



HAL
open science

Social, Structured, and Semantic Search

Raphaël Bonaque, Bogdan Cautis, François Goasdoué, Ioana Manolescu

► **To cite this version:**

Raphaël Bonaque, Bogdan Cautis, François Goasdoué, Ioana Manolescu. Social, Structured, and Semantic Search. [Research Report] RR-8797, Inria Saclay Ile de France. 2015, 38 p. hal-01218116

HAL Id: hal-01218116

<https://inria.hal.science/hal-01218116v1>

Submitted on 20 Oct 2015

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Social, Structured, and Semantic Search

Raphaël Bonaque, Bogdan Cautis, François Goasdoué, Ioana
Manolescu

**RESEARCH
REPORT**

N° 8797

October 2015

Project-Team OAK

ISRN INRIA/RR--8797--FR+ENG

ISSN 0249-6399



Social, Structured, and Semantic Search

Raphaël Bonaque^{*†}, Bogdan Cautis^{†*}, François Goasdoué^{‡*},
Ioana Manolescu^{*†}

Project-Team OAK

Research Report n° 8797 — October 2015 — 38 pages

Abstract: Social content such as blogs, tweets, news etc. is a rich source of interconnected information. We identify a set of requirements for the meaningful exploitation of such rich content, and present a new data model, called **S3**, which is the first to satisfy them. **S3** captures *social* relationships between users, and between users and content, but also the *structure* present in rich social content, as well as its *semantics*. We provide the first top-*k* keyword search algorithm taking into account the social, structured, and semantic dimensions and formally establish its termination and correctness. Experiments on real social networks demonstrate the efficiency and qualitative advantage of our algorithm through the joint exploitation of the social, structured, and semantic dimensions of **S3**.

Key-words: semantic web, social network, database, top-k, RDF

* INRIA

† Université Paris-Sud

‡ Université Rennes 1

**RESEARCH CENTRE
SACLAY – ÎLE-DE-FRANCE**

1 rue Honoré d'Estienne d'Orves
Bâtiment Alan Turing
Campus de l'École Polytechnique
91120 Palaiseau

Recherche sur du contenu structuré, social et sémantique

Résumé : Les contenus sociaux comme les blogs, les tweets, les journaux en ligne etc. sont une source riche d'informations liées. Nous identifions dans ce rapport un ensemble de conditions nécessaires à une exploration pertinente de ce contenu riche et introduisons un nouvel modèle de données, **S3**, qui est le premier à les satisfaire. **S3** capte les relations *sociales* entre les utilisateurs et les contenus mais aussi la *structure* et la *sémantique* de ces derniers. Nous proposons aussi le premier algorithme de recherche top k qui prend en compte les dimensions structurelles, sociales et sémantiques et donnons une preuve formelle de sa correction et de sa terminaison. Une évaluation expérimentale sur des vrais réseaux sociaux valide l'efficacité et la qualité de notre approche sur l'exploration conjointe des dimensions structurelles, sociales et sémantiques de **S3**.

Mots-clés : web sémantique, réseaux sociaux, base de données, recherche top k, RDF

1 Introduction

The World Wide Web (or Web, in short) was designed for users to interact with each other by means of pages interconnected with hyperlinks. Thus, the Web is the earliest inception of an *online* social network (whereas “real-life” social networks have a much longer history in social sciences). However, the technologies and tools enabling large-scale online social exchange have only become available recently. A popular model of such exchanges features: *social network users*, who may be connected to one another, *data items*, and the possibility for users to *tag* data items, i.e., to attach to an item an annotation expressing the user’s view or classification of the item. Variants of this “user-item-tag” (UIT) model can be found e.g., in [15, 19, 30]. In such contexts, a user, called *seeker*, may ask a query, typically as a set of keywords. The problem then is to find the best query answers, taking into account both the relevance of items to the query, and the social proximity between the seeker and the items, based also on tags. Today’s major social networks e.g., Facebook [5], all implement some UIT variant. We identify a set of basic requirements which UIT meets:

R0. UIT models *explicit social connections* between users, e.g., u_1 is a friend of u_0 in Figure 1, to which we refer throughout this paper unless stated otherwise. It also captures *user endorsement (tags)* of data items, as UIT search algorithms *exploit both the user endorsement and the social connections* to return items most likely to interest the seeker, given his social and tagging behavior.

To fully exploit the content shared in social settings, we argue that the model used for such data (and, accordingly, the query model) must also satisfy the requirements below:

R1. The current wealth of publishing modes (through social networks, blogs, interlinked Web pages etc.) allows many different relations between items. For example, document d_1 *replies to* document d_0 (think for instance of opposite-viewpoint articles in a heated debate), while document d_2 *comments on* the paragraph of d_0 identified by the URI $d_{0.3.2}$. The model must capture **relations between items**, in particular since **they may lead to implicit relations between users**, according to their manipulations of items. For instance, the fact that u_2 posted d_1 as a reply to d_0 , posted by u_0 , entails that u_2 at least read d_0 , and thus some form of exchange has taken place between u_0 and u_2 ; if one looked for *explicit* social connections only, we would wrongly believe that u_0 and u_2 have no relation to each other.

R2. Items shared in social media often have a rich structured content. For instance, the article d_0 comprises many sections, and paragraphs, such as the one identified by the URI $d_{0.3.2}$. **Document structure must be reflected in the model** in order to return *useful* document fragments as query results, instead of a very large document or a very small snippet of a few words (e.g., exactly the search keywords). Document structure also helps discern when users have *really* interacted through content. For instance, u_3 has interacted with u_0 , since u_3 comments on the fragment $d_{0.3.2}$ of u_0 ’s article d_0 . In contrast, when user u_4 tags with “university” the fragment $d_{0.5.1}$ of d_0 , disjoint from $d_{0.3.2}$, u_4 may not even have read the same text as u_3 , thus the two likely did not interact.

R3. Item and tag semantics must be modeled. Social Web data encapsulates users’ knowledge on a multitude of topics; ontologies, either general such as DBPedia or Google’s Knowledge Base, or application-specific, can be leveraged to *give query answers which cannot be found without relying on semantics*. For instance, assume u_1 looks for information about *university graduates*: document d_1 states that u_2 holds a M.S. degree. Assume a knowledge base specifies that *a M.S. is a degree* and that *someone having a degree is a graduate*. The ability to return as result the snippet of d_1 most relevant to the query is directly conditioned by the ability to exploit the ontology (and the content-based interconnections along the path: u_1 friend of u_0 , u_0 posted d_0 , d_1 replied to d_0).

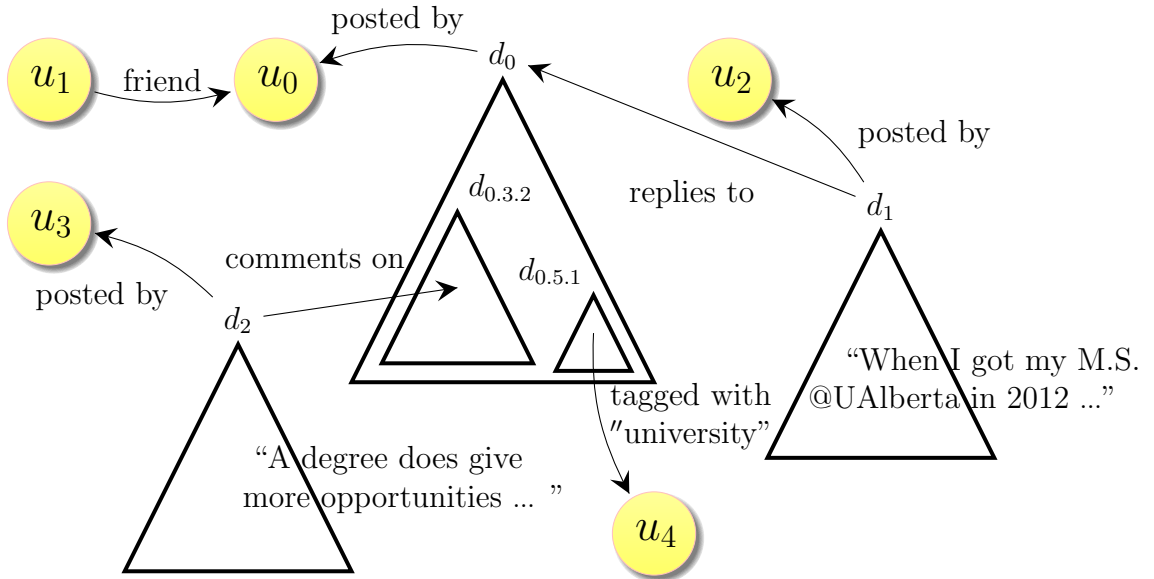


Figure 1: Motivating example.

R4. In many contexts, tagging may apply to tags themselves, e.g., in annotated corpora, where an annotation (tag) obtained from an analysis can further be annotated with provenance details (when and how the annotation was made) or analyzed in its turn. Information from higher-level annotations is obviously still related to the original document. The model should allow expressing **higher-level tags**, to exploit their information for query answering.

R5. The data model and queries should have **well-defined semantics**, to precisely characterize computed results, ensure correctness of the implementation, and allow for optimization.

R6. The model should be **generic** (not tied to a particular social network model), **extensible** (it should allow easy extension or customization, as social networks and applications have diverse and rapidly evolving needs), and **interoperable**, i.e., it should be possible to get richer / more complete answers by integrating different sources of social connections, facts, semantics, or documents. This ensures in particular independence from any proprietary social network viewpoint, usefulness in a variety of settings, and a desirable form of “monotonicity”: the more content is added to the network, the more its information value increases.

This work makes the following contributions.

1. We present **S3**, a novel *data model* for structured, semantic-rich content exchanged in social applications; it is the first model to meet the requirements **R0** to **R6** above.
2. We revisit *top-k social search for keyword queries*, to retrieve the most relevant *document fragments* w.r.t. the social, structural, and semantic aspects captured by **S3**. We identify a set of *desirable properties of the score function* used to rank results, provide a *novel query evaluation algorithm* called **S3_k** and *formally establish its termination and correctness*; the algorithm intelligently exploits the score properties to stop *as early as possible*, to return answers fast, with little evaluation effort. **S3_k** is the first to formally guarantee a specific result in a structured, social, and semantic setting.

U URIs	L literals	\mathcal{K} keywords	$Ext(k)$ extension of k
Ω users	D documents	T tags	I graph instance

Table 1: Main data model notations.

3. We implemented $S3_k$ based on a concrete score function (extending traditional ones from XML keyword search) and experimented with *three real social datasets*. We demonstrate the *feasibility* of our algorithm, and its *qualitative advantage* over existing approaches: it finds relevant results that would be missed by ignoring any dimension of the graph.

An $S3$ instance can be exploited in many other ways: through structured XML and/or RDF queries as in [8], searching for users, or focusing on annotations as in [3]; one could also apply graph mining etc. In this paper, we first describe the data model, and then revisit the top-k document search problem, since it is the most widely used (and studied) in social settings.

In the sequel, Section 2 presents the $S3$ data model, while Section 3 introduces a notion of generic score and instantiates it through a concrete score. Section 4 describes $S3_k$, we present experiments in Section 5, then discuss related works in Section 6 and conclude.

2 Data model

We now describe our model integrating social, structured, and semantic-rich content into a *single weighted RDF graph*, and based on a small set of *$S3$ -specific RDF classes and properties*. We present weighted RDF graphs in Section 2.1, and show how they model social networks in Section 2.2. We add to our model structured documents in Section 2.3, and tags and user-document interactions in Section 2.4; Section 2.5 introduces our notion of social paths. Table 1 recaps the main notations of our data model.

URIs and literals We assume given a set U of Uniform Resource Identifiers (URIs, in short), as defined by the standard [28], and a set of literals (constants) denoted L , disjoint from U .

Keywords We denote by \mathcal{K} the set of all possible *keywords*: it contains all the URIs, plus the stemmed version of all literals. For instance, stemming replaces “graduation” with “graduate”.

2.1 RDF

An *RDF graph* (or *graph*, in short) is a set of *triples* of the form $\mathbf{s} \mathbf{p} \mathbf{o}$, stating that the *subject* \mathbf{s} has the *property* \mathbf{p} and the value of that property is the *object* \mathbf{o} . In relational notation (Figure 2), $\mathbf{s} \mathbf{p} \mathbf{o}$ corresponds to the tuple (\mathbf{s}, \mathbf{o}) in the binary relation \mathbf{p} , e.g., u_1 hasFriend u_0 corresponds to $\text{hasFriend}(u_1, u_0)$. We consider every triple is *well-formed* [25]: its subject belongs to U , its property belongs to U , and its object belongs to \mathcal{K} .

Notations We use \mathbf{s} , \mathbf{p} , \mathbf{o} to denote a subject, property, and respectively, object in a triple. Strings between quotes as in “*string*” denote literals.

RDF types and schema The property type built in the RDF standard is used to specify to which *classes* a resource belongs. This can be seen as a form of resource typing.

A valuable feature of RDF is RDF Schema (RDFS), which allows enhancing the resource descriptions provided by RDF graphs. An RDF Schema declares *semantic constraints* between the classes and the properties used in these graphs, through the use of four RDF built-in properties. These constraints can model:

Constructor	Triple	Relational notation
Class assertion	$s \text{ type } o$	$o(s)$
Property assertion	$s \text{ p } o$	$p(s, o)$

Constructor	Triple	Relational notation
Subclass constraint	$s \prec_{sc} o$	$s \subseteq o$
Subproperty constraint	$s \prec_{sp} o$	$s \subseteq o$
Domain typing constraint	$s \leftrightarrow_d o$	$\Pi_{\text{domain}}(s) \subseteq o$
Range typing constraint	$s \leftrightarrow_r o$	$\Pi_{\text{range}}(s) \subseteq o$

Figure 2: RDF (top) and RDFS (bottom) statements.

