [?]. Otra vía se centra en el uso de modelos, incluyendo decodificadores y en algunos casos redes neuronales convolucionales [?]. También se emplea la construcción de la consulta formal mediante un compilador, especialmente cuando se ha utilizado un lenguaje intermedio entre el lenguaje natural y el formal de consulta[?].

1.3. Extracción de conocimiento mediante enfoques basados en técnicas de *prompt engineering* mediante LLMs

Con la creciente atención dada a los modelos de lenguaje a gran escala, estos se han convertido en un componente esencial en el procesamiento del lenguaje natural. A medida que aumenta el tamaño de los modelos preentrenados, también está cambiando gradualmente su uso. A diferencia de modelos como BERT [?] y T5 [?], que requieren un proceso de entrenamiento con una pequeña cantidad de datos, modelos como GPT-3 [?] requieren un diseño de un texto de entrada para generar resultados deseados. El reciente modelo de *ChatGPT* [?], que emplea Aprendizaje por Reforzamiento para la Retroalimentación Humana (RLHF) [?], simplifica el diseño de textos de entrada de calidad, lo que permite una mejor utilización de la capacidad de ZSL de modelos preentrenados a gran escala de manera conversacional. Debido a la sólida capacidad de dichos en la generación de código [?] y al hecho de que los modelos de generación de código suelen requerir una gran cantidad de datos anotados para producir buenos resultados [?], un modelo de generación de código de ZSL se considera fundamental.

Para la tarea específica de generar código de un lenguaje de consulta formal como *SQL* y *Cypher* se han utilizado distintos enfoques basados en varias de las principales técnicas de *prompt engineering*.

1.3.1. Zero-Shot Learning

Esta técnica se enfoca en la capacidad de un modelo para comprender y generar código en un lenguaje de consulta sin requerir ejemplos específicos de entrenamiento en ese lenguaje en particular. En otras palabras, el

modelo puede realizar esta tarea "desde cero", sin conocimiento previo del lenguaje. Algunos estudio interesantes se han realizado principalmente en la tarea de traducir lenguaje natural a lenguaje *SQL* [?] [?], los cuales permiten inferir la calidad mínima que estos modelos pueden alcanzar en la realización de dicha tarea [?].

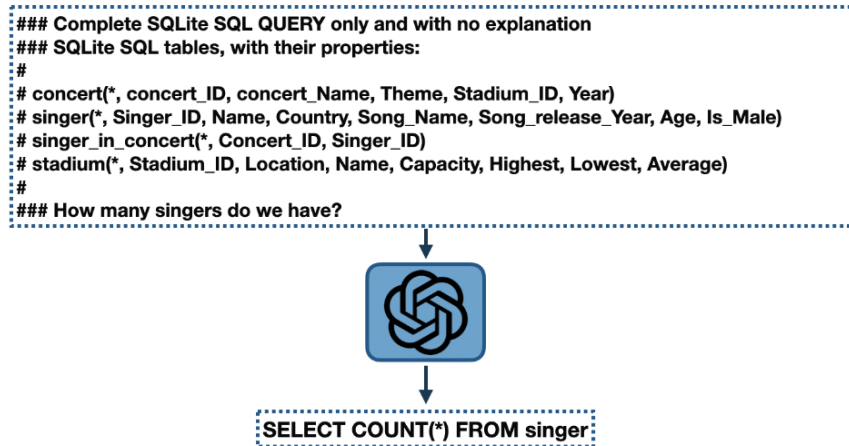


Figura 1.4: Ejemplo del flujo de trabajo de experimentos basados en ZSL para la traducción de lenguaje natural a código en *SQL*.

1.3.2. Few-Shot Learning

En contraste, el enfoque de Few-Shot Learning (FSL) se basa en la idea de que el modelo tiene acceso a un pequeño número de ejemplos (pocos ejemplos) en el lenguaje de consulta deseado para mejorar su capacidad de generar código en ese lenguaje. Esto puede ser especialmente útil cuando se necesita una adaptación rápida a un nuevo lenguaje o contexto. Tal y como muestran algunos resultados experimentales, este enfoque puede tener resultados superiores a varios modelos basados en *fine-tuning* [?].

Un ejemplo notable de esta aplicación es el modelo Codex [?], que ha demostrado ser un fuerte baseline en el *benchmark* Spider sin ninguna fase de entrenamiento. Además, se ha observado que proporcionar un pequeño número de ejemplos en el dominio en el prompt permite a Codex superar a los modelos de estado del arte cuyos parámetros han sido ajustados con pocos ejemplos de pocos dominios [?].

Además, se ha propuesto un marco de selección de demostraciones ODIS [?] que utiliza tanto ejemplos fuera de dominio como ejemplos ge-

nerados sintéticamente en el dominio para construir demostraciones. Este enfoque ha demostrado ser efectivo en comparación con los métodos de línea de base que se basan en una única fuente de datos [?].

En cuanto a los conjuntos de datos, existen varios conjuntos de datos de texto a SQL que han sido propuestos para evaluar el rendimiento de los modelos LLM. Algunos de estos conjuntos de datos incluyen CoSQL, TableQA, DuSQL, CHASE, y BIRD-SQL [?] [?] [?] [?] [?] [?].

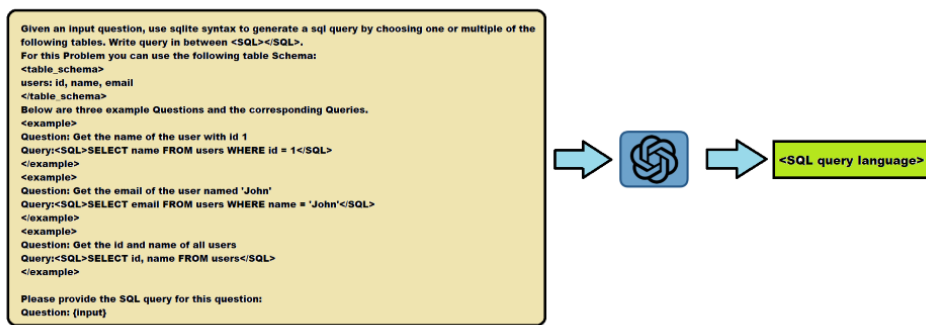


Figura 1.5: Ejemplo del flujo de trabajo de experimentos basados en FSL para la traducción de lenguaje natural a código en SQL.

1.3.3. Chain-of-Thought Prompting

El enfoque de la cadena de pensamiento (*Chain of Thought, CoT*) [?] en la conversión de texto a SQL ha demostrado ser una estrategia prometedor para mejorar la capacidad de los modelos de lenguaje de gran tamaño (LLMs) para realizar razonamientos complejos. Este enfoque se ha utilizado en varias investigaciones recientes para mejorar la capacidad de los LLMs para realizar tareas de razonamiento complejo, como la conversión de texto a SQL [?].

Un estudio propuso un nuevo paradigma para la generación de consultas SQL a partir de texto, llamado *Divide-and-Prompt*, que divide la tarea en sub tareas y luego aborda cada sub tarea a través de la cadena de pensamiento. Se presentaron tres métodos basados en la indicación para mejorar la capacidad de los LLMs para generar consultas SQL a partir de texto. Los experimentos mostraron que estas indicaciones guían a los LLMs para generar consultas SQL a partir de texto con mayor precisión de ejecución [?].

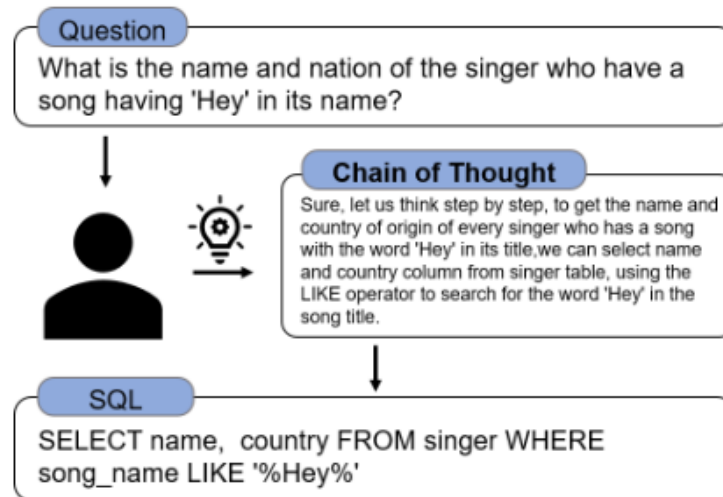


Figura 1.6: Ejemplo del flujo de trabajo de experimentos basados en CoT para la traducción de lenguaje natural a código en SQL.

1.3.4. Fine-Tuning

El entrenamiento (*fine-tuning*) de un modelo es un proceso que se utiliza para mejorar el rendimiento de un modelo de aprendizaje automático en una tarea específica. En el contexto de la conversión de texto a SQL, el ajuste fino puede ser utilizado para mejorar la precisión y la eficacia de los modelos de lenguaje de gran tamaño (LLMs) en la generación de consultas SQL a partir de texto. Se han realizado estudios [?]ropuso un método de *fine-tuning* para los modelos LLMs utilizando un conjunto de datos de texto a SQL, donde el método principal consistía en entrenar el modelo en el conjunto de datos de entrenamiento y luego ajustar el modelo en un conjunto de datos que contenía consultas SQL generadas por humanos. Este método demostró ser efectivo para mejorar la precisión de los modelos LLMs en la generación de consultas SQL a partir de texto [?].

Este enfoque constituye una poderosa herramienta en la conversión de lenguaje natural a SQL. Sin embargo, no es una solución mágica, ya que pocas organizaciones tienen conjuntos de datos de entrenamiento NL-2-SQL disponibles de manera inmediata. Algunos expertos consideran que las mejores arquitecturas se podrían lograr combinando modelos ajustados con agentes RAG (*Retrieval Augmented Generation*) [?].

