

# The AWS Well-Architected Framework: A Focus on Machine Learning

**Siva Kumar Mamillapalli**

[siva.mamill@gmail.com](mailto:siva.mamill@gmail.com)

## Abstract

In recent years, machine learning has transitioned from the realm of research and development to becoming widely adopted, driven by the proliferation of data sources and scalable cloud computing resources. AWS customers now leverage AI/ML across diverse applications such as call center operations, personalized recommendations, fraud detection, social media content moderation, audio and video analysis, product design, and identity verification. Industries benefiting from AI/ML include insurance, healthcare, manufacturing, finance, media, and telecom. Machine learning, with its ability to uncover patterns in data through algorithms, empowers users significantly, emphasizing the importance of responsible deployment. AWS is dedicated to developing AI and ML services that are fair and accurate, providing tools and guidance for building responsible AI and ML applications. This paper outlines proven best practices for designing and continuously improving ML workloads, offering guidance and architectural principles applicable across cloud platforms while including specific resources for implementing these practices on AWS.

**Keywords:** AWS, AI/ML, Responsible AI, Well Architected Framework, GenAI, Data Modeling, ML Lens

## 1. Introduction

The AWS Well-Architected Framework helps you assess the benefits and risks associated with decisions made while building workloads on AWS. By using this framework, you gain insight into operational and architectural best practices for designing and managing cloud workloads. It offers a consistent approach to measuring your operations and architectures against best practices, helping you identify areas for improvement.

The quality of your ML models' input data is crucial for producing accurate results. As data evolves over time, continuous monitoring is essential to detect, correct, and address accuracy and performance issues. This may involve retraining your model using the most up-to-date and refined data.

Unlike application workloads that follow specific instructions to solve a problem, ML workloads rely on algorithms that learn from data through an iterative and ongoing process. The ML Lens enhances and extends the Well-Architected Framework to address the unique nature of ML workloads.

## **2. Literature Review**

### **2.1 Well-Architected Framework pillars**

The AWS Well-Architected Framework outlines best practices for designing and managing cloud workloads, structured around six pillars:

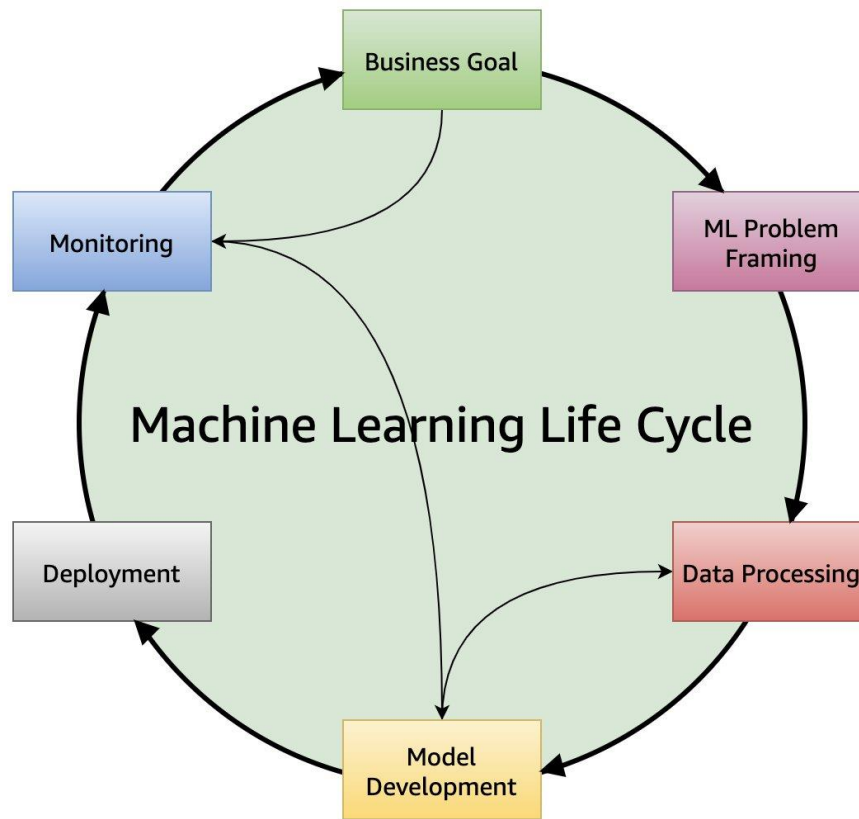
1. **Operational Excellence:** Focuses on running, monitoring, and improving workloads for business value. Key areas: organization, preparation, operation, evolution.
2. **Security:** Protects systems, data, and assets through risk management and mitigation. Key areas: foundations, identity management, detection, protection, incident response, application security.
3. **Reliability:** Ensures workloads can recover from disruptions and consistently perform. Key areas: foundations, architecture, change management, failure management.
4. **Performance Efficiency:** Optimizes resource use to meet requirements and adapt to demand. Key areas: selection, review, monitoring, trade-offs.
5. **Cost Optimization:** Focuses on cost-effective operations and maximizing ROI. Key areas: financial management, usage awareness, resource efficiency, demand and supply management.
6. **Sustainability:** Aims to reduce environmental impact, especially energy consumption. Key areas: region selection, software/architecture alignment, data, hardware, and processes.

