

HADAD: A Lightweight Approach for Optimizing Hybrid Complex Analytics Queries

Anonymous Author(s)

ABSTRACT

Hybrid complex analytics workloads typically include (i) data management tasks (joins, selections, etc.), easily expressed using relational algebra (RA)-based languages, and (ii) complex analytics tasks (regressions, matrix decompositions, etc.), mostly expressed in linear algebra (LA) expressions. Such workloads are common in many application areas, including scientific computing, web analytics, and business recommendation. Existing solutions for evaluating hybrid analytical tasks – ranging from LA-oriented systems, to relational systems (extended to handle LA operations), to hybrid systems – either optimize data management and complex tasks separately, exploit RA properties only while leaving LA-specific optimization opportunities unexploited, or focus heavily on low-level physical optimization, and leaving semantic query optimization opportunities unexplored. Additionally, they are not able to exploit precomputed (materialized) results to avoid computing again (part of) a given mixed (LA and/or RA) computation.

In this paper, we take a major step towards filling this gap by proposing a unified semantic query optimization framework for such a hybrid setting. We describe HADAD, an extensible lightweight approach for optimizing hybrid complex analytics queries, based on a common abstraction that facilitates a unified reasoning: a relational model endowed with integrity constraints, which can be used to express the properties of the two computation formalisms. Our solution can be naturally and portably applied on top of cross RA-LA (or LA) platforms; while their existing code need not be rewritten. An extensive empirical evaluation shows that HADAD yields significant performance gains on diverse workloads, from LA-centered ones to hybrid ones.

CCS CONCEPTS

- Data management systems → Semantic query optimization.

KEYWORDS

linear algebra, query rewriting, integrity constraints, chase

ACM Reference Format:

Anonymous Author(s). 2021. HADAD: A Lightweight Approach for Optimizing Hybrid Complex Analytics Queries. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 62 pages. <https://doi.org/10.1145/nnnnnnnnnnnnnn>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference'17, July 2017, Washington, DC, USA

© 2021 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnnnnnnnnn>

1 INTRODUCTION

Modern analytical tasks typically include (i) data management tasks (e.g., joins, filters) to perform pre-processing steps, including feature selection, transformation, and engineering [20, 24, 37, 45, 47], tasks that are easily expressed using RA-based languages, as well as (ii) complex analytics tasks (e.g., regressions, matrices decompositions), which are mostly expressed using LA operations [36]. To perform such analytical tasks, data scientists can choose from a variety of systems, tools, and languages. Languages/libraries such as R [8] and NumPy [7], as well as LA systems such as SystemML [21], TensorFlow [13] and MLLib [10] treat matrices and LA operations as first-class citizens: they offer a rich set of built-in LA operations. However, it can be difficult to express pre-processing tasks in these systems. Further, expression rewrites, based on equivalences that hold due to well-known LA properties, are not exploited in some of these systems, leading to missed optimization opportunities.

Many works propose integrating RA and LA processing, where both algebraic styles can be used together [9, 23, 27, 35, 38, 40, 49]. [9] offers calling LA packages through user defined functions (UDFs), where libraries such as NumPy are embedded in the host language. Others suggest to extend RDBMS to treat LA objects as first-class citizens by using built-in functions to express LA operations [35, 40]. However, LA operations' semantics remain hidden behind these functions, where the optimizers treat as black-boxes. SPORES [49] and SPOOF [23] optimize LA expressions by converting them into RA, optimizing the latter, and then converting the result back to an (optimized) LA expression. They only focus on optimizing LA pipelines containing operations that can be expressed in RA. The restriction is that LA properties of complex operations such as inverse, matrix-decompositions are entirely unexploited. Morpheus [27] speeds up LA pipelines over large joins by pushing computation into each joined table, thereby avoiding expensive materialization. LARA [38] focuses on low-level optimization by exploiting data layouts (e.g., column-wise) to choose LA operators' physical implementations.

A limitation of such approaches is that they lack high-level reasoning about LA properties and rewrites, which can drastically enhance the pipelines' performance [48]. Further, they do not support semantic query optimization, which includes exploiting integrity constraints and materialized views and can bring enormous performance advantages in hybrid RA-LA and even plain LA settings.

We propose HADAD, an extensible lightweight framework for providing semantic query optimization (including views-based, integrity constraint-based, and LA property-based rewriting) on top of both pure LA and hybrid RA-LA platforms, with no need to modify their internals. At the core of HADAD lies a common abstraction: *relational model with integrity constraints*, which enables reasoning in hybrid settings. Moreover, it makes it very easy to extend HADAD's semantic knowledge of LA operations by simply declaring appropriate constraints, with no need to change HADAD

code. As we show, constraints are sufficiently expressive to declare (and thus allow HADAD to exploit) more properties of LA operations than previous work could consider.

Last but not least, our holistic, cost-based approach enables to judiciously apply for each query *the best available optimization*. For instance, given the computation $M(NP)$ for some matrices M , N and P , we may rewrite it into $(MN)P$ if its estimated cost is smaller than that of the original expression, or we may turn it into MV if a materialized view V stores exactly the result of (NP) .

HADAD capitalizes on a framework previously introduced in [15] for rewriting queries across many data models, using materialized views, in a polystore setting that does not include the LA model. The novelty of HADAD is to extend the benefits of rewriting and views optimizations to pure LA and hybrid RA-LA computations, which are crucial for ML workloads.

Contributions. The paper makes the following contributions:

- ❶ We propose an *extensible lightweight approach to optimize hybrid complex analytics queries*. Our approach can be implemented on top of existing systems without modifying their internals; it is based on a powerful intermediate abstraction that supports reasoning in hybrid settings, namely *a relational model with integrity constraints*.
- ❷ We formalize the problem of *rewriting computations using previously materialized views in hybrid settings*. To the best of our knowledge, ours is the first work that brings views-based rewriting under integrity constraints in the context of LA-based pipelines and hybrid analytical queries.
- ❸ We provide *formal guarantees* for our solution in terms of soundness and completeness.
- ❹ We conduct an extensive set of empirical experiments on typical LA- and hybrid-based expressions, which show the benefits of HADAD.

Outline. The rest of this paper is organized as follows: §2 highlights HADAD’s optimizations that go beyond the state of the art based on *real-world scenarios*, §3 formalizes the query optimization problem in the context of a hybrid setting, §5 provides an end-to-end overview of our approach. §6 presents our novel reduction of the rewriting problem into one that can be solved by existing techniques from the relational setting. §7 describes our extension to the query rewriting engine, integrating two different cost models, to help prune out inefficient rewritings as soon as they are enumerated. We formalize our solution’s guarantees in §8 and present the experiments in §9. We discuss related work and conclude in §10.

2 HADAD OPTIMIZATIONS

We highlight below examples of performance-enhancing opportunities that are exploited by HADAD and not being addressed by LA-oriented and cross RA-LA existing solutions.

LA Pipeline Optimization. Consider the Ordinary-Least Regression (OLS) pipeline: $(X^T X)^{-1} (X^T y)$, where X is a square matrix of size $10K \times 10K$ and y is a vector of size $10K \times 1$. Suppose available a materialized view $V = X^{-1}$. HADAD rewrites the pipeline to $(V(V^T(X^T y)))$, by exploiting the LA properties $(CD)^{-1} = D^{-1}C^{-1}$, $(CD)E = C(DE)$ and $(D^T)^{-1} = (D^{-1})^T$ as well as the view V . The rewriting is more efficient than the original pipeline since it avoids computing the expensive inverse operation. Moreover, it optimizes

the matrix chain multiplication order to minimize the intermediate result size. This leads to a $150\times$ speed-up on MLLib [42]. Current popular LA-oriented systems [7, 8, 13, 21, 42] are not capable of exploiting such rewrites, due to the lack of systematic exploration of standard LA properties and views.

Hybrid RA-LA Optimization. Cross RA-LA platforms such as Morpheus [27], SparkSQL [17] and others [34, 38] can greatly benefit from HADAD’s cross-model optimizations, which can find these platforms’ rewrites but also additional ones that they miss.

Factorization of LA Operations over Joins. For instance, Morpheus implements a powerful optimization that *factorizes an LA operation on a matrix M obtained by joining tables R and S and casting the join result as a matrix. Factorization pushes the LA operation before the join, to operate on R and S , cast as matrices*.

Consider a specific instantiation of factorization: $\text{colSums}(MN)$, where matrix M has size $20M \times 120$ and N has size 120×100 ; both matrices are dense. The LA operation is colSums , which sums up the elements in each column, returning the vector of these sums (the operation is common in ML-relevant algorithms such as K-means-clustering [41]). On this pipeline, Morpheus applies a multiplication factorization rule to push the multiplication by N down to R and S : it computes RN and SN , then concatenates the resulting matrices to obtain MN . This intermediate result is materialized as a dense $20M \times 100$ matrix of size $\approx 16\text{GB}$ (Morpheus is prototyped on R, which uses double precision). Finally, colSums is applied to the intermediate result, reducing it to a vector of size 1×100 .

HADAD can help Morpheus do much better, by pushing the colSums operator to R and S (instead of the multiplication with N), then concatenating the resulting vectors. This leads to much smaller intermediate results, since the combined size of vectors $\text{colSums}(R)$ and $\text{colSums}(S)$ is only 1×120 .

To this end, HADAD rewrites the pipeline to $\text{colSums}(MN)$ by exploiting the LA property $\text{colSums}(AB) = \text{colSums}(A)B$ and applying its cost estimator, which favors rewritings with a small intermediate result size. Evaluating this HADAD-produced rewriting, Morpheus’s multiplication pushdown rule no longer applies, while the colSums pushdown rule is now enabled, leading to $125\times$ speed-up.

Pushing Selection from LA Analysis to RA Preprocessing. Consider another hybrid example on a Twitter dataset [12]. The JSON dataset contains tweet ids, extended tweets, entities including hashtags, filter-level (e.g., sensitive level), media, URL, and tweet text, etc. We implemented it on SparkSQL (with SystemML [21]).

In the preprocessing stage, our SparkSQL query constructs a tweet-hashtag filter-level matrix N of size $2M \times 1000$, for all tweets posted from “USA” mentioning “covid”, where rows are tweets, columns are hashtags, and values are filter-levels.

N is then loaded into SystemML, where rows with filter-level less than 4 are selected. The result undergoes an Alternating Least Square (ALS) [44] computation in order to predict, given a set of tweets and hashtags, the sensitivity level of these tweets.

A core building block of the ALS computation is the LA pipeline $(uv^T - N)v$. In our example, u is a tweet feature vector (of size $2M \times 1$) and v is a hashtag feature vector (1000×1).

We have two materialized views available: V_1 stores the tweet id and text as a text datasource in Solr, and V_2 stores tweet id, hashtag id, and filter-level for all tweets posted from “USA”, and is materialized on disk as CSV file. The rewriting modifies the preprocessing of N by introducing V_1 and V_2 ; it also pushes the filter-level selection from the LA pipeline into the preprocessing stage. To this end, it rewrites $(uv^T - Ny)$ to $uv^T v - Nv$, which is more efficient for two reasons. First, N is ultra sparse (0.00018% non-zero), which renders the computation of Nv extremely efficient. Second, SystemML evaluates the chain $uv^T v$ efficiently, computing $v^T v$ first, which results in a scalar, instead of computing uv^T , which results in a dense matrix of size $2M \times 1000$ (HADAD’s cost model realizes this). Without the rewriting help from HADAD, SystemML is unable to exploit its own efficient operations for lack of awareness of the distributivity property of vector multiplication over matrix addition, $Av + Bv = (A + B)v$. The rewriting achieves 14x speed-up.

HADAD detects and applies all the above-mentioned optimizations combined. It captures RA-, LA-, and cross-model optimizations precisely because it reduces all rewrites to a single setting in which they can synergize: relational rewrites under integrity constraints.

3 PROBLEM STATEMENT

We consider a set of *value domains* \mathcal{D}_i , e.g., \mathcal{D}_1 denotes integers, \mathcal{D}_2 denotes real numbers, \mathcal{D}_3 strings, etc, and two basic data types: *relations* (sets of tuples) and *matrices* (bi-dimensional arrays). Any attribute in a tuple or cell in a matrix is a value from some \mathcal{D}_i . We assume a matrix can be implicitly converted into a relation (the order among matrix rows is lost), and the opposite conversion (each tuple becomes a matrix line, in some order that is unknown, unless the relation was explicitly sorted before the conversion).

We consider a hybrid language \mathcal{L} , comprising a set R_{ops} of (unary or binary) RA operators; concretely, R_{ops} comprises the standard relational selection, projection, and join. We also consider a set L_{ops} of LA operators, comprising: unary (e.g., inversion and transposition) and binary (e.g., matrix product) operators. The full set L_{ops} of LA operations we support is detailed in §6.1. A *hybrid expression* in \mathcal{L} is defined as follows:

- any value from a domain \mathcal{D}_i , any matrix, and any relation, is an expression;
- (RA operators): given some expressions E, E' , $ro_1(E)$ is also an expression, where $ro_1 \in R_{ops}$ is a unary relational operator, and E ’s type matches ro_1 ’s expected input type. The same holds for $ro_2(E, E')$, where $ro_2 \in R_{ops}$ is a binary relational operator (i.e., the join);
- (LA operators): given some expressions E, E' which are either numeric matrices or numbers (which can be seen as degenerate matrices of 1×1), and some real number r , the following are also expressions: $lo_1(E)$ where $lo_1 \in L_{ops}$ is a unary operator, and $lo_2(E, E')$ where $lo_2 \in L_{ops}$ is a binary operator (again, provided that E, E' match the expected input types of the operators).

Clearly, an important set of *equivalence rules* hold over our hybrid expressions, well-known respectively in the RA and the LA literature. These equivalences lead to *alternative evaluation strategies* for each expression. Further, we assume given a (possibly empty) set of

materialized views $\mathcal{V} \in \mathcal{L}$, which have been previously computed over some inputs (matrices and/or relations), and whose results are directly available (e.g., as a file on disk). Detecting when a materialized view can be used instead of evaluating (part of) an expression is another important source of alternative evaluation strategies.

Given an expression E and a *cost model* that assigns a cost (a real number) to an expression, we consider the problem of **identifying the most efficient rewrite** derived from E by: (i) exploiting RA and LA equivalence rules, and/or (ii) replacing part of an expression with a scan of a materialized view equivalent to that expression.

Below, we detail our approach, the equivalence rules we capture, and two alternative cost models we devised for this hybrid setting. Importantly, our solution (based on a *relational encoding with integrity constraints*) capitalizes on the framework previously introduced in [15], where it was used to *rewrite queries using materialized views in a polystore setting*, where the data, views, and query cover a variety of data models (relational, JSON, XML, etc.). Those queries can be expressed in a combination of query languages, including SQL, JSON query languages, XQuery, etc. **The ability to rewrite such queries using heterogeneous views directly and fully transfers to HADAD**: thus, instead of a relation, we could have the (tuple-structured) results of an XML or JSON query; views materialized by joining an XML document with a JSON one and a relational database could also be reused. **The novelty of our work is to extend the benefits of rewriting and view-based optimization to LA computations, crucial for ML workloads**. In §6, we focus on capturing matrix data and LA computations in the relational framework, along with relational data naturally; this enables our novel, holistic optimization of hybrid expressions.

4 PRELIMINARIES

We recall conjunctive queries [26], integrity constraints [14], and query rewriting under constraints [32]; these concepts are at the core of our approach.

4.1 Conjunctive Query and Constraints

A conjunctive query (or simply CQ) Q is an expression of the form $Q(\bar{x}) :- R_1(\bar{y}_1), \dots, R_n(\bar{y}_n)$, where each R_i is a predicate (relation) of some finite arity, and $\bar{x}, \bar{y}_1, \dots, \bar{y}_n$ are tuples of variables or constants. Each $R_i(\bar{y}_i)$ is called a relational atom. The expression $Q(\bar{x})$ is the *head* of the query, while the conjunction of relational atoms $R_1(\bar{y}_1), \dots, R_n(\bar{y}_n)$ is its *body*. All variables in the head are called *distinguished*. Also, every variable in \bar{x} must appear at least once in $\bar{y}_1, \dots, \bar{y}_n$. Different forms of constraints have been studied in the literature [14]. In this work, we use Tuple Generating Dependencies (**TGDs**) and Equality Generating Dependencies (**EGDs**), stated by formulas of the form $\forall x_1, \dots, x_n \phi(x_1, \dots, x_n) \rightarrow \exists z_1, \dots, z_k \psi(y_1, \dots, y_m)$, where $\{z_1, \dots, z_k\} = \{y_1, \dots, y_m\} \setminus \{x_1, \dots, x_n\}$. The constraint’s *premise* ϕ is a possibly empty conjunction of relational atoms over variables x_1, \dots, x_n and possibly constants. The constraint’s *conclusion* ψ is a non-empty conjunction of atoms over variables y_1, \dots, y_m and possibly constants, atoms that are relational ones in the case of **TGDs** or equality atoms – of the form $w = w'$ – in the case of **EGDs**. For instance, consider a relation *Review(paper, reviewer, track)* listing reviewers of papers submitted to a conference’s tracks, and a relation *PC(member, affiliation)* listing the affiliation of every program committee member [29]. The fact that a

paper can only be submitted to a single track is captured by the following EGD: $\forall p \forall r \forall t \forall r' \forall t' \text{Review}(p, r, t) \wedge \text{Review}(p, r', t') \rightarrow t = t'$. We can also express that papers can be reviewed only by PC members by the following TGD: $\forall p \forall r \forall t \text{Review}(p, r, t) \rightarrow \exists a \text{PC}(r, a)$.

4.2 Provenance-Aware Chase & Back-Chase

A key ingredient leveraged in our approach is relational query rewriting using views, in the presence of constraints. The state-of-the-art method for this task, called Chase & Backchase, was introduced in [30] and improved in [32], as the Provenance-Aware Chase & Back-Chase (**PACB** in short). At the core of these methods is the idea to *model views as constraints*, in this way reducing the view-based rewriting problem to constraints-only rewriting. Specifically, for a given view V defined by a query, the constraint V_{IO} states that *for every match of the view body against the input data, there is a corresponding (head) tuple in the view output*, while the constraint V_{OI} states the converse inclusion, i.e., *each view output tuple is due to a view body match*. From a set \mathcal{V} of view definitions, PACB therefore derives a set of view constraints $C_V = \{V_{IO}, V_{OI} \mid V \in \mathcal{V}\}$.

Given a source schema σ with a set of integrity constraints I , a set \mathcal{V} of views defined over σ , and a conjunctive query Q over σ , the **rewriting problem** thus becomes: find every reformulation query ρ over the schema of view names \mathcal{V} that is equivalent to Q under the constraints $I \cup C_V$.

EXAMPLE 4.1. For instance, if $\sigma = \{R, S\}$, $I = \emptyset$, $\tau = \{\mathcal{V}\}$ and we have a view V materializing the join of relations R and S , $V(x, y) :- R(x, z), S(z, y)$, the pair of constraints capturing V is the following:

$$\begin{aligned} V_{IO} : \quad & \forall x \forall z \forall y R(x, z) \wedge S(z, y) \rightarrow V(x, y) \\ V_{OI} : \quad & \forall x \forall y V(x, y) \rightarrow \exists z R(x, z) \wedge S(z, y). \end{aligned}$$

Given the query $Q(x, y) :- R(x, z), S(z, y)$, PACB finds the reformulation $\rho(x, y) :- V(x, y)$. Algorithmically, this is achieved by:

(i) chasing Q with the constraints $I \cup C_V^{IO}$, where $C_V^{IO} = \{V_{IO} \mid V \in \mathcal{V}\}$; intuitively, this enriches (extends) Q with all the consequences that follow from its atoms and the constraints $I \cup C_V^{IO}$.

(ii) restricting the chase result to only the \mathcal{V} -atoms; the result is called the universal plan U .

(iii) annotating each atom of the universal plan U with a unique ID called a provenance term.

(iv) chasing U with the constraints in $I \cup C_V^{OI}$, where $C_V^{OI} = \{V_{OI} \mid V \in \mathcal{V}\}$, and annotating each relational atom a introduced by these chase steps with a provenance formula¹ $\pi(a)$, which gives the set of U -subqueries whose chasing led to the creation of a ; the result of this phase, called the backchase, is denoted B .

(v) matching Q against B and outputting as rewritings the subsets of U that are responsible for the introduction (during the backchase) of the atoms in the image $h(Q)$ of Q ; these rewritings are read off directly from the provenance formula $\pi(h(Q))$.

In our example, I is empty, $C_V^{IO} = \{V_{IO}\}$, and the result of the chase in phase (i) is $Q_1(x, y) :- R(x, z), S(z, y), V(x, y)$. The universal plan obtained in (ii) by restricting Q_1 to the schema of view names is $U(x, y) :- V(x, y)^{p_0}$, where p_0 denotes the provenance term

¹Provenance formulas are constructed from provenance terms using logical conjunction and disjunction.

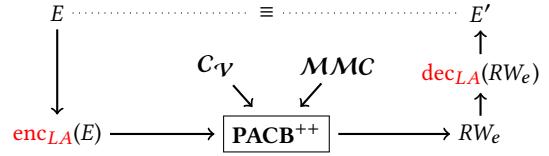


Figure 1: Outline of Our Reduction

of atom $V(x, y)$. The result of backchasing U with C_V^{OI} in phase (iv) is $B(x, y) :- V(x, y)^{p_0}, R(x, z)^{p_0}, S(z, y)^{p_0}$. Note that the provenance formulas of the R and S atoms (a simple term, in this example) are introduced by chasing the view V . Finally, in phase (v) we find one match image given by h from Q 's body into the R and S atoms from B 's body. The provenance formula $\pi(h(Q))$ of the image h is p_0 , which corresponds to an equivalent rewriting $\rho(x, y) :- V(x, y)$.

5 HADAD OVERVIEW

We outline here our approach as an extension to [15] for solving the rewriting problem introduced in §3.

Hybrid Expressions and Views. A hybrid expression (whether asked as a query, or describing a materialized view) can be purely relational (RA), in which case we assume it is specified as a **conjunctive query** [26]. Other expressions are purely LA ones; we assume that they are defined in a dedicated LA language such as R [8], DML [21], etc., using LA operators from our set L_{ops} (see §6.1), commonly used in real-world ML workloads. Finally, a hybrid expression can combine RA and LA, e.g., an RA expression (resulting in a relation) is treated as a matrix input by an LA operator, whose output may be converted again to a table and joined further, etc.

Our approach is based on a reduction to a relational model. Below, we show how to bring our hybrid expressions - and, most specifically, their LA components - under a relational form (the RA part of each expression is already in the target formalism).

Encoding into a Relational Model. Let E be an LA expression (query) and \mathcal{V} be a set of materialized views. We reduce the LA-views based rewriting problem to the relational rewriting problem under integrity constraints, as follows (see Figure 1). First, we encode *relationally* E , \mathcal{V} , and the set L_{ops} of LA operators. Note that the relations used in the encoding are *virtual* and *hidden*, i.e., invisible to both the application designers and users. They only serve to support query rewriting via relational techniques.

These virtual relations are accompanied by a set of relational integrity constraints $enc_{LA}(LA_{prop})$ that reflect a set LA_{prop} of LA properties of the supported L_{ops} operations. For instance, we model the *matrix addition* operation using a relation $\text{add}_M(M, N, R)$, denoting that R is the result of $M+N$, together with a set of constraints stating that add_M is a *functional* relation that is *commutative*, *associative*, etc. These constraints are **Tuple Generating Dependencies (TGDs)** or **Equality Generating Dependencies (EGDs)** [14], which are a generalization of key and foreign key dependencies. We detail our relational encoding in §6.

Reduction from LA-based to Relational Rewriting. Our reduction translates the declaration of each view $V \in \mathcal{V}$ to constraints $enc_{LA}(V)$ that reflect the correspondence between V 's input data and its output. Separately, E is also encoded as a relational query $enc_{LA}(E)$ over the relational encodings of L_{ops} and its matrices.

Operation	Encoding	Operation	Encoding	Operation	Encoding
Matrix scan	$name(M, n)$	Inversion	$inv_M(M, R)$	Cells sum	$sum(M, s)$
Multiplication	$multi_M(M, N, R)$	Scalar Multiplication	$multi_{MS}(s, M, R)$	Row sum	$rowSums(M, R)$
Addition	$add_M(M, N, R)$	Determinant	$det(M, R)$	Colsums	$colSums(M, R)$
Division	$div_M(M, N, R)$	Trace	$trace(M, s)$	Direct sum	$sum_D(M, N, R)$
Hadamard product	$multi_E(M, N, R)$	Exponential	$exp(M, R)$	Direct product	$dirct_D(M, N, R)$
Transposition	$tr(M, R)$	Adjoints	$adj(M, R)$	Diagonal	$diag(M, R)$

Table 1: Snippet of the \mathcal{VRM} Schema

Now, the reformulation problem is reduced to a purely relational setting, as follows. We are given a relational query $enc_{LA}(E)$ and a set $C_V = enc_{LA}(V_1) \cup \dots \cup enc_{LA}(V_n)$ of relational integrity constraints encoding the views \mathcal{V} . We add as further input a set of relational constraints $enc_{LA}(LA_{prop})$, we called them Matrix-Model Encoding constraints, or \mathbf{MMC} in short. We must find the rewritings RW_r^i expressed over the relational views C_V and \mathbf{MMC} , for some integer k and $1 \leq i \leq k$, such that each RW_r^i is equivalent to $enc_{LA}(E)$ under these constraints $(C_V \cup \mathbf{MMC})$. Solving this problem yields a *relationaly encoded rewriting* RWe expressed over the (virtual) relations used in the encoding; a final *decoding* step is needed to obtain E' , the rewriting of (LA or, more generally, hybrid) E using the views \mathcal{V} .

The challenge in coming up with the reduction consists in designing an encoding, i.e., one in which rewritings found by (i) encoding relationally, (ii) solving the resulting relational rewriting problem, and (iii) decoding a resulting rewriting over the views, is guaranteed to produce an equivalent expression E' (see §6).

Relational Rewriting Using Constraints. To solve the relational rewriting problem under constraints, the engine of choice is Provenance Aware Chase&Backchase (PACB) [32]. Our PACB rewriting engine (*PACB⁺⁺* hereafter) has been extended to utilize the *Pruned_{prov}* algorithm sketched and discussed in [31, 32], which prunes inefficient rewritings during the search phase, based on a simple cost model (see §7).

Decoding of the Relational Rewriting. For the selected relational reformulation RWe by *PACB⁺⁺*, a *decoding step* $dec(RWe)$ is performed to translate RWe into the native syntax of its respective underlying store/engine (e.g., R, DML, etc.).

6 REDUCTION TO THE RELATIONAL MODEL

Our internal model is relational, and it makes prominent use of expressive integrity constraints. This framework suffices to describe the features and properties of most data models used today, notably including relational, XML, JSON, graph, etc [15, 16].

Going beyond, in this section, we present a novel way to *reason relationally about LA primitives/operations* by treating them as uninterpreted functions with black-box semantics, and adding *constraints that capture their important properties*. First, we give an overview of a wide range of LA operations that we consider in (§6.1). Then, in (§6.2), we show how matrices and their operations can be represented (*encoded*) using a set of *virtual relations*, part of a schema we call \mathcal{VRM} (for *Virtual Relational Encoding of Matrices*), together with the integrity constraints \mathbf{MMC} that capture the LA properties of these operations. Regardless of matrix data's

physical storage, we only use \mathcal{VRM} to encode LA expressions and views relationally to reason about them. (§6.3) exemplifies relational rewritings obtained via our reduction.

6.1 Matrix Algebra

We consider a wide range of matrix operations [18, 39], which are common in real-world machine learning algorithms [3]: element-wise multiplication (i.e., Hadamard-product) ($multi_E$), matrix-scalar multiplication ($multi_{MS}$), matrix multiplication ($multi_M$), addition (add_M), division (div_M), transposition (tr), inversion (inv_M), determinant (det), trace ($trace$), diagonal ($diag$), exponential (exp), adjoints (adj), direct sum (sum_D), direct product ($dirct_D$), summation (sum), rows and columns summation ($rowSums$, $colSums$, respectively), QR decomposition (QR), Cholesky decomposition (CHO), LU decomposition (LU), and pivoted LU decomposition (LUP).

6.2 VREM Schema and Relational Encoding

To model LA operations on the \mathcal{VRM} relational schema (part of which appears in Table 1), we also rely on a set of integrity constraints \mathbf{MMC} , which are encoded using relations in \mathcal{VRM} . We detail the encoding below.

6.2.1 Base Matrices and Dimensionality Modeling. We denote by $M_{k \times z}(\mathcal{D})$ a matrix of k rows and z columns, whose entries (values) come from a domain \mathcal{D} , e.g., the domain of real numbers \mathbb{R} . For brevity we just use $M_{k \times z}$. We define a virtual relation $name(M, n) \in \mathcal{VRM}$ attaching a unique ID M to any matrix identified by a name denoted n (which may be e.g. of the form `"/data/M.csv"`). This relation (shown at the top left in Table 1) is accompanied by an EGD key constraint $I_{name} \in \mathbf{MMC}_m$, where $\mathbf{MMC}_m \subset \mathbf{MMC}$, and I_{name} states that *two matrices with the same name n have the same ID*:

$$I_{name}: \forall M \forall N \ name(M, n) \wedge name(N, n) \rightarrow M = N$$

Note that the matrix ID in $name$ (and all the other virtual relations used in our encoding) are not IDs of individual matrix objects: rather, each identifies an *equivalence class* (induced by value equality) of expressions. That is, two expressions are assigned the same ID iff they yield value-based-equal matrices. In Table 1, we use M and N to denote an operation's input matrices' IDs and R for the resulting matrix ID, and s for scalar's input and output.

The dimensions of a matrix are captured by a $size(M, k, z)$ relation, where k and z are the number of rows, resp. columns and M is an ID. An EGD constraint $I_{size} \in \mathbf{MMC}_m$ holds on the $size$ relation, stating that the ID determines the dimensions:

$$I_{size}: \forall M \forall k_1 \forall z_1 \forall k_2 \forall z_2 \\ size(M, k_1, z_1) \wedge size(M, k_2, z_2) \rightarrow k_1 = k_2 \wedge z_1 = z_2$$

The identity I and zero O matrices are captured by the $\text{Zero}(O)$ and $\text{Identity}(I)$ relations, where O and I denote their IDs, respectively. They are accompanied by EGD constraints $\mathcal{I}_{iden}, \mathcal{I}_{zero} \in \mathbf{MMC}_m$, stating that zero matrices with the same sizes have the same IDs, and this also applies for identity matrices with the same size:

$$\mathcal{I}_{zero}: \forall O_1 \forall O_2 \forall k \forall z$$

$$\text{Zero}(O_1) \wedge \text{size}(O_1, k, z) \wedge \text{Zero}(O_2) \wedge \text{size}(O_2, k, z) \rightarrow O_1 = O_2$$

$$\mathcal{I}_{iden}: \forall I_1 \forall I_2$$

$$\text{Identity}(I_1) \wedge \text{size}(I_1, k, k) \wedge \text{Identity}(I_2) \wedge \text{size}(I_2, k, k) \rightarrow I_1 = I_2$$

6.2.2 Encoding Matrix Algebra Expressions. LA operations are encoded into dedicated relations, as shown in Table 1. We now illustrate the encoding of an LA expression on the \mathbf{VRM} schema.

EXAMPLE 6.1. Consider the LA expression $E: ((MN)^T)$, where the two matrices $M_{100 \times 1}$ and $N_{1 \times 10}$ are stored as “M.csv” and “N.csv”, respectively. The encoding function $\text{enc}_{LA}(E)$ takes as argument the LA expression E and returns a conjunctive query whose: (i) body is the relational encoding of E using \mathbf{VRM} (see below), and (ii) head has one distinguished variable, denoting the equivalence class of the result. For instance:

```
enc(((MN)^T) =
  Let enc(MN) =
    Let enc(M) = Q0(M):- name(M, "M.csv");
    enc(N) = Q1(N):- name(N, "N.csv");
    R1 = freshId()
    in
    Q2(R1):- multi_M(M, N, R1) ∧ Q0(M) ∧ Q1(N);
    R2 = freshId()
    in
    Q(R2):- tr(R1, R2) ∧ Q2(R1);
```

In the above, nesting is dictated by the syntax of E . From the inner (most indented) to the outer, we first encode M and N as small queries using the name relation, then their product (to whom we assign the newly created identifier R_1), using the multi_M relation and encoding the relationship between this product and its inputs in the definition of $Q_2(R_1)$. Next, we create a fresh ID R_2 used to encode the full E (the transposed of Q_2) via relation tr , in the query $Q(R_2)$. For brevity, we omit the matrices’ size relations in this example and hereafter. Unfolding $Q_2(R_1)$ in the body of Q yields:

```
Q(R2) :- tr(R1, R2) ∧ multi_M(M, N, R1) ∧
          Q0(M) ∧ Q1(N);
```

Now, by unfolding Q_0 and Q_1 in Q , we obtain the final encoding of $((MN)^T)$ as a conjunctive query Q :

```
Q(R2) :- tr(R1, R2) ∧ multi_M(M, N, R1) ∧
          name(M, "M.csv") ∧ name(N, "N.csv");
```

6.2.3 Encoding LA Properties as Integrity Constraints. Figure 2 shows some of the constraints $\mathbf{MMC}_{LA_{prop}} \subset \mathbf{MMC}$, which capture textbook LA properties [18, 39] of our LA operations (§6.1). The TGDs (1), (2) and (3) state that matrix addition is commutative, matrix transposition is distributive with respect to addition, and the transposition of the inverse of matrix M is equivalent to the inverse of the transposition of M , respectively. We also express that the *virtual relations are functional* by using EGD key constraints. For example, the following $\mathcal{I}_{multi_M} \in \mathbf{MMC}_{LA_{prop}}$ constraint

$$\forall M \forall N \forall R \text{ add}_M(M, N, R) \rightarrow \text{add}_M(N, M, R) \quad (1)$$

$$\forall M \forall N \forall R_1 \forall R_2 \text{ add}_M(M, N, R_1) \wedge \text{tr}(R_1, R_2) \rightarrow$$

$$\exists R_3 \exists R_4 \text{ tr}(M, R_3) \wedge \text{tr}(N, R_4) \wedge \text{add}_M(R_3, R_4, R_2) \quad (2)$$

$$\forall M \forall R_1 \forall R_2 \text{ inv}_M(M, R_1) \wedge \text{tr}(R_1, R_2) \rightarrow$$

$$\exists R_3 \text{ tr}(M, R_3) \wedge \text{inv}_M(R_3, R_2) \quad (3)$$

Figure 2: MMC Constraints Capturing Basic LA Properties

$$\begin{aligned} & \forall M \forall N \forall R_1 \forall R_2 \forall R_3 \forall R_4 \\ & \text{name}(M, "M.csv") \wedge \text{name}(N, "N.csv") \wedge \text{tr}(N, R_1) \wedge \\ & \text{tr}(M, R_2) \wedge \text{inv}_M(R_2, R_3) \wedge \\ & \text{add}_M(R_1, R_3, R_4) \rightarrow \text{name}(R_4, "V.csv") \end{aligned}$$

Figure 3: Relational Encoding of view V

states that multi_M is functional, that is the products of pairwise equal matrices are equal.

$$\begin{aligned} \mathcal{I}_{multi_M} : & \forall M \forall N \forall R_1 \forall R_2 \\ & \text{multi}_M(M, N, R_1) \wedge \text{multi}_M(M, N, R_2) \rightarrow R_1 = R_2 \end{aligned}$$

Other properties [18, 39] of the LA operations we consider are similarly encoded; due to space constraints, we delegate them to Appendix A [11].

6.2.4 Encoding LA Views as Constraints. We translate each view definition V (defined in LA language such as R, DML, etc) into relational constraints $\text{enc}_{LA}(V) \in \mathcal{C}_V$, where \mathcal{C}_V is the set of relational constraints used to capture the views \mathcal{V} . These constraints show how the view’s inputs are related to its output over the \mathbf{VRM} schema. Figure 3 illustrates the encoding as a TGD constraint of the view $V: (N)^T + (M^T)^{-1}$ stored in a file “ $V.csv$ ” and computed based on the matrices N and M (e.g., stored as “ $N.csv$ ” and “ $M.csv$ ”, respectively).

6.2.5 Encoding Matrix Decompositions. Matrix decompositions play a crucial role in many LA computations. For instance, for every symmetric positive definite matrix M there exists a unique Cholesky Decomposition (CD) of the form $M = LL^T$, where L is a lower triangular matrix. We model CD, as well as other well-known decompositions (LU, QR, and Pivoted LU or PLU) as a set of virtual relations \mathbf{VRM}_{dec} , which we add to \mathbf{VRM} . For instance, to CD we associate a relation $\text{CHO}(M, L)$, which denotes that L is the output of the CD decomposition for a given matrix M whose ID is M . CHO is a functional relation, meaning every symmetric positive definite matrix has a unique CD decomposition. This functional aspect is captured by an EGD, conceptually similar to the constraint \mathcal{I}_{multi_M} (§6.2.3). The property $M = LL^T$ is captured as a TGD constraint $\mathcal{I}_{cho} \in \mathbf{MMC}_{LA_{prop}}$:

$$\begin{aligned} \mathcal{I}_{cho} : & \forall M \text{ type}(M, "S") \rightarrow \exists L_1 \exists L_2 \text{ cho}(M, L_1) \wedge \\ & \text{type}(L_1, "L") \wedge \text{tr}(L_1, L_2) \wedge \text{multi}_M(L_1, L_2, M) \end{aligned} \quad (4)$$

The atom $\text{type}(M, "S")$ indicates the type of matrix M , where the constant “ S ” denotes a matrix that is *symmetric positive definite*; similarly, $\text{type}(L_1, "L")$ denotes that the matrix L_1 is a lower triangular matrix. For each base matrix, its type (if available) (e.g., symmetric, upper triangular, etc.) is specified as TGD constraint. For example, we state that a certain matrix M (and any other matrix value-equal to M) is symmetric positive definite as follows:

$$\forall M \text{ name}(M, "M.csv") \rightarrow \text{type}(M, "S") \quad (5)$$

EXAMPLE 6.2. Consider a view $V=N + LL^T$, where $L = \text{cho}(M)$ and M is a symmetric positive definite matrix encoded as in (5). Let E be the LA expression $M + N$. The reader realizes easily that V can be used to answer E directly, thanks to the specific property of the CD decomposition (4), and since $M + N = N + M$, which is encoded in (1). However, at the syntactic level, V and E are very dissimilar. Knowledge of (1) and (4) and the ability to reason about them is crucial in order to efficiently answer E based on V .

The output matrix of CD decomposition is a lower triangular matrix L , which is not symmetric positive definite, meaning that CD decomposition can not be applied again on L . For other decompositions, such as $QR(M) = [Q, R]$ decomposition, where M is a square matrix, Q is an orthogonal matrix [39] and R is an upper triangular matrix, there exists a QR decomposition for the orthogonal matrix Q such that $QR(Q) = [Q, I]$, where I is an identity matrix and $QR(R) = [I, R]$. We say the *fixed point* of the QR decomposition is $QR(I) = [I, I]$. These properties of the Q decompositions are captured with the following constraints, which are part of \mathbf{MMC}_{LAprop} :

$$\begin{aligned} & \forall M \forall n \forall k \text{ name}(M, n) \wedge \text{size}(M, k, k) \rightarrow \exists Q \exists R \\ & QR(M, Q, R) \wedge \text{type}(Q, "O") \wedge \text{type}(R, "U") \\ & \wedge \text{multi}_M(Q, R, M) \end{aligned} \quad (6)$$

$$\begin{aligned} & \forall Q_1 \text{ type}(Q_1, "O") \rightarrow \exists I \ QR(Q_1, Q_1, I) \wedge \text{identity}(I) \\ & \wedge \text{multi}_M(Q_1, I, Q_1) \end{aligned} \quad (7)$$

$$\begin{aligned} & \forall R_1 \text{ type}(R_1, "U") \rightarrow \exists I \ QR(R_1, I, R_1) \wedge \text{identity}(I) \\ & \wedge \text{multi}_M(I, R_1, R_1) \end{aligned} \quad (8)$$

$$\forall I \text{ identity}(I) \rightarrow QR(I, I, I) \quad (9)$$

Known LA properties of the other matrix decompositions (LU and PLU) are similarly encoded in Appendix A.

