

Infrastructures et Stockage et de Traitement de Données - IST

Leonard Cseres | May 27, 2025

Serverless

- Zero server ops (auto-scaling, no provisioning)
- No compute costs when idle
- Stateless
- Asynchronous, concurrent, easy to parallelize

Structure

- Initialization code: code outside the handler method
- Local storage: /tmp
- event: information about the event
- context: func. name, req. id, deployment params, remaining time, env

Deployment Parameters Memory, CPU, Timeout, Concurrency, Environment variables

Pricing

- Resources consumed (memory × CPU)
- Number of invocations
- Ingress/egress traffic

Lifecycle



- Outside of these phases the platform freezes the execution environment
- After invocation, the platform will keep the execution environment for some time for optimizing
- The function must be stateless, but the reuse of the exec. env. can be used for caching
 - Objects outside the handler function
 - /tmp
 - Background procs. or callbacks initiated by the function and did not complete when it ended resume if the platform reuses the execution environment

Permissions

(1) Policies allowing other AWS services to invoke the function

```
Version: "2012-10-17"
Statement:
- Effect: "Allow"
  Principal:
    Service: "apigateway.amazonaws.com"
  Action: "lambda:InvokeFunction"
  Resource: "*"

```

(2) Policies allowing the function to access other services

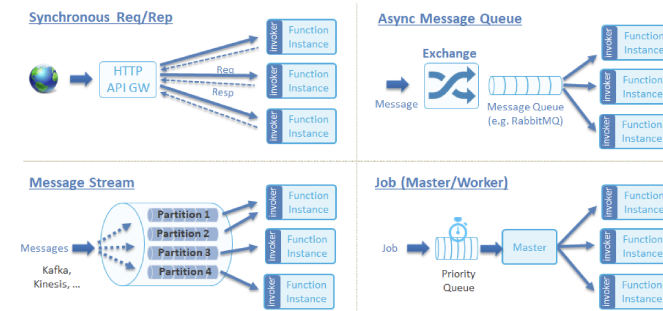
```
Version: "2012-10-17"
Statement:
- Effect: "Allow"
  Action: ["s3:GetObject", "s3:ListBucket"]
  Resource:
    ["arn:aws:s3:::your-bucket-name/*",
     "arn:aws:s3:::your-bucket-name"]

```

S3 Object Lambda

- Use cases: redaction, conversion, augmentation, compression, resizing, watermarking, custom row level authorization

S3 Access Point Alias for a bucket with custom permissions



S3 Object Lambda Access Point Enhanced S3 access point

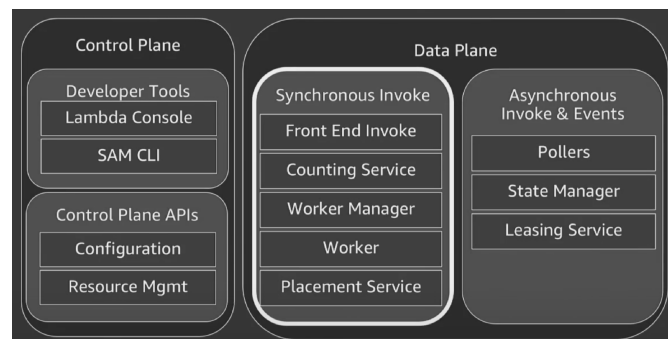
Configured with:

- Standard S3 access point
- AWS Lambda transform function
- (optional) IAM policy to restrict access to this access point

Cold Start Problem Mitigations

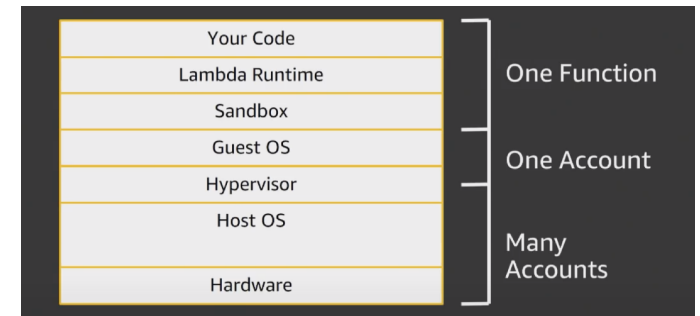
- Not only in FaaS, also in auto-scaling
- Reduce dependencies and optimize function initialization
- Fixed allocation of a set of VMs to run the function

