

pemo

# Data Strategy Document - Detailed

by Mehmet Akif Kucuk

2025-06-12

Thursday

---

<b>In a nutshell;</b>	<b>3</b>
<b>Current Data Structure of pemo</b>	<b>4</b>
<b>Proposed High-Level Data Structure for pemo</b>	<b>5</b>
<b>Proposed Detailed Data Structure for pemo</b>	<b>6</b>
<b>1. Data Governance Foundation</b>	<b>8</b>
1.1. Data Inventory and Asset Management (Table-Level) - Data Catalog L1	8
1.2. Data Classification and Sensitivity Mapping (Column-Level) Data Catalog L2	8
<b>2. Compliance and Regional Requirements</b>	<b>9</b>
2.1. Data Residency Policy	9
2.2. Data Retention and Lifecycle Management	9
<b>3. Access Control Security</b>	<b>10</b>
3.1 Role-Based Access Management	10
<b>4. Data Architecture and Performance</b>	<b>10</b>
4.1. Database Optimization	10
4.2. Data Pipeline Architecture	11
4.3. Curated Data and User Access Optimization	11
<b>5. Data Visualization</b>	<b>11</b>
5.1. Data Visualization	11
<b>6. Data Quality and Monitoring</b>	<b>12</b>
6.1. Data Quality and Monitoring	12
<b>7. External Integrations and Partnerships</b>	<b>12</b>
7.1. Third-Party Data Sharing Controls	12
<b>8. Data Backup and Disaster Recovery</b>	<b>13</b>
<b>9. Monitoring Success Metrics and KPIs</b>	<b>14</b>

## In a nutshell;

**Risk Mitigation:** Eliminates potential millions in UAE Data Regulatory/SAMA fines and prevents data breaches

**Revenue Growth:** Enables \$2-5M additional revenue through better data utilization and faster decision-making

**Cost Optimization:** Reduces data-related costs by 30-50% through precalculation, automation and efficiency

**Curated Data:** Maintains curated data by following medallion structure to drive faster, accurate and consistent decisions

