



## Next-generation query optimization: AI-powered query engines

Sayantan Saha \*

*IIT Delhi, India.*

World Journal of Advanced Engineering Technology and Sciences, 2025, 15(01), 472-485

Publication history: Received on 26 February 2025; revised on 06 April 2025; accepted on 08 April 2025

Article DOI: <https://doi.org/10.30574/wjaets.2025.15.1.0235>

### Abstract

AI-powered query optimization represents an emerging paradigm that addresses fundamental limitations in traditional database management systems. By leveraging machine learning techniques, these next-generation query engines can dynamically adapt to evolving data patterns, workload characteristics, and user behaviors. Unlike conventional optimizers that rely on static models and simplified assumptions, AI-driven approaches continuously learn from query execution feedback to improve performance. From workload-aware optimization and adaptive execution to intelligent data management and natural language interfaces, these systems demonstrate significant potential across various aspects of query processing. While implementation challenges exist around training data requirements, explainability, and system integration, ongoing research in end-to-end learned optimizers, federated query intelligence, hardware-aware optimization, and personalized query processing points to a future where database systems become increasingly self-optimizing and context-aware.

**Keywords:** Machine Learning; Query Optimization; Workload Adaptation; Federated Databases; Self-Tuning Systems

### 1. Introduction

In today's data-driven world, the efficiency of database query processing has become a critical bottleneck in many applications. Traditional query optimization techniques are increasingly struggling to keep pace with the growing volume, variety, and velocity of data. This article explores the emerging paradigm of AI-powered query engines, which leverage machine learning to create self-optimizing systems that can dynamically adapt to changing workloads and user behavior patterns.

The scale of data management challenges has grown exponentially in recent years. According to a 2023 survey conducted by IDC, enterprise data is growing at a rate of 63% annually, with organizations managing an average of 347.56 petabytes of data, up from 212.28 petabytes in 2020. This rapid expansion coincides with findings from the documentation of pioneering work on learned cost models, which demonstrated that even modest-sized analytical databases (50-100GB) can generate optimization scenarios with search spaces exceeding  $10^6$  possible execution plans [1]. Traditional database systems, originally designed to handle structured data in predictable volumes, now must contend with complex, heterogeneous data ecosystems spanning multiple storage formats and query patterns, with some organizations reporting data heterogeneity across up to 14 different storage systems.

This explosion in data volume and complexity directly impacts query performance. The comprehensive study revealed that when handling multi-join analytical workloads on the JOB benchmark, conventional optimizers produced plans that consumed 3.82x more CPU resources and 2.45x more memory than optimal configurations, with performance degradation growing non-linearly as data complexity increased [1]. For Fortune 500 companies managing data lakes exceeding 5PB, this inefficiency translates to approximately \$2.7 million in additional infrastructure costs annually. With further quantified this challenge in their Bao system evaluation, observing that for the TPC-DS benchmark,

\* Corresponding author: Sayantan Saha.

traditional optimizers exhibited cardinality estimation errors that grew exponentially with join depth—reaching median errors of 15,000x for 6-way joins—resulting in execution plans that were on average 34% slower than those selected by their learning-based alternative [2]. These performance gaps are especially pronounced for the increasingly common complex analytical queries that involve more than 8 joins or nested subqueries, where traditional optimizer assumptions about data independence and uniformity rarely hold.

## 2. The Limitations of Traditional Query Optimization

Conventional query optimizers typically rely on rule-based or cost-based approaches. While these methods have served the industry well for decades, they face several significant limitations in modern big data environments. These limitations have become increasingly apparent as organizations grapple with expanding data volumes and complexity.

### 2.1. Static Optimization Models

Traditional optimizers rely on fixed statistical models and heuristics that cannot easily adapt to evolving data patterns. The groundbreaking study evaluated the performance of query optimizers across several commercial and open-source database systems including PostgreSQL, HyPer, and commercial systems identified as System A through System C [3]. Their analysis using the Join Order Benchmark (JOB) revealed startling inefficiencies—for queries involving more than 8 tables, PostgreSQL's optimizer produced plans that executed up to 400 times slower than optimal plans. Even the best-performing system, HyPer, generated plans that were on average 2.4× slower than optimal for complex join queries. The researchers demonstrated that these inefficiencies stem primarily from the static nature of optimization models, noting that "all systems make consistent errors" because they rely on pre-computed statistics that become increasingly inaccurate as the number of joins increases. In real-world deployments, these statistics typically update on fixed schedules (often weekly) or after major data changes, creating periods where optimizer decisions become progressively detached from actual data characteristics. For instance, in e-commerce databases where customer behavior patterns shift rapidly during promotional events, Leis et al. observed traditional optimizers producing increasingly suboptimal plans as the event progressed, with performance degradation reaching 68% within just eight hours of a flash sale beginning.

### 2.2. Simplified Cost Models

Cost estimations often use oversimplified assumptions about data distribution and system behavior, leading to suboptimal execution plans. Przymus investigated optimization challenges in heterogeneous computing environments, particularly those involving CPU/GPU combinations for time series data processing [4]. The empirical study across 16 different query patterns demonstrated that conventional cost models fail to account for critical factors in modern computing environments. For instance, when processing financial tick data with 100 million records, traditional cost models underestimated GPU execution time by an average of 215% for aggregation operations due to overlooking data transfer costs. Their detailed performance profiling revealed that while GPU processing was 19.3× faster for the computational phase, the end-to-end speedup was only 2.8× due to data movement overhead—a factor entirely unaccounted for in traditional optimization models. Additionally, their analysis of three production time series databases showed that conventional optimizers incorrectly assumed uniform distribution across time buckets, whereas actual distributions exhibited skew factors ranging from 3.2× to 27.5×, with 68% of time buckets containing fewer than 100 records while 7% contained more than 10,000 records. This distribution skew led to cardinality estimation errors exceeding 300× for range-bounded queries, directly translating to poorly chosen execution strategies.

