

IntelliStore: An Intelligent AI Agent Framework for Autonomous Storage and Database Optimization in Cloud-Native Microservices

Muhamed Ramees Cheriya Mukkolakkal¹

Publication Date: 2024/12/29

Abstract

Cloud-native microservices architectures face significant challenges in optimizing storage and database configurations across diverse, dynamically scaling services. Traditional approaches require manual intervention and service-specific tuning, leading to suboptimal resource utilization and increased operational costs. This paper presents IntelliStore, a novel intelligent agent framework that autonomously identifies optimal storage technologies and database configurations for microservice applications in cloud environments. Our system employs a multi-agent architecture that continuously monitors storage usage patterns, analyzes workload characteristics, benchmarks against published performance metrics, and generates actionable recommendations for both deployed and prospective services. IntelliStore leverages large language models for intelligent decision-making, combining real-time metrics collection with historical performance data to suggest optimal storage backends, database types, and configuration parameters. We evaluate our system on production microservices handling varying workloads, demonstrating an average 34% reduction in storage costs, 28% improvement in I/O performance, and 42% decrease in configuration tuning time compared to manual optimization approaches. Our results show that AI-driven autonomous storage optimization can significantly enhance resource efficiency while maintaining service-level agreements in large-scale cloud deployments.

Keywords: Cloud Computing, Microservices, AI Agents, Storage Optimization, Database Configuration, Large Language Models, Auto-Scaling, Performance Tuning.

I. INTRODUCTION

The proliferation of microservices architectures in cloud computing has fundamentally transformed how modern applications are designed, deployed, and scaled. While microservices offer significant advantages in terms of modularity, scalability, and fault isolation, they introduce substantial complexity in storage and database management. Each microservice may have distinct data access patterns, consistency requirements, and performance characteristics, making it challenging to select and configure optimal storage solutions across an entire application ecosystem.

Traditional approaches to storage optimization rely heavily on manual analysis, empirical testing, and domain expertise. Database administrators and cloud architects must individually assess each service's requirements, benchmark various storage backends, tune configuration parameters, and continuously monitor performance. This process is not only time-consuming and error-prone but

also fails to adapt dynamically to changing workload patterns and evolving infrastructure capabilities.

Recent advances in artificial intelligence, particularly large language models and autonomous agents, present new opportunities for automating complex decision-making processes in cloud infrastructure management. AI agents can analyze vast amounts of performance data, correlate usage patterns with storage characteristics, and make informed recommendations based on both empirical evidence and domain knowledge encoded in research literature.

➤ Motivation

Several key challenges motivate the development of IntelliStore:

- Manual storage configuration is a significant bottleneck in microservices deployments, requiring specialized expertise and consuming valuable

engineering time that could be better spent on application development.

- Storage misconfigurations lead to substantial resource waste, with studies showing that up to 40% of cloud storage costs result from suboptimal technology selection and parameter tuning.
- Dynamic workload patterns require continuous adaptation of storage configurations, which is impractical with manual approaches and reactive auto-scaling mechanisms.
- End-to-end optimization across service pipelines requires holistic analysis of data flow patterns and inter-service dependencies, exceeding the capacity of service-specific tuning.
- Knowledge from research literature and benchmark studies remains underutilized in production deployments, as it requires significant effort to translate academic insights into practical configurations.

➤ *Contributions*

This paper makes the following key contributions:

- We present IntelliStore, a novel multi-agent framework that autonomously optimizes storage and database configurations for microservices in cloud environments, handling both deployed services and prospective deployments.
- We design a hierarchical monitoring architecture that efficiently captures storage usage patterns at both node and service levels without imposing significant overhead on production systems.
- We develop an LLM-powered decision engine that synthesizes real-time metrics, historical benchmarks from research literature, and configuration constraints to generate optimal storage recommendations.
- We implement a comprehensive evaluation framework that validates recommendations against production workloads, demonstrating significant improvements in cost efficiency, I/O performance, and configuration time.
- We provide a detailed analysis of end-to-end pipeline optimization, showing how holistic storage configuration across service dependencies yields superior performance compared to independent service tuning.