1.4. Consideraciones generales

A pesar del uso comprobado de dichos enfoques, todavía existe una escasez de estudios orientados a la traducción de lenguaje natural a código de consulta a BDOG como por ejemplo *Cypher* [?], por lo tanto, es visible la necesidad de elaborar experimentos enfocados en dicha tarea, tomando como bases las ideas anteriormente expuestas. Debido a esto, resulta imprescindible desarrollar las primeras experimentaciones de la tarea en cuestión utilizando el enfoque ZSL, el cual en promedio permite obtener una cota inferior de qué tan efectivos pueden ser los LLMs con respecto a una tarea específica [].

Capítulo 2

Propuesta de Solución

En el presente capítulo se abordará la metodología seguida para diseñar el experimento propuesto en este trabajo ???. Primeramente, se expondrá un marco teórico que formaliza la definición del problema a tratar, esto con el objetivo de presentar los conocimientos base tenidos en cuenta para los enfoques probados. Luego, se detallan los primeros acercamientos desechados ??, argumentando las deficiencias de estos a la hora de arrojar resultados consistentes para la tarea que se desea desarrollar. Finalmente, se detalla la metodología definitiva a implementar, teniendo en cuenta la experiencia obtenida de las anteriores y mostrando su robustez para el análisis experimental [].

De forma general, el componente común para cada vía de solución constituye la presencia de un Gran Modelo de Lenguaje, pues representan los modelos más recientes utilizados para la tarea en cuestión; además, ofrecen resultados alentadores para el caso de traducción a lenguaje *SQL* según lo visto en la sección 1.3. Por lo tanto, tiene sentido probar su eficacia para traducir a código en *Cypher*, ya que ambos presentan similitudes como lenguajes formales declarativos para consultar bases de datos. Dicho modelo será analizado como una “caja negra” capaz de hacer tareas de traducción de lenguaje natural a una consulta semánticamente equivalente en el lenguaje *Cypher*.

Para cada vía de solución se deberá considerar el despliegue de un sistema de gestión de bases de datos para alguna BDOG, ya que es en este componente donde se almacenará la información a extraer por consultas en un lenguaje orientado a este tipo de almacenamiento. En el caso particular de este trabajo, se considerará el uso de *Neo4J*, con el cual se puede interactuar a partir del lenguaje *Cypher* ya mencionado. Por esto, es importante

considerar la implementación de un módulo intermedio para interactuar con una instancia del sistema de gestión *Neo4J*.

El enfoque de *prompt engineering* a utilizar será ZSL, por lo tanto, los textos de entrada que se le darán al modelo para la generación de código *Cypher* no contendrán ejemplos de pares de lenguaje natural con su correspondiente traducción al lenguaje de consulta objetivo. Por lo tanto, se tomarán algunas ideas experimentadas en el estado del arte para *SQL* vistas en la sección ??.

2.1. Definición del problema *Text-to-Cypher*

Dentro del ámbito de las bases de datos y las consultas que las acceden, existe una tarea particularmente compleja: traducir preguntas formuladas en lenguaje natural a un lenguaje de consulta estructurado como *Cypher*. Dada una pregunta Q , formulada en lenguaje cotidiano, y un esquema de base de datos S , el cual se compone a partir del tuplo N, A, R , donde encontramos múltiples nodos N (que representan instancias de una entidad en la base de datos), atributos C (tanto en nodos como en relaciones) y relaciones entre pares de nodos R . La problemática subyacente en el proceso de convertir *Text-to-Cypher* se centra en generar una consulta en lenguaje *Cypher* Y que sea equivalente y responda adecuadamente a la pregunta inicial Q realizada por un usuario humano.

2.1.1. LLM para resolver *Text-to-Cypher*

Con el auge de la inteligencia artificial y el aprendizaje automático, la tarea de convertir texto en lenguaje natural a código en *Cypher* ha sido recientemente abordada a través de técnicas modernas. En trabajos más recientes, algunos investigadores, [?] [?], han abogado por formular esta tarea como un desafío de generación. Utilizando lo que se conoce como "*prompts*." indicaciones P , es posible dirigir y guiar a un gran modelo de lenguaje M en esta labor. Este modelo, una vez entrenado, puede estimar una distribución de probabilidad sobre posibles consultas *Cypher* Y . De esta forma, el modelo es capaz de generar, paso a paso y token por token, una consulta apropiada.

La fórmula subyacente para generar la consulta Y se estructura como:

$$P_M(Y|P, S, Q) = \prod_{i=1}^{|Y|} P_M(Y_i|P, S, Q, Y < i)$$

Para simplificar, $Y < i$ se refiere al fragmento inicial o prefijo de la consulta *Cypher* que se está construyendo. Mientras que $P_M(Y_i|*)$ denota la probabilidad condicional asociada con la generación del i -ésimo token, considerando factores como el prefijo existente $Y < i$, la indicación P , el esquema S y la pregunta original Q .

Uno de los hallazgos más reveladores en el campo reciente es el concepto de aprendizaje en contexto (ICL, por sus siglas en inglés) [?], en el cual, grandes modelos de lenguaje pueden adaptarse y aprender de unos pocos ejemplos presentados en un contexto específico. Esta estrategia, defendida por varios investigadores [?] [?] [?], ha mostrado que los LLMs pueden abordar y dominar una amplia variedad de tareas complejas con una cantidad limitada de datos. Sin embargo, hay un equilibrio que mantener: agregar más ejemplos conlleva un aumento en los costos, tanto en términos de mano de obra para preparar esos ejemplos como en términos de costos de procesamiento y tokens al interactuar con APIs avanzadas como la de OpenAI [?]. En este estudio, el foco recae en trabajar eficientemente con indicaciones al modelo de lenguaje sin requerir ejemplos adicionales.

2.2. Enfoques considerados

El desarrollo de un sistema con un LLM (*Large Language Model*) para hacer inferencias de *Zero-Shot* y generar *Cypher* implica varios enfoques y desafíos. A continuación, se detallan los enfoques considerados y las dificultades encontradas.

Primero, se exploraron modelos de código abierto. El primer modelo probado fue Alpaca Lora 7-B [?]. Sin embargo, la calidad de las inferencias no fue la adecuada. En algunos casos, las respuestas eran *Cypher* no compilable con errores de sintaxis. En otros casos, el modelo generaba otro lenguaje de consulta que no era *Cypher*, como por ejemplo *SQL*. Posteriormente, se probó con GPT4A11 [?] y con Vicuna-7B [?], obteniendo resultados similares. Estos modelos, aunque útiles, no proporcionaban la precisión y la generación de *Cypher* que se requería para el sistema.

Además, se consideró la posibilidad de entrenar un propio modelo de lenguaje grande con muchos parámetros para tareas específicas. Este modelo podría ser capaz de generar *Cypher* sin haber sido entrenado específicamente para eso. Sin embargo, esta opción fue descartada debido a la inviabilidad total. Entrenar un modelo de lenguaje grande requiere una gran cantidad de recursos de computación [?]. Además, la tarea de entrenar un

modelo para generar *Cypher* sin haber sido entrenado específicamente para eso es extremadamente desafiante. Aunque técnicamente posible, requeriría una inversión significativa de tiempo y recursos, así como un conocimiento profundo de los modelos de lenguaje y el aprendizaje profundo.

En resumen, el desarrollo de un sistema con un LLM para hacer inferencias de *Zero-Shot* y generar *Cypher* implica una serie de desafíos. Estos incluyen la elección del modelo correcto, la necesidad de entrenar el modelo para generar *Cypher* específicamente, y la inversión significativa de recursos y tiempo. A pesar de estos desafíos, el uso de LLMs ofrece la promesa de generar *Cypher* de manera eficiente y precisa, lo que puede ser de gran utilidad en una variedad de aplicaciones.

2.3. Propuesta de solución diseñada

En el desarrollo de sistemas de bases de datos, la interacción eficiente y la gestión de esquemas son fundamentales para la manipulación y el mantenimiento de datos. En este contexto, se ha desarrollado una arquitectura de software compuesta por varios componentes interconectados diseñados para interactuar con una base de datos de grafos, con especial atención en *Neo4J*, un sistema de manejo de bases de datos basado en grafos.

Uno de los componentes clave de esta arquitectura es el *GraphContractor*, un módulo diseñado para facilitar la interacción con la base de datos. Este componente actúa como intermediario entre la aplicación y la base de datos, manejando la lógica necesaria para establecer conexiones, ejecutar consultas y manejar los resultados. La modularidad del *GraphContractor* permite su reutilización y fácil mantenimiento, además de proporcionar una capa de abstracción que simplifica las operaciones de la base de datos para los desarrolladores.

Para la generación de esquemas de la base de datos, se ha creado el *SchemaMaker*, una herramienta automatizada que se encarga de construir los esquemas necesarios para *Neo4J*. Esta herramienta juega un rol crucial en la estructuración de la base de datos, ya que define la organización de nodos, relaciones, propiedades y restricciones. El *SchemaMaker* garantiza que la base de datos esté correctamente configurada para cumplir con los requisitos del dominio y las necesidades de la aplicación, asegurando así la integridad y coherencia de los datos.

En la interfaz entre el lenguaje natural y la base de datos, se ha integrado el modelo de lenguaje GPT-4, utilizado como una "caja negra" para la traducción de consultas en lenguaje natural a *Cypher*, el lenguaje de con-

sulta para *Neo4J*. La capacidad de GPT-4 para comprender y generar texto hace posible que los usuarios realicen consultas complejas sin necesidad de conocer la sintaxis específica de *Cypher*. Para mejorar la precisión y relevancia de las traducciones, se ha elaborado una plantilla de prompt que se nutre de la salida del SchemaMaker. Esta plantilla guía al modelo de lenguaje proporcionando contexto y estructura, lo que permite que GPT-4 genere consultas *Cypher* más precisas y eficientes.

Finalmente, se ha desarrollado el DBSeeder, un componente encargado de poblar la base de datos de *Neo4J* con datos iniciales o de prueba. Utilizando consultas *Cypher* generadas por el modelo de lenguaje GPT-4, el DBSeeder trabaja en conjunto con el GraphContractor para insertar nodos, atributos y relaciones en la base de datos. Esta funcionalidad es especialmente valiosa en las etapas de desarrollo y prueba, donde se requiere de una base de datos poblada para validar el diseño y la lógica de la aplicación. Además, el DBSeeder está diseñado para operar dentro de un contenedor de *Docker*, lo que ofrece ventajas significativas en términos de portabilidad, escalabilidad y aislamiento del entorno de desarrollo.

