APPENDIX

A1: Statistics of the generated instances (Table : NumOfRows)

20K Products :	
Vendor: 203; Product1: 10,000; Person: 10210; ProductFeature: 10,519; ProductFeatureProduct1: 212,342; ProductType: 329; ProductTypeProduct1: 40,000; Offer: 400,000;	Product: 20,000 Product2: 10,000 Review: 2000000; ReviewC: 200,000 ProductFeatureProduct: 423,222 ProductFeatureProduct2: 210,880 ProductTypeProduct: 80,000 ProductTypeProduct: 40,000 Producer: 405
200K Prooducts :	
Vendor: 2,027; Person: 102,596; ProductFeature: 47,884; ProductFeatureProduct1: 1,941,862; ProductType: 2,011; ProductTypeProduct1: 500,000; Offer: 400,000;	Product: 200,000 Review: 2,000,000; ReviewC: 2,000,000 ProductFeatureProduct: 3,885,664 ProductFeatureProduct2: 1,943,802 ProductTypeProduct: 1,000,000 ProductTypeProduct2: 500,000 Producer: 3,956
2M Products :	
Vendor: 19, 969; Person: 1, 024, 622; ProductFeature: 94, 259; ProductFeatureProduct1: 19, 285, 093; ProductType: 3, 949; ProductTypeProduct1: 5, 000, 000; Offer: 40, 000, 000;	Product: 2,000,000 Review: 1,9999,800; ReviewC: 19,999,800 ProductFeatureProduct: 3,8571,526 ProductFeatureProduct2: 19,286,433 ProductTypeProduct: 10,000,000 ProductTypeProduct2: 5,000,000 Producer: 39,530

A2: Pre-computed hints

Next we show the source hints used for the experiment. The top block contains the empty joins hints, the middle block contains the data redundancy hints, and the bottom one contains the materialized views. Here, T_i denotes T is a table/relation in the VDB Σ for data source S_i , and $\mathbb{D} \in \{s, m, l\}$ denotes one of the three instances considered in Appendix A1.

```
\begin{split} & \text{EFJ1: Product1}_1 \bowtie_{\text{Product1}_1.nr} = \text{Product2}_{2.nr} \ \text{Product2}_2 =_{\mathbb{D}} \ \emptyset \\ & \text{EFJ2: ProductFeatureProduct1}_1 \bowtie_{\text{Product1}_1.nr} = \text{Product2}_{2.nr} \ \text{Product3}_{2.nr} \ \text{Product4}_{2.nr} \ \text{Product4}_{2.nr}
```

The statics for the materialized views, over each of the three data instances, are listed in the table below:

20K Products :	ts: 200K Prooducts:		2M Products :	
op1: 200, 067; op2: 199, pfpf1: 212, 342; pfpf2: 210 ppd1: 10, 000; ppd2: 10,	,880 pfpf1: 1,941,862;	pfpf2: 1943802	op1: 20,005,677; pfpf1: 19,285,093; ppd1: 1,000,000;	op2: 19, 994, 323 pfpf2: 19, 286, 433 ppd2: 1,000,000

A3: Query evaluation

A4: Analysis of the experimental results

In this section, we analyze the impact hints-based optimizations over the queries in our experiment. Next, we use q_1 – q_{12} to denote the translation from Q1–Q12 performed by Ontop (hence, with classical OBDA optimizations) in the federation situation.

Empty federated joins. The pre-computed empty-federated joins have an impact on queries q_1 - q_6 . For instance, Ontop translation q_3 , containing 5 terminal federated joins (joins of unions) and 6 union operations, is optimized into a query without federated joins and a single union. Similarly, translation q_4 , containing 8 terminal federated joins and 11 union operations, is optimized into a query without federated joins and 3 unions.