6.2.6 Encoding LA-Oriented System Rewrite Rules. Most LA-oriented systems [7, 8] execute an incoming expression (LA pipeline) *as-is*, that is: run operations in a sequence, whose order is dictated by the expression syntax. Such systems do not exploit basic LA properties, e.g., reordering a chain of multiplied matrices in order to reduce the intermediate size. SystemML [21] is the only system that models *some* LA properties as static rewrite rules. It also comprises a set of *rewrite rules* which modify the given expressions to avoid large intermediates for aggregation and statistical operations such as $\text{rowSums}(M)$, $\text{sum}(M)$, etc. For example, SystemML uses rule:

$$\text{sum}(MN) = \text{sum}(\text{colSums}(M)^T \odot \text{rowSums}(N)) \quad (i)$$

to rewrite $\text{sum}(MN)$ (summing all cells in the matrix product) where \odot is a matrix element-wise multiplication, to avoid actually computing MN and materializing it; similarly, it rewrites $\text{sum}(M^T)$ into $\text{sum}(M)$, to avoid materializing M^T , etc. However, the performance benefits of rewriting depend on the rewriting power (or, in other words, on *how much the system understands the semantics of the incoming expression*), as the following example shows.

$$\begin{array}{ll} RW_1 : (M^{-1})^T + N^T & RW_2 : (M^T)^{-1} + N^T \\ RW_3 : N^T + (M^{-1})^T & RW_4 : (N^T)^{-1} + N^T \\ RW_5 : (N + M^{-1})^T & \end{array}$$

Figure 4: Equivalent rewritings of the pipeline Q_p .

EXAMPLE 6.3. Consider the LA expression $E=((M^T)^k(M+N)^T)$, where M and N are square matrixes, and expression $E'=\text{sum}(E)$, which computes the sum of all cells in E . The expression E' can be rewritten to $RW_1 : \text{sum}(E'')$, where E'' is:

$$\text{sum}(\text{colSums}(M+N)^T \odot \text{rowSums}(M^k))$$

Failure to exploit the properties $LA_{prop1} : M^T N^T = (MN)^T$ and $LA_{prop2} : (M^n)^T = (M^T)^n$ prevents from finding rewriting RW_1 . E' admits the alternative rewriting

$RW_2 : \text{sum}((\text{colSums}((M^T)^k))^T \odot (\text{colSums}(M+N)^T))$ which can be obtained by directly applying the rewrite rule (i) above and $\text{rowSums}(M^T) = \text{colSums}(M)^T$, without exploiting the properties LA_{prop1} and LA_{prop2} . However, RW_2 introduces more intermediate results than RW_1 .

To fully exploit the potential of rewrite rules (for statistical or aggregation operations), they should be accompanied by sufficient knowledge of, and reasoning on, known properties of LA operations.

To bring such fruitful optimization to other LA-oriented systems lacking support of such rewrite rules, we have incorporated SystemML's rewrite rules into our framework, encoding them as a set of integrity constraints over the virtual relations in the schema \mathcal{VRM} , denoted $\mathbf{MMC}_{StatAgg}$ ($\mathbf{MMC}_{StatAgg} \subset \mathbf{MMC}$). Thus, these rewrite rules can be exploited together with other LA properties. For instance, the rewrite rule (i) is modeled by the following integrity constraint $I_{sum} \in \mathbf{MMC}_{StatAgg}$:

$$\begin{aligned} & \forall M \forall N \forall R \text{ multi}_M(M, N, R) \wedge \text{sum}(R, s) \rightarrow \\ & \exists R_1 \exists R_2 \exists R_3 \exists R_4 \text{ colSums}(M, R_1) \wedge \text{tr}(R_1, R_2) \\ & \wedge \text{rowSums}(N, R_3) \wedge \text{multi}_E(R_2, R_3, R_4) \wedge \text{sum}(R_4, s) \end{aligned}$$

We refer the reader to Appendix B for a full list of SystemML's encoded rewrite rules.

6.3 Relational Rewriting Using Constraints

With the set of views constraints \mathcal{CV} and $\mathbf{MMC} = \mathbf{MMC}_m \cup \mathbf{MMC}_{LAprop} \cup \mathbf{MMC}_{StatAgg}$, we rely on $PACB^{++}$ to rewrite a given expression under integrity constraints. We exemplify this below, and detail $PACB^{++}$'s inner workings in Section 7.

The view V shown in Figure 3 can be used to *fully* rewrite (return the answer for) the pipeline $Q_p : (M^{-1} + N)^T$ by exploiting the TGDs (1), (2) and (3) listed in Figure 2, which describe the following three LA properties, denoted $LA_{prop_i} : M+N = M+N ; ((M+N))^T = (M)^T + (N)^T$ and $((M)^{-1})^T = ((M)^T)^{-1}$. The relational rewriting RW_0 of Q_p using the view V is $RW_0(R_4) :- \text{name}(R_4, "V.csv")$. In this example, RW_0 is the only *views-based* rewriting of Q_p . However, five other rewritings exist (shown in Figure 4), which reorder its operations just by exploiting the set LA_{prop_i} of LA properties.

Rewritings RW_0 to RW_5 have different evaluation costs. We discuss next how we estimate which among these alternatives (including evaluating Q_p directly) is likely the most efficient.

7 CHOICE OF AN EFFICIENT REWRITING

We introduce our cost model (§7.1), which can take two different sparsity estimators (§7.2). Then, we detail our extension to the PACB rewriting engine based on the *Prune_{prov}* algorithm (§7.3) to prune out inefficient rewritings.

7.1 Cost Model

We estimate the cost of an expression E , denoted $\gamma(E)$, as the sum of the intermediate result sizes if one evaluates E “as stated”, in the syntactic order dictated by the expression. Real-world matrices may be *dense* (most or all elements are non-zero) or *sparse* (a majority of zero elements). The latter admit more economical representations that do not store zero elements, which our intermediate result size measure excludes. To estimate the number of non-zeros (*nnz*, in short), we incorporated two different sparsity estimators from the literature (discussed in §7.2) into our framework.

EXAMPLE 7.1. Consider $E_1 = (MN)M$ and $E_2 = M(NM)$, where we assume the matrices $M_{50K \times 100}$ and $N_{100 \times 50K}$ are dense. The total cost of E_1 is $\gamma(E_1) = 50K \times 50K$ and $\gamma(E_2) = 100 \times 100$.

7.2 LA-based Sparsity Estimators

We outline below two existing *sparsity estimators* [21, 49] that we have incorporated into our framework to estimate *nnz*.²

7.2.1 Naïve Metadata Estimator. The naïve metadata estimator [22, 49] derives the sparsity of the output of LA expression solely from the base matrices’ sparsity. This incurs no runtime overhead since metadata about the base matrices, including the *nnz*, columns and rows are available before runtime in a specific metadata file. The most common estimator is the *worst-case estimator* [22], which we use in our framework.

7.2.2 Matrix Non-zero Count (MNC) Estimator. The MNC estimator [46] exploits matrix structural properties such as single non-zero per row, or columns with varying sparsity, for efficient, accurate, and general sparsity estimation; it relies on count-based histograms that exploit these properties. We have also adopted this framework into our approach, and compute histograms about the *base* matrices offline. However, the MNC framework still needs to derive and construct histograms for *intermediate results* online (during rewriting cost estimation). We study this overhead in (§9).

7.3 Rewriting Pruning: PACB⁺⁺

We extended the PACB rewriting engine with the *Prune_{prov}* algorithm sketched and discussed in [31, 32], to eliminate inefficient rewritings during the rewriting search phase. The naïve PACB algorithm generates all minimal (by join count) rewritings before choosing a *minimum-cost* one. While this suffices on the scenarios considered in [15, 32], the settings we obtain from our LA encoding stress-test the naïve algorithm, as commutativity, associativity, etc. blow up the space of alternate rewritings exponentially. Scalability considerations forced us to further optimize naïve PACB to find only *minimum-cost* rewritings, aggressively pruning the others during the generation phase. We illustrate *Prune_{prov}* and our improvements next.

Prune_{prov} Minimum-Cost Rewriting. Recall from (§4) that the minimal rewritings of a query Q are obtained by first finding the set \mathcal{H} of all matches (i.e., containment mappings) from Q to the

result B of backchasing the universal plan U . Denoting with $\pi(S)$ the provenance formula of a set of atoms S , PACB computes the DNF form D of $\bigvee_{h \in \mathcal{H}} \pi(h(Q))$. Each conjunct c of D determines a subquery $sq(c)$ of U which is guaranteed to be a rewriting of Q .

The idea behind cost-based pruning is that, whenever the naïve PACB backchase would add a provenance conjunct c to an existing atom a ’s provenance formula $\pi(a)$, *Prune_{prov}* does so more conservatively: if the cost $\gamma(sq(c))$ is larger than the minimum cost threshold T found so far, then c will never participate in a minimum-cost rewriting and need not be added to $\pi(a)$. Moreover, atom a itself need not be chased into B in the first place if all its provenance conjuncts have above-threshold cost.

EXAMPLE 7.2. Let $E = M(NM)$, where we assume for simplicity that $M_{50K \times 100}$ and $N_{100 \times 50K}$ are dense. Exploiting the associativity of matrix-multiplication $(MN)M = M(NM)$ during the chase leads to the following universal plan U annotated with provenance terms:

$$\begin{aligned} U(R_2) : & name(M, "M.csv")^{p_0} \wedge size(M, 50000, 100)^{p_1} \wedge \\ & name(N, "N.csv")^{p_2} \wedge size(N, 100, 50000)^{p_3} \wedge \\ & multi_M(M, N, R_1)^{p_4} \wedge multi_M(R_1, M, R_2)^{p_5} \wedge \\ & multi_M(N, M, R_3)^{p_6} \wedge multi_M(R_3, R_2)^{p_7} \end{aligned}$$

Now, consider in the back-chase the associativity constraint C :

$$\begin{aligned} & \forall M \forall N \forall R_1 \forall R_2 \\ & multi_M(M, N, R_1) \wedge multi_M(R_1, M, R_2) \rightarrow \\ & \exists R_4 multi_M(N, M, R_4) \wedge multi_M(M, R_4, R_2) \end{aligned}$$

There exists a containment mapping h embedding the two atoms in the premise P of C into the U atoms whose provenance annotations are p_4 and p_5 . The provenance conjunct collected from P ’s image is $\pi(h(P)) = p_4 \wedge p_5$.

Without pruning, the backchase would chase U with the constraint C , yielding U' which features additional $\pi(h(P))$ -annotated atoms $multi_M(N, M, R_4)^{p_4 \wedge p_5} \wedge multi_M(M, R_4, R_2)^{p_4 \wedge p_5}$

E has precisely two matches h_1, h_2 into U' . $h_1(E)$ involves the newly added atoms as well as those annotated with p_0, p_1, p_2, p_3 . Collecting all their provenance annotations yields the conjunct $c_1 = p_0 \wedge p_1 \wedge p_2 \wedge p_3 \wedge p_4 \wedge p_5$. c_1 determines the U -subquery $sq(c_1)$ corresponding to the rewriting $(MN)M$, of cost $(50K)^2$.

$h_2(E)$ ’s image yields the provenance conjunct $c_2 = p_0 \wedge p_1 \wedge p_2 \wedge p_3 \wedge p_6 \wedge p_7$, which determines the rewriting $M(NM)$ that happens to be the original expression E of cost 100^2 .

The naïve PACB would find both rewritings, cost them, and drop the former in favor of the latter.

With pruning, the threshold T is the cost of the original expression 100^2 . The chase step with C is never applied, as it would introduce the provenance conjunct $\pi(h(P))$ which determines U -subquery $sq(\pi(h(P))) = multi_M(M, N, R_1)^{p_4} \wedge multi_M(R_1, M, R_2)^{p_5}$

of cost $(50K)^2$ exceeding T . The atoms needed as image of E under h_1 are thus never produced while backchasing U , so the expensive rewriting is never discovered. This leaves only the match image $h_2(E)$, which corresponds to the efficient rewriting $M(NM)$.

Our improvements on Prune_{prov}. Whenever the pruned chase step is applicable and applied for each TGD constraint, the original algorithm searches for all *minimal-rewritings* \mathcal{RW} that can be found “so far”, then it costs each $rw \in \mathcal{RW}$ to find the “so far” *minimum-cost* one rwe and adjusts the threshold T to the cost of rwe .

²Solving the problem of sparsity estimation is beyond the scope of this paper

However, this strategy can cause *redundant* costing of $rw \in \mathcal{R}^*W$ whenever the pruned chase step is applied again for another constraint. Therefore, in our modified version of *Prune_{prov}*, we keep track of the rewriting costs already estimated, to prevent such redundant work. Additionally, the search for *minimal-rewritings* “so far” (matches of the query Q into the evolving universal plan instance U' , see §4) whenever the pruned chase step is applied is modeled as a query evaluation of Q against U' (viewed as a symbolic/canonical database [14]). This involves repeatedly evaluating the same query plan. However, the query is evaluated over evolutions of the same instance. Each pruned chase step *adds a few new tuples* to the evolving instance, corresponding to atoms introduced by the step, while most of the instance is unchanged. Therefore, instead of evaluating the query plan from scratch, we employ incremental evaluation as in [32]. The plan is kept in memory along with the populated hash tables, and whenever new tuples are added to the evolving instance, we push them to the plan.

8 GUARANTEES ON THE REDUCTION

We detail the conditions under which we guarantee that our approach is *sound* (i.e., generates only equivalent, cost-optimal rewritings), and *complete* (i.e., finds all equivalent cost-optimal rewritings).

Let $\mathcal{V} \subseteq \mathcal{L}$ be a set of materialized view definitions, where \mathcal{L} is the language of hybrid expressions described in §3.

Let LA_{prop} be a set of properties of the LA operations in L_{ops} that admits relational encoding over \mathcal{VRM} . We say that LA_{prop} is *terminating* if it corresponds to a set of TGDs and EGDs with terminating chase (this holds for our choice of LA_{prop}).

Denote with γ a cost model for expressions from \mathcal{L} . We say that γ is *monotonic* if expressions are never assigned a lower cost than their subexpressions (this is true for both models we used).

We call $E \in \mathcal{L}$ $(\gamma, LA_{prop}, \mathcal{V})$ -optimal if for every $E' \in \mathcal{L}$ that is (LA_{prop}, \mathcal{V}) -equivalent to E we have $\gamma(E') \geq \gamma(E)$.

Let $Eq^\gamma(LA_{prop}, \mathcal{V})(E)$ denote the set of all $(\gamma, LA_{prop}, \mathcal{V})$ -optimal expressions that are (LA_{prop}, \mathcal{V}) -equivalent to E .

We denote with $HADAD(LA_{prop}, \mathcal{V}, \gamma)$ our parameterized solution based on relational encoding followed by PACB++ rewriting and next by decoding all the relational rewritings generated by the cost-based pruning PACB++ (recall Figure 1). Given $E \in \mathcal{L}$, $HADAD(LA_{prop}, \mathcal{V}, \gamma)(E)$ denotes all expressions returned by $HADAD(LA_{prop}, \mathcal{V}, \gamma)$ on input E .

THEOREM 8.1 (SOUNDNESS). *If the cost model γ is monotonic, then for every $E \in \mathcal{L}$ and every $rw \in HADAD(LA_{prop}, \mathcal{V}, \gamma)(E)$, we have $rw \in Eq^\gamma(LA_{prop}, \mathcal{V})(E)$.*

THEOREM 8.2 (COMPLETENESS). *If γ is monotonic and LA_{prop} is terminating, then for every $E \in \mathcal{L}$ and every $rw \in Eq^\gamma(LA_{prop}, \mathcal{V})(E)$, we have $rw \in HADAD(LA_{prop}, \mathcal{V}, \gamma)(E)$.*

9 EXPERIMENTAL EVALUATION

We evaluate HADAD to answer research questions below about our approach:

- **§9.1.1, §9.1.2, §9.2.1 and §9.2.2:** Can HADAD find rewrites with/without views that lead to a greater performance improvement than original pipelines without modifying the internals of the existing systems? Are the identified rewrites found by state-of-the-art platforms?

- **§9.1.3 and §9.2.3:** Is HADAD’s optimization overhead compensated by the performance gains in execution?

We evaluate our approach, first on **LA-based pipelines** (§9.1), then on **hybrid expressions** (§9.2). Due to space constraints, we only discuss a subset of our results here and delegate other results to Appendix C, D, E, and F.

Experimental Environment. We used a single node with an Intel(R) Xeon(R) CPU E5-2640 v4 @ 2.40GHz, 40 Cores (hyper-threading), 123GB RAM, disk read speed 616 MB/s, and disk write speed 455 MB/s. We run on OpenJDK Java 8 VM. As for **LA systems/libraries**, we used **R 3.6.0**, **Numpy 1.16.6** (**python 2.7**), **TensorFlow 1.4.1**, **Spark 2.4.5 (MLlib)**, **SystemML 1.2.0** For hybrid experiments, we use **Morpheus** [4] and **SparkSQL** [17].

Systems Configuration Tuning. We discuss here the most important installation and configuration details. We use a JVM-based linear algebra library for SystemML as recommended in [48], at the optimization level 4. Additionally, we enable multi-threaded matrix operations in a single node. We run Spark in single node setting and OpenBLAS (built from the sources as detailed in [5]) to take advantage of its accelerations [48]. SparkMLlib’s datatypes do not support many basic LA operations, such as scalar-matrix multiplication, matrix element-wise multiplication, etc. To support them, we use the Breeze Scala library [2], convert MLlib’s datatypes to Breeze types and express the basic LA operations in Spark. The driver memory allocated for Spark and SystemML is 115GB. To maximize TensorFlow performance, we compile it from sources. For all systems/libraries, we set the number of **cores** to 24; all systems use **double precision** numbers.

9.1 LA-based Experiments

In this experiments, we study the performance benefits of our approach on LA-based pipelines as well as our optimization overhead.

Datasets. We used several **real-world, sparse matrices**, for which Table 4 lists the dimensions and the sparsity (S_X) (i) dielFilterV3real (**DFV** in short) is an analysis of a microwave filter with different mesh qualities [28]; (ii) 2D_54019_highK (**2D_54019** in short) is a 2D semiconductor device simulation [28]; (iii) we used several subsets of an Amazon books review dataset [1] (in JSON), and similarly (iv) subsets of a Netflix movie rating dataset [6]. The latter two were easily converted into matrices where columns are items and rows are customers [49]; we extracted smaller subsets of all real datasets to ensure the various computations applied on them fit in memory (e.g., Amazon/(AS) denotes the small version of the Amazon dataset). We also used a set **synthetic, dense matrices**, described in Table 5.

LA benchmark. We select a set \mathcal{P} of 57 **LA pipelines** used in prior studies and/or frequently occurring in real-world LA computations, as follows:

- **Real-world pipelines** include: a chain of matrix self products used for reachability queries and other graph analytics [46] (**P1.29** in Table 2); expressions used in Alternating Least Square Factorization (ALS) [49] (**P2.25** in Table 14); Poisson Nonnegative Matrix Factorization (PNMF) (**P1.13** in Table 15) [49]; Nonnegative Matrix Factorization (NMF) (**P1.25** and **P1.26** in Table 15) [48]; complex predicate for image

No.	Expression	No.	Expression	No.	Expression
P1.1	$(MN)^T$	P1.2	$A^T + B^T$	P1.3	$C^{-1}D^{-1}$
P1.4	$(A + B)v_1$	P1.5	$((D)^{-1})^{-1}$	P1.6	$\text{trace}(s_1 D)$
P1.7	$((A)^T)^T$	P1.8	$s_1 A + s_2 A$	P1.9	$\det(D^T)$
P1.10	$\text{rowSums}(A^T)$	P1.11	$\text{rowSums}(A^T + B^T)$	P1.12	$\text{colSums}((MN))$
P1.13	$\text{sum}(MN)$	P1.14	$\text{sum}(\text{colSums}((N^T M^T)))$	P1.15	$(MN)M$
P1.16	$\text{sum}(A^T)$	P1.17	$\det(CDC)$	P1.18	$\text{sum}(\text{colSums}((A)))$
P1.19	$(C^T)^{-1}$	P1.20	$\text{trace}(C^{-1})$	P1.21	$(C + D^{-1})^T$
P1.22	$\text{trace}((C + D)^{-1})$	P1.23	$\det((CD)^{-1}) + D$	P1.24	$\text{trace}((CD)^{-1}) + \text{trace}(D)$
P1.25	$M \odot (N^T / (MNN^T))$	P1.26	$N \odot (M^T / (M^T MN))$	P1.27	$\text{trace}(D(CD)^T)$
P1.28	$A \odot (A \odot B + A)$	P1.29	$CCCC$	P1.30	$NM \odot NMR^T$

Table 2: LA Benchmark Pipelines (Part 1)

No.	Expression	No.	Expression	No.	Expression
P2.1	$\text{trace}(C + D)$	P2.2	$\det(D^{-1})$	P2.3	$\text{trace}(D^T)$
P2.4	$s_1 A + s_1 B$	P2.5	$\det((C + D)^{-1})$	P2.6	$C^T(D^T)^{-1}$
P2.7	$DD^{-1}C$	P2.8	$\det(C^T D)$	P2.9	$\text{trace}(C^T D^T + D)$
P2.10	$\text{rowSums}(MN)$	P2.11	$\text{sum}(A + B)$	P2.12	$\text{sum}(\text{rowSums}(N^T M^T))$
P2.13	$((MN)M)^T$	P2.14	$((MN)M)N$	P2.15	$\text{sum}(\text{rowSums}(A))$
P2.16	$\text{trace}(C^{-1}D^{-1}) + \text{trace}(D)$	P2.17	$(((C + D)^{-1})^T)((D^{-1})^{-1})C^{-1}C$	P2.18	$\text{colSums}((A^T + B^T)$
P2.19	$(C^T D)^{-1}$	P2.20	$(M(NM))^T$	P2.21	$(D^T D)^{-1}(D^T v_1)$
P2.22	$\exp((C + D)^T)$	P2.23	$\det(C) * \det(D) * \det(C)$	P2.24	$(D^{-1}C)^T$
P2.25	$(u_1 v_2^T - X)v_2$	P2.26	$\exp((C + D)^{-1})$	P2.27	$(((C + D)^T)^{-1})D)C$

Table 3: LA Benchmark Pipelines (Part 2)

Name	Rows n	Cols m	Nnz $\ X\ _0$	S_X
DFV	1M	100	8050	0.0080%
2D_54019	50K	100	3700	0.0740%
Amazon/(AS)	50K	100	378	0.0075%
Amazon/(AM)	100K	100	673	0.0067%
Amazon/(AL ₁)	1M	100	6539	0.0065%
Amazon/(AL ₂)	10M	100	11897	0.0011%
Amazon/(AL ₃)	100K	50K	103557	0.0020%
Netflix/(NS)	50K	100	69559	1.3911%
Netflix/(NM)	100K	100	139344	1.3934%
Netflix/(NL ₁)	1M	100	665445	0.6654%
Netflix/(NL ₂)	10M	100	665445	0.0665%
Netflix/(NL ₃)	100K	50K	15357418	0.307%

Table 4: Overview of Used Real Datasets.

Name	Rows n	Cols m	Name	Rows n	Cols m
<i>Syn₁</i>	50K	100	<i>Syn₆</i>	20K	20K
<i>Syn₂</i>	100	50K	<i>Syn₇</i>	100	1
<i>Syn₃</i>	1M	100	<i>Syn₈</i>	50K	1
<i>Syn₄</i>	5M	100	<i>Syn₉</i>	100K	1
<i>Syn₅</i>	10K	10K	<i>Syn₁₀</i>	100	100

Table 5: Syntactically Generated Dense Datasets

Matrix Name	Used Data
<i>A</i> and <i>B</i>	AM, AL ₁ , AL ₂ , NM, NL ₁ , NL ₂ , diefFilter, <i>Syn₃</i> or <i>Syn₄</i>
<i>C</i> and <i>D</i>	<i>Syn₅</i> or <i>Syn₆</i>
<i>M</i>	AS, NS, <i>Syn₁</i> , or 2D_54019
<i>N</i>	<i>Syn₂</i>
<i>R</i>	<i>Syn₁₀</i>
<i>X</i>	AL ₃ or NL ₃
<i>v₁, v₂</i> and <i>u₁</i>	<i>Syn₇</i> , <i>Syn₈</i> and <i>Syn₉</i> , respectively.

Table 6: Matrices used for each matrix name in a pipeline

masking [46] (**P1.30** in Table 15); recommendation computation [46] (**P1.30** in Table 15); finally, Ordinary Least Squares Regression (OLS) [48] (**P2.21** in Table 14).

- **Synthetic pipelines** were also generated, based on a set of basic matrix operations (inverse, multiplication, addition, etc.), and a set of combination templates, written as a Rule-Iterated Context-Free Grammar (RI-CFG) [43]. Expressions thus generated include **P2.16**, **P2.16**, **P2.23**, **P2.24** in Table 15.

Methodology. In (§9.1.1), we show the performance benefits of our approach to LA-oriented systems/tools mentioned above using a set $\mathcal{P}^{-Opt} \subset \mathcal{P}$ of 38 pipelines in Table 15 and Table 14, and the matrices in Table 6. The performance of these pipelines can be improved *just by exploiting LA properties* (in the absence of views). For TensorFlow and NumPy, we present the results only for dense matrices, due to limited support for sparse matrices. In (§9.1.2), we show how our approach improves the performance of 30 pipelines from \mathcal{P} , denoted \mathcal{P}^{Views} , using pre-materialized views. Finally, in (§9.1.3), we study our rewriting performance and *optimization overhead* for the set $\mathcal{P}^{Opt} = \mathcal{P} \setminus \mathcal{P}^{-Opt}$ of 19 pipelines that are already optimized.

9.1.1 Effectiveness of LA Rewriting (No Views). For each system, we run the original pipeline and our rewriting 5 times; we report the average of the last 4 running times. We exclude the data loading time. For fairness, we ensured SparkMLlib and SystemML compute the entire pipeline (despite their lazy evaluation mode).

Figure 5 illustrates the original pipeline execution time Q_{exec} and the selected rewriting execution time RW_{exec} for **P1.1**, **P1.3**, **P1.4**, and **P1.15**, including the rewriting time RW_{find} , using the MNC cost model. For each pipeline, the used datasets are on top of the figure. For brevity in the figures, we use **SM** for SystemML, **NP** for NumPy, **TF** for Tensorflow, and **SP** for MLlib. For **P1.1** (see Figure 5(a)), both matrices are dense. The speed-up (1.5 \times to 4 \times) comes from rewriting $(MN)^T$ (intermediate result size to $(50K)^2$) into $N^T M^T$, much cheaper since both N^T and M^T are of size $50K \times 100$.

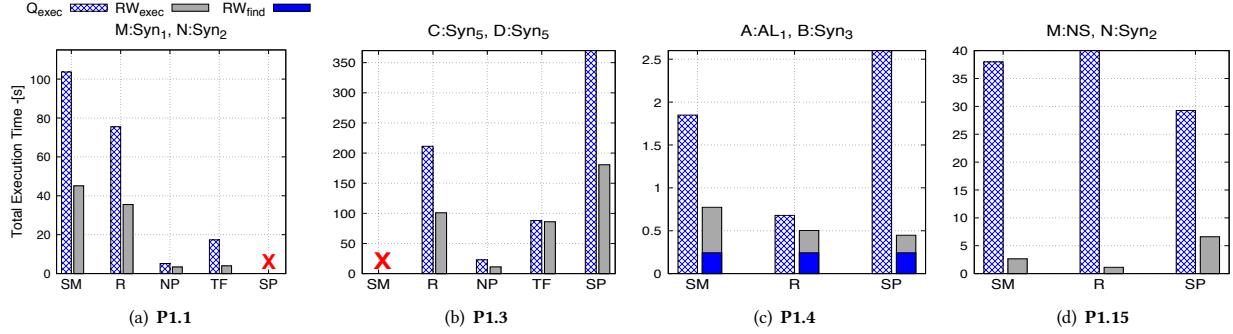


Figure 5: P1.1, P1.3, P1.4, and P1.15 evaluation before and after rewrite

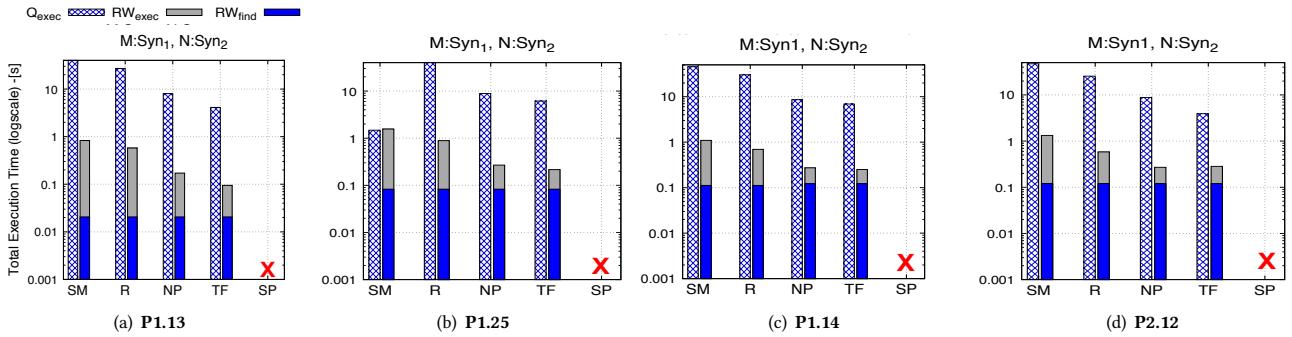


Figure 6: P1.13, P1.25, P1.14 and P2.12 evaluation before and after rewrite

We exclude MLlib from this experiment since it failed to allocate memory for the intermediate matrix (Spark/MLlib limits the maximum size of a dense matrix). As a variation (not plotted in the Figure), we ran the same pipeline with the ultra-sparse *AS* matrix (0.0075% non-zeros) used as *M*. The Q_{exec} and RW_{exec} time are very comparable using SystemML, because we avoid large dense intermediates. In R, this scenario lead to a *c* since the multiplication operator tries to densify the matrix *M*. To avoid it, we cast *M* during load time to a dense matrix type. Thus, the speed-up achieved is the same as if *M* and *N* were both dense. If, instead, *NS* (1.3860% non-zeros) plays the role of *M*, our rewrite achieves $\approx 1.8\times$ speed-up for SystemML.

For **P1.3** (Figure 5(b)), the speed-up comes from rewriting $C^{-1}D^{-1}$ to $(DC)^{-1}$. Interestingly, TensorFlow is the only system that applies this optimization by itself. SystemML timed-out (>1000 secs) for both original pipeline and its rewriting.

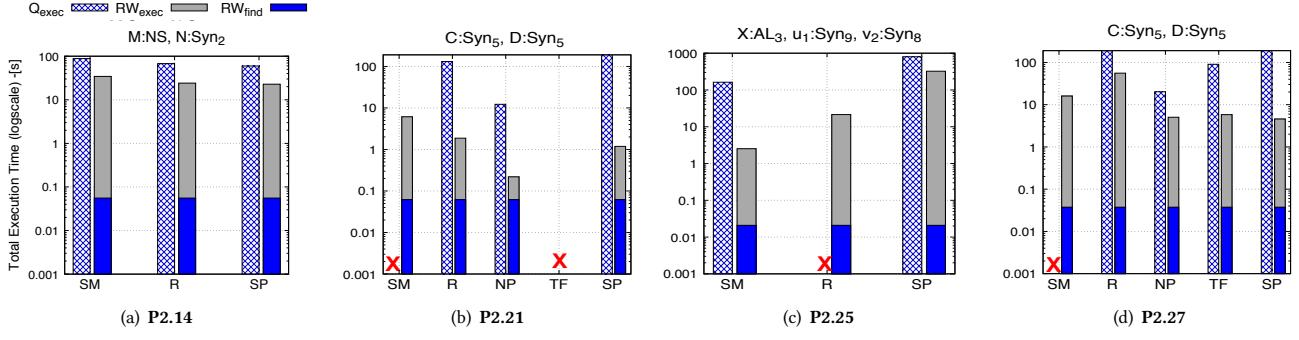
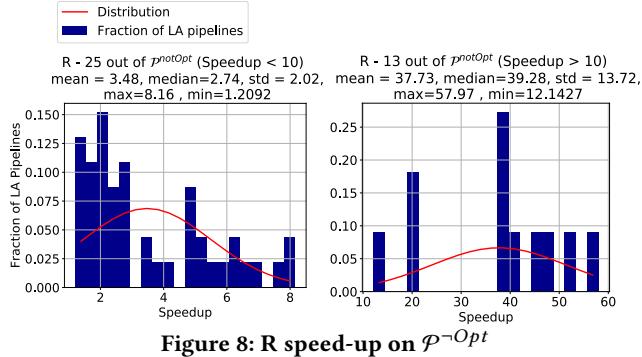
For pipeline **P1.4** (Figure 5(c)), we rewrite $(A + B)v_1$ to $Av_1 + Bv_1$. Adding a sparse matrix *A* to a dense matrix *B* results into materializing a dense intermediate of size $1M \times 100$. Instead, $Av_1 + Bv_1$ has fewer non-zeros in the intermediate results, and Av_1 can be computed efficiently since *A* is sparse. The MNC sparsity estimator has a noticeable overhead here. We run the same pipeline, where the dense *Syn₄* matrix plays both *A* and *B* (not shown in the Figure). This leads to speed-up of up to $9\times$ for MLlib, which does not natively support matrix addition, thus we convert its matrices to Breeze types in order to perform it (as in [48]).

P1.15 (Figure 5(d)) is a matrix chain multiplication. The naïve left-to-right evaluation plan $(MN)M$ computes an intermediate matrix of size $O(n^2)$, where *n* is 50K. Instead, the rewriting $M(NM)$

only needs an $O(m^2)$ intermediate matrix, where *m* is 100, and is much faster. To avoid MLlib memory failure on **P1.15**, we use the distributed matrix of type **BlockMatrix** for both matrices. While *M* thus converted has the same sparsity, Spark views it as being of a dense type (multiplication on BlockMatrix is considered to produce dense matrices) [10]. SystemML does optimize the multiplication order if the user does not enforce it. Further (not shown in the Figure), we ran **P1.15** with *AS* in the role of *M*. This is $4\times$ faster in SystemML since with an ultra sparse *M*, multiplication is more efficient. This is not the case for MLlib which views it as dense. For R, we again had to densify *M* during loading to prevent crashes.

Figures 6(a) and 6(b) study **P1.13** and **P1.25**, two real-world pipelines involved in ML algorithms, using the MNC cost model; note the log-scale *y* axis. Rewriting **P1.13** to $\text{sum}(t(\text{colSums}((M)) * \text{rowSums}(N)))$ yields a speed-up of $50\times$; while SystemML has this rewrite as a static rule, it did not apply it. Our rewrite allowed SystemML and the others to benefit from it. Not shown in the Figure, we re-ran this with *M* ultra sparse (using *AS*) and SystemML: the rewrite did not bring benefits, since *MN* is already efficient. *In this experiment and subsequently, whenever MLlib is absent, this is due to its lack of support for LA operations* (here, sum of all cells in a matrix) on BlockMatrix. For **P1.25**, the important optimization is selecting the multiplication order in MNN^T (Figure 6(b)). SystemML is efficient here, due to its dedicated operator *tsmm* for transpose-self matrix multiplication and *mmchain* for matrix multiply chains.

Figures 6(c) and 6(d) shows up to $42\times$ rewriting speed-up achieved by turning **P1.14** and **P2.12** into $\text{sum}(t(\text{colSums}((M)) * \text{rowSums}(N)))$. This exploits several properties: (i) $(MN)^T = N^T M^T$, (ii) $\text{sum}(M^T) =$

Figure 7: P2.14, P2.21, P2.25 and P2.27 evaluation before and after rewriting using the views V_{exp} Figure 8: R speed-up on P^{-Opt}

$\text{sum}(M)$, (iii) $\text{sum}(\text{row}/\text{colSums}((M))) = \text{sum}(M)$, and (iv) $\text{sum}(MN) = \text{sum}(\text{t}(\text{colSums}((M))) * \text{rowSums}(N))$. SystemML captures (ii), (iii), and (iv) as static rewrite rules, however, it is unable to exploit these performance-saving opportunities since it is unaware of (i). Other systems lack support for more or all of these properties.

Figure 8 shows the distribution of the significant rewriting speed-up on P^{-Opt} running on R, and using the MNC-based cost model. For clarity, we split the distribution into two figures: on the left, 25 P^{-Opt} pipelines with speed-up lower than 10 \times ; on the right, the remaining 13 with greater speed-up. Among the former, 87% achieved at least 1.5 \times speed-up. The latter are sped up by 10 \times to 60 \times . **P1.5** is an extreme case here (not plotted): it is sped up by about 1000 \times , simply by rewriting $((D)^{-1})^{-1}$ into D.

9.1.2 Effectiveness of view-based LA rewriting. We have defined a set V_{exp} of 12 views that pre-compute the result of some expensive operations (multiplication, inverse, determinant, etc.) which can be used to answer our P^{Views} pipelines, and materialized them on disk as CSV files. The experiments outlined below used the naïve cost model; all graphs have a log-scale y axis.

Discussion. For **P2.14** (Figure 7(a)), using the view $V_4 = NM$ by and the multiplication associativity leads to up to 2.8 \times speed-up.

Figure 7(b) shows the gain due to the view $V_1 = D^{-1}$, for the ordinary-least regression (OLS) pipeline **P2.21**. It has 8 rewritings, 4 of which use V_1 ; they are found thanks to the properties $(CD)^{-1} = D^{-1}C^{-1}$, $(CD)E = C(DE)$ and $(D^T)^{-1} = (D^{-1})^T$ among others. The cheapest rewriting is $V(V^T(D^Tv_1))$, since it introduces small intermediates due to the optimal matrix chain multiplication order. This rewrite leads to 70 \times , 55 \times and 150 \times speed-ups on R, NumPy and MLLib, respectively; TensorFlow is omitted as matmul operator does not support matrix-vector multiplication. It could run this

pipeline by converting the matrices to NumPy, whose performance we already report separately. On SystemML, the original pipeline timed out (> 1000 seconds).

Pipeline **P2.25** (Figure 7(c)) benefits from a view V_5 , which pre-computes a dense intermediate vector multiplication result; then, rewriting based on the property $(A + B)v = Av + Bv$ leads to a 65 \times speed-up in SystemML. For MLLib, as discussed before, to avoid memory failure, we used BlockMatrix types, for all matrices and vectors, thus they were treated as dense. In R, the original pipeline triggers a memory allocation failure, which the rewriting avoids. Figure 7(d) shows that for **P2.27** exploiting the views $V_2 = (D+C)^{-1}$ and $V_3 = DC$ leads to speed-ups of 4 \times to 41 \times on different systems. Properties enabling rewriting here are $C + D = D + C$, $(D^T)^{-1} = (D^{-1})^T$ and $(CD)E = C(DE)$.

9.1.3 Rewriting Performance and Overhead. We now study the running time RW_{find} of our rewriting algorithm, and the *rewrite overhead* defined as $RW_{find}/(Q_{exec} + RW_{find})$, where Q_{exec} is the time to run the pipeline “as stated”. We ran each experiment 100 times and report the average of the last 99 times. The global trends are as follows. (i) For a fixed pipeline and set of data matrices, the overhead is slightly higher using the MNC cost model, since histograms are built during optimization. (ii) For a fixed pipeline and cost model, sparse matrices lead to a higher overhead simply because Q_{exec} tends to be smaller. (iii) Some (system, pipeline) pairs lead to a low Q_{exec} when the system applies internally the same optimization that HADAD finds “outside” of the system.

