

Efficient ML Model Training over Silos: to Factorize or to Materialize

ABSTRACT

This study aims to address the challenge of training machine learning (ML) models using data from disparate sources. The existing approaches for this problem include materialization and factorization. Materialization involves integrating all source tables before loading the join result into an ML tool, while factorization enables pushing down linear algebra operators to underlying source tables to enhance efficiency. However, existing solutions fall short in two areas when dealing with data sources outside a single database. First, they lack adequate support for the complex relationships between source tables and the desired target table, which are typically defined by schema mappings and row matching. Second, there is no systematic method for distinguishing between cases when materialization or factorization is more efficient.

To address these challenges, we propose *Ilargi*, a system that executes the algorithmic operations of a linear ML model across disparate source tables. *Ilargi* covers a broad range of typical dataset relationships in ML use cases, representing them in matrix format, which leads to a unified execution of data transformation operations and linear algebra operations. With the aid of first-order logic-based pruning rules and an ML-based cost estimation module, *Ilargi* can automatically choose between materialization and factorization with 90-96% accuracy depending on the dataset. Through intensive experiments, we demonstrate that *Ilargi* effectively improves the time-wise efficiency of model training.

1 INTRODUCTION

When the training data of a machine learning model comes from different sources, is it faster to train the model over the source tables or the materialized join result?

Training data from different sources. The goal of ML methods is to automatically extract information from training data for accurate predictions over unseen data [22, 57]. In real world applications, data is often not stored in a central database or file system, but spread over different data silos. Take, for instance, drug risk prediction: features can reside in datasets collected from clinics, hospitals, pharmacies, and laboratories distributed geographically [11]. Another example is training models for keyboard stroke prediction: training requires data from millions of phones [37].

Data integration. When training data is spread over different sources, the common practice is to integrate them into one training dataset, as a preprocessing step of ML pipeline [70, 62]. Fig. 1 illustrates an example of predicting the mortality (binary classification) of patients based on tables, S_1 and S_2 , maintained by different departments in the same hospital. Data integration (DI) systems facilitate interoperability among multiple, heterogeneous sources, and provide a unified view for users [36, 33]. A fundamental principle underlying DI systems is the description of data sources and their relationships [24]: *i*) mappings between different source schemata, i.e., schema matching and mapping [63, 30] and *ii*) linkages between data instances, i.e., data matching (also known as record linkage

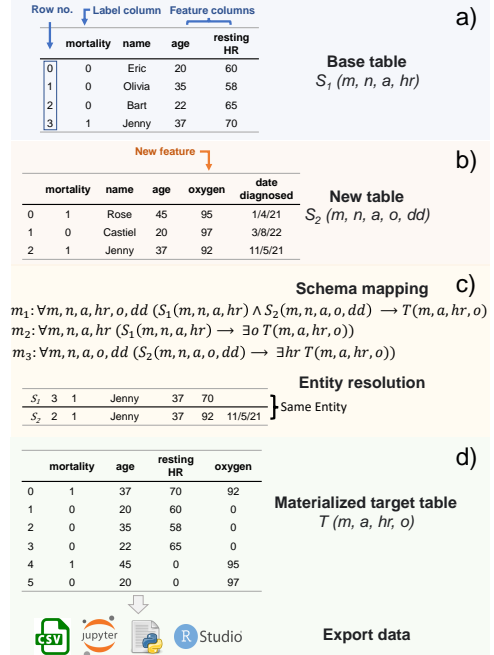


Figure 1: Data integration tasks before model training

or entity resolution) [12]. We refer to such information obtained through data integration as *data integration metadata*.

Materialization as common practice. With the DI metadata, the exact choice of merging source datasets may depend on the source data volume, integration scenarios, and technology stack, etc. A data scientist or ML expert trains models with popular ML frameworks such as scikit-learn, PyTorch, TensorFlow, and Keras. It is straightforward to use Python libraries such as Pandas to merge datasets, when we only have a few CSV files and simple column and row matching. With large-scale datasets stored in databases or data lakes [32, 35], merging the datasets may require dedicated transformation scripts or data integration tools and services, e.g., Informatica [41], Pentaho [60], Talend [71], SSIS [69]. Data science processes are exploratory and interactive [62, 70, 66], possibly demanding more data preparation and integration tasks such as feature engineering and selection [72, 34], normalization and binding [43], ontology matching [27, 6], data cleaning [52, 65, 8]. All these solutions for generating the training dataset, have one thing in common: source datasets need to be merged and materialized, as shown in Figure 1d.

Factorization as a new paradigm. Pushing ML models over joins is coined as the problem of factorized learning and has been intensively studied [48, 49, 67, 44, 14, 3, 42, 53, 17]. In a nutshell, factorized learning can be seen as an extended application of data dependency and redundancy studies over relational models. To efficiently train linear ML models over multiple tables of a single database, *factorized learning* [48] has been proposed, also known

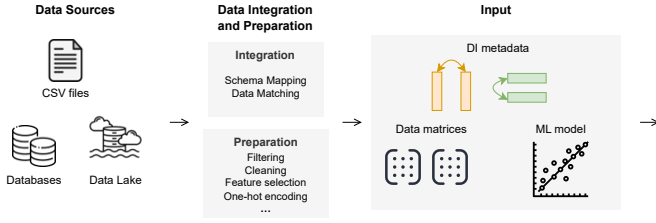


Figure 2: Ilargi's approach overview.

as *learning over factorized joins* [67]. The training process of an ML model requires complex arithmetic computations of linear algebra (LA) operators. Given an ML model and joinable tables of a database, a general approach [14] is proposed to reformulate LA operators and push down their computation to each table (see Sec. 2.1). Compared to materialization, factorized learning does not affect model training accuracy but improves time-wise efficiency.

Research challenges. It is natural to ask whether factorized learning can help with ML model training over *multiple* sources instead of a single database. The main challenges are two-fold:

(C1) Typical integration scenarios of multiple sources are not covered. When two input tables have a primary key-foreign key (PK-FK) constraint, the result of the join often contains redundant data. Factorized learning [14] proposes to execute the computation of LA operators over input tables instead of their join results. The approach also supports multi-table joins over PK-FK constraints or natural joins over shared join columns. However, as illustrated in Table 1, when datasets come from different sources, there are more possibilities for how to join them, depending on the ML applications. We are still missing a comprehensive and formal analysis of possible relationships between source datasets for ML use cases.

(C2) Factorized learning does not always yield speedups in integration scenarios. Factorized learning achieves speedup when joining tables leads to redundancy in the training dataset. However, data sources may have overlapping values, e.g., a patient's personal details are recorded in two hospitals. That is, source tables may contain redundant data. The question arises: given a downstream ML model and source datasets, is it faster to train over the materialized result, or to factorize by computing LA operations within each data source? The choice is not obvious, as we will show in Sec. 5.

Our solution. In this work, we focus on a common scenario: training data comes from multiple sources. To train ML models efficiently, we address the above two challenges. i) To enable factorization over sources, we formalize possible integration scenarios and define matrix-based representations for DI metadata, such that we can develop LA rewriting rules over source datasets. ii) We leverage materialization and factorization and choose the more efficient option. We reveal that the choice depends on three factors: data (incl. data integration metadata), ML algorithms (LA operators), and hardware configuration (e.g., number of CPU cores). As shown in Fig. 2, we propose an approach that *represents* dataset relationships as matrices, *computes* operations of data transformation and model training, and *optimizes* the workflow by choosing between materialization and factorization. We summarize our contributions as follows.

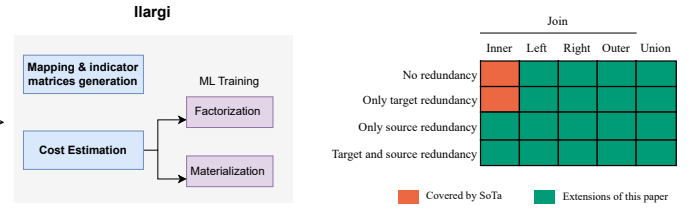


Figure 3: Ilargi vs. SoTA [14].

- **Problem formalization** (Sec. 2). We analyze the scenarios of integrating multiple source datasets for a training dataset with a widely used mapping formalism, i.e., tuple-generating dependencies [7, 28]. We identify research gaps for extending factorization from a single database to multiple sources as shown in Fig. 3.
- **Representation & computation** (Sec. 3). We propose *Ilargi*, a system that tensorizes DI metadata, i.e., represents schema mappings and data matching as matrices. *Ilargi* provides a unified computation of data transformation and linear algebra operations, which enables factorization over source datasets.
- **Optimization** (Sec. 4). We identify the relevant factors affecting the choice of materialization and factorization. We propose an estimator based on a tree boosting algorithm, which combines these factors and effectively predicts the decision boundary of materialization and factorization with up to 96% accuracy.
- **Experiments** (Sec. 5). Through an extensive evaluation using datasets with diverse characteristics, we show that our cost estimation method is effective in telling apart cases of factorization and materialization over both synthetic and real datasets.

2 DATA INTEGRATION AND FACTORIZATION

We first explain the existing strategies of factorization. To represent schema-level relationships in typical integration scenarios for ML use cases, we discuss our formal framework based on first-order logic. We analyze the gaps and reveal the core research questions.

2.1 Preliminary: Factorization

The training process of an ML model requires complex arithmetic computations. Materialization and factorization are the two possibilities for conducting LA computations of a linear ML model over separated tables from different sources. In Sec. 1 we have discussed common practices for materialization, e.g., through joins in SQL transformation scripts or Python libraries.

Factorization. Given an ML model and joinable tables of a database, factorized learning [48, 14] is the process of reformulating the execution rules of LA operations of the ML model and pushing down the computation to each table. Compared to materialization, factorized learning does not affect model training accuracy, as the arithmetic computation results of the original operators will remain the same. Similar to [14], in this work, we focus on linear regression, logistic regression, K-Means, and Gaussian Non-Negative Matrix Factorization. Gradient descent is used in both linear regression and logistic regression algorithms. For simplicity, we use linear regression to explain factorized learning.

