

Secure and Cost-Efficient Deployment of Data-Intensive AI Workloads in Cloud Platforms

Balasubramanian Bava Jagannathan

Abstract: Cloud infrastructure remains the primary deployment platform for data-hungry AI pipelines, with elastic compute and managed storage allowing rapid provisioning at scale. Yet engineering production-grade deployments remains a poorly solved problem. The objectives on performance, cost, security, and operational readiness are tightly coupled with each other, but the existing deployment frameworks optimize them separately and only satisfy one objective at a time. End-to-end AI pipelines entail diverse data ingestion/transformation, feature generation, model training, batch inference, online serving, and continuous monitoring, which exhibit heterogeneous resource utilization and scaling. A monolithic deployment strategy cannot simultaneously meet the needs of various components. We propose a framework that proactively realizes data locality, elastic resource allocation, governance-aware isolation, and observability readiness at design time as opposed to applying these concepts post-deployment. It generates candidate deployment plans in placement and scaling dimensions from stage-level workload characteristics, filters them using policy and observability feasibility gates, and emits run-time readiness artifacts for auditability and reliable run-time operations. The problem is framed as a constrained multi-objective optimization. The parameters of interest are the tail latency, total cloud cost, and a surrogate for operational risk, which incorporates the exposure surface and the blast radius. The trade-offs between data localization, elasticity, and governance are investigated, and it is shown that joint planning can reveal deployment options overlooked by several performance- and cost-first baselines.

Keywords: *Cloud Deployment, Data-Intensive Workloads, Artificial Intelligence Pipelines, Governance-Aware Optimization, Operational Observability*

I. Introduction

Cloud platforms provide the elastic compute, the managed storage, and the operational tooling necessary to build modern data pipelines and AI. Productionizing data and AI workloads is more than just provisioning infrastructure. In addition, the same pipeline must satisfy latency SLOs, control cloud spending, limit data exposure, and gracefully handle demand variability and partial system failures across its data lifecycle [1]. When pipelines include multiple services (e.g., object stores, stream processors, GPU clusters, and real-time serving endpoints), the relationships among these properties are not easy to reason about individually.

Data movement is another potential cause of fragility. Training workloads typically involve numerous passes over large datasets. The high throughput of online inference requires low-latency access to the feature stores and model artifacts. Batch inference needs to finish before a deadline at the lowest transfer cost. Each stage has a different data access pattern. It interacts in unpredictable ways with the price of cloud resources, the topology of the node network, and the caching behavior if not profiled at a fine granularity [2]. Furthermore, workload heterogeneity is a challenge, with GPU-bound training and CPU-bound preprocessing and memory-bound serving all occurring in the same pipeline.

Security and governance boundaries may also have implications: in some regulated or sensitive-data scenarios, least-privilege access, audit logging, and minimizing data-plane exposure may be required properties of the architecture rather than hardening tasks performed on an imposed architecture. Security as an orthogonal layer of concern leads to a functioning but non-compliant architecture. The blast radius is increased with overprivileged service accounts being compromised by an attacker. Misconfigurations rather than cryptographic flaws are responsible for the majority of cloud data exposure incidents [7].

Reliability is not an accident. Projects that do not have standardized telemetry, SLO monitors, and incident response runbooks will amass unacceptable operational debt, regardless of their performance characteristics. Automated RCA and proactive anomaly detection may help reduce incident response time [9], but they require observability infrastructure to be deployed from the start, rather than in response to incidents.

Despite the progress in cloud ecosystems, data engineering, and MLOps, we lack systems guidance for deployments that anticipate data locality and access patterns, elastic resource allocations, security-aware isolation and governance, and operations reliability. Further, any one of these optimizations typically results in plans that are infeasible under the other constraints. To fill this gap, this paper proposes a deployment model, a formal problem definition, and an evaluation approach

based on the performance, cost, and operational risk dimensions.

