# OntoTune: Ontology-Driven Learning for Query Optimization with Convolutional Models

1st Songhui Yue
*Charleston Southern University*
syue@csunive.edu

2nd Yang Shao
*Charleston Southern University*
yshao@student.csuniv.edu

3th Sean Hayes
*Charleston Southern University*
shayes@csuniv.edu

*Abstract*—**Query optimization has been studied using machine learning, reinforcement learning, and, more recently, graph-based convolutional networks. Ontology, as a structured, information-rich knowledge representation, can provide context, particularly in learning problems. This paper presents OntoTune, an ontology-based platform for enhancing learning for query optimization. By connecting SQL queries, database metadata, and statistics, the ontology developed in this research is promising in capturing relationships and important determinants of query performance. This research also develops a method to embed ontologies while preserving as much of the relationships and key information as possible, before feeding it into learning algorithms such as tree-based and graph-based convolutional networks. A case study shows how OntoTune's ontology-driven learning delivers performance gains compared with database system default query execution.**

*Index Terms*—**Ontology, Query Optimization, Feature Embedding, Convolutional Neural Network**

## I. INTRODUCTION

Database query optimization has been studied for decades, starting from fixed heuristic models, to learning models [1], [2], and the latter highly relies on deep reinforcement learning models [3], and other deep learning models [4], such as a transformer model [5]. Those studies are performed on query datasets that are composed of different types [6], [7]. These queries often arise from template families (e.g., particular table-usage patterns or the presence of DISTINCT/COUNT). The database workload contains the tables and attributes, and their related information includes table and attribute schemas, index availability, data distributions and skew, and the quality of cardinality estimates. All these types of information, along with the execution environment, when organized into features, can determine the execution speed of a query set in a particular database system, such as PostgreSQL, a widely used open-source system.

Modern workloads such as the Jobs/Join-Order Benchmark (JOB) [8] and StackOverflow (so_pg) [9] expose a persistent challenge: long-tail queries and environment-sensitive behavior that fixed-cost models fail to capture [7]. Learning-based approaches help, but they can be brittle: outcomes depend on batch size, random seed, hot vs. cold start, and on hardware, statistics, and DBMS configuration (GUCs) [10]. From a modeling standpoint, ML predictors exhibit high variance across training sets [11] and suffer when deployed under distribution shift, yielding large errors on data that differ from the training distribution [11], [12]. These factors make experiments unstable, hard to reproduce, and difficult to explain.

In this study, we propose OntoTune, an ontology-driven learning platform for query optimization. Our view in this study is simple: an ontology and a concrete knowledge graph (KG) provide the missing semantic and provenance layers for learned optimization. We name the objects and relations that matter to plan selection, and we record runs with their full context. Queries, plan trees and operators, planner arms (i.e., bundles of PostgreSQL GUC switches), runtime metrics, environments, sampling policies, and model predictions/rewards all become first-class, queryable entities. The ontology (TBox) [13] gives the schema, and the KG (ABox) [14] stores instances.

Ontotune is a PostgreSQL-based prototype. The system follows a similar implementation architecture to Bao (a learned component that sits on top of the native optimizer) [7], to which we add a query-related metadata and system configuration extraction module via the PG extension, and a workflow on the training side that utilizes extracted data as feature providers. By preserving the query metadata, DB statistics, plan trees and operators, runtime and buffer metrics, and relevant Grand Unified Configurations (GUCs) (such as `enable_*` and `work_mem`), experiments are traceable and comparable across configurations.

For the utilization of the ontology and KG, we adopt a CNN-based predictor [15] trained on an ontology-derived feature matrix that fuses SQL, plan, and context information. On the datasets considered, the CNN is our strongest variant. For online arm selection (bandit-style reinforcement learning), we apply a 1-complement reward transform to reduce the autoselection of seldom-used arms. Because the representation lives at the ontology/KG layer, tree- and graph-based models (TreeConv [7], GCN [16], GAT [17]) can plug into the same data pipeline without re-engineering; we treat them as forward-looking components for future exploration. Case studies illustrate both regimes: when the learned policy outperforms the PostgreSQL baseline and when gains disappear (e.g., under cold starts or highly exploratory settings).