Concretely, for the P^{-Opt} pipelines, on the dense and sparse matrices listed in Table 6, using the naïve cost model, 64% of the RW_{find} times are under 25ms (50% are under 20ms), and the longest is about 200ms. Using the MNC estimator, 55% took less than 20ms, and the longest (outlier) took about 300ms. Among the 39 P^{-Opt} pipelines, SystemML finds efficient rewritings for a set of 6, denoted P_{SM}^{-Opt} , while TensorFlow optimizes a different set of 11, denoted P_{TF}^{-Opt} . On these subsets, where HADAD’s optimization is redundant, using dense matrices, the overhead is very low: with the MNC model, 0.48% to 1.12% on P_{SM}^{-Opt} (0.64% on average), and 0.0051% to 3.51% on P_{TF}^{-Opt} (1.38% on average). Using the naïve estimator slightly reduces this overhead, but across P^{-Opt} , this model misses 4 efficient rewritings. On sparse matrices, the overhead is at most 4.86% with the naïve estimator and up to 5.11% with the MNC one.

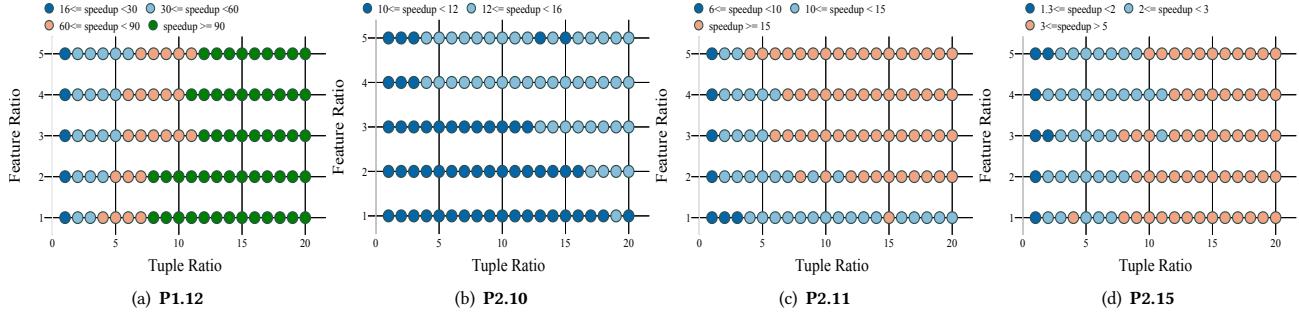


Figure 9: Speed-ups of Morpheus (with HADAD rewrites) over Morpheus (without HADAD rewrites) for pipelines P1.12, P2.10, P2.11 and P2.15 on synthetic data for a PK-FK join.

Among the already-optimal pipelines \mathcal{P}^{Opt} , 70% involve expensive operations such as inverse, determinant, matrix exponential, leading to rather long Q_{exec} times. Thus, the rewriting overhead is less than 1% of the total time, on all systems, using sparse or dense matrices, and the naïve or the MNC-based cost models. For the other \mathcal{P}^{Opt} pipelines with short Q_{exec} , mostly matrix multiplications chains already in the optimal order, on *dense* matrices, the overhead reaches 0.143% (SparkMLlib) to 9.8% (TensorFlow) using the naïve cost model, while the MNC cost model leads to an overhead of 0.45% (SparkMLlib) up to 10.26% (TensorFlow). On *sparse* matrices, using the naïve and MNC cost models, the overhead reaches up to 0.18% (SparkMLlib) to 1.94% (SystemML), and 0.5% (SparkMLlib) to 2.61% (SystemML), respectively.

9.2 Hybrid (LA and RA) Experiments

We now study the benefits of rewriting on hybrid expressions combining relational and linear algebra. In (§9.2.1), we show the performance benefits of our approach to Morpheus [4]. In (§9.2.2), we evaluate our micro-hybrid benchmark on SparkSQL with SystemML. (§9.2.3) discusses our optimization overhead.

9.2.1 Morpheus Experiments. In this experiment, we study the performance-saving opportunity that HADAD adds to Morpheus. Morpheus is a cross RA-LA platform for pushing LA computations through joins, where LA pipelines can be automatically "factorized" using its framework of algebraic rewrite rules. The rules rewrite LA operations over the join output's matrix (table) into LA operations over the base tables' matrices, thereby avoiding data redundancy introduced by the join, which leads to runtime inefficiencies due to redundancy in the computations. For example, (based on [4]), consider a data scientist using ML to predict if a customer will move to a competitor. She joins a table *Customers* with columns (features) such as income, employerID, etc. with another table *Employers* such as revenue, city, etc. The column vectors in the base tables are viewed as matrices, *C* and *E*, respectively. The materialized join output is viewed as a matrix *M*. Factorized LA rewrites LA operation over *M*, e.g., *rowSums(M)* into LA operations over *C* and *E*. Thus, the user writes LA pipeline on one table *M* but under the cover Morpheus rewrites the pipeline to operate on *C* and *E*. However, based on equivalences that hold due to LA properties, LA pipelines' rewrites are not exploited by Morpheus, leading to missed optimization opportunities, as we outlined below.

Methodology. We used the same experiments setup introduced in [4] for generating synthetic datasets for PK-FK join of **R** and **S** tables. The quantities varied are the *tuple ratio* (n_S/n_R) and feature ratio (d_R/d_S), where n_S and n_R are the number of rows and d_R and d_S are the number of columns (features) in **R** and **S**, respectively. We fixed $n_R = 1M$ and $d_S = 20$, this helps quantify the amount of redundancy introduced by a PK-FK join on **R** and **S**. The join's output matrix $\mathbf{M}(n_S \times (d_R + d_S))$ is always dense. We evaluate a set $\mathcal{P}_M^{-opt} \subset \mathcal{P}^{-opt}$ of 8 pipelines and their selected rewritings found by HADAD (using naïve cost-model) on Morpheus.³

Discussion. **P1.9: colSums(MN)** is the example from §2, with **M** the output (viewed as matrix) of joining tables **R** and **S** generated as described above. **N** is a $ncol(\mathbf{M}) \times 100$ dense matrix. HADAD's rewriting yields up to $125\times$ speed-up (see Figure 9(a)).

Figure 9(b) shows up to $15\times$ speed-up for **P1.10: rowSums(NM)**, where the size of **N** is $100 \times nrow(\mathbf{M})$. This is due to HADAD's rewriting: *NrowSums(M)*, which enables Morpheus to push the *rowSums* operator to **R** and **S** instead of computing the large intermediate matrix multiplication.

P1.11: sum(N+M) is run *as-is* by Morpheus since it does not factorize element-wise operations, e.g., addition. However, HADAD rewrites **P1.11** into *sum(N)+sum(M)*, which avoids the (large and dense) intermediate result of the element-wise matrix addition. The HADAD rewriting enables Morpheus to execute *sum(M)* by pushing *sum* to **R** and **S**, for up to $20\times$ speed-up (see Figure 9(c)).

Morpheus evaluates **P1.12: sum(rowSums(M))** by pushing the *rowSums* operator to **R** and **S**. HADAD finds the rewriting *sum(M)*, which enables Morpheus to push the *sum* operation instead, achieving up to $4.5\times$ speed-up (see Figure 9(d)).

Since Morpheus does not exploit the associative property of matrix multiplication, it cannot reorder multiplication chains to avoid large intermediate results, which lead to runtime exception in R (Morpheus' backend). For example, for the chain $N^T NM$, when the size of **M** is $1M \times 20$ and the size of **N** is $20 \times 1M$, the size of the $N^T N$ intermediate result is $1M \times 1M$, which R cannot handle. HADAD exploits associativity and selects the rewriting $N^T(NM)$ of intermediate result size (100×20).

³We select a set of pipelines where Morpheus can rewrite (factorize) them. This includes pipelines that contain aggregate operations, transpose and matrix multiplication. For the inverse, we only evaluated P2.21 (OLS), which is the realistic one.

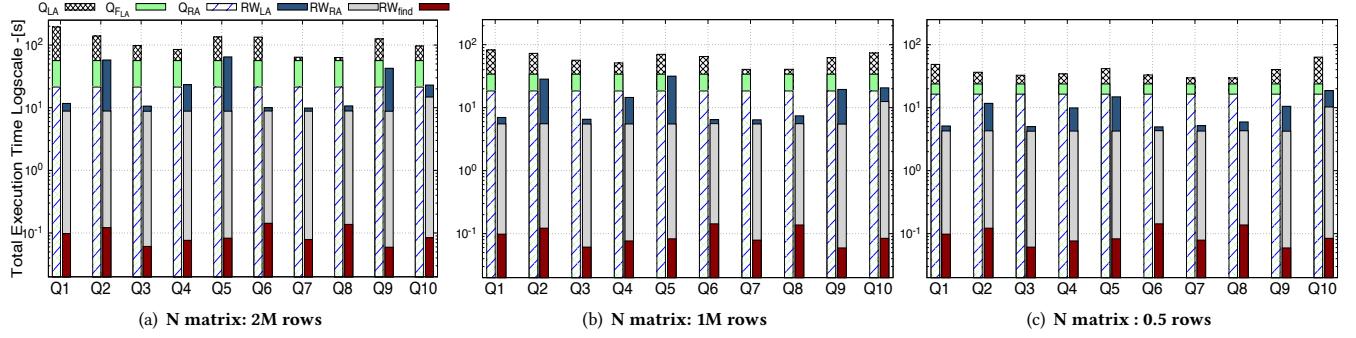


Figure 10: Micro-Hybrid Benchmark Twitter Dataset

Interestingly, for pipelines that involve transpose operator, Morpheus applies its special rewrite rules that replace an operation on \mathbf{M}^T with an operation on \mathbf{M} before pushing the operation to the base tables. For example, for P1.14: $\text{sum}(\text{colSums}(\mathbf{M}^T \mathbf{N}^T))$ and P2.12: $\text{sum}(\text{rowSums}(\mathbf{M}^T \mathbf{N}^T))$, Morpheus rewrites the pipelines to $\text{sum}((\text{rowSums}(\mathbf{N} \mathbf{M})))$ and $\text{sum}((\text{colSums}(\mathbf{N} \mathbf{M})))$, respectively, and then applies its factorized matrix-multiplication rewrite rule on $\mathbf{N} \mathbf{M}$. HADAD rewrites both pipelines to $\text{sum}(\mathbf{N} \mathbf{M})$, enabling again Morpheus to apply its factorized rewrite rule on $\mathbf{N} \mathbf{M}$ and achieving speed-up ranging from 1.3 up to 1.5 \times . The reason is that for both rewritings (found by Morpheus and HADAD), Morpheus's factorized matrix multiplication rewrite rule dominates the execution time.

9.2.2 Micro-hybrid Benchmark Experiments. In this experiment, we create a *micro-hybrid benchmark* on Twitter[12] and MIMIC [33] datasets to empirically study HADAD's rewriting benefits in a hybrid setting. The benchmark comprises ten different queries combining relational and linear algebra expressions.

Twitter Dataset Preparation. We obtain from Twitter API [12] 16GB of tweets (in JSON). We extract the structural parts of the dataset, which include user and tweet information, and store them in tables **User** (**U**) and **Tweet** (**T**), linked via PK-FK relationships. The dataset is detailed in §2. The tables **User** and **Tweet** as well as **TweetJSON** (**TJ**) are stored in Parquet format.

Twitter Queries and Views. Queries consist of two parts: (i) RA preprocessing (Q_{RA}) and (ii) LA analysis (Q_{LA}). In the Q_{RA} part, queries construct two matrices: \mathbf{M} and \mathbf{N} . The matrix \mathbf{M} ($2M \times 12$; *dense*) is the output of joining **T** and **U**. The construction of matrix \mathbf{N} is described in §2⁴. We fix the Q_{RA} part across all queries and vary the Q_{LA} part using a set of LA pipelines detailed below.

In addition to the views defined in §2, we define three hybrid RA-LA materialized views: V_3 , V_4 and V_5 , which store the result of applying `rowSums`, `colSums` and matrix multiplication operations over base tables **T** and **U** (viewed as matrices).

Importantly, rewritings based on these views can only be found by *exploiting together LA properties and Morpheus' rewrites rules* (we incorporated them in our framework as a set of integrity constraints). The full list of queries and views is in [11].

Discussion. For the first run of all queries (Figure 10(a)), matrix \mathbf{N} ($2M \times 1000$) is constructed from **TJ** for all tweets that mention “covid” and are posted from “USA”, and \mathbf{M} ($2M \times 12$) is constructed

No.	Expression
P3.1	<code>rowSums(XM) + (uv^T + N^T)v</code>
P3.2	<code>ucolSums((XM)^T) + N</code>
P3.3	$((\mathbf{N} + \mathbf{X})\mathbf{v})\text{colSums}(\mathbf{M})$
P3.4	$\text{sum}(\mathbf{C} + \mathbf{N}\text{rowSums}(\mathbf{X}\mathbf{M}))\mathbf{v}$
P3.5	$\text{ucolSums}(\mathbf{M}\mathbf{X}) + \mathbf{N}$
P3.6	$\text{rowSums}((\mathbf{M}\mathbf{X})^T) + (uv^T + \mathbf{N})\mathbf{v}$
P3.7	$X\mathbf{N}\mathbf{v} + \text{rowSums}(\mathbf{M})^T$
P3.8	$\mathbf{N} \odot \text{trace}(\mathbf{C} + v\text{colSums}(\mathbf{M}\mathbf{X})\mathbf{C})$
P3.9	$X \odot \text{sum}(\text{colSums}(\mathbf{C})^T + \text{rowSums}(\mathbf{M})) + \mathbf{N}$
P3.10	$\mathbf{N} \odot \text{sum}((\mathbf{X} + \mathbf{C})\mathbf{M})$

Table 7: LA pipelines used in micro-hybrid benchmark

as described above. Both matrices are then loaded in SystemML to be used for the analysis part (varied using the set of pipelines in Table 7). Before evaluating an LA pipeline, the queries filter \mathbf{N} 's rows with *filter-level* less than 4 (medium). Now, we analyze HADAD's rewrites and optimizations for each query in detail in the following paragraphs:

Q1: The rewriting *modifies the pre-processing* of \mathbf{N} by introducing V_1 and V_2 ; it also *pushes the filter-level selection* into the pre-processing part. Now, for the LA pipeline P3.1 (see Table 7), HADAD applies several optimizations: (i) it rewrites $(uv^T + \mathbf{N}^T)\mathbf{v}$ to $uv^T\mathbf{v} + \mathbf{N}^T\mathbf{v}$, where u and v are synthetic vectors of size 12×1 and $2M \times 1$. First, \mathbf{N} is ultra sparse, which makes the computation of $\mathbf{N}^T\mathbf{v}$ extremely efficient. Second, SystemML evaluates $uv^T\mathbf{v}$ efficiently in one go without intermediates, taking advantage of `tsmm` operator (discussed earlier) and `mmchain` for matrix multiply chains, where the best way to evaluate it computes $\mathbf{v}^T\mathbf{v}$ first, which results in a scalar, instead of computing uv^T , which results in a dense matrix of size $1000 \times 2M$. Alone, SystemML is unable to exploit its own efficient operations for lack of awareness of the LA property $A\mathbf{v} + B\mathbf{v} = (\mathbf{A} + \mathbf{B})\mathbf{v}$; (ii) HADAD also rewrites `rowSums(XM)` into XV_4 , where $V_4 = \text{rowSums}(\mathbf{T}) + K\text{rowSums}(\mathbf{U})$, by exploiting the property $\text{rowSums}(XM) = X\text{rowSums}(\mathbf{M})$ together with Morpheus's rewrite rule: $\text{rowSums}(\mathbf{M}) \rightarrow \text{rowSums}(\mathbf{T}) + K\text{rowSums}(\mathbf{U})$, where \mathbf{M} is the join's output table (matrix) of **T** and **U** base tables (feature matrices). The rewriting achieves up to 16.5 \times speed-up.

Q2: From hereafter, the rewriting of the pre-processing part of \mathbf{N} is similar to the one in Q1. The speed-up of 2.5 \times in this query comes from rewriting the pre-processing part and turning P3.2:

⁴Matrix \mathbf{N} is represented in Matrix Market Format (MTX) since it is sparse.

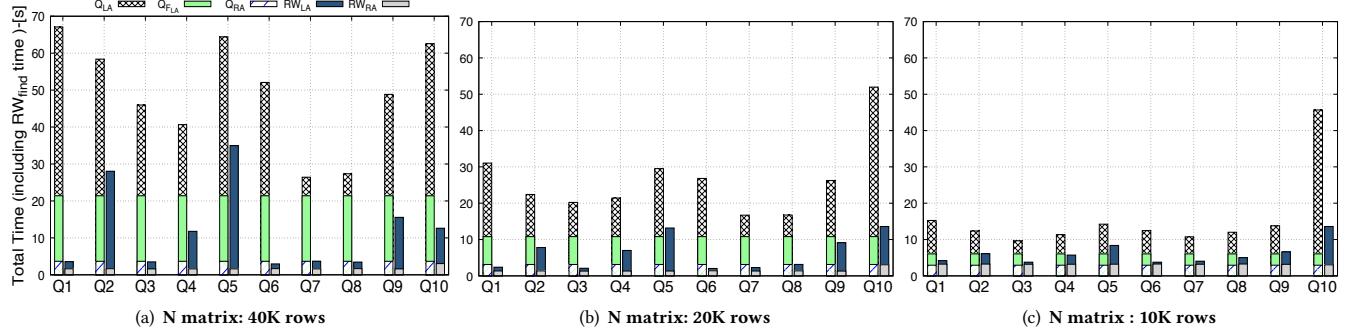


Figure 11: Micro-Hybrid Benchmark MIMIC

ucolSums($(XM)^T$) $+N$ to $u(XV_4)^T+N$, where u and X are synthetic matrices of size $2M \times 1$ and $1000 \times 2M$, respectively. The rewriting is achieved by exploiting : $\text{colSums}((XM)^T) = \text{rowSums}(XM)^T$ and $\text{rowSums}(XM) = X\text{rowSums}(M)$ along with Morpheus's rewrite rule: $\text{rowSums}(M) \rightarrow \text{rowSums}(T) + K\text{rowSums}(U)$. Both the rewriting and the original LA pipeline still introduce an unavoidable large dense intermediate of size $2M \times 1000$.

Q3: In addition to the rewriting of the pre-processing part, for **P3.3**, HADAD (i) avoids introducing the large dense intermediate ($N+X; 2M \times 1000$) by distributing the multiplication of v (1000×1) and consider the sparsity of N , which again makes the multiplication with v more efficient; (ii) directly rewrites $\text{colSums}(M)$ to $V_3 = [\text{colSums}(T), \text{colSums}(K)U]$ by utilizing Morpheus's rewrite rule: $\text{colSums}(M) \rightarrow [\text{colSums}(T), \text{colSums}(K)U]$. The obtained rewriting achieves $9.2\times$ speed-up.

Q4: For this query, the speed-up of $3.6\times$ attributes to rewriting the pre-processing part and optimizing **P3.4** by distributing the sum operation ($\text{sum}(A+B) = \text{sum}(A)+\text{sum}(B)$), which prevents materializing a dense intermediate. The matrix C here is a synthetic dense matrix of size $2M \times 1000$. Interestingly, SystemML includes this rewrite rule but did not apply it during optimization. Moreover, HADAD rewrites $\text{rowSums}(XM)$ to XV_4 (see **Q1**), which then enables SystemML to further optimize the chain of NXV_4v , where the sizes of X and v are $1000 \times 2M$ and 1×1000 , respectively.

Q5: HADAD's rewriting speeds-up this query by $2.3\times$. It rewrites $\text{u}colSums(MX)$ in **P3.5** to uV_3X (see **Q3** for V_3 's definition). The view is exploited by HADAD due to utilizing $\text{colSums}(MX) = \text{colSums}(M)X$ together with the Morpheus's rewrite rule (shown in **Q3**) for pushing the colSums to the base tables (matrices) U and T . The found rewriting enables SystemML to optimize the matrix-chain multiplication by computing V_3X first, which results in 1×1000 matrix instead of computing uV_3 which results in $2M \times 12$. The rewriting and the original pipeline still introduce unpreventable dense intermediate of size $2M \times 1000$.

Q6: The LA pipeline (**P3.6**) in this query is a variation of **P3.1**. In addition to distributing the multiplication of v , HADAD rewrites $\text{rowSums}((M X)^T)$ to $(V_3 X)^T$ (see **Q3** for V_3 's definition) by exploiting $\text{rowSums}((MX)^T) = \text{colSums}(MX)^T$ and $\text{colSums}(MX) = \text{colSums}(M)X$, all together with Morpheus's rewrite rule as illustrated in **Q3**. The obtained rewriting (along with the rewriting of the pre-processing part) achieves speed-up of $13.4\times$.

Q8: In addition to rewriting the pre-processing part, HADAD applies several optimizations for the analysis part (**P3.8**). First, it distributes the trace operation to prevent materializing a dense intermediate, which SystemML does not apply. Second, it exploits the view V_3 (see **Q3**) by utilizing $\text{colSums}(MX) = \text{colSums}(M)X$, which then enables SystemML to further optimize the matrix chain multiplication by computing $v((V_3 X)C)$, where the size of v , X and C are $20K \times 1$, $12 \times 20K$ and $20K \times 20K$, respectively. The element-wise multiplication with N in the rewiring and the original pipeline is extremely efficient since N is ultra-sparse. The rewriting speeds-up the query by $5.9\times$.

Q9: Beside rewriting the pre-processing part, the speed-up of $3\times$ also attributes to turning $\text{sum}(\text{colSums}(C)^T + \text{rowSums}(M))$ in **P3.9** to $\text{sum}(V_5)$, where $V_5 = [CT, CKU]$. The view V_5 is utilized by exploiting the property $\text{sum}(CM) = \text{sum}(\text{colSums}(C)^T \odot \text{rowSums}(M))$ together with Morpheus's rewrite rule: $CM \rightarrow [CT, CKU]$. The result of an element-wise multiplication with the dense matrix X in the rewriting and the original pipelines is a dense intermediate.

Q10: HADAD's rewrite leads to a speed-up of $3.9\times$. In addition to rewriting the pre-processing part, HADAD rewrites $\text{sum}((X + C)M)$ **P3.10** to $\text{sum}(XM) + \text{sum}(V_5)$, where $V_5 = [CT, CKU]$. The view V_5 is utilized by exploiting the property $(X + C)M = XM + CM$ together with Morpheus's rewrite rule: $CM \rightarrow [CT, CKU]$ (see **Q9**). This optimization goes beyond SystemML's optimization since it does not consider distributing the multiplication of M , which then enables exploiting the view and distributing the sum operation to avoid introducing a large dense intermediate.

For the second and third runs of the benchmark (see Figures 10(b) and 10(c)), we vary the text-search selection condition ("Trump" and "US election") in the pre-processing part where the number of N 's rows are $1M$ and $0.5M$, respectively. Accordingly, we change the size of synthetic matrices to ensure matrices' dimensions compatibility across the pipelines. Importantly, the rewriting time across all queries in the benchmark is very negligible compared to the original queries execution time.

MIMIC Dataset Preparation. MIMIC dataset [33] comprises health data for patients. The total size of the dataset is 46.6 GB, and it consists of : (i) all charted data for all patients and their hospital admission information, ICU stays, laboratory measurements, caregivers' notes, and prescriptions; (ii) the role of caregivers (e.g., MD stands for "medical doctor"), (iii) lab measurements (e.g., ABG stands for "arterial blood gas") and (iv) diagnosis related groups (DRG) codes

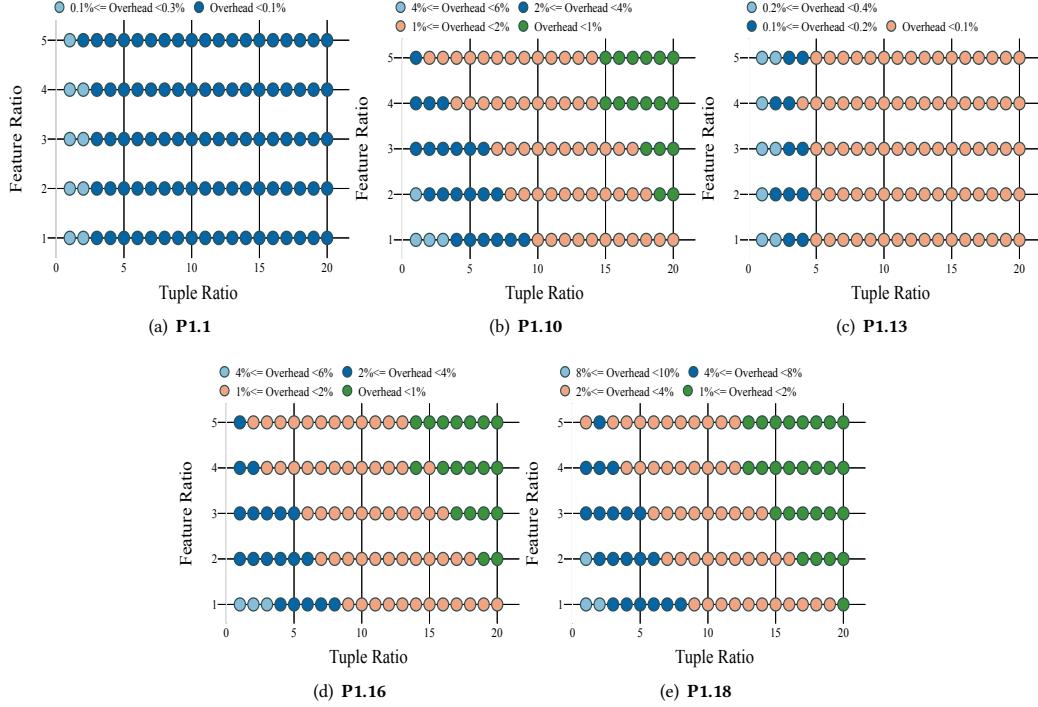


Figure 12: HADAD RW_{find} overhead as a percentage (%) of the total time ($Q_{exec} + RW_{fond}$) for pipelines P1.1, P1.10, P1.13, P1.16 and P1.18 running on Morpheus

descriptions. We use subset of the dataset, which includes **Patients** (P), **Admission** (A), **Service** (S), and **Callout** (C) tables. We convert tables’ categorical features (columns) to numeric using one-hot encoding [?].

MIMIC Queries and Views. Similar to the Twitter’s dataset benchmark, queries consist of two parts: (i) preprocessing (Q_{RA}) and (ii) analysis (Q_{LA}). In the Q_{RA} part, the queries construct two main matrices: M and N. The matrix M (40K×82; *dense*) is the join’s output table (matrix) of P and A. The matrix N (40K×30K; *ultra-sparse*) is *patient-service outcome* (e.g., cancelled (1), serving (2), etc) matrix, constructed from joining C and S for all patients who are in “CCU” care unit. We fix the Q_{RA} part across all queries and vary the Q_{LA} part using a set of LA pipelines in Table 7, a well as N and M matrices. For views, we define three cross RA-LA materialized views: V_1 , V_2 and V_3 , which store the result of applying `rowSums`, `colSums` and matrix multiplication operations over P and A base tables (matrices), respectively. These views can only be found by *exploiting together LA properties and Morpheus rewrites’ rules* in similar fashion as we detailed in Twitter’s benchmark experiment.

Discussion The results exhibit similar trends to the Twitter’s benchmark. The first run of the benchmark is shown in Figure 11(a); both matrices M and N are loaded in SystemML to be used for the analysis part (varied using the set of pipelines in Table 7). Before evaluating an LA pipeline, all queries filter N’s rows with *outcome* is equal to 2. HADAD applies the same set of optimizations as described in Twitter’s benchmark. For the second and third runs of the queries (see Figures 11(b) and 11(c)), we construct the N matrix

for patients who are in “TSICU” and “MICU” care units, where N’s rows are 20K and 10K, respectively.

9.2.3 Rewriting Time Overhead. For Morpheus, using the same experiment setup illustrated in §9.2.1, the rewriting time’s overhead for a set $\mathcal{P}_M^{opt} \subset \mathcal{P}^{Opt}$ of 5 pipelines (running on Morpheus) is very negligible compared to the pipelines’ execution time. For pipelines that contain matrix multiplication expressions such as P1.1 and P1.13, the rewriting time is generally less than 0.1% of the total time, as shown in Figure 12. However, for the other pipelines P1.10, P1.16, and P1.18, which contain only aggregate operations, the rewriting time is up to 9% of the total time (when the data size is very small (0.32GB) and the computation is extremely efficient) and less than 1% (when the data size is large (19.2GB) and the computation is expensive) as shown in Figure 12. As for the micro-hybrid benchmark, the RW_{find} reaches 78ms and up to 145 ms.

9.3 Experiments Summary

We have shown that HADAD brings significant performance-saving across LA-oriented and cross RA-LA platforms without the need to modify their internals. It improves their performance by order of magnitudes on typical LA-based and hybrid pipelines. Moreover, as we confirm experimentally, the time spent searching for rewritings is a small fraction of the query execution time hence a worthwhile investment. In addition, the rewriting overhead of pipelines that are already in the optimal form is very negligible compared to the original pipeline execution time in the presence of sparse/dense matrices and using naïve and MNC-based cost models.

10 RELATED WORK AND CONCLUSION

LA Systems/Libraries. SystemML [21] offers high-level R-like linear algebra language and applies some logical LA pattern-based rewrites and physical execution optimizations, based on cost estimates for the latter. SparkMLlib [42] provides LA operations and built-in function implementations of popular ML algorithms, such as linear regression, etc. on Spark RDDs. R [8] and NumPy [7] are two of the most popular computing environments for statistical data analysis, widely used in academia and industry. They provide a high-level abstraction that can simplify the programming of numerical and statistical computations, by treating matrices as first-class citizens and by providing a rich set of built-in LA operations. However, LA properties in most of these systems remain unexploited, which makes them *miss opportunities to use their own highly efficient operators* (recall hybrid-scenario in §2). Our experiments (§9.1) show that LA pipelines evaluation in these systems can be sped up, by more than 10×, by our rewriting using (i) LA properties and (ii) materialized views.

Bridging the Gap: RA and LA. There has been a recent increase in research for unifying the execution of RA and LA expressions [9, 27, 35, 40]. A key limitation of these approaches is that *the semantics of LA operations remains hidden behind built-in functions or UDFs, preventing performance-enhancing rewrites; as shown in §9.2.1.*

Closer to our work, SPORES [49], SPOOF [23], LARA [38] and RAVEN [34]. SPORES and SPOOF optimize LA expressions, by converting them into RA, optimizing the latter, and then converting the result back to an (optimized) LA expression. They are restricted to a small set of selected LA operations (the ones that can be expressed in RA), while we support significantly more (§6.1), and model properties allowing to optimize with them. LARA relies on a declarative domain-specific language for collections and matrices, which can enable optimization (i.e., *selection pushdown*) across the two algebraic abstractions. It heavily focuses on low-level optimization such as exploiting the choice of data layouts for physical LA operators’ optimization. RAVEN [34] takes a step forward by providing intermediate representation to enhance in-database model inferencing performance. It transforms classical ML models into equivalent neural networks (NN) to leverage highly optimized ML engines on CPU/GPU. [19, 25] study the expressive power of cross RA-LA [19] / LA [25] query languages. In contrast to HADAD, as all aforementioned solutions do not reason with constraints, they provide no capabilities for *holistic* semantic query optimizations including RA/LA views-based and LA pure rewritings; such optimizations can bring large performance saving as shown in (§9)

We see HADAD as complementary to all of these works, where it can be naturally and portably applied on top of these platforms; while their existing code need not be rewritten.

Conclusion. HADAD is an extensible lightweight framework for optimizing hybrid analytics queries, based on the powerful intermediate abstraction of a *a relational model with integrity constraints*. HADAD extends the capability of [15] with a reduction from LA (or LA view)-based rewriting to relational rewriting under constraints. It enables a full exploration of rewrites using a large set of LA operations, with no modification to the execution platform. Our experiments show significant performance gains on various LA and hybrid workloads across popular LA and cross RA-LA platforms.

REFERENCES

- [1] Amazon Review Data. <https://nijianmo.github.io/amazon/index.html>, Accessed August, 2020.
- [2] Breeze Wiki. <https://github.com/scalanlp/breeze/wiki>, Accessed June, 2020.
- [3] Kaggle Survey. <https://www.kaggle.com/kaggle-survey-2019>, Accessed June, 2020.
- [4] Morpheus. <https://github.com/lchen001/Morpheus>, Accessed December, 2020.
- [5] Native Blas in SystemDS. <https://apache.github.io/systemds/native-backend>, Accessed June, 2020.
- [6] Netflix Movie Rating. <https://www.kaggle.com/netflix-inc/netflix-prize-data>, Accessed August, 2020.
- [7] NumPy. <https://numpy.org/>, Accessed June, 2020.
- [8] Project R. <https://www.r-project.org/other-docs.html>, Accessed June, 2020.
- [9] Python/NumPy in Monetdb. <https://tinyurl.com/1l1jy21v>, Accessed June, 2020.
- [10] SparkMLlib. <https://spark.apache.org/mllib>, Accessed June, 2020.
- [11] Technical Report. https://github.com/hadad-paper/HADAD_SIGMOD2021.
- [12] Twitter API. <https://developer.twitter.com/en/docs>, Accessed January, 2021.
- [13] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. Tensorflow: A System for Large-Scale Machine Learning. In *USENIX*, pages 265–283, 2016.
- [14] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [15] R. Alotaibi, D. Bursztyn, A. Deutscher, I. Manolescu, and S. Zampetakis. Towards Scalable Hybrid Stores: Constraint-based Rewriting to the Rescue. In *SIGMOD*, pages 1660–1677, 2019.
- [16] R. Alotaibi, B. Cautis, A. Deutscher, M. Latrache, I. Manolescu, and Y. Yang. ESTOCADA: Towards Scalable Polystore Systems. *PVLDB*, pages 2949–2952, 2020.
- [17] M. Armbrust, R. S. Xin, C. Lian, Y. Huai, D. Liu, J. K. Bradley, X. Meng, T. Kaftan, M. J. Franklin, A. Ghodsi, et al. SparkSQL: Relational Data Processing in Spark. In *SIGMOD*, pages 1383–1394, 2015.
- [18] S. Axler. *Linear Algebra Done Right*. Springer, 2015.
- [19] P. Barceló, N. Higuera, J. Pérez, and B. Subercaseaux. On the Expressiveness of LARA: A Unified Language for Linear and Relational Algebra. In *ICDT*, pages 6:1–6:20, 2020.
- [20] D. Baylor, E. Breck, H.-T. Cheng, N. Fiedel, C. Y. Foo, Z. Haque, S. Haykal, M. Ispir, V. Jain, L. Koc, et al. Tfx: A Tensorflow-based Production-Scale Machine Learning Platform. In *SIGKDD*, pages 1387–1395, 2017.
- [21] M. Boehm, M. W. Dusenberry, D. Eriksson, A. V. Evfimievski, F. M. Manshadi, N. Pansare, B. Reinwald, F. R. Reiss, P. Sen, A. C. Surve, et al. SystemML: Declarative Machine Learning on Spark. *PVLDB*, pages 1425–1436, 2016.
- [22] M. Boehm, A. V. Evfimievski, N. Pansare, and B. Reinwald. Declarative Machine Learning Classification of Basic Properties and Types. *arXiv:1605.05826*, 2016.
- [23] M. Boehm, B. Reinwald, D. Hutchison, P. Sen, A. V. Evfimievski, and N. Pansare. On Optimizing Operator Fusion Plans for Large-Scale Machine Learning in SystemML. *PVLDB*, pages 1755–1768, 2018.
- [24] J.-H. Böse, V. Flunkert, J. Gasthaus, T. Januschowski, D. Lange, D. Salinas, S. Schelter, M. Seeger, and Y. Wang. Probabilistic Demand Forecasting at Scale. *PVLDB*, pages 1694–1705, 2017.
- [25] R. Brijder, F. Geerts, J. V. den Bussche, and T. Weerwag. On the Expressive Power of Query Languages for Matrices. *ACM Trans. Database Syst.*, pages 15:1–15:31, 2019.
- [26] A. K. Chandra and P. M. Merlin. Optimal Implementation of Conjunctive Queries in Relational Databases. In *ACM symposium on Theory of computing*, pages 77–90, 1977.
- [27] L. Chen, A. Kumar, J. Naughton, and J. M. Patel. Towards Linear Algebra Over Normalized Data. *PVLDB*, pages 1214–1225, 2017.
- [28] T. A. Davis and Y. Hu. The University of Florida Sparse Matrix Collection. *TOMS*, 2011.
- [29] A. Deutscher. FOL Modeling of Integrity Constraints (Dependencies). In *Encyclopedia of Database Systems, Second Edition*. 2018.
- [30] A. Deutscher, L. Popa, and V. Tannen. Query Reformulation with Constraints. In *SIGMOD*, pages 65–73, 2006.
- [31] I. Ileana. *Query Rewriting Using Views : a Theoretical and Practical Perspective*. Theses, Télécom ParisTech, Oct. 2014.
- [32] I. Ileana, B. Cautis, A. Deutscher, and Y. Katsis. Complete Yet Practical Search for Minimal Query Reformulations Under Constraints. In *SIGMOD*, pages 1015–1026, 2014.
- [33] A. Johnson et al. MIMIC-III. <http://www.nature.com/articles/sdata201635>, 2016.
- [34] K. Karanassos, M. Interlandi, D. Xin, F. Psallidas, R. Sen, K. Park, I. Popivanov, S. Nakandal, S. Krishnan, M. Weimer, et al. Extending Relational Query Processing with ML Inference. In *CIDR*, 2020.
- [35] D. Kernert, F. Köhler, and W. Lehner. Bringing Linear Algebra Objects to Life in a Column-Oriented in-Memory Database. In *In Memory Data Management and Analysis*, pages 44–55. Springer, 2013.
- [36] A. Kumar, M. Boehm, and J. Yang. Data Management in Machine Learning: Challenges, Techniques, and Systems. In *SIGMOD*, pages 1717–1722, 2017.

- [37] A. Kumar, R. McCann, J. Naughton, and J. M. Patel. Model Selection Management Systems: The Next Frontier of Advanced Analytics. In *SIGMOD*, pages 17–22, 2016.
- [38] A. Kunft, A. Katsifodimos, S. Schelter, S. Breß, T. Rabl, and V. Markl. An Intermediate Representation for Optimizing Machine Learning Pipelines. *PVLDB*, pages 1553–1567, 2019.
- [39] K. Kuttler. *Linear Algebra: Theory and Applications*. The Saylor Foundation, 2012.
- [40] S. Luo, Z. J. Gao, M. Gubanov, L. L. Perez, and C. Jermaine. Scalable Linear Algebra on a Relational Database System. In *ICDE*, pages 1224–1238, 2018.
- [41] J. MacQueen et al. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, pages 281–297, 1967.
- [42] X. Meng, J. Bradley, B. Yavuz, E. Sparks, S. Venkataraman, D. Liu, J. Freeman, D. Tsai, M. Amde, S. Owen, et al. MLlib: Machine Learning in Apache Spark. *The Journal of Machine Learning Research*, pages 1235–1241, 2016.
- [43] M. Milani, S. Hosseinpour, and H. Pehlivan. Rule-based Production of Mathematical Expressions. *Mathematics*, 6:254, 11 2018.
- [44] T. M. Mitchell. *Machine Learning*. McGraw-Hill, 1997.
- [45] D. Sculley, G. Holt, D. Golovin, E. Davydov, T. Phillips, D. Ebner, V. Chaudhary, M. Young, J.-F. Crespo, and D. Dennison. Hidden Technical Debt in Machine Learning Systems. In *NeurIPS*, pages 2503–2511, 2015.
- [46] J. Sommer, M. Boehm, A. V. Evfimievski, B. Reinwald, and P. J. Haas. MNC: Structure-Exploiting Sparsity Estimation for Matrix Expressions. In *SIGMOD*, pages 1607–1623, 2019.
- [47] E. R. Sparks, A. Talwalkar, D. Haas, M. J. Franklin, M. I. Jordan, and T. Kraska. Automating Model Search for Large Scale Machine Learning. In *ACM Symposium on Cloud Computing*, pages 368–380, 2015.
- [48] A. Thomas and A. Kumar. A Comparative Evaluation of Systems for Scalable Linear Algebra-based Analytics. *PVLDB*, pages 2168–2182, 2018.
- [49] Y. R. Wang, S. Hutchison, J. Leang, B. Howe, and D. Suciu. SPORES: Sum-Product Optimization via Relational Equality Saturation for Large Scale Linear Algebra. *PVLDB*, pages 1919–1932, 2020.

