



Degree Project in Computer Science

Second cycle, 30 credits

IcedHops: reducing read and write latency in an Iceberg-backed offline feature store

Integrating HopsFS and Pylceberg Python library to reduce read and write latency on Hopsworks offline feature store, with a comparative analysis of alternative solutions

SEBASTIANO MENEGHIN

IcedHops: reducing read and write latency in an Iceberg-backed offline feature store

Integrating HopsFS and Pylceberg Python library to reduce read and write latency on Hopsworks offline feature store, with a comparative analysis of alternative solutions

SEBASTIANO MENEGHIN

Master's Programme, ICT Innovation, 120 credits

Date: February 27, 2025

Supervisors: Sina Sheikholeslami, Fabian Schmidt, Davit Bzhalava

Examiner: Vladimir Vlassov

School of Electrical Engineering and Computer Science

Host company: Hopsworks AB

Swedish title: IcedHops: minska läs- och skrivlatens i en Iceberg-stödd offlinebutik för funktioner

Swedish subtitle: Integrering av HopsFS och Pylceberg Python-bibliotek för att minska läs- och skriva fördröjningen i den offline-lagringen Hopsworks, med en jämförande analys av alternativa lösningar.

Abstract

The growing need for efficient data management in Machine Learning (ML) workflows has led to the widespread adoption of feature stores, centralized data platforms that supports feature engineering, model training and prediction inference. The Hopsworks' feature store has demonstrated outperformance compared to its alternatives, leveraging Apache Hudi and Spark for offline data storage, but suffers from high write and read latency, even for small quantities of data (1GB or less). This thesis explores the potential of Apache Iceberg as an alternative table format to reduce latencies, developing "IcedHops", an integration of HopsFS (Hopsworks HDFS distribution) and PyIceberg Python library.

The research begin with an evaluation of potential system integration alternatives, documenting the advantages and limitations of each approach. Then, IcedHops is implemented and evaluated, benchmarking it against the existing Spark-based solution and an alternative Delta Lake implementation (delta-rs). Extensive experiments were conducted across varying table sizes and CPU configurations to assess write and read performance. Results show that IcedHops significantly reduces write latency – from 40 to 140 times lower than the legacy system – and read latency – from 55% to 60 times lower than the legacy system. Compared to delta-rs, IcedHops demonstrates reduced write latency for large tables – up to 7 times lower – and equal read latency, but exhibits lower scaling benefits with additional CPU cores – 20% less than delta-rs.

These findings confirm that alternatives to Spark-based pipelines in small-scale scenarios are possible and are worth of further investigations, and the system implemented will be included in the Hopsworks feature store. Furthermore, this thesis work and results finally provides a baseline for future work about additional open table formats, alternative languages to mitigate Python's performance overhead, and strategies to improve resource utilization in data management platforms.

Keywords

Machine Learning, Feature Store, Spark, Apache Iceberg, Delta Lake, Python, Read/Write Latency, Open Table Formats

Sammanfattning

Det växande behovet av effektiv datahantering i arbetsflöden för maskininlärning (ML) har lett till en utbredd användning av feature stores – centraliserade dataplattformar som stöder feature engineering, modellträning och inferens. Hopsworks feature store har visat bättre prestanda jämfört med sina alternativ och använder Apache Hudi och Spark för offline-datalagring. Dock lider systemet av hög skriv- och läslatens, även för små datamängder (1 GB eller mindre). Denna avhandling undersöker potentialen hos Apache Iceberg som ett alternativt tabellformat och integrerar det med HopsFS (Hopsworks HDFS-distribution) samt PyIceberg Python-biblioteket för att minska latensen.

Forskningen inleds med en utvärdering av potentiella systemintegrationsalternativ, där fördelar och begränsningar med varje metod dokumenteras. Därefter implementeras och utvärderas en PyIceberg-baserad arkitektur, vilken jämförs med den befintliga Spark-baserade lösningen samt en alternativ Delta Lake-implementering (delta-rs). Omfattande experiment genomfördes med varierande tabellstorlekar och CPU-konfigurationer för att bedöma skriv- och läsprestandorna. Resultaten visar att PyIceberg avsevärt minskar skrifftördröjningen – från 40 till 140 gånger lägre än det äldre systemet – och läsfördröjningen – från 55% till 60 gånger lägre än det äldre systemet. Jämfört med delta-rs uppvisar PyIceberg minskad skrifftördröjning för stora tabeller – upp till sju gånger lägre – och liknande läsfördröjning, men har sämre skalningsfördelar vid ökning av CPU-kärnor (20% mindre än delta-rs).

Dessa resultat bekräftar att alternativ till Spark-baserade pipelines i småskaliga scenarier är möjliga och värda ytterligare undersökningar. Det implementerade systemet kommer att integreras i Hopsworks feature store. Dessutom utgör denna avhandling en baslinje för framtida forskning kring ytterligare öppna tabellformat, alternativa programmeringsspråk för att hantera Pythons prestandabegränsningar samt strategier för att förbättra resursutnyttjandet i datahanteringsplattformar.

Nyckelord

Maskininlärning, Feature Store, Spark, Apache Iceberg, Delta Lake, Python, Läs- och skrivilatens, Open Table Formats

Sommario

La crescente necessità di piattaforme per una efficienza gestione dei dati per applicazioni di *Machine Learning (ML)* ha portato a un'ampia diffusione dei *feature store*, piattaforme dati centralizzate che supportano *feature engineering*, addestramento di modelli e inferenza. Il *feature store* di Hopsworks ha dimostrato prestazioni superiori rispetto alle sue alternative, sfruttando Apache Hudi e Spark per il suo *offline feature store*. Tuttavia, questo sistema soffre di un'elevata latenza di scrittura e lettura, anche per piccole quantità di dati (1GB o inferiori). Questa tesi esplora l'uso di Apache Iceberg come *table format* alternativo per ridurre tali latenze, sviluppando "IcedHops", una integrazione tra HopsFS (distribuzione Hopsworks di HDFS) e la libreria Python PyIceberg.

Questa ricerca inizia con un'analisi delle possibili strategie di integrazione, documentando i vantaggi e i limiti di ciascun approccio. Successivamente, IcedHops viene sviluppato ed evaluato, confrontandolo con la soluzione esistente basata su Spark e con un'alternativa basata su Delta Lake (delta-rs). Esperimenti approfonditi sono stati condotti su tabelle di diverse dimensioni e diverse configurazioni CPU per misurare le prestazioni di scrittura e lettura. I risultati mostrano che IcedHops riduce significativamente la latenza di scrittura – da 40 a 140 volte inferiore rispetto al sistema legacy – e la latenza di lettura – dal 55% a 60 volte inferiore. Rispetto a delta-rs, IcedHops offre una latenza di scrittura inferiore per le tabelle più grandi – fino a 7 volte minore – e prestazioni di lettura equivalenti, ma mostra minori vantaggi di scalabilità con l'aumento dei CPU cores (20% ridotti rispetto a delta-rs).

Questi risultati confermano che alternative ad architetture basate su Spark, per gestire dati su piccola scala, esistono e sono più efficienti, e dunque il sistema sviluppato verrà integrato nel *feature store* di Hopsworks. Inoltre, questa tesi fornisce una base per futuri studi su nuovi *open table format* e strategie da adottare per ottimizzare l'uso delle risorse nelle piattaforme di gestione dei dati.

Parole Chiave

Machine Learning, Feature Store, Apache Iceberg, Delta Lake, PyIceberg, Latenza in Lettura/Scrittura, Open Table Formats.

vi | Sommario

Acknowledgments

“The only true voyage of discovery would be not to visit strange lands but to possess other eyes., – Marcel Proust

The completion of this master thesis was only possible thanks to the invaluable support of many, to whom I am deeply grateful.

Many thanks to Jim Dowling, how offered me the opportunity to work on this project and Davit Bzhalava, who supported me along the way. I also had the great pleasure of collaborating with the employees of Hopsworks AB, who generously shared their patience and expertise.

Many thanks to my KTH examiner, Vladimir Vlassov, and my supervisors, Fabian Schmidt and Sina Sheikholeslami, for their guidance in conducting thorough research and for their invaluable feedback.

Heartfelt thanks to my parents and relatives, who have supported me, who have inspired me, who have believed in me, giving me a blank canvas with colors they had never used before. Your wisdom, guidance, and the values you instilled in me have shaped who I am today – and the person I strive to become in the future.

Lastly, I cannot begin to express my gratitude to the friends who have stood by my side through adventures and challenges. Michele, Daniele, and Davide, for countless years, you offered me advice and challenged my ideas, guiding me towards a clearer understanding of myself. Alfonso, Giacomo, Luca, and Virginia, you made feel like home, at every step of this journey, no matter where we were. Giovanni, you have filled my mind with 'whys' and 'hows' – questions that helped me grow at every answer, though some still await resolution. Giulia, and all the incredible people from MESA, you have shown me how essential it is to have a cause to believe in and care about. May the future bring us even more time to learn, discover, and laugh – together.

Stockholm, February 2025

Sebastiano Meneghin

Contents

1	Introduction	1
1.1	Background	4
1.2	Problem	6
1.2.1	Research Questions	7
1.3	Purpose	8
1.4	Goals	8
1.5	Ethics and Sustainability	9
1.6	Research Methodology	10
1.7	Research Limitation	11
1.8	Structure of the thesis	12
2	Background and Related work	13
2.1	Data storage	14
2.1.1	Block storage vs. File storage vs. Object storage	14
2.1.2	Hadoop Distributed File System	17
2.1.3	HopsFS	18
2.1.4	Cloud object stores, an alternative	20
2.2	Data management	20
2.2.1	Brief history of Data Base Management Systems	21
2.2.2	Data lakehouse architecture	23
2.2.3	Data lakehouse comparison	25
2.3	Query engine	29
2.3.1	Apache Spark	29
2.3.2	Apache Kafka	30
2.3.3	Catalogs	31
2.3.4	Duck DB	32
2.3.5	Arrow Flight	33
2.4	Application - Hopsworks	33
2.4.1	Machine Learning Operations (MLOps)	33

2.4.2	Hopsworks AI Data Platform	34
2.5	System architectures	35
2.5.1	Legacy system - Hudi - writing	36
2.5.2	Legacy system - Hudi - reading	37
2.5.3	IcedHops - writing	37
2.5.4	IcedHops - reading	37
2.5.5	delta-rs - writing	39
2.5.6	delta-rs - reading	39
3	Method	41
3.1	System integration	41
3.1.1	Integration process	42
3.1.2	Requirements	43
3.1.3	Development environment	44
3.2	System evaluation - Legacy vs. IcedHops	44
3.2.1	Evaluation process - RQ2	44
3.2.2	Industrial use case	46
3.2.3	Experimental data	46
3.2.4	Experimental design	48
3.2.5	Experimental environment	49
3.2.6	Evaluation framework	50
3.2.7	Reliability and validity	51
3.3	System evaluation - IcedHops vs. Delta Lake	51
3.3.1	Evaluation process - RQ3	52
3.3.2	Evaluation framework	53
4	Implementation	55
4.1	Integration design and usage	55
4.1.1	Catalog choice	56
4.1.2	Query engine choice	57
4.1.3	Usage	58
4.2	Experimental setup	59
5	Results and Analysis	63
5.1	Major results	63
5.1.1	Write experiments	64
5.1.2	Read experiments	67
5.1.3	Legacy pipeline write latency breakdown	70
5.1.4	In-memory resources usage	72
5.2	Results analysis and discussion	72

5.2.1	Discussion on major results	72
5.2.2	Considerations on the legacy system	74
5.2.3	Considerations on the PyIceberg library	75
6	Conclusions and Future work	77
6.1	Conclusions	77
6.2	Limitations	79
6.3	Future work	80
References		81
A	Technology comparisons	89
B	System architectures	92
C	Write experiments results	95
D	Read experiments results	104
E	Legacy pipeline write latency breakdown results	113

List of Figures

1.1	Sustainable Development Goals supported by this thesis	10
2.1	Data stack abstraction	13
2.2	Hadoop Distributed File System architecture	19
2.3	Architecture of data lakehouse	24
2.4	GitHub stars of Open Table Formats (OTFs) repositories	26
2.5	Feature store in an MLOps pipeline	35
2.6	Legacy system - Hudi - write process	36
2.7	Legacy system - Hudi - read process	37
2.8	IcedHops - write process	38
2.9	IcedHops - read process	38
2.10	delta-rs - write process	39
2.11	delta-rs - read process	39
3.1	System integration process	43
3.2	System evaluation process - Legacy vs. IcedHops	45
3.3	System evaluation process - IcedHops vs. delta-rs	52
5.1	Histogram of the write experiment - Latency - 1 CPU core . .	65
5.2	Histogram of the write experiment - Throughput - 1 CPU core	66
5.3	Histogram of the read experiment - Latency - 1 CPU core . .	68
5.4	Histogram of the read experiment - Throughput - 1 CPU core .	69
5.5	Histogram of the write on legacy pipeline - Time breakdown - 1 core	71
B.1	Legacy system - Write process - Magnified diagram	93
B.2	Legacy system - Read process - Magnified diagram	94
C.1	Histogram of the write experiment - Latency - 1 CPU core . .	96
C.2	Histogram of the write experiment - Latency - 2 CPU cores .	97
C.3	Histogram of the write experiment - Latency - 4 CPU cores .	98

C.4	Histogram of the write experiment - Latency - 8 CPU cores . . .	99
C.5	Histogram of the write experiment - Throughput - 1 CPU core . . .	100
C.6	Histogram of the write experiment - Throughput - 2 CPU cores . . .	101
C.7	Histogram of the write experiment - Throughput - 4 CPU cores . . .	102
C.8	Histogram of the write experiment - Throughput - 8 CPU cores . . .	103
D.1	Histogram of the read experiment - Latency - 1 CPU core . . .	105
D.2	Histogram of the read experiment - Latency - 2 CPU cores . . .	106
D.3	Histogram of the read experiment - Latency - 4 CPU cores . . .	107
D.4	Histogram of the read experiment - Latency - 8 CPU cores . . .	108
D.5	Histogram of the read experiment - Throughput - 1 CPU core . .	109
D.6	Histogram of the read experiment - Throughput - 2 CPU cores . .	110
D.7	Histogram of the read experiment - Throughput - 4 CPU cores . .	111
D.8	Histogram of the read experiment - Throughput - 8 CPU cores . .	112
E.1	Histogram of write on legacy pipeline - Latency breakdown - 1 CPU core	114
E.2	Histogram of write on legacy pipeline - Latency breakdown - 2 CPU cores	115
E.3	Histogram of write on legacy pipeline - Latency breakdown - 4 CPU cores	116
E.4	Histogram of write on legacy pipeline - Latency breakdown - 8 CPU cores	117

List of Tables

2.1	Data storage features comparison	15
2.2	Data storage pros and cons comparison	17
5.1	Write experiments results expressed as latency	65
5.2	Write experiments results expressed as throughput	66
5.3	Read experiments results expressed as latency	68
5.4	Read experiments results expressed as throughput	69
5.5	Writes on legacy pipeline - Time breakdown	71
A.1	Comparison of data architectures	90
A.2	Comparison of data lakehouses	91
C.1	Write experiment - Latency - 1 CPU core	96
C.2	Write experiment - Latency - 2 CPU cores	97
C.3	Write experiment - Latency - 4 CPU cores	98
C.4	Write experiment - Latency - 8 CPU cores	99
C.5	Write experiment - Throughput - 1 CPU core	100
C.6	Write experiment - Throughput - 2 CPU cores	101
C.7	Write experiment - Throughput - 4 CPU cores	102
C.8	Write experiment - Throughput - 8 CPU cores	103
D.1	Read experiment - Latency - 1 CPU core	105
D.2	Read experiment - Latency - 2 CPU cores	106
D.3	Read experiment - Latency - 4 CPU cores	107
D.4	Read experiment - Latency - 8 CPU cores	108
D.5	Read experiment - Throughput - 1 CPU core	109
D.6	Read experiment - Throughput - 2 CPU cores	110
D.7	Read experiment - Throughput - 4 CPU cores	111
D.8	Read experiment - Throughput - 8 CPU cores	112
E.1	Write on legacy pipeline - Latency breakdown - 1 CPU core .	114

E.2	Write on legacy pipeline - Latency breakdown - 2 CPU cores .	115
E.3	Write on legacy pipeline - Latency breakdown - 4 CPU cores .	116
E.4	Write on legacy pipeline - Latency breakdown - 8 CPU cores .	117

Listings

3.1	Experimental environment details	50
4.1	Instanciate Iceberg catalog with SQLite	58
4.2	Writing with IcedHops	59
4.3	Reading with IcedHops	59
4.4	Measuring latency using Timeit	60
4.5	Measuring latency using the time difference	61

List of acronyms and abbreviations

ACID	Atomicity, Consistency, Isolation and Durability
AI	Artificial Intelligence
API	Application Programming Interface
AWS	Amazon Web Services
BI	Business Intelligence
CC	Concurrency Control
CDC	Change Data Capture
CIDR	Conference on Innovative Data Systems Research
CoW	Copy on Write
CPU	Central Processing Unit
CRUD	Create Read Update Delete
D	Deliverable
DBMS	Data Base Management System
ELT	Extract Load Transform
ETL	Extract Transform Load
G	Goal
GCS	Google Cloud Storage
HDD	Hard Disk Drive
HDFS	Hadoop Distributed File System
HMS	Hive Metastore
HopsFS	Hopsworks' HDFS distribution
IN	Industrial Need
JDBC	Java DataBase Connectivity
JVM	Java Virtual Machine
LAN	Local Area Network
LOC	Lines Of Code

ML	Machine Learning
MLOps	Machine Learning Operations
MoR	Merge on Read
NAS	Network Attached Storage
NN	Neural Network
OLAP	On-Line Analytical Processing
OTF	Open Table Format
PA	Project Assumption
PC	Personal Computer
RAM	Random Access Memory
RDD	Resilient Distributed Dataset
RQ	Research Question
SAN	Storage Area Network
SDG	Sustainable Development Goal
SDK	Software Development Kit
SF	Scale Factor
SQL	Structured Query Language
SSD	Solid State Drive
SSH	Secure Shell protocol
TPC	Transaction Processing Performance Council
VM	Virtual Machine

Chapter 1

Introduction

Data lakehouses' adoption reached in 2024 a critical mass, with more than 50% organizations running their analytics on these platform, and projection up to 67% within 2027 [1]. A data lakehouse is a modern data architecture that creates a single platform by combining the key benefits of data lakes (large repositories of raw data in its original form) and data warehouses (organized sets of structured data). Thus, data lakehouse is a cost-effective architecture [2] that bridges those architectures by providing the scalability of the data lakes together with analytical computations and **Atomicity, Consistency, Isolation and Durability (ACID)** properties of data warehouses [3]. This surge is strongly related to the acceleration of AI adoption [4], which put pressure on storage solutions and data processing capabilities [1]. In support to this, data lakehouse systems include data partitioning, reducing the query computational needs, support schemas evolution and provide "time travel" capabilities, enabling data versioning over time [5].

From the first appearance of data lakehouses [6], three primary applications arose [7]:

1. **Apache Hudi:** introduced by Uber in 2017, then supported by Alibaba, Bytedance, Uber and Tencent [8].
2. **Apache Iceberg:** open sourced by Netflix * in 2018, now used by Airbnb, Apple, Expedia, LinkedIn, Lyft [9].
3. **Delta Lake:** open sourced by Databricks in 2019 [10], supported now by Databricks and Microsoft.

*Netflix first implementation at <https://github.com/Netflix/iceberg>

Apache Iceberg (from now on, Iceberg) was originally developed as an alternative to Apache Hive (from now on, Hive), a data warehouse system that enables querying of large datasets stored in Hadoop using a **Structured Query Language (SQL)**-like language. Iceberg is designed to manage bigger datasets with frequent updates and deletions, and to support schema evolution, going thus beyond Hive's capabilities [11]. This is possible thanks to the Iceberg metadata management layer, which enables data warehouse-like capabilities via cloud object storage. In 2024, Iceberg has consolidated its position in 2024 [12], and has been recently included in arguably all the major big data vendors' solutions [13].

Apache Hudi (from now on, Hudi), Iceberg and Delta Lake were initially designed to fit a single engine, respectively Hive, Trino and Apache Spark (from now on, Spark) [7]. These three solutions have been improved to better scalability and flexibility needs, each integrating new engners and features [14]. Despite being beneficial, this evolution highlights a potential limitation: systems developed in different ways performe differently, and depending on the use case, even widespread solutions might not always be the optimal solution. For example, integrating the system with Spark, a query and compute engine [15], is a good solution when processing massive datasets (1 TB+) on the cloud. However, it is yet unclear how good this solution could perform when processing smaller datasets (1 GB to 100 GB) [16].

Such Spark limitations of Spark are also brought to light by technologies such as the Python library Polars [17] and the **Data Base Management System (DBMS)** DuckDB [18], that combined to process data locally, in small quantities, showed much better performance than a Spark Cluster. This solution further avoids the high expenses and computation time lying behind Spark usage in the same scenario. Thus, similar solutions employed in small-scale scanarios (from 1GB to 100GB) could greatly increase performance and reduce costs [19] [20].

Remaining in the data science domain, Python has perhaps become the de facto standard programming language. Python is arguably the most popular general-purpose programming language [21, 22, 23], and the go-to option for **Machine Learning (ML)** and **Artificial Intelligence (AI)** applications [24]. Among the most commonly used libraries by **ML** application developers, we find NumPy and Pandas [22], while PyTorch and TensorFlow are the first two for **Neural Network (NN)** libraries. Accessing data lakehouse solutions via a Python client would be the favorite option for most of **ML** and **AI** developers, so they do not have to resort to alternatives, such as Spark and its **Application Programming Interfaces (APIs)** (e.g., PySpark).

Therefore, supporting Iceberg with a native Python client would be directly beneficial for Hopsworks AB, the host company of this master thesis. Hopsworks AB develops a homonymous **AI Lakehouse for ML**. This software includes a multi-user data platform, called feature store, that enables storage and access of reusable features *. It also implements resource-efficient and point-in-time joins between datasets, adding historical retrieval features to saved data [25].

This project seeks to reduce latency (in seconds), thus increasing throughput (in rows/second), of reading and writing data on Iceberg tables, as offline feature store in Hopsworks. At time of writing, the writing pipeline is Spark-based and uses Hudi as **Open Table Format (OTF)**, while this thesis seeks to employ Iceberg without Spark-dependencies. If the hypothesis that a faster non-Spark alternative is possible is proved, Hopsworks AB will consider this alternative as an integration or replacement of their current open-source feature store implementation. This could considerably improve the experience of Python users working with small-scale datasets (1 GB to 100 GB). More broadly, this thesis will discuss the feasibility of Spark alternatives in small-scale use scenarios [26].

The primary contributions of this work are:

- The evaluation of the possible technologies for the integrations between Iceberg and Hopsworks feature store, and the selection of the best one according to defined requirements, described in Section 3.1.2. IcedHops employs the technologies most fitting the requirements, consisting of SQLCatalog as catalog, DuckDB as query engine, and PyArrowFileIO as FileIO.
- The results of the experiments conducted on IcedHops and the Hudi legacy system, showing differences in latency and throughput for read and write operations. These experiments were conducted fifty times, with different dataset sizes and **Central Processing Unit (CPU)** settings. IcedHops, accessing Iceberg from a Python client, reduced write latency by 40 to 140 times, and read latency by 15 to 60 times in the smaller tables, compared to what was experienced using the old system.
- The results of the comparison between IcedHops and the recently implemented delta system [26], on which the same experiments above were conducted in a parallel thesis work. IcedHops has write latency from 5 to 7 times lower compared to what was experienced using the

*Definition from the company's website at <https://www.hopsworks.ai/>

delta-rs system, and comparable read latency. The delta-rs system showed better support for multiple **CPU** cores setting.

These findings provide a significant contribution to the data management industry, supporting existing studies on the constraints of utilising Spark with small-scale volumes of data. Furthermore, the reproducibility of the experiments conducted adds great value to this work. Thanks to the well-defined environment, in Section 4.2 and the available code ^{*}, these experiments can be used as starting or final point for further exploration and testing in this field.

1.1 Background

There are three crucial components for a comprehensive understanding of this work: the evolution of data infrastructure to data lakehouses, the role of Spark as a data management tool, and the rise of Python as the most used programming language in the data science field.

The term "data lakehouse" was used by Databricks in 2020 [6] to characterize a new architectural standard developing across the industry. This novel paradigm integrates the data lake's capacity to store and manage unstructured data with the **ACID** features characteristic of data warehouses. Data warehouses became the standard in the 1990s and early 2000s [27], facilitating firms in deriving business intelligence insights from data coming from various structured data sources. The architectural issues of this technology became evident at the end of the 2010s, with the rise of Big Data, characterized by increasing volumes, variety, and velocity of data, including significant amounts of unstructured information [28]. Data lakes were thus introduced, as a more flexible and scalable solution, serving as a central repository for all data. They enable the development of more intricate built-upon architectures, including data warehouses for **Business Intelligence (BI)** and **ML** pipelines. While being more appropriate for unstructured data, this architecture entails several complications and expenses associated with the need for duplicated data (data lake and data warehouse) and complex **Extract Load Transform (ELT)/Extract Transform Load (ETL)** pipelines. Data lakehouse solutions addressed the challenges of data lakes by combining the best of both worlds: the flexibility and scalability of data lakes with the data governance, reliability, and performance of data warehouses. This is achieved

^{*}Code available at: <https://github.com/SebastianoMeneghin/apache-iceberg-as-offline-feature-store>

by integrating data management and performance capabilities directly into open data formats like Parquet [29]. Three pivotal technologies facilitated this paradigm:

- A metadata layer, providing data lineage and facilitating efficient data discovery and access.
- An optimized query engine, leveraging techniques like **Random Access Memory (RAM)/Solid State Drive (SSD)** caching and advanced query optimization.
- A user-friendly **API**, simplifying data access and integration with **ML** and **AI** applications.

Uber first open-sourced this architectural design with Hudi in 2017 [8], followed by Netflix in 2018, with Iceberg [9], and ultimately by Databricks, with Delta Lake in 2020 [10].

Spark is a distributed computing platform designed to facilitate large-scale, data-intensive applications [30]. Developed as an improvement over Hadoop MapReduce (from now on just MapReduce), Spark addresses its limitations, such as high latency due to disk I/O and limited support for iterative algorithms. Spark achieves significantly higher performance by leveraging in-memory computing, eliminating the need for frequent disk access to store intermediate results. This, coupled with its DAG execution model, **Resilient Distributed Datasets (RDDs)**, and support for multiple processing models (batch, stream, **ML**, graph), has established Spark as the de facto standard for many data-intensive workloads. While Spark offers a powerful and versatile platform, it may not always be the most suitable choice for all scenarios. For instance, Apache Flink [31], specifically designed for stream and real-time processing, excels in low-latency applications where Spark may exhibit higher latency. Similarly, for small-scale datasets (from 1 GB to 100 GB), the overhead associated with launching and managing a Spark cluster can be substantial. In such cases, in-memory **DBMS** like DuckDB [18] and Polars [17] offer compelling alternatives. These solutions, optimized for smaller datasets, provide high-performance **On-Line Analytical Processing (OLAP)** capabilities and DataFrame operations within an embedded environment, delivering significantly faster results compared to initiating a Spark cluster for equivalent tasks. This project investigates the feasibility of employing alternatives to Spark for efficient data processing on small-scale datasets, exploring their potential advantages in terms of performance.

Python reigns supreme as the programming language of choice for data scientists [32]. Its user-friendliness, high-level abstraction, and emphasis on clear code expression initially attracted a large and enthusiastic community [22]. This supportive community, in turn, fueled the development of a vast ecosystem of libraries and APIs specifically designed for data science tasks. Over three decades, Python has become the de facto standard in this domain, with popular libraries like NumPy, Pandas, PyTorch, and TensorFlow, forming the backbone of countless data science projects. Python is regarded as the most popular programming language based on the volume of search results for the query (+"*<language> programming*") across 25 distinct search engines *. The TIOBE index [21] and 2024 GitHub report [23] underscore the current trends, clearly demonstrating Python's ascendancy, also evident from the following milestones:

- Python surpasses C in 2021.
- Python surpasses Java in 2022.
- Python surpasses JavaScript in 2024.

These results emphasize the need for offering Python APIs, especially for programmers and data scientists, to augment interaction and broaden the framework's possibilities.

1.2 Problem

The Hopsworks feature store [33] first used Hudi for its offline feature store, since in 2017 it was the first data lakehouse to be open-sourced. Hopsworks AB thus actively implemented new technologies for their software, to cater growing customer's needs. In the legacy system, Spark serves as the query engine, executing queries (read, write, or delete) on the offline feature store. This system demonstrated that a write operation, even on a tiny dataset (1 GB or less), takes one or more minutes to finalise.

This adversely affects Hopsworks' standard use case, which operates between testing scenarios on tiny data volumes (from 1 GB to 10 GB) and production scenario on higher volumes, despite still with relative small volumes (from 10 GB to 100 GB).

*Evaluation methodology defined at https://www.tiobe.com/tiobe-index/programminglanguages_definition/

The fundamental hypothesis of this study is that the prolonged transaction time is a problem determined by Spark. This has prompted Hopsworks to implement Spark alternatives [16] for data ingestion in their Hudi system, and to search for alternatives to the Hudi-Spark tandem [26]. Iceberg provides support for several alternative to Spark [14], and Iceberg tables can be accessed directly from Python, via the PyIceberg library *. However, the Hopsworks feature store has not been yet integrated with Iceberg, and its underlying file system, **Hopsworks' HDFS distribution (HopsFS)** [34], and several implementations are possible [11].

1.2.1 Research Questions

This research project aims to assess and compare the performance of the legacy system, relying on Hudi and Spark, against IcedHops, a newly developed systems providing an alternative to Spark, and a comprehensive comparison of the latters. Those systems features the PyIceberg library *, which operates on Iceberg tables, and the Rust delta-rs library †, which operates on Delta Lake tables, in both cases hosted on **HopsFS**. To do this, an integration of Iceberg with the Hopsworks feature store must be designed. Accordingly, this study tackles the three following **Research Questions (RQs)**:

- RQ1: Is a working integration of PyIceberg and HopsFS in the Hopsworks offline feature store possible?
- RQ2: What are the differences in read and write latency and throughput on the Hopsworks offline feature store, between the existing legacy system and IcedHops?
- RQ3: What are the differences in read and write latency and throughput on the Hopsworks offline feature store, between IcedHops and the delta-rs alternative?

It is fundamental to notice that, while the measured performance is expected to be comparable if PyIceberg and/or delta-rs are included into the Hopsworks client in the future, they are officially not functioning via the offline feature store **API** but only using the same file system, **HopsFS**.

*PyIceberg repository accessible at <https://github.com/apache/iceberg-python>

†Project repository available at <https://github.com/delta-io/delta-rs>

1.3 Purpose

This thesis project aims to diminish read and write latency (seconds) and thus enhance data throughput (rows/second) for operations on the Hopsworks offline feature store. This research will evaluate the performance of the existing legacy pipeline, which utilises Spark for writing, against IcedHops on a small-scale dataset by assessing variations in latency and throughput during read and write operations. If IcedHops is shown to be a more efficient option, Hopsworks AB will contemplate including this pipeline into their application. The same assessment will be conducted, within the same data scale domain, between PyIceberg and delta-rs pipelines.

The overall ramifications of this thesis are far larger, given Spark's prominence within the open-source community, which has had almost 3000 contributors during its existence [35]. Opting using IcedHops or delta-rs instead of Spark provides developers with a wider array of options for managing small-scale data (1 GB - 100 GB). This thesis work also provides an industrial use case, that can be used to understand which alternative to Spark could fit the best a specific developer need.

1.4 Goals

This project aims to reduce latency and hence enhance data throughput for reading and writing on Iceberg tables inside **HopsFS**. The achievement of this objective is tight to a specified list of **Goals (Gs)**, here described. These are connected to the collection of **RQs**, delineating a distinct framework entailing each project milestone.

1. **Gs** aimed to answer RQ1:

G1: Understand PyIceberg library tools, and identify which are needed to interact with **HopsFS**.

G2: Implement interaction between PyIceberg and **HopsFS**.

2. **Gs** aimed to answer RQ2:

G3: Design the experiments to evaluate the performance difference between the legacy access to Apache Hudi and the IcedHops PyIceberg-based access to Apache Iceberg, on **HopsFS**.

G4: Perform the designed experiments.

G5: Visualize the experiments' results, focusing on an effective comparison of performances.

G6: Examine and describe the findings in a dedicated Section of the thesis report.

3. **Gs** aimed to answer RQ3:

G7: Investigate related work on delta-rs library-based access to Delta Lake, on **HopsFS**.

G8: Visualize experiments' results of IcedHops and delta-rs alternatives to the legacy system, focusing on an effective comparison of performances.

