

Pay-as-you-go Approximate Join Top-k Processing for the Web of Data – Technical Report

Andreas Wagner
AIFB, Karlsruhe Institute of
Technology
Karlsruhe, Germany
a.wagner@kit.edu

Veli Bicer
IBM Research, Smarter Cities
Technology Centre
Dublin, Ireland
velibice@ie.ibm.com

Thanh Duc Tran
AIFB, Karlsruhe Institute of
Technology
Karlsruhe, Germany
ducThanh.tran@kit.edu

ABSTRACT

In recent years, the amount of RDF data on the Web has drastically increased. For an effective search over such a large Web of data, ranking of results is crucial. To allow efficient query processing of ranked queries, *top-k join processing* has been proposed, which aims at computing k top-ranked results, without complete result materialization. However, Web search poses novel challenges. Most notably, users are frequently willing to sacrifice result accuracy in favor of result computation time. Thus, there is a strong need to *approximate the top-k results*. Previous work on approximate top- k processing, however, is not applicable for the join top- k setting – it solely targets the *selection top-k problem*. Further, existing approaches require *offline computed score statistics* to be available. Unfortunately, many important kinds of Web search queries, e.g., keyword or spatial queries, commonly *query-/user-dependent ranking* is employed. Thus, one only has score information at runtime. In this paper, we address these shortcomings and propose a novel approximate top- k join framework, well-suited for Web search queries and ranking. For this, all necessary score statistics are learned via a *pay-as-you-go training at runtime*. We study our approach in a theoretical manner, and *formally show its efficiency as well as effectiveness*. Further, we conducted experiments on benchmarks comprising synthetic as well as real-world Web data/-queries. Our results are very promising, as we could achieve up to 65% time savings, while maintaining a high precision/recall.

1. INTRODUCTION

With the proliferation of the Web of data, RDF (Resource Description Format) has gained a significant amount of attention. In fact, RDF has become an accepted standard for publishing data on the Web.¹ Adhering to a set of *triples* $\{\langle s, p, o \rangle\}$, RDF data forms a data graph, cf. Fig. 1-a. Every triple describes a particular entity (the *subject*) s through a *predicate/object* pair: p, o .

Web Search Ranking. For web-scale data, queries frequently produce a large number of results (*bindings*), thereby making *ranking* a key factor for an effective search. Each result has a ranking score associated, measuring its relevance for the user/information need. In a Web search setting, however, ranking functions of

ten need to *reflect on query constraints or user characteristics*, for judging a binding's relevance.

Example 1. Find movies with highest ratings, featuring an actress “Audrey Hepburn”, and playing close to Rome, cf. Fig. 1.

Exp. 1 would require a ranking function to incorporate the movie rating, quality of keyword matches for “Audrey Hepburn”, and distance of the movie's location to Rome. While one may assume that a higher rating value is preferred by any user and query, *scores for keyword and location constraint dependent on query and/or user characteristics*. For instance, in order to rank a triple t for “Audrey Hepburn”, a function may measure the edit distance between that keyword and triple t 's attribute value. Notice, given another keyword (e.g., only “Audrey”), t 's attribute value would yield a different score. Further, depending on the user's geographic knowledge of Italy, she may have different notions of “closeness” to Rome, e.g., distance ≤ 100 km. These three ranking criteria could be aggregated via a summation, Fig. 1-c. Ranking functions of that kind are commonly referred to as *query-/user-dependent ranking* [1, 7, 37], and are employed for many important Web search queries, e.g., keyword, spatial or temporal queries, cf. [4, 8, 9, 22, 23].

Join Top-k Processing. Most RDF query processing approaches focus on efficiency gains via query optimization, indexing or join techniques – *without taking ranking scores into account* [16].

On the other hand, work on relational databases has shown that result computation efficiency may be drastically improved by rank-aware query processing [18]. Rank-aware approaches (also known as *top-k query processing*) aim at computing k top-ranked bindings without full materialization of the result set. That is, after computing a certain number of bindings, the algorithm allows to terminate early, as it can guarantee that no binding with higher score, than those already found, exists [18].

Web search queries commonly feature one or more joins to combine multiple triples into a result [16]. Focusing on such queries, we target a particular flavor of top- k processing: the *join top-k problem*. Here, a complete binding is obtained by joining single triples, and the score of a complete binding is an aggregation of its comprised triples' scores [18].

Two recent papers proposed join top- k for the Web of data [25, 39]. However, these works miss to target common Web search behavior/information needs. That is, search on Web data is mostly performed by end-users having information needs that do not require a high result accuracy or completeness. In fact, while systems may compute a large number of top-ranked results, users oftentimes only investigate a small fraction of those results. For instance, in Exp. 1 a system could return the top-50 movie results. However, a user may just visit one or two results, until she discovers a movie of interest. Thus, users likely “miss” relevant results, because of their search behavior and information need. This drives the need for ap-

¹<http://webdatacommons.org>

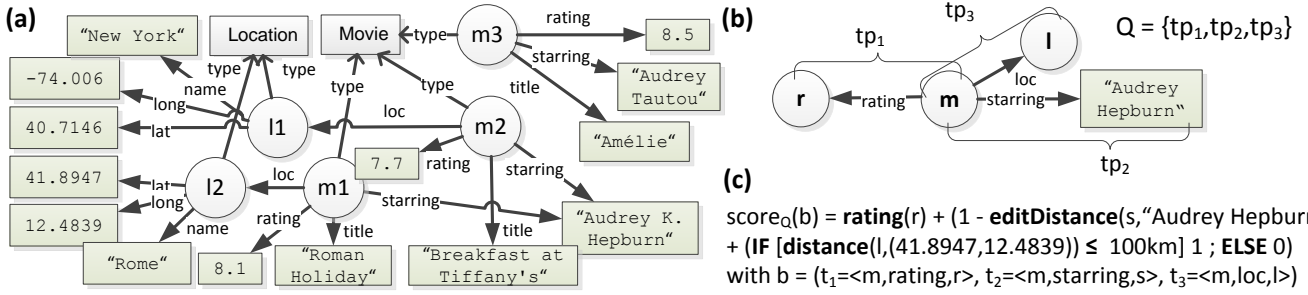


Figure 1: (a) Data about the movies “Roman Holiday”, “Breakfast at Tiffany’s”, and “Amélie”. (b) Query graph asking for a movie starring “Audrey Hepburn”. (c) Scoring function aggregating (via a summation) scores for triple pattern bindings (bold): movie ratings, edit distance w.r.t. “Audrey Hepburn”, and distance of the movie’s location to Rome (lat: 41.8947, long: 12.4839) ≤ 100 km. While the first constraint is query- and user-independent, the keyword constraint is query-dependent, and the location ranking is user-dependent, due to user-defined notion of “closeness” (here: ≤ 100 km).

proximated Web search query processing: *a system should be able to trade off result accuracy/completeness for computation time.*

Approximate Join Top-k Processing. Approximate top-k strategies, allowing false positive/negative top-k bindings, have been proposed before [2, 3, 28, 35, 38]. Unfortunately, these approaches are not well-suited Web search:

(P1) *Binding Probability.* RDF data adheres to a fine-grained triple model, and frequently many joins are required to form a query. However, existing approaches are not targeted at join top-k, but instead aim at the *selection top-k problem* [2, 3, 28, 35, 38]. Selection top-k is a very different task, as it computes top-ranked results, where each result is a “single” entity. In particular, no joins of multiple triples/entities are considered [18]. Applying such works for approximate join top-k is not straightforward, because those approaches do not capture the “binding probability” of a partial result: *Example 2.* Say a partial binding $b = (*, \langle m_2, \text{starring}, \text{"Audrey K. Hepburn"} \rangle, *)$ in Fig. 1-a matches the *starring* constraint, tp_2 , in Fig. 1-b, but still has to be joined with bindings for its unevaluated patterns (tp_1 and tp_3). The binding probability estimates how likely b will join with bindings for tp_1 and tp_3 , thereby contributing to one or more complete bindings.

(P2) *Score Probability.* To determine the probability for a partial binding to have a “sufficiently” large score to lead to a top-k result, score probability distributions are employed:

Example 3. Partial binding b in Exp. 2 has a known/fixed score for its triple $\langle m_2, \text{starring}, \text{"Audrey K. Hepburn"} \rangle$. However, bindings for pattern tp_1 and tp_3 are unknown, and their scores are modeled by means of a probability distribution.

Web search often uses an user/query-dependent ranking, e.g., for keyword or spatial queries. For these ranking functions, *scores are only known at runtime*, see Exp. 1. Existing approaches, however, require ranking scores to be available *at offline time*, in order to estimate and store the score probability distributions, e.g., via histograms [2, 3, 28, 35, 38]. Thus, for user/query-dependent ranking, previous works would need to *load all matching triples for each pattern into memory, compute their scores, and construct the necessary score statistics/distributions at runtime*. Notice, this must be done every time a query is issued. Such a procedure, however, would nullify the advantages of top-k processing, as materialization of bindings from disk is commonly the most expensive task in query processing.

This problem is exacerbated by the fact that score probability distributions must also capture aggregated scores from multiple patterns. For instance, in Exp. 3 we would need to estimate a distribution over scores for bindings of $tp_1 \bowtie tp_3$. As tp_3 is ranked via an user-dependent functions (Exp. 1), we again have no offline score information about the scores for $tp_1 \bowtie tp_3$. So, one would either

need to materialize/sample the join $tp_1 \bowtie tp_3$ at runtime, or (given a summation as score aggregation) compute the desired score distribution as a convolution of the two probability distributions for tp_1 and tp_3 . However, the former is associated with prohibitive costs, and the latter restricts the score aggregation to be a summation as well as imposes a harsh independence assumption.

In order to address problem P1 and P2, we propose an approximate join top-k strategy well-suited for the Web of data.

With regard to binding probabilities, *we reuse existing works on selectivity estimation* [24, 31, 32, 36]. Traditionally, selectivity estimation has been exploited to estimate the result set size for a given query, and to optimize join processing accordingly. In the following, we show how to make use of such solutions in order to approximate the binding probability.

For score probabilities, we propose a flexible approach based on Bayesian inferencing. In particular, we support any kind of ranking function, including query-/user-dependent functions. For this, *we use scores seen during join processing as “training examples” to iteratively learn score probability distribution at runtime.*

Contributions. Our contributions are as follows:

- While previous works solely aimed at top-k approximation for a relational setting, this is the first work towards an approximate join top-k processing for Web search.
- Our approach is lightweight, i.e., we do not require offline knowledge about ranking functions or score distributions. Every score statistic needed can be learned at runtime following the pay-as-you-go paradigm.
- We provide a theoretical analyses of our approach, and formally show its efficiency as well as effectiveness. In particular, we prove our distribution learning to have a constant space complexity, and a runtime complexity bounded by the result size. Furthermore, we give bounds for the approximation error of our overall approach.
- We implemented our approach and conducted experiments on two RDF benchmarks. Evaluation results are promising, as we could achieve time savings of up to 65%, while still having a high precision/recall.

Outline. We outline preliminaries in Sect. 2, and present the approximate top-k join in Sect. 3. In Sect. 4, we discuss evaluation results. Last, we give an overview over related works in Sect. 5, and conclude with Sect. 6.

2. PRELIMINARIES

Data Model. Given a set of edge labels ℓ , a RDF data is a labeled graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \ell)$, where $\mathcal{V} = \mathcal{V}_E \cup \mathcal{V}_A$ with entity nodes as \mathcal{V}_E and attribute nodes as \mathcal{V}_A . *Triples* are given by edges $\mathcal{E} = \{\langle s, p, o \rangle\}$, with $s \in \mathcal{V}_E$ as subject, $p \in \ell$ as predicate, and $o \in \mathcal{V}_E \cup \mathcal{V}_A$ as object. An example is depicted in Fig. 1-a.

Query and Result Model. We abstract from a particular type of queries, e.g., keyword or spatial queries, and model queries as basic graph patterns (BGPs), thereby capturing a key fragment of SPARQL.² A BGP query Q constitutes a directed labeled graph $Q = (\mathcal{V}^Q, \mathcal{E}^Q)$, with $\mathcal{V}^Q = \mathcal{V}_V^Q \cup \mathcal{V}_C^Q$ as disjoint union of variable, \mathcal{V}_V^Q , and constant nodes, \mathcal{V}_C^Q . Edges \mathcal{E}^Q are called *triple patterns*, with each pattern adhering to $tp = \langle s, p, o \rangle$, where $s \in \mathcal{V}_V^Q$, $p \in \mathcal{L}$, and $o \in \mathcal{V}_V^Q \cup \mathcal{V}_C^Q$. For example, pattern $\langle m, \text{starring}, \text{"Audrey Hepburn"} \rangle$ has one variable, m , one constant, "Audrey Hepburn", as object, and *starring* as predicate, Fig. 1-b. As shorthand we write Q as a *set of its triple patterns*: $Q = \{tp_i\}_i$.

A *binding* b for a query Q is a vector (t_1, \dots, t_n) of triples, such that each triple t_i *matches* (defined below) exactly one pattern tp_i in Q , and triples in b form a connected subgraph of \mathcal{G} . Via the matching of patterns in Q to triples, b *binds* variables to nodes in the data. Formally, for binding b there is a function $\mu_b : \mathcal{V}_V^Q \mapsto \mathcal{V}$, mapping every variable in Q to an entity/attribute value node in \mathcal{V} .

During query processing *partial* bindings, which feature some patterns with no matching triples, may occur. We refer to such a pattern, say tp_i , as *unevaluated*, and write $*$ in b 's i -th position: $(t_1, \dots, t_{i-1}, *, t_{i+1}, \dots, t_n)$. For a partial binding b we denote its evaluated pattern as $Q(b) \subseteq Q$, and its unevaluated pattern as $Q^u(b) = Q \setminus Q(b)$. b is *complete*, if all pattern have been evaluated.

A binding b *comprises* a binding b' , if any matched triple t_i in b' is also contained in b at position i . We say b' *contributes* to b .

Example 4. In Fig. 1-b, $Q = \{tp_1, tp_2, tp_3\}$. Partial binding $b_{31} = (*, *, t_{31} = \langle m_1, \text{loc}, l_1 \rangle)$ in Fig. 2-a matches pattern tp_3 , while $Q^u(b_{31}) = \{tp_1, tp_2\}$ are unevaluated. b_{31} binds variable m and l to entity m_1 and l_1 , respectively. Complete binding $b = (t_{12}, t_{21}, t_{31})$ comprises partial binding b_{31} , and b_{31} contributes to b .

Ranking Function. For quantifying the relevance of a binding b w.r.t. query Q , we employ a *ranking function*: $score_Q : \mathcal{B}^Q \mapsto \mathbb{R}$, with \mathcal{B}^Q as set of all partial/complete bindings for Q . More precisely, $score_Q(b)$ is given by an aggregation over b 's triples: $score_Q(b) = \bigoplus_{t \in b} score_Q(t)$, with \bigoplus as monotonic aggregation function. Function $score_Q$ is specific for Q , and the user who issued Q . Thus, we allow for user/query-dependent ranking of triple patterns bindings. A ranking function for our example is in Fig. 1-c. Note, $score_Q$ could be defined as *part of the query*, e.g., by means of an "extended" ORDER BY clause in SPARQL.

Sorted Access. For every pattern tp_i in query Q , a *sorted access* sa_i retrieves matching triples in descending score order. Previous works on join top- k over Web data discussed efficient sorted access implementations for RDF stores [25, 39]. Let us present simple approaches for our running example, Fig. 2-a:

Example 5. Given the keyword pattern $tp_2 = \langle m, \text{starring}, \text{"Audrey Hepburn"} \rangle$, a sorted access must materialize all triples, which have a value that contains "Audrey" or "Hepburn". After materialization, these triples are sorted with descending similarity w.r.t. that keyword – measured via, e.g., edit distance. Thus, sorted access sa_2 loads three triples comprising "Audrey" or "Hepburn", and sorts them according to their edit distance to "Audrey Hepburn". On the other hand, for pattern $\langle m, \text{loc}, l \rangle$, an R-tree on the attribute pair (lat, long) may be used [14]. This offline computed index yields two hits: l_1 and l_2 . While l_1 is an exact match, thus, ranking its triple t_{31} with max. score 1, l_2 is more distant from Rome. Note, l_2 and triple t_{32} , respectively, is only loaded if needed, i.e., if join-2 pulls on sa_3 a second time. An index for attribute *rating* can be constructed offline. Here, triples are sorted with descending attribute value. Sorted access sa_1 can simply iterate over this index, and materialize a triple each time join-1 pulls.

