

---

# AN INVESTIGATION INTO THE SUITABILITY OF GRAPH DATABASE TECHNOLOGY IN THE ANALYSIS OF SPATIO-TEMPORAL DATA

---

**S.D. Baker Effendi      A.B. van der Merwe**

Department of Computer Science  
Stellenbosch University  
Stellenbosch, South Africa  
{20056354@, abvdm@cs}.sun.ac.za

November 11, 2019

## ABSTRACT

Large quantities of spatio-temporal data are captured everyday, whether by large web-based companies for social data or by other industries such as those concerning disaster relief or marine data analysis. This increases the need for backend systems to provide realtime query response times while scaling well (in terms of storage and performance) with increasing quantities of structured or semi-structured, multi-dimensional data. Traditionally, relational database solutions have been used and, with technologies such as PostGIS, this solution has been well adapted in this regard. However, the use of graph database technology has been rising in popularity and development and has been found to handle graph-like data much more effectively.

This work is motivated by the need to effectively store multi-dimensional, interconnected data and to investigate whether or not graph database technology is better suited to address this when compared to the traditional relational database approach. Three database technologies will be investigated using this dataset namely: PostgreSQL, JanusGraph, and TigerGraph. The dataset used is the Yelp challenge dataset and the evaluation is based on how each database performs under data analysis scenarios similar to those found on an enterprise level.

**Keywords** Graph database · spatio-temporal data · NoSQL · Yelp dataset

## 1 Introduction

The use of graph database technology as well as the volumes of spatio-temporal data have increased in recent years [1]. Graphs are composed of vertices (nodes) and edges (relations) joining two vertices together. In this paper, nodes are referred to as vertices and relations as edges. In a graph database, a vertex could be a user, business, vehicle, or city and edges could be a friendship, review, road, or location affiliation. Due to the relational nature of graphs, they can be used to model a variety of real-world phenomena such as physical, biological, social and information systems [2].

Graph databases form part of a group of database technologies known as Not Only SQL (NoSQL) databases. This, and the aforementioned characteristics of graphs, make graph databases a flexible and scalable solution for many enterprise level problems. It is important to note that large streams of social media data are logged daily, a significant amount of which contain geotags and timestamps [3].

Although many categories of data can implicitly be interpreted in a graph structure, many traditional relational databases lack the architecture and high-level query languages to effectively model and manipulate this structure. Graph databases are designed to represent data as attributes on vertices and edges. They are often schema-less (no fixed data structure) and their attributes are described as key-value pairs. A major strength of graph databases are that they excel in complex join-style queries where relational databases often handle these inefficiently for large tables [4]. In contrast, graph databases perform poorly in full graph aggregate queries (but fairly well when working in a local subgraph) whereas relational databases are efficient in this regard.

This research aims to clarify the suitability of which database technology is best suited to handle large-scale, spatio-temporal data. In particular, the investigation will cover the query speed, expressiveness of the available query languages, and computational demand on the open source graph database, JanusGraph, the open source object-relational database system, PostgreSQL, and the enterprise level graph analytics platform, TigerGraph.

The Yelp dataset will be used for our investigation. There is an annual challenge where students have a chance to conduct research using this dataset. The dataset is relevant to this research as it holds, among others, spatiotemporal properties. Hundreds of academic papers have been written using this dataset and winners of the challenge receive a prize based on their technical depth and rigor and other criteria as described in [5].

Our research aims will be achieved by; (i) modelling and importing the Yelp dataset into the three previously mentioned databases, (ii) benchmarking them against one another using various analysis that represent real-world contexts to fairly represent the type of querying and data manipulation the dataset would typically undergo, and (iii) comparing the suitability of JanusGraph with that of Tiger-

Graph, by investigating their graph query language implementations.

## 2 Problem Definition

The research in this paper compares the problem of modelling, visualizing, and querying large-scale spatio-temporal data using traditional relational database approaches versus a more modern graph database one. Query performance, storage and computational cost, and ease and efficiency of simulation of real-world applications will be investigated. The Yelp dataset not only contains spatio-temporal properties, but also social and linguistic aspects.

One of the core issues is how to store spatio-temporal data efficiently. Coordinates are comprised of latitude and longitude and time adds a third dimension. The accessing times of secondary storage is an issue when housing large volumes of data. The use of B-Trees, R-Trees, and a *geograph*<sup>1</sup> are three techniques investigated for indexing uni- and multidimensional data efficiently.

Effectively querying a graph topology is another problem which one needs to deal with. Basic graph pattern matching is a popular technique used to extract data from graph databases, but there are a range of graph query languages which implement this and other, more complex, graph patterns. Section 4 addresses three graph querying languages namely; Gremlin, Cypher, and GSQL. The solution to how traditional relational operators, such as union and difference, are implemented in graph pattern matching are also covered under the aforementioned section. The conciseness of querying graphs versus using traditional SQL is addressed in the conclusion (Section 8). To effectively measure the difficulty of learning and writing these queries, one requires experiments involving developers new to these languages and measuring their progress, but this is beyond the scope of our investigation.

Due to the rise of popularity in web-based applications, the ease of incorporating these three databases will be investigated with a Flask [6] back-end and Angular [7] driven front-end. This will show how appropriate it is in real-world applications and production settings to use graph database technologies and any challenges encountered during implementation versus using a classical SQL database back-end.

## 3 Literature Review

The following literature review aims to present and clearly define the terms and definitions related to database technologies. This section begins by defining the following principles; ACID compliance and the CAP theorem. The relational model will be explained to show how classic relational database technology came to be. The shortfalls

<sup>1</sup>A grid-based geospatial system where vertices are grids and edges connect other vertices to these grid vertices to indicate their physical location.

of relational databases will then be used to introduce what NoSQL technology is, and how this contrasts with the many variations of NoSQL databases including key-value, column family, and document stores. Based on our literature review, we will summarize the contributions of our research to spatio-temporal data mining, analysis, and visualization.

### 3.1 ACID Compliance

From an article written in *Database Guide* [8]; Atomicity, Consistency, Isolation, Durability (ACID) are the properties that can guarantee that transaction based databases perform reliably. These properties concern themselves with how databases recover after failures such as a transaction, system, or media failure.

**Atomicity** A single transaction could be comprised of one or more steps, and if one of them fails then the transaction will only be partially successful. Atomicity guarantees that if one step fails, the whole transaction fails and the state of the database will be rolled back to a state prior to the transaction.

**Consistency** Data may need to conform to a variety of rules and constraints, especially with regard to relational databases. If any data does not conform to these rules or the given schema, it will not be consistent with the rest of the database. Allowing this data to become a part of the database implies that we cannot guarantee our predefined rules. This is one of the properties that NoSQL databases relax on in order to provide higher availability and partition tolerance. A consistent database will ensure that any data being inserted into the database will follow all the predefined rules or else the transaction will be rejected.

**Isolation** To avoid transaction conflicts, each transaction needs to be performed in isolation. This means that no transaction will affect another, and if necessary, a given transaction will take precedence over another in case of a conflict.

**Durability** For a database to be reliable in the event of any of the aforementioned failures, it should be the case that if a transaction has been committed, it needs to be able to guarantee that the transaction has been stored permanently. Durability along with atomicity work together to keep a database reliable when faced with a failure either during or directly after a transaction has been committed.

A database that is ACID-compliant is robust against data being corrupted during a failure and guarantees that only successful transactions are processed.

### 3.2 CAP Theorem

Gilbert & Lynch [9] explain that, in distributed systems, there is often a trade-off between consistency, availability, and network partition tolerance (CAP). Eric Brewer presented this theorem in the context of geographically separated datacenters supporting web services, implicating that a multi-node database cannot ensure all three properties under CAP. Following this web service context, the three properties are described as follows:

- **Consistency** means that the server will return the correct response to each request.
- **Availability** will guarantee that every request will receive a response.
- **Partition tolerance** refers to the system implementation allowing the database to be split into multiple groups unable to communicate with one another.

Gilbert & Lynch, when referring to these properties, say:

“CAP states that any protocol implementing an atomic read/write register cannot guarantee both safety and liveness in a system prone to partitions.”

The CAP theorem illustrates the trade-off between safety and liveness when considering an unreliable system. For the safety property to hold, every point in each step of a transaction should be atomically consistent. For the liveness property to hold, a desirable outcome should result as long as the execution continues for long enough.

### 3.3 Relational Databases

Edgar F. Codd, an IBM fellow, is normally accredited for proposing the relational model [10]. The relational model provides a solution to the problem of making a distinction between the logical view and physical representation of the data being stored. This was the main shortcoming with the database management systems of the late sixties. Other issues included programmers having to follow large pointer chains to access data and having to redesign programs whenever storage layout structure had to be reconfigured (such as the introduction of indexing).

The relational model solved these issues by providing a clear boundary between the logical and physical representation of data, created a simple model that could be understood by all programmers, and introduced high level language concepts.

The model's structure stores data in tuples with relations represented as tables. Each tuple has attributes, primary and candidate keys. Each attribute has a domain which is the data type of the attribute's values. The tuples became rows and attributes the columns for each row under their respective tables. Furthermore, the degree refers to the number of columns in a table whereas cardinality the

number of rows in a table. An illustration of this can be seen in figures 1 and 2.

The high level language concepts introduced are the algebraic operators we see in SQL today such as SELECT, UNION, JOIN, etc. These operators allow users to convert relations (tables) into other relations (resulting in tables as output).

Since the 1970's, relational database management systems have been the leading database technology with a variety of capabilities and language dialects. This was until NoSQL obtained traction among large Internet corporations and was implemented in the form of distributed, non-relational databases in the 2000s [4].

### 3.4 NoSQL Databases

NoSQL databases were created to address the shortfalls in classic relational databases when used for data of large volume, variety, and velocity [11]. Many of these issues were introduced with the increasing popularity of the Internet and systems hosted in the cloud. Relational databases scale poorly over multiple nodes and on large volumes of data, which is why SQL databases could no longer be used as the data solution for major Internet companies such as Facebook, Amazon, and Google.

NoSQL database technologies are non-relational and use non-SQL languages to manipulate and query the data. NoSQL is designed to run on multiple nodes (distributed systems), scale horizontally, and some even implemented technologies such as massively parallel processing (MPP) [12] or in-memory processing (such as H-store [13]). NewSQL is a distributed SQL technology which implements NoSQL features to achieve ACID-compliance (alongside partition tolerance). They will not be referred to as relational database technology when mentioned in this paper [11].

A main challenge for non-relational database systems is the conflict between high availability in distributed systems and remaining transaction based and ACID-complaint (see Subsection 3.1). This results in NoSQL DBMS to typically choose availability and partition tolerance over consistency (see Subsection 3.2) and have become BASE (Basically Available, Soft-state, Eventually consistent)-type systems to compensate [14]. Not all NoSQL databases follow this trend and there are some vendors such as Neo4j and OrientDB that are fully ACID compliant [8].

NoSQL databases fall under a broad scope and, in this paper, they will be classified under four categories; key-value stores, document stores, column family (or wide-column) stores, and graph databases.

**Key-Value Store** The basis of a key-value store is that the data is organized in the form of key-value pairs. The key is typically a numeric or string value whereas the value can be any data object or collection of data objects (where a list would be represented by multiple entries of the same

key). Each document has a unique identifier with a list of attributes forming the key-value paired data. A simple example of this can be seen in Figure 3.

Castellano [15] explains that key-value stores are designed to be lightning fast, simple, and unstructured. They expose three operations namely; PUT, GET, and DELETE, with some implementations adding an additional operation, SEARCH, that matches keys or key-value pairs given a specified search expression and key namespace. These databases are decentralized in nature so struggles to provide the transactional guarantees of ACID.

**Document Store** As the name suggests, document store databases store their data in the form of documents. Typically, the document formats are JSON, XML, PDF, etc. Unlike key-value stores, document stores are semi-structured and both keys and values are fully searchable [11]. Document stores are still schema-less and are well suited for data that is dissimilar such as those which do not fit well in a table with set columns [16] or require many "nulls". These databases, like key-value stores, do not perform well if the data is highly relational and requires some kind of normalization.

**Column Family Store** Manoj [16] explains column families as more structured data stores in that they store data in columns. They are a hybrid row/column store but store their data in distributed architectures instead of tables. Data is stored in a column-family (analogous to a table in SQL), but column-families have no relation to one another unlike in relational models. Column-families are less flexible than the previous key-value and document store databases as one will have to predefine a column family's attributes and are grouped together under keyspaces. A column-family model is illustrated in Figure 4.

**Graph** Graph databases, the focus of this paper, applies a graph structure to store and manipulate data. Data can be stored as key-pair attributes on both the vertices or edges of the database (as can be seen in Figure 5). Edges can be unidirectional or bidirectional which may add more information about the kind of relation between two edges. Graph databases are the only NoSQL databases in the four classifications of this section that concern themselves with relations and can be visualized easily in a more human-friendly manner [11].

Graph databases are strong at finding patterns and revealing information about the relationship between vertices rather than aggregate queries on the data itself.

Within the category of graph databases we find three main subcategories; property (or attributed) graphs, hypergraphs, and resource description framework (RDF) triples [2]. The type of graph databases focused on in this investigation are property graphs.

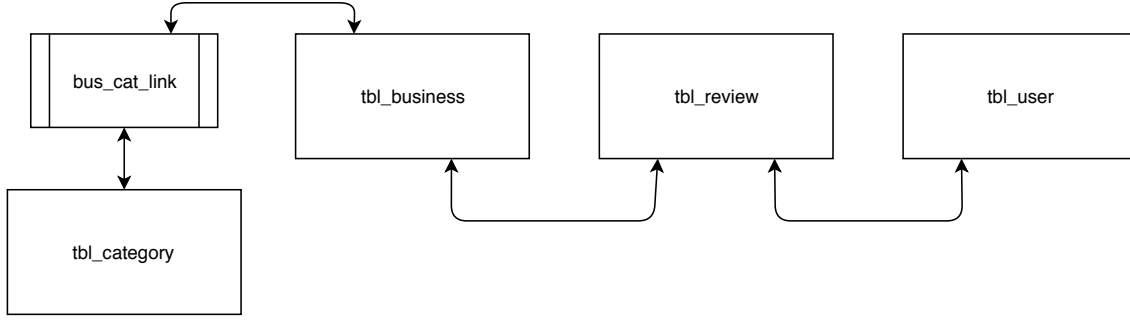


Figure 1: A simplified representation of the Yelp dataset in a relational database context illustrating tables and their relations.

