

Towards Hybrid Link Traversal: Challenges and Research Directions for Heterogeneous Access Interfaces

Ruben Eschauzier¹, Ruben Taelman¹

¹Department of Electronics and Information Systems, Ghent University – imec

Abstract

Some cool abstract goes here!

Keywords

Link Traversal-based Query Processing, Hybrid Link Traversal, SPARQL, Dataspaces, Query Optimization

This paper addresses (Issue #2) of the W3C Dataspaces Community Group.

1. Introduction

While centralizing data into warehouses or lakes can benefit query engine performance, it is antithetical to the core principles of *dataspaces*. Dataspaces promote a decentralized ecosystem where data remains at the source, managed by independent participants under agreed-upon governance models. This architectural shift introduces significant technical challenges for query engines, which must now discover and query data across a vast network of sources, each enforcing its own usage policies.

Centralized querying requires collecting large volumes of proprietary or sensitive data in a single source, which conflicts with the requirements for fine-grained sovereignty and minimized data replication. Conversely, traditional *federated* approaches [1, 2, 3] support decentralization but typically assume a static federation of a small number (10–100) of uniform, expressive endpoints (e.g., SPARQL endpoints) [4]. Enforcing complex usage policies [5] on such heavyweight interfaces significantly impacts performance [6], and these engines struggle to scale to the thousands of sources characteristic of a dataspace.

In contrast, realistic dataspace environments involve a massive number of permissioned sources exposing data through highly *heterogeneous interfaces* [7, 8]. Rather than a uniform layer of SPARQL endpoints, a dataspace participant may expose data via generic HTTP documents, constrained APIs, derived views citevan2024granular, or specialized query services. Due to the scale and autonomy of the dataspace, the complete set of relevant sources and their specific interface capabilities is possibly unknown during query time. Such an environment requires a *hybrid* federation approach capable of querying across diverse data interfaces , discovered at runtime, without prior centralization.

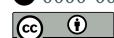
Link Traversal-based Query Processing (LTQP) is an approach that meets these needs by discovering data sources during execution via hypermedia links found in earlier results. Starting from seed references (e.g., a self-description or catalog entry), it follows links asynchronously. While LTQP is traditionally defined as an approach for querying over federated *hypermedia* documents by following links at runtime, data within a dataspace is often exposed through diverse access interfaces that may not support standard hypermedia traversal and querying.

We propose *Hybrid LTQP*, a generalization of link traversal that supports the diverse access methods found in dataspaces—including SPARQL endpoints, materialized views, and Triple Pattern Fragments (TPF) [9]. To leverage these varied levels of expressivity, a hybrid engine must offload computation to the server whenever possible. For instance, a SPARQL endpoint can execute complex sub-queries locally. By leveraging internal indexing and local data knowledge, the endpoint can produce results far more

The Third International Workshop on Semantics in Dataspaces, co-located with the Extended Semantic Web Conference, June 01, 2025, Portorož, Slovenia

✉ ruben.eschauzier@ugent.be (R. Eschauzier); ruben.taelman@ugent.be (R. Taelman)

>ID 0000-0002-6475-806X (R. Eschauzier); 0000-0001-5118-256X (R. Taelman)



© 2025 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

efficiently than a client-side LTQP engine. Offloading these tasks to capable endpoints significantly reduces overall query latency. However, building such a system introduces several unresolved challenges. In this paper, we review the state-of-the-art in hybrid querying, analyze related work, and outline open issues for future research.

2. Related Work

While Hybrid LTQP has previously been discussed in literature, the exact definition of hybrid link traversal querying varies between papers.

2.1. Hybrid Traversal - Local vs Remote

A branch of hybrid link traversal querying uses precomputed indexes, either located within the decentralized environment or computed locally to speed-up query execution or source discovery. The first of such approaches, uses precomputed indexes to rerank sources based on relevance to the query's triple pattern and joins, and retrieves them to execute queries over. [10]. Other works [11] propose to cache information on sources to help in prioritization of sources.

