

Data at Marfeel

Addressing complexity at scale with the latest technologies

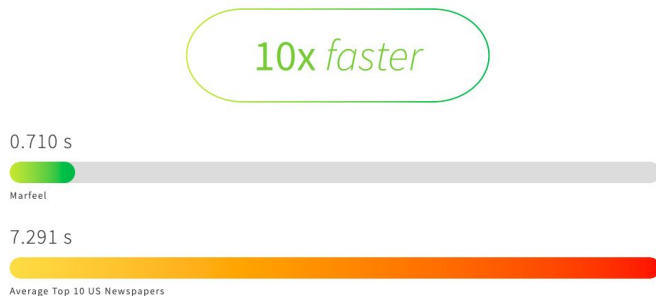
Alessandro Pregnolato
Head of Data



What does **Marfeel** do?



Optimize.
Engage.
Monetize.

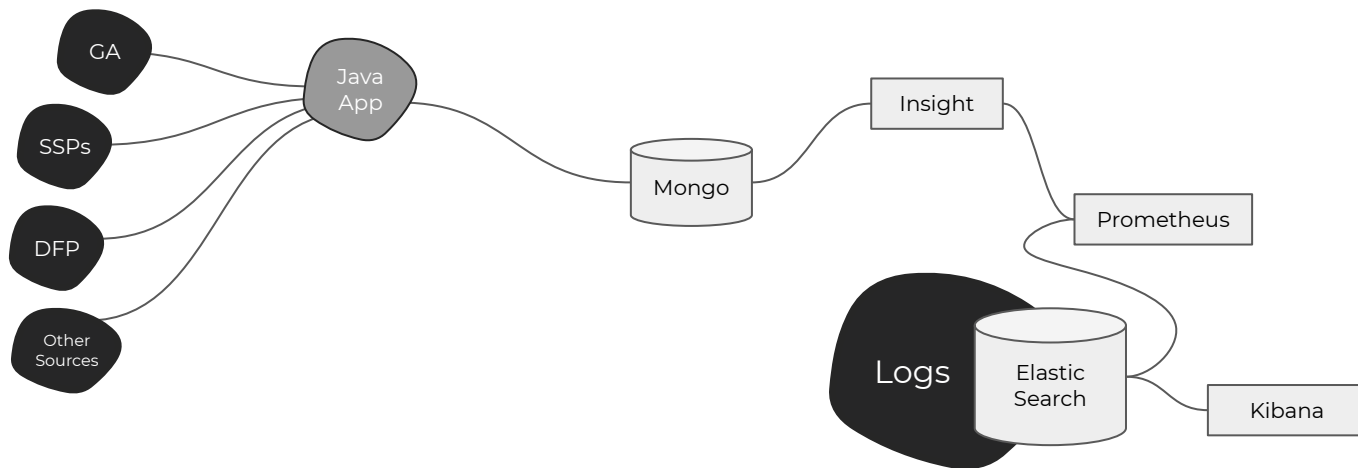


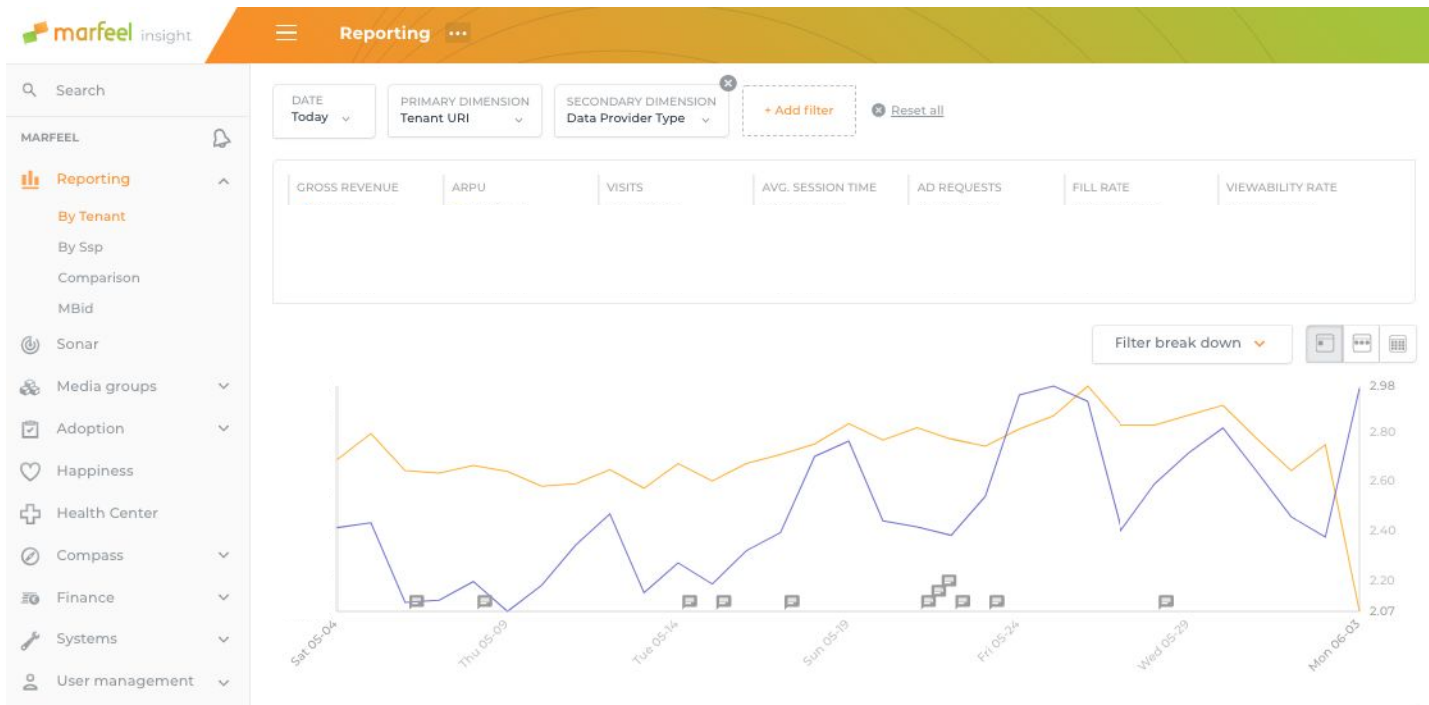
- A lightning fast, responsive mobile web
- A sophisticated monetization layer

...& **more** to come

Tonnes of data, each month.

- **700M** Visits
- **1.75Bn** Page views
- **4 Bn** Ad requests
- **20Tb** of logs data





Monitoring

SYS_nginx_consumer_tenant_5XX_errors_tooMany (11 active)	
SYS_nginx_consumer_tenant_stale_item_tooMany (21 active)	
SYS_nginx_consumer_tenant_stale_section_home_tooMany (39 active)	
SYS_nginx_consumer_tenant_stale_section_tooMany (15 active)	
DA_cluster_requests_outsideClusters_last4hours (6 active)	
GTB_JS_MSG (2 active)	