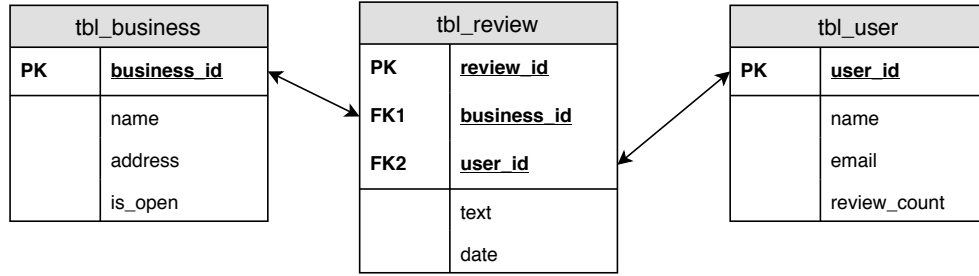


Figure 2: A closeup on the foreign key and primary key structure within the business, user, and review tables.

Businesses	
Key	Attributes
0	Name: "McDonald's" Categories: "Fast-food" Categories: "Takeaway" Stars: 2.5 Open: true
1	Name: "KFC" Categories: "Fast-food" Categories: "Restaurant" Stars: 3.0 Open: true

Figure 3: An example of business entries stored in a key-value store.

Yelp Keyspace	
Users Column-Family	Businesses Column-Family
RowID: 0 FirstName: "David" Email: "david@sun.ac.za" ReviewCount: 8	RowID: 0 Name: "McDonald's" Categories: "Fast-food" Categories: "Takeaway" Stars: 2.5 Open: true
RowID: 2 FirstName: "Kyle" Email: "kyle@gmail.com" ReviewCount: 3	RowID: 4 Name: "KFC" Categories: "Fast-food" Categories: "Restaurant" Stars: 3.0 Open: true

Figure 4: An example of user and business column-families.

### 3.5 Summary

From the literature above, we face the following questions:

**What would make a non-relational solution more suitable than a relational solution?**

Moniruzzaman & Hossain [11] elaborate the shortcomings of relational database technologies in terms of performance when distributed over geographically diverse datacenters due to their strong emphasis on consistency while maintaining ACID-compliance. As emphasized by Chen [2] & Makris et. al. [1], due to the trend of services hosted in the cloud and velocity of incoming data and requests,

a horizontally scalable, distributed system would be best suited to address these issues – which would suggest a NoSQL solution.

**When compared among other NoSQL databases, what makes a graph database solution stand out?**

Spatial data is often highly relational due to the ties between vertices representing physical locations and their association with other entities in the database. This means the data will be closely related and that trend will follow in how the queries are written. Chen [2] explains that

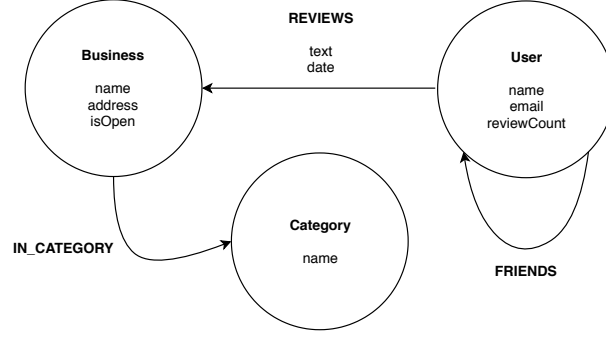


Figure 5: A simplified representation of the Yelp dataset in a property graph database context illustrating vertices, edges, and their attribute keys.

the queries in a graph database are designed to reveal hidden trends among the relationships in the data (especially when visualized) rather than within the data itself adding additional value in this implementation.

Time adds an ordinal relation among all entries in the data that hold this property. Graph databases are the most relational among the NoSQL databases and have the ability to ask complex questions on complex data. These complex questions would be multi join-style queries which are typically expensive, but due to a graph’s relational architecture, graph databases should be better suited in comparison to other NoSQL database solutions handling semi-structured data [4].

#### Why compare a relational solution to a graph database and not another NoSQL solution?

Relational databases should come close to what is required in our setting due to their ability to manipulate data based on table relations, but their lack of scalability may impact overall benchmark performance negatively when tested on large volumes of the kind of data we consider, especially due to how join-style queries scale with table size [4].

Makris, et. al. [1] compares how MongoDB (with GeoJSON) and PostgreSQL (with PostGIS) perform against one another on spatial data and show that PostgreSQL outperforms its NoSQL rival. Considering our discussion in subsection 3.4, it does not come as a surprise, given the lack of suitability of relational queries in a document store databases. In general, complex queries perform poorly in databases such as MongoDB. One important observation that this paper does point out is how well PostgreSQL performs, despite the volume of data. This suggests that the scalability issue of relational databases was not as prominent in this comparison as much as how well these complex queries were handled.

#### Why spatio-temporal data?

Spatio-temporal data is constantly being generated by GPS-equipped devices everyday [3][17]. This data has important applications in epidemiologic [18], marine, disaster

relief, and social data investigation [19]. Few papers address a graph database solution for spatio-temporal data, but due to its importance and abundance it would be important to seek an appropriate technology to store, model, and visualize these datasets.

**In conclusion** The findings in this paper would answer which technology, graph or relational, would perform the best not only in benchmark performance but also in suitability in terms of querying, modelling, and visualizing of spatio-temporal data. This answer will be useful if one would like to perform analysis on not only a spatio-temporal dataset, but also investigate relationships in the data with a suitable database, and perform/write efficient queries.

## 4 Graph Querying Techniques and Languages

The following section explores the three popular graph querying languages – two of which are used in the databases being investigated in this report – and properties and techniques of graph querying. Graph query languages have a notion of traversing the graph and accumulating results from pattern matching. This makes it very different to traditional SQL and is the main factor in contributing to the learning curve – how to conceptualize graph expressions. These languages are still fairly new with some of the latest additions being GraphQL by Facebook (released in 2015) [20] and GSQL by TigerGraph (released in 2017) [21]. As techniques of graph traversal becomes more refined, one can expect a few more query languages to be developed.

### 4.1 Languages

**Gremlin** Gremlin is a functional, data-flow query language running in its own virtual machine [22]. Due to how Gremlin is compiled, it can be written in both an imperative and declarative manner and also be embedded in a host-language [23]. Imperatively querying with Gremlin involves explicitly stating the traversal pattern, whereas

declarative allows the traverser to decide. The benefits of using an embedded query language is that security risks from string concatenation and sanitizing input is handled automatically.

Gremlin integrates into multiple vendor’s products (in our case, JanusGraph) and benefits from Turing-completeness. Each transaction is local to a given thread which means that a transaction is automatically created in that thread without having to explicitly call a create method. Transactions are atomic and support rollbacks.

**Cypher** In an effort to create an easier graph language to learn, Cypher (see the openCypher project [24]) was created as a very SQL-like query language. Unfortunately, while it is easier to pick up, it is not Turing-complete [25].

It should be noted that whether the fact that native Cypher is not Turing-complete should be considered as a drawback, is a problem dependent consideration. For most ordinary querying, not having a Turing-complete query language poses few to no major drawbacks at all. Examples of algorithms, commonly used on an enterprise level, which cannot be implemented in native Cypher, are given by PageRank/Label Propagation style algorithms [25]. The workaround is that Cypher provides calls and algorithms libraries to conduct these. Many important algorithms affected by these limitation are implemented in the Neo4j Graph Algorithms library<sup>2</sup>. An example of using this library is given in Figure 7. A problem with a hard-coded PageRank-style algorithm is given by the limitation of not having full control over the number of iterations and terminating conditions. Ultimately, one should consult the documentation before deciding if a Turing-complete graph query language is required.

Use of ASCII-art helps to create more intuitive and easy to visualize queries, e.g. edges are denoted by `-->` or `--[...]->` and vertices by the use of parenthesis as in `(b:Business)`.

Cypher is a declarative language, so there is less control over how the traverser goes about each step. This means that Cypher has less flexibility and could perform worse than Gremlin. Another criticism in terms of performance is that Cypher compiles into Gremlin which is then executed by the TinkerPop engine [27] and this overhead should still be reduced. At the time of writing, Neo4j has been gathering more attention around Cypher, so the performance issues may improve significantly in the near future.

Cypher has the advantage of being able to express complex traversals in a simple and more intuitive manner than Gremlin. Due to the open-source nature of Cypher and, how closely it is linked to Gremlin, it benefits from the same portability advantages by being able to integrate with multiple vendors. TinkerPop 3 supports Cypher<sup>3</sup> so any

TinkerPop 3 enabled graph database can be queried with either language and thus benefits from the ability that either query language can be used depending on if a high-level traversal or simple query is required. An example of a useful query on the dataset we consider are given in Figure 8.

**GSQL** GSQL is described as a SQL-like language which is the conceptual descendent from technologies such as Gremlin, Hadoop MapReduce, SQL, Cypher and SPARQL [28]. GSQL, like Gremlin, is a Turing-complete graph query language that accumulates data along a traversal. One limitation that Gremlin has compared to GSQL, is that Gremlin cannot simultaneously group two tables by separate group-by attributes<sup>4</sup>. GSQL achieves this by providing the ability to define multiple grouping accumulators and can use these accumulators to accumulate data based on varying criteria within the same steps.

Accumulators are variables which accumulate information over a graph traversal and come in two major groups namely, scalar and collection. Scalar accumulators such as `OrAccum` or `SumAccum` store a single value. Collection accumulators such as `ListAccum` or `SetAccum` store a set of values or, as is the case with `ListAccum`, can nest accumulators. An example of a GSQL query using accumulators can be seen in Figure 9.

In GSQL, there is a large emphasis on creating a language that enables massively parallel processing on queries. The vertex and edge blocks in GSQL queries indicate independent computations separated by incoming vertices or edges referred to as guarding conditions. These blocks are pieced together by the output of one block being the input to another. Control flow can be handled by if-then-else or while statements, allowing for subsequent blocks using dynamically calculated input.

As with Gremlin, there is an emphasis on a strong, functional programming style. One has the ability to define named, parameterized queries which is analogous to creating a function with arguments. These parameterized queries can then be called by other queries enabling the re-use of code. As is the case with Gremlin and Cypher, TigerGraph allows for the conversion of Cypher to GSQL for those migrating from their competitor, Neo4j [29].

## 4.2 Graph Pattern Matching

Graph pattern matching is an example of *declarative* (descriptive) querying. Basic graph patterns follow the structure of the graph to query. A basic graph pattern for a property graph<sup>5</sup> is a graph where variables appear on the edges and vertices. A *match* for a basic graph pattern is

<sup>2</sup><https://neo4j.com/docs/graph-algorithms/3.5/labs-algorithms/>

<sup>3</sup><https://github.com/opencypher/cypher-for-gremlin>

<sup>4</sup>Note that Gremlin is able to group two tables by separate group-by attributes, but it would need to do this while using the store step between groupings due to the dataflow architecture of the language.

<sup>5</sup>Which is the type of graph being investigated in this paper.

```

g.V().has("User", "user_id", "qUL3CdRRF1vedNvaq06rIA")
  .outE("REVIEWS").has("stars", gt(3)).order().by("date", desc)
  .as("stars", "text")
  .inV().has("location", geoWithin(Geoshape.circle(35.15, -80.79, 5)))
  .as("business_id")
.select("stars", "text", "business_id")
.by("stars").by("text").by("business_id")

```

Figure 6: Returns all the reviews by the specified user ID within a 5km radius of 35°15N 80°79W ordered by date descending. The ability to query and index spatio-temporal properties is provided by JanusGraph’s integration with a search engine such as Elasticsearch.

```

CALL algo.pageRank("Page", "LINKS", {
  iterations:20,
  dampingFactor:0.85,
  write: true,
  writeProperty:"pagerank"
})

```

Figure 7: An example from the Neo4j documentation [26] on using the PageRank algorithm from the Neo4j Graph Algorithms library.

```

CALL apoc.periodic.iterate(
  "MATCH (p1:User) RETURN p1",
  "MATCH (p1)-[:REVIEWS]->(r1)
  -->()<--
  (r2)<-[:REVIEWS]-(p2)
  WHERE id(p1) < id(p2)", {
    batchSize:100,
    parallel:true,
    iterateList:true
  }
);

```

Figure 8: Returns all users who have reviewed the same businesses as a given user. The query iterates through all users. Note how the return match verifies that the user p1 does not match user p2.

then mapped against the graph being queried. The variables in the basic graph pattern subgraph is matched to selected values or constants in the original graph and returned as a result [30]. An example of the pattern produced by Figure 6 can be seen in Figure 10.

Complex graph patterns extend on basic graph patterns by including the traditional relational operators used for sets such as union, difference, optional, and filter. These operators are described next.

**Projection** A projection returns a subset of data from the accumulated results of the pattern match. An example of this is to return only the stars from reviews between a user and business and exclude the text and edge IDs.

**Join** The join of two basic graph patterns corresponds to the function of a natural join in classic relational query languages such as SQL. Since the output of a basic graph pattern is the result of the variables specified on the graph pattern, the output of a join between two basic graph patterns is the union of their output variables.

**Union and difference** The union of two basic graph patterns is satisfied when one pattern or the other satisfies the pattern match. The difference of two basic graph patterns where the set of matches in the one are not in the set of matches in the other.

**Optional** Optional works much the same as *join*, but instead of discarding the results from the evaluation which cannot be joined, the results from both matches are kept. This allows data with incomplete or unavailable properties to remain in the output.