**Competitive Advantage:** Provides real-time insights with <1 hour time-to-insight for critical business questions

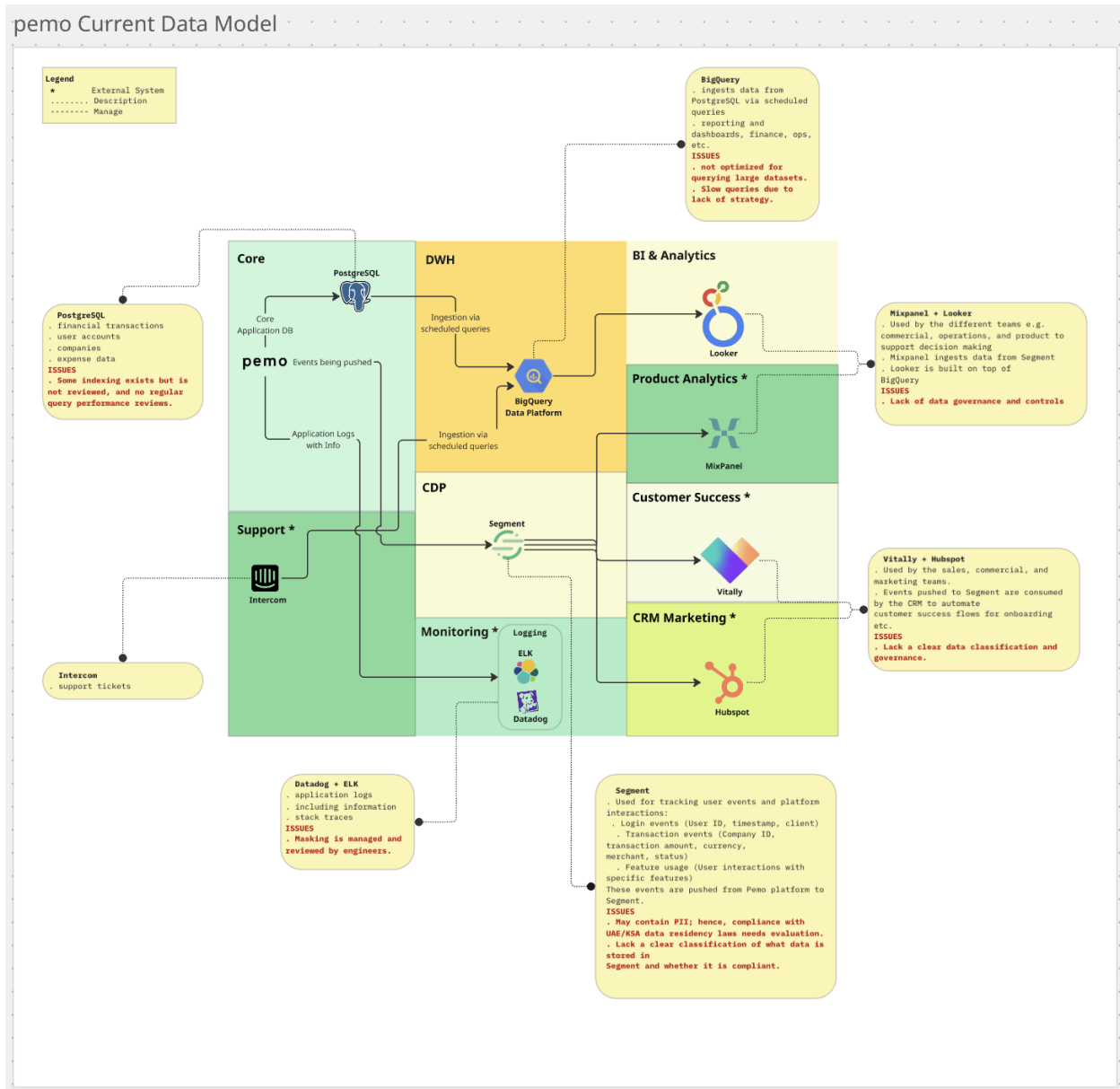
**Business Continuity:** Ensures <4 hour recovery time during disasters, protecting against extended downtime

