

# Property Graph Schema Optimization for Domain-Specific Knowledge Graphs

Chuan Lei<sup>1</sup>, Rana Alotaibi<sup>2</sup>, Abdul Quamar<sup>1</sup>, Vasilis Efthymiou<sup>1</sup>, Fatma Özcan<sup>1</sup>

IBM Research - Almaden<sup>1</sup>, University of California at San Diego<sup>2</sup>

chuan.lei|vasilis.efthymiou@ibm.com, ahquamar|fozcan@us.ibm.com, ralotaib@eng.ucsd.edu

## ABSTRACT

Enterprises are creating domain-specific knowledge graphs by curating and integrating their business data from multiple sources. The data in these knowledge graphs can be described using ontologies, which provide a semantic abstraction to define the content in terms of the entities and the relationships of the domain. The rich semantic relationships in an ontology contain a variety of opportunities to reduce edge traversals and consequently improve the graph query performance. Although there has been a lot of effort to build systems that enable efficient querying over knowledge graphs, the problem of schema optimization for query performance has been largely ignored in the graph setting. In this work, we show that graph schema design has significant impact on query performance, and then propose optimization algorithms that exploit the opportunities from the domain ontology to generate efficient property graph schemas. To the best of our knowledge, we are the first to present an ontology-driven approach for property graph schema optimization. We conduct empirical evaluations with two real-world knowledge graphs from medical and financial domains. The results show that the schemas produced by the optimization algorithms achieve up to 2 orders of magnitude speed-up compared to the baseline approach.

## 1 INTRODUCTION

Domain-specific knowledge graphs are playing an increasingly important role to derive business insights in many enterprise applications such as customer engagement, fraud detection, network management, etc [38]. One distinct characteristic of these enterprise knowledge graphs, compared to the open-domain knowledge graphs like DBpedia [31], Freebase [11], and YAGO2 [44], is their deep domain specialization. The domain specialization is typically captured by an ontology which provides a semantic abstraction to describe the entities and their relationships of the data in the knowledge graphs. A few widely used domain-specific ontologies include Unified Medical Language System (UMLS)<sup>1</sup>

and SNOMED Clinical Terms<sup>2</sup> in the medical domain, Financial Industry Business Ontology (FIBO)<sup>3</sup> and Financial Report Ontology (FRO)<sup>4</sup> in the financial domain, and many more in various other domains<sup>5</sup>. The ontology is often used to drive the creation of a knowledge graph by ingesting and transforming raw data from multiple sources into standard terminologies. The curated knowledge graphs allow users to express their queries in standard vocabularies, which promotes more interoperable and effective enterprise applications and services for specific domains [17, 22].

There are two popular approaches to store and query knowledge graphs: RDF data model and SPARQL query language [42] or property graph model and graph query languages such as Gremlin [45] and Cypher [25]. An important difference between RDF and property graphs is that RDF regularizes the graph representation as a set of triples, which means that even literals are represented as graph vertices. Such artificial vertices make it hard to express graph queries in a natural way. The property graph model instead uses vertices to represent entities and edges to represent the relationships between them, with each specified using key-value properties pairs [49]. For this reason, property graph systems are rapidly gaining popularity for graph storage and retrieval. Examples include Neo4j [6], Apache JanusGraph [5], Azure Cosmos DB [2], Amazon Neptune [1], to name a few. Many graph applications (e.g., community detection, centrality analysis, and link prediction) heavily rely on the performance of graph queries over the property graph systems. Many techniques have been proposed for optimizing the query performance, system scalability, and transaction support for these systems [13, 34, 36, 48]. However the problem of property graph schema optimization has been largely ignored, which is also critical to graph query performance.

In this paper, we tackle the *property graph schema optimization problem* for domain-specific knowledge graphs. Our goal is to create an optimized schema<sup>6</sup> based on a given

<sup>2</sup><http://www.snomed.org/>

<sup>3</sup><https://spec.edmouncil.org/fibo/>

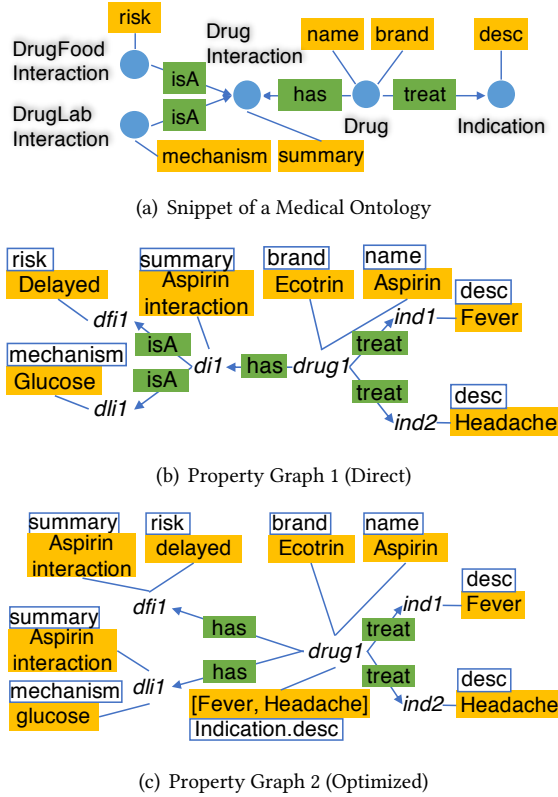
<sup>4</sup><http://www.xbrlsite.com/2015/fro/us-gaap/xbrl/Ontology/Overview.html>

<sup>5</sup><https://lod-cloud.net/>

<sup>6</sup>We use the terms property graph schema, graph schema, and schema interchangeably.

<sup>1</sup><https://www.nlm.nih.gov/research/umls/index.html>

ontology, such that the corresponding property graph can efficiently support various types of graph queries (e.g., pattern matching, path finding, or aggregation queries) with better query performance. The raw data is loaded directly as a property graph that conforms to the optimized schema<sup>7</sup>. One straightforward way to create a property graph schema from an ontology is to directly map each ontology concept to a schema node, and to map each ontology relationship to a schema edge, analogous to ER diagram to relational schema mapping. However, we argue that the graph query performance varies vastly for different property graphs with the same data but corresponding to different schemas, and the rich semantic information in the ontology provides unique opportunities for schema optimization. We illustrate this using two examples from the medical domain.



**Figure 1: Motivating Example.**

*Example 1 (Pattern matching query).* Consider the ontology in Figure 1(a), *summary* is a property of *DrugInteraction* concept, which is connected to *DrugFoodInteraction* and

<sup>7</sup>A property graph schema may not be logically equivalent to a given ontology. Capturing the full expressivity of ontologies (e.g., negation, role inclusion, transitivity) in the form of a property graph schema is an unexplored and challenging problem, which is beyond the scope of this work.

*DrugLabInteraction* concepts via inheritance (*isA*) relationships. Figures 1(b) and 1(c) show two alternative property graphs conforming to two different schemas with several vertices and edges. In Figure 1(b), the vertex *di1* (i.e., an instance of *DrugInteraction*) leads to both *dfi1* and *dli1*. In Figure 1(c), *drug1* directly connects to *dfi1* and *dli1* vertices. For any query that requires edge traversals from *drug1* to either *dfi1* or *dli1* or both, the property graph 2 clearly requires less number of edge traversals. A pattern matching query interested in *Drug* and the associated *risk* of *DrugFoodInteraction* achieves 2 orders of magnitude performance gains on the optimized property graph (23ms) compared to the property graph 1 (3245ms).

*Example 2 (Aggregation query).* In Figure 1(a), *Drug* concept is also connected to *Indication* concept via a *treat* (1:M) relationship. In this case, we observe that if we replicate certain properties accessible via a 1:M relationship, edge traversals can be avoided. Figure 1(c) shows that the vertex *drug1* has an additional property, which is a list of descriptions replicated from the property *desc* of *ind1* and *ind2*. An aggregation query (COUNT) on the *desc* of *Indication* treated by *Drug* runs 8 times faster on this optimized property graph (78ms) than the property graph 1 (627ms). In this case, avoiding the edge traversals is extremely beneficial, especially when the number of edges between these two types of vertices is large.

These two examples show that edge traversal is one of the dominant factors affecting graph query performance, and having an optimized schema can greatly improve query performance. We can reduce edge traversals by merging nodes or replicating data. However, this needs to be done carefully, as the resulting knowledge graph needs to preserve its semantics information. Fortunately, the rich semantic relationships in an ontology provide a variety of opportunities to reduce graph traversals. To generate an optimized graph schema, we need to identify and exploit these opportunities in the ontology, and design different techniques to utilize them accordingly. As illustrated in the examples, certain optimization techniques require data replication resulting in space overheads. Hence, the schema optimization has to trade off between the query performance and the space consumption of the resulting property graph.

**Our proposed approach.** To the best of our knowledge, we are the first to address the problem of property graph schema optimization to improve graph query performance. In addition to the ontology, our approach also takes into account the space constraints, if any, and additional information such as data distribution and workload summaries<sup>8</sup>. We propose a set of rules that are designed to optimize the graph query performance with respect to different types of

<sup>8</sup>We refer to the access frequency of concepts, relationships and properties as workload summaries which will be formally defined later.

relationships in the ontology. When there is a space constraint, we estimate the cost-benefit of applying these rules to each individual relationship by leveraging the additional data distribution and workload information. We propose two algorithms, concept-centric and relation-centric, which incorporate the cost-benefit scores to produce an optimized property graph schema. Our approach can seamlessly handle updates to the property graph, as long as its schema remains unchanged.

**Contributions.** The contributions of this paper can be summarized as follows:

1. We introduce an ontology-driven approach for property graph schema optimization.
2. We design a set of rules that reduce the edge traversals by exploiting the rich semantic relationships in the ontology, resulting in better graph query performance.
3. We propose concept-centric and relation-centric algorithms that harness the proposed rules to generate an optimized property graph schema from an ontology, under space constraints. The concept-centric algorithm utilizes the centrality analysis of concepts, and the relation-centric algorithm uses a cost-benefit model.

4. Our experiments show that our ontology-driven approach effectively produces optimized graph schemas for two real-world knowledge graphs from medical and financial domains. The queries over the optimized property graphs achieve up to 2 orders of magnitude performance gains compared to the graphs resulting from the baseline approach.

The rest of the paper is organized as follows. Section 2 introduces the basic concepts, formulates the problem, and provides an overview of our ontology-driven approach. Section 3 describes our optimization rules for different types of relationships in an ontology. Section 4 explains the algorithms to produce optimized property graph schema. We provide our experimental results in Section 5, review related work in Section 6, and finally conclude in Section 7.

## 2 PRELIMINARIES & APPROACH OVERVIEW

### 2.1 Preliminaries

An ontology describes a particular domain and provides a structured view of the data. Specifically, it provides an expressive data model for the concepts that are relevant to that domain, the properties associated with the concepts, and the relationships between concepts.

*Definition 1 (Ontology ( $O$ )).* An ontology  $O(C, R, P)$  contains a set of concepts  $C = \{c_n | 1 \leq n \leq N\}$ , a set of data properties  $P = \{p_m | 1 \leq m \leq M\}$ , and a set of relationships between the concepts  $R = \{r_k | 1 \leq k \leq K\}$ .

An ontology is typically described in OWL [7], wherein a concept is defined as a *class*, a property associated with a concept is defined as a *DataProperty* and a relationship between a pair of concepts is defined as an *ObjectProperty*. Each DataProperty  $p_i \in P_n$  represents a characteristic of a concept  $c_n \in C$  and  $P_n \subseteq P$  represents the set of DataProperties associated with the concept  $c_n$ . Each ObjectProperty  $r_k = (c_s, c_d, t)$  is associated with a source concept  $c_s \in C$ , also referred to as the domain of the ObjectProperty, a destination concept  $c_d \in C$ , also referred to as the range of the ObjectProperty, and a type  $t$ . The type  $t$  can be either a functional (i.e., 1:1, 1:M, M:N), an inheritance (a.k.a *isA*) or a union/membership relationship<sup>9</sup>. In this paper, we use the ontology as a semantic data model of a knowledge graph.

