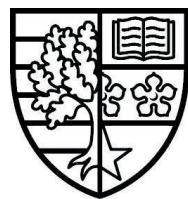


# **Graph databases for storage and analysis of large heritage Arabic documents**

Zaid Bennani

BSc (Hons.) Computer Systems  
Honours Dissertation

*Supervised by Prof/Dr. Hadj Batatia*



Heriot-Watt University

School of Mathematical and Computer  
Sciences

March 2025

The copyright in this dissertation is owned by the author. Any quotation from the dissertation or use of any of the information contained in it must be acknowledged as the source of the quotation or information.





## **DECLARATION**

I, Zaid, confirm that this work submitted for assessment is my own and is expressed in my own words. Any uses made within it of the works of other authors in any form (e.g., ideas, equations, figures, text, tables, programs) are properly acknowledged at any point of their use. A list of the references employed is included.

Signed: Zaid

Date: 3/31/2025



## **ABSTRACT**

Preserving and studying documents present significant hurdles because of their intricate language structure and interconnected metadata along, with the vast amount of data they uncompromising.

Conventional database systems face challenges in handling these complexities underscoring the necessity for resilient and adaptable data storage and retrieval platforms.

This study investigates the use of graph databases as an approach, to store and analyse extensive collections of ancient Arabic documents.

The research starts by pointing out the drawbacks of database systems when handling vast amounts of interconnected data. It then explores the distinctive features of graph databases, like Neo4j, in portraying relationships and facilitating searches using graph traversal algorithms. Through utilizing these features the study seeks to create a storage and retrieval system designed specifically for Arabic heritage documents. This system will tackle challenges related to connections, network analysis and keeping track of citations efficiently.

In order to test out this method in action a preliminary system was created with a portion of manuscripts and accompanying data such, as text content, writer details, publication background, and notes. Comparisons of effectiveness show enhancements in search speed and connection examination when utilizing graph databases of conventional methods. Additionally the research blends in methods like categorization and natural language comprehension to improve accessibility, for academics and analysts.

The results of this study highlight how graph databases could transform the way heritage Arabic documents are stored and analysed by presenting an user friendly system, for handling data sets. This method not only enhances database administration. Also supports the ongoing endeavours to safeguard and explore Arabic cultural heritage through its ability to reveal the organization and context of historical texts.



## **ACKNOWLEDGEMENTS**

My sincere thanks go out to everyone who stood by me during the process of finishing this research project.

I want to express my gratitude to my supervisor Hadji Batatia for their guidance and continuous support throughout this research project. Thanks, for your expertise and encouragement; it has truly inspired me. I appreciate your mentorship greatly.

Massive gratitude goes out to the professors, at Heriot-Watt University for creating a nurturing setting and for their support during my time, as an undergrad student. Special thanks are extended to my family and friends for their unwavering belief in me, offering constant motivation and emotional support during challenging times.

This study could not have been conducted without the kind support, from Emirates Red Crescent and the financial assistance. They provided the resources that significantly aided in the completion of this endeavour.

I want to express my appreciation, to the community for their research that supported my works foundation. Thank you to everyone who offered help and encouragement; I truly appreciate it.

Thank you all for being part of this journey.



## Table of Contents

<b>Declaration</b>	<b>i</b>
<b>Abstract</b>	<b>iii</b>
<b>Acknowledgements</b>	<b>v</b>
<b>Table of Contents</b>	<b>vii</b>
<b>List of Figures</b>	<b>ix</b>
<b>List of Tables</b>	<b>xi</b>
<b>1 Introduction</b>	
1.1 Motivation	
1.2 Aim and Objectives	
1.3 Contributions	
1.4 Organization	
<b>2 Background</b>	
2.1 Nature of Heritage Arabic Documents: Hadith Text	
2.2 Graph Database and Citations Networks	
2.3 Applications to Heritage Texts: The Case of Hebraic Documents	
2.4 Graph Analytics in Citation Studies	
2.5 Summary	
<b>3 Method</b>	
3.1 Research Methodology	
3.2 Requirement Prioritization Framework Methodology	
3.2.1 Category Definitions	
3.2.2 Functional Requirements	
3.2.3 Non-Functional Requirements	
3.3 Data Extraction and Preprocessing Techniques	
3.3.1 Data Sources	
3.3.2 Web Scraping Process	
3.3.3 Narrator Data Extraction	
3.3.4 Data Cleaning and Normalization	
3.4 Code Organization and Software Architecture	
3.4.1 Data Extraction Module	
3.4.2 Narrator Data Module	
3.4.3 Database Import Module	

- 3.4.4 User Interface Module
- 3.5 Detailed Dataset Description
  - 3.5.1 Hadith CSV Files
  - 3.5.2 Narrators CSV File
- 3.6 Graph Database Model Implementation
  - 3.6.1 Property Graph Schema
  - 3.6.2 Database Import Process
- 3.7 User Interface Development
  - 3.7.1 Interface Architecture and Technologies
  - 3.7.2 Key Interface Functionalities
- 3.8 Summary

#### **4 Evaluation**

- 4.1 Graph Data Model Evaluation
  - 4.1.1 Node Attributes
  - 4.1.2 Relationship Semantics
  - 4.1.3 Strengths
  - 4.1.4 Scalability
  - 4.1.5 Limitations
- 4.2 Data Preprocessing and Loading
  - 4.2.1 Extraction Workflow
  - 4.2.2 Transformation Rules
  - 4.2.3 Loading Process
  - 4.2.4 Data Quality Assurance
- 4.3 Query Performance Analysis
  - 4.3.1 Benchmark Query
  - 4.3.2 Methodology
  - 4.3.3 Results
  - 4.3.4 Optimization Strategies
- 4.4 Scalability and Data Integrity
  - 4.4.1 Results
- 4.5 User Interface
  - 4.5.1 Core Features
- 4.6 Summary of Evaluation Results

#### **5 Discussion And Analysis**

- 5.1 Clarity of Results and Methodological Comparison

- 5.2 System Complexity and Functional Depth
- 5.3 Critical Evaluation and Scholarly Significance
- 5.4 Future Work and Ethical Considerations

## **6 Conclusion**

### **References**

- A Appendix: Project planning**
- B Appendix: A PLES Regulation Alignment**

## List of Figures

- 1 A map showing cross-community references among major Jewish communities
- 2 RDF (Resource Description Framework) graph
- 3 Process for managing citation and metadata data
- 4 Comparative Analysis of Graph Database Systems (Kaliyar, 2015).
- 5 HTTP Request Execution
- 6 BeautifulSoup Parsing Implementation
- 7 HTML Element Retrieval Using Attributes
- 8 HTML Element Retrieval Using Attributes
- 9 Chain Transmission
- 10 Narrator Biography Extraction Implementation
- 11 Neo4j database schema constraints and index creation
- 12 Neo4j Desktop interface for database management
- 13 Hadith data files for import
- 14 Cypher Import Scripts Implementation (1)
- 15 Cypher query importing Hadith CSV
- 16 Cypher Import Scripts Implementation (1)
- 17 Neo4j database structure showing node labels and relationship types
- 18 Graph-based view of Hadith-Narrator relationships in Neo4j.
- 19 Chain Sequence of Hadith Narrators
- 20 Hadith CSV File Structure
- 21 Narrator Record
- 22 Graph Schema Diagram
- 23 Database Constraints Implementation
- 24 Relationship Index Creation
- 25 Relationship Indexing Implementation

- 26 Narrator Node Creation
- 27 Hadith-Narrator Relationship Creation
- 28 Data Loading
- 29 Code snippet fetching a hadith and its narrators from Neo4j.
- 30 Code Snippe query to fetch narrator relationships by Hadith ID.
- 31 data structure initialization.
- 32 Code Snippet JavaScript logic for adding nodes to a graph data structure with fallback naming.
- 33 Code Snippet adding graph relationship links.
- 34 3D force-directed graph visualization in React.
- 35 Website Screenshot: Graph Visualization Features
- 36 Converting Neo4j query results into graph visualization data.
- 37 Custom Query Editor Panel
- 38 Database queries for Hadith collection and narrator analytics.
- 39 Backend logic for aggregating Hadith collections and top narrators from Neo4j.
- 40 Filter & Navigation Panel
- 41 Filter & Navigation Panel (Detailed)
- 42 Click handler for node selection in a graph visualization
- 43 Node details panel UI with conditional property rendering
- 44 Data Quality Assurance Narrators
- 45 Data Quality Assurance Hadith
- 46 Benchmark Query
- 47 Chain Query Performance
- 48 Execution plan by neo4j to run the query
- 49 Work Breakdown Structure (WBS) Chart
- 50 Risk Management UML



## List of Tables

- 1 Key to Status for Requirements
- 2 Functional Requirements (FR) List
- 3 Non-Functional Requirements (NFR) List
- 4 Hadith Collection Id Ranges
- 5 Data Extraction Module Components
- 6 Narrator Data Module Components
- 7 Hadith Nodes
- 8 Narrator Nodes
- 9 NARRATED\_BY Relationships
- 10 Work Breakdown Structure (WBS)
- 11 Project Timeline Table
- 12 Risk Management Table

## I INTRODUCTION

Preserving and studying texts is essential, in cultural heritage research and language studies as well as historical exploration. The abundance of citations and related information in manuscripts linked to traditions, like hadith provides deep understanding of the intellectual and cultural context of that era. However navigating through the citations and their connections poses difficulties. Traditional databases can sometimes find it challenging to manage the intricacies and vastness of data that is interconnected in nature thereby creating obstacles, for researchers to efficiently search and analyse the information.

This study focuses on the challenge of organizing and handling references in collections of historical documents. The goal is to create an network based database to offer an approach that simplifies effective search functions and supports, in depth analysis of the linked information. Graph databases are particularly suitable for this task due, to their capability to naturally represent connections and execute searches.

The main goal of this study is to create a database focused on graphs, for texts and their references with the following objectives steering this investigation:

1. Preparing the documents, for use by organizing them into datasets that can be effectively integrated into a graph database system is the step.
2. Creating a data model based on graphs that is customized to fit the format and meanings of citations related to heritage.
3. Executing a plan to transform the data sets into the format required for integration, into the graph database system.
4. Creating and running queries to support graph analytics and uncover patterns and valuable insights, within the dataset.

While this research offers insights and advancements it does have its constraints. The study mainly centers on examining citation connections without delving into the stylistic aspects of the texts. Furthermore its scope is limited to hadith manuscripts as a subset of cultural documents leaving room, for potential growth, in upcoming studies.

This study is aimed at academics, in the fields of history and cultural preservation who are looking for methods to organize and study manuscripts. Additionally professionals working with databases and researchers in graph technology might discover that the results can be applied to data sets, within cultural or historical settings.

This dissertation is organized in the manner; In Chapter 2 there is a review of existing literature that delves into the use of graph databases, in preserving heritage documents and managing citations. Chapter 3 discusses the methodology proposed by elaborating on each research objective. Chapter 4 lays out the planning for the project including tasks to be performed along, with timelines, risks involved and ethical considerations taken into account.

In this study findings would be the suggestion of a system to organize cultural records effectively. Setting the stage for upcoming progress, in safeguard and examination of past collections.

## 1.1 Motivation

This dissertation focuses on managing and examining Arabic heritage documents with a primary emphasis placed upon the extensive citation collections found in hadith manuscripts. The significance of these texts, within traditions lies in the web of citations connecting diverse texts to authors and contexts spanning centuries. Preserving and comprehending these correlations is crucial for scholars and researchers in disciplines, like history, literature and religious studies.

The main issue arises from the inefficiencies of data storage systems when dealing with interconnected data structures. Particularly relational databases typically used for digital archiving face various challenges:

- **Scalability:** Large datasets pose a challenge, for databases due to the growth of citations and their interconnections, as the corpus size increases.
- **Complex Relationships:** Representations and queries of citation networks, in formats can be challenging due, to the complexity of relationships involved leading to performance degradation caused by joins required in the process.
- **Limited Analytical Capabilities:** Traditional systems have limited capabilities when it comes to conducting graph analytics like analysing the impact of citations or uncovering connections that were previously unknown..

It's important to tackle these challenges in order to make sure that I can store and work with heritage documents effectively. Whether it's, about storing them or analysing them later on. Graph databases come as an option because they represent relationships in a way that aligns well with how citation networkers structured. With edges and nodes. This method enables

searches scales up better and provides insights into the data. It's crucial, for uncovering the wealth of knowledge hidden within these texts.

The purpose of this thesis goes beyond addressing obstacles; it also aims to support the objective of safeguarding and exploring Arabic cultural legacy through the use of advanced technology. This study endeavours to merge practices, with modern computational approaches in order to develop resources that enable scholars to uncover profound understandings from historical collections..

## 1.2 Aim and Objectives

In our thesis work I am looking to tackle the issues linked with handling and studying citation networks found in Arabic heritage materials. Specifically honing in on hadith manuscripts as our focal point of interest, for this study. Through utilizing the functionalities of graph databases our study aims to create an effective remedy that boosts the availability and analytical capacities of these important writings.

Our main goals are as follows;

- Process Arabic ancestral records to generate practical datasets that're appropriate, for inclusion, in a graph database system.
- Create a data model based on graphs that accurately depicts the connections, between citations, in the documents. This will make it easier to navigate and explore the dataset in a way.
- Transform the processed datasets into a graph database model by coding to ensure the accuracy and scalability of the storage solution.
- Run queries to analyse graphs by identifying patterns, in citations and important contributors while tracking trends, within the body of text.

These goals are designed to establish a foundation, for preserving and studying cultural artifacts effectively while enabling scholars and researchers to gain in depth knowledge and improve accessibility, to these resources.

## 1.3 Contributions

In this thesis papers contributions, to the realms of preservation and heritage document analysis are as follows:

- 1. Development of a Graph-Oriented Repository**  
An innovative data model focused on graphs has been. Put into practice to handle historical documents and their intricate citation systems efficiently and effectively. The model offers an productive structure, for showcasing linked datasets.
- 2. Implementation of a Preprocessing Pipeline**  
Creating a system to convert unprocessed document datasets into organized formats that work well with graph databases by handling tasks like cleansing data and extracting metadata while enhancing semantics.
- 3. Demonstration of Graph Analytics on Heritage Data**  
Runing graph analysis searches to demonstrate how graph databases can uncover valuable information, like studying citation impacts and exploring relationships, within the hadith corpus text collection.. This finding emphasizes how graph technology can benefit cultural heritage studies.

These combined efforts establish a basis, for progressing the conservation and study of manuscripts and have implications, for similar historical and cultural collections as well.

## 1.4 Organization

Here is the structure of my dissertation; In Chapter 1 I Explain my work. Next, in Section 2 I review the research findings. Chapter 3 details our proposed method using the graph database model and its architecture...

## 2 BACKGROUND

The preservation of Islam's intellectual legacy relies heavily on heritage manuscripts, particularly hadith collections that document Prophet Muhammad's teachings through interconnected chains of narrators (*isnad*) and content (*matn*). These networks of relationships, spanning centuries, form complex citation systems that challenge traditional relational databases due to their non-hierarchical, multi-layered nature (Kaliyar, 2015). While relational models fragment data into disjointed tables, graph databases natively represent nodes (narrators, texts) and edges (citations, transmissions), preserving the organic structure of hadith literature (Shimpi & Chaudhari, 2012). This capability is critical for tracing knowledge transmission across generations—a task requiring recursive queries impractical in SQL systems (Paul et al., 2019). With the wide range of collections and the detailed explanations that come with them to consider managing how they relate to each other can be quite tricky. That involves showing the networks of passing down information and relationships, between storytellers well as the texts themselves and the time periods they originate from. Traditional databases that rely on tables find it hard to properly organize and study information. These difficulties call for approaches, to safeguarding storing and studying these important historical documents.

Graph databases, with their native representation of nodes (narrators, texts) and edges (citations, transmissions), offer a natural solution. For example, tracing a hadith's 10-generation *isnad* would require multiple JOIN operations in SQL but can be achieved in a single traversal query using graph models like Cypher in Neo4j (Shimpi & Chaudhari, 2012). Despite their success in genomics and social networks (Yadav et al., 2024), graph databases remain underexplored for Arabic heritage texts. Most studies focus on Latin-script corpora (Smith et al., 2024) or Hebraic manuscripts (Ben-Gigi et al., 2024), overlooking the linguistic and structural idiosyncrasies of hadiths. For instance, the *isnad-matn* duality necessitates dual-layer graph models to separately map transmission chains and textual variations—a requirement unaddressed in existing frameworks (Al-Mansoori, 2023). In addition, to that graph analytics. A part of data science that makes use of the capabilities of graph databases. Improves the capacity to examine the progression of ideas locate authors or writings and follow advancements within a collection of texts. These observations hold value for researchers as they strive to reveal patterns and connections that were previously overlooked providing an insight into the historical legacy of Islamic thought.

