

p e m o

Data Strategy Document - Non Tech

by Mehmet Akif Kucuk

2025-06-12

Thursday

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

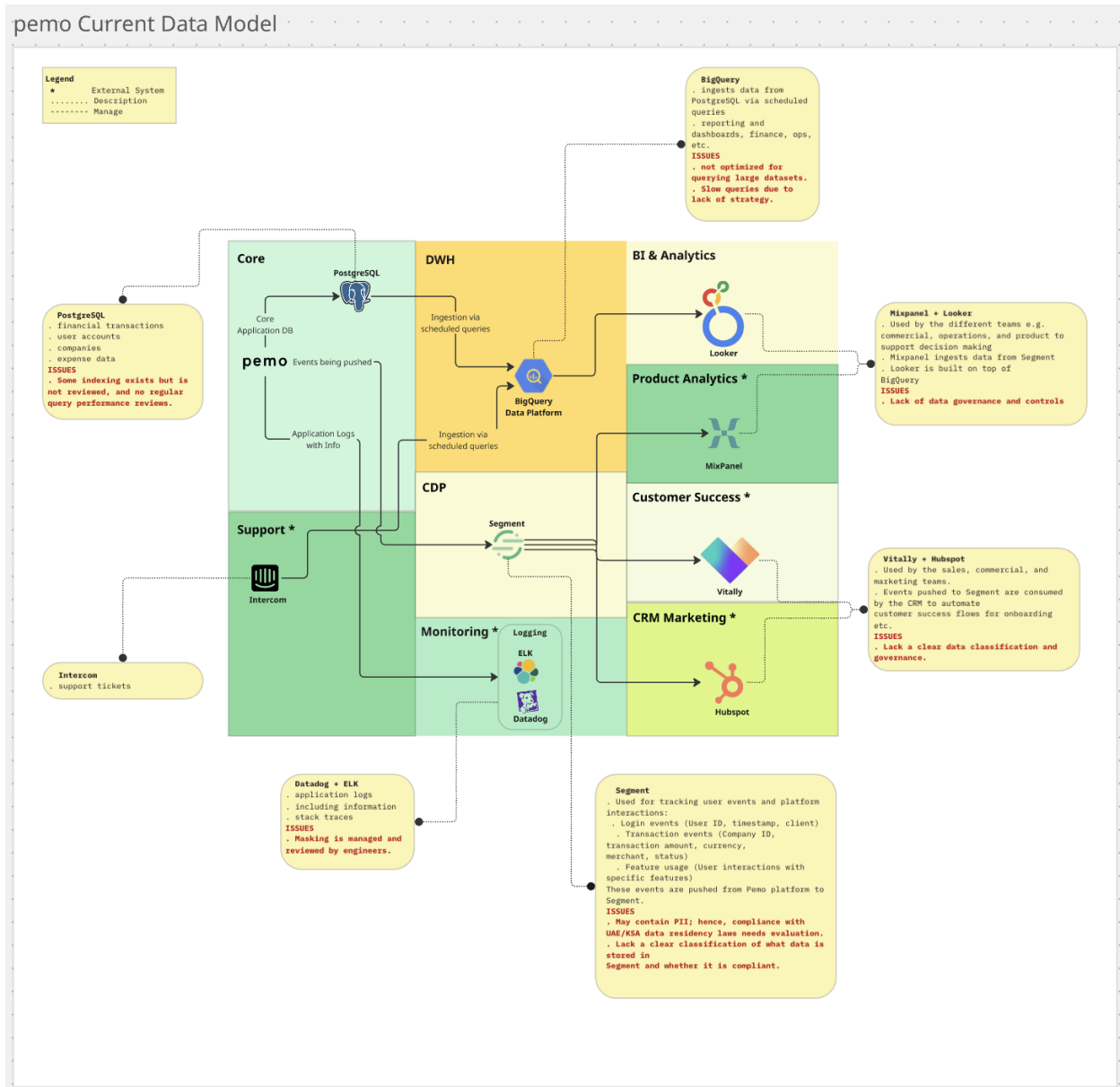
Business Continuity: Ensures <4 hour recovery time during disasters, protecting against extended downtime

