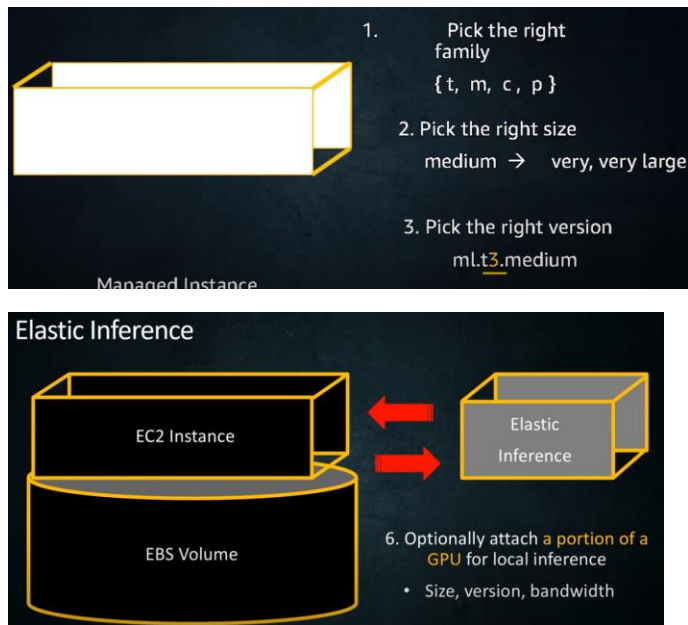


## Domain X: Misc



## X.1 SageMaker Deep Dive

### X.1.1 Fully Managed Notebook Instances with Amazon SageMaker



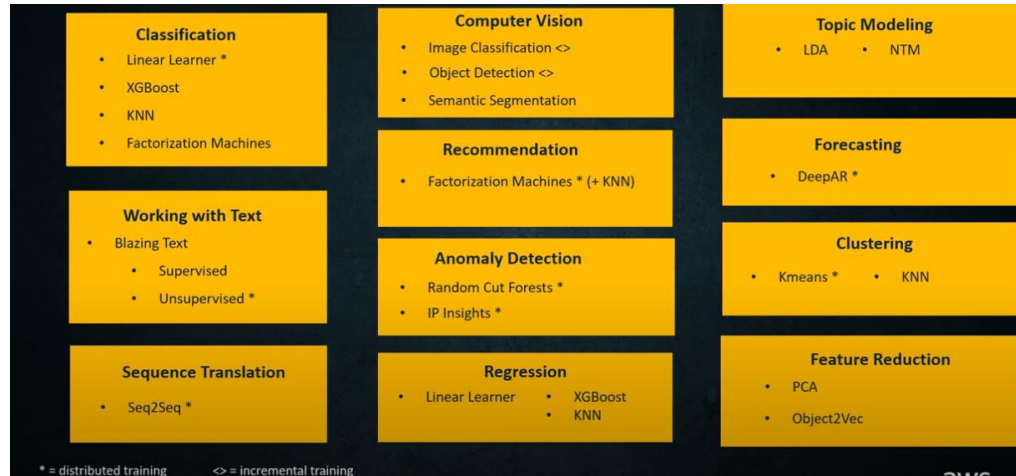
#### Elastic Inference

Elastic Inference is a service that allows attaching a portion of a GPU to an existing EC2 instance.<sup>2</sup> This approach is particularly useful when running inference locally on a notebook instance.<sup>2</sup> By selecting an appropriate Elastic Inference configuration based on size, version, and bandwidth, users can accelerate their inference tasks without needing a full GPU.<sup>2</sup>

#### Use Cases for Elastic Inference

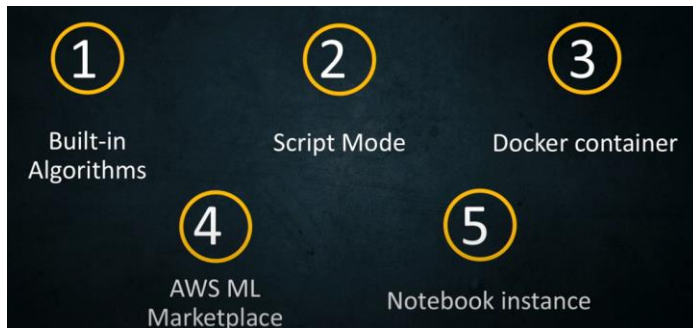
- You need to run inference tasks locally on your notebook instance.
- Your workload benefits from GPU acceleration but doesn't require a full GPU.
- You want to optimize cost by only paying for the portion of GPU resources used.

