

GRAPHITI: Bridging Graph and Relational Database Queries

YANG HE, Simon Fraser University, Canada

RUIJIE FANG, University of Texas at Austin, USA

ISIL DILLIG, University of Texas at Austin, USA

YUEPENG WANG, Simon Fraser University, Canada

This paper presents an automated reasoning technique for checking equivalence between graph database queries written in Cypher and relational queries in SQL. To formalize a suitable notion of equivalence in this setting, we introduce the concept of *database transformers*, which transform database instances between graph and relational models. We then propose a novel verification methodology that checks equivalence modulo a given transformer by reducing the original problem to verifying equivalence between a pair of SQL queries. This reduction is achieved by embedding a subset of Cypher into SQL through syntax-directed translation, allowing us to leverage existing research on automated reasoning for SQL while obviating the need for reasoning simultaneously over two different data models. We have implemented our approach in a tool called GRAPHITI and used it to check equivalence between graph and relational queries. Our experiments demonstrate that GRAPHITI is useful both for verification and refutation and that it can uncover subtle bugs, including those found in Cypher tutorials and academic papers.

CCS Concepts: • **Software and its engineering** → **Automatic programming; Software verification; Formal software verification**; • **Theory of computation** → **Program verification**.

Additional Key Words and Phrases: Program Verification, Equivalence Checking, Relational Databases, Graph Databases.

ACM Reference Format:

Yang He, Ruijie Fang, Isil Dillig, and Yuepeng Wang. 2025. GRAPHITI: Bridging Graph and Relational Database Queries. *Proc. ACM Program. Lang.* 1, 1, Article 1 (January 2025), 44 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

Over the past decades, *graph* databases have garnered significant attention from both industry and academia, offering more flexible data models with different trade-offs compared to *relational* databases. As a result, developers are increasingly interested in migrating relational database applications to graph databases [52], and systems like Apache Age [1] aim to incorporate graph components into relational databases to help with this transition.

Nevertheless, transitioning between relational and graph databases often requires developers to convert queries from one model to the other, and, as evidenced by numerous posts on online forums [3, 38, 47], translating between relational and graph queries can be quite challenging due to misunderstandings of joins versus relationships, aggregation semantics, and other complexities. In fact, we have identified multiple incorrect translations in existing literature where queries claimed to be equivalent were, in reality, non-equivalent [30, 37, 40]. This underscores a growing need for rigorous reasoning about the equivalence of graph and relational queries.

Although significant progress has been made in verifying the equivalence of relational database queries [10–12, 55], to the best of our knowledge, no existing work addresses the equivalence verification problem between relational and graph queries. A key challenge in this area arises from

Authors' Contact Information: Yang He, Simon Fraser University, Burnaby, Canada, yha244@sfu.ca; Ruijie Fang, University of Texas at Austin, Austin, TX, USA, ruijief@cs.utexas.edu; Isil Dillig, University of Texas at Austin, Austin, TX, USA, isil@cs.utexas.edu; Yuepeng Wang, Simon Fraser University, Burnaby, Canada, yuepeng@sfu.ca.

2025. ACM 2475-1421/2025/1-ART1
<https://doi.org/XXXXXXX.XXXXXXX>

the distinct data models of relational and graph databases. Relational databases organize data in tables, with queries operating on rows and columns through well-defined operations like joins and aggregations. In contrast, graph databases represent data as nodes and edges, with queries typically expressing relationships and traversals through graph structures. This fundamental difference in data representation complicates both the definition of equivalence between graph and relational queries and the verification process itself.

This paper takes a first step towards developing automated reasoning techniques that can be used to check equivalence between relational queries written in SQL and graph database queries implemented in Cypher [19], the most popular graph database query language. Our formal reasoning technique is built on a novel method that embeds a subset of the Cypher language into SQL, building on the insight that paths in a graph database instance correspond to joins of rows in a relational database. This observation not only allows us to translate Cypher queries to SQL in a syntax-directed way but also facilitates checking equivalence between Cypher and SQL queries. At a high level, our approach hinges on three crucial components: (1) a formal foundation for defining equivalence between graph and relational database instances (over arbitrary schemas); (2) a correct-by-construction technique for translating graph database queries to relational queries (over a specific schema); and (3) a novel verification methodology that leverages (2) to establish equivalence between Cypher and SQL queries that operate over any arbitrary schema. We next explain each of our contributions in more detail.

Formal foundation for graph and relational database equivalence. As mentioned earlier, a key challenge in reasoning about equivalence between graph and relational queries is the lack of a straightforward mapping between the data models. To address this, we introduce the concept of *database transformer*, adapted from prior work on schema mappings [17, 35, 57], which allows transforming a database instance from one data model (graphs) to an equivalent instance in another model (relational databases). This transformation forms the foundation for defining equivalence between graph and relational database queries.

Correct-by-construction transpilation. Building on this notion of database transformer, we introduce the concept of a *standard database transformer (SDT)* as the default correctness specification. In simple terms, the SDT provides a set of transformation rules to map graph elements (nodes and edges) into relational tables, maintaining the structure and semantics of the graph within the relational model. For example, nodes in a graph schema are transformed into tables in the relational schema, where attributes of the node become columns, and edges are represented as relationships between these tables with foreign keys. The resulting relational schema is referred to as the *induced relational schema*. Our method then defines syntax-directed transpilation rules to convert any Cypher query into a SQL query over this induced schema. The core insight is that Cypher path queries, which perform pattern matching over subgraphs, can be mapped to relational joins. However, the actual translation is tricky due to the fact that Cypher supports flexible, multi-step pattern matching over graph structures, which demands careful handling to ensure that pattern matching in Cypher—whether simple or complex—is accurately represented as joins in SQL, preserving the semantics of the original graph query.

Equivalence checking methodology for arbitrary schema. While the transpilation method described above can generate an equivalent SQL query, the equivalence is *modulo* the SDT. However, in practice, the target relational database often uses a different schema (rather than the *induced relational schema*), so the syntax-directed transpilation method alone is insufficient. To address this gap, we propose a verification methodology that checks equivalence between graph and relational queries modulo any database transformer.

The key insight of our approach is that, instead of directly reasoning about equivalence between Cypher and SQL queries—which would require complex SMT encodings that combine graph and relational structures—we reduce the problem to checking equivalence between SQL queries over different schemas. This reduction allows us to leverage existing techniques and tools for SQL equivalence checking [24, 55], avoiding the need for handling the intricate challenge of reasoning over two fundamentally different data models at the same time.

Figure 1 illustrates the proposed approach for checking equivalence between Cypher and SQL queries. Our method takes four inputs: (1) a Cypher query Q_G , (2) an SQL query Q_R , (3) schemas for the graph and relational databases, and (4) a correctness specification Φ in the form of a database transformer. To establish equivalence between Q_G and Q_R modulo Φ , the method first derives the induced relational schema and the SDT. It then applies a correct-by-construction transpilation technique to generate a SQL query Q'_R that is provably equivalent to Q_G modulo the SDT (but not modulo Φ). In the final step, the method computes a semantic "diff" between the induced and target relational schemas, constructing a *residual database transformer* to align the two relational instances. Finally, an off-the-shelf SQL equivalence checker is then used to check equivalence between Q'_R and Q_R .

Overall, our proposed methodology has two key advantages. First, when the user just wants to translate a Cypher query to SQL but does not care about the underlying relational schema (or if the desired schema is the same as the induced relational schema), our method can be used to perform correct-by-construction transpilation. Second, given any arbitrary correctness specification (in the form of a database transformer), our method can leverage a combination of the proposed transpilation approach and existing automated reasoning tools for SQL to reason about equivalence between any pair of Cypher and SQL queries.

We have implemented the proposed approach in a new tool called GRAPHITI for reasoning about equivalence between Cypher and SQL queries and conducted an extensive experimental evaluation of GRAPHITI on 410 benchmarks. These include 45 benchmarks sourced from public platforms such as StackOverflow, tutorials, and academic papers; 160 benchmarks translated from SQL by students with Cypher experience; and 205 benchmarks translated using ChatGPT. In the first evaluation, we combine GRAPHITI with the VERIEQL [24] bounded model checker for SQL, performing equivalence verification on all 410 query pairs. This reveals equivalence violations in 34 benchmarks, including 3 from the wild, 4 from manual translations, and 27 from GPT-generated translations. In the second evaluation, we pair GRAPHITI with the deductive SQL verifier MEDIATOR [55] for full-fledged verification of an aggregation-free subset. This enables unbounded equivalence verification between Cypher and SQL queries, where both tables and graphs can have arbitrary sizes. Here, about 80% of supported queries are verified as equivalent in a push-button manner. Finally, in the third experiment, we show that GRAPHITI can generate SQL queries that are competitive with manually-written ones in terms of execution efficiency.

We have implemented the proposed approach in a new tool called GRAPHITI for reasoning about equivalence between Cypher and SQL queries and conducted an extensive experimental evaluation of GRAPHITI on 410 benchmarks. These include 45 benchmarks sourced from public platforms such as StackOverflow, tutorials, and academic papers; 160 benchmarks translated from SQL by students with Cypher experience; and 205 benchmarks translated using ChatGPT. In the first evaluation, we combine GRAPHITI with the VERIEQL [24] bounded model checker for SQL, performing equivalence verification on all 410 query pairs. This reveals equivalence violations in 34 benchmarks, including 3 from the wild, 4 from manual translations, and 27 from GPT-generated translations. In the second evaluation, we pair GRAPHITI with the deductive SQL verifier MEDIATOR [55] for full-fledged verification of an aggregation-free subset. This enables unbounded equivalence verification between Cypher and SQL queries, where both tables and graphs can have arbitrary sizes. Here, about 80% of supported queries are verified as equivalent in a push-button manner. Finally, in the third experiment, we show that GRAPHITI can generate SQL queries that are competitive with manually-written ones in terms of execution efficiency.

Contributions. To summarize, this paper makes the following key contributions:

- We propose the first technique for reasoning about equivalence between graph and relational queries, based on a formal definition of equivalence modulo database transformer.
- We introduce the concept of *standard database transformer*, which acts as the default correctness specification for equivalence between Cypher and SQL queries.

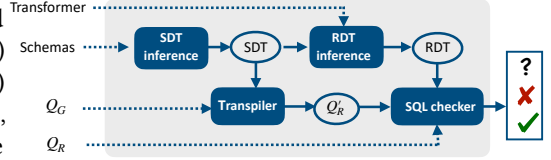


Fig. 1. Overview of approach.

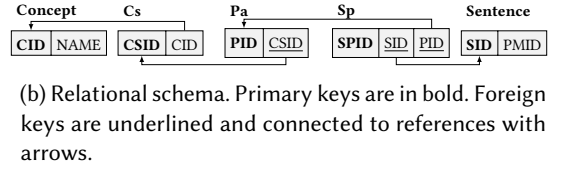
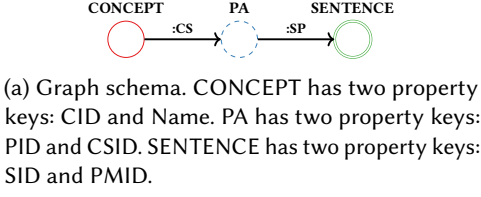
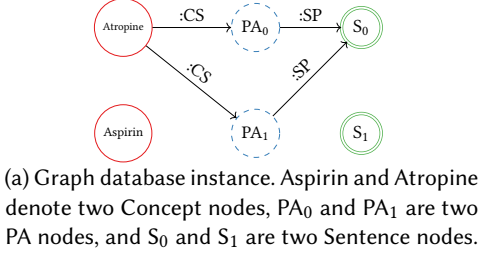


Fig. 2. A pair of relational and graph schemas.



(b) Relational database instance.

CID	NAME
1	Atropine
2	Aspirin

Concept

CID	CSID
1	0
1	1

Cs

PID	CSID
0	0
1	1

Pa

SPID	SID	PID
0	0	0
1	0	1

Sp

SID	PMID
0	0
1	0

Sentence

Fig. 3. Example graph and relational database instances.

- We develop a sound and complete transpilation technique that translates a subset of Cypher queries into equivalent SQL queries, guaranteeing that both queries produce the same results under *standard database transformer*.
- We show how to leverage the proposed transpilation approach to reduce the equivalence checking problem (modulo *any* database transformer) into the problem of checking equivalence between a pair of SQL queries, allowing us to leverage existing work on automated reasoning for SQL.
- We implement these ideas in a tool called GRAPHITI and conduct an empirical evaluation on 410 benchmarks, showing that our approach can be used for both verification and falsification.

2 Motivating Example

In this section, we motivate the problem addressed in this work through an example from prior work [30] that studies SQL analytics for graphs. The queries in this section pertain to a real-world biomedical research database [39].

Incorrect translation. Figure 2a shows a graph schema that contains three types of nodes called CONCEPT, PA¹, and SENTENCE. The CS relationship (represented by an edge) links concepts to predication arguments, and the SP relationship (also represented as an edge) links predication arguments to sentences. On the other hand, Figure 2b shows the relational representation of the graph model. As shown in Figure 2b, the relational schema contains five tables called Concept, Cs, Pa, Sp, Sentence that correspond to both the nodes and edges of the graph schema in Figure 2a. For some intuition about the correspondence between the two databases, Figure 3 shows sample instances of both database schemas that contain the same entries “Atropine” and “Aspirin”.

Next, consider the SQL and Cypher queries shown in Figures 4a and 4c respectively. The SQL query aims to find all concepts that link to a concept c_1 with $CID = 1$ and their corresponding frequencies of connected paths to c_1 through the join of Cs, Pa, and Sp tables. Similarly, the Cypher query first finds all sentences linked to c_1 through the CS - PA - SP path and then counts the frequencies of paths from those sentences to all concepts. According to Lin et al. [30], these queries are intended to be equivalent; however, they are actually *not* equivalent due to the subtle differences in Cypher and SQL semantics. At a high level, both queries explore how certain concepts (starting

¹Here, PA stands for *predication argument*

```

SELECT c2.CID, Count (*) FROM Cs AS c2, Pa AS p2, Sp AS s2
WHERE s2.PID = p2.PID AND p2.CSID = c2.CSID AND s2.SID IN (
  SELECT s1.SID FROM Cs AS c1, Pa AS p1, Sp AS s1
  WHERE s1.PID = p1.PID AND p1.CSID = c1.CSID AND c1.CID = 1 )
GROUP BY CID

```

(a) SQL query

c2.CID	Count(*)
1	2

(b) The result of SQL query.

```

MATCH (c1:CONCEPT {CID:1})-[:CS]->(p1:PA)-[:SP]->(s:SENTENCE)
WITH s
MATCH (s:SENTENCE)-[:SP]->(p2:PA)-[:CS]->(c2:CONCEPT)
RETURN c2.CID, Count (*)

```

(c) Cypher query

c2.CID	Count(*)
1	4

(d) The result of Cypher query.

Fig. 4. A pair of SQL and Cypher queries with their execution results.

```

CONCEPT(cid, name)→Concept(cid, name)      CONCEPT(cid, _), CS(cid, csid, cid, pid), PA(pid, csid)→Cs(cid, csid)
PA(pid, csid)→Pa(pid, csid)                  PA(pid, _), SP(spид, sid, pid, pid, sid), SENTENCE(sid, _)→Sp(spид, sid, pid)
SENTENCE(sid, pmid)→Sentence(sid, pmid)

```

Fig. 5. Database transformer for our example. All variables are implicitly universally quantified.

with CID = 1) are linked to other concepts via their occurrence in shared sentences. They do this by traversing relationships (either explicitly in the graph or via joins in the relational database) and aggregating the results to identify how frequently these connections occur. However, the two queries actually differ in how they compute the frequencies and end up providing different results on two database instances that are meant to contain the same data.

In particular, when run on the database instances from Figure 3, the SQL query produces the table shown in Figure 4b whereas the Cypher query produces the one in Figure 4d. These results agree on the CID's of the related concepts; however they differ on the *frequencies* of the relationships, as is evident from the entries in the Count column of these tables.² As this example illustrates, queries that *appear* to be ostensibly equivalent can have subtle differences in their semantics, motivating the need for automated reasoning tools that can be used to expose semantic differences between graph and relational queries. In the remainder of this section, we elucidate some important aspects and design choices behind our proposed approach.

Need for database transformers. In order to conclude that the SQL and Cypher queries are semantically different, we need to reason about how they behave when executed on the *same data*. However, because graph and relational databases have such different data models, the input database instances are never identical. Thus, in order to reason about query equivalence, we first need a mechanism for defining *data equivalence*. In our framework, this is done through the concept of *database transformer*, which takes as input a graph database instance D over a certain schema Ψ and produces a relational database instance D' over a relational schema Ψ' . Our concept of database transformer is adapted from prior work on *schema mappings* [17, 35], generalized to model the correspondence between graph and relational databases. For our running example, the correspondence between the two database instances is given by the transformer shown in Figure 5. Intuitively, each rule describes how each table in the relational database can be generated based on nodes and edges in the graph. For example, the rule $\text{CONCEPT}(\text{cid}, _), \text{CS}(\text{cid}, \text{csid}, \text{cid}, \text{pid}), \text{PA}(\text{pid}, \text{csid}) \rightarrow \text{Cs}(\text{cid}, \text{csid})$ specifies if there is an edge CS connecting two nodes CONCEPT and PA in the graph, and the first two properties of CS are cid and csid, then there is a row (cid, csid) in the Cs table. Here, the last two attributes of CS serve as foreign keys referencing to the source and target nodes of the CS edge.

Induced relational schema and default transformer. As mentioned in Section 1, our approach to reasoning about equivalence between SQL and Cypher queries relies on first transpiling the given

²An equivalent Cypher query of the SQL query in Figure 4a is shown in Appendix C.

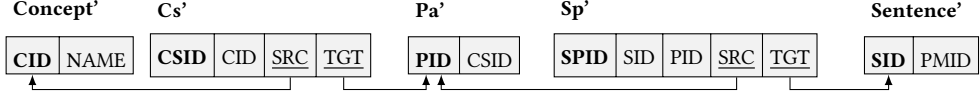


Fig. 6. Induced relational schema. Primary keys are in bold. Foreign keys are underlined and connected to references with arrows.

Concept' (cid, name) \rightarrow Concept (cid, name) Concept' (cid, _), Cs' (cid, csid, cid, pid), Pa' (pid, csid) \rightarrow Cs (cid, csid)
 Pa' (pid, csid) \rightarrow Pa (pid, csid) Pa' (pid, _), Sp' (spid, sid, pid, pid, sid), Sentence' (sid, _) \rightarrow Sp (spid, sid, pid)
 Sentence' (sid, PMID) \rightarrow Sentence (sid, PMID)

Fig. 8. Residual database transformer. All variables are universally quantified.

Cypher query to a SQL query over the *induced relational schema*, which corresponds to a natural relational representation of the graph database. In particular, Figure 6 shows the induced relational schema for Figure 2a and is obtained by (a) translating both node and edge types in the graph schema into tables, and (b) translating incidence and adjacency information between node and edge types into *functional dependencies*. For instance, nodes of type CONCEPT are mapped to the Concept table, edges of type CS are mapped to the Cs table, and so on. To handle functional dependencies, we need to introduce foreign keys in the corresponding tables. Compared to the relational schema in Figure 2b, the induced relational schema preserves the attributes while using additional attributes (i.e. SRC and TGT) as foreign keys to represent functional dependencies. For example, for edges of type CS, the source node is of type CONCEPT; thus, the induced relational schema has the SRC attribute as a foreign key to CID of the Concept table. Similarly, the TGT attribute is considered a foreign key to the PA table. Given the original graph database schema, our method constructs a so-called *standard database transformer* Φ_{std} that can be used to convert any instance of the given graph schema to a relational database over its induced schema.

Syntax-directed transpilation. Intuitively, the standard database transformer establishes a one-to-one correspondence between elements of the graph database and the corresponding entries in the induced relational database. Hence, we can use syntax-directed translation to directly transpile the Cypher query to a SQL query over the induced schema. Here, the transformer Φ_{std} fixes the mapping between “atomic elements” of both databases; thus, atomic queries over nodes and edges in Cypher can be translated to atomic queries over SQL tables. This forms the base case for an inductive transpilation scheme, leveraging the key insight that Cypher path queries correspond to relational joins. For example, a Cypher path query like **MATCH** (u)-[:CS]->(v) translates to a SQL join query between the Concept, Cs, and Pa tables, since the database transformer specifies that the CS edge type maps to the Cs table in the relational schema. Such a transpilation scheme is conceptually simple, but as our final transpilation result in Figure 7 shows, one must take significant care to address the different types of path patterns in Cypher syntax, in addition to translating compositions of path queries with other Cypher operators such as aggregation. Specifically, the first pattern matching in Figure 4c is translated to T1 while the WITH clause propagates the intermediate results to T2. Similarly, the second pattern matching is translated to T3. Due to the shared node s :SENTENCE in these two path patterns, we join T2 and T3 as T4. The final RETURN clause is translated to GroupBy because of the aggregation expression.

```

WITH T1 AS ( SELECT c1.CID AS c1_CID, ..., s.SID AS s_SID
FROM Concept AS c1 JOIN CS AS r1
      JOIN PA AS p1 JOIN SP AS r2 JOIN Sentence AS s
ON c1.CID = 1 AND c1.CID = r1.SRC AND ... AND r2.TGT = s.SID ),
T2 AS ( SELECT s.SID FROM T1 ),
T3 AS ( SELECT s.SID AS s_SID, ..., c2.CID AS c2_CID
FROM Sentence AS s JOIN SP AS r3
      JOIN PA AS p2 JOIN CS AS r4 JOIN Concept AS c2
ON s.SID = r3.TGT AND ... AND r4.SRC = c2.CID ),
T4 AS ( SELECT * FROM T2 JOIN T3 ON T2.s_SID = T3.s_SID )
SELECT T4.c2_CID, Count (*) FROM T4 GROUP BY T4.c2_CID
  
```

Fig. 7. Transpilation result for the Cypher query.

Checking equivalence. While our approach allows correct-by-construction transpilation from Cypher to SQL, there are several reasons why this is not sufficient. First, our approach does not guarantee that the transpilation result is the most efficient, so even though one could use existing SQL query optimizers to further optimize the query, the user may want to write their hand-optimized SQL query. Second, the user may want to use a different relational schema rather than our default version. Third, while our transpilation rules allow translating Cypher to SQL, they do not address the reverse direction. Motivated by these shortcomings, our method leverages the proposed transpilation algorithm to perform verification between any given pair of Cypher and SQL queries by utilizing existing tools for SQL. In particular, the key idea is to infer a *residual database transformer* that specifies the relationship between the relational database over the induced schema and the target relational database (as specified by the user-provided database transformer). For our running example, Figure 8 shows the residual transformer that can be used to convert instances of the induced relational schema from Figure 6 to instances of the desired schema. Given such a residual schema, we can use an existing SQL equivalence checker, such as VERIEQL [24], to refute equivalence between these queries and obtain the counterexample shown in Figure 3.

3 Preliminaries

3.1 Background on Graph Databases

A graph database instance is a *property graph*, which contains nodes and edges carrying data. Typically, nodes model entities, and edges model relationships. Each node or edge in the graph stores data represented as pairs of property keys and values. Additionally, each node or edge is assigned a *label*, which describes the kind of entity or relationship it models. A well-formed property graph should conform to a graph schema, which is formalized below.

Definition 3.1 (Node/edge type). A node type t_{node} is a tuple (l, K_1, \dots, K_n) where l is the label of the node (e.g., Actor) and K_1, \dots, K_n are the property keys for that node type (e.g., name, dob, etc.). An edge type t_{edge} is also a tuple $(l, t_{\text{src}}, t_{\text{tgt}}, K_1, \dots, K_m)$ where l is a label (e.g., ACTS_IN), t_{src} and t_{tgt} are the types of the source and target nodes respectively, and K_1, \dots, K_m are the property keys (e.g., role) for that edge type.

For each node type $t_{\text{node}} = (l, K_1, \dots, K_n)$, we define $\text{label}(t_{\text{node}})$ to give the label l of t_{node} , and we assume that K_1 is the *default property key* for t_{node} , which is a key with a globally unique value, similar to a primary key in a relational database. We define $\text{keys}(t_{\text{node}}) = \{K_1, \dots, K_n\}$ to yield the set of property keys associated with t_{node} . For an edge type $t_{\text{edge}} = (l', t_{\text{src}}, t_{\text{tgt}}, K'_1, \dots, K'_m)$ we similarly define $\text{label}(t_{\text{edge}}) = l'$ and $\text{keys}(t_{\text{edge}}) = \{K'_1, \dots, K'_m\}$, with key K'_1 being the default property key. Additionally, we define $\text{dstType}(t_{\text{edge}}) = t_{\text{tgt}}$ and $\text{srcType}(t_{\text{edge}}) = t_{\text{src}}$.

Definition 3.2 (Graph database schema). A graph database schema Ψ_G is a pair (T_N, T_E) where T_N is a set of node types and T_E is a set of edge types.

For each graph database schema $\Psi_G = (T_N, T_E)$, we assume that the label of each node or edge type uniquely defines it: $\forall t_1, t_2 \in T_N \cup T_E. \text{label}(t_1) \neq \text{label}(t_2)$. Thus, we can use types and labels interchangeably. Additionally, we assume that all property keys are unique inside a given schema Ψ_G , i.e., there are no name clashes between arbitrary pairs of property keys between different types.

Definition 3.3 (Graph database). An instance of a graph database schema $\Psi_G = (T_N, T_E)$ is a tuple $G = (N, E, P, T)$ where N is a set of nodes, $E \subseteq N \times N$ is a set of edges, $P : (N \cup E) \times \mathbf{Keys} \rightarrow \mathbf{Values}$, and $T : N \cup E \rightarrow T_N \cup T_E$ gives the type of a node $n \in N$ or an edge $e \in E$.

Query	Q	$::= R \mid \text{OrderBy}(R, k, b) \mid \text{Union}(Q, Q) \mid \text{UnionAll}(Q, Q)$
Return Query	R	$::= \text{Return}(C, \bar{E}, \bar{k})$
Clause	C	$::= \text{Match}(PP, \phi) \mid \text{Match}(C, PP, \phi) \mid \text{OptMatch}(C, PP, \phi) \mid \text{With}(C, \bar{X}, \bar{X})$
Path Patt.	PP	$::= NP \mid NP, EP, PP$
Node Patt.	NP	$::= (X, l)$
Edge Patt.	EP	$::= (X, l, d)$
Expression	E	$::= k \mid v \mid \text{Cast}(\phi) \mid \text{Agg}(E) \mid E \oplus E$
Predicate	ϕ	$::= \top \mid \perp \mid E \odot E \mid \text{IsNull}(E) \mid E \in \bar{v} \mid \text{Exists}(PP) \mid \phi \wedge \phi \mid \phi \vee \phi \mid \neg \phi$

$X \in \text{Node/Edge Names} \quad l \in \text{Labels} \quad k \in \text{Property Keys} \quad v \in \text{Values} \quad b \in \text{Bools}$
 $\text{Agg} \in \{\text{Count, Avg, Sum, Min, Max}\} \quad d \in \{\rightarrow, \leftarrow, \leftrightarrow\}$

Fig. 9. Featherweight Cypher syntax where \oplus, \odot correspond to arithmetic and logical operators respectively.

