

Cost Models for Big Data Query Processing: Learning, Retrofitting, and Our Findings

Tarique Siddiqui^{1,2}, Alekh Jindal¹, Shi Qiao¹, Hiren Patel¹, Wangchao Le¹

¹Microsoft ²University of Illinois, Urbana-Champaign

ABSTRACT

Query processing over big data is ubiquitous in modern clouds, where the system takes care of picking both the physical query execution plans *and* the resources needed to run those plans, using a cost-based query optimizer. A good cost model, therefore, is akin to better resource efficiency and lower operational costs. Unfortunately, the production workloads at Microsoft show that costs are very complex to model for big data systems. In this work, we investigate two key questions: (i) can we *learn* accurate cost models for big data systems, and (ii) can we *integrate* the learned models within the query optimizer. To answer these, we make three core contributions. First, we exploit workload patterns to learn a large number of individual cost models and combine them to achieve high accuracy and coverage over a long period. Second, we propose extensions to Cascades framework to pick optimal resources, i.e., number of containers, during query planning. And third, we integrate the learned cost models within the Cascade-style query optimizer of SCOPE at Microsoft. We evaluate the resulting system, CLEO, in a production environment using both production and TPC-H workloads. Our results show that the learned cost models are 2 to 3 orders of magnitude more accurate, and 20× more correlated with the actual runtimes, with a large majority (70%) of the plan changes leading to substantial improvements in latency as well as resource usage.

ACM Reference Format:

Tarique Siddiqui, Alekh Jindal, Shi Qiao, Hiren Patel, Wangchao Le. 2020. Cost Models for Big Data Query Processing: Learning, Retrofitting, and Our Findings. In *In Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data (SIGMOD'20)*, June 14–19, 2020, Portland, OR, USA. ACM, NY, NY, USA, 15 pages. <https://doi.org/10.1145/3318464.3380584>

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

SIGMOD'20, June 14–19, 2020, Portland, OR, USA

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-6735-6/20/06...\$15.00

<https://doi.org/10.1145/3318464.3380584>

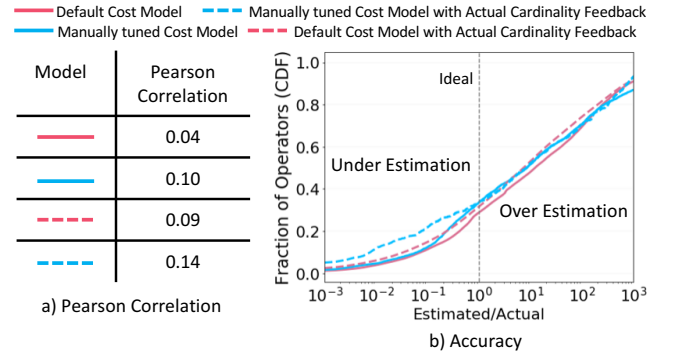


Figure 1: Impact of manual tuning and cardinality feedback on cost models in SCOPE

1 INTRODUCTION

There is a renewed interest in cost-based query optimization in big data systems, particularly in modern cloud data services (e.g., Athena [4], ADLA [36], BigSQL [20], and BigQuery [18]) that are responsible for picking both the query execution plans and the resources (e.g., number of containers) needed to run those plans. Accurate *cost models* are therefore crucial for generating efficient combination of plan and resources. Yet, the traditional wisdom from relational databases is that cost models are less important and fixing cardinalities automatically fixes the cost estimation [28, 30]. The question is whether this also holds for the new breed of big data systems. To dig deeper, we analyzed one day's worth of query logs from the big data infrastructure (SCOPE [10, 45]) at Microsoft. We feed back the actual runtime cardinalities, i.e., the ideal estimates that any cardinality estimator, including learned models [13, 24, 39, 42] can achieve. Figure 1 compares the ratio of cost estimates with the actual runtimes for two cost models in SCOPE: 1) a default cost model, and 2) a manually-tuned cost model that is partially available for limited workloads. The vertical dashed-line at 10⁰ corresponds to an ideal situation where all cost estimates are equal to the actual runtimes. Thus, the closer a curve is to the dashed line, the more accurate it is.

The dotted lines in Figure 1(b) show that fixing cardinalities reduces the over-estimation, but there is still a wide gap between the estimated and the actual costs, with the Pearson correlation being as low as 0.09. This is due to the complexity of big data systems coupled with the variance in cloud environments [38], which makes cost modeling incredibly difficult. Furthermore, any improvements in cost modeling

need to be consistent across workloads and over time since performance spikes are detrimental to the reliability expectations of enterprise customers. Thus, accurate cost modeling is still a challenge in SCOPE like big data systems.

In this paper, we explore the following two questions:

(1) Can we learn accurate, yet robust cost models for big data systems? This is motivated by the presence of massive workloads visible in modern cloud services that can be harnessed to accurately model the runtime behavior of queries. This helps not only in dealing with the various complexities in the cloud, but also specializing or *instance optimizing* [32] to specific customers or workloads, which is often highly desirable. Additionally, in contrast to years of experience needed to tune traditional optimizers, learned cost models are potentially easy to update at a regular frequency.

(2) Can we effectively integrate learned cost models within the query optimizer? This stems from the observation that while some prior works have considered learning models for predicting query execution times for a given physical plan in traditional databases [2, 5, 17, 29], none of them have integrated learned models within a query optimizer for selecting physical plans. Moreover, in big data systems, resources (in particular the number of machines) play a significant role in cost estimation [41], making the integration even more challenging. Thus, we investigate the effects of learned cost models on query plans by *extending* the SCOPE query optimizer in a *minimally invasive* way for *predicting costs* in a *resource-aware* manner.

Our key ideas are as follows. We note that the cloud workloads are quite diverse in nature, i.e., there is no representative workload to tune the query optimizer, and hence there is no single cost model that fits the entire workload, i.e., *no-one-size-fits-all*. Therefore, we learn a large collection of smaller-sized cost models, one for each common subexpressions that are typically abundant in production query workloads [21, 22]. While this approach results in specialized cost models that are very accurate, the models do not cover the entire workload: expressions that are not common across queries do not have models. The other extreme is to learn a cost model per operator, which covers the entire workload but sacrifices the accuracy with very general models. Thus, *there is an accuracy-coverage trade-off that makes cost modeling challenging*. To address this, we define the notion of cost model *robustness* with three desired properties: (i) high accuracy, (ii) high coverage, and (iii) high retention, i.e., stable performance for a long-time period before retraining. We achieve these properties in two steps: First, we bridge the accuracy-coverage gap by learning additional mutually enhancing models that improve the coverage as well as the accuracy. Then, we learn a combined model that automatically corrects and combines the predictions from multiple individual models, providing accurate and stable predictions for a sufficiently long window (e.g., more than 10 days).

We implemented our ideas in a Cloud LEarning Optimizer (CLEO) and integrated it within SCOPE. CLEO uses a feedback loop to periodically train and update the learned cost models

within the Cascades-style top-down query planning [19] in SCOPE. We extend the optimizer to invoke the learned models, instead of the default cost models, to estimate the cost of candidate operators. However, in big data systems, the cost depends heavily on the resources used (e.g., number of machines for each operator) by the optimizer [41]. Therefore, we extend the Cascades framework to explore resources, and propose mechanisms to explore and derive optimal number of machines for each stage in a query plan. Moreover, instead of using handcrafted heuristics or assuming fixed resources, we leverage the learned cost models to find optimal resources as part of query planning, thereby *using learned models for producing both runtime as well as resource-optimal plans*.

In summary, our key contributions are as follows.

- (1) We motivate the cost estimation problem from production workloads at Microsoft, including prior attempts for manually improving the cost model (Section 2).
- (2) We propose machine learning techniques to learn highly *accurate* cost models. Instead of building a generic cost model for the entire workload, we learn a large collection of smaller specialized models that are resource-aware and highly accurate in predicting the runtime costs (Section 3).
- (3) We describe the accuracy and coverage trade-off in learned cost models, show the two extremes, and propose additional models to bridge the gap. We combine the predictions from individual models into a *robust model* that provides the best of both accuracy and coverage over a long period (Section 4).
- (4) We describe integrating CLEO within SCOPE, including periodic training, feedback loop, model invocations during optimization, and novel extensions for finding the optimal resources for a query plan (Section 5).
- (5) Finally, we present a detailed evaluation of CLEO, using both the production workloads and the TPC-H benchmark. Our results show that CLEO improves the correlation between predicted cost and actual runtimes from 0.1 to 0.75, the accuracy by 2 to 3 orders of magnitude, and the performance for 70% of the changed plans (Section 6). In our technical report [12], we further describe practical techniques to address performance regressions in our production settings.

2 MOTIVATION

In this section, we give an overview of SCOPE, its workload and query optimizer, and motivate the cost modeling problem from production workloads at Microsoft.

2.1 Overview of SCOPE

SCOPE [10, 45] is the big data system used for internal data analytics across the whole of Microsoft to analyze and improve its various products. It runs on a hyper scale infrastructure consisting of hundreds of thousands of machines, running a massive workload of hundreds of thousands of jobs per day that process exabytes of data. SCOPE exposes a job service interface where users submit their analytical queries and the system takes care of automatically provisioning resources and running queries in a distributed environment.

SCOPE query processor partitions data into smaller subsets and processes them in parallel. The number of machines

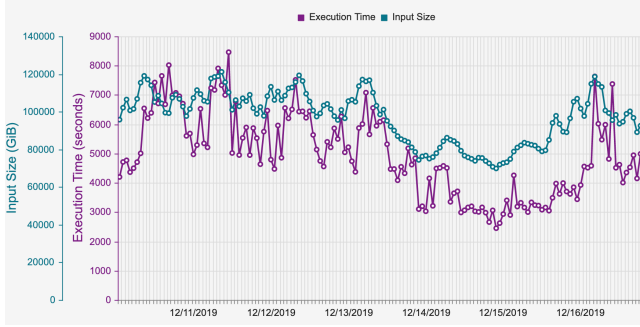


Figure 2: 150 instances of an hourly recurring job that extracts facts from a production clickstream.

running in parallel (i.e., degree of parallelism) depends on the number of partitions of the input. When no specific partitioning is required by upstream operators, certain physical operators (e.g., Extract and Exchange (also called Shuffle)), decide partition counts based on data statistics and heuristics. The sequence of intermediate operators that operate over the same set of input partitions are grouped into a stage — all operators in a stage run on the same set of machines. Except for selected scenarios, Exchange operator is commonly used to re-partition data between two stages.

