# Buenas Prácticas en Arquitecturas de Analítica en la Nube
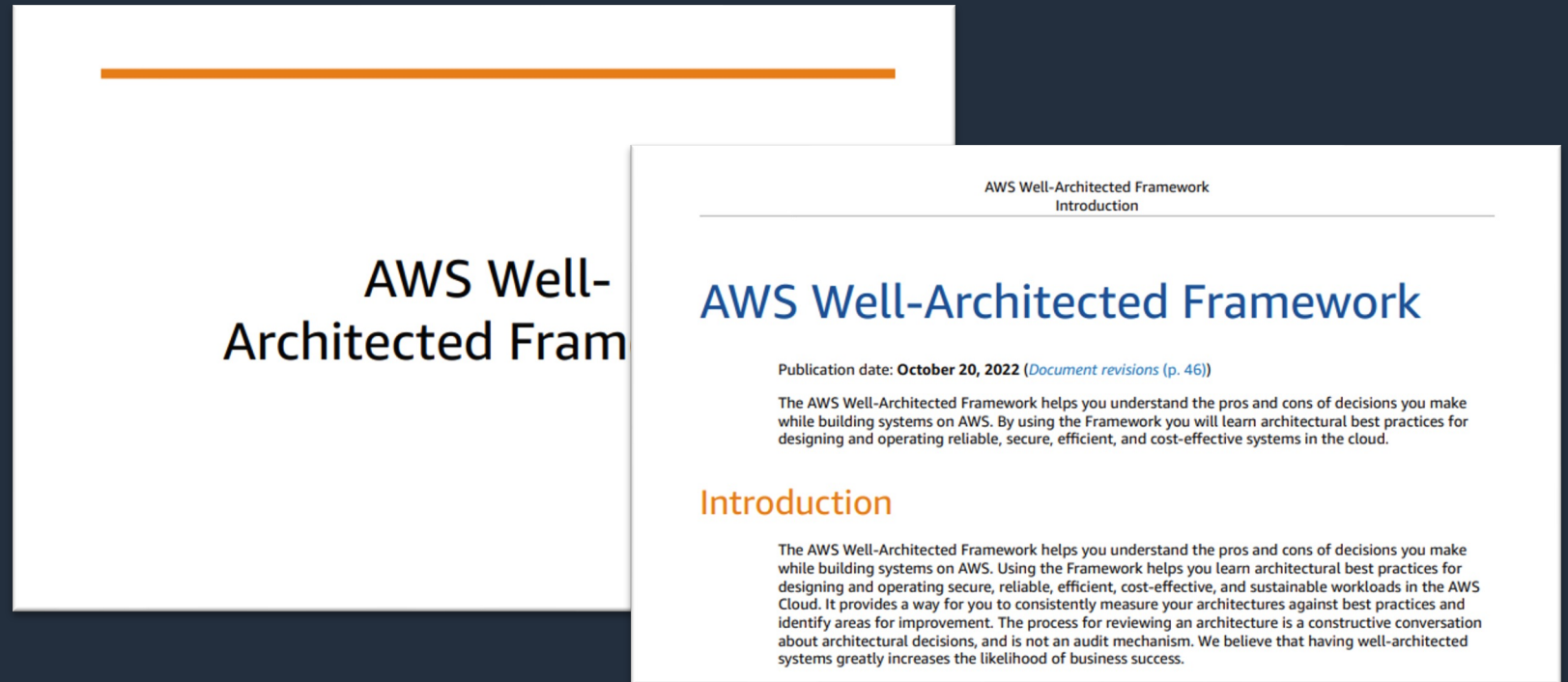
Carlos Paez – Solutions Architect

