# From Development to Deployment: Streamlining MLOps with Monoliths, Microservices, and Amazon SageMaker

## Karanbir Singh

*Senior Software Engineer, Salesforce, San Francisco, United States*

---***---

**Abstract -** *In the era of AI, Machine learning models are integral to modern software applications. From spam detectors to self-driving cars, intelligent machine learning models are making their mark. However, the process of transitioning from model development to deployment poses significant challenges. This article aims to explore the model deployment process in detail and compare different deployment strategies such as Deploy as Monolithic, Deploy as Microservices, and Deploy using Amazon SageMaker. It also shed light on how Microservices and Amazon SageMaker can streamline and enhance Machine Learning Operations (MLOps). Additionally, it highlights relevant tools and practices that complement these approaches.*

***Key Words*: MLOps , Microservices, Monoliths, Amazon SageMaker, Kubernetes, Artificial Intelligence, Scalability, Model Development Lifecycle**

## 1. INTRODUCTION

Developing a machine learning model is just the beginning of its lifecycle. To deliver value in a production environment, the model needs to be deployed efficiently, scaled to meet demand, and maintained over time. This process, commonly referred to as Machine Learning Operations (MLOps), encompasses the activities that ensure the model performs as expected when integrated into a broader application system.

This article outlines the key stages of model development and explores the architectural choices for deploying models in a production environment, focusing on Kubernetes and Amazon SageMaker.

### 1.1. Model Development: A Brief Overview

The process of model development typically involves the following steps:

- **Data Analysis**: Understanding the data, cleaning it, and preparing it for training is the foundation of any machine learning model. This phase involves identifying patterns, relationships, and anomalies in the dataset.

- **Algorithm Selection**: Depending on the use case, a suitable algorithm is chosen. This could range from simple regression models to complex deep learning architectures.

- **Training**: The model is trained on the dataset to learn the underlying patterns and relationships. The goal is to optimize the model to generalize well to new, unseen data.

- **Model Evaluation**: After training, the model is evaluated to ensure it is neither underfitting (failing to capture important patterns) nor overfitting (capturing noise in the data as if it were important).

- **Efficiency Computation**: The model's efficiency is computed in terms of its performance metrics, such as accuracy, precision, recall, or F1 score, depending on the use case.

## 2. Architectural Approaches for Model Deployment

Understanding how the model will be used in production as well as target audience is essential to guiding architectural choices for deployment

### 2.1. Monolithic Approach

The monolithic approach involves deploying the entire system as a single, unified unit. This means that all components—whether they pertain to the user interface, business logic, or machine learning models—are tightly coupled and deployed together as a single application.

### 2.1.1. Example Use Case: Car Dealership Application

Let's consider a medium-sized car dealership that aims to provide personalized car recommendations to its customers based on their preferences. Since the target audience is relatively small, and the system does not require handling multiple versions of the application or models simultaneously, a monolithic architecture can be a practical choice.
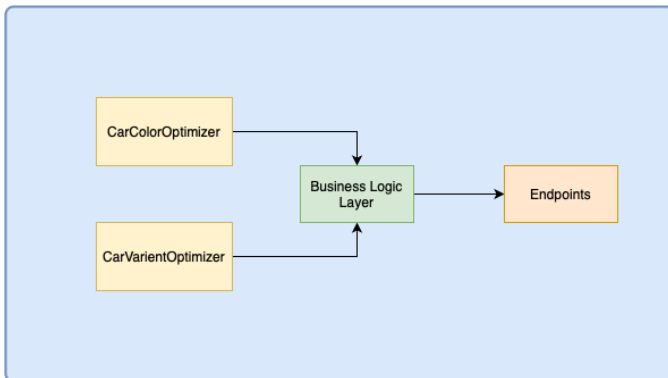
**Figure 1.** *Architecture of the Monolith Deployment for the Car Customization Application.*

In this case, the machine learning model that powers the car recommendation engine is embedded directly within the application, along with the user interface and business logic. All the components are bundled together into a single deployable unit, making it easier to manage and deploy.

**Pros:**

● **Easy to maintain**: Since all components are part of the same application, maintaining the system is straightforward. Updates, bug fixes, and new features can be implemented in a unified manner, without the need to coordinate changes across multiple services.

● **Quick deployment**: The entire system can be deployed at once, reducing the complexity of the deployment phase. This can be especially beneficial in environments where frequent updates or multiple versions are not required.

**Cons:**

● **Lack of scalability**: One of the key challenges of the monolithic approach is scalability. If any part of the system needs to scale to handle increased demand (e.g., the recommendation engine during a promotional event), the entire application must be scaled as a whole. This can lead to inefficient resource usage.

