# Serverless in AI Insights

Oct. 2022
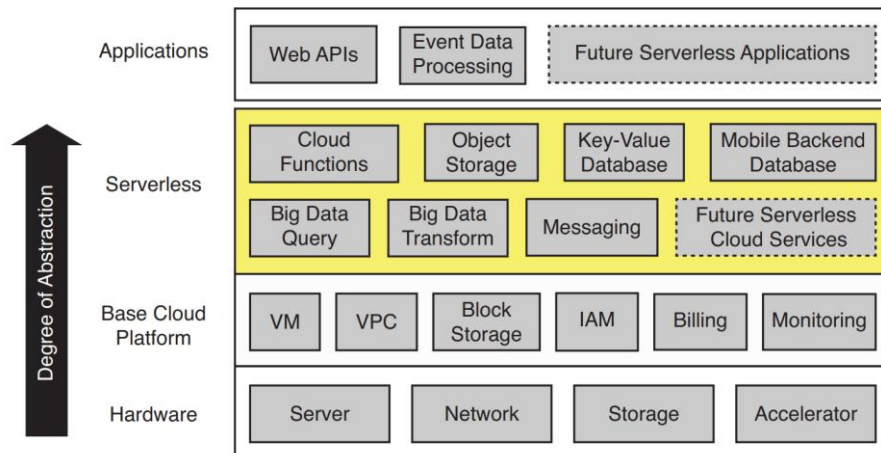
# Outline

- Definition: serverless, serverless service, FaaS

- Serverless ML architecture

- Google Colab Notebook and AWS SageMaker severless

- AWS serverless service

# Definition

- Serverless: no server management, further simplified cloud computing

- Serverless service:  components in the serverless computing paradigm

- FaaS (function-as-a-service): no services management, building block for serverless applications
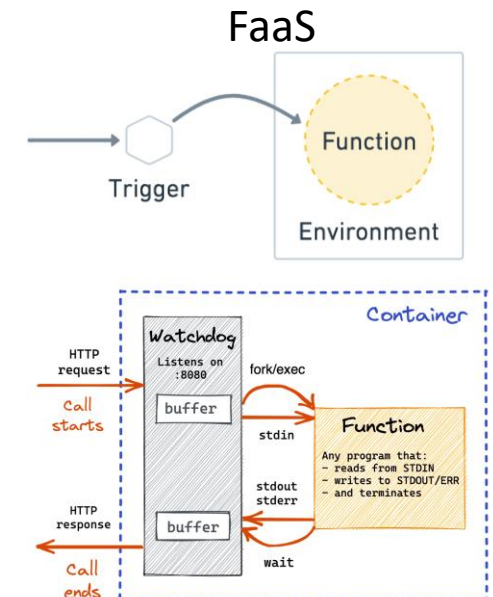
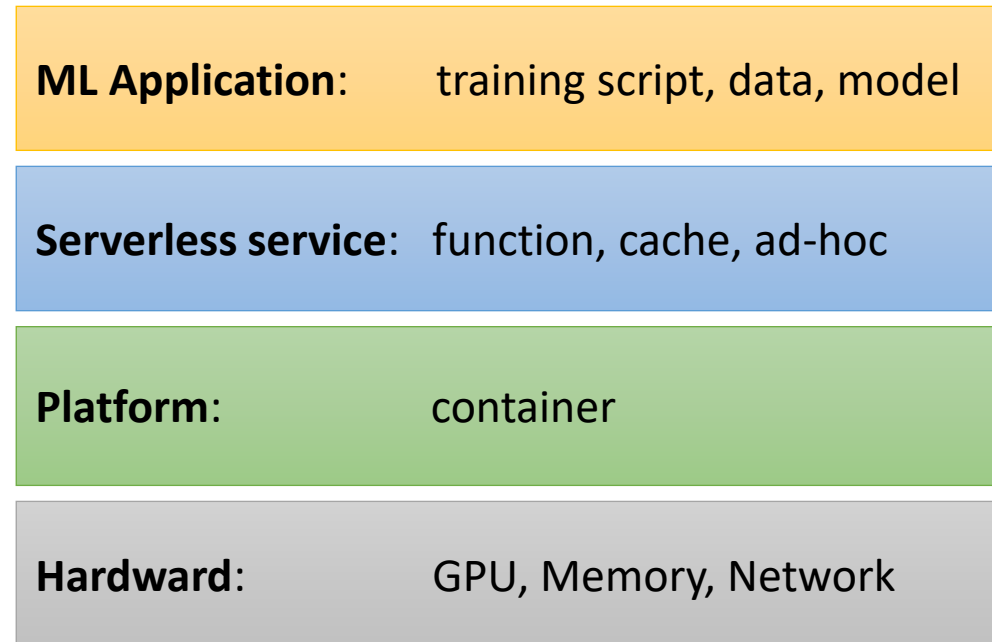Berkeley View on Serverless Computing



