# Multi-layer Optimizations for End-to-End Data Analytics

**Amir Shaikhha**
University of Oxford
United Kingdom

**Maximilian Schleich**
University of Oxford
United Kingdom

**Alexandru Ghita**
University of Oxford
United Kingdom

**Dan Olteanu**
University of Oxford
United Kingdom

## Abstract

We consider the problem of training machine learning models over multi-relational data. The mainstream approach is to first construct the training dataset using a feature extraction query over input database and then use a statistical software package of choice to train the model. In this paper we introduce Iterative Functional Aggregate Queries (IFAQ), a framework that realizes an alternative approach. IFAQ treats the feature extraction query and the learning task as one program given in the IFAQ's domain-specific language, which captures a subset of Python commonly used in Jupyter notebooks for rapid prototyping of machine learning applications. The program is subject to several layers of IFAQ optimizations, such as algebraic transformations, loop transformations, schema specialization, data layout optimizations, and finally compilation into efficient low-level C++ code specialized for the given workload and data.

We show that a Scala implementation of IFAQ can outperform mlpack, Scikit, and TensorFlow by several orders of magnitude for linear regression and regression tree models over several relational datasets.

***CCS Concepts*** • **Computing methodologies** → *Supervised learning by regression*; • **Information systems** → *Database management system engines*; • **Software and its engineering** → *Domain specific languages*.

***Keywords*** In-Database Machine Learning, Multi-Query Optimization, Query Compilation.

## 1 Introduction

The mainstream approach in supervised machine learning over relational data is to specify the construction of the data matrix followed by the learning task in scripting languages such as Python, MATLAB, Julia, or R using software environments such as Jupyter notebook. These environments call libraries for query processing, e.g., Pandas [34] or Spark-SQL [4], or database systems, e.g., PostegreSQL [37]. The materialized training dataset then becomes the input to a statistical package, e.g., mlpack [14], scikit-learn [41], TensorFlow [1], and PyTorch [40]) that learns the model. There

are clear advantages to this approach: it is easy to use, allows for quick prototyping, and does not require substantial programming knowledge. Although it uses libraries that provide efficient implementations for their functions (usually by binding to efficient low-level C code), they miss optimization opportunities for the end-to-end relational learning pipeline. In particular, they fail to exploit the relational structure of the underlying data, which was removed by the materialization of the training dataset. The effect is that the efficiency on one machine is often severely limited, with common deployments having to rely on expensive and energy-consuming clusters of machines to perform the learning task.

The thesis of this paper is that this performance limitation can be overcome by systematically optimizing the end-to-end relational learning pipeline as a whole. We introduce Iterative Functional Aggregate Queries, or IFAQ for short, a framework that is designed to uniformly process and automatically optimize the various phases of the relational learning pipeline. IFAQ takes as input a program that fully specifies both the data matrix construction and the model training in a dynamically-typed language called D-IFAQ. This language captures a fragment of scripting languages such as Python that is used for rapid prototyping of machine learning models. This is possible thanks to the iterative and collection-oriented constructs provided by D-IFAQ.

An IFAQ program is optimized by multiple compilation layers, whose optimizations are inspired by techniques developed by the data management, programming languages, and high-performance computing communities. IFAQ performs automatic memoization, which identifies code fragments that can be expressed as batches of aggregate queries. This optimization further enables non-trivial loop-invariant code motion opportunities. The program is then compiled down to a statically-typed language, called S-IFAQ, which benefits from further optimization stages, including loop transformations such as loop fusion and code motion. IFAQ investigates optimization opportunities at the interface between the data matrix construction and learning steps and interleaves the code for both steps by pushing the data-intensive computation from the second step past the joins of the first step, inspired by recent query processing techniques [2, 38]. The
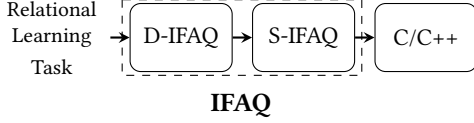
**Figure 1.** A relational learning task is expressed in D-IFAQ, transformed into an optimized S-IFAQ expression, and compiled to efficient C++ code.

outcome is a highly optimized code with no separation between query processing and machine learning. Finally, IFAQ compiles the optimized program into low-level C++ code that further exploits data-layout optimizations. IFAQ can outperform equivalent solutions by orders of magnitude.

The contributions of this paper are as follows:

- Section 2 introduces the IFAQ framework, which comprises languages and compiler optimizations that are designed for efficient end-to-end relational learning pipelines.
- As proof of concept, Section 3 demonstrates IFAQ for two popular models: linear regression and regression tree.
- Section 4 shows how to systematically optimize an IFAQ program in several stages.
- Section 5 benchmarks IFAQ, mlpack, TensorFlow, and scikit-learn. It shows that IFAQ can train linear regression and regression tree models over two real datasets orders-of-magnitude faster than its competitors. In particular, IFAQ learns the models faster than it takes the competitors to materialize the training dataset. Section 5 further shows the performance impact of individual optimization layers.

## 2 Overview

Figure 1 depicts the IFAQ workflow. Data scientists use software environments such as Jupyter Notebook to specify relational learning programs. In our setting, such programs are written in a dynamically-typed language called D-IFAQ and subject to high-level optimizations (Section 4.1). Given the schema information of the dataset, the optimized program is translated into a statically-typed language called S-IFAQ (Section 4.2). If there are type errors, they are reported to the user. IFAQ performs several optimizations inspired by database query processing techniques (Section 4.3). Finally, it synthesizes appropriate data layouts, resulting in efficient low-level C/C++ code (Section 4.4).

### 2.1 IFAQ Core Language

The grammar of the IFAQ core language is given in Figure 2. This functional language support the following data types. The first category consists of numeric types as well as categorical types (e.g., boolean values, string values, and other custom enum types). Furthermore, for categorical types the other alternative is to one-hot encode them, the type of which is represented as $\mathbb{R}^n_{T_1}$. This type represents an array of $n$ real numbers each one corresponding to a particular value in the domain of $T_1$. In the one-hot encoding of a value of type

$$
\begin{array}{rcl}
p & ::= & e \mid \mathbf{x} \leftarrow e \; \texttt{while(e)} \; \{ \; \mathbf{x} \leftarrow e \} \; \mathbf{x} \\[4pt]
e & ::= & e + e \mid e * e \mid -e \mid uop(e) \mid e \; bop \; e \mid c \\[4pt]
& \mid & \displaystyle\sum_{x \in e} e \mid \mathop{\lambda}_{x \in e} e \mid \{\{\overrightarrow{e \to e}\}\} \mid [\![\overrightarrow{e}]\!] \mid \texttt{dom}(e) \mid e(e) \\[10pt]
& \mid & \{\overrightarrow{x = e}\} \mid \langle x = e \rangle \mid e.x \mid e[e] \\[4pt]
& \mid & \texttt{let} \; x = e \; \texttt{in} \; e \mid x \mid \texttt{if} \; e \; \texttt{then} \; e \; \texttt{else} \; e \\[4pt]
c & ::= & `id` \mid "id" \mid n \mid r \mid \texttt{false} \mid \texttt{true} \\[4pt]
T & ::= & S \mid \{\overrightarrow{x : T}\} \mid \langle \overrightarrow{x : T} \rangle \mid \texttt{Map[}T, \; T\texttt{]} \mid \texttt{Set[}T\texttt{]} \\[4pt]
S & ::= & B \mid C \\[4pt]
C & ::= & \texttt{string} \mid \texttt{enum} \mid \texttt{bool} \mid \texttt{Field} \\[4pt]
B & ::= & \mathbb{Z} \mid \mathbb{R} \mid \mathbb{R}^n_C
\end{array}
$$

**Figure 2.** Grammar of IFAQ core language.

$T_1$, only the $i^{th}$ value is 1 and the rest are zero. However, a value of type denoting the one-hot encoding of $\mathbb{R}^n_{T_1}$, can take arbitrary real numbers at each position.

The second category consists of record types, which are similar to *struct*s in C; the record values contain various fields of possibly different data types. The partial records, where some fields have no values are referred to as *variants*.

The final category consists of collection data types such as (ordered) sets and dictionaries. Database relations are represented as dictionaries (mapping elements to their multiplicities). In D-IFAQ, the elements of collections can have different types. However, in S-IFAQ the elements of collections should have the same data type. In order to distinguish the variables with collection data type with other types of variables, we denote them as **x** and $x$, respectively.

The top-level program consists of several initialization expressions, followed by a loop. This enables IFAQ to express iterative algorithms such as gradient descent optimization algorithms. The rest of the constructs of this language are as follows: (i) various binary and unary operations on both scalar and collection data structures,[1] (ii) let binding for assigning the value of an expression to a variable, (iii) conditional expressions (iv) $\sum_{x \in e_1} e_2$: the summation operator in order to iterate over the elements of a collection ($e_1$) and perform a stateful computation using a particular addition operator (e.g., numerical addition, set union, bag union),[2] (v) $\lambda_{x \in e_1} e_2$: constructing a dictionary where the key domain is $e_1$ and the value for each key $x$ is constructed using $e_2$, (vi) constructing sets ($[\![e]\!]$) and dictionaries ($\{\{e \to e\}\}$) given a list of elements, (vii) dom($e$): getting the domain of a dictionary, (viii) $e_0(e_1)$: retrieving the associatiated value to the given key in a dictionary, (ix) constructing records ($\{\overrightarrow{x = e}\}$) and variants ($\langle x = e \rangle$), (x) statically (e.f) and dynamically (e[f]) assessing the field of a record or a variant.

---

[1] More specifically ring-based operations.

[2] This operator could be every monoid-based operator. Thus, even computing the minimum of two numbers is an addition operator.