- subclass relationships, which we denote by \prec_{sc} ; for instance, any *M.S.Degree* is also a *Degree*;
- subproperty relationships, denoted \prec_{sp} ; for instance, *workingWith* someone also means being *acquaintedWith* him;
- typing of the first attribute (or domain) of a property, denoted \leftrightarrow_d , e.g., the domain of *hasDegreeFrom* is a *Graduate*;
- typing of the second attribute (or range) of a property, denoted \leftrightarrow_r , e.g., the range of *hasDegreeFrom* is an *University*.

Figure 2 shows the constraints we use, and how to express them. In this figure, domain and range denote respectively the first and second attributes of a property. The figure also shows the relational notation for these constraints, which in RDF are interpreted under the open-world assumption [1], i.e., as *deductive constraints*. For instance, if a graph includes the triples *hasFriend* \leftrightarrow_d *Person* and u_1 *hasFriend* u_0 , then the triple u_1 type *Person* holds in this graph even if it is not explicitly present. This *implicit* triple is due to the \leftrightarrow_d constraint in Figure 2.

Saturation RDF entailment is the RDF reasoning mechanism that allows making explicit all the implicit triples that hold in an RDF graph \mathcal{G} . It amounts to repeatedly applying a set of normative immediate *entailment* rules (denoted \vdash_{RDF}^i) on \mathcal{G} : given some triples explicitly present in \mathcal{G} , a rule adds some triples that directly follow from them. For instance, continuing the previous example,

$$u_1 \text{ hasFriend } u_0, \text{ hasFriend } \leftrightarrow_r \text{ Person } \vdash_{\text{RDF}}^i u_0 \text{ type Person}$$

Applying immediate entailment \vdash_{RDF}^i repeatedly until no new triple can be derived is known to lead to a unique, finite fixpoint graph, known as the *saturation* (a.k.a. closure) of \mathcal{G} . RDF entailment is part of the RDF standard itself: the answers to a query on \mathcal{G} must take into account all triples in its saturation, since *the semantics of an RDF graph is its saturation* [25].

In the following, we assume, without loss of generality, that all RDF graphs are saturated; many saturation algorithms are known, including incremental [9] or massively parallel ones [24].

Weighted RDF graph Relationships between documents, document fragments, comments, users, keywords etc. naturally form a graph. We encode each edge from this graph by a *weighted RDF triple* of the form (s, p, o, w) , where (s, p, o) is a regular RDF triple, and $w \in [0, 1]$ is termed the *weight* of the triple. Any triple whose weight is not specified is assumed to be of weight 1.

We define the saturation of a weighted RDF graph as the saturation derived *only from its triples whose weight is 1*. Any entailment rule of the form $a, b \vdash_{\text{RDF}}^i c$ applies only if the weight

of a and b is 1; in this case, the entailed triple c also has the weight 1. We restrict inference in this fashion to distinguish triples which certainly hold (such as: “a M.S. is a degree”, “ u_1 is a friend of u_0 ”) from others whose weight is computed, and carries a more quantitative meaning, such as “the similarity between d_0 and d_1 is 0.5”¹.

Graph instance I and S3 namespace We use I to designate the weighted RDF instance we work with. The RDF Schema statements in I allow a semantic interpretation of keywords, as follows:

Definition 2.1 (Keyword extension). *Given an S3instance I and a keyword $k \in \mathcal{K}$, the extension of k , denoted $Ext(k)$, is defined as follows:*

- $k \in Ext(k)$
- for any triple of the form $b \text{ type } k$, $b \prec_{sc} k$ or $b \prec_{sp} k$ in I, we have $b \in Ext(k)$.

For example, given the keyword *degree*, and assuming that M.S. \prec_{sc} degree holds in I, we have $M.S. \in Ext(\text{degree})$. *The extension of k does not generalize it*, in particular it does not introduce any loss of precision: whenever k' is in the extension of k , the RDF schema in I ensures that k' is an *instance*, or a *specialization* (particular case) of k . This is in coherence with the principles behind the RDF schema language².

For our modeling purposes, we define below a small set of RDF classes and properties used in I; these are shown prefixed with the S3 namespace. The next sections show how I is populated with triples derived from the users, documents and their interactions.

2.2 Social network

We consider a set of social network users $\Omega \subset U$, i.e., each user is identified by a URI. We introduce the special RDF class S3:user, and for each user $u \in \Omega$, we add: $u \text{ type } S3:user \in I$.

To model the relationships between users, such as “friend”, “co-worker” etc., we introduce the special property S3:social, and model any concrete relationship between two users by a triple whose property specializes S3:social. Alternatively, one may see S3:social as the *generalization of all social network relationships*.

Weights are used to encode the strength w of each relationship going from a user u_1 to a user u_2 : $u_1 \text{ S3:social } u_2 \ w \in I$. As customary in social network data models, the higher the weight, the closer we consider the two users to be.

Extensibility Depending on the application, it may be desirable to consider that two users satisfying some condition are involved in a social interaction. For instance, if two people have worked the same year for a company of less than 10 employees (such information may be in the RDF part of our instance), they must have *worked together*, which could be a social relationship. This is easily achieved with a query that retrieves all such user pairs (in SPARQL or in a more elaborate language [8] if the condition also carries over the documents), and builds a $u \text{ workedWith } u'$ triple for each such pair of users. Then it suffices to add these triples to the instance, together with the triple: $\text{workedWith } \prec_{sp} \text{ S3:social}$.

¹One could generalize this to support inference over triples of any weight, leading to e.g., “ u_1 is of type Person with a weight of 0.5”, in the style of probabilistic databases.

²One could also allow a keyword $k' \in Ext(k)$ which is only close to (but not a specialization of) k , e.g., “student” in $Ext(\text{“graduate”})$, at the cost of a loss of precision in query results. We do not pursue this alternative here, as we chose to follow standard RDF semantics.

2.3 Documents and fragments

We consider that content is created under the form of structured, tree-shaped *documents*, e.g., XML, JSON, etc. A document is an unranked, ordered tree of *nodes*. Let N be a set of node names (for instance, the set of allowed XML element and attribute names, or the set of node names allowed in JSON). Any node has a *URI*. We denote by $D \subset U$ the set of all node URIs. Further, each node has a *name* from N , and a *content*, which we view as *a set of keywords* from \mathcal{K} : we consider each text appearing in a document has been broken into words, stop words have been removed, and the remaining words have been stemmed to obtain our version of the node’s text content. For example, in Figure 1, the text of d_1 might become {"M.S.", "UAlberta", "2012"}.

We term any subtree rooted at a node in document d a *fragment* of d , implicitly defined by the URI of its root node. The set of fragments (nodes) of a document d is denoted $Frag(d)$. We may use f to refer interchangeably to a fragment or its URI. If f is a fragment of d , we say d is an *ancestor* of f .

To simplify, *we use document and fragment interchangeably*; both are identified by the URI of their unique root node.

Document-derived triples We capture the *structural relationships* between documents, fragments and keywords through a set of RDF statements using S3-specific properties. We introduce the RDF class S3:doc corresponding to the documents, and we translate:

- each $d \in D$ into the I triple d type S3:doc;
- each document $d \in D$ and fragment rooted in a node n of d into n S3:partOf d ;
- each node n and keyword k appearing in the content of n into n S3:contains k ;
- each node n whose name is m , into n S3:nodeName m .

Example 2.1. *Based on the sample document shown in Figure 1, the following triples are part of I:*

$$\begin{array}{ll} d_{0.3.2} \text{ S3:partOf } d_{0.3} & d_1 \text{ S3:contains "M.S."} \\ d_{0.3} \text{ S3:partOf } d_0 & d_1 \text{ S3:nodeName text} \end{array}$$

The following constraints, part of I, model the natural relationships between the S3:doc class and the properties introduced above:

$$\begin{array}{ll} \text{S3:partOf} \leftrightarrow_d \text{ S3:doc} & \text{S3:partOf} \hookrightarrow_r \text{ S3:doc} \\ \text{S3:contains} \leftrightarrow_d \text{ S3:doc} & \text{S3:nodeName} \leftrightarrow_d \text{ S3:doc} \end{array}$$

which read: the relationship S3:partOf connects pairs of fragments (or documents); S3:contains describes the content of a fragment; and S3:nodeName associates names to fragments.

Fragment position We will need to assess how closely related a given fragment is to one of its ancestor fragments. For that, we use a function $pos(d, f)$ which returns the *position* of fragment f within document d . Concretely, pos can be implemented for instance by assigning Dewey-style IDs to document nodes, as in [16, 20]. Then, $pos(d, f)$ returns the list of integers (i_1, \dots, i_n) such that the path starting from d ’s root, then moving to its i_1 -th child, then to this node’s i_2 -th child etc. ends in the root of the fragment f . For instance, in Figure 1, $pos(d_{0.3.2}, d_0)$ may be (3, 2).

Example 2.2. *Considering again Figure 1, sample outputs of the pos function are:*

$$\begin{array}{l} pos(d_0, d_{0.3.2}) = [3, 2] \\ pos(d_1, d_2) = [] \end{array}$$

2.4 Relations between structure, semantics, users

We now show how dedicated S3 classes and properties are used to encode all kinds of connections between users, content, and semantics in a single S3 instance.

Tags A typical user action in a social setting is to *tag* a data item, reflecting the user’s opinion that the item is related to some concept or keyword used in the tag. We introduce the special class S3:relatedTo to *account for the multiple ways in which a user may consider that a fragment is related to a keyword*. We denote by T the set of all tags.

For example, in Figure 1, u_4 tags $d_{0.5.1}$ with the keyword “university”, leading to the triples:

$$\begin{array}{ll} \text{a type S3:relatedTo} & \text{a S3:hasSubject } d_{0.5.1} \\ \text{a S3:hasKeyword "university"} & \text{a S3:hasAuthor } u_4 \end{array}$$

In this example, a is a *tag* (or annotation) resource, encapsulating the various tag properties: its content, who made it, and on what. The tag subject (the value of its S3:hasSubject property) is either a document or another tag. The latter allows to express *higher-level annotations*, when an annotation (tag) can itself be tagged.

A tag may lack a keyword, i.e., it may have no S3:hasKeyword property. Such no-keyword tags model *endorsement* (support), such as **like** on Facebook, **retweet** on Twitter, or **+1** on Google+.

Tagging may differ significantly from one social setting to another. Thus, in a restaurant rating site, a tag ranges from \star (terrible) to $\star\star\star\star$ (excellent); in a collaborative question answering site, users tag questions with one of the existing discussion topics, e.g., “clock” and “cpu-speed” for a question related to CPU overclocking etc. Tags may also be produced by programs, e.g., a natural language processing (NLP) tool may recognize a text fragment related to a person. Just like the S3:social property can be specialized to model arbitrary social connections between users, subclasses of S3:relatedTo can be used to model different kinds of tags. For instance, assuming a_2 is a tag produced by a NLP software, this leads to the I triples:

$$\begin{array}{l} a_2 \text{ type NLP:recognize} \\ \text{NLP:recognize} \prec_{sc} \text{ S3:relatedTo} \end{array}$$

User actions on documents Users *post* (or *author*, or *publish*) content, modeled by the dedicated property S3:postedBy. Some of this content may be *comments* on (or *replies* / *answers* to) other fragments; this is encoded via the property S3:commentsOn. When user u posts document c , which comments on document d and possibly cites part of it, each fragment copied (cited as such) from d into c is now part of c and thus has a new URI.

Example 2.3. In Figure 1, d_2 is posted by u_3 , as a comment on $d_{0.3.2}$, leading to the following I triples:

$$d_2 \text{ S3:postedBy } u_3 \quad d_2 \text{ S3:commentsOn } d_{0.3.2}$$

As before, we view any concrete relation between documents e.g., *answers to*, *retweets*, *comments on*, *is an old version of* etc. as a specialization (sub-property) of S3:commentsOn; the corresponding connections lead to implicit S3:commentsOn triples, as explained in Section 2.1. Similarly, forms of authorship connecting users to their content are modeled by specializing S3:postedBy. This allows integrating (querying together) many social networks over partially overlapping sets of URIs, users and keywords.

Inverse properties As syntactic sugar, to simplify the traversal of connections between users and documents, we introduce a set of *inverse properties*, denoted respectively S3:postedBy,

Class	Semantics
S3:user	the users (the set of its instances is Ω)
S3:doc	the documents (the set of its instances is D)
S3:relatedTo	generalization of item “tagging” with keywords (the set of all instances of this class is T : the set of tags)
Property	Semantics
S3:postedBy	connects users to the documents they posted
S3:commentsOn	connects a comment with the document it is about
S3:partOf	connects a fragment to its parent nodes
S3:contains	connects a document with the keyword(s) it contains
S3:nodeName	asserts the name of the root node of document
S3:hasSubject	specifies the subject (document or tag) of a tag
S3:hasKeyword	specifies the keyword of a tag
S3:hasAuthor	specifies the poster of a tag
S3:social	generalization of social relationships in the network

Table 2: Classes and properties in the S3 namespace.

