# Azure OpenAI Batch API Accelerator
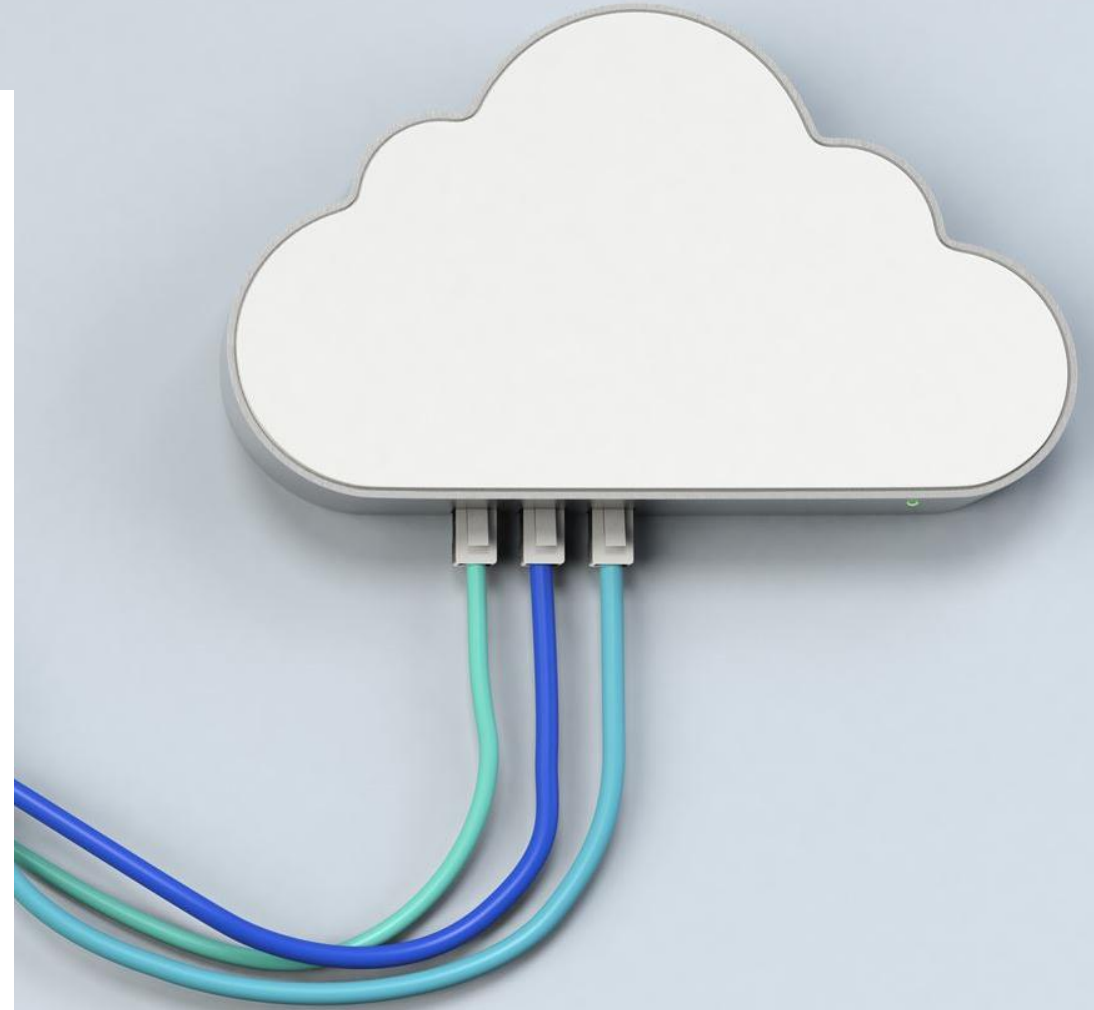
DJ Dean – Principal Cloud Solution Architect

Amit Mukherjee – Prin. Technical Specialist – GenAI

October 2nd, 2024

Microsoft

# Why We Need the AOAI Batch API

- The Batch API significantly reduces costs for non-time sensitive workloads.

- Implementing multi-threaded/multi-process functionality can be difficult and time-consuming.

- The Batch API provides significant performance benefits when processing large quantities of files.