**Customer Trust:** Demonstrates strong data protection, essential for customer retention and growth

# Current Data Structure of pemo

## pemo Current Data Model Frame

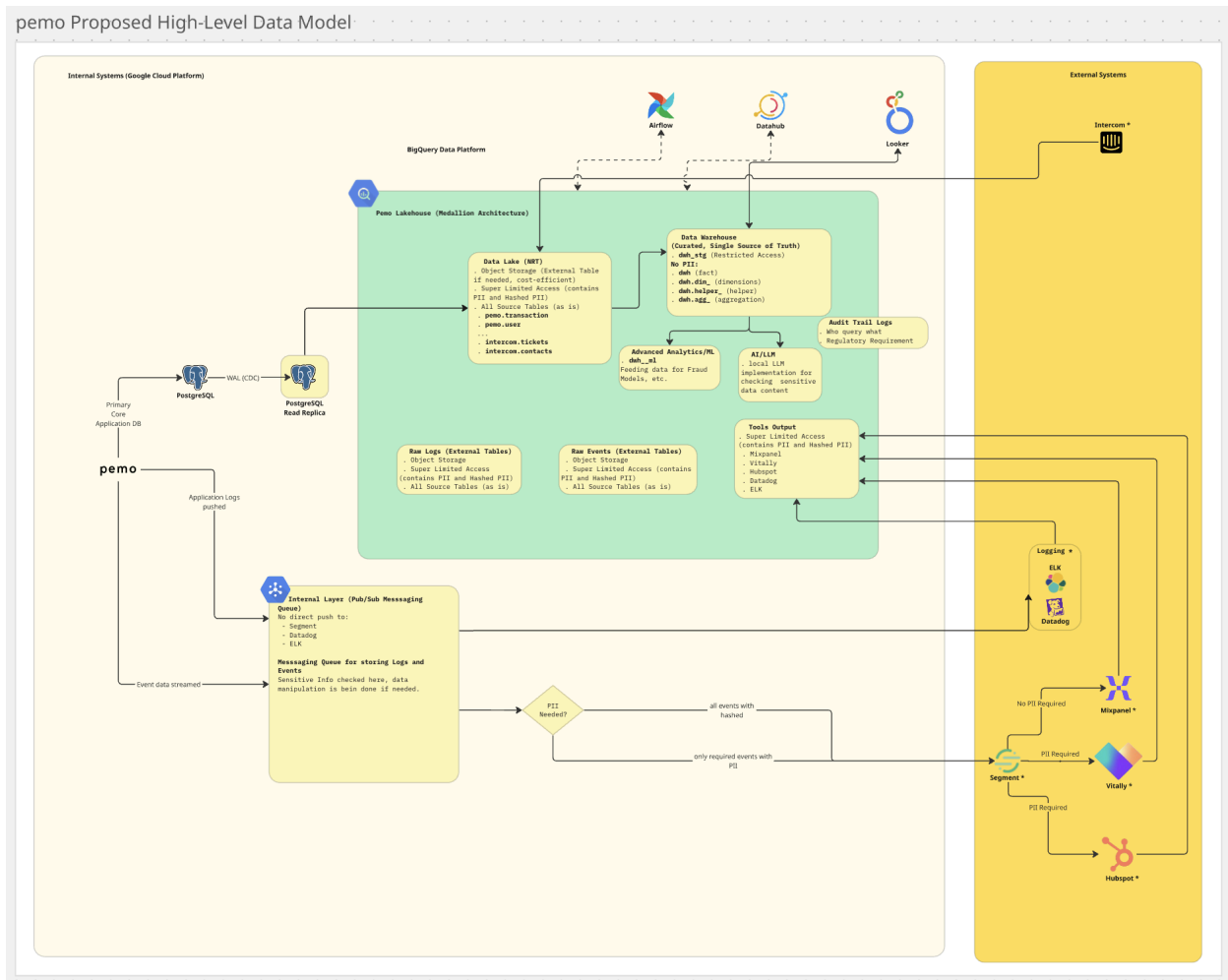
[https://miro.com/app/board/uXjVlr6Vhbw=/?share\\_link\\_id=264565662366](https://miro.com/app/board/uXjVlr6Vhbw=/?share_link_id=264565662366)



# Proposed High-Level Data Structure for pemo

## pemo Proposed High Level Data Model Frame

[https://miro.com/app/board/uXjVlr6Vhbw=/?share\\_link\\_id=264565662366](https://miro.com/app/board/uXjVlr6Vhbw=/?share_link_id=264565662366)

