# A Framework for Federated SPARQL Query Processing over Heterogeneous Linked Data Fragments

Lars Heling
lars.heling@kit.edu
Karlsruhe Institute of Technology
Karlsruhe, Germany

Maribel Acosta
maribel.acosta@rub.de
Ruhr-University Bochum
Bochum, Germany

## ABSTRACT

Linked Data Fragments (LDFs) refer to Web interfaces that allow for accessing and querying Knowledge Graphs on the Web. These interfaces, such as SPARQL endpoints or Triple Pattern Fragment servers, differ in the SPARQL expressions they can evaluate and the metadata they provide. Client-side query processing approaches have been proposed and are tailored to evaluate queries over individual interfaces. Moreover, federated query processing has focused on federations with a single type of LDF interface, typically SPARQL endpoints. In this work, we address the challenges of SPARQL query processing over federations with heterogeneous LDF interfaces. To this end, we formalize the concept of federations of Linked Data Fragment and propose a framework for federated querying over heterogeneous federations with different LDF interfaces. The framework comprises query decomposition, query planning, and physical operators adapted to the particularities of different LDF interfaces. Further, we propose an approach for each component of our framework and evaluate them in an experimental study on the well-known FedBench benchmark. The results show a substantial improvement in performance achieved by devising these interface-aware approaches exploiting the capabilities of heterogeneous interfaces in federations.

## 1 INTRODUCTION

The increasing number and size of Knowledge Graphs published as Linked Data led to the development of different interfaces to support querying Knowledge Graphs on the Web [7, 12, 15, 27]. These interfaces mainly differ in their expressivity, server availability, and client cost as shown in Figure 1. The Linked Data Fragment (LDF) framework provides a uniform way to describe these interfaces regarding their querying expressivity and the metadata they provide [13, 27]. The query expressivity of these interfaces ranges from triple patterns in Triple Pattern Fragment (TPF) servers to the full fragment of SPARQL in SPARQL endpoints. To support efficient querying, these developments also drove the research in the area of client-side SPARQL query processing tailored to the individual interfaces [3, 16, 26]. Therefore, most of the existing approaches focus on querying data from a single dataset or through a federation of sources with the same interfaces. However, the problem of evaluating SPARQL queries over heterogeneous federations of such LDF interfaces has not gained much attention thus far. To devise efficient querying plans in this scenario, it is not sufficient to combine and fuse existing solutions because the capabilities and limitations of the interfaces have to be taken into account altogether. An effective solution, thus, requires to re-define the notions of the main tasks of federated engines – i.e query decomposition,
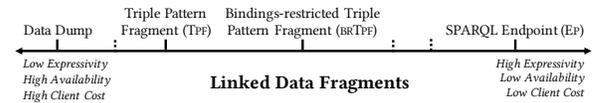


Figure 1: Linked Data Fragment Spectrum (based on [27]).

planning, and execution – to integrate the different capabilities of the interfaces in a single query plan.

In this work, we formalize federations of LDF interfaces and propose a framework that serves as a foundation for devising efficient solutions for querying heterogeneous federations. At the core of the framework, we present novel concepts for query decomposition, planning, and execution in heterogeneous federations. These components include the definitions of (i) interface-compliant sub-expressions and the semantics of their evaluation, (ii) interface-aware query planning, and (iii) polymorphic physical operators that implement different execution strategies according to the contacted interfaces. To accompany our theoretical contributions, we propose simple yet novel approaches to query heterogeneous LDF federations. Each approach addresses the particularities of the interfaces and is designed to reduce the query execution times and the load on the members of the federations. Our results show the effectiveness of our framework and illustrate how leveraging the interfaces' capabilities in a single plan can substantially improve query execution. In summary, the contributions of this work are

- a general definition of Linked Data Fragment (LDF) federations,
- a framework for querying heterogeneous LDF federations addressing query decomposition, planning, and physical operators,
- a practical solution for each component of the framework, and
- an experimental evaluation of a prototypical implementation of the solutions on heterogeneous federations.

The remainder of this work is structured as follows. Section 2 presents a motivating example, and in Section 3, we present our definition of federations of LDFs. Our framework is presented in Section 4 and evaluated in Section 5. We discuss related work in Section 6 and conclude our work in Section 7.

## 2 MOTIVATING EXAMPLE

As a motivating example, consider the query to *get American Presidents, the political party they are a member of as well as their predecessors and successors* shown in Listing 1.

**Listing 1: Example query. Prefixes as in http://prefix.cc**

```
SELECT * WHERE {
    ?x wdt:P39 wd:Q11696 .              # tp₁
    ?x wdt:P102 ?party  .              # tp₂
    ?y owl:sameAs ?x .                 # tp₃
    ?y dbo:predecessor ?predecessor .  # tp₄
    ?y dbo:successor ?successor  .  }  # tp₅
```

Let us assume, we want to evaluate the query over a federation that consists of the SPARQL endpoint of Wikidata[1] and the Triple Pattern Fragment (TPF) server [27] of DBpedia[2]. The Wikidata endpoint provides solutions to triple patterns $tp_1, tp_2, tp_3$ and the DBpedia TPF server to $tp_3, tp_4$, and $tp_5$. As the members of the federation implement different Linked Data Fragment (LDF) interfaces, we denote such a federation *heterogeneous*. In this example, we are not able to apply an existing query decomposition approach from query processing over SPARQL endpoint federations, as these approaches do not consider the capabilities of the LDF interfaces. For instance, FedX [24] would group triple patterns $tp_4$ and $tp_5$ into a subquery which is not *compliant* with the DBpedia TPF interface. On the contrary, a naive decomposition that evaluates the query on the triple pattern level at the relevant sources would be *compliant* with the interfaces in the federation, since they are all able to evaluate triple patterns. However, such a decomposition leads to inefficient query plans on SPARQL endpoints as they produce an excessive number of requests on the server and thus, lead to long execution times. In our example, this approach would fail to evaluate the subexpression ($tp_1$ AND $tp_2$) at the Wikidata endpoint. Their individual evaluation at the endpoint leads to an overhead in requests and intermediate results transferred that could be avoided. The example illustrates the challenges that arise in heterogeneous federations and motivates our research to address those challenges. In this work, we propose a framework that is tailored to leverage the capabilities of the different interfaces in heterogeneous federations. Based on this framework, our implementation reduces the number of requests by almost 25% leading to a tenfold decrease in query execution time over the naive approach for the example query.

# 3　FEDERATIONS OF LINKED DATA FRAGMENT SERVICES

Following existing works [13, 27], we introduce a formalization of Linked Data Fragment (LDF) interfaces based on the SPARQL expressions they are able to evaluate and the metadata they provide. Verborgh et al. [27] define a Linked Data Fragment (LDF) for an RDF graph $G$ as a tuple consisting of a URI, a selector function, a set of RDF triples that are the result of applying the selector function over $G$, metadata in the form of a set of RDF triples, and a set of hypermedia controls. Based on this work, Hartig et al. [13] propose a formal framework for comparing LDF interfaces in terms of expressiveness, complexity, and performance when evaluating SPARQL queries over different interfaces. The concept of LDF interfaces by Hartig et al. comprises the following: i) a notion of a *server language* to differentiate between different capabilities of LDF interfaces, and ii) an evaluation function in which an LDF interface provides a set of SPARQL solution mappings upon

a request. We adapt the server language definition from [13] to be based on the SPARQL expressions an LDF interface can evaluate. Therefore, we first revise the SPARQL expressions considered in the literature.

Let the sets of RDF terms $U$, $B$, and $L$ be pairwise disjoint sets of URIs, blank nodes, and Literals, and $V$ be a set of variables disjoint from $U$, $B$, and $L$. A triple $(s, p, o) \in (U \cup B) \times U \times (U \cup B \cup L)$ is called an RDF triple. A *set* of RDF triples is an RDF graph $G$ and the universe of RDF graphs is denoted as $\mathcal{G}$. Following the notation by Peréz et al. [18] and Schmidt et al. [23], SPARQL expressions are constructed using the operators AND, UNION, OPTIONAL, FILTER, and VALUES and can be defined recursively as follows.

*Definition 3.1 (SPARQL Expression).* A SPARQL expression is an expression that is recursively defined as follows.

(1) A triple pattern $tp \in (U \cup V) \times (U \cup V) \times (U \cup L \cup V)$ is a SPARQL expression [18],
(2) if $P_1$ and $P_2$ are SPARQL expressions, then the expressions ($P_1$ AND $P_2$), ($P_1$ UNION $P_2$) and ($P_1$ OPTIONAL $P_2$) are SPARQL expressions (*conjunctive expression, union expression, optional expression*) [18],
(3) if $P$ is a SPARQL expression and $R$ is a SPARQL filter condition, then the expression $P$ FILTER $R$ is a SPARQL expression (*filter expression*) [18],
(4) if $P$ is a SPARQL expression and $D$ is a SPARQL values datablock, the expression $P$ VALUES $D$ is a SPARQL expression (*values expression*), and
(5) if $P$ is SPARQL expression and $S \subset V$ is a set of variables, the expression $\text{SELECT}_S(P)$ is an expression (*select query*) [23].

Furthermore, we denote the universe of SPARQL expressions as $\mathcal{P}$ and $vars(P) \subset V$ as the set of variables in the expression $P$. We can define the interface languages of different LDF interfaces by means of the fragment of SPARQL expressions they can evaluate.

*Definition 3.2 (Interface Language).* Let $\mathcal{L}$ be the universe of interface languages, an interface language $L \in \mathcal{L}$ is the fragment of SPARQL expressions that an interface can evaluate.

Moreover, we denote $P \in L$ if a SPARQL expression $P$ is a SPARQL expression that is part of an interface language $L$. Common interface languages can, thus, be defined in the following way.

- $L_{\text{CoreSparql}}$: Any SPARQL expression defined in Def. 3.1.
- $L_{\text{BGP}}$: Conjunctive expressions: ($P_1$ AND $P_2$) where $P_1$ and $P_2$ are either conjunctive expressions or triple patterns,
- $L_{\text{Tp}}$: Triple patterns.
- $L_{\text{Tp+Values}}$: Triple patterns and values expressions of the form $P$ VALUES $D$, where $P$ is a triple pattern.

The definition of interface languages based on SPARQL expression also allows for defining containment relations between the different languages according to their expressiveness.

*Definition 3.3 (Interface Language Containment).* Let $L_1$ and $L_2$ be two interface languages, we say that $L_1$ is contained in $L_2$, if all SPARQL expressions in $L_1$ are also in $L_2$:

$$L_1 \subseteq L_2, \text{if } \forall P \in L_1 \Rightarrow P \in L_2.$$

For example, we can state the following containment relations for the previously introduced languages: $L_{\text{Tp}} \subseteq L_{\text{BGP}} \subseteq L_{\text{CoreSparql}}$,

or $L_{\text{Tp}} \subseteq L_{\text{Tp+Values}}$. With this formalism, we can define the languages of common LDF interfaces and compare them according to their expressiveness. Triple Pattern Fragment (TPF) servers support querying triple patterns ($L_{\text{Tp}}$) [27], bindings-restricted TPF server support triple patterns and Values expressions ($L_{\text{Tp+Values}}$) [12], and SPARQL endpoint support any expression ($L_{\text{CoreSparql}}$).