Rewriting rules for arithmetic operations. Given two tables S and R connected by PK-FK relationship, they can be joined to obtain the target table T , i.e., $T \leftarrow S \bowtie R$. Given a LA operator over target

Table 1: Four example integration scenarios.

No.	Dataset Relationship	Schema mappings
1	Full outer join	$m_1 : \forall m, n, a, hr, o, dd (S_1(m, n, a, hr) \wedge S_2(m, n, a, o, dd) \rightarrow T(m, a, hr, o))$ $m_2 : \forall m, n, a, hr (S_1(m, n, a, hr) \rightarrow \exists o T(m, a, hr, o))$ $m_3 : \forall m, n, a, o, dd (S_2(m, n, a, o, dd) \rightarrow \exists hr T(m, a, hr, o))$
2	Inner join	$m_1 : \forall m, n, a, hr, o, dd (S_1(m, n, a, hr) \wedge S_2(m, n, a, o, dd) \rightarrow T(m, a, hr, o))$
3	Left join	$m_1 : \forall m, n, a, hr, o, dd (S_1(m, n, a, hr) \wedge S_2(n, a, o, dd) \rightarrow T(m, a, hr, o))$ $m_2 : \forall m, n, a, hr (S_1(m, n, a, hr) \rightarrow \exists o T(m, a, hr, o))$
4	Union	$m_2 : \forall m, n, a, hr, o (S_1(m, n, a, hr, o) \rightarrow T(m, a, hr, o))$ $m_3 : \forall m, n, a, hr, o, dd (S_2(m, n, a, hr, o, dd) \rightarrow T(m, a, hr, o))$

table T , and joinable tables S and R , the rewriting rules of factorizing the LA operator and creating new LA operations over S and R are proposed in the framework *Morpheus* [14]. After transforming S and R into their matrix form, their row-matching relationship can be defined as a matrix. This is referred to as *indicator matrix*, we elaborate its definition in Sec. 3.1. Morpheus covers the common LA operators used in linear ML models. The LA operators that benefit from the speed-up brought by factorization, can be divided into four groups, element-wise scalar operators (binary arithmetic operations \odot such as $+$, $-$, $*$, $/$ and \wedge , transpose T^T , etc.), aggregation (sum up matrix elements by row/column/all), matrix multiplication, and matrix inversion. Consider the linear regression. It contains two LA operators that may benefit from factorization, transpose T^T , and left matrix multiplication TX .

Algorithm 1 Linear regression using Gradient Descent ([14])

Input: T, y, w, γ
for $i \in 1 : n$ **do**
 $w = w - \gamma(T^T((Tw) - Y))$
end for

We showcase the rewriting rules of factorizing two operators from Morpheus: transpose T^T (x is a scalar, and the transpose is implemented using a flag), and left matrix multiplication TX (X is a model or weight vector). I_1 and I_2 are the indicator matrices of S and R . c_S , c_R and c_T are the column numbers of tables S, R, T .

$$T^T \odot x \rightarrow (T \odot x)^T$$

$$TX \rightarrow I_1(SX[1 : c_S,]) + I_2(RX[c_R + 1 : c_T,])$$

Research gap. Existing factorized learning works [48, 49, 67, 44, 14, 3, 42, 53, 17] are applicable over multiple tables joinable with foreign key constraints (a special case of inclusion dependencies) or join dependencies. Inclusion dependencies and join dependencies are developed in the context of schema design given a single database schema. To describe relationships of multiple source datasets, we introduce the widely used schema mapping formalism next.

2.2 Unifying Formalism for DI & Factorization

Schema mappings lay at the heart of data integration. Schema mappings are semantic descriptions of the contents of data sources and the relationships between data source schemas and the target schema. One of the most commonly used mapping languages is *source-to-target tuple generating dependencies (s-t tgd)* [7, 28], which are also known as Global-Local-as-View (GLAV) assertions [51].

Tgd definition. Let S and T be a source relational schema and a target relational schema sharing no relation symbols. A *schema mapping* M between S and T is a triple $M = \langle S, T, \Sigma \rangle$, where Σ is a set of dependencies over $\langle S, T \rangle$. Σ can be expressed as logical formulas over source and target schemas. An s-t tgd is a first-order sentence in the form of $\forall x (\varphi(x) \rightarrow \exists y \psi(x, y))$, where $\varphi(x)$ is a conjunction of atomic formulas over the source schema S , and $\psi(x, y)$ is a conjunction of atoms over the target schema T .

EXAMPLE 2.1. Consider the example of integrating S_1 and S_2 in Fig 1 via a full outer join. The dataset relationship can be expressed by three tgds, i.e., m_1 , m_2 , and m_3 . We represent mapped attributes with the same variable names, e.g., $S_1.m$ and $S_2.m$. The tgd m_1 specifies that the overlapped rows of S_1 and S_2 are added to T (\wedge denotes a natural join between S_1 and S_2); m_2 and m_3 indicate that all rows of S_1 and S_2 will be transformed to generate new tuples in T , respectively. Among the three tgds, it is the union relationship. The three tgds together describe that the instances in T are obtained via a full outer join between the datasets S_1 and S_2 .

Tgds for unifying formal framework. Besides schema mappings, tgds can be used to describe inclusion dependencies and join dependencies [29] for existing factorization works. Next, we employ tgds to describe the dataset relationships in different scenarios where source datasets are merged into the training dataset.

2.3 Integration Scenario for Training Dataset

An *integration scenario* [5] is a triple of source schemas S , target schema T , and their schema mappings M that are typically specified by tgds. In this work, we consider the schema of the training dataset as the target schema T , and use an integration scenario to describe the dataset relationships between the source and training datasets.

Existing factorization works, including state-of-the-art [14], tackle inner joins, as elaborated in Sec. 6. Our approach covers more ways of merging source datasets, i.e., left joins, full outer joins, and unions. In Table 1, we list example integration scenarios covering these four basic types of dataset relationships.¹ Understanding the dataset relationships is useful for many ML use case where training data comes from silos. Next, for clarity, we explain the four types of dataset relationships and the schema mappings in Table 1, through the use case of feature augmentation.

Use case. *Feature augmentation* is the exploratory process of finding new datasets and selecting features that help improve the ML model

¹By composing the basic scenarios in Table 1, we obtain more complicated integration scenarios. For instance, given three datasets S_1 , S_2 and S_3 , a data scientist can first perform a left join over S_1 and S_2 , then an inner join between the join result and S_3 . With user-defined functions (UDFs), the basic integration scenarios can be extended to fuzzy joins or rule-based transformations, which is more common in practice. Such extensions are straightforward, once the basic integration scenarios are understood.

performance [18, 26, 54]. Figure 1 shows an example: starting from a base table S_1 , a data scientist augments the features by introducing the table S_2 and selecting the new feature o (*oxygen*). Example 1 (full outer join) is explained in Example 2.1.

Example 2 (inner join) represents the integration scenario where only overlapped rows in two sources will be transformed, i.e., an inner join between S_1 and S_2 followed by a projection on columns m, a, hr, o . It can be used to describe the feature augmentation processes where fewer missing values are preferred.

Example 3 (left join) shows a left join between S_1 and S_2 . Compared to Example 1, we slightly change the schema of S_2 by dropping the label column m . Example 3 describes another typical feature augmentation scenario for supervised learning: only the base table S_1 contains the label column. Thus, when adding features from the new table S_2 , only rows overlapped with S_1 will be selected.

Example 4 (union) is a special case of Example 1, where S_1 and S_2 do not share any rows. We modify the schemas of S_1 and S_2 such that they share the same set of feature columns which are mapped to the target schema T . Example 4 can represent the scenario when a new table is selected to bring more data samples.

2.4 Open Problems & Challenges

Representation and execution. Data sources are autonomous, and their schemata are designed independently [24]. It is non-trivial to extend existing factorization work [14] to aforementioned integration scenarios, due to the following gaps.

i) *Joins*. As shown in Table 1, the source tables can be combined in different manners depending on ML application. Only inner-joins are discussed in [14]. On the logical level, the inner-join of multiple source tables over the same join column, can be expressed as a single tgds, while the rest of the three cases may require the union of multiple tgds as in Table 1.

ii) *Row-overlap*. Inclusion dependencies, e.g., foreign key constraints, indicate a containment relationship between instances coming from two tables. In contrast, two joinable or unionable columns from two sources (expressed as the same variables in the left-hand-side of an s-t tgd) may or may not have instance overlaps.

iii) *Column-overlap*. Morpheus has only discussed representation for row matching and does not have a declarative representation for complex scenarios of column overlap. For instance, tables are joined with a single join column, in PK-FK table joins or multi-table joins in [14]. As shown in Table 1, in integration scenarios, it is possible for two source tables to share more than one common attribute as the result of schema matching.

To summarize, the dataset relationships in integration scenarios are more complicated than PK-FK or multi-table joins in the problem setting of Morpheus [14], which leads to the below question.

Research question Q1

How to represent and utilize integration metadata, e.g., schema mapping and data matching, to enable factorization over different integration scenarios?

Optimization. As outlined in Sec. 2.1, we can compute LA operations in two manners: materialization or factorization. To choose the more efficient strategy, the key is to understand *redundancy*

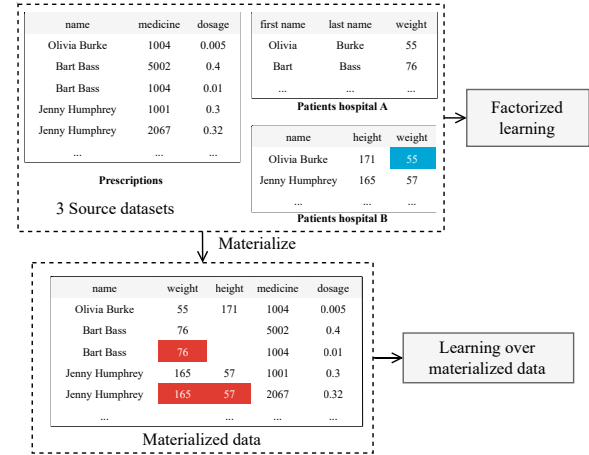


Figure 4: Running example: source redundancy (blue) vs. target redundancy (red).

and *sparsity*.