Cada uno de estos componentes representa un eslabón en la cadena de herramientas que permitirán interactuar con la base de datos de grafos de manera más intuitiva y automatizada. La integración de tecnologías avanzadas como GPT-4 en el proceso no solo mejora la accesibilidad para los usuarios finales sino que también agiliza el ciclo de desarrollo, ofreciendo un enfoque moderno y eficiente en la gestión de bases de datos de grafos como *Neo4J*.

2.3.1. Interacción con una base de datos *Neo4J*

El componente GraphContractor actúa como un facilitador o intermediario entre el usuario y la base de datos *Neo4J*. Su propósito es simplificar las tareas de conexión y ejecución de consultas contra la base de datos, manejando internamente los detalles de la comunicación y posibles excepciones.

Al instanciar GraphContractor, se le proporciona una URL, junto con un nombre de usuario y contraseña para la autenticación. Luego, esta herramienta intenta establecer una conexión con la instancia de la base de datos *Neo4J* en la URL especificada. Si la conexión es exitosa, el GraphContractor estará listo para ejecutar consultas. Si la conexión falla, por ejemplo, debido a problemas con la red o las credenciales de acceso, se informará al usuario al respecto con un mensaje informativo.

Una vez que el GraphContractor está conectado a la base de datos, es

posible realizar consultas a una base de datos objetivo. Para dicha tarea, este acepta una cadena de texto que representa la consulta *Cypher* a ejecutar. Al efectuar dicha funcionalidad con una consulta de *Cypher* válida, GraphContractor la ejecutará en la base de datos y devolverá los resultados. Si ocurre algún error durante la ejecución de la consulta, como una sintaxis incorrecta de *Cypher* o un problema de conexión, el error se capturará y se presentará al usuario, ofreciendo retroalimentación inmediata.

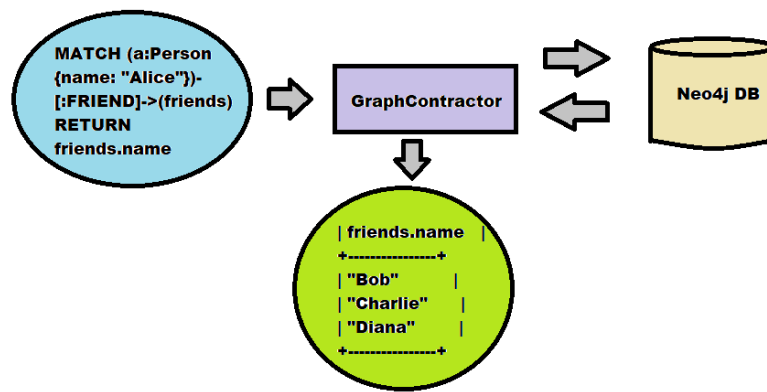


Figura 2.1: Flujo de funcionamiento del componente GraphContractor.

En resumen, GraphContractor encapsula la complejidad de la gestión de la conexión a la base de datos y la ejecución de consultas, proporcionando una interfaz simplificada para interactuar con *Neo4J*. Esto permite a los usuarios centrarse en la lógica de sus consultas y manejo de datos, en lugar de en los detalles subyacentes de la implementación de la base de datos.

2.3.2. Extracción de información de una base de datos *Neo4J*

En el ámbito de la extracción de información de bases de datos orientadas a grafos, como *Neo4J*, se ha desarrollado un componente denominado KnowledgeBase. Este componente actúa como una interfaz avanzada para la interacción con dichas bases de datos, aprovechando las capacidades del componente GraphContractor. La funcionalidad principal de KnowledgeBase radica en su habilidad para encapsular consultas predefinidas en el lenguaje *Cypher*, el cual es específico para bases de datos *Neo4J*.

El diseño de KnowledgeBase tiene como objetivo principal facilitar la extracción de información de manera eficiente y transparente. Para lograr

esto, utiliza consultas en *Cypher* que están integradas dentro de sus funcionalidades. Estas consultas son ejecutadas a través de una instancia interna de *GraphContractor*, proporcionando una capa de abstracción que simplifica las interacciones del usuario con la base de datos.

Sus funcionalidades principales implementadas fueron:

- **Inicialización y Configuración:** El proceso de inicialización establece una conexión vital con la base de datos y prepara el componente para operaciones de consulta. Esto implica la integración con un sistema central que maneja las interacciones con la base de datos.
- **Verificación de Existencia de Entidades:** Una funcionalidad central permite verificar la presencia de entidades específicas en la base de datos. Se emplean etiquetas y propiedades para formular consultas que determinan la existencia de la entidad.
- **Comprobación de Atributos en Entidades:** Otra capacidad importante es determinar si una entidad posee un atributo específico. Esta función es clave para validar la integridad y completitud de los datos.
- **Extracción y Análisis de Entidades y Atributos:** La extracción de etiquetas de entidades provee una visión general de los tipos de datos almacenados. Adicionalmente, se realiza un análisis detallado de los atributos asociados con cada entidad y relación, incluyendo el tipo de dato y los rangos de valores.
- **Inferencia del Tipo de Dato:** La habilidad para identificar el tipo de dato de un valor proporcionado es esencial para el manejo adecuado de los datos, permitiendo realizar conversiones y validaciones de tipo cuando sea necesario.
- **Identificación de Claves y Atributos:** Se realizan operaciones para extraer claves de entidades y relaciones. Esto proporciona información detallada sobre los campos disponibles en diferentes tipos de nodos y enlaces.
- **Análisis de Relaciones entre Entidades:** Identificar y catalogar las relaciones entre distintas entidades es crucial para comprender las interacciones y conexiones dentro de la base de datos.

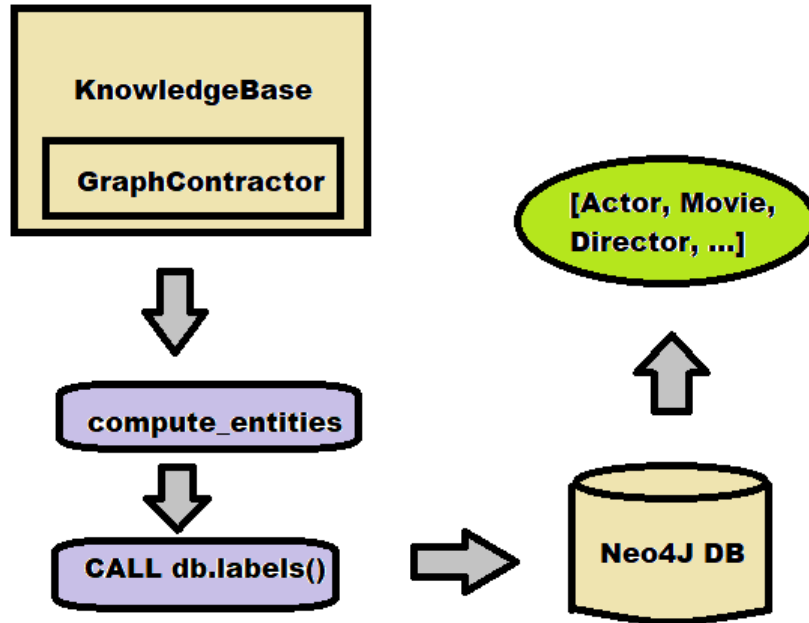


Figura 2.2: Flujo de funcionamiento del componente KnowledgeBase.

2.3.3. Almacenamiento de información en una base de datos *Neo4J*

El componente DBSeeder es una herramienta diseñada para cargar datos en una base de datos en forma de grafo. Su propósito es automatizar el proceso de toma de datos estructurados y su inserción en la base de datos, creando nodos y relaciones entre ellos según se define en los datos de entrada. La interacción de este con una instancia de una base de datos *Neo4J* es a través de una instancia del componente KnowledgeBase visto en la sección 2.3.2.

Al inicializar esta herramienta, se le proporcionan dos piezas de información esenciales: un conjunto de conocimientos que describe la estructura de la base de datos y una ruta a un archivo que contiene los datos a ser sembrados en la base de datos. Estos datos de entrada son esenciales para guiar el proceso de sembrado.

Una vez configurada, la herramienta tiene la capacidad de procesar los datos de entrada. Esto se realiza leyendo cada línea del archivo de datos, donde cada línea representa un conjunto de información que debe ser transformada en elementos dentro de la base de datos. La herramienta ana-

liza cada línea para comprender y aislar las partes que corresponden a entidades y las relaciones entre ellas.

Para cada conjunto de datos, la herramienta verifica si los elementos ya existen en la base de datos. Si no es así, procede a crear nuevos nodos que representan entidades y luego establece relaciones entre estos nodos, basándose en la relación especificada en los datos.

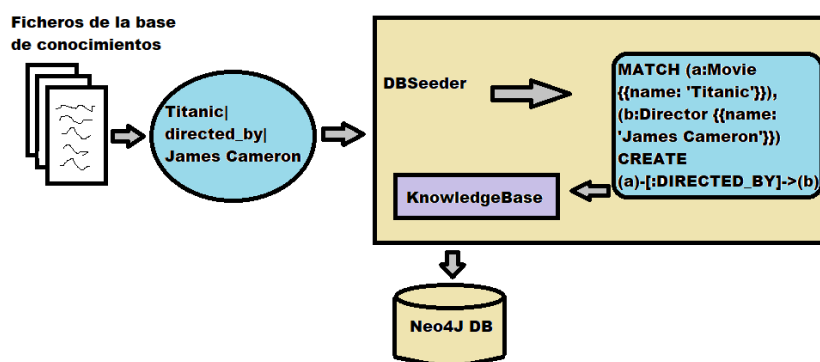


Figura 2.3: Flujo de funcionamiento del componente DBSeeder.

