Demonstration of Accelerating Machine Learning Inference Queries with Correlative Proxy Models

Zhihui Yang¹, Yicong Huang², Zuozhi Wang², Feng Gao¹, Yao Lu³, Chen Li², X. Sean Wang⁴

¹Zhejiang Lab, Hangzhou, China; ²UC Irvine, CA, USA;

³Microsoft Research, WA, USA; ⁴Fudan University, Shanghai, China.

¹{zhyang14,gaof}@zhejianglab.com, ²{yicongh1,zuozhiw,chenli}@ics.uci.edu,

³luyao@microsoft.com, ⁴xywangcs@fudan.edu.cn

ABSTRACT

We will demonstrate a prototype query-processing engine, which utilizes correlations among predicates to accelerate machine learning (ML) inference queries on unstructured data. Expensive operators such as feature extractors and classifiers are deployed as user-defined functions (UDFs), which are not penetrable by classic query optimization techniques such as predicate push-down. Recent optimization schemes (e.g., Probabilistic Predicates or PP) build a cheap proxy model for each predicate offline, and inject proxy models in the front of expensive ML UDFs under the independence assumption in queries. Input records that do not satisfy query predicates are filtered early by proxy models to bypass ML UDFs. But enforcing the independence assumption may result in sub-optimal plans. We use correlative proxy models to better exploit predicate correlations and accelerate ML queries. We will demonstrate our query optimizer called CORE, which builds proxy models online, allocates parameters to each model, and reorders them. We will also show end-to-end query processing with or without proxy models.

PVLDB Reference Format:

Zhihui Yang, Yicong Huang, Zuozhi Wang, Feng Gao, Yao Lu, Chen Li, X. Sean Wang. Demonstration of Accelerating Machine Learning Inference Queries with Correlative Proxy Models. PVLDB, 15(12): 3734 - 3737, 2022. doi:10.14778/3554821.3554887

1 INTRODUCTION

Consider an example workflow illustrated in Figure 1, where input tweets are processed by a machine learning (ML) user-defined function (UDF) Geotagger (\mathcal{F}_1) followed by a predicate state = 'CA' (σ_1), and another ML UDF Sentiment (\mathcal{F}_2) followed by a predicate sentiment = positive (σ_2). It enables downstream visualization and statistics, such as word cloud. ML queries are costly due to the expensive ML UDFs; improving the efficiency for ML inference has been a recent research focus. In our example, classic query optimization techniques such as predicate push-down cannot help much because σ_1 and σ_2 are stuck behind their corresponding ML UDFs regardless of their selectivity.

Recent works [5, 8] propose to rewrite the query and insert light-weight filters before the expensive ML UDFs, thus forming a

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit https://creativecommons.org/licenses/by-nc-nd/4.0/ to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 15, No. 12 ISSN 2150-8097. doi:10.14778/3554821.3554887

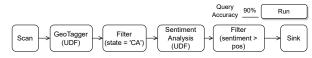


Figure 1: An example workflow for tweet analysis.

proxy model [11]. Figure 2 demonstrates an example plan with two proxy models $\hat{\sigma}_1$ and $\hat{\sigma}_2$; they quickly discard input records that are unlikely to satisfy the predicates and thus improve the query performance with an acceptable accuracy. In [8], an independence assumption is made to train filters and assemble these filters. But, query predicates are often correlated in many applications. In our example, sentiments may vary in different states – the sentiment in California can be different from that in Texas. The query optimization (QO) in [8] overestimates the reduction when building the filters and thus yields sub-optimal plans.

In our recent work [13], we proposed an optimizer called "CORE" that better exploits predicate correlations. By relaxing the independence assumption among different predicates, a proxy model is specific not only to a predicate but also its input relation, i.e., prefix σ 's and $\hat{\sigma}$'s, as well as parameter choices of prefix $\hat{\sigma}$'s. In Figure 2, $\hat{\sigma}_2$ learns upon filtering the raw input by $\hat{\sigma}_1 \wedge \sigma_1$. Note that \mathcal{F}_1 is a row processor and does not filter as σ_1 and $\hat{\sigma}_1$ do. Enumerating and building proxy models with different orders and parameter choices offline result in infeasible building and storage costs. CORE builds proxy models *online* to avoid exhaustive offline filter construction.