# Customer Problem statement

- How to process tens of thousands of files to extract information

- How do track errors at the time of processing

- How to automate entire process

- How to make cost effective

AOAI Bath Accelerator

Automated Batch Job Submission and Creation

Multi-threaded Async Processing to Reduce Overall Processing Time

Automated Error Tracking

Multi-directory Hierarchy Support

Configurable Micro-batch Support
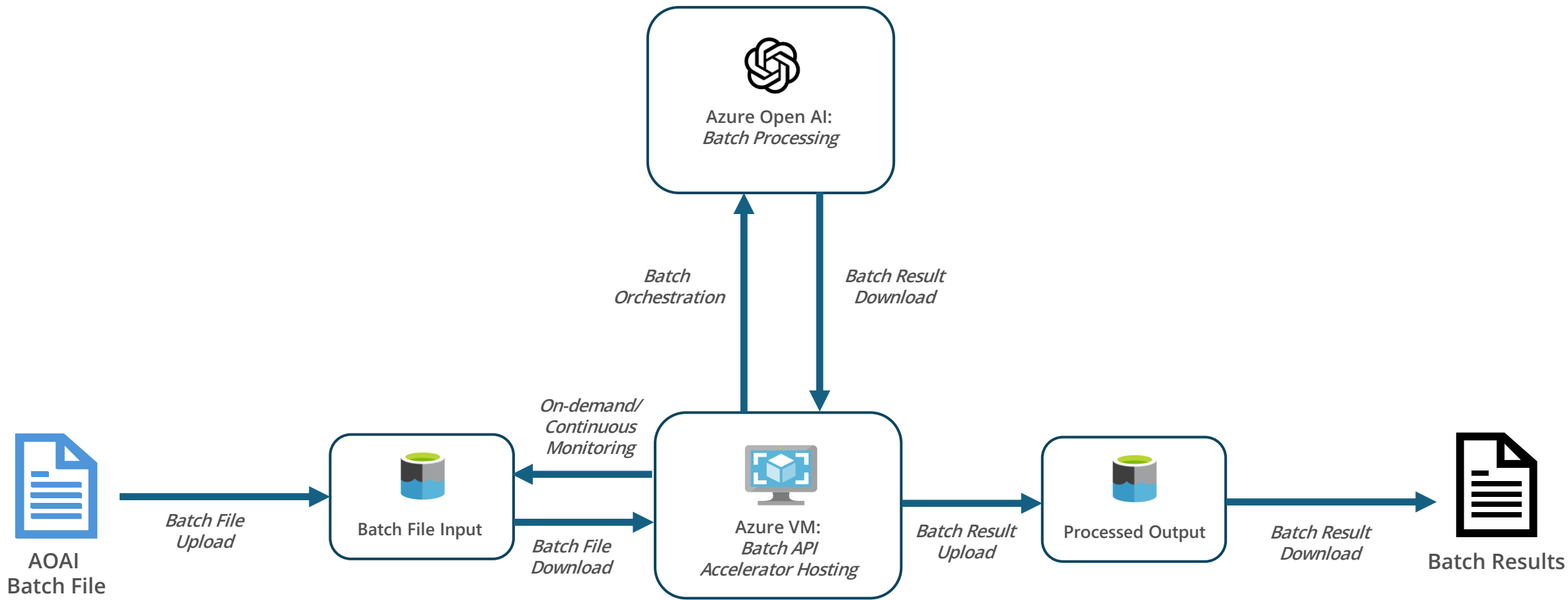
Automated Post-job Cleanup

