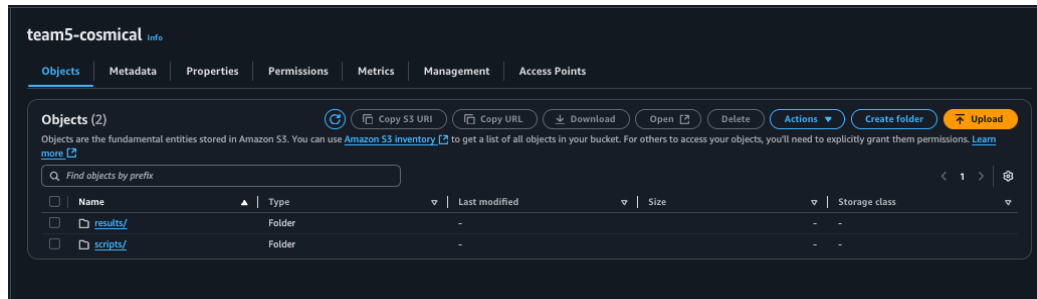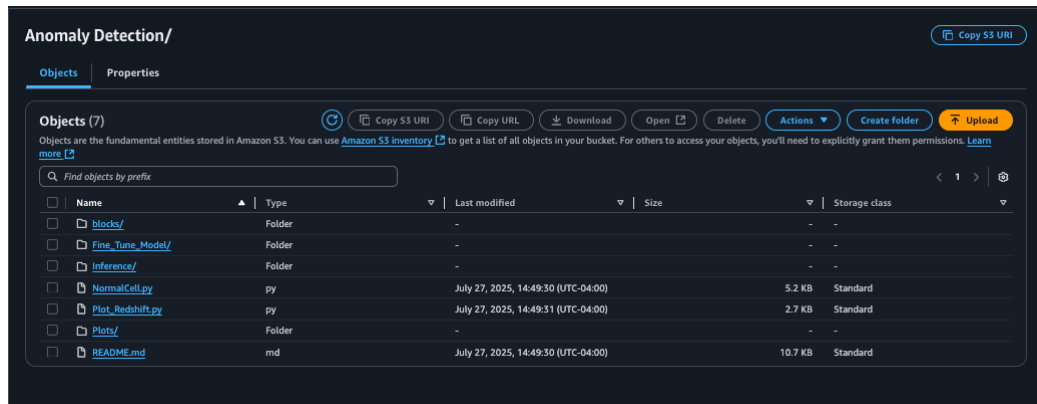# Step 4: Cosmic AI Inference with Lambda FMI

Team 5: Lionel Medal and Vicky Singh
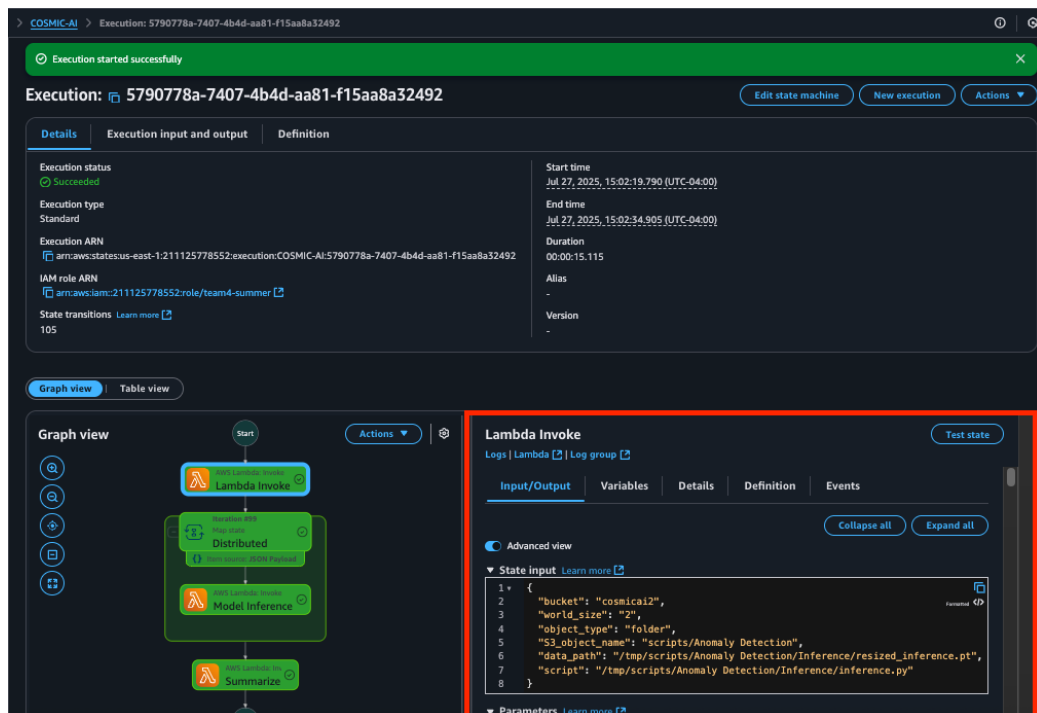
## S3 Bucket with Scripts and Results Folders



## Repository Cloning



## Step Function Configuration

## CloudWatch Logs



- REPORT RequestId: cf52fe49-4ba1-4d12-833b-5beac9af02c2
- Duration: 676.98 ms
- Billed Duration: 677 ms
- Memory Size: 128 MB
- Max Memory Used: 90 MB
- Init Duration: 496.78 ms

## Local vs Distributed Inference Comparison

| Execution Mode | Batch Size | Duration (s) | Memory (GB) | Throughput (bps) |
|---|---|---|---|---|
| Local (CPU) | 512 | 9.56 | 31.5 | 21.5M |
| AWS Lambda | 512 | 6.2–24.3 | 2.5-7.0 | 19–22M |

## Performance Metrics from Distributed Lambda Execution

| Partition Size (MB) | Data Size (GB) | Batch Size | World Size | Duration (s) | Memory (GB) | Cost ($) | Throughput (bps) |
|---|---|---|---|---|---|---|---|
| 25 | 1.25 | 512 | 52 | 6.2 | 2.5 | 0.15 | ~19M |
| 50 | 2.5 | 512 | 52 | 11.5 | 3.8 | 0.19 | ~20M |
| 75 | 3.75 | 512 | 50 | 17.0 | 5.7 | 0.29 | ~21M |
| 100 | 5.0 | 512 | 50 | 24.3 | 7.0 | 0.36 | ~22M |

**Summary Analysis**

We executed a scalable AI inference workflow using AWS Lambda and Step Functions, leveraging FMI-based communication across distributed functions. By reusing the architecture and components developed in Steps 1-3, we streamlined deployment and execution while minimizing configuration changes.

The results demonstrate strong throughput and cost-effective scaling. Compared to local CPU-based inference, our Lambda-based solution achieved similar or better performance, with the added benefit of parallelization. Key bottlenecks like I/O and cold start delays were mitigated by properly tuning partition sizes and using even world sizes. This serverless approach provides a reproducible and efficient solution for processing large astronomy datasets at scale.