

A Prospective Analysis of Security Vulnerabilities within Link Traversal-Based Query Processing

Ruben Taelman, Ruben Verborgh

IDLab, Department of Electronics and Information Systems, Ghent University – imec, {firstname.lastname}@ugent.be

Abstract

The societal and economic consequences surrounding Big Data-driven platforms have increased the call for decentralized solutions. However, retrieving and querying data in more decentralized environments requires fundamentally different approaches, whose properties are not yet well understood. Link-Traversal-based Query Processing (LTQP) is a technique for querying over decentralized data networks, in which a client-side query engine discovers data by traversing links between documents. Since decentralized environments are potentially unsafe due to their non-centrally controlled nature, there is a need for client-side LTQP query engines to be resistant against security threats aimed at the query engine's host machine or the query initiator's personal data. As such, we have performed an analysis of potential security vulnerabilities of LTQP. This article provides an overview of security threats in related domains, which are used as inspiration for the identification of 10 LTQP security threats. This list of security threats forms a basis for future work in which mitigations for each of these threats need to be developed and tested for their effectiveness. With this work, we start filling the unknowns for enabling query execution over decentralized environments. Aside from future work on security, wider research will be needed to uncover missing building blocks for enabling true data decentralization.

Keywords

Linked Data, RDF, Link Traversal, SPARQL, Security, Solid

1. Introduction

Contrary to the Web's initial design as a *decentralized* ecosystem, the Web has grown to be a very centralized place, as large parts of the Web are currently made up of a few large Big Data-driven centralized platforms [1]. This large-scale centralization has led to a number of problems related to personal information abuse, and other economic and societal problems. In order to solve these problems, there are calls to go back to the original vision of a decentralized Web. The leading effort to achieve this decentralization is Solid [1]. Solid proposes a radical *decentralization* of data across *personal data vaults*, where everyone is in full control of its own personal data vault. This vault can contain any number of documents, where its owner can determine who or what can access what parts of this data. In contrast to the current state of the Web where data primarily resides in a small number of huge data sources, Solid leads to a Web where data is spread over a huge number of data sources.

Our focus in this article is not on decentralizing data, but on finding data after it has been decentralized, which can be done via *query processing*. The issue of query processing over data has been

6th Workshop on Storing, Querying and Benchmarking Knowledge Graphs (QuWeDa) at ISWC 2022, virtual

* Corresponding author.

✉ ruben.taelman@ugent.be (R. Taelman); ✉ ruben.verborgh@ugent.be (R. Verborgh)



2022 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).



CEUR Workshop Proceedings (CEUR-WS.org)

primarily tackled from a Big Data standpoint so far. However, if decentralization efforts such as Solid will become a reality, we need to be prepared for the need to query over a huge number of data sources. For example, decentralized social networking applications will need to be able to query over networks of friends containing hundreds or thousands of data documents. As such, we need new query techniques that are specifically designed for such levels of decentralization. A promising technique to achieve this is Link-Traversal-based Query Processing (LTQP) [2, 3]. LTQP is able to query over a set of documents that are connected to each other via *links*. An LTQP engine typically starts from one or more documents, and *traverses* links between them in a crawling-manner to resolve the given query.

Since LTQP is still a relative young area of research, there are still a number of open problems that need to be tackled, notably result completeness and query termination [2]. Aside from these known issues, we also state the importance of *security*. Security is a highly important and well-investigated topic in the context of Web applications [4, 5], but it has not yet been investigated in the context of LTQP. As such, **we investigate in this article security issues related to LTQP engines**, which may threaten the integrity of the user’s data, machine, and user experience, but also lead to privacy issues if personal data is unintentionally leaked. Specifically, we focus on data-driven security issues that are inherent to LTQP due to the fact that it requires a query engine to follow links on the Web, which is an uncontrolled, unpredictable and potentially unsafe environment. Instead of analyzing a single security threat in-depth, we perform a broader high-level analysis of multiple security threats.

Since LTQP is still a relatively new area of research, its real-world applications are currently limited. As such, we can not learn from security issues that arose in existing systems. Instead of waiting for – potentially unsafe– widespread applications of LTQP, we draw inspiration from related domains that *are already well-established*. Specifically, we draw inspiration from the domains of crawling and Web browsers in Section 2, and draw links to what impact these known security issues will have on LTQP query engines. In Section 3, we introduce a guiding use case that will be used to illustrate different threats. After that, we discuss our method of categorizing vulnerabilities in Section 4. Next, we list 10 data-driven security vulnerabilities related to LTQP in Section 5, which are derived from known vulnerabilities in similar domains, and through analysis of the LTQP implementation within the Comunica query engine [6]. Finally, we discuss the future of LTQP security and conclude in Section 6.