## X.1.2 SageMaker Built-in Algorithms



Task Category	Algorithms	Supervised/UnSupervised
Classification	<ul style="list-style-type: none"> <li>Linear Learner (distributed)</li> <li>XGBoost</li> <li>KNN</li> <li>Factorization Machines</li> </ul>	Supervised
Regression	<ul style="list-style-type: none"> <li>Linear Learner</li> <li>XGBoost</li> <li>KNN</li> </ul>	Supervised
Computer Vision	<ul style="list-style-type: none"> <li>Object Detection (incremental)</li> <li>Semantic Segmentation</li> </ul>	Supervised
Working with Text	<ul style="list-style-type: none"> <li>BlazingText</li> </ul>	Supervised / Unsupervised
Sequence Translation	<ul style="list-style-type: none"> <li>Seq2Seq (distributed)</li> </ul>	Supervised
Recommendation	<ul style="list-style-type: none"> <li>Factorization Machines (distributed)</li> <li>KNN</li> </ul>	Supervised
Anomaly Detection	<ul style="list-style-type: none"> <li>Random Cut Forests (distributed)</li> <li>IP Insights (distributed)</li> </ul>	Unsupervised / Semi-supervised
Topic Modeling	<ul style="list-style-type: none"> <li>LDA</li> <li>NTM</li> </ul>	Unsupervised
Forecasting	<ul style="list-style-type: none"> <li>DeepAR (distributed)</li> </ul>	Supervised
Clustering	<ul style="list-style-type: none"> <li>K-means (distributed)</li> <li>KNN</li> </ul>	Unsupervised
Feature Reduction	<ul style="list-style-type: none"> <li>PCA</li> <li>Object2Vec</li> </ul>	Unsupervised / Semi-supervised

