

# Overview of Foundry Data Pipelines

## Understanding the Foundations of Palantir's Data Flow Architecture

- **Definition:** A Foundry data pipeline represents the controlled flow of data from raw sources through transformations into curated, high-quality datasets ready for analytics and machine learning.
- **Core Principles:** Key characteristics include ownership, reliability, quality, and maintainability—ensuring data integrity and consistency across the enterprise.
- **Integration Objective:** To deliver a digital, unified view of organizational reality by combining diverse systems under a common schema for shared use.
- **Strategic Role:** Pipelines serve as the backbone for operational decision-making, governance, collaboration, and scalable analytics within Foundry.

# Why Data Pipelines Matter

Driving Business Value Through Reliable Data Flows



## Operational Decision-Making

Pipelines provide timely, reliable data to support mission-critical business processes, ensuring leaders make decisions with current, trusted insights.



## Data Quality Assurance

Continuous validation and standardization across stages preserve accuracy and consistency, protecting downstream analytics from corruption.



## Scalability and Governance

Built-in mechanisms in Foundry manage growing data volumes while enforcing security and compliance standards automatically.



## Collaboration Enablement

By connecting both technical and non-technical users, pipelines foster a shared understanding of data assets and empower cross-functional work.

# Sample Data Pipeline: Flight Example

## An End-to-End Walkthrough in Foundry

- **Pipeline Flow:** Raw data ingestion → basic cleanup → business transformation → ontology mapping → application layer. Each stage refines and contextualizes data.
- **Example Inputs:** Includes datasets such as passengers\_preprocessed, flight\_alerts\_raw, and status\_mapping\_raw feeding into transformation logic.
- **Transformation Steps:** Cleansing, joining, aggregating, and enriching data to produce analytics-ready tables, e.g., average delay metrics per route.
- **Key Operations:** Filtering, casting, joining, aggregating, and enriching ensure consistent, validated output aligned to downstream analytics.

# When to Build a Pipeline

## Decision Framework for Data Engineering in Foundry



### Recurring Data Updates

A pipeline is justified when source data changes regularly—daily, hourly, or in real time—requiring consistent ingestion and transformation.



### Multiple Consumers

When transformed data is needed by several teams or applications, pipelines provide a shared, scalable delivery mechanism.



### Complex Logic and Governance

Sophisticated business rules, validation, or versioning demand the structured control that pipelines offer.



### Production and Collaboration Needs

Pipelines support operational decision-making and multi-user contribution, unlike one-off analyses or prototypes.

# Pipeline Stages Overview

## The Multi-Stage Architecture of Foundry Pipelines



### Stage 1: Data Connection

Ingests raw data from databases, APIs, or storage systems into Foundry while ensuring authentication, scheduling, and initial validation.



### Stage 2: Datasource Project

Applies schema alignment, deduplication, and basic cleanup to raw datasets to create standardized, reliable sources.



### Stage 3: Transform Project

Implements business logic—joins, aggregations, enrichments—to produce canonical datasets for reuse across the platform.



### Stage 4 & 5: Ontology and Workflow

Maps data to business entities (Ontology) and enables applications, dashboards, and actions for end users (Workflow).

# Data Connection Stage

## Establishing Trusted Ingestion into Foundry

- **Purpose:** To securely connect external systems—databases, APIs, cloud storage—and ingest raw data into Foundry with proper authentication and schedules.
- **Common Connectors:** Supports Oracle, MySQL, PostgreSQL, S3, Azure Blob, Kafka, and ERP/CRM platforms like SAP and Salesforce for broad enterprise integration.
- **Best Practices:** Use incremental syncs to minimize transfer cost, enable retry logic, and monitor health metrics for reliable data flow.
- **Governance:** Document dependencies and credentials management to maintain security and auditability across data ingress points.

# Datasource Project Stage

## Standardizing and Cleaning Raw Data

- **Purpose:** Transforms raw ingested data into standardized, schema-aligned datasets that serve as trusted building blocks for downstream projects.
- **Core Activities:** Remove system artifacts, rename columns, cast data types, handle missing values, and apply baseline validation checks.
- **Transformation Patterns:** Schema standardization, type casting, deduplication, basic filtering, and null handling ensure structural integrity.
- **Best Practice:** Create one Datasource Project per logical data source, ensuring modularity and ease of debugging and versioning.