2.2 Recurring Workloads

SCOPE workloads primarily consist of **recurring jobs**. A recurring job in SCOPE is used to provide periodic (e.g., hourly, six-hourly, daily, etc.) analytical result for a specific application functionality. Typically, a recurring job consists of a script template that accepts different input parameters similar to SQL modules. Each *instance* of the recurring job runs on different input data, parameters and have potentially different statements. As a result, each instance is different in terms of input/output sizes, query execution plan, total compute hour, end-to-end latency, etc. Figure 2 shows 150 instances of an hourly recurring job that extracts facts from a production clickstream. Over these 150 instances, we can see a big change in the total input size and the total execution time, from 69,859 GiB to 118,625 GiB and from 40 mins and 50 seconds to 2 hours and 21 minutes respectively. Note that a smaller portion of SCOPE workload is ad-hoc as well. Our analysis from four of the production clusters shows 7%–20% ad-hoc jobs on a daily basis. However, compared to ad-hoc jobs, recurring jobs represent long term business logic with critical value, and hence the focus of several prior research works [1, 6–8, 15, 21–23, 27, 42] and also the primary focus for performance improvement in this paper.

2.3 Overview of SCOPE Optimizer

SCOPE uses a cost-based optimizer based on the Cascades Framework [19] for generating the execution plan for a given query. Cascades [19] transforms a logical plan using multiple tasks: (i) Optimize Groups, (ii) Optimize Expressions, (iii) Explore Groups, and (iv) Explore Expressions, and (v) Optimize Inputs. While the first four tasks search for candidate plans via transformation rules, our focus in this work is essentially

on the Optimize Inputs tasks, where the cost of a physical operator is estimated. The cost of an operator is modeled to capture its *runtime latency*, estimated using a combination of data statistics and hand-crafted heuristics developed over many years. Cascades performs optimization in a top-down fashion, where physical operators higher in the plan are identified first. The exclusive (or local) costs of physical operators are computed and combined with costs of children operators to estimate the total cost. In some cases, operators can have a *required property* (e.g., sorting, grouping) from its parent that it must satisfy, as well as can have a *derived property* from its children operators. In this work, we optimize *how partition counts are derived* as it is a key factor in cost estimation for massively parallel data systems. Overall, our goal is to improve the cost estimates with minimal changes to the Cascades framework. We next analyze the accuracy of current cost models in SCOPE.

2.4 Cost Model Accuracy

The solid red line in Figure 1 shows that the cost estimates from the default cost model range between an under-estimate of 100× to an over-estimate of 1000×, with a Pearson correlation of just 0.04. As mentioned in the introduction, this is because of the difficulty in modeling the highly complex big data systems. Current cost models rely on hand-crafted heuristics that combine statistics (e.g., cardinality, average row length) in complex ways to estimate each operator’s execution time. These estimates are usually way off and get worse with constantly changing workloads and systems in cloud environments. Big data systems, like SCOPE, further suffer from the widespread use of custom user code that ends up as black boxes in the cost models.

One could consider improving a cost model by considering newer hardware and software configurations, such as machine SKUs, operator implementations, or workload characteristics. SCOPE team did attempt this path and put in significant efforts to improve their default cost model. This alternate cost model is available for SCOPE queries under a flag. We turned this flag on and compared the costs from the improved model with the default one. Figure 1b shows the alternate model in solid blue line. We see that the correlation improves from 0.04 to 0.10 and the ratio curve for the manually improved cost model shifts a bit up, i.e., it reduces the over-estimation. However, it still suffers from the wide gap between the estimated and actual costs, again indicating that cost modeling is non-trivial in these environments.

Finally, as discussed in the introduction and shown as dotted lines in Figure 1(b), fixing cardinalities to the perfect values, i.e., that best that any cardinality estimator [13, 24, 39, 42] could achieve, does not fill the gap between the estimated and the actual costs in SCOPE-like systems.

3 LEARNED COST MODELS

In this section, we describe how we can leverage the common sub-expressions abundant in big data systems to learn a large set of smaller-sized but highly accurate cost models.

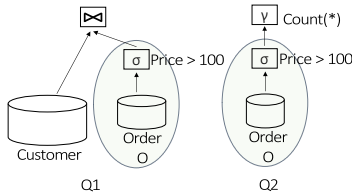


Figure 3: Illustrating common subexpressions.

3.1 Specialized Cost Models

As described in Section 2.2, shared cloud environments often have a large portion of recurring analytical queries (or jobs), i.e., the same business logic is applied to newer instances of the datasets that arrive at regular intervals (e.g., daily or weekly). Due to shared inputs, such recurring jobs often end up having one or more common subexpressions across them. For instance, the SCOPE query processing system at Microsoft has more than 50% of jobs as recurring, with a large fraction of them appearing daily [23], and as high as 60% having common subexpressions between them [21, 22, 42]. Common subexpression patterns have also been reported in other production workloads, including Spark SQL queries in Microsoft’s HDInsight [37], SQL queries from risk control department at Ant Financial Services Group [46], and iterative machine learning workflows [43].

Figure 3 illustrates a common subexpression, consisting of a scan followed by a filter, between two queries. We exploit these common subexpressions by learning a large number of *specialized* models, one for each unique *operator-subgraph* template representing the subexpression. An operator-subgraph template covers the root physical operator (e.g., Filter) and all prior (descendants) operators (e.g., scan) from the root operator of the subexpression. However, parameters and inputs in operator-subgraphs *can vary over time*, and are used as features for the model (along with other logical and physical features) as discussed in Section 3.3.

The operator-subgraph templates essentially *capture the context of root operator*, i.e., *learn the behavior of root physical operator conditioned on the operators beneath it in the query plan*. This is helpful because of two reasons. First, the execution time of an operator depends on whether it is running in a pipelined manner, or is blocked until the completion of underneath operators. For example, the latency of a hash operator running on top of a filter operator is typically smaller compared to when running over a sort operator. Similarly, the grouping and sorting properties of operators beneath the root operator can influence the latency of root operator [19].

Second, the estimation errors (e.g., of cardinality) grow quickly as we move up the query plan, with each intermediate operator building upon the errors of children operators. The operator-subgraph models mitigates this issue partially since the intermediate operators are fixed and the cost of root operator depends only on the leaf level inputs. Moreover, when the estimation errors are systematically off by certain factors (e.g., 10x), the subgraph models can adjust the weights such that the predictions are close to actual

Loss Function	Median Error
Median Absolute Error	246%
Mean Absolute Error	62%
Mean Squared Error	36%
Mean Squared-Log Error	14%

Table 1: Median error using 5-fold CV over the production workload for regression loss functions

(discussed subsequently in Section 3.4). This is similar to adjustments learned explicitly in prior cardinality estimation work [39]. These adjustments generalize well since recurring jobs share similar schemas and the data distributions remain relatively stable, even as the input sizes change over time. Accurate cardinality estimations are, however, still needed in cases where simple adjustment factors do not exist [42].

Next, we discuss the learning settings, feature selection, and our choice of learning algorithm for operator-subgraphs.

3.2 Learning Settings

Target variable. Given an operator-subgraph template, we learn the *exclusive cost* of the root operator as our target. At every intermediate operator, we predict the exclusive cost of the operator conditioned on the subgraph below it. The exclusive cost is then combined with the costs of the children subgraphs to compute the total cost of the sub-graph, similar to how default cost models combine the costs.

Loss function. As the loss function, we use mean-squared log error between the predicted exclusive cost (p) and actual exclusive latency (a) of the operator: $\frac{\sum_n (\log(p+1) - \log(a+1))^2}{n}$, here 1 is added for mathematical convenience. Table 1 compares the average median errors using 5-fold cross validation (CV) of mean-squared log error with other commonly used regression loss functions, using elastic net as the learning model (described subsequently Section 3.4). We note that not taking the log transformation makes learning more sensitive to extremely large differences between actual and predicted costs. However, large differences often occur when the job’s running time itself is long or even due to outlier samples because of machine or network failures (typical in big data systems). We, therefore, minimize the relative error (since $\log(p+1) - \log(a+1) = \log(\frac{p+1}{a+1})$), that reduces the penalty for large differences. Moreover, our chosen loss function helps penalize under-estimation more than over-estimation, since under estimation can lead to under allocation of resources which is typically decided based on cost estimates. Finally, log transformation implicitly ensures that the predicted costs are always positive.

3.3 Feature Selection

It is expensive to extract and combine a large number of features every time we predict the cost of an operator — a typical query plan in big data systems can involve 10s of physical operators, each of which can have multiple possible candidates. Moreover, large feature sets require more number of training samples, while many operator-subgraph instances typically have much fewer samples. Thus, we perform an offline analysis to identify a small set of useful features.

Feature	Description
Input Cardinality (I)	Total Input Cardinality from children operators
Base Cardinality (B)	Total Input Cardinality at the leaf operators
Output Cardinality (C)	Output Cardinality from the current operator
AverageRowLength (L)	Length (in bytes) of each tuple
Number of Partitions (P)	Number of partitions allocated to the operator
Input (IN)	Normalized Inputs (ignored dates, numbers)
Parameters (PM)	Parameters

Table 2: Basic Features

Category	Features
Input or Output data	\sqrt{I} , \sqrt{B} , L^*I , L^*B , $L^*\log(B)$, $L^*\log(I)$, $L^*\log(C)$
Input \times Output	B^*C , I^*C , $B^*\log(C)$, $I^*\log(C)$, $\log(I)^*\log(C)$, I^*C , $\log(B)^*\log(C)$
Input or Output per partition	I/P , C/P , I^*L/P , C^*L/P , $\sqrt{I/P}$, $\sqrt{C/P}$, $\log(I)/P$

Table 3: Derived Features

For selecting features, we start with a set of basic statistics that are frequently used for estimating costs of operators in the default cost model. These include the cardinality, the average row length, and the number of partitions. We consider three kinds of cardinalities: 1) base cardinality: the total input cardinality of the leaf operators in the subgraph, 2) input cardinality: the total input cardinality from the children operators, and 3) output cardinality: the output cardinality of the operator-subgraph. We also consider normalized inputs (ignoring dates and numbers) and parameters to the job that typically vary over time for the recurring jobs.

