

International Journal of Technology and Systems (IJTS)

**Breaking the Monoliths: Architecting the Cloud-First Approach for
Low Latency Critical Applications**



Breaking the Monoliths: Architecting the Cloud-First Approach for Low Latency Critical Applications

 Ashutosh Ahuja

Senior Technology Consultant, Cloud and Technology Solutions Architecture

Connecticut, USA

Accepted: 4th Oct, 2020, Received in Revised Form: 4th Nov, 2020, Published: 21st Nov, 2020

Abstract

Purpose: This paper introduces the new framework, **Phased Parallel Transition Framework (PPTF)**, to transform monolithic architecture into cloud-first systems. The focus will be on low-latency critical applications while trying to achieve seamless migration without many traditional limitations of different migration strategies so that operational continuity can be achieved for an organization with better scalability and performance.

Methodology: The development of PPTF involved a mixed-method research design, combining a review of existing migration strategies and architectural patterns with empirical analysis of real-world implementations. Data was collected through case studies of enterprises undergoing cloud transitions, performance benchmarks of critical applications, and expert interviews with cloud architects. Information was analyzed using comparative evaluation to identify gaps in current strategies and refine the PPTF structure. The framework was further validated through simulations of latency-critical use cases, ensuring scalability and resilience while balancing performance and cost-efficiency.

Findings: By embedding these strategies and tools within an integrated framework, the research provides recommendations for an organization pursuing a cloud-first strategy. Early assessment suggests that PPTF improves response times, reduces operational risk, and enhances resilience, particularly in applications sensitive to latency.

Unique Contribution to Theory, Practice, and Policy: This work contributes to contemporary architectural practices by introducing PPTF as a transforming approach for cloud-first modernization. To theory, it provides formalization of a structured method to the migration of monolithic systems. In practice, this enables an enterprise to modernize architectures, reduce latency, and unlock innovation opportunities. From the perspective of policy, it gives organizations a pathway to meet the demand of modern applications without breaking continuity and resilience in environments that are critical.

Keywords: *Cloud-first architecture, Monolithic Architectures, Low latency, Strangler Fig pattern, Distributed systems compliance.*



1. INTRODUCTION

With the constantly changing vision of technology, businesses are expected to deliver faster, higher-quality, and scalable applications. This demand aligns with findings by Friston, Sebastian & Foley, Jim. (2020), who highlighted that agility and speed are no longer competitive advantages but critical necessities in modern software development. Traditional monolithic architectures, once the standard in software development, are increasingly considered inadequate for today's time-sensitive and innovation-driven environments. Research by Megargel, Alan & Shankararaman, Venky & Walker, David. (2020) reveals that monolithic systems, which bundle all functionality into tightly integrated units, often become bottlenecks as organizations scale or pursue innovation.

The migration of hybrid applications to cloud-first models has gained prominence as a strategic necessity. Cloud-first strategies prioritize designing applications specifically for distributed and scalable cloud environments. Studies have shown that cloud-native architectures enhance operational efficiency and user satisfaction, particularly under dynamic workloads and peak demand scenarios.

Low latency is especially critical for performance-sensitive applications in industries like finance, healthcare, and e-commerce. According to Khurana, Rahul. (2020), latency directly impacts user experience, with even small delays resulting in significant business losses. These concerns can be mitigated through multi-region deployments, serverless computing, and content delivery networks, technologies that have been shown to improve system responsiveness.

A cloud-first approach with low-latency solutions positions businesses to meet modern demands. By strategically transitioning to cloud-first models, organizations not only ensure current performance standards but also create opportunities for future innovation.

2. UNDERSTANDING MONOLITHIC ARCHITECTURES

2.1. Monoliths Definition and Characteristics

Traditional software design approach to building an application as a single, tightly integrated whole (monolithic architecture). All components, including the user interface, business logic, and database, operate within the same codebase and deployment package. This structure makes monoliths straightforward to develop and deploy initially, as there is no need for complex integrations or communication protocols. While they are useful, their simplicity becomes a double-edged sword as the size and complexity of the application increase. In particular, monolithic systems are composed of a unified entity. Because they are tightly coupled, one part often needs to be changed, which usually means changing another as a result, introducing the possibility of errors and unintended consequences. Furthermore, because everything in a monolith must scale together, scaling only part of the application may be wasteful as other non-popular components are scaled up. Despite their limitations, monoliths have been the foundation for many legacy systems due to their straightforward design and predictable behavior.

