

# H&M Database System



Group 3:

Zhikun Chen, Jinglin He, Xinlei He, Wan Ting Hsu, Sepideh Kiumarsi, Yuexin Li, Yuting Wang

## Background

H&M is the world's second-largest clothing retailer

- Operates in more than 74 countries with more than 5,000 stores worldwide
- Closed 5% of stores in 2021 because of COVID

Shift from traditional retail settings to eCommerce

## Problem

Lots of competition within the online retail space

- H&M needs to differentiate itself

Different consumer expectations in eCommerce that H&M needs to adapt to

Create the best online shopping experience for shoppers

## Our Proposed Solution

To help H&M more easily access consumer data we have created an organized database that can be used to understand their customers and ultimately develop product recommendation systems

# Original Data Description

- **Data Source:** Kaggle
- **Data Name:** H&M Personalized Fashion Recommendation datasets
  - Article csv file
  - Customer csv file
  - Transaction csv file

- **Data Provider:** H&M Group

- **Data Access Link:**

<https://www.kaggle.com/competitions/h-and-m-personalized-fashion-recommendations/data>

## Key data information:

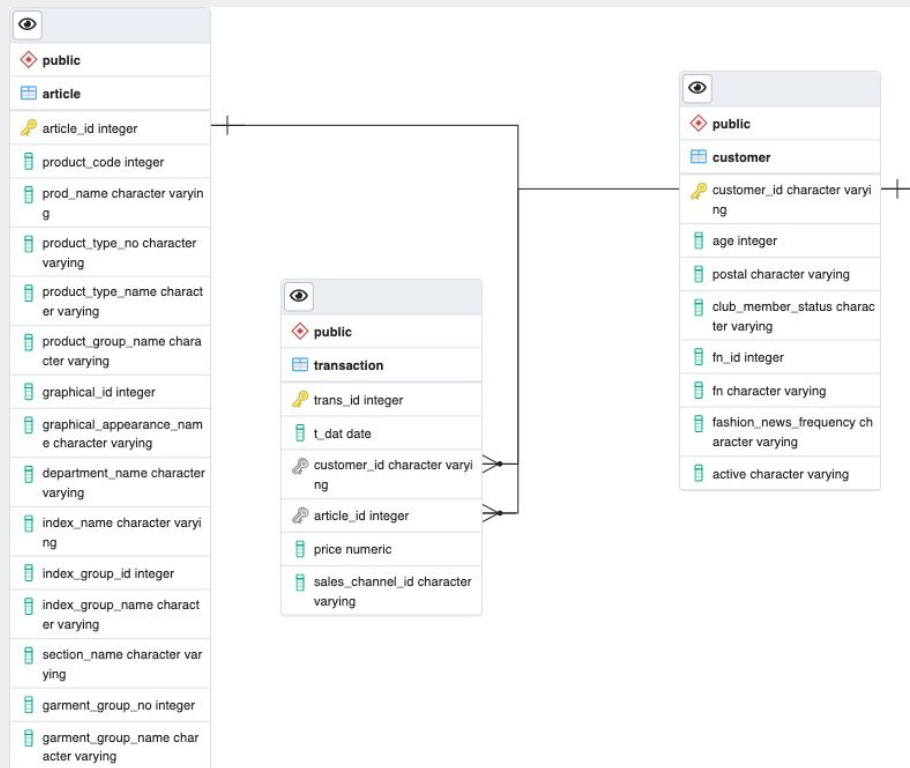
- **Article file**
  - Article id, product name, product type, graphical appearance, color, department etc.
- **Customer file**
  - Customer id, customer active status, customer membership status, age.
- **Transaction file**
  - Transaction date, customer id, article article (products purchased), price, sales channel

**H&M Personalized Fashion Recommendations**

Provide product recommendations based on previous purchases

# Normalization - 1NF

- Each table cell should contain a single value.
- Each record needs to be unique.



