# Explaining Results of Path Queries on Graphs
## Single-Path Results for Context-Free Path Queries

Jelle Hellings

Department of Computer Science, University of California, Davis
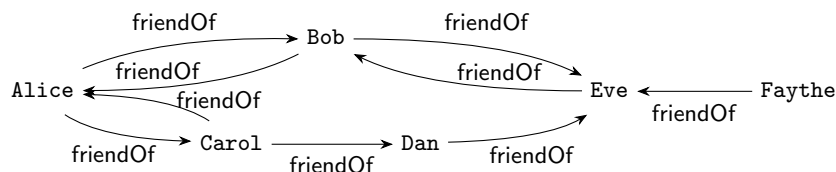Davis, CA 95616-8562, USA
jhellings@ucdavis.edu

**Abstract.** Many graph query languages use, at their core, path queries that yield node pairs that are connected by a path of interest. For the end-user, such node pairs only give limited insight as to *why* this query result is obtained, as the pair does not directly identify the underlying path of interest. To address this limitation of path queries, we propose the *single-path semantics*, which evaluates path queries to, for each node pair $(m, n)$, a single path from $m$ to $n$ satisfying the conditions of the query. To put our proposal in practice, we provide an efficient algorithm for evaluating *context-free path queries*, a particular powerful type of path queries, using the single-path semantics. Additionally, we perform a short evaluation of our techniques that shows that the single-path semantics is practically feasible, even when query results grow large.

**Keywords:** Graph Queries · Path Results · Context-Free Path Queries

## 1 Introduction

The graph data model is one of the most versatile and natural data models in use: graph-structured data is everywhere and examples can be found in family trees, social networks, process models, gene networks, XML data, and RDF data [1,2,9,12,30]. As an example, consider the small social network visualized in Figure 1 in which nodes represent peoples and edges represent the relationships between people.

A central step in the analysis of such graph data is the ability to *query* the data for relationships of interest. For this purpose, many different query languages have been developed, including XPATH for querying XML data [8,9,11],



**Fig. 1.** A typical example of graph data: a social network relating peoples.

SPARQL for querying RDF data [19,30], the graph query languages GX-
PATH [24], CYPHER [28], and GREMLIN [29], and formal verification languages
such as PDL, KAT, CTL, and LTL [12,22,23]. At their core, these graph query
languages depend on *path queries* that can be used to express *indirect relation-
ships* that can be derived from the data [2]. Examples of such path queries are the
well-known regular path queries [5] and the context-free path queries [18,20,23,31].
Unfortunately, path queries are typically evaluated to only a set of node pairs
$(m, n)$ that are connected by a path of interest, which gives little insight in the
way pairs $(m, n)$ are obtained, limiting their capabilities for graph analytics.

*Example 1.* Let $\mathfrak{G}$ be the social network visualized in Figure 1. The path query
indirectFriendOf = friendOf$^+$, expressed by a regular expression, will return the
derived relationship between pairs $(m, n)$ such that $m$ is a friend of $n$, or a
friend-of-a-friend of $n$, or a friend-of-a-friend-of-a-friend of $n$, and so on. The
pair (Alice, Eve) is in the result of evaluating indirectFriendOf on graph $\mathfrak{G}$.
Unfortunately, Alice cannot use this result to determine whom of her friends can
help her to get in contact with Eve, and Alice will have to further analyze the
underlying graph data.

Example 1 illustrates the need to answer path queries with the underlying
paths of interest inspected by these queries. To address this need, we propose the
*single-path semantics* for evaluating path queries: using the single-path semantics,
a path query will evaluate to a shortest path connecting node pair $(m, n)$ (for
each node pair in the query result).

*Example 2.* Consider the setting of Example 1. Evaluation of indirectFriendOf us-
ing the single path semantics can result in the path "Alice friendOf Bob friendOf
Eve", from which Alice can derive a way to contact Eve. As the single path
semantics requires a single and shortest path between node pairs, the path
"Alice friendOf Carol friendOf Dan friendOf Eve" cannot be in the output.

The need for the single-path semantics extends beyond the above toy example.
Not only can single-path semantics provide more relevant information to end-users,
the single-path semantics can also aid in graph analytics and data exploration,
and can be used to provide *data provenance* for traditional path queries [10],
this by providing paths that show why a path query includes a certain node pair
in its output. Furthermore, in the large-scale graph data setting in which many
complex path queries are evaluated, there is a need for tools to support query
debugging [21], for which the single-path semantics can also be of use.

In this paper, we deal with the issues outlined by proposing the single-path
semantics. We focus our study on the context-free path queries, as these are
a particular powerful type of path queries that cannot only express all typical
path queries (e.g. [5,20]), but also have applications in model checking [23],
bio-informatics [31], and parser construction [17]. In specific, we formalize the
*single-path semantics*, introduce the algorithm MINIMIZESETGG that provides
efficient evaluation of context-free path queries on graphs using the single-path
semantics, and evaluate the performance of MINIMIZESETGG in practice. Our

results show promise, as MINIMIZESETGG can easily answer queries whose results contains tens-of-millions of paths, even if these paths have considerable lengths.

## 2    Preliminaries

First, we introduce the terminology and notation used throughout this paper.

Let $\Sigma$ be a set of symbols. We call a sequence $s = \sigma_1 \ldots \sigma_n$ of symbols, $\sigma_1, \ldots, \sigma_n \in \Sigma$, a string over $\Sigma$. We write $|s| = n$ to denote the length of $s$. The empty string is denoted by $\epsilon$ and we usually treat individual symbols as strings of length one. The concatenation of two strings $s_1$ and $s_2$ is denoted by $s_1 \circ s_2$. We denote the set of all strings over $\Sigma$ by $\Sigma^*$. A *language* over $\Sigma$ is a (possibly infinite) set of strings over $\Sigma$.

A *graph* is a triple $\mathfrak{G} = (\mathcal{V}, \Sigma, \delta)$, in which $\mathcal{V}$ is a finite set of nodes, $\Sigma$ is a finite set of alphabet symbols used as edge labels, and $\delta \subseteq \mathcal{V} \times \Sigma \times \mathcal{V}$ is a finite set of labeled edges. To simplify presentation, we assume that $\mathcal{V}$ and $\Sigma$ do not overlap ($\mathcal{V} \cap \Sigma = \emptyset$). A *path* in $\mathfrak{G}$ is a sequence $\pi = m_1\ \sigma_1\ m_2 \ldots m_{n-1}\ \sigma_{n-1}\ m_n$ such that, for every $m_i\ \sigma_i\ m_{i+1}$ in the sequence, $1 \leq i < n$, we have $(m_i, \sigma_i, m_{i+1}) \in \delta$. We write $m_1 \pi m_n$ to indicate that $\pi$ is a path starting at node $m_1$ and ending at node $m_n$. We write $|\pi| = n - 1$ to denote the length of $\pi$ and we write $\mathrm{trace}(\pi) = \sigma_1 \ldots \sigma_{n-1}$ to denote the trace of $\pi$, the string represented by the sequence of edge labels in $\pi$.

