

Multimodal Deep Learning at Scale -Learning from Catalogs at Mirakl



Arthur Delaitre Data Scientist, Mirakl



Sang-Hoon Yoon Data Scientist, Mirakl Could not come to present his work :



Milton Minervino Data Scientist, Mirakl

What is a marketplace ?

Marketplace: a platform where multiple providers sell products and services















- Mirakl is the leader marketplace SaaS platform that empowers both B2B and B2C organizations to launch and grow an enterprise marketplace at scale.
- Our Data Science team aims to develop and industrialize AI features to improve user experience (automatic categorization and mapping, catalog cleaning, customer care support, ...)



Main topics

1. Catalog data

A large and diverse source of data ideal for multimodal learning

2. Catalog use cases

Key use cases: Categorization and Duplicates finding

- 3. Categorize with multimodal product embeddings Multimodal product embeddings: product2vec
- 4. Curate the catalog by removing duplicates A reliable and coherent product catalog

5. Key takeaways



Catalog data

A large and diverse source of data ideal for multimodal learning



What is a catalog in a marketplace ?

In Mirakl, the catalog is a database containing all products data of the marketplace

Multiple sellers can upload their products : **Mirakl Catalog Manager** ensures that everything is represented as a single, coherent and reliable catalog.

Image: Bluetooth Handsfree Headsets (CAT_25930) Apple iPhone 13 mini 128GB - Pink - Unlocked APPLE Cellular Phones / Cell Phone Batteries (CAT_11111) Apple iPhone 13 mini 128GB - Blue - Unlocked APPLE Cellular Phones / Cell Phone Cables (CAT_25930) Apple iPhone 13 mini 128GB - Blue - Unlocked APPLE Cellular Phones / Cell Phones & Smartphones SKU: 15923255 D3/31/2022 11:30 D4/08/2022 10:18 Tech Outlet Center Published Image: Cell Phone Cables (CAT_25930) Image: Cell Phone Signal Boosters (CAT_496383) Apple iPhone 13 Pro 256GB - Graphite - Unlocked APPLE Cellular Phones / Cell Phones & Smartphones SKU: 15923255 DPC-A: 999003559762 Am 12/23/2021 10:46 12/23/2021 12:47 Tech Outlet Center Published Image: Cell Phone Signal Boosters (CAT_29000) Apple iPhone 13 Pro 256GB - Graphite - Unlocked APPLE Cellular Phones / Cell Phones & Smartphones SKU: 15922540 DPC-A: 999003559782 AM 12/23/2021 10:46 12/23/2021 12:47 Tech Outlet Center Published Image: Cell Phone Signal Boosters (CAT_29000) Apple iPhone 13 Pro 256GB - Graphite - Unlocked - Refurbished APPLE Cellular Phones / Cell Phones & Smartphones SKU: 15922540 DPC-A: 999003559782 AM 12/23/2021 10:46 12/23/2021 12:47 Tech Outlet Center Published APPLE Cellular Phones / Cell Phones & Smartphones Cell Phones & Smartphones DPC-A: 9990035597588 AM 12/2	E E Cellular Phones (CAT_26336)	Product	۱ţ	Product ID	Creation date $\downarrow\uparrow$	Update date 🛛 🕸	Provider(s)	Status
 Bluetooth Speakerphones (CAT_25940) Cell Phone Batteries (CAT_11111) Cell Phone Cables (CAT_25930) Cell Phone Cases (CAT_30177) Cell Phone Signal Boosters (CAT_496383) Cell Phones & Smartphones (CAT_29000) Cell Phones & Smartphones (CAT_29000) 	Bluetooth Handsfree Headsets (CAT_25938)	•	Apple iPhone 13 mini 128GB - Pink - Unlocked - Refurbished APPLE Cellular Phones / Cell Phones & Smartphones	SKU: 16035702 UPC-A: 999003631017	03/31/2022 11:30 AM	04/08/2022 10:18 AM	Tech Outlet Center	Published
Image: Cell Phone Batteries (CAT_1111) Image: Cell Phone Cables (CAT_25930) Image: Cell Phone Cases (CAT_30177) Image: Cell Phone Signal Boosters (CAT_496383) Image: Cell Phones & Smartphones (CAT_29000) Image: Cell Phones & Smartphones (CAT_29000)	Bluetooth Speakerphones (CAT_25940)							
Image: Cell Phone Cables (CAT_25930) Image: Cell Phone Cases (CAT_30177) Image: Cell Phone Signal Boosters (CAT_496383) Image: Cell Phones & Smartphones (CAT_29000) Image: Cell Phone & Smartphones (CAT_	Cell Phone Batteries (CAT_11111)	0	Apple iPhone 13 mini 128GB - Blue - Unlocked - Refurbished APPLE Cellular Phones / Cell Phones & Smartphones	SKU: 15923255 UPC-A: 999003559762	12/23/2021 10:46 AM	12/23/2021 1:47 PM	Tech Outlet Center	Published
E Cell Phone Cases (CAT_30177) Cell Phone Signal Boosters (CAT_496383) Apple iPhone 13 Pro 256GB - Graphite - Unlocked - Refurbished APPLE Cellular Phones / Cell Phones & Smartphones SKU: 15922540 UPC-A: 999003559588 12/23/2021 10:46 AM 12/23/2021 12:17 Tech Outlet Center Published	Cell Phone Cables (CAT_25930)							
Image: Cell Phone Signal Boosters (CAT_496383) Apple iPhone 13 Pro 256GB - Graphite - Unlocked - Refurbished APPLE SKU: 15922540 UPC-A: 999003559588 12/23/2021 10:46 AM 12/23/2021 12:17 Tech Outlet Center Published Cell Phones & Smartphones (CAT_29000) Cellular Phones / Cell Phones & Smartphones Cellular Phones & Smartphones SKU: 15922540 UPC-A: 999003559588 12/23/2021 12:17 Tech Outlet Center Published	Cell Phone Cases (CAT_30177)							
Image: Cell Phones & Smartphones (CAT_29000) APPLE Cellular Phones / Cell Phones & Smartphones	Cell Phone Signal Boosters (CAT_496383)		Apple iPhone 13 Pro 256GB - Graphite - Unlocked - Refurbished APPLE Cellular Phones / Cell Phones & Smartphones	SKU: 15922540 UPC-A: 999003559588	12/23/2021 10:46 AM	12/23/2021 12:17 PM	Tech Outlet Center	Published
	Cell Phones & Smartphones (CAT_29000)							