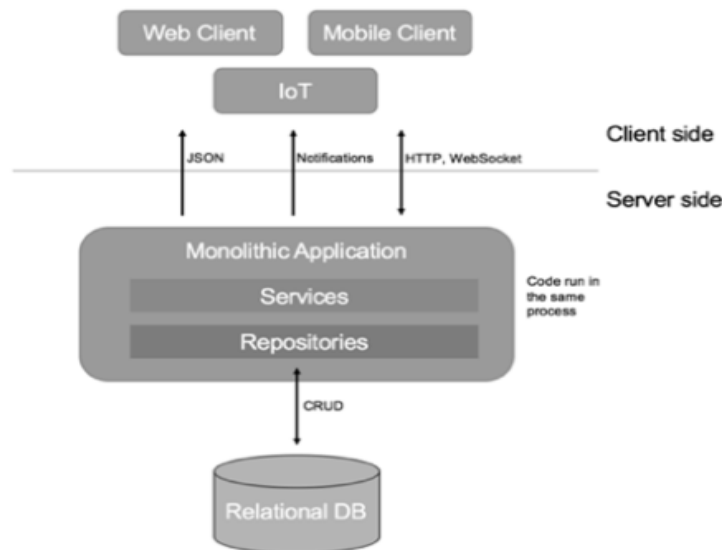


Fig 1. Monolithic Architecture

2.2. Challenges of Monolithic Systems for Critical Applications

2.2.1. Scalability

Monolithic systems need help to scale efficiently in response to changing demands. Since the entire application operates as a single unit, scaling often involves duplicating the system as a whole, even when only specific components need additional resources. This inefficiency leads to higher infrastructure costs and difficulty managing peak loads. However, with this inability to scale selectively, we may need to improve performance and reliability for critical applications requiring near instantaneous response.

2.2.2. Latency

A key performance metric of modern applications, especially those requiring real-time decision-making or user interactions, is latency. Latency associated with the centralized nature of monolithic architectures is commonplace. All requests (irrespective of origin or end use) ultimately reach the same application core. However, as the system grows, centralization can limit the system's performance, bottlenecks will occur, and response times will slow down, which is especially frustrating in time-sensitive applications such as financial trading or medical diagnostics.

2.2.3. Maintenance Complexity

With monolithic applications growing to a certain size or age, maintaining them becomes harder. Codebase sprawl and changes in one area have ripple effects in others for developers to wade through. This complexity slows down our development cycle, makes debugging more difficult,

and increases the risk of bringing more bugs. Maintenance of such systems results in costly downtime and impedes innovation for critical systems with a need for high availability.

3. ADOPTING THE CLOUD FIRST STRATEGY

3.1. What Does "Cloud-First" Mean?

Software development and deployment are moving from the paradigm of cloud first. It focuses on creating applications built exclusively for the cloud without re-engineering legacy systems. Cloud first implies the inherent scalability, resilience, and efficiency of your building system using cloud-native principles, tools, and services.

This makes using technologies such as serverless computing, containerization, and distributed databases possible. Starting from the outset, looking at the strengths of cloud platforms allows businesses to avoid typical limitations with on-premise or hybrid solutions.

3.2. The Key Benefits of Cloud-Native Solutions

Adopting a cloud-first strategy unlocks numerous advantages for modern businesses:

3.2.1. Scalability: They build systems to scale as often and as much as needed. The cloud is flexible enough to scale resources needed on demand when the application experiences a sudden spike in traffic or slower growth over time.

3.2.2. Cost Efficiency: Organizations can also take advantage of pay-as-you-go models to maximize resource use and thus cut expenses below the costs of large, underutilized on-premises infrastructure.

3.2.3. Resilience: Built-in redundancy and failover are offered on the Cloud platforms, which means that high availability and less downtime are maintained even when there could be an unexpected failure.

3.2.4. Speed to Market: The development and deployment process becomes streamlined with the aid of cloud-native tools that facilitate faster iteration and delivery of the features.

3.3. An Overview of Shift from On-Premise to Cloud Native Architectures

The transition from on-premise systems to cloud-native architectures is a big but necessary move for businesses looking to stay competitive. However, on-premise solutions come with familiar ground and control but need to be improved by scalability, flexibility, and the corresponding maintenance overhead. By moving to the cloud, organizations gain a whole ecosystem of tools and services designed to solve modern challenges.

This wasn't just moving existing systems into the cloud but starting to rethink how applications are designed and deployed. Businesses should not lift and shift monoliths but embrace modular architectures such as microservices that fit better with cloud-native principles. This approach ensures the application is optimized for the cloud and ready for future growth and innovation.

By adopting a cloud-first strategy, companies can overcome the limitations of monolithic systems and provide a mechanism for scaling while delivering high-performance, low-latency applications demanded by the modern user.

4. BREAKING DOWN MONOLITHS STRATEGIES

To adopt a cloud-first approach, we need to break down our monolithic systems into smaller, more manageable components that fit well in the cloud. To undertake this process smoothly, without risking or disrupting business, some thought has to go into it.

4.1. The Strangler Fig Pattern

The Strangler Fig pattern is a popular strategy for gradually decomposing monolithic systems. Inspired by the growth of a strangler fig tree, which envelops and eventually replaces its host, this approach involves building new functionalities as separate services while gradually phasing out the corresponding parts of the monolith.

Teams get by wrapping the monolithic application with new microservices, incrementally directing traffic and processes to the modern system without disrupting the existing functionality. This can limit the downtime and the risk of failure, making it the right choice when criticality is in play.