G9: Examine and describe the findings in a dedicated Section of the thesis report.

To reach these **Gs**, several **Deliverables (Ds)** will be created:

D1: Experiment results on the difference in performance between the legacy access to Apache Hudi and the IcedHops PyIceberg-based access to Apache Iceberg, on **HopsFS**. This **D** is related to completing **Gs3–5**.

D2: Experiment results on the difference in performance between IcedHops PyIceberg-based access to Apache Iceberg and the delta-rs library-based access to Delta Lake, on **HopsFS**. This **D** is related to completing **G7** and **G8**.

D3: This thesis report. It offers additional detail on the implementation, design choices, performance, and evaluation of the results. This **D** is a comprehensive report of all thesis work, incorporating the analysis specified in **G6** and **G9**.

1.5 Ethics and Sustainability

This research project focusses on software, namely the development of more efficient data-intensive processing pipelines applicable in **ML** and **NN** training. The Green Software Foundation * states that software may be "part of the climate problem or part of the climate solution." Green software is described as software that minimises its environmental effect by using fewer physical

*Foundation's website available at <https://greensoftware.foundation/>



Figure 1.1: Illustrations of the **SDG** supported by this thesis.

resources and less energy, while optimising energy use to employ lower-carbon sources [36]. In the realm of **ML** and **NN** training, minimising training duration—and therefore the read and write latencies associated to the dataset—has been proved to significantly decrease carbon emissions [37, 38]. This thesis minime latency and thus enhance data throughput for reading and writing on Iceberg tables on **HopsFS**. This objective adheres to the essential principles of green software by reducing **CPU** utilisation time relative to the prior system. Minimising **CPU** utilisation time lowers energy consumption, resulting in a reduced carbon impact. This as a consequence, enhance energy efficiency of data-intensive computer pipelines, which are extensively used in **ML** and **NN** training.

For the same reason, this project supports two **Sustainable Development Goals (SDGs)** *, namely Affordable and Clean Energy (Goal 7), and Industry Innovation and Infrastructure (Goal 9). In particular, it tackles objectives 7.3, "Double the improvement in energy efficiency", and 9.4, "Upgrade all industries and infrastructures for sustainability".

Furthermore, the experiments created to satisfy **G2–3** has been carefully designed to minimize the number of trails necessary for statistical relevance, pursuing resource effiency. All the data used for this thesis are thoroughly explained, and all the scripts made available, to best address reproducibility of this project.

1.6 Research Methodology

This thesis is built from a few **Industrial Needs (INs)**, provided by Hopsworks AB, and a few **Project Assumptions (PAs)**, validated through a literature study. The Hopsworks's **INs** are:

***SDGs** website available at <https://sdgs.un.org/>

IN1 : The Hopsworks offline feature store, supported by the legacy pipeline, exhibits high latency (over one minute) and low throughput in write operations for small-scale data (from 1 GB to 100 GB). These performances indicate the potential of employing Spark alternatives in a small-scale data scenario.

IN2 : Hopsworks actively looks to customer needs and software's integration capabilities. Improving the read and write operations performance on their feature store and/or integrate their software with additional table formats, like Iceberg, is benefitted by all Hopsworks feature store users.

The **PAs** will be validated in Chapter 2. Those **PAs** are:

PA1 : Python is the most popular programming language and the most used in data science workflows. **ML** and **AI** developers prefer Python tools to work. Thus high-performance Python libraries will typically be preferred over **Java Virtual Machine (JVM)** or other alternatives.

PA2 : Spark has been proved to perform worse than other options when processing data in small-scale scenarios (from 1 GB to 100 GB). Alternatives to the Spark-based system could strongly improve reading and writing operations on the Hopsworks feature store, in different ways.

This thesis work fulfill the **INs** following a system implementation and evaluation guided by its **Gs**. Initially IcedHops, an integration between PyIceberg and **HopsFSs** will be implemented [34], followed by test to validate the approach. Then, an evaluation structure will be designed and used to compare the performances of the current legacy system, IcedHops and the delta-rs pipeline developed in the related work. These experiments will involve datasets of varying sizes and evaluate critical performance metrics such as read and write latency (in seconds), and throughput (in rows/second). Those two metrics were chosen as they most affect the computational time of accessing **OTFs**, thus being fair comparators between pipelines.

1.7 Research Limitation

The project is conducted in collaboration with Hopsworks AB, and as such the implementation will focus on working with their feature store and related system, supported by **HopsFS**. The considerations drawn in this thesis provide an insight into Spark limitations, and on which tools perform better in different data scales. However, those results cannot be generalized for any system.

1.8 Structure of the thesis

Chapter 2 provides to readers the basic information needed to understand the layered data stack analyzed in this study. Additionally, it presents the legacy and novel system architectures that will be used in the experiments. Chapter 3 delineates the methodologies for the system integration and the system evaluations. Chapter 4 illustrates the decisions made during the system selection phase and describes the design of the experiments. Chapter 5 presents the results of the experiments, focussing on the differences between the described pipelines, with various CPU settings. The chapter includes a discussion Section describing the major findings and implications of this thesis. Chapter 6 eventually summarises the contributions and discoveries of this thesis, highlights its limitations, and sets the discussion on prospective future research.

Chapter 2

Background and Related work

This project aims to improve performances of systems consisting of a layered data stack, able to handle large volumes of structured and unstructured data, at a high velocity, also defined as Big Data [39]. A generic data stack handles storing, management and retrieval of the data, offering access to those to the applications on top. However, there is no a single data stack for all cases, but different architectures used on different needs [40, 41, 9]. The following Sections of this report defines the data stack for this use case, which is illustrated in Figure 2.1.

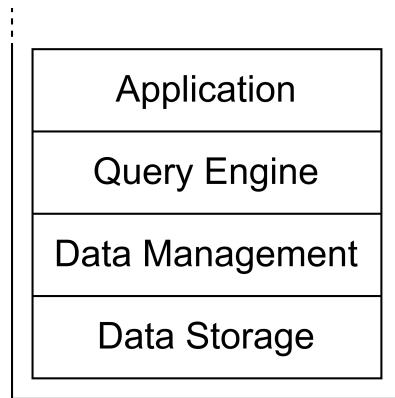


Figure 2.1: Data stack abstraction for this project.

The data stack and this chapter are divided into four Sections:

1. **Data Storage:** handles the physical storage of data. This layer determines how data is physically stored, including aspects like centralization or distribution, on-premise or cloud deployment, and storage formats such as files, objects, or blocks.

2. **Data Management:** handles the organization, governance, and lifecycle of data. The data management layer may provide features such as **ACID** properties, data versioning, support for open data formats, and the ability to store and manage both structured and unstructured data effectively.
3. **Query Engine:** handles the execution of data queries. This layer is responsible for efficiently accessing, retrieving, and writing data based on user requests. Key features of a query engine may include caching mechanisms, highly scalable architectures, and support for diverse programming languages through **APIs**.
4. **Application:** a system that utilizes the capabilities of the underlying data stack to achieve specific objectives. In this project, the focus will be on the Hopsworks feature store software.

Following the data stack section, the architectures of the legacy system, IcedHops, and the delta-rs system architectures are explained in Section 2.5, showing how the technologies are reflected within the pipelines measured during this thesis' experiments.

2.1 Data storage

This Section describes first what data storage is and which typologies of storage exist. Thus, this project's data storage layer, **HopsFS**, is presented, starting from its evolution from **Hadoop Distributed File System (HDFS)**, going into the details about its tools, and ending with possible alternatives, the cloud object storages.

2.1.1 Block storage vs. File storage vs. Object storage

Data can be stored and organized in physical storages, such as **Hard Disk Drives (HDDs)** or **SSDs**. The three main type of data storage are (1) Block storage, (2) File storage, and (3) Object storage briefly compared in Table 2.1. Each typology is describe in the following paragraphs, and their pros and cons are summarized in Table 2.2. This Subsection is a re-elaboration of three articles from major cloud providers (Amazon, Google, and IBM) [42, 43, 44] according to the author's understanding.

Table 2.1: Data storage features comparison. Table inspired by major cloud providers articles [42, 43, 44].

Characteristics	Block Storage	File Storage	Object Storage
Performance	High	High	Low
Scalability	Low	Low	High
Cost	High	High	Low

Block Storage

Block storage is a data storage method that divides data into discrete blocks of fixed size, each assigned a unique identifier. These blocks are stored independently on a storage system, such as a **Storage Area Network (SAN)** or within a cloud environment.

This decentralized approach offers several key advantages. Firstly, it enables high performance with fast read/write speeds and low latency, crucial for demanding applications like databases and virtual machine environments. Secondly, block storage provides direct, low-level access to storage volumes, similar to physical disks, granting users and applications granular control over data organization and management. This flexibility allows for a wide range of use cases, including powering virtual machine environments, supporting high-performance databases, and enabling efficient file sharing.

While offering significant benefits, block storage also presents certain limitations. It typically requires specialized hardware and infrastructure, potentially leading to higher costs compared to other storage options. Furthermore, while it offers a degree of scalability, expanding beyond certain limits can become complex and costly. Despite these considerations, block storage remains a vital technology for modern IT environments, enabling high performance, flexibility, and agility in data management.

File Storage

File storage is a hierarchical data organization method that stores data in files, which are organized into folders within a structure of directories and subdirectories. Files are characterized by extensions (e.g., ".txt", ".png", ".csv"), defining how the data is organized and accessed. This system simplifies locating and retrieving individual files when their exact paths are known, making it intuitive and user-friendly.

This structure is particularly beneficial for managing structured data and is widely used in **Personal Computer (PC)** and **Network Attached Storage (NAS)** devices. It enables centralized file sharing on **Local Area Network (LAN)** and supports common file-level protocols, ensuring compatibility across Windows and Linux systems. Storing data on a separate NAS device or in the cloud also enhances data protection and disaster recovery, with options to replicate data across multiple geographic locations for added security.

However, as the volume of files grows, scaling becomes challenging. Locating files in a large hierarchy can be time-consuming, and scaling often requires investing in additional or higher-capacity hardware. Cloud-based file storage services mitigate these challenges by offering scalable, off-site storage managed by service providers. These services eliminate hardware maintenance costs and provide flexible, subscription-based models that adapt to varying storage and performance needs.

File storage remains popular for applications requiring simplicity and centralized access, such as file sharing, personal storage, and cloud-based platforms like Dropbox and Google Drive. While other storage solutions may be better suited for managing massive datasets or unstructured data, file storage's accessibility, affordability, and ease of use ensure its ongoing relevance.

Object Storage

Object storage is a flat data storage method that organizes data into self-contained objects, each containing metadata that describes attributes like size, creation date, and unique identifiers. This metadata not only defines the data but also enables efficient querying and retrieval of large datasets. This makes object storage particularly well-suited for managing unstructured data, such as videos, images, and other media files that do not fit neatly into traditional hierarchical systems.

The flat structure of object storage eliminates complex hierarchies like folders and directories, simplifying organization and improving scalability. This structure allows object storage systems to replicate data across multiple regions, enhancing accessibility and fault tolerance in case of hardware failures. As a result, users benefit from faster data access in different parts of the world and robust disaster recovery options.

However, object storage has limitations. Objects are immutable, meaning they cannot be directly altered once created. Any changes require the creation of a new object. Additionally, object storage does not support transactional

operations, as it lacks mechanisms like file locking, making it unsuitable for applications requiring frequent updates or real-time data changes. It also has slower writing performance compared to file or block storage solutions.

Overall, object storage is an excellent choice for use cases requiring high scalability, such as social networks, video streaming platforms, and cloud-based services. Its flat structure and metadata-driven design are ideal for managing large, static datasets. However, other storage options are preferred when high performance is required for frequently changing files or when transactional consistency is critical.

Table 2.2: Data storage pros and cons comparison. Table inspired by major cloud providers articles [42, 43, 44].

Storage Typology	Pros	Cons
Block	High performance High reliability Easy updates	Lacks metadata Not easily searchable High cost
File	Easy on small-scale User-friendly User-manageable File-level locking	Inefficient on unstructured data Limited scalability
Object	Ideal on unstructured data Cost-effective Highly scalable Efficient advanced retrieval	No file-level locking Low performance No data updates

2.1.2 Hadoop Distributed File System

HDFS, a distributed file system *, is designed to store and process massive datasets efficiently. It leverages a cluster of commodity hardware, allowing for cost-effective scalability and high availability. Unlike traditional file systems, **HDFS** prioritizes high-throughput data access, making it well-suited for applications that process large volumes of data (higher than 100 GB), such as log analysis, data warehousing, and machine learning [45].

At the core of **HDFS** lies a master-slave architecture. A single Namenode acts as the central control point, managing the file system namespace, tracking

*HDFS official guide available at https://hadoop.apache.org/docs/r1.2.1/hdfs_user_guide.html

file metadata, and controlling client access to files. Multiple Datanodes serve as worker nodes, each responsible for storing and managing a portion of the data within the cluster. **HDFS** divides files into large blocks, which are then replicated across multiple Datanodes to ensure data redundancy and fault tolerance. This distributed storage approach enhances data availability and minimizes the risk of data loss due to hardware failures. Figure 2.2 presents a simplified visual representation of the Namenode read/write operations and Datanode orchestrating operations in **HDFS**.

The Namenode maintains a comprehensive record of the file system namespace, including file locations, block mappings, and access permissions. This metadata is crucial for efficient data retrieval and processing. By strategically placing data replicas across different nodes, **HDFS** minimizes data movement, optimizing performance and reducing network congestion. This aligns with the core principle that "moving computation is cheaper than moving data". **HDFS** is designed for applications that primarily require write-once-read-many access patterns. This design choice simplifies data management and optimizes performance for batch processing tasks, where data is typically written once and then read multiple times for analysis or processing.

Furthermore, **HDFS** relaxes some of the strict requirements defined by POSIX, such as low-latency access, to prioritize high-throughput data transfer. This allows **HDFS** to efficiently handle large-scale data processing tasks, making it a cornerstone of many big data applications.

2.1.3 HopsFS

HopsFS, introduced at the 15th USENIX Conference in 2017 [34], represents a significant evolution of **HDFS**. **HopsFS** addresses the critical scalability and metadata management limitations inherent to its predecessor, decentralizing metadata management by leveraging RonDB *, a distributed NewSQL database, which allows metadata to be stored and managed across multiple Namenodes. This architecture supports Namenode replication and dynamic scaling, significantly enhancing throughput and operational efficiency.

HopsFS encapsulates file system operations as distributed transactions, using advanced NewSQL features such as partition pruning and write-ahead caching. These techniques enable faster metadata retrieval and scalable operations, as seen in experiments where **HopsFS** outperformed **HDFS** by 16 to 37 times in throughput for real-world workloads [33]. Despite its scalability

*Details available at <https://www.rondb.com>

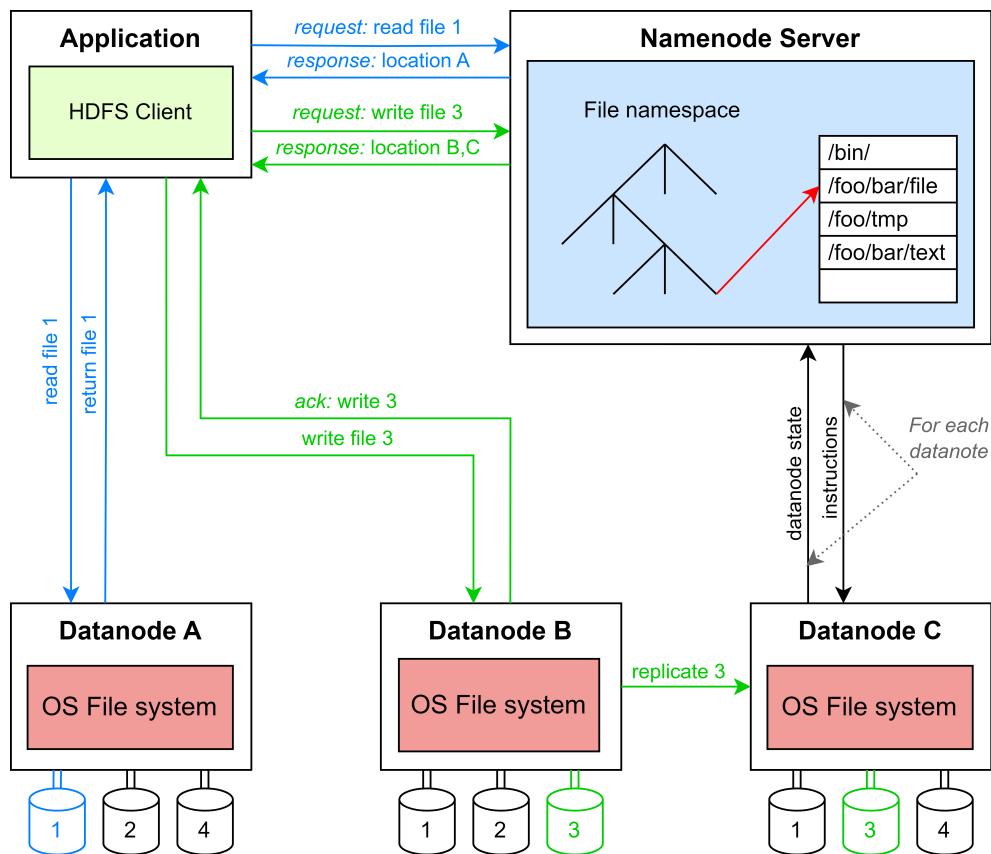


Figure 2.2: Hadoop Distributed File System (HDFS) architecture displaying in different colors basic operations: read (blue), write (green) and Namenode-Datanodes management messages (black). Note: for representation simplicity, files are not segmented into blocks and a single Namenode-Datanode message exchange is pictured. Diagram inspired by the Data-intensive Computing lectures at KTH by Prof. A. H. Payberah. Course website available at <https://www.kth.se/student/kurser/kurs/ID2221?l=en>.

and performance advantages, **HopsFS** remains backward-compatible with **HDFS**, serving as a drop-in replacement. While **HDFS** is widely adopted as a cornerstone of big data applications, **HopsFS** provides a forward-looking solution for environments where metadata scalability and performance are critical. Both systems demonstrate the strengths of distributed file systems in managing massive datasets, but **HopsFS** builds on the foundation of **HDFS** by addressing its limitations and pushing the boundaries of what distributed file systems can achieve.

2.1.4 Cloud object stores, an alternative

The advent of cloud computing has revolutionized data storage, with object storage emerging as a dominant paradigm [46]. Pioneered by [Amazon Web Services \(AWS\)](#) with its S3 service in 2006, cloud object storage services have rapidly proliferated, offered by major providers like [Google Cloud Storage \(GCS\)](#) and Microsoft Azure. This widespread adoption stems from the inherent advantages of cloud-based solutions, particularly their scalability and cost-efficiency, as explained in 2.1.1. Cloud object storage empowers users to dynamically scale their storage capacity on-demand, paying only for the resources consumed. This pay-as-you-go model eliminates the need for significant upfront investments in hardware and infrastructure, significantly reducing operational costs. Furthermore, cloud providers leverage economies of scale to unparalleled levels of availability and scalability that are difficult to achieve for most organizations.

[HDFS](#), initially released in 2006, evolved alongside the rise of cloud object storage. While [HDFS](#) has been widely adopted for on-premise deployments, the increasing popularity of cloud services has shifted the landscape. Cloud object storage solutions like [AWS S3](#), [GCS](#), and Azure Blob Storage have gained significant traction, with applications prioritizing support for these platforms due to their widespread adoption. While [HDFS](#) and its advancements, such as [HopsFS](#), still have their place in specific use cases, the convenience and scalability offered by cloud object storage services have made them the preferred choice for many organizations. However, in the latest year, hybrid solution have been created, to extend cloud object storage with typical file-based system advantages, such as [HopsFS-S3](#) [47].

2.2 Data management

This Section introduces the concept of data mananagement layer, starting from its origin and evolution in to today's technologies. This dissertation includes advantages and limitations of each evolutive step, and describe the functions of [DBMS](#). what a data storage is and which typologies of storages exist. The Subsections 2.2.2 focuses on the architecture of data lakehouses, explaing in details each component involed. Then, in Subsections 2.2.3, three data lakehouse frameworks are compared, namely Hudi, employed in the legacy version of Hopsworks feature store, Iceberg and Delta Lake, the two alternatives evaluated.

2.2.1 Brief history of Data Base Management Systems

Since the 2010s, with the advent of Big Data, the data volume, variety, and production velocity have increased exponentially [28, 48]. While on one side, this proved to be of enormous value, on the other this posed several challenges [49] on data architectures, which had to evolve to cope with these new needs. Data lakehouse frameworks, like Hudi, Iceberg and Delta Lake [8, 9, 10] are the last step of this evolution. However, to truly understand these tools, it is needed to start from the beginning of DBMS evolution.

Before big data, companies already wanted to get insights from their data, automating the workflow from the data sources to point of access to this data. Here is where ETL and relational databases first came into use. An ETL pipeline consists of three steps:

1. **Extracts** data from APIs or other various data sources.
2. **Transforms** data according to one or more goal. Often this means removing absent fields, standardizes to a specific format to match the database schema, and validates the data.
3. **Loads** it into a relational database (e.g., MySQL).

This workflow structure enabled companies to generate BI insights and data reports based on organizational data. However, a key limitation of this approach is its restricted ability to perform analytical queries that require joining multiple tables and working on several data dimensions. These types of queries, while executed less frequently than simpler queries, are essential for strategic decision-making (e.g., identifying the customer segment with the highest profitability over the past year, with the capabilities of drilling down on products, marketing and sales information).

To address the increasing demand for analytical queries, more advanced DBMs replaced traditional relational databases, optimizing performance for business-oriented analytical workloads. These systems, known as OLAPs, introduced specialized storage and query execution strategies tailored for large-scale analytical processing. The most notable example of an OLAP system is the data warehouse, which revolutionized how organizations handle structured data analysis. A data warehouse workflow, enables companies to process and analyze significantly larger datasets than traditional databases. Unlike transactional databases, which optimize for rapid insert/update operations, data warehouses focus on read-optimized queries by structuring data into columnar formats, reducing scan times for analytical queries. They

maintain core relational database properties such as **ACID** transactions and data versioning, ensuring data integrity and consistency for complex business intelligence applications.

Over time, the exponential growth of unstructured data (e.g., images, videos, logs, and sensor data), posed new challenges for companies aiming to leverage such information. Traditional data warehouses struggled to accommodate this data due to their rigid schema requirements and high storage costs. Moreover, they were not designed to support **AI/ML** workflows that rely on diverse, unstructured datasets. To overcome these limitations, organizations adopted a new paradigm known as data lakes. Data lakes leverage cost-efficient object storage (Section 2.1.1) and a schema-on-read approach, where raw data is loaded first and transformed only when needed. This shift from the traditional **ETL** model to **ELT** offers greater flexibility in handling diverse data types. For example, businesses can use data lakes to store raw IoT sensor readings and later refine them for predictive maintenance models, once the needed data transformation will be designed. Despite reducing storage costs and improving flexibility, data lakes introduced new complexities. Unlike structured databases, querying data lakes directly for **BI** reports is impractical, as they lack indexing and transactional consistency. Additionally, maintaining separate storage solutions for structured (data warehouses) and unstructured (data lakes) data increases operational complexity and costs. A major drawback of data lakes is the risk of becoming "data swamps", repositories filled with ungoverned, low-quality data that provide little business value. Recognizing these challenges, organizations sought a hybrid solution combining the best features of data warehouses and data lakes. This led to the emergence of the data lakehouse architecture. First described by Databricks in 2020 at [Conference on Innovative Data Systems Research \(CIDR\)](#) [3], data lakehouses integrate structured and unstructured data storage capabilities of data lakes, while maintaining the governance, **ACID** compliance, and indexing capabilities of data warehouses. The data lakehouse approach resolves many pain points associated with separate data warehouses and data lakes. It enables organizations to use a single storage for all types of data, while supporting high-performance **SQL** queries and **ML** workloads. Furthermore, by leveraging open file formats like Apache Parquet, ORC and Avro, data lakehouses ensure interoperability with various analytics engines, avoiding the risk of getting data locked into a proprietary format [50]. All these capabilities and features make lakehouses an ideal solution for enterprise-scale data processing. Table A.1 provides a comparison of data warehouses, data

lakes, and data lakehouses, highlighting key takeaways for each technology and summarizing what described so far.

2.2.2 Data lakehouse architecture

Data lakehouse architectures combines the desirable attributes of data warehouses and data lakes, mitigating the challenges encountered by both these technologies, and eliminate the need of a two-tier data stacks to run varying analytical workloads. The key components of data lakehouses are at the bases of any other architecture seen so far, but their design focus on reducing components coupling, providing more agility and choice when architecting such platform [50]. The data lakehouse components, graphically presented in Figure 2.3, are:

- **Data storage:** is where data files land after ingestion from various systems, usually a Cloud object store, as described in Section 2.1.
- **File formats:** hold the actual raw data and are physically stored in the data storage. They are open-source file formats, such Apache Parquet or JSON, and are typically column-oriented.
- **Table format:** also referred as OTF, acts as metadata layer on the top of the file formats, that abstracts the underlying physical data structure. They offer API access to the query and storage engines.
- **Storage engine:** handles data management tasks, aimed at optimizing efficiency of queries over the data. It performs tasks as data compaction, indexing, and partitioning.
- **Catalog:** sometimes described as metastore, it enables efficient search and discovery. The catalog keeps track of information about each table (name, column names, data types) and a reference to metadata for each table (table format).
- **Query engine:** is responsible for processing data, performing read and write operations leveraging the table format API. They are further explained in Section 2.3.

OTF are open standards, and by design support interoperability throughout the stack, as visible in Figure 2.3. Thus, depending on the integration capabilities of the table format, there are several implementation options for data storage, file formats, storage engine, catalog and query engine.

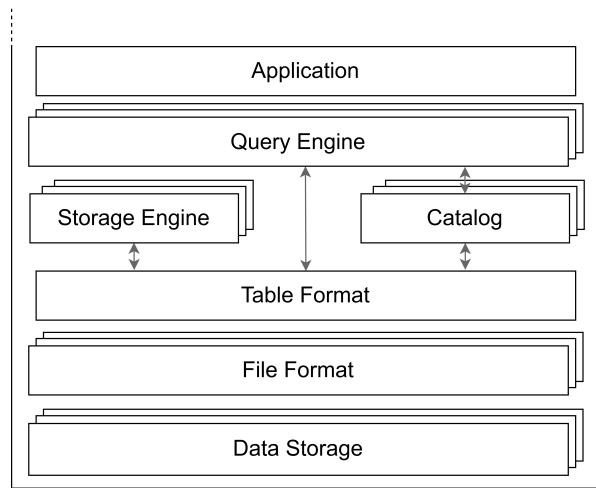


Figure 2.3: Abstraction of architecture of a data lakehouse. Inspired by "Open table formats in perspective" article on OneHouse Website, available at <https://www.onehouse.ai/blog/open-table-formats-and-the-open-data-lakehouse-in-perspective>

Additionally, recent introduction of tools like Apache XTable * demonstrates the trend towards a universal compatibility between OTFs.

Open Table Formats

OTFs provide an abstraction layer on top of data lakes, enabling database-like functionalities, enhancing data management capabilities significantly [3, 51]. One of the cornerstone features of OTFs is support for full **Create Read Update Delete (CRUD)** operations. The ability to perform updates and deletes sets data lake house data storage apart from traditional file-based storages, where such operations are cumbersome and inefficient. Performance and scalability are other notable features that OTFs bring to the table. These formats are designed to excel in Big Data environments, where data volumes are massive and continue to grow. OTFs could support various optimization techniques, such as indexing, partitioning, and caching, to expedite data retrieval and processing. This not only improves query performance but also ensures that the system can scale horizontally to accommodate increasing data loads without a significant degradation in performance. As a result, organizations can manage their data ecosystems more effectively, making data-driven insights more accessible and actionable. Transactional support

*XTable repository available at <https://xtable.apache.org/>

with **ACID** compliance is another key feature of **OTFs**. This ensures that all data transactions are processed reliably, maintaining data integrity and consistency across the board. This is particularly important in scenarios where multiple transactions occur simultaneously or when the system needs to recover from partial failures. **OTFs** guarantee that each transaction is completed successfully or fully rolled back, providing an essential level of data reliability and trustworthiness for critical business operations. This comprehensive functionality allows for flexible and complex data workflows, and ensures that data lakes and warehouses can be updated in real time, reflecting the most current state of information.

2.2.3 Data lakehouse comparison

The three main data lakehouse frameworks investigated in this project are Hudi [8], Iceberg [9] and Delta Lake [10]. Their popularity has increased proportionally with the popularity of the data lakehouse architecture, as visible by the growing community of each of these technologies in Figure 2.4, becoming the de-facto standards for data lakehouse implementations [52]. Some alternatives are newly being developed, such Apache Paimon *, but are still at early stages of their development and will not be investigated in this thesis project. The following subsections present each of these technologies, specifying their development history, key features, integration capabilities, **API** offered and in which use cases they are most suitable for.

Apache Hudi

Apache Hudi, open-sourced by Uber in 2017 [8], is an open-source framework that addresses the challenges of low-latency data ingestion and incremental processing. Hudi enables efficient record updates and deletions in data lakes, thus eliminating the need to rewrite entire datasets [53]. The framework supports **Change Data Capture (CDC)** for efficient updates and deletes. Hudi offers two primary write optimization modes: **Copy on Write (CoW)** for high read performance and **Merge on Read (MoR)** for balancing both write and read performance.

Hudi's metadata layer is structured around a timeline-based architecture, maintaining a history of all table operations, thus providing version control. The commit timeline records actions such as inserts, updates, deletes, and compactions, enabling time travel by referencing specific commit instants.

*GitHub repository available at <https://github.com/apache/paimon/>.

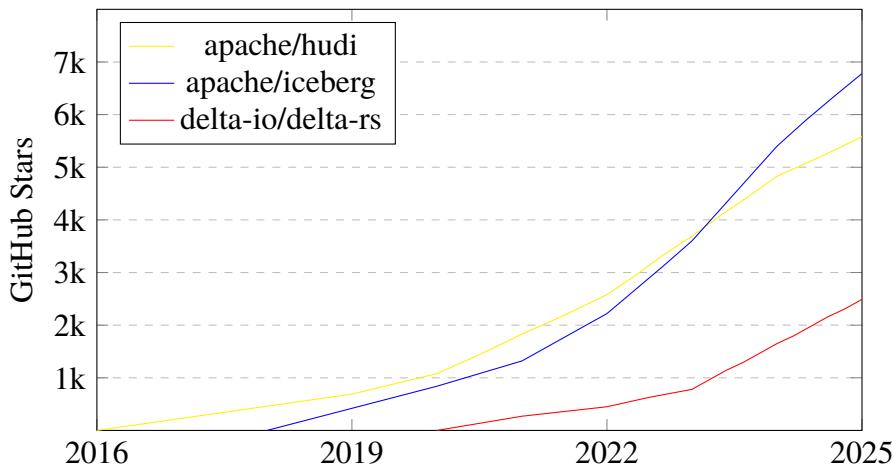


Figure 2.4: Trend of GitHub stars of following repositories <https://github.com/apache/iceberg>, <https://github.com/apache/hudi>, <https://github.com/delta-io/delta-rs>.

The metadata table optimizes file listing by indexing data files, partitions, and record locations, improving query performance. Hudi also maintains delta logs (**MoR**) to track incremental changes before they are compacted into columnar storage. The file groups and file slices structure organizes base and log files, ensuring efficient data versioning.

Hudi is optimized for real-time analytics through low-latency streaming ingestion, and integrates seamlessly with Spark and Flink for data ingestion and processing. Hudi also supports reading data from Hive, Impala, and Presto. The framework incorporates multimodal indexing, such as Bloom filters and record-level indexing, and employs both Multi-Version Concurrency Control (**CC**) and Optimistic **CC** for transaction management. Typical use cases for Hudi include streaming ingestion with frequent updates and deletes, incremental data processing, and real-time updates in domains like IoT, fintech, and log data processing [54]. While Hudi excels in real-time ingestion, its metadata overhead can be significant for large-scale analytical queries, potentially resulting in slower query performance compared to other formats like Iceberg and Delta Lake [55].

Hudi's core implementation is in Java, and it integrates deeply with Apache Spark, offering a Spark datasource **API** for reading and writing Hudi tables. This makes Java and Scala the primary languages for Hudi development and usage. While a standalone Python **Software Development Kit (SDK)** for Hudi doesn't exist, Python users can still interact with Hudi tables through PySpark,

leveraging the Spark API. Although community-driven efforts like hudi-rs ^{*}, a native Rust implementation with Python bindings, are emerging, offering read support for some technologies, Hudi is currently still JVM-centric.

