# Sample Sequence Diagram of E2E Architecture

| Snowflake | Databricks | GitHub | Streamlit |
|-----------|------------|--------|-----------|

Snowflake → Databricks: Grab all raw, production premiums and claims data.

Databricks → Databricks: Triggers a Job Run

**Note (Databricks):**
Filter to correct business segments.
Drop invalid records, adjust outlier records. \n Aggregate and merge.

**Note (Databricks):**
Identify key information in new columns.
Create new feature columns for segmentation.
Generate transformations (grouping, percentiles, etc.).

**Note (Databricks):**
Calculate on-level premium.
Trend losses.
Develop losses.

**Note (Databricks):**
Remove duplicate/redundant segments.
Exclude non-credible segments.
Identify impactful segments (good & bad).

Databricks → Databricks: Job Logging

Databricks → Snowflake: Append-Only Data Storage

**Note (Snowflake):** Outputs from Databricks Job run are appended to Snowflake tables.

Databricks → GitHub: Regression Testing

GitHub → GitHub: Pull Request Conditions

GitHub → GitHub: Check Conditions Pass

GitHub → Streamlit: Update Dashboard

Snowflake → Streamlit: Pull Latest Gold Layer Data

Streamlit → Streamlit: End-User Dashboard

**Note (Streamlit):** Advanced Dashboard Logic.

| Snowflake | Databricks | GitHub | Streamlit |
|-----------|------------|--------|-----------|

# Sample Flow Chart of Medallion Architecture

**Bronze Layer**

Snowflake Raw Data → Bronze Notebook → Parquet Storage

**Silver Layer**

Silver Notebook → Processed Data → Parquet Storage

**Gold Layer**

Gold Notebook → Final Business Logic → Snowflake Append-Only Tables

**Orchestration**

GitHub Actions → Databricks Workflow

Trigger Bronze

Log Status

Trigger Silver

Log Status

Trigger Gold

Log Status

# Idempotency

Idempotency is necessary for application testing.

**∞-loop**

The output remains the same for each execution given the same input.

Idempotency is typically an essential characteristic of data engineering pipelines.

static input → Function → static output

---

## Non-Idempotent

```
def func(float_val):
    float_val += random_int()
    return float_val
```

## Idempotent

```
def func(float_val):
    return float_val
```

# Integration Snapshot Testing (Analytics Engineering)

Saved from a previous run of the parent() function.

```
def parent(df_parquet):
    var_1 = child_1(df_parquet)
    var_2 = child_2(var_1)
    var_3 = child_3(var_2)
```

Expected Output

Actual Output

Comparison Function (Passes?)

Dependency Injection (PARQUET)

Parent

(Yes/No)

Child

Child

Child

# Integration Unit Testing

The expected output is derived using a secondary, trusted calculation method
**(i.e.  manual, LLM, alternative software).**

Expected Output

Comparison Function (Passes?)

Dependency Injection (PARQUET)

Parent

Actual Output

(Yes/No)

Child

Child

Child

# Alternative Testing

"Good software is well-tested software."

**1. SMOKE TESTING**

Does it run?

**2. PERFORMANCE TESTING**

Is it fast?

**3. UI TESTING**

Does the UI operate as expected?

**4. PROPERTY TESTING**

Does it fit within the bounds of an expected output?

# Snapshot Testing

- Probably not conceptually validated (using a previous out as a test case).
- **Fickle:** Needs to be changed often.
- Must be updated with each functional change.
- Verifies that only non-functional changes have occurred (i.e., refactors).
- Should only be updated with peer review.

# Unit Testing

- Conceptually validates a function is performing correctly.
- Ideally, these are written for critical building blocks of a software.
- Helps developers sleep at night.

**Asynchronous**

Read

Write

Read

Read

Write

Write

Write

Asynchronous reads and writes run operations
on separate threads decrease the runtime of
the data pipeline.

