

# veDB-HTAP: a Highly Integrated, Efficient and Adaptive HTAP System

Jianjun Chen, Li Zhang, Yu Xie, Wei Ding, Lixun Cao, Ye Liu, Yonghua Ding, Fangshi Li, Ke Wu, Haibo Xiu\*, Kui Wei, Le Cai, Rui Chang, Yuxiang Chen, Yuanjin Lin, Shangyu Luo, Jianfeng Qian, Xu Wang, Zikang Wang, Jian Zhang, Mingyi Zhang, Shicai Zeng, Jason Sun, Lei Zhang, Rui Shi, Pengwei Zhao

ByteDance Inc \*Duke University

## ABSTRACT

In this paper, we describe veDB-HTAP, a highly integrated, efficient, and adaptive HTAP system recently built in ByteDance. veDB-HTAP adopts a highly integrated system architecture by leveraging the Secondary Engine mechanism provided by MySQL and provides a seamless query processing experience across OLTP and OLAP engines. In addition, we introduce a cost-based and machine-learning-based smart query router that significantly outperforms the rule-based query router used in ByteHTAP, a precursor of veDB-HTAP. A key design principle of veDB-HTAP is the collaboration and adaptability of major system components, including query planning, query execution, and unified storage. Our adaptive query execution can be classified into two categories: 1) adaptive execution that dynamically collects and utilizes runtime statistics for better query performance; 2) utilizing runtime resource information to achieve a high quality of service even under heavy workloads. The experiments show that veDB-HTAP can achieve more than 3× speedup for TPC-H while consuming only one-third of the resources compared to ByteHTAP.

## CCS CONCEPTS

• Information systems → DBMS engine architectures.

## KEYWORDS

HTAP systems, cloud-native databases

### ACM Reference Format:

Jianjun Chen, Li Zhang, Yu Xie, Wei Ding, Lixun Cao, Ye Liu, Yonghua Ding, Fangshi Li, Ke Wu, Haibo Xiu\*, Kui Wei, Le Cai, Rui Chang, Yuxiang Chen, Yuanjin Lin, Shangyu Luo, Jianfeng Qian, Xu Wang, Zikang Wang, Jian Zhang, Mingyi Zhang, Shicai Zeng, Jason Sun, Lei Zhang, Rui Shi, Pengwei Zhao. 2025. veDB-HTAP: a Highly Integrated, Efficient and Adaptive HTAP System. In *Proceedings of the 2025 International Conference on Management of Data (SIGMOD '25)*. ACM, New York, NY, USA, 14 pages. <https://doi.org/XXXXXX.XXXXXXX>

Dr. Jianjun Chen is the corresponding author, [jianjun.chen@bytedance.com](mailto:jianjun.chen@bytedance.com).

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

SIGMOD '25, June 22–27, 2025, Berlin, Germany

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/XXXXXX.XXXXXXX>

## 1 INTRODUCTION

As ByteDance’s business grows, we have observed many applications that require both transactional processing and analytical processing over fresh data with strong data consistency requirements. To meet those business requirements, we built ByteHTAP [27] system, which was launched in the middle of 2021.

To provide context for this paper, we begin with a brief overview of ByteHTAP system. ByteHTAP employs an architecture of *separate engines over a shared storage*, integrating a ByteDance’s existing OLTP (TP) system (i.e., veDB [15, 72]) and an open-source OLAP (AP) system (i.e., Apache Flink [3]). ByteHTAP provides a unified API for query processing, facilitating seamless interaction with both OLTP and OLAP systems. A rule-based query router within the proxy layer automatically routes queries to the appropriate engine—either OLTP or OLAP—based on predefined rules. The Metadata Service provides centralized metadata access for the OLAP query optimizer and storage nodes. ByteHTAP guarantees strong data consistency by providing consistent data snapshots for its queries. The Metadata Service relays a globally committed Log Sequence Number (LSN) to the OLAP engine. Each query is assigned a read LSN derived from the globally committed LSN, ensuring that the query operates on a consistent data snapshot.

ByteHTAP has garnered extensive adoption within ByteDance since its release, which also exposes multiple important limitations. Additionally, we aim at providing HTAP capability on ByteDance’s public cloud (i.e., Volcano Engine [6]) as an extension to our veDB product [15], besides supporting internal customers. Many use cases are expected to involve upgrades from our existing veDB MySQL customers, who wish to leverage enhanced OLAP query-processing capabilities without modifying their applications.

In response to a myriad of new customer requirements, we identified three substantial limitations of ByteHTAP that motivated the design of *veDB-HTAP*:

- **Subtle semantic differences between OLTP and OLAP engines.**

- First, even though ByteHTAP’s OLTP engine uses MySQL read committed (RC) isolation level as default, its OLAP engine only supports snapshot isolation which often uses a slightly older snapshot than the most recently committed one in the OLTP engine due to log replication delay. Hence, the same query executed by both the OLTP and OLAP engines may produce different results. Therefore, ByteHTAP does not support session-level data consistency. Queries directed to the OLAP engine may not reflect earlier updates made by DML operations within the same session.

- Second, ByteHTAP adopts Apache Flink as its OLAP query engine, which uses its own SQL dialects and lacks certain MySQL dialects support. The system’s unified storage also needs enhancement to support certain MySQL dialects.
- **Simple rule-based query routing in the proxy has many limitations.** The rule-based query router in ByteHTAP’s proxy layer determines query routing based solely on predefined rules, without considering the performance characteristics or syntax support of OLTP and OLAP engines. This can lead to query failures if a query is sent to an engine that does not support its syntax. To address this, we propose a more sophisticated query router that leverages a cost-based query optimizer for in-depth plan analysis while maintaining MySQL compatibility. Additionally, we are exploring machine learning (ML)-based approaches to further enhance routing accuracy.
- **Limitations in Flink’s query engine and resource scheduler to support large-scale OLAP workloads.** Although ByteHTAP has made significant improvements over Flink and achieved up to 25% performance improvements on the TPC-H benchmark compared to Flink’s open source version [27], it did not satisfy the requirements of our new use cases, which demanded a more efficient and elastic OLAP system. For example, ByteDance’s application-platform-as-a-service team wants to build a scalable platform-as-a-service platform on top of ByteHTAP and they need to support thousands of developers with tens of millions of tables. In addition, they want great OLAP query performance over complex queries, i.e., seconds of query latency, that may consist of many joins and other complex operators. Furthermore, they want our system to be elastic and cost-effective, and it can be used in a multi-tenant environment.

In summary, we aim to build a highly integrated HTAP system that alleviates users from concerns about underlying engine differences by providing consistent querying semantics across its OLTP and OLAP engines. In addition, we seek to improve query routing to ensure each query is directed to the appropriate engine for execution. Furthermore, we strive to develop an efficient and elastic OLAP query engine capable of handling complex OLAP workloads in a multi-tenant environment.

Based on the aforementioned considerations, we propose veDB-HTAP, a system designed to serve as a generic HTAP solution capable of accommodating both large-scale and light-weight HTAP processing needs, including: (1) catering to large-scale customers with multi-tenant demands, and providing support for up to tens of millions of tables while ensuring high performance and scalability, and (2) offering solutions for light-workloads OLTP clients, and prioritizing cost efficiency and excellent MySQL compatibility.

We released veDB-HTAP to internal customers near the end of 2023 and plan to expand its availability to external customers on ByteDance’s Volcano Engine soon. Since its release, several customers have onboarded and we are actively migrating many of ByteHTAP customers to veDB-HTAP. Customers of veDB-HTAP have reported significant improvement in performance and resource utilization compared to ByteHTAP. veDB-HTAP’s support of MySQL Read Committed isolation level also simplifies the migration of our customers’ workloads. In summary, our key contributions are as follows:

- veDB-HTAP transforms system architecture of ByteHTAP into a highly integrated system. In veDB-HTAP, queries are always sent to the OLTP engine and we utilize the Secondary Engine mechanism [10] provided by MySQL to direct query to OLAP engine when necessary. veDB-HTAP also enables a seamless querying experience across OLTP and OLAP engines by providing a unified Read Committed isolation level. In addition, a cost-based smart query routing is implemented and used in production. We also designed and implemented a query routing prototype based on machine learning techniques, which showed some interesting and promising results in our preliminary study presented in this paper.
- We demonstrate how to build a high-performance and elastic OLAP query engine in veDB-HTAP to replace the Flink engine in ByteHTAP. A key aspect of our design strategy for achieving high performance and elasticity in our OLAP query system is to ensure that the major system components, including query planning, query execution, and the unified storage, operate in a collaborative and adaptive manner, which will be illustrated with concrete examples within the paper. Overall, our adaptive query execution can be classified into two categories: 1) adaptive execution that dynamically collects and utilizes runtime stats for better query performance; 2) utilizes runtime resource informations for better quality of service. The experiments show that veDB-HTAP can achieve more than 3× speedup while consuming 1/3 fewer resources than ByteHTAP.

The rest of the paper is organized as follows: Section 2 gives an overview of related work. Section 3 describes the overall architecture of veDB-HTAP as well as the implementations of its key components. Section 4 shows how veDB-HTAP supports a highly integrated HTAP processing through architecture change and smart query routing. Section 5 focuses on how veDB-HTAP’s OLAP query engine achieves high performance and elasticity. In Section 6, we provide some empirical measurements of veDB-HTAP. Finally, Section 7 concludes our work.

## 2 RELATED WORK

In this section, we first briefly review early HTAP databases and techniques in industry and research literature. Then, we discuss the latest HTAP systems and techniques in detail. Lastly, we present works in the area of resource management and workload isolation.

HTAP systems and techniques have been extensively proposed and studied over the years [26, 37, 39, 52, 55, 62, 66]. For example, SAP Hana [35, 44, 51, 71] and WildFire [18–20, 46] are dedicated commercial HTAP systems. In addition, many existing systems are extended to support HTAP functionalities, including Oracle’s Database In-Memory Option [42, 56], the BLU Acceleration [64] and IDAA [24] in IBM DB2, Google’s F1 Lightning [79], and the HTAP-related components in Microsoft SQL Server 2016 [43] and Greenplum [47]. Lastly, HTAP techniques are also discussed in some research works [48, 65, 70].

In recent years, HTAP systems have evolved to integrate with new architectures and technologies, leading to the emergence of innovative HTAP databases and prototypes.

For instance, TiDB [36] is an open-source distributed HTAP system built on top of TiKV and TiFlash. TiKV serves as a row-based

storage engine for transactional workloads, while TiFlash functions as a column-based storage engine for analytical queries. Starting with TiDB v7.0.0, TiFlash introduced support for disaggregated storage and compute [13]. TiDB features its own TiDB Server, which supports MySQL’s protocol and common syntax, though not all MySQL features are currently supported [12]. In contrast, veDB-HTAP employs the MySQL Secondary Engine to provide strong MySQL compatibility.

Alibaba PolarDB [74, 80, 81] is a cloud-native relational database. PolarDB-IMCI [74] extends PolarDB with in-memory column indexes and a column-based execution engine to speed up the processing for complex OLAP queries. It executes OLAP queries in its MySQL Read-Only (RO) nodes to separate OLTP workloads from OLAP workloads. Similar to veDB-HTAP, PolarDB-IMCI adopts a disaggregated storage architecture and also separates OLTP and OLAP processing. However, veDB-HTAP’s architecture utilizes MySQL Secondary Engine extension [10] and can plug in any OLAP engine with its modular design. veDB-HTAP has a powerful proprietary MPP OLAP engine that can execute queries across different servers while PolarDB-IMCI can only execute queries on a single MySQL RO node. PolarDB-SCC [81] adopts several techniques to reduce the data freshness gap between its Read-Write (RW) and RO nodes. Similarly, veDB-HTAP also provides Read Committed isolation and utilizes techniques such as multi-level LSNs to reduce the gap of data freshness between its OLTP and OLAP engines.