Another approach is to store entire RDF graphs in local stores to quickly execute parts of the query. For this approach, non-blocking indexed-based operators [12] are used to leverage the locally computed indexes for faster query execution. A major limitation is the freshness of the data in the store, which can quickly go stale for high velocity data sources. Various methods, like freshness (or coherence)-aware query processing exist [13], however it remains a limitation.

2.2. Hybrid Traversal - Closed Systems

Other approaches integrate into the dataspace itself. ESPRESSO [14] is such an approach, it is a framework around the Solid dataspace specifically that defines several apps in the dataspace to perform certain *optimization* tasks. In the Solid environment, data is stored into personal data vaults. These vaults use the Linked Data Platform [15] to arrange the data into a folder-like structure. To avoid traversing the entire pod to query it, ESPRESSO uses the *Brewmaster* application to locally create indexes of the pod and uses the *CoffeeFilter* application to search its contents. In addition, ESPRESSO uses an overlay network to distribute end-user queries to relevant data resources across Solid servers, where each Solid server is mapped to a federated database node in the overlay network.

Unlike approaches based on link traversal or service discovery within Solid, ESPRESSO relies on an explicitly configured overlay network, with no mechanism for dynamically discovering participating servers or query endpoints. Furthermore, currently ESPRESSO only supports distributed keyword-based search, it does not support SPARQL queries over the sources.

While these approaches do hybrid link traversal, they significantly deviate from our interpretation of hybrid traversal. This in contrast to the setting described in this paper where the discovered sources themselves are heterogeneous, serendipitously discovered, and cannot be assumed to be present for every source within the dataspace.

2.3. Hybrid Link Traversal over Heterogeneous Sources

State-of-the-art LTQP engines, most notably Comunica [16, 17], demonstrate hybrid capabilities by dynamically federating over discovered SPARQL endpoints, TPF interfaces, and hypermedia documents. However, these systems are currently limited by a *lowest common denominator* strategy: they fail to leverage the full expressivity of sophisticated interfaces. For instance, upon discovering a SPARQL endpoint, existing engines restrict interaction to simple *triple pattern* lookups.

This forces the client to retrieve high volumes of intermediate results and perform joins locally, effectively negating the optimization and latency benefits of server-side execution.

3. Challenges with Hybrid Link Traversal

In this section, we will outline specific identified problems with hybrid link traversal making implementation difficult.

3.1. Establishing Exclusive Groups

In classical federated query processing, *exclusive groups* are sets of triple patterns within a query that can *only* be answered by a single source [1]. Instead of retrieving results for individual triple patterns and joining them locally, the engine can delegate the entire group as a single sub-query (e.g., a Basic Graph Pattern) to the remote endpoint. This technique significantly reduces the number of remote requests and processed intermediate results [18, 1].

The identification of exclusive groups is a strict requirement for this sub-query delegation. If a query engine cannot verify that a group of patterns belongs *exclusively* to a single source, it cannot safely dispatch them as a conjunctive sub-query (BGP). Doing so risks *incomplete results* because the endpoint will execute a server-side join. If the endpoint lacks data for one of the patterns, or if valid join partners reside on a different server, the server-side join execution will filter out bindings that could have been successfully joined with data from another source. Thus, the exclusive grouping strategy fails to produce the valid distributed join results.

Federation engines like *FedX* [1], *HiBISCuS* [19], and *FedUp* [18] rely on either prior knowledge on the content of the federation members, or issue queries (such as ASK) to determine these groups.

For Hybrid Link Traversal to use SPARQL endpoints to compute joins locally, the engine should also identify exclusive groups. The problem is the unbounded nature of LTQP [20]. As during LTQP the query engine does not have access to all sources it will query over, exclusive groups can never be established. Even if we move to a finite web [20], the engine would have to first traverse all reachable [20] sources before exclusive groups can be determined. This is antithetical to the streaming nature of LTQP engines [16, 21] and will delay time to first result unacceptably. Thus, traditional approaches for determining exclusive groups are not applicable.

3.1.1. Research Avenues

Discoverable Indexes: Indexes can be effectively used to catalogue the data sources present in a dataspace. These indexes could be deployed as services, similar to approaches like ESPRESSO [14], pointing towards available source URIs and potentially exposing summary statistics. Using such indexes to quickly determine relevant sources, LTQP engines can apply traditional federation techniques to establish exclusive groups and perform source selection.