Serverless Service

| Service | Programming Interface | Cost Model |
|---------|----------------------|------------|
| Cloud Functions | Arbitrary code | Function execution time |
| BigQuery/Athena | SQL-like query | The amount of data scanned by the query |
| DynamoDB | puts() and gets() | Per put() or get() request + storage |
| SQS | enqueue/dequeue events | per-API call |

Table 3: Examples of serverless computing services and their corresponding programming interfaces and cost models. Note that for the serverless compute offerings described here: BigQuery, Athena, and cloud functions, the user pays separately for storage (e.g., in Google Cloud Storage, AWS S3, or Azure Blob Storage).
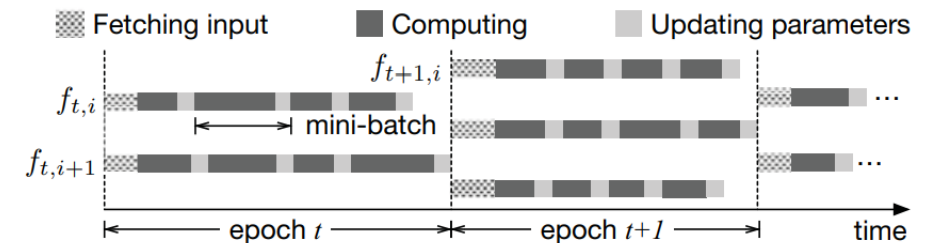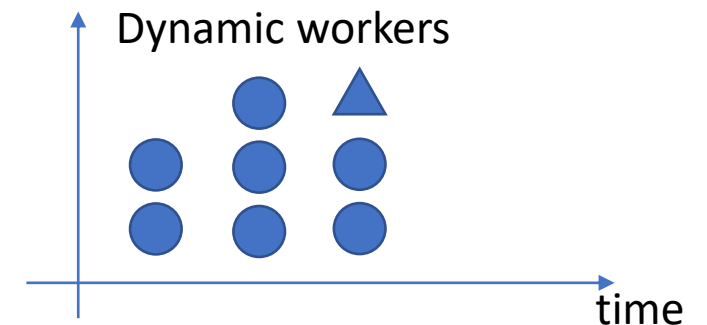
FaaS

# Serverless ML Architecture

- Training
  - Multi-dimensional split on training tasks (e.g, epoch, data, preprocessing)
  - Dynamically allocate workers(services/resources) and coordinate states between workers
  - High performance function and cache service, and ad-hoc data processing service

