**ENVS 5726** Guest Lecture

*Building Data Pipelines to Support Mineral Exploration Activity*

Abdel Alfahham
alfahham@sas.upenn.edu

*Founding Data Engineer*
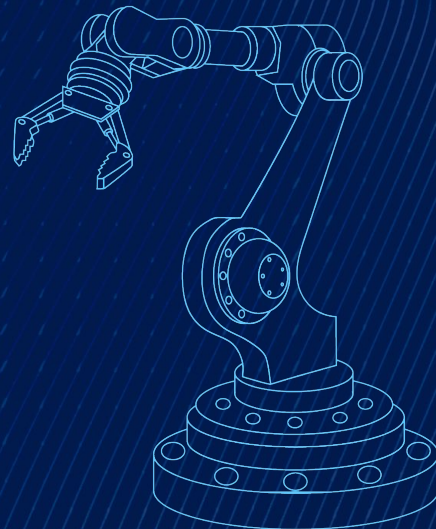
# TABLE OF CONTENTS

# 01

# Introduction

## Summary and Objective

# Career Timeline

- BS Environmental Sciences    [2009 - 2013]
- Environmental Consulting    [2013 - 2018]
- Penn MSAG    [2016 - 2019]
- Data Engineer    [2019 - Current]

- Self Taught Python (coding in general)
- Intro to SQL in Environmental Consulting
- Intro to GIS/QGIS/CAD in Environmental Consulting

**Founding Data Engineer**
ecue ai · Full-time
Jul 2024 – Present · 1 yr 6 mos
New York, New York, United States · Remote

**Software Engineer**
Grata · Full-time
Sep 2022 – Aug 2024 · 2 yrs
New York, New York, United States · Remote

**KoBold Metals**
Full-time · 2 yrs 11 mos
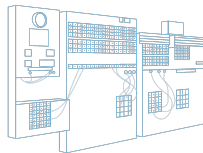San Francisco Bay Area · Remote

**Data Engineer**
Nov 2019 – Sep 2022 · 2 yrs 11 mos

# Lecture Objective

1. Provide a *detailed but practical overview* on the various components of an end-to-end *data platform*.

2. Introduce critical *data engineering, data infrastructure and software development concepts.*

3. Demonstrate the *impact of data engineering in earth science.*

# Background

- During this lecture I will be discussing and demonstrating a hypothetical pipeline that I created to share the knowledge I have learned from the past 7 years or so as a data engineer.

- The main objective is to share the fundamentals and give jumping off points for individuals that are interested in pursuing data engineering or data science roles.

## 02

# Geophysical Surveys

Ground EM Fundamentals

# Geophysics in Mineral Exploration

**Geophysical surveys allows exploration teams to identify promising areas without having to dig or drill, saving time and money in the early stages of mineral exploration.**

# Types of Geophysical Surveys

Exploration

| Drill Permit | Geophysics | Drill Target Selected | Drilling | Drill Results |

| Electromagnetics (EM) | Gravity | Induced Polarization (IP) | Magnetics |

# Electromagnetic (EM) Surveys

1. An artificial electromagnetic signal from the **transmitter (Tx)** is generated and it penetrates deep into the earth.
2. When the signal hits **conductors** (conductive layers of rock or minerals like nickel, copper, and other base metals) it changes.
3. A **receiver (Rx)** picks up these changed signals and the information is saved.



4.8 kW Transmitter

• Portable Transmitter (Tx) unit for Surface, Borehole and Underground Pulse-EM surveys



Induction Coil (dB/dt)

Induction sensors measure the rate of change of the B-field. This derivative value allows you to see targets at depth exhibiting a large range of conductivities including weakly conductive zinc mineralization.



Primary EM field

Tx

Rx

Surface

Modified primary field

Eddy currents

Secondary field

Conductor

# Electromagnetic (EM) Surveys

**What it measures:**

- Ground EM surveys detect how well electrical current flows through subsurface materials.
- Highly conductive materials (like ore deposits) conduct electricity easily, while resistive materials don't

**How it works:**

- An artificial electromagnetic signal is transmitted into the ground, and the response is measured; conductive materials create stronger secondary signals that the receiver detects

**Why it matters:**

- Conductivity patterns reveal the location, size, and type of subsurface targets; anomalies indicate potential mineral deposits, geological structures.

# Revisiting Electromagnetics (EM)

- **Strong Conductors and Anomalies:** Conductors and anomalies are prominent features on an EM survey, these are zones that conduct electricity particularly well. *Note: depending on the type of mineral the company is searching for, they may be looking for low conductivity.*
- **Shape and Depth of Anomalies:** The shape and depth of a conductive anomaly can help with modelling structures underground.
- **Context of Geological Features:** EM survey results are often used in conjunction with other geophysical data such as rock chip sampling, this allows an explorer to be further informed about the potential. For example, an anomaly along a well-known deposit belt or occurrence can serve to increase the confidence of prospectivity in a project.
- **Comparison with previous data:** If there is historical data and older survey results for the region, it can be compared with the new EM data to give the company a better understanding of its project.

# 03

# Data Pipelines

Automation, Scaling, Design

# Basic Components of Data Platform

Let's Start the Demo Now!