To complement the definition of LDF interfaces, we introduce the concept of *interface metadata*. For a SPARQL expression $P$, an interface may provide interface-specific metadata $M(P)$ describing the data obtained from the RDF graph. The interface metadata may range from simple statistics such as the number of expected results to more elaborate metadata describing statistics, provenance, and licensing information. Similar to [27], we assume the metadata provided for a given expression $P$ to be an RDF graph, that is $M : \mathcal{P} \to \mathcal{G}$. Examples for common interface metadata are:

- SPARQL endpoints $M_{\text{Ep}}$: $M(P) = \emptyset$, $\forall P \in L_{\text{CoreSparql}}$,
- Triple Pattern Fragments $M_{\text{TPF}}$: $M(P)$ is an RDF graph that contains an estimate of the number of triples that match the expression $P$, $\forall P \in L_{\text{Tp}}$,
- Bindings-restricted Triple Pattern Fragments $M_{\text{BRTPF}}$: $M(P)$ is an RDF graph that contains an estimate of the number of triples that match the expression $P$, $\forall P \in L_{\text{Tp+Values}}$.

Besides enabling a more fine-grained distinction between LDF interfaces, the metadata may impact the potential querying strategies employed by a client. Finally, combining interface language and metadata, we define a Linked Data Fragment interface as follows.

*Definition 3.4 (Linked Data Fragment Interface).* A Linked Data Fragment interface is a 2-tuple $f = (L_f, M_f)$, where

- $L_f \in \mathcal{L}$, the interface language,
- $M_f : \mathcal{P} \to \mathcal{G}$, the interface metadata for an expression $P$.

Conceptually, we distinguish LDF *interfaces* which define the interface language and metadata, and LDF *services*, which are Web servers that implement a specific interface.

*Definition 3.5 (Linked Data Fragment Service).* A Linked Data Fragment service $c \in U$ is a Web service that supports the evaluation of SPARQL expressions and provides metadata according to the LDF interface $int(c) = (L_c, M_c)$ that it implements.

We reuse the function $ep : U \to \mathcal{G}$ [6] that maps an LDF service to the RDF graph $ep(c)$ available at the service. The evaluation of a SPARQL expression $P$ over an LDF service $c$ is then given as

$$[\![P]\!]_c := \begin{cases} [\![P]\!]_{ep(c)}, & \text{if } P \in L_c. \\ \emptyset, & \text{otherwise.} \end{cases} \tag{1}$$

Note the difference in the subscript $c$ and $ep(c)$ to distinguish between solution mappings produced by an LDF service $[\![\cdot]\!]_c$ regarding its interface language and the solution mappings for evaluating any expression over the graph available at the LDF service $[\![\cdot]\!]_{ep(c)}$. Combining these previous definitions, we define the concept of federations of LDF services as follows.

*Definition 3.6 (Federation of Linked Data Fragment Services).* A Federation of Linked Data Fragment services is a 3-tuple $F = (C, int, ep)$, where

- $C = \{c_1, \ldots, c_n\} \subset U$, a set of URIs for LDF services,

- *int*, a function that maps an LDF service to its interface,
- *ep*, a function that maps each LDF service to the graph available at that service.

Federations in which all LDF services implement the same LDF interfaces are called *homogeneous*, and *heterogeneous* otherwise. For practical reasons, in the remainder of this work, we just consider graphs in the federation without blank nodes and focus on federations in which all members are at least able to evaluate triple patterns of any form: $L_{\text{TP}} \subseteq L_c, \forall c \in C$. [3]

*Example 3.7.* We can define the federation from our motivating example as $F_{ex} = (\{c_1, c_2\}, int, ep)$ with $c_1 = \texttt{wikidata:sparql}$[1], $c_2 = \texttt{dbpedia:tpf}$[2], $int(c_1) = (L_{\text{CoreSparql}}, M_{\text{Ep}})$, $int(c_2) = (L_{\text{Tp}}, M_{\text{TPF}})$, $ep(c_1) = G_{Wikidata}$ and $ep(c_2) = G_{DBpedia}$.

Following the notation by Acosta et al. [2], we denote the evaluation of a SPARQL expression over a federation of LDF interfaces $F$ as $[\![\cdot]\!]_F$ and define the semantics in the following way.

*Definition 3.8 (Set Semantics of SPARQL Query Processing over LDF Service Federations).* Given a SPARQL expression $P$ and a federation $F = (C, int, ep)$, the result set of evaluating $P$ over $F$ is given as

$$[\![P]\!]_F := [\![P]\!]_G, \text{ with } G = \bigcup_{\forall c \in C} ep(c)$$

## 4 FEDERATED QUERY PROCESSING OVER HETEROGENEOUS FEDERATIONS

In the presence of heterogeneous LDF service federations, novel challenges arise that cannot be addressed by existing approaches. Therefore, we propose a framework for heterogeneous federations to address the central components of federated query processing: (§4.1) query decomposition, (§4.2) query planning, and (§4.3) physical operators. Furthermore, for each component of the framework, we propose an approach aiming to obtain efficient query plans.

### 4.1 Query Decomposition

The goal of query decomposition is grouping the query into subexpressions such that the evaluation of the subexpressions over the members of the federation minimizes execution time while ensuring that all expected answers are produced. Existing decomposition approaches assume that all federation members are able to evaluate any SPARQL expression. Since this assumption is not valid in heterogeneous federations, we propose *interface-compliant* query decompositions and their evaluation over such federations.

Given a given SPARQL query $P$ and a federation $F = (C, int, ep)$, query decomposition aims to group a query into subexpressions such that they can be answered by the relevant sources in the federation. The first step to achieve this goal is source selection. That is, select the relevant sources $r(tp_i)$ for all triple patterns $tp_i$ in $P$, with $r(tp_i) = \{c \in C \mid [\![tp_i]\!]_{ep(c)} \neq \emptyset\}$. Because we require all LDF services to at least evaluate triple patterns of any form ($L_{\text{TP}} \subseteq L_c, \forall c \in C$), in principle, the relevant sources can be selected by evaluating each triple pattern at each service. Typically, the capabilities of the services allow more efficient source selection

---

[3]This means that we do not include data dumps, even though they are also considered Linked Data Fragments in other works [13, 27].

strategy implementations, such as ASK queries for SPARQL end-points, leveraging the metadata $(u, \text{void:triples}, cnt) \in M(tp_i)$ for TPF and brTPF servers, or using pre-computed data catalogues.

Once the relevant sources are identified, the query engine decomposes the query into subexpressions to be evaluated at the services in the federations. If there exists a triple pattern $tp$ in a basic graph pattern (BGP) $P$ with no relevant source $r(tp) = \emptyset$, the evaluation of $P$ over the federation is the empty set. In the following, we focus on query decompositions for BGPs where all triple patterns have at least one relevant source. For simplicity, we extend notation and consider a BGP $P = (tp_1 \text{ AND } \ldots \text{ AND } tp_n)$ also as a set of $n$ triple patterns: $P = \{tp_1, \ldots, tp_n\}$.

*Definition 4.1 (Query Decomposition).* Given a BGP $P$ and an LDF service federation $F = (C, int, ep)$, a query decomposition $D(P, F) = \{d_1, \ldots, d_m\}$ is a set of tuples $d_i = (SE_i, S_i)$ where

- $SE_i$ is a subexpression of $P$, and
- $S_i \subseteq C$ a non-empty the subset of services over which $SE_i$ is evaluated, such that $\bigcup\limits_{i=1,\ldots,n} SE_i = P$.

Because the query decomposition as such does not consider the interface language of the services, it is possible that for a valid query decomposition $D$: $\exists d_i \in D : \exists c \in S_i : SE_i \notin L_c$. Considering the query from our motivating example, a valid decomposition would be $D(P, F) = \{((tp_1 \text{ AND } tp_2), \{c_1\}), ((tp_3 \text{ AND } tp_4 \text{ AND } tp_5), \{c_2\})\}$, even though $c_2$ is a TPF server that can only evaluate triple patterns. Therefore, we introduce an evaluation function $\theta$ for the interface-compliant evaluation of SPARQL expressions.

*Definition 4.2 (Interface-compliant Evaluation of an Expression).* Given a BGP $P$ and an LDF service $c \in U$, the interface-compliant evaluation of $P$ over $c$ is given as follows.

$$\theta_c(P) := \begin{cases} [\![P]\!]_c & \text{if } P \in L_c. & (2) \\ [\![P_1]\!]_c \bowtie \ldots \bowtie [\![P_l]\!]_c & \text{otherwise.} & (3) \end{cases}$$

For some $P_1, \ldots, P_l$ with $P = (P_1 \text{ AND } \ldots \text{ AND } P_l)$ in Equation (3), such that the following conditions hold:

- Equivalence: $\theta_c(P) \equiv [\![P]\!]_{ep(c)}$
- Compliance: $P_i \in L_c, \ \forall P_i$ in $\{P_1 \ldots P_l\}$

The intuition of the interface-compliant evaluation is as follows. If the BGP $P$ is in the language of the service, then $P$ can be evaluated directly at the service. Otherwise, the original expression $P$ is split into subexpressions, such that each subexpression is in the language of the service. As a result, joining the solutions of evaluating the individual subexpressions at the service yields the same solutions as evaluating $P$ over the graph of the service $ep(c)$. We denote the number of subexpressions in the interface-compliant evaluation as $|\theta_c(P)|$. A compliant evaluation of $SE = (tp_3 \text{ AND } tp_4 \text{ AND } tp_5)$ from our previous example at the DBpedia TPF server $c_2$ would be $\theta_c(SE) = [\![tp_3]\!]_{c_2} \bowtie [\![tp_4]\!]_{c_2} \bowtie [\![tp_5]\!]_{c_2}$. With this notion, we define the interface-compliant evaluation of a query decomposition.

*Definition 4.3 (Interface-compliant Evaluation of a Query Decomposition).* Given a query decomposition $D(P, F)$ for the BGP $P$ and federation $F$, the evaluation of $P$ following the query decomposition $\theta_{D(P,F)}(P)$ is given as the conjunction ($\bowtie$) of the subexpressions $SE_i$ evaluated at all ($\cup$) services in $S_i$:

$$\theta_{D(P,F)}(P) := \bowtie_{(SE_i, S_i) \in D(P,F)} (\cup_{c_j \in S_i} \theta_{c_j}(SE_i))$$

After defining query decompositions and the evaluation of such decompositions that is compliant with the interfaces of the LDF services in the federations, the problem of finding a *suitable* query decomposition for a given query arises. The common goal of query decomposition approaches is finding a decomposition that yields complete answers according to the assumed semantics, while the cost of executing the decomposition by the query engine is minimized [10, 28]. However, these approaches do not explicitly measure the expected answer completeness of query decompositions. For instance, in [28] the answer completeness is encoded implicitly in the query decomposition cost by considering the number of non-selected endpoints: if fewer relevant endpoints are contacted according to a decomposition, its cost is higher and vice versa. Extending existing approaches, we propose the concept of query decomposition *density* as a measure to estimate and compare the expected answer completeness of different decompositions. In contrast to [28], our density measure not only considers the non-selected endpoints but also how triple patterns are grouped into subexpressions that are evaluated jointly at the services.

*Query Decomposition Density.* The query decomposition density is a proxy for the expected answer completeness of a decomposition. We define density as a relative measure with respect to a decomposition that guarantees answer completeness, i.e., the *atomic* decomposition. The atomic decomposition evaluates every single triple pattern in a subexpression at all relevant sources and thus, guarantees answer completeness.