HEIS_NOTIFICATIONS_subscribers_not_increasing (2 active)	
MPL_CTR_drop_yday_vs_lastWeek (1 active)	
MPL_collector_down (1 active)	
MPL_traffic_not_marfeelized (5 active)	
ALOT_mblddr_timeout_above20Percent_last1hour (0 active)	
CDN_Mobile_First_Detector_Tenant_Above_Transition_Threshold_last_week (0 active)	
CDN_Mobile_First_Detector_Tenant_Below_Threshold_last_week (0 active)	
DA_adExchange_executors_enabledForinactiveAdExchange (0 active)	
DA_adExchange_queue_failedActions_isAbove1Percent_last1hour (0 active)	

Tenant m.eldiariodechihuahua.mx has ads.txt but doesn't contain Marfeel correct data

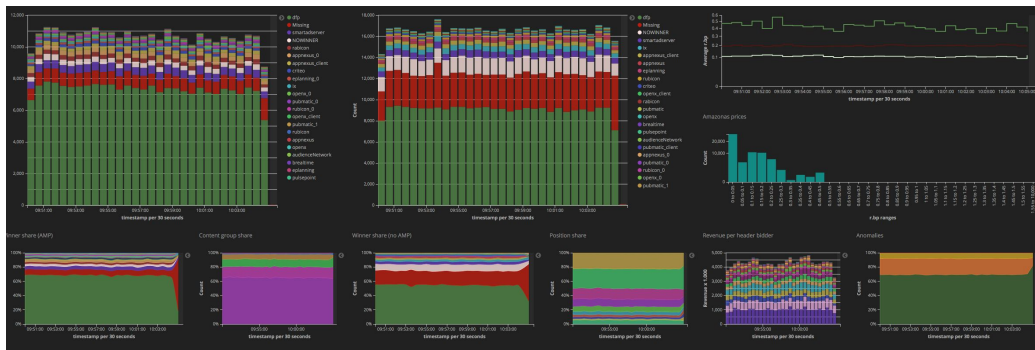
① GSA-1066

Tenant economiahoy.mx has ads.txt but doesn't contain Marfeel correct data

① GSA-1063

Tenant prima.fr has ads.txt but doesn't contain Marfeel correct data

① GSA-1062



Some great achievements so far...