We use the notation $P(n, k)$ to give the value of a property key k in node n and analogously define $P(e, k)$ for an edge e . We also use the notation $G \triangleright \Psi_G$ to denote that G is an instance of schema Ψ_G , and we refer to any subgraph of G as a *property graph*.

3.2 Query Language for Graph Databases

To formalize our method, we focus on a subset of the popular Cypher database query language [19]. This subset, which we refer to as “Featherweight Cypher”, is presented in Figure 9.³ A query Q is either a union of return queries R or an order-by statement following one or more such return queries. Each return query R takes as input a clause C , a list of expressions \bar{E} , and a list of property key names \bar{k} . Intuitively, the return query shapes a list of graphs into a table. Each clause in the return statement is a *Match*, representing a pattern match over an input property graph, and the patterns are specified using the path pattern PP . We do not include Cypher features such as unbounded-length path queries and graph reachability primitives (e.g., `shortestPath`), which are not expressible in the core SQL fragment considered in this paper.

Example 3.4. Consider the following Cypher query

MATCH (n:EMP)-[:WORK_AT]->(m:DEPT) **RETURN** m.dname **AS** name, **Count**(n) **AS** num

that returns a table containing department names and the number of employees. We can represent it using our featherweight Cypher syntax as follows

$\text{Return}(\text{Match}([(n, \text{EMP}), (e, \text{WORK_AT}, \rightarrow), (m, \text{DEPT})], \top), [m.dname, \text{Count}(n.id)], [name, num])$

Here, the match clause retrieves all paths of length one from EMP nodes to DEPT nodes connected by an edge of type WORK_AT. Then, the return clause reshapes the set of matched paths into a table with two columns: name and num. The name column is populated with values corresponding to the property key dname of the DEPT node, and the num column is populated by the count of n.id, where m and n refer to the source and target nodes of the matched edge, respectively.

3.3 Relational Databases

Definition 3.5 (Relational schema). A relational database schema is a pair $\Psi_R := (S, \xi)$ where $S: \mathcal{R} \rightarrow [\mathcal{A}]$ is a mapping from a set of relation names \mathcal{R} to a list of attributes, and ξ is an *integrity constraint*. We assume that all attribute names in a schema are unique.

We represent an integrity constraint as a conjunction of three types of *atomic* constraints:

- (1) **Primary key constraints:** A primary key constraint $\text{PK}(R) = a$ specifies that attribute a is the *primary key* for relation R — i.e., there cannot be multiples tuples of R that agree on a .

³The denotational semantics is formally described in Appendix A.

Query	Q	$::= R \mid \Pi_L(Q) \mid \sigma_\phi(Q) \mid \rho_R(Q) \mid Q \cup Q \mid Q \bowtie Q \mid Q \otimes Q$ $\mid \text{GroupBy}(Q, \bar{E}, L, \phi) \mid \text{With}(\bar{Q}, \bar{R}, Q) \mid \text{OrderBy}(Q, a, b)$
Attribute List	L	$::= E \mid \rho_a(E) \mid L, L$
Attribute Expr	E	$::= a \mid v \mid \text{Cast}(\phi) \mid \text{Agg}(E) \mid E \oplus E$
Predicate	ϕ	$::= b \mid E \odot E \mid \text{IsNull}(E) \mid E \in \bar{v} \mid \bar{E} \in Q \mid \phi \wedge \phi \mid \phi \vee \phi \mid \neg \phi$
Join Op	\otimes	$::= \times \mid \bowtie_\phi \mid \bowtie_{\neg \phi} \mid \bowtie_{\phi} \mid \bowtie_{\neg \phi}$

$R \in \text{Relation Names}$ $a \in \text{Attr Names}$ $v \in \text{Values}$ $b \in \text{Bools}$ $\text{Agg} \in \{\text{Count}, \text{Avg}, \text{Sum}, \text{Min}, \text{Max}\}$

Fig. 10. Featherweight SQL syntax; \oplus and \odot represent arithmetic and logical operators respectively

Transformer	Φ	$::= P, \dots, P \rightarrow P \mid \Phi \Phi$
Predicate	P	$::= E(t, \dots, t)$
Term	t	$::= c \mid v \mid _$

$E \in \text{Table Names} \cup \text{Node Labels} \cup \text{Edge Labels}$ $c \in \text{Constants}$ $v \in \text{Variables}$

Fig. 11. Syntax of the database transformer.

- (2) **Foreign key constraints:** A foreign key constraint $\text{FK}(R.a) = R'.a'$ specifies that the attribute a in relation R is a *foreign key* corresponding to attribute a' in relation R' . That is, the values stored in attribute a of relation R must be a subset of the values stored in attribute a' of R' .
- (3) **Not-null constraints:** A not-null constraint $\text{NotNull}(R, a)$ specifies that the value stored at attribute a of relation R must not be Null.

Definition 3.6 (Relational database instance). A relational database instance R is a collection of tuples $\{r_1, \dots, r_m\}$, where each $r_i \in R$ is of the form $(a_1 : v_1, \dots, a_n : v_n)$. Here, a_1, \dots, a_n are attributes and v_1, \dots, v_n are values. We let $\text{Attrs}(r_i)$ return the list of attributes a_1, \dots, a_n in sequence. We use the notation $r_i.a$ to denote the value stored in attribute a of tuple r_i .

As with graph databases, we use the notation $R \triangleright \Psi_R$, to denote that R is instance of Ψ_R .

SQL query language. In this paper, we consider relational database queries written in SQL. Figure 10 shows the subset of SQL that we use in our formalization. At a high level, this language extends core relational algebra (e.g., projection Π , selection σ , renaming ρ , joins \otimes , set and bag unions \cup, \bowtie) to incorporate `GroupBy`, `OrderBy`, and `With` clauses. It can express a representative fragment of SQL queries that are commonly used in practice. The semantics of these SQL operators are standard and formally defined by prior work such as He et al. [24].

4 Problem Statement

In this section, we first describe the language for database transformers and then formally define the equivalence checking problem between graph and relational databases.

4.1 Language for Database Transformers

In this section, we present a small domain-specific language (DSL), shown in Figure 11, for expressing database transformers. Following prior work [57], our DSL generalizes the standard concept of *schema mapping* [17, 35] for relational databases to a more flexible form. In particular, a database transformer in this DSL is expressed as a set of first-order formulas of the form $P_1, \dots, P_n \rightarrow P_0$, where each P_i is a predicate that represents a database element, such as a table in a relational database or a node or edge in a graph database. Each predicate is of the form $E(t_1, \dots, t_n)$ where E corresponds a table name, node labels, or edge labels, and each t_j is a term (variable or a constant), with $_$ denoting a fresh variable that is not used. All variables are implicitly universally quantified but the quantifiers are omitted in the syntax for brevity. Intuitively, the formula $P_1, \dots, P_n \rightarrow P_0$ expresses that, if predicates P_1, \dots, P_n hold over a database instance D , then predicate P_0 holds over another database instance D' .

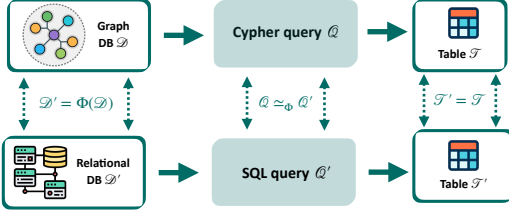


Fig. 12. Definition of equivalence between Cypher and SQL queries where Φ is the user-provided schema transformer. A pair of Cypher and SQL queries are equivalent if they produce the same table (modulo renaming) whenever they are executed on a pair of database instances satisfying Φ .

Semantics. To define the semantics of our transformer DSL, we first represent a database transformer Φ as a set of universally quantified first-order logic formulas, denoted as $\llbracket \Phi \rrbracket$. The idea behind the semantics of our DSL is to represent each database instance as a set of ground predicates and then check whether these predicates entail the first-order logic formula $\llbracket \Phi \rrbracket$ under the Herbrand semantics.

To make this discussion more precise, we introduce a function C that maps a database instance to a set of ground predicates representing its structure and contents. Formally, C is defined as follows:

$$C(D) = \{E(t_1, \dots, t_n) \mid E \in D\}$$

The mapping of elements in D to predicates depends on whether D is a relational or graph database:

- **For relational database instance D :** If R is a table in D with a set of records $\{(a_1, \dots, a_n)\}$, then R is converted to the following set of ground predicates:

$$\{R(a_1, \dots, a_n) \mid (a_1, \dots, a_n) \in R\}$$

Here, R represents the table name, and each a_i is a constant representing a value in the record.

- **For graph database instance D :** If $N(l, a_1, \dots, a_n)$ denotes a node with label l and values a_1, \dots, a_n of property keys K_1, \dots, K_n , nodes are converted to the following set of ground facts:

$$\{l(a_1, \dots, a_n) \mid \text{node } N(l, a_1, \dots, a_n) \in D\}$$

Similarly, if $E(l, s, t, a_1, \dots, a_n)$ denotes an edge with label l that connects nodes s and t and has property values a_1, \dots, a_n , then edges are converted to predicates as follows:

$$\{l(a_1, \dots, a_n, s, t) \mid \text{edge } E(l, s, t, a_1, \dots, a_n) \in D\}$$

Given a database instance D , $C(D)$ yields a set of ground predicates representing the structure and contents of D . We can now define the semantics of the transformer Φ as follows:

$$\Phi(D) = D' \iff C(D) \cup C(D') \models \llbracket \Phi \rrbracket$$

where the notation $S \models \varphi$ indicates that the set S of ground predicates is a Herbrand model of φ .

Example 4.1. Consider the graph and relational database instances G, R from Figures 3a and 3b respectively. For the transformer Φ shown in Figure 5, we have $\Phi(G) = R$.

4.2 Equivalence Checking Problem

In this section, we formally define what it means for a pair of graph and relational queries to be equivalent modulo a database transformer Φ , expressed in the DSL of Section 4.1. To this end, we first introduce some necessary definitions and notations.

Definition 4.2 (Database query). A database query Q over schema Ψ takes as input a database instance D such that $D \models \Psi$ and yields a table. We denote the semantics of Q as $\llbracket Q \rrbracket_D$.

The definition above is sufficiently general to describe both SQL and Cypher queries, as both query languages return a table.

Algorithm 1 Methodology for checking equivalence between Cypher and SQL queries

```

1: procedure CHECKEQUIVALENCE( $\Psi_G, Q_G, \Psi_R, Q_R, \Phi$ )
   Input: Graph and relational schemas  $\Psi_G, \Psi_R$ , Cypher query  $Q_G$ , SQL query  $Q_R$ , transformer  $\Phi$ 
   Output:  $\top$  for equivalence or  $\perp$  indicating failure
2:    $(\Phi_{\text{sd}}, \Psi'_R) \leftarrow \text{INFERSDT}(\Psi_G)$ 
3:    $Q_{R'} \leftarrow \text{TRANPILE}(Q_G, \Phi_{\text{sd}}, \Psi'_R)$ 
4:   return REDUCETOSQL( $\Psi_R, Q_R, \Psi'_R, Q_{R'}, \Phi, \Phi_{\text{sd}}$ )

```

Definition 4.3 (Database equivalence). Let D, D' be database instances over schemas Ψ, Ψ' respectively and let Φ be a database transformer that can be used to convert instances of Ψ to instances of Ψ' . Then, D is said to be equivalent to D' modulo Φ , denoted $D \sim_\Phi D'$, if $\Phi(D) = D'$.

According to the above definition, a graph database instance G is equivalent to a relational database instance R if the contents of R can be obtained from G by applying the transformer Φ to G . Next, to define equivalence between graph and relational queries, we need to define what it means for the query outputs to be the same. Since queries in both Cypher and SQL return tables, we need a notion of equivalence between tables.

Definition 4.4 (Table equivalence). Two tables T and T' are said to be equivalent, denoted $T \equiv T'$, if there exists a bijective mapping π from columns of T to those of T' such that, for each tuple $r \in T$ with multiplicity n , there exists a unique tuple $r' \in T'$ with multiplicity n where $\forall a \in \text{Attrs}(r). r.a = r'.\pi(a)$.

In other words, our notion of table equivalence disregards the order of attributes as well as their names, which allows for a more robust notion of query equivalence.⁴

Definition 4.5 (Graph-relational query equivalence). Let Ψ_G and Ψ_R be graph and relational schemas respectively, and Φ be a transformer from Ψ_G to Ψ_R . A query Q over Ψ_G is *equivalent* to Q' over Ψ_R modulo Φ , denoted $Q \simeq_\Phi Q'$, iff:

$$\forall G, R. (G \triangleright \Psi_G \wedge R \triangleright \Psi_R \wedge G \sim_\Phi R) \Rightarrow \llbracket Q \rrbracket_G \equiv \llbracket Q' \rrbracket_R$$

In other words, Q, Q' are considered equivalent if they produce the same tables (modulo renaming/re-ordering of columns) when executed on a pair of database instances G, R satisfying $G \sim_\Phi R$. Figure 12 visualizes this definition of equivalence between graph databases and relational databases.

Definition 4.6 (Equivalence checking problem). Given graph and relational queries Q, Q' and a transformer Φ from graph schema Ψ_G to relational schema Ψ_R , the equivalence checking problem is to decide whether $Q \simeq_\Phi Q'$.

5 Equivalence Checking Algorithm

In this section, we present our algorithm, summarized in Algorithm 1, for checking equivalence between Cypher and SQL queries. As shown in Algorithm 1, our algorithm consists of three steps:

- (1) **Schema and transformer inference:** Given the graph database schema Ψ_G , our algorithm first invokes the `INFERSDT` function to infer both the induced relational schema $\Psi'_R = (S, \xi)$ as well as the standard database transformer Φ_{sd} .

⁴If a query includes an `OrderBy` clause, list semantics will be applied, where the result is an ordered list of tuples. In this case, the equivalence of two tables T and T' is defined such that there exists a bijective mapping π between their attributes, and for each pair of tuples $r \in T$ and $r' \in T'$ at the same index, it holds that $\forall a \in \text{Attrs}(r). r.a = r'.\pi(a)$.

$$\begin{array}{c}
\frac{t_{\text{node}} = (l, K_1, \dots, K_n) \quad \xi = (\text{PK}(R_l) = K_1) \quad \Phi = \{l(K_1, \dots, K_n) \rightarrow R_l(K_1, \dots, K_n)\}}{t_{\text{node}} \hookrightarrow (\{R_l \mapsto (K_1, \dots, K_n)\}, \xi, \Phi)} \text{(Node)} \\
\\
\frac{\begin{array}{l} t_{\text{edge}} = (l, t_{\text{src}}, t_{\text{tgt}}, K_1, \dots, K_m) \quad \text{label}(t_{\text{src}}) = s \quad \text{label}(t_{\text{tgt}}) = t \\ \xi = \text{PK}(R_l) = R_l.K_1 \wedge \text{FK}(R_l.\text{fk}_s) = \text{PK}(R_s) \wedge \text{FK}(R_l.\text{fk}_t) = \text{PK}(R_t) \\ \Phi = \{l(K_1, \dots, K_m, \text{fk}_s, \text{fk}_t) \rightarrow R_l(K_1, \dots, K_m, \text{fk}_s, \text{fk}_t)\} \end{array}}{t_{\text{edge}} \hookrightarrow (\{R_l \mapsto (K_1, \dots, K_m, \text{fk}_s, \text{fk}_t)\}, \xi, \Phi)} \text{(Edge)} \\
\\
\frac{T_1 \hookrightarrow (S_1, \xi_1, \Phi_1) \quad T_1 \hookrightarrow (S_2, \xi_2, \Phi_2)}{T_1 \uplus T_2 \hookrightarrow (S_1 \uplus S_2, \xi_1 \wedge \xi_2, \Phi_1 \cup \Phi_2)} \text{(Set)} \quad \frac{(T_N \uplus T_E) \hookrightarrow (S, \xi, \Phi)}{(T_N, T_E) \hookrightarrow (S, \xi, \Phi)} \text{(Schema)}
\end{array}$$

Fig. 13. Rules for the INFERSDT procedure. R'_l denotes the table name of R_l in the induced relational schema.

- (2) **Syntax-directed transpilation:** Next, our algorithm uses the inferred database transformer and integrity constraints to transpile the Cypher query into an SQL query Q'_R that is guaranteed to be equivalent to Q_G modulo the standard Φ_{sdt} .
- (3) **Checking SQL equivalence:** Finally, the algorithm computes a residual database transformer Φ_{rdt} that can be used to convert instances of Ψ'_R into instances of Ψ_R and checks equivalence between SQL queries Q_R and Q'_R modulo the residual database transformer Φ_{rdt} relating a pair of database schemas.

Discussion. An alternative approach to solving the equivalence checking problem would be to directly reason about equivalence between the graph query Q_G and the relational query Q_R . While such an approach would be more direct, we adopt the methodology illustrated in Figure 1 for several reasons. First, in order to directly verify equivalence between graph and relational database queries, we need suitable SMT encodings of *both* graphs and relations, which makes the resulting constraint solving problem harder compared to the alternative. Second, the reduction to relational equivalence checking allows us to leverage a variety of existing tools that have been developed for SQL, including testing tools [6], bounded model checkers [24], and deductive verifiers [55].

5.1 Induced Relational Schema and Standard Transformer Inference

We first discuss the INFERSDT procedure for inferring the induced relational schema as well as the standard database transformer. This procedure is formalized in Figure 13 using inference rules of the form: $\Psi_G \hookrightarrow (S, \xi, \Phi_{\text{sdt}})$ where Ψ_G corresponds to elements of the graph schema (nodes, edges, and subgraphs), $\Psi_R = (S, \xi)$ is the induced relational schema for Ψ_G , and Φ_{sdt} is the standard database transformer that can be used to convert instances of Ψ_G into instances of Ψ_R .

Induced relational schema. Intuitively, the induced relational schema is the “closest” relational representation of the graph database schema $\Psi_G = (T_N, T_E)$ that represents each node and edge type as a relational table. As shown in the Node rule, for each node type $t_{\text{node}} = (l, K_1, \dots, K_n) \in T_N$ with label l and default property key K_1 , we introduce a table R_l with attributes K_1, \dots, K_n in the induced relational schema. We also use the default property key K_1 as the primary key of the corresponding table and generate an integrity constraint $\text{PK}(R_l) = K_1$. Similarly, as shown in the Edge rule, for each edge type $t_{\text{edge}} = (l, t_{\text{src}}, t_{\text{tgt}}, K_1, \dots, K_m) \in T_E$ with default property key K_1 , we introduce a table R_l with attributes $K_1, \dots, K_m, \text{fk}_s, \text{fk}_t$. Here, K_1 is also the primary key of table R_l , and fk_s, fk_t are foreign keys which reference to primary keys of the tables corresponding to source and target nodes. Thus, the integrity constraint is $\text{PK}(R_l) = R_l.K_1 \wedge \text{FK}(R_l.\text{fk}_s) = \text{PK}(R_s) \wedge \text{FK}(R_l.\text{fk}_t) = \text{PK}(R_t)$, where R_s, R_t are the tables corresponding to the source and target nodes.

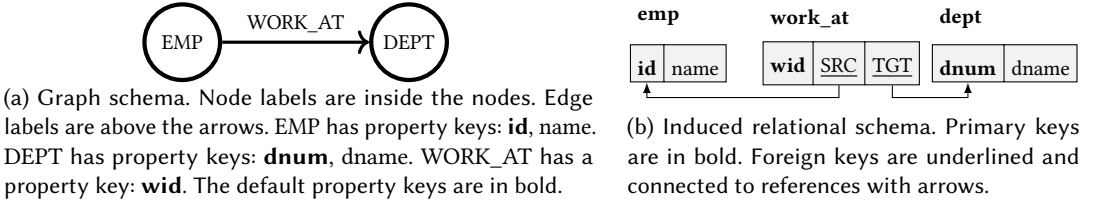


Fig. 14. Example of a graph schema and its induced relational schema.

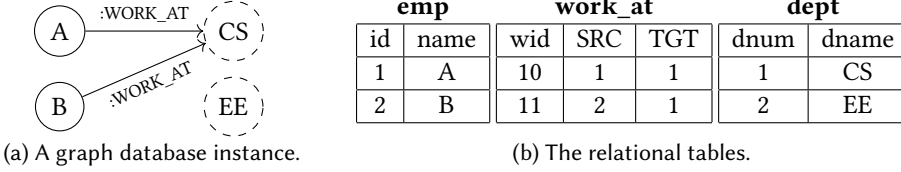


Fig. 15. Example graph database and its corresponding relational database over the induced relational schema.

Example 5.1. Consider the graph schema Ψ_G shown in Figure 14a. Its induced relational schema Ψ_R is visualized in Figure 14b.

Standard database transformer. The standard database transformer (SDT) is expressed in the same language as a general database transformer from Section 4.1, which transforms instances of a graph schema Ψ_G to instances of a relational schema Ψ_R . Intuitively, the standard database transformer for a graph schema Ψ_G converts each node and edge type to a separate table, and all occurrences of that element type in a graph database G become tuples in the corresponding table of the relational database. Specifically, as shown in the Node rule, for each node of type $t_{\text{node}} = (l, K_1, \dots, K_n)$, we generate a formula $l(K_1, \dots, K_n) \rightarrow R_l(K_1, \dots, K_n)$ which transforms a predicate $R_l(v_1, \dots, v_n)$ representing a graph element to a predicate $R'_l(v_1, \dots, v_n)$ representing a tuple in the relational database. Similarly, we generate a formula $l(K_1, \dots, K_m, \text{fk}_s, \text{fk}_t) \rightarrow R_l(K_1, \dots, K_m, \text{fk}_s, \text{fk}_t)$ for each edge type in the graph schema.

Example 5.2. Consider the graph database G visualized in Figure 15a. The SDT Φ_{sdT} transforms G to the relational database shown in Figure 15b.

5.2 Syntax-Directed Transpilation

Building upon the standard database transformer (SDT) introduced earlier, we now turn to the core task of translating Cypher queries into corresponding SQL queries over the induced schema. This process presents several challenges, as it involves mapping graph-based operations in Cypher to the relational model in SQL, while ensuring that the original query semantics are preserved. Specifically, Cypher includes features like pattern matching over subgraphs and optional matches, which do not have direct equivalents in SQL. Additionally, Cypher aggregates data over matched subgraphs, whereas SQL aggregates over grouped tuples. Finally, ensuring consistent mappings of graph nodes and edges across the query is critical to maintaining the integrity of references throughout the process. Our syntax-directed transpilation approach addresses these challenges by converting Cypher queries into SQL queries over the induced relational schema. The key idea is that path patterns in Cypher can be represented by relational joins in SQL. For instance, Cypher's match clauses are translated into SQL inner joins, and optional match clauses map to outer joins. This approach ensures that the pattern matching semantics of Cypher queries are faithfully represented within the relational model, maintaining the integrity of the original graph-based operations. We

$\frac{\neg\text{hasAgg}(\bar{E}) \quad \Phi_{\text{sdt}}, \Psi_R \vdash C \xrightarrow{\text{clause}} X, Q \quad \Phi_{\text{sdt}}, \Psi_R \vdash E_i \xrightarrow{\text{expr}} E'_i \quad 1 \leq i \leq \bar{E} }{\Phi_{\text{sdt}}, \Psi_R \vdash \text{Return}(C, \bar{E}, \bar{k}) \xrightarrow{\text{query}} \Pi_{\rho_{\bar{k}}(\bar{E}')} (Q)} \quad (\text{Q-Ret})$	$\frac{\text{hasAgg}(\bar{E}) \quad \Phi_{\text{sdt}}, \Psi_R \vdash E_i \xrightarrow{\text{expr}} E'_i \quad 1 \leq i \leq \bar{E} \quad \Phi_{\text{sdt}}, \Psi_R \vdash C \xrightarrow{\text{clause}} X, Q \quad \bar{A} = \text{filter}(\lambda x. \neg\text{IsAgg}(x), \bar{E}')}{\Phi_{\text{sdt}}, \Psi_R \vdash \text{Return}(C, \bar{E}, \bar{k}) \xrightarrow{\text{query}} \text{GroupBy}(Q, \bar{A}, \rho_{\bar{k}}(\bar{E}'), \tau)} \quad (\text{Q-Agg})$
$\frac{\Phi_{\text{sdt}}, \Psi_R \vdash Q \xrightarrow{\text{query}} Q' \quad \Phi_{\text{sdt}}, \Psi_R \vdash k \xrightarrow{\text{expr}} a}{\Phi_{\text{sdt}}, \Psi_R \vdash \text{OrderBy}(Q, k, b) \xrightarrow{\text{query}} \text{OrderBy}(Q', a, b)} \quad (\text{Q-OrderBy})$	
$\frac{\Phi_{\text{sdt}}, \Psi_R \vdash Q_1 \xrightarrow{\text{query}} Q'_1 \quad \Phi_{\text{sdt}}, \Psi_R \vdash Q_2 \xrightarrow{\text{query}} Q'_2}{\Phi_{\text{sdt}}, \Psi_R \vdash \text{Union}(Q_1, Q_2) \xrightarrow{\text{query}} Q'_1 \cup Q'_2} \quad (\text{Q-Union})$	$\frac{\Phi_{\text{sdt}}, \Psi_R \vdash Q_1 \xrightarrow{\text{query}} Q'_1 \quad \Phi_{\text{sdt}}, \Psi_R \vdash Q_2 \xrightarrow{\text{query}} Q'_2}{\Phi_{\text{sdt}}, \Psi_R \vdash \text{UnionAll}(Q_1, Q_2) \xrightarrow{\text{query}} Q'_1 \uplus Q'_2} \quad (\text{Q-UnionAll})$

Fig. 16. Translation rules for queries.

now describe a core subset of our syntax-directed transpilation rules, with further details available in Appendix B.

Translating queries. The translation rules for queries, illustrated in Figure 16, use judgments of the form $\Phi_{\text{sdt}}, \Psi_R \vdash Q \xrightarrow{\text{query}} Q'$, where a Cypher query Q maps to an SQL query Q' given SDT Φ_{sdt} and induced database transformer Ψ_R . Among these rules, handling Return requires particular attention, as it involves distinguishing between cases with and without aggregate functions. If there are no aggregation functions in \bar{E} , the Q-Ret rule produces a straightforward translation: the Cypher query $\text{Return}(C, \bar{E}, \bar{k})$ becomes a simple SQL projection $\Pi_{\rho_{\bar{k}}(\bar{E}')} (Q)$. Here, Q is the translated result of the Cypher clause C , and \bar{E}' represents the translated expressions of \bar{E} , with all attributes renamed to \bar{k} . In contrast, when aggregation functions appear in \bar{E} , the translation shifts to the Q-Agg rule, which requires generating a GroupBy query. This is necessary because SQL uses GroupBy to manage aggregation by partitioning rows based on non-aggregated columns. In the Cypher query $\text{Return}(C, \bar{E}, \bar{k})$, the non-aggregation expressions \bar{A} act as grouping keys, while the aggregated expressions ensure the correct computation of results for each group.