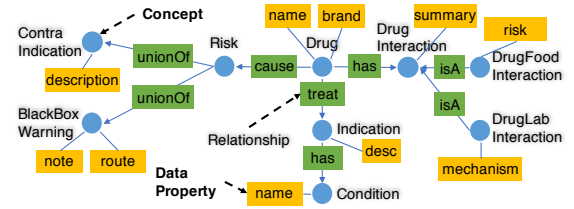


Figure 2: Medical Ontology.

We adopt the widely used property graph model from [41].

*Definition 2 (Property Graph ( $\mathcal{PG}$ )).* A property graph  $\mathcal{PG}(V, E)$  is a directed multi-graph with vertex set  $V$  and edge set  $E$ , where each node  $v \in V$  and each edge  $e \in E$  has data properties consisting of multiple attribute-value pairs.

Similar to a relational database schema that describes tables, columns, and relationships of a relational database, the property graph schema is critical for creating high-quality domain-specific graphs. A property graph instantiated from a property graph schema provides agile and robust knowledge services with correctness, coverage, and freshness [38].

A property graph schema  $\mathcal{PGS}$  can be specified in a data definition language such as Neo4j’s Cypher [25], Tiger-Graph’s GSQL [21], or GraphQL SDL [27]. They all define notions of node types and edge types, as well as property types that are associated with a node type or with an edge type. We adopt Cypher due to its popularity, but our proposed techniques are independent of the aforementioned languages. Table 1 provides the notations used in this paper.

### 2.2 Approach Overview

Given an ontology  $O$  providing a semantic abstraction of the input data, the problem of *property graph schema optimization* is to generate a property graph schema that produces

<sup>9</sup>Even if inheritance and union are not ObjectProperties, we simplify the notation for presentation purposes.

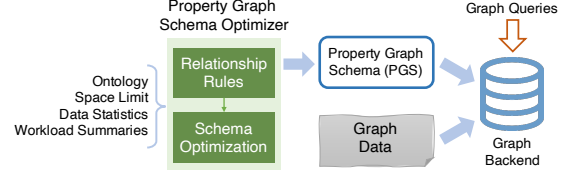
**Table 1: Notations.**

Notations	Definitions
$O$	an ontology
$c_i$	$c_i \in C$ : a concept in an ontology
$r_i$	$r_i \in R$ : a relationship in an ontology
$c_i.P_i$	all data properties associated to $c_i$
$c_i.inE$	all incoming relationships of $c_i$
$c_i.outE$	all outgoing relationships of $c_i$
$c_i.R_i$	$c_i.R_i = c_i.inE \cup c_i.outE$
$r_i.src$	the source concept of $r_i$
$r_i.dst$	the destination concept of $r_i$
$r_i.type$	the relationship type of $r_i$ (i.e., 1:1, union, inheritance, 1:M, or M:N)
$\mathcal{PGS}$	a property graph schema
$vs_i$	$vs_i \in VS$ : a schema vertex
$vs_i.PS_i$	all property schema of $vs_i$
$es_i$	an edge schema defined in $\mathcal{PGS}$
$es_i.type$	the edge type of $e_i$
$\mathcal{PG}$	a property graph
$V_i$	$V_i \in V$ all instance vertices of $vs_i$
$v_{i,j}$	$v_{i,j} \in V_i$ , an instance vertex of $vs_i$
$v_{i,j}.p_k$	a property of $v_{i,j}$
$e_k$	$e_k = (v_{src}, v_{dst}) \in E, v_{src}, v_{dst} \in V$

the best query performance for various graph queries (e.g., pattern matching, path finding, or aggregation queries). Optimizing the property graph might entail data replication and hence increased memory footprint. In real knowledge graph applications, especially in a multi-tenant setting, there is a limit on the amount of memory that we can trade for query performance. Hence, any practical solution needs incorporate a space constraint while producing an optimized property graph schema.

Figure 3 provides an overview of our property graph schema optimization approach. The property graph schema optimizer takes as input an ontology and optionally a space limit, data statistics, as well as workload summaries<sup>10</sup>. It utilizes a set of rules designed for different types of relationships to produce an optimized property graph schema. The raw graph data is then loaded into a graph database (e.g., Neo4j or JanusGraph) conforming to the optimized schema. At query time, users can directly express graph queries against this instantiated property graph corresponding to the optimized schema.

<sup>10</sup> Access frequencies of concepts, relationships, and data properties in an ontology

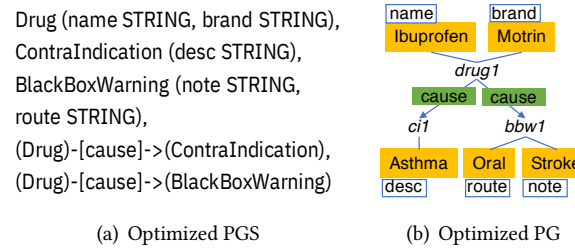


**Figure 3: Approach Overview.**

### 3 RELATIONSHIP RULES

Graph queries often involve multi-hop traversal or vertex attribute lookup/analytics on property graphs. As shown in the motivating examples, edge traversals over a graph are vital to the overall query performance. Hence, we focus on the rich semantic relationships in an ontology and propose a set of novel rules for different types of relationships. These rules minimize edge traversals and consequently improve graph query performance.

**Union Rule.** In an ontology, a union relationship ( $r_{un} = (c_i, c_j)$ ) contains a union concept ( $c_i$ ) and a member concept ( $c_j$ ). Each instance of a union concept is an instance of one of its member concepts, and each instance of a member concept is also an instance of the union concept. Figure 2 shows that *BlackBoxWarning* and *ContraIndication* are two member concepts of a union concept *Risk*. A graph query accessing an instance of *Risk* is equivalent to accessing the instances of either *BlackBoxWarning*, or *ContraIndication*, or both. In other words, if we create a property graph directly from the ontology shown in Figure 2, then the queries starting from any vertices of either *BlackBoxWarning* or *ContraIndication* concepts have to traverse through some vertex of *Risk* in order to reach the vertices of *Drug*. This leads to unnecessary edge traversal.



**Figure 4: Union Relationship.**

Hence we propose a union rule to alleviate this issue. The union rule first creates a union node  $vs_i$  (based on the corresponding  $c_i$  in  $O$ ) and its member node  $vs_j$  (based on the corresponding  $c_j$  in  $O$ ) in the property graph schema. Then the member node  $vs_j$  is connected to the other nodes that connect to the union node  $vs_i$  in the property graph schema

---

**Algorithm 1** Union Rule (union)

---

**Input:** A union relationships  $r_{un}$

```

1:  $vs_i \leftarrow r_{un}.src$  // the union concept of  $r_{un}$ 
2:  $vs_j \leftarrow r_{un}.dst$  // the member concept of  $r_{un}$ 
3: for each  $r \in vs_i.ES_i$  do
4:   if  $\neg(r \text{ of type union})$  then
5:      $vs_j.ES_j \leftarrow vs_j.ES_j \cup r$ 

```

---

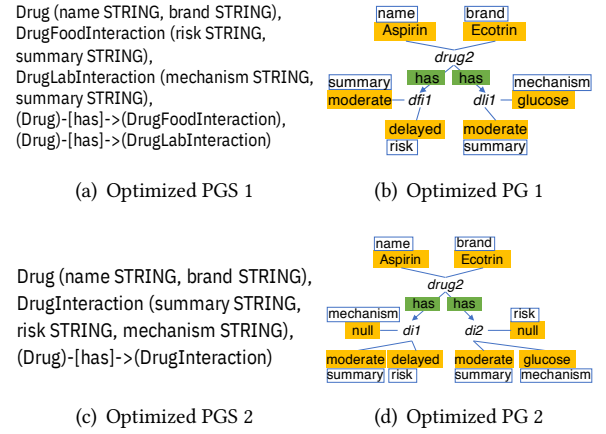
(Algorithm 1). Figures 4(a) and 4(b) show the property graph schema and the corresponding property graph after applying the union rule to the above example. In the optimized property graph, retrieving the drugs (e.g., *Ibuprofen*) that cause *Asthma* requires only a single edge traversal, instead of 2 in the property graph directly instantiated from the ontology.

**Inheritance Rule.** An inheritance relationship ( $r_{ih} = (c_i, c_j)$ ) contains a parent concept ( $c_i$ ) and a child concept ( $c_j$ ). Similar to the union rule, we create a parent node  $vs_i$  (corresponding to  $c_i$ ) and its child node  $vs_j$  (corresponding to  $c_j$ ) in the property graph schema. Unlike a union concept, a parent concept in the inheritance relationship may have instances that are not present in any of its children concepts.

- (1) Connect the child node  $vs_j$  directly to the nodes that are connected to its parent node  $vs_i$ , and attach all data properties  $vs_i.P_i$  of  $vs_i$  to the child node  $vs_j$  in the schema;
- (2) Connect the parent node directly to the nodes that are connected to its child node, and attach all data properties  $vs_j.P_j$  of  $vs_j$  to the parent node  $vs_i$  in the schema;
- (3) Or connect the parent  $vs_i$  and child  $vs_j$  nodes with an edge of type *isA*.

In the first two cases, edge traversals can be avoided in the property graph conforming to the property graph schema. Figure 2 shows that *DrugFoodInteraction* and *DrugLabInteraction* are two children concepts of *DrugInteraction*. Applying the inheritance rule to these concepts can lead to two alternative optimized property graph schemas shown in Figure 5. Figures 5(a) and 5(b) demonstrate the first scenario where the data properties (*summary*) of the parent concept *DrugInteraction* are directly attached to two children concepts *DrugFoodInteraction* and *DrugLabInteraction*. Figures 5(c) and 5(d) depict the second scenario where the data properties *risk* and *mechanism* of two respective children concepts are now attached to the parent concept *DrugInteraction*.

However, attaching the data properties ( $c_i.P_i$ ) from the parent concept to the child concept incurs data replication as  $c_i.P_i$  is shared among all children concepts (Figure 5(b)). If the number of data properties shared by the children concepts is large, the data replication can introduce significant space overhead. On the other hand, when the data properties ( $c_j.P_j$ )



**Figure 5: Inheritance Relationship.**

from the children concepts are replicated to their parent concept ( $c_i$ ),  $c_i$  may end up with a large number of data properties (Figure 5(d)). However, these data properties may not exist in many instance vertices of  $c_i$ . Consequently, the instance vertices of  $c_i$  may consume unnecessary space. To remedy the above two issues, we propose to exploit the Jaccard similarity [32] between  $c_i.P_i$  and  $c_j.P_j$  to decide the best strategy for the inheritance relationship:

$$JS(c_i.P_i, c_j.P_j) = |c_i.P_i \cap c_j.P_j| / |c_i.P_i \cup c_j.P_j|. \quad (1)$$

As described in Algorithm 2, if  $JS(c_i.P_i, c_j.P_j)$  is greater than a threshold  $\theta_1$ , it indicates that the child concept  $c_j$  shares a lot of data properties with its parent concept  $c_i$ . Intuitively, this means that  $c_j$  has only few properties in addition to the ones of  $c_i$ . In this case, moving  $c_j.P_j$  from the child concept to  $c_i$  incurs less space overhead compared to the other way. Similarly, if  $JS(c_i.P_i, c_j.P_j)$  is less than a threshold  $\theta_2$  ( $\theta_2 \leq \theta_1$ ), the child concept  $c_j$  has little in common with its parent  $c_i$ . Intuitively, this means that  $c_j$  has many additional properties compared to  $c_i$ . Therefore, it is more cost effective to make the data properties of the parent concept  $c_i.P_i$  available at  $c_j$ . In either case, the inheritance rule avoids edge traversals in the resulting property graph.

