

COVID-19 Data Pipeline - Technical Documentation

Executive Summary

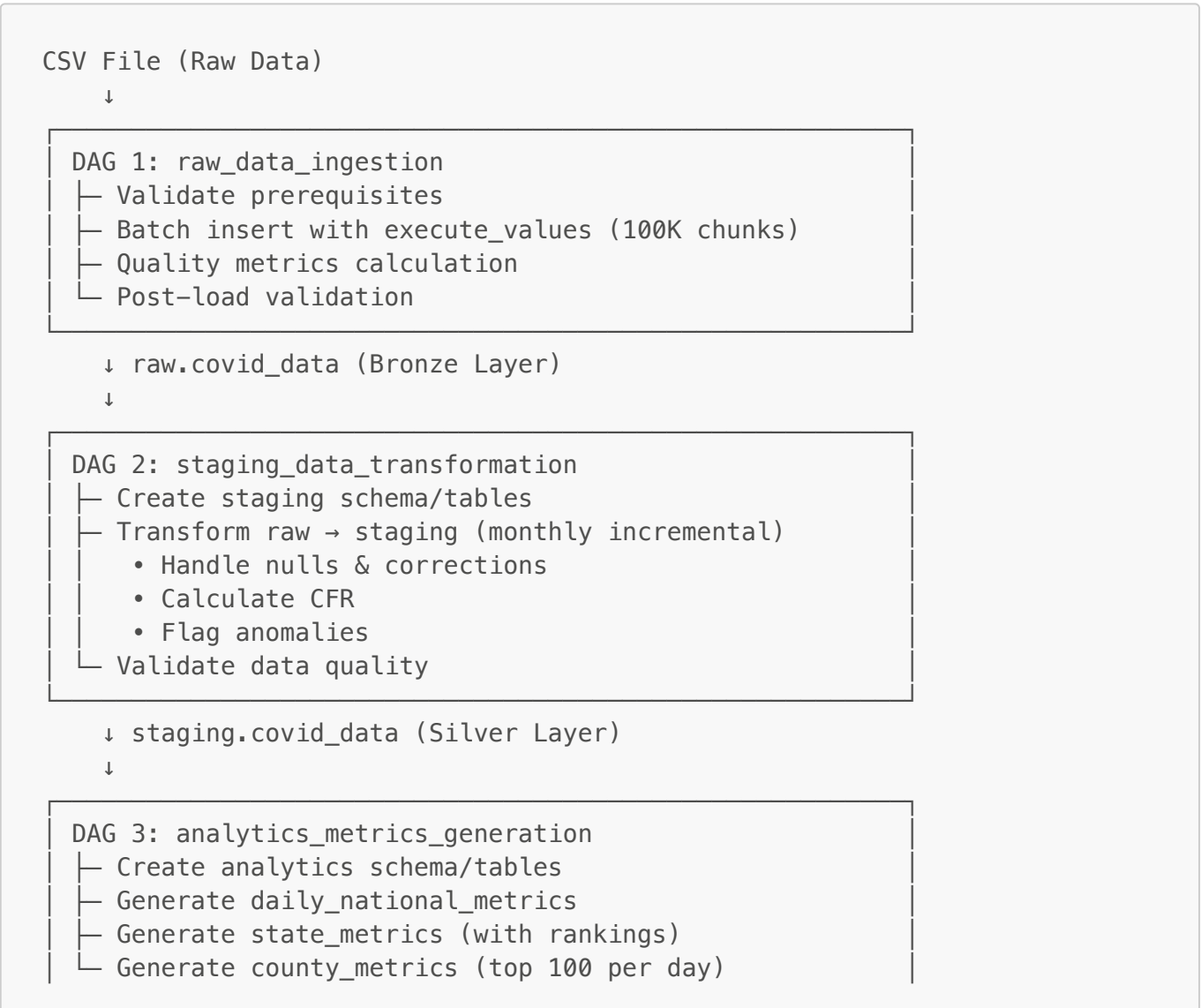
This document describes a COVID-19 data processing pipeline built with Apache Airflow and PostgreSQL. The system implements a **Medallion Architecture** (Bronze → Silver → Gold) with three distinct data layers: Raw, Staging, and Analytics. The pipeline processes millions of COVID-19 records through automated ETL workflows, providing business-ready metrics for executive dashboards.

