

# The Cloud Cost-Optimization Flywheel: A Systematic Approach to Reducing Infrastructure Waste without Compromising Delivery Velocity

Janardhan Reddy Kasireddy

Reveal Global Consulting, USA



**ABSTRACT:** With increasing cloud infrastructure spending, organizations confront cost inefficiencies in security, compute, storage, and data platforms. The flywheel methodology provides an approach to continuously optimize costs using the cycle of observing, attributing, correcting, guardrailing, and validating the use and spending of cloud resources. The framework addresses the challenge of creating sustainable optimizations or short-term cost savings that went out of control due to a lack of governance policies. Optimization opportunities exist in various parts of infrastructure, such as the encryption service, which causes too many API calls due to caching issues; non-production environments that run continuously but are rarely used; databases with long support lifetimes; computing clusters that are not utilized efficiently; and storage services that retain obsolete data with no policies for managing data lifecycles. Experience shows that cost reduction and performance measures are synergistic and not inversely correlated; large waste reductions are consistent with increased velocity, system throughput, and delivery responsiveness. The transition from sporadic cost reduction to everyday engineering discipline demands a cultural foundation of constant questioning, systematic waste elimination, and a process for institutionalizing optimization that maintains the cost of suboptimal behavior through automation and governance policies.

**KEYWORDS:** Cloud Cost Optimization, Flywheel Methodology, Infrastructure Waste Reduction, Performance Engineering, Storage Lifecycle Management

## I. INTRODUCTION

The cost of cloud infrastructure can become substantial, particularly for large organizations. Flexera's 2024 State of the Cloud Report states that cloud spend optimization is the most important issue for enterprises. The report also estimates that businesses waste 28% of cloud spending on unused or underutilized resources [1]. Most of this waste comes from overprovisioned and underutilized resources, and due to a lack of visibility into the usage trends of cloud resources. Unlike a typical CAPEX model, cloud expenditure consists of a series of small recurring transactions across security, compute, storage, and data platform stacks, which means that cloud spending does not usually happen as a sudden spike. The report also found that organizations are struggling with managing multi-cloud complexity, as 89% of

enterprises have adopted a multi-cloud strategy, and tracking and optimizing cloud costs across distinct cloud providers and pricing models becomes more difficult.

Historically, cloud cost optimization was viewed as a finance problem and not an engineering problem. Atlas Systems stated that cloud cost optimization requires a holistic approach that balances performance and cost, where successful optimization initiatives integrate engineering disciplines and financial control frameworks [2]. This separation of domains results in a tradeoff between cost optimization and speed of operation to the extent that unnecessary costs are incurred by engineering teams that hesitate to engage in additional optimization due to the fear of performance degradation or delays in delivering infrastructure to end users. The company asserts that continuous monitoring, right-sizing, automation, and accountability form the four pillars of cloud cost optimization. More recently, empirical research has demonstrated that systematic cost optimization can diminish waste without sacrificing system reliability and performance when treated as a discipline of engineering rather than merely as a budget constraint.

The proliferation of cloud services has fundamentally transformed how organizations provision and consume computing resources. However, this transformation has introduced new challenges in cost governance and financial accountability. Traditional IT budgeting approaches, designed for predictable capital expenditures and depreciation schedules, prove inadequate for managing the dynamic, consumption-based pricing models inherent to cloud computing. Organizations find themselves navigating a complex landscape where hundreds or thousands of individual services, each with distinct pricing structures, combine to generate monthly bills that can fluctuate significantly based on usage patterns, architectural decisions, and operational practices.

### **1.1 Economic Impact and Business Case for Cloud Cost Optimization**

The financial implications of cloud waste extend far beyond immediate budget concerns, representing a significant drag on organizational profitability and competitive positioning. When organizations waste substantial percentages of their cloud spending on unused or underutilized resources, they effectively reduce the return on investment for their entire cloud migration and modernization initiatives. This waste manifests across multiple dimensions, including idle compute instances left running during off-hours, oversized database instances provisioned for peak loads that occur infrequently, storage volumes accumulating obsolete data without lifecycle management, and network resources maintaining connectivity for decommissioned applications.

The business case for systematic cloud cost optimization becomes compelling when examining the cumulative impact across typical enterprise environments. Organizations operating substantial cloud footprints often discover that comprehensive optimization initiatives can reduce overall cloud spending by substantial margins while simultaneously improving operational performance and system reliability. These savings translate directly to improved profit margins, increased budget availability for innovation initiatives, and enhanced financial flexibility for responding to market opportunities. The opportunity cost of unoptimized cloud infrastructure proves equally significant, as capital locked in wasteful spending cannot be allocated to strategic initiatives, including new product development, market expansion, or competitive differentiation efforts.

Return on investment calculations for cloud optimization programs typically demonstrate favorable economics with rapid payback periods. Initial optimization efforts targeting high-impact domains such as non-production environment scheduling, storage lifecycle management, and right-sizing oversized instances often achieve payback within weeks or months of implementation. These quick wins generate momentum and organizational confidence, enabling expansion into more sophisticated optimization domains requiring deeper technical analysis and longer implementation cycles. The compounding nature of optimization efforts means that organizations implementing systematic approaches experience progressively increasing returns as capabilities mature and optimization scope expands across additional infrastructure layers.

Industry benchmarks for cloud efficiency reveal substantial variation in how effectively organizations manage their cloud spending relative to actual business value delivered. High-performing organizations demonstrate cloud efficiency ratios indicating that a much higher percentage of their cloud spending directly supports productive workloads and business outcomes, with minimal waste from idle resources, oversized instances, or obsolete data retention. Conversely, organizations lacking systematic optimization approaches exhibit efficiency ratios indicating substantial room for improvement through relatively straightforward interventions. The gap between high performers and laggards in cloud cost efficiency often represents millions in annual savings opportunities for large enterprises operating significant cloud infrastructure.