Redundancy. Normalization is one of the most fundamental database concepts. When designing a relational table schema, to store less redundant data, a good database designer normalizes the table through lossless join decomposition, leveraging functional dependencies, join dependencies, and inclusion dependencies. Factorized learning works the same way: when joinable tables have PK-FK relationship, the join result (i.e., the target table) may have more data redundancy. We refer to such redundancy as *target redundancy*. This is the reason why, when executing the same LA operator, the overhead of executing it on the target table is larger than executing it on the source tables. However, the redundancy analysis in integration scenarios is more complicated than in the cases of a single database. First, depending on the schema mapping and entity resolution results, source tables may have overlapping columns and rows. We refer to such redundancy as *source redundancy*. Fig. 4 shows an example: for the patient *Olivia Burke*, their record of weight exist in two hospital databases. As shown in Fig. 3, we aim to tackle both types of redundancies.

Sparsity. By *sparsity*, we refer to elements with values of zero in the matrices that participate in linear algebra computation. A matrix element is zero, either the data value is zero, or transformed to zero during data preprocessing, e.g., one-hot encoding or null value replacement. In integration scenarios, an instance in the target table may have missing values, if no corresponding values exist in the source table. In Sec. 4, we will reveal that the sparsity plays an important role in choosing between materialization and factorization.

Research question Q2

Given data sources and an ML algorithm, how to choose between materialization and factorization?

3 ILARGI: ML OVER DISPARATE SOURCES

For research question Q1, we propose matrix-based representations of DI metadata, and rewrite rules for factorizing ML models. For Q2, we elaborate on cost estimation in Sec. 4.

Overview. Fig. 5 illustrates the workflow of Ilargi with three inputs: source datasets in matrices, a user-defined ML algorithm (e.g.,

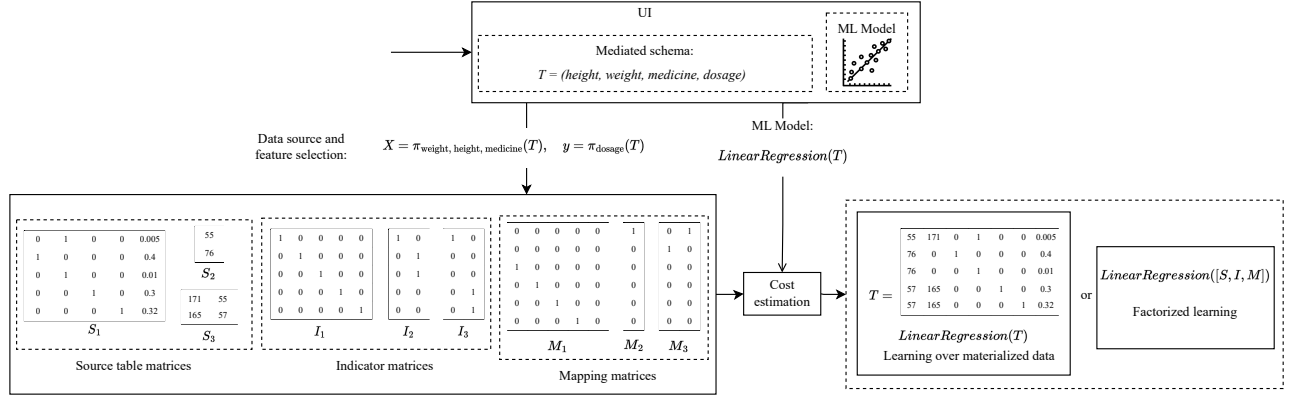


Figure 5: Ilargi workflow and intermediate results of running example.

Python scripts), and metadata about the source datasets and their data integration metadata, e.g., row matches and column matches. The output of Ilargi is the trained ML model. Ilargi has three main functions. First, given the input source datasets and their metadata, Ilargi generates their matrix-based representation (Sec.3.1). Such matrix-based representations enable a unified computation of data transformation and linear algebra operations (Sec.3.2). Second, with the metadata, the estimator decides to factorize or materialize based on the cost estimation (Sec. 4). Finally, the ML model is trained in the chosen strategy.

Running example. Fig. 4 illustrates a running example extended from Fig. 1. A pharmacist tries to build a linear regression model to predict medicine dosage. Data is from three source tables S_1 , S_2 , and S_3 with patients' weight, height, and medicine (ID) as features.

3.1 DI Metadata as Matrices

Prior to our approach, source tables are preprocessed and transformed into matrices. Fig. 5 shows the matrix form of three source tables of the running example.

Mapping matrices. We represent schema mappings with a set of mapping matrices, denoted as M . Given a source table S_k and target table T , a mapping matrix $M_k \in M$ of size $c_T \times c_k$, describes how the columns of S_k mapped to the columns of T .

Definition 1 (Mapping matrix). Mapping matrices between source tables S_1, S_2, \dots, S_n and target table T are a set of binary matrices $M = \{M_1, \dots, M_n\}$. M_k ($k \in [1, n]$) is a matrix with the shape $c_T \times c_k$, where

$$M_k[i, j] = \begin{cases} 1, & \text{if } j^{\text{th}} \text{ column of } S_k \text{ is mapped to the } i^{\text{th}} \text{ column of } T \\ 0, & \text{otherwise} \end{cases}$$

Intuitively, in $M_k[i, j]$ the vertical coordinate i represents the target table column while the horizontal coordinate j represents the mapped source table column. A value of 1 in M_k specifies the existence of column correspondences between S_k and T , while the value 0 shows that the current target table attribute has no corresponding column in S_k . Fig. 5 shows the mapping matrices M_1, M_2, M_3 for source tables S_1, S_2, S_3 , respectively.

Indicator matrices. We use the *indicator matrix* [14] to preserve the row matching between each source table S_k and the target table T . An indicator matrix I_k of size $c_T \times c_k$ describes how the rows of source table S_k map to the rows of target table T , as in Fig. 5.

Definition 2 (Indicator matrix [14]). Indicator matrices between source tables S_1, S_2, \dots, S_n and target table T are a set of row vectors $I = \{I_1, \dots, I_n\}$. I_k ($k \in [1, n]$) is a row vector of size r_T , where

$$I_k[i, j] = \begin{cases} 1 & \text{if the } j\text{-th row of } S_k \text{ is mapped to the } i\text{-th row of } T \\ 0 & \text{otherwise} \end{cases}$$

Implementation. Mapping and indicator matrices are simple, easy-to-implement representations for schema mapping and row matching. A key observation is that they are often sparse. After trying out straightforward implementation following the above definitions, and alternative data structures in Python, our current implementation of mapping and indicator matrices is in sparse matrix format, SciPy Compressed Sparse Row Format (CSR)², which copes with the matrix sparsity problem, and improves performance in later steps of factorization/materialization.

3.2 Rewriting Rules for Factorization

With data and metadata represented in matrices, next, we explain how to rewrite a linear algebra over target schema to linear algebras over source schemas. The rewriting rules, although for arithmetic operations, share similar principles of view-based query rewriting [23, 19]. Here we use the example of LA operator *left matrix multiplication* (LMM). The full set of LA rewrite rules based on mapping and indicator matrices is elaborated in our technique report [25].

Left Matrix Multiplication. Given a matrix X with the size $c_T \times c_X$, the LMM of T and X is denoted as TX . The result of LMM between our target table matrix and another matrix X is a matrix of size $r_T \times c_X$. Our rewrite of LMM goes as follows.

$$TX \rightarrow I_1 S_1 M_1^T X + \dots + I_n S_n M_n^T X$$

We first compute the local result $I_k S_k M_k^T X$ ($k \in [1, n]$) for each source table, which are assembled for the final results.

Implementation-level optimization. To reduce computation overhead, our implementation adopts the optimal matrix multiplication order algorithm [20] from the version included in NumPy³ for SciPy matrices, computed using matrix dimensions.

Why mapping and indicator matrices? In Sec. 2.4 we mentioned the gaps regarding representation and execution. As discussed in Sec. 2.3, at the logical level, with schema mappings we bridge the

²https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.csr_matrix.html

³https://numpy.org/doc/stable/reference/generated/numpy.linalg.multi_dot.html

gaps, i.e., we can merge tables more than just inner joins, and have multiple matched columns, not necessarily follow set containment relationships. At the computation level, with schema mapping and data matching encoded in matrices, the local data transformation is also computed in linear algebra, e.g., $I_k S_k M_k^T$ in the above example. It facilitates unified, LA-based rewrite rules of factorization.

Answer to Q1

We represent DI metadata, e.g., schema mapping and row matchings as matrices, which brings a unified computation with linear algebra operations for factorization over sources.

Applicability, assumptions, and extensibility. Fig. 3 summarizes Ilargi’s applicability compared to SOTA solution [14]. Ilargi can handle both normalized and denormalized tabular data. To prepare the training data, feature and label columns are from source datasets, which means that for each column in the target table, we can find its corresponding column from at least one source schema. Similar to [14], our input source table matrices, e.g., S_1, S_2, S_3 in Fig. 5, are source datasets after preprocessing, e.g., in NumPy arrays. The schema mapping matrices are built between the schemas of source table matrices and the target schema of the training dataset. Our running example and current implementation replace null values to zero, and apply one-hot-encoding for categorical variables. It is straightforward to extend the implementation to replace the null values to mean/median values or user-specified values, and other methods for handling categorical variables [43].

4 COST ESTIMATION IN ILARGI

We illustrate the research question Q2 intuitively in Fig. 6. We explain Ilargi’s optimizer depicted in Fig. 8, which performs cost estimation given source datasets, DI metadata, the ML model, and the hardware configuration.

4.1 Schema Mappings as Pruning Rules

Formal framework. We use schema mappings to specify the schema-wise dataset relationships in our system. A *schema mapping* M between source schemas S and target schema T is a triple $M = \langle S, T, \Sigma \rangle$, where Σ is a set of dependencies over $\langle S, T \rangle$. The dependencies Σ can be expressed as tgds, as shown in Table 1. We refer to an instance I over S as the source instance, and an instance J over T as the target instance. Intuitively, if storing the source instance I takes more space compared to J , the optimizer considers it more source redundancy and recommends materialization, and vice versa. During implementation and preliminary experiments, we observed that factorization brings extra computation overhead. Thus, if the sizes of source instances and target instances are the same, our optimizer will recommend materialization.