These optimizations pave the road to other optimizations that can be applied later, for instance w.r.t. materialized views. For instance, consider query q_2 , containing 17 terminal federated joins and 14 union operations. Such translation is optimized into a query with 6 federated

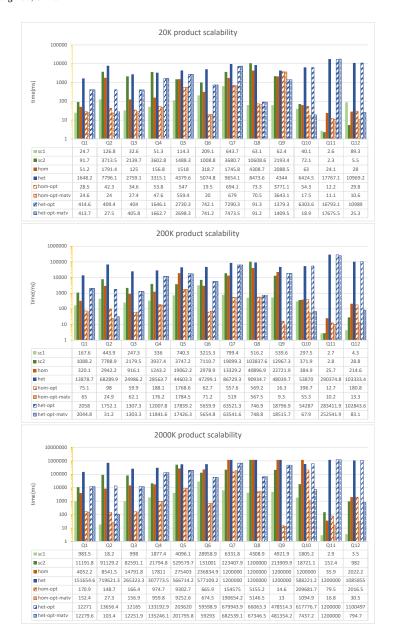


Figure 6: The figure reports the results (time in ms, average of 10 executions) of the SQL evaluation, over the three data instances 20K, 200K and 2M products. The timeout for each execution is set to 1200 seconds.

joins and a single union. Such optimized query, that we here call q_2 , has the following structure:

```
\begin{array}{l} (u_1(\mathsf{ProductFeature}_3) \bowtie_{\theta_1} u_2(\mathsf{ProductFeatureProduct}_1) \bowtie_{\theta_2} u_3(\mathsf{Product}_1) \bowtie_{\theta_3} u_4(\mathsf{Produce}_4)) \cup \\ (v_1(\mathsf{ProductFeature}_3) \bowtie_{\varphi_1} v_2(\mathsf{ProductFeatureProduct}_2) \bowtie_{\varphi_2} v_3(\mathsf{Product}_2) \bowtie_{\varphi_3} v_4(\mathsf{Produce}_4)) \end{array}
```

We will see in the remainder of this section that this query can be further optimized by exploiting the pre-computed materialized views.

These optimizations lead to a considerable reduction of the query evaluation time. For the 200K instance and the hom configuration, the evaluation times of q_2 , q_3 and q_4 go from 2942.2 ms, 916.1 ms and 1243.2 ms to 98 ms, 59.9 ms and 188.1 ms, respectively. Note that these executions are even faster than those for the "easiest" configuration, sc1.

Data redundancy. The pre-computed data redundancy hint DRS has an impact on translations q_6-q_9 . After optimizing with respect to this hint, in fact, the evaluation of these three queries becomes as fast as for the sc1 situation. This improvement the hint allows for the

removal of several redundant operations. Consider, for instance, SQL query q_9 , containing 5 terminal federated joins and 4 union operations. Optimizing such translation with respect to data redundancy leads to a query without federated joins and unions.

Materialized views. The pre-computed materialized views have an impact on translations q_2 , q_{10} and q_{12} . For instance, after applying the optimizations based on empty federated joins hints, query q'_2 above is further simplified into

$$q_2'': (\mathsf{u}(S^\mathsf{M}.ppd1) \bowtie_\theta \mathsf{v}(S^\mathsf{M}.pfpf1)) \cup (\mathsf{u}'(S^\mathsf{M}.ppd2) \bowtie_\varphi \mathsf{v}'(S^\mathsf{M}.pfpf2))$$

Observe that, for such query, no computation of federated joins or access to inefficient sources is performed. This simplification results, for the 2M instance and for the hom and het settings, in a further reduction of the evaluation times respectively from 8541.5 ms and 719621.3 ms to 27.3 ms and 103.4 ms.

Query q_{10} , instead, is transformed from the following algebra expression:

$$\begin{array}{l} ((\mathsf{u}_1(\mathsf{Product1}_1) \bowtie_{\theta_1} \mathsf{v}_1(\mathsf{Offer}_4) \bowtie_{\phi_1} \mathsf{w}_1(\mathsf{Vendor}_4))) \cup \\ ((\mathsf{u}_2(\mathsf{Product2}_2) \bowtie_{\theta_2} \mathsf{v}_2(\mathsf{Offer}_4) \bowtie_{\phi_2} \mathsf{w}_2(\mathsf{Vendor}_4))) \end{array}$$

into the query

$$q_{10}^{\prime\prime}: (\mathsf{u}(S^{\mathsf{M}}.op1)\bowtie_{\theta^{\prime}}\mathsf{v}(\mathsf{Vendor}_4)) \cup (\mathsf{u}^{\prime}(S^{\mathsf{M}}.op2)\bowtie_{\phi^{\prime}}\mathsf{v}^{\prime}(\mathsf{Vendor}_4))$$

This simplification results, for the 2M instance and for the hom and het settings, in a reduction of the evaluation times respectively from 209681.7 ms and 617776.8 ms to 1094.9 ms and 7437.2 ms.