Figures: Mirakl Catalog Management





Number of products*Number of categories310 M140 k

Number of marketplaces Number of sellers

+300

+200 k



The kind of data we have

Product data contain:

- **images** (product pictures)
- texts
- tabular data (category, size, storage capacity, ...)



Different product data from sellers



The kind of data we have

- 🗆 🔚 Cellular Phones (CAT_26336)
 - 📱 Bluetooth Handsfree Headsets (CAT_25938)
 - Cell Phone Signal Boosters (CAT_496383)
 - Cell Phones & Smartphones (CAT_29000)



Title: iPhone 13 6,1" 5G 128 Go Double SIM **Description**: This is an iPhone ... **Color**: Midnight **Dimensions**: 146,7 x 71,5 x 7,65 mm Screen: 6.1" OLED **Processor**: A15 Bionic Chip **Camera resolution**: 12MP (Ultra Wide) **Weight**: 173g Storage: 128GB



The kind of data we have





Title: iPhone 13 6,1" 5G 128 Go Double SIM **Description**: Le dernier smartphone d'Apple ... Color: Gold 🥂 **Dimensions**: 146,7 x 71,5 x 7,65 mm Screen: NaN / **Processor**: A15 Bionic Chip **Camera resolution**: 12MP (Ultra Wide) Weight: 6,10 ounces Storage: 128-512 GB



Recurring issues with product data

The recurring issues:

- Data quality
 - Product with missing/ corrupted data
 - Sellers may mislabel
- Multiple standards
 - Different metric systems
 - Different languages
 - Heterogeneous category trees

These product data issues prevent marketplaces from proposing a reliable catalog



Catalog use cases Categorization and Duplicates finding



Use case 1 : Product Categorization

Some products on the marketplace are assigned to a wrong category



Goal: predict the category with product data such as image, texts, color, size, etc



Use case 2 : Product duplicates finding



Goal: find the product data referring to the same product to clean the catalog



Categorization

Categorize with multimodal product embeddings



Categorization model architecture



Multimodal categorization model that corrects mislabeled product categories

Requirements:

- Robustness: stable with missing input
- Multilingual support: 5 languages

What we did:

- Data augmentation with dropout
- Use of multilingual pre-trained embeddings
- Multi-task learning and hierarchical classification

Ability to compute product2vec embeddings



Representation of products: product2vec





Leveraging the product2vec

Achievements with product2vec

- Categorization on very small catalogs
- Similar products detection
- Category embeddings
- Visual attribute prediction
- Use product2vec as an input for ML catalog models



Engineering considerations



Training challenges



Training dataset creation & preprocessing:

- Distributed on multiple clusters
- Dataset size exceeding memory capacity



άW

Training on terabytes of data using TFRecords

Why TFRecords ?