Example 5.3. Consider the following Cypher query and the SDT from Example 5.2:

`Return(Match([(n, EMP), (e, WORK_AT, →), (m, DEPT)], τ), [m.dname, Count(n.id)], [name, num])`

Since there is a Count aggregation in the return query, we apply the Q-Agg rule to translate it to a GroupBy query in SQL. Specifically, we first apply the C-Match1 rule to translate the match clause `Match([(n, EMP), (e, WORK_AT, →), (m, DEPT)], τ)` into a SQL query Q (explained in Example 5.4). Then we translate the returned Cypher expressions `m.dname` and `Count(n.id)` to their corresponding SQL expressions `m.dname` and `Count(n.id)`. Among these expressions, we find the one that does not contain aggregations, namely `m.dname`, and use it as the grouping key for GroupBy. Since there is no filtering based on the aggregated results, the Having clause for GroupBy is not generated. Therefore, the translated SQL query is `GroupBy(Q, [m.dname], $\rho_{[\text{name}, \text{num}]}([m.dname, \text{Count}(n.id)])$, τ)`.

Translating Clauses. Unlike translating entire queries, translating individual clauses requires tracking additional information about the node and edge variables used within the clauses. This is necessary to ensure that multiple occurrences of the same variable across different clauses are translated to refer to the same tuple in SQL. As shown in Figure 17, our translation judgments are of the form $\Phi_{\text{sdt}}, \Psi_R \vdash C \xrightarrow{\text{clause}} X, Q'$, meaning a Cypher clause C is translated into a SQL query Q' , and X is the set of all used node and edge variables.

Our translation rules for the Match and OptMatch clauses are based on the observation that graph pattern matching in Cypher can be emulated using sequences of join operations in SQL.

$$\begin{array}{c}
\frac{\Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} \mathcal{X}, Q \quad \Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi'}{\Phi_{\text{sdt}}, \Psi_R \vdash \text{Match}(PP, \phi) \xrightarrow{\text{clause}} \mathcal{X}, \sigma_{\phi'}(Q)} \text{ (C-Match1)} \quad \frac{\Phi_{\text{sdt}}, \Psi_R \vdash C \xrightarrow{\text{clause}} \mathcal{X}, \Pi_L(Q) \quad \mathcal{X}' = \mathcal{X} \setminus \bar{Y} \cup \bar{Z}}{\Phi_{\text{sdt}}, \Psi_R \vdash \text{With}(C, \bar{Y}, \bar{Z}) \xrightarrow{\text{clause}} \mathcal{X}', \Pi_{\rho_{L[\bar{Z}/\bar{Y}]}(L)}(Q)} \text{ (C-With)} \\
\\
\frac{\Phi_{\text{sdt}}, \Psi_R \vdash C \xrightarrow{\text{clause}} \mathcal{X}_1, Q_1 \quad \Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} \mathcal{X}_2, Q_2 \quad \Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi' \quad \text{fresh } T_1, T_2 \quad \phi'' = \phi' \wedge \wedge_{(X:I) \in \mathcal{X}_1 \cap \mathcal{X}_2} T_1.\text{PK}(R_I) = T_2.\text{PK}(R_I) \text{ where } I(\dots) \rightarrow R_I(\dots) \in \Phi_{\text{sdt}}}{\Phi_{\text{sdt}}, \Psi_R \vdash \text{Match}(C, PP, \phi) \xrightarrow{\text{clause}} \mathcal{X}_1 \cup \mathcal{X}_2, \rho_{T_1}(Q_1) \bowtie_{\phi''} \rho_{T_2}(Q_2)} \text{ (C-Match2)} \\
\\
\frac{\Phi_{\text{sdt}}, \Psi_R \vdash C \xrightarrow{\text{clause}} \mathcal{X}_1, Q_1 \quad \Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} \mathcal{X}_2, Q_2 \quad \Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi' \quad \text{fresh } T_1, T_2 \quad \phi'' = \phi' \wedge \wedge_{(X:I) \in \mathcal{X}_1 \cap \mathcal{X}_2} T_1.\text{PK}(R_I) = T_2.\text{PK}(R_I) \text{ where } I(\dots) \rightarrow R_I(\dots) \in \Phi_{\text{sdt}}}{\Phi_{\text{sdt}}, \Psi_R \vdash \text{OptMatch}(C, PP, \phi) \xrightarrow{\text{clause}} \mathcal{X}_1 \cup \mathcal{X}_2, \rho_{T_1}(Q_1) \bowtie_{\phi''} \rho_{T_2}(Q_2)} \text{ (C-OptMatch)}
\end{array}$$

Fig. 17. Translation rules for clauses.

$$\begin{array}{c}
\frac{I(K_1, \dots, K_n) \rightarrow R_I(K_1, \dots, K_n) \in \Phi_{\text{sdt}}}{\Phi_{\text{sdt}}, \Psi_R \vdash (X, I) \xrightarrow{\text{pattern}} \{(X : I)\}, \rho_X(R_I)} \text{ (PT-Node)} \\
\\
\frac{\Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} \mathcal{X}, Q' \quad (X_3, I_3) = \text{head}(PP) \quad I_1(\dots) \rightarrow R_{I_1}(\dots) \in \Phi_{\text{sdt}} \quad I_2(\dots, \text{fk}_s, \text{fk}_t) \rightarrow R_{I_2}(\dots, \text{fk}_s, \text{fk}_t) \in \Phi_{\text{sdt}} \quad I_3(\dots) \rightarrow R_{I_3}(\dots) \in \Phi_{\text{sdt}} \quad \phi = (R_{I_2}.\text{fk}_s = \text{PK}(R_{I_1})) \quad \phi' = (R_{I_2}.\text{fk}_t = \text{PK}(R_{I_3})) \quad \xi(\Psi_R) \Rightarrow \phi \wedge \phi'}{\Phi_{\text{sdt}}, \Psi_R \vdash (X_1, I_1), (X_2, I_2, d_2), PP \xrightarrow{\text{pattern}} \{(X_1 : I_1), (X_2 : I_2)\} \cup \mathcal{X}, \rho_{X_1}(R_{I_1}) \bowtie_{\phi} \rho_{X_2}(R_{I_2}) \bowtie_{\phi'} Q'} \text{ (PT-Path)}
\end{array}$$

Fig. 18. Translation rules for path patterns.

Specifically, the join operations for Match are inner joins, while those for OptionalMatch are left outer joins. Intuitively, a Match clause returns no results if there is no matching pattern, mirroring the behavior of inner joins where unmatched tuples are discarded. In contrast, an OptionalMatch clause returns null for missing matches, similar to how outer joins include unmatched rows with null values.

For example, consider the clause $\text{Match}(C, PP, \phi)$. The translation rule C-Match2 first translates the preceding clause C into a subquery Q_1 and the path pattern PP into another subquery Q_2 . It also collects the sets of node and edge variables used in C and PP , denoted as \mathcal{X}_1 and \mathcal{X}_2 , respectively. For each common variable X with label I , the translation generates a join predicate $T_1.K_1 = T_2.K_1$. This ensures that occurrences of the same variable in different parts of the clause are correctly matched by joining on their primary keys, effectively referring to the same tuple in the SQL translation.

The C-With rule handles the With clause by translating $\text{With}(C, \bar{Y}, \bar{Z})$ into a renaming operation in SQL. It generates a query $\Pi_{\rho_{L[\bar{Z}/\bar{Y}]}(L)}(Q)$, which projects and renames columns, replacing the old names \bar{Y} with the new names \bar{Z} .

Example 5.4. Consider the Cypher clause $\text{Match}([(n, \text{EMP}), (e, \text{WORK_AT}), \rightarrow), (m, \text{DEPT})], \top)$ and the Φ_{sdt} from Example 5.2. Based on C-Match1, we use the PT-Path rule to collect all node and edge variables in the pattern, namely $\mathcal{X} = \{(n : \text{EMP}), (e : \text{WORK_AT}), (m : \text{DEPT})\}$, and translate it to a SQL query $\rho_n(\text{emp}) \bowtie_{n.\text{id}=e.\text{SRC}} \rho_e(\text{work_at}) \bowtie_{e.\text{TGT}=m.\text{dnum}} \rho_m(\text{dept})$. Thus, the Cypher clause is translated to $\sigma_{\top}(\rho_n(\text{emp}) \bowtie_{n.\text{id}=e.\text{SRC}} \rho_e(\text{work_at}) \bowtie_{e.\text{TGT}=m.\text{dnum}} \rho_m(\text{dept}))$.

Example 5.5. Consider the Cypher clause $\text{OptMatch}(C, PP, \phi)$ where C is the Match clause in Example 5.4, $PP = [(m, \text{DEPT})]$ and the Φ_{sdt} from Example 5.2. Based on C-Match1 and C-OptMatch, we know $\mathcal{X}_1 = \{(n : \text{EMP}), (e : \text{WORK_AT}), (m : \text{DEPT})\}$ and $\mathcal{X}_2 = \{(m : \text{DEPT})\}$. C and PP are translated to $Q_1 = \sigma_{\top}(\rho_n(\text{emp}) \bowtie_{n.\text{id}=e.\text{SRC}} \rho_e(\text{work_at}) \bowtie_{e.\text{TGT}=m.\text{dnum}} \rho_m(\text{dept}))$ and $Q_2 = \sigma_{\phi}(\rho_m(\text{dept}))$, respectively. Since there is a shared node (m, DEPT) between C and PP , and the primary key of dept is dnum , the Cypher clause is translated to $\rho_{T_1}(Q_1) \bowtie_{T_1.\text{dnum}=T_2.\text{dnum}} \rho_{T_2}(Q_2)$.

Algorithm 2 Algorithm for inferring the residual of database transformers and SDT's

```

1: procedure REDUCETOSQL( $\Psi_R, Q_R, \Psi'_R, Q'_R, \Phi, \Phi_{\text{sdt}}$ )
   Input: Database transformer  $\Phi$ , standard database transformer  $\Phi_{\text{sdt}}$ 
   Output:  $\top$  for equivalence or  $\perp$  indicating failure
2:    $\sigma \leftarrow \{P_1 \mapsto P_0 \mid P_1(\dots) \rightarrow P_0(\dots) \in \Phi_{\text{sdt}}\}$ 
3:    $\Phi_{\text{rdt}} \leftarrow \Phi[\sigma]$ 
4:   return CheckSQL( $\Psi_R, Q_R, \Psi'_R, Q'_R, \Phi_{\text{rdt}}$ )

```

Translating patterns. Figure 18 illustrates the rules for translating patterns, using judgments of the form $\Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} \mathcal{X}, Q'$. This notation indicates that a Cypher pattern PP translates to an SQL query Q' , with all node and edge variables in the pattern represented by \mathcal{X} .

The translation of patterns follows an inductive structure, with two key rules. The base case handles a single node pattern using the PT-Node rule, which maps the node variable X to its corresponding table R'_l in the relational schema derived from SDT Φ_{sdt} , renaming it to X . The inductive case, represented by the PT-Path rule, addresses more complex patterns where a new node (X_2, l_2) expands an existing sub-pattern PP . The rule identifies the connecting node (X_3, l_3) in PP and determines the appropriate joins. It locates the tables $R'_{l_1}, R'_{l_2}, R'_{l_3}$ corresponding to nodes X_1, X_2 , and X_3 , and constructs join predicates to connect the foreign keys fk_s and fk_t in the edge table R'_{l_2} to the primary keys in the source and target node tables, respectively. This approach aligns with the observation that Cypher's pattern matching naturally corresponds to a series of SQL join operations.

Example 5.6. Given the standard database transformer Φ in Example 5.2, let us focus on the path pattern $[(n, \text{EMP}), (e, \text{WORK_AT}, \rightarrow), (m, \text{DEPT})]$. According to the PT-Path rule, we first need to apply the PT-Node rule to get the variables $\{(m : \text{DEPT})\}$ and query $\rho_m(\text{DEPT})$ from the node pattern (m, DEPT) . Based on these results, we can use the PT-Path rule to collect all variables $\mathcal{X} = \{(n : \text{EMP}), (e : \text{WORK_AT}), (m : \text{DEPT})\}$ and obtain the SQL query $\rho_n(\text{emp}) \bowtie_{n.\text{id}=e.\text{SRC}} \rho_e(\text{work_at}) \bowtie_{e.\text{TGT}=m.\text{dnum}} \rho_m(\text{dept})$.

THEOREM 5.7 (SOUNDNESS OF TRANSLATION). *Let Ψ_G be a graph schema and Q be a Cypher query over Ψ_G . Let Ψ_R be the induced relational schema of Ψ_G , and Φ_{sdt} be the standard database transformer from Ψ_G to Ψ_R . If $\Phi_{\text{sdt}}, \Psi_R \vdash Q \xrightarrow{\text{query}} Q'$, then Q' is equivalent to Q modulo Φ_{sdt} , i.e., $Q \simeq_{\Phi_{\text{sdt}}} Q'$.*

THEOREM 5.8 (COMPLETENESS OF TRANSLATION). *Let Ψ_G be a graph schema and Ψ_R be the induced relational schema of Ψ_G . Given any Cypher query Q over Ψ_G accepted by the grammar shown in Figure 9, there exists a SQL query Q' over Ψ_R such that $\Phi_{\text{sdt}}, \Psi_R \vdash Q \xrightarrow{\text{query}} Q'$.*

5.3 Reduction to SQL Equivalence Checking

The final step of our algorithm utilizes the transpiled query to reduce our original problem to checking equivalence between a pair of SQL queries. As shown in Algorithm 2, the REDUCETOSQL procedure first infers the residual database transformer Φ_{rdt} through a simple syntactic substitution: Since every clause of the SDT is of the form $P_1(\dots) \rightarrow P_0(\dots)$, we can obtain the residual transformer simply by substituting every occurrence of P_1 in Φ by P_0 . Finally, since the residual transformer Φ_{rdt} specifies how to convert instances of the induced schema to those of the desired schema, we can use an existing tool for checking SQL equivalence by utilizing Φ_{rdt} . As stated by the following theorems, the original Cypher query is equivalent to the given SQL query if and only if Q_R and Q'_R are equivalent modulo Φ_{rdt} .

Table 1. Statistics of Cypher and SQL queries in the benchmarks. Sizes are the number of AST nodes.

Dataset	#	SQL Size				Cypher Size				Transformer Size			
		Min	Max	Avg	Med	Min	Max	Avg	Med	Min	Max	Avg	Med
StackOverflow	12	15	74	32.5	28	20	149	54.9	41	1	6	3.3	4
Tutorial	26	5	76	25.8	22	12	77	31.2	28	1	17	8.3	5
Academic	7	27	59	46.9	54	45	121	75.0	66	5	7	6.7	7
VeriEQL	60	15	143	42.2	39	29	174	65.8	61	1	10	6.7	10
Mediator	100	9	63	18.6	13	20	114	33.8	28	1	11	5.1	4
GPT-Translate	205	5	143	28.2	25	12	171	46.0	38	1	17	5.9	5
Total	410	5	143	28.2	25	12	174	45.7	38	1	17	5.9	5

THEOREM 5.9 (SOUNDNESS). *Let $\text{CheckSQL}(\Psi_R, Q, \Psi'_R, Q', \Phi_{\text{rdt}})$ be a sound procedure for equivalence checking of SQL queries Q, Q' over relational schemas Ψ_R, Ψ'_R connected by RDT Φ_{rdt} . Given a Cypher query Q_G over graph schema Ψ_G , a SQL query Q_R over relational schema Ψ_R , and their database transformer Φ , if $\text{CHECKEQUIVALENCE}(\Psi_G, Q_G, \Psi_R, Q_R, \Phi)$ returns \top , it holds that $Q_G \simeq_\Phi Q_R$.*

THEOREM 5.10 (COMPLETENESS). *Let $\text{CheckSQL}(\Psi_R, Q, \Psi'_R, Q', \Phi_{\text{rdt}})$ be a complete procedure for equivalence checking of SQL queries Q, Q' over schemas Ψ_R, Ψ'_R connected by RDT Φ_{rdt} . Given a Cypher query Q_G over graph schema Ψ_G , a SQL query Q_R over relational schema Ψ_R , and their database transformer Φ , if $Q_G \simeq_\Phi Q_R$, then $\text{CHECKEQUIVALENCE}(\Psi_G, Q_G, \Psi_R, Q_R, \Phi)$ returns \top .*

6 Evaluation

In this section, we describe three experiments to evaluate GRAPHITI. Because GRAPHITI's verification methodology reduces the Cypher-SQL equivalence checking problem to pure SQL, our results depend on what backend is used for SQL equivalence checking. Thus, in our first experiment, we evaluate GRAPHITI using the VERIEQL [24] bounded model checker as its backend, and in our second experiment, we use a deductive verifier called MEDIATOR [55] as the backend. Finally, we also evaluate the quality of GRAPHITI's transpilation results.

Benchmarks. We evaluate GRAPHITI on 410 pairs of SQL and Cypher queries (see Table 1) from the following sources:

- **StackOverflow:** We identified 12 StackOverflow posts where users inquire about translating a SQL query to Cypher or vice versa. All of these posts contain a description of the schemas and SQL/Cypher queries that are intended to be semantically equivalent.
- **Tutorial.** We identified 26 tutorial examples, including from the official Neo4j guide [36], that explain how a SQL query can be implemented using the Cypher query language.
- **Academic.** We identified 7 examples from academic papers [4, 30] that contain relational queries and their corresponding version in Cypher.
- **VeriEQL.** We collected 60 benchmarks from the VERIEQL paper [24] by randomly sampling 20 queries from each of its three datasets and asking various people with at least 3 months of Cypher experience to write an equivalent Cypher query.
- **Mediator.** We collected 100 benchmarks from the MEDIATOR evaluation set [55]. Each MEDIATOR benchmark, consisting of an SQL query pair (Q_1, Q_2) over schemas Ψ_1 and Ψ_2 , was translated into pairs (G_1, Q_2) and (G_2, Q_1) where G_1, G_2 are Cypher queries over graph schemas Ψ'_1 and Ψ'_2 , and Ψ_1 and Ψ_2 are the induced relational schemas for Ψ'_1 and Ψ'_2 .
- **GPT-Translate.** Given that large language models like GPT are increasingly used by people for coding and transpilation tasks, we included GPT-generated Cypher queries to assess GRAPHITI's ability to detect errors in automated translations. Specifically, we used GPT to transpile SQL queries from the previous five categories, yielding an additional 205 benchmarks.

Table 2. Results of bounded equivalence checking.

Dataset	#	# Non-Equiv	Avg Checked Bound	Avg Refutation Time (s)
StackOverflow	12	1	9.2	0.6
Tutorial	26	1	7.7	56.2
Academic	7	1	2.5	5.4
VeriEQL	60	4	7.2	8.5
Mediator	100	0	33.2	N/A
GPT-Translate	205	27	18.7	25.9
Total	410	34	19.6	23.4

Database transformers. Since the induced schema of graph databases may differ from the schema of relational databases, GRAPHITI requires a database transformer to describe the relationship between the graph and relational schemas. To evaluate GRAPHITI across all pairs of SQL and Cypher queries, one of the authors constructed a database transformer for each query pair based on their schema descriptions. We observe that writing these database transformers is not difficult. As shown in Table 1, each database transformer consists of an average of 5.9 rules, which takes approximately one minute to write.

Machine configuration. All of the experiments are conducted on a laptop with an Intel Core i7-8750H processor and 32GB physical memory running the Debian 12 operating system.

6.1 Evaluation of GRAPHITI with BMC backend

In this section, we present the results of the evaluation in which we use the VERIEQL bounded model checker as GRAPHITI’s SQL equivalence checking backend. VERIEQL is a bounded model checker that requires a hyperparameter specifying the size bound of symbolic tables. However, since it is difficult to estimate a suitable bound a priori, we set a 10-minute time limit and gradually increase the bound until either a counterexample is found or the time-limit is reached. For each refuted benchmark, GRAPHITI uses the relational counterexamples produced by VERIEQL to construct a counterexample over the graph schema.

The results of this evaluation are presented in Table 2. Here, the column labeled “# Non-Equiv” shows the number of benchmarks proven to be *not* equivalent, and the last column shows the average time to find a counterexample. The column labeled “Avg Checked Bound” shows the average size of symbolic tables (measured in terms of the number of rows) when the 10 minute time limit is reached. As shown in Table 2, GRAPHITI refutes equivalence for 34 out of the 410 benchmarks, taking an average of 23.4 seconds to find a counterexample. For the remaining 376 benchmarks, GRAPHITI performs bounded verification, demonstrating that there is no counterexample for database instances with symbolic tables of average size 19.6.

Uncovered bugs. We have manually inspected all 34 bugs uncovered by GRAPHITI and confirmed that all counterexamples produced by the tool correspond to true positives. As expected, GPT-generated queries have a higher probability of being incorrect compared to the human-written queries. In particular, GRAPHITI finds a counterexample to equivalence for 13% of the queries transpiled by GPT. This experiment shows that GPT may introduce semantic bugs when converting SQL to Cypher, and GRAPHITI can effectively identify these bugs. As developers increasingly rely on large language models for assistance with coding-related tasks, we believe this demonstrates GRAPHITI’s practical value for developers relying on LLM-generated queries.

Perhaps more surprisingly, GRAPHITI also finds incorrect translations among benchmarks in the StackOverflow, Tutorial, Academic, and VeriEQL categories, all of which involve queries

Table 3. Results of full equivalence verification.

Dataset	#	# Supported	# Verified	# Unknown	Avg Time (s)
StackOverflow	12	1	1	0	1.0
Tutorial	26	1	1	0	0.2
Academic	7	0	0	0	N/A
VeriEQL	60	0	0	0	N/A
Mediator	100	100	77	23	20.5
GPT-Translate	205	94	73	21	23.5
Total	410	196	152	44	16.8

transpiled by experts. Most surprisingly, GRAPHITI uncovers a bug in an example from a Neo4j tutorial [36] that is intended to help developers familiar with SQL to learn Cypher. This tutorial contains several pairs of SQL and Cypher queries that are intended to be equivalent, but, for one of these examples, GRAPHITI finds a counterexample on which the query results are actually different. We refer the interested reader to Appendix D for a case study explaining some of the uncovered bugs, including the example from the Neo4j tutorial.

False negatives. Since VERIEQL is a bounded model checker, benchmarks that are not refuted by GRAPHITI within the 10-minute time limit *may* still contain bugs. To assess how frequently this occurs, we sampled 50 pairs of queries that were not refuted by GRAPHITI and manually inspected whether the translation was correct. For 48 of the 50 manually-inspected benchmarks, we found that the translation is indeed correct, but for 2 benchmarks (both from the GPT-Translate category), the translation is incorrect but GRAPHITI fails to find a counterexample within the 10 minute time limit. Thus, while GRAPHITI with the VERIEQL backend does not provide theoretical soundness guarantees, we find that it is useful for finding bugs and has a low chance of missing incorrect translations (4% according to our manual inspection results).

Key finding: Using a bounded model checker backend, GRAPHITI identifies 27 bugs among 205 SQL queries transpiled to Cypher using GPT. More surprisingly, among the 205 manually-written query pairs that are meant to be equivalent, GRAPHITI also identifies 7 inconsistencies, including 3 benchmarks from the wild and 4 benchmarks from manual translations.

6.2 Evaluation of GRAPHITI with Deductive Verifier

In this section, we present the results of a second experiment wherein we evaluate GRAPHITI with MEDIATOR as its backend. As mentioned earlier, MEDIATOR is an SMT-based deductive verifier for reasoning about SQL applications over different schemas. Unlike VERIEQL, MEDIATOR can perform full-fledged verification; however, it supports a limited subset of SQL without aggregations or outer joins. Additionally, unlike VERIEQL, MEDIATOR cannot disprove equivalence by generating counterexamples. Hence, when using GRAPHITI with MEDIATOR as its backend, GRAPHITI can either prove equivalence or it returns “Unknown”. For performing this experiment, we also use a time limit of 10 minutes per benchmark.

The results of this experiment are summarized in Table 3. As shown in the “# Supported” column, about half of the benchmarks (196 out of 410) fall inside the fragment of SQL supported by MEDIATOR, so we conduct our evaluation on this subset. Overall, GRAPHITI can verify 77.6% of these 196 benchmarks, with an average running time of 16.8 seconds. Since the syntax-directed transpilation performed by GRAPHITI takes negligible time, most of the verification time is dominated by SMT queries for discharging the generated verification conditions.

Qualitative analysis. We manually inspected 44 benchmarks that cannot be verified by GRAPHITI with the MEDIATOR backend. Among these 44 benchmarks, two of them are in fact *refuted* by

Table 4. Execution time of transpiled and manually-written SQL queries.

Dataset	#	Avg Exec Time (s)		% Transpiled Faster	% Trans Slower (1x, 1.1x]	% Trans Slower (1.1x, 1.2x]	% Trans Slower (1.2x, +∞)
		Transpiled	Manual				
StackOverflow	12	1.9	1.8	41.7%	8.3%	50.0%	0.0%
Tutorial	26	2.3	1.7	19.2%	11.5%	46.2%	23.1%
Academic	7	2.4	3.2	71.4%	0.0%	14.3%	14.3%
Total	45	2.2	2.0	33.3%	8.9%	42.2%	15.6%

GRAPHITI with VERIEQL, so they should *not* be verified. Among the 42 remaining benchmarks, MEDIATOR fails to complete verification within the 10 minute time limit for 14 of these, and, for the final 28 benchmarks, it terminates but returns “Unknown”. To gain insight about failure cases, note that MEDIATOR needs to infer an *inductive bisimulation invariant* between the two queries [55]. However, for queries involving long join chains, the corresponding bisimulation invariant can be complex. This can either lead to expensive SMT queries, thereby causing time-outs, or the required invariant might fall outside the inference capabilities of MEDIATOR.

Key finding: Among the 196 SQL queries supported by MEDIATOR, GRAPHITI can prove equivalence between roughly 80% of (Cypher, SQL) query pairs using the MEDIATOR backend.

6.3 Evaluation of Transpilation

While the primary goal of this work is to enable checking equivalence between graph and relational queries, an additional benefit of our method is that it can transpile graph database queries to relational queries over the induced schema. To assess the practical effectiveness, we conduct an experiment to evaluate how well our method transpiles graph queries into *efficient* SQL queries.

Efficiency of transpilation. First, we evaluate how long GRAPHITI takes to transpile each Cypher query to a SQL query over the induced schema. GRAPHITI can successfully transpile all 410 queries, and the average, median, and maximum transpilation times are 6.3, 3.0, and 180.2 milliseconds respectively. Hence, we can conclude that transpilation is very fast in practice.