Note that the similarity score of a parent concept and a child concept remains unchanged even if new data properties are added to one or both concepts as a result of applying other rules. The reason is that the Jaccard similarity is computed based on the given ontology, as it represents the semantic similarity between two concepts with an inheritance relationship. Hence we calculate the Jaccard similarity score for all inheritance relationships before applying any rules.

**One-to-one Rule.** A 1:1 relationship ( $r_{1:1} = (c_i, c_j)$ ) indicates that an instance of  $c_i$  can only relate to one instance of

---

**Algorithm 2** Inheritance Rule (inheritance)

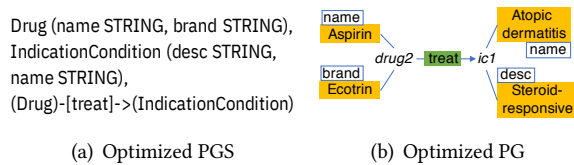
---

**Input:** An inheritance relationship  $r_{ih}$

```
1:  $vs_i \leftarrow r_{ih}.src$  // Parent concept
2:  $vs_j \leftarrow r_{ih}.dst$  // Child concept
3:  $jsim \leftarrow JS(vs_i.PS_i, vs_j.PS_j)$  // Jaccard similarity of  $r_{ih}$ 
4: if  $jsim > \theta_1$  then
5:    $vs_i.P_i \leftarrow vs_i.PS_i \cup vs_j.PS_j$ 
   //  $ES_{ih}$  is the set of inheritance relationships
6:    $vs_i.ES_i \leftarrow (vs_i.ES_i \cup vs_j.ES_j) \setminus r_{ih}$ 
7: else if  $jsim < \theta_2$  then
8:    $vs_j.PS_j \leftarrow vs_j.PS_j \cup vs_i.PS_i$ 
9:    $vs_j.ES_j \leftarrow (vs_j.ES_j \cup vs_i.ES_i) \setminus r_{ih}$ 
```

---

$c_j$  and vice versa (e.g., *Indication* and *Condition* in Figure 2). Two concepts ( $c_i$  and  $c_j$ ) of a 1:1 relationship can be represented as one combined node  $vs_{ij}$  in the optimized schema, which is similar to joining two tables in relational databases where one row in one table is linked with only one row in another table and vice versa. If two tables are merged, then a join can be saved when two tables are queried together. The other tables can still join with the merged table through their respective relationships. Namely, the 1:1 rule preserves the original semantics and does not lead to any information loss. Any query accessing instance vertices of  $c_i$  and  $c_j$  can be satisfied by looking up the merged instance vertex of  $c_{ij}$ . In Figure 6(a), *IndicationCondition* is the merged concept with two data properties, *name* and *note*, attached. Hence the edge traversal (e.g., from *Drug* to *Condition* in Figure 2) is avoided and the number of instance vertices (i.e., space consumption) is reduced as well. Algorithm 3 shows the one-to-one rule, which is straightforward to follow.



**Figure 6: 1:1 Relationship.**

---

**Algorithm 3** 1:1 Rule (oneToOne)

---

**Input:** A 1:1 relationship  $r_{1:1}$

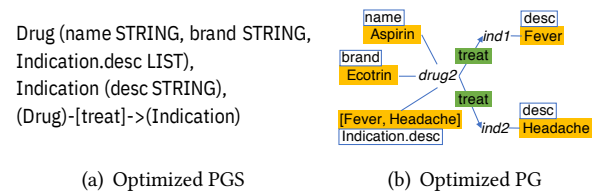
```
1:  $vs_i \leftarrow r_{1:1}.src$ 
2:  $vs_j \leftarrow r_{1:1}.dst$ 
3:  $vs_{i,j} \leftarrow \emptyset$ 
4:  $vs_{i,j}.ES_{i,j} \leftarrow (vs_i.ES_i \cup vs_j.ES_j) \setminus r_{1:1}$ 
5:  $vs_{i,j}.PS_{i,j} \leftarrow vs_i.PS_i \cup vs_j.PS_j$ 
```

---

**One-to-many Rule.** A 1:M relationship ( $r_{1:M} = (c_i, c_j)$ ) indicates that an instance of  $c_i$  can potentially refer to several instances of  $c_j$ . In other words, in a 1:M relationship, an instance of  $c_i$  allows zero, one, or many corresponding instances of  $c_j$ . However, an instance of  $c_j$  cannot have more than one corresponding instance of  $c_i$ .

To better support the aggregation (e.g., COUNT, SUM, AVG, etc.) and neighborhood (1-hop) lookup functions in graph queries, we first create two nodes  $vs_i$  and  $vs_j$  corresponding to  $c_i$  and  $c_j$  in the optimized schema. Then we propagate each data property  $vs_j.P_j$  of  $vs_j$  as a property of type *LIST* to the other node  $vs_i$  (Fig. 7(a)). The aggregation and neighborhood lookup functions can directly leverage these localized list properties instead of traversing through the edges of the 1:M relationships. This is similar to denormalization technique in relational databases where data replication is added to one or more tables in order to avoid costly joins. As depicted in Figure 7(b), *Indication.desc* is a data property of *drug2* consisting of a list of descriptions (i.e., [*Fever*, *Headache*]) that saves the aggregation queries edge traversals to the other instance vertices (e.g., *ind1* and *ind2*). The potential savings can be substantial when there are many edges between instance vertices of two concepts such as *Drug* and *Indication*.

However, the newly introduced property of type *LIST* introduces additional space overheads, which can be expensive depending on the data distribution. Therefore, choosing the appropriate set of data properties from each 1:M relationship to propagate is critical with respect to both query performance and space consumption. We will describe algorithms to choose the data properties to merge in Section 4.2. Algorithm 4 corresponds to the one-to-many rule.



**Figure 7: 1:M Relationship.**

---

**Algorithm 4** 1:M Rule (oneToMany)

---

**Input:** A 1:M relationship  $r_{1:M}$

```
1:  $vs_i \leftarrow r_{1:M}.src$ 
2:  $vs_j \leftarrow r_{1:M}.dst$ 
3: for each  $p \in vs_j.PS_j$  do
4:    $vs_i.PS_i.addAsList(p)$ 
```

---

**Many-to-many Rule.** An  $M:N$  relationship ( $r_{M:N} = (c_i, c_j)$ ) indicates that an instance of  $c_i$  can have several corresponding instances of  $c_j$ , and vice versa. An  $M:N$  relationship is essentially equivalent to two  $1:M$  relationships, namely,  $r_{1:M} = (c_i, c_j)$  and  $r_{1:M} = (c_j, c_i)$ . Therefore, the many-to-many rule is identical to the one-to-many rule, except that the property propagation is done for both directions. Namely, in the optimized schema, a data property of the node  $vs_i$  corresponding to  $c_i$  in  $O$  is propagated as a property of type *LIST* to the node  $vs_j$  corresponding to  $c_j$  in  $O$ , and vice versa. Hence applying the many-to-many rule leads to the same potential gains for queries with aggregate or neighborhood (1-hop) lookup functions at the cost of introducing additional space consumption.

In summary, all proposed rules reduce the number of edge traversals which improve graph query performance. Moreover, these rules can be utilized in graph systems using different storage backends. The potential benefits could be more significant when the storage backend changes from in-memory to disk as edge traversals may incur additional disk I/Os. However, *union*, *inheritance*, *one-to-many*, and *many-to-many* rules may incur space overheads. In Section 4, we describe our property graph schema optimization algorithms, trading off performance gain and space overhead.

## 4 PROPERTY GRAPH SCHEMA OPTIMIZATION

In this section, we first introduce a property graph schema optimization algorithm in an ideal scenario (i.e., no space constraints). Then, we describe our concept-centric and relation-centric algorithms that harness the proposed rules and a cost-benefit model to generate an optimized property graph schema for a given space constraint.

### 4.1 Optimization Without Space Constraints

To produce an optimized property graph schema, we need to determine how to utilize the proposed rules described in Section 3. A straightforward approach is to iteratively apply these rules in order and generate the property graph schema.

Specifically, Algorithm 5 takes as input an ontology  $O$  and first computes the Jaccard similarity scores for all inheritance relationships (Lines 1-2). Then, it iteratively applies the appropriate rule to each relationship in the ontology (Lines 3-16). At the end of each iteration, it checks if the ontology converges (Line 17). Finally when no more rule applies, a property graph schema is generated (Lines 18-19). In fact, these rules can be applied in any order, and the generated property graph schema is always the same.

---

#### Algorithm 5 Ontology to PGS without Space Limits

---

**Input:** Ontology  $O = (C, R, P)$

**Output:** A property graph schema  $\mathcal{PGS}$

```

// Compute Jaccard similarity for each inheritance relationship and get all parent concepts
1: for each  $r \in R$  of type inheritance do
2:    $r.js \leftarrow \text{computeJS}(r)$ 
3: repeat
4:    $O_{prev} \leftarrow O$ 
5:   for each  $r \in R$  do
6:     switch  $r.type$  do
7:       case 1:1
8:          $O \leftarrow \text{oneToOne}(O, r)$ 
9:       case 1:M
10:         $O \leftarrow \text{oneToMany}(O, r)$ 
11:      case M:N
12:         $O \leftarrow \text{manyToMany}(O, r)$ 
13:      case union
14:         $O \leftarrow \text{union}(O, r)$ 
15:      case inheritance
16:         $O \leftarrow \text{inheritance}(O, r)$ 
17: until  $O = O_{prev}$ 
18:  $\mathcal{PGS} \leftarrow \text{generatePGS}(O)$ 
19: return  $\mathcal{PGS}$ 

```

---

**THEOREM 3.** *Applying the union, inheritance, 1:M and M:N rules in any order produces a unique  $\mathcal{PGS}$ , if there is no space constraint.*

The detailed proof can be found in Appendix A.

### 4.2 Schema Optimization With Space Constraints

While the naïve approach harnesses all potential optimization opportunities aggressively, it incurs space overheads from *union*, *inheritance*, *1:M*, and *M:N* rules. In cases where the number of such relationships is large in the ontology, this can be expensive with respect to the space consumption, especially in a cluster setting, where many large-scale property graphs co-exist. Hence our goal is to produce an optimized property graph schema for a given space limit. The quality and the space consumption of an optimized property graph schema are measured based on the total benefit and cost (i.e., space consumed) by applying the rules (given by Equations 3-5 in Section 4.2.2).

**Definition 4 (Optimal Property Graph Schema).** Let  $\mathcal{PGS}$  be the set of all property graph schemas, such that  $\forall \mathcal{PGS}' \in \mathcal{PGS}$  we have  $\text{Cost}(\mathcal{PGS}') \leq S$ , where  $S$  is a given space budget.  $\mathcal{PGS}_{opt} \in \mathcal{PGS}$  is an optimal property graph schema if  $\nexists \mathcal{PGS}' \in \mathcal{PGS}$  such that  $\text{Benefit}(\mathcal{PGS}') > \text{Benefit}(\mathcal{PGS}_{opt})$ .

Finding an optimal property graph schema is exponential in the number of concepts and relationships in the ontology, which is practically infeasible. Hence, we need to design efficient heuristics to produce a near-optimal property graph schema. To achieve this goal, we propose two property graph schema optimization algorithms that leverage additional information such as data and workload characteristics.

**Data characteristics** contain the basic statistics about each concept, data property, and relationship specified in the given ontology. The statistics include the cardinality of data instances of each concept and relationship, as well as the data type of each data property. The data characteristics allow us to identify and prioritize the more beneficial relationships when applying *union*, *inheritance*, *one-to-many* and *many-to-many* rules, such that the space can be used more efficiently.