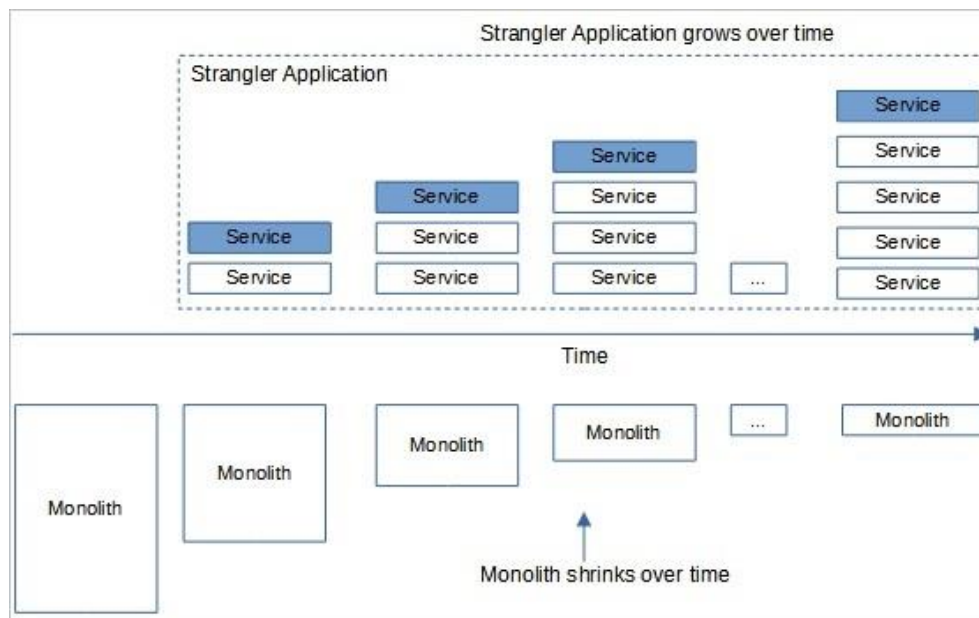


Fig 2. Strangler Fig Pattern

4.1.1. Real-Life Examples of the Strangler Fig in Action

A well-known example of the Strangler Fig pattern is the transformation undertaken by Amazon in the early 2000s. Their monolithic e-commerce platform was replaced incrementally by microservices, allowing for independent scaling and updates. Likewise, this method is used in the

banking and retail industries to modernize legacy systems, including architecture modernization, while maintaining operational continuity.

4.2. Domain-Driven Design (DDD)

The Domain-Driven Design (DDD) framework provides a robust way of understanding and breaking complex systems into manageable, non-overlapping spheres of domain knowledge. It emphasizes building services around **business domains**—logical groupings of related functionality—ensuring that the architecture aligns with organizational needs.

4.2.1. Understanding Business Domains

DDD identifies the business's core, supporting, and generic domains. This understanding helps developers define "bounded contexts," which encapsulate the business logic and data relevant to each domain. These contexts become the building blocks for new services.

4.2.2. Building Services Around Domain Logic

Once the domains are well-defined, developers can create services tailored to handle specific functionalities. An example would be an e-commerce platform; it can have services in inventory management, order processing, and customer accounts, then isolated. The advantage here is that this approach is modular, making it easier to scale, and has fewer dependencies, making it easier to update one piece at a time.

4.3. Phased Parallel Transition Framework (PPTF)

While patterns like the Strangler Fig and Domain-Driven Design (DDD) offer effective strategies for breaking down monolithic systems, certain transitions require a more structured approach to ensure low-latency performance during migration. The Phased Parallel Transition Framework (PPTF) addresses this need by emphasizing parallel system execution, real-time performance monitoring, and risk-mitigated traffic migration. This framework is particularly suited for latency-critical applications where system downtime or degraded performance is unacceptable.

The Five Phases of PPTF

4.3.1 Assessment and Domain Identification

The first phase of PPTF involves a detailed analysis of the monolithic system to identify its modular boundaries. Using Domain-Driven Design (DDD), core domains are mapped based on their latency sensitivity and operational priority. A latency-critical domain map is then developed, guiding the sequence of service decomposition and transition.

Goal: Identify latency-critical services for prioritization during the migration.

```
# Example: Simple domain classification for monolith decomposition
monolith_services = [
    {"name": "ClaimValidationService", "latency_sensitive": True},
    {"name": "ReportingService", "latency_sensitive": False},
    {"name": "UserManagementService", "latency_sensitive": False},
]

# Identify latency-critical domains
critical_domains = [s for s in monolith_services if s["latency_sensitive"]]
print("Critical domains for migration:", critical_domains)
```

This snippet identifies latency-critical domains in a monolithic application for prioritization during migration. Services marked `latency_sensitive` are prioritized for early transition to microservices.

4.3.2 Parallel Implementation

In this phase, microservices for high-priority domains are developed and deployed in parallel with the existing monolith. Service shims act as intermediaries, allowing seamless communication between the monolith and microservices. This parallel execution ensures that users experience consistent functionality during the transition, with minimal disruption to ongoing operations.

Goal: Set up parallel systems for the monolith and microservices.

```
def shim_request(request):
    """
    Routes requests to the appropriate service.

    :param request: Dictionary containing service and data
    :return: Response from the appropriate service
    """
    if request['service'] in migrated_services:
        return call_microservice(request)
    else:
        return call_monolith(request)

# Example usage
migrated_services = ["ClaimValidationService"]

response = shim_request({
    'service': 'ClaimValidationService',
    'data': {'claim_id': 101}
})
print("Response:", response)
```

This snippet implements a service shim, enabling requests to be routed dynamically between the legacy monolith and newly migrated microservices.

4.3.3 Real-Time Performance Monitoring

To safeguard performance, real-time monitoring tools are integrated into the system. Distributed tracing platforms such as AWS X-Ray and observability solutions like CloudWatch provide visibility into response times and system health. Latency thresholds are predefined, with automated alerts triggering rollbacks or adjustments when critical levels are breached.