*Definition 4.4 (Atomic Decomposition).* Given a federation $F = (\{c_1, \ldots, c_k\}, int, ep)$ and a BGP $P = (tp_1 \text{ AND } \ldots \text{ AND } tp_n)$, the atomic decomposition is given as

$$D^*(P, F) = \{(tp_1, r(tp_1)), \ldots, (tp_n, r(tp_n))\}.$$

LEMMA 4.5. *The evaluation of $P$ following $D^*(P, F)$ yields complete answers, that is:*

$$\theta_{D^*(P,F)}(P) = [\![P]\!]_F \qquad (4)$$

PROOF. We provide a direct proof by assuming the left-hand side in Eq. (4). Since we require all services to be able to evaluate triple patterns and $D^*$ is composed of triple patterns only, the evaluation of $D^*(P, F)$ is given by Def. 4.2 and Def. 4.3 as

$$\theta_{D^*(P,F)}(P) := \bowtie_{(SE_i, S_i) \in D^*(P,F)} (\cup_{c_j \in S_i} [\![SE_i]\!]_{c_j}) \qquad (5)$$

with $SE_i = tp_i$. By Def. 4.4, $S_i$ corresponds to the relevant sources of $tp_i$, which is given by $r(tp_i) = \{r_{i1}, \ldots, r_{im}\}$. Next, we expand the Eq. (5) with $S_i$ in the following way.

$$([\![tp_1]\!]_{r_{11}} \cup \cdots \cup [\![tp_1]\!]_{r_{1l}}) \bowtie \ldots \bowtie ([\![tp_n]\!]_{r_{n1}} \cup \cdots \cup [\![tp_n]\!]_{r_{no}}) \quad (6)$$

Next, we show that we can evaluate all triples patterns at all sources (relevant and non-relevant). By definition, we have that the evaluation of a triple pattern over a non-relevant source is the empty set: $[\![tp_i]\!]_c = \emptyset, \ \forall c \notin r(tp_i)$. Further, since $([\![tp_i]\!]_{r_{ij}} \cup \emptyset) = [\![tp_i]\!]_{r_{ij}}$, we can expand Eq. (6) to

$$([\![tp_1]\!]_{c_1} \cup \cdots \cup [\![tp_1]\!]_{c_k}) \bowtie \ldots \bowtie ([\![tp_n]\!]_{c_1} \cup \cdots \cup [\![tp_n]\!]_{c_k}) \quad (7)$$

According to Eq. (1) and the fact that triple patterns are in the interface language of all services, we have $[\![tp_i]\!]_{c_j} = [\![tp_i]\!]_{ep(c_j)}$ and can rewrite Eq. (7) as

$$([\![tp_1]\!]_{ep(c_1)} \cup \cdots \cup [\![tp_1]\!]_{ep(c_k)}) \bowtie \ldots \bowtie ([\![tp_n]\!]_{ep(c_1)} \cup \cdots \cup [\![tp_n]\!]_{ep(c_k)}) \quad (8)$$

Because we assume set semantics, the following equality holds[4]

$$([\![tp_i]\!]_{ep(c_1)} \cup \cdots \cup [\![tp_i]\!]_{ep(c_k)}) = [\![tp_i]\!]_{\bigcup_{c \in C} ep(c)} \quad (9)$$

With Eq. (9) we can reformulate Eq. (8) as

$$[\![tp_1]\!]_{\bigcup_{c \in C} ep(c)} \bowtie \ldots \bowtie [\![tp_n]\!]_{\bigcup_{c \in C} ep(c)} \quad (10)$$

and according to Def. 3.8 and Definition 4 in [23], we have the following equality:

$$[\![tp_1]\!]_{\bigcup_{c \in C} ep(c)} \bowtie \ldots \bowtie [\![tp_n]\!]_{\bigcup_{c \in C} ep(c)}$$
$$= [\![tp_1]\!]_F \bowtie \ldots \bowtie [\![tp_n]\!]_F$$
$$= [\![tp_1 \text{ AND } \ldots \text{ AND } tp_n]\!]_F$$
$$= [\![P]\!]_F$$

$\square$

Another type of structure that preserves completeness are exclusive groups [24], which are subexpressions of a query that can only be answered by a single source. They are defined as follows.

*Definition 4.6 (Exclusive Group).* Given a federation $F = (C, int, ep)$, a BGP $X$ is called an exclusive group, if for all triple patterns $tp_i \in X$ there exists only one relevant source $c_X \in C$:

$$X = \{tp_i \mid r(tp_i) = \{c_X\}\}$$

We represent query decompositions by *decomposition graphs* and compute the relative density with respect to the decomposition graph of the atomic decomposition $D^*(P, F)$ as a measure of completeness. More edges in the graph of a given decomposition yield a higher density and, thus, expected answer completeness.

*Definition 4.7 (Query Decomposition Graph).* Let $D(P, F)$ be a query decomposition for the BGP $P$ and federation $F = (C, int, ep)$. The decomposition graph $G_{D(P,F)} = (V, E)$ of $D(P, F)$ is:
The set of vertices $V = \{tp_i \in P\} \cup \{r(tp_i) \mid \forall tp_i \in P\}$,
The set of edges $E \subseteq V \times V$ are given by the following rules.

Rule I  Add an edge between a triple pattern $tp_i \in P$ and a relevant source $r_{ij} \in r(tp_i)$, if $tp_i$ is part of a subexpression $SE$ that is evaluated at $r_{ij}$: $\exists (SE, S) \in D(P, F)$ with $tp_i \in SE \wedge r_{ij} \in S$.
Rule II  Add an edge between two triple patterns $tp_i$ and $tp_j$, if they do not co-occur in a subexpression $SE$ in $D$: $(tp_i, tp_j) \in E$, if $\nexists SE \in D(P, F): tp_i \in SE \wedge tp_j \in SE$.

Rule III  Add an edge between two triple patterns $tp_i$ and $tp_j$, if they are part of the same exclusive group $X$: $(tp_i, tp_j) \in E$, if $tp_i \in X \wedge tp_j \in X$.
Rule IV  Add an edge between all triple patterns, if the decomposition is composed of just a single subexpression to be evaluated at one source : $D(P, F) = \{(SE, S)\} \wedge |S| = 1$.

The rules for adding edges to the graph are designed in such a way that the maximum number of edges is present for the decomposition graph of the atomic decomposition $G_{D^*(P,F)} = (V^*, E^*)$. This is because each triple pattern is connected to each relevant source (Rule I) and there is an edge between each pair of triple patterns (Rule II). If a decomposition contacts fewer sources, the decomposition graph will have fewer edges according to Rule I. Further, if more triple patterns are grouped together into subexpressions in a query decomposition, its graph will also have fewer edges according to Rule II. The rationale of this rule is that grouping triple patterns could potentially miss solution mappings that are only produced by joining data from two different sources. The remaining rules are introduced to handle the following exceptions. Rule III handles exclusive groups: triple patterns of exclusive groups can be grouped into a single subexpression without negatively impacting the answers completeness. Finally, Rule IV handles the following cases: if the decomposition only has a single subexpression that is evaluated at a single source, it does not have an impact on the completeness how these triple patterns are grouped into subexpressions of the decomposition. In contrast to Rule III, in Rule IV even though $SE$ is just evaluated at a single source, $SE$ does not need to be an exclusive group and can have other relevant sources that are not in $S$. By these rules, we can measure the density of a decomposition graph relative to the maximum number of possible edges as given by the atomic decomposition graph.

*Definition 4.8 (Density of a Query Decomposition).* Given a query decomposition $D(P, F)$ and the corresponding graph $G_{D(P,F)} = (V, E)$, the density $density(D(P, F))$ is computed as:

$$density(D(P, F)) = \frac{|E|}{|E^*|} \in [0, 1].$$

THEOREM 4.9. *The evaluation of a query decomposition $D(P, F)$ over a federation $F$ yields complete answers, if $density(D(P, F)) = 1$:*

$$density(D(P, F)) = 1 \implies \theta_{D(P,F)}(P) = [\![P]\!]_F \quad (11)$$

PROOF. We prove the implication in Eq. (11) by contradiction. We assume $density(D(P, F)) = 1$ and $\theta_{D(P,F)}(P) \neq [\![P]\!]_F$. According to Def. 4.8, $density(D(P, F)) = 1$ holds only if the decomposition graph of $D(P, F)$ has the same number of edges as the decomposition graph of the atomic decomposition: $|E| = |E^*|$. Following Def. 4.4 and Def. 4.7, in $E^*$ there is an edge between each triple pattern and its relevant sources (Rule I) and an edge between every pair of triple patterns (Rule II). The maximum number of edges is

$$|E^*| = \underbrace{\sum_{tp_i \in P} r(tp_i)}_{\text{Rule I}} + \underbrace{0.5 \cdot |P| \cdot (|P| - 1)}_{\text{Rule II}}.$$

Since we prove completeness, we focus on the case when a decomposition yields fewer answers: $\theta_{D(P,F)}(P) \subset [\![P]\!]_F$. This can occur in two cases:

---

[4]We prove this equality by contradiction. Consider $G = \bigcup_{c \in C} ep(c)$. Assume that there exists a solution mapping $\mu$ s.t. $\mu \in ([\![tp_i]\!]_{ep(c_1)} \cup \cdots \cup [\![tp_i]\!]_{ep(c_k)})$ and $\mu \notin [\![tp_i]\!]_G$. This means that the evaluation of a subexpression over some source, e.g. $[\![tp_i]\!]_{ep(c_j)}$, is producing additional answers w.r.t. the evaluation over the union of all RDF graphs. This could only happen if $ep(c_j) \nsubseteq G$, however, this contradicts the definition of $G$. Now assume that $\mu \in [\![tp_i]\!]_G$ but $\mu \notin ([\![tp_i]\!]_{ep(c_1)} \cup \cdots \cup [\![tp_i]\!]_{ep(c_k)})$. Without loss of generality, assume that $\mu$ was produced from matching an RDF triple $t \in G$ s.t. $t \notin ep(c)$ for all services $c \in C$ in the federation. This is again a contradiction with the definition of $G$.

CASE 1: A part of the query is not evaluated at a relevant source. Without loss of generality, consider that a triple pattern $tp_i \in P$ is not evaluated at a relevant source $c_j$ and $[\![tp_i]\!]_{c_j}$ contributes to the answers of $P$. In this case, the decomposition graph $G_{D(P,F)}$ does not have an edge $(tp_i, c_j)$ according to Rule I and, therefore, $|E| < |E^*|$. This contradicts the assumption that $density(D(P,F)) = 1$.

CASE 2: Triple patterns with several relevant sources are grouped into subexpressions. Consider the solution mapping $\mu \in [\![P]\!]_F$, with $\mu = \{\mu_1 \cup \mu_2 \mid \mu_1 \in [\![tp_1]\!]_{c_1} \wedge \mu_2 \in [\![tp_2]\!]_{c_2}, \mu_1 \sim \mu_2\}$, and without loss of generality assume that $r(tp_1) = r(tp_2) = \{c_1, c_2\}$. Such a solution mapping $\mu$ does not exist in $\theta_{D(P,F)}(P)$ in the case that the two triple patterns are evaluated jointly at the source $c_1$ and $c_2$, that is

$$((tp_1 \text{ AND } tp_2), \{c_1, c_2\}) \in D(P,F)$$

In this case, the edge $(tp_1, tp_2)$ does not exist in $E$ according to Rule II but the edge exists in $E^*$ because

$$(tp_1, \{c_1, c_2\}), (tp_2, \{c_1, c_2\}) \in D^*(P,F)$$

Therefore, we have $|E| < |E^*|$ which contradicts the assumption that $density(D(P,F)) = 1$.

$\square$

We can prove that a decomposition density of 1 implies answer completeness, however, the inverse (i.e., $density(D(P,F)) = 1 \impliedby \theta_{D(P,F)}(P) = [\![P]\!]_F$ ) cannot be guaranteed. For example, there might be a triple pattern $tp_1$ with two relevant sources $c_1$ and $c_2$ with just source $c_1$ contributing to the final answers. A decomposition $D(P,F)$ where $tp_1$ is not evaluated at $c_2$ might still yield complete answers but $density(D(P,F)) < 1$ according to Rule I. Therefore, the decomposition density is a measure for the *expected* completeness based on the assumptions that answer completeness is negatively affected while i) contacting fewer relevant sources, and ii) grouping triple patterns that can be evaluated at several sources into subexpressions. Estimating the true completeness more accurately would require additional information on the data provided by the LDF services than just the relevant sources. Such additional information could be used to improve the effectiveness of our measure, for example by weighting the edges in the decomposition graph according to their importance. However, such an extension is out of the scope of this work.