# Transform Project Stage

## Applying Business Logic and Deriving Value

- **Purpose:** Implements complex business logic to produce reusable, canonical datasets that drive analytics, KPIs, and machine learning models.
- **Core Activities:** Join multiple datasources, calculate derived metrics, apply aggregations, and create domain-specific datasets for reuse.
- **Common Patterns:** Dimensional modeling, denormalization, Slowly Changing Dimensions, and data enrichment improve accessibility and performance.
- **Governance and Testing:** Each Transform Project requires controlled changes, validation, and impact analysis to maintain data quality and lineage.

# Ontology Project Stage

## Mapping Data to Business Entities

- **Purpose:** Defines semantic mappings between datasets and real-world entities, transforming technical data into business-relevant objects, properties, and links.
- **Key Components:** Objects represent entities, properties define attributes, and links establish relationships such as Employee–Department or Product–Order.
- **Business Value:** Provides a unified semantic layer enabling business users to explore data intuitively through familiar domain concepts.
- **Security and Actions:** Supports dynamic security, fine-grained permissions, and write-back actions that empower operational workflows.

# Workflow Project Stage

## Building End-User Applications on Top of the Ontology

- **Purpose:** Develop operational applications, dashboards, and reports leveraging the Ontology's semantic layer and data actions.
- **Key Tools:** Workshop for operational workflows, Contour for analytics, Quiver for time-series insights, and Reports for distribution.
- **Capabilities:** Integrates actions for write-back to systems, interactive dashboards for stakeholders, and automated analytics pipelines.
- **Outcome:** Transforms curated data into business impact by enabling decision-making, monitoring, and real-time collaboration.

# Recommended Project Structure & Naming

## Organizing Foundry Pipelines for Scalability and Clarity

- **Folder Hierarchy:** Organize projects by stage—Raw, Datasource, Transform, Ontology, Applications, and Code Repositories—to ensure modularity and traceability.
- **Naming Conventions:** Use descriptive, domain-based names with consistent casing (e.g., finance\_revenue\_monthly). Avoid ambiguous abbreviations.
- **Governance Benefits:** Structured naming improves searchability, access control, and onboarding of new collaborators.
- **Practical Setup:** Begin with workspace creation, define permissions, configure connections, and progressively build transformations.

# Prototyping in Foundry

## Rapid Exploration Before Production Deployment

- **Purpose:** Allows engineers to experiment with logic, test transformations, and validate datasets safely before moving to production.
- **Tools:** Pipeline Builder for visual prototyping, Notebooks for interactive PySpark analysis, and Dataset Preview for real-time inspection.
- **Best Practices:** Work in personal branches, use sample data, document assumptions, and validate early with business users to prevent rework.
- **Outcome:** Accelerates innovation cycles while maintaining isolation and data integrity across teams.

# Transforming Data in Foundry

## Applying Logic, Validation, and Enrichment at Scale

- **Transformation Types:** Includes schema alignment, filtering, joins, aggregations, and advanced logic such as window functions or JSON parsing.
- **Performance Best Practices:** Push filters early, optimize joins, partition effectively, and cache reused datasets to improve efficiency.
- **Data Quality:** Validate inputs, handle nulls and duplicates, and add quality metrics and alerts for transparency and governance.
- **Code Quality:** Follow modular, idempotent, and testable design principles with descriptive naming and documentation.

# Maintaining Pipelines in Production

## Ensuring Reliability, Performance, and Data Quality



### Health Monitoring

Implement health checks for job success, data freshness, schema stability, and row count validation.



### Performance Optimization

Track build duration, resource usage, and bottlenecks. Optimize joins, partitions, and scheduling frequency.



### Version Control & Deployment

Use branches for changes, conduct peer reviews, and test before merging to maintain stability.



### Troubleshooting

Analyze logs, compare historical runs, and leverage data lineage visualization to identify root causes.

# Pipeline Types Overview

## Batch, Incremental, and Streaming Pipelines in Foundry

- **Batch Pipelines:** Recompute full datasets at each run—simple, consistent, ideal for small-to-medium data volumes and high-complexity logic.
- **Incremental Pipelines:** Process only new or changed data, enabling faster performance and scalability for large datasets with frequent updates.
- **Streaming Pipelines:** Handle real-time data with sub-minute latency, suitable for event-driven or monitoring applications requiring immediacy.
- **Trade-offs:** Batch is simple but slow; Incremental balances efficiency; Streaming delivers speed with higher operational complexity.