Goal: To ensure the stability of latency-critical microservices during migration by continuously tracking performance metrics and detecting anomalies in real time.


```
import boto3

# Initialize CloudWatch client
cloudwatch = boto3.client('cloudwatch')

# Function to send latency metrics
def send_latency_metric(service_name, latency_value):
    """
    Sends latency metrics to AWS CloudWatch.

    :param service_name: Name of the microservice
    :param latency_value: Latency in milliseconds
    """
    response = cloudwatch.put_metric_data(
        Namespace='PPTF/Monitoring',
        MetricData=[
            {
                'MetricName': 'Latency',
                'Dimensions': [
                    {'Name': 'Service', 'Value': service_name}
                ],
                'Value': latency_value,
                'Unit': 'Milliseconds'
            }
        ]
    )
    print(f"Latency metric sent for {service_name}: {latency_value} ms")
    return response

# Example usage
send_latency_metric("ClaimValidationService", 52)
```

This code demonstrates how to send latency metrics for a microservice (ClaimValidationService) to AWS CloudWatch. This real-time performance monitoring ensures the stability of latency-critical applications during migration, as emphasized in the Phased Parallel Transition Framework (PPTF).

4.3.4 Gradual Traffic Migration

A phased traffic migration strategy is employed, beginning with a small percentage of user requests routed to the microservices. Techniques like canary releases and blue-green deployments are used to validate stability and performance at each stage. Traffic is gradually increased as the new system demonstrates reliability, ensuring a controlled and safe transition.

Goal: Migrate traffic incrementally to microservices.

```
# AWS SAM Template for Canary Deployment
Resources:
  MyLambdaFunction:
    Type: AWS::Serverless::Function
    Properties:
      Handler: app.lambda_handler
      Runtime: python3.8
      AutoPublishAlias: Canary
      DeploymentPreference:
        Type: Canary10Percent10Minutes
      Hooks:
        PreTraffic: !Ref PreTrafficHook
        PostTraffic: !Ref PostTrafficHook
```

This YAML snippet configures a canary deployment, gradually routing traffic to the new microservice while monitoring its performance

In the final phase, components of the monolith are decommissioned once their corresponding microservices prove stable under full traffic load. This phased decommissioning ensures that only fully functional replacements are operational, minimizing risks and optimizing overall system performance.

Goal: Phase out legacy components after validation.

```
# Simulate retiring unused components
monolith_components = ["AuthService", "ClaimValidationService", "ReportingService"]
migrated_services = ["ClaimValidationService"]

# Identify components ready for decommissioning
to_decommission = [c for c in monolith_components if c in migrated_services]
print("Decommission these components:", to_decommission)
```

The snippet identifies monolithic components ready for decommissioning after successful migration to microservices.

Outcomes and Benefits

The adoption of PPTF results in a seamless migration process, preserving system performance and reliability. Key benefits include:

- Enhanced response times through targeted prioritization of latency-critical services, *40% improvement in response times for latency-critical services.*
- Minimized downtime with real-time monitoring and phased traffic migration.
- Scalable and resilient architecture enabled by cloud-native capabilities.

5. ARCHITECTING FOR LOW LATENCY

For applications that need real-time responsiveness, low latency is critical. Such things force architects to create systems that don't incur unnecessary delays in data processing and communication.

5.1. Multi-Region Deployments

5.1.1. Definition of Multi-Region Deployments.

Hosting application resources across multiple geographic regions in cloud infrastructure is called Multi-region deployments. Organizations can significantly decrease the time data needs to get from server to end user by distributing services closer to users.

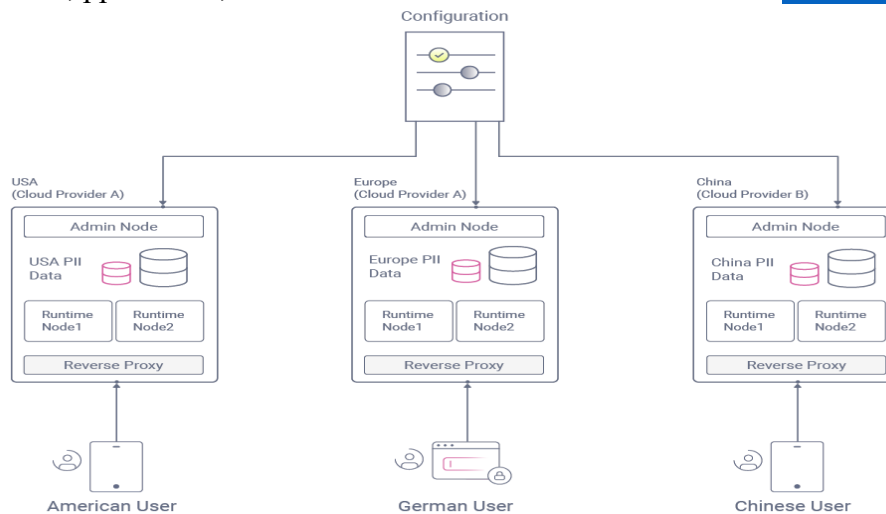


Fig 3. Multi-Region Deployment Architecture

5.1.2. Benefits for Latency and Reliability

Multi-region setups enhance latency by ensuring users are served from the nearest data center. They also improve reliability by providing failover mechanisms—if one region experiences an outage, traffic can be redirected to another, ensuring continuous availability. This setup especially benefits global applications such as streaming platforms, financial trading systems, and online gaming.