This demonstration will show an online query processing engine on top of CORE, which builds proxy models using a small portion of the input data, and executes the optimized plan on the remaining data. Users will be able to interact with the system in various ways, including submitting new queries and comparing performance with or without proxy models. We will also show a correlation score for a new query, and show end-to-end execution of a new query. We can observe performance improvements by up to 63% compared to [8] and by up to 80% compared to running the workflow as it is. In addition, CORE builds proxy models online for a new query and leverages a branch-and-bound search process to reduce the building costs. We will demonstrate CORE including building proxy models, allocating accuracy parameters to proxy models, and reordering with proxy models. Users will be able to interact with the query optimizer by specifying an order or deciding accuracy parameters for proxy models without waiting for the optimizer to finish to accelerate the QO phase.



Figure 2: An optimized query plan with proxy models.

2 THE DEMONSTRATION SYSTEM

We build a prototype on top of Texera [10], an open source system to support cloud-based collaborative data analytics using workflows. We implement a library of ML operators by deploying them as UDFs, and utilize Texera to support Java and Python functions. Table 1 shows a partial list of operators provided in the system. These operators depict row manipulators; they produce one output row per input row. Using these operators, developers are able to construct a workflow using a Web UI. They can also interact with the system by choosing to execute the workflow as it is, using PPs, or using CORE.

Table 1: A partial list of ML modules provided in the system.

Module Name	Description
Entity Recognition	Label sequences of words such as person or
	company names.
Sentiment Analysis	Predict the sentiment of the text (i.e., positive
	or negative).
Pos Tagger	Assign parts of speech to each word (e.g., a
	noun or a verb).
Object Detection	Identify objects in an image or video (e.g., a
	dog or a cup).
Activity Recognition	Predict the activity of a person (e.g., applying
	lipstick).

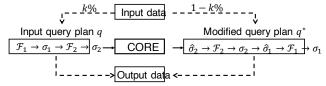


Figure 3: Given a query plan q, CORE generates an optimized plan q^* by applying proxy models. Part of the input data (k%) is used for building proxy models, and the remaining data is processed by q^* .

Figure 3 shows the architecture of the CORE system. Its input is a workflow that includes multiple ML UDFs, such as the example in Figure 1. We follow the setup of the prior work, such as NoScope [5] and PP [8], and leverage the query interface in [8]. Users are able to specify a global target accuracy A, which sets a trade-off goal between an acceptable accuracy and query-processing speedup. In Figure 1, $\mathcal{A} = 90\%$. CORE optimizes the input workflow by building proxy models online. A proxy model is specific to a predicate $c\phi v$, where c is a predicate column, ϕ is a comparison (e.g., > or =), and v is a constant value. A workflow can have one or more predicate clauses in conjunction: $\wedge c\phi v$. CORE builds a proxy model for each predicate online, considers proxy models' combinations, allocates their accuracy parameters, and injects them into the modified query plan q* (Figure 2). A small portion of the input data (e.g., k%) is used to build proxy models, and the remaining data is processed by the optimized plan q^* .

Key technical challenges: We build the demonstration system to answer the following technical questions to reduce the overhead of building proxy models online. More details can be found in [13].

- Building proxy models online. Enumerating and building proxy models offline result in infeasible building and storage costs. To build $\hat{\sigma}$ online, the demonstration system generates its labeled sample L by pulling initial records from the input, filtering these records by its input relation d, and labeling L using its predicate σ . Sample L is divided into the training set, the testing set, and the validation set. We re-sample the training data to ensure a label balance. The classifier M is trained on the training set using light-weight classification algorithms, such as a linear SVM and a shallow NN. During training, we leverage a grid-search cross-validation on the F1-score to decide the best hyper-parameters. We derive the accuracy versus reduction curve R using the validation set.
- Allocating parameters for proxy models. The system leverages a hill-climbing search to find an optimal accuracy allocation for proxy models. The main challenge is that building proxy models online is time-consuming because (i) there are an exponential number of candidates ô's, and (ii) generating a labeled sample and training a classifier are computationally costly. The system reuses previously materialized samples and trained classifiers to reduce labeling costs and training costs, respectively.
- Reordering proxy models. Building all proxy models for different orders can be computationally expensive. We leverage a branch-and-bound search to prune candidate plans. Specifically, we compute a lower bound and an upper bound of costs $\sum C$ for a specific order of proxy models. During the search, we tighten the lower and upper bounds of $\sum C$ as we collect more information, such as selectivity and reduction, and prune candidate plans to reduce the optimization overhead. Additionally, we adopt a fine-grained search tree to improve the search process further.