The relationship between cloud spending and business growth metrics provides additional context for understanding optimization importance. Organizations experiencing rapid business growth naturally expect corresponding increases in

cloud infrastructure costs as they scale to serve expanding customer bases, process increasing data volumes, and support growing user populations. However, the rate of cloud cost growth should correlate reasonably with business growth indicators such as revenue, transaction volume, or user count. When cloud costs grow disproportionately faster than business metrics, this divergence signals inefficiency requiring investigation and remediation. Conversely, successful optimization enables organizations to achieve business growth with sub-linear cloud cost growth, improving unit economics and competitive positioning.

Table 1: Cloud Cost Optimization Challenges and Strategic Priorities [1, 2]

Challenge Category	Description	Strategic Response
Wasted Resources	Overprovisioned capacity and unused resources across the infrastructure	Continuous monitoring and right-sizing initiatives
Multi-Cloud Complexity	Tracking costs across diverse cloud platforms and pricing models	Integrated financial governance frameworks
Visibility Limitations	Insufficient insight into cloud usage patterns and spending drivers	Comprehensive cost attribution mechanisms
Organizational Alignment	Separation between financial constraints and engineering practices	Accountability structures linking technical and business objectives

## II. THE FLYWHEEL METHODOLOGY: A STRUCTURED FRAMEWORK FOR SUSTAINABLE COST OPTIMIZATION

As the cloud cost optimization problem is not about identifying the costs associated with the cloud but about establishing sustainable change, the management of cloud computing resources requires the establishment of a systematic approach for provisioning inefficiencies, resource allocation strategies, and performance monitoring systems [3]. The central recommendation for time-agnostic proactive optimization is the need for complex control loops based on real-time observability, dynamic resource allocation, and proactive auto-scaling policies. Without them, organizations risk falling into the reactive cycle of responding to cost alerts immediately by cutting costs, which can lead to pipeline breaks, performance regressions, or stakeholder dissatisfaction arising in pipelines that experience a reduction in resources. To be most effective at managing cloud computing resources, the study found, organizations should determine their resource usage, use predictive analytics to forecast future capacity, and apply governance policy to both prevent resource sprawl and meet service level expectations.

The flywheel methodology, a six-phased feedback loop for incremental improvements, addresses this issue. Work done by IEEE, for example, on finding optimal virtual machine placement and migration strategies in cloud computing environments indicates that such optimization strategies must consider workload characteristics, resource availability, energy consumption, and service level agreement requirements [4]. The work shows a reduction in operational expenses and improved resource utilization via optimal placement algorithms that dynamically allocate virtual machines based on the state of the system and expected future resource requirements. The first stage is called observation, which detects unexpected costs and identifies which services are to blame. In terms of cost, it is the task of the pattern recognition and deviation detection phase to distinguish between normal operational cost and waste. This research indicates that machine learning algorithms applied to cloud resource management can detect optimization potential. By recognizing anomalous resource usage patterns and forecasting future resource requirements, these machine learning approaches outperform standard threshold-based monitoring techniques significantly.

Attribution, step two in this method, maps spending patterns to behaviors. Attribution is critical but difficult. Surface-level cost measures often hide inefficiencies in system behavior. The IEEE report states that systems must account for the interdependencies of resources, where the decision upon one resource may have further cascade effects on other resources. For instance, inefficient application layer design or configuration may conceal the cost of encryption services as a security expense. Attribution answers the question of why a particular amount of money is spent by connecting resource usage to application behavior, user behavior, or system configuration settings, which could be the cause of consumption.

The third stage, the correction stage, is the stage in which changes are made towards addressing the root behavioral problems rather than the symptom behaviors, as superficial solutions that do not change behavior reintroduce the costs in other ways. It is important that the correction focuses on reconfiguring the systems, processes, and configurations that resulted in the wasteful patterns. The guardrails phase automates corrections and institutionalizes them through policy. Human vigilance cannot maintain improvement gains as teams rotate, knowledge spreads, and priorities shift

over time. Optimization efforts can use automated guardrails to stop old problems from coming back, making improvements last in the system. Validation, the fifth phase, measures improvements that are more than just cost. Performance metrics, throughput measurements, and reliability readings can show how the optimization step is succeeding, while the repeat step, or phase, ensures that optimization is not a project but a discipline, compounding value over time as things improve.

### **The Flywheel Methodology: A Structured Framework for Sustainable Cost Optimization**

As the cloud cost optimization problem is not about identifying the costs associated with the cloud but about establishing sustainable change, the management of cloud computing resources requires the establishment of a systematic approach for provisioning inefficiencies, resource allocation strategies, and performance monitoring systems [3] (in the International Journal of Engineering Research and Development). The central recommendation for time-agnostic proactive optimization is the need for complex control loops based on real-time observability, dynamic resource allocation, and proactive auto-scaling policies. Without them, organizations risk falling into the reactive cycle of responding to cost alerts immediately by cutting costs, which can lead to pipeline breaks, performance regressions, or Stakeholder dissatisfaction arising in pipelines that experience a reduction in resources. To be most effective at managing cloud computing resources, the study found, organizations should determine their resource usage, use predictive analytics to forecast future capacity, and apply governance policy to both prevent resource sprawl and meet service level expectations.

The flywheel methodology, a six-phased feedback loop for incremental improvements, addresses this issue. Work done by IEEE, for example, on finding optimal virtual machine (VM) placement and VM migration strategies in cloud computing environments indicates that such optimization strategies must consider workload characteristics, resource availability, energy consumption, and service level agreement requirements [4]. The work shows a reduction in operational expenses and improved resource utilization via optimal placement algorithms that dynamically allocate virtual machines based on the state of the system and expected future resource requirements. The first stage is called observation, which detects unexpected costs and identifies which services are to blame. In terms of cost, it is the task of the pattern recognition and deviation detection phase to distinguish between normal operational cost and waste. This research indicates that machine learning algorithms applied to cloud resource management can detect optimization potential. By recognizing anomalous resource usage patterns and forecasting future resource requirements, these machine learning approaches outperform standard threshold-based monitoring techniques by a factor of 10.