A *path query* $q$ is specified by a language $\mathcal{L}$ that contains all traces of paths of interest, e.g., via a regular expression (RPQs) or via a context-free grammar (CFPQs). The evaluation of $q$ on graph $\mathfrak{G}$ using the standard *relational semantics* simply consists of all node pairs that are connected by paths whose trace is in $\mathcal{L}$. To denote the evaluation of $q$ on $\mathfrak{G}$, we write $[\![q]\!]_{\mathfrak{G}}$, and we have $[\![q]\!]_{\mathfrak{G}} = \{(m, n) \mid \exists$ path $m\pi n$ in $\mathfrak{G}$ with $\mathrm{trace}(\pi) \in \mathcal{L}\}$.

*Example 3.* In Example 1, the query indirectFriendOf was expressed by the regular expression friendOf$^+$. This regular expression represents the language

$$\mathcal{L} = \{\mathsf{friendOf}, \mathsf{friendOf} \circ \mathsf{friendOf}, \mathsf{friendOf} \circ \mathsf{friendOf} \circ \mathsf{friendOf}, \ldots\}.$$

We have $(\texttt{Alice}, \texttt{Eve}) \in [\![\mathsf{indirectFriendOf}]\!]_{\mathfrak{G}}$, with $\mathfrak{G}$ the graph in Figure 1, as there exists a path $\pi = $ "$\texttt{Alice}$ friendOf $\texttt{Bob}$ friendOf $\texttt{Eve}$" with $\mathrm{trace}(\pi) \in \mathcal{L}$.

A *grammar* is a triple $\mathscr{C} = (\mathcal{N}, \Sigma, \mathcal{P})$, in which $\mathcal{N}$ is a set of non-terminals, $\Sigma$ is a finite set of alphabet symbols, and $\mathcal{P}$ is a set of production rules. We require that $\mathcal{N}$ and $\Sigma$ do not overlap ($\mathcal{N} \cap \Sigma = \emptyset$). The set of production rules, $\mathcal{P}$, consists of production rules of the form $\textsc{a} \mapsto \textsc{b}\ \textsc{c}$ or $\textsc{a} \mapsto \sigma$, in which $\textsc{a}, \textsc{b}, \textsc{c} \in \mathcal{N}$ and $\sigma \in \Sigma$.

Each non-terminal in $\mathcal{N}$ represents a language over $\Sigma$: the production rules in $\mathcal{P}$ describe how to produce strings out of non-terminals via rewrite steps. To illustrate this, consider a string $s = s_1 \circ \textsc{a} \circ s_2$ in which $s_1, s_2 \in (\mathcal{N} \cup \Sigma)^*$ and $\textsc{a} \in \mathcal{N}$. If there exists a production rule $(\textsc{a} \mapsto s') \in \mathcal{P}$, then we can rewrite $s$ into $s_1 \circ s' \circ s_2$ by applying the rewrite $\textsc{a} \mapsto s'$. We write $s \rightarrow_{\mathcal{P}}^* s'$ if $s$ can be

rewritten into $s'$ using production rules in $\mathcal{P}$, and we write $s \to_{\mathcal{P}}^+ s'$ if $s \to_{\mathcal{P}}^* s'$ and at least one rewrite step is necessary to rewrite $s$ into $s'$.

The *language* of non-terminal $\text{A} \in \mathcal{N}$ is defined by $\mathcal{L}(\mathscr{C}; \text{A}) = \{s \in \Sigma^* \mid \text{A} \to_{\mathcal{P}}^* s\}$. Given a grammar with non-terminal $\text{A}$, we simply write $\text{A}$ to denote the path query based on the language $\mathcal{L}(\mathscr{C}; \text{A})$.

*Example 4.* Consider the grammar $\mathscr{C} = (\mathcal{N}, \Sigma, \mathcal{P})$ in which $\mathcal{N} = \{\text{A}\}$, $\Sigma = \{\textsf{friendOf}\}$, and $\mathcal{P} = \{\text{A} \to \textsf{friendOf}, \text{A} \to \text{A A}\}$. The language $\mathcal{L}(\mathscr{C}; \text{A})$ is equivalent to the language $\mathcal{L}$ of Example 3. Hence, we have $(\texttt{Alice}, \texttt{Eve}) \in [\![\text{A}]\!]_{\mathfrak{G}}$, in which $\mathfrak{G}$ is the graph visualized in Figure 1.

## 3   The single-path semantics

in Section 2, we already introduced the typical *relational semantics* of path queries. Unfortunately, the step toward *path-based semantics*—in which a path query yields paths $m\pi n$ instead of node pairs $(m,n)$—is not straightforward. Even in basic situations, the resulting set of paths can already be unbounded in size, making it impossible to simply evaluate to such a set:

*Example 5.* Consider Example 1 and the graph $\mathfrak{G}$ visualized in Figure 1. This graph is cyclic, as there is a path "$\texttt{Alice}$ $\textsf{friendOf}$ $\texttt{Carol}$ $\textsf{friendOf}$ $\texttt{Dan}$ $\textsf{friendOf}$ $\texttt{Eve}$ $\textsf{friendOf}$ $\texttt{Bob}$ $\textsf{friendOf}$ $\texttt{Alice}$". Hence, we can make paths of arbitrary lengths that match the query $\texttt{indirectFriendOf}$, and the set of all paths matching the query is unbounded in size.

Restricting the paths considered in the evaluation, e.g., to simple paths, assures that the set of paths considered is finite. Unfortunately, changing the paths considered during evaluation defeats the purpose of path-based semantics as a data provenance and debugging tool for normal path queries. Moreover, it is well-known that such restrictions make query evaluation prohibitive expensive [2,3,7,25]. Restricting the number of paths in the result, e.g., to a single path per node pair, will also assure a finite result. Unfortunately, as Example 5 already shows, individual paths in such a finite result set can still have a practically unbounded length. To address these issues, we choose to return a single as-short-as-possible path for each node-pair $(m,n)$:

**Definition 1.** *Let $q$ be a path query specified by language $\mathcal{L}$ and let $\mathfrak{G}$ be a graph. The evaluation of $q$ on $\mathfrak{G}$ using the* single-path semantics, *denoted by* $\text{single}(q|_{\mathfrak{G}})$, *yields, for every $(m,n) \in [\![q]\!]_{\mathfrak{G}}$, a single shortest path $m\pi n$ in $\mathfrak{G}$ such that $\text{trace}(\pi) \in \mathcal{L}$. (Hence, for every other path $m\pi'n$ in $\mathfrak{G}$ with $\text{trace}(\pi') \in \mathcal{L}$, we have $|\pi| \le |\pi'|$.)*

Toward evaluating context-free path queries using the single-path semantics, we proceed in three steps. First, in Section 3.1, we show that all paths of interest of a context-free path query can be represented by a grammar. Then, in Section 3.2, we propose MINIMIZESET, an algorithm for computing a shortest string in a context-free language. Finally, in Section 3.3, we combine these results and show how to evaluating context-free path queries using the single-path semantics.