Partial bindings retrieved from sorted accesses are combined via joins. That is, an equi-join combines two (or more) inputs, each of them either being another join or a sorted access. This way, multiple joins form a tree. For instance, three sorted accesses are combined via two joins in Fig. 2-a.

It is important to notice: if a join input i produces a partial binding b , its set of unevaluated triple patterns, $Q^u(b)$, is the same for any binding from that input. Thus, we say input i has a set of unevaluated patterns $Q^u(i)$ associated. For instance, input i_1 has $Q^u(i_1) = \{tp_2, tp_3\}$ as unevaluated patterns, Fig. 2-a.

Problem. Our goal is to produce k high-ranked, complete query bindings that may differ from the true top- k results in terms of false positives/negatives. These approximations aim at saving computation time and input depth (#bindings pulled from sorted accesses). For this, we use a top- k test: for a given partial query binding, we estimate its probability for contributing to the final top- k results, and discard partial bindings having only a minor probability.

In our work, we exploit *conjugate priors* from the field of Bayesian statistics for learning necessary probability distributions.

Bayesian Inference. Let Θ be a set of parameters. One can model *prior beliefs* about these parameters in the form of probabilities: $\Theta \sim P(\Theta | \alpha)$, with $\Theta \in \Theta$ [15]. α is a vector of *hyperparameters* allowing to parametrize the prior distribution. Suppose we observe relevant data $\mathbf{x} = \{x_1, \dots, x_n\}$ w.r.t. Θ , where each $x_i \sim P(x_i | \Theta)$. Then, the dependency between observations \mathbf{x} and prior parameters Θ can be written as $P(\mathbf{x} | \Theta)$. Using the Bayes theorem we can estimate a *posterior* probability, capturing parameters Θ conditioned on observed events \mathbf{x} . In simple terms, a posterior distribution models how likely parameters Θ are, in light of the seen data \mathbf{x} , and the prior beliefs [15]:

$$P(\Theta | \mathbf{x}, \alpha) \propto P(\mathbf{x} | \Theta) \cdot P(\Theta | \alpha) = \frac{P(\mathbf{x} | \Theta) \cdot P(\Theta | \alpha)}{\sum_{\Theta} P(\mathbf{x} | \Theta) P(\Theta)} \quad (1)$$

Example 6. For the *rating* pattern tp_1 in Fig-2-a, scores are based on rating values. So, we can compute sufficient statistics, i.e. mean $\bar{x}_1 = 8.1$ and variance $s_1^2 = 0.16$, for these scores at offline time, cf. $stat_1$ in Fig-2-b. Such statistics represent prior beliefs about the true distribution, which is capturing only those scores observed for bindings of tp_1 that are part of a complete binding. Only triple t_{11} and t_{12} contribute to complete bindings, thus, only their scores should be modeled via a distribution. We update the prior beliefs using scores samples \mathbf{x} observed during query processing, thereby learning the true (posterior) distribution as we go.

As we are interested in estimating an *unobserved* event x^* , we need to calculate the *posterior predictive distribution*, i.e., the distribution of new events given observed data \mathbf{x} :

$$P(x^* | \mathbf{x}, \alpha) = \sum_{\Theta} P(x^* | \Theta) P(\Theta | \mathbf{x}, \alpha) \quad (2)$$

An important kind of Bayesian priors are the *conjugate priors* [15]:

Definition 1 (Conjugate Prior). A prior distribution family \mathcal{D} for a parameter set Θ is called *conjugate* iff $P(\Theta) \in \mathcal{D} \Rightarrow P(\Theta | \mathbf{x}) \in \mathcal{D}$.

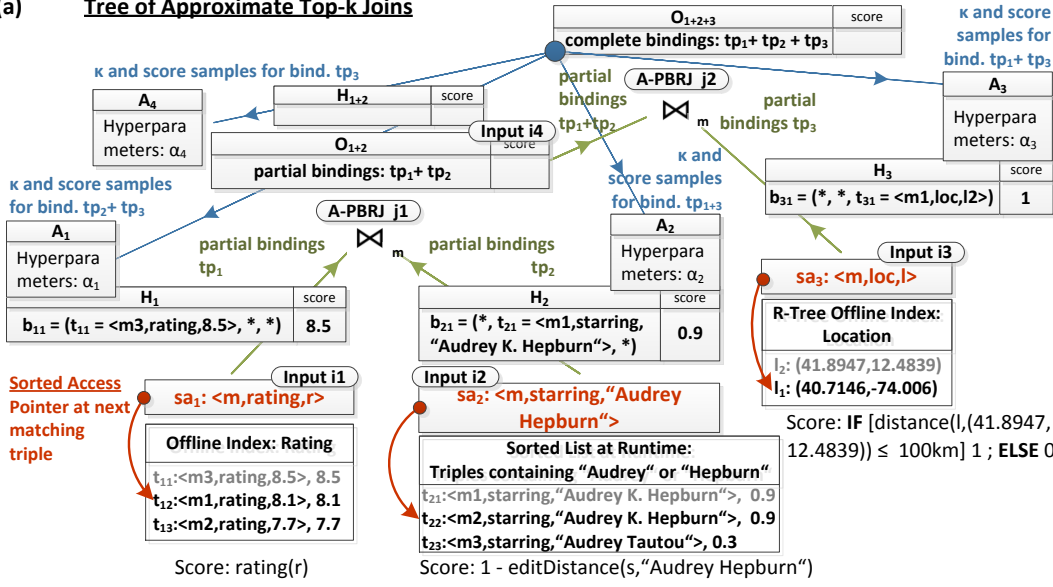
Intuitively, conjugate priors require the posterior and prior distribution to belong to the same distribution family. In other words, these priors provide an "computational convenience" as they give a closed-form of the posterior distribution [15]. Thus, posterior computation is easy and efficient.

3. APPROXIMATE TOP-K JOIN

We now present an *approximated* top- k join processing algorithm, allowing a trade-off between result computation efficiency and accuracy. For this, we extend the well-known Pull/Bound Rank Join [34] with a novel *probabilistic component* \mathcal{PC} .

²<http://www.w3.org/TR/rdf-sparql-query/>

(a) Tree of Approximate Top-k Joins



(b) Sufficient Statistics

Stat₁: Offline computed sample mean and variance for **rating ranking**
 $(\bar{x}_1, s_1^2) = (8.1, 0.16)$

Stat₂: Runtime computed sample mean and variance for **edit distance ranking**
 $(\bar{x}_2, s_2^2) = (0.7, 0.12)$

Stat₃: Runtime computed mean and variance for **location ranking** based on uniform distribution
 $(\bar{x}_3, s_3^2) =$
 $\left(\frac{(0+1)^2}{2}, \frac{(1-0)^2}{12} \right) =$
 $(0.25, 0.08)$

Figure 2: (a) Tree of approximate rank joins, combining three sorted accesses (one for each triple pattern in Fig. 1-b). Two information flows occur between operators in the tree: partial bindings (green), and score samples (blue). (b) Sufficient statistics calculated based on scores observed at indexing time (stat₁) and runtime (stat₂ and stat₃), respectively.

In particular, we will show that our instantiation of *PC* is *efficient and effective*. For the former, we show our score distribution training to have a constant space complexity, and a runtime complexity bounded by the result size (Lemma 1 and 2). For the latter, we discuss the quality of learned score distributions in Thm. 1 and 2, and provide an upper bound for approximation errors in Thm. 3.

3.1 Approximate Rank Join Framework

We define an approximate Pull/Bound Rank Join (A-PBRJ) comprising three parts: a pulling strategy *PS*, a bounding strategy *BS*, and a probabilistic component *PC*. *PS* determines the next join input to pull from [34]. The bounding strategy *BS* gives an upper bound, β , for the maximal possible score of unseen results for that join [34]. Last, we use *PC* to estimate a probability for a partial binding to contribute to the final top-*k* result.

Approximate Pull/Bound Rank Join. The A-PBRJ is depicted in Algo. 1. Following [34], on line 5 we check whether output buffer *O* comprises *k* complete bindings, and if there are unseen bindings with higher scores – measured via bound β . If both conditions hold, the A-PBRJ terminates and reports *O*. Otherwise, *PS* selects an input *i* to pull from (line 6), and produces a new partial binding *b* from the sorted access on input *i*, line 7. After materialization, we update β using bounding strategy *BS*.

Example 7. By means of strategy *PS*, join *j*₂ decides to first pull on sorted access *sa*₃, and materialize partial binding *t*₃₁, see Fig. 2-a. Then, join *j*₂ pulled on input *i*₄ (join *j*₁), which in turn pulled on its input *i*₁ (*sa*₁) loading binding *t*₁₁, and afterwards on input *i*₂ (*sa*₂) marginalizing *t*₂₁. *t*₁₁ \bowtie *t*₂₁ on variable *m* fails as *m*₃ \neq *m*₁.

In line 9, *PC* estimates the probability for partial binding *b* leading to a complete top-*k* binding: the top-*k* test. If *b* fails this top-*k* test *b* is pruned. That is, we do not attempt to join it, and do not insert it in *H*_{*i*}, which holds “seen” bindings from input *i*. Otherwise, if this test holds, *b* is further processed (lines 10 - 15). That is, we join *b* with seen bindings from the other input *j*, and add results to *O*. Further, *b* is inserted into buffer *H*_{*i*}, line 11. For learning the necessary probability distributions, *PC* trains on seen bindings/scores in *O*, line 13. Notice, we continuously train *PC* throughout the query processing – every time “enough” new bindings are in *O*, line 12. *PC* requires parameter κ for its pruning decision. κ holds

Algorithm 1: Approx. Pull/Bound Rank Join (A-PBRJ).

Params: Pulling strategy *PS*, bounding strategy *BS*, probabilistic component *PC*.

Index : Sorted access *sa_i* and *sa_j* for input *i* and *j*, resp..

Buffer : Output buffer *O*. *H_i* and *H_j* for “seen” bindings pulled from *sa_i* and *sa_j*. Hyperparameter buffer *A*.

Input : Query *Q*, result size *k*, and top-*k* test threshold τ .

Output : Approximated top-*k* result.

```

1 begin
2   PC.initialize()
3    $\beta \leftarrow -\infty$ 
4    $\kappa \leftarrow -\infty$ 
5   while |O| < k or  $\min_{b' \in O} \text{score}_Q(b') < \beta$  do
6     i  $\leftarrow$  PS.input()
7     b  $\leftarrow$  next triple pattern binding from sai
8      $\beta \leftarrow$  BS.update(b)
9     // top-k test
10    if PC.probabilityTopK(b,  $\kappa$ ) >  $\tau$  then
11      Hj  $\bowtie$  {b} and add each binding to O
12      Add b to Hi
13      if #new bindings b in O  $\geq$  training threshold then
14        PC.train(b)
15        Retain only k top-ranked bindings in O
16      if |O|  $\geq$  k then  $\kappa \leftarrow \min_{b' \in O} \text{score}_Q(b')$ 
17  // return approximated top-k results
18  return O

```

the the smallest currently known top-*k* score (line 15). Note, κ is initialized as $-\infty$, line 4. See Fig. 2-a for a tree of A-PBRJs.

Choices for *PS* and *BS*. Multiple works proposed bounding, e.g., [12, 17, 26, 34], and pulling strategies, e.g., [17, 27]. Commonly, the corner bound [17] is employed as bounding strategy:

Definition 2 (Corner Bound). For a join operator, we maintain *u_i* and *l_i* for each input *i*. *u_i* is the highest score observed from *i*, while *l_i* is the lowest observed score on *i*. If input *i* is exhausted, *l_i* is set to $-\infty$. The bound for scores of unseen join results is $\beta =$

$\max\{u_1 \oplus l_2, u_2 \oplus l_1\}$.

In example Fig. 2-a, join j_1 currently has $\beta = \max\{8.5 + 0.9, 0.9 + 8.5\}$, with $u_1 = l_1 = 8.5$, and $u_2 = l_2 = 0.9$.

On the other hand, the *corner-bound-adaptive strategy* [17] is frequently used as pulling strategy \mathcal{PS} :

Definition 3 (Corner-Bound-Adaptive Pulling). *The corner-bound-adaptive pulling strategy chooses the input i such that: $i = 1$ iff $u_1 \oplus l_2 > u_2 \oplus l_1$, and $i = 2$ otherwise. If $u_1 \oplus l_2 = u_2 \oplus l_1$, the input with the smaller # unseen partial bindings is chosen.*

For instance, as $8.5 + 0.9 = 0.9 + 8.5$ in join j_1 , with both inputs having 2 unseen partial bindings, either input may be selected.

3.2 Probabilistic Component

The probabilistic component implements the top- k test, which follows a simple intuition: For a partial binding b , a subset of query patterns are evaluated, $Q(b)$, while some others are not, $Q^u(b)$. Thus, scores (bindings) of unevaluated patterns are uncertain. However, partial bindings for already evaluated patterns restrict the space of likely/possible scores (bindings) for unevaluated patterns.

We model these binding and score uncertainties via two probabilities: (1) Probability for partial binding b contributing to one or more complete bindings (*binding probability*). (2) Probability distribution over scores of complete bindings (*score probability*).

Binding Probability. To address the former probability, we use a selectivity estimation function *sel*. Simply put, given a query Q , *sel*(Q) estimates the probability that there is at least one binding for Q , see [24, 31, 32, 36]. For example, selectivity of pattern $tp_3 = \langle m, \text{loc}, l \rangle$ is *sel*(tp_3) = $\frac{2}{3}$, because out of the three movie entities only two have a *loc* predicate, cf. Fig. 1-a.

By means of function *sel*, we define a *complete binding indicator* for partial binding b :

$$\mathbb{1}(Q^u(b) \mid b) := \begin{cases} 1 & \text{if } \text{sel}(Q^u(b) \mid b) > 0 \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

Intuitively, $\mathbb{1}(Q^u(b) \mid b)$ models whether matching triples for b 's remaining unevaluated patterns can exist, given variable assignments dictated by b .

That is, $Q^u(b) \mid b$ is a set of patterns $\{\bar{tp}_i\}$, such that pattern $tp_i \in Q^u(b)$, and each variable v in tp_i that is bound by b is replaced with its assignment in b , $\mu_b(v)$, resulting in a new pattern \bar{tp}_i .

Example 8. Consider Fig. 2-a and partial binding $b_{11} = (t_{11} = \langle m_3, \text{rating}, 8.5 \rangle, *, *)$. Here, $Q^u(b_{11}) \mid b_{11}$ is $\{\langle m_3, \text{starring}, \text{"Audrey Hepburn"} \rangle, \langle m_3, \text{loc}, l \rangle\}$, which is obtained by replacing variable m in tp_2 and tp_3 with its assignment in b_{11} ($\mu_{b_{11}}(m) = m_3$).

It is important to notice that the complete binding indicator (Eq. 3) is independent of a particular *sel* implementation – any selectivity estimation for BGP queries may be used.