There are other commercial and prototype database systems that are enhanced to support HTAP workloads. Amazon Redshift [16] provides some HTAP capabilities by facilitating both the in-place querying of data in the OLTP services by using Redshift’s Federated Query, and the seamless copying and synchronization of data to Redshift by using Glue Elastic Views [68]. SingleStoreDB [63], formerly known as MemSQL [28], utilizes the cloud-native architecture and takes advantage of the elastic infrastructure. It supports transactions that perform scans on column stores and perform the seek on row stores, and uses indexes to speed up point reads and writes. Pixels-Turbo [21] is a hybrid query engine that processes queries in a scalable VM cluster and invokes cloud functions to accelerate the processing of unpredictable workload spikes. Diva [40] attempts to make MVCC systems more HTAP-friendly by decoupling three design concerns in MVCC DBMSs: separating version searching and version cleaning while isolating recovery from MVCC. This design improves the system’s performance while keeping memory and storage footprints, I/O, and runtime overhead lower than vanilla systems.

In addition, HTAP-related techniques are also widely discussed [61]. Many works focus on the storage of HTAP systems. LASER [67] shows how to use a Real-Time LSM-Tree to build a lifecycle-aware storage engine for HTAP systems. AIDE [45] is an analytic offload engine aiding MVCC databases for offloading query plans to computational storage, which aims to reduce resource contention and accelerate query operations by eliminating random and redundant data access in storage.

In the evaluation of an HTAP system, [53] proposed a method from the perspectives of performance and freshness. A concept called *throughput frontier* is introduced to capture the performance of an HTAP system, which visualizes both transactional and analytical throughput in a graph. A hybrid benchmark called HATtrick

is also proposed, which incorporates both throughput frontier and freshness as metrics.

Lastly, some works have discussed resource management and workload isolation in multi-tenant database systems [17, 25, 34, 57, 60, 73, 83, 85]. SQLVM [58, 59] provides an abstraction for performance isolation in a multi-tenant RDBMS based on the promise of reservation of resources. It proposes a fine-grained resource scheduling mechanism to meet each tenant’s reservations. It also provides a novel metering logic to audit the promise of resources. Another work [32] focuses on the problem of CPU sharing among tenants co-located at a server. It discusses how to use SQLVM [58] as an abstraction of CPU reservations to provide resources to tenants without any restriction on the tenant’s workloads. In comparison, veDB-HTAP adopts an adaptive approach to handle resource contention for both CPU and memory.

### 3 SYSTEM ARCHITECTURE AND IMPORTANT COMPONENTS

In this section, we first give a brief overview of MySQL Secondary Engine. Then, we talk about veDB-HTAP’s architecture and important components.

#### 3.1 MySQL Secondary Engine Overview

MySQL has provided built-in support for pluggable storage engines. Customized storage engines such as [75] have been developed and combined with existing storage engines for extended functionality. Coupled with additional query processing capability for workload off-loading or acceleration, this type of alternative storage engine is referred to as Secondary Engine by MySQL Heatwave [10]. MySQL can send different query plans to appropriate storage engines for execution based on the data processing characteristics of different storage engines and the cost, thereby leveraging multiple storage engines for optimal workload performance.

In integrating the Secondary Engine into MySQL, an additional phase/stage is added at the end of the optimizer to decide whether to send the query plan to the primary or secondary engine for execution. Users can specify whether to use a specific storage engine through the Secondary Engine keyword defined in the CREATE TABLE or ALTER TABLE command for existing MySQL tables. As an example, the following command

```
|| create table orders (
||   oid int not null primary key,
||   cust_id int not null, o_date date
|| ) secondary_engine = RAPID
```

would create a table with all columns also stored in the secondary engine registered in the system as RAPID, which is the query processing engine used in Heatwave [10]. We can plugin any customized query processing engine such as our own OLAP engine in the same fashion. Once the table has been created with the Secondary Engine support, loading the table will automatically load both MySQL’s default storage engine as well as the Secondary Storage. For existing tables that have been altered with Secondary Engine support, a background process will start to replicate the data from the existing MySQL storage to the newly added Secondary Storage.

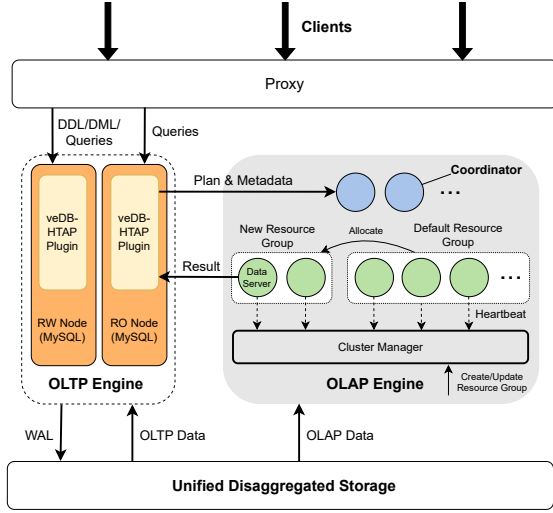


Figure 1: An illustration of the veDB-HTAP architecture.

### 3.2 veDB-HTAP System Overview

As shown in Figure 1, veDB-HTAP inherits ByteHTAP’s cloud-native architecture, with the following major changes.

**3.2.1 veDB-HTAP Plugin.** As shown in Figure 1, we extended the MySQL OLTP engine within veDB-HTAP using the veDB-HTAP Plugin, our proprietary Secondary Engine extension, on both the Read-Write (RW) and Read-Only (RO) nodes. With this new architecture, veDB-HTAP’s proxy consistently routes user queries to the OLTP engine, which determines the appropriate engine to use based on query complexity and other properties. In contrast, ByteHTAP sends queries to both OLTP and OLAP engines. In veDB-HTAP, DDLs and DMLs can only be sent to the RW node, while queries can be directed to both the RW and RO nodes.

**3.2.2 Adopt a proprietary high-performance and elastic MPP OLAP engine.** ByteHTAP utilizes Flink as the OLAP engine, which has been widely used in our company for ETL processing in streaming applications. However, it has some limitations in being used in large scale OLAP processing, e.g., resource scheduling/managing for high QPS workload, and lack of vectorized query engine. We have been improving the Flink engine for a while and it can support many existing customers currently.

In veDB-HTAP, we want to significantly extend its OLAP capabilities to further improve query performance in terms of QPS and query latency and reduce the cost of customers through great multi-tenancy support. For example, one of our largest customers manages thousands of shared tenants with a diverse range of job SLAs and priorities. Enabling resource sharing and isolation among a large number of tenants within a single veDB-HTAP cluster is essential for cost reduction. Since we have built a distributed query engine from scratch in our Krypton project [29], it is natural to reuse the Krypton’s distributed query processing engine in veDB-HTAP, which consists of a set of *Coordinators* for distributed plan generation and optimization, and a set of *Data Servers* (DS) for

distributed plan execution (see Figure 1). We have made significant enhancements to the OLAP engine and its Columnar Store to make the OLAP system highly efficient and elastic to support veDB-HTAP customers.

To support better resource sharing and resource isolation for multi-tenancy, veDB-HTAP introduces a centralized *Cluster Manager* to allow users to create resource groups among Data Servers. For example, customers can create a shared resource group to hold many low-priority tenants with lenient SLA and dedicated resource groups for tenants with stringent SLA. This can be done by sending *create/update resource groups* commands to the cluster manager with the corresponding min/max CPU cores and min/max memory sizes. The Cluster Manager will calculate the initial number of Data Servers to be allocated from a default resource group to the dedicated resource groups. When the workload on a resource group changes, Data Servers may be allocated to or removed from this group on the fly.

## 4 HIGH INTEGRATION BETWEEN OLTP AND OLAP ENGINES

In this section, we describe how veDB-HTAP seamlessly integrates its OLAP engine with its OLTP engine through veDB-HTAP Plugin to provide extendable OLAP query processing capability, MySQL Read Committed transaction isolation level, and smart routing.

### 4.1 Extensible OLAP Engine Integration Using veDB-HTAP Plugin

As shown in Figure 1, veDB-HTAP Plugin reuses the parser and analyzer from MySQL to generate a MySQL-compatible query plan and performs simple optimizations. It then sends the partially optimized query plan, together with required metadata, to a *coordinator* for distributed plan generation and optimization. After optimization, the colocated query scheduler divides the plan into a set of *fragments* in a distributed fashion and sends them to a set of *Data Servers* for execution, which will read the data from the unified storage and periodically send heartbeats to the cluster manager. Whenever a *block* (i.e., a batch of rows) is produced by its local fragment executor, it will be shuffled to its destination Data Server(s) for future processing. The final result blocks will be buffered at a dedicated Data Server and fetched by veDB-HTAP Plugin.

An important design decision is to have the OLAP query optimizer take a partially optimized logical plan from the OLTP planner instead of SQL text. This offers not only tighter compatibility between OLTP and OLAP query planners but also better efficiency in terms of reducing redundant optimization work. For example, the expression handling, such as data type deducing, casting, and constant folding optimization, are done by the OLTP planner. Some complex OLAP-oriented optimization such as join reorder, cost-based remote/local distribution enforcement required for distributed query processing, as well as common plan subtree sharing, are all handled in the OLAP optimizer, which is specifically designed for such tasks. In addition, we have defined an interface to translate an OLTP logical plan into an OLAP logical plan to allow future extensions to connect to different OLTP engines beyond MySQL.

Another important design decision is that we choose to physically decouple Coordinators (OLAP Query Optimizers) and Data Servers from the MySQL nodes to achieve better isolation between OLTP and OLAP workloads and better scalability and elasticity. As OLAP query optimization tasks can be resource-intensive, running them on separate processes minimizes the impact to OLTP workloads. As the metadata is stored inside MySQL system tables, our OLAP optimizer does not need to store any metadata, which makes it stateless and easily scaled out.

## 4.2 Support Read-Committed Isolation Level

In ByteHTAP, the OLAP engine supports consistent reads over a snapshot, which might miss the latest committed data in the OLTP engine due to log shipping and processing delays. Thus, ByteHTAP fails to maintain read consistency between the OLAP and OLTP engines. For instance, a query directed to the OLAP engine might not reflect the outcomes of a user's recent DML operations. veDB-HTAP addresses this shortcoming by implementing a Read-Committed (RC) transaction isolation level that spans both OLTP and OLAP engines. This upgrade guarantees that the effects of a write operation will be visible to subsequent read operations, irrespective of whether these operations occur within the same or different sessions.

Such support requires that a query in the OLAP engine must read and only read data committed prior to the query's arrival at the OLTP engine. This necessitates: (1) acquiring aligned read snapshot across OLTP and OLAP engines; (2) guaranteeing the query execution after the OLAP engine receives the corresponding latest data snapshot.

**Acquire aligned read snapshot on read-only (RO) nodes:** The data is replicated to the OLAP engine by applying logical log records. The veDB-HTAP Plugin assigns a logical log LSN to a query as its read snapshot. However, logical log is unavailable on RO nodes of the OLTP engine, which becomes a problem when a query is sent to a RO node to execute. To address this limitation, veDB-HTAP incorporates control information of logical log records in physical commit log records. The veDB-HTAP Plugin has been enhanced to process these physical log records, enabling it to construct OLTP snapshots and retrieve the corresponding logical log LSN. With this improvement, RO nodes can send queries to the OLAP engine using the most recently committed LSN as the snapshot, ensuring compliance with RC isolation level.