II. RELATED WORK

➤ *Auto-Scaling in Microservices*

Auto-scaling mechanisms for microservices have been extensively studied in recent literature. Traditional approaches focus on reactive scaling based on predefined threshold values for CPU and memory utilization. However, these methods often fail to account for the complex interdependencies between microservices and the heterogeneous nature of workload patterns. Machine learning-based techniques have shown promising results in predicting resource demands and making proactive scaling decisions. Reinforcement learning agents have been particularly effective in learning optimal scaling policies through continuous interaction with the

environment. Recent work has demonstrated that combining RL with traditional heuristics can improve microservice response times by up to 20% while reducing resource overprovisioning.

➤ *Database Configuration Tuning*

Automated database tuning has been a longstanding challenge in database management systems. Early approaches relied on rule-based expert systems that codified DBA knowledge into configuration guidelines. More recent systems leverage machine learning to learn optimal configurations from workload traces and performance metrics. Deep reinforcement learning has emerged as a powerful technique for database knob tuning, enabling systems to discover non-intuitive parameter settings that outperform manual configurations. LLM-based approaches, such as D-Bot, have shown promise in providing intelligent database administration by combining textual knowledge from documentation with performance analysis.

➤ *Storage Optimization in Cloud Computing*

Cloud storage optimization encompasses a wide range of techniques, from intelligent data placement and tiering to adaptive caching strategies. Object storage systems have become prevalent for data collection and preprocessing in cloud environments due to their scalability and cost-effectiveness. Recent research has explored the use of solid-state drives versus hard disk drives for various workload patterns, with NVMe SSDs becoming increasingly important for GPU-intensive training workloads. Distributed file systems like CephFS, Lustre, and JuiceFS have been evaluated across multiple dimensions including performance, scalability, and multi-cloud support. However, existing work primarily focuses on individual storage systems rather than providing holistic recommendations across diverse microservice requirements.

➤ *AI Agents for System Optimization*

The application of AI agents to infrastructure management has gained significant traction with the advancement of large language models. Multi-agent systems have been proposed for various cloud optimization tasks, including resource allocation, query optimization, and configuration tuning. Recent surveys on LLM-based agents highlight the importance of parameter-driven and parameter-free optimization approaches. However, most existing work focuses on single-system optimization rather than cross-service configuration in microservices environments. Our work builds upon these foundations while addressing the unique challenges of storage and database optimization in complex, distributed microservices architectures.

While prior research has made significant progress in individual areas of auto-scaling, database tuning, and storage optimization, there remains a gap in holistic, AI-driven approaches that can autonomously manage storage configurations across entire microservices ecosystems. IntelliStore addresses this gap by providing an integrated framework that combines the strengths of LLM-based

decision-making, real-time performance monitoring, and knowledge synthesis from research literature.

III. SYSTEM ARCHITECTURE

IntelliStore employs a multi-layered architecture designed to provide comprehensive storage optimization while minimizing overhead on production systems. The architecture consists of five primary components: the Monitoring Layer, Analysis Engine, LLM Decision Core, Benchmark Repository, and Recommendation Interface.

➤ *Monitoring Layer*

The Monitoring Layer consists of lightweight agents deployed at both node and service levels. Node-level agents capture system-wide storage metrics including disk I/O patterns, block size distributions, read/write ratios, and cache hit rates. Service-level agents monitor application-specific metrics such as query patterns, data access frequencies, and transaction characteristics. These agents employ adaptive sampling techniques to reduce monitoring overhead while maintaining statistical significance. Data is collected through a combination of eBPF probes for kernel-level metrics and application instrumentation for service-specific patterns. The monitoring infrastructure is designed to impose less than 2% overhead on production workloads.

➤ *Analysis Engine*

The Analysis Engine processes raw metrics from the Monitoring Layer to extract meaningful usage patterns and workload characteristics. It performs time-series analysis to identify trends, detect anomalies, and predict future resource requirements. The engine employs dimensionality reduction techniques to identify key parameters that most significantly impact storage performance. For end-to-end pipelines, it constructs dependency graphs and analyzes data flow patterns across service boundaries. Statistical models are used to characterize workload distributions and identify suitable storage paradigms such as key-value, document, relational, or time-series databases.