El proceso de llenado en la base de datos asegura que no se introduzcan duplicados y que los datos se estructuren correctamente en la base de datos de acuerdo con sus reglas y definiciones. Al finalizar, la base de datos debería reflejar una red de nodos interconectados que representan tanto las entidades como las relaciones definidas en el archivo de datos original. Este proceso es fundamental para preparar la base de datos para su uso en aplicaciones que requieren acceso a datos relacionales y estructurados en forma de grafo.

2.3.4. Selección del modelo

La selección de GPT-4 para la generación de código a partir de lenguaje natural utilizando aprendizaje *Zero-Shot* (ZSL) está justificada por varias razones fundamentadas en investigaciones y comparaciones técnicas recientes. GPT-4 es un avance significativo respecto a sus predecesores, construido sobre la arquitectura de GPT-3 pero alcanzando nuevos niveles de

rendimiento y escala []. Este modelo mejora en la corrección factual de las respuestas y reduce las "alucinaciones", donde el modelo comete errores de hecho o razonamiento, obteniendo un 40% más de precisión que GPT-3.5 en las pruebas de rendimiento factual internas de OpenAI [].

El modelo GPT-4 se basa en la arquitectura Transformer, que utiliza mecanismos de atención para procesar texto, y se ha mejorado con una mezcla de expertos (MoE) para lograr un modelo con aproximadamente 1,76 billones de parámetros, un orden de magnitud mayor que GPT-3 []. Además, estudios recientes han mostrado que GPT-4 supera a GPT-3.5 en aprendizaje zero-shot en casi todas las tareas evaluadas, lo que incluye una variedad de dominios de razonamiento como deductivo, inductivo, abductivo, analógico, causal y multi-salto, a través de tareas de preguntas y respuestas [].

Además, el modelo GPT-4 emplea técnicas de *fine-tuning* y Aprendizaje por Refuerzo con Retroalimentación Humana (RLHF), lo que le permite ser un modelo multimodal robusto capaz de procesar entradas textuales y visuales y generar salidas basadas en texto []. Este enfoque ha demostrado ser eficaz en la mejora de las capacidades de razonamiento de los modelos de lenguaje grandes (LLMs), lo que lo hace especialmente adecuado para tareas complejas que requieren razonamiento, como la traducción de lenguaje natural a lenguaje de consulta formal [].

Cuadro 2.1 Rendimiento de GPT-4 en referencias académicas. []

	GPT-4	GPT-3.5	LM SOTA	SOTA
MMLU	86.4%	70.0%	70.7%	75.2%
HellaSwag	95.3%	85.5%	84.2%	85.6%
AI2 Reasoning Challenge (ARC)	96.3%	85.2%	85.2%	86.5%
WinoGrande	87.5%	81.6%	85.1%	85.1%
HumanEval	67.0%	48.1%	26.2%	65.8%
DROP	80.9%	64.1%	70.8%	88.4%
GSM-8K	92.0%	57.1%	58.8%	87.3%

En comparación con otros modelos como PALM, Chinchilla, LaMDA, LLaMA y Gopher, que también han evaluado sus habilidades de razonamiento, GPT-4 se destaca por su capacidad mejorada de aprendizaje zero-shot y por el uso de estrategias de prompting refinadas para mejorar aún más su rendimiento en tareas de razonamiento, lo que lo convierte en una elección

prometedora para la generación de código [1].

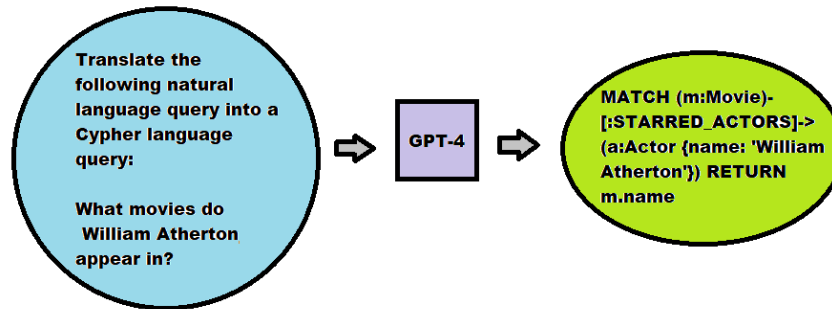


Figura 2.4: Ejemplo básico sobre cómo utilizar GPT-4 para traducir lenguaje natural en lenguaje de consulta *Cypher*.

2.3.5. Diseño de la información entrada al LLM

Tal y como se mostró en la sección anterior, el uso de GPT-4 para la tarea de traducción no resulta complicado, pues este es utilizado como una gran “caja negra” capaz de realizar operaciones que se le indiquen en un texto (*prompt*) de entrada. Debido a esto, también es importante mencionar que no cualquier entrada es efectiva o suficiente para obtener el resultado esperado [?]. Al proceso de diseñar una entrada de calidad para que el modelo realice una tarea específica exitosamente se denomina *prompt engineering* [7].

Dado que la hipótesis central de este trabajo se basa en el uso del aprendizaje *Zero-Shot* (ZSL), el texto de entrada al modelo GPT-4 no puede contener ejemplos de cómo traducir una consulta en lenguaje humano a lenguaje *Cypher*, es decir, no deberá reflejar contenido demostrativo de la tarea a realizar, lo cual se justifica por la misma definición del enfoque ZSL [?]. Además, como parte de la información de entrada al modelo para este tipo de tareas, es común añadir una descripción de la estructura de la base de datos a consultar [?] [?], lo cual se conoce como esquema de la base de datos [?].

Por lo mencionado anteriormente, el texto de entrada al modelo deberá contener:

- **La tarea a ejecutar:** Se le describirá al modelo la tarea a realizar, mencionando los datos que recibirá de entrada y el formato en que se desea obtener la respuesta.

- **Esquema de la base de datos:** Se especificará el contenido de la base de datos objetivo en forma de grafo, mencionando las entidades, relaciones (especificando si son en una sola o ambas direcciones entre un par de entidades) y atributos presentes en la misma (correspondientes tanto a las entidades como a las relaciones entre estas).
- **Consulta en lenguaje natural:** En este caso, se añadirá la consulta en lenguaje natural humano a traducir.

Para la obtención del esquema de la base de datos de tipo *Neo4J* se diseñó el componente *SchemaMaker*. Con el fin de elaborar la descripción de la estructura de la fuente de datos objetivo, este recibe los nombres de las entidades, las relaciones y los atributos presentes en la misma. Dicha información la obtiene auxiliándose del componente *KnowledgeBase* analizado en la sección 2.3.2, mediante el cual se realizan las consultas en lenguaje *Cypher* pertinentes a la instancia de la fuente de datos *Neo4J* en cuestión. A continuación se muestra un ejemplo del funcionamiento de la herramienta *SchemaMaker*:

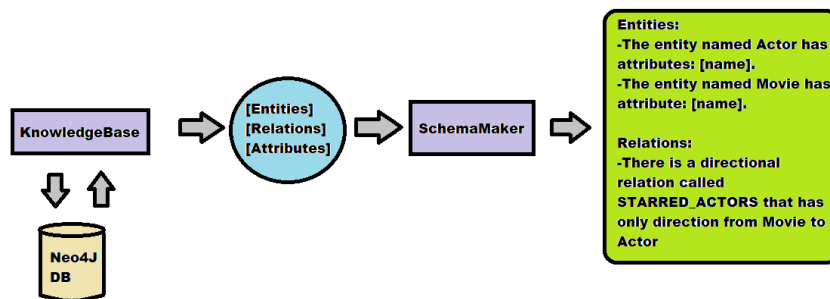


Figura 2.5: Ejemplo del proceso que se realiza para obtener un formato verbalizado de una base de datos *Neo4J* con el uso del *SchemaMaker*.

Finalmente, el texto de entrada para el modelo a utilizar mencionado en la sección 2.3.4 se integraría de la siguiente manera al proceso de traducción de lenguaje natural a lenguaje *Cypher*:

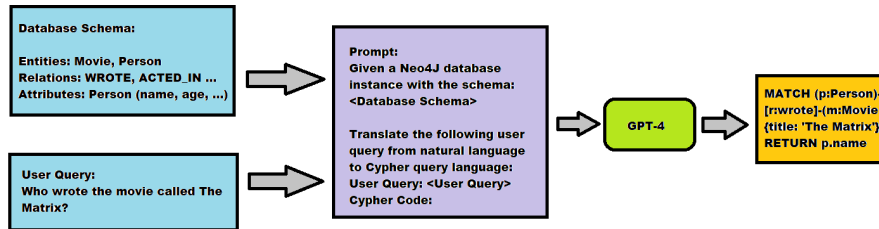


Figura 2.6: Flujo de entrada y salida en el proceso de traducción de lenguaje natural a código en *Cypher*.

Tal y como se muestra en la imagen anterior, junto con la entrada de una consulta en lenguaje natural se elabora un texto de entrada al modelo utilizando además, el esquema de la base de datos *Neo4J* a consultar, la cual como ya se mencionó en esta subsección, es producida por el SchemaMaker.

2.3.6. Caso de estudio

A partir de los contenidos abordados en esta sección, resulta importante mostrar la arquitectura general del sistema diseñado, así como el flujo de funcionamiento de la misma.

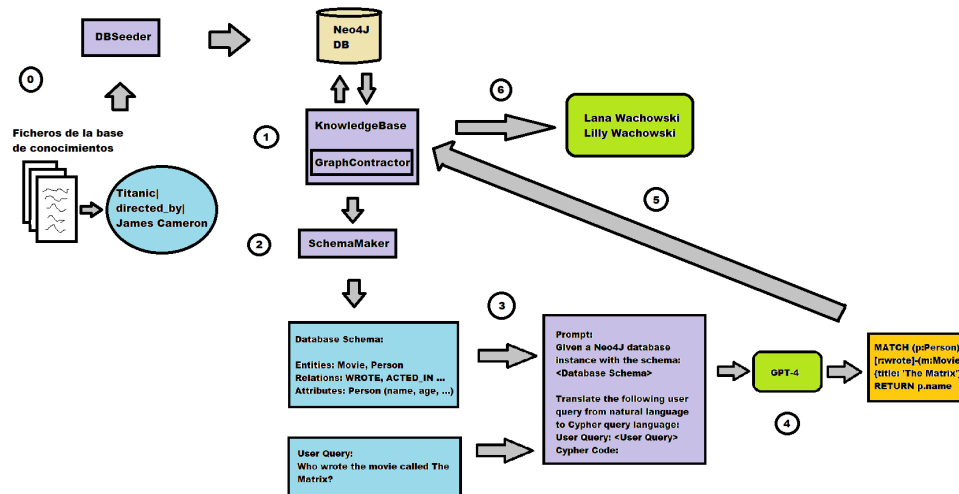


Figura 2.7: Arquitectura y funcionamiento de la propuesta de solución.