ENVS 5726

# Basic Components of Data Pipeline

# Basic Components of Data Pipeline
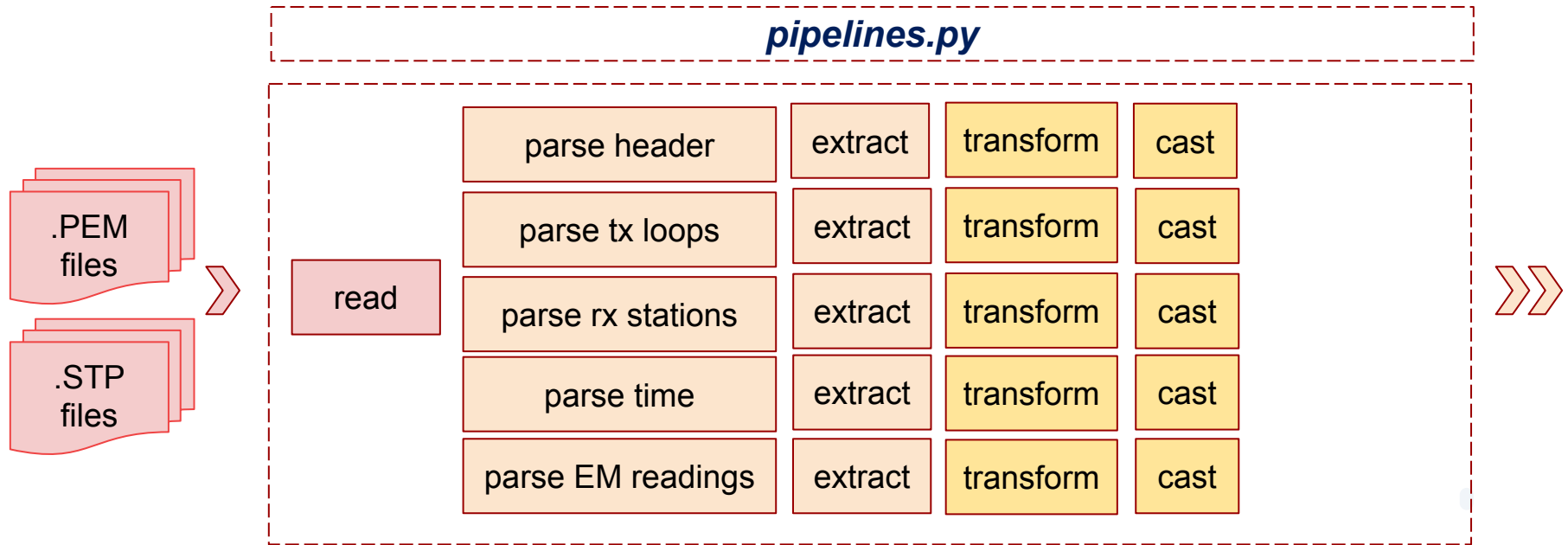