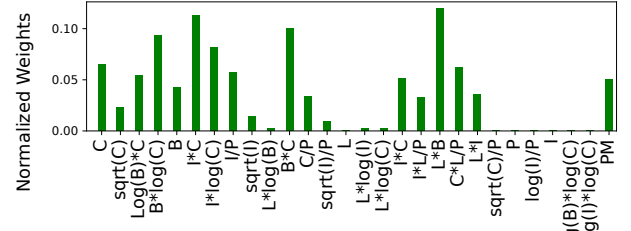
We further combine basic features to create additional derived features to capture the behavior of operator implementations and other heuristics used in default cost models. We start a large space of possible derivations by applying (i) logarithms, square root, squares, and cubes of basic features, (ii) pairwise products among basic features and derived features listed in (i), and (iii) cardinality features divided by the number of partitions (i.e. machine). Given this set of candidate features, we use a variant of elastic net [47] model to select a subset of useful features that have at least one non-zero weight over all subgraph models.

Table 2 and Table 3 depict the selected basic and derived features with non-zero weights. We group the derived features into three categories (i) input or output data, capturing the amount of data read, or written, (ii) the product of input and output, covering the data processing and network communication aspects, and finally (iii) per-machine input or output, capturing the partition size.

Further, we analyze the influence of each feature. While the influence of each feature varies over different subgraph models, Figure 4 shows the aggregated influence over all subgraph models of each feature. Given a total of K non zero features and N subgraph models, with w_{in} as the weight of feature i in model n , we measure the influence of feature i using normalized weight nw_i , as $nw_i = \frac{\sum_N abs(w_{in})}{\sum_K \sum_N abs(w_{kn})}$.

3.4 Choice of learning model

For learning costs over operator-subgraphs, we considered a number of variants of linear-regression, SVM, decision tree and their ensembles, as well as neural network models. On 5-fold cross-validation over our production workload, the following models give more accurate results compared to the

**Figure 4: Features weights (Op-Subgraph model)**

default cost model: (i) Neural network. 3-layers, hidden layer size = 30, solver = adam, activation = relu, l2 regularization = 0.005, (ii) Decision tree: depth = 15, (iii) Random forest number of trees = 20, depth = 5, (iv) FastTree Regression (a variant of Gradient Boosting Tree): number of trees = 20, depth = 5, and (v) Elastic net: $\alpha = 1.0$, fit intercept=True, l1 ratio=0.5.

We observe that the strict structure of subgraph template helps reduce the complexity, making the simpler models, e.g., linear- and decision tree-based regression models, perform reasonably well with the chosen set of features. A large number of operator-subgraph templates have fewer training samples, e.g., more than half of the subgraph instances have < 30 training samples for the workload described in Section 2. In addition, because of the variance in cloud environments (e.g., workload fluctuations, machine failures, etc.), training samples can have noise both in their features (e.g., inaccurate statistics estimates) and the class labels (i.e., execution times of past queries). Together, both these factors lead to over-fitting, making complex models such as neural network as well as ensemble-based models such as gradient-boost perform worse.

Elastic net [47], a L_1 and L_2 regularized linear regression model, on the other hand, is relatively less prone to overfitting. In many cases, the number of candidate features (ranging between 25 to 30) is as many as the number of samples, while only a select few features are usually relevant for a given subgraph. The relevant features further tend to differ across subgraph instances. Elastic net helps perform *automatic feature selection*, by selecting a few relevant predictors for each subgraph independently. Thus, we train all subgraphs with the same set of features, and let elastic net select the relevant ones. Another advantage of elastic net model is that it is intuitive and easily interpretable, like the default cost models which are also weighted sums of a number of statistics. This is an important requirement for effective debugging and analysis of production jobs.

Table 4 depicts the Pearson correlation and median error of the five machine learning models over the production workload. We see that operator-subgraphs models trained using elastic net can make sufficiently accurate cost predictions (14% median error), with a high correlation (more than 0.92) with the actual runtime, a substantial improvement over the default cost model (median error of 258% and Pearson correlation of 0.04). In addition, elastic net models are

fast to invoke during query optimization, and have low storage footprint that indeed makes it feasible for us to learn specialized models for each possible subgraph.

4 ROBUSTNESS

We now discuss how we learn *robust cost models*. As defined in Section 1, robust cost models cover the entire workload with high accuracy for a substantial time period before requiring retraining. We first explain the coverage and accuracy tradeoff for the operator-subgraph model. Then, we discuss the other extreme, namely an *operator* model, and introduce additional models to bridge the gap between the two extremes. Finally, we discuss how we combine predictions from individual models to achieve robustness.

4.1 Accuracy-Coverage Tradeoff

The operator-subgraph model presented in Section 3 is highly specialized. As a result, this model is likely to be highly accurate. Unfortunately, the operator-subgraph model does not cover subgraphs that are not repeated in the training dataset, i.e., it has limited coverage. For example, over 1 day of Microsoft production workloads, operator-subgraphs have learned models for only 54% of the subgraphs. Note that we create a learned model for a subgraph if it has at least 5 occurrences over the single day worth of training data. Thus, it is difficult to predict costs for arbitrary query plans consisting of subgraphs never seen in training dataset.

The other extreme: Operator model. In contrast to the operator-subgraph model, we can learn a model for each physical operator, similar to how traditional query optimizers model the cost. The operator models estimate the execution latency of a query by composing the costs of individual operators in a hierarchical manner akin to how default cost models derive query costs. As a result, operator models can predict the cost of any query in the workload, including those previously unseen in the training dataset. However, similar to traditional cost model, operator models also suffer from poor accuracy since the behavior of an operator changes based on what operations appear below it. Furthermore, the estimation errors in the features or statistics at the lower level operators of a query plan are propagated to the upper level operators that significantly degrades the final prediction accuracy. On 5-fold cross-validation over 1 day of Microsoft production workloads, operator models results in 42% median error and 0.77 Pearson correlation, which although better than the default cost model (258% median error and 0.04 Pearson correlation), is relatively lower compared to that of operator-subgraph models (14% median error and 0.92 Pearson correlation). Thus, there is an accuracy and coverage tradeoff when learning cost models, with operator-subgraph and operator models being the two extremes of this tradeoff.

4.2 Bridging the Gap

We now present additional models that fall between the two extreme models in terms of the accuracy-coverage trade-off.

Model	Correlation	Median Error
Default	0.04	258%
Neural Network	0.89	27%
Decision Tree	0.91	19 %
Fast-Tree regression	0.90	20%
Random Forest	0.89	32%
Elastic net	0.92	14%

Table 4: Correlation and error w.r.t. actual runtime for the operator-subgraphs

Operator-input model. An improvement to per-operator is to learn a model for all jobs that share similar inputs. Similar inputs also tend to have similar schema and similar data distribution even as the size of the data changes over time, thus operator models learned over similar inputs often generalize over future job invocations. In particular, we create a model for each operator and input template combination. An input template is a normalized input where we ignore the dates, numbers, and parts of names that change over time for the same recurring input, thereby allowing grouping of jobs that run on the same input schema over different sessions. Further, to partially capture the context, we featurize the intermediate subgraph by introducing two additional features: 1) the number of logical operator in the subgraph (CL) and 2) the depth of the physical operator in the sub-graph (D). This helps in distinguishing subgraph instances that are extremely different from each other.

Operator-subgraphApprox model. While operator-subgraph models exploit the overlapping and recurring nature of big data analytics, there is also a large number (about 15-20%) of subgraphs that are similar but are not exactly the same. To bridge this gap, we relax the notion of subgraph similarity, and learn one model for all subgraphs that have the same inputs and the same *approximate* underlying subgraph. We consider two subgraphs to be approximately same if they have the same physical operator at the root, and consist of the *same frequency of each logical operator in the underneath subgraph (ignoring the ordering between operators)*. This relaxed similarity criteria allows grouping of similar subgraphs without substantially reducing the coverage. Overall, operator-subgraphApprox model is a hybrid of the operator-subgraph and operator-input models: it achieves much higher accuracy compared to operator or operator-input models, and more coverage compared to operator-subgraph model.

Table 5 depicts the Pearson correlation, the median accuracy using 5-fold cross-validation, as well as the coverage of individual cost models using elastic net over production workloads. As we move from more specialized to more generalized models (i.e., operator-subgraph to operator-subgraphApprox to operator-input to operator), we see that the model accuracy decreases while the coverage over the workload increases. We provide more details on the importance of features for each model in our technical report [12].

4.3 The Combined Model

Given multiple learned models, with varying accuracy and coverage, the strawman approach is to select a learned model

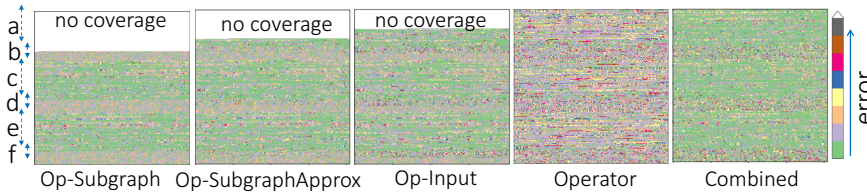


Figure 5: Heatmap of errors over 42K operators from production jobs. Each point in the heatmap depicts the error of learned cost with respect to the actual runtime

Model	Correlation	Median Error
Default	0.04	258%
Neural Network	0.79	31%
Decision Tree	0.73	41 %
FastTree Regression	0.84	19%
Random Forest	0.80	28%
Elastic net	0.68	64%

Table 6: Correlation and error w.r.t actual runtimes for the Combined Model

in decreasing order of their accuracy over training data, starting from operator-subgraphs to operator-subgraphApprox to operator-inputs to operators. However, as discussed earlier, there are subgraph instances where more accurate models result in poor performance over test data due to over-fitting.