2. Related Work

This section lists relevant related work in the topics of LTQP and security.

2.1. Link-Traversal-Based Query Processing

More than a decade ago, Link-Traversal-based Query Processing (LTQP) [3, 2] was been introduced as an alternative query paradigm for enabling query execution over document-oriented interfaces. These documents are usually Linked Data [7] serialized using any RDF serialization. RDF is suitable for LTQP and decentralization because of its global semantics, which allows queries to be written independently of the schemas of specific documents. In order to execute these queries, LTQP processing occurs over live data, and discovers links to other documents via the *follow-your-nose principle* during query execution. This is in contrast to the typical query execution over centralized database-oriented inter-

faces such as SPARQL endpoints, where data is assumed to be loaded into the endpoint beforehand, and no additional data is discovered during query execution.

Concretely, LTQP typically starts off with an input query and a set of seed documents. The query engine then dereferences all seed documents via an HTTP GET request, discovers links to other documents inside those documents, and recursively dereferences those discovered documents. Since document discovery can be a very long (or infinite) process, query execution happens during the discovery process based on all the RDF triples that are extracted from the discovered documents. This is typically done by implementing these processes in an iterative pipeline [8]. Furthermore, since this approach can lead to a large number of discovered documents, different reachability criteria [9] have been introduced to restrict what links to follow for a given query.

So far, LTQP research in the area of security has been limited. One work has indicated the importance of *trustworthiness* [10] during link traversal, as people may publish false or contradicting information, which would need to be avoided or filtered out during query execution. Another work mentioned the need for LTQP engines to adhere to `robots.txt` files [11] in order to not lead to unintentional denial of service attacks of data publishers. Given the focus of our work on data-driven security vulnerabilities in LTQP engines, we only consider this issue of *trustworthiness* further in this work, and omit the security vulnerabilities from a data publisher's perspective.

2.2. Vulnerabilities Of RDF Query Processing

Research involving the security vulnerabilities of RDF query processing has been primarily focused on injection attacks within Web applications that internally send SPARQL queries to a SPARQL endpoint [4, 5, 12, 13]. So far, no research has been done on vulnerabilities specific to RDF federated querying or link traversal.

2.3. Linked Data Access Control

Kirrane et al. [14] surveyed the existing approaches for achieving access control in RDF, for both authentication and authorization. The authors mention that only a minority of those works apply specifically to the document-oriented nature of Linked Data. They do however mention that non-Linked-Data-specific approaches could potentially be applied to Linked Data in future work. To the best of our knowledge, no security vulnerabilities have yet been identified for any of these.

2.4. Web Crawlers

Web crawling [15] is a process that involves collecting information on the Web by following links between pages. Web crawlers are typically used for Web indexing to aid search engines. Focused crawling [16] is a special form of Web crawling that prioritizes certain Web pages, such as Web pages about a certain topic, or domains for a certain country. LTQP can therefore be considered as an area of focused crawling where the priority lies in achieving query results.

One related work in this area involves abusing crawlers to initiate attacks on other Web sites [17]. This may cause performance degradation on the attacked Web site, or could even cause the crawling agent to be blocked by the server. These attacks involve convincing the crawler to follow a link to a third-party Web site that exploits a certain vulnerability, such as an SQL injection. Additionally, this work describes a type of attack that allows vulnerable Web sites to be used for improving the PageRank [18] of an attacker-owned Web site via forged backlinks.

Some other works focus on mitigation of so-called *crawler traps* [19,20] or *spider traps*. These are sets of URLs that cause an infinite crawling process, which can either be intentional or accidental. Such crawler traps can have multiple causes:

- Links between dynamic pages that are based on URLs with query parameters;
- Infinite redirection loops via using the HTTP 3xx range;
- Links to search APIs;
- Infinitely paged resources, such as calendars;
- Incorrect relative URLs that continuously increase the URL length.

Crawler traps are mostly discovered through human intervention when many documents in a single domain are discovered. Recently, a new detection technique was introduced [21] that attempts to measure the *distance* between documents, and rejects links to documents that are too similar.