**Reduce delay of query execution caused by waiting for the corresponding latest data snapshot:** When a query requires the Read Committed (RC) isolation level, the veDB-HTAP Plugin may automatically delay its execution to ensure the required data is ready in the OLAP engine. To minimize this delay, two strategies are employed: (1) a carefully designed pipelined waiting mechanism and (2) use of a fine-grained commit-LSN examination. The pipelined waiting mechanism allows queries to progress through the plan optimizer and be dispatched to DS nodes for execution, even if the storage nodes have not yet applied the necessary log records. A secondary check occurs when DS nodes send queries to the storage nodes, ensuring that the required data is ready. This approach overlaps the delay caused by log application with the

network latency and query plan optimization time, effectively reducing or even eliminating the waiting period. In the fine-grained commit-LSN examination, the veDB-HTAP Plugin utilizes both instance-level and table-level commit LSNs to determine whether a query can proceed without delay. This method builds upon similar principles to enhance efficiency and responsiveness as [81].

For applications with a stringent query latency requirement, veDB-HTAP provides a compatible mode, called *AP Snapshot mode*, which allows a query to be executed immediately in the OLAP engine with its currently available logical LSN.

## 4.3 Smart Routing

In order to execute a query correctly and efficiently, ByteHTAP used a smart query router to automatically decide which engine to use, mainly based on pre-defined rules over query syntaxes. However, this approach is limited by the difficulty of enumerating all possible rules, as well as the inability of existing rules to cover all potential queries.

In veDB-HTAP, all queries enter the system via the OLTP engine. We can leverage its query optimizer to do the cost estimation to make more accurate, cost-based decisions to accommodate a wider range of query. Since OLAP queries typically incur higher costs than OLTP queries, we establish a cost threshold to differentiate between them. If a query plan's cost is lower than the pre-defined threshold, it will be executed by the OLTP engine; otherwise, it will be forwarded to the OLAP engine through veDB-HTAP Plugin. The production results showed that cost-based routing outperformed rule-based routing in the majority of cases. Moreover, the smart router also ensures that queries sent to the OLAP engine accommodate the syntax and LSN requirements.

**4.3.1 Limitations in cost-based routing.** However, a query's cost may not always align with its actual execution time, which can lead to incorrect routing decisions, i.e., sending a query to an engine that spends more time than the other one. We can use three examples to illustrate the problem, using the schema and data from the TPC-H 100G benchmark.

*Example 4.1. A Top-N query with ORDER BY and LIMIT clauses. The cost of the query in the OLTP engine is  $1.0 \times 10^3$ , with an execution time of 20ms. The query's execution time is 650ms in the OLAP engine.*

```
|| SELECT l_orderkey, l_linenumber FROM lineitem
|| WHERE l_shipdate <= DATE '1993-6-01'
|| ORDER BY l_orderkey, l_linenumber
|| LIMIT 11500, 1700;
```

*Example 4.2. A query joins seven tables, with an index on the column c\_phone. Its cost from the OLTP engine is  $2.2 \times 10^5$ , where it finishes in 70ms. The OLAP engine takes 3.54s to finish.*

```
|| SELECT SUM(l1.l_extendedprice) AS total_extended_price
|| FROM customer c, nation n, orders o, lineitem l1, orders o1,
||      partsupp p, orders o2
|| WHERE c.c_phone IN ('15-741-346-9870', '15-741-346-9871',
||                    '15-741-346-9872', '15-741-346-9873')
|| AND n.name <> 'egypt' AND l1.l_quantity > 5
|| AND o.o_custkey = c.c_custkey AND n.n_nationkey = c.c_nationkey
|| AND l1.l_orderkey = o.o_orderkey AND o1.o_custkey = o.o_custkey
|| AND o1.o_custkey = o2.o_custkey AND l1.l_partkey = p.ps_partkey;
```



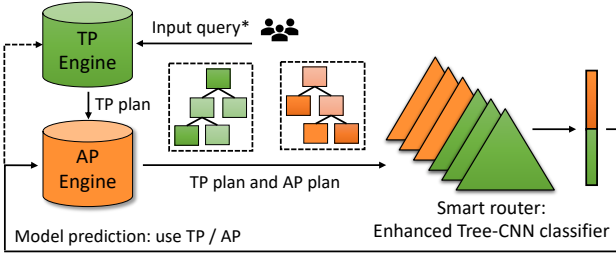


Figure 2: Workflow of the Learning-based Smart Router.

*Example 4.3.* A query joins three tables and has the *SUBSTRING* function in the *WHERE* clause. Its cost from the OLTP engine is  $5.2 \times 10^3$ , and the engine executes it in 5.80s. The OLAP engine completes in 310ms.

```
SELECT COUNT(*) FROM customer, nation, orders
WHERE SUBSTRING(c_phone, 1, 2) IN ('20', '40', '22', '30', '39',
'42', '21')
AND c_mktsegment = 'machinery'
AND n_name = 'egypt' AND o_orderstatus = 'p'
AND o_custkey = c_custkey AND n_nationkey = c_nationkey;
```

The examples above show that the cost obtained from the OLTP engine does not always align with the runtime performance. In Example 4.1, the query's cost is  $1.0 \times 10^3$ , and the OLTP engine takes 20ms to execute. However, in Example 4.2, despite a  $220\times$  higher cost than Example 4.1, the query finishes in only 70ms, which is just  $3.5\times$  slower than Example 4.1. On the contrary, the query's cost in Example 4.3 is about  $1/42$  of the cost in Example 4.2. However, its execution time in OLTP engine is conversely  $74\times$  higher than the time in Example 4.2.

These discrepancies highlight two key challenges for cost-based query routing:

- The estimated cost may not be always reliable enough to reflect the true performance of execution. A high cost may not always lead to slow execution, and a low cost may not guarantee a fast execution. Therefore, using a single static cost threshold for routing may produce incorrect choices of engines.
- It is difficult to decide the cost threshold for routing as the cost obtained from the OLTP engine cannot estimate the performance of the OLAP engine. In the examples above, we cannot find a single threshold to route all three queries correctly. We have also considered to include the plan cost from OLAP engine. However, the OLTP cost and OLAP cost cannot be directly compared as they are generated by different systems and have different characteristics.

**4.3.2 Learning-based routing.** Inspired by recent advancements in learning-based query optimization [30, 33, 49, 77, 82, 86, 87] and its successful application in real-world industrial systems [76], we propose a learning-based smart router. Based on an enhanced tree-based convolutional neural network (Tree-CNN) [50, 54], our model learns from historical query executions to predict which engine would execute a query more efficiently.

**Workflow of smart router.** The Figure 2 shows the workflow of the smart router, which is integrated with the workflow of veDB-HTAP and cost-based routing. Firstly, a query is sent to the OLTP engine and an OLTP plan is generated. If the query's cost is below

Training Size	5%	10%	20%	50%	100%
Learning-based	86.2%	87.9%	91.6%	92.4%	94.3%
Cost-based	80.2%	80.1%	80.0%	80.0%	79.1%

Table 1: Routing Accuracies with Different Training Sizes

a low threshold, it will be executed in the OLTP engine directly. We set this threshold conservatively to avoid the overhead of the learning-based routing for relatively simple queries. The remaining more complex queries will be sent to the OLAP engine. Then, the OLAP engine generates the OLAP plan and sends both OLTP and OLAP plans to the smart router. The smart router takes those plans as inputs, and predicts which engine will execute the query more efficiently. Lastly, the query is routed to the corresponding engine according to the prediction.

**Model structure.** Tree-CNN has been proven effective for performance prediction for queries [49, 76, 86]. Our work differs from previous studies in two major aspects: First, our model is more practical and more efficient to train because it only needs to select a better engine from two candidates according to the optimized plans, instead of predicting the actual query performance for all candidate plans. Second, our model processes and compares plans generated by two distinct engines with different characteristics. This introduces additional challenges compared to processing plans from a single engine.

Our model begins by formatting the input plans into tree structures and then encoding them using a series of convolutional layers. This encoding process transforms each plan tree into a vector representation that encapsulates key patterns and relationships within each plan. We also apply an attention mechanism to each convolutional layer to integrate runtime factors, which allows the model to prioritize the most critical aspects of the encoded features based on the current runtime conditions. Then, the results from the convolutional layers are processed through dynamic pooling and fully connected layers, which produces a final score to indicate which engine will have a better performance for the input query.

**Experimental setup and results.** To train and evaluate our model, we created customized workloads based on the queries collected from ByteDance's online applications. These queries range from simple single-table queries to complex queries involving up to seven-way joins. The workloads also include queries that the cost-based routing is difficult to handle and may produce incorrect routing decisions, such as *Top-N* queries with an *ORDER BY*, *LIMIT*, or *OFFSET* clause.

In our experiment, we used a workload of 7,033 queries, 90% of which were used for training and the remaining 10% for testing. We compared the routing accuracy of the pure learning-based routing with the pure cost-based routing. For the cost-based-routing test, an exhaustive search was first performed to identify the optimal cost threshold that maximized the routing accuracy for the training set. We also trained both methods with different percentages of the training data to test the impact of training sizes on the routing accuracy. The experiments were conducted on a cluster, with the details of its configuration provided in Section 6.

Table 1 shows the experimental results. The learning-based router consistently outperforms the cost-based router across all

training sizes, and the performance improvement is up to 15.0%. Additionally, the routing accuracy of the learning-based method consistently improves as the training size increases, which demonstrates the effectiveness of our training.

To the best of our knowledge, our work presents the first deep-learning-based query router for a production HTAP system. However, several challenges remain to be addressed before the learning-based router can be fully adopted in production applications with complex workloads. For example, a key challenge is determining when and how to retrain the model to accommodate evolving workloads. We are actively working on resolving these issues.

## 5 HIGHLY EFFICIENT AND ADAPTIVE QUERY PROCESSING

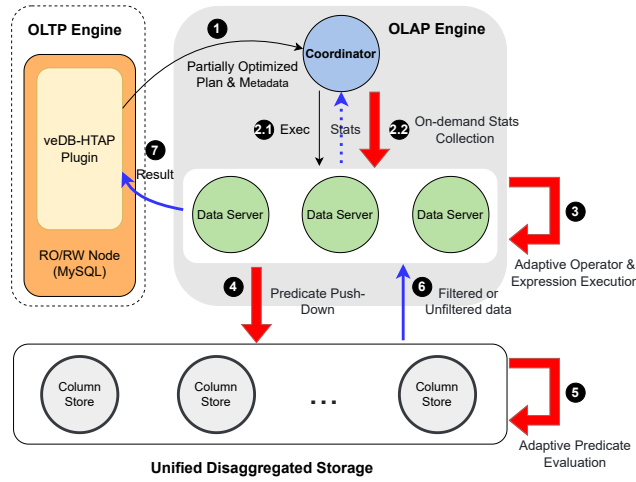


Figure 3: Adaptive Query Execution Workflow

As the ByteHTAP [27] users grow, we have noticed a few limitations in the query performance and resources utilization of the Flink-based AP engine. For example, for TPC-H 100G benchmark, we typically need 128 CPU cores to achieve 140s total running time.<sup>1</sup> In comparison, a well-known HTAP system on the market (referred to as HTAP-X) only needs 32 CPU cores to achieve 80s total running time for the same dataset. This has motivated us to develop our in-house massively-parallel-processing query engine.