**Filter** The filter operator restricts the matches over which the traversal is performed. In practice, these filtering criteria vary in complexity with the ability to search over regular expressions (when querying string data), between dates (when querying temporal data), and over a radius (when querying spatial data).

### 4.3 Navigational Queries

Angles, et. al. [30] describes navigational queries as queries where the length of the traversal is potentially arbitrary such as *path queries*. Path queries are the most basic navigational queries, where one is only interested in the results accumulated when traversing from a source to a destination. Path queries are useful when looking at friend-of-a-friend relations between users in social networks and find applications in route-finding [31].

An example of one such query can be seen in Figure 11 and 12.

Navigational queries that try to match no-repeated-node or edge paths problems are typically NP-complete. Due to this, it is often necessary to add additional limitations on the pattern to be matched or use imperative querying



```

CREATE QUERY getNearbyBusinesses(DOUBLE lat, DOUBLE lon, DOUBLE distKm) FOR GRAPH YelpGraph {
  SetAccum<STRING> @@vSet;
  @@vSet += getNearbyGridId(distKm, lat, lon);
  Grids = to_vertex_set(@@vSet, "Geo_Grid");

  bus =
    SELECT b
    FROM Grids:s-(Business_Geo:e)-Business:b
    WHERE geoDistance(lat, lon, e.LATITUDE, e.LONGITUDE) <= distKm;
  PRINT bus;
}

```

Figure 9: A GSQL query that returns business vertices nearby the specified coordinates which lie within the radius distKm.

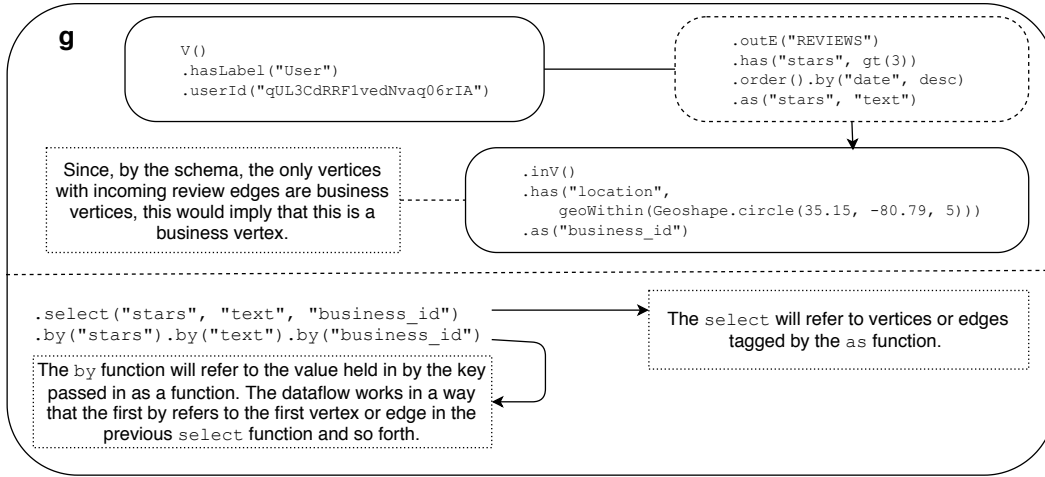


Figure 10: An illustration of the graph pattern produced by Figure 6. The top half shows the graph traversal pattern and the bottom half how the output is projected from the matched vertices and edges.

```

MATCH (u:User)-[:friends*]->f
WHERE f <> u
RETURN f

```

Figure 11: Simple friend-of-a-friend query written in Cypher.

```

g.V().hasLabel("User").as("u")
.out("FRIENDS")
.hasLabel("User").except("u")

```

Figure 12: Simple friend-of-a-friend query written in Gremlin.

techniques. Another common path traversal includes shortest paths from one vertex to another or path existence queries. It is clear that path traversals may be complex and so it is important to bound queries using constructs such as repeat...times(x), in Gremlin, to limit the search space.

## 5 Design and Architecture

The design of the software, architecture, and technologies used to create the web-application and databases are described in the following section. The analysis being performed by the web-application, how it performs each analysis, and how the backend queries each database is discussed.

### 5.1 Database Technologies

This section covers the underlying technologies used to implement the databases being benchmarked. This is important to note in terms of how portable each technology is, how difficult the configuration is before each database can be used for a given project, or limitations on performance or hardware requirements due to software requirements. While each database technology has varying capabilities, in terms of being supported for a given operating system, they each provide support for being deployed in a containerized environment.

### 5.1.1 PostgreSQL

PostgreSQL is written in the C programming language, is an object-oriented relational database, and is queried via SQL commands. PostgreSQL, like many relational databases, is ACID-compliant and robust to transactional failures. It is built to be extensible with a variety of extensions one can install for additional features such as UUID or spatial indexing. PostgreSQL can deploy on a variety of major operating systems such as Windows, Mac, Linux (Redhat, Debian, and a few others), Solaris, and BSD.

PostgreSQL is a mature product and has a large amount of support from its open-source community. Postgres is flexible in that it supports a variety of data types and allows the definition of own types. Founder of NewsBlur<sup>6</sup>, Samuel Clay, mentions using Postgres for multiple years for storing millions of sites and subscriptions. Canonical<sup>7</sup> founder, Mark Shuttleworth, explains that, while using Postgres during the development of Launchpad, finding it “robust, fast, and professional in every regard” [32].

Many of these features and opinions of using PostgreSQL in production environments on this scale is why PostgreSQL was a reputable relational database to benchmark against.

### 5.1.2 JanusGraph

JanusGraph is a highly scalable graph database that is ready to be clustered between multiple machines. It is a transactional database which supports ACID-compliance and eventual consistency [33]. It is written in Java and thus platform independent. JanusGraph is a project under The Linux Foundation and is forked from the Titan project as a continuation of the vision in creating an open-source, scalable, highly concurrent graph database. There is support for those wishing to migrate from Titan in order to benefit from the bug fixes and additional features now supported via JanusGraph [34].

JanusGraph is largely based on the Apache tech-stack making use of technologies such as Apache TinkerPop<sup>8</sup>, Lucene, Cassandra, Hadoop, and more. This adds to complexity when configuring JanusGraph as, for each technology plugged in, there may be configuration necessary. Gremlin is the native language through which JanusGraph is queried but, as mentioned in Section 4.1, can be extended for Cypher queries. JanusGraph benefits from optional support for advanced search capabilities and having no-single point of failure [35].

JanusGraph can store graph data via three supporting backends; Apache Cassandra, Apache HBase, and Apache Berkeley. The CAP Theorem (Section 3.2) should be referred to when considering which of the three backends to use – this is illustrated in Figure 13.

<sup>6</sup><https://newsblur.com/>

<sup>7</sup><https://canonical.com/>

<sup>8</sup>Thus makes use of the property graph model.

Examples of companies who have deployed JanusGraph in production include Netflix, Redhat, Uber, and IBM [37]. The professional support, documentation, and fact that one is able to leverage all these advanced features for free is why JanusGraph is one of the graph databases used in this investigation. The configuration of JanusGraph used in this report is with Apache Cassandra as the storage backend and Elasticsearch as the search engine for spatial and temporal query support.

### 5.1.3 TigerGraph

TigerGraph is an enterprise level graph analytics platform developed in the C++ programming language. TigerGraph was developed with hindsight from projects such as Apache TinkerPop and Neo4j and provides features such as native parallel graph processing and fast offline batch loading [38] [39]. Unlike JanusGraph, TigerGraph was developed from scratch in order to effectively create the next generation of graph database technology. TigerGraph won Strata Data’s “Most Disruptive Startup” Award for its return in this decision [40].

Some of the use cases explicitly mentioned by TigerGraph<sup>9</sup> are in geospatial and time series analysis. This lends itself as a promising database technology for this investigation. TigerGraph is queried using their GSQL querying language (see Section 4.1) where queries are optimized via an installation process where a REST endpoint is also generated in the process. Like JanusGraph, TigerGraph can be deployed on multi-machine clusters but this is limited to the enterprise version of this product. TigerGraph uses Apache Zookeeper for cluster management and Apache Kafka for message queuing.

GraphStudio is a web interface which is packaged along with TigerGraph which provides an interface to write, install and visualize queries, design and export one’s graph schema, and monitor database performance. This makes use of an Nginx web server [29].

For all intents and purposes, the developer edition is more than capable to perform the investigation of this paper. There is an enterprise version that allows for additional features such as multi-machine clustering.

## 5.2 Web-application Simulation

The web-application, called Providentia<sup>10</sup>, is used to queue analysis in a pipeline on which each benchmark is to be run, server performance measured, and accumulated results be displayed. This is deployed on target hardware and will import a subset of the data determined by a configurable

<sup>9</sup><https://www.tigergraph.com/solutions/>

<sup>10</sup>The name of the web-application is a nod to JanusGraph and Titan’s theme of Roman mythology. Providentia is associated with provision and forethought [41]. This was thought to be fitting due to the nature of our experiments designed to find the best database for storing and modeling our particular data.

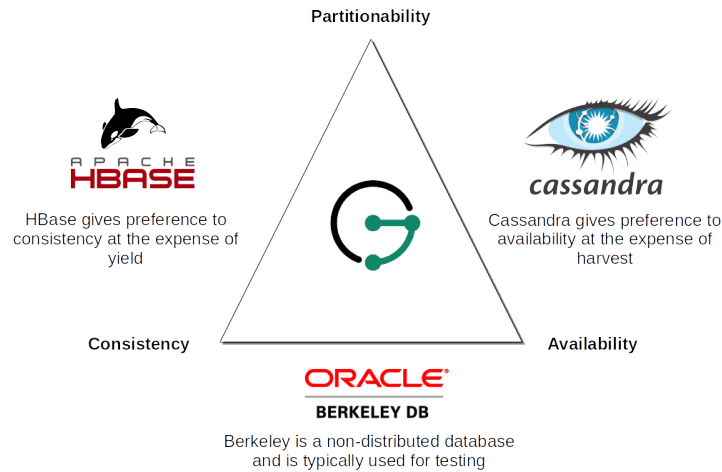


Figure 13: The CAP theorem illustrated using JanusGraph’s three supporting storage backends. This diagram is largely inspired from Chapter 1 in Titan’s documentation and is adjusted for JanusGraph [36].

percentage. Then one will be able to use the web-based interface to perform all necessary benchmarking tasks.

The architecture and how each technology communicates is illustrated in Figure 14. The databases are containerized using Docker<sup>11</sup>.

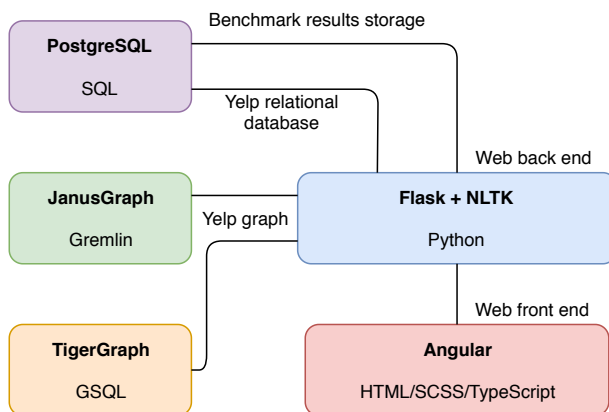


Figure 14: The architecture of Providentia.

JanusGraph has some Java specific features that add limitations when making use of embedded Gremlin in Python. The limitations are, when trying to make use of mixed indexing search predicates such as spatial queries, that one may only do this via Java or a superset language thereof. The workaround for this was to make use of the Gremlin Translator which takes a Gremlin query as a string and interprets it on the server side.

<sup>11</sup><https://www.docker.com/>

The first motivation towards using a Python backend is that the text in the reviews can be analysed using NLTK for easily implemented sentiment analysis. The second is that a simple REST API can be easily and quickly designed and deployed using the Flask framework. Angular was subjectively chosen as the front end framework as it allows for fast front end web development. All benchmarking results are stored in a separate database within PostgreSQL.

The front end of Providentia allows a user to query each database, test the sentiment classifier, add benchmarking jobs and to review the performance and results of each analysis. Each job is run serially to avoid too much interference and competition between each database for resources. At intervals the CPU performance and memory consumption is measured and stored in PostgreSQL. The server performance and results of an analysis can be viewed together to validate that the outputs are the same and how each database utilizes the server’s resources.

### 5.3 Data Analysis

Using each of the databases, a number of data analysis jobs are performed on the data. This section describes what each analysis aims to do and how they align to a typical real world use case. Each of these jobs have some kind of spatiotemporal aspect to test the accessibility of the data to demonstrate how well the given database handles the data.

These analysis are run over different percentages of the data loaded in each database and the performance is then measured and discussed in Section 7. This section goes into more detail about the types of queries written and how well each language expresses each query.

### 5.3.1 Sentiment Analysis

One application for graph analytic platforms is in machine learning. In order to further build context around each of the kernels mentioned in the next section, a simple binary sentiment classifier is used to classify the text of reviews as representing either a positive or negative sentiment. Sentiment analysis is a class of natural language processing where subjective information is extracted from a given text [42].

Although this investigation does not explore the applications of machine learning models on spatiotemporal data, it will explore how it could reveal interesting information alongside patterns among the data. All natural language processing is done using the Natural Language Toolkit (NLTK)<sup>12</sup> for Python.

#### 5.3.1.1 Machine Learning Model

NLTK’s Naïve Bayes classifier<sup>13</sup> was used as the model to train and classify sentiment. Naïve Bayes is a probabilistic machine learning model which has proven very effective in text classification. A limitation of this model, on a binary classification problem, includes only being able to perform linear separability. This should not be an issue regarding the use of adjectives as feature vectors, since the most informative adjectives seem to be verbose and particular to a specific sentiment e.g. words such as horrible, disgusting, perfect, and wonderful [43].