To illustrate, Figure 5 depicts the heat-map representation of the accuracy and coverage of different individual models, over more than 42K operator instances from our production workloads. Each point in the heat-map represents the error of the predicted cost with respect to the actual runtime: the more green the point, the less the error. The empty region at the top of the chart depicts that the learned model does not cover those subgraph instances. We can see that the predictions are mostly accurate for operator-subgraph models, while Operator models have relatively more error (i.e., less green) than the operator-subgraph models. Operator-input have more coverage with marginally more error (less green) compared to operator-subgraphs. However, for regions b, d, and f, as marked on the left side of the figure, we notice that operator-input performs better (more blue and green) than operator-subgraph. This is because operators in those regions have much fewer training samples that results in over-fitting. Operator-input, on the other hand, has more training samples, and thus performs better. Thus, it is difficult to decide a rule or threshold that always select the best performing model for a given subgraph instance.

Learning a meta-ensemble model. We introduce a meta-ensemble model that uses the predictions from specialized models as meta features, along with the following extra features: (i) cardinalities (I, B, C), (ii) cardinalities per partition (I/P, B/P, C/P), and (iii) number of partitions (P) to output a more accurate cost. Table 6 depicts the performance of different machine learning models that we use as a meta-learner on our production workloads. We see that the FastTree regression [14] results in the most accurate predictions. FastTree regression is a variant of the gradient boosted regression trees [16] that uses an efficient implementation of

Model	Correlation	Median Error	Coverage
Default	0.04	258%	100%
Op-Subgraph	0.92	14%	54%
Op-Subgraph Approx	0.89	16 %	76%
Op-Input	0.85	18%	83%
Operator	0.77	42%	100%
Combined	0.84	19%	100%

Table 5: Performance of learned models w.r.t to actual runtimes

the MART gradient boosting algorithm [33]. It builds a series of regression trees (estimators), with each successive tree fitting on the residual of trees that precede it. Using 5-fold cross-validation, we find that the maximum of only 20 regression trees with mean-squared log error as the loss function and the sub-sampling rate of 0.9 is sufficient for optimal performance. As depicted in Figure 5, using FastTree Regression as a meta-learner has three key advantages.

First, FastTree regression can effectively characterize the space where each model performs well. The regression trees recursively split the space defined by the predictions from individual models and features, creating fine-grained partitions such that prediction in each partition are highly accurate.

Second, FastTree regression performs better for operator instances where individual models perform worse. For example, some of the red dots found in outlier regions b, d and f of the individual model heat-maps are missing in the combined model. This is because FastTree regression makes use of sub-sampling to build each of the regression trees in the ensemble, making it more resilient to overfitting and noise in execution times of prior queries. Similarly, for region a where subgraph predictions are missing, FastTree regression creates an extremely large number of fine-grained splits using extra features and operator predictions to give even better accuracy compared to Operator models.

Finally, the combined model is naturally capable of covering all possible plans, since it uses Operator model as one of the predictors. Further, as depicted in Figure 5, the combined model is comparable to that of specialized models, and almost always has better accuracy than Operator model. The combined model is also flexible enough to incorporate additional models or features, or replace one model with another. However, on also including the default cost model, it did not result in any improvement on SCOPE. Next, we discuss how we integrate the learned models within the query optimizer.

5 OPTIMIZER INTEGRATION

In this section, we describe our two-fold integration of CLEO within SCOPE. First, we discuss our end-to-end feedback loop to learn cost models and generate predictions during optimization. Then, we discuss how we extend the optimizer for supporting resource-exploration using learned models.

5.1 Integrating Learned Cost Models

The integration of learned models with the query optimizer involves the following three components.

Instrumentation and Logging. Big data systems, such as SCOPE, are already instrumented to collect logs of query

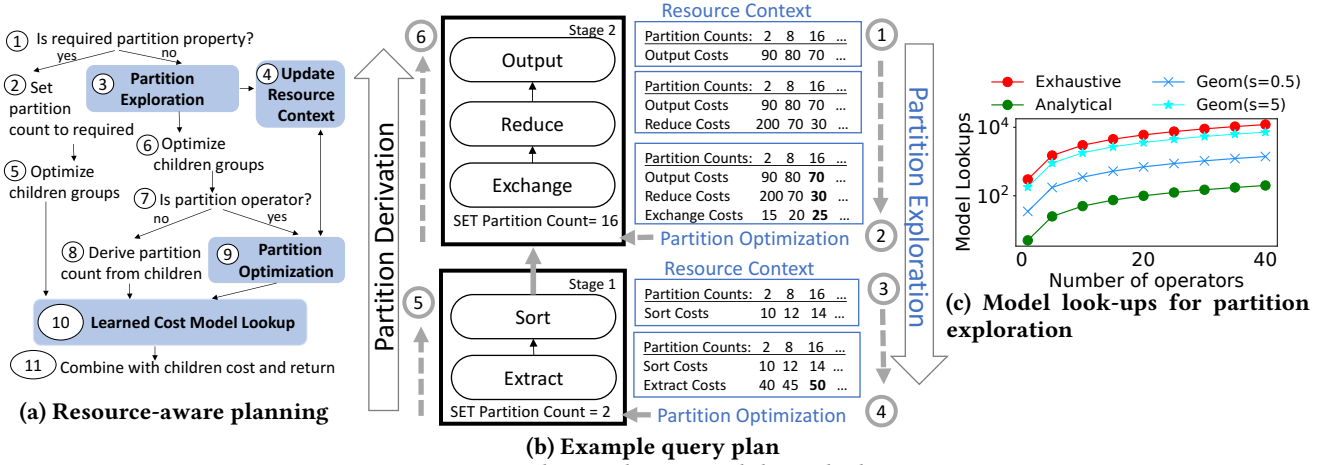


Figure 6: Integrating learned cost models with the query optimizer

plan statistics such as cardinalities, estimated costs, as well as runtime traces for analysis and debugging. For uniquely identifying each operator and the subplan, query optimizers annotate each operator with a *signature* [7], a 64-bit hash value that can be recursively computed in a bottom-up fashion by combining (i) the signatures of children operators, (ii) hash of current operator’s name, and (iii) hash of operator’s logical properties. We extend the optimizer to compute three more signatures, one for each individual sub-graph mode, and extract additional statistics (i.e., features) that were missing for CLEO. Since all signatures can be computed simultaneously in the same recursion, and that there are only 25 – 30 features (most of which were already extracted), the additional overhead is minimal ($\leq 10\%$), as we describe in Section 6. This overhead includes both the logging and the model lookup that we discuss subsequently.

Training and Feedback. Given the logs of past query runs, we learn each of the four individual elastic net-based models independently and in parallel using our SCOPE-based parallel model trainer. Experimentally, we found that a training window of two days and a training frequency of every ten days results in acceptable accuracy and coverage (Section 6). We then use the predictions from individual models on the next day queries to learn the combined FastTree regression model. Since individual models can be trained in parallel, the training time is not much, e.g., it takes less than 45 minutes for training 25K models from over 50,000 jobs using a cluster of 200 nodes. Once trained, we serialize the models and feedback them to the optimizer. The models can be served either from a text file, using an additional compiler flag, or using a web service that is backed by a SQL database.

Look-up. All models relevant for a cluster are loaded upfront by the optimizer, into a hash map with keys as *signatures* of models, to avoid expensive lookup calls during optimization. When loaded simultaneously, all 25K models together take about 600 MB of memory, which is within an acceptable range. Finally, for cost estimation, we modify the Optimize Input phase of Cascade optimizer to invoke learned models. Figure 6a highlights the the key steps that we added in blue.

Essentially, we replace the calls to the default cost-models with the learned model invocations (Step 10 in Figure 6a) to predict the exclusive cost of an operator, which is then combined with costs of children operators, similar to how default cost models combine costs. Moreover, since there is an operator model and a combined model for every physical operator, CLEO can cover all possible query plans, and can even generate a plan unseen in training data. All the features that learned models need are available during query optimization. However, we do not reuse the partition count derived by the default cost model, rather we try to find a more optimal partition count (steps 3, 4, 9) as it drives the query latency. We discuss the problem with existing partition count selection and our solution below.

5.2 Resource-aware Query Planning

The degree of parallelism (i.e., the number of machines or containers allocated for each operator) is a key factor in determining the runtime of queries in massively parallel databases [41], which implicitly depends on the partition count. This makes partition count as an important feature in determining the cost of an operator (as noted in Figures 4).

Unfortunately, in existing optimizers, the partition count is not explored for all operators, rather partitioning operators (e.g., Exchange for stage 2 in Figure 6b) set the partition count for the entire stage based on their local statistics [44]. The above operators on the same stage simply derive the partition count set by the partitioning operator. For example, in stage 2 of Figure 6b, Exchange sets a partition count to 2 as it results in its smallest local cost (i.e., 15). The operators above Exchange (i.e., Reduce and Output) derive the same partition count, resulting in a total cost of 305 for the entire stage. However, we can see that the partition count of 16 results in a much lower overall cost of 125, though it’s not locally optimal for Exchange. Thus, *not optimizing the partition count for the entire stage results in a sub-optimal plan.*

To address this, we explore the partition counts during query planning, by making query planning resource-aware. Figures 6a and 6b illustrate our resource-aware query planning approach. We introduce the notion of a *resource-context*,

within the optimizer-context, for tracking costs of partitions across operators in a stage. Furthermore, we add a *partition exploration* step, where each physical operator attaches a list of learned costs for different partition counts to the resource context (step 3 in Figure 6a). For example, in Figure 6b, the resource-context for stage 2 shows the learned cost for different partition counts for each of the three operators. On reaching the stage boundary, the partitioning operator Exchange performs *partition optimization* (step 9 in Figure 6a) to set its local partition count to 16, that results in the lowest total cost of 125 across for the stage. Thereafter, the higher level operators simply derive the selected partition count (line 8 in Figure 6a), like in the standard query planning, and estimate their local costs using learned models. Note that when a partition comes as required property [19] from upstream operators, we set the partition count to the required value without any exploration (Figure 6a step 2).

SCOPE currently does not allow varying other resources, therefore we focus only on partition counts in this work. However, the resource-aware query planning with the three new abstractions to Cascades framework, namely the resource-context, partition-exploration, and partition-optimization, is general enough to incorporate additional resources such as memory sizes, number of cores, VM instance types, and other infrastructure level decisions to jointly optimize for both plan and resources. Moreover, our proposed extensions can also be applied to other big data systems such as Spark, Flink, and Calcite that use variants of Cascades optimizers and follow a similar top-down query optimization as SCOPE.