In contrast to Flink’s batch/stream-based data processing engine, which usually aims for streaming applications, our new engine is designed with optimal analytical workload performance in mind by employing a push-based and vectorized execution engine. To maximize the CPU utilization, the execution engine adopts a coroutine-based asynchronous scheduler that can efficiently perform user-space context switch with minimal overhead. As shown later, for the above TPC-H 100G benchmark, we can achieve approximately 35s total run time with only 32 CPU cores.

<sup>1</sup>The 22 queries are submitted by a single stream

One key design of the new query processing engine is to have the query optimizer and executor work together so the performance-aware features can be adaptively performed based on the compile time and runtime information, and, therefore, achieve optimal query performance in practice. Different from existing work [78], which mainly focuses on dynamically changing query plans during query execution, our adaptive query execution framework can dynamically change the execution behavior inside the query executor without the need of changing the query plans.

Figure 3 shows the high level design of our adaptive query execution workflow. During the planning phase, the Coordinator can collect the missing column stats on-demand to generate optimal distributed plans (Section 5.1.1). Some of these stats, such as table row count, can be attached to plan and sent to the Data Server for better runtime decision. The Data Server can also collect necessary runtime stats along with its resource utilization metrics to help determine its execution behavior, such as hash table selection, disk spill decision, etc (Sections 5.1.2 to 5.2.1). In addition, the underlying storage can execute the push-down predicates in an adaptive fashion based on the resource utilizations on storage nodes to ensure the system stability and quality of service (Section 5.2.2). In summary, our adaptive query execution can be classified into two categories: (1) adaptive execution that dynamically collects and utilizes runtime stats for better query performance; (2) utilizes runtime resource informations for better quality of service.

### 5.1 Runtime Stats-Based Adaptiveness

In this section, we will first explain how stats are collected on-demand in veDB-HTAP and then showcase how these stats are used at runtime to improve the execution efficiency.

**5.1.1 On-Demand Stats Collection.** In cost-based query optimizers, the estimation of operator stats is of paramount importance in generating good plans. Typically, databases gather base table stats through stats collecting utilities, such as analyze-table commands executed by database users or administrators. However, this method often demands that users be familiar with the workloads and comprehend how query optimizer works. Motivated by this pain point, we introduced our lightweight and highly effective on-demand stats collecting feature, which aims to relieve users from cumbersome stats-related operations. In addition, when the base table data is updated to a certain extent where our HTAP system can readily obtain such information, we can leverage the cache eviction policy to expel stale stats and recollect them. There are a couple key challenges in this approach:

- The initial time needed to obtain on-demand stats during query planning contributes to the overall query execution time. To solve this problem, we offer an asynchronous option to the on-demand stats approach where a reasonable wait-time can be configured. If the on-demand stats job does not finish within the wait-time, the optimizer proceeds without stats. Nevertheless, the stats collecting job continues to retrieve the stats and store them in the cache for the benefit of subsequent queries.
- On-demand stats collection may increase system overhead and interfere with other critical user workloads running in parallel. To address this challenge, the first solution is to make on-demand

stats lightweight. Based on our customers' workloads and benchmark studies, we identified three most essential stats for our optimizer: 1. Base table row count, 2. Filter factor on base table filters. 3. Column NDV. Base table row counts can be collected easily when a table is populated. For base table filter factor, we employ a sampling-based approach, which asynchronously issues a scan request with filter predicates on a small portion of the table to obtain a sampled filter factor. We believe this is a more efficient and lightweight approach than a histogram-based method and provides better estimation accuracy, especially in complex predicate cases such as compound predicates or LIKE predicate. For the column NDV, we extended our on-demand sampling approach to issue an asynchronously approximate count distinct request on target columns. In addition, we also utilize dedicated resource work group for these stats collecting jobs to prevent resource contention from normal user queries.

**5.1.2 Adaptive Hash Table Selection.** Hash tables are the fundamental data structures used for GROUP BY and JOIN operators. It provides constant lookup time  $O(1)$  for inserting a key-value pair and retrieving the value from a given key. If the hash table can store  $ht_{capacity}$  elements in its memory block, and  $ht_{ndv}$  elements are currently inserted, the *load factor* of the hash table  $l$  can be defined as follows:

$$l = ht_{ndv} / ht_{capacity} \quad (1)$$

We call  $ht_{capacity}$  the capacity and  $ht_{ndv}$  the size of the hash table.

The hash table lookup performance usually stays stable as long as its load factor stays within its maximum value  $l_{max}$ , and degrades exponentially once it exceeds this threshold [41] due to collisions.<sup>2</sup> Most existing work on hash tables focus on either improving the lookup performance [8, 22, 23] or enhancing the memory efficiency [2, 9] by designing hash tables with high  $l_{max}$ . However, it is hard to find an optimal solution that satisfies both criteria, which motivates us to propose an adaptive hash table selection strategy that dynamically selects the most memory-efficient hash table based on the estimated hash table capacity at runtime with great lookup performance.

To estimate the hash table size, we utilize the HyperLogLog (HLL) algorithm [38] to calculate the number of distinct values (ndv) on the build side at runtime (i.e.,  $ht_{ndv}$  in (1)). Because the HLL algorithm is able to estimate cardinalities of  $> 10^9$  with a typical accuracy (standard error) of 2%, using 1.5 KB of memory, this approach works well in practice with minimal computation overhead. Once the estimated hash table size  $ht_{ndv}$  is calculated, based on (1), the hash table capacity  $ht_{capacity}$  can be expressed as:

$$ht_{capacity} = next_{pow2}(ht_{ndv} / l_{max}) \quad (2)$$

where  $next_{pow2}$  is a function to get the next power of two and  $l_{max}$  is the hash table's suggested maximum load factor.

In veDB-HTAP, we currently maintain two hash table implementations, with a higher  $l_{max}$  and a lower  $l_{max}$  separately. In order to achieve high  $l_{max}$ , additional metadata needs to be stored in the hash table, which causes extra indirections during lookup. Our experimental results have shown that, only when under the same

capacity, the one with higher  $l_{max}$  exhibits worse lookup performance (around 10%) than the other. In all other cases, it consumes less memory and shows better lookup performance. Based on above findings we provide a simple but effective strategy to adaptively select the hash table implementation. In specific, before hash build, we calculate the memory capacity used for each type of hash table according to (2). If they are equal, we select the one with better lookup performance, i.e., the one with lower  $l_{max}$ , otherwise, we always select the one with smaller capacity, i.e., the one with higher  $l_{max}$ . We will provide more experimental results in Section 6.

**5.1.3 Adaptive Runtime Filter.** Runtime filtering has been a popular optimization feature available in the industry for decades [4, 5, 11, 84]. The basic idea is to decrease the overall execution time by reducing data transmission and computing in the join operation by filtering out input data that does not match the join in advance for the probe side. One key challenge in runtime filtering is to decide if a filter on a join condition should be pushed down. When the filtering condition is not selective, pushing down filter may cause performance degradation due to the overhead introduced by runtime filter.

A common solution to this issue is to perform the cost evaluation during query optimization that heavily relies on the accuracy of join cardinality estimation, which can be inaccurate for complicate queries. In contrast, we push the cost estimation from query optimizer down to the query executor to achieve more intelligent runtime filtering decision. The actual join cardinality is collected during the hash build, and the runtime filter decision is made right after the hash build is finished. In specific, when building a hash table for each incoming build-side stream, we first collect the actual number of distinct values (ndv) of the hash table, denoted as  $ht_{ndv}$ , and then estimate the bloom filter ndv for the whole table, denoted as  $bf_{ndv}$ , based on  $ht_{ndv}$  and the number of streams. Assuming the row count and ndv of the prob side (left) table are  $target_{row}$  and  $target_{ndv}$ , respectively, the saved CPU cycles from runtime filtering can be expressed as:

$$cost_{saving} = target_{row} \cdot (1 - bf_{ndv} / target_{ndv}) \cdot \delta, \quad (3)$$

where  $\delta$  is a cost saving coefficient that can be interpreted as the CPU cycles spent on transmitting a single row over the network and joining that row on the given hash map. Note that both  $target_{row}$  and  $target_{ndv}$  come from the query optimizer.

Meanwhile, there is also an overhead of building the bloom filter, which can be expressed as:

$$cost_{overhead} = bf_{ndv} \cdot \beta, \quad (4)$$

where  $\beta$  is the overhead coefficient that can be interpreted as the CPU cycles spent on building a single entry in the bloom filter.

If the bloom filter is beneficial for the join operator, then based on equations (3) and (4), we should have the following condition satisfied:

$$cost_{overhead} < cost_{saving} \quad (5)$$

The executor uses the above cost model to evaluate whether the runtime filter should be generated and pushed down to the storage node. Different from the work in [78], which only uses the runtime stats such as row counts on both sides to make a rough decision, we estimate the bloom filter ndv in a more accurate way. As shown

<sup>2</sup>We limit our discussion to hash tables using open addressing method [41] for collision resolution where all hash table entries are stored in the bucket array itself, as most modern hash table is implemented in this way to achieve low memory overhead and good lookup performance.



later in our experiments, compared to the case where the runtime filter is always turned on, this collaborative runtime filter strategy can bring over 20% performance boost to certain TPC-H queries.

## 5.2 Resource-Aware Adaptiveness

Next, we explain how the executor dynamically changes the execution behavior based on runtime resource informations such as CPU usage, memory availability, etc, to improve both query performance and quality of service.

**5.2.1 Adaptive Disk Spill.** Our multi-tenant Data Server clusters can overcommit memory for cost efficiency, as concurrently running queries may not require their full memory allocation simultaneously. If a Data Server encounters memory shortages, disk spilling acts as a safeguard mechanism to ensure the successful execution of queries.

Most disk spill implementations [4, 5, 7, 69] employ preset memory thresholds at the query level or the operator level for *spillable* operators, such as aggregation and join, so that disk spill can take place whenever the consumed memory exceeds the thresholds. For example, the optimizer can estimate an operator’s memory usage and compare it with the threshold to make the spill decision during planning. Similarly, the executor can also compare the current operator’s memory consumption with this threshold to make the spill decision at runtime.

However, in real-world scenarios, selecting the right query or operator level threshold is a challenging task for several reasons. First, workloads are dynamic, and their memory consumption is unpredictable. Second, accurate cardinality estimation is difficult for the optimizer, especially with complex queries. Third, the threshold must balance preventing severe memory shortages on a Data Server with fully utilizing available resources.

This motivates us to propose an adaptive disk spill approach based on the runtime memory utilization of Data Servers. Two memory utilization thresholds, referred to as the *soft limit* and *hard limit*, are introduced. Before the memory utilization reaches the *soft limit*, no disk spill is required, and queries can allocate memory as needed. When the memory consumption exceeds the *hard limit*, all running and incoming queries will be terminated to prevent the system from out-of-memory (OOM).

When the memory usage falls between these two limits, indicating that the system memory is under pressure but not critically low, disk spill will be triggered. However, instead of spilling all running queries, in our adaptive approach, disk spill will only be selectively applied to a subset of queries and their execution fragments. These candidate fragments are chosen based on their runtime memory usage, with the goal of bringing the overall memory consumption down below the *soft limit*. If only spilling these queries are not fast enough to achieve this goal, additional queries may be considered to spill in a more aggressive way. This adaptive approach minimizes the performance impact of disk spilling to the overall system throughput by targeting only a subset of queries, and provides a fast response to ensure system stability at critical resource pressures.