Caching: Similar to discoverable indexes, caches provide prior knowledge about the data sources present in a dataspace. Provided the cache maintains a degree of freshness and completeness, the query engine can perform a ‘simulated traversal’ offline over the cached data to identify relevant sources. Any gaps in this knowledge can then be supplemented by live traversal prior to or during execution.

Authoritativeness Assumptions: Establishing exclusive groups in a decentralized dataspace requires a notion of *source authority* [22]. Under standard Link Traversal, any reachable source can assert triples about any subject (e.g., a third-party source asserting $\langle P1 \rangle \text{ foaf : knows } \langle P2 \rangle$), creating a “wild west” of data where no single source can be deemed the exclusive authority for a topic. This lack of authority prevents the query engine from determining if a set of patterns can be safely delegated to a single endpoint.

A possible solution is restricting the scope of validity such that a data source is the *sole authority* for triples where the subject corresponds to the source’s own URI (or a URI within its controlled namespace). By enforcing this *Subject Authority Constraint*, the query engine can statically determine that all triples matching $\langle \text{SourceA} \rangle \ ?p \ ?o$ reside exclusively at *Source A*. This allows the engine to safely treat such

patterns as an exclusive group, enabling the delegation of complex sub-queries (e.g., BGPs) to hybrid sources like SPARQL endpoints without the risk of incomplete results.

3.2. Integrating Dynamically Delegated Sub-queries into the Query Plan

When exclusive groups are identified and subqueries are dynamically pushed to expressive server-side interfaces, the resulting intermediate results must fit into the existing query plan.

In the traditional optimize-then-execute approach, used by Comunica-link-traversal [16, 17] and SQUIN [21] (except in [23]), the query plan is constructed before execution and evaluated over dynamically discovered sources.

Thus, an exclusive group can only be integrated if the precomputed plan already produces an intermediate result that exactly matches the triple patterns in that group.

Example 3.1. As an example, consider a query plan with the join structure

$$(T_1 \bowtie T_2) \bowtie (T_3 \bowtie T_4).$$

If the engine discovers an exclusive group consisting of the triple patterns $\{T_2, T_3\}$, this group cannot be integrated into the plan. The plan never produces an intermediate result of the form $T_2 \bowtie T_3$, since these patterns occur in different subtrees. As a result, the exclusive group does not fit into the existing join structure and cannot be pushed as a single subquery.

3.2.1. Research Avenues

A clear path forward is the use of adaptive query processing techniques. These techniques should enable low-cost plan switching without discarding intermediate results, as the large number of sources in a dataspace would otherwise lead to plan thrashing. Algorithms such as Eddies [24], SteMs [25] (as used in [23]), and STAIRs [26] are viable candidates due to their tuple-level adaptivity.

Despite their potential, the literature lacks a comprehensive analysis of adaptive frameworks within the context of LTQP. Specifically, there is little guidance on designing routing strategies, the logic determining how tuples are forwarded between operators, or the overhead costs associated with each.

Eddies offer a constrained form of adaptivity. In this model, the join operators are instantiated at the start of the query; the engine only adapts the *order* in which tuples visit these operators. Because these join operators (e.g., Symmetric Hash Joins) accumulate internal state as they process data, the engine is effectively locked into the existing operator instances.

Example 3.2. Consider a query with triple patterns $\{T_1, T_2, T_3, T_4\}$. If an Eddy starts by routing results from $T_1 \bowtie T_2$, the join state is physically stored within that specific operator's hash tables. If the engine later discovers that $T_2 \bowtie T_3$ is more selective due to newly traversed data, it can change the routing sequence, but it cannot easily dissolve the $T_1 \bowtie T_2$ state to form a $T_2 \bowtie T_3$ join without significant re-initialization overhead.

In contrast, **SteMs** (State Modules) decouple state from the join operation entirely. Since SteMs do not hold fixed join states, any tuple can be routed to any SteM in any order. However, this flexibility introduces a *recomputation tax*: if a SteM is probed before its join partners have arrived, the intermediate result is not captured and must be re-derived once more data is discovered via traversal.