We format the database and query context as an ontology/knowledge graph (KG) in this application for four reasons: First, interoperability [18]: a minimal, task-driven vocabulary avoids over-engineering while standardizing names
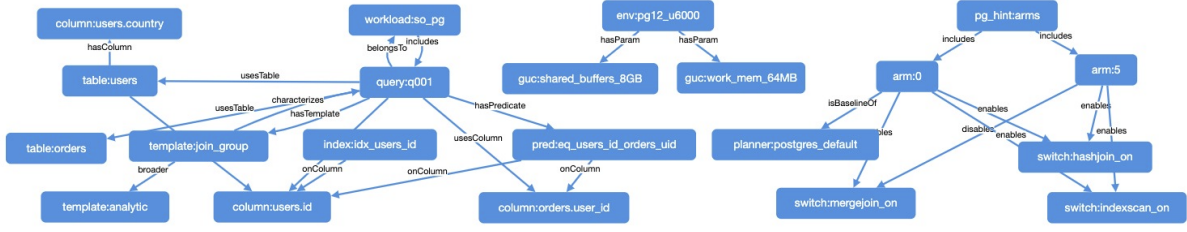
Fig. 1. A Partial Knowledge Graph for DB Query Optimization (with Ontology and Instance Examples)

and relations. Second, explanation [19]: the KG ties each run to its query, arm, plan, and metrics, making case analysis straightforward. Third, reproducibility: environment/GUC snapshots, plan fingerprints, and sampling choices (hot/cold start, batch regimes, exploration) are part of the data, not side notes. Figure 1 sketches the overall design: the ontology defines the core classes; the KG records per-run evidence and links among them. Fourthly, the ontology+KG can provide structured representations that feed into AI models, such as LLMs [10], to improve context understanding for further query optimization and transfer learning.

The contributions of this study are summarized as follows:

- A minimal ontology (TBox) and a practical knowledge graph (ABox) for learning-based query optimization, covering queries, plans, operators, arms (GUC bundles), environments, metrics, predictions, and rewards.
- OntoTune: a PostgreSQL-based platform that mirrors prior learned-optimizer integration but adds end-to-end provenance capability: plan extraction, metric collection, and explicit snapshots of environment and GUC settings.
- The case study combines CNN into the reinforcement learning process through designing a reward-cost complement transform trick to avoid the autoselection of seldomly-used arms between different batches. CNNs provide our strongest results so far on selected datasets.

The subsequent sections are structured as follows: Section II introduces the ontology and KG schema and how we populate and embed them into feature matrices. Section III describes the implementation of the learning pipeline. Section IV presents case studies and evaluation under fixed configurations. Section V discusses strengths and limitations. Section VI reviews related work, and Section VII concludes with future directions.

## II. METHODOLOGY

The methodlogies in this study includes the ontology design, data extraction, embeding, and the learning pipleine.

### A. Query-level Ontology

As demonstrated in Fig 1, we name the objects and relations that matter to learned plan selection and make every run connected with the query, plan, and configurations. The figure is composed of nodes that contain two parts: the left side of the colon contains concepts, corresponding to the ontology TBox, and the right side contains entities, corresponding to the KG.



Fig. 2. OntoTune's system architecture: developed based on Bao's workflow [7].

The core classes include: {Query, Plan, PlanNode, Table, Column, Index, Arm, Environment, Execution, Reward}. Typical relations are: (Query, hasPlan, Plan), (Plan, hasNode, PlanNode), (Execution, useArm, Arm).

A PG extension instantiates the ABox per run (plan trees, per-operator stats, runtime/buffer metrics, GUC snapshots), so each execution ties a concrete (Query, Plan, Arm, Environment) to observed configurations and predictions.