### 3.1   Representing the paths of interest of a path query

Let $\mathfrak{G} = (\mathcal{V}, \Sigma, \delta)$ be a graph and $(m, n) \in \mathcal{V}$ a pair of nodes. There is a close correspondence between labeled graphs and finite automata and we can easily interpret $(\mathfrak{G}, m, n)$ as a finite automata with initial state $m$ and final state $n$. The language of this finite automata is $\mathcal{L}(\mathfrak{G}; m, n) = \{\mathrm{trace}(\pi) \mid m\pi n \text{ is a path in } \mathfrak{G}\}$. It is well-known that the intersection of a finite automaton and a grammar can be represented by another context-free grammar:

**Lemma 1 (Bar-Hillel et al. [4]).** *Let* $\mathscr{C} = (\mathcal{N}, \Sigma, \mathcal{P})$ *be a grammar, let* $\mathfrak{G} = (\mathcal{V}, \Sigma, \delta)$ *be a graph, let* $\mathrm{A} \in \mathcal{N}$, *and let* $m, n \in \mathcal{V}$. *The language* $\mathcal{L}(\mathscr{C}; \mathrm{A}) \cap \mathcal{L}(\mathfrak{G}; m, n)$ *can be represented by a grammar.*

Lemma 1 guarantees that there is a finite representation of the set of all strings in $\mathcal{L}(\mathscr{C}; \mathrm{A}) \cap \mathcal{L}(\mathfrak{G}; m, n)$, each such string representing the trace of a path $m\pi n$ in $\mathfrak{G}$ with $\mathrm{trace}(\pi) \in \mathcal{L}(\mathscr{C}; \mathrm{A})$. Unfortunately, there can be several paths with the same trace, complicating the derivation of the underlying paths. To improve on this, we show the existence of *graph-annotated grammars* that directly represent the set of paths instead of their traces:

**Definition 2.** *Let* $\mathscr{C} = (\mathcal{N}, \Sigma, \mathcal{P})$ *be a grammar and let* $\mathfrak{G} = (\mathcal{V}, \Sigma, \delta)$ *be a graph. We denote triples* $(\mathrm{A}, m, n) \in \mathcal{N} \times \mathcal{V}^2$ *by* $\mathrm{A}|_{mn}$. *An annotated grammar over* $(\mathscr{C}, \mathfrak{G})$ *is a grammar* $\mathscr{C}|_{\mathfrak{G}} = (\mathcal{N}|_{\mathfrak{G}}, \Sigma, \mathcal{P}|_{\mathfrak{G}})$ *in which*

1. $\mathcal{N}|_{\mathfrak{G}} = \{\mathrm{A}|_{mn} \in \mathcal{N} \times \mathcal{V}^2 \mid \mathcal{L}(\mathscr{C}; \mathrm{A}) \cap \mathcal{L}(\mathfrak{G}; m, n) \neq \emptyset\}$;
2. $\mathcal{P}|_{\mathfrak{G}} = P_\Sigma \cup P_\mathcal{N}$ *with* $P_\Sigma = \{\mathrm{A}|_{mn} \mapsto \sigma \mid (m, \sigma, n) \in \delta \wedge (\mathrm{A} \mapsto \sigma) \in \mathcal{P}\}$ *and* $P_\mathcal{N} = \{\mathrm{A}|_{mn} \mapsto \mathrm{B}|_{mo} \ \mathrm{C}|_{on} \mid (\mathrm{A} \mapsto \mathrm{B} \ \mathrm{C}) \in \mathcal{P}\}$.

The notation $\mathrm{A}|_{mn}$ denotes a node-annotated non-terminal: any string produced from rewriting this non-terminal is a trace of a path $m\pi n$. As rewriting $\mathrm{A}|_{mn}$ eventually leads to rewrite steps using production rules in $P_\Sigma$, which represent single edges in $\mathfrak{G}$, the path $\pi$ can be derived by keeping track of these node-annotations. Notice that $|\mathcal{N}|_{\mathfrak{G}}| \leq |\mathcal{N}||\mathcal{V}|^2$, $|P_\Sigma| \leq |\mathcal{P}||\delta|$, and $|P_\mathcal{N}| \leq |\mathcal{P}||\mathcal{V}|^3$. We illustrate these annotated grammars with an example:

*Example 6.* Let $\mathfrak{G}$ be the graph visualized in Figure 1 and $\mathscr{C}$ the grammar of Example 4. We construct the annotated grammar $\mathscr{C}|_{\mathfrak{G}} = (\mathcal{N}|_{\mathfrak{G}}, \Sigma, \mathcal{P}|_{\mathfrak{G}})$. For brevity, we refer to each person by the first letter of their name. We have

$$\mathcal{N}|_{\mathfrak{G}} = \{\mathrm{Q}|_{mn} \mid m, n \in \{\mathrm{A}, \mathrm{B}, \mathrm{C}, \mathrm{D}, \mathrm{E}\}\} \cup \{\mathrm{Q}|_{\mathrm{F}n} \mid n \in \{\mathrm{A}, \mathrm{B}, \mathrm{C}, \mathrm{D}, \mathrm{E}\}\}.$$

We have $\mathcal{P}|_{\mathfrak{G}} = P_\Sigma \cup P_\mathcal{N}$, in which $P_\Sigma$ represents all edges in $\mathfrak{G}$ and $P_\mathcal{N}$ represents all ways in which paths in $\mathfrak{G}$ can be combined. E.g., $\mathrm{Q}|_{\mathrm{AB}} \in P_\Sigma$, as $(\mathtt{Alice}, \mathsf{friendOf}, \mathtt{Bob})$ is and edge in $\mathfrak{G}$, and $(\mathrm{Q}|_{\mathrm{AB}} \mapsto \mathrm{Q}|_{\mathrm{AD}} \ \mathrm{Q}|_{\mathrm{DB}}) \in P_\mathcal{N}$, as there is a $\mathsf{friendOf}$-labeled path from $\mathtt{Alice}$ to $\mathtt{Dan}$ and another $\mathsf{friendOf}$-labeled path from $\mathtt{Dan}$ to $\mathtt{Bob}$. To produce a path from Alice to Eve, we use $\mathscr{C}|_{\mathfrak{G}}$:

$$\mathrm{Q}|_{\mathtt{AliceEve}} \rightarrow^*_{\mathcal{P}|_{\mathfrak{G}}} \{\text{Rewrite } \mathrm{Q}|_{\mathtt{AliceEve}} \mapsto \mathrm{Q}|_{\mathtt{AliceCarol}} \ \mathrm{Q}|_{\mathtt{CarolEve}}\}$$

$$\mathrm{Q}|_{\mathtt{AliceCarol}} \ \mathrm{Q}|_{\mathtt{CarolEve}} \rightarrow^*_{\mathcal{P}|_{\mathfrak{G}}} \{\text{Rewrite } \mathrm{Q}|_{\mathtt{CarolEve}} \mapsto \mathrm{Q}|_{\mathtt{CarolDan}} \ \mathrm{Q}|_{\mathtt{DanEve}}\}$$

$$\mathrm{Q}|_{\mathtt{AliceCarol}} \ \mathrm{Q}|_{\mathtt{CarolDan}} \ \mathrm{Q}|_{\mathtt{DanEve}} \rightarrow^*_{\mathcal{P}|_{\mathfrak{G}}} \{\text{Rewrite } \mathrm{Q}|_{\mathtt{AliceCarol}} \mapsto \mathsf{friendOf}, \dots\}$$

$$\mathsf{friendOf} \circ \mathsf{friendOf} \circ \mathsf{friendOf}.$$

The annotation in each node-annotated non-terminal carry information that can be used to map strings in $\mathscr{C}|_{\mathfrak{G}}$ to paths in the underlying graph $\mathfrak{G}$. E.g., in this rewrite, we derived a path from `Alice` to `Eve` in graph $\mathfrak{G}$ of length three, namely the path "`Alice friendOf Carol friendOf Dan friendOf Eve`".

