

Building a Financial Data Warehouse: A lesson in empathy

Solmaz Shahalizadeh
Data Analysis Lead, Shopify
[@solmaz_sh](https://twitter.com/solmaz_sh)

@solmaz_sh

Data Team Lead @ Shopify

Statistics/ Machine Learning in Cancer Research

Analytics and development in Banking



what is this talk about?

1. How did our data journey start at Shopify?
2. Why we needed a massive change?
3. Change is difficult, how did we push forward?
4. After all of this, where are we now? What next?
5. Takeaways

[Ways to sell](#)[Pricing](#)[Blog](#)[More](#)[Log in](#)

Shopify is everything you need to **sell anywhere**

Start your free 14-day trial today!

Email address

[Get started](#)



ONLINE STORE



Overview

Today

Sales

	Total	Orders
Today	\$199.98	2
Yesterday	\$119.97	3
Last 7 Days	\$829.87	13
Last 30 Days	\$3.9k	62
Last 90 Days	\$10.3k	164



Visitors

	Total	Unique
Today	1337	1322
Yesterday	20k	17.5k
Last 7 Days	152k	105k
Last 30 Days	565.8k	358.9k
Last 90 Days	1.8m	992.6k



Conversions

		Today
Added to Cart		
0.33% +0.33%	Reached Checkout	0.33% -0.5%



Top Products

<input type="checkbox"/> Muscles Leggings	4 sold
<input type="checkbox"/> Muscles Leggings	4 sold
<input type="checkbox"/> Muscles Leggings	4 sold

Traffic Sources

	Today
Direct	88%
995 visitors	
Referrals	7%
74 visitors	
Search Engines	5%
60 visitors	

Social Referrals

Today

Facebook	8%
	10 visitors

0%
0 visitors0%
0 visitors

Top Countries

Today

United States	78%	190
Canada	7%	17
Russia	4%	10
United States	78%	190
Canada	7%	17
Russia	4%	10
Canada	7%	17
Russia	4%	10

Top Referrers

Today

m.facebook.com	54%	20
i.facebook.com	16%	6
www.facebook.com	11%	4
iris.josephnogucci.com	5%	2
josephnogucci.us5.list-manage.com	3%	1
www.facebook.com	3%	1
iris.josephnogucci.com	3%	1
josephnogucci.us5.list-manage.com	3%	1

Top Search Terms

Today

nogucci	25%	1
joseph nogucci	25%	1
josephnogucci	25%	1
prosperity buddha...	25%	1

About Shopify

Shopify is a leading cloud-based, multichannel commerce platform designed for small and medium-sized businesses. Merchants can use the software to design, set up and manage their stores across multiple sales channels, including web, mobile, social media, marketplaces, brick-and-mortar locations, and pop-up shops. The platform also provides a merchant with a powerful back-office and a single view of their business. The Shopify platform was engineered for reliability and scale, using enterprise-level technology made available to businesses of all sizes. Shopify currently powers over 200,000 businesses in approximately 150 countries, including: Tesla Motors, Budweiser, Red Bull, LA Lakers, the New York Stock Exchange, GoldieBlox, and many more.

Key Facts

SHOPIFY PLATFORM RELEASED
2006

MERCHANTS IN
~150 Countries

EMPLOYEES
900+

APPS IN OUR APP STORE
1000+

MERCHANTS
200,000+

THEMES IN OUR THEME STORE
100+

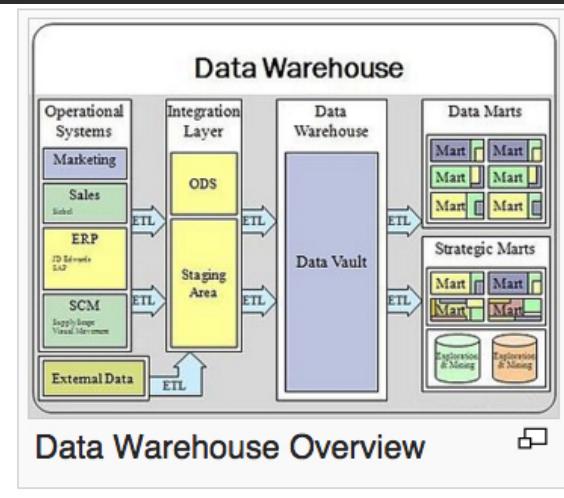
TOTAL SALES ON SHOPIFY
\$12 Billion

EXPERTS IN OUR NETWORK
680+

Where did our data journey start?

