



# E-commerce products

Taxonomy and extended similarity measuring between products

Team 4 - Frytki: Szymon Rećko, Mateusz Sperkowski, Patryk Tomaszewski, Kinga Ułasik

# What is e-commerce?

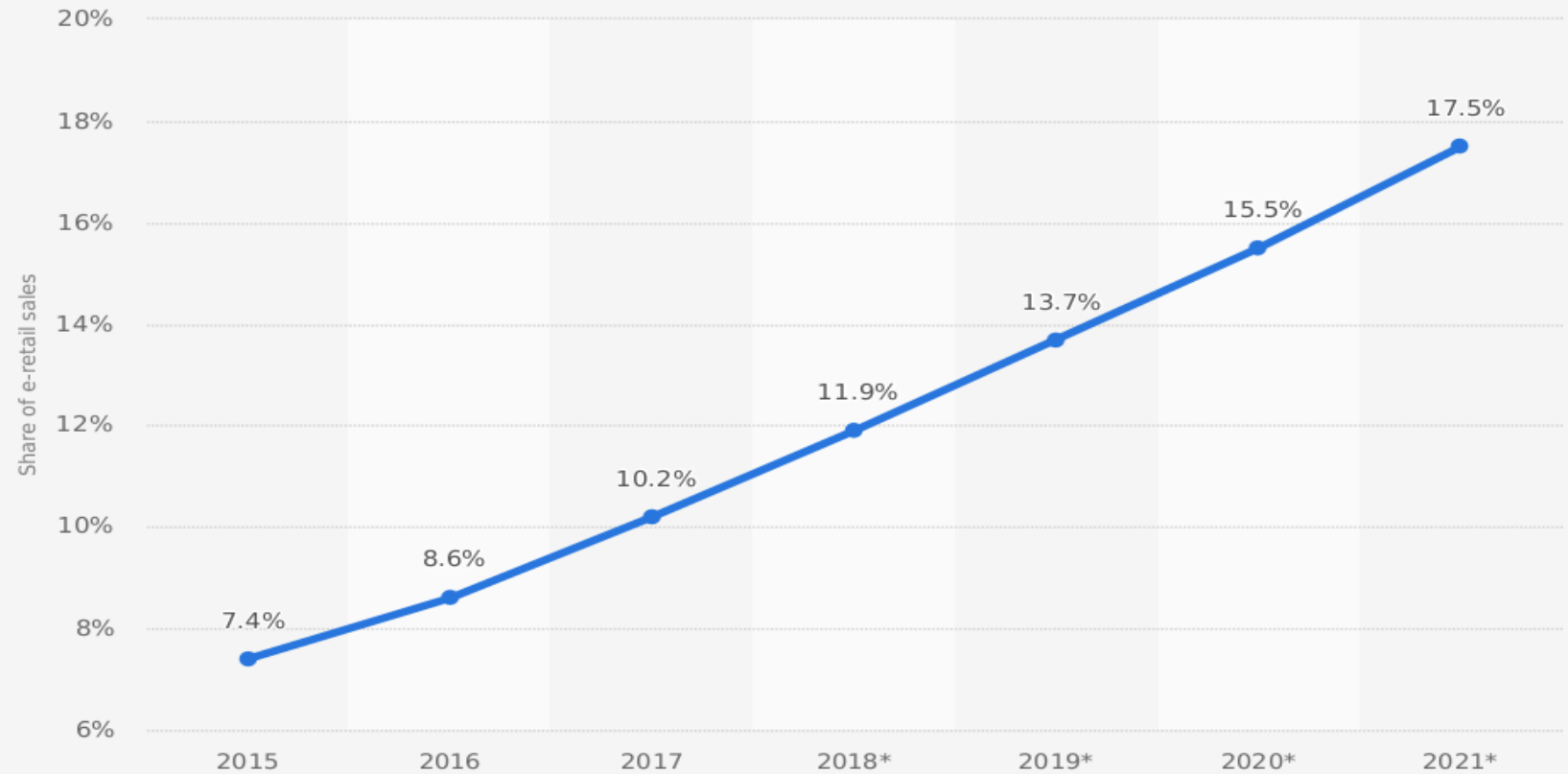
- **E-commerce** is the activity of electronically buying or selling products on online services or over the Internet.
- Top 5 e-commerce companies in the world:



- The ones we might know and use:
  - Allegro (18th largest)
  - Etsy (20th largest)
  - Zalando (21st largest)

# Rise of e-commerce

**E-commerce share of total global retail sales from 2015 to 2021**



**Sources**

eMarketer; Website (retailtechnews.com)  
© Statista 2018

**Additional Information:**

Worldwide; eMarketer; 2015 to 2017

## Recommendation systems

- “People don’t know what they want until you show it to them” - Steve Jobs
- A fix for the lack of help by the shops staff or an annoying ad?
- Most commonly evaluates users shopping history and viewing behaviour to recommend things to buy.
- Cross-selling vs **upselling**.