Using induction, we can prove that graph-annotated grammars can always be used as illustrated in Example 6:

**Proposition 1.** *Let $\mathscr{C} = (\mathcal{N}, \Sigma, \mathcal{P})$ be a grammar, let $\mathfrak{G} = (\mathcal{V}, \Sigma, \delta)$ be a graph, let $\mathscr{C}|_{\mathfrak{G}} = (\mathcal{N}|_{\mathfrak{G}}, \Sigma, \mathcal{P}|_{\mathfrak{G}})$ be the annotated grammar over $(\mathscr{C}, \mathfrak{G})$, let $m\pi n$ be a path in $\mathfrak{G}$, and let $\mathrm{A} \in \mathcal{N}$ be a non-terminal. We have $\mathrm{trace}(\pi) \in \mathcal{L}(\mathscr{C}; \mathrm{A})$ if and only if we can derive $\pi$ from $\mathrm{A}|_{mn} \in \mathcal{N}|_{\mathfrak{G}}$.*

We note that annotated grammars can, on their own, be used in *interactive* data exploration tools in which users can explorer the query results by zooming in on certain paths in the dataset, e.g., for graph analysis and query debugging.

### 3.2   Deriving shortest strings of a grammar

Next, we propose an efficient way to compute a shortest string in the language defined by a grammar. Mclean et al. [27] already proved that a shortest string can be computed effective given a grammar, but did not provide a practical algorithm for computing shortest strings. Toward such an algorithm, we introduce rewrites using *simple* production rules:

**Definition 3.** *Let $\mathcal{P}$ be a set of production rules. We define $\mathrm{heads}(\mathcal{P}) = \{\mathrm{A} \mid (\mathrm{A} \mapsto s) \in \mathcal{P}\}$ and we define the set of non-terminals derivable from $\mathrm{A}$ using the production rules in $\mathcal{P}$ by $\langle \mathrm{A} \rangle_{\mathcal{P}} = \{\mathrm{B} \in \mathcal{N} \mid \exists s_1 \exists s_2 \; \mathrm{A} \rightarrow^+_{\mathcal{P}} s_1 \circ \mathrm{B} \circ s_2\}$.*

*A set of production rules $\mathcal{P}$ is* non-recursive *if, for every $\mathrm{A} \in \mathrm{heads}(\mathcal{P})$, we have $\mathrm{A} \notin \langle \mathrm{A} \rangle_{\mathcal{P}}$. A set of production rules $\mathcal{P}$ is* deterministic *if, for every $\mathrm{A} \in \mathrm{heads}(\mathcal{P})$, there exists exactly one production rule $(\mathrm{A} \mapsto s) \in \mathcal{P}$. Finally, a set of production rules $\mathcal{P}$ is* effective *if $\mathrm{A} \in \mathrm{heads}(\mathcal{P})$ implies that there exists a string $s \in \Sigma^*$ such that $\mathrm{A} \rightarrow^*_{\mathcal{P}} s$. We refer to a set of production rules that is non-recursive, deterministic, and effective as* simple.

A simple set of production rules $\mathcal{P}$ over alphabet $\Sigma$ can be used to rewrite every non-terminal $\mathrm{A} \in \mathrm{heads}(\mathcal{P})$ into a unique string $\mathrm{ustring}_{\mathcal{P}}(\mathrm{A})$ over $\Sigma$ in a straightforward manner, as $\mathcal{P}$ does not provide any choices during such a rewrite.

*Example 7.* Consider Example 6. For brevity, we restrict ourselves to Alice, Carol, Dan, and Eve. With respect to these four people, the following set of production rules in the annotated grammar is deterministic non-recursive:

$$\mathrm{Q}|_{\mathrm{AC}} \mapsto \mathsf{friendOf}, \quad \mathrm{Q}|_{\mathrm{CA}} \mapsto \mathsf{friendOf}, \quad \mathrm{Q}|_{\mathrm{CD}} \mapsto \mathsf{friendOf}, \quad \mathrm{Q}|_{\mathrm{DE}} \mapsto \mathsf{friendOf},$$

$$\mathrm{Q}|_{\mathrm{AD}} \mapsto \mathrm{Q}|_{\mathrm{AC}} \; \mathrm{Q}|_{\mathrm{CD}}, \quad \mathrm{Q}|_{\mathrm{AE}} \mapsto \mathrm{Q}|_{\mathrm{AD}} \; \mathrm{Q}|_{\mathrm{DE}}, \quad \mathrm{Q}|_{\mathrm{CE}} \mapsto \mathrm{Q}|_{\mathrm{CA}} \; \mathrm{Q}|_{\mathrm{AE}}.$$

We can use simple production rules derived from a grammar to represent shortest strings in that grammar:

**Lemma 2.** *Let $\mathscr{C} = (\mathcal{N}, \Sigma, \mathcal{P})$ be a grammar. There exists a simple set of production rules $\mathcal{P}' \subseteq \mathcal{P}$ such that, for every non-terminal $\mathrm{A} \in \mathcal{N}$ with $\mathcal{L}(\mathscr{C}; \mathrm{A}) \neq \emptyset$, ustring$_{\mathcal{P}'}(\mathrm{A})$ is a shortest string in $\mathcal{L}(\mathscr{C}; \mathrm{A})$.*

We say that a set of production rules that satisfies the conditions of Lemma 2 is *minimizing*. Unfortunately, not every simple set of production rules is minimizing:

*Example 8.* Consider Example 7. The provided simple set of production rules $\mathcal{P}'$ is not minimizing: we have $|\text{ustring}_{\mathcal{P}'}(\mathrm{Q}|_{\mathrm{CE}})| = 4$, while a shorter string of length two exists. By replacing the production rule for $\mathrm{Q}|_{\mathrm{CE}}$ in $\mathcal{P}'$ by $\mathrm{Q}|_{\mathrm{CE}} \mapsto \mathrm{Q}|_{\mathrm{CD}} \; \mathrm{Q}|_{\mathrm{DE}}$, we obtain a minimizing set of production rules.

Using a minimizing set of production rules, it is straightforward to produce shortest strings for $\mathrm{A} \in \text{heads}(\mathcal{P})$. Moreover, the way to obtain these shortest strings, by rewriting $\mathrm{A}$, also provides complete information on how these shortest strings can be obtained from the original grammar. Next, we propose the MINIMIZESET algorithm to construct a minimizing set of production rules. The pseudo-code of this algorithm can be found in Figure 2, *left*.

