



Decoding Passion: Discovering User Interest Profiles with Hybrid AI

AI Hub @ Allegro

Anna Mikołajczyk
Anna Pintara

Data Scientist
Data Scientist

allegro

Discovering user interests profiles

Our hobbies



belly dance

healthy lifestyle

board games



motherhood

skiing

contemporary dance

Hobby as a gateway to inspiration

Moving from individual product data to a hobby insight.

BEHAVIORAL TRACE



Item #329487

Running Shoes

"User bought shoes."

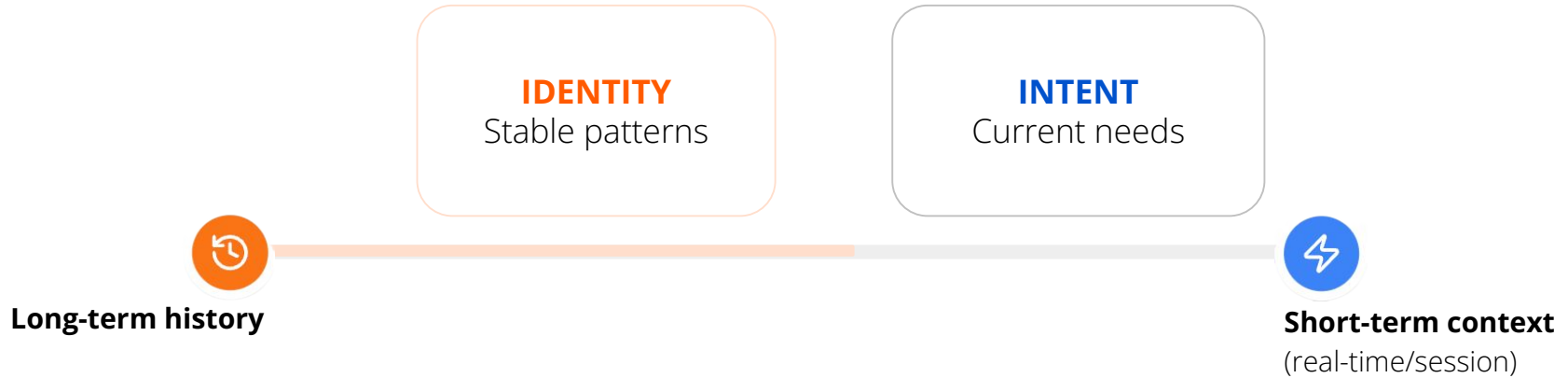
HOBBY INSIGHT



RUNNING ENTHUSIAST

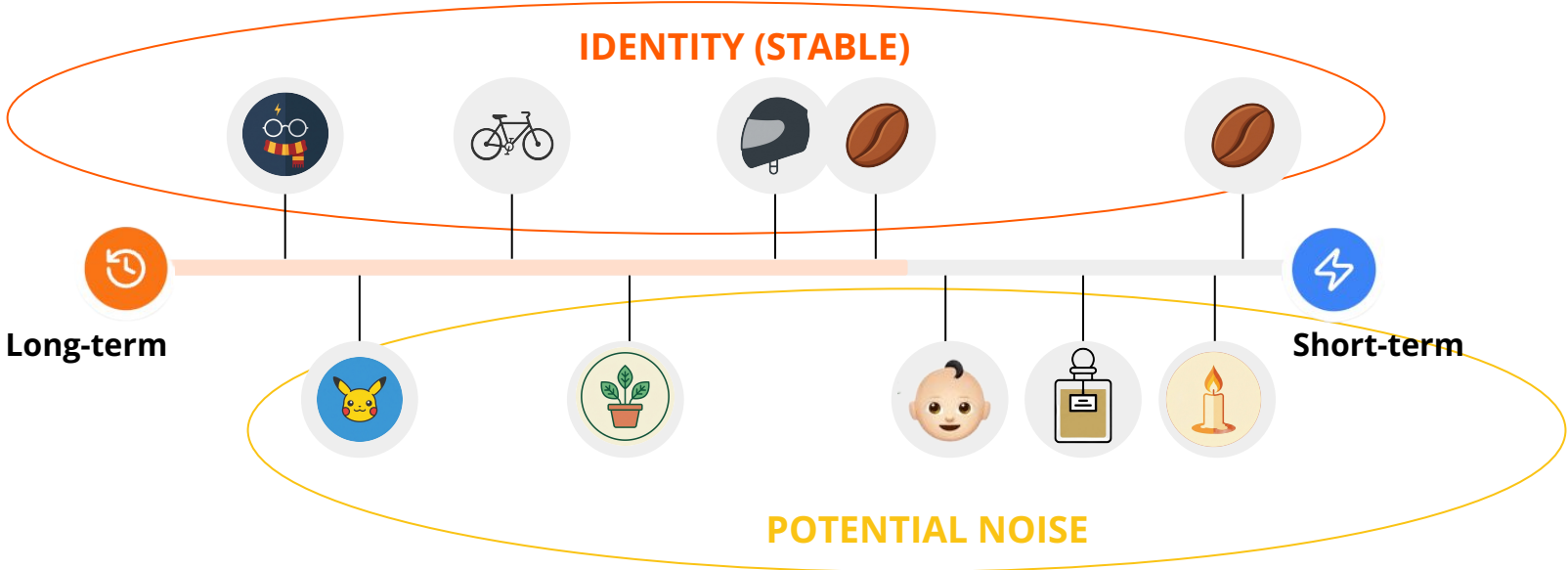
The Horizon of Understanding

The longer we know you the more we understand users **identity**.