## Our Project

- Automatic methods for measuring similarity between products on multilevel dimensions
- Taxonomy
- Extracting crucial information from descriptions and titles

# **BERT-based similarity learning for product matching**

**Janusz Tracz<sup>1</sup>, Piotr Wójcik<sup>1</sup>, Kalina Jasinska-Kobus<sup>1, 2</sup>,  
Riccardo Belluzzo<sup>1</sup>, Robert Mroczkowski<sup>1</sup>, Ireneusz Gawlik<sup>1, 3</sup>**

<sup>1</sup> ML Research at Allegro.pl

<sup>2</sup> Poznan University of Technology

<sup>3</sup> AGH University of Science and Technology

`{janusz.tracz, piotr.wojcik, kalina.kobus, riccardo.belluzzo,  
robert.mroczkowski, ireneusz.gawlik}@allegro.pl`

# BERT-based similarity learning for product matching

## Generalized zero-shot multi-class classification

Categorizing instance into multiple classes, some of which have not been part of the model's training data.

### Available data

- Title
- Description
- Category
- List of attributes (attribute name, value, and unit)

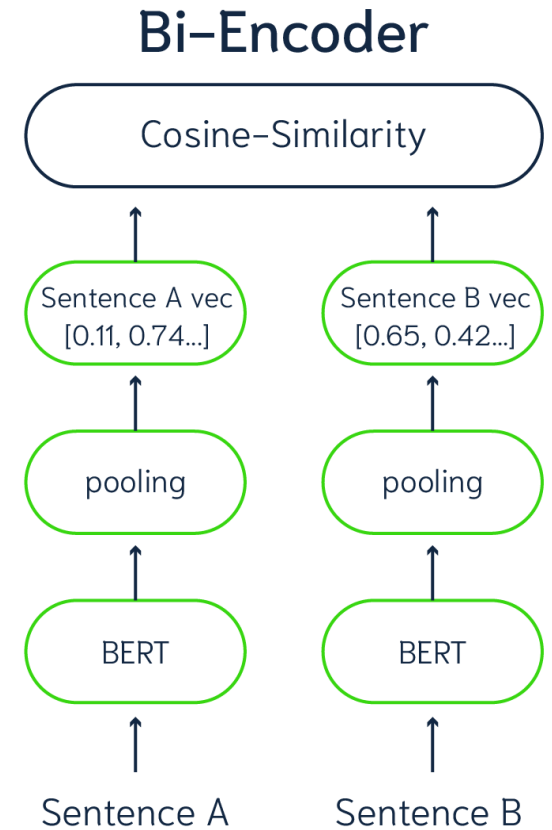
# BERT-based similarity learning for product matching

## Bi-Encoder architecture

- Products are represented as text
- A previously trained transformer is applied
- A distance between embeddings is calculated as similarity between instances

Used transformer: BERT

Used distance: cosine distance





# BERT-based similarity learning for product matching

## Textual representation

Each of the products was represented as a concatenation of:

- title
- attribute values
- attribute units

Description and attribute names were omitted as they deteriorated model performance.

Additionally, an assumption was made that only products from the same category were compared.

# BERT-based similarity learning for product matching

## Similarity learning with triplet loss objective

Training data consists of triples in the form of

$$(o, p^+, p^-)$$

with the elements being the anchor, a matching product, and a non-matching product.

Then, the transformer is adjusted to minimize the following loss function:

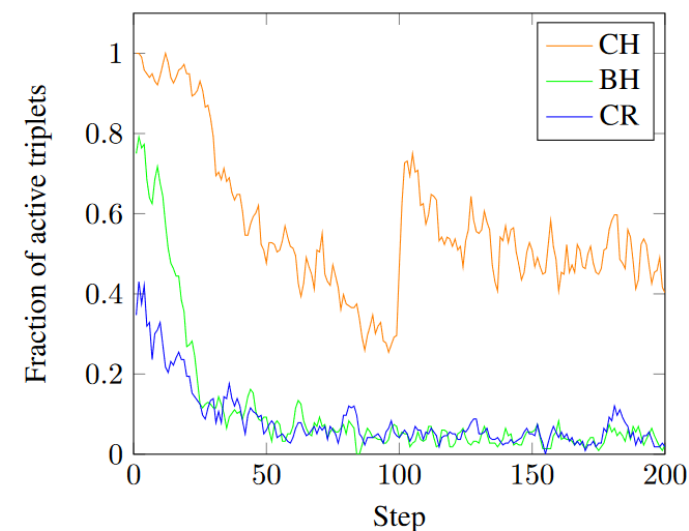
$$\mathcal{L}(o, p^+, p^-) = \max(0, m + d(\mathcal{E}_\theta(o), \mathcal{E}_\phi(p^+)) - d(\mathcal{E}_\theta(o), \mathcal{E}_\phi(p^-)))$$