II. Background and Related Work

A. Cloud-Based Data Analytics and AI Workload Management

The resource access patterns of production AI workloads at cloud scale are very different from research workloads in a small cluster. In a large-scale analysis of a production MLaaS system with over 7.5 million task instances, we find that these workloads are highly skewed: 77% of task instances are submitted by the 5% of users who each run over 17,500 instances [3]. However, because the median run time of a program was 23 minutes, with a variance of four orders of magnitude, queuing delay became a major factor in the end-to-end completion time of a large percentage of jobs rather than raw compute throughput.

Additionally, our trace showed a disproportionate number of scheduling delays for short-running instances: 9% of the short-running instances spent more than 50% of their total completion time in the scheduling queue compared to only 3% of the long-running instances. [3]. For example, we observed that the 90th percentile queuing time to use the V100 machines used for the high-end workload was 13,709 seconds, compared to a 360-second median for the low-end workload. For example, stage-level profiling captures not only compute intensity but also affinity for different types of GPUs and contention when scheduling a stage. Likewise, GPU locality enforcement that co-runs worker instances on machines with high-speed interconnects resulted in over a 10× speedup for some large-scale classification training workloads [3], providing further evidence that placement effects not captured by cost and throughput models can be substantial.

Intensive work on GPU-based training has shown that communication overheads impose non-trivial dependencies on both the network topology and the batch size, and that distributed training topology, rather than simply per-node resource allocation, must be a factor in deployment considerations to achieve predictable performance [2].

B. Resource Management and Scaling for AI Inference

Resource allocation for inference serving differs from training. Proteus, a high-throughput inference serving system, was evaluated for a cluster of 20 CPU workers, 10 GTX 1080 Ti GPU workers, and 10 V100 GPU workers showed that static allocation at deployment time consistently underperforms with diurnal variations in demand [5]. Compared to a non-scaling baseline (maximizing throughput on model part variants with low precision), Proteus improved throughput by 60% and

reduced SLO violations by >10× due to the MILP-based dynamic model placement and adaptive batching. Its maximum accuracy degradation at peak demand was 4.85%, which is 2.8× lower than the greedy heuristic baseline, and 3.2× lower than the partially dynamic baseline [5].

Proteus's adaptive batching mechanism also directly addresses the tension between throughput and tail-latency for SLOs. Proteus avoids SLO violations that would exist when reactive techniques are used since it executes batches just as the first request in its queue is about to timeout. Proteus avoids having the MILP solver be on the critical path for serving an incoming query by solving the resource allocation problem offline every 30 seconds, in a manner similar to the one in this work.

Clockwork, on the other hand, tackles the same inference-serving problem but from a different angle. Instead of dynamically scaling resource allocation, it attempts to eliminate performance variance by scheduling only at a centralized controller [6]. Clockwork shows up to 100× lower variance in tail latency from a single INFER action per GPU, at only a small drop in throughput. Over six hours and 208 million requests with 4026 model instances replaying a real network trace, only 58 requests had a mispredicted timing. Notably, no mispredicted request missed our 100 ms SLO [6]. In fact, the 99th percentile error on the prediction of INFER was only 144 microseconds over, and 55 microseconds under the real execution time, which indicates that predictable execution at scale is possible when scheduling is not done in lower layers.

C. Cloud Security, Governance, and Data Exposure

Security in cloud-native AI applications is from protecting the service boundary but also involves identity management, managing data access patterns and service-level isolation boundaries [7]. Minimum viable governance frameworks for machine learning pipelines are best applied in the service's deployment toolchain rather than in an after-the-fact audit mode [8]. If security is delayed until post-deployment review, then misconfigurations progress through to production, expanding the production blast radius for a potential future incident. Instead, security must be modeled as a feasibility constraint on deployment decisions, rather than a score variable to be balanced against latency, cost, and other objectives.