Key Statistics:

- **Partitioning Strategy:** Monthly (84 partitions covering 2020-2026)
- **Update Frequency:** Monthly catchup processing for new records and one time for inserting raw data from csv file
- **Data Quality:** Automated validation with anomaly detection

Architecture Overview

Pipeline Flow



↓ analytics.* (Gold Layer)
↓

- Streamlit Dashboard

- Executive Overview
 - National Daily Metrics
 - State Comparison
 - County Hotspots

Technology Stack

Component	Technology	Version
Orchestration	Apache Airflow	2.11.0
Database	PostgreSQL	15+
Language	Python	3.11+
Data Processing	Pandas	2.x
Visualization	Streamlit + Plotly	Latest
Containerization	Docker Compose	Latest

DAG 1: Raw Data Ingestion

Configuration

```
dag_id: "raw_data_ingestion"  
schedule_interval: None (manual trigger one time)  
catchup: False
```

Task Flow

```
validate_prerequisites → ingest_csv_to_postgres → post_load_quality_check
```

Task 1: validate_prerequisites

Purpose: Ensure system readiness and create necessary database structures

Operations:

- 1. Verify CSV file exists at `/opt/airflow/data/covid-data.csv`
- 2. Check database connectivity
- 3. Create `raw` schema if not exists

4. Create `raw.covid_data` table (partitioned by `report_date`)
5. Generate 84 monthly partitions (2020-01 through 2026-12)
6. Create default partition for out-of-range dates

Key Feature - Monthly Partitioning:

```
-- Example partitions created:
raw.covid_data_2020_01  -- 2020-01-01 to 2020-02-01
raw.covid_data_2020_02  -- 2020-02-01 to 2020-03-01
...
raw.covid_data_2026_12  -- 2026-12-01 to 2027-01-01
raw.covid_data_default  -- Catch-all for other dates
```

Performance Impact: Query filtering by date only scans relevant partitions (partition pruning), achieving 10-20x performance improvement.

Task 2: ingest_csv_to_postgres

Purpose: Load raw CSV data into PostgreSQL with quality tracking

Processing Strategy:

- **Chunk Size:** 100,000 rows per iteration
- **Batch Size:** 5,000 rows per database insert
- **Method:** `execute_values()` from `psycopg2` (3-10x faster than `executemany`)

Data Quality Checks (Per Chunk):

```
Quality Metrics Calculated:
├ Null values in primary keys (report_date, county_fips)
├ Negative values in new_cases/deaths (tracked but allowed for
  corrections)
├ Duplicate records
└ Quality thresholds: null_pct < 10%, negative_pct < 5%
```

Transformation Steps:

1. Schema validation (ensure all required columns present)
2. Duplicate removal (keep last occurrence)
3. Filter out rows with null primary keys
4. Column renaming (CSV format → database format)
5. NaT/NaN conversion to SQL NULL

Insert Pattern (UPSERT):

```
INSERT INTO raw.covid_data (...) VALUES %s
ON CONFLICT (report_date, county_fips)
```

```
DO UPDATE SET
    positive_new_cases = EXCLUDED.positive_new_cases,
    ...
```

This ensures idempotency - the pipeline can be rerun safely.

Task 3: post_load_quality_check

Purpose: Validate successful ingestion

Validations:

- Total row count > 0
- No NULL values in primary key columns
- Date range verification (min/max dates)

Output Example:

```
Total rows in database: 2672600
Date range: 2020-01-21 to 2022-04-29
✓ Quality check passed
✓ Summary: {'nulls': 279710, 'negatives': 51184, 'duplicates': 258154}
```

DAG 2: Staging Data Transformation

Configuration

```
dag_id: "staging_data_transformation"
schedule_interval: @monthly
catchup: True (processes historical data)
start_date: 2020-01-30
```

Task Flow

```
create_staging_schema → transform_raw_to_staging → validate_staging_data
```

Monthly Incremental Processing

Key Design Decision: This DAG processes data **monthly** using the execution date:

```
execution_date = "2020-02-15" # Example Airflow execution date
start_date = "2020-02-01" # First day of month
end_date = "2020-03-01" # First day of next month
```

```
# Process all data in February 2020
WHERE report_date >= '2020-02-01' AND report_date < '2020-03-01'
```

Why Monthly?