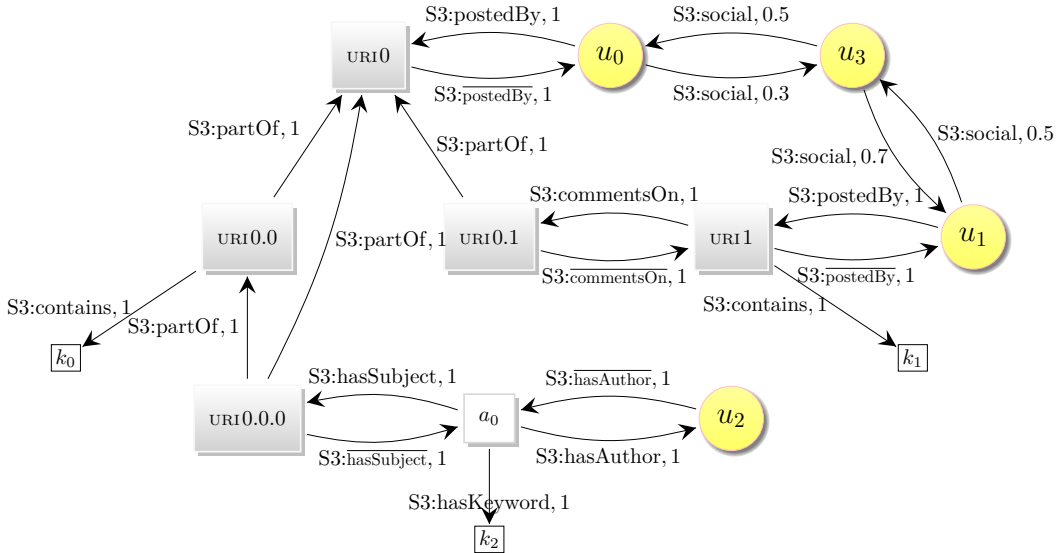


Figure 3: Sample S3 instance I.

$S3:\overline{commentsOn}$, $S3:\overline{hasSubject}$ and $S3:\overline{hasAuthor}$, with the straightforward semantics: $s \overline{p} o \in I$ iff $o p s \in I$ where \overline{p} is the inverse property of p . For instance, $u_0 S3:\overline{friend} u_1$ in Figure 1.

Table 2 summarises the above S3 classes and properties, while Figure 3 illustrates an I instance.

2.5 Social paths

We define here social paths on I (established either through explicit social links, or through user interactions). We call **network edges** those I edges encapsulating quantitative information on the links between user, documents and tags, i.e., *the set of edges whose properties are in the*

namespace `S3` other than `S3:partOf`, and whose subjects and objects are either users, documents, or tags.

Definition 2.2 (Network edges). *We define:*

$$I_{\text{net}} = \{\mathbf{s} \ \mathbf{p} \ \mathbf{o} \ \mathbf{w} \mid \mathbf{p} \in \text{S3} \setminus \{\text{S3:partOf}\}, \mathbf{s}, \mathbf{o} \in (\Omega \cup D \cup T)\}$$

For instance, in Figure 3, `u1 S3:social u3 0.5` and `URI0 S3:postedBy u1` are network edges; `URI0.0 S3:contains k0` and `URI0.1 S3:partOf URI0` are not. The intuition behind the exclusion of `S3:partOf` is that *structural relations between fragments, or between fragments and keywords, merely describe data content and not an interaction*. However, if two users comment on the same fragment, or one comments on a fragment of a document posted by the other (e.g., `u2` and `u0` in Figure 1), this is form of social interaction.

When two users interact with unrelated fragments of a same document, such as `u3` and `u4` on disjoint subtrees of `d0`, this does not establish a social link between `u3` and `u4`, since they may not even have read the same text³. We introduce:

Definition 2.3 (Document vertical neighborhood). *Two documents are vertical neighbors if one of them is a fragment of the other. The function `neigh`: $U \rightarrow 2^U$ returns the set of vertical neighbors of an URI.*

In Figure 3, `URI0` and `URI0.0.0` are vertical neighbors, so are `URI0` and `URI0.1`, but `URI0.0.0` and `URI0.1` are not. In the sequel, due to the strong connections between nodes in the same vertical neighborhood, **we consider (when describing and exploiting social paths) that a path entering through any of them can exit through any other**; a vertical neighborhood acts like a single node only and exactly from the perspective of a social path⁴. We can now define social paths:

Definition 2.4 (Social path). *A social path (or simply a path) in I is a chain of network edges such that the end of each edge and the beginning of the next one are either the same node, or vertical neighbors.*

We may also designate a path simply by *the list of nodes it traverses*, when the edges taken along the path are clear. In Figure 3, `u2` $\xrightarrow{\text{u2 S3:hasAuthor a0 1}}$ `a0` $\xrightarrow{\text{a0 S3:hasSubject URI0.0.0 1}}$ `URI0.0.0` \dashrightarrow `URI0` $\xrightarrow{\text{URI0 S3:postedBy u0 1}}$ `u0` is an example of such a path (the dashed line: `URI0.0.0` \dashrightarrow `URI0`, is not an edge in the path but a connection between vertical neighbors, `URI0.0.0` been the end of an edge and `URI0` the beginning of the next edge). Also, in this Figure, there is no social path going from `u2` to `u1` avoiding `u0`, because it is not possible to move from `URI0.1` to `URI0.0.0` through a vertical neighborhood.

Social path notations The set of *all social paths from a node x (or one of its vertical neighbours) to a node y (or one of its vertical neighbors)* is denoted $x \rightsquigarrow y$. The length of a path p is denoted $|p|$. The restriction of $x \rightsquigarrow y$ to paths of length exactly n is denoted $x \rightsquigarrow_n y$, while $x \rightsquigarrow_{\leq n} y$ holds the paths of at most n edges.

Path normalization To harmonize the weight of each edge in a path depending on its importance, we introduce path normalization, which modifies the weights of a path's edge as follows. Let n be the ending point of a social edge in a path, and e be the *next* edge in this path. The normalized weight of e for this path, denoted $e.n_w$, is defined as:

³To make such interactions count as social paths would only require simple changes to the path normalization introduced below.

⁴In other contexts, e.g., to determine their relevance w.r.t. a query, vertical neighbors are considered separately.

$$e.n_w = e.w / \sum_{e' \in \text{out}(\text{neigh}(n))} e'.w$$

where $e.w$ is the weight of e , and $\text{out}(\text{neigh}(n))$ the set of network edges outgoing from any vertical neighbor of n . This normalizes the weight of e w.r.t. the weight of edges outgoing from any vertical neighbor of n . Observe that $e.n_w$ depends on n , however e does not necessarily start in n , but in any of its vertical neighbors. Therefore, $e.n_w$ indeed depends on the path (which determines the vertical neighbor n of e 's entry point).

In the following, we assume all social paths are normalized.

Example 2.4. In Figure 3, consider the path:

$$p = u_0 \xrightarrow{u_0 \text{ S3:postedBy URI0 } 1} \text{URI0} \dashrightarrow \text{URI0.0.0} \xrightarrow{\text{URI0.0.0 S3:hasSubject } a_0 } a_0$$

Its first edge is normalized by the edges leaving u_0 : one leading to URI0 (weight 1) and the other leading to u_3 (weight 0.3). Thus, its normalised weight is $1/(1 + 0/3) = 0.77$.

Its second edge exits URI0.0.0 after a vertical neighborhood traversal $\text{URI0} \dashrightarrow \text{URI0.0.0}$. It is normalized by the edges leaving $\text{neigh}(\text{URI0})$, i.e., all the edges leaving a fragment of URI0 . Its normalised weight is $1/(1 + 1 + 1 + 1) = 0.25$.

S3 meets the requirements from Section 1, as follows. Genericity, extensibility and interoperability (**R6**) are guaranteed by the reliance on the Web standards RDF (Section 2.1) and XML/JSON (Section 2.3). These enable specializing the **S3** classes and properties, e.g., through application-dependent queries (see Extensibility in Section 2.2). Our document model (Section 2.3) meets requirement **R2**; the usage of RDF (Section 2.1) ensures **R3**, while the relationships introduced in Section 2.4 satisfy **R1** as well as **R4** (higher-level tags). For what concerns **R5** (formal semantics), the data model has been described above; we consider queries next.

3 Querying an S3 instance

Users can search **S3** instances through keyword queries; the answer consists of the k top-score fragments, according to a joint structural, social, and semantic score. Section 3.1, defines queries and their answers. After some preliminaries, we introduce a *generic score*, which can be instantiated in many ways, and a set of *feasibility conditions* on the score, which suffice to ensure the termination and correctness of our query answering algorithm (Section 3.3). We present our concrete score function in Section 3.4.

3.1 Queries

S3 instances are queried as follows:

Definition 3.1 (Query). A query is a pair (u, ϕ) where u is a user and ϕ is a set of keywords.

We call u the *seeker*. We define the top- k answers to a query as the k documents or fragments thereof with the highest scores, further satisfying the following constraint: the presence of a document or fragment at a given rank precludes the inclusion of its vertical neighbors at lower ranks in the results⁵.

As customary, top- k answers are ranked using a score function: $s(q, d)$ returns for a document d and query q a value in \mathbb{R} , based on the graph I . The top k results are recursively defined as the top $k - 1$ best results, plus the best among the documents which are neither fragments nor ancestors of the $k - 1$ best results.

⁵This assumption is standard in XML keyword search, e.g., [4].

Definition 3.2 (Query answer). A top- k answer to the query q using the score s , denoted $T_{k,s}(q)$, is defined inductively over k as follows:

- $T_{0,s}(q) = \emptyset$
- if $k > 0$, $T_{k,s}(q)$ contains exactly $T_{k-1,s}(q)$ plus one document from $\underset{d \in D \setminus \text{neigh}(T_{k-1,s}(q))}{\text{argmax}}(s(q, d))$.

Here, argmax is not necessarily unique: there might be ties. If a group of documents have identical scores, any of them can be in the query answer (ties are broken in some non-deterministic way). For instance, if the graph consists of four documents scored $(d_1, 0.8)$, $(d_2, 0.6)$, $(d_3, 0.6)$, $(d_4, 0.5)$ and we ask for the top 2 results, $\{d_1, d_2\}$ as well as $\{d_1, d_3\}$ are answers.

3.2 Connecting query keywords and documents

Answering queries over I requires finding best-scoring documents, based on the *direct and indirect connections* between documents, the seeker, and search keywords. The connection can be direct, for instance, when the document contains the keyword, or indirect, when a document is connected by a chain of relationships to a search keyword k , or to some keyword from k 's extension.

We denote the **set of direct and indirect connections between a document d and a keyword k** by $\text{con}(d, k)$. It is a set of three-tuples $(\text{type}, \text{frag}, \text{src})$ such that:

- $\text{type} \in \{\text{S3:contains}, \text{S3:relatedTo}, \text{S3:commentsOn}\}$ is the **type** of the connection,
- $f \in \text{Frag}(d)$ is the **fragment** of d (possibly d itself) due to which d is involved in this connection,
- $\text{src} \in \Omega \cup D$ (users or documents) is the **source** (origin) of this connection (see below).

Below we describe the possible situations which create connections. Let d, d' be documents or tags, and f, f' be fragments of d and d' , respectively⁶. Further, let k, k' be keywords such that $k' \in \text{Ext}(k)$, and $\text{src} \in \Omega \cup D$ be a user or a document.

Documents connected to the keywords of their fragments If the fragment f contains a keyword k , i.e., $f \text{ S3:contains } k \in I$, then:

$$(\text{S3:contains}, f, d) \in \text{con}(d, k)$$

which reads: “ d is connected to k through a S3:contains relationship due to f ”. This connection holds even if f contains not k itself, but some $k' \in \text{Ext}(k)$. For example, in Figure 1, if the keyword “university” appears in the fragment whose URI is $d_{2.7.5}$, then $\text{con}(d_2, \text{“university”})$ includes $(\text{S3:contains}, d_{2.7.5}, d_2)$. Observe that a given k' and f may lead to many connections, if k' specializes several keywords and/or if f has many ancestors.

Connections due to tags For every tag a of the form

$$\begin{array}{ll} \text{a type S3:relatedTo} & \text{a S3:hasSubject f} \\ \text{a S3:hasAuthor src} & \text{a S3:hasKeyword k'} \end{array}$$

⁶We here slightly extend notations, since tags do not have fragments: if d is a tag, we consider that its only fragment is d .

$con(d, k)$ includes (S3:relatedTo, f, src). In other words, whenever a fragment f of d is tagged by a source src with a specialization of the keyword k , this leads to a S3:relatedTo connection between d and k due to f , whose source is the tag author src . For instance, the tag a of u_4 in Figure 1 creates the connection (S3:relatedTo, $d_{0.5.1}, u_4$) between d_0 and “university”.

More generally, if a tag a on fragment f has *any type of connection* (not just S3:hasKeyword) to a keyword k due to source src , this leads to a connection (S3:relatedTo, f, src) between d and k . The intuition is that the tag adds its connections to the tagged fragment and, transitively, to its ancestors. (As the next section shows, the importance given to such connections decreases as the distance between d and f increases.)

If the tag a on f is a simple endorsement (it has no keyword), the tag inherits d ’s connections, as follows. Assume d has a connection of type $type$ to a keyword k : then, a also has a $type$ connection to k , whose source is src , the tag author. The intuition is that when src endorses (likes, +1s) a fragment, src agrees with its content, and thus connects the tag, to the keywords related to that fragment and its ancestors. For example, if a user u_5 endorsed d_0 in Figure 1 through a no-keyword tag a_5 , the latter tag is related to “university” through: (S3:relatedTo, $d_{0.5.1}, u_5$).