Attribution, step two in this method, maps spending patterns to behaviors. Attribution is critical but difficult. Surface-level cost measures often hide inefficiencies in system behavior. The IEEE report states that systems must account for the interdependencies of resources, where the decision upon one resource may have further cascade effects on other resources. For instance, inefficient application layer design or configuration may conceal the cost of encryption services as a security expense. Attribution answers the question of why a particular amount of money is spent by connecting resource usage to application behavior, user behavior, or system configuration settings, which could be the cause of consumption.

The third stage, the correction stage, is the stage in which changes are made towards addressing the root behavioral problems rather than the symptom behaviors, as superficial solutions that do not change behavior reintroduce the costs in other ways. It is important that the correction focuses on reconfiguring the systems, processes, and configurations that resulted in the wasteful patterns. The last phase, known as guardrails, automates corrections and institutionalizes them through policy. Human vigilance cannot maintain improvement gains as teams rotate, knowledge spreads, and priorities shift over time. Optimization efforts can use automated guardrails to stop old problems from coming back, making improvements last in the system. Validation, the fifth phase, measures improvements that are more than just cost. Performance metrics, throughput measurements, and reliability readings can show how the optimization step is succeeding, while the repeat step, or phase, ensures that optimization is not a project but a discipline, compounding value over time as things improve.

#### **2.1 Implementing the Flywheel: Organizational Readiness and Change Management**

Successfully implementing the flywheel methodology requires more than technical capability and analytical tools. Organizations must cultivate organizational readiness spanning cultural attitudes, stakeholder alignment, skills development, and governance structures that enable sustained optimization efforts. The transition from viewing cloud cost management as a reactive financial exercise to embracing it as a proactive engineering discipline necessitates fundamental shifts in how teams perceive their responsibilities, measure success, and allocate effort across competing priorities.

Organizational readiness assessment should evaluate multiple dimensions before launching comprehensive optimization initiatives. Leadership commitment represents perhaps the most critical prerequisite, as sustained optimization requires executive sponsorship willing to invest resources, empower teams with appropriate authority, and maintain focus through inevitable challenges and setbacks. Without visible leadership support, optimization initiatives risk being deprioritized when competing demands arise or when initial efforts encounter resistance from teams comfortable with existing practices. Leadership must articulate clear expectations that cost efficiency constitutes a core engineering responsibility rather than an optional enhancement pursued only when convenient.

Stakeholder alignment across organizational boundaries proves equally essential for optimization success. Cloud cost optimization inherently involves multiple stakeholders, including engineering teams responsible for infrastructure design and implementation, finance organizations tracking budgets and variances, procurement teams negotiating vendor agreements and pricing structures, and business units consuming cloud resources to deliver customer value. These diverse stakeholders often operate with different priorities, incentives, and success metrics, creating potential conflicts that can derail optimization efforts. Effective stakeholder alignment establishes a shared understanding of optimization objectives, agreed-upon metrics for measuring progress, and clear accountability structures defining roles and responsibilities.

Engineering teams must develop new competencies spanning cost analysis, resource optimization, and financial modeling that traditionally fell outside typical technical roles. Cloud architects and developers need visibility into the cost implications of their architectural decisions, understanding how choices regarding service selection, instance sizing, data storage strategies, and application design patterns directly impact infrastructure spending. This cost awareness should inform design reviews, architecture decisions, and implementation choices rather than being considered only after deployment when optimization opportunities become more constrained. Organizations can accelerate competency development through training programs, knowledge sharing sessions, cost optimization workshops, and mentoring relationships pairing experienced practitioners with team members new to optimization disciplines.

Change management strategies must address both technical process changes and cultural transformations required for sustainable optimization. Introducing new workflows for cost monitoring, approval gates for infrastructure provisioning, and governance policies constraining certain practices inevitably encounters resistance from teams accustomed to previous approaches prioritizing speed and flexibility above cost considerations. Effective change management acknowledges these concerns, provides context explaining why changes are necessary, demonstrates how new approaches enable rather than constrain effective engineering, and celebrates early successes, building momentum for broader adoption. Communication plans should regularly share optimization wins, recognize team contributions, and reinforce connections between cost efficiency and organizational success.

Governance frameworks institutionalize optimization practices through policies, standards, and accountability mechanisms that persist beyond individual initiatives or personnel changes. Well-designed governance balances necessary controls, preventing wasteful practices against flexibility, enabling innovation and rapid experimentation. Governance policies might include requirements for business justification when provisioning resources exceeding certain cost thresholds, automated approval workflows routing high-cost requests through appropriate review channels, mandatory tagging standards enabling cost attribution to specific teams or projects, and regular review cycles evaluating resource utilization and identifying optimization opportunities. These governance mechanisms should feel supportive rather than bureaucratic, enabling teams to move quickly while maintaining appropriate cost discipline.

Metrics and reporting systems provide visibility into optimization progress and ongoing cost efficiency. Organizations should establish baseline measurements capturing current cost efficiency levels before optimization initiatives begin, enabling quantification of improvement over time. Regular reporting cadences communicate optimization results to stakeholders, track progress toward cost reduction targets, identify emerging cost trends requiring attention, and maintain organizational focus on efficiency objectives. Dashboards and visualization tools help stakeholders at different organizational levels access relevant cost information, from executive summaries highlighting total spending and savings trends to detailed technical views enabling engineers to optimize specific resources or workloads.