FaaS platform



- Control plane: developer-facing API
- Data plane: underlying infrastructure

1. Call hits the LB
2. Frontend Invoke (entrypoint, retrieves metadata)
3. Counting Service (count invocations and concurrent functions)
4. Worker Manager (manages the instances, sandboxes, lifecycle)
5. Worker (compute instance, running the sandboxes)
6. Placement Service (optimizes the intelligent placement of the function, minimize the cold start, ensure optimal utilization, monitors health of workers)



Functions Security Functions are isolated from each other

- **Environment isolation:** customize the environment without affecting other
- **Security isolation:** data is not shared
- **Performance isolation:** “noisy neighbors” should not infer with other

Sandboxes are built from the same tech. as Docker

- One Sandbox executes several function invocations in serial
- A Sandbox is never reused across several functions

Hundreds or thousands of Guest OS's may run on a physical host

- Guest OS's are shared within an account
- A Guest OS is never reused across accounts

Data

- **OLAP:** Online analytical processing
- **OLTP:** Online transaction processing

Data Warehouse (1990s) Central repository of key data ingested from OLTP

- Data is well-integrated, highly structured, highly curated, and highly trusted

Hadoop & Big Data (2006) Data has high velocity, high volume, high variety

Data Lake (2011) “Store everything just in case”. Raw data, without optimizations. Data is unstructured, semi-structured, and structured. Usage of S3.

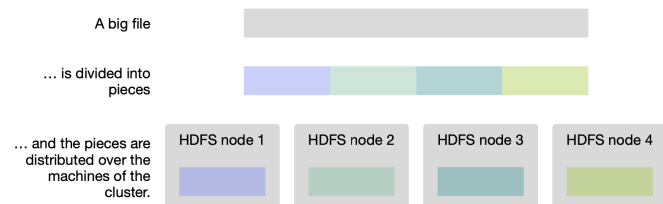
Data Lakehouse (2020) Data lake with a data warehouse. Optimized binary format (Parquet). Concurrent reads/writes (Delta Lake).

Data mesh (2021) Each company division has its own independent data lake. More agility.

MapReduce (2004) Programming model for parallel processing of large data sets

HDFS Hadoop Distributed File System

- Files are distributed across multiple nodes
- Replicated across multiple nodes 3+ times
- Chunks of 64MB
- **Move the processing to the data**

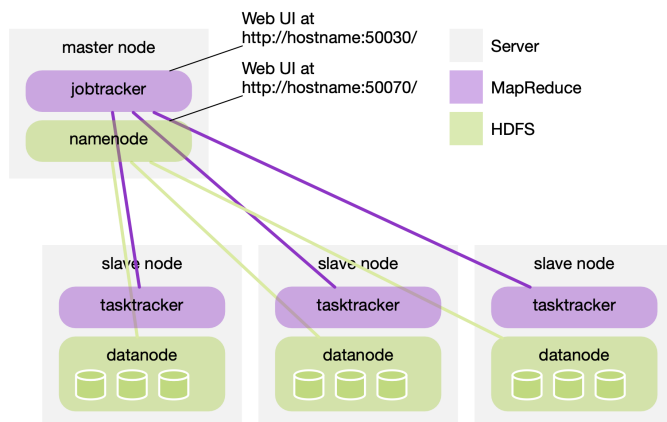


HDFS design decisions

- Files stored as chunks
- Reliability through replication
- Single master to coordinate access, keep metadata
- No data caching
- Simplify the API
- HDFS namenode**
 - manages the file system (file structure, metadata)
 - coordinates file operations (no data is moved through the namenode)
 - maintains overall health (block re-replication & rebalancing, GC)
- HDFS datanode**: manages the chunks

Hadoop 1.x cluster

- Per cluster:
 - One Namenode (NN): master node for HDFS
 - One Jobtracker (JT): master node for job submission
- Per slave machine:
 - One Tasktracker (TT): contains multiple task slots
 - One Datanode (DN): serves HDFS data blocks



Binary file formats

- Avro (Hadoop)
- Thrift (Facebook)
- Protocol Buffers (Google)
- Column-oriented binary formats: (Optimized for analytics, efficient storage and queries)
 - Parquet (Twitter & Cloudera)
 - ORC (Hadoop)

Parquet *Column-based* format - files are organized by column, rather than by row, which saves storage space and speeds up analytics queries

- Row groups: Each group contains column chunks (contiguous in file)