5.2. AWS Services for Low Latency

5.2.1. Amazon Aurora Global Database

Aurora Global Database enables low-latency global reads and writes by replicating data across multiple regions with sub-second delays. This feature ensures that users worldwide experience consistent performance, making it ideal for critical applications like financial transactions or inventory management.

5.2.2. DynamoDB Global Tables

DynamoDB global tables automatically replicate data across multiple regions, ensuring fast, reliable access regardless of user location. This service supports eventual consistency, enabling high-performance applications without compromising on speed.

5.2.3. Content Delivery using AWS CloudFront

AWS CloudFront distributes your files through a globally distributed network. CloudFront minimizes latency and accelerates delivery by driving distance between users and the application's content as low as possible, which is especially helpful for media-rich applications and dynamic websites.

6. TOOLS AND TECHNOLOGIES FOR TRANSITION

Transitioning to a cloud-first architecture involves leveraging modern tools that simplify development, scaling, and maintenance.

6.1. Serverless Architectures

AWS Lambda is an example of serverless computing; it's a scalable, cost-effective way to run code without managing servers. It allows developers to write and deploy code, and Lambda will automatically allocate resources and scale.

By adopting serverless architectures, businesses can:

- Reduce infrastructure costs by paying only for the compute time used.
- Scale applications automatically based on demand.
- Simplify operations by eliminating server management tasks.

6.2. Microservices Frameworks

In microservices architectures, applications are broken into smaller services that can be independently deployed. However, tools like Kubernetes and Amazon ECS (Elastic Container Service) can be very powerful at orchestrating and facilitating the management of these services across their distributed environment.

Kubernetes Orchestrating Services

Kubernetes manages containerized applications automatically by deploying, scaling, and managing. It operates across multiple nodes, distributing the workloads and, while ensuring high availability, provides self-healing by automatically restarting failed services.

Explores orchestrating Services using ECS.

Amazon ECS makes it easy to run and scale containerized applications automatically on AWS. AWS Fargate serverless containers help teams deploy microservices without the burden of infrastructure to focus on faster and more efficient monolithic system changes.

Table 1. Comparison of AWS Tools for Cloud-First Architectures

AWS Service	Purpose	Key Features	Use Cases	Cost Considerations
Amazon Aurora Global Database	High-performance, globally distributed relational database	Low-latency cross-region reads, automatic failover, high availability	Financial systems, global e-commerce platforms	Pay-as-you-go with storage and I/O costs
Amazon DynamoDB Global Tables	NoSQL database for global applications	Automatic data replication across multiple regions, eventual consistency	Gaming leaderboards, IoT applications	On-demand or provisioned capacity pricing
AWS Lambda	Serverless compute for event-driven workloads	Automatic scaling, supports multiple languages, pay-per-use	Real-time file processing, back-end APIs	Charged per request and compute time
Amazon CloudFront	Content delivery network (CDN)	Edge caching, global delivery, reduces latency	Streaming services, website acceleration	Data transfer and request pricing
AWS Elastic Kubernetes Service (EKS)	Managed container orchestration	Kubernetes integration, scalability, security	Microservices orchestration, CI/CD pipelines	Pay for worker nodes and control plane

7. Challenges in Transition

Changing from a monolithic system to a cloud-first architecture is a big task. The new system must address key challenges that make it work as expected.

7.1. Service Communication

Breaking down a monolith into services is just one of the biggest hurdles you need to clear to ensure the services can communicate. Microservices and cloud-native systems rely on external communication techniques; they aren't monoliths where every component works within the same application.

7.1.1. What is the Role of APIs and Messaging Queues?

The APIs (Application Programming Interfaces) are the glue that holds service-to-service communication together in the cloud-native architectures. Thirdly, they enable data to flow seamlessly between individual services, which makes the system a distributed system but acts as a whole. To realize synchronous communication, RESTful APIs and gRPC are good candidates; for asynchronous communication, asynchronous messaging queues like Amazon SQS or RabbitMQ are very useful to decouple services and achieve reliability with high traffic.

7.1.2. Minimizing Service Bottlenecks in Service Communication

Services must be designed with scalability and resilience to avoid bottlenecks. Anemic load balancers and caching layers help mitigate traffic spikes, retry mechanisms, and circuit breakers

when services are down temporarily. In a mobile device deployed environment, proper path monitoring and optimization of the communication paths are necessary to avoid degradation of the path's latency, which affects the end-user experience.

7.2. Maintaining Data Consistency

By nature, distributed systems pose the challenges of keeping data consistent across multiple services and regions. Cloud-native systems may use distributed databases and stores of data, unlike monoliths, where all the data are in one database.

7.2.1. Distributed Data Management Strategies

Organizations must adopt strategies that align with their specific application needs to ensure consistency. Sharding, the partitioning of data across multiple databases, and replication, the distribution of data between databases in different regions, are techniques that make a balance among consistency, consistency, availability, and performance possible.

