



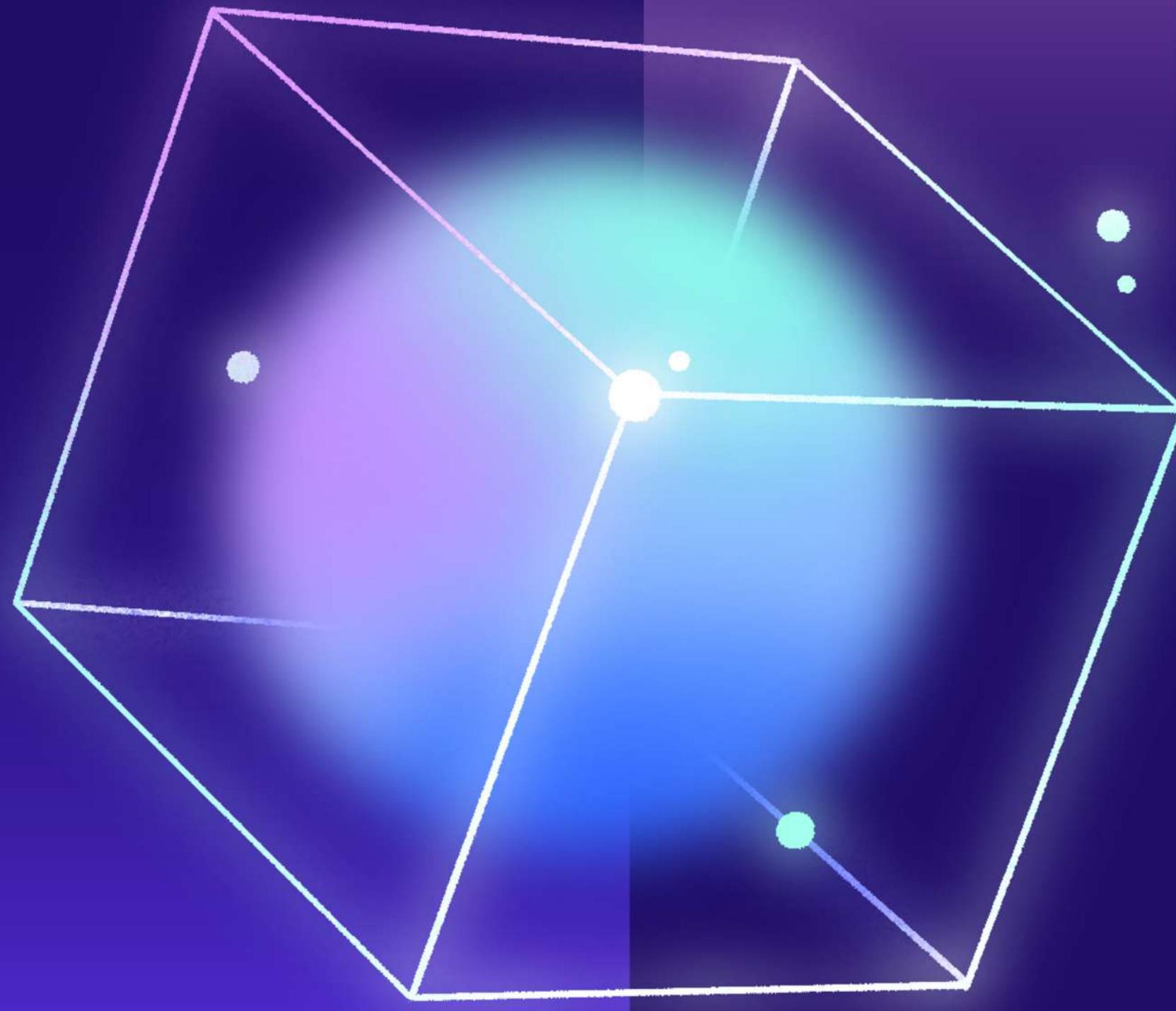
END-TO-END ETL PIPELINES

...

Pelcharskyi Artur

3 December, 2025

WHAT IS ETL?



ETL stands for **Extract, Transform, Load**, which is a process used in data management and analytics to move data from various sources into a centralized system, usually a data warehouse.

In essence, ETL ensures that data from multiple sources becomes consistent, accurate, and ready for analysis, enabling better decision-making.

•••

ETL STEPS

- **Extract:** This is the first step, where data is collected from different sources such as databases, APIs, or files. The goal is to gather raw data, even if it's in different formats.
- **Transform:** Here, the extracted data is cleaned, formatted, and processed to fit the target system. This can include filtering, aggregating, converting data types, or applying business rules.
- **Load:** Finally, the transformed data is loaded into the target system, like a database or data warehouse, where it can be used for reporting, analysis, or machine learning.