A single source of truth



Monitoring & alerts on most KPIs



A data-driven culture (to some extent)

"We don't know much about our tenants"

"I cannot count articles published per day"

"We could segment by tenants' attributes such as vertical, content type (news/evergreen), keywords/tags, topics (ML), wordcount, images, video, etc."

"We could create audiences"

"Not enough flexibility"

"I'd like to create my own visualizations and dashboards"

"I need different granularities"

"I'd like to cross this with content and tenant data"

"Cannot compare tenants or YoY"

"Cannot export data"

"Activation and QBR reports are very limited"

"I cannot join collections, nor cross them with other data"

"A big proportion of this data is not being used"

"We cannot look at yearly trends because there's no historical"

"I cannot perform complex operations (weighted averages, running totals, etc.)"

"The tools are dictating the events modeling"

What's a Data-driven culture?

Five building blocks of a data-driven culture

"Having clean, high-quality data, from a central source, and with clear metadata, is ineffective if staff can't access it"

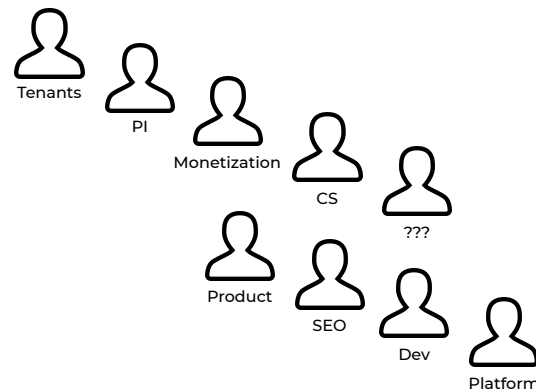
[Carl Anderson](#), [Michael Li](#)



1. Single source of truth
2. Data dictionary
3. Broad data access
4. Data Literacy
5. Data-driven decision making

... means NO BARRIERS

- Technical (DWH Modeling & SQL)
- Functional (Business knowledge)



- We can only report to our tenants ***Traffic Metrics*** consistent with their own data (*Google Analytics*)
- We can only report to our tenants ***Revenue Metrics*** from *SSPs & ADX*
- These source provide limited granularity
- Granular, accurate data requires access to paid tools (such as Google AdManager Premium, GA 360, etc.) whose cost is prohibitive at our scale



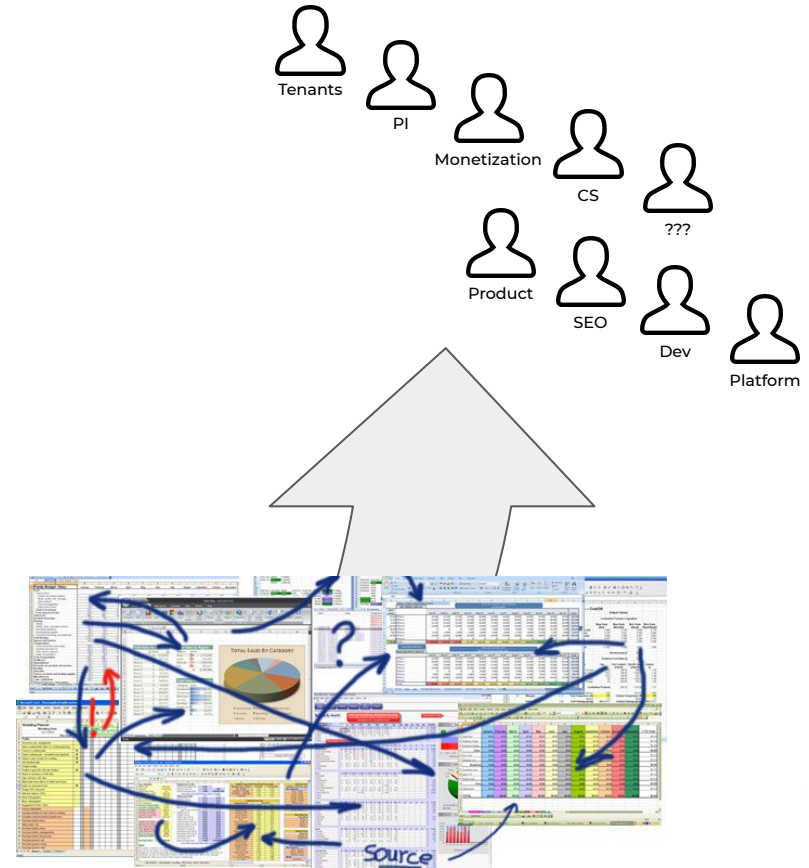
Implications

- **High-level, accurate data** → fit for *reporting* but not for *analytics*
- **Granular, approximated data** → fit for *analytics* but not for *reporting*
- Exploiting the available data currently requires such a degree of **technical and business knowledge** that's unreasonable to expect from our stakeholders