- Aligns with business reporting cycles
- Optimizes partition utilization
- Enables efficient backfilling with catchup=True
- Reduces processing time (11 min vs hours for full reprocessing)

Task 1: create_staging_schema

Creates staging layer structures:

```
CREATE TABLE staging.covid_data (
  id BIGSERIAL,
  report_date DATE NOT NULL,
  county_fips VARCHAR(16) NOT NULL,
  county_name VARCHAR(64) NOT NULL,
  state_name VARCHAR(32) NOT NULL,

  -- Cleaned metrics
  positive_new_cases INT NOT NULL DEFAULT 0,
  positive_total INT NOT NULL DEFAULT 0 CHECK (positive_total >= 0),
  death_new_count INT NOT NULL DEFAULT 0,
  death_total INT NOT NULL DEFAULT 0 CHECK (death_total >= 0),

  -- Calculated fields
  case_fatality_rate NUMERIC(5,2),
  is_anomaly BOOLEAN DEFAULT FALSE,

  -- Metadata
  data_source VARCHAR(32),
  processed_at TIMESTAMP NOT NULL DEFAULT NOW(),
  raw_ingested_at TIMESTAMP NOT NULL,

  PRIMARY KEY (report_date, county_fips)

  CHECK (case_fatality_rate >= 0 AND case_fatality_rate <= 100)

) PARTITION BY RANGE (report_date);
```

Key Differences from Raw Layer:

- VARCHAR with specific lengths (performance optimization)
- NOT NULL constraints on cleaned fields
- CHECK constraints on totals (must be non-negative)
- Additional calculated fields (CFR, anomaly flag)

Task 2: transform_raw_to_staging

Purpose: Clean, validate, and enrich data

Transformation Logic:

```
-- Handle Nulls
COALESCE(county_name, 'Unknown')

-- Preserve Negative Corrections (Important!)
positive_new_cases can be negative → Correction to total
death_new_count can be negative → Correction to total
-- Only replace NULL with 0

-- Enforce Non-Negative Totals
CASE WHEN positive_total < 0 THEN 0 ELSE positive_total END

-- Calculate Case Fatality Rate
CASE
  WHEN positive_total > 0 THEN
    ROUND((death_total::NUMERIC / positive_total * 100), 2)
  ELSE 0
END as case_fatality_rate

-- Anomaly Detection
is_anomaly = TRUE IF:
  • death_total > positive_total (inconsistent)
  • CFR > 50% (unrealistic)
  • positive_total < 0 (impossible)
  • death_total < 0 (impossible)
```

Critical Design Decision - Negative Values:

In COVID reporting, negative **new_cases** or **new_deaths** represent **corrections** to previously reported totals:

```
Day 1: total_cases = 1000, new_cases = 100
Day 2: total_cases = 1050, new_cases = 50
      (Discovered 50 of the 1000 were miscounted)
      new_cases should be: 1050 - 1000 = 50
      But raw data shows: new_cases = -50 (correction)
```

Therefore: We preserve negative daily counts but enforce non-negative cumulative totals.

Task 3: validate_staging_data

Validation Checks:

- Data consistency (deaths ≤ total cases)

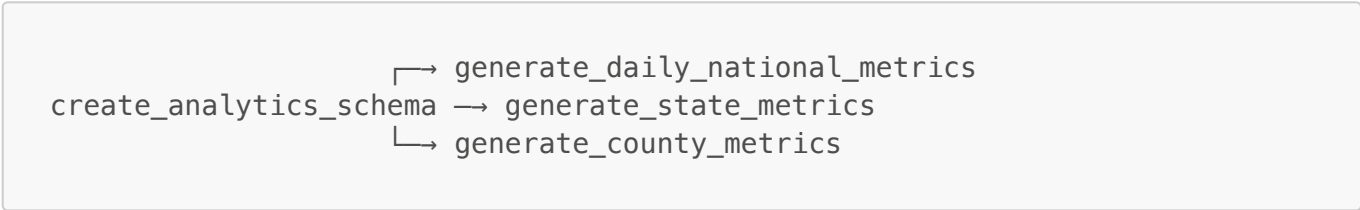
- Anomaly rate < 10%
- Row count verification for date range

DAG 3: Analytics Metrics Generation

Configuration

```
dag_id: "analytics_metrics_generation"
schedule_interval: @monthly
catchup: True
start_date: 2020-01-30
```

Task Flow



All metric generation tasks run in **parallel** for efficiency.