**Access frequencies** provide an abstraction of the workload in terms of how each concept, relationship, and data property accessed by each query in the workload. We use  $AF(c_i \xrightarrow{r_k} c_j.P_j)$  to indicate the frequency of queries (the number of queries) that access a data property in  $c_j.P_j$  from the concept  $c_i$  through the relationship  $r_k$ . The high frequency of a relationship indicates its relative importance among all relationships in the given ontology. Hence it is more imperative to apply the above proposed rules to these relationships with high frequency.

In case of no prior knowledge about access frequency, we assume that it follows a uniform distribution. Our approach can also handle updates (i.e., insert, delete, and modify) to the property graph if they do not incur any schema changes. If the accumulated updates change the data distributions, then we can apply the rules locally to the affected part of the ontology. Note that data statistics changes can invalidate certain rule applied earlier, or can trigger new rules, especially inheritance and union rules. We can make local adjustments to accommodate these changes. Minimizing such transformation overheads is left as future work.

**4.2.1 Concept-Centric Algorithm.** As described in Section 2, an ontology describes a particular domain and provides a concept-centric view over domain-specific data. Intuitively, some concepts are more critical to the domain, and have more relationships with the other concepts [40]. We expect these key concepts to be queried more frequently than other concepts, which is confirmed in [39]. This leads to our concept-centric algorithm that exploits the structural information in an ontology to identify key concepts which we believe are more likely to be accessed more often. Hence, this algorithm is useful when no workload summary is available.

To determine these key concepts, we utilize centrality analysis over the ontology to rank all concepts according to their respective centrality score. The centrality analysis is based on the commonly used PageRank algorithm [14] as

its underlying assumption, more important websites likely to receive more links from other websites, is similar to our intuition of key concepts. Our modified PageRank algorithm, called *OntologyPR* (Algorithm 6), determines the centrality score of each concept in an ontology. Compared to PageRank, we further introduce weights for both in and out degrees of concepts in determining their centrality scores.

---

**Algorithm 6** Ontology PageRank Algorithm  
(OntologyPR)

---

**Input:**  $O = (C, R, P)$

**Output:**  $O = (C, R, P)$

```

1:  $C_{un} \leftarrow$  empty set
2: for each  $r \in R$  of type union do
3:    $c_i \leftarrow r.src$  // the union concept of  $r$ 
4:    $c_j \leftarrow r.dst$  // the member concept of  $r$ 
5:    $C_{un}.add(c_i)$ 
6:    $c_j.R_j \leftarrow (c_j.R_j \cup c_i.R_i) \setminus r$ 
7:  $O.remove(C_{un})$ 
8: for each  $r \in R$  do
9:   if  $r$  is of type inheritance then
10:     $R_{ih}.add(r)$ 
11:     $O.remove(r)$ 
12:   else
13:     $O.add(r')$  // add a reverse relation  $r'$ 
14: pageRank( $O$ ) // PageRank on the modified ontology
15:  $O.add(R_{ih})$  // add inheritance relationships back
16: updatePR( $O$ ) // update PageRank score for inheritance concepts
17: return  $O$  //  $O$  associated with PageRank scores

```

---

**Inheritance.** To cater for inheritance relationships, we remove these relationships from the ontology while running the initial PageRank algorithm. This allows us to calculate the page ranks of a concept based on the links from other concepts that are not children of the same concept. After computing the page rank values of all concepts, we re-attach these relationships and update the page ranks of each concept by doing a depth-first traversal over its inheritance relationships to find the parent with the highest page rank. If this value is higher than the current page rank of the concept, we use this value as the new page rank of the concept. This enables a child concept to inherit the page rank of its parent. The intuition is that a child concept inherits all its other properties from the same chain of concepts and hence would have a similar estimate of centrality.

**Unions.** The union concept in the ontology represents a logical membership of two or more concepts. Any incoming edge to a union concept can therefore be considered as pointing to at least one of the member concepts of the

union. Similarly each outgoing edge can be considered as emanating from at least one of the member concepts.

To handle union concepts, the *OntologyPR* algorithm iterates over all incoming and outgoing edges to/from the union concept. For each incoming edge to the union concept, we create new edges between the source concept and each of the member concepts of the union. For each outgoing edge, similarly, we create new edges between the destination and each of the member concepts of the union. Thus the page rank mass is appropriately distributed to/from the member nodes of the union. Finally, the union node itself is removed from the graph as its contribution towards centrality analysis has already been accounted for by the new edges to/from the member concepts of the union.

**Out-degree of Concepts.** In the default PageRank algorithm, the weight distribution of the page rank is proportional to the in-degree of a node as it receives page rank values from all its neighbors that point to it. In other words nodes with a high in-degree would tend to have a higher page rank than nodes with a low in-degree. However, for a domain ontology, we observe that both in-degree and out-degree are equally important in terms of the key concept. Hence, we introduce a reverse edge in the ontology, essentially making the graph equivalent to an undirected graph. Then, the *OntologyPR* algorithm uses this modified ontology as an input to determine the centrality score of each concept.

Using *OntologyPR*, we associate PageRank scores with each concept in the ontology. To accurately capture the relative importance of the concepts, we further leverage the *data characteristics* and *access frequency* information to rank all concepts. Namely, the ranking score for a concept is defined as follows.

$$Score(c_i) = \frac{c_i.pr \times AF(c_i)}{Size(c_i)} \quad (2)$$

where  $c_i.pr$  denotes the PageRank score of  $c_i$ ,  $AF(c_i)$  denotes the access frequency of  $c_i$  including accessing all data properties of  $c_i$ , and  $Size(c_i)$  denotes the size of  $c_i$  including all data properties of  $c_i$ .

Based on Equation 2, our concept-centric algorithm (Algorithm 7) first sorts all concepts in a descending order of their respective scores (Lines 1-2). Then, it iterates through each concept  $c$  (Lines 3-8). For each concept, the algorithm utilizes the *applyRules* procedure to apply all rules (Section 3) to the relationships connecting to  $c$ . During this process, the algorithm updates the space limit as it is consumed by the rules. Once the space is fully exhausted, the algorithm terminates (Lines 7-8) and returns the optimized property graph schema (Line 10).

**Complexity Analysis.** The *OntologyPR* is the dominant procedure in Algorithm 7, and its time complexity is  $O((|R| + |C|)k)$ , where  $|R|$  is the number of relationships,  $|C|$  is the

---

#### Algorithm 7 Concept-Centric Algorithm

---

**Input:** Ontology  $O = (C, R, P)$ , space limit  $S$

**Output:** A property graph schema  $\mathcal{PGS}$

```

1:  $O \leftarrow \text{ontologyPR}(O)$ 
2:  $C_{srt} \leftarrow \text{sort}(C)$ 
3: for each  $c \in C_{srt}$  do
4:   for each  $r \in c.R$  do
5:      $S' \leftarrow S$ 
6:      $O, S \leftarrow \text{applyRules}(r, S')$ 
7:     if  $S < 0$  then
8:       break
9:  $\mathcal{PGS} \leftarrow \text{generatePGS}(O)$ 
10: return  $\mathcal{PGS}$ 
```

---

number of concepts, and  $k$  is the maximum number of iterations. The time complexity of sorting concepts is  $O(|C|\log|C|)$ . Finally, the time complexity of applying rules to the sorted concepts is  $O(|R|)$ . Thus, the overall time complexity of Algorithm 7 is  $O((|R| + |C|)k)$ .

**4.2.2 Relation-Centric Algorithm.** Intuitively, the concept-centric algorithm prioritizes the relationships of the key concepts in an ontology by leveraging information such as access frequency, data characteristics, and structural information from the ontology. However, the relationship selection is limited to each concept locally. Namely, the concept-centric algorithm does not have a global optimal ordering among all relationships in the ontology. To address this issue, we propose the relation-centric algorithm based on a cost-benefit model for each type of relationships described as follows.

**Cost Benefit Models.** The union rule, introduced in Section 3, connects the member concept directly to all concepts that are connected to the union concept. Then, the benefit of applying this rule to a union relationship  $r$  is the access frequency of  $r$ , and the cost is the number of edges that we copy from the union concept to the member concept. Formally:

$$\begin{aligned} \text{Benefit}(r) &= AF(c_i \xrightarrow{r} c_j) \\ \text{Cost}(r) &= \sum_{r' \in (c_i.R_i \setminus R_{un})} |r'|, \end{aligned} \quad (3)$$

where  $c_i$  denotes the union concept and  $|r'|$  denotes the number of edges between the instance vertices of  $c_i$  and the ones of a neighborhood concept<sup>11</sup> of  $c_i$ .

The benefit of applying the inheritance rule to an inheritance relationship is the access frequency of that relationship multiplied by the Jaccard similarity between  $c_i.P_i$  and  $c_j.P_j$ . Depending on that similarity, the cost of inheritance rule can be either the number of new edges attached to the parent,

---

<sup>11</sup>The neighborhood concepts do not include the member concepts of  $c_i$ .

or the number of new edges attached to the child. Formally:

$$\begin{aligned} \text{Benefit}(r) &= AF(c_i \xrightarrow{r} c_j.P_j) \times JS(c_i, c_j) \\ \text{Cost}(r) &= \begin{cases} \sum_{p \in c_j.P_j} |c_j| \times p.type + \sum_{r \in (c_j.R_j \setminus R_{ih})} |r|, & \text{if } \theta_1 < JS(c_i, c_j) \\ \sum_{p \in c_i.P_i} |c_i| \times p.type + \sum_{r \in (c_i.R_i \setminus R_{ih})} |r|, & \text{if } JS(c_i, c_j) < \theta_2, \end{cases} \end{aligned} \quad (4)$$

where  $JS(c_i, c_j)$  denotes the Jaccard similarity between  $c_i.P_i$  and  $c_j.P_j$ ,  $p.type$  indicates the data type size of  $p$  (e.g., the size of INT, DOUBLE, STRING, etc.),  $\sum_{p \in c_j.P_j} |c_j| \times p.type$  ( $\sum_{p \in c_i.P_i} |c_i| \times p.type$ ) denotes the space overheads incurred by propagating  $c_j.P_j$  ( $c_i.P_i$ ) to  $c_i$  ( $c_j$ ), and  $\sum_{r \in (c_j.R_j \setminus R_{ih})} |r|$  ( $\sum_{r \in (c_i.R_i \setminus R_{ih})} |r|$ ) denotes the space overhead incurred by connecting the neighbors of  $c_i$  ( $c_j$ ) to  $c_j$  ( $c_i$ ).

Similarly, the cost-benefit model for one-to-many rule, leveraging both data characteristics and access frequency information, is described as:

$$\begin{aligned} \text{Benefit}(r) &= AF(c_i \xrightarrow{r} c_j.p) \\ \text{Cost}(r) &= |r| \times p.type, \end{aligned} \quad (5)$$

where  $|r| \times p.type$  denotes the space overhead incurred by replicating  $p$  as a data property of type *LIST* to  $c_i$ .

As described in Section 3, each  $M:N$  relationship is equivalent to two  $1:M$  relationships. Thus, we first convert each  $M:N$  relationship in the ontology into two  $1:M$  relationships, and then use Equation 5 to decide the cost-benefit for each of them. Potentially some of the original  $M:N$  relationships could be optimized for only one direction. This increases the flexibility of applying many-to-many rule such that more frequently accessed data properties can be propagated to the other end of the relationship.

With the cost and benefit scores, our goal is to select a subset of relationships in the ontology that maximize the total benefit within the given space limit. We map our **relationship selection problem** to the **0/1 Knapsack Problem**, which is NP-hard [50].

**PROPOSITION 1 (REDUCTION).** *If both benefit and cost of a relationship are positive, then every instance of the relationship selection problem can be reduced to a valid instance of the 0/1 Knapsack problem.*

The proof of Proposition 1 can be found in Appendix B.