For our experiments we reused work from [31, 32]. The authors employed indexes for triple patterns with two constants: SP, PS, SO, OS, PO, and OP. Each index maps a $\langle val_1, val_2 \rangle$ pair to its cardinality, i.e., the number of its matching triples in the data. For instance, $\langle m_1, \text{starring} \rangle$ would map to 1 in Fig. 1-a. However, the binding indicator only requires a selectivity estimation function to make a boolean decision: either $\text{sel}(Q^u(b) \mid b) > 0$ or not. Thus, not all six indexes are necessary: SP, PO, and SO are sufficient. A rudimentary implementation of $\text{sel}(Q^u(b) \mid b)$ returns 1 iff

$$\forall_{\langle s, p, o \rangle \in Q^u(b) \mid b, \text{ where } s \in \mathcal{V}_E \text{ SP.card}(\langle s, p \rangle) > 0 \quad \wedge \quad (4a)$$

$$\forall_{\langle s, p, o \rangle \in Q^u(b) \mid b, \text{ where } o \in \mathcal{V}_E \cup \mathcal{V}_A \text{ PO.card}(\langle p, o \rangle) > 0 \quad \wedge \quad (4b)$$

$$\forall_{\langle s, p, o \rangle \in Q^u(b) \mid b, \text{ where } s \in \mathcal{V}_E \wedge o \in \mathcal{V}_E \cup \mathcal{V}_A \text{ SO.card}(\langle s, o \rangle) > 0 \quad (4c)$$

and 0 otherwise, with *card* as cardinality function.

Example 9. For $Q^u(b_{11}) \mid b_{11}$ in Exp. 8, $\mathbb{1}(Q^u(b_{11}) \mid b_{11}) = 0$, because selectivity for both patterns is 0. In our implementation probe in SP for $\langle m_3, \text{loc} \rangle$ returns “pair does not exist”. Thus, our *sel*($Q^u(b_{11}) \mid b_{11}$) correctly returns 0.

Score Probability. Given a partial binding b , let X_b^s be a random variable for modeling scores of complete bindings comprising b . Further, scores for bindings of b 's unevaluated patterns, $Q^u(b)$, are captured via $X_{Q^u(b)}^s$. Then, it holds:

$$P(X_b^s \geq x) = P(X_{Q^u(b)}^s \geq \delta(x, b)) \quad (5)$$

where $\delta(x, b) := x - \text{score}_Q(b)$. LHS of Eq. 5 may be simplified to the RHS, because partial binding b already has a “certain” score, $\text{score}_Q(b)$, and only the scores for its unevaluated patterns are unknown. So, $\delta(x, b)$ may be seen as “delta” between b 's score and a desired score x . For instance, b_{31} has a score of 1, however, scores for bindings to tp_1 and tp_2 are uncertain, and modeled via $X_{Q^u(b_{31})}^s$.

Top- k Test. Having introduced above probabilities, we can define the top- k test. For a partial binding b , we are interested in b 's probability of contributing to a complete binding in the top- k results. Thus, we use (1) the binding probability, i.e., the complete binding indicator, to determine whether b might contribute to any complete binding, and (2) the score probability to estimate scores for complete bindings comprising b :

$$P(X_b^s \geq \kappa) \cdot \mathbb{1}(Q^u(b) \mid b) > \tau \Leftrightarrow \quad (6a)$$

$$P(X_{Q^u(b)}^s \geq \delta(\kappa, b)) \cdot \mathbb{1}(Q^u(b) \mid b) > \tau \quad (6b)$$

with κ as the smallest currently known top- k score (see Algo. 1 line 15), and $\tau \in [0, 1]$ as test threshold.

Discussion. Threshold τ provides a key instrument for Web search systems, as it always to adjust result accuracy, depending on user and information need. For instance, a system may decide to compute top-50 results in total, and increase τ every time a new top-ranked binding could be reported. Generally speaking, τ should not be conceived as a constant, but rather as a function in reported top-ranked results. This way, systems may reflect on the typical user behavior, who frequently only visit a small subset of the computed results.

Second, note that κ in Eq. 6 refers to the smallest currently known top- k binding score. Thus, as long as no k complete bindings have been found, κ is set to $-\infty$ (Algo. 1, line 4), and the score probability is always 1. So, the only reason for a partial binding b to be pruned is, if it fails the complete binding indicator. That is, if b is not expected to contribute to any complete binding.

Last, each top- k test computation causes costs in the form of probability estimations. We will provide empirical results regarding those costs in our evaluation. On the other hand, recent work introduced a “cost-aware” rank join, which schedules sorted and random accesses based on their associated costs [27] – this line of work can be directly applied here. In fact, the top- k test may be treated as “one more” access in their optimization problem [27].

3.3 Score Distribution Learning

With regard to the binding probability, we only require a selectivity estimation function. For this, *all necessary indexes can be computed offline*. We leave the selectivity estimation problem to orthogonal work, and focus on score probabilities in this paper.

Estimation of score probabilities is rather difficult: (1) In a join top- k setting, we are only interested in scores for complete bindings. Consider join j_1 in Fig. 2-a: we are not interested in “any” aggregation of two scores for bindings of tp_1 and tp_2 , e.g., $\text{score}_Q(t_{11}) + \text{score}_Q(t_{21}) = 8.5 + 0.9$. Instead, we wish to capture only scores of valid join results, e.g., $\text{score}_Q(t_{11}) + \text{score}_Q(t_{23})$. (2) We allow

| (a) Predictive Distributions | (b) Priors |
|---|---|
| Input 11 $P(X_{Q^u(i_1)}^s)$ with $Q^u(i_1) = \{tp_2, tp_3\}$ | $\text{stat}_2 \oplus \text{stat}_3$ (0.7 + 0.25, 0.12 + 0.08) = (0.95, 0.20) |
| Input 12 $P(X_{Q^u(i_2)}^s)$ with $Q^u(i_2) = \{tp_1, tp_3\}$ | $\text{stat}_1 \oplus \text{stat}_3$ (8.1 + 0.24, 0.16 + 0.08) = (8.35, 0.24) |
| Input 13 $P(X_{Q^u(i_3)}^s)$ with $Q^u(i_3) = \{tp_1, tp_2\}$ | $\text{stat}_1 \oplus \text{stat}_2$ (8.1 + 0.7, 0.16 + 0.12) = (8.8, 0.28) |
| Input 14 $P(X_{Q^u(i_4)}^s)$ with $Q^u(i_4) = \{tp_3\}$ | stat_3 (0.25, 0.08) |

Figure 3: (a) Four predictive score distributions, one for each join input, together with their priors. For instance, $X_{Q^u(i_1)}^s$ models scores for bindings to $tp_2 \bowtie tp_3$, which are comprised in complete bindings. (b) Priors are based on sufficient statistics. Note, aggregation function \oplus in Fig. 1-c is a summation. For instance, $\text{stat}_1 \oplus \text{stat}_3$ becomes (8.1 + 0.24, 0.16 + 0.08).

for user/query-dependent ranking. For those functions, no offline score statistics can be constructed. At the same time, for other ranking functions, e.g., based on the `rating` predicate, we do have offline score statistics. Thus, we aim for a flexible framework allowing to incorporate any offline information.

3.3.1 Overview

We estimate probabilities for $X_{Q^u(b)}^s$ by approximating its true distribution with $X_{Q^u(i)}^s$ with i as the input from which b was pulled, and $Q^u(i)$ as set of i 's unevaluated patterns:

$$P(X_{Q^u(b)}^s) \approx P(X_{Q^u(i)}^s)$$

Example 10. For instance, in Fig. 2-a we approximate $P(X_{Q^u(b_{11})}^s)$ with $P(X_{Q^u(i_1)}^s)$, as binding b_{11} is produced by input i_1 . $X_{Q^u(i_1)}^s$ captures scores for bindings to patterns $\{tp_2, tp_3\}$.

So, instead of learning a distribution for every partial binding b , we train a score distribution $X_{Q^u(i)}^s$ for each input i . For our example we learn four distributions, cf. Fig. 3.

Relying on $X_{Q^u(i)}^s$, we provide an implementation of $\mathcal{PC}.\text{initialize}()$, $\mathcal{PC}.\text{train}()$ in Algo. 2, and $\mathcal{PC}.\text{probabilityTopK}()$ in Algo. 3 – as required by our A-PBRJ framework in Algo. 1.

3.3.2 Pay-as-you-go Distribution Learning

We use conjugate priors, which allow a pay-as-you-go distribution learning for $X_{Q^u(i)}^s$. With runtime ranking scores, we have no information about the true distribution of $X_{Q^u(i)}^s$. In such a scenario, a common assumption is to use a Gaussian score distribution. Thus, we employ a Gaussian conjugate prior (Eq. 7a). However, we also outline how to extend our approach to other distributions, if offline knowledge about $X_{Q^u(i)}^s$ is available (Sect. 3.3.3).

Using a Gaussian conjugate prior with unknown mean and unknown variance, prior and posterior distribution can be decomposed as: $P(\mu, \sigma^2 | \alpha) = P(\mu | \sigma^2, \alpha) \cdot P(\sigma^2 | \alpha)$ [15]. The mean μ follows a Gaussian, Eq. 7b, and the variance σ^2 a inverse-Gamma distribution, Eq. 7c. Hyperparameters $\alpha_0 = (\mu_0, \eta_0, \sigma_0^2, \nu_0)$ parameterize both distributions, where μ_0 is prior mean with quality η_0 , and σ_0^2 is prior variance with quality ν_0 [15]:

$$X_{Q^u(i)}^s \sim \text{normal}(\mu, \sigma^2) \quad (7a)$$

$$\mu | \sigma^2 \sim \text{normal}\left(\mu_0, \frac{\sigma^2}{\eta_0}\right) \quad (7b)$$

$$\sigma^2 \sim \text{inverse-gamma}(0.5 \cdot \nu_0, 0.5 \cdot \nu_0 \sigma_0^2) \quad (7c)$$

Prior Distribution. Prior initialization is implemented by means of the $\mathcal{PC}.\text{initialize}()$ method in Algo. 1 (line 2). For each input i we specify a prior distribution for $X_{Q^u(i)}^s$ via prior hyperparameters α_0 . For α_0 we require sufficient score statistics in the form of a sample mean, $\bar{x} = \frac{1}{n} \sum_{x_i \in \mathbf{x}} x_i$, and a sample variance $s^2 = \frac{1}{(n-1)} \sum_{x_i \in \mathbf{x}} (x_i - \bar{x})^2$, with \mathbf{x} as sample. There are multiple ways to obtain the necessary score samples, depending on what kinds of ranking functions are used:

Example 11. Fig. 2-b depicts three sufficient statistics based on information from the sorted accesses: (1) Offline information in the case of sa_1 . That is, scores are known before runtime, thus, $\bar{x}_1 = 8.1$ and $s_1^2 = 0.16$ can be computed offline. (2) Online information for access sa_2 . Recall, the list of matching triples for keywords “Audrey” and “Hepburn” must be fully materialized. So, $\bar{x}_2 = 0.7$ and $s_2^2 = 0.12$ may be computed from runtime score samples. (3) Last, given access sa_3 , we have neither offline scores, nor a fully materialized list of triples (sa_3 loads a triple solely upon a pull request). In lack of more information, we assume each score to be equal likely, i.e., a uniform distribution. With min. score as 0 and max. score as 1: $\bar{x}_3 = 0.25$ and $s_3^2 = 0.08$.

Algorithm 2: $\mathcal{PC}.\text{train}()$

Params: Weight $w \geq 1$ for score sample \mathbf{x} .

Buffer : Buffer **A** storing hyperparameters.

Input : Complete bindings **B** \subseteq **O**, and join j .

```

1 begin
  // train hyperparameters for each input
2 foreach input  $i$  in join  $j$  do
  // load prior hyperparameters for input  $i$ 
3    $\alpha_n = (\mu_n, \eta_n, \sigma_n^2, \nu_n) \leftarrow \mathbf{A}_i$ 
  // get scores of bindings of unevaluated patterns
4   foreach complete binding  $b \in \mathbf{B}$  do
5     get  $b'$  comprised in  $b$  for patterns  $Q^u(i)$ 
6     add  $\text{score}_Q(b')$  to score sample  $\mathbf{x}$ 
  // sample mean and variance
7    $\bar{x} \leftarrow \text{mean}(\mathbf{x}) = \frac{1}{n} \sum x_i$ 
8    $s^2 \leftarrow \text{var}(\mathbf{x}) = \frac{1}{(n-1)} \sum (x_i - \bar{x})^2$ 
  // posterior hyperparameters
9    $\nu_{n+1} \leftarrow \nu_n + w$ 
10   $\eta_{n+1} \leftarrow \eta_n + w$ 
11   $\mu_{n+1} \leftarrow \frac{\eta_n \mu_n + w \bar{x}}{\eta_{n+1}}$ 
12   $\sigma_{n+1}^2 \leftarrow \frac{1}{\nu_{n+1}} \cdot (\nu_n \sigma_n^2 + (w-1)s^2 + \frac{\eta_n w}{\eta_{n+1}} \cdot (\bar{x} - \mu_n)^2)$ 
  // store hyperparameters for input  $i$  in join  $j$ 
13   $\mathbf{A}_i \leftarrow \alpha_{n+1} = (\mu_{n+1}, \eta_{n+1}, \sigma_{n+1}^2, \nu_{n+1})$ 

```

Similar to [15], we initialize hyperparameters α_0 with: μ_0 as sample mean, σ_0^2 as sample variance, and $\eta_0 = \nu_0$ as sample size/quality. For every input, we aggregate necessary sample means and variances for μ_0 and σ_0^2 , respectively. The aggregation function for these statistics must be the same as for the ranking function, \oplus . For example, given input i_1 where $Q^u(i_1) = \{tp_2, tp_3\}$, we sum up (aggregate) statistics stat_2 and stat_3 : $\bar{x}_2 + \bar{x}_3$ for μ_0 , and $s_2^2 + s_3^2$ for σ_0^2 . Similarly, priors for the remaining three distributions can be calculated, see Fig. 3. Further, η_0 and ν_0 are used to quantify the prior quality. For instance, stat_1 and stat_2 are exact statistics, while stat_3 relies on a uniform distribution. So, a prior may be weighted depending on its employed statistics. In other words, weighting reflects the “trustworthiness” of the prior.

Posterior Distribution. Having estimated a prior distribution, we continuously update the distribution with seen scores during query processing. Based on [15], the training procedure $\mathcal{PC}.\text{train}()$

is given in Algo. 2.

Intuitively, each time new complete bindings (Algo. 1, line 12) are produced, all distributions $X_{Q^{(i)}}^s$ could be trained. That is, complete binding scores are used to update hyperparameters from the previous n -th training iteration, α_n , resulting in new posterior hyperparameters, α_{n+1} . For this, we use standard training (Algo. 2, lines 11-12) as discussed in [15]:

$$\mu_{n+1} = \frac{\eta_n \mu_n + w \bar{x}}{\eta_{n+1}} \quad (8a)$$

$$\sigma_{n+1}^2 = \frac{1}{\nu_{n+1}} \cdot \left(\nu_n \sigma_n^2 + (w-1)s^2 + \frac{\eta_n w}{\eta_{n+1}} \cdot (\bar{x} - \mu_n)^2 \right) \quad (8b)$$

In simple terms, the prior mean μ_n is updated with the new sample mean \bar{x} , Eq. 8a, and the prior variance σ_n^2 is updated with the sample variance s^2 , Eq. 8b. Each input “extracts” its own score sample \mathbf{x} (Algo. 2, lines 5-6), from which the sample mean and variance is computed. This is because every $X_{Q^{(i)}}^s$ models scores for different “unevaluated” patterns.

Prior hyperparameters are weighted via η_n and ν_n . Further, for each hyperparameter update, a parameter w is used as weight, which indicates the quality of the score samples \mathbf{x} . In our experiments, we set w as sample size. Finally, new hyperparameters α_{n+1} are stored (Algo. 2, line 13), and used as prior for the next training.