2.5. Web Browsers

Web browsers enable users to visualize and interact with Web pages. This interaction is closely related to LTQP, with the main difference that LTQP works autonomously, while Web browsers are user-driven. Silic et al. [22] analyzed the architectures of modern Web browsers, determined the main vulnerabilities, and discuss how these issues are handled. They list the following main threats for Web browsers:

1. **System compromise:** Malicious arbitrary code execution with full privileges on behalf of the user. For example, exploits in the browser or third-party plugins caused by bugs. These types of attacks are mitigated through automatic updates once exploits become known.
2. **Data theft:** Ability to steal local network or system data. For example, a Web page includes a sub-resource to URLs using the file scheme (`file://`), which are usually blocked.
3. **Cross domain compromise:** Code from a Fully Qualified Domain Name (FQDN) executes code (or reads data) from another FQDN. For example, a malicious domain could extract authentication cookies from your bank's website you are logged into. This is usually blocked through the same-origin policy, but can be explicitly allowed through Cross-Origin Resource Sharing (CORS) (<https://fetch.spec.whatwg.org/#http-cors-protocol>).
4. **Session hijacking:** Session tokens are compromised through theft or session token prediction. For example, cross-domain request forgery (CSRF) [23] is a type of attack that involves an attacker forcing a user logged in on another Web site to perform an action without their consent. Web browsers do not protect against these attacks, but they are typically handled by Web frameworks via the Synchronizer Token Pattern [24].
5. **User interface compromise:** Manipulating the user interface to trick the user into performing an action without their knowledge. For example, placing an invisible button in front of another button. This category also includes CPU and memory hogging to block the user from taking any further actions. Web browsers have limited protections for these types of attacks that involve placing limitations on user interface manipulations.

3. Use Case

In this section, we introduce a use case that will be used to illustrate the security threats discussed throughout this article.

We assume a Web with public and private information, which may for instance be achieved via personal data vaults following the principles of the Solid ecosystem [1]. This data vault is in full control of the owner, and they can host any kind of file in here, such as Linked Data files.

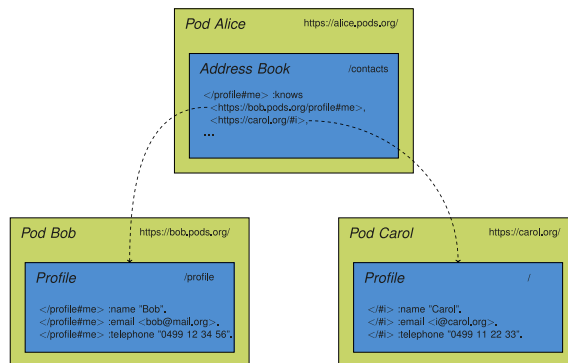


Fig. 1: Overview of the address book use case in which Alice has an address book with links to the profiles of Carol and Bob, which contain further details.

For this use case, we assume the existence of three people (Alice, Bob, and Carol), each having their own personal data vault. Alice uses her vault to store an address book containing the people she knows. Instead of storing contact details directly in the address book, she stores *links* to the profiles of her contacts (Bob and Carol). Bob and Carol can then self-define their own contact details. Fig. 1 shows an illustration of this setup.

The LTQP paradigm is well-suited to handle query execution over such setups. If Alice for instance would like to obtain the names of all her contacts, she could initiate a query starting from her address book as seed document, and the query engine would follow the links to her contacts, and obtain the names from their respective profiles. Some documents may require authentication before they can be accessed, for which Alice's query engine makes use of Alice's identity. In all threats throughout this article, we assume that Carol has malicious intentions that Alice is unaware of.

In this use case, two main roles can be identified. The first is the role of data publisher, which is taken up by Alice, Bob, and Carol through their personal data vaults. The second is the role of the query initiator, which here applies to Alice, as she issues a query over her contacts.

4. Classification Of Security Vulnerabilities

In this section, we first introduce the background on classifying security vulnerabilities in software. After that, we introduce a classification method specifically for the LTQP domain, to assess the validity of our work.

4.1. Background

Security vulnerabilities in software can be classified using many different methods [25,26]. Generic classification methods often result in very large taxonomies, which are shown to result in practical problems [25] due to their size and complexity.

Seacord et al. [25] claim that classification methods must be based on engineering analysis of the problem domain, instead of being too generic. For this, they suggest the use of domain-specific attributes for classifying security vulnerabilities for each domain separately. Furthermore, they introduce the following terminology for security vulnerabilities, by building upon earlier formal definitions of vulnerabilities [26]:

Security Flaw A defect in a software application or component that, when combined with the necessary conditions, can lead to a software vulnerability.

Vulnerability A set of conditions that allows violation of an explicit or implicit security policy.

Exploit A technique that takes advantage of a security vulnerability to violate an explicit or implicit security policy.

Mitigation Techniques to prevent or limit exploits against vulnerabilities.

For the remainder of this article, we will make use of this terminology, and we adopt a method hereafter for classifying software vulnerabilities specific to the LTQP domain as recommended by Seacord et al. [25].

4.2. Classification Method

Our classification method considers the listing of several security *vulnerabilities*. Each vulnerability has one or more possible *exploits*, which may take advantage of this vulnerability. The different properties of each exploit are shown in Table 1.