# 3 Applications

As proof of concept, we show in this section how linear regression and regression tree models can be expressed in D-IFAQ. We assume that the model is learned over a training dataset $\mathbf{Q}$ with features $a_1, \ldots, a_n$ and label $a_{n+1}$, which is the output of a query over a database with relations $R_1, \ldots, R_m$.

Further applications that can benefit from IFAQ include supervised learning problems (e.g., logistic regression, factorization machines, and support vector machines), as well as unsupervised learning problems (e.g., k-means clustering, principal component analysis, and learning Bayesian networks via Chow-Liu trees). IFAQ further supports different optimization algorithms, e.g., batch or stochastic gradient descent, coordinate descent, and alternating minimization.

**Linear Regression:** A linear regression model is given by:

$$LR(x) = \sum_{f \in \mathbf{F}} \theta(f) * x[f]$$

where $x$ is a record with fields $\mathbf{F} = \{a_1, \ldots, a_n\}$, and $\theta$ is a dictionary with these fields as its keys. The $\theta$ dictionary defines the parameters of the model, and the output of $LR(x)$ is a prediction for the label given input $x$.

Linear regression models require categorical features to be one-hot encoded. Thus, for a categorical feature $a_i$, the field $x[a_i]$ encodes a vector of one indicator value for each value in the domain of $a_i$ (Section 2.1). The corresponding parameter element $\theta(a_i)$ encode a vector of parameters which is of the same dimension as $x[a_i]$. We assume without loss of generality that $x[a_1]$ only takes value 1 and then $\theta(a_1)$ is the intercept of the model.

Given the training dataset $\mathbf{Q}$, the error of fitting $LR$ to $\mathbf{Q}$ is given by the sum of the least squares loss function:[3]

$$J(\theta) = \frac{1}{2|\mathbf{Q}|} \sum_{x \in \text{dom}(\mathbf{Q})} \Big( \sum_{f \in \mathbf{F}} \theta(f) * x[f] - x[a_{n+1}] \Big)^2$$

We minimize $J(\theta)$ using batch gradient descent (BGD) optimization, which repeatedly updates each parameter $\theta(f)$ by learning rate $\alpha$ in the direction of the partial derivative of $J(\theta)$ w.r.t. $\theta(f)$ until convergence. This is represented in the following D-IFAQ program ($\mathbf{Q}$ is expressed in IFAQ as a query over a multi-relational database):

```
let F = [[a₁, a₂, …, aₙ]] in
```
$$\theta \leftarrow \theta_0$$
```
while( not converged ) {
```
$$\theta = \lambda_{f_1 \in \mathbf{F}} \Big( \theta(f_1) -$$
$$\frac{\alpha}{|\mathbf{Q}|} \sum_{x \in \text{dom}(\mathbf{Q})} \mathbf{Q}(x) * \big( \sum_{f_2 \in \mathbf{F}} \theta(f_2) * x[f_2] - x[a_{n+1}] \big) * x[f_1] \Big)$$
```
}
```
$$\theta$$

---

[3]It is common to penalize $J(\theta)$ with a regularization term. IFAQ supports regularization as well, but we omit it from the examples to avoid clutter.

This program is inefficient, because it computes for each BGD iteration all pairwise products $x[f_1] * x[f_2]$. These products, however, are invariant, and can be computed once outside the while loop. This is the main high-level optimization performed by IFAQ for this case (Section 4.1). The rewriting would gather a collection of these products, which can be computed once as database aggregate queries. The following BGD iterations can then be computed directly over these aggregates, and do not require a pass over $\mathbf{Q}$. The collection of these aggregates represents the *non-centered covariance matrix* (or *covar matrix* for short). The aggregates can then be further optimized using ideas from the database query processing, in particular they can be computed directly over the input database and without materializing $\mathbf{Q}$ (Section 4.3).

We next provide a running example which is used to explain the optimization rewritings in the following sections.

**Example 3.1.** In retail forecasting scenarios, the goal is to learn a model which can predict future sales for items at different stores. We assume the retailer has a database with three relations: **S**ales(*item, store, units*), Sto**R**es(*store, city*), **I**tems(*item, price*). The goal is to learn a model that predicts $u$ with the features $\mathbf{F} = \{i, s, c, p\}$. The training dataset is given by the result of the join of the three relations: $\mathbf{Q} = \mathbf{S} \bowtie \mathbf{R} \bowtie \mathbf{I}$.

We learn the model with BGD using the D-IFAQ program above. To avoid clutter, we focus on the core computation of the program and we make two simplifications: (1) we assume $\frac{\alpha}{|\mathbf{Q}|} = 1$, and (2) we hide the term for $x[a_{n+1}]$. Then the inner-loop expression is given by:

$$\theta = \lambda_{f_1 \in \mathbf{F}} \Big( \theta(f_1) - \sum_{x \in \text{dom}(\mathbf{Q})} \mathbf{Q}(x) * \big( \sum_{f_2 \in \mathbf{F}} \theta(f_2) * x[f_2] \big) * x[f_1] \Big)$$

where $\mathbf{F} = [['i', 's', 'c', 'p']]$.

**Decision Tree.** We next consider learning decision trees using the CART algorithm [7]. With our optimizations, we can learn both classification and regression trees efficiently. In the following, we focus on regression trees.

A regression tree model is a (binary) tree with inner nodes representing conditional control statements to model decisions and their consequences. Given an element $x$ in the dataset $\mathbf{Q}$, the condition for a feature $f$ is of the form $x[f]$ op $t$, which is denoted as $c(f, \text{op}, t)$ here. For categorical features (e.g., city), $t$ may be a set of categories and op denotes inclusion. For continuous features (e.g., price), $t$ is a real number and op is inequality. Leaf nodes represent predictions for the label. For regression trees, the prediction is the average of the label values in the fragment of the training dataset that satisfies all control statements on the root to leaf path.

For a given node $N$, let $\delta$ encode the conjunction of all conditions along the path from the root to $N$. The CART algorithm recursively seeks for the condition that minimises
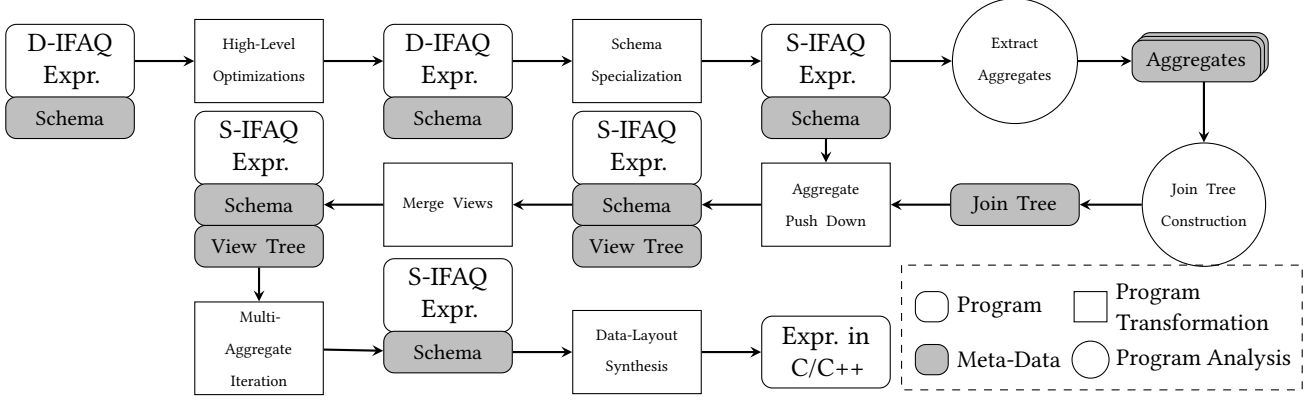
**Figure 3.** Transformation steps for an expression written in IFAQ.

the following optimization problem for a given cost function:

$$\min_{f \in F} \min_{t \in T_f} \mathsf{cost}(\mathbf{Q}, \delta \wedge c(f, \mathsf{op}, t)) + \mathsf{cost}(\mathbf{Q}, \delta \wedge c(f, !\mathsf{op}, t)).$$

where $!\mathsf{op}$ denotes the negation of $\mathsf{op}$, $T_f$ is the set of all possible thresholds for $f$. Once this condition is found, a new node with condition $c(f, \mathsf{op}, t)$ is constructed and the algorithm recursively constructs the next node for each child.

For regression trees, the cost is given by the variance, which is represented as the following D-IFAQ expression:

$$\mathsf{cost}(\mathbf{Q}, \delta') = \sum_{x \in \mathsf{dom}(\mathbf{Q})} \mathbf{Q}(x) * x[a_{n+1}]^2 * \delta' -$$
$$\frac{1}{\sum\limits_{x \in \mathsf{dom}(\mathbf{Q})} \mathbf{Q}(x) * \delta'} \Big( \sum_{x \in \mathsf{dom}(\mathbf{Q})} \mathbf{Q}(x) * x[a_{n+1}] * \delta' \Big)^2$$

In this formula, $\delta'$ can be $\delta \wedge c(f, \mathsf{op}, t)$ or $\delta \wedge c(f, !\mathsf{op}, t)$.

In contrast to linear regression, all aggregates depend on node-specific information ($\delta'$). It is thus not possible to hoist and compute them only once for all recursions of the CART algorithm. Nevertheless, all other optimizations presented in the next section are applicable.

## 4 Optimizations

This section details the IFAQ optimizations. Figure 3 overviews the transformations applied at different stages.