### **2.2 Well-Architected machine learning lifecycle**

The ML lifecycle is an iterative, cyclic process that provides best practices across phases to ensure a successful machine learning project. It includes the following phases:

- Business goal identification
- ML problem framing
- Data processing (collection, preprocessing, feature engineering)
- Model development (training, tuning, evaluation)
- Model deployment (inference, prediction)
- Model monitoring

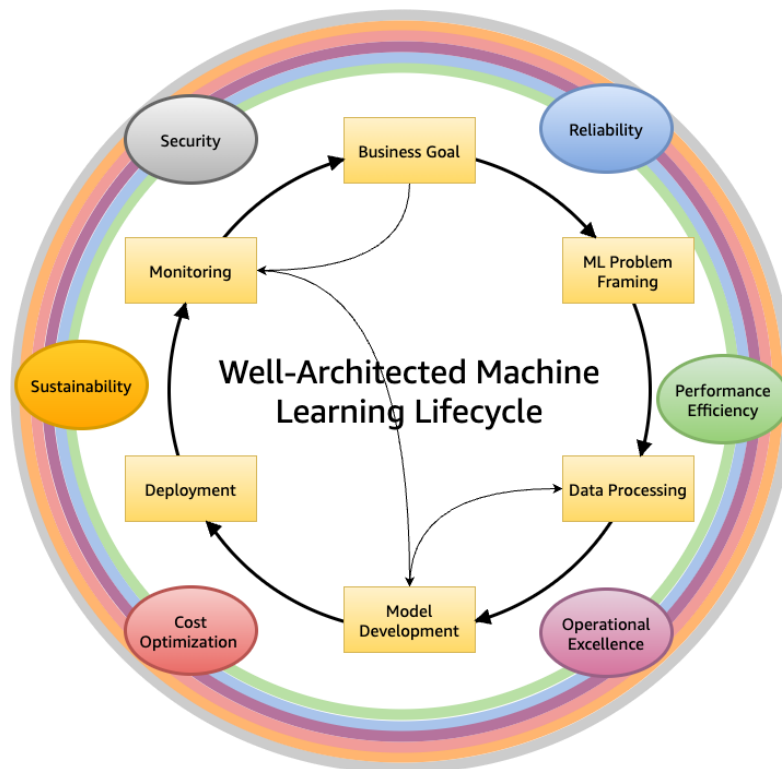
These phases may not be strictly sequential and can involve feedback loops, as shown in Figure 1, to revisit earlier stages.