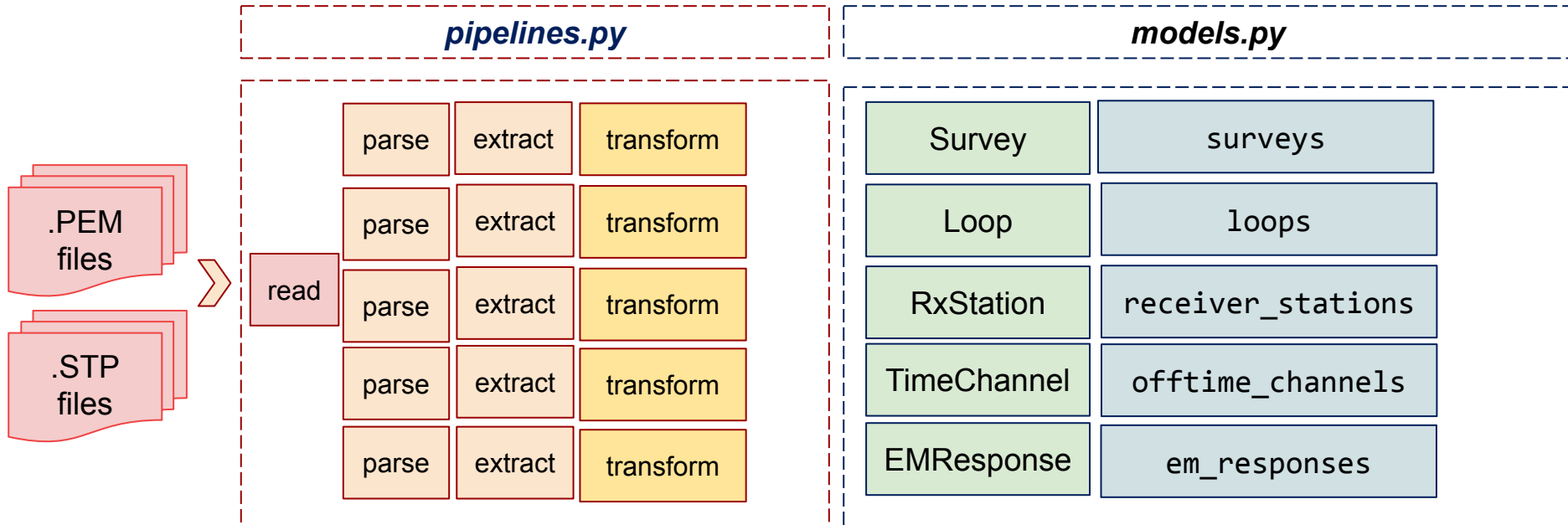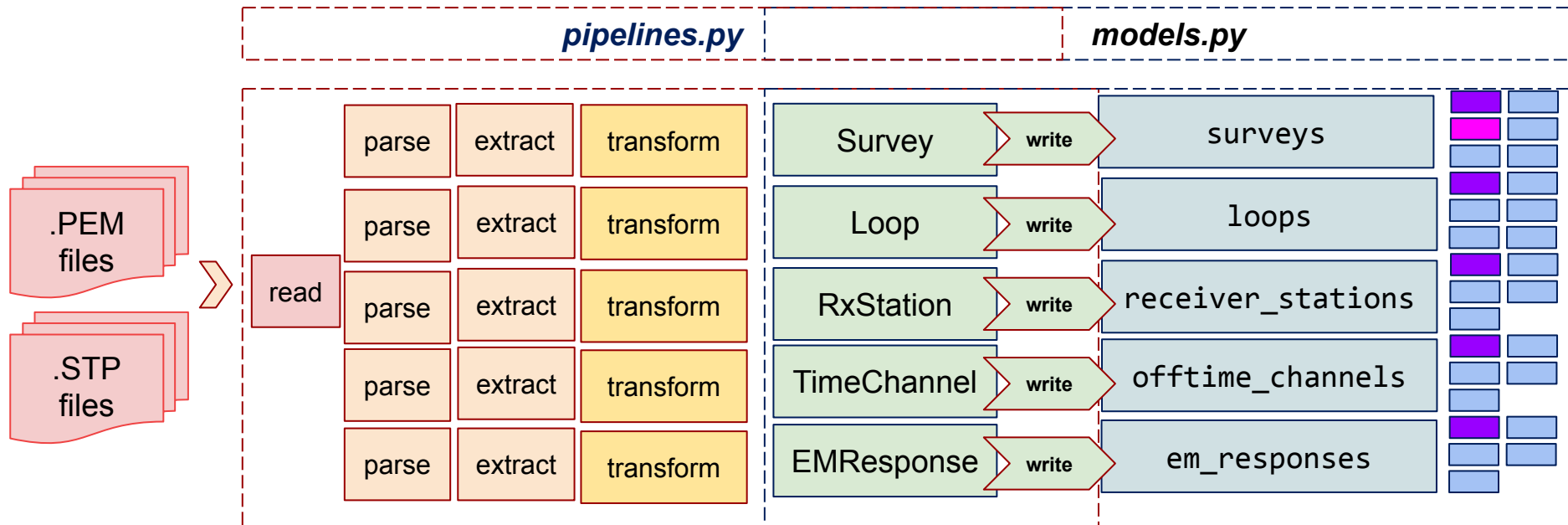