**Synchronous**

Read

Read

Write

Write

**Chunking**

**Time Complexity:** O(n)
**Space Complexity:** O(1)

**Parallelization**

**Time Complexity:** O(1)
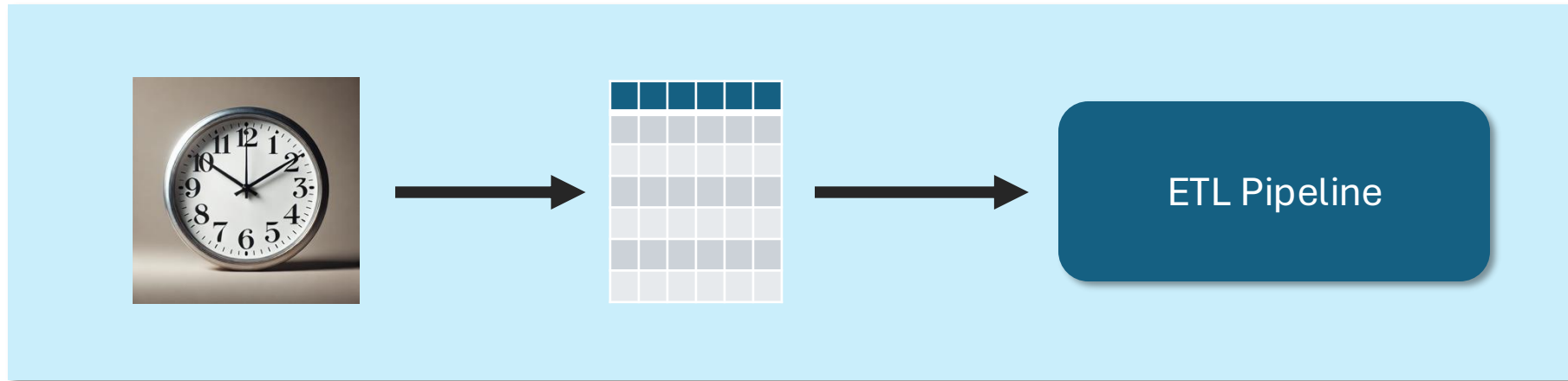**Space Complexity:** O(n)

# Programming Language Trade-Offs

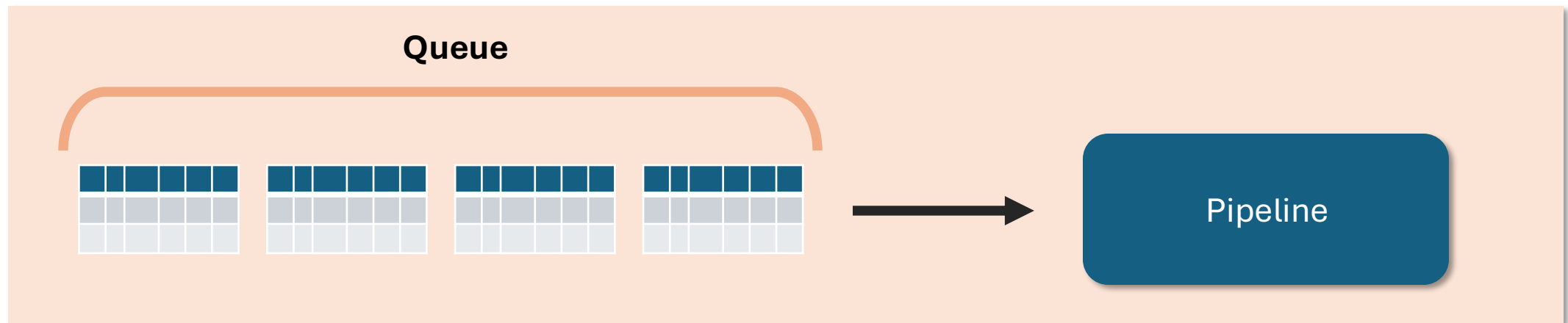| Use Case | Pandas | Polars | PySpark | SQL |
|---|---|---|---|---|
| Small datasets (<10M rows) | ✅ Best | ✅ Faster | ❌ Overkill | ✅ Good for querying |
| Large datasets (>10M rows, fits in RAM) | ❌ Memory issues | ✅ Efficient | ✅ Good | ✅ Good |
| Big Data (TB-scale, distributed) | ❌ Impossible | ❌ Limited | ✅ Best | ✅ Best |
| Parallel processing | ❌ Single-threaded | ✅ Multi-threaded | ✅ Distributed | ✅ Query optimizations |
| Complex ETL | ✅ Simple | ✅ Efficient | ✅ Distributed pipelines | ✅ SQL transformations |
| ML/Statistical modeling | ✅ Best for ML | ✅ Works | ❌ Spark ML (limited) | ❌ Not ideal |
| Cloud-based processing | ❌ Local only | ❌ Local only | ✅ Cloud (Databricks, EMR) | ✅ Cloud-native (BigQuery, Snowflake) |