- Metadata at file end: chunk positions, min/max stats, dictionary filtering
- Supports sorting on one column for efficient filtering
- Encoding: Plain, Dictionary, Run-Length Encoding (RLE), Delta encoding
- Compression: gzip, Snappy, LZO
- Data types: BOOLEAN, INT32/64/96, FLOAT, DOUBLE, BYTE_ARRAY
- Logical types: STRING, ENUM, UUID, DECIMAL, DATE, TIME, TIMESTAMP, JSON, LIST, MAP
- Tools: pqr to inspect schema, content, metadata, convert to CSV/JSON

Hive Metastore & AWS Glue Data Catalog

- Hive Metastore: Metadata DB for tables, columns, types, file locations, partitions
- SerDe (Serializer/Deserializer): Reads/writes tables in various formats (Avro, ORC, Parquet, CSV, JSON, Regex, custom)
- AWS Glue Data Catalog: Hive-compatible, shared by AWS services (Athena, Glue, Redshift, etc.)
 - Tables point to folders (not files); new files = new data
 - If a file misses a column, value is null
 - Populated by DDL, UI, or Glue crawlers (auto schema inference)

Partitioning & Best Practices

- Partitioning: Split table data into subfolders (e.g., by year/month/day)
- Partition keys become virtual columns
- Queries using partition keys only scan relevant files
- Hive-style: Folders named key=value (recommended)
- Glue crawlers auto-detect partitions and add keys to schema

Data Pipelines & Transformations

- ETL: Extract, Transform, Load
- Typical steps: Ingest → Transform → Store → Catalog → Consume
- Ingestion: Batch (one-time, scheduled), Streaming (real-time)
- Transformations: Format conversion (e.g., CSV→Parquet), standardization, quality checks, partitioning, denormalization, cataloging
- Tools: AWS DMS (Database Migration Service), Kinesis (streaming), Glue (ETL), Elastic MapReduce (Hadoop/Spark)

AWS Glue

- Crawler: Discovers schema, writes to catalog
- Job: ETL (Python, Spark, or visual editor)
- Run: On demand, event, or schedule

Data Consumption

- AWS Athena: Serverless SQL queries on data lake (CSV, JSON, Parquet, ORC, Avro, logs, etc.)
 - Based on Presto/Trino
 - Requires tables in Glue Data Catalog
 - Pay per data scanned
- Amazon Redshift/Spectrum: Data warehouse/lakehouse
- QuickSight: Visualization
- Sagemaker: Machine learning

Presto / Trino

- Distributed SQL query engine (originally by Facebook)
- Interactive queries at petabyte scale
- Used by Athena, Netflix, others
- Coordinator/worker architecture

Delta Lake

- Stores table data as Parquet files + `__delta_log/` with JSON commit files (transaction log)
- Supports versioning, time travel, concurrent access

- Checkpoints (every 10 commits) for fast recovery
- Open source, supported by Spark, Presto, Trino, Hive, etc.

```
-- Create a table
CREATE TABLE table_name (
  id INT,
  name STRING
);

-- Create an external table with CSV format
CREATE EXTERNAL TABLE table_name (
  col1 INT,
  col2 STRING
)
ROW FORMAT SERDE
  'org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe'
WITH SERDEPROPERTIES ('field.delim' = ',')
STORED AS INPUTFORMAT 'org.apache.hadoop.mapred.TextInputFormat'
OUTPUTFORMAT
  'org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat'
LOCATION 's3://your-bucket/path/'
TBLPROPERTIES ('classification' = 'csv');

-- Create a Parquet table with partitions
CREATE EXTERNAL TABLE table_name (
  col1 STRING,
  col2 DOUBLE
)
PARTITIONED BY (year STRING, month STRING)
STORED AS PARQUET
LOCATION 's3://your-bucket/path/'
TBLPROPERTIES ('classification' = 'parquet');

-- Select data
SELECT * FROM table_name;

-- Filter data
SELECT * FROM table_name WHERE col1 = 'value';

-- Insert data
INSERT INTO table_name VALUES (1, 'example');

-- Create a view
CREATE VIEW view_name AS
SELECT col1, col2 FROM table_name;

-- Drop a table
DROP TABLE table_name;

-- Repair partitions (for partitioned tables)
MSCK REPAIR TABLE table_name;

-- Add a partition manually
ALTER TABLE table_name ADD PARTITION (year='2025', month='01')
LOCATION 's3://your-bucket/path/year=2025/month=01/';

-- Show partitions
SHOW PARTITIONS table_name;

-- Describe table structure
DESCRIBE table_name;
```