### X.1.3 SageMaker Training types



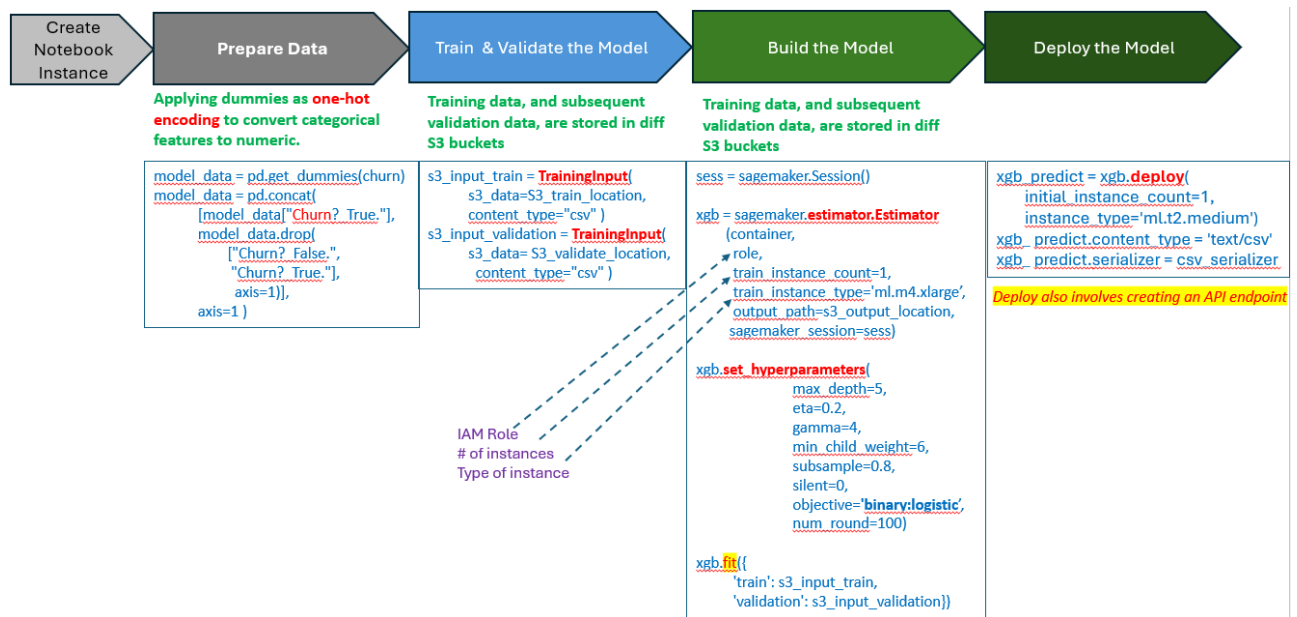
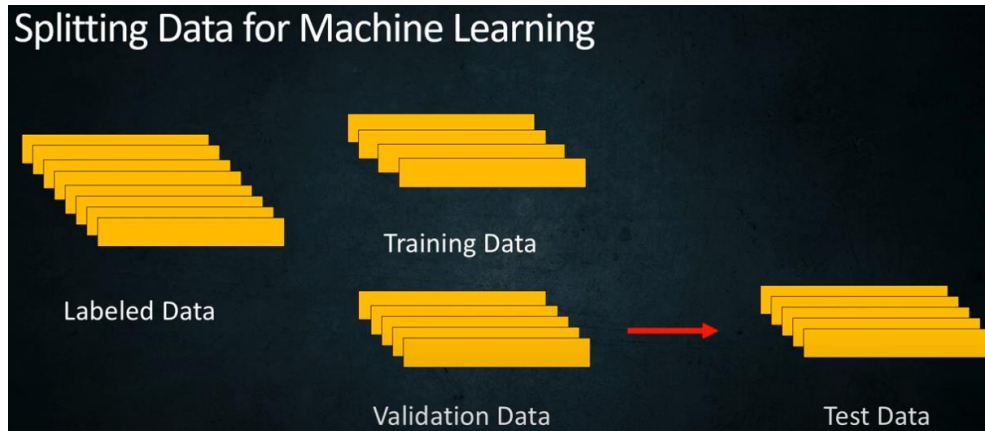
Training Type	Description	When to Use
<b>1. Built-in Algorithms</b>	Pre-configured algorithms provided by Amazon SageMaker, optimized for performance and ease of use	<ul style="list-style-type: none"><li>• Working on common ML tasks (e.g., classification, regression)</li><li>• When you need a quick start without deep ML expertise</li></ul>
<b>2. Script Mode</b>	Custom training scripts using popular ML frameworks (e.g., TensorFlow, PyTorch, Scikit-learn)	<ul style="list-style-type: none"><li>• You have existing scripts in popular ML frameworks</li><li>• For customizing model architecture while leveraging SageMaker's infrastructure</li></ul>
<b>3. Docker Container</b>	Custom Docker containers with your own algorithms or environments	<ul style="list-style-type: none"><li>• Need complete control over training env.</li><li>• Custom or proprietary algorithms</li><li>• For complex, multi-step training pipelines</li></ul>
<b>4. AWS ML Marketplace</b>	Pre-built algorithms and models from third-party vendors available through the AWS Marketplace	<ul style="list-style-type: none"><li>• Need industry-specific or specialized models</li><li>• When you want to explore alternative solutions without building from scratch</li></ul>
<b>5. Notebook Instance</b>	Interactive development and training using Jupyter notebooks on managed instances	<ul style="list-style-type: none"><li>• During the initial stages of model development</li><li>• When you need an interactive environment for debugging and visualization</li></ul>

#### Key Considerations:

- Skill Level: Built-in Algorithms and Marketplace for beginners, Script Mode and Containers for more advanced users
- Customization Needs: From low (Built-in) to high (Containers)
- Development Speed: Notebooks for rapid prototyping, Built-in for quick deployment, Containers for complex but reproducible setups
- Scale: Consider moving from Notebooks to other options as your data and model complexity grow.

