

Appendix

A1: Scalability of the generated instances of the 10 tables

20 K Products (s) :	
Vendor : 203	Product : 20000
Person : 10210	Review : 200000
ProductFeature : 10519	ProductFeatureProduct : 423222
ProductType : 329	ProductTypeProduct : 80,000
Offer : 400000	Producer : 405
200K Products (m) :	
Vendor : 2027	Product : 200000
Person : 102,596	Review : 2000000
ProductFeature : 47884	ProductFeatureProduct : 3885664
ProductType : 2011	ProductTypeProduct : 1000000
Offer : 400000	Producer : 3956
2000K Products (l) :	
Vendor : 19969	Product : 2000000
Person : 1024622	Review : 19999800
ProductFeature : 94259	ProductFeatureProduct : 38571526
ProductType : 3949	ProductTypeProduct : 10000000
Offer : 40000000	Producer : 39530

A2: Pre-computed inter-source hints

EFJ1 : $\text{Product1}_1 \bowtie_{\text{Product1}_1.nr \equiv \mathbb{D}} \text{Product2}_2.nr \text{ Product2}_2 = \emptyset$
EFJ2 : $\text{ProductFeatureProduct1}_1 \bowtie_{\text{ProductFeatureProduct1}_1.product \equiv \mathbb{D}} \text{Product2}_2.nr \text{ Product2}_2 = \emptyset$
EFJ3 : $\text{Product1}_1 \bowtie_{\text{Product1}_1.nr \equiv \mathbb{D}} \text{ProductFeatureProduct2}_1.product \text{ ProductFeatureProduct2}_1 = \emptyset$
DRS : $\text{ReviewC}_1(...) \equiv \mathbb{D} \text{ Review}_5(...)$
$op1(...) \leftarrow \pi_{...}(\text{Offer}_4 \bowtie_{\text{Offer}_4.product = \text{Product1}_1.nr} \text{Product1}_1)$ $op2(...) \leftarrow \pi_{...}(\text{Offer}_4 \bowtie_{\text{Offer}_4.product = \text{Product2}_2.nr} \text{Product2}_2)$ $pfpf1(...) \leftarrow \pi_{...}(\text{ProductFeature}_3 \bowtie_{\text{ProductFeature}_3.nr = \text{ProductFeatureProduct1}_1.product.feature} \text{ProductFeatureProduct1}_1)$ $pfpf2(...) \leftarrow \pi_{...}(\text{ProductFeature}_3 \bowtie_{\text{ProductFeature}_3.nr = \text{ProductFeatureProduct2}_2.product.feature} \text{ProductFeatureProduct2}_2)$ $ppd1(...) \leftarrow \pi_{...}(\text{Product1}_1 \bowtie_{\text{Product1}_1.nr = \text{Producer}_4.product} \text{Producer}_4)$ $ppd2(...) \leftarrow \pi_{...}(\text{Product2}_2 \bowtie_{\text{Product2}_2.nr = \text{Producer}_4.product} \text{Producer}_4)$

Here, T_i denotes T is a table/relation in the VDB Σ for data source S_i (or S'_i), and \mathbb{D} is the instance for each instance of the OBDF specification.

A3: Concrete query evaluation results

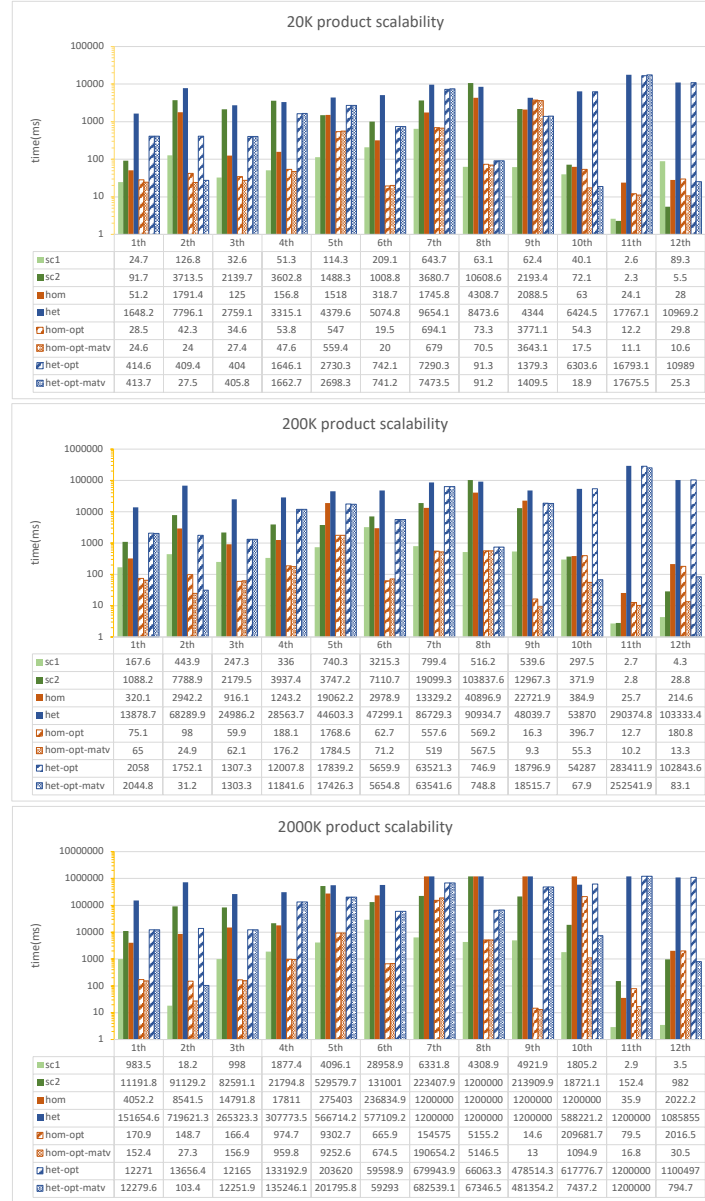


Fig. 5. Milliseconds used of evaluating the SQL queries over different configurations, where 120000000 (ms) denotes timeout (i.e., the query cannot be evaluated in 20min).