*Example 4.10.* Let us consider the BGP $P = (tp_1 \text{ AND } tp_2 \text{ AND } tp_3 \text{ AND } tp_4)$ from the SPARQL query of the motivating example in Section 2 and the federation $F_{ex} = (\{c_1, c_2\}, int, ep)$. The relevant sources are $r(tp_1) = \{c_1\}$, $r(tp_2) = \{c_1\}$, $r(tp_3) = \{c_1, c_2\}$, $r(tp_4) = \{c_2\}$. The atomic query decomposition is $D^*(P,F) = \{(tp_1, \{c_1\}), (tp_2, \{c_1\}), (tp_3, \{c_1, c_2\}), (tp_4, \{c_2\})\}$ and the corresponding graph is shown in Fig. 2a. In $P$, the triple patterns $tp_1$ and $tp_2$ form an exclusive group, as they are both only answerable by service $c_1$. Therefore, we can combine them in a single subexpression without reducing the expected completeness in $D_1(P,F) = \{((tp_1 \text{ AND } tp_2), \{c_1\}), (tp_3, \{c_1, c_2\}), (tp_4, \{c_1\})\}$. The corresponding graph shown in Fig. 2b is identical to $G_{D^*(P,G)}$ and thus its expected completeness is $density(D_1(P,F)) = \frac{11}{11} = 1$. Alternatively, we can choose to evaluate $tp_3$ only at service $c_2$ with $D_2(P,F) =$

$\{((tp_1 \text{ AND } tp_2), \{c_1\}) \, ((tp_3 \text{ AND } tp_4), \{c_2\})\}$ (Fig. 2c) or evaluate $tp_3$ at service $c_1$ with $D_3(P,F) = \{((tp_1 \text{ AND } tp_2 \text{ AND } tp_3), \{c_1\}), (tp_4, \{c_2\})\}$ (Fig. 2d). Since both corresponding decomposition graphs have fewer edges than the graph of $G_{D^*(P,F)}$, we expect fewer answers because: $density(D^*(P,F)) > density(D_2(P,F)) > density(D_3(P,F))$.

*Query Decomposition Cost.* The example illustrates how query decompositions can have different levels of expected completeness. Ideally, one would always choose the atomic query decomposition to guarantee complete answers. However, there are also costs associated with the evaluation of a decomposition that are induced by the amount of transferred data for intermediate results during query execution as well as the number of services that need to be contacted. In federations of SPARQL endpoints, both goals are achieved by i) decomposing the query into as few subexpressions as possible, and ii) reducing the number of endpoints contacted by selecting just those sources which are likely to contribute to the final answer of the query. In contrast, when facing heterogeneous federations of LDF services, the languages of the LDF services need to be considered as well. The reason is that the interface-compliant evaluation might yield additional costs in cases when subexpressions cannot be evaluated by a service as a whole. There might be several interface-compliant evaluations for an expression because the original expression could be split in different ways into subexpressions that can be evaluated by the service. We denote an interface-compliant evaluation of an expression $P$ with the minimal number of subexpressions as $\theta_c^*(P)$, which is the evaluation of $P$ that requires separating the expression into the fewest subexpression to be interface-compliant. Note that $|\theta_c^*(P)| = 1$, if $P \in L_c$.

We propose a lower bound for query decomposition cost that considers the number of services contacted and the number of subexpressions in an interface-compliant evaluation of the decomposition. In particular, this lower bound combines: (1) The number of sources $|S|$ to be contacted per subexpression. (2) The number of additional subexpressions ($|\theta_c^*(SE)| - 1$) required for an interface–compliant evaluation for each subexpression and all corresponding sources.

*Definition 4.11 (Cost of a Query Decomposition).* The cost of evaluating a query decomposition $D(P,F)$ is given by

$$cost(D(P,F)) = \sum_{(SE,S) \in D(P,F)} |S| + \sum_{(SE,S) \in D(P,F) \wedge \forall c \in S} (|\theta_c^*(SE)| - 1).$$

Note that the proposed query decomposition cost provides a lower bound for evaluating a decomposition while computing the exact cost requires knowledge about the technical configurations of the services in the federation. For instance, obtaining solutions from TPF servers might require several requests for paginating the results, while a single request might suffice on a SPARQL endpoint.

*Example 4.12.* Let us consider the decomposition $D_2(P,F)$ from Example 4.10 and the subexpression $SE_1 = (tp_1 \text{ AND } tp_2 \text{ AND } tp_3)$ to be evaluated at source $S_1 = \{c_1\}$. In contrast to its density, the cost of evaluating $D_2(P,F)$ depends on the LDF interface $c_1$ implements. If $c_1$ is a SPARQL endpoint, i.e. $int(c_1) = (L_{\text{CORESPARQL}}, M_{\text{EP}})$, the evaluation $[\![SE_1]\!]_{c_1}$ is interface-compliant and thus $|\theta_{c_1}^*(SE_1)| = 1$. However, if $c_1$ is a TPF server, i.e. $int(c_1) = (L_{\text{TP}}, M_{\text{TPF}})$, the interface-compliant evaluation of $SE_1$ would require evaluating

$D^*(P,F) = \{(tp_1,\{c_1\}),(tp_2,\{c_1\}),(tp_3,\{c_1,c_2\}),(tp_4,\{c_2\})\}$

(a) Graph $G_{D^*(P,F)}$

$D_1(P,F) = \{((tp_1 \text{ AND } tp_2),c_1),(tp_3,\{c_1,c_2\}),(tp_4,\{c_2\})\}$

(b) Graph $G_{D_1(P,F)}$

$D_1(P,F) = \{((tp_1 \text{ AND } tp_2),c_1),(tp_3,\{c_1,c_2\}),(tp_4,\{c_2\})\}$

(c) Graph $G_{D_2(P,F)}$

$D_1(P,F) = \{((tp_1 \text{ AND } tp_2),c_1),(tp_3,\{c_1,c_2\}),(tp_4,\{c_2\})\}$
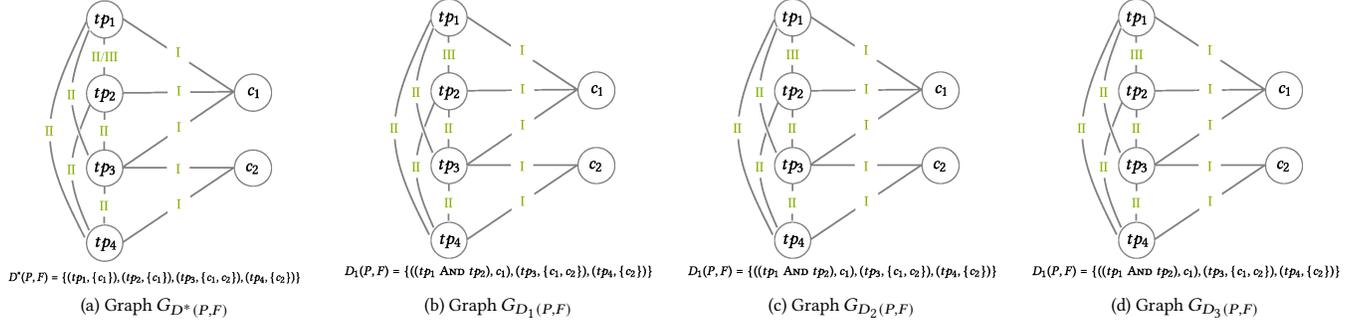
(d) Graph $G_{D_3(P,F)}$

**Figure 2: Query decomposition graphs for decompositions from Example 4.10. The rules for adding edges are indicated in green.**

the triple patterns individually with $\theta^*_{c_1}(SE_1) = [\![tp_1]\!]_{c_1} \bowtie [\![tp_2]\!]_{c_1} \bowtie [\![tp_3]\!]_{c_1}$ and thus $|\theta^*_{c_1}(SE_1)| = 3$. Hence, the evaluation at the TPF server requires two additional subexpressions to be evaluated. This may lead to higher execution costs as there are potentially more intermediate results to be transferred and the service needs to be contacted at least two additional times.

Finally, we can combine both the density and cost of a query decomposition into the query decomposition problem which aims to obtain a query decomposition that maximizes the expected answer completeness while minimizing the execution cost.

*Definition 4.13 (Query Decomposition Problem).* Given a BGP $P$ and a federation $F = (C, int, ep)$, the query decomposition problem is finding a query decomposition $D(P,F)$ that minimizes the execution cost while maximizing its density:

$$D(P,F) = \arg\max density(D(P,F)) \wedge \arg\min cost(D(P,F))$$

Note that this problem is a multi-objective optimization problem, where there might not be a single best solution but rather a set of optimal trade-off solutions, i.e. Pareto-optimal solutions.

*Example 4.14.* Consider two alternative example federations which differ in the LDF interfaces of their services $c_1$ and $c_2$:

$F_1 = (\{c_1, c_2\}, int, ep)$ with $int(c_1) = int(c_2) = (L_{\text{CORESPARQL}}, M_{\text{EP}})$.
$F_2 = (\{c_1, c_2\}, int, ep)$ with $int(c_1) = (L_{\text{TP}}, M_{\text{TPF}})$ and $int(c_2) = (L_{\text{CORESPARQL}}, M_{\text{EP}})$.

The density and cost for the query decompositions are given in the following table, where the best values are indicated in bold.

| | $F_1$ | | | | $F_2$ | | | |
|---|---|---|---|---|---|---|---|---|
| | $D^*$ | $D_1$ | $D_2$ | $D_3$ | $D^*$ | $D_1$ | $D_2$ | $D_3$ |
| $\sum \lvert S\rvert$ | 5 | 4 | 2 | 2 | 5 | 4 | 2 | 2 |
| $\sum(\lvert\theta^*_c(SQ)\rvert - 1)$ | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 2 |
| *cost* | 5 | 4 | **2** | **2** | 5 | 5 | **3** | 4 |
| *comp* | **1** | **1** | $\frac{9}{11}$ | $\frac{8}{11}$ | **1** | **1** | $\frac{9}{11}$ | $\frac{8}{11}$ |

The decomposition cost in the example shows how both the number of subexpressions and the number of sources they are evaluated at ($\sum \lvert S\rvert$) as well as the capabilities of the interface $\sum(\lvert\theta^*_c(SE)\rvert-$

1) have an impact on the overall cost. Further, it shows the trade-off between the two conflicting objectives density and cost. In both federations, the decompositions that yield the highest density also have the highest cost and vice versa. Approaches to solving the query decomposition problem need to determine solutions that yield a suitable (depending on the use case) trade-off between the number of answers and execution cost. According to Def. 4.8, two main factors impact on the density. First, the triple patterns should be evaluated at as many relevant sources as possible (Rule I). Second, the more fine-grained the subexpressions for triple patterns that have several relevant sources in common, the higher the density (Rule II). Similarly, the costs of decompositions originate from two main aspects. First, contacting fewer sources with larger subexpressions will reduce costs and, second, decomposing the query into subexpressions that are interface-compliant will reduce the cost. One way of pruning sources without affecting answer completeness is to determine the relevant sources that do not contribute to the final answers of the query [20]. However, this can be very challenging for queries with triple patterns that contain terms from common ontologies (e.g., RDF/S, OWL), as they can be answered by many of the sources in the federation. For this purpose, some approaches rely on pre-computed statistics/catalogues [17, 20, 21] and/or the query capabilities of SPARQL endpoints, such as Ask queries [28]. We propose a query decomposition approach that can be combined with a heuristic-based source pruning method and can be applied for any heterogeneous federation.