The horizon of understanding (illustration)

Distinguishing stable identity from transactional noise across 730 days.



The Definition of Hobby



Lifestyle

PARENTING, GARDENING



Activity

MUSIC, FISHING



Profession

HAIRDRESSER,
INFLUENCER



Lifestage

RENOVATING, FESTIVALS



Style

RETRO, MINIMALIST



Cultural affinity

POTTERHEADS, LEGO

The approach: hybrid pipeline

STEP 01



HOBBY GENERATION

Assigning topics to millions of products using textual data.

TECHNOLOGY STACK

Generative AI



STEP 02

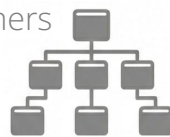


HOBBIES STRUCTURATION

Organizing raw tags into coherent, hierarchical tree-based structure.

TECHNOLOGY STACK

Clustering
+ Sentence Transformers
+ Generative AI



STEP 03



USERS PROFILING

Aggregating history to find the strongest interests using behavioral dynamics.

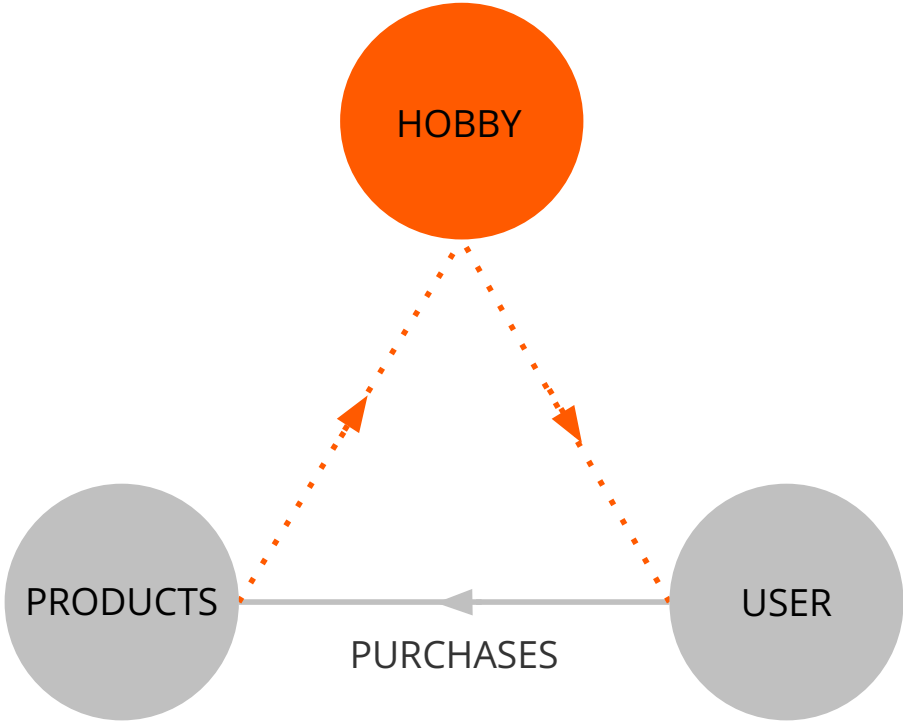
TECHNOLOGY STACK

RFM Scoring



WHAT ARE THE HOBBY SIGNALS?

How do we get hobbies?



Text is all we need

Allegro / Home and Garden / Equipment / Rugs and Carpets / Doormats / F

Brand: Grupoerik | Condition: new

Doormat Harry Potter Platform 9 3/4 Doormats 60x40 cm

4.98 51 ratings and 11 reviews | 10 people purchased recently



Harry Potter Platform 9 3/4

Make your entrance not only clean, but also stylish. This doormat is more than just a practical piece - it's a way to show off the character of your home right from the doorstep! The wipers are made of natural coconut fibers, with a non-slip rubber base. Not only do they look great, but they also keep things organized. They are suitable for both indoor and outdoor use.