Quality of transpiled queries. Next, we also set out to evaluate the quality of the transpiled queries by comparing the execution time of manually written SQL queries against GRAPHITI’s transpilation result. However, performing this evaluation is challenging for two reasons: First, we only have access to the “ground truth” Cypher and SQL queries for some benchmark categories. Second, we do not have database instances that these queries are meant to be executed on. To deal with the first challenge, we conduct this evaluation only on those benchmarks from the StackOverflow, Tutorial, and Academic categories for which we are given the original Cypher query and its SQL equivalent (or vice versa). To deal with the second challenge, we generate mock database instances and assess query efficiency on them. For each benchmark with SQL query Q_R over schema Ψ_R and Cypher query Q_G , we use GRAPHITI to transpile Q_G into SQL query Q'_R over the induced schema Ψ'_R . We then create databases R and R' over Ψ_R and Ψ'_R , respectively, ensuring $\Phi_{\text{rdt}}(R') = R$. To account for execution time variability, we start with 10,000 tuples in each table of R and iteratively increase the table size by 10x, up to 1 million, choosing the largest size where manually written SQL queries run within 10 seconds. This approach results in 1 million tuples for 36 benchmarks and between 10,000 and 1 million for the remaining 9 benchmarks. Finally, we measure the execution times of Q'_R on R' and Q_R on R . Table 4 summarizes the results: for 33.3% of the benchmarks, the transpiled queries are faster than the manually written queries. For the remaining benchmarks, 8.9% exhibit a slowdown of no more than 1.1x, 42.2% exhibit a slowdown between 1.1x and 1.2x, and 15.6% exceed 1.2x. These results indicate that using GRAPHITI to perform automated transpilation

could be useful in real-world scenarios that necessitate graph-to-relational query conversion, such as for legacy systems, resource-constrained environments, or data integration.

Key finding: GRAPHITI can transpile Cypher queries to SQL queries in milliseconds. The execution time of transpiled SQL queries is faster than manually-written queries on 33.3% of benchmarks and within 1.2x slowdown on 51.1% of benchmarks.

7 Related Work

In this section, we discuss prior work that is mostly related to our techniques for equivalence verification between Cypher and SQL queries.

Automated reasoning for SQL. Despite the undecidability of checking equivalence between SQL queries [49], there has been much prior work on automated reasoning for relational queries. We can categorize existing work into three classes. The first class targets a decidable subset of SQL and proposes decision procedures for that subset. Examples of work in this category include [2, 5, 21]. Approaches in the second category propose sound but incomplete algorithms for an undecidable subset of SQL; examples of work in this space include [10, 12, 61, 62]. The third category performs bounded verification to find bugs in SQL queries; examples include COSETTE [11], QEX [51], and VERIEQL [24]. There is also prior work on verifying relational database applications that involve both queries and updates [55]. Our proposed approach reduces the verification problem between graph and relational queries to checking equivalence between a pair of relational queries; as such, it can leverage any future advances in this area.

Migration between database instances. There is prior work on migrating data between different schemas, including [18, 27, 32, 53, 57–59]. The most related to this paper is DYNAMITE [57], which automates data migration between graph and relational databases. However, DYNAMITE is only useful for migrating the *contents* of the database and cannot be used for transpiling queries. There are also prior papers that address the query transpilation problem in the context of SQL [13, 14, 56]. While our notion of database transformer is inspired by prior work [35, 57], to the best of our knowledge, this paper is the first to formalize the transpilation procedure from Cypher to SQL queries and leverage it for formal equivalence checking.

Data representation refactoring. There is a related line of work on *data representation refactoring*, which aims to refactor programs or specifications from one data representation to another [8, 9, 15, 20, 28, 41–43, 54]. For instance, Solidare [41] refactors smart contracts between different ADTs. QBS [8] converts Java programs operating over collections to SQL queries. Since graph and relational schemas can be viewed as different data representations, the query transpilation problem in this work can be viewed as a form of data representation refactoring. However, none of the existing techniques addresses the query transpilation problem between graph and relational data. Additionally, this work can be viewed as presenting a novel methodology for verifying data representation refactoring: Rather than directly going from representation R to representation R' , our idea is to introduce an auxiliary representation R'' that simplifies the problem by enabling syntax-directed translation. To the best of our knowledge, such a methodology based on a layer of indirection has not previously been explored in this context.

Graph database query languages. There has been a long line of prior work on graph databases and semantic foundations of graph query languages [1, 16, 19, 22, 23, 46, 50]. The unifying insight behind many of such works is that a graph database schema may be viewed as a graphical representation of the Entity-Relationship Diagram (ER Diagram) of a relational database schema. In part due to this unifying insight, these graph query languages are both semantically and syntactically similar

to Cypher. Because of this similarity, we believe our proposed methodology can be adapted fairly easily to graph database query languages other than Cypher.

Testing database queries. Another related line of related work focuses on testing database queries. Work in this space includes differential testing [25, 33, 60] and metamorphic testing [7, 26, 31, 44] to detect bugs in database management systems, mutation-based testing to grade programming assignments involving SQL queries [6], and provenance-based techniques for explaining wrong SQL queries [34]. In contrast, GRAPHITI transpiles graph database queries into equivalent SQL queries and uses existing automated reasoning tools for SQL equivalence checking. Thus, GRAPHITI is complementary and can benefit from advances in testing SQL queries.

Transpiling Cypher queries. There are a few existing tools that can translate Cypher queries to queries in SQL-like languages [29, 48]. The most relevant to this paper is `OPENCYPHERTRANSPILER` [29], which first transforms a Cypher query into a logical plan and renders it as a relational query. However, it supports only a limited subset of Cypher and lacks soundness guarantees for translated queries.⁵ In contrast, GRAPHITI ensures soundness during transpilation and supports a broader subset of Cypher queries. Another tool, `KUZU` [48], can execute graph queries on relational databases with a Cypher interface. It compiles Cypher queries into an intermediate representation similar to a relational database’s logical plan. However, `KUZU` does not transpile Cypher into SQL but instead directly executes the intermediate representation on the database.

8 Limitation

The current version of GRAPHITI is focused on a specific subset of SQL and Cypher, which does not yet cover all modern features, such as variable-length pattern matching in Cypher. However, considering the lack of prior research on reasoning about equivalence between graph and relational database queries, we believe our selected fragments offer a strong foundation for advancing this area of study. While some SQL and Cypher queries lie outside the scope of our current subset, evaluations on a diverse set of benchmarks, including real-world queries, demonstrate that this subset is expressive enough for practical use cases. Future work can further extend the transpilation rules and backend equivalence verifiers to support additional features.

9 Conclusion and Future Work

In this paper, we proposed automated reasoning techniques between graph and relational database queries. Specifically, we first proposed a formal definition of equivalence between graph and relational queries and used it as the basis of a correct-by-construction transpilation strategy for converting Cypher queries to SQL queries over a so-called *induced relational schema*. We then showed how our translation approach can be used to check equivalence between graph and SQL queries over *arbitrary* schema by leveraging existing automated reasoning techniques for SQL. We have also evaluated our implementation, GRAPHITI, on equivalence checking tasks involving real-world Cypher and SQL queries and showed that GRAPHITI can be useful for (a) uncovering subtle bugs in Cypher queries that are meant to be equivalent to a reference SQL implementation, and (b) verifying full equivalence between Cypher and SQL queries.

Looking ahead, we plan to explore the development of a graphical interface for specifying database transformers between graph and relational databases, inspired by prior work on schema mapping visualization [45]. This interface would aim to further reduce the manual effort required by users to verify equivalence between graph and relational queries, providing a more intuitive and user-friendly approach. We see this as a promising avenue for future research.

⁵The detailed evaluation of `OPENCYPHERTRANSPILER` can be found in Appendix E.

References

- [1] AGE. 2024. Apache AGE. <https://age.apache.org>.
- [2] Alfred V. Aho, Yehoshua Sagiv, and Jeffrey D. Ullman. 1979. Equivalences Among Relational Expressions. *SIAM J. Comput.* 8, 2 (1979), 218–246. <https://doi.org/10.1137/0208017>
- [3] Amazon. 2024. What's the Difference Between a Graph Database and a Relational Database? https://aws.amazon.com/compare/the-difference-between-graph-and-relational-database/?nc1=h_ls.
- [4] Abdelkrim Boudaoud, Houari Mahfoud, and Azeddine Chikh. 2022. Towards a Complete Direct Mapping from Relational Databases to Property Graphs. In *Proceedings of the International Conference on Model and Data Engineering (MEDI)*. 222–235. https://doi.org/10.1007/978-3-031-21595-7_16
- [5] Ashok K. Chandra and Philip M. Merlin. 1977. Optimal Implementation of Conjunctive Queries in Relational Data Bases. In *Proceedings of the ACM Symposium on Theory of Computing (STOC)*. 77–90. <https://doi.org/10.1145/800105.803397>
- [6] Bikash Chandra, Bhupesh Chawda, Biplab Kar, K. V. Maheshwara Reddy, Shetal Shah, and S. Sudarshan. 2015. Data generation for testing and grading SQL queries. *VLDB J.* 24, 6 (2015), 731–755. <https://doi.org/10.1007/S00778-015-0395-0>
- [7] Tsong Yueh Chen, S. C. Cheung, and Siu-Ming Yiu. 2020. Metamorphic testing: a new approach for generating next test cases. *CoRR* abs/2002.12543 (2020). arXiv:2002.12543 <https://arxiv.org/abs/2002.12543>
- [8] Yanju Chen, Yuepeng Wang, Maruth Goyal, James Dong, Yu Feng, and Isil Dillig. 2022. Synthesis-powered optimization of smart contracts via data type refactoring. *Proc. ACM Program. Lang.* 6, OOPSLA2 (2022), 560–588. <https://doi.org/10.1145/3563308>
- [9] Alvin Cheung, Armando Solar-Lezama, and Samuel Madden. 2013. Optimizing database-backed applications with query synthesis. In *ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. ACM, 3–14. <https://doi.org/10.1145/2491956.2462180>
- [10] Shumo Chu, Brendan Murphy, Jared Roesch, Alvin Cheung, and Dan Suciu. 2018. Axiomatic Foundations and Algorithms for Deciding Semantic Equivalences of SQL Queries. *Proc. VLDB Endow.* 11, 11 (2018), 1482–1495. <https://doi.org/10.14778/3236187.3236200>
- [11] Shumo Chu, Chenglong Wang, Konstantin Weitz, and Alvin Cheung. 2017. Cosette: An Automated Prover for SQL. In *Proceedings of the Biennial Conference on Innovative Data Systems Research (CIDR)*. <http://cidrdb.org/cidr2017/papers/p51-chu-cidr17.pdf>
- [12] Shumo Chu, Konstantin Weitz, Alvin Cheung, and Dan Suciu. 2017. HoTSQL: proving query rewrites with univalent SQL semantics. In *Proceedings of the SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. 510–524. <https://doi.org/10.1145/3062341.3062348>
- [13] Carlo Curino, Hyun Jin Moon, Alin Deutsch, and Carlo Zaniolo. 2013. Automating the database schema evolution process. *VLDB J.* 22, 1 (2013), 73–98. <https://doi.org/10.1007/s00778-012-0302-x>
- [14] Carlo Curino, Hyun Jin Moon, and Carlo Zaniolo. 2008. Graceful database schema evolution: the PRISM workbench. *Proc. VLDB Endow.* 1, 1 (2008), 761–772. <https://doi.org/10.14778/1453856.1453939>
- [15] Benjamin Delaware, Clément Pit-Claudel, Jason Gross, and Adam Chlipala. 2015. Fiat: Deductive Synthesis of Abstract Data Types in a Proof Assistant. In *Proceedings of the ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*. 689–700. <https://doi.org/10.1145/2676726.2677006>
- [16] Alin Deutsch, Nadime Francis, Alastair Green, Keith Hare, Bei Li, Leonid Libkin, Tobias Lindaaker, Victor Marsault, Wim Martens, Jan Michels, Filip Murlak, Stefan Plantikow, Petra Selmer, Oskar van Rest, Hannes Voigt, Domagoj Vrgoc, Mingxi Wu, and Fred Zemke. 2022. Graph Pattern Matching in GQL and SQL/PGQ. In *International Conference on Management of Data (SIGMOD)*. 2246–2258. <https://doi.org/10.1145/3514221.3526057>
- [17] Ronald Fagin, Laura M. Haas, Mauricio A. Hernández, Renée J. Miller, Lucian Popa, and Yannis Velegrakis. 2009. Clio: Schema Mapping Creation and Data Exchange. In *Conceptual Modeling: Foundations and Applications - Essays in Honor of John Mylopoulos (Lecture Notes in Computer Science, Vol. 5600)*. 198–236. https://doi.org/10.1007/978-3-642-02463-4_12
- [18] Yu Feng, Ruben Martins, Jacob Van Geffen, Isil Dillig, and Swarat Chaudhuri. 2017. Component-based synthesis of table consolidation and transformation tasks from examples. In *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. 422–436. <https://doi.org/10.1145/3062341.3062351>
- [19] Nadime Francis, Alastair Green, Paolo Guagliardo, Leonid Libkin, Tobias Lindaaker, Victor Marsault, Stefan Plantikow, Mats Rydberg, Petra Selmer, and Andrés Taylor. 2018. Cypher: An Evolving Query Language for Property Graphs. In *Proceedings of the 2018 International Conference on Management of Data (SIGMOD)*. 1433–1445. <https://doi.org/10.1145/3183713.3190657>
- [20] Xi Ge, Quinton L. DuBose, and Emerson R. Murphy-Hill. 2012. Reconciling manual and automatic refactoring. In *International Conference on Software Engineering (ICSE)*. 211–221. <https://doi.org/10.1109/ICSE.2012.6227192>
- [21] Todd J. Green. 2011. Containment of Conjunctive Queries on Annotated Relations. *Theory Comput. Syst.* 49, 2 (2011), 429–459. <https://doi.org/10.1007/S00224-011-9327-6>

- [22] Harry Halpin and James Cheney. 2011. Dynamic Provenance for SPARQL Updates Using Named Graphs. In *Workshop on the Theory and Practice of Provenance (TaPP)*. <https://doi.org/10.1145/2567948.2577357>
- [23] Olaf Hartig and Jorge Pérez. 2018. Semantics and Complexity of GraphQL. In *Proceedings of the World Wide Web Conference on World Wide Web (WWW)*. 1155–1164. <https://doi.org/10.1145/3178876.3186014>
- [24] Yang He, Pinhan Zhao, Xinyu Wang, and Yuepeng Wang. 2024. VeriEQL: Bounded Equivalence Verification for Complex SQL Queries with Integrity Constraints. *Proc. ACM Program. Lang.* 8, OOPSLA1 (2024), 1071–1099. <https://doi.org/10.1145/3649849>
- [25] Ziyue Hua, Wei Lin, Luyao Ren, Zongyang Li, Lu Zhang, Wenpin Jiao, and Tao Xie. 2023. GDsmith: Detecting Bugs in Cypher Graph Database Engines. In *Proceedings of the SIGSOFT International Symposium on Software Testing and Analysis (ISSTA)*. 163–174. <https://doi.org/10.1145/3597926.3598046>
- [26] Yuancheng Jiang, Jiahao Liu, Jinsheng Ba, Roland HC Yap, Zhenkai Liang, and Manuel Rigger. 2024. Detecting Logic Bugs in Graph Database Management Systems via Injective and Surjective Graph Query Transformation. In *Proceedings of the International Conference on Software Engineering (ICSE)*. <https://doi.org/10.1145/3597503.3623307>
- [27] Zhongjun Jin, Michael R. Anderson, Michael J. Cafarella, and H. V. Jagadish. 2017. Foofah: Transforming Data By Example. In *Proceedings of the ACM International Conference on Management of Data (SIGMOD)*. 683–698. <https://doi.org/10.1145/3035918.3064034>
- [28] Shadaj Laddad, Conor Power, Mae Milano, Alvin Cheung, and Joseph M. Hellerstein. 2022. Katara: synthesizing CRDTs with verified lifting. *Proc. ACM Program. Lang.* 6, OOPSLA2 (2022), 1349–1377. <https://doi.org/10.1145/3563336>
- [29] Jerry Liang. 2025. openCypher Transpiler. <https://github.com/microsoft/openCypherTranspiler>.
- [30] Chunbin Lin, Benjamin Mandel, Yannis Papakonstantinou, and Matthias Springer. 2016. Fast In-Memory SQL Analytics on Relationships between Entities. *CoRR* abs/1602.00033 (2016). arXiv:1602.00033 <http://arxiv.org/abs/1602.00033>
- [31] Qiuyang Mang, Aoyang Fang, Boxi Yu, Hanfei Chen, and Pinjia He. 2024. Testing Graph Database Systems via Equivalent Query Rewriting. In *Proceedings of the International Conference on Software Engineering (ICSE)*. <https://doi.org/10.1145/3597503.3639200>
- [32] Ruben Martins, Jia Chen, Yanju Chen, Yu Feng, and Isil Dillig. 2019. Trinity: An Extensible Synthesis Framework for Data Science. *Proc. VLDB Endow.* 12, 12 (2019), 1914–1917. <https://doi.org/10.14778/3352063.3352098>
- [33] William M. McKeeman. 1998. Differential Testing for Software. *Digit. Tech. J.* 10, 1 (1998), 100–107.
- [34] Zhengjie Miao, Sudeepa Roy, and Jun Yang. 2019. Explaining Wrong Queries Using Small Examples. In *Proceedings of the International Conference on Management of Data (SIGMOD)*. 503–520. <https://doi.org/10.1145/3299869.3319866>
- [35] Renée J. Miller, Laura M. Haas, and Mauricio A. Hernández. 2000. Schema Mapping as Query Discovery. In *Proceedings of the International Conference on Very Large Data Bases (VLDB)*. 77–88. <http://www.vldb.org/conf/2000/P077.pdf>
- [36] Neo4j. 2024. *Cypher for SQL Users*. <https://neo4j.com/docs/getting-started/cypher-intro/cypher-sql/> Accessed: 2024-11-07.
- [37] Neo4j. 2024. Neo4j Docs. https://neo4j.com/docs/getting-started/cypher-intro/cypher-sql/#_return_customers_without_existing_orders.
- [38] Neo4j. 2024. Transition from Relational to Graph Database. <https://neo4j.com/docs/getting-started/appendix/graphdb-concepts/graphdb-vs-rdbms>.
- [39] NLM. 2024. SemMedDB Database Details - Version 2.0. https://lhncbc.nlm.nih.gov/ii/tools/SemRep_SemMedDB_SKR/dbinfo20.html.
- [40] Stack Overflow. 2024. How could i use this SQL on cypher(neo4j). <https://stackoverflow.com/questions/43438214/how-could-i-use-this-sql-on-cypherneo4j>.
- [41] Shankara Pailoor, Yuepeng Wang, and Isil Dillig. 2024. Semantic Code Refactoring for Abstract Data Types. *Proc. ACM Program. Lang.* 8, POPL (2024), 816–847. <https://doi.org/10.1145/3632870>
- [42] Shankara Pailoor, Yuepeng Wang, Xinyu Wang, and Isil Dillig. 2021. Synthesizing data structure refinements from integrity constraints. In *Proceedings of the International Conference on Programming Language Design and Implementation (PLDI)*. 574–587. <https://doi.org/10.1145/3453483.3454063>
- [43] Kia Rahmani, Kartik Nagar, Benjamin Delaware, and Suresh Jagannathan. 2021. Repairing serializability bugs in distributed database programs via automated schema refactoring. In *ACM SIGPLAN International Conference on Programming Language Design and Implementation (PLDI)*. 32–47. <https://doi.org/10.1145/3453483.3454028>
- [44] Manuel Rigger and Zhendong Su. 2020. Finding bugs in database systems via query partitioning. *Proc. ACM Program. Lang.* 4, OOPSLA (2020), 211:1–211:30. <https://doi.org/10.1145/3428279>
- [45] George G. Robertson, Mary Czerwinski, and John E. Churchill. 2005. Visualization of mappings between schemas. In *Proceedings of the Conference on Human Factors in Computing Systems (CHI)*. 431–439. <https://doi.org/10.1145/1054972.1055032>
- [46] Marko A. Rodriguez. 2015. The Gremlin graph traversal machine and language (invited talk). In *Proceedings of the Symposium on Database Programming Languages (DBPL)*. 1–10. <https://doi.org/10.1145/2815072.2815073>

- [47] StackOverflow. 2024. Comparison of relational databases and graph databases. <https://stackoverflow.com/questions/13046442/comparison-of-relational-databases-and-graph-databases>.
- [48] Kuzu Team. 2025. Kuzu. <https://kuzudb.com/>.
- [49] Boris A Trakhtenbrot. 1950. Impossibility of an algorithm for the decision problem in finite classes. In *Doklady Akademii Nauk SSSR* 70. 569–572.
- [50] Oskar van Rest, Sungpack Hong, Jinha Kim, Xuming Meng, and Hassan Chafi. 2016. PGQL: a property graph query language. In *Proceedings of the International Workshop on Graph Data Management Experiences and Systems (GRADES)*. 7. <https://doi.org/10.1145/2960414.2960421>
- [51] Margus Veanes, Nikolai Tillmann, and Jonathan de Halleux. 2010. Qex: Symbolic SQL Query Explorer. In *International Conference on Logic for Programming, Artificial Intelligence, and Reasoning (LPAR)*. 425–446. https://doi.org/10.1007/978-3-642-17511-4_24
- [52] Chad Vicknair, Michael Macias, Zhendong Zhao, Xiaofei Nan, Yixin Chen, and Dawn Wilkins. 2010. A comparison of a graph database and a relational database: a data provenance perspective. In *Proceedings of the Annual Southeast Regional Conference*. 42. <https://doi.org/10.1145/1900008.1900067>
- [53] Chenglong Wang, Alvin Cheung, and Rastislav Bodik. 2017. Synthesizing highly expressive SQL queries from input-output examples. In *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. 452–466. <https://doi.org/10.1145/3062341.3062365>
- [54] Meng Wang, Jeremy Gibbons, Kazutaka Matsuda, and Zhenjiang Hu. 2013. Refactoring pattern matching. *Sci. Comput. Program.* 78, 11 (2013), 2216–2242. <https://doi.org/10.1016/J.SCICO.2012.07.014>
- [55] Yuepeng Wang, Isil Dillig, Shuvendu K. Lahiri, and William R. Cook. 2018. Verifying equivalence of database-driven applications. *Proc. ACM Program. Lang.* 2, POPL (2018), 56:1–56:29. <https://doi.org/10.1145/3158144>
- [56] Yuepeng Wang, James Dong, Rushi Shah, and Isil Dillig. 2019. Synthesizing database programs for schema refactoring. In *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. 286–300. <https://doi.org/10.1145/3314221.3314588>
- [57] Yuepeng Wang, Rushi Shah, Abby Criswell, Rong Pan, and Isil Dillig. 2020. Data Migration using Datalog Program Synthesis. *Proc. VLDB Endow.* 13, 7 (2020), 1006–1019. <https://doi.org/10.14778/3384345.3384350>
- [58] Navid Yaghmazadeh, Christian Klinger, Isil Dillig, and Swarat Chaudhuri. 2016. Synthesizing transformations on hierarchically structured data. In *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. 508–521. <https://doi.org/10.1145/2908080.2908088>
- [59] Navid Yaghmazadeh, Xinyu Wang, and Isil Dillig. 2018. Automated Migration of Hierarchical Data to Relational Tables using Programming-by-Example. *Proc. VLDB Endow.* 11, 5 (2018), 580–593. <https://doi.org/10.1145/3187009.3177735>
- [60] Yingying Zheng, Wensheng Dou, Yicheng Wang, Zheng Qin, Lei Tang, Yu Gao, Dong Wang, Wei Wang, and Jun Wei. 2022. Finding bugs in Gremlin-based graph database systems via Randomized differential testing. In *Proceedings of the SIGSOFT International Symposium on Software Testing and Analysis (ISSTA)*. 302–313. <https://doi.org/10.1145/3533767.3534409>
- [61] Qi Zhou, Joy Arulraj, Shamkant B. Navathe, William Harris, and Jinpeng Wu. 2022. SPES: A Symbolic Approach to Proving Query Equivalence Under Bag Semantics. In *International Conference on Data Engineering (ICDE)*. 2735–2748. <https://doi.org/10.1109/ICDE53745.2022.00250>
- [62] Qi Zhou, Joy Arulraj, Shamkant B. Navathe, William Harris, and Dong Xu. 2019. Automated Verification of Query Equivalence Using Satisfiability Modulo Theories. *Proc. VLDB Endow.* 12, 11 (2019), 1276–1288. <https://doi.org/10.14778/3342263.3342267>

A Semantics of Cypher Queries

Query semantics. Figure 19 presents the formal semantics of featherweight Cypher using standard functional combinators such as map, filter, foldl, head, and zip. In accordance with the standard Cypher semantics [19], we view a query Q as a function from a graph database instance to a table; hence, $\llbracket Q \rrbracket$ takes as input a property graph G and outputs a table. As shown in Figure 19, the semantics of Union, UnionAll, and OrderBy in Cypher are analogous to their SQL counterparts. The semantics of a $\text{Return}(C, \bar{E}, \bar{k})$ query is denoted $\llbracket \text{Return}(C, \bar{E}, \bar{k}) \rrbracket_G$: Given a property graph G , it first obtains a list of subgraphs by evaluating clause C and then transforms the result into a table using arguments \bar{E} and \bar{k} . Here, \bar{E} specifies a list of expressions to be evaluated against the subgraphs, each of which forms a row of the table. \bar{k} specifies the corresponding list of column names of the return value.

Clause semantics. The semantics of a Cypher clause C is denoted as $\llbracket C \rrbracket$: A clause transforms an input property graph G into a list of subgraphs, where each subgraph corresponds to the result of pattern matching against the input graph G . There are four different types of clauses in Cypher: First, the clause $\text{Match}(PP, \phi)$ returns all subgraphs of G that match the pattern PP and that satisfy predicate ϕ . Second, the clause $\text{Match}(C, PP, \phi)$ evaluates the clause C and the match clause $\text{Match}(PP, \phi)$ on the input graph G and merges their results. Third, the optional match clause OptMatch is similar to the previous match clause, except that it will use a null value for missing parts of the pattern. Finally, the fourth type of clause, namely $\text{With}(C, \bar{X}, \bar{Y})$, first applies the nested clause C to the input graph G and then renames nodes and edges as specified by the old and new names \bar{X} and \bar{Y} respectively. We further illustrate the semantics of OptMatch through an example.

Example A.1. Figure 20a shows a graph database instance, where employee A works at the CS department, but employee B has no department. Consider the following Cypher query

```
Return(OptMatch(Match([(n, EMP)],  $\top$ ), [(n, EMP), (e, WORK_AT,  $\rightarrow$ ), (m, DEPT)],  $\top$ ),
      [n.name, m.dname], [n.name, m.dname])
```

Running this query on the graph database yields the result in Figure 20b. The idea is that since there is no outgoing edge from the node representing employee B, its corresponding "department" entry in the result table is Null.

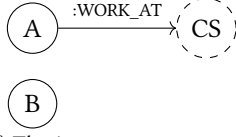
Semantics of path patterns. Given a graph database instance G and path pattern PP , $\llbracket PP \rrbracket_G$ returns a list of subgraphs of G that match the given pattern. A path pattern can have two forms. A *node pattern* NP is a pair (X, l) where l is the type of the node and X is a variable for the matched node that can be used in the rest of the query. In this case, the output of pattern matching is a list of single-node subgraphs, with each subgraph containing a node of type l . A path pattern PP can also be of the form NP, EP, PP where EP is an edge pattern (X, l, d) where l is the type of the edge, d is its direction, and X is a variable for the matched edge that can be used in the rest of the query. For example, the pattern NP_1, EP, NP_2 would match all edges that match the edge pattern EP and where additionally the source and target nodes match the node patterns NP_1 and NP_2 .