Our experiments in Section 6 show that the resource-aware query planning not only generates better plans in terms of latency, but also leads to resource savings. However, the challenge is that estimating the cost for every single partition count for each operator in the query plan can explode the search space and make query planning infeasible. We discuss our approach to address this next.

5.3 Efficient Resource Exploration

We now discuss two techniques for efficiently exploring the partition counts in CLEO, without exploding the search space.

Sampling-based approach. Instead of considering every single partition count, one option is to consider a uniform sample over the set of all possible containers for the tenant or the cluster. However, the relative change in partition is more interesting when considering its influence on the cost, e.g., a change from 1 to 2 partitions influences the cost more than a change from 1200 to 1210 partitions. Thus, we sample partition counts in a geometrically increasing sequence, where a sample x_{i+1} is derived from previous sample x_i using: $x_{i+1} = \lceil x_i + x_i/s \rceil$ with $x_0 = 1$, $x_1 = 2$. Here, s is a skipping coefficient that decides the gap between successive samples. A large s leads to a large number of samples and more accurate predictions, but at the cost of higher model look-ups and prediction time.

Analytical approach. We reuse the individual learned models to directly model the relationship between the partition

Cluster	Day	Total Jobs	Recurring Jobs	Recurring Templates	Total Sub-Expr.	Common Sub-Expr.	Recurring Sub-Expr.	Ad-hoc Sub-Expr.
Cluster1	Day1	64796	52400	17662	3546087	2874485	484909	186693
	Day2	62065	50055	16888	3284706	2633981	468643	182082
	Day3	68479	55767	18143	3657239	2946294	493745	217200
Cluster2	Day1	52952	49452	6741	2523173	1948962	488531	85680
	Day2	36776	33167	6044	1595925	1238273	285043	72609
	Day3	40464	36429	6867	2232402	1624081	522783	85538
Cluster3	Day1	30277	25929	6273	1337034	1031231	238666	67137
	Day2	29120	25002	5917	1292181	1017849	211155	63177
	Day3	29667	25306	6281	1326588	1023611	235668	67309
Cluster4	Day1	14562	13145	2277	530045	416972	61492	51581
	Day2	16000	14132	2696	519914	408397	53819	57698
	Day3	18641	17040	2606	536294	420282	49434	66578
Overall		463799	397824	98395	22381588	17584418	3593888	1203282

Figure 7: Workload consisting of almost 0.5 million jobs from 4 different production clusters over 3 days.

count and the cost of an operator. The key insight here is that only features where partition count is present are relevant for partition exploration, while the rest of the features can be considered constants since their values are fixed during partition exploration. Thus, we can express operator cost as follows: $cost \propto \frac{(\theta_1 * I + \theta_2 * C + \theta_3 * I * C)}{P} + \theta_c * P$ where I , C , and P refer to input cardinality, cardinality and partition count respectively. During optimization, we know I , C , $I * C$, therefore: $cost \propto \frac{\theta_P}{P} + \theta_c * P$. Extending the above relationship across all operators (say n in number) in a stage, the relationship can be modeled as: $cost \propto \frac{\sum_{i=1}^n \theta_{P_i}}{P} + \sum_{i=1}^n \theta_{C_i} * P$

Thus, during partition exploration, each operator calculates θ_P and θ_C and adds them to resource-context, and the partitioning operator selects the optimal partition count by optimizing the above function. There are three possible scenarios: (i) $\sum_{i=1}^n \theta_{P_i}$ is positive while $\sum_{i=1}^n \theta_{C_i}$ is negative: we can have the maximum number of partitions for the stage since there is no overhead of increasing the number of partitions, (ii) $\sum_{i=1}^n \theta_{P_i}$ is negative while $\sum_{i=1}^n \theta_{C_i}$ is positive: we set the partition count to minimum as increasing the partition count increases the cost, (iii) $\sum_{i=1}^n \theta_{P_i}$ and $\sum_{i=1}^n \theta_{C_i}$ are either both positive or both negative: we can derive the optimal partition count by differentiating the cost equation with respect to P . Overall, for m physical operator, and the maximum possible partition count of P_{max} , the analytical model makes $5 \cdot m \cdot \log_{\frac{s+1}{s}} P_{max}$ cost model look-ups. Figure 6c shows the number of model look-ups for sampling and analytical approaches as we increase the number of physical operators from 1 to 40 in a plan. While the analytical model incurs a maximum of 200 look-ups, the sampling approach can incur several thousands depending on the skipping coefficients.

6 EXPERIMENTS

In this section, we present an evaluation of our learned optimizer CLEO. For fairness, we feed the same statistics (e.g., cardinality, average row length) to learned models that are used by the SCOPE default cost model. Our goals are five fold: (i) to compare the prediction accuracy of our learned cost models over all jobs as well as over only ad-hoc jobs across

	All jobs				Ad-hoc jobs			
	Correlation	Median Error	95%tile Error	Coverage	Correlation	Median Error	95%tile Error	Coverage
Default	0.12	182%	12512%	100%	0.09	204%	17791%	100%
Op-Subgraph	0.86	9%	56%	65%	0.81	14%	57%	36%
Op-Subgraph Approx	0.85	12 %	71%	82%	0.80	16 %	79%	64%
Op-Input	0.81	23%	90%	91%	0.77	26%	103%	79%
Operator	0.76	33%	138%	100%	0.73	42%	186%	100%
Combined	0.79	21%	112%	100%	0.73	29%	134%	100%

Table 7: Breakdown of accuracy and coverage of each learned model for all jobs and ad-hoc jobs separately on Cluster1.

Cluster	Default (all jobs)		Learned (all jobs)		Learned (ad-hoc jobs)	
	Correlation	Median Accuracy	Correlation	Median Accuracy	Correlation	Median Accuracy
Cluster 1	0.12	182%	0.79	21%	0.73	29%
Cluster 2	0.08	256%	0.77	33%	0.75	40%
Cluster 3	0.15	165%	0.83	26%	0.81	38%
Cluster 4	0.05	153%	0.74	15%	0.72	26%

Table 8: Pearson Correlation and Median accuracy of default and combined learned model over all jobs and ad-hoc jobs on each cluster.

multiple clusters, (ii) to test the coverage and accuracy of learned cost models over varying test windows, (iii) to compare the CLEO cost estimates with those from CardLearner, (iv) to explore why perfect cardinality estimates are not sufficient for query optimization, (iv) to evaluate the effectiveness of sampling strategies and the analytical approach proposed in Section 5.2 in finding the optimal resource (i.e., partition count), and (v) to analyze the performance of plans produced by CLEO with those generated from the default optimizer in CLEO using both the production workloads and the TPC-H benchmark (v) to understand the training and runtime overheads when using learned cost models.

Workload. As summarized in Figure 7, we consider a large workload trace from 4 different production clusters comprising of 423 virtual clusters, with each virtual cluster roughly representing a business unit within Microsoft, and consisting a total of ≈ 0.5 million jobs over 3 days, that ran with a total processing time of ≈ 6 million hours, and use a total of ≈ 1.4 billion containers. The workload exhibits variations in terms of the load (e.g., more than $3\times$ jobs on Cluster1 compared to Cluster4), as well as in terms of job properties such as average number of operators per job (50s in Cluster1 compared to 30s in Cluster4) and the average total processing time of all containers per job (around 17 hours in Cluster1 compared to around 5 hours in Cluster4). The workload also varies across days from a 30% decrease to a 20% increase of different job characteristics on different clusters. Finally, the workload consists of a mix of both recurring and ad-hoc jobs, with about 7% – 20% ad-hoc jobs on different clusters and different days.

6.1 Accuracy

We first compare the accuracy and correlation of learned models to that of the default cost model for each of the clusters. We use the elastic net model for individual learned model and the Fast-Tree regression for the combined model. Cross-validation results for other ML models are presented in Tables 6, 7 and 8, and the CDF distribution of predictions are shown in our technical report [12]. We learn the models on day 1 and day 2, and predict on day 3 of the workload

as described in Section 5.1. Table 8 shows the Pearson correlation and median accuracy of default cost model and the combined model for all jobs and only ad-hoc jobs separately for each of the clusters on day 3. We further show the breakdown of results for each of the individual models on cluster 1 jobs in Table 7.

All jobs. We observe that learned models result between $8\times$ to $10\times$ better accuracy and between $6\times$ to $14\times$ better correlation compared to the default cost model across the 4 clusters. For operator-subgraph, the performance is highly accurate (9% median error and .86 correlation), but at lower coverage of 65%. This is because there are a large number of specialized models, highly optimized for specific sub-graph instances. As the coverage increases from operator-subgraphApprox to operator-input to operator, the accuracy decreases. Overall, the combined model is able to provide the best of both worlds, i.e., accuracy close to those of individual models and 100% coverage like that of the operator model. The 50th and the 95th percentile errors of the combined model are about $10\times$ and $1000\times$ better than the default SCOPE cost model. These results show that it is possible to learn highly accurate cost models from the query workload.

Only ad-hoc Jobs. Interestingly, *the accuracy over ad-hoc jobs drop slightly but it is still close to those over all jobs (Table 8 and Table 7)*. This is because: (i) ad-hoc jobs can still have one or more similar subexpressions as other jobs (e.g., they might be scanning and filtering the same input before doing completely new aggregates), which helps them to leverage the subgraph models learned from other jobs. This can be seen in Table 7 – the sub-graph learned models still have substantial coverage of subexpression on ad-hoc jobs. For example, the coverage of 64% for Op-Subgraph Approx model means that 64% of the sub-expressions on ad-hoc jobs had matching Op-Subgraph Approx learned model. (ii) Both the operator and the combined models are learned on a per-operator basis, and have much lower error (42% and 29%) than the default model (182%). This is because the query processor, the operator implementations, etc. still remain the same, and their behavior gets captured in the learned cost models. Thus, even if there is no matching sub-expression in ad-hoc jobs, the operator and the combined models still result in better prediction accuracy.