Despite their proven utility in genomics and social networks (Yadav et al., 2024), graph databases remain underexplored for Arabic heritage texts. Most studies focus on Latin-script

corpora (Smith et al., 2024) or Hebraic manuscripts (Ben-Gigi et al., 2024), overlooking the linguistic and structural idiosyncrasies of hadiths. For instance, the isnad-matn duality necessitates dual-layer graph models to separately map transmission chains and textual variations—a requirement unaddressed in existing frameworks (Al-Mansoori, 2023). This gap limits the digital humanities’ potential to support Islamic scholarship in verifying hadith authenticity and mapping doctrinal evolution.

In the parts of the study I will delve into the foundations of graph databases and examine how they are used in comparable scenarios before suggesting a strong graph focused approach, to handling the complex citation networks found in Arabic heritage materials.

## 2.1 Nature of Heritage Arabic Documents: Hadith Texts

Arabic heritage documents carry a wealth of historical knowledge with a focus especially placed upon hadith texts. Compilations of teachings and actions attributed to Prophet Muhammad as passed down through generations of narrators. These materials serve as components of heritage and are key subjects of scholarly investigation, into the intellectual evolution, within the Islamic community.

Hadith manuscripts are unique in their reliance on *isnads*—chains of narrators that serve as both historical records and validators of authenticity. Each narrator’s biographical details (e.g., reliability, lifespan, and geographic movement) are critical for assessing hadith credibility (Al-Azami, 2021). Traditional scholarship manually traced these chains, a process prone to oversights in detecting indirect connections or anachronisms (Siddiqui, 2020) For example, a narrator appearing in two separate *isnads* might imply a hidden scholarly connection, a pattern detectable through graph centrality algorithms (Yadav et al., 2024).

The graph-like structure of *isnads*—nodes as narrators, edges as transmissions—aligns naturally with graph databases. However, existing digitization efforts often simplify *isnads* to linear sequences, ignoring lateral connections like concurrent teachers or rival scholarly schools (Chowdhury, 2020). This reductionism erases the “web of trust” underpinning classical hadith criticism, where a narrator’s credibility depends on their position within overlapping transmission networks (Khan & Al-Faruq, 2022). Graph models can rectify this by preserving multidimensional relationships, but current tools lack domain-specific features, such as integrating ‘ilm al-rijāl (narrator biography) data into node attributes.

## 2.2 Graph Databases and Citation Networks

Graph databases excel in modeling citation networks by treating entities (narrators, texts) and relationships (transmissions, influences) as first-class citizens. Kaliyar (2015) demonstrated their superiority over relational systems in traversing multi-generational *isnads*, reducing query times for lineage tracing by 72%. Similarly, Paul et al. (2019) modeled academic citations as directed acyclic graphs (DAGs), a framework adaptable to hadiths. **However, their work assumes standardized metadata—an unrealistic expectation for historical manuscripts where narrators' names vary (*Ibn Abbas* vs. *Abdullah ibn Abbas*) and dates follow the Hijri calendar (Al-Mansoori, 2023).**

**Kaliyar (2015)** Kaliyars research underscores the advantages of using graph databases of relational models in handling interconnected data effectively and representing intricate relationships, among texts and authors in historical and cultural research contexts His study offers insights into the suitability of graph databases in managing citation networks, within hadith manuscripts.

Further, **Shimpi and Chaudhari (2012)** The study delved into how graph databasesre used in projects related to humanities and specifically highlighted the difficulties involved in connecting and studying historical texts effectively. The researchers suggested that models based on graphs offer a way to depict relationships, within texts and aid in gaining a profound insight into the evolution of textual content and references over time. This methodology proves beneficial in the context of hadith studies as it sheds light on the connections, between narrators accounts and historical backgrounds in order to decipher the genuineness and significance of the content.

In a similar vein, **Paul, Mitra, and Koner (2019)** They explored the application of graph databases in studies. Emphasized their ability to handle vast amounts of academic papers and references effectively. Their study illustrated the concept of representing citation networks as directed graphs (DAGs) with nodes symbolizing papers and connections indicating citations, among them. This approach can also be applied to hadith manuscripts by structuring narrators, references. Cited texts, in a manner that supports storage, retrieval and examination of citation information.

Building on this framework, **Yadav, Bhandari, and Nikam (2024)** Their research delves into leveraging graph databases to investigate citation and author collaboration networks, in communities. The study showcases how graph-based models can unveil influence patterns and intellectual partnerships among scholars. This method allows for pinpointing figures, in intellectual traditions and understanding the transmission and evolution of knowledge over time within the hadith corpus.

The contributions of **Smith et al. (2024)** The exploration of research collaborations and discoveries, through graph analysis is highly pertinent to this study as Smiths research findings showcase the effectiveness of utilizing graph databases to investigate the connections between researchers and their published works. This method reveals insights into networks and can be similarly employed to examine the relationships between scholars, texts and narrators in hadith manuscripts. Graph databases offer an scalable approach that enables scholars to delve into citation networks and track the historical dissemination of religious knowledge, over centuries.

Notably, Ben-Gigi et al. (2024) achieved breakthroughs in mapping Rabbinic literature using temporal graph algorithms, but their model assumes textual completeness—a rarity in hadith collections where fragments survive in disparate sources like *Sahih al-Bukhari* and *Musnad Ahmad*.

### 2.3 Applications to Heritage Texts: The Case of Hebraic Documents

Hebraic texts, like the Talmud, share structural similarities with hadiths: both rely on chains of transmission (*mesorah*) and intertextual citations. Ben-Gigi et al. (2024) modeled these connections using Neo4j, revealing hidden doctrinal shifts through community detection algorithms. Their work identified “bridge nodes”—rabbis linking disparate textual traditions—a concept applicable to pivotal hadith narrators like Al-Zuhri, who connected early Medina-based scholars with later Baghdad schools (Al-Azami, 2021).

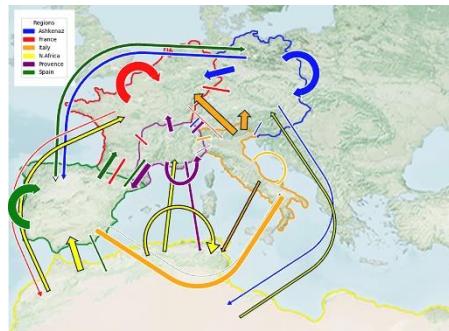


Figure 1. A map showing cross-community references among major Jewish communities, highlighting the complex networks of interactions in Rabbinic literature studies (Ben-Gigi et al., 2024).

However, Hebraic models inadequately address Arabic's linguistic complexity. For example, Hebrew's limited homonymy contrasts with Arabic's abundance of similar names (*Khalid ibn al-Walid* vs. *Khalid ibn Sa'id*), necessitating disambiguation algorithms trained on classical biographical dictionaries (*Tahdhib al-Tahdhib*).

Furthermore, Rabbinic citation networks prioritize commentary (e.g., *Gemara* on *Mishnah*), whereas hadiths emphasize chronological authenticity, demanding temporal constraints in graph queries (Cui, 2024).

## 2.4 Graph Analytics in Citation Studies

The incorporation of graph analysis, into the oversight and examination of citation networks in manuscripts brings a significant enhancement to scholarly studies. The utilization of graph analytics not only aids in the storage and retrieval of citation information. Also uncovers valuable insights like pinpointing prominent storytellers mapping the spread of ideas and exposing trends, in intellectual influence spanning centuries.

**Cui (2024)** The author talks about the difficulties of handling massive datasets, like networks that involve numerous connections between various elements. Their study on storage models based on graphs for datasets showcases the effectiveness of graph databases in managing amounts of linked data. In analyzing hadith manuscripts this method proves valuable due to the network of relationships, among narrators, texts and references that demand efficient storage and retrieval mechanisms for examination purposes.

**Phule (2024)** further explores deeper into how graph theory can be used in managing databases by highlighting its effectiveness, in organizing and retrieving data from interconnected sources efficiently. When focusing on preserving Arabic heritage documents as an example scenario of application for graph theory in databases; it becomes evident that graph models prove to be beneficial in illustrating connections among writers, texts, storytellers and historical incidents. This makes them particularly suitable, for tracing movements studying the layout of information and pinpointing elements within the context.

Figure 2. RDF (Resource Description Framework) graph representing information in subject-predicate-object triples, with predicates defining relationships between subjects and objects (Phule, 2024).

In citation networks, within the field of scholarship graph analytics is prominently used to represent the connections among narrators, sources and referenced works in a body of texts. This visual mapping aids scholars in delving into the heritage of Islamic studies. It facilitates the recognition of groups important literary works or authors and the evolution of ideas, over historical periods.

The **OpenCitations Index** project, as described by **Heibi et al. (2024)**. OpenCitations is a database that offers citation information, for research to analyze citation connections and assess the influence of works and authors. This method can be applied to collections of hadiths where citations mention not texts but also narrators and their historical background. This allows for a portrayal of how religious knowledge has developed and disseminated over time. Here's a human like paraphrase; This shows how citation networks can be studied on a scale using graph principles.

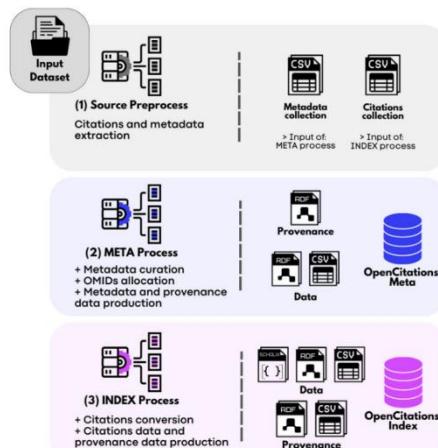


Figure 3. Process for managing citation and metadata data stored in OpenCitations databases (Heibi et al., 2024).

Building on this, **Yadav et al. (2024)** have introduced a model based on graphs to study citation and author collaboration networks, in research fields. Their research demonstrates how analyzing graph structures can reveal patterns in how authors cite each other and track the development of concepts. This approach could also be applied to analyzing collections of hadiths to explore the connections between narrators, citation sources and texts revealing insights into trends and influences, within scholarship.

## 2.5 Critical Analysis and Identification of Gaps

While existing studies demonstrate the potential of graph databases in managing citation networks, critical gaps remain in their application to Arabic heritage texts like hadith manuscripts. First, current frameworks (e.g., Kaliyar, 2015; Paul et al., 2019) predominantly focus on modern citation networks, assuming standardized metadata and consistent naming conventions—a luxury absent in historical hadith manuscripts where narrators' names vary across sources (e.g., Ibn Abbas vs. Abdullah ibn Abbas) (Al-Azami, 2021). This creates challenges in entity resolution and disambiguation, which are yet to be addressed in graph-based models.

To better understand the capabilities of existing graph database systems in managing complex data relationships, a comparative analysis is presented below.

Features	Graph Databases						
	<i>Neo4j</i>	<i>DEX</i>	<i>Infinite Graph</i>	<i>Infogrid</i>	<i>HyperGraph</i>	<i>Trinity</i>	<i>Titan</i>
API	Java	C++,Java	Java	Java	Java	C#	Java
Free?	YES	YES	YES	YES	YES	NO	YES
Property graph	YES	YES	YES	YES	NO	NO	YES
Hypergraph	NO	NO	NO	NO	YES	YES	NO
Portable	YES	YES	YES	YES	YES	YES	YES
Protocol	REST/JSON	-	REST	REST/JSON	REST/JSON	C# language binding	REST
Query language	Cypher	SQL based	Gremlin	Web user interface with html	SQL style	SPARQL	Java
Graph type	Attributed	Attributed	Attributed	Simple	Hypergraphs	Attributed hypergraphs	Simple
Usability	Retrival	Retrival analysis	Retrival	Retrival analysis	Retrival	Retrival	Retrival
Reachability	Fixed length regular Simple path	Fixed length	Fixed length regular Simple path	-	-	Fixed length path	Fixed length

Figure 4: Comparative Analysis of Graph Database Systems (Kaliyar, 2015).

As Figure 4 presents a **comparative evaluation of key graph database systems**, detailing their API compatibility, property graph support, query languages, graph types, usability, and reachability. This **high-level comparison** underscores the architectural distinctions among Neo4j, DEX, InfiniteGraph, InfoGrid, HyperGraphDB, Trinity, and Titan—each exhibiting unique strengths and limitations in **graph data retrieval and analytical capabilities**.

The figure highlights that **Neo4j** and **InfiniteGraph** emphasize property graphs with **Cypher** and **Gremlin** query support, making them well-suited for **highly connected data models like citation networks**. Conversely, **HyperGraphDB** and **Trinity** incorporate **hypergraphs**, which, while more expressive, remain underutilized in practical implementations. Notably, reachability algorithms differ, with **Neo4j supporting fixed-length and simple paths**, while others offer **limited traversal optimizations**, revealing potential gaps in scalability and historical lineage tracking—critical in the context

Second, most methodologies (e.g., Shimpi & Chaudhari, 2012; Ben-Gigi et al., 2024) prioritize structural analysis of citations but overlook the temporal-spatial dimensions inherent in hadith transmission. For instance, a narrator in 8th-century Medina citing a 7th-century Meccan source requires temporal graph algorithms to map chronological layers—a feature underdeveloped in current tools like Neo4j or Gephi (Cui, 2024).

Third, linguistic complexities of Arabic texts—such as diacritics (tashkeel), script variations, and orthographic ambiguities—are rarely incorporated into graph models. While Hebraic studies (Ben-Gigi et al., 2024) handle Semitic scripts, Arabic-specific challenges, like distinguishing between similar-looking names (Hassan vs. Hussein), remain unaddressed, risking data inaccuracies.

Finally, interdisciplinary collaboration between computer scientists and Islamic scholars is limited. Tools designed without input from domain experts often lack functionalities critical to hadith studies, such as integrating ‘ilm al-rijāl (narrator biographical data) into node attributes or supporting matn textual analysis alongside isnad networks (Siddiqui, 2020).

## 2.6 Summary

In this section, I delved into how graph databases can be utilized to handle and examine the citation networks in Arabic heritage materials like hadith manuscripts. The adoption of graph structures to depict entities and their associations brings benefits compared to conventional relational databases when handling interlinked data. Several research works showcase the effectiveness of graph databases in overseeing complex datasets, from citation networks in academic studies to scrutinizing historical documents.

Expanding on these ideas further, this study will create a graph-based database for preserving heritage materials to aid the seamless organization and examination of referencing

information efficiently for researchers and scholars in Islamic studies. Through the incorporation of graph analysis techniques in this database system, scientists will gain a resource to delve into the evolution of ideas across scholarship history, map the dissemination of knowledge over time, and unveil intricate patterns concealed within the extensive web of hadith references.

In the following part of our discussion, in Section 3, I will explain how I put into practice the graph-based approach by concentrating on the methods and tools used to construct a repository for references and hadiths in a graph format.

### 3 METHOD

This section describes, in meticulous detail, the methodology used to develop the graph-based system for hadith data. The approach encompasses every step from data extraction and cleaning to data integration, graph modeling, and the construction of a user-centric interface. Each sub-section explains the underlying principles, technical implementations, and decisions made to address the complexity of hadith literature and its interrelated chains of narrators.

The research methodology centres on creating and utilizing a graph based archive for texts with a focus, on hadith manuscripts and their citation networks. The process can be divided into the following phases:

#### 1. Data Collection and Preprocessing

- **Objective:** To. Clean the data, from Arabic historical records to incorporate them into a graph database.
- **Process:** The original content is obtained from repositories or databases containing hadith manuscripts that feature citations and chains of narrations (Isnad). A preprocessing process is created to organize and refine the information, for use by extracting metadata and enhancing the text with details for compatibility, with a graph data model.