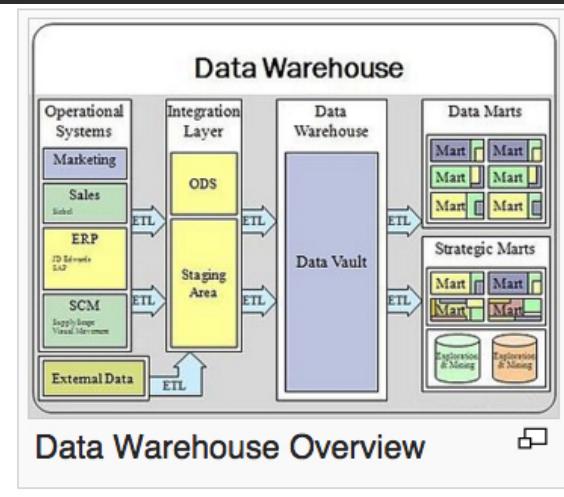
What is a Data Warehouse?

In computing, a **data warehouse** (DW or DWH), also known as an **enterprise data warehouse** (EDW), is a system used for **reporting** and **data analysis**. DWs are central repositories of integrated data from one or more disparate sources. They store current and historical data and are used for creating analytical reports for knowledge workers throughout the enterprise. Examples of reports could range from annual and quarterly comparisons and trends to detailed daily sales analyses.



What is a Data Warehouse?

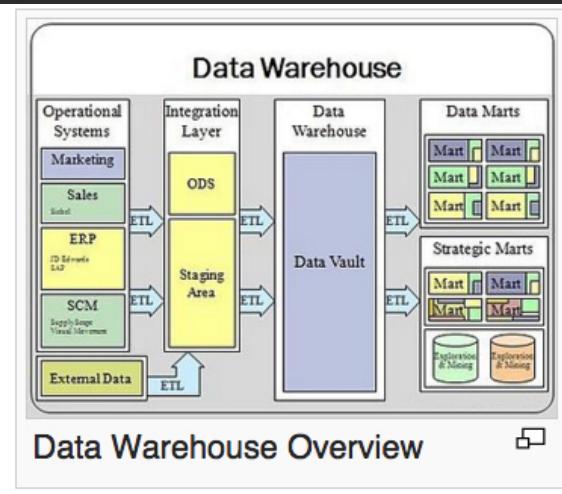
In computing, a **data warehouse** (DW or DWH), also known as an **enterprise data warehouse** (EDW), is a system used for **reporting** and **data analysis**. DWs are central repositories of integrated data from one or more disparate sources. They store current and historical data and are used for creating analytical reports for knowledge workers throughout the enterprise. Examples of reports could range from annual and quarterly comparisons and trends to detailed daily sales analyses.



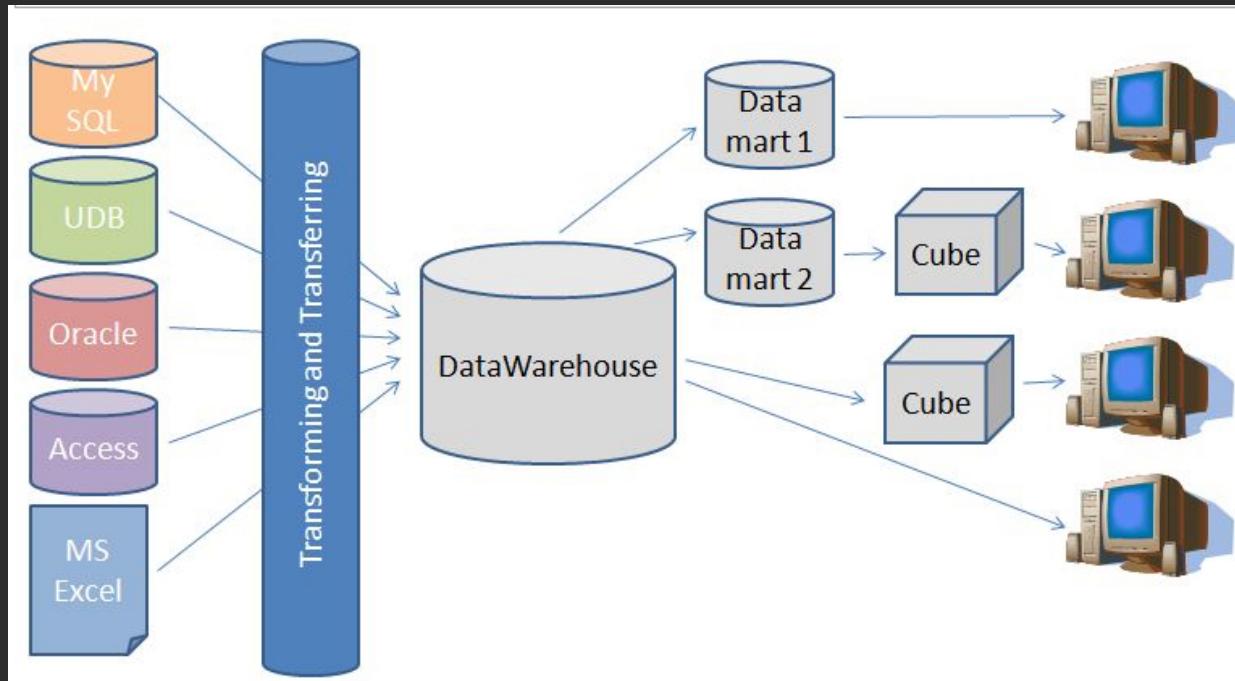
What is a Data Warehouse?