# Voice of Customer

ontada®

"Ontada is at the unique position of serving providers, patients and life science partners with data-driven insights. **We leverage the Azure OpenAI batch API to process tens of millions of unstructured documents efficiently, enhancing our ability to extract valuable clinical information. What would have taken months to process now takes just a week. This significantly improves evidence-based medicine practice and accelerates life science product R&D**. Partnering with Microsoft, we are advancing AI-driven oncology research, aiming for breakthroughs in personalized cancer care and drug development."
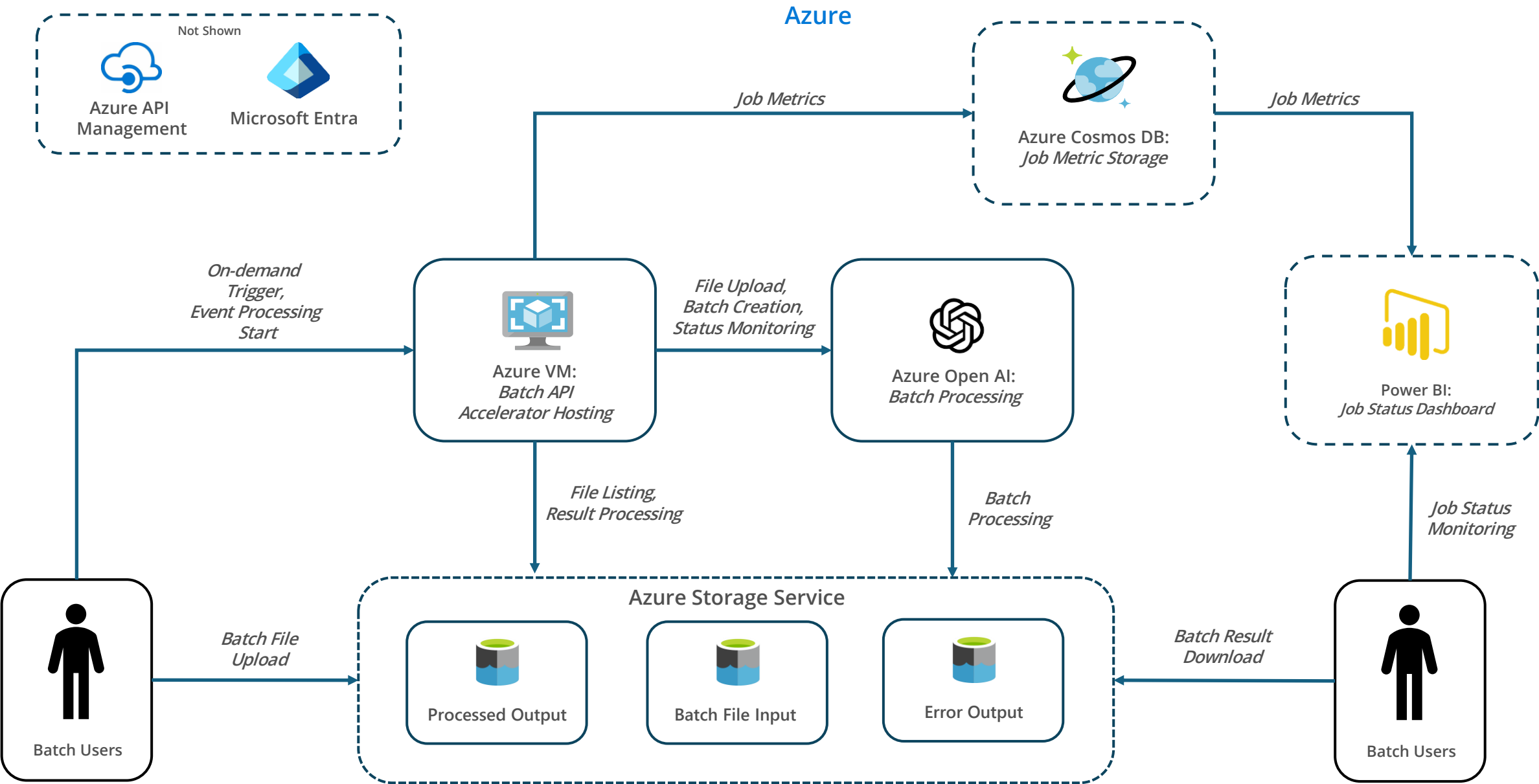
Sagran Moodley, Chief Innovation and Technology Officer, Ontada