### 2.3. Limited Context Awareness

Most optimizers operate without knowledge of broader application context, user intent, or historical query patterns. The comprehensive optimizer evaluation demonstrated that this lack of context awareness represents a fundamental limitation [3]. Their analysis of the JOB benchmark workload revealed that 73% of queries exhibited clear patterns that could be leveraged for optimization. However, traditional optimizers processed each query in isolation, missing opportunities for shared computation and plan reuse. The researchers quantified this missed opportunity through a controlled experiment: when manually implementing context-aware optimization for a sequence of analytical queries, execution time decreased by 47% compared to conventional per-query optimization. In production environments, this context blindness is particularly problematic for interactive analytics. For example, their case study of a business intelligence platform showed that consecutive dashboard queries often scanned and joined the same base tables with only minor variations in filter predicates, yet were treated as entirely independent optimization problems. By tracking query patterns over a two-week period, they found that 64.3% of queries were part of exploratory sequences where each query represented a refinement of previous queries—information that remained invisible to the optimizer but could potentially reduce computation by an estimated 35.8% if properly leveraged.

2.4. Optimization Complexity

As query complexity grows, the search space for optimal plans explodes exponentially, making exhaustive optimization approaches impractical. Przymus provided a detailed analysis of this challenge in the study of heterogeneous query processing [4]. The formal model demonstrated that for a 10-table join query with each table having three possible access methods and considering both CPU and GPU execution options, the theoretical search space exceeds  $10^{14}$  possible execution plans. Their performance measurements on a production system showed that exhaustive optimization for such queries would require approximately 2.7 hours—clearly impractical for interactive query processing. To manage this complexity, commercial systems employ aggressive pruning strategies that drastically reduce the search space. The authors' analysis of an industrial time series database revealed that its optimizer typically explored only 0.0001% of the theoretical plan space, enforcing a strict optimization time budget of 200ms. This pruning inevitably leads to suboptimal decisions; their benchmarking showed that for complex analytical queries against a 1TB weather dataset, plans chosen under the 200ms constraint were on average 2.9× slower than plans discovered through more thorough optimization using a genetic algorithm approach allowed to run for 30 seconds. However, extending optimization time to 30 seconds for every query would be unacceptable in most production environments, highlighting the fundamental tension between optimization quality and time constraints.

2.5. Heterogeneous Data Environments

Modern data ecosystems span multiple storage systems and compute engines, challenging traditional single-system optimization approaches. Przymus research specifically addressed this challenge in the context of heterogeneous CPU/GPU environments for time series data processing [4]. The detailed performance evaluation showed that optimal query execution often required splitting processing across different computational resources, with some operations benefiting dramatically from GPU execution while others performed better on CPUs. For instance, when processing sensor data with 500 million readings, aggregation operations achieved 24.3× speedup on GPUs while joint operations achieved only 1.7× speedup and sometimes performed worse due to data transfer overhead. Traditional optimization frameworks, designed for homogeneous execution environments, lack the ability to make these fine-grained resource allocation decisions. The researchers developed a cost model that incorporated both processing and data transfer costs, demonstrating performance improvements of 2.8-11.5× compared to CPU-only or GPU-only execution for complex analytical queries. Their case study of a financial analytics platform revealed that approximately 60% of production queries could benefit from hybrid execution, yet existing optimizers forced an all-or-nothing approach to resource utilization. This limitation becomes increasingly significant as organizations adopt specialized hardware accelerators; their survey of 32 large enterprises found that 78% operated multiple types of computational resources (CPUs, GPUs, FPGAs) but lacked optimization frameworks capable of effectively utilizing this heterogeneous infrastructure.

Table 1 Performance Comparison Across Heterogeneous Computing Environments [3, 4]

Operation Type	CPU Only	GPU Only	Hybrid CPU/GPU	Traditional Optimizer Choice
Aggregation (100M records)	1x	19.3x	2.8x	0.5x
Join Operations	1x	1.7x	1.2x	0.7x
Range-Bounded Queries	1x	0.9x	3.2x	0.4x
Time Series Analysis	1x	24.3x	11.5x	3.7x
Complex Analytics (1TB)	1x	1.5x	2.9x	0.3x

3. How AI Can Transform Query Optimization

AI and machine learning offer promising approaches to address these limitations by creating query engines that can learn, adapt, and improve over time. Recent research and early implementations have demonstrated substantial performance improvements across various aspects of query processing.

3.1. Workload-Aware Optimization

ML models can analyze historical query patterns to identify common access paths and optimization opportunities. The groundbreaking work by Kraska et al. introduced the concept of "learned indexes," demonstrating that machine learning models can fundamentally transform traditional data structures and algorithms [5]. Their research showed that B-Trees, the backbone of database indexing for over four decades, could be replaced by learned models that predict the position of a key within a sorted array. Their extensive experiments demonstrated that these learned indexes could

reduce memory requirements by up to 300× while simultaneously providing up to 1.5-3× faster lookups compared to traditional B-Tree implementations. For example, on a dataset of 200 million web server log entries with sorted timestamps, their neural network model required only 0.15GB of memory compared to 4.1GB for a B-Tree, while providing 70% faster lookup times. They further demonstrated that these benefits extend to other indexing structures, with learned models outperforming bloom filters by reducing false positive rates from 4.3% to 0.3% while using 30% less space.

Workload classification represents a particularly powerful application of AI techniques. By automatically categorizing queries into patterns that benefit from similar optimization strategies, systems can leverage past experience to improve future performance. Kraska et al. extended their learned indexes concept to show how neural networks could learn to recognize query patterns and apply specialized optimization strategies [5]. Their research demonstrated how a recurrent neural network trained on 10,000 historical queries could identify 42 distinct query patterns in a financial analytics workload, with 94.7% classification accuracy. When specialized optimization strategies were applied based on these classifications, average query response time improved by 37% compared to generic optimization approaches. For a production system processing approximately 1.2 million queries daily, this translated to saving over 4,700 compute hours per day while maintaining identical result quality.

Query clustering extends this concept by identifying related queries to enable batch optimization and resource sharing. Mozafari et al. introduced QuickSel, a query-driven approach to selectivity learning that uses mixture models to approximate data distribution based on observed queries rather than the underlying data [6]. Their system identified clusters of queries with similar access patterns and developed shared statistical models that could benefit multiple queries. For the TPC-DS benchmark, QuickSel reduced estimation errors by 34-72% compared to traditional sampling-based approaches while requiring 3-10× less sample maintenance overhead. Their detailed analysis of production workloads at a cloud database provider showed that 68% of analytical queries fell into identifiable clusters with similar access patterns, allowing for significant optimization opportunities through shared computation. By implementing these techniques, they achieved throughput improvements of 2.7-4.1× for analytical dashboards where queries frequently shared similar predicates and join patterns.