Para una mejor comprensión de la figura 2.3.6 se añadieron un conjunto de índices que resaltan las distintas fases por las que pasa el sistema implementado. A continuación se enumeran y explican en qué consisten cada una de estas etapas:

1. **Paso 0:** Este representa el proceso mediante el cual se construye una instancia de una base de datos *Neo4J* a partir de un conjunto de ficheros de texto. Esta tarea se puede llevar a cabo mediante el componente *DBSeeder*, el cual puede programarse con funcionalidades específicas de acuerdo al formato de los ficheros iniciales y los datos que estos contienen. Como se mencionó en la sección 2.3.3, el componente *DBSeeder* se apoya del componente *GraphContractor* internamente para realizar peticiones a la base de datos con el objetivo de añadir nuevos registros de información, traduciendo la creación de entidades, relaciones y atributos en instrucciones de *Cypher* ejecutadas sobre una base de datos *Neo4J* objetivo.
2. **Paso 1:** En esta fase, el componente *KnowledgeBase* extrae la información referente a las entidades, relaciones y atributos de una base de datos *Neo4J* existente. Para ello, hace uso internamente de una instancia de un *GraphContractor* 2.3.1.
3. **Paso 2:** En este paso, se utiliza la herramienta *SchemaMaker* 2.3.5, el

cual recibe la información de la base de datos procedente del componente KnowledgeBase y produce una descripción verbalizada y legible en lenguaje natural sobre la estructura de la instancia de *Neo4J* objetivo.

4. **Paso 3:** En este momento del proceso general, se recibe una consulta descrita en lenguaje natural humano sobre una solicitud de datos de la base de datos objetivo. Luego esta es utilizada junto con el esquema de la base de datos obtenido del SchemaMaker para conformar un texto de entrada al modelo GPT-4 con las características mencionadas en la sección 2.3.5.
5. **Paso 4:** Una vez obtenido el texto de entrada para realizar una inferencia con el modelo GPT-4, se procede a hacer una ejecución de este, produciendo así como salida un texto que representa un código en lenguaje *Cypher*.
6. **Paso 5:** En esta etapa, se utiliza la salida del modelo GPT-4 para ser ejecutada sobre la base de datos *Neo4J* objetivo mediante la herramienta KnowledgeBase.
7. **Paso 6:** Finalmente, se obtiene la respuesta procedente de la base de datos *Neo4J* con la información solicitada.

Con lo anteriormente explicado se expone un ejemplo de caso de uso donde se tienen como entradas al sistema una instancia de una base de datos *Neo4J* y una consulta en lenguaje humano para obtener como salida final el conjunto de datos extraídos de dicha base de datos objetivo correspondientes a la solicitud textual dada sobre estos.

Capítulo 3

Detalles de Implementación

En este capítulo se exponen los detalles sobre la implementación de los componentes que intervinieron en el proyecto, resaltando los principales aspectos de programación tenidos en cuenta.

El sistema implementado fue desarrollado con el *IDE VSCode* [?], el cual fue configurado para utilizar el lenguaje *Python* (versión 3.9) [?] y contar con las facilidades que esta herramienta ofrece para detectar errores de sintaxis y de dependencias en tiempo de compilación. Por otro lado, el sistema operativo utilizado para implementar el proyecto fue *Ubuntu-22.04*, con el cual se facilitó el proceso de instalación de bibliotecas para el lenguaje de programación utilizado.

En general, no se optó por utilizar una arquitectura de software específica, ya que el objetivo del proyecto se centra fundamentalmente en utilizar el sistema resultante para evaluar los experimentos orientados a responder la hipótesis de estudio y no en desarrollar una aplicación a ser llevada a producción para emplearse por usuarios humanos.

3.1. Despliegue de una instancia de una base de datos *Neo4J*

Para desplegar una instancia de una base de datos en forma de grafo de tipo *Neo4J* nos apoyamos en *Docker* [?], con el cual se desplegó un contenedor que pudiese ejecutar el sistema de gestión de *Neo4J* y ser accesible ante peticiones de una aplicación mediante el protocolo *Bolt* [?].

Una vez instalado *Docker* se utilizó el siguiente comando:
este comando ejecuta un contenedor Docker con la imagen de *Neo4J*,

```

1 docker run \
2 --name testneo4j \
3 -p 7474:7474 -p 7687:7687 \
4 -d \
5 -v $HOME/neo4j/data:/data \
6 -v $HOME/neo4j/logs:/logs \
7 -e NEO4J_AUTH=neo4j/testpassword \
8 neo4j

```

Figura 3.1: Comando utilizado para desplegar una base de datos *Neo4J*.

expone los puertos 7474 y 7687, guarda los datos y los registros en el *host*, y establece las credenciales de autenticación para la interfaz de usuario de *Neo4J*:

- **--name testneo4j**: Esto asigna el nombre testneo4j al contenedor que se está.
- **-p 7474:7474 -p 7687:7687**: Estas opciones mapean los puertos 7474 y 7687 del contenedor a los mismos puertos del *host*. Esto permite acceder a los servicios que se ejecutan en estos puertos dentro del contenedor desde el *host*.
- **-d**: Esta opción hace que el contenedor se ejecute en segundo plano (modo *detached*).
- **-v \$HOME/neo4j/data : /data -v \$HOME/neo4j/logs:/logs**: Estas opciones montan los directorios *\$HOME/neo4j/data* y *\$HOME/neo4j/logs* del *host* en los directorios */data* y */logs* del contenedor, respectivamente. Esto permite que los datos y los registros generados por el contenedor se guarden en el *host*.
- **-e NEO4J_AUTH=neo4j/testpassword**: Esta opción establece la variable de entorno *NEO4J_AUTH* en el contenedor con el valor *neo4j/testpassword*. Esto se utiliza para configurar la autenticación en *Neo4J*.
- **neo4j**: Esta es la imagen de la que se está creando el contenedor. En este caso, se está utilizando la imagen oficial de *Neo4J*.

3.2. GraphContractor

Para la implementación del componente GraphContractor se diseñó una clase para interactuar con una instancia de base de datos *Neo4J*. La clase hereda de la clase Graph del módulo *py2neo*, que proporciona una interfaz de alto nivel para interactuar con bases de datos *Neo4J*.

```

1 class GraphContractor(Graph):
2     """
3         Graph Contractor class for interacting with a Neo4J DB
4         instance
5     """
6     def __init__(self, url, name, password):
7         try:
8             self.graph = Graph(url, auth=(name, password))
9
10        except Exception as e:
11            print(e)
12            print('Error connecting to the database(Remember VPN)')
13
14    def make_query(self, query: str):
15        try:
16            return self.graph.run(query).data()
17        except BaseException as e:
18            print(e)
19            return str(e)

```

Figura 3.2: Implementación de la clase GraphContractor.

La clase GraphContractor tiene un método `__init__` que se utiliza para inicializar una nueva instancia de la clase. Este método toma tres argumentos: `url`, `name` y `password`, que se utilizan para establecer una conexión con la base de datos *Neo4J*. Esta clase también tiene un método `make_query` que se utiliza para ejecutar consultas en la base de datos *Neo4J*. Este método toma una consulta en formato de cadena y la ejecuta en la base de datos.

3.3. KnowledgeBase

El componente KnowledgeBase fue implementado en una clase cuyo constructor recibe una instancia de una entidad de tipo GraphContractor:

```

1 class KnowledgeBase:
2     def __init__(self, graph: GraphContractor) -> None:
3         self.graph = graph
4
5     def entity_exists(self, label, property_name, property_value):
6         ...
7     def entity_has_attribute(self, label, property_name,
8         entity_name):
9         ...
10    def compute_entities(self):
11        ent = self.graph.graph.run('CALL db.labels()')
12        entities = QueryUtils._unfold_graph_resp(ent)
13        return entities
14
15    def compute_attributes(self, entities, relations):
16        ...
17    def compute_relations(self, entities):
18        ...
19    def _infer_data_type(self, value):
20        ...
21    def get_type_min_max_entity_attribute(self, entity,
22        attribute_name):
23        ...
24    def get_type_min_max_relation_attribute(self, relation,
25        attribute_name):
26        ...
27    def get_keys_of_label(self, label):
28        ...
29    def get_keys_of_relation(self, relation):
30        ...

```

Figura 3.3: Implementación de la clase KnowledgeBase.

Esta herramienta presenta un conjunto de métodos auxiliares para di-

versas tareas que impliquen la extracción de información de una base de datos, donde cada uno internamente utiliza la instancia del `GraphContractor` proporcionado. Por ejemplo, tal como se muestra en la figura 3.3, para obtener el conjunto de entidades de la base de datos se ejecuta un código de *Cypher* correspondiente sobre la instancia de *Neo4J* a consultar.

3.4. DBSeeder

La clase `DBSeeder` en el código proporcionado en la figura 3.4 se utiliza para poblar (*seed*) una base de datos con información de una base de datos base y una base de conocimientos de tipo *Neo4J*.

```

1 class DBSeeder:
2     def __init__(self, kb: KnowledgeBase, db_base_file_path: str) ->
        None:
3         self.kb = kb
4         self.db_base_file_path = db_base_file_path
5
6     def seed_db(self):
7         ...
8     def create_entities_relations_attributes(self, entity1,
        relation_type, entity2):
9         ...

```

Figura 3.4: Implementación de la clase `DBSeeder`.