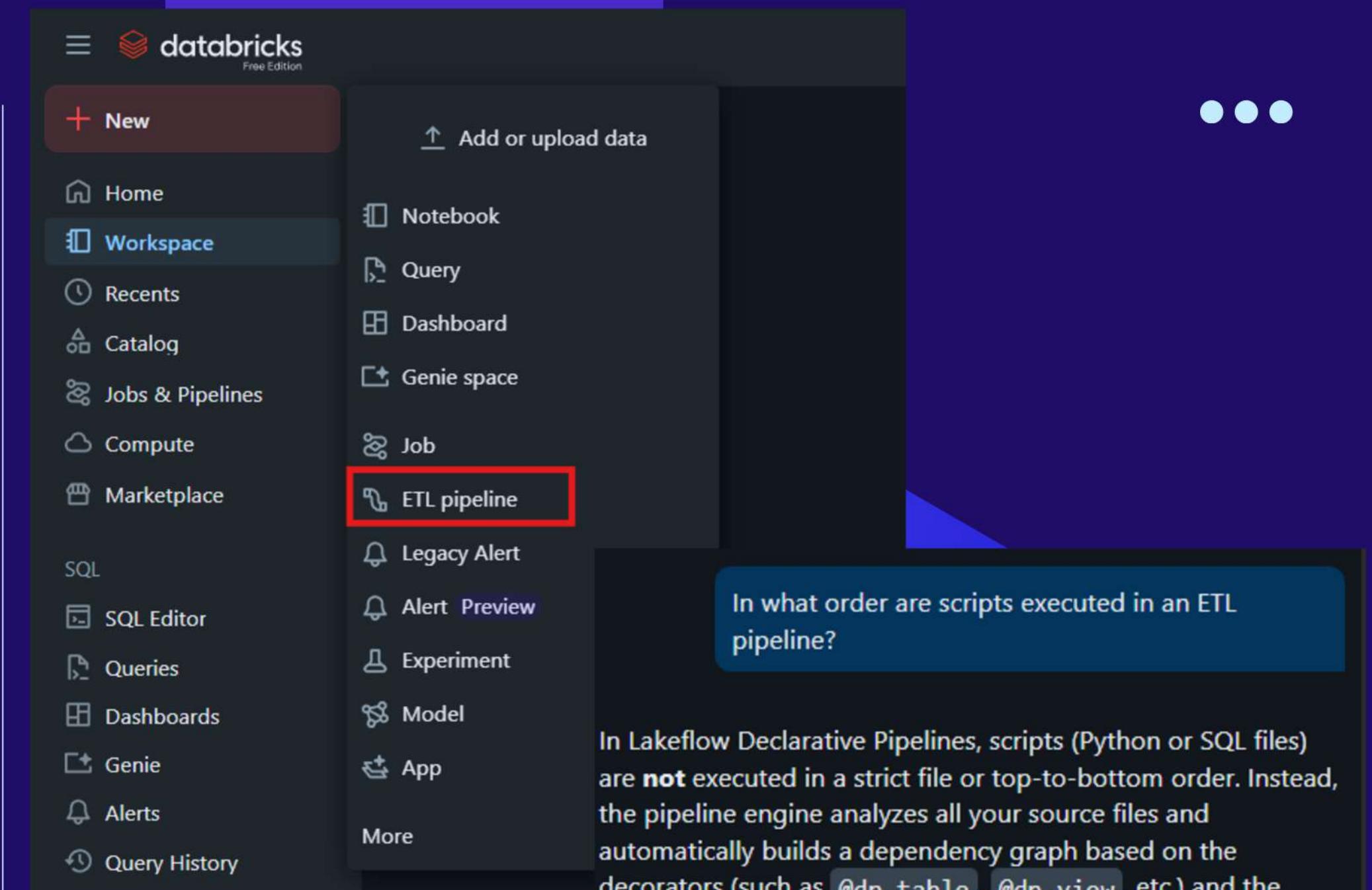
•••

ETL PIPELINES IN DATABRICKS

The ETL Pipeline is part of the Pipelines module in Databricks. This module helps structure code to be modular, scalable, and easy to maintain. Databricks also enables working efficiently with individual tables that are updated within a pipeline.

A **Pipeline** is not the same as a sequence of **Jobs**. In Jobs, we work with notebooks that can be connected in a chain, like **Notebook1** -> **Notebook2** ->

In a Pipeline, the focus is on table dependencies, which are automatically resolved. (More details in the demo.)



The screenshot shows the Databricks workspace sidebar. The 'Jobs & Pipelines' section is expanded. The 'ETL pipeline' option is highlighted with a red box. A callout bubble to the right asks, "In what order are scripts executed in an ETL pipeline?" Below it, a detailed explanation is provided: "In Lakeflow Declarative Pipelines, scripts (Python or SQL files) are **not** executed in a strict file or top-to-bottom order. Instead, the pipeline engine analyzes all your source files and automatically builds a dependency graph based on the decorators (such as `@dp.table`, `@dp.view`, etc.) and the references between tables, views, and flows. Transformations are executed in the order required to satisfy these dependencies, ensuring that upstream tables/views are computed before downstream consumers, regardless of file or function order 1. This means you can organize your code across multiple files, and Lakeflow Declarative Pipelines will always resolve and execute transformations according to the dependency graph, not the order of files or functions." At the bottom right, there are like, dislike, and share icons, and a note: '> 1 citation'.