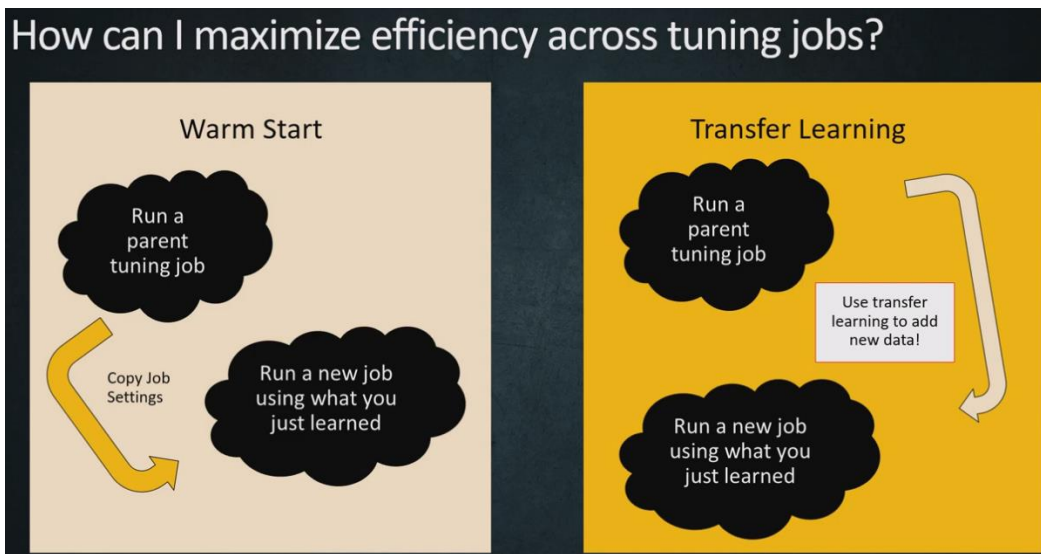
## X.1.4 Train Your ML Models with Amazon SageMaker

### Splitting Data for ML

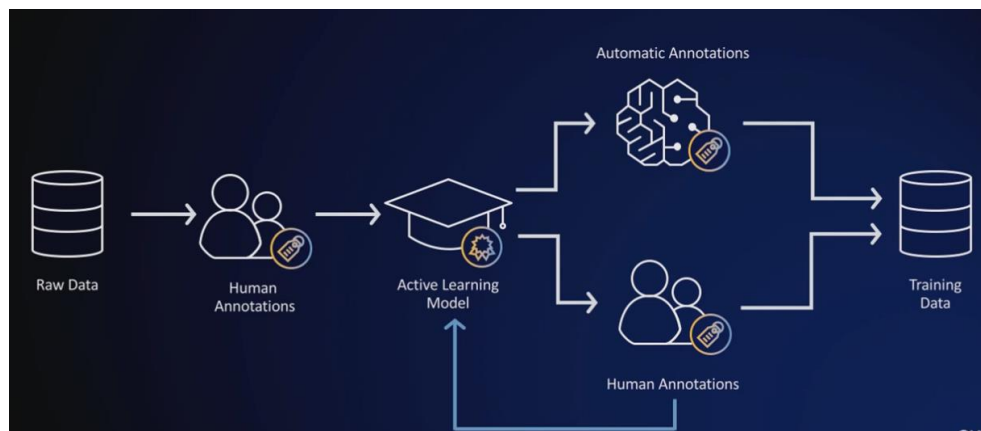


## X.1.5 Tuning Your ML Models with Amazon SageMaker

### Maximizing Efficiency across tuning jobs



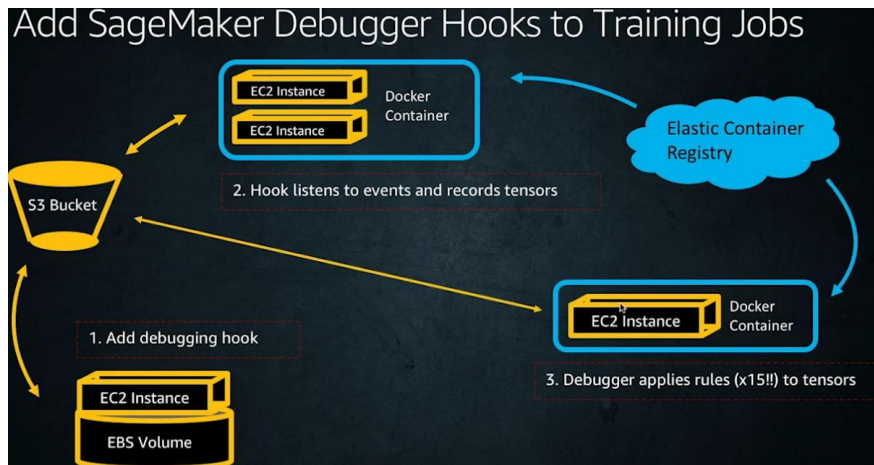
## X.1.6 Tuning Your ML Models with Amazon SageMaker