El método constructor `__init__` recibe dos argumentos: `kb`, que es una instancia de `KnowledgeBase`, y `db_base_file_path`, que es la ruta del archivo de la base de datos base. Por otro lado, el método `seed_db` se encarga de almacenar información en la base de datos objetivo. Este método abre el el conjunto de ficheros que contienen información de la base de datos para su lectura y luego itera sobre cada línea del archivo, llamando al método `create_entities_relations_attributes`, el cual verifica si las entidades `entity1` y `entity2` existen en la base de datos. Si no existen, se crean. Luego, se crea una relación entre el par de entidades dadas y la relación especificada.

3.5. SchemaMaker

En la figura 3.6 se muestra la implementación de la clase SchemaMaker, la cual tiene el método estático compute_schema_description.

```

1 class SchemaMaker:
2     @staticmethod
3     def compute_schema_description(entities, relations, attributes):
4         schema_description = ""
5         schema_description += f"Entities: {entities}\n"
6         for entity in entities:
7             entity_attrs = attributes[entity]
8             if len(entity_attrs) > 0:
9                 schema_description += f"The entity named {entity} has
                                the attributes: {[attr[0] for attr in
                                entity_attrs]}\n"
10        for relation in relations:
11            for ent1, ent2, is_double_sense in relations[relation]:
12                if is_double_sense:
13                    schema_description += f"There is a relation
                                called {relation} between the entitites
                                {ent1} and {ent2}. The relation {relation}
                                can be used in both senses.\n"
14                    continue
15                    schema_description += f"There is a directional
                                relation called {relation} that has only
                                direction from {ent1} to {ent2}.\n"
16                if len(attributes[relation]) > 0:
17                    schema_description += f"The relation {relation} has
                                attributes {attributes[relation]}\n"
18        return schema_description

```

Figura 3.5: Implementación de la clase Schema Maker.

Dicho método toma tres argumentos: entities, relations y attributes, los cuales utiliza para generar una descripción del esquema de una base de datos. Este comienza inicializando una cadena vacía schema_description y luego agrega información sobre las entidades y relaciones. Primero, agrega una línea que indica las entidades presentes en el esquema. Luego, para

cada entidad, si tiene atributos, agrega una línea que indica los atributos de esa entidad. Después de manejar las entidades, el método pasa a las relaciones. Para cada relación, si es de doble sentido, agrega una línea que indica que existe una relación bidireccional entre las dos entidades involucradas. Si no es de doble sentido, agrega una línea que indica que existe una relación unidireccional desde la primera entidad a la segunda. Finalmente, si la relación tiene atributos, agrega una línea que indica los atributos de la relación.

3.6. GPT-4

Para implementar una estructura capaz de realizar inferencias a partir del modelo GPT-4 se utilizó la biblioteca *Langchain* y se definió una función `get_model` que se utiliza para inicializar un modelo de lenguaje basado en el tipo de modelo especificado. Fue necesario además, el uso de una `API_KEY` de *OpenAI* con el objetivo de acceder a los modelos disponibles [?].

La función `get_model` toma dos argumentos: `model_type` y `model_name`. La variable `model_type` puede ser `chat` o cualquier otro tipo de modelo, y `model_name` es el nombre del modelo específico que se va a utilizar. Para el caso específico de este trabajo el tipo de modelo fue `chat` y el nombre utilizado fue `gpt-4`, mientras que la temperatura del modelo utilizada fue 0,7.

Primero, la función inicializa una plantilla de *prompt* utilizando de acuerdo a si el modelo a utilizar es de tipo `chat` o generativo. La plantilla de *prompt* se inicializa con un texto predefinido que describe la tarea del modelo y los marcadores de posición para el lenguaje de consulta, el tipo de base de datos, el esquema y la consulta. Luego, la función inicializa un modelo de lenguaje basado en `model_type`. Ambos modelos se inicializan con un nombre de modelo y una temperatura. Finalmente, la función inicializa una instancia de un modelo capaz de hacer inferencias a partir de un texto de entrada predefinido con variables.

```

1 from langchain.prompts import PromptTemplate, ChatPromptTemplate
2 from langchain.chat_models import ChatOpenAI
3 from langchain import LLMChain
4 from langchain import OpenAI
5
6 # template for the model
7 template = """
8 You are an agent capable of transforming natural language queries
   to queries in the query language {query_language}. Your task
   is: Given a database schema of type {database_type} and a query
   written in human natural language, return only the code to
   answer that query in the query language {query_language} and
   respect the relations directions.
9
10 The database schema is: {schema}
11
12 The natural language query is: {query}
13
14 The code in the query language {query_language} is:
15
16 """.strip()
17
18 def get_model(model_type, model_name):
19     # Init prompt template
20     prompt = ChatPromptTemplate.from_template(template=template) if
       model_type == "chat" else PromptTemplate(template=template,
       input_variables=[
21         "query_language", "database_type", "schema", "query"])
22
23     # Init llm
24     llm = ChatOpenAI(model=model_name, temperature=0.7) if
       model_type == "chat" else OpenAI(temperature=0.7)
25
26     # Init chain
27     llm_chain = LLMChain(prompt=prompt, llm=llm)
28
29     return llm_chain

```

Figura 3.6: Implementación para utilizar el modelo GPT-4.

Capítulo 4

Análisis Experimental

En este capítulo se presentan los marcos experimentales utilizados para evaluar la efectividad del sistema propuesto en el capítulo 2 para la traducción de una consulta en lenguaje natural al lenguaje de consulta formal *Cypher*. Cada enfoque utilizado consistió en el uso de un conjunto de tuplas que contenían común una consulta en lenguaje natural de ejemplo a traducir hacia un segundo elemento correspondiente con un objetivo a medir en la traducción.

Todos los experimentos fueron ejecutados en un servidor privado virtual (VPS) con sistema operativo *Ubuntu-20.04*, memoria RAM de 16Gb, una CPU AMD basada en la arquitectura *x86_64*, con 8 núcleos y una velocidad de 2649.998 MHz y con un ancho de banda de 16Mb/s para la comunicación con servicios como la API de *OpenAI*.

El primer sistema de evaluación fue sobre el *benchmark MetaQA 4.1*, el cual constituye el principal conjunto de datos de evaluación para la tarea *Text-to-Cypher* vista en la sección 2.1. En este caso se utilizó la versión clásica, donde los pares de evaluación consistían en una consulta en lenguaje natural con su correspondiente respuesta en la base de datos.

4.1. Evaluación sobre el *benchmark MetaQA* [?]

MetaQA [?] es un conjunto de datos diseñado para la tarea de razonamiento de múltiples pasos (*multi-hop*) en respuesta a preguntas. Está compuesto por entidades, relaciones y preguntas en lenguaje natural relacionadas con películas. Cada nodo en el grafo de conocimientos representa una entidad (como una película, actor o director), y las aristas representan relaciones entre las entidades. El conjunto de datos también incluye pre-

guntas a tres niveles de complejidad (*1-hop*, *2-hop* y *3-hop*), con cada nivel requiriendo razonamiento sobre un número creciente de aristas en la base de datos en forma de grafos analizada para responder correctamente a las preguntas. A continuación se muestra un ejemplo de la distribución de dicho conjunto de datos:

Cuadro 4.1 Distribución de los conjuntos de datos del *benchmark MetaQA*.

	1-hop	2-hop	3-hop
Train	96,106	118,980	114,196
Dev	9,992	14,872	14,274
Test	9,947	14,872	14,274

En este estudio solo se utilizarán los datos referentes a los conjuntos de evaluación (Test) para cada uno de los grupos especificados, ya que el modelo empleado es un gran modelo de lenguaje mediante la técnica *Zero-Shot*, por lo que no es necesario hacer un proceso de entrenamiento al mismo para realizar la tarea en cuestión, ya que se desea analizar la capacidad de inferencia del mismo sin haber sido entrenado específicamente para esta.

Las principales métricas de evaluación utilizadas fueron el número de consultas que al ser traducidas a *Cypher* y ser ejecutadas ejecutadas sobre la base de conocimiento daban una respuesta idéntica a la respuesta objetivo (*correct*), así como el porcierto de dichas consultas acertadas (*correct%*) sobre el total de consulta (*n*), número de consultas compiladas con éxito y su porcierto correspondiente (*compiled* y *compiled%* respectivamente), algunas relacionadas con los recursos consumidos para el experimento como el costo monetario (*cost (USD)*), tiempo de ejecución de la evaluación en segundos (*elapsed_seconds*). También, se calculó eficacia del modelo con respecto a cada tipo de consulta en el conjunto de evaluación y medidas clásicas como la precisión, el recobrado y la medida *F1* para evaluar la calidad de la extracción de información. La métrica *correct%* fue la utilizada para comparar el resultado del modelo empleado sobre otros resultados en el estado del arte ?? Además, implícitamente, al evaluar la efectividad del modelo sobre las consultas de los conjuntos de evaluación de *1-hop*, *2-hop* y *3-hop*, se evalúa la eficacia del modelo sobre consultas que requieren de una relación, dos relaciones y hasta tres relaciones de conexión respectivamente para encontrar la respuesta a la consulta.

Para la preparación del conjunto de datos se insertaron los elementos correspondientes a la base de conocimientos en una instancia de *Neo4J* con

ayuda del componente DBSeeder visto en la sección 2.3.3. Luego se tomaron los conjuntos de prueba (*Test*) para *1-hop*, *2-hop* y *3-hop* y para cada par de evaluación se ejecutó el procedimiento descrito en el listado 2.3.6.