Example 12. Consider input i_1 , and say $\eta_0 = \nu_0 = 1$, its prior is $\alpha_0 = (0.95, 1, 0.20, 1)$, cf. Fig. 3. We observe scores $\mathbf{x} = \{x_1, x_2\}$ from $\mathbf{B} = \{(t_{12}, t_{21}, t_{31}), (t_{13}, t_{22}, t_{32})\}$, with $w = |\mathbf{x}| = 2$, $x_1 = 1.9 = \text{score}_Q(t_{21}) + \text{score}_Q(t_{31})$, and $x_2 = 0.9 = \text{score}_Q(t_{22}) + \text{score}_Q(t_{32})$. Then, $s^2 = 0.26$, $\bar{x} = 1.4$, and posterior hyperparameters are:

$$\eta_1 = \nu_1 = 1 + 2 = 3$$

$$\sigma_1^2 = \frac{1}{3} \cdot \left(0.2 + (2-1) \cdot 0.26 + \frac{(1.4 - 0.95)^2}{3} \right) = 0.26$$

$$\mu_1 = \frac{(0.95 + 2 \cdot 1.4)}{3} = 1.25$$

After each such update only posterior hyperparameters are stored, thereby making the learning highly space efficient:

Lemma 1 (Distribution Learning Space Complexity). *Given an A-PRBJ operator j , at any time during query processing, we require $O(1)$ of space for score distribution learning.*

Proof Sketch. Given an A-PRBJ operator, for every of its inputs a single parameter vector (hyperparameters α) is stored after a training iteration (Algo. 2, line 13). As each vector, α , features a fixed number of parameters, the space consumption remains constant, i.e., independent of #training iterations. ■

Further, for the learning time complexity we can show:

Lemma 2 (Distribution Learning Time Complexity). *Given an A-PRBJ operator j , and \mathbf{B} complete bindings, score learning time complexity is bounded by $O(|\mathbf{B}|)$.*

Proof Sketch. Given an A-PRBJ operator, and a set of new results \mathbf{B} : A score sample, \mathbf{x} , is constructed (Algo. 2, lines 5-6) with $O(|\mathbf{B}|)$ complexity. Mean and variance is computed from \mathbf{x} in $O(|\mathbf{x}|)$ time. However, as $|\mathbf{x}| \leq |\mathbf{B}|$, it holds that $O(|\mathbf{x}|) \in O(|\mathbf{B}|)$. Notice, computation of mean and variance could also be done while collecting the sample in lines 5-6. Further, hyperparameters are updated via \mathbf{x} in constant time (Algo. 2, lines 9-12). Thus, every training iteration costs overall $O(|\mathbf{B}|)$. ■

Predictive Distribution. In Algo. 3, we provide an implementation of the $\mathcal{PC}.\text{probabilityTopK}()$ method, see Algo. 1 line 9. At any point during query processing, one may need to perform a top- k test. For this, our approach allows to always give a distribution for

$X_{Q^{(i)}}^s$ based on the currently known hyperparameters α_n (Algo. 3, line 2). As hyperparameters are continuously trained, the distributions improve over time.

More specifically, we use the posterior predictive distribution. This distribution estimates probabilities for *new* scores, based on observed scores, and prior distribution. As shown in [15], this distribution can be easily obtained in a closed form as non-standardized Student’s t -distribution with ν_n degrees of freedom, see Algo. 3, line 3. Being able to compute the posterior predictive in such an easy manner is an advantage of a Gaussian conjugate prior – for other priors this estimation may be more complex [15].

Using the posterior predictive distribution, we may give an approximation of $P(X_{Q^{(b)}}^s)$ via $P(X_{Q^{(i)}}^s)$, Algo. 3, line 4. Last, in line 5 we compute the binding probability (Eq. 3) by means of the selectivity estimation function. Now we have obtained binding as well as score probability, and can return a probability for b contributing to the top- k results, line 6 (Eq. 6).

Algorithm 3: $\mathcal{PC}.\text{probabilityTopK}()$

Buffer : Buffer \mathbf{A} storing hyperparameters.

Input : Partial bindings b , input i , and join j .

Output : Probability that b will result in one (or more) final top- k bindings.

```

1 begin
  // load hyperparameters for input  $i$  at join  $j$ 
2  $\alpha_n = (\mu_n, \eta_n, \sigma_n^2, \nu_n) \leftarrow \mathbf{A}_i$ 
  // posterior predictive distribution based on current  $\alpha_n$ 
3  $X_{Q^{(i)}}^s \sim t_{(\nu_n)} \left( x \mid \mu_n, \frac{\sigma_n^2(\eta_n+1)}{\eta_n} \right)$ 
  // approximate score probability
4  $p_S \leftarrow P(X_{Q^{(b)}}^s \geq \delta(\kappa, b)) \approx P(X_{Q^{(i)}}^s \geq \delta(\kappa, b))$ 
  // compute binding probability
5  $p_B \leftarrow \mathbb{1}(Q^u(b) \mid b)$ 
  // probability that  $b$  contributes to top- $k$ 
6 return  $p_S \cdot p_B$ 
```

3.3.3 Discussion

Refined Conjugate Priors. Whenever one has offline information about the true distribution of $X_{Q^{(i)}}^s$ (or good approximation for it), we may replace the Gaussian conjugate prior in Eq. 7 with a more suitable one. For this, only minor changes in the training (Algo. 2), and predictive distribution estimation (Algo. 3, line 3) are required. Both these tasks are well-known problems in the literature [15]. No further modifications are needed – the top- k test works with any valid score distribution for $X_{Q^{(i)}}^s$.

A wide variety of discrete/continuous conjugate priors are known, thus, in the best case, there is a conjugate prior for the true distribution of $X_{Q^{(i)}}^s$. Otherwise, if no matching conjugate prior exists, we can exploit a mixture of multiple conjugate priors: $\sum_i w_i P_i(\theta)$ with each $P_i(\theta)$ being a conjugate prior, and w_i as weights such that $\sum_i w_i = 1$ and $0 < w_i < 1$. Notice, [5] showed that any distribution from the exponential family could be approximated (arbitrarily close) by means of a mixed conjugate prior.

Maintenance. We require maintenance of binding and score probabilities. Binding probabilities are estimated via a selectivity estimation function, Eq. 3. Maintenance of these statistics varies with the specific selectivity estimation implementation, and is out of scope for this work.

With regard to score probabilities, we train posterior distributions during query processing. Thus, only for prior distributions sufficient statistics are needed. These statistics may differ depending on the ranking functions employed. For example, given the

rating ranking in Fig. 1-c, sufficient statistics can be computed before runtime in the form of a sample mean and variance: $stat_1 = (\bar{x}_1 = 8.1, s_1^2 = 0.16)$, Fig. 2-b. In fact, one may even store further distribution characteristics, e.g., distribution skewness or symmetry. This way, more refined conjugate priors could be estimated (see paragraph above). In contrast, for user-/query-dependent ranking, e.g., the keyword constraint “Audrey Hepburn” in Fig. 1-c, scores are unknown before runtime. Thus, no sufficient statistic can be stored and/or maintained. However, in such a case, a minimal and maximal score may be kept. For instance, for the distance ranking constraint, we would store a minimal and maximal score as 0 and 1, respectively. This way, we may assume a uniform score distribution as naive prior, and compute mean and variance as: $stat_3 = (\bar{x}_3 = 0.25, s_3^2 = 0.08)$, cf. Fig. 2-b and Exp. 11.

3.4 Theoretical Analysis

In this section, we discuss theoretical aspects concerning the *effectiveness of the A-PBRJ operator*. That is, we discuss the quality of the learned score distributions (Thm. 1 and 2), and provide bounds for the approximation error, cf. Thm. 3.

Distribution Quality. The aggregation function \oplus for our ranking function $score_Q$ could be any monotonic function. However, when we restrict the aggregation to a summation (see Fig. 1-c), we can formally show that a Gaussian distribution/conjugate prior in Eq. 7 is a good approximation for the true distribution of $X_{Q^u(i)}^s$. Notice, many common aggregations employ summations, e.g., TF-IDF inspired functions may be represented by summations [38]. For such a summation-based aggregation function it holds:

Theorem 1. *Given a query $Q = \{tp_k\}$ and $X_{Q^u(i)}^s = \sum_{tp_k \in Q^u(i)} X_{tp_k}^s$, the Central Limit Theorem (CLT) holds:*

$$\frac{\sum_k (X_{tp_k}^s - \mu_k)}{\sqrt{\sum_k \sigma_k^2}} \xrightarrow{n \rightarrow \infty} \text{normal}(0, 1)$$

with n as the number of patterns in $Q^u(i)$, and $X_{tp_k}^s$ as random variable modeling scores of bindings for pattern tp_k . Further, μ_k and σ_k^2 stand for the finite mean and variance of $X_{tp_k}^s$.

Informal Proof Sketch. As we do not have knowledge about the ranking functions, $score_Q$, or the distributions for $X_{tp_k}^s$, we can only outline a very informal sketch of proof in the following.

Our argumentation is based on two assumptions: (A1) Recall, we define a *separate* ranking function $score_Q$ for every triple pattern tp_k in query Q , cf. Sect. 2. In particular, each function computes its score solely by considering the triple binding of “its own” pattern, tp_k . Thus, assuming scores from different ranking functions to be independently distributed is a reasonable simplification. Formally, for every two pattern tp_1 and tp_2 in Q , we assume: $X_{tp_1}^s \perp X_{tp_2}^s$. (A2) We assume each random variable $X_{tp_k}^s$ to have a finite mean μ_k and variance σ_k^2 . Notice, most common distributions feature a finite mean and variance. So, this assumption does not restrict $score_Q$ and $X_{tp_k}^s$, respectively.

In its simplest form, the Central Limit Theorem is only applicable to i.i.d. random variables [13, p. 329]. However, in the Lindeberg Theorem this restriction is lifted, i.e., every variable $X_{tp_k}^s$ may adhere to a different distribution [13, p. 330]:

Lemma 3 (Lindeberg Condition). *If*

$$\lim_{n \rightarrow \infty} \frac{1}{s_n^2} \sum_k E \left((X_{tp_k}^s - \mu_k)^2 \right) \cdot \mathbb{1}(|X_{tp_k}^s - \mu_k| \geq \varepsilon s_n) = 0$$

holds, where $s_n^2 = \sum_k \sigma_k^2$, then the Central Limit Theorem in Thm. 1 holds.

Further, it is known that [6, p. 368]:

Lemma 4. *If each random variable $X_{tp_k}^s$ is uniformly bounded, and $\lim_{n \rightarrow \infty} s_n = \infty$, then the Lindeberg Condition in Lemma 3 holds.*

As our $score_Q$ function is bounded in $[0, 1]$ (defined in Sect. 2), every variable $X_{tp_k}^s$ is also bounded: $P(0 \leq X_{tp_k}^s \leq 1) = 1$. Further, the variance σ_k^2 can be expected to be > 0 for each $X_{tp_k}^s$, because $X_{tp_k}^s$ models ranking scores. That is, ranking scores are supposed to vary between different results, in order to assist an users in discovering results of interest. Therefore, it holds that $s_n = \sum_k \sigma_k^2 \rightarrow \infty$ with $n \rightarrow \infty$.

In simple terms, Thm. 1 states that the true distribution of $X_{Q^u(i)}^s$ converges (in the number of patterns in $Q^u(i)$) to a Gaussian distribution. Now, the question is: “how fast” does $X_{Q^u(i)}^s$ converge to a Gaussian distribution? For this convergence it holds:

Theorem 2 (Berry-Esseen Theorem). *Let $\rho_k = E(|X_{tp_k}^s|^3) < \infty$ be the third absolute normalized moment of $X_{tp_k}^s$. Then, it holds [13, p. 355]:*

$$\sup_x |F(x) - \phi(x)| \leq C \cdot \frac{\sum_k \rho_k}{(\sum_k \sigma_k^2)^{\frac{3}{2}}}$$

with $\phi(x)$ as standard Gaussian CDF, and $F(x)$ as exact CDF of $X_{Q^u(i)}^s$. Further, tp_k , μ_k , and σ_k^2 are defined as in Thm. 1.

C is a constant in Thm. 2, and is currently estimated as $0.4097 \leq C \leq 0.5600$ [13, p. 355]. Thus, intuitively, Thm. 2 gives an absolute bound on how close the true distribution of $X_{Q^u(i)}^s$ is to a Gaussian.

Approximation Error. Let X_i^e denote a random variable for the error introduced by pruning from input i . This error may be measured as the number of pruned partial bindings from i , which would have contributed to the final top- k result. Then, it holds that:

Lemma 5. *Random variable X_i^e follows a binomial distribution such that: $X_i^e \sim \text{bin}(c_i, \varepsilon + \tau)$. c_i stands for the # bindings pulled from input i , in order to produce the top- k results. Further, τ is the error threshold from Eq. 6, and ε is a small constant ≥ 0 .*

Proof Sketch. From a given input i we pull c_i partial bindings. Every of these bindings could be pruned “wrongfully” by the top- k test in Eq. 6, either because the binding probability was falsely estimated as 0, or because the score probability was smaller than the threshold τ . Let the probability for the former be bounded by a constant $\varepsilon \geq 0$, while the probability for the latter is known to be $\leq \tau$. Thus, the overall probability for a partial bindings to be wrongfully pruned in Eq. 6 is $\leq \varepsilon + \tau$. Further, pulling c_i partial bindings from input i may be conceived as c_i trials, where a wrongly pruned binding is a “hit”. Therefore, we can model X_i^e by means of binomial distribution, with c_i as number of trials, and $\varepsilon + \tau$ as “success” probability. Formally, $X_i^e \sim \text{bin}(c_i, \varepsilon + \tau)$, cf. Lemma 5.

Note, ε is a small error introduced by the binding indicator function $\mathbb{1}(Q^u(b) | b)$, see also Eq. 3. This error depends on the accuracy of the selectivity estimation, however, as the binding indicator only requires a binary decision, its induced error is frequently very small. In fact, our simplistic implementation in Eq. 4 is exact, i.e., leads to $\varepsilon = 0$.

Given a tree of one or more A-PBRJ operators having a total of n inputs, let X^e capture the overall error, i.e., number of false positives/negatives in the top- k results. We can show that X^e also adheres to a binomial distribution:

Theorem 3. *$X^e \sim \text{bin}(\sum_{i=1}^n c_i, \varepsilon + \tau)$, with input depth c_i , threshold τ , and ε , as defined in Lemma 5.*

Proof Sketch. Given a tree of joins having n inputs: $\{i_1, \dots, i_n\}$. Let every input i_j pull c_j partial bindings, in order for the join tree to produce the desired k top-ranked results. Further, the error (# wrongly pruned partial bindings) for each input i_j is modeled via variable $X_i^e \sim \text{bin}(c_i, \varepsilon + \tau)$, cf. Lemma 5. False positives/negatives results comprised in the final top- k bindings are caused by wrongly pruned partial bindings. More precisely, for a given input i_j , every wrongfully pruned partial binding could lead to a false positive/negative top- k result. Thus, errors from the individual inputs “sum up” to a total error – captured by X^e . In other words, random variable X^e is a summation over the random variables X_i^e . Further, errors made in the inputs are independent of each other, i.e., every pair of variables, X_i^e and X_j^e , is independent: $X_i^e \perp X_j^e$. Thus, X^e is again a binomial distribution with $\sum_{i=1}^n c_i$ trials and “success” probability $\varepsilon + \tau$: $X^e \sim \text{bin}(\sum_{i=1}^n c_i, \varepsilon + \tau)$ ■

Exploiting Thm. 3, we can give an expected error as a function of threshold τ : $E(X^e) = \sum_{i=1}^n c_i \cdot (\varepsilon + \tau)$.

4. EVALUATION

We conducted experiments for (1) *analyzing the efficiency and effectiveness of the A-PBRJ operator* in Sect. 4.2, and (2) *inspecting the behavior of our probabilistic component* in Sect. 4.3. By means of the former, we illustrate the overall performance of our approach, when compared with the exact PBRJ. The latter provides insights into overhead and accuracy of the probabilistic component.

4.1 Evaluation Setting

Data. We used two RDF benchmarks based on synthetic/real-world data: (1) The SP² benchmark features synthetic DBLP data comprising information about computer science bibliographies [33]. (2) The DBpedia SPARQL benchmark (DBPSB), which holds real-world data extracted from the DBpedia knowledge base [29]. For both benchmarks we used datasets of 10M triples each.