D. Operational Intelligence and Observability

To detect and analyze problems, telemetry should be collected in a systematic way from the instant of service deployment. An automated root cause analysis (RCACopilot) system deployed across production services handling 150 billion messages a day showed that

multi-source diagnostics (error logs, exception stack traces, and socket metrics) can be more effective than single-source diagnostics. [9]. When provided with the summarized diagnostic information and using large language model-based chain-of-thought reasoning, the system achieves a Micro-F1 of 0.766 and a Macro-F1 of 0.533, outperforming all baselines, including fine-tuned models and prompting GPT-4 without structured context [9]. The system has been deployed on over 30 operational teams, and the time to collect diagnostic information for an incident varies from 15 to 841 seconds, depending on the complexity of the system involved.

Multi-modal AI techniques that use execution logs, network traffic, and API usage metrics for cloud-native security monitoring lead to a 35% reduction in detection time and a 22% increase in the F1-score for detecting attacks compared to single-modality baselines [10]. The

multi-modality framework achieved a precision of 95% versus 82% for the single-modality, and a recall of 93% versus 78% for single-modality frameworks [10]. These results indicate that observability depth, the correlation of telemetry across heterogeneous tooling determines the operational intelligence that can be produced during and after an incident.

E. Gap Summary

All the above works target workload characterization, inference serving, governance policy, and operational intelligence independently, but none co-designs data locality, scaling policies, governance constraints, and observability requirements as a joint objective during the deployment planning stage. However, the conceptual framework presented in subsequent sections addresses this gap.

Sub-Area	Key Focus	Key Findings / Metrics
Cloud-Based AI Workloads	Task distribution and scheduling	77% tasks from 5% users; median runtime 23 min; 90th percentile queue delay up to 13,709 sec
GPU Training Optimization	Placement and communication overhead	GPU locality yields >10× speedup; topology impacts performance
Inference Resource Management	Dynamic vs static allocation	Throughput +60%; SLO violations reduced >10×; accuracy drop limited to 4.85%
Centralized Scheduling	Predictability in inference	100× lower latency variance; <100 ms SLO maintained
Security & Governance	Deployment-time enforcement	Misconfigurations dominate exposure risks
Observability & RCA	Multi-source telemetry	Micro-F1: 0.766; detection latency reduced by 35%; precision up to 95%

Table 1. Overview of Background Studies and Key Findings in Cloud-Based AI Workloads [2, 3, 6, 7, 9, 10].

III. System Model and Problem Formulation

A. Cloud Environment Model

We model the cloud as a set of service classes of compute instances (serverless, container and VM), managed storage (object store, block store and distributed cache), a network fabric, an egress pricing scheme, an identity and access management plane and an observability stack (logs, metrics and distributed traces). The defining metrics for each class of service are unit cost and latency scaling, maximum throughput, and isolation boundary.

B. Workload Model

W is an ordered set of stages: $S = \{s_1, s_2, \dots, s_n\}$. Each stage corresponds logically to a stage in a data pipeline, such as ingest, transform, featureization, train, validate, serve, and monitor. Further, production workload traces confirm the stage-level heterogeneity such as for CTR prediction where the GPU utilization is different for training and inference. The median GPU usage for

inference CTR prediction workloads is less than 0.1 despite a high CPU usage [3]. Each stage s_k has its profile ϕ_k containing execution mode, input data volume and access pattern, compute intensity and preferred accelerator type, latency SLO, burstiness coefficient, and data sensitivity class.

C. Deployment Plan

A deployment plan Π consists of a placement p_k , a resource allocation and scaling policy r_k , a data access mapping a_k , and an isolation and access-control policy g_k . At each replica group k , the placement p_k specifies a region, an availability zone, and a service tier. The resource allocation and scaling policy r_k specifies an instance type, a minimum and maximum replica count, and a set of scaling policies. The data access mapping a_k specifies one of direct access, regional cache, selective replication, or partitioned locality. The isolation and access-control policy g_k specifies least-privilege, network segmentation, and per-stage credential scoping.