### How to automate

Put a check to see if Accuracy falls below a % (e.g. > 80%), invoke Human in the loop

## X.1.6 Add Debugger to Training Jobs in Amazon SageMaker



### How it works

1. Add debugging hook:
  - An EC2 instance with an attached EBS volume is used to initiate the process.
  - The debugging hook is added to the training job configuration.
2. Hook listens to events and records tensors:
  - Docker containers running on EC2 instances are used for the training job.
  - The hook listens for specific events during the training process and records tensor data.
3. Debugger applies rules to tensors:
  - Another EC2 instance with a Docker container is used for debugging.
  - The debugger applies predefined rules (mentioned as "x15!!" in the image) to the recorded tensor data.

### Benefits of debugger

DeadRelu, ExplodingTensor,  
PoorWeightInitialization,  
SaturatedActivation,  
VanishingGradient,  
WeightUpdateRatio, AllZero,  
ClassImbalance, Confusion,  
LossNotDecreasing, Overfit,  
Overtraining, SimilarAcrossRuns,  
TensorVariance, UnchangedTensor,  
CheckInputImages,  
NLPSequenceRatio, TreeDepth

entry\_point = 'mnist.py'  
and 1P SM algos

Rules: Built-in & BYO

No Change Needed

Visualization

1. **Comprehensive Built-in Rules/Algorithms:** The debugger offers a wide range of built-in rules to detect common issues in machine learning models, such as:
  - DeadRelu, ExplodingTensor, PoorWeightInitialization
  - SaturatedActivation, VanishingGradient
  - WeightUpdateRatio, AllZero, ClassImbalance
  - Confusion, LossNotDecreasing, Overfit
  - Overtraining, SimilarAcrossRuns
  - TensorVariance, UnchangedTensor
  - CheckInputImages, NLPSequenceRatio, TreeDepth



2. **Customizable** (BYO - Bring Your Own): Users can create and add their own custom debugging rules.
3. **Easy Integration**: The entry point is `'mnist.py'` and it works with SageMaker's built-in algorithms (1P SM algos), suggesting easy integration with existing SageMaker workflows.
4. **No Code Changes Required**: The "No Change Needed" text implies that adding debugging capabilities doesn't require modifying the existing model code.
5. **Visualization**: The debugger provides visualization capabilities, as indicated by the image on the right, which appears to show a tensor or weight distribution.
6. **Real-time Monitoring**: The variety of rules suggests that the debugger can monitor various aspects of model training in real-time, helping to identify issues as they occur.

### X.1.7 Deployment using SageMaker

Deployment Strategy	Description	When to Use
Blue/Green Deployment with Linear Traffic Shifting	<b>Gradually shift traffic</b> from the old version (blue) to the new version (green) over time	<ul style="list-style-type: none"> <li>• When you need <b>fine-grained control over the traffic shift</b></li> <li>• For critical applications <b>requiring minimal risk</b></li> <li>• When you have the resources to run two full environments simultaneously</li> </ul>
Canary Deployment	Release a new version to a <b>small subset of users</b> before rolling it out to the entire infrastructure	<ul style="list-style-type: none"> <li>• When you want to <b>test in production with real users</b></li> <li>• For <b>early detection of issues</b> before full deployment</li> <li>• When you have a diverse user base</li> </ul>
A/B Testing	<b>Run two versions simultaneously</b> and compare their performance based on metrics	<ul style="list-style-type: none"> <li>• When <b>you want to test specific features</b> or changes</li> <li>• When you need to optimize based on user behavior or business metrics</li> </ul>
Rolling Deployment	<b>Gradually replace instances</b> of the old version with the new version	<ul style="list-style-type: none"> <li>• When you have limited resources and can't run two full environments</li> <li>• For <b>applications that can handle mixed versions</b></li> <li>• When you <b>need to minimize downtime</b></li> </ul>