3 DEMONSTRATION SCENARIOS

3.1 Correlated Workflows

In the demonstration, we provide three datasets with text, images, and videos, an operator library as illustrated in Table 1, and several pre-constructed workflows over the datasets with different correlation among the query predicates. These workflows retrieve texts, images, and videos that match given filter conditions, which are conjunctions of multiple clauses. Each filter is an equality condition on an ML-generated label column. Users will also be able to construct a workflow using the operator library.

Twitter text dataset. It contains 2M tweets from January 2017 to September 2017 in the United States. Each tweet is a string with a maximum of 140 characters. We demonstrate several workflows with various NLP modules such as entity recognition and sentiment analysis to analyze tweets.

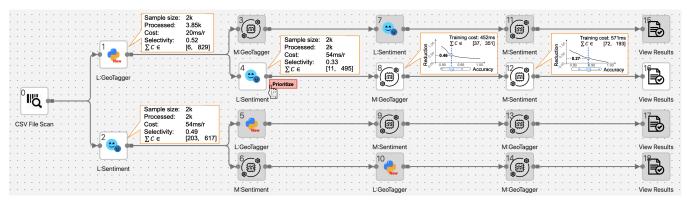


Figure 4: A fine-grained tree generated by CORE demonstrating runtime metrics of five nodes and the selected order Geotagger → Sentiment.

COCO image dataset. COCO is a public dataset with 123K images and 80 object classes such as "person," "bicycle," and "dog." Each image is labeled with multiple objects for their class labels and bounding box positions. We demonstrate workflows with different levels of correlation among the query predicates.

UCF101 video dataset. UCF101 contains 13K videos collected from YouTube. Each video is labeled with one of 101 action categories such as "applying lipstick" and "baby crawling." We demonstrate several workflows with object detection and activity recognition models to retrieve videos with specific labels.

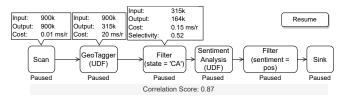


Figure 5: A paused workflow demonstrating its correlation score and the runtime metrics of the first three operators.

We show the correlation score of the query predicates on the interface. We compute the correlation score using its formulation provided by CORDS [4]. When users start running the workflow, the system collects and shows each operator's runtime metrics, such as the number of processed records, cost, and selectivity. Users may pause the workflow, observe each operator's runtime metrics, and resume the workflow to gain valuable insights [12]. In Figure 5, the user paused the workflow and observed the first three operators' runtime metrics. The cost of the Geotagger is 20ms per record, and the selectivity of the state = 'CA' filter is 0.52.

3.2 Visualizing Query Optimization with CORE

In this scenario, we use pre-constructed workflows over the Twitter dataset to demonstrate the CORE query optimizer including building proxy models online, allocating accuracy to each proxy model, and reordering with proxy models. We show the optimization overhead of the workflow and how users can interact with the query optimizer.

Given a workflow, CORE converts candidate query plans in the space \mathbb{H} to a fine-grained tree. Figure 4 shows the fine-grained tree

for the workflow in Figure 1. There are two types of tree nodes. (i) An *L*-node represents a pair of ML UDF and its filter, and it is to generate labeled samples. (ii) An *M*-node is to train its corresponding classifiers and derive an accuracy-versus-reduction relationship *R*. The root of the tree is the CSV File Scan operator in the workflow, which reads initial records from the input. Leaf nodes are View Results operators, which store results of the workflow.