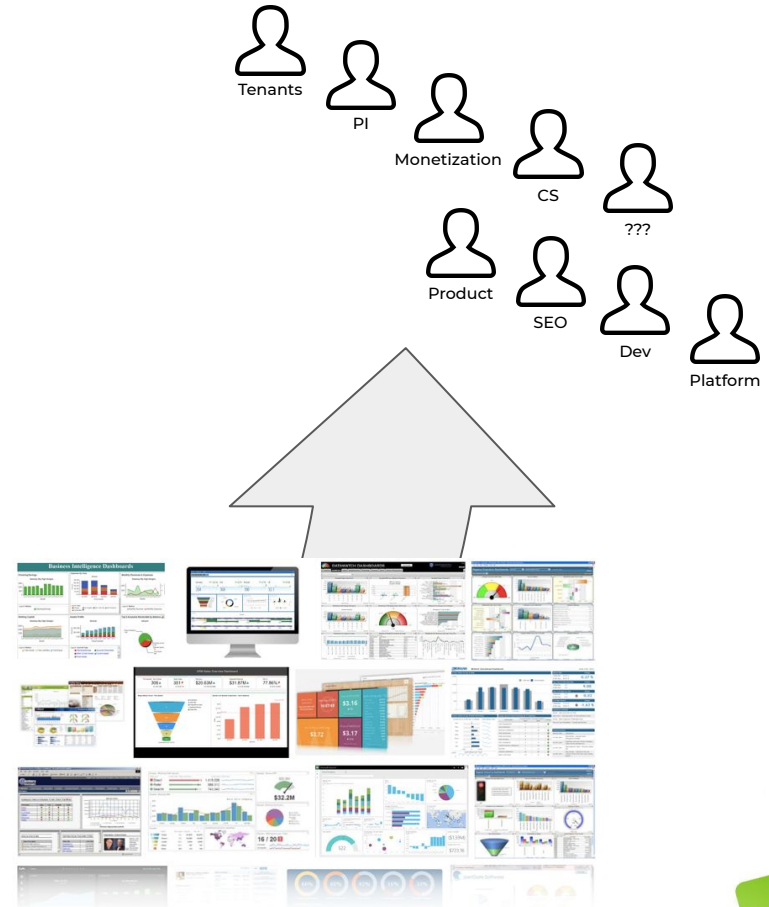


What to do?