## X.2 From the Exam Guide

### X.2.1 Domain 1: Data Preparation for Machine Learning (ML)

#### Data formats and ingestion mechanisms

	Format	Type	Description	Advantages	Common Use Cases
ROW BASED	CSV	Text	Simple tabular format	<ul style="list-style-type: none"><li>• Human-readable</li><li>• Widely supported</li><li>• Easy to generate</li></ul>	<ul style="list-style-type: none"><li>• Simple data exchange</li><li>• Small to medium datasets</li></ul>
	JSON	<ul style="list-style-type: none"><li>• Semi-structured</li><li>• Text</li></ul>	Flexible format	<ul style="list-style-type: none"><li>• Human-readable</li><li>• Supports nested structures</li><li>• Language-independent</li></ul>	<ul style="list-style-type: none"><li>• Web APIs</li><li>• Configuration files</li><li>• Document databases</li></ul>
	Apache Avro	Binary	Data serialization	<ul style="list-style-type: none"><li>• Compact serialization</li><li>• Language-independent</li></ul>	<ul style="list-style-type: none"><li>• Data serialization</li><li>• RPC protocols</li><li>• <b>Hadoop</b> data storage</li></ul>
	RecordIO	Binary	SageMaker-specific format for efficient data loading	<ul style="list-style-type: none"><li>• <b>Optimized for SageMaker</b></li><li>• Supports large datasets</li></ul>	<ul style="list-style-type: none"><li>• SageMaker model training</li><li>• Large-scale ML datasets</li></ul>
COLUMNAR	Apache Parquet	Binary	Optimized format	<ul style="list-style-type: none"><li>• Efficient compression</li><li>• Fast query performance</li><li>• Schema evolution support</li></ul>	<ul style="list-style-type: none"><li>• Big data analytics</li><li>• Data warehousing</li><li>• Machine learning datasets</li></ul>
	Apache ORC (Optimized Row Columnar)	Binary	Optimized for <b>Hadoop</b>	<ul style="list-style-type: none"><li>• High compression ratio</li><li>• Fast data processing</li></ul>	<ul style="list-style-type: none"><li>• Hive data storage</li><li>• Big data processing</li><li>• <b>Analytics</b> workloads</li></ul>

## Core AWS data sources

Feature	S3	EFS	FSx for NetApp ONTAP
Best for	Large, infrequent-change data	Shared, frequent-change data	High-performance, multi-protocol
Latency	Higher	Low	Lowest
Scalability	Virtually unlimited	Up to petabytes	Up to hundreds of petabytes
Cost	Lowest	Moderate	Highest
ML Use Case	Training data, model artifacts	Distributed training, real-time	High-performance computing, Windows ML
Storage Type	Object storage	File storage	High-performance file storage
Access Pattern	Good for sequential access	Good for random access	Excellent for all access patterns
Shared Access	Not native	Native	Native
Protocols	S3 API, HTTP/S	NFS	NFS, SMB, iSCSI

## Domain 2: ML Model Development

### Common regularization techniques

Technique	Description	Benefits	Best Used When
Dropout	Randomly "drops out" a proportion of neurons during training	<ul style="list-style-type: none"><li>• Reduces overfitting</li><li>• Improves generalization</li><li>• Acts as an ensemble method</li><li>• Prevents co-adaptation of features</li></ul>	<ul style="list-style-type: none"><li>• Large neural networks</li><li>• Limited training data</li><li>• Complex tasks with risk of overfitting</li></ul>
Weight Decay (L2)	Adds a penalty term to the loss function based on the squared magnitude of weights	<ul style="list-style-type: none"><li>• Prevents large weights</li><li>• Improves generalization</li><li>• Stabilizes learning</li><li>• Helps with feature selection</li></ul>	<ul style="list-style-type: none"><li>• Most neural networks</li><li>• When you want to keep all features but reduce their impact</li><li>• When dealing with multicollinearity</li></ul>
L1 Regularization	Adds a penalty term to the loss function based on the absolute value of weights	<ul style="list-style-type: none"><li>• Encourages sparsity in the model</li><li>• Feature selection (drives some weights to zero)</li><li>• Robust to outliers</li><li>• Computationally efficient for sparse data</li></ul>	<ul style="list-style-type: none"><li>• When feature selection is important</li><li>• Dealing with high-dimensional data</li><li>• When you want a sparse model</li></ul>