# Pandas VS. Polars



⚡ **Pandas vs Polars Performance Scaling** ⚡

# Batch or Stream Processing

**Batch**



**Stream**

# Vectorization

PyArrow
Numpy
Pandas

# CI/CD

**INTERNET COMPUTER**

Cloud Provider as
State Machine
Replication for
Byzantine Fault
Tolerance

Active Data Centers    Boundary Node Data Centers

Upcoming Data Centers

# INTERNET COMPUTER

Subnets consist of 13 or more nodes. This allows for variable decentralization.

Latency is < 2 seconds.

Egress (query calls) doesn't need to go through consensus.
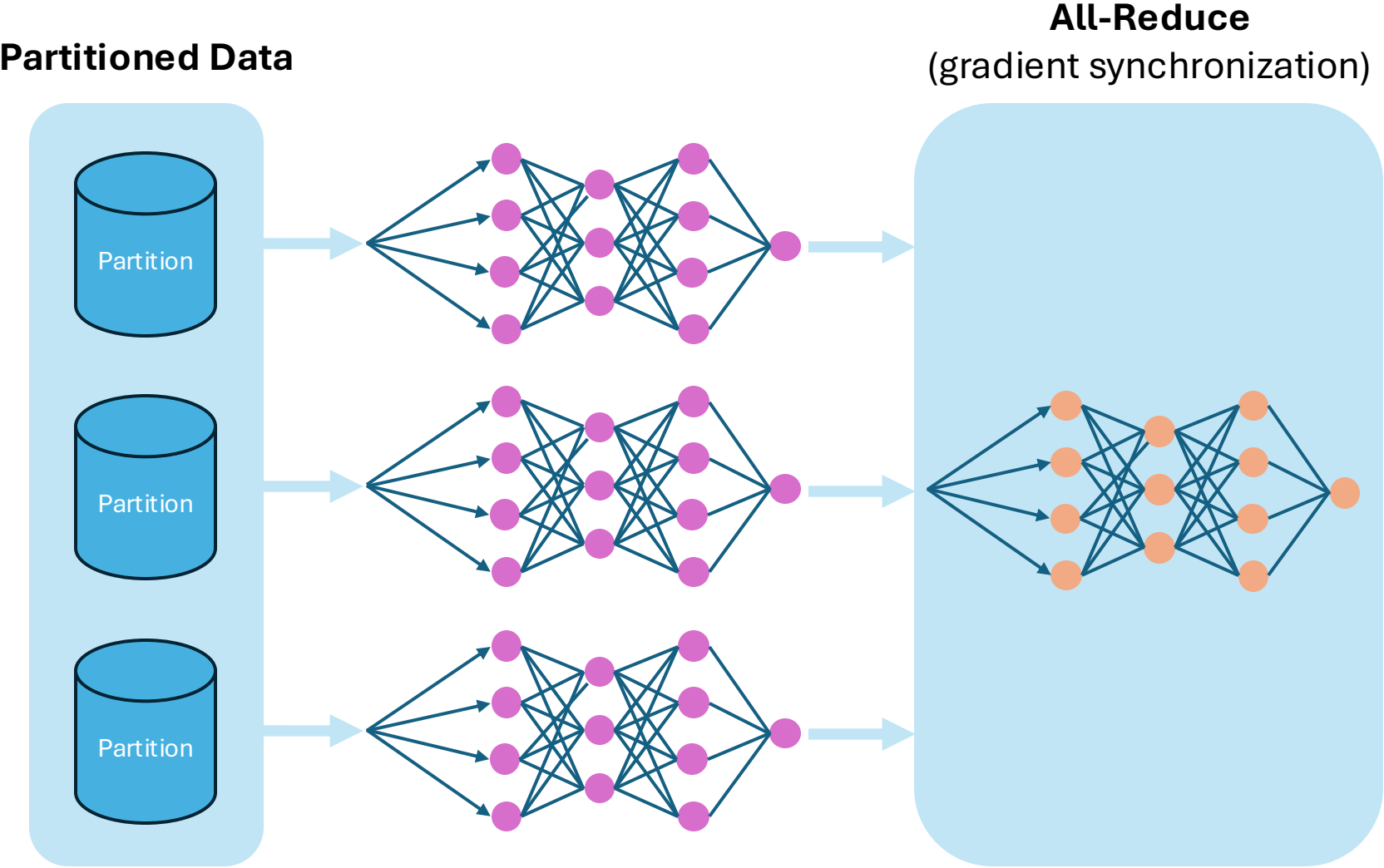
**Reverse gas:** Developers pay for gas fees, and end-users can use the blockchain without needing to pay in ICP (like a normal website).