D. Objective Functions and Constraints

The deployment problem can be formalized as follows:

Minimize: $f(\Pi) = (C(\Pi), L(\Pi), R(\Pi))$

where $C(\Pi)$ is the total cloud cost including compute, storage, and network transfer; $L(\Pi)$ is the 95th percentile tail latency of the online inference stages and the end-to-end latency of the batch and training workloads; and $R(\Pi)$ is a risk proxy defined as the product of exposure surface and blast radius. The feasible region is defined by three feasibility constraints:

(C1) Governance feasibility: Each stage's access-control configuration satisfies least-privilege, bounded exposure, and audit-readiness requirements.

(C2) Observability feasibility: Includes telemetry configuration (logs, metrics, traces) above the coverage thresholds needed for reporting SLOs and helping incident response.

(C3) SLO feasibility: The expected $L(\Pi)$ under the design-load scenarios meets declared SLOs for all latency-sensitive stages.

The framework determines Pareto-efficient plans with respect to (C, L, R) subject to $(C1) - (C3)$.

Component	Description	Key Attributes
Cloud Environment	Compute, storage, network, IAM, observability stack	Cost, latency scaling, throughput, isolation
Workload Model (W)	Pipeline stages (ingest → monitor)	Stage profile φ : data volume, compute type, SLO, burstiness
Deployment Plan (Π)	Placement, scaling, access, governance	Region, instance type, scaling rules, access policy
Objective Functions	Optimization targets	Cost (C), Latency (L), Risk (R)
Constraints	Feasibility requirements	Governance (C1), Observability (C2), SLO (C3)

Table 2. System Model Components and Formal Definitions of Deployment Framework [2, 3, 7.

IV. Proposed Framework

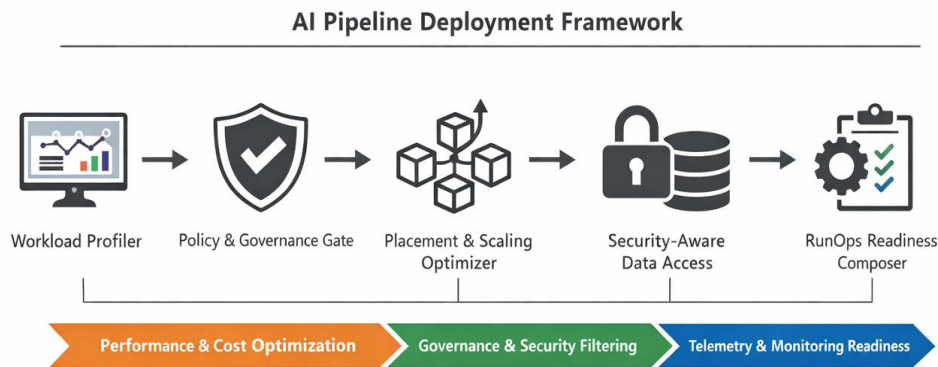


Figure 1. Workflow of the Proposed AI Pipeline Deployment Framework

A. Workload Profiler

The Workload Profiler extracts a stage-level profile $\varphi = \{\varphi_1, \varphi_2, \dots, \varphi_n\}$ from the workload specification and telemetry from completed or shadow executions. Production traces show that most tasks in a distributed cluster are of a repeated nature: for example, over 65% of the diverse tasks are executed at least five times, and their run-time can be predicted from group tag, user name, and resource requests [3]. The duration prediction error is less than 25% for 78% of the cases considered, which is a

cutoff determined by other scheduling works to be sufficient for high-quality placement decisions [3].

We classify each stage as one of four archetypes: large-sequential for bulk scans for pre-training, random-read for lookup-heavy feature extraction during online inference, stream-append for streaming ingest with ordering guarantees (e.g. data ingest), and mixed for read- and write-intensive iterative access patterns used during validation and monitoring stages of model training. These archetypes determine the data-locality optimizations evaluated in the optimization phase. Two additional

insights captured in the production traces (GPU sharing with fractional GPU allocation) can improve profiling. On average and under the worst-case peak hour load, GPU sharing reduces the number of GPUs required by 50% and 73%, respectively [3]. Hence, the profiler can focus on placement scenarios that tolerate sharing within SLOs.