Queries. We employed queries from the SP² benchmark [33] and the DBPSB benchmark [29]. As both benchmarks comprise SPARQL queries, we translated the queries to our query model (basic graph patterns). Queries featuring no basic graph patterns were discarded, i.e., we omitted 12 and 4 queries in DBPSB and SP², respectively. This lead to 13 SP² queries and 120 queries for DBPSB – a comprehensive load of 133 queries in total.

The DBPSB queries were generated from 8 “seed queries” as proposed in [29]: each seed query comprises a special variable, which is randomly assigned a constant. Note, each constant is chosen such that the resulting query has a non-empty result. With every instantiation of that variable, a new benchmark query is given. Overall, we instantiated each of the 8 seed queries with 15 random bindings, resulting in a total of 120 queries. Queries varied in # pattern $\in [2, 9]$ and result sizes $\in [1, 5390436]$. Additional query statistics are in Table 1, and a complete query listing is given in the appendix, Sect. 8.

| | SP ² Queries | DBPSB Queries |
|-----------------------|-------------------------|---------------|
| # Queries | 13 | 120 |
| # Triple pattern | [2, 9] | [2, 4] |
| Mean(#Triple pattern) | 5 | 2.8 |
| Var(#Triple pattern) | 6.4 | 0.6 |
| # Results | [1, 5.4M] | [1, 50] |
| Mean(#Results) | 590K | 3.9 |
| Var(#Results) | 2.1B | 51.6 |

Table 1: Query statistics for the SP² and DBPSB benchmark.

Systems. We randomly generated bushy query plans. For a given query, all systems rely on the same plan. We implemented three

systems, solely differing in their join operator: (1) A system with *join-sort* operator, JS, which does not employ top- k processing, but instead produces all results, and then sorts them. (2) An *exact* rank join operator, PBRJ, featuring the corner-bound (Def. 2), and the corner-bound-adaptive pulling strategy in Def. 3. (3) Last, we implemented our *approximate* operator, A-PBRJ, see Sect. 3. All systems have indexes for random and sorted access.

Aiming at Web search ranking, we assume *no offline score information* for the A-PBRJ operator. Thus, we employ a Gaussian conjugate prior (with unknown mean and variance) for score probability learning. Training and top- k test implementation follows Algo. 2 and 3, see Sect. 3.3.2. Further, we used sufficient statistics based on a uniform distribution over $[0, 1]$, as discussed in Exp. 11 for sorted access sa_3 . Prior weights ν_0 and η_0 are both 1. Weight w in Algo. 2 is the sample size, $|x|$. Scores for single triple pattern bindings are random (see below), and complete binding scores are computed as summation. For binding probabilities, we aimed at a simple implementation as presented in Eq. 4, Sect. 3.2. The necessary SP, PO, and SO indexes were loaded into memory.

We implemented all systems in Java 6. Experiments were run on a Linux server with two Intel Xeon 5140 CPUs at 2.33GHz, 48GB memory (16GB assigned to the JVM), and a RAID10 with IBM SAS 148GB 10K rpm disks. Before each query execution, all operating system caches were cleared. The presented values are averages collected over five runs.

Hypothesis (H1): We expect that JS is outperformed by PBRJ, as it computes all results for a query. Further, we expect A-PBRJ to outperform JS and PBRJ, by trading effectiveness for efficiency.

Parameters. We used parameters as follows: We vary the number of results to be computed: $k \in \{1, 5, 10, 20\}$.

Hypothesis (H2): We predict efficiency to decrease in parameter k for A-PBRJ and PBRJ. Effectiveness should not be affected.

We chose triple pattern binding scores, $score_Q(t)$, at random with distribution $d \in \{u, n, e\}$ (uniform, normal, and exponential distribution). By means of varying distributions, we aim at an abstraction from a particular ranking function, and examine performance for different “classes” of functions. We employed standard parameters for all distributions, and normalized scores to be in $[0, 1]$.

Hypothesis (H3): A-PBRJ’s efficiency and effectiveness is not influenced by the choice of the ranking function and its score distribution, respectively.

We used top- k test thresholds $\tau \in [0, 0.8]$ for inspecting the trade-off between computation efficiency and effectiveness.

Hypothesis (H4): We expect efficiency of A-PBRJ to be increasing in τ , while its effectiveness will be decreasing in τ .

Metrics. We measure efficiency via: (1) Number of sorted and random accesses. (2) Time spend for top- k result computation. (3) Max. memory needed for buffering intermediate results.

As effectiveness metrics we use: (1) Precision: fraction of approximated top- k results being exact top- k results. (2) Recall: fraction of exact top- k results being reported as approximate results. Notice, precision and recall have identical values, as both share the same denominator k . We therefore discuss only precision results in the following. Further, precision is given as average over our query load, the so-called macro-precision. (3) Rank distance: approximate vs. exact top- k rank: $\frac{1}{k} \sum_{i=1, \dots, k} |rank^*(b) - rank(b)|$, with $rank^*(b)$ and $rank(b)$ as approximated and exact rank for binding b [19]. (4) Score error: approximate vs. exact top- k score: $\frac{1}{k} \sum_{b=1, \dots, k} |score_Q^*(b) - score_Q(b)|$, with $score_Q^*(b)$ and $score_Q(b)$ as approximated and exact score for binding b [38].

4.2 Evaluation: A-PBRJ

First, we the performance of the A-PBRJ system in terms of its

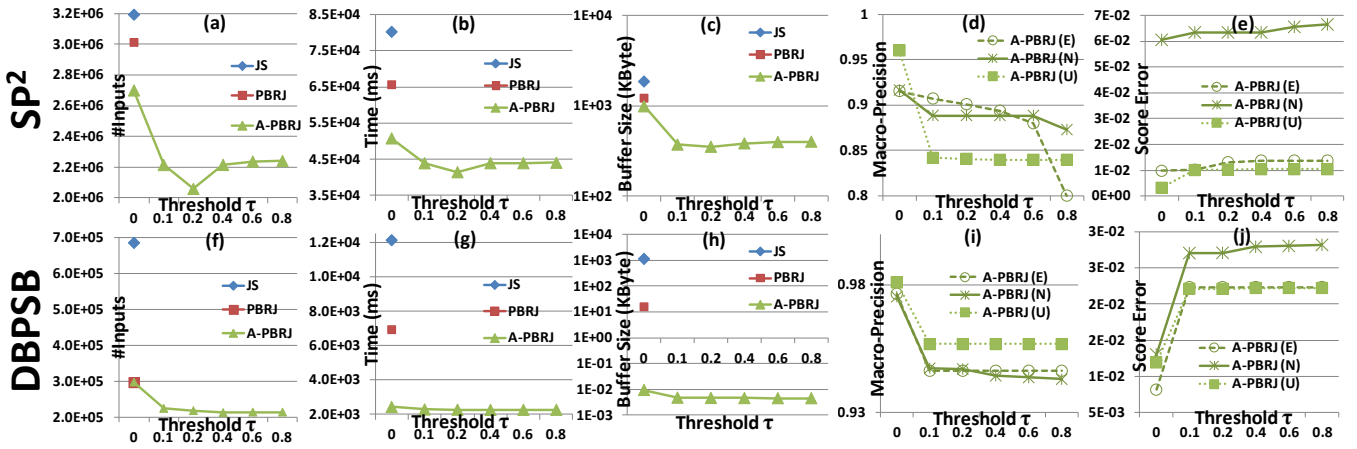


Figure 4: Average results over all queries from SP² and DBPSB: (a) + (f) Efficiency: #inputs vs. threshold τ . (b) + (g) Efficiency: time vs. threshold τ . (c) + (h) Efficiency: buffer sizes vs. threshold τ . (d) + (i) Effectiveness: macro-precision vs. threshold τ , given varying score distributions. (e) + (j) Effectiveness: score error vs. threshold τ , given varying score distributions.

effectiveness and efficiency.

Efficiency: Overall Results. Efficiency results are depicted in Fig. 4 (a)+(f), (b)+(g), and (c)+(h), as functions in threshold τ for SP² and DBPSB. As expected in hypothesis H1, we observed #inputs and computation time to be decreasing in τ , cf. Fig. 4 (a)+(f) and (b)+(g). For SP² (DBPSB), A-PBRJ needed up to 25% (23%) less inputs w.r.t. baseline PBRJ, and 30% (67%) regarding JS. We explain these gains with pruning of partial bindings, thereby omitting “unnecessary” joins and join attempts. Thus, A-PBRJ produces overall fewer partial and complete bindings, as determined by pruning of the top- k test. Fewer #inputs translated into average time savings of 35% (65%) as compared to PBRJ, and 47% (80%) w.r.t. JS, given SP² (DBPSB). Further, we noticed smaller buffers, holding less intermediate results, to contribute to these time savings, due to more efficient probing and maintenance of their hash tables.

Interestingly, we saw an increase in #inputs for $\tau \in [0.2, 0.4]$ in SP², and $\tau \in [0.4, 0.8]$ in DBPSB, Fig. 4 (a)+(f). For instance, comparing $\tau = 0.2$ and $\tau = 0.4$ in SP², A-PBRJ read 8% more inputs. DBPSB was less effected – we noticed a marginal increase of 2% for $\tau = 0.4$ vs. $\tau = 0.6$. We explain the increase in both benchmarks with a too “aggressive” pruning, i.e., too many partial bindings have been wrongfully pruned, leading to more inputs being read in order to produce the desired k results. More precisely, we inspected SP² benchmark runs for $\tau = 0.2$ vs. $\tau = 0.4$, and noticed a large number of partial bindings throughout the tree were pruned “too early”, i.e., they would have led to a larger or even a complete binding. In both cases, β thresholds (Algo. 1, line 8), capturing the upper-bound for unseen join results, were not accurate any more, as they depended on seeing these pruned bindings. In fact, $\tau \in [0.6, 0.8]$ was even more aggressive, however, the ratio between pruned bindings and read inputs was high enough to compensate for the extra inputs, Fig. 4 (a). Note, such wrongfully pruned bindings not necessarily have an effect on result quality. That is, we frequently observed these bindings to not contribute to the exact top- k result. Overall, we saw a “sweet spot” at $\tau \approx 0.2$ for SP² and DBPSB – here, we noted pruning to be fairly accurate, i.e., only few partial bindings were wrongfully pruned. This is also reflected in high precision (recall) values for both benchmarks given $\tau \approx 0.2$: 88% (95%) in SP² (DBPSB) – as discussed in the following.

With regard to computation time for SP² and DBPSB queries, we noticed similar effects as for the #inputs, cf. Fig. 4 (b)+(g). In particular, the “sweet spot” at $\tau \approx 0.2$ is also reflected here. That is, for $\tau = 0.2$ A-PBRJ needs 40% (68%) less time than PBRJ. In compari-

son to JS, A-PBRJ is 51% (82%) faster in SP² (DBPSB). We explain such drastic time savings with the less #inputs being read, due partial bindings pruned by our top- k test. For instance, we were able to prune up to 40% (90%) of the inputs, given SP² (DBPSB). This lead to significantly less join attempts, and also smaller buffers (explained below). In fact, we noted that binding probabilities played a crucial role, as they oftentimes allowed reliable discovery (pruning) of partial bindings with little/no probability of leading to complete bindings. Thus, similar to a random access, binding probabilities allowed to “probe” unevaluated patterns.

Last, considering buffer sizes in SP² (DBPSB), we saw 60% (95%) less space consumption w.r.t. PBRJ, and 75% (99%) in contrast to JS, cf. Fig. 4 (c)+(h). Note, Fig. 4 (c)+(h) shows average buffers size per join in KByte, while employing a dictionary encoding for nodes in the data graph. Smaller buffers in A-PBRJ are possible due to fewer intermediate results and less inputs being loaded. Interestingly, as these buffers are based on hash tables, we observed smaller buffers to also have a positive effect on runtime behavior. That is, less maintenance was necessary, and they allowed for more efficient probing, due to less entries. So, buffer sizes also contributed to the above discussed time savings.

Efficiency: Sensitivity Analysis. Let us discuss efficiency sensitivity for different parameters – as expressed in hypotheses H2-H4.

As expressed by hypothesis H2, we observed #inputs and time to increase in k for A-PBRJ and PBRJ. For instance, comparing $k = 1$ and $k = 20$, A-PBRJ needed a factor of 1.2 (5.7) more time, given SP² (DBPSB). Similarly, 1.2 (6.8) times more inputs were consumed by A-PBRJ for SP² (DBPSB). Such a behavior is expected (H2), as more inputs/join attempts are required to produce a larger result. PBRJ leads to a similar performance decrease in SP². For instance, for $k = 1$ vs. $k = 20$ a factor of 1.3 (1.2) more inputs (time) are needed. Given DBPSB queries, an interesting observation is that PBRJ achieves a 50 – 60% smaller performance decrease than A-PBRJ for $k = 1$ vs. $k = 20$. We explain this with many DBPSB queries having a cardinality ≤ 10 , e.g., Q_1 where cardinality is 1. For such queries, A-PBRJ can not prune “more” bindings for a larger k , as the algorithm first needs to compute k complete results, which fails due to their small cardinality. Thus, A-PBRJ is more sensitive to k w.r.t. queries with cardinality $\leq k$.

Furthermore, we can confirm our hypothesis H3 with regard to system efficiency: *we did not find correlations between system performance and score distributions*. In other words, score distributions (ranking functions) had no impact on performance of the

A-PBRJ, given SP² or DBPSB. For SP², we noticed A-PBRJ to outperform PBRJ by 22% (42%) for e distribution, 21% (35%) given u distribution, and 8% (14%) for n distribution w.r.t. #inputs (time). For DBPSB queries, A-PBRJ resulted in different gains over PBRJ: 27% (65%) for e distribution, 23% (64%) given u distribution, and 21% (64%) for n distribution w.r.t. #inputs (time).

Last, with regard to parameter τ , we noted A-PBRJ efficiency to increase with $\tau \in [0, 2]$ for SP² (DBPSB) – partially confirming hypothesis H4. However, as outlined above, too aggressive pruning let to “inverse” effects. An important observation is, however, that our approach was already able to achieve performance gains with a very small $\tau < 0.1$. Here, partial bindings were pruned primarily due to their low binding probability. In fact, A-PBRJ could even save time for $\tau = 0$: 26% (60%) with SP² (DBPSB). We inspected queries leading to such saving, and saw that many of their partial bindings had a binding probability of 0. Thus, they were pruned even with $\tau = 0$. We argue that this a strong advantage of A-PBRJ: even for low error thresholds (leading to minor effectiveness decrease), we could achieve significant efficiency gains.

Effectiveness: Overall Results. Fig. 4 (d)+(i), and (e)+(j) depict effectiveness evaluation results for varying score distributions as average over all queries. More precisely, (d)+(i) give macro-precision, and (e)+(j) show score error – both as functions in τ . For space reasons, figures depicting the rank distance were omitted. Instead, we describe key observations in the text.

We observed high precision values $\in [0.8, 0.95]$ for both benchmarks and over all queries, cf. Fig. 4 (d)+(i). More precisely, we saw best results for a small $\tau < 0.1$, and the uniform distribution. However, given $\tau < 0.1$, all distributions led to very similar precision results $\in [0.9, 0.95]$ and $[0.95, 0.98]$ for SP² and DBPSB, respectively. This confirms our hypothesis H3: *A-PBRJ’s effectiveness is not affected by a particular score distribution*. We explain these good approximations with accurate score distributions (discussed in Sect. 4.3), and reliable binding probabilities. In fact, our simplistic implementation in Eq. 4 led to *exact* binding probabilities, thereby pruning partial bindings *with certainty*. This is because Eq. 4 checks necessary (but not sufficient) conditions, which a partial bindings has to fulfill, in order to contribute to a complete binding. This resulted in accurate pruning of many partial bindings – for some queries in DBPSB up to 97% of their total inputs. Such DBPSB queries featured highly selective patterns, and had only a small result cardinality ≤ 10 , thereby allowing for an highly effective pruning via binding probabilities.