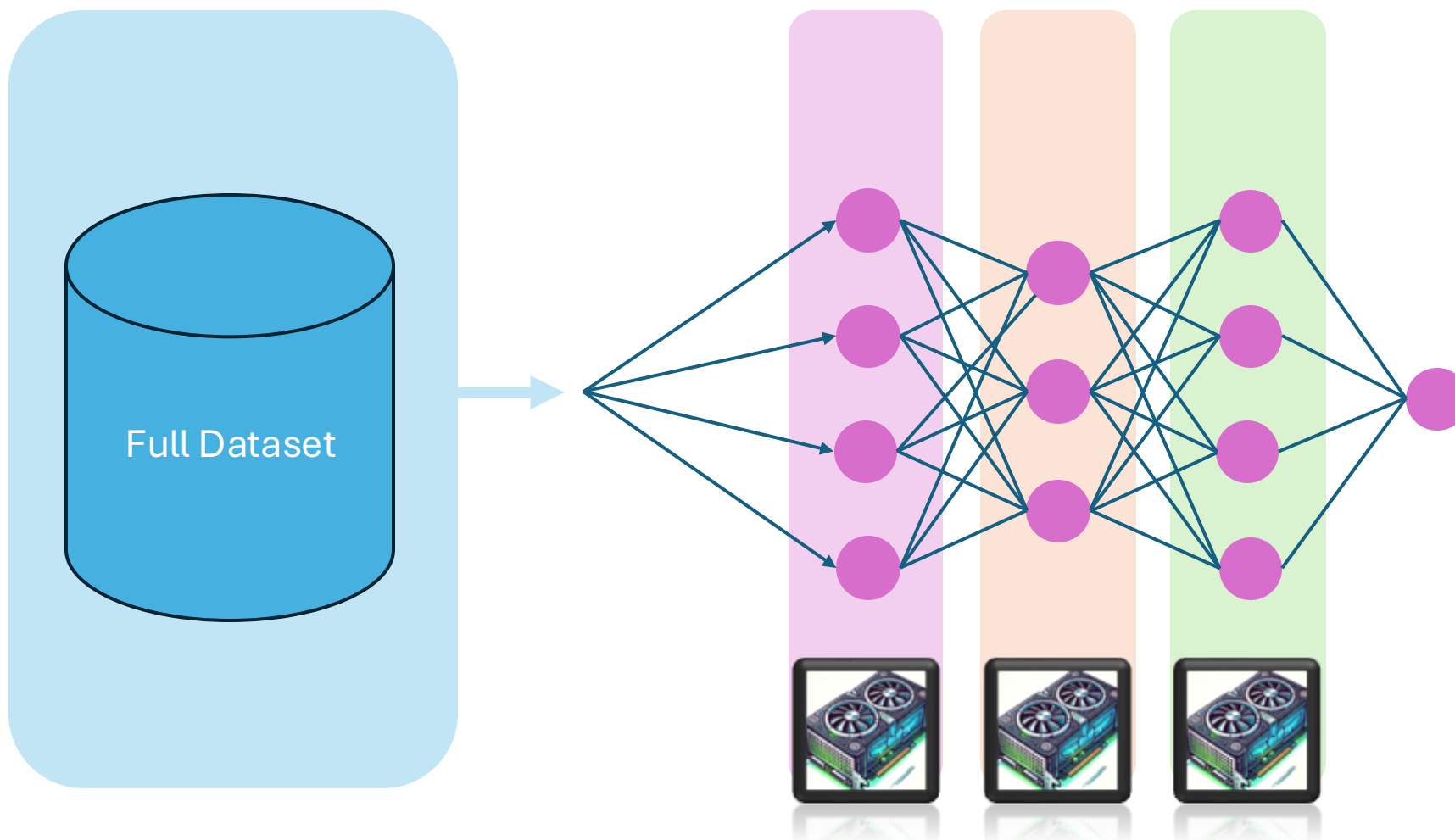
Software becomes immune to cyber hacks or datacenter outages due to any cause of failure.
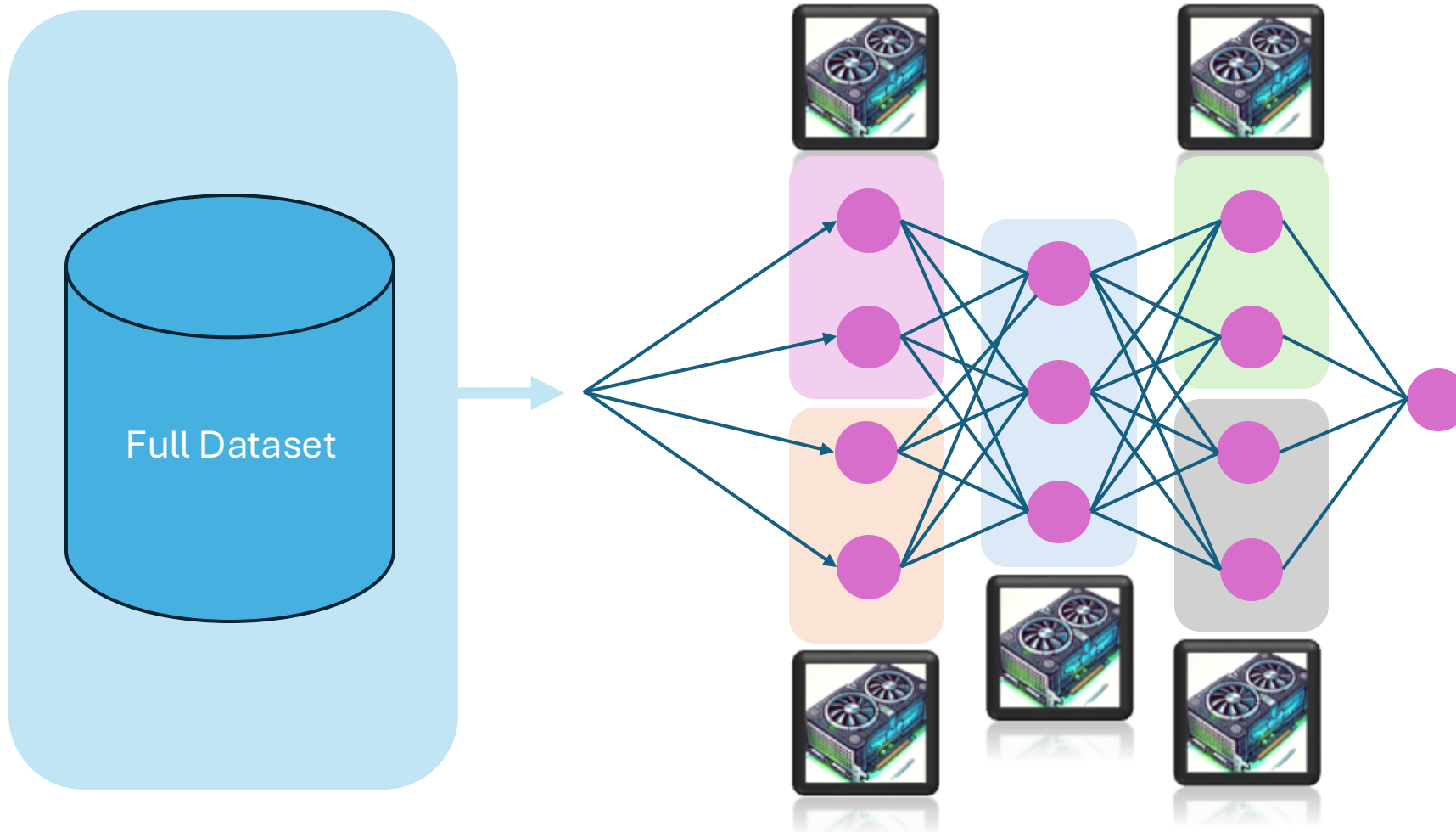
# Model Parallel Training

# Tensor Parallel Training

# RAG: Retrieval Augmented Generation

# Example DL Project