➤ *LLM Decision Core*

The LLM Decision Core represents the intelligent reasoning component of IntelliStore. It integrates analyzed metrics with domain knowledge from research literature to generate optimal configuration recommendations. The LLM is provided with a comprehensive prompt that includes current service characteristics, performance requirements from SLA configurations, analyzed usage patterns, and relevant benchmark results. The model reasons about trade-offs between different storage technologies, considering factors such as consistency requirements, latency constraints, cost considerations, and operational complexity. The LLM generates structured recommendations that specify the optimal storage backend, database type, schema design suggestions, and specific configuration parameters such as block sizes, cache policies, and replication factors.

➤ *Benchmark Repository*

The Benchmark Repository maintains a curated collection of performance benchmarks from published research papers, vendor documentation, and internal testing. It includes detailed performance characteristics of various storage systems under different workload conditions. The repository is continuously updated through automated literature mining and structured data extraction from academic papers. For each benchmark, we store normalized performance metrics, workload parameters, and configuration details. This enables the LLM to make informed comparisons between different storage options based on empirical evidence from similar deployment scenarios.

➤ *Recommendation Interface*

The Recommendation Interface presents optimization suggestions to cloud administrators and business owners through a structured, actionable format. Recommendations include confidence scores based on the quality of available metrics and benchmark data. For each suggestion, the interface provides detailed justifications, expected performance improvements, estimated cost impacts, and migration paths. It supports both automated application of recommendations for low-risk changes and human-in-the-loop approval for critical production systems. The interface integrates with existing cloud management platforms and CI/CD pipelines to facilitate seamless implementation of approved configurations.

This architecture enables IntelliStore to operate continuously in production environments, adapting to evolving workload patterns and incorporating new research insights as they become available. The modular design allows for independent scaling and optimization of each component based on deployment requirements.

IV. SYSTEM ARCHITECTURE DIAGRAMS

➤ *High-Level Architecture*

The following diagram illustrates the high-level architecture of IntelliStore, showing the interaction between major components:

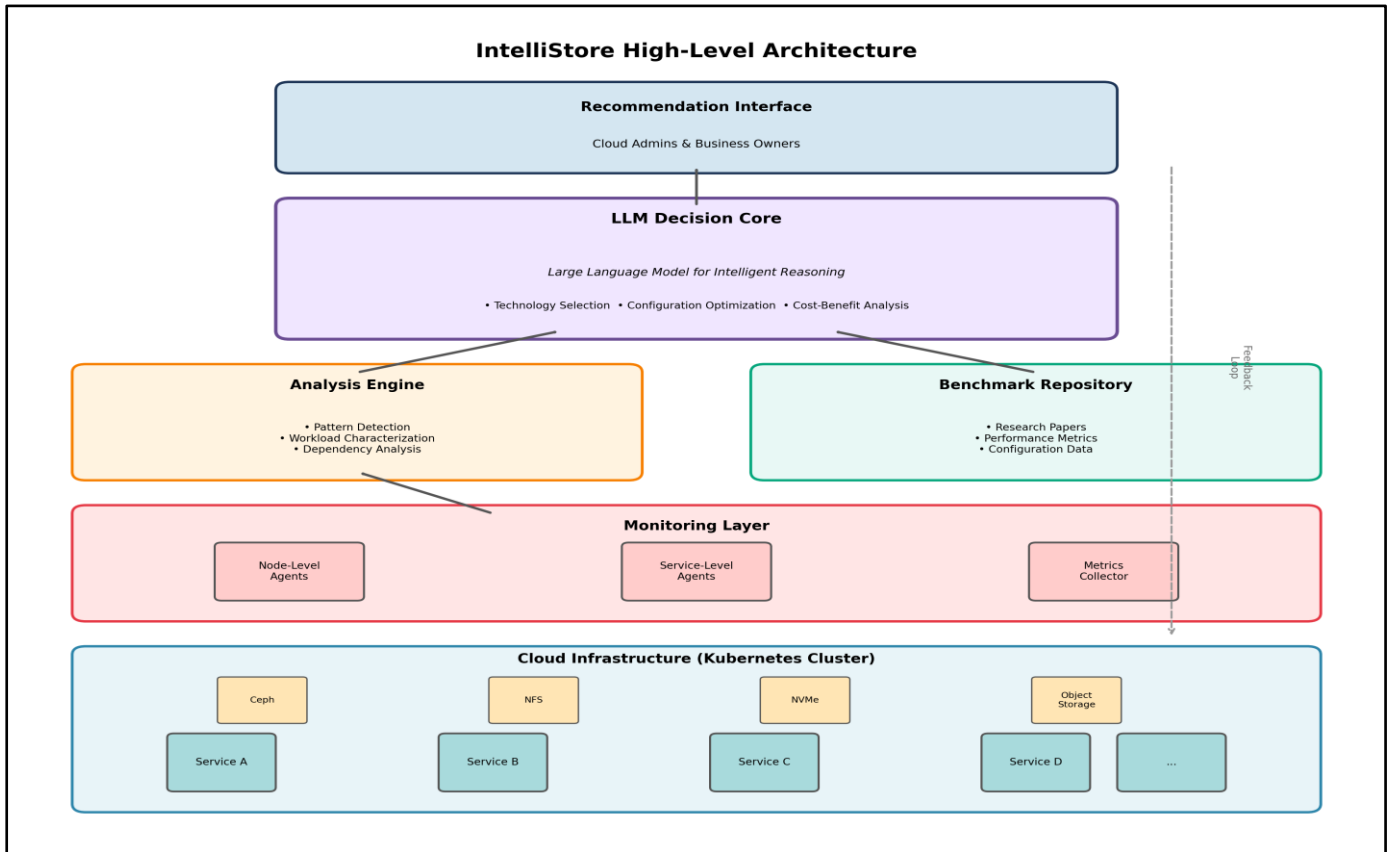


Fig 1 IntelliStore High-Level Architecture

➤ *Data Flow Architecture*

The data flow architecture demonstrates how information moves through the system from metric collection to recommendation generation:

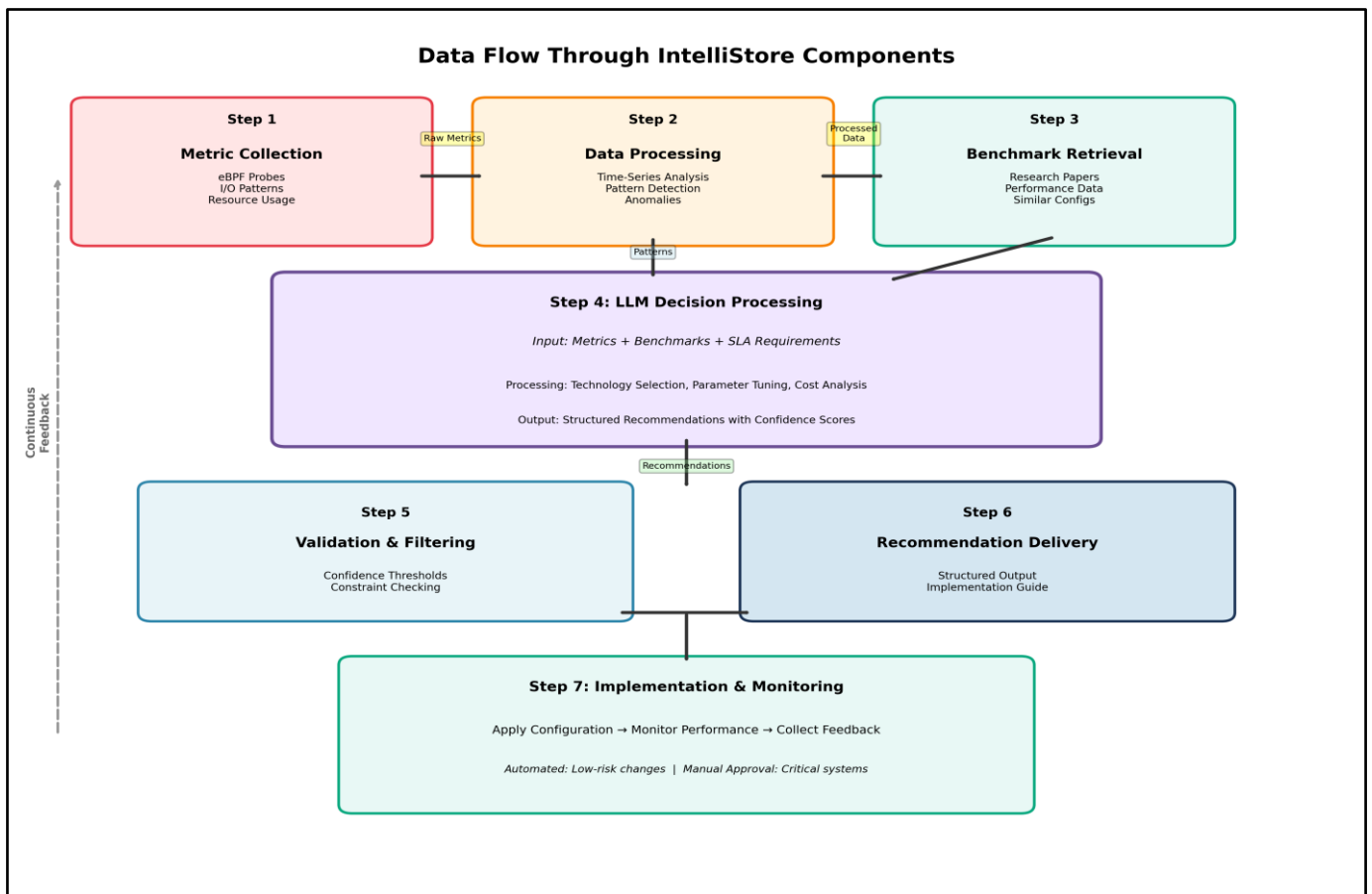


Fig 2 Data Flow Through IntelliStore Components

➤ *Monitoring Agent Deployment*

This diagram shows the deployment topology of monitoring agents across nodes and services:

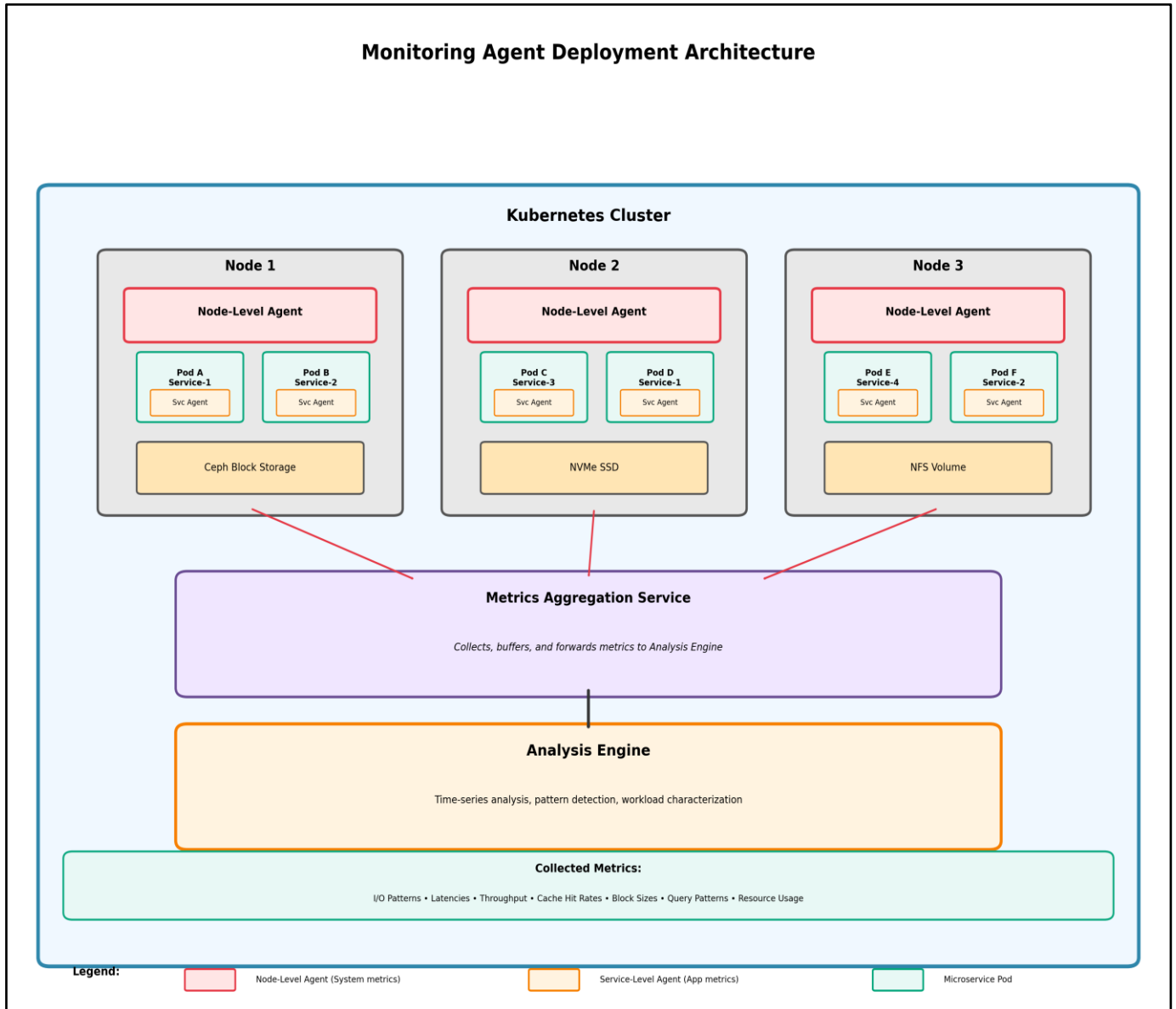


Fig 3 Monitoring Agent Deployment Architecture

V. IMPLEMENTATION DETAILS

➤ *Monitoring Agent Implementation*

The monitoring agents are implemented using a combination of eBPF programs for kernel-level tracing and user-space collectors for application metrics. For storage I/O monitoring, we utilize eBPF kprobes attached to block layer functions to capture read/write operations, latency distributions, and access patterns. The agents employ adaptive sampling rates that increase during periods of unusual activity and decrease during stable operation to minimize overhead. Collected metrics are buffered locally and transmitted to the Analysis Engine through a lightweight protocol based on Protocol Buffers for efficient serialization.

➤ *LLM Integration and Prompting Strategy*

The LLM Decision Core utilizes a carefully crafted prompting strategy that structures the decision-making process into discrete reasoning steps. The prompt includes

sections for workload analysis, storage technology comparison, configuration parameter reasoning, and cost-benefit analysis. We employ few-shot learning with representative examples from the Benchmark Repository to guide the LLM toward consistent, high-quality recommendations. The system uses temperature sampling with a relatively low value to balance creativity with deterministic decision-making. Generated recommendations are parsed through a structured output format that ensures all required fields are present and properly validated.

➤ *Configuration Parameter Optimization*

For existing deployments, IntelliStore can perform in-place configuration optimization. When a service is already using a particular storage backend such as Ceph, the system analyzes current performance metrics and identifies tunable parameters that would yield the greatest improvement. For example, in Ceph deployments, the system may recommend adjusting block sizes based on

the predominant I/O pattern, modifying OSD thread pool sizes for better concurrency, or tuning cache parameters to match working set sizes. The optimization process uses a combination of analytical models and empirical testing in staging environments to validate recommendations before production deployment.

➤ *End-to-End Pipeline Analysis*

For microservices organized in processing pipelines, IntelliStore performs holistic optimization that considers data flow patterns across service boundaries. The Analysis Engine constructs a directed graph representing data dependencies and analyzes metrics across the entire pipeline. This enables identification of bottlenecks that may not be apparent when examining individual services in isolation. For example, if an upstream service generates data with certain characteristics, downstream services can be configured to optimize for those patterns. The LLM considers these inter-service dependencies when generating recommendations, potentially suggesting different storage technologies for different pipeline stages based on their specific roles and access patterns.

VI. EXPERIMENTAL EVALUATION

➤ *Experimental Setup*

We evaluated IntelliStore on a production Kubernetes cluster with 50 nodes running 120 microservices across three distinct application pipelines.

Table 1 Performance Comparison Between Baseline and IntelliStore Configurations

Metric	Baseline	IntelliStore
Storage Cost Reduction	0%	34%
I/O Latency Improvement (p95)	12.4ms	8.9ms (-28%)
Throughput Improvement	8,500 IOPS	11,200 IOPS (+32%)
Configuration Time Reduction	24 hours	14 hours (-42%)
SLA Compliance	96.2%	99.1%

➤ *Case Study: Ceph Block Size Optimization*

One particularly interesting result came from optimizing a microservice using Ceph block storage. The baseline configuration used the default 4MB block size. IntelliStore's analysis of I/O patterns revealed that 78% of operations were small random reads averaging 64KB. The LLM recommended reducing the block size to 256KB and adjusting the cache tier configuration. This resulted in a 45% reduction in read latency and a 22% improvement in cache hit rate. The configuration change was validated in a staging environment before production deployment, confirming the predicted improvements with minimal risk.