In order to quantify “how bad” false positive/negative results are, we employed the score error and rank distance metric. Score error results depicted in Fig. (e)+(j). We observed that rank distance (score error) was $\in [0.01, 0.02]$ ($[0.07, 0.11]$) for a small $\tau < 0.1$, over all distributions and both benchmarks. We explain this we our high precision (recall), i.e., A-PBRJ led to only few false positive/negative top- k results given $\tau = 0$. Further, as expected in H4, both metrics increased in τ , due to more pruning and false positives/negatives. For instance, for $\tau = 0$ vs. $\tau = 0.8$ we noted an increase by a factor of 3.3 (2.6) for rank distance in SP² (DBPSB). Overall, rank distance and score error results were very promising: we saw an average score error of 0.03 (0.02) over all τ and queries, given SP² (DBPSB).

Effectiveness: Sensitivity Analysis. Concerning the sensitivity effectiveness w.r.t. parameter k , we can confirm the initial hypothesis H2: *k does not impact the A-PBRJ’s effectiveness*. Given SP², we saw A-PBRJ to be fairly stable in different values for parameter k . For instance, macro-precision was in $[0.87, 0.89]$. Other metrics showed similar, minor fluctuations. Also for DBPSB, we noted only minor effectiveness fluctuations, e.g., macro-precision varied

around 7% w.r.t. different k . We explain this good behavior with (1) quality of learned score distributions for SP² (discussed later), and (2) query characteristics of DBPSB. As for the latter, we noted that DBPSB mostly contained very selective queries having a cardinality ≤ 10 . Thus, for such queries the risk of pruning “the wrong” partial bindings did not increase in k . Note, SP² featured very different queries – many SP² queries had a large cardinality $\gg 100$. Thus, we argue the stable performance of A-PBRJ for SP² to be an indicator for a reliable top- k test and score/binding probabilities.

As predicted in hypothesis H3, we observed A-PBRJ to not be influenced by varying score distributions, Fig. 4 (d)+(i), and (e)+(j). Given SP², we saw a macro-precision (over all queries and values for τ) of 0.86 for u distribution, 0.88 for e distribution, and 0.89 for n distribution. Also for the DBPSB benchmark, we observed only minor changes in macro-precision: 0.96 for u distribution, 0.95 for e as well as n distribution. We explain this with a good score distribution quality (see Sect. 4.3), leading to reliable top- k tests – independent of the actual triple scores. Note, we made similar observations for score error and rank distance metric. One exception was, however, the score error for distribution n , Fig. 4 (e)+(j). This distribution led to an higher score error than, e.g., the u distribution by a factor of 6.9 (1.13) for SP² (DBPSB). We analyzed these outlier queries in SP² and DBPSB: scores from wrongly pruned top- k results were higher for n than for other distributions. This caused an increase in the score error.

With regard to the effectiveness of A-PBRJ vs. parameter τ , we noticed that metrics over both benchmarks decreased with increasing τ . For instance, macro-precision decreased for $\tau = 0$ vs. $\tau = 0.8$ with 11% (5%), given SP² (DBPSB). Such a behavior is expected (H4), as chances of pruning “the wrong” bindings increase with higher τ values. Thus, while leading to efficiency gains (discussed above), a higher value for τ causes effectiveness losses. Overall, this confirms hypotheses H4 and H1 that *A-PBRJ trades off effectiveness for efficiency, as dictated by threshold τ* .

4.3 Evaluation: Probabilistic Component

In this section, we analyze the performance of the probabilistic component in terms of its efficiency and effectiveness. More specifically, as binding probabilities are estimated via a given selectivity estimation framework based on previous work [31, 32], we focus on learning and computation of score probabilities.

| | SP ² | DBPSB | Dist. | SP ² | DBPSB |
|--------------|-----------------|---------|-------|-----------------|-------|
| Time (ms) | [1, 300] | [1, 24] | e | 0.34 | 0.04 |
| # Samples | [1, 4M] | [1, 50] | n | 0.34 | 0.01 |
| Avg. #Sample | 400K | 4 | u | 0.31 | 0.02 |

(a) Efficiency: average learning time and # learning samples.

(b) Effectiveness: average p -value from the goodness-of-fit test.

Table 2: Efficiency and effectiveness of score distribution learning.

Efficiency. First, we analyze the overhead introduced by score distribution learning. For this, we measured the time needed for hyperparameter training, cf. Algo. 2 in Sect. 3.2. We set the training threshold, i.e., the # new bindings after which a new training procedure is triggered, to 1 (Algo. 1, line 12). Table 2-a gives average training times and # samples for distribution training.

We observed average learning times $\in [1, 300]$ ms ($[1, 24]$ ms) over all score distributions, queries and thresholds τ , given the SP² (DBPSB) benchmark. We noted the driving factor to be the overall query selectivity. That is, SP² queries often had a large cardinality $\gg 100$, which led to a high number (up to 4,227,732) of training samples. In contrast, DBPSB featured highly selective queries, resulting in few training iterations – only up to 50 score

samples were available on average. Overall, this explains the additional training time (factor 12.5) needed for SP² queries, when compared to DBPSB. Note, one may easily cope with high cardinality queries by: (1) setting a larger training threshold, or (2) stop distribution learning, if distribution “quality” does not improve any more. However, such optimizations are left to future work.

Second, we measured the extra time required for performing a top- k test, see Algo. 3 in Sect. 3.2. On average over both benchmarks and all parameters, a top- k test needed 4.3K ns. This time comprises a selectivity estimation lookup for the binding probability, and the score probability computation. In contrast, a sorted (random) access took 26.8K ns (1.7M ns) on average. Thus, a top- k test is fairly cheap in comparison to a sorted/random access.

Effectiveness. For judging learning effectiveness, we captured how well the trained distributions “fits” the observed complete binding scores. More precisely, we applied the well-known Kolmogorov-Smirnov test [13], which measures, via a p -value in $[0, 1]$, whether a sample comes from the population of a specific distribution. Table 2-b shows p -values as averages over both benchmarks.

We observed drastic differences between p -values for SP² and DBPSB. SP² results were very promising, and reflect that learned distributions accurately capture the true scores of complete bindings. *This confirms hypothesis H3, as our distribution learning could achieve high-quality approximations: a p -value of 0.34 for distribution e as well as n , and 0.31 for u distribution.* With regard to DBPSB, we could not train good distributions, i.e., we measured poor p -values: 0.04 for distribution e , 0.01 given n , and 0.02 for u distribution. We explain this with the few training samples available in DBPSB queries, due to their high selectivity. As discussed above, SP² queries featured much more score samples, cf. Table 2-a. However, the interesting observation is that the overall approach, A-PBRJ, was hardly affected. In fact, A-PBRJ achieved a very high precision $\in [0.95, 0.98]$ for DBPSB. This is because the score probabilities are only relevant for the top- k test, if k complete bindings have been computed (see discussion in Sect. 3.2). However, for many DBPSB queries, their cardinality was $\leq k$. Thus, the low quality distributions had little to no effect.

Overall, we can conclude that the probabilistic component trains score distributions in an efficient and effective manner, if sufficient score samples are available.

5. RELATED WORK

Early work on rank-aware query processing aimed at the selection top- k problem [10, 11]. Here, the goal is to find k top-ranked results, where each result is an entity with n attributes, that is ranked according to m criteria, defined on those attributes [18].

To foster efficient result computation, *approximate selection top- k* techniques have been proposed [2, 3, 28, 35, 38]. [38] used score statistics to predict the highest possible complete score of a partial binding. Partial results are discarded, if they are not likely to contribute to a top- k result. Focusing on distributed top- k queries, [28] employed histograms to predict aggregated score values over a space of data sources. Anytime measures for selection top- k have been introduced by [2, 3]. For this, the authors used offline score information, e.g., histograms, to predict complete binding scores at runtime. In [35], a framework for approximate top- k processing under budgetary constraints, as well as algorithms for scheduling sorted and random access have been presented.

In contrast, we target the join top- k problem. In this setting, scores are assigned to single triples, and a complete result is obtained by joining these triples. The score of a complete binding is an aggregation of the scores for its comprised triples [18]. A large body of work aimed at an *exact join top- k* processing, e.g., [12, 17,

20, 21, 27, 30, 34]. In particular, a rank-join based on the A* algorithm, J*, was proposed in [30]. [17] introduced the hash rank-join algorithm (HRJN), which was further addressed in [21]. [20] investigated the join top- k problem given a non-monotonic aggregation function. In [34], the authors developed an algorithm template, the Pull/Bound Rank Join (PBRJ), which covers previous work on rank join approaches. Using the PBRJ, the authors introduced a novel bound for the case that a join features more than two inputs, and a triple has multiple scores. Extending [34], efficiency and costs aspects were further discussed in [12]. A cost-aware scheduling strategy for random and sorted accesses in rank-joins was presented in [27]. Recently, two works adapted top- k join processing to RDF data and SPARQL queries [25, 39].

As result accuracy is not of high importance for Web search, we aim at an *approximate join top- k* processing. However, to the best of our knowledge, there is no work addressing this problem. Further, approximate selection top- k techniques are not directly applicable. This is because the selection top- k problem solely ranks “single” entities, and does not consider joins. In a join top- k setting, multiple triples are combined via joins. This leads to many partial bindings that never contribute to a top- k binding, because they do not satisfy the join conditions. Thus, an approximate join top- k should not only capture a score probability, but also a probability for satisfying the query constraints. However, approximate selection top- k strategies solely judge the likelihood of a partial binding leading to a top- k result by means of score probabilities [2, 3, 28, 35, 38].

Moreover, existing approximate top- k approaches heavily rely on “offline ranking”. That is, scores must be known before runtime for computing statistics, e.g., histograms [2, 3, 28, 35, 38], or suitable/approximated score distributions [38]. However, we target a Web search context, where systems oftentimes exploit user-/query-dependent ranking functions. Given such a ranking, scores are only known at runtime. In fact, our approach not only supports such an “online ranking”, but allows a flexible integration of any kind of ranking function and offline available score information.

6. CONCLUSION

In this paper, we introduced an approximate join top- k algorithm, A-PBRJ, well-suited for the Web of data. For this, we extended the well-known PBRJ framework with a novel probabilistic component. This component allows us to estimate the probability of a partial query binding (1) to lead to a complete binding, and (2) to have a score higher than the smallest currently known top- k score. We conducted a theoretical analyses showing that our approach features an efficient and effective instantiation of the probabilistic component. Furthermore, we implemented and evaluated the A-PBRJ system by means of two state-of-the-art Web data benchmarks. Our results are promising, as we could achieve times savings of up to 65% over the baselines, while maintaining a high precision/recall.

7. REFERENCES

- [1] R. Agrawal, R. Rantzaou, and E. Terzi. Context-sensitive ranking. In *SIGMOD*, 2006.
- [2] B. Arai, G. Das, D. Gunopulos, and N. Koudas. Anytime measures for top- k algorithms. In *VLDB*, 2007.
- [3] B. Arai, G. Das, D. Gunopulos, and N. Koudas. Anytime measures for top- k algorithms on exact and fuzzy data sets. *VLDB Journal*, 18:407–427, 2009.
- [4] A. Baid, I. Rae, J. Li, A. Doan, and J. Naughton. Toward scalable keyword search over relational data. *VLDB*, 2010.

[5] J. M. Bernardo and A. F. M. Smith. *Bayesian Theory*. John Wiley & Sons, 1994.

[6] P. Billingsley. *Probability and Measure*. Wiley-Interscience, 3 edition, 1995.

[7] S. Chaudhuri, G. Das, V. Hristidis, and G. Weikum. Probabilistic information retrieval approach for ranking of database query results. *TODS*, 31(3):1134–1168, 2006.

[8] Y.-Y. Chen, T. Suel, and A. Markowetz. Efficient query processing in geographic web search engines. In *SIGMOD*, 2006.

[9] G. Cong, C. S. Jensen, and D. Wu. Efficient retrieval of the top-k most relevant spatial web objects. *VLDB*, 2009.

[10] R. Fagin. Combining fuzzy information from multiple systems. *JCSS*, 58:83–99, 1999.

[11] R. Fagin, A. Lotem, and M. Naor. Optimal aggregation algorithms for middleware. In *SIGMOD*, 2001.

[12] J. Finger and N. Polyzotis. Robust and efficient algorithms for rank join evaluation. In *SIGMOD*, 2009.

[13] A. Gut. *Probability: A Graduate Course*. Springer, 2012.

[14] A. Guttman. R-trees: a dynamic index structure for spatial searching. *SIGMOD Rec.*, 14(2):47–57, 1984.

[15] P. D. Hoff. *A First Course in Bayesian Statistical Methods*. Springer, 2009.

[16] K. Hose, R. Schenkel, M. Theobald, and G. Weikum. Database foundations for scalable RDF processing. In *RR*, 2011.

[17] I. F. Ilyas, W. G. Aref, and A. K. Elmagarmid. Supporting top-k join queries in relational databases. *VLDB Journal*, 13:207–221, 2004.

[18] I. F. Ilyas, G. Beskales, and M. A. Soliman. A survey of top-k query processing techniques in relational database systems. *ACM Comput. Surv.*, 40(4):11:1–11:58, 2008.

[19] M. Kendall and J. D. Gibbons. *Rank Correlation Methods*. Oxford University Press, 1990.

[20] B. Kimelfeld and Y. Sagiv. Incrementally Computing Ordered Answers of Acyclic Conjunctive Queries. In *NGITS*, volume 4032 of *Lecture Notes in Computer Science*, pages 141–152. 2006.

[21] C. Li, K. C.-C. Chang, I. F. Ilyas, and S. Song. RankSQL: query algebra and optimization for relational top-k queries. In *SIGMOD*, 2005.

[22] F. Li, K. Yi, and W. Le. Top-k queries on temporal data. *The VLDB Journal*, 19(5):715–733, 2010.

[23] Y. Luo, X. Lin, W. Wang, and X. Zhou. Spark: top-k keyword query in relational databases. In *SIGMOD*, 2007.

[24] A. Maduko, K. Anyanwu, A. Sheth, and P. Schliekelman. Graph summaries for subgraph frequency estimation. In *ESWC*, 2008.

[25] S. Magliacane, A. Bozzon, and E. Della Valle. Efficient execution of top-k SPARQL queries. In *ISWC*, 2012.

[26] N. Mamoulis, M. L. Yiu, K. H. Cheng, and D. W. Cheung. Efficient top-k aggregation of ranked inputs. *TODS*, 32, 2007.

[27] D. Martinenghi and M. Tagliasacchi. Cost-Aware Rank Join with Random and Sorted Access. *TKDE*, 24(12):2143–2155, 2012.

[28] S. Michel, P. Triantafyllou, and G. Weikum. KLEE: a framework for distributed top-k query algorithms. In *VLDB*, 2005.

[29] M. Morsey, J. Lehmann, S. Auer, and A.-C. Ngonga Ngomo. DBpedia SPARQL Benchmark – Performance Assessment with Real Queries on Real Data. In *ISWC*, 2011.

[30] A. Natsev, Y.-C. Chang, J. R. Smith, C.-S. Li, and J. S. Vitter. Supporting Incremental Join Queries on Ranked Inputs. In *VLDB*, 2001.

[31] T. Neumann and G. Moerkotte. Characteristic sets: Accurate cardinality estimation for RDF queries with multiple joins. In *ICDE*, 2011.

[32] T. Neumann and G. Weikum. Scalable join processing on very large RDF graphs. In *SIGMOD*, 2009.

[33] M. Schmidt, T. Hornung, G. Lausen, and C. Pinkel. SP²Bench: A SPARQL Performance Benchmark. In *ICDE*, 2009.

[34] K. Schnaitter and N. Polyzotis. Evaluating rank joins with optimal cost. In *PODS*, 2008.