Apache Iceberg

Apache Iceberg, open-sourced by Netflix in 2018 [9], is an open-source framework that addresses the limitations of Hive tables by emphasizing scalability and correctness for large-scale analytics. By separating metadata from data, Iceberg supports efficient snapshot-based isolation and query planning [11, 56]. The framework provides full ACID transactions, ensuring atomicity, consistency, isolation, and durability. Iceberg allows for schema evolution, enabling the addition, removal, and renaming of columns, without breaking existing queries. It also supports partition evolution, allowing changes to partitioning schemes without necessitating data rewrites.

Iceberg’s metadata layer consists of multiple components that enable efficient time travel and point-in-time query planning. The manifest list acts as an index, pointing to multiple manifest files, each of which tracks a subset of data files. These manifest files store information about data file locations, partitioning, and statistics (e.g., min/max values). The snapshot metadata records changes over time, referencing manifest lists and enabling time travel by tracking table versions. Iceberg also maintains a table metadata file, which holds high-level properties, schema, partitioning details, and references to the latest snapshot.

Iceberg is designed to scale to petabyte-scale datasets with optimized metadata management using manifest files. Its architecture enables data warehouse-like functionality, leveraging cloud object storage for efficient data access. Furthermore, Iceberg supports collaboration across multiple applications with transactionally consistent access. The framework integrates with query engines like Spark, Trino, Presto, and Dremio and supports Flink for both reading and writing. Iceberg excels in read performance, particularly for tables with large numbers of partitions [54]. It is ideal for analytical workloads that involve large-scale datasets and frequent schema or partition evolution. However, Iceberg’s write performance may not be as optimized for streaming scenarios compared to other formats.

Iceberg’s core is written in Java, and it is heavily integrated with Apache Spark [12, 56]. The most common way to interact with Iceberg is through SparkSQL or other query engines like Trino and Presto, which support

^{*}hudi-rs repository available at <https://github.com/apache/hudi-rs>.

Iceberg natively, making Java and Scala the dominant languages in these environments. However, Iceberg offers more than just Spark bindings. PyIceberg [57] provides a direct Python interface for read and – from release 0.6.0 in 2024 – write operations, making Iceberg accessible to Python developers without requiring PySpark. Furthermore, query engines like Trino, accessible from Python via libraries like PyHive, allow Python users to query Iceberg tables through **SQL**. Iceberg community has recently developed a Rust implementation, `iceberg-rust`^{*}, which provides read access to Iceberg tables, further expanding its accessibility to a wider range of language ecosystems.

Delta Lake

Delta Lake, open-sourced by Databricks in 2019 [10], is an open-source framework that enhances the reliability and performance of data lakes, facilitating a seamless transition between batch and streaming use cases. Often considered the first data lakehouse, Delta Lake employs a transaction log to record all changes to data, ensuring consistent views and write isolation, which supports concurrent data operations [58]. The framework offers **ACID** transactions, as well as features such as unified batch and streaming processing, indexing, and schema enforcement.

Delta Lake’s metadata layer is built around the Delta Log, which records every transaction in JSON files, enabling time travel. The transaction log maintains a sequential history of operations, while checkpoint files (stored in Parquet format) periodically summarize log entries for faster metadata access. The protocol metadata defines the table’s versioning, schema, and supported features, ensuring compatibility across different readers and writers. It also incorporates metadata-informed data skipping during merge operations and supports streaming through change data feeds.

Delta Lake tightly integrates with Spark, thus it is particularly well-suited for real-time streaming, batch processing, and **ML** pipelines requiring data versioning. While its open-source adoption continues to grow, Delta Lake remains deeply integrated with the Databricks ecosystem [52, 59]. Additionally, its focus on schema and partition evolution is less pronounced compared to frameworks like Iceberg, making Delta Lake less suitable for situations where data schemas are frequently changing.

Delta Lake core is written in Java and Scala [58], providing native support within the Spark ecosystem through the delta package. This allows Java and Scala developers to interact directly with Delta tables. Python support is

^{*}iceberg-rust available at <https://github.com/apache/iceberg-rust>.

provided via PySpark, enabling Python developers to leverage the Spark **API** for Delta Lake operations. However, Delta Lake's reach extends beyond the **JVM**. With the release of Delta Kernel [60], a Java library providing low-level access, and the development of delta-rs ^{*}, a Rust-based implementation, Delta Lake has become accessible to a wider audience. The library delta-rs allows interaction with Delta tables without Spark or **JVM** dependencies, and its Python bindings makes this particularly attractive to the Python data science community.

In Table A.2 is presented a comparation between key features of the three data lakehouse frameworks.

2.3 Query engine

This Section describes the technologies used to query, cache, and process data in this project. The category of query engine includes mainly Apache Spark and DuckDB, but technologies operating at the same abstraction level of those are here presented, namely Apache Kafka, Arrow Flight and Catalogs.

2.3.1 Apache Spark

Apache Spark is an open-source distributed computing framework designed for large-scale data processing [15]. It builds upon MapReduce, a distributed programming model developed by Google for handling massive datasets [61], later adapted into Hadoop MapReduce by Yahoo! engineers [62]. Spark enhances this model using **RDDs**[63], a disitrbutued memory abstraction that enables lazy in-memory computation, diffently by on-disk MapReduce's computation, that is tracked through lineage graphs. This increases fault tolernace [63] and allow to manage even bigger scale computations.

Spark supports various workloads beyond batch processing, leveraging an in-memory execution model for efficiency. Its core components include SparkSQL for querying structured data, Spark Streaming for real-time processing, MLlib for scalable **ML**, and GraphX for large-scale graph processing. With support for iterative computations, Spark is particularly well-suited for **ML** and graph analytics. It also provides **APIs** for Scala, Java, Python (via PySpark), and R, making it accessible to a broad range of users. Despite its advantages, Spark has limitations. Its in-memory execution requires significant RAM, and **JVM**-based execution can lead to

^{*}delta-rs repository available at <https://github.com/delta-io/delta-rs>.

performance overhead due to garbage collection. Tuning performance involves fine-tuning configurations, which can be complex. Additionally, Spark Streaming's micro-batch processing introduces higher latency compared to other frameworks, making it not the best solution for small-scale datasets processing [20].

2.3.2 Apache Kafka

Apache Kafka is a robust, open-source distributed platform for handling data streaming, that has become a cornerstone of modern data architectures [64]. Kafka supports high throughput data ingestion and processing, making it exceptionally well-suited for handling massive volumes of data in real, given its ability to manage and process continuous streams of data with very low latency. Kafka's architecture, which emphasizes fault tolerance, scalability, and durability, allows it to handle the demands of mission-critical systems that require continuous data flow and processing. The key components of Kafka architecture are:

1. **Producer:** an application that publishes data, labeling into specific topic that group similar messages.
2. **ZooKeeper:** the responsible for managing the Kafka cluster, including broker(s) information and topic message tracking, tracked with an offset for each topic.
3. **Broker:** a processing node, or server, that handles data for specific topics. It receives messages from producers and delivers them to consumers upon request, enabling asynchronous communication. When for a single topic there are multiple brokers, one is elected as a leader, and the others replicate its content.
4. **Consumer:** an application that subscribes to specific topics to receive and process data. Multiple consumers can subscribe to the same topic.

Kafka enables applications to behave as producers and consumers, without the need of developing any synchronization protocol. This enables producers to reach high throughput, as they can broadcast messages without waiting for any acknowledgements or availability signal from the consumers. Due to its distributed architecture, a Kafka ecosystem can be optimised according to specific needs, allowing several brokers, producers and consumers to coexist.

2.3.3 Catalogs

The term catalog is used in the data domain in multiple contexts, in each of those with a different definitions. In the data lakehouse context, thus in this project context, a catalog is a technical catalog also called metastore [11]. As explained in Section 2.2.2, a catalog plays an important role in tracking tables and their metadata. It is, at a minimum, the source of truth for a table's current metadata location. This is in contrast to a federated catalog, which is a central portal for data discovery and collaboration, since it tracks datasets across multiple data stores. It focuses on business needs such as data governance and documentation, and sets standards for authentication and authorisation [65]. A federated catalog may point to tables from Cassandra, Postgres, Hive, and other systems.

Made clarity over the catalog definition regarding this project, there are several technologies providing such tool. When data lakehouses were first created [8, 11], the **Hive Metastore (HMS)** was the most popular catalog. **HMS** is the technical catalog for Hive tables, but it can actually belong to both catalog categories, since it may also track **Java DataBase Connectivity (JDBC)** tables, and other datasets. However, scale challenges led developers and big tech companies to develop their own more scalable tools, which contributed to the development of data lakehouses technology itself, such as the case of Iceberg [12].

Data lakehouse technologies support different Catalogs, as shown in Table ???. Depending on the specific business need or on the available cloud resources and ecosystem, possible "pluggable" catalogs are:

- **SQL Catalog:** is a C library, with Python bindings (SQLite), that provides lightweight disk-base database, that allows accessing the database with nonstandard variants of **SQL**, where metadata can be saved. This does not require a separate server process, thus it is great in embedded applications, but it is not suitable for large-scale applications.
- **AWS Glue:** provides a Data Catalog that serves as a managed metastore within the **AWS** ecosystem. It integrates well with other **AWS** services, thus it is a common choice for data lakehouses built on **AWS**.
- **DynamoDB:** it is primarily a No **SQL** databases that due to its scalability and performance, can be used as a metastore. This is less common compared to other options.
- **JDBC:** per sé is not a catalog, but it's the standard Java **API** for

connecting to relational databases. In this context, **JDBC** would be used to connect to and interact with database systems that stores metadata.

- **Nessie**: acts as a versioned metastore, providing Git-like capabilities for managing data lake tables
- **Databricks Unity**: is Databricks' unified governance solution, thus it is a federated catalog, which however includes also technical catalog capabilities. It provides a centralized metastore specifically designed for the Databricks Lakehouse Platform, thus is the most common option for Delta Lake-backed lakehouses [60].

As data lakehouse frameworks grew to support more and more languages and engines, pluggable catalogs started causing some practical problems, related to compatibility. It proved indeed difficult, for commercial offerings, to support many different clients and catalog features. To overcome these problems, some developers, like Iceberg's community [65], created also a REST catalog protocol, a common **API** for interacting with any catalog. The advantages of such protocol are multiple: (1) one client implementation, for new languages or engines, can support any catalog; (2) secure table sharing is enabled, using credential vending or remote signing; (3) the amount of failures is reduced, since server-side deconfliction and retries are supported.

When implementing a new data lakehouse, the catalog decision depends on the specific integrations and feature proper of each catalog option, also considering how this can be configured on the top of the previously added layer, as shown in Table 2.3.

2.3.4 Duck DB

DuckDB [18] is an open-source, embeddable, **OLAP DBMS** designed for efficient processing of small-scale datasets (from 1GB to 100GB) within the same process as the application using it. This embedded, in-process operation, inspired by SQLite's success, simplifies deployment and eliminates the overhead of inter-process communication, leading to high responsiveness and low latency. DuckDB requires no external dependencies and compiles into a single amalgamation file, enhancing portability across major operating systems and architectures, including web browsers via DuckDB-Wasm. It offers **API** for various languages, such Java, C, C++, Go, Rust and Python and supports complex **SQL** queries, including window functions and **ACID** properties through Multi-Version **CC**.

Data is stored in persistent, single-file databases, and secondary indexes enhance query performance. DuckDB's columnar-vectorized query execution engine, a key design choice for **OLAP** workloads, processes data in batches (vectors), significantly improving performance compared to traditional row-by-row processing systems. While optimized for analytical queries processing significant portions of datasets, DuckDB is not designed for massive data volumes (1TB or more) that require disk-based processing, as its core processing relies on in-memory operations. However, its extensible architecture allows for adding features like support for Parquet, JSON, and cloud storage protocols as extensions.

2.3.5 Arrow Flight

Arrow Flight is an high-performance framework designed for efficient data transfer over networks, mostly utilising Arrow tables [66]. This protocol facilitates the transmission of large data volumes stored in a specific format, such as Arrow tables, without the necessity of serialisation or deserialisation for transfer. This significantly accelerates the data transfer, making Arrow Flight highly efficient. Arrow Flight is engineered for cross-platform compatibility and supports many programming languages, including C++, Python, and Java. The protocol further facilitates parallelism, enhancing transmission speeds by using numerous nodes in parallel networks. The Arrow Flight protocol is constructed upon gRPC, facilitating standardisation and simplifying the construction of connectors.

2.4 Application - Hopsworks

This Section describes the application layer, the uppermost layer presented in Figure 2.1, which takes advantage of the data stack described above. The software described is the Hopsworks feature store, which this project contributes to. This software is part of the broader **Machine Learning Operations (MLOps)** platform offered by Hopsworks AB, the company that hosted this master thesis.

2.4.1 Machine Learning Operations (MLOps)

MLOps are a full set of practices related to the development and automation of **ML** workflows. Through the **MLOps** lenses, a **ML** workflow is logically separated in smaller steps, and considered from a data, code and model

perspective. Differently from a classical software application, where only the code needs versioning, in the context of **ML** applications, it is key that all three of data, code and model are versioned. Perhaps, different data might be used by models in different moments, as well as new model might be trained on the same data, to truly understand their improvements. The novel challenges are thus related to data validation, **ML** artifacts versioning, **ML** workflow orchestration, and infrastructure management, given higher computation and storage capabilities needed by those workflows [4, 39].

The need of the described features saw new solutions emerging, like the Hopsworks AI data platform [67]. In the context of data versioning, the specific solution is called feature store [68], which consists in a centralized fast-access storage for both real-time and batch data.

A simple architecture following **MLOps** principle is presented in Figure 2.5. As first step, data are gathered either streaming (real-time) or batch data sources. A feature pipeline process those the data, performing model-independent transformations [69], and saves them in the feature store. A training pipeline runs now model-dependent transformations on features and labels retrieved from the feature store, and train the respective model, saving the trained model in the model registry. The last pipeline, the inference pipeline, extract at its need a specific model from the model registry and access the features, over which a inference will be conducted, from the feature store. The output of this pipeline, which normally is embedded directly in the consuming application, are the predictions of the model over the selected features, altogether with logs describing the performance of the model according to a pre-selected measure. All the pipeline are decoupled and could thus work asynchronously, making the whole **ML** workflow scalable, maintainable and effective.

2.4.2 Hopsworks AI Data Platform

As briefly said above, the feature Store is a key data layer in an **ML** application built over the **MLOps** principles. The feature store enables feature reusability, as well as multi-user collaboration and a centralized authenticated access to the feature. The Hopsworks feature store organizes features in feature groups. Those are mutable collection of features, over which developers can perform **CRUD** operations, accessing them via the Hopsworks **APIs**.

The Hopsworks feature store supports both batch data sources and real-time data streaming. This hybrid system is possible thanks to the dual implementation of an offline and an offline feature store. The offline feature

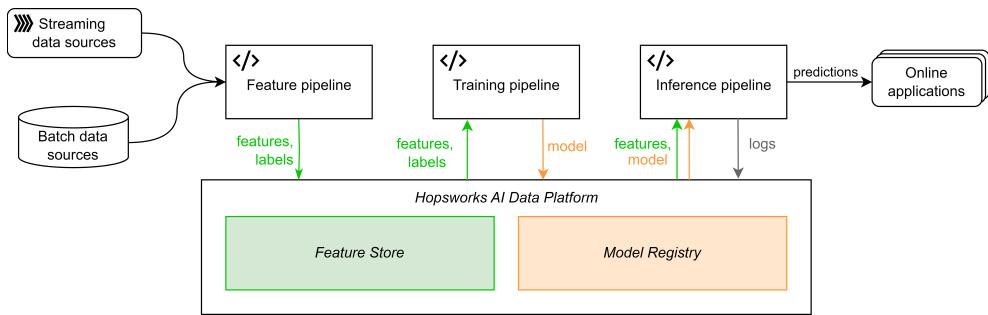


Figure 2.5: **MLOps** pipeline using a feature store and a model registry. Diagram inspired by Hopsworks documentation available at <https://www.hopsworks.ai/dictionary/feature-store>.

store is a column-based storage suited for batch data that is updated with a low frequency (every few hours at maximum frequency). The online feature store, on the other hand, is a real-time row-based, key-value data storage based on RonDB, thus enabling low latency and real-time processing. To maintain consistency between those two stacks, the Hopsworks feature store has a unique point of entry for data, which is Kafka, explained in detail in [2.3.2](#). This guarantees that each message is delivered to both storages, which acts as consumers subscribed to the Kafka topics.

The model registry works under the same principles of the feature store, providing a centralized access and a standard deployment for all the models saved in it. It also saves information about model performance along time, model schemas and training feature logs, to ease model reproducibility and tracking.

2.5 System architectures

This Section describes the architectures of the legacy system, based on Hudi, of IcedHops, developed during this thesis work and based on Iceberg, and the system implemented in a related thesis work [\[26\]](#), based on Delta Lake. Those are the systems that will be run and measured in the experimental part of this thesis work. This Section is divided into six Subsections, according to the system and the operation run over it. For each Subsection, a chart presents the operation protocol step by step.

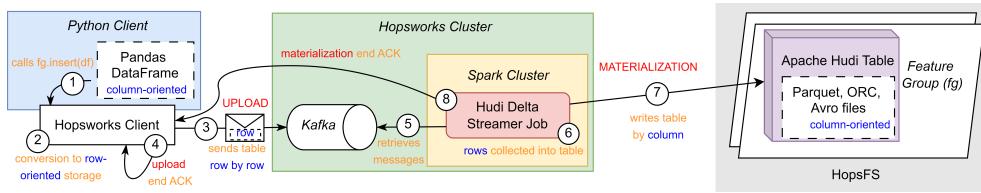


Figure 2.6: Legacy Hudi-based system writing a DataFrame (Pandas) to the Hopworks offline feature store via a Python client.

Each stage is assigned a numerical identifier. The transformation of the table format is highlighted in blue, illustrating the conversion from columns to rows and subsequently back to columns. The upload phase encompasses steps one through four, whereas the materialization phase concludes at step eight. The diagram was created based on individual interviews with Hopworks AB developers working on the Hopworks online and offline feature store and the related work [26] has been used as inspiration.

2.5.1 Legacy system - Hudi - writing

Figure 2.6 * illustrates how data is written from the client to the offline feature store in the legacy Hopworks system. This process is primarily divided into two synchronous phases: upload and materialization. During the upload phase, the input Pandas DataFrame is transformed into individual rows and transmitted to Kafka sequentially, one row at a time. Once the upload concludes, the client receives a notification. Meanwhile, a Spark job—Hudi Delta Streamer—operates asynchronously within the cluster from the moment the Hopworks cluster is initialized. This job continuously fetches messages from Kafka and, upon gathering an entire table, proceeds to write the data on the table in an Apache Hudi table sitting on top of a HopsFS system, using a columnar format. Once the materialization is completed, The Python client will receive a notification upon completion.

Similar to the pipeline, the upload and materialization phases operate asynchronously. During the experimental phase of the thesis, the materialization function was invoked to assess the overall process latency without factoring in the Hudi Delta Streamer's data retrieval interval. This approach allowed the system to trigger materialization on demand rather than waiting for the scheduled period. As a result, the experiments were able to obtain precise measurements of the total process latency.

*For enhanced visualization, refer Figure B.1.

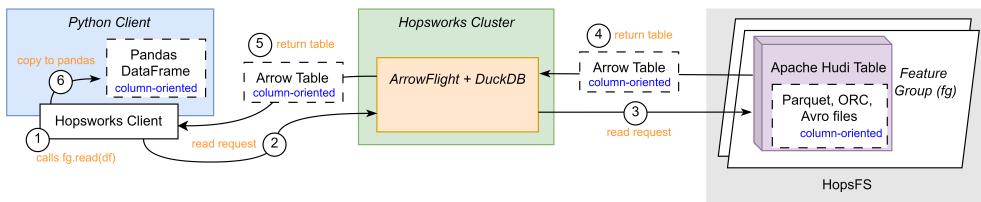


Figure 2.7: Legacy system reading a table from the Hopsworks offline feature store and loading it into the Python client’s local memory. The process is streamlined using Arrow Tables that avoid table conversion. Diagram inspired by the Hopsworks feature store paper [33] and by the related work [26].

2.5.2 Legacy system - Hudi - reading

Figure 2.7 * shows the offline feature store. Conversely to the legacy system writing process, the reading process uses already a Spark alternative, which consists of a combination of an Arrow Flight server and a DuckDB instance. This was implemented by Hopsworks AI in order to directly send the data in columnar format, without the need of converting them into row-based format.

2.5.3 IcedHops - writing

Figure 2.8 shows IcedHops writes on an Iceberg table instanced on top of HopsFS. If first requests metadata information about the existing table to Iceberg Catalog, implemented using SQLite in the example. Once the JSON file containing the metadata is received, it reads the table location and request to read the table. IcedHops streamlines the process without passing from a server instance (Spark), removing the middle-tier from the process.

2.5.4 IcedHops - reading

Figure 2.9 shows how IcedHops reads on an Iceberg table instanced on top of HopsFS. If first requests metadata information about the existing table to Iceberg Catalog, implemented using SQLite in the example. Once the JSON file containing the metadata is received, it reads the table location and proceed to write (append), by column, on the Iceberg Table. IcedHops streamlines the process without passing from a server instance (Arrow Flight), removing the middle-tier from the process.

*For enhanced visualization, refer Figure B.2.

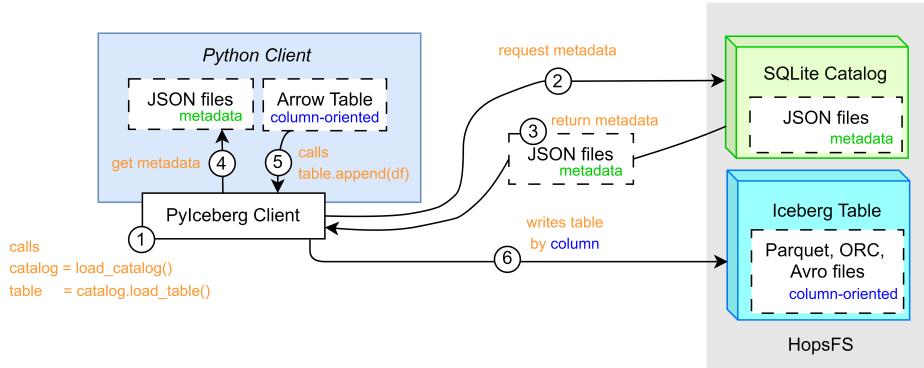


Figure 2.8: IcedHops writing an Arrow Table from a Python client to an Iceberg Table stored on [HopsFS](#).

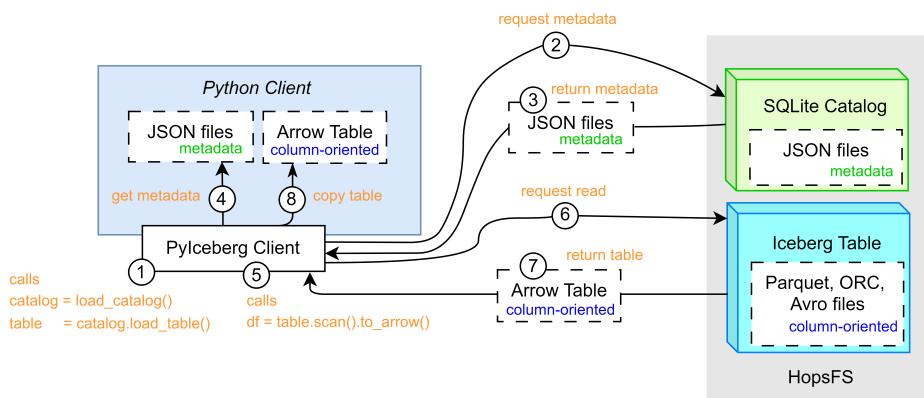


Figure 2.9: IcedHops reading an Iceberg Table stored on HopsFS and loading it into memory.

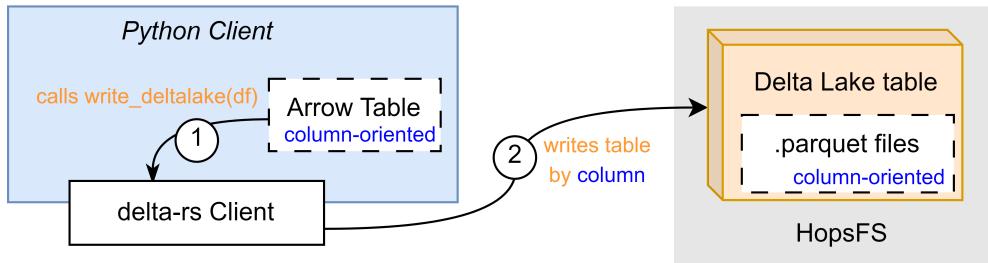


Figure 2.10: Delta-rs library writing an Arrow Table from a Python client to a Delta Lake table stored on [HopsFS](#). Diagram inspired by the related work [26].

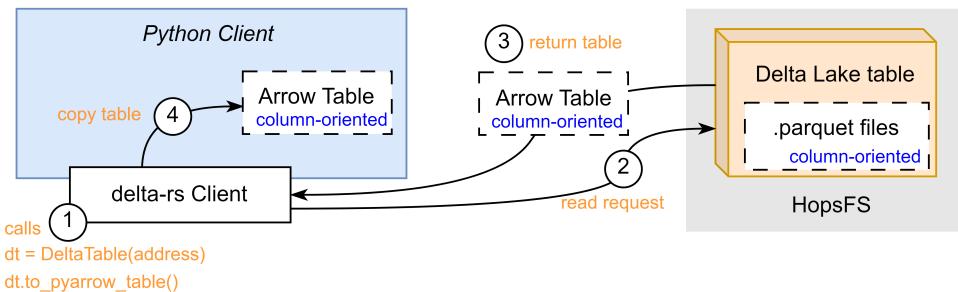


Figure 2.11: Delta-rs library reading a Delta Lake table stored in [HopsFS](#) and loading it into memory. Diagram inspired by the related work [26].

2.5.5 delta-rs - writing

Figure 2.10 shows how the delta-rs library writes on a Delta Lake table instanced on top of [HopsFS](#). The delta-rs library streamlines the process without passing from a server instance (Spark), removing the middle-tier from the process. In this system, the only file format supported is Parquet, while the other systems supported also ORC and Avro.

2.5.6 delta-rs - reading

Figure 2.11 shows how the delta-rs library reads a Delta Lake table instanced on top of [HopsFS](#). The delta-rs library streamlines the process without passing from a server instance (Arrow Flight), removing the middle-tier from the process. In this system, the only file format supported is Parquet, while the other systems supported also ORC and Avro.

Chapter 3

Method

This chapter defines three methodologies that will be applied sequentially in this project, answering the three **RQs** defined in Section 1.2.1. Section 3.1 defines the system integration process that outputs part of **D3**, i.e., the integration detail. This output will enable the Legacy vs. Icedhops system evaluation defined in Section 3.2, which will output **D1**, i.e., the results of the experiments. The IcedHops experiment results (**D1-partial**) will enable the IcedHops vs. Delta Lake system evaluation, defined in Section 3.3, which will output **D2**, i.e., the results of the comparative experiments. The analysis of **Ds1–2** will be delivered in **D3**, i.e. this comprehensive thesis report.

It has to be noticed that the system evaluation processes, described in Sections 3.2-3.3, are inspired by the system evaluation process of a related work [26]. That related work was hosted by the same company which hosted this thesis work, Hopsworks AB, and it focused on answering **RQs** on some intent similar to the ones described in Section 1.2.1. The related work enabled write and read operations on Delta Lake tables stored on **HopsFS**, and investigated performance differences between that system and the current Hopsworks legacy system. Following similar system evaluation processes will not only allow to answer this project's **RQs**, but will also enable the reader to fairly compare the three technologies investigated in these two thesis works. For the same reason, it eases the comparative process described in Section 3.3, and will make its results more fair and consistent.

3.1 System integration

This Section explains the method and principles used for the system integration process, between PyIceberg and **HopsFS**. This Section is divided

into three Subsections: Integration process, describing the activities to be conducted to integrate PyIceberg and **HDFS**; Requirements; and Development environment, detailing the tools and resources that will be used during the integration process.

3.1.1 Integration process

The integration development process will follow an iterative in Figure 3.1. This project will require numerous interactions with **HopsFS** maintainers (i.e., the industrial supervisors), since PyIceberg library has to interact with **HopsFS**. This review process creates the need for a feedback loop, allowing the system to fit all the stakeholder requirements, described in Subsection 3.1.2. This process answers **RQ1**, to which **Gs1–2** are associated, as described in Section 1.4. The relationships between each process activity and **Gs** are here explained:

1. **Set requirements collaboratively:** this activity is key for the completion of **Gs1–2**, as it is an initial system analysis, performed together with the industrial supervisors, who are knowledgeable on Hopsworks' infrastructure. This task sets the project requirements and investigates what needs to be integrated at a high-level abstraction.
2. **Analyze system and tools:** this activity solves **G1** performing low-level code analysis on system and PyIceberg function calls, understanding what needs to be developed to integrate **HopsFS** and PyIceberg. This activity is the first of an iterative loop accounting for the requirements fulfillment.
3. **Develop integration:** this activity partially solves **G2**, designing and developing a code solution from the information gathered during the analyses above.
4. **Test integration:** this activity partially solves **G2**, verifying the solution validity with integration tests. Failed integration tests or unfulfilled requirements will trigger a new loop iteration. Otherwise, the integration process will be concluded.

This process will produce a part of **D3**, which is the code implementation for IcedHops, the PyIceberg and **HopsFS** integration, including the consideration about the specific tools selected or discarded.

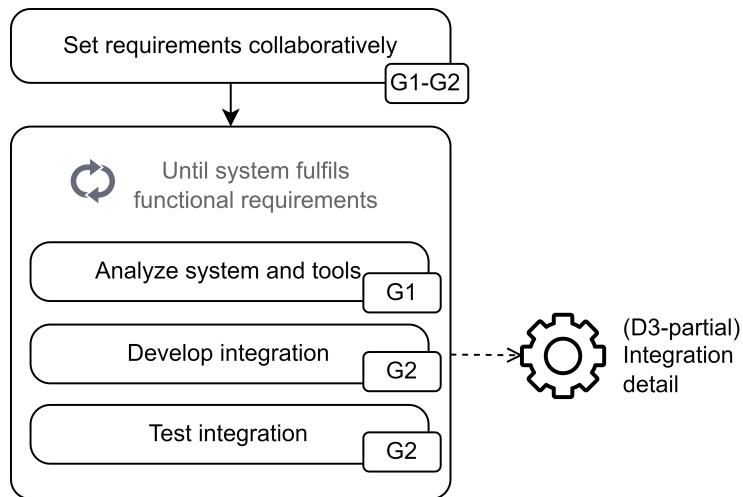


Figure 3.1: Diagram of the system integration process answering RQ1. Each activity is associated to specific Gs. The process produces the integration detail (D3 partial). The loop iterates until the functional requirements, defined in Section 3.1.2, are fulfilled.

3.1.2 Requirements

A series of requirements are defined in agreement with Hopsworks AB, to create a solution that could be later used by the company, in a production environment. The **functional requirements** are:

1. **Write Iceberg Tables:** the solution should allow to write Iceberg tables on **HopsFS** via the PyIceberg library.
2. **Read Iceberg Tables:** the solution should allow to read Iceberg tables on **HopsFS** via the PyIceberg library.

The **non-functional requirements** are:

1. **Consistency:** the solution should be consistent with the current open-source codebase.
2. **Maintainability:** the solution should minimize the need for maintenance and support.
3. **Scalability:** the solution should handle read or write operations on Iceberg tables, up to 100 GB, to support small-scale dataset scenarios.

3.1.3 Development environment

The system implementation will be developed using the following technologies:

- **Computing resources:** the system integration will be developed on a **Virtual Machine (VM)** accessed via **Secure Shell protocol (SSH)** from a computer terminal. Local development will be avoided, due to low reproducibility and complexity of mounting **HopsFS** on a local machine.
- **Code versioning and shared development:** GitHub will be used for versioning and sharing the developed solution.

3.2 System evaluation - Legacy vs. IcedHops

This Section explains the method and principles used for the system evaluation processes, measuring and comparing the performance (latency, in seconds, and throughput, in rows/second) of reading and writing Hudi and Iceberg tables on **HopsFS**. Those operations are conducted respectively on the current legacy system and on IcedHops, integrated in this thesis work.