Semantics of expressions. The semantics of an expression E is denoted as $\llbracket E \rrbracket_{G,gs}$ where G is an input graph database instance and gs is a list of subgraphs (of G) denoting the result of performing a path pattern matching on G . $\llbracket E \rrbracket_{G,gs}$ evaluates the expression E on the subgraphs gs and yields a value. An expression E can be a property key k , a value v , aggregation expressions, a $\text{Cast}(\phi)$ expression, and arithmetics $E \oplus E$. For example, a property key k is evaluated by looking up the value of this key in the first subgraph of the list of matched subgraphs gs .⁶ *Aggregation*

⁶Our semantics guarantees that gs is a list of length 1 when passed to $\llbracket k \rrbracket$.

$\llbracket Q \rrbracket :: \text{Graph} \rightarrow \text{Table}$	
$\llbracket \text{OrderBy}(R, k, b) \rrbracket_G$	$= \text{foldl}(\lambda xs. \lambda_xs ++ [\text{MinTuple}(E, b, \llbracket R \rrbracket_G - xs)], [], \llbracket R \rrbracket_G) \text{ where}$ $\text{MinTuple}(E, b, xs) = \text{foldl}(\lambda x. \lambda y. \text{ite}(\text{Cmp}(E, b, x, y), x, y), \text{head}(xs), xs),$ $\text{Cmp}(E, b, x_1, x_2) = (\text{lookup}(x_1, k) < \text{lookup}(x_2, k)) = b$
$\llbracket \text{Union}(Q_1, Q_2) \rrbracket_G$	$= \text{dedup}(\llbracket \text{UnionAll}(Q_1, Q_2) \rrbracket_G)$
$\llbracket \text{UnionAll}(Q_1, Q_2) \rrbracket_G$	$= \llbracket Q_1 \rrbracket_G ++ \llbracket Q_2 \rrbracket_G$
$\llbracket R \rrbracket :: \text{Graph} \rightarrow \text{Table}$	
$\llbracket \text{Return}(C, \bar{E}, \bar{k}) \rrbracket_G$	$= \text{ite}(\neg \text{hasAgg}(\bar{E}), \text{map}(\lambda g. \text{map}(\lambda (E, k). (k, \llbracket E \rrbracket_{G, [g]}), \text{zip}(\bar{E}, \bar{k})), \llbracket C \rrbracket_G),$ $\text{map}(\lambda gs. \text{map}(\lambda (E, k), (k, \llbracket E \rrbracket_{G, gs}), \text{zip}(\bar{E}, \bar{k})), \text{Groups})) \text{ where}$ $\text{Groups} = \text{map}(\lambda V. \text{filter}(\lambda g. \llbracket \bar{A} \rrbracket_{G, [g]} = V, \llbracket C \rrbracket_G), \text{dedup}(\bar{V})),$ $\bar{A} = \text{filter}(\lambda E. \neg \text{IsAgg}(E), \bar{E}),$ $\bar{V} = \text{map}(\lambda g. \llbracket \bar{A} \rrbracket_{G, [g]}, \llbracket C \rrbracket_G)$
$\llbracket C \rrbracket :: \text{Graph} \rightarrow [\text{Graph}]$	
$\llbracket \text{Match}(PP, \phi) \rrbracket_G$	$= \text{filter}(\lambda g. \llbracket \phi \rrbracket_{G, [g]} = \top, \llbracket PP \rrbracket_G)$
$\llbracket \text{Match}(C, PP, \phi) \rrbracket_G$	$= \text{filter}(\lambda g. \llbracket \phi \rrbracket_{G, [g]} = \top, \text{map}(\lambda g_1. \text{map}(\lambda g_2. \text{merge}(g_1, g_2), \llbracket PP \rrbracket_G), \llbracket C \rrbracket_G))$
$\llbracket \text{OptMatch}(C, PP, \phi) \rrbracket_G$	$= \text{foldl}(\lambda gs. \lambda g. gs ++ \text{ite}(v_1(g) = 0, v_2(g), v_1(g)), [], \llbracket C \rrbracket_G) \text{ where}$ $v_1(g) = \text{filter}(\lambda x. \llbracket \phi \rrbracket_{G, [x]} = \top,$ $\text{map}(\lambda g'. \text{merge}(g, g'), \text{filter}(\lambda g'. g \cap g'' \neq \emptyset, \llbracket PP \rrbracket_G))),$ $v_2(g) = [\text{merge}(g, \text{Nullify}(\text{head}(\llbracket PP \rrbracket_G)))]$
$\llbracket \text{With}(C, \bar{X}, \bar{Y}) \rrbracket_G$	$= \text{map}(\lambda (N, E, P, T). (N, E, P[\bar{X} \mapsto \bar{Y}], T[\bar{X} \mapsto \bar{Y}]), \llbracket C \rrbracket_G)$
$\llbracket PP \rrbracket :: \text{Graph} \rightarrow [\text{Graph}]$	
$\llbracket NP \rrbracket_G$	$= \text{map}(\lambda n. (\{n\}, \emptyset, \{(X, k) \mapsto P(n, k) \mid k \in \text{keys}(T(n))\}, \{X \mapsto (I, \text{keys}(T(n)))\}),$ $[n \in N \mid \text{label}(T(n)) = I]) \text{ where } (N, E, P, T) = G, (X, I) = NP$
$\llbracket NP, EP, PP' \rrbracket_G$	$= \text{foldl}(\lambda gs. \lambda g'. gs ++ \text{map}(\lambda g. \text{merge}(g, g'),$ $\text{filter}(\lambda g''. g' \cap g'' \neq \emptyset, \text{Subgraphs}(G, [NP, EP, NP']))) , [], \llbracket PP' \rrbracket_G)$ $\text{where } NP' = \text{head}(PP')$
$\llbracket E \rrbracket :: \text{Graph} \rightarrow [\text{Graph}] \rightarrow \text{Value}$	
$\llbracket k \rrbracket_{G, gs}$	$= \text{lookup}(\text{head}(gs), k)$
$\llbracket v \rrbracket_{G, gs}$	$= v$
$\llbracket \text{Cast}(\phi) \rrbracket_{G, gs}$	$= \text{ite}(\llbracket \phi \rrbracket_{G, gs} = \text{Null}, \text{Null}, \text{ite}(\llbracket \phi \rrbracket_{G, gs} = \top, 1, 0))$
$\llbracket \text{Count}(E) \rrbracket_{G, gs}$	$= \text{ite}(\bigwedge_{g \in gs} \llbracket E \rrbracket_{G, [g]} = \text{Null}, \text{Null}, \text{foldl}(+, 0, \text{map}(\lambda g. \text{ite}(\llbracket E \rrbracket_{G, [g]} = \text{Null}, 0, 1), gs)))$
$\llbracket \text{Sum}(E) \rrbracket_{G, gs}$	$= \text{ite}(\bigwedge_{g \in gs} \llbracket E \rrbracket_{G, [g]} = \text{Null}, \text{Null},$ $\text{foldl}(+, 0, \text{map}(\lambda g. \text{ite}(\llbracket E \rrbracket_{G, [g]} = \text{Null}, 0, \llbracket E \rrbracket_{G, [g]}), gs)))$
$\llbracket \text{Avg}(E) \rrbracket_{G, gs}$	$= \llbracket \text{Sum}(E) \rrbracket_{G, gs} / \llbracket \text{Count}(E) \rrbracket_{G, gs}$
$\llbracket \text{Min}(E) \rrbracket_{G, gs}$	$= \text{ite}(\bigwedge_{g \in gs} \llbracket E \rrbracket_{G, [g]} = \text{Null}, \text{Null},$ $\text{foldl}(\text{min}, +\infty, \text{map}(\lambda g. \text{ite}(\llbracket E \rrbracket_{G, [g]} = \text{Null}, +\infty, \llbracket E \rrbracket_{G, [g]}), gs)))$
$\llbracket \text{Max}(E) \rrbracket_{G, gs}$	$= \text{ite}(\bigwedge_{g \in gs} \llbracket E \rrbracket_{G, [g]} = \text{Null}, \text{Null},$ $\text{foldl}(\text{max}, -\infty, \text{map}(\lambda g. \text{ite}(\llbracket E \rrbracket_{G, [g]} = \text{Null}, -\infty, \llbracket E \rrbracket_{G, [g]}), gs)))$
$\llbracket E_1 \oplus E_2 \rrbracket_{G, gs}$	$= \llbracket E_1 \rrbracket_{G, gs} \oplus \llbracket E_2 \rrbracket_{G, gs}$
$\llbracket \phi \rrbracket :: \text{Graph} \rightarrow [\text{Graph}] \rightarrow \text{Bool} \cup \{\text{Null}\}$	
$\llbracket \top \rrbracket_{G, gs}$	$= \top$
$\llbracket \perp \rrbracket_{G, gs}$	$= \perp$
$\llbracket E_1 \odot E_2 \rrbracket_{G, gs}$	$= \llbracket E_1 \rrbracket_{G, gs} \odot \llbracket E_2 \rrbracket_{G, gs}$
$\llbracket \text{IsNull}(E) \rrbracket_{G, gs}$	$= \text{ite}(\llbracket E \rrbracket_{G, gs} = \text{Null}, \top, \perp)$
$\llbracket E \in \bar{v} \rrbracket_{G, gs}$	$= \bigvee_{v \in \bar{v}} \llbracket E \rrbracket_{G, gs} = v$
$\llbracket \text{Exists}(PP) \rrbracket_{G, gs}$	$= \bigvee_{g \in \llbracket PP \rrbracket_G} \bigwedge_{K \in \bar{K}} \llbracket K \rrbracket_{G, gs} = \llbracket K \rrbracket_{G, [g]} \text{ where } \bar{K} = [\text{PK}(\text{head}(PP)), \text{PK}(\text{last}(PP))]$
$\llbracket \phi_1 \wedge \phi_2 \rrbracket_{G, gs}$	$= \llbracket \phi_1 \rrbracket_{G, gs} \wedge \llbracket \phi_2 \rrbracket_{G, gs}$
$\llbracket \phi_1 \vee \phi_2 \rrbracket_{G, gs}$	$= \llbracket \phi_1 \rrbracket_{G, gs} \vee \llbracket \phi_2 \rrbracket_{G, gs}$
$\llbracket \neg \phi \rrbracket_{G, gs}$	$= \neg \llbracket \phi \rrbracket_{G, gs}$

Fig. 19. Semantics of Cypher queries. Here, $\text{ite}(\phi, E_1, E_2)$ denotes the if-then-else expression that evaluates to E_1 when ϕ is true and E_2 otherwise. $\text{dedup}(L)$ removes all duplicated elements from list L . $\text{hasAgg}(\bar{E})$ checks if there is an aggregation expression in \bar{E} . $\text{merge}(G_1, G_2)$ merges graphs G_1 and G_2 into one graph. $\text{Nodes}(G)$ returns all nodes of graph G . $\text{Nullify}(G)$ yields a graph isomorphic to G with all values set to be Null. $\text{Subgraphs}(G, [NP_1, EP, NP_2])$ returns all subgraphs of G that match the simple pattern NP_1, EP, NP_2 .



(a) The input property graph.

n.name	m.dname
A	CS
B	Null

(b) The output of a Cypher query.

Fig. 20. An input property graph and the output of the example Cypher query.

$$\begin{array}{c}
\frac{}{\Phi_{\text{sdt}}, \Psi_R \vdash k \xrightarrow{\text{expr}} k} \text{ (E-Prop)} \quad \frac{}{\Phi_{\text{sdt}}, \Psi_R \vdash v \xrightarrow{\text{expr}} v} \text{ (E-Value)} \quad \frac{\Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi'}{\Phi_{\text{sdt}}, \Psi_R \vdash \text{Cast}(\phi) \xrightarrow{\text{expr}} \text{Cast}(\phi')} \text{ (E-Pred)} \\
\frac{\Phi_{\text{sdt}}, \Psi_R \vdash E \xrightarrow{\text{expr}} E'}{\Phi_{\text{sdt}}, \Psi_R \vdash \text{Agg}(E) \xrightarrow{\text{expr}} \text{Agg}(E')} \text{ (E-Agg)} \quad \frac{\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \xrightarrow{\text{expr}} E'_1 \quad \Phi_{\text{sdt}}, \Psi_R \vdash E_2 \xrightarrow{\text{expr}} E'_2}{\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \oplus E_2 \xrightarrow{\text{expr}} E'_1 \oplus E'_2} \text{ (E-Arith)}
\end{array}$$

Fig. 21. Translation rules for expressions.

$$\begin{array}{c}
\frac{}{\Phi_{\text{sdt}}, \Psi_R \vdash \top \xrightarrow{\text{pred}} \top} \text{ (P-True)} \quad \frac{}{\Phi_{\text{sdt}}, \Psi_R \vdash \perp \xrightarrow{\text{pred}} \perp} \text{ (P-False)} \quad \frac{\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \xrightarrow{\text{expr}} E'_1 \quad \Phi_{\text{sdt}}, \Psi_R \vdash E_2 \xrightarrow{\text{expr}} E'_2}{\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \odot E_2 \xrightarrow{\text{pred}} E'_1 \odot E'_2} \text{ (P-Logic)} \\
\frac{\Phi_{\text{sdt}}, \Psi_R \vdash E \xrightarrow{\text{expr}} E'}{\Phi_{\text{sdt}}, \Psi_R \vdash \text{IsNull}(E) \xrightarrow{\text{pred}} \text{IsNull}(E')} \text{ (P-IsNull)} \quad \frac{\Phi_{\text{sdt}}, \Psi_R \vdash E \xrightarrow{\text{expr}} E'}{\Phi_{\text{sdt}}, \Psi_R \vdash E \in \bar{v} \xrightarrow{\text{pred}} E' \in \bar{v}} \text{ (P-In)} \quad \frac{\Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi'}{\Phi_{\text{sdt}}, \Psi_R \vdash \neg \phi \xrightarrow{\text{pred}} \neg \phi'} \text{ (P-Not)} \\
\frac{\begin{array}{c} \circ \in \{\wedge, \vee\} \\ \Phi_{\text{sdt}}, \Psi_R \vdash \phi_1 \xrightarrow{\text{pred}} \phi'_1 \\ \Phi_{\text{sdt}}, \Psi_R \vdash \phi_2 \xrightarrow{\text{pred}} \phi'_2 \end{array}}{\Phi_{\text{sdt}}, \Psi_R \vdash \phi_1 \odot \phi_2 \xrightarrow{\text{pred}} \phi'_1 \odot \phi'_2} \text{ (P-AndOr)} \quad \frac{\begin{array}{c} (X_1, l_1) = \text{head}(PP) \quad (X_2, l_2) = \text{last}(PP) \\ l_1(\dots) \rightarrow R_{l_1}(\dots) \in \Phi_{\text{sdt}} \quad l_2(\dots) \rightarrow R_{l_2}(\dots) \in \Phi_{\text{sdt}} \\ \Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} \mathcal{X}, Q \quad \bar{a} = [\text{PK}(R_{l_1}), \text{PK}(R_{l_2})] \end{array}}{\Phi_{\text{sdt}}, \Psi_R \vdash \text{Exists}(PP) \xrightarrow{\text{pred}} \bar{a} \in \Pi_{\bar{a}}(Q)} \text{ (P-Exists)}
\end{array}$$

Fig. 22. Translation rules for predicates.

expressions involve the Count, Sum, Avg, Min, Max operators and their semantics are similar to SQL, where an aggregation operator is used to calculate an output value over the return values of the sub-expression. Finally, the $\text{Cast}(\phi)$ expression simply casts a predicate ϕ to 0, 1, or Null.

Semantics of predicates. Given a predicate ϕ , property graph G , and list of subgraphs gs , $\llbracket \phi \rrbracket_{G,gs}$ defines how to evaluate ϕ on G, gs . Similar to SQL, we take the three-valued logic interpretation of a boolean predicate ϕ , meaning that the evaluation result can contain Null, and we can perform boolean arithmetics involving \top , \perp , Null as terms. In particular, $\perp \wedge \text{Null} = \perp$ and $\top \vee \text{Null} = \top$; otherwise, the evaluation result is Null if any of its arguments are Null.

B Transpilation of Cypher Predicates and Expressions

As shown in Figure 21 and Figure 22, most of the rules for translating expressions and predicates are straightforward. The only rule that is tricky is P-Exists for translating Exists predicates. Specifically, given a predicate $\text{Exists}(PP)$ in Cypher, the P-Exists rule first translates the pattern PP to a SQL query Q and then translates $\text{Exists}(PP)$ to an **IN** predicate in SQL, namely $\bar{a} \in \Pi_{\bar{a}}(Q)$. Recall from the Cypher semantics that $\text{Exists}(PP)$ checks if there exists a match of pattern PP in the graph, so the translated SQL query uses an **IN** predicate to perform the existence checking. Here, the attribute list \bar{a} is obtained from the primary keys of relations corresponding to the first and last nodes in the path pattern. The predicate $\bar{a} \in \Pi_{\bar{a}}(Q)$ checks that \bar{a} used in the parent query is in the set of results returned by the subquery Q .

Example B.1. Given the standard database transformer Φ_{sdt} in Example 5.2, let us consider the Cypher predicate $\text{Exists}([(n, \text{EMP}), (e, \text{WORK_AT}, \rightarrow), (m, \text{DEPT})])$. Based on the P-Exists rule, it is translated to a SQL predicate

$$[n.\text{id}, m.\text{dnum}] \in \Pi_{[n.\text{id}, m.\text{dnum}]} (\rho_n(\text{emp}) \bowtie_{n.\text{id}=e.\text{SRC}} \rho_e(\text{work_at}) \bowtie_{e.\text{TGT}=m.\text{dnum}} \rho_m(\text{dept}))$$

C An Equivalent Cypher Query of Motivating Example

As discussed in Section 2, the Cypher query in Figure 4c is proven non-equivalent to the SQL query in Figure 4a. The issue arises because the original query filters $s: \text{SENTENCE}$ using the **WITH** clause, but this does not enforce filtering at the level of individual matches in the second **MATCH** clause. As a result, the second **MATCH** explores more paths than intended, leading to double counting. The correct Cypher query shown below moves the filtering into an **EXISTS** condition, ensuring each **SENTENCE** node is considered only if it satisfies the first pattern, which prevents the discrepancy.

```
MATCH (s:SENTENCE)<-[r3:SP]-(p2:PA)<-[r4:CS]-[c2:CONCEPT]
WHERE EXISTS { MATCH (c1:CONCEPT {CID: 1})-[r1:CS]->(p1:PA)-[r2:SP]->(s:SENTENCE) }
RETURN c2.CID, Count (*)
```

D Qualitative Analysis of Manually-Written Buggy Queries

We manually inspected the buggy queries in Section 6.1 to gain intuition about common problems. At a high level, there are several root causes of bugs we identified in the manually-written benchmarks. These include:

- (1) Using nested match instead of an existential pattern, as shown in Section 2.
- (2) Misusing path patterns for **OPTIONAL MATCH**. For example, let us look at the following benchmark adapted from the Neo4j tutorial⁷. We removed the **OrderBy** clauses and renamed node and edge labels to avoid confusion. Given the SQL query

```
SELECT P.ProductName, Sum (OD.UnitPrice * OD.Quantity) AS Volume FROM Customers AS C
LEFT JOIN Orders AS O ON C.CustomerID = O.CustomerID
LEFT JOIN OrderDetails AS OD ON O.OrderID = OD.OrderID
LEFT JOIN Products AS P ON OD.ProductID = P.ProductID
WHERE C.CompanyName = 'Drachenblut Delikatessen' GROUP BY P.ProductName
```

The tutorial provides a corresponding Cypher query

```
MATCH (C:Customer {CompanyName:'Drachenblut Delikatessen'})
OPTIONAL MATCH (P:Product)<-[OD:OrderDetails]-(O:Order)<-[:Purchased]-(C)
RETURN P.ProductName, Sum (OD.UnitPrice * OD.Quantity) AS Volume
```

However, this Cypher query is not equivalent to the SQL query. The underlined optional match clause uses a path pattern containing three nodes and two edges, which does not consider the missing **Product** as queried by a left outer join in SQL. The difference is more straightforward to understand if we inspect the transpiled SQL query

```
WITH T0 AS (SELECT C.CustomerID AS C_CustomerID
FROM Customer AS C WHERE C.CompanyName = 'Drachenblut Delikatessen'),
T1 AS (SELECT P.ProductName AS P_ProductName, OD.Quantity AS OD_Quantity,
OD.UnitPrice AS OD_UnitPrice, C.CustomerID AS C_CustomerID,
FROM Product AS P JOIN OrderDetails AS OD JOIN Purchased AS PD JOIN Customer
AS C ON P.ProductID = OD.TGT AND OD.SRC = PD.TGT AND PD.SRC = C.CustomerID),
T2 AS (SELECT T1.P_ProductName AS ProductName, T1.O_Quantity AS Quantity, T1.O_UnitPrice AS UnitPrice
FROM T0 LEFT JOIN T1 ON T0.C_CustomerID = T1.C_CustomerID)
SELECT ProductName, Sum (UnitPrice * Quantity) AS Volume FROM T2 GROUP BY ProductName;
```

Here, the underlined clause uses inner joins instead of left outer joins in the original SQL query.

⁷<https://neo4j.com/docs/getting-started/cypher-intro/cypher-sql>

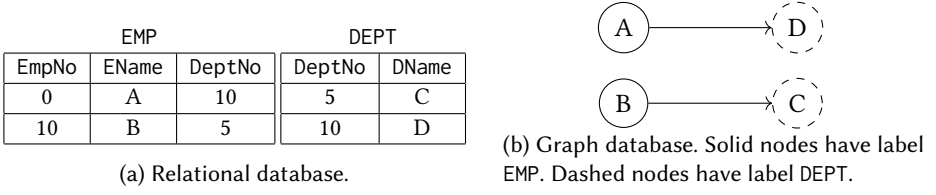


Fig. 23. A counterexample that demonstrates non-equivalence.

Table 5. Transpilation Results of OPENCYPHERTRANSPILER.

Dataset	#	# Unsupported	# SynErr	# Incorrect	# Correct
StackOverflow	12	8	0	0	4
Tutorial	26	14	0	0	12
Academic	7	6	0	0	1
VeriEQL	60	44	1	0	15
Mediator	100	100	0	0	0
GPT-Translate	205	112	1	2	90
Total	410	284	2	2	122

- (3) Misusing different nodes or edges with the same label. For instance, consider the following SQL query collected from the VeriEQL dataset

```
SELECT t0.EmpNo, t0.DeptNo, t1.DeptNo AS DeptNo0, FROM (
  SELECT EmpNo, EName, DeptNo, DeptNo + EmpNo AS f9 FROM EMP WHERE EmpNo = 10
) AS t0 JOIN (SELECT DeptNo, Name, DeptNo + 5 AS f2 FROM DEPT) AS t1
ON t0.EmpNo = t1.DeptNo AND t0.f9 = t1.f2
```

A graduate student wrote a corresponding Cypher query as follows:

```
MATCH (t0:EMP {EmpNo: 10})-[:WORK_AT]->(t1:DEPT)
WHERE t1.DeptNo + t0.EmpNo = t1.DeptNo + 5
RETURN t0.EmpNo, t1.DeptNo, t1.DeptNo AS DeptNo0
```

However, this Cypher query is not equivalent to the SQL query because it fails to introduce a DEPT node different from t1. Consequently, the query matches an incorrect t1 and uses an incorrect filtering predicate. GRAPHITI disproves the equivalence between the SQL and Cypher queries and provides a counterexample as shown in Figure 23. The counterexample contains a relational and a graph database holding the same data in different formats. Running the SQL query on the relational database in Figure 23a returns a table with one row (10, 5, 10), whereas running the Cypher query on the graph database in Figure 23b returns empty result.

E Comparing GRAPHITI's Transpiler with OPENCYPHERTRANSPILER

In this section, we compare the transpiler of GRAPHITI with OPENCYPHERTRANSPILER [29] on all 410 benchmarks. The results of running OPENCYPHERTRANSPILER on these benchmarks are presented in Table 5. As shown in the table, 284 out of 410 (69%) queries fall outside the Cypher fragment supported by OPENCYPHERTRANSPILER. Among the remaining 126 benchmarks, OPENCYPHERTRANSPILER ill-translate 2 Cypher queries into syntactically invalid SQL queries. Furthermore, a manual inspection reveals that 2 other Cypher queries are transpiled into semantically incorrect SQL queries. These findings demonstrate that OPENCYPHERTRANSPILER serves as a “best-effort” tool without any soundness guarantees about the result of transpilation and cannot preserve semantic equivalence during transpilation. In contrast, GRAPHITI can correctly transpile all 410 Cypher queries used in our evaluation.

Here are three sample Cypher queries that OPENCYPHERTRANSPILER cannot transpile correctly.

- (1) A Cypher query that cannot be translated by OPENCYPHERTRANSPILER.

```

MATCH (N:PERSON)
WITH N.ZIPCODE AS ZIP, Count(*) AS POPULATION
WHERE POPULATION > 3
RETURN ZIP, POPULATION

```

OPENCYPHERTRANSPILER does not support expressions such as **Count**(*) or **Avg**(*).

- (2) A Cypher query where OPENCYPHERTRANSPILER generates ill-formed SQL query with undefined variables.

```

MATCH (X:USR), (U:PIC), (V:PIC, (W:PIC)
WHERE X.USRUID = U.PICUID AND X.USRUID = V.PICUID AND X.USRUID = W.PICUID AND
    W.PRCSIZE = V.PRCSIZE AND U IS NOT NULL AND V IS NULL
RETURN DISTINCT X.USRUID AS XID, X.USRNAME AS XNAME

```

OPENCYPHERTRANSPILER translates the above Cypher query to the following:

```

SELECT DISTINCT __X_USRUID AS XID, ... FROM (
    ... WHERE ... AND ((U) IS NOT NULL) AND ...
) AS _proj

```

However, this SQL query uses an undefined table name “U”.

- (3) A Cypher query where OPENCYPHERTRANSPILER produces a semantically incorrect transpilation result.

```

MATCH (e1:EMPLOYEES)
OPTIONAL MATCH (e2:EMPLOYEEUNI)-[:IS]->(e1)
RETURN e2.UNIQUE_ID, e1.NAME

```

OPENCYPHERTRANSPILER translates the above Cypher query to the following:

```

SELECT UNIQUE_ID, NAME FROM ( SELECT * FROM EMPLOYEEUNI JOIN IS ON ...) AS T1
LEFT JOIN EMPLOYEES AS T2 ON ...

```

The correct SQL query should be as follows:

```

SELECT UNIQUE_ID, NAME FROM ( SELECT * FROM EMPLOYEEUNI JOIN IS ON ...) AS T1
RIGHT JOIN EMPLOYEES AS T2 ON ...