**5.2.2 Adaptive Predicates Push-Down.** veDB-HTAP’s unified storage system is engineered to not only facilitate columnar data generation and scanning but also support a broad spectrum of push-down

computations, such as filters and aggregates, which can significantly accelerate query processing by minimizing data transfer across the network, albeit at the expense of higher resource consumption. Besides push-down processing and data scans across thousands of data partitions on a storage node, the unified storage engine also need to execute concurrent background tasks such as log replication and apply, Delta Store flush, compaction, and garbage collection on those data partitions. A crucial aspect of the unified storage engine is its ability to segregate the performance effects of these distinct tasks.

To enhance cost efficiency, the engine optimally serves a multi-tenant workload at elevated resource utilization levels. To protect it from ever getting overloaded, the multi-tenant unified storage engine coordinates with other system components, such as Data Servers and the replication framework, allowing for the selective delay, avoidance, rejection, or parameter value changing of operations based on resource utilization while guaranteeing the correctness of query executions.

One of those approaches is adaptive predicates push-down. Initially, the Data Server issues a scan request with predicates attached to the storage engine. Once receiving the request, the storage engine performs the scan and predicates evaluation as expected. During scan processing, the storage engine checks its CPU utilization periodically. If it exceeds certain threshold, the rest scanned data will be returned to the Data Server without any predicate evaluation. On the Data Server side, the scan operator will selectively apply the predicates to the returned data to ensure correctness.

## 6 PERFORMANCE STUDY

In this section, a set of experiments are conducted to evaluate our system. Except for the elasticity experiment, all the experiments are conducted on a cluster of six machines with configurations listed in Table 2. Note the cluster contains four Data Servers and one Coordinator that is co-located on one of Data Servers. In addition, the cluster uses one read-write (RW) node and one read-only (RO) node to handle HTAP SQL queries.

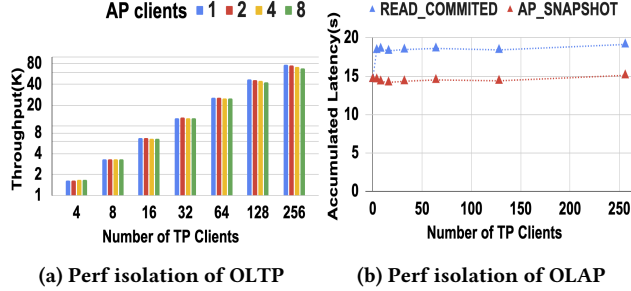
To evaluate the HTAP capabilities of our highly integrated system, we conduct experiments using a hybrid OLTP and OLAP workload based on the CH-benCHmark [31]. This benchmark combines the unmodified TPC-C [1] workload with modified TPC-H queries. The experiments are performed on CH-benCHmark datasets with 100 warehouses. Additionally, to assess the efficiency and adaptability of our new MPP OLAP engine in handling complex query processing, we use the widely adopted TPC-H [14] benchmark with data volumes of 100GB and 1TB. We also compare the execution times of individual queries against the published results of a well-known HTAP system on the market, referred to as HTAP-X.

### 6.1 Hybrid OLTP and OLAP Workloads

The performance isolation between OLTP and OLAP is critical for HTAP systems. We test the performance isolation on the RW node with mixed workloads from CH-BenCHmark. It enables running both OLTP and OLAP queries on a set of shared tables in one database, and the data is firstly loaded from the OLTP side and then replicated to the OLAP side. As Figure 4 (a) shows, when the number of OLTP clients is less than 128, OLTP’s throughput is almost not

**Table 2: cluster configurations.**

RW/RO Node	16v CPU, 1 NUMA node 64GB DRAM
Data Server/Coordinator Nodes	8v CPU, 1 NUMA node 32GB DRAM
Storage	Unified Storage
OS	Debian 4.14.81
Network	25Gbps Ethernet

**Figure 4: Performance isolation on HTAP workloads**

affected as the number of OLAP clients increases. However, when we use 256 OLTP clients, as increasing OLAP clients to issue analytical queries, there is a small OLTP throughput degradation (<10%) since CPU resources on the RW node become a bottleneck.

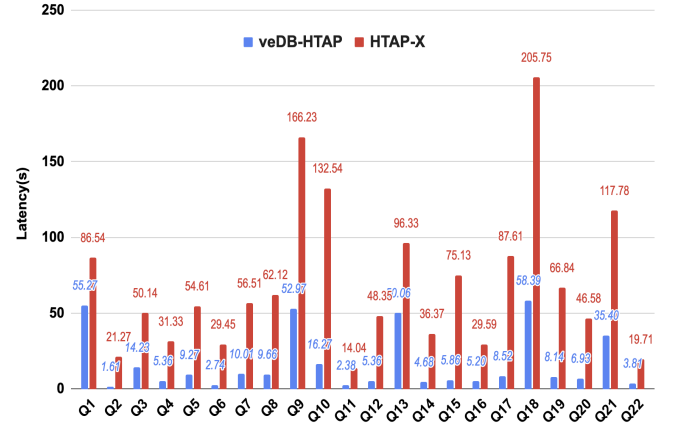
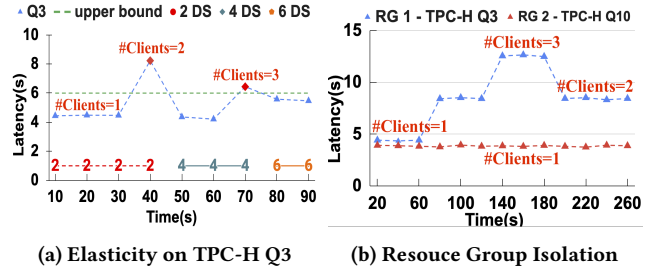
Next, we conducted similar experiments by submitting analytical queries to the RO node while keeping the TP workload on the RW node. We verified that there was no OLTP throughput degradation in such configurations. In production, we recommend customers to run the veDB-HTAP Plugin on MySQL RO nodes if they want to avoid potential performance degradation over OLTP workloads.

Figure 4 (b) shows the performance impact on OLAP queries when we increase the number of OLTP clients with different transaction isolation levels. As we increase the number of TP clients from 4 to 256, the accumulated AP query latency only increases slightly, which demonstrates great performance isolation in our system. In addition, we can observe that the total query latency in Read Committed (RC) isolation is about 4 seconds higher than the AP snapshot mode (for each query, we notice there is about 100-200ms latency delay), where AP queries can be evaluated over current data snapshots like in ByteHTAP. This performance degradation is expected since AP queries need to wait for the data snapshot with the latest committed LSN in the RC mode.

## 6.2 OLAP Performance

Achieving advanced OLAP performance is one of the foremost motivations of veDB-HTAP. In this section, we evaluate the OLAP performance of our system, and compare it with HTAP-X's latest published results, using equivalent cluster resources as reported by HTAP-X. We present the query latency results for all 22 queries in TPC-H 1T benchmark.

Figure 5 shows the result for TPC-H 1TB, where veDB-HTAP's total latency is 372 seconds, which is about 4× speedup compared

**Figure 5: TPC-H 1T performance****Figure 6: Elasticity and Isolation**

with HTAP-X's 1534 seconds of latency. Except Q1, all other queries' running times on veDB-HTAP are at least 2× speedup than HTAP-X because of better query plans and more efficient query execution. We also conduct a comparison with ByteHTAP on TPC-H 1T. veDB-HTAP can achieve about 3× speedup (372 seconds vs 1034 seconds) with less than one-third computing resources. In the TPC-H 1T experiment, we observe that disk spill happens for memory-intensive query Q18. Without the disk spill functionality, TPC-H 1T won't succeed under our experiment setup with constrained memory configuration.

In summary, veDB-HTAP delivers a significant improvement in OLAP performance over ByteHTAP, and demonstrates substantially superior performance compared to HTAP-X.

## 6.3 Elasticity and Isolation

To evaluate the elasticity of our system, we use Q3 in the TPC-H 100GB benchmark with an upper-bound latency requirement of 6 seconds, under different resource group configurations. As shown in Figure 6 (a), when increasing #Clients (the number of query streams) from 1 to 2, the query latency also increases and exceeds the upper bound of the latency requirement. Meanwhile, all DS nodes' CPUs are fully utilized and the system is overloaded. After adding two more DS nodes to the resource group, the query latency drops down in a few seconds. We repeat the above process: add another two DS nodes into the resource group after increasing

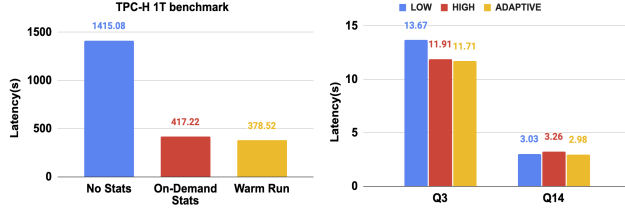


Figure 7: On-demand stats collection

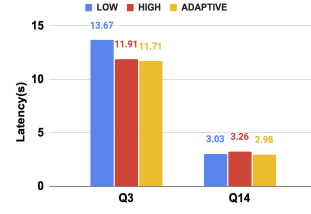


Figure 8: Adaptive hash table selection

#Clients from 2 to 3. Similar reductions in latency are observed, demonstrating that our system scales out efficiently and rapidly.

To demonstrate the resource isolation capability, we run TPC-H Q3 and TPC-H Q10 on two different resource groups concurrently, where each resource group contains two DS nodes. As shown in Figure 6 (b), increasing the #Clients of TPC-H Q3 on resource group 1 has no impact for TPC-H Q10 on resource group 2.

#### 6.4 Adaptive Query Processing

In this section, we conduct multiple experiments to evaluate the effectiveness of adaptive query processing techniques in veDB-HTAP including runtime stats based and resource aware adaptiveness in Section 5.

**On-Demand Stats Collection** To assess the efficacy of the on-demand stats, we carry out experiments on the TPC-H 1T benchmark, as depicted in Figure 7. Initially, we execute all 22 queries without gathering any stats and some queries like Q5, Q8, and Q13 can't even be completed due to suboptimal plans. The remaining queries are finished in 1415 seconds. Subsequently, we sequentially run all 22 queries with the on-demand stats feature. All 22 queries are completed in 417 seconds and we notice that the plans are progressively improved starting from Q4. In the third run, since all the collected stats are saved in the cache, we attain the optimal performance of 379 seconds. This experiment illustrates how our system can progressively achieve great performance without any human intervention.

**Adaptive Hash Table Selection** To evaluate the effectiveness of adaptive hash table selection, we conduct experiments on TPC-H 1T benchmark with different hash table implementations. As Figure 8 shows, neither the high nor the low load hash table is optimal for both queries. For TPC-H Q3, the high load factor hash table selection can achieve about 13% performance improvement compared with the low load factor hash table because its allocated capability is smaller in this query. However, choosing the high load factor hash table causes 7% performance degradation for TPC-H Q14 than the low load factor hash table, since the high load factor hash table implementation has a worse lookup performance while their allocated capacity is same.

By adopting the adaptive hash table selection method, we can achieve the optimal performance and memory usage for both queries. Note that in our design, there is an overhead in calculating the estimated NDV. However, this experiment shows the gain from our approach outweighs the overhead of calculating the NDV.