4.1.1. Resultados

Los resultados obtenidos para cada métrica analizada para cada conjunto de evaluación se muestran en la siguiente figura:

Cuadro 4.2 Resultados de ejecutar GPT-4 en los conjuntos de datos de prueba de *MetaQA* para *hop1*, *hop2* y *hop3*.

	n	compiled	correct	compiled %	correct %
hop 1	9947	9386	7613	94.36	76.53
hop 2	14872	13733	6462	92.34	43.45
hop 3	14274	13221	4430	92.62	31.03

En la tabla 4.2 se muestran los resultados de eficacia del modelo GPT-4 para traducir consultas a *Cypher* tal que puedan ser utilizadas para extraer información de la base de datos objetivo. En este caso, una consulta generada por el sistema utilizado fue considerada eficaz si al ejecutar el código de *Cypher* sobre la base de datos correspondiente, dicho resultado coincide con los datos de respuesta esperados asociados a cada par de los conjuntos de evaluación. Para aquellas consultas cuyo código de *Cypher* correspondiente requería de la presencia de una relacion específica entre dos entidades en cuestión se tuvo relevante resultado del 76,53% de acierto. Por otro lado, aquellas consultas que requerían de la generación de una consulta con dos y hasta tres relaciones tuvieron como resultados unos discretos 43,45% y 31,03% respectivamente, lo que nos indica la deficiencia de este modelo para responder expresiones en lenguaje natural complejas que requieran acceder a la información de más una relación entre dos entes de la base de datos en forma de grafo.

Es importante resaltar de los resultados anteriores la efectividad del sistema para generar consultas de *Cypher* compilables, donde para cada lote de evaluación se obtuvieron valores sobre el 92%, lo que confirma que se tuvo éxito en la inmensa mayoría de las veces que el modelo trató de traducir la consulta inicial en lenguaje natural a solamente un texto conteniendo un código de *Cypher*. En los casos donde dicho suceso no fue posible, se debió principalmente a que el modelo generó un texto adicional describiendo la consulta, ofrecía más de una consulta o simplemente había un error sin-

táctico en el código.

Cuadro 4.3 Resultados para las consultas del lote *hop 1*.

Tipo de consulta	Correctas	Total	Efectividad
actor_to_movie	650	879	73.94
director_to_movie	446	553	80.65
movie_to_actor	672	1105	60.81
movie_to_director	984	1301	75.63
movie_to_genre	946	1143	82.76
movie_to_language	251	294	85.37
movie_to_tags	611	846	72.22
movie_to_writer	909	1091	83.31
movie_to_year	1122	1420	79.01
tag_to_movie	236	411	57.42
writer_to_movie	786	904	86.94

Cuadro 4.4: Resultados para las consultas del lote *hop 2*.

Tipo de consulta	Correctas	Total	Efectividad
actor_to_movie_to_director	488	929	52.53
director_to_movie_to_director	86	164	52.44
director_to_movie_to_language	160	193	82.90
writer_to_movie_to_writer	481	763	63.04
actor_to_movie_to_genre	470	823	57.11
director_to_movie_to_genre	380	533	71.30
actor_to_movie_to_actor	585	971	60.25
writer_to_movie_to_actor	414	838	49.40
actor_to_movie_to_writer	417	834	50.00
movie_to_director_to_movie	501	1081	46.35
actor_to_movie_to_year	186	985	18.88
writer_to_movie_to_genre	550	844	65.17
director_to_movie_to_actor	300	483	62.11
movie_to_actor_to_movie	254	1180	21.53
writer_to_movie_to_year	45	1020	4.41
director_to_movie_to_year	138	594	23.23
director_to_movie_to_writer	107	402	26.62
movie_to_writer_to_movie	225	896	25.11
Continuación en la siguiente página			

Cuadro 4.4 – continuación desde la página anterior

Tipo de consulta	Correctas	Total	Efectividad
writer_to_movie_to_director	380	763	49.80
writer_to_movie_to_language	133	250	53.20
actor_to_movie_to_language	162	326	49.70

Cuadro 4.5: Resultados para las consultas del lote *hop 3*.

Tipo de consulta	Correctas	Total	Efectividad
movie_to_director_to_movie_to_language	132	616	21.43
movie_to_director_to_movie_to_actor	147	1069	13.75
movie_to_actor_to_movie_to_language	377	840	44.88
movie_to_writer_to_movie_to_year	256	833	30.73
movie_to_actor_to_movie_to_director	609	1166	52.23
movie_to_director_to_movie_to_genre	334	1045	31.96
movie_to_writer_to_movie_to_director	200	917	21.81
movie_to_actor_to_movie_to_year	304	1132	26.86
movie_to_actor_to_movie_to_writer	510	1182	43.15
movie_to_actor_to_movie_to_genre	441	1148	38.41
movie_to_director_to_movie_to_writer	243	1151	21.11
movie_to_writer_to_movie_to_genre	323	835	38.68
movie_to_writer_to_movie_to_actor	172	873	19.70
movie_to_director_to_movie_to_year	268	1099	24.39
movie_to_writer_to_movie_to_language	114	368	30.98

En las tablas 4.8, 4.1.1 y 4.1.1 se muestran para cada lote de prueba, el número de consultas de cada tipo presente, donde a su vez se reflejan aquellas en cuyo caso el sistema propuesto obtuvo un resultado correcto. Es posible notar que a medida que aumenta la complejidad de la consulta, la efectividad del modelo decrece, pues se hace cada vez más difícil generar consultas de *Cypher* válidas que contengan hasta cuatro entidades a solicitar, ya que este proceso implica predecir correctamente las direcciones de las relaciones entre un mayor número de entidades, y basta con que una dirección en alguna relación falle para que la consulta ofrezca un resultado incorrecto.

Cuadro 4.6 Precisión, recobrado y medida F1 para cada lote de evaluación.

	Precisión	Recobrado	F1
hop 1	85.94	80.48	83.12
hop 2	32.24	59.28	41.77
hop 3	50.47	39.19	44.12

En la tabla 4.6 se muestran la precisión, recobrado y medida $F1$ para cada lote de evaluación. Para calcular los mismos se tuvieron en cuenta la cantidad de resultados positivos correctamente identificados (verdaderos positivos), aquellos identificados como correctos pero que no se encontraban en la respuesta objetivo (falsos positivos) y aquellos que estaban presentes en la consulta objetivo pero que no se obtuvieron en la consulta generada por el modelo propuesto al ser ejecutada sobre la base de datos en cuestión (falsos negativos). Estas medidas representan en una mejor manera la capacidad del modelo para la extracción de información, ya que a pesar de que para ciertas consultas de prueba no se obtuvieron todos los resultados esperados, es útil conocer qué tan distante está la respuesta dada de la esperada, por lo cual, la decisión de calcular la precisión, recobrado y medida $F1$ es correcta ya que las mismas son una correcta forma de implícitamente evaluar dicha medida de correctitud [?].

Cuadro 4.7 Costo monetario y tiempo de ejecución del experimento.

	cost (USD)	elapsed_seconds
hop 1	139.33	57587.00
hop 2	220.66	79240.58
hop 3	218.62	92191.07

En la tabla 4.7 es posible ver reflejados los recursos monetarios y de tiempo consumidos por la realización del experimento en el *VPS* utilizado 4. Como se muestra, la ejecución del modelo *GPT-4* a partir de la *API* de *OpenAI* resulta costoso y requiere de condiciones ideales de ejecución, como por ejemplo una conexión a *Internet* estable para poder acceder a la misma.

Cuadro 4.8 Comparación de los resultados de otros modelos respecto al *benchmark MetaQA*.

Método	1-hop	2-hop	3-hop
SOTA	97.50	98.80	94.80
zero-shot	24.75	6.37	9.72
zero-shot-cot	18.41	12.86	21.89
zero-shot+graph	91.69	46.82	19.40
zero-shot-cot+graph	86.16	47.36	19.29
zero-shot+graph+change-order	95.20	40.48	20.17
zero-shot-cot+graph+change-order	95.87	47.71	23.95
zero-shot Cypher Generation	30.00	10.00	13.00
GPT-4 zero-shot Cypher Generation	76.53	43.45	31.03
one-shot Cypher Generation	99.00	77.00	96.00

La tabla 4.8 refleja el resultado del sistema implementado comparado con otros enfoques utilizados sobre *MetaQA*. En cada columna de la tabla relacionada con *1-hop*, *2-hop* y *3-hop* se reflejan los valores porcentuales de acierto de ejecución de dichas vías de solución propuestas sobre el conjunto de evaluación (Test) correspondiente. La primera fila contiene el mejor resultado para cada conjunto con respecto al estado del arte, las seis filas representan el resultado de utilizar GPT-3 (code-davinci-003) para la tarea de extracción de información de la base de datos sin utilizar lenguaje *Cypher* como paso intermedio. En las filas 8 y 10 se reflejan los resultados para GPT-3 utilizando *Cypher* como vía para extraer información de una base de datos *Neo4J* utilizando los enfoques *Zero-Shot* y *One-Shot*. Finalmente, la fila 9 contiene los resultados referentes para cada conjunto del modelo propuesto.

De acuerdo con el estudio más reciente realizado por Guo et al. [?], la propuesta de sistema de traducción de esta investigación supera el mejor resultado que se tenía para la traducción de lenguaje natural a lenguaje *Cypher* utilizando aprendizaje *Zero-Shot* sobre el *benchmark MetaQA* y donde el modelo utilizado fue GPT-3, sin embargo, sus capacidades de extracción de conocimiento a partir de *Cypher* quedan todavía lejos de los mejores resultados del estado del arte para dicha tarea.

4.2. Discusiones

Después de aplicar GPT-4 para traducir consultas a Cypher, es pertinente destacar tanto las fortalezas como las deficiencias del sistema. Los resultados revelan una eficiencia notable en la generación de código *Cypher* compilable, alcanzando un 92% de éxito como promedio en las pruebas realizadas. Este alto grado de precisión indica que el modelo es eficaz en la creación de consultas sin errores sintácticos ni semánticos, lo que es esencial para su aplicación práctica en entornos de bases de datos como *Neo4J*.

Sin embargo, a pesar de esta eficacia en la compilación, el modelo demostró limitaciones en su capacidad para generar consultas correctas a medida que aumentaba la complejidad de las relaciones entre entidades. Se observó un descenso significativo en la precisión, pasando de un 76,53% en consultas simples (*1-hop*) a 43,45% y 31,03% en consultas más complejas (*2-hop* y *3-hop*). Esto sugiere que, aunque GPT-4 es competente en la traducción de consultas sencillas, su rendimiento se reduce considerablemente con consultas que involucran múltiples relaciones entre entidades.