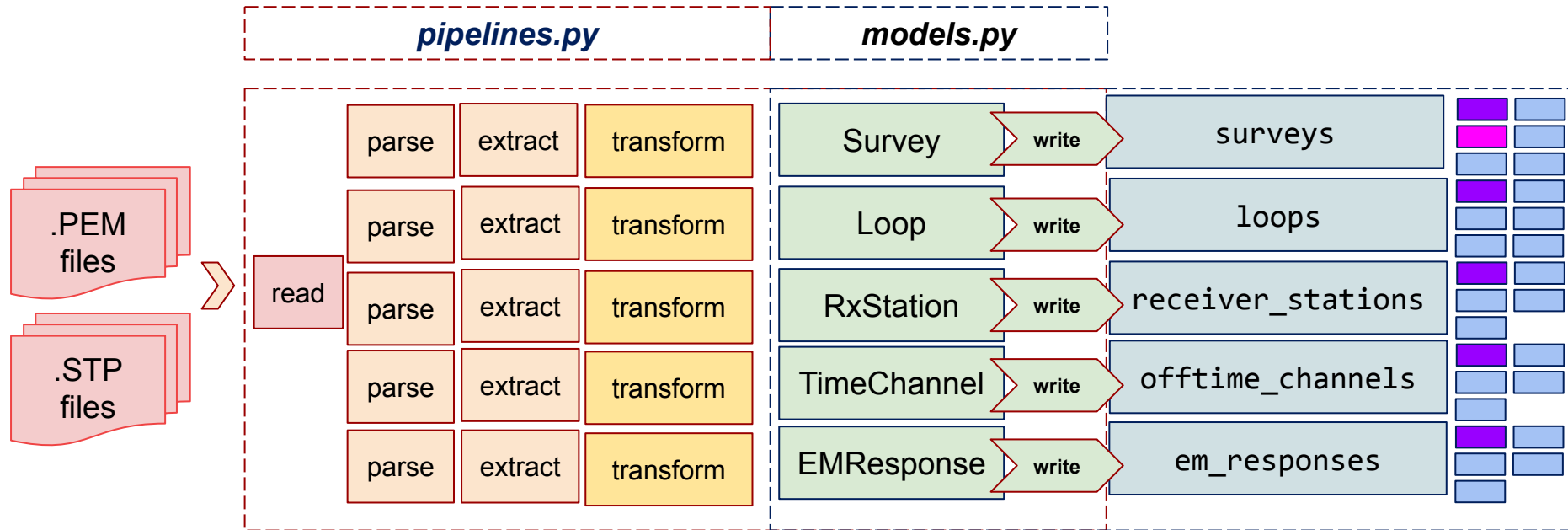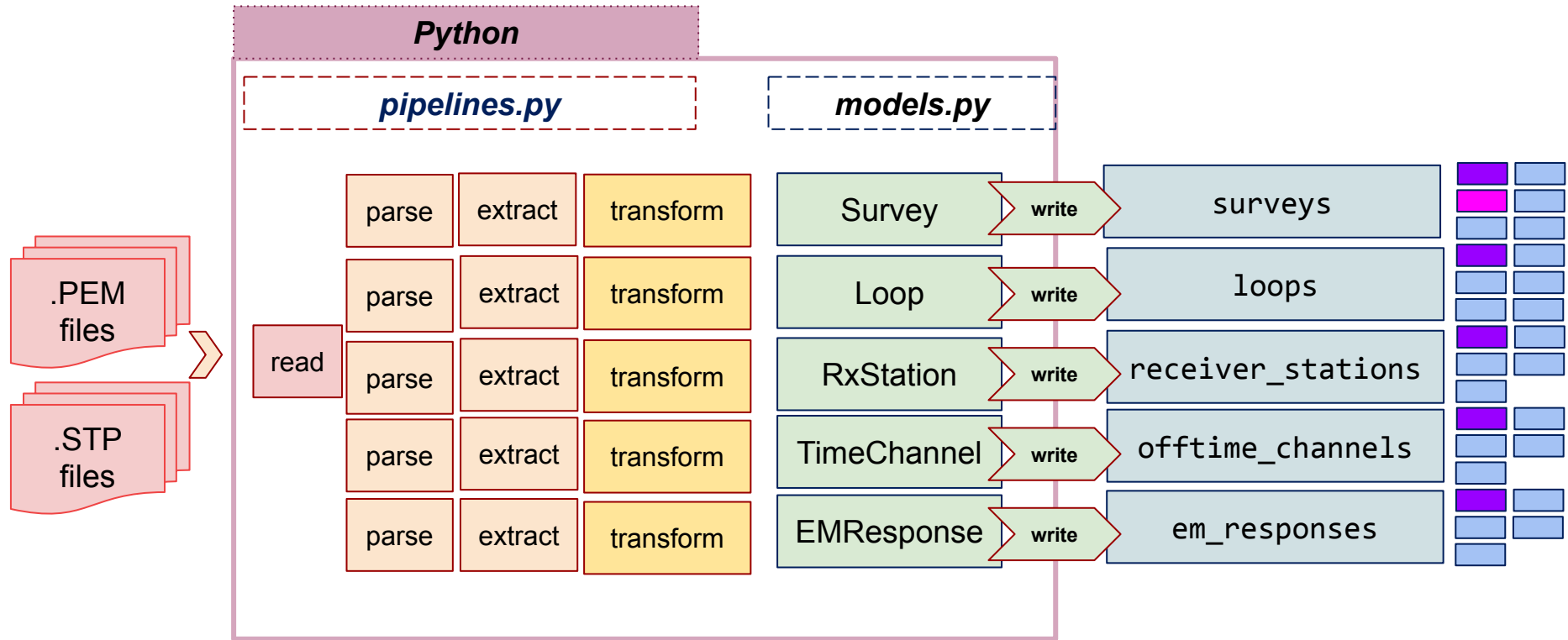
# Basic Components of Data Pipeline