**Query Decomposition Approach.** We propose an approach that does not rely on specific statistics about the members and has two central goals: (1) maximize the density by evaluating all triple patterns at the relevant sources that contribute to the final answers, and (2) reduce the execution cost by obtaining subexpressions that leverage the capabilities of the services as much as possible. Furthermore, we add an optional source pruning step to further decrease cost by reducing the number of sources contacted. The decomposer is outlined in Algorithm 1. Its inputs are a BGP $P$ and a federation $F = (C, int, ep)$. First, the algorithm creates the atomic decomposition by iterating over each triple patterns $tp$ in $P$, determines the set of relevant sources as $S$, and adds $(tp, S)$ to the decomposition $D$ (Line 2 - Line 5). Next, the relevant sources per triple pattern can be pruned in Line 6. This pruning step is not

(a) Join Ordering and Union Expressions

(b) Interface-compliant subquery plans

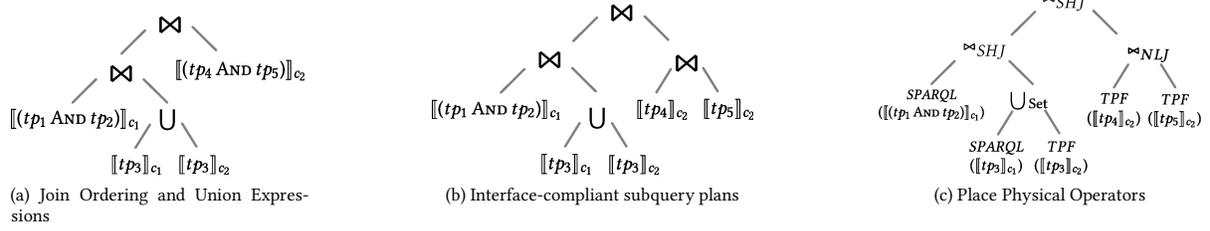(c) Place Physical Operators

**Figure 3: Query planning steps for the query decomposition** $D_4(P, F_{ex}) = \{((tp_1 \text{ AND } tp_2), \{c_1\}), (tp_3, \{c_1, c_2\}), ((tp_4 \text{ AND } tp_5), \{c_2\})\}$.

---

**Algorithm 1:** Interface-aware Query Decomposer

**Input:** BGP P = $\{tp_1, \ldots, tp_n\}$, Federation $F = (C, int, ep)$
1  $D = \emptyset$
2  **for** $tp \in P$ **do**
3     $S = \text{relevantSources}(tp)$
4     $D = D \cup \{(tp, S)\}$
5  **end**
6  $D = \text{pruneSources}(D)$
7  **do**
8     $updated = False$
9     **for** $\forall (SE_i, S_i), (SE_j, S_j) \in D \land SE_i \neq SE_j$ **do**
10       **if** $|vars(SE_i) \cap vars(SE_j)| > 0$
11       $\land |S_i \cup S_j| = 1$
12       $\land (SE_i \text{ AND } SE_j) \in L_c, \forall c \in S_i$ **then**
13          $D = D \setminus \{(SE_i, S_i), (SE_j, S_j)\}$
14          $D = D \cup \{((SE_i \text{ AND } SE_j), S_i)\}$
15          $updated = True$
16          **break**
17    **end**
18 **while** $updated$
19 **return** $D$

---

required, however, it allows for reducing the decomposition cost by i) reducing the number of sources to be contacted, and ii) allowing to group more triple patterns into subexpressions in the following steps. The source pruning approach is interchangeable and we detail our source pruning heuristic in the next paragraph. After pruning the sources, the algorithm tries to merge as many subexpressions in the decomposition $D$ as possible. All possible combinations of subexpressions $(SE_i, S_i)$ and $(SE_j, S_j)$ are considered and merged if they fulfill the following three conditions:

Condition I    Both subexpressions have variables in common: $|vars(SE_i) \cap vars(SE_j)| > 0$. (Line 10)

Condition II    Both subexpressions have exactly one source in common: $|S_i \cup S_j| = 1$. (Line 11)

Condition III    The common source $c$ can evaluate the conjunction of both expressions: $(SE_i \text{ AND } SE_j) \in L_c$. (Line 12)

If two subexpressions fulfill all conditions, the individual subexpressions are removed from the decomposition $D$ and their conjunction is added to $D$. This process is repeated until no more subexpressions can be merged ($updated = False$). A central property of the query decomposition generated by the algorithm is the fact that the evaluation of all subexpressions is compliant with all corresponding sources. That is, $\forall (SE, S) \in D(P, F) : \forall c \in S : \theta_c(SE) = [\![SE]\!]_c$. As a result, the interface-compliant evaluation (Def. 4.3) of all decomposition generate by Algorithm 1 is given as

$$\theta_{D(P,F)}(P) := \bowtie_{(SE_i, S_i) \in D(P,F)} (\cup_{c_j \in S_i} [\![SE_i]\!]_{c_j}).$$

Note that this property does not require the query planner to find the subexpression minimizing evaluation $\theta_c^*(SE_i)$.

**Source Pruning Approach.** We propose a heuristic that leverages the atomic decomposition graph $G_{D^*(P,F)} = (V^*, E^*)$ and does not rely on data statistics. Our approach iterates over the source vertices $c_i \in V^*$ by non-increasing out-degree (i.e. starting with the most *popular* source). For each triple pattern $tp_j$ connected to $c_i$ ($(c_i, tp_j) \in E^*$), the edges to all other sources are removed for $tp_j$: $E^* = E^* \setminus \{(c_k, tp_j) \in E^* \mid \forall c_k \neq c_i\}$. In addition, the relevant sources for triple patterns with the same common subject are not pruned to maximize completeness. The rationale for this is the observation that RDF datasets typically follow entity-centric descriptions, where the URI of an entity appears in the subject of triples in the authoritative dataset. For example, triples with subject dbr:Berlin are all part of the DBpedia dataset.

## 4.2 Query Planner

The main tasks of the query planner are finding an efficient logical plan and placing physical operators such that the execution time of the query plan is minimized. For both tasks, common cost-based query planners leverage statistics on the data of the members in the federation. In heterogeneous federations, however, the query planning approaches cannot always rely on the same level of statistics from all sources and need to be adjusted to the statistics available at the individual sources. For instance, obtaining fine-grained statistics might require access to the entire dataset of a source for efficient computation [14] or require the services to be able to execute complex SPARQL expressions, such as aggregate queries. Furthermore, in the case that the interface language of an LDF service does not support the evaluation of a subexpression from the decomposition, the planner needs to obtain an efficient subplan for evaluating the subexpression over that service. In this section, we first discuss the steps necessary to obtain efficient query plans, and thereafter, we propose a query planner for query decompositions that respects the interface restrictions in heterogeneous federations.

*Join Ordering with Union Expressions.* The query planner determines a join ordering for the subexpressions in a decomposition that minimizes the number of intermediate results. The challenge lies in estimating the size of intermediate results from subexpressions and joins. This is particularly difficult in heterogeneous interfaces due to two factors. First, the methods to estimate cardinalities depend on the interface languages and the metadata supported by the interfaces. For example, determining the cardinality of a subexpression comprised of two triple patterns could be achieved

by a COUNT query if the interface, e.g. a SPARQL endpoint, supports the evaluation of such expressions. However, this could be an expensive operation on the server and thus time-consuming for the client. Moreover for other interfaces, such as TPF servers, this would not be possible and the cardinality would need to be estimated according to the metadata of the triple patterns. If available, statistical data on the data distribution could be used alternatively to estimate the number of intermediate results [8, 21]. Second, federated plans comprise union operators to combine data from alternative relevant sources. In this case, estimating the number of intermediate results from a union operator that will contribute to a join is more difficult due to the different data distributions in each source. Therefore, the planner must devise appropriate join orderings in the presence of unions from different sources. Fig. 3a shows a join ordering with unions for a query decomposition from the query and federation of our motivating example: $D_4(P, F_{ex}) = \{((tp_1 \text{ AND } tp_2), \{c_1\}), (tp_3, \{c_1, c_2\}), ((tp_4 \text{ AND } tp_5), \{c_2\})\}$.

*Interface-compliant Subexpression Plans.* If a decomposer does not provide decompositions in which the subexpressions $SE$ are interface-compliant, the query planner additionally needs to find subplans that evaluate the subexpression in an interface-compliant manner. In those cases, the query planner needs to break down $SE$ into subexpressions that minimize the cost of the interface-compliant evaluation $\theta_c^*(SE)$. Since the resulting interface-compliant evaluation consists of several joins, the query planner also needs to determine the join ordering for $\theta_c^*(SE)$. For example, if the service is a TPF server, this would require first splitting the subexpression into its individual triple patterns and thereafter, finding an appropriate join ordering. The latter could rely on existing query planning approaches for TPF servers [3, 27]. Fig. 3b shows the interface-compliant evaluation for $[\![(tp_4 \text{ AND } tp_5)]\!]_{c_2}$ over the DBpedia TPF server ($c_2$) for decomposition $D_4(P, F_{ex})$. The evaluation is given by $\theta_{c_2}^*(tp_4 \text{ AND } tp_5) = [\![tp_4]\!]_{c_2} \bowtie [\![tp_5]\!]_{c_2}$ and it introduces an additional join operation in the query plan.

*Placing Physical Operators.* Finally, the query planner selects physical operators to obtain an executable physical query plan. This includes placing access operators that retrieve the solution mappings from the services as well as physical join and union operators to process the intermediate results. The access operators transform the subexpressions into requests that can be processed by the corresponding LDF services. Ideally, the access operators leverage the querying capabilities of the interfaces such that the results are obtained efficiently. For example, traditional federated query engines for SPARQL endpoints require only access operators that adhere to the SPARQL protocol to get solution mappings from the endpoints. In heterogeneous federations, however, appropriate access operators for each LDF interface in the federation need to be implemented and placed accordingly by the planner. Moreover, physical join operators that implement different join strategies, such as symmetric hash join or bind join, need to be placed effectively as they incur different costs. Finally, the planner needs to place the appropriate physical union operators in the plan that respects the semantics of the query language. Fig. 3c shows an example of a physical plan for decomposition $D_4(P, F_{ex})$, where service $c_1$ is a Wikidata SPARQL endpoint and service $c_2$ a the DBpedia TPF server.

