# Project Step 4
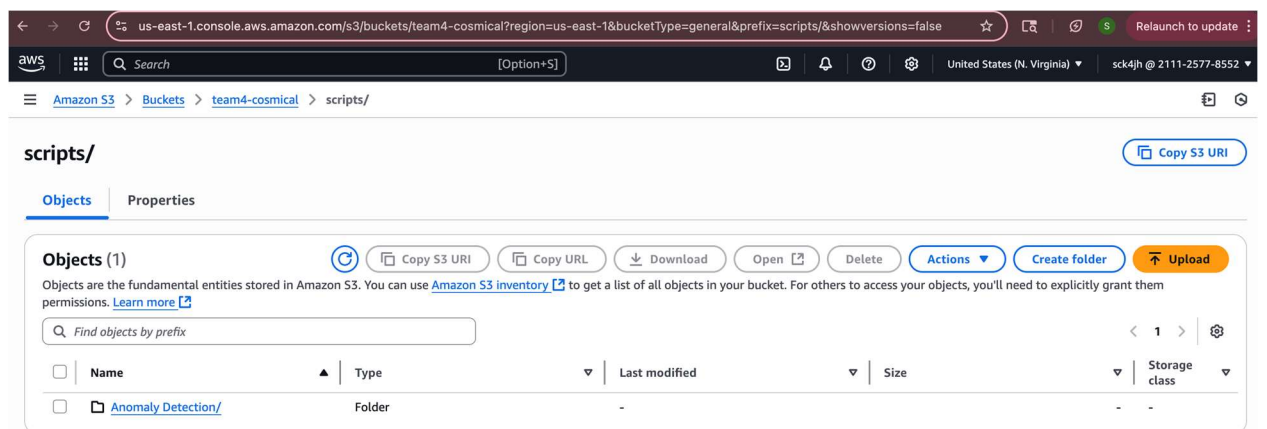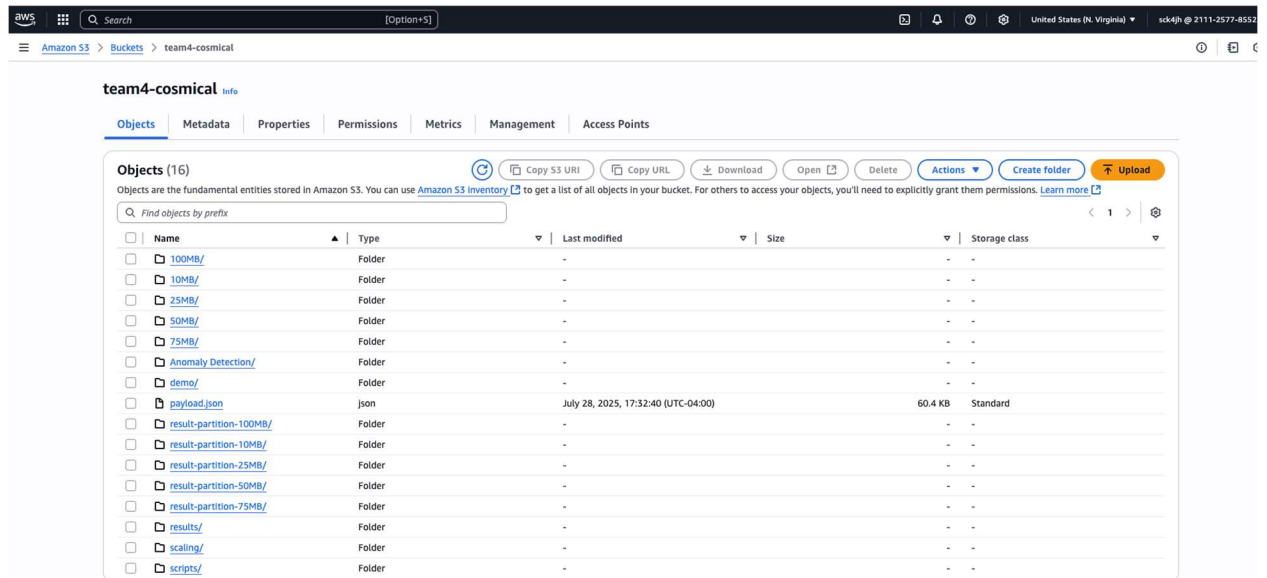
Sam Knisely, Darreion Bailey

1. Create an S3 bucket with "results" and "scripts" folders
2. Clone the AI for Astronomy repository (https://github.com/mstaylor/AI-for-Astronomy.gitLinks to an external site.)
3. Copy the Anomaly Detection folder to your S3 bucket under the scripts path

- The below screenshot shows our S3 bucket, team4-cosmical, with the results and scripts folders and the cloned AI for Astronomy repository. We also coped the AI for Astronomy folder into the scripts folder as seen in the second screenshot.