The MINIMIZESET algorithm works rather intuitively. Let $\mathscr{C} = (\mathcal{N}, \Sigma, \mathcal{P})$ be a grammar. Production rules of the form $(\mathrm{A} \to \sigma) \in \mathcal{P}$, $\sigma \in \Sigma$, produce the shortest possible strings: if $(\mathrm{A} \to \sigma) \in \mathcal{P}$, then $\sigma$ is a shortest string in $\mathcal{L}(\mathscr{C}; \mathrm{A})$. If such productions rules exist for $\mathrm{A}$, then we choose one of them for the minimizing set of production rules (Line 3). Next, we process non-terminals $\mathrm{A}$ for which we have determined the length $cost(\mathrm{A})$ of the shortest strings in $\mathcal{L}(\mathscr{C}; \mathrm{A})$. We do so on increasing string length by using a min-priority queue *new* (Line 7). We process $\mathrm{A}$ by checking, for each production rule $(\mathrm{C} \mapsto \mathrm{A} \; \mathrm{B}) \in \mathcal{P}$ or $(\mathrm{C} \mapsto \mathrm{B} \; \mathrm{A}) \in \mathcal{P}$, whether using this production rule will allow us to rewrite $\mathrm{C}$ into a shorter string than the currently-found string with length $cost[\mathrm{C}]$ (Line 10 and Line 12). We do so by rewriting—in this production rule—$\mathrm{A}$ to a string of length $cost(\mathrm{A})$ (Line 14).

**Theorem 1.** *Let $\mathscr{C} = (\mathcal{N}, \Sigma, \mathcal{P})$ be grammar. Execution of MINIMIZESET$(\mathscr{C})$ yields a minimizing set of production rules $\mathcal{P}'$ for $\mathscr{C}$ in $\mathcal{O}(|\mathcal{N}|(|\mathcal{N}| \log|\mathcal{N}| + |\mathcal{P}|))$. Using $\mathcal{P}'$, a set $R$ of shortest strings $s_{\mathrm{A}}$ in $\mathcal{L}(\mathscr{C}; \mathrm{A})$, $\mathrm{A} \in \mathcal{N}$, can be constructed in $\mathcal{O}(L)$, in which $L = \sum\{|s_{\mathrm{A}}| \mid s_{\mathrm{A}} \in R\}$ is the total length of these shortest strings.*

*Proof (sketch).* The main *while*-loop maintains the following invariants:

1. The set $\{\mathcal{P}'[\mathrm{A}] \mid \mathrm{A} \in \mathcal{P}'\}$ is simple.
2. If $\mathrm{A} \in \mathcal{P}'$ and $\mathcal{P}'[\mathrm{A}] = (\mathrm{A} \mapsto \mathrm{B} \; \mathrm{C})$, then $cost[\mathrm{A}] \geq cost[\mathrm{B}] + cost[\mathrm{C}]$, $cost[\mathrm{A}] > cost[\mathrm{B}]$, and $cost[\mathrm{A}] > cost[\mathrm{C}]$.
3. If $\mathrm{A} \in \mathcal{P}'$ and $s$ is a shortest string in $\mathcal{L}(\mathscr{C}; \mathrm{A})$, then $|s| \leq cost[\mathrm{A}]$.
4. Let $m$ be the priority of the last element removed from *new*. No new element is inserted in *new* with priority less than or equal to $m$.
5. Let $m$ be the priority of the last element removed from *new*. For every $\mathrm{A} \in \mathcal{N}$ and every shortest string $s$ in $\mathcal{L}(\mathscr{C}; \mathrm{A})$ with $|s| \leq m$, we have $cost[\mathrm{A}] = |s|$.

As each non-terminal is added to *new* at most once, the MINIMIZESET algorithm terminates. At termination, Invariants 1–5 guarantee that the resulting set of production rules is minimizing.

**Algorithm** MINIMIZESET($\mathscr{C} = (\mathcal{N}, \Sigma, \mathcal{P})$):

1:  $\mathcal{P}', cost :=$ empty mapping, empty mapping.
2:  $new$ is a min-priority queue.
3:  **for all** $(\text{A} \mapsto \sigma) \in \mathcal{P}$ **do**
4:      **if** $\text{A} \notin cost$ **then**
5:          $cost[\text{A}], \mathcal{P}'[\text{A}] := 1, (\text{A} \mapsto \sigma)$.
6:          add A to $new$ with priority 1.
7:  **while** $new \neq \emptyset$ **do**
8:      Take A with minimum priority in $new$.
9:      Remove A from $new$.
10:     **for all** $(\text{C} \mapsto \text{A B}) \in \mathcal{P}$ with $\text{B} \in cost$ **do**
11:         PRODUCE($\text{C} \mapsto \text{A B}$).
12:     **for all** $(\text{C} \mapsto \text{B A}) \in \mathcal{P}$ with $\text{B} \in cost$ **do**
13:         PRODUCE($\text{C} \mapsto \text{B A}$).
14: **return** $\{\mathcal{P}'[\text{A}] \mid \text{A} \in \mathcal{P}'\}$.

**Procedure** PRODUCE($\text{D} \mapsto \text{E F}$):

15: **if** $\text{D} \notin cost$ **then**
16:     $cost[\text{D}] := cost[\text{E}] + cost[\text{F}]$.
17:     $\mathcal{P}'[\text{D}] := \text{D} \mapsto \text{E F}$.
18:     Add D to $new$ with priority $cost[\text{E}] + cost[\text{F}]$.
19: **else if** $cost[\text{D}] > cost[\text{E}] + cost[\text{F}]$ **then**
20:     $cost[\text{D}] := cost[\text{E}] + cost[\text{F}]$.
21:     $\mathcal{P}'[\text{D}] := \text{D} \mapsto \text{E F}$.
22:     Lower priority of $\text{D} \in new$ to $cost[\text{E}] + cost[\text{F}]$.

**Algorithm** MINIMIZESETGG($\mathscr{C}, \mathfrak{G}$):

1:  $\mathcal{P}', cost :=$ empty mapping, empty mapping.
2:  $new$ is a min-priority queue.
3:  **for all** $(\text{A} \mapsto \sigma) \in \mathcal{P}$ and $(m, \sigma, n) \in \delta$ **do**
4:      **if** $\text{A}|_{mn} \notin cost$ **then**
5:          $cost[\text{A}|_{mn}], \mathcal{P}'[\text{A}|_{mn}] := 1, (\text{A}|_{mn} \mapsto \sigma)$.
6:          Add $\text{A}|_{mn}$ to $new$ with priority 1.
7:  **while** $new \neq \emptyset$ **do**
8:      Take $\text{A}|_{mn}$ with minimum priority in $new$.
9:      Remove $\text{A}|_{mn}$ from $new$.
10:     **for all** $(\text{C} \mapsto \text{A B}) \in \mathcal{P}$ with $\text{B}|_{no} \in cost$ **do**
11:         PRODUCEGG($\text{C}|_{mo} \mapsto \text{A}|_{mn} \text{B}|_{no}$).
12:     **for all** $(\text{C} \mapsto \text{B A}) \in \mathcal{P}$ with $\text{B}|_{om} \in cost$ **do**
13:         PRODUCEGG($\text{C}|_{on} \mapsto \text{B}|_{om} \text{A}|_{mn}$).
14: **return** $\{\mathcal{P}'[\text{A}|_{mn}] \mid \text{A}|_{mn} \in \mathcal{P}'\}$.