In computing, a **data warehouse** (DW or DWH), also known as an **enterprise data warehouse** (EDW), is a system used for **reporting** and **data analysis**. DWs are central repositories of integrated data from one or more disparate sources. They store current and historical data and are used for creating analytical reports for knowledge workers throughout the enterprise.

Examples of reports could range from annual and quarterly comparisons and trends to detailed daily sales analyses.



ETL systems overview



1st generation data warehouse

ETL:

- E: xtract from Production MySQL and other sources using in-house ruby etl
- !T: very little (almost no) T in the ETL
- L: oad into a Vertica database

* a copy of “any” operational table in the warehouse

1st generation reporting

In-house reporting (SQL and coffeescript):

- Clean the data at the report level using SQL
- Define the metric in the report
- use D3 for visualizations

* JIT metric definition

KPI for an e-commerce SaaS company

Customers

GMV (Gross Merchandise Volume)

Revenue

KPI for an e-commerce SaaS company

Customers: *by cohort, by monthly recurring revenue, by partnership*

GMV: *by payment gateway, by industry, by monthly recurring revenue*

Revenue: *by cohort, by subscription type, by partnership*



shop details



billing information



partnership information

A black and white photograph of four large, cylindrical concrete silos standing side-by-side in a field. They are partially hidden behind a line of bushes and small trees. The silos have horizontal bands and some vertical markings. The sky is clear.

hard to connect silos

image link <https://goo.gl/aOkbXU>

Different business areas, different definitions

- What is the geographical location of a customer?
- What do we mean when we talk about cohorts?
- Did the customer churn?
- ...

Things that are called the same thing, do not mean the same thing

A black and white cat is looking down at a green lizard on a wooden floor. The cat's head is on the left, and the lizard is on the right, both facing each other. The lizard has its mouth open. The background shows a wooden floor and a portion of a colorful rug.

Difficult to talk
about data

image link <https://goo.gl/VJwFYx>

Exploratory analysis for all:

While analysts idea of question!= business idea of question:

- 1) write ad-hoc SQL query
- 2) assume analyst 100% understands the business process
- 3) assume the ad-hoc cleaning and conforming of values is reproducible
- 4) extract data and sending data dump
- 5) upon seeing the results, business user realizes this is not what they wanted

A photograph of Keanu Reeves sitting cross-legged on a grassy field, looking down at a small object in his hands. A giant panda is lying on the grass next to him, facing him. The scene is outdoors with green grass in the background.

sad users & analysts

image link <https://goo.gl/zdatv6>

“The world as we have created it is a process of our thinking. It cannot be changed without changing our thinking.”

- Albert Einstein

Dimensional modelling

Dimensional Modelling is a design technique for databases intended to support end-user queries in a data warehouse. It is oriented around understandability and performance.

- Wikipedia

Focus on a business process:

- What are the dimensions?
- What are the measures?
- What is the grain?

Focus on modelling a process

Not a specific source tables but a business process

e.g: processing a financial transaction in exchange for a product as part of GMV

Define dimensions, facts, grain

Business Process: processing financial transaction

Dimensions: *payment gateway that processed the transaction, Industry of the shop, monthly recurring revenue (subscription amount) of the shop*

Measures: \$ amount of transaction

Define dimensions, facts, grain

Business Process: processing financial transaction

Dimensions: *payment gateway that processed the transaction, Industry of the shop, monthly recurring revenue (subscription amount) of the shop*

Measures: \$ amount of transaction

Remove anything without analytical value: name of the customer, exact address of the customer, etc.