3.2.1 Evaluation process - RQ2

This evaluation process will follow a sequential approach described in Figure 3.2. Each step of this process is related to one of the **Gs3–6** associated with the **RQ2** in Section 1.4, to which this process answers. The relationships between each process activity and **Gs** are here explained:

1. **Design experiments:** this activity maps perfectly to **G3**, designing the experiments that will be conducted to evaluate the performance difference in performance between the current legacy access to Apache Hudi compared to the IcedHops PyIceberg-based access to Iceberg Tables in **HopsFS**.
2. **Perform experiments:** this activity maps perfectly to **G4**, using the integration detail (**D3-partial**) to develop and conduct the designed experiments on the analyzed systems. Here, data is collected as latency, expressed in seconds.
3. **Transform data according to metrics:** this activity is requisite to fulfill **G5**, and **G8** of **RQ3**. The activity is conducted because throughput is not

directly measured, but it is computed from latency following the formula here below:

$$\text{Throughput (rows/second)} = \frac{\text{Number of rows (rows)}}{\text{Latency (seconds)}}$$

4. **Visualize results:** this activity maps perfectly to **G5**, visualizing the experiments' result according latency, measured in seconds, and throughput, measured in rows/second. This activity also generates **D1**, the experiment results complemented with tables and histograms, presented in Chapter **5**.
5. **Analyze results:** this activity maps perfectly to **G6**, analyzing and interpreting the results delivered in **D1**. This activity contributes to **D3**, generating the analysis of experimental results, presented in Chapter **5**.

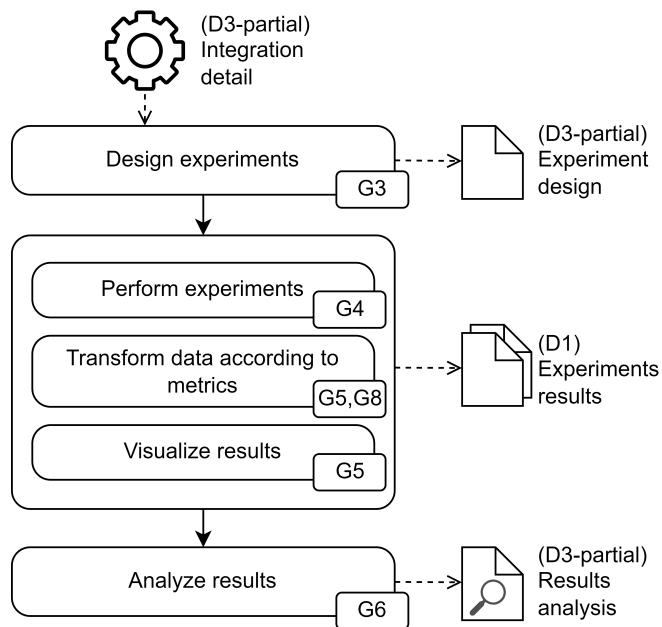


Figure 3.2: Diagram of the system evaluation process answering **RQ2**. Each activity is associated to specific **G**. The process produces two **Ds**, the experiments results (**D1**) and a results analysis (**D3-partial**).

3.2.2 Industrial use case

Several choices must be made for a system evaluation: which data will be used, which environment will run the experiments, and which metrics will be used to evaluate the system. Those choices depends on the scenarios for which the system is created. Thus, this Subsection describes the typical use case for this system. Following, the other Subsections describe which were the decision taken accordingly. While conducting their research in Hopsworks, the author outlined a typical industrial use case for the Hopsworks feature store, by reading internal documentation and discussing with several employees on their customer needs and trends. The use case is described by:

- **Table size:** most of the Hopsworks' customers' workloads were limited (from 1M to 100M rows), with only few clients needing support for massive workloads (more than 1B rows). Thus, this project opted to improve performance for the smaller workloads (from 100k to 100M rows). The dataset selected are presented in Section 3.2.3.
- **Type of data:** the Hopsworks feature store works only with structured data (e.g., numbers, strings), thus experiment design and selected datasets embody this scenario.
- **Rows over storage size:** In the experimental part of this thesis, in order to have a reliable unit measure for table size, the number of rows will be over the storage size (bytes). Perhaps, in the structured data domain of this use case, storage size (bytes) is neither linked to a table structure nor to a storage structure, arguably making this unit measure not reliable (i.e, table with a lot of rows and few columns, and a table with few rows and a lot of columns occupies the same memory).
- **Client configuration:** the client configuration is modeled to reflect typical customers' clients' computational and storage capabilities. Thus, the configuration are limited between one and eight CPU cores, RAM just sufficient for system needs and common SSD storages. The experimental environment is further detailed in Section 3.2.5.

3.2.3 Experimental data

The datasets that will be used to perform read and write experiments come from TPC-H benchmark suite *. TPC-H is a decision support benchmark

*Benchmark suite website available at <https://www.tpc.org/tpch/>

by **Transaction Processing Performance Council (TPC)**, that consists in a collection of business-oriented queries in specific industry sectors [70], which became the de-facto standard for experiments on data storage system. Perhaps, it has been used in related studies [18, 71, 26].

The TPC-H benchmark contains eight tables, and any part of the data can be generated via the TPC-H data generation tool *. The two tables that will be used are the SUPPLIER and the LINEITEM, respectively, the smallest (10k rows) and largest (60M rows) tables. The size (number of rows) of a table depends on the **Scale Factor (SF)**, that can be varied to progressively change in the table size. The SUPPLIER table has seven columns, while the LINEITEM table has sixteen. This difference influences the average size of memory each row occupies. Below for each table their columns are listed, specifying which data type they store.

- SUPPLIER

- S_SUPPKEY : identifier
- S_NAME : fixed text, size 25
- S_ADDRESS : variable text, size 40
- S_NATIONKEY : identifier
- S_PHONE : fixed text, size 15
- S_ACCTBAL : decimal
- S_COMMENT : variable text, size 101

- LINEITEM

- L_ORDERKEY : identifier
- L_PARTKEY : identifier
- L_SUPPKEY : identifier
- L_LINENUMBER : integer
- L_QUANTITY : decimal
- L_EXTENDEDPRICE : decimal
- L_DISCOUNT : decimal
- L_TAX : decimal

*Available at https://www.tpc.org/tpc_documents_current_versions/current_specifications5.asp

- L_RETURNFLAG : fixed text, size 1
- L_LINESTATUS : fixed text, size 1
- L_SHIPDATE : date
- L_COMMITDATE : date
- L_RECEIPTDATE : date
- L_SHIPINSTRUCT : fixed text, size 25
- L_SHIPMODE : fixed text, size 10
- L_COMMENT : variable text, size 44

Considering the different structure of the two tables used (i.e., number of columns and data types) comparison across different tables cannot be done using the selected metrics, i.e., latency (seconds) and throughput(rows/second). For this reason, the system evaluation will only consider same tables on different configuration. This project used five table variations to benchmark the system integration. SF was varied to obtain a table at each significant order of magnitude from 10k to 60M rows. These are the tables:

1. *supplier_sf1*: size = 10000 rows
2. *supplier_sf10*: size = 100000 rows
3. *supplier_sf100*: size = 1000000 rows
4. *lineitem_sf1*: size = 6000000 rows
5. *lineitem_sf10*: size = 60000000 rows

Lastly, despite no information are given about the row and storage size (bytes) ratio, in accordance to Section 3.2.2, this Section presents comprehensive information on the data, including how to retrieve it and how it is composed. These gives to the reader the ability of calculating the storage occupancy of each table, depending on their preferred unit of measure.

3.2.4 Experimental design

The experiments aim to highlight the differences between the current legacy system and IcedHops. Two experimental pipelines were designed to isolate the performance of those two systems:

1. **Legacy pipeline:** is the current Hopsworks feature store which stores data in Hudi tables. This system uses a pipeline based on Kafka, and Spark to write data on the Hudi tables, saved on **HopsFS**. The pipeline uses a Spark alternative, DuckDB, and Arrow Flight to read data. This is described in Sections [2.5.1-2.5.2](#).
2. **IcedHops:** is the system implemented in Chapter [4](#), which stores data in Iceberg Tables. This systems uses a pipeline based on PyIceberg, SQLite, and Arrow Flight to both write to and read from the Iceberg tables, saved on **HopsFS**. This is described in Sections [2.5.3-2.5.4](#).

The experiments will verify how the performance will change based on different **CPU** resources provided: one, two, four, and eight cores, according to the typical Hopsworks use case (Section [3.2.2](#)). Each time, the experimental environment will be modified, creating a new Jupyter server where the host the experiments. The data used for experiments, as described in Section [3.2.3](#), will come from two different tables, modified according to a **SF**, for a total of five times for each table. For the writing experiments conducted on the legacy system, different parts of the whole writing process will be measured, to verify how different parts of the legacy system will scale in the different environments described above. Those two parts are the upload and the materialization, which are explained in Section [2.5.1](#). This will also help defining whether and what part of the legacy pipeline is a bottleneck.

In conclusion, the experiments conducted will be a total of two (pipelines) times four (**CPU** configurations) times five (tables) times two (read and write operations), thus 80 experiments, performed each fifty times, to ensure statistical significance to the results.

3.2.5 Experimental environment

The experimental environment consists of a physical machine in Hopsworks AB' offices, virtualized to enable remote shared development. The **CPU** details of the machine are present in Listing [3.1](#), noting that only eight cores at maximum were dedicated during the experiments. The machine mounts about 5.4 TBs of **SSD** memory, allowing for fast read and write speed, respectively 2.7 GB/s, and 1.9 GB/s (measured with a simple *dd* bash command). The experimental environment will be set up with a Jupyter Server of different CPU cores, depending on the experiment. The Jupyter server is allocated by default with 2048MB, which will be adjusted during the experiments, according to the needs (see Section [5.1.4](#)).

Notice that, despite this experimental environment is virtualized in isolation, it runs on shared resources, thus experiments result might vary depending on the machine's load. To mitigate this, all experiments will be run when the machine load is low (less than 50% of **CPU** and **RAM** usage).

Listing 3.1: Output of a *lscpu* bash command on the machine.

Architecture :	x86_64
CPU op-mode(s) :	32-bit , 64-bit
Address sizes :	48 bits physical , 48 bits virtual
Byte Order :	Little Endian
CPU(s) :	32
On-line CPU(s) list :	0-31
Vendor ID :	AuthenticAMD
Model name :	AMD Ryzen Threadripper PRO 5955WX 16-Cores
CPU family :	25
Model :	8
Thread(s) per core :	2
Core(s) per socket :	16
Socket(s) :	1
Stepping :	2
Frequency boost :	enabled
CPU max MHz :	7031.2500
CPU min MHz :	1800.0000
BogoMIPS :	7985.56
Virtualization features :	
Virtualization :	AMD-V
Caches (sum of all) :	
L1d:	512 KiB (16 instances)
L1i:	512 KiB (16 instances)
L2:	8 MiB (16 instances)
L3:	64 MiB (2 instances)

3.2.6 Evaluation framework

The system evaluation framework will evaluate the different system on three key aspects:

1. **Functional requirements:** the success or failure of an experiment run will determine whether those are matched or not, and how much. This

is an additional step, since only in the system integration process, the process is stopped and restarted from the beginning if the functional requirements are not met. These requirements are described in Section [3.1.2](#).

2. **Non-functional requirements:** Consistency and maintainability are mainly addressed during integration. On the other hand, scalability is measured during the system evaluation experiments designed below. The metric used for measuring this requirement is the throughput, as defined in [RQ1](#).
3. **How does the legacy pipeline compare to IcedHops?** this question answers directly [RQ2](#), measuring the throughput of the pipelines defined in Section [3.2.4](#). Results are then compared using a visual approach.

3.2.7 Reliability and validity

Experiments result are significant according to their reliability and validity. Each experiment will be performed fifty times to ensure the reliability of the its results on the system performance. This number was agreed to balance the results' consistency and resource efficiency (i.e., time and computing resources), as described in Section [1.5](#). Secondly, due to the complex nature of the pipelines used, the data distribution of results might vary from one experiment to the other. Thus, a bootstrapping technique will be employed to restore the validity of the data collected. Data will be resampled with substitutions a thousand times, then for each experiment an average measure, altogether with a confidence interval.

3.3 System evaluation - IcedHops vs. Delta Lake

This Section firstly explains the steps of this system evaluation processes, aimed at measuring and comparing the performance (latency, in seconds, and throughput, in rows/second) of reading and writing Iceberg and Delta Lake tables on [HopsFS](#). Those experiments will be conducted respectively on IcedHops, integrated in this thesis work, and on the delta-rs-based system, implemented in a related work [\[26\]](#). Then, it presents the evaluation framework used by in this process.

3.3.1 Evaluation process - RQ3

This evaluation process will follow a sequential approach described in Figure 3.3. Each step of this process is related to one of the **Gs**7–9 associated with the **RQ3** in Section 1.4. The relationships between each process activity and **Gs** are here explained:

1. **Analyze related work results:** this activity maps perfectly to **G7**, analyzing the results coming from the related work taken as input by this process.
2. **Visualize results:** this activity maps perfectly to **G8**. It takes as input also PyIceberg experiments results (**D1-partial**), visualizing thus the experiments' result according to latency and throughput. This activity generates **D2**, the comparative experiments results complemented with tables and histograms, presented in Chapter 5.
3. **Analyze results:** this activity maps perfectly to **G9**, analyzing and interpreting the results delivered in **D2**. This activity contributes to **D3**, generating the comparative analysis of the experiment results, presented in Chapter 5.

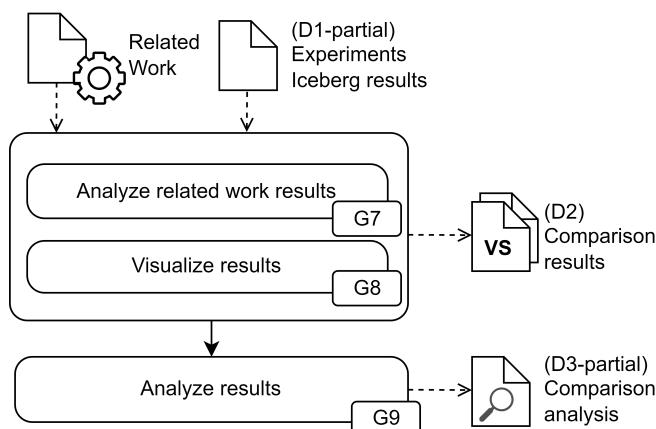


Figure 3.3: Diagram of the system evaluation process answering **RQ3**. Each activity is associated to specific **G**. The process produces two **Ds**, the comparative experiments results (**D2**) and a comparative results analysis (**D3-partial**).

3.3.2 Evaluation framework

This system evaluation framework is similar to the system evaluation framework described in Section 3.2.6, but will not focus on functional requirements, as they will be already assessed in previous Section, for IcedHops, and they have been already assessed in the related work, for delta-rs. Thus, this evaluation framework will evaluate the different system on two key aspects:

1. **Non-functional requirements:** Consistency and maintainability are mainly addressed during integration, while scalability is measured during the system evaluation experiments. The metric used for measuring this requirement is the throughput, as defined in RQ2.
2. **How does IcedHops compare to delta-rs?** this question answers directly RQ3, measuring the throughput of IcedHops defined in Section 3.2.4 and the delta-rs pipeline [26]. Results are then compared using a visual approach.

Chapter 4

Implementation

This Chapter follows first the system integration process, defined in Section 3.1, describing which integration choices were taken, which components were selected and their usage. Thus it explains how was conducted the experimental part of this thesis, presented in Section 3.2.

4.1 Integration design and usage

The first step of the system integration process consisted of analyzing the Hopsworks system and the PyIceberg tools, as described in Section 3.1 and it was subsequently divided in several activities. First was drawn a full picture of all the technologies behind the Hopsworks online and offline feature store, as well as their integration capabilities with third parties software, by reading the development documentation and by interviewing the colleagues in Hopsworks AB. Secondly, the Hopsworks feature store API's code performing write and read operations on Hudi tables on HopsFS was studied and broken-down in detail. Then was conducted a thorough research within scientific and commercial domains about Iceberg, PyIceberg, their capabilities, limitations and existing integrations with other file systems and cloud storages. Lastly, this research was completed looking at alternatives to PyIceberg and alternatives to Iceberg, seeking methods and implementations already developed that could have helped the integration of IcedHops. This information has been summarized in Sections 2.2–2.5.

All those activities, while accounting for more than two months of research, were key since permitted to identify which software parts had to be integrated and which pieces of technology had to be plugged-in to satisfy the requirements described in Section 3.1.2, thus to be able to write on and

read from Iceberg Table hosted on [HopsFS](#), via the PyIceberg library. These activities outlined the need of a catalog, a query engine, and a FileIO to be plugged-in. While the formers are fundamental components of any Data lakehouse architecture, as explained in Section [2.2.2](#), the latter is a pluggable module for performing **CRUD** operations on files, specifically required by PyIceberg. To verify the integrability of each of those technology with PyIceberg and HopsFS, and evaluated their performances, several Jupyter notebooks were created *. While only the notebooks using the working technologies with the correct methods and workarounds are presented in the repository, this second step was reached to the production of approximatively 2500 [Lines Of Code \(LOC\)](#).

PyIceberg, being a rather recently developed library, does not features yet all the integrations of older Iceberg libraries [\[56\]](#), such the Java [API](#). In addition, despite [HopsFS](#) expose the same methods of [HDFS](#), the environment where HopsFS was mounted, described in Section [3.2.5](#), did not allow to use some catalogs. Sections [4.1.1-4.1.2](#) describe the choices taken for both catalog (SQL Catalog) and query engine (DuckDB), describing the reason behind those choices and the problem encountered. Regarding the FileIO, since there was a single compatible option (PyArrowFileIO), this component could not have been subject of any design comparison. Lastly, Section [4.1.3](#) describe how instanciate the integrated system, and how to perform operations over it.

4.1.1 Catalog choice

The discared Catalogs are:

- **REST**: it is supported by [HopsFS](#). However, since this would have need to develop the interface from scratch, thus was discarded as not fulfilling the maintainability not-functional requirement, described in Section [3.1.2](#).
- **AWS Glue**: it is supported by [HDFS](#), but it did not pass the integration test with [HopsFS](#), thus was discared. Additionally, since it is a proprietary solution of [AWS](#), this would have lowered the reproducibility of the experiments later conducted, due to the additional costs of this solution.
- **AWS DynamoDB**: it is supported by [HopsFS](#). Was however discarded, for the same reproducibility reason explained above.

*Code available at: <https://github.com/SebastianoMeneghin/apache-iceberg-as-offline-feature-store>

- **HMS**: it is supported by **HopsFS**. **HMS** is however a complex tool, developed to be tightly integrated with MapReduce and Spark environment, and perhaps perform its best in large-scale data scenarios. This was discarded since not matching with purpose of avoiding Spark and its environment, and the industrial use case described in Section 3.2.2. Furthermore, this did not fulfill the maintainability not-functional requirement.

The **chosen** catalog is:

- **SQL Catalog**: it is supported by **HopsFS**, and it could be instantiated on a SQLite database supported by PyIceberg. This was the choice for the system integration, as it is an open-source catalog, the most light-weight option among the alternatives, and needs few lines of code to be used, fulfilling both the not-functional and the functional requirements. This solution suits perfectly small-scale scenarios such this thesis' use case, but it does not fit a large-scale scenario. Furthermore, the specific SQLite database is not built for concurrency.

4.1.2 Query engine choice

In the choice of the query engine, all the candidates are known to be suitable for both **HopsFS**, since this all of these engines are supported by Hopsworks AI Data Platoform, which uses HopsFS are data storage layer. The discarded query engines are:

- **AWS Athena and Snowflake**: both were discarded since they are proprietary solutions. This would have lowered the reproducibility of the experiments later conducted, due to the additional costs of this solution.
- **Spark**: this was discarded since it direclty violates the purpose of this project, i.e. create a Spark alternative to read and write data on **OTF** stored on **HopsFS**.
- **Presto, Trino, Flink**: were discarded as designed to perform their best in large-scale data scenarios, thus it did not suit the industrial use case described in Section 3.2.2. Additionally, it did not fulfill the maintainability not-functional requirement, describe in Section 3.1.2.

The **chosen** query engine is:

- **DuckDB**: was the choice for the system integration, as it is a portable open-source **OLAP-DBMS**, which proved to be the best performing engine in related work on small-scale scenarios [18, 16].

4.1.3 Usage

Once selected PyArrowFileIO as FileIO, SQLite as support to SQL Catalog, and DuckDB as query engine, all the libraries and dependencies are directly managed by the installation of the PyIceberg library, using the command `pip install pyiceberg[pyarrow, duckdb, sqlite]`, as described on PyIceberg documentation [56]. This integration supports all PyIceberg methods, but this Section will focus only on the methods used for the experiments. Listing 4.1 describe how to instantiate an Iceberg catalog using SQLite, to enable metadata management on **HopsFS** (or **HDFS**), and how to create a namespace and a table within the namespace. Following this, an example of write operation on **HopsFS** (or **HDFS**) is described in Listing 4.2, while Listing 4.3 provides an example of read operation on **HopsFS** (or **HDFS**).

Listing 4.1: Instantiating an Iceberg catalog using SQLite

```
from pyiceberg.catalog.sql import SqlCatalog

catalog = ("default", **{
    "uri" : "sqlite:///catalog_path",
    "warehouse" : "hdfs_path",
    "hdfs.host" : "hdfs.host"})
catalog.create_namespace("ns")
table = catalog.create_table(
    "ns.table",
    schema=your_df.schema,
    location="hdfs_path",)
```

Listing 4.2: Writing a DataFrame with IcedHops on an Iceberg Table stored on HopsFS (or HDFS).

```
import pandas as pd
from pyiceberg.catalog import load_catalog

df = pd.DataFrame({ "num": [1, 2, 3],
                     "letter": ["a", "b", "c"]})
catalog = load_catalog()
table = catalog.load_table("ns.table")
table.append(df)
```

Listing 4.3: Reading a table with IcedHops from an Iceberg Table stored on HopsFS (or HDFS).

```
from pyiceberg.catalog import load_catalog

catalog = load_catalog()
table = catalog.load_table("ns.table")
df = table.scan().to_arrow()
```

4.2 Experimental setup

As defined in Section 3.2.4, experiments run different system configurations with five tables fifty times per experiment. The run time of the experiments, i.e., the read and write latency, was measured using two different approaches. This decision is motivated by the need to measure accurate results, while it was impossible to use the most accurate method for measuring across all systems. The initial technique utilizes Python’s `timeit` function, which ensures the process executing the code runs in isolation, providing a reliable approximation of operation latency. As shown in Listing 4.4, `timeit` is implemented by specifying a `SETUP_CODE` that executes before the experiment and a `TEST_CODE` that is measured during execution, with the resulting time (in seconds) returned by the `timeit` function. This method was chosen because `timeit` offers a straightforward interface for executing and timing a concise code snippet, and was used to conduct the experiments about the read operations of legacy system, described in Section 2.5.2.

Listing 4.4: Timeit usage to measure the time to read from an Iceberg table stored on HopsFS.

```

import timeit
SETUP_CODE= '''
from pyiceberg.catalog import load_catalog '''


TEST_CODE= '''
catalog = load_catalog()
table   = catalog.load_table("ns.table")
df      = table.scan().to_arrow() '''


# Measure the execution runtime
write_result = timeit.timeit(setup = SETUP_CODE,
                             stmt    = TEST_CODE,
                             number = 1)

```

Running in an isolated Python environment, the timeit approach had two limitations: (1) it was not possible to breakdown the two main steps performed by the legacy system when performing a write operation, (2) it was not possible to separate the timings of the creation of a new catalog and of the read operations in IcedHops, since timeit cannot access variables declared elsewhere (this catalog creation/deletion was necessary at every step, to the machine to perform caching). Thus, for those two cases, see Sections 2.5.1 and 2.5.4, a second approach was implemented. This approach consisted in recording the time before and after the script run, using the function time, from the Python standard library called time, and an example of this approach is showed in Listing 4.5. The usage of two different approaches was not considered problematic during the experiments, as some trials revealed that the latency measured by the two methods was equal within a 95% confidence interval.

Listing 4.5: A simple time difference approach the time to read from an Iceberg table stored on HopsFS.

```
import time
from pyiceberg.catalog import load_catalog

catalog = load_catalog()

before = time.time()
table = catalog.load_table("ns.table")
df    = table.scan().to_arrow()
after = time.time()

reading_time = after - before
```

Chapter 5

Results and Analysis

This Chapter is the output of both system evaluation processes defined in Section 3.2–3.3. It starts with Section 5.1, which presents the results of both the experiments and the comparison in tables, histograms, and written descriptions. Then, Section 5.2 complements the chapter by analyzing and discussing the major findings.

5.1 Major results

This Section presents the main results of the eighty experiments performed as defined in Section 3.2.4, and the forty experiments performed in the related work [26], conducted with the same experimental design and environment. The experiments are grouped into subsections according to the measured operation, i.e., read or write. Each Subsection presents histograms and tables to visualize the results using both metrics, i.e., latency (measured in seconds) and throughput (measured in rows/seconds), defined in RQ2.

It must be noted that the results are visualized using a log scale for clarity, as results differing from more than one significant figure would not be clearly interpreted using a linear scale, in a histogram representation. For each measurement, a 95% confidence interval was calculated using the bootstrapping technique mentioned in Section 3.2.7. For better readability, this interval is not reported in the histograms of this Section, but it is reported in Appendices C–D, with a histogram and a table for each experiment expressed in both metrics.

Of the two metrics, only latency was measured during the experiments, while throughput was calculated with the formula present in Section 3.2.1. Latency and throughput are inversely related by a constant factor, as all

experiments were conducted with fixed-size tables. This implies that halving the latency results in a doubling of throughput, while reducing latency to a quarter results in a quadrupling of throughput. Since the results are described according to both metrics, this introduces redundancy. To minimize repetition, trends will primarily be discussed in terms of latency, with throughput trends highlighted only when they exhibit significantly different behavior.

5.1.1 Write experiments

Figures 5.1–5.2, and Tables 5.1–5.2 are showing the outcome of the write experiments conducted. The outcomes are presented in latency in Figure 5.1 and Table 5.1, while those are presented in throughput in Figure 5.2 and Table 5.2. The experiments were conducted on each of the two systems defined in Section 3.2.4, as well as on the delta-rs rust-based system described in the related work [26]. The five tables of different sizes being written were defined in Section 3.2.3.

Both histograms, i.e., Figures 5.1–5.2, report the results of the experiment performed with one CPU core. Instead, Tables 5.1–5.2 also present a calculated percentage of improvement (decrease in the case of latency, increase in the case of throughput) of the metric as the CPU cores increase.

Legacy pipeline vs. IcedHops

The latency on write operations measured using IcedHops is more than fourty times lower than the latency measured using the legacy pipeline, in all the experiments. This trend is more prominent for bigger tables (6M and 60M rows), where latency measured using IcedHops is less than one hundredth of the latency measured using the legacy system.

IcedHops vs. delta-rs

The latency on write operations measured using IcedHops is lower than the latency measured using the delta-rs pipeline, in all the experiments. This trend is more prominent for the biggest table (60M rows), where latency measured using IcedHops is around seven times lower than the latency measured using the delta-rs pipeline, while it is lesser prominent the smaller are the tables, with approximatively equal latency for the smallest table (10K rows).

Table 5.1: Write experiment results expressed as latency. Experiments performed with multiple **CPU** cores are expressed as latency percentage decrease compared to the one **CPU** core experiment.

Pipeline	Number of rows	1 CPU core latency (seconds)	2 CPU cores (% decrease)	4 CPU cores (% decrease)	8 CPU cores (% decrease)
Hudi Legacy	10K	50.23794	-0.98	-2.06	-1.97
	100K	59.55605	-0.37	0.07	-1.23
	1M	112.17773	3.22	3.00	2.49
	6M	511.84908	7.52	5.84	7.01
	60M	2716.20939	13.81	13.63	14.41
Iceberg IcedHops	10K	1.20406	0.75	0.73	-1.18
	100K	1.26524	5.19	0.83	1.32
	1M	1.71770	-0.22	5.01	-0.08
	6M	3.57649	0.28	0.88	1.23
	60M	24.60989	4.48	2.46	4.01
Delta Lake	10K	1.25088	-0.95	2.72	-8.70
	100K	1.36800	4.44	2.30	5.54
	1M	9.38231	9.19	10.27	11.59
	6M	19.75149	17.47	17.86	20.39
	60M	177.19458	24.29	29.99	31.21

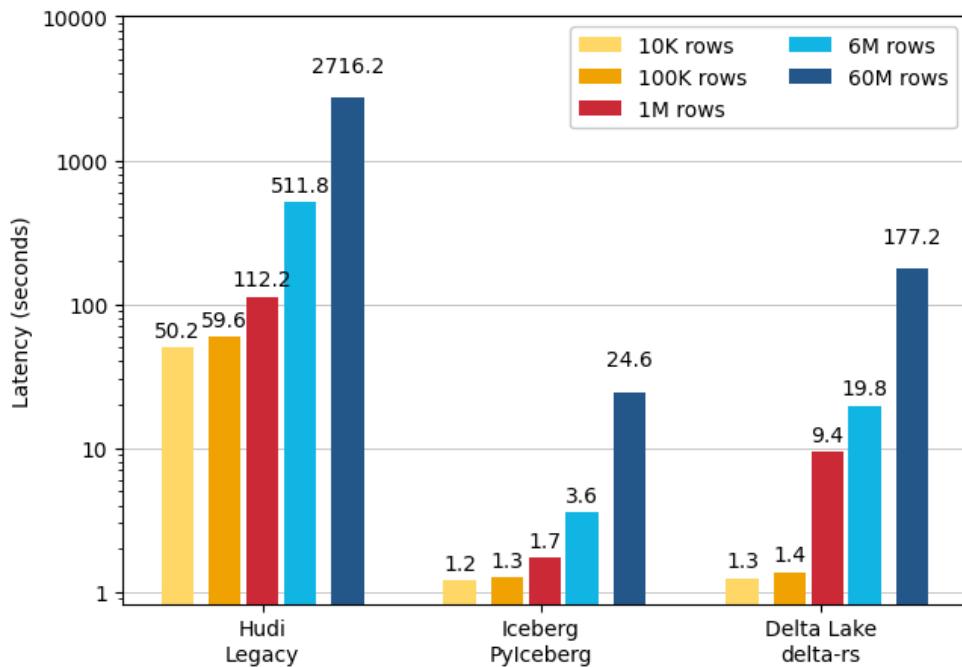


Figure 5.1: Histogram in log-scale of the write experiment results expressed as latency. The experiment was performed with one **CPU** core.

Table 5.2: Write experiment results expressed as throughput. Experiments performed with multiple **CPU** cores are expressed as throughput percentage increase compared to the one **CPU** core experiment.

Pipeline	Number of rows	1 CPU core throughput (k rows/second)	2 CPU cores (% increase)	4 CPU cores (% increase)	8 CPU cores (% increase)
Hudi Legacy	10K	0.199 05	-0.97	-2.02	-1.93
	100K	1.679 09	-0.37	0.07	-1.21
	1M	8.914 43	3.33	3.10	2.56
	6M	11.722 21	8.13	6.20	7.54
	60M	22.089 61	16.02	15.78	16.84
Iceberg IcedHops	10K	8.305 21	0.75	0.74	-1.16
	100K	79.036 58	5.47	0.84	1.33
	1M	582.172 95	-0.22	5.28	-0.08
	6M	1 677.621 61	0.28	0.89	1.25
	60M	2 438.044 52	4.69	2.52	4.18
Delta Lake delta-rs	10K	7.994 35	-0.94	2.80	-8.00
	100K	73.099 24	4.65	2.35	5.87
	1M	106.583 58	10.12	11.45	13.11
	6M	303.774 53	21.16	21.75	25.61
	60M	338.610 80	32.08	42.83	45.37

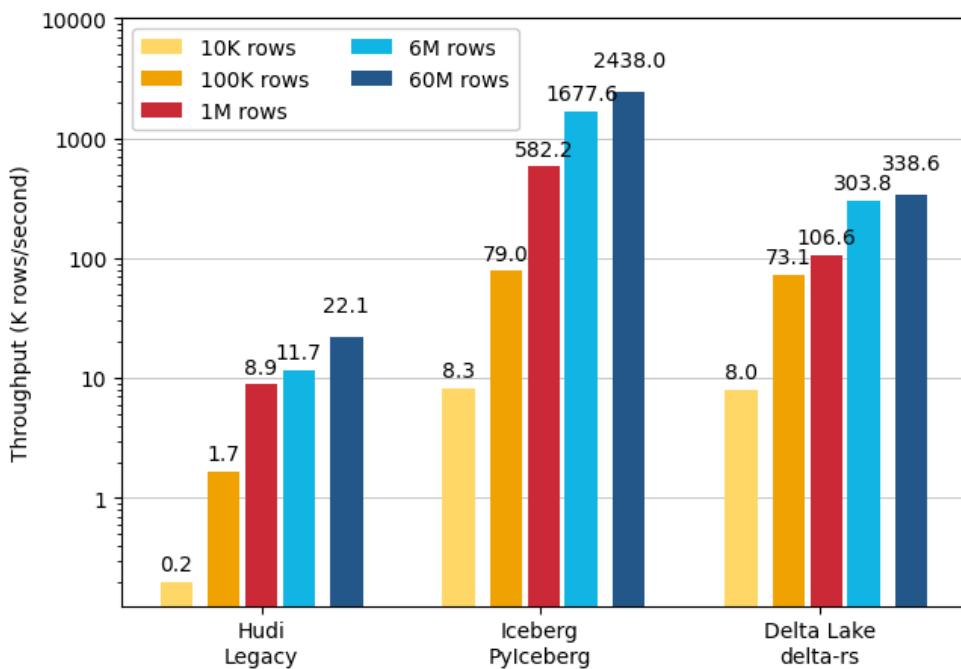


Figure 5.2: Histogram in log-scale of the write experiment results expressed as throughput. The experiment was performed with one **CPU** core.