Option #1 A troop of data Monkeys

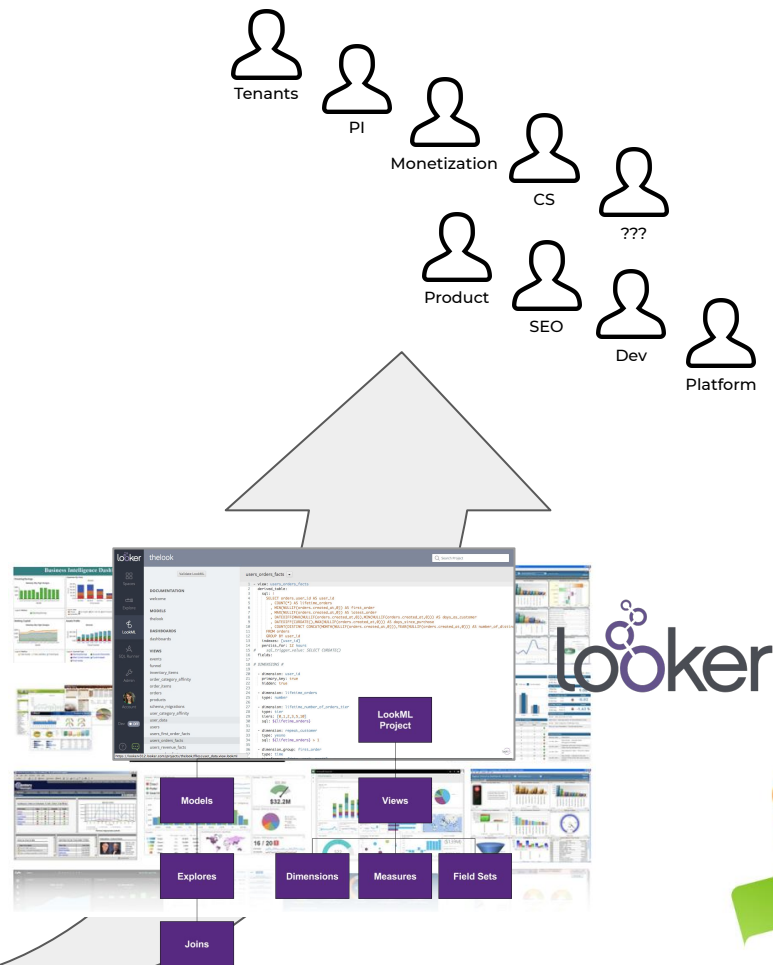


Option # 2 An army of BI developers



Option #3 Technology

- A logical layer in which to store table relationships and business rules
- Flexible access to the DWH data (writes SQL for you)



BI Layer - Tech Comparison

	Tableau	Qlikview	Looker
Latency	Low (in-memory)	Low (in-memory)	Mid (from DB)
Loading & Pre-processing	Required	Required	Not Required
Maintenance	Mid (if managed)	Mid (if managed)	Low (centralized)
Development & Deployment time	High (dashboards)	High (dashboards)	Low (data model only)
Logical Model	N	N	Y
SQL Engine	N	N	Y
ETL Layer	Y	Y	Not required
OLAP Layer	Y	Y	Limited
Visualization capability	High	Mid/High	Mid
Mobile Adaptiveness	Low	Mid	High
Learning Curve (Developers)	High	High	Mid
Learning Curve (Business Users)	High	High	Low
Flexibility	Low	Low	High
Price	Mid	Mid	Mid
Pros	Fast (In-memory) OLAP Layer Advanced Visualization	Fast (In-memory) OLAP Layer	True Self-Service Embeds business logic Only one model to build Restricts data interactions Unlimited scaling
Cons	Rigid Requires Dashboards Development Limited self-exploration Limited scaling		Not as powerful Slower (Relies mostly on DB)

Looker supported DB

- Amazon Aurora
- Amazon Redshift
- Apache Spark 1.5+
- Apache Spark 2.0
- Aster Data
- Clustrix
- Exasol
- Google BigQuery Legacy SQL
- Google BigQuery Standard SQL
- Google Cloud PostgreSQL
- Google Cloud SQL
- IBM Netezza
- MariaDB
- MemSQL
- Microsoft Azure PostgreSQL
- Microsoft Azure SQL Data Warehouse
- Microsoft Azure SQL Database
- Microsoft SQL Server 2005
- Microsoft SQL Server 2008+
- MySQL
- Oracle
- PostgreSQL
- PrestoDB
- Qubole Presto
- Qubole Presto Service
- SAP HANA
- Snowflake
- Teradata
- Vector
- Vertica 7.1+

(Discarded all non-distributed & high-end corporate solutions)

DWH Layer - Tech Comparison

	AWS Redshift	Snowflake	Google Big Query	Clickhouse
Speed	Mid/high	Mid/high	Mid/high	High (?)
Maintenance	Mid	Low	Mid	Mid/High
Dynamic resizing	Limited	Y	Y	N
Concurrency	Low	High	Mid	Mid
Indexes	Sort/Dist Key	Self-tuning	Self-tuning	Sort Key (primary only)
Real-time Ingestion	AWS Kinesis Firehose	Snowpipe	Y	Y
Complex Types	N	Json/XML	Nested Struct Types	Array
Join on Array/Nested DS	N (UDF?)	Y (JOIN on Json/XML)	Hive-like (Explode)	Hive-like (Explode)
Approximated Calculations	N	N	Y	Y
Transactions	Y	Y	Y	N
Replication	N	Y	Y	Y
Fault Resistance	Backup	Distributed (replication)	Distributed (replication)	Distributed (replication)
Subqueries	Y	Y	Y	N
Window Functions	Y	Y	Y	N
UDF (Python, JS, etc.)	Y	Y	Y	N
Connectivity	Extensive	Extensive	Extensive	JDBC/ODBC only
Tableau Connectivity	Y	Y	Y	JDBC/ODBC only
Looker Connectivity	Y	Y	Y	N
LogStash Output	S3 only	S3 only	Y	N
Google Analytics integration	N	N	Native	N
Cost	Mid	Mid/High	Mid/High	Low
Pros	Highly Tunable, on AWS, widely adopted, previous experience	Self-tuning, fully elastic, high-concurrency, Json/XML support, <u>cheap storage</u>	Fully managed, linear self-scaling, high-concurrency, Json/XML support, Logstash/GA integration	Very Fast. Open Source
Cons	Storage and computing are coupled Not great at handling concurrency	Speed? Price?	Pricing model	High maintenance Non-standard SQL,, No UDF, No Window Functions, No Looker connectivity



Pros

- AWS, established, widely adopted
- Highly Tunable
- It works

Concerns

- Storage and computing are coupled (Spectrum doesn't quite cut it)
- Not great at handling concurrency
- Didn't evolve much since 2013. Outdated (?)