**Adaptive Runtime Filter** Since our system uses a disaggregated storage architecture, optimizations, such as runtime filters

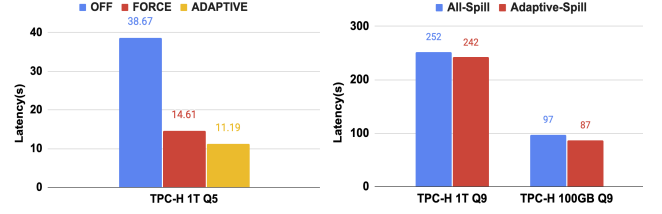


Figure 9: Adaptive runtime filter

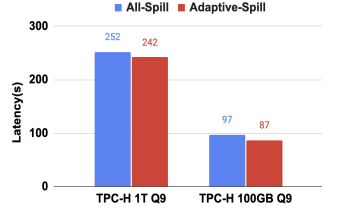


Figure 10: Adaptive disk spill

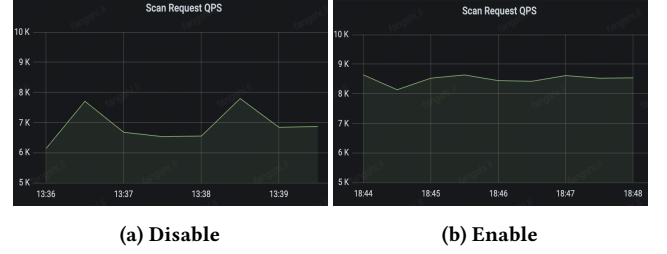


Figure 11: Adaptive Predicates Push-Down

and predicate pushdown, can significantly reduce data transfer costs from storage to Data Servers through the network. Our experiment shows that we achieve around 41.36% performance improvement on the TPC-H 100GB benchmark with the help of the runtime filter, and even more improvement on the TPC-H 1T workload.

Since building and applying runtime filters during query execution can incur additional CPU overhead on both Data Servers and storage nodes, it's essential to choose the right joins with high selectivity as the runtime filter candidates in an adaptive manner as we described in Section 5. We conduct experiments to evaluate the adaptive runtime filter on TPC-H Q5, which contains a mix of high and low-selectivity joins on a TPC-H 1T. Figure 9 shows we can achieve about 62.2% performance improvement from 38.67 seconds to 14.61 seconds by forcing runtime filters on all join conditions compared with the case without using runtime filters at all. With the adaptive runtime filter, we can further improve the latency to 11.19 seconds by choosing only highly selective joins, i.e., not using runtime filters over the join between the tables SUPPLIER and LINEITEM to avoid unnecessary runtime filter overhead.

**Adaptive Disk Spill** To demonstrate the effectiveness of adaptive disk spill, we concurrently run 1 TPC-H 1T Q9 and 8 TPC-H 100G Q9 on the cluster that has 128GB total memory. We compare the performance across three setups, namely, (1) no-spill, which simply disables the disk spill feature; (2) all-spill, which spills all the queries when the soft limit is reached; and (3) adaptive-spill, which selectively spills queries as needed in an adaptive manner as explained in Section 5.2.1. As expected, the no-spill setup can't finish because of limited memory capacity. Figure 10 illustrates the results for the all-spill and adaptive-spill setups, where adaptive-spill provides 4% and 12% performance improvements over all-spill for TPC-H 1T Q9 and TPC-H 100G Q9, respectively. This is because that TPC-H 1T Q9 will always be spilled in the given environment due to insufficient system memory, but for TPC-H 100G Q9, if only a small portion of these 8 queries are spilled (together with the

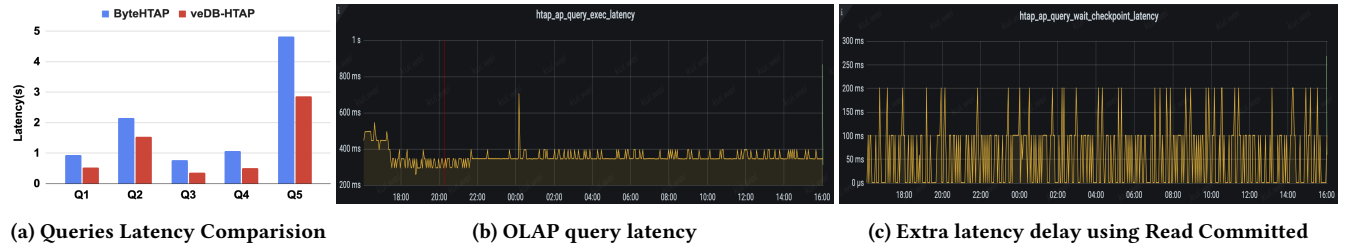


Figure 12: Production workloads

spilled TPC-H 1T Q9), then the rest may not need to be spilled at all, which explains the 12% overall performance gain.

**Adaptive Predicates Push-Down** In Figure 11, we present the outcomes of an experiment with a query workload consisting of two distinct query sets, one with some predicates that can be pushed down to the storage engine, and the other without any predicates. We want to demonstrate how adaptive predicates push-down can effectively improve the QPS of a mixed query workload, particularly in scenarios when CPU is heavily loaded in the multi-tenant storage. Note that since the shared storage nodes can not scale out as easily as the Data Servers since they contain states, it is very important to prevent storage nodes get overloaded in the case of heavy workloads.

As shown in Figure 11(a), in the absence of adaptive predicates push-down, query execution under a heavy workload in the storage can result in significant variability in the storage scan QPS. In contrast, as illustrated in Figure 11 (b), the average QPS when enabling this feature is about 20% higher and has less QPS variation occurring. This is mainly because with adaptive predicates push-down enabled, the storage will disable evaluating some of the pushed down predicates in an adaptive manner to avoid overload its CPUs, which can be reallocated to the second set of queries, thereby increasing the overall QPS.

## 6.5 Performance of Production Workloads

Since ByteHTAP became generally available in the middle of 2021, it has been deployed widely to serve many ByteDance internal customers, including Douyin, Lark, Finance, etc. Since veDB-HTAP was released near the end of 2023, we are working on migrating existing customers from ByteHTAP to veDB-HTAP for performance and resource utilization improvement.

In this section, we studied veDB-HTAP’s performance on several real-world workloads in our production environment. Figure 12 (a) shows the performance improvement based on the measurement of some representative queries from one early adopter after migrating from ByteHTAP to veDB-HTAP. This customer achieves about 2× performance improvement, by using only about 20% CPU resource for veDB-HTAP compared with ByteHTAP to serve their production workload. In Figure 12 (b), another new customer of veDB-HTAP requires read committed transaction isolation with a database over 10 Terabytes in size. Figure 12 (b) shows over a 24-hour period, the OLAP query latency consistently remained below 1 second while the system handled thousands of QPS from concurrent OLTP queries. In addition, Figure 12 (c) illustrates the additional latency incurred when waiting for the latest snapshot

under the Read Committed isolation level, which remained below 20ms.

As veDB-HTAP continues to be adopted by more customers, we will gain deeper insights into the system and share our findings in the future.

## 7 CONCLUSIONS

veDB-HTAP is designed to be a highly integrated, efficient and adaptive HTAP system. veDB-HTAP adopts a highly integrated system architecture and provides a seamless query processing experience across OLTP and OLAP engines with a newly built cost-based and ML-based smart query router. One key design principle of veDB-HTAP is to make major system components, including query planning, query execution, and unified storage, collaborative and adaptive. The experiments show that veDB-HTAP can achieve more than 3× speedup for TPC-H while consuming only one-third of the resources compared to ByteHTAP.

## ACKNOWLEDGMENTS

We would like to extend our thanks to all people who made contributions to the design and development of this veDB-HTAP system: Xiaofeng Bao, Peijun Hao, Ron Hu, Shuaiyu Jin, Yupeng Jin, Gang Liao, Bingyu Liu, Zixiong Liu, Yintong Ma, Yuanyuan Nie, Siyuan Ren, Xin Sui, Mengqi Wang, Peipei Wang, Lingyan Yin, Li Zhang, Lutong Zhang, Ruyi Zhang, Xiahao Zhang, Jie Zhou, Yihao Hu.

## REFERENCES

- [1] 2022. *TPC-C Specification*. Retrieved February 23, 2022 from [http://www.tpc.org/tpc\\_documents\\_current\\_versions/pdf/tpc-c\\_v5.11.0.pdf](http://www.tpc.org/tpc_documents_current_versions/pdf/tpc-c_v5.11.0.pdf)
- [2] 2024. *Abseil Flat HashMap*. Retrieved October 25, 2024 from <https://abseil.io/docs/cpp/guides/container>
- [3] 2024. *Apache Flink*. Retrieved March 12, 2024 from <https://flink.apache.org>
- [4] 2024. *Apache Impala*. Retrieved February 26, 2024 from <https://impala.apache.org/>
- [5] 2024. *Apache Spark*. Retrieved February 23, 2024 from <https://spark.apache.org>
- [6] 2024. *ByteDance Volcano Engine*. Retrieved March 6, 2024 from <https://www.volcengine.com>
- [7] 2024. *Clickhouse*. Retrieved February 26, 2024 from <https://clickhouse.com/>
- [8] 2024. *ClickHouse Hash Table*. Retrieved October 25, 2024 from <https://clickhouse.com/blog/hash-tables-in-clickhouse-and-zero-cost-abstractions>
- [9] 2024. *F14 Hash Table*. Retrieved October 25, 2024 from <https://github.com/facebook/folly/blob/main/folly/container/F14.md>
- [10] 2024. *MySQL Heatwave*. Retrieved February 23, 2024 from <https://dev.mysql.com/doc/heatwave/en/mys-hw-analytics.html>
- [11] 2024. *Presto*. Retrieved February 23, 2024 from <https://prestodb.io>
- [12] 2024. *TiDB MySQL Compatibility*. Retrieved March 6, 2024 from <https://docs.pingcap.com/tidb/stable/mysql-compatibility>
- [13] 2024. *TiFlash Disaggregated Storage and Compute Architecture and S3 Support*. Retrieved March 6, 2024 from <https://docs.pingcap.com/tidb/stable/tiflash-disaggregated-and-s3>
- [14] 2024. *TPC-H*. Retrieved February 11, 2024 from <http://www.tpc.org/tpch/>