Finally, **STAIRs** (Storage, Transformation and Access for Intermediate Results) attempt to bridge this gap through state migration. This allows a query engine to physically move join states between different operator instances. While this maximizes flexibility, state migration is a resource-intensive operation. Determining the optimal threshold for migration, and the number of required migrations to incorporate the diverse precomputed sub-queries discovered during traversal is an open challenge.

3.3. Service Discovery in Hybrid Query Environments

A fundamental challenge in Hybrid Link Traversal Querying is the gap between resource identifiers and the aggregated services that host them. For standard LTQP, the reachability of data is determined by the recursive dereferencing of URIs, adhering to the “follow-your-nose” principle where a URI identifies a resource and provides its description. However, SPARQL endpoints function as aggregators: they contain data about resources but do not necessarily share the resource’s URI namespace, nor are they always explicitly linked from the resource descriptions themselves. Consequently, a query engine may possess a URI for a resource, but lack the URI pointing towards the SPARQL endpoint that contains that resource, rendering the data unreachable despite being publicly accessible.

3.3.1. Research Avenues

Explicit Linkage Mechanisms An avenue for mitigating the service discovery gap is by using explicit links between resources and their hosting endpoints. This approach relies on extending the “follow-your-nose” principle to include service metadata within the resource description itself. Existing methods such as the Vocabulary of Interlinked Datasets (VOID) can embed pointers directly in the RDF data. For instance, a triple taking the form $\langle r, \text{void:sparqlEndpoint}, s \rangle$ allows a traversal engine to immediately identify that the resource r is queryable via the endpoint s . However, this relies on adoption of this approach while, for example, currently only 33% percent of endpoints expose VOID descriptions [27].

Service Registries and Indexing. Where direct linkage is absent, research directs towards the usage of external service registries and catalogs. Unlike linkage-based discovery, this approach treats the mapping between URI namespaces and SPARQL endpoints as a separate knowledge base. Within a dataspace, services such as ESPRESSO [14] can aggregate these connections, effectively bridging the gap between identifiers and query services. However, relying on purely centralized indexes is antithetical to the goals of a decentralized dataspace. Consequently, future directions should investigate distributed alternatives, such as Peer-to-Peer (P2P) discovery mechanisms, to prevent the reliance on single points of failure.

3.4. Alternative View Deduplication

In decentralized dataspaces, a single dataset can often be accessed through multiple architectural views. Consider a personal data store (PDS) following the Solid protocol [28]. Here, data is structured as a hierarchical resource tree where documents are interconnected via RDF predicates. A Link Traversal Query Processing (LTQP) engine navigates this store by dereferencing these URIs on-the-fly.

Data providers may also expose a SPARQL endpoint over the same PDS to improve query efficiency. Ideally, an LTQP engine should incorporate this endpoint into its query plan. However, because endpoint discovery is itself a byproduct of the traversal process, an engine may dereference several hypermedia documents before it encounters the metadata pointing to the endpoint.

This creates a synchronization problem: the engine ingests the same data from two different interfaces, leading to duplicated results and wasted computational resources. In the specific context of Solid, a client might assume that any (public) URI sharing the Pod’s root namespace is covered by the discovered endpoint. However, in general dataspace environments, this assumption is often unsubstantiated; alternative views may only cover specific subsets of data, such as historical archives versus high-frequency data.

3.4.1. Research Avenues

View Overlap Detection To determine an appropriate deduplication strategy, a query engine must first assess the data coverage of the discovered access interfaces. One approach from the literature involves issuing queries to both the local data store and the remote endpoint to calculate *coherence* [13],

a metric indicating the extent to which the two sources agree. However, this method is restricted to queryable interfaces and incurs significant overhead, as each coherence check requires additional HTTP requests. An alternative is to rely on server-side metadata that explicitly denotes which underlying resources a given access interface aggregates. Unfortunately, this relies heavily on the widespread adoption of specific metadata vocabularies by data providers, which cannot easily be guaranteed in a decentralized, autonomous dataspace.