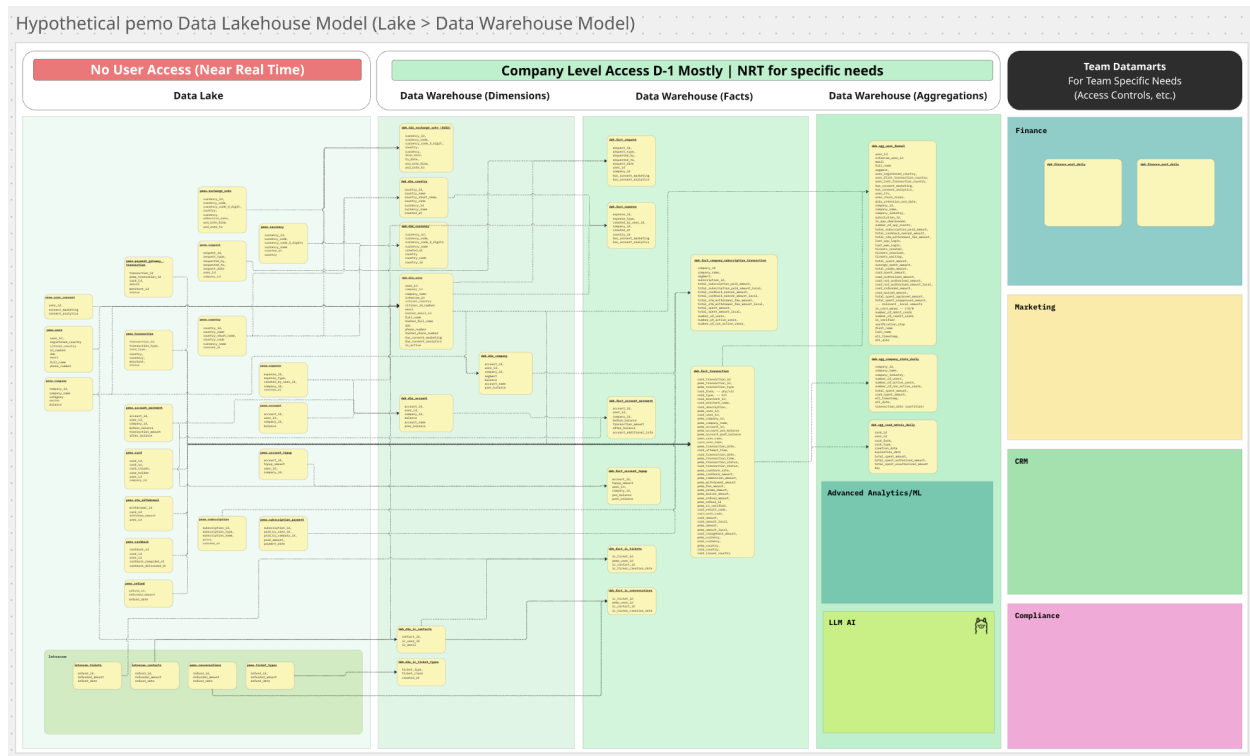




# Hypothetical pemo Data Lakehouse Model

## *Hypothetical pemo Data Lakehouse Model (Lake > Data Warehouse Model) Frame*

[https://miro.com/app/board/uXjVlr6Vhbw=?share\\_link\\_id=264565662366](https://miro.com/app/board/uXjVlr6Vhbw=?share_link_id=264565662366)



# 1. Data Governance Foundation

## 1.1. Data Inventory and Asset Management (Table-Level) - Data Catalog L1

**Objective:** Provide clear answers to "Where does data come from, where does it go, and how is it transferred?" This enables proactive risk identification.

Source System	Src Table/Stream	Target System	Target Table/Stream	Migration Method	Frequency	Data Residency	Target Retention	Notes
PostgreSQL	users	BigQuery	dataset.users	Scheduled Query	Daily	me-central2	5 years	ID No, email
Event Stream	payments	BigQuery	dataset.payments	Datastream	Saatlik	me-central2	3 years	Card no
Log System	app_logs	Datadog	-	Cloud Function	Real Time	-	90 days	user, ip
DWH	transaction, user company	DWH	agg_company_sta ts_daily	Transform ETL	NRT every 10m	me-central2	No Retention	Precalculated, very fast result

<b>Current State</b>	No central inventory and catalog exists
<b>Business Impact</b>	<ul style="list-style-type: none"><li>- Enables identification of risky data flows across all internal and external systems</li><li>- Simplifies compliance audits</li><li>- Protects company from fines and builds customer trust</li></ul>
<b>Cost Impact</b>	<ul style="list-style-type: none"><li>- Requires initial time and resource investment</li><li>- However, potential fines and reputation damage from data loss or breaches are significantly higher</li></ul>
<b>Technical Notes</b>	<i>Utilize dbt and SQL Mesh for source-target mapping</i>

## 1.2. Data Classification and Sensitivity Mapping (Column-Level) Data Catalog L2

**Objective:** Identify sensitive information within data and ensure adequate protection measures.

Table/Stream	Column Name	Description	Sensitivity Level	DLP Found	Manipulation	Data Residency	Retention	Notes
dataset.users	customer_id	ID Number	PII	Yes	Hash	europa-west1	5 years	Not hashed currently
dataset.events	card_number	Credit Card No	Confidential	Yes	Custom	europa-west1	5 years	Already custom manipulated
app_logs	ip_address	IP Address	Confidential	No	Masking	-	90 days	Already masked to Segment