- [15] 2024. *veDB for MySQL*. Retrieved March 6, 2024 from <https://www.volcengine.com/product/vedb-mysql>
- [16] Nikos Armatatzoglou, Sanuj Basu, Naga Bhanoori, Mengchu Cai, Naresh Chainani, Kiran Chinta, Venkatraman Govindaraju, Todd J Green, Monish Gupta, Sebastian Hillig, et al. 2022. Amazon Redshift re-invented. In *Proceedings of the 2022 International Conference on Management of Data*. 2205–2217.
- [17] Pankaj Arora, Surajit Chaudhuri, Sudipto Das, Junfeng Dong, Cyril George, Ajay Kalhan, Arnd Christian König, Willis Lang, Changsong Li, Feng Li, et al. 2023. Flexible Resource Allocation for Relational Database-as-a-Service. *Proceedings of the VLDB Endowment* 16, 13 (2023), 4202–4215.
- [18] Ronald Barber, Christian Garcia-Arellano, Ronen Grosman, Guy Lohman, C Mohan, Rene Muller, Hamid Pirahesh, Vijayshankar Raman, Richard Sidle, Adam Storm, et al. 2019. Wiser: A highly available HTAP DBMS for iot applications. In *2019 IEEE International Conference on Big Data (Big Data)*. IEEE, 268–277.
- [19] Ronald Barber, Christian Garcia-Arellano, Ronen Grosman, Rene Mueller, Vijayshankar Raman, Richard Sidle, Matt Spilchen, Adam J Storm, Yuanyuan Tian, Pinar Tözün, et al. 2017. Evolving Databases for New-Gen Big Data Applications. In *CIDR*.
- [20] Ronald Barber, Matt Huras, Guy Lohman, C Mohan, Rene Mueller, Fatma Özcan, Hamid Pirahesh, Vijayshankar Raman, Richard Sidle, Oleg Sidorkin, et al. 2016. Wildfire: Concurrent blazing data ingest and analytics. In *Proceedings of the 2016 International Conference on Management of Data*. 2077–2080.
- [21] Haoqiong Bian, Tiannan Sha, and Anastasia Ailamaki. 2023. Using Cloud Functions as Accelerator for Elastic Data Analytics. *Proceedings of the ACM on Management of Data* 1, 2 (2023), 1–27.
- [22] Altan Birler, Tobias Schmidt, Philipp Fent, and Thomas Neumann. 2024. Simple, Efficient, and Robust Hash Tables for Join Processing. In *Proceedings of the 20th International Workshop on Data Management on New Hardware (Santiago, AA, Chile) (DaMoN '24)*. Association for Computing Machinery, New York, NY, USA, Article 4, 9 pages. <https://doi.org/10.1145/3662010.3663442>
- [23] Maximilian Böther, Lawrence Benson, Ana Klimovic, and Tilmann Rabl. 2023. Analyzing Vectorized Hash Tables across CPU Architectures. *Proc. VLDB Endow.* 16, 11 (July 2023), 2755–2768. <https://doi.org/10.14778/3611479.3611485>
- [24] Dennis Butterstein, Daniel Martin, Knut Stolze, Felix Beier, Jia Zhong, and Lingyun Wang. 2020. Replication at the speed of change: a fast, scalable replication solution for near real-time HTAP processing. *Proceedings of the VLDB Endowment* 13, 12 (2020), 3245–3257.
- [25] Wei Cao, Feifei Li, Gui Huang, Jianghang Lou, Jianwei Zhao, Dengcheng He, Mengshi Sun, Yingqiang Zhang, Sheng Wang, Xueqiang Wu, et al. 2022. Polardb-x: An elastic distributed relational database for cloud-native applications. In *2022 IEEE 38th International Conference on Data Engineering (ICDE)*. IEEE, 2859–2872.
- [26] Jianjun Chen, Yu Chen, Zhibiao Chen, Ahmad Ghazal, Guoliang Li, Sihao Li, Weijie Ou, Yang Sun, Mingyi Zhang, and Minqi Zhou. 2019. Data management at huawei: Recent accomplishments and future challenges. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*. IEEE, 13–24.
- [27] Jianjun Chen, Yonghua Ding, Ye Liu, Fangshi Li, Li Zhang, Mingyi Zhang, Kui Wei, Lixun Cao, Dan Zou, Yang Liu, Lei Zhang, Rui Shi, Wei Ding, Kai Wu, Shangyu Luo, Jason Sun, and Yuming Liang. 2022. ByteHTAP: Bytedance's HTAP System with High Data Freshness and Strong Data Consistency. *Proc. VLDB Endow.* 15, 12 (sep 2022), 3411–3424. <https://doi.org/10.14778/3554821.3554832>
- [28] Jack Chen, Samir Jindel, Robert Walzer, Rajkumar Sen, Nika Jimsheleishvili, and Michael Andrews. 2016. The MemSQL Query Optimizer: A modern optimizer for real-time analytics in a distributed database. *Proceedings of the VLDB Endowment* 9, 13 (2016), 1401–1412.
- [29] Jianjun Chen, Rui Shi, Heng Chen, Li Zhang, Ruidong Li, Wei Ding, Liya Fan, Hao Wang, Mu Xiong, Yuxiang Chen, Benchao Dong, Kuankuan Guo, Yuanjin Lin, Xiao Liu, Haiyang Shi, Peipei Wang, Zikang Wang, Yemeng Yang, Junda Zhao, Dongyan Zhou, Zhikai Zuo, and Yuming Liang. 2023. Krypton: Real-time Serving and Analytical SQL Engine at ByteDance. *Proc. VLDB Endow.* 16, 12 (2023), 3528–3542. <https://doi.org/10.14778/3611540.3611545>
- [30] Tianyi Chen, Jun Gao, Hedui Chen, and Yaofeng Tu. 2023. LOGER: A learned optimizer towards generating efficient and robust query execution plans. *Proceedings of the VLDB Endowment* 16, 7 (2023), 1777–1789.
- [31] Richard Cole, Florian Funke, Leo Giakoumakis, Wey Guy, Alfons Kemper, Stefan Krompass, Harumi Kuno, Raghunath Nambiar, Thomas Neumann, Meikel Poess, et al. 2011. The mixed workload CH-benCHmark. In *Proceedings of the Fourth International Workshop on Testing Database Systems*. 1–6.
- [32] Sudipto Das, Vivek R. Narasayya, Feng Li, and Manoj Syamala. 2013. CPU sharing techniques for performance isolation in multi-tenant relational database-as-a-service. *Proceedings of the VLDB Endowment* 7, 1 (2013), 37–48. <https://doi.org/10.14778/2732219.2732223>
- [33] Bolin Ding, Rong Zhu, Jingren Zhou, et al. 2024. Learned Query Optimizers. *Foundations and Trends® in Databases* 13, 4 (2024), 250–310.
- [34] Aaron J Elmore, Sudipto Das, Divyakant Agrawal, and Amr El Abbadi. 2011. Zephyr: live migration in shared nothing databases for elastic cloud platforms. In *Proceedings of the 2011 ACM SIGMOD International Conference on Management of data*. 301–312.
- [35] Franz Färber, Sang Kyun Cha, Jürgen Primsch, Christof Bornhövd, Stefan Sigg, and Wolfgang Lehner. 2011. SAP HANA database: data management for modern business applications. *Proceedings of the VLDB Endowment* 4, 4 (2011), 45–51.
- [36] Dongxu Huang, Qi Liu, Qiu Cui, Zhuhe Fang, Xiaoyu Ma, Fei Xu, Li Shen, Liu Tang, Yuxing Zhou, Menglong Huang, et al. 2020. TiDB: a Raft-based HTAP database. *Proceedings of the VLDB Endowment* 13, 12 (2020), 3072–3084.
- [37] Murtadha Al Hubail, Ali Alsuliman, Michael Blow, Michael Carey, Dmitry Ly-chagin, Ian Maxon, and Till Westmann. 2019. Couchbase analytics: NoETL for scalable NoSQL data analysis. *Proceedings of the VLDB Endowment* 12, 12 (2019), 2275–2286.
- [38] Matti Karppa and Rasmus Pagh. 2022. HyperLogLogLog: Cardinality Estimation With One Log More. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (Washington DC, USA) (KDD '22)*. Association for Computing Machinery, New York, NY, USA, 753–761. <https://doi.org/10.1145/3534678.3539246>
- [39] Alfons Kemper and Thomas Neumann. 2011. HyPer: A hybrid OLTP&OLAP main memory database system based on virtual memory snapshots. In *2011 IEEE 27th International Conference on Data Engineering*. IEEE, 195–206.
- [40] Jongbin Kim, Jaeseon Yu, Jaechan Ahn, Sooyong Kang, and Hyungsoo Jung. 2022. Diva: Making MVCC Systems HTAP-Friendly. In *Proceedings of the 2022 International Conference on Management of Data*. 49–64.
- [41] Donald E. Knuth. 1998. *The art of computer programming, volume 3: (2nd ed.) sorting and searching*. Addison Wesley Longman Publishing Co., Inc., USA.
- [42] Tirthankar Lahiri, Shasank Chavan, Maria Colgan, Dinesh Das, Amit Ganesh, Mike Gleeson, Sanket Hase, Allison Holloway, Jesse Kamp, Teck-Hua Lee, et al. 2015. Oracle database in-memory: A dual format in-memory database. In *2015 IEEE 31st International Conference on Data Engineering*. IEEE, 1253–1258.
- [43] Per-Åke Larson, Adrian Birka, Eric N Hanson, Weiyun Huang, Michal Nowakiewicz, and Vassilis Papadimos. 2015. Real-time analytical processing with SQL server. *Proceedings of the VLDB Endowment* 8, 12 (2015), 1740–1751.
- [44] Juchang Lee, SeungHyun Moon, Kyu Hwan Kim, Deok Hoe Kim, Sang Kyun Cha, and Wook-Shin Han. 2017. Parallel replication across formats in SAP HANA for scaling out mixed OLTP/OLAP workloads. *Proceedings of the VLDB Endowment* 10, 12 (2017), 1598–1609.
- [45] Kitaek Lee, Insoon Jo, Jaechan Ahn, Hyuk Lee, Hwang Lee, Woong Sul, and Hyungsoo Jung. 2023. Deploying Computational Storage for HTAP DBMSs Takes More Than Just Computation Offloading. *Proceedings of the VLDB Endowment* 16, 6 (2023), 1480–1493.
- [46] Chen Luo, Pinar Tözün, Yuanyuan Tian, Ronald Barber, Vijayshankar Raman, and Richard Sidle. 2019. Umzi: Unified Multi-Zone Indexing for Large-Scale HTAP. In *EDBT*. 1–12.
- [47] Zhenghua Lyu, Huan Hubert Zhang, Gang Xiong, Gang Guo, Haozhou Wang, Jinbao Chen, Asim Praveen, Yu Yang, Xiaoming Gao, Alexandra Wang, et al. 2021. Greenplum: A Hybrid Database for Transactional and Analytical Workloads. In *Proceedings of the 2021 ACM SIGMOD International Conference on Management of Data*. 2530–2542.
- [48] Darko Makreshanski, Jana Giceva, Claude Barthels, and Gustavo Alonso. 2017. BatchDB: Efficient isolated execution of hybrid OLTP+ OLAP workloads for interactive applications. In *Proceedings of the 2017 ACM International Conference on Management of Data*. 37–50.
- [49] Ryan Marcus, Parimarjan Negi, Hongzi Mao, Nesime Tatbul, Mohammad Alizadeh, and Tim Kraska. 2021. Bao: Making learned query optimization practical. In *Proceedings of the 2021 International Conference on Management of Data*. 1275–1288.
- [50] Ryan Marcus, Parimarjan Negi, Hongzi Mao, Chi Zhang, Mohammad Alizadeh, Tim Kraska, Olga Papaemmanouil, and Nesime Tatbul. 2019. Neo: a learned query optimizer. *Proceedings of the VLDB Endowment* 12, 11 (2019), 1705–1718.
- [51] Norman May, Alexander Böhm, and Wolfgang Lehner. 2017. SAP HANA—The Evolution of an In-Memory DBMS from Pure OLAP Processing Towards Mixed Workloads. *Datenbanksysteme für Business, Technologie und Web (BTW 2017)* (2017).
- [52] John Meehan, Nesime Tatbul, Stan Zdonik, Cansu Aslantas, Ugur Cetintemel, Jiang Du, Tim Kraska, Samuel Madden, David Maier, Andrew Pavlo, et al. 2015. S-Store: Streaming Meets Transaction Processing. *Proceedings of the VLDB Endowment* 8, 13 (2015).
- [53] Elena Milkai, Yannis Chronis, Kevin P. Gaffney, Zhihan Guo, Jignesh M. Patel, and Xiangyao Yu. 2022. How Good is My HTAP System?. In *Proceedings of the 2022 ACM SIGMOD International Conference on Management of Data*. 1810–1824.
- [54] Lili Mou, Ge Li, Lu Zhang, Tao Wang, and Zhi Jin. 2016. Convolutional neural networks over tree structures for programming language processing. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 30.
- [55] Barzan Mozafari, Jags Ramnarayan, Sudhir Menon, Yogesh Mahajan, Soubhik Chakraborty, Hemant Bhanawat, and Kishor Bachhav. 2017. SnappyData: A Unified Cluster for Streaming, Transactions and Interactive Analytics.. In *CIDR*.
- [56] Niloy Mukherjee, Shasank Chavan, Maria Colgan, Dinesh Das, Mike Gleeson, Sanket Hase, Allison Holloway, Hui Jin, Jesse Kamp, Kartik Kulkarni, et al. 2015. Distributed architecture of oracle database in-memory. *Proceedings of the VLDB*