Client-Side Deduplication Regardless of whether the data coverage of discovered access interfaces is known, the query engine must perform deduplication on the client side. This requires tracking the *provenance* (i.e., the source) of bindings during query processing to filter out duplicates based on the computed or assumed coverage of alternative interfaces. A naive approach is to maintain a “seen” set of previously produced bindings for each triple pattern, discarding any newly discovered matches that already exist in the set. However, retaining a complete history of bindings is highly memory-intensive and scales poorly for large queries. While probabilistic data structures like Bloom filters can drastically reduce this memory footprint, their inherent false-positive rates risk accidentally discarding unique, valid results.

If the engine knows which sources are covered by an alternative access interface, it can attempt further optimizations. First, upon discovering the interface, any relevant links that have been discovered but not yet dereferenced can be pruned from the traversal queue. Second, the engine can better manage the lifecycle of the “seen” sets. Once the engine can guarantee that all data capable of producing duplicate results has been processed, the associated “seen” sets can be safely discarded to free memory. However, the mechanism for reliably detecting this point during execution remains an open research area.

4. Conclusion

Hybrid querying over heterogeneous sources has the potential to significantly reduce query latency and network overhead, by combining the dynamic discovery of link traversal with the computational efficiency of expressive interfaces like SPARQL endpoints,

However, as outlined in this paper, realizing a fully functional Hybrid LTQP engine introduces substantial technical hurdles. The unbounded nature of the web makes it difficult to establish exclusive groups for safe sub-query delegation without relying on prior knowledge, discoverable indexes, or strict authoritativeness assumptions. Furthermore, query engines must adopt highly adaptive processing techniques to seamlessly integrate dynamically discovered materialized sub-queries without thrashing the query plan or losing intermediate results. We also highlight the ongoing need for robust service discovery mechanisms to bridge the gap between resource identifiers and the endpoints aggregating them. Finally, we identify a challenge with alternative views and the inherent risk of result duplication due to them and identify two research directions for this challenge.

Ultimately, transitioning from a “lowest common denominator” traversal strategy to a fully hybrid approach is essential to leverage the diverse access interface expected in dataspaces. Addressing these open research avenues will be critical in developing query engines capable of scaling to the realistic demands of a decentralized data ecosystem.

Declaration on Generative AI

During the writing of this paper, the author(s) used DeepL and GPT-4o in order to: Grammar, translation and spelling check. After using these tool(s)/service(s), the author(s) reviewed and edited the content as needed and take(s) full responsibility for the publication’s content.