linkedin.com/in/carlospaez/

# Well-Architected Framework
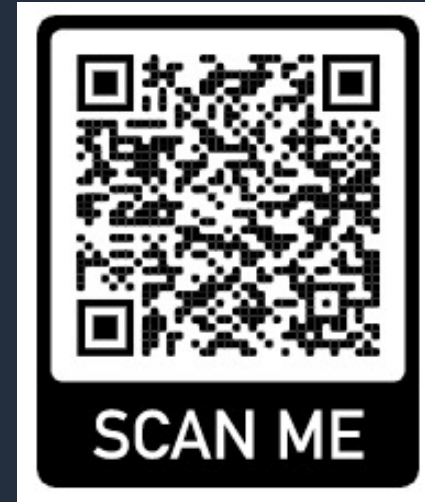
aws

# Well-Architected Framework

The AWS Well-Architected Framework helps you understand the pros and cons of decisions you make while building systems on AWS. By using the Framework you will learn architectural best practices for designing and operating reliable, secure, efficient, and cost-effective systems in the cloud.

- **Pillars**
  Design principles

- **Lenses**
  Best Practices

AWS Well-Architected Framework

**AWS Well-Architected Framework**

Publication date: **October 20, 2022** (*Document revisions* (p. 46))

The AWS Well-Architected Framework helps you understand the pros and cons of decisions you make while building systems on AWS. By using the Framework you will learn architectural best practices for designing and operating reliable, secure, efficient, and cost-effective systems in the cloud.

## Introduction

The AWS Well-Architected Framework helps you understand the pros and cons of decisions you make while building systems on AWS. Using the Framework helps you learn architectural best practices for designing and operating secure, reliable, efficient, cost-effective, and sustainable workloads in the AWS Cloud. It provides a way for you to consistently measure your architectures against best practices and identify areas for improvement. The process for reviewing an architecture is a constructive conversation about architectural decisions, and is not an audit mechanism. We believe that having well-architected systems greatly increases the likelihood of business success.

https://docs.aws.amazon.com/wellarchitected/latest/framework/welcome.html

aws

Well-Architected
Data Analytics Lens

SCAN ME

Data Analytics
AWS Well-Architected F

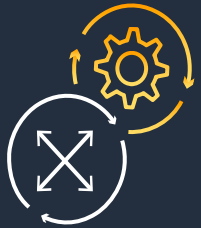Data Analytics Lens AWS Well-Architected Framework

**Data Analytics Lens: AWS Well-Architected Framework**
Copyright © 2022 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

Amazon's trademarks and trade dress may not be used in connection with any product or service that is not Amazon's, in any manner that is likely to cause confusion among customers, or in any manner that disparages or discredits Amazon. All other trademarks not owned by Amazon are the property of their respective owners, who may or may not be affiliated with, connected to, or sponsored by Amazon.

aws

# Analytics Lens

Defines the principles and key architectural elements for designing Analytics Workload based on the 6 pillars of Well-Architected Framework.



| Operational Excellence | Security | Reliability | Performance Efficiency | Cost Optimization | Sustainability |

https://docs.aws.amazon.com/wellarchitected/latest/analytics-lens/analytics-lens.html

aws

# Operational Excellence

Modernize deployment of the analytics jobs and applications

Build financial accountability models for data and workload usage

Monitor the health of the analytics pipelines

https://docs.aws.amazon.com/wellarchitected/latest/analytics-lens/operational-excellence.html

aws

# Monitor the health of the analytics pipelines

# Security

- Classify and protect data

- Control the data access

- Control the access to workload infrastructure

https://docs.aws.amazon.com/wellarchitected/latest/analytics-lens/security.html

aws

# Control the access to workload infrastructure



Centralized permissions

1. Administrator sets up user permissions on lake resources: databases, tables, and columns

Data Steward

Lake Formation

Data Catalog | Access Control

Amazon S3 Data Lake Storage

3. Lake Formation unifies metadata and data access permissions. It **authorizes** access to resources

Amazon Athena

Amazon Redshift

AWS Glue

Amazon EMR

Data Analyst

2. User access data via integrated services

aws

# Reliability

Design resiliency for analytics workload

Govern data and metadata changes

https://docs.aws.amazon.com/wellarchitected/latest/analytics-lens/reliability.html

aws

# Design resiliency for analytics workload

AWS Regions are comprised of multiple AZs for high availability, high scalability, and high fault tolerance. Applications and data are replicated in real time and consistent in the different AZs.