Observation of source redundancy. Consider the union example in Table 1. Following the table unionability definition in [64, 59, 46], here we say the instances from two source tables S_1 and S_2 are union-compatible if: 1) S_1 and S_2 have the same number of feature columns; 2) each corresponding pair of feature columns have the same domain. Both S_1 and S_2 have the same set of columns $\{m, a, hr, o\}$ that also exist in the target table T , which are later used as features. Note that a source table may have non-feature

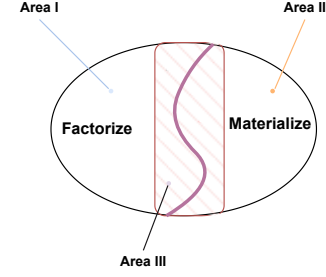


Figure 6: An abstraction of different decision areas

Table 2: Notations used in the paper

Symbol	Meaning
S_k	The k -th source table
T	Target table
I_k	The indicator matrix for S_k
M_k	The mapping matrix for S_k
j_T	The join type of T (inner/left/outer/union)
r_k/r_T	Number of rows in S_i/T
c_k/c_T	Number of columns in S_i/T
$m_k/m_T/m_X/m_w$	Number of nonzero elements in $S_i/T/X/w$

columns, e.g., $S_2.dd$, which will be dropped during preprocessing. Each matched attribute in the common feature set, e.g., m , have the same semantics and domain. With or without row overlap between S_1 and S_2 , the size of source instances is large or the same as target instances, i.e., source redundancy. The optimizer will recommend materialization.

Pruning rules. Union is a typical example for area II in Fig. 6. When source tables are unionable, materialization is more efficient. In the logical level, we can describe such relationships as tgds whose left-hand side (LHS) contains only predicates from a single source. In Example 4 $LHS(m_2) = S_1, LHS(m_3) = S_2$. We derived the following rule for favoring materialization, solely based on schema mappings.

PROPOSITION 3. *If an integration scenario can be expressed by one or a set of tgds whose LHS only contains predicates from a single source table, materialization is more efficient than factorization.*

The bigger picture. Tgds and equality generating dependencies (egds) are the two major types of database dependencies [28]. With tgds and egds we can formalize relationships within a database or across different databases of data sources, e.g., join dependencies, inclusion dependencies, functional dependencies, and schema mappings [55, 1, 31]. We can express functional dependencies in edges, and develop pruning rules for area II, e.g., identify target redundancies through functional dependencies. We demonstrate such a possibility with an example in [25]. We choose tgds, as they provide a high-level abstraction of dataset relationships, and are a convenient tool to analyze the algorithmic complexity [13, 61] and reasoning tasks [45]. It is a wide-open research field on logic-based formalism for model training and inference over data sources.

Implementation and limitations. The tgd-based pruning rules are implemented in Python, extensible to other query and programming languages, e.g., SQL, Java. Such schema-only rules can be checked quickly, but cannot cover all the cases, as the relative comparison between target redundancy and source redundancy also

Table 3: LA operator computations cost comparison

Operation	Materialization	Factorization
$T \odot x$		
$T^T \odot x$		
$f(T)$		
$f(T^T)$		
rowSum(T)	m_T	$\sum_{k=1}^K m_k$
rowSum(T^T)		
colSum(T)		
colSum(T^T)		
TX	$c_X \cdot m_T + r_T \cdot m_X$	$\sum_{k=1}^K c_X \cdot m_k + r_k \cdot m_X$
$T^T X$	$c_X \cdot m_T + c_T \cdot m_X$	$\sum_{k=1}^K c_X \cdot m_k + c_k \cdot m_X$
XT	$r_X \cdot m_T + c_T \cdot m_X$	$\sum_{k=1}^K r_X \cdot m_k + c_k \cdot m_X$
XT^T	$r_X \cdot m_T + r_T \cdot m_X$	$\sum_{k=1}^K r_X \cdot m_k + r_k \cdot m_X$

depends on the row matching, i.e., area III in Fig. 6. We propose a more general method next.

4.2 Complexity Analysis of LA Operations

We conduct a comparative analysis of materialization and factorization regarding the computational complexity of LA operations. For materialization, we consider the total number of computations involved in performing an LA operation on the materialized table. For factorization, we consider the number of operations involved in performing an LA operation on each source table. Below we continue to use LMM for the explanation.

Complexity analysis in state-of-the-art. In [14], a straightforward complexity analysis is provided. For each LA operator, the decision of factorization and materialization is based on the arithmetic computation complexity, i.e., the size comparison between joinable tables and materialized table. For instance, given a matrix T with the shape $r_T \times c_T$, and a matrix X with the shape $c_T \times c_X$ ($r_X = c_T$), the complexity of their left matrix multiplication TX is $c_X \cdot c_T \cdot r_T$. Thus, in [14] the complexity of computing TX for

materialization is $c_X \cdot c_T \cdot r_T$, while for factorization is $\sum_{k=1}^K c_X \cdot c_k \cdot r_k$.

Complexity analysis in Ilargi. Table 3 summarizes the cost comparison between materialization and factorization cases. At the logical level, the choice between factorization and materialization depends on the relative source and target redundancy, as discussed in Sec 2.4. However, when it comes to lower-level arithmetic operation computation, the sparsity of the matrices also matters. A *sparse matrix* is a matrix whose number of nonzero elements is negligible compared to the number of zeros, while a *dense matrix* has more non-zero elements. Given two matrices A and B , the naive matrix multiplication complexity is $O(r_A \cdot c_A \cdot c_B)$ as in [14]. In our approach, the matrices are stored in sparse format. Thus, the complexity of matrix multiplication AB is $O(c_B \cdot m_A + r_A \cdot m_B)$, where m_A , m_B are the numbers of nonzero elements in matrix A and B [39]. The matrix form of the source and target tables can be sparse or dense. As discussed in Sec. 3.1, our mapping and indicator matrices are sparse. Therefore, We compute the computation

cost as in Table 3. We also note the complexity of transposed operations. As described in Section 3.2, $T^T X$ is actually defined as $(X^T T)^T$ and XT^T is actually defined as $(TX^T)^T$. Here, we ignore the complexity added by the transposes, as these are not executed on data matrices. Notably, the differentiation between sparse and dense matrices is more about an intuition rather than a rigorous measurement [39]. Take the example of TX for materialization, and we have $m_T \leq r_T \cdot c_T$. When the target table matrix is dense, the non-zero element number m_T is close to the size of T , i.e., $r_T \cdot c_T$. The time for multiplying T and X is at most $O(r_T \cdot c_T \cdot c_X)$.

Computing the sparsity of T . Some elements in the complexity analysis can be inferred directly from the source tables, while others require estimation. Elements which can be directly inferred include r_X , c_X , r_k , c_k and m_k , as these are stored in SciPy objects. The elements r_T and c_T cannot be inferred from the source tables directly, but can be inferred from the indicator matrices and mapping matrices, respectively. The element that requires the most complicated computations is m_T , which depends upon the source tables' sparsities, the indicator matrix and the mapping matrix. In order to compute m_T , we compute T using the materialization computation $T \leftarrow I_1 S_1 M_1^T + \dots + I_K S_K M_K^T$. Then, we compute the sparsity directly from the result of this computation.

4.3 Complexity Analysis of ML Models

With the above complexity analysis of LA operations, we estimate the complexity of ML models. The complexity analysis on the model level is model-specific, depending on which linear algebras are involved in the model. The speedup brought by factorization is only applicable for a subset of LA operators [14]. Thus, here the discussion at the model level also focuses on such LA operators that bring speedups. We explain the complexity analysis of linear regression, with details of remaining models in the technical report [25].

Linear regression. First, we analyze the complexity of linear regression in Algorithm 1 (Sec. 2.1) in the materialized case, which is dominated by two matrix multiplication operations, i.e., $T^T(Tw)$. For the first matrix multiplication Tw : we denote the shape of weights vector w as $c_T \times 1$. We make the common assumption that w is dense, then $m_w = r_w \times c_w$. We denote the result of Tw as X , which is an intermediate result of linear regression algorithm. The size of X is $r_T \times 1$, i.e., $r_X = r_T$, $c_X = 1$. We assume X is also dense, thus, $m_X = r_X \times c_X$. Now we define the complexity of linear regression in the materialized case based on Tw and $T^T X$.

$$O_{\text{materialized}}(T) = \underbrace{c_w \cdot m_T + r_T \cdot m_w}_{Tw} + \underbrace{m_T + c_T \cdot m_X}_{T^T X}$$

Next, we define the complexity of the factorized case.

$$O_{\text{factorized}}(T) = \underbrace{\sum_{k=1}^K (c_w \cdot m_k + r_k \cdot m_w)}_{Tw} + \underbrace{\sum_{k=1}^K (m_k + c_k \cdot m_X)}_{T^T X}$$

Complexity ratio. We define a variable *complexity ratio* to indicate whether materialization or factorization leads to more redundancy.

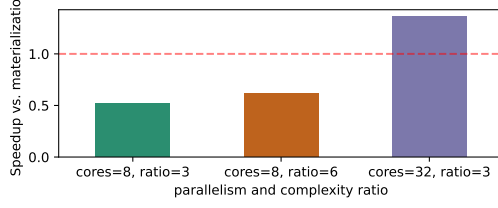


Figure 7: The speedups of training Linear Regression model with different the numbers of CPU cores

The complexity ratio is measured as the ratio of the materialization complexity divided by the factorization complexity.

$$ratio = \frac{O_{materialization}(T)}{O_{factorization}(T)}$$

4.4 ML-based Estimator

The third factor. The complexity ratio incorporates two factors, data including DI metadata, and ML algorithms. Our preliminary experiments revealed hardware configurations as a third vital factor. Fig. 7 showcases how the choice varies with changing parallelism w.r.t the numbers of CPU cores. Interestingly, factorization outperforms materialization when the parallelism increases to 32, even at a complexity ratio of 3. This result is attributed to factorization’s ability to decouple source table records, thus enabling parallel computational tasks and leveraging extensive hardware parallelism.