Table 2: Flywheel Methodology Phases and Optimization Principles [3,4]

Phase	Core Function	Optimization Principle
Observation	Pattern recognition and anomaly detection	Distinguish waste signals from expected operational costs
Attribution	Causal analysis of spending patterns	Trace expenditures to underlying system behaviors
Correction	Behavioral modification at the root cause level	Address sources rather than symptoms of inefficiency
Guardrails	Automation and policy enforcement	Institutionalize improvements through systematic controls
Validation	Multi-dimensional impact assessment	Measure outcomes across cost, performance, and reliability
Repetition	Continuous improvement cycles	Compound value through iterative optimization

### **III. SECURITY AND COMPUTE LAYER OPTIMIZATIONS: ADDRESSING HIGH-IMPACT INEFFICIENCIES**

Often overlooked in cost analysis, the security layer can reveal surprising inefficiencies unrelated to its security properties. For example, in our thorough analysis of the AWS Key Management Service operations, we found that organizations tend to use poorly optimized encryption workflows in the Cloud, causing a collection of API calls and unnecessary cache misses on the application layer [5]. Most users of KMS are not aware that for every encryption operation, KMS creates a data key. Therefore, high-volume encryption operations can lead to meaningful API costs for KMS. A correct understanding of the distinction between customer master keys and data keys allows data keys to be cached and reused for multiple encryption operations with security guarantees, but many organizations do not try to take advantage of this. They treat each encryption operation as needing new key material, but the crypto architecture allows data keys to be reused within the security domain of the key management service or application. This analysis shows that envelope encryption combined with data key caching drastically reduces the required number of API calls to KMS without compromising security (the customer master key never leaves the boundary of the KMS service, and the data keys can be cached for the duration of some workflow or application session).

Current research into encryption strategies and key management practices indicates that organizations are attempting to balance security and operational considerations [6]. The study investigates the merits of client-side, server-side, and hybrid encryption strategies and concludes that the key difference between a low-cost, efficient encryption strategy and a secure encryption strategy is key management. The results demonstrate that misunderstandings about the cloud provider's key management security model contribute to at least some of these problems. An effective key-management system will utilize several principles, including hierarchical key management, key rotation strategies tailored to the threat model, and specialized caching schemes to strike a balance between performance and the security properties of the full-scale key management plan. In addition to the API costs and overhead, encryption incurs additional costs, which include the cost of cryptographic computation, using network bandwidth to fetch encrypted keys, and the additional overhead introduced by cryptographic data structures. This shows the need to build holistic mechanisms.

Another area in non-critical environments where that default behavior causes a lot of waste is compute resource management. Lower environments are often provisioned all the time, yet are only really used during business hours. Development and test environments often run during normal working hours. Non-production environments may continue to run after hours, on weekends, and on holidays, adding to costs. Scheduling automation can turn off and remove non-production idle resources, without changing development processes or testing patterns. Resources are created when the team arrives and are destroyed when they leave. During working hours, both development and testing are in full swing, but outside those hours, there is no cost. Organizations often find that savings from shifting non-production compute resources are substantial. Cluster lifecycle management in non-production clusters frequently assumes a persistent cluster model where cluster resources are at least initialized and running for the service's lifetime as a long-running, bursty, periodic workload. Deploying workloads on on-demand clusters, with fresh clusters created for each workload execution, establishes stricter operational standards for cluster usage and prevents wasted compute resources, especially when the cost of cluster initialization is not a significant concern for non-production workloads.

#### **3.1 Advanced Compute Optimization Strategies and Workload Characterization**

Beyond basic scheduling automation and lifecycle management, advanced compute optimization requires a deep understanding of workload characteristics, resource consumption patterns, and performance requirements that enable

sophisticated right-sizing and purchasing strategy decisions. Organizations often provision compute resources based on peak capacity requirements or worst-case scenarios, resulting in substantial overprovisioning during typical operating conditions. Effective optimization balances the need for adequate capacity during demand spikes against the cost of maintaining excess capacity during normal operations.

Workload characterization involves systematic analysis of resource utilization patterns across temporal dimensions, including daily cycles, weekly patterns, seasonal variations, and business-driven demand fluctuations. Different workload types exhibit distinct resource consumption profiles requiring tailored optimization approaches. Batch processing workloads that execute periodically on predictable schedules can leverage spot instances or burstable instance types that provide cost advantages compared to continuously running dedicated instances. Interactive applications serving user requests require consistent availability but may experience predictable traffic patterns, enabling scheduled capacity adjustments during known low-traffic periods. Analytics workloads processing large datasets benefit from memory-optimized or compute-optimized instance types matching their specific resource requirements rather than general-purpose instances, potentially mismatched to actual needs.

Instance family selection significantly impacts both cost and performance outcomes. Cloud providers offer diverse instance families optimized for different workload characteristics, including compute-intensive processing, memory-intensive data analysis, storage-intensive databases, graphics-intensive rendering, and general-purpose balanced workloads. Organizations frequently default to general-purpose instance families for convenience, despite specialized instance types potentially offering superior price-performance characteristics for specific workloads. Systematic evaluation of workload requirements against available instance families often reveals opportunities for substantial cost reduction or performance improvement through better matching of instance characteristics to actual workload needs. Purchasing strategy optimization encompasses decisions regarding on-demand instances, reserved instances, savings plans, and spot instances, each offering different trade-offs between cost, commitment, and availability guarantees. On-demand instances provide maximum flexibility but the highest per-hour costs, making them appropriate for unpredictable workloads or temporary infrastructure. Reserved instances and savings plans offer significant discounts in exchange for commitment to consistent usage levels over extended periods, making them cost-effective for steady-state workloads with predictable capacity requirements. Spot instances provide the deepest discounts by utilizing spare cloud capacity, but can be interrupted with minimal notice, making them suitable for fault-tolerant workloads that can handle interruptions gracefully.