➤ *Pipeline Optimization Results*

For end-to-end pipeline optimization, we observed particularly strong results. In a data processing pipeline consisting of ingestion, transformation, and aggregation services, IntelliStore recommended heterogeneous storage configurations tailored to each stage. The ingestion service was configured with high-throughput object storage, the transformation service used local NVMe SSDs for intermediate results, and the aggregation service employed a time-series database optimized for

The cluster included a mix of CPU-optimized and memory-optimized instance types with various storage backends including Ceph block storage, NFS shared volumes, and local NVMe SSDs. Workloads ranged from high-throughput data ingestion services to latency-sensitive API endpoints and batch processing pipelines. We compared IntelliStore recommendations against manually configured baselines established by experienced DevOps engineers over a six-month period.

➤ *Performance Metrics*

Our evaluation focused on four key metrics:

- Storage cost efficiency measured as dollars per terabyte-month with consideration for performance tiers
- I/O performance measured through 95th percentile read and write latencies and throughput in IOPS
- Configuration time measured as engineer-hours required to research, test, and deploy storage configurations
- SLA compliance measured as the percentage of requests meeting latency targets

VII. RESULTS

Table 1 Summarizes the Performance Improvements Achieved by IntelliStore Across our Test Microservices:

analytical queries. This holistic approach yielded a 52% reduction in end-to-end processing time compared to a homogeneous storage configuration, demonstrating the value of pipeline-aware optimization.

VIII. DISCUSSION

➤ *Key Insights*

Our evaluation of IntelliStore revealed several important insights. First, the value of combining real-time metrics with research literature cannot be overstated. The LLM's ability to reason about trade-offs based on published benchmarks enabled it to make recommendations that considered both theoretical performance characteristics and empirical evidence from similar deployments. Second, the importance of end-to-end pipeline analysis became evident in our results. Individual service optimization often leads to suboptimal global configurations when inter-service data flows are not considered. Third, the system's ability to adapt to changing workload patterns proved crucial for maintaining optimal configurations over time.

➤ *Limitations and Challenges*

Despite promising results, IntelliStore faces several limitations. The quality of recommendations depends heavily on the availability and accuracy of monitoring data. Services with insufficient instrumentation may receive suboptimal suggestions. The system's reliance on LLMs introduces potential variability in recommendations, though we mitigated this through structured prompting and validation. Migration to new storage backends can be complex and risky, requiring careful planning and testing. Our current implementation focuses on storage configuration but does not address all aspects of database schema design, which remains partially manual. Finally, the benchmark repository requires continuous curation to remain current with evolving storage technologies.

➤ *Generalizability*

While our evaluation focused on Kubernetes-based microservices, the core principles of IntelliStore generalize to other cloud-native architectures. The monitoring and analysis components could be adapted to serverless platforms, edge computing environments, or hybrid cloud deployments. The LLM-based decision engine is agnostic to specific technologies and can reason about any storage system for which adequate documentation and benchmarks exist. However, practical deployment would require platform-specific integration work and validation of recommendations against platform constraints.

IX. FUTURE WORK

➤ *Several Promising Directions Exist for Extending IntelliStore:*

- Integrating automated schema design capabilities that leverage the LLM's understanding of data access patterns to suggest optimal database schemas and indexing strategies
- Developing predictive models that can anticipate future storage requirements based on historical trends and upcoming feature deployments
- Implementing automated A/B testing frameworks that can safely validate recommendations in production by gradually migrating traffic between configurations
- Extending the system to handle multi-cloud and hybrid cloud environments with considerations for data sovereignty, network costs, and cross-region replication
- Developing specialized agents for specific storage technologies that can provide deeper optimization beyond general recommendations
- Incorporating cost optimization beyond storage efficiency to include compute and network costs associated with different storage configurations
- Exploring the use of reinforcement learning to continuously refine recommendation quality based on observed outcomes in production deployments

X. CONCLUSION

This paper presented IntelliStore, an intelligent AI agent framework for autonomous storage and database optimization in cloud-native microservices environments. Our system addresses the growing complexity of storage configuration in modern distributed systems by combining real-time monitoring, LLM-powered decision-making, and knowledge synthesis from research literature. Through comprehensive evaluation on production workloads, we demonstrated significant improvements in storage costs, I/O performance, configuration time, and SLA compliance compared to manual optimization approaches.