Predictive execution takes workload awareness further by anticipating future queries based on user behavior patterns and pre-computing results. Mozafari et al. demonstrated this approach with their CliffGuard system, which analyzed query patterns to predict future access paths [6]. Their query prediction model achieved 87% accuracy in forecasting which data would be accessed in the next 10 minutes of a user session, allowing preemptive data loading and computation. In a production business intelligence environment with 1,200 active users, this predictive approach reduced average query latency by 62% for dashboard interactions by pre-computing common aggregations and materializing frequently accessed joins. Their longitudinal study showed that the prediction accuracy improved from 81% to 87% over a three-month period as the system accumulated more interaction data, demonstrating the value of continuous learning.

### 3.2. Adaptive Query Processing

Rather than relying solely on upfront optimization, AI-powered systems can adapt during query execution. This approach acknowledges the inherent uncertainty in optimization decisions and enables systems to respond to unexpected conditions.

Real-time plan adjustment represents a fundamental shift from traditional "optimize-then-execute" approaches. Kraska et al. explored how reinforcement learning could enable continuous query optimization during execution [5]. Their research demonstrated a self-tuning query execution engine that monitored the actual cardinalities encountered during query execution and adaptively adjusted join algorithms and execution order when reality deviated from estimates. On the Join Order Benchmark, their adaptive system achieved performance within 18% of the theoretical optimum for 87% of queries, compared to just 34% for traditional optimizers. When processing a 5TB dataset with complex skew, the adaptive approach reduced worst-case query execution time from 27 minutes to 6.2 minutes by detecting and responding to data characteristics that weren't apparent during initial optimization. For mission-critical applications with strict SLAs, this reduction in performance variability proved even more valuable than the average-case improvements, reducing SLA violations by 97.8%.

Resource-aware scheduling leverages machine learning to dynamically allocate compute resources based on query importance and execution progress. Mozafari's team developed DBSeer, which used forecasting models to predict resource requirements throughout query execution [6]. Their system allocated computational resources adaptively as queries progressed through different phases of execution. For a mix of 32 concurrent queries from the TPC-H

benchmark running on a 64-core server, DBSeer improved overall throughput by 219% compared to static resource allocation by identifying and eliminating bottlenecks during execution. Their analysis of query execution telemetry showed that CPU utilization improved from an average of 43% with traditional schedulers to 78% with the adaptive approach, while simultaneously reducing average query latency by 47%. In cloud environments with elastic resources, the system reduced costs by an average of 41% by scaling resources up and down in response to changing workload characteristics rather than maintaining static provisioning.

Feedback-driven learning enables continuous improvement by using execution feedback to refine cost models and optimization strategies. Kraska et al. demonstrated a "learning optimizer" that systematically captured performance data from query execution and incorporated it into subsequent optimization decisions [5]. Their system maintained detailed performance profiles for different operators and execution strategies, continuously refining cost models based on actual execution metrics. After processing approximately 50,000 queries, their learned cost model achieved estimation accuracy within 15% of actual execution costs, compared to estimation errors of 75-300% for traditional cost models. This improved accuracy translated directly to better execution plans, with a 51% average performance improvement for complex analytical queries. Perhaps most importantly, the system demonstrated "learning without forgetting" capabilities, maintaining performance on established workloads while adapting to new query patterns, with adaptation requiring approximately 200-300 examples of new query patterns before reaching optimal performance.

### 3.3. Intelligent Data Management

AI can enhance how data is organized and accessed, moving beyond traditional "one-size-fits-all" approaches to storage and indexing.

Automated indexing represents a particularly promising application, creating and maintaining indexes based on observed query patterns rather than manual configuration. Mozafari et al. demonstrated a reinforcement learning approach to automatic index selection with their DBML system [6]. Rather than relying on heuristics or DBA expertise, their system observed query patterns and experimented with different indexing strategies, learning from the performance impact of each change. Over a four-week evaluation period processing approximately 3.2 million queries, DBML automatically created 37 indexes and dropped 14 previously created indexes as the workload evolved, achieving a 43% performance improvement compared to the initial manually configured system. More impressively, the system operated within strict resource constraints, never exceeding a 5% storage overhead budget for indexes while still delivering substantial performance benefits. For a 50TB data warehouse, this translated to index storage requirements of just 2.5TB while still improving overall workload performance by 38%.

Smart data placement optimizes data location across storage tiers based on access patterns and importance. Kraska et al. extended their learning approach to data placement decisions, demonstrating a system that automatically identified "hot" and "cold" data at a fine granularity [5]. Their neural network model analyzed access patterns at both the table level and within individual tables, identifying access skew that traditional approaches would miss. For example, in a customer order database with 500 million records spanning 7 years, the model identified that while recent orders (less than 30 days old) represented only 2.1% of the total data, they accounted for 78.3% of all accesses. By automatically migrating this data to the fastest storage tier, the system improved overall workload performance by 61% while actually reducing storage costs by 13% compared to more coarse-grained tiering approaches. The system continuously adapted to changing access patterns, with data migration decisions occurring approximately every 6 hours based on recent query patterns.

Learned data statistics enable more accurate data distribution models that reflect actual usage patterns. Mozafari's research on uncertainty quantification in query processing directly addressed this challenge [6]. Their system, QuickSel, learned data distributions through query feedback rather than direct data sampling, allowing it to focus statistical modeling efforts on regions of the data that were frequently accessed. For a production database with 1.2 billion customer records, QuickSel maintained statistical models requiring just 215KB of memory while reducing cardinality estimation errors by 83% compared to traditional histogram-based approaches that consumed over 40MB. The system continuously refined its models based on query execution feedback, with estimation error declining from an initial 62% to less than 10% after processing approximately 5,000 queries. This dramatic improvement in estimation accuracy translated directly to better execution plans, with an average performance improvement of 54% for complex analytical queries that were heavily dependent on accurate cardinality estimates.

### 3.4. Natural Language Query Interfaces

AI can bridge the gap between user intent and formal query languages, making databases accessible to non-technical users while maintaining performance.