One problem faced by sentiment analysis is that of negation which may make certain adjectives less information e.g. good vs not good [44]. Nevertheless, the model performs well in terms of separating the two classes and, given the training and validation data (see Section 5.3.1.2), has shown to have good generalization performance. The model used in the experiments produce the results discussed in Section 7 and has a training and validation accuracy of 94.15% and 91.35% respectively.

#### 5.3.1.2 Training Process

The data used to train and validate the Naïve Bayes classifier was extracted from the Yelp dataset. The data is extracted with the following assumption: 1 star reviews hold negative sentiment and 5 star reviews positive sentiment. Training and validation data are split by a 70% training and 30% validation ratio.

The extracted data is then tagged as follows: 5000 reviews with a 1 star rating are tagged as “negative” and 5000 reviews with a 5 star rating as “positive”.

The data pre-processing and training of the Naïve Bayes model used in this investigation is based on the process described by Munir [45]. Using NLTK, one can separate

different parts of speech within a given sentence or paragraph. First, a bag of words is created from both classes of the training data. For a given review, punctuation is removed, words are tokenized, stop words removed, and only adjectives are kept and these adjectives are added to our bag. A frequency distribution is then used to determine the most informative words for a given label.

The top 5000 occurring adjectives are kept in the bag. These informative words can be seen as 5000 attributes where a review has the value “true” under an attribute found in the text and “false” when a review does not have that attribute contained within its text. This will be how a feature set will be constructed for a review. A Python dictionary is constructed with these 5000 attributes as keys and their value as the respective state as a boolean. This is then the first element of a tuple and the tag is the second. A trivial example of a feature set can be seen in Figure 15.

All of the reviews are next processed into these feature sets and tagged. They are then given to the model to learn. Once the model has been trained, a review is classified by breaking it up into a feature set, given to the model to classify and the predicted class label is then returned.

### 5.3.2 Kernels

Each of the following kernels represent some kind of user story for data analysis. Each one has some spatio-temporal constraint which will be applied on the respective queries. Each of these kernels will be benchmarked 15 times each such that the mean and standard deviation can be considered when comparing response times.

#### 5.3.2.1 Kate’s Restaurant Recommendation

**Description** This analysis selects a user near the beginning of the dataset named Kate. This user has a number of reviews for restaurants in the Las Vegas area. A subset of reviews which hold a strictly greater than 3 star rating by Kate are selected, sorted by date descending, and limited by 10 reviews per user in order to take the most relevant ones. These businesses are then selected, filtered by category “Restaurants”, and other users who have an equal or greater than star rating are then selected as the recommending users.

Now assume that Kate has relocated to a new area. All businesses which have been rated strictly larger than 3 stars by the recommending users are then selected as restaurants to recommend to Kate in the new area. The text in these reviews are checked for sentiment and the percentage positive sentiment is displayed alongside the average star rating for each recommended restaurant. This sentiment vs. average star rating is used as a metric to analyse in terms of asking the question: How reliable is the star rating versus the actual sentiment found in the text?

The purpose of the first part of this kernel is to test the performance of a 1-hop graph traversal pattern. This hop

<sup>12</sup><https://www.nltk.org/>

<sup>13</sup>`nltk.NaiveBayesClassifier`

Review Text					
I just ate the most yummy pizza at this restaurant! The setting was wonderful with outdoor seating on a patio with a beautiful view. Service was incredible!					
Feature Set					
Horrible	Dirty	Worst	Wonderful	Yummy	Perfection
false	false	false	true	true	false
Tag	Positive				

Figure 15: An example of a review tagged as positive with it's feature set.

is demonstrated in category filtering and finding users with mutual sentiment for a given review. This type of situation is faced by many recommendation technologies and this is quite a basic technique for recommendation. The additional challenge is the relocation of Kate and seeing how responsive the database is to the spatial and temporal aspect which is the second part of this kernel. The accumulated list of users is split into a separate query for each user to test the ability of each database technology to perform concurrent reads on subsets of data which is a strength of NoSQL databases.

**SQL** Listing 1 works as follows: First, the reviews, user, and business tables are joined together to find out which businesses Kate reviewed. Once all these businesses are obtained, all users who rated these businesses above 3 stars are obtained by another review join. Finally, the categories table (and its linking table) is joined to the businesses to only include restaurants.

```
SELECT DISTINCT OtherReviews.user_id
FROM users
JOIN review KateReviews
  ON users.id = KateReviews.user_id
  AND users.id = "qUL3CdRRF1vedNvaq06rIA"
  AND KateReviews.stars > 3
JOIN business KateBus
  ON KateReviews.business_id = KateBus.id
JOIN review OtherReviews
  ON OtherReviews.user_id != KateReviews.user_id
  AND OtherReviews.business_id = KateReviews.business_id
JOIN bus_2_cat Bus2Cat
  ON OtherReviews.business_id = Bus2Cat.business_id
JOIN category Categories
  ON Bus2Cat.category_id = Categories.id
  AND Categories.name = "Restaurants"
```

Listing 1: A SQL query that returns all the IDs of all the users who have rated restaurants Kate has been to above 3 stars.

Listing 2 begins by joining the review and business tables to match all businesses reviewed by the given user. These businesses are filtered by those within the Las Vegas area (with a radius of 50km). The reviews are ordered by date and limited to return only 10, as to keep only the latest reviews from the user.

```
SELECT review.stars, review.text, review.
  business_id
FROM review
JOIN business
  ON review.business_id = business.id
  AND review.user_id = "...
  AND ST_DWithin(location,
    ST_MakePoint(-80.79, 35.15)::geography,
    5000)
  AND review.stars > 3
ORDER BY review.date DESC
LIMIT 10
```

Listing 2: A SQL query that returns the star rating, text and business ID for restaurants a user has reviewed above 3 stars.

**Gremlin** The query in Listing 3 works by obtaining the vertex representing Kate and declaring the referencing for this vertex using the as step. The review edges leaving Kate are then filtered by those with a star rating above 3. The inV() step simply refers to the businesses who have the incoming review edge. The in("REVIEWS") step is different from the inE step as it skips across the edge and references the user vertices directly – which are then selected if they are not Kate's vertex. These users are referenced and the category vertices are checked to filter out businesses which aren't restaurants. The users referenced earlier are then selected, duplicates are removed and their ID values are selected for the return.

```
g.V().has("User",
  "user_id", "qUL3CdRRF1vedNvaq06rIA")
  .as("kate")
  .outE("REVIEWS").has("stars", gt(3))
  .inV().in("REVIEWS").where(neq("kate"))
  .as("users")
  .out("REVIEWS").out("IN_CATEGORY")
  .has("name", "Restaurants")
  .select("users").dedup()
  .values("user_id").fold()
```

Listing 3: A Gremlin query that returns all the IDs of all the users who have rated restaurants Kate has been to above 3 stars.

Listing 4 starts again at the user vertex. The reviews are once again filtered by rating and then ordered by date. The reviews are then referenced for later twice as this can be used to create a dictionary-style output later. The vertices

with these incoming review edges are implicitly business vertices again and are selected by their location using the mixed spatial index from ElasticSearch. These businesses are referenced for later. The reviews are selected and limited to 10 and all the prior references are selected. The dictionary-style output is produced using a combination of the select and by step<sup>14</sup>. The output will now look like a list of dictionary objects with the keys; stars, text, and business\_id which are easily serialized and deserialized over sockets or a REST API.

```
g.V().has("User", "user_id", "...")
  .outE("REVIEWS").has("stars", gt(3))
  .order().by("date", desc)
  .as("stars", "text")
  .inV().has("location",
    geoWithin(
      Geoshape.circle(35.15,-80.79, 5)
    )).as("business_id")
  .select("stars").limit(10)
  .select("stars", "text", "business_id")
  .by("stars").by("text").by("business_id")
```

Listing 4: A Gremlin query that returns the star rating, text and business ID for restaurants a user has reviewed above 3 stars.

**GSQL** GSQL queries work by first declaring variables and seeds, then writing the query. The GSQL query in Listing 5 declares a SetAccum which does not allow duplicates and the three seeds; categories, businesses, and PSet. First, businesses which are categorized as restaurants are selected, then businesses which were rated by the user – which will be Kate. The users who reviewed the intersection of these two business sets are then accumulated based on the constraints that their ratings were above 3 stars and that those users aren't Kate. The results of the accumulated users is then printed which returns a JSON string of the variable's contents when the REST endpoint is queried.

```
CREATE QUERY getSimilarUsersBasedOnRestaurants(
  VERTEX<User> p) FOR GRAPH MyGraph {
  SetAccum<STRING> @@userIds;
  categories = { Category.* };
  businesses = { Business.* };
  PSet = { p };

  Restaurants =
  SELECT b
  FROM businesses:b-(In_Category)->Category:c
  WHERE c.id == "Restaurants";

  PRatedBusinesses =
  SELECT b
  FROM PSet-(Reviews)->Business:b
  WHERE r.STARS > 3;

  PRatedRestaurants =
```

```
PRatedBusinesses INTERSECT Restaurants;

PeopleRatedSameBusinesses =
  SELECT tgt
  FROM PRatedRestaurants:m
  -(reverse_Reviews:r)->
  User:tgt
  WHERE tgt != p AND r.STARS > 3
  ACCUM @@userIds += tgt.id;

  PRINT @@userIds;
}
```

Listing 5: A GSQL query that returns all the IDs of all the users who have rated restaurants Kate has been to above 3 stars.

Listing 6 begins by declaring a custom type which has a tuple format. This allows one to record only the properties one wants to, similar to the select and by steps in Listing 3. A HeapAccum is constructed with a max size of 1000, and is used as it can order the accumulator by the properties in the custom type – in this case the date property. Using the geo-grid, all grid IDs are obtained using the getNearbyGridId method, which is then converted into a vertex set by matching the grid IDs to their respective vertices in the graph. The businesses connected to these grid vertices is obtained and businesses categorized as restaurants are then intersected as before Listing 5. The top 10 tuples from the heap are popped and accumulated in a list which is then printed.

```
CREATE QUERY getRecentGoodReviewsNearUser(
  Vertex<User> p) FOR GRAPH MyGraph {
  TYPEDEF tuple<DATETIME reviewDate,
    STRING businessId, INT stars,
    STRING text> restAndReview;

  DOUBLE lat = 35.15;
  DOUBLE lon = -80.79;
  INT distKm = 5;
  HeapAccum<restAndReview>
    (10, reviewDate DESC) @@busAndReviews;
  ListAccum<restAndReview> @@finalReviews;
  businesses = { Business.* };
  users = { User.* };
  PSet = { p };
  Grids = to_vertex_set(
    getNearbyGridId(distKm, lat, lon),
    "Geo_Grid");

  NearbyBusinesses =
  SELECT b
  FROM Grids:s-(Business_Geo:e)-Business:b
  WHERE geoDistance(lat, lon,
    e.LATITUDE, e.LONGITUDE) <= distKm;

  Restaurants =
  SELECT b
  FROM businesses:b-(In_Category)->Category:c
  WHERE c.id == "Restaurants";

  NearbyRestaurants =
  NearbyBusinesses INTERSECT Restaurants;
```

<sup>14</sup>select will select the prior references and the by step will select the attribute to display from the references in select respectively.



```

NearbyRestReviewsByP =
  SELECT u
  FROM NearbyRestaurants:b
  -(reverse_Reviews:tgt)->
  User:u
  WHERE tgt.STARS > 3 AND u == p
  ACCUM @@busAndReviews +=
    restAndReview(tgt.REVIEW_DATE, b.id,
                  tgt.STARS, tgt.TEXT);
  FOREACH i IN RANGE[0, 9] DO
    @@finalReviews += @@busAndReviews.pop();
  END;

  PRINT @@finalReviews;
}

```

Listing 6: A GSQL query that returns the star rating, text and business ID for restaurants a user has reviewed above 3 stars.

### 5.3.2.2 Review Trends in Phoenix 2018

**Description** This analysis goes deeper into observing the trend of various characteristics of reviews versus their star ratings. This is a common analysis performed on the Yelp dataset [46], but the version in this investigation selects a subset of reviews only within the 2018 year in the Phoenix area. The spatiotemporal boundaries placed on this subset may reveal hidden trends to be considered in future work.

Reviews are extracted first by location (which results in a much smaller subset than extracting by date first) then by date. The reviews are then separated by star rating. For each star rating, the characteristics of “funny”, “useful”, and “cool” are accumulated and the text is classified as either positive or negative. For each star rating these are normalized and placed next to one another to see the characteristic of a review from each star group.

**SQL** This kernel is the least complex of the three as it has a single join with a spatio-temporal constraint. Listing 7 returns selected characteristics on reviews where the date year is 2018 and reviewed businesses are within the Phoenix area.

```

SELECT text, review.stars, cool, funny, useful
FROM business
JOIN review ON business.id = review.business_id
AND ST_DWithin(
  location,
  ST_MakePoint(-112.56, 33.45)::geography,
  50000)
AND date_part("year", date) = 2018)

```

Listing 7: A SQL query that returns all the review text and ratings for businesses within 50km of the Phoenix area during 2018.

**Gremlin** One caveat of using mixed indexes on dates via the Gremlin Translator is highlighted in the query for this kernel. Usually, since Gremlin is designed to be embedded, one should make use of objects when appropriate. Since JanusGraph is being queried from Python<sup>15</sup>, with no support for mixed query specific parameters, date related parameters need to be parsed using static methods from the Instant class in Java. This can be seen in Listing 8 when filtering reviews by date. Alternatively, one could also use the filter step as is done in Listing 11.

Out of the set of businesses within the Phoenix area and set of all the reviews in 2018, the businesses set would be the smaller of the two. This is important when using a dataflow language since the whole subset will be accumulated before moving to successive functions in the query. Due to this characteristic of Gremlin, businesses are accumulated before the reviews.

```

g.V().has("Business",
  "location", geoWithin(
    Geoshape.circle(33.45, -112.56, 50)))
.inE("REVIEWS")
.has("date", between(
  Instant.parse("2018-01-01T00:00:00.00Z"),
  Instant.parse("2018-12-31T23:59:59.99Z")
)).valueMap()

```

Listing 8: A Gremlin query that returns all the review text and ratings for businesses within 50km of the Phoenix area during 2018.