- Serving
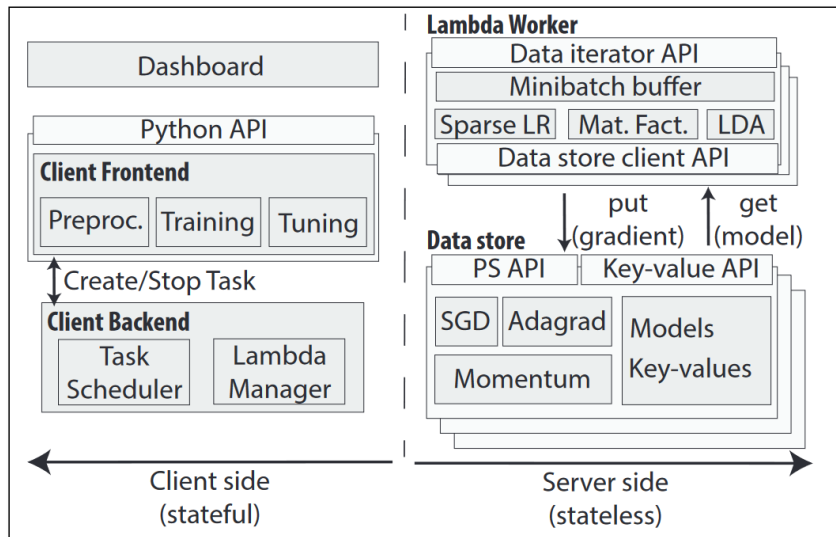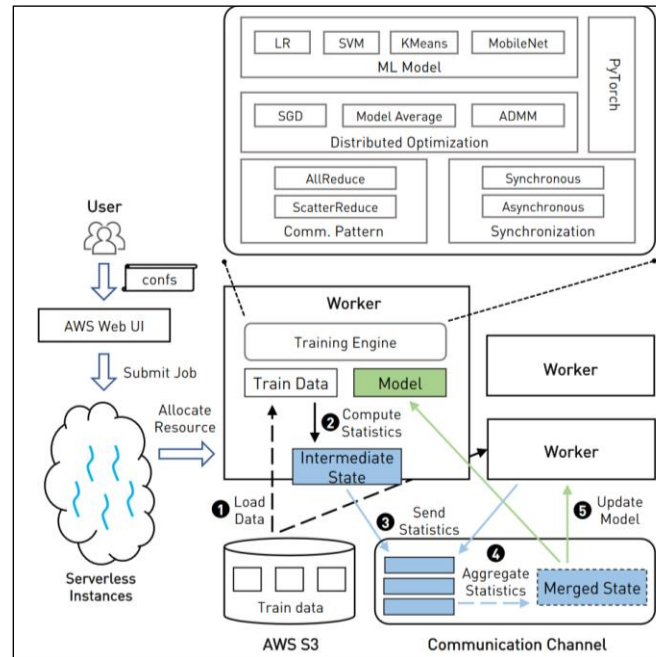  - Smart warm-up and batch, meet SLO with minimal cost

**ML Application**: training script, data, model

**Serverless service**: function, cache, ad-hoc

**Platform**: container

**Hardware**: GPU, Memory, Network

Elastic and Intelligent

Dynamic workers

time

Fetching input     Computing     Updating parameters

$f_{t,i}$

$f_{t+1,i}$

mini-batch

$f_{t,i+1}$

epoch $t$          epoch $t+1$          time

# Reference: Serverless ML Framework

- CIRRUS, 2019: https://andrewmzhang.com/assets/pdf/p13-Carreira.pdf
- LambdaML, 2021, https://arxiv.org/abs/2105.07806
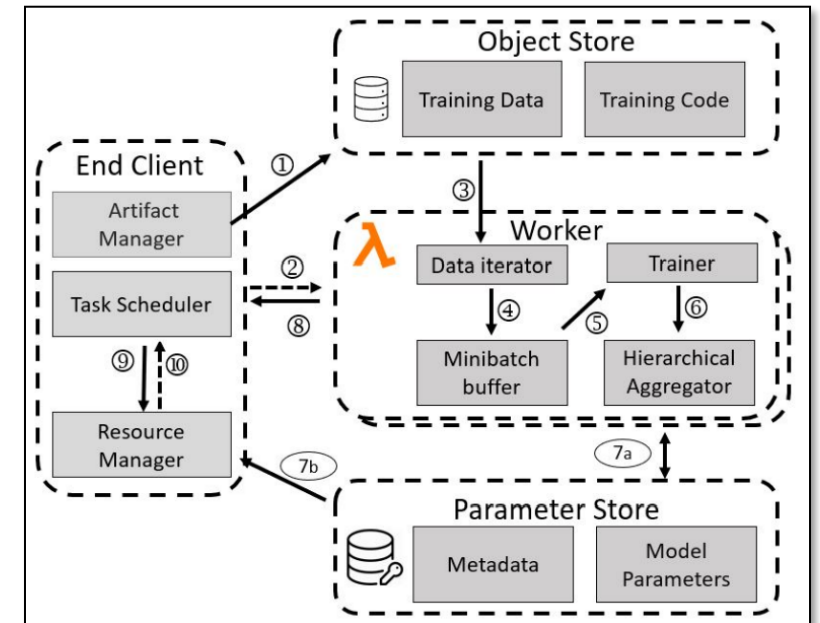- SMLT, 2022, https://arxiv.org/abs/2205.01853

CIRRUS



LambdaML



SMLT

# Design Notes

- Fast storage service is critical to exchange gradients and parameters
- Data parallelism is the first step, add advanced services to parse and coordinate model for model parallelism
- Leverage vGPU features, dynamically assign GPU memory in a worker as training goes
- Dynamically schedule the number of workers without training interruption
- Leverage open source fass project (e.g., fission), build ML oriented container (gpu, rdma), avoid general FaaS limits.