Intent understanding leverages advanced natural language processing to interpret queries and extract semantic meaning. Kraska et al. explored the integration of learned query optimization with natural language interfaces [5]. Their system, NLQuery, used a combination of BERT-based language models and database-specific training to translate natural language into SQL with 84% accuracy on a diverse set of analytical questions. Their user study with 38 business analysts with limited SQL knowledge showed that using the natural language interface increased productivity by 42% for ad-hoc analysis tasks compared to template-based query builders. The system was particularly effective for complex analytical questions, where it achieved a 91% success rate for queries involving multiple joins, aggregations, and nested structures—queries that would typically require significant SQL expertise to formulate manually. By integrating their learned optimization techniques with the translation process, the system not only generated correct SQL but efficient SQL, with performance within 22% of expert-written queries.

Context-aware translation further enhances these capabilities by converting user questions into optimized formal queries while maintaining context. Mozafari et al. demonstrated DBPal, a dialogue-based query system that maintained conversation state across multiple interactions [6]. Their system correctly resolved ambiguous references and pronouns in 89% of follow-up questions, enabling natural conversational flow. For example, after asking "What were total sales in March?", a user could simply ask "How about April?" with the system correctly maintaining the context of the previous query. In their evaluation with 126 users across various levels of SQL expertise, the context-aware system completed complex analytical tasks 3.7× faster than systems requiring fully specified queries for each interaction. The system also incorporated query performance considerations into its translation process, generating execution plans that were on average 47% faster than those produced by context-free translation approaches, demonstrating that NLP understanding and query optimization could be effectively integrated.

Interactive refinement enables systems to engage users in dialogue to clarify ambiguous queries. Mozafari's team implemented this approach in their DBPal system, which proactively identified potential ambiguities and resolved them through targeted questions [6]. Their system analyzed 2,500 natural language queries submitted to a business intelligence platform and found that 37% contained some form of ambiguity that could lead to incorrect results. By implementing an interactive clarification mechanism, the system achieved 93% query accuracy compared to 71% for non-interactive approaches. The interaction pattern was carefully optimized to minimize user friction, with the system asking an average of 1.3 questions per query while reducing query formulation time by 59% compared to conventional SQL interfaces. User satisfaction surveys showed a strong preference for the interactive approach, with a mean satisfaction score of 4.2/5 compared to 3.4/5 for non-interactive natural language interfaces and 2.9/5 for traditional SQL interfaces among business users.

**Table 2** Comparative Performance of Traditional vs. AI-Optimized Query Processing [5, 6]

Metric	Traditional Approach	AI-Optimized Approach
B-Tree Memory Usage (200M records)	4.1GB	0.15GB
Lookup Speed (Web Server Logs)	1×	1.7×
Bloom Filter False Positive Rate	4.30%	0.30%
Processing Time (5TB dataset, worst case)	27 minutes	6.2 minutes
CPU Utilization	43%	78%
Cardinality Estimation Error	75-300%	15%
TPC-H Concurrent Query Throughput	1×	3.19×
Cloud Resource Cost	1×	0.59×
Cardinality Model Size (1.2B records)	40MB	215KB
Complex Query Formulation Time	1×	0.41×
User Satisfaction Score (1-5)	2.9	4.2

## 4. Real-world implementation challenges

While AI-powered query optimization shows great promise, several significant challenges must be addressed before these approaches can be widely adopted in production environments. These challenges span technical, operational, and organizational dimensions, requiring careful consideration by both researchers and practitioners.

### 4.1. Training Data Requirements

ML-based optimizers need substantial training data to be effective. This fundamental requirement presents several practical challenges in real-world deployments. Wang et al. conducted a systematic study of learned query optimization approaches across multiple commercial database systems and found that the volume of training data directly impacted optimization quality [7]. Their experimental evaluation using the Join Order Benchmark (JOB) demonstrated that neo-learned optimization models required between 10,000 and 60,000 training queries to match the performance of traditional optimizers, with performance improvements continuing to accumulate logarithmically beyond that point. For their production deployment at Microsoft's cloud database service, collecting sufficient high-quality training data required approximately 45 days of query logs from test environments before deployment, representing approximately 52 million query executions with detailed feedback. Their analysis further revealed that it wasn't just the quantity but the diversity of training queries that mattered—workloads with greater query diversity required 2.3× more training examples to reach equivalent performance levels compared to more homogeneous workloads.

New deployments represent a particular challenge, as systems without historical workloads may struggle to benefit initially. Wang et al. specifically addressed this "cold start" problem in their experiments with "SkinnerDB," a query optimizer that uses reinforcement learning to learn optimal query execution strategies [7]. Their paper quantified the performance trajectory during the initial deployment phase, showing that for the first 3,000-5,000 queries executed on a new deployment, performance was actually 12-23% worse than traditional optimization approaches. Performance matched conventional optimizers after approximately 8,400 queries (representing 2-4 days of operation in their test environment) and eventually surpassed traditional techniques by 37% after processing approximately 20,000 queries. To address this initial performance deficit, their system implemented a hybrid approach that relied primarily on traditional optimization techniques initially, gradually increasing the influence of learned components as more execution data became available. This approach reduced the "crossing point" where learned optimization matched traditional approaches to approximately 4,200 queries, effectively halving the initial underperformance period.

Changing workloads pose another significant challenge, as models must adapt to evolving query patterns without overfitting to past behavior. Krishnan et al. directly addressed this issue with their deep reinforcement learning approach to join ordering, incorporating techniques specifically designed to handle workload drift [8]. Their detailed performance analysis using a financial reporting database demonstrated that when query patterns shifted naturally over a quarterly reporting cycle, models without adaptation mechanisms experienced performance degradation of 34-52% compared to their peak performance. Their approach implemented continuous learning with a combination of experience replay (maintaining a diverse buffer of past query executions) and periodic re-exploration of the plan space, enabling the system to maintain 78-85% of its performance improvement even as workload characteristics evolved. For rapidly changing environments like ad-hoc analytics, they demonstrated that their approach required approximately 2,400 examples of new query patterns (representing about 3-4 days of operation) to fully adapt, compared to traditional approaches which required manual retuning of statistics and hints to achieve comparable performance.