En cuanto a las deficiencias del sistema, se identificaron varios aspectos que no se abordaron en el estudio. Uno de los más críticos fue la incapacidad del modelo para evaluar consultas anidadas y funciones de agregación, lo cual limita su aplicabilidad en escenarios de análisis de datos más complejos. Asimismo, la ausencia de un análisis multidominio impidió una evaluación adecuada de la capacidad de generalización del modelo, un factor crucial para determinar su eficacia en diferentes contextos y bases de datos. Además, el formato de las respuestas generadas por el modelo fue bastante básico, lo que plantea un área de mejora para futuras versiones, especialmente en aplicaciones que requieren un análisis de datos más detallado y avanzado. También, es posible notar que el tamaño del esquema de la base de datos utilizada no fue lo suficientemente extensa como para no caber en el texto de entrada limitado del modelo GPT-4, por lo cual se hace imprescindible probar como se comportaría el modelo para dichos casos extremos y posibles formas de resolverlo.

Conclusiones

Después de aplicar GPT-4 para la traducción de consultas al lenguaje *Cypher*, este estudio presenta conclusiones relevantes tanto en términos de fortalezas como de deficiencias. Destaca la eficiencia del modelo en generar código *Cypher* compilable, con un impresionante 92% de éxito en las pruebas, lo que subraya su competencia en la creación de consultas sin errores sintácticos o semánticos. Esta alta precisión es crucial para su aplicación práctica en bases de datos como *Neo4J*.

Sin embargo, el modelo exhibe limitaciones significativas al manejar consultas más complejas. Mientras que en consultas sencillas (*1-hop*) la precisión es del 76,53%, en consultas más complejas (*2-hop* y *3-hop*) esta precisión disminuye drásticamente a 43,45% y 31,03%, respectivamente. Esto indica que aunque GPT-4 es eficaz en traducciones simples, su rendimiento se ve comprometido en escenarios que involucran múltiples relaciones entre entidades.

Además, se identificaron varias áreas críticas no abordadas en el estudio, como la incapacidad del modelo para manejar consultas anidadas y funciones de agregación. Esto limita su utilidad en análisis de datos más complejos. La falta de un análisis multidominio también plantea preguntas sobre la capacidad de generalización del modelo, un factor esencial para determinar su eficacia en diferentes contextos y bases de datos. Otro aspecto a mejorar es el formato básico de las respuestas generadas, que no satisface necesidades de análisis de datos más detallado y avanzado. Además, se señala que el tamaño limitado del esquema de la base de datos utilizada no puso a prueba la capacidad del modelo para manejar esquemas más grandes, una limitación importante para su aplicación práctica.

A pesar de estos desafíos, el sistema de traducción propuesto supera al modelo GPT-3 en la traducción de lenguaje natural a *Cypher* en el *benchmark* MetaQA, aunque todavía no alcanza los mejores resultados en la extracción de conocimiento usando *Cypher*.

Gracias a las medidas de precisión, recuperación y medida *F1* utilizadas

para evaluar la capacidad del modelo para la extracción de información, considerando los verdaderos positivos, falsos positivos y falsos negativos, demuestran que, a pesar de no obtener todos los resultados esperados en ciertas consultas, es importante entender la distancia entre la respuesta dada y la esperada.

Finalmente, es evidente que la eficacia del modelo disminuye con el aumento de la complejidad de las consultas, especialmente en aquellas que requieren la predicción correcta de las direcciones de relaciones entre un mayor número de entidades. Este estudio deja claro que, mientras que el uso de un *grand* modelo de lenguaje con la técnica de aprendizaje *Zero-Shot* muestra una eficiencia notable en ciertos aspectos, aún hay un camino considerable por recorrer para mejorar su rendimiento en escenarios más complejos y variados.

Recomendaciones

Para futuros trabajos en la aplicación de grandes modelos de lenguajes en la traducción de consultas a *Cypher*, se recomienda abordar varias áreas clave para mejorar la eficacia y versatilidad del modelo. Estas recomendaciones incluyen:

- **Mejorar la Comprensión de Consultas Complejas:** Es esencial perfeccionar la capacidad del modelo para manejar consultas con múltiples relaciones entre entidades (*2-hop* y *3-hop*), que actualmente presentan una disminución significativa en la precisión. Esto podría implicar un entrenamiento adicional específico para estos tipos de consultas o la implementación de algoritmos más sofisticados para la comprensión de relaciones complejas.
- **Gestión de Consultas Anidadas y Funciones de Agregación:** Desarrollar sistemas de evaluación que contengan consultas anidadas y funciones de agregación, ampliando su aplicabilidad en análisis de datos avanzados y complejos.
- **Ampliación del Esquema de la Base de Datos:** Probar el modelo con esquemas de bases de datos más extensos y complejos permitiría evaluar y mejorar su capacidad de manejar casos más cercanos a escenarios del mundo real. Proponer una metodología para resolver dicho problema a partir de la adición de módulos adicionales de preprocesamiento de la consulta de entrada que permitan reducir el tamaño de la descripción de la base de datos, mostrando solamente los aspectos más relevantes para la consulta en lenguaje humano a responder.
- **Análisis Multidominio:** Realizar pruebas en múltiples dominios y tipos de bases de datos podría ayudar a evaluar y mejorar la capacidad de generalización del modelo, lo que es crucial para su eficacia en diferentes contextos.

- **Mejorar del Formato de Respuestas Generadas:** Elaborar consultas de prueba que generen salidas con formatos complejos, lo cual evaluará la capacidad del sistema de devolver datos de la manera especificada.
- **Mejorar las estadísticas relacionadas con las consultas incorrectas durante la evaluación:** Trabajar en desarrollar algoritmos y técnicas que permitan determinar, para una salida de un gran modelo de lenguaje en la tarea de este trabajo, cuántas entidades se detectaron correctamente en el código de *Cypher* generado, contabilizar y separar por grupos bien determinados las consultas de evaluación cuya respuesta falló debido a cuestiones semánticas y sintácticas, lo cual permitirá saber que estructuras del lenguaje *Cypher* son inherentemente complejas de traducir.

Al implementar estas recomendaciones, futuros trabajos podrán superar las limitaciones actuales del modelo GPT-4 en la traducción de consultas a *Cypher* y ampliar su aplicabilidad en una variedad de entornos y situaciones prácticas.

Bibliografía

- [1] Amazon. What is structured data? <https://aws.amazon.com/es/what-is/structured-data/>, 2023. (Citado en la página 11).
- [2] Web Archive. Creación de una base de conocimiento. <https://web.archive.org/web/20180315025341/http://es.ccm.net/faq/2158-organizacion-crear-una-base-de-conocimientos>, 2023. (Citado en la página 11).
- [3] Nghi D. Q. Bui, Hung Le, Yue Wang, Junnan Li, Akhilesh Deepak Gotmare, and Steven C. H. Hoi. Codetf: One-stop transformer library for state-of-the-art code llm, 2023. Accedido el: 05 de diciembre, 2023. (Citado en la página 12).
- [4] Naihao Deng, Yulong Chen, and Yue Zhang. Recent advances in text-to-sql: A survey of what we have and what we expect. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 2166–2187, Online, October 2022. Presented at the 29th International Conference on Computational Linguistics, October 12-17, 2022. (Citado en la página 13).
- [5] Meta GenAI. Llama 2: Open foundation and fine-tuned chat models. Technical report, Simons Foundation, member institutions, 2023. (Citado en la página 12).
- [6] SAP insights. ¿qué es el modelado de datos? <https://www.sap.com/latinamerica/products/technology-platform/datasphere/what-is-data-modeling.html>, 2023. (Citado en la página 11).
- [7] Sam Hays Michael Sandborn Carlos Olea Henry Gilbert Ashraf El-nashar Jesse Spencer-Smith Douglas C. Schmidt Jules White, Qu-

- chen Fu. A prompt pattern catalog to enhance prompt engineering with chatgpt. *ArXiv*, 2023. cs.SE. (Citado en las páginas 12 y 39).
- [8] Peng Li, Tianxiang Sun, Qiong Tang, Hang Yan, Yuanbin Wu, Xuanjing Huang, and Xipeng Qiu. Large code generation models are better few-shot information extractors. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, volume 1, pages 15339–15353, Association for Computational Linguistics, 2023. Accedido el: 05 de diciembre, 2023. (Citado en la página 12).
- [9] Fundación MAPFRE. ¿cuánta información se genera y almacena en el mundo? <https://www.fundacionmapfre.org/blog/cuanta-informacion-se-genera-y-almacena-en-el-mundo/>, 2023. (Citado en la página 11).
- [10] Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. A comprehensive overview of large language models. *ArXiv*, 2023. (Citado en la página 12).
- [11] OpenAI. Gpt-4 technical report. Technical report, Simons Foundation, member institutions, 2023. (Citado en la página 12).
- [12] O. M. Parkhi, S. M. Ali, M. Elgammal, and C. K. I. Williams. Zero-shot learning - a comprehensive evaluation of the good, the bad and the ugly. *ArXiv*, 2017. cs.LG. (Citado en la página 12).
- [13] Neo4J Official Site. Neo4j graph database: The mosted trusted database for intelligent applications. <https://neo4j.com/product/neo4j-graph-database/>, 2023. (Citado en la página 12).
- [14] Neo4J Official Site. Query a neo4j database using cypher. <https://neo4j.com/docs/getting-started/cypher-intro/>, 2023. (Citado en la página 12).
- [15] Neo4J Official Site. Query a neo4j database using cypher. <https://neo4j.com/docs/getting-started/cypher-intro/>, 2023. (Citado en la página 13).
- [16] IBM Official Website. Lenguaje de consulta estructurada (sql). <https://www.ibm.com/docs/es/db2woc?topic=reference-sql>, 2023. (Citado en la página 12).

- [17] Wikipedia. Graph database. https://en.wikipedia.org/wiki/Graph_database, 2023. (Citado en la página 11).
- [18] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A survey of large language models. *ArXiv*, 2023. (Citado en la página 12).