# Serverless notebook training, Google Colab

## What is Colab?

Colab, or "Colaboratory", allows you to write and execute Python in your browser, with

- Zero configuration required
- Access to GPUs free of charge
- Easy sharing

Serverless: **Colab allows you to instantly spring up a Jupyter Notebook completely serverless in the browser**. That means you don't have to worry about provisioning hardware, your Python version and path or if you're on Windows, MacOS or even a phone!

# Serverless notebook training, gradient notebook

# Disaggregated data, Resource Utils, Workflow

# AWS SageMaker serverless Inference (FaaS)

- SageMaker serverless inference



Serverless still require config
- memory selection
- max concurrency

Real-time requires config
- compute/GPU selection
- Autoscaling policy

- 4 modes of ML model deployment
  - <u>Serverless inference</u> (sparse request, intermittent, infrequent)
  - <u>Real-time inference</u> (constant request, mili-sec response time)
  - Asynchronous
  - Batch transform
- Use case
  - Fraud detection (serverless inference) → Business grow, consistent traffic (real-time inference)

https://aws.amazon.com/blogs/machine-learning/deploying-ml-models-using-sagemaker-serverless-inference-preview/

# One step further, ML as a service (serverless inference)

- Large language model, e.g., Nvidia NeMo API, Cohere.ai

## What is NeMo LLM?

NVIDIA NeMo LLM is a service that provides a fast path to customizing and using large language models trained on several frameworks. Developers can deploy enterprise AI applications using NeMo LLM on private and public clouds.

They can also experience Megatron 530B—one of the largest language models—through the cloud API or experiment via the LLM service.

### Integrate large language models into your builds

We've made an API that can be used in different libraries that fit every stack. No matter your level of developer experience, Cohere makes it easy to build machine learning into your application with our Python, Node, and Go SDKs.

Explore Docs

| Python | Node | Go | cURL | CLI |
| --- | --- | --- | --- | --- |

```
0   import cohere
1   from cohere.classify import Example
2   co = cohere.Client('{apiKey}')
3   response = co.classify(
4     model='large',
5     inputs=["this movie was great", "this movie was bad"],
6     examples=[Example("love this movie", "Positive Review"),
7   print('The confidence levels of the labels are: {}'.format(
```

Reference:
https://www.nvidia.com/en-us/gpu-cloud/nemo-llm-service/
https://cohere.ai/
https://sambanova.ai/products/dataflow-as-a-service/

# AWS serverless service (applications)



**Amazon EMR Serverless**
Run big data applications using open-source frameworks without managing clusters and servers

**Step 1: Create your application**
Choose the open-source framework and version you want to use.

**Step 2: Submit jobs**
Submit jobs to your application through APIs or EMR Studio. You can also submit jobs using workflow orchestration services like Apache Airflow or Amazon Managed Workflows for Apache Airflow.

**Step 3: Debug jobs**
Use familiar open-source tools such as Spark UI and Tez UI to monitor and debug jobs.

**Data Pipelines on EC2**

**Data Pipelines on EMR Serverless**

## Amazon Redshift Serverless

Get insights from data in seconds without having to manage data warehouse infrastructure

Amazon Redshift Serverless makes it easier to run and scale analytics without having to manage your data warehouse infrastructure. Developers, data scientists, and analysts can work across databases, data warehouses, and data lakes to build reporting and dashboarding applications, perform real-time analytics, share and collaborate on data, and build and train machine learning (ML) models. Go from large amounts of data to insights in seconds. Amazon Redshift Serverless automatically provisions and intelligently scales data warehouse capacity to deliver fast performance for even the most demanding and unpredictable workloads, and you pay only for what you use. Just load data and start querying right away in Amazon Redshift Query Editor or in your favorite business intelligence (BI) tool and continue to enjoy the best price performance and familiar SQL features in an easy-to-use, zero administration environment.