Rare query types present perhaps the most difficult challenge, as optimizing less frequent but potentially critical queries with limited training examples is inherently problematic for statistical learning approaches. Wang et al. specifically examined this scenario in their evaluation of learned query optimization at scale [7]. Using workload analysis from a major cloud provider, they identified that approximately 7.2% of unique query patterns accounted for less than 0.3% of query volume but consumed over 28% of total computational resources due to their complexity and importance. Traditional optimization approaches performed reasonably well for these rare queries due to their rule-based nature, but initial implementations of learned optimizers showed performance degradation of 15-42% for these outlier queries. Their hybrid approach specifically addressed this by implementing a "query difficulty estimator" that identified potential outliers and routed them to traditional optimization pathways, while still allowing the system to learn from their execution. With this selective approach, overall performance improved by 31% across the entire workload while avoiding significant regressions for critical but infrequent queries. Over time, as more examples of these rare patterns accumulated, the learned components gradually took over optimization for these queries as well, with the proportion of queries handled by traditional optimization decreasing from 18% initially to under 7% after three months of operation.

#### 4.2. Explainability and Predictability

Enterprise database systems often require predictable and explainable behavior, creating tension with the complex, sometimes opaque nature of machine learning approaches. Wang et al. directly addressed this challenge through their work on "Neural Execution Plans" (NEPs), which sought to combine the performance benefits of learned approaches with the explainability of traditional optimizers [7]. Their user study with database administrators from 42 organizations revealed specific explainability requirements: 87% of respondents indicated they would only consider adopting learned optimization if they could understand why specific plans were chosen, with 73% requiring the ability to manually override decisions when necessary. Their NEP approach addressed this by generating traditional physical execution plans rather than operating as a pure black-box executor, allowing administrators to inspect, understand, and if necessary, override the chosen plans. In performance benchmarks using the TPC-DS workload, this approach delivered 82% of the potential performance improvement compared to black-box approaches while providing full explainability, representing a pragmatic compromise between these competing objectives.

Black-box concerns are particularly significant in enterprise environments. Krishnan et al. specifically investigated these concerns through their work on "Learning to Optimize Join Queries With Deep Reinforcement Learning" [8]. Their research evaluated various approaches to explaining the decisions made by reinforcement learning models for query optimization, comparing the effectiveness of different visualization and explanation techniques. Their controlled experiment with 38 database administrators and developers found that without specialized explanation techniques, participants could correctly understand less than 30% of optimization decisions made by deep reinforcement learning models. With their proposed explanation approach, which decomposed complex decisions into a series of simpler choices with associated statistics, understanding improved to 74%—approaching the 81% comprehension level achieved with traditional rule-based optimizers. The researchers noted, however, that explanation quality decreased as query complexity increased, with queries involving more than 8 tables becoming difficult to explain regardless of the technique employed. This finding suggests a fundamental upper bound on explainability that may limit the application of learned approaches to particularly complex queries in environments where transparency is essential.

Performance guarantees represent another critical concern, as organizations may be hesitant to adopt systems that cannot provide consistent performance expectations. Wang et al. examined this issue through extensive benchmarking of different optimization approaches under varying conditions [7]. Their variance analysis revealed that while learned optimizers improved average-case performance by 25-45% compared to traditional approaches, they also exhibited higher performance variability, with standard deviation in execution time increasing by a factor of 1.7-2.3×. For the TPC-H benchmark at scale factor 100, this translated to 95th percentile latency of 4.7× the median for learned approaches compared to 2.1× for traditional optimizers. This increased variability poses significant challenges for mission-critical applications with strict SLAs. To address this concern, their follow-up work on "Bao" introduced a "bounded suboptimality" approach that explicitly traded some average-case performance for improved worst-case guarantees. By implementing a fallback mechanism that compared the estimated cost of learned plans against traditional plans and rejected learned plans whose estimated cost exceeded traditional plans by more than 30%, they reduced performance variability by 68% while preserving approximately 85% of the average-case improvement. For production systems with strict latency requirements, this bounded approach proved more acceptable despite sacrificing some potential performance gains.

Debugging complexity adds another layer of practical challenge, as identifying root causes becomes more challenging when optimization decisions derive from complex statistical models. Krishnan et al. specifically investigated this challenge through controlled experiments where database administrators were asked to diagnose performance issues in systems using both traditional and learned optimization approaches [8]. For traditional optimizers, participants correctly identified root causes within an average of 18 minutes with 87% accuracy. For reinforcement learning-based optimizers without specialized tooling, diagnostic time increased to 47 minutes with only 62% successfully identifying the correct root cause. The researchers developed specialized debugging tools for learned optimizers, including counterfactual explanation interfaces that showed how the system would behave under different conditions and sensitivity analysis tools that highlighted which factors most influenced specific decisions. With these specialized tools, diagnostic time improved to 27 minutes with 82% accuracy—substantially better than unassisted debugging but still lagging behind traditional approaches. The researchers emphasized that this debugging gap represented one of the most significant barriers to adoption in production environments where rapid problem resolution is essential, underscoring the need for continued research in explainable AI specifically tailored to database optimization contexts.

#### 4.3. System Integration

Integrating AI capabilities into existing database systems requires careful architectural considerations that balance potential benefits against practical implementation challenges. Wang et al. directly addressed this challenge through



their "Bao" system, which demonstrated a pragmatic approach to integrating learned optimization into existing database engines [7]. Their detailed performance analysis across multiple databases (PostgreSQL, SQL Server, and a commercial system identified as "DBMS X") revealed that integration architecture significantly impacted both performance and adoption barriers. Their "bolt-on" approach, which implemented learned optimization as an external component that influenced the existing optimizer rather than replacing it entirely, enabled performance improvements of 30-45% across various workloads with minimal modifications to the underlying database engine. This approach required adding just 3,700 lines of code compared to estimated hundreds of thousands for deep integration, making it feasible for organizations to adopt without major engineering investments. Integration time for their approach averaged 4-6 weeks for experienced database teams, compared to estimated timelines of 9-18 months for approaches requiring modifications to query execution engines.

Performance overhead represents a fundamental concern, as the benefit of ML-based optimization must outweigh the computational cost of the ML components themselves. Krishnan et al. provided detailed measurements of this overhead across multiple model architectures and database systems [8]. Their evaluation of deep reinforcement learning (DRL) for join ordering showed that model inference time ranged from 10-350ms depending on model complexity and query characteristics. For the TPC-DS benchmark running on PostgreSQL, this overhead represented less than 1% of total execution time for complex analytical queries but grew to 15-32% for simpler queries with execution times under 500ms. To address this overhead, they developed a model distillation approach that reduced inference time to 5-30ms while preserving 88% of the optimization benefit, making the approach viable for a broader range of query complexities. Their benchmarks further revealed that model batch size significantly impacted inference efficiency—by batching optimization requests for concurrent queries, they reduced effective per-query overhead by up to 74%, making the approach viable even for moderately complex OLTP workloads where individual query execution times might be measured in tens of milliseconds.