A L_{ops} OPERATIONS PROPERTIES (LA_{prop}) CAPTURED AS INTEGRITY CONSTRAINTS

Table 8: L_{ops} Operations Properties (LA_{prop}) Captured as Integrity Constraints

LA Property	Relational Encoding as Integrity Constraints
Addition of Matrices	
$M + N = N + M$	$\forall M, N, R \text{ add}_M(M, N, R) \rightarrow \text{add}_M(N, M, R)$
$(M + N) + D = M + (N + D)$	$\forall M, N, D, R_1, R_2 \text{ add}_M(M, N, R_1) \wedge \text{add}_M(R_1, D, R_2) \rightarrow \exists R_3 \text{ add}_M(N, D, R_3) \wedge \text{add}_M(M, R_3, R_2)$
$c(M + N) = cM + cN$	$\forall c, M, N, R_1, R_2 \text{ add}_M(M, N, R_1) \wedge \text{multi}_{MS}(c, R_1, R_2) \rightarrow \exists R_3, R_4 \text{ multi}_{MS}(c, M, R_3) \wedge \text{multi}_{MS}(c, N, R_4) \wedge \text{add}_M(R_3, R_4, R_2)$
$(c + d)M = cM + dM$	$\forall c, d, s, M, R_1 \text{ adds}(c, d, s) \wedge \text{multi}_{MS}(s, M, R_1) \rightarrow \exists R_2, R_3 \text{ multi}_{MS}(c, M, R_2) \wedge \text{multi}_{MS}(d, M, R_3) \wedge \text{add}_M(R_2, R_3, R_1)$
$M + 0 = M$	$\rightarrow \exists \text{Zero}(O)$ $\forall M, n, O \text{ name}(M, n) \wedge \text{Zero}(O) \rightarrow \text{add}_M(M, O, M)$ $\forall O \text{ Zero}(O) \rightarrow \text{add}_M(O, O, O)$
Product of Matrices	
$(MN)D = M(ND)$	$\forall M, N, D, R_1, R_2 \text{ multi}_M(M, N, R_1) \wedge \text{multi}_M(R_1, D, R_2) \rightarrow \exists R_3 \text{ multi}_M(N, D, R_3) \wedge \text{multi}_M(M, R_3, R_2)$
$M(N + D) = MN + MD$	$\forall M, N, D, R_1, R_2 \text{ add}_M(N, D, R_1) \wedge \text{multi}_M(M, R_1, R_2) \rightarrow \exists R_3, R_4 \text{ multi}_M(M, N, R_3) \wedge \text{multi}_M(M, D, R_4) \wedge \text{add}_M(R_3, R_4, R_2)$
$(M + N)D = MD + MD$	$\forall M, N, D, R_1, R_2 \text{ add}_M(M, N, R_1) \wedge \text{multi}_M(R_1, D, R_2) \rightarrow \exists R_3, R_4 \text{ multi}_M(M, D, R_3) \wedge \text{multi}_M(N, D, R_4) \wedge \text{add}_M(R_3, R_4, R_2)$
$d(MN) = (dM)N$	$\forall d, M, N, R_1, R_2 \text{ multi}_M(M, N, R_1) \wedge \text{multi}_{MS}(d, R_1, R_2) \rightarrow \exists R_3 \text{ multi}_{MS}(d, M, R_3) \wedge \text{multi}_M(R_3, N, R_2)$
$c(dM) = (cd)M$	$\forall c, dM, R_1, R_2 \text{ multi}_{MS}(d, M, R_1) \wedge \text{multi}_{MS}(c, R_1, R_2) \rightarrow \exists s \text{ multi}_S(c, d, s) \wedge \text{multi}_{MS}(s, M, R_2)$
$I_k M = M = MI_z$	$\forall M, n, k, z \text{ name}(M, n), \text{size}(M, k, z) \rightarrow \exists I_1 \text{ Identity}(I_1), \text{size}(I_1, k, k)$ $\forall M, n, k, z \text{ name}(M, n), \text{size}(M, k, z) \rightarrow \exists I_1 \text{ Identity}(I_1), \text{size}(I_1, z, z)$ $\forall M, I_1, n, k, z \text{ name}(M, n), \text{size}(M, k, z), \text{Identity}(I_1), \text{size}(I_1, k, k) \rightarrow \text{multi}_M(I_1, M, M)$ $\forall M, I_1, n, k, z \text{ name}(M, n), \text{size}(M, k, z), \text{Identity}(I_1), \text{size}(I_1, z, z) \rightarrow \text{multi}_M(M, I_1, M)$
Transposition of Matrices	
$(MN)^T = (N)^T(M)^T$	$\forall M, N, R_1, R_2 \text{ multi}_M(M, N, R_1) \wedge \text{tr}(R_1, R_2) \rightarrow \exists R_3, R_4 \text{ tr}(M, R_3) \wedge \text{tr}(N, R_4) \wedge \text{multi}_M(R_4, R_3, R_2)$
$(M + N)^T = (M)^T + (N)^T$	$\forall M, N, R_1, R_2 \text{ add}_M(M, N, R_1) \wedge \text{tr}(R_1, R_2) \rightarrow \exists R_3, R_4 \text{ tr}(M, R_3) \wedge \text{tr}(N, R_4) \wedge \text{add}_M(R_3, R_4, R_2)$
$(cM)^T = c(M)^T$	$\forall c, M, R_1, R_2 \text{ multi}_{MS}(c, M, R_1) \wedge \text{tr}(R_1, R_2) \rightarrow \exists R_3 \text{ tr}(M, R_3) \wedge \text{multi}_{MS}(c, R_3, R_2)$
$((M)^T)^T = M$	$\forall n, M \text{ name}(M, n) \rightarrow \exists R_1 \text{ tr}(M, R_1) \wedge \text{tr}(R_1, M)$
$(I)^T = I$, where I is identity matrix	$\forall I_1 \text{ Identity}(I_1) \rightarrow \text{tr}(I_1, I_1)$
$(O)^T = O$, where O is zero matrix	$\forall O_1^i \text{ Zero}(O_1^i) \rightarrow \text{tr}(O_1^i, O_1^i)$
Inverses of Matrices	
$((M)^{-1})^{-1} = M$	$\forall n, M \text{ name}(M, n) \rightarrow \exists R_1 \text{ inv}_M(M, R_1) \wedge \text{inv}_M(R_1, M)$
$(MN)^{-1} = (N)^{-1}(M)^{-1}$	$\forall M, N, R_1, R_2 \text{ multi}_M(M, N, R_1) \wedge \text{inv}_M(R_1, R_2) \rightarrow \exists R_3, R_4 \text{ inv}_M(M, R_3) \wedge \text{inv}_M(N, R_4) \wedge \text{multi}_M(R_4, R_3, R_2)$
$((M)^T)^{-1} = ((M)^{-1})^T$	$\forall M, R_1, R_2 \text{ tr}(M, R_1) \wedge \text{inv}_M(R_1, R_2) \rightarrow \exists R_3 \text{ inv}_M(M, R_3) \wedge \text{tr}(R_3, R_2)$
$((kM))^{-1} = k^{-1}(M)^{-1}$	$\forall k, M, R_1, R_2 \text{ multi}_{MS}(k, M, R_1) \wedge \text{inv}_M(R_1, R_2) \rightarrow \exists R_3, s \text{ inv}_S(k, s) \wedge \text{inv}_M(M, R_3) \wedge \text{multi}_{MS}(s, R_3, R_2)$
$M^{-1}M = I = MM^{-1}$	$\forall M, R_1, R_2 \text{ inv}_M(M, R_1) \wedge \text{multi}_M(R_1, M, R_2) \rightarrow \text{Identity}(R_2)$ $\forall M, R_1, R_2 \text{ inv}_M(M, R_1) \wedge \text{multi}_M(M, R_1, R_2) \rightarrow \text{Identity}(R_2)$

Table 9: L_{ops} Operations Properties (LA_{prop}) Captured as Integrity Constraints

LA Property	Relational Encoding as Integrity Constraints
Determinant of Matrices	
$\det(MN) = \det(M) * \det(N)$	$\forall M, N, R_1, d \text{ multi}_M(M, N, R_1) \wedge \det(R_1, d) \rightarrow \exists d_1, d_2 \det(M, d_1) \wedge \det(N, d_2) \wedge \text{multi}_S(d_1, d_2, d)$
$\det((M)^T) = \det(M)$	$\forall M, R_1, d \text{ tr}(M, R_1) \wedge \det(R_1, d) \rightarrow \det(M, d)$
$\det((M)^{-1}) = (\det(M))^{-1}$	$\forall M, R_1, d \text{ inv}_M(M, R_1) \wedge \det(R_1, d) \rightarrow \exists d_1 \det(M, d_1) \wedge \text{inv}_S(d_1, d)$
$\det((cM)) = c^k \det(M)$	$\forall M, c, k, d^O \text{ size}(M, k, k) \wedge \text{multi}_{MS}(c, M, d^O) \rightarrow \exists s_1, s_2 \text{ pow}(c, k, s_1) \wedge \det(M, s_2) \wedge \text{multi}_S(s_1, s_2, d^O)$
$\det((I)) = 1$	$\forall I_1, d^O \text{ Identity}(I_1) \wedge \det(I_1, d^O) \rightarrow d^O = 1$
Adjoint of Matrices	
$\text{adj}(M) = R^T$	$\forall M, R_1 \text{ adj}(M, R_1) \rightarrow \exists R_2 \text{ cof}(M, R_2) \wedge \text{tr}(R_2, R_1)$
$\text{adj}(M)^T = \text{adj}(M^T)$	$\forall M, R_1, R_2 \text{ adj}(M, R_1) \wedge \text{tr}(R_1, R_2) \rightarrow \exists R_3 \text{ tr}(M, R_3) \wedge \text{adj}(R_3, R_2)$
$\text{adj}(M)^{-1} = \text{adj}(M^{-1})$	$\forall M, R_1, R_2 \text{ adj}(M, R_1) \wedge \text{inv}_M(R_1, R_2) \rightarrow \exists R_3 \text{ inv}_M(M, R_3) \wedge \text{adj}(R_3, R_2)$
$\text{adj}(MN) = \text{adj}(NM)$	$\forall M, R_1, R_2 \text{ multi}_M(M, N, R_1) \wedge \text{adj}(R_1, R_2) \rightarrow \exists R_3 \text{ multi}_M(N, M, R_3) \wedge \text{adj}(R_3, R_2)$
Trace of Matrices	
$\text{trace}(M + N) = \text{trace}(M) + \text{trace}(N)$	$\forall M, N, R_1, s_1 \text{ add}_M(M, N, R_1) \wedge \text{trace}(R_1, s_1) \rightarrow \exists s_2, s_3 \text{ trace}(M, s_2) \wedge \text{trace}(N, s_3) \wedge \text{add}_S(s_2, s_3, s_1)$
$\text{trace}(MN) = \text{trace}(NM)$	$\forall M, N, R_1, s_1 \text{ multi}_M(M, N, R_1) \wedge \text{trace}(R_1, s_1) \rightarrow \exists R_2 \text{ multi}_M(N, M, R_2) \wedge \text{trace}(R_2, s_1)$
$\text{trace}(M^T) = \text{trace}(M)$	$\forall M, R_1, s_1 \text{ tr}(M, R_1) \wedge \text{trace}(R_1, s_1) \rightarrow \text{trace}(M, s_1)$
$\text{trace}(cM) = c\text{trace}(M)$	$\forall M, R_1, c, s_1 \text{ multi}_{MS}(c, M, R_1) \wedge \text{trace}(R_1, s_1) \rightarrow \exists s_2 \text{ trace}(M, s_2) \wedge \text{multi}_S(c, s_1, s_2)$
$\text{trace}(I_k) = k$	$\forall I_O, k, s_1 \text{ Identity}(I_O) \wedge \text{size}(I_O, k, k) \wedge \text{trace}(I_O, s_1) \rightarrow s_1 = k$
Direct Sum	
$(M \oplus N) + (C \oplus D) = (M + C) \oplus (N + D)$	$\forall M, N, R_1, C, D, R_2, R_3 \text{ sum}_D(M, N, R_1) \wedge \text{sum}_D(C, D, R_2) \wedge \text{add}_M(R_1, R_2, R_3) \rightarrow \exists R_4, R_5 \text{ add}_M(M, C, R_4) \wedge \text{add}_M(N, D, R_5)$
$(M \oplus N) + (C \oplus D) = (MC) \oplus (ND)$	$\forall M, N, R_1, C, D, R_2, R_3 \text{ sum}_D(M, N, R_1) \wedge \text{sum}_D(C, D, R_2) \wedge \text{multi}_M(R_1, R_2, R_3) \rightarrow \exists R_4, R_5 \text{ multi}_M(M, C, R_4) \wedge \text{multi}_M(N, D, R_5)$
$c(M \oplus N) = (cM \oplus cN)$	$\forall M, N, R_1, c, R_2 \text{ sum}_D(M, N, R_1) \wedge \text{multi}_{MS}(c, R_1, R_2) \rightarrow \exists R_3, R_4 \text{ multi}_{MS}(M, c, R_3) \wedge \text{multi}_{MS}(N, c, R_4) \wedge \text{sum}_D(R_3, R_4, R_2)$
Exponential of Matrices	
$\exp(0) = I$	$\forall O, R_1 \text{ Zero}(O) \wedge \exp(O, R_1) \rightarrow \text{Identity}(R_1)$
$\exp(M^T) = \exp(M)^T$	$\forall M, R_1, R_2 \text{ tr}(M, R_1) \wedge \exp(R_1, R_2) \rightarrow \exists R_3 \exp(M, R_3) \wedge \exp(R_3, R_2)$

Table 10: Matrix Decompositions Properties Captured as Integrity Constraints

Decomposition Property	Relational Encoding as Integrity Constraints
Cholesky Decomposition (CD)	
$\text{CHO}(M) = L$ such that $M = LL^T$, where M is symmetric positive definite	$\forall M \text{ type}(M, "S") \rightarrow \exists L_1 \exists L_2 \text{ CHO}(M, L_1) \wedge \text{type}(L_1, "L") \wedge \text{tr}(L_1, L_2) \wedge \text{multi}_M(L_1, L_2, M)$
QR Decomposition	
$QR(M) = [Q, R]$ such that $M = QR$	$\forall M \forall n \forall k \text{ name}(M, n) \wedge \text{size}(M, k, z) \rightarrow \exists Q_1, R_1$ $QR(M, Q_1, R_1) \wedge \text{type}(Q_1, "O") \wedge \text{size}(Q_1, k, k) \wedge \text{type}(R_1, "U") \wedge \text{size}(R_1, k, z) \wedge \text{multi}_M(Q_1, R_1, M)$ $\forall Q_1 \text{ type}(Q_1, "O") \wedge \text{size}(Q_1, k, k) \rightarrow \exists I_1 QR(Q_1, Q_1, I_1) \wedge \text{identity}(I_1) \wedge \text{size}(I_1, k, k) \wedge \text{multi}_M(Q_1, I_1, Q_1)$ $\forall R_1 \text{ type}(R_1, "U") \wedge \text{size}(R_1, k, z) \rightarrow \exists I_1 QR(R_1, I_1, R_1) \wedge \text{identity}(I_1) \wedge \text{size}(I_1, k, k) \wedge \text{multi}_M(I_1, R_1, R_1)$ $\forall I_1 \text{ identity}(I_1) \rightarrow QR(I_1, I_1, I_1)$
LU Decomposition	
$LU(M) = [L, U]$ such that $M = LU$	$\forall M \forall n \forall k \text{ name}(M, n) \wedge \text{size}(M, k, z) \rightarrow \exists L_1, U_1$ $LU(M, L_1, U_1) \wedge \text{type}(L_1, "L") \wedge \text{size}(L_1, k, z) \wedge \text{type}(U_1, "U") \wedge \text{size}(U_1, z, z) \wedge \text{multi}_M(L_1, U_1, M)$ $\forall L_1 \text{ type}(L_1, "L") \wedge \text{size}(L_1, k, z) \rightarrow \exists I_1 LU(L_1, L_1, I_1) \wedge \text{identity}(I_1) \wedge \text{multi}_M(L_1, I_1, L_1) \wedge \text{size}(I_1, z, z)$ $\forall U_1 \text{ type}(U_1, "U") \wedge \text{size}(U_1, z, z) \rightarrow \exists I_1 LU(U_1, I_1, U_1) \wedge \text{identity}(I_1) \wedge \text{size}(I_1, z, z) \wedge \text{multi}_M(I_1, U_1, U_1)$ $\forall I_1 \text{ identity}(I_1) \rightarrow LU(I_1, I_1, I_1)$
Pivoted LU Decomposition	
$LUP(M) = [L, U, P]$ such that $PM = LU$, where M is a square matrix	$\forall M \forall n \forall k \text{ name}(M, n) \wedge \text{size}(M, k, z) \rightarrow \exists L_1, U_1, P_1, R_1$ $LUP(M, L_1, U_1, P_1) \wedge \text{type}(L_1, "L") \wedge \text{type}(U_1, "U") \wedge \text{type}(P_1, "P") \wedge \text{multi}_M(L_1, U_1, R_1) \wedge \text{multi}_M(P_1, M, R_1)$ $\forall L_1 \text{ type}(L_1, "L") \wedge \text{size}(L_1, k, z) \rightarrow \exists I_1, I_2 LUP(L_1, L_1, I_1, I_2) \wedge \text{identity}(I_1) \wedge \text{identity}(I_2) \wedge \text{size}(I_1, z, z) \wedge \text{size}(I_2, k, k) \wedge \text{multi}_M(L_1, I_1, L_1) \wedge \text{multi}_M(I_2, L_1, L_1)$ $\forall U_1 \text{ type}(U_1, "U") \rightarrow \exists I_1 LUP(U_1, I_1, U_1, I_1) \wedge \text{identity}(I_1) \wedge \text{multi}_M(I_1, U_1, U_1)$ $\forall I_1 \text{ identity}(I_1) \rightarrow LU(I_1, I_1, I_1)$

B SYSTEMML REWRITE RULES ENCODED AS INTEGRITY CONSTRAINTS

Table 11: SystemML Algebraic Aggregate Rewrite Rules Captured as Integrity Constraints

SystemML Algebraic Simplification Rule	Integrity Constraints $\mathcal{MMC}_{StatAgg}$
UnnecessaryAggregates	
$\text{sum}(t(M)) \rightarrow \text{sum}(M)$	$\forall M, R_1, s \text{ tr}(M, R_1), \text{sum}(R_1, s) \rightarrow \text{sum}(M, s)$
$\text{sum}(\text{rev}(M)) \rightarrow \text{sum}(M)$	$\forall M, R_1, s \text{ rev}(M, R_1), \text{sum}(R_1, s) \rightarrow \text{sum}(M, s)$
$\text{sum}(\text{rowSums}(M)) \rightarrow \text{sum}(M)$	$\forall M, R_1, s \text{ rowSums}(M, R_1), \text{sum}(R_1, s) \rightarrow \text{sum}(M, s)$
$\text{sum}(\text{colSums}(M)) \rightarrow \text{sum}(M)$	$\forall M, R_1, s \text{ colSums}(M, R_1), \text{sum}(R_1, s) \rightarrow \text{sum}(M, s)$
$\text{min}(\text{rowMin}(M)) \rightarrow \text{min}(M)$	$\forall M, R_1, s \text{ rowMin}(M, R_1), \text{min}(R_1, s) \rightarrow \text{min}(M, s)$
$\text{min}(\text{colMin}(M)) \rightarrow \text{min}(M)$	$\forall M, R_1, s \text{ colMin}(M, R_1), \text{min}(R_1, s) \rightarrow \text{min}(M, s)$
$\text{max}(\text{colMax}(M)) \rightarrow \text{max}(M)$	$\forall M, R_1, s \text{ colMax}(M, R_1), \text{max}(R_1, s) \rightarrow \text{max}(M, s)$
$\text{max}(\text{rowMax}(M)) \rightarrow \text{max}(M)$	$\forall M, R_1, s \text{ rowMax}(M, R_1), \text{max}(R_1, s) \rightarrow \text{max}(M, s)$
pushdownUnaryAggTransposeOp	
$\text{rowSums}(t(M)) \rightarrow t(\text{colSums}(M))$	$\forall M, R_1, R_2 \text{ tr}(M, R_1) \wedge \text{rowSums}(R_1, R_2) \rightarrow \exists R_3 \text{ colSums}(M, R_3) \wedge \text{tr}(R_3, R_2)$
$\text{colSums}(t(M)) \rightarrow t(\text{rowSums}(M))$	$\forall M, R_1, R_2 \text{ tr}(M, R_1) \wedge \text{colSums}(R_1, R_2) \rightarrow \exists R_3 \text{ rowSums}(M, R_3) \wedge \text{tr}(R_3, R_2)$
$\text{rowMean}(t(M)) \rightarrow t(\text{colMean}(M))$	$\forall M, R_1, R_2 \text{ tr}(M, R_1) \wedge \text{colMean}(R_1, R_2) \rightarrow \exists R_3 \text{ rowMean}(M, R_3) \wedge \text{tr}(R_3, R_2)$
$\text{colMean}(t(M)) \rightarrow t(\text{rowMean}(M))$	$\forall M, R_1, R_2 \text{ tr}(M, R_1) \wedge \text{rowMean}(R_1, R_2) \rightarrow \exists R_3 \text{ colMean}(M, R_3) \wedge \text{tr}(R_3, R_2)$
$\text{rowVar}(t(M)) \rightarrow t(\text{colVar}(M))$	$\forall M, R_1, R_2 \text{ tr}(M, R_1) \wedge \text{rowVar}(R_1, R_2) \rightarrow \exists R_3 \text{ colVar}(M, R_3) \wedge \text{tr}(R_3, R_2)$
$\text{colVar}(t(X)) \rightarrow t(\text{rowVar}(X))$	$\forall M, R_1, R_2 \text{ tr}(M, R_1) \wedge \text{colVar}(R_1, R_2) \rightarrow \exists R_3 \text{ rowVar}(M, R_3) \wedge \text{tr}(R_3, R_2)$
$\text{rowMax}(t(M)) \rightarrow t(\text{colMax}(M))$	$\forall M, R_1, R_2 \text{ tr}(M, R_1) \wedge \text{rowMax}(R_1, R_2) \rightarrow \exists R_3 \text{ colMax}(M, R_3) \wedge \text{tr}(R_3, R_2)$
$\text{colMax}(t(M)) \rightarrow t(\text{rowMax}(M))$	$\forall M, R_1, R_2 \text{ tr}(M, R_1) \wedge \text{colMax}(R_1, R_2) \rightarrow \exists R_3 \text{ rowMax}(M, R_3) \wedge \text{tr}(R_3, R_2)$
$\text{rowMin}(t(M)) \rightarrow t(\text{colMin}(M))$	$\forall M, R_1, R_2 \text{ tr}(M, R_1) \wedge \text{rowMin}(R_1, R_2) \rightarrow \exists R_3 \text{ colMin}(M, R_3) \wedge \text{tr}(R_3, R_2)$
$\text{colMin}(t(M)) \rightarrow t(\text{rowMin}(M))$	$\forall M, R_1, R_2 \text{ tr}(M, R_1) \wedge \text{colMin}(R_1, R_2) \rightarrow \exists R_3 \text{ rowMin}(M, R_3) \wedge \text{tr}(R_3, R_2)$
simplifyTraceMatrixMult	
$\text{trace}(MN) \rightarrow \text{sum}(M \odot t(N))$	$\forall M, N, R_1, r \text{ multi}_M(M, N, R_1) \wedge \text{trace}(R_1, r) \rightarrow \exists R_3, R_4 \text{ tr}(N, R_3) \wedge \text{multi}_E(M, R_3, R_4) \wedge \text{sum}(R_4, r)$
simplifySumMatrixMult	
$\text{sum}(MN) \rightarrow \text{sum}(t(\text{colSums}(M)) \odot \text{rowSums}(N))$	$\forall M, N, R_1, r \text{ multi}_M(M, N, R_1) \wedge \text{sum}(R_1, r) \rightarrow \exists R_2, R_3, R_4, R_5 \text{ colSums}(M, R_2) \wedge \text{tr}(R_2, R_3) \wedge \text{rowSums}(N, R_4) \wedge \text{multi}_E(R_3, R_4, R_5), \text{sum}(R_5, r)$
$\text{colSums}(MN) \rightarrow \text{colSums}(M)N$	$\forall M, N, R_1, R_2 \text{ multi}_M(M, N, R_1) \wedge \text{colSums}(R_1, R_2) \rightarrow \exists R_3 \text{ colSums}(M, R_3) \wedge \text{multi}_M(R_3, N, R_2)$
$\text{rowSums}(MN) \rightarrow M\text{rowSums}(N)$	$\forall M, N, R_1, R_2 \text{ multi}_M(M, N, R_1) \wedge \text{rowSums}(R_1, R_2) \rightarrow \exists R_3 \text{ rowSums}(N, R_3) \wedge \text{multi}_M(M, R_3, R_2)$
$\text{colSums}(M) \rightarrow M$ if x is row vector	$\forall M, n, i \text{ name}(M, n) \wedge \text{size}(M, "1", j) \rightarrow \text{colSums}(M, M)$
$\text{colMean}(M) \rightarrow M$ if x is row vector	$\forall M, n, j \text{ name}(M, n) \wedge \text{size}(M, "1", j) \rightarrow \text{colSums}(M, M)$
$\text{colVar}(M) \rightarrow M$ if x is row vector	$\forall M, n, j \text{ name}(M, n) \wedge \text{size}(v, "1", j) \rightarrow \text{colVar}(M, M)$
$\text{colMax}(M) \rightarrow M$ if x is row vector	$\forall M, n, j \text{ name}(M, n) \wedge \text{size}(M, "1", j) \rightarrow \text{colMax}(M, M)$
$\text{colMin}(M) \rightarrow M$ if x is row vector	$\forall M, n, j \text{ name}(M, n) \wedge \text{size}(M, "1", j) \rightarrow \text{colMin}(M, M)$
$\text{colSums}(M) \rightarrow \text{sum}(M)$ if x is col vector	$\forall M, i, R_1 \text{ colSums}(M, R_1) \wedge \text{size}(M, i, "1") \rightarrow \text{sum}(M, R_1)$
$\text{colMean}(M) \rightarrow \text{mean}(M)$ if x is col vector	$\forall M, i, R_1 \text{ colMean}(M, R_1) \wedge \text{size}(M, i, "1") \rightarrow \text{mean}(M, R_1)$
$\text{colMax}(X) \rightarrow \text{max}(M)$ if x is col vector	$\forall M, i, R_1 \text{ colMax}(M, R_1) \wedge \text{size}(M, i, "1") \rightarrow \text{max}(M, R_1)$
$\text{colMin}(M) \rightarrow \text{min}(X)$ if x is col vector	$\forall M, i, R_1 \text{ colMin}(M, R_1) \wedge \text{size}(M, i, "1") \rightarrow \text{min}(M, R_1)$
$\text{colVar}(M) \rightarrow \text{var}(M)$ if x is col vector	$\forall M, i, R_1 \text{ colVar}(M, R_1) \wedge \text{size}(M, i, "1") \rightarrow \text{var}(M, R_1)$

SystemML Algebraic Simplification Rule	Integrity Constraints $\mathcal{MMC}_{StatAgg}$
simplifyRowWiseAgg	
rowSums(M) -> M if x is col vector	$\forall M, n, i \ name(M, n) \wedge size(M, i, "1") \rightarrow \text{rowSums}(M, M)$
rowMean(M) -> M if x is col vector	$\forall M, n, i \ name(M, n) \wedge size(M, i, "1") \rightarrow \text{rowMean}(M, M)$
rowVar(M) -> M if x is col vector	$\forall M, n, i \ name(M, n) \wedge size(M, i, "1") \rightarrow \text{rowVar}(M, M)$
rowMax(M) -> M if x is col vector	$\forall M, n, i \ name(M, n) \wedge size(M, i, "1") \rightarrow \text{rowMax}(M, M)$
rowMin(M) -> M if x is col vector	$\forall M, n, i \ name(M, n) \wedge size(M, i, "1") \rightarrow \text{rowMin}(M, M)$
rowSums(M) -> sum(M) if x is row vector	$\forall M, j, R_1 \ \text{rowSums}(M, R_1) \wedge size(M, "1", j) \rightarrow \text{sum}(M, R_1)$
rowMean(M) -> mean(M) if x is row vector	$\forall M, j, R_1 \ \text{rowMean}(M, R_1) \wedge size(M, "1", j) \rightarrow \text{mean}(M, R_1)$
rowMax(M) -> max(M) if x is row vector	$\forall M, j, R_1 \ \text{rowMax}(M, R_1) \wedge size(M, "1", j) \rightarrow \text{max}(M, R_1)$
rowMin(X) -> min(M) if x is row vector	$\forall M, j, R_1 \ \text{rowMin}(M, R_1) \wedge size(M, "1", j) \rightarrow \text{min}(M, R_1)$
rowVar(X) -> var(M) if x is row vector	$\forall M, j, R_1 \ \text{rowVar}(M, R_1) \wedge size(M, "1", j) \rightarrow \text{var}(M, R_1)$
pushdownSumOnAdd	
sum(M+N) -> sum(A)+sum(B)	$\forall M, N, s \ \text{add}_M M, N, s_1 \wedge \text{sum}(M, s_1) \rightarrow \exists s_2, s_3 \ \text{sum}(M, s_2) \wedge \text{sum}(N, s_3) \wedge \text{add}_s(s_2, s_3, s_1)$
ColSumsMVMMult	
colSums(M*N) -> t(M)N	$\forall M, N, R_1, R_2, i \ size(N, i, "1") \wedge \text{multi}_E(M, N, R_1) \wedge \text{colSums}(R_1, R_2) \ \exists R_3 \ \text{tr}(M, R_3) \wedge \text{multi}_M(R_3, N, R_2)$
rowSums(M*M) -> Mt(N)	$\forall M, N, R_1, R_2, j \ size(N, "1", j) \wedge \text{multi}_E(M, N, R_1) \wedge \text{rowSums}(R_1, R_2) \ \exists R_3 \ \text{tr}(N, R_3) \wedge \text{multi}_M(M, R_3, R_2)$

C \mathcal{P}^{-Opt} AND \mathcal{P}^{Views} PIPELINES REWRITES

No.	Rewrite	No.	Rewrite	No.	Rewrite
P1.1	$N^T M^T$	P1.2	$(A + B)^T$	P1.3	$(DC)^{-1}$
P1.4	$Av_1 + Bv_1$	P1.5	D	P1.6	$s_1 \text{trace}(D)$
P1.7	A	P1.8	$(s_1 + s_2)A$	P1.9	$\det(D)$
P1.10	$\text{colSums}(A)^T$	P1.11	$\text{colSums}(A + B)^T$	P1.12	$\text{colSums}(M)N$
P1.13	$\text{sum}(\text{colSums}(M)^T * \text{rowSums}(N))$	P1.14	$\text{sum}(\text{colSums}(M)^T * \text{rowSums}(N))$	P1.15	$M(NM)$
P1.16	$\text{sum}(A)$	P1.17	$\det(C) * \det(D) * \det(C)$	P1.18	$\text{sum}(A)$
P1.25	$M \odot (N^T / (M(NN^T)))$				

Table 12: \mathcal{P}^{-Opt} Pipelines (Part 1) Rewrites

No.	Rewrite	No.	Rewrite	No.	Rewrite
P2.1	$\text{trace}(C) + \text{trace}(D)$	P2.2	$1/\det(D)$	P2.3	$\text{trace}(D)$
P2.4	$s_1(A + B)$	P2.5	$1/\det((C + D))$	P2.6	$(D^{-1}C)^T$
P2.7	C	P2.8	$\det(C) * \det(D)$	P2.9	$\text{trace}(DC) + \text{trace}(D)$
P2.10	$M\text{rowSums}N$	P2.11	$\text{sum}(A) + \text{sum}(B)$	P2.12	$\text{sum}(\text{colSums}(M)^T * \text{rowSums}(N))$
P2.13	$(M(NM))^T$	P2.14	$(M(NM))N$	P2.15	$\text{sum}(A)$
P2.16	$\text{trace}((DC)^{-1}) + \text{trace}D$	P2.17	$(((C + D)^{-1})^T)D$	P2.18	$\text{rowSums}(A + B)^T$
P2.25	$u_1 v_2^T v_2 - X v_2$				

Table 13: \mathcal{P}^{-Opt} Pipelines (Part 2) Rewrites

No.	Expression	No.	Expression	No.	Expression
V_1	$(D)^{-1}$	V_2	$(C^T)^{-1}$	V_3	NM
V_4	$u_1 v_2^T$	V_5	DC	V_6	$A + B$
V_7	C^{-1}	V_8	$C^T D$	V_9	$(D + C)^{-1}$
V_{10}	$\det(CD)$	V_{11}	$\det(DC)$	V_{12}	$(DC)^T$

Table 14: The set of views V_{exp}

No.	Rewrite	No.	Rewrite	No.	Rewrite
P1.2	$(V_6)^T$	P1.3	$V_7 V_1$	P1.4	$(V_6)v_1$
P1.11	$\text{colSums}(V_6)^T$	P1.15	$M(V_3)$	P1.17	$V_{10} * \det(C)$
P1.19	V_2	P1.20	$\text{trace}(V_7)$	P1.21	$(C + V_1)^T$
P1.22	$\text{trace}(V_9)$	P1.24	$\text{trace}(V_1 V_7) + \text{trace}(D)$	P1.29	$V_5 CCC$
P1.30	$V_3 \odot V_3 R^T$	P2.2	$\det(V_1)$	P2.4	$s_1(V_6)$
P2.5	$\det(V_9)$	P2.6	$(V_1 C)^T$	P2.9	$\text{trace}(V_{12}) + \text{trace}(D)$
P2.11	$\text{sum}(V_6)$	P2.13	$(MV_3)^T$	P2.14	$MV_3 N$
P2.16	$\text{trace}(V_7 V_1) + \text{trace}D$	P2.17	$(V_9^T)D$	P2.18	$\text{rowSums}(V_6)^T$
P2.20	$(MV_3)^T$	P2.21	$V_1(V_1^T (D^T v_1))$	P2.25	$V_4 v_1 - X v_1$
P1.23	$\det((V_7 V_1) + D)$	P2.26	$\exp(V_9)$	P2.27	$V_9^T V_5$