Design goal. Our goal is to create an estimator capable of making a binary decision—factorization or materialization—based on data characteristics, computational complexity, and hardware. Our approach, therefore, creates a more efficient machine learning training pipeline over data sources, optimizing the overall efficiency.

Ilargi’s tree boosting estimator. Among the various statistical models, tree boosting models are highly regarded in many cost estimators [16, 2] due to their explainability and prediction speed. Moreover, tree boosting has the capability to identify non-linear relationships among features. Given these advantages, we propose a tree-boosting estimator leveraging XGBoost [15]. The main challenge lies in identifying features that are relevant to the decision-making process between materialization and factorization.

Selection of hardware features. Our estimator is generalizable and portable, focusing on macro-level features instead of micro-level. Micro-level features (e.g., CPU cache speeds and memory latency), while valuable, can tie the model to specific hardware, limiting portability. Instead, we incorporate macro-level features like CPU cores and memory bandwidth. Shown in Sec. 5, combined with memory read/write costs, we can effectively estimate maximum memory access costs for both materialization and factorization.

Cost estimation pipeline. Fig. 8 illustrates the workflow of our proposed estimator with the example of linear regression.

1. Feature engineering. The estimator calculates the complexity ratio in Sec.4.3⁴, and theoretical memory read and write quantities for each operator within the input ML algorithm to be trained. Since factorization and materialization have different complexities and memory I/O amounts, we normalize these features by dividing by

⁴Hyperparameters such as the rank r for GNMF and the number of clusters k for KMeans are factored into the cost estimation. However, other hyperparameters like the learning rate γ for linear and logistic regression do not affect the cost estimation.

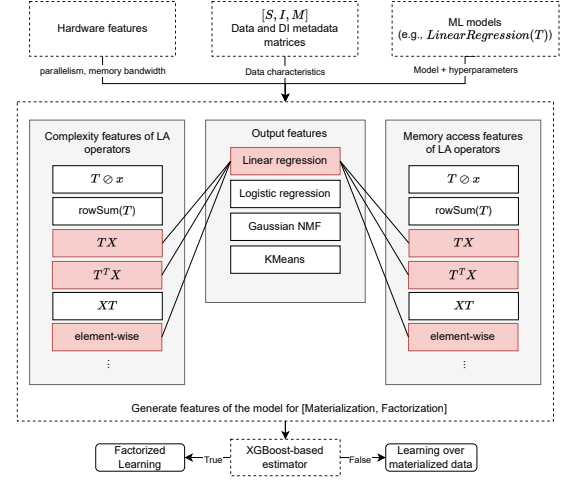


Figure 8: Workflow of the estimator.

total amounts, further dividing these normalized features by parallelism and memory bandwidth. We document and explain the used features in our technical report [25].

2. Train XGBoost & predict and recommend strategy. The key to a good performance of our estimator is creating a comprehensive training dataset for XGBoost, which spans a wide array of data characteristics and input models. We elaborate on our empirical approach and its effectiveness in Sec. 5.4.

Limitations and applicability. While the estimator offers benefits such as ease of construction and ability to capture non-linear features, it requires prior training. This requirement implies the need for a diverse dataset collection. Additionally, the features are sensitive to operator combinations in supported ML algorithms, reducing its predictive accuracy for models with unseen operators. Moreover, the estimator assumes the presence of operator information in the given ML algorithm supporting factorization, which potentially limits the types of models the estimator can support. Once the estimator is trained, it is used for binary classification, i.e., to determine whether to materialize or to factorize.

Answer to Q2

1. Schema mappings can be used as pruning rules.
2. The choice of factorization and materialization, depends on the datasets and their relationships, ML algorithms and their LA operators, and hardware configuration.

5 EVALUATION

5.1 Datasets and Setup

Synthetic datasets. To assess the accuracy and generalizability of the proposed estimator, we needed to collect runtime samples from a diverse range of integration scenarios, varying in redundancy, sparsity, and complexity. To create such scenarios, we have designed a synthetic data generator. The generator takes as input the dimensions of the source and target tables, target table sparsity, and the type of join. The output is a pair of source tables that, when joined, produce a target table with the specified characteristics. These tables are in matrices and then used to capture the runtime

Parameter	Values	Description
r_T	100k, 500k, 1M	Target table rows
c_T	10-50	Target table columns
j_T	left, inner, outer	Jointype
$sparsity$	0.0 - 0.9	Fraction of 0 values in T
Scalar complexity ratio	0.5 - 9.5	-
S_1 Column ratio	0.1 - 1.0	Fraction of columns from T in S_k
S_2 Column ratio	0.1 - 1.0	

Table 4: Parameters used for synthetic dataset generation.

of LA operations and ML models on both materialized and factorized data. Using the parameter ranges shown in Table 4, we created approximately 1800 datasets that allow us to accurately capture relations between data characteristics and operator speedup.

Real datasets. We use the seven real-world datasets from Morphheus and its extensions [50, 68, 14, 53], as provided by Project Hamlet⁵ [50, 68]. Each dataset consists of two or more tables connected in a star schema with PK-FK relationships. nnz indicates the number of nonzero elements. We include the characteristics of these datasets in the technical report [25].

Hardware & Software. We run our experiments with 16 cores and 48 cores of AMD EPYC 7H12 CPU. The system uses Ubuntu 20.04.4 LTS as the OS. Each experiment is run on one CPU. We use Python 3.10.4, SciPy 1.8.0 and NumPy 1.22.4. To enable parallelized sparse matrix multiplication, we choose MKL as the backend of NumPy. The number of iterations of model training is 20.

5.2 Speedup of Factorized LA Operators

In this section, we examine the speedups achieved by factorized LA operators, in comparison to their materialized counterparts.

Scalar operators. Fig. 9 illustrates the speedups with respect to varying complexity ratios and the join types. Speedups increase as the complexity ratio grows. Among the results, scalar multiplication operators achieve significantly higher speedups compared to rowSums and colSums. Notably, the rowSum and transposed colSums ($colSum(T^T)$) do not exhibit any speedup when compared to their materialized counterparts. In integration scenarios with different join types, factorization over inner joins delivers the most significant speedups compared to other types.

Fig. 10 presents speedups of factorization over different target table sizes and parallelism. As the size of the target table increases, all operators achieve higher speedups to varying degrees, with scalar multiplications benefiting the most. In contrast, speedups of factorization slightly decrease as the parallelism increases. One might expect that as the number of available cores for parallel computation increases, the speedup achieved through factorization would also increase. However, it is important to note that the decreased speedup might be due to the fact that increasing parallelism can also benefit materialized execution, thereby reducing the relative speedup of factorized execution.

Matrix operators. Fig. 11 displays the speedups of factorized matrix operators. LMM (TX) and transposed RMM (XT^T) exhibit noticeable speedups in all four configurations. In contrast, RMM (XT)

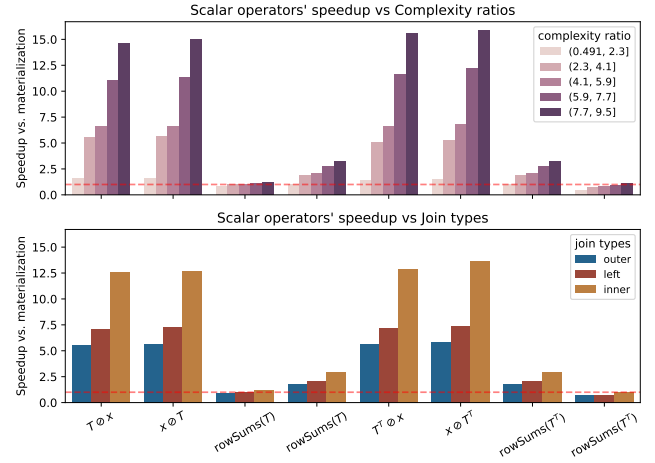


Figure 9: Speedups of factorized scalar operators w.r.t parallelism and complexity ratio. Element-wise operators deliver more significant speedups than reduction operators (e.g., rowSums).

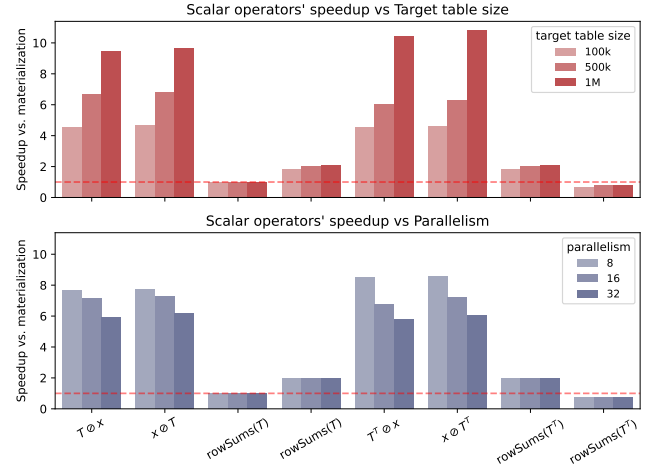


Figure 10: Speedups of factorized scalar operators w.r.t target table size and join type. Factorized scalar operators exhibit higher scalability on large datasets.

and transposed LMM ($T^T X$) are slower than materialization. All matrix operators deliver higher speedups with increasing parallelism and complexity ratios. Similar to the evaluation for scalar operators, matrix operators demonstrate better performance when applied to inner joins. However, in the evaluation that varies data sizes, the speedup of LMM declines as the target table size increases, indicating the LMM operator does not scale.

Analysis. The study reveals the impact of data characteristics on the speedup of factorization over materialization at the operator level. For instance, some operators (e.g., scalar multiplications, LMM) perform better with factorization, while others (e.g., rowSum, RMM) perform worse due to varying utilization of parallel resources and memory access efficiency.