### B. Embedding

For a context (template) $\tau$, let the column universe $\mathcal{U}_\tau$ enumerate referenced Table/Column instances. We build a matrix

$$X_\tau \in \mathbb{R}^{R \times C}, \quad C = |\mathcal{U}_\tau|,$$

where rows are feature channels drawn from: (i) template/SQL indicators and buckets; (ii) plan-level cost/row shares; (iii) operator-type cost shares aggregated to referenced columns;

TABLE I
REPRESENTATIVE FEATURES USED IN THE ONTOTUNE FEATURE MATRIX.

| Feature name | Meaning / How computed |
|---|---|
| `tpl_has_distinct` | Query uses DISTINCT |
| `tpl_need_sort_for_merge` | MergeJoin requires explicit sort (derived from plan/heuristics); Binary. |
| `tpl_group_by_cols_bucket_{0..3}` | One-hot over bucketized #GROUP BY columns; exactly one bucket set to 1 (others 0). |
| `tpl_rows_bucket_{0..2}` | One-hot over bucketized estimated result size; |
| `sql_has_window` | Window functions present; parsed into `template_features` then broadcast. |
| `sql_has_like` | LIKE patterns present; binary. |
| `sql_num_join_bucket_{0..2}` | One-hot over bucketized join count (e.g., 0, 1–2, $\geq$3). |
| `sql_num_subquery_bucket_{0..1}` | One-hot over subquery count (0 vs. $\geq$1). |
| `plan_cost_share` | Per-table cost share normalized by total plan cost, then broadcast to that table's columns. |
| `plan_rows_share` | Per-table rows share normalized by total estimated rows, broadcast similarly. |
| `col_cost_from_scan_share` | Fraction of total plan cost attributed to scans mapped onto each referenced column. |
| `col_cost_from_agg_share` | Fraction of cost from Aggregate nodes assigned to their argument/output columns. |
| `plan_cost_op_HashJoin_share` | Total cost share of HashJoin nodes over the whole plan; normalized to $[0, 1]$. |
| `plan_cost_op_Sort_share` | Total cost share of Sort nodes; normalized to $[0, 1]$. |

and (iv) column traits (numeric/indexed/in-where/in-join/in-orderby). Broadcasting from plan nodes $p$ to columns $u \in \mathcal{U}_\tau$ follows:

$$\text{plan\_cost\_share}(u) = \sum_{p:\, u \in S(p)} \frac{\text{cost}(p)}{\sum_q \text{cost}(q)},$$
$$\text{plan\_rows\_share}(u) = \sum_{p:\, u \in S(p)} \frac{\text{rows}(p)}{\sum_q \text{rows}(q)}. \tag{1}$$

We reform the data from SQL/plan/context to well-typed properties that are broadcasted into the matrix embedding and the epresentative properties can be found in Table I. The same extraction pipeline can export adjacency/incidence for Tree/Graph models.

The basic SQL-based ontology information is passed from the PG extension along with the plan, including analysis of a particular table attribute, such as whether it is numeric, whether it forms an index, and whether it needs sorting. When the values of those attributes are synthesized to the table level, they will decide the table's state in the matrix, as shown in Fig 3 (d).

### C. OntoTune's System Architecture

As shown in Fig 2, we use PG 12 as the basis for prototyping the OntoTune, as it is easier to make comparisons with the prior research work. We use C to implement the extension. We put the extraction of the DB statistics on the extension side, including the information from pg_statistics, pg_class, and pg_namespace, as well as getting the information from GUC. The cost information is from the plans that are also sent from the extension side.

Before the execution, we store the SQLs in SQLite3, which are later used to extract extra query-level ontology information using SQLGLOT, a Python library for parsing SQL statements. The "run_query" script is to run the batches of queries and invoke the training between two batches. OntoServer provides

the remote procedures, including recording the rewards of each query and making predictions for a specific query.

### III. IMPLEMENTATION

To adopt CNN into the online reinforcement learing process, it is important to tackle a challenge that, when meeting with less samples, for example, for arm 1, which normally performs the worst, it won't be selected many times in the cold start running set, where an arm selection is decided by the PG's cost estimation, thus, in the second round, it tends to predict it to be the best arm to use.

We implemented a method to avoid this case (although for some cases, it is worth the exploration with a tailored design): by having an inverse of the predicted value – turn cost as reward and objective after training turns to getting the most significant reward, and finally, make a reverse to get the cost. This way, this process walks around the effect that the trained model gives the unseen arms the smallest cost, because the model after the inverse will just get the largest number to be the fastest (after 1-x), and the smallest number from the unseen arms is dropped.