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

Current Data Structure of pemo

pemo Current Data Model Frame

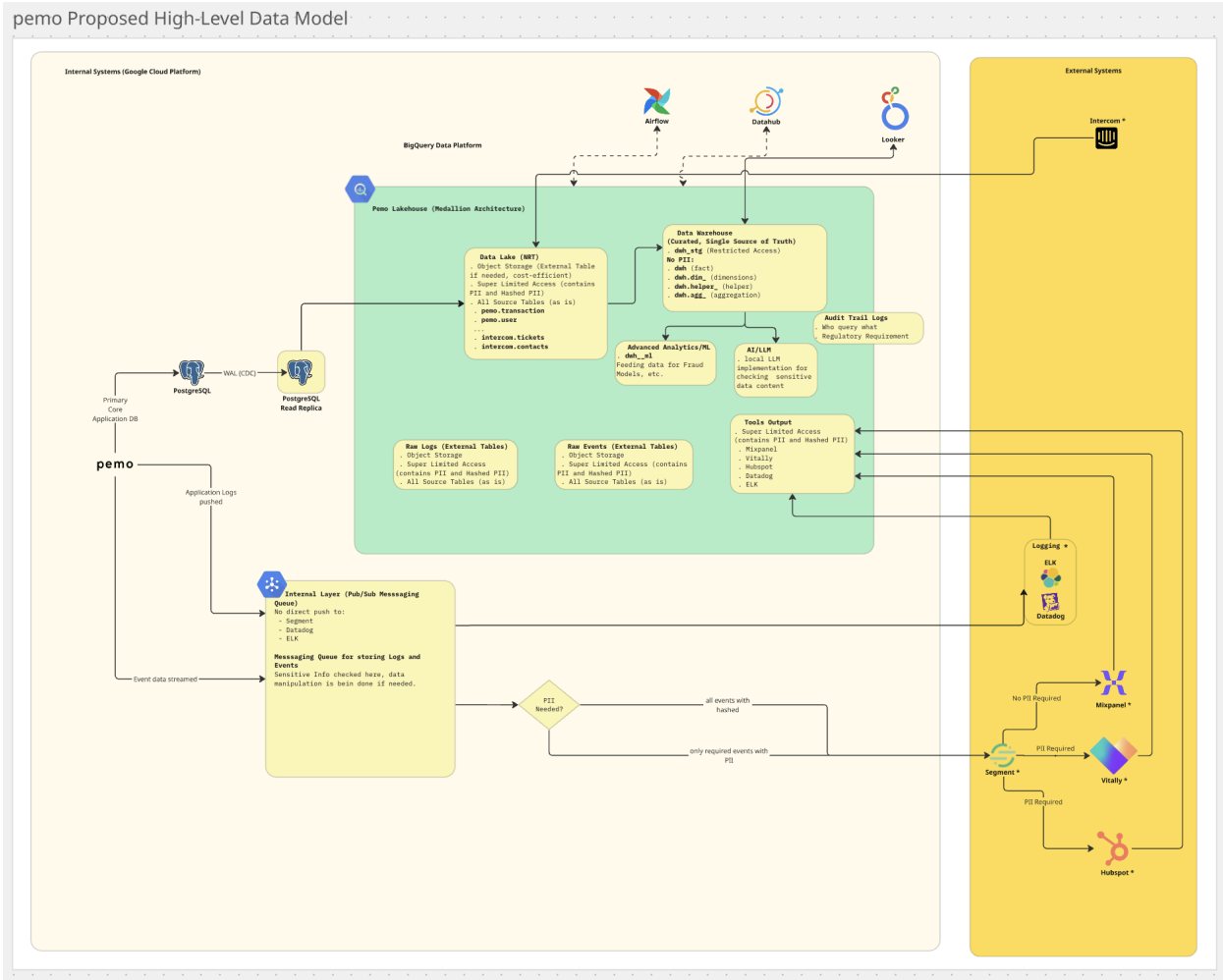
https://miro.com/app/board/uXjVlr6Vhbw=/?share_link_id=264565662366