Connections due to comments When a comment on f is connected to a keyword, this also connects any ancestor d of f to that keyword; the connection source carries over, while the type of d ’s connection is S3:commentsOn. For instance, in Figure 1, since d_2 is connected to “university” through (S3:contains, $d_{2.7.5}, d_2$) and since d_2 is a comment on $d_{0.3.2}$, it follows that d_0 is also related to “university” through (S3:commentsOn, $d_{0.3.2}, d_2$).

Formally, whenever d ’s S3:commentsOn f and there exist $type', frag', src$ such that $(type', frag', src) \in con(d', k)$, we have:

$$(S3:commentsOn, f, src) \in con(d, k)$$

3.3 Generic score model

We introduce a set of proximity notions, based on which we state the conditions to be met by a score function, for our query evaluation algorithm to compute a top-k query answer.

Path proximity We consider a measure of proximity *along one path*, denoted \overrightarrow{prox} , between 0 and 1 for any path, such that:

- $\overrightarrow{prox}(()) = 1$, i.e., the proximity is maximal on an empty path (in other words, from a node to itself),
- for any two paths p_1 and p_2 , such that the start point of p_2 is in the vertical neighborhood of the end point of p_1 :

$$\overrightarrow{prox}(p_1 || p_2) \leq \min(\overrightarrow{prox}(p_1), \overrightarrow{prox}(p_2)),$$

where $||$ denotes path concatenation. This follows the intuition that proximity along a concatenation of two paths is at most the one along each of these two components paths: proximity can only decrease as the path gets longer.

Social proximity associates to two vertices connected by at least one social path, a comprehensive measure over *all the paths* between them. We introduce such a global proximity notion, because different paths traverse different nodes, users, documents and relationships, all of which may impact the relation between the two vertices. Considering all the paths gives a *qualitative* advantage to our algorithm, since it enlarges its knowledge to the types and strength of all connections between two nodes.

Definition 3.3 (Social proximity). *The social proximity measure $prox : (\Omega \cup D \cup T)^2 \rightarrow [0, 1]$, is an aggregation along all possible paths between two users, documents or tags, as follows:*

$$prox(a, b) = \oplus_{path}(\{\overrightarrow{prox}(p), |p|, p \in a \rightsquigarrow b\}),$$

where $| \cdot |$ is the number of vertices in a path and \oplus_{path} is a function aggregating a set of values from $[0, 1] \times \mathbb{N}$ into a single scalar value.

Observe that the set of all paths between two nodes may be infinite, if the graph has cycles; this is often the case in social graphs. For instance, in Figure 3, a cycle can be closed between $(u_0, \text{URI0}, u_0)$. Thus, in theory, the score is computed over a potentially infinite set of paths. However, in practice, our algorithm works with *bounded social proximity* values, relying only on paths of a bounded length:

$$prox^{\leq n}(a, b) = \oplus_{path}(\{\overrightarrow{prox}(p), |p|, p \in a \rightsquigarrow_{\leq n} b\})$$

Based on the proximity measure, and the connections between keywords and documents introduced in Section 3.2, we define:

Definition 3.4 (Generic score). *Given a document d and a query $q = (u, \phi)$, the score of d for q is:*

$$score(d, (u, \phi)) = \oplus_{gen}(\{(k, type, pos(d, f), prox(u, src)) \\ |k \in \phi, (type, f, src) \in con(d, k)\})$$

where \oplus_{gen} is a function aggregating a set of (keyword, relationship type, importance of fragment f in d , social proximity) tuples into a value from $[0, 1]$.

Importantly, the above score *reflects the semantics, structure, and social content of the S3 instance*, as follows.

First, \oplus_{gen} aggregates over the keywords in ϕ . Recall that tuples from $con(d, k)$ account not only for k but also for keywords $k' \in Ext(k)$. This is how semantics is injected into the score.

Second, the score of d takes into account the relationships between fragments f of d , and keywords k , or $k' \in Ext(k)$, by using the sequence $pos(d, f)$ (Section 2.3) as an indication of the structural importance of the fragment within the document. If the sequence is short, the fragment is likely a large part of the document. Document structure is therefore taken into account here both *directly* through pos , and *indirectly*, since the con tuples also propagate relationships from fragments to their ancestors (Section 3.2).

Third, the score takes into account the social component of the graph through $prox$: this accounts for the relationships between the seeker u , and the various parties (users, documents and tags), denoted src , due to which f may be relevant for k .

Feasibility properties For our query answering algorithm to converge, the generic score model must have some properties which we describe below.

1. Relationship with path proximity This refers to the relationship between path proximity and score. First, the score should only *increase* if one adds *more paths* between a seeker and a data item. Second, the contribution of the paths of length $n \in \mathbb{N}$ to the social proximity can be expressed using the contributions of shorter “prefixes” of these paths, as follows. We denote by $ppSet^n(a, b)$ the set of the path proximity values for all paths of length n going from a to b :

$$ppSet^n(a, b) = \{\overrightarrow{prox}(p) \mid p \in a \rightsquigarrow_n b\}$$

Then, the first property is that there exists a function U_{prox} with values in $[0, 1]$, taking as input (i) the bounded social proximity for path of length at most $n - 1$, (ii) the proximity along paths of length n , and (iii) the length n , and such that:

$$prox^{\leq n}(a, b) = prox^{\leq n-1}(a, b) + U_{prox}(prox^{\leq n-1}(a, b), ppSet^n(a, b), n)$$

2. Long paths attenuation The influence of social paths should decrease as they get longer; intuitively, the farther away two items are, the weaker their connection and thus their influence on the score. More precisely, there exists a bound $B_{prox}^{>n}$ tending to 0 as n grows, and such that:

$$B_{prox}^{>n} \geq prox - prox^{\leq n}$$

3. Score soundness The score of a document should be positively correlated with the social proximity from the seeker to the document fragments that are relevant for the query.

Denoting $score_{[g]}$ the score where the proximity function $prox$ is replaced by a continuous function g having the same domain $(\Omega \cup D \cup T)^2$, $g \mapsto score_{[g]}$ must be monotonically increasing and continuous for the uniform norm.

4. Score convergence This property bounds the score of a document and shows how it relates to the social proximity. It requires the existence of a function B_{score} which takes a query $q = (u, \phi)$ and a number $B \geq 0$, known to be an upper bound on the social proximity between the seeker and any source: for any d , query keyword k , and $(type, f, src) \in con(d, k)$, we know that $prox(u, src) \leq B$. B_{score} must be positive, and satisfy, for any q :

- for any document d , $score(d, q) \leq B_{score}(q, B)$;
- $\lim_{B \rightarrow 0} (B_{score}(q, B)) = 0$ (tends to 0 like B).

We describe a concrete *feasible score*, i.e., having the above properties, in the next section.

3.4 Concrete score

We start by instantiating \overrightarrow{prox} , $prox$ and $score$.

Social proximity Given a path p , we define $\overrightarrow{prox}(p)$ as the product of the normalized weights (recall Section 2.5) found along the edges of p . We define our concrete social proximity function $prox(a, b)$ as a weighted sum over all paths from a to b :

$$prox(a, b) = C_\gamma \times \sum_{p \in a \rightsquigarrow b} \frac{\overrightarrow{prox}(p)}{\gamma^{|p|}}$$

where $\gamma > 1$ is a scalar coefficient, and $C_\gamma = \frac{\gamma-1}{\gamma}$ is introduced to ensure that $prox \leq 1$. Recall that by Definition 3.3, $prox$ requires a \oplus_{path} aggregation over the (social proximity, length) pairs of the paths between the two nodes. Hence, this concrete social proximity corresponds to choosing:

$$\oplus_{path}(S) = C_\gamma \times \sum_{(sp, len) \in S} \frac{sp}{\gamma^{len}}$$

where (sp, len) is a (social proximity, length) pair from its input.

Example 3.1. Social proximity Let us consider in Figure 3 the social proximity from u_0 to URI0 , using the $\overrightarrow{\text{prox}}$ and \oplus_{path} previously introduced. An edge connects u_0 directly to URI0 , leading to the normalized path p :

$$p = u_0 \xrightarrow{u_0 \text{ S3:postedBy URI0 } \frac{1}{1+0.3}} \text{URI0}$$

which accounts for a partial social proximity:

$$\text{prox}^{\leq 1}(u_0, \text{URI0}) = \frac{\overrightarrow{\text{prox}}(p)}{\gamma^{|p|}} = \frac{1/(1+0.3)}{\gamma^1}$$

Score function We define a simple concrete **S3 score function** which, for a document d , is the product of the scores of each query keyword in d . The score of a keyword is summed over all the connections between the keyword and the document. The weight for a given connection and keyword only depends on the *social distance between the seeker and the sources of the keyword*, and the *structural distance between the fragment involved in this relation and d* , namely the length of $\text{pos}(d, f)$. Both distances decrease exponentially as the path length grows. Formally:

Definition 3.5 (S3_k score). Given a query (u, ϕ) , the S3_k score of a document d for the query is defined as:

$$\text{score}(d, (u, \phi)) = \prod_{k \in \phi} \left(\sum_{(type, f, src) \in \text{con}(d, k)} \eta^{|\text{pos}(d, f)|} \times \text{prox}(u, src) \right)$$

for some damping factor $\eta < 1$.

Recall from Definition 3.4 that an aggregation function \oplus_{gen} combines the contributions of (keyword, relationship type, importance, social proximity) tuples in the score. The above definition corresponds to the following \oplus_{gen} aggregator:

$$\oplus_{\text{gen}}(S) = \prod_{k \in \phi} \left(\sum_{\substack{rel, \text{prox} \\ \exists type, (k, type, rel, \text{prox}) \in S}} \eta^{|\text{rel}|} \times \text{prox} \right)$$

Note that if we ignore the social aspects and restrict ourselves to top- k search on documents (which amounts to $\text{prox} = 1$), \oplus_{gen} gives the best score to the lowest common ancestor (LCA) of the nodes containing the query keywords. Thus, our score extends typical XML IR works, e.g., [4] (see also Section 6).

Obviously, there are many possible ways to define \oplus_{gen} and \oplus_{path} , depending on the application. In particular, different types of connections may not be accounted for equally; our algorithm only requires a *feasible score* (with the feasibility properties).

Theorem 3.1 (Score feasibility). *The S3_k score function (Definition 3.5) has the feasibility properties (Section 3.3).*

Proof. (Theorem 3.1)

1. Relationship with path proximity The S3_k score function satisfies this constraint by choosing U_{prox} as follows:

$$\begin{aligned} & U_{\text{prox}}(\text{prox}^{\leq n-1}(a, b), \text{ppSet}^n(a, b), n) \\ &= \sum_{\text{pprox} \in \text{ppSet}^n(a, b)} \frac{\text{pprox}}{\gamma^n} \end{aligned}$$

2. Long paths attenuation We show that long paths have a limited influence in the social proximity by exhibiting the bound based on our concrete functions. We start by stating the following lemma.

Lemma 3.1. *For any $i \in \mathbb{N}$, $u \in \mathbf{I}$, and set of paths P_i of length exactly i starting in u ,*

$$\sum_{p \in P_i} \overrightarrow{prox}(p) \leq 1.$$

This is done by induction on $i \in \mathbb{N}$, using the expression of \overrightarrow{prox} and the path normalization:

- If $i = 0$, then there is by convention only one (empty) path of length 0 from u , which has a \overrightarrow{prox} of 1 (computed as a product over an empty set).
- If $i > 1$, then any path in P_i can be decomposed as a path of length $i - 1$ starting in u followed by a single edge:

$$\begin{aligned} \sum_{p \in P_i} \overrightarrow{prox}(p) &= \sum_{p \in P_i} \left(\prod_{e \text{ edge} \in p} e.n_w \right) \\ &= \sum_{\substack{p' \text{ prefix of } p \in P_i \\ |p'|=i-1}} \left(\prod_{e \text{ edge} \in p'} e.n_w \right) \\ &\quad \times \sum_{\substack{e' \text{ edge starting at the end of } p' \\ \text{the concatenation of } p' \text{ and } e' \text{ is in } P_i}} e'.n_w \end{aligned}$$

By the definition of normalized weights for any p' :

$$\sum_{e' \text{ edge starting at the end of } p'} e'.n_w \leq 1,$$
 therefore:

$$\begin{aligned} \sum_{p \in P_i} \overrightarrow{prox}(p) &\leq \\ &\sum_{\substack{p' \text{ prefix of } p \in P_i \\ |p'|=i-1}} \left(\prod_{e \text{ edge} \in p'} e.n_w \times 1 \right) \end{aligned}$$

As $P_{i-1} = \{p' \text{ prefix of } p \in P_i, |p'| = i - 1\}$ is a set of paths of length $i - 1$ starting from u , the inductive hypothesis implies that:

$$\sum_{p \in P_i} \overrightarrow{prox}(p) \leq 1$$

From the lemma follows that $prox - prox^{\leq n} \leq \sum_{i > n} \frac{1}{\gamma^n}$. Because $\gamma > 1$, we have thus shown that $B_{prox}^{>x} = \sum_{i > n} \frac{1}{\gamma^n}$ tends to 0 as n grows to infinity.

3. Score soundness We show that the $\mathbf{S3}_k$ score increases with $prox$. Recall (Definition 3.4) that the score of a document c is:

$$score(c, (u, \phi)) = \prod_{k \in \phi} \left(\sum_{(type, f, src) \in con(c, k)} \eta^{pos(c, f)} \times prox(u, src) \right)$$

Because all quantities above are positive, and due to the continuity and monotonicity of addition and multiplication, $score(c, (u, \phi))$ is monotonous and increasing with $prox$. That is: for any src , if we replace any value of $prox(u, src)$ with an greater one, $score(d, (u, \phi))$ increases.