We consider $K$ candidate arms $a \in \{0, \ldots, K-1\}$ for each incoming query under context $\tau$ with features $X_\tau$. Let $y_{\tau,a} > 0$ denote the observed runtime cost (e.g., milliseconds) when arm $a$ is executed in context $\tau$. To stabilize the long-tailed distribution of runtimes, we apply a log–min–max scaler. We regress a *high-is-good* reward and select arms by minimizing the recovered cost with optimism, candidate filtering and gated exploration.

*a) Scaling and complement.:* Let $y_{\tau,a} > 0$ be runtime. Apply log–min–max scaling

$$\phi(y) = \frac{\log(1+y) - \ell_{\min}}{\ell_{\max} - \ell_{\min}} \in [0, 1], \tag{2}$$
$$\phi^{-1}(c) = \exp\big(c(\ell_{\max} - \ell_{\min}) + \ell_{\min}\big) - 1. \tag{3}$$

**Mapping between description and SQL**

| | |
|---|---|
| Target | select t1.name |
| from | table1 t1 |
| where | t1.condition > number1 |

(a) — select t1.name / from table1 t1 / where t1.condition > number1 / group_by / sorting name asc

(b)

(c) — inTarget(attribute1) boolean / inFrom(table1) boolean / inWhere(table1) boolean / inWhere(attribute1) boolean

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| inTarget | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| inFrom | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 |
| table_arributes (1/0) | table_1 | table_2 | table_... | table_1 .attibute1 | table_1 .attibute2 | ... | table_2 .attibute1 | table_2 .attibute2 | ... | table_x .attribute... |
| inWhere | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 |
| isIndex | 0 | 0 | 0 | 1 | 0 | 9 | 1 | 0 | 0 | 0 |
| isNumber | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 |
| isSorted | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |

(d)

Fig. 3. Partial Embedding of the OntoTune Ontology: for an SQL query

*b) Reward–Cost Complement Transform:* Define the complement

$$J(x) = 1 - x, \qquad (4)$$

and the (scaled) reward

$$r_{\tau,a} = J\big(\phi(y_{\tau,a})\big) = 1 - \phi(y_{\tau,a}) \in [0,1]. \qquad (5)$$

*c) Model and loss.:* A CNN $f_\theta$ predicts reward from $(X_\tau, a)$:

$$\hat{r}_{\tau,a} = f_\theta(X_\tau, a) \in [0,1]. \qquad (6)$$

Recover scaled/real cost for selection:

$$\hat{c}_{\tau,a} = J(\hat{r}_{\tau,a}) = 1 - \hat{r}_{\tau,a}, \qquad \tilde{y}_{\tau,a} = \phi^{-1}(\hat{c}_{\tau,a}). \qquad (7)$$

Train with MSE in reward space:

$$\mathcal{L}(\theta) = \frac{1}{|D|} \sum_{(\tau,a)\in D} \big(f_\theta(X_\tau, a) - r_{\tau,a}\big)^2. \qquad (8)$$

*d) Scoring and candidate filtering.:* With per-context counts $n_{\tau,a}$ and optimism hyperparameter $\beta \geq 0$,

$$s_{\tau,a} = \tilde{y}_{\tau,a} - \frac{\beta}{\sqrt{\max(1, n_{\tau,a})}}. \qquad (9)$$

Keep the best $K_{\text{top}}$ arms after removing a banned tail $B_\tau$ and enforcing a minimum sample threshold $n_{\min}$:

$$A_\tau^{\text{cand}} = \Big(\text{TopK}_a\big(s_{\tau,a}; K_{\text{top}}\big) \setminus B_\tau\Big) \setminus \{a : n_{\tau,a} < n_{\min}\}. \qquad (10)$$

*e) Gated $\varepsilon$-greedy and final decision.:* Let $N_\tau$ be the frequency of $\tau$, and $\texttt{estC}_\tau$ a normalized batch cost:

$$\varepsilon_\tau = \max\Big(\varepsilon_{\min}, \frac{\varepsilon_0}{\sqrt{N_\tau}}\Big) \cdot \frac{1}{1 + \texttt{estC}_\tau}. \qquad (11)$$

Choose

$$a_\tau^\star = \begin{cases} \arg\min_{a \in A_\tau^{\text{cand}}} s_{\tau,a}, & \text{with prob. } 1 - \varepsilon_\tau, \\ \text{Uniform}\big(A_\tau^{\text{cand}}\big), & \text{with prob. } \varepsilon_\tau. \end{cases} \qquad (12)$$

## IV. CASE STUDY AND EXPERIMENTS

CNN is used to capture features from the matrix formed by the ontology embedding and to perform regression on which arm's plan will have the best execution time – reward prediction. The purpose of this experiment is to serve as a case study to demonstrate how the OntoTune platform can provide insights into what should be added or emphasized next to make the platform better suited to support query optimization research.