### Open Source frameworks for SageMaker script mode - TensorFlow vs PyTorch

Feature	TensorFlow	PyTorch
Type	Open-source ML framework	Open-source ML framework
Specialization	General-purpose, excels in production deployment	Flexible, great for research and prototyping
Distributed Training	Supported via Horovod or parameter servers	Supported via PyTorch Distributed
GPU Acceleration	Fully supported	Fully supported
Model Serving	Native support in SageMaker	Native support in SageMaker
Automatic Model Tuning	Supported	Supported

## Domain 4: ML Solution Monitoring, Maintenance, and Security

### Design principles for ML lenses relevant to monitoring

Principle	Key Points
1. Continuous Monitoring	<ul style="list-style-type: none"><li>• Real-time monitoring with CloudWatch</li><li>• SageMaker Model Monitor for quality checks</li><li>• Set up alerts for key metrics</li></ul>
2. Automated Remediation	<ul style="list-style-type: none"><li>• Auto-scaling policies for endpoints</li><li>• Automated model retraining triggers</li><li>• AWS Lambda for automated responses</li></ul>
3. Data Quality Assurance	<ul style="list-style-type: none"><li>• Monitor input data drift</li><li>• Implement data validation checks</li><li>• Use Amazon Athena for ad-hoc queries</li></ul>
4. Model Performance Tracking	<ul style="list-style-type: none"><li>• Track accuracy, latency, throughput</li><li>• A/B testing for model comparisons</li><li>• SageMaker Experiments for version logging</li></ul>
5. Explainability and Interpretability	<ul style="list-style-type: none"><li>• SageMaker Clarify for bias detection</li><li>• SHAP values for interpretability</li><li>• Maintain model cards for documentation</li></ul>
6. Security and Compliance	Encryption at rest and in transit /IAM roles / audit with AWS CloudTrail
7. Cost Optimization	Monitor and optimize resource utilization/ Auto-scaling /Use Spot Instances
8. Scalability and Elasticity	Horizontal scaling/Multi-model endpoints for efficiency/Caching strategies
9. Fault Tolerance and High Availability	Multi AZs/Circuit breakers and fallbacks/Use multi-model endpoints
10. Operational Excellence	IaC with CloudFormation/ AWS Step Functions for ML workflows

### How to use AWS CloudTrail to log, monitor, and invoke re-training activities

Aspect	Description	Key Points
Logging	Record API calls and events	...
Monitoring	Track ML-related activities	...
Re-training Triggers	Use events to initiate re-training	<ul style="list-style-type: none"><li>• Set up CloudWatch Events rules based on CloudTrail logs</li><li>• Trigger Lambda functions for automated re-training</li><li>• Integrate with Step Functions for complex workflows</li></ul>

**Monitoring and observability tools to troubleshoot latency and performance issues (for example, AWS X-Ray, Amazon CloudWatch Lambda Insights, Amazon CloudWatch Logs Insights)**

Tool	Key Features	Use Cases	Benefits for Troubleshooting
AWS X-Ray	<ul style="list-style-type: none"> <li>• Distributed tracing</li> <li>• Service map visualization</li> <li>• Trace analysis</li> <li>• Integration with many AWS services</li> </ul>	<ul style="list-style-type: none"> <li>• End-to-end request tracking</li> <li>• Identifying bottlenecks</li> <li>• Analyzing service dependencies</li> </ul>	<ul style="list-style-type: none"> <li>• Visualize application's component interactions</li> <li>• Pinpoint exact location of performance issues</li> <li>• Understand downstream impact of issues</li> </ul>
CloudWatch Lambda Insights	ONLY FOR Lambda functions		
CloudWatch Logs Insights	<ul style="list-style-type: none"> <li>• Log query and visualization</li> <li>• Built-in and custom queries</li> </ul>	<ul style="list-style-type: none"> <li>• ...</li> </ul>	<ul style="list-style-type: none"> <li>• ....</li> </ul>