Element-wise scalar operators over factorization show notable speedups because they can leverage high parallelism: each element's result does not depend on other elements. Interestingly, rowSum and colSum (reduction operators) show different speedups with factorization due to synchronization primitives needed for aggregations. In addition, given our implementation's row-major

⁵Project Hamlet datasets: <https://adalabucsd.github.io/hamlet.html>

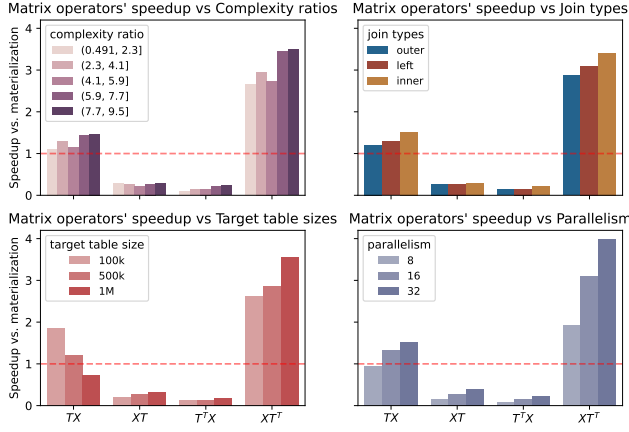


Figure 11: Speedups of factorized matrix operators w.r.t parallelism, complexity ratio, target size, and join type. LMM and transposed RMM show higher speedups with factorization, while RMM and transposed LMM with materialization.

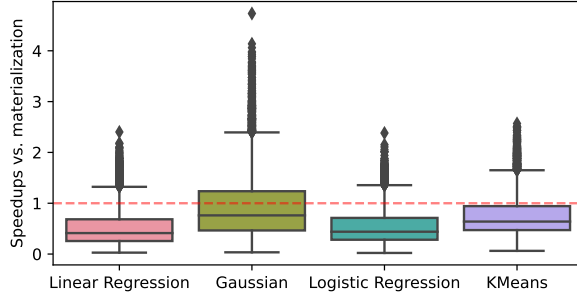


Figure 12: Speedups of factorization across all configurations. Although factorization often presents speedups over materialization, it's not consistently beneficial. It's crucial to understand this variable performance to apply the right strategy for different configurations.

memory layout, colSum's row-major access yields more efficient memory I/O than rowSum. Idem for LMM and transposed RMM.

Takeaways:

The speedups of LA operators using factorization fluctuate considerably depending on data characteristics, parallelism, and memory I/O efficiency.

5.3 Speedups of factorized ML-model training

The speedups of individual LA operators exhibit significant variation across different data and parallelism settings. Some scalar operators can achieve a speedup of over 10x, while certain matrix operators like RMM may underperform compared to their materialized counterparts. Consequently, ML model training, composed of combinations of these operators, may exhibit different speedups.

As shown in Fig. 12, speedups on various data configurations vary significantly, ranging from 0x to 5x, highlighting the need for an accurate estimator. Fig. 13 demonstrates that speedups for linear regression (other ML algorithms exhibited similar behavior) increase with parallelism and complexity ratio, supporting our theory in Sec. 4.3. That is, these factors are key for deciding between factorization and materialization.

Analysis. The speedup variability is attributed to the *mix* of operators in a model, with scalar operators and matrix multiplications

having a significant impact. For instance, the linear regression model, rich in highly accelerated scalar multiplications, shows a higher average speedup than GaussianNMF, which includes slower RMM operators.

Takeaways: The significant variability in speedup (ranging from 0.2x to 5x) depending on data and parallelism configurations calls for an accurate estimator. Thus, the estimator should consider key features such as data characteristics, operator complexities, and hardware.

5.4 Effectiveness of ML-based Estimator

We elaborate on training configurations of our tree boosting based estimator in Sec. 4.4. We compare our solution with baselines over synthetic data and validate our estimator using real-world data.

5.4.1 Configuration.

Training and test data. Based on the ML algorithm speedup evaluation, we extracted features and generated training data with labels (True for speedup when factorizing and False for the opposite). We further categorize the labels into 4 categories: True-Positive, True-Negative, False-Positive, and False-Negative. For instance, True-Positive prediction means that the estimator predicts factorization is faster correctly, while a False-Positive is the opposite. The data was split into training (80%) and test (20%) sets.

Evaluation. We now evaluate the accuracy and F1-score of our estimator model and compare it with other estimators using a synthetic test set. Accuracy reflects the model's ability to accurately decide between factorization and materialization.

Besides quality metrics, we assess estimators' effectiveness by comparing their overall speedup over materialization on the test set, focusing on algorithms trained with parallelism of 32. The observed speedups with the estimator approach, achieving a maximum speedup of 1.24x (as observed in our experiments), underscore the advantages of our ML-based estimator over threshold-based estimators in practical ML workloads, demonstrating its potential to enhance the overall performance of ML tasks.

Five baselines for comparison. We compared the performance of our proposed estimator in Sec. 4.4 with the following five baselines.

- i) To confirm the significance of hardware information, we train another tree boosting model using the same dataset, but *without hardware features*.
- ii) A second baseline is applying *logistic regression* instead of tree boosting. As a linear model, it may not effectively capture complex feature interactions like the tree-based model.



Figure 13: Speedups of factorized linear regression training with varying parallelism and complexity ratio. The decision boundary depends on a combination of complexity and hardware features.

Table 5: Accuracy, F1-score and Overall Speedups of estimators. Our tree boosting estimator shows superior predictive quality and effectiveness over synthetic test data

	Accuracy	F1-score	Overall Speedups
Tree Boosting (ours)	0.971	0.936	1.24
Tree Boosting w/o hardware	0.712	0.560	1.19
Logistic Regression	0.776	0.620	1.20
Complexity Ratio	0.447	0.420	1.16
tuple ratio + feature ratio	0.490	0.426	1.16
Selectivity	0.244	0.362	1.08

iii) The third baseline is using only the *complexity ratio* in Sec. 4.3. A complexity ratio of 1 marks the complexity borderline between materialization and factorization.

iv) The *tuple ratio* and *feature ratio*, metrics used in the state-of-the-art [14], quantify the target table's redundancy relative to the joinable tables. Morpheus [14] suggests a threshold of 5 for tuple ratio, and 1 for feature ratio.

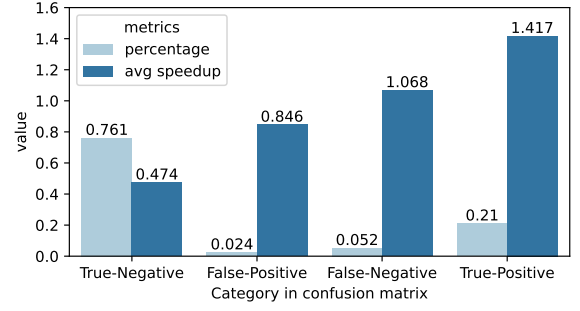
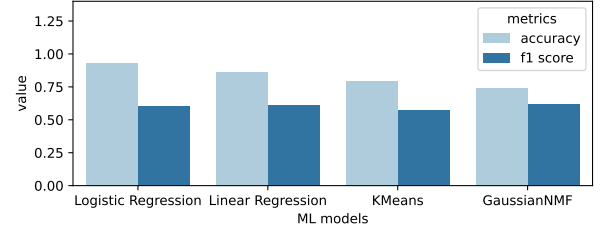
v) Finally, we consider the selectivity of joins, a traditional metric that measures the data proportion filtered out during join operations. Selectivity larger than 1 suggests that the result table has more data than the number of tuples of the largest source table, implying potential benefit using factorization.

5.4.2 Evaluation on synthetic data. Table 5 compares accuracy and F1-score of different estimators. The tree boosting model incorporating hardware features demonstrates 96.9% accuracy. Its high F1-score further emphasizes its balance between precision and recall, minimizing false predictions. In addition, the overall speedup with our estimator achieves 1.24, surpassing the baselines.

Performance of baselines. The tree boosting estimator without hardware features still achieves relatively high accuracy and overall speedups but presents a significant F1-score drop, verifying our hypothesis on the importance of hardware information in deciding between factorization or materialization. Finally, threshold-based estimators fall short in both accuracy and F1-score. The complexity ratio outperforms the tuple ratio suggested in earlier research across various threshold configurations. These results highlight our tree boosting model's effectiveness in determining whether to employ factorization or materialization.

Other metrics. Fig. 14 provides a confusion matrix for predictions from our tree boosting estimator, showing over 20% of test instances can benefit from factorized learning. We also analyze the average speedup for each confusion matrix category. Correctly predicted instances, on average, achieve a 1.42x speedup with factorization. Furthermore, when comparing the average speedups of False-Positive and True-Negative predictions, we find an interesting trend. Even when factorization is incorrectly applied due to mispredictions, the model training process can still achieve around 85% of the expected performance. This implies that our estimator is highly robust, maintaining reasonable performance even in the face of errors. Comparing the average speedups of True-Negative instances, our estimator shows a high positive correlation between accurate predictions and speedups.

Ablation study for generalizability. To evaluate the generalizability of our performance estimator, we utilized a leave-one-out cross-validation method, specific to different ML models to

**Figure 14: Confusion matrix of prediction result in test set. Over 20% instances gain speedup using factorized learning.****Figure 15: The predictive quality of our estimator under a "leave-one-out" scenario. Test set only contains the left-out model. For instance, the estimator for Logistic Regression has not been trained with any datapoints that we have gathered for logistic regression.**

be trained. In this approach, we systematically exclude each ML model from the training set, one at a time. The excluded model is then used as the test set to evaluate the predictive quality of the performance estimator trained without any information from that specific model. This way, we assess how well our estimator can generalize its predictions to unseen data related to the excluded model.

Fig. 15 indicates a decrease in both accuracy and F1-score when certain models are excluded from training data. This is due to missing operator features in the training set. For example, excluding the KMeans model, which heavily uses the dense scalar operator, creates a predictive bias due to an unbalanced feature distribution. Despite these constraints, the estimator remains functional with an approximate 80% accuracy rate. This suggests that even with certain predictive biases, our estimator is resilient and generalizable in optimizing ML workloads with new ML models, offering reasonably accurate predictions across different scenarios.