Table 15: \mathcal{P}^{Views} Pipelines Rewrites

D ADDITIONAL RESULTS: \mathcal{P}^{Opt} PIPELINES - NAÏVE-BASED COST MODEL

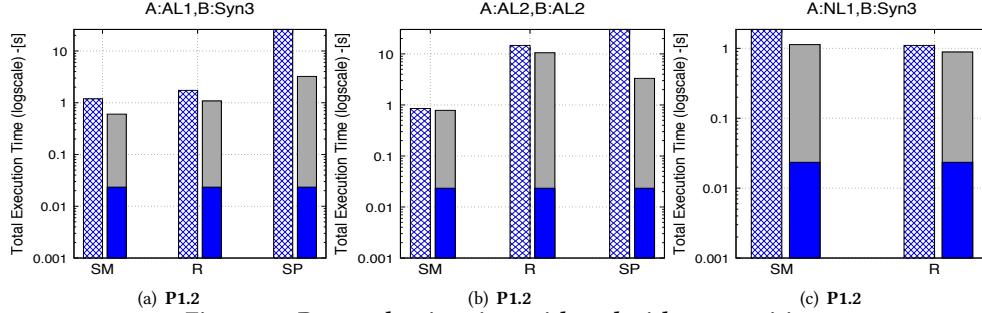


Figure 13: P1.2 evaluation time with and without rewriting

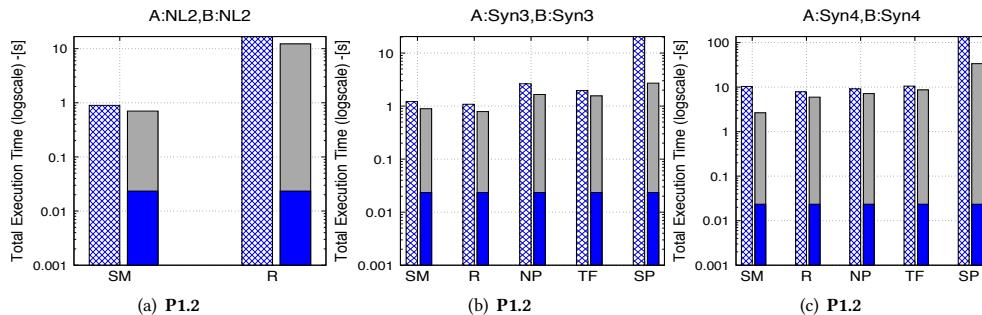


Figure 14: P1.2 evaluation time with and without rewriting

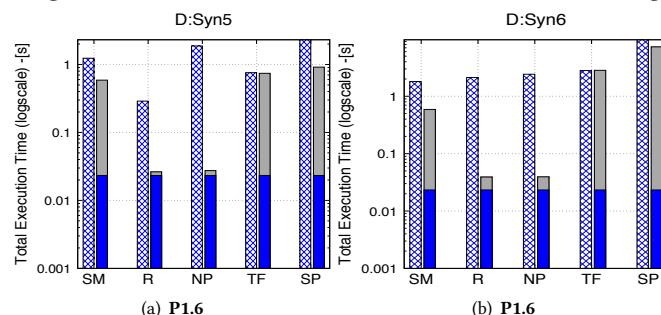


Figure 15: P1.6 evaluation time with and without rewriting

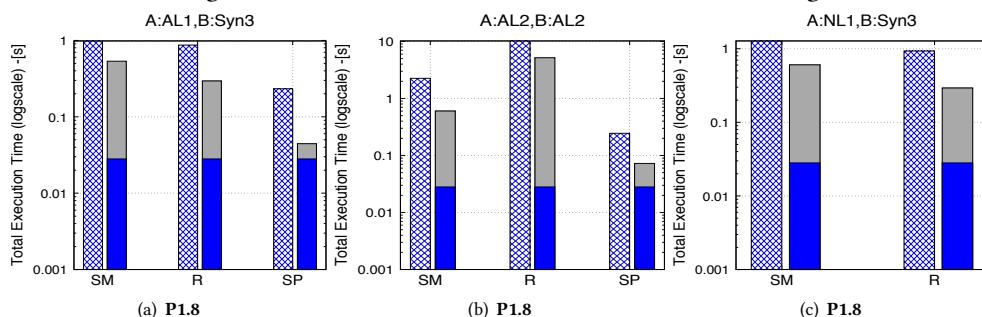
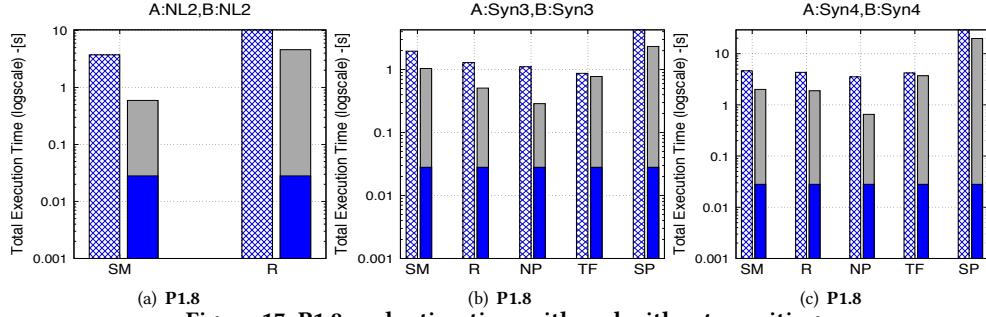
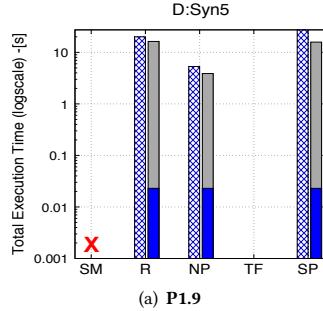
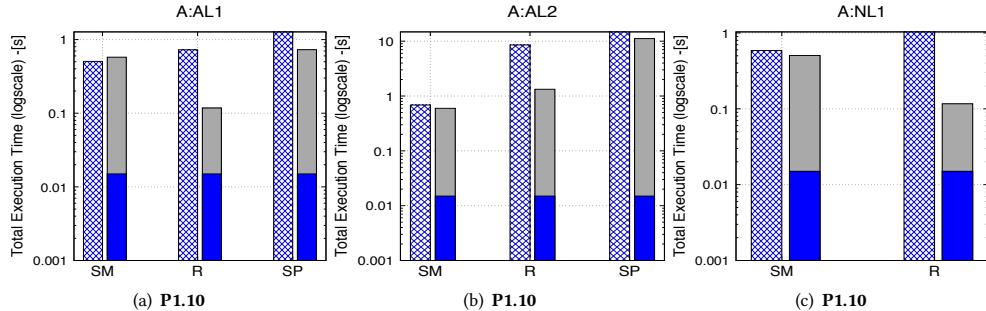
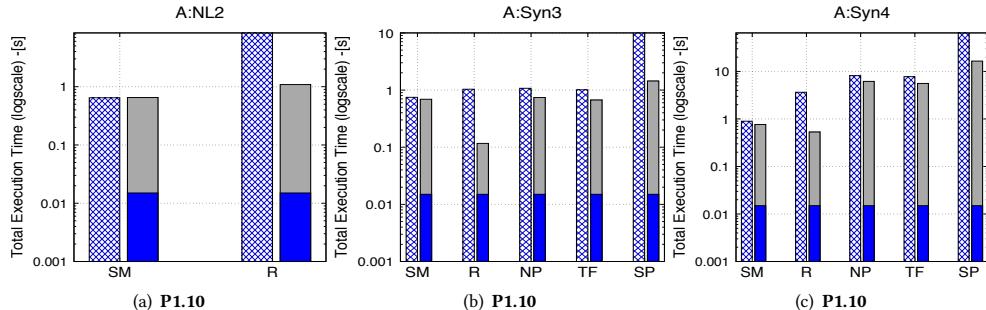
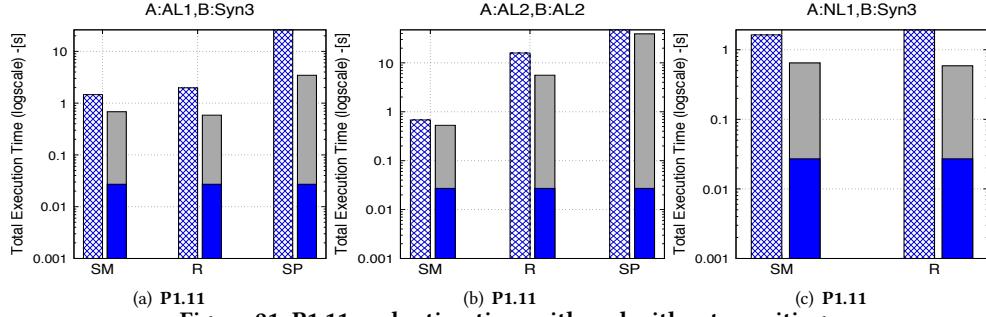
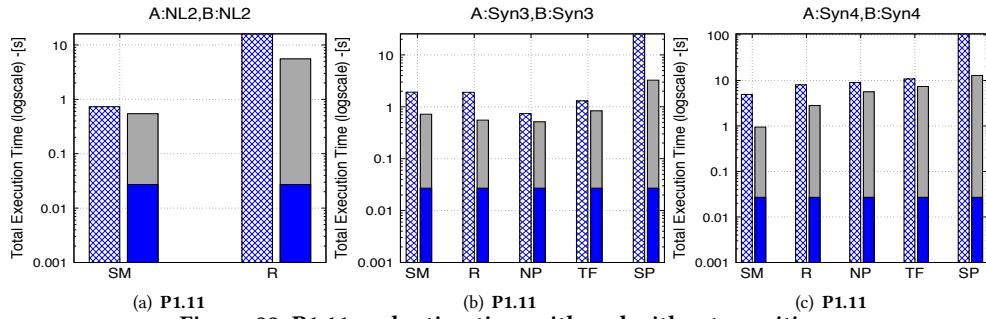
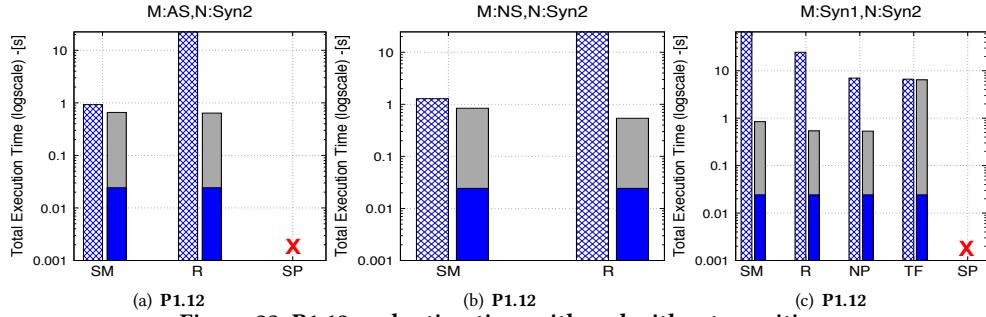
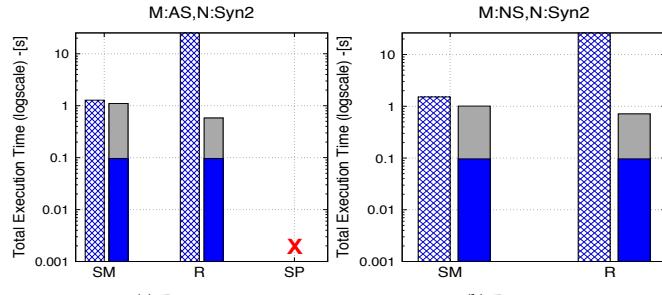
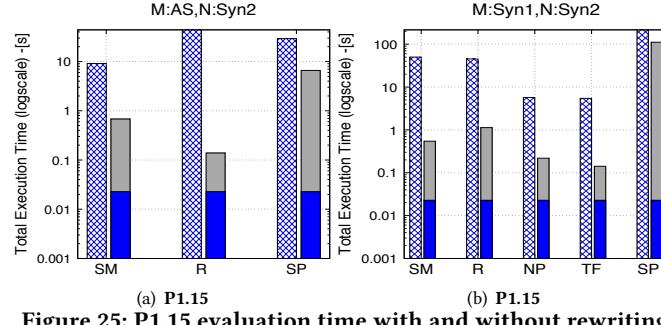
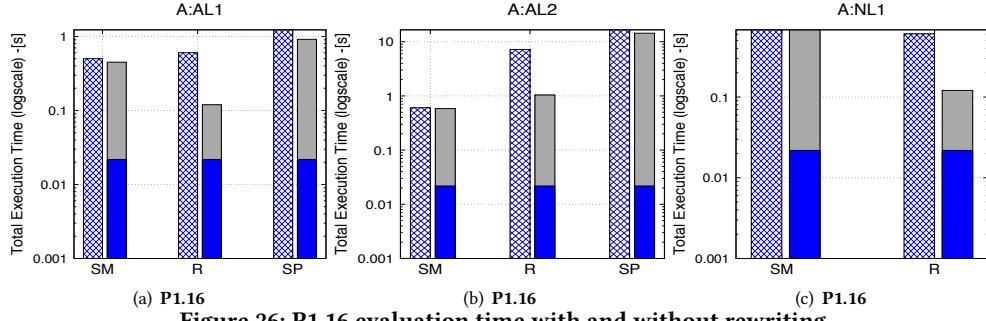
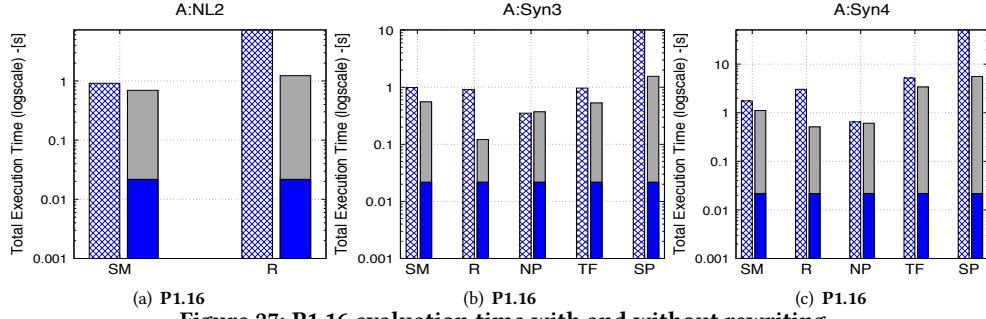
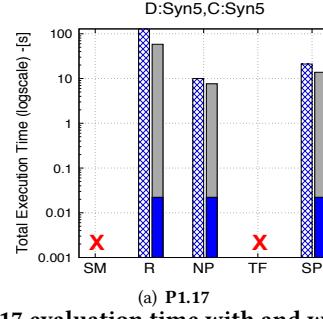
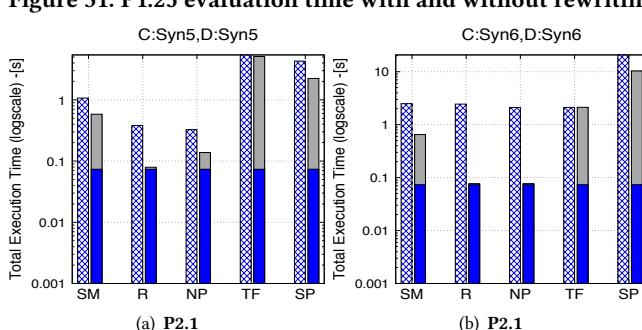
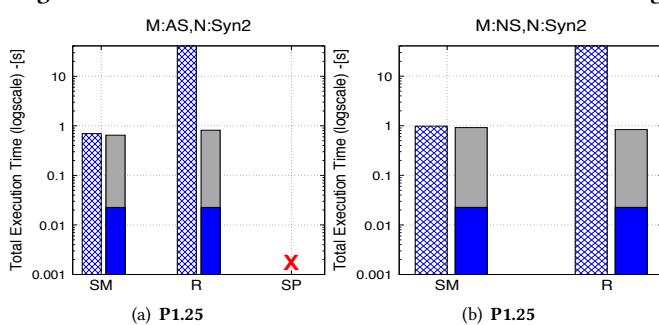
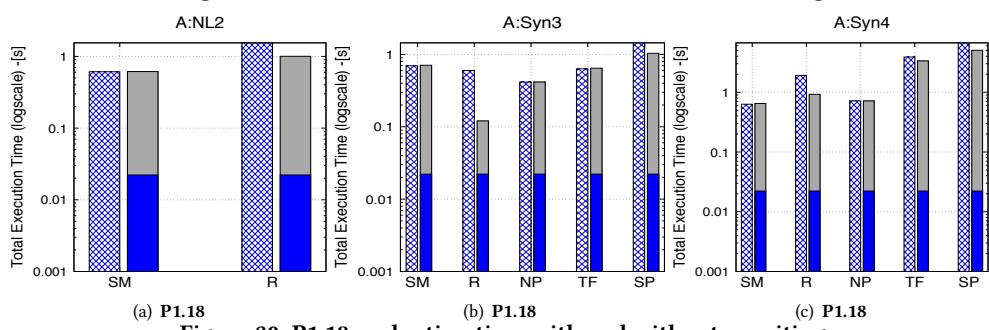
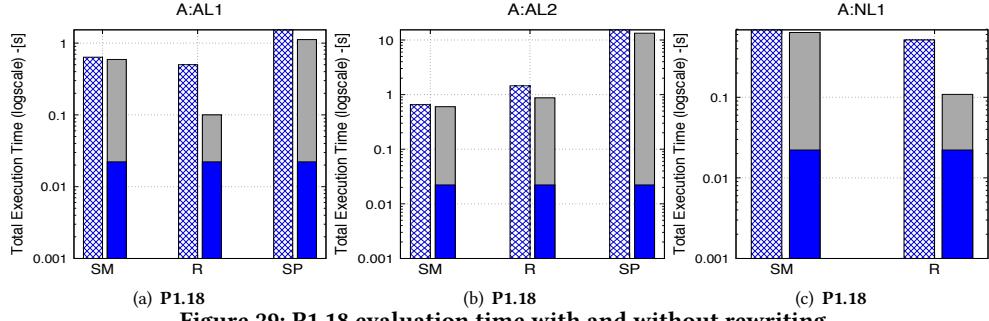


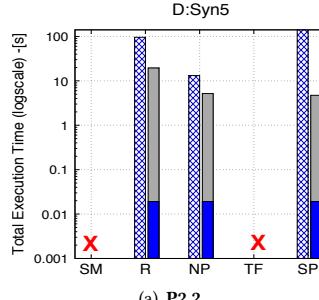
Figure 16: P1.8 evaluation time with and without rewriting

**Figure 17: P1.8 evaluation time with and without rewriting****Figure 18: P1.9 evaluation time with and without rewriting****Figure 19: P1.10 evaluation time with and without rewriting****Figure 20: P1.10 evaluation time with and without rewriting**

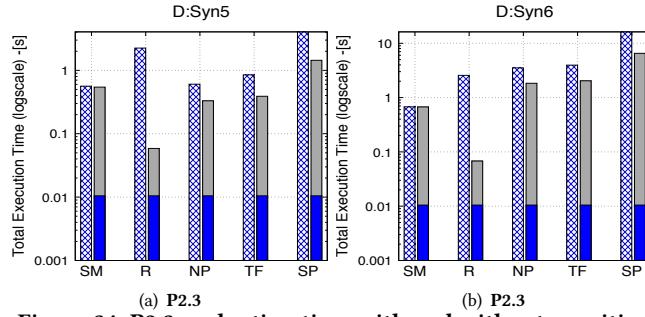
**Figure 21: P1.11 evaluation time with and without rewriting****Figure 22: P1.11 evaluation time with and without rewriting****Figure 23: P1.12 evaluation time with and without rewriting****Figure 24: P1.14 evaluation time with and without rewriting**

**Figure 25: P1.15 evaluation time with and without rewriting****Figure 26: P1.16 evaluation time with and without rewriting****Figure 27: P1.16 evaluation time with and without rewriting****Figure 28: P1.17 evaluation time with and without rewriting**



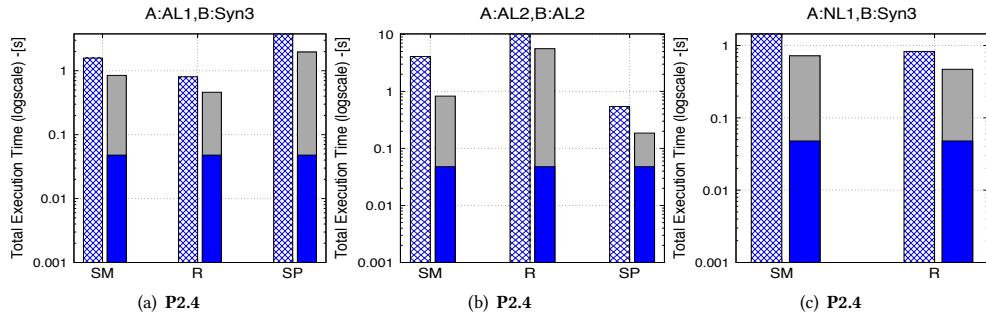


(a) P2.2

Figure 33: P2.2 evaluation time with and without rewriting

(a) P2.3

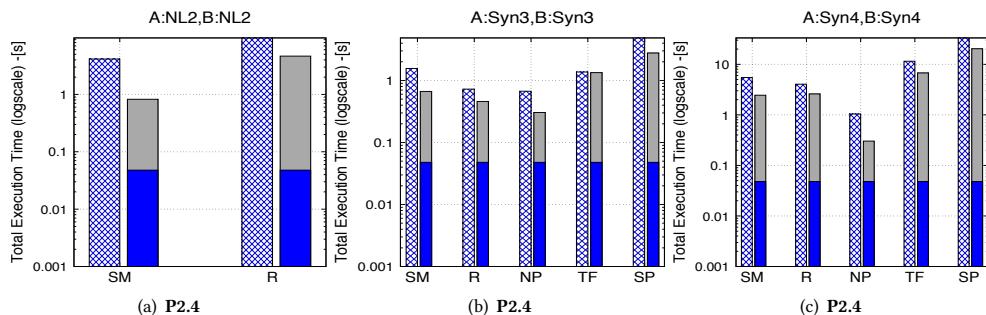
(b) P2.3

Figure 34: P2.3 evaluation time with and without rewriting

(a) P2.4

(b) P2.4

(c) P2.4

Figure 35: P2.4 evaluation time with and without rewriting

(a) P2.4

(b) P2.4

(c) P2.4

Figure 36: P2.4 evaluation time with and without rewriting

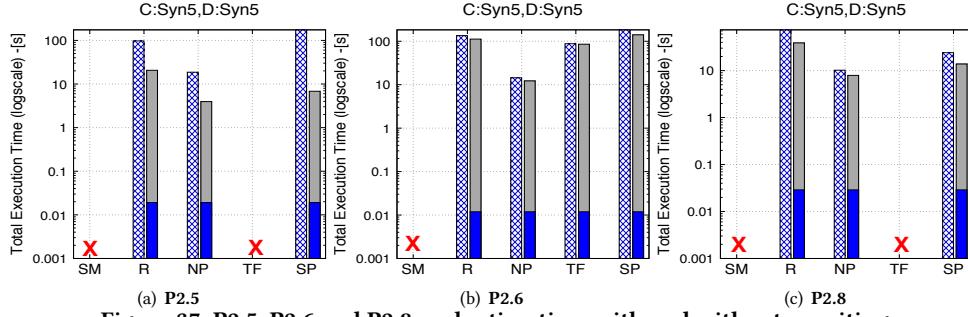


Figure 37: P2.5, P2.6 and P2.8 evaluation time with and without rewriting

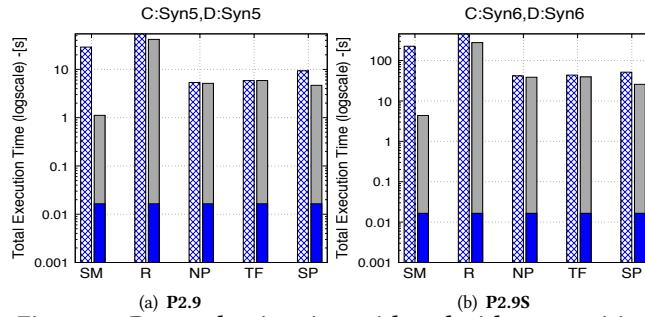


Figure 38: P2.9 evaluation time with and without rewriting

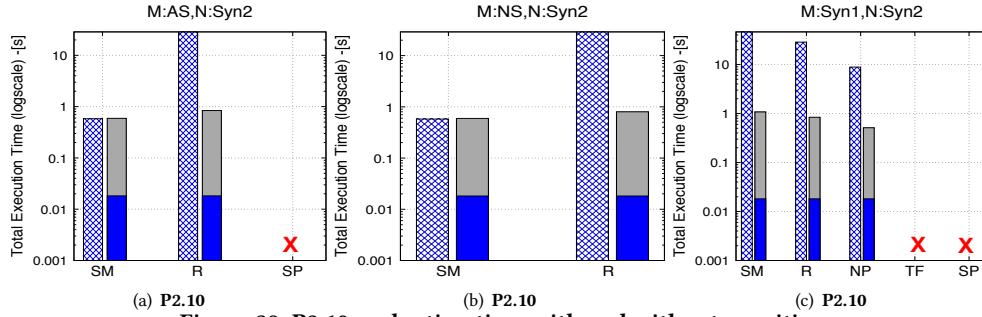


Figure 39: P2.10 evaluation time with and without rewriting

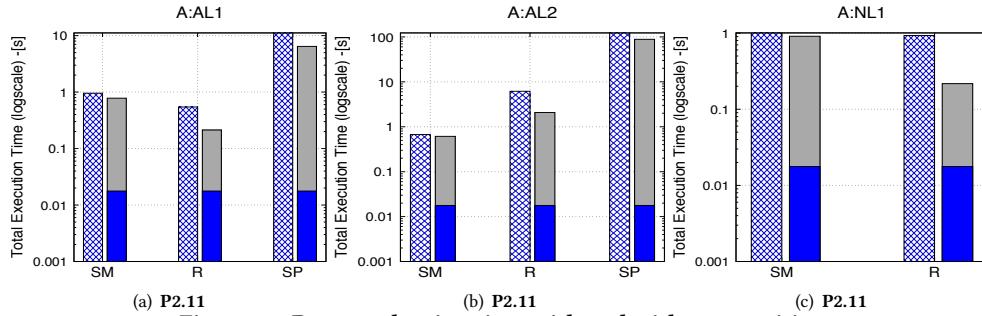
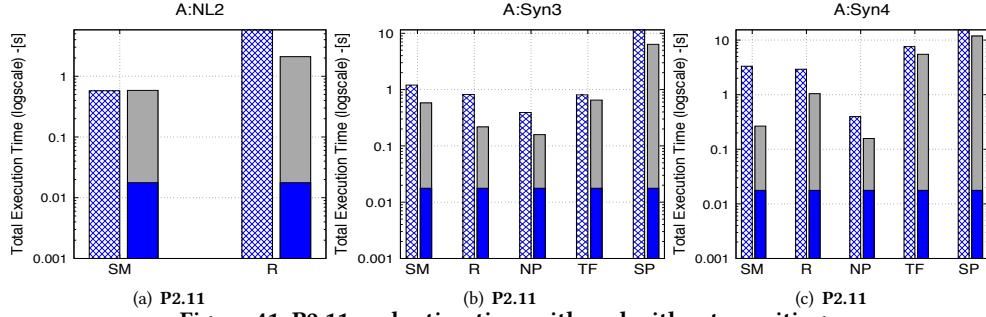
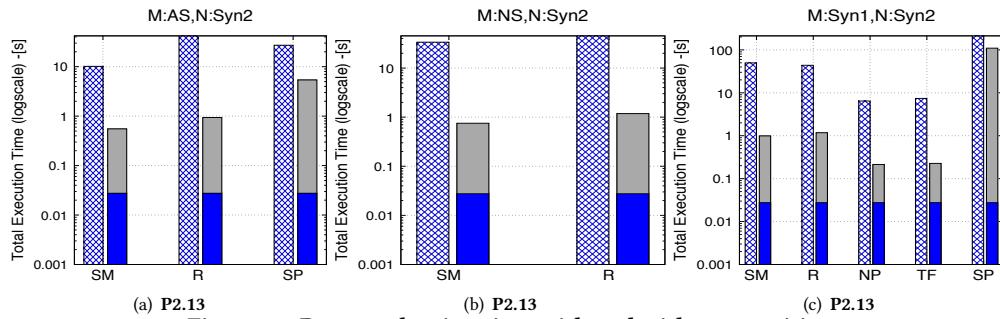
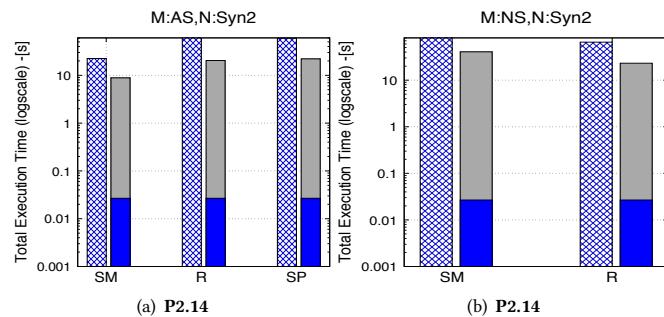
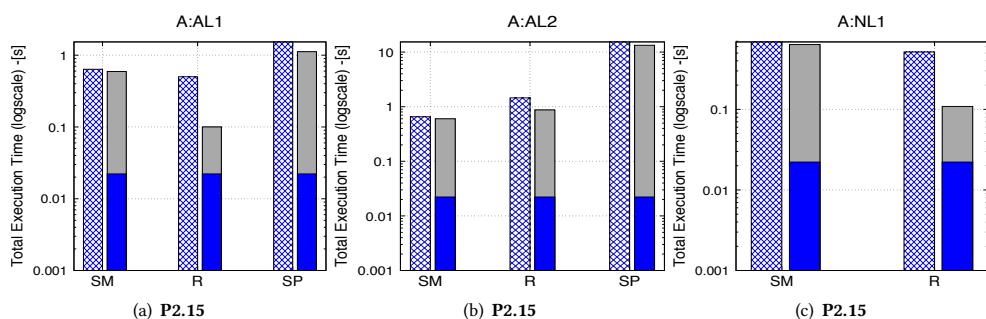
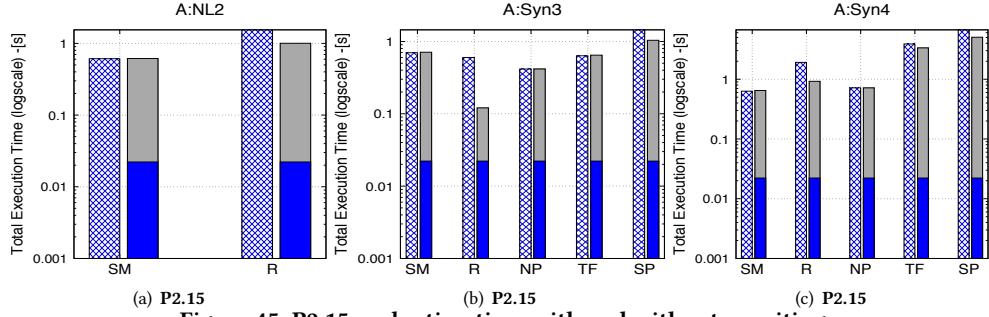
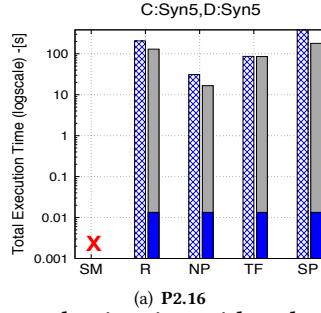
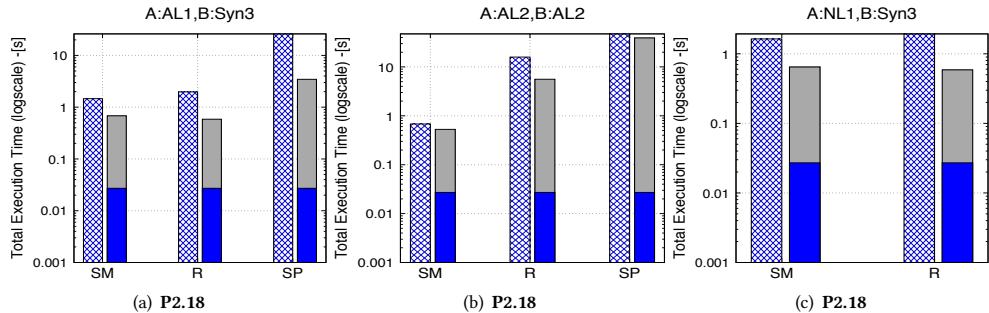
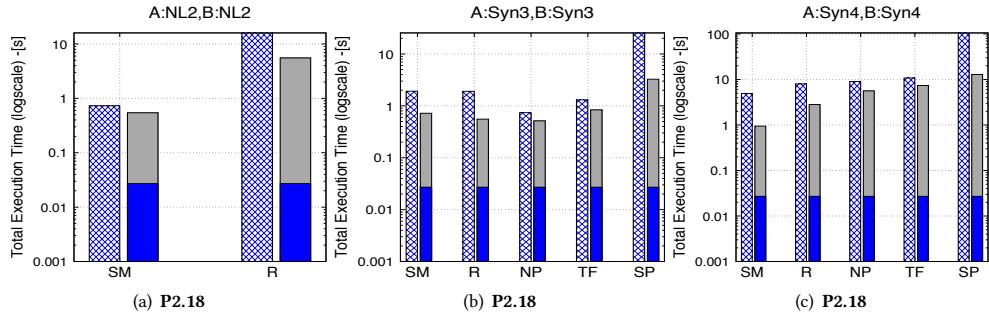


Figure 40: P2.11 evaluation time with and without rewriting

**Figure 41: P2.11 evaluation time with and without rewriting****Figure 42: P2.13 evaluation time with and without rewriting****Figure 43: P2.14 evaluation time with and without rewriting****Figure 44: P2.15 evaluation time with and without rewriting**

**Figure 45: P2.15 evaluation time with and without rewriting****Figure 46: P2.16 evaluation time with and without rewriting****Figure 47: P2.18 evaluation time with and without rewriting****Figure 48: P2.18 evaluation time with and without rewriting**

E ADDITIONAL RESULTS: $\mathcal{P}^{\neg Opt}$ PIPELINES - MNC-BASED COST MODEL

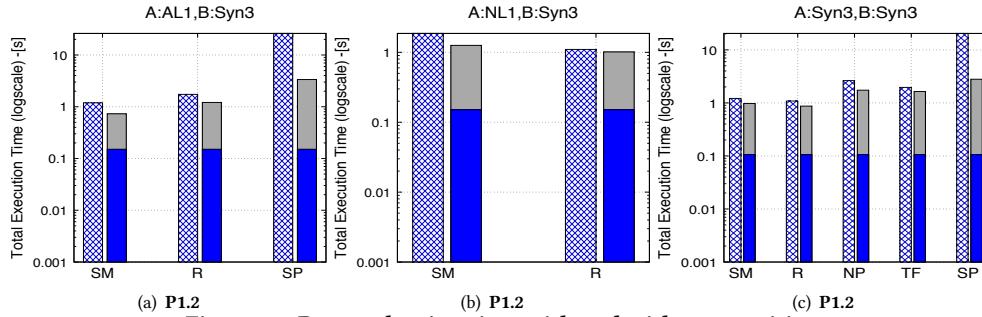


Figure 49: P1.2 evaluation time with and without rewriting

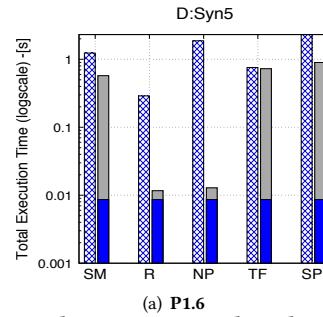


Figure 50: P1.6 evaluation time with and without rewriting

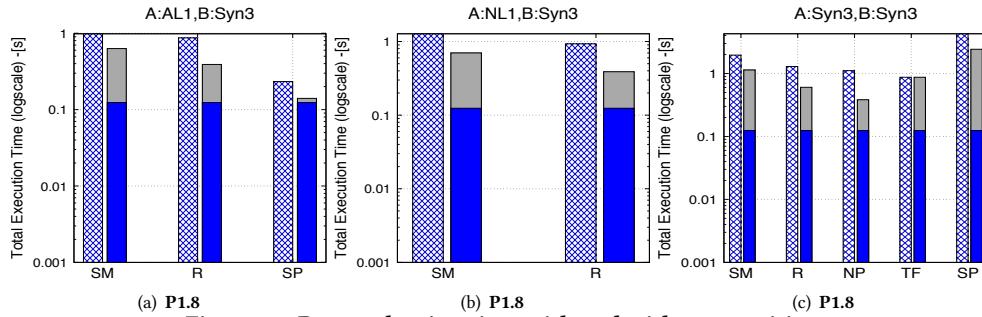


Figure 51: P1.8 evaluation time with and without rewriting

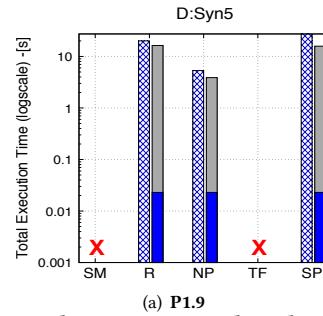
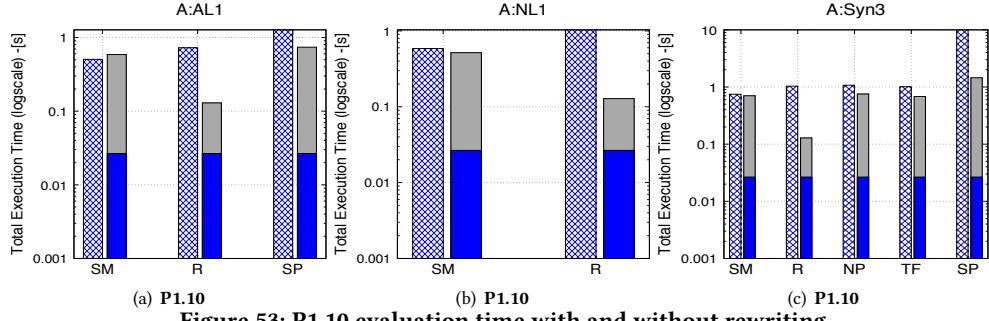
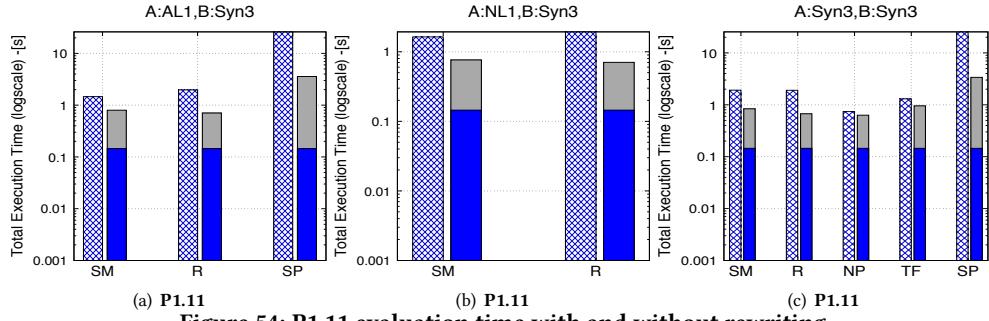
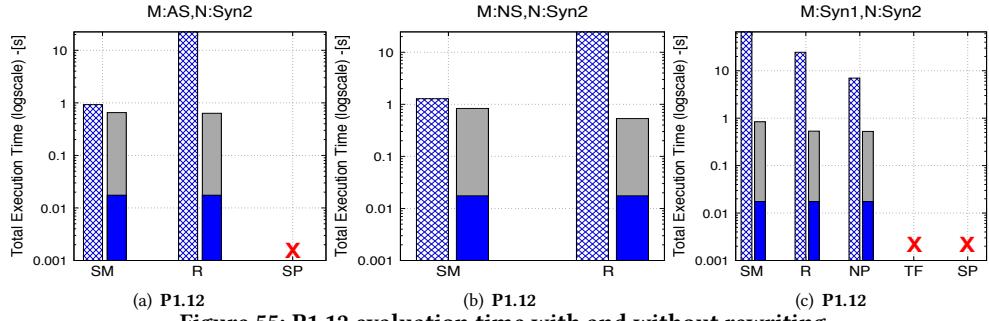
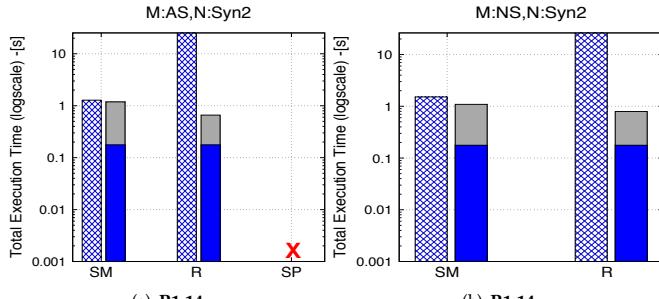
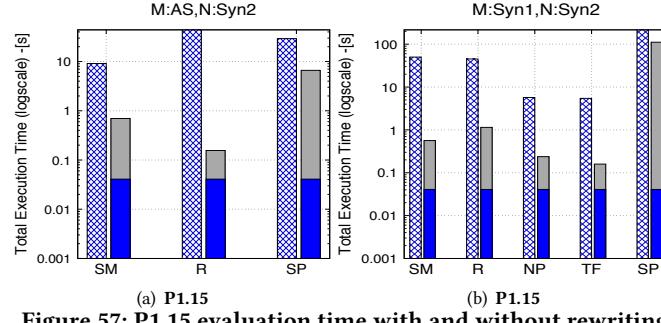
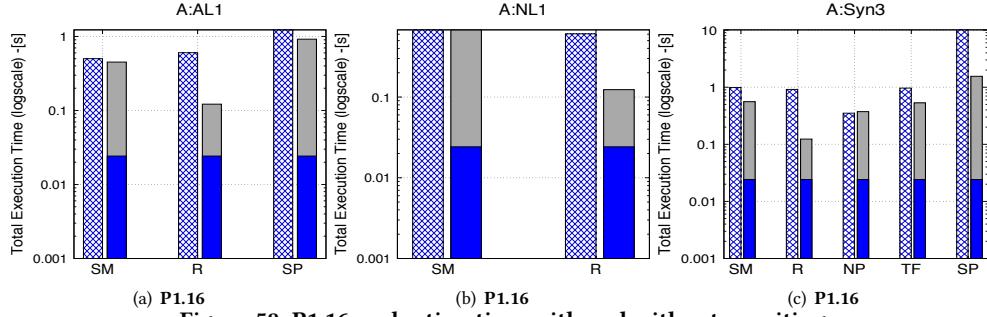
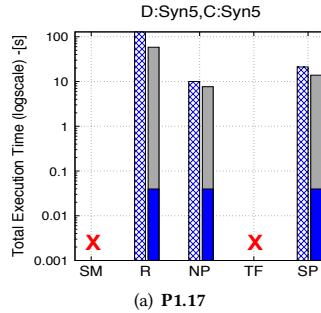
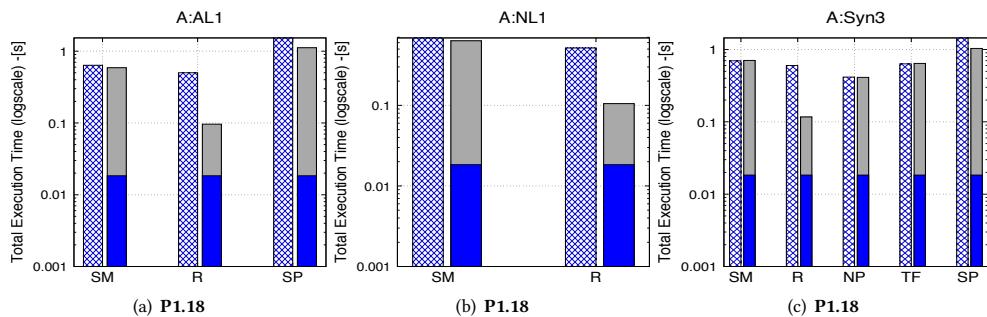
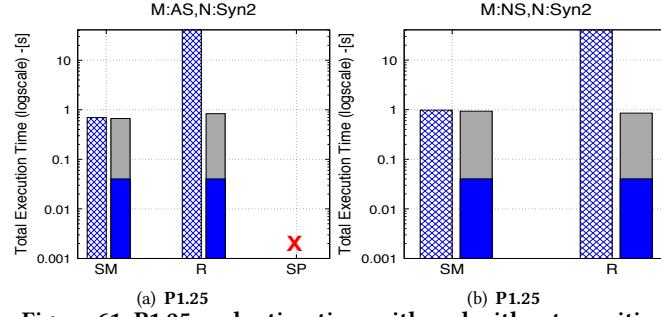
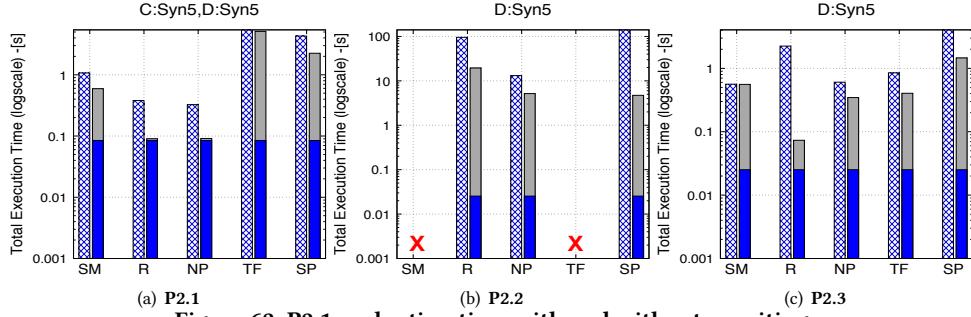
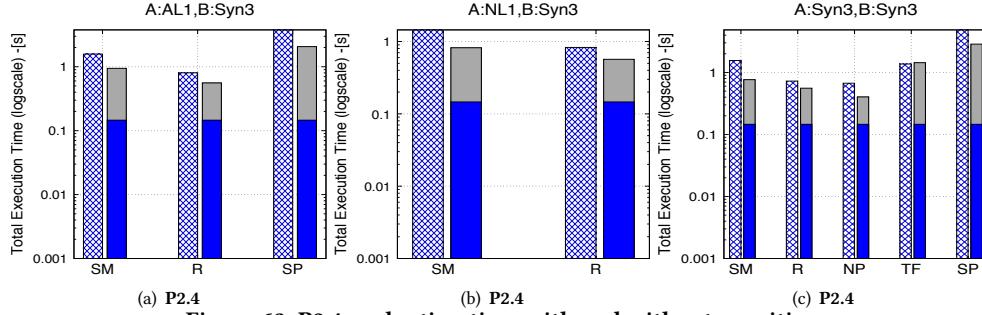
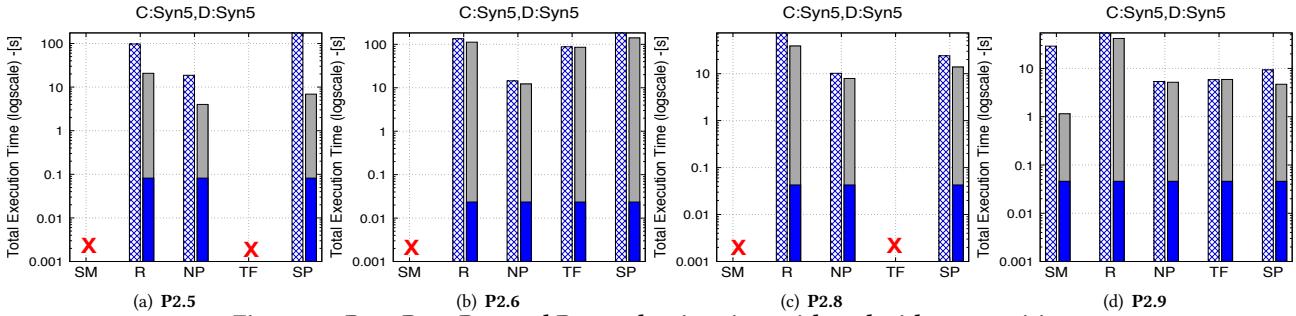
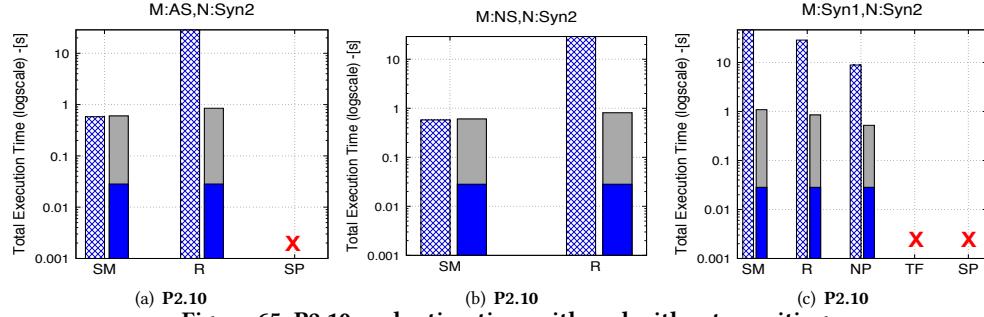
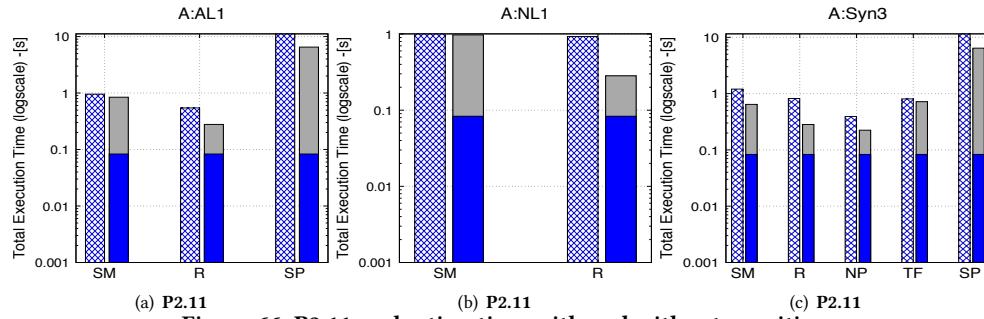
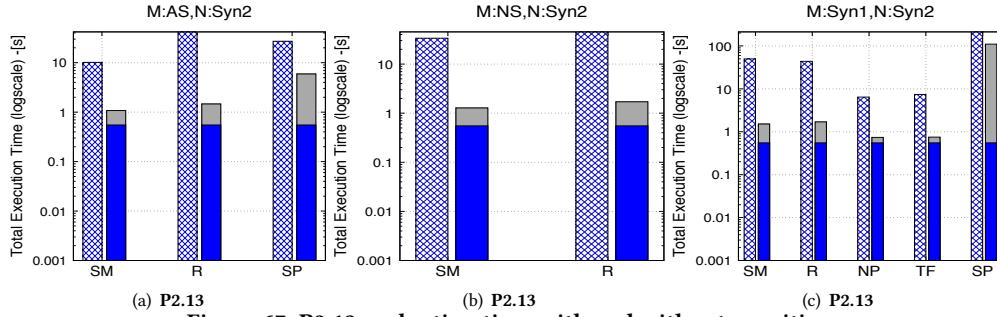
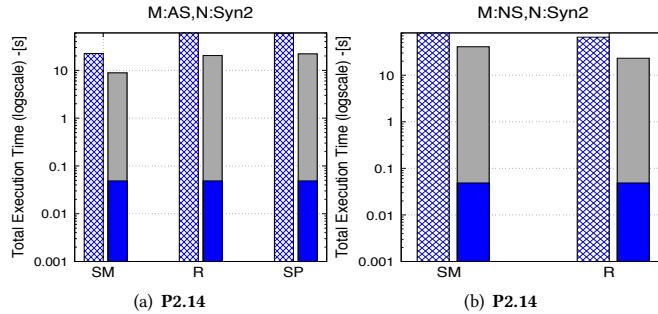