```

F Proofs

THEOREM F.1 (SOUNDNESS OF TRANSLATION (5.7)). *Let Ψ_G be a graph schema and Q be a Cypher query over Ψ_G . Let $\Psi_R = (S, \xi)$ be the induced relational schema of Ψ_G , and Φ_{sdt} be the standard database transformer from Ψ_G to Ψ_R . If $\Phi_{\text{sdt}}, \Psi_R \vdash Q \xrightarrow{\text{query}} Q'$, then Q' is equivalent to Q modulo Φ_{sdt} , i.e., $Q \simeq_{\Phi_{\text{sdt}}} Q'$.*

PROOF. Prove by structural induction on Q .

- (1) Base case: $Q = \text{Return}(C, \bar{E}, \bar{k})$.

By the inductive hypothesis, we have $\Phi_{\text{sdt}}(\llbracket C \rrbracket_G) = \llbracket Q' \rrbracket_R$ where $\Phi_{\text{sdt}}, \Psi_R \vdash C \xrightarrow{\text{clause}} X, Q'$. Suppose that $|\bar{E}| = m$ and $|\llbracket Q \rrbracket_R| = |\llbracket C \rrbracket_G| = n$. Let us discuss the predicate $\neg\text{hasAgg}(\bar{E})$ in two cases.

- (a) If $\neg\text{hasAgg}(\bar{E}) = \top$, then, for any expression $e \in \bar{E}$, $\text{IsAgg}(e) = \perp$. By Figure 19, we know

$$\begin{aligned}
 \llbracket \text{Return}(C, \bar{E}, \bar{k}) \rrbracket_G &= \text{map}(\lambda g. \text{map}(\lambda (E, k). (k, \llbracket E \rrbracket_{G, [g]}), \text{zip}(\bar{E}, \bar{k})), \llbracket C \rrbracket_G) \\
 &= [\\
 &\quad [(k_1, \llbracket E_1 \rrbracket_{G, [g_1]}), \dots, (k_m, \llbracket E_m \rrbracket_{G, [g_1]})], \\
 &\quad \dots, \\
 &\quad [(k_1, \llbracket E_1 \rrbracket_{G, [g_n]}), \dots, (k_m, \llbracket E_m \rrbracket_{G, [g_n]})] \\
 &]
 \end{aligned}$$

Also, by the SQL semantics [24], we know

$$\begin{aligned} \llbracket \Pi_{\rho_{\bar{k}}(\bar{E})}(Q') \rrbracket_R &= \text{map}(\lambda t. \text{map}(\lambda(E', k). (k, \llbracket E' \rrbracket_{R, [t]}), \text{zip}(\bar{E}', \bar{k})), \llbracket Q' \rrbracket_R) \\ &= [\\ &\quad (k_1 : \llbracket E'_1 \rrbracket_{R, [t_1]}, \dots, k_m : \llbracket E'_m \rrbracket_{R, [t_1]}), \\ &\quad \dots, \\ &\quad (k_1 : \llbracket E'_1 \rrbracket_{R, [t_n]}, \dots, k_m : \llbracket E'_m \rrbracket_{R, [t_n]}), \\ &\quad] \end{aligned}$$

There exists a bijective mapping $\pi = \{E_j \mapsto E'_j \mid 1 \leq j \leq m\}$ and $\pi^{-1} = \{E'_j \mapsto E_j \mid 1 \leq j \leq m\}$. Let $g_i \in \llbracket C \rrbracket_G$ be a graph derived from G and $t_i \in \llbracket Q' \rrbracket_R$ be a tuple derived from R s.t. $\Phi_{\text{sdt}}(g_i) = t_i$. Then, we know $\llbracket E_j \rrbracket_{G, [g_i]} = \llbracket E'_j \rrbracket_{R, [t_i]}$ by Lemma F.4. This case is proved because the following formula holds

$$\forall G, R. G \triangleright \Psi_G \wedge R \triangleright \Psi_R \wedge G \sim_{\Phi_{\text{sdt}}} R \Rightarrow \Phi_{\text{sdt}}(\llbracket C \rrbracket_G) = \llbracket Q \rrbracket_R$$

- (b) If $\neg \text{hasAgg}(\bar{E}) = \perp$, then, for some expression $E \in \bar{E}$, $\text{IsAgg}(E) = \top$.

Suppose that $n_v = |\text{dedup}(\bar{V})| \leq n = |\bar{V}|$. Then, by the definition of Groups in Figure 19, we know

$$\llbracket \text{Return}(C, \bar{E}, \bar{k}) \rrbracket_G = \text{map}(\lambda gs. \text{map}(\lambda(E, k). (k, \llbracket E \rrbracket_{G, gs}), \text{zip}(\bar{E}, \bar{k})), \text{Groups})$$

where Groups = $[gs_1, \dots, gs_{n_v}]$ and gs_i is a list of graphs consisting of at least one graph derived from G .

Also, by the SQL semantics [24], we know

$$\begin{aligned} \llbracket \text{GroupBy}(Q', \bar{A}', \rho_{\bar{k}}(\bar{E}'), \top) \rrbracket &= \text{map}(\lambda xs. \text{map}(\lambda(E, k). (k, \llbracket E' \rrbracket_{R, xs}), \text{zip}(\bar{E}', \bar{k})), \\ &\quad \text{filter}(\lambda xs. \llbracket \top \rrbracket_{R, xs} = \top, \text{Ts})) \\ &= \text{map}(\lambda xs. \text{map}(\lambda(E, k). (k, \llbracket E' \rrbracket_{R, xs}), \text{zip}(\bar{E}', \bar{k})), \text{Ts}) \end{aligned}$$

where Ts denote the grouping list derived by Q' .

There exists a bijective mapping $\pi = \{E_j \mapsto E'_j \mid 1 \leq j \leq m\}$ and $\pi^{-1} = \{E'_j \mapsto E_j \mid 1 \leq j \leq m\}$. Then, $\pi(\bar{E}) = \pi([E_1, \dots, E_m]) = [E'_1, \dots, E'_m] = \bar{E}'$ and $\pi(\bar{A}) = \text{filter}(\lambda E. \neg \text{hasAgg}(E), \pi(\bar{E})) = \text{filter}(\lambda E. \neg \text{hasAgg}(E), \bar{E}) = \bar{A}'$. Since $\Phi_{\text{sdt}}(\llbracket C \rrbracket_G) = \llbracket Q \rrbracket_R$ and $\pi(\bar{A}) = \bar{A}'$, we always find n_v unique results while evaluating \bar{A} on $\llbracket C \rrbracket_G$ and \bar{A}' on $\llbracket Q \rrbracket_R$ (namely, $|\text{Groups}| = |\text{Ts}| = n_v$). Let $\text{Ts} = [ts_1, \dots, ts_m]$ and ts_i corresponds to $gs_i \in \text{Groups}$. Then $\Phi_{\text{sdt}}(\text{Groups}) = \text{Ts}$ and, therefore, $\Phi_{\text{sdt}}(\llbracket \text{Return}(C, \bar{E}, \bar{k}) \rrbracket_G) = \llbracket \text{GroupBy}(Q', \bar{A}', \rho_{\bar{k}}(\bar{E}'), \top) \rrbracket$.

Thus, Theorem F.1 in this case is proved.

- (2) Inductive case: $Q = \text{OrderBy}(R_1, k, b)$.

Specifically, R_1 denotes the Return statement of Cypher. By the inductive hypothesis, we have

$\llbracket R_1 \rrbracket_G = \llbracket Q' \rrbracket_R$ where $\Phi_{\text{sdt}}, \Psi_R \vdash R_1 \xrightarrow{\text{query}} Q'$. Also, by Lemma F.4, we have $\llbracket k \rrbracket_{G, gs} = \llbracket k \rrbracket_{R, \tau}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash k \xrightarrow{\text{expr}} k$ and k is property key. By Figure 19, we know

$$\llbracket \text{OrderBy}(R_1, k, b) \rrbracket_G = \text{foldl}(\lambda xs. \lambda_xs ++ [\text{MinTuple}(k, b, \llbracket R_1 \rrbracket_G - xs)], [], \llbracket R_1 \rrbracket_G)$$

Also, by the SQL semantics [24], we know

$$\llbracket \text{OrderBy}(Q', k, b) \rrbracket_G = \text{foldl}(\lambda xs. \lambda_xs ++ [\text{MinTuple}(k, b, \llbracket Q' \rrbracket_{\tau} - xs)], [], \llbracket Q' \rrbracket_{\tau})$$

Since OrderBy in Cypher and OrderBy in SQL will not change the order of columns, we can reuse the bijective mapping π from the inductive hypothesis for this case. Further, as

$\Phi_{\text{sdt}}(\llbracket R_1 \rrbracket_G) = \llbracket Q' \rrbracket_{\mathcal{T}}$ and the functions MinTuple and MinTuple are equivalent by their definitions, $\Phi_{\text{sdt}}(\llbracket \text{OrderBy}(R_1, k, b) \rrbracket_G) = \llbracket \text{OrderBy}(Q', k, b) \rrbracket_G$

(3) Inductive case: $Q = \text{Union}(Q_1, Q_2)$.

By the inductive hypothesis, we have $\llbracket Q_1 \rrbracket_{G,gs} \simeq_{\Phi_{\text{sdt}}} \llbracket Q'_1 \rrbracket_R$ and $\llbracket Q_2 \rrbracket_{G,gs} \simeq_{\Phi_{\text{sdt}}} \llbracket Q'_2 \rrbracket_R$ where $\Phi_{\text{sdt}}, \Psi_R \vdash Q_1 \xrightarrow{\text{query}} Q'_1$ and $\Phi_{\text{sdt}}, \Psi_R \vdash Q_2 \xrightarrow{\text{query}} Q'_2$. Similarly, since the order of columns will not be changed in this case, we can inherit the bijective mapping π from the inductive hypothesis for this case. By Figure 19 and the definitions of dedup and Distinct , $\llbracket \text{Union}(Q_1, Q_2) \rrbracket_G = \text{dedup}(\llbracket \text{UnionAll}(Q_1, Q_2) \rrbracket_G) \simeq_{\Phi_{\text{sdt}}} \text{Distinct}(\llbracket Q'_1 \uplus Q'_2 \rrbracket_R) = \llbracket Q'_1 \cup Q'_2 \rrbracket_R$.

(4) Inductive case: $Q = \text{UnionAll}(Q_1, Q_2)$.

By the inductive hypothesis, we have $\llbracket Q_1 \rrbracket_{G,gs} \simeq_{\Phi_{\text{sdt}}} \llbracket Q'_1 \rrbracket_R$ and $\llbracket Q_2 \rrbracket_{G,gs} \simeq_{\Phi_{\text{sdt}}} \llbracket Q'_2 \rrbracket_R$ where $\Phi_{\text{sdt}}, \Psi_R \vdash Q_1 \xrightarrow{\text{query}} Q'_1$ and $\Phi_{\text{sdt}}, \Psi_R \vdash Q_2 \xrightarrow{\text{query}} Q'_2$. Also, since the order of columns will not be changed in this case, we can inherit the bijective mapping π from the inductive hypothesis for this case. By Figure 19, $\llbracket \text{UnionAll}(Q_1, Q_2) \rrbracket_G = \llbracket Q_1 \rrbracket_G ++ \llbracket Q_2 \rrbracket_G \simeq_{\Phi_{\text{sdt}}} \llbracket Q'_1 \rrbracket_R ++ \llbracket Q'_2 \rrbracket_R = \llbracket Q'_1 \uplus Q'_2 \rrbracket_R$.

□

LEMMA F.2. Let Ψ_G be a graph schema and C be a Cypher clause. Let Ψ_R be a relational schema and Φ_{sdt} be the standard database transformer from Ψ_G to Ψ_R . If $\Phi_{\text{sdt}}, \Psi_R \vdash C \xrightarrow{\text{clause}} \mathcal{X}, Q$, it holds that⁸

$$\forall G, R. G \triangleright \Psi_G \wedge R \triangleright \Psi_R \wedge G \sim_{\Phi_{\text{sdt}}} R \Rightarrow \Phi_{\text{sdt}}(\llbracket C \rrbracket_G) = \llbracket Q \rrbracket_R$$

PROOF. Prove by structural induction on C .

(1) Base case: $C = \text{Match}(PP, \phi)$.

By Lemma F.3, we have $\Phi_{\text{sdt}}(\llbracket PP \rrbracket_G) = \llbracket Q \rrbracket_R$ where $\Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} \mathcal{X}, Q$. Also, by Lemma F.5, we have $\llbracket \phi \rrbracket_{G,gs} = \llbracket \phi' \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi'$, gs is a list of graph derived from G and \mathcal{T} is a list of tuples derived from R . By Figure 19,

$$\begin{aligned} \Phi_{\text{sdt}}(\llbracket \text{Match}(PP, \phi) \rrbracket_G) &= \Phi_{\text{sdt}}(\text{filter}(\lambda g. \llbracket \phi \rrbracket_{G,[g]} = \top, \llbracket PP \rrbracket_G)) \\ &= \text{filter}(\lambda t. \llbracket \phi' \rrbracket_{\Phi_{\text{sdt}}(G),[t]} = \top, \Phi_{\text{sdt}}(\llbracket PP \rrbracket_G)) \\ &= \text{filter}(\lambda t. \llbracket \phi' \rrbracket_{R,[t]} = \top, \llbracket Q \rrbracket_R) \\ &= \llbracket \sigma_{\phi'}(Q) \rrbracket_R \end{aligned}$$

(2) Inductive case: $C = \text{MATCH}(C_1, PP, \phi)$.

By the inductive hypothesis, we have $\Phi_{\text{sdt}}(\llbracket C_1 \rrbracket_G) = \llbracket Q_1 \rrbracket_R$ where $\Phi_{\text{sdt}}, \Psi_R \vdash C_1 \xrightarrow{\text{clause}} \mathcal{X}_1, Q_1$.

By Lemma F.3, we have $\Phi_{\text{sdt}}(\llbracket PP \rrbracket_G) = \llbracket Q_2 \rrbracket_R$ where $\Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} \mathcal{X}_2, Q_2$. Also, by Lemma F.5, we have $\llbracket \phi \rrbracket_{G,gs} = \llbracket \phi' \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi'$, gs is a list of graph derived from G and \mathcal{T} is a list of tuples derived from R .

By the definition of merge for graph, we know, for any two graphs g_1 and g_2 , there is no duplicate node, edge, property key in $\text{merge}(g_1, g_2)$. Similarly, by the definition of merge for tuple, we know the function merge has the same functionality as merge . Therefore, $\Phi_{\text{sdt}}(\text{merge}(g_1, g_2)) = \text{merge}(t_1, t_2)$ iff $\Phi_{\text{sdt}}(g_1) = t_1 \wedge \Phi_{\text{sdt}}(g_2) = t_2$.

For any Cypher clause C_1 and pattern PP , let us discuss $\mathcal{X}_1 \cap \mathcal{X}_2$ in two cases.

⁸We slightly abuse Φ_{sdt} over lists to apply Φ_{sdt} to each element in list $\llbracket C \rrbracket_G$.

(a) If $\mathcal{X}_1 \cap \mathcal{X}_2 = \emptyset$, then we know $\phi'' = \phi'$ by Figure 17 and

$$\begin{aligned}
 \Phi_{\text{sdt}}(\llbracket \text{MATCH}(C_1, PP, \phi) \rrbracket_G) &= \Phi_{\text{sdt}}(\text{filter}(\lambda g. \llbracket \phi \rrbracket_{G, [g]} = \top, \text{map}(\lambda g_1. \text{map}(\lambda g_2. \text{merge}(g_1, g_2), \\
 &\quad \llbracket PP \rrbracket_G), \llbracket C_1 \rrbracket_G))) \\
 &= \text{filter}(\lambda t. \llbracket \phi' \rrbracket_{\Phi_{\text{sdt}}(G), [t]} = \top, \text{map}(\lambda t_1. \text{map}(\lambda t_2. \text{merge}(t_1, t_2), \\
 &\quad \Phi_{\text{sdt}}(\llbracket PP \rrbracket_G), \Phi_{\text{sdt}}(\llbracket C_1 \rrbracket_G))) \\
 &= \text{filter}(\lambda t. \llbracket \phi' \rrbracket_{R, [t]} = \top, \text{map}(\lambda t_1. \text{map}(\lambda t_2. \text{merge}(t_1, t_2), \\
 &\quad \llbracket Q_1 \rrbracket_R), \llbracket Q_2 \rrbracket_R)) \\
 &= \text{filter}(\lambda t. \llbracket \phi' \rrbracket_{R, [t]} = \top, \llbracket Q_1 \times Q_2 \rrbracket_R) \\
 &= \text{filter}(\lambda t. \llbracket \phi' \rrbracket_{R, [t]} = \top, \llbracket \rho_{T_1}(Q_1) \times \rho_{T_2}(Q_2) \rrbracket_R) \\
 &= \llbracket \sigma_{\phi'}(\rho_{T_1}(Q_1) \times \rho_{T_2}(Q_2)) \rrbracket_R \\
 &= \llbracket \rho_{T_1}(Q_1) \bowtie_{\phi'} \rho_{T_2}(Q_2) \rrbracket_R \\
 &= \llbracket \rho_{T_1}(Q_1) \bowtie_{\phi''} \rho_{T_2}(Q_2) \rrbracket_R
 \end{aligned}$$

by Figure 19 where T_1 and T_2 are fresh names.

(b) If $\mathcal{X}_1 \cap \mathcal{X}_2 \neq \emptyset$, then we know $\phi'' = \phi' \wedge \wedge_{(X:I) \in \mathcal{X}_1 \cap \mathcal{X}_2} T_1.\xi_{pk}(\Lambda(I)) = T_2.\xi_{pk}(\Lambda(I))$ for merging two overlapping graphs T_1 and T_2 by Figure 17 and

$$\begin{aligned}
 \Phi_{\text{sdt}}(\llbracket \text{MATCH}(C_1, PP, \phi) \rrbracket_G) &= \Phi_{\text{sdt}}(\text{filter}(\lambda g. \llbracket \phi \rrbracket_{G, [g]} = \top, \text{map}(\lambda g_1. \text{map}(\lambda g_2. \text{merge}(g_1, g_2), \\
 &\quad \llbracket PP \rrbracket_G), \llbracket C_1 \rrbracket_G))) \\
 &= \text{filter}(\lambda t. \llbracket \phi' \rrbracket_{\Phi_{\text{sdt}}(G), [t]} = \top, \text{map}(\lambda t_1. \text{map}(\lambda t_2. \text{merge}(t_1, t_2), \\
 &\quad \Phi_{\text{sdt}}(\llbracket PP \rrbracket_G), \Phi_{\text{sdt}}(\llbracket C_1 \rrbracket_G))) \\
 &= \text{filter}(\lambda t. \llbracket \phi' \rrbracket_{R, [t]} = \top, \text{map}(\lambda t_1. \text{map}(\lambda t_2. \text{merge}(t_1, t_2), \\
 &\quad \llbracket Q_1 \rrbracket_R), \llbracket Q_2 \rrbracket_R)) \\
 &= \text{filter}(\lambda t. \llbracket \phi' \rrbracket_{R, [t]} = \top, \\
 &\quad \llbracket \rho_{T_1}(Q_1) \bowtie_{\wedge_{(X:I) \in \mathcal{X}_1 \cap \mathcal{X}_2} T_1.\xi_{pk}(\Lambda(I)) = T_2.\xi_{pk}(\Lambda(I))} \rho_{T_2}(Q_2) \rrbracket_R) \\
 &= \llbracket \sigma_{\phi'}(\rho_{T_1}(Q_1) \bowtie_{\wedge_{(X:I) \in \mathcal{X}_1 \cap \mathcal{X}_2} T_1.\xi_{pk}(\Lambda(I)) = T_2.\xi_{pk}(\Lambda(I))} \rho_{T_2}(Q_2)) \rrbracket_R \\
 &= \llbracket \rho_{T_1}(Q_1) \bowtie_{\phi' \wedge \wedge_{(X:I) \in \mathcal{X}_1 \cap \mathcal{X}_2} T_1.\xi_{pk}(\Lambda(I)) = T_2.\xi_{pk}(\Lambda(I))} \rho_{T_2}(Q_2) \rrbracket_R \\
 &= \llbracket \rho_{T_1}(Q_1) \bowtie_{\phi''} \rho_{T_2}(Q_2) \rrbracket_R
 \end{aligned}$$

by Figure 19 where T_1 and T_2 are fresh names.

Thus, Lemma F.2 in this case is proved.

(3) Inductive case: $C = \text{OptMatch}(C_1, PP, \phi)$.

By the inductive hypothesis, we have $\Phi_{\text{sdt}}(\llbracket C_1 \rrbracket_G) = \llbracket Q_1 \rrbracket_R$ where $\Phi_{\text{sdt}}, \Psi_R \vdash C_1 \xrightarrow{\text{clause}} \mathcal{X}_1, Q_1$.

By Lemma F.3, we have $\Phi_{\text{sdt}}(\llbracket PP \rrbracket_G) = \llbracket Q_2 \rrbracket_R$ where $\Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} \mathcal{X}_2, Q_2$. Also, by

Lemma F.5, we have $\llbracket \phi \rrbracket_{G, gs} = \llbracket \phi' \rrbracket_{R, \mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi'$, gs is a list of graph derived from G and \mathcal{T} is a list of tuples derived from R . Further, by the definition of Nullify, we know $\text{Nullify}(\text{head}(\llbracket PP \rrbracket_G))$ return a graph which is isomorphic to the first graph of $\llbracket PP \rrbracket_G$ but with Null values for all its property keys. Therefore, $\Phi_{\text{sdt}}(\text{Nullify}(\text{head}(\llbracket PP \rrbracket_G)))$ is equal to a unique tuple T_{Null} used in the semantics of SQL [24].

For any graph $g_i \in \llbracket C \rrbracket_G$, let us discuss the predicate $|v_1(g_i)| = 0$ in two cases.

(a) If the predicate $|v_1(g_i)| = 0$ holds, then, for any graph $g'_j \in \llbracket PP \rrbracket_G$, $\llbracket \phi \rrbracket_{G, [\text{merge}(g_i, g'_j)]} = \perp$ does not hold. By Figure 19,

$$\begin{aligned}
 \Phi_{\text{sdt}}(\text{ite}(|v_1(g_i)| = 0, v_2(g_i), v_1(g_i))) &= \Phi_{\text{sdt}}(\llbracket v_2(g_i) \rrbracket) \\
 &= \llbracket \text{merge}(\Phi_{\text{sdt}}(g_i), \Phi_{\text{sdt}}(\text{Nullify}(\text{head}(\llbracket PP \rrbracket_G)))) \rrbracket \\
 &= \llbracket \text{merge}(t_i, T_{\text{Null}}) \rrbracket \\
 &= [t_i] \times [T_{\text{Null}}]
 \end{aligned}$$

where $t_i \in \mathcal{T}$ is a corresponding tuple of g_i s.t. $\Phi_{\text{sdt}}(g_i) = t_i$.

(b) If the predicate $|v_1(g_i)| = 0$ does not hold, then, there exists a graph $g'_j \in \llbracket PP \rrbracket_G$ s.t. $\llbracket \phi \rrbracket_{G, [\text{merge}(g_i, g'_j)]} = \top$. By Figure 19,

$$\begin{aligned} \Phi_{\text{sdt}}(\text{ite}(|v_1(g_i)| = 0, v_2(g_i), v_1(g_i))) &= \Phi_{\text{sdt}}(v_1(g_i)) \\ &= \text{filter}(\lambda x. \llbracket \phi' \rrbracket_{\Phi_{\text{sdt}}(G), [x]} = \top, \text{map}(\lambda t'. \text{merge}(t_i, t'), \\ &\quad \text{filter}(\lambda t''. t \cap t'' \neq \emptyset, \Phi_{\text{sdt}}(\llbracket PP \rrbracket_G))) \\ &= \text{filter}(\lambda x. \llbracket \phi' \rrbracket_{R, [x]} = \top, \text{map}(\lambda t'. \text{merge}(t_i, t'), \\ &\quad \text{filter}(\lambda t''. t \cap t'' \neq \emptyset, \llbracket Q_2 \rrbracket_R))) \\ &= \sigma_{\phi'}([t_i] \bowtie_{\phi'} \llbracket Q_2 \rrbracket_R) \\ &= [t_i] \bowtie_{\phi' \wedge \phi''} \llbracket Q_2 \rrbracket_R \end{aligned}$$

where $t_i \in \mathcal{T}$ is a corresponding tuple of g_i s.t. $\Phi_{\text{sdt}}(g_i) = t_i$, and $\phi'' = \wedge_{(X:l) \in \mathcal{X}_1 \cap \mathcal{X}_2} T_1.\xi_{pk}(\Lambda(l)) = T_2.\xi_{pk}(\Lambda(l))$.

Thus, $\Phi_{\text{sdt}}(\text{ite}(|v_1(g_i)| = 0, v_2(g_i), v_1(g_i))) = [t_i] \bowtie_{\phi' \wedge \wedge_{(X:l) \in \mathcal{X}_1 \cap \mathcal{X}_2} T_1.\xi_{pk}(\Lambda(l)) = T_2.\xi_{pk}(\Lambda(l))} \llbracket Q_2 \rrbracket_R$ and

$$\begin{aligned} \Phi_{\text{sdt}}(\llbracket \text{OptMatch}(C_1, PP, \phi) \rrbracket_G) &= \Phi_{\text{sdt}}(\text{foldl}(\lambda g.s. \lambda g. gs ++ \text{ite}(|v_1(g_i)| = 0, v_2(g_i), v_1(g_i)), \\ &\quad [], \llbracket C_1 \rrbracket_G)) \\ &= \text{foldl}(\lambda xs. \lambda x. xs ++ \Phi_{\text{sdt}}(\text{ite}(|v_1(g_i)| = 0, v_2(g_i), v_1(g_i))), \\ &\quad [], \Phi_{\text{sdt}}(\llbracket C_1 \rrbracket_G)) \\ &= \text{foldl}(\lambda xs. \lambda x. xs ++ [t_i] \bowtie_{\phi' \wedge \phi''} \llbracket Q_2 \rrbracket_R, [], \llbracket Q_1 \rrbracket_R) \\ &= \llbracket Q_1 \bowtie_{\phi' \wedge \phi''} Q_2 \rrbracket_R \\ &= \llbracket \rho_{T_1}(Q_1) \bowtie_{\phi' \wedge \phi''} \rho_{T_2}(Q_2) \rrbracket_R \end{aligned}$$

(4) Inductive case: $C = \text{With}(C_1, \bar{Y}, \bar{Z})$.