# BERT-based similarity learning for product matching

## Batch construction strategy

In the article, different strategies were considered when selecting the negative match to minimize inactive triplets:

- randomly from a category (CR)
- most similar product from the non-matching products in the sampled batch (BH)
- most similar product from the non-matching products in the entire category (CH)



(a) Active triplet fraction for HerBERT initialised model for different negative item selection strategies.

# BERT-based similarity learning for product matching

## Results

	Available matches	Products
CULTURE	300K	800K
ELECTRONICS	200K	400K
BEAUTY	300K	200K

Table 1: Datasets used for our experiments.

	CULTURE	ELECTRONICS	BEAUTY
BOW	0.8863	0.8032	0.7687
HerBERT-NFT	0.8206	0.6716	0.5542
HerBERT	0.9550	0.8580	0.9064
eComBERT-NFT	0.8208	0.6755	0.6127
eComBERT	<b>0.9777</b>	<b>0.8840</b>	<b>0.9219</b>




Table 2: Test accuracy per each dataset. NFT stands for *non-finetuned*.

---

# MULTILINGUAL TRANSFORMERS FOR PRODUCT MATCHING – EXPERIMENTS AND A NEW BENCHMARK IN POLISH

---

A PREPRINT

**Michał Moźdzzonek<sup>1</sup>**,  **Anna Wróblewska<sup>1</sup>**,  **Sergiy Tkachuk<sup>2</sup>**,  **Szymon Łukasik<sup>2,3</sup>**

<sup>1</sup>Faculty of Mathematics and Information Science, Warsaw University of Technology, Warsaw, Poland

Email: {michal.mozdzonek,anna.wroblewska1}@pw.edu.pl

<sup>2</sup>Systems Research Institute, Polish Academy of Sciences

ul. Newelska 6, 01-447 Warsaw, Poland

Email: {stkachuk,slukasik}@ibspan.waw.pl

<sup>3</sup>Faculty of Physics and Applied Computer Science, AGH University of Science and Technology

al. Mickiewicza 30, 30-059 Kraków, Poland

Email: slukasik@agh.edu.pl

# mBERT and XML- RoBERTa transfer learning

- Product matching problem
- Transfer learning for data in different languages
- Web Data Commons dataset (4 categories, sizes: small, medium, large)
- Own Polish product matching dataset
- Running pre-trained models and performance comparison