B. Policy and Governance Gate

The Policy and Governance Gate encodes the organizational and regulatory constraints that any compliant deployment plans must satisfy. The governance rules enforced by the Policy Gate are least-privilege credential binding, bounded exposure surface, data residency compliance, and audit readiness, which is the requirement that all data-plane access events must be immutably logged in a centralized manner [7].

This observability requirement is defined by a gate that expects minimum telemetry in the form of structured logs (e.g., request IDs, latency, and error codes), metrics (resource-related and SLO rather than implementation-focused counters), and distributed traces (the critical path through each latency-sensitive stage). Operational experience has shown (in deployment) that multi-source diagnostic collection (error logs, exception stack traces, and system metrics correlated together) provides better root-cause analysis than any of these sources alone [9]. Plans that cannot be instrumented to the required coverage threshold due to the constraints of the available service model are pruned before optimization.

C. Placement And Scaling Optimizer

The Placement and Scaling Optimizer generates candidate deployment plans by composing stage-level placement options, scaling policies and data-locality strategies. Placement candidates are enumerated based on the Cartesian product of service tiers and region/zone pairs that are compatible with the data-residency constraints returned by the governance gate. Candidates for scaling policies include static provisioning, rule-based autoscaling policies, and predictive autoscaling based on burstiness coefficients computed by the profiler. Data-locality policies include direct access, cached access, selective replication, and partitioned locality.

Each candidate plan's (C, L, R) triplet is scored with analytical models of the cloud provider's pricing, network topology, and the sample execution profiles. In particular, the latency model divides the stage's latency into compute and data transfer latencies based on the data size, access pattern, and locality strategy of data transfer. The model evaluates the risk based on the exposure surface and the

blast radius generated under each placement and access-control configuration and returns the Pareto frontier and a weighted-optimal recommendation based on specified priority weights.

Based on data collected from running dynamic inference serving systems, we show that MILP allocation outperforms greedy heuristics; e.g., at peak load, the drop in accuracy is at most 4.85% with dynamic allocation versus 13.7% with greedy allocation. Additionally, worst-case SLO violations are a factor of $4.3\times$ lower with dynamic allocation. This motivates the use of exact multi-objective search on the candidate plan space instead of a heuristic approximation.

D. Security-Aware Data Access Patterns

The data plane should be exposed as little as possible, such that direct access is used instead of replication if latency requirements can be satisfied without caching. Caching is only scoped to the minimal data required by each stage with automatic expiry. Selective replication is only available when profiling clearly shows that transfer latency materially impacts SLO attainment. Per-stage credential scoping limits the blast radius of any service identity that is compromised [7]. The rules guarantee that the optimizer will not accept a plan that reduces latency or cost by increasing the exposure surface beyond the governance feasibility threshold.

E. RunOps Readiness Composer

The RunOps Readiness Composer takes the selected deployment plan and emits operational artifacts to make the deployment observable, auditable, and recoverable by design: Telemetry baseline to each stage's SLO and criticality class, SLO monitoring with alerting thresholds and burn-rate rules for the SLO error budgets, audit hooks for logging data-plane access events in a centralized audit store, and operational runbook pointers to remediation playbooks for each stage.

When structured, the multi-source telemetry dataset is of action-oriented use. A model that fused execution logs, exception traces and socket metrics achieved a Micro-F1 score of 0.766 and a Macro-F1 score of 0.533, bettering all other baseline models based on each source alone and direct-prompting [9]. Multi-modal telemetry fusion in cloud security monitoring further improved the F1-score by 22% and the detection latency by 35% compared to single-modality models [10], confirming that such upfront choices in the design of telemetry architecture determine the operational intelligence available to investigators during an incident.