article_id	product_code	prod_name	product_type_no	product_type_name	product_group_name	graphical_appearance_no	graphical_appearance_name	colour_group_code	colour_group_name	perceived_colour_value_id	perceived_colour
106773015	106775	Strap top	253	Wool top	Garment Upper body	1010016	Solid	9	Black	4	Dark
106773044	106775	Strap top	253	Wool top	Garment Upper body	1010016	Solid	10	White	3	Light
106773061	106775	Strap top (1)	253	Wool top	Garment Upper body	1010017	Stripe	11	Off White	1	Dusty Light
110060001	110065	CP T-shirt (dark)	306	Bra	Underwear	1010016	Solid	9	Black	4	Dark
110060002	110065	CP T-shirt (dark)	306	Bra	Underwear	1010016	Solid	10	White	3	Light
110060011	110065	CP T-shirt (dark)	306	Bra	Underwear	1010016	Solid	12	Light Beige	1	Dusty Light
111560001	111565	20 den 1p Stockings	304	Underwear Tights	Socks & Tights	1010016	Solid	9	Black	4	Dark
111560003	111565	20 den 1p Stockings	302	Socks	Socks & Tights	1010016	Solid	13	Beige	2	Medium Dusty
111560001	111566	Shape Up 30 den 1p Tights	273	Leggings/Tights	Garment Lower body	1010016	Solid	9	Black	4	Dark
111560001	111593	Support 40 den 1p Tights	304	Underwear Tights	Socks & Tights	1010016	Solid	9	Black	4	Dark
111600001	111609	200 den 1p Tights	304	Underwear Tights	Socks & Tights	1010016	Solid	9	Black	4	Dark
110719046	110719	SWATCHBERRY OC	252	Sweater	Garment Upper body	1010001	All over pattern	7	Grey	1	Dusty Light
110719052	110719	SWATCHBERRY OC	252	Sweater	Garment Upper body	1010001	All over pattern	71	Light Blue	1	Dusty Light
114420001	114425	Alma BANDEAU 2-g	306	Bra	Underwear	1010017	Stripe	10	White	3	Light
114420001	114425	Alma BANDEAU 1-g	306	Bra	Underwear	1010016	Solid	6	Light Blue	1	Dusty Light