Autoscaling configuration optimization ensures that infrastructure scales appropriately in response to demand changes without overprovisioning during scale-up events or scaling down too aggressively and impacting performance. Many organizations configure autoscaling policies too conservatively, scaling up quickly at the first sign of increased demand but scaling down slowly or not at all, resulting in excess capacity during normal operations. Effective autoscaling policies balance responsiveness to demand increases with appropriate scale-down behavior, removing unneeded capacity once demand subsides. Scaling metrics should reflect actual workload requirements rather than simplistic CPU utilization thresholds that may not accurately indicate capacity needs for memory-intensive or I/O-intensive workloads. Container orchestration platforms, including Kubernetes, introduce additional optimization opportunities through bin-packing efficiency, resource requests and limits configuration, horizontal pod autoscaling, and cluster autoscaling. Efficient bin-packing maximizes the number of containers scheduled on each node, reducing the total node count required to run workloads and consequently lowering compute costs. Properly configured resource requests and limits ensure that containers receive adequate resources without reserving excessive capacity that remains unused. Horizontal pod autoscaling adjusts the number of container replicas based on observed metrics, while cluster autoscaling adds or removes nodes based on overall cluster utilization, creating a multi-layered optimization approach matching capacity to actual demand at both the container and node levels.

Table 3: Security Layer Optimization Strategies and Key Management Practices [5, 6]

Optimization Domain	Inefficiency Pattern	Corrective Strategy
Encryption Workflows	Excessive API calls from poor caching implementation	Envelope encryption with data key reuse
Key Management Architecture	Unnecessary cryptographic operations from conservative approaches	Hierarchical key structures with appropriate rotation policies
Security Service Usage	Application-layer behaviors generating redundant operations	Caching mechanisms preserving security while reducing overhead
Cost Components	Direct API charges plus computational and storage overhead	Holistic optimization addressing multiple cost dimensions

#### **IV. PLATFORM MODERNIZATION AND PERFORMANCE ENGINEERING AS COST CONTROL MECHANISMS**

Infrastructure modernization, while often occurring to access new features and security controls, is also a cost optimization measure. A Lumenalta study on cloud modernization strategies found that organizations committed to platform modernization projects saw a range of cost, agility, and performance improvements [7]. Cloud modernization can refer to application refactoring, infrastructure optimization, process automation, or architectural redesign (from monolithic architecture to microservices architecture). If buyers introduce changes to their existing applications and cloud environments through modernization, they can realize savings from better usage consolidation, elimination of legacy licensing and maintenance fees, utilization of more attractive cloud-native services with better price-performance profiles, and automation to reduce labor overhead. Accordingly, the study proposes that successful modernization requires plans that treat the trade-off between short-term cost reduction opportunities and long-term architectural durability systematically, and that organizations tend to obtain better solutions by pursuing incremental modernization paths that deliver measurable immediate value rather than completely transforming the system in one step.

Their dual nature is clearest in database platform management, where vendor support lifecycles create business triggers for database version upgrades. Each year, as a database instance runs beyond vendor standard support windows, extended support charges are assessed until a database version upgrade is performed and extended proportionally to the size of the infrastructure. These features and the migration of the database instances in each environment to the current version eliminate the additional support costs, as well as allow the database stack to be used with performance and feature enhancements. In this example, it shows how infrastructure can become more expensive over time as the vendor's support policy changes and older versions roll over into different price brackets. The case study shows that upgrade cycles are necessary not only for security and technical reasons but also for cost control.

When done properly, performance engineering can provide the highest form of cost optimization by delivering a cost reduction, throughput improvement, or a combination of the two factors, depending on the priorities of the organization. Previous studies on optimizing resource consumption in big data processing have documented underutilized computational clusters due to configuration, workload, and resource allocation problems [8]. The work also studies dynamic provisioning of resources, scheduling workloads, and optimizing data locality as techniques that can help increase cluster utilization and decrease overall resource spending. It makes the point that effective monitoring is critical for any evidence-based optimization approach, as it allows for the detection of resource bottlenecks and anti-patterns in queries, as well as for the quantitative validation of configuration changes. Performance interventions can be done in an instance right-sizing manner if consumption information is available for resource reallocation, in the case of parallelism if workloads can be balanced among resources, or by configuring executors for workload characteristics. Observability data allows companies to increase their cluster utilization from low to very high levels, meaning that they are doing considerably more work per unit of infrastructure provisioned. Sometimes this does not reduce the cost of provisioned infrastructure, such as when throughput is the main bottleneck, where costs are unchanged while job completion time is massively reduced. These changes can lead to overall lower costs and are commonly a combination of many cost reduction factors. The eventual cost reduction can be important and result from a collection of vectors of improvement through a complete optimization program that includes better purchasing regimes, improved capacity planning, improvements in internal tuning, optimized processor architecture, and improved platform versions. These

reductions are not driven by a single intervention but by the application of modernization, tuning, and cloud-native approaches at the platform level.

#### **4.1 Data Platform Optimization and Query Performance Tuning**

Data platforms, including data warehouses, data lakes, and analytics engines, represent significant cost centers in modern cloud environments, with spending driven by factors including compute resources for query execution, storage volumes for data retention, and data transfer between services and regions. Organizations operating substantial data platforms often discover that optimization initiatives targeting query performance, data organization, and resource allocation can dramatically reduce costs while simultaneously improving query response times and analytical capabilities.

Query performance optimization begins with understanding how queries interact with underlying data structures and execution engines. Poorly written queries can consume excessive computational resources by scanning entire datasets when selective filtering would suffice, performing unnecessary joins or aggregations, or executing inefficient operations that overwhelm available memory and spill to disk. Query optimization techniques include rewriting queries to leverage partition pruning and predicate pushdown, adding appropriate indexes or materialized views to accelerate common access patterns, and restructuring data models to align with typical query patterns. Organizations should establish query performance monitoring by identifying expensive queries consuming disproportionate resources, enabling targeted optimization efforts focusing on high-impact improvements.