Change of performance as the CPU cores increase

The experiments with increased **CPU** cores showed different results across the different pipeline technologies. For the legacy system, significant latency reductions were only observed for the largest table (60M rows), with a decrease of approximately 14%. Smaller tables (10K and 100K rows) actually saw slight increases in latency, and even the larger tables (1M and 6M rows) only improved by a maximum of 7%. IcedHops saw some benefit from additional **CPU** cores, but the improvement was limited, with only a maximum of 5% decrease in latency for some tables (100K, 60M). The most substantial gains were seen with delta-rs, where the larger tables (6M and 60M rows) experienced considerable improvements of 20-30%, and smaller tables (100K and 1M rows) only saw modest latency decreases of 5-10%.

5.1.2 Read experiments

Figures 5.3–5.4, and Tables 5.3–5.4 are showing the outcome of the write experiments conducted. The outcomes are presented in latency in Figure 5.3 and Table 5.3, while those are presented in throughput in Figure 5.4 and Table 5.4. The experiments were conducted on each of the two systems defined in Section 3.2.4, as well as on the delta-rs rust-based system described in the related work [26]. The five tables of different sizes being read were defined in Section 3.2.3.

Both histograms, i.e., Figures 5.3–5.4, report the results of the experiment performed with one **CPU** core. Instead, Tables 5.3–5.4 also present a calculated percentage of improvement (decrease in the case of latency, increase in the case of throughput) of the metric as the **CPU** cores increase.

Legacy pipeline vs. IcedHops

The latency of read operation measured using IcedHops results from 55% up to sixty times lower than the latency measured using the legacy pipeline. The most prominent reduction (sixty times) is obtained in the table containing 100k rows, while the less prominent (55%) is obtained with the largest table (60M rows). Regarding the other tables, IcedHops has around fifteen times lower read operation latency than the legacy system.

IcedHops vs. delta-rs

The latency on read operations measured using IcedHops is completely comparable with the latency measured using the delta-rs pipeline, if not for

Table 5.3: Read experiment results expressed as latency. Experiments performed with multiple **CPU** cores are expressed as latency percentage decrease compared to the one **CPU** core experiment.

Pipeline	Number of rows	1 CPU core latency (seconds)	2 CPU cores (% decrease)	4 CPU cores (% decrease)	8 CPU cores (% decrease)
Hudi Legacy	10K	0.63144	1.05	-0.76	0.65
	100K	2.65043	-0.49	0.40	-0.45
	1M	8.59296	-0.16	-1.76	2.80
	6M	33.52580	0.44	0.23	0.26
	60M	33.69031	0.18	0.09	1.63
Iceberg IcedHops	10K	0.01723	-6.36	-72.94	1.53
	100K	0.04687	7.89	8.62	8.95
	1M	0.43018	37.37	45.75	46.04
	6M	2.12331	38.96	52.50	58.34
	60M	21.81955	38.07	53.66	59.49
Delta Lake delta-rs	10K	0.05345	22.90	18.50	19.29
	100K	0.05763	0.73	3.78	5.18
	1M	0.53878	56.62	64.81	67.69
	6M	1.94947	53.41	72.76	74.50
	60M	22.98186	50.35	75.72	87.19

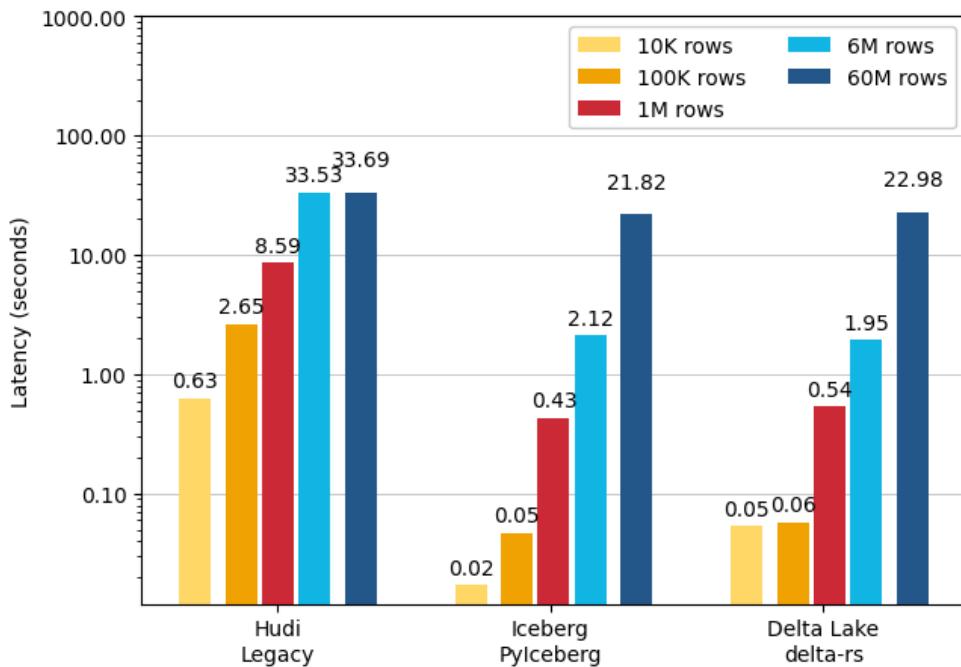


Figure 5.3: Histogram in log-scale of the read experiment results expressed as latency. The experiment was performed with one **CPU** core.

Table 5.4: Read experiment results expressed as throughput. Experiments performed with multiple **CPU** cores are expressed as throughput percentage increase compared to the one **CPU** core experiment.

Pipeline	Number of rows	1 CPU core throughput (k rows/second)	2 CPU cores (% increase)	4 CPU cores (% increase)	8 CPU cores (% increase)
Hudi Legacy	10K	15.83673	1.06	-0.75	0.66
	100K	37.72979	-0.48	0.40	-0.45
	1M	116.37431	-0.16	-1.73	2.88
	6M	178.96662	0.44	0.23	0.26
	60M	1780.92764	0.18	0.09	1.65
Iceberg IcedHops	10K	580.48003	-5.98	-42.18	1.56
	100K	2133.40362	8.57	9.43	9.83
	1M	2324.62963	59.66	84.33	85.33
	6M	2825.77691	63.83	110.53	140.02
	60M	2749.82753	61.47	115.82	146.85
Delta Lake	10K	187.08171	29.70	22.69	23.89
	100K	1735.08031	0.74	3.93	5.47
	1M	1856.03325	130.51	184.19	209.48
	6M	3077.75801	114.64	267.09	292.19
	60M	2610.75520	101.42	311.87	680.55

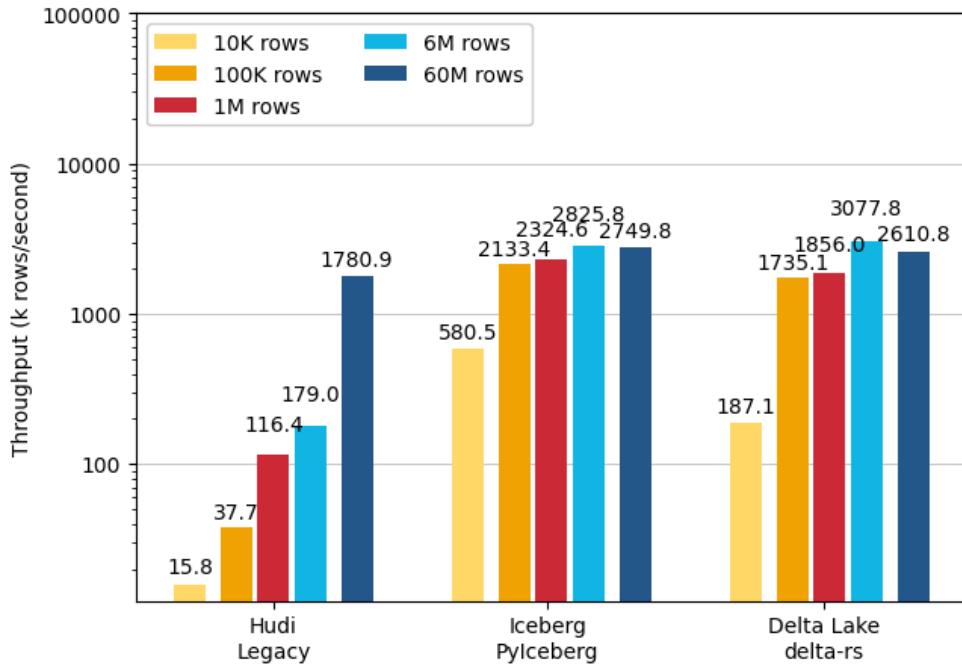


Figure 5.4: Histogram in log-scale of the read experiment results expressed as throughput. The experiment was performed with one **CPU** core.

the smallest table (10K), where IcedHops latency was three times lower than delta-rs pipeline latency. For the other tables, The latency on write operations measured using IcedHops is lower than the latency measured using the delta-rs pipeline, in all the experiments. This trend is more prominent for the biggest table (60M rows), where latency measured using IcedHops is around seven times lower than the latency measured using the delta-rs pipeline, while it is lesser prominent the smaller are the tables, with approximatively equal latency for the smallest table (10K rows).

Change of performance as the CPU cores increase

The experiments with increased **CPU** cores showed different results across the different pipeline technologies. In the legacy pipeline, increasing the number of **CPU** cores had minimal impact on latency reduction. The observed improvements remained below 1% for most table sizes, with some cases even showing slight increases. IcedHops saw some benefit from utilizing multiple **CPU** cores. While the smallest table (10K rows) showed a slight increase in latency, all other tables benefited from latency reductions. For tables larger than 100K rows, the reduction was more pronounced, reaching around 39% for mid-sized and large tables (1M, 6M, and 60M rows). Delta-rs showed the most substantial gains in read latency with additional CPU cores. For small tables (10K and 100K rows), the letancy reductions remained below 5%, but for tables larger than 100K rows, the reduction was significant, reaching around 55% for mid-sized table (1M rows) and around 50% for larger tables (6M and 60M rows). The trend continued with additional cores, further enhancing performance as table sizes increased.

5.1.3 Legacy pipeline write latency breakdown

Table 5.5 and Figure 5.5 present an analysis of the write latency in the legacy pipeline, dividing it into two phases: upload time and materialization time, as detailed in Section 2.5.1. This breakdown is applied to all five tables introduced in Section 3.2.3. 5.5 illustrates the results obtained from the experiment using a single **CPU** core, whereas Table 5.5 provides both the data from this experiment and the computed percentage reduction in latency as the number of **CPU** cores increases. Regarding the impact of upload time on overall write latency, it accounts for only a small fraction (approximately 5%) when dealing with smaller tables (containing 10K or 100K rows). However, as table sizes increase, this contribution follows a nearly linear growth trend. This is indeed highlighted in the results experimennts using the tables of 100K,

Table 5.5: Contributions to the write latency of the upload and materialization steps in the legacy pipeline. Experiments performed with multiple CPU cores are expressed as latency percentage decrease compared to the one CPU core experiment.

Number of rows	1 CPU core		2 CPU cores		4 CPU cores		8 CPU cores	
	latency (seconds)	(% decrease)	upl.	mat.	upl.	mat.	upl.	mat.
10K	2.50	47.73	4.50	-1.28	4.53	-2.46	4.65	-2.32
100K	3.67	55.91	6.36	-0.79	5.56	-0.30	6.28	-1.68
1M	22.59	89.60	17.51	-0.37	14.88	0.01	16.49	-1.03
6M	244.60	267.25	15.83	-0.10	13.47	-1.15	15.09	-0.39
60M	2438.22	278.16	15.34	0.43	15.17	0.20	15.96	0.86

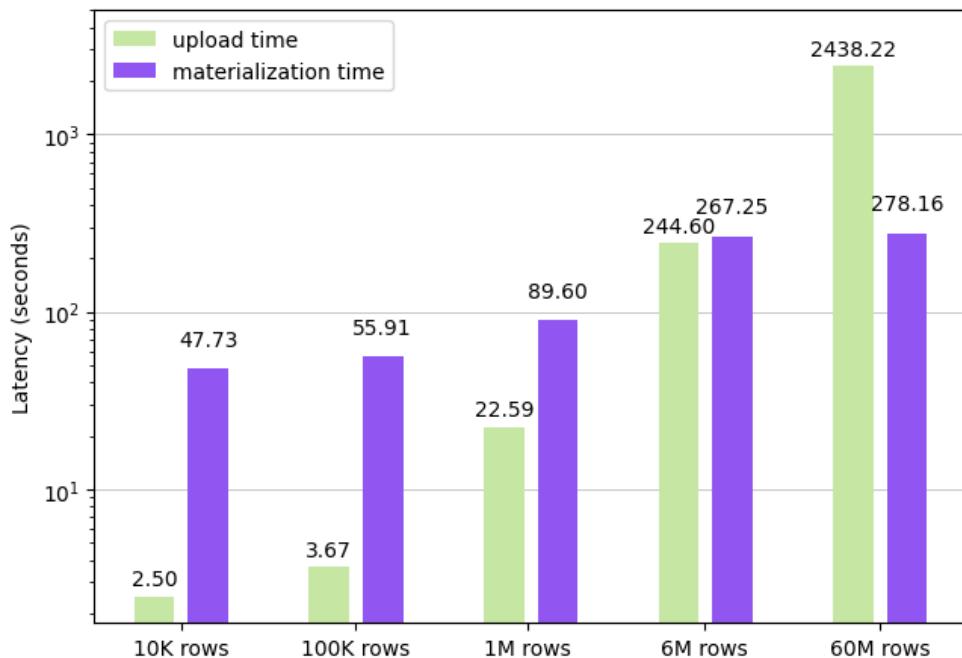


Figure 5.5: Histogram in log-scale displaying the contributions to the write latency of the upload and materialization steps in the legacy pipeline. The experiment was performed with one CPU core.

6M, and 60M rows. This radically changes the proportion between the upload and materialized time contributions, making the upload time 90% of the total write latency for the largest table (60M rows). Differently, the materialization time contribution makes up 95% of the total write latency for the smallest table (10K rows), but then its absolute value does not increase by more than a significant figure even if the table size is increased by three significant figures.

Analyzing the impact of multiple **CPU** cores, the upload time benefits significantly from increased core availability, particularly for larger tables, where latency decreases by 15% compared to just 4% for smaller tables. In contrast, the materialize time remains largely unaffected, showing only minor fluctuations, with latency varying by 1-2%.

5.1.4 In-memory resources usage

The resources of the experimental environment, defined in Section 3.2.5, were adjusted to match computational needs. The write operations were demanding on the available **RAM** resources, requiring up to 24 GBs to operate with the larger tables (6M and 60M rows). Thus, on were allocated 32768 MBs of **RAM**, to avoid slowing down operations.

5.2 Results analysis and discussion

This Section discusses the major results, described in detail in Section 5.1, explaining what they mean, which are their implication for Hopsworks AB, and for the research area. Following, this Section contains the author's considerations over the legacy system and IcedHops, which was developed while conducting this thesis research. Additional consideration regarding the delta-rs library can be found in the related work's discussion [26].

5.2.1 Discussion on major results

The results presented in Section 5.1 disclose a substantial difference in the performance (both latency and throughput) of the legacy system, and IcedHops. Despite the two systems were tested on two different tasks, i.e. write and read operations, IcedHops integrated in this thesis shown from 15 to 140 times lower latency than the legacy system, in all experiments conducted using tables from 10K to 6M rows. Furthermore, the results reveal that IcedHops and the delta-rs system differ generally much less, with IcedHops outperforming the delta-rs system only in write operations with table bigger

than 1M rows, where the system latency was up to seven times lower than its counterpart. These findings confirm the hypothesis that alternatives to a Spark-based system exist, and that those are preferable when operating on small tables (from 10K to 60M rows). Furthermore, the large performance differences observed indicate substantial room for optimization in the legacy system, on both read and write operations.

Going more in depth in the results of the write experiments, IcedHops has a latency from 40 to 140 times lower than the legacy system on all the tables sizes. This system performed almost identically to the delta-rs-based system on smaller tables (10k and 100k rows), while outperformed it on larger tables (1M, 6M and 60M rows), with a latency up to seven times lower than its counterpart. In the experiments using multiple **CPU** cores, the delta-rs system experienced latency decreases up to 30%, while IcedHops latency decreased of maximum 5%. However, despite the better scalability results, the delta-rs system still underperform compared to IcedHops, for all the table sizes and **CPU** settings used in this thesis research. The above findings clearly demonstrate that IcedHops, in the industrial use case defined in Section 3.2.2, is a better alternative to legacy Spark-based system, while also confirming the need for Spark alternatives when operating on small tables (from 10K to 60M rows). Those alternatives could not only increase the performances, but also reduce computational costs, system complexity, and deployability, since managing a Spark cluster for small size data is expensive and more complex, and requires as well **JVM** dependencies, which cannot be deployed from any Python environment manager, such as pip, or conda.

Considering now the read experiments, the results show smaller performance differences, between the legacy system and IcedHops. When considering the largest table (60M rows), IcedHops exhibits only 55% latency reduction, however for all the other tables IcedHops has a latency from 15 to 56 times lower than the legacy system latency. The delta-rs pipeline follows a very similar trend, despite slightly underperforming IcedHops, with a maximum of 10% latency increase. The reason behind this different trend between IcedHops and legacy system's performance, compared to what seen for the write operation, is that the legacy system, despite it is already using an alternative to Spark, cannot efficiently scale down resources when reading small size tables. The Spark alternative employs a combination of Arrow Flight and DuckDB, which is instanciated within the Hopsworks cluster with a predefined amount of dedicated resources. Furthermore, in the experiments using multiple **CPU** cores, the delta-rs system scales better than IcedHops, experiencing latency decreases of respectively 55% and 39%, while the legacy

system is not impacted at all. The above outcomes shows that Hopsworks AB has still room for improvement even for the read operation, where employing IcedHops or delta-rs system would increase the system overall flexibility to different scales of data and to multiple **CPU** cores settings, at the price of moving computation on the client side.

All the above findings impact significantly Hopsworks AI, directly affecting the Hopsworks feature store, which is their main product. The results recommend the adoption of one of those two new systems (i.e. IcedHops, delta-rs), over the current system, using Apache Hudi and Spark to implement an offline feature store, or the extend this system with one or both the new systems. This would allow the users of the feature store, depending on their use case (data scale), their computational resources (**CPUs**), and their habits, to choose the preferred datalakehouse format to use. This would resonate with the Hopsworks AB objective of having products accessible from and integrable with multiple tools, and could be the starting point of possible further integration with currently developing technologies, such Apache XTable *. For a complete replacement of the legacy system across different use cases, further experiments should be conducted. The author envisions that beyond a certain data threshold, the Spark-based system may still outperform the new solutions due to its scalable distributed nature.

Perhaps the industrial use case described in Section 3.2.2 limits the results and findings to specific data scales, computational resources, and operations. Modifying one or more of these variables could lead to different outcomes, meaning the findings cannot be fully generalized. Additionally, since this thesis was conducted in collaboration with the company developing the evaluated product, some degree of bias is inevitable. Nonetheless, the results align with research expectations, reinforcing the idea that Spark alternatives should be considered in scenarios of small size (10K to 60M rows) data processing. Additionally, these findings highlight the need for further research and experimentation on Spark alternatives, to which thousands developers already contributes everyday.

5.2.2 Considerations on the legacy system

The legacy system, described in Sections 2.5.1-2.5.2, consists of a multi-layered architecture that employs Spark as compute engine and Hudi as **OTF**. This system is modeled to support multiple users, each having different use cases and resources. This could create several bottlenecks, depending on

*XTable repository available at <https://xtable.apache.org/>

the scenario, that goes far beyond the high latency of instantiating a Spark cluster. An example for this issue is using Kafka, which is used as single data entry-point in the Hopsworks feature store. Using Kafka as single entry-point ensures data consistency between the online feature store, working with streaming data, and the offline feature store, working on batch data. However, this creates a performance bottleneck for this infrastructure as data increases, needing to make a piece the full table in single rows, to save each of those as a single message. This was clearly visible in the experiments conducted in this thesis work, described in Figure 5.5 and Table 5.5, where Kafka accounted for more than 95% of the latency of the write operations. Now that these experiments outlined this problem, more research is needed to find an effective solution. Some approaches might reduce the overall latencies: (1) implement micro-batching mechanism, altogether with using columns-based format for data exchange; (2) enable client multi-processing, to support multiple settings with multiple CPU cores; (3) refine trade-off between partitions, topic replicas and brokers.

5.2.3 Considerations on the PyIceberg library

IcedHops, the system integration between Iceberg and HopsFS, described in Section 4.1, was possible thanks to the PyIceberg library, developed by the Iceberg open-source community for Python access to Iceberg Table. It was possible to create a Spark alternative system thanks to the high level of integrability of the library and the overarching OTF (compared to its counterparts, as highlighted in A.2), and thanks to several other online and book sources documenting this framework and library. The author's recommendation on this regard are: (1) the Iceberg and PyIceberg maintainers should revise their current documentation and integrate it with new findings, thus developers can directly access all the information needed from Iceberg webpage, without the need of buying a book, or of relying on documentation created by software companies that integrates Iceberg in their products; (2) intensively investigate solution for efficient local resources management, that might make Iceberg libraries better adaptable to changing local environment, as delta-rs shown to be able to, during the experiments described in Section 5.1.

Chapter 6

Conclusions and Future work

The conclusions of this thesis are presented in this Chapter. The RQs are revisited and discussed in relation to the experimental findings in Section 6.1, highlighting the contributions of this work. Section 6.2 details the limitations of the project, while Section 6.3 explores potential future research directions.

6.1 Conclusions

This thesis work tackled two RQs defined in Section 1.2.1. These were:

- RQ1: Is a working integration of PyIceberg and HopsFS in the Hopsworks offline feature store possible?
- RQ2: What are the differences in read and write latency and throughput on the Hopsworks offline feature store, between the existing legacy system and IcedHops?
- RQ3: What are the differences in read and write latency and throughput on the Hopsworks offline feature store, between IcedHops and the delta-rs alternative?

The work conducted in this thesis answered RQ1 by performing a system integration. Thus, RQ2 was answered by evaluating IcedHops against the legacy system, and RQ3 by evaluating IcedHops against the delta-rs alternative.

The first step was integrating HopsFS with the PyIceberg library. To complete this task, the technologies supported by HopsFS and PyIceberg were evaluated and the promising ones tested, according to the functional and non-functional requirements, until the best technology set was selected. Each

was tested in a dedicated script ^{*}, within the experimental environment. This contributed to the creation of the integration details, described in Section 4.1, where the reasoning and findings behind the author’s choice for SQLCatalog, DuckDB and PyArrowFileIO are explained. This can be used by developers to districte in similar complex system architectures, together with the thorough data lakehouse comparison tackled in the Section 2.2.3. Additionally, most of the significance of this work lies in the creation of a production-ready solution, which overcame system dependencies and use case limitations, while fulfilling all specified requirements. The integration of IcedHops into the production Hopsworks feature store shortly after the thesis’s release perhaps serves as a powerful testament to its practical value and real-world applicability.

The second and third steps were comparing the legacy system with IcedHops, and the latter with the delta-rs alternative. Experiments were designed and conducted to measure the performance of read and write operations on those system, according to two metrics, i.e., latency (in seconds) and throughput (in rows/second), accounting for 350 hours of testing. The results presented in Section 5.1 showed that IcedHops, accessing Iceberg tables on HopsFS, has a latency at least fifteen times lower than the legacy system, for both read and write operations, for the tables contained from 10K to 6M rows. The write experiments demonstrate that IcedHops significantly outperforms the legacy system in all cases, with write latency more than forty times lower across the board. This trend becomes even more striking for larger tables (6M and 60M rows), where IcedHops achieves write latencies that respectively are 140 and 110 times lower of those observed in the legacy system. These experiments were key to identify a critical bottleneck in current Hopsworks feature store architecture, where Kafka accounted for 95% of the write latency. In the other hand, the decrease of performance gap with the growth in size of the tables suggests that for even larger tables a Spark-based system might perform better, but more experiments are needed to verify this hypothesis. A separate comparison between IcedHops and delta-rs reveals that while IcedHops consistently achieves lower latencies, the extent of this advantage depends on table size. For the smallest table (10K rows), the two pipelines perform almost identically, whereas for the largest table (60M rows), ICedHops exhibits a performance edge of nearly seven times lower latency. Similarly, read performance experiments highlight considerable advantages of IcedHops over the legacy system. The latency reduction ranges from 55% for the largest table (60M rows) to as much as sixty times lower for the

^{*}Code available at: <https://github.com/SebastianoMeneghin/apache-iceberg-as-offline-feature-store>

100K rows table, with a consistent trend across all table sizes. This observed reduction in the performance gap is likely a result of the legacy system's use of a Spark alternative for reading, which already provides some improvements, but does not optimize resource usage as effectively as IcedHops, especially at smaller data scales. In contrast, a comparison between IcedHops and delta-rs shows that their read performance is nearly identical, except for the smallest table (10K rows), where IcedHops achieves three times lower latency. This indicates that while IcedHops is superior for write operations, its advantage diminishes in read operations, with delta-rs performing similarly. Experiments with increasing CPU cores reveal notable differences in scalability between the three systems. The legacy system shows minimal gains from additional CPU resources, while IcedHops significantly benefits from multi-core execution only on reading operations, reaching 39% performance gains with the larger tables (6M and 60M rows), and a modest 5% for writing operations. The delta-rs, however, scales significantly better, achieving 20-30% performance gains in write operations with the larger tables (6M and 60M rows), as well as latency reductions of 55% for the 1M rows table and 50% for the larger tables (6M and 60M rows) in read operations.

Overall, these findings reinforce the recommendation to adopt IcedHops in the defined use case, as it provides a substantial improvement over the legacy system, either as a replacement or as an extension for Iceberg table users. At the same time, the comparison between IcedHops and delta-rs suggests that IcedHops represents overall a better option for single-core machines, particularly for write operations, but delta-rs might be preferable in scenarios where multi-core scalability is crucial.

6.2 Limitations

The limitations of this study primarily stem from constraints in resources, both in terms of time and computational resources, and from in scope, which is closely tied to the defined industrial use case, described in Section 3.2.2.

This research focuses on reducing the read and write latency of the Hopsworks offline feature store, which dictated the choice of technology for implementation and experimentation. As a result, the findings are intrinsically influenced by the specific tools and infrastructure used, particularly given the collaboration with the company developing this technology, lowering their generality. Furthermore, the thesis' scope was shaped by the predefined use case, described in Section 3.2.2, which determined the CPU and data loads used in the experiments. While this helped establish a clear boundary for the

research contributions, it also limits the broader applicability of the results, needing further research to validate these findings in different contexts.

Despite access to substantial computational resources, their availability was restricted due the shared nature of the development environment, where other critical workloads had to be prioritized. Additionally, time constraints limited the number of experiments to reach a 95% confidence interval. All experiments were run fifty times, which significantly increased the time required to perform all experiments, reaching 350 hours.

6.3 Future work

The findings and limitations of this thesis open multiple avenues for future research. The scope of this work was primarily constrained by the specific industrial use case and the available computational resources. Relaxing one or more of these constraints could provide further insights and help generalize the findings across broader applications.

Future experiments should investigate the system's behavior with significantly larger datasets (1B rows, 10B rows, or more). The current study focused on tables up to 60M rows, but exploring a broader data scale could help determine whether there is a threshold where a Spark-based system begins to outperform the proposed alternatives due to its distributed nature. Similarly, additional experiments on different data formats, workloads, and storage solutions—such as [AWS S3](#) or [GCS](#)—could further validate the results. This thesis was conducted in collaboration with Hopsworks AB, focusing on optimizing the Hopsworks offline feature store. While the findings are valuable within the specific use case context, described in Section 3.2.2, testing the proposed approaches on other feature store implementations—such as Databricks or Feast—would help assess their broader applicability. As well, new libraries supporting the [OTF](#) evaluated in this thesis, like iceberg-rust, should be investigated, as they might further improve performances and scalability, better exploiting local resources and various computational settings.

Lastly, while this thesis focused on PyIceberg and delta-rs as alternatives to the legacy Spark-based system, other [OTFs](#) are currently under development, such as Apache Paimon, and performing similar experiments and comparisons with these technologies would broaden the knowledge on the possible alternatives to Spark.

References

- [1] S. A. Jackson, “State of the Data Lakehouse 2025,” <https://www.dremio.com/>. [Page 1.]
- [2] “Dremio Report Highlights Surge in Data Lakehouse Adoption for Enhanced Cost Efficiency and Analytics,” <https://www.bigdatawire.com/this-just-in/new-dremio-report-highlights-surge-in-data-lakehouse-adoption-for-enhanced-cost-efficiency-and-analytics/>, 2024. [Page 1.]
- [3] M. Armbrust, A. Ghodsi, R. Xin, and M. Zaharia, “Lakehouse: A New Generation of Open Platforms that Unify Data Warehousing and Advanced Analytics,” 2021. [Pages 1, 22, and 24.]
- [4] M. , “The state of AI in early 2024,” <https://www.mckinsey.com/capabilities/quantumblack/our-insights/the-state-of-ai>, May 2024. [Pages 1 and 34.]
- [5] D. Croci, “Data Lakehouse, beyond the hype,” Dec. 2022. [Page 1.]
- [6] “What Is a Lakehouse?” <https://www.databricks.com/blog/2020/01/30/what-is-a-data-lakehouse.html>. [Pages 1 and 4.]
- [7] “Apache Hudi vs Delta Lake vs Apache Iceberg - Data Lakehouse Feature Comparison,” <https://www.onehouse.ai/blog/apache-hudi-vs-delta-lake-vs-apache-iceberg-lakehouse-feature-comparison>. [Pages 1 and 2.]
- [8] P. Rajaperumal, “Uber Engineering’s Incremental Processing Framework on Hadoop,” <https://www.uber.com/en-EC/blog/hoodie/>, Mar. 2017. [Pages 1, 5, 21, 25, and 31.]
- [9] U. Team, “Iceberg Architecture Examples,” <https://www.upsolver.com/blog/iceberg-architecture-examples>, Feb. 2024. [Pages 1, 5, 13, 21, 25, and 27.]

- [10] M. Armbrust, T. Das, L. Sun, B. Yavuz, S. Zhu, M. Murthy, J. Torres, H. van Hovell, A. Ionescu, A. Łuszczak, M. Świtakowski, M. Szafrański, X. Li, T. Ueshin, M. Mokhtar, P. Boncz, A. Ghodsi, S. Paranjpye, P. Senster, R. Xin, and M. Zaharia, “Delta lake: High-performance ACID table storage over cloud object stores,” *Proc. VLDB Endow.*, vol. 13, no. 12, pp. 3411–3424, Aug. 2020. doi: 10.14778/3415478.3415560 [Pages 1, 5, 21, 25, and 28.]
- [11] T. Shiran, J. Hughes, and A. Merced, *Apache Iceberg: The Definitive Guide*. O’Reilly, Mar. 2024. ISBN 978-1-09-814861-4 [Pages 2, 7, 27, and 31.]
- [12] Dani, “Apache Iceberg: The Hadoop of the Modern Data Stack?” <https://blog.det.life/apache-iceberg-the-hadoop-of-the-modern-data-stack-c83f63a4ebb9>, Dec. 2024. [Pages 2, 27, and 31.]
- [13] L. Clark, “Big data vendors embrace Apache Iceberg,” https://www.theregister.com/2024/10/14/apache_iceberg_feature_announcements/. [Page 2.]
- [14] “The (Ongoing) Evolution Of Table Formats,” <https://www.montecarlodata.com/blog-the-evolution-of-table-formats/>, Sep. 2024. [Pages 2 and 7.]
- [15] M. Zaharia, R. S. Xin, P. Wendell, T. Das, M. Armbrust, A. Dave, X. Meng, J. Rosen, S. Venkataraman, M. J. Franklin, A. Ghodsi, J. Gonzalez, S. Shenker, and I. Stoica, “Apache Spark: A unified engine for big data processing,” *Commun. ACM*, vol. 59, no. 11, pp. 56–65, Oct. 2016. doi: 10.1145/2934664 [Pages 2 and 29.]
- [16] A. Khazanchi, “Faster reading with DuckDB and arrow flight on hopsworks : Benchmark and performance evaluation of offline feature stores,” Master’s thesis, KTH Royal Institute of Technology / KTH, School of Electrical Engineering and Computer Science (EECS) / KTH, School of Electrical Engineering and Computer Science (EECS), 2023. [Pages 2, 7, and 58.]
- [17] R. Vink, “I wrote one of the fastest DataFrame libraries - Ritchie Vink,” <https://www.ritchievink.com/blog/2021/02/28/i-wrote-one-of-the-fastest-dataframe-libraries/>. [Pages 2 and 5.]