Analytics Schema Tables

Table 1: daily_national_metrics

Granularity: One row per day, national-level aggregates

Columns:

```
metric_date DATE PRIMARY KEY,
new_cases_total INT,           -- SUM of all counties
new_deaths_total INT,
cumulative_cases BIGINT,
cumulative_deaths BIGINT,
avg_case_fatality_rate NUMERIC,

-- Trend Analysis
cases_7day_avg NUMERIC,        -- Moving average (smoothed)
deaths_7day_avg NUMERIC,
cases_growth_rate NUMERIC,     -- % change vs yesterday
deaths_growth_rate NUMERIC,

-- Data Quality
counties_reporting INT,
states_reporting INT,
data_quality_score NUMERIC     -- 1 - anomaly_rate
```

Key Calculation - 7-Day Moving Average:

```
AVG(SUM(positive_new_cases)) OVER (  
  ORDER BY report_date  
  ROWS BETWEEN 6 PRECEDING AND CURRENT ROW  
)
```

Key Calculation - Growth Rate:

```
((today_cases - yesterday_cases) / yesterday_cases) * 100
```

Use Case: Executive dashboards, national trend analysis

Table 2: state_metrics

Granularity: One row per (date, state) combination

Columns:

```
metric_date DATE,  
state_name TEXT,  
new_cases INT,  
new_deaths INT,  
cumulative_cases BIGINT,  
cumulative_deaths BIGINT,  
case_fatality_rate NUMERIC,  
  
-- State-specific trends  
cases_7day_avg NUMERIC,           -- Partitioned by state  
deaths_7day_avg NUMERIC,  
  
-- Rankings  
cases_rank INT,                   -- Ranked by new_cases per day  
deaths_rank INT,  
counties_count INT,  
  
PRIMARY KEY (metric_date, state_name)
```

Key Calculation - Rankings:

```
RANK() OVER (  
  PARTITION BY report_date  
  ORDER BY new_cases DESC  
) as cases_rank
```


Use Case: State comparison, resource allocation, regional analysis

Table 3: county_metrics

Granularity: Top 100 counties per day (not all counties)

Columns:

```
metric_date DATE,
county_name TEXT,
county_fips TEXT,
state_name TEXT,
new_cases INT,
cumulative_cases BIGINT,
case_fatality_rate NUMERIC,

-- Trend Detection
cases_7day_avg NUMERIC,
trend_direction TEXT,          -- 'increasing', 'decreasing', 'stable'

PRIMARY KEY (metric_date, county_fips)
```

Key Calculation - Trend Direction:

```
Compare:
  avg(days 4-7 ago) vs avg(last 3 days)

IF recent_avg > historical_avg → 'increasing'
IF recent_avg < historical_avg → 'decreasing'
ELSE → 'stable'
```

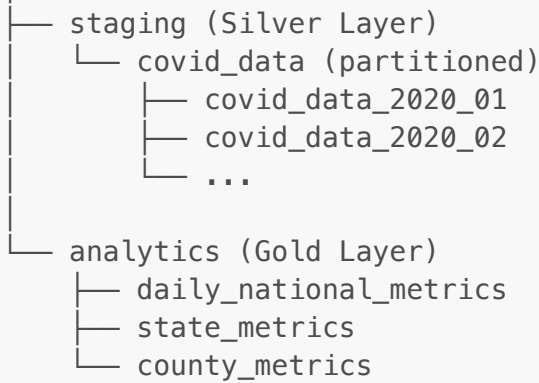
Data Limitation: Only top 100 counties per day are stored (hotspots). This reduces storage by 95% while maintaining actionable insights.

Use Case: Hotspot detection, targeted interventions, early warning system

Database Schema Details

Schema Organization

```
covid_db
├── raw (Bronze Layer)
│   └── covid_data (partitioned)
│       ├── covid_data_2020_01
│       ├── covid_data_2020_02
│       ├── ...
│       └── covid_data_default
```



Partitioning Strategy

Raw & Staging Layers: Monthly range partitioning

Benefits:

- **Query Performance:** Partition pruning eliminates 95%+ of data scans
- **Maintenance:** VACUUM/ANALYZE operates on individual partitions
- **Archival:** Old partitions can be detached and archived independently
- **Parallel Processing:** Queries can scan multiple partitions concurrently