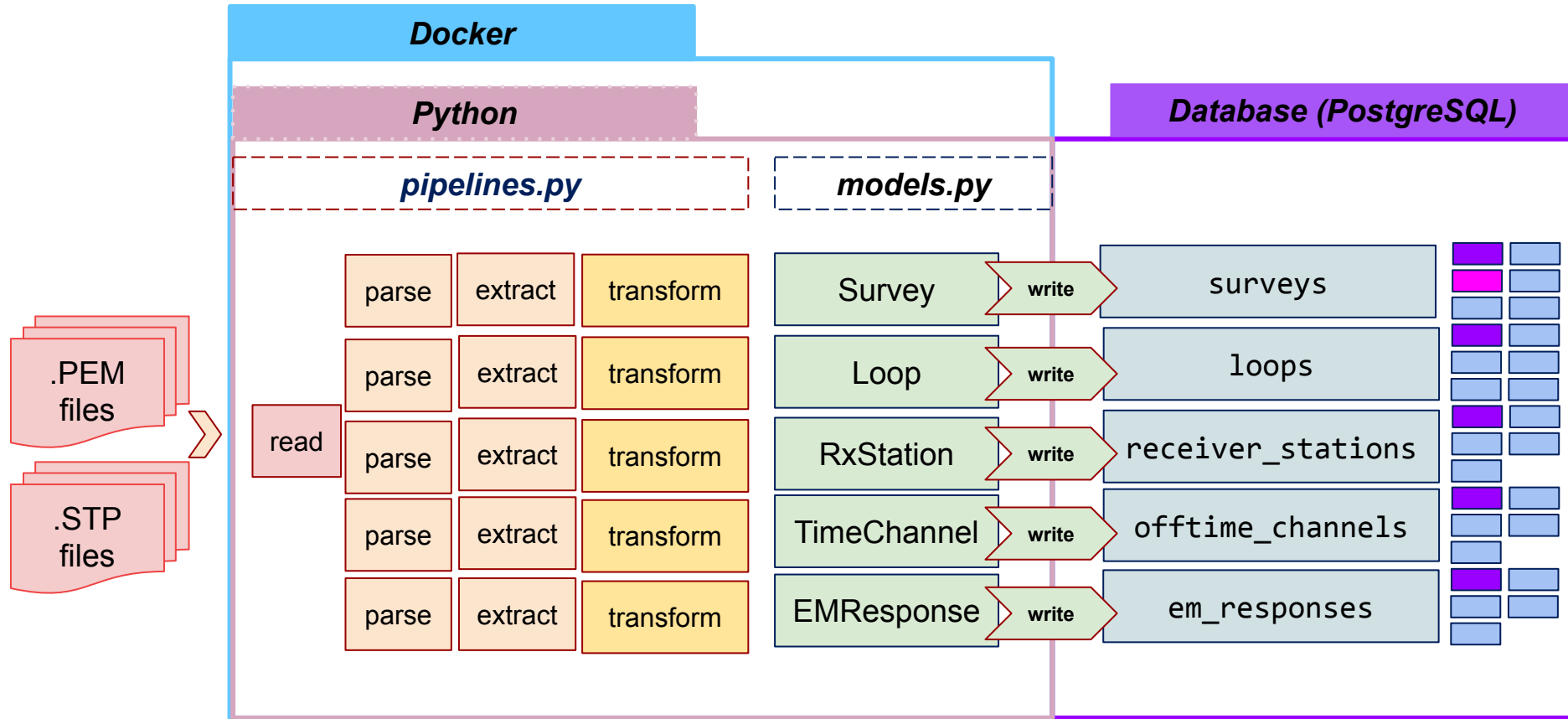# Basic Components of Data Pipeline

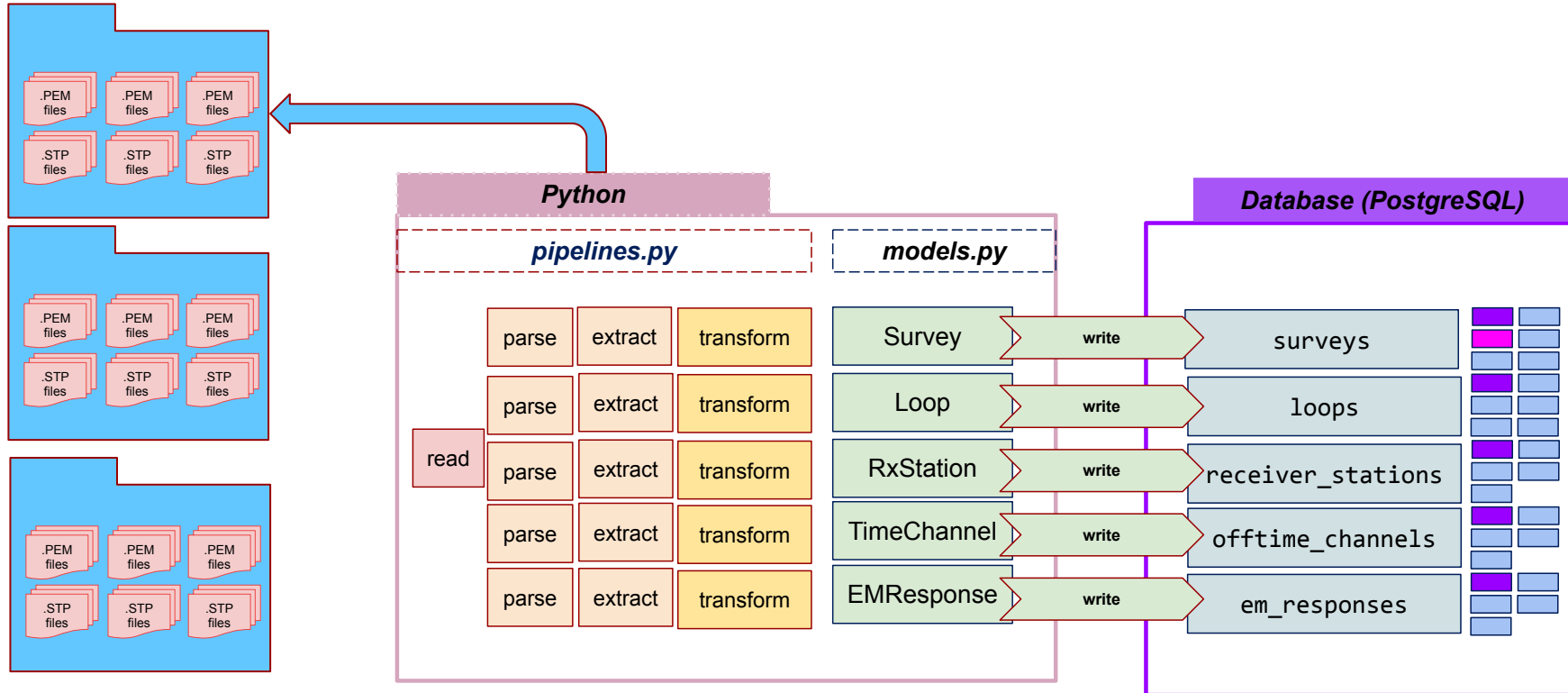# Basic