The key contributions of this work include a scalable monitoring architecture that captures storage usage patterns with minimal overhead, an LLM-based decision engine that generates optimal recommendations by reasoning about trade-offs and constraints, and a holistic approach to pipeline optimization that considers inter-service dependencies. Our results show that AI-driven automation can significantly enhance resource efficiency while reducing the operational burden on DevOps teams.

As cloud infrastructure continues to grow in complexity and scale, intelligent automation tools like IntelliStore will become increasingly essential for managing the operational challenges of modern distributed systems. By leveraging the reasoning capabilities of large language models and the wealth of knowledge available in research literature, we can build systems that not only react to current conditions but proactively optimize for future requirements, ultimately enabling more efficient and cost-effective cloud deployments.

REFERENCES

- [1]. Delimitrou, C., & Kozyrakis, C. (2014). Quasar: resource-efficient and QoS-aware cluster management. In Proceedings of ACM ASPLOS.
- [2]. Ma, L., et al. (2018). Query-based workload forecasting for self-driving database management systems. In Proceedings of ACM SIGMOD.
- [3]. Gan, Y., et al. (2019). An Open-Source Benchmark Suite for Microservices and Their Hardware-Software Implications for Cloud & Edge Systems. In Proceedings of ACM ASPLOS.
- [4]. Gias, A.U., et al. (2019). ATOM: Model-Driven Autoscaling for Microservices. In Proceedings of IEEE ICDCS.
- [5]. Kannan, R.S., et al. (2019). GrandSLAM: Guaranteeing SLAs for Jobs in Microservices Execution Frameworks. In Proceedings of ACM EuroSys.
- [6]. Mahgoub, A., et al. (2020). OPTIMUS-CLOUD: Heterogeneous Configuration Optimization for Distributed Databases in the Cloud. In Proceedings of USENIX ATC.
- [7]. Mirhosseini, A., et al. (2021). Parslo: A Gradient Descent-based Approach for Near-optimal Partial

- SLO Allotment in Microservices. In Proceedings of IEEE ICDE.
- [8]. Chow, K., et al. (2022). DeepRest: Deep Resource Estimation for Interactive Microservices. In Proceedings of ACM EuroSys.
 - [9]. Hossen, M.R., et al. (2022). Practical Efficient Microservice Autoscaling with QoS Assurance. In Proceedings of ACM HPDC.
 - [10]. Luo, S., et al. (2022). The power of prediction: microservice auto scaling via LSTM. In Proceedings of the ACM Symposium on Cloud Computing.
 - [11]. Park, J., et al. (2022). Improved Q Network Auto-Scaling in Microservice Architecture. *Applied Sciences*, 12(3), 1206.
 - [12]. Hu, X., et al. (2023). LLM As DBA. In Proceedings of the ACM SIGMOD International Conference on Management of Data.
 - [13]. Sarma, S.K., et al. (2023). An Auto-Scaling Approach for Microservices in Cloud Computing Environments. *Journal of Grid Computing*, 21(4), 1-25.
 - [14]. Abdullah, M., et al. (2024). Auto-Scaling Techniques in Cloud Computing: Issues and Research Directions. *Sensors*, 24(17), 5551.
 - [15]. JuiceFS Team (2024). LLM Storage Selection & Detailed Performance Analysis of JuiceFS. Technical Report.
 - [16]. Omdia (2024). Optimizing AI Data Storage Management. Technical Analysis Report.
 - [17]. Zhang, L., et al. (2024). A Survey of LLM × DATA: Bridging LLMs and Data Management. arXiv preprint arXiv:2505.18458.
 - [18]. Zhang, Y., et al. (2024). Blueprinting the Cloud: Unifying and Automatically Optimizing Cloud Data Infrastructures with BRAD. In Proceedings of VLDB.