Component	Function	Key Contributions
Workload Profiler	Extracts stage-level characteristics	Predicts runtime (<25% error in 78% cases); identifies workload archetypes

Policy & Governance Gate	Enforces compliance constraints	Ensures least-privilege, audit readiness, data residency
Placement & Scaling Optimizer	Generates and evaluates plans	Uses multi-objective optimization; reduces SLO violations by 4.3×
Data Access Controller	Defines locality strategies	Balances latency vs exposure; limits replication
RunOps Readiness Composer	Generates operational artifacts	Enables observability, auditability, RCA readiness

Table 3. Functional Components of the Proposed Deployment Framework and Their Roles [3, 5, 7, 9, 10].

V. Experimental Evaluation Methodology

A. Workload Scenarios

Offline training is a multi-epoch deep learning job that performs very large sequential reads on a partition of the object store while the job checkpoints periodically. The trade-offs are between a data-locality strategy and storage cost. A second focus is batch inference, in which a pre-trained model is applied to a batch of inputs within a deadline-bound fashion: the trade-off has been studied most often in the context of an SLO on completion time vs. parallelism, though the third use case is online inference: request serving with bursty arrivals and p95 tail latency SLOs. The general trade-off is between the number of replicas and expected peak demand latency.

When running production workloads on heterogeneous GPU clusters, instance run-times vary four orders of magnitude when instance-based training and inference jobs share a cluster and 79% of total GPU utilization are via gang-scheduling [3]. In contrast, batch inference jobs using pre-trained models have predictable instance run-times and can be scheduled with deadlines if they need to be re-run [3].

B. Baselines

Three baselines are considered. In compute-centric placement, stages are assigned to the lowest cost compute node without considering data locality, governance, and observability requirements. Data-locality-only placement minimizes data transfer latency using aggressive data placement and replication. However, it does not model

governance exposure and cost. Rules-based autoscaling applies manual governance and observability steps after scaling. It uses rules that have predefined (> 0) thresholds for CPU and memory utilization.

C. Metrics

Common metrics for online inference include the SLO violation ratio (the number of SLO violations over the number of queries), p95 tail latency [5], and evidence completeness (the fraction of required governance artifacts that are machine-generated and version controlled). Observability coverage is the percentage of pipeline stages instrumented with logs, metrics, and traces above a minimum threshold. The readiness metrics are independent of performance and cost so that the feasibility gates do not get relaxed due to optimizations before deep debugging is possible.

D. Protocol

For each workload scenario we profile with the Workload Profiler, we compute candidate plans using the Placement and Scaling Optimizer and filter them using the Policy and Governance Gate. We then deploy the resulting plan and the three base plans to our experimental infrastructure, and run controlled load experiments using the workload scenario's load profile. We collect workload performance, cost, and readiness telemetry for analysis over the (C, L, R) design space. For systems too expensive to deploy at full scale, we fall back on trace-driven simulation, validated by our previous systems work [6], and the small-scale empirical measurements we obtain.

Category	Aspect	Details
Workload Types	Training, Batch, Online Inference	Heterogeneous GPU/CPU workloads; bursty and deadline-driven
Baselines	Comparison strategies	Compute-centric, locality-only, rules-based autoscaling
Performance Metrics	Latency, throughput	p50/p95/p99 latency; completion time
Cost Metrics	Resource usage	Compute, storage, egress costs
Readiness Metrics	Observability & governance	Coverage %, audit completeness, SLO readiness
Evaluation Method	Deployment & simulation	Controlled experiments + trace-driven simulation

Table 4. Experimental Setup, Baselines, and Evaluation Metrics for Workload Analysis [3, 5, 6].

VI. Discussion

A. Integrated Decision-Making Under Constraints

The primary output of this problem formulation is that for many instances, the latency/cost-optimal deployment plans are not feasible in the presence of governance or observability constraints and cannot be deployed in production in regulated or sensitive-data environments. Filtering through governance and observability constraints prior to scoring for latency or cost yields an actionable Pareto frontier.