Google BigQuery

Pros

- Google (strong relationship)
- Fully managed, linear self-scaling, high-concurrency
- Json/XML support, GA integration

Concerns

- Awkward pricing model - Pay per query (flat rates start from 10K per month)
- Quite Hadoop-like. More complex to use?



Pros

- Open Source
- Allegedly very fast
- Some prestigious adopters (E.G. CloudFlare)

Concerns

- High-maintenance (concerns about the *Total Cost of Ownership*)
- Even if it was cheaper, do we need to process SO much data?
- Is it worth the trade off?
 - Non-standard SQL
 - No Subqueries
 - No Analytic functions
 - No UDF
 - No Looker connectivity



Pros

- Separates storage and computing. Storage is cheap
- Handles concurrency very well
- Semi-structured data support (Json, XML)
- Virtual warehouses (pay per usage, predictable cost)
- Lots of advanced, handy functionality

Concerns

- Will it be fast enough?
- Cost



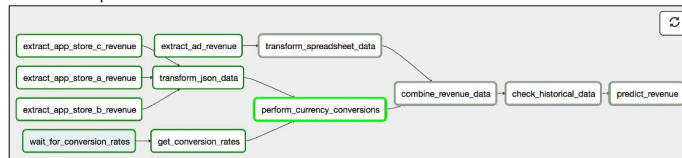
Apache Airflow

- Open Source, great community
- Well established
- Extremely versatile
- Powerful
- Distributed

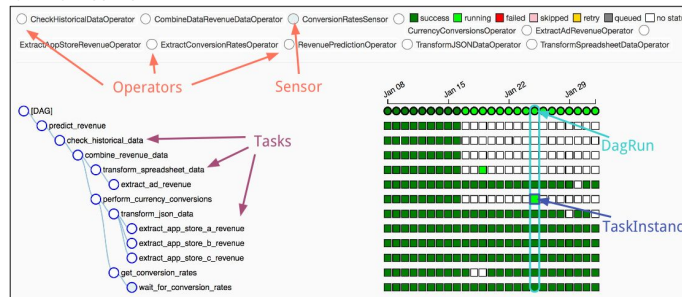
...& it could do some of the heavy-lifting