References

- [1] A. Schwarte, P. Haase, K. Hose, R. Schenkel, M. Schmidt, 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 10, Springer, 2011, pp. 601–616.
- [2] M. Schmidt, O. Görlitz, P. Haase, G. Ladwig, A. Schwarte, T. Tran, 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 10, Springer, 2011, pp. 585–600.
- [3] L. Heling, M. Acosta, Federated sparql query processing over heterogeneous linked data fragments, in: Proceedings of the ACM Web Conference 2022, 2022, pp. 1047–1057.
- [4] M.-H. Dang, J. Aimonier-Davat, P. Molli, O. Hartig, H. Skaf-Molli, Y. Le Crom, Fedshop: a benchmark for testing the scalability of sparql federation engines, in: International Semantic Web Conference, Springer, 2023, pp. 285–301.
- [5] B. Esteves, H. J. Pandit, V. Rodríguez-Doncel, Odrl profile for expressing consent through granular access control policies in solid, in: 2021 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW), IEEE, 2021, pp. 298–306.
- [6] A. Padia, T. Finin, A. Joshi, et al., Attribute-based fine grained access control for triple stores, in: 3rd Society, Privacy and the Semantic Web-Policy and Technology workshop, 14th International Semantic Web Conference, 2015.
- [7] A. Halevy, M. Franklin, D. Maier, Principles of dataspace systems, in: Proceedings of the twenty-fifth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems, 2006, pp. 1–9.
- [8] E. Curry, Dataspaces: fundamentals, principles, and techniques, in: Real-time Linked Dataspace: Enabling Data Ecosystems for Intelligent Systems, Springer, 2019, pp. 45–62.
- [9] R. Verborgh, M. Vander Sande, O. Hartig, J. Van Herwegen, L. De Vocht, B. De Meester, G. Haezendonck, P. Colpaert, Triple pattern fragments: a low-cost knowledge graph interface for the web, Journal of Web Semantics 37 (2016) 184–206.
- [10] G. Ladwig, T. Tran, Linked data query processing strategies, in: International Semantic Web Conference, Springer, 2010, pp. 453–469.
- [11] M. M. Sabri, A hybrid framework for online execution of linked data queries, in: Proceedings of the 24th International Conference on World Wide Web, 2015, pp. 515–519.
- [12] G. Ladwig, T. Tran, Sihjoin: Querying remote and local linked data, in: Extended Semantic Web Conference, Springer, 2011, pp. 139–153.
- [13] J. Umbrich, A hybrid framework for querying linked data dynamically (2012).
- [14] M. Ragab, Y. Savateev, R. Moosaei, T. Tiropanis, A. Poulovassilis, A. Chapman, G. Roussos, Espresso: a framework for empowering search on decentralized web, in: International Conference on Web Information Systems Engineering, Springer, 2023, pp. 360–375.
- [15] S. Speicher, J. Arwe, A. Malhotra, Linked data platform 1.0, 2015. URL: <https://www.w3.org/TR/ldp/>, w3C Recommendation.
- [16] R. Taelman, J. Van Herwegen, M. Vander Sande, R. Verborgh, Comunica: a modular sparql query engine for the web, in: International Semantic Web Conference, Springer, 2018, pp. 239–255.
- [17] R. Taelman, R. Verborgh, Link traversal query processing over decentralized environments with structural assumptions, in: International Semantic Web Conference, Springer, 2023, pp. 3–22.
- [18] J. Aimonier-Davat, B. Nédelec, M.-H. Dang, P. Molli, H. Skaf-Molli, Fedup: querying large-scale federations of sparql endpoints, in: Proceedings of the ACM Web Conference 2024, 2024, pp. 2315–2324.
- [19] M. Saleem, A.-C. Ngonga Ngomo, Hibiscus: Hypergraph-based source selection for sparql endpoint federation, in: European semantic web conference, Springer, 2014, pp. 176–191.
- [20] O. Hartig, J.-C. Freytag, Foundations of traversal based query execution over linked data, in: Proceedings of the 23rd ACM conference on Hypertext and social media, 2012, pp. 43–52.

- [21] O. Hartig, Squin: a traversal based query execution system for the web of linked data, in: Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data, 2013, pp. 1081–1084.
- [22] B. Bogaerts, B. Ketsman, Y. Zeboudj, H. Aamer, R. Taelman, R. Verborgh, Distributed subweb specifications for traversing the web, Theory and Practice of Logic Programming 24 (2024) 394–420.
- [23] O. Hartig, M. T. Özsü, Walking without a map: Ranking-based traversal for querying linked data, in: International Semantic Web Conference, Springer, 2016, pp. 305–324.
- [24] R. Avnur, J. M. Hellerstein, Eddies: Continuously adaptive query processing, in: Proceedings of the 2000 ACM SIGMOD international conference on Management of data, 2000, pp. 261–272.
- [25] V. Raman, A. Deshpande, J. M. Hellerstein, Using state modules for adaptive query processing, in: Proceedings 19th International Conference on Data Engineering (Cat. No. 03CH37405), IEEE, 2003, pp. 353–364.
- [26] A. Deshpande, J. M. Hellerstein, Lifting the burden of history from adaptive query processing, in: Proceedings of the Thirtieth international conference on Very large data bases-Volume 30, 2004, pp. 948–959.
- [27] C. Buil-Aranda, A. Hogan, J. Umbrich, P.-Y. Vandenbussche, Sparql web-querying infrastructure: ready for action?, in: International Semantic Web Conference, Springer, 2013, pp. 277–293.
- [28] A. V. Sambra, E. Mansour, S. Hawke, M. Zereba, N. Greco, A. Ghanem, D. Zagidulin, A. Aboulnaga, T. Berners-Lee, Solid: a platform for decentralized social applications based on linked data, MIT CSAIL & Qatar Computing Research Institute, Tech. Rep. 2016 (2016).