The problem of finding an optimal order of proxy models and allocating their accuracy parameters simultaneously is NP-hard. CORE searches over the tree using a branch-and-bound algorithm to prune candidate plans. Our demo shows each node's runtime metrics during the search process. For an L-node, we show the number of records in the labeled sample, the number of processed records, cost, and selectivity. For an M-node, we show the training cost and the accuracy-versus-reduction curve. CORE tightens the lower and upper bounds for different query plans as we collect information, such as selectivity and reduction. The demo shows each candidate's lower and upper bounds of $\sum C$.

Figure 4 illustrates the runtime metrics of five nodes, including three L-nodes and two M-nodes. For the runtime metrics of node 1, the Geotagger processed 3.85K records but only kept the first 2K records in its labeled sample. The lower and upper bounds of $\sum C$ for the order Geotagger \rightarrow Sentiment became 6 and 829 after collecting the cost and selectivity of the Geotagger. Based on the runtime metrics of M-node 8, the accuracy-versus-reduction curve is shown. The lower and upper bounds of Geotagger \rightarrow Sentiment became 37 and 352, respectively, after collecting the 0.46 reduction. We also demonstrate pruned candidate plans. Their corresponding nodes in the fine-grained tree are gray, such as nodes 3, 5, and 6.

The user may interact with CORE by specifying an order or deciding an accuracy parameter for a proxy model without waiting for the optimizer to reduce the QO overhead. Specifically, the user can select a query plan from candidates to build proxy models. In Figure 4, the user clicks the Priority button and node 4 to select the order Geotagger \rightarrow Sentiment. After observing the accuracy-versus-reduction curve from an M-node, the user can specify an accuracy parameter for a node's proxy model by choosing an accuracy value on the curve. For node 8 in Figure 4, the choice is at a 0.46 reduction.

3.3 End-to-End Query Processing

We use the pre-constructed workflows over the Twitter dataset to demonstrate the system. A user is able to submit a new workflow and specify a query accuracy A. The demonstration system supports two other modes to run a workflow besides using CORE. (i) ORIG runs the original workflow as it is, and (ii) PP builds a light-weight proxy model for each predicate, and injects them early in a plan. PP decides the accuracy parameter for each proxy model using a dynamic programming algorithm with an independence assumption of predicates. Users are able to run and compare the workflow using ORIG, PP, and CORE.

The demonstration system uses a small portion of the input data to generate an optimal plan, and processes the remaining data using the optimal plan. The demo shows the percentage of the initial input data used to build proxy models, the total time of running the workflow including the optimization overhead, the QO cost percentage (i.e., the percentage of the QO time over the total processing time), and the execution cost. Users may compare the performance of running the workflow with different modes. In general, both PP and CORE improve the performance of workflows, and CORE achieves more improvement on correlated workflows than PP.

4 RELATED WORK

Operator reordering in database optimization. [1] studies how to order correlated predicates in streaming systems. It uses a greedy algorithm for selection ordering and collects samples at runtime to estimate selectivity. Our QO gives an optimal solution and uses a branch-and-bound search to quickly prune plans in the space of proxy models. [9] studies various optimization techniques of complex user-defined functions on map-reduce-style big data systems, such as predicate simplification and UDF semantic inference. These techniques are orthogonal to our solution. [2] provides approximate answers to queries by running queries on a small sampling subset of data. Our approach provides approximate answers by exploiting the accuracy of ML inference predicates.

Optimization with Proxy models (a.k.a. cascaded filters). Proxy models have been studied for decades to accelerate ML inference. Jones et al. [11] cascade weak classifiers as proxy models to speed-up face detection in images. Recently proxy models have been applied in big-data systems to accelerate ML inference-based analysis tasks [3, 5, 6, 8]. NoScope [5] inserts a cheap specialized model before expensive DNNs to accelerate selection video queries. Certain classes of video queries [7] including selection without guarantees [3] and selection with statistical guarantees [6] are optimized using proxy models. Probabilistic predicates (PPs) [8] optimize various workloads by inserting multiple offline-built proxy models before expensive ML UDFs with an assumption of independence

between predicates. Unlike [3, 5, 6], PP and our proposed CORE leverage general proxy models for various workloads. CORE follows this line of work and further relaxes the independence assumption.