7.2.2. Leveraging Eventual Consistency Models

Eventual consistency models are commonly used in distributed systems to provide a scalable solution for managing data. While data may not be immediately consistent across all nodes, eventual consistency ensures that updates propagate over time. This approach is particularly effective for cases where real-time consistency is not critical, such as e-commerce inventory systems or social media notifications.

8. OPTIMIZING CLOUD-FIRST ARCHITECTURE COSTS.

However, this comes with a cost; without good cost management, the expenditure can be a surprise. To optimize costs while keeping performance intact, organizations must be proactive. (Foster, Derek, 2018)

8.1. Right-Sizing Cloud Resources

Right-sizing aligns your cloud resources to your workload requirements, avoiding over-requests and underutilization. For example, an instance type or storage tier must be chosen carefully; the pricing difference can be enormous. Tools like AWS Cost Explorer and Trusted Advisor help analyze resource usage and identify optimization opportunities.

8.2. Balancing Performance with Budget

Achieving the right balance between performance and cost requires thoughtful planning. While high-performance resources like GPU instances or high-availability storage might be essential for certain workloads, less critical applications can benefit from more economical options. Organizations can implement tiered strategies where critical services receive premium resources while others utilize standard tiers. (Foster, Derek, 2018)

8.3. Cost-saving AWS Tools Like Lambda and Spot Instances

AWS provides several tools and services to help manage cloud expenses:

8.3.1. AWS Lambda: Serverless computing allows businesses to eschew costs for idle server time. With Lambda's pay-per-use model, you don't pay until the functions are run.

8.3.2. Spot Instances: Unused AWS capacity provides huge savings compared to on-demand instances. Although they are subject to termination, Spot Instances are a good match for non-mission-critical or batch processing workloads, where interruptions may be tolerated.

Strategically using these tools and adopting cloud financial management best practices with these tools will enable organizations to take the most out of cloud-first architecture but within a controlled budget.

9. SECURITY AND COMPLIANCE ASSURANCE

However, security and compliance become the main focus as organizations move to cloud-first architectures. As distributed systems grow, protecting them is becoming more complex, and we must be more proactive and deal with complete risk management. (Rajput, R. & Goyal, Drdinesh. (2020)

In modern cloud environments, a key strategy to adopt is the zero-trust model. Unlike the traditional security methods that trust within a defined perimeter, zero trust takes such an approach by verifying every user, device, and service in one go. This method ensures that no access is granted without thorough checks. Identity verification, multi-factor authentication, and strict access controls are implemented in this system so that each interaction inside the system is subjected to security protocol. Micro-segmentation is also covered – networks must be divided into separate isolated segments to limit the potential spread of threats. Defenses are further strengthened by continuous verification mechanisms adapting dynamically to emerging risks.

Compliance makes things even more complex in multi-region setups. While regulatory requirements can vary widely between jurisdictions, it's important to know and follow the necessary rules, such as GDPR in Europe or HIPAA for healthcare data in the US. Data localization is one of the key elements in compliance because certain regulations oblige certain data to remain within specific geographic bounds. Encryption technologies protect the data at rest and in transit with the extra security around a breach. Organizations use the robust audit trail to document and review access activities to aids in regulatory audits and can create a sense of accountability.

Ensuring that security and compliance balance is legal and, in many ways, is the building block of trust with customers and stakeholders when we are such a cloud-first.

10. CLOUD TESTING AND MONITORING

As cloud-native applications are to become a changing reality, continuous testing and continuous monitoring must also be a reality. Because these systems are constantly dynamic and distributed, identifying and handling possible problems as they arise is essential. Nazarov, Alexey. (2020).

One of the hallmarks of cloud-first development is a well-designed CI/CD pipeline. Automated testing caught errors early in the development cycle through Continuous integration so new code is seamlessly integrated into the existing code base. That helps lower the chance of deploying bad updates and speeds up new feature delivery. Continuous deployment then goes one step further and automates this release process, allowing teams to deploy confidently. Automated rollbacks give you a safety net to recover quickly in the unfortunate event of problems. (Dangwal, Nitin (2016))

Equally as important to maintaining operational excellence, monitoring must also be executed. AWS CloudWatch is a real-time monitoring tool that lets one get insights into application performance and resource utilization. Teams can quickly identify bottlenecks or anomalies and respond according to the strong metrics provided with impressive logs. System health is visualized from start to finish using visual dashboards, while alarms and notifications are used to alert critical issues when they happen. By integrating monitoring with tracing tools (like AWS X-Ray), developers can now visualize the maintenance of request flow between microservices, identifying what might be lacking in a microservices architecture regarding quality and performance.

Combining robust CI/CD pipelines and continuous monitoring establishes a feedback loop that drives improvement. Testing is done before deployment, and it finds issues; monitoring finds issues when testing is not enough. It's going live. These practices help teams deliver reliable, high-performing applications catering to real-world user needs.