[35] M. Shmueli-Scheuer, C. Li, Y. Mass, H. Roitman, R. Schenkel, and G. Weikum. Best-Effort Top-k Query Processing Under Budgetary Constraints. In *ICDE*, 2009.

[36] M. Stocker, A. Seaborne, A. Bernstein, C. Kiefer, and D. Reynolds. SPARQL basic graph pattern optimization using selectivity estimation. In *WWW*, 2008.

[37] A. Telang, C. Li, and S. Chakravarthy. One Size Does Not Fit All: Toward User- and Query-Dependent Ranking for Web Databases. *TKDE*, 24(9):1671–1685, 2012.

[38] M. Theobald, G. Weikum, and R. Schenkel. Top-k query evaluation with probabilistic guarantees. In *VLDB*, 2004.

[39] A. Wagner, T. T. Duc, G. Ladwig, A. Harth, and R. Studer. Top-k linked data query processing. In *ESWC*, 2012.

8. APPENDIX

In this section, we present the query load that was used during our experiments. Queries for the SP² benchmark are based on [33], while the DBPSB benchmark queries are generated from seed queries in [29]. All queries are given in RDF N3³ notation.

Listing 1: Prefixes used for SP² and DBPSB queries.

```

1 @prefix rdf:
2   <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
3 @prefix rdfs:
4   <http://www.w3.org/2000/01/rdf-schema#> .
5 @prefix dc:
6   <http://purl.org/dc/elements/1.1/> .
7 @prefix dcterms:
8   <http://purl.org/dc/terms/> .
9 @prefix xs:
10  <http://www.w3.org/2001/XMLSchema#> .
11 @prefix bench:
12  <http://localhost/vocabulary/bench/> .
13 @prefix foaf:
14  <http://xmlns.com/foaf/0.1/> .
15 @prefix swrc:
16  <http://swrc.ontoware.org/ontology#> .
17 @prefix dbpedia:
18  <http://dbpedia.org/ontology/> .
19 @prefix dbpedia:prop:
20  <http://dbpedia.org/property/> .
21 @prefix dbpedia:res:
22  <http://dbpedia.org/resource/> .
23 @prefix skos:
24  <http://www.w3.org/2004/02/skos/core> .
25 @prefix yago:
26  <http://dbpedia.org/class/yago/> .

```

Listing 2: Queries for SP² benchmark [33].

³<http://www.w3.org/TeamSubmission/n3/>

```

1  ### 1
2  ?journal dc:title "Journal 1 (1940)"^^xs:string .
3  ?journal dcterms:issued ?yr .
4  ?journal rdf:type bench:Journal .
5
6  ### 2
7  ?inproc dcterms:partOf ?proc .
8  ?inproc bench:booktitle ?booktitle .
9  ?inproc swrc:pages ?page .
10 ?inproc dc:title ?title .
11 ?inproc rdfs:seeAlso ?ee .
12 ?inproc foaf:homepage ?url .
13 ?inproc dcterms:issued ?yr .
14 ?inproc dc:creator ?author .
15 ?inproc rdf:type bench:Inproceedings .
16
17 ### 3
18 ?article rdf:type bench:Article .
19 ?article swrc:pages ?value .
20
21 ### 4
22 ?article rdf:type bench:Article .
23 ?article swrc:month ?value .
24
25 ### 5
26 ?article rdf:type bench:Article .
27 ?article swrc:isbn ?value .
28
29 ### 6
30 ?article1 dc:creator ?author1 .
31 ?author1 foaf:name ?name1 .
32 ?article1 swrc:journal ?journal .
33 ?article2 swrc:journal ?journal .
34 ?article2 dc:creator ?author2 .
35 ?author2 foaf:name ?name2 .
36 ?article1 rdf:type bench:Article .
37 ?article2 rdf:type bench:Article .
38
39 ### 7
40 ?article dc:creator ?person .
41 ?person foaf:name ?name .
42 ?person2 foaf:name ?name .
43 ?inproc dc:creator ?person2 .
44 ?article rdf:type bench:Article .
45 ?inproc rdf:type bench:Inproceedings .
46
47 ### 8
48 ?document dcterms:issued ?yr .
49 ?document dc:creator ?author .
50 ?author foaf:name ?name .
51 ?document rdf:type ?class .
52 ?class rdfs:subClassOf foaf:Document .
53
54 ### 9
55 ?doc dc:title ?title .
56 ?bag2 ?member2 ?doc .
57 ?doc2 dcterms:references ?bag2 .
58 ?doc rdf:type ?class .
59 ?class rdfs:subClassOf foaf:Document .
60
61 ### 10
62 ?erdoes foaf:name "Paul Erdoes"^^xs:string .
63 ?document dc:creator ?erdoes .
64 ?document dc:creator ?author .
65 ?document2 dc:creator ?author .
66 ?document2 dc:creator ?author2 .
67 ?author2 foaf:name ?name .
68 ?erdoes rdf:type foaf:Person .
69
70 ### 11
71 ?person rdf:type foaf:Person .
72 ?subject ?predicate ?person .
73

```

```

74 ### 12
75 ?article dc:creator ?person1 .
76 ?person1 foaf:name ?name .
77 ?person2 foaf:name ?name .
78 ?inproc dc:creator ?person2 .
79 ?inproc rdf:type bench:Inproceedings .
80 ?article rdf:type bench:Article .
81
82 ### 13
83 ?erdoes foaf:name "Paul Erdoes"^^xs:string .
84 ?document dc:creator ?erdoes .
85 ?document dc:creator ?author .
86 ?document2 dc:creator ?author .
87 ?document2 dc:creator ?author2 .
88 ?author2 foaf:name ?name .
89 ?erdoes rdf:type foaf:Person .

```

Listing 3: Queries for DBPSB benchmark [29].

```

1  ### 1
2  ?var5 rdf:type dbpedia:Person .
3  ?var5 foaf:page ?var8 .
4  ?var5 dbpedia:thumbnail ?var4 .
5  ?var5 rdfs:label "Thaksin Shinawatra"@nn .
6
7  ### 2
8  ?var5 dbpedia:thumbnail ?var4 .
9  ?var5 rdf:type dbpedia:Person .
10 ?var5 rdfs:label
11     "\u0420\u0438\u0448\u0435\u0435,
12     \u0428\u0438\u043D\u0430\u0432\u0430\u0442\u0440\u0430"@ru .
13 ?var5 foaf:page ?var8 .
14
15 ### 3
16 ?var5 dbpedia:thumbnail ?var4 .
17 ?var5 rdf:type dbpedia:Person .
18 ?var5 rdfs:label
19     "Amadeo, quinto
20     Duque de Aosta"@es .
21 ?var5 foaf:page ?var8 .
22
23 ### 4
24 ?var5 dbpedia:thumbnail ?var4 .
25 ?var5 rdf:type dbpedia:Person .
26 ?var5 rdfs:label "Godeberta"@en .
27 ?var5 foaf:page ?var8 .
28
29 ### 5
30 ?var5 dbpedia:thumbnail ?var4 .
31 ?var5 rdf:type dbpedia:Person .
32 ?var5 rdfs:label "Thaksin Shinawatra"@nl .
33 ?var5 foaf:page ?var8 .
34
35 ### 6
36 ?var5 dbpedia:thumbnail ?var4 .
37 ?var5 rdf:type dbpedia:Person .
38 ?var5 rdfs:label
39     "\u827E\u9A30\u00B7
40     \u4F0A\u683C\u8A00"@zh .
41 ?var5 foaf:page ?var8 .
42
43 ### 7
44 ?var5 dbpedia:thumbnail ?var4 .
45 ?var5 rdf:type dbpedia:Person .
46 ?var5 rdfs:label "Vlad\u00EDmir Karpets"@es .
47 ?var5 foaf:page ?var8 .
48
49 ### 8
50 ?var5 dbpedia:thumbnail ?var4 .
51 ?var5 rdf:type dbpedia:Person .
52 ?var5 rdfs:label "Daniel Pearl"@en .
53 ?var5 foaf:page ?var8 .
54

```



```

55 ### 9
56 ?var5 dbpedia:thumbnail ?var4 .
57 ?var5 rdf:type dbpedia:Person .
58 ?var5 rdfs:label
59     "\u030DE\u030EA\u030FC\u030FB\u030EB
60     \u030A4\u030FC\u030BA\u030FB\u030C9\u030EB
61     \u030EC\u030A2\u030F3"@ja .
62 ?var5 foaf:page ?var8 .
63
64 ### 10
65 ?var5 dbpedia:thumbnail ?var4 .
66 ?var5 rdf:type dbpedia:Person .
67 ?var5 rdfs:label "Walter Hodge"@en .
68 ?var5 foaf:page ?var8 .
69
70 ### 11
71 ?var5 dbpedia:thumbnail ?var4 .
72 ?var5 rdf:type dbpedia:Person .
73 ?var5 rdfs:label "Damian Wayne"@en .
74 ?var5 foaf:page ?var8 .
75
76 ### 12
77 ?var5 dbpedia:thumbnail ?var4 .
78 ?var5 rdf:type dbpedia:Person .
79 ?var5 rdfs:label
80     "\u0417\u0430\u0431
81     \u0435\u0432\u0441\u0441\u0430
82     \u0438\u0439, \u0410\u0430\u0430\u0437
83     \u0438\u0438\u0438\u0438\u0438"@ru .
84 ?var5 foaf:page ?var8 .
85
86 ### 13
87 ?var5 dbpedia:thumbnail ?var4 .
88 ?var5 rdf:type dbpedia:Person .
89 ?var5 rdfs:label
90     "\u0413\u0438
91     \u0438\u0430\u0443\u0443\u0440\u0438,
92     \u0418\u0438\u0438\u0438\u0438\u0438\u0438
93     \u0438\u0438\u0438 \u0417\u0440\u0443\u0443\u0440
94     \u0430\u0438\u0438\u0438\u0438\u0438\u0438
95     \u0447"@ru .
96 ?var5 foaf:page ?var8 .
97
98 ### 14
99 ?var5 dbpedia:thumbnail ?var4 .
100 ?var5 rdf:type dbpedia:Person .
101 ?var5 rdfs:label "Francis Atterbury"@en .
102 ?var5 foaf:page ?var8 .
103
104 ### 15
105 ?var5 dbpedia:thumbnail ?var4 .
106 ?var5 rdf:type dbpedia:Person .
107 ?var5 rdfs:label "Damian Wayne"@es .
108 ?var5 foaf:page ?var8 .
109
110 ### 16
111 ?var4 dbpedia:prop:birthPlace
112     "Vigny, Val d'Oise"@en .
113 ?var4 dbpedia:birthDate ?var6 .
114 ?var4 foaf:name ?var8 .
115 ?var4 dbpedia:deathDate ?var10 .
116
117 ### 17
118 ?var4 dbpedia:prop:birthPlace
119     "Salisbury, England"@en .
120 ?var4 dbpedia:birthDate ?var6 .
121 ?var4 foaf:name ?var8 .
122 ?var4 dbpedia:deathDate ?var10 .
123
124 ### 18
125 ?var4 dbpedia:prop:birthPlace
126     "Bailey in the city of Durham"@en .
127 ?var4 dbpedia:birthDate ?var6 .

```

```

128 ?var4 foaf:name ?var8 .
129 ?var4 dbpedia:deathDate ?var10 .
130
131 ### 19
132 ?var4 dbpedia:prop:birthPlace
133     "Vasilievskaya,
134     Tambov Governorate"@en .
135 ?var4 dbpedia:birthDate ?var6 .
136 ?var4 foaf:name ?var8 .
137 ?var4 dbpedia:deathDate ?var10 .
138
139 ### 20
140 ?var4 dbpedia:prop:birthPlace
141     dbpedia:Waltham%2C_Massachusetts .
142 ?var4 dbpedia:birthDate ?var6 .
143 ?var4 foaf:name ?var8 .
144 ?var4 dbpedia:deathDate ?var10 .
145
146 ### 21
147 ?var4 dbpedia:prop:birthPlace
148     dbpedia:Valencia%2C_Spain .
149 ?var4 dbpedia:birthDate ?var6 .
150 ?var4 foaf:name ?var8 .
151 ?var4 dbpedia:deathDate ?var10 .
152
153 ### 22
154 ?var4 dbpedia:prop:birthPlace
155     dbpedia:Halifax%2C_West_Yorkshire .
156 ?var4 dbpedia:birthDate ?var6 .
157 ?var4 foaf:name ?var8 .
158 ?var4 dbpedia:deathDate ?var10 .
159
160 ### 23
161 ?var4 dbpedia:prop:birthPlace dbpedia:Sucre .
162 ?var4 dbpedia:birthDate ?var6 .
163 ?var4 foaf:name ?var8 .
164 ?var4 dbpedia:deathDate ?var10 .
165
166 ### 24
167 ?var4 dbpedia:prop:birthPlace
168     dbpedia:L%C3%BACar .
169 ?var4 dbpedia:birthDate ?var6 .
170 ?var4 foaf:name ?var8 .
171 ?var4 dbpedia:deathDate ?var10 .
172
173 ### 25
174 ?var4 dbpedia:prop:birthPlace
175     dbpedia:%C3%89tampes .
176 ?var4 dbpedia:birthDate ?var6 .
177 ?var4 foaf:name ?var8 .
178 ?var4 dbpedia:deathDate ?var10 .
179
180 ### 26
181 ?var4 dbpedia:prop:birthPlace
182     dbpedia:Montgomery_County%2C_Maryland
183     .
184 ?var4 dbpedia:birthDate ?var6 .
185 ?var4 foaf:name ?var8 .
186 ?var4 dbpedia:deathDate ?var10 .
187
188 ### 27
189 ?var4 dbpedia:prop:birthPlace
190     "Berkeley, Gloucestershire"@en .
191 ?var4 dbpedia:birthDate ?var6 .
192 ?var4 foaf:name ?var8 .
193 ?var4 dbpedia:deathDate ?var10 .
194
195 ### 28
196 ?var4 dbpedia:prop:birthPlace
197     dbpedia:Papal_States .
198 ?var4 dbpedia:birthDate ?var6 .
199 ?var4 foaf:name ?var8 .
200 ?var4 dbpedia:deathDate ?var10 .