**Figure 1: Machine Learning Lifecycle**

Here's a brief overview of each phase:

- I. **Business Goal:** The organization should define the problem and its business value, measuring success against specific objectives and criteria.
- II. **ML Problem Framing:** The business problem is reframed as an ML problem, identifying what to predict (label/target) and how to optimize performance and error metrics.
- III. **Data Processing:** Data must be collected, prepared, and processed through feature engineering to create usable input for model training.
- IV. **Model Development:** Involves building, training, tuning, and evaluating models, including creating CI/CD pipelines for automating these processes in staging and production.
- V. **Deployment:** After training, tuning, and evaluation, the model is deployed to production for making predictions and inferences.
- VI. **Monitoring:** A monitoring system ensures the model maintains performance by detecting and addressing issues early.



**Figure 2: Well-Architected ML lifecycle**

## 2.3 Well-Architected machine learning design principles

Well-Architected ML design principles guide the creation of efficient, secure, and scalable ML workloads in the cloud. These principles include:

- I. **Assign Ownership:** Ensure the right skills, resources, and accountability to boost productivity.
- II. **Provide Protection:** Implement security controls to safeguard model data, algorithms, and services.
- III. **Enable Resiliency:** Ensure fault tolerance and recoverability with version control, traceability, and explainability.
- IV. **Enable Reusability:** Use modular components to improve reliability, productivity, and cost efficiency.
- V. **Enable Reproducibility:** Apply version control across components for model governance and auditing.
- VI. **Optimize Resources:** Conduct trade-off analysis to achieve the best outcomes with available resources.
- VII. **Reduce Cost:** Identify cost-saving opportunities through automation and optimization.
- VIII. **Enable Automation:** Leverage CI/CD, pipelining, and continuous training to enhance agility, performance, and cost efficiency.
- IX. **Enable Continuous Improvement:** Continuously monitor, analyze, and evolve the workload.
- X. **Minimize Environmental Impact:** Set sustainability goals and use efficient hardware and software to reduce the environmental footprint.



development and deployment. Additionally, the lens serves as a valuable tool for regularly evaluating and improving existing workloads, enabling you to identify and resolve potential issues early in the process.

## 4. References

1. Mishra, Abhishek. Machine learning in the AWS cloud: Add intelligence to applications with Amazon SageMaker and Amazon Rekognition. John Wiley & Sons, 2019.
2. Amershi et al., "Software Engineering for Machine Learning: A Case Study," 2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP), Montreal, QC, Canada, 2019, pp. 291-300, doi: 10.1109/ICSE-SEIP.2019.00042, 2021
3. AWS Prescriptive Guidance: Planning for successful MLOps, Oct 2022  
<https://docs.aws.amazon.com/pdfs/prescriptive-guidance/latest/ml-operations-planning/ml-operations-planning.pdf>,
4. AWS Architecture Center, Mar 2018: [Reference Architecture Examples and Best Practices](#)
5. Architecture Best Practices for ML, Jun 2021: [https://aws.amazon.com/architecture/machine-learning/?cards-all.sort-by=item.additionalFields.sortDate&cards-all.sort-order=desc&awsf.content-type=\\*all&awsf.methodology=\\*all](https://aws.amazon.com/architecture/machine-learning/?cards-all.sort-by=item.additionalFields.sortDate&cards-all.sort-order=desc&awsf.content-type=*all&awsf.methodology=*all)
6. Amazon AI Fairness and Explainability, Jan 2023 AWS reinvent:  
<https://pages.awscloud.com/rs/112-TZM-766/images/Amazon.AI.Fairness.and.Explainability.Whitepaper.pdf>
7. Tools for building GenAI in AWS , May 2021 :  
<https://pages.awscloud.com/rs/112-TZM-766/images/Amazon.AI.Fairness.and.Explainability.Whitepaper.pdf>
8. Wittig, Andreas, and Michael Wittig. Amazon Web Services in Action: An in-depth guide to AWS. Simon and Schuster. Jan 2021