**Rightsizing instance - SageMaker Inference Recommender vs AWS Compute Optimizer)**

Tool	Purpose	Key Features	Benefits
SageMaker Inference Recommender	Optimize ML model deployment	<ul style="list-style-type: none"> <li>• Automated benchmarking</li> <li>• Instance type recommendations</li> <li>• Performance vs. cost analysis</li> </ul>	<ul style="list-style-type: none"> <li>• Improved inference performance</li> <li>• Cost optimization for ML workloads</li> </ul>
AWS Compute Optimizer	Optimize EC2 instance types	<ul style="list-style-type: none"> <li>• ML-powered recommendations</li> <li>• Right-sizing suggestions</li> </ul>	Performance/Savings boost

## Appendix:

### Analytics Tools Summary

Service	Description	Primary ML Use Case
Amazon Athena	Serverless query for S3	Ad-hoc analysis of ML datasets
Amazon EMR	Managed big data platform	Large-scale data processing for ML
AWS Glue	Serverless data integration service	ETL for ML data preparation
AWS Glue DataBrew	Visual data preparation tool	Feature engineering and data cleaning
AWS Glue Data Quality	Automated data quality checks	Ensuring ML data quality and consistency
Amazon Kinesis	Real-time data streaming platform	Stream processing for ML applications
Amazon Kinesis Data Firehose	Real-time streaming data delivery	Ingesting streaming data for ML models
AWS Lake Formation	Centralized data lake service	Building secure data lakes for ML
Amazon Managed Service for Apache Flink	Serverless Apache Flink applications	Real-time data processing for ML
Amazon OpenSearch Service	Distributed search and analytics	Log analytics and ML model monitoring
Amazon QuickSight	Business intelligence service	Visualizing ML insights and predictions
Amazon Redshift	Data warehousing service	Large-scale data analysis for ML

### AWS Secrets Manager

AWS Secrets Manager is a secrets management service that helps you protect access to your applications, services, and IT resources.

#### Key Features:

- Secure Storage: Encrypts and stores secrets (e.g., passwords, API keys)
- Rotation: Automates the rotation of secrets
- Fine-grained Access Control: Uses IAM policies to control access
- Auditing: Integrates with AWS CloudTrail for auditing
- Cross-Region Replication: Supports replication of secrets across regions

AWS Storage Gateway

AWS Storage Gateway is a hybrid storage service that enables on-premises applications to seamlessly use AWS cloud storage.

By using AWS Storage Gateway, organizations can seamlessly integrate their on-premises IT environments with AWS cloud storage, enabling hybrid cloud use cases and facilitating cloud migration strategies

Machine Learning:

Service	Primary Function	Key Features/Use Cases
Amazon Augmented AI (A2I)	Human review of ML predictions	<ul style="list-style-type: none"><li>Improves ML model accuracy</li><li>Customizable human review workflows</li><li>Integrates with SageMaker and other AWS services</li></ul>
Amazon Bedrock	Foundation model service	<ul style="list-style-type: none"><li>Access to pre-trained foundation models</li><li>Customization and fine-tuning capabilities</li><li>Secure and scalable deployment</li></ul>
Amazon CodeGuru	<ul style="list-style-type: none"><li>Automated code reviews</li><li>application performance recommendations</li></ul>	<ul style="list-style-type: none"><li>Identifies code defects and vulnerabilities</li><li>Provides performance optimization suggestions</li><li>Supports Java and Python</li></ul>
Amazon Comprehend	Natural Language Processing (NLP)	<ul style="list-style-type: none"><li>Entity recognition</li><li>Sentiment analysis</li><li>Topic modeling</li><li>Language detection</li></ul>
Amazon Comprehend Medical	NLP for healthcare and life sciences	<ul style="list-style-type: none"><li>Medical entity extraction</li><li>Protected PHI detection</li></ul>
Amazon DevOps Guru	ML-powered cloud operations	<ul style="list-style-type: none"><li>Anomaly detection in operational data</li><li>Root cause analysis</li><li>Proactive issue resolution recommendations</li></ul>