- Endowment* 8, 12 (2015), 1630–1641.
- [57] Vivek Narasayya and Surajit Chaudhuri. 2022. Multi-tenant cloud data services: state-of-the-art, challenges and opportunities. In *Proceedings of the 2022 International Conference on Management of Data*. 2465–2473.
  - [58] Vivek Narasayya, Sudipto Das, Manoj Syamala, Badrish Chandramouli, and Surajit Chaudhuri. 2013. Sqlvm: Performance isolation in multi-tenant relational database-as-a-service. In *CIDR* 2013.
  - [59] Vivek Narasayya, Sudipto Das, Manoj Syamala, Surajit Chaudhuri, Feng Li, and Hyunjung Park. 2013. A demonstration of SQLVM: performance isolation in multi-tenant relational database-as-a-service. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*. 1077–1080.
  - [60] Vivek Narasayya, Ishai Menache, Mohit Singh, Feng Li, Manoj Syamala, and Surajit Chaudhuri. 2015. Sharing buffer pool memory in multi-tenant relational database-as-a-service. *Proceedings of the VLDB Endowment* 8, 7 (2015), 726–737.
  - [61] Krishna Kantikiran Pasupuleti, Boris Klotz, Vijaykrishnan Nagarajan, Ananthakiran Kandukuri, and Nipun Agarwal. 2022. High Availability Framework and Query Fault Tolerance for Hybrid Distributed Database Systems. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 3451–3460.
  - [62] Andrew Pavlo, Gustavo Angulo, Joy Arulraj, Haibin Lin, Jiexi Lin, Lin Ma, Prashanth Menon, Todd C Mowry, Matthew Perron, Ian Quah, et al. 2017. Self-Driving Database Management Systems. In *CIDR*, Vol. 4. 1.
  - [63] Adam Prout, Szu-Po Wang, Joseph Victor, Zhou Sun, Yongzhu Li, Jack Chen, Evan Bergeron, Eric Hanson, Robert Walzer, Rodrigo Gomes, and Nikita Shamgunov. 2022. Cloud-Native Transactions and Analytics in SingleStore. In *Proceedings of the 2022 ACM SIGMOD International Conference on Management of Data*. 2340–2352.
  - [64] Vijayshankar Raman, Gopi Attaluri, Ronald Barber, Naresh Chainani, David Kalmuk, Vincent Kulandaisamy, Jens Leenstra, Sam Lightstone, Shaorong Liu, Guy M Lohman, et al. 2013. DB2 with BLU acceleration: So much more than just a column store. *Proceedings of the VLDB Endowment* 6, 11 (2013), 1080–1091.
  - [65] Aunn Raza, Periklis Chrysogelos, Angelos Christos Anadiotis, and Anastasia Ailamaki. 2020. Adaptive HTAP through elastic resource scheduling. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 2043–2054.
  - [66] Hemant Saxena, Lukasz Golab, Stratos Idreos, and Ihab F Ilyas. 2021. Real-Time LSM-Trees for HTAP Workloads. *arXiv preprint arXiv:2101.06801* (2021).
  - [67] Hemant Saxena, Lukasz Golab, Stratos Idreos, and Ihab F Ilyas. 2023. Real-time LSM-trees for HTAP workloads. In *2023 IEEE 39th International Conference on Data Engineering (ICDE)*. IEEE, 1208–1220.
  - [68] Mohit Saxena, Benjamin Sowell, Daiyan Alamgir, Nitin Bahadur, Bijay Bisht, Santosh Chandrachood, Chhitti Keswani, G Krishnamoorthy, Austin Lee, Bohou Li, et al. 2023. The Story of AWS Glue. *Proceedings of the VLDB Endowment* 16, 12 (2023), 3557–3569.
  - [69] Raghav Sethi, Martin Traverso, Dain Sundstrom, David Phillips, Wenlei Xie, Yutian Sun, Nezhir Yegitbasi, Haozhun Jin, Eric Hwang, Nileema Shingte, and Christopher Berner. 2019. Presto: SQL on Everything. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*. 1802–1813. <https://doi.org/10.1109/ICDE.2019.00196>
  - [70] Sijie Shen, Rong Chen, Haibo Chen, and Binyu Zang. 2021. Retrofitting High Availability Mechanism to Tame Hybrid Transaction/Analytical Processing. In *15th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 21)*. 219–238.
  - [71] Vishal Sikka, Franz Färber, Wolfgang Lehner, Sang Kyun Cha, Thomas Peh, and Christof Bornhövd. 2012. Efficient transaction processing in SAP HANA database: the end of a column store myth. In *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*. 731–742.
  - [72] Jason Sun, Haoxiang Ma, Li Zhang, Huicong Liu, Haiyang Shi, Shangyu Luo, Kai Wu, Kevin Bruhwiler, Cheng Zhu, Yuanyuan Nie, et al. 2023. Accelerating Cloud-Native Databases with Distributed PMem Stores. In *2023 IEEE 39th International Conference on Data Engineering (ICDE)*. IEEE, 3043–3057.
  - [73] Lalitha Viswanathan, Alekh Jindal, and Konstantinos Karanasos. 2018. Query and resource optimization: Bridging the gap. In *2018 IEEE 34th International Conference on Data Engineering (ICDE)*. IEEE, 1384–1387.
  - [74] Jianying Wang, Tongliang Li, Haoze Song, Xinjun Yang, Wenchao Zhou, Feifei Li, Baoyue Yan, Qianqian Wu, Yukun Liang, ChengJun Ying, Yujie Wang, Baokai Chen, Chang Cai, Yubin Ruan, Xiaoyi Weng, Shibin Chen, Liang Yin, Chengzhong Yang, Xin Cai, Hongyan Xing, Nanlong Yu, Xiaofei Chen, Dapeng Huang, and Jianling Sun. 2023. PolarDB-IMCI: A Cloud-Native HTAP Database System at Alibaba. *Proc. ACM Manag. Data* 1, 2, Article 199 (jun 2023), 25 pages. <https://doi.org/10.1145/3589785>
  - [75] Louis Woods, Istvan Zsolt, and Gustavo Alonso. 2014. Ibox: An Intelligent Storage Engine with Support for Advanced SQL Off-loading. *Proceedings of the VLDB Endowment* 7, 11 (2014), 963–974. <https://doi.org/10.14778/2732967.2732972>
  - [76] Ziniu Wu, Ryan Marcus, Zhengchun Liu, Parimarjan Negi, Vikram Nathan, Pascal Pfeil, Gaurav Saxena, Mohammad Rahman, Balakrishnan Narayanaswamy, and Tim Kraska. 2024. Stage: Query Execution Time Prediction in Amazon Redshift. In *Companion of the 2024 International Conference on Management of Data*. 280–294.
  - [77] Xianghong Xu, Zhibing Zhao, Tieying Zhang, Rong Kang, Luming Sun, and Jianjun Chen. 2023. COOOL: A Learning-To-Rank Approach for SQL Hint Recommendations. *5th International Workshop on Applied AI for Database Systems and Applications* (2023).
  - [78] Maryann Xue, Yingyi Bu, Abhishek Somani, Wenchen Fan, Ziqi Liu, Steven Chen, Herman van Hovell, Bart Samwel, Mostafa Mokhtar, RK Korlapati, et al. 2024. Adaptive and Robust Query Execution for Lakehouses at Scale. *Proceedings of the VLDB Endowment* 17, 12 (2024), 3947–3959.
  - [79] Jiacheng Yang, Ian Rae, Jun Xu, Jeff Shute, Zhan Yuan, Kelvin Lau, Qiang Zeng, Xi Zhao, Jun Ma, Ziyang Chen, et al. 2020. F1 Lightning: HTAP as a Service. *Proceedings of the VLDB Endowment* 13, 12 (2020), 3313–3325.
  - [80] Xinjun Yang, Yingqiang Zhang, Hao Chen, Feifei Li, Bo Wang, Jing Fang, Chuan Sun, and Yuhui Wang. 2024. PolarDB-MP: A Multi-Primary Cloud-Native Database via Disaggregated Shared Memory. In *Proceedings of the 2024 ACM SIGMOD International Conference on Management of Data*. 295–308.
  - [81] Xinjun Yang, Yingqiang Zhang, Hao Chen, Chuan Sun, Feifei Li, and Wenchao Zhou. 2023. PolarDB-SCC: A Cloud-Native Database Ensuring Low Latency for Strongly Consistent Reads. *Proc. VLDB Endow.* 16, 12 (aug 2023), 3754–3767. <https://doi.org/10.14778/3611540.3611562>
  - [82] Zongheng Yang, Wei-Lin Chiang, Sifei Luan, Gautam Mittal, Michael Luo, and Ion Stoica. 2022. Balsa: Learning a Query Optimizer Without Expert Demonstrations. In *Proceedings of the 2022 International Conference on Management of Data (Philadelphia, PA, USA) (SIGMOD/PODS '22)*. Association for Computing Machinery, New York, NY, USA, 931–944. <https://doi.org/10.1145/3514221.3517885>
  - [83] Zhifeng Yang, Quanqing Xu, Shanyan Gao, Chuanhui Yang, Guoping Wang, Yuzhong Zhao, Fanyu Kong, Hao Liu, Wanhong Wang, and Jinliang Xiao. 2023. OceanBase Paetica: A Hybrid Shared-Nothing/Shared-Everything Database for Supporting Single Machine and Distributed Cluster. *Proceedings of the VLDB Endowment* 16, 12 (2023), 3728–3740. <https://doi.org/10.14778/3611540.3611560>
  - [84] Chaohun Zhan, Maomeng Su, Chuangxian Wei, Xiaoqiang Peng, Liang Lin, Sheng Wang, Zhe Chen, Feifei Li, Yue Pan, Fang Zheng, and Chengliang Chai. 2019. AnalyticDB: real-time OLAP database system at Alibaba cloud. *Proc. VLDB Endow.* 12, 12 (aug 2019), 2059–2070. <https://doi.org/10.14778/3352063.3352124>
  - [85] Ning Zhang, Junichi Tatemura, Jignesh Patel, and Hakan Hacigumus. 2014. Re-evaluating designs for multi-tenant OLTP workloads on SSD-based/O subsystems. In *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data*. 1383–1394.
  - [86] Rong Zhu, Wei Chen, Bolin Ding, Xingguang Chen, Andreas Pfadler, Ziniu Wu, and Jingren Zhou. 2023. Lero: A learning-to-rank query optimizer. *Proceedings of the VLDB Endowment* 16, 6 (2023), 1466–1479.
  - [87] Rong Zhu, Lianggui Weng, Bolin Ding, and Jingren Zhou. 2024. Learned Query Optimizer: What is New and What is Next. In *Companion of the 2024 International Conference on Management of Data*. 561–569.