Components of Data Pipeline

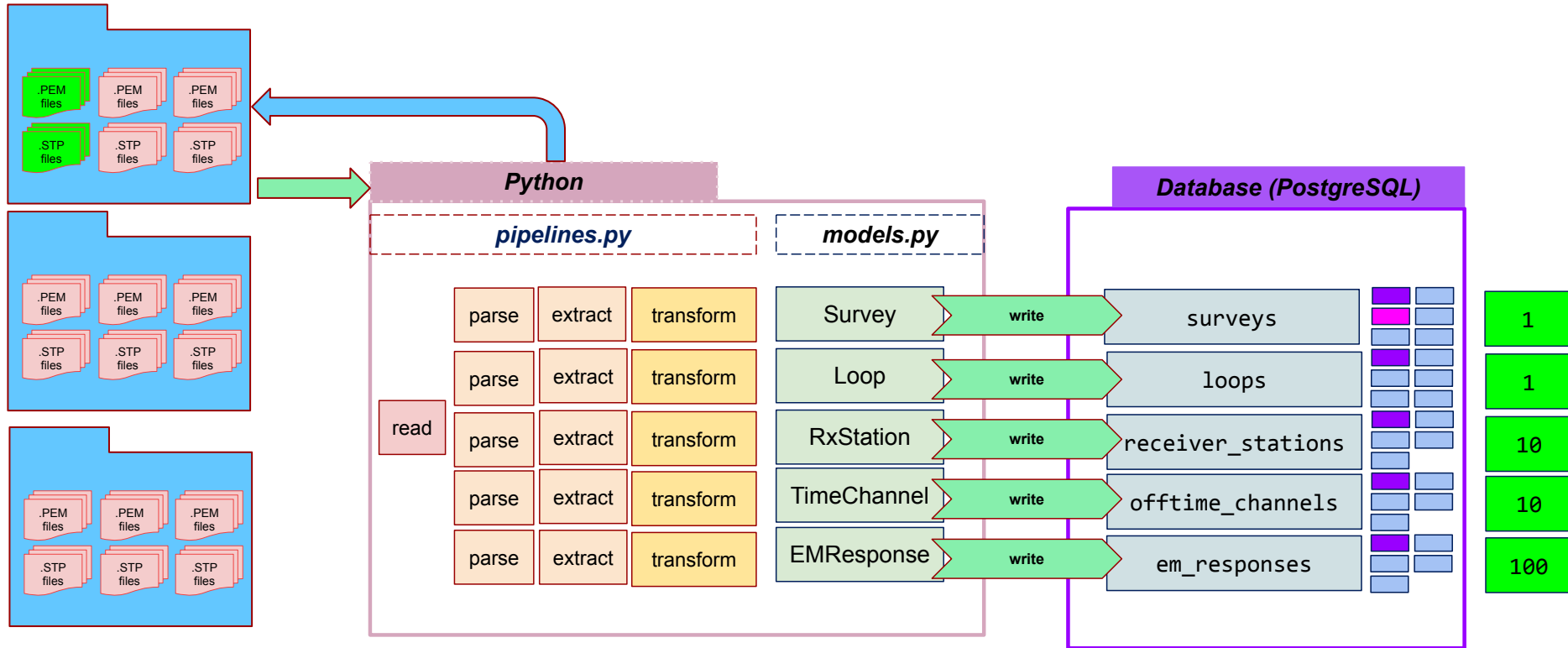# Basic Components of Data Pipeline