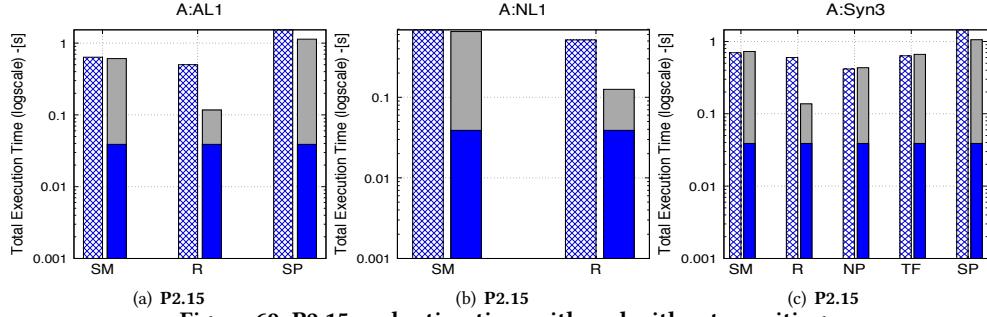
Figure 52: P1.9 evaluation time with and without rewriting

**Figure 53: P1.10 evaluation time with and without rewriting****Figure 54: P1.11 evaluation time with and without rewriting****Figure 55: P1.12 evaluation time with and without rewriting****Figure 56: P1.14 evaluation time with and without rewriting**

**Figure 57: P1.15 evaluation time with and without rewriting****Figure 58: P1.16 evaluation time with and without rewriting****Figure 59: P1.17 evaluation time with and without rewriting****Figure 60: P1.18 evaluation time with and without rewriting**

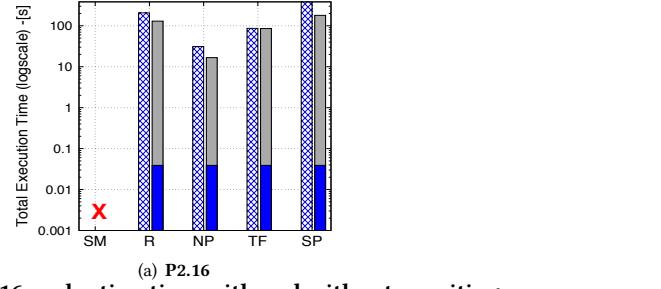
**Figure 61: P1.25 evaluation time with and without rewriting****Figure 62: P2.1 evaluation time with and without rewriting****Figure 63: P2.4 evaluation time with and without rewriting****Figure 64: P2.5, P2.6, P2.8 and P2.9 evaluation time with and without rewriting**

**Figure 65: P2.10 evaluation time with and without rewriting****Figure 66: P2.11 evaluation time with and without rewriting****Figure 67: P2.13 evaluation time with and without rewriting****Figure 68: P2.14 evaluation time with and without rewriting**



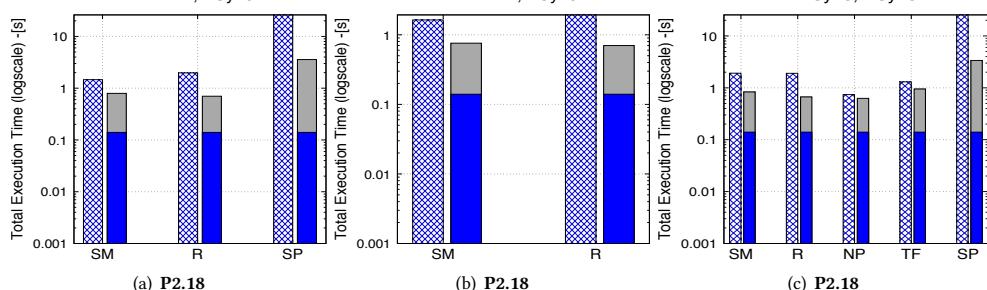
(a) P2.15 (b) P2.15 (c) P2.15

Figure 69: P2.15 evaluation time with and without rewriting



(a) P2.16

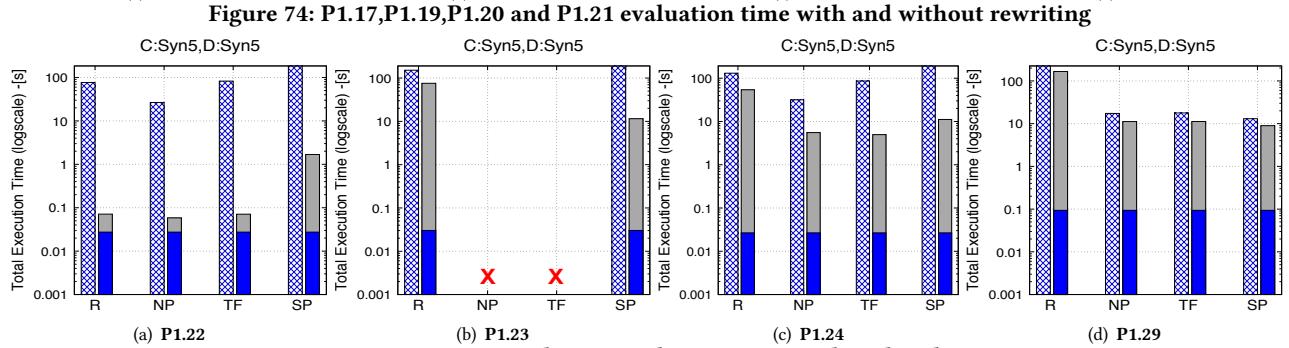
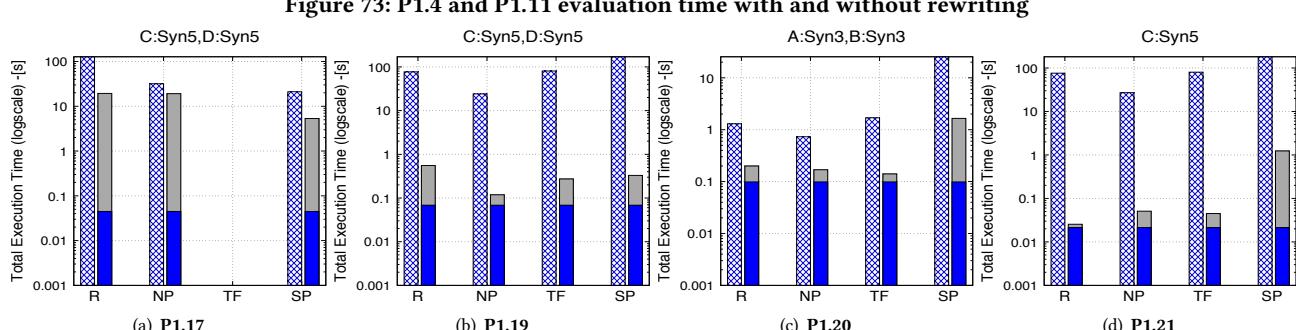
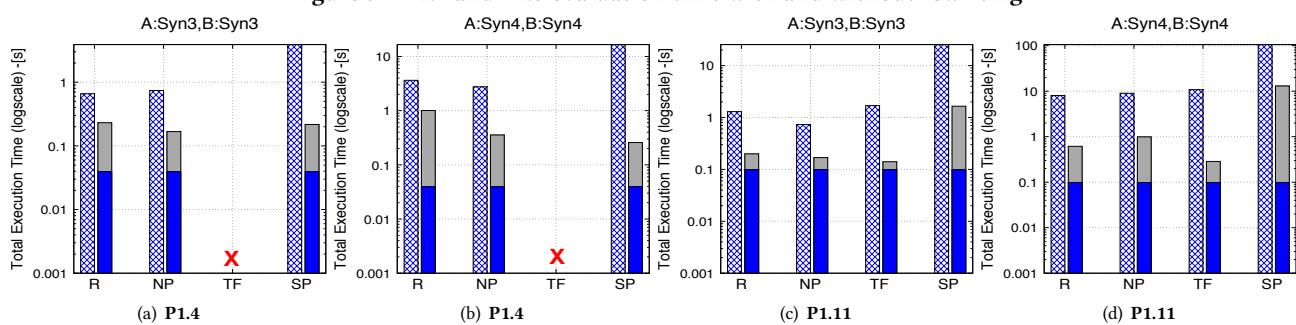
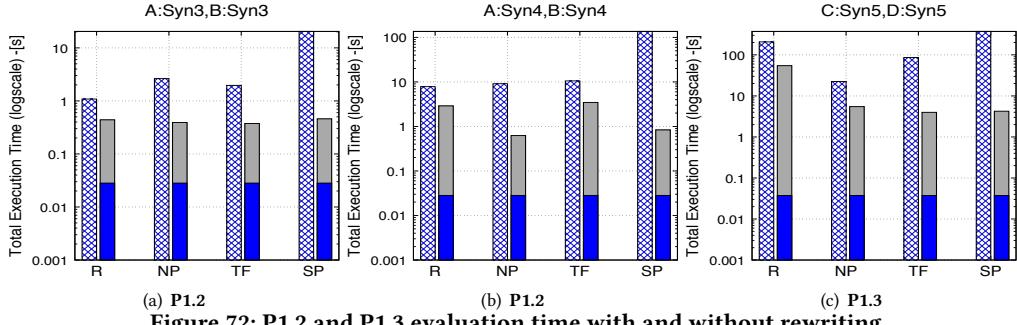
Figure 70: P2.16 evaluation time with and without rewriting

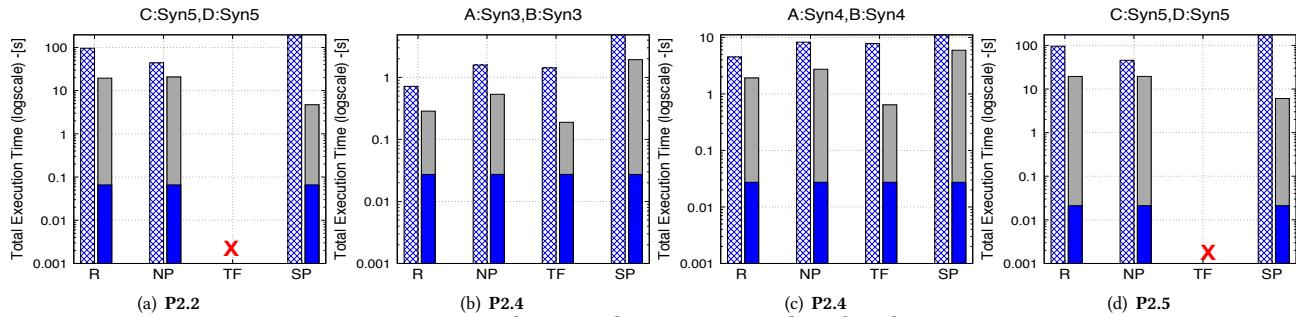
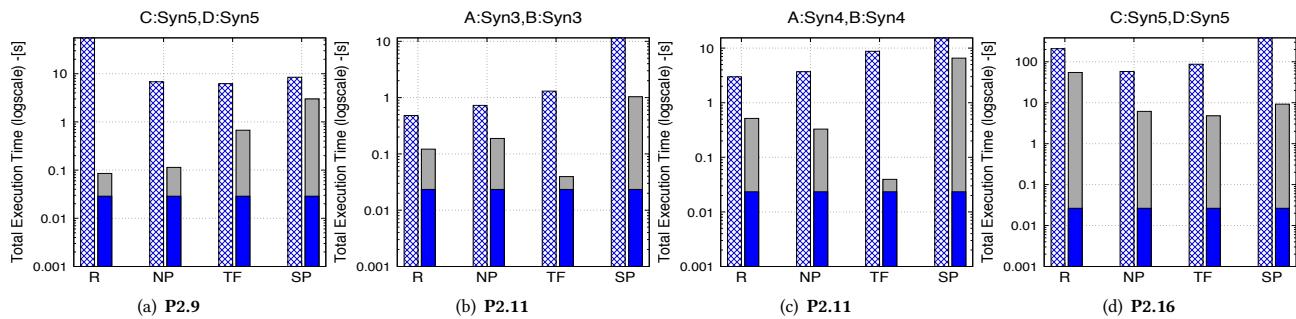


(a) P2.18 (b) P2.18 (c) P2.18

Figure 71: P2.18 evaluation time with and without rewriting

F ADDITIONAL RESULTS: \mathcal{P}^{Views} PIPELINES



**Figure 76: P2.2,P2.4 and P2.5 evaluation time with and without rewriting****Figure 77: P2.9,P2.11 and P2.16 evaluation time with and without rewriting**

G HYBRID TWITTER BENCHMARK CONSTRAINTS

```

1 #EGDs
2 name(id1, n1), name(id2, n1) -> id2 = id1;
3 tr(id1, idt1), tr(id1, idt2) -> idt2 = idt1;
4 in(id1, idt1), in(id1, idt2) -> idt2 = idt1;
5 ins(id1, idt1), ins(id1, idt2) -> idt1 = idt2;
6 in(id1, idr1), in(idr1, idr3) -> idr3 = id1;
7 multi(idA, idB, idr1), multi(idA, idB, idr2) -> idr2 = idr1;
8 multiS(idA, idB, idr1), multiS(idA, idB, idr2) -> idr2 = idr1;
9 det(idA, d1), det(idA, d2) -> d2 = d1;
10 trace(idA, id1), trace(idA, id2) -> id2 = id1;
11 add(idA, idB, idr1), add(idA, idB, idr2) -> idr1 = idr2;
12 size(idA, i1, j1), size(idA, i2, j2) -> j2 = j1, i2 = i1;
13 name(id1, n1), name(id2, n1) -> id2 = id1;
14 name(mID1, n1), name(mID1, n2) -> n1 = n2;
15 adds(idA, idB, idr1), adds(idA, idB, idr2) -> idr1 = idr2;
16 sum(id1, idt1), sum(id1, idt2) -> idt2 = idt1;
17 rowSums(idA, id1), rowSums(idA, id2) -> id2 = id1;
18 colSums(idA, id1), colSums(idA, id2) -> id2 = id1;
19 user(EID, FOLC1, FRIC1, LC1, P1, V1), user(EID, FOLC2, FRIC2, LC2, P2, V2) -> FOLC2 = FOLC1, FRIC2 = FRIC1, LC2 = LC1, P2 = P1, V2 = V1;
20 constructMatrix2(rID, mID, v11, v12), constructMatrix2(rID, mID, v21, v22) -> v21 = v11, v22 = v12;
21 constructMatrix4(rID, mID, v11, v12, v13, v14), constructMatrix4(rID, mID, v21, v22, v23, v24) -> v21 = v11, v22 = v12, v23 = v13, v24 = v14;
22 constructMatrix11(rID, mID, v11, v12, v13, v14, v15, v16, v17, v18, v19, v110, v211), constructMatrix11(rID, mID, v21, v22, v23, v24, v25, v26, v27, v28, v29, v210, v211) ->
v11 = v21, v12 = v22, v13 = v23, v14 = v24, v15 = v25, v16 = v26, v17 = v27, v18 = v28, v19 = v29, v110 = v210, v111 = v211;
23 constructMtxMatrix(i, j, v, id)(i, j1, v1, id), constructMtxMatrix(j2, j2, v2, id) -> i2 = i1, j2 = j1, v2 = v1;
24 indicator(id1, id2, idr31), indicator(id1, id2, idr32) -> idr32 = idr31;
25 deriveFrom(id1, id2, idr31), deriveFrom(id1, id2, idr32) -> idr32 = idr31;
26 concat(id1, id2, idr1), concat(id1, id2, idr2) -> idr2 = idr1;
27 subsettweeT(id1, id2, id3, id5, id6, id7, idr1),
subsettweeT(id1, id2, id3, id4, id5, id6, id7, idr2) -> idr2 = idr1;
28 subsetuseT(id1, id2, id3, id4, id5, id6, idr1),
subsetuseT(id1, id2, id3, id4, id5, id6, idr2) -> idr2 = idr1;
29 tweet(id1), tweet(id2) -> id2 = id1;
30 childsp(x, y1, z, d), childsp(x, y2, z, d) -> y2 = y1;
31 valsP(x, v1), valsP(x, v2) -> v2 = v1;
32 childsjV2(x, y1, z, d), childsjV2(x, y2, z, d) -> y2 = y1;
33
34 #TGDs
35 add(idA, idB, idr1) -> add(idB, idA, idr1);
36 multi(id1, id2, idr1), rowSums(id1, idr2) -> rowSums(id2, idr3), multi(id1, idr3, idr2);
37 add(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
38 multi(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idB, idC, idr3), multi(idA, idr3, idr2);
39 add(idA, idB, idr1), tr(idr1, idT) -> tr(idA, idAT), tr(idB, idBT), add(idAT, idBT, idT);
40 tr(idA, idB, idAT), multi(idB, idAT, idr2) -> multi(idA, idB, idr1), tr(idr1, idr2);
41 constructMtxMatrix(i, j, v, id), name(id, n), matrixFilterLT(id, x, idr1) ->
name(idr1, n), constructMtxMatrix(i, j, v, idr1), LT(v, x);
42
43 subsettweeT(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
subsetuseT(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
44 indicator(mID2, mID3, idK), deriveFrom(idS1, idR1, id), rowSums(id, idr1) ->
rowSums(idS1, idr2), rowSums(idR1, idr3), multi(idK, idr3, idr4), add(idr2, idr4, idr1);
45
46 user(EID, FOLC, FRIC, LC, P, V),
47 tweetR(EID, FC, QC, RC, F, PS, R),
48 constructMatrix11(mID1, rWID, FOLC, FRIC, LC, P, V, FC, QC, RC, F, PS, R) ->
name(mID1, n),
name(mID2, "tweet"),
name(mID3, "user");
49 subsettweeT(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
subsetuseT(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
50 indicator(mID2, mID3, idK), deriveFrom(idS1, idR1, id), rowSums(id, idr1) ->
rowSums(idS1, idr2), rowSums(idR1, idr3), multi(idK, idr3, idr4), add(idr2, idr4, idr1);
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86

```

Listing G1: Forward Constraints Used for Q1

```

1 #EGDs
2 name(id1, n1), name(id2, n1) -> id2 = id1;
3 tr(id1, idt1), tr(id1, idt2) -> idt2 = idt1;
4 in(id1, idt1), in(id1, idt2) -> idt2 = idt1;
5 ins(id1, idt1), ins(id1, idt2) -> idt1 = idt2;
6 in(id1, idr1), in(idr1, idr3) -> idr3 = id1;
7 multi(idA, idB, idr1), multi(idA, idB, idr2) -> idr2 = idr1;
8 multiS(idA, idB, idr1), multiS(idA, idB, idr2) -> idr2 = idr1;
9 det(idA, d1), det(idA, d2) -> d2 = d1;
10 trace(idA, id), trace(idA, id2) -> id2 = id1;
11 add(idA, idB, idr1), add(idA, idB, idr2) -> idr1 = idr2;
12 size(idA, i1, j1), size(idA, i2, j2) -> j2 = j1, i2 = i1;
13 name(id1, n1), name(id2, n1) -> id2 = id1;
14 name(mID1, n1), name(mID1, n2) -> n1 = n2;
15 addS(idA, idB, idr1), addS(idA, idB, idr2) -> idr1 = idr2;
16 sum(id1, idt1), sum(id1, idt2) -> idt2 = idt1;
17 rowSums(idA, id1), rowSums(idA, id2) -> id2 = id1;
18 colSums(idA, id1), colSums(idA, id2) -> id2 = id1;
19 user(EID, FOLC1, FRIC1, LC1, P1, V1), user(EID, FOLC2, FRIC2, LC2, P2, V2) -> FOLC2 = FOLC1, FRIC2 = FRIC1, LC2 = LC1, P2 = P1, V2 = V1;
20 constructMatrix2(rID, mID, v11, v12, v13, v14), constructMatrix2(rID, mID, v21, v22, v23, v24) -> v21 = v11, v22 = v12, v23 = v13, v24 = v14;
21 constructMatrix4(rID, mID, v11, v12, v13, v14), constructMatrix4(rID, mID, v11, v12, v13, v14, v15, v16, v17, v18, v19, v110, v111), constructMatrix4(rID, mID, v21, v22, v23, v24, v25, v26, v27, v28, v29, v210, v211) ->
v11 = v21, v12 = v22, v13 = v23, v14 = v24, v15 = v25, v16 = v26, v17 = v27, v18 = v28, v19 = v29, v110 = v210, v111 = v211;
22 constructMtxMatrix(i, j, v, id)(i, j1, v1, id), constructMtxMatrix(i2, j2, v2, id) -> i2 = i1, j2 = j1, v2 = v1;
23 indicator(id1, id2, idr31), indicator(id1, id2, idr32) -> idr32 = idr31;
24 deriveFrom(id1, id2, idr31), deriveFrom(id1, id2, idr32) -> idr32 = idr31;
25 concat(id1, id2, idr1), concat(id1, id2, idr2) -> idr2 = idr1;
26 subsetweet(id1, id2, id3, id4, id5, id6, id7, idr1), subsetweet(id1, id2, id3, id4, id5, id6, id7, idr2) -> idr2 = idr1;
27 subsetuser(id1, id2, id3, id4, id5, id6, idr1), subsetuser(id1, id2, id3, id4, id5, id6, idr2) -> idr2 = idr1;
28 tweet(id1), tweet(id2) -> id2 = id1;
29 childsp(x, y1, z, d), childsp(x, y2, z, d) -> y2 = y1;
30 valsp(x, v1), valsp(x, v2) -> v2 = v1;
31 childsjV2(x, y1, z, d), childsjV2(x, y2, z, d) -> y2 = y1;
32
33 #TGDs
34 add(idA, idB, idr1) -> add(idB, idA, idr1);
35 multi(id1, id2, idr1), rowSums(idr1, idr2) -> rowSums(id2, idr3), multi(id1, idr3, idr2);
36 add(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
37 multi(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idB, idC, idr3), multi(idA, idr3, idr2);
38 add(idA, idB, idr1), tr(idr1, idT) -> tr(idA, idAT), tr(idB, idBT), add(idAT, idBT, idT);
39 tr(idA, idB, idr1), tr(idB, idAT, idAT, idr2) -> multi(idA, idB, idr1), tr(idr1, idr2);
40 constructMtxMatrix(i, j, v, id), name(id, n), matrixFilterLT(id, x, idr1) ->
name(id1, n), constructMtxMatrix(i, j, v, idr1), LT(v, x);
41 tr(id, id1), colSums(idr1, idr2) -> rowSums(id, idr1), tr(idr3, idr2);
42
43 subsetweet(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
44 subsetuser(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
45 indicator(mID2, mID3, idK), deriveFrom(idS1, idR1, id), rowSums(id, idr1) ->
rowSums(idS1, idr1), rowSums(idR1, idr3), multi(idR3, idr4), add(idR2, idr4, idr1);
46
47 user(EID, FOLC, FRIC, LC, P, V),
48 tweetR(EID, FC, QC, RC, F, PS, R),
49 constructMatrix11(mID1, rwID, FOLC, FRIC, LC, P, V, FC, QC, RC, F, PS, R) ->
name(mID1, n),
50 name(mID2, "tweet"),
51 name(mID3, "user"),
52 subsetweet(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
53 subsetuser(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
54 indicator(mID2, mID3, idr3),
55 deriveFrom(idS1, idR1, mID1);
56 name(idS, "tweet"),
57 name(idR, "user"),
58 indicator(idS, idR, idK),
59 subsetweet(idS, "FC", "QC", "RC", "F", "PS", "R", idS1),
60 subsetuser(idR, "FOLC", "FRIC", "LC", "P", "V", idR1),
61 rowSums(idS1, idr1),
62 rowSums(idR1, idr2),
63 multi(idK, idr2, idr3),
64 add(idr1, idr3, idr4) -> name(idr4, "VV1");
65
66 tweet(id),
67 childsp(id, "tweetID", "o", tID),
68 childsp(id, "flevel", "o", fID),
69 childsp(id, "entities", "o", eID),
70 childsp(eID, "hashtag", "o", hID),
71 childsp(hID, "hID", "o", hcID),
72 childsp(id, "place", "o", pID),
73 childsp(pID, "code", "o", cID) -> V1(tID, hcID, fID, cID);
74
75 tweet(id),
76 childsp(id, "tweetID", "o", tID),
77 childsp(id, "text", "o", text) -> V2(idV2),
78 childsjV2(idV2, "tweetID", "o", tID),
79 childsjV2(idV2, "text", "o", text);
80
81
82
83
84
85
86
87

```

Listing G2: Forward Constraints Used for Q2

```

1 #EGDs
2 name(id1, n1), name(id2, n1) -> id2 = id1;
3 tr(id1, idt1), tr(id1, idt2) -> idt2 = idt1;
4 in(id1, idt1), in(id1, idt2) -> idt2 = idt1;
5 ins(id1, idt1), ins(id1, idt2) -> idt1 = idt2;
6 in(id1, idr1), in(idr1, idr3) -> idr3 = id1;
7 multi(idA, idB, idr1), multi(idA, idB, idr2) -> idr2 = idr1;
8 multiS(idA, idB, idr1), multiS(idA, idB, idr2) -> idr2 = idr1;
9 det(idA, d1), det(idA, d2) -> d2 = d1;
10 trace(idA, id), trace(idA, id2) -> id2 = id1;
11 add(idA, idB, idr1), add(idA, idB, idr2) -> idr1 = idr2;
12 size(idA, i1, j1), size(idA, i2, j2) -> j2 = j1, i2 = i1;
13 name(id1, n1), name(id2, n1) -> id2 = id1;
14 name(mID1, n1), name(mID1, n2) -> n1 = n2;
15 addS(idA, idB, idr1), addS(idA, idB, idr2) -> idr1 = idr2;
16 sum(id1, idt1), sum(id1, idt2) -> idt2 = idt1;
17 rowSums(idA, id1), rowSums(idA, id2) -> id2 = id1;
18 colSums(idA, id1), colSums(idA, id2) -> id2 = id1;
19 user(EID, FOLC1, FRIC1, LC1, P1, V1), user(EID, FOLC2, FRIC2, LC2, P2, V2) -> FOLC2 = FOLC1, FRIC2 = FRIC1, LC2 = LC1, P2 = P1, V2 = V1;
20 constructMatrix2(rID, mID, v11, v12, v13, v14), constructMatrix2(rID, mID, v21, v22, v23, v24) -> v21 = v11, v22 = v12, v23 = v13, v24 = v14;
21 constructMatrix4(rID, mID, v11, v12, v13, v14, v15, v16, v17, v18, v19, v110, v111), constructMatrix4(rID, mID, v21, v22, v23, v24, v25, v26, v27, v28, v29, v210, v211) ->
v11 = v21, v12 = v22, v13 = v23, v14 = v24, v15 = v25, v16 = v26, v17 = v27, v18 = v28, v19 = v29, v110 = v210, v111 = v211;
22 constructMtxMatrix11(i, j, v, id)(i, j1, v1, id), constructMtxMatrix11(i2, j2, v2, id) -> i2 = i1, j2 = j1, v2 = v1;
23 indicator(id1, id2, idr31), indicator(id1, id2, idr32) -> idr32 = idr31;
24 deriveFrom(id1, id2, idr31), deriveFrom(id1, id2, idr32) -> idr32 = idr31;
25 concat(id1, id2, idr1), concat(id1, id2, idr2) -> idr2 = idr1;
26 subsetweet(id1, id2, id3, id4, id5, id6, id7, idr1),
27 subsetweet(id1, id2, id3, id4, id5, id6, id7, idr2) -> idr2 = idr1;
28 subsetweet_r(id1, id2, id3, id4, id5, id6, idr1),
29 subsetweet_r(id1, id2, id3, id4, id5, id6, idr2) -> idr2 = idr1;
30 subsetuse_r(id1, id2, id3, id4, id5, id6, idr2) -> idr2 = idr1;
31 tweet(id1), tweet(id2) -> id2 = id1;
32 childsp(x, y1, z, d), childsp(x, y2, z, d) -> y2 = y1;
33 valsp(x, v1), valsp(x, v2) -> v2 = v1;
34 valspjV2(x, y1, z, d), valspjV2(x, y2, z, d) -> y2 = y1;
35
36 #TGDs
37 add(idA, idB, idr1) -> add(idB, idA, idr1);
38 multi(id1, id2, idr1), rowSums(idr1, idr2) -> rowSums(id2, idr3), multi(id1, idr3, idr2);
39 add(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
40 multi(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idB, idC, idr3), multi(idA, idr3, idr2);
41 constructMtxMatrix1(i, j, v, id), name(id, n), matrixFilterLT(id, x, idr1) ->
42 name(idr1, n), constructMtxMatrix1(i, j, v, idr1), LT(v, x);
43 add(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
44
45 subsetweet(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
46 subsetuser(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
47 indicator(mID2, mID3, idK), deriveFrom(idS1, idR1, id), colSums(id, idr1) ->
48 colSums(idK, idr2), multi(idr2, idR1, idr3), colSums(idS1, idr4), concat(idr4, idr3, idr1);
49
50
51 user(EID, FOLC, FRIC, LC, P, V),
52 tweetR(EID, FC, QC, RC, F, PS, R),
53 constructMatrix11(mID1, rwID, FOLC, FRIC, LC, P, V, FC, QC, RC, F, PS, R) ->
54 name(mID1, n),
55 name(mID2, "tweet"),
56 name(mID3, "user"),
57 subsetweet(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
58 subsetweet(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
59 subsetuser_r(idR, "FOLC", "FRIC", "LC", "P", "V", idR1),
60 indicator(mID2, mID3, idr3),
61 deriveFrom(idS1, idR1, mID1),
62 name(idS, "tweet"),
63 name(idR, "user"),
64 indicator(idS, idR, idK),
65 subsetweet(idS, "FC", "QC", "RC", "F", "PS", "R", idS1),
66 subsetweet(idR, "FOLC", "FRIC", "LC", "P", "V", idR1),
67 rowSums(idS1, idr1),
68 rowSums(idR1, idr2),
69 multi(idK, idr2, idr3),
70 add(idr1, idr3, idr4) -> name(idr4, "VV1");
71
72 tweet(id),
73 childsp(id, "tweetID", "o", tID),
74 childsp(id, "flevel", "o", fID),
75 childsp(id, "entities", "o", eID),
76 childsp(eID, "hashtag", "o", hID),
77 childsp(hID, "hID", "o", hcID),
78 childsp(id, "place", "o", pID),
79 childsp(pID, "code", "o", cID) -> V1(tID, hcID, fID, cID);
80
81 tweet(id),
82 childsp(id, "tweetID", "o", tID),
83 childsp(id, "text", "o", text) ->
84 V2(tdV2),
85 childspjV2(idV2, "tweetID", "o", tID),
86 childspjV2(idV2, "text", "o", text);