Catalog number: FGE0006

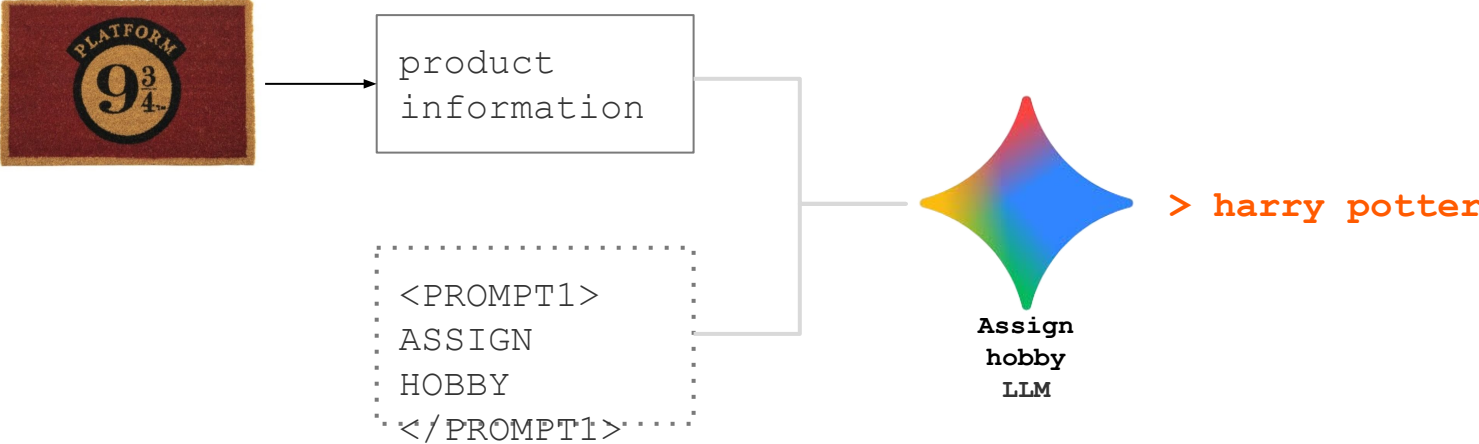
Dimension: 60x40 cm

Thickness: 1,5 cm

Brand, producer: Harry Potter Grupo Erik

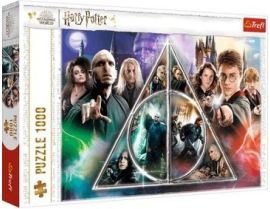



- 1 Name
- 2 Description
- 3 Category path
- 4 Images
- 5 Parameters

Hobby generation

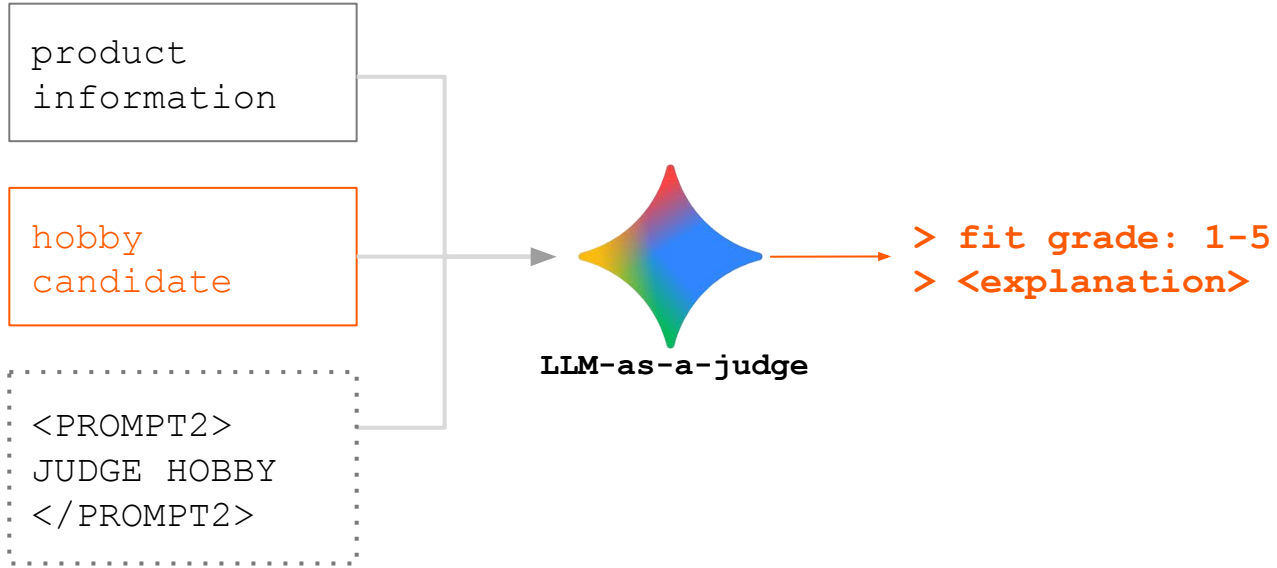


We use Large Language Model to interpret product textual info. Each item receives a **hobby/topic tag**.


Can we trust an LLM?

	<p>Name</p> <p><i>Puzzle 1000 Deathly Hallows Trefl</i></p>	<p>Hobby</p> <p>harry potter</p>		<p>Name</p> <p>Anti-smog mask size S</p>	<p>Hobby</p> <p>honey extraction</p>
	<p>Name</p> <p><i>Born2Be Black Women's Boots 38 Transitional Rubber</i></p>	<p>Hobby</p> <p>festival wear</p>		<p>Name</p> <p>SILOO ball medal hanger</p>	<p>Hobby</p> <p>achievements</p>

Hobby judge to the rescue




Hobby judge verdict

Name	Hobby	Fit grade	Explanation
	honey extraction	1	No clear connection.
	achievements	1	Abstract concept, not a hobby.

How many hobbies per product?