# Pipeline Comparison and Use Cases

## Selecting the Right Type for Your Data Workload

- **Batch Pipelines:** Recomputes entire datasets at each run; best for periodic reports or complex recalculations with manageable data volumes.
- **Incremental Pipelines:** Processes only deltas; ideal for large datasets with frequent updates like IoT feeds or transaction logs.
- **Streaming Pipelines:** Processes continuous data in real time for immediate responses such as fraud detection or operational alerts.
- **Key Metrics:** Batch = simplicity, Incremental = efficiency, Streaming = immediacy; trade-offs balance latency, complexity, and cost.

# Data Lineage in Foundry

## Visualizing and Understanding Data Flow

- **Purpose:** Provides end-to-end visibility of data flow from source to output, mapping all upstream and downstream dependencies.
- **Capabilities:** Supports impact analysis, permission checks, and schedule visualization to trace and understand pipeline dependencies.
- **Use Cases:** Used for debugging, compliance documentation, and optimization by revealing redundant or outdated data paths.
- **Value:** Enhances trust and transparency by showing exactly how each dataset is produced and consumed across the enterprise.

# Build Schedules Application

## Automating and Monitoring Pipeline Executions



### Purpose

Manages the scheduling and orchestration of pipeline builds across environments to ensure timely data availability.



### Capabilities

Provides centralized dashboards for viewing schedule status, success/failure rates, and historical performance metrics.



### Best Practices

Stagger schedules to prevent resource contention, configure failure notifications, and document interdependencies.

### Outcome

Ensures reliable, predictable data delivery and proactive monitoring for business-critical pipelines.

# Data Health Application

## Monitoring Quality and Reliability Across Pipelines

- **Purpose:** Centralizes visibility of data quality, completeness, and freshness across the entire pipeline ecosystem.
- **Key Features:** Displays health metrics, historical trends, and alerts; integrates with build schedules for proactive monitoring.
- **Recommended Checks:** Track row counts, null rates, schema consistency, and data recency for inputs, intermediates, and outputs.
- **Value:** Enables quick detection of anomalies, ensuring stakeholders trust and rely on data-driven operations.

# Foundry Branching

## Safe and Collaborative Development Workflow



### Purpose

Allows developers to experiment safely in isolated branches without impacting production pipelines or data.



### Workflow

Create branch → make and test changes → submit proposal → review → merge to Main; ensures accountability and quality.



### Benefits

Supports collaboration, peer review, and rollback while maintaining a unified end-to-end environment.



### Governance

Inactive branches auto-close, reducing clutter and ensuring controlled, compliant development practices.

# Development Best Practices

## Building Reliable and Scalable Foundry Pipelines

- **Start Simple:** Begin with batch pipelines, validate logic, and evolve toward incremental or streaming as needs mature.
- **Prototype and Test Early:** Use feature branches, sample data, and health checks to validate transformations before deployment.
- **Documentation and Reviews:** Document business logic, decisions, and code changes. Enforce peer reviews for quality assurance.
- **Continuous Monitoring:** Track metrics, maintain health checks, and refine pipeline performance iteratively.

# Data Quality & Performance Optimization

Ensuring Accuracy, Efficiency, and Reliability

- **Validate Early:** Run data quality checks at every stage—inputs, transformations, and outputs—to detect issues before propagation.
- **Optimize Joins & Filters:** Push filters upstream and use broadcast joins or partitioning to reduce compute costs and improve speed.
- **Monitor Continuously:** Track job duration, row counts, schema drift, and freshness metrics to maintain steady performance.
- **Alert and Iterate:** Automate alerts for anomalies and refine transformations regularly to sustain long-term reliability.

# Collaboration & Governance

## Empowering Teams While Ensuring Compliance

- **Branching and Reviews:** Encourage collaboration via branching workflows and mandatory peer reviews to ensure quality and accountability.
- **Clear Naming & Structure:** Use descriptive project and dataset names with consistent hierarchies to aid discoverability and maintainability.
- **Access Controls:** Apply least-privilege principles and granular permissions at dataset and project levels for security and compliance.
- **Knowledge Sharing:** Maintain documentation, wikis, and shared notebooks to preserve institutional knowledge across teams.

# Conclusion & Key Takeaways

## Mastering the Foundations of Palantir Foundry Data Pipelines

- **Foundational Mastery:** Understanding the complete data pipeline lifecycle—from ingestion to application—empowers scalable, trustworthy analytics.
- **Operational Excellence:** Reliability, maintainability, and data quality must remain at the heart of every design decision.
- **Collaboration and Governance:** Cross-functional teamwork and strict data governance drive consistent outcomes and compliance.
- **Continuous Learning:** Evolving alongside Foundry's ecosystem ensures data engineers remain adaptive and future-ready.