We adopt the fully polynomial time approximation scheme (FPTAS) [50] for our relation selection problem, which guarantees that the benefit of the optimized property graph schema  $\text{Benefit}(\mathcal{PGS})$  is within  $1-\epsilon$  ( $\epsilon > 0$ ) bound to the benefit of the optimal property graph schema  $\text{Benefit}(\mathcal{PGS}_{opt})$ .

Algorithm 8 takes as inputs an ontology and the space limit. Similar to Algorithm 5, it computes the Jaccard similarity scores for all inheritance relationships (Lines 1-2). Then it computes the cost and benefit for each relationship in the

---

#### Algorithm 8 Relation-Centric Algorithm

---

**Input:**  $O = (C, R, P)$ , space limit  $S$

**Output:** A property graph schema  $\mathcal{PGS}$

```
// Compute Jaccard similarity for each inheritance relationship and get all parent concepts
1: for each  $r \in R$  of type inheritance do
2:    $r.js \leftarrow \text{computeJS}(r)$ 
3:  $\text{Benefit}, \text{Cost} \leftarrow \emptyset$ 
4: for each  $r_i \in R$  do
5:    $\text{Benefit}[i] \leftarrow \text{Benefit}(r_i)$ 
6:    $\text{Cost}[i] \leftarrow \text{Cost}(r_i)$ 
7:  $R_{opt} \leftarrow \text{knapsack}(R, \text{Benefit}, \text{Cost}, S)$ 
8: for each  $r_i \in R_{opt}$  do
9:    $O \leftarrow \text{applyRules}(r_i)$ 
10:  $\mathcal{PGS} \leftarrow \text{generatePGS}(O)$ 
11: return  $\mathcal{PGS}$ 
```

---

ontology  $O$  using Equations 3, 4, and 5 (Lines 3-6). Next, the FPTAS algorithm is used to select the near-optimal subset of relationships  $R_{opt}$  with the given space limit  $S$  (Line 7). In *applyRules* procedure, the algorithm applies the corresponding rules;  $r \in R_{opt}$  (Lines 8-9). Lastly, an optimized property graph schema is generated (Lines 10-11).

**Complexity Analysis.** The FPTAS *knapsack* is the dominant procedure in Algorithm 8, and its time complexity is  $O(|R|^2 \lfloor |R|/\epsilon \rfloor)$  [50], where  $|R|$  is the number of relationships and  $\epsilon \in (0, 1]$ . The rest of Algorithm 8 is linear to  $|R|$ . Thus, the time complexity of Algorithm 8 is  $O(|R|^3)$ .

## 5 EXPERIMENTAL EVALUATION

In this section, we present experiments to evaluate the effectiveness of our property graph schema design algorithms, and compare the query performance of different property graphs generated by different algorithms.

### 5.1 Experimental Setup

**Infrastructure.** We implemented our approach in Java with JDK 1.8.0 running on Ubuntu 14.04 with 16-core 3.4 GHz CPU and 128 GB of RAM. We choose two popular graph database systems, Neo4j [6] and JanusGraph [5], as our graph backends. We executed each experiment ten times and here we report their average.

**Data sets.** To evaluate the effectiveness of our system on different application domains, we use the following two data sets and their corresponding ontologies.

1. Financial data set (*FIN*) [43] includes data from two main sources: Securities and Exchange Commission (SEC) [8] and Federal Deposit Insurance Corporation [3]. The size of the data set is approximately 53 GB. The corresponding financial ontology contains 28 concepts, 96 properties, and 138

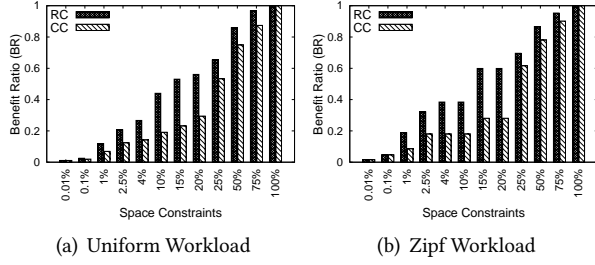


Figure 8: Varying Space Constraints (MED).

relationships (4 union, 69 inheritance, and 30 one-to-many relationships). It contains financial entities, financial metrics, lender, borrower, investment relationships, the officers of the companies as well as their relationships, etc.

2. Medical data set (*MED*) contains medical knowledge that is used to support evidence-based clinical decision and patient education. The total size of this data set is around 12 GB. The corresponding medical ontology consists of 43 concepts, 78 properties, and 58 relationships (11 inheritance, 5 one-to-one, 30 one-to-many, and 12 many-to-many relationships).

**Methodology and metrics.** To evaluate the quality of the property graph schemas produced by our algorithms, we vary the space limit and the Jaccard similarity thresholds for inheritance relationships with two different workload summaries (uniform and Zipf). Specifically, we show how effectively *PGSG* leverages the given space limit, how robust *PGSG* is to various workloads, and how sensitive *PGSG* is to different similarity thresholds. *PGSG* chooses the property graph schema with a higher total benefit score from relation-centric (*RC*) and concept-centric (*CC*) algorithms. We measure the quality of a property graph schema by  $BR = \frac{B_{SC}}{B_{NSC}}$ , where  $B_{NSC}$  is the total benefit score of the property graph schema generated by Algorithm 5 without any space constraint, and  $B_{SC}$  indicates the total benefit score achieved by either *RC* or *CC* algorithm.

To verify the graph query performance, we express most graph queries in both Cypher [25] and Gremlin [4], including path, reachability, and graph analytical queries. Among these query types, we construct a variety of query workloads conforming to different workload distributions over both financial and medical data sets. The details of these query workloads are described in Section 5.3. We use latency as the metric to measure these graph queries. Latency is measured in milliseconds as the total time of all queries in a workload executed in sequential order. We also use the number of edge traversals required in a query as the second metric. It directly reveals the computational savings achieved by our optimized property graph schema. Lastly, we evaluate the efficiency

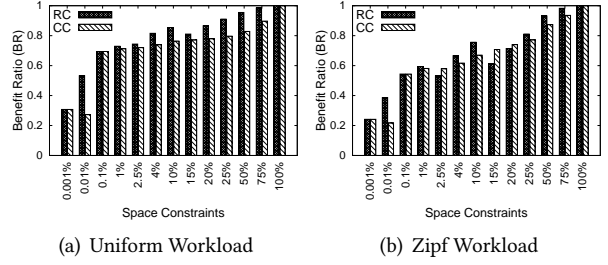


Figure 9: Varying Space Constraints (FIN).

of our concept-centric and relation-centric algorithms with different space constraints.

## 5.2 Property Graph Schema Quality

**Varying Space Constraint.** In Figures 8 and 9, we focus on the quality of the property graph schema produced by our concept-centric (*CC*) and relation-centric (*RC*) algorithms compared to our method without space constraints *NSC* (Algorithm 5). We choose two commonly seen workload summaries, uniform and Zipf distributions. The Zipf workload gives more access to the key concepts in the ontology. Namely, the access frequencies of concepts in the ontology follow either uniform or Zipf distribution. And the skew factor of Zipf distribution is set to 1. We first use *NSC* to produce an optimal property graph schema  $PGS_{NSC}$  without any space constraint, and then compute the total benefit score  $B_{NSC}$  achieved by  $PGS_{NSC}$  as well as the total amount of space  $S_{NSC}$  needed by  $PGS_{NSC}$ . We also compute the total amount of space  $S_{DIR}$  needed by the direct mapping algorithm from the given ontology. The space used by *NSC* is approximately 29GB for *MED* and 106GB for *FIN*, respectively. The total amount of space needed by the direct mapping algorithm  $S_{DIR}$  is 12GB for *MED* and 53GB for *FIN*, respectively. We, then, vary the space constraint from  $S_{DIR}$  to  $S_{NSC}$ , such that the range of the Y-axis in Figures 8 and 9 is from 0 to 1. Figures 8 and 9 show results from *MED* and *FIN* data sets respectively.

In Figure 8, we observe that *RC* consistently outperforms *CC* with both uniform and Zipf workloads. The reason is that *RC* has a global ordering of all relationships, and the global ordering is near-optimal with respect to the given space constraint due to the adopted approximate Knapsack algorithm. On the contrary, *CC* suffers from a rather local optimal ordering with respect to each concept. Hence, it misses the opportunity to utilize the space for more beneficial relationships. Moreover, we observe that with approximately 20% of the maximum space constraint, *RC* is able to produce high-quality property graph schemas which achieve above 50% of the total benefit. In other words, both algorithms can

effectively utilize the rather limited space. Lastly, both *RC* and *CC* produce the same property graph schema as  $PGS_{NSC}$  when the space constraint reaches 100%, which substantiates Theorem 3.

Similarly, *RC* outperforms *CC* In Figure 9, as *CC* utilizes the space for one concept at a time, missing the opportunities for more beneficial relationships in the ontology. We also observe that both algorithms, with uniform and Zipf workloads, have a couple of drops when the space constraint increases. The reason is primarily due to the complexity of *FIN* ontology. Given that the inheritance relationships are more dominant in *FIN*, the given space may be exhausted quickly by certain inheritance relationships. Again, *RC* and *CC* produce the same property graph schema as  $PGS_{NSC}$  with 100% space constraint.

**Varying Jaccard Similarity.** In Figure 10, we show the sensitivity of both *CC* and *RC* with respect to the Jaccard similarity thresholds ( $\theta_1$  and  $\theta_2$ ). In this experiment, we choose *FIN* ontology because it consists of multiple inheritance relationships. Uniform and Zipf workload distributions are used to examine the robustness of our *CC* and *RC* algorithms. Note that the space constraint in this experiment is set to  $(S_{NSC} - S_{DIR})/2$  under each specific Jaccard similarity threshold. The reason is that the cost (space overhead) of the same inheritance relationship can vary (Eq. 4) depending on the similarity threshold. Consequently, the space consumption of the optimal property graph changes under different thresholds. As shown in Figure 10, both *CC* and *RC* are robust under different similarity thresholds. In the worst case, they achieve more than 70% of the maximum benefit score under 50% space constraint. This shows that when the cost-benefit of an inheritance relationship changes due to a different threshold, both *CC* and *RC* can adjust accordingly by choosing different and more beneficial relationships to optimize. Hence, the total benefit scores achieved by both algorithms are relatively stable.

In summary, *CC* and *RC* produce high quality property graph schemas under various settings. They work effectively with any given space constraints. Moreover, *RC* always produces a near-optimal property graph schema and outperforms *CC* in most cases. Our property graph schema generator still leverages both algorithms to choose the property graph schema with the highest benefit score under any space constraints.

### 5.3 Graph Query Execution

In this section, we focus on the graph query execution performance over the property graphs created by our ontology-driven approach. We use both *MED* and *FIN* data sets to conduct our experiments. First, we create a micro benchmark to empirically examine whether the property graph

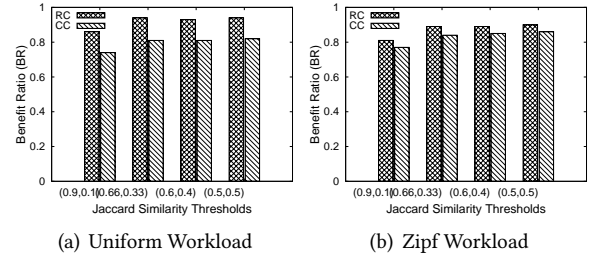


Figure 10: Varying Jaccard Thresholds (FIN).

schema from our approach can actually benefit a set of graph primitives including simple pattern matching, vertex property lookup, and aggregation on vertices. Second, we study the overall execution time for a given graph query workload by mixing the above graph primitives. We run the graph queries, expressed in Cypher and Gremlin, on Neo4j and JanusGraph, respectively. Note that our goal is not to compare the performance between two systems, rather to show that our schema optimization results in query performance improvements irrespective of the backend.