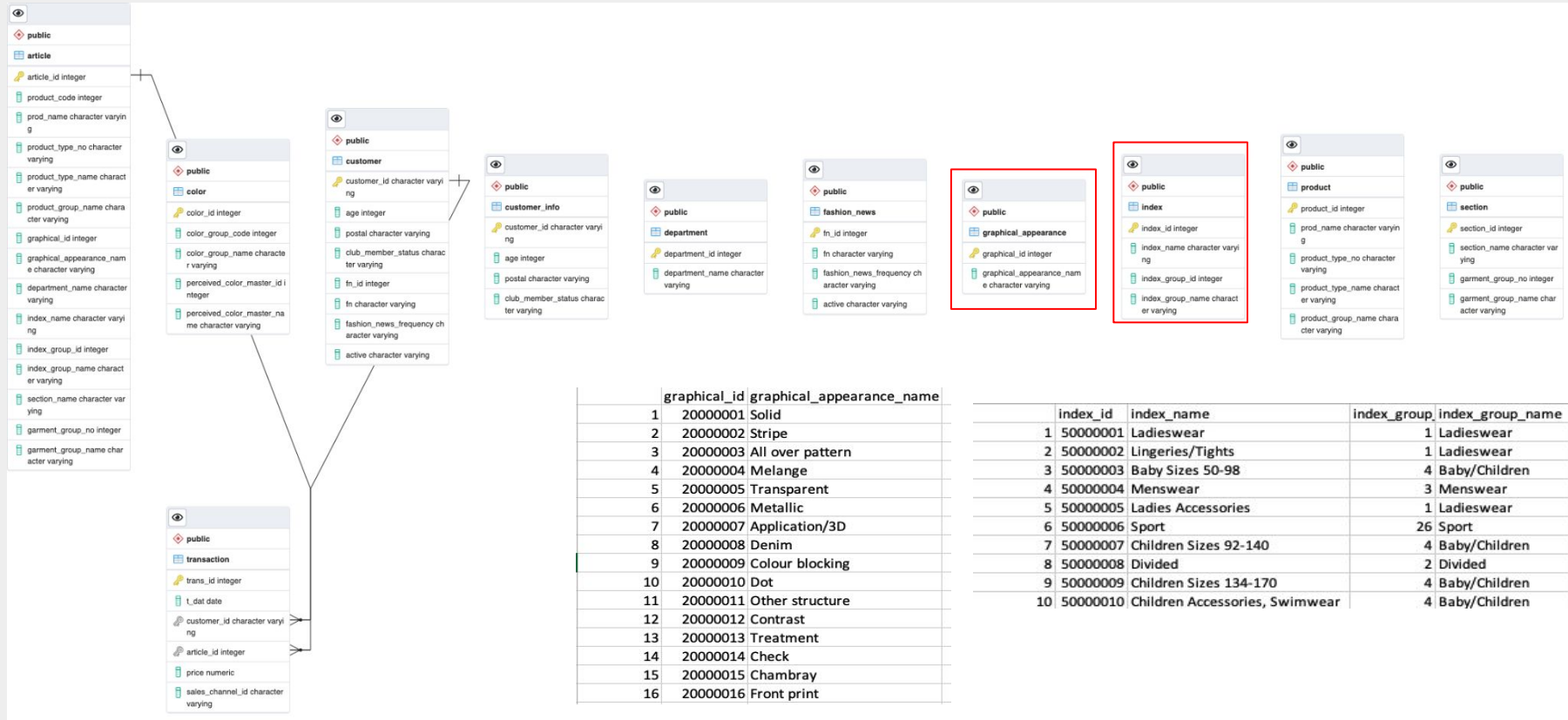
customer_id	FN	Active	club_member_status	fashion_news_frequency	age
00000dbacae5abe5e23885899a1fa44253a17956c6d1c3d25f88aa139fdcf657			ACTIVE	NONE	49
0000423b00ade91418ccfa3b26c6af3dd342b51fd051eec9c12fb36984420fa			ACTIVE	NONE	25
000058a12d5b43e67d225668fa1f8d618c13dc232df0cad8ffe7ad4a1091e318			ACTIVE	NONE	24
00005ca1c9ed5f5146b52ac8639a40ca9d57aeff4d1bd2c5fe1b1ca5dff07c43e			ACTIVE	NONE	54
00006413d8573cd20ed7128e53b7b13819fe5cfd2d801fe7f0c26dd8d65a85a	1.0	1.0	ACTIVE	Regularly	52
000064249685c11552da43ef22a5030f35a1471723d5b02dd9d22452b1f5a6					
0000757967448a6cb83efb3ea7a3fb9d418ac7ad2379d8cd0c725276a467a2a			ACTIVE	NONE	20
00007d2de826758b65a93dd24ce629ed66842531df6699338c5570910a014cc2	1.0	1.0	ACTIVE	Regularly	32

t_dat	customer_id	article_id	price	sales_channel_id
9/20/18	000058a12d5b43e67d225668fa1f8d618c13dc232df0cad8ffe7ad4a1091e318	663713001	0.050830508	2
9/20/18	000058a12d5b43e67d225668fa1f8d618c13dc232df0cad8ffe7ad4a1091e318	541518023	0.030491525	2
9/20/18	00007d2de826758b65a93dd24ce629ed66842531df6699338c5570910a014cc2	505221004	0.015237288	2
9/20/18	00007d2de826758b65a93dd24ce629ed66842531df6699338c5570910a014cc2	685687003	0.016932203	2
9/20/18	00007d2de826758b65a93dd24ce629ed66842531df6699338c5570910a014cc2	685687004	0.016932203	2
9/20/18	00007d2de826758b65a93dd24ce629ed66842531df6699338c5570910a014cc2	685687001	0.016932203	2
9/20/18	00007d2de826758b65a93dd24ce629ed66842531df6699338c5570910a014cc2	505221001	0.020322034	2
9/20/18	00083cda041544b2fbb0e0d2905ad17da7cf1007526fb4c73235dccbcb132280	688873012	0.030491525	1
9/20/18	00083cda041544b2fbb0e0d2905ad17da7cf1007526fb4c73235dccbcb132280	501323011	0.053372881	1