6.2 Robustness

We now look at the robustness (as defined in Section 1) of learned models in terms of accuracy, coverage, and retention over a month long test window on cluster 1 workloads.

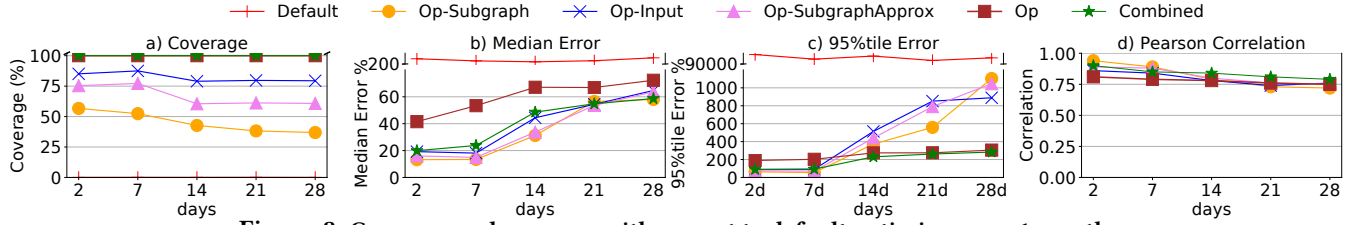


Figure 8: Coverage and accuracy with respect to default optimizer over 1 month

Coverage over varying test window. Figure 8a depicts the coverage of different subgraph models as we vary the test window over a duration of 1 month. The coverage of per-operator and combined model is always 100% since there is one model for every physical operator. The coverage of per-subgraph models, strictest among all, is about 58% after 2 days, and decreases to 37% after 28 days. Similarly, the coverage of per-subgraphApprox ranges between 75% and 60%. The per-operator-input model, on the other hand remains stable between 78% and 84%.

Error and correlation over varying test window. Figures 8b and 8c depict the median and 95%ile error percentages respectively over a duration of one month. While the median error percentage of learned models improves on default cost model predictions by 3-15x, the 95%ile error percentage is better by over three orders of magnitude. For specialized learned models, the error increases slowly over the first two weeks and then grows much faster due to decrease in the coverage. Moreover, the 95%ile error of the subgraph models grows worse than their median error. Similarly, in Figure 8d, we see that predicted costs from learned models are much more correlated with the actual runtimes, having high Pearson correlation (generally between 0.70 and 0.96), compared to default cost models that have a very small correlation (around 0.1). A high correlation over a duration month of 1 month shows that learned models can better discriminate between two candidate physical operators. Overall, based on these results, we believe re-training every 10 days should be acceptable, with a median error of about 20%, 95% error of about 200%, and Pearson correlation of around 0.80.

Robustness of the combined model. From Figure 8, we note that the combined model (i) has 100% coverage, (ii) matches the accuracy of best performing individual model at any time (visually illustrated in Figure 5), while having a high correlation (> 0.80) with actual runtimes, and (iii) gives relatively stable performance with graceful degradation in accuracy over longer run. Together, these results show that the combined model in CLEO is indeed robust.

6.3 Impact of Cardinality

Comparison with Cardlearner. Next, we compare the accuracy and the Pearson correlation of CLEO and the default cost model with CardLearner [42], a learning-based cardinality estimation that employs a Poisson regression model to improve the cardinality estimates, but uses the default cost model to predict the cost. For comparison, we considered jobs from a single virtual cluster from cluster 4, consisting of

about 900 jobs. While CLEO and the default cost model use the cardinality estimates from the SCOPE optimizer, we additionally consider a variant (CLEO + CardLearner) where CLEO uses the cardinality estimates from CardLearner. Overall, we observe that the median error of default cost with CardLearner (211%) is better than that of default cost model alone (236%) but much still worse than that of CLEO's (18%) and CLEO + CardLearner (13%). Figure 9 depicts the CDF of estimated and actual costs ratio, where we can see that by learning costs, CLEO significantly reduces both the under-estimation as well as over-estimation, while CardLearner only marginally improves the accuracy for both default cost model and CLEO. Similarly, the Pearson correlation of CardLearner's estimates (.01) is much worse than that of CLEO's (0.84) and CLEO + CardLearner (0.86). These results are consistent with our findings from Section 2 where we show that fixing the cardinality estimates is not sufficient for filling the wide gap between the estimated and the actual costs in SCOPE-like systems. Interestingly, we also observe that the Pearson correlation of CardLearner is marginally less than that of the default cost model (0.04) in spite of better accuracy, which can happen when a model makes both over and under estimation over different instances of the same operator.

Why cardinality alone is not sufficient? To understand why fixing cardinality alone is not enough, we performed the following two micro-experiments on a subset of jobs, consisting of about 200 K subexpressions from cluster 4.

1. The need for a large number of features. First, starting with perfect cardinality estimates (both input and output) as only features, we fitted an elastic net model to find / tune the best weights that lead to minimum cost estimation error (using the loss function described in Section 3.4). Then, we incrementally added more and more features and also combined feature with previously selected features, retraining the model with every addition. Figure 12 shows the decrease in cost model error as we cumulatively add features from left to right. We use the same notations for features as defined in Table 2 and Table 3.

We make the following observations.

- (1) When using only perfect cardinalities, the median error is about 110%. However, when adding more features, we observe that the error drops by more than half to about 40%. This is because of the more complex environments and queries in big data processing that are very hard to cost with just the cardinalities (as also discussed in Section 2.3).
- (2) Deriving additional statistics by transforming (e.g., sqrt, log, squares) input and output cardinalities along with other

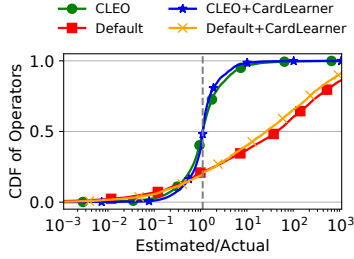


Figure 9: Comparison with CardLearner

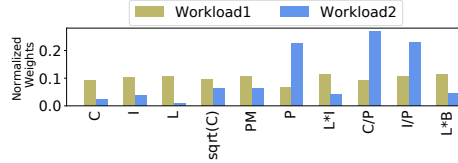


Figure 10: Hash join operator having different weights over different sets of sub-expressions.

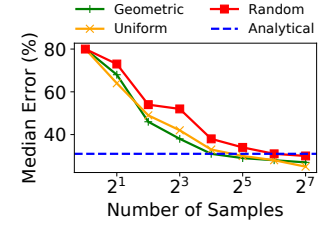


Figure 11: Partition Exploration Accuracy vs. Efficiency

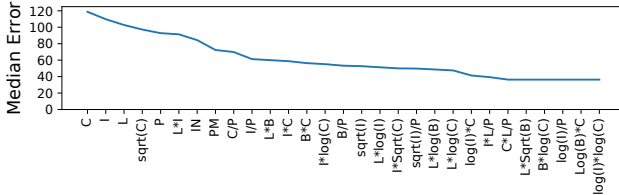


Figure 12: Impact on median error as we cumulatively add features from left to right. The first two features are perfect output (C) and input (I) cardinalities.

features helps in reducing the error. However, it is hard to arrive at these transformations in hand coded cost models.

(3) We also see that features such as parameter (PM), input (IN), and partitions (P) are quite important, leading to sharp drops in error. While partition count is good indicator of degree of parallelism (DOP) and hence the runtime of job; unfortunately query optimizers typically use a fixed partition count for estimating cost as discussed in Section 5.2.

(4) Finally, there is a pervasive use of unstructured data (with schema imposed at runtime) as well as custom user code (e.g., UDFs) that embed arbitrary application logic in big data applications. It is hard to come up with generic heuristics, using just the cardinality, that effectively model the runtime behavior of all possible inputs and UDFs.

2. Varying optimal weights. We now compare the optimal weights of features of physical operators when they occur in different kinds of sub-expressions. We consider the popular hash join operator and identified two sets of sub-expressions in the same sample dataset: (i) hash join appears on top of two scan operators, and (ii) hash join appears on top of two other join operators, which in turn read from two scans.

Figure 10 shows the weights of the top 10 features of the hash join cost model for the two sets. We see that the partition count is more influential for set 2 compared to set 1. This is because there is more network transfer of data in set 2 than in set 1 because of two extra joins. Setting the partition count that leads to minimum network transfer is therefore important. On the other hand, for jobs in set 1, we notice that hash join typically sets the partition count to the same value as that of inputs, since that minimizes repartitioning. Thus, partition count is less important in set 1. To summarize, even when cardinality is present as a raw or derived feature, its relative importance is instance specific (heavily influenced by the partitioning and the number of machines) and hard to capture in a static cost model.

6.4 Efficacy of Partition Exploration

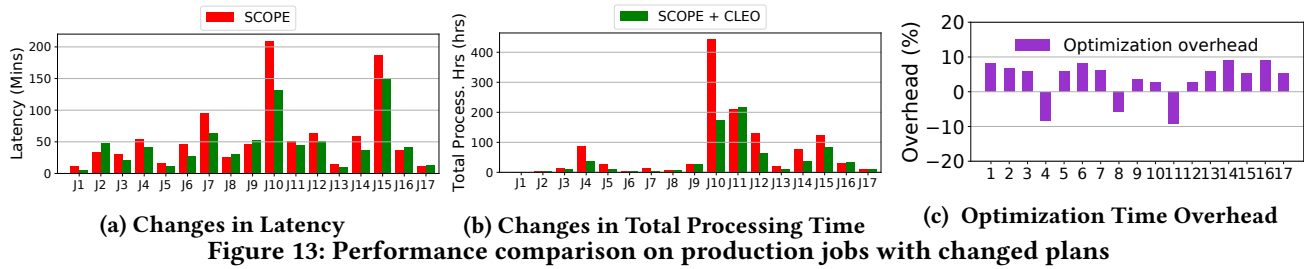
In this section, we explore the effectiveness of partition exploration strategies proposed in Section 5.2 in selecting the partition count that leads to lowest cost for learned models. We used a subset of 200 sub-expression instances from the production workload from cluster 1, and exhaustively probe the learned models for all partition counts from 0 to 3000 (the maximum capacity of machines on a virtual cluster) to find the most optimal cost. Figure 11 depicts the median error in cost accuracy with respect to the optimal cost for (i) three sampling strategies: random, uniform, and geometric as we vary the number of samples of partition counts (ii) the analytical model (dotted blue line) that selects a single partition count.