**Procedure** PRODUCEGG($\text{D}|_{uw} \mapsto \text{E}|_{uv} \text{F}|_{vw}$):

15: **if** $\text{D}|_{uw} \notin cost$ **then**
16:     $cost[\text{D}|_{uw}] := cost[\text{E}|_{uv}] + cost[\text{F}|_{vw}]$.
17:     $\mathcal{P}'[\text{D}|_{uw}] := \text{D}|_{uw} \mapsto \text{E}|_{uv} \text{F}|_{vw}$.
18:     Add $\text{D}|_{uw}$ to $new$ with priority $cost[\text{E}|_{uv}] + cost[\text{F}|_{vw}]$.
19: **else if** $cost[\text{D}|_{uw}] > cost[\text{E}|_{uv}] + cost[\text{F}|_{vw}]$ **then**
20:     $cost[\text{D}|_{uw}] := cost[\text{E}|_{uv}] + cost[\text{F}|_{vw}]$.
21:     $\mathcal{P}'[\text{D}|_{uw}] := \text{D}|_{uw} \mapsto \text{E}|_{uv} \text{F}|_{vw}$
22:     Lower priority of $\text{D}|_{uw} \in new$ to $cost[\text{E}|_{uv}] + cost[\text{F}|_{vw}]$.

**Fig. 2.** On the *left*, the MINIMIZESET algorithm that constructs a minimizing set of production rules for the grammar $\mathscr{C}$. On the *right*, the MINIMIZESETGG algorithm that constructs a minimizing set of production rules for the annotated grammar $\mathscr{C}|_{\mathfrak{G}}$, of which only the necessary parts are implicitly constructed.

To obtain the stated complexity, we represent *costs* as an array holding $|\mathcal{N}|$ integers. The costs used in *cost* and *new* are integers in the range $1, \ldots, 2^{|\mathcal{N}|-1}$, which we can represent using $\log(2^{|\mathcal{N}|}) = |\mathcal{N}|$ bits. The initialization steps perform $\mathcal{O}(|\mathcal{P}|)$ steps. The *while*-loop will, in the worst case, visit every non-terminal once. For each of these non-terminals, one insertion into and one removal from the priority queue *new* is performed. The inner *for*-loops will visit every production rule twice, causing at most $2|\mathcal{P}|$ decrease key operations on priority queue *new*. When using a Fibonacci heap for a priority queue holding at most $e$ elements, each insert and removal costs $\mathcal{O}(\log e)$ and each decrease key operation costs an amortized $\mathcal{O}(1)$ heap operations [16]. Hence, a total of $\mathcal{O}(|\mathcal{N}|\log|\mathcal{N}| + |\mathcal{P}|)$ heap operations are performed. Taking the size of the integers representing priorities into account, the heap operations cost $\mathcal{O}(|\mathcal{N}|(|\mathcal{N}|\log|\mathcal{N}| + |\mathcal{P}|))$.

### 3.3  Deriving shortest paths for path query results

Using the above results, we can already answer context-free path queries under the single-path semantics, this by applying MINIMIZESET on an annotated grammar. Unfortunately, this approach has high overhead due to the explicit construction and storing of the annotated grammar. Luckily, during the execution of MINIMIZESET, the relevant parts of $\mathscr{C}|_{\mathfrak{G}}$ can be implicitly derived from $\mathscr{C}$ and $\mathfrak{G}$. We obtain the MINIMIZESETGG algorithm by integrate these implicit derivation steps into MINIMIZESET. The resulting pseudo-code can be found in Figure 2, *right*. We conclude:
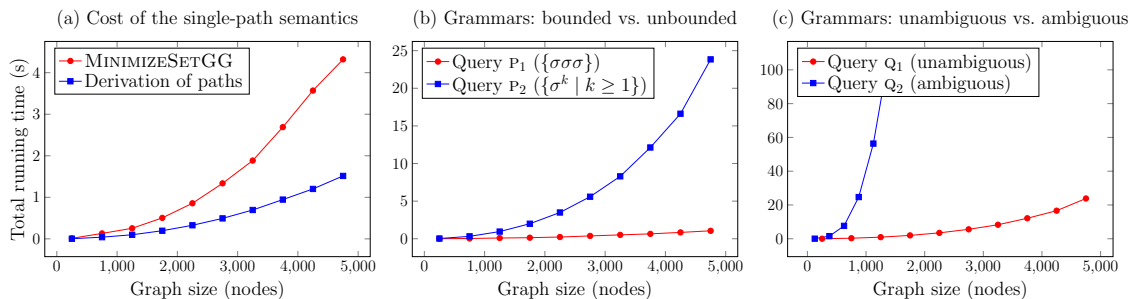
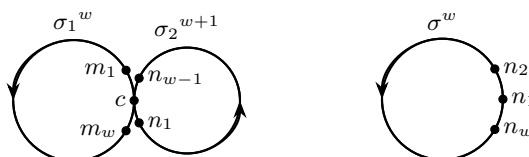**Fig. 3.** Measurements on the performance of MINIMIZESETGG.



**Fig. 4.** On the *left*, the double-cyclic graph: two cycles, one having $w - 1$ edges labeled with $\sigma_1$, and one having $w$ edges labeled with $\sigma_2$. The two cycles are connected via a shared node $c$. On the *right*, the cyclic graph having $w$ nodes labeled with $\sigma$.

**Theorem 2.** *Let $\mathscr{C} = (\mathcal{N}, \Sigma, \mathcal{P})$ be a grammar, $\mathfrak{G} = (\mathcal{V}, \Sigma, \delta)$ be a graph, and $\textsc{a} \in \mathcal{N}$ a context-free path query. We can evaluate $\mathrm{single}(\textsc{a}|_{\mathfrak{G}})$ using* MINI-MIZESETGG *in $\mathcal{O}(|\mathcal{N}||\mathcal{V}|^2(|\mathcal{N}||\mathcal{V}|^2 \log(|\mathcal{N}||\mathcal{V}|^2) + |\mathcal{P}|(|\mathcal{V}|^3 + |\delta|)) + L)$, in which $L = \sum\{|\pi| \mid \pi \in \mathrm{single}(q|_{\mathfrak{G}})\}$ is total length of the shortest paths in the result.*

## 4   Empirical Evaluation

To show that the path-based semantics for context-free path are viable in practice, we implemented the MINIMIZESETGG algorithm of Figure 2, *right*, and the straightforward path derivation algorithm in `C++`14. Open-source code of the full `C++`14 implementation of the data structures, algorithms, and supporting tooling used can be found at https://www.jhellings.nl/projects/cfpqpaths/. Using this implementation, we ran three different experiments to study the behavior of the MINIMIZESETGG algorithm. The programs were compiled and run on a workstation with an Intel Core i5-4670 CPU, running at a maximum of $3.8$ GHz, and with $16$ GiB of main memory. In each of our experiments, we test with synthetic graphs that are designed specifically to test extreme-case behavior of the algorithms. Visualizations of these graphs can be found in Figure 4.

*Cost of the single-path semantics.* As the first experiment, we study the cost of evaluating context-free path queries using the single-path semantics. To put MINIMIZESETGG to the test, we run these experiments with the grammar

$$\textsc{q} \mapsto \textsc{a } \textsc{q}', \qquad \textsc{q}' \mapsto \textsc{q } \textsc{b}, \qquad \textsc{q} \mapsto \textsc{a } \textsc{b}, \qquad \textsc{a} \mapsto \sigma_1, \qquad \textsc{b} \mapsto \sigma_2,$$

which is context-free and cannot be expressed by a regular language. We use the double-cyclic graphs of Figure 4, as this combination of query and graphs produces very large path. More specifically, we have proven that, in this case, the longest shortest paths have a size that is *quadratic* in the size of the graph, whereas for all single-symbol grammars and regular grammars the maximum size is only *linear* in the size of the graph (details omitted due to space limitations). The results of the experiment can be found in Figure 3(a). As is clear from the results, single-path evaluation is practically feasible: the query Q, evaluated on a double-cyclic graph of 4750 nodes, yields a set of $11 \cdot 10^6$ distinct paths, of which the longest (non-simple) path has $11 \cdot 10^6$ edges, and the average path has $5.6 \cdot 10^6$ edges. Hence, the query result is large. Still, MINIMIZESETGG finished in only 4.3 s and the longest path was constructed in 1.5 s. Hence, even for queries and graphs that produce very large results, the query costs are reasonable.