4. Configure the Step Function input payload with your bucket name, world size, and correct paths

   - The step function invoke statement can be modified as seen below with bucket names, world sizes, batch sizes, paths to different partitions of data, and paths to the scripts, data, and output results.



5. Execute the step function and monitor the CloudWatch logs at /aws/lambda/cosmic-executor

6. Examine the results in your S3 bucket's results folder

We saved results in the below format. Each folder contains each rank and the combined JSON.



7. Compare the distributed inference performance with the local execution from Step 3

| Metric | Local Execution | AWS Execution | Notes |
|---|---|---|---|
| Execution Time/Batch (s) | 4.23 | 2.46 | AWS ~42% faster |
| Throughput (Mbps) | 161.98 | 34.44 | Local ~4.7x higher |
| CPU Time (s) | ~12.7 | ~12.5 | AWS slightly better |
| Memory (MB) | ~25337 | ~3762 | AWS ~85% less memory |
| Sample/sec (normalized) | ~121 | ~211 | AWS ~75% faster |

- These three steps were completed to perform the benchmarking below.

Rubric Steps:

1. Test Baseline - Experiment with multiple batch sizes
   - The baseline group tested was the 50MB partitioned data with a world size of 10. The batch size testing yielded the below results.



Average Sample/sec vs. Batch Size (50MB data, world size 10)

2. Test various dataset access - Benchmark execution time, cost, and throughput
   - We tested benchmarked execution time, cost, and throughput using various data sizes, data partitions, and batch sizes. See the graphs is step 1 above and step 4 below.
3. Test at scale – Validate Parallel Execution

4. Performance Measurement - Create a Table to show the performance in "execution time vs. batch size" "execution time vs. input partitions" "cost vs. dataset by varying sizes"

**CAI Varying Partition Sizes**
**Throughput (bps) vs. Data Size (GB)**

# CAI Varying Partition Sizes
## Inference Time (s) vs. Data Size (GB)



Table 1: Execution Time vs. Batch Size

| Batch Size | Execution Time (s/batch) | Throughput (Mbps) | Memory Usage (GB) |
|---|---|---|---|
| 32.0 | 0.18 | 178.2 | 1.5 |
| 64.0 | 0.3 | 216.6 | 1.3 |
| 128.0 | 0.59 | 231.7 | 1.5 |
| 256.0 | 1.08 | 236.7 | 2.5 |
| 512.0 | 2.19 | 227.6 | 3.8 |

**Table 2: Execution Time vs. Input Partitions**

| Partition Size (MB) | Avg Execution Time (s/batch) | Throughput (Gbps) | Memory Usage (GB) | Num Concurrent Jobs |
|---|---|---|---|---|
| 25.0 | 1.48 | 3.9 | 1.7 | 480.0 |
| 50.0 | 1.08 | 3.6 | 2.5 | 240.0 |
| 75.0 | 1.32 | 3.3 | 3.5 | 160.0 |
| 100.0 | 1.67 | 2.5 | 4.3 | 120.0 |

**Table 3: AWS Lambda Cost Summary**

| Partition Size (MB) | Num Lambda Invocations | Total Duration (s) | Memory (GB) | Cost ($) |
|---|---|---|---|---|
| 25.0 | 491.0 | 3633.4 | 2.8 | 83.27 |
| 50.0 | 245.0 | 2646.0 | 4.0 | 43.2266 |
| 75.0 | 163.0 | 3227.4 | 5.9 | 51.7402 |
| 100.0 | 122.0 | 4074.8 | 7.0 | 58.0096 |

Cost insights:

1. Smaller partitions (25MB) have more parallel invocations but lower memory requirements

2. Larger partitions (100MB) have fewer invocations but require more memory

3. The 50MB partition offers a good balance between parallelism and resource usage

4. Cost increases roughly linearly with dataset size for all partition sizes

Summary - Provide observation and analysis of the overall results for the project on scalable and efficient CosmicAI infrastructure using Serverless Cloud Computing

The results showed that running CosmicAI on serverless infrastructure was significantly faster and more efficient than running it locally. Inference time dropped from 4.23 seconds to 2.46 seconds, throughput increased by about 75%, and memory usage dropped by over 80%. The main bottleneck in the serverless setup was reading data from S3, which was slower than reading

from local disk. The most efficient setup used 50 MB data chunks, which balanced speed, memory use, and cost well. Smaller chunks caused too much overhead, while larger chunks slowed down processing. Overall, cost scaled predictably with dataset size, making it easy to estimate expenses at scale. The results confirmed that serverless is a strong solution for scaling CosmicAI, with the only major limitation being cloud storage speed.