**Microbenchmark Using Graph Primitives.** With both *MED* and *FIN* data sets, we compare the query performance of the property graph created by the optimized graph schema (*OPT*) to the baseline property graph created by a direct mapping of the ontology (*DIR*). The following parameter settings are used to produce *OPT*: Jaccard similarity thresholds  $\theta_1 = 66\%$ ,  $\theta_2 = 33\%$ , and space constraint 0.5 ( $S_{NSC} - S_{DIR}$ ). All queries ( $Q_1$ - $Q_{12}$ ) are first expressed against *DIR* and then rewritten into the semantically equivalent queries over *OPT*. These queries are constructed according to the query patterns in [12]. We show a few representative queries used in the microbenchmark below.

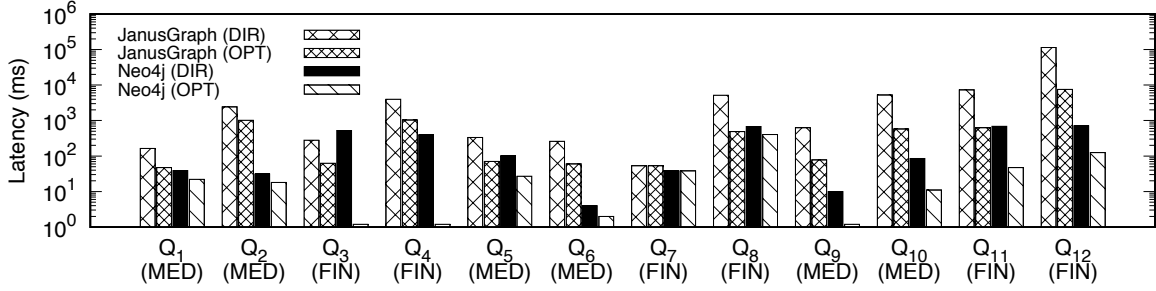
```
Q1: MATCH (d:Drug)-[p:cause]->(r:Risk)<-
      [p2:unionOf]-(ci:ContraIndication)
RETURN d.name
```

```
Q3: MATCH (aa:AutonomousAgent)<-[r1:isA]-
      (p:Person)<-[r2:isA]-(cp:ContractParty)
RETURN aa
```

```
Q5: MATCH (dl:DrugLabInteraction)-[r:isA]->
      (di:DrugInteraction)
RETURN di.summary
```

```
Q7: MATCH (n:Corporation)
RETURN n.hasLegalName
```

```
Q9: MATCH p=(d:Drug)-[r:hasDrugRoute]->
      (dr:DrugRoute)
```



**Figure 11: Microbenchmark - Pattern Matching ( $Q_1$ - $Q_4$ ), Property Lookup ( $Q_5$ - $Q_8$ ), Aggregation ( $Q_9$ - $Q_{12}$ ).**

```
RETURN dr.drugRouteId, size(COLLECT(
d.brand)) AS numberOfDrugBrands
```

```
Q11: MATCH p=(con:Contract)-[r:isManagedBy]->
(corp:Corporation)
RETURN size(COLLECT(con.hasEffectiveDate)) AS
numberOfEffectiveDates
```

As shown in Figure 11, the results are unequivocal. The optimized schema has significant advantages over the direct mapping schema for all types of queries. The graph pattern matching queries ( $Q_1$ - $Q_4$ ) report all matches of a sub-graph with 3 vertices and 2 edges in the property graph. Query execution times with our approach are at least 2.4 times faster than the direct mapping schema. The number of edge traversals on *DIR* is always 2 as the query is specified with 2 edges connecting 3 vertices. On the other hand, our property graph only requires at most 1 edge traversal as some of the neighbor vertices have been already merged with the starting vertices.

$Q_5$ - $Q_8$  are vertex property lookup queries. Both  $Q_5$  and  $Q_8$  are interested in a property of a vertex of a parent concept, and the starting vertex is a vertex of a child concept.  $Q_6$  starts from a vertex and looks for a property of its neighbor vertex. *OPT* has the property of type *List* with the starting vertex, and is able to return the result without any edge traversal.  $Q_7$  looks for a property of the starting vertex. In this case, *OPT* and *DIR* have identical query performance as no edge traversal is required. Hence *OPT* takes advantage of having the property of the parent concept available at the starting vertex, and consequently returns the result without any edge traversals. Therefore, the query runs more than an order of magnitude slower on the property graph of *DIR* than the one on *OPT* in the worst case.

$Q_9$ - $Q_{12}$  are graph aggregation queries that involve traversal from one vertex to the other. They count the number of neighbors of the starting vertex. On average, the query execution time is an order of magnitude faster for *OPT* approach compared to *DIR*. Again, the reason is that the aggregation

on the neighbor vertices can be instantaneously returned from the starting vertex. The above results suggest that using the proposed ontology-driven approach can bring significant benefits to a variety of graph queries.

Lastly, we observe that the performance gain on Neo4j is more substantial compared to JanusGraph (e.g.,  $Q_3$ ,  $Q_4$ ,  $Q_9$ , etc.). This shows that disk-based graph systems (e.g., Neo4j) benefits much more from our techniques, as the optimized schema requires significantly less disk I/O. Namely, the graph system loads less number of vertices and edges into memory. We expect such benefit to become even greater when the size of the property graph increases.

**Table 2: Microbenchmark - # Edge Traversals.**

Graph Queries	# Edge Traversals		Graph Queries	# Edge Traversals	
	DIR	OPT		DIR	OPT
$Q_1$	21,608	6,072	$Q_7$	0	0
$Q_2$	288,142	115,014	$Q_8$	493,588	0
$Q_3$	36,272	0	$Q_9$	67,397	0
$Q_4$	510,460	97,614	$Q_{10}$	429,636	15,327
$Q_5$	38,768	0	$Q_{11}$	524,265	0
$Q_6$	32,586	0	$Q_{12}$	110,4756	548,262

In addition, Table 2 reveals that *OPT* substantially reduces the number of edge traversals required in most queries, which leads to significant computational savings and performance gains. In several cases (e.g.,  $Q_3$ ,  $Q_6$ ), edge traversals can be completely avoided as the queried information is available locally within the starting vertices. On the other hand, the performance gains of certain queries (e.g.,  $Q_5$ ,  $Q_8$ ,  $Q_{12}$ ) are not as significant as others, even though the number of edge traversals with *OPT* is much smaller than the one with *DIR*. The reason is that the costs of lookup and return operations are non-trivial in both *DIR* and *OPT*, which can be observed from the latency of these queries in Fig. 11.

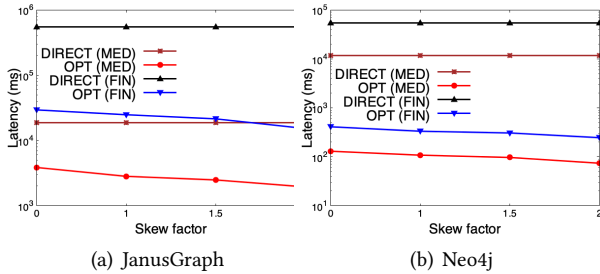
**Graph Query Workload Performance.** To evaluate the runtime performance of the property graph schema generated by our approach, we first generate two query workloads,

including both uniform and Zipf distributions in terms of the access frequency of the concepts in the ontology. We vary the Zipf’s skew factor from 0 (i.e., uniform distribution) to 2 (highly skewed). All query workloads consist of 15 queries of mixed types (i.e., pattern matching, lookup, and aggregation), similar to the ones used in the microbenchmark. The space limit is set to 20% of the space consumed by *NSC* (i.e., 15.4GB for *MED* and 80GB for *FIN*). The Jaccard similarity thresholds are  $\theta_1 = 66\%$  and  $\theta_2 = 33\%$ . The optimized schemas ( $OPT_{MED}$  and  $OPT_{FIN}$ ) are produced by the best performing algorithm of *RC* and *CC*.

**Table 3: Benefit Ratio w.r.t  $B_{NSC}$ .**

Skew Factor	MED				FIN			
	0	1	1.5	2	0	1	1.5	2
RC	56%	59%	62%	71%	67%	71%	74%	88%
CC	30%	43%	50%	63%	65%	74%	80%	88%

Table 3 shows the quality of the property graph schema produced by *RC* and *CC* compared to the one without space constraints *NSC*. We observe that both *RC* and *CC* correctly prioritize the most cost-effective relationships when the workloads are highly skewed. *RC* performs better than *CC* over *MED*, because *MED* has more data properties per concepts and *RC* makes more flexible decisions in terms of which relationships to optimize. On the other hand, *CC* performs better than *RC* over *FIN* as it successfully selects few concepts that are frequently accessed by the highly skewed workloads.



**Figure 12: Total Query Latency (MED & FIN).**

We compare our optimized schemas to the direct mapping schemas ( $DIRECT_{MED}$ ,  $DIRECT_{FIN}$ ) on both JanusGraph and Neo4j. The total query latency measures the performance on these property graphs corresponding to different schemas. Fig. 12 shows the total query latency in log scale. Both  $OPT_{MED}$  and  $OPT_{FIN}$  offer significant performance boosts to the graph query workloads on both JanusGraph and Neo4j. In Fig. 12(a), we observe that the total query latency on the optimized schema, on average, is around 7 and 26 times faster than the direct mapping one over *MED* and

*FIN*, respectively. The winning margin is substantially bigger (i.e., 129 and 176 times faster) on Neo4j (Fig. 12(b)). The total query latency on both optimized schema is approximately 2 orders of magnitude faster than the direct mapping. Moreover, we also observe that the total query latency decreases with increasing skew factor. Both  $OPT_{MED}$  and  $OPT_{FIN}$  achieve the lowest latency when the workload distributions are highly skewed. This indicates that the most frequently accessed concepts and relationships in the workloads are chosen to be optimized given the space limit. Based on these results, we verify that the designed rules for different types of relationships in the ontology are effective in terms of reducing edge traversals and consequently improving the graph query performance. Furthermore, we demonstrate that our approach can effectively utilize the given space constraint by leveraging data distribution and workload summaries.

## 5.4 Efficiency of Property Graph Schema Algorithms

Finally, we study the execution time of our concept-centric and relation-centric algorithms (Table 4). First, we observe that both *CC* and *RC* produce an optimized property graph schema in less than one second with different space constraints (shown in Table 4 as percentages of the space consumed by Algorithm 5). The optimization time of both algorithms is negligible compared to an exhaustive search approach, which even failed to produce an optimal schema for *MED* after 3 hours. Second, neither of the algorithms is sensitive to the space constraint, since both algorithms have a polynomial time complexity with respect to the number of concepts and relationships in the given ontology. Third, *RC* is consistently faster than *CC*, and the performance difference is more significant in *FIN*. This is due to the cost of *ontologyPR* procedure being dominant in *CC*. It usually takes more iterations to converge when the ontology (i.e., *FIN*) is more complex.

**Table 4: Efficiency of RC & CC (Time in ms).**

Space Constraint	MED			FIN		
	25%	50%	75%	25%	50%	75%
RC	23	23	26	192	188	193
CC	34	36	36	373	344	372

## 6 RELATED WORK

Schema optimization for improving query performance has been studied in the database community for decades [18, 24, 35, 51]. In recent years, the emergence of many large-scale knowledge graphs has drawn attention for schema optimization. In this section, we present important works in this field, highlighting the main differences to our approach.

**Schema Optimization in RDBMS/NoSQL.** Extensive work is available for schema design problem in relational database systems [10, 15, 20, 24, 30, 51]. RDBMSs provide a clean separation between logical and physical schemas. The logical schema includes a set of table definitions and determines a physical schema consisting of a set of base tables [10, 24, 51]. The physical layout of these base tables is then optimized with auxiliary data structures such as indexes and materialized views for the expected workload [10, 30]. Typically, the physical design often involves identifying candidate physical structures and selects a good subset of these candidates [20]. NoSE [35] is introduced to recommend schemas for NoSQL applications. Its cost-based approach utilizes a binary integer programming formulation to generate a schema based on the conceptual data model from the application.