↑ Add or upload data

+ New

Home

Workspace

Recents

Catalog

Jobs & Pipelines

Compute

Marketplace

SQL

SQL Editor

Queries

Dashboards

Genie

Alerts

Query History

ETL pipeline

Legacy Alert

Alert Preview

Experiment

Model

App

More

In what order are scripts executed in an ETL pipeline?

In Lakeflow Declarative Pipelines, scripts (Python or SQL files) are **not** executed in a strict file or top-to-bottom order. Instead, the pipeline engine analyzes all your source files and automatically builds a dependency graph based on the decorators (such as `@dp.table`, `@dp.view`, etc.) and the references between tables, views, and flows. Transformations are executed in the order required to satisfy these dependencies, ensuring that upstream tables/views are computed before downstream consumers, regardless of file or function order 1. This means you can organize your code across multiple files, and Lakeflow Declarative Pipelines will always resolve and execute transformations according to the dependency graph, not the order of files or functions.

Like Dislike Share > 1 citation

WORLD BANK ETL PIPELINE

Extract



Transform

...

DEMO

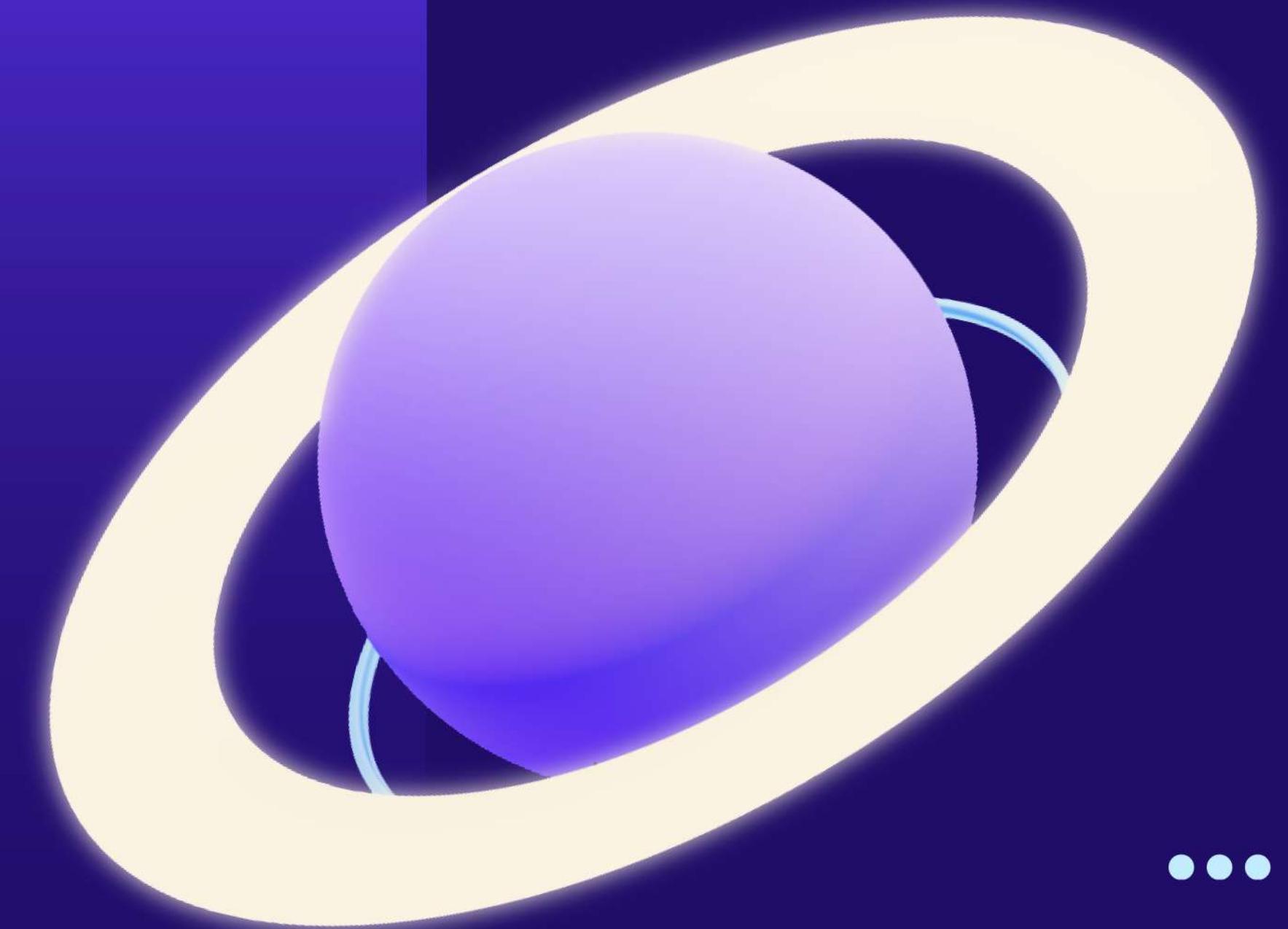
[REPO LINK](#)



Q/A

...

**THANK YOU
FOR YOUR
ATTENTION**



...