---

**Algorithm 2:** Query Planning Algorithm

**Input:** Decomposition $D(P, F) = \{(SE_1, S), \ldots, (SE_n, S_n)\}$

1  List $L$
2  **for** $(SE_i, S_i) \in D(P, F)$ **do**
3  | $card_i = \text{estimateCardinality}(SE_i, S_i)$
4  | $L.append((SE_i, S_i, card_i))$
5  **end**
6  $L = \text{sort}(L, card_i)$ // Sort $L$ by non-decreasing $card_i$
7  $d_1 = L.get(0)$
8  $L.remove(0)$
9  $T_1 = \text{AccessPlan}(d_1)$
10 **while** $|L| > 0$ **do**
11 | $d_2 = L.get(0)$
12 | **for** $i = 1; i < |L|; i + +$ **do**
13 | | $(SE_i, S_i, card_i) = L.get(i)$
14 | | **if** $|vars(T_1) \cap vars(SE_i)| > 0$ **then**
15 | | | $d_2 = (SE_i, S_i, card_i)$
16 | | | $L.remove(i)$
17 | | | **break**
18 | **end**
19 | $T_2 = \text{AccessPlan}(d_2)$
20 | $O = \text{getPhysicalOperator}(T_1, T_2)$
21 | $T_1 = \text{JoinPlan}(T_1, T_2, O)$
22 **end**
23 **return** $T_1$

---

**Query Planning Approach.** We now present a heuristic-based query planner for heterogeneous federations. In particular, it relies on decomposition obtained by Algorithm 1. First, we present the overall planning approach and, thereafter, we present details of our prototypical implementation. The query planner is outlined in Algorithm 2. It starts by estimating the cardinality of each subexpression in the decomposition (Line 3) and creates a list $L$ in which the subexpressions are sorted by non-decreasing cardinality (Line 6). The query planner starts building the query plan with the subexpression $d_1$ with the lowest cardinality and creates the corresponding access plan $T_1$ (Line 9)[5]. It iterates over the remaining subexpressions in $L$ and determines the next subexpression to join $T_1$ with. This is either a remaining subexpression with the lowest cardinality and a common variable (Line 15) or if there is no join remaining in the BGP, it is the subexpression with the lowest cardinality (Line 11). Once the subexpression $d_2$ is selected, the access plan $T_2$ for $d_2$ is created (Line 19) and the appropriate physical join operator $O$ is determined (Line 20). Finally, $T_1$ becomes the JoinPlan of $T_1$ and $T_2$ (Line 21). When $L$ is empty, the final plan $T_1$ is returned (Line 23). After presenting the generic planning approach, we now provide details on the specific steps in our prototypical implementation. The current implementation focuses on the three well-known LDF interfaces: Triple Pattern Fragments (TPF), Bindings-Restricted Triple Pattern Fragments (brTPF), and SPARQL endpoints. Further, it relies on the properties of decompositions generated by our interface-aware query decomposer presented in Algorithm 1. That is, each subexpression $SE_i$ is interface-compliant for all sources in $S_i$. For each service $c \in S_i$, estimateCardinality (Line 3) obtains the estimated cardinality $card_i^c$ for the subexpression $SE_i$ at the service $c$ in line with the interface language and the metadata of $c$. As evaluating $SE_i$ at several sources reflects a union operation, it then sums up those individual cardinalities to obtain

---

[5]The access plan for $d_i = (SE_i, S_i)$ refers to the union of evaluating subexpression $SE_i$ at each source in $S_i$.

the total cardinality of $SE_i$ at all sources: $card_i = \sum_{c \in S_i} card_i^c$. If $SE_i$ is a triple pattern and the source is a brTPF or a TPF server, we request the triple pattern and use the $\texttt{void:count}$ in the meta-data as the cardinality estimation. If $SE_i$ is a BGP or a triple pattern and the source is a SPARQL endpoint, we use a COUNT query to estimate the cardinality. Further, we estimate the join cardinality of two subexpressions $SE_i$ and $SE_j$ as the minimum of their cardinalities. Next, we implement appropriate access operators for all three interfaces. Since all subexpressions are compliant with the interface, we do not need to first obtain an interface-compliant evaluation in the $\texttt{AccessPlans}$. Finally, we determine the physical join operator according to the estimated number of requests to execute the join. We distinguish between two different common join strategies: symmetric hash join and bind join. The reason to use the number of requests to determine the join strategy is two-fold: i) the number of requests directly have an effect on the execution time, and ii) fewer requests lead to a reduced load on the services in the federation. Thus, we compare the number of requests necessary when placing a bind join or a symmetric hash join and choose the operator that yields fewer requests. The request estimations depend on the implementation of the physical join operator, which we detail in the following section.

## 4.3 Physical Operators

The heterogeneity of LDF interfaces in a federation introduces challenges but also opens opportunities for implementing novel physical operators. Access operators to retrieve answers from LDF services need to be implemented in efficient ways reducing the load on the LDF services and the time for obtaining results to improve query execution time. For instance, TPF servers have a *page size* configuration that limits the number of answers that are returned upon a requested triple pattern. Additionally, many public SPARQL endpoints are configured with fair use policies that can lead to zero or incomplete query results [25]. Consequently, implementations of access operators for SPARQL endpoints should not overload the SPARQL endpoints and adhere to the usage policies. Yet, physical join operators can be designed to simultaneously handle different LDF interfaces and follow different join strategies depending on the capabilities of the underlying services. We call these kinds of operators *polymorphic* and present a novel Polymorphic Bind Join tailored to TPF, brTPF, and SPARQL interfaces.

**Polymorphic Bind Join.** The Polymorphic Bind Join (PBJ) implements a Nested Loop Join algorithm that is able to adjust its join strategy according to the LDF interface. It simultaneously executes a tuple- and block-based nested loop join according to the supported interface language. Our current implementation supports the languages $L_{\text{TP}}$, $L_{\text{TP+VALUES}}$ and $L_{\text{CORESPARQL}}$. By leveraging the capabilities of each service, PBJ reduces the number of requests when accessing more capable sources using the block-based approach. In particular, PBJ is designed for cases where the inner relation is either an access operator or the union of access operators. For each LDF interface $f$, a block size $B_f$ is defined. During the execution, the operator keeps a reservoir per service that is filled by tuples from the outer relation. When the reservoir reaches the block size $B_f$ of the corresponding LDF interface, the bindings from the

reservoir are requested at the services. For example, when querying a TPF server in a nested loop join, each solution mapping of the outer relation is used to instantiate and resolve the triple pattern of the inner relation, hence, $B_{\text{TPF}} = 1$. However, as the interface languages of brTPF servers and SPARQL endpoints support SPARQL values expressions, the PBJ changes its operation accordingly by requesting a triple pattern or a subexpression with several bindings. The number of bindings that can be sent to a brTPF server $B_{\text{BRTPF}}$ depends on the server configuration [12]. For SPARQL endpoints, $B_{\text{EP}}$ is not limited, yet too many values may lead to long runtimes at the endpoint and potentially incomplete results.[6]

The proposed query planner selects a Symmetric Hash Join (SHJ) or Polymorphic Bind Join (PBJ) in $\texttt{getPhysicalOperator}$ (Line 20) depending on the estimated number of requests. The number of requests to execute the SHJ or PBJ depends on the sub-plans $T_1$ and $T_2$. If $T_1$ is an $\texttt{AccessPlan}$, the number of requests to obtain the tuples of $T_1$ are determined by its cardinality $card_{T_1}$ and the interfaces over which $T_1$ is evaluated. Otherwise, if $T_1$ is a $\texttt{JoinPlan}$, no additional requests are necessary to obtain the tuples for $T_1$. For the first case, the requests $R_{acc}(T_1)$ depend on the maximum number of tuples that can be obtained per requests from the corresponding LDF service, which we denote as $Max_{\text{EP}}$, $Max_{\text{BRTPF}}$, and $Max_{\text{TPF}}$.[7]

$$R_{acc}(T) = \sum_{\substack{c \in S \wedge \\ int(c)=\text{EP}}} \left\lceil \frac{card_T^c}{Max_{\text{EP}}} \right\rceil + \sum_{\substack{c \in S \wedge \\ int(c)=\text{BRTPF}}} \left\lceil \frac{card_T^c}{Max_{\text{BRTPF}}} \right\rceil + \sum_{\substack{c \in S \wedge \\ int(c)=\text{TPF}}} \left\lceil \frac{card_T^c}{Max_{\text{TPF}}} \right\rceil$$

As a result, we can compute the number of request for the SHJ as the sum of the requests for the two sub-plans:

$$R_{SHJ}(T_1, T_2) = R_{acc}(T_1) + R_{acc}(T_2)$$

For the PBJ, we need to determine the number of requests that need to be performed in the inner relation $R_{bind}(T_1, T_2)$ , which depends on the cardinality $card_{T_1}$ of the outer relation $T_1$ and the block sizes for the services in the inner relation:

$$R_{bind}(T_1, T_2) = \sum_{\substack{c \in S_2 \wedge \\ int(c)=\text{EP}}} \left\lceil \frac{card_{T_1}}{B_{\text{EP}}} \right\rceil + \sum_{\substack{c \in S_2 \wedge \\ int(c)=\text{BRTPF}}} \left\lceil \frac{card_{T_1}}{B_{\text{BRTPF}}} \right\rceil + \sum_{\substack{c \in S_2 \wedge \\ int(c)=\text{TPF}}} \left\lceil \frac{card_{T_1}}{B_{\text{TPF}}} \right\rceil$$

The overall number of requests for the PBJ is

$$R_{PBJ}(T_1, T_2) = R_{acc}(T_1) + R_{bind}(T_1, T_2).$$

## 5 EXPERIMENTAL EVALUATION

We evaluate a prototypical implementation of the interface-compliant query decomposer, query planner, and polymorphic bind join. The goal is to investigate the impact of the components on the performance when querying heterogeneous federations of LDF interfaces.

**Datasets and Queries.** We use the well-known FedBench benchmark [22] which is comprised of 9 datasets and tailored to assess the performance of federated SPARQL querying strategies. We use a total of 25 queries from Cross Domain (CD1-7), Life Science (LS1-7) and Linked Data (LD1-11) in our evaluation.

---

[6]In our implementation, we set $B_{\text{BRTPF}} = 30$ and $B_{\text{EP}} = 50$, to reduce the requests while not overloading the endpoint.
[7]In our implementation, we set $Max_{\text{BRTPF}} = 100$ [12] and $Max_{\text{TPF}} = 100$ [27], and $Max_{\text{EP}} = 10000$ (most common value reported at https://sparqles.ai.wu.ac.at/).
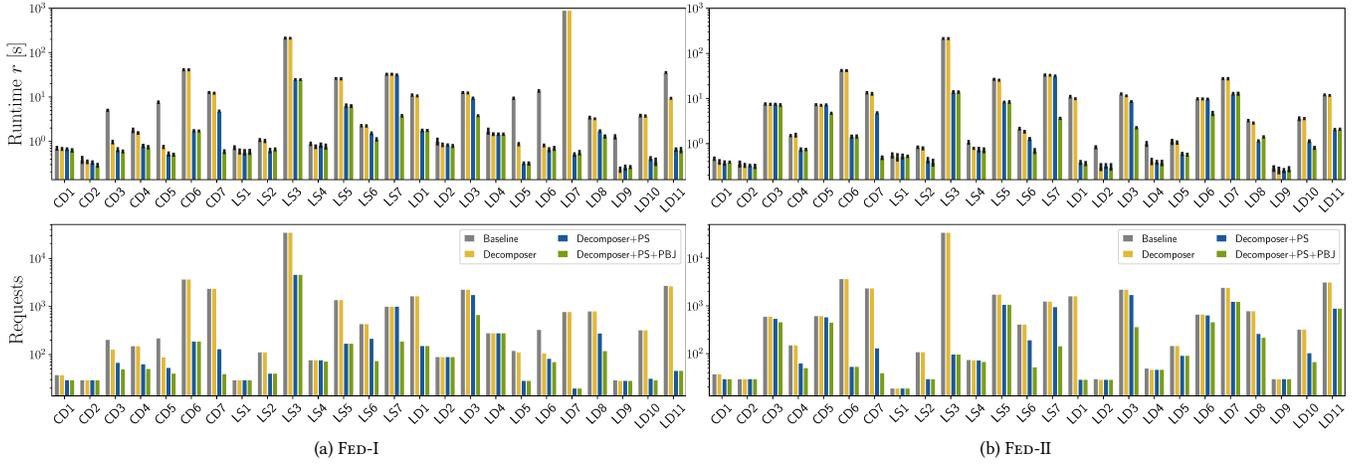
**Figure 4: Average runtimes [s] (log-scale) and total number of requests (log-scale) for each query and both federations.**

**Table 1: Heterogeneous Federations: Fed-I, Fed-II. SPARQL endpoints indicated in bold and brTPF servers in italic.**

|          | DBpedia | NYTimes | LinkedMDB | Jamendo | GeoNames | SWDF | KEGG | Drugbank | ChEBI |
|----------|---------|---------|-----------|---------|----------|------|------|----------|-------|
| Fed-I  | **SPARQL** | *brTPF* | *brTPF* | TPF | **SPARQL** | TPF | *brTPF* | TPF | **SPARQL** |
| Fed-II | TPF | *brTPF* | *brTPF* | **SPARQL** | TPF | **SPARQL** | *brTPF* | **SPARQL** | TPF |

**Federations.** We evaluate our approach on two heterogeneous federations Fed-I and Fed-II shown in Table 1 to study the performance in different scenarios. The central difference between the federations is that in Fed-I the three largest datasets are accessible via SPARQL endpoints while in Fed-II they are accessible via TPF servers. The other datasets are accessible via TPF or brTPF servers.