The reported results are confirmed by three independent runs using the GitHub release (https://github.com/songhui01/OntoTune.git) and the stated configuration. Empirically, the model sometimes selects the best arms for long-tail queries in some datasets and

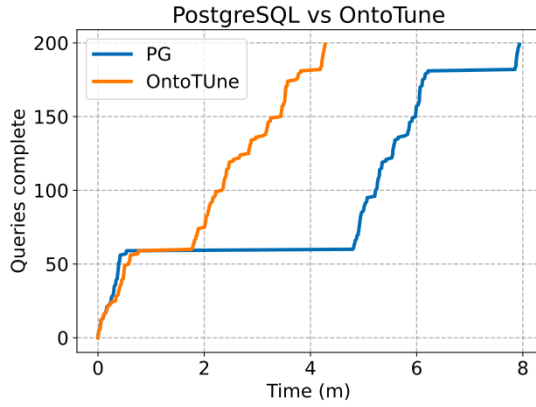Fig. 4. Queries vs Time for PostgreSQL and OntoTune with Random Set (42).



Fig. 5. Queries vs Time for PostgreSQL and OntoTune with Random Set (49).

configurations, but in another dataset, it can take longer due to selecting a very expensive arm.

### A. Configurations

Unless otherwise stated, all experiments were executed on Ubuntu 24.04 with PostgreSQL 12. The workload is the StackOverflow (so_pg) dataset. The host machine is equipped with an Intel Core i5-class CPU and an NVIDIA Blackwell GPU (20 GB VRAM).

We fixed the following GUC settings to make results comparable across runs: `shared_buffers=2GB`, `work_mem=4MB`, `maintenance_work_mem=64MB`, `effective_io_concurrency=2`, `max_worker_processes=8`, `max_parallel_workers_per_gather=1`, `max_parallel_workers=1`, and `effective_cache_size=2GB`. All other parameters remained at PostgreSQL defaults.

### B. Results

On StackOverflow (so_pg) under the fixed configuration, OntoTune outperforms PostgreSQL for one representative random seed (Fig. 4): the orange curve rises more steeply and completes the 200-query batch several minutes earlier, with fewer long plateaus. For another seed (Fig. 5), OntoTune underperforms: after an initial rise, it exhibits extended stalls and finishes noticeably later than PostgreSQL. These trends were consistent across three independent repetitions for each seed.

### C. Analysis

The win in Fig. 4 is consistent with our design: the learned policy leverages prior executions to identify better arms to shorten completion time by avoiding long stalls. The curve shows that the method performs well when meeting the two long-tailing queries. The loss in Fig. 5 occurs when early exploration or misranking admits an expensive arm into the candidate set, as a result, a few costly queries dominate wall-time and flatten the curve.
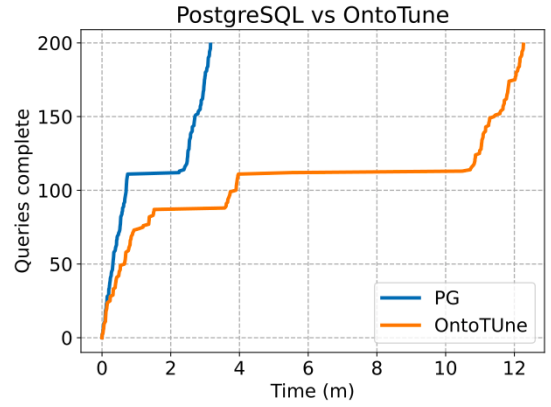
A plausible explanation for Fig. 5 is that the system was still in the learning phase; with sufficient time and samples, Onto-Tune is expected to recover and improve. To make such improvements reliable, we plan to enrich the ontology/knowledge graph with tail-aware nodes and relations (e.g., spill risk, operator-level memory pressure, join-selectivity drift) so that root causes of long tails can be systematically mitigated via arm choice or dynamic configuration [10].