<b>Current State</b>	No central column-level catalog
<b>Business Impact</b>	<ul style="list-style-type: none"><li>- Ensures compliance with UAE &amp; SAMA regulations</li><li>- Protects brand security</li><li>- Reduces potential lawsuit and penalty risks</li></ul>
<b>Cost Impact</b>	Late discovery of violations can result in very high penalties



<b>Technical Notes</b>	Google DLP, Tokenization, Masking, Hashing via Dataflow
------------------------	---

## 2. Compliance and Regional Requirements

### 2.1. Data Residency Policy

**Objective:** Ensure data remains within required geographical boundaries per local laws.

*Maintain and follow the rules defined here:*

Region	Data Types	Storage Location	Transfer Rules	Local Requirements
EU	All Personal Data	europa-west1	No US transfer	GDPR Article 44-49
UAE	Customer Data	me-central2	Local processing	UAE Data Protection Law
KSA	Financial Data	Local only	No cloud export	SAMA Cyber Security Framework

<b>Current State</b>	Core applications are aligned with data residency requirements
<b>Business Impact</b>	Prevents business shutdowns due to non-compliance
<b>Cost Impact</b>	Additional regional infrastructure costs
<b>Technical Notes</b>	Explore GCP's multi-region solutions with strong compliance alignment

### 2.2. Data Retention and Lifecycle Management

**Objective:** Define how long data should be stored and subsequent actions, considering Compliance as well.

*Check retention period in table and column level mapping (1.1, 1.2) and apply deletion:*

<b>Current State</b>	No central data retention system
<b>Business Impact</b>	<ul style="list-style-type: none"> <li>- Eliminates unnecessary data burden</li> <li>- Ensures compliance (e.g., "emails must be stored for 5 years")</li> <li>- Reduces risk of old data breaches</li> </ul>
<b>Cost Impact</b>	<ul style="list-style-type: none"> <li>- Archived data reduces storage costs</li> <li>- Manual SQL deletion is risky → Automation essential</li> </ul>
<b>Technical Notes</b>	<ul style="list-style-type: none"> <li>- BigQuery TTL (Time-to-Live) configuration</li> <li>- "Soft Delete + Retention Policy" → archive first, permanent delete later</li> </ul>

## 3. Access Control Security

### 3.1 Role-Based Access Management

**Objective:** Control who accesses which data at what level, minimizing unnecessary access.

Role	Dataset	Access Level	Notes
Compliance Officer	dwh.raw_events	Full Access	Original PII visible
Data Engineer	dwh.raw_events	Masked	Masked PII visible
Analyst	dwh.raw_events	Anonymized	Unencodable Data Visible
UAE Analyst	dwh.raw_events	UAE-Only	Regional Records Only
Marketing Analyst	segment.event_metrics	Marketing	Role to show who has access to external systems as well

<b>Current State</b>	Poor level of access controls
<b>Business Impact</b>	- Only necessary personnel access data → Enhanced security - Simplified audit processes
<b>Cost Impact</b>	Access confusion leads to data breach risks
<b>Technical Notes</b>	- BigQuery IAM Management with on-the-fly hashing (5-10% CPU cost) - Column and Row Level Access Controls

## 4. Data Architecture and Performance

### 4.1. Database Optimization

**Objective:** Achieve faster, accurate results at lower cost through precalculation strategies.

<b>Current State</b>	Unoptimized databases, queries
<b>Business Impact</b>	Quick, accurate, and consistent data access
<b>Cost Impact</b>	Increased costs due to precalculation processes and storage
<b>Technical Notes</b>	Partitioning, indexing frequently used columns, materialized views, clustering, and merge operations

## 4.2. Data Pipeline Architecture

**Objective:** Improve current data transfer methods using Near Real-Time ingestion with Change Data Capture technology.

Near Real Time Data ingestion by using the power of Change Data Capture technology