5.4.3 Feature Importance. Training the boosting tree model allows us to extract node information gain for ranking feature importance. Table 6 shows the top-10 features by information gain. Complexity features related to LMM, one of the most common operators in model training, have significantly high information gain. Operators like rowSums and colSums, which are not easily parallelized, are also important to determining speedup. Apart from complexity features, we observed a notable influence of memory I/O features, particularly those related to dense operators. Even though our raw data is stored in a sparse format, computations frequently generate dense intermediate results. This can significantly impact overall performance due to high I/O costs.

The complexity ratio, although not the most critical feature, is important in the estimator along with memory I/O features. Both materialization and factorization complexities are intrinsically

linked to the complexity ratio. Conversely, the tuple ratio ranks lower, implying a lesser impact on performance estimation.

Our feature importance analysis offers valuable insights for designing a performance estimator for factorized learning, emphasizing the importance of considering hardware features. Incorporating these in the estimation process can improve prediction accuracy, given the computational demands of factorized learning and the influence of hardware features on its efficiency.

5.4.4 Validation on real-world data. We further evaluated our estimator using real-world data to validate its practical usability in realistic ML workloads. The predictive quality of our estimator is presented in Table 7. For a fair comparison, we use the same datasets employed in Morpheus [14].

Despite certain differences in data characteristics (e.g., target table sizes, sparsity) between the real-world and synthetic data, our estimator continues to demonstrate remarkable predictive quality in comparison to other estimators. Even with a slight drop in predictive quality, our estimator still manages to achieve an overall speedup of 1.24 when applied to model training tasks using real-world data sets. This performance is close to the theoretical peak value of 1.27. This outcome indicates that dataset variations do not significantly impact the effectiveness of our estimator.

In addition, we investigate the accuracy for each dataset in Fig. 16. The accuracy values do not show large variations across different datasets. Through evaluations conducted on real-world datasets, our estimator has demonstrated consistent effectiveness, underscoring its robustness and reliability across a range of data characteristics. In summary, by leveraging our performance estimator to make informed decisions on whether to employ factorization or materialization, we can enhance the overall speedup of ML training workloads across data sources.

Table 6: Top 10 information gain of features in the estimator.

Top 10 features	Info. Gain
rowSums mem. write	30.96
LMM complexity with factorization	22.69
Dense scalar operation complexity	20.53
LMM complexity with materialization	8.68
colSums complexity with materialization	6.30
Complexity ratio	6.21
Total mem. write with materialization	5.88
Dense MM operation complexity	5.63
Dense MM mem. write	5.77
Feature Ratio	5.45

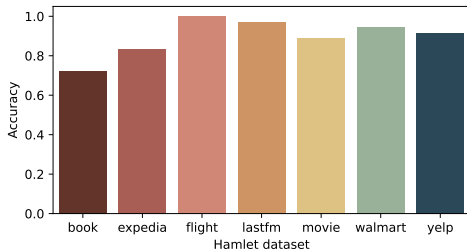


Figure 16: Accuracy w.r.t different real-world datasets.

Table 7: Accuracy, F1-score and Overall Speedup of estimators over real-world data. Our proposed estimator maintains a performance similar to the evaluation conducted with the test set, while other estimators demonstrate varying degrees of performance degradation.

	Accuracy	F1-score	Overall Speedups
Tree Boosting (ours)	0.905	0.821	1.22
Logistic Regression	0.651	0.45	1.11
Complexity Ratio	0.409	0.387	1.09
Tuple Ratio	0.464	0.211	0.94
Selectivity	0.591	0.237	0.97

Takeaways:

- i) Our tree-boosting performance estimator outperforms other estimators in terms of accuracy, F1-score, and overall speedup on both synthetic and real-world datasets, thereby enhancing ML training workloads over source tables.
- ii) Our estimator exhibits robust predictive power on unseen ML algorithms and datasets.

6 RELATED WORK

Factorized learning. Early efforts required *manual* design for factorizing specific ML algorithms. Orion [48] supports linear and logistic regression, which are later extended to decision trees, feature ranking, and Naive Bayes in Santoku [49]. These two works mainly tackle inner joins via PK-FK relationship, while F [67] has extended factorization of linear regression to natural joins. As the successors of F, AC/DC [44] added support to categorical variables, and LMFAO supports more ML algorithms such as decision trees. In a more *automated* manner, Morpheus [14] proposed a general factorization framework enabling the automation of linear algebra rewriting. By reordering of multiplication and exploiting pre-computed results, HADAD [3] creates rewrites for systems like Morpheus and further speeds up the computation. Later, Trinity [42] extends Morpheus in such a way that Factorized LA logic is written *once*, but can be reused in multiple programming languages and LA tools. MorpheusFI [53] added support to non-linear feature interactions. Non-linear models such as Gaussian Mixture Models and Neural Networks are supported in [17].

Ilargi's main contribution is not on extending the ML algorithm. We expand the dataset relationships to various integration scenarios, cf. Table 1. We take a more critical inspection of the speedup of factorization compared to materialization by utilizing DI metadata. We also discovered a third factor, hardware configuration, which has an impact on the decision boundary. We are currently expanding the implementation of Ilargi to more ML models, e.g., decision trees.

Tensor-based ML and data management systems. A large body of work has studied the unification of linear and relational algebra operations at the level of the execution engine [56, 9]. SystemDS [10] models data as tensors, and performs optimizations for end-to-end ML pipelines. TQP [38] executes relational operators such as joins and projections in tensor algebra, while multiple works leverage GPUs for relational joins [4, 21, 40]. Finally, Hummingbird [58, 47] compiles traditional ML models into tensor-based runtimes designed for deep learning models. Ilargi is complementary to these approaches as its input source datasets also come in matrix form. The key novelty of Ilargi lays in the modeling of data integration

metadata (e.g., schema mappings traditionally expressed in first-order logic) in a matrix representation.

7 CONCLUSION

In this paper, we address the challenge of training machine learning models over disparate data sources. We introduce a formal framework based on tgds that links existing factorization approaches based on inclusion dependencies and join dependencies with data integration scenarios specified by schema mappings. Ilargi leverages matrices to encode schema mapping and row matching, enabling factorization over sources. Moreover, Ilargi utilizes schema mappings as pruning rules, and employs a tree boosting based estimator, which leverages complexity ratios and hardware configuration to choose between materialization and factorization. Our experiments demonstrate the effectiveness and efficiency of our estimator, which covers more scenarios than state-of-the-art solutions and performs correct predictions with 90–96% accuracy depending on the dataset.