Attribute	Values
Attacker	Data publisher, ...
Victim	LTQP engine, query initiator, data publisher, ...
Impact	Incorrect query results, system crash, ...
Difficulty	Easy, medium, hard

Table 1: Exploit properties specific to LTQP, with several possible values for each attribute.

5. Data-Driven Vulnerabilities

As shown before in Subsection 2.2, most research on identifying security vulnerabilities within RDF query processing focuses on the query itself as a means of attacking, mostly through injection techniques. Since LTQP engines also accept queries as input, these existing techniques will therefore also apply to LTQP engines.

In this work, we acknowledge the importance of these vulnerabilities, but we instead place our attention onto a new class of vulnerabilities that are specific to LTQP engines as a consequence of the open and uncontrolled nature of data on the Web. Concretely, we consider two main classes of security vulnerabilities to LTQP engines:

1. **Query-driven:** vulnerabilities that are caused by modifying queries that are the input to certain query engines.
2. **Data-driven:** vulnerabilities that are caused by the presence, structuring, or method of publishing data on the Web.

To the best of our knowledge, all existing work on security vulnerabilities within RDF querying has focused on query-driven vulnerabilities. Given its importance for LTQP engines, we focus on data-driven vulnerabilities for the remainder of this work.

Hereafter, we explain and classify each vulnerability using the classification method from Section 4. For each vulnerability, we provide at least one possible example of an *exploit* based on our use case.

Unless mentioned otherwise, we do not make any assumptions about specific forms or semantics of LTQP, which can influence which links are considered. The only general assumption we make is that we have an LTQP query engine that follows links in any way, and executes queries over the union of the discovered documents.

5.1. Unauthoritative Statements

A consequence of the open-world assumption [27] where anyone can say anything about anything, is that both valid and invalid (and possibly malicious) things can be said. When a query engine is traversing the Web, it is therefore possible that it can encounter information that impacts the query results in an undesired manner. This information could be *untrusted* [10, 28], *contradicting*, or *incorrect*. Without mitigations to this vulnerability, query results from an LTQP can therefore never be really trusted, which brings the practical broad use of LTQP into question.

Exploit: producing untrusted query results by adding unauthoritative triples

Given our use case, Carol could for instance decide to add one additional triple to her profile, such as: `<https://bob.pods.org/profile#me> :name "Dave"`. She would therefore indicate that Bob's name is "Dave". This is obviously false, but she is "allowed" to state this under the open world assumption. However, this means that if Alice would naively query for all her friend's names via LTQP, she would have two names for Bob appear in her results, namely "Bob" and "Dave", where this second result may be undesired.

Attacker Data publisher (Carol)

Victim Query results from the LTQP engine of Alice

Impact Untrusted query results

Difficulty Easy (adding triples to an RDF document)

5.2. Intermediate Result And Query Leakage

This vulnerability assumes the existence of a *hybrid* LTQP query engine that primarily traverses links, but can exploit database-oriented interfaces such as SPARQL endpoints if they are detected in favour of a range of documents. Furthermore, we assume a range of documents that require authentication, as their contents are not accessible to everyone. Query engines typically decompose queries into smaller sub-queries, and join these intermediate results together afterwards. In the case of a hybrid LTQP engine, intermediate results that are obtained from the traversal process from non-public documents could be joined with data from a discovered SPARQL endpoint. An attacker could therefore set up an interface that acts as a SPARQL endpoint, but is in fact a **malicious interface that intercepts intermediate results** from LTQP engines.

Exploit: capturing intermediary results via malicious SPARQL endpoint

Based on our use case, Carol could include a triple with a link to the SPARQL endpoint at `http://attacker.com/sparql`. If Alice makes use of a hybrid LTQP engine with an adaptive query planner, this

internal query planner could decide to make use of this malicious endpoint once it has been discovered. Depending on the query planner, this could mean that non-public intermediate results from the traversal process such as Bob's telephone are used as input to the malicious SPARQL endpoint. Other query planning algorithms could even decide to send the full original SPARQL query to the malicious endpoint. Depending on the engine and its query plan, this could give the attacker knowledge of intermediate results, or even the full query. This vulnerability enables attackers to obtain insights into user behaviour, which is a privacy concern. A more critical problem is when private data is being leaked that normally exists behind access control, such as bank account numbers.

Attacker SPARQL endpoint publisher (Carol)

Victim Intermediary results of the LTQP engine of Alice

Impact Leakage of (intermediary) query results

Difficulty Medium (setting up a malicious SPARQL endpoint)

5.3. Session Hijacking

In this vulnerability, we assume the presence of some form of authentication (such as WebID-OIDC [29]) that leads to an active authenticated session. This vulnerability is similar to that of Web browsers, where the session token can be compromised through theft or session token prediction. Such a vulnerability could lead to cross-domain request forgery (CSRF) [23] attacks, where an attacker forces the user to perform an action while authenticated without the user's consent.