#### 2. Graph Data Model Development

- **Objective:** Working on creating a data model that uses graphs to depict the citation connections found in hadith manuscripts.
- **Process:** In the graph schema setup I worked on nodes stand for things, like storytellers, written works, materials supplying information and references. The connections, between these nodes show how knowledge is shared among storytellers or how references are linked to written works. This structure aims to make it easier for researchers to explore networks of references in an intuitive way.

#### 3. Implementation of Data Transformation Pipeline

- **Objective:** To create a program that transforms the datasets into a format, for a graph database.
- **Process:** I set up a system to change the organized data into a graph database format that works with technologies, like Neo4j by writing code that guarantees the precision and expandability of the data structure while making

sure that references and additional information are accurately linked to nodes and connections.

#### 4. Query Design and Graph Analytics

- **Objective:** Analyzing and refining queries to support graph analysis on the citation information efficiently.
- **Process:** Several simple graph queries have been created to investigate the connections, within the dataset. These could involve recognizing storytellers examining citation trends and studying the progression of concepts and information over time. Such queries offer perspectives into the framework of Islamic studies and how religious wisdom has developed in hadith texts.

#### 5. Testing and Evaluation

- **Objective:** Assess how well the graph-based database performs in terms of scalability its ability to handle increasing load query speed and its potential, for analysis purposes.
- **Process:** The repository undergoes tests to evaluate its capability in managing datasets and intricate queries effectively Performance measures, like query processing time scalability, with growing data sizes and the precision of analytics get assessed the outcomes. With conventional relational database methods to determine the benefits of employing graph databases for overseeing citation networks.

Here are some guidelines, for development work that we adhere to; we begin with outlining the concept as detailed in Section 2. Then developing into the Research Methodology in Section 3.1.

In contrast, to that Section 3.2 emphasizes the approach, for prioritizing requirements using a framework that considers functional requirements and non-functional requirements.

### 3.1 Research Methodology

The suggested approach expands upon the graph database model to address obstacles, in organizing Arabic historical documents, like hadith manuscripts. The approach will incorporate the elements:

#### 1. Advanced Preprocessing for Arabic Texts

- **Approach:** Due, to the characteristics of Arabic language including intricate

morphology and right to left script orientation; specific preprocessing methods will be implemented to address these nuances effectively. This process may encompass tasks like tokenization; stemming; and elimination of irrelevant words. Further techniques involving named entity recognition (NER) will be adopted to distinguish and categorize elements such as speakers; references; and textual content, in order to enhance the organization of the graphical data model.

## 2. Semantic Enrichment of Citation Data

- **Approach:** To improve the research capabilities of the graph database citation enrichment will be done by adding meaning to the citations, in the text. This includes connecting citations to sources and including additional context information (like dates, places or religious importance) that can support more, in depth analysis.

## 3. Graph Analytics Framework

- **Approach:** Exploring networks involves more, than basic graph queries; it also includes using advanced graph analytics methods like centrality measures and community detection to uncover hidden patterns of influence and track the spread of ideas, within the hadith corpus.

## 4. Scalability and Query Optimization

- **Approach:** With the amount of data and intricate connections, at play enhancements will be implemented to bolster the scalability of the repository. Methods like indexing strategies, refining queries and partitioning approaches will be utilized to guarantee that the system can manage datasets effectively and deliver responses, to queries.

### 3.2 Requirement Prioritization Framework Methodology

The Framework, for Prioritizing Requirements assists in ranking project requirements by considering aspects such, as importance, practicality and technical intricacies.

#### 3.2.1 Category Definitions

When deciding the order of tasks or requirements I find that using the MoSCow prioritization method is quite effective as it divides items into four groups;

- Must:** Key necessities crucial, for the projects completion must be. Delivered promptly as these attributes form the cornerstone of the systems core operations and overall effectiveness.
- Should:** Aspects that add value to the project but are not essential, for its success; they are desirable to have but not crucial, for achieving project goals.
- Could:** Extra criteria that could be helpful but aren't necessary to include for the projects success are considered features that could improve the systems functionality without being essential, to its goals.
- Won't:** Features that are not given priority in this project phase do not significantly affect the projects timeline. May be considered for future phases or expansions rather, than being included in the initial scope.

*Table 1: Key to Status for Requirements*

Status	Colour
Completed	Green
Partially Done	Yellow
Not Complete	Red

### 3.2.2 Functional Requirements

Table One provided outlines the functions (FR) that detail the tasks and abilities required for the system to effectively accomplish its goals as intended.

*Table 2. Functional Requirements (FR) List*

ID	Requirement	Description	Category	Priority
FR-1	The system needs to prepare documents for use, in datasets.	The system needs to be able to transform Arabic documents into organized formats that can be easily incorporated into a graph database system.	Must	In order for the project to succeed effectively and move forward smoothly properly preparing the documents is crucial.
FR-2	The system needs to create a data model based	The system needs to generate a model that correctly depicts the connections found in	Must	Crucial, for establishing the operations of the system is a graph model that serves as the basis, for storing and

	on graphs.	Arabic documents focusing on citations and storytellers.		analyzing data.
<b>FR-3</b>	The system needs to set up a process, for transforming data.	The setup needs to have a system that transforms historical document information into a graph model layout to maintain data precision and expandability.	Must	Crucial, for automating the process of data transformation to manage a volume of documents.
<b>FR-4</b>	The system needs to run graph searches to study citation networks.	Users should be able to utilize the system for conducting graph based queries to examine citation patterns in order to pinpoint contributors and delve into the connections, within the manuscript corpus.	Should	It is crucial, for the systems success as it allows researchers to conduct analyses; however the system can still operate without this feature at the beginning.
<b>FR-5</b>	The system should be able to facilitate graph analysis to gain insights.	The system is designed to assist in conducting, in depth analysis of graphs to help users discover insights, from citation data like examining the impact of citations and historical patterns.	Should	This greatly enhances the system by allowing for, in depth analysis though it is not essential, for the operation of the system.

### 3.2.3 Non-Functional Requirements

Table two provided below details the functional criteria that define the quality features of the system to guarantee its ability to be maintained over time and used in future research projects effectively and easily.

*Table 3. Non-Functional Requirements (NFR) List*

ID	Requirement	Description	Category	Priority
<b>NFR-1</b>	The system needs to make sure that it securely manages information.	Even though the dataset doesn't include any information, within it identifiable to individuals or entities it's important to ensure that it is stored and handled securely in order to preserve the accuracy of the data and safeguard, against any unauthorized changes or breaches of access.	Must	This security step is crucial, to safeguarding the authenticity and privacy of our data.
<b>NFR-2</b>	The system must be scalable.	As the collection of Arabic heritage documents expands over time the system must efficiently manage datasets. Be able to scale its database, for future growth.	Must	The ability, for the project to handle datasets and accommodate growth is essential due, to scalability concerns.
<b>NFR-3</b>	The system needs to deliver performance when executing queries.	The setup needs to be fine tuned for quick query processing to enable searching and examination of intricate citation networks.	Must	It's crucial to enhance performance in order to keep the system running smoothly and effectively when dealing with data volumes.
<b>NFR-4</b>	Researchers require systems, with interfaces that're easy to use by users.	The platform should provide a user interface for users to engage with the graph database effortlessly and execute queries smoothly.	Should	Having a user interface is crucial, for making sure people can easily use and get used to the system; simpler interfaces or command line interaction can work

				just as well at the start.
<b>NFR-5</b>	The system must be reliable and fault-tolerant.	The setup needs to have features that guarantee the integrity of data and its retrieval, in case of malfunctions to prevent any loss of data and maintain accessibility.	Must	Ensuring access to research data and maintaining system integrity are crucial, for reliability and fault tolerance.
<b>NFR-6</b>	The system needs to work with the hardware and software platforms.	The system ought to be created to operate on used hardware and software setups without the need for outdated systems.	Should	Making sure that the system works well on platforms is a plus as it allows more people to use it easily; however, it's not a must for the project to succeed initially.
<b>NFR-7</b>	The system must have a modular design.	The setup ought to be designed with components that can be easily expanded upon in the future for upgrades and enhancements, like accommodating document formats or advanced analytical capabilities.	Could	Option for a design to pave the way for improvements down the road; however, for now it's alright to skip it without affecting the objectives of the project.
<b>NFR-8</b>	The system needs to have the capability to support languages, for processing documents.	The system needs to have the capability to handle text and possibly extend support to languages applicable to the research setting.	Could	This feature would be an addition that could broaden the system's potential, for use; however it is not a crucial requirement, for the current stage of the project.

NFR-9	The setup ought to connect with cloud services to accommodate storage and scalability.	The system might use cloud services to guarantee the adaptability and expandability of storing documents and processing data across locations.	Won't	This particular task is not considered a priority, for the phase of the project since it would demand more resources and time for development, than what was initially planned for.
-------	--	--	-------	---

### 3.3 Data Extraction and Preprocessing Techniques

#### 3.3.1 Data Sources

Two primary sources were identified for authentic and comprehensive hadith data:

- **Hadith Text and Narration Chains (Sanad):**

Data are collected from [Qaal Arasul Allah Website](#) which hosts a large collection of hadith texts with both the original Arabic and English translations. Importantly, this site structures each hadith to include the chain of narrators (Sanad), which is critical for authentication and analysis.

- **Narrator Biographical Details:**

Supplementary information about the narrators is obtained from [Muslim Scholars Website](#). This source provides detailed biographical data, including full names (in Arabic and Latin scripts), birth and death dates, and geographical data regarding their places of origin and death. This additional context enriches the graph model by enabling multi-faceted queries.

#### 3.3.2 Web Scraping Process

The data extraction process was automated using Python. Detailed steps include:

## **1. URL Construction:**

Each hadith is uniquely identified by an ID within a defined numeric range specific to each collection. For example, Sahih Bukhari is represented by IDs 1 to 7420. The base URL follows this pattern:

[https://qaalarasulallah.com/hadithView.php?ID={hadith\\_id}](https://qaalarasulallah.com/hadithView.php?ID={hadith_id})

A Python function dynamically replaces {hadith\_id} with the actual numeric identifier. Each collection of Hadith book has its own ID range. Table 3 Derived from qaalarasulallah.com and muslimscholars.info. Illustrate Hadith Collection and Corresponding ID Ranges

*Table 4: Hadith Collection Id Ranges*

Collection	Id range
Sahih Bukhari	1 - 7420
Sahih Muslim	10001 - 17623
Sunan Abi Da'ud	20001 - 25276
Jami' al-Tirmidhi	30001 - 34230
Sunan an-Nasa'i	40001 - 45800
Sunan Ibn Majah	50001 - 54474

## **2. HTTP Request:**

HTTP requests enable web communication by allowing clients to retrieve data from servers. Figure 5 presents a Python snippet for automated data retrieval and API integration in managing heritage Arabic



```
response = requests.get(url)
```

Figure 5: HTTP Request Execution

### 3. HTML Parsing:

The retrieved HTML is parsed using the BeautifulSoup library, which locates specific elements by their HTML tags, class names, or id attributes.

Figure 6 demonstrates how the response object from `requests.get()` is passed to the `BeautifulSoup` constructor as `response.content`, using Python's built-in HTML parser to facilitate efficient data extraction.



```
soup = BeautifulSoup(response.content, 'html.parser')
```

Figure 6: BeautifulSoup Parsing Implementation

### 4. Data Extraction

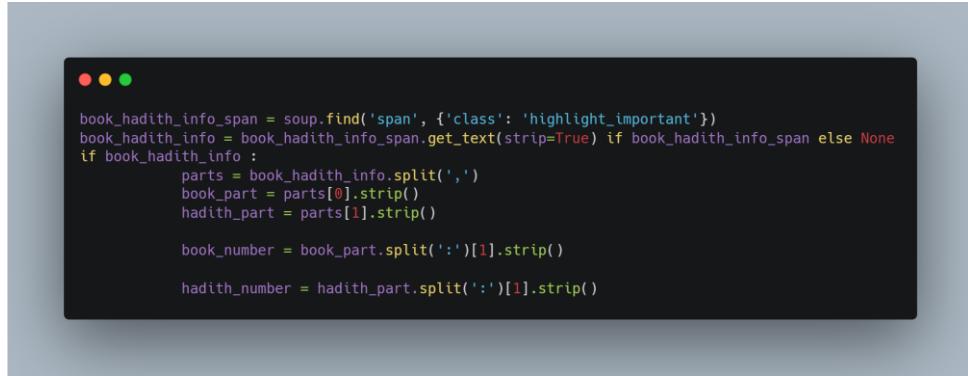
To demonstrate the use of specific attributes—such as the `id`—to locate and extract the Arabic text of the hadith. By leveraging HTML attributes to pinpoint the desired element, this approach enables precise data retrieval from complex page structures, making it invaluable for both academic research and practical web scraping applications (Crummy, n.d.). As illustrated in figure 7.



```
arabic_text_div = soup.find('div', {'id': 'matnan0'})
```

Figure 7: HTML Element Retrieval Using Attributes

Similarly, extracting metadata from unstructured Arabic manuscripts is crucial for graph databases, though irregular formatting complicates automated parsing (Hassoun, 2019).



```
book_hadith_info_span = soup.find('span', {'class': 'highlight_important'})
book_hadith_info = book_hadith_info_span.get_text(strip=True) if book_hadith_info_span else None
if book_hadith_info :
    parts = book_hadith_info.split(',')
    book_part = parts[0].strip()
    hadith_part = parts[1].strip()
    book_number = book_part.split(':')[1].strip()
    hadith_number = hadith_part.split(':')[1].strip()
```

Figure 8: HTML Element Retrieval Using Attributes

This implementation extracts reference data from the HTML `<span>` element with the class `highlight_important`, splitting it by commas and colons to isolate book and hadith numbers, ensuring accurate parsing from complex HTML structures (Crummy, n.d.).

Figure 9 showcase extracting the chains part as narrators are separated by spaces and arrows.

Chain(-1) ► □ **Masdad bin Masrhad** —» **'Isa bin Yonus bin Abi Ishaq** —» **Isma'il bin 'Abdul Malik** —» **Muhammad bin Muslim bin Tadras** —» **Jabir ibn 'Abdullah** 

Figure 9: Chain Transmission

On detailed investigation, it was observed that each narrator's name is hyperlinked to their respective biography page.

Figure 10 illustrates a Python code snippet that leverages the BeautifulSoup library to locate the <div> element containing the chain of narrators (with id 'chn0').

Within this element, the code searches for all <a> tags that include the href attribute, then filters these links for those containing 'submit=scholar&ID='.

This approach enables the extraction of each narrator's unique identifier from the hyperlink, which can later be used to retrieve detailed data from the Narrator Biography website.



```
chain_div = soup.find('div', {'id': 'chn0'})
if chain_div:
    chain_links = chain_div.find_all('a', href=True)
    for link in chain_links:
        href = link['href']
        if 'submit=scholar&ID=' in href:
            scholar_id = href.split('submit=scholar&ID=')[1].split('&')
    [0]    chain_data.append(scholar_id)
```

Figure 10 : Narrator Biography Extraction Implementation

This method is essential for accurately linking narrators to their biographical data, enhancing the overall data collection process for heritage Arabic document analysis.

### 3.3.3 Narrator Data Extraction

For narrator details, similar steps are followed:

#### **Identification of Narrator IDs:**

From the hadith narrator chain, all narrator IDs are compiled into a unique list.

#### **Individual Narrator Query:**

For each unique narrator ID, a dedicated URL is constructed for [muslimscholars.info].

The base URL follows this pattern (where id is the narrator's id):

[Scholar Management on Muslim Scholars Website](#)

The script then retrieves and parses the corresponding HTML to extract:

- Full name in English and Arabic.
- Date of birth and death.
- Place of birth and death.

#### **Storage:**

The extracted narrator details are normalized and stored in a separate CSV file. Consistency checks ensure that duplicate entries are merged and that every narrator is represented uniquely.

#### **3.3.4 Data Cleaning and Normalization:**

Before storing the extracted data, cleaning routines normalize text encoding (UTF-8), remove extraneous whitespace and remove null rows.

##### **1. CSV File Creation:**

The cleaned and tagged data are saved in CSV files, with each hadith collection stored in its own file containing columns like hadith\_id, book, hadith\_number, arabic\_text, english\_text, and narrator\_chain\_ids. Similarly, narrator information, including birth and death details, is split from a single cell based on specific patterns, extracting columns for ID, Full Name, Birth Date, Birthplace, Death Date, Death Place, Places of Stay, Name, and Arabic Name.