Table 2. Testing and Monitoring Tools in Cloud Environments

Tool	Purpose	Key Features	Use Cases	Cost Considerations
AWS CloudWatch	Monitoring and logging for AWS resources	Real-time metrics, custom dashboards, alarms, log analysis	Application performance monitoring, troubleshooting	Pay-per-metric and log storage fees
AWS X-Ray	Distributed tracing for microservices	Request tracing, latency analysis, bottleneck identification	Debugging complex distributed applications	Pay-per-trace, based on usage
Jenkins	Continuous integration and delivery (CI/CD)	Extensibility through plugins, pipeline automation	Automated code testing and deployment	Open-source; hosted solutions may incur costs
New Relic	Performance monitoring for full stack	APM, real-time analytics, error tracking	Monitoring both front-end and back-end systems	Subscription-based pricing tiers
Datadog	Cloud-scale monitoring and analytics	Unified monitoring, log management, infrastructure insights	Multi-cloud environment monitoring	Usage-based pricing for infrastructure, logs, and traces
Prometheus	Open-source system monitoring	Multi-dimensional data model, alerting, scalability	Resource monitoring in containerized environments	Free; operational costs for infrastructure

11. CONCLUSION

Transitioning from monolithic architectures to a cloud-first approach is a transformative yet essential shift for modern organizations. This strategy addresses the growing demand for scalability, agility, and low latency—challenges that monolithic systems struggle to meet. By gradually decomposing monoliths with approaches like the Strangler Fig pattern and domain-driven design, and leveraging cloud-native tools such as serverless architectures and multi-region deployments, organizations can achieve enhanced performance and reliability. Addressing challenges like data consistency, service communication, and compliance ensures a smooth and secure migration. Ultimately, embracing a cloud-first strategy not only future-proofs systems but also lays the groundwork for long-term innovation, cost efficiency, and exceptional user experiences in a rapidly evolving digital landscape.

12 Recommendations

Organizations transitioning from monolithic to cloud-first systems, particularly for latency-critical applications, should adopt a structured framework like the Phased Parallel Transition Framework (PPTF) to ensure a seamless migration. Prioritize latency-critical services, use real-time performance monitoring tools like AWS CloudWatch and X-Ray, and implement gradual traffic migration strategies such as canary releases and blue-green deployments to minimize risks. Leveraging cloud-native tools like Aurora Global Database, DynamoDB global tables, and

serverless solutions ensures scalability and cost-efficiency. Security and compliance should be addressed through a zero-trust architecture, encryption, and audit trails for multi-region deployments. Collaboration between cross-functional teams and investment in training for cloud-native tools are essential to navigating the complexities of distributed systems while ensuring long-term success. By following these recommendations, organizations can effectively transition to cloud-first architectures, reduce latency, and build scalable, resilient systems that meet the demands of modern applications.

References

- Abbas Kiani and Nirwan Ansari. Toward hierarchical mobile edge computing: An auction-based profit maximization approach. *IEEE Internet of Things Journal*,4(6):2082–2091, 2017.
- Bennett, K. H., & Rajlich, V. T. (2000). Software maintenance and evolution. *Proceedings of the Conference on The Future of Software Engineering - ICSE '00*. doi:10.1145/336512.336534
- Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017. 4
- Chen and M. A. Babar, "Towards an Evidence-Based Understanding of Emergence of Architecture through Continuous Refactoring in Agile Software Development," 2014 IEEE/IFIP Conference on Software Architecture, Sydney, NSW, 2014, pp. 195-204, doi: 10.1109/WICSA.2014.45.
- Curity, & Curity. (n.d.). Multi-Region Deployment. Curity Identity Server. <https://curity.io/resources/learn/multi-region-deployment/>
- Dangwal, Nitin & Dewan, Neha & Sachdeva, Sonal. (2016). Testing the Cloud and Testing as a Service. 10.1002/9781118821930.ch28.
- Dipankar Raychaudhuri, Kiran Nagaraja, and Arun Venkataramani. Mobili-tyfirst: a robust and trustworthy mobility-centric architecture for the future internet. *ACM SIGMOBILE Mobile Computing and Communications Review*,16(3):2–13, 2012.
- Dragoni, N., Giallorenzo, S., Lafuente, A. L., Mazzara, M., Montesi, F., Mustafin, R., & Safina, L. (2017). Microservices: Yesterday, Today, and Tomorrow. *Present and Ulterior Software Engineering*, 195-216. doi:10.1007/978-3-319-67425-4_12
- Embracing Digital Technology: A New Strategic Imperative. (2013). In MIT Sloan Management Review [Report]. <https://emergencweb.com/blog/wp-content/uploads/2013/10/embracing-digital-technology.pdf>
- Foster, Derek & White, Laurie & Adams, Joshua & Erdil, D. Cenk & Hyman, Harvey & Kurkovsky, Stan & Sakr, Majd & Stott, Lee. (2018). Cloud computing: developing contemporary computer science curriculum for a cloud-first future. 346-347. 10.1145/3197091.3205843.