**KEEP
CALM
AND
EDUCATE**

2nd generation data warehouse

E: in-house Ruby extracts to



T: Heavy transforms in



L: Load into Amazon Redshift



Building the ETL framework in Spark

spark out of memory error

Web Videos News Shopping Images More ▾ Search tools

About 1,300,000 results (0.33 seconds)

out of memory - spark java.lang.OutOfMemoryError: Java ...
stackoverflow.com/.../spark-java-lang-outofmemoryerror-java-heap-spac... ▾
Jan 15, 2014 - I have a few suggestions: If your nodes have 6g, then use 6g rather than 4g, `spark.executor.memory=6g`. Make sure you're using all the **memory** by ...

Tuning - Spark 1.5.1 Documentation
spark.apache.org/docs/latest/tuning.html ▾
Determining Memory Consumption; Tuning Data Structures; Serialized RDD you find that your JVM is garbage-collecting frequently or running **out of memory**, ...

Memory Issues in while accessing files in Spark - Cloudera ...
<https://community.cloudera.com/t5/.../Spark/Memory...Spark/.../18250> ▾
Sep 4, 2014 - Hi, I am working on a spark cluster with more than 20 worker nodes and each node with a memory of 512 MB. ... Below are the error messages and the corresponding code. I believe that's what's running **out of memory**.

Troubleshooting Hive on Spark - Cloudera
www.cloudera.com/.../en/.../admin_hos_troubleshooting.html ▾ Cloudera ▾
Important : Hive on Spark is included in CDH 5.4 but is not currently ... Problem: **Out-of-memory error**: You might get an **out-of-memory error** similar to the ...

Apache Spark User List - Re: Out of memory on large RDDs
apache-spark-user-list.1001560.n3.nabble.com/Re-Out-of-memory-on-la... ▾
Mar 11, 2014 - 7 posts · 5 authors

```
org.apache.spark.SparkException: Job aborted due to stage failure: Task 0.0:0 failed 4 times.  
Driver stacktrace:  
    at org.apache.spark.scheduler.DAGScheduler.org$apache$spark$scheduler$DAGScheduler.  
    at org.apache.spark.scheduler.DAGScheduler$$anonfun$abortStage$1.apply(DAGScheduler.scala:1015)  
    at org.apache.spark.scheduler.DAGScheduler$$anonfun$abortStage$1.apply(DAGScheduler.scala:1015)  
    at scala.collection.mutable.ResizableArray$class.foreach(ResizableArray.scala:59)  
    at scala.collection.mutable.ArrayBuffer.foreach(ArrayBuffer.scala:47)  
    at org.apache.spark.scheduler.DAGScheduler.abortStage(DAGScheduler.scala:1015)  
    at org.apache.spark.scheduler.DAGScheduler$$anonfun$handleTaskSetFailed$1.apply(DAGScheduler.scala:1015)  
    at org.apache.spark.scheduler.DAGScheduler$$anonfun$handleTaskSetFailed$1.apply(DAGScheduler.scala:1015)  
    at scala.Option.foreach(Option.scala:236)  
    at org.apache.spark.scheduler.DAGScheduler.handleTaskSetFailed(DAGScheduler.scala:1015)  
    at org.apache.spark.scheduler.DAGSchedulerEventProcessActor$$anonfun$receive$2.apply(DAGSchedulerEventProcessActor.scala:1015)  
    at akka.actor.ActorCell.receiveMessage(ActorCell.scala:498)  
    at akka.actor.ActorCell.invoke(ActorCell.scala:456)  
    at akka.dispatch.Mailbox.processMailbox(Mailbox.scala:237)  
    at akka.dispatch.Mailbox.run(Mailbox.scala:219)  
    at akka.dispatch.ForkJoinExecutorConfigurator$AkkaForkJoinTask.exec(AbstractDispatcher.java:393)  
    at scala.concurrent.forkjoin.ForkJoinTask.doExec(ForkJoinTask.java:260)  
    at scala.concurrent.forkjoin.ForkJoinPool$WorkQueue.runTask(ForkJoinPool.java:133)  
    at scala.concurrent.forkjoin.ForkJoinPool.runWorker(ForkJoinPool.java:1979)  
    at scala.concurrent.forkjoin.ForkJoinWorkerThread.run()
```

The Little Warehouse That Couldn't, Or: How We Learned to Stop Worrying and Move to Spark

[Yandu Oppacher](#) (Shopify)

Monday, June 15

5:45 PM – 6:00 PM

Imperial Ballroom (Level 2)

[SLIDES PDF](#)

[VIDEO](#)