We note that the analytical model, although approximate, gives more accurate results compared to sampling-based approaches until the sample size of about 15 to 20, thereby requiring much fewer model invocations. Further, for a sample size between 4 to 20, we see that the geometrically increasing sampling strategy leads to more accurate results compared to uniform and random approaches. This is because it picks more samples when the values are smaller, where costs tend to change more strongly compared to higher values. During query optimization, each sample leads to five learned cost model predictions, four for individual models and one for the combined model. Thus, for a typical plan in big data system that consists of 10 operators, the sampling approach requires $20 * 5 * 10 = 1000$ model invocations, while the analytical approach requires only $5 * 10 = 50$ invocations for achieving the same accuracy. This shows that the analytical approach is practically more effective when we consider both efficiency and accuracy together. Hence, CLEO uses the analytical approach as the default partition exploration strategy.

6.5 Performance

We split the performance evaluation of CLEO into three parts: (i) the runtime performance over production workloads, (ii) a case-study on the TPC-H benchmark [34], explaining in detail the plan changes caused by the learned cost models, and (iii) the overheads incurred due to the learned models. For these evaluations, we deployed a new version of the SCOPE runtime with CLEO optimizer (i.e., SCOPE+CLEO) on the production clusters and compared the performance of this runtime with the default production SCOPE runtime. We used the analytical approach for partition exploration.

6.5.1 Production workload. Since production resources are expensive, we selected a subset of the jobs, similar to prior



work on big data systems [42], as follows. We first recompiled all the jobs from a single virtual cluster from cluster 4 with CLEO and found that 182 (22%) out of 845 jobs had plan changes when partition exploration was turned off, while 322 (39%) had plan changes when we also applied partition exploration. For execution, we selected jobs that had at least one change in physical operator implementations, e.g., for aggregation (hash vs stream grouping), join (hash vs merge), addition or removal of grouping, sort or exchange operators. We picked 17 such jobs and executed them with and without CLEO over the same production data, while redirecting the output to a dummy location, similar to prior work [1].

Figure 13a shows the end-to-end latency for each of the jobs, compared to their latency when using the default cost model. We see that the learned cost models improve latency in 70% (12 jobs) cases, while they degrade latency in the remaining 30% cases. Overall, the average improvement across all jobs is 15.35%, while the cumulative latency of all jobs improves by 21.3%. Interestingly, CLEO was able to improve the end-to-end latency for 10 out of 12 jobs with less degree of parallelism (i.e., the number of partitions). This is in contrary to the typical strategy of scaling out the processing by increasing the degree of parallelism, which does not always help. Instead, resource-awareness plan selection can reveal more optimal combinations of plan and resources. To understand the impact on resource consumption, Figure 13b shows the total processing time (CPU hours). Learned cost models reduce the Overall, we see the total processing time reducing by 32.2% on average and 40.4% cumulatively across all 17 jobs— a significant operational cost saving in big data computing environments.

Thus, the learned cost models could reduce the overall costs while still improving latencies in most cases. Below, we dig deeper into the plans changes using the TPC-H workload.

6.5.2 TPC-H workload. We generated TPC-H dataset with a scale factor of 1000, i.e., total input of 1TB. We ran all 22 queries 10 times, each time with randomly chosen different parameters, to generate the training dataset. We then trained our cost models on this workload and feedback the learned models to re-run all 22 queries. We compare the performance with- and without the feedback as depicted in Figure 14. For each observation, we take the average of 3 runs. Overall, 6 TPC-H queries had plan changes when using resource-aware cost model. Out of these, 4 changes (Q8, Q9, Q16, Q20) improve latency as well as total processing time, 1 change improves only the latency (Q11), and 1 change (Q17) leads

to performance regression in both. We next discuss the key changes observed in the plans.

1. More optimal partitioning. In Q8, for Part (100 partitions) $\bowtie_{partKey}$ Lineitem (200 partitions) and in Q9 for Part (100 partitions) $\bowtie_{partKey}$ (Lineitem \bowtie Supplier) (250 partitions), the default optimizer performs the merge join over 250 partitions, thereby re-partitioning both Part and the other join input into 250 partitions over the Part key column. Learned cost models, on other hand, perform the merge join over 100 partitions, which requires re-partitioning only the other join inputs and not the Part table. Furthermore, re-partitioning 200 or 250 partitions into 100 partitions is cheaper, since it involves partial merge [9], compared to re-partitioning them into 250 partitions. In Q16, for the final aggregation and top-k selection, both the learned and the default optimizer repartition the 250 partitions from the previous operator’s output on the final aggregation key. While the default optimizer re-partitions into 128 partitions, the learned cost models pick 100 partitions. The aggregation cost has almost negligible change due to change in partition counts. However, repartitioning from 250 to 100 turns out to be substantially faster than re-partitioning from 250 to 128 partitions.

2. Skipping exchange (shuffle) operators. In Q8, for Part (100 partitions) $\bowtie_{partKey}$ Lineitem (200 partitions), the learned cost model performs join over 100 partitions and thus it skips the costly Exchange operator over the Part table. On the other hand, the default optimizer creates two Exchange operator for partitioning each input into 250 partitions.

3. More optimal physical operator: For both Q8 and Q20, the learned cost model performs join between Nations and Supplier using merge join instead of hash join (chosen by default optimizer). This results in an improvement of 10% to 15% in the end-to-end latency, and 5% to 8% in the total processing time (cpu hours).

4. Regression due to partial aggregation. For Q17, the learned cost models add local aggregation before the global one to reduce data transfer. However, this degrades the latency by 10% and total processing time by 25% as it does not help in reducing data. Currently, learned models do not learn from their own execution traces. We believe that doing so can potentially resolve some of the regressions.

6.5.3 Training and Runtime Overheads. We now describe the training and run-time overheads of CLEO. It takes less than 1 hour to analyze and learn models for a cluster running about 800 jobs a day, and less than 4 hours for training over 50K jobs instances at Microsoft. We use a parallel model

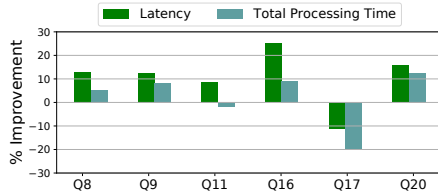


Figure 14: Performance change with SCOPE + CLEO for TPC-H queries (higher is better)

trainer that leverages SCOPE to train and validate models independently and in parallel, which significantly speeds up the training process. For a single cluster of about 800 jobs, CLEO learns about 23K models which when simultaneously loaded takes about 600 MB of memory. About 80% of the memory is taken by the individual models while the remaining is used by the combined model. The additional memory usage is not an issue for big data environments, where an optimizer can typically have 100s of GBs of memory. Finally, we saw between 5-10% increase in the optimization time for most of the jobs when using learned cost models, which involves the overhead of signature computation of sub-graph models, invoking learned models, as well as any changes in plan exploration strategy due to learned models. Figure 13c depicts the overhead in the optimization time for each of the 17 production jobs we executed. Since the optimization time is often in orders of few hundreds of milliseconds, the overhead incurred here is negligible compared to the overall compilation time (orders of seconds) of jobs in big data systems as well as the potential gains we can achieve in the end-to-end latency (orders of 10s of minutes).

7 RELATED WORK

Machine learning has been used for estimating the query execution time of a given physical plan in centralized databases [2, 17, 29]. In particular, the operator and sub-plan-level models in [2] share similarities with our operator and operator-subgraph model. However, we discovered the coverage-accuracy gap between the two models to be substantially large. To bridge this gap, we proposed additional mutually enhancing models and then combined the predictions of these individual models to achieve the best of accuracy and coverage. There are other works on query progress indicators [11, 31] that use the run time statistics from the currently running query to tell how much percentage of the work has been completed for the query. Our approach, in contrast, uses compile time features to make the prediction before the execution starts.

Cardinality is a key input to cost estimation and several learning and feedback-driven approaches [1, 39, 42]. However, these works have either focused only on recurring or strict subgraphs [1, 42], or learn only the ratio between the actual and predicted cardinality [39] that can go wrong in many situations, e.g., partial coverage results in erroneous estimates due to mixing of disparate models. Most importantly, as we discuss in Section 2, fixing cardinalities alone do not always lead to accurate costs in big data systems. There are other factors such as resources (e.g., partitions) consumed,

operator implementations (e.g., custom operators), and hardware characteristics (e.g., parallel distributed environment) that could determine cost. In contrast to cardinality models, CLEO introduces novel learning techniques (e.g., multiple models, coupled with an ensemble) and extensions to optimizer to robustly model cost. That said, cardinality is still an important feature (see Figure 4), and is also key to deciding partition counts, memory allocation at runtime, as well as for speculative execution in the job manager. A more detailed study on cardinality estimates in big data systems is an interesting avenue for future work.

Several works find the optimal resources given a physical query execution plan [3, 35, 40]. They either train a performance model or apply non-parametric Bayesian optimization techniques with a few sample runs to find the optimal resource configuration. However, the optimal execution plan may itself depend on the resources, and therefore in this work, we jointly find the optimal cost and resources. Nevertheless, the ideas from resource optimization work can be leveraged in our system to reduce the search space, especially if we consider multiple hardware types.

Finally, several recent works apply machine learning techniques to improve different components of a data system [25, 32]. The most prominent being learned indexes [26], which overfits a given stored data to a learned model that provides faster lookup as well as smaller storage footprint. [32] takes a more disruptive approach, where the vision is to replace the traditional query optimizers with one built using neural networks. In contrast, our focus in this paper is on improving the cost-estimation of operators in big data systems and our goal is to integrate learned models into existing query optimizers in a minimally invasive manner.