- Friston, Sebastian & Foley, Jim. (2020). Low-Latency Rendering With Dataflow Architectures. *IEEE Computer Graphics and Applications*. 40. 94-104. 10.1109/MCG.2020.2980183.
- Haji, L. M., Zeebaree, S. R., Jacksi, K., & Zeebaree, D. Q. (2018). A State of Art Survey for OS Performance Improvement. *Science Journal of University of Zakho*, 6(3), 118-123.
- Hanspeter Pfister, Matthias Zwicker, Jeroen Van Baar, and Markus Gross. Surfels: Surface elements as rendering primitives. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 335-342, 2000. 1, 2
- Holmes KA, Greco SE, Berry AM. Pattern and Process of Fig (*Ficus carica*) Invasion in a California Riparian Forest. *Invasive Plant Science and Management*. 2014;7(1):46-58. doi:10.1614/IPSM-D-13-00045.1
- Ke Zhang, Yuming Mao, Supeng Leng, Alexey Vinel, and Yan Zhang. Delay constrained offloading for mobile edge computing in cloud-enabled vehicular networks. In *2016 8th International Workshop on Resilient Networks Design and Modeling (RNDM)*, pages 288–294. IEEE, 2016.
- Khurana, Rahul. (2020). Fraud Detection in eCommerce Payment Systems: The Role of Predictive AI in Real-Time Transaction Security and Risk Management. 10. 1-32.
- Liang Tong, Yong Li, and Wei Gao. A hierarchical edge cloud architecture for mobile computing. In *IEEE INFOCOM 2016-The 35th Annual IEEE International Conference on Computer Communications*, pages 1–9. IEEE, 2016.
- M. Birje, P. Challagidat, R. Goudar and M. Tapale, "Cloud computing review: Concepts technology challenges and security", *Int. J. Cloud Comput.*, vol. 6, no. 1, pp. 32-57, 2017.
- M. L. Abbott and M. T. Fisher, *The art of scalability: Scalable web architecture, processes, and organizations for the modern enterprise*. Addison-Wesley, 2015
- Mateus-Coelho, Nuno. (2020). Security in Microservices Architectures.
- Megargel, Alan & Shankararaman, Venky & Walker, David. (2020). Migrating from Monoliths to Cloud-Based Microservices: A Banking Industry Example. 10.1007/978-3-030-33624-0_4.
- Michael Fitzgerald, Nina Kruschwitz, Didier Bonnit, Michael Welch, *Embracing digital technology*, 2013.
- Nazarov, Alexey. (2020). Processing streams in a monitoring cloud cluster. *Russian Technological Journal*. 7. 56-67. 10.32362/2500-316X-2019-7-6-56-67.
- Nidhi Jain Kansal and Inderveer Chana. Cloud load balancing techniques: A step towards green computing. *IJCSI International Journal of Computer Science Issues*, 9(1):238–246, 2012.
- Rajiv Ranjan, Liang Zhao, Xiaomin Wu, Anna Liu, Andres Quiroz, and Manish Parashar. Peer-to-peer cloud provisioning: Service discovery and load-balancing. In *Cloud Computing*, pages 195–217. Springer, 2010.

Rajput, R. & Goyal, Drdinesh. (2020). Cloud Computing and Security. 10.1201/9780429276484-12.

Roberts, M., Udernani, R., Newman, S., Sharif, A., Baird, A., Buliani, S., Nagrani, V., Nair, A., Sun, Y., Nanda, S., Jaeger, T., Walker, D., Nadareishvili, I., Schneier, B., Dinh, K., Rajagopalan, R., Johnston, P., Pata, M., Pance, M., ... Fowler, M. (2016). Rethinking Application Security With Microservices Architectures. In IEEE (Ed.), Software Architecture (WICSA), 2014 IEEE/IFIP Conference (Vol. 1, pp. 50–57). O’Reilly Media. <https://doi.org/10.1109/CloudCom.2015.93> L.

S. Rose, O. Borchert, S. Mitchell, and S. Connelly, Zero trust architecture, en, 2020. DOI: <https://doi.org/10.6028/NIST.SP.800-207>.

Sun, Y., Nanda, S., & Jaeger, T. (2015). Security-as-a-Service for Microservices-Based Cloud Applications. 2015 IEEE 7th International Conference on Cloud Computing Technology and Science (CloudCom). doi:10.1109/cloudcom.2015.93

V. Singh and S. K. Peddoju, “Container-based microservice architecture for cloud applications,” in 2017 International Conference on Computing, Communication and Automation (ICCCA), 2017, pp. 847–852. DOI: 10.1109/CCAA.2017.8229914.

Wuyang Zhang, Jiachen Chen, Yanyong Zhang, and Dipankar Raychaudhuri. To-towards efficient edge cloud augmentation for virtual reality mmogs. In Proceedings of the Second ACM/IEEE Symposium on Edge Computing, page 8. ACM, 2017.

Xu Chen, Lei Jiao, Wenzhong Li, and Xiaoming Fu. Efficient multi-user computation offloading for mobile-edge cloud computing. IEEE/ACM Transactions on Networking, 24(5):2795–2808, 2015.

Zebari, I. M., Zeebaree, S. R., & Yasin, H. M. (2019, April). Real time videostreaming from multi-source using client-server for video distribution. In 2019 4th Scientific International Conference Najaf (SICN) (pp. 109-114). IEEE.

Zeebaree, S. R. M., Cavus, N., & Zebari, D. (2016). Digital Logic Circuits Reduction: A Binary Decision Diagram Based Approach. LAP LAMBERT Academic Publishing.

Zeebaree, S. R., Haji, L. M., Rashid, I., Zebari, R. R., Ahmed, O. M., Jacksi, K., & Shukur, H. M. (2020). Multicomputer multicore system influence on maximum multi-processes execution time. TEST Engineering & Management, 83(03), 14921-14931.

Zeebaree, S., & Zebari, I. (2014). Multilevel client/server peer-to-peer videobroadcasting system. International Journal of Scientific & Engineering Research, 5(8), 260-265



©2020 by the Authors. This Article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>)