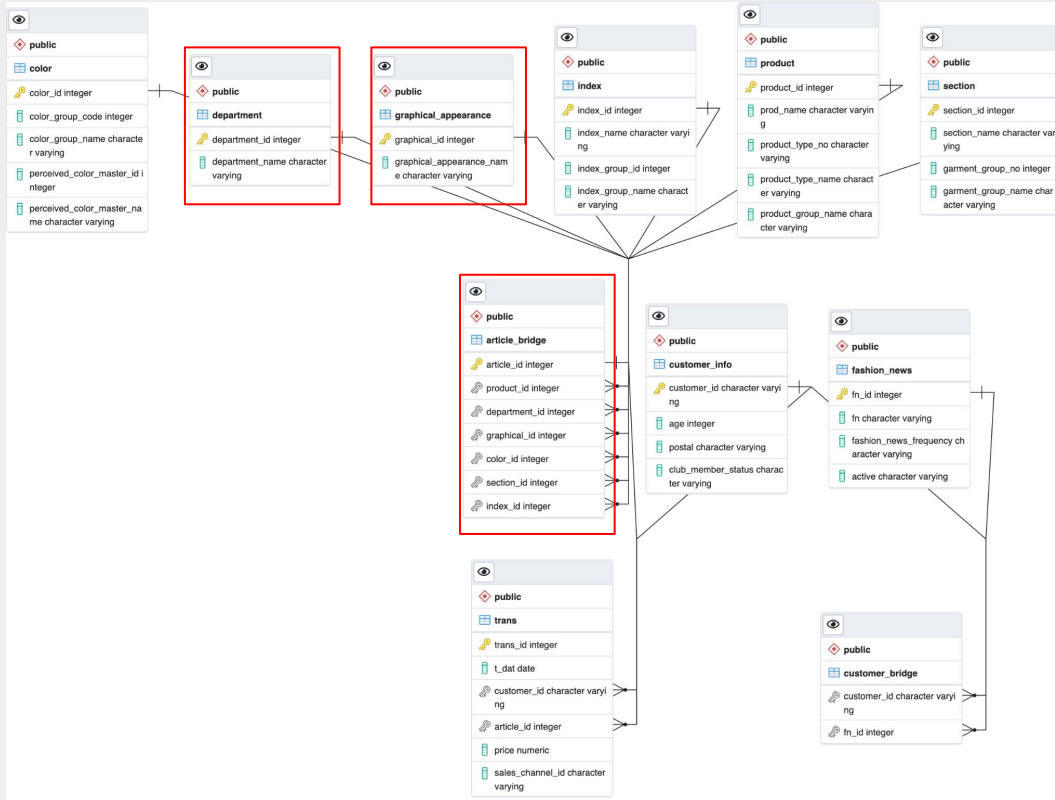
# Normalization - 2NF

- The table must be already in 1 NF and all non-key columns of the tables must depend on the PRIMARY KEY
- The partial dependencies are removed and placed in a separate table