● **Difficult to update individual components**: As the system grows, making updates to individual components can become more challenging. For instance, if the machine learning model needs to be retrained and updated, it may require redeploying the entire application, leading to potential downtime and increased deployment complexity.

While the monolithic approach may work well for smaller projects or systems with limited scalability requirements,

it can struggle as the system grows or when components need to be independently updated or scaled.

## 2.2. Microservices Approach

In the software industry, microservices architecture is a method of structuring applications as a collection of smaller, independent services that communicate with each other via APIs [1]. Each service operates as a standalone component, allowing for greater flexibility in development, deployment, and scaling. This approach contrasts with monolithic architecture, where all components are bundled together and deployed as a single unit.

### 2.2.1. Example Use Case: Car Company Application in a Private Cloud

Building on our previous example, let's consider that a car company wants to create an application for customizing car options. However, due to concerns about data exposure, they prefer to build the system in their private cloud. In this scenario, the microservices architecture is a suitable choice, as it allows for independent scaling and management of services while maintaining control over data security within their private infrastructure.

For instance, the car company might need to deploy different versions of the application for various car models and variants, such as electric vehicles, luxury sedans, or SUVs. Each variant may require different machine learning models for recommendation engines, user customization, or predictive maintenance. The ability to deploy and manage these models as separate microservices ensures that updates and scaling can occur independently, without affecting the entire system.
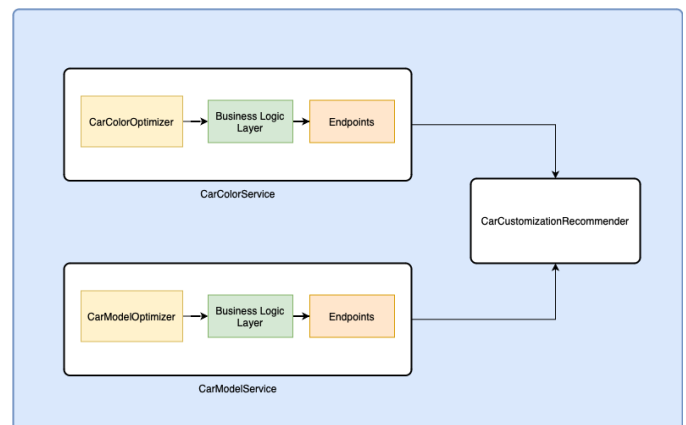


**Figure 2.** *Architecture of the Microservices Deployment for the Car Customization Application*

Deploying the system in a private cloud adds an extra layer of security, as all data remains within the company's

controlled environment. This is especially important when handling sensitive customer data, proprietary algorithms, or intellectual property related to car designs and features. By using microservices in a private cloud, the car company can balance scalability with security.

## 2.2.2. Microservices with Kubernetes

To efficiently manage this architecture, Kubernetes can be employed to orchestrate and scale containerized microservices. Kubernetes is a powerful tool for automating the deployment, scaling, and operation of application containers. [2] When used in conjunction with microservices in a private cloud, Kubernetes provides robust control over the system's infrastructure while ensuring that the application remains secure and scalable.

**Pros:**

- **Quick Deployment**: Kubernetes automates much of the deployment process, enabling rapid and consistent releases across the private cloud environment. This allows the car company to roll out new features or model updates without risking disruptions to the entire system.

- **Auto-scaling**: Kubernetes can automatically scale individual microservices based on traffic and load. This ensures efficient resource utilization, allowing the system to handle varying levels of demand—such as during a new car launch—without compromising performance.

- **Ease of Maintenance**: Kubernetes abstracts away much of the complexity involved in maintaining containerized applications. It offers built-in tools for monitoring, logging, and service discovery, which simplify the management of microservices over time.

Kubernetes enables the deployment of machine learning models and microservices in a highly scalable, maintainable, and efficient manner. For environments that require multiple models or services to be managed, updated, and scaled independently, Kubernetes excels. Additionally, Kubernetes offers advanced features such as rolling updates, service discovery, and self-healing, ensuring that the application remains available and resilient, even during updates or unexpected failures.

By deploying microservices within a private cloud and orchestrating them with Kubernetes, the car company can achieve a balance between scalability, security, and operational efficiency. This approach ensures that sensitive data remains protected while still benefiting from the flexibility and robustness of modern cloud-native technologies.

## 2.3. Amazon SageMaker Approach

Amazon SageMaker is a fully managed service from AWS that allows data scientists and developers to build, train, and deploy machine learning models at scale. SageMaker simplifies the entire machine learning lifecycle by providing a set of tools and infrastructure that automates much of the process. [3]

**Key Features of Amazon SageMaker:**

- **Managed Infrastructure**: Amazon SageMaker abstracts away the need to manage underlying infrastructure. It computes resources, handles scaling, and ensures that the environment is configured optimally for the model's requirements.