## V. Discussion

Our study positions OntoTune as a platform rather than a single model: it encodes queries, plans, and execution context into an ontology/knowledge-graph layer and learns over this representation. This separation between the semantic layer and the learner enables CNNs (and, in principle, graph/TreeConv variants) to operate on a consistent, query-level abstraction. Compared with Bao—which learns features directly on the plan tree—our approach explicitly transforms SQL/plan/context signals into ontology relations and features. The case study shows feasibility: a convolutional regressor can select effective arms and, in favorable settings, avoid bad plans that produce long stalls.

One limitation of the study is that the performance can be sensitive to environment and parameter changes (e.g., GUCs, statistics refresh), leading to instability even when the selected arms are usually reasonable. The ontology currently covers most—but not all—categories we need; tail-aware attributes are still incomplete. Because learning requires many batches of executed queries, our experiments used relatively small subsets, and long-horizon benefits were not thoroughly tested. These are threats to validity and define the boundaries of what the present evidence supports.

In the long run, instability can be mitigated with environment-aware transfer learning (offline pretraining followed by light online adaptation), and to generalize across workloads by mining templates and using ontology-graph similarity for retrieval and warm starts. On the representation side, the TBox/ABox can be broadened with tail-focused nodes and relations, then instrument a closed loop that (i) automatically

updates the knowledge base from experimental feedback, (ii) schedules/evaluates new runs, and (iii) deploys safety fallbacks when predicted risk or uncertainty is high. Together, these steps aim to make ontology-driven learning a robust and systematic path toward practical query optimization.

## VI. RELATED WORK

Learned query optimization has progressed from plan-tree encoders to end-to-end selection policies. Bao learns over operator trees and uses a tree-convolutional network to rank alternative knobs/arms, demonstrating practical gains under controlled settings [7]. Subsequent systems extend this direction with alternative learners and search strategies (e.g., learned re-ranking and tighter candidate exploration). Lero further explores model-driven selection by embedding query/plan signals and learning when a nondefault choice yields lower latency [1]. These works share the premise that a learned model can exploit execution history to guide per-query decisions.

Our work differs in representation: instead of encoding the physical plan solely as a tree, we explicitly transform SQL/plan/context information into an ontology/knowledge-graph and derive a matrix (and optional graph) embedding from that semantic layer. This design separates the what (semantics and provenance) from the how (choice of learner), allowing CNNs—as shown in our case study—and, in future, GCN/GAT or TreeConv, to operate on one consistent, query-level abstraction.

Orthogonal to the above, LLM-based tuning has been explored for database configuration and query improvement, e.g., Lambda-Tune [10] and GPTuner [20]. These methods focus on generating or selecting configurations and often do not expose full metrics/provenance for ablation. Our platform is complementary: the ontology captures relations among queries, plans, environments, and outcomes, enabling controlled ablations and explanations (e.g., which relations/features correlate with long-tail stalls). Pairing LLMs with ontology/KG for template discovery in this study can be a natural next step.

## VII. CONCLUSION

We presented OntoTune, an ontology-driven platform that records queries, plans, arms, environments, and outcomes as first-class objects and turns them into a reproducible knowledge base. From this layer, we derive a consistent embedding—matrix features (and an optional adjacency matrix) suitable for convolutional and graph models. A CNN case study demonstrated the end-to-end workflow, including data pipeline and arm selection, and showed that the learned policy can leverage prior executions to avoid disk- or memory-heavy plans and improve completion time under a fixed PostgreSQL configuration.

This work is a first step. The current ontology is intentionally minimal and results can still be sensitive to environment and parameter changes. Going forward, we will expand tail-aware features in the ontology/KG, apply environment-aware transfer learning for stability, and explore GCN/GAT/TreeConv as complementary learners. We also plan to pair the ontology with LLM-assisted template discovery and configuration suggestions, and to close the loop by automatically updating the knowledge base from feedback and scheduling new runs. The goal of this study is a systematic, explainable, and reproducible path to learned query optimization.

## REFERENCES

[1] R. Zhu, W. Chen, B. Ding, X. Chen, A. Pfadler, Z. Wu, and J. Zhou, "Lero: A learning-to-rank query optimizer," *arXiv preprint arXiv:2302.06873*, 2023.