<b>Current State</b>	No real time ingestion flow
<b>Business Impact</b>	<ul style="list-style-type: none"><li>- Current and real-time data → Better decision making</li><li>- Enhanced performance and reporting quality</li></ul>
<b>Cost Impact</b>	<ul style="list-style-type: none"><li>- Near real-time structures are more expensive and complex</li><li>- However, delayed data also delays decisions</li></ul>
<b>Technical Notes</b>	<ul style="list-style-type: none"><li>- CDC → EventStream → BigQuery (external or internal tables)</li><li>- GCP Dataflow for real-time pipeline setup</li><li>- Asynchronous flow separation as scale increases</li></ul>

## 4.3. Curated Data and User Access Optimization

**Objective:** Enable end users to access data quickly and accurately.

Warehouse layer objects may run frequently based on time criticality (costly but necessary)

No user access to initial ingestion layer; access limited to Data Warehouse layer and beyond

<b>Current State</b>	No data management in Data Warehouse level
<b>Business Impact</b>	<ul style="list-style-type: none"><li>- Golden-level curated data is being flagged and accessed</li><li>- Enhanced performance and reporting quality</li></ul>
<b>Cost Impact</b>	Some near-real-time data causes cost but also brings faster decision making
<b>Technical Notes</b>	dbt → BigQuery Transformation Management

## 4.4. Data Generation and Push for Advanced Analytics/Machine Learning

**Objective:** Empower fraud detection systems with rapid and accurate data access for real-time analytics and machine learning model training.

Warehouse layer objects are scheduled to run frequently to support near-real-time fraud detection

<b>Current State</b>	Limited data orchestration at the Data Warehouse level for fraud-specific use cases.
<b>Business Impact</b>	<ul style="list-style-type: none"><li>- Curated, high-quality fraud detection data enables faster identification of suspicious activities.</li><li>- Improved model accuracy enhances fraud prevention, reducing financial losses.</li></ul>

<b>Cost Impact</b>	Investment in optimized pipelines may yield long-term cost efficiencies.
<b>Technical Notes</b>	Data pipelines feed structured data into machine learning models (e.g., anomaly detection algorithms).

## 4.5. LLM Implementations

**Objective:** Leverage Large Language Models (LLMs) to securely explore and extract insights from sensitive data while maintaining compliance and data privacy.

Warehouse layer objects are designed to anonymize and preprocess sensitive data before feeding into LLMs, with strict access controls.

<b>Current State</b>	Limited automation for anonymization and preprocessing of sensitive datasets.
<b>Business Impact</b>	Enhanced data privacy measures build trust with stakeholders and ensure regulatory compliance.
<b>Cost Impact</b>	Processing sensitive data with LLMs incurs higher computational costs due to anonymization and encryption requirements.
<b>Technical Notes</b>	Use batch processing for non-time-critical tasks to balance cost and performance.

## 5. Data Visualization

### 5.1. Data Visualization

**Objective:** Ensure reporting and visualization uses curated data

Curated sources will be mapped in Data Catalog to be used by visualization tool, this layer will be clean, heavy aggregations to be tracked

<b>Current State</b>	Querying BigQuery via Looker
<b>Business Impact</b>	Access beautiful insights with a curated data on the back-end
<b>Cost Impact</b>	Cost will be reduced after proper reporting catalog
<b>Technical Notes</b>	Looker might connect other layers, OLAP layer if needed

## 6. Data Quality and Monitoring

### 6.1. Data Quality and Monitoring

**Objective:** Ensure data accuracy, completeness, and consistency across all systems.

Quality Dimension	Measure	Target	Monitoring	Alert Threshold
Completeness	% of non-null values	>95%	Daily	<90%
Timeliness	Data freshness	<15mins	Continuous	<30mins
Consistency	Source Target Match	100%	Daily	>99%

<b>Current State</b>	No data quality checks and monitoring being implemented
<b>Business Impact</b>	Reduce time spent fixing data issues
<b>Cost Impact</b>	Minimal implementation cost for checks
<b>Technical Notes</b>	Python, Spark Framework or SQL with dashboard visualization