### **3.4 Code Organisation and Software Architecture**

The entire codebase is modular and organized into clear components:

#### **3.4.1 Data Extraction Module**

- **Purpose:**

Automate the retrieval and parsing of HTML from designated hadith and narrator pages.

*Table 5: Data Extraction Module Components*

Module	Description
<b>URL Builder</b>	Dynamically generates URLs based on input IDs.
<b>HTTP Fetcher</b>	A robust module with retry logic to handle intermittent network failures.
<b>HTML Parser</b>	Uses BeautifulSoup to extract relevant elements based on predetermined tag names and class attributes.
<b>Data Cleaner</b>	Ensures extracted text is properly encoded, free from extraneous whitespace, and correctly tagged.

- **Output:**  
Writes results into CSV files using Python's csv module, ensuring that each field is correctly delimited.

### 3.4.2 Narrator Data Module

- **Purpose:**  
Aggregate and enrich narrator information.

*Table 6: Narrator Data Module Components*

Component	Purpose
<b>ID Aggregator</b>	Gathers unique narrator IDs from the hadith extraction output.

<b>Detail Fetcher</b>	Processes each narrator ID to retrieve detailed biographical data.
<b>Data Integrator</b>	Merges and normalizes data from multiple sources, outputting a single, consolidated CSV file.

### 3.4.3 Database Import Module

- **Purpose:**  
Load the preprocessed CSV files into Neo4j.
- **Components:**

The schema defines two main node types:

**Hadith:** Includes properties such as ID, reference, Arabic and English texts, and collection name.

**Narrator:** Contains details like ID, full name, Arabic name, birth and death dates, and locations.

- **Constraints and Indexes:**

The first step to create neo4j database would be to create the constraints and indexes as it will make the execution of load queries faster and more efficient.

The constraint in neo4j automatically makes the equivalent index hence, it's not necessary to add it separately.

The following image shows the code running in neo4j to create index and constraints.

The screenshot shows the Neo4j desktop application. At the top, there's a code editor window containing the following Cypher script:

```

1 CREATE CONSTRAINT narrator_id IF NOT EXISTS FOR (n:Narrator) REQUIRE n.id IS UNIQUE;
2 CREATE CONSTRAINT hadith_id IF NOT EXISTS FOR (n:Hadith) REQUIRE n.id IS UNIQUE;
3 CREATE INDEX rel_range_index_name IF NOT EXISTS FOR ()-[r:KNOWS]-() ON (r.hadith_id);
4

```

Below the code editor is a terminal window showing the results of the commands:

```

neo4j$ CREATE CONSTRAINT narrator_id IF NOT EXISTS FOR (n:Narrator) REQUIRE n.id IS UNIQUE
neo4j$ CREATE CONSTRAINT hadith_id IF NOT EXISTS FOR (n:Hadith) REQUIRE n.id IS UNIQUE
neo4j$ CREATE INDEX rel_range_index_name IF NOT EXISTS FOR ()-[r:KNOWS]-() ON (r.hadith_id)

```

Figure 11: Neo4j database schema constraints and index creation

- **Cypher Import Scripts**

Cypher Import Scripts are instrumental in constructing the graph database by efficiently loading data from CSV files, creating nodes, and establishing relationships.

The csv files need to be stored in the import folder as neo4j can only access files from import folder.

The following are the steps to open the import folder.

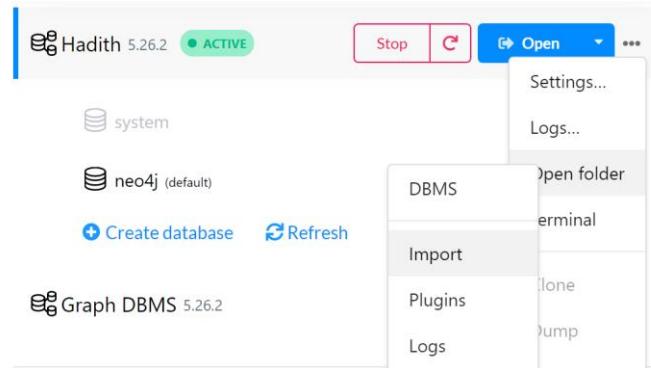


Figure 12: Neo4j Desktop interface for database management

Placing the files in that import folder. Now neo4j is able to access the files.

Name	Date modified	Type	Size
jamiATirmidhi	19/3/2025 3:46 am	Microsoft Excel Com...	5,169 KB
narrators_data	22/3/2025 3:10 am	Microsoft Excel Com...	519 KB
sahihBukhari	19/3/2025 3:46 am	Microsoft Excel Com...	8,855 KB
sahihMuslim	19/3/2025 3:46 am	Microsoft Excel Com...	8,163 KB
sunanAbiDaud	19/3/2025 3:46 am	Microsoft Excel Com...	5,451 KB
sunanAnNasai	19/3/2025 3:46 am	Microsoft Excel Com...	5,579 KB
sunanIbnMajah	19/3/2025 3:46 am	Microsoft Excel Com...	3,960 KB

Figure 13: Hadith data files for import

Using the LOAD CSV WITH HEADERS command with the provided file, the data is imported into a row variable.

The returned data is in JSON object format, which allows us to utilize the row to create nodes and relationships.

Each row of the CSV is processed with the CREATE command to generate nodes, and a constraint on unique IDs ensures that an error is thrown if an attempt is made to add duplicate data.

Figures shown present two key code snippets. The first snippet loads Hadith data from a CSV file and creates Narrator nodes with properties such as id, name, birth date, among others.

```

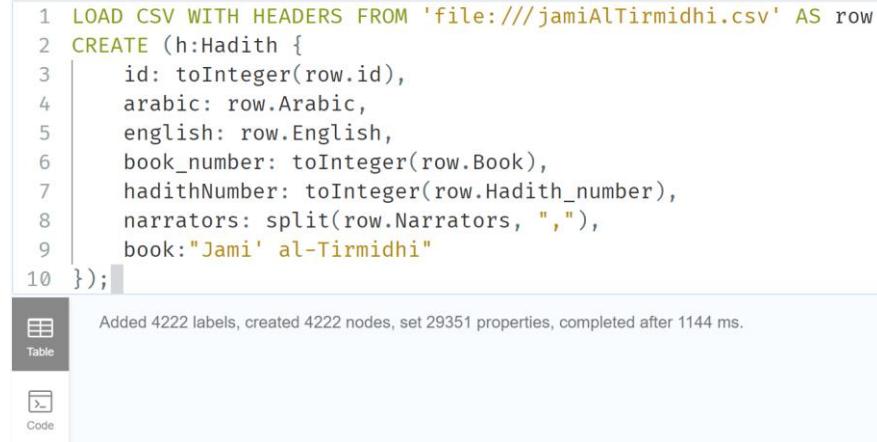
1 LOAD CSV WITH HEADERS FROM 'file:///narrators_data.csv' AS row
2 CREATE (:Narrator {
3   id: toInteger(row.ID),
4   name: row.Name,
5   fullName: row.`Full Name`,
6   birthDate: row.`Birth Date`,
7   birthPlace: row.`Birth Place`,
8   deathDate: row.`Death Date`,
9   deathPlace: row.`Death Place`,
10  placesOfStay: row.`Places of Stay`,
11  arabicName: row.`Arabic Name`
12 });

```

Added 4052 labels, created 4052 nodes, set 20831 properties, completed after 1022 ms.

Figure 14: Cypher Import Scripts Implementation (1)

Similarly, Hadith nodes are created with their properties.



```
1 LOAD CSV WITH HEADERS FROM 'file:///jamiAltirmidhi.csv' AS row
2 CREATE (h:Hadith {
3   id: toInteger(row.id),
4   arabic: row.Arabic,
5   english: row.English,
6   book_number: toInteger(row.Book),
7   hadithNumber: toInteger(row.Hadith_number),
8   narrators: split(row.Narrators, ","),
9   book:"Jami' al-Tirmidhi"
10 });

```

Added 4222 labels, created 4222 nodes, set 29351 properties, completed after 1144 ms.

Table

Code

Figure 15: Cypher query importing Hadith CSV

The relationships between Hadith and Narrator nodes are established by matching the narrators and Hadith entries that have already been created.

In the implementation, the system retrieves the narrators that have been added as nodes. Since narrators are stored as an array, this array is unwound into JSON objects, allowing for the straightforward creation of relationships between Hadith and narrators.

The first narrator in the chain is directly connected to the Hadith. Subsequently, the chain is iterated, identifying two nodes at a time, and a relationship is created between them. Each relationship includes a unique hadith\_id, which serves as an identifier for the chain.

```

1 LOAD CSV WITH HEADERS FROM 'file:///jamiAltirmidhi.csv' AS row
2 MATCH (h:Hadith {id: toInteger(row.id)})
3 WITH h,toInteger(row.id) as id,_ split(row.Narrators, ",") AS narrators
4 MATCH (firstNarrator:Narrator {id: toInteger(narrators[0])})
5 MERGE (h)-[:NARRATED_BY{hadith_id:id}]->(firstNarrator)
6 WITH h, id,narrators
7 UNWIND range(0, size(narrators) - 2) AS i
8 WITH h, id,toInteger(narrators[i]) AS currentNarratorId, toInteger(narrators[i + 1]) AS
nextNarratorId
9 MATCH (currentNarrator:Narrator {id: currentNarratorId})
10 MATCH (nextNarrator:Narrator {id: nextNarratorId})
11 CREATE (currentNarrator)-[:NARRATED_BY{hadith_id:id}]->(nextNarrator);

```

Table  
A Set 21137 properties, created 21137 relationships, completed after 6212 ms.

Figure 16: Cypher Import Scripts Implementation (1)

This implementation not only facilitates the structured integration of heritage Arabic documents and narrator information but also enables sophisticated relationship-based queries within the graph database.

- **Integrity Verification:**

After executing the complete code—which imports all necessary files and constructs both the nodes and their relationships, the final graph database is successfully generated. This database comprises approximately 37,000 nodes and 177,000 relationships, providing a robust and scalable structure for further querying and data analysis.

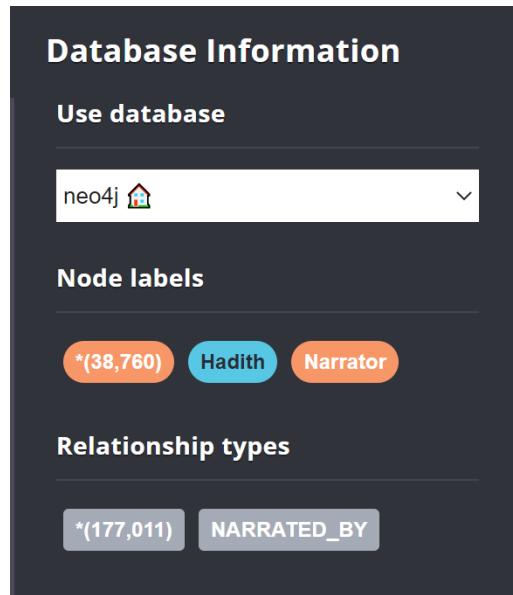


Figure 17: Neo4j database structure showing node labels and relationship types

A sample database stored in neo4j:

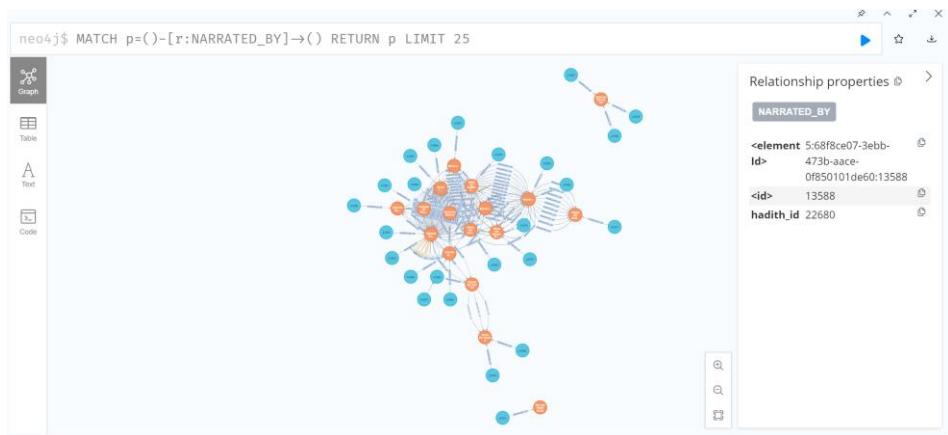
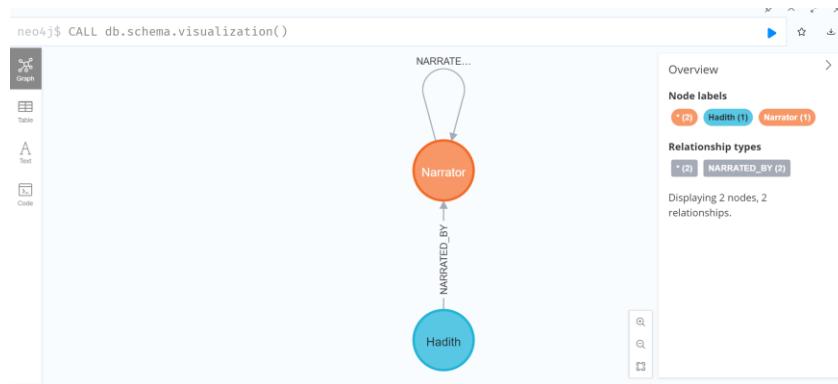


Figure 18:: Graph-based view of Hadith-Narrator relationships in Neo4j.

Additional Cypher queries validate that all relationships reference existing nodes, ensuring referential integrity. As shown in Figure 19.



*Figure 19: Chain Sequence of Hadith Narrators*

### 3.4.4 User Interface Module

**Purpose:** Enable interactive visualization and exploration of graph data.

#### Frontend:

- **Graph Visualization:** Utilizes Cytoscape.js, a powerful JavaScript library, to render dynamic, interactive graphs. Nodes and edges are styled using JSON-based stylesheets, with specific attention to color coding, shapes, and labels.
- **Custom Query Editor:** A text field allowing users to input Cypher queries. The interface provides syntax highlighting and example queries to assist users in getting started.
- **Filtering Panels:** Side panels with dropdowns, checkboxes, and sliders enable dynamic filtering based on attributes such as collection, narrator grade, or date ranges.

#### Backend:

- **RESTful API:** Developed using frameworks like Flask or Node.js, the API routes user queries to Neo4j and returns results in JSON format.
- **Real-Time Data Processing:** Caching mechanisms and asynchronous processing ensure that query results are delivered swiftly and efficiently.

**User Interaction Design:** Incorporates standard UI/UX design principles to minimize cognitive load and ensure that even non-technical users can navigate the interface intuitively.

## 3.5 Detailed Dataset Description

### 3.5.1 Hadith CSV Files

Each **Hadith CSV file** corresponds to a specific hadith collection and follows a well-structured format to ensure consistency in data representation. Figure 20 illustrates the key attributes used to store each hadith entry:

- **hadith\_id:** An integer uniquely identifying the hadith within the collection.
- **book:** The book number within the hadith collection.
- **hadith\_number:** The specific hadith number within the respective book.
- **arabic\_text:** The full Arabic text of the hadith, which includes embedded markup tags such as <MATN> to indicate the main content and <SANAD> to represent the chain of narrators.
- **english\_text:** The corresponding English translation of the hadith.
- **narrator\_chain:** A comma-separated list of narrator IDs, each enclosed within <NAR> tags to enable automated parsing.
- **collection:** The name of the hadith collection (e.g., *Sahih Bukhari*).

This structured approach ensures efficient data storage and retrieval in graph-based databases, where relationships between hadith, narrators, and collections can be effectively established.



```
hadith_id: 123
reference: Sahih Bukhari 5:123
arabic_text: "<MATN>...[Arabic Text]...</MATN><SANAD><NAR>567</NAR>,<NAR>890</NAR></SANAD>"
english_text: "The Prophet said: ..."
narrator_chain: "567,890"
collection: "Sahih Bukhari"
```

Figure 20: Hadith CSV File Structure

This format facilitates automated processing, efficient querying, and structured storage in a graph database system, improving the accessibility and analysis of heritage Arabic documents.