By the inductive hypothesis, we have $\Phi_{\text{sdt}}(\llbracket C_1 \rrbracket_G) = \llbracket \Pi_L(Q) \rrbracket_R$ where $\Phi_{\text{sdt}}, \Psi_R \vdash C_1 \xrightarrow{\text{clause}} \mathcal{X}, \Pi_L(Q)$ and L is an alias list. By the definition of $P[\bar{Y} \mapsto \bar{Z}]$ and $T[\bar{Y} \mapsto \bar{Z}]$, we know, for any variable $(Y_i, l_i) \in \mathcal{X}$ and $Y_i \in \bar{Y}$, $P(Y_i, k) = P(Z_i, k)$ and $T(Y_i) = T(Z_i)$ where k is a property key from the node or edge of the label l , i.e., $k \in \text{Keys}(l)$. Therefore, this process is equal to the renaming operation in SQL, i.e., $\rho_{L[\bar{Z}/\bar{Y}]}$. By Figure 19,

$$\begin{aligned} \Phi_{\text{sdt}}(\llbracket \text{With}(C_1, \bar{Y}, \bar{Z}) \rrbracket_G) &= \Phi_{\text{sdt}}(\text{map}(\lambda (N, E, P, T). (N, E, P[\bar{X} \mapsto \bar{Y}], T[\bar{X} \mapsto \bar{Y}]), \llbracket C_1 \rrbracket_G)) \\ &= \llbracket \Pi_{\rho_{L[\bar{Z}/\bar{Y}]}(Q)} \rrbracket_R \end{aligned}$$

□

LEMMA F.3. Let Ψ_G be a graph schema and PP be a Cypher pattern. Let Ψ_R be a relational schema and Φ_{sdt} be the standard database transformer from Ψ_G to Ψ_R . If $\Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} \mathcal{X}, Q$, it holds that⁹

$$\forall G, R. G \triangleright \Psi_G \wedge R \triangleright \Psi_R \wedge G \sim_{\Phi_{\text{sdt}}} R \Rightarrow \Phi_{\text{sdt}}(\llbracket PP \rrbracket_G) = \llbracket Q \rrbracket_R$$

PROOF. Prove by structural induction on PP .

(1) Base case: $PP = NP$.

Suppose that $G = (N, E, P, T)$ and $NP = (X, l)$. By Figure 19, there exists a list of graphs $\llbracket NP \rrbracket_G = \text{map}(\lambda n. (\{n\}, \emptyset, \{(X, k) \mapsto P(n, k) \mid k \in \text{keys}(T(n))\}, \{X \mapsto (l, \text{keys}(T(n)))\}, [n \in N \mid \text{label}(T(n)) = l]))$ and a list of tuples $\mathcal{T} = \rho_X(R(\Lambda(l)))$ s.t. $\forall i. \Phi_{\text{sdt}}(g_i) = t_i$ where $g_i \in \llbracket NP \rrbracket_G$ is a graph derived from G and $t_i \in \mathcal{T}$ is a tuple derived from the relational table $\Lambda(l)$. Therefore, $\Phi_{\text{sdt}}(\llbracket NP \rrbracket_G) = \rho_X(R(\Lambda(l))) = \llbracket \rho_X(\Lambda(l)) \rrbracket_R$.

⁹We slightly abuse Φ_{sdt} over lists to apply Φ_{sdt} to each element in list $\llbracket PP \rrbracket_G$.

(2) Inductive case: $PP = NP_1, EP_2, PP'$.

Suppose that $NP = (X_1, l_1)$, $EP = (X_2, l_2, d_2)$, $\text{head}(PP') = NP_3 = (X_3, l_3)$. By the inductive hypothesis, we know $\Phi_{\text{sdt}}(\llbracket NP_1 \rrbracket_G) = \llbracket \rho_{X_1}(\Lambda(l_1)) \rrbracket_R$ and $\Phi_{\text{sdt}}(\llbracket NP_3 \rrbracket_G) = \llbracket \rho_{X_3}(\Lambda(l_3)) \rrbracket_R$, and $\Phi_{\text{sdt}}(\llbracket PP' \rrbracket_G) = \llbracket Q' \rrbracket_R$ where $\Phi_{\text{sdt}}, \Psi_R \vdash PP' \xrightarrow{\text{pattern}} \mathcal{X}, Q'$.

Let us first discuss $\text{SubGraphs}(G, [NP_1, EP_2, NP_3])$. By its definition, we are supposed to get a list of subgraphs of G matching the pattern $[NP_1, EP_2, NP_3]$, and $\Phi_{\text{sdt}}(\text{SubGraphs}(G, [NP_1, EP_2, NP_3])) = \llbracket \rho_{X_1}(\Lambda(l_1)) \bowtie_{\phi_1} \rho_{X_2}(\Lambda(l_2)) \bowtie_{\phi_2} \rho_{X_3}(\Lambda(l_3)) \rrbracket_R$ where $\phi_1 = \text{link}(\xi, \Lambda(l_1), \Lambda(l_2))$ and $\phi_2 = \text{link}(\xi, \Lambda(l_2), \Lambda(l_3))$ by the definition of link in Figure 18.

We assume that variables in NP and EP are unique. Therefore, for any graphs g'' derived from NP_1, EP_2, NP_3 and g' derived from PP' , we are only able to merge them to a graph without isolated nodes and the overlapping node for g'' and g' must be the node derived from NP_3 ; otherwise, we cannot return such a merged graph as an element of $\text{SubGraphs}(G, [NP_1, EP_2, NP_3])$. Furthermore, for any graph $g' \in \llbracket PP' \rrbracket_G$, we have

$$\begin{aligned} & \Phi_{\text{sdt}}(\text{map}(\lambda g. \text{merge}(g, g'), \text{filter}(\lambda g'', g'' \cap g' \neq \emptyset, \text{SubGraphs}(G, [NP_1, EP_2, NP_3]))) \\ &= \Phi_{\text{sdt}}(\text{map}(\lambda g. \text{merge}(g, g'), \text{filter}(\lambda g'', g'' \cap g' \neq \emptyset, [g'_1, \dots, g'_n]))) \\ &= \text{map}(\lambda t. \text{merge}(t, t'), \text{filter}(\lambda t'', \llbracket \phi_3 \rrbracket_{D, [t''] \times [t']}, [t'_1, \dots, t'_n]))) \\ &= \llbracket [t'_1, \dots, t'_n] \bowtie_{\phi_3} [t'] \rrbracket_R \\ &= \llbracket \rho_{X_1}(\Lambda(l_1)) \bowtie_{\phi_1} \rho_{X_2}(\Lambda(l_2)) \bowtie_{\phi_2} \rho_{X_3}(\Lambda(l_3)) \bowtie_{\phi_3} [t'] \rrbracket_R \\ &= \llbracket \rho_{X_1}(\Lambda(l_1)) \bowtie_{\phi_1} \rho_{X_2}(\Lambda(l_2)) \bowtie_{\phi_2} [t'] \rrbracket_R \end{aligned}$$

where $\Phi_{\text{sdt}}(g') = t'$ and $\phi_3 = \text{link}(\xi, \Lambda(l_3), \Lambda(l_3))$.

$$\begin{aligned} \Phi_{\text{sdt}}(\llbracket NP_1, EP_2, PP' \rrbracket_G) &= \Phi_{\text{sdt}}(\text{foldl}(\lambda gs. \lambda g'. gs ++ \text{map}(\lambda g. \text{merge}(g, g'), \text{filter}(\lambda g'', \\ & \quad g'' \cap g' \neq \emptyset, \text{SubGraphs}(G, [NP_1, EP_2, NP_3]))) , [], \llbracket PP' \rrbracket_G)) \\ &= \text{foldl}(\lambda xs. \lambda t'. xs ++ \llbracket \rho_{X_1}(\Lambda(l_1)) \bowtie_{\phi_1} \rho_{X_2}(\Lambda(l_2)) \bowtie_{\phi_2} [t'] \rrbracket_R, \\ & \quad [], \llbracket Q' \rrbracket_R) \\ &= \llbracket \rho_{X_1}(\Lambda(l_1)) \bowtie_{\phi_1} \rho_{X_2}(\Lambda(l_2)) \bowtie_{\phi_2} Q' \rrbracket_R \end{aligned}$$

□

LEMMA F.4. Let Ψ_G be a graph schema and E be a Cypher expression. Let Ψ_R be a relational schema and Φ_{sdt} be the standard database transformer from Ψ_G to Ψ_R . If $\Phi_{\text{sdt}}, \Psi_R \vdash E \xrightarrow{\text{expr}} E'$, it holds that

$$\forall G, R. G \triangleright \Psi_G \wedge R \triangleright \Psi_R \wedge G \sim_{\Phi_{\text{sdt}}} R \wedge (\forall i. \Phi_{\text{sdt}}(g_i) = t_i) \Rightarrow \llbracket E \rrbracket_{G, gs} = \llbracket E' \rrbracket_{R, \mathcal{T}}$$

where $g_i \in gs$ is a graph derived from G , and $t_i \in \mathcal{T}$ is a tuple derived from R .

PROOF. Prove by structural induction on E .

(1) Base case: $E = k$.

Lemma F.4 in this case is proved because $\llbracket k \rrbracket_{G, gs} = \text{lookup}(\text{head}(gs), k) = \text{lookup}(\text{head}(\mathcal{T}), k) = \llbracket k \rrbracket_{R, \mathcal{T}}$ by Figure 19.

(2) Base case: $E = v$.

Lemma F.4 in this case is proved because $\llbracket v \rrbracket_{G, gs} = v = \llbracket v \rrbracket_{R, \mathcal{T}}$ by Figure 19.

(3) Inductive case: $E = \text{Cast}(\phi)$.

By Lemma F.5, we have $\llbracket \phi \rrbracket_{G, gs} = \llbracket \phi' \rrbracket_{R, \mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi'$. By Figure 19, $\llbracket \text{Cast}(\phi) \rrbracket_{G, gs} = \text{ite}(\llbracket \phi \rrbracket_{G, gs} = \text{Null}, \text{Null}, \text{ite}(\llbracket \phi \rrbracket_{G, gs} = \top, 1, 0)) = \text{ite}(\llbracket \phi' \rrbracket_{R, \mathcal{T}} = \text{Null}, \text{Null}, \text{ite}(\llbracket \phi' \rrbracket_{R, \mathcal{T}} = \top, 1, 0)) = \llbracket \text{Cast}(\phi') \rrbracket_{R, \mathcal{T}}$.

(4) Inductive case: $E = \text{Count}(E_1)$.

By the inductive hypothesis, we have $\llbracket E_1 \rrbracket_{G,gs} = \llbracket E'_1 \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \xrightarrow{\text{expr}} E'_1$. By Figure 19,

$$\begin{aligned} \llbracket \text{Count}(E_1) \rrbracket_{G,gs} &= \text{ite}(\wedge_{g \in gs} \llbracket E_1 \rrbracket_{G,[g]} = \text{Null}, \text{Null}, \\ &\quad \text{foldl}(+, 0, \text{map}(\lambda g. \text{ite}(\llbracket E_1 \rrbracket_{G,[g]} = \text{Null}, 0, 1), gs))) \\ &= \text{ite}(\wedge_{t \in \mathcal{T}} \llbracket E'_1 \rrbracket_{R,[t]} = \text{Null}, \text{Null}, \\ &\quad \text{foldl}(+, 0, \text{map}(\lambda t. \text{ite}(\llbracket E'_1 \rrbracket_{R,[t]} = \text{Null}, 0, 1), \mathcal{T}))) \\ &= \llbracket \text{Count}(E'_1) \rrbracket_{R,\mathcal{T}} \end{aligned}$$

(5) Inductive case: $E = \text{Sum}(E_1)$.

By the inductive hypothesis, we have $\llbracket E_1 \rrbracket_{G,gs} = \llbracket E'_1 \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \xrightarrow{\text{expr}} E'_1$. By Figure 19,

$$\begin{aligned} \llbracket \text{Sum}(E_1) \rrbracket_{G,gs} &= \text{ite}(\wedge_{g \in gs} \llbracket E_1 \rrbracket_{G,[g]} = \text{Null}, \text{Null}, \\ &\quad \text{foldl}(+, 0, \text{map}(\lambda g. \text{ite}(\llbracket E_1 \rrbracket_{G,[g]} = \text{Null}, 0, \llbracket E_1 \rrbracket_{G,[g]}), gs))) \\ &= \text{ite}(\wedge_{t \in \mathcal{T}} \llbracket E'_1 \rrbracket_{R,[t]} = \text{Null}, \text{Null}, \\ &\quad \text{foldl}(+, 0, \text{map}(\lambda t. \text{ite}(\llbracket E'_1 \rrbracket_{R,[t]} = \text{Null}, 0, \llbracket E'_1 \rrbracket_{R,[t]}), \mathcal{T}))) \\ &= \llbracket \text{Sum}(E'_1) \rrbracket_{R,\mathcal{T}} \end{aligned}$$

(6) Inductive case: $E = \text{Avg}(E_1)$.

By the inductive hypothesis, we have $\llbracket \text{Sum}(E_1) \rrbracket_{G,gs} = \llbracket \text{Sum}(E'_1) \rrbracket_{R,\mathcal{T}}$ and $\llbracket \text{Count}(E_1) \rrbracket_{G,gs} = \llbracket \text{Count}(E'_1) \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash E \xrightarrow{\text{expr}} E'$. By Figure 19,

$$\begin{aligned} \llbracket \text{Sum}(E_1) \rrbracket_{G,gs} &= \llbracket \text{Sum}(E_1) \rrbracket_{G,gs} / \llbracket \text{Count}(E_1) \rrbracket_{G,gs} \\ &= \llbracket \text{Sum}(E'_1) \rrbracket_{R,\mathcal{T}} / \llbracket \text{Count}(E'_1) \rrbracket_{R,\mathcal{T}} \\ &= \llbracket \text{Sum}(E'_1) \rrbracket_{R,\mathcal{T}} \end{aligned}$$

(7) Inductive case: $E = \text{Min}(E_1)$.

By the inductive hypothesis, we have $\llbracket E_1 \rrbracket_{G,gs} = \llbracket E'_1 \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \xrightarrow{\text{expr}} E'_1$. By Figure 19,

$$\begin{aligned} \llbracket \text{Min}(E_1) \rrbracket_{G,gs} &= \text{ite}(\wedge_{g \in gs} \llbracket E_1 \rrbracket_{G,[g]} = \text{Null}, \text{Null}, \\ &\quad \text{foldl}(\text{min}, +\infty, \text{map}(\lambda g. \text{ite}(\llbracket E_1 \rrbracket_{G,[g]} = \text{Null}, +\infty, \llbracket E_1 \rrbracket_{G,[g]}), gs))) \\ &= \text{ite}(\wedge_{t \in \mathcal{T}} \llbracket E'_1 \rrbracket_{R,[t]} = \text{Null}, \text{Null}, \\ &\quad \text{foldl}(\text{min}, +\infty, \text{map}(\lambda t. \text{ite}(\llbracket E'_1 \rrbracket_{R,[t]} = \text{Null}, +\infty, \llbracket E'_1 \rrbracket_{R,[t]}), \mathcal{T}))) \\ &= \llbracket \text{Min}(E'_1) \rrbracket_{R,\mathcal{T}} \end{aligned}$$

(8) Inductive case: $E = \text{Max}(E_1)$.

By the inductive hypothesis, we have $\llbracket E_1 \rrbracket_{G,gs} = \llbracket E'_1 \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \xrightarrow{\text{expr}} E'_1$. By Figure 19,

$$\begin{aligned} \llbracket \text{Max}(E_1) \rrbracket_{G,gs} &= \text{ite}(\wedge_{g \in gs} \llbracket E_1 \rrbracket_{G,[g]} = \text{Null}, \text{Null}, \\ &\quad \text{foldl}(\text{max}, -\infty, \text{map}(\lambda g. \text{ite}(\llbracket E_1 \rrbracket_{G,[g]} = \text{Null}, -\infty, \llbracket E_1 \rrbracket_{G,[g]}), gs))) \\ &= \text{ite}(\wedge_{t \in \mathcal{T}} \llbracket E'_1 \rrbracket_{R,[t]} = \text{Null}, \text{Null}, \\ &\quad \text{foldl}(\text{max}, -\infty, \text{map}(\lambda t. \text{ite}(\llbracket E'_1 \rrbracket_{R,[t]} = \text{Null}, -\infty, \llbracket E'_1 \rrbracket_{R,[t]}), \mathcal{T}))) \\ &= \llbracket \text{Max}(E'_1) \rrbracket_{R,\mathcal{T}} \end{aligned}$$

(9) Inductive case: $E = E_1 \oplus E_2$.

By the inductive hypothesis, we have $\llbracket E_1 \rrbracket_{G,gs} = \llbracket E'_1 \rrbracket_{R,\mathcal{T}}$ and $\llbracket E_2 \rrbracket_{G,gs} = \llbracket E'_2 \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \xrightarrow{\text{expr}} E'_1$ and $\Phi_{\text{sdt}}, \Psi_R \vdash E_2 \xrightarrow{\text{expr}} E'_2$. By Figure 19, $\llbracket E_1 \oplus E_2 \rrbracket_{G,gs} = \llbracket E_1 \rrbracket_{G,gs} \oplus \llbracket E_2 \rrbracket_{G,gs} = \llbracket E'_1 \rrbracket_{R,\mathcal{T}} \oplus \llbracket E'_2 \rrbracket_{R,\mathcal{T}} = \llbracket E'_1 \oplus E'_2 \rrbracket_{R,\mathcal{T}}$.

□

LEMMA F.5. Let Ψ_G be a graph schema and ϕ be a Cypher predicate. Let Ψ_R be a relational schema and Φ_{sdt} be the standard database transformer from Ψ_G to Ψ_R . If $\Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi'$, it holds that

$$\forall G, R. G \triangleright \Psi_G \wedge R \triangleright \Psi_R \wedge G \sim_{\Phi_{\text{sdt}}} R \wedge (\forall i. \Phi_{\text{sdt}}(g_i) = t_i) \Rightarrow \llbracket \phi \rrbracket_{G,gs} = \llbracket \phi' \rrbracket_{R,\mathcal{T}}$$

where $g_i \in gs$ is a graph derived from G , and $t_i \in \mathcal{T}$ is a tuple derived from R .

PROOF. Prove by structural induction on ϕ .

- (1) Base case: $\phi = \top$.

Lemma F.5 in this case is proved because $\llbracket \top \rrbracket_{G,gs} = \top = \llbracket \top \rrbracket_{R,\mathcal{T}}$ by Figure 19.

- (2) Base case: $\phi = \perp$.

Lemma F.5 in this case is proved because $\llbracket \perp \rrbracket_{G,gs} = \perp = \llbracket \perp \rrbracket_{R,\mathcal{T}}$ by Figure 19.

- (3) Inductive case: $\phi = E_1 \odot E_2$.

By Lemma F.4, we have $\llbracket E_1 \rrbracket_{G,gs} = \llbracket E'_1 \rrbracket_{R,\mathcal{T}}$ and $\llbracket E_2 \rrbracket_{G,gs} = \llbracket E'_2 \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \xrightarrow{\text{expr}} E'_1$ and $\Phi_{\text{sdt}}, \Psi_R \vdash E_2 \xrightarrow{\text{expr}} E'_2$. By Figure 19, $\llbracket E_1 \odot E_2 \rrbracket_{G,gs} = \llbracket E_1 \rrbracket_{G,gs} \odot \llbracket E_2 \rrbracket_{G,gs} = \llbracket E'_1 \rrbracket_{R,\mathcal{T}} \odot \llbracket E'_2 \rrbracket_{R,\mathcal{T}} = \llbracket E'_1 \odot E'_2 \rrbracket_{R,\mathcal{T}}$.

- (4) Inductive case: $\phi = \text{IsNull}(E)$.

By Lemma F.4, we have $\llbracket E \rrbracket_{G,gs} = \llbracket E \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash E \xrightarrow{\text{expr}} E'$. By Figure 19, $\llbracket \text{IsNull}(E) \rrbracket_{G,gs} = \text{ite}(\llbracket E \rrbracket_{G,gs} = \text{Null}, \top, \perp) = \text{ite}(\llbracket E' \rrbracket_{R,\mathcal{T}} = \text{Null}, \top, \perp) = \llbracket \text{IsNull}(E') \rrbracket_{R,\mathcal{T}}$.

- (5) Inductive case: $\phi = E \in \bar{v}$.

By Lemma F.4, we have $\llbracket E \rrbracket_{G,gs} = \llbracket E' \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash E \xrightarrow{\text{expr}} E'$. By Figure 19, $\llbracket E \in \bar{v} \rrbracket_{G,gs} = \bigvee_{v \in \bar{v}} \llbracket E \rrbracket_{G,gs} = \bigvee_{v \in \bar{v}} \llbracket E' \rrbracket_{R,\mathcal{T}} = \llbracket E' \in \bar{v} \rrbracket_{R,\mathcal{T}}$.

- (6) Inductive case: $\phi = \text{Exists}(PP)$.

By Lemma F.3, we have $\Phi_{\text{sdt}}(\llbracket PP \rrbracket_G) = \llbracket Q \rrbracket_R$ where $\Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} \mathcal{X}, Q$. Also, by Lemma F.4, we have $\llbracket K \rrbracket_{G,gs} = \llbracket K \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash K \xrightarrow{\text{expr}} K$ and K is property key. Suppose that $\text{head}(PP) = (X_1, l_1)$ and $\text{last}(PP) = (X_2, l_2)$. By the definition of PKs and ξ_{pk} , $\llbracket \bar{K} \rrbracket_{G,gs} = \llbracket [\text{PK}(\text{head}(PP)), \text{PK}(\text{last}(PP))] \rrbracket_{G,gs} = \llbracket [\xi_{pk}(\Lambda(l_1)), \xi_{pk}(\Lambda(l_2))] \rrbracket_{R,\mathcal{T}} = \llbracket \bar{K} \rrbracket_{R,\mathcal{T}}$ iff $\forall i. \Phi_{\text{sdt}}(g_i) = t_i$ for any $g_i \in gs$ and $t_i \in \mathcal{T}$. Furthermore, $\text{map}(\lambda g. \llbracket \bar{K} \rrbracket_{G,[g]}, \llbracket PP \rrbracket_G) = \llbracket \bar{K} \rrbracket_{G,[g_1]}, \dots, \llbracket \bar{K} \rrbracket_{G,[g_n]} = \llbracket \bar{K} \rrbracket_{R,[t_1]}, \dots, \llbracket \bar{K} \rrbracket_{R,[t_n]} = \text{map}(\lambda t. \llbracket \bar{K} \rrbracket_{R,[t]}, \llbracket Q \rrbracket_R) = \llbracket \Pi_{\bar{K}}(Q) \rrbracket_R$ where $n = |\llbracket PP \rrbracket_G| = |\llbracket Q \rrbracket_R|$.

By Figure 19,

$$\begin{aligned} \llbracket \text{Exists}(PP) \rrbracket_{G,gs} &= \bigvee_{g \in \llbracket PP \rrbracket_G} \bigwedge_{K \in \bar{K}} \llbracket K \rrbracket_{G,gs} = \llbracket K \rrbracket_{G,[g]} \\ &= \bigvee_{g \in \llbracket PP \rrbracket_G} \bigwedge_{K \in \bar{K}} \llbracket K \rrbracket_{R,\mathcal{T}} = \llbracket K \rrbracket_{G,[g]} \\ &= \bigvee_{g \in \llbracket PP \rrbracket_G} \llbracket \bar{K} \in \llbracket \bar{K} \rrbracket_{G,[g]} \rrbracket_{R,\mathcal{T}} \\ &= \bigvee_{t \in \text{map}(\lambda g. \llbracket \bar{K} \rrbracket_{G,[g]}, \llbracket PP \rrbracket_G)} \llbracket \bar{K} \in t \rrbracket_{R,\mathcal{T}} \\ &= \text{foldl}(\lambda y. \lambda y. y \vee \llbracket \bar{K} \in y \rrbracket_{R,\mathcal{T}}, \perp, \text{map}(\lambda g. \llbracket \bar{K} \rrbracket_{G,[g]}, \llbracket PP \rrbracket_G)) \\ &= \text{foldl}(\lambda y. \lambda y. y \vee \llbracket \bar{K} \in y \rrbracket_{R,\mathcal{T}}, \perp, \llbracket \Pi_{\bar{K}}(Q) \rrbracket_R) \\ &= \llbracket \bar{K} \in \Pi_{\bar{K}}(Q) \rrbracket_{R,\mathcal{T}} \end{aligned}$$

- (7) Inductive case: $\phi = \phi_1 \wedge \phi_2$.

By the inductive hypothesis, we have $\llbracket \phi_1 \rrbracket_{G,gs} = \llbracket \phi'_1 \rrbracket_{R,\mathcal{T}}$ and $\llbracket \phi_2 \rrbracket_{G,gs} = \llbracket \phi'_2 \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash \phi_1 \xrightarrow{\text{pred}} \phi'_1$ and $\Phi_{\text{sdt}}, \Psi_R \vdash \phi_2 \xrightarrow{\text{pred}} \phi'_2$. By Figure 19, $\llbracket \phi_1 \wedge \phi_2 \rrbracket_{G,gs} = \llbracket \phi_1 \rrbracket_{G,gs} \wedge \llbracket \phi_2 \rrbracket_{G,gs} = \llbracket \phi'_1 \rrbracket_{R,\mathcal{T}} \wedge \llbracket \phi'_2 \rrbracket_{R,\mathcal{T}} = \llbracket \phi'_1 \wedge \phi'_2 \rrbracket_{R,\mathcal{T}}$.

- (8) Inductive case: $\phi = \phi_1 \vee \phi_2$.

By the inductive hypothesis, we have $\llbracket \phi_1 \rrbracket_{G,gs} = \llbracket \phi'_1 \rrbracket_{R,\mathcal{T}}$ and $\llbracket \phi_2 \rrbracket_{G,gs} = \llbracket \phi'_2 \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash \phi_1 \xrightarrow{\text{pred}} \phi'_1$ and $\Phi_{\text{sdt}}, \Psi_R \vdash \phi_2 \xrightarrow{\text{pred}} \phi'_2$. By Figure 19, $\llbracket \phi_1 \vee \phi_2 \rrbracket_{G,gs} = \llbracket \phi_1 \rrbracket_{G,gs} \vee \llbracket \phi_2 \rrbracket_{G,gs} = \llbracket \phi'_1 \rrbracket_{R,\mathcal{T}} \vee \llbracket \phi'_2 \rrbracket_{R,\mathcal{T}} = \llbracket \phi'_1 \vee \phi'_2 \rrbracket_{R,\mathcal{T}}$.

(9) Inductive case: $\phi = \neg \phi_1$.

By the inductive hypothesis, we have $\llbracket \phi_1 \rrbracket_{G,gs} = \llbracket \phi'_1 \rrbracket_{R,\mathcal{T}}$ where $\Phi_{\text{sdt}}, \Psi_R \vdash \phi_1 \xrightarrow{\text{pred}} \phi'_1$. By Figure 19, $\llbracket \neg \phi_1 \rrbracket_{G,gs} = \neg \llbracket \phi_1 \rrbracket_{G,gs} = \neg \llbracket \phi'_1 \rrbracket_{R,\mathcal{T}} = \llbracket \neg \phi'_1 \rrbracket_{R,\mathcal{T}}$.