REFERENCES

- [1] Serge Abiteboul, Richard Hull, and Victor Vianu. *Foundations of databases*. Vol. 8. Addison-Wesley Reading, 1995.
- [2] Andrew Adams et al. "Learning to optimize halide with tree search and random programs". In: *ACM Trans. Graph.* 38.4 (2019), 121:1–121:12.
- [3] Rana Alotaibi et al. "HADAD: A Lightweight Approach for Optimizing Hybrid Complex Analytics Queries". In: *Proceedings of the 2021 International Conference on Management of Data*. SIGMOD/PODS '21: International Conference on Management of Data. Virtual Event China: ACM, June 9, 2021, pp. 23–35.
- [4] Rasmus Resen Amossen and Rasmus Pagh. "Faster Join-Projects and Sparse Matrix Multiplications". In: *ICDT 2009*. St. Petersburg, Russia: Association for Computing Machinery, 2009, 121–126.
- [5] Patricia C Arocena et al. "The iBench integration metadata generator". In: *Proceedings of the VLDB Endowment* 9.3 (2015), pp. 108–119.
- [6] David Aumüller et al. "Schema and ontology matching with COMA++". In: *Proceedings of the 2005 ACM SIGMOD international conference on Management of data*. 2005, pp. 906–908.
- [7] C. Beeri and M. Y. Vardi. "A proof procedure for data dependencies". In: *JACM* 31.4 (1984), pp. 718–741.
- [8] Felix Biessmann et al. "DataWig: Missing Value Imputation for Tables." In: *J. Mach. Learn. Res.* 20.175 (2019), pp. 1–6.
- [9] Matthias Boehm, Matteo Interlandi, and Chris Jermaine. "Optimizing Tensor Computations: From Applications to Compilation and Runtime Techniques". In: *Companion of the 2023 International Conference on Management of Data*. SIGMOD '23. Seattle, WA, USA: Association for Computing Machinery, 2023, 53–59.
- [10] Matthias Boehm et al. "SystemDS: A Declarative Machine Learning System for the End-to-End Data Science Lifecycle". In: (2020).
- [11] Jacqueline M Bos et al. "Prediction of clinically relevant adverse drug events in surgical patients". In: *PloS one* 13.8 (2018), e0201645.
- [12] David Guy Brizan and Abdullah Uz Tansel. "A. survey of entity resolution and record linkage methodologies". In: *Communications of the IIMA* 6.3 (2006), p. 5.
- [13] Balder ten Cate and Phokion G. Kolaitis. "Structural characterizations of schema-mapping languages". In: *Database Theory - ICDT 2009, 12th International Conference, St. Petersburg, Russia, March 23–25, 2009, Proceedings*. Ed. by Ronald Fagin. Vol. 361. ACM International Conference Proceeding Series. ACM, 2009, pp. 63–72.
- [14] L. Chen et al. "Towards Linear Algebra over Normalized Data". In: *PVLDB* 10.11 (2017).
- [15] Tianqi Chen and Carlos Guestrin. "XGBoost: A Scalable Tree Boosting System". In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13–17, 2016*. Ed. by Balaji Krishnapuram et al. ACM, 2016, pp. 785–794.
- [16] Tianqi Chen et al. "TVM: An Automated End-to-End Optimizing Compiler for Deep Learning". In: *13th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2018, Carlsbad, CA, USA, October 8–10, 2018*. Ed. by Andrea C. Arpaci-Dusseau and Geoff Voelker. USENIX Association, 2018, pp. 578–594.
- [17] Zhaoyue Cheng et al. "Efficient Construction of Nonlinear Models over Normalized Data". In: *37th IEEE International Conference on Data Engineering, ICDE 2021, Chania, Greece, April 19–22, 2021*. IEEE, 2021, pp. 1140–1151.
- [18] N. Chepurko et al. "ARDA: automatic relational data augmentation for machine learning". In: *VLDB* 13.9 (2020), pp. 1373–1387.
- [19] Rada Chirkova, Jun Yang, et al. "Materialized views". In: *Foundations and Trends® in Databases* 4.4 (2012), pp. 295–405.
- [20] Thomas H Cormen et al. *Introduction to algorithms*. MIT press, 2022.
- [21] Shaleen Deep, Xiao Hu, and Paraschos Koutis. "Fast Join Project Query Evaluation Using Matrix Multiplication". In: *SIGMOD* 2020. 2020, 1213–1223.
- [22] Marc Peter Deisenroth, A Aldo Faisal, and Cheng Soon Ong. *Mathematics for machine learning*. Cambridge University Press, 2020.
- [23] Alin Deutsch, Lucian Popa, and Val Tannen. "Query reformulation with constraints". In: *ACM SIGMOD Record* 35.1 (2006), pp. 65–73.
- [24] A. Doan, A. Halevy, and Z. Ives. *Principles of data integration*. Elsevier, 2012.
- [25] *Efficient ML Model Training over Silos: to Factorize or to Materialize*. Technical report. https://github.com/amademicnobody12/sigmod24-310/blob/main/technical_report.pdf. 2023.
- [26] Mahdi Esmailoghli, Jorge-Arnulfo Quiané-Ruiz, and Ziawach Abedjan. "CO-COA: COrelation COefficient-Aware Data Augmentation." In: *EDBT*. 2021, pp. 331–336.
- [27] Jérôme Euzenat, Pavel Shvaiko, et al. *Ontology matching*. Vol. 18. Springer, 2007.
- [28] R. Fagin. "Tuple-Generating Dependencies". In: *Encyclopedia of Database Systems*. Ed. by LING LIU and M. TAMER ÖZSU. Boston, MA: Springer US, 2009, pp. 3201–3202.
- [29] Ronald Fagin and M Vardi. "The theory of data dependencies". In: *Mathematics of Information Processing* 34 (1986), p. 19.
- [30] Ronald Fagin et al. "Clio: Schema mapping creation and data exchange". In: *ER*. Springer, 2009, pp. 198–236.
- [31] Ronald Fagin et al. "Data exchange: Semantics and query answering". In: *Database Theory—ICDT 2003: 9th International Conference Siena, Italy, January 8–10, 2003 Proceedings* 9. Springer, 2003, pp. 207–224.
- [32] Grace Fan et al. "Table Discovery in Data Lakes: State-of-the-Art and Future Directions". In: *Companion of the 2023 International Conference on Management of Data*. SIGMOD '23. Seattle, WA, USA: Association for Computing Machinery, 2023, 69–75.
- [33] Behzad Golshan et al. "Data Integration: After the Teenage Years". In: *Proceedings of the 36th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2017, Chicago, IL, USA, 2017*, pp. 101–106.
- [34] Isabelle Guyon et al., eds. *Feature Extraction - Foundations and Applications*. Vol. 207. Studies in Fuzziness and Soft Computing. Springer, 2006.
- [35] Rihan Hai et al. "Data Lakes: A Survey of Functions and Systems". In: *IEEE Transactions on Knowledge and Data Engineering* (2023), pp. 1–20.
- [36] Alon Halevy, Anand Rajaraman, and Joann Ordille. "Data Integration: The Teenage Years". In: *Proceedings of the 32nd International Conference on Very Large Data Bases*. VLDB '06. Seoul, Korea: VLDB Endowment, 2006, 9–16.
- [37] Andrew Hard et al. "Federated learning for mobile keyboard prediction". In: *arXiv preprint arXiv:1811.03604* (2018).
- [38] Dong He et al. "Query Processing on Tensor Computation Runtimes". In: vol. 15. VLDB Endowment, 2022, 2811–2825.
- [39] Ellis Horowitz and Sartaj Sahni. "Fundamentals of data structures". In: (1982).
- [40] Zichun Huang and Shimin Chen. "Density-Optimized Intersection-Free Mapping and Matrix Multiplication for Join-Project Operations". In: vol. 15. VLDB Endowment, 2022, 2244–2256.
- [41] Informatica. <https://www.informatica.com/es/products/data-integration/powercenter.html>.
- [42] David Justo et al. "Towards a Polyglot Framework for Factorized ML". In: *Proceedings of the VLDB Endowment* 14.12 (July 1, 2021), pp. 2918–2931.
- [43] John D Kelleher, Brian Mac Namee, and Aoife D'arcy. *Fundamentals of machine learning for predictive data analytics: algorithms, worked examples, and case studies*. MIT press, 2020.
- [44] Mahmoud Abo Khamis et al. "AC/DC: In-Database Learning Thunderstruck". In: *Proceedings of the Second Workshop on Data Management for End-To-End Machine Learning*. SIGMOD/PODS '18: International Conference on Management of Data. Houston TX USA: ACM, June 15, 2018, pp. 1–10.
- [45] Phokion G. Kolaitis et al. "Nested dependencies: structure and reasoning". In: *Proceedings of the 33rd ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS'14, Snowbird, UT, USA, June 22–27, 2014*. Ed. by Richard Hull and Martin Grohe. ACM, 2014, pp. 176–187.
- [46] Christos Koutras et al. "Valentine: Evaluating Matching Techniques for Dataset Discovery". In: *2021 IEEE 37th International Conference on Data Engineering (ICDE)*. 2021, pp. 468–479.
- [47] Dimitrios Koutsoukos et al. "Tensors: An Abstraction for General Data Processing". In: *Proc. VLDB Endow.* 14.10 (2021), 1797–1804.
- [48] A. Kumar, J. Naughton, and J. M. Patel. "Learning generalized linear models over normalized data". In: *SIGMOD*. 2015, pp. 1969–1984.
- [49] Arun Kumar et al. "Demonstration of Santoku: optimizing machine learning over normalized data". In: *Proceedings of the VLDB Endowment* 8.12 (Aug. 2015), pp. 1864–1867.

- [50] Arun Kumar et al. "To join or not to join? thinking twice about joins before feature selection". In: *Proceedings of the 2016 International Conference on Management of Data*. 2016, pp. 19–34.
- [51] Maurizio Lenzerini. "Data integration: A theoretical perspective". In: *PODS*. ACM. 2002, pp. 233–246.
- [52] Peng Li et al. "CleanML: A Benchmark for Joint Data Cleaning and Machine Learning [Experiments and Analysis]". In: *CoRR abs/1904.09483* (2019). arXiv: 1904.09483.
- [53] Side Li, Lingjiao Chen, and Arun Kumar. "Enabling and Optimizing Non-linear Feature Interactions in Factorized Linear Algebra". In: *Proceedings of the 2019 International Conference on Management of Data*. SIGMOD/PODS '19: International Conference on Management of Data. Amsterdam Netherlands: ACM, June 25, 2019, pp. 1571–1588.
- [54] Jiabin Liu et al. "Feature augmentation with reinforcement learning". In: *2022 IEEE 38th International Conference on Data Engineering (ICDE)*. IEEE. 2022, pp. 3360–3372.
- [55] David Maier. *The Theory of Relational Databases*. Computer Science Press, 1983.
- [56] Nantia Makrynioti and Vasilis Vassalos. "Declarative Data Analytics: A Survey". In: *IEEE Transactions on Knowledge and Data Engineering* 33.6 (June 2021), pp. 2392–2411.
- [57] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- [58] Supun Nakandala et al. "A Tensor Compiler for Unified Machine Learning Prediction Serving". In: *OSDI 2020*. USA: USENIX Association, 2020.
- [59] Fatemeh Nargesian et al. "Table union search on open data". In: *Proceedings of the VLDB Endowment* 11.7 (2018), pp. 813–825.
- [60] Pentaho. <https://www.hitachivantara.com/en-us/products/pentaho-platform/data-integration-analytics.html>.
- [61] Reinhard Pichler and Sebastian Skritek. "The Complexity of Evaluating Tuple Generating Dependencies". In: *Proceedings of the 14th International Conference on Database Theory*. ICDT '11. Uppsala, Sweden: Association for Computing Machinery, 2011, 244–255.
- [62] Neoklis Polyzotis et al. "Data Lifecycle Challenges in Production Machine Learning: A Survey". In: *SIGMOD Rec.* 47.2 (2018), pp. 17–28.
- [63] Erhard Rahm and Philip A Bernstein. "A survey of approaches to automatic schema matching". In: *the VLDB Journal* 10.4 (2001), pp. 334–350.
- [64] Raghu Ramakrishnan, Johannes Gehrke, and Johannes Gehrke. *Database management systems*. Vol. 3. McGraw-Hill New York, 2003.
- [65] Theodoros Rekatsinas et al. "HoloClean: Holistic Data Repairs with Probabilistic Inference". In: *Proceedings of the VLDB Endowment* 10.11 (2017).
- [66] Yuji Roh, Geon Heo, and Steven Euijong Whang. "A Survey on Data Collection for Machine Learning: A Big Data - AI Integration Perspective". In: *IEEE Transactions on Knowledge and Data Engineering* 33.4 (2021), pp. 1328–1347.
- [67] Maximilian Schleich, Dan Olteanu, and Radu Ciucanu. "Learning Linear Regression Models over Factorized Joins". In: *Proceedings of the 2016 International Conference on Management of Data*. SIGMOD '16. New York, NY, USA: Association for Computing Machinery, June 14, 2016, pp. 3–18.
- [68] Vraj Shah, Arun Kumar, and Xiaojin Zhu. "Are key-foreign key joins safe to avoid when learning high-capacity classifiers?" In: *Proceedings of the VLDB Endowment* 11.3 (2017), pp. 366–379.
- [69] SSIS. <https://learn.microsoft.com/en-us/sql/integration-services/sql-server-integration-services?view=sql-server-ver16>.
- [70] Michael Stonebraker and Ihab F. Ilyas. "Data Integration: The Current Status and the Way Forward". In: *IEEE Data Eng. Bull.* 41.2 (2018), pp. 3–9.
- [71] Talend. <https://www.talend.com/products/integrate-data/>.
- [72] Lei Yu and Huan Liu. "Efficient feature selection via analysis of relevance and redundancy". In: *The Journal of Machine Learning Research* 5 (2004), pp. 1205–1224.