# mBERT and XML- RoBERTa transfer learning

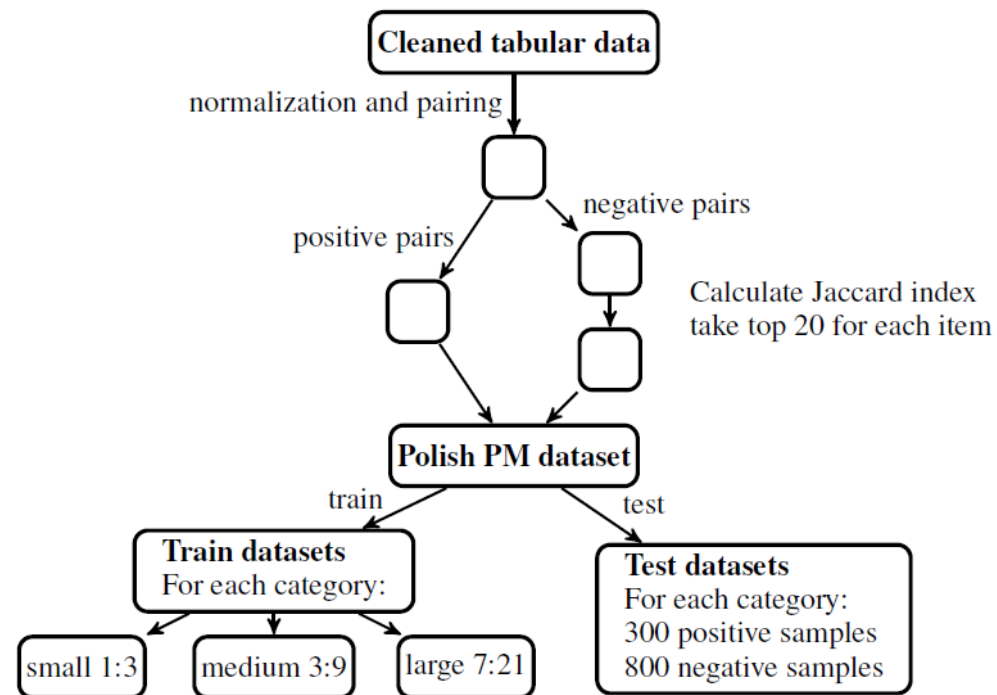


Figure 2: The process of creating the Polish PM datasets. In each training set, the ratio of positive to negative samples is 1:3.

# mBERT and XML- RoBERTa transfer learning

- Selecting the title column only and concatenating it with token markers
- HuggingFace Transformers library
- Two types of models: : mBERT and XLM-RoBERTa
- Pre-trained models on Wikipedia articles in about 100 languages
- Models run on both WDC and Polish datasets



# mBERT and XML- RoBERTa transfer learning

Table 7: F1 scores for models trained on English WDC datasets. Mean value and standardized error (confidence level 95%) for each dataset were calculated from 4 samples. For further information on how standardized error was calculated see Section 5.

dataset type	dataset size	mBERT	XML-RoBERTa	Ditto (reported in Li et al. [2020])	WDC-Deepmatcher (reported in Peeters et al.)
Cameras	small	<b>82.13</b> ( $\pm 4.70$ )	81.96( $\pm 7.75$ )	80.89	68.59
	medium	87.86( $\pm 2.04$ )	<b>88.11</b> ( $\pm 4.22$ )	88.09	76.53
	large	90.88( $\pm 2.28$ )	<b>92.36</b> ( $\pm 0.76$ )	91.23	87.19
	xlarge	-	-	<b>93.78</b>	89.21
Computers	small	<b>86.43</b> ( $\pm 3.69$ )	81.10( $\pm 13.40$ )	80.76	70.55
	medium	<b>90.13</b> ( $\pm 1.89$ )	88.69( $\pm 2.19$ )	88.62	77.82
	large	92.48( $\pm 2.33$ )	<b>93.71</b> ( $\pm 0.77$ )	91.70	89.55
	xlarge	-	-	<b>95.45</b>	90.80
Shoes	small	<b>79.20</b> ( $\pm 7.89$ )	74.98( $\pm 13.36$ )	75.89	73.86
	medium	<b>84.11</b> ( $\pm 3.40$ )	81.30( $\pm 8.21$ )	82.66	79.48
	large	90.28( $\pm 2.36$ )	<b>91.26</b> ( $\pm 2.09$ )	88.07	90.39
	xlarge	-	-	90.10	<b>92.61</b>
Watches	small	<b>87.31</b> ( $\pm 1.64$ )	83.78( $\pm 4.38$ )	85.12	66.32
	medium	<b>91.17</b> ( $\pm 4.21$ )	89.50( $\pm 3.69$ )	91.12	79.31
	large	93.52( $\pm 2.63$ )	93.62( $\pm 0.67$ )	<b>95.69</b>	91.28
	xlarge	-	-	<b>96.53</b>	93.45

# mBERT and XML- RoBERTa transfer learning

Table 6: F1 scores for mBERT and XLM-RoBERTa trained on Polish datasets. Mean value and standardized error (confidence level 95%) for each dataset were calculated from **4** samples. For further information on how standardized error was calculated see Section [5](#).

dataset type	dataset size	mBERT	XLM-RoBERTa
Household chemistry (pl. chemia)	small	<b>85.73</b> ( $\pm 1.89$ )	83.15( $\pm 4.15$ )
	medium	<b>90.78</b> ( $\pm 3.03$ )	89.03( $\pm 5.96$ )
	large	<b>93.25</b> ( $\pm 1.77$ )	92.52( $\pm 1.77$ )
Drinks (pl. napoje)	small	<b>85.17</b> ( $\pm 1.61$ )	84.43( $\pm 7.16$ )
	medium	<b>88.98</b> ( $\pm 2.63$ )	88.44( $\pm 2.88$ )
	large	89.39( $\pm 2.12$ )	<b>89.93</b> ( $\pm 3.99$ )
All	small	<b>85.73</b> ( $\pm 1.96$ )	84.67( $\pm 9.03$ )
	medium	<b>90.78</b> ( $\pm 1.13$ )	88.63( $\pm 2.79$ )
	large	91.41( $\pm 3.17$ )	<b>91.61</b> ( $\pm 1.39$ )