A4: Analysis of the hints for SQL query optimization/simplification

Here, we analyze the performance improved by the inter-source hints in the ODBF settings. For easy describing, u_j , v_k and w_l are used to denote the expressions like π_a , σ_b and $\pi_a\sigma_b$.

Query q_1 - q_6 benefit from the hints of pre-computed empty federated joins, i.e., EFJ1–EFJ3. Such rewriting reduces the computation of federated joins in the SQL queries significantly, such as rewriting the federated joins like:

$$(u_1(\text{Product1}_1) \cup v_1(\text{Product2}_2)) \bowtie_{\theta} (u_2(\text{Product1}_1) \cup v_2(\text{Product2}_2))$$

into $u(\text{Product1}_1) \cup v(\text{Product2}_2)$, and delivers more computation to local sources to be evaluated in a parallel way. For example, after rewriting terminal federated joins in the query using the approach described from line 11–17 of algorithm $\text{unfold}_{\text{ODBF}}$, the SQL translations q_3 and q_4 of Q_3 and Q_4 are rewritten into the queries with the format $sql_1 \cup sql_2$, where sql_1 is a query only referring source S_1 and sql_2 only referring source S_2 (S'_2), and the SQL translation q_2 of Q_2 is rewritten into q'_2 :

$$(u_1(\text{ProductFeature}_3) \bowtie_{\theta_1} u_2(\text{ProductFeatureProduct}_1) \bowtie_{\theta_2} u_3(\text{Product1}_1) \bowtie_{\theta_3} u_4(\text{Producer}_4)) \cup (v_1(\text{ProductFeature}_3) \bowtie_{\varphi_1} v_2(\text{ProductFeatureProduct}_2) \bowtie_{\varphi_2} v_3(\text{Product2}_2) \bowtie_{\varphi_3} v_4(\text{Producer}_4))$$

Take the 200K product scalability and hom configuration as an example. Empty federated join based rewriting makes the evaluation times of q_2 , q_3 and q_4 reduced respectively from 2942.2 ms, 916.1 ms and 1243.2 ms to 98 ms, 59.9 ms and 188.1 ms (even more faster than in the sc1 situation).

The SQL translations q_6 – q_9 of Q_6 – Q_9 benefit from redundancy based rewriting. Such rewriting removes the equivalent/containment redundancy w.r.t. the pre-computed data redundancy in the union expressions, such as removing one element from the union $(v(\text{ReviewC}_1) \bowtie_{\theta} A) \cup (u(\text{Review}_5) \bowtie_{\theta} A)$ by taking data distribution and heterogeneity into consideration, so as to reduce the computation of federated joins, such as rewriting the join:

$$(u_1(\text{ReviewC}_1) \cup v_1(\text{Review}_5)) \bowtie_{\theta} (u_2(\text{ReviewC}_1) \cup v_2(\text{Review}_5))$$

as $u(\text{ReviewC}_1) \cup v(\text{Review}_5)$, and the access of inefficient data sources. For example, after redundancy based rewriting, q_9 is rewritten as a query only referring source S_2 (S'_2). After removing redundancy, the evaluation of these three queries is as fast as in the sc1 situation.

The SQL query translation q_2 , q_{10} and q_{12} of the SPARQL queries Q_2 , Q_{10} and Q_{12} benefit from the materialized views. Here q'_{10} and q_{10} are the same, and q'_{12} and q_{12} are the same. For example, after applying all the hints, q_2 is further rewritten as:

$$q''_2 : (u(S^M.ppd1) \bowtie_{\theta} v(S^M.pfpf1)) \cup (u'(S^M.ppd2) \bowtie_{\varphi} v'(S^M.pfpf2))$$

which is a query over S^M . No computation of federated joins and accessing of inefficient sources is needed. In the 2M product scalability, for hom and het, it makes the query evaluation time further reduced respectively from 8541.5 ms and 719621.3 ms to 27.3 ms and 103.4 ms. Moreover, q_{10} is the query:

$$((u_1(\text{Product1}_1) \bowtie_{\theta_1} v_1(\text{Offer}_4) \bowtie_{\varphi_1} w_1(\text{Vendor}_4))) \cup ((u_2(\text{Product2}_2) \bowtie_{\theta_2} v_2(\text{Offer}_4) \bowtie_{\varphi_2} w_2(\text{Vendor}_4)))$$

Based on the materialization $op1$ and $op2$ in S^M , q_{10} is finally rewritten as the query $q''_{10} : (u(S^M.op1) \bowtie_{\theta'} v(\text{Vendor}_4)) \cup (u'(S^M.op2) \bowtie_{\varphi'} v'(\text{Vendor}_4))$. Again take 2M products scalability as an example. In the federation setting, i.e., hom and het, the rewriting reduced the evaluation times respectively from 209681.7 ms and 617776.8 ms to 1094.9 ms and 7437.2 ms. Although in the heterogeneous situation, S'_4 is a set of CSV files, it does not effect the overall performance much, since the source only needs to evaluate simple selecting queries.

In summary, the overall evaluation results indicate the rationality and importance of our research as well as the high-efficiency and power of the proposed solution in improving the overall all query answering performance.