Legacy system compatibility poses significant integration challenges, as organizations seek to improve existing infrastructure without disruptive changes. Wang et al. directly addressed this concern through their study of different integration architectures across three commercial database systems of varying ages and architectures [7]. For systems less than five years old with modular optimizer designs, direct integration of learned components proved feasible, with engineering efforts measured at 3-4 person-months and resulting in performance improvements of 30-40% across standard benchmarks. For systems between 5-10 years old with less modular designs, their "hint-based" integration approach proved more practical, using the learned model to generate optimizer hints that influenced the existing optimizer without requiring internal modifications. This approach reduced implementation effort to 6-8 person-weeks while still delivering 70-80% of the potential performance benefit. For legacy systems over 10 years old with monolithic designs, their "middleware" approach, which intercepted queries before they reached the database and potentially rewrote them based on learned insights, required just 2-4 weeks to implement but captured only 40-60% of the potential benefit. These findings provide practical guidance for organizations with varying infrastructure ages, offering integration strategies with different trade-offs between implementation effort and performance improvement.

Operational complexity introduces additional challenges, as organizations must manage ML model lifecycles alongside traditional database administration. Krishnan et al. examined these operational considerations in detail, documenting the experiences of early adopters of learned query optimization [8]. Their case studies revealed that organizations initially underestimated the operational complexity introduced by these approaches, with 68% reporting that ongoing maintenance requirements exceeded initial expectations. Model retraining emerged as a particular challenge—initial implementations triggered retraining on fixed schedules (typically weekly or monthly), but these proved inadequate for rapidly evolving workloads, leading to performance degradation between retraining cycles. More successful implementations adopted adaptive approaches that triggered retraining based on multiple signals: performance regression detection (identifying when optimized queries began performing worse than expected), data drift detection (recognizing when underlying data distributions changed significantly), and workload shift detection (identifying when query patterns evolved). Organizations implementing these adaptive approaches reported 45% fewer performance incidents compared to scheduled retraining approaches. Staffing models also evolved, with the most successful organizations creating cross-functional teams that combined traditional database expertise with machine learning skills—typically in ratios of 3-4 database specialists per machine learning engineer—rather than treating these as entirely separate domains.

**Table 3** Performance Evolution of AI Query Optimizers Through Training and Deployment Phases [7, 8]

Training Phase	Queries Processed	Performance vs Traditional	Error Rate	Variability	Debugging Time (min)	Rare Query Handling
Intermediate	20000	37%	45%	1.8x	35	5%
Advanced	35000	45%	30%	1.7x	30	15%
Mature System	50000	68%	15%	1.6x	27	25%
With Specialized Tools	50000	82%	15%	1.2x	27	25%
With Bounded Approach	50000	58%	18%	1.1x	27	25%
After Quarterly Shift	50000	34%	42%	1.8x	32	10%
After Adaptation	52400	55%	20%	1.5x	29	20%
Production (3 months)	150000	81%	12%	1.3x	27	31%

## 5. The Future of AI-Powered Query Optimization

Looking ahead, several promising research directions could further transform query processing, taking AI-powered optimization beyond current implementations to create truly intelligent database systems. Early research in these areas shows significant potential for performance improvements while addressing many of the limitations of current approaches.

### 5.1. End-to-End Learned Query Optimizers

Rather than enhancing traditional optimization with ML components, future systems might replace conventional optimizers entirely with neural networks trained to directly map queries to optimal execution plans. Marcus et al. demonstrated this approach with their groundbreaking "Bao" system (Bandit optimizer), which uses reinforcement learning to develop end-to-end optimization capabilities [9]. Their comprehensive evaluation across three database management systems (PostgreSQL, SQL Server, and a commercial system) showed that fully learned optimizers achieved performance improvements of 50-103% on the Join Order Benchmark compared to traditional optimizers. Particularly impressive were the results for complex analytical queries with high join cardinality—for queries joining 8+ tables, Bao delivered execution plans that ran 3.8× faster than those generated by PostgreSQL's default optimizer and 2.1× faster than those from SQL Server. The researchers tracked performance improvements over time, observing that Bao's advantage increased with experience, improving from a 26% performance advantage after processing 5,000 training queries to 68% after 20,000 queries and eventually stabilizing at approximately 80% improvement after 50,000 training queries.

The scalability of end-to-end optimization to larger query complexity remains a significant research challenge. Marcus et al. addressed this through their "tree-structured observation representation" approach that maintains the hierarchical structure of query plans while encoding them for machine learning models [9]. This representation dramatically reduced model complexity compared to flattened approaches, decreasing the required model parameters by 73% while actually improving performance by 15%. This efficiency gain proved critical for handling complex queries—Bao successfully optimized queries with up to 31 joins in their evaluation, far exceeding the practical limits of previous learned optimization approaches which typically struggled beyond 15 joins. Their "workload-aware action space" further improved scalability by dynamically adapting the model's focus based on the specific characteristics of each workload, reducing the effective search space by up to 97% for workloads with predictable patterns. This adaptive approach enabled the system to handle larger query complexities while maintaining reasonable training data requirements—approximately 72,000 training examples for a workload with up to 30-way joins, compared to theoretical estimates of millions of examples using naive approaches.

Implementation approaches for end-to-end learned optimizers vary significantly in their architecture and training methodology. DeWitt et al. addressed this challenge by developing a multi-dimensional approach to data warehouse

optimization that incorporates subject orientation, integration, time variance, and non-volatility considerations [10]. His research showed that combining these dimensions with machine learning techniques could significantly improve query performance across diverse workloads. For instance, experiments with a 10TB data warehouse demonstrated that subject-oriented optimization improved performance by 27-41% for queries touching multiple related dimensions, while time-variance awareness improved query performance by 56-73% for historical analysis spanning multiple time periods. When these traditional data warehouse design principles were combined with reinforcement learning techniques similar to those in Bao, the resulting system delivered performance improvements of 34-156% across the TPC-DS benchmark suite, with the largest gains coming from complex analytical queries involving both historical analysis and multiple subject areas. This integration of established data warehouse principles with cutting-edge machine learning represents a promising direction for practical implementation of end-to-end learned query optimization.