Proposed Data Structure for pemo

pemo Proposed High Level Data Model Frame

https://miro.com/app/board/uXjVlr6Vhbw=?share_link_id=264565662366

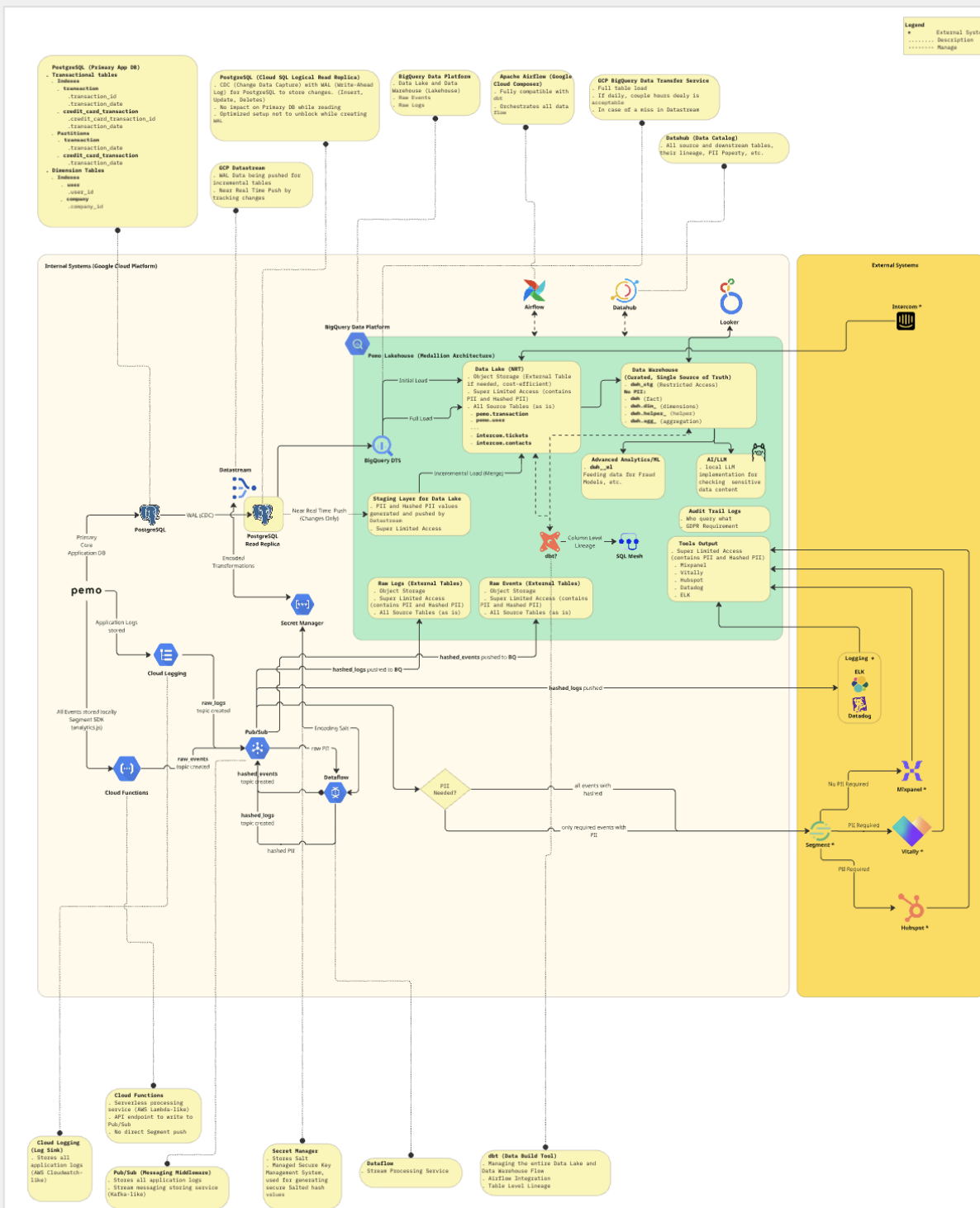


Proposed Data Structure for pemo

pemo Proposed Detailed Data Model Frame

https://miro.com/app/board/uXjVIr6Vhbww/?share_link_id=264565662366

pemo Proposed Detailed Data Model



1. Data Governance Foundation

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

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

Current State	No central inventory and catalog exists
Business Impact	<ul style="list-style-type: none">- Enables identification of risky data flows across all internal and external systems- Simplifies compliance audits- Protects company from fines and builds customer trust

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

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

Current State	No central column-level catalog
Business Impact	<ul style="list-style-type: none">- Ensures compliance with UAE & SAMA regulations- Protects brand security- Reduces potential lawsuit and penalty risks

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	europe-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

Current State	Core applications are aligned with data residency requirements
Business Impact	Prevents business shutdowns due to non-compliance

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:

Current State	No central data retention system
Business Impact	<ul style="list-style-type: none">- Eliminates unnecessary data burden- Ensures compliance (e.g., "emails must be stored for 5 years")- Reduces risk of old data breaches

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

Current State	Poor level of access controls
Business Impact	<ul style="list-style-type: none">- Only necessary personnel access data → Enhanced security- Simplified audit processes

4. Data Architecture and Performance

4.1. Database Optimization

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

Current State	Unoptimized databases, queries
Business Impact	Quick, accurate, and consistent data access with less cost

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

Current State	No real time ingestion flow
Business Impact	- Current and real-time data → Better decision making - Enhanced performance and reporting quality

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

Current State	No data management in Data Warehouse level
Business Impact	- Golden-level curated data is being flagged and accessed - Enhanced performance and reporting quality

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

Current State	Limited data orchestration at the Data Warehouse level for fraud-specific use cases.
Business Impact	- Curated, high-quality fraud detection data enables faster identification of suspicious activities. - Improved model accuracy enhances fraud prevention, reducing financial losses.

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.

Current State	Limited automation for anonymization and preprocessing of sensitive datasets.
Business Impact	Enhanced data privacy measures build trust with stakeholders and ensure regulatory compliance.

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

Current State	Querying BigQuery via Looker
Business Impact	Access beautiful insights with a curated data on the back-end

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
Timeliness	Data freshness	<15mins	Continuous	<30mins
Consistency	Source Target Match	100%	Daily	>99%

Current State	No data quality checks and monitoring being implemented
Business Impact	Reduce time spent fixing data issues

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

Current State	No sensitive data classification checks, especially for logs and events
Business Impact	<ul style="list-style-type: none"> - Prevents PII data from going to wrong destinations - Provides audit comfort

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

Current State	No centralized backup or disaster recovery plan documented
Business Impact	<ul style="list-style-type: none">- Ensures business continuity during outages, minimizing downtime and financial loss- Builds customer trust by protecting sensitive data

9. Monitoring Success Metrics and KPIs

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