	Name	Hobby
	<p><i>Puzzle 1000 Deathly Hallows Trefl</i></p>	<p>harry potter</p> <p>problem solving</p> <p>family activity</p> <p>fantasy fandom</p>

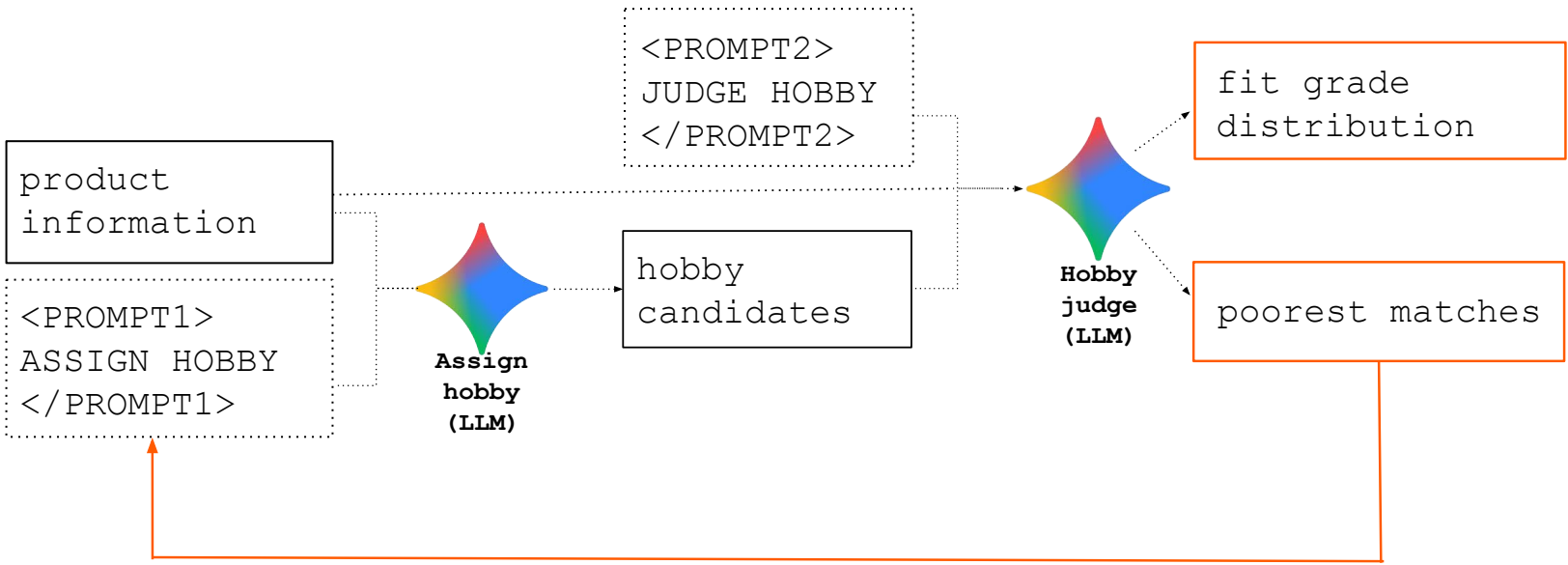
Which input data is relevant?

	Name	Category Path	Description
	<p><i>Born2Be Black Women's Boots 38 Transitional Rubber</i></p>	<p><i>Allegro</i></p> <ul style="list-style-type: none"> > <i>Fashion</i> > <i>Clothes,</i> > <i>Footwear,</i> > <i>Accessories</i> > <i>Footwear</i> > <i>Women's</i> > <i>Wellingtons</i> 	<p><small>Born2Be Czarne Kaloze damskie 38 Przejeściowy Guma Kod produktu: 398643 Kolor: czarny Materiał: Guma Misek buta: Okrągły Ociepianie: Neouteployony Rodzaj obcasa: Płaski Rozmiarówka: Typ większe Sezon: Przejeściowy Styl: Casual Technic stopy: Neutralny Wkładka/wyciółka: Materiał włókiennicy Zapięcie: Niesamozatwarty TABELA ROZMIARÓW</small></p> <p><small>Rozmiar 38</small></p> <ul style="list-style-type: none"> [1] Długość wkładki 22,70 cm [2] Szerokość wkładki 7,50 cm [3] Wysokość obcasa 2,5 cm [6] Wysokość buta 36,5 cm [7] Obwód cholewki 38,00 cm [8] Obwód łydki 32,00 cm Długość stopy 22,20 cm <p><small>Rozmiar 37</small></p> <ul style="list-style-type: none"> [1] Długość wkładki 22,90 cm [2] Szerokość wkładki 7,50 cm [3] Wysokość obcasa 2,5 cm [6] Wysokość buta 36,5 cm [7] Obwód cholewki 38,00 cm [8] Obwód łydki 33,00 cm Długość stopy 22,40 cm <p><small>Rozmiar 36</small></p> <ul style="list-style-type: none"> [1] Długość wkładki 24,00 cm [2] Szerokość wkładki 8,00 cm [3] Wysokość obcasa 2,5 cm [6] Wysokość buta 36,5 cm [7] Obwód cholewki 38,00 cm [8] Obwód łydki 33,00 cm Długość stopy 23,30 cm <p><small>Rozmiar 39</small></p> <ul style="list-style-type: none"> [1] Długość wkładki 24,20 cm [2] Szerokość wkładki 8,00 cm [3] Wysokość obcasa 2,5 cm [6] Wysokość buta 36,5 cm [7] Obwód cholewki 38,00 cm [8] Obwód łydki 35,00 cm Długość stopy 23,70 cm <p><small>Rozmiar 40</small></p> <ul style="list-style-type: none"> [1] Długość wkładki 24,70 cm [2] Szerokość wkładki 8,50 cm [3] Wysokość obcasa 2,5 cm [6] Wysokość buta 37 cm [7] Obwód cholewki 40,00 cm [8] Obwód łydki 35,00 cm Długość stopy 24,20 cm