Experiments on production inference serving have shown that static resource allocation to optimize accuracy and ignore throughput may lead the overall throughput to be half of the query demand at diurnal peaks at best [5]. In contrast, optimizing for throughput only (regardless of the model's accuracy) incurs at most an accuracy drop of 20% relative to demand [5]. These results are generally not sustainable for production settings where both accuracy SLOs and throughput SLOs must be satisfied. The Pareto frontier structure makes this trade-off more explicit than if there were only one potential operating point.

B. Trade-Off Characterization

Three canonical trade-offs appear across all workloads. The first is locality vs. cost/exposure, which states that although a replication can reduce stage-level latency and intra-region transfers by increasing data locality, this improves locality at the cost of increased storage and exposure. The governance gate disallows plans with too much replication (above some exposure threshold), preferring caching or sharding. The elasticity/tail latency tension occurs because aggressive scale-in minimizes per second compute cost but incurs increased p95 and p99 tail latency for bursty requests. Production evidence confirms predictive scheduling (task run time distributions with prediction error < 25% in 78% of cases) lowers average task completion time 63-77% compared to first-in-first-out baselines [3]. Third, we want observability depth versus overhead: instrumentation can reduce incident response time but incurs the overhead cost of ingesting and storing additional data. RunOps Readiness Composer limits the coverage of telemetry to the minimum for each stage's SLO and class of criticality to avoid over-instrumentation.

C. Observability as a Deployment Prerequisite

Observability coverage is not a post-deployment enhancement but a feasibility gate. Strong operational evidence exists that root cause analysis systems that fuse and reason over structured diagnostic information from multiple systems achieve a state-of-the-art Micro-F1 of 0.766, compared to 0.103 from fine-tuned models from working with alerts only, lacking any contextual or structured diagnostic information. If, for example,

execution logs, network traffic and API usage are fused a precision of 95% and a recall of 93% can be achieved, compared with a precision of 82% and a recall of 78%, respectively, for each modality. Furthermore, fusion reduces the latency from 1.2 seconds to 0.78 seconds [10]. Hence, all performance differences can be attributed to the choices made at deployment time on the telemetry architecture, which supports the framework's assertion that observability infrastructure has to be planned.

D. Limitations and Future Work

Cost- and latency-models in the framework are constructed based on properties of workloads and pricing schemes, which results in degradation of the framework's prediction performance for workloads with high access pattern irregularity and for services with rapidly varying pricing schemes. In future work, we will respond to data by calibrating the inputs with billing actuals on a regular cadence and by dynamically updating model parameters. We could also relax the governance policies to soft preferences (e.g., prefer low exposure when both a high exposure and a low exposure plan are feasible). There are currently no standardized end-to-end benchmarks for Data+AI pipeline deployments, and there is a need for reproducible workload generators and reference telemetry baselines for verifying future framework extensions and allowing for comparisons between frameworks.

Conclusion

Production deployments of data-hungry AI pipelines are challenged by the need to meet performance, cost, governance, and observability targets, which are typically optimized separately within existing systems. By incorporating governance feasibility and observability coverage as hard constraints in the planning formulation, the proposed framework treats deployability as a first-class design aim. This guarantees that all plans on the Pareto frontier are both efficient and policy-compliant. The Workload Profiler captures resource access patterns and scaling behavior at the stage level using evidence-based profiling methods. The Policy and Governance Gate denies the progression of plans that do not meet least-privilege, data residency, or audit-readiness requirements before the cost and latency scoring. The Placement and Scaling Optimizer employs exact multi-objective search to output Pareto-efficient plans over the deployment cost, latency, and operational risk space. The RunOps Readiness Composer emits machine-generated telemetry configurations, SLO monitors and audit hooks to make deployments observable and recoverable from birth. The evaluation protocol covers training and batch and online inference scenarios through a set of performance, cost, and readiness-related metrics that provide a thorough view of deployment quality. Future work should extend to cost model calibration, soft

governance preference modeling, and end-to-end benchmarks for reproducible cross-framework comparisons.