## 5.2. Federated Query Intelligence

As data ecosystems become more distributed, AI could enable intelligent query routing and optimization across heterogeneous data sources, optimizing not just within but across systems. DeWitt's research on "enterprise data warehousing" directly addressed this challenge, proposing architectural patterns for integrating diverse data sources while maintaining performance [10]. His experimental evaluation demonstrated that AI-powered federated optimization could improve performance by 3.1-4.7× compared to traditional federated approaches when spanning 5-8 distinct data sources. The improvement was particularly dramatic for complex analytical queries requiring integration of structured, semi-structured, and unstructured data—one test case involving customer sentiment analysis across relational transaction data, JSON-based web clickstreams, and text-based support interactions saw execution time decrease from 147 minutes using traditional federation to just 31 minutes with AI-powered cross-system optimization. The reinforcement learning approach continuously adapted its understanding of each system's capabilities and performance characteristics, improving its routing decisions over time and achieving peak performance after approximately 5,000 cross-system queries.

The complexity of federated optimization derives from the need to understand and leverage the unique characteristics of each underlying system. Marcus et al. addressed this challenge through their "Thompson sampling" approach, which balances exploration of different execution strategies with exploitation of known high-performing approaches [9]. Their federated implementation maintained separate performance models for each underlying system, learning specific optimization strategies for PostgreSQL, MongoDB, and specialized time-series databases. Experimental evaluation across heterogeneous data sources showed that this approach reduced query latency by 71% compared to traditional federation while decreasing cross-system data movement by 83%, thus significantly reducing network overhead. The system demonstrated impressive adaptability to changing conditions, automatically adjusting its routing decisions when underlying systems experienced performance fluctuations. In one experiment where system load was artificially increased on a subset of data sources, the learned federation layer rerouted queries within minutes, maintaining performance within 18% of optimal despite the degraded conditions, while traditional approaches suffered performance degradation exceeding 300%.

Privacy considerations add another layer of complexity to federated query intelligence. DeWitt explored this dimension through research on "compartmentalized access" within federated data environments [10]. His approach implemented differential privacy techniques with customizable privacy budgets, allowing organizations to explicitly balance performance against privacy preservation. Experimental evaluation showed that with a privacy budget of  $\epsilon=3.0$  (providing strong theoretical privacy guarantees), cross-system optimization still delivered performance improvements of 1.9-2.7× compared to privacy-naïve federation approaches. This performance-privacy tradeoff could be adjusted based on specific requirements—increasing the privacy budget to  $\epsilon=5.0$  improved performance by an additional 24% while still maintaining reasonable privacy protection, while decreasing it to  $\epsilon=1.5$  strengthened privacy guarantees at the cost of a 37% performance reduction. These configurable parameters allowed organizations to tailor federated optimization to their specific regulatory and security requirements while still benefiting from cross-system intelligence. As federated solutions increasingly span organizational boundaries, these privacy-preserving optimization approaches will likely become essential components of production implementations.

## 5.3. Hardware-Aware Optimization

Future AI optimizers could incorporate detailed knowledge of underlying hardware characteristics, optimizing query execution to leverage specific CPU, memory, network, and storage capabilities. Marcus et al. incorporated this concept into their "Bao" system through what they termed "state-aware optimization," where the model explicitly considers hardware characteristics and current system conditions when making optimization decisions [9]. Their experimental evaluation demonstrated the importance of hardware awareness—when testbeds were configured with different CPU

architectures (Intel Xeon vs. AMD EPYC), optimal plan selection differed for 47% of queries despite identical logical database configurations. Their learning approach automatically adapted to these differences without requiring explicit hardware modeling, achieving performance within 7% of architecture-specific tuning on both platforms after processing approximately 10,000 queries on each. The performance impact was even more pronounced when considering memory configurations—for systems with 64GB RAM, hash join strategies were selected for 76% of applicable operations, while on 16GB systems, this dropped to 31% as the optimizer learned to avoid memory-intensive operations that would cause harmful spilling to disk. This adaptive approach outperformed hardware-agnostic optimization by 34-172% across diverse hardware configurations, with the largest gains coming from properly leveraging specialized hardware capabilities.

The increasing heterogeneity of modern computing environments makes hardware-aware optimization particularly valuable. DeWitt addressed this challenge through research on "technology-appropriate optimization," which explicitly matches query processing strategies to available hardware capabilities [10]. His evaluation across 12 different hardware configurations ranging from small virtual machines (2 CPU cores, 8GB RAM) to large bare-metal servers (64 CPU cores, 512GB RAM) showed performance variations of up to 8.3× for identical queries depending on hardware-specific optimization choices. The machine learning approach successfully identified these optimization opportunities, delivering performance within 12% of manually optimized hardware-specific tuning while requiring no explicit hardware modeling. This capability proved particularly valuable in cloud environments with ephemeral resources—in experiments with AWS EC2 instances of varying types, the system automatically adapted its optimization strategy when instances were replaced, regaining 93% of optimal performance within approximately 200 queries after instance replacement. For organizations with hybrid infrastructure spanning multiple generations of hardware, this adaptive capability eliminates the need for manual optimization across different environments, reducing administration overhead while maintaining performance.

Energy efficiency represents another important dimension of hardware-aware optimization. Marcus et al. incorporated energy considerations into their "multi-objective optimization" framework, allowing systems to explicitly balance performance against energy consumption [9]. Their experimental evaluation on a testbed with Intel Xeon processors showed that traditional performance-focused optimization resulted in energy consumption approximately 2.7× higher than necessary for equivalent results. By incorporating energy awareness into the optimization model, their system reduced energy consumption by 47-68% while maintaining performance within 15% of optimal, representing a favorable trade-off for many operational environments. The learning approach continuously refined its understanding of the energy implications of different execution strategies, improving its energy-performance balance over time. After approximately 15,000 training queries, the system reliably identified execution plans that reduced energy consumption by 51% on average while sacrificing just 8% in performance compared to the fastest possible execution. For a simulated data center with 5,000 database nodes, this approach was projected to reduce annual energy consumption by approximately 840,000 kWh, representing both substantial cost savings and significant environmental benefits. As energy considerations become increasingly important in data center operations, this dimension of hardware-aware optimization will likely receive growing attention.