Algorithm 1: $S3_k$ – Top- k algorithm.

Input : a query $q = (u, \phi)$
Output: the best k answers to q over an $S3$ instance I , $T_{k,s}(q)$

```

1  $candidates \leftarrow []$  // initially empty list
2  $discarded \leftarrow \emptyset$ 
3  $borderPath \leftarrow []$ 
4  $allProx \leftarrow \delta_u$  //  $\delta_u[v] = \begin{cases} 1 & \text{if } v = u \\ 0 & \text{otherwise} \end{cases}$ 
5  $threshold \leftarrow \infty$  // Best possible score of a document not yet explored, updated in
   ComputeCandidatesBounds
6  $n \leftarrow 0$ 
7 while not StopCondition( $candidates$ ) do
8    $n \leftarrow n + 1$ 
9   ExploreStep()
10  ComputeCandidatesBounds()
11  CleanCandidatesList()
12 return  $candidates[0, k - 1]$ 

```

4. Score convergence This criterion states that *if* we are given a bound B on the proximity between the seeker and any source, *then* we can exhibit an upper bound on a document's score, namely B_{score} , which tends to 0 when B does.

Because $\eta < 1$, a document's $S3_k$ score for a query can be bound as follows:

$$score(c, (u, \phi)) \leq \prod_{k \in \phi} \left(\sum_{(type, f, src) \in con(c, k)} prox(u, src) \right)$$

Therefore, if we have B such that $prox(u, src) \leq B$ for any keyword query $k \in \phi$, and for any src such that $(type, f, src) \in con(c, k)$, it follows directly that:

$$score(c, (u, \phi)) \leq \prod_{k \in \phi} B \times |con(c, k)|$$

Based on this, we set $B_{score}(q, B) = \prod_{k \in \phi} B \times \max_{c \in I} |con(c, k)|$. This $B_{score}(q, B)$ clearly tends to 0 as B does.

The above concludes the proof that the $S3_k$ score (Definition 3.5) satisfies the feasibility constraints and thus, is a concrete score.

□

4 Query answering algorithm

In this section, we describe our *Top-k* algorithm called $S3_k$, which computes the answer to a query over an $S3$ instance using our $S3_k$ score, and formally state its correctness.

$q = (u, \phi)$	Query: seeker u and keyword set ϕ
k	Result size
n	Number of iterations of the main loop of the algorithm
$candidates$	Set of documents and/or fragments which are candidate query answers at a given moment
$discarded$	Set of documents and/or fragments which have been ruled out of the query answer
$borderPath[v]$	Paths from u to v explored at the last iteration ($u \rightsquigarrow_n v$)
$allProx[v]$	Bounded social proximity ($prox^{\leq n}$) between the seeker u and a node v , taking into account all the paths from u to v known so far
$connect[c]$	Connections between the seeker and the candidate c : $connect[c] = \{(k, type, pos(d, f), src) k \in \phi, (type, f, src) \in con(c, k)\}$
$threshold$	Upper bound on the score of the documents not visited yet

Table 3: Main variables used in our algorithms.

4.1 Algorithm

The main idea, outlined in Algorithm 1, is the following. The instance is explored starting from the seeker and going to other vertices (users, documents, or resources) at increasing distance. At the n -th iteration, the I vertices explored are those connected to the seeker by at least a path of length at most n . We term *exploration border* the set of graph nodes reachable by the seeker through a path of length exactly n . Clearly, the border changes as n grows.

During the exploration, documents are collected in a set of *candidate* answers. For each candidate c , we maintain a score interval: its *currently known lowest possible score*, denoted $c.lower$, and its *highest possible score*, denoted $c.upper$. These scores are updated as new paths between the seeker and the candidates are found. Candidates are kept sorted *by their highest possible score*; the k first are the answer to the query when the algorithm stops, i.e., when no candidate document outside the current first k can have an *upper bound* above the *minimum lower bound* within the top k ranks.

Further, the search algorithm relies on three tables:

- $borderPath$ is a table storing, for a node v in I , the set of paths of length n between u (the seeker) and v , where n is the current distance from u that the algorithm has traversed.
- $allProx$ is a table storing, for a node v in I , the proximity between u and v taking into account all the paths known so far from u to v . Initially, its value is 0 for any $v \neq u$.
- $connect$ is a table storing for a candidate node c the set of the connections (Section 3.2) discovered so far between the seeker and c .

These tables are updated during the search. While they are defined on all the I nodes, we only compute them gradually, for the nodes on the exploration border.

Termination condition Of course, search should not explore the whole graph, but instead stop as early as possible, while returning the correct result. To check whether the algorithm can terminate, To this aim, we maintain during the search an upper bound on the score of score

Algorithm 2: Algorithm StopCondition

Input : *candidates* set
Output: **true** if *candidates*[0, *k* - 1] is $T_{k,s}(q)$, **false** otherwise

- 1 **if** $\exists d, d' \in \text{candidates}[0, \dots, k - 1], d \in \text{neigh}(d')$ **then**
- 2 | **return false**
- 3 $\text{min_topk_lower} \leftarrow \infty$
- 4 **foreach** $c \in \text{candidates}[0, \dots, k - 1]$ **do**
- 5 | $\text{min_topk_lower} \leftarrow \min(\text{min_topk_lower}, c.\text{lower})$
- 6 $\text{max_non_topk_upper} \leftarrow \text{candidates}[k].\text{upper}$
- 7 **return** $\max(\text{max_non_topk_upper}, \text{threshold}) \leq \text{min_topk_lower}$ // Boolean result

of all documents unexplored so far, named *threshold*. Observe that we do not need to return the exact score of our results, and indeed we may never narrow down the (lower bound, upper bound) intervals to single numbers; we just need to make sure that no document unexplored so far is in among the top *k*. Algorithm 2 outlines the procedure to decide whether the search is complete: when (i) the candidate set does not contain documents such that one is a fragment of another, and (ii) no document can have a better score than the current top *k*.

Graph exploration Algorithm 3 describes one search step (iteration), which visits nodes at a social distance *n* from the seeker. For the ones that are documents or tags, the **GetDocuments** algorithm (see hereafter) looks for related documents that can also be candidate answers (these are added to *candidates*); *discarded* keeps track of related documents with scores too low for them to be candidates. The *allProx* table is also updated using the U_{prox} function, whose existence follows from the first score feasibility property (Section 3.3), to reflect the knowledge acquired from the new exploration border (*borderPath*). Observe that Algorithm 3 computes $prox^{\leq n}(u, src)$ iteratively using the first feasibility property; at iteration *n*, $allProx[src] = prox^{\leq n}(u, src)$.

Computing candidate bounds The **ComputeCandidateBounds** algorithm (Algorithm 4) maintains during the search the lower and upper bounds of the documents in candidates as well as the best possible score of unexplored documents. A candidate's *lower* bound is computed as its score where its social proximity to the user⁷ is approximated by its bounded version, based only on the paths explored so far:

$$\oplus_{gen}(\{(kw, type, pos(d, f), allProx[src]) \mid kw \in \phi, (type, f, src) \in con(d, kw)\})$$

This is a lower bound because, *during exploration, a candidate can only get closer to the seeker* (as more paths are discovered).

A candidate's *upper* bound is computed as its score, where the social proximity to the user is replaced by the sum between the bounded proximity and the function $B_{prox}^{>n}(u, src)$, whose existence follows from the long path attenuation property (Section 3.3). The latter is guaranteed to offset the difference between the bounded and actual social proximity:

$$\oplus_{gen}(\{(kw, type, pos(d, f), allProx[src] + B_{prox}^{>n}(u, src)) \mid kw \in \phi, (type, f, src) \in con(d, kw)\})$$

⁷The actual (exact) social proximity requires a complete traversal of the graph; our algorithms work with approximations thereof.

Algorithm 3: Algorithm ExploreStep

Update: *borderPath* and *allProx*

```

1 if  $n = 1$  then
2    $\lfloor$  borderPath  $\leftarrow$  out( $\{u\}$ )
3 else
4   foreach  $v \in I$  do
5      $\lfloor$  newBorderPath[ $v$ ]  $\leftarrow$   $\emptyset$ 
6   foreach  $p \in$  borderPath do
7     foreach network edge  $e$  in out(neigh( $p.end$ )) do
8        $\lfloor$   $m \leftarrow e.target$ 
9         if  $m$  is a document or a tag then
10           $\lfloor$  GetDocuments( $m$ )
11           $\lfloor$  newBorderPath[ $m$ ].add( $p|e$ )
12    $\lfloor$  borderPath  $\leftarrow$  newBorderPath
13 foreach  $v \in I$  do
14    $\lfloor$  newAllProx[ $v$ ]  $\leftarrow$  allProx[ $v$ ] +  $U_{prox}(allProx[v],$ 
15    $\lfloor$   $\{\overrightarrow{prox}(p), p \in borderPath[v]\}, n)$ 
16 allProx  $\leftarrow$  newAllProx

```

Algorithm 4: Algorithm ComputeCandidateBounds

Input : loop counter n

Update: *upper* and *lower* bounds of the *candidates*, *threshold*

```

1 foreach  $c \in candidates$  do
2    $\lfloor$   $c.lower \leftarrow \oplus_{gen}(\{(kw, t, p, allProx[src])|$ 
3    $\lfloor$   $(kw, t, p, src) \in connect[c]\})$ 
4    $\lfloor$   $c.upper \leftarrow \oplus_{gen}(\{(kw, t, p, allProx[src] +$ 
5    $\lfloor$   $B_{prox}^{>n}(u, src))|(kw, t, p, src) \in connect[c]\})$ 
6  $threshold \leftarrow B_{score}(q, B_{prox}^{>n})$ 

```

The above bounds rely on $con(d, k)$, the set of all connections between a candidate d and a query keyword k (Section 3.2); clearly, the set is not completely known when the search starts. Rather, connections accumulate gradually in the *connect* table (Algorithm *GetDocuments*), whose tuples are used as approximate (partial) $con(d, k)$ information in *ComputeCandidateBounds*.

Finally, *ComputeCandidateBounds* updates the relevance threshold using the known bounds on *score* and *prox*. The new bound estimates the best possible score of the unexplored documents⁸; it tends to 0 as n grows, due to the score convergence feasibility property (Section 3.3).

Cleaning the candidate set Algorithm 5 removes from *candidates* some documents that cannot be in the answer, i.e., those for which k candidates with better scores are sure to exist, as well as those having a candidate neighbor with a better score. In the former case, this translates into removing every document d whose score is less than that of k other candidates which are neither (i) neighbors pairwise nor (ii) neighbors with other candidates with scores greater than that of d .

⁸See the Threshold Correctness Lemma (4.2).

Algorithm 5: Algorithm CleanCandidatesList

Update: *candidates* and *discarded*

```

1 lower_bounds_list  $\leftarrow$  [] // Initially empty list
2 max_rank = k
3 foreach c  $\in$  candidates do
4   if  $|lower\_bounds\_list| < max\_rank$  or c.lower  $\geq lower\_bounds\_list[max\_rank]$ 
   then
5     lower_bound_list.insertSort(c.lower)
6   else if c.upper  $< lower\_bounds\_list[max\_rank]$  then
7     candidates.remove(c)
8     discarded.add(c)
9   if No vertical neighbor of c may have a score equal or higher than c then
10    foreach f  $\in$  neigh(c) do
11      candidates.remove(f) // If f was a candidate
12      discarded.add(f)
13  else if  $|lower\_bounds\_list| < max\_rank$  then
14    max_rank += 1

```

Algorithm 6: Algorithm GetDocuments

Input : Document or tag *x*
Update: *candidates* and *discarded*

```

1 foreach d such that there exists a chain of triples from x to d in I using only S3:partOf,
  S3:commentsOn, S3:commentsOn, S3:hasSubject and S3:hasSubject labels do
2   if d is a document and d  $\notin$   $(candidates \cup discarded)$  then
3     connect[d]  $\leftarrow$   $\emptyset$ 
4     kw_index  $\leftarrow$  0
5     while kw_index  $<$  length( $\phi$ ) do
6       kw  $\leftarrow$   $\phi[kw\_index]$ 
7       if con(kw, d) is empty then //A keyword from the query is missing
8         kw_index  $\leftarrow$  length( $\phi$ )
9         discarded.add(d)
10      else
11        kw_index += 1
12        foreach  $(type, frag, source) \in con(kw, d)$  do
13          connect[d].add(
14            (kw, type, pos(d, frag), source)
15          )
16    candidates.add(d)

```

Getting candidate documents Algorithm 6 checks whether every unexplored document reachable from a given document or tag through a chain of triples labelled S3:partOf, S3:commentsOn,

Algorithm 7: Algorithm AnytimeTermination**Output:** A list of candidates

```

1 return_list ← []
2 while length(return_list) < k and candidates ≠ ∅ do
3   d ← candidates[0] //candidates is kept sorted by upper bound
4   return_list.add(d)
5   candidates.remove(neigh(d))
6 return return_list

```

S3:commentsOn, S3:hasSubject, or S3:hasSubject (i.e., any edge that may connect two documents or tags), is a candidate answer. If yes, it is added to *candidates* and the necessary information to estimate its score, derived from *con*, is recorded in *connect*.

Anytime termination Algorithm 7 can be used to make our query answering algorithm (Algorithm 1) *anytime*. It outputs the *k* best candidate documents known so far, based on their current upper bound score.

4.2 Correctness of the algorithm

The theorems below state the correctness of our algorithm for *any score function having the feasibility properties* identified in Section 3.3.