Exploit: triggering unintended operations on SPARQL endpoint behind access control

For example, we assume that Alice has a flawed SPARQL endpoint running at `http://my-endpoint.com/sparql`, which requires Alice's session for accepting read and write queries. Alice's query engine may have Alice's session stored by default for when she wants to query against her own endpoint. If Carol knows this, she could add a malicious triple with a link to `http://my-endpoint.com/sparql?query=DELETE * WHERE { ?s ?p ?o }` in her profile. While the SPARQL protocol only allows update queries via HTTP POST, Alice's flawed query engine could implement this incorrectly so that update queries are also accepted via HTTP GET. If Alice executes a query over her address book, the query engine could dereference this link with her session enabled, which would cause her endpoint to be cleared. This vulnerability is however not specific to SPARQL endpoints, but may occur on any type of Web API that allows modifying data via HTTP GET requests.

Attacker Data publisher (Carol)

Victim Alice's stored data

Impact Removal or modification of Alice's stored data

Difficulty Easy (adding a malicious link to flawed endpoint)

5.4. Cross-Site Data Injection

This vulnerability concerns ways by which attackers can inject data or links into documents. For instance, HTTP GET parameters are often used to parameterize the contents of documents. If such parameters are not properly validated or escaped, they can be used by attackers to include malicious data or links.

Exploit: injecting untrusted links via flawed trusted API

For example, assuming Alice executes a query over a page from Carol, and a compromised API `http://trusted.org/?name` that dynamically creates RDF responses based on the `?name` HTTP GET parameters. In this case, the API simply has a Turtle document template into which the name is filled in as a literal value, but it does not do any escaping. We assume Alice decides to fully trust all links from `http://trusted.org/` to other pages, but only trusts information directly on Carol's page or links to other trusted domains. If Carol includes a link to `<http://trusted.org/?name=Bob>`. `<> rdfs:seeAlso <http://hacker.com/invalid-data>`. `<> foaf:name "abc"`, then this would cause the API to produce a Turtle document that contains a link to `http://hacker.com/invalid-data`, which would lead to unwanted data to be included in the query results.

Attacker Data publisher (Carol)

Victim Query results from the LTQP engine of Alice

Impact Untrusted query results

Difficulty Easy (adding triples to an RDF document)

5.5. Arbitrary Code Execution

Advanced crawlers such as the Googlebot [30] allow JavaScript logic to be executed for a limit duration, since certain HTML pages are built dynamically via JavaScript at the client-side. In this vulnerability, we assume a similar situation for LTQP, where Linked Data pages may also be created client-side via an expressive programming language such as JavaScript. This would in fact already be applicable to HTML pages that dynamically produce JSON-LD script tags or RDFa in HTML via JavaScript. In order to query over such dynamic Linked Data pages, a query engine must initiate a process similar to Googlebot's JavaScript execution phase. Such a process does however open the door to potentially major security vulnerabilities if malicious code is being read and executed by the query engine during traversal.

Exploit: manipulate local files via overprivileged JavaScript execution

For example, we assume that Alice's LTQP query engine executes JavaScript on HTML pages before extracting its RDFa and JSON-LD. Furthermore, this LTQP engine has a security flaw that allows executed JavaScript code to access and manipulate the local file system. Carol could include a malicious piece of JavaScript code in her profile that makes use of this flaw to upload all files on the local file system to the attacker, and deletes all files afterwards so that she can hold Alice's data for ransom.

Attacker Data publisher (Carol)

Victim Files on machine in which Alice's query engine runs

Impact Removal or modification of files on Alice's machine

Difficulty Easy (adding JavaScript code to a document)

5.6. Link Traversal Trap

LTQP by nature depends on the ability of iteratively following links between documents. It is however possible that such **link structures cause infinite traversal paths** and make the traversal engine get trapped, either intentional or unintentional, just like crawler traps. Given this reality, LTQP query

engines must be able to detect such traps. Otherwise, query engines could never terminate, and possibly even produce infinite results.

Exploit: forming a link cycle

A link cycle is a simple form of link traversal trap that could be formed in different ways. First, at the application-level, Carol's profile could contain a link path to document X, and document X could contain a link path back to Carol's profile. Second, at the HTTP protocol-level, Carol's server could return for her profile's URL an (HTTP 3xx) redirect chain to URL X, and URL X could contain a redirect chain back to the URL of her profile. Third, at the application level, a cycle structure could be simulated via virtual pages that always link back to similar pages, but with a different URL. For example, the Linked Open Numbers [31] project generates a long virtual sequence of natural numbers, which could produce a bottleneck when traversed by an LTQP query engine.