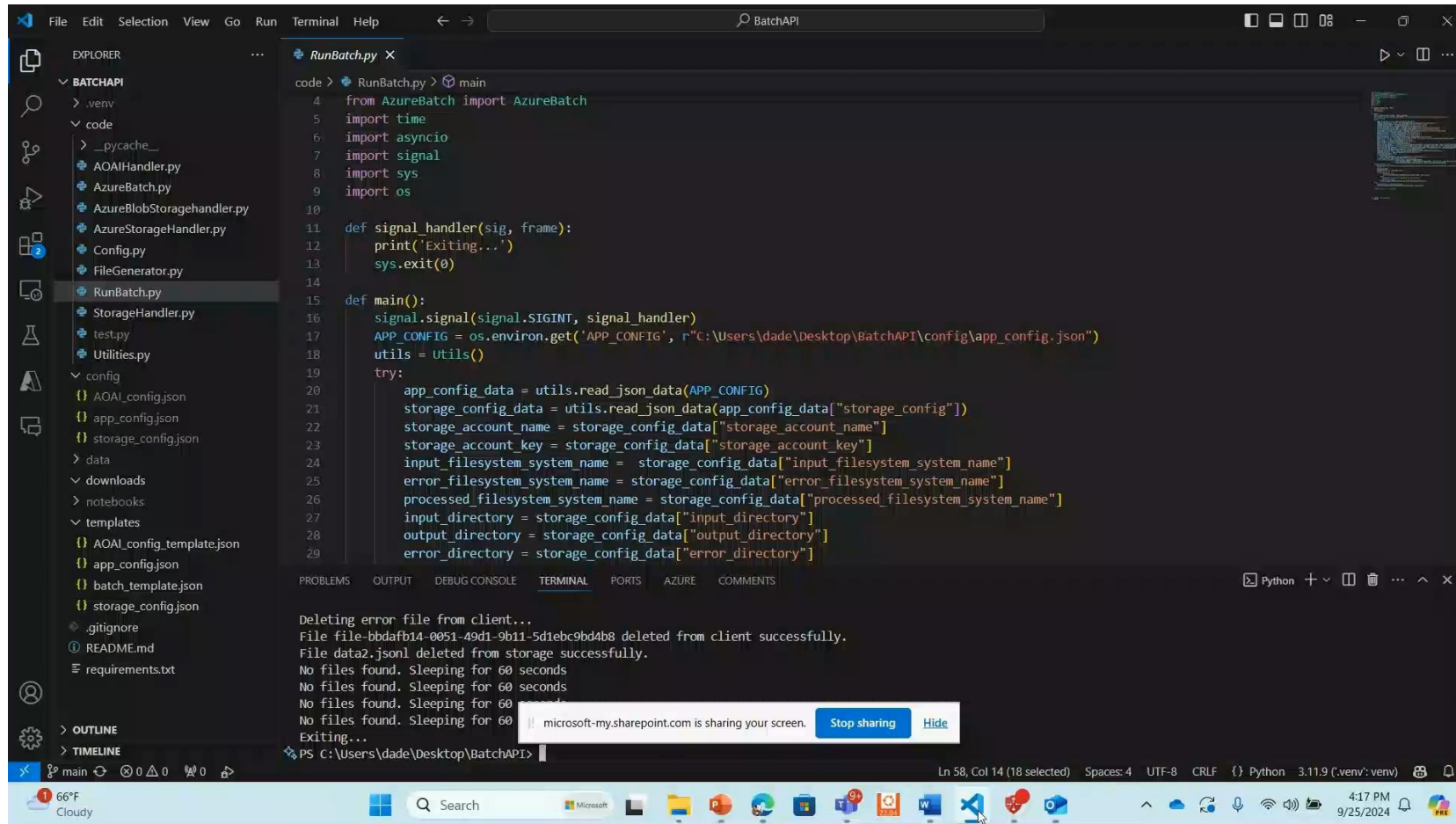
# AOAI Batch Accelerator Overview

## Azure

# High-level Architecture (Detailed) | Data Flow and Consumption

## Azure

Not Shown

Azure API Management

Microsoft Entra

Azure Cosmos DB:
*Job Metric Storage*

Job Metrics

Job Metrics

On-demand Trigger, Event Processing Start

Azure VM:
*Batch API Accelerator Hosting*

File Upload, Batch Creation, Status Monitoring

Azure Open AI:
*Batch Processing*

Power BI:
*Job Status Dashboard*

File Listing, Result Processing

Batch Processing

Job Status Monitoring

## Azure Storage Service

Processed Output

Batch File Input

Error Output

Batch File Upload

Batch Result Download

Batch Users

Batch Users

# Demo