```

Listing G3: Forward Constraints Used for Q3

```

1 #EGDs
2 name(id1, n1), name(id2, n1) -> id2 = id1;
3 tr(id1, idt1), tr(id1, idt2) -> idt2 = idt1;
4 in(id1, idt1), in(id1, idt2) -> idt2 = idt1;
5 ins(id1, idt1), ins(id1, idt2) -> idt1 = idt2;
6 in(id1, idr1), in(idr1, idr3) -> idr3 = id1;
7 multi(idA, idB, idr1), multi(idA, idB, idr2) -> idr2 = idr1;
8 multiS(idA, idB, idr1), multiS(idA, idB, idr2) -> idr2 = idr1;
9 det(idA, d1), det(idA, d2) -> d2 = d1;
10 trace(idA, id), trace(idA, id2) -> id2 = id1;
11 add(idA, idB, idr1), add(idA, idB, idr2) -> idr1 = idr2;
12 size(idA, i1, j1), size(idA, i2, j2) -> j2 = j1, i2 = i1;
13 name(id1, n1), name(id2, n1) -> id2 = id1;
14 name(mID1, n1), name(mID1, n2) -> n1 = n2;
15 addS(idA, idB, idr1), addS(idA, idB, idr2) -> idr1 = idr2;
16 sum(id1, idt1), sum(id1, idt2) -> idt2 = idt1;
17 rowSums(idA, id1), rowSums(idA, id2) -> id2 = id1;
18 colSums(idA, id1), colSums(idA, id2) -> id2 = id1;
19 user(EID, FOLC1, FRIC1, LC1, P1, V1), user(EID, FOLC2, FRIC2, LC2, P2, V2) -> FOLC2 = FOLC1, FRIC2 = FRIC1, LC2 = LC1, P2 = P1, V2 = V1;
20 constructMatrix2(rID, mID, v11, v12, v13, v14), constructMatrix2(rID, mID, v21, v22, v23, v24) -> v21 = v11, v22 = v12, v23 = v13, v24 = v14;
21 constructMatrix4(rID, mID, v11, v12, v13, v14), constructMatrix4(rID, mID, v11, v12, v13, v14, v15, v16, v17, v18, v19, v110, v111), constructMatrix11(rID, mID, v21, v22, v23, v24) -> v21 = v11, v22 = v12, v23 = v13, v24 = v14;
22 constructMatrix11(rID, mID, v11, v12, v13, v14, v15, v16, v17, v18, v19 = v29, v110 = v210, v111 = v211);
23 v11 = v21, v12 = v22, v13 = v23, v14 = v24, v15 = v25, v16 = v26, v17 = v27, v18 = v28, v19 = v29, v110 = v210, v111 = v211;
24 constructMtxMatrix(i, j, v, id)(i, j1, v1, id), constructMtxMatrix(j2, j2, v2, id) -> j2 = j1, v2 = v1;
25 indicator(id1, id2, idr31), indicator(id1, id2, idr32) -> idr32 = idr31;
26 deriveFrom(id1, id2, idr31), deriveFrom(id1, id2, idr32) -> idr32 = idr31;
27 concat(id1, id2, idr1), concat(id1, id2, idr2) -> idr2 = idr1;
28 subsetweet(id1, id2, id3, id4, id5, id6, id7, idr1), subsetweet(id1, id2, id3, id4, id5, id6, id7, idr2) -> idr2 = idr1;
29 subsetweet(id1, id2, id3, id4, id5, id6, id7, idr1), subsetweet(id1, id2, id3, id4, id5, id6, idr1), subsetweet(id1, id2, id3, id4, id5, id6, idr2) -> idr2 = idr1;
30 subsetweet(id1, id2, id3, id4, id5, id6, idr1), subsetweet(id1, id2, id3, id4, id5, id6, idr2) -> idr2 = idr1;
31 subsetweet(id1, id2, id3, id4, id5, id6, idr1), subsetweet(id1, id2, id3, id4, id5, id6, idr2) -> idr2 = idr1;
32 tweet(id1), tweet(id2) -> id2 = id1;
33 childsp(x, y1, z, d), childsp(x, y2, z, d) -> y2 = y1;
34 valsp(x, v1), valsp(x, v2) -> v2 = v1;
35 childsjV2(x, y1, z, d), childsjV2(x, y2, z, d) -> y2 = y1;
36
37 #TGDs
38 add(idA, idB, idr1) -> add(idB, idA, idr1);
39 multi(id1, id2, idr1), rowSums(idr1, idr2) -> rowSums(id2, idr3), multi(id1, idr3, idr2);
40 add(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
41 multi(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idB, idC, idr3), multi(idA, idr3, idr2);
42 constructMtxMatrix(i, j, v, id), name(id, n), matrixFilterLT(id, x, idr1) ->
43 name(idr1, n), constructMtxMatrix(i, j, v, idr1), LT(v, x);
44 add(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
45
46
47 subsetweet(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
48 subsetweet(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
49 indicator(mID2, mID3, idK), deriveFrom(idS1, idR1, id), rowSums(id, idr1) ->
50 rowSums(idS1, idr2), rowSums(idR1, idr3), multi(idK, idr3, idr4), add(idr2, idr4, idr1);
51
52
53 user(EID, FOLC, FRIC, LC, P, V),
54 tweetR(EID, FC, QC, RC, F, PS, R),
55 constructMatrix11(mID1, rwdID, FOLC, FRIC, LC, P, V, FC, QC, RC, F, PS, R) ->
56 name(mID1, n),
57 name(mID2, "tweet"),
58 name(mID3, "user"),
59 subsetweet(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
60 subsetweet(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
61 indicator(mID2, mID3, idr3),
62 deriveFrom(idS1, idR1, mID1);
63
64 name(idS, "tweet"),
65 name(idR, "user"),
66 indicator(idS, idR, idK),
67 subsetweet(idS, "FC", "QC", "RC", "F", "PS", "R", idS1),
68 subsetweet(idR, "FOLC", "FRIC", "LC", "P", "V", idR1),
69 rowSums(idS1, idr1),
70 rowSums(idR1, idr2),
71 multi(idK, idr2, idr3),
72 add(idr1, idr3, idr4) -> name(idr4, "VV1");
73
74 tweet(id),
75 childsp(id, "tweetID", "o", tID),
76 childsp(id, "flevel", "o", fID),
77 childsp(id, "entities", "o", eID),
78 childsp(eID, "hashtag", "o", hID),
79 childsp(hID, "hID", "o", hcID),
80 childsp(id, "place", "o", pID),
81 childsp(pID, "code", "o", cID) -> V1(tID, hcID, fID, cID);
82
83 tweet(id),
84 childsp(id, "tweetID", "o", tID),
85 childsp(id, "text", "o", text) ->
86 V2(idV2),
87 childsjV2(idV2, "tweetID", "o", tID),
88 childsjV2(idV2, "text", "o", text);
89

```

Listing G4: Forward Constraints Used for Q4

```

1 #EGDs
2 name(id1, n1), name(id2, n1) -> id2 = id1;
3 tr(id1, idt1), tr(id1, idt2) -> idt2 = idt1;
4 in(id1, idt1), in(id1, idt2) -> idt2 = idt1;
5 ins(id1, idt1), ins(id1, idt2) -> idt1 = idt2;
6 in(id1, idr1), in(idr1, idr3) -> idr3 = id1;
7 multi(idA, idB, idr1), multi(idA, idB, idr2) -> idr2 = idr1;
8 multiS(idA, idB, idr1), multiS(idA, idB, idr2) -> idr2 = idr1;
9 det(idA, d1), det(idA, d2) -> d2 = d1;
10 trace(idA, id), trace(idA, id2) -> id2 = id1;
11 add(idA, idB, idr1), add(idA, idB, idr2) -> idr1 = idr2;
12 size(idA, i1, j1), size(idA, i2, j2) -> j2 = j1, i2 = i1;
13 name(id1, n1), name(id2, n1) -> id2 = id1;
14 name(mID1, n1), name(mID1, n2) -> n1 = n2;
15 addS(idA, idB, idr1), addS(idA, idB, idr2) -> idr1 = idr2;
16 sum(id1, idt1), sum(id1, idt2) -> idt2 = idt1;
17 rowSums(idA, id1), rowSums(idA, id2) -> id2 = id1;
18 colSums(idA, id1), colSums(idA, id2) -> id2 = id1;
19 user(EID, FOLC1, FRIC1, LC1, P1, V1), user(EID, FOLC2, FRIC2, LC2, P2, V2) -> FOLC2 = FOLC1, FRIC2 = FRIC1, LC2 = LC1, P2 = P1, V2 = V1;
20 constructMatrix2(rID, mID, v11, v12, v13, v14), constructMatrix2(rID, mID, v21, v22, v23, v24) -> v21 = v11, v22 = v12, v23 = v13, v24 = v14;
21 constructMatrix4(rID, mID, v11, v12, v13, v14), constructMatrix4(rID, mID, v16, v17, v18, v19, v110, v111), constructMatrix4(rID, mID, v21, v22, v23, v24, v25, v26, v27, v28, v29, v210, v211) ->
v11 = v21, v12 = v22, v13 = v23, v14 = v24, v15 = v25, v16 = v26, v17 = v27, v18 = v28, v19 = v29, v110 = v210, v111 = v211;
22 constructMtxMatrix(i, j, v, id)(i, j1, v1, id), constructMtxMatrix(i, j2, v2, id) -> i2 = i1, j2 = j1, v2 = v1;
23 indicator(id1, id2, idr31), indicator(id1, id2, idr32) -> idr32 = idr31;
24 deriveFrom(id1, id2, idr31), deriveFrom(id1, id2, idr32) -> idr32 = idr31;
25 concat(id1, id2, idr1), concat(id1, id2, idr2) -> idr2 = idr1;
26 subsetweet(id1, id2, id3, id4, id5, id6, id7, idr1),
27 subsetweet(id1, id2, id3, id4, id5, id6, id7, idr2) -> idr2 = idr1;
28 subsetweet_r(id1, id2, id3, id4, id5, id6, idr1),
29 subsetweet_r(id1, id2, id3, id4, id5, id6, idr2),
30 subsetuser(id1, id2, id3, id4, id5, id6, idr2) -> idr2 = idr1;
31 subsetuser_r(id1, id2, id3, id4, id5, id6, idr2) -> idr2 = idr1;
32 tweet(id1), tweet(id2) -> id2 = id1;
33 childsp(x, y1, z, d), childsp(x, y2, z, d) -> y2 = y1;
34 valsp(x, v1), valsp(x, v2) -> v2 = v1;
35 childsjV2(x, y1, z, d), childsjV2(x, y2, z, d) -> y2 = y1;
36
37 #TGDs
38 add(idA, idB, idr1) -> add(idB, idA, idr1);
39 multi(id1, id2, idr1), colSums(idr1, idr2) -> colSums(id1, idr3), multi(idr3, id2, idr2);
40 multi(id1, id2, idr1), rowSums(idr1, idr2) -> rowSums(id2, idr3), multi(id1, idr3, idr2);
41 add(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
42 multi(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idB, idC, idr3), multi(idA, idr3, idr2);
43 namecMtxMatrix(i, j, v, id), namec(d, n), matrixFilterLT(id, x, idr1) ->
namec(idr1, n), constructMtxMatrix(i, j, v, idr1), LT(v, x);
44 add(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
45
46 subsetweet(mID2, "FC", "QC", "F", "PS", "R", idS1),
47 subsetweet(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
48 indicator(mID2, mID3, idK), deriveFrom(idS1, idR1, id), colSums(id, idr1) ->
colSums(idK, idr2), multi(idr2, idR1, idr3), colSums(idS1, idr4), concat(idr4, idr3, idr1);
49
50
51
52 user(EID, FOLC, FRIC, LC, P, V),
53 tweetR(EID, FC, QC, RC, F, PS, R),
54 constructMatrix11(mID1, rwd1, FOLC, FRIC, LC, P, V, FC, QC, RC, F, PS, R) ->
name(mID1, n),
55 name(mID2, "tweet"),
56 name(mID3, "user"),
57 subsetweet(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
58 subsetweet(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
59 indicator(mID2, mID3, idr3),
60 deriveFrom(idS1, idR1, mID1);
61 name(idS, "tweet"),
62 name(idR, "user"),
63 indicator(idS, idR, idK),
64 subsetweet(idS, "FC", "QC", "RC", "F", "PS", "R", idS1),
65 subsetweet(idR, "FOLC", "FRIC", "LC", "P", "V", idR1),
66 rowSums(idS1, idr1),
67 rowSums(idR1, idr2),
68 multi(idK, idr2, idr3),
69 add(idr1, idr3, idr4) -> name(idr4, "VV1");
70
71
72 tweet(id),
73 childsp(id, "tweetID", "o", tID),
74 childsp(id, "flevel", "o", fID),
75 childsp(id, "entities", "o", eID),
76 childsp(eID, "hashtag", "o", hID),
77 childsp(hID, "hID", "o", hcID),
78 childsp(id, "place", "o", pID),
79 childsp(pID, "code", "o", cID) -> V1(tID, hcID, fID, cID);
80
81
82 tweet(id),
83 childsp(id, "tweetID", "o", tID),
84 childsp(id, "text", "o", text) ->
V2(idV2),
85 childsjV2(idV2, "tweetID", "o", tID),
86 childsjV2(idV2, "text", "o", text);
87

```

Listing G5: Forward Constraints Used for Q5

```

1 #EGDs
2 name(id1, n1), name(id2, n1)-> id2 = id1;
3 tr(id1, idt1), tr(id1, idt2)-> idt2 = idt1;
4 in(id1, idt1), in(id1, idt2)-> idt2 = idt1;
5 ins(id1, idt1), ins(id1, idt2)-> idt1 = idt2;
6 in(id1, idr1), in(idr1, idr3)-> idr3 = id1;
7 multi(idA, idB, idr1), multi(idA, idB, idr2)-> idr2 = idr1;
8 multi(idA, idB, idr1), multi(idA, idB, idr2)-> idr2 = idr1;
9 det(idA, d1), det(idA, d2)-> d2 = d1;
10 trace(idA, id1), trace(idA, id2)-> id2 = id1;
11 add(idA, idB, idr1), add(idA, idB, idr2)-> idr1 = idr2;
12 size(idA, i1, j1), size(idA, i2, j2)-> j2 = j1, i2 = i1;
13 name(id1, n1), name(id2, n1)-> id2 = id1;
14 name(mID1, n1), name(mID1, n2)-> n1 = n2;
15 addS(idA, idB, idr1), addS(idA, idB, idr2)-> idr1 = idr2;
16 sum(id1, idt1), sum(id1, idt2)-> idt2 = idt1;
17 rowSums(idA, id1), rowSums(idA, id2)-> id2 = id1;
18 colSumS(idA, id1), colSumS(idA, id2)-> id2 = id1;
19 user(EID, FOLC1, FRIC1, LC1, P1, V1), user(EID, FOLC2, FRIC2, LC2, P2, V2)-> FOLC2 = FOLC1, FRIC2 = FRIC1, LC2 = LC1, P2 = P1, V2 = V1;
20 constructMatrix2(rID, mID, v11, v12, v13, v14), constructMatrix2(rID, mID, v21, v22)-> v21 = v11, v22 = v12;
21 constructMatrix4(rID, mID, v11, v12, v13, v14), constructMatrix4(rID, mID, v21, v22, v23, v24)-> v21 = v11, v22 = v12, v23 = v13, v24 = v14;
22 constructMatrix11(rID, mID, v11, v12, v13, v14, v15, v16, v17, v18, v19, v110, v111), constructMatrix11(rID, mID, v21, v22, v23, v24)-> v21 = v11, v22 = v12, v23 = v13, v24 = v14;
23 v11 = v21, v12 = v22, v13 = v23, v14 = v24, v15 = v25, v16 = v26, v17 = v27, v18 = v28, v19 = v29, v110 = v210, v111 = v211;
24 constructMtxMatrix(i, j, v, id)(i1, j1, v1, id), constructMtxMatrix(i2, j2, v2, id)-> i2 = i1, j2 = j1, v2 = v1;
25 indicator(id1, id2, id3), indicator(id1, id2, id3)-> idr32 = idr31;
26 deriveFrom(id1, id2, id3), deriveFrom(id1, id2, id3)-> idr32 = idr31;
27 concat(id1, id2, idr1), concat(id1, id2, idr2)-> idr2 = idr1;
28 subset_tweet(id1, id2, id3, id4, id5, id6, id7, idr1);
29 subset_tweet(id1, id2, id3, id4, id5, id6, id7, idr2)-> idr2 = idr1;
30 subset_user(id1, id2, id3, id4, id5, id6, idr1);
31 subset_user(id1, id2, id3, id4, id5, id6, idr2)-> idr2 = idr1;
32 tweet(id1), tweet(id2)-> id2 = id1;
33 childsp(x, y1, z, d), childsp(x, y2, z, d)-> y2 = y1;
34 valsP(x, v1), valsP(x, v2)-> v2 = v1;
35 childsjV2(x, y1, z, d), childsjV2(x, y2, z, d)-> y2 = y1;
36
37 #TGDs
38 add(idA, idB, idr1)-> add(idB, idA, idr1);
39 multi(id1, id2, idr1), colSumS(id1, idr3), multi(idr3, id2, idr2);
40 multi(id1, id2, idr1), rowSums(id1, idr2)-> rowSums(id2, idr3), multi(id1, idr3, idr2);
41 add(idA, idB, idr1), multi(idC, idr2)-> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
42 multi(idA, idB, idr1), multi(idr1, idC, idr2)-> multi(idB, idr3, idr2), multi(idA, idr3, idr2);
43 constructMatrixXMatrix(i, j, v, id), name(id, n), matrixFilterLT(id, x, idr1)->
44 name(idr1, n), constructMatrixXMatrix(i, j, v, idr1), LT(v, x);
45 add(idA, idB, idr1), multi(idr1, idC, idr2)-> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
46
47 subset_tweet(mID2, "FC", "QC", "RC", "PS", "R", idS1),
48 subset_user(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
49 indicator(mID2, mID3, idK), deriveFrom(idS1, idR1, id), colSumS(id, idr1)->
50 colSumS(idK, idr2), multi(idr2, idR1, idr3), colSumS(idS1, idr4), concat(idr4, idr3, idr1);
51
52
53 user(EID, FOLC, FRIC, LC, P, V),
54 tweetR(EID, FC, QC, RC, F, PS, R),
55 constructMatrix11(mID1, rwID, FOLC, FRIC, LC, P, V, FC, QC, RC, F, PS, R)->
56 name(mID1, n),
57 name(mID2, "tweet"),
58 name(mID3, "user"),
59 subset_tweet(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
60 subset_user(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
61 indicator(mID2, mID3, idR3),
62 deriveFrom(idS1, idR1, mID1),
63 name(idS, "tweet"),
64 name(idR, "user"),
65 indicator(idS, idR, idK),
66 subset_tweet(idS, "FC", "QC", "RC", "F", "PS", "R", idS1),
67 subset_user(idR, "FOLC", "FRIC", "LC", "P", "V", idR1),
68 rowSums(idS1, idr1),
69 rowSums(idR1, idr2),
70 multi(idK, idr2, idr3),
71 add(idr1, idr3, idr4)-> name(idr4, "VV1");
72
73 tweet(id),
74 childsp(id, "tweetID", "o", tID),
75 childsp(id, "flevel", "o", fID),
76 childsp(id, "entities", "o", eID),
77 childsp(eID, "hashtag", "o", hID),
78 childsp(hID, "hID", "o", hcID),
79 childsp(id, "place", "o", pID),
80 childsp(pID, "code", "o", cID)-> V1(tID, hcID, fID, cID);
81
82 tweet(id),
83 childsp(id, "tweetID", "o", tID),
84 childsp(id, "text", "o", text)->
85 V2(dV2),
86 childsjV2(idV2, "tweetID", "o", tID),
87 childsjV2(idV2, "text", "o", text);

```

Listing G6: Forward Constraints Used for Q6

```

1 #EGDs
2 name(id1, n1), name(id2, n1) -> id2 = id1;
3 tr(id1, idt1), tr(id1, idt2) -> idt2 = idt1;
4 in(id1, idt1), in(id1, idt2) -> idt2 = idt1;
5 ins(id1, idt1), ins(id1, idt2) -> idt1 = idt2;
6 in(id1, idr1), in(idr1, idr3) -> idr3 = id1;
7 multi(idA, idB, idr1), multi(idA, idB, idr2) -> idr2 = idr1;
8 multiS(idA, idB, idr1), multiS(idA, idB, idr2) -> idr2 = idr1;
9 det(idA, d1), det(idA, d2) -> d2 = d1;
10 trace(idA, id), trace(idA, id2) -> id2 = id1;
11 add(idA, idB, idr1), add(idA, idB, idr2) -> idr1 = idr2;
12 size(idA, i1, j1), size(idA, i2, j2) -> j2 = j1, i2 = i1;
13 name(id1, n1), name(id2, n1) -> id2 = id1;
14 name(mID1, n1), name(mID1, n2) -> n1 = n2;
15 addS(idA, idB, idr1), addS(idA, idB, idr2) -> idr1 = idr2;
16 sum(id1, idt1), sum(id1, idt2) -> idt2 = idt1;
17 rowSums(idA, id1), rowSums(idA, id2) -> id2 = id1;
18 colSums(idA, id1), colSums(idA, id2) -> id2 = id1;
19 user(EID, FOLC1, FRIC1, LC1, P1, V1), user(EID, FOLC2, FRIC2, LC2, P2, V2) -> FOLC2 = FOLC1, FRIC2 = FRIC1, LC2 = LC1, P2 = P1, V2 = V1;
20 constructMatrix2(rID, mID, v11, v12, v13, v14), constructMatrix2(rID, mID, v21, v22, v23, v24) -> v21 = v11, v22 = v12, v23 = v13, v24 = v14;
21 constructMatrix4(rID, mID, v11, v12, v13, v14), constructMatrix4(rID, mID, v11, v12, v13, v14) -> v21 = v11, v22 = v12, v23 = v13, v24 = v14;
22 constructMatrix11(rID, mID, v11, v12, v13, v14, v15, v16, v17, v18, v19, v110, v111), constructMatrix11(rID, mID, v21, v22, v23, v24, v25, v26, v27, v28, v29, v210, v211) -> v11 = v21, v12 = v22, v13 = v23, v14 = v24, v15 = v25, v16 = v26, v17 = v27, v18 = v28, v19 = v29, v110 = v210, v111 = v211;
23 constructMtxMatrix(i, j, v, id)(i, j1, v1, id), constructMtxMatrix(j2, j2, v2, id) -> i2 = i1, j2 = j1, v2 = v1;
24 indicator(id1, id2, idr3), indicator(id1, id2, idr3) -> idr32 = idr31;
25 deriveFrom(id1, id2, idr3), deriveFrom(id1, id2, idr3) -> idr32 = idr31;
26 concat(id1, id2, idr1), concat(id1, id2, idr2) -> idr2 = idr1;
27 subsettweeT(id1, id2, id3, id4, id5, id6, id7, idr1), subsettweeT(id1, id2, id3, id4, id5, id6, id7, idr2) -> idr2 = idr1;
28 subsetuseT(id1, id2, id3, id4, id5, id6, idr1), subsetuseT(id1, id2, id3, id4, id5, id6, idr2) -> idr2 = idr1;
29 tweet(id1), tweet(id2) -> id2 = id1;
30 childsp(x, y1, z, d), childsp(x, y2, z, d) -> y2 = y1;
31 valsp(x, v1), valsp(x, v2) -> v2 = v1;
32 childsjV2(x, y1, z, d), childsjV2(x, y2, z, d) -> y2 = y1;
33
34
35
36
37 #TGDs
38 add(idA, idB, idr1) -> add(idB, idA, idr1);
39 multi(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idB, idC, idr3), multi(idA, idr3, idr2);
40 multi(id1, id2, idr1), colSums(id1, idr2) -> colSums(id1, idr2), multi(idr3, id2, idr2);
41 multi(id2, idr1), rowSums(id1, idr2) -> rowSums(id2, idr3), multi(id1, idr3, idr2);
42 add(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
43 construcMtxMatrix(i, j, v, id), name(id, n), matrixFilterLT(id, x, idr1) ->
44 name(idr1, n), construcMtxMatrix(i, j, v, idr1), LT(v, x);
45 add(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
46
47 subsettweeT(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
48 subsetuseT(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
49 indicator(mID2, mID3, idK), deriveFrom(idS1, idR1, id), colSums(id, idr1) ->
50 colSums(idK, idr2), multi(idr2, idR1, idr3), colSums(idS1, idr4), concat(idr4, idr3, idr1);
51
52
53 user(EID, FOLC, FRIC, LC, P, V),
54 tweetR(EID, FC, QC, RC, F, PS, R),
55 constructMatrix11(mID1, rwd1, FOLC, FRIC, LC, P, V, FC, QC, RC, F, PS, R) ->
56 name(mID1, n),
57 name(mID2, "tweet"),
58 name(mID3, "user"),
59 subsettweeT(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
60 subsetuseT(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
61 indicator(mID2, mID3, idr3),
62 deriveFrom(idS1, idR1, mID1);
63
64 name(idS, "tweet"),
65 name(idR, "user"),
66 indicator(idS, idR, idK),
67 subsettweeT(idS, "FC", "QC", "RC", "F", "PS", "R", idS1),
68 subsetuseT(idR, "FOLC", "FRIC", "LC", "P", "V", idR1),
69 rowSums(idS1, idr1),
70 rowSums(idR1, idr2),
71 multi(idK, idr2, idr3),
72 add(idr1, idr3, idr4) -> name(idr4, "VV1");
73
74 tweet(id),
75 childsp(id, "tweetID", "o", tID),
76 childsp(id, "flevel", "o", fID),
77 childsp(id, "entities", "o", eID),
78 childsp(eID, "hashtag", "o", hID),
79 childsp(hID, "hID", "o", hcID),
80 childsp(id, "place", "o", pID),
81 childsp(pID, "code", "o", cID) -> V1(tID, hcID, fID, cID);
82
83 tweet(id),
84 childsp(id, "tweetID", "o", tID),
85 childsp(id, "text", "o", text) ->
86 V2(idV2),
87 childsjV2(idV2, "tweetID", "o", tID),
88 childsjV2(idV2, "text", "o", text);
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
559
560
561
562
563
564
565
566
567
568
569
569
570
571
572
573
574
575
576
577
578
579
579
580
581
582
583
584
585
586
587
587
588
589
589
590
591
592
593
594
595
596
597
597
598
599
599
600
601
602
603
604
605
606
606
607
608
609
609
610
611
612
613
613
614
615
615
616
617
617
618
618
619
619
620
620
621
621
622
622
623
623
624
624
625
625
626
626
627
627
628
628
629
629
630
630
631
631
632
632
633
633
634
634
635
635
636
636
637
637
638
638
639
639
640
640
641
641
642
642
643
643
644
644
645
645
646
646
647
647
648
648
649
649
650
650
651
651
652
652
653
653
654
654
655
655
656
656
657
657
658
658
659
659
660
660
661
661
662
662
663
663
664
664
665
665
666
666
667
667
668
668
669
669
670
670
671
671
672
672
673
673
674
674
675
675
676
676
677
677
678
678
679
679
680
680
681
681
682
682
683
683
684
684
685
685
686
686
687
687
688
688
689
689
690
690
691
691
692
692
693
693
694
694
695
695
696
696
697
697
698
698
699
699
700
700
701
701
702
702
703
703
704
704
705
705
706
706
707
707
708
708
709
709
710
710
711
711
712
712
713
713
714
714
715
715
716
716
717
717
718
718
719
719
720
720
721
721
722
722
723
723
724
724
725
725
726
726
727
727
728
728
729
729
730
730
731
731
732
732
733
733
734
734
735
735
736
736
737
737
738
738
739
739
740
740
741
741
742
742
743
743
744
744
745
745
746
746
747
747
748
748
749
749
750
750
751
751
752
752
753
753
754
754
755
755
756
756
757
757
758
758
759
759
760
760
761
761
762
762
763
763
764
764
765
765
766
766
767
767
768
768
769
769
770
770
771
771
772
772
773
773
774
774
775
775
776
776
777
777
778
778
779
779
780
780
781
781
782
782
783
783
784
784
785
785
786
786
787
787
788
788
789
789
790
790
791
791
792
792
793
793
794
794
795
795
796
796
797
797
798
798
799
799
800
800
801
801
802
802
803
803
804
804
805
805
806
806
807
807
808
808
809
809
810
810
811
811
812
812
813
813
814
814
815
815
816
816
817
817
818
818
819
819
820
820
821
821
822
822
823
823
824
824
825
825
826
826
827
827
828
828
829
829
830
830
831
831
832
832
833
833
834
834
835
835
836
836
837
837
838
838
839
839
840
840
841
841
842
842
843
843
844
844
845
845
846
846
847
847
848
848
849
849
850
850
851
851
852
852
853
853
854
854
855
855
856
856
857
857
858
858
859
859
860
860
861
861
862
862
863
863
864
864
865
865
866
866
867
867
868
868
869
869
870
870
871
871
872
872
873
873
874
874
875
875
876
876
877
877
878
878
879
879
880
880
881
881
882
882
883
883
884
884
885
885
886
886
887
887
888
888
889
889
890
890
891
891
892
892
893
893
894
894
895
895
896
896
897
897
898
898
899
899
900
900
901
901
902
902
903
903
904
904
905
905
906
906
907
907
908
908
909
909
910
910
911
911
912
912
913
913
914
914
915
915
916
916
917
917
918
918
919
919
920
920
921
921
922
922
923
923
924
924
925
925
926
926
927
927
928
928
929
929
930
930
931
931
932
932
933
933
934
934
935
935
936
936
937
937
938
938
939
939
940
940
941
941
942
942
943
943
944
944
945
945
946
946
947
947
948
948
949
949
950
950
951
951
952
952
953
953
954
954
955
955
956
956
957
957
958
958
959
959
960
960
961
961
962
962
963
963
964
964
965
965
966
966
967
967
968
968
969
969
970
970
971
971
972
972
973
973
974
974
975
975
976
976
977
977
978
978
979
979
980
980
981
981
982
982
983
983
984
984
985
985
986
986
987
987
988
988
989
989
990
990
991
991
992
992
993
993
994
994
995
995
996
996
997
997
998
998
999
999
1000
1000
1001
1001
1002
1002
1003
1003
1004
1004
1005
1005
1006
1006
1007
1007
1008
1008
1009
1009
1010
1010
1011
1011
1012
1012
1013
1013
1014
1014
1015
1015
1016
1016
1017
1017
1018
1018
1019
1019
1020
1020
1021
1021
1022
1022
1023
1023
1024
1024
1025
1025
1026
1026
1027
1027
1028
1028
1029
1029
1030
1030
1031
1031
1032
1032
1033
1033
1034
1034
1035
1035
1036
1036
1037
1037
1038
1038
1039
1039
1040
1040
1041
1041
1042
1042
1043
1043
1044
1044
1045
1045
1046
1046
1047
1047
1048
1048
1049
1049
1050
1050
1051
1051
1052
1052
1053
1053
1054
1054
1055
1055
1056
1056
1057
1057
1058
1058
1059
1059
1060
1060
1061
1061
1062
1062
1063
1063
1064
1064
1065
1065
1066
1066
1067
1067
1068
1068
1069
1069
1070
1070
1071
1071
1072
1072
1073
1073
1074
1074
1075
1075
1076
1076
1077
1077
1078
1078
1079
1079
1080
1080
1081
1081
1082
1082
1083
1083
1084
1084
1085
1085
1086
1086
1087
1087
1088
1088
1089
1089
1090
1090
1091
1091
1092
1092
1093
1093
1094
1094
1095
1095
1096
1096
1097
1097
1098
1098
1099
1099
1100
1100
1101
1101
1102
1102
1103
1103
1104
1104
1105
1105
1106
1106
1107
1107
1108
1108
1109
1109
1110
1110
1111
1111
1112
1112
1113
1113
1114
1114
1115
1115
1116
1116
1117
1117
1118
1118
1119
1119
1120
1120
1121
1121
1122
1122
1123
1123
1124
1124
1125
1125
1126
1126
1127
1127
1128
1128
1129
1129
1130
1130
1131
1131
1132
1132
1133
1133
1134
1134
1135
1135
1136
1136
1137
1137
1138
1138
1139
1139
1140
1140
1141
1141
1142
1142
1143
1143
1144
1144
1145
1145
1146
1146
1147
1147
1148
1148
1149
1149
1150
1150
1151
1151
1152
1152
1153
1153
1154
1154
1155
1155
1156
1156
1157
1157
1158
1158
1159
1159
1160
1160
1161
1161
1162
1162
1163
1163
1164
1164
1165
1165
1166
1166
1167
1167
1168
1168
1169
1169
1170
1170
1171
1171
1172
1172
1173
1173
1174
1174
1175
1175
1176
1176
1177
1177
1178
1178
1179
1179
1180
1180
1181
1181
1182
1182
1183
1183
1184
1184
1185
1185
1186
1186
1187
1187
1188
1188
1189
1189
1190
1190
1191
1191
1192
1192
1193
1193
1194
1194
1195
1195
1196
1196
1197
1197
1198
1198
1199
1199
1200
1200
1201
1201
1202
1202
1203
1203
1204
1204
1205
1205
1206
1206
1207
1207
1208
1208
1209
1209
1210
1210
1211
1211
1212
1212
1213
1213
1214
1214
1215
1215
1216
1216
1217
1217
1218
1218
1219
1219
1220
1220
1221
1221
1222
1222
1223
1223
1224
1224
1225
1225
1226
1226
1227
1227
1228
1228
1229
1229
1230
1230
1231
1231
1232
1232
1233
1233
1234
1234
1235
1235
1236
1236
1237
1237
1238
1238
1239
1239
1240
1240
1241
1241
1242
1242
1243
1243
1244
1244
1245
1245
1246
1246
1247
1247
1248
1248
1249
1249
1250
1250
1251
1251
1252
1252
1253
1253
1254
1254
1255
1255
1256
1256
1257
1257
1258
1258
1259
1259
1260
1260
1261
1261
1262
1262
1263
1263
1264
1264
1265
1265
1266
1266
1267
1267
1268
1268
1269
1269
1270
1270
1271
1271
1272
1272
1273
1273
1274
1274
1275
1275
1276
1276
1277
1277
1278
1278
1279
1279
1280
1280
1281
1281
1282
1282
1283
1283
1284
1284
1285
1285
1286
1286
1287
1287
1288
1288
1289
1289
1290
1290
1291
1291
1292
1292
1293
1293
1294
1294
1295
1295
1296
1296
1297
1297
1298
1298
1299
1299
1300
1300
1301
1301
1302
1302
1303
1303
1304
1304
1305
1305
1306
1306
1307
1307
1308
1308
1309
1309
1310
1310
1311
1311
1312
1312
1313
1313
1314
1314
1315
1315
1316
1316
1317
1317
1318
1318
1319
1319
1320
1320
1321
1321
1322
1322
1323
1323
1324
1324
1325
1325
1326
1326
1327
1327
1328
1328
1329
1329
1330
1330
1331
1331
1332
1332
1333
1333
1334
1334
1335
1335
1336
1336
1337
1337
1338
1338
1339
1339
1340
1340
1341
1341
1342
1342
1343
1343
1344
1344
1345
1345
1346
1346
1347
1347
1348
1348
1349
1349
1350
1350
1351
1351

```

```

1 #EGDs
2 name(id1, n1), name(id2, n1)-> id2 = id1;
3 tr(id1, idt1), tr(id1, idt2)-> idt2 = idt1;
4 in(id1, idt1), in(id1, idt2)-> idt2 = idt1;
5 ins(id1, idt1), ins(id1, idt2)-> idt1 = idt2;
6 in(id1, idr1), in(idr1, idr3)-> idr3 = id1;
7 multi(idA, idB, idr1), multi(idA, idB, idr2)-> idr2 = idr1;
8 multiS(idA, idB, idr1), multiS(idA, idB, idr2)-> idr2 = idr1;
9 det(idA, d1), det(idA, d2)-> d2 = d1;
10 trace(idA, id), trace(idA, id2)-> id2 = id1;
11 add(idA, idB, idr1), add(idA, idB, idr2)-> idr1 = idr2;
12 size(idA, i1, j1), size(idA, i2, j2)-> j2 = j1, i2 = i1;
13 name(id1, n1), name(id2, n1)-> id2 = id1;
14 name(mID1, n1), name(mID1, n2)-> n1 = n2;
15 addS(idA, idB, idr1), addS(idA, idB, idr2)-> idr1 = idr2;
16 sum(id1, idt1), sum(id1, idt2)-> idt2 = idt1;
17 rowSums(idA, id1), rowSums(idA, id2)-> id2 = id1;
18 colSums(idA, id1), colSums(idA, id2)-> id2 = id1;
19 user(EID, FOLC1, FRIC1, LC1, P1, V1), user(EID, FOLC2, FRIC2, LC2, P2, V2)-> FOLC2 = FOLC1, FRIC2 = FRIC1, LC2 = LC1, P2 = P1, V2 = V1;
20 constructMatrix2(rID, mID, v11, v12, v13, v14), constructMatrix2(rID, mID, v21, v22, v23, v24)-> v21 = v11, v22 = v12, v23 = v13, v24 = v14;
21 constructMatrix4(rID, mID, v11, v12, v13, v14), constructMatrix4(rID, mID, v21, v22, v23, v24)-> v21 = v11, v22 = v12, v23 = v13, v24 = v14;
22 constructMatrix11(rID, mID, v11, v12, v13, v14, v15, v16, v17, v18, v19, v110, v111), constructMatrix11(rID, mID, v21, v22, v23, v24, v25, v26, v27, v28, v29, v210, v211)->
v11 = v21, v12 = v22, v13 = v23, v14 = v24, v15 = v25, v16 = v26, v17 = v27, v18 = v28, v19 = v29, v110 = v210, v111 = v211;
23 constructMtxMatrix(i, j, v, id)(i, j1, v1, id), constructMtxMatrix(j2, j2, v2, id)-> i2 = i1, j2 = j1, v2 = v1;
24 indicator(id1, id2, idr3), indicator(id1, id2, idr3)-> idr32 = idr31;
25 deriveFrom(id1, id2, idr3), deriveFrom(id1, id2, idr32)-> idr32 = idr31;
26 concat(id1, id2, idr1), concat(id1, id2, idr2)-> idr2 = idr1;
27 subsetUser(id1, id2, id3, id4, id5, id6, id7, idr1), subsetUser(id1, id2, id3, id4, id5, id6, id7, idr2)-> idr2 = idr1;
28 subsetUser(id1, id2, id3, id4, id5, id6, idr1), subsetUser(id1, id2, id3, id4, id5, id6, idr2)-> idr2 = idr1;
29 subsetUser(id1, id2, id3, id4, id5, id6, idr1), subsetUser(id1, id2, id3, id4, id5, id6, idr2)-> idr2 = idr1;
30 tweet(id1), tweet(id2)-> id2 = id1;
31 childsp(x, y1, z, d), childsp(x, y2, z, d)-> y2 = y1;
32 valSp(x, v1), valSp(x, v2)-> v2 = v1;
33 childsjV2(x, y1, z, d), childsjV2(x, y2, z, d)-> y2 = y1;
34
35
36
37 #TGDs
38 add(idA, idB, idr1)-> add(idB, idA, idr1);
39 multi(id1, id2, idr1), colSums(idr1, idr2)-> colSums(id1, idr3), multi(idr3, id2, idr2);
40 multi(idA, idB, idr1), multi(idr1, idC, idr2)-> multi(idB, idC, idr3), multi(idA, idr3, idr2);
41 add(idA, idB, idr1), multi(idr1, idC, idr2)-> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
42 add(idA, idB, idr1), trace(idr1, t)-> trace(idA, idr2), trace(idB, idr3), addS(idr2, idr3, t);
43
44 constructMtxMatrix(i, j, v, id), name(id, n), matrixFilterLT(id, x, idr1)->
name(id1, n), constructMtxMatrix(i, j, v, idr1), LT(v, x);
45 add(idA, idB, idr1), multi(idr1, idC, idr2)-> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
46
47
48 subsetUser(mID2, "FC", "QC", "RC", "F", "PS", "R"), subsetUser(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
49 indicator(mID2, mID3, idK), deriveFrom(idS1, idR1, id), colSums(id, idR1)->
colSums(idK, idR2), multi(idr2, idR1, idR3), colSums(idS1, idR4), concat(idr4, idr3, idR1);
50
51
52
53
54 user(EID, FOLC, FRIC, LC, P, V),
55 tweetR(EID, FC, QC, RC, F, PS, R),
56 constructMatrix11(mID1, rwID, FOLC, FRIC, LC, P, V, FC, QC, RC, F, PS, R)->
name(mID1, n),
57 name(mID2, "tweet"),
58 name(mID3, "user"),
59 subsetUser(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
60 subsetUser(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
61 indicator(mID2, mID3, idR3),
62 deriveFrom(idS1, idR1, mID1),
63
64
65 name(idS, "tweet"),
66 name(idR, "user"),
67 indicator(idS, idR, idK),
68 subsetUser(idS, "FC", "QC", "RC", "F", "PS", "R", idS1),
69 subsetUser(idR, "FOLC", "FRIC", "LC", "P", "V", idR1),
70 rowSums(idS1, idR1),
71 rowSums(idR1, idR2),
72 multi(idK, idR2, idR3),
73 add(idR1, idR3, idR4)-> name(idR4, "VV1");
74
75
76 tweet(id),
77 childsp(id, "tweetID", "o", tID),
78 childsp(id, "flevel", "o", fID),
79 childsp(id, "entities", "o", eID),
80 childsp(cID, "hashtag", "o", hID),
81 childsp(hID, "hID", "o", hcID),
82 childsp(id, "place", "o", pID),
83 childsp(pID, "code", "o", cID)-> V1(tID, hcID, fID, cID);
84
85
86 childsp(id, "tweetID", "o", tID),
87 childsp(id, "text", "o", text)->
V2(idV2),
88 childsjV2(idV2, "tweetID", "o", tID),
89 childsjV2(idV2, "text", "o", text);