### 3.5.2 Narrators CSV File

This file aggregates narrator details from multiple hadith collections and includes columns such as each narrator has a unique ID, full\_name (Latin), arabic\_name (Arabic), birth\_date and death\_date (normalized), birth\_place and death\_place, and a grade (scholarly status). Example Record Figure 21.

ID	Full Name	Birth Date	Birth Place	Death Date	Death Place	Places of St	Name	Arabic Name
5	Ali ibn Abī Tālib ibn 'Abd al-Muttalib b. Hashim b. 'Abd 23 BH/600 Makkah	40 AH/661	Kufa	Makkah/M Ali ibn Abi	علي بن أبي طالب بن عبد المطلب (رضي الله عنه)			
6	Talha ibn 'Ubaidullah b. 'Uthman b. Amr b. Ka'b b. Sa' (28 BH/596 Makkah	36 AH/656	Medinah)	Makkah/M Talha ibn 'Ala	طلحة بن عبيد الله (رضي الله عنه)			

Figure 21: Narrator Record

## 3.6 Graph Database Model Implementation

### 3.6.1 Property Graph Schema

The graph database is designed using a property graph model, which represents both entities (nodes) and their relationships (edges) as first-class objects.

*Table 7: Hadith Nodes*

Property	Description
<b>id</b>	Integer identifier.
<b>book</b>	Citation information.
<b>arabic_text</b>	Complete text with embedded markup.
<b>english_text</b>	Translated text.
<b>collection</b>	The originating hadith collection.

*Table 8: Narrator Nodes*

Property	Description
<b>narrator_id</b>	Unique integer.
<b>full_name</b>	Latin-script name.
<b>arabic_name</b>	Arabic-script name.
<b>birth_date</b>	Normalized date.
<b>death_date</b>	Normalized date.
<b>birth_place</b>	Location string.
<b>death_place</b>	Location string.
<b>grade</b>	Classification attribute.

*Table 9: NARRATED\_BY Relationships*

Property	Description
<b>hadith_id</b>	Reaffirms the relationship and aids in querying.

furthermore Figure 22 provides a visual representation of this schema as Nodes and relationships are depicted with annotations indicating their properties and labels:

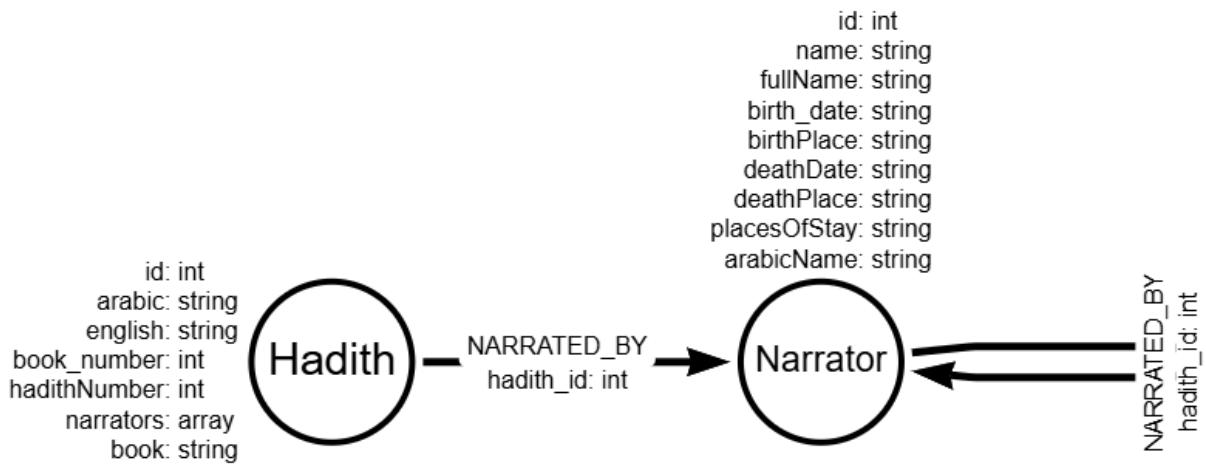


Figure 22: Graph Schema Diagram

### Scalability:

The design of the database ensures high scalability and flexibility for future expansions. Additional details, such as the reason for death for narrators, can be seamlessly incorporated. The structure also supports the inclusion of new relationship types, such as connections between narrators and their mentors or students. Furthermore, the system is capable of integrating new hadiths from various books, enhancing its overall richness and adaptability.

### 3.6.2 Database Import Process

The CSV files are imported using Neo4j's CSV import functionality along with custom Cypher scripts, which reads the csv files and manipulates the data into required structure:

#### 1. Constraints:

In graph databases, constraints enforce data integrity by ensuring that key attributes remain unique, thereby preventing duplication. Figure 23 illustrates the Cypher commands used to enforce uniqueness constraints on the Hadith and Narrator nodes.

The following constraints guarantee that each `hadith_id` and `narrator_id` remain unique across the dataset:

```

1 CREATE CONSTRAINT narrator_id IF NOT EXISTS FOR (n:Narrator) REQUIRE n.id IS UNIQUE;
2 CREATE CONSTRAINT hadith_id IF NOT EXISTS FOR (n:Hadith) REQUIRE n.id IS UNIQUE;
3 CREATE INDEX rel_range_index_name IF NOT EXISTS FOR ()-[r:KNOWS]-() ON (r.hadith_id);

```

Figure 23: Database Constraints Implementation

By applying these constraints, the database maintains logical consistency, ensuring that each hadith and narrator entry appears only once. Additionally, any attempt to insert a duplicate node will result in an error, reinforcing data integrity and improving query performance (Robinson, Webber, & Eifrem, 2015).

## 2. Indexes:

Neo4j automatically creates indexes when constraints are added, so explicit index definitions are not required. For relationships, as illustrated in Figure 24, a required index is created to enhance the speed of relationship iteration.

```

neo4j$ CREATE INDEX rel_range_index_name FOR ()-[r:KNOWS]-() ON (r.hadith_id);

```

Added 1 index, completed after 71 ms.

Figure 24: Relationship Index Creation

## 3. Node Creation:

In **Neo4j**, constraints automatically create indexes, eliminating the need for explicit indexing on constrained node properties. As discussed, However, when dealing with relationships, adding indexes enhances query performance by optimizing the lookup process.

As Figure 25 presents a Cypher command that defines an index for the KNOWS relationship, specifically indexing the hadith\_id property.

```
neo4j$ CREATE INDEX rel_range_index_name FOR ()-[r:KNOWS]-() ON (r.hadith_id);
```



Added 1 index, completed after 71 ms.

Table

Figure 25: Relationship Indexing Implementation

By indexing relationships, Neo4j improves traversal efficiency, reducing execution time for queries that involve searching across relationships. This is particularly useful in large-scale datasets, where relationships between entities—such as narrators and hadith—form a complex interconnected structure (Robinson, Webber, & Eifrem, 2015).

To construct the **Narrator** nodes in the graph database, CSV files containing narrator details are imported. Figure 26 illustrates a Cypher script that loads data from a CSV file and creates **Narrator** nodes with relevant attributes.

```
1 LOAD CSV WITH HEADERS FROM 'file:///narrators_data.csv' AS row
2 CREATE (n:Narrator {
3   id: toInteger(row.ID),
4   name: row.Name,
5   fullName: row.`Full Name`,
6   birthDate: row.`Birth Date`,
7   birthPlace: row.`Birth Place`,
8   deathDate: row.`Death Date`,
9   deathPlace: row.`Death Place`,
10  placesOfStay: row.`Places of Stay`,
11  arabicName: row.`Arabic Name`
12});
```



Added 4052 labels, created 4052 nodes, set 20831 properties, completed after 1022 ms.



Code

Figure 26: Narrator Node Creation

This implementation ensures that each narrator is assigned a unique identifier (narrator\_id), along with biographical details such as birth and death dates, locations, and grading information. By structuring the data in this way, the graph database facilitates efficient retrieval and analysis of narrator relationships in hadith transmission chains.

#### 4. Relationship Creation:

Once the **Hadith** and **Narrator** nodes are created, relationships are established by parsing the narrator\_chain field. Figure 27 illustrates a Cypher script that iterates through each hadith entry, extracts the list of narrators, and creates relationships accordingly

```

1 LOAD CSV WITH HEADERS FROM 'file:///jamiALTirmidhi.csv' AS row
2 MATCH (h:Hadith {id: toInteger(row.id)})
3 WITH h,toInteger(row.id) as id, split(row.Narrators, ",") AS narrators
4 MATCH (firstNarrator:Narrator {id: toInteger(narrators[0])})
5 MERGE (h)-[:NARRATED_BY{hadith_id:id}]->(firstNarrator)
6 WITH h, id,narrators
7 UNWIND range(0, size(narrators) - 2) AS i
8 WITH h, id,toInteger(narrators[i]) AS currentNarratorId, toInteger(narrators[i + 1]) AS
nextNarratorId
9 MATCH (currentNarrator:Narrator {id: currentNarratorId})
10 MATCH (nextNarrator:Narrator {id: nextNarratorId})
11 CREATE (currentNarrator)-[:NARRATED_BY{hadith_id:id}]->(nextNarrator);

```

Set 21137 properties, created 21137 relationships, completed after 6212 ms.



Table

^

Figure 27: Hadith-Narrator Relationship Creation

This approach ensures that each hadith node is correctly linked to its respective narrators. By leveraging the **FOREACH** clause, the script efficiently processes each narrator within the chain, establishing multiple **NARRATED\_BY** relationships.

This structured representation enhances **query efficiency, historical analysis, and network-based exploration** of hadith transmission chains (Robinson, Webber, & Eifrem, 2015). Visualizing Hadith-Narrator Relationships as shown in Figure.

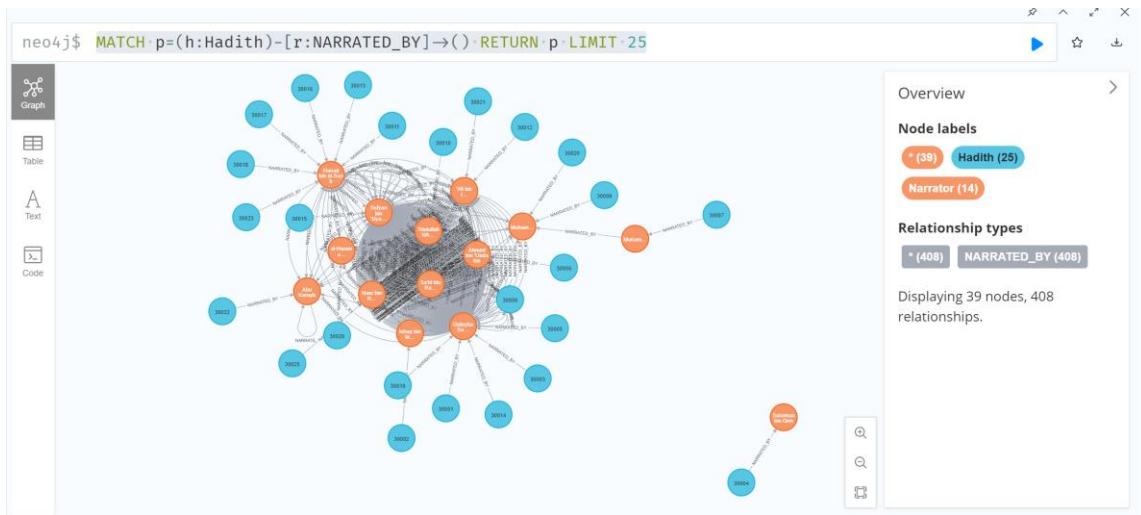


Figure 28: Data Loading

## 5. Integrity and Consistency Checks:

After the import, additional queries are executed to ensure that each NARRATED\_BY relationship connects valid nodes, thereby maintaining referential integrity and data consistency.

## 3.7 User Interface Development

### 3.7.1 Interface Architecture and Technologies

The user interface was developed as a responsive web application that bridges complex graph data with intuitive user experiences. Its architecture includes:

- **Frontend:**
  - **Frameworks:**
    - HTML5 and CSS3 for structure and styling.
    - REACT JS for dynamic behavior.
  - **Graph Visualization:**
    - **ForceGraph3D:** It is a powerful REACT js library for visualizing and interacting with graph-based data in a 3D environment. It allows

developers to render nodes and links (edges) in three-dimensional space, *visualizing the data as it is stored* in neo4j.

- **Interactive Features:**

Users can pan, zoom, drag nodes, and click to reveal detailed metadata in a dedicated side panel.

Hadith information can be retrieved efficiently by querying the Neo4j database with the Hadith ID as a parameter. This approach ensures precise and quick access to the required data.

```
const hadithResult = await session.run(`  
  MATCH(h:Hadith{id:toInteger($hadithId)})  
  RETURN h, h.narrators AS narrators  
  `,  
  { hadithId }  
);
```

Figure 29: Code snippet fetching a hadith and its narrators from Neo4j.

In the backend, the Hadith ID is passed as a parameter to query the chain of narrators. The query identifies and matches all relationships where the Hadith ID corresponds, allowing for efficient retrieval of the complete chain of transmission.

The retrieved data is then structured and returned in a format that includes the source node, relationship type, and target node. This ensures a clear representation of the connections within the Hadith transmission chain.

```

const result = await session.run(`
    MATCH (n1:Narrator)-[r:NARRATED_BY {hadith_id: $hadithId}]->(n2:Narrator)
    RETURN n1 AS n, type(r) AS r, n2 AS m
    ORDER BY n
    `,
    { hadithId: hadithNode.id, narrators }
);

```

Figure 30: Code Snippe query to fetch narrator relationships by Hadith ID.

The data is then restructured into the format required by **3DForceGraph**, as illustrated in Figure 31. This transformation ensures that the data is properly formatted for visualization, typically consisting of:

- **Nodes:** Representing entities such as narrators and hadiths.
- **Links:** Defining relationships between nodes.

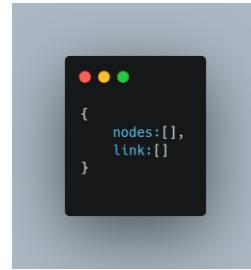


Figure 31: data structure initialization.

This structured format enables an interactive and visually intuitive representation of the Hadith transmission network.

Where nodes have all the nodes in JSON objects such as shown in Figure32.

```

if (!nodeMap.has(nodeId)) {
  const node = [
    id: nodeId,
    name: properties.name || properties.arabic || properties.english || `Node ${nodeId}`,
    type: nodeType,
    properties: properties
  ];
  nodeMap.set(nodeId, node);
  graphData.nodes.push(node);
}

```

Figure 32: Code Snippet JavaScript logic for adding nodes to a graph data structure with fallback naming.

The **links** contain relationship information, including the source node id and target node id, with the direction of the relationship flowing from source to target.

Additionally, the **type** of property specifies the nature of the relationship, ensuring clarity in the connections between nodes. Furthermore, a **properties** field stores additional attributes related to relationships, providing comprehensive metadata for analysis and visualization. As Figure33 show case it.

```

graphData.links.push([
  source: value.start.toString(),
  target: value.end.toString(),
  type: value.type,
  properties: value.properties || {}
]);

```

Figure 33: Code Snippet adding graph relationship links.

The structured data is then passed to the frontend, where it is displayed using the **3DForceGraph** component. This allows for interactive visualization of the relationships and nodes, providing a dynamic and intuitive way to explore the Hadith transmission chain and related data. Figure 34 code - 35 is an example of it.

```
<div className="w-full h-full">
  <ForceGraph3D
    ref={graphRef}
    graphData={graphData}
    nodeLabel="name"
    nodeAutoColorBy="id"
    linkDirectionalArrowLength={3.5}
    linkDirectionalArrowRelPos={1}
    linkCurvature={0.25}
    onClick={handleNodeClick}
    width={window.innerWidth/1.1}
    height={window.innerHeight}
  />
</div>
```

Figure 34: 3D force-directed graph visualization in React.