# Basic Components of Data Pipeline
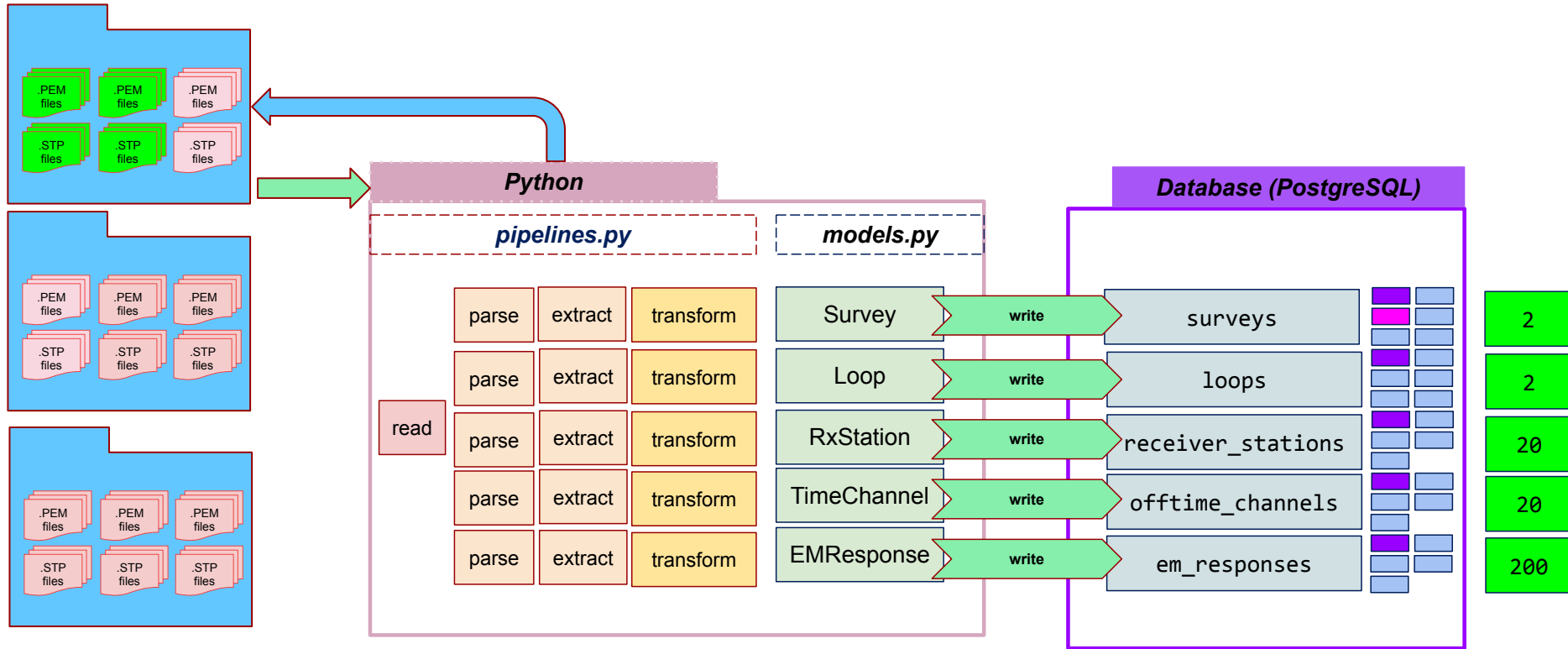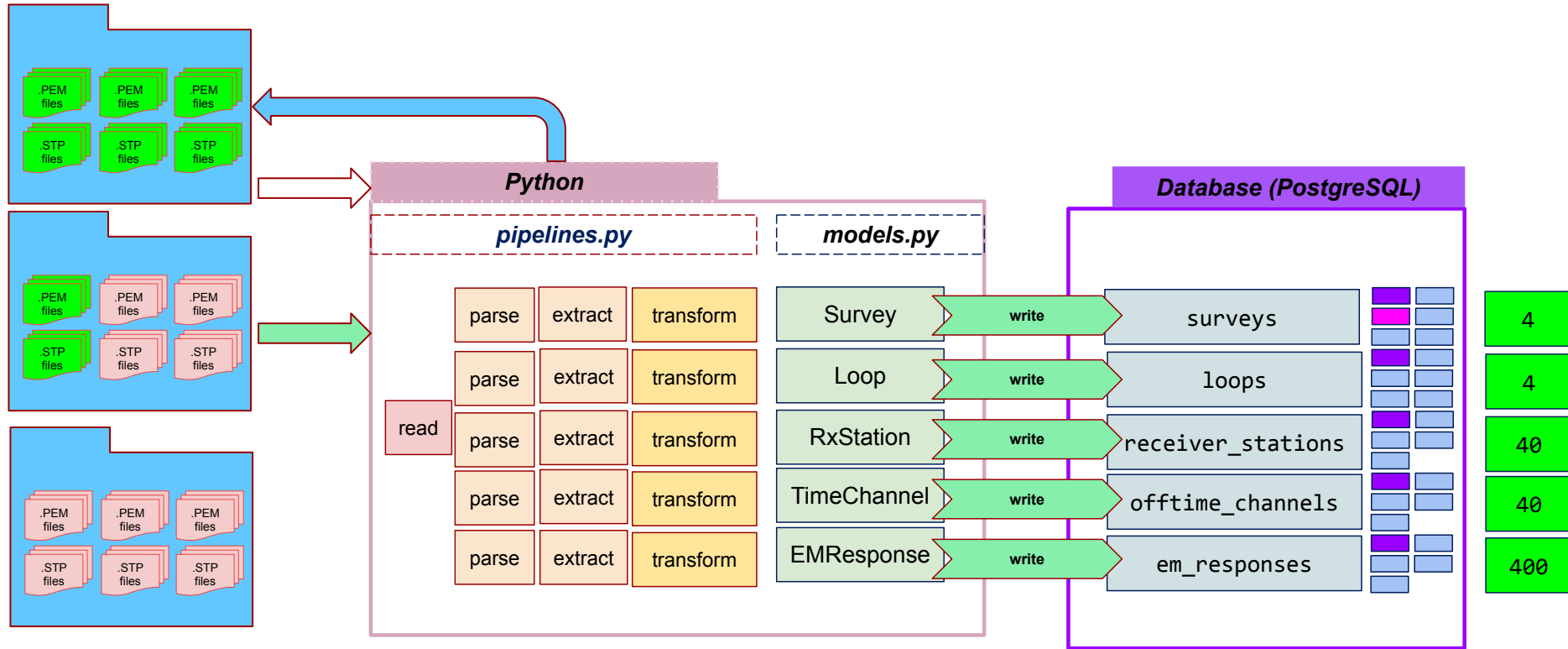
Scaling Data Pipeline