if Snowflake turned out to be too slow, or expensive

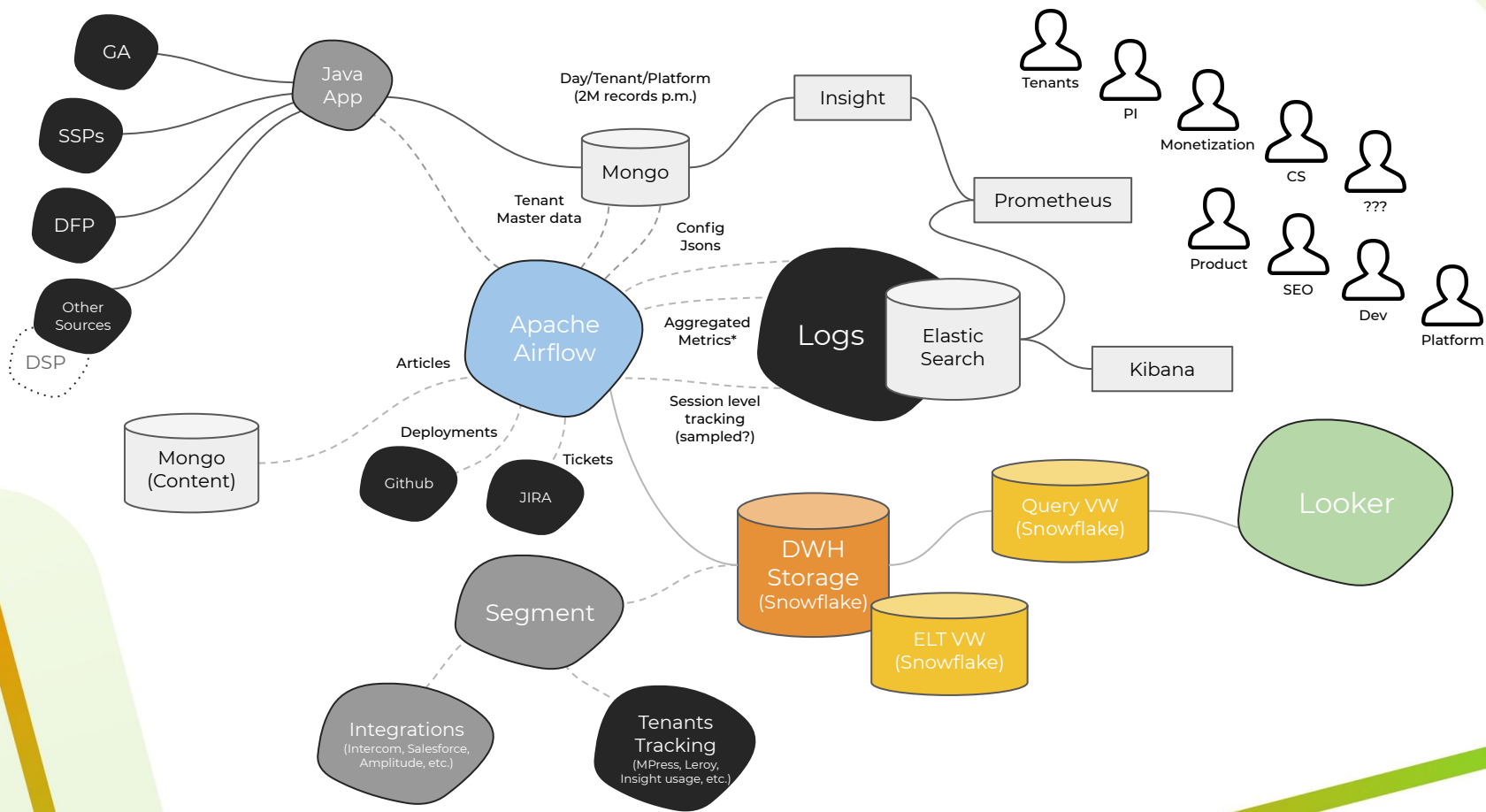
UI DAG Graph View



UI DAG Tree View



Hi-level data architecture proposal



Before proceeding, we must validate that:

- **Looker** did a good job at:
 - Resolving the complexity and fragmentation of our data sources
 - Removing (most) barriers to Data Accessibility & Literacy by modeling the required technical and business knowledge into its logical layer
- **Snowflake** could scale (both in terms of performance and costs)



Proof of Concept (PoC)



CE



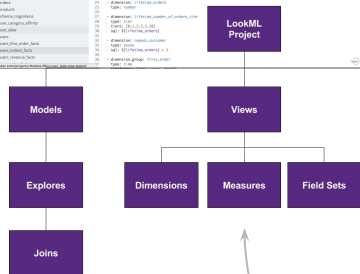
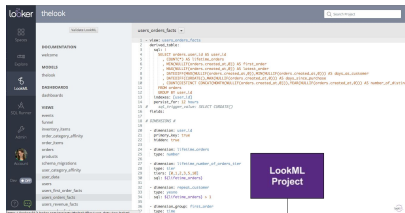
Product



Monetization

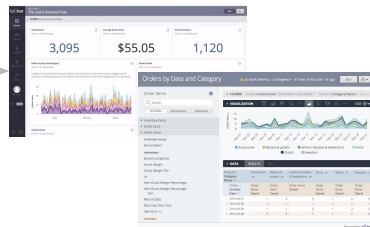


PI



Looker Model

Logical layer

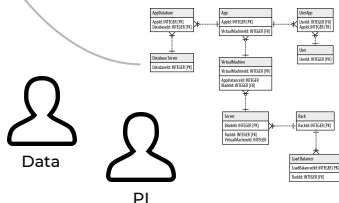


Looker dashboard and data exploration POC



Data Warehouse

Analytics data model



Data

Da



PI

Staging and processing (ELT)

Raw Json Data

- person "Dusty"
 - "person": "Dusty"
 - "favorite movies": ["Twelve number eleven", "He goed the bad and the ugly", "Kickass"]
- person "Dustin"
 - "person": "Dustin"
 - "favorite movies": ["Love free or Die Hard", "Jurassic Park", "Shamashak Redemption"]
- person "Oliver"
 - "person": "Oliver"
 - "favorite movies": ["Sheepless in Seattle", "Driving Miss Daisy", "Love Actually"]
- person "Mung"
 - "person": "Mung"
 - "favorite movies": ["Space Balls", "Saving Private Ryan", "40 First Dates"]