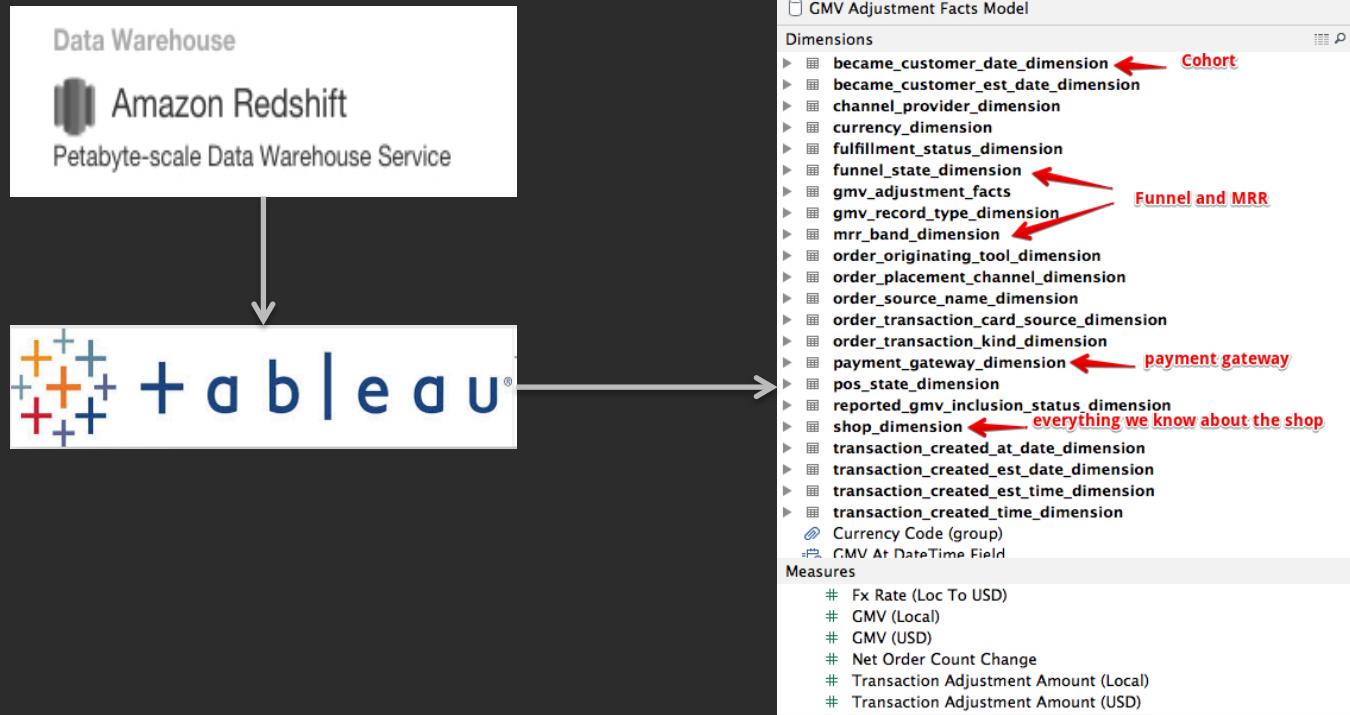
ABOUT YANDU

Yandu spent seven years working on the relational engine of IBM's Cognos BI and reporting tool, working with dozens of vendors and helping to rebuild their 15 year old engine. For the past year and a half he's been working on Shopify's data engineering team to help build a new framework for dimensionally modeling Shopify's data.



Shopify's commerce platform now powers over 150K merchants and continues to grow. The huge volume and variety of our data pushed our homegrown reporting and warehousing system to the edge. As the maintenance and performance costs became too much we moved to HDFS and Spark, using both python and scala to transform our vast amounts of operational data into dimensional models to provide better insight. Find out the lessons we've learned moving our entire organization onto fully conformed facts and dimensions, the clusters we've cratered, the walls we've hit, and what we did to overcome them to build our Starscream framework.

2nd generation reporting



KPI for an e-commerce SaaS company

Customers: *by cohort, by monthly recurring revenue, by partnership*

Sales: *by payment gateway, by industry, by monthly recurring revenue*

Revenue: *by cohort, by subscription type, by partnership*



shop details



billing information



partnership information

A walk in the data garden

Insightful reports

Self-service analytics (Tableau/ ipython)

New Business Analyst role

Faster feature selection (ML)



image link <https://goo.gl/MtVWcJ>

Takeaways

- Separation of business process/data model from ETL implementation
- Dimensional modelling for business process/ data models
- Clear data definitions to enable self service analytics
- Reproducible data models for reliable (and faster) predictive analysis

Questions?

We are hiring!

<http://shopify.com/careers>

Thank you

@solmaz_sh