Data organization strategies significantly impact both query performance and storage costs. Partitioning large tables by frequently filtered columns, such as date ranges, enables queries to scan only relevant partitions rather than entire tables, dramatically reducing data volumes processed and consequently lowering compute costs. Clustering or sorting data within partitions by commonly filtered or joined columns further improves query performance by enabling more efficient data access patterns. Compression reduces storage costs and can improve query performance by reducing I/O requirements, though compression algorithms must be selected considering the trade-offs between compression ratio, decompression computational overhead, and compatibility with predicate pushdown optimization.

Data lifecycle management for analytical platforms extends beyond simple retention policies to include considerations regarding hot versus cold data access patterns. Recently ingested data typically experiences frequent query access, justifying storage in high-performance tiers optimized for low-latency access. As data ages, query frequency typically declines, enabling migration to lower-cost storage tiers accepting higher access latency in exchange for reduced storage costs. Some data platforms support automatic tiering based on access patterns, while others require explicit policies defining when data transitions between tiers. Historical data required for compliance or occasional analysis but rarely queried can be archived to the lowest-cost storage tiers, dramatically reducing costs while maintaining data availability for the infrequent cases when access is necessary.

Workload management and resource allocation ensure that diverse workloads sharing data platform resources receive appropriate priority and capacity allocation without interfering with each other. Interactive queries supporting business users require low latency and should be prioritized over batch processing jobs that can tolerate longer execution times. Resource allocation mechanisms, including workload management queues, query concurrency limits, and memory allocation policies, prevent individual queries or workloads from monopolizing platform resources and degrading performance for other users. Proper workload management configuration balances competing objectives of resource utilization efficiency, query performance predictability, and fair resource allocation across different user populations and use cases.

Materialized views and aggregate tables trade increased storage costs for reduced query execution costs by pre-computing and storing query results for frequently executed queries or common aggregation patterns. When queries repeatedly perform expensive operations such as complex joins, aggregations across large datasets, or calculations requiring substantial computational effort, materializing these results enables subsequent queries to retrieve pre-computed values rather than re-executing expensive operations. Organizations must balance the storage costs of maintaining materialized views against the compute cost savings from avoiding repeated expensive query execution, considering factors including query frequency, computational complexity, and data update patterns affecting materialization maintenance costs.

Table 4: Platform Modernization and Performance Engineering Benefits [7, 8]

Modernization Dimension	Traditional Challenge	Optimization Outcome
Infrastructure Evolution	Legacy systems with licensing overhead	Cloud-native services with superior price-performance ratios
Database Platform Management	Extended support charges for outdated versions	Version currency eliminates avoidable recurring costs
Cluster Resource Utilization	Suboptimal efficiency from configuration issues	Enhanced throughput through right-sizing and tuning
Workload Distribution	Imbalanced processing and bottlenecks	Dynamic provisioning and locality optimization

## **V. STORAGE AND NETWORK HYGIENE: ELIMINATING ACCUMULATED TECHNICAL DEBT**

Although storage costs are low relative to other public cloud services on a per-gigabyte basis, costs can add up quickly when organizations do not optimize their lifecycle management. For example, studies of organizations that use cloud storage management show that organizations store excess data by keeping data that is no longer needed, duplicating data, and failing to migrate data from a high-tier to lower tiers when it is infrequently accessed [9]. This technique includes the practice of data classification taxonomies, automated lifecycle management policies, tiering, and deduplication, all of which reduce storage costs and offer sufficient availability in the production environment. Understanding data access patterns, prioritizing and automating data movement between storage classes based on usage frequency, and applying retention and deletion controls can create a more efficient, cost-effective, and compliant storage environment, moving away from the misconception of the need for indefinite data retention. Systems employing a systematic approach to storage lifecycle management can save important costs, reduce backup windows, simplify disaster recovery, and more efficiently organize the data repositories they manage.

Checkpoint and recovery data is used during workflow execution, but after that can persist, becoming waste unless system-defined or user-requested cleanup is performed. A few case studies show what happens when checkpoint data is allowed to build up with no cleanup performed by either the system or a user (even if the majority of the data is no longer useful in the event of recovery). These costs show that a proper maintenance process to delete unused checkpoints and a smart retention policy can save a lot of money each year. People who think that storage is a good deal will regard these expenses as surprising, but the real perception here is that storage looks cheap until volumes increase big enough that material cost matters: lifecycle management is a first-class engineering challenge.

As explained by Snowflake's research on storage lifecycle policies, modern data platforms in the cloud require advanced data retention policies that balance cost management and governance policies [10]. Data lifecycle policies automatically manage data through its lifecycle, from the point of generation and usage to the point of archiving or deletion when it is no longer needed. Policies that automatically shift data between tiers based on usage patterns to improve performance and reduce costs (i.e., frequently accessed data on high-performance storage, infrequently accessed data on lower-cost storage, including archival storage) are also available to organizations. The research defined optimal lifecycle policies based on data classification, compliance, query performance, and cost optimization objectives. An effective automated lifecycle management system can help organizations save significant amounts of money by reducing storage costs while maintaining appropriate levels of data access for legitimate data access patterns. A variety of policies can be applied depending on the data type, retention period, and business need.

There are other areas within the network sphere where hygiene can reduce waste. One of these is the decommissioning of virtual private network (VPN) endpoints. VPN endpoints created for a certain connectivity need are often left up long after the associated applications have ceased, projects have been completed, and network topologies have changed. Network auditing can reveal underutilized endpoints, including those that incur charges but no longer serve an active connectivity function. Systematic reclamation of such resources results in faster gains while simplifying the architecture and operational workload. This is a classic example of infrastructure leftovers: resources that were set aside for a good reason but are no longer needed and keep costing money. In general, network consolidation efforts can yield meaningful reductions in the number of endpoints due to the removal of orphaned resources and rationalization of connectivity patterns, resulting in cost and architecture simplification.