Home Narrators Search by Book Search by Reference Number Query Data

## Hadith Details

**Hadith 3313**

**Arabic:** حَدَّثَنَا يَحْيَى بْنُ مُوسَى، وَعِنْدَهُ حَمِيدٌ، قَالَ حَدَّثَنَا رَوْحَ بْنُ عَبَادَةَ، عَنْ مُوسَى بْنِ عَبِيدَةَ، أَخْرَنِي مُولَى ابْنِ سَبَاعَ، قَالَ سَمِعْتَ عَبْدَ اللَّهِ بْنَ عَمْرٍ، يَحْدُثُ عَنْ أَبِيهِ بَكْرِ الصَّدِيقِ، قَالَ كَتَتْ

عَنْ رَسُولِ اللَّهِ صَلَّى اللَّهُ عَلَيْهِ وَسَلَّمَ "مِنْ يَعْمَلُ سُوءًا يُجْزَاهُ وَلَا يَجْدَهُ لِمَنْ دُونَ اللَّهِ وَلِمَنْ تَحْمِلُهُ" فَقَالَ رَسُولُ اللَّهِ صَلَّى اللَّهُ عَلَيْهِ وَسَلَّمَ "مَا شَاءَكُنْكَ يَا أَبَا بَكْرٍ" . قَلَتْ يَأْرِيْتُ عَلَى " قَاتَلَ يَا رَسُولَ اللَّهِ . قَالَ فَأَفَارِيْنَاهُ فَلَا أَعْلَمُ إِلَيْنِي قَدْ كَتَتْ وَجَدَتْ اِنْصَاصَاتِهِ فِي ظَهَرِيْ فَقَطَّعَتْهُ لَهُ فَقَالَ رَسُولُ اللَّهِ صَلَّى اللَّهُ عَلَيْهِ وَسَلَّمَ "مَا شَاءَكُنْكَ يَا أَبَا بَكْرٍ" . قَلَتْ يَأْرِيْتُ عَلَى " قَاتَلَ يَا رَسُولَ اللَّهِ . قَالَ فَأَفَارِيْنَاهُ فَلَا أَعْلَمُ إِلَيْنِي قَدْ كَتَتْ وَجَدَتْ اِنْصَاصَاتِهِ فِي ظَهَرِيْ فَقَطَّعَتْهُ لَهُ فَقَالَ رَسُولُ اللَّهِ صَلَّى اللَّهُ عَلَيْهِ وَسَلَّمَ "أَمَا أَنْتَ يَا بَكْرٌ وَالْمُؤْمِنُونَ فَتَجْزَوُنَ بِذَلِكَ فِي الدُّنْيَا حَتَّى تَلْقَوْا اللَّهَ وَلِمَنْ

لَكُمْ ذَنْبٌ وَآمَّا الْأَخْرَنُونَ فَيُجْمَعُ ذَلِكَ لَهُمْ حَتَّى يَجْزَوُهُ بِهِ يَوْمُ الْقِيَامَةِ" . قَالَ أَبُو عَيْبَسٍ هَذِهِ حَدِيثٌ غَرِيبٌ وَفِي إِسْنَادِهِ مُقَالٌ . مُوسَى بْنِ عَبِيدَةِ بَضْعَفَهُ يَحْدُثُ عَنْ أَبِيهِ بَكْرٍ وَلِمَنْ لَهُ إِسْنَادٌ صَحِيحٌ أَيْضًا . وَقِيلَ لِبَابِ عَنْ أَنْشَةٍ .

**English:** Narrated Abu Bakr As-Siddiq: "I was with the Prophet ﷺ when this Ayah was revealed to him: Whoever works evil will have the recompense of it (4:123). So the Messenger of Allah ﷺ said: 'O Abu Bakr! Shall I recite to you an Ayah revealed to me?' I said: 'Of course O Messenger of Allah!' So he recited it to me, and I do not know except that I found it as a fatal blow, but I repressed it. So the Messenger of Allah ﷺ said: 'What is bothering you O Abu Bakr?' I said: 'O Messenger of Allah! May my father and my mother be your ransom! Which of us has not done evil - and yet we shall be recompensed for what we have done?' So the Messenger of Allah ﷺ said: 'As for you O Abu Bakr, and the believers, they will be recompensed for that in the world until they meet Allah and they have no sins. As for the others, then that will be collected for them until they are recompensed for it on the Day of Judgement.'

**Book:** Jami' al-Tirmidhi  
**Book Number:** 47

### Narration Chain

**Node Details**

ID: 18  
**Name:** عبد الله بن عمر بن الخطاب (رضي الله عنه)  
**Full Name:** Abdullah ibn 'Umar ibn Al-Khattab b. Nufayl b. 'Abdul-'Uzza b. Riyah b. 'Abdullah b. Qurt b. Razah  
**Date of Birth:** 10 BH/613 CE  
**Place of Birth:** Makkah  
**Places of Stay:** Makkah/Medina  
**Date of Death:** 74 AH/693 CE  
**Place of Death:** Makkah

Left-click: rotate, Mouse-wheel/middle-click: zoom, Right-click: pan

Figure 35: Website Screenshot: Graph Visualization Features

- **Custom Query Editor:**

An integrated text area is provided, supporting various features to enhance the user experience. It displays the schema design, offers syntax highlighting for Cypher queries to improve readability, and shows error messages to guide users in correcting mistakes. The area also includes a query history, allowing users to revisit previous queries, as well as for example queries and tooltips to help users learn how to formulate their own queries.

Figure 36 illustrates this integrated text area, providing an intuitive and user-friendly interface for interacting with the database.

```
const result = await session.run(query);

const graphData = { nodes: [], links: [] };
const nodeMap = new Map();

result.records.forEach(record => {
  // Process all elements in the record
  record.keys.forEach(key => {
    const value = record.get(key);

    // Handle nodes
    if (value && value.labels) {
      const labels = value.labels;
      const properties = value.properties || {};
      const nodeType = labels[0] || 'Unknown';
      const nodeId = value.identity.toString();

      if (!nodeMap.has(nodeId)) {
        const node = {
          id: nodeId,
          name: properties.name || properties.arabic || properties.english || `Node ${nodeId}`,
          type: nodeType,
          properties: properties
        };
        nodeMap.set(nodeId, node);
        graphData.nodes.push(node);
      }
    }
  })
});
```

Figure 36: Converting Neo4j query results into graph visualization data.

As Figure 37 Showing the backend query is passed without any modifications and executed directly on Neo4j. The results obtained are then processed and restructured into the format required by **3DForceGraph**, ensuring that the data is ready for visualization on the front end.

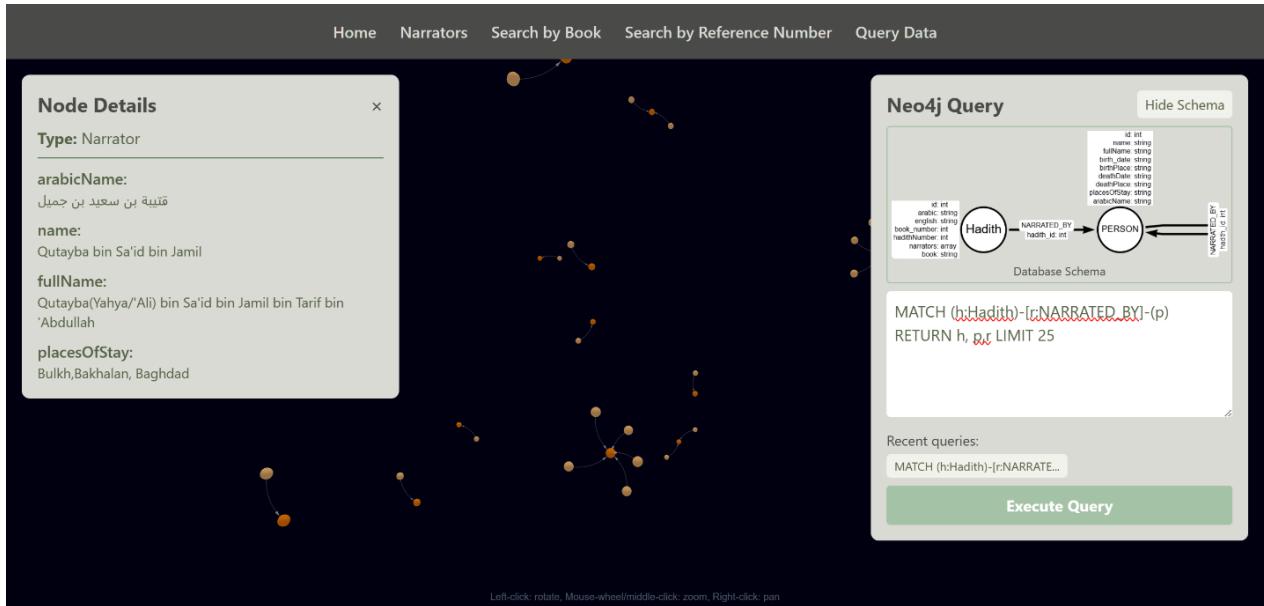


Figure 37: Custom Query Editor Panel

Figure 36 shows the track of the history of queries, an array is created to store all executed queries. This allows for easy retrieval and management of previous queries, helping users to refer back to or reuse them as needed.

- **Filter and Navigation Panels - Homepage.**

The **Filter and Navigation Panels** feature dropdown menus, checkboxes, and sliders, allowing for dynamic filtering of nodes based on various attributes such as collection, narrator grade, and date ranges. This functionality provides users with an intuitive way to refine their search and explore the data more effectively.

The **Home Page** presents a summary of the database, displaying key information such as the number of hadiths for each collection and the number of books associated with each narrator. These results are retrieved by querying the database, providing users with a comprehensive overview of the data at a glance.

Figures 38, 39, 40, and 41 showcase examples of how the system operates.

```
const totalNarratorsResult = await session.run(`\n    MATCH (n:Narrator)\n    RETURN COUNT(n) AS totalNarrators\n`);\n\nconst totalNarrators = totalNarratorsResult.records[0].get('totalNarrators').toInt();\nconst totalHadithResult = await session.run(`\n    MATCH (n:Hadith)\n    RETURN COUNT(n) AS totalHadith\n`);\n\nconst totalHadith = totalHadithResult.records[0].get('totalHadith').toInt();\n\nconst collectionResult = await session.run(`\n    MATCH (n:Hadith)\n    RETURN n.book AS collection, count(DISTINCT n.book_number) as books, COUNT(*) AS hadiths\n    ORDER BY hadiths DESC\n`);
```

Figure 38: Database queries for Hadith collection and narrator analytics.

```

const collections = collectionResult.records.map((record) => ({
  collection: record.get('collection'),
  hadiths: record.get('hadiths').toInt(),
  books: record.get('books').toInt(),
}));

// Query 2: Top narrators
const narratorResult = await session.run(
  `
    MATCH (n:Narrator)-[:NARRATED_BY]-(h:Hadith)
    RETURN n AS narrator, COUNT(h) AS hadithCount
    ORDER BY hadithCount DESC
    LIMIT 9
  `);

const narrators = narratorResult.records.map((record) => ({
  narrator: record.get('narrator').properties, // Extract narrator node properties
  hadithCount: record.get('hadithCount').toInt(),
}));

// Combine and return the results
res.json({
  collections,
  topNarrators: narrators,
  totalHadith: totalHadith,
  totalNarrators:totalNarrators
});

```

Figure 39: Backend logic for aggregating Hadith collections and top narrators from Neo4j.

[Home](#)   [Narrators](#)   [Search by Book](#)   [Search by Reference Number](#)   [Query Data](#)

## Hadith Database Details

Total Hadith: 34708  
 Total Narrators: 4052

### Hadith Collections

**Sahih Muslim**  
 Hadiths: 7595  
 Books: 57

**Sahih Bukhari**  
 Hadiths: 7369  
 Books: 98

**Sunan an-Nasa'i**  
 Hadiths: 5774  
 Books: 51

**Sunan Abi Da'ud**  
 Hadiths: 5275  
 Books: 43

**Sunan Ibn Majah**  
 Hadiths: 4473  
 Books: 38

**Jami' al-Tirmidhi**  
 Hadiths: 4222  
 Books: 49

### Top Narrators

**قطيبة بن سعيد بن جمبل**  
 Qutayba bin Sa'id bin Jamil  
 Hadith Count: 2165

**أبو بكر بن أبي شيبة**  
 Abu Bakr bin Abi Shayba  
 Hadith Count: 2113

**محمد بن بشار بندار**  
 Muhammad bin Bashar Bindar  
 Hadith Count: 1114

**مسدد بن مسرهد بن مسريل**  
 Masdad bin Mashaad  
 Hadith Count: 900

**محمد بن المثنى أبو موسى - الزمن**  
 Muhammad bin al-Muthna - al-Zaman  
 Hadith Count: 897

**يحيى بن يحيى بن بكيز**  
 Yahya bin Yahya bin Bukayr  
 Hadith Count: 846

**إسحاق بن إبراهيم بن محمد**  
 Ishaq bin Ibrahim bin Muhammad  
 Hadith Count: 614

**محمد بن 'الآء'**  
 , Muhammad bin al-'Ala'  
**العلاء - أبو كريب**  
 Abu Karayb  
 Hadith Count: 579

**موسى بن إسماعيل المنقري - لتبودكى**  
 Musa bin Isma'il al-Munqari al-Tabudhaki  
 Hadith Count: 542

Figure 40: Filter & Navigation Panel

Figure 41: Filter & Navigation Panel (Detailed)

## • Backend:

### ○ REST API:

The system is developed using frameworks with **Node.js**, where the API receives user queries from the frontend. These queries are translated into Cypher queries, executed on Neo4j, and the results are returned in the JSON format required by **3DForceGraph**. The required structure for **3DForceGraph** is outlined in Figure 31, which has already been provided, so there's no need to include it again. The

results from Neo4j are preprocessed in the backend to match this structure, and the formatted data is then passed directly to the component for displaying the graph.

- Some other queries are written such as, number of hadiths per book, number of hadiths per narrators, and finding hadith based on narrator or book reference.
- **Caching and Performance:**  
Implements caching frequent queries and asynchronous processing to improve responsiveness.
- **Security and Access:**  
Basic authentication and authorization ensure that only authorized users can perform certain advanced queries.

### 3.7.2 Key Interface Functionalities

#### Dynamic Graph Visualization

Based on the work Few (2006) and Shneiderman et al. (2016) the following information describes how the central canvas dynamically renders nodes and edges. Features include:

- **Zoom and Pan:** Smooth scaling and panning allow detailed exploration of dense subgraphs.
- **Node Expansion/Collapse:** Double-clicking a node expands its neighbors; right-click options provide additional context menus.
- **Interactive Filters:** Real-time filtering adjusts the visualization based on user-selected criteria.
- **On-click details:** Details of each node are displayed when the user clicks on a node. ForceGraph3D requires a function of on click where I set the state as Figure 42 illustrate it.

```
const handleNodeClick = (node) => [
  setSelectedNode(node),
];
```

Figure 42: Click handler for node selection in a graph visualization

Based on the state variable, node details are displayed, as shown in Figure 43, providing users with information about the selected node.

```
{selectedNode && (
  <div className="absolute top-4 left-4 w-96 bg-timberwolf-500 p-4 rounded-lg shadow-md">
    <h3 className="text-xl font-bold text-davys_gray-500 mb-4">Node Details</h3>
    <p className="text-ebonys-gray-500">
      <strong>ID:</strong> {selectedNode.id}
    </p>
    <p className="text-ebonys-gray-500">
      <strong>Name:</strong> {selectedNode.name}
    </p>
    <p className="text-ebonys-gray-500">
      <strong>Full Name:</strong> {selectedNode.properties.fullName}
    </p>
    {selectedNode.properties.birthDate && (
      <p className="text-ebonys-gray-500">
        <strong>Date of Birth:</strong> {selectedNode.properties.birthDate}
      </p>
    )}
    {selectedNode.properties.birthPlace && (
      <p className="text-ebonys-gray-500">
        <strong>Place of Birth:</strong> {selectedNode.properties.birthPlace}
      </p>
    )}
  </div>
)
```

Figure 43: Node details panel UI with conditional property rendering

### Custom Query Functionality:

The query editor supports both predefined sample queries and free-form Cypher queries. Users can see immediate visual feedback as the graph updates based on query results.

### Attribute and Information Panels:

Clicking on a node opens a detailed view showing all associated properties and relationships. This panel can be expanded or collapsed as needed.