Get started with Amazon Redshift Serverless (1:28)

Reference: https://aws.amazon.com/blogs/big-data/announcing-amazon-emr-serverless-preview-run-big-data-applications-without-managing-servers/

# Serverless services on AWS

Modern applications are built serverless-first, a strategy that prioritizes the adoption of serverless services, so you can increase agility throughout your application stack. We've developed serverless services for all three layers of your stack: compute, integration, and data stores. Consider getting started with these services:

## Compute

### AWS Lambda

AWS Lambda is an event-driven, pay-as-you-go compute service that lets you run code without provisioning or managing servers.

### AWS Fargate

AWS Fargate is a serverless compute engine that works with Amazon Elastic Container Service (ECS) and Amazon Elastic Kubernetes Service (EKS).

### Use cases

**Web applications**    Data processing    Batch processing    Event ingestion

## Application integration

### Amazon EventBridge

Amazon EventBridge is a serverless event bus that lets you build event-driven applications at scale across AWS and existing systems.

### AWS Step Functions

AWS Step Functions is a visual workflow orchestrator that makes it easy to sequence multiple AWS services into business-critical applications.

### Amazon SQS

Amazon Simple Queue Service (SQS) is a message queuing service enabling you to decouple and scale microservices, distributed systems, and serverless applications.

### Amazon SNS

Amazon Simple Notification Service (SNS) is a fully managed messaging service for both application-to-application (A2A) and application-to-person (A2P) communication.

### Amazon API Gateway

Amazon API Gateway is a fully managed service that makes it easy to create and publish APIs at any scale.

### AWS AppSync

AWS AppSync is a fully managed service that accelerates application development with scalable GraphQL APIs.

## Data store

### Amazon S3

Amazon Simple Storage Service (Amazon S3) is an object storage service designed to store and protect any amount of data.

### Amazon DynamoDB

Amazon DynamoDB is a key-value and document database service, delivering single-digit millisecond performance at any scale.

### Amazon RDS Proxy

Amazon RDS Proxy is a managed database proxy for Amazon Relational Database Service (RDS) that makes applications more scalable and secure.

### Amazon Aurora Serverless

Amazon Aurora Serverless is a MySQL and PostgreSQL-compatible relational database that automatically scales capacity based on your application's needs.

Reference: https://aws.amazon.com/serverless/

# AWS serverless framework (developer tools)

- **AWS Serverless Application Model (AWS SAM)** is an open-source framework for building serverless applications. It provides shorthand syntax to express functions, APIs, databases, and event source mappings.

- **AWS Cloud Development Kit (AWS CDK)** is an open source software development framework to define your cloud application resources using familiar programming languages.

- **Serverless Framework** - The Serverless Framework consists of an open source CLI and a hosted dashboard. Together, they provide you with full serverless application lifecycle management.

- **Chalice** is a framework for writing serverless apps in Python. It allows you to quickly create and deploy applications that use AWS Lambda.

- **Arc.codes** provides everything you need to build massively scalable serverless apps with low code, clear and terse config, and zero ceremony.

- **Claudia.js** makes it easy to deploy Node.js projects to AWS Lambda and API Gateway.

Reference: https://aws.amazon.com/serverless/getting-started/

# Recent serverless paper highlights (SOCC'21)

1. Llama: A Heterogeneous & Serverless Framework for Auto-Tuning Video Analytics Pipelines (Standford)

2. Faa$T: A Transparent Auto-Scaling Cache for Serverless Applications (Microsoft)

3. Atoll: A Scalable Low-Latency Serverless Platform (UW and Google)

4. Kraken : Adaptive Container Provisioning for Deploying Dynamic DAG applications in Serverless platforms (PSU)

5. Speedo: Fast dispatch and orchestration of serverless workflows (IIT)

# Feedback

1. Concrete feature, well-defined Interface, direct reuse

2. Serverless essential service (data-intensive application)

3. Fine-grained resource scheduling (Google reference, plaque)

4. User view (security, cost saving, bill model)

5. Platform view (serverless vs PaaS, resource utilization, pro/cons, profits)