## 5.1 Backup, Disaster Recovery, and Data Protection Optimization

Backup and disaster recovery infrastructure represents another significant cost center that organizations can optimize without compromising data protection and business continuity objectives. Traditional backup approaches often retain multiple full backups across extended retention periods, consuming substantial storage capacity despite much of the backed-up data remaining unchanged between backup cycles. Modern backup strategies leverage incremental and differential backup techniques that store only changed data blocks, dramatically reducing storage requirements while maintaining the ability to restore to any desired point in time within the retention window.

Backup retention policies should align with actual recovery requirements, regulatory compliance mandates, and business continuity objectives rather than defaulting to indefinite retention or arbitrary retention periods disconnected from legitimate needs. Many organizations discover that they maintain backup data far longer than necessary, retaining daily backups for months or years when business requirements only mandate retention for weeks. Graduated retention policies that maintain recent backups at higher frequency and granularity while transitioning to lower frequency and longer intervals for historical backups balance recovery flexibility against storage costs. For example, organizations might retain daily backups for the most recent month, weekly backups for the prior quarter, and monthly backups for the prior year, providing adequate recovery options while minimizing storage consumption compared to retaining daily backups indefinitely.

Backup storage tier selection significantly impacts costs, with different storage classes offering trade-offs between access speed and storage pricing. Recent backups supporting short recovery time objectives should reside in storage tiers optimized for rapid access, accepting higher storage costs in exchange for minimal recovery latency. Older backups accessed infrequently for compliance or historical analysis can migrate to lower-cost archival storage tiers, where substantially reduced storage pricing justifies accepting longer retrieval times when restoration becomes necessary. Some backup solutions support automatic tiering that migrates backup data between storage classes based on age or access patterns, optimizing costs without requiring manual intervention.

Disaster recovery infrastructure maintained for business continuity often incurs substantial ongoing costs through continuously running standby resources, data replication bandwidth consumption, and redundant infrastructure in secondary regions or availability zones. Organizations should evaluate disaster recovery requirements considering factors including acceptable recovery time objectives, recovery point objectives representing tolerable data loss, and the business impact of extended service unavailability. Not all systems require the most stringent disaster recovery capabilities, and tailoring disaster recovery strategies to actual requirements enables cost optimization while maintaining adequate protection for critical systems.

Cross-region data transfer represents a significant cost component for disaster recovery strategies involving active data replication between geographically distributed regions. Cloud providers typically charge for data egress from source regions and sometimes for ingress to destination regions, making continuous replication of large datasets expensive. Organizations can optimize these costs through selective replication, focusing on critical data requiring active disaster recovery capabilities while employing alternative approaches such as periodic backup-based recovery for less critical systems. Compression and deduplication reduce the data volume requiring transfer, lowering bandwidth costs while maintaining adequate protection levels.

Snapshot-based backup strategies for block storage volumes and database instances provide cost-effective alternatives to traditional backup approaches for certain use cases. Snapshots capture point-in-time storage state and leverage incremental storage techniques that only consume capacity for changed blocks, making them substantially more storage-efficient than full backup copies. However, snapshot retention costs can accumulate over time, particularly when numerous snapshots are retained or when data churn rates are high, generating substantial changed block volumes. Organizations should establish snapshot retention policies balancing recovery granularity against storage costs, deleting obsolete snapshots no longer required for recovery purposes.

## VI. MONITORING, OBSERVABILITY, AND COST ATTRIBUTION

Effective cloud cost optimization depends fundamentally on comprehensive monitoring, observability, and cost attribution capabilities that provide visibility into spending patterns, resource utilization, and the relationship between infrastructure consumption and business outcomes. Without adequate visibility, organizations struggle to identify optimization opportunities, validate the impact of optimization interventions, and maintain cost discipline as infrastructure evolves. Modern cloud environments generate vast quantities of monitoring data spanning resource metrics, application performance indicators, cost allocation tags, and usage patterns that organizations must collect, aggregate, analyze, and act upon.

Cost allocation and chargeback mechanisms establish accountability by attributing infrastructure spending to specific teams, projects, applications, or business units consuming resources. Effective cost allocation requires comprehensive tagging strategies that label resources with metadata identifying their organizational ownership, business purpose, environment classification, and cost center allocation. Automated tagging policies enforced at resource provisioning time ensure consistent tag application, while tag compliance monitoring identifies untagged or incorrectly tagged resources requiring remediation. With proper cost allocation in place, organizations can generate detailed spending reports showing which teams or projects consume what resources, enabling informed discussions about resource utilization efficiency and optimization priorities.

Anomaly detection capabilities identify unusual spending patterns that may indicate waste, misconfiguration, or security incidents requiring investigation. Machine learning algorithms can learn normal spending patterns and alert when actual spending deviates significantly from expected baselines, enabling rapid response to cost anomalies before they accumulate into substantial waste. Anomaly detection proves particularly valuable for identifying issues such as runaway autoscaling consuming excessive resources, misconfigured batch jobs executing more frequently than intended, or compromised credentials being exploited for cryptocurrency mining or other malicious activities that generate substantial unexpected costs.

Resource utilization monitoring provides the empirical foundation for right-sizing decisions by revealing actual resource consumption patterns compared to provisioned capacity. Organizations frequently discover that instances provisioned based on initial capacity estimates or inherited configurations substantially exceed actual requirements, creating right-sizing opportunities. Comprehensive utilization monitoring spanning CPU, memory, disk I/O, network bandwidth, and application-specific metrics enables confident right-sizing decisions supported by data rather than guesswork. Monitoring should capture both average utilization levels and peak utilization to ensure that right-sizing maintains adequate capacity for demand spikes while eliminating excess capacity during normal operations.