Attacker Data publisher (Carol)

Victim Query process of Alice's query engine

Impact Unresponsiveness of Alice's query engine

Difficulty Easy

5.7. System Hogging

The *user interface compromise* vulnerability for Web browsers includes attacks involving CPU and memory hogging through (direct or indirect) malicious code execution or by exploiting software flaws. Such vulnerabilities also exist for LTQP query engines, especially regarding the use of different RDF serializations, and their particularities with respect to parsing.

Exploit: producing infinite RDF documents

For example, RDF serializations such as Turtle [32] are implicitly designed as to allow streaming serialization and deserialization. JSON-LD even explicitly allows this through its Streaming JSON-LD note [33]. Due to this streaming property, RDF documents of infinite size can be generated, since serializations place no limits on their document sizes. Valid use cases exist for publishers to generate infinite RDF documents, which can be streamed to query engines. Query engines with non-streaming or flawed streaming parsers, can lead to CPU and memory issues. Furthermore, similar issues can occur due to very long or infinite IRIs or literals inside documents. Other attacks could exist that specifically target known flaws in RDF parsers that cause CPU or memory issues.

Attacker Data publisher (Carol)

Victim Machine in which Alice's query engine runs

Impact Unresponsiveness or crashing of Alice's query engine or machine

Difficulty Easy

5.8. Document Corruption

Since the Web is not a centrally controlled system, it is possible that documents are incorrectly formatted, either intentional or unintentional. RDF formats typically prescribe a restrictive syntax, which require parsers to emit an error when it encounters illegal syntax. When an LTQP engine discovers and parses a large number of RDF documents, possibly in an uncontrolled manner, it is undesired that a

syntax error in just a single RDF document can cause the whole query process to terminate with an error. Furthermore, the phenomenon of *Link Rot* [34] can lead to links going dead (HTTP 404) at any point in time, while finding a link to a URL that produces a 404 response should not always cause the query engine to terminate.

Exploit: publishing an invalid RDF document

For example, Carol could decide to introduce a syntax error in her profile document, or she could simply remove it to produce a 404 response. This would could cause Alice’s queries over her friends from that point on to fail.

Attacker Data publisher (Carol)

Victim Alice’s query engine

Impact Crashing of Alice’s query engine

Difficulty Easy

5.9. Cross-Query Execution Interaction

Query engines of all forms typically make use of caching techniques to improve performance of query execution. LTQP query engines can leverage caching techniques for document retrieval. Within a single query execution, or across multiple query executions, the documents may be reused, which could reduce the overall number of HTTP requests. Such forms of caching can lead to vulnerabilities based on information leaking across different query executions. We therefore make the assumption of caching-enabled LTQP engines in this vulnerability.

Exploit: timing attack to determine prior knowledge

A first exploit of this vulnerability is an attack that enables Carol to gain knowledge about whether or not Bob’s profile has been requested before by Alice. We assume that Alice’s engine issues a query over a document from Carol listing all her pictures. We also assume that Bob’s profile contains a link to Carol’s profile. If Carol includes a link from her pictures document to Bob’s profile, and Bob’s profile already links to Carol’s profile, then the query engine could fetch these three documents in sequence (Carol’s pictures, Bob’s profile, Carol’s profile). Since Carol’s pictures and profile are in control of Carol, she could perform a timing attack [35] to derive how long the Alice’s query engine took to process Bob’s profile. Since HTTP delays typically form the bottleneck in LTQP, Carol could thereby derive if Bob’s profile was fetched from a cache or not. This would enable Carol to gain knowledge about prior document lookups, which could for example lead to privacy issues with respect to the user’s interests.

Attacker Data publisher (Carol)

Victim Privacy about Alice’s document usage

Impact Alice’s document usage becomes known to Carol

Difficulty Hard

Exploit: unauthenticated cache reuse

A second exploit assumes the presence of a software flaw inside Alice’s LTQP query engine that makes document caches ignore authorization information. This example is also a form of the *Intermediate Result and Query Leakage* vulnerability that was explained before, for which we assume the existence of a *hybrid* LTQP query engine. If Alice queries a private file containing her passwords

from a server using its authentication key, this can cause this passwords file to be cached. If Carol has a query endpoint that is being queried by Alice, and Carol is aware of the location of Alice's passwords, then she could maliciously introduce a link to Alice's passwords file. Even if the query was not executed with Alice's authentication key, the bug in Alice's query engine would cause the passwords file to be fetched in full from the cache, which could cause parts of it to be leaked to Carol's query endpoint.