- [18] M. Raasveldt and H. Mühleisen, “DuckDB: An Embeddable Analytical Database,” in *Proceedings of the 2019 International Conference on Management of Data*, ser. SIGMOD ’19. New York, NY, USA: Association for Computing Machinery, Jun. 2019. doi: 10.1145/3299869.3320212. ISBN 978-1-4503-5643-5 pp. 1981–1984. [Pages 2, 5, 32, 47, and 58.]
- [19] “Updates to the H2O.ai db-benchmark! – DuckDB,” <https://duckdb.org/2023/11/03/db-benchmark-update.html>. [Page 2.]
- [20] “DataFrames at Scale Comparison: TPC-H,” <https://docs.coiled.io/blog/tpch.html>. [Pages 2 and 30.]
- [21] “TIOBE Index,” <https://www.tiobe.com/tiobe-index/>. [Pages 2 and 6.]
- [22] “Stack Overflow Developer Survey 2023,” https://survey.stackoverflow.co/2023/?utm_source=social-share&utm_medium=social&utm_campaign=dev-survey-2023. [Pages 2 and 6.]
- [23] G. Staff, “Octoverse: AI leads Python to top language as the number of global developers surges,” Oct. 2024. [Pages 2 and 6.]
- [24] “Python Machine Learning | Data | Print,” <https://www.packtpub.com/en-us/product/python-machine-learning-9781789955750?type=print>. [Page 2.]
- [25] A. Pettersson, “Resource-efficient and fast Point-in-Time joins for Apache Spark : Optimization of time travel operations for the creation of machine learning training datasets,” Master’s thesis, KTH, School of Electrical Engineering and Computer Science (EECS) / KTH, School of Electrical Engineering and Computer Science (EECS), 2022. [Page 3.]
- [26] G. Manfredi, *Reducing Read and Write Latency in a Delta Lake-backed Offline Feature Store : Adding Support for HDFS and HopsFS in the Delta-Rs Rust Library to Reduce Read and Write Latency on the Delta Lake-backed Hopsworks Offline Feature Store*. KTH Royal Institute of Technology, 2024. [Pages 3, 7, 35, 36, 37, 39, 41, 47, 51, 53, 63, 64, 67, and 72.]
- [27] S. Chaudhuri and U. Dayal, “An overview of data warehousing and OLAP technology,” *ACM SIGMOD Record*, vol. 26, no. 1, pp. 65–74, Mar. 1997. doi: 10.1145/248603.248616 [Page 4.]

- [28] “The 80% Blind Spot: Are You Ignoring Unstructured Organizational Data?” <https://www.forbes.com/councils/forbestechcouncil/2019/01/29/the-80-blind-spot-are-you-ignoring-unstructured-organizational-data/>. [Pages 4 and 21.]
- [29] “Dremel made simple with Parquet,” https://blog.x.com/engineering/en_us/a/2013/dremel-made-simple-with-parquet. [Page 5.]
- [30] M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker, and I. Stoica, “Spark: Cluster Computing with Working Sets.” [Page 5.]
- [31] P. Carbone, A. Katsifodimos, S. Ewen, V. Markl, S. Haridi, and K. Tzoumas, “Apache Flink™: Stream and Batch Processing in a Single Engine.” [Page 5.]
- [32] G. van Rossum, “Python tutorial,” no. R 9526, Jan. 1995. [Page 6.]
- [33] J. de la Rúa Martínez, F. Buso, A. Kouzoupis, A. A. Ormenisan, S. Niazi, D. Bzhalava, K. Mak, V. Jouffrey, M. Ronström, R. Cunningham, R. Zangis, D. Mukhedkar, A. Khazanchi, V. Vlassov, and J. Dowling, “The Hopsworks Feature Store for Machine Learning,” in *Companion of the 2024 International Conference on Management of Data*, ser. SIGMOD/PODS ’24. New York, NY, USA: Association for Computing Machinery, Jun. 2024. doi: 10.1145/3626246.3653389. ISBN 9798400704222 pp. 135–147. [Pages 6, 18, and 37.]
- [34] S. Niazi, M. Ismail, S. Haridi, J. Dowling, S. Grohsschmiedt, and M. Ronström, “HopsFS: Scaling Hierarchical File System Metadata Using NewSQL Databases.” [Pages 7, 11, and 18.]
- [35] “The Apache Spark Open Source Project on Open Hub,” <https://openhub.net/p/apache-spark>. [Page 8.]
- [36] “What is Green Software?” <https://greensoftware.foundation/articles/what-is-green-software>, Oct. 2021. [Page 10.]
- [37] D. Patterson, J. Gonzalez, Q. Le, C. Liang, L.-M. Munguia, D. Rothchild, D. So, M. Texier, and J. Dean, “Carbon Emissions and Large Neural Network Training,” Apr. 2021. [Page 10.]

- [38] D. Patterson, J. Gonzalez, U. Hözle, Q. Le, C. Liang, L.-M. Munguia, D. Rothchild, D. R. So, M. Texier, and J. Dean, “The Carbon Footprint of Machine Learning Training Will Plateau, Then Shrink,” *Computer*, vol. 55, no. 7, pp. 18–28, Jul. 2022. doi: 10.1109/MC.2022.3148714 [Page 10.]
- [39] “(PDF) Big Data Processing Stacks,” *ResearchGate*, Oct. 2024. doi: 10.1109/MITP.2017.6 [Pages 13 and 34.]
- [40] M. Frampton, *Complete Guide to Open Source Big Data Stack*. Berkeley, CA: Apress, 2018. ISBN 978-1-4842-2148-8 978-1-4842-2149-5 [Page 13.]
- [41] S. Sakr, “Big Data Processing Stacks,” *IT Professional*, vol. 19, no. 1, pp. 34–41, Jan. 2017. doi: 10.1109/MITP.2017.6 [Page 13.]
- [42] “Block vs File vs Object Storage - Difference Between Data Storage Services - AWS,” <https://aws.amazon.com/compare/the-difference-between-block-file-object-storage/>. [Pages 14, 15, and 17.]
- [43] “How Object vs Block vs File Storage differ,” <https://cloud.google.com/discover/object-vs-block-vs-file-storage>. [Pages 14, 15, and 17.]
- [44] “Object vs. File vs. Block Storage: What’s the Difference? | IBM,” <https://www.ibm.com/think/topics/object-vs-file-vs-block-storage>, Aug. 2024. [Pages 14, 15, and 17.]
- [45] A. Kala Karun and K. Chitharanjan, “A review on hadoop — HDFS infrastructure extensions,” in *2013 IEEE Conference on Information & Communication Technologies*, Apr. 2013. doi: 10.1109/CICT.2013.6558077 pp. 132–137. [Page 17.]
- [46] J. Wu, L. Ping, X. Ge, Y. Wang, and J. Fu, “Cloud Storage as the Infrastructure of Cloud Computing,” in *2010 International Conference on Intelligent Computing and Cognitive Informatics*, Jun. 2010. doi: 10.1109/ICICCI.2010.119 pp. 380–383. [Page 20.]
- [47] M. Ismail, S. Niazi, G. Berthou, M. Ronström, S. Haridi, and J. Dowling, “HopsFS-S3: Extending Object Stores with POSIX-like Semantics and more (industry track),” in *Proceedings of the 1st International Middleware Conference Industrial Track*. Delft Netherlands: ACM,

- Dec. 2020. doi: 10.1145/3429357.3430521. ISBN 978-1-4503-8201-4 pp. 23–30. [Page 20.]
- [48] H. E. Pence, “What is Big Data and Why is it Important?” *Journal of Educational Technology Systems*, vol. 43, no. 2, pp. 159–171, Dec. 2014. doi: 10.2190/ET.43.2.d [Page 21.]
- [49] Y. Demchenko, Z. Zhao, P. Grossi, A. Wibisono, and C. de Laat, “Addressing Big Data challenges for Scientific Data Infrastructure,” in *4th IEEE International Conference on Cloud Computing Technology and Science Proceedings*, Dec. 2012. doi: 10.1109/CloudCom.2012.6427494 pp. 614–617. [Page 21.]
- [50] D. Mazumdar, J. Hughes, and J. B. Onofre, “The Data Lakehouse: Data Warehousing and More,” Oct. 2023. [Pages 22 and 23.]
- [51] “Data Lakehouse: A survey and experimental study,” *Information Systems*, vol. 127, p. 102460, Jan. 2025. doi: 10.1016/j.is.2024.102460 [Page 24.]
- [52] P. Jain, P. Kraft, C. Power, T. Das, I. Stoica, and M. Zaharia, “Analyzing and Comparing Lakehouse Storage Systems,” 2023. [Pages 25 and 28.]
- [53] “Apache Hudi | Technical Documentation,” <https://hudi.apache.org/docs/overview>. [Pages 25 and 91.]
- [54] O. Katz, “Hudi vs Iceberg vs Delta Lake: Data Lake Table Formats Compared,” Apr. 2021. [Pages 26 and 27.]
- [55] “Apache Iceberg vs Delta Lake,” <https://www.starburst.io/blog/iceberg-vs-delta-lake/>. [Page 26.]
- [56] “Apache Iceberg | Technical Documentation,” <https://iceberg.apache.org/docs/nightly/>. [Pages 27, 56, 58, and 91.]
- [57] “PyIceberg,” <https://py.iceberg.apache.org/>. [Page 28.]
- [58] “Delta Lake | Technical Documentation,” <https://docs.delta.io/latest/index.html>. [Pages 28 and 91.]
- [59] S. Asad, “Apache Iceberg vs Delta Lake vs Apache Hudi: A Comprehensive Comparison,” Jan. 2025. [Page 28.]

- [60] “Announcing Delta Lake 3.0 with New Universal Format and Liquid Clustering,” <https://www.databricks.com/blog/announcing-delta-lake-30-new-universal-format-and-liquid-clustering>, Wed, 06/28/2023 - 23:53. [Pages 29 and 32.]
- [61] J. Dean and S. Ghemawat, “MapReduce: Simplified data processing on large clusters,” *Commun. ACM*, vol. 51, no. 1, pp. 107–113, Jan. 2008. doi: 10.1145/1327452.1327492 [Page 29.]
- [62] D. Borthakur, “The Hadoop Distributed File System: Architecture and Design,” 2007. [Page 29.]
- [63] M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauley, M. J. Franklin, S. Shenker, and I. Stoica, “Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing.” [Page 29.]
- [64] J. Kreps, N. Narkhede, and J. Rao, “Kafka: A Distributed Messaging System for Log Processing.” [Page 30.]
- [65] R. Blue, “Catalogs and the REST catalog,” 2024. [Pages 31 and 32.]
- [66] wesm, “Introducing Apache Arrow Flight: A Framework for Fast Data Transport,” <https://arrow.apache.org/blog/2019/10/13/introducing-arrow-flight/>, Oct. 2019. [Page 33.]
- [67] “Hopsworks - The Real-time AI Lakehouse,” <https://www.hopsworks.ai/>. [Page 34.]
- [68] “Meet Michelangelo: Uber’s Machine Learning Platform,” <https://www.uber.com/en-IN/blog/michelangelo-machine-learning-platform/>, Sep. 2017. [Page 34.]
- [69] “The Big Dictionary of MLOps,” <https://www.hopsworks.ai/mlops-dictionary>, 2024. [Page 34.]
- [70] “TPC-H Benchmark: Decision support, standard verification, revision v3.0.1,” https://www.tpc.org/tpc_documents_current_versions/current_specifications5.asp, 2022. [Page 47.]
- [71] A. Behm, S. Palkar, U. Agarwal, T. Armstrong, D. Cashman, A. Dave, T. Greenstein, S. Hovsepian, R. Johnson, A. Sai Krishnan, P. Leventis,

- A. Luszczak, P. Menon, M. Mokhtar, G. Pang, S. Paranjpye, G. Rahn, B. Samwel, T. Van Bussel, H. Van Hovell, M. Xue, R. Xin, and M. Zaharia, “Photon: A Fast Query Engine for Lakehouse Systems,” in *Proceedings of the 2022 International Conference on Management of Data*. Philadelphia PA USA: ACM, Jun. 2022. doi: 10.1145/3514221.3526054. ISBN 978-1-4503-9249-5 pp. 2326–2339. [Page 47.]
- [72] B. Inmon, “Five Steps to a successful data lakehouse,” vol. 2021. [Page 90.]

Appendix A

Technology comparisons

This appendix contains the comparative tables for the technologies presented in Section 2.2.1–2.2.3, respectively between the different DBMs developed along their history, and the different OTFs presented in this thesis report.

Table A.1: Comparison of key features of Data Warehouses, Data Lakes, and Data Lakehouses [72]. Technologies presented and discussed in details in Section 2.2.1.

Feature	Data Warehouse	Data Lake	Data Lakehouse
<i>Primary Use</i>	Business Intelligence, SQL Analytics	Unstructured data storage, AI/ML workflows	Unified storage for structured/unstructured data, SQL + AI/ML processing
<i>Data Type</i>	Structured	Unstructured	All data types
<i>Query Performance</i>	Optimized for SQL -based queries, fast indexed access	Slow for structured queries, requires extensive data preparation	High performance for both SQL and ML workloads
<i>Schema Enforcement</i>	Strict schema-on-write	Schema-on-read	Flexible schema with enforcement and evolution support
ACID support	Fully supported	Not supported	Fully supported
<i>Scalability</i>	Limited scalability due to high storage costs	High scalability, low-cost storage	Highly scalable with structured query optimizations
<i>Storage Cost</i>	High due to proprietary formats	Low due to object storage	Moderate, optimized for efficiency
<i>Governance</i>	Strong governance, access control, and security	Weak governance, risk of data swamps	Fine-grained governance with access control and security
<i>Key Takeaways</i>	Best for structured data and BI queries. High-performance SQL analytics.	Ideal for cost-effective big data storage. Difficult to manage and query efficiently	Combines the best of warehouses and lakes. Supports both SQL and AI/ML use cases efficiently.

Table A.2: Comparison of key features of Apache Hudi, Apache Iceberg, Delta Lake [53, 56, 58]. Technologies presented and discussed in details in Section 2.2.3.

Feature	Apache Hudi	Apache Iceberg	Delta Lake
<i>Metadata Implementation</i>	Tabular	Hierarchical	Tabular
<i>Time Travel</i>	snapshots	transaction log	incremental commits
<i>Caching</i>	No	Yes	Yes
<i>Optimization</i>	CoW, MoR	CoW, MoR	CoW
<i>Schema Evolution</i>	Limited	Full	Partial
<i>Partition Evolution</i>	Explicit	Hidden	Explicit
<i>Concurrency Control</i>	Optimistic, Multi-version	Optimistic	Optimistic
<i>Storage Engines</i>	HDFS, AWS S3, GCS , Azure Blob Storage		
<i>File Formats</i>	Parquet Avro ORC	Parquet Avro ORC	Parquet
<i>Catalogs</i>	Hive AWS Glue	Hive AWS Glue JDBC REST	Hive AWS Glue Unity
<i>Query Engines</i>	Spark Flink Presto Hive Impala	Spark Flink Presto Athena Trino Snowflake	Spark Presto Athena Redshift Snowflake
<i>APIs</i>	Java (Spark) Python (Spark) Rust Scala (Spark,Flink)	Java (Spark) Python Rust	Java (Spark,Flink) Python (Spark) Rust Scala (Spark)

Appendix B

System architectures

This appendix reports the legacy architecture diagrams, described in Section [2.5](#), here increased in size to improve readability.

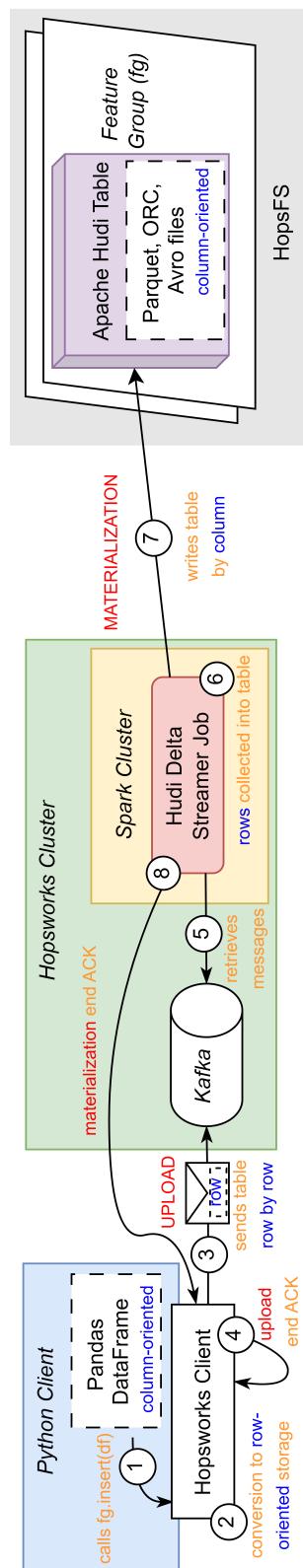


Figure B.1: Legacy system writing a Pandas DataFrame from a Python client to the Hopsworks offline feature store. This image was magnified to enhance visualization.

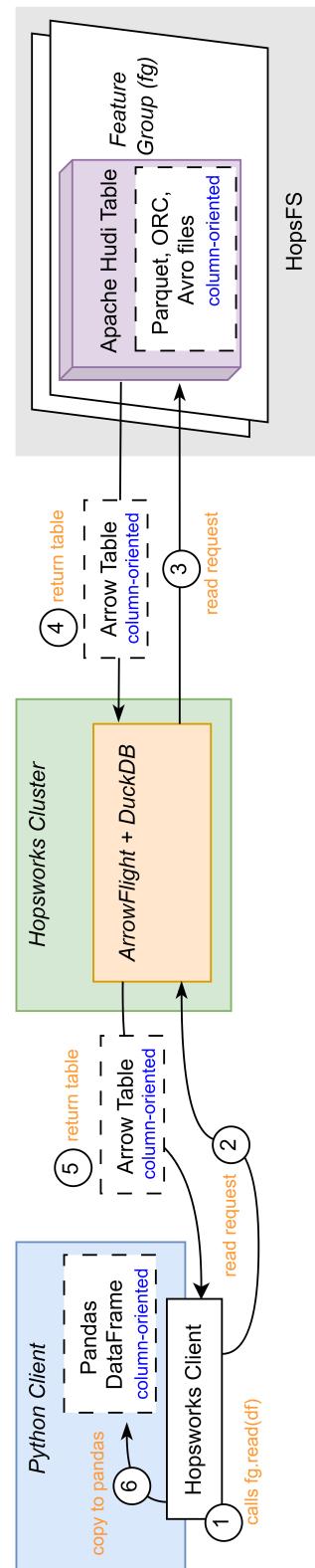


Figure B.2: Legacy system reading a table from the Hopsworks offline feature store and loading it into the Python client's local memory. This image was magnified to enhance visualization.

Appendix C

Write experiments results

This appendix contains all the graphs and the tables related to the write experiments conducted, and of the comparison with the related work. Results are reported first expressed as latency (measured in seconds) and then as throughput (measured in rows/seconds). Latency was directly measured, while throughput was computed from the latency and table size.

Table C.1: Write experiment results expressed as latency. The experiment was performed with one **CPU** core.

Pipeline	Number of rows	Latency (seconds)	Latency (seconds) 95% Confidence Interval	
			low	high
Hudi Legacy	10K	50.237 94	49.497 73	50.959 89
	100K	59.556 05	58.911 90	60.197 33
	1M	112.177 73	111.398 58	112.961 77
	6M	511.849 08	510.652 80	512.889 00
	60M	2 716.209 39	2 700.807 03	2 732.348 88
Iceberg IcedHops	10K	1.204 06	1.169 30	1.234 95
	100K	1.265 24	1.221 50	1.299 14
	1M	1.717 70	1.662 75	1.792 99
	6M	3.576 49	3.535 57	3.614 05
	60M	24.609 89	24.453 04	24.785 83
Delta Lake delta-rs	10K	1.250 88	1.238 22	1.264 26
	100K	1.368 00	1.337 73	1.389 56
	1M	9.382 31	9.240 54	9.541 70
	6M	19.751 49	19.386 52	20.111 50
	60M	177.194 58	174.583 94	180.046 79

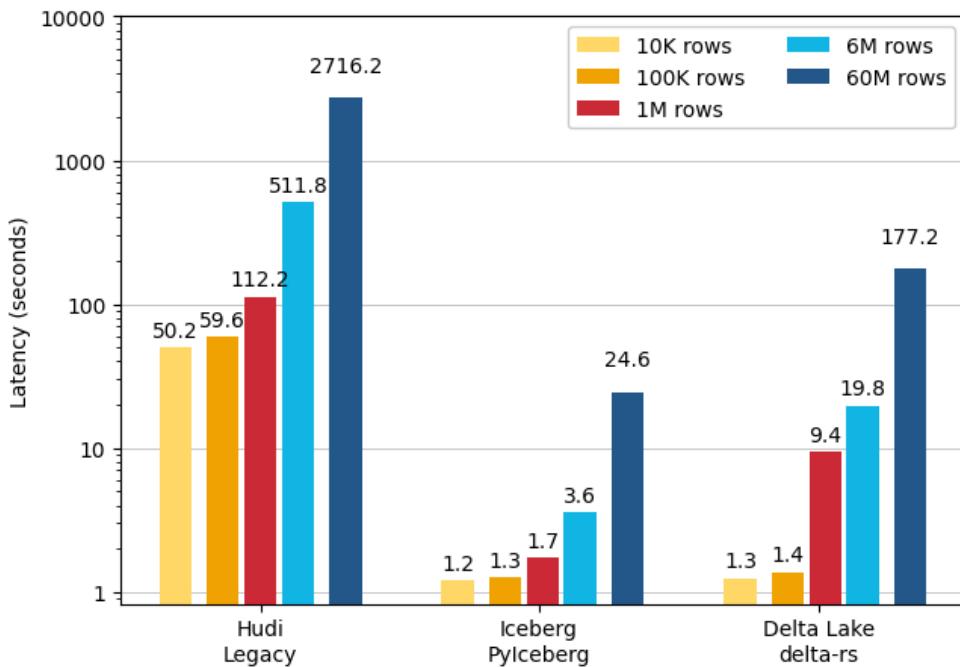


Figure C.1: Histogram in log-scale of the write experiment results expressed as latency. The experiment was performed with one **CPU** core.

Table C.2: Write experiment results expressed as latency. The experiment was performed with two **CPU** cores.

Pipeline	Number of rows	Latency (seconds)	Latency (seconds)	
			95% Confidence Interval low	high
Hudi Legacy	10K	50.72779	50.11508	51.31597
	100K	59.77647	59.10732	60.48278
	1M	108.56462	108.01696	109.06968
	6M	473.38310	472.30635	474.45390
	60M	2341.08359	2334.15850	2347.68142
Iceberg IcedHops	10K	1.19507	1.13429	1.26015
	100K	1.19962	1.14105	1.25251
	1M	1.72148	1.69402	1.74827
	6M	3.56635	3.50274	3.64406
	60M	23.50784	23.46474	23.54621
Delta Lake delta-rs	10K	1.26280	1.25128	1.27656
	100K	1.30721	1.27933	1.32999
	1M	8.52018	8.34308	8.70174
	6M	16.30152	15.92007	16.66990
	60M	134.15567	131.76812	136.58794

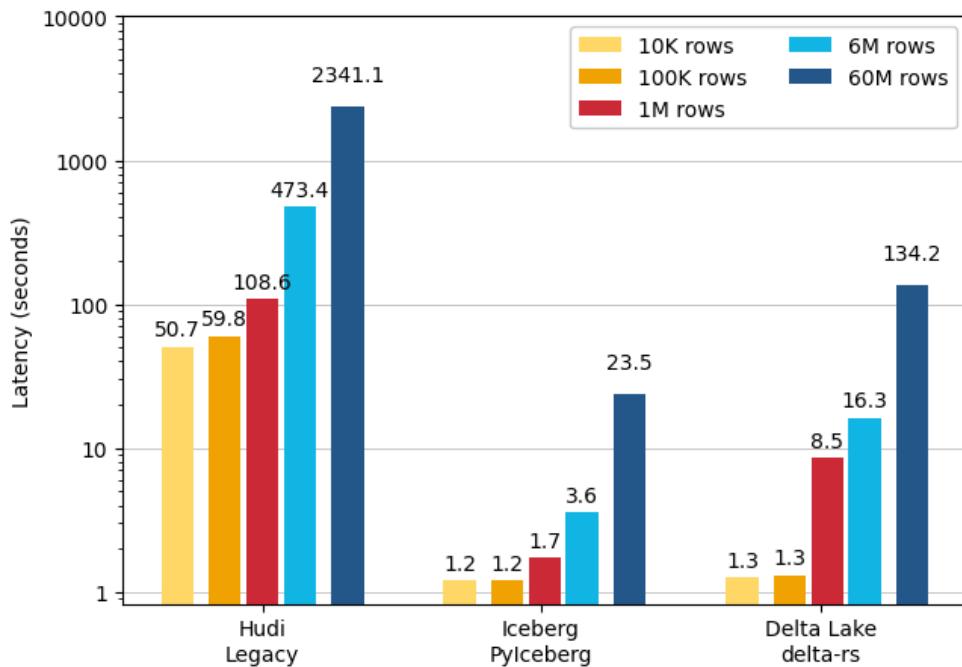


Figure C.2: Histogram in log-scale of the write experiment results expressed as latency. The experiment was performed with two **CPU** cores.

Table C.3: Write experiment results expressed as latency. The experiment was performed with four **CPU** cores.

Pipeline	Number of rows	Latency (seconds)	Latency (seconds) 95% Confidence Interval	
			low	high
Hudi Legacy	10K	51.271 46	50.609 85	51.917 87
	100K	59.514 01	58.855 12	60.137 37
	1M	108.807 23	108.291 44	109.369 79
	6M	481.958 70	480.984 72	482.900 19
	60M	2 346.108 73	2 337.175 80	2 355.212 33
Iceberg IcedHops	10K	1.195 28	1.161 00	1.226 59
	100K	1.254 75	1.218 73	1.285 07
	1M	1.631 61	1.586 39	1.672 95
	6M	3.544 88	3.514 41	3.571 31
	60M	24.005 35	23.897 06	24.153 70
Delta Lake delta-rs	10K	1.216 84	1.202 38	1.233 06
	100K	1.336 53	1.323 13	1.350 43
	1M	8.418 63	8.245 95	8.603 09
	6M	16.223 45	15.831 16	16.585 56
	60M	124.061 39	121.162 37	126.547 71

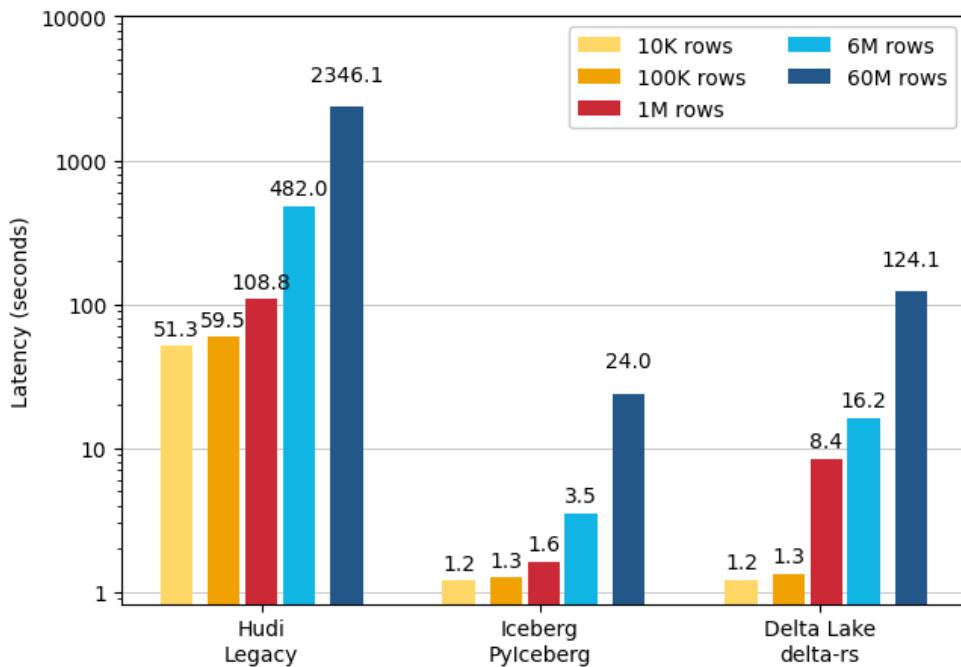


Figure C.3: Histogram in log-scale of the write experiment results expressed as latency. The experiment was performed with four **CPU** cores.

Table C.4: Write experiment results expressed as latency. The experiment was performed with eight **CPU** cores.

Pipeline	Number of rows	Latency (seconds)	Latency (seconds)	
			95% Confidence Interval low	high
Hudi Legacy	10K	51.225 98	50.587 30	51.826 30
	100K	60.288 08	59.778 63	60.766 77
	1M	109.381 79	108.911 21	109.852 72
	6M	475.972 84	474.807 88	477.121 64
	60M	2 324.807 40	2 319.160 36	2 331.084 60
Iceberg IcedHops	10K	1.218 22	1.159 14	1.290 76
	100K	1.248 58	1.219 85	1.276 29
	1M	1.719 08	1.704 10	1.731 50
	6M	3.532 48	3.473 37	3.609 68
	60M	23.623 13	23.560 93	23.700 51
Delta Lake delta-rs	10K	1.359 65	1.248 49	1.570 41
	100K	1.292 20	1.265 57	1.311 20
	1M	8.294 98	8.143 43	8.460 28
	6M	15.724 10	15.239 56	16.176 34
	60M	121.892 14	119.518 99	124.134 25

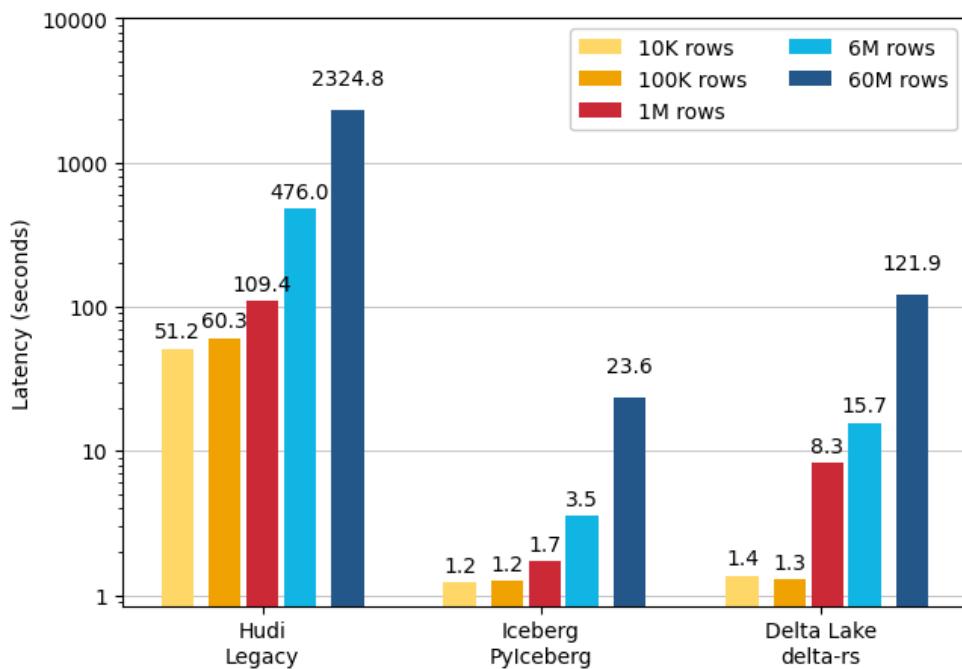


Figure C.4: Histogram in log-scale of the write experiment results expressed as latency. The experiment was performed with eight **CPU** cores.

Table C.5: Write experiment results expressed as throughput. The experiment was performed with one **CPU** core.