Custom metrics and application instrumentation extend monitoring beyond infrastructure-level resource consumption to capture business-relevant indicators that connect technical resource usage to business outcomes. Application-level metrics such as transaction volumes, user activity levels, data processing throughput, and business event rates provide context for understanding whether infrastructure spending scales appropriately with business activity. When infrastructure costs grow disproportionately faster than business metrics, this divergence signals inefficiency requiring investigation. Conversely, optimization initiatives that reduce costs while maintaining or improving business metric trends demonstrate genuine efficiency improvements rather than cost reduction achieved through unacceptable performance degradation.

Dashboard and visualization tools make monitoring data accessible to diverse stakeholders across the organization, from executives requiring high-level spending summaries to engineers needing detailed technical metrics for optimization efforts. Executive dashboards highlight total spending trends, major cost drivers, savings from optimization initiatives, and progress toward cost reduction targets, providing the visibility needed for informed strategic decisions. Technical dashboards enable engineers to drill down into specific resources, services, or workloads, analyzing detailed utilization patterns and identifying optimization opportunities. Self-service analytics capabilities empower teams to explore cost data, generate custom reports, and answer their own questions without depending on centralized teams, accelerating optimization efforts through distributed decision-making.

Forecasting and budgeting capabilities leverage historical spending patterns and planned infrastructure changes to project future costs, enabling proactive capacity planning and budget management. Accurate forecasting helps organizations avoid budget surprises, plan for anticipated spending increases associated with business growth or new initiatives, and identify when actual spending trajectories diverge from projections in ways requiring investigation. Budget alerts notify stakeholders when spending approaches or exceeds predefined thresholds, enabling timely intervention before cost overruns become severe. Forecasting models should incorporate both historical trends and planned changes, such as upcoming product launches, marketing campaigns driving user growth, or infrastructure migrations that may significantly impact spending patterns.

## VII. CONCLUSION: COST OPTIMIZATION AS AN ENGINEERING CULTURE AND CONTINUOUS PRACTICE

Cloud cost optimization is not just about operational or financial optimization. It is an engineering discipline about simultaneously optimizing cost, operational, and delivery performance. In cloud cost optimization, the flywheel methodology provides a repeatable framework for transforming episodic cost optimization activities into a compounding cycle of continuous iterative improvement. It includes security, compute, storage, and platform

infrastructures. The observed interventions, including encryption service optimization, non-production environment scheduling, cluster lifecycle management, cloud platform modernization, storage cleanup, and network optimization, provide evidence that a meaningful waste reduction can be achieved without compromises and is possible by eliminating inefficiencies that impact both financial and operational performance. Organizations that target high-leverage areas for remediation, including key management service call patterns, ephemeral compute lifecycle management, database version currency, performance engineering, and storage retention thresholds, can realize value quickly and gain confidence and traction to drive more common cost optimization. Such an outcome requires imparting cultural practices in engineering teams to challenge defaults, trace costs to root causes, remediate behavior rather than symptoms, and automate governance solutions to avoid progress being lost through personnel churn or changing focus. To do such tasks well, organizations don't need special tools or access, just a disciplined approach to the system's observability, error attribution, remediation, operations, validation, and iteration boundaries as the scale of the infrastructure and workloads that they are automating grows. Organizations that build this discipline can achieve meaningful economic efficiency gains by increasing throughput, which shows that, when done right, cost efficiency improves technical and business metrics.

## REFERENCES

- [1] Tanner Luxner, "Cloud computing trends: Flexera 2024 State of the Cloud Report," Flexera Blog, 2024. [Online]. Available: <https://www.flexera.com/blog/finops/cloud-computing-trends-flexera-2024-state-of-the-cloud-report/>
- [2] Kiruthika Ramesh, "Cloud Cost Optimization to Reduce Waste & Maximize Value," Atlas Systems Blog, 2025. [Online]. Available: <https://www.atlassystems.com/blog/cloud-cost-optimization>
- [3] Titilayo Deborah Olorunyomi et al., "Developing Sustainable Cloud Governance and Cost Management Strategies for Financial Operations," International Journal of Engineering Research and Development, 2024. [Online]. Available: <https://ijerd.com/paper/vol20-issue11/2011478485.pdf>
- [4] Long Wang, et al., "An Iterative Optimization Framework for Adaptive Workflow Management in Computational Clouds," IEEE, 2013. [Online]. Available: <https://ieeexplore.ieee.org/document/6680948>
- [5] Mihir Popat, "Mastering cloud security with AWS KMS: The ultimate guide to data encryption and key management," Medium, 2024. [Online]. Available: <https://mihirpopat.medium.com/mastering-cloud-security-with-aws-kms-the-ultimate-guide-to-data-encryption-and-key-management-b0a136b8c3b1>
- [6] Moses Blessing, et al., "Cloud Encryption Strategies and Key Management," ResearchGate, 2024. [Online]. Available: [https://www.researchgate.net/publication/383660212\\_Cloud\\_Encryption\\_Strategies\\_and\\_Key\\_Management](https://www.researchgate.net/publication/383660212_Cloud_Encryption_Strategies_and_Key_Management)
- [7] Lumenalta, "What is cloud modernization? (Updated 2025)," [Online]. Available: <https://lumenalta.com/insights/what-is-cloud-modernization>
- [8] Praveen Kumar, Vijay Singh Rathore, "Improvising and Optimizing Resource Utilization in Big Data Processing," ResearchGate, 2016. [Online]. Available: <https://www.researchgate.net/publication/314940260>
- [9] Yaser Mansouri, et al., "Data storage management in cloud environments: Taxonomy, survey, and future directions," ResearchGate, 2017. [Online]. Available: <https://www.researchgate.net/publication/321735551>
- [10] Kaushal Jain, "Optimize Storage Costs and Simplify Compliance with Storage Lifecycle Policies, Now Generally Available," Snowflake, 2025. [Online]. Available: <https://www.snowflake.com/en/blog/storage-lifecycle-policies-ga/>