8 CONCLUSION

Accurate cost prediction is critical for resource efficiency in big data systems. At the same time, modeling query costs is incredibly hard in these systems. In this paper, we present techniques to learn cost models from the massive cloud workloads. We recognize that cloud workloads are highly heterogeneous in nature and *no one model fits all*. Instead, we leverage the common subexpression patterns in the workload and learn specialized models for each pattern. We further describe the accuracy and coverage trade-off of these specialized models and present additional mutual enhancing models to bridge the gap. We combine the predictions from all of these individual models into a *robust model* that provides the best of both accuracy and coverage over a sufficiently long period of time. A key part of our contribution is to integrate the learned cost models with existing query optimization frameworks. We present details on integration with SCOPE, a Cascade style query optimizer, and show how the learned models could be used to efficiently find both the query and resource optimal plans. Overall, applying machine learning to systems is an active area of research, and this paper presents a practical approach for doing so deep within the core of a query processing system.

REFERENCES

- [1] S. Agarwal, S. Kandula, N. Bruno, M.-C. Wu, I. Stoica, and J. Zhou. Re-optimizing data-parallel computing. In *Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation*, pages 21–21. USENIX Association, 2012.
- [2] M. Akdere, U. Çetintemel, M. Riondato, E. Upfal, and S. B. Zdonik. Learning-based query performance modeling and prediction. In *Data Engineering (ICDE), 2012 IEEE 28th International Conference on*, pages 390–401. IEEE, 2012.
- [3] O. Alipourfard, H. H. Liu, J. Chen, S. Venkataraman, M. Yu, and M. Zhang. Cherrypick: Adaptively unearthing the best cloud configurations for big data analytics. In *NSDI*, volume 2, pages 4–2, 2017.
- [4] AWS Athena. <https://aws.amazon.com/athena/>.
- [5] F. R. Bach and M. I. Jordan. Kernel independent component analysis. *Journal of machine learning research*, 3(Jul):1–48, 2002.
- [6] N. Bruno, S. Agarwal, S. Kandula, B. Shi, M. Wu, and J. Zhou. Recurring job optimization in scope. In *SIGMOD*, pages 805–806, 2012.
- [7] N. Bruno, S. Jain, and J. Zhou. Continuous cloud-scale query optimization and processing. *Proceedings of the VLDB Endowment*, 6(11):961–972, 2013.
- [8] N. Bruno, S. Jain, and J. Zhou. Recurring Job Optimization for Massively Distributed Query Processing. *IEEE Data Eng. Bull.*, 36(1):46–55, 2013.
- [9] N. Bruno, Y. Kwon, and M.-C. Wu. Advanced join strategies for large-scale distributed computation. *Proceedings of the VLDB Endowment*, 7(13):1484–1495, 2014.
- [10] R. Chaiken, B. Jenkins, P. Larson, B. Ramsey, D. Shakib, S. Weaver, and J. Zhou. SCOPE: easy and efficient parallel processing of massive data sets. *PVLDB*, 1(2):1265–1276, 2008.
- [11] S. Chaudhuri, V. Narasayya, and R. Ramamurthy. Estimating progress of execution for sql queries. In *Proceedings of the 2004 ACM SIGMOD international conference on Management of data*, pages 803–814. ACM, 2004.
- [12] CLEO: Technical Report. <http://arxiv.org/abs/2002.12393>.
- [13] A. Dutt, C. Wang, A. Nazi, S. Kandula, V. Narasayya, and S. Chaudhuri. Selectivity Estimation for Range Predicates Using Lightweight Models. *PVLDB*, 12(9):1044–1057, 2019.
- [14] FastTree. <https://www.nuget.org/packages/Microsoft.ML.FastTree/>.
- [15] A. D. Ferguson, P. Bodik, S. Kandula, E. Boutin, and R. Fonseca. Jockey: guaranteed job latency in data parallel clusters. In *EuroSys*, pages 99–112, 2012.
- [16] J. H. Friedman. Stochastic gradient boosting. *Computational statistics & data analysis*, 38(4):367–378, 2002.
- [17] A. Ganapathi, H. Kuno, U. Dayal, J. L. Wiener, A. Fox, M. Jordan, and D. Patterson. Predicting multiple metrics for queries: Better decisions enabled by machine learning. In *Data Engineering, 2009. ICDE'09. IEEE 25th International Conference on*, pages 592–603. IEEE, 2009.
- [18] Google BigQuery. <https://cloud.google.com/bigquery>.
- [19] G. Graefe. The Cascades framework for query optimization. *IEEE Data Eng. Bull.*, 18(3):19–29, 1995.
- [20] IBM BigSQL. <https://www.ibm.com/products/db2-big-sql>.
- [21] A. Jindal, K. Karanasos, S. Rao, and H. Patel. Selecting Subexpressions to Materialize at Datacenter Scale. In *VLDB*, 2018.
- [22] A. Jindal, S. Qiao, H. Patel, Z. Yin, J. Di, M. Bag, M. Friedman, Y. Lin, K. Karanasos, and S. Rao. Computation Reuse in Analytics Job Service at Microsoft. In *SIGMOD*, 2018.
- [23] S. A. Jyothi, C. Curino, I. Menache, S. M. Narayanamurthy, A. Tumanov, J. Yaniv, R. Mavlyutov, I. Goiri, S. Krishnan, J. Kulkarni, et al. Morpheus: Towards automated slos for enterprise clusters. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, pages 117–134, 2016.
- [24] A. Kipf, T. Kipf, B. Radke, V. Leis, P. Boncz, and A. Kemper. Learned cardinalities: Estimating correlated joins with deep learning. *CIDR*, 2019.
- [25] T. Kraska, M. Alizadeh, A. Beutel, E. Chi, J. Ding, A. Kristo, G. Leclerc, S. Madden, H. Mao, and V. Nathan. Sagedb: A learned database system. *CIDR*, 2019.
- [26] T. Kraska, A. Beutel, E. H. Chi, J. Dean, and N. Polyzotis. The case for learned index structures. In *Proceedings of the 2018 International Conference on Management of Data*, pages 489–504. ACM, 2018.
- [27] C. Lei, Z. Zhuang, E. A. Rundensteiner, and M. Y. Eltabakh. Redoop infrastructure for recurring big data queries. *PVLDB*, 7(13):1589–1592, 2014.
- [28] V. Leis, A. Gubichev, A. Mirchev, P. Boncz, A. Kemper, and T. Neumann. How good are query optimizers, really? *Proceedings of the VLDB Endowment*, 9(3):204–215, 2015.
- [29] J. Li, A. C. König, V. Narasayya, and S. Chaudhuri. Robust estimation of resource consumption for sql queries using statistical techniques. *Proceedings of the VLDB Endowment*, 5(11):1555–1566, 2012.
- [30] G. Lohman. Is query optimization a NP -hard problem. In *Proc. Workshop on Database Query Optimization*, volume 13. Oregon Graduate Center Comp. Sci. Tech. Rep, 2014.
- [31] G. Luo, J. F. Naughton, C. J. Ellmann, and M. W. Watzke. Toward a progress indicator for database queries. In *Proceedings of the 2004 ACM SIGMOD international conference on Management of data*, pages 791–802. ACM, 2004.
- [32] R. Marcus, P. Negi, H. Mao, C. Zhang, M. Alizadeh, T. Kraska, O. Papaemmanouil, and N. Tatbul. Neo: A learned query optimizer. *arXiv preprint arXiv:1904.03711*, 2019.
- [33] MART. <http://statweb.stanford.edu/~jhf/MART.html>.
- [34] M. Poess and C. Floyd. New tpc benchmarks for decision support and web commerce. *ACM Sigmod Record*, 29(4):64–71, 2000.
- [35] K. Rajan, D. Kakadia, C. Curino, and S. Krishnan. Perforator: eloquent performance models for resource optimization. In *Proceedings of the Seventh ACM Symposium on Cloud Computing*, pages 415–427. ACM, 2016.
- [36] R. Ramakrishnan, B. Sridharan, J. R. Douceur, P. Kasturi, B. Krishnamachari-Sampath, K. Krishnamoorthy, P. Li, M. Manu, S. Michaylov, R. Ramos, et al. Azure data lake store: a hyperscale distributed file service for big data analytics. In *Proceedings of the 2017 ACM International Conference on Management of Data*, pages 51–63. ACM, 2017.
- [37] A. Roy, A. Jindal, H. Patel, A. Gosalia, S. Krishnan, and C. Curino. SparkCruise: Handsfree Computation Reuse in Spark. *PVLDB*, 12(12):1850–1853, 2019.
- [38] J. Schäd, J. Dittrich, and J.-A. Quiané-Ruiz. Runtime Measurements in the Cloud: Observing, Analyzing, and Reducing Variance. *PVLDB*, 3(1-2):460–471, 2010.
- [39] M. Stillger, G. M. Lohman, V. Markl, and M. Kandil. Leo-db2's learning optimizer. In *VLDB*, volume 1, pages 19–28, 2001.
- [40] S. Venkataraman and others. Ernest: Efficient performance prediction for large-scale advanced analytics. In *NSDI*, pages 363–378, 2016.
- [41] L. Viswanathan, A. Jindal, and K. Karanasos. Query and Resource Optimization: Bridging the Gap. In *ICDE*, pages 1384–1387, 2018.
- [42] C. Wu, A. Jindal, S. Amizadeh, H. Patel, W. Le, S. Qiao, and S. Rao. Towards a Learning Optimizer for Shared Clouds. *PVLDB*, 12(3):210–222, 2018.
- [43] D. Xin, S. Macke, L. Ma, J. Liu, S. Song, and A. Parameswaran. HELIX: Holistic Optimization for Accelerating Iterative Machine Learning. *PVLDB*, 12(4):446–460, 2018.
- [44] Z. Yint, J. Sun, M. Li, J. Ekanayake, H. Lin, M. Friedman, J. A. Blakeley, C. Szyperski, and N. R. Devanur. Bubble execution: resource-aware reliable analytics at cloud scale. *Proceedings of the VLDB Endowment*, 11(7):746–758, 2018.
- [45] J. Zhou, N. Bruno, M.-C. Wu, P.-A. Larson, R. Chaiken, and D. Shakib. SCOPE: parallel databases meet MapReduce. *VLDB J.*, 21(5):611–636, 2012.
- [46] Q. Zhou, J. Arulraj, S. Navathe, W. Harris, and D. Xu. Automated Verification of Query Equivalence Using Satisfiability Modulo Theories. *PVLDB*, 12(11):1276–1288, 2019.
- [47] H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *Journal of the royal statistical society: series B (statistical methodology)*, 67(2):301–320, 2005.