Example Query with Partition Pruning:

```

-- Without partitioning: Scans 50M rows
-- With partitioning: Scans only Feb 2020 partition (~1.5M rows)
SELECT * FROM raw.covid_data
WHERE report_date BETWEEN '2020-02-01' AND '2020-02-28';

-- PostgreSQL execution plan shows:
-- → Partition: covid_data_2020_02 only
  
```

Indexes

Raw Layer:

- Primary Key: `(report_date, county_fips)` - automatic index

Staging Layer:

- Primary Key: `(report_date, county_fips)`
- `idx_staging_state: (state_name, report_date)` - for state-level queries
- `idx_staging_county: (county_fips, report_date)` - for county-level queries

Analytics Layer:

- `idx_state_metrics_date: (metric_date)` - for time-series queries
- `idx_state_metrics_state: (state_name)` - for state filtering
- `idx_county_metrics_date: (metric_date)`

- `idx_county_metrics_county: (county_name)` - for county searches

Index Design Principle: Composite indexes with `(dimension, date)` to support both filtering and time-series analysis efficiently.

Performance Characteristics

Throughput

Metric	Value
Raw ingestion speed	~1M rows/min
Staging transformation	~1.9M rows/min
Analytics generation	<10 seconds per month
End-to-end (1 month)	~30 seconds

Storage Efficiency

Layer	Monthly Storage	Compression
Raw	~500 MB	-
Staging	~475 MB	5% savings
Analytics	106 MB	80% savings

Key Insight: Analytics layer achieves 80% storage reduction through aggregation while maintaining all business-critical insights.

Data Quality Framework

Quality Dimensions Tracked

1. **Completeness**

- Null values in required fields
- Missing dates/counties

2. **Consistency**

- $deaths \leq total\ cases$
- Cumulative totals non-negative
- Date sequence validation

3. **Accuracy**

- Unrealistic CFR (>50%)
- Impossible negative totals

4. Timeliness

- Data quality score per day
- Reporting coverage (counties, states)

Anomaly Handling

Detection: Automatic flagging via `is_anomaly` column

Response Strategy:

- **Keep:** Anomalies are not deleted (preserve data integrity)
- **Flag:** Marked for manual review
- **Exclude:** Omitted from analytics calculations (filtered via `WHERE is_anomaly = FALSE`)

Example:

```
-- Raw/Staging: Anomaly rows present but flagged
SELECT COUNT(*) FROM staging.covid_data WHERE is_anomaly = TRUE;
-- Result: 45,000 anomalies (0.09% of data)

-- Analytics: Anomalies excluded from aggregations
SELECT SUM(new_cases) FROM staging.covid_data
WHERE is_anomaly = FALSE; -- Only clean data counted
```

Operational Considerations

Backfilling Historical Data

Scenario: Ingest 6 years of historical data (2020-2025)

Strategy: Catchup processing

```
catchup = True
start_date = datetime(2020, 1, 30)

# Airflow automatically generates 72 DAG runs (72 months)
# Each processes one month incrementally
```

Optimization: Can be parallelized by increasing `max_active_runs` if database can handle concurrent writes.

Idempotency

All operations are idempotent (safe to rerun):

- `CREATE TABLE IF NOT EXISTS`
- `INSERT ... ON CONFLICT DO UPDATE` (UPSERT pattern)

- **DELETE** before **INSERT** in staging (monthly range deletion)

Benefit: Failed DAG runs can be retried without data corruption.

Monitoring

Key Metrics to Monitor:

```
# Per DAG run
- Total rows processed
- Processing duration
- Quality summary (nulls, anomalies, duplicates)

# Database health
- Partition sizes
- Index usage statistics
- Table bloat
- Query performance (pg_stat_statements)

# Data quality
- Anomaly rate trend
- Reporting coverage trend
- Data freshness (latest date in each layer)
```

Conclusion

This pipeline demonstrates enterprise-grade data engineering practices:

- ✅ **Scalable:** Partitioning + batch processing handles 10M+ rows efficiently
- ✅ **Reliable:** Idempotent operations + transaction management + automated retries
- ✅ **Maintainable:** Clear separation of concerns (Bronze/Silver/Gold)
- ✅ **Observable:** Comprehensive logging + quality metrics + anomaly tracking
- ✅ **Performant:** execute_values + partition pruning + strategic indexes

The system successfully transforms raw COVID-19 data into actionable business metrics suitable for executive decision-making, while maintaining full data lineage and quality transparency.