



# T-SQL → Spark SQL & PySpark Code Converter Tool

Individual Intern Project, Summer 2024

## **Contributors:**

Ray Wang with grateful  
guidance from  
Numair Ul-Ghani, Brian Lui

# Context

For  
Microsoft  
Internally  
and  
Our  
Customers

## Problem



How can we rapidly enable the conversion of T-SQL into PySpark and Spark SQL to **save time** and **ensure correctness**?



## Opportunity



Can we better equip team members with custom tools using LLMs like gpt-4o to automate this process?

*“Helping customers quickly convert their T-SQL is a capability we’d love to offer!” – TSP/CSA Data & AI*

## Key Outcomes

Less friction in  
migrations

Lower Spark SQL &  
PySpark barrier

Increased uptake of  
Fabric/Databricks

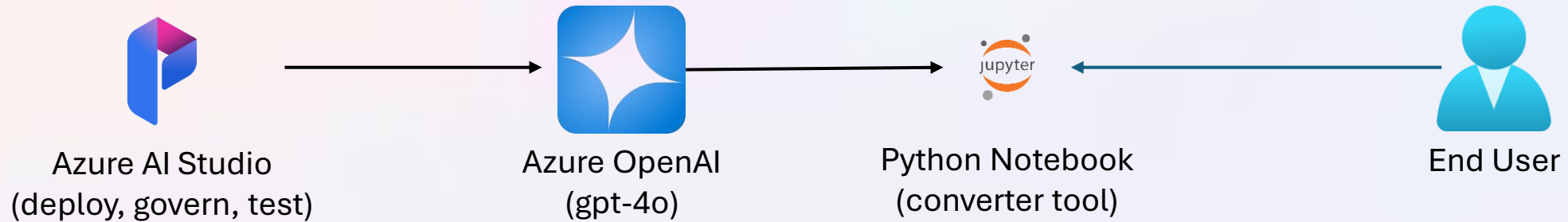
# Conversion Steps and Customized Components

*How is this different from plain old gpt-4o?*

## Code blocks in a shareable Jupyter Notebook:

1. Install dependencies and import environment variables
2. Ingest T-SQL code
  - A. Supports single string paste and bulk .sql file. formats
3. Validate T-SQL for correctness
4. List dependencies on tables, views, and other database objects
5. **Define custom conversion logic** (tuning for Databricks / Fabric)
  - A. Few-shot prompting and custom context reminders of syntactical differences
  - B. Detailed, guard railed system message with temperature set to 0.1
6. Translation to Spark SQL
  - A. Works best for basic T-SQL / dialect conversion for DDLs, etc.
7. Translation to PySpark
  - A. Works best for stored procedures / complex logic)

# Technical Resources and Requirements



## Azure:

- Active Azure subscription and Azure OpenAI resource with deployed gpt-4o model, and sufficient tokens-per-minute provisioned to avoid rate limiting (recommended: 130K TPM)

## Local computer-side:

- Visual Studio Code (or another IDE with Jupyter Notebook support)
- Python 3.11
- Git Bash to pull the repository (optional, can also download as .zip)
  - .ipynb notebook that can execute in any Jupyter environment
  - .env file to enable easy linking with Azure resources (template included)

# Benchmarking

	Key Constraints	Integrability	Accuracy and Performance
<b>This Converter (custom gpt-4o)</b>	<b>Context length: 128K tokens</b> <i>(supports many stored procs)</i>	Supports <b>multiple statements &amp; .sql files</b>	1 <sup>st</sup> / Best (output requires no debugging most of the time)
<b>ChatGPT Plus (gpt-4o)</b>	<b>Context length: 4K tokens</b> <i>at a time (prone to truncation)</i>	Requires manual pasting in of SQL extracts	2 <sup>nd</sup> / Okay (output valid some of the time, data not secured)
<b>Databricks Assistant (gpt-4)</b>	<b>Context length: 500 tokens</b> <i>(usually stops part-way during periods of high demand)</i>	Automatically debugs and can offer conversions	2 <sup>nd</sup> / Okay (helpful for debugging, but insufficient for conversion)
<b>Copilot (gpt-4o)</b>	<b>Context length: ~2K chars</b> <i>(script usually doesn't fit)</i>	Requires manual pasting in of SQL extracts	3 <sup>rd</sup> / Poor (lazily truncates logic)

# Success Story: Major Public Transit Agency

## Context

Agency is migrating from Synapse to Databricks with CSA assistance.

## Action

Converted DDLs, DMLs, and Stored Procs to Spark SQL and PySpark with the code converter and Databricks Assistant debugging

## Result

Completed main ETL procedures in 4 weeks with just 1 CSA + intern



Q&A