References

- [1] Syed Nyamtulla and Dr. Dharendra Kumar Tripathi, "Serverless vs Traditional Cloud Architectures: Performance and Cost Evaluation of AI/ML Workloads in HPC Environments," International Research Journal of Engineering & Applied Sciences, 2025. [Online]. Available: <https://www.irjeas.org/wp-content/uploads/admin/volume13/V13I14/IRJEAS04V13I4017.pdf>
- [2] Weizheng Xu et al., "Parallelizing DNN training on GPUs: Challenges and opportunities." Companion Proceedings of the Web Conference 2021. [Online]. Available: <https://dl.acm.org/doi/pdf/10.1145/3442442.3452055>
- [3] Qizhen Weng et al., "{MLaaS} in the wild: Workload analysis and scheduling in {Large-Scale} heterogeneous {GPU} clusters," 19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22). 2022. [Online]. Available: <https://www.usenix.org/system/files/nsdi22-paper-weng.pdf>
- [4] Juncheng Gu et al., "Tiresias: A {GPU} cluster manager for distributed deep learning," 16th USENIX Symposium on Networked Systems Design and Implementation (NSDI 19). 2019. [Online]. Available: <https://www.usenix.org/system/files/nsdi19-gu.pdf>
- [5] Sohaib Ahmad et al., "Proteus: A high-throughput inference-serving system with accuracy scaling." Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1. 2024. [Online]. Available: <https://dl.acm.org/doi/pdf/10.1145/3617232.3624849>
- [6] Arpan Gujarati et al., "Serving {DNNs} like clockwork: Performance predictability from the bottom up," 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20). 2020. [Online]. Available: <https://www.usenix.org/system/files/osdi20-gujarati.pdf>
- [7] Maria Papaioannou et al., "A survey on security threats and countermeasures in the internet of medical things (IoMT)," Transactions on Emerging Telecommunications Technologies 33.6 (2022): e4049. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/ett.4049>
- [8] IBRAHEEM ADEBAYO ADEREMI et al., "Explainable AI for Water Quality Monitoring: A Systematic Review of Transparency, Interpretability, and Trust." IEEE Sensors Reviews (2025). [Online]. Available: <https://ieeexplore.ieee.org/document/11112533>
- [9] Yinfang Chen et al., "Automatic root cause analysis via large language models for cloud incidents," Proceedings of the Nineteenth European Conference on Computer Systems. 2024. [Online]. Available: <https://dl.acm.org/doi/pdf/10.1145/3627703.3629553>
- [10] Falope Samson, "Multi-Modal AI for Serverless Cloud Security." (2026). [Online]. Available: <https://www.researchgate.net/profile/Falope-Samson/publication/403569954>
- [11] Deepak Narayanan et al., "Efficient large-scale language model training on GPU clusters using Megatron-LM," Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis. 2021. [Online]. Available: <https://dl.acm.org/doi/pdf/10.1145/3458817.3476209>
- [12] Sandra Wachter et al., "Why a right to explanation of automated decision-making does not exist in the general data protection regulation." International data privacy law 7.2 (2017): 76-99. [Online]. Available: <https://academic.oup.com/idpl/article-abstract/7/2/76/3860948?redirectedFrom=PDF>
- [13] Hongzi Mao et al., "Park: An open platform for learning-augmented computer systems." Advances in Neural Information Processing Systems 32 (2019). [Online]. Available: <https://proceedings.neurips.cc/paper/2019/file/f69e505b08403ad2298b9f262659929a-Paper.pdf>
- [14] Jasmin Bogatinovski et al., "Artificial Intelligence for IT Operations (AIOps)," Workshop White Paper, arXiv preprint arXiv:2101.06054 (2021). [Online]. Available: <https://arxiv.org/pdf/2101.06054>
- [15] JOHN OUSTERHOUT et al., "The RAMCloud storage system." ACM Transactions on Computer Systems (TOCS) 33.3 (2015): 1-55. [Online]. Available: <https://dl.acm.org/doi/pdf/10.1145/2806887>