# Normalization - 3NF

## Finalized ERD



- Non-Primary key columns shouldn't depend on the other non-Primary key columns
- There is no transitive functional dependency

article_id	colour_id	department	graphical_id	index_id	product_id	section_id	articlebridge'
108775015	41000001	31000001	20000001	50000001	10000001	60000001	108775015
108775044	41000002	31000001	20000001	50000001	10000001	60000001	108775044
108775051	41000003	31000001	20000002	50000001	10000002	60000001	108775051
110065001	41000001	31000002	20000001	50000002	10000003	60000002	110065001
110065002	41000002	31000002	20000001	50000002	10000003	60000002	110065002
110065011	41000004	31000002	20000001	50000002	10000003	60000002	110065011
111565001	41000001	31000003	20000001	50000002	10000004	60000003	111565001
111586001	41000001	31000003	20000001	50000002	10000006	60000003	111586001
111593001	41000001	31000003	20000001	50000002	10000007	60000003	111593001
111609001	41000001	31000003	20000001	50000002	10000008	60000003	111609001
116379047	41000009	31000001	20000001	50000001	10000011	60000001	116379047
118458004	41000011	31000001	20000004	50000004	10000012	60000005	118458004
118458029	41000010	31000001	20000004	50000004	10000012	60000005	118458029
118458038	41000013	31000001	20000004	50000004	10000012	60000005	118458038

department_id	department_name	graphical_id	graphical_appearance_name
31000001	Jersey Basic	20000001	Solid
31000002	Clean Lingerie	20000002	Stripe
31000003	Tights basic	20000003	All over pattern
31000004	Baby basics	20000004	Melange
31000005	Casual Lingerie	20000005	Transparent
31000007	Jersey	20000006	Metallic
31000008	EQ & Special Collections	20000007	Application/3D
31000009	Hair Accessories	20000008	Denim
		20000009	Colour blocking

# ETL Process

1. Clean tables by manipulating the data to only the attributes we want for the table. As you can see there are 105,542 rows for the color table.

```
{r}
colour <- articles %>% select(colour_group_code, colour_group_name, perceived_colour_value_id, perceived_colour_value_name, perceived_colour_master_id,
perceived_colour_master_name)
colour
```

Description: df [105,542 x 6]

colour_group_code	colour_group_name	perceived_colour_value_id	perceived_colour_value_name	perceived_colour_master_id	perceived_colour_master_name
9	Black	4	Dark	5	Black
10	White	3	Light	9	White
11	Off White	1	Dusty Light	9	White
9	Black	4	Dark	5	Black
10	White	3	Light	9	White
12	Light Beige	1	Dusty Light	11	Beige
9	Black	4	Dark	5	Black
13	Beige	2	Medium Dusty	11	Beige
9	Black	4	Dark	5	Black
9	Black	4	Dark	5	Black

1-10 of 105,542 rows

2. Find only the distinct rows. We eliminated the duplicate rows. Now we have 296 distinct color combination rows.

```
{r}
colour <- unique(colour)
colour
```

Description: df [296 x 6]

colour_group_code	colour_group_name	perceived_colour_value_id	perceived_colour_value_name	perceived_colour_master_id
9	Black	4	Dark	5
10	White	3	Light	9
11	Off White	1	Dusty Light	9
12	Light Beige	1	Dusty Light	11
13	Beige	2	Medium Dusty	11
7	Grey	1	Dusty Light	12
71	Light Blue	1	Dusty Light	2
6	Light Grey	1	Dusty Light	12
73	Dark Blue	4	Dark	2
8	Dark Grey	4	Dark	12

1-10 of 296 rows | 1-6 of 6 columns

3. Create a primary key (ID) for each table.

```
{r}
colour = colour %>% mutate(colour_id = sub("A", "4", row_number() + 1000000))
colour
```

Description: df [296 x 7]

colour_group_code	colour_group_name	perceived_colour_value_id	perceived_colour_value_name	perceived_colour_master_id	perceived_colour_master_name	colour_id
9	Black	4	Dark	5	Black	41000001
10	White	3	Light	9	White	41000002
11	Off White	1	Dusty Light	9	White	41000003
12	Light Beige	1	Dusty Light	11	Beige	41000004
13	Beige	2	Medium Dusty	11	Beige	41000005
7	Grey	1	Dusty Light	12	Grey	41000006
71	Light Blue	1	Dusty Light	2	Blue	41000007
6	Light Grey	1	Dusty Light	12	Grey	41000008
73	Dark Blue	4	Dark	2	Blue	41000009
8	Dark Grey	4	Dark	12	Grey	41000010