**GSQL** Since only selected characteristics of a review are desired, a tuple is created at the beginning of Listing 9. The businesses within the Phoenix area are selected first, then reviews where the date part is 2018. The review tuples are accumulated into a ListAccum.

```

CREATE QUERY getReviewsFromPhoenix2018()
FOR GRAPH MyGraph {
  TYPEDEF tuple<STRING text, INT stars,
    INT cool, INT funny,
    INT useful> review;

  DOUBLE lat = 33.45;
  DOUBLE lon = -112.56;
  INT distKm = 50;
  ListAccum<review> @@reviewList;
  Grids = to_vertex_set(
    getNearbyGridId(distKm, lat, lon),
    "Geo_Grid");

  NearbyBusinesses =
    SELECT b
    FROM Grids:s-(Business_Geo:e)-Business:b
    WHERE geoDistance(lat, lon,
      e.LATITUDE, e.LONGITUDE) <= distKm;

  ReviewsForBusinesses =
    SELECT b

```

<sup>15</sup>Where ideally, it would be within a JVM language which has access to the JanusGraph specific classes and functions, e.g. Groovy.

```

FROM NearbyBusinesses:b
  - (reverse_Reviews:r)-
    User
WHERE YEAR(r.REVIEW_DATE) == 2018
ACCUM @@reviewList +=
  review(r.TEXT, r.STARS,
    r.COOL, r.FUNNY, r.USEFUL);

PRINT @@reviewList;
}

```

Listing 9: A GSQL query that returns all the review text and ratings for businesses within 50km of the Phoenix area during 2018.

### 5.3.2.3 Ranking Las Vegas by Friends' Sentiment

**Description** The purpose of this analysis is the ability to aggregate relations from depth 1 – 2 of a graph pattern while maintaining spatio-temporal constraints. The user story of this kernel is that a user in the dataset would like to travel to Las Vegas over the Nov – Dec period. Instead of asking from each of the hundreds of direct friends to thousands of mutual friends, the sentiment from their reviews written during the Nov – Dec period (irrespective of year) in the Las Vegas area will be analysed.

Both friends and mutual friends will be aggregated and all reviews written during the Nov – Dec period will be filtered. These reviews will be filtered by the spatial constraint of whether they are connected to businesses within 30km of the Las Vegas center. The remaining reviews will have their text data extracted and returned for analysis. Using the sentiment classifier, a percentage positive sentiment will be generated and this will be the result of the data analysis.

**SQL** Listing 10 begins by joining reviews and businesses before joining the two depths of friend relations. Julie's ID must be checked for at depth two as to only include mutual friends. Businesses are constrained by location and reviews by their month value under the date attribute. All reviews where the user ID matches the IDs of friends at depth 1 or 2 are selected and the text and star ratings of these reviews is returned.

```

SELECT DISTINCT R.text, R.stars FROM review R
JOIN business B ON R.business_Id = B.id
INNER JOIN friends F2 ON R.user_id = F2.
  friend_id
INNER JOIN friends F1 ON F2.user_id = F1.
  friend_id
WHERE F1.user_id = "7weuSPSSqYLUFga6IYP4pg"
AND F2.user_id <> "7weuSPSSqYLUFga6IYP4pg"
AND (R.user_id = F1.user_id
OR R.user_id = F1.friend_id)
AND ST_DWithin(
  B.location,
  ST_MakePoint(-115.14, 36.16)::geography,
  30000)
AND (date_part("month", R.date) >= 11

```

```

AND date_part("month", R.date) <= 12)

```

Listing 10: A SQL query that returns all the review text from reviews written by friends and mutual friends for businesses within 30km of the Las Vegas center.

**Gremlin** Starting with Julie's user vertex at the beginning of Listing 11, direct friends and mutual friends are accumulated as f1 and f2 respectively. First, duplicate users are removed, then any instance of Julie's vertex is removed from the union of these two accumulated groups of friends.

Once the traversal needs to be filtered by the temporal aspect, the same issue of having to call multiple functions due to the caveat in Listing 8 is encountered. The month value is extracted from each of these dates and filtered by their month values accordingly and referenced for later use. These reviews are then constrained by the locations of the business vertices by which they are connected. These reviews are then selected and their text and star rating data returned.

```

g.V().has("User",
  "user_id", "7weuSPSSqYLUFga6IYP4pg")
.as("julie")
.out("FRIENDS").as("f1")
.out("FRIENDS").as("f2")
.union(select("f1"), select("f2"))
.dedup().where(neq("julie"))
.outE("REVIEWS").filter{
  it.get().value("date")
    .atZone(ZoneId.of("-07:00"))
    .toLocalDate().getMonthValue() >= 11
  &&
  it.get().value("date")
    .atZone(ZoneId.of("-07:00"))
    .toLocalDate().getMonthValue() <= 12
}.as("text").as("stars")
.inV().has("location", geoWithin(
  Geoshape.circle(36.16, -115.14, 30)))
.select("text", "stars")
.by("text").by("stars")

```

Listing 11: A Gremlin query that returns all the review text from reviews written by friends and mutual friends for businesses within 30km of the Las Vegas center.

**GSQL** In a similar fashion to the Gremlin variant, the users from both depths of the friend relations are aggregated before removing Julie's user vertex. The SetAccum is used so duplicate user vertices should not be an issue. The intersection of the nearby businesses and businesses reviewed by the accumulated friend users are then used to extract the review data. Another SetAccum, @@reviews, is used to extract the text and star rating data from each review, where the reviews are constrained by their month values. The accumulated review data in @@reviews is then returned.

```

CREATE QUERY getFriendReviewsInArea(
  VERTEX<User> p, DOUBLE lat, DOUBLE lon)

```



```

FOR GRAPH MyGraph {
  TYPEDEF tuple<STRING text, INT stars> review
  ;
  SetAccum<review> @@reviews;
  SetAccum<VERTEX> @@F1F2;

  INT distKm = 30;
  users = { User.* };
  PSet = { p };
  Grids = to_vertex_set(
    getNearbyGridId(distKm, lat, lon),
    "Geo_Grid");

  NearbyBusinesses =
    SELECT b
    FROM Grids:s-(Business_Geo:e)-Business:b
    WHERE geoDistance(lat, lon,
      e.LATITUDE, e.LONGITUDE) <= distKm;

  F1 =
    SELECT f
    FROM PSet-(Friends)-User:f
    ACCUM @@F1F2 += f;

  F2 =
    SELECT f
    FROM F1-(Friends)-User:f
    ACCUM @@F1F2 += f;

  @@F1F2.remove(p);

  FReviewedBusinesses =
    SELECT b
    FROM users:f-(Reviews)-Business:b
    WHERE @@F1F2.contains(f);

  NearbyFBusiness =
    NearbyBusinesses
    INTERSECT
    FReviewedBusinesses;

  GetTheReviews =
    SELECT b
    FROM NearbyFBusiness:b
    -(reverse_Reviews:tgt)-
    User:u
    WHERE MONTH(tgt.REVIEW_DATE) >= 11
      AND MONTH(tgt.REVIEW_DATE) <= 12
      AND @@F1F2.contains(u)
    ACCUM @@reviews +=
      review(tgt.TEXT, tgt.STARS);

  PRINT @@reviews;
}

```

Listing 12: A GSQL query that returns all the review text from reviews written by friends and mutual friends for businesses within 30km of the Las Vegas center.

## 6 Implementation

The following section describes how the systems were setup, data processed, modelled and describes notable indexes used in each technology.

### 6.1 Benchmark Setup

This subsection describes the hardware used for benchmarking as well as the technical details of the dataset and preprocessing performed.

#### 6.1.1 Hardware Platform

All experiments were run on the following two machines listed in Table 1, that made it possible to consider how the technologies utilize multiple cores and perform with different storage limitations.

#### 6.1.2 Dataset

The dataset used is from the Yelp Dataset Challenge [5]. Due to the enormity of the dataset<sup>16</sup> only a subset of the data is used in increments to observe how well each database scales. These increments are percentage based and are applied per file e.g. 10% would mean 10% of `business.json`, `user.json`, and `review.json` applied separately. The dataset is stored as non-valid JSON and first had to be preprocessed and converted into valid JSON.

During preprocessing many attributes not used in the analysis or benchmarking were removed to save on import time and storage costs. Only businesses, users, and reviews were used from the dataset. This resulted in a  $\pm 11.39\%$  reduction in uncompressed storage size<sup>17</sup>, significant reduction in complexity and improvement in consistency among attributes. The result of this preprocessing can be seen in Table 2.

### 6.2 Schema Design

The following subsections describe the different indexing techniques used on the spatio-temporal attributes of the data and graphically present the schemas used in each database.

#### 6.2.1 Indexing

When storing a database considered to be large, it will necessitate that the data be stored on secondary storage, since it would most likely not fit in memory. This slows down the data access speeds considerably. When specific records need to be retrieved among large volumes of data, an intelligent method of organizing the data needs to be implemented. We can do this by narrowing our search

<sup>16</sup>Totaling in size around 8.69 gigabytes in uncompressed format.

<sup>17</sup>This results in a close to 1 gigabyte reduction, which is a significant improvement.

Table 1: The specifications of the two machines used to benchmark database performance.

Machine	CPUs	vCPUs	Base Clock	Memory	Storage Media	OS
Setup 1	6	12	3.4GHz	32GB	SSD	Debian 10
Setup 2	8	16	2.27GHz	32GB	SSD (across a network)	Ubuntu 18.04

Table 2: Data used from the Yelp dataset after preprocessing.

Business		User		Review	
Attribute	Data Type	Attribute	Data Type	Attribute	Data Type
business_id	string	user_id	string	review_id	string
name	string	name	string	user_id	string
address	string	review_count	integer	business_id	string
city	string	yelping_since	timestamp	stars	integer
state	string	friends	string array	date	timestamp
postal code	string	useful	integer	text	string
latitude	float	funny	integer	useful	integer
longitude	float	cool	integer	funny	integer
stars	float	fans	integer	cool	integer
review_count	integer	average_stars	float		
is_open	integer				
categories	string array				

space to a small subset where our target data lies – these small subsets are our *indexes*.

There are two broad classes [47] of retrievals methods used when gathering data, namely:

- Sequential, e.g. when we retrieve from our reviews all records between June 2018 and November 2018.
- Random, e.g. when we retrieve from our users records containing information about J. Doe.

The way we search for our data using these two classes is guided by our indexes in order to improve the performance of our search. The method of indexing differs depending on the data type being used.

Traditionally, databases only stored primitive data types, but now support various others types such as IP, timestamps, arrays, UUID, and JSON. Spatial data is typically two-dimensional and this cannot be efficiently indexed with an B-tree but, for example, we would use an R-tree. TigerGraph make use of a supposedly more efficient method of querying spatial data that fits graph architecture appropriately, called a geograph, the performance of which will be compared in our experiments. In the following paragraphs the underlying index structures used for our databases will be discussed.

**B-trees** B-trees can be considered as a generalization of binary search trees [47]. Unlike binary trees, more than two paths may leave a given node depending on the outcome of the query at a node, e.g. at node 0 if  $x > 0$  traverse to node 1,  $x = 0$  traverse to node 2,  $x > 1$  traverse to node 3.

Typical binary trees may become unbalanced after some number of insert and deletion operations, but B-trees always remain balanced. All leaves in a B-tree have the same depth and any search operation among  $n$  records will never visit more than  $1 + \log_d n$  nodes.

B-trees are popular for indexing one-dimensional data and are the default indexing method for many databases, and in our use case, are used for not only primary keys but also temporal indexing. One of the important uses of B-trees is the efficiency gain in sequential and range queries, further optimized by clustering each record in the data by their date fields.

**R-Trees** The nature of spatial data being two-dimensional reveals a shortcoming in using B-trees as the method of efficiently indexing coordinates. Most successful methods of indexing multi-dimensional data have following B-tree-like structures [48] and, in a similar fashion, guided the search to a smaller space.

Traditional R-trees implement indexing by guiding the search toward bounded (hyper) rectangles enclosing the multi-dimensional spatial object. This allows us to query over arbitrary regions such as the nearest restaurants within a 5km radius of a given point without doing a full scan.

Disadvantages of R-trees are that they are slow to update and create a significant redundancy in terms of data storage [49].

ElasticSearch (the search engine used for our JanusGraph configuration) and PostGIS are two examples of technologies that implement R-trees in the database technologies being benchmarked this paper.

**Geograph** The geograph is a grid-based “indexing”<sup>18</sup> solution used by TigerGraph which naturally fits the graph architecture and saving on data storage costs. The idea is that two-dimensional coordinates are mapped to a given grid ID where a grid is represented as a graph vertex. Any vertex associated at that point is then linked by an edge.

A grid can be of any size but setting this size may be dependent on the distribution of points in the dataset. This allows queries to leverage the massively parallel processing (MPP) techniques implored by TigerGraph which create fast updates in contrast to R-trees. The mapping from coordinates to grid ID works in a way such that one can still do searches over an arbitrary region without scanning the whole graph.

A disadvantage of this approach is an uneven distribution of vertices linked to each grid, but this can be managed by manually configuring grid sizes.

## 6.2.2 Relational Design

Figure 16 is the design of the Yelp dataset modelled in a relational database, specifically PostgreSQL. The *location* attribute is indexed with an R-tree using the PostGIS<sup>19</sup> extension.

The decision not to extend *city* and *state* attributes into separate tables was taken, otherwise more joins would be required for an attribute that is never queried. *Location* is the only purely spatial attribute – the main attribute in the benchmarking. It may be faster to simply index state as an attribute due to the low cardinality of city and state paired in a single table. City is indexed and clustered such that records in the same city are physically stored together which, alongside location, should help retrieval speeds from a spatial query perspective. PostgreSQL makes use of B-trees and hash indexes for native data types [50].

The review table is commonly used in queries and holds the most interesting attribute in terms of temporal insight. Since temporal information is one-dimensional and ordinal, the *date* attribute is indexed and clustered.