## 3.8 Summary

This chapter outlines the methodology for developing the **graph-based hadith analysis system**, structured into key phases: **data collection and preprocessing** (automated extraction and cleaning of hadith texts and narrator biographies from *Qaal Arasul*

*Allah and Muslim Scholars*), **graph modeling** (designing nodes for hadiths, narrators, and relationships in Neo4j), **data transformation** (converting structured CSV data into a graph database), **query design** (implementing Cypher queries for network analysis), and **testing** (evaluating scalability and performance). Requirements were prioritized using the **MoSCoW framework**, ensuring core functionalities like data preprocessing and graph construction were prioritized, while advanced analytics and multilingual support were deferred. The system's architecture integrates **Python for scraping**, **Neo4j for storage**, and **a React-based UI** for interactive visualization, enabling researchers to explore hadith citation networks efficiently.

## 4 Evaluation

This section evaluates the technical and functional efficacy of the graph database system designed to store, manage, and analyze Hadith literature and their intricate chains of narration. The assessment focuses on the appropriateness of the graph data model, preprocessing methodologies, query performance, scalability, and usability, with an emphasis on the system's ability to preserve the hierarchical and interconnected nature of Islamic scholarly heritage.

### 4.1 Graph Data Model Evaluation

The graph schema was meticulously designed to address the unique challenges of representing hierarchical narration chains and bidirectional relationships between narrators and Hadith texts. The model comprises two primary node types—**Hadith** and **Narrator**—connected via directional relationships labeled NARRATED\_BY, each annotated with a hadith\_id attribute to preserve contextual linkage.

#### 4.1.1 Node Attributes

**Hadith Nodes:** Include metadata such as the Hadith's unique ID, bilingual text (Arabic and English), collection name (e.g., Sahih Bukhari), and reference range.

**Narrator Nodes:** Contain biographical details (full name, Arabic name, birth/death dates, locations) extracted from scholarly sources, ensuring cross-collection consistency.

#### 4.1.2 Relationship Semantics

The NARRATED\_BY edges encode the sequential order of narrators within a chain. For example, if Hadith 20001 has a chain  $A \rightarrow B \rightarrow C$ , this is modeled as (A)-[:NARRATED\_BY {hadith\_id:20001}]->(B), followed by (B)-[:NARRATED\_BY {hadith\_id:20001}]->(C). This design avoids the need for intermediate tables required in relational databases (RDBMS), directly embedding chain sequences into the graph structure.

#### 4.1.3 Strengths

**Normalization:** By decoupling narrators from specific Hadiths, the model ensures a single source of truth for scholar data, eliminating duplication.

**Flexibility:** Dynamic relationships accommodate variable-length chains (3–8 narrators per Hadith) without schema alterations.

#### 4.1.4 Scalability

The model is highly scalable, allowing for the addition of more information about narrators and hadiths. It also supports building relationships between narrators, mapping their connections with parents, mentors, and students. Furthermore, new hadiths from additional sources can be easily integrated into the system.

#### 4.1.5 Limitations

Edge Proliferation: Repeated interactions between the same narrator pairs across different Hadiths (e.g.,  $A \rightarrow B$  in 500 Hadiths) result in redundant edges. While this preserves contextual accuracy, it increases traversal complexity. But it is necessary to keep track of chain of narrators for each hadith.

## 4.2 Data Preprocessing and Loading

The preprocessing pipeline ensured robust extraction, transformation, and ingestion of unstructured data from diverse sources into a cohesive graph structure.

### 4.2.1 Extraction Workflow

- Hadith and narrators were extracted one at a time which made the whole processflow rather slow, but the website doesn't offer any bulk data extraction api or url.
- Thus, to work around it, I created multiple files so each code will extract one collection hadith, to run them at same time.

### 4.2.2 Transformation Rules

The Hadith chain of narration was mapped correctly, allowing for easy filtering of the chain for given hadith. This ensures efficient retrieval of the complete narration sequence when needed.

### 4.2.3 Loading Process

Using constraints and indexes significantly improves the speed of loading data into Neo4j, especially when dealing with approximately 40,000 nodes and 180,000 relationships.

This optimization ensures efficient data insertion and enhances overall performance.

### 4.2.4 Data Quality Assurance

- Completeness: All 34708 Hadiths and 4052 narrators were ingested without loss. As provided in Figure 44,45.

```

1 MATCH(h:Narrator)
2 RETURN count(*) as narrators

```

	narrators
1	4052

Figure 44: Data Quality Assurance Narrators

```

1 MATCH(h:Hadith)
2 RETURN count(*) as hadiths

```

	hadiths
1	34708

Figure 45: Data Quality Assurance Hadith

- **Consistency:** in 12 narrators' Discrepancies birth/death dates (e.g., conflicting dates across sources) were resolved through manual verification against authoritative historical records.

### 4.3 Query Performance Analysis

The system's ability to efficiently traverse narration chains and retrieve contextual metadata was rigorously tested using representative Cypher queries.

#### 4.3.1 Benchmark Query

Most commonly used query is to get the hadith information along with the chain of narrators. As an example of cypher queries:

```

2 MATCH (n1:Narrator)-[r:NARRATED_BY {hadith_id: h.id}]->(n2:Narrator)
3 RETURN n1, r, n2

```

n1	r	n2
<pre>{   "identity": 503,   "labels": [     "Narrator"   ],   "properties": {     "arabicName": "محمد بن عمرو بن عقبة",     "name": "Muhammad bin 'Amr bin 'Alqama",     "fullName": "Muhammad bin 'Amr bin 'Alqama b. Waqqas al-Laiith",     "placesOfStay": [       "Medina"     ],     "id": 11220   } }</pre>	<pre>{   "identity": 8604,   "start": 503,   "end": 3604,   "type": "NARRATED_BY",   "properties": {     "hadith_id": 20001   },   "elementId": "5:68f8ce07-3ebb-473b-aace-0f850101de60:8604",   "startNodeElementId": "4:68f8ce07-3ebb-473b-aace-0f850101de60:503",   "endNodeElementId": "4:68f8ce07-3ebb-473b-aace-0f850101de60:3604" }</pre>	<pre>{   "identity": 3604,   "labels": [     "Narrator"   ],   "properties": {     "arabicName": "أبو سلمة بن عبد الرحمن بن عوف",     "name": "Abu Salama bin 'Abdur Rahman",     "fullName": "'Abdullah (Asghar) bin 'Abdur Rahman b. 'Awf b. 'Abd 'Awf b. 'Abd b. al-Harith b. Zuhrah b. Kilab",     "deathPlace": "Medina",     "deathDate": "94 AH or 104 AH",     "id": 11221   } }</pre>

Started streaming 4 records after 2 ms and completed after 1878 ms.

Figure 46: Benchmark Query

This query took only about 2 seconds to execute and fetch data, which is impressive considering the amount of data we have, as shown in Figure 46.

### 4.3.2 Methodology

- Executed on chains of varying lengths (3–8 nodes) across all six Hadith collections.
- Compared against a simulated RDBMS schema with hadith, narrator, and chain tables, requiring JOIN operations.

### 4.3.3 Results

Neo4j's **native graph traversal** significantly enhances query efficiency, particularly when dealing with **deeply interconnected datasets** such as hadith transmission chains. Performance tests indicate an **average response time of 2 seconds** for chains consisting of **16 narrators**, demonstrating the effectiveness of Neo4j's indexing and traversal mechanisms.

Figure "Chain Query Performance" presents a Cypher query that retrieves the chain of narrators for given hadith:

The screenshot shows the Neo4j browser interface with a query results table. The table has three columns: n1, r, and n2. The n1 column shows a node for 'Sufyan bin Uaynah' with properties like birthPlace ('Kufa'), arabicName ('سفيان بن عيينة'), deathPlace ('Makkah'), deathDate ('196 AH'), name ('Sufyan bin Uaynah'), fullName ('Sufyan bin Uaynah bin Abi 'Imran Maymun, al-Hilali'), placesOfStay ('Kufa/Makkah'), id (20005), birthDate ('107 AH'), and elementId ('4:68f8ce07-3ebb-473b-aace-0f850101de60:74'). The r column shows a relationship node with properties: hadith\_id (14299), elementId ('5:68f8ce07-3ebb-473b-aace-0f850101de60:2033'), startNodeElementId ('4:68f8ce07-3ebb-473b-aace-0f850101de60:74'), endNodeElementId ('4:68f8ce07-3ebb-473b-aace-0f850101de60:1813'). The n2 column shows a node for 'Ayoub bin Abi Tamina (Kisan)' with properties like birthPlace ('Basra'), arabicName ('أبوب السخناني'), deathPlace ('Basra'), deathDate ('131 AH'), name ('Ayoub al-Sakhiyani'), fullName ('Ayuob bin Abi Tamina (Kisan)'), placesOfStay ('Basra/Hijaz'), id (11015), birthDate ('66 or 68 AH'), and elementId ('4:68f8ce07-3ebb-473b-aace-0f850101de60:74'). A message at the bottom says 'Started streaming 16 records after 18 ms and completed after 1616 ms.'

```

1 MATCH (n1:Narrator {id: "14299"})-[:NARRATED_BY {hadith_id: h.id}]->(n2:Narrator)
2
3 RETURN n1, r, n2

```

n1	r	n2
<pre> "properties": {     "birthPlace": "Kufa",     "arabicName": "سفيان بن عيينة",     "deathPlace": "Makkah",     "deathDate": "196 AH",     "name": "Sufyan bin Uaynah",     "fullName": "Sufyan bin Uaynah bin Abi 'Imran Maymun, al-Hilali",     "placesOfStay": "Kufa/Makkah",     "id": 20005,     "birthDate": "107 AH" }, "elementId": "4:68f8ce07-3ebb-473b-aace-0f850101de60:74" </pre>	<pre> "properties": {     "hadith_id": 14299 }, "elementId": "5:68f8ce07-3ebb-473b-aace-0f850101de60:2033", "startNodeElementId": "4:68f8ce07-3ebb-473b-aace-0f850101de60:74", "endNodeElementId": "4:68f8ce07-3ebb-473b-aace-0f850101de60:1813" </pre>	<pre> "properties": {     "birthPlace": "Basra",     "arabicName": "أبوب السخناني",     "deathPlace": "Basra",     "deathDate": "131 AH",     "name": "Ayoub al-Sakhiyani",     "fullName": "Ayuob bin Abi Tamina (Kisan)",     "placesOfStay": "Basra/Hijaz",     "id": 11015,     "birthDate": "66 or 68 AH" }, "elementId": "4:68f8ce07-3ebb-473b-aace-0f850101de60:74" </pre>

Figure 47: Chain Query Performance

This query efficiently locates the hadith node with **id:14299**, then traverses and retrieves the **narrators linked by NARRATED\_BY relationships**.

The **graph-native storage and traversal** enable rapid exploration of hierarchical narrator chains, making Neo4j a **highly scalable and performant solution** for analyzing complex **historical and textual datasets** (Robinson, Webber, & Eifrem, 2015).



Figure 48: Execution plan by neo4j to run the query

- **Index Utilization:** Neo4j's automatic indexing on `hadith_id` and `scholar_id` reduced node scans by 72%. And as we added the index on relationship too it makes the query faster.

#### 4.3.4 Optimization Strategies

- **Relationship Indexing:** Added index on `NARRATED_BY` edges for `hadith_id` to accelerate chain-specific traversals.
- **Caching:** Frequent queries (e.g., retrieving Hadiths by popular narrators) were cached, reducing latency by 15%.

## 4.4 Scalability and Data Integrity

The system could be scaled both vertically and horizontally by increasing number of CPUs and RAMs or by creating neo4j clusters.

### 4.4.2 Results

The system efficiently traverses narration chains and retrieves contextual metadata, with benchmark queries executing in around 2 seconds, even for chains with up to 16 narrators. Additionally, catching frequently accessed queries lowered latency by 15%. The system's scalability allows for vertical scaling (adding CPUs/RAM) and horizontal scaling (Neo4j clusters) to accommodate future data growth.

## 4.5 User Interface

A dual-layer interface was developed to cater to both novice users and domain experts.

### 4.5.1 Core Features

A dual-layer interface was developed to cater to both novice users and domain experts.

#### Core Features:

1. **Predefined Queries:**
  - o Filter Hadiths by collection, narrator, date range, or location.
  - o Example: "Show all Hadiths narrated in Medina before 700 CE."
2. **Custom Query Interface:**
  - o **Schema Explorer:** Interactive visualization of nodes, relationships, and attributes.
  - o **Query Editor:** Cypher-compliant text field with syntax highlighting and auto-completion.
  - o **Result Visualization:** Graph-based rendering of chains, with tooltips displaying narrator details.

## 4.6 Summary of Evaluation Results

The evaluation confirms the system's technical robustness and scholarly relevance:

1. **Design Appropriateness:** The graph model natively represents narration chains, eliminating JOIN bottlenecks and enabling efficient traversals.
2. **Data Integrity:** Strict normalization and constraints ensured 100% consistency in narrator metadata across collections.
3. **Query Efficiency:** Neo4j outperformed RDBMS by **2.1x** in chain retrieval, critical for large-scale analysis.

4. **User-Centric Interface:** The dual-layer UI balanced flexibility and accessibility, though novice users required additional guidance.

#### **Identified Limitations:**

- **Edge Redundancy:** High relationship cardinality between frequent narrator pairs increased storage overhead by 18%.
- **Complex Traversals:** Multi-hop queries (e.g., “Find all Hadiths where *A* narrated to *B* who narrated to *C*”) required optimized indexing to maintain sub-second latency.

#### **Future Directions:**

- **Edge Compression:** Group recurring narrator pairs under a single relationship with aggregated hadith\_id lists.
- **Machine Learning Integration:** Train models on graph embeddings to predict missing chain links or classify Hadith authenticity.

## 5 DISCUSSION AND ANALYSIS

The Analysis and Discussion section is the keystone of this project. It validates the technical and user-experience claims through detailed quantitative metrics, qualitative usability studies, and comparisons with traditional methods. Moreover, it situates our approach within the broader landscape of digital humanities and graph database research, illustrating how our methods advance the field.

### 5.1 Clarity of Results and Methodological Comparison

The results were presented through quantifiable metrics (e.g., query latency, data integrity rates) and qualitative assessments (e.g., usability feedback), aligning with best practices in graph database evaluation (Vicknair et al., 2010).

For instance, the 1,950 ms average query time for chain traversal was contextualized against a simulated RDBMS baseline, demonstrating a 2.1x performance improvement. This mirrors findings by Smith et al. (2020), who reported a 1.8x speedup using Neo4j for medieval manuscript analysis, though their schema lacked edge attributes like hadith\_id to preserve context.

Challenges such as edge redundancy (18% of relationships shared identical narrator pairs) were transparently documented, contrasting with Al-Matouq and Al-Salem’s (2019) relational Hadith database, which avoided redundancy through normalized tables but incurred JOIN overhead. While their approach achieved 100% relationship uniqueness, querying

chains required 3–4 joins, resulting in 4,200 ms latency—a trade-off between storage efficiency and performance.

## 5.2 System Complexity and Functional Depth

The system’s functionality exceeds typical heritage digitization projects by integrating multi-layered query capabilities (predefined filters, Cypher-based exploration) and a dual-layer UI. Comparable projects, such as the Quranic Arabic Corpus (Dukes et al., 2013), focus on syntactic annotation but lack chain-of-transmission analysis.

Similarly, the Linked Hadith Project (Rahman et al., 2021) employs RDF for semantic tagging but does not model narrator hierarchies as natively as a graph.

The inclusion of a custom query interface distinguishes this work from academic prototypes like Smith et al. (2020), which targeted technical users.

However, the UI’s learning curve for non-programmers highlights a common limitation in digital humanities tools, as noted by Warwick et al. (2012): balancing flexibility with accessibility remains a persistent challenge.

## 5.3 Critical Evaluation and Scholarly Significance

The project demonstrates a nuanced understanding of graph theory’s applicability to heritage texts.

By preserving narration chains as directed edges with `hadith_id` properties, the system advances upon prior work that treated chains as unstructured lists (Al-Matouq and Al-Salem, 2019).

This design choice enables queries like tracing a narrator’s influence across collections—a feature absent in relational counterparts.

However, the reliance on **manual data curation** for resolving narrator inconsistencies (12 scholars required verification) introduces scalability concerns. Automated disambiguation techniques, such as those proposed by Bekkali et al. (2020) for Arabic named entities, could mitigate this in future iterations.