1-10 of 296 rows

4. Generate bridge tables with primary keys from related tables.

# Analytical Procedure

## Why do we interact with data?

- Data is not always in the format as wish
- Employees who have different job responsibilities have different requirements for data analysis results

## Customer Needs

- Fast, easy and convenient
- Getting a better understanding of the customers profile and demand for marketing campaign
- Different requirements for data

## Results and Benefits

- Ops team: To know better about the fashion and trends in the market and make appropriate designing and manufacturing decisions
- C-level: To obtain business insights and support decision-making procedure through data analysis.



# Tools & Techniques

We have multiple front-end and back-end interaction tools for our SQL database, such as Python, PGAdmin, and Metabase. Each tool has its own feature.

- **Python** is a front end programming language to help with data visualization, ETL process before data ingestion, and data analytics. Psycopg2 is the library we used for connection to our SQL server. We can use Hadoop for HDFS and Spark for ML analytics.
- **pgAdmin** is the platform for our SQL server. Our clients or our product manager can directly interact with our database using query language on pgAdmin.
- **Metabase** is the automated data visualization tool for our SQL database. Our data security is ensured with limited views for different levels of analyst and manager. Different positions within the company can view different visualizations and it is always updated.

# Database Interaction Demo

```
Query Editor  Query History

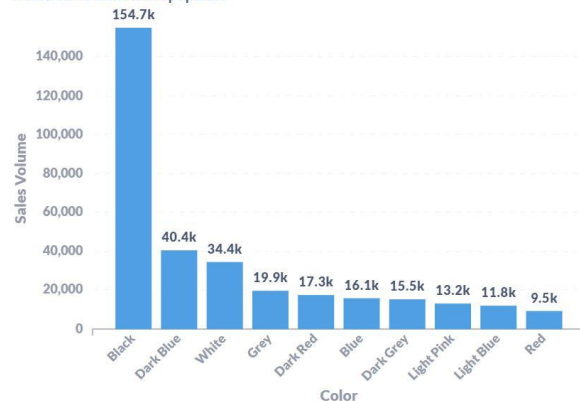
28 Question4: Which color is the most popular?
29 SELECT c.color_group_name, c.perceived_color_master_name, COUNT(c.color_group_name) AS number_of_color
30 FROM article_bridge AS a
31 INNER JOIN color AS c
32 ON a.color_id = c.color_id
33 INNER JOIN trans AS t
34 ON a.article_id = t.article_id
35 GROUP BY c.color_group_name, c.perceived_color_master_name
36 ORDER BY number_of_color DESC;
37
```

	color_group_name character varying	perceived_color_master_name character varying	number_of_color bigint
1	Black	Black	154721
2	Dark Blue	Blue	40383
3	White	White	34417
4	Grey	Grey	19897
5	Dark Red	Red	17283
6	Blue	Blue	16104
7	Dark Grey	Grey	15549
8	Light Pink	Pink	13205

## 5310 Group Project ⓘ

Our analytics • Edited in 4 hours by you

Which color is the most popular?



<http://160.39.250.97:3000/public/dashboard/e1077600-e5b2-485d-b939-e730658f1e3f>

Which product is the most popular?



# Conclusion

- Gain deeper business insights from the results generated by data analytics
  - more targeted marketing strategy and campaign based on the data analytical results.
- Gain better insight of customer behaviors, and helps H&M decision making process
  - eg. most popular color and graphical print among customers
- Provide information and experience for ETL procedure implementation in other department (e.g accounting, HR, public relations), improving company operation efficiency.



Thank you!