*Grammars: bounded vs. unbounded.* In the second experiment, we take a more in-depth look of the cost of context-free path query evaluation. In practice, many path queries are *bounded* in the sense that only paths of a limited length are inspected in the graph. E.g., to ask for friends-of-friends in a social network, one only has to inspect paths of length two. Some path queries, however, are *unbounded*, as context-free path queries can use recursion. This is of use, e.g., to query for pairs of indirect friends (Example 1). As unbounded queries can yield much larger result sets than bounded queries, we inspect the impact of the type of queries on the running time of MINIMIZESETGG. For this experiment, we use the queries $P_1$ (bounded) and $P_2$ (unbounded):

$$P_1 \mapsto S\ B \qquad\qquad B \mapsto S\ S \qquad\qquad S \mapsto \sigma;$$
$$P_2 \mapsto S\ P_2 \qquad\qquad P_2 \mapsto \sigma \qquad\qquad S \mapsto \sigma.$$

The language described by $P_1$ is $\mathcal{L}_1 = \{\sigma\sigma\sigma\}$, and the language described by $P_2$ is $\mathcal{L}_2 = \{\sigma^k \mid k \geq 1\}$. We use the cycle graphs with $w$ nodes of Figure 4, on which query $P_1$ will evaluate to a very sparse result set of $w$ paths, whereas query $P_2$ will evaluate to a very dense result set of $w^2$ paths. We measured the running time of MINIMIZESETGG for both queries. The results of the experiment can be found in Figure 3(b). As is clear from the results, the performance of single-path evaluation depends largely on the size of the query result. On the one hand, query evaluation for $P_1$, a bounded query yielding a small result set, finished within a second on all graphs. On the other hand, query evaluation for $P_2$, an unbounded query yielding large result sets, produced a result set $22 \cdot 10^6$ paths on the larges graph and did so in 22 s. We notice that the size of the result set is the limiting factor here: in the previous experiment, we already demonstrated that MINIMIZESETGG can easily deal with very large paths constructed by complex context-free path queries.

*Grammars: unambiguous vs. ambiguous.* In the third and final experiment, we look at the impact of the design of context-free path queries on the cost of their evaluation. This experiment is inspired by well-known results from parsing

and compiler construction (see, e.g., [17]): for grammars that are deterministic and unambiguous, e.g., LL($k$) or LR($k$) grammars, simple high-performance parsers with a linear running time exist. For non-deterministic and for ambiguous grammars, such high-performance parsers do not exist, however. The MINIMIZE-SETGG algorithm we propose works on all grammars, even grammars that are non-deterministic and ambiguous. This raises the question whether the type of grammars impacts the overall performance. To answer this question, we construct two equivalent queries $Q_1$ (unambiguous) and $Q_2$ (ambiguous):

$$Q_1 \mapsto S\ Q_1 \qquad\qquad Q_1 \mapsto \sigma \qquad\qquad S \mapsto \sigma;$$
$$Q_2 \mapsto Q_2\ Q_2 \qquad\qquad Q_2 \mapsto \sigma.$$

Both queries specify the language $\mathcal{L} = \{\sigma^k \mid k \geq 1\}$. As in the previous experiment, we use cycle graphs. On these cycle graphs, both queries will evaluate to very dense results sets. We measured the running time of MINIMIZESETGG for both queries. The results of the experiment can be found in Figure 3(c). As is clear from the results, evaluation of the unambiguous query $Q_1$ is magnitudes faster than evaluation of the ambiguous query $Q_2$, even though MINIMIZESETGG does not yet optimize for deterministic or unambiguous grammars. The reason for this is simple: for any shortest path in the graph, $Q_2$ has many different ways to derive the trace of this path, whereas $Q_1$ only has a single derivation. Consequently, MINIMIZESETGG will have to inspect many more choices while evaluating $Q_2$. Still, we believe that further optimizations for deterministic and unambiguous grammars are possible, a direction we leave open for future work.

## 5   Related Work

There is an abundant literature on graph queries, formal languages, and context-free grammars. There is only limited work toward answering graph queries with paths, however. Likewise, there is only limited work on the related problem of deriving shortest strings from grammars. Next, we give a brief overview.

As stated before, path-based semantics have only gained limited attention. For the regular expressions, Barceló et al. [6] introduced the extended regular path queries that have path variables for output. The main focus of Barceló et al. is, however, on the use of path variables for expressivity purposes, and path-based results are only studied in limited details. Recent work by Hofman et al. [21] provides an alternative to use path-based query semantics for debugging: to gain more insight in the behavior of regular path queries with respect to the expected behavior, Hofman et al. propose a technique based on separability. Although this approach addresses query debugging, it does not lift the other limitations of the traditional query semantics used to evaluate path queries. In practical graph database systems, path-based results can already be used in some limited settings [2]. E.g., SPARQL can return RDF graphs via `CONSTRUCT` queries [19], which can be used to encode fixed-size paths; whereas Gremlin can enumerate graph traversal steps (which can encode paths) via the `.path()` step,

which comes at prohibitive high costs [29]. In the setting of model checking using CTL [12], path-based query semantics are widely used to produce witnesses and counterexamples that show why the graph does or does not meet the conditions expressed by the CTL formulae. Unfortunately, model checking languages lack the expressive power found in most path query languages used to query graph databases. This sharply contrasts our work, as we show that path-based results are viable both in theory and in practice, this even for complex context-free path queries. Hence, to the best of our knowledge, our work is the first to systematically formalize and study path results for complex graph query languages.

Barrett et al. [7] studies variations of the single-path semantics we propose in this work. They do so from a complexity-theoretical standpoint, however, by classifying the complexity of query evaluation using variations of our single-path semantics. E.g., they show that the single-path semantics is feasible for regular path queries and context-free path queries, but becomes unfeasible when only simple paths are to be returned. As their focus is on classifying the complexity of evaluation, Barret et al. do not provide practical algorithms for the evaluation of path queries using the single-path semantics. We improve on this work by providing the algorithm MINIMIZESETGG, an efficient algorithm for evaluating context-free path queries on graphs using the single-path semantics.

Finally, we have shown that the evaluation of context-free path queries on graphs using the single-path semantics can be reduced to the derivation of a shortest string from a grammar. Mclean et al. [27] proved that such a shortest string could be computed effective given a grammar, but failed to give a practical algorithm for doing so. We improve on these results by providing the algorithm MINIMIZESET, an efficient algorithm for computing the shortest string in a grammar. Other works, e.g. [13,15,14,26], provide ways to enumerate strings in a grammar, but these algorithms cannot effectively be used to quickly find the shortest such string.