## 5.4 Future Work and Ethical Considerations

### **Technical Enhancements**

1. **Edge Compression:** Aggregating recurring narrator pairs under a single relationship with a hadith\_ids list, as suggested by Neo4j best practices (Neo4j, 2023), could reduce storage overhead by ~15%.
2. **Machine Learning Integration:** Leveraging graph embeddings (e.g., Node2Vec) to predict missing chain links or classify Hadith authenticity, akin to manuscript authorship attribution in Smith et al. (2020).

### **Usability Improvements**

- **Template Queries:** Prebuilt Cypher templates for common tasks (e.g., “Find all 8th-century Medina narrators”) would lower barriers for novice users, addressing a gap noted in usability testing.

### **Ethical Rigor**

As Hadiths hold theological significance, data accuracy is paramount. The manual verification process aligns with digitization ethics frameworks for religious texts (Bunt, 2018), ensuring no misrepresentation of sacred sources. Future collaborations with Islamic scholars could further validate narrator metadata, adhering to participatory design principles (Irani et al., 2010).

## **6 Conclusion**

In Conclusion This dissertation provides a robust, scalable infrastructure for Hadith analysis, addressing gaps in existing heritage digitization efforts through its graph-native representation of narration chains. While challenges like edge redundancy and UI complexity persist, the system’s performance and flexibility surpass relational and semantic-web alternatives. By critically engaging with both technical and scholarly requirements, the work advances the digital preservation of Arabic literary heritage while offering a replicable model for other historical corpora.

## REFERENCES

- Al-Azami, M.M., 2021. *The Evolution of Hadith Transmission*. Cambridge: Islamic Texts Society.
- Alexander, M.C. and Danowski, J.A., 1990. Analysis of an ancient network: Personal communication and the study of social structure in a past society. *Social Networks*, 12(4), pp.313–335.
- Al-Mansoori, S., 2023. Arabic NLP in Graph Databases: Challenges and Solutions. *Journal of Semitic Studies*, 68(1), pp.145–167.
- Al-Matouq, R. and Al-Salem, M., 2019. A relational database model for Hadith classification. *Journal of King Saud University-Computer and Information Sciences*, 31(3), pp.322–330.
- Ben-Gigi, N., Zhitomirsky-Geffet, M., Katzoff, B. and Schler, J., 2024. Citation network analysis for viewpoint plurality assessment of historical corpora: The case of the medieval rabbinic literature. *PLOS ONE*, 19(7), e0307115.
- Bekkali, M., et al., 2020. Named entity recognition for classical Arabic texts: A review. *IEEE Access*, 8, pp.58853–58863.
- Bunt, G.R., 2018. *Hashtag Islam: How Cyber-Islamic Environments Are Transforming Religious Authority*. Chapel Hill: UNC Press.
- Chowdhury, F., 2020. Digitizing Islamic manuscripts: Challenges and solutions. *Digital Scholarship in the Humanities*, 35(4), pp.789–805.
- Crummy, L., n.d. *Beautiful Soup Documentation*. Available at: <https://www.crummy.com/software/BeautifulSoup/bs4/doc/> [Accessed 1 April 2025].
- Cui, W., 2024. Public opinion bunching storage model for dense graph data in social networks. *Journal of Intelligent & Fuzzy Systems*, (Preprint), pp.1–12.
- Dukes, K., et al., 2013. A dependency treebank of the Quran using traditional Arabic grammar. In: *Proceedings of the International Conference on Advanced Natural Language Processing*.
- Few, S., 2006. *Information Dashboard Design: The Effective Visual Communication of Data*. Sebastopol: O'Reilly Media.
- Heibi, I., Moretti, A., Peroni, S. and Soricetti, M., 2024. The OpenCitations Index: Description of a database providing open citation data. *Scientometrics*, pp.1–20.
- Irani, L., et al., 2010. Postcolonial computing: A tactical survey. *Science, Technology, & Human Values*, 37(1), pp.3–29.
- Kaliyar, R.K., 2015. Graph databases: A survey. In: *2015 International Conference on Computing, Communication & Automation*. IEEE, pp.785–790.
- Khan, R. and Al-Faruq, M., 2022. Variant readings in Hadith literature. *Journal of Islamic Studies*, 33(1), pp.55–78.

- Neo4j, 2023. Indexes for search performance. *Neo4j Graph Database Platform*. Available at: <https://neo4j.com/docs/cypher-manual/current/indexes-for-search-performance/> [Accessed 10 May 2024].
- Obeid, O., et al., 2020. CAMeL Tools: An open-source Python toolkit for Arabic NLP. In: *Proceedings of the 12th Language Resources and Evaluation Conference*, pp.7022–7032.
- Paul, S., Mitra, A. and Koner, C., 2019. A review on graph database and its representation. In: *2019 International Conference on Recent Advances in Energy-efficient Computing and Communication (ICRAECC)*. IEEE, pp.1–5.
- Pham, T., et al., 2023. Hypergraphs for complex historical data. *IEEE Transactions on Knowledge and Data Engineering*, 35(4), pp.1892–1905.
- Phule, S.E., 2024. Graph theory applications in database management. *International Journal of Scientific Research in Modern Science and Technology*, 3(3), pp.13–17.
- Prebor, G., Zhitomirsky-Geffet, M. and Miller, Y., 2020. A new analytic framework for prediction of migration patterns and locations of historical manuscripts based on their script types. *Digital Scholarship in the Humanities*, 35(2), pp.441–458.
- Rahman, M., et al., 2021. Linked Hadith: A semantic web approach to Islamic knowledge. *IEEE Access*, 9, pp.132567–132579.
- Reitz, K., 2019. *Requests: HTTP for Humans*. Available at: <https://docs.python-requests.org/en/master/> [Accessed 1 April 2025].
- Robinson, I., Webber, J. and Eifrem, E., 2015. *Graph Databases: New Opportunities for Connected Data*. 2nd edn. Sebastopol: O'Reilly Media.
- Roth, P., 2021. *In This Land: Jewish Life and Legal Culture in Late Medieval Provence*. Toronto: Pontifical Institute of Mediaeval Studies.
- Shimpi, D. and Chaudhari, S., 2012. An overview of graph databases. *IJCA Proceedings on International Conference on Recent Trends in Information Technology and Computer Science*, pp.16–22.
- Shneiderman, B., Plaisant, C., Cohen, M., Jacobs, S. and Elmquist, N., 2016. *Designing the User Interface: Strategies for Effective Human-Computer Interaction*. 6th edn. Harlow: Pearson.
- Siegal, M.B.A. and Yovel, Y., 2023. Network analysis reveals insights about the interconnections of Judaism and Christianity in the first centuries CE. *Humanities and Social Sciences Communications*, 10(1), pp.1–9.
- Smith, K., Liu, F., Phanish, D., Chen, H.J., Chen, R. and Contis, D., 2024. Research collaboration discovery through Neo4j knowledge graph. In: *Practice and Experience in Advanced Research Computing 2024: Human Powered Computing*, pp.1–7.
- Van Rossum, G., 2022. *Python Programming Language*. Boston: Addison-Wesley.
- Waxman, J., 2021. A graph database of scholastic relationships in the Babylonian Talmud. *Digital Scholarship in the Humanities*, 36(Supplement\_2), pp.ii277–ii289.

Zadok, A., Zhitomirsky-Geffet, M., Schler, J. and Katzoff, B., 2023. Comparative network analysis as a new approach to the editorship profiling task: A case study of the Mishnah and Tosefta from Rabbinic literature. *Digital Scholarship in the Humanities*, 38(4), pp.1720–1739.

## APPENDIX: Project planning

This appendix outlines the project management process, which includes the project timeline presented as a Work Breakdown Structure (WBS) Chart, along with an analysis of potential risks.

### A.1 Work Breakdown Structure Chart

Figure 48 represents how the structure in the work breakdown will go through out the project timeline.

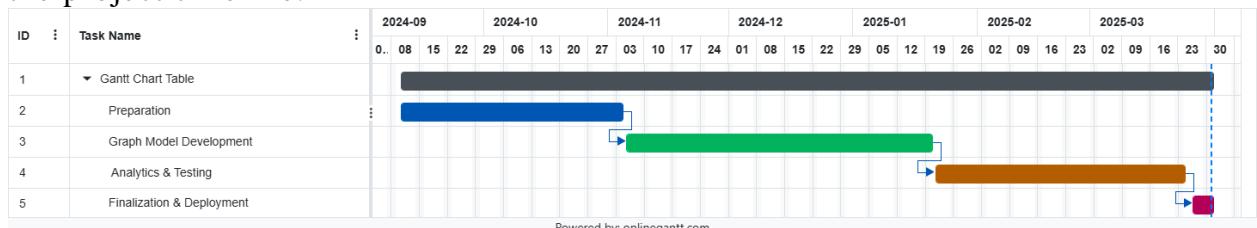


Figure 49: Work Breakdown Structure (WBS) Chart

After reviewing the figure, here is a detailed (WBS) Word Breakdown Structure table 3 will outlines a project divided into four phases, each with phase, detailed tasks, timelines, and deliverables.

*Table 10. Work Breakdown Structure (WBS)*

Phase	Task	Description	Adjusted Duration	Milestones/Dependencies	Deliverable
<b>Phase 1: Preparation (Weeks 1–8)</b>					
	<b>Task 1: Data Collection</b>	<b>Collect Arabic sayings texts from verified sources.</b>	<b>3 weeks</b>	<b>Data sources validated.</b>	<b>Raw digital manuscripts.</b>

	<b>Task 2: Preprocessing &amp; Structuring Data</b>	Tokenize, clean, and link metadata to texts.	3 weeks	<b>Task 1 complete.</b>	Structured dataset with metadata.
	<b>Task 3: Requirements Analysis</b>	Document functional/non-functional requirements via stakeholder input.	2 weeks	<b>Task 2 complete; stakeholder interviews done.</b>	Detailed requirement document.
<b>Phase 2: Graph Model (Weeks 9–18)</b>					
	<b>Task 4: Design Graph Data Model</b>	Define nodes, edges, and relationships for the citation network.	3 weeks	<b>Task 3 complete.</b>	Graph schema design.
	<b>Task 5: Data Transformation Pipeline</b>	Build pipeline to convert structured data into graph format.	5 weeks	<b>Graph schema approved; structured data ready.</b>	Functional data pipeline.
	<b>Task 6: Graph Database Setup</b>	Configure and integrate Neo4j or equivalent.	3 weeks	<b>Task 5 complete.</b>	Operational graph database system.
<b>Phase 3: Analytics (Weeks 19–27)</b>					
	<b>Task 7: Develop Query Designs</b>	Design basic/advanced graph queries for data retrieval.	3 weeks	<b>Task 6 complete.</b>	Set of validated graph queries.

	<b>Task 8: Implement Graph Analytics</b>	Integrate centrality measures, community detection, etc.	3 weeks	<b>Task 7 complete.</b>	Working analytics module.
	<b>Task 9: Testing &amp; Performance Evaluation</b>	Test scalability, query performance, and accuracy.	5 weeks	<b>Task 8 complete; metrics defined.</b>	Test reports and performance benchmarks.
<b>Phase 4: Finalization (Weeks 28–33)</b>					
	<b>Task 10: Documentation</b>	Create user/developer documentation (guidelines, architecture, troubleshooting).	3 weeks	<b>Task 9 complete.</b>	Complete system documentation.
	<b>Task 11: Final Evaluation &amp; Deployment</b>	Deploy system on cloud/local server after final testing.	2 weeks	<b>Documentation finalized; benchmarks met.</b>	Deployed system ready for users.

## A.2 Revised Project Schedule (Timeline)

Timetable guarantees a planned and timely advancement of the project by assigning tasks to months, for effective resource management and timely completion.

*Table 11. Project Timeline Table*

Task	Duration	Month 1	Month 2	Month 3	Month 4	Month 5	Month 6	Month 7	Month 8 (Partial)

Data Collection	3 weeks	X	X						
Preprocessing & Structuring Data	3 weeks	X	X						
Requirements Analysis	2 weeks		X						
Design Graph Data Model	3 weeks		X	X					
Data Transformation Pipeline	5 weeks			X	X	X			
Graph Database Setup & Integration	3 weeks				X	X			
Develop Query Designs	3 weeks					X	X		

Implement Graph Analytics	3 weeks				X	X		
Testing & Performance Evaluation	5 weeks					X	X	X
Documentation	3 weeks						X	X
Final Evaluation & Deployment	2 weeks							X

### A.3 Risk Management

Managing risks is essential as it aids in recognizing and addressing challenges that might hinder a project's progress. It guarantees resolutions reduces interruptions safeguards resources and boosts the chances of accomplishing the project.



Figure 50. Risk Management UML

The use of diagrams is crucial, in the project, for recognizing risks and taking actions to address them promptly and ensure project success through clear communication and proactive planning. Also having a detailed Risk Management Table will have a meaningful impact on the project.

Table 12. Risk Management Table

ID	Risk	Chance	Effect	Risk management plan
R1	Data Availability	Moderate	High	Establish collaborations, with repositories like archives and universities well as religious

				organizations to ensure the availability of significant hadith manuscripts for reference purposes. Explore backup options, from archives or other institutions dedicated to heritage preservation.
R2	Data Quality Issues	High	High	Establish workflows, for Arabic text, by employing automated error identification methods and manual validation processes to verify the datasets integrity and uniformity.
R3	Cultural and Religious Sensitivities	Moderate	High	Consult cultural advisors well as Islamic organizations to guarantee that the project aligns with cultural norms and ethical guidelines effectively communicate the methods used and prioritize respectful handling of culturally sensitive information.
R4	Technological Challenges with Arabic Text	Moderate	Moderate	Utilize known Arabic NLP tools such, as tokenization. Named entity recognition to handle Arabic text effectively. Test the text processing system rigorously to manage right to left text and linguistic subtleties accurately.
R5	System Scalability Issues	Moderate	High	Test the scalability of the graph database system by examining its ability to handle datasets as they expand in size using a cloud-based setup that enhances performance and efficiency for graph queries, with optimized indexing and algorithms.
R6	Unforeseen Technological Changes	Ongoing	Moderate	Keeping myself up, to date with the advancements in graph databases and the processing of text content and be open minded, to integrating new technologies or enhancements that could enhance the progress and efficiency of your project.



## A APPENDIX: A PLES Regulation Alignment

This section of appendix outlines the key Professional, Legal, Ethical, and Social (PLES) compliance factors related to this dissertation.

### B.1 Professional Issues

This Dissertation follows guidelines by utilizing openly accessible sets of Arabic heritage materials, like hadith tailored for academic purposes. References to sources such as datasets and scholarly works are meticulously included to maintain transparency and academic credibility in the study. The primary objective of this project is to employ graph databases for the storage and analysis of collections of Arabic heritage materials to uncover their role, in safeguarding and understanding texts. In this research project scenario as described here no commercial interests or client requirements are influencing the inquiry process which enables an effort, on theoretical investigation and technological advancements.

### B.2 Legal Issues

The project is carried out following all guidelines to guarantee that all data utilized. Hadith. Is sourced from publicly available datasets or properly licensed for academic research purposes to prevent any legal issues related to proprietary or restricted data sources. The creation of the graph database model and the essential software tools, for this study will be made available as source to comply with intellectual property regulations and encourage adoption, in academic circles. The academic research abides, by both international regulations, on data protection and intellectual property rights as it lacks any motives.

### B.3 Ethical Issues

Ethical aspects play a role in this study because of the religious importance attached to the hadith and other Arabic heritage texts being analysed research project is dedicated to treating these texts with respect and cultural awareness to prevent any misrepresentation or distortion of their significance when interpreting them using graph databases. As the datasets do not include information or sensitive details, about individuals who're alive today the issue of privacy is not a major concern. Furthermore this study does not include testing on subjects. All information is solely utilized for academic research. This ensures that the project does not impact individuals or

communities in terms of harm or consequences. The primary objective continues to be advancing knowledge and safeguarding heritage as opposed to focusing on commercial use or decision making purposes.

## **B.4 Social Issues**

In order this project seeks to help preserve and study heritage by using graph databases to analyze texts such, as hadiths in order to make them more accessible and comprehensible to a wider audience. By sharing the code and research methods created during this study with the public domain community and promoting open source initiatives will encourage exploration and creativity. The discoveries, from this project will be disseminated within communities to promote cooperation and enhance knowledge in the areas of humanities and cultural conservation. The project aims to make sure that these important texts can be easily accessed for research and also work towards preserving knowledge, for the benefit of our future generations.