Attacker Data publisher (Carol)

Victim Alice's private data

Impact Alice's private data is leaked

Difficulty Easy (if cache is flawed)

5.10. Document Priority Modification

Different techniques are possible to determine the priority of documents [36] during query processing. If queries do not specify a custom ordering, this prioritization will impact the ordering of query results. Some of these techniques are purely graph-based, such as PageRank [18], and can therefore suffer from purely data-driven attacks. This vulnerability involves attacks that can influence the priority of documents, and thereby maliciously influence what query results come in earlier or later.

Exploit: malicious PageRank prioritization of documents

One possible exploit is similar to the attack to modify priorities within crawlers [17]. We assume that Alice issues a query that returns grocery stores in the local area, which is executed via a LTQP query engine that makes use of PageRank to prioritize documents. Furthermore, we assume a highly-scoring, but vulnerable API that accepts HTTP GET parameters that can be abused to inject custom URLs inside the API responses. If Carol aims to increase the ranking of her grocery store within Alice's query for better visibility, then she could exploit this vulnerable API. Concretely, Carol could place links from the grocery store's page to this vulnerable API using GET parameters that would cause it to link back to Carol's grocery store. Such an attack would lead to a higher PageRank for Carol's grocery store, and therefore an earlier handling and result representation of Carol's grocery store.

Attacker Data publisher (Carol)

Victim Order of Alice's query results

Impact Carol's page is ranked higher

Difficulty Medium

6. Conclusions

With this prospective analysis, we have illustrated the importance of more security-oriented research in the domain on LTQP and the general handling of decentralized environments such as Solid [1], especially in the presence of data behind authentication. In future work, work is needed to determine mitigation strategies for them, which may be inspired by existing techniques in related domains. Furthermore, research will be needed to test the impact of these mitigations on implementations, analyze their performance impact, introduce more performant techniques and algorithms, and introduce and apply attack models to test their effectiveness. Since our analysis of security vulnerabilities is by no means exhaustive, additional research efforts are needed to uncover and predict potential security vulnerabilities

in LTQP. Such future research—with our work as a first step—is crucial for enabling a decentralized Web which we can query securely.

Acknowledgements

Ruben Taelman is a postdoctoral fellow of the Research Foundation – Flanders (FWO) (1274521N).

References

- [1] Verborgh, R.: Re-decentralizing the Web, for good this time. In: Seneviratne, O. and Hendler, J. (eds.) *Linking the World's Information: A Collection of Essays on the Work of Sir Tim Berners-Lee* (2020).
- [2] Hartig, O.: An Overview on Execution Strategies for Linked Data Queries. *Datenbank-Spektrum*. 13, 89–99 (2013).
- [3] Hartig, O., Bizer, C., Freytag, J.-C.: Executing SPARQL Queries over the Web of Linked Data. In: Bernstein, A., Karger, D.R., Heath, T., Feigenbaum, L., Maynard, D., Motta, E., and Thirunarayan, K. (eds.) *Proceedings of the 8th International Semantic Web Conference*. pp. 293–309 (2009).
- [4] Halfond, W.G., Viegas, J., Orso, A., others: A classification of SQL-injection attacks and countermeasures. In: *Proceedings of the IEEE international symposium on secure software engineering*. pp. 13–15. IEEE (2006).
- [5] Orduña, P., Almeida, A., Aguilera, U., Laiseca, X., López-de-Ipiña, D., Goiri, A.G.: Identifying security issues in the semantic web: Injection attacks in the semantic query languages. *Actas de las {VI} Jornadas Científico-Técnicas en Servicios Web y {SOA}*. 51, 4529–4542 (2010).
- [6] Taelman, R., Van Herwegen, J., Vander Sande, M., Verborgh, R.: Comunica: a Modular SPARQL Query Engine for the Web. In: *Proceedings of the 17th International Semantic Web Conference* (2018).
- [7] Berners-Lee, T.: *Linked Data*. <https://www.w3.org/DesignIssues/LinkedData.html> (2006).
- [8] Hartig, O.: SQUIN: a Traversal Based Query Execution System for the Web of Linked Data. In: *Proceedings of the ACM SIGMOD International Conference on Management of Data*. pp. 1081–1084. ACM (2013).
- [9] Hartig, O.: SPARQL for a Web of Linked Data: Semantics and Computability. In: *Proceedings of the 9th international conference on The Semantic Web: research and applications*. pp. 8–23. Berlin, Heidelberg (2012).
- [10] Verborgh, R., Taelman, R.: Guided Link-Traversal-Based Query Processing. Presented at the May (2020).
- [11] Umbrich, J., Hose, K., Karnstedt, M., Harth, A., Polleres, A.: Comparing data summaries for processing live queries over linked data. *World Wide Web*. 14, 495–544 (2011).
- [12] Yang, X., Chen, Y., Zhang, W., Zhang, S.: Exploring injection prevention technologies for security-aware distributed collaborative manufacturing on the Semantic Web. *The International Journal of Advanced Manufacturing Technology*. 54, 1167–1177 (2011).
- [13] Asghar, H., Anwar, Z., Latif, K.: A deliberately insecure RDF-based Semantic Web application framework for teaching SPARQL/SPARUL injection attacks and defense mechanisms. *computers & security*. 58, 63–82 (2016).
- [14] Kirrane, S., Mileo, A., Decker, S.: Access control and the resource description framework: A survey. *Semantic Web*. 8, 311–352 (2017).
- [15] Shkapenyuk, V., Suel, T.: Design and implementation of a high-performance distributed web crawler. In: *Proceedings 18th International Conference on Data Engineering*. pp. 357–368. IEEE (2002).
- [16] Novak, B.: A survey of focused web crawling algorithms. *Proceedings of SIKDD*. 5558, 55–58 (2004).