In principle, our approach is similar to the logical schema design in RDBMSs, which defers the physical design to the underlying graph systems. Other than that, our approach is different from the above methods since the data modeling for graphs is inherently different from the relational data model. Specifically, the graph structure results in more expressive data models than those produced using relational databases, allowing the formation of graph queries (e.g., reachability, path finding, pattern matching) in a very intuitive fashion. Moreover, our approach exploits the rich semantic information available in an ontology to drive the schema optimization, which is not considered by any of the previous works.

**Schema Optimization in Knowledge Graphs.** In the last few years, RDF has been growing significantly for expressing graph data. A variety of schemas have been proposed for physically storing graph data in both centralized and distributed settings [9, 13, 16, 26, 29, 33, 34, 36, 37]. Some of these works focus on optimizing RDF data storage and SPARQL queries based on either workload statistics [33, 34, 36, 37] or heuristics [48]. A fundamental difference to those works is that we neither re-load the data to follow a new schema, nor build new indices, nor optimize the queries on-the-fly (e.g., join reordering). Instead, inspired by database literature as stated above, we provide an optimized property graph schema design before loading the data, and directly instantiate a property graph conforming to this schema on a graph database. Graph queries are then executed on the graph database, where graph query optimization techniques can be further utilized. Other works [9, 13, 16, 26] attempt to transform RDF data into relational data and provide SPARQL views over relational schemas, leveraging the many years of experience in RDBMS schema optimization. Unlike those approaches, we do not create views over a property graph. Instead, we directly express property graph queries in Cypher or Gremlin over optimized property graphs. Whether a graph

database internally uses views or not for query optimization is orthogonal to this work.

Recently, works such as [28, 45, 46] address a similar problem in the context of property graphs. GRFusion [28] focuses on filling the gap between the relational and the graph models rather than optimizing the graph schema to achieve better query performance. Szárnyas et al. [46] propose to use incremental view maintenance for property graph queries. However, their approach can only support a subset of property graph queries by using nested relational algebra. SQL-Graph [45] and Db2 Graph [47] introduce a physical schema design that combines relational storage for adjacency information with JSON storage for vertex and edge attributes. It also translates Gremlin queries into SQL queries in order to leverage relational query optimizers. However SQLGraph and Db2 Graph also focus on physical schema design which only targets on the relational databases. The query translator is limited to Gremlin queries with no side effects. Our ontology-driven approach is different for the following reasons. First, our approach produces a high-quality schema applicable to any graph system compatible with property graph model and Gremlin or Cypher queries. Second, we exploit the rich semantic information in an ontology to guide the schema design. Last but not least, our approach can further leverage these techniques to decide how the property graph should be stored on different storage backends.

Materialized views [19, 23] are also introduced to answer graph pattern queries. Views are either given as inputs or generated based on query workloads. Then a subset of views are chosen to answer a query. Hence the optimized schema generated from our approach can be considered as a view on the original property graph, which can be consumed by their technique.

## 7 CONCLUSIONS

To the best of our knowledge, our ontology-driven approach is the first to address the property graph schema optimization problem for domain-specific knowledge graphs. Our approach takes advantages of the rich semantic information in an ontology to drive the property graph schema optimization. The produced schemas gain up to 3 orders of magnitude graph query performance speed-up compared to a direct mapping approach in two real-world knowledge graphs.

## REFERENCES

- [1] Amazon neptune. <https://aws.amazon.com/neptune/>, March 2020.
- [2] Azure cosmos db. <https://azure.microsoft.com/en-us/services/cosmos-db/>, March 2020.
- [3] Federal deposit insurance corporation. <https://www.fdic.gov/regulations/resources/call/index.html>, March 2020.
- [4] Gremlin query language. <https://tinkerpop.apache.org/gremlin.html>, March 2020.

- [5] Janusgraph: Distributed graph database. <http://janusgraph.org/>, March 2020.
- [6] The neo4j graph platform. <https://neo4j.com/>, March 2020.
- [7] Owl 2 web ontology language document overview. <https://www.w3.org/TR/owl2-overview/>, March 2020.
- [8] Securities and exchange commission. <https://www.sec.gov/dera/data/financial-statement-data-sets.html>, March 2020.
- [9] D. J. Abadi, A. Marcus, S. Madden, and K. Hollenbach. Sw-store: a vertically partitioned DBMS for semantic web data management. *VLDB J.*, 18(2):385–406, 2009.
- [10] S. Agrawal, S. Chaudhuri, and V. R. Narasayya. Automated selection of materialized views and indexes in sql databases. In *VLDB*, pages 496–505, 2000.
- [11] K. D. Bollacker, C. Evans, P. Paritosh, T. Sturge, and J. Taylor. Freebase: a collaboratively created graph database for structuring human knowledge. In *ACM SIGMOD*, pages 1247–1250, 2008.
- [12] A. Bonifati, W. Martens, and T. Timm. An analytical study of large SPARQL query logs. *Proc. VLDB Endow.*, 11(2):149–161, 2017.
- [13] M. A. Bornea, J. Dolby, A. Kementsietsidis, K. Srinivas, P. Dantressangle, O. Udrea, and B. Bhattacharjee. Building an efficient RDF store over a relational database. In *ACM SIGMOD*, pages 121–132, 2013.
- [14] S. Brin and L. Page. The anatomy of a large-scale hypertextual web search engine. In *WWW*, pages 107–117, 1998.
- [15] N. Bruno and S. Chaudhuri. Automatic physical database tuning: A relaxation-based approach. In *ACM SIGMOD*, pages 227–238, 2005.
- [16] E. I. Chong, S. Das, G. Eadon, and J. Srinivasan. An efficient sql-based RDF querying scheme. In *VLDB*, pages 1216–1227, 2005.
- [17] V. Christophides, V. Efthymiou, and K. Stefanidis. *Entity Resolution in the Web of Data*. Synthesis Lectures on the Semantic Web: Theory and Technology. Morgan & Claypool Publishers, 2015.
- [18] E. F. Codd. A relational model of data for large shared data banks. *Commun. ACM*, 13(6):377–387, June 1970.
- [19] J. M. F. da Trindade, K. Karanasos, C. Curino, S. Madden, and J. Shun. Kaskade: Graph views for efficient graph analytics. In *ICDE*, pages 193–204, 2020.
- [20] D. Dash, N. Polyzotis, and A. Ailamaki. Cophy: A scalable, portable, and interactive index advisor for large workloads. *PVLDB*, 4(6):362–372, 2011.
- [21] A. Deutsch, Y. Xu, M. Wu, and V. Lee. Tigergraph: A native MPP graph database. *CoRR*, abs/1901.08248, 2019.
- [22] X. L. Dong and D. Srivastava. *Big Data Integration*. Synthesis Lectures on Data Management. Morgan & Claypool Publishers, 2015.
- [23] W. Fan, X. Wang, and Y. Wu. Answering pattern queries using views. *IEEE Trans. Knowl. Data Eng.*, 28(2):326–341, 2016.
- [24] S. Finkelstein, M. Schkolnick, and P. Tiberio. Physical database design for relational databases. *ACM Trans. Database Syst.*, 13(1):91–128, 1988.
- [25] N. Francis, A. Green, P. Guagliardo, et al. Cypher: An evolving query language for property graphs. In *ACM SIGMOD*, pages 1433–1445, 2018.
- [26] S. Harris and N. Shadbolt. SPARQL query processing with conventional relational database systems. In *WISE*, pages 235–244, 2005.
- [27] O. Hartig and J. Hidders. Defining schemas for property graphs by using the graphql schema definition language. In *Proceedings of the 2Nd Joint International Workshop on Graph Data Management Experiences & Systems (GRADES) and Network Data Analytics (NDA)*, GRADES-NDA’19, pages 6:1–6:11, 2019.
- [28] M. S. Hassan, T. Kuznetsova, H. C. Jeong, W. G. Aref, and M. Sadoghi. Extending in-memory relational database engines with native graph support. In *EDBT*, pages 25–36, 2018.
- [29] J. Huang, D. J. Abadi, and K. Ren. Scalable sparql querying of large rdf graphs. *PVLDB*, 4:1123–1134, 2011.
- [30] H. Kimura, G. Huo, A. Rasin, S. Madden, and S. B. Zdonik. Coradd: Correlation aware database designer for materialized views and indexes. *PVLDB*, 3(1-2):1103–1113, 2010.
- [31] J. Lehmann, R. Isele, M. Jakob, et al. Dbpedia - A large-scale, multilingual knowledge base extracted from wikipedia. *Semantic Web*, 2015.
- [32] J. Leskovec, A. Rajaraman, and J. D. Ullman. *Mining of Massive Datasets*. Cambridge University Press, New York, NY, USA, 2nd edition, 2014.
- [33] A. Maduko, K. Anyanwu, A. P. Sheth, and P. Schliekelman. Estimating the cardinality of RDF graph patterns. In *WWW*, pages 1233–1234, 2007.
- [34] M. Meimaris, G. Papastefanatos, N. Mamoulis, and I. Anagnostopoulos. Extended characteristic sets: Graph indexing for SPARQL query optimization. In *ICDE*, pages 497–508, 2017.
- [35] M. J. Mior, K. Salem, A. Aboulmaga, and R. Liu. Nose: Schema design for nosql applications. In *ICDE*, pages 181–192, May 2016.
- [36] T. Neumann and G. Moerkotte. Characteristic sets: Accurate cardinality estimation for RDF queries with multiple joins. In *ICDE*, pages 984–994, 2011.
- [37] T. Neumann and G. Weikum. The RDF-3X engine for scalable management of RDF data. *VLDB J.*, 19(1):91–113, 2010.
- [38] N. Noy, Y. Gao, A. Jain, A. Narayanan, A. Patterson, and J. Taylor. Industry-scale knowledge graphs: Lessons and challenges. *Commun. ACM*, 62(8):36–43, July 2019.
- [39] A. Quamar, C. Lei, D. Miller, et al. An ontology-based conversation system for knowledge bases. In *ACM SIGMOD*, pages 361–376, 2020.
- [40] A. Quamar, F. Özcan, and K. Xirogiannopoulos. Discovery and creation of rich entities for knowledge bases. In *ExploreDB*, 2018.
- [41] I. Robinson, J. Webber, and E. Eifrem. *Graph Databases*. O’Reilly Media, Inc., 2013.
- [42] S. Sakr and G. Al-Naymat. Relational processing of RDF queries: a survey. *SIGMOD Record*, 38(4):23–28, 2009.
- [43] J. Sen, F. Özcan, A. Quamar, G. Stager, A. R. Mittal, M. Jammi, C. Lei, D. Saha, and K. Sankaranarayanan. Natural language querying of complex business intelligence queries. In *SIGMOD*, pages 1997–2000, 2019.
- [44] F. M. Suchanek, G. Kasneci, and G. Weikum. Yago: A large ontology from wikipedia and wordnet. *Semantic Web*, 2008.
- [45] W. Sun, A. Fokoue, K. Srinivas, A. Kementsietsidis, G. Hu, and G. T. Xie. Sqlgraph: An efficient relational-based property graph store. In *ACM SIGMOD*, pages 1887–1901, 2015.
- [46] G. Szárnyas. Incremental view maintenance for property graph queries. *arXiv preprint arXiv:1712.04108*, 2017.
- [47] Y. Tian, E. L. Xu, W. Zhao, et al. IBM db2 graph: Supporting synergistic and retrofittable graph queries inside IBM db2. In *SIGMOD*, pages 345–359, 2020.
- [48] P. Tsialiamanis, L. Sidiropoulos, I. Fundulaki, V. Christophides, and P. A. Boncz. Heuristics-based query optimisation for SPARQL. In *EDBT*, pages 324–335, 2012.
- [49] O. van Rest, S. Hong, J. Kim, X. Meng, and H. Chafi. PGQL: a property graph query language. In *Graph Data-management Experiences and Systems*, page 7, 2016.
- [50] V. V. Vazirani. *Approximation Algorithms*. Springer-Verlag, Berlin, Heidelberg, 2001.
- [51] D. C. Zilio, J. Rao, S. Lightstone, et al. Db2 design advisor: Integrated automatic physical database design. In *VLDB*, pages 1087–1097, 2004.