5 CONCLUSIONS

This paper demonstrates a novel query optimizer, CORE, to accelerate ML inference queries. It improves state-of-the-art techniques by relaxing the independence assumption among different predicates.

ACKNOWLEDGMENTS

This work was partially supported by the National Key R&D Program of China (No. 2020AAA0103903), the NSFC (No. 61732004) and the USA NSF award IIS-2107150. We want to thank the Wide-Area Joint Computing team at Zhejiang Lab for their contributions to the development of the system.

REFERENCES

- [1] Shivnath Babu, Rajeev Motwani, Kamesh Munagala, Itaru Nishizawa, and Jennifer Widom. 2004. Adaptive Ordering of Pipelined Stream Filters. In Proceedings of the ACM SIGMOD International Conference on Management of Data, June 13-18, 2004. ACM, Paris, France, 407-418.
- [2] Surajit Chaudhuri, Bolin Ding, and Srikanth Kandula. 2017. Approximate Query Processing: No Silver Bullet. In Proceedings of the 2017 ACM International Conference on Management of Data, SIGMOD Conference 2017, May 14-19, 2017. ACM, Chicago, IL, USA, 511-519.
- [3] Kevin Hsieh, Ganesh Ananthanarayanan, Peter Bodík, Shivaram Venkataraman, Paramvir Bahl, Matthai Philipose, Phillip B. Gibbons, and Onur Mutlu. 2018. Focus: Querying Large Video Datasets with Low Latency and Low Cost. In 13th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2018, October 8-10, 2018. USENIX Association, Carlsbad, CA, USA, 269-286.
- [4] Ihab F. Ilyas, Volker Markl, Peter J. Haas, Paul Brown, and Ashraf Aboulnaga. 2004. CORDS: Automatic Discovery of Correlations and Soft Functional Dependencies. In Proceedings of the ACM SIGMOD International Conference on Management of Data, June 13-18, 2004. ACM, Paris, France, 647-658.
- [5] Daniel Kang, John Emmons, Firas Abuzaid, Peter Bailis, and Matei Zaharia. 2017. NoScope: Optimizing Deep CNN-Based Queries over Video Streams at Scale. PVLDB 10, 11 (2017), 1586–1597.
- [6] Daniel Kang, Edward Gan, Peter Bailis, Tatsunori Hashimoto, and Matei Zaharia. 2020. Approximate Selection with Guarantees using Proxies. Proc. VLDB Endow. 13, 11 (2020), 1990–2003.
- [7] Yao Lu, Aakanksha Chowdhery, and Srikanth Kandula. 2016. Optasia: A relational platform for efficient large-scale video analytics. In Proceedings of the Seventh ACM Symposium on Cloud Computing. 57–70.
- [8] Yao Lu, Aakanksha Chowdhery, Srikanth Kandula, and Surajit Chaudhuri. 2018. Accelerating Machine Learning Inference with Probabilistic Predicates. In Proceedings of the 2018 International Conference on Management of Data, SIGMOD Conference 2018, June 10-15, 2018. ACM, Houston, TX, USA, 1493-1508.
- [9] Astrid Rheinländer, Ulf Leser, and Goetz Graefe. 2017. Optimization of Complex Dataflows with User-Defined Functions. ACM Comput. Surv. 50, 3 (2017), 38:1– 28:20
- [10] Texera. 2021. Texera Website. https://github.com/Texera/texera/.
- [11] Paul A. Viola and Michael J. Jones. 2001. Rapid Object Detection using a Boosted Cascade of Simple Features. In 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2001), with CD-ROM, 8-14 December 2001. IEEE Computer Society, Kauai, HI, USA, 511-518.
- [12] Zuozhi Wang, Avinash Kumar, Shengquan Ni, and Chen Li. 2020. Demonstration of interactive runtime debugging of distributed dataflows in Texera. *Proceedings* of the VLDB Endowment 13, 12 (2020), 2953–2956.
- [13] Zhihui Yang, Zuozhi Wang, Yicong Huang, Yao Lu, Chen Li, and X. Sean Wang. 2022. Optimizing Machine Learning Inference Queries with Correlative Proxy Models. CoRR abs/2201.00309 (2022).