**Implementation.** We implemented a prototypical federated query engine for heterogeneous federations that implements the proposed query planner, decomposer, source pruning (PS), and polymorphic bind join (PBJ) operator. Our implementation is based on CROP [?] and implemented in Python 2.7.13. As Baseline, we use execute the query plans from our query planner for the atomic decompositions. The decomposer, source pruning, and PBJ are disabled in the Baseline. Note that, while Comunica [26] can query heterogeneous interfaces, its performance is currently not competitive as it does not implement query decomposition, source pruning, or polymorphic join operators. Therefore, we do not consider Comunica in our evaluation. We use the `Server.js` v2.2.3[8] and original Java brTPF server implementation [12] to deploy the TPF and brTPF servers with HDT [?] backends. We used Virtuoso v07.20.3229 with the default `virtuoso.ini` (cf. supplemental material). All LDF services and the client were executed on a single Debian Jessie server (2x16 core Intel(R) Xeon(R) E5-2670 2.60GHz CPU; 256GB RAM) to avoid network latency. The timeout was set to 900 seconds. After a warm-up run, the queries were executed five times. The source code, experimental results, and additional material are provided in the supplemental material of this submission.

---

[8]https://github.com/LinkedDataFragments/Server.js

**Table 2: Average total runtime $\sum r$, number of requests $\sum req.$, and answers $\sum ans.$ per run as well as the mean decomposition completeness $\overline{comp}$ and decomposition cost $\overline{cost}$[9].**

|        |                   | $\sum r$ | $\sum req.$ | $\sum ans.$ | $\overline{comp}$ | $\overline{cost}$ |
|--------|-------------------|----------|-------------|-------------|-------------------|-------------------|
| Fed-I  | Baseline          | 1337.87  | 54452       | 13534       | **1.0**           | 1.0               |
|        | Decomposer        | 1274.15  | 53958       | 13534       | **1.0**           | 0.95              |
|        | Decomposer+PS     | 93.66    | 9645        | 13171       | 0.77              | **0.55**          |
|        | Decomposer+PS+PBJ | 54.69    | 7271        | 13171       | 0.77              | **0.55**          |
| Fed-II | Baseline          | 433.37   | 57671       | 13578       | **1.0**           | 1.0               |
|        | Decomposer        | 425.15   | 57662       | 13578       | **1.0**           | 0.9               |
|        | Decomposer+PS     | 116.45   | 9040        | 13171       | 0.77              | **0.53**          |
|        | Decomposer+PS+PBJ | 69.38    | 6121        | 13171       | 0.77              | **0.53**          |

**Metrics.** We evaluated the performance by the following metrics: (i) *Runtime*: Elapsed time spent by the engine evaluating a query. (ii) *Number of Requests*: Total number of requests submitted to the LDF services during the query execution. (iii) *Number of Answers*: Total number of answers produced. (iv) *Dieffiency*: Continuous efficiency as the answers are produced over time [5].

## 5.1 Experimental Results

We start providing an overview of the performance of the different components. In Fig. 4a and Fig. 4b the mean runtimes and number of requests are shown per query for Fed-I and Fed-II. The values are also summarized in Table 2. Considering the impact of the individual components, the results show that enabling the decomposer without pruning the sources and no PBJ (Decomposer), only provides a slight improvement in the runtime over the Baseline, even though all queries yield the same number of requests or less. This is because, without source pruning, only exclusive groups can be merged by the decomposer. The results when adding the source pruning approach (Decomposer+PS) show that pruning sources considerably reduces both the runtime and the number of requests for the majority of queries. The reasons for the improvement are two-fold: i) the decomposer can create more and larger

subexpressions, and ii) fewer services are contacted during the execution of the query plan. Finally, with the polymorphic join operator (Decomposer+PS+PBJ), we observe the lowest overall runtimes and number of requests in both federations. In Fed-I, executing all queries with Decomposer+PS+PBJ is more than 34 times faster than the Baseline and 6 times faster in Fed-II. The results show that our interface-aware federated query approaches, that adjust to the specifics of heterogeneous interfaces, can greatly improve the performance in terms of runtime. Simultaneously, it reduces the load on the servers by requiring fewer requests. The results show that the interfaces present in the federation (Fed-I vs. Fed-II) substantially impact the querying performance when not considering the interfaces' capabilities (Baseline). Yet, our interface-aware solution (Decomposer+PS+PBJ) enables similar performance results regardless of the interfaces.

**Query Decomposition.** The results show the effectiveness of the proposed *density* measure as a proxy for completeness and *cost* measures as means to assess the expected execution cost of query decompositions. In Table 2, we can observe that, in both federations, the decomposer without source pruning yields complete answers with $\overline{density}$ = 1.0, since only exclusive groups are merged (Rule III). The cost can only be slightly reduced (Fed-I: $\overline{cost}$ = 0.95 and Fed-II: $\overline{cost}$ = 0.9)[9]. However, adding the source pruning (Decomposer+PS) enables decompositions with about half the cost. Contacting fewer services reduces the cost but also leads to a reduction in the expected completeness ($\overline{density}$ = 0.77) and to fewer answers ($\sum ans$) that are being produced. 97% of all answers are still produced when sources are pruned.[10] These results show that the improvement achieved by the decomposer in its ability to leverage the interfaces' capabilities depends on the source pruning.

**Polymorphic Bind Join.** The results in both Fig. 4 and Table 2 reveal that, in the two federations, adding the Polymorphic Bind Join (Decomposer+PS+PBJ) reduces the number of requests by more than 25% and, as a consequence, reduces the overall runtimes. We investigate the dieefficiency to better understand the impact of the PBJ. In Fig. 5, we show the dieefficiency plots for four example queries. The plot for LS3 in Fig. 5a shows that the performance of the PBJ is similar to a regular NLJ in case it cannot leverage the capabilities of the interfaces, e.g., LS3 where only TPFs are contacted. However, if the capabilities of the services can be leveraged, the PBJ allows for producing the results at a higher rate as shown for queries LD3 (Fig. 5b) and LS6 (Fig. 5c). In the latter, all answers are produced at once. Nonetheless, the semi-blocking nature of the PBJ can also have a detrimental effect on dieefficiency and runtime as observed for query LS8 in Fig. 5d. Here, the production of the answers is delayed because the block size of the inner relation (which consumes data from a brTPF server) is not reached until all results of the outer relation in the PBJ are produced. Future work could study an adaptive PBJ with variable block sizes, determined according to the expected number of tuples of the outer relation.

Summarizing our experimental evaluation, the results show the effectiveness of our interface-aware techniques for query decomposition, planning, and physical operators. Furthermore, the results illustrate that our techniques can cope with heterogeneous

federations that are composed of different combinations of interfaces.

**Limitations.** The central assumption of our framework is the access to high-level information about the federation (e.g., interfaces and relevant sources) while fine-grained statistics (e.g., data distributions) are not available. Therefore, the proposed framework components are limited to devise approximate solutions to the problems of query decomposition, planning, and execution. While our experimental results show substantial improvements in the Fed-Bench benchmark, these improvements might not hold in other federations. Yet, our framework is a foundation for federated query processing in heterogeneous federations and the components can be refined in the case that additional statistics are available. For instance, by weighting edges in the decomposition graph according to probability of sources contributing to the answers of a query.

# 6 RELATED WORK

Query processing over homogeneous federations of SPARQL endpoints has been broadly studied and existing approaches address different challenges. For instance, [1, 24] leverage requests during runtime to obtain efficient query plans, while [8, 11, 17, 19, 20] implement cost models that rely on pre-computed statistics, and [4] focuses on runtime adaptivity. Furthermore, approaches that specifically study query decomposition have been proposed. Vidal et al. [28] propose a formalization of the query decomposition problem in a way such that it can be mapped to the vertex coloring problem. Vidal et al. [28] propose the heuristic *Fed-DSATUR* to solve the problem. Similar, Endris et al. [10] formalize the query decomposition problem for federated SPARQL querying and present a decomposition approach that relies on RDF Molecule Templates, which represent metadata obtained by executing SPARQL queries over endpoints. Different from our work, these approaches assume all members in the federation to be SPARQL endpoints, and thus, the proposed solutions rely on their querying capabilities.

Additional Linked Data Fragment (LDF) interfaces and corresponding SPARQL clients have been proposed. They range from less expressive interfaces, such as (Bindings-Restricted) Triple Pattern Fragments ([12]) [27], to more expressive interfaces such as SaGe [15] or smart-KG [7]. To study the expressiveness of LDF interfaces, Hartig et al. [13] propose the Linked Data Fragment Machines as a formal framework that includes client demand, server demand, and communication cost when executing queries over these interfaces. Similar to their work, we also formalize the concept of a server language to distinguish the capabilities of different interfaces in the federation. Yet, our work goes beyond individual interfaces and studies the problem of heterogeneous LDF federations.

Lastly, a few approaches have addressed the problem of heterogeneous interfaces. Comunica [26] is a client able to query heterogeneous LDF federations. But, in contrast to our work, Comunica does not support interface-aware query decomposition and handles the query execution on a triple pattern level, even if different interfaces are present. Moreover, the physical join operators, such as the nested loop join, do not adapt to the different interfaces. In a

---

[9]Cost values are normalized: $cost(D(P, F))/cost(D^*(P, F))$

[10]In Fed-I, the Baseline does not yield all answers, due to a timeout in LD7.

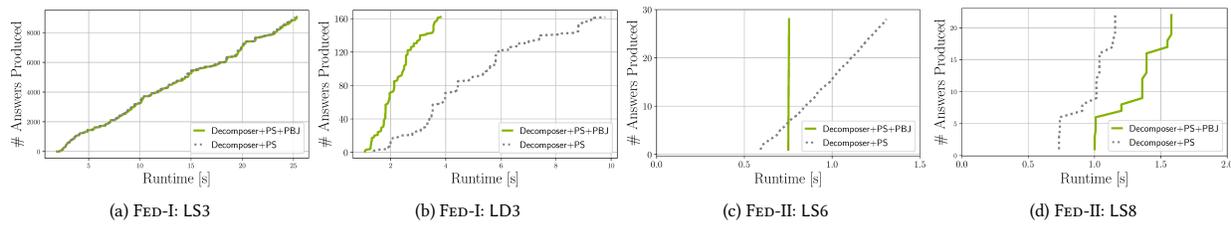(a) FED-I: LS3      (b) FED-I: LD3      (c) FED-II: LS6      (d) FED-II: LS8

**Figure 5: Example dieffiency plots for the approach with the PBJ (green) and without the PBJ (dotted).**

recent paper, Cheng and Hartig [9] study query plans in heterogeneous federations. Similar to our work, they conceptualize different interfaces, and federation members implementing those interfaces. They focus on a formal language for logical query plans over such federations but, in contrast to our work, they do not propose specific solutions to devise such plans and derive physical plans to be evaluated by an engine. Montoya et al. [16] propose a client to query replicas of datasets via heterogeneous interfaces (brTPF server and SPARQL endpoints) to exploit their characteristics. Different from our work, they focus on different interfaces for single datasets but do no investigate federated querying.