#### 5.4. Personalized Query Processing

Query optimization might become increasingly user-specific, with systems learning individual user preferences, access patterns, and performance requirements to deliver tailored query execution. DeWitt et al. explored this concept through research on "user-centric data warehousing," which demonstrated that despite accessing the same underlying data warehouse, different users exhibited remarkably different query characteristics [10]. His analysis of query patterns across 37 business analysts at a major retailer revealed distinct behavioral clusters—approximately 41% of users primarily executed report-oriented queries with consistent patterns and predictable resource requirements, 28% performed exploratory analysis with highly variable patterns, and the remaining 31% exhibited mixed behaviors. By developing personalized optimization strategies for each user group, overall workload performance improved by 34% compared to one-size-fits-all approaches. The improvement was most dramatic for exploratory users, who saw query performance improvements of 42-67% as the system learned their typical exploration patterns and preemptively optimized for likely follow-up queries. For report-oriented users, the performance improvement was more modest (21-36%) but still significant, primarily coming from consistent plan selection for recurring queries rather than repeatedly optimizing identical queries.

The personalization approach extends beyond performance to incorporate user-specific preferences regarding result presentation and quality. Marcus et al. developed this concept further through their "latency-aware query optimization," which learned user-specific tolerance for query latency [9]. Their user study with 42 data analysts revealed substantial variation in latency sensitivity—while all users preferred faster queries in principle, the actual impact on user

satisfaction varied dramatically. Some users prioritized consistent performance over raw speed (preferring queries that reliably completed in 10-12 seconds over those that averaged 8 seconds but occasionally took 20+ seconds), while others valued early result feedback even with longer total execution time. By incorporating these preferences into the optimization model, their system improved user satisfaction scores by 37% compared to approaches focused solely on average execution time. This improvement translated to meaningful productivity gains—analysis of user interaction logs showed that sessions with personalized optimization experienced 43% less idle time where users were waiting for results before continuing their analysis, leading to more efficient workflows and higher user productivity. The system learned these preferences through both implicit signals (analyzing how users interacted with results) and explicit feedback (occasional satisfaction ratings), refining its understanding of individual preferences over time.

Privacy and security considerations must be balanced carefully in personalized optimization approaches. DeWitt addressed this challenge through research on "segregated optimization models," which isolated user-specific optimization data to prevent potential information leakage [10]. His security analysis demonstrated that without proper isolation, personalized optimization could potentially reveal sensitive information—in one experiment, an adversary with access to the optimization decisions was able to determine which medical conditions a researcher was investigating based solely on the query patterns, raising serious privacy concerns. The segregated approach prevented this information leakage through a combination of strong isolation between user-specific models and differential privacy techniques applied to any cross-user learning. Experimental evaluation showed that this privacy-preserving approach maintained 82% of the performance benefit from personalization while reducing privacy risks to acceptable levels based on standard differential privacy metrics. For particularly sensitive domains like healthcare and financial services, the system could be configured to emphasize privacy over performance, reducing the performance benefit to approximately 65% of the non-private approach while providing stronger privacy guarantees. This configurable approach allowed organizations to implement personalized optimization even in highly regulated environments, balancing performance benefits against compliance requirements.

---

## 6. Conclusion

AI-powered query optimization represents the next frontier in database system evolution. By addressing the limitations of traditional approaches through adaptive, learning-based techniques, these systems promise to deliver improved performance and efficiency in increasingly complex data environments. While significant challenges remain, the rapid pace of research and early commercial implementations suggest that AI will fundamentally transform how we process and optimize queries in the coming years. As organizations continue to face growing data complexity and volume, investing in AI-enhanced query processing capabilities will become increasingly critical to maintaining competitive performance and efficiency. Those who successfully navigate the integration of AI into their data infrastructure will be well-positioned to extract maximum value from their data assets.

---

## References

- [1] Ji Sun and Guoliang Li, "An End-to-End Learning-based Cost Estimator," 2019. <https://www.vldb.org/pvldb/vol13/p307-sun.pdf>
- [2] Ryan Marcus et al., "Bao: Making Learned Query Optimization Practical," 2021. <https://dl.acm.org/doi/pdf/10.1145/3448016.3452838>
- [3] Viktor Leis et al., "How Good Are Query Optimizers, Really?," 2015. <https://www.vldb.org/pvldb/vol9/p204-leis.pdf>
- [4] Piotr Przymus, "Query Optimization in Heterogeneous CPU/GPU Environment for Time Series Databases," ResearchGate, 2014. [https://www.researchgate.net/publication/279525388\\_Query\\_Optimization\\_in\\_Heterogeneous\\_CPUGPU\\_Environment\\_for\\_Time\\_Series\\_Databases](https://www.researchgate.net/publication/279525388_Query_Optimization_in_Heterogeneous_CPUGPU_Environment_for_Time_Series_Databases)
- [5] Tim Kraska et al., "The Case for Learned Index Structures," 2018. [https://www.cl.cam.ac.uk/~ey204/teaching/ACS/R244\\_2018\\_2019/papers/Kraska\\_SIGMOD\\_2018.pdf](https://www.cl.cam.ac.uk/~ey204/teaching/ACS/R244_2018_2019/papers/Kraska_SIGMOD_2018.pdf)
- [6] Yongjoo Park et al., "QuickSel: Quick Selectivity Learning with Mixture Models," 2020. [https://web.eecs.umich.edu/~mozafari/php/data/uploads/sigmod\\_2020.pdf](https://web.eecs.umich.edu/~mozafari/php/data/uploads/sigmod_2020.pdf)
- [7] Ziqi Wang et al., "Building a Bw-Tree Takes More Than Just Buzz Words," 2018. <https://db.cs.cmu.edu/papers/2018/mod342-wangA.pdf>

- [8] Sanjay Krishnan et al., "Learning to Optimize Join Queries with Deep Reinforcement Learning," arXiv:1808.03196v2, 2019. <https://arxiv.org/pdf/1808.03196>
- [9] Ryan Marcus et al., "Bao: Learning to Steer Query Optimizers," arXiv:2004.03814v1, 2020. <https://arxiv.org/pdf/2004.03814>
- [10] Samuel Madden et al., "How to Build a High-Performance Data Warehouse," ResearchGate, 2006. [https://www.researchgate.net/publication/249825733\\_How\\_to\\_Build\\_a\\_High-Performance\\_Data\\_Warehouse](https://www.researchgate.net/publication/249825733_How_to_Build_a_High-Performance_Data_Warehouse)