DATA+AI

SUMMIT 2022

- Save texts and images as **binary files**.
- A Dataset can be composed of multiple files.
- Allows for optimized data fetch without having to load everything in memory.







File specification and Dataset creation

Training in less than 3 hours

- For a dataset of ~20 M products
- On a g4dn.8xlarge (single GPU)

Training challenges





ONNX & ONNX Runtime

ONNX - open standard format defining a common set of computational functions:

- ONNX represents deep learning models in a wide variety of frameworks
- **ONNX Runtime** provides tools to optimize the ONNX graph
- Dynamic quantization to reduce latency and model size



With ONNX Runtime and dynamic quantization, the use of huge models like transformers (e.g. BERT models) in production is possible with CPUs



Achievements with categorization model

Mirakl categorizes automatically products into the right categories

We obtained the cornerstone of catalog ML use cases: product2vec embeddings



Finding duplicate products

A reliable and coherent product catalog



Finding duplicates use case

Goal : find product data from sellers that refer to the same product





Notes:

- Catalogs may be very **dense**: **thousands** of similar products in a category
- Sellers may describe the same product very differently

Product2vec embeddings are not precise and robust enough for this task

Model architecture

Goal : Take both product data as input, infer whether or not they refer to the same product



Focus on duplicated products search

Reference set



Cannot use product embeddings (not precise enough, either too many false positives or no positives at all)



Can use an **algorithm that compares each product low-level data to tell if they relate to the same product**. Filtering and pairwise comparisons



\rightarrow We need a robust architecture



Engineering considerations



Inference challenges



0

Filtering step

Goal : Drastically decrease the number of candidate pairs so that a multimodal deep learning model can be used on a limited set.

- 1. Compute basic features on images:
 - Histograms
 - Texture
 - Aspect ratio
- 2. Derive some features for pairs of images using vectorized computing
- 3. Apply a Decision Tree
 - 1M predictions per second

Consequence : From 1B pairs to 1M



DATA+A

SUMMIT 2022



Decisior

tree

Inference challenges

Image + Texts







Spark optimization

Optimize data frames joins:

- Use of **spark.sql.autoBroadcastJoinThreshold** to tune joins or deactivate it if needed.
- Carefully **partition** files and Dataframe, and **persist** Data **On disk only** for large DataFrames (do not forget to unpersist).

Optimize data queries:

• Use of **Z-ordering** in delta tables to speed up data queries







Spark optimization

Tips

- Add a non serializable model into spark context to predict with Spark UDF/pandas UDF
- Singleton class avoids loading object at every UDF call and accelerates executions

Be careful when using custom UDF, they can introduce memory leaks



rom pyspark import SparkFiles

class Singleton(type):

_instances = {}

def __call__(cls, *args, **kwargs):

if cls not in cls._instances:

cls._instances[cls] = super(Singleton, cls).__call__(*args, **kwargs)
return cls._instances[cls]

class <u>CatalogEstimator</u>(*metaclass=<u>Singleton</u>*):

model = None

def __init__(self, remote_model_path):

File_name = get_file_name(remote_model_path) Estimator.model = self._load_model(SparkFiles.get(file_name))

@staticmethod

def _load_model(model_path):

return catalog_model.load_model(model_path)

@staticmethod

def predict(remote_model_path, sentence):

spark.sparkContext.addFile("/tmp/model.bin")

@pandas_udf(T.StringType())

def predict_pandas_udf(s: pd.Series) -> pd.Series:

s = s.apply(*lambda x*: CatalogEstimator.predict("/tmp/model.bin", x))

return s

spark_DF.select('id', predict_pandas_udf(col("body")).display())

DATA+AI SUMMIT 2022

Conclusion - key takeaways

- Multimodal product embeddings are the cornerstone of Machine Learning on Catalog Data
- Before pulling out the big guns, check if pragmatic and frugal solutions can resolve your O(n²) problems

 Optimization on Spark and inference time reduction with ONNX are paramount to scale your pipeline



DATA+AI SUMMIT 2022

Thank you