Call for experiments

- **How many** hobbies?
- Which **input data**?
- **Prompt modification** impact?
- **Language** choice?
- ... any other

How we managed to compare experiments



How big is our hobby base?

~115k

generated
hobbies

~10%

removed
by hobby judge

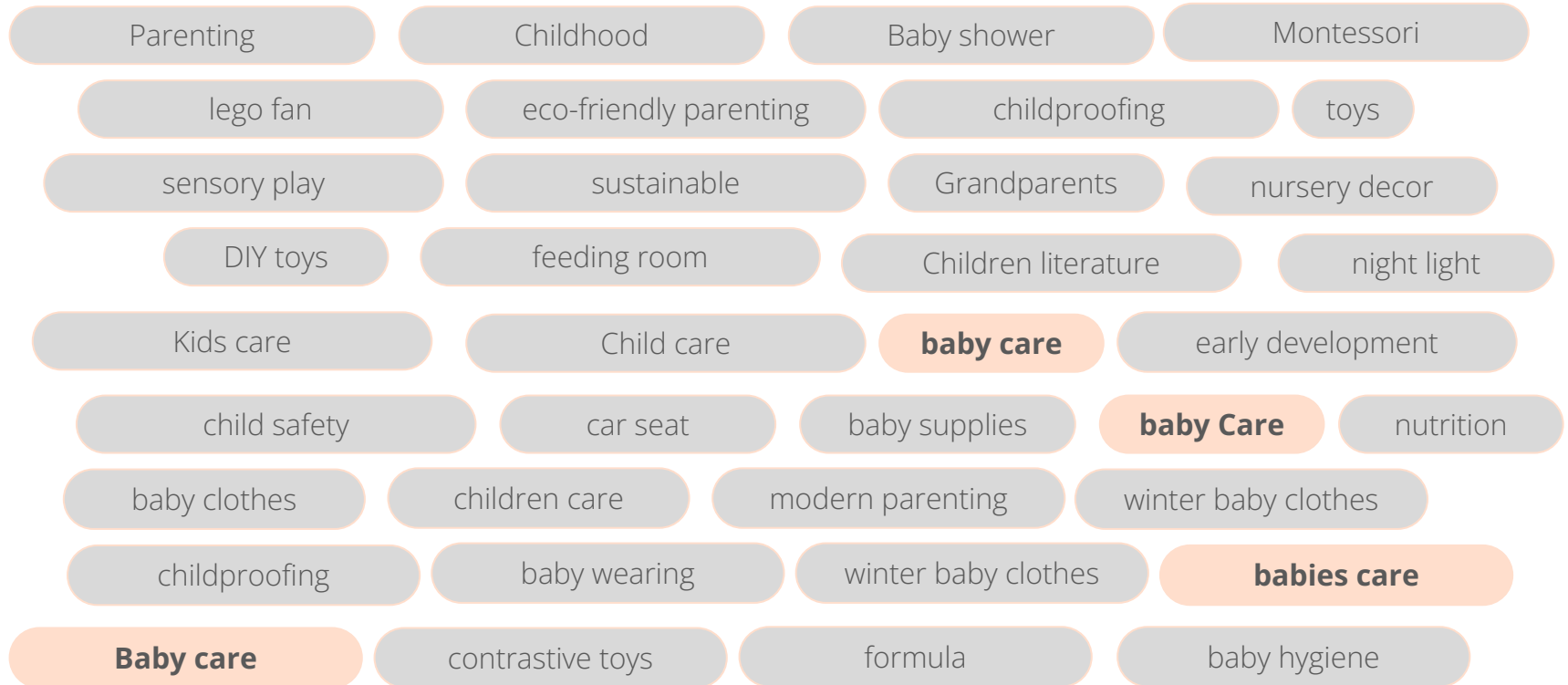
~100k

remaining
hobby candidates

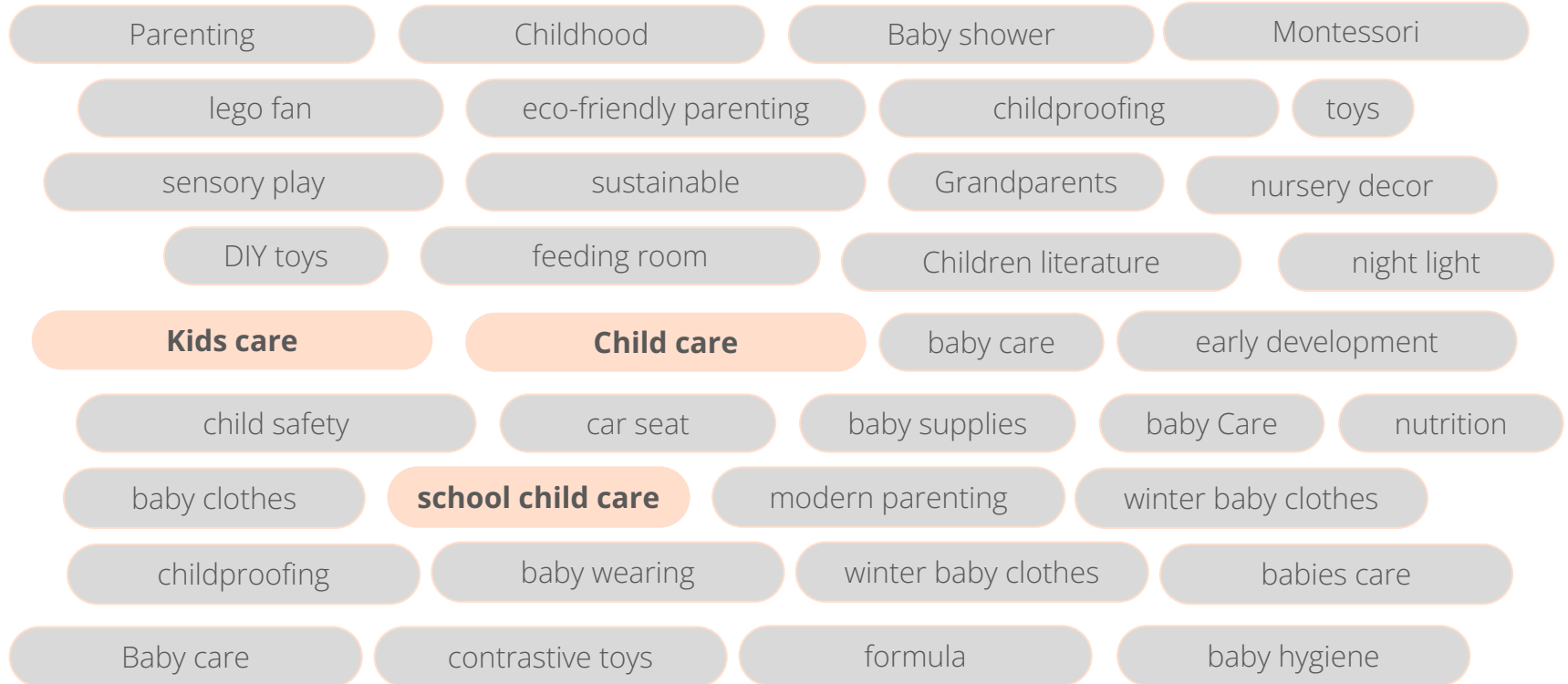
Highly unstructured data: 350 children topics



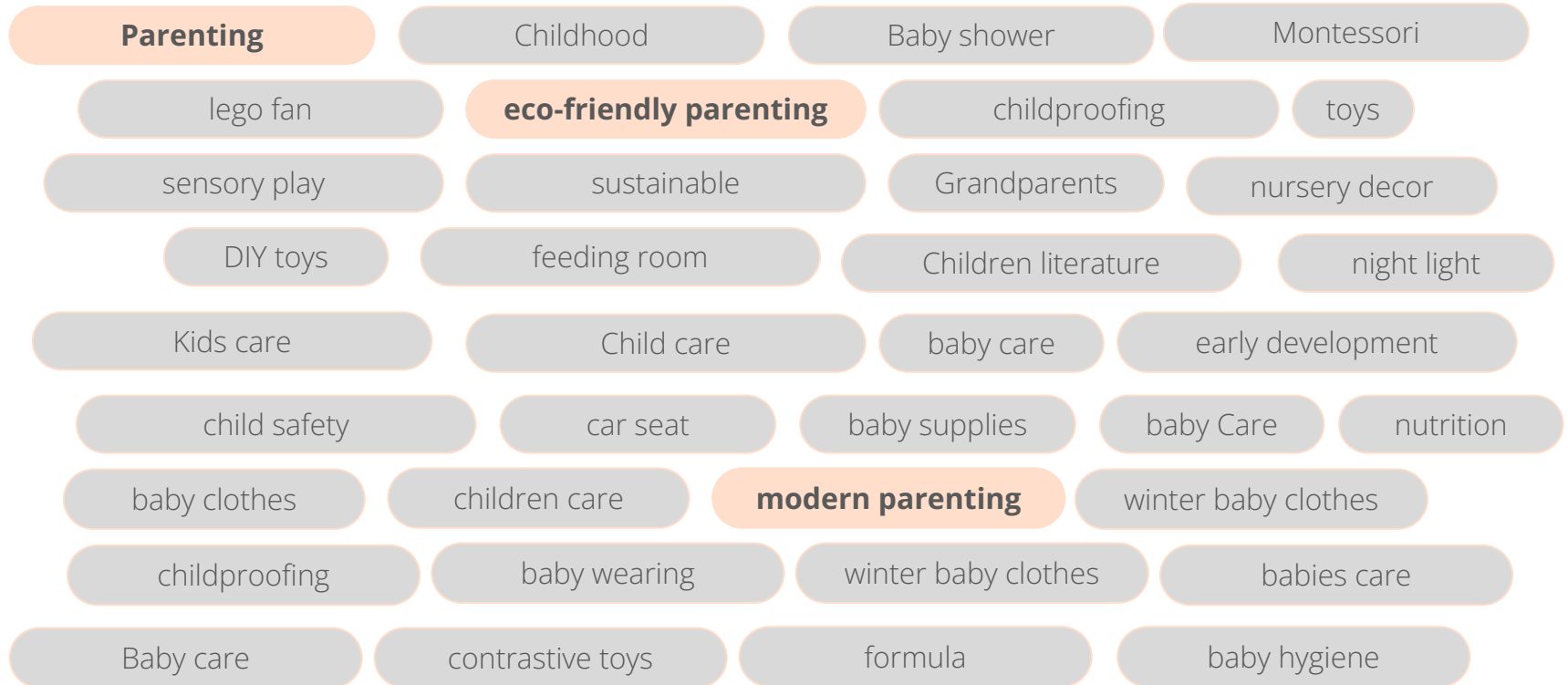
Highly unstructured data: Inconsistent forms