Lemma 4.1 (Bounds correctness and convergence). *At the end of each iteration of Algorithm 1, for every candidate c in $candidates$: $c.upper \geq score(c, q) \geq c.lower$. Furthermore, $\lim_{n \rightarrow +\infty} (c.upper) = score(c, q)$ and $\lim_{n \rightarrow +\infty} (c.lower) = score(c, q)$ where n is the number of iterations.*

Proof. (Lemma 4.1) *c.upper* and *c.lower* are only modified in Algorithm 4, so at the n -th iteration after the execution of Algorithm 4, we have:

$$\begin{aligned}
 c.upper &= \oplus_{gen} (\{(k, t, p, allProx(src)) \mid (k, t, p, src) \in connect[c]\}) \\
 c.lower &= \oplus_{gen} (\{(k, t, p, allProx(src) + B_{prox}^{>n}(u, src)) \mid \\
 &\quad (k, t, p, src) \in connect[c]\})
 \end{aligned}$$

Since documents are added to *candidates* only through Algorithm 6 we know that $(k, t, p, src) \in connect[c] \Leftrightarrow k \in \phi, (t, f, src) \in con(c, k), p = pos(f, c)$.

Because we also know that $allProx[src] = prox^{\leq n}(u, src)$ we have:

$$\begin{aligned}
 c.upper &= score_{[prox \leq n]}(c, q) \\
 c.lower &= score_{[prox \leq n + B_{prox}^{>n}]}(c, q)
 \end{aligned}$$

where $score_{[f]}$ denotes our score function (Section 3.3) using f as the social proximity function. By definition, $prox^{\leq n} \leq prox$ and by the second feasibility property we know that $prox \leq prox^{\leq n} + B_{prox}^{>n}$ and $\lim_{n \rightarrow +\infty} (prox^{\leq n} - prox) = 0 =$

$\lim_{n \rightarrow +\infty} (prox^{\leq n} + B_{prox}^{>n} - prox)$. The third feasibility property ensures that $g \mapsto score_{[g]}$ is monotonically increasing and continuous and therefore that $c.lower \leq score(c, q) \leq c.upper$ and that $\lim_{n \rightarrow +\infty} (c.lower) = \lim_{n \rightarrow +\infty} (c.upper) = score(c, q)$. \square

Lemma 4.2 (Threshold correctness). *At any time, for any document d not in *candidates* nor in *discarded*: $threshold \geq score(d, q)$.*

Proof. (Lemma 4.2) Let d be a document neither in *candidates* nor in *discarded*, at an iteration $n > 0$ ⁹. Because d is not in *candidates* nor in *discarded*, there is no path of length smaller than $n - 1$ from u to a source for $con(d, k), k \in \phi$ (otherwise, this path would have been found by Algorithm 6 and d would have been added to *candidates* or *discarded*). Let src be such a source. Because $prox^{\leq n}(u, src) = 0$, and due to the long path attenuation property on the score (Section 3.3), we know that $prox(u, src) \leq B_{prox}^{>n}$. Therefore, thanks to the fourth score feasibility property, $score(d, q) \leq B_{score}(q, B_{prox}^{>n})$. This is exactly the *threshold* value set by Algorithm 4 at line 6. \square

We partition *candidates* into *groups* of documents. Two documents in *candidates* are in the same group if and only if *candidates* comprises a list of nodes, containing both documents, such that each node is a vertical neighbor of the next one in the list. Observe that in particular candidates that are not part of the same document are never in the same group. We call *score of a group* the highest score of its elements, and *lower bound of a group* the highest lower bound of its elements.

Lemma 4.3. *The scores of the k groups with the best scores can only increase, and only documents with a lower score than the score of the k -th group or having vertical neighbors in the k groups with the best scores can be in *discarded*.*

Proof. (Lemma 4.3)

Observe that the lemma's claim is true when *discarded* is empty. We show that it is also true when documents are added to *discarded* or when the k best groups change.

Unless a document's score is null and gets discarded by GetDocuments (line 10), a document can only go to *discarded* by leaving *candidates* due to CleanCandidatesList (algorithm 5).

The document with the best score in each group only has neighbors with lower scores than itself and therefore cannot be removed by lines 11-12. The only other way of discarding documents from candidates (lines 7-8) removes documents having an upper bound lower than the lower bound of the k -th group with the best lower bound. Because lower bounds are always smaller than the scores the best elements of the k groups with the best scores, they cannot be removed this way either. Therefore, at each call of CleanCandidatesList, each document with the best score in the k groups with the best scores remains in *candidates*, guaranteeing that the scores of the k best group can only increase and that adding new documents to *discarded* respects the lemma's claim.

As for the change of the k best groups, if a new group introduced in *candidates* removes the current k -th best group from its k -th position then the score of the new group is better than that

⁹The theorem also holds if $n = 0$, as $threshold = \infty$.

of the removed one. Hence, documents which were in *discarded* because they were a neighbor of some element of the removed group, are now in *discarded* because they have a lower score than the new k -th best group (documents that had a lower score than the previous k best groups continue to do so with the new k ones). \square

Theorem 4.1 (Stop correctness). *When a stop condition is met, the first k elements in candidates are a query answer.*

Proof. (Theorem 4.1)

Let us show by induction on $i \in [0, k]$ that if the stop condition is met then the i documents with the best scores among the first k elements of *candidates* form a $T_{i, \text{score}}(q)$: a top- i answer to the query.

For $i = 0$, the 0 highest scored elements of the k first candidates is the empty set and we purposefully expand the notion of query answer such that: $T_{0, \text{score}}(q) = \emptyset$.

For $i > 0$, given that the $i - 1$ documents with the best scores among the first k elements of *candidates* form a $T_{i-1, \text{score}}(q)$, we have to show that adding a i -th document with the i -th best score among the first k elements of *candidates*, to this $T_{i-1, \text{score}}(q)$ forms a $T_{i, \text{score}}(q)$. In order to do so, let d be a such a document. By the definition of query answer, we only have to show that d is in $D \setminus \text{neigh}(T_{k-1, \text{score}}(q))$ and has the maximum score ($\text{score}(d, q)$) on that set.

d is not in $\text{neigh}(T_{k-1, \text{score}}(q))$ since it follows from the stop condition that $\nexists d, d' \in \text{candidates}[0, \dots, k-1]$, $d \in \text{neigh}(d')$ (lines 1-2) and from the induction hypothesis that $T_{k-1, \text{score}}(q) \subset \text{candidates}[0, \dots, k-1]$. Therefore, d is in $D \setminus \text{neigh}(T_{k-1, \text{score}}(q))$.

Let us show now that d maximizes the score over $D \setminus \text{neigh}(T_{k-1, \text{score}}(q))$, by showing it does so over the intersection of $D \setminus \text{neigh}(T_{k-1, \text{score}}(q))$ with the following partition of D :

- the k first elements of *candidates*: by construction d has the best score among the k first elements of *candidates* which are not in $T_{k-1, \text{score}}(q)$

- the other elements of *candidates*: if d' is a document in *candidates*, which is not ranked among the first k , then:

$d'.\text{upper} \leq \text{candidates}[k].\text{upper}$ because *candidates* is kept sorted by upper bound

$\text{candidates}[k].\text{upper} \leq \min_{c \in \text{candidates}[0, k-1]}(c.\text{lower}) \leq d.\text{lower}$ because of the stop condition (line 7)

$\text{score}(d', q) \leq d'.\text{upper}$ and $d.\text{lower} \leq \text{score}(d, q)$ by bounds correctness (Lemma 4.1).

Hence $\text{score}(d', q) \leq \text{score}(d, q)$.

- *discarded*: From Lemma 4.3, it follows that if d' is a document in *discarded* then either it has a lower score than d or it has a neighbor d'' in *candidates* such that $\text{score}(d'', q) \geq \text{score}(d', q)$. In the latter case, if this neighbor d'' is not in $T_{k-1, \text{score}}(q)$ then $\text{score}(d'', q) \leq \text{score}(d, q)$ and $\text{score}(d', q) \leq \text{score}(d, q)$. In any case there exists no d' in *discarded* $\setminus \text{neigh}(T_{k-1, \text{score}}(q))$ with a better score than d .

- $D \setminus (\text{candidates} \cup \text{discarded})$: from the Threshold correctness (Lemma 4.2), we know that documents not in *candidates* nor in *discarded* have a score lower than *threshold*, thus the line 7 in Algorithm 2 and bounds correctness (Lemma 4.1) guarantee that $\text{score}(d, q) \geq \text{threshold}$.

We have shown that the k documents with the best score among the first k elements of *candidates*, i.e., the first k elements of *candidates*, are a query answer. \square

We say the tie of two equal-score documents d, d' is *breakable* if examining a set of paths of *bounded length* suffices to decide their scores are equal. (In terms of our score feasibility properties, this amounts to $B_{prox}^{>n} = 0$ for some n). Our generic score function (Definition 3.5) does not guarantee all ties are breakable. However, any finite-precision number representation eventually brings the lower and upper bounds on d and d' 's scores too close to be distinguished, de facto breaking ties.

Theorem 4.2 (Correctness with breakable ties).

If there exists a query answer of size k and all ties are breakable then Algorithm 1 returns a query answer of size k .

Proof. (Theorem 4.2) If the algorithm terminates then it follows from Theorem 4.1 that the algorithm returns an answer. Now, to exhibit a contradiction, let us suppose that it does not terminate.

If there exists a query answer of size k then *candidates* will eventually contain at least k groups; otherwise some documents from this query answer would have gone to *discarded*, which by Lemma 4.3 is impossible if *candidates* doesn't contain at least k groups with better scores than these of the k documents from this query answer.

Because it is impossible to add new documents to *candidates* when all documents have been explored, after some time, the size of *candidates* will not change and no document will leave it anymore.

By Lemma 4.3, at least one document in each of the k best groups remains in *candidates*. Consider an element d with the best score in one the k best groups at a time where *candidates* has reached its final size. Because of the bound correctness and convergence lemma (Lemma 4.1), and the fact that all ties are breakable, eventually $\forall d' \in \text{neigh}(d)$ $d.lower \geq d'.upper$ holds. Therefore, the CleanCandidatesList algorithm, at lines 11-12, remove documents from *candidates* if any such d' is in candidates. Hence, d has no neighbors in *candidates* and is the only document in its group.

The k best groups are therefore k documents without neighbors in *candidates*, with the k best scores in *candidates*, and by bound convergence the best lower bounds eventually. This automatically triggers the stop condition and forces the algorithm to terminate, which contradicts our hypothesis.

We have shown that if there exists a query answer and all ties are breakable then the algorithm terminates and thus returns an answer (Theorem 4.1). \square

Theorem 4.3 (Anytime correctness). *Using anytime termination, Algorithm 1 eventually returns a query answer.*

Proof. (Theorem 4.3) Because (i) there is a finite number of documents that can be candidates and (ii) the simple convergence of bounds (Lemma 4.1) is actually a uniform convergence: for any query q and positive value ϵ , after some number of iterations, the lower and upper bounds of any candidate can differ from its score of at most ϵ .

If we choose for ϵ a value smaller than half the minimum positive difference between two document scores, then eventually any two candidates with different scores will have their upper bounds ordered in the same way as their score. In particular, after some time, $\text{argmax}(c.upper) = \text{argmax}(score(q, c))$ and therefore the anytime termination algorithm (Algorithm 7) produces a query answer. \square

It is worth noting that in our experiments (Section 5), the threshold-based termination condition was always met, thus we never needed to wait for convergence of the lower and upper bound scores, in order to find the top- k answers.

5 Implementation & experiments

We describe experiments creating and querying S3 instances. We present the datasets in Section 5.1, while Section 5.3 outlines our implementation and some optimizations we brought to the search algorithm. We report query processing times in Section 5.4, study the quality of our returned results in Section 5.5, then we conclude.

5.1 Datasets, queries, and systems

Datasets We built three datasets, I_1 , I_2 , and I_3 , based respectively on content from Twitter, Vodkaster and Yelp.

The instance I_1 was constructed starting from tweets obtained through the public streaming Twitter API, and based on the Tweepy library [27]. Over a one-day interval (from May 2nd 2014 16h44 GMT to May 3rd 2014 12h44 GMT), we gathered roughly one million tweets. From every tweet that is not a retweet, we create a document having three nodes: (i) *text* from the *text* field of the tweet, we extract named entities and words (with the help of the twitter NLP tools library [18]) and match them against a general-purpose ontology we created from DBpedia, as we explain below (ii) *date*, extracted from the *created_at* tweet field and (iii) a *geo* node: if the tweet includes a human readable location, i.e., recognizable keywords in the *place* field of the tweet (be it from a *name* or a *full_name* property) we insert it in this node.

The RDF graph of our instance is built from four DBpedia datasets, namely: Mapping-based Types, Mapping-based Properties, Persondata and DBpedia Lexicalizations Dataset. These were chosen as they were the most likely to contain concepts (names, entities etc.) occurring in tweets.

Within the *text* fields, we replace each word w for which a triple of the form u `http://xmlns.com/foaf/0.1/name w` holds in the DBpedia knowledge base, by the respective URI u .

When a tweet t' authored by an user u is a *retweet* of another tweet t , and further if t' adds a hashtag h , then we add to I_{d_1} the following triples:

$$\begin{array}{ll} \text{a type S3:relatedTo} & \text{a S3:hasSubject } t \\ \text{a S3:hasKeyword } h & \text{a S3:hasAuthor } u \end{array}$$

Thus, a retweet that adds multiple keywords leads to the creating several tags. If a tweet t'' is a *reply* to another tweet t (as identified by the *in_reply_to_status_id* field), we consider that t'' is a comment on t . If t is present in our dataset¹⁰, we add the corresponding S3:commentsOn triple in I_1 .