A business category and a user’s friends are many-to-many relationships. The linking tables *bus\_2\_cat* and *friends* tables handle these relationships.

## 6.2.3 Graph Design

The JanusGraph design can be seen in Figure 17, whereas the TigerGraph is given in Figure 18. The motivation for the difference in the two designs is mainly due to the databases handling spatial data differently.

JanusGraph indexes attributes on nodes and edges using either composite indexes [51], which index native data types on equality conditions, or mixed indexes which leverages

an indexing backend on more complex data types or for complex search predicates, e.g. fuzzy search on strings [52].

Figure 17 show which attributes are indexed using composite indexes and which use mixed indexes – making use of the indexing backend. As in Section 6.2.2, *location* is indexed using R-trees with ElasticSearch’s geo-search predicates. *Date* is indexed using ElasticSearch for equality conditions using the `java.time.Instant` class – this uses a Bkd-tree indexing implementation [53].

TigerGraph puts less of a focus on indexing and more on writing efficient and fast queries. One notable difference between the structure in the JanusGraph implementation and the TigerGraph implementation in Figure 18 is that there are no indexes and the extension of the *location* attribute as the *Business\_Geo* edge and *Geo\_Grid* vertex. The code leveraged for this design idea was inspired from the TigerGraph geospatial webinar [49] and C++ code on the TigerGraph “ecosys”<sup>20</sup>.

# 7 Results and Discussion

This section discusses how each database performed on implementing the kernels mentioned in Section 5.3.2 and discusses the effectiveness of each query language when producing a query to extract the required data for each kernel.

## 7.1 Kate’s Restaurant Recommendation

Figures 19 shows linear growth in the response time of PostgreSQL whereas JanusGraph and TigerGraph remain fairly horizontal over increasing volumes of data. Due to the complexity of this query it comes as no surprise that the graph databases far outperform their relational counterpart. TigerGraph outperforms both PostgreSQL and JanusGraph in terms of query response time and shows a very high consistency as there is almost no standard deviation around the mean.

Table 3 shows which restaurants would be recommended to Kate. One can see that both positive sentiment and star average over reviews for highly regarded restaurants compare well and remain consistent.

The reviews returned by the query were typically well above 3 stars and, since the Naïve Bayes performed well when predicting unseen text data, it comes as no surprise that most of the reviews were tagged as positive. Further analysis could look at how positive sentiment and star rating would compare for inconsistently performing restaurants, but a special measure would need to be put in place to determine how “inconsistency” is measured.

<sup>18</sup>Mentioned in quotations due to TigerGraph not implementing indexes, but rather optimizing data access by the notion of installing queries.

<sup>19</sup><https://postgis.net/>

<sup>20</sup>[https://github.com/tigergraph/ecosys/tree/master/guru\\_scripts/geospatial\\_search](https://github.com/tigergraph/ecosys/tree/master/guru_scripts/geospatial_search)

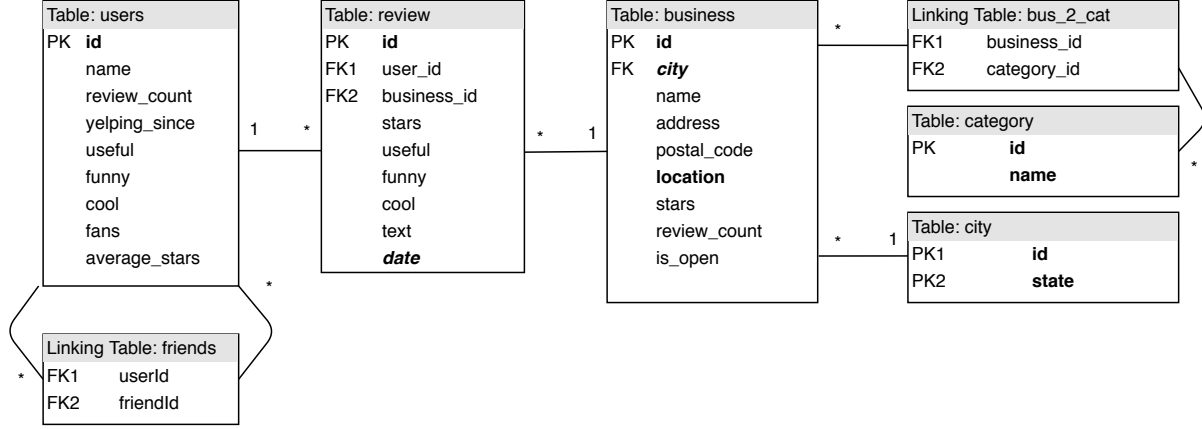


Figure 16: A UML diagram of the relational design of the Yelp dataset. Indexed attributes are in bold whereas clustered attributes are in italics.

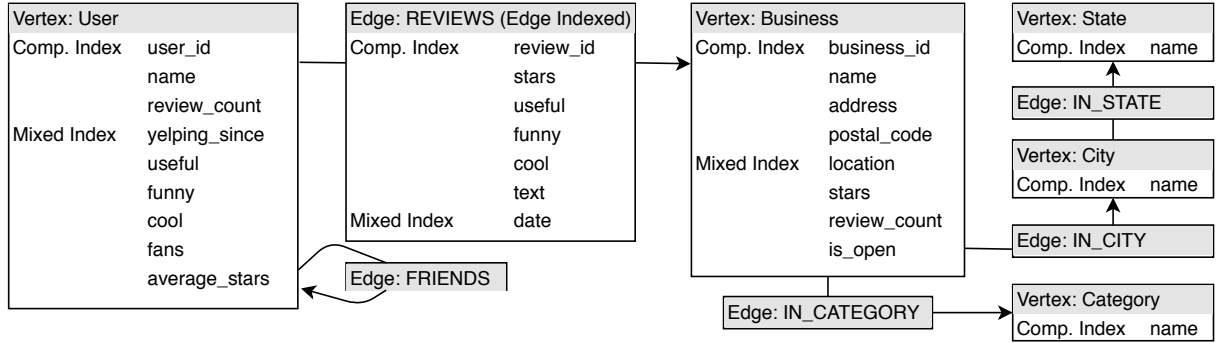


Figure 17: A UML diagram of the graph design for JanusGraph.

Table 3: The result of analysis on the review data of recommended restaurants. Only results with 5 reviews or more are displayed.

Businesses in Phoenix 2018		
Name	Pos Sentiment	Star Average
Paco's Tacos & Tequila	92.9825%	4.5833
Oak Steakhouse Charlotte	100.0%	5.0
The Cheesecake Factory	100.0%	4.4615
Block & Grinder	90.0%	4.3333
Best Wok	100.0%	4.2

## 7.2 Review Trends in Phoenix 2018

Figure 20 shows a phenomena of interest where JanusGraph performs poorly in relation to the other two database technologies with high deviation around the mean. PostgreSQL scales horizontally for this kernel which is most likely be due to the query being the simplest of the three kernels. TigerGraph outperforms both.

The result of a subgraph produced by this query can be seen in Figure 21.

The general trends shown in the review data for Phoenix during 2018 have the following characteristics:

- More critical, lower scoring reviews tend to be longer and most useful.
- Reviews with 3 or 4 stars seem to be the funniest.
- Reviews with 4 or 5 stars tend to be the coolest.

The percentage positive sentiment, when scored relatively, is ordered consistently with the average star rating. This validates the performance of the binary sentiment classifier in that one almost does not need to see the star rating and can rely on text data alone when considering a broad spectrum of reviews.

This analysis could be performed over varied year brackets and different areas to see if performance is consistent or not. The implication of this could lead to experimenting with more sophisticated machine learning models on the dataset to be more precising in that it could potentially predict the star rating as is done in [54] and [55].

## 7.3 Ranking Las Vegas by Friends' Sentiment

Figure 27 shows both graph database technologies outperform PostgreSQL as PostgreSQL shows linearly growth as it did in Figures 19. For the experiments on 13% of the dataset, JanusGraph shows a lot of deviation around the

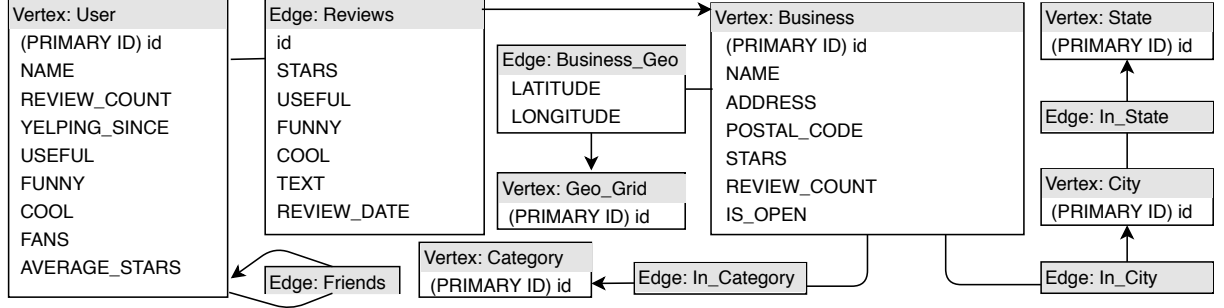


Figure 18: A UML diagram of the graph design for TigerGraph.

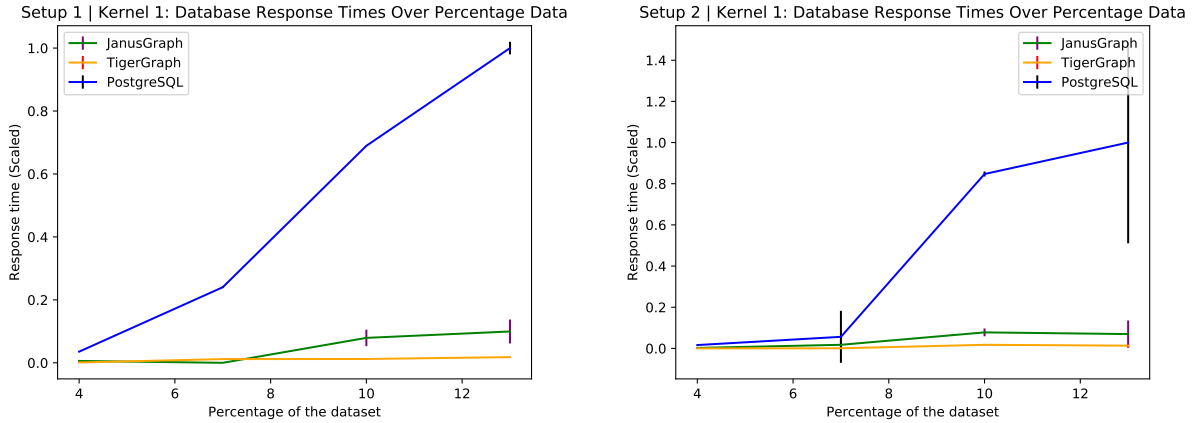


Figure 19: Database response times over varying percentages of the dataset for setup 1 and 2 for the kernel: “Kate’s Restaurant Recommendation”. The error bars display standard deviation.

mean. This may be due to all the moving parts on which JanusGraph is implemented on and it’s multi-level caching implementation behaving poorly at this size of the dataset.

The result of this analysis was focused more on the performance rather than the data extracted. The results of the data analysis on this kernel does not show anything more interesting about the review data than what was already discussed in Section 7.2. One notable difference in this kernel however, is that it produces a much more complex query and the databases perform accordingly.

The result of this kernel is simply to show that the results may vary and can be correlated depending on the relationships between different data points.

Table 4: The result of analysis on the review data of Julie’s friends.

Las Vegas Sentiment vs Star Average	
Positive Sentiment (%)	Star Average
70.7767	3.8917

## 7.4 Queries

**SQL** SQL is a mature and well supported querying language which makes it simple to implement a solution. The caveat of this simplicity is that the resulting solution may long and convoluted for complex queries – such as the one produced in Listing 1. The SQL queries produced for these kernels have a good balance between readability and expressiveness but, as complexity grew, so did size and queries began to lose the readability aspect.

SQL handles temporal data well and, in the PostgreSQL dialect, comes well supported with functions to operate on various temporal data types. This level of support provides ease of programming when implementing a SQL-based solution to a dataset with spatio-temporal properties.

**Gremlin** Gremlin was found to produce the most concise queries of the three languages. The limitation of Gremlin is that, if one makes use of the mixed indexing search predicates, one may be limited to programming languages with support from these drivers to have the embedded Gremlin functionality. In the context of the technologies implemented in this investigation, a JVM language would be better suited as the backend for a JanusGraph data storage

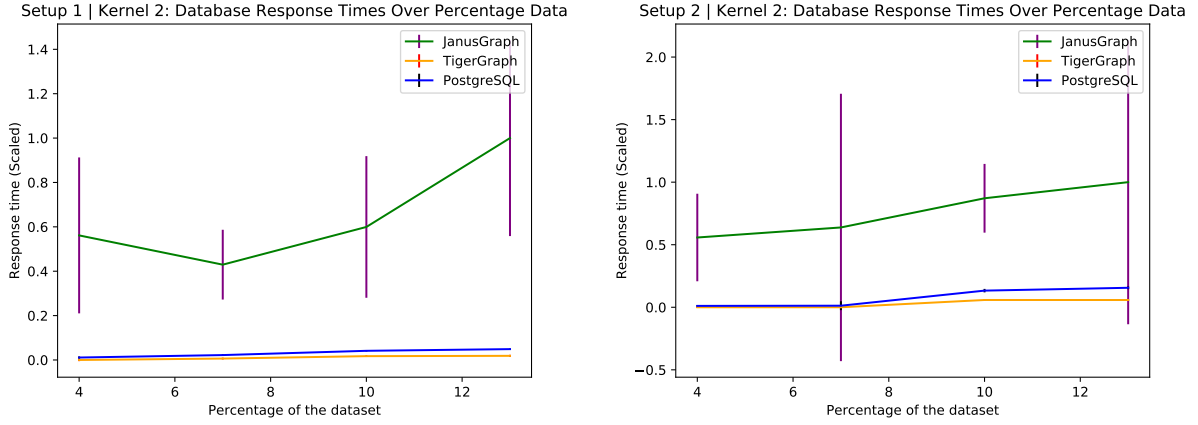


Figure 20: Database response times over varying percentages of the dataset for setup 1 and 2 for the kernel: “Review Trends in Phoenix 2018”. The error bars display standard deviation.