[2] H. Gadde, "Intelligent query optimization: Ai approaches in distributed databases," *International Journal of Advanced Engineering Technologies and Innovations*, vol. 1, no. 2, pp. 650–691, 2024.

[3] J. Ortiz, M. Balazinska, J. Gehrke, and S. S. Keerthi, "Learning state representations for query optimization with deep reinforcement learning," in *Proceedings of the Second Workshop on Data Management for End-To-End Machine Learning*, 2018, pp. 1–4.

[4] Z. Yang, "Machine learning for query optimization," Ph.D. dissertation, University of California, Berkeley, 2022.

[5] Y. Zhao, G. Cong, J. Shi, and C. Miao, "Queryformer: A tree transformer model for query plan representation," *Proceedings of the VLDB Endowment*, vol. 15, no. 8, pp. 1658–1670, 2022.

[6] J. Heitz and K. Stockinger, "Join query optimization with deep reinforcement learning algorithms," *arXiv preprint arXiv:1911.11689*, 2019.

[7] R. Marcus, P. Negi, H. Mao, N. Tatbul, M. Alizadeh, and T. Kraska, "Bao: Making learned query optimization practical," in *Proceedings of the 2021 International Conference on Management of Data*, 2021, pp. 1275–1288.

[8] V. Leis, A. Gubichev, A. Mirchev, P. Boncz, A. Kemper, and T. Neumann, "How good are query optimizers, really?" *Proceedings of the VLDB Endowment*, vol. 9, no. 3, pp. 204–215, 2015.

[9] R. Marcus, "Stack (so_pg) dataset for postgresql," https://rmarcus.info/stack.html, 2021, postgreSQL 12/13 archives and 5000+ queries.

[10] V. Giannakouris and I. Trummer, "λ-tune: Harnessing large language models for automated database system tuning," *Proceedings of the ACM on Management of Data*, vol. 3, no. 1, pp. 1–26, 2025.

[11] B. Ding, S. Das, R. Marcus, W. Wu, S. Chaudhuri, and V. R. Narasayya, "Ai meets ai: Leveraging query executions to improve index recommendations," in *Proceedings of the 2019 International Conference on Management of Data*, 2019, pp. 1241–1258.

[12] L. Ma, B. Ding, S. Das, and A. Swaminathan, "Active learning for ml enhanced database systems," in *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, 2020, pp. 175–191.

[13] B. Ben Mahria, I. Chaker, and A. Zahi, "A novel approach for learning ontology from relational database: from the construction to the evaluation," *Journal of Big Data*, vol. 8, no. 1, p. 25, 2021.

[14] U. Simsek, E. Kärle, K. Angele, E. Huaman, J. Opdenplatz, D. Sommer, J. Umbrich, and D. Fensel, "A knowledge graph perspective on knowledge engineering," *SN Computer Science*, vol. 4, no. 1, p. 16, 2023.

[15] A. Ma'arif, W. Rahmaniar, H. I. K. Fathurrahman, A. Z. K. Frisky *et al.*, "Understanding of convolutional neural network (cnn): A review." *International Journal of Robotics & Control Systems*, vol. 2, no. 4, 2022.

[16] S. Zhang, H. Tong, J. Xu, and R. Maciejewski, "Graph convolutional networks: a comprehensive review," *Computational Social Networks*, vol. 6, no. 1, pp. 1–23, 2019.

[17] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.

[18] S. Lee, S. Ahn, and Y. Seo, "Relation modeling on knowledge graph for interoperability in recommender systems," in *Proceedings of the 38th ACM/SIGAPP Symposium on Applied Computing*, 2023, pp. 751–758.

[19] E. Rajabi and K. Etminani, "Knowledge-graph-based explainable ai: A systematic review," *Journal of information science*, vol. 50, no. 4, pp. 1019–1029, 2024.

[20] J. Lao, Y. Wang, Y. Li, J. Wang, Y. Zhang, Z. Cheng, W. Chen, M. Tang, and J. Wang, "Gptuner: An llm-based database tuning system," *ACM SIGMOD Record*, vol. 54, no. 1, pp. 101–110, 2025.