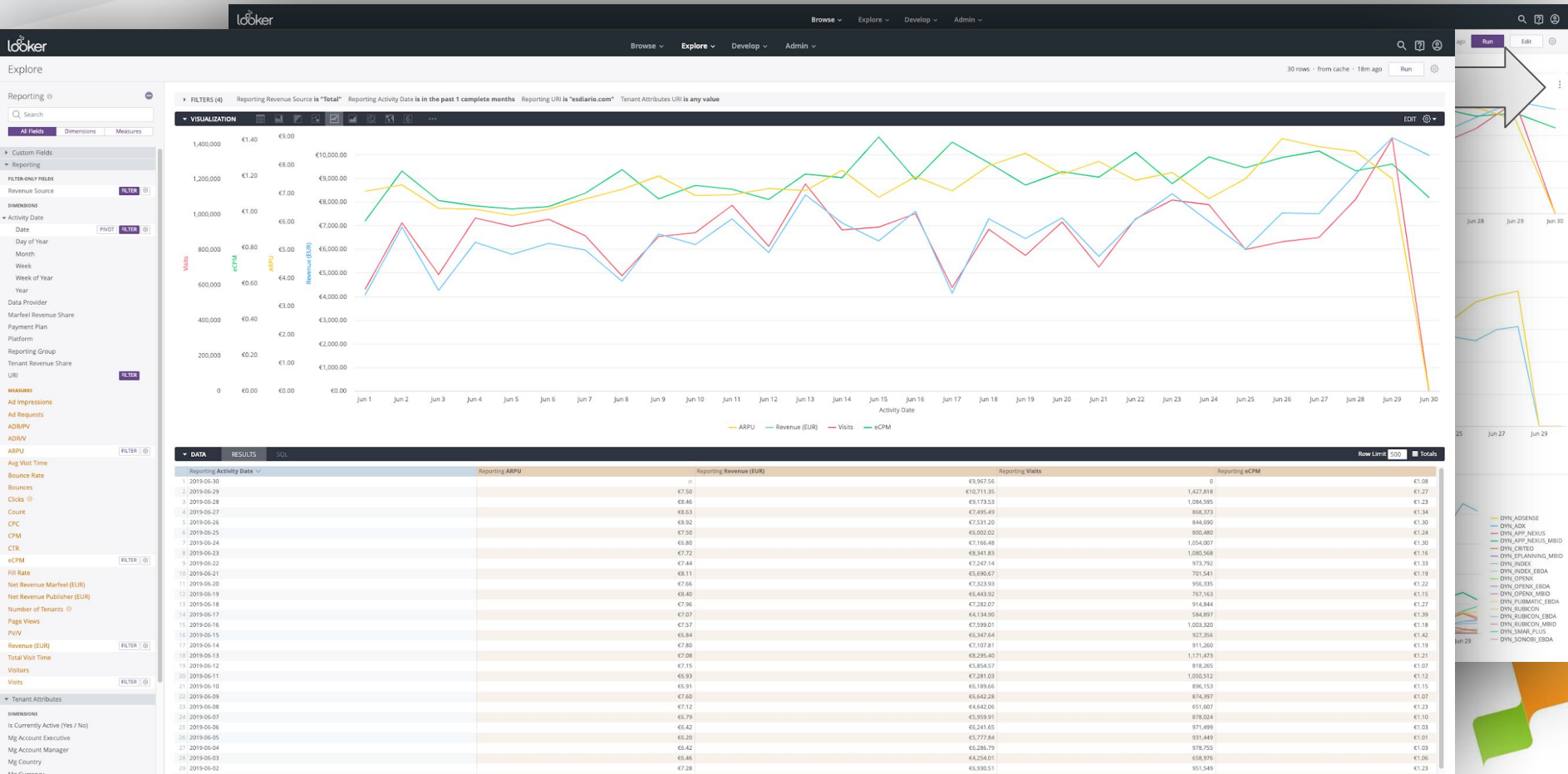


Pageviews & Visits (varys)



GA Connect, SSP, ADX, etc.

THE RESULT



Looking back...

In just a few years...

When	Where	Data Sources	ETL	DWH	BI Layer
2012	Softonic	Mostly backend DB	<i>Pentaho Data Integration</i>	SQL Server + <i>Hadoop</i>	<i>Qlikview</i> + <i>Tableau</i>
2014	King	Mostly tracking events	Java / Jenkins	<i>Hadoop</i> + Exasol (then BigQuery)	Qlikview (then Looker)
2016	Typeform	Mostly backend DB	Pentaho Data Integration (then Apache Airflow)	Redshift	Tableau (then Looker)
2019	Marfeel	Multiple APIs & microservices	Apache Airflow (Python)	Snowflake	Looker
20??	?				

Thank you.

Alessandro Pregnolato
Head of Data