AWS Availability Zone (AZ)

AWS Region

Transit — AZ

AZ — AZ

Transit — AZ

**A Region** is a physical location in the world where we have multiple **Availability Zones.**

Datacenter

Datacenter

Datacenter

**Availability Zones** consist of one or more discrete data centers, each with redundant power, networking, and connectivity, housed in separate facilities.

aws

# Performance Efficiency

Choose the best-performing compute solution

Choose the best-performing storage solution

Choose the best-performing file format and partitioning

https://docs.aws.amazon.com/wellarchitected/latest/analytics-lens/performance-efficiency.html

aws

# Choose the best-performing compute solution

**Amazon Redshift Serverless**

Get insights from data in seconds ●

Experience consistently high performance ●

Get started with no modifications ●

Pay for what you use ●

Automatic scaling ●

Compute provisioning ●

Automated patching ●

Automatic failover ●

Advanced monitoring ●

Backup and recovery ●

Routine maintenance ●

Security and industry compliance ●

## YOU
focus on insights

## aws
takes care of the rest

aws

# Cost Optimization

Choose cost-effective compute and storage solutions based on workload usage patterns

Manage cost over time

Use optimal pricing models based on infrastructure usage patterns

https://docs.aws.amazon.com/wellarchitected/latest/analytics-lens/cost-optimization.html

aws

# Use optimal pricing models based on infrastructure usage patterns

| S3 Intelligent-Tiering | S3 Standard | S3 Standard-IA | **S3 Glacier Instant Retrieval** | S3 Glacier Flexible Retrieval | S3 Glacier Deep Archive |
|---|---|---|---|---|---|

## AWS Region ≥ 3 Availability Zones

| **Changing access patterns** | **Frequently accessed data** | **Infrequently accessed data** | **Rarely accessed data** | **Archive data** | **Long term archive data** |
|---|---|---|---|---|---|
| • Milliseconds access | • Milliseconds access | • Milliseconds access | • Milliseconds access | • Retrieval options from minutes to hours | • Retrieval in hours |
| • No retrieval charge | • No retrieval charge | • Per-GB retrieval charge | • Per-GB retrieval charge | • **Free bulk retrievals** | |
| • **Archive Instant Access tier** | | | | | |

aws

# Sustainability

To be released soon..

aws

# Well-Architected Analytics Review

aws

# Well-Architected Review Process

- Create a plan to fix high risk issues
- Identify a significant workload
- Prepare for review
- Review architecture
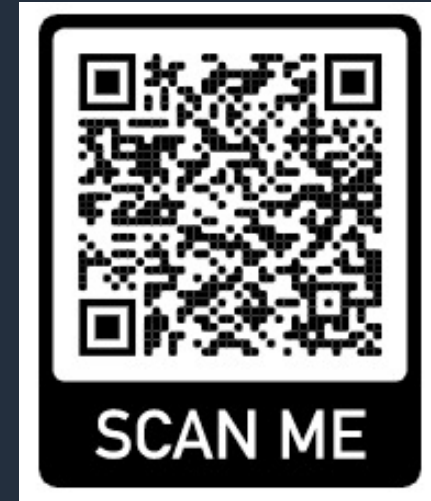- Review results and improvement plan

aws

# Resumen

- Usar Well-Architected Framework – Analytics Lens como guía
- Aprovechar el review como self-assessment de cargas de trabajo
- Adecuar el análisis dependiendo del
  - Escenario
  - Casos de uso
  - Componentes
  - NFR

aws

# Resources

- Well-Architected Framework ([link](#))
- Data Analytics Lens white paper ([link](#))
- Data Classification white paper ([link](#))

- Analytics on AWS ([link](#))
- Modern Data Architecture on AWS ([link](#))

- AWS Global Infrastructure ([link](#))
- AWS S3 Storage Classes ([link](#))



Carlos Paez – Solutions Architect
linkedin.com/in/carlospaez/

aws

# Q&A

aws

¡Gracias!

aws