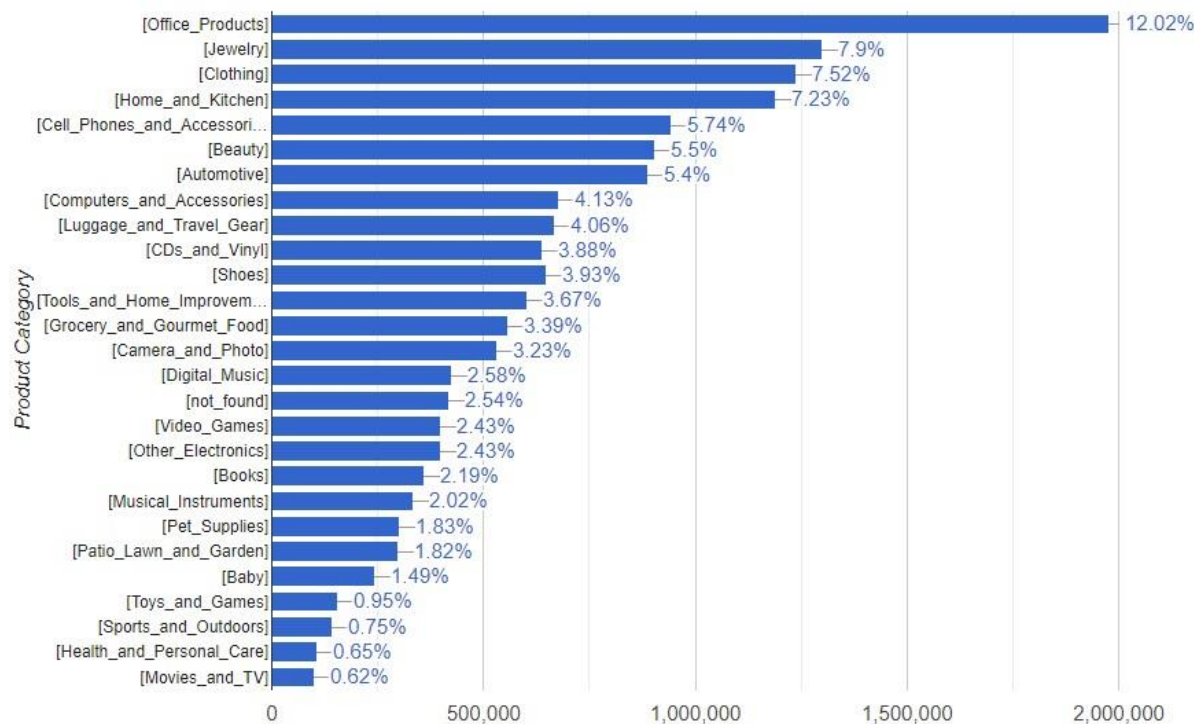
The slide features a white background with two large teal geometric shapes. On the left, a teal triangle points towards the center. On the right, a teal trapezoid is positioned, also pointing towards the center. The text 'Solution Concept' is centered between these two shapes.

# Solution Concept

# Dataset

## "Web Data Commons - Training Dataset and Gold Standard for Large-Scale Product Matching"

- 16 million English-language offers sourced from a wide array of 79 thousand websites.
- Includes product categorization based on Amazon product data and TF-IDF scores for 26 product categories.
- Each offer was assigned to one of 26 categories



# Dataset

## GOLD STANDARD

Category	# positive pairs	# negative pairs	% title	% description	% brand	% price	% specTableContent
Computers	300	800	100	82	42	11	22
Cameras	300	800	100	73	25	3	7
Watches	300	800	100	71	15	1	7
Shoes	300	800	100	70	8	1	2
<b>All</b>	<b>1200</b>	<b>3200</b>	<b>100</b>	<b>74</b>	<b>23</b>	<b>4</b>	<b>10</b>

# Proposed approach

## Products data processing

- Exploratory Data Analysis
- Select and/or augment product attributes that will be concatenated to create input for model.

## Implementation of modified product similarity pipeline

- Fine-tune BERT-based model (RoBERTa, DistilBERT) or different embedding model, like BGE (#1 on HF MTEB) using bi-encoder framework.
- Experiment with loss/distance functions.
- Performance tests and evaluation

The image features a white background with two large teal geometric shapes. On the left, a large teal parallelogram is positioned. On the right, a teal triangle points towards the center. Centered between these shapes is the text "THANK YOU FOR ATTENTION" in a dark gray, sans-serif font, arranged in two lines.

THANK YOU FOR  
ATTENTION