Pipeline	Number of rows	Throughput (k rows/second)	Throughput (k rows/second) 95% Confidence Interval	
			low	high
Hudi Legacy	10K	0.199 05	0.196 23	0.202 03
	100K	1.679 09	1.661 20	1.697 45
	1M	8.914 43	8.852 55	8.976 78
	6M	11.722 21	11.698 44	11.749 67
	60M	22.089 61	21.959 13	22.215 58
Iceberg IcedHops	10K	8.305 21	8.097 48	8.552 14
	100K	79.036 58	76.974 09	81.866 64
	1M	582.172 95	557.726 36	601.411 81
	6M	1 677.621 61	1 660.186 13	1 697.038 46
	60M	2 438.044 52	2 420.738 21	2 453.683 12
Delta Lake delta-rs	10K	7.994 35	7.909 77	8.076 11
	100K	73.099 24	71.965 38	74.753 63
	1M	106.583 58	104.803 17	108.218 78
	6M	303.774 53	298.336 72	309.493 33
	60M	338.610 80	333.246 70	343.674 23

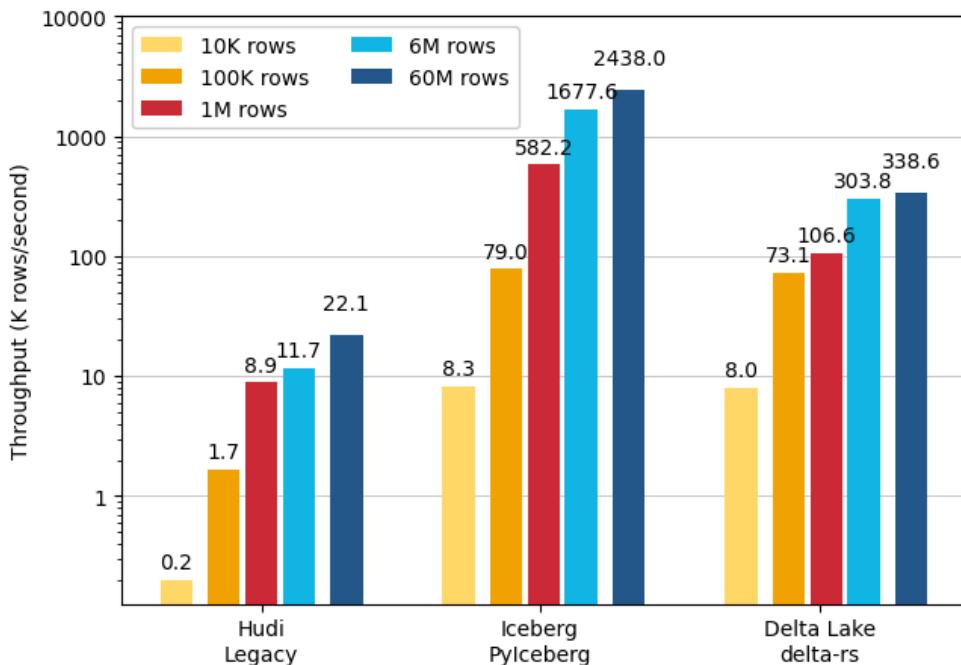


Figure C.5: Histogram in log-scale of the write experiment results expressed as throughput. The experiment was performed with one **CPU** core.

Table C.6: Write experiment results expressed as throughput. The experiment was performed with two **CPU** cores.

Pipeline	Number of rows	Throughput (k rows/second)	Throughput (k rows/second) 95% Confidence Interval	
			low	high
Hudi Legacy	10K	0.19713	0.19487	0.19954
	100K	1.67290	1.65336	1.69184
	1M	9.21110	9.16845	9.25781
	6M	12.67472	12.64612	12.70362
	60M	25.62916	25.55713	25.70520
Iceberg IcedHops	10K	8.36770	7.93556	8.81610
	100K	83.36001	79.83944	87.63880
	1M	580.89514	571.99564	590.31126
	6M	1682.39446	1646.51435	1712.94234
	60M	2552.34016	2548.18101	2557.02817
Delta Lake delta-rs	10K	7.91888	7.83356	7.99183
	100K	76.49869	75.18833	78.16603
	1M	117.36839	114.91950	119.85982
	6M	368.06392	359.93015	376.88277
	60M	447.24164	439.27743	455.34535

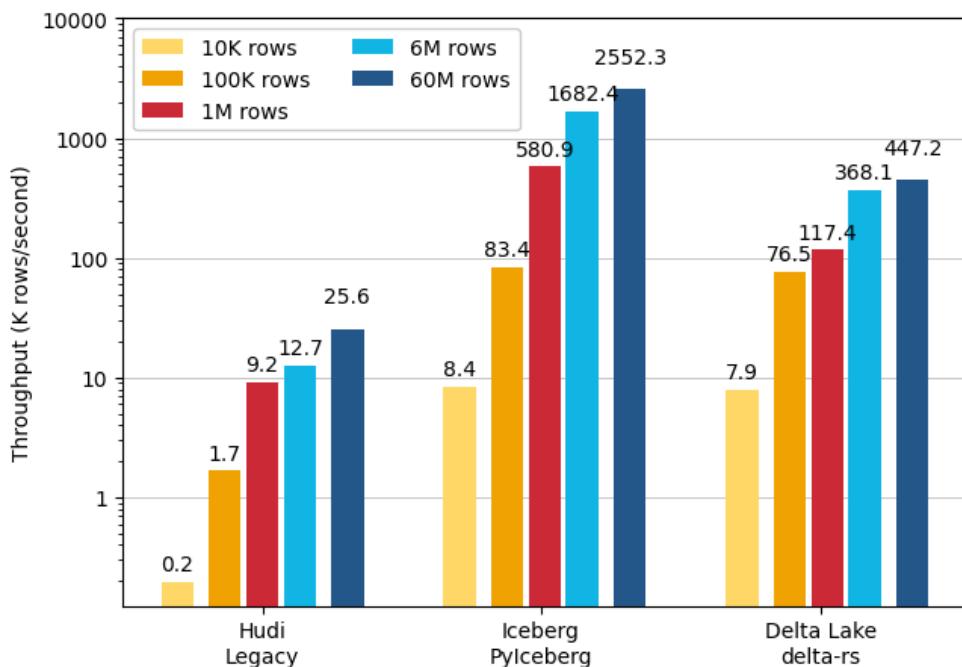


Figure C.6: Histogram in log-scale of the write experiment results expressed as throughput. The experiment was performed with two **CPU** cores.

Table C.7: Write experiment results expressed as throughput. The experiment was performed with four **CPU** cores.

Pipeline	Number of rows	Throughput (k rows/second)	Throughput (k rows/second) 95% Confidence Interval	
			low	high
Hudi Legacy	10K	0.195 04	0.192 61	0.197 58
	100K	1.680 28	1.662 86	1.699 09
	1M	9.190 57	9.143 29	9.234 34
	6M	12.449 20	12.424 93	12.474 41
	60M	25.574 26	25.475 41	25.672 01
Iceberg IcedHops	10K	8.366 27	8.152 67	8.613 26
	100K	79.697 15	77.816 50	82.052 81
	1M	612.891 38	597.747 52	630.363 05
	6M	1 692.583 43	1 680.055 53	1 707.255 44
	60M	2 499.443 20	2 484.091 84	2 510.768 77
Delta Lake delta-rs	10K	8.218 03	8.109 89	8.316 84
	100K	74.820 47	74.050 66	75.578 59
	1M	118.784 11	116.237 36	121.271 61
	6M	369.835 13	361.760 41	378.999 43
	60M	483.631 53	474.129 47	495.203 24

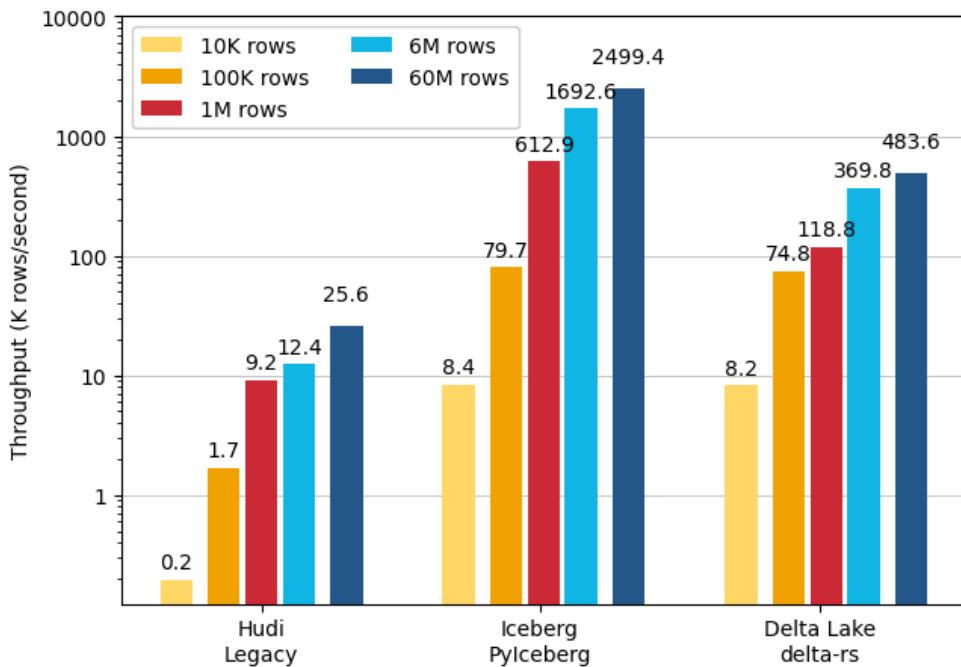


Figure C.7: Histogram in log-scale of the write experiment results expressed as throughput. The experiment was performed with four **CPU** cores.

Table C.8: Write experiment results expressed as throughput. The experiment was performed with eight **CPU** cores.

Pipeline	Number of rows	Throughput (k rows/second)	Throughput (k rows/second) 95% Confidence Interval	
			low	high
Hudi Legacy	10K	0.195 21	0.192 95	0.197 68
	100K	1.658 70	1.645 64	1.672 84
	1M	9.142 29	9.103 10	9.181 79
	6M	12.605 76	12.575 41	12.636 69
	60M	25.808 59	25.739 09	25.871 43
Iceberg IcedHops	10K	8.208 73	7.747 35	8.627 10
	100K	80.090 89	78.352 19	81.977 34
	1M	581.707 82	577.533 99	586.819 41
	6M	1 698.521 31	1 662.196 40	1 727.430 78
	60M	2 539.884 07	2 531.591 50	2 546.589 00
Delta Lake delta-rs	10K	7.354 81	6.367 77	8.009 70
	100K	77.387 41	76.265 73	79.015 66
	1M	120.554 78	118.199 43	122.798 45
	6M	381.579 81	370.912 17	393.712 18
	60M	492.238 46	483.347 68	502.012 28

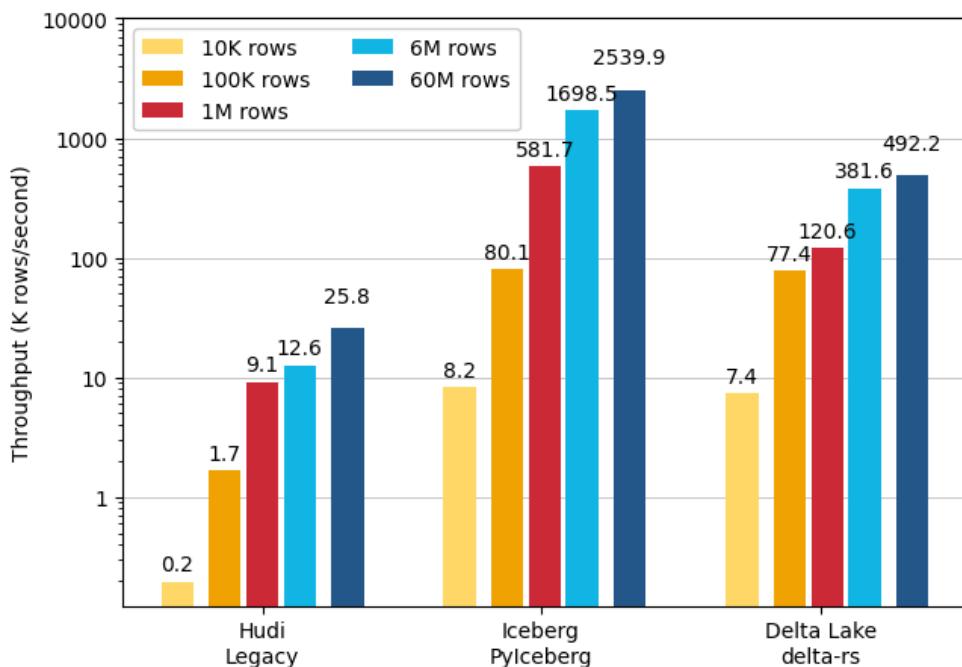


Figure C.8: Histogram in log-scale of the write experiment results expressed as throughput. The experiment was performed with eight **CPU** cores.

Appendix D

Read experiments results

This appendix contains all the graphs and the tables related to the read experiments conducted and the comparison with the related work. Results are first reported expressed as latency (measured in seconds) and then as throughput (measured in rows/seconds). Latency was directly measured, while throughput was computed from the latency and table size.

Table D.1: Read experiment results expressed as latency. The experiment was performed with one **CPU** core.

Pipeline	Number of rows	Latency (seconds)	Latency (seconds) 95% Confidence Interval	
			low	high
Hudi Legacy	10K	0.63144	0.62360	0.64086
	100K	2.65043	2.64277	2.65880
	1M	8.59296	8.32432	8.88965
	6M	33.52580	33.22121	33.83881
	60M	33.69031	33.34263	34.03532
Iceberg IcedHops	10K	0.01723	0.01646	0.01799
	100K	0.04687	0.04283	0.05107
	1M	0.43018	0.41987	0.44190
	6M	2.12331	2.10768	2.13933
	60M	21.81955	21.69884	21.93705
Delta Lake delta-rs	10K	0.05345	0.03917	0.08056
	100K	0.05763	0.05505	0.06055
	1M	0.53878	0.52444	0.55318
	6M	1.94947	1.92855	1.96958
	60M	22.98186	22.83994	23.15794

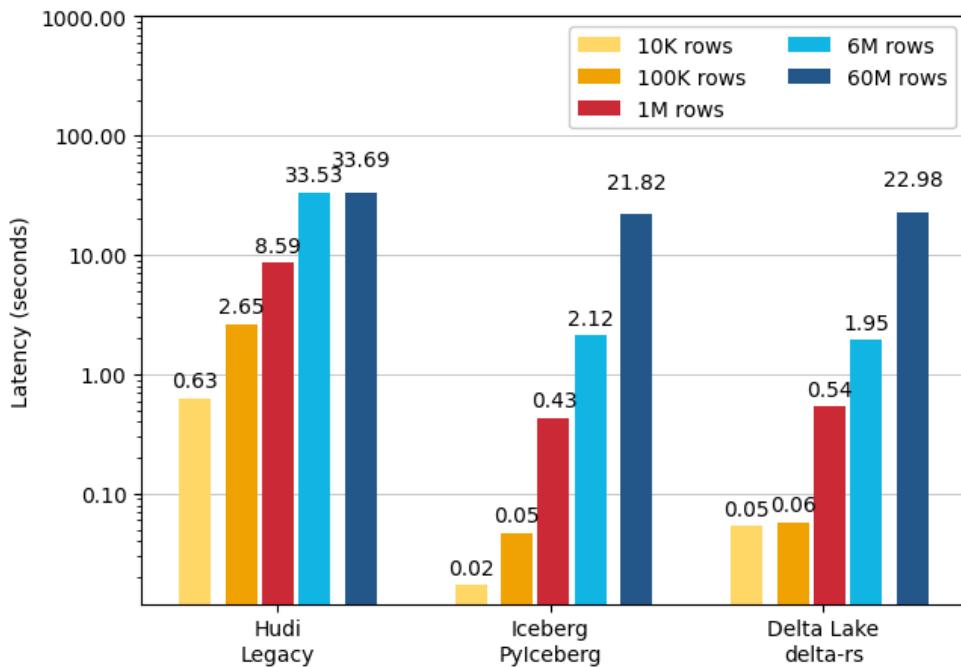


Figure D.1: Histogram in log-scale of the read experiment results expressed as latency. The experiment was performed with one **CPU** core.

Table D.2: Read experiment results expressed as latency. The experiment was performed with two **CPU** cores.

Pipeline	Number of rows	Latency (seconds)	Latency (seconds) 95% Confidence Interval	
			low	high
Hudi Legacy	10K	0.624 82	0.621 81	0.628 02
	100K	2.663 32	2.656 15	2.670 70
	1M	8.606 87	8.301 38	8.979 21
	6M	33.377 40	33.096 84	33.706 81
	60M	33.628 80	33.286 17	34.017 80
Iceberg IcedHops	10K	0.018 32	0.017 35	0.019 31
	100K	0.043 17	0.041 43	0.045 18
	1M	0.269 43	0.262 17	0.276 64
	6M	1.296 06	1.230 03	1.343 85
	60M	13.512 94	13.410 47	13.624 77
Delta Lake delta-rs	10K	0.041 21	0.039 31	0.043 85
	100K	0.057 21	0.051 10	0.066 94
	1M	0.233 73	0.225 13	0.242 63
	6M	0.908 26	0.899 00	0.918 35
	60M	11.409 75	11.273 01	11.596 08

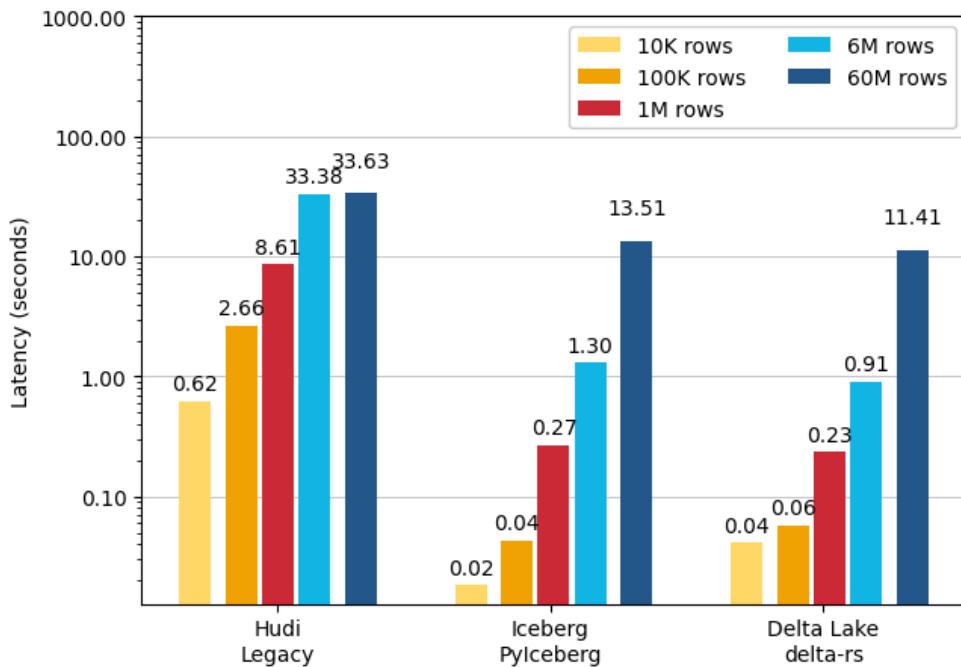


Figure D.2: Histogram in log-scale of the read experiment results expressed as latency. The experiment was performed with two **CPU** cores.

Table D.3: Read experiment results expressed as latency. The experiment was performed with four **CPU** cores.

Pipeline	Number of rows	Latency (seconds)	Latency (seconds) 95% Confidence Interval	
			low	high
Hudi Legacy	10K	0.636 24	0.623 55	0.659 45
	100K	2.639 86	2.633 50	2.646 47
	1M	8.744 36	8.522 35	9.000 47
	6M	33.447 49	33.185 73	33.740 60
	60M	33.661 53	33.263 12	34.090 78
Iceberg IcedHops	10K	0.029 79	0.015 68	0.059 14
	100K	0.042 83	0.041 13	0.044 74
	1M	0.233 37	0.218 55	0.248 06
	6M	1.008 57	1.001 23	1.016 16
	60M	10.110 29	10.040 90	10.183 32
Delta Lake delta-rs	10K	0.043 57	0.039 31	0.051 02
	100K	0.055 46	0.053 72	0.057 86
	1M	0.189 58	0.150 76	0.255 24
	6M	0.531 06	0.507 52	0.569 20
	60M	5.579 90	5.545 08	5.618 30

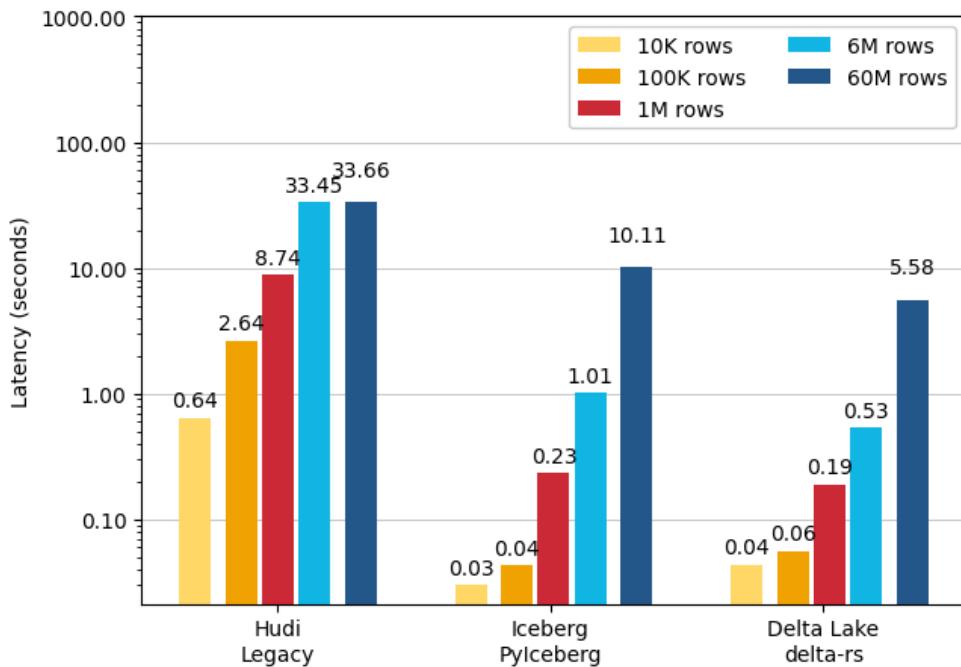


Figure D.3: Histogram in log-scale of the read experiment results expressed as latency. The experiment was performed with four **CPU** cores.

Table D.4: Read experiment results expressed as latency. The experiment was performed with eight **CPU** cores.

Pipeline	Number of rows	Latency (seconds)	Latency (seconds) 95% Confidence Interval	
			low	high
Hudi Legacy	10K	0.627 33	0.622 54	0.632 50
	100K	2.662 35	2.653 22	2.672 16
	1M	8.352 14	8.141 63	8.618 32
	6M	33.440 13	33.169 78	33.752 54
	60M	33.142 44	32.900 73	33.405 44
Iceberg IcedHops	10K	0.016 96	0.015 98	0.018 10
	100K	0.042 68	0.040 96	0.044 78
	1M	0.232 12	0.229 27	0.235 31
	6M	0.884 65	0.875 63	0.894 54
	60M	8.839 20	8.758 94	8.933 90
Delta Lake delta-rs	10K	0.043 14	0.038 83	0.051 29
	100K	0.054 65	0.052 99	0.056 65
	1M	0.174 09	0.169 83	0.178 29
	6M	0.497 07	0.485 62	0.510 95
	60M	2.944 31	2.859 72	3.053 01

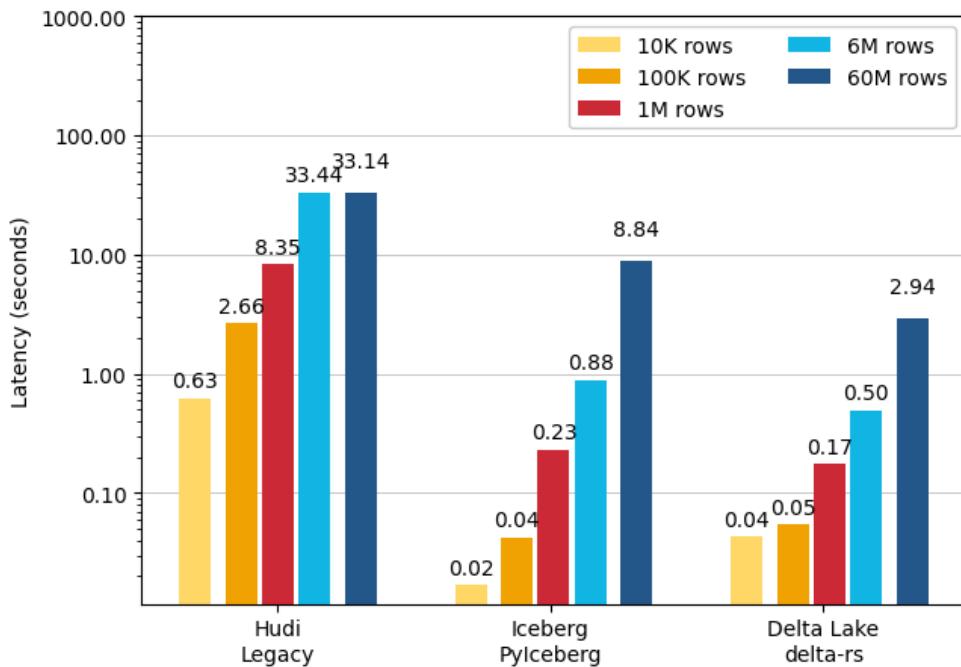


Figure D.4: Histogram in log-scale of the read experiment results expressed as latency. The experiment was performed with eight **CPU** cores.

Table D.5: Read experiment results expressed as throughput. The experiment was performed with one **CPU** core.

Pipeline	Number of rows	Throughput (k rows/second)	Throughput (k rows/second) 95% Confidence Interval	
			low	high
Hudi Legacy	10K	15.836 73	15.603 92	16.035 91
	100K	37.729 79	37.611 00	37.839 07
	1M	116.374 31	112.490 38	120.129 97
	6M	178.966 62	177.311 21	180.607 49
	60M	1 780.927 64	1 762.874 34	1 799.498 02
Iceberg IcedHops	10K	580.480 03	555.595 31	607.645 51
	100K	2 133.403 62	1 958.119 12	2 334.613 16
	1M	2 324.629 63	2 262.967 50	2 381.662 01
	6M	2 825.776 91	2 804.612 49	2 846.730 95
	60M	2 749.827 53	2 735.099 13	2 765.125 11
Delta Lake delta-rs	10K	187.081 71	124.129 41	255.274 25
	100K	1 735.080 31	1 651.502 61	1 816.685 79
	1M	1 856.033 25	1 807.743 07	1 906.811 84
	6M	3 077.758 01	3 046.330 18	3 111.145 48
	60M	2 610.755 20	2 590.903 90	2 626.976 89

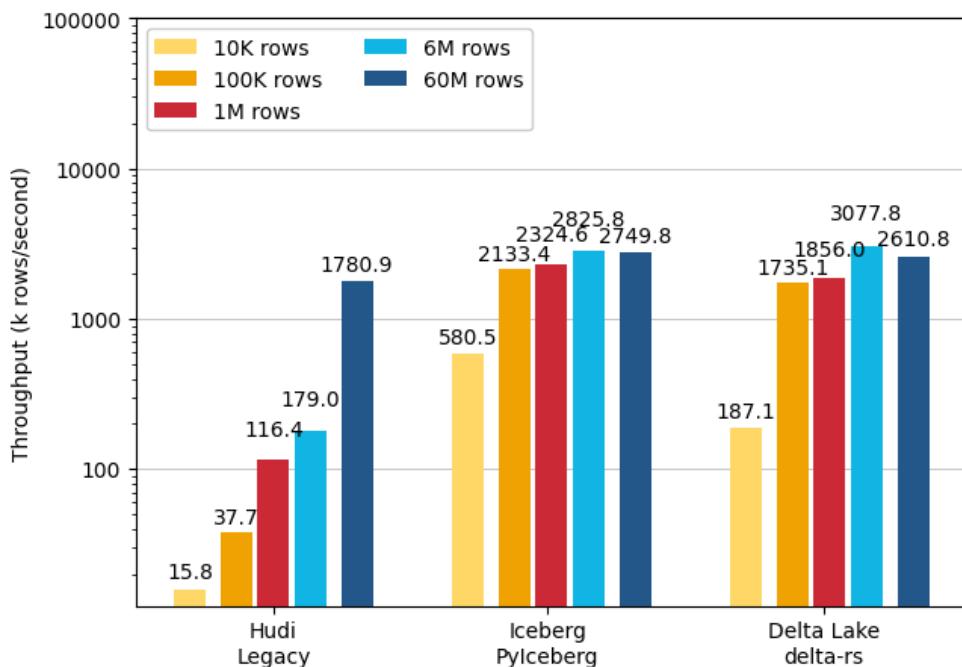


Figure D.5: Histogram in log-scale of the read experiment results expressed as throughput. The experiment was performed with one **CPU** core.

Table D.6: Read experiment results expressed as throughput. The experiment was performed with two **CPU** cores.

Pipeline	Number of rows	Throughput (k rows/second)	Throughput (k rows/second) 95% Confidence Interval	
			low	high
Hudi Legacy	10K	16.004 50	15.922 95	16.082 20
	100K	37.547 12	37.443 31	37.648 42
	1M	116.186 27	111.368 42	120.461 87
	6M	179.762 35	178.005 56	181.286 17
	60M	1 784.185 08	1 763.782 51	1 802.550 61
Iceberg IcedHops	10K	545.758 37	517.757 02	576.367 42
	100K	2 316.163 25	2 213.607 44	2 413.481 99
	1M	3 711.541 85	3 614.794 02	3 814.373 02
	6M	4 629.432 93	4 464.785 44	4 877.936 41
	60M	4 440.187 28	4 403.743 16	4 474.117 01
Delta Lake delta-rs	10K	242.640 97	228.056 10	254.401 00
	100K	1 747.921 59	1 493.971 60	1 957.125 61
	1M	4 278.422 15	4 121.532 10	4 441.780 08
	6M	6 606.023 33	6 533.473 80	6 674.052 28
	60M	5 258.661 43	5 174.163 60	5 322.447 83

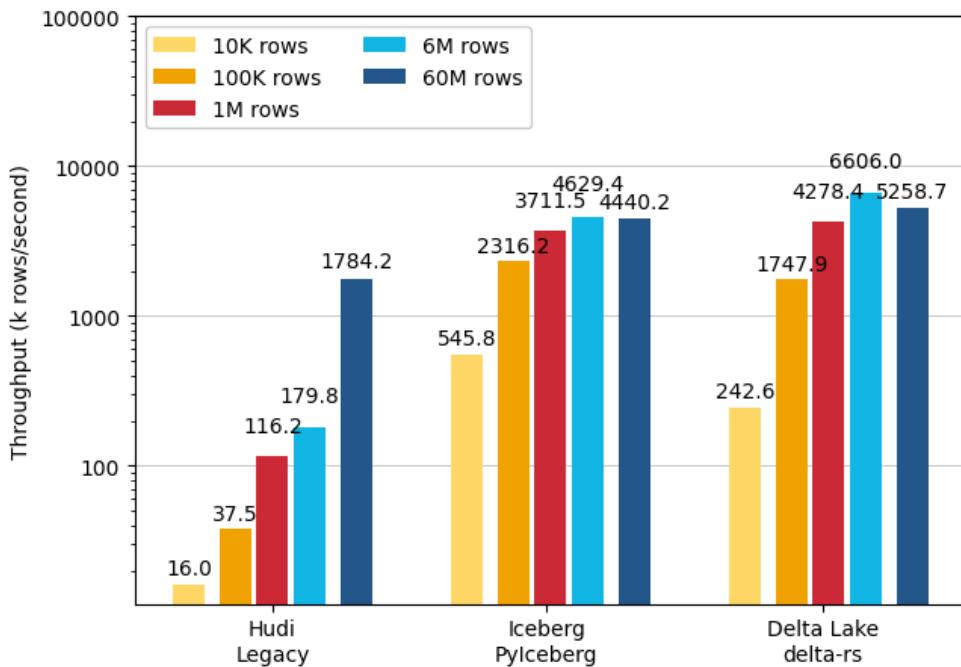


Figure D.6: Histogram in log-scale of the read experiment results expressed as throughput. The experiment was performed with two **CPU** cores.

Table D.7: Read experiment results expressed as throughput. The experiment was performed with four **CPU** cores.