## 7. External Integrations and Partnerships

### 7.1. Third-Party Data Sharing Controls

**Objective:** Monitor and control what data goes to external systems.

External System	Data Shared	Sharing ID	Data Content Controlled by	Sharing Approved by	Shared Time
segment	User Events	share_sgm_01	Mehmet Kucuk	Compliance Officer	2025-09-01 12:30
datadog	Application Log	share_dd_01	Akif Kucuk	Engineering Lead	2025-09-01 13:30

- **Additional Internal layer** before sending, check whether the content has PII by Sensitive Data Exploration, content control person recorded, an, no direct push to external systems

<b>Current State</b>	No sensitive data classification checks, especially for logs and events
<b>Business Impact</b>	<ul style="list-style-type: none"> <li>- Prevents PII data from going to wrong destinations</li> <li>- Provides audit comfort</li> </ul>
<b>Cost Impact</b>	<ul style="list-style-type: none"> <li>- Internal layers can be costly</li> <li>- However, sensitive data exposure creates reputation and financial loss</li> </ul>
<b>Technical Notes</b>	<ul style="list-style-type: none"> <li>- Internal Layer for Logs and Events</li> <li>- Log flow: Cloud Logging → Pub/Sub → Dataflow → Transformed Pub/Sub → Parquet → BigQuery → External Tools</li> <li>- Event flow: Cloud Function → Internal Endpoint → Pub/Sub → Dataflow → Hashed Pub/Sub → BigQuery</li> <li>- Log all external data transfers with timestamps</li> <li>- Use only anonymous IDs (hashed + salted) for external data</li> <li>- Store external data temporarily for audit queries</li> </ul>

## 8. Data Backup and Disaster Recovery

**Objective:** Ensure critical data is securely backed up and quickly restored during failures, while meeting compliance requirements.

System	Asset	Backup Frequency	Storage Location	Recovery Time Obj (RTO)	Compliance Notes
PostgreSQL	pemo.user	Daily full, hourly incremental	me-central2, europe-west1	<4h	SAMA, GDPR
bigquery	dwh.transaction	Daily snapshot	share_dd_01	<24h	UAE Data Protection, SAMA

<b>Current State</b>	No centralized backup or disaster recovery plan documented
<b>Business Impact</b>	<ul style="list-style-type: none"> <li>- Ensures business continuity during outages, minimizing downtime and financial loss</li> <li>- Builds customer trust by protecting sensitive data</li> </ul>
<b>Cost Impact</b>	<ul style="list-style-type: none"> <li>- Initial setup and ongoing storage costs</li> <li>- Avoids much higher costs from data loss, regulatory fines, or reputation damage</li> </ul>
<b>Technical Notes</b>	Follow GCP best practices including Cloud SQL Automated Backups and BigQuery Auto Snapshots

## 9. Monitoring Success Metrics and KPIs

<b>Data Access Metrics</b>	<ul style="list-style-type: none"><li>- Who accessed which tables</li><li>- Mostly queried tables</li><li>- Mostly querying user</li><li>- Who tried to access restricted tables systems</li></ul>
<b>Operational Metrics</b>	<ul style="list-style-type: none"><li>- Data Pipeline Uptime: &gt;99.5%</li><li>- Query Performance: &lt;3 second average response</li><li>- Data Quality Score: &gt;95%</li><li>- Compliance Audit Results: Zero findings</li></ul>
<b>Business Metrics</b>	<ul style="list-style-type: none"><li>- Time to Insight: &lt;1 hour for critical business questions</li><li>- Decision Making Speed: 50% faster data-driven decisions</li><li>- Revenue Impact: \$2-5M additional revenue from better data usage</li><li>- Cost Optimization: 30-50% reduction in data-related costs</li></ul>
<b>Risk Metrics</b>	<ul style="list-style-type: none"><li>- Security Incidents: Zero data breaches</li><li>- Compliance Violations: Zero regulatory fines</li><li>- Data Recovery Time: &lt;4 hours for critical data</li><li>- Vendor SLA Compliance: &gt;99% across all vendors</li></ul>