```

Listing G8: Forward Constraints Used for Q8

```

1 #EGDs
2 name(id1, n1), name(id2, n1) -> id2 = id1;
3 tr(id1, idt1), tr(id1, idt2) -> idt2 = idt1;
4 in(id1, idt1), in(id1, idt2) -> idt2 = idt1;
5 ins(id1, idt1), ins(id1, idt2) -> idt1 = idt2;
6 in(id1, idr1), in(idr1, idr3) -> idr3 = id1;
7 multi(idA, idB, idr1), multi(idA, idB, idr2) -> idr2 = idr1;
8 multiS(idA, idB, idr1), multiS(idA, idB, idr2) -> idr2 = idr1;
9 det(idA, d1), det(idA, d2) -> d2 = d1;
10 trace(idA, id), trace(idA, id2) -> id2 = id1;
11 add(idA, idB, idr1), add(idA, idB, idr2) -> idr1 = idr2;
12 size(idA, i1, j1), size(idA, i2, j2) -> j2 = j1, i2 = i1;
13 name(id1, n1), name(id2, n1) -> id2 = id1;
14 name(mID1, n1), name(mID1, n2) -> n1 = n2;
15 addS(idA, idB, idr1), addS(idA, idB, idr2) -> idr1 = idr2;
16 sum(id1, idt1), sum(id1, idt2) -> idt2 = idt1;
17 rowSums(idA, id1), rowSums(idA, id2) -> id2 = id1;
18 colSums(idA, id1), colSums(idA, id2) -> id2 = id1;
19 user(EID, FOLC1, FRIC1, LC1, P1, V1), user(EID, FOLC2, FRIC2, LC2, P2, V2) -> FOLC2 = FOLC1, FRIC2 = FRIC1, LC2 = LC1, P2 = P1, V2 = V1;
20 constructMatrix2(rID, mID, v11, v12, v13, v14), constructMatrix2(rID, mID, v21, v22, v23, v24) -> v21 = v11, v22 = v12, v23 = v13, v24 = v14;
21 constructMatrix4(rID, mID, v11, v12, v13, v14), constructMatrix4(rID, mID, v16, v17, v18, v19, v110, v111), constructMatrix4(rID, mID, v21, v22, v23, v24, v25, v26, v27, v28, v29, v210, v211) ->
v11 = v21, v12 = v22, v13 = v23, v14 = v24, v15 = v25, v16 = v26, v17 = v27, v18 = v28, v19 = v29, v110 = v210, v111 = v211;
22 constructMatrix11(rID, mID, v11, v12, v13, v14, v15, v16, v17, v18, v19, v110, v111), constructMatrix11(rID, mID, v21, v22, v23, v24, v25, v26, v27, v28, v29, v210, v211) ->
v11 = v21, v12 = v22, v13 = v23, v14 = v24, v15 = v25, v16 = v26, v17 = v27, v18 = v28, v19 = v29, v110 = v210, v111 = v211;
23 constructMtxMatrix(i, j, v, id), constructMtxMatrix(j, f2, v2, id) -> i2 = i1, j2 = j1, v2 = v1;
24 indicator(id1, id2, idr31), indicator(id1, id2, idr32) -> idr32 = idr31;
25 deriveFrom(id1, id2, idr31), deriveFrom(id1, id2, idr32) -> idr32 = idr31;
26 concat(id1, id2, idr1), concat(id1, id2, idr2) -> idr2 = idr1;
27 subsetweet(id1, id2, id3, id4, id5, id6, id7, idr1),
28 subsetweet(id1, id2, id3, id4, id5, id6, id7, idr2) -> idr2 = idr1;
29 subsetweet(id1, id2, id3, id4, id5, id6, idr1),
30 subsetweet(id1, id2, id3, id4, id5, id6, idr2) -> idr2 = idr1;
31 subsetweet(id1, id2, id3, id4, id5, id6, idr2) -> idr2 = idr1;
32 tweet(id1), tweet(id2) -> id2 = id1;
33 childsp(x, y1, z, d), childsp(x, y2, z, d) -> y2 = y1;
34 valsp(x, v1), valsp(x, v2) -> v2 = v1;
35 childsjV2(x, y1, z, d), childsjV2(x, y2, z, d) -> y2 = y1;
36
37 #TGDs
38 colSums(id1, idr1), tr(idr1, idr2), rowSums(id2, idr3), multi(idr2, idr3, idr4), sum(idr4, s) ->
39 multi(id1, id2, idr5), sum(idr5, s);
40 multi(id1, id2, idr1), colSums(idr1, idr2) -> colSums(id1, idr3), multi(idr3, id2, idr2);
41 multi(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idB, idC, idr3), multi(idA, idr3, idr2);
42 constructMtxMatrix(i, j, v, id), name(id, n), matrixFilterLT(id, x, idr1) ->
43 name(idr1, n), constructMtxMatrix(i, j, v, idr1), LT(v, x);
44 add(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
45
46 subsetweet(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
47 subsetweet(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
48 indicator(mID2, mID3, idK) deriveFrom(idS1, idR1, id),
49 multi(idX, id, idr1) ->
50 multi(idX, idK, idr2),
51 multi(idr2, idR1, idr3),
52 multi(idX, idS1, idr4),
53 concat(idr4, idr3, idr1);
54
55 user(EID, FOLC, FRIC, LC, P, V),
56 tweetREID, FC, QC, RC, F, PS, R,
57 constructMatrix11(mID1, rwID, FOLC, FRIC, LC, P, V, FC, QC, RC, F, PS, R) ->
58 name(mID1, n),
59 name(mID2, "tweet"),
60 name(mID3, "user"),
61 subsetweet(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
62 subsetweet(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
63 indicator(mID2, mID3, idr3),
64 deriveFrom(idS1, idR1, mID1);
65
66 name(idS, "tweet"),
67 name(idR, "user"),
68 indicator(idS, idR, idK),
69 subsetweet(idS, "FC", "QC", "RC", "F", "PS", "R", idS1),
70 subsetweet(idR, "FOLC", "FRIC", "LC", "P", "V", idR1),
71 name(idC, "C"),
72 multi(idC, idK, idr1),
73 multi(idr1, idR1, idr3),
74 multi(idC, idS1, idr4),
75 concat(idr4, idr3, idr5) -> name(idr5, "VV1");
76
77 tweet(id),
78 childsp(id, "tweetID", "o", tID),
79 childsp(id, "flevel", "o", fID),
80 childsp(id, "entities", "o", eID),
81 childsp(eID, "hashtag", "o", hID),
82 childsp(hID, "hID", "o", hcID),
83 childsp(id, "place", "o", pID),
84 childsp(pID, "code", "o", cID) -> V1(tID, hcID, fID, cID);
85
86 tweet(id),
87 childsp(id, "tweetID", "o", tID),
88 childsp(id, "text", "o", text) ->
89 V2(idV2),
90 childsjV2(idV2, "tweetID", "o", tID),
91 childsjV2(idV2, "text", "o", text);

```

Listing G9: Forward Constraints Used for Q9

```

1 #EGDs
2 name(id1, n1), name(id2, n1) -> id2 = id1;
3 tr(id1, idt1), tr(id1, idt2) -> idt2 = idt1;
4 in(id1, idt1), in(id1, idt2) -> idt2 = idt1;
5 ins(id1, idt1), ins(id1, idt2) -> idt1 = idt2;
6 in(id1, idr1), in(idr1, idr3) -> idr3 = id1;
7 multi(idA, idB, idr1), multi(idA, idB, idr2) -> idr2 = idr1;
8 multiS(idA, idB, idr1), multiS(idA, idB, idr2) -> idr2 = idr1;
9 det(idA, d1), det(idA, d2) -> d2 = d1;
10 trace(idA, id1), trace(idA, id2) -> id2 = id1;
11 add(idA, idB, idr1), add(idA, idB, idr2) -> idr1 = idr2;
12 size(idA, i1, j1), size(idA, i2, j2) -> j2 = j1, i2 = i1;
13 name(id1, n1), name(id2, n1) -> id2 = id1;
14 name(mID1, n1), name(mID1, n2) -> n1 = n2;
15 addS(idA, idB, idr1), addS(idA, idB, idr2) -> idr1 = idr2;
16 sum(id1, idt1), sum(id1, idt2) -> idt2 = idt1;
17 rowSums(idA, id1), rowSums(idA, id2) -> id2 = id1;
18 colSums(idA, id1), colSums(idA, id2) -> id2 = id1;
19 user(EID, FOLC1, FRIC1, LC1, P1, V1), user(EID, FOLC2, FRIC2, LC2, P2, V2) -> FOLC2 = FOLC1, FRIC2 = FRIC1, LC2 = LC1, P2 = P1, V2 = V1;
20 constructMatrix2(rID, mID, v11, v12, v13, v14), constructMatrix2(rID, mID, v21, v22, v23, v24) -> v21 = v11, v22 = v12, v23 = v13, v24 = v14;
21 constructMatrix4(rID, mID, v11, v12, v13, v14), constructMatrix4(rID, mID, v17, v18, v19, v110, v111), constructMatrix4(rID, mID, v21, v22, v23, v24, v25, v26, v27, v28, v29, v210, v211) ->
v11 = v21, v12 = v22, v13 = v23, v14 = v24, v15 = v25, v16 = v26, v17 = v27, v18 = v28, v19 = v29, v110 = v210, v111 = v211;
22 constructMatrix11(rID, mID, v11, v12, v13, v14, v15, v16, v17, v18, v19, v110, v111), constructMatrix11(rID, mID, v21, v22, v23, v24, v25, v26, v27, v28, v29, v210, v211) ->
v11 = v21, v12 = v22, v13 = v23, v14 = v24, v15 = v25, v16 = v26, v17 = v27, v18 = v28, v19 = v29, v110 = v210, v111 = v211;
23 constructMtxMatrix(i, j, v, id), constructMtxMatrix(j, f2, v2, id) -> i2 = i1, j2 = j1, v2 = v1;
24 indicator(id1, id2, idr31), indicator(id1, id2, idr32) -> idr32 = idr31;
25 deriveFrom(id1, id2, idr31), deriveFrom(id1, id2, idr32) -> idr32 = idr31;
26 concat(id1, id2, idr1), concat(id1, id2, idr2) -> idr2 = idr1;
27 subsetweet(id1, id2, id3, id4, id5, id6, id7, idr1),
28 subsetweet(id1, id2, id3, id4, id5, id6, id7, idr2) -> idr2 = idr1;
29 subsetweet_r(id1, id2, id3, id4, id5, id6, idr1),
30 subsetweet_r(id1, id2, id3, id4, id5, id6, idr2) -> idr2 = idr1;
31 subsetuser(id1, id2, id3, id4, id5, id6, idr2) -> idr2 = idr1;
32 tweet(id1), tweet(id2) -> id2 = id1;
33 childsp(x, y1, z, d), childsp(x, y2, z, d) -> y2 = y1;
34 valsp(x, v1), valsp(x, v2) -> v2 = v1;
35 childsjV2(x, y1, z, d), childsjV2(x, y2, z, d) -> y2 = y1;
36
37 #TGDs
38 colSums(id1, idr1), tr(idr1, idr2), rowSums(id2, idr3), multi(idr2, idr3, idr4), sum(idr4, s) ->
39 multi(id1, id2, idr5), sum(idr5, s);
40 multi(id1, id2, idr1), colSums(idr1, idr2) -> colSums(id1, idr3), multi(idr3, id2, idr2);
41 multi(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idB, idC, idr3), multi(idA, idr3, idr2);
42 constructMtxMatrix(i, j, v, id), name(id, n), matrixFilterLT(id, x, idr1) ->
43 name(idr1, n), constructMtxMatrix(i, j, v, idr1), LT(v, x);
44 add(idA, idB, idr1), multi(idr1, idC, idr2) -> multi(idA, idC, c1), multi(idB, idC, c2), add(c1, c2, idr2);
45
46 subsetweet(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
47 subsetweet(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
48 indicator(mID2, mID3, idK) deriveFrom(idS1, idR1, id),
49 multi(idX, id, idr1) ->
50 multi(idX, idK, idr2),
51 multi(idr2, idR1, idr3),
52 multi(idX, idS1, idr4),
53 concat(idr4, idr3, idr1);
54
55 user(EID, FOLC, FRIC, LC, P, V),
56 tweetREID, FC, QC, RC, F, PS, R,
57 constructMatrix11(mID1, rWID, FOLC, FRIC, LC, P, V, FC, QC, RC, F, PS, R) ->
58 name(mID1, n),
59 name(mID2, "tweet"),
60 name(mID3, "user"),
61 subsetweet(mID2, "FC", "QC", "RC", "F", "PS", "R", idS1),
62 subsetweet(mID3, "FOLC", "FRIC", "LC", "P", "V", idR1),
63 indicator(mID2, mID3, idr3),
64 deriveFrom(idS1, idR1, mID1);
65
66 name(idS, "tweet"),
67 name(idR, "user"),
68 indicator(idS, idR, idK),
69 subsetweet(idS, idR, idK),
70 subsetweet(idR, "FOLC", "FRIC", "LC", "P", "V", idR1),
71 name(idC, "C"),
72 multi(idC, idK, idr1),
73 multi(idR1, idR1, idr3),
74 multi(idC, idS1, idr4),
75 concat(idr4, idr3, idr5) -> name(idr5, "VV1");
76
77 tweet(id),
78 childsp(id, "tweetID", "o", tID),
79 childsp(id, "flevel", "o", fID),
80 childsp(id, "entities", "o", eID),
81 childsp(eID, "hashtag", "o", hID),
82 childsp(hID, "hID", "o", hcID),
83 childsp(id, "place", "o", pID),
84 childsp(pID, "code", "o", cID) -> V1(tID, hcID, fID, cID);
85
86 tweet(id),
87 childsp(id, "tweetID", "o", tID),
88 childsp(id, "text", "o", text) ->
89 V2(idV2),
90 childsjV2(idV2, "tweetID", "o", tID),
91 childsjV2(idV2, "text", "o", text);

```

Listing G10: Forward Constraints Used for Q10

H HYBRID TWITTER BENCHMARK QUERIES AND VIEWS

```

import scala.math._
import org.apache.sysml.api.mlcontext._
import org.apache.sysml.api.mlcontext.ScriptFactory._
import org.apache.sysml.api.mlcontext.MatrixFormat._
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.SparkConf
import org.apache.spark.sql._
import scala.io.Source
import scala.tools.nsc.io._
import scala.collection.immutable._
import org.apache.spark.storage.StorageLevel._

//Preprocessing M
val sqlContext = new org.apache.spark.sql.SQLContext(sc)
val User = sqlContext.read.format("csv").option("header", "true").load("User.csv");
User.createOrReplaceTempView("User")
val Tweet = sqlContext.read.format("csv").option("header", "true").load("Tweet.csv");
Tweet.createOrReplaceTempView("Tweet")
val M = spark.sql("""SELECT followers_count,friends_count, listed_count,protected,verified,
                     favorite_count,quote_count, reply_count,retweet_count,favorited,
                     possibly_sensitive,retweeted
                     FROM User as U, Tweet as T
                     WHERE T.id = U.id""")
constructMatirx (M2,"M.csv");

//Preprocessing (N)
val TweetJSON = spark.read.json("tweets/")
TweetJSON.createOrReplaceTempView("TweetJSON")
val N = spark.sql("""SELECT TJ.id, hID, TJ.filter_level
                     FROM TweetJSON AS TJ
                     LATERAL VIEW EXPLOD (entities.hashtags.id) AS hID
                     WHERE text LIKE '%covid%' AND
                     TJ.place.country_code ='US'""")
constructMTXMatirx (2380000,1000,N,"N mtx");

//Analysis
val ml = new MLContext(spark)
ml.setConfigProperty("sysml.cp.parallel.ops","true")
ml.setConfigProperty("sysml.optlevel","4")
ml.setConfigProperty("sysml.localtmpdir","systemml")
//DML Query
val dmlTextQuery =
s"""
| M = read("M.csv", format="csv", header=FALSE, sep=",");
| N = readMM("N mtx");
| V1 = read("V1.csv", format="csv", header=FALSE, sep ",");
| U1 = read("U1.csv", format="csv", header=FALSE, sep ",");
| X = read("X.csv", format="csv", header=FALSE, sep ",");
| NF = ifelse(N<=4,N,0)
| print(as.scalar(NF[20000,400]))
| res = rowSums(X%*%M) + ((U1%*%t(V1)) + t(NF))%*%V1
| while(FALSE){}
|   print(as.scalar(res[500,1]))
| }
"""
"".stripMargin
val dmlScript = dml(dmlTextQuery)
val result = ml.execute(dmlScript)

```

Listing H1: Q1

```

import scala.math._
import org.apache.sysml.api.mlcontext._
import org.apache.sysml.api.mlcontext.ScriptFactory._
import org.apache.sysml.api.mlcontext.MatrixFormat._
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.SparkConf
import org.apache.spark.sql._
import scala.io.Source
import scala.tools.nsc.io._
import scala.collection.immutable._
import org.apache.spark.storage.StorageLevel._

//Preprocessing M
val sqlContext = new org.apache.spark.sql.SQLContext(sc)
val User = sqlContext.read.format("csv").option("header", "true").load("User.csv");
User.createOrReplaceTempView("User")
val Tweet = sqlContext.read.format("csv").option("header", "true").load("Tweet.csv");
Tweet.createOrReplaceTempView("Tweet")
val M = spark.sql("""SELECT followers_count,friends_count, listed_count,protected,verified,
                     favorite_count,quote_count, reply_count,retweet_count,favorited,
                     possibly_sensitive,retweeted
                  FROM User as U, Tweet as T
                  WHERE T.id = U.id""")
constructMatirx (M2,"M.csv");

//Preprocessing (N)
val TweetJSON = spark.read.json("tweets/")
TweetJSON.createOrReplaceTempView("TweetJSON")
val N = spark.sql("""SELECT TJ.id, hID, TJ.filter_level
                  FROM TweetJSON AS TJ
                  LATERAL VIEW EXPLOD (entities.hashtags.id) AS hID
                  WHERE text LIKE '%covid%' AND
                  TJ.place.country_code ='US'""")
constructMTXMatirx (2380000,1000,N,"N mtx");

//Analysis
val ml = new MLContext(spark)
ml.setConfigProperty("sysml.cp.parallel.ops","true")
ml.setConfigProperty("sysml.optlevel","4")
ml.setConfigProperty("sysml.localtmpdir","systemml")
//DML Query
val dmlTextQuery =
s"""
| M = read("M.csv", format="csv", header=FALSE, sep=",")
| N = readMM("N mtx")
| X = read("X.csv", format="csv", header=FALSE, sep=",");
| U = read("U.csv", format="csv", header=FALSE, sep ",");
| NF = ifelse(N<=4,N,0)
| print(as.scalar(NF[20000,400]))
| res =U%*%colSums(t(X%*%M))+NF
| while(FALSE){}
|   print(as.scalar(res[1,80]))
| }
"""
.stripMargin
val dmlScript = dml(dmlTextQuery)
val result = ml.execute(dmlScript)

```

Listing H2: Q2

```

import scala.math._
import org.apache.sysml.api.mlcontext._
import org.apache.sysml.api.mlcontext.ScriptFactory._
import org.apache.sysml.api.mlcontext.MatrixFormat._
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.SparkConf
import org.apache.spark.sql._
import scala.io.Source
import scala.tools.nsc.io._
import scala.collection.immutable._
import org.apache.spark.storage.StorageLevel._

//Preprocessing M
val sqlContext = new org.apache.spark.sql.SQLContext(sc)
val User = sqlContext.read.format("csv").option("header", "true").load("User.csv");
User.createOrReplaceTempView("User")
val Tweet = sqlContext.read.format("csv").option("header", "true").load("Tweet.csv");
Tweet.createOrReplaceTempView("Tweet")
val M = spark.sql("""SELECT followers_count,friends_count, listed_count,protected,verified,
                     favorite_count,quote_count, reply_count,retweet_count,favorited,
                     possibly_sensitive,retweeted
                     FROM User as U, Tweet as T
                     WHERE T.id = U.id""")
constructMatirx (M2,"M.csv");

//Preprocessing (N)
val TweetJSON = spark.read.json("tweets/")
TweetJSON.createOrReplaceTempView("TweetJSON")
val N = spark.sql("""SELECT TJ.id, hID, TJ.filter_level
                     FROM TweetJSON AS TJ
                     LATERAL VIEW EXPLOD (entities.hashtags.id) AS hID
                     WHERE text LIKE '%covid%' AND
                     TJ.place.country_code ='US'""")
constructMTXMatirx (2380000,1000,N,"N mtx");

//Analysis
val ml = new MLContext(spark)
ml.setConfigProperty("sysml.cp.parallel.ops","true")
ml.setConfigProperty("sysml.optlevel","4")
ml.setConfigProperty("sysml.localtmpdir","systemml")
//DML Query
val dmlTextQuery =
s"""
| M = read("M.csv", format="csv", header=FALSE, sep=",");
| N = readMM("N mtx");
| X = read("X.csv", format="csv", header=FALSE, sep=",");
| V = read("V.csv", format="csv", header=FALSE, sep=",");
| NF = ifelse(N<=4,N,0)
| print(as.scalar(NF[20000,400]))
| res = ((NF+X)%*%V)%*% colSums(M)
| while(FALSE){}
|   print(as.scalar(res[500,1]))
| }
"""
.stripMargin
val dmlScript = dml(dmlTextQuery)
val result = ml.execute(dmlScript)

```

Listing H3: Q3

```

import scala.math._
import org.apache.sysml.api.mlcontext._
import org.apache.sysml.api.mlcontext.ScriptFactory._
import org.apache.sysml.api.mlcontext.MatrixFormat._
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.SparkConf
import org.apache.spark.sql._
import scala.io.Source
import scala.tools.nsc.io._
import scala.collection.immutable._
import org.apache.spark.storage.StorageLevel._

//Preprocessing M
val sqlContext = new org.apache.spark.sql.SQLContext(sc)
val User = sqlContext.read.format("csv").option("header", "true").load("User.csv");
User.createOrReplaceTempView("User")
val Tweet = sqlContext.read.format("csv").option("header", "true").load("Tweet.csv");
Tweet.createOrReplaceTempView("Tweet")
val M = spark.sql("""SELECT followers_count,friends_count, listed_count,protected,verified
                     favorite_count,quote_count, reply_count,retweet_count,favorited,
                     possibly_sensitive,retweeted
                  FROM User as U, Tweet as T
                  WHERE T.id = U.id""")
constructMatirx (M2,"M.csv");

//Preprocessing (N)
val TweetJSON = spark.read.json("tweets/")
TweetJSON.createOrReplaceTempView("TweetJSON")
val N = spark.sql("""SELECT TJ.id, hID, TJ.filter_level
                     FROM TweetJSON AS TJ
                     LATERAL VIEW EXPLOD (entities.hashtags.id) AS hID
                     WHERE text LIKE '%covid%' AND
                     TJ.place.country_code ='US'""")
constructMTXMatirx (2380000,1000,N,"N mtx");

//Analysis
val ml = new MLContext(spark)
ml.setConfigProperty("sysml.cp.parallel.ops","true")
ml.setConfigProperty("sysml.optlevel","4")
ml.setConfigProperty("sysml.localtmpdir","systemml")
//DML Query
val dmlTextQuery =
s"""
| M = read("M.csv", format="csv", header=FALSE, sep=",")
| N = readMM("N mtx")
| X = read("X.csv", format="csv", header=FALSE, sep ",");
| U = read("U.csv", format="csv", header=FALSE, sep ",");
| v = read("v.csv", format="csv", header=FALSE, sep ",");
| NF = ifelse(N<=4,N,0)
| print(as.scalar(NF[20000,400]))
| res = sum(U+NF%*%colSums(X%*%M)%*%v)
| while(FALSE){}
|   print(res)
| }
"""
"".stripMargin
val dmlScript = dml(dmlTextQuery)
val result = ml.execute(dmlScript)

```

Listing H4: Q4

```

import scala.math._
import org.apache.sysml.api.mlcontext._
import org.apache.sysml.api.mlcontext.ScriptFactory._
import org.apache.sysml.api.mlcontext.MatrixFormat._
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.SparkConf
import org.apache.spark.sql._
import scala.io.Source
import scala.tools.nsc.io._
import scala.collection.immutable._
import org.apache.spark.storage.StorageLevel._

//Preprocessing M
val sqlContext = new org.apache.spark.sql.SQLContext(sc)
val User = sqlContext.read.format("csv").option("header", "true").load("User.csv");
User.createOrReplaceTempView("User")
val Tweet = sqlContext.read.format("csv").option("header", "true").load("Tweet.csv");
Tweet.createOrReplaceTempView("Tweet")
val M = spark.sql("""SELECT followers_count,friends_count, listed_count,protected,verified,
                     favorite_count,quote_count, reply_count,retweet_count,favorited,
                     possibly_sensitive,retweeted
                  FROM User as U, Tweet as T
                  WHERE T.id = U.id""")
constructMatirx (M2,"M.csv");

//Preprocessing (N)
val TweetJSON = spark.read.json("tweets/")
TweetJSON.createOrReplaceTempView("TweetJSON")
val N = spark.sql("""SELECT TJ.id, hID, TJ.filter_level
                  FROM TweetJSON AS TJ
                  LATERAL VIEW EXPLOD (entities.hashtags.id) AS hID
                  WHERE text LIKE '%covid%' AND
                  TJ.place.country_code ='US'""")
constructMTXMatirx (2380000,1000,N,"N mtx");

//Analysis
val ml = new MLContext(spark)
ml.setConfigProperty("sysml.cp.parallel.ops","true")
ml.setConfigProperty("sysml.optlevel","4")
ml.setConfigProperty("sysml.localtmpdir","systemml")
//DML Query
val dmlTextQuery =
s"""
| M = read("M.csv", format="csv", header=FALSE, sep=",")
| N = readMM("N mtx")
| U = read("U.csv", format="csv", header=FALSE, sep=",");
| X = read("X.csv", format="csv", header=FALSE, sep ",");
| NF = ifelse(N<=4,N,0)
| print(as.scalar(NF[20000,400]))
| res = U%*%colSums(M%*%X) + NF
| while(FALSE){}
|   print(res)
| }
"""
"".stripMargin
val dmlScript = dml(dmlTextQuery)
val result = ml.execute(dmlScript)

```

Listing H5: Q5

```

import scala.math._
import org.apache.sysml.api.mlcontext._
import org.apache.sysml.api.mlcontext.ScriptFactory._
import org.apache.sysml.api.mlcontext.MatrixFormat._
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.SparkConf
import org.apache.spark.sql._
import scala.io.Source
import scala.tools.nsc.io._
import scala.collection.immutable._
import org.apache.spark.storage.StorageLevel._

//Preprocessing M
val sqlContext = new org.apache.spark.sql.SQLContext(sc)
val User = sqlContext.read.format("csv").option("header", "true").load("User.csv");
User.createOrReplaceTempView("User")
val Tweet = sqlContext.read.format("csv").option("header", "true").load("Tweet.csv");
Tweet.createOrReplaceTempView("Tweet")
val M = spark.sql("""SELECT followers_count,friends_count, listed_count,protected,verified
                     favorite_count,quote_count, reply_count,retweet_count,favorited,
                     possibly_sensitive,retweeted
                  FROM User as U, Tweet as T
                  WHERE T.id = U.id""")
constructMatirx (M2,"M.csv");

//Preprocessing (N)
val TweetJSON = spark.read.json("tweets/")
TweetJSON.createOrReplaceTempView("TweetJSON")
val N = spark.sql("""SELECT TJ.id, hID, TJ.filter_level
                     FROM TweetJSON AS TJ
                     LATERAL VIEW EXPLOD (entities.hashtags.id) AS hID
                     WHERE text LIKE '%covid%' AND
                     TJ.place.country_code ='US'""")
constructMTXMatirx (2380000,1000,N,"N mtx");

//Analysis
val ml = new MLContext(spark)
ml.setConfigProperty("sysml.cp.parallel.ops","true")
ml.setConfigProperty("sysml.optlevel","4")
ml.setConfigProperty("sysml.localtmpdir","systemml")
//DML Query
val dmlTextQuery =
s"""
| M = read("M.csv", format="csv", header=FALSE, sep=",");
| N = readMM("N mtx");
| V1 = read("V1.csv", format="csv", header=FALSE, sep ",");
| U1 = read("U1.csv", format="csv", header=FALSE, sep ",");
| X = read("X.csv", format="csv", header=FALSE, sep ",");
| NF = ifelse(N<=4,N,0)
| print(as.scalar(NF[20000,400]))
| res = rowSums(t(M%*%X)) + (((U1%*%V1) ) + NF)%*%V1)
| while(FALSE){}
|   print(res[1000,1])
| }
"""
"".stripMargin
val dmlScript = dml(dmlTextQuery)
val result = ml.execute(dmlScript)

```

Listing H6: Q6

```

import scala.math._
import org.apache.sysml.api.mlcontext._
import org.apache.sysml.api.mlcontext.ScriptFactory._
import org.apache.sysml.api.mlcontext.MatrixFormat._
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.SparkConf
import org.apache.spark.sql._
import scala.io.Source
import scala.tools.nsc.io._
import scala.collection.immutable._
import org.apache.spark.storage.StorageLevel._

//Preprocessing M
val sqlContext = new org.apache.spark.sql.SQLContext(sc)
val User = sqlContext.read.format("csv").option("header", "true").load("User.csv");
User.createOrReplaceTempView("User")
val Tweet = sqlContext.read.format("csv").option("header", "true").load("Tweet.csv");
Tweet.createOrReplaceTempView("Tweet")
val M = spark.sql("""SELECT followers_count,friends_count, listed_count,protected,verified,
                     favorite_count,quote_count, reply_count,retweet_count,favorited,
                     possibly_sensitive,retweeted
                  FROM User as U, Tweet as T
                  WHERE T.id = U.id""")
constructMatirx (M2,"M.csv");

//Preprocessing (N)
val TweetJSON = spark.read.json("tweets/")
TweetJSON.createOrReplaceTempView("TweetJSON")
val N = spark.sql("""SELECT TJ.id, hID, TJ.filter_level
                  FROM TweetJSON AS TJ
                  LATERAL VIEW EXPLOD (entities.hashtags.id) AS hID
                  WHERE text LIKE '%covid%' AND
                  TJ.place.country_code ='US'""")
constructMTXMatirx (2380000,1000,N,"N mtx");

//Analysis
val ml = new MLContext(spark)
ml.setConfigProperty("sysml.cp.parallel.ops","true")
ml.setConfigProperty("sysml.optlevel","4")
ml.setConfigProperty("sysml.localtmpdir","systemml")
//DML Query
val dmlTextQuery =
s"""
| M = read("M.csv", format="csv", header=FALSE, sep=",")
| N = readMM("N mtx")
| X = read("X.csv", format="csv", header=FALSE, sep=",");
| U = read("U.csv", format="csv", header=FALSE, sep ",");
| NF = ifelse(N<=4,N,0)
| print(as.scalar(NF[20000,400]))
| res= X%*%NF%*%U + rowSums(t(M))
| while(FALSE){}
| print(res[2,30])
|
"""
"".stripMargin
val dmlScript = dml(dmlTextQuery)
val result = ml.execute(dmlScript)

```

Listing H7: Q7

```

import scala.math._
import org.apache.sysml.api.mlcontext._
import org.apache.sysml.api.mlcontext.ScriptFactory._
import org.apache.sysml.api.mlcontext.MatrixFormat._
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.SparkConf
import org.apache.spark.sql._
import scala.io.Source
import scala.tools.nsc.io._
import scala.collection.immutable._
import org.apache.spark.storage.StorageLevel._

//Preprocessing M
val sqlContext = new org.apache.spark.sql.SQLContext(sc)
val User = sqlContext.read.format("csv").option("header", "true").load("User.csv");
User.createOrReplaceTempView("User")
val Tweet = sqlContext.read.format("csv").option("header", "true").load("Tweet.csv");
Tweet.createOrReplaceTempView("Tweet")
val M = spark.sql("""SELECT followers_count,friends_count, listed_count,protected,verified
                     favorite_count,quote_count, reply_count,retweet_count,favorited,
                     possibly_sensitive,retweeted
                  FROM User as U, Tweet as T
                  WHERE T.id = U.id""")
constructMatirx (M2,"M.csv");

//Preprocessing (N)
val TweetJSON = spark.read.json("tweets/")
TweetJSON.createOrReplaceTempView("TweetJSON")
val N = spark.sql("""SELECT TJ.id, hID, TJ.filter_level
                     FROM TweetJSON AS TJ
                     LATERAL VIEW EXPLOD (entities.hashtags.id) AS hID
                     WHERE text LIKE '%covid%' AND
                     TJ.place.country_code ='US'""")
constructMTXMatirx (2380000,1000,N,"N mtx");

//Analysis
val ml = new MLContext(spark)
ml.setConfigProperty("sysml.cp.parallel.ops","true")
ml.setConfigProperty("sysml.optlevel","4")
ml.setConfigProperty("sysml.localtmpdir","systemml")
//DML Query
val dmlTextQuery =
s"""
| M = read("M.csv", format="csv", header=FALSE, sep=",")
| N = readMM("N mtx")
| C = read("C.csv", format="csv", header=FALSE, sep ",");
| V = read("V.csv", format="csv", header=FALSE, sep ",");
| X = read("X.csv", format="csv", header=FALSE, sep ",");
| NF = ifelse(N<=4,N,0)
| print(as.scalar(NF[20000,400]))
| res = NF*trace(C+ V%*%colSums(M%*%X)%*%C)
| while(FALSE){}
|   print(res[2,30])
| }
"""
"".stripMargin
val dmlScript = dml(dmlTextQuery)
val result = ml.execute(dmlScript)

```

Listing H8: Q8

```

import scala.math._
import org.apache.sysml.api.mlcontext._
import org.apache.sysml.api.mlcontext.ScriptFactory._
import org.apache.sysml.api.mlcontext.MatrixFormat._
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.SparkConf
import org.apache.spark.sql._
import scala.io.Source
import scala.tools.nsc.io._
import scala.collection.immutable._
import org.apache.spark.storage.StorageLevel._

//Preprocessing M
val sqlContext = new org.apache.spark.sql.SQLContext(sc)
val User = sqlContext.read.format("csv").option("header", "true").load("User.csv");
User.createOrReplaceTempView("User")
val Tweet = sqlContext.read.format("csv").option("header", "true").load("Tweet.csv");
Tweet.createOrReplaceTempView("Tweet")
val M = spark.sql("""SELECT followers_count,friends_count, listed_count,protected,verified,
                     favorite_count,quote_count, reply_count,retweet_count,favorited,
                     possibly_sensitive,retweeted
                  FROM User as U, Tweet as T
                  WHERE T.id = U.id""")
constructMatirx (M2,"M.csv");

//Preprocessing (N)
val TweetJSON = spark.read.json("tweets/")
TweetJSON.createOrReplaceTempView("TweetJSON")
val N = spark.sql("""SELECT TJ.id, hID, TJ.filter_level
                  FROM TweetJSON AS TJ
                  LATERAL VIEW EXPLOD (entities.hashtags.id) AS hID
                  WHERE text LIKE '%covid%' AND
                  TJ.place.country_code ='US'""")
constructMTXMatirx (2380000,1000,N,"N mtx");

//Analysis
val ml = new MLContext(spark)
ml.setConfigProperty("sysml.cp.parallel.ops","true")
ml.setConfigProperty("sysml.optlevel","4")
ml.setConfigProperty("sysml.localtmpdir","systemml")
//DML Query
val dmlTextQuery =
s"""
| M = read("M.csv", format="csv", header=FALSE, sep=",")
| N = readMM("N mtx")
| C = read("C.csv", format="csv", header=FALSE, sep=",");
| X = read("X.csv", format="csv", header=FALSE, sep ",");
| NF = ifelse(N<=4,N,0)
| print(as.scalar(NF[20000,400]))
| res = X*sum(t(colSums(C)*rowSums(M))+NF
| while(FALSE){}
|   print(res[2,30])
| }
"""
.stripMargin
val dmlScript = dml(dmlTextQuery)
val result = ml.execute(dmlScript)

```

Listing H9: Q9

```

import scala.math._
import org.apache.sysml.api.mlcontext._
import org.apache.sysml.api.mlcontext.ScriptFactory._
import org.apache.sysml.api.mlcontext.MatrixFormat._
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.SparkConf
import org.apache.spark.sql._
import scala.io.Source
import scala.tools.nsc.io._
import scala.collection.immutable._
import org.apache.spark.storage.StorageLevel._

//Preprocessing M
val sqlContext = new org.apache.spark.sql.SQLContext(sc)
val User = sqlContext.read.format("csv").option("header", "true").load("User.csv");
User.createOrReplaceTempView("User")
val Tweet = sqlContext.read.format("csv").option("header", "true").load("Tweet.csv");
Tweet.createOrReplaceTempView("Tweet")
val M = spark.sql("""SELECT followers_count,friends_count, listed_count,protected,verified,
                     favorite_count,quote_count, reply_count,retweet_count,favorited,
                     possibly_sensitive,retweeted
                  FROM User as U, Tweet as T
                  WHERE T.id = U.id""")
constructMatirx (M2,"M.csv");

//Preprocessing (N)
val TweetJSON = spark.read.json("tweets/")
TweetJSON.createOrReplaceTempView("TweetJSON")
val N = spark.sql("""SELECT TJ.id, hID, TJ.filter_level
                  FROM TweetJSON AS TJ
                  LATERAL VIEW EXPLOD (entities.hashtags.id) AS hID
                  WHERE text LIKE '%covid%' AND
                  TJ.place.country_code ='US'""")
constructMTXMatirx (2380000,1000,N,"N mtx");

//Analysis
val ml = new MLContext(spark)
ml.setConfigProperty("sysml.cp.parallel.ops","true")
ml.setConfigProperty("sysml.optlevel","4")
ml.setConfigProperty("sysml.localtmpdir","systemml")
//DML Query
val dmlTextQuery =
s"""
| M = read("M.csv", format="csv", header=FALSE, sep=",")
| N = readMM("N mtx")
| C = read("C.csv", format="csv", header=FALSE, sep=",");
| X = read("X.csv", format="csv", header=FALSE, sep ",");
| NF = ifelse(N<=4,N,0)
| print(as.scalar(NF[20000,400]))
| res = NF*sum((X+C)%*%M)
| while(FALSE){}
|   print(res[2,30])
| }
"""
.stripMargin
val dmlScript = dml(dmlTextQuery)
val result = ml.execute(dmlScript)

```

Listing H10: Q10

```

T<-read.csv("Tweet.csv")
U<-read.csv("User.csv")
K <-sparseMatrix(i=1:nrow(U), j=S[, "id"], x=1)
T1 <- as.matrix(subset(T, select = c("favorite_count", "quote_count", "reply_count",
                                     "retweet_count", "favorited", "possibly_sensitive", "retweeted" )))
U1 <-as.matrix(subset(U, select = c("followers_count", "friends_count", "listed_count", "protected", "verified" )))
V3<-rowSums(T1)+K%*%rowSums(U1)
write.csv(as.matrix(V3), "V3.csv", row.names = FALSE)

```

Listing H11: V3

```

T<-read.csv("Tweet.csv")
U<-read.csv("User.csv")
K <-sparseMatrix(i=1:nrow(U), j=S[, "id"], x=1)
T1 <- as.matrix(subset(T, select = c("favorite_count", "quote_count", "reply_count",
                                     "retweet_count", "favorited", "possibly_sensitive", "retweeted" )))
U1 <-as.matrix(subset(U, select = c("followers_count", "friends_count", "listed_count", "protected", "verified" )))
V4<-cbind(colSums(T1),colSums(K)%*%U1)
write.csv(as.matrix(V4), "V4.csv", row.names = FALSE)

```

Listing H12: V4

```

T<-read.csv("Tweet.csv")
U<-read.csv("User.csv")
C<-read.csv("C.csv")
K <-sparseMatrix(i=1:nrow(U), j=S[, "id"], x=1)
T1 <- as.matrix(subset(T, select = c("favorite_count", "quote_count", "reply_count",
                                     "retweet_count", "favorited", "possibly_sensitive", "retweeted" )))
U1 <-as.matrix(subset(U, select = c("followers_count", "friends_count", "listed_count", "protected", "verified" )))
V4<-cbind(C%*%T1,(C%*%K)%*%U1)
write.csv(as.matrix(V5), "V5.csv", row.names = FALSE)

```

Listing H13: V5