```

| | | | | | |
|-----|--|--|--|-----|---|
| 200 | | | | 273 | ?var3 skos:broader ?var4 . |
| 201 | ### 29 | | | 274 | ?var3 rdfs:label ?var6 . |
| 202 | ?var4 dbpedia:prop:birthPlace | | | 275 | |
| 203 | dbpediares:City_of_London . | | | 276 | ### 42 |
| 204 | ?var4 dbpedia:birthDate ?var6 . | | | 277 | ?var4 rdfs:label "IX\u00F3\u00F5"@de . |
| 205 | ?var4 foaf:name ?var8 . | | | 278 | ?var3 skos:broader ?var4 . |
| 206 | ?var4 dbpedia:deathDate ?var10 . | | | 279 | ?var3 rdfs:label ?var6 . |
| 207 | | | | 280 | |
| 208 | ### 30 | | | 281 | ### 43 |
| 209 | ?var4 dbpedia:prop:birthPlace | | | 282 | ?var4 rdfs:label "(1083) Salvia"@de . |
| 210 | "Houghton, Norfolk, England"@en . | | | 283 | ?var3 skos:broader ?var4 . |
| 211 | ?var4 dbpedia:birthDate ?var6 . | | | 284 | ?var3 rdfs:label ?var6 . |
| 212 | ?var4 foaf:name ?var8 . | | | 285 | |
| 213 | ?var4 dbpedia:deathDate ?var10 . | | | 286 | ### 44 |
| 214 | | | | 287 | ?var4 rdfs:label |
| 215 | ### 31 | | | 288 | "(1296) Andr\u00E9"@de . |
| 216 | ?var4 rdfs:label "(372) Palma"@de . | | | 289 | ?var3 skos:broader ?var4 . |
| 217 | ?var3 skos:broader ?var4 . | | | 290 | ?var3 rdfs:label ?var6 . |
| 218 | ?var3 rdfs:label ?var6 . | | | 291 | |
| 219 | | | | 292 | ### 45 |
| 220 | ### 32 | | | 293 | ?var1 rdf:type yago:ChristianLGBTPeople . |
| 221 | ?var4 rdfs:label "(11554) Asios"@de . | | | 294 | ?var1 foaf:givenName ?var2 . |
| 222 | ?var3 skos:broader ?var4 . | | | 295 | |
| 223 | ?var3 rdfs:label ?var6 . | | | 296 | ### 46 |
| 224 | | | | 297 | ?var1 rdf:type yago:DefJamRecordingsArtists . |
| 225 | ### 33 | | | 298 | ?var1 foaf:givenName ?var2 . |
| 226 | ?var4 rdfs:label "(3080) Moisseiev"@de . | | | 299 | |
| 227 | ?var3 skos:broader ?var4 . | | | 300 | ### 47 |
| 228 | ?var3 rdfs:label ?var6 . | | | 301 | ?var1 rdf:type yago:IndianFilmActors . |
| 229 | | | | 302 | ?var1 foaf:givenName ?var2 . |
| 230 | ### 34 | | | 303 | |
| 231 | ?var4 rdfs:label "(1273) Helma"@de . | | | 304 | ### 48 |
| 232 | ?var3 skos:broader ?var4 . | | | 305 | ?var1 rdf:type yago:EnglishKeyboardists . |
| 233 | ?var3 rdfs:label ?var6 . | | | 306 | ?var1 foaf:givenName ?var2 . |
| 234 | | | | 307 | |
| 235 | ### 35 | | | 308 | ### 49 |
| 236 | ?var4 rdfs:label | | | 309 | ?var1 rdf:type yago:GuitarPlayers . |
| 237 | "(119878) 2002 CY224"@en . | | | 310 | ?var1 foaf:givenName ?var2 . |
| 238 | ?var3 skos:broader ?var4 . | | | 311 | |
| 239 | ?var3 rdfs:label ?var6 . | | | 312 | ### 50 |
| 240 | | | | 313 | ?var1 rdf:type yago:FilipinoFemaleModels . |
| 241 | ### 36 | | | 314 | ?var1 foaf:givenName ?var2 . |
| 242 | ?var4 rdfs:label | | | 315 | |
| 243 | "039A\u578B\u6F5C | | | 316 | ### 51 |
| 244 | \u6C34\u8266"@ja . | | | 317 | ?var1 rdf:type yago:BluesBrothers . |
| 245 | ?var3 skos:broader ?var4 . | | | 318 | ?var1 foaf:givenName ?var2 . |
| 246 | ?var3 rdfs:label ?var6 . | | | 319 | |
| 247 | | | | 320 | ### 52 |
| 248 | ### 37 | | | 321 | ?var1 rdf:type yago:AmericanSongwriters . |
| 249 | ?var4 rdfs:label | | | 322 | ?var1 foaf:givenName ?var2 . |
| 250 | "(4444) \u042D | | | 323 | |
| 251 | \u0448\u0435\u0440"@ru . | | | 324 | ### 53 |
| 252 | ?var3 skos:broader ?var4 . | | | 325 | ?var1 rdf:type yago:FrenchJazzViolinists . |
| 253 | ?var3 rdfs:label ?var6 . | | | 326 | ?var1 foaf:givenName ?var2 . |
| 254 | | | | 327 | |
| 255 | ### 38 | | | 328 | ### 54 |
| 256 | ?var4 rdfs:label | | | 329 | ?var1 rdf:type yago:EnglishJazzComposers . |
| 257 | "(3834) Zappafrank"@es . | | | 330 | ?var1 foaf:givenName ?var2 . |
| 258 | ?var3 skos:broader ?var4 . | | | 331 | |
| 259 | ?var3 rdfs:label ?var6 . | | | 332 | ### 55 |
| 260 | | | | 333 | ?var1 rdf:type yago:HarveyMuddCollegeAlumni . |
| 261 | ### 39 | | | 334 | ?var1 foaf:givenName ?var2 . |
| 262 | ?var4 rdfs:label "(2612) Kathryn"@de . | | | 335 | |
| 263 | ?var3 skos:broader ?var4 . | | | 336 | ### 56 |
| 264 | ?var3 rdfs:label ?var6 . | | | 337 | ?var1 rdf:type yago:Bassist109842629 . |
| 265 | | | | 338 | ?var1 foaf:givenName ?var2 . |
| 266 | ### 40 | | | 339 | |
| 267 | ?var4 rdfs:label "(290) Bruna"@de . | | | 340 | ### 57 |
| 268 | ?var3 skos:broader ?var4 . | | | 341 | ?var1 rdf:type yago:Curate109983572 . |
| 269 | ?var3 rdfs:label ?var6 . | | | 342 | ?var1 foaf:givenName ?var2 . |
| 270 | | | | 343 | |
| 271 | ### 41 | | | 344 | ### 58 |
| 272 | ?var4 rdfs:label "(438) Zeuxo"@de . | | | 345 | ?var1 rdf:type yago:GreekFemaleModels . |

```

346 ?var1 foaf:givenName ?var2 .
347
348 ### 59
349 ?var1 rdf:type yago:FilipinoReligiousLeaders .
350 ?var1 foaf:givenName ?var2 .
351
352 ### 60
353 ?var4 skos:subject
354     dbpediares:Category:1004_deaths .
355 ?var4 foaf:name ?var6 .
356
357 ### 61
358 ?var4 skos:subject
359     dbpediares:Category:
360         11th_century_in_England .
361 ?var4 foaf:name ?var6 .
362
363 ### 62
364 ?var4 skos:subject
365     dbpediares:Category:1067_deaths .
366 ?var4 foaf:name ?var6 .
367
368 ### 63
369 ?var4 skos:subject
370     dbpediares:Category:1107_births .
371 ?var4 foaf:name ?var6 .
372
373 ### 64
374 ?var4 skos:subject
375     dbpediares:Category:%C5%A0koda_trams .
376 ?var4 foaf:name ?var6 .
377
378 ### 65
379 ?var4 skos:subject
380     dbpediares:Category:
381         %C3%81guilas_Cibae%C3%B1as_players .
382 ?var4 foaf:name ?var6 .
383
384 ### 66
385 ?var4 skos:subject
386     dbpediares:Category:1255_births .
387 ?var4 foaf:name ?var6 .
388
389 ### 67
390 ?var4 skos:subject
391     dbpediares:Category:0s_BC_births .
392 ?var4 foaf:name ?var6 .
393
394 ### 68
395 ?var4 skos:subject
396     dbpediares:Category:
397         1130_disestablishments .
398 ?var4 foaf:name ?var6 .
399
400 ### 69
401 ?var4 skos:subject
402     dbpediares:Category:
403         .32_S%26W_Long_firearms .
404 ?var4 foaf:name ?var6 .
405
406 ### 70
407 ?var4 skos:subject
408     dbpediares:Category:1009_deaths .
409 ?var4 foaf:name ?var6 .
410
411 ### 71
412 ?var4 skos:subject
413     dbpediares:Category:1144_deaths .
414 ?var4 foaf:name ?var6 .
415
416 ### 72
417 ?var4 skos:subject
418     dbpediares:Category:1239_deaths .

```

```

419 ?var4 foaf:name ?var6 .
420
421 ### 73
422 ?var4 skos:subject
423     dbpediares:Category:1070s_deaths .
424 ?var4 foaf:name ?var6 .
425
426 ### 74
427 ?var4 skos:subject
428     dbpediares:Category:1105 .
429 ?var4 foaf:name ?var6 .
430
431 ### 75
432 ?var3 dbpedia:influenced
433     dbpediares:Ram%C3%B3n_Emeterio_Betances .
434 ?var3 foaf:page ?var4 .
435 ?var3 rdfs:label ?var6 .
436
437 ### 76
438 ?var3 dbpedia:influenced dbpediares:Rob_Corddry .
439 ?var3 foaf:page ?var4 .
440 ?var3 rdfs:label ?var6 .
441
442 ### 77
443 ?var3 dbpedia:influenced
444     dbpediares:Parakrama_Niriella .
445 ?var3 foaf:page ?var4 .
446 ?var3 rdfs:label ?var6 .
447
448 ### 78
449 ?var3 dbpedia:influenced
450     dbpediares:Alexander_VI .
451 ?var3 foaf:page ?var4 .
452 ?var3 rdfs:label ?var6 .
453
454 ### 79
455 ?var3 dbpedia:influenced dbpediares:Iqbal .
456 ?var3 foaf:page ?var4 .
457 ?var3 rdfs:label ?var6 .
458
459 ### 80
460 ?var3 dbpedia:influenced
461     dbpediares:Al-Maqrizi .
462 ?var3 foaf:page ?var4 .
463 ?var3 rdfs:label ?var6 .
464
465 ### 81
466 ?var3 dbpedia:influenced
467     dbpediares:Clarence_Irving_Lewis .
468 ?var3 foaf:page ?var4 .
469 ?var3 rdfs:label ?var6 .
470
471 ### 82
472 ?var3 dbpedia:influenced
473     dbpediares:Ibn_Khaleel .
474 ?var3 foaf:page ?var4 .
475 ?var3 rdfs:label ?var6 .
476
477 ### 83
478 ?var3 dbpedia:influenced
479     dbpediares:David_Friedl%C3%A4nder .
480 ?var3 foaf:page ?var4 .
481 ?var3 rdfs:label ?var6 .
482
483 ### 84
484 ?var3 dbpedia:influenced
485     dbpediares:John_Warnock .
486 ?var3 foaf:page ?var4 .
487 ?var3 rdfs:label ?var6 .
488
489 ### 85
490 ?var3 dbpedia:influenced
491     dbpediares:Vladimir_Lenin .

```

```

492 ?var3 foaf:page ?var4 .
493 ?var3 rdfs:label ?var6 .
494
495 ### 86
496 ?var3 dbpedia:influenced
497     dbpediares:Niall_McLaren .
498 ?var3 foaf:page ?var4 .
499 ?var3 rdfs:label ?var6 .
500
501 ### 87
502 ?var3 dbpedia:influenced
503     dbpediares:David_J._Farber .
504 ?var3 foaf:page ?var4 .
505 ?var3 rdfs:label ?var6 .
506
507 ### 88
508 ?var3 dbpedia:influenced
509     dbpediares:Fran_Lebowitz .
510 ?var3 foaf:page ?var4 .
511 ?var3 rdfs:label ?var6 .
512
513 ### 89
514 ?var3 dbpedia:influenced
515     dbpediares:Kathleen_Raine .
516 ?var3 foaf:page ?var4 .
517 ?var3 rdfs:label ?var6 .
518
519 ### 90
520 ?var0 rdfs:label "The Subtle Knife"@en .
521 ?var0 rdf:type ?var1 .
522
523 ### 91
524 ?var0 rdfs:label "Patrioter"@sv .
525 ?var0 rdf:type ?var1 .
526
527 ### 92
528 ?var0 rdfs:label "Scar Tissue (libro)"@es .
529 ?var0 rdf:type ?var1 .
530
531 ### 93
532 ?var0 rdfs:label
533     "Jason Bournes ultimatum"@nn .
534 ?var0 rdf:type ?var1 .
535
536 ### 94
537 ?var0 rdfs:label "Jane Eyre"@fr .
538 ?var0 rdf:type ?var1 .
539
540 ### 95
541 ?var0 rdfs:label
542     "Gone with the Wind (livro)"@pt .
543 ?var0 rdf:type ?var1 .
544
545 ### 96
546 ?var0 rdfs:label "Alkumets\u00E4"@fi .
547 ?var0 rdf:type ?var1 .
548
549 ### 97
550 ?var0 rdfs:label
551     "Flight from the Dark"@en .
552 ?var0 rdf:type ?var1 .
553
554 ### 98
555 ?var0 rdfs:label "Mongol Empire"@en .
556 ?var0 rdf:type ?var1 .
557
558 ### 99
559 ?var0 rdfs:label "Kultahattu"@fi .
560 ?var0 rdf:type ?var1 .
561
562 ### 100
563 ?var0 rdfs:label
564     "Aseiden k\u00E4ytt\u00F6"@fi .

```

```

565 ?var0 rdf:type ?var1 .
566
567 ### 101
568 ?var0 rdfs:label "Marrow (novel)"@en .
569 ?var0 rdf:type ?var1 .
570
571 ### 102
572 ?var0 rdfs:label "The Acid House"@en .
573 ?var0 rdf:type ?var1 .
574
575 ### 103
576 ?var0 rdfs:label "\u06B3\u04E49\u08BBA"@zh .
577 ?var0 rdf:type ?var1 .
578
579 ### 104
580 ?var0 rdfs:label "Dawn of the Dragons"@en .
581 ?var0 rdf:type ?var1 .
582
583 ### 105
584 ?var2 rdf:type dbpedia:Person .
585 ?var2 rdfs:label
586     "Ab\u016B l-Hasan Ban\u012Bsadr"@de .
587 ?var2 foaf:page ?var4 .
588
589 ### 106
590 ?var2 rdf:type dbpedia:Person .
591 ?var2 rdfs:label
592     "Abdul Rahman of Negeri Sembilan"@en .
593 ?var2 foaf:page ?var4 .
594
595 ### 107
596 ?var2 rdf:type dbpedia:Person .
597 ?var2 rdfs:label "A.W. Farwick"@en .
598 ?var2 foaf:page ?var4 .
599
600 ### 108
601 ?var2 rdf:type dbpedia:Person .
602 ?var2 rdfs:label "Abdullah G\u00F6l"@sv .
603 ?var2 foaf:page ?var4 .
604
605 ### 109
606 ?var2 rdf:type dbpedia:Person .
607 ?var2 rdfs:label "Aaron Pe\u00F1a"@en .
608 ?var2 foaf:page ?var4 .
609
610 ### 110
611 ?var2 rdf:type dbpedia:Person .
612 ?var2 rdfs:label "Abby Lockhart"@fr .
613 ?var2 foaf:page ?var4 .
614
615 ### 111
616 ?var2 rdf:type dbpedia:Person .
617 ?var2 rdfs:label "Abd al-Latif"@en .
618 ?var2 foaf:page ?var4 .
619
620 ### 112
621 ?var2 rdf:type dbpedia:Person .
622 ?var2 rdfs:label "Abdel Halim Khaddam"@fr .
623 ?var2 foaf:page ?var4 .
624
625 ### 113
626 ?var2 rdf:type dbpedia:Person .
627 ?var2 rdfs:label "Abdur Rahman Khan"@fr .
628 ?var2 foaf:page ?var4 .
629
630 ### 114
631 ?var2 rdf:type dbpedia:Person .
632 ?var2 rdfs:label
633     "A\u0142\u0142a Kudriawcewa"@pl .
634 ?var2 foaf:page ?var4 .
635
636 ### 115
637 ?var2 rdf:type dbpedia:Person .

```

```

638 ?var2 rdfs:label "Aaron Raper"@en .
639 ?var2 foaf:page ?var4 .
640
641 ### 116
642 ?var2 rdf:type dbpedia:Person .
643 ?var2 rdfs:label "A.L. Williams"@en .
644 ?var2 foaf:page ?var4 .
645
646 ### 117
647 ?var2 rdf:type dbpedia:Person .
648 ?var2 rdfs:label "Abdul Kadir Khan"@sv .
649 ?var2 foaf:page ?var4 .
650
651 ### 118
652 ?var2 rdf:type dbpedia:Person .
653 ?var2 rdfs:label "A. J. Pierzynski"@en .
654 ?var2 foaf:page ?var4 .
655
656 ### 119
657 ?var2 rdf:type dbpedia:Person .
658 ?var2 rdfs:label "Abdullah Ahmad Badawi"@pl .
659 ?var2 foaf:page ?var4 .
660
661 ### 120
662 ?var2 rdf:type dbpedia:Person .
663 ?var2 rdfs:label
664 "Abdelbaset Ali Mohmed Al Megrahi"@en .
665 ?var2 foaf:page ?var4 .

```