Highly unstructured data: synonyms

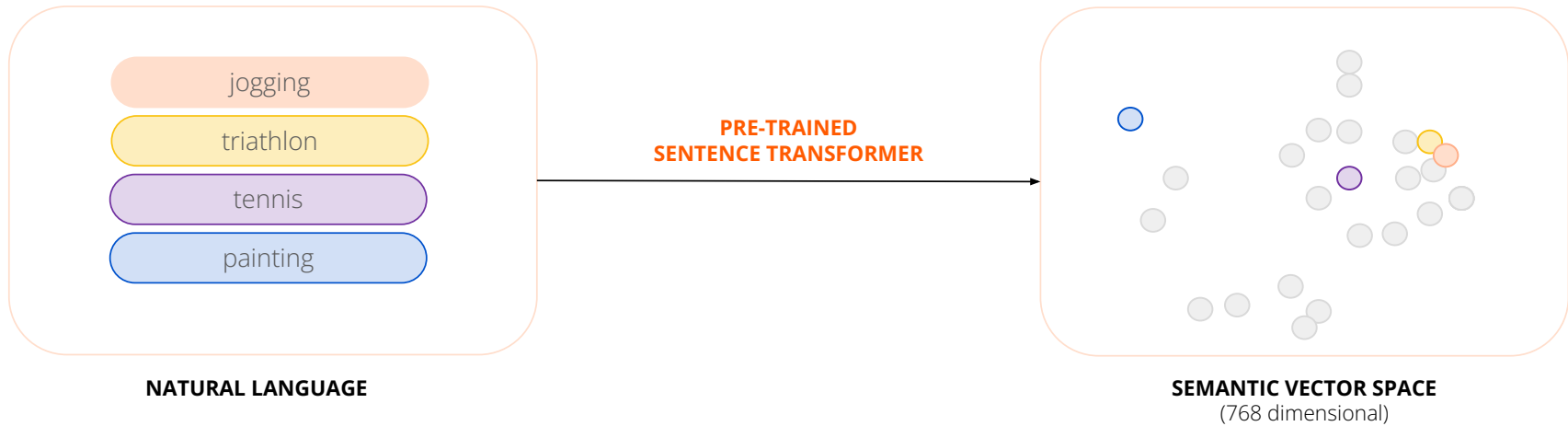


Highly unstructured data: inconsistent granularities



HOW TO ORGANIZE TOPIC CHAOS?

From natural language to mathematical language



Embeddings reflect semantic meaning, so we expect:

- > jogging to be closer to triathlon than tennis and painting
- > painting far away from all three other as not related to sport