## 7 CONCLUSION AND FUTURE WORK

We formalize the concept of federations of Linked Data Fragment services and present the challenges that querying approaches over heterogeneous federations face. In particular, we present a theoretical framework and practical solutions for query decomposition, query planning, and physical operators tailored to heterogeneous LDF federations. In our experimental study, we evaluated a prototypical implementation of our proposed solutions. The results show a substantial improvement in performance achieved by devising interface-aware strategies to exploit the capabilities of TPF, brTPF, and SPARQL endpoints during federated query processing. Future work may focus on extending the proposed framework to other LDF interfaces and studying how state-of-the-art query decomposition, planning, and source pruning approaches from federated SPARQL engines can be applied to heterogeneous federations.

## REFERENCES

[1] Ibrahim Abdelaziz, Essam Mansour, Mourad Ouzzani, Ashraf Aboulnaga, and Panos Kalnis. 2017. Lusail: A System for Querying Linked Data at Scale. *Proc. VLDB Endow.* 11, 4 (2017), 485–498. https://doi.org/10.1145/3186728.3164144

[2] Maribel Acosta, Olaf Hartig, and Juan F. Sequeda. 2019. Federated RDF Query Processing. In *Encyclopedia of Big Data Technologies*. https://doi.org/10.1007/978-3-319-63962-8_228-1

[3] Maribel Acosta and Maria-Esther Vidal. 2015. Networks of Linked Data Eddies: An Adaptive Web Query Processing Engine for RDF Data. In *The Semantic Web - ISWC 2015 - 14th International Semantic Web Conference, Bethlehem, PA, USA, October 11-15, 2015, Proceedings, Part I*. 111–127. https://doi.org/10.1007/978-3-319-25007-6_7

[4] Maribel Acosta, Maria-Esther Vidal, Tomas Lampo, Julio Castillo, and Edna Ruckhaus. 2011. ANAPSID: An Adaptive Query Processing Engine for SPARQL Endpoints. In *The Semantic Web - ISWC 2011 - 10th International Semantic Web Conference, Bonn, Germany, October 23-27, 2011, Proceedings, Part I*. 18–34. https://doi.org/10.1007/978-3-642-25073-6_2

[5] Maribel Acosta, Maria-Esther Vidal, and York Sure-Vetter. 2017. Dieffiency Metrics: Measuring the Continuous Efficiency of Query Processing Approaches. In *The Semantic Web - ISWC 2017 - 16th International Semantic Web Conference, Vienna, Austria, October 21-25, 2017, Proceedings, Part II*. 3–19. https://doi.org/10.1007/978-3-319-68204-4_1

[6] Carlos Buil Aranda, Marcelo Arenas, and Óscar Corcho. 2011. Semantics and Optimization of the SPARQL 1.1 Federation Extension. In *The Semantic Web: Research and Applications - 8th Extended Semantic Web Conference, ESWC 2011, Heraklion, Crete, Greece, May 29 - June 2, 2011, Proceedings, Part II*. 1–15. https://doi.org/10.1007/978-3-642-21064-8_1

[7] Amr Azzam, Javier D. Fernández, Maribel Acosta, Martin Beno, and Axel Polleres. 2020. SMART-KG: Hybrid Shipping for SPARQL Querying on the Web. In *WWW '20: The Web Conference 2020, Taipei, Taiwan, April 20-24, 2020*. 984–994. https://doi.org/10.1145/3366423.3380177

[8] Angelos Charalambidis, Antonis Troumpoukis, and Stasinos Konstantopoulos. 2015. SemaGrow: optimizing federated SPARQL queries. In *Proceedings of the 11th International Conference on Semantic Systems, SEMANTICS 2015, Vienna, Austria, September 15-17, 2015*. 121–128. https://doi.org/10.1145/2814864.2814886

[9] Sijin Cheng and Olaf Hartig. 2020. FedQPL: A Language for Logical Query Plans over Heterogeneous Federations of RDF Data Sources (Extended Version). arXiv:2010.01190 [cs.DB]

[10] Kemele M. Endris, Mikhail Galkin, Ioanna Lytra, Mohamed Nadjib Mami, Maria-Esther Vidal, and Sören Auer. 2017. MULDER: Querying the Linked Data Web by Bridging RDF Molecule Templates. In *Database and Expert Systems Applications - 28th International Conference, DEXA 2017, Lyon, France, August 28-31, 2017, Proceedings, Part I*. 3–18. https://doi.org/10.1007/978-3-319-64468-4_1

[11] Olaf Görlitz and Steffen Staab. 2011. SPLENDID: SPARQL Endpoint Federation Exploiting VOID Descriptions. In *Proceedings of the Second International Workshop on Consuming Linked Data (COLD2011), Bonn, Germany, October 23, 2011*. http://ceur-ws.org/Vol-782/GoerlitzAndStaab_COLD2011.pdf

[12] Olaf Hartig and Carlos Buil Aranda. 2016. Bindings-Restricted Triple Pattern Fragments. In *On the Move to Meaningful Internet Systems: OTM 2016 Conferences - Confederated International Conferences: CoopIS, C&TC, and ODBASE 2016, Rhodes, Greece, October 24-28, 2016, Proceedings*. 762–779. https://doi.org/10.1007/978-3-319-48472-3_48

[13] Olaf Hartig, Ian Letter, and Jorge Pérez. 2017. A Formal Framework for Comparing Linked Data Fragments. In *The Semantic Web - ISWC 2017 - 16th International Semantic Web Conference, Vienna, Austria, October 21-25, 2017, Proceedings, Part I*. 364–382. https://doi.org/10.1007/978-3-319-68288-4_22

[14] Lars Heling and Maribel Acosta. 2020. Estimating Characteristic Sets for RDF Dataset Profiles Based on Sampling. In *The Semantic Web - 17th International Conference, ESWC 2020, Heraklion, Crete, Greece, May 31-June 4, 2020, Proceedings*. 157–175. https://doi.org/10.1007/978-3-030-49461-2_10

[15] Thomas Minier, Hala Skaf-Molli, and Pascal Molli. 2019. SaGe: Web Preemption for Public SPARQL Query Services. In *The World Wide Web Conference, WWW 2019, San Francisco, CA, USA, May 13-17, 2019*. 1268–1278. https://doi.org/10.1145/3308558.3313652

[16] Gabriela Montoya, Christian Aebeloe, and Katja Hose. 2018. Towards Efficient Query Processing over Heterogeneous RDF Interfaces. In *Emerging Topics in Semantic Technologies - ISWC 2018 Satellite Events [best papers from 13 of the workshops co-located with the ISWC 2018 conference]*. 39–53. https://doi.org/10.3233/978-1-61499-894-5-39

[17] Gabriela Montoya, Hala Skaf-Molli, and Katja Hose. 2017. The Odyssey Approach for Optimizing Federated SPARQL Queries. In *The Semantic Web - ISWC 2017 - 16th International Semantic Web Conference, Vienna, Austria, October 21-25, 2017, Proceedings, Part I*. 471–489. https://doi.org/10.1007/978-3-319-68288-4_28

[18] Jorge Pérez, Marcelo Arenas, and Claudio Gutiérrez. 2009. Semantics and complexity of SPARQL. *ACM Trans. Database Syst.* 34, 3 (2009), 16:1–16:45. https://doi.org/10.1145/1567274.1567278

[19] Bastian Quilitz and Ulf Leser. 2008. Querying Distributed RDF Data Sources with SPARQL. In *The Semantic Web: Research and Applications, 5th European Semantic Web Conference, ESWC 2008, Tenerife, Canary Islands, Spain, June 1-5, 2008, Proceedings*. 524–538. https://doi.org/10.1007/978-3-540-68234-9_39

[20] Muhammad Saleem and Axel-Cyrille Ngonga Ngomo. 2014. HiBISCuS: Hypergraph-Based Source Selection for SPARQL Endpoint Federation. In *The Semantic Web: Trends and Challenges - 11th International Conference,*

*ESWC 2014, Anissaras, Crete, Greece, May 25-29, 2014. Proceedings.* 176–191. https://doi.org/10.1007/978-3-319-07443-6_13

[21] Muhammad Saleem, Alexander Potocki, Tommaso Soru, Olaf Hartig, and Axel-Cyrille Ngonga Ngomo. 2018. CostFed: Cost-Based Query Optimization for SPARQL Endpoint Federation. In *Proceedings of the 14th International Conference on Semantic Systems, SEMANTICS 2018, Vienna, Austria, September 10-13, 2018.* 163–174. https://doi.org/10.1016/j.procs.2018.09.016

[22] Michael Schmidt, Olaf Görlitz, Peter Haase, Günter Ladwig, Andreas Schwarte, and Thanh Tran. 2011. FedBench: A Benchmark Suite for Federated Semantic Data Query Processing. In *The Semantic Web - ISWC 2011 - 10th International Semantic Web Conference, Bonn, Germany, October 23-27, 2011, Proceedings, Part I.* 585–600. https://doi.org/10.1007/978-3-642-25073-6_37

[23] Michael Schmidt, Michael Meier, and Georg Lausen. 2010. Foundations of SPARQL query optimization. In *Database Theory - ICDT 2010, 13th International Conference, Lausanne, Switzerland, March 23-25, 2010, Proceedings.* 4–33. https://doi.org/10.1145/1804669.1804675

[24] Andreas Schwarte, Peter Haase, Katja Hose, Ralf Schenkel, and Michael Schmidt. 2011. FedX: Optimization Techniques for Federated Query Processing on Linked Data. In *The Semantic Web - ISWC 2011 - 10th International Semantic Web Conference, Bonn, Germany, October 23-27, 2011, Proceedings, Part I.* 601–616.

https://doi.org/10.1007/978-3-642-25073-6_38

[25] Arnaud Soulet and Fabian M. Suchanek. 2019. Anytime Large-Scale Analytics of Linked Open Data. In *The Semantic Web - ISWC 2019 - 18th International Semantic Web Conference, Auckland, New Zealand, October 26-30, 2019, Proceedings, Part I.* 576–592. https://doi.org/10.1007/978-3-030-30793-6_33

[26] Ruben Taelman, Joachim Van Herwegen, Miel Vander Sande, and Ruben Verborgh. 2018. Comunica: A Modular SPARQL Query Engine for the Web. In *The Semantic Web - ISWC 2018 - 17th International Semantic Web Conference, Monterey, CA, USA, October 8-12, 2018, Proceedings, Part II.* 239–255. https://doi.org/10.1007/978-3-030-00668-6_15

[27] Ruben Verborgh, Miel Vander Sande, Olaf Hartig, Joachim Van Herwegen, Laurens De Vocht, Ben De Meester, Gerald Haesendonck, and Pieter Colpaert. 2016. Triple Pattern Fragments: A low-cost knowledge graph interface for the Web. *J. Web Semant.* 37-38 (2016), 184–206. https://doi.org/10.1016/j.websem.2016.03.003

[28] Maria-Esther Vidal, Simón Castillo, Maribel Acosta, Gabriela Montoya, and Guillermo Palma. 2016. On the Selection of SPARQL Endpoints to Efficiently Execute Federated SPARQL Queries. *Trans. Large Scale Data Knowl. Centered Syst.* 25 (2016), 109–149. https://doi.org/10.1007/978-3-662-49534-6_4