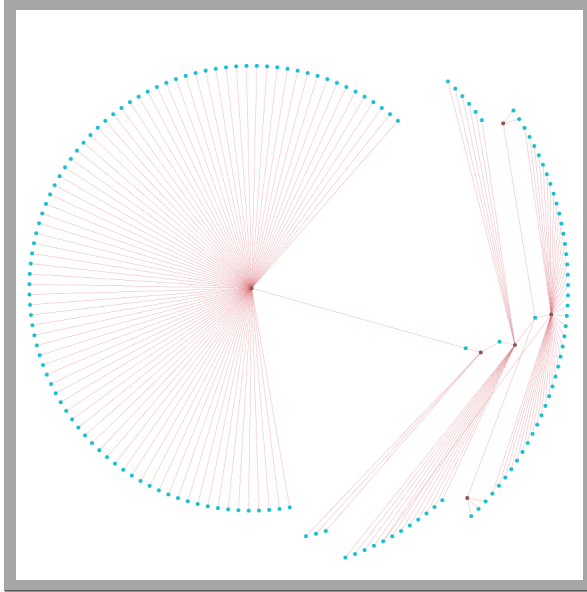


Figure 21: A subset of the graph produced by TigerGraph on the result of the query for all reviews in the Phoenix area in 2018. Maroon edges represent reviews. Blue vertices are the users and brown vertices represent the businesses.

solution. The Gremlin queries produced for these kernels were found to be readable in terms of describing the data flow of the traversal within a graph topology context. One may not enjoy Gremlin’s referencing steps going back and forth within a query using the `as` and `select` steps but, after some experience with Gremlin, this will no longer be an issue.

The ability of certain steps allowing one to skip across edges to refer to vertices directly is part of why Gremlin is able to produce such concise queries. The performance of these queries heavily relies on the data flow produced

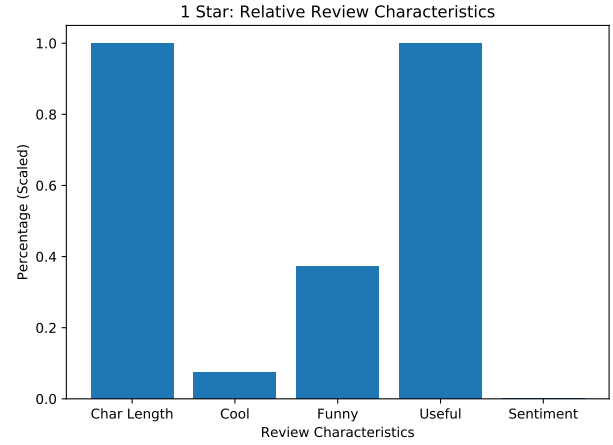


Figure 22: Relative characteristics of 1 star reviews over reviews of restaurants in Phoenix 2018.

by the ordering of such steps. This makes it important to use `filter` steps and be conscious of the ordering of each step.

Gremlin is well supported and has an extensive documentation<sup>21</sup> but is a vastly different querying language when compared to SQL. The implication of this is that there is a small but significant learning curve involved. Effective imperative Gremlin queries will most likely only be written after some experience. Fortunately, Gremlin supports declarative querying which allows a user new to Gremlin to write effective queries with little experience.

**GSQL** The GSQL queries produced by each kernel resulted in the queries with the most vertical space of the three query languages. This is necessary for segmenta-

<sup>21</sup><http://tinkerpop.apache.org/docs/current/reference/>

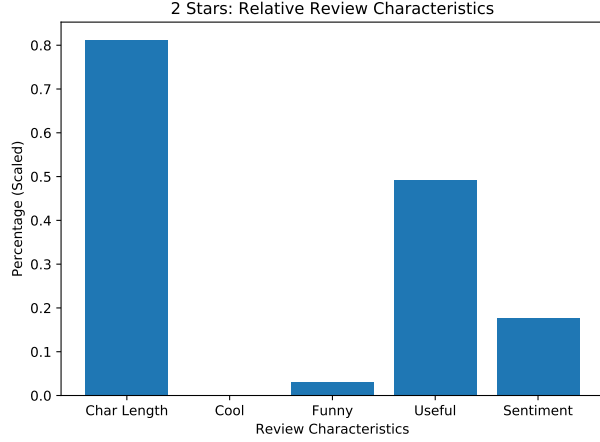


Figure 23: Relative characteristics of 2 star reviews over reviews of restaurants in Phoenix 2018.

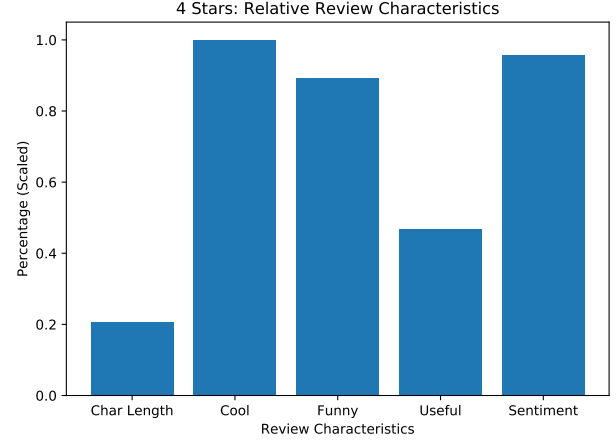


Figure 25: Relative characteristics of 4 star reviews over reviews of restaurants in Phoenix 2018.

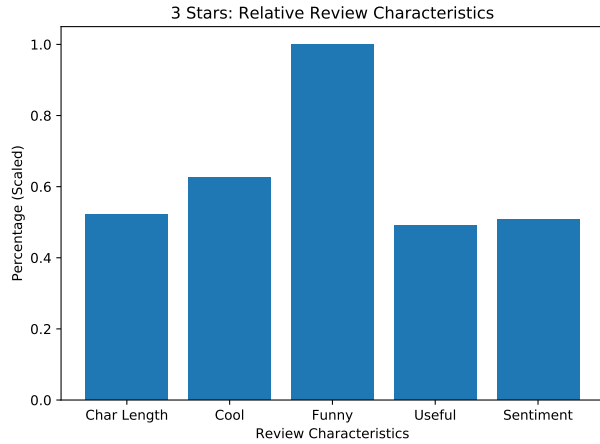


Figure 24: Relative characteristics of 3 star reviews over reviews of restaurants in Phoenix 2018.

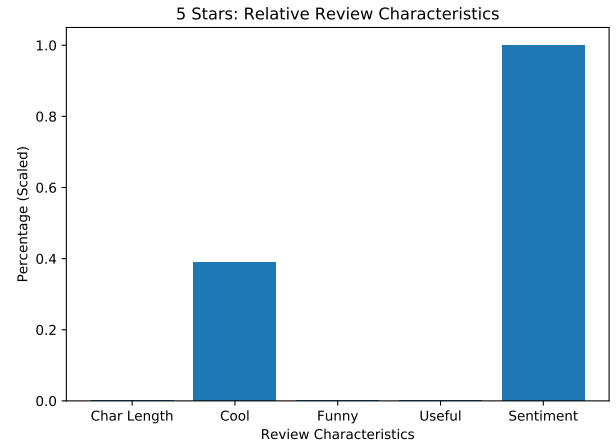


Figure 26: Relative characteristics of 5 star reviews over reviews of restaurants in Phoenix 2018.

tion of the query which is used for parallel graph traversals. The result of this segmentation and vertical space has made each query extremely readable and expressive. The conservative use of ASCII art and use of keywords from SQL provides a good balance between query visualization and familiarity. The development of queries using GraphStudio reduces the learning curve significantly as queries are developed in statically typed, compiler driven context – only allowing one to install a query once all errors have been addressed.

GSQL is well suited for spatial queries using the geo-grid approach – which integrates well within a graph topology – and temporal queries with a selection of built-in function for manipulating temporal data types – as with SQL. By segmenting the query, the compiler is able to determine what can automatically be executed in parallel which adds to the fast response times of TigerGraph when compared to the other two databases.

The Rest++ API allows one to write a parameterized query once and access it anywhere without having to worry about driver issues other than being able to communicate with a REST API. This was a particular pleasure in the post-query development of connecting a web application backend to communicate with TigerGraph.

## 8 Conclusion and Future Work

In this paper, we analyzed and compared the response times of three database technologies with respect to handling interconnected spatio-temporal data. The technologies compared were two open source database technologies, PostgreSQL and JanusGraph, and one enterprise level technology, TigerGraph. The linear growth in the relational model was clearly illustrated in the results whereas the graph database solutions scaled more horizontally. This alone is an advantage NoSQL databases have over traditional relational models when querying large volumes of



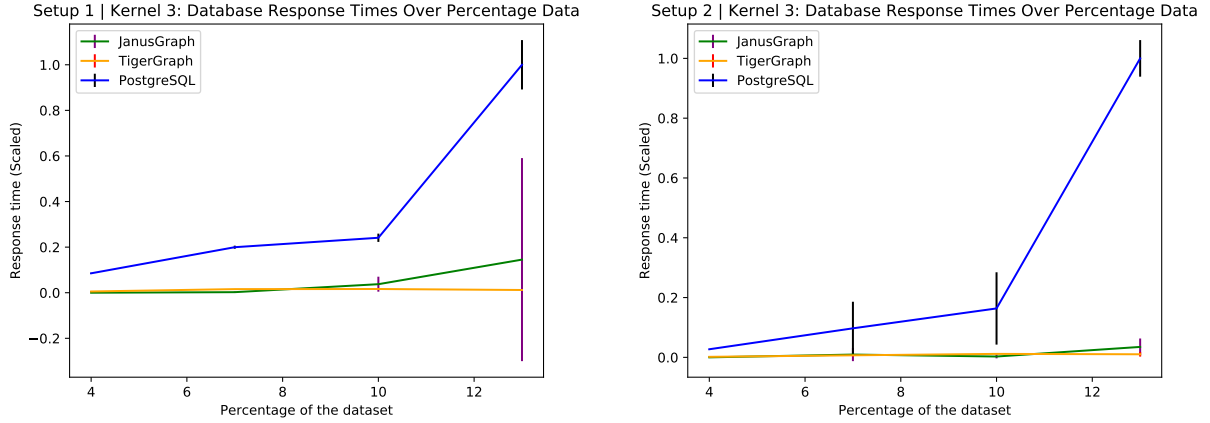


Figure 27: Database response times over varying percentages of the dataset for setup 1 and 2 for the kernel: “Ranking Las Vegas by Friends’ Sentiment”. The error bars display standard deviation.

data. These three systems were evaluated by employing a set of spatio-temporal queries similar to those that would be found in real world scenarios when analysing data in a dataset such as the Yelp Challenge Dataset.

The results show that graph database technology has been shown to outperform PostgreSQL in all three of the kernels. This result is partially due to the fact that the kernels produce complex queries due to the interconnected nature of the data. This dataset produced a dense graph which graph databases have the ability to perform effective traversals over when compared to multi-join style queries produced by the relational implementation. The spatio-temporal multi-dimensional aspect has shown to be supported well in all of the databases and evident by the response times of the queries with constraints of this nature.

**Benchmark Results** The results in Section 7 strongly suggest that graph database technology and, specifically TigerGraph, provide the faster response times when querying this dataset applied with spatio-temporal constraints. One notable observation is how inconsistent JanusGraph performed and this may be due to JanusGraph’s caching implementation. JanusGraph maintains multiple levels of caching both on the transaction level and database level. This excludes the storage backend’s caching – in this case Cassandra. The cache has an expiration time and, since these experiments were run serially but chosen randomly, the JanusGraph specific jobs were run out of order and the cache could have expired at this time.

Another potential reason for the inconsistent gradient between the means in each result may be due to the fact that the dataset adds unpredictable levels of complexity – in terms of how connected the data is – at the end of an import for a given percentage. The horizontal scaling for the graph databases suggests that the impact of this is minimal.

Nevertheless, as the queries became more and more complicated, the graph databases maintained a horizontal scale whereas PostgreSQL grew linearly in these cases. When the queries were not complicated, as was the case in Section 7.2, PostgreSQL responded nearly as fast as TigerGraph and outperforming JanusGraph.

**Visualization** Graph databases have an advantage in data visualization with tools such as Cytoscape<sup>22</sup> or GraphStudio which is built into TigerGraph. This can be important in use cases involving further analysis on data patterns such as those found on social datasets such as the Yelp dataset. Relational databases can have stored data processed and re-shaped into a graph structure but this requires extra overhead and configuration as it is not a native topology of this technology.

**Migration and Product Maturity** Graph database technology is currently under rapid development, each vendor with their own API and query language. Each graph query language is designed to express graph traversals in a more graph-oriented approach. The learning curve from SQL may pose an issue when migrating but languages such as Cypher or GSQL make this minimal by applying concepts from SQL to the graph traversal context.