## A PROOF SKETCH OF THEOREM 3

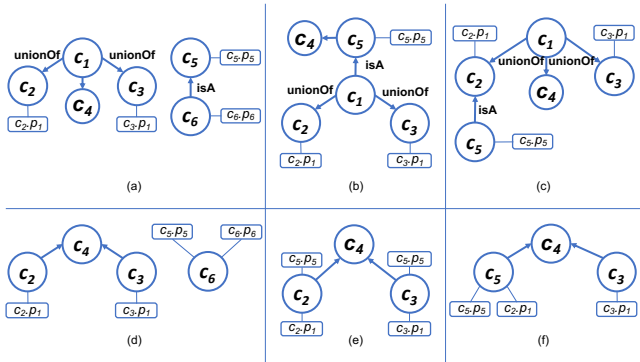
PROOF. Let  $O = (C, R, P)$  be an ontology given as input to Algorithm 5, and let  $O_{out} = (C_{out}, R_{out}, P_{out})$  be the resulting ontology, which is used in Line 18 to produce the output  $\mathcal{PGS}$ . Proving Theorem 3 is equivalent to proving

that applying the rules for any  $R' \subseteq R$  in any order will yield the same result  $O_{out}$ . The theorem trivially holds when  $|R'| = 0$  ( $O_{out} = O$ ), and when  $|R'| = 1$  (only one rule can be triggered).

**Base case.**  $|R'| = 2$ , i.e., for any two relationships, applying the rules in any order yields the same result. Since we only have two relationships, only two rules will be triggered if the relationships are of different types, or one rule will be triggered twice if the two relationships are of the same type. Therefore, we need to prove that applying each pair of rules in any order will yield the same results, examining every possible scenario for each rule.

Specifically, we need to prove that the following pairs of rules are order-independent: (i) union rule and inheritance rule, (ii) inheritance rule and 1:M rule, (iii) union rule and 1:M rule, (iv) inheritance rule and M:N rule, (v) union rule and M:N rule, and (vi) 1:M rule and M:N rule.

(i) *Union and Inheritance.* To prove that union and inheritance rules are order-independent, we examine all the cases in which those two rules may be triggered in the same graph, as shown in Figure 13(a), (b), (c). We assume that the Jaccard similarity between the two concepts connected with an inheritance relationship is less than  $\theta_2$  (see Algorithm 2), so the inheritance rule is triggered and the properties of the parent concept are copied to the child concept. It is straightforward to apply the following observations to the case in which the Jaccard similarity is greater than  $\theta_1$  as well. Figure 13 contains more than two relationships, but only two relationships are sufficient to prove the case<sup>12</sup>. The additional relationships shown are for illustration purpose only.



**Figure 13: Union and Inheritance Rules Independence.**

In the trivial case of Figure 13(a), the source and destination concepts of the union and inheritance relationships are not inter-connected. If we apply the union rule first, we will end up with the left part of Figure 13(d), leaving the right

<sup>12</sup>Consider only the relationships  $(c_1, c_2)$ ,  $(c_6, c_5)$  for Figure 13(a),  $(c_1, c_2)$ ,  $(c_1, c_5)$  for Figure 13(b), and  $(c_1, c_2)$ ,  $(c_5, c_2)$  for Figure 13(c).

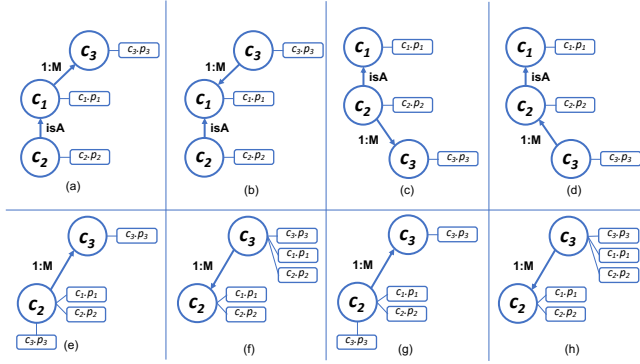
part of Figure 13(a) unchanged, and if we apply the inheritance rule first, we end up with the right part of Figure 13(d), leaving the left part of Figure 13(a) unchanged. In both cases, applying the second rule generates the graph of Figure 13(d).

The case shown in Figure 13(b) is more complex, where the same concept ( $c_1$ ) corresponds to a union concept and a child concept. Applying the union rule first, we remove  $c_1$  and connect its member concepts  $c_2$  and  $c_3$  to  $c_5$  through inheritance relationships. Note that those inheritance relationships come with the same Jaccard value as the original one connecting  $c_1$  to  $c_5$ , which we have assumed to be less than  $\theta_2$ . Then, the inheritance rule is triggered, removing  $c_5$ , copying its properties to its new children  $c_2$  and  $c_3$ , and connecting them to  $c_4$ , as shown in Figure 13(e). If we apply inheritance first, instead of union, then we first remove  $c_5$ , copy its properties to  $c_1$  and connect  $c_1$  to  $c_4$ . Then, applying the union rule, we remove  $c_1$  and connect the member concepts  $c_2$  and  $c_3$  to  $c_4$ , again resulting in the graph of Figure 13(e). The same observations hold for the case in which  $c_1$  corresponds to a parent concept and a union concept.

In a similar way, we can show that union and inheritance rules are order-independent in the case of Figure 13(c), in which the same concept ( $c_2$ ) corresponds to a member concept and a parent concept. If we apply the union rule first, we remove  $c_1$  and connect the member concepts  $c_2$  and  $c_3$  to  $c_4$ . Then, applying the inheritance rule, we remove  $c_2$ , copy its properties to  $c_5$ , and connect  $c_4$  to  $c_5$ , resulting in the graph of Figure 13(f). If we apply the inheritance rule first, we remove  $c_2$ , copy its properties to  $c_5$ , and connect  $c_1$  to  $c_5$  through a union relationship. Finally, we apply the union rule and remove  $c_1$ , connecting  $c_4$  to  $c_5$  and  $c_3$ , also resulting in the graph of Figure 13(f).

(ii) *Inheritance and 1:M.* We follow a similar strategy to prove that inheritance and 1:M rules are order-independent, enumerating all possible cases in which those two rules may be triggered in the same graph, as shown in Figure 14(a), (b), (c), (d). This time, as well as in all the remaining cases (iii) - (vi), the proof is simpler, since there is no alternative intermediate graph involved, if we follow one rule first or another. The only difference is in the set of properties attached to each concept. Again, we assume that the Jaccard similarity between the two concepts connected with an inheritance relationship is less than  $\theta_2$ , so the inheritance rule is triggered and the properties of the parent concept are copied to the child concept.

We skip the trivial case in which the inheritance and 1:M relationships are not related, and start with the case depicted in Figure 14(a), where the parent concept  $c_1$  is also the source concept of an 1:M relationship. If we apply inheritance first, then we copy the properties of  $c_1$  to  $c_2$ , remove  $c_1$  and connect  $c_2$  to  $c_3$  through a 1:M relationship. Then, we apply the 1:M rule and copy  $c_3$ 's properties to  $c_2$ , resulting in the graph



**Figure 14: Inheritance and 1:M Rules Independence.**

of Figure 14(e). If we apply the 1:M rule first, then we first copy the properties of  $c_3$  to  $c_1$  and then we apply inheritance to copy the properties of  $c_1$  (also including the properties of  $c_3$ ) to  $c_2$ , remove  $c_1$  and connect  $c_2$  to  $c_3$  through a 1:M relationship, resulting again in the graph of Figure 14(e).

In the case of Figure 14(b), the parent concept ( $c_1$ ) is now also the destination of an 1:M relationship. If we apply inheritance first, then we copy the properties of  $c_1$  to  $c_2$ , remove  $c_1$  and connect  $c_3$  to  $c_2$  through a 1:M relationship. Then, we apply the 1:M rule and copy  $c_2$ 's properties to  $c_3$ , resulting in the graph of Figure 14(f). If we apply the 1:M rule first, then we first copy the properties of  $c_1$  to  $c_3$  and then we apply inheritance to copy the properties of  $c_1$  to  $c_2$ , remove  $c_1$  and connect  $c_3$  to  $c_2$  through a 1:M relationship. Finally, we apply 1:M rule again (remember that Algorithm 5 iterates until convergence) and and copy the properties of  $v_2$  to  $v_3$ , again resulting in the graph of Figure 14(f).

In Figure 14(c),  $c_2$  is a child and a source concept of a 1:M relationship. In short, if we apply inheritance first, we remove  $c_1$  and copy its properties to  $c_2$  and then we apply 1:M and also copy the properties of  $c_3$  to  $c_1$ , resulting in Figure 14(g). If we apply 1:M first, we copy the properties of  $c_3$  to  $c_2$  and then we apply inheritance to copy the properties of  $c_1$  to  $c_2$  and remove  $c_1$ , again resulting in Figure 14(g).

Finally, in Figure 14(d),  $c_2$  is a child and a destination concept of a 1:M relationship. If we apply inheritance first,

we remove  $c_1$  and copy its properties to  $c_2$  and then we apply 1:M and copy the properties of  $c_2$  (including the properties of  $c_1$ ) to  $c_3$ , resulting in the graph of Figure 14(h). If we apply 1:M first, we copy the properties of  $c_2$  to  $c_3$  and then we apply inheritance to copy the properties of  $c_1$  to  $c_2$  and remove  $c_1$ . Again, we need to trigger the 1:M rule once more to copy the properties of  $c_2$ , now also including the properties of  $c_1$ , to  $c_3$  and get the graph of Figure 14(h).

For the remaining pairs of rules (iii) – (vi), we can follow the same strategy and prove that they are order-independent for all possible cases.

**Induction hypothesis.** Applying the rules in any order for any  $R' \subseteq R$ , where  $|R'|=n$ , always results in the same  $O'$ .

Then, applying the rules in any order for any  $R'' \subseteq R$ , such that  $|R''| = n+1$  and  $R' \subset R''$ , will always result in the same  $O''$ , since there is only one additional relationship in  $R''$  compared to  $R'$ , and only one possible rule corresponding to this new relationship can be triggered.  $\square$

## B 0/1 KNAPSACK PROBLEM REDUCTION

**PROOF.** Given an instance of 0/1 Knapsack problem, our reduction produces the following instance of relationship selection: the cost  $Cost(r_i)$  of relationship  $r_i$  is set to  $w_i$ , and the benefit  $Benefit(r_i)$  of relationship  $r_i$  is set to  $b_i$  as well. We set the space limit  $S$  to  $W$ . Clearly this reduction runs in polynomial time.

If we started with a YES instance of 0/1 Knapsack, then we claim that the reduction produces a YES instance of relationship selection. Suppose there exists a subset  $T \subseteq X$  for which  $\sum_{i \in T} b_i = B$  is maximized and  $\sum_{i \in T} w_i \leq W$ . Then selecting the relationship in  $T$  has total benefit  $B$  and weight no greater than  $W$ , so the instance of relationship selection produced by the reduction is a YES instance.

If the reduction produces a YES instance of relationship selection, then we claim that  $(X, B)$  is a YES instance of 0/1 Knapsack. Let  $T \subseteq X$  be the selected relationships, whose total benefit is  $B$  and whose total cost is at most  $W$ . In other words, we have  $\sum_{i \in T} b_i = B$  and  $\sum_{i \in T} w_i \leq W$ . We conclude that  $(X, B)$  is a YES instance of Knapsack problem as required.  $\square$