For what concerns the social network, we set Ω_{I_1} as the set of IDs of the users having posted some tweets, and we create links between users as follows. We assign to every pair of users (a, b) a similarity value $u_{\sim}(a, b)$ between 0 and 1. u_{\sim} is a weighted sum of two Jaccard similarity

¹⁰Sometimes our corpus contains a re-tweet of an original tweet that we did not capture. This is unavoidable unless one has access to the full Twitter history.

I ₁ (Twitter)	
Users	492,244
S3:social edges	17 544 347
Documents	467,710
Fragments (non-root)	1,273,800
Tags	609,476
Keywords	28,126,940
Tweets	999,370
Retweets	85%
Reply to users' status	6.9%
String-keyword associations extracted from DBpedia	3,301,425
S3:social edges per user having any (average)	317
Nodes (without keywords)	2 972 560
Edges (without keywords)	24 554 029

I ₂ (Vodkaster)		I ₃ (Yelp)	
Users	5,328	Users	366,715
S3:social edges (vdk:follow)	94,155	S3:social edges (yelp:friend)	3,868,771
Documents (movie comments)	330,520	Documents (reviews)	2,064,371
Fragments (non-root)	529,432	Keywords	59,614,201
Keywords	3,838,662	Businesses	61,184
Movies	20,022		

Figure 4: Statistics on our instances.

coefficients:

$$\begin{aligned}
 u_{\sim}(a, b) = & \\
 & t \times \frac{|\text{keywords in comments from both } a \text{ and } b|}{|\text{keywords in comments from } a| + |\text{keywords in comments from } b|} \\
 & + (1 - t) \times \frac{|\text{keywords posted by both } a \text{ and } b|}{|\text{keywords posted by } a| + |\text{keywords posted by } b|}
 \end{aligned}$$

Whenever this similarity is above a certain threshold we create an edge with weight u_{\sim} between the two users. Based on experiments with the data, we set $t = 0.5$ and the $u_{\sim}(a, b)$ threshold for creating a link between a and b to 0.1.

The instance I₂ uses data from Vodkaster, a French social network dedicated to movies. The original social network data comprised *follower* relations between the users and a list of comments on the movies, in French, along with their author. Whenever user u follows user v , we included $u \text{ vdk:follow } v$ 1 in I₂, where *vdk:follow* is a custom subproperty of S3:social expressing the act of following someone in Vodkaster ($\text{vdk:follow} \prec_{\text{sp}} \text{S3:social}$). The first comment on a movie was translated in I₂ as an original document; each additional comment on this film was then considered a comment on the first. The textual content of each comment was stemmed using Snowball algorithms [26] and each (stemmed) sentence was made a fragment of the comment.

The instance I₃ is based on Yelp [29], a crowd-sourced reviews website about local businesses. This dataset contains a list of textual reviews of businesses, and the friend list of each user.

As for I_2 , we consider that the first review of a business is commented on by the subsequent reviews of the same business and we introduce a dedicated subproperty of S3:social: yelp:friend to express the friendship on Yelp: if user u is friend with user v then we include u yelp:friend v 1 in I_3 . As for I_1 , we use DBpedia to gain additional semantic on the extracted keywords.

Table 4 shows the main features of the three quite different data instances. I_1 is by far the largest. I_2 was not matched with a knowledge base since its content is in French; I_2 and I_3 have no tags.

Queries For each instance we created workloads of 100 queries, based on three independent parameters:

- f , the keyword frequency: either *rare*, denoted ‘-’ (among the 25% least frequent in the document set), or *common*, denoted ‘+’ (among the 25% most frequent)
- l , the number of keywords in the query: 1 or 5
- k , the expected number of results: 5 or 10

This lead to a total of 8 workloads, identified by $qset_{f,l,k}$, for each dataset. To further analyze the impact of varying k , we added 10 more workloads for I_1 , where $f \in \{+, -\}$, $l = 1$, and $k \in [1, 5, 10, 50]$ (used in Figure 8). *We stress here that injecting semantics in our workload queries, by means of keyword extensions (Definition 2.1), increased their size on average by 50%.*

Systems Our algorithms were fully implemented in Python 2.7; the code has about 6K lines. We stored some data tables in PostgreSQL 9.3, while others were built in memory, as we explain shortly. All our experiments were performed on a 4 cores Intel Xeon E3-1230 V2 @3.30GHz with 16Go of RAM, running Debian 8.1.

No existing system directly compares with ours, as we are the first to consider fine-granularity content search with semantics in a social network. To get at least a rough performance comparison, we used the top- k social search system described in [15], based on the UIT (user, item tag) model, and referred to as **TopkS** from now on. We used the Java-based code provided by the authors of [15]. The data model of TopkS is rather basic, since its documents (*items*) have no internal structure nor semantics; tags are all independent, i.e., there is no semantic connection between them. Further, (*user, user, weight*) tuples reflect weighted links between users. TopkS computes a social score of each item, and separately a content-based score; the overall item score is then obtained as $\alpha \times$ social score $+(1 - \alpha) \times$ content score, where α is a parameter of TopkS.

We adapted our instances into TopkS’s simpler data model. From I_1 , we created I'_1 as follows: (i) the relations between users were kept with their weight, (ii) every tweet was merged with all its retweets and replies into a single item, and (iii) every keyword k in the content of a tweet that is represented by item i posted by user u led to introducing the (user, item, tag) triple (u, i, k) . To obtain I'_2 and I'_3 : (i) the relations between the users are kept with their weight; (ii) each film or business becomes an item (iii) each word extracted from a user’s review on a item becomes a tuple in the (user, item, tag) relation.

5.2 Data Layout

The instance is stored on disk using PostgreSQL. We describe below our data layout; the underlined attributes denote a (clustering) primary key. For the keyword extension part we store the following relations:

- `text_association(name, id)`, with an additional on name and another one in id;

- `DBpediaURI_association(URI, id)`;
- `keyword_extension(id, extended_id)` with an additional index on `id`;

In the following, by the *components* of the `S3` graph we mean the maximal subgraphs of documents and tags connected by triples labelled by `S3:partOf`, `S3:commentsOn`, `S3:commentsOn`, `S3:hasSubject`, or `S3:hasSubject`.

The documents and the connections they participate to are compiled into the following set of tables:

- `network_edges(source, target, weight, type of relation)`, with an additional index on `source`;
- `keyword_container(fragment or tag, component, keyword)` stores the keywords contained in the documents and tags, grouped by component;
- `structure(root, JSON structure)` stores for each document tree the URI of the root of the tree and the position of the other URIs as a JSON expression;
- `component(document, root of the document if exists, component)`, with an additional index on the component attribute; this table associates to each document the root of the document tree and the component it is in;
- `raw_docs(URI, raw)` associates documents with their URIs; this table is only used to return actual documents to the querier (it is not used in the query processing phase);
- `original(URI, optional external identifier, type)`, where *type* is one of user, document root, fragment, or tag, with two additional indexes, on external identifier and on type. This table helps finding where the URIs in the database come from and what they are, in term of representation outside the algorithm. It is not used in the query processing phase, but serves to connect it with the application using it.

5.3 Implementation and optimizations

We briefly discuss our implementation, focusing on optimizations w.r.t. the conceptual description in Section 4.

The first optimization concerns the computation of *prox*, required for the score (Definition 3.5). While the score involves connections between documents and keywords found on any path, in practice `S3k` explores paths (and nodes) increasingly far from the seeker, and stores such paths in *borderPath*, which may grow very large and hurt performance. To avoid storing *borderPath*, we compute for each explored vertex *v* the weighted sum over all paths of length *n* from the seeker to this vertex:

$$borderProx(v, n) = \sum_{p \in u \rightsquigarrow v, |p|=n} \frac{prox(p)}{\gamma^n}$$

and compute *prox* directly based on this value.

Furthermore, Algorithm `GetDocuments` considers documents reachable from *x* through edges labeled `S3:partOf`, `S3:commentsOn`, `S3:commentsOn`, `S3:hasSubject` or `S3:hasSubject`. Reachability by such edges defines a *partition* of the documents into *connected components*. Further, by construction of *con* tuples (Section 3.2), connections carry over from one fragment to another, across such edges. Thus, a fragment matches the query keywords iff its component matches it, leading to an efficient pruning procedure: we compute and store the partitions, and test that each keyword (or extension thereof) is present in every component (instead of fragment). Partition

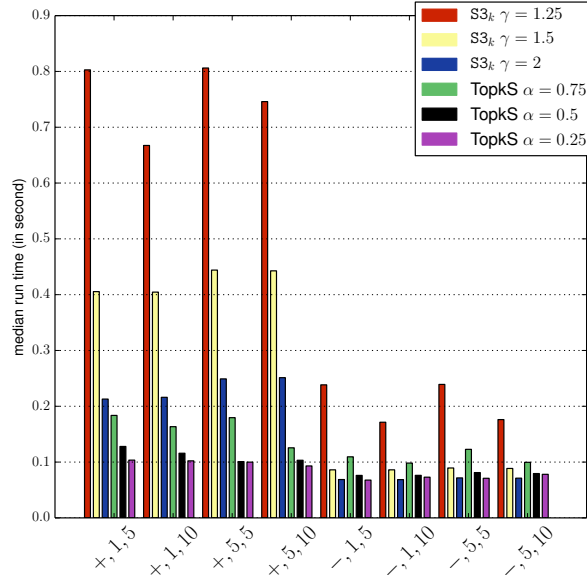


Figure 5: Query answering times on I_1 (Twitter).

maintenance is easy when documents and tags are added, and more expensive for deletions, but luckily these are rarer.

The query answering algorithm creates a boolean vector indexed by the components, initially all false. Further, it creates (in memory) the *allProx* table and two sparse matrices which are computed only once: *distance matrix* that encodes the graph of network edges in I (accounting for the vertical neighborhood), and a *component* matrix storing the component of each fragment or tag. This simplifies Algorithm 3, since computing *allProx* and finding new components to explore can be implemented using matrix and vector operations. For instance, the new distance vector *borderProx* w.r.t. the seeker at step $n + 1$ can be obtained by multiplying the distance matrix with the previous distance vector from step n .

Last but not least, the search algorithm can be *parallelized*, in two independent ways. First, within an iteration we discover new documents in parallel by splitting the search across components. Second, an iteration can start executing before the current one is finished: as long as *borderProx* is available in the current iteration, one can start computing the next *borderProx* using the (fixed) distance matrix. More precisely, **ExploreStep** can be seen as consisting of two main blocks: (i) computing the new *borderProx* using the distance matrix and the previous *borderProx* (lines 1- 12 except for line 10); (ii) computing *allProx* using the new *borderProx* and the previous *allProx* (lines 13-16) plus the call to **GetDocuments** (line 10). The latter algorithm can also be seen as consisting of two parts: (iii) identifying the newly discovered components (line 1), respectively (iv) testing the documents they contain (the remaining lines). In practice, we use 8 concurrent threads each running a task of one of the forms (i)-(iv) above, and synchronize them with a custom scheduler. This has divided the query answering time on average by a factor of 2.

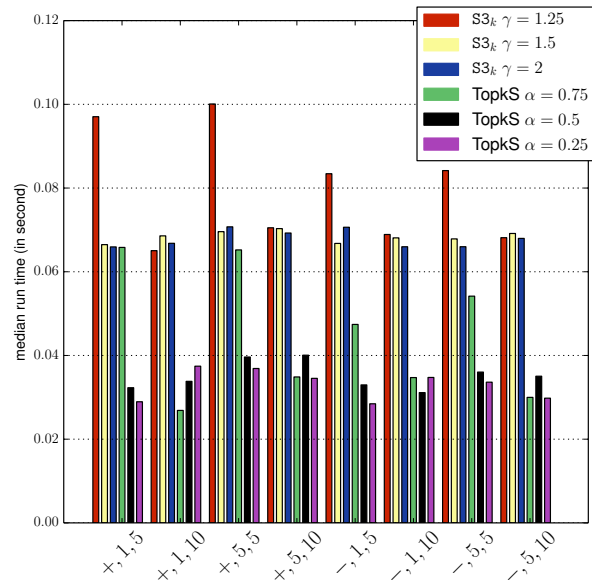


Figure 6: Query answering times on I_2 (Vodkaster).

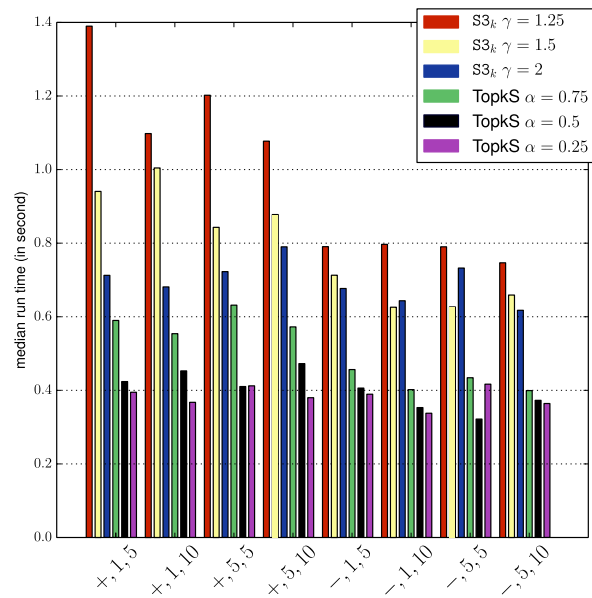


Figure 7: Query answering times on I_3 (Yelp).