□

THEOREM F.6 (COMPLETENESS OF TRANSLATION (5.8)). *Let Ψ_G be a graph schema and Ψ_R be the induced relational schema of Ψ_G . Given any Cypher query Q over Ψ_G accepted by the grammar shown in Figure 9, there exists a SQL query Q' over Ψ_R such that $\Phi_{\text{sdt}}, \Psi_R \vdash Q \xrightarrow{\text{query}} Q'$.*

PROOF. Prove by structural induction on Q .

(1) Base case: $Q = \text{Return}(C, \bar{E}, \bar{k})$.

By Lemma F.7, there exist a variable set \mathcal{X}' and a SQL query Q' such that $\Phi_{\text{sdt}}, \Psi_R \vdash C \xrightarrow{\text{pred}} \mathcal{X}', Q'$. Also, by Lemma F.9, for any expression $E \in \bar{E}$ there exists a SQL expression E' such that $\Phi_{\text{sdt}}, \Psi_R \vdash E \xrightarrow{\text{pred}} E'$. Let us discuss the predicate $\neg \text{hasAgg}(\bar{E})$ in two cases.

(a) If $\neg \text{hasAgg}(\bar{E}) = \top$, then for any expression $E \in \bar{E}$, $\text{IsAgg}(E) = \perp$. According to the Q-Ret rule in Figure 16, there exists a SQL query $Q'' = \Pi_{\rho_{\bar{E}}(\bar{E}')} (Q')$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \text{Return}(C, \bar{E}, \bar{k}) \xrightarrow{\text{query}} Q''$.

(b) If $\neg \text{hasAgg}(\bar{E}) = \perp$, then there exist some expression $E \in \bar{E}$ s.t. $\text{IsAgg}(E) = \top$. Therefore, this Return statement in Cypher should be translated in GroupBy in SQL. Let $\bar{A} = \text{filter}(\lambda x. \neg \text{IsAgg}(x), \bar{E}')$ be the list of SQL expressions containing no aggregation function. According to the Q-Agg rule in Figure 16, there exists a SQL query $Q'' = \text{GroupBy}(Q, \bar{A}, \rho_{\bar{k}}(\bar{E}'), \top)$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \text{Return}(C, \bar{E}, \bar{k}) \xrightarrow{\text{query}} Q''$.

Thus, Theorem 5.8 in this case is proved.

(2) Inductive case: $Q = \text{OrderBy}(R_1, k, b)$.

Specifically, R_1 denotes the Return statement of Cypher. By the inductive hypothesis, we have $\Phi_{\text{sdt}}, \Psi_R \vdash R_1 \xrightarrow{\text{query}} Q'_1$. Also, by Lemma F.9, there exists a SQL expression $E' = k$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash k \xrightarrow{\text{expr}} E'$. By Lemma F.10, there exists a SQL predicate $b' = b$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash b \xrightarrow{\text{expr}} b'$. According to the Q-OrderBy rule in Figure 16, there exists a SQL query $Q'_1 = \text{OrderBy}(Q'_1, E', b')$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \text{OrderBy}(R_1, k, b) \xrightarrow{\text{expr}} Q'_1$.

(3) Inductive case: $Q = \text{Union}(Q_1, Q_2)$.

By the inductive hypothesis, there exist two queries Q'_1 and Q'_2 such that $\Phi_{\text{sdt}}, \Psi_R \vdash Q_1 \xrightarrow{\text{query}} Q'_1$ and $\Phi_{\text{sdt}}, \Psi_R \vdash Q_2 \xrightarrow{\text{query}} Q'_2$. According to the Q-Union rule in Figure 16, there exists a SQL query $Q' = Q'_1 \cup Q'_2$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \text{Union}(Q_1, Q_2) \xrightarrow{\text{expr}} Q'$.

(4) Inductive case: $Q = \text{UnionAll}(Q_1, Q_2)$.

By the inductive hypothesis, there exist two queries Q'_1 and Q'_2 such that $\Phi_{\text{sdt}}, \Psi_R \vdash Q_1 \xrightarrow{\text{query}} Q'_1$ and $\Phi_{\text{sdt}}, \Psi_R \vdash Q_2 \xrightarrow{\text{query}} Q'_2$. According to the Q-Union rule in Figure 16, there exists a SQL query $Q' = Q'_1 \uplus Q'_2$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \text{UnionAll}(Q_1, Q_2) \xrightarrow{\text{expr}} Q'$.

□

LEMMA F.7. Let Ψ_G be a graph schema, $\Psi_R = (S, \xi)$ be the induced relational schema of Ψ_G , and Λ be the corresponding schema mapping. Given any Cypher clause C accepted by the grammar shown in Figure 9, there exists a SQL query Q over Ψ_R such that $\Phi_{\text{sdt}}, \Psi_R \vdash C \xrightarrow{\text{clause}} \mathcal{X}, Q$.

PROOF. Prove by structural induction on C .

- (1) Base case: $C = \text{Match}(PP, \phi)$.

By Lemma F.8, there exist a variable set \mathcal{X} and a SQL query Q such that $\Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} \mathcal{X}, Q$. Also, by Lemma F.10, there exists a SQL predicate ϕ' such that $\Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi'$. According to the C-Match1 rule in Figure 17, there exists a SQL query $Q' = \sigma_{\phi'}(Q)$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \text{Match}(PP, \phi) \xrightarrow{\text{clause}} \mathcal{X}, Q'$.

- (2) Inductive case: $C = \text{Match}(C_1, PP, \phi)$.

By the inductive hypothesis, there exist a variable set \mathcal{X}_1 and a SQL query Q_1 such that $\Phi_{\text{sdt}}, \Psi_R \vdash C_1 \xrightarrow{\text{clause}} \mathcal{X}_1, Q_1$. By Lemma F.8, there exist a variable set \mathcal{X}_2 and a SQL query Q_2 such that $\Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{clause}} \mathcal{X}_2, Q_2$. Also, by Lemma F.10, there exists a SQL predicate ϕ' such that $\Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi'$. Let T_1 and T_2 be two fresh name for the output tables of Q_1 and Q_2 . Then we have a new predicate $\phi'' = \phi' \wedge \bigwedge_{(X:I) \in \mathcal{X}_1 \cap \mathcal{X}_2} T_1.\xi_{pk}(\Lambda(I)) = \xi_{pk}(\Lambda(I))$ to join Q_1 and Q_2 . According to the C-Match2 rule in Figure 17, there exists a SQL query $Q' = \sigma_{\phi'}(Q)$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \text{Match}(C_1, PP, \phi) \xrightarrow{\text{clause}} \mathcal{X}_1 \cup \mathcal{X}_2, \rho_{T_1}(Q_1) \bowtie_{\phi''} \rho_{T_2}(Q_2)$.

- (3) Inductive case: $C = \text{OptMatch}(C_1, PP, \phi)$.

By the inductive hypothesis, there exist a variable set \mathcal{X}_1 and a SQL query Q_1 such that $\Phi_{\text{sdt}}, \Psi_R \vdash C_1 \xrightarrow{\text{clause}} \mathcal{X}_1, Q_1$. By Lemma F.8, there exist a variable set \mathcal{X}_2 and a SQL query Q_2 such that $\Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{clause}} \mathcal{X}_2, Q_2$. Also, by Lemma F.10, there exists a SQL predicate ϕ' such that $\Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi'$. Let T_1 and T_2 be two fresh name for the output tables of Q_1 and Q_2 . Then we have a new predicate $\phi'' = \phi' \wedge \bigwedge_{(X:I) \in \mathcal{X}_1 \cap \mathcal{X}_2} T_1.\xi_{pk}(\Lambda(I)) = \xi_{pk}(\Lambda(I))$ to join Q_1 and Q_2 . According to the C-Match2 rule in Figure 17, there exists a SQL query $Q' = \sigma_{\phi'}(Q)$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \text{Match}(C_1, PP, \phi) \xrightarrow{\text{clause}} \mathcal{X}_1 \cup \mathcal{X}_2, \rho_{T_1}(Q_1) \bowtie_{\phi''} \rho_{T_2}(Q_2)$.

- (4) Inductive case: $C = \text{With}(C_1, \bar{Y}, \bar{Z})$.

By the inductive hypothesis, we have $\Phi_{\text{sdt}}, \Psi_R \vdash C_1 \xrightarrow{\text{clause}} \mathcal{X}_1, \Pi_L(Q_1)$ where L denotes a list of columns. According to the C-With rule in Figure 17, there exists a SQL query $Q'_1 = \Pi_{\rho_L(\bar{Z}/\bar{Y})}(Q_1)$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \text{Match}(C_1, PP, \phi) \xrightarrow{\text{clause}} \mathcal{X}_1 \setminus \bar{Y} \cup \bar{Z}, Q'_1$.

□

LEMMA F.8. Let Ψ_G be a graph schema, $\Psi_R = (S, \xi)$ be the induced relational schema of Ψ_G , and Λ be the corresponding schema mapping. Given any Cypher pattern PP accepted by the grammar shown in Figure 9, there exists a SQL query Q over Ψ_R such that $\Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} \mathcal{X}, Q$.

PROOF. Prove by structural induction on PP .

- (1) Base case: $PP = NP = (X, l)$.

According to the PT-Node rule in Figure 18, there exist a variable set $\mathcal{X} = \{(X, l)\}$ and a SQL query $Q = \rho_X(\Lambda(l))$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash NP \xrightarrow{\text{pattern}} \mathcal{X}, Q$.

- (2) Inductive case: $PP = NP_1, EP_1, PP_1$.

Let $NP_1 = (X_1, l_1)$ be the first node pattern in PP , $EP = (X_2, l_2, d_2)$ be the first edge pattern in PP , and $\text{head}(PP_1) = (X_3, l_3)$ be the first node pattern in PP_1 .

By the inductive hypothesis, there exist a variable set $X_1 = \{(X_1, l_1)\}$ and a SQL query $Q_1 = \rho_{X_1}(\Lambda(l_1))$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash NP_1 \xrightarrow{\text{pattern}} X_1, Q_1$, a variable set $X_2 = \{(X_2, l_2)\}$ and a SQL query $Q_2 = \rho_{X_2}(\Lambda(l_2))$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash NP_2 \xrightarrow{\text{pattern}} X_2, Q_2$, and a variable set X_3 and a SQL query Q_3 such that $\Phi_{\text{sdt}}, \Psi_R \vdash PP_1 \xrightarrow{\text{pattern}} X_3, Q_3$. Let $\phi_1 = \text{link}(\xi, \Lambda(l_1), \Lambda(l_2))$ be a predicate to join the relations $\Lambda(l_1)$ and $\Lambda(l_2)$, and $\phi_2 = \text{link}(\xi, \Lambda(l_2), \Lambda(l_3))$ to join the relations $\Lambda(l_2)$ and $\Lambda(l_3)$. According to the PT-Path rule in Figure 18, there exists a variable set $X = X_1 \cup X_2 \cup X_3$ and a SQL expression $Q = Q_1 \bowtie_{\phi_1} Q_2 \bowtie_{\phi_2} Q_3$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} X, Q$. \square

LEMMA F.9. Let Ψ_G be a graph schema, $\Psi_R = (S, \xi)$ be the induced relational schema of Ψ_G , and Λ be the corresponding schema mapping. Given any Cypher expression E accepted by the grammar shown in Figure 9, there exists a SQL expression E' such that $\Phi_{\text{sdt}}, \Psi_R \vdash E \xrightarrow{\text{expr}} E'$.

PROOF. Prove by structural induction on E .

- (1) Base case: $E = k$.

According to the E-Prop rule in Figure 21, there exists a SQL expression $E' = k$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash k \xrightarrow{\text{expr}} E'$.

- (2) Base case: $E = v$.

According to the E-Value rule in Figure 21, there exists a SQL expression $E' = v$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash v \xrightarrow{\text{expr}} E'$.

- (3) Inductive case: $E = \text{Cast}(\phi)$.

By Lemma F.10, there exists a SQL predicate ϕ' such that $\Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi'$. According to the E-Pred rule in Figure 21, there exists a SQL expression $E' = \text{Cast}(\phi')$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \text{Cast}(\phi) \xrightarrow{\text{expr}} E'$.

- (4) Inductive case: $E = \text{Agg}(E_1)$ where $\text{Agg} \in \{\text{Count}, \text{Max}, \text{Min}, \text{Avg}, \text{Sum}\}$.

By the inductive hypothesis, there exists a SQL expression E'_1 such that $\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \xrightarrow{\text{expr}} E'_1$. Also, for each aggregation function in Cypher there exists an equivalent SQL counterpart. According to the E-Agg rule in Figure 21, there exists a SQL expression $E' = \text{Agg}(E'_1)$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \text{Agg}(E_1) \xrightarrow{\text{expr}} E'$.

- (5) Inductive case: $E = E_1 \oplus E_2$.

By the inductive hypothesis, there exists two SQL expressions E'_1 and E'_2 such that $\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \xrightarrow{\text{expr}} E'_1$ and $\Phi_{\text{sdt}}, \Psi_R \vdash E_2 \xrightarrow{\text{expr}} E'_2$. According to the E-Arith rule in Figure 21, there exists a SQL expression $E' = E'_1 \oplus E'_2$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \oplus E_2 \xrightarrow{\text{expr}} E'_1 \oplus E'_2$. \square

LEMMA F.10. Let Ψ_G be a graph schema, $\Psi_R = (S, \xi)$ be the induced relational schema of Ψ_G , and Φ be the corresponding database transformer. Given any Cypher predicate ϕ accepted by the grammar shown in Figure 9, there exists a SQL predicate ϕ' such that $\Phi_{\text{sdt}}, \Psi_R \vdash \phi \xrightarrow{\text{pred}} \phi'$.

PROOF. Prove by structural induction on ϕ .

- (1) Base case: $\phi = \top$.

According to the P-True rule in Figure 22, there exists a SQL predicate $\phi' = \top$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \top \xrightarrow{\text{pred}} \phi'$.

- (2) Base case: $\phi = \perp$.

According to the P-False rule in Figure 22, there exists a SQL predicate $\phi' = \perp$ such that

$$\Phi_{\text{sdt}}, \Psi_R \vdash \perp \xrightarrow{\text{pred}} \phi'.$$

- (3) Inductive case: $\phi = E_1 \odot E_2$.

By Lemma F.9, there exists two SQL expressions E'_1 and E'_2 such that $\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \xrightarrow{\text{expr}} E'_1$ and $\Phi_{\text{sdt}}, \Psi_R \vdash E_2 \xrightarrow{\text{expr}} E'_2$. According to the P-Logic rule in Figure 22, there exists a SQL predicate $\phi' = E'_1 \odot E'_2$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash E_1 \odot E_2 \xrightarrow{\text{pred}} \phi'$.

- (4) Inductive case: $\phi = \text{IsNull}(E)$.

By Lemma F.9, there exists a SQL expression E' such that $\Phi_{\text{sdt}}, \Psi_R \vdash E \xrightarrow{\text{expr}} E'$. According to the P-IsNull rule in Figure 22, there exists a SQL predicate $\phi' = \text{IsNull}(E')$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \text{IsNull}(E) \xrightarrow{\text{pred}} \phi'$.

- (5) Inductive case: $\phi = E \in \bar{v}$.

By Lemma F.9, there exists a SQL expression E' such that $\Phi_{\text{sdt}}, \Psi_R \vdash E \xrightarrow{\text{expr}} E'$. According to the P-In rule in Figure 22, there exists a SQL predicate $\phi' = E' \in \bar{v}$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash E \in \bar{v} \xrightarrow{\text{pred}} \phi'$.

- (6) Inductive case: $\phi = \text{Exists}(PP)$.

By Lemma F.8, there exists a SQL query Q such that $\Phi_{\text{sdt}}, \Psi_R \vdash PP \xrightarrow{\text{pattern}} X, Q$. Let (X_1, l_1) and (X_2, l_2) be the head node and the tail node of the path pattern PP , and $\bar{a} = [\xi_{pk}(\Lambda(l_1)), \xi_{pk}(\Lambda(l_2))]$ be a list of primary keys of those two nodes by the definition of Λ . According to the P-Exist rule in Figure 22, there exists a SQL predicate $\phi' = \bar{a} \in \Pi_{\bar{a}}(Q)$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \text{Exists}(PP) \xrightarrow{\text{pred}} \phi'$.

- (7) Inductive case: $\phi = \phi_1 \circ \phi_2$ where $\circ \in \{\wedge, \vee\}$.

By the inductive hypothesis, there exist two SQL predicates ϕ'_1 and ϕ'_2 such that $\Phi_{\text{sdt}}, \Psi_R \vdash \phi_1 \xrightarrow{\text{pred}} \phi'_1$ and $\Phi_{\text{sdt}}, \Psi_R \vdash \phi_2 \xrightarrow{\text{pred}} \phi'_2$. According to the P-AndOr rule in Figure 22, there exists a SQL predicate $\phi' = \phi'_1 \circ \phi'_2$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \phi_1 \circ \phi_2 \xrightarrow{\text{pred}} \phi'$.

- (8) Inductive case: $\phi = \neg\phi_1$.

By the inductive hypothesis, there exists a SQL predicate ϕ'_1 such that $\Phi_{\text{sdt}}, \Psi_R \vdash \phi_1 \xrightarrow{\text{pred}} \phi'_1$. According to the P-Not rule in Figure 22, there exists a SQL predicate $\phi' = \neg\phi'_1$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash \neg\phi_1 \xrightarrow{\text{pred}} \phi'$.

□

LEMMA F.11. *Let G be a graph database instance over graph schema Ψ_G , Φ_{sdt} be a standard database transformer between Ψ_G and Ψ_R , $\Psi_{R'}$ be the induced relational schema, and $D' \triangleright \Psi_{R'}$ be a relational database such that $D' = \Phi_{\text{sdt}}(G)$. Given a database transformer Φ from Ψ_G to Ψ_R , a residual database transformer Φ_{rdt} and a relational database instance $D \triangleright \Psi_R$ such that $D = \Phi(G)$, it holds that $D' \sim_{\Phi_{\text{rdt}}} D$.*

PROOF. By Algorithm 2, we know the residual database transformer is obtained by a syntactic substitution. Let us discuss the substitution into two cases.

- (1) For any node $N(l, a_1, \dots, a_n)$ with label l and values a_1, \dots, a_n of property keys K_1, \dots, K_n , there exists a corresponding clause $l(a_1, \dots, a_n) \rightarrow l'(a_1, \dots, a_n)$ in Φ_{sdt} . Therefore, σ and Φ_{rdt} are updated by $\sigma \cup \{l \mapsto l'\}$ and $\Phi_{\text{rdt}} \cup \Phi[l \mapsto l']$, respectively. Let $l(a_1, \dots, a_n) \rightarrow R_l(a'_1, \dots, a'_n)$ be a formula of Φ that represent the equivalence between the node N of l label and the table of R_l relation by the semantics of Φ . Then $l(a_1, \dots, a_n) \rightarrow R_l(a'_1, \dots, a'_n)[l \mapsto l'] = l'(a_1, \dots, a_n) \rightarrow R_l(a'_1, \dots, a'_n)$ is a clause of Φ_{rdt} that denote the equivalence between

the table of l' relation in Ψ'_R and the table of R_l relation in Ψ_R . Formally, for any database instance $D_{l'} \in D'$ that corresponds to nodes $N_l \in G$ and database instance $D_{R_l} \in D$, $\Phi_{\text{rdt}}(D_{l'}) = D_{R_l}$

- (2) For any edge $E(l, s, t, a_1, \dots, a_n)$ with label l that connects nodes s and t and has property values a_1, \dots, a_n , there exists a corresponding clause $s(\dots), l(s, t, a_1, \dots, a_n), t(\dots) \rightarrow s'(\dots), l'(a_1, \dots, a_n, s, t), t'(\dots)$ in Φ_{sdt} . Similarly, σ and Φ_{rdt} are updated by $\sigma \cup \{l \mapsto l'\}$ and $\Phi_{\text{rdt}} \cup \Phi[l \mapsto l']$, respectively. Let $s(\dots), l(s, t, a_1, \dots, a_n), t(\dots) \rightarrow R_s(\dots), R_l(a_1, \dots, a_n, s, t), R_t(\dots)$ be a formula of Φ by the semantics of Φ . Also, since s and t are two labels for nodes, then by the case 1 we have $s'(\dots) \rightarrow R_s(\dots)$ and $t'(\dots) \rightarrow R_t(\dots)$ from Φ_{rdt} (namely, $\Phi[s \mapsto s']$ and $\Phi[t \mapsto t']$). Then $\Phi[s \mapsto s', l \mapsto l', t \mapsto t'] = s'(\dots), l'(s, t, a_1, \dots, a_n), t'(\dots) \rightarrow R_s(\dots), R_l(a_1, \dots, a_n, s, t), R_t(\dots)$. Formally, $\Phi_{\text{rdt}}(D_{l'}) = D_{R_l}$ where $D_{l'} \in D'$ and $D_{R_l} \in D$.

Thus, we know $D' \triangleright \Psi'_R \wedge D \triangleright \Psi_R \wedge \Phi_{\text{rdt}}(D') = D \Rightarrow D' \sim_{\Phi_{\text{rdt}}} D$ where Φ_{rdt} is derived from Algorithm 2. □

LEMMA F.12. Let Ψ_G be a graph schema, $\Psi_{R'}$ be the induced relational schema. Given a graph query Q_G over Ψ_G and a SQL query $Q_{R'}$ over $\Psi_{R'}$ such that $\Phi_{\text{sdt}}, \Psi_R \vdash Q_G \xrightarrow{\text{query}} Q_{R'}$, let Φ be a database transformer from Ψ_G to Ψ_R , and Φ_{rdt} be the residual database transformer between Ψ'_R and Ψ_R . If $Q_G \simeq_{\Phi} Q_R$, it holds that $Q_{R'} \simeq_{\Phi_{\text{rdt}}} Q_R$.

PROOF. By the definition 4.5, if $Q_G \simeq_{\Phi} Q_R$ holds, i.e.,

$$\forall G, D. (G \triangleright \Psi_G \wedge D \triangleright \Psi_R \wedge G \sim_{\Phi} D) \Rightarrow \llbracket Q_G \rrbracket_G \equiv \llbracket Q_R \rrbracket_D$$

then $\Phi(G) = D \Leftrightarrow G \sim_{\Phi} D$ holds. Therefore, we can infer that $\llbracket Q_G \rrbracket_G \equiv \llbracket Q_R \rrbracket_D$ is true. Furthermore, by Lemma F.11, we know $\Phi_{\text{rdt}}(D') = D$ for any $D' \triangleright \Psi_{R'}$ where D' is a relational database instance over $\Psi_{R'}$. Also, by Lemma F.1, we know $Q_G \simeq_{\Phi_{\text{sdt}}} Q_{R'} \Rightarrow \llbracket Q_G \rrbracket_G \equiv \llbracket Q_{R'} \rrbracket_{D'}$ where $\Phi_{\text{sdt}}(G) = D'$. Then we can infer that $\llbracket Q_R \rrbracket_D \equiv \llbracket Q_G \rrbracket_G \equiv \llbracket Q_{R'} \rrbracket_{D'}$. Therefore, this Lemma is proved by that $Q_{R'} \simeq_{\Phi'} Q_R$ holds, i.e.,

$$\forall D', D. (D' \triangleright \Psi_{R'} \wedge D \triangleright \Psi_R \wedge D' \sim_{\Phi'} D) \Rightarrow \llbracket Q_{R'} \rrbracket_{D'} \equiv \llbracket Q_R \rrbracket_D$$
□

THEOREM F.13 (SOUNDNESS(5.9)). Let $\text{CheckSQL}(\Psi_R, Q, \Psi'_R, Q', \Phi_{\text{rdt}})$ be a sound procedure for equivalence checking of SQL queries Q, Q' over relational schemas Ψ_R, Ψ'_R connected by RDT Φ_{rdt} . Given a Cypher query Q_G over graph schema Ψ_G , a SQL query Q_R over relational schema Ψ_R , and their database transformer Φ , if $\text{CHECKEQUIVALENCE}(\Psi_G, Q_G, \Psi_R, Q_R, \Phi)$ returns \top , it holds that $Q_G \simeq_{\Phi} Q_R$.

PROOF. By Lemma F.1, it holds that $Q_G \simeq_{\Phi_{\text{sdt}}} Q_{R'}$, i.e.,

$$\forall G, D'. (G \triangleright \Psi_G \wedge D' \triangleright \Psi_{R'} \wedge G \sim_{\Phi_{\text{sdt}}} D') \Rightarrow \llbracket Q_G \rrbracket_G \equiv \llbracket Q_{R'} \rrbracket_{D'}$$

Since $G \sim_{\Phi_{\text{sdt}}} D'$ holds, we have $\llbracket Q_G \rrbracket_G \equiv \llbracket Q_{R'} \rrbracket_{D'}$. Further, if $\text{CHECKEQUIVALENCE}(\Psi_G, Q_G, \Psi_R, Q_R, \Phi)$ returns \top , then CheckSQL also returns \top , i.e., $\llbracket Q_{R'} \rrbracket_{D'} \equiv \llbracket Q_R \rrbracket_D$ where $\Phi_{\text{rdt}}(D') = D$ by the definition of Φ_{rdt} . Also, we have $\Phi(G) = D$ by the definition of Φ and $Q_G \simeq_{\Phi_{\text{sdt}}} Q_{R'} \Rightarrow \llbracket Q_G \rrbracket_G \equiv \llbracket Q_{R'} \rrbracket_{D'}$ by Lemma 5.7. Then we have $\llbracket Q_G \rrbracket_G \equiv \llbracket Q_{R'} \rrbracket_{D'} \equiv \llbracket Q_R \rrbracket_D$.

Apparently, $G \sim_{\Phi} D$ holds by our assumption. Thus, $Q_G \simeq_{\Phi} Q_R$ is proved, i.e.,

$$\forall G, D. (G \triangleright \Psi_G \wedge D \triangleright \Psi_R \wedge G \sim_{\Phi} D) \Rightarrow \llbracket Q_G \rrbracket_G \equiv \llbracket Q_R \rrbracket_D$$
□

THEOREM F.14 (COMPLETENESS(5.10)). *Let $\text{CheckSQL}(\Psi_R, Q, \Psi'_R, Q', \Phi_{\text{rdt}})$ be a complete procedure for equivalence checking of SQL queries Q, Q' over schemas Ψ_R, Ψ'_R connected by RDT Φ_{rdt} . Given a Cypher query Q_G over graph schema Ψ_G , a SQL query Q_R over relational schema Ψ_R , and their database transformer Φ , if $Q_G \simeq_\Phi Q_R$, then $\text{CHECKEQUIVALENCE}(\Psi_G, Q_G, \Psi_R, Q_R, \Phi)$ returns \top .*

PROOF. If $Q_G \simeq_\Phi Q_R$ holds, then $Q_{R'} \simeq_{\Phi_{\text{rdt}}} Q_R$ holds by Lemma F.12. Further, we know that if $Q_{R'} \simeq_{\Phi_{\text{rdt}}} Q_R$, then $\text{CheckSQL}(\Psi_{R'}, Q_{R'}, \Psi_R, Q, \Phi')$ returns \top . Therefore, the procedure $\text{CHECKSQLEQUIVALENCE}(\Psi_G, Q_G, \Psi_R, Q_R, \Phi)$ returns \top . □