Security and reliability used to be an issue when considering migrating to graph database technology but both of these vendors can be configured to use encrypted communication and can be robust to failures. For example, TigerGraph’s Rest++ API can be encrypted, integrated with Single Sign-On, and require authorization with LDAP authentication. JanusGraph transactions can be configured to be ACID-compliant when using BerkeleyDB but this is not generally the case with Cassandra or HBase. TigerGraph and Neo4j are ACID-compliant so it is in the

<sup>22</sup><https://cytoscape.org/>



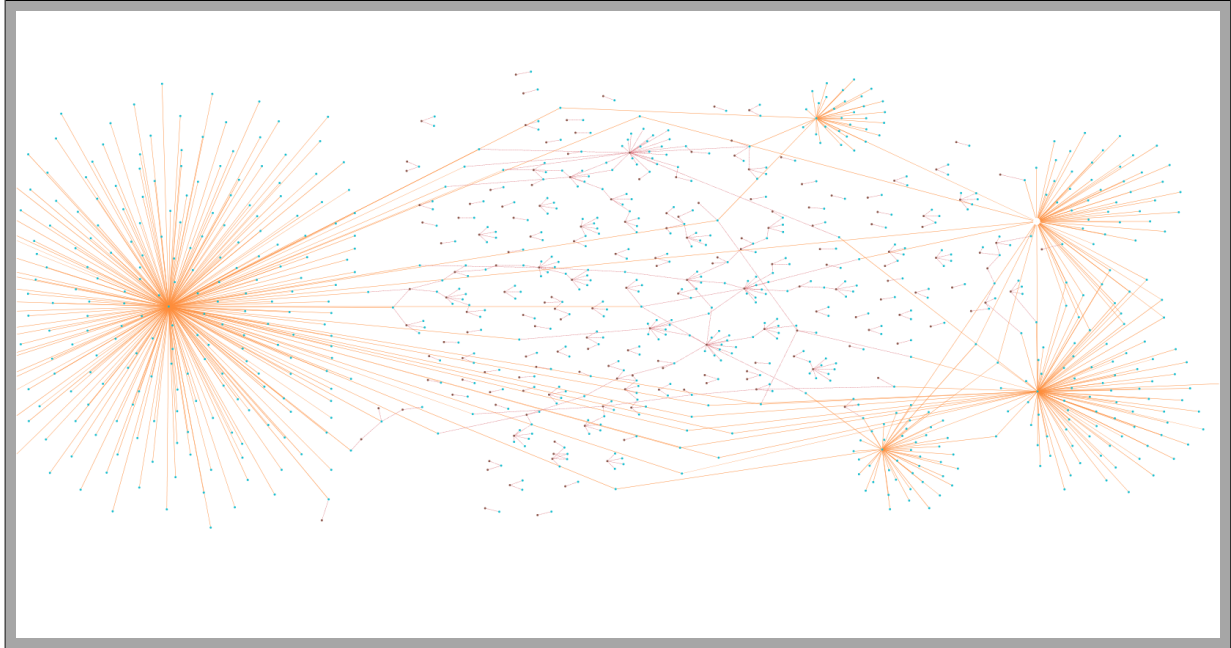


Figure 28: A subset of the graph produced by TigerGraph on the result of the query for this kernel. Orange edges represent friend relations and maroon edges represent reviews. Blue vertices represent users and brown vertices represent businesses. The white center of the cluster on the top right is Julie and once can see the center of the giant cluster on the left is a mutual friend of Julie's.

position to compete with relational databases with regards to reliable transactions. The implication of this is that graph database technology can provide the same level of robustness and security as relational database technology can so one should not have to sacrifice either of these two properties when migrating.

**Data Structure** Before pre-processing, the Yelp dataset contains additional attributes many of which are missing or partial and this makes the dataset semi-structured. Since a relational database schema is fixed, this would increase the number of tables in a relational database for each potential attribute that could have a relationship to any other data entry. Graph databases are schemaless and are well-suited to handle such unstructured or semi-structured data.

**Query Languages** Of the two graph querying languages, GSQL was found to be the easiest to learn and implement and, with an imperative and statically typed language, many developers may find GSQL very familiar. The Rest++ API feature was found to greatly enhance the post-query process due to any other applications only having to query a parameterized HTTP endpoint. GSQL also adds a lot of flexibility in terms of how the data is formatted and structured in the HTTP response which helps for seamless deserializing of the JSON result.

**Future Work** This paper investigated the response times and, to some degree, how effective each query language was in producing queries to return the necessary data. Neither the storage efficiency nor performance capabilities for varied limitations on hardware were measured in any sophisticated way. Investigating these issues for large-scale spatio-temporal data would only add to the findings in this paper in terms of how suitable graph database technology is.

Only one dataset was used in this investigation but other spatio-temporal datasets of varying quantities should be benchmarked. Doing so would create a more robust conclusion to the suitability of graph database technology for storing and querying spatio-temporal data.

It would be worthwhile to add not only more graph database technologies, but other NoSQL, SQL, and NewSQL solutions. One could, for example, investigate the performance of Cassandra with Elasticsearch and measure whether the graph abstraction provided by JanusGraph is truly beneficial or not. This would add to a more complete view on how well suited each database solution is for spatio-temporal data, or large-scale data in general. There are new relational database technologies which employ auto-sharding and other techniques to help traditional SQL technologies scale horizontally such as MySQL Cluster<sup>23</sup> and the Citus<sup>24</sup> extension for PostgreSQL. Using

<sup>23</sup><https://www.mysql.com/products/cluster/>

<sup>24</sup><https://www.citusdata.com/product>

these technologies could help relational database technologies compete with NoSQL technologies when facing large-scale data in general.

## References

- [1] Antonios Makris, Konstantinos Tserpes, Giannis Spiliopoulos, and Dimosthenis Anagnostopoulos. Performance evaluation of mongodb and postgresql for spatio-temporal data. In *EDBT/ICDT Workshops*, 2019.
- [2] Yaowen Chen et al. *Comparison of Graph Databases and Relational Databases When Handling Large-Scale Social Data*. PhD thesis, University of Saskatchewan, 2016.
- [3] Luke Sloan and Jeffrey Morgan. Who tweets with their location? understanding the relationship between demographic characteristics and the use of geoservices and geotagging on twitter. *PloS one*, 10(11):e0142209, 2015.
- [4] Jaroslav Pokorný. Nosql databases: a step to database scalability in web environment. *International Journal of Web Information Systems*, 9(1):69–82, 2013.
- [5] Yelp Inc. Yelp Dataset Challenge. <https://www.yelp.com/dataset/challenge/>, Jan-Dec 2019.
- [6] Armin Ronacher and contributors. Flask. <https://palletsprojects.com/p/flask/>, 2015.
- [7] Google and contributors. Angular. <https://angular.io/>, 2010.
- [8] Ian. What does ACID mean in Database Systems? <https://database.guide/>, June 2016.
- [9] Seth Gilbert and Nancy Lynch. Perspectives on the cap theorem. *Computer*, 45(2):30–36, 2012.
- [10] Edgar F Codd. Relational database: a practical foundation for productivity. In *Readings in Artificial Intelligence and Databases*, pages 60–68. Elsevier, 1989.
- [11] ABM Moniruzzaman and Syed Akhter Hossain. Nosql database: New era of databases for big data analytics-classification, characteristics and comparison. *arXiv preprint arXiv:1307.0191*, 2013.
- [12] Alin Deutsch, Yu Xu, Mingxi Wu, and Victor Lee. Tigergraph: A native mpp graph database. *arXiv preprint arXiv:1901.08248*, 2019.
- [13] Robert Kallman, Hideaki Kimura, Jonathan Natkins, Andrew Pavlo, Alexander Rasin, Stanley Zdonik, Evan PC Jones, Samuel Madden, Michael Stonebraker, Yang Zhang, et al. H-store: a high-performance, distributed main memory transaction processing system. *Proceedings of the VLDB Endowment*, 1(2):1496–1499, 2008.
- [14] Deka Ganesh Chandra. Base analysis of nosql database. *Future Generation Computer Systems*, 52:13–21, 2015.
- [15] Luca Castellano. Distributed, transactional key-value store, May 19 2015. US Patent 9,037,556.
- [16] V Manoj. Comparative study of nosql document, column store databases and evaluation of cassandra. *International Journal of Database Management Systems*, 6(4):11, 2014.
- [17] Haoyu Tan, Wuman Luo, and Lionel M Ni. Clost: a hadoop-based storage system for big spatio-temporal data analytics. In *Proceedings of the 21st ACM international conference on Information and knowledge management*, pages 2139–2143. ACM, 2012.
- [18] Jaymie R Meliker, Melissa J Slotnick, Gillian A AvRuskin, Andrew Kaufmann, Geoffrey M Jacquez, and Jerome O Nriagu. Improving exposure assessment in environmental epidemiology: Application of spatio-temporal visualization tools. *Journal of Geographical Systems*, 7(1):49–66, 2005.
- [19] A. Govardhan K. Venkateswara Rao and K.V. Chalapathi Rao. Spatiotemporal data mining: Issues, tasks and applications. *International Journal of Computer Science and Engineering Survey*, 3(1):39, 2012.
- [20] The GraphQL Foundation. GraphQL. <https://graphql.org/>, 2019.
- [21] TigerGraph. Tigergraph docs : Tigergraph v1.0. <https://doc-archive.tigergraph.com/1.0/TigerGraph.html>, 2017.
- [22] Apache. The Gremlin Graph Traversal Machine and Language. <https://tinkerpop.apache.org/gremlin.html>, 2019.
- [23] Apache. TinkerPop3 Documentation. <http://tinkerpop.apache.org/docs/3.3.0/reference/>, 2019.
- [24] Neo4j Inc. openCypher. <https://www.opencypher.org>, 2018.
- [25] Yu Xu. Modern Graph Query Language – GSQL. <https://www.kdnuggets.com/>, 2019.
- [26] Inc. Neo4j. The pagerank algorithm. <https://neo4j.com/docs/graph-algorithms/current/algorithms/page-rank/>, 2019.
- [27] George Anadiotis. Back to the future: Does graph database success hang on query language? <https://www.zdnet.com/>, 2018.
- [28] *Querying Graph Databases with the GSQL Query*, 2018.
- [29] Martin Heller. TigerGraph review: A graph database designed for deep analytics. <https://www.infoworld.com/>, 2018.
- [30] Renzo Angles, Marcelo Arenas, Pablo Barceló, Aidan Hogan, Juan L Reutter, and Domagoj Vrgoc. Foundations of modern graph query languages. *CoRR, abs/1610.06264*, 2016.
- [31] Riko Jacob Christopher L. Barrett and Madhav V. Marathe. Formal-language-constrained path problems. *SAIM J. Comput.*, 30(3):809–837, 2000.

- [32] The PostgreSQL Global Development Group. PostgreSQL About. <https://www.postgresql.org/about/>, 2019.
- [33] JanusGraph Authors. JanusGraph. <https://janusgraph.org/>, 2019.
- [34] JanusGraph Authors. Migrating from Titan. <https://docs.janusgraph.org/advanced-topics/migrating/>, 2019.
- [35] JanusGraph Authors. Introduction: The Benefits of JanusGraph. <https://docs.janusgraph.org/>, 2019.
- [36] Titan Authors. Chapter 1. The Benefits of Titan. [http://s3.thinkaurelius.com/docs/titan/1.0.0/benefits.html#\\_titan\\_and\\_the\\_cap\\_theorem](http://s3.thinkaurelius.com/docs/titan/1.0.0/benefits.html#_titan_and_the_cap_theorem), 2015.
- [37] Oleksandr Porunov. JanusGraph GitHub. <https://github.com/JanusGraph/janusgraph/blob/master/README.md#powered-by-janusgraph>, 2019.
- [38] Florin Rusu & Zhiyi Huang. Benchmarking Graph Analytic Systems: TigerGraph, Neo4j, Neptune, JanusGraph, and ArangoDB. arXiv:1907.07405, 2019.
- [39] Trip Beernink. Private communication. 2019.
- [40] Lila Razzaqui. TigerGraph Wins Strata Data’s “Most Disruptive Startup” Award. <https://www.globenewswire.com/>, 2018.
- [41] MCMXCV–MMXIX Encyclopedia Mythica. Providentia. <https://pantheon.org/articles/p/providentia.html>, 1997.
- [42] Shashank Gupta. Sentiment Analysis: Concept, Analysis and Applications. <https://towardsdatascience.com/>, 2018.
- [43] Irina Rish et al. An empirical study of the naive bayes classifier. *IJCAI 2001 workshop on empirical methods in artificial intelligence*, 3(22):41–46, 2001.
- [44] Eduardo Blanco and Dan Moldovan. Some issues on detecting negation from text. In *Twenty-Fourth International FLAIRS Conference*, 2011.
- [45] Samira Munir. Basic Sentiment Analysis using NLTK. <https://towardsdatascience.com/>, 2019.
- [46] Zhiwei Zhang. Machine Learning and Visualization with Yelp Dataset. <https://medium.com/>, 2017.
- [47] Douglas Comer. Ubiquitous b-tree. *ACM Computing Surveys (CSUR)*, 11(2):121–137, 1979.
- [48] DM Gavril. *R-tree index optimization*. University of Maryland, Center for Automation Research, Computer Vision ..., 1994.
- [49] Xinyu Chang. Graph Gurus Episode 8: Location, Location, Location. [https://www.youtube.com/watch?v=gPF\\_SXibDxw&t=887s](https://www.youtube.com/watch?v=gPF_SXibDxw&t=887s), 2019.
- [50] Jason Harris. PostgreSQL Vs. MySQL: Differences In Performance, Syntax, And Features. <https://blog.panoply.io/postgresql-vs.-mysql>, 2018.
- [51] JanusGraph. Indexing for Better Performance. <https://docs.janusgraph.org/v0.2/basics/index-performance/>, 2019.
- [52] JanusGraph. Index Backend. <https://docs.janusgraph.org/v0.2/index-backend/>, 2019.
- [53] Nick Knize. Numeric and Date Ranges in Elasticsearch: Just Another Brick in the Wall. <https://www.elastic.co/>, 1 2017.
- [54] Ch Sarath Chandra Reddy, K Uday Kumar, J Dheeraj Keshav, Bakshi Rohit Prasad, and Sonali Agarwal. Prediction of star ratings from online reviews. In *TENCON 2017-2017 IEEE Region 10 Conference*, pages 1857–1861. IEEE, 2017.
- [55] Dagmar Monett and Hermann Stolte. Predicting star ratings based on annotated reviews of mobile apps. In *2016 Federated Conference on Computer Science and Information Systems (FedCSIS)*, pages 421–428. IEEE, 2016.