### 4.1 High-Level Optimizations

The high-level optimizations applied to the D-IFAQ expressions are: normalization, loop scheduling, factorization, static memoization, and loop-invariant code motion. The first two transformations are preprocessing steps for factorization. The impact of static memoization becomes positive once it is combined with loop-invariant code motion.

**Normalization.** This transformation brings the expressions into a normalized form. Similar to other algebraic-based systems (e.g., logical circuits) the normal form is sum of products. Thus, this transformation applies distributivity and pushes the products inside summation (Figure 4a).

**Example 4.1.** The inner-loop of our running example:
$$\theta = \lambda_{f_1 \in \mathbf{F}} \Big( \theta(f_1) - \sum_{x \in \mathsf{dom}(\mathbf{Q})} (\mathbf{Q}(x) * \sum_{f_2 \in \mathbf{F}} \theta(f_2) * x[f_2]) * x[f_1] \Big)$$
is normalized by pushing $x[f_1]$ into the sum over $f_2$:
$$\theta = \lambda_{f_1 \in \mathbf{F}} \Big( \theta(f_1) - \sum_{x \in \mathsf{dom}(\mathbf{Q})} \mathbf{Q}(x) * \sum_{f_2 \in \mathbf{F}} \theta(f_2) * x[f_2] * x[f_1] \Big)$$

**Loop Scheduling.** This transformation reorders nested loops (i.e., nested summations) such that the outer loop iterates over a smaller collection (Figure 4b). Pushing larger loops inside allows to factorize the free variables of the outer smaller loop outside the more expensive inner loop.

**Example 4.2.** We next reorder the loops:
$$\theta = \lambda_{f_1 \in \mathbf{F}} \Big( \theta(f_1) - \sum_{f_2 \in \mathbf{F}} \sum_{x \in \mathsf{dom}(\mathbf{Q})} \mathbf{Q}(x) * \theta(f_2) * x[f_2] * x[f_1] \Big)$$

**Factorization.** Now that the loops are correctly ordered, we factorize the multiplication operations outside summation (Figure 4c) so that they are no longer performed in an expensive loop. This results in less arithmetic operations, which in turn leads to better performance.

**Example 4.3.** We next factorize the inner-loop expression:
$$\theta = \lambda_{f_1 \in \mathbf{F}} \Big( \theta(f_1) - \sum_{f_2 \in \mathbf{F}} \theta(f_2) * \sum_{x \in \mathsf{dom}(\mathbf{Q})} \mathbf{Q}(x) * x[f_2] * x[f_1] \Big)$$

**Static Memoization.** There are repetitive computations inside a loop that cannot be easily hoisted outside. This is because they depend on variables defined inside the loop. One case stands out: the free variables are bound in loops ranging over a statically-known finite domain. In this case we can *memoize* all values as a dictionary mapping from the domain of the free variables to the values (Figure 4d).

**Example 4.4.** Static memoization results in:
$$\mathsf{let}\ \mathbf{M} = \lambda_{f_1 \in \mathbf{F}} \lambda_{f_2 \in \mathbf{F}} \sum_{x \in \mathsf{dom}(\mathbf{Q})} \mathbf{Q}(x) * x[f_2] * x[f_1]\ \mathsf{in}$$
$$\theta = \lambda_{f_1 \in \mathbf{F}} \Big( \theta(f_1) - \sum_{f_2 \in \mathbf{F}} \theta(f_2) * \mathbf{M}(f_2)(f_1) \Big)$$

**Loop-Invariant Code Motion.** After static memoization, the computations that can be shared are no longer dependent on the variables defined in the loop and can be hoisted outside of the loop: A loop-invariant let binding is moved outside the loop (Figure 4e).

**Example 4.5.** The result of loop-invariant code motion is:

```
let F = [['i', 's', 'c', 'p']] in
let M = λ   λ   ∑      Q(x) * x[f₁] * x[f₂] in
       f₁∈F f₂∈F x∈dom(Q)
θ ← θ₀
while( not converged ){
  θ = λ  (θ(f₁) − ∑  θ(f₂) * M(f₁)(f₂))
     f₁∈F         f₂∈F
}
θ
```

Here $M$ is the covar matrix (Section 3) that is automatically introduced by static memoization and hoisted outside the loop thanks to loop-invariant code motion. However, for the cases where the number of features is more than the number of elements, none of this would happen. This is because the loop scheduling optimization is not applied, and as a result static memoization and code motion will not be applied too.

## 4.2 Schema Specialization

At this stage, the dictionaries with statically-known keys of type Field (e.g., relations with the list of attributes provided as schema) are converted into records. Before checking for such dictionaries, partial evaluation transformations (Figure 4f) such as loop unrolling are performed. Also, the dynamic field accesses are converted to static field accesses. Figure 4g gives the schema specialization rules.

**Example 4.6.** Schema specialization converts the dictionary encoding the $\theta$ parameter into a record $\theta$ and all dynamic field accesses into static ones. The loop responsible for constructing the dictionary $M$ is unrolled into the creation of the record $M$. The transformed S-IFAQ program is:

```
let M = {i = {..., c =  ∑     Q(x)*x.i*x.c, ...}, ...} in
                      x∈dom(Q)
θ ← θ₀
  while( not converged ){
    θ = {i = θ.i − (... + θ.c * M.i.c + θ.p * M.i.p), ...}
  }
  θ
```

## 4.3 Aggregate Query Optimizations

These optimizations focus on the query processing aspect of IFAQ programs and aim at interleaving the aggregates and joins to achieve faster computation and smaller result.

**Example 4.7.** In our running example, $\mathbf{Q}$ is the result of joining the input relations. The computation of $\mathbf{Q}$ is represented as the following S-IFAQ expression:

```
let Q =
   ∑        ∑        ∑      (
x_s∈dom(S) x_r∈dom(R) x_i∈dom(I)
   let k = {i = x_s.i, s = x_s.s, c = x_r.c, p = x_i.p} in
   {{k →S(x_s)* R(x_r)* I(x_i)* (x_s.i==x_i.i)* (x_s.s==x_r.s)}}
)
```

We next explain how to compute the covar matrix $M$ over the query defining $\mathbf{Q}$ without materializing $\mathbf{Q}$.

**Extract aggregates.** This pass analyzes the input S-IFAQ expression and extracts a batch of aggregates from it. These aggregates are then used in S-IFAQ expressions in lieu of those that rely on the materialization of the query result. In our example, $M$ captures such aggregates.

**Join Tree Construction.** This pass takes the join defining $\mathbf{Q}$ and produces a *join tree*, where relations are nodes and an edge between two nodes is annotated with the variables on which the nodes of this edge join. The join order is computed using the state-of-the-art query optimization techniques [25] and IFAQ assumes it is given as input. The join tree is used to factorize the computation of the aggregates in $M$.

**Example 4.8.** For the join of relations $\mathbf{S}$, $\mathbf{R}$, and $\mathbf{I}$ in our running example, we may consider the join tree $\mathbf{R}\overset{i}{-}\mathbf{S}\overset{s}{-}\mathbf{I}$ where $\mathbf{S}$ is the root and $\mathbf{R}$ and $\mathbf{I}$ are its children.

**Aggregate Pushdown.** This pass decomposes each aggregate into a view per edge in the join tree, resulting in a *view tree*. Each view is used to partially push down aggregates past joins and to allow the sharing of common views across the aggregate batch. The produced intertwining of aggregates and joins are injected back into the input S-IFAQ expression.

**Example 4.9.** Let us focus on the elements $M_{c,p}$ and $M_{c,c}$ of the (nested) record $M$ that defines aggregates:

```
let M_{c,p} =   ∑     Q(x) * x.c * x.p
             x∈dom(Q)
let M_{c,c} =   ∑     Q(x) * x.c * x.c
             x∈dom(Q)
let M = {c={..., p=M_{c,p}, c=M_{c,c},...} ,...} in ...
```

Aggregate pushdown for $M_{c,p}$ yields the view tree: $V_R - M_{c,p} - V_I$. The views $V_R$ and $V_I$ compute the aggregates for x.c and x.p while iterating over the relations $\mathbf{R}$ and $\mathbf{I}$:

```
let V_R =   ∑     R(x_r) * {{{s=x_r.s} →x_r.c}} in
          x_r∈dom(R)
let V_I =   ∑     I(x_i) * {{{i=x_i.i} →x_i.p}} in
          x_i∈dom(I)
  let M_{c,p}=   ∑     S(x_s) * V_R({s=x_s.s}) * V_I({i=x_s.i})
             x_s∈dom(S)
```

A similar situation happens for the view tree $V'_R - M_{c,c} - V'_I$:

```
let V'_R =   ∑     R(x_r) * {{{s=x_r.s} →x_r.c * x_r.c}} in
           x_r∈dom(R)
let V'_I =   ∑     I(x_i) * {{{i=x_i.i} →1}} in
           x_i∈dom(I)
let M_{c,c} =   ∑     S(x_s) * V'_R({s=x_s.s}) * V'_I({i=x_s.i})
             x_s∈dom(S)
```

Computing these batches of aggregates requires multiple scans over each of the relations, the performance of which can be even worse than materializing the result of join. Next,

| | | |
|---|---|---|
| $e_1 * (e_2 + e_3)$ | $\rightsquigarrow$ | $e_1 * e_2 + e_1 * e_3$ |
| $e_1 * \sum_{x \in e_2} e_3$ | $\rightsquigarrow$ | $\sum_{x \in e_2} (e_1 * e_3)$ |
| $e_1 * (-e_2)$ | $\rightsquigarrow$ | $-(e_1 * e_2)$ |
| $-\sum_{x \in e_2} e_3$ | $\rightsquigarrow$ | $\sum_{x \in e_2} -e_3$ |