5.4 Query answering times

Figures 5 – 7 show the running times of $S3_k$ on our three instances. We used different values of the γ social proximity damping factor (Section 3.4), and of α for TopkS. For each workload, we plot the average time (over its 100 queries). *All runs terminated by reaching the threshold-based stop condition* (Algorithm 2).

A first thing to notice is that while all running times are comparable, TopkS runs consistently faster. This is mostly due to the different proximity functions: our *prox*, computed from all possible paths, has a much broader scope than TopkS, which explores and uses only one (shortest) path. In turn, as we show later, we return a significantly *different* set of results, due to *prox*'s broader scope and to considering document structure and semantics.

Decreasing the γ in $S3_k$ reduces the running time. This is expected, as γ gives more weight to nodes far from the seeker, whose exploration is costly. Similarly, *increasing* α in TopkS forces to look further in the graph, and affects negatively its performance.

The influence of k is more subtle. When the number of candidates is low and the exploration of the graph is not too costly, higher k values allow to include most candidates among the k highest-scoring ones. This reduces the exploration needed to refine their bounds enough to clarify their relative ranking. In contrast, if the number of candidates is important and the exploration costly, a small k value significantly simplifies the work. This can be seen in Figure 8 where, with frequent keywords, increasing k does not affect the 3 fastest quartiles but significantly slows down the slowest quartile, since the algorithm has to look further in the graph.

The same figure also shows that rare-keyword workloads (whose labels start by $-$) are faster to evaluate than the frequent-keyword ones (workload labels starting with $+$). This is because finding rare keywords tends to require exploring longer paths. Social damping at the end of such paths is high, allowing to decide that possible matches found even farther from the seeker will not make it into the top- k . In contrast, matches for frequent keywords are found soon, while it is still possible that nearby exploration may significantly change their relative scores. In this case, more search and computations are needed before the top- k elements are identified.

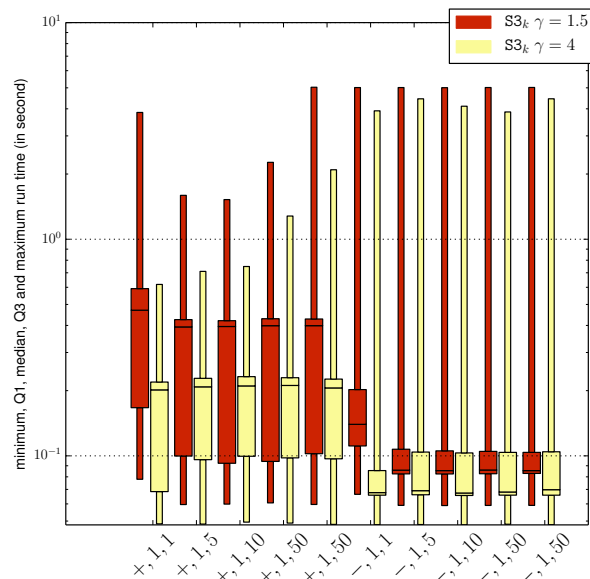
5.5 Qualitative comparison

We compare now the answers of our $S3_k$ algorithm and those of TopkS from a *qualitative* angle. $S3_k$ follows links between documents to access further content, while TopkS does not; we term *graph reachability* the fraction of candidates reached by our algorithm which are not reachable by the TopkS search. Further, while $S3_k$ takes into account semantics by means of semantic extension (Definition 2.1), TopkS only relies on the query keywords. We call *semantic reachability* the ratio between the number of candidates examined by an algorithm *without* expanding the query, and the number of candidates examined *with* query expansion. Observe that some $S3_k$ candidates may be ignored by TopkS due to the latter's lack of support for *both* semantics and connections between documents. Finally, we report two measures of distance between the results of the two algorithms. The first is the *intersection size* i.e., the fraction of $S3_k$ results that TopkS also returned. The second, L_1 , is based on Spearman's well-known *foot rule* distance between lists [6]. Modulo normalisation to ensure that $L_1(\tau, \tau) = 1$, $L_1(\tau_1, \tau_2) = 0$ if they share no elements and 0.5 if they have the same elements in reverse order, L_1 is, defined as:

$$L_1(\tau_1, \tau_2) = 2(k - |\tau_1 \cap \tau_2|)(k + 1) + \sum_{i \in \tau_1 \cap \tau_2} |\tau_1(i) - \tau_2(i)| - \sum_{\substack{\tau \in \{\tau_1, \tau_2\} \\ i \in \tau \setminus (\tau_1 \cap \tau_2)}} \tau(i)$$

where $\tau_j(i)$ is the rank of item i in the list τ_j .

The averages of these 4 measures over the 8 workloads on each instance appear in Figure 9. The ratios are low, and show that different candidates translate in different answers (the low L_1

Figure 8: Query answering times on I_1 when varying k .

Measure \ Instance	I_1	I_2	I_3
Graph reachability	12%	23%	41%
Semantic reachability	83%	100%	78%
L_1	8%	10%	4%
Intersection size	13.7%	18.4%	5.6%

Figure 9: Relations between the $S3_k$ and TopkS answers.

stands witness for this). Few $S3_k$ results can be attained by an algorithm such as TopkS, which ignores semantics and relies only on the shortest path between the seeker and a given candidate.

5.6 Experiment conclusion

Our experiments have demonstrated first the ability of the $S3$ data model to *capture very different social applications*, and to query them meaningfully, accounting for their structure and enriching them with semantics. Second, we have shown that $S3_k$ query answering can be quite efficient, even though considering all paths between the seeker and a candidate answer slows it down w.r.t. simpler algorithms, which rely on a shortest-path model. We have experimentally verified the expected impact of the social damping factor γ and of the result size k on running time. Third, and most importantly, we have shown that taking into account in the relevance model the social, structured, and semantic aspects of the instance bring a *qualitative gain*, enabling meaningful results that would not have been reachable otherwise.

6 Related work

Prior work on *keyword search in databases* spreads over different research directions:

Top-k search in a social environment uses UIT models [19, 30, 15] we outlined in Section 1. Top-k query results are found based on a score function accounting for the presence of each keyword in the tags of a candidate item, and a simple social distance based on the length of the social edge paths; query answering algorithms are inspired from the general top-k framework of [7]. As documents are considered atomic, and relations between them are ignored, requirements **R1**, **R2** and **R4** are not met. Further, the lack of semantics also prevents **R5**. Recent developments tend to focus on performance and scalability, or the integration of more attributes such as locality or temporality [5, 14], without meeting the abovementioned requirements.

Semi-structured document retrieval based on keywords, rely mostly on the *Least Common Ancestors* approach, by which a set of XML nodes containing the requested keywords are resolved into one result item, their common ancestor node [21, 4]. This field pioneered by [10], encompassed by our model, generalizes LCA constraints but lacks both social and semantics, and thus meets only **R2**. Other recent developments in this area, including more flexible and comprehensive reasoning patterns, have been presented in [2] but have the same limitations. IR-style search in relational databases [11, 12] considers data relations through key-foreign key pairs, but ignores text structure, semantics, and lacks a social dimension.

Semantic search on full-text documents, either via RDF [13, 23] or a semantic similarity measure [22], allows to query interconnected, semantic rich unstructured textual documents or entities, thus meeting **R1**, **R5** and **R6**. Efforts to consider XML structure in such semantics-rich models [17, 8] also enable **R2**.

All the aforementioned models can be seen as partial views over the **S3** model we devised, and they could easily be transcribed into it modulo some minor variations; for instance, Facebook’s GraphSearch [5] is a restricted form of SPARQL query one could ask over an **S3** instance. Slight adaptations may be needed for social contexts tolerating *similarity* between keywords that goes beyond the strict specialization relation (in RDF sense) we consider. We have hinted in Section 2 how this could be included.

7 Conclusion

We devised the **S3** data model for structured, semantic-rich content exchanged in social applications. We also provided the **S3_k** top-k keyword search algorithm, which takes into account the social, structural and semantical aspects of **S3**. Finally, we demonstrated the practical interest of our approach through experiments on three real social networks.

Next, we plan to extend **S3_k** to a massively parallel in-memory computing model to make it scale further. We also consider generating user-centric knowledge bases to be used in **S3_k**, to further adapt results to the user’s semantic perspective.

References

- [1] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [2] C. Aksoy, A. Dimitriou, and D. Theodoratos. Reasoning with patterns to effectively answer XML keyword queries. *The VLDB Journal*, 2015.

-
- [3] P. Buneman, E. V. Kostylev, and S. Vansummeren. Annotations are relative. In *Proceedings of the 16th International Conference on Database Theory*, pages 177–188. ACM, 2013.
- [4] L. J. Chen and Y. Papakonstantinou. Supporting top-k keyword search in XML databases. In *ICDE*, 2010.
- [5] M. Curtiss, I. Becker, T. Bosman, S. Doroshenko, L. Grijincu, T. Jackson, S. Kunnatur, S. Lassen, P. Pronin, S. Sankar, et al. Unicorn: A system for searching the social graph. *PVLDB*, 2013.
- [6] R. Fagin, R. Kumar, and D. Sivakumar. Comparing top k lists. *SIAM Journal on Discrete Mathematics*, 2003.
- [7] R. Fagin, A. Lotem, and M. Naor. Optimal aggregation algorithms for middleware. *Journal of Computer and System Sciences*, 66(4), 2003.
- [8] F. Goasdoué, K. Karanasos, Y. Katsis, J. Leblay, I. Manolescu, and S. Zampetakis. Growing triples on trees: an XML-RDF hybrid model for annotated documents. *VLDB Journal*, 2013.
- [9] F. Goasdoué, I. Manolescu, and A. Roatiş. Efficient query answering against dynamic RDF databases. In *EDBT*, 2013.
- [10] L. Guo, F. Shao, C. Botev, and J. Shanmugasundaram. XRANK: Ranked keyword search over XML documents. In *SIGMOD*, 2003.
- [11] V. Hristidis, L. Gravano, and Y. Papakonstantinou. Efficient IR-style keyword search over relational databases. In *VLDB*, 2003.
- [12] V. Hristidis, H. Hwang, and Y. Papakonstantinou. Authority-based keyword search in databases. *ACM TODS*, 33(1), 2008.
- [13] W. Le, F. Li, A. Kementsietsidis, and S. Duan. Scalable keyword search on large RDF data. *IEEE TKDE*, 26(11), 2014.
- [14] Y. Li, Z. Bao, G. Li, and K.-L. Tan. Real time personalized search on social networks. In *ICDE*, 2015.
- [15] S. Maniu and B. Cautis. Network-aware search in social tagging applications: instance optimality versus efficiency. In *CIKM*, 2013.
- [16] P. E. O’Neil, E. J. O’Neil, S. Pal, I. Cseri, G. Schaller, and N. Westbury. Ordpaths: Insert-friendly XML node labels. In *SIGMOD*, 2004.
- [17] M. Paradies, S. Malaika, J. Siméon, S. Khatchadourian, and K.-U. Sattler. Entity matching for semistructured data in the cloud. In *ACM SAC*, 2012.
- [18] A. Ritter, S. Clark, Mausam, and O. Etzioni. Named entity recognition in tweets: An experimental study. In *EMNLP*, 2011.
- [19] R. Schenkel, T. Crecelius, M. Kacimi, S. Michel, T. Neumann, J. X. Parreira, and G. Weikum. Efficient top-k querying over social-tagging networks. In *SIGIR*, 2008.
- [20] I. Tatarinov, S. D. Viglas, K. Beyer, J. Shanmugasundaram, E. Shekita, and C. Zhang. Storing and querying ordered XML using a relational database system. In *SIGMOD*, 2002.

- [21] M. Theobald, H. Bast, D. Majumdar, R. Schenkel, and G. Weikum. TopX: efficient and versatile top-k query processing for semistructured data. *The VLDB Journal*, 17(1), 2008.
- [22] M. Theobald, R. Schenkel, and G. Weikum. Efficient and self-tuning incremental query expansion for top-k query processing. In *SIGIR*, 2005.
- [23] T. Tran, H. Wang, S. Rudolph, and P. Cimiano. Top-k exploration of query candidates for efficient keyword search on graph-shaped (RDF) data. In *ICDE*, 2009.
- [24] J. Urbani, S. Kotoulas, J. Maassen, F. van Harmelen, and H. E. Bal. WebPIE: A web-scale parallel inference engine using MapReduce. *J. Web Sem.*, 10, 2012.
- [25] Resource Description Framework. <http://www.w3.org/RDF>.
- [26] Snowball stemming library for python. <https://pypi.python.org/pypi/PyStemmer>.
- [27] Tweepy library. <http://www.tweepy.org/>.
- [28] Uniform Resource Identifier. <http://tools.ietf.org/html/rfc3986>.
- [29] Yelp Dataset Challenge. http://www.yelp.com/dataset_challenge.
- [30] S. A. Yahia, M. Benedikt, L. V. Lakshmanan, and J. Stoyanovich. Efficient network aware search in collaborative tagging sites. *PVLDB*, 2008.



**RESEARCH CENTRE
SACLAY – ÎLE-DE-FRANCE**

1 rue Honoré d'Estienne d'Orves
Bâtiment Alan Turing
Campus de l'École Polytechnique
91120 Palaiseau

Publisher
Inria
Domaine de Voluceau - Rocquencourt
BP 105 - 78153 Le Chesnay Cedex
inria.fr

ISSN 0249-6399