- **Built-in Algorithms and Frameworks**: SageMaker supports a wide range of built-in algorithms and pre-configured environments for popular frameworks like TensorFlow, PyTorch, and Scikit-learn. This reduces the time needed to set up the development environment.

- **Automatic Model Tuning**: SageMaker provides automated hyperparameter tuning to optimize model performance. It uses machine learning to search for the best set of hyperparameters, making it easier to achieve optimal results without manual intervention.

- **One-Click Deployment**: Once a model is trained, SageMaker allows for one-click deployment, where the model can be launched as an endpoint that automatically scales based on traffic. This simplifies the process of making the model accessible to production applications.

- **Multi-Model Endpoints**: SageMaker supports multi-model endpoints, enabling multiple models to be hosted on a single endpoint. This can reduce costs by consolidating resources and simplifying the architecture.

- **Monitoring and Logging**: SageMaker integrates with AWS CloudWatch to monitor deployed models, providing insights into performance and operational metrics. This is crucial for maintaining models in production and ensuring they continue to perform well over time.

**Pros:**

- **End-to-End Management**: SageMaker covers the entire machine learning lifecycle, from data preparation to deployment, reducing the need for multiple tools and integrations.

- **Scalability**: Amazon SageMaker automatically scales the infrastructure to meet demand, ensuring that models can handle varying levels of traffic.

- **Ease of Use**: SageMaker's interface and integrations with other AWS services make it user-friendly, even for teams without extensive DevOps experience.

**Cons:**

- **Dependency on AWS Ecosystem**: While SageMaker provides a seamless experience within the AWS ecosystem, it may not be the best choice for organizations using multi-cloud strategies or those seeking to avoid vendor lock-in.

- **Cost**: SageMaker's managed service can incur significant costs, particularly for large-scale deployments or models that require constant tuning and monitoring.

## 3. Other Relevant Tools and Practices

In addition to Kubernetes and Amazon SageMaker, several tools and practices complement MLOps pipelines:

- **Pickling a Model**: Pickling is a process for serializing a Python object, such as a trained machine learning model, so it can be saved to disk and later reloaded. This technique is commonly used for model persistence, allowing models to be easily loaded and deployed in production environments.

- **Jenkins**: Jenkins is an open-source automation server that facilitates continuous integration and continuous delivery (CI/CD). It can be used to automate the deployment of machine learning models by integrating with Kubernetes and other tools to streamline the build, test, and deployment processes.

- **Docker**: Docker containers are essential to modern MLOps pipelines. Docker enables the packaging of applications, including their dependencies, into isolated containers that can be deployed consistently across different environments. Both Kubernetes and Amazon SageMaker leverage Docker containers for model deployment.

## 4. Conclusion

Deploying machine learning models in production requires careful consideration of the architectural approach. While monolithic systems may suffice for small projects, Kubernetes and Amazon SageMaker offer the scalability and flexibility needed for larger, more complex deployments. Kubernetes, in particular, stands out for its

ability to automate deployment and scaling, making it an ideal choice for teams with containerization expertise. Amazon SageMaker simplifies the machine learning lifecycle through a fully managed service, making it a great option for teams seeking to accelerate deployment with minimal infrastructure management.

The integration of complementary tools like Jenkins and Docker further enhances the efficiency and reliability of model deployment pipelines. As machine learning continues to drive innovation, MLOps practices will play a critical role in ensuring that models deliver consistent, real-world value.

## REFERENCES

[1] Fowler, M. (2014). *Microservices: A definition of this new architectural term.* Retrieved from https://martinfowler.com/articles/microservices.html

[2] Kubernetes Documentation. (n.d.). *Kubernetes documentation.* Retrieved from https://kubernetes.io/docs/

[3] Amazon Web Services. (n.d.). *Amazon SageMaker developer guide.* Retrieved from https://docs.aws.amazon.com/sagemaker/

## BIOGRAPHIES

Karanbir Singh is an expert in MLOps with a specialization in tools like Kubernetes and Amazon SageMaker. Karanbir holds a MS in Software Engineering from San Jose State University, United States. With a strong educational background in Computer Science and Engineering, they have honed their skills in cutting-edge AI techniques. Currently working at Salesforce, Karanbir is dedicated to making complex technical concepts accessible and practical for a wide audience. Their work reflects a deep commitment to advancing the field of artificial intelligence and empowering others with knowledge.