Pipeline	Number of rows	Throughput (k rows/second)	Throughput (k rows/second) 95% Confidence Interval	
			low	high
Hudi Legacy	10K	15.71731	15.16413	16.03729
	100K	37.88087	37.78617	37.97227
	1M	114.35939	111.10533	117.33856
	6M	179.38567	177.82731	180.80062
	60M	1782.45022	1760.00676	1803.79967
Iceberg IcedHops	10K	335.65743	169.10433	637.82194
	100K	2334.62210	2235.35439	2431.53216
	1M	4285.05339	4031.24413	4575.64254
	6M	5949.03530	5904.55474	5992.63025
	60M	5934.54677	5891.98628	5975.55876
Delta Lake delta-rs	10K	229.53446	196.00845	254.37802
	100K	1803.24782	1728.35674	1861.44303
	1M	5274.74103	3917.94620	6633.05053
	6M	11298.14809	10541.11329	11822.24655
	60M	10752.87337	10679.37868	10820.39355

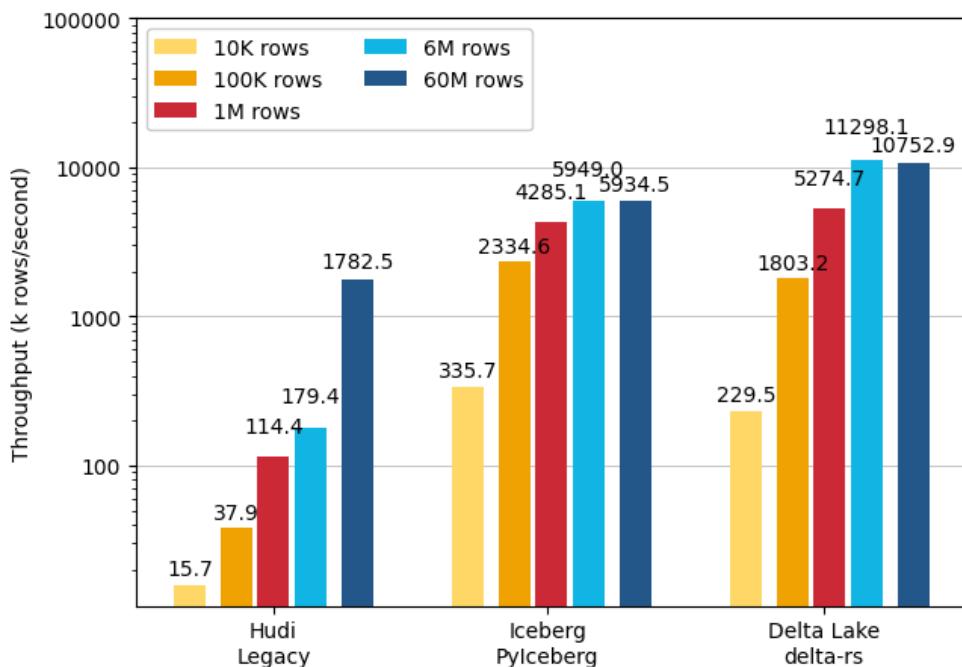


Figure D.7: Histogram in log-scale of the read experiment results expressed as throughput. The experiment was performed with four **CPU** cores.

Table D.8: Read experiment results expressed as throughput. The experiment was performed with eight **CPU** cores.

Pipeline	Number of rows	Throughput (k rows/second)	Throughput (k rows/second) 95% Confidence Interval	
			low	high
Hudi Legacy	10K	15.940 66	15.810 39	16.063 27
	100K	37.560 87	37.422 90	37.689 99
	1M	119.729 81	116.031 88	122.825 54
	6M	179.425 15	177.764 38	180.887 57
	60M	1 810.367 75	1 796.114 54	1 823.667 90
Iceberg IcedHops	10K	589.518 81	552.402 08	625.871 90
	100K	2 343.108 33	2 233.189 92	2 441.607 41
	1M	4 308.143 74	4 249.775 75	4 361.607 98
	6M	6 782.335 84	6 707.348 51	6 852.171 57
	60M	6 787.943 65	6 715.992 13	6 850.140 93
Delta Lake delta-rs	10K	231.783 14	194.952 45	257.555 78
	100K	1 829.910 95	1 765.222 67	1 887.200 60
	1M	5 744.056 36	5 608.985 88	5 888.367 78
	6M	12 070.777 05	11 742.934 47	12 355.358 82
	60M	20 378.273 02	19 652.710 07	20 981.075 17

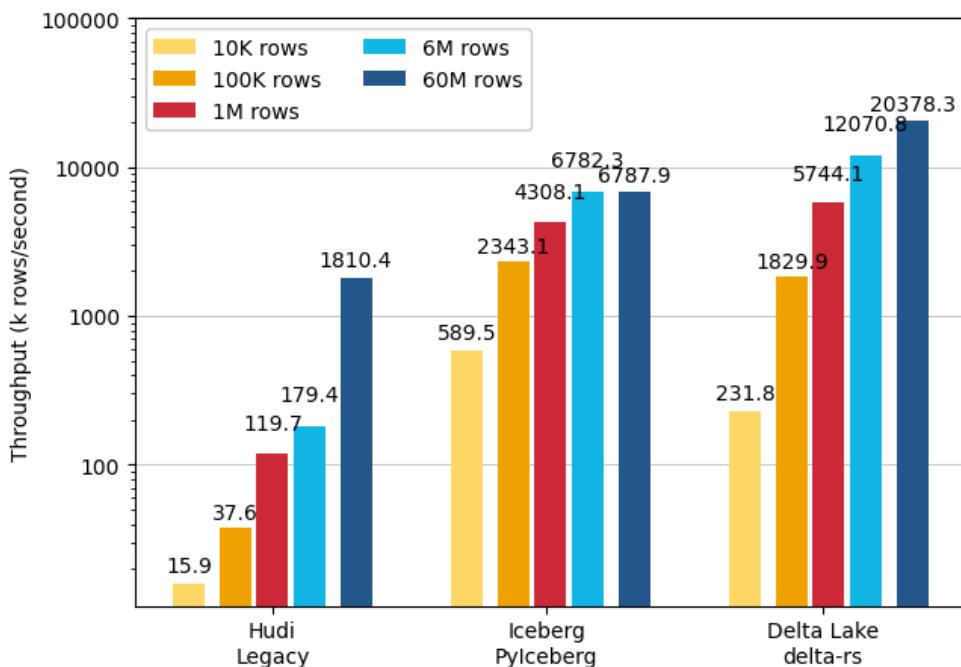


Figure D.8: Histogram in log-scale of the read experiment results expressed as throughput. The experiment was performed with eight **CPU** cores.

Appendix E

Legacy pipeline write latency breakdown results

This appendix contains all the graphs and the tables related to the write latency breakdown, between upload and materialization steps, for the Hudi-based legacy pipeline. Results are reported expressed as latency (measured in seconds), which was measured directly.

Table E.1: Contributions to the write latency of the upload and materialization steps in the legacy pipeline. The experiment was performed with one **CPU** core.

Step	Number of rows	Latency (seconds)	Latency (seconds) 95% Confidence Interval	
			low	high
upload materialize	10K	2.4865	2.3896	2.6261
		47.7262	47.0445	48.4031
upload materialize	100K	3.6684	3.6310	3.7098
		55.9005	55.2494	56.5541
upload materialize	1M	22.5934	22.4496	22.7349
		89.5754	88.8286	90.3049
upload materialize	6M	244.6123	244.0368	245.1905
		267.2490	266.4287	268.1549
upload materialize	60M	2437.7840	2422.8704	2453.6746
		278.0504	276.2340	280.0921

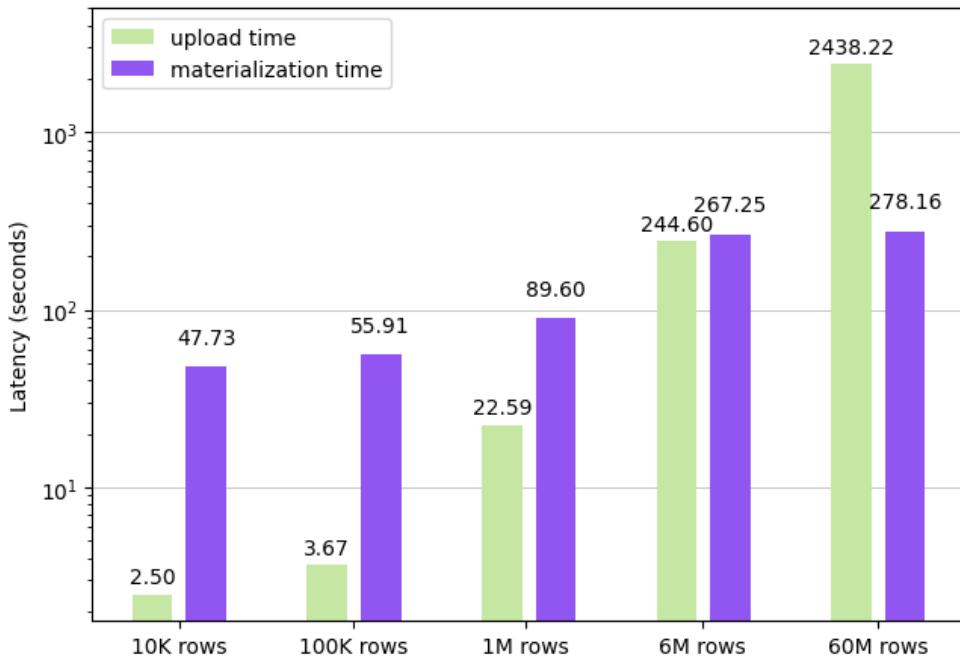


Figure E.1: Histogram in log-scale displaying the contributions to the write latency of the upload and materialization steps in the legacy pipeline. The experiment was performed with one **CPU** core.

Table E.2: Contributions to the write latency of the upload and materialization steps in the legacy pipeline. The experiment was performed with two CPU cores.

Step	Number of rows	Latency (seconds)	Latency (seconds) 95% Confidence Interval	
			low	high
upload materialize	10K	2.3873	2.3276	2.4466
		48.3305	47.7020	48.9923
upload materialize	100K	3.4348	3.4008	3.4671
		56.3367	55.5626	57.1129
upload materialize	1M	18.6349	18.5673	18.7104
		89.9267	89.4012	90.4514
upload materialize	6M	205.8854	205.2177	206.4984
		267.5079	266.6853	268.3512
upload materialize	60M	2064.1357	2057.6396	2070.4450
		276.9608	275.7156	278.2849

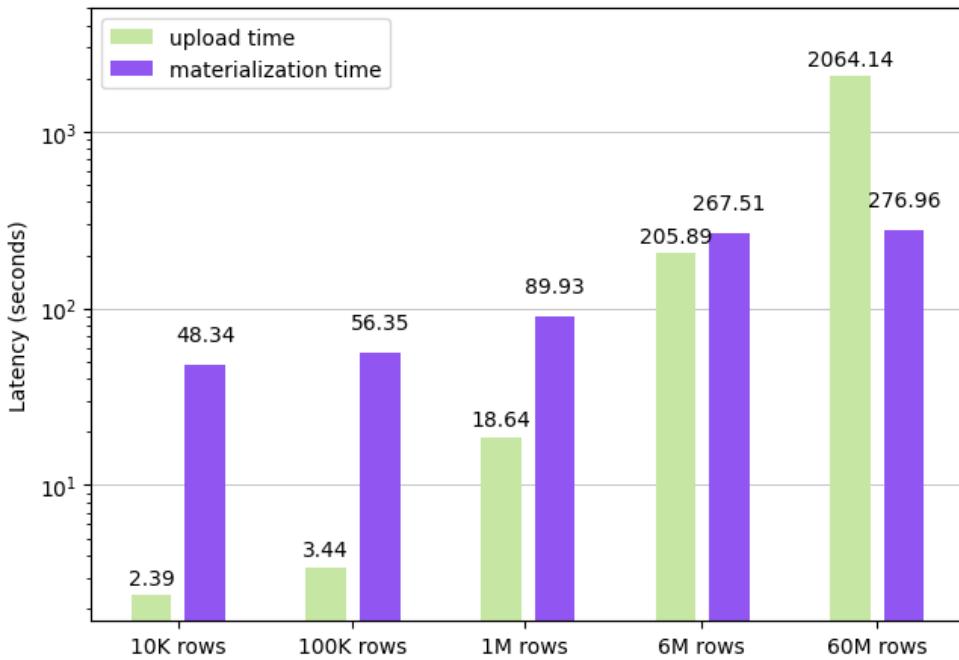


Figure E.2: Histogram in log-scale displaying the contributions to the write latency of the upload and materialization steps in the legacy pipeline. The experiment was performed with two CPU cores.

Table E.3: Contributions to the write latency of the upload and materialization steps in the legacy pipeline. The experiment was performed with four CPU cores.

Step	Number of rows	Latency (seconds)	Latency (seconds)	
			95% Confidence Interval low	high
upload materialize	10K	2.3846	2.3299	2.4335
		48.9061	48.2436	49.5470
upload materialize	100K	3.4650	3.4245	3.5071
		56.0524	55.3682	56.6822
upload materialize	1M	19.2296	19.1455	19.3161
		89.5864	89.0209	90.1313
upload materialize	6M	211.6758	211.0694	212.2839
		270.3233	269.5967	270.9895
upload materialize	60M	2068.5260	2060.3358	2077.2837
		277.6001	276.0065	279.1456

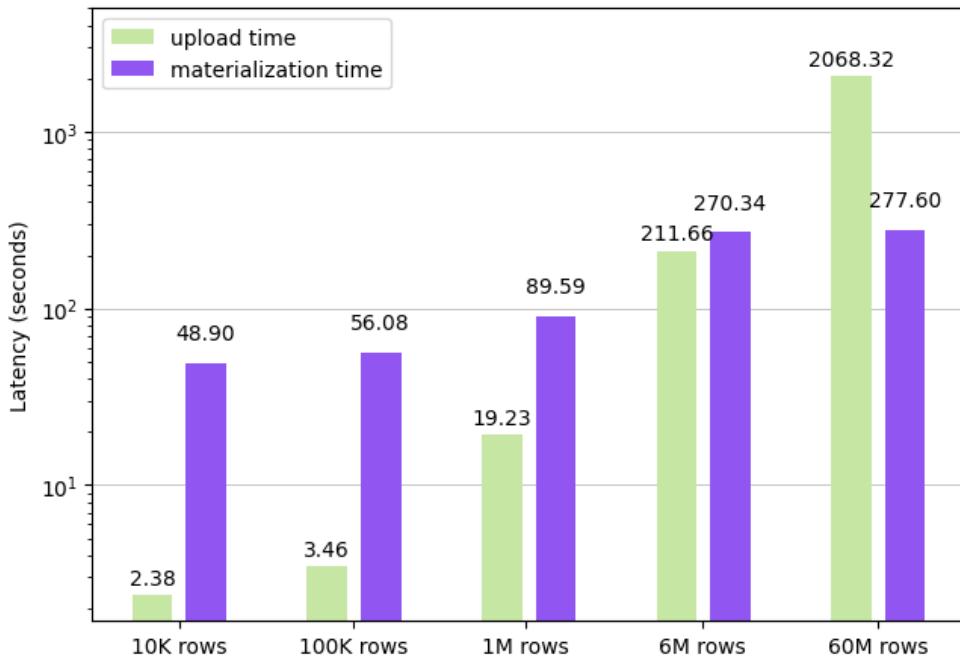


Figure E.3: Histogram in log-scale displaying the contributions to the write latency of the upload and materialization steps in the legacy pipeline. The experiment was performed with four CPU cores.

Table E.4: Contributions to the write latency of the upload and materialization steps in the legacy pipeline. The experiment was performed with eight **CPU** cores.

Step	Number of rows	Latency (seconds)	Latency (seconds)	
			95% Confidence Interval low	high
upload materialize	10K	2.3815	2.3304	2.4358
		48.8485	48.1979	49.4467
upload materialize	100K	3.4392	3.4081	3.4700
		56.8428	56.3177	57.3685
upload materialize	1M	18.8642	18.7808	18.9532
		90.5153	90.0306	90.9718
upload materialize	6M	207.6646	207.1606	208.2090
		268.2752	267.3569	269.2456
upload materialize	60M	2049.1371	2043.5991	2055.4782
		275.7636	274.2773	274.2773

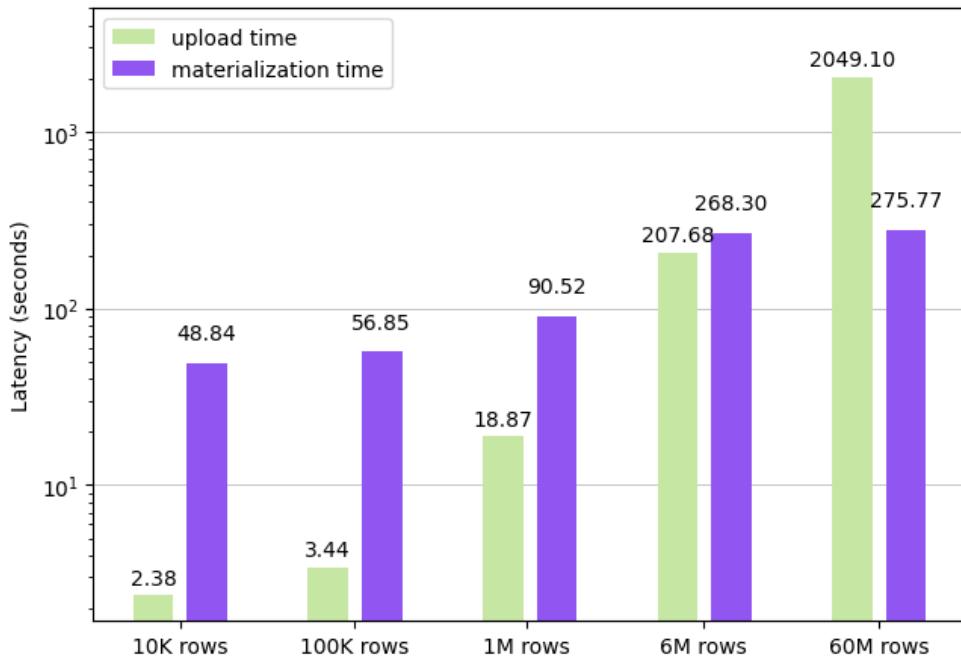


Figure E.4: Histogram in log-scale displaying the contributions to the write latency of the upload and materialization steps in the legacy pipeline. The experiment was performed with eight **CPU** cores.

€€€€ For DIVA €€€€

```
{  
    "Author1": { "Last name": "Meneghin",  
    "First name": "Sebastiano",  
    "Local User Id": "u16dn7xj",  
    "E-mail": "meneghin@kth.se",  
    "organisation": {"L1": "School of Electrical Engineering and Computer Science",  
    }  
    },  
    "Cycle": "2",  
    "Course code": "DA258X",  
    "Credits": "30.0",  
    "Degree1": {"Educational program": "Master's Programme, ICT Innovation, 120 credits"  
    , "programcode": "TIVNM"},  
    "Degree": "Master's Programme, Distributed Systems and Data Mining for Big Data, 120 credits"  
    , "subjectArea": "Computer Science"},  
    "Title": {  
        "Main title": "IcedHops: reducing read and write latency in an Iceberg-backed offline feature store",  
        "Subtitle": "Integrating HopsFS and Pylceberg Python library to reduce read and write latency on Hopsworks offline feature store, with a  
        comparative analysis of alternative solutions",  
        "Language": "eng"},  
        "Alternative title": {  
            "Main title": "IcedHops: minska läs- och skrivlatens i en Iceberg-stödd offlinebutik för funktioner",  
            "Subtitle": "Integrering av HopsFS och Pylceberg Python-bibliotek för att minska läs- och skrivförrörningen i den offline-lagringen Hopsworks, med  
            en jämförande analys av alternativa lösningar.",  
            "Language": "swe"},  
    },  
    "Supervisor1": { "Last name": "Sheikholeslami",  
    "First name": "Sina",  
    "Local User Id": "u1znylhw",  
    "E-mail": "sinash@kth.se",  
    "organisation": {"L1": "School of Electrical Engineering and Computer Science",  
    "L2": "Computer Science"},  
    },  
    "Supervisor2": { "Last name": "Schmidt",  
    "First name": "Fabian",  
    "Local User Id": "u1mrsz0u",  
    "E-mail": "fschm@kth.se"},  
    "Supervisor3": { "Last name": "Bzhalava",  
    "First name": "Davit",  
    "E-mail": "davit@hopsworks.ai",  
    "Other organisation": "Hopsworks AB"},  
    "Examiner1": { "Last name": "Vlassov",  
    "First name": "Vladimir",  
    "Local User Id": "u19yb2c8",  
    "E-mail": "vladv@kth.se",  
    "organisation": {"L1": "School of Electrical Engineering and Computer Science",  
    "L2": "Computer Science"},  
    },  
    "Cooperation": { "Partner_name": "Hopsworks AB"},  
    "National Subject Categories": "10201, 10205, 10206",  
    "Other information": {"Year": "2025", "Number of pages": "1,119"},  
    "Copyrightleft": "copyright",  
    "Series": {"Title of series": "TRITA – EECS-EX", "No. in series": "2025:0000"},  
    "Opponents": {"Name": "Arvid Wennerström"},  
    "Presentation": {"Date": "2025-02-25 15:00"},  
    "Language": "eng",  
    "Room": "via Zoom https://kth-se.zoom.us/j/ddddddddddd",  
    "Address": "Isafjordsgatan 22 (Kistagången 16)",  
    "City": "Stockholm"},  
    "Number of lang instances": "3",  
    "Abstract[eng]": "€€€€  
The growing need for efficient data management in Machine Learning (ML) workflows has led to the widespread adoption of feature stores, centralized data platforms that supports feature engineering, model training and prediction inference. The Hopsworks' feature store has demonstrated outperformance compared to its alternatives, leveraging Apache Hudi and Spark for offline data storage, but suffers from high write and read latency, even for small quantities of data (1GB or less). This thesis explores the potential of Apache Iceberg as an alternative table format to reduce latencies, developing "IcedHops", an integration of HopsFS (Hopsworks HDFS distribution) and Pylceberg Python library. The research begin with an evaluation of potential system integration alternatives, documenting the advantages and limitations of each approach. Then, IcedHops is implemented and evaluated, benchmarking it against the existing Spark-based solution and an alternative Delta Lake implementation (delta-rs). Extensive experiments were conducted across varying table sizes and CPU configurations to assess write and read performance. Results show that IcedHops significantly reduces write latency – from 40 to 140 times lower than the legacy system – and read latency – from 55% to 60 times lower than the legacy system. Compared to delta-rs, IcedHops demonstrates reduced write latency for large tables – up to 7 times lower – and equal read latency, but exhibits lower scaling benefits with additional CPU cores – 20% less than delta-rs. These findings confirm that alternatives to Spark-based pipelines in small-scale scenarios are possible and are worth of further investigations, and",  
}
```

the system implemented will be included in the Hopsworks feature store. Furthermore, this thesis work and results finally provides a baseline for future work about additional open table formats, alternative languages to mitigate Python's performance overhead, and strategies to improve resource utilization in data management platforms.

€€€€,

"Keywords[eng]": €€€€

Machine Learning, Feature Store, Spark, Apache Iceberg, Delta Lake, Python, Read/Write Latency, Open Table Formats €€€€,

"Abstract[swe]": €€€€

Det växande behovet av effektiv datahantering i arbetsflöden för maskininlärning (ML) har lett till en utbredd användning av feature stores – centraliseraade dataplattformar som stöder feature engineering, modellträning och inferens. Hopsworks feature store har visat bättre prestanda jämfört med sina alternativ och använder Apache Hudi och Spark för offline-datalagring. Dock lider systemet av hög skriv- och läslatens, även för små datamängder (1 GB eller mindre). Denna avhandling undersöker potentialen hos Apache Iceberg som ett alternativt tabellformat och integrerar det med HopsFS (Hopsworks HDFS-distribution) samt Pylceberg Python-biblioteket för att minska latensen.

Forskningsenheten inleds med en utvärdering av potentiella systemintegrationsalternativ där fördelar och begränsningar med varje metod dokumenteras. Därefter implementeras och utvärderas en Pylceberg-baserad arkitektur, vilken jämförs med den befintliga Spark-baserade lösningen samt en alternativ Delta Lake-implementering (delta-rs). Omfattande experiment genomförs med varierande tabellstorlekar och CPU-konfigurationer för att bedöma skriv- och läsprestanda. Resultaten visar att Pylceberg avsevärt minskar skrífördröjningen – från 40 till 140 gånger lägre än det äldre systemet – och läsfördröjningen – från 55% till 60 gånger lägre än det äldre systemet. Jämfört med delta-rs upptäcks Pylceberg minskad skrífördröjning för stora tabeller – upp till sju gånger lägre – och liknande läsfördröjning, men har sämre skalningsfördelar vid ökning av CPU-kärnor (20% mindre än delta-rs).

Dessa resultat bekräftar att alternativ till Spark-baserade pipelines i småskaliga scenarier är möjliga och värda ytterligare undersökningar. Det implementerade systemet kommer att integreras i Hopsworks feature store. Dessutom utgör denna avhandling en baslinje för framtida forskning kring ytterligare öppna tabellformat, alternativa programmeringsspråk för att hantera Pythons prestandabegränsningar samt strategier för att förbättra resursutnyttjandet i datahanteringsplattformar. €€€€,

"Keywords[swe]": €€€€

Maskininlärning, Feature Store, Spark, Apache Iceberg, Delta Lake, Python, Läs- och skrivelatens, Open Table Formats €€€€,

"Abstract[ita]": €€€€

La crescente necessità di piattaforme per una efficiente gestione dei dati per applicazioni di *Machine Learning (ML)* ha portato a un'ampia diffusione dei *feature store*, piattaforme dati centralizzate che supportano *feature engineering*, addestramento di modelli e inferenza. Il *feature store* di Hopsworks ha dimostrato prestazioni superiori rispetto alle sue alternative, sfruttando Apache Hudi e Spark per il suo *offline feature store*. Tuttavia, questo sistema soffre di un'elevata latenza di scrittura e lettura, anche per piccole quantità di dati (1GB o inferiori). Questa tesi esplora l'uso di Apache Iceberg come *table format* alternativo per ridurre tali latenze, sviluppando "IcedHops", una integrazione tra HopsFS (distribuzione Hopsworks di HDFS) e la libreria Python Pylceberg.

Questa ricerca inizia con un'analisi delle possibili strategie di integrazione, documentando i vantaggi e i limiti di ciascun approccio. Successivamente, IcedHops viene sviluppato ed evaluato, confrontandolo con la soluzione esistente basata su Spark e con un'alternativa basata su Delta Lake (delta-rs). Gli esperimenti approfonditi sono stati condotti su tabelle di diverse dimensioni e diverse configurazioni CPU per misurare le prestazioni di scrittura e lettura. I risultati mostrano che IcedHops riduce significativamente la latenza di scrittura – da 40 a 140 volte inferiore rispetto al sistema legacy – e la latenza di lettura – dal 55% a 60 volte inferiore. Rispetto a delta-rs, IcedHops offre una latenza di scrittura inferiore per le tabelle più grandi – fino a 7 volte minore – e prestazioni di lettura equivalenti, ma mostra minori vantaggi di scalabilità con l'aumento dei CPU cores (20% ridotti rispetto a delta-rs).

Questi risultati confermano che alternative ad architetture basate su Spark, per gestire dati su piccola scala, esistono e sono più efficienti, e dunque il sistema sviluppato verrà integrato nel *feature store* di Hopsworks. Inoltre, questa tesi fornisce una base per futuri studi su nuovi *open table format* e strategie da adottare per ottimizzare l'uso delle risorse nelle piattaforme di gestione dei dati. €€€€,

"Keywords[ita]": €€€€

Machine Learning, Feature Store, Apache Iceberg, Delta Lake, Pylceberg, Latenza in Lettura/Scrittura, Open Table Formats. €€€€,

}

acronyms.tex

```
%%% Local Variables:
%%% mode: latex
%%% TeX-master: t
%%% End:
% The following command is used with glossaries-extra
\setabbreviationstyle[acronym]{long-short}
% The form of the entries in this file is \newacronym{label}{acronym}{phrase}
% or \newacronym[options]{label}{acronym}{phrase}
% see "User Manual for glossaries.sty" for the details

\newacronym{KTH}{KTH}{KTH Royal Institute of Technology}
\newacronym{ACID}{ACID}{Atomicity, Consistency, Isolation and Durability}
\newacronym{AI}{AI}{Artificial Intelligence}
\newacronym{ML}{ML}{Machine Learning}
\newacronym{BI}{BI}{Business Intelligence}
\newacronym[shortplural={RDDs}, firstplural={Resilient Distributed Datasets (RDDs)}]{RDD}{RDD}{Resilient Distributed Dataset}
\newacronym{OLAP}{OLAP}{On-Line Analytical Processing}
\newacronym{ELT}{ELT}{Extract Load Transform}
\newacronym{ETL}{ETL}{Extract Transform Load}
\newacronym{HDFS}{HDFS}{Hadoop Distributed File System}
\newacronym{JVM}{JVM}{Java Virtual Machine}
\newacronym[shortplural={INs}, firstplural={Industrial Needs (INs)}]{IN}{IN}{Industrial Need}
\newacronym[shortplural={PAs}, firstplural={Project Assumptions (PAs)}]{PA}{PA}{Project Assumption}
\newacronym[shortplural={APIs}, firstplural={Application Programming Interfaces (APIs)}]{API}{API}{Application Programming Interface}
\newacronym{OLTP}{OLTP}{On-Line Transaction Processing}
\newacronym[shortplural={DBMs}, firstplural={Data Base Management Systems}]{DBMS}{DBMS}{Data Base Management System}
\newacronym[shortplural={Gs}, firstplural={Goals}]{G}{G}{Goal}
\newacronym[shortplural={RQs}, firstplural={Research Questions}]{RQ}{RQ}{Research Question}
\newacronym[shortplural={Ds}, firstplural={Deliverables}]{D}{D}{Deliverable}
\newacronym{CRUD}{CRUD}{Create Read Update Delete}
\newacronym[shortplural={SDGs}, firstplural={Sustainable Development Goals}]{SDG}{SDG}{Sustainable Development Goal}
\newacronym{AWS}{AWS}{Amazon Web Services}
\newacronym{GCS}{GCS}{Google Cloud Storage}
\newacronym[shortplural={HDDs}, firstplural={Hard Disks Drives}]{HDD}{HDD}{Hard Disk Drive}
\newacronym[shortplural={SSDs}, firstplural={Solid State Drives}]{SSD}{SSD}{Solid State Drive}
\newacronym[shortplural={PCs}, firstplural={Personal Computers}]{PC}{PC}{Personal Computer}
\newacronym[shortplural={OSes}, firstplural={Operating Systems (OSes)}]{OS}{OS}{Operating System}
\newacronym[shortplural={DFSes}, firstplural={Distributed File Systems (DFSes)}]{DFS}{DFS}{Distributed File System}
\newacronym{BPMN}{BPMN}{Business Process Model and Notation}
\newacronym{TLS}{TLS}{Transport Layer Security}
\newacronym{SSH}{SSH}{Secure Shell protocol}
\newacronym{VM}{VM}{Virtual Machine}
\newacronym{CoC}{CoC}{Conquer of Completion}
\newacronym{SF}{SF}{Scale Factor}
\newacronym{CPU}{CPU}{Central Processing Unit}
\newacronym{GPU}{GPU}{Graphical Processing Unit}
\newacronym{HopsFS}{HopsFS}{Hopsworks' \glsentryshort{HDFS} distribution}
\newacronym{LocalFS}{LocalFS}{Local File System}
\newacronym{TPC}{TPC}{Transaction Processing Performance Council}
\newacronym{RAM}{RAM}{Random Access Memory}
\newacronym{RPC}{RPC}{Remote Procedural Call}
\newacronym{CIDR}{CIDR}{Conference on Innovative Data Systems Research}
\newacronym{MLOps}{MLOps}{Machine Learning Operations}
\newacronym{LOC}{LOC}{Lines Of Code}
\newacronym{NN}{NN}{Neural Network}
\newacronym{SAN}{SAN}{Storage Area Network}
\newacronym{NAS}{NAS}{Network Attached Storage}
\newacronym{LAN}{LAN}{Local Area Network}
\newacronym{WAN}{WAN}{Wide Area Network}
\newacronym[shortplural={OTFs}, firstplural={Open Table Formats}]{OTF}{OTF}{Open Table Format}
\newacronym{CDC}{CDC}{Change Data Capture}
\newacronym{CoW}{CoW}{Copy on Write}
\newacronym{MoR}{MoR}{Merge on Read}
\newacronym{CC}{CC}{Concurrency Control}
\newacronym{SDK}{SDK}{Software Development Kit}
\newacronym{HMS}{HMS}{Hive Metastore}
\newacronym{JDBC}{JDBC}{Java DataBase Connectivity}
\newacronym{SQL}{SQL}{Structured Query Language}
```