## 6   Conclusions and Future Work

To address the limitations of the traditional semantics for evaluating path queries, such as the regular path queries and the context-free path queries, we proposed the single-path semantics. This path-based semantics is not only useful for end-users, but also enables new directions in the design of graph query languages and enables new tools for graph analytics, data exploration, data provenance, and debugging of complex path queries. To show the practical viability of the single-path semantics, we also propose algorithms that evaluate context-free path queries using the single-path semantics. Our initial results are promising: our experimental evaluation shows that queries can be evaluated using the single-path semantics with little effort, even in cases where the path-based query results are very large. Based on our initial results, we see several avenues for the further study of evaluating queries with path-based semantics:

1. The algorithms in our paper are *bottom-up* and are tuned toward evaluating a query over the entire graph. In many practical applications, the end-user is

only interested in a part of the graph, e.g., paths that originate or end at a certain node. For such applications, we are interested in the development of *top-down* and *goal-oriented* algorithms.

2. Our measurements showed that the cost of evaluating a context-free path query depends heavily on the structure of the grammar used by the query: evaluating different grammars that express the same query can have widely different costs. This raises an interesting query optimization question: can we automatically optimize grammars to reduce the cost of evaluation?

3. Furthermore, it is open whether simpler, more efficient, query evaluation algorithms exist for restricted classes of context-free grammars (e.g., deterministic grammars or unambiguous grammars [17]). It is not directly clear if such algorithms exist: deterministic and unambiguous grammars will still face ambiguity and non-deterministic choices in their evaluation on graphs, as complex graphs can have many paths with the same traces.

# References

1. Alon, U.: An Introduction to Systems Biology: Design Principles of Biological Circuits. Chapman and Hall/CRC (2006)
2. Angles, R., Arenas, M., Barceló, P., Hogan, A., Reutter, J., Vrgoč, D.: Foundations of modern query languages for graph databases. ACM Comput. Surv. **50**(5), 68:1–68:40 (2017)
3. Arenas, M., Conca, S., Pérez, J.: Counting beyond a yottabyte, or how SPARQL 1.1 property paths will prevent adoption of the standard. In: Proceedings of the 21st International Conference on World Wide Web. pp. 629–638. ACM (2012)
4. Bar-Hillel, Y., Perles, M.A., Shamir, E.: On formal properties of simple phrase structure grammars. Zeitschrift für Phonetik, Sprachwissenschaft und Kommunikationsforschung **14**, 143–172 (1961)
5. Barceló, P.: Querying graph databases. In: Proceedings of the 32nd Symposium on Principles of Database Systems. pp. 175–188. ACM (2013)
6. Barceló, P., Libkin, L., Lin, A.W., Wood, P.T.: Expressive languages for path queries over graph-structured data. ACM Trans. Database Syst. **37**(4), 31:1–31:46 (2012)
7. Barrett, C., Jacob, R., Marathe, M.: Formal-language-constrained path problems. SIAM J. Comput. **30**(3), 809–837 (2000)
8. Berglund, A., Boag, S., Chamberlin, D., Fernández, M.F., Kay, M., Robie, J., Siméon, J.: XML path language (XPath) 2.0 (second edition). Tech. rep., W3C (2010), http://www.w3.org/TR/2010/REC-xpath20-20101214/
9. Bray, T., Paoli, J., Sperberg-McQueen, C.M., Maler, E., Yergeau, F., Cowan, J.: Extensible markup language (XML) 1.1 (second edition). Tech. rep., W3C (2006), http://www.w3.org/TR/2006/REC-xml11-20060816
10. Cheney, J., Chiticariu, L., Tan, W.C.: Provenance in databases: Why, how, and where. Foundations and Trends in Databases **1**(4), 379–474 (2009)
11. Clark, J., DeRose, S.: XML path language (XPath) version 1.0. Tech. rep., W3C (1999), http://www.w3.org/TR/1999/REC-xpath-19991116/
12. Clarke, E.M., Grumberg, O., Peled, D.: Model Checking. The MIT Press (1999)
13. Dömösi, P.: Unusual algorithms for lexicographical enumeration. Acta Cybern. **14**(3), 461–468 (2000)

14. Dong, Y.: Linear algorithm for lexicographic enumeration of CFG parse trees. Science in China Series F: Information Sciences **52**(7), 1177–1202 (2009)
15. Florêncio, C.C., Daenen, J., Ramon, J., den Bussche, J.V., Dyck, D.V.: Naive infinite enumeration of context-free languages in incremental polynomial time. J. Univers. Comput. Sci. **21**(7), 891–911 (2015)
16. Fredman, M.L., Tarjan, R.E.: Fibonacci heaps and their uses in improved network optimization algorithms. J. ACM **34**(3), 596–615 (1987)
17. Grune, D., Jacobs, C.J.: Parsing Techniques: A Practical Guide. Springer-Verlag New York, 2nd edn. (2008)
18. Harel, D., Pnueli, A., Stavi, J.: Propositional dynamic logic of nonregular programs. J. Comput. Syst. Sci. **26**(2), 222–243 (1983)
19. Harris, S., Seaborne, A.: SPARQL 1.1 query language. Tech. rep., W3C (2013), http://www.w3.org/TR/2013/REC-sparql11-query-20130321
20. Hellings, J.: Conjunctive context-free path queries. In: Proceedings of the 17th International Conference on Database Theory (ICDT 2014). pp. 119–130 (2014)
21. Hofman, P., Martens, W.: Separability by short subsequences and subwords. In: 18th International Conference on Database Theory. Leibniz International Proceedings in Informatics (LIPIcs), vol. 31, pp. 230–246. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik (2015)
22. Kozen, D.: Kleene algebra with tests. ACM Trans. Program. Lang. Syst. **19**(3), 427–443 (1997)
23. Lange, M.: Model checking propositional dynamic logic with all extras. J. Appl. Log. **4**(1), 39–49 (2006)
24. Libkin, L., Martens, W., Vrgoč, D.: Querying graph databases with XPath. In: Proceedings of the 16th International Conference on Database Theory. pp. 129–140. ACM (2013)
25. Losemann, K., Martens, W.: The complexity of evaluating path expressions in SPARQL. In: Proceedings of the 31st ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems. pp. 101–112. ACM (2012)
26. Mäkinen, E.: On lexicographic enumeration of regular and context-free languages. Acta Cybern. **13**(1), 55–61 (1997)
27. Mclean, M.J., Johnston, D.B.: An algorithm for finding the shortest terminal strings which can be produced from non-terminals in context-free grammars. In: Combinatorial Mathematics III, Lecture Notes in Mathematics, vol. 452, pp. 180–196. Springer Berlin Heidelberg (1975)
28. Robinson, I., Webber, J., Eifrem, E.: Graph Databases: New Opportunities for Connected Data. O'Reilly Media, Inc., 2nd edn. (2015)
29. Rodriguez, M.A.: The gremlin graph traversal machine and language (invited talk). In: Proceedings of the 15th Symposium on Database Programming Languages. pp. 1–10. ACM, New York, NY, USA (2015)
30. Schreiber, G., Raimond, Y.: RDF 1.1 primer. Tech. rep., W3C (2014), http://www.w3.org/TR/2014/NOTE-rdf11-primer-20140624
31. Sevon, P., Eronen, L.: Subgraph queries by context-free grammars. J. Integr. Bioinform. **5**(2), 157–172 (2008)