- [17] Zarras, A., Maggi, F.: Hiding Behind the Shoulders of Giants: Abusing Crawlers for Indirect Web Attacks. In: 2017 15th Annual Conference on Privacy, Security and Trust (PST). pp. 355–35509. IEEE (2017).
- [18] Page, L., Brin, S., Motwani, R., Winograd, T.: The PageRank citation ranking: Bringing order to the web. Stanford InfoLab (1999).
- [19] Crawler traps: how to identify and avoid them. <https://www.contentkingapp.com/academy/crawler-traps/> (2020).
- [20] Heydon, A., Najork, M.: Mercator: A scalable, extensible web crawler. *World Wide Web*. 2, 219–229 (1999).
- [21] David, B., Delong, M., Filiol, E.: Detection of crawler traps: formalization and implementation –defeating protection on internet and on the TOR network. *Journal of Computer Virology and Hacking Techniques*. 1–14 (2021).
- [22] Šilić, M., Krolo, J., Delač, G.: Security vulnerabilities in modern web browser architecture. In: The 33rd International Convention MIPRO. pp. 1240–1245 (2010).
- [23] Barth, A., Jackson, C., Mitchell, J.C.: Robust defenses for cross-site request forgery. In: Proceedings of the 15th ACM conference on Computer and communications security. pp. 75–88 (2008).
- [24] Alur, D., Crupi, J., Malks, D.: Core J2EE patterns: best practices and design strategies. Prentice Hall Professional (2003).
- [25] Seacord, R.C., Householder, A.D.: A structured approach to classifying security vulnerabilities. CARNEGIE-MELLON UNIV PITTSBURGH PA SOFTWARE ENGINEERING INST (2005).
- [26] Fithen, W.L., Hernan, S.V., O’Rourke, P.F., Shinberg, D.A.: Formal modeling of vulnerability. *Bell Labs technical journal*. 8, 173–186 (2004).
- [27] Drummond, N., Shearer, R.: The open world assumption. In: eSI Workshop: The Closed World of Databases meets the Open World of the Semantic Web (2006).
- [28] Umbrich, J., Hogan, A., Polleres, A., Decker, S.: Link traversal querying for a diverse web of data. *Semantic Web*. 6, 585–624 (2015).
- [29] WebID-OIDC Authentication Spec. Solid, <https://github.com/solid/webid-oidc-spec> (2019).
- [30] Google: Understand the JavaScript SEO basics. <https://developers.google.com/search/docs/guides/javascript-seo-basics> (2021).
- [31] Vrandečić Denny, Krötzsch, M., Rudolph, S., Lösch, U.: Leveraging non-lexical knowledge for the linked open data web. 5th Review of April Fool’s day Transactions. 18–27 (2010).
- [32] Prud’hommeaux, E., Carothers, G., Machina, L.: JSON-LD 1.1 Processing Algorithms and API. W3C, <https://www.w3.org/TR/turtle/> (2014).
- [33] Taelman, R.: Streaming JSON-LD. W3C, <https://www.w3.org/TR/json-ld11-streaming/> (2020).
- [34] Fetterly, D., Manasse, M., Najork, M., Wiener, J.L.: A large-scale study of the evolution of Web pages. *Software: Practice and Experience*. 34, 213–237 (2004).
- [35] Dhem, J.-F., Koeune, F., Leroux, P.-A., Mestre, P., Quisquater, J.-J., Willems, J.-L.: A practical implementation of the timing attack. In: International Conference on Smart Card Research and Advanced Applications. pp. 167–182 (1998).
- [36] Hartig, O., Özsu, M.T.: Walking Without a Map: Ranking-Based Traversal for Querying Linked Data. In: Proceedings of the 13th International Semantic Web Conference. pp. 305–324 (2016).