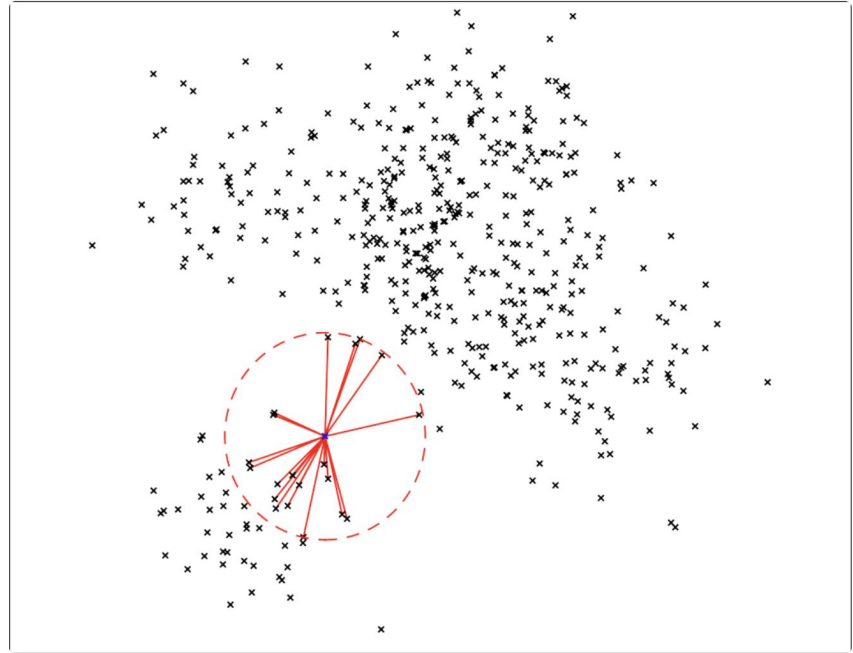
Nearest Neighbor Search

What is it?

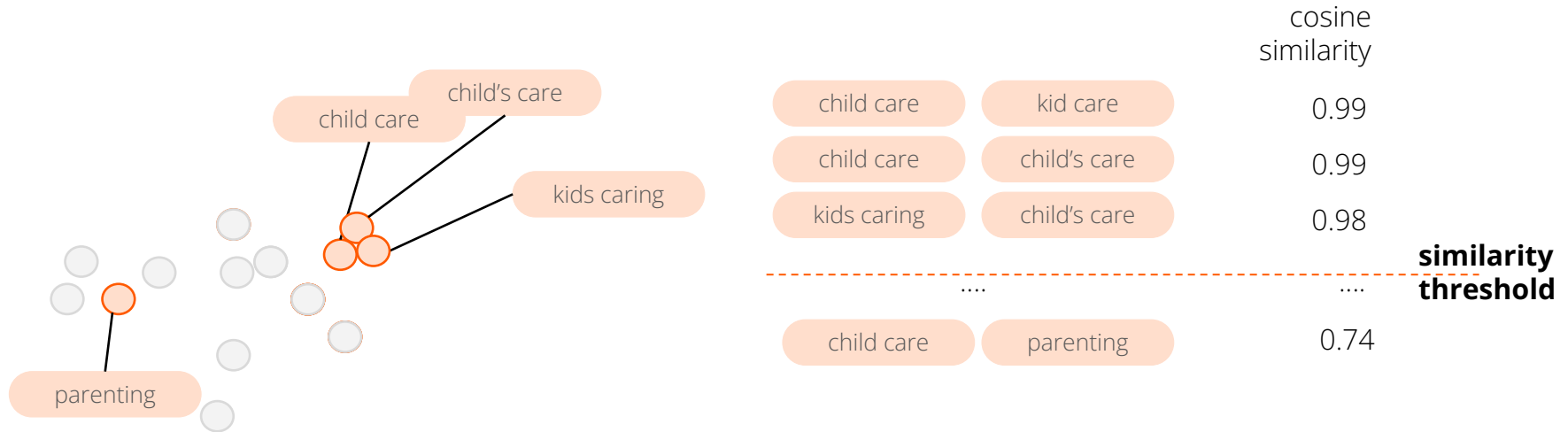
- Finds closest data points

Core Idea:

- calculate distances with all data points
- smallest distances = **neighbors**

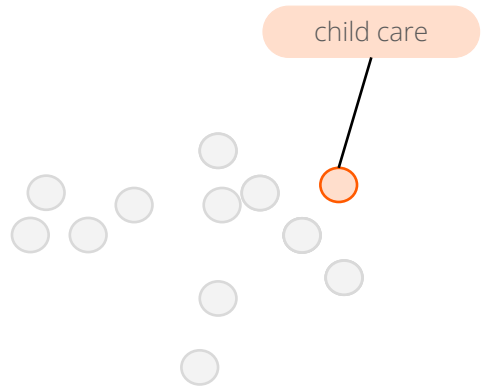
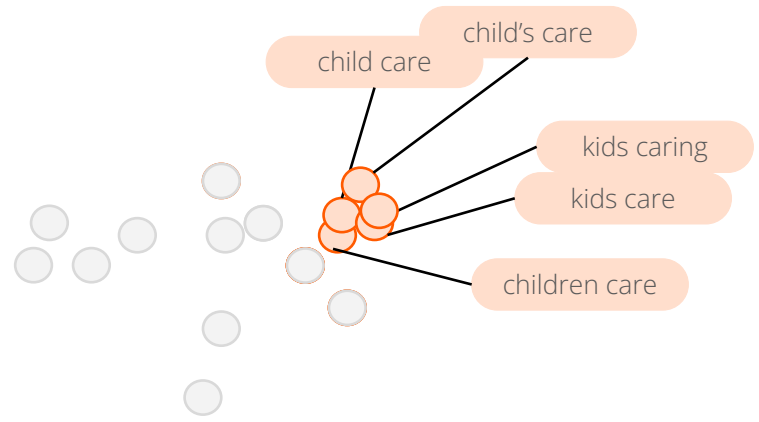


Hobbies pre-processing: synonyms & multiple forms



We want to get rid of the neighbors that are too close to each other (cosine similarity above threshold) by **merging them into only one!**

Hobbies pre-processing: synonyms & multiple forms



Pre-processed hobby base