ENVS 5726

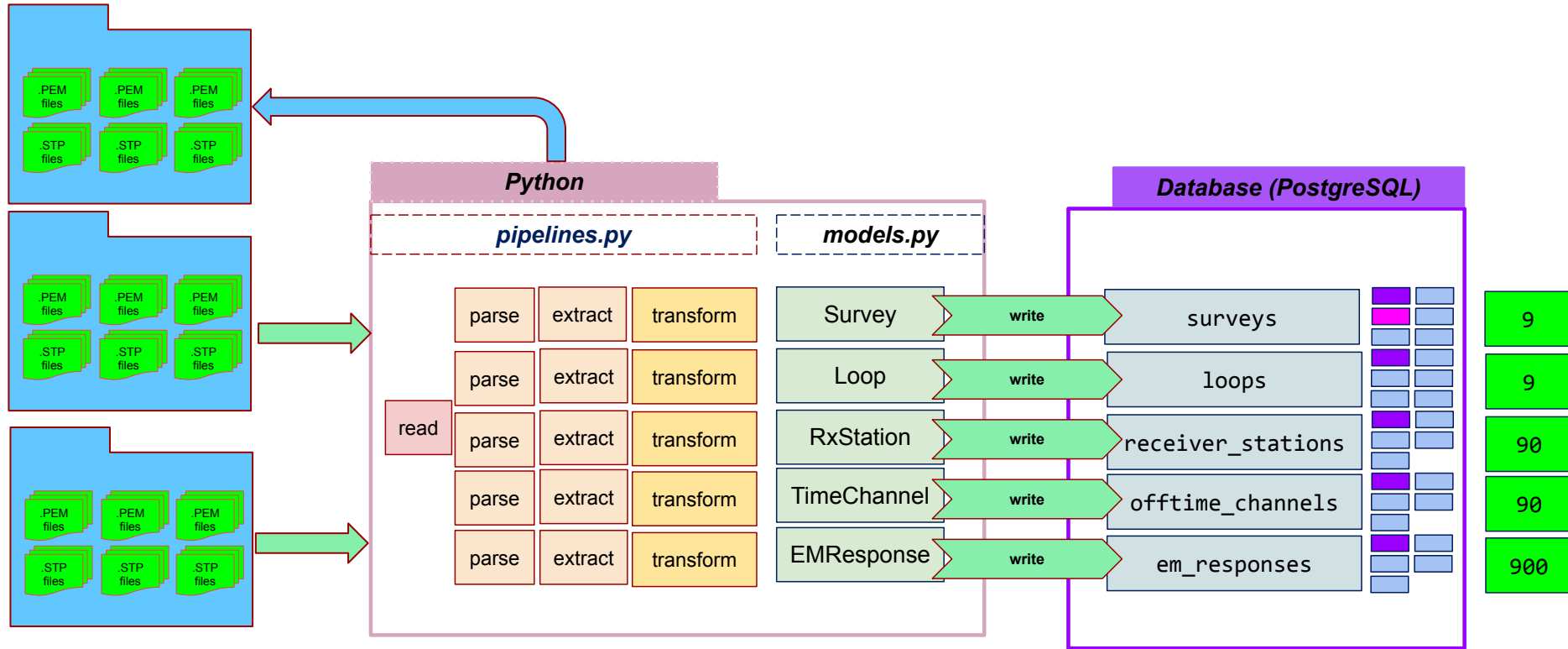# Scaling Data Pipeline

# Scaling Data Pipeline

ENVS 5726

04

# Demo Check-In

# Demo time!

# Data Engineering Pipelines

- **Scalable Task Design with Deduplication Checks:** Each save function (save_survey, save_loops, save_receiver_stations) includes upfront checks that query for existing records before inserting. This prevents duplicate data ingestion if the pipeline re-runs, returning early with counts of existing records which is a critical pattern for reliable ETL pipelines.

- **Batch Processing with Scalable Looping Patterns:** The code processes large datasets efficiently by iterating through lists of parsed records (loops, stations, measurements) and committing them in controlled batches. Rather than individual commits per record, functions like save_receiver_stations loop through all stations, add them to the session, and commit once thus reducing database round-trips and enabling the pipeline to scale to thousands of measurements per file.

- **Rich Data Extraction with Structured Parsing:** The parsing tasks (e.g., parse_measurements, parse_stations) use regex patterns and line-by-line iteration to extract complex, nested data structures from unstructured file content. Each function returns well-typed dictionaries with metadata context (line names, component types), ensuring downstream storage functions receive clean, validated data ready for database insertion.

# Data Engineering
# Models and Schema

- **Robust Data Integrity with Cascades:** The models explicitly define parent-child relationships using *cascade="all, delete-orphan"*. This ensures that when a parent record (like a survey) is deleted, all associated children (loops, stations, measurements) are automatically cleaned up, preventing orphaned records in the database.

- **Explicit Database Constraints and Indexing:** proactively enforces data quality and performance through __table_args__. It defines UniqueConstraint to prevent logical duplicates (e.g., ensuring a specific loop_point_number only appears once per survey) and creates specific Index entries (e.g., on easting and northing) to speed up lookups.

- **Intentional Typing and Spatial Awareness:** The models go beyond basic types by integrating specialized fields like Enum for fixed categories (e.g., ComponentEnum) and PostGIS Geometry columns with specific SRIDs (Spatial Reference System Identifiers). This ensures the database enforces valid values and correctly handles complex spatial data rather than relying solely on application logic.

05

Questions and Conclusions

# Conclusion

Knowledge about underlying data and user objective is critical to design a resilient and effective data pipeline.

A thoughtful schema design can go a long way (quality, resilience, and reliability).

Off-the-shelf orchestration tools are critical for pipeline observability.

Learning and applying the fundamentals (earth science, data engineering, software development) is critical despite LLMs.

Get feedback from users and stakeholders, and incorporate that feedback into the logic.

# Discussion

# Questions

fin