**(a)** Normalization Rules

| | | |
|---|---|---|
| $\sum_{x \in e_1} \sum_{y \in e_2} e_3$ | $\rightsquigarrow$ | $\sum_{y \in e_2} \sum_{x \in e_1} e_3$ |
| | *(if $|e_1| > |e_2|$)* | |

**(b)** Loop Scheduling Rules

| | | |
|---|---|---|
| $e_1 * e_2 + e_1 * e_3$ | $\rightsquigarrow$ | $e_1 * (e_2 + e_3)$ |
| $\sum_{x \in e_2} (e_1 * e_3)$ | $\rightsquigarrow$ | $e_1 * \sum_{x \in e_2} e_3$ |
| | *(if $x \notin fvs(e_1)$)* | |

**(c)** Factorization Rules

| | | |
|---|---|---|
| $\sum_{x \in e_1} \Gamma(\sum_{y \in e_2} e_3)$ | $\rightsquigarrow$ | let $z = \lambda_{x \in e_1} \sum_{y \in e_2} e_3$ in $\sum_{x \in e_1} \Gamma(z(x))$ |

**(d)** Static Memoization Rules

| | | |
|---|---|---|
| $\sum_{x \in e_1} ($ let $y = e_2$ in $e_3)$ | $\rightsquigarrow$ | let $y = e_2$ in $\sum_{x \in e_1} e_3$ |
| | *(if $x \notin fvs(e_2)$)* | |
| x ← e₁ <br> while(e₂) <br>   let y = e₃ in <br>   x ← e₄ <br> x | $\rightsquigarrow$ | let y = e₃ in <br> x ← e₁ <br> while(e₂) <br>   x ← e₄ <br> x |
| | *(if $x \notin fvs(e_3)$)* | |

**(e)** Loop-Invariant Code Motion Rules

| | | |
|---|---|---|
| $\sum_{x \in [[e_1,...,e_n]]} \Gamma(x)$ | $\rightsquigarrow$ | $\Gamma(e_1) + ... + \Gamma(e_n)$ |
| $\{\{e_1 \rightarrow e_2\}\} + \{\{e_1 \rightarrow e_3\}\}$ | $\rightsquigarrow$ | $\{\{e_1 \rightarrow e_2 + e_3\}\}$ |
| $\{\{e_1 \rightarrow e_2\}\} + \{\{e_3 \rightarrow e_4\}\}$ | $\rightsquigarrow$ | $\{\{e_1 \rightarrow e_2, e_3 \rightarrow e_4\}\}$ |

**(f)** Partial Evaluation Rules

| | | |
|---|---|---|
| $e_1['f']$ | $\rightsquigarrow$ | $e_1.f$ |
| $\{\{...,'f_i' \rightarrow e_i,...\}\}$ | $\rightsquigarrow$ | $\{...,f_i = e_i,...\}$ |
| $e_1(e_2)$ | $\rightsquigarrow$ | $e_1[e_2]$ |
| | *(if $e_1$ is transformed)* | |
| $\lambda_{x \in [[...,'f_i',...]]} \Gamma(e_1[x])$ | $\rightsquigarrow$ | $\{...,f_i = \Gamma(e_1.f_i),...\}$ |

**(g)** Schema Specialization Rules

| | | |
|---|---|---|
| let x = $\sum_{z \in e_1} e_2$ in <br> let y = $\sum_{z \in e_1} e_3$ in <br> $\Gamma(x, y)$ | $\rightsquigarrow$ | let t = <br> $\sum_{z \in e_1} \{x = e_2, y = e_3\}$ <br> in $\Gamma(t.x, t.y)$ |

**(h)** Loop Fusion Rule

| | | |
|---|---|---|
| let x = e₀ in $\Gamma(x)$ | $\rightsquigarrow$ | $\Gamma(e_0)$ |
| let x = e₀ in e₁ | $\rightsquigarrow$ | e₁ |
| | *(if $x \notin fvs(e_1)$)* | |
| let x = <br>   let y = e₀ in e₁ <br> in e₂ | $\rightsquigarrow$ | let y = e₀ in <br> let x = e₁ <br> in e₂ |
| let x = e₀ in <br> let y = e₀ in <br> $\Gamma(x, y)$ | $\rightsquigarrow$ | let x = e₀ in <br> $\Gamma(x, x)$ |

**(i)** Generic Optimization Rules

**Figure 4.** Transformation Rules of IFAQ. We use $\Gamma(e_1)$ on the left-hand-side to denote a context in which $e_1$ is used, and $\Gamma(e_2)$ on the right-hand-side represents the same context where each occurrence of $e_1$ is substituted by $e_2$.

we show how to share the computation across these aggregates by fusing the summations over the same collection.
**Merge Views.** This transformation consolidates the views generated in the previous transformation. All views computed at a node in the join tree will have the same key (the join variables shared between the node and its parent, or empty for the root node). These views are merged into a single view with the same key as the merged views and with all distinct aggregates of the merged views.
**Multi-Aggregate Iteration.** This transformation creates one summation over a relation for all aggregates of the view to be computed at that relation node. This computation seeks into the views computed at child nodes to fetch aggregate values used to compute the aggregates in the view. This transformation also shares computation and behaves similarly to horizontal loop fusion (Figure 4h) on S-IFAQ expressions.

**Example 4.10.** By performing horizontal loop fusion, one can merge $V_R$ and $V_R'$ into $W_R$, as well as $V_I$ and $V_I'$ into $W_I$:

$$
\text{let } W_R = \sum_{x_r \in \text{dom}(\mathbf{R})} \mathbf{R}(x_r) * \\
\{\{\{s=x_r.s\} \rightarrow \{v_R=x_r.c, v_R'=x_r.c * x_r.c\}\}\} \text{ in}
$$
$$
\text{let } W_I = \sum_{x_i \in \text{dom}(\mathbf{I})} \mathbf{I}(x_i) * \\
\{\{\{i=x_i.i\} \rightarrow \{v_I=x_i.p, v_I'=1\}\}\} \text{ in}
$$
$$
\text{let } M_{cc,pc} = \sum_{x_s \in \text{dom}(\mathbf{S})} \mathbf{S}(x_s) * \Big( \\
\text{let } w_R = W_R(\{s=x_s.s\}) \text{ in} \\
\text{let } w_I = W_I(\{i=x_s.i\}) \text{ in} \\
\{m_{c,p} = w_R.v_R * w_I.v_I, m_{c,c} = w_R.v_R' * w_I.v_I'\} \Big) \text{ in}
$$
$$
\text{let } M_{c,p} = M_{cc,pc}.m_{c,p} \text{ in let } M_{c,c} = M_{cc,pc}.m_{c,c}
$$

Rather than iterating multiple times over each relation, here we thus iterate over each relation only once.
**Dictionary To Trie.** This transformation pass converts relations and intermediate views from dictionaries into tries organized by join attributes. Hence, the summations iterating over the full domain of a relation or view, are converted into nested summations over the domain of individual attributes.

6

One key benefit of this transformation is more opportunities for factorizing the computation as well as hoisting the computation outside the introduced nested summations.

**Example 4.11.** As the intermediate views $W_R$ and $W_I$ have single fielded records as their key, changing from a dictionary to a trie does not have any impact on them. Thus, we focus only on the code for computing $M_{cc,pc}$. Rather than iterating over the domain of the key of the dictionary representing relation $S$, we have to iterate in the key hierarchy of the converted trie structure, named as $S'$. This trie data-structure is a nested dictionary containing the domain of the field $s$ at its first level, and the domain of field $i$ at its second level. The corresponding S-IFAQ expression is:

$$
\begin{aligned}
&\texttt{let } M_{cc,pc} = \\
&\quad \sum_{x_s \in \text{dom}(S')} \\
&\qquad \sum_{x_i \in \text{dom}(S'(x_s))} \\
&\qquad S'(x_s)(x_i) * \big( \\
&\qquad\qquad \texttt{let } w_R = W_R(\{\texttt{s}=x_s.\texttt{s}\}) \texttt{ in} \\
&\qquad\qquad \texttt{let } w_I = W_I(\{\texttt{i}=x_i.\texttt{i}\}) \texttt{ in} \\
&\qquad\qquad \{m_{c,p} = w_R.v_R * w_I.v_I, \\
&\qquad\qquad\quad m_{c,c} = w_R.v_R' * w_I.v_I'\} \big) \texttt{ in } \ldots
\end{aligned}
$$

This expression can be further transformed using loop-invariant code motion (Figure 4e):

$$
\begin{aligned}
&\texttt{let } M_{cc,pc} = \\
&\quad \sum_{x_s \in \text{dom}(S')} \\
&\qquad \texttt{let } w_R = W_R(\{\texttt{s}=x_s.\texttt{s}\}) \texttt{ in} \\
&\qquad \sum_{x_i \in \text{dom}(S'(x_s))} \\
&\qquad S'(x_s)(x_i) * \big( \\
&\qquad\qquad \texttt{let } w_I = W_I(\{\texttt{i}=x_i.\texttt{i}\}) \texttt{ in} \\
&\qquad\qquad \{m_{c,p} = w_R.v_R * w_I.v_I, \\
&\qquad\qquad\quad m_{c,c} = w_R.v_R' * w_I.v_I'\} \big) \texttt{ in } \ldots
\end{aligned}
$$

Finally, the multiplication operands $w_R.v_R$ and $w_R.v_R$ can be factored out of the inner summation:

$$
\begin{aligned}
&\texttt{let } M_{cc,pc} = \\
&\quad \sum_{x_s \in \text{dom}(S')} \\
&\qquad \texttt{let } w_R = W_R(\{\texttt{s}=x_s.\texttt{s}\}) \texttt{ in} \\
&\qquad \{m_{c,p} = w_R.v_R, m_{c,c} = w_R.v_R'\} * \\
&\qquad \sum_{x_i \in \text{dom}(S'(x_s))} \\
&\qquad S'(x_s)(x_i) * \big( \\
&\qquad\qquad \texttt{let } w_I = W_I(\{\texttt{i}=x_i.\texttt{i}\}) \texttt{ in} \\
&\qquad\qquad \{m_{c,p} = w_I.v_I, m_{c,c} = w_I.v_I'\} \big) \texttt{ in } \ldots
\end{aligned}
$$

Next, we move to more low-level optimizations including the physical representation of the data structures.

### 4.4 Data-Layout Synthesis

**Static Record Representation.** The generated code for records can be dictionaries from the field name to the corresponding value. This is in essence what happens in the query interpreters of database systems through using data

**Table 1.** Characteristics of the Retailer and Favorita datasets.

|  | Retailer | Favorita |
|---|---|---|
| Tuples/Size of Database | 87M(1.5GB) | 125M(2.5GB) |
| Tuples/Size of Join Result | 86M(17GB) | 127M(2.6GB) |
| Relations/Continuous Attrs | 5 / 35 | 5 / 6 |

dictionaries. However, as S-IFAQ uses code generation, this representation can be improved by generating static definitions for records (e.g., classes in Scala and structs in C/C++).
**Immutable to Mutable.** As S-IFAQ is a functional language, its data structures are immutable. Even though such data structures improve the reasoning power, their runtime performance has room for improvement. Especially, when summation produces collections, it can lead to the production of many unused intermediate collections. This can be implemented more efficiently by appending the elements into a mutable collection. As the summation construct is a tail recursive function, there are well-known optimization techniques [61] that can be applied to it; an inital empty mutable collection at the beginning is allocated, and then at each iteration an in-place update is performed.
**Scalar Replacement and Single-Field-Record Removal.** In many cases, the intermediate records are not *escaped* from their scope, meaning that their fields can be treated as local variables. This optimization is called scalar replacement, and checking its applicability is achieved using escape analysis. Furthermore, we have observed many cases where one is using a record with a single field. These records can be completely removed and substituted by their single field.
**Memory Management.** One major drawback of using languages with managed runtimes, such as JVM-based languages, is the lack of control over the memory management. By using a low-level language such as C and C++ as the target language, one can have more fine-grained control on memory management. This way, one can make sure that the remaining records from the previous optimization, are mostly stack allocated rather than being allocated on heap.
**Dictionary to Array.** S-IFAQ treats input relations as dictionaries from records to their multiplicity. However, in most cases the multiplicity is one, and it is more efficient to represent them as arrays. Furthermore, by statically setting the multiplicities to one, the compiler sees more opportunities for constant folding and other optimizations.
**Sorted Dictionary.** In many cases, one iterates over a set of keys (coming from the domain of a dictionary) and looks for the associated values with that key in other dictionaries. If all these dictionaries are sorted on that key, the process of getting the associated value for a key becomes faster; instead of looking for a key in the whole domain, it can ignore the already iterated domain.
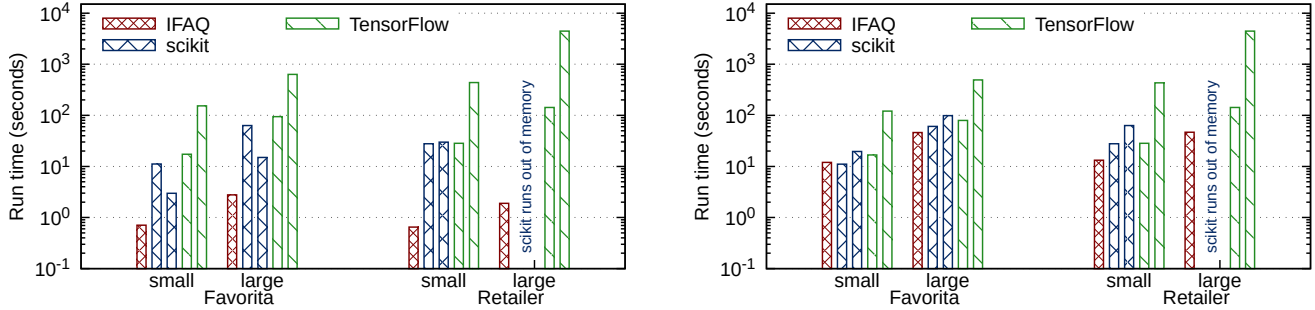
**Figure 5.** Performance comparison of IFAQ, TensorFlow, and scikit-learn for learning linear regression models (left) and regression trees (right) over the two variants of Favorita and Retailer. For TensorFlow and scikit-learn, the performance is separated into the time it takes to compute the training dataset (left bar) and the time for learning the model (right bar).

## 5  Experimental Results

In this section, we benchmark IFAQ, TensorFlow and scikit-learn for learning linear regression and regression tree models using two real-world datasets. We then micro benchmark individual IFAQ optimizations.

**Datasets.** We consider two real-world datasets, used in retail forecasting scenarios (Table 1): (1) *Favorita* [17] is a publicly available Kaggle dataset; and (2) *Retailer* is is a dataset from a US retailer [48]. Both datasets have a fact table with information about sales and respectively inventory for different stores, dates, and products. The other tables provide information about the products and stores. We consider common retail-forecasting models that predict future sales for Favorita, and future inventory demand for Retailer. For each dataset, we learn the models over the natural join of all relations, which include the sales and respectively inventory for all dates except the last month. The sales and respectively inventory for the last month is used as the test dataset to measure model accuracy. We also consider a smaller variant for each dataset whose size is 25% of the join result.

**Experimental Setup.** All experiments were performed on an Intel(R) Core(TM) i7-4770/3.40GHz/32GB/Ubuntu 18.04. We use g++ 6.4.0 for compiling the generated C++ code using the O3 optimization flag. For all the experiments, we compute the average of four subsequent runs. We do not consider the time to load the database into RAM and assume that all relations are indexed by their join attributes.

We learn the models over all continuous attributes for Favorita and Retailer. We learn regression trees up to depth four (i.e., max 31 nodes).

The input to IFAQ is a program that performs batch gradient descent optimization for linear regression models or the CART algorithm for learning regression trees. IFAQ automatically optimizes the code. For linear regression, it employs the full suite of optimizations, including the memoization and hoisting of the covar matrix. For regression trees, the aggregates cannot be hoisted, but they still benefit from the lower level optimizations, including loop fusion and data-layout

synthesis. The optimized code interleaves the computation of the training dataset with the learning task.

TensorFlow learns the linear regression model with the predefined LinearRegressor estimator. We report the time to perform one epoch with a batch size of 100,000 instances. This batch size gave the best performance/accuracy trade-off. The regression trees are learned with the BoostedTrees estimator. Since TensorFlow and scikit-learn do not support the query processing task, we use Pandas DataFrames to materialize the training dataset [35]. We compare the accuracy of all systems by comparing the root-mean-square-error (RMSE) over the test dataset. For linear regression, we compare against the closed-form solution of scikit-learn or MADlib [23] when scikit-learn fails.

Scikit-learn requires that the training dataset is represented in memory, which can lead to out-of-memory errors. In our experiments, TensorFlow also runs out of memory for learning over the full Retailer dataset. In this case, we wrote the data to a file and used the CSV parser to iteratively load mini-batches of data during learning. We do not include the time for writing the dataset to the file.

We attempted to benchmark against mlpack [14], an efficient C++ machine learning library, but it runs out-of-memory on all our experiments. One explanation is that mlpack copies the input data to compute its transpose. In our experiments, mlpack failed for as little as 5% of the Favorita training dataset (100MB on disk).

**Benchmarking End-To-End Learning.** Figure 5 shows the results of the end-to-end benchmarks for learning linear regression and regression tree models in IFAQ, scikit-learn, and TensorFlow. The runtime for the latter two is the sum of two components, depicted in Figure 6 by two bars: (1) the left bar is the time to materialize the project-join query that defines the training dataset; (2) the right bar is the time to learn the model. IFAQ learns the model directly over the input database in one computation step.
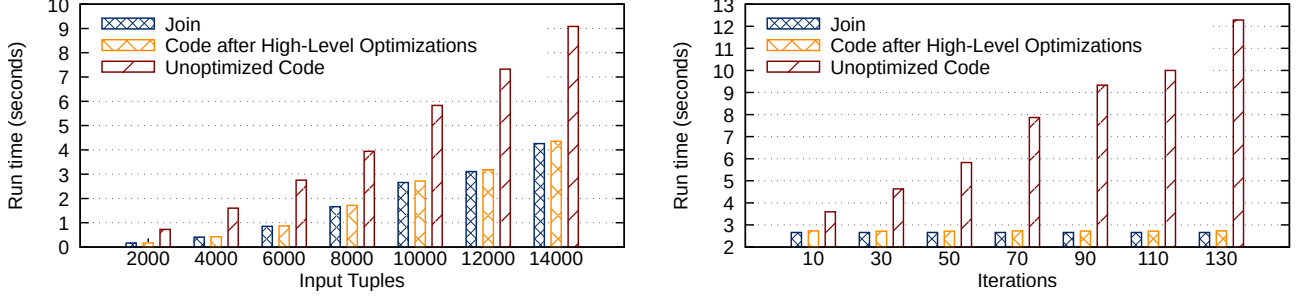
**Figure 6.** Impact of high-level optimizations for learning the linear regression model using BGD by varying the number of input tuples for 50 iterations (left) and varying the number iterations for 10,000 input tuples (right). The join computation is identical for both unoptimized and optimized code and shown by a separate bar.
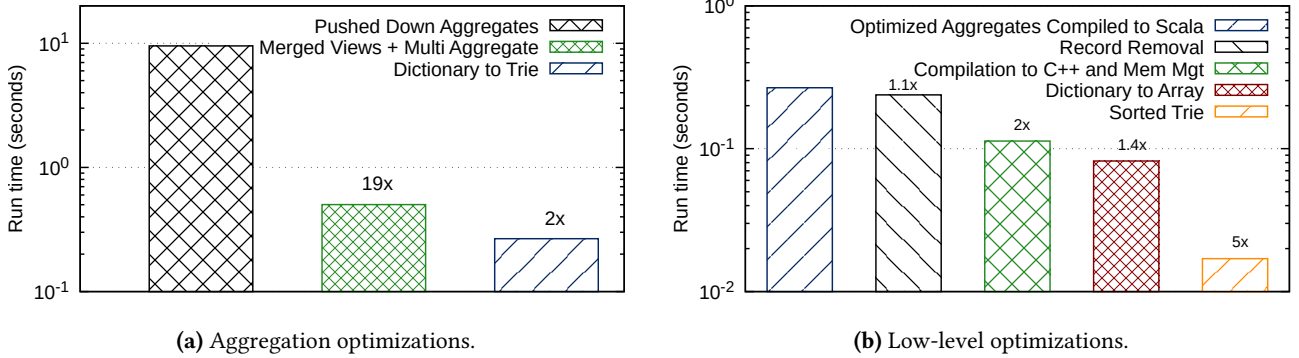


**(a)** Aggregation optimizations.

**(b)** Low-level optimizations.

**Figure 7.** Impact of aggregation and low-level optimizations for the covar matrix computation for 1,000,000 input tuples.

IFAQ significantly outperforms both scikit-learn and TensorFlow, whenever they do not fail. The end-to-end computation in IFAQ is consistently faster than it takes the competitors to materialize the training dataset. It is remarkable that IFAQ outperforms them even though it runs in one thread whereas the competitors use 8 threads (one per CPU core) for learning the model.

The reason for the performance improvement is due to the suite of optimizations in IFAQ. First, IFAQ represents the core computation of the learning algorithms as aggregate queries and computes them directly over the input database, which is significantly smaller than the join result (Table 1). It then uses several optimizations to further improve the computation of these aggregates (Section 4). We evaluate the impact of individual optimizations below.

For linear regression, the RMSE for IFAQ is within 1% of the closed form solution. TensorFlow computes a single epoch, and the RMSE is worse than IFAQ's (e.g., for Retailer small it is 3% higher). TensforFlow would thus need more epochs to achieve the accuracy of IFAQ. Scikit-learn and IFAQ learn very similar regression trees (using CART variants) so the accuracies are very close.

**Categorical Attributes.** The previous experiments only report on runtime performance for datasets with continuous features only. All aforementioned systems, including IFAQ, can also support categorical features via one-hot encoding. This translates into 87 features for Retailer and 526 features for Favorita, one per category of each categorical attribute and one per continuous attribute. We found that this encoding leads to an out-of-memory error for scikit-learn. Furthermore, IFAQ becomes much slower as it needs to generate quadratically many aggregates in the number of features. We leave the more efficient handling of categorical attributes, e.g., by extending IFAQ with sparse tensors as in LMFAO [47], as future work.

**Compilation Overhead.** The gcc compilation of the generated C++ code for Retailer is on average in 4.3 and 8.3 seconds for linear regression and respectively regression trees, while for Favorita this is 9.7 and respectively 2.4 seconds.

**Individual Optimizations.** We evaluate the effect of individual optimizations for learning the linear regression model using BGD on a subset of the Favorita dataset. We start with the impact of high-level optimizations using an interpreter for D-IFAQ. Then, we continue with the impact of the aggregate query optimizations on an interpreter for S-IFAQ.

Finally, we show the impact of lower level optimizations on a compiler for S-IFAQ which generates Scala or C++ code.

To evaluate the effect of individual high-level optimizations, we start with an unoptimized program that first materializes the join and then learns the model on this materialized dataset. This is what existing ML tools do. For linear regression, applying high-level optimizations results in hoisting the covar matrix out of the while loop (Section 4.1).

Figure 6 shows the impact of high-level optimizations by showing the performance of unoptimized and optimized programs. Both these programs need to materialize the join result. However, the unoptimized program needs to perform aggregations at each iteration on the training dataset, whereas in the optimized program, most of these data-intensive computations are hoisted outside the loop. The two graphs show that most of the computation for the optimized program is dedicated to the join computation. In addition, the increase in the number of iterations has a negligible impact on the performance of the optimized program.

For the following optimization layers, we focus on the data-intensive computation, i.e., computing the covar matrix.

The impact of aggregate optimizations (Section 4.3) on the performance of computing the covar matrix is shown in Figure 7a: (i) merging views and multi-aggregate iteration have a significant impact on the performance thanks to horizontal loop fusion; (ii) converting dictionaries to tries also improves the performance thanks to the sharing of computation enabled by factorization.

Finally, Figure 7b shows the impact of low-level optimizations (Section 4.4) on the computation of the covar matrix. As the code is more optimized, we can process a significantly larger amount of data. The leftmost bar shows the code after aggregate optimizations and compilation in Scala. The following two optimizations have the highest impact: (i) generating C++ code with efficient memory management has 2x performance improvement, which requires less heap allocation that the Scala code; and (ii) using an already sorted trie rather than a hash-table trie data-structure gives 5x speedup.

## 6 Related Work

**High Performance Computing.** There is high demand for efficient matrix processing in numerical and scientific computing. BLAS [16] exposes a set of low-level routines for common linear algebra primitives used in higher-level libraries including LINPACK, LAPACK, and ScaLAPACK for parallel processing. Highly optimized BLAS implementations are provided for dense linear algebra by hardware vendors such as Intel or AMD and code generators such as ATLAS [63] and for sparse linear algebra by Combinatorial BLAS [8]. HPAT [59] compiles a high-level scripting language into high-performance parallel code. IFAQ can reuse techniques developed in all the aforementioned tools.

**Numerical and Linear Algebra DSLs.** The compilers of high-level linear algebra DSLs, such as Lift [56], Opt [15], Halide [45], Diderot [12], and OptiML [57], focus on generating efficient parallel code from the high-level programs.

Spiral [42] provides a domain-specific compiler for synthesizing Digital Signal Processing kernels, e.g., Fourier transforms. It is based on the SPL [65] language that expresses recursion and mathematical formulas. The LGen [55] and SLinGen [54] systems target small-scale computations involving fixed-size linear algebra expressions common in graphics and media processing applications. They use two-level DSLs, namely LL to perform tiling decisions and Σ-LL to enable loop level optimizations. The generated output is a C function that includes intrinsics to enable SIMD vector extensions.

TACO [30] generates efficient low-level code for compound linear algebra operations on dense and sparse matrices. Workspaces [29] further improve the performance of the generated kernels by increasing sharing, and enabling optimizations such as loop-invariant code motion. It has similarities to our static memoization optimization (Section 4.1).

APL [26] is the pioneer of array languages, the design of which inspired functional array languages such as SAC [21], Futhark [24], and $\widetilde{F}$ [51]. The key feature of all these languages is the support for fusion [3, 13, 50, 58], which is essential for efficient low-level code generation.

Finally, there are mainstream programming languages used by data scientists. The R programming language [44] is widely used by statisticians and data miners. It provides a standard language for statistical computing that includes arithmetic, array manipulation, object oriented programming and system calls. Similarly, MATLAB provides a wide range of toolboxes for data analytics tasks.

These DSLs do not express the hybrid query processing and machine learning workloads that IFAQ supports. Next, we highlight techniques used to overcome this issue.

**Functional Languages for Data Processing.** Apart from the linear-algebra-based DSLs, there are languages for data processing based on the nested relational model [46] and the monad calculus and monoid comprehension [5, 6, 9, 18, 22, 60, 62, 64]. Functional collection programming abstractions existing in the standard library of the mainstream programming languages such as Scala, Haskell, and recently Java 8 are also based on these languages. A similar abstraction is used in distributed data processing frameworks such as Spark [66], as well as as an intermediate representation in Weld [39]. Similar to IFAQ, DBToaster [31] uses a bag-based collection, yet for incremental view maintenance. A key challenge for such collection-based languages is efficient code generation, which has been investigated both by the DB [39, 52, 53] and PL [28, 33, 49] communities. The IFAQ languages are also inspired by the same model, with more similarities to monoid comprehension; the summation operator can be considered

as a monoid comprehension. However, they are more expressive and allow for a wider range of optimizations and generation of efficient low-level code for aggregate batches. **In-Database Machine Learning.** The common approach to learning models over databases is to first construct the training dataset using a query engine, and then learn the model over the materialized training dataset using a machine learning framework. Python Pandas or database systems such as PostgreSQL and SparkSQL [4] are the most common tools used for the first step. The second step commonly uses scikit-learn [41], TensorFlow [1], PyTorch [40], R [44], MLlib [36], SystemML [20], or XGBoost [11]. On-going efforts avoid the expensive data export/import at the interface between the two steps, e.g., MLlib over SparkSQL, and the Python packages over Pandas. MADlib [23], Bismarck [19], and GLADE PF-OLA [43] define ML tasks as user-defined aggregate functions inside database systems so that the ML tasks share the same processing environment with the query engine. The two steps are nevertheless treated as black boxes and executed after the training dataset is materialized.

Alternative approaches avoid the materialization of the training dataset. Morpheus [10, 32] reformulates in-database machine learning tasks in terms of linear algebra operations on top of R [10] and NumPy [32]. It supports a limited class of schema definitions (i.e., key-foreign key star or chain joins) and learns models expressible in linear algebra. IFAQ is the latest development of a sustained effort to train efficiently machine learning models over arbitrary joins. Its predecessors are: F [38, 48] for linear regression; AC/DC [27] for polynomial regression and factorization machines; and LM-FAO [47] for models whose data-intensive computation can be expressed as batches of aggregates. IFAQ introduces: high-level optimizations, a systematic treatment of the various optimizations through compilation stages, and languages expressive enough to capture the full ML application.

## 7 Conclusion and Future Work

In this paper we introduced IFAQ, an optimization framework for relational learning applications, which includes the construction of the data matrix via feature extraction queries over multi-relational data and the subsequent model training step. IFAQ takes as input such an application expressed in a dynamically-typed language that captures a fragment of scripting languages such as Python. IFAQ then applies multiple layers of optimizations that are inspired by techniques developed by the databases, programming languages, and high-performance computing communities. As proof of concept, we showed that IFAQ outperforms mainstream approaches such as TensorFlow and Scikit by orders of magnitude for linear regression and regression tree models.

This work opens exciting avenues of future research. New technical developments include: a compilation approach to capture LMFAO's efficient support for categorical variables with multi-root join trees and group-by aggregates [47]; support for parallelization and many-core architectures; and an investigation of the trade-off between runtime performance and size of generated C++ code for models with high degree and many parameters (e.g., factorization machines). We would also like to improve the usability of IFAQ as follows: build an IFAQ library of optimization algorithms and ML models beyond the simple ones discussed in this paper and including boosting trees, random forests, and neural networks; generate optimized code for model selection over different subsets of the given variables; allow IFAQ to work directly on Jupyter notebooks that specify the construction of the data matrix and the model training; and investigate whether the IFAQ compilation techniques can be incorporated into popular data science tools such as Scikit and TensorFlow.

## References

[1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. 2016. TensorFlow: A System for Large-Scale Machine Learning.. In *Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation (OSDI'16)*. USENIX Association, USA, 265–283.

[2] Mahmoud Abo Khamis, Hung Q. Ngo, and Atri Rudra. 2016. FAQ: Questions Asked Frequently. In *Proceedings of the 35th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems (PODS '16)*. Association for Computing Machinery, New York, NY, USA, 13–28.

[3] Johan Anker and Josef Svenningsson. 2013. An EDSL approach to high performance Haskell programming. In *ACM Haskell Symposium*. ACM, New York, NY, USA, 1–12.

[4] Michael Armbrust, Reynold S. Xin, Cheng Lian, Yin Huai, Davies Liu, Joseph K. Bradley, Xiangrui Meng, Tomer Kaftan, Michael J. Franklin, Ali Ghodsi, and Matei Zaharia. 2015. Spark SQL: Relational Data Processing in Spark. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data (SIGMOD '15)*. ACM, New York, NY, USA, 1383–1394.

[5] Val Breazu-Tannen, Peter Buneman, and Limsoon Wong. 1992. Naturally embedded query languages. In *ICDT'92*. Springer Berlin Heidelberg, Berlin, Heidelberg, 140–154.

[6] Val Breazu-Tannen and Ramesh Subrahmanyam. 1991. Logical and computational aspects of programming with sets/bags/lists. In *Automata, Languages and Programming*. Springer Berlin Heidelberg, Berlin, Heidelberg, 60–75.

[7] L. Breiman, J. Friedman, R. Olshen, and C. Stone. 1984. *Classification and Regression Trees*. Wadsworth and Brooks, Monterey, CA.

[8] Aydın Buluç and John R Gilbert. 2011. The Combinatorial BLAS: Design, Implementation, and Applications. *Int. J. High Perform. Comput. Appl.* 25, 4 (Nov. 2011), 496–509.

[9] Peter Buneman, Shamim Naqvi, Val Tannen, and Limsoon Wong. 1995. Principles of Programming with Complex Objects and Collection Types. *Theor. Comput. Sci.* 149, 1 (Sept. 1995), 3–48.

[10] Lingjiao Chen, Arun Kumar, Jeffrey Naughton, and Jignesh M Patel. 2017. Towards linear algebra over normalized data. *Proceedings of the VLDB Endowment* 10, 11 (2017), 1214–1225.

[11] Tianqi Chen and Carlos Guestrin. 2016. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '16)*. Association for Computing Machinery, New York, NY, USA, 785–794.

[12] Charisee Chiw, Gordon Kindlmann, John Reppy, Lamont Samuels, and Nick Seltzer. 2012. Diderot: A Parallel DSL for Image Analysis and Visualization. In *Proceedings of the 33rd ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI'12)*. ACM, 111–120.

[13] Koen Claessen, Mary Sheeran, and Bo Joel Svensson. 2012. Expressive Array Constructs in an Embedded GPU Kernel Programming Language. In *Proceedings of the 7th Workshop on Declarative Aspects and Applications of Multicore Programming (DAMP '12)*. ACM, NY, USA, 21–30.

[14] Ryan R. Curtin, Marcus Edel, Mikhail Lozhnikov, Yannis Mentekidis, Sumedh Ghaisas, and Shangtong Zhang. 2018. mlpack 3: a fast, flexible machine learning library. *J. Open Source Software* 3 (2018), 726. Issue 26.

[15] Zachary Devito, Michael Mara, Michael Zollhöfer, Gilbert Bernstein, Jonathan Ragan-Kelley, Christian Theobalt, Pat Hanrahan, Matthew Fisher, and Matthias Niessner. 2017. Opt: A Domain Specific Language for Non-Linear Least Squares Optimization in Graphics and Imaging. *ACM Trans. Graph.* 36, 5, Article Article 171 (Oct. 2017), 27 pages.

[16] J. J. Dongarra, Jeremy Du Croz, Sven Hammarling, and I. S. Duff. 1990. A Set of Level 3 Basic Linear Algebra Subprograms. *ACM Trans. Math. Softw.* 16, 1 (March 1990), 1–17.

[17] Corporacion Favorita. 2017. Corp. Favorita Grocery Sales Forecasting: Can you accurately predict sales for a large grocery chain? https://www.kaggle.com/c/favorita-grocery-sales-forecasting/

[18] Leonidas Fegaras and David Maier. 2000. Optimizing Object Queries Using an Effective Calculus. *TODS* 25, 4 (Dec. 2000), 457–516.

[19] Xixuan Feng, Arun Kumar, Benjamin Recht, and Christopher Ré. 2012. Towards a Unified Architecture for in-RDBMS Analytics. In *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data (SIGMOD '12)*. ACM, New York, NY, USA, 325–336.

[20] Amol Ghoting, Rajasekar Krishnamurthy, Edwin Pednault, Berthold Reinwald, Vikas Sindhwani, Shirish Tatikonda, Yuanyuan Tian, and Shivakumar Vaithyanathan. 2011. SystemML: Declarative machine learning on MapReduce. In *ICDE'11*. IEEE, 231–242.

[21] Clemens Grelck and Sven-Bodo Scholz. 2006. SAC—A Functional Array Language for Efficient Multi-threaded Execution. *Int. Journal of Parallel Programming* 34, 4 (2006), 383–427.

[22] Torsten Grust and MarcH. Scholl. 1999. How to Comprehend Queries Functionally. *Journal of Intelligent Information Systems* 12, 2-3 (1999), 191–218.

[23] Joseph M Hellerstein, Christoper Ré, Florian Schoppmann, Daisy Zhe Wang, Eugene Fratkin, Aleksander Gorajek, Kee Siong Ng, Caleb Welton, Xixuan Feng, Kun Li, et al. 2012. The MADlib analytics library: or MAD skills, the SQL. *Proceedings of the VLDB Endowment* 5, 12 (2012), 1700–1711.

[24] Troels Henriksen, Niels GW Serup, Martin Elsman, Fritz Henglein, and Cosmin E Oancea. 2017. Futhark: purely functional GPU-programming with nested parallelism and in-place array updates. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation*. ACM, 556–571.

[25] Y. E. Ioannidis and Younkyung Kang. 1990. Randomized Algorithms for Optimizing Large Join Queries. In *Proceedings of the 1990 ACM SIGMOD International Conference on Management of Data (SIGMOD '90)*. ACM, New York, NY, USA, 312–321.

[26] Kenneth E Iverson. 1962. A Programming Language. In *Proceedings of the May 1-3, 1962, spring joint computer conference*. ACM, 345–351.

[27] Mahmoud Abo Khamis, Hung Q. Ngo, XuanLong Nguyen, Dan Olteanu, and Maximilian Schleich. 2018. AC/DC: In-Database Learning Thunderstruck. In *Proceedings of the Second Workshop on Data Management for End-To-End Machine Learning (DEEM'18)*. ACM, New York, NY, USA, Article 8, 10 pages.

[28] Oleg Kiselyov, Aggelos Biboudis, Nick Palladinos, and Yannis Smaragdakis. 2017. Stream Fusion, to Completeness. In *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages (POPL 2017)*. ACM, New York, NY, USA, 285–299.

[29] Fredrik Kjolstad, Peter Ahrens, Shoaib Kamil, and Saman Amarasinghe. 2019. Tensor Algebra Compilation with Workspaces. In *Proceedings of the 2019 IEEE/ACM International Symposium on Code Generation and Optimization (CGO 2019)*. IEEE Press, Piscataway, NJ, USA, 180–192.

[30] Fredrik Kjolstad, Shoaib Kamil, Stephen Chou, David Lugato, and Saman Amarasinghe. 2017. The Tensor Algebra Compiler. *Proc. ACM Program. Lang.* 1, OOPSLA, Article 77 (Oct. 2017), 29 pages.

[31] Christoph Koch, Yanif Ahmad, Oliver Kennedy, Milos Nikolic, Andres Nötzli, Daniel Lupei, and Amir Shaikhha. 2014. DBToaster: Higher-order delta processing for dynamic, frequently fresh views. *VLDBJ* 23, 2 (2014).

[32] Side Li, Lingjiao Chen, and Arun Kumar. 2019. Enabling and Optimizing Non-linear Feature Interactions in Factorized Linear Algebra. In *Proceedings of the 2019 International Conference on Management of Data*. ACM, 1571–1588.

[33] Geoffrey Mainland, Roman Leshchinskiy, and Simon Peyton Jones. 2013. Exploiting Vector Instructions with Generalized Stream Fusion. In *Proceedings of the 18th ACM SIGPLAN International Conference on Functional Programming (ICFP'13)*. ACM, New York, NY, USA, 37–48.

[34] Wes McKinney. 2011. pandas: a foundational Python library for data analysis and statistics. In *Python for High Performance and Scientific Computing*, Vol. 14. 1–9.

[35] Wes McKinney. 2012. *Python for data analysis: Data wrangling with Pandas, NumPy, and IPython.* " O'Reilly Media, Inc.".

[36] Xiangrui Meng, Joseph Bradley, Burak Yavuz, Evan Sparks, Shivaram Venkataraman, Davies Liu, Jeremy Freeman, DB Tsai, Manish Amde, Sean Owen, Doris Xin, Reynold Xin, Michael J. Franklin, Reza Zadeh, Matei Zaharia, and Ameet Talwalkar. 2016. MLlib: Machine Learning in Apache Spark. *The Journal of Machine Learning Research* 17, 1 (2016), 1235–1241.

[37] Bruce Momjian. 2001. *PostgreSQL: introduction and concepts*. Vol. 192. Addison-Wesley New York.

[38] Dan Olteanu and Maximilian Schleich. 2016. Factorized Databases. *SIGMOD Rec.* 45, 2 (Sept. 2016), 5–16.

[39] Shoumik Palkar, James J Thomas, Anil Shanbhag, Deepak Narayanan, Holger Pirk, Malte Schwarzkopf, Saman Amarasinghe, Matei Zaharia, and Stanford InfoLab. 2017. Weld: A Common Runtime for High Performance Data Analytics. In *Conference on Innovative Data Systems Research (CIDR)*.

[40] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic Differentiation in PyTorch. In *NIPS 2017 Autodiff Workshop: The Future of Gradient-based Machine Learning Software and Techniques*.

[41] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in Python. *Journal of machine learning research* 12, Oct (2011), 2825–2830.

[42] Markus Püschel, José M. F. Moura, Jeremy R. Johnson, David Padua, Manuela M. Veloso, Bryan W. Singer, Jianxin Xiong, Franz Franchetti, Aca Gacic, Yevgen Voronenko, Kang Chen, Robert W. Johnson, and Nicholas Rizzolo. 2005. SPIRAL: Code generation for DSP transforms. *Proc. IEEE* 93, 2 (2005), 232–275.

[43] Chengjie Qin and Florin Rusu. 2015. Speculative approximations for terascale distributed gradient descent optimization. In *Proceedings of the Fourth Workshop on Data analytics in the Cloud*. ACM, 1.

[44] R Core Team. 2013. *R: A Language and Environment for Statistical Computing*. R Foundation for Stat. Comp., www.r-project.org.

[45] Jonathan Ragan-Kelley, Connelly Barnes, Andrew Adams, Sylvain Paris, Frédo Durand, and Saman Amarasinghe. 2013. Halide: A Language and Compiler for Optimizing Parallelism, Locality, and Recomputation in Image Processing Pipelines. In *Proceedings of the 34th ACM*

*SIGPLAN Conference on Programming Language Design and Implementation (PLDI'13)*. ACM, New York, NY, USA, 519–530.

[46] Mark A Roth, Herry F Korth, and Abraham Silberschatz. 1988. Extended algebra and calculus for nested relational databases. *ACM Transactions on Database Systems (TODS)* 13, 4 (1988), 389–417.

[47] Maximilian Schleich, Dan Olteanu, Mahmoud Abo Khamis, Hung Q. Ngo, and XuanLong Nguyen. 2019. A Layered Aggregate Engine for Analytics Workloads. In *Proceedings of the 2019 International Conference on Management of Data (SIGMOD '19)*. ACM, New York, NY, USA, 1642–1659.

[48] Maximilian Schleich, Dan Olteanu, and Radu Ciucanu. 2016. Learning Linear Regression Models over Factorized Joins. In *Proceedings of the 2016 International Conference on Management of Data (SIGMOD '16)*. ACM, New York, NY, USA, 3–18.

[49] Amir Shaikhha, Mohammad Dashti, and Christoph Koch. 2018. Push versus Pull-Based Loop Fusion in Query Engines. *Journal of Functional Programming* 28 (2018), e10.

[50] Amir Shaikhha, Andrew Fitzgibbon, Simon Peyton Jones, and Dimitrios Vytiniotis. 2017. Destination-passing Style for Efficient Memory Management. In *Proceedings of the 6th ACM SIGPLAN International Workshop on Functional High-Performance Computing (FHPC 2017)*. ACM, New York, NY, USA, 12–23.

[51] Amir Shaikhha, Andrew Fitzgibbon, Dimitrios Vytiniotis, and Simon Peyton Jones. 2019. Efficient differentiable programming in a functional array-processing language. *Proceedings of the ACM on Programming Languages* 3, ICFP (2019), 97.

[52] Amir Shaikhha, Yannis Klonatos, and Christoph Koch. 2018. Building Efficient Query Engines in a High-Level Language. *ACM Transactions on Database Systems* 43, 1, Article 4 (April 2018), 45 pages.

[53] Amir Shaikhha, Yannis Klonatos, Lionel Parreaux, Lewis Brown, Mohammad Dashti, and Christoph Koch. 2016. How to Architect a Query Compiler. In *Proceedings of the 2016 International Conference on Management of Data (SIGMOD'16)*. ACM, New York, NY, USA, 1907–1922.

[54] Daniele G. Spampinato, Diego Fabregat-Traver, Paolo Bientinesi, and Markus Püschel. 2018. Program Generation for Small-scale Linear Algebra Applications. In *Proceedings of the 2018 International Symposium on Code Generation and Optimization (CGO 2018)*. ACM, New York, NY, USA, 327–339.

[55] Daniele G. Spampinato and Markus Püschel. 2014. A Basic Linear Algebra Compiler. In *Proceedings of Annual IEEE/ACM International Symposium on Code Generation and Optimization (CGO '14)*. Association for Computing Machinery, New York, NY, USA, 23–32.

[56] Michel Steuwer, Christian Fensch, Sam Lindley, and Christophe Dubach. 2015. Generating Performance Portable Code Using Rewrite Rules: From High-level Functional Expressions to High-performance OpenCL Code. In *Proceedings of the 20th ACM SIGPLAN International Conference on Functional Programming (ICFP 2015)*. ACM, New York, NY, USA, 205–217.

[57] Arvind Sujeeth, HyoukJoong Lee, Kevin Brown, Tiark Rompf, Hassan Chafi, Michael Wu, Anand Atreya, Martin Odersky, and Kunle Olukotun. 2011. OptiML: An Implicitly Parallel Domain-Specific Language for Machine Learning. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11) (ICML '11)*. 609–616.

[58] Bo Joel Svensson and Josef Svenningsson. 2014. Defunctionalizing Push Arrays. In *Proceedings of the 3rd ACM SIGPLAN Workshop on Functional High-performance Computing (FHPC '14)*. ACM, NY, USA, 43–52.

[59] Ehsan Totoni, Todd A Anderson, and Tatiana Shpeisman. 2017. HPAT: high performance analytics with scripting ease-of-use. In *Proceedings of the International Conference on Supercomputing*. ACM, 9.

[60] Phil Trinder. 1992. Comprehensions, a Query Notation for DBPLs. In *Proc. of the 3rd DBPL workshop (DBPL3)*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 55–68.

[61] Philip Wadler. 1984. Listlessness is better than laziness: Lazy evaluation and garbage collection at compile-time. In *Proceedings of the 1984 ACM Symposium on LISP and functional programming*. ACM, 45–52.

[62] Philip Wadler. 1990. Comprehending Monads. In *Proceedings of the 1990 ACM Conference on LISP and Functional Programming (LFP '90)*. ACM, New York, NY, USA, 61–78.

[63] Clint Whaley and Jack Dongarra. 1998. Automatically Tuned Linear Algebra Software. In *Proceedings of the 1998 ACM/IEEE Conference on Supercomputing (SC '98)*. IEEE Computer Society, USA, 1–27.

[64] Limsoon Wong. 2000. Kleisli, a functional query system. *Journal of Functional Programming* 10, 1 (2000), 19–56.

[65] Jianxin Xiong, Jeremy Johnson, Robert Johnson, and David Padua. 2001. SPL: A Language and Compiler for DSP Algorithms. In *Proceedings of the ACM SIGPLAN 2001 Conference on Programming Language Design and Implementation (PLDI'01)*. ACM, New York, NY, USA, 298–308.

[66] Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, and Ion Stoica. 2012. Resilient Distributed Datasets: A Fault-tolerant Abstraction for In-memory Cluster Computing. In *Proceedings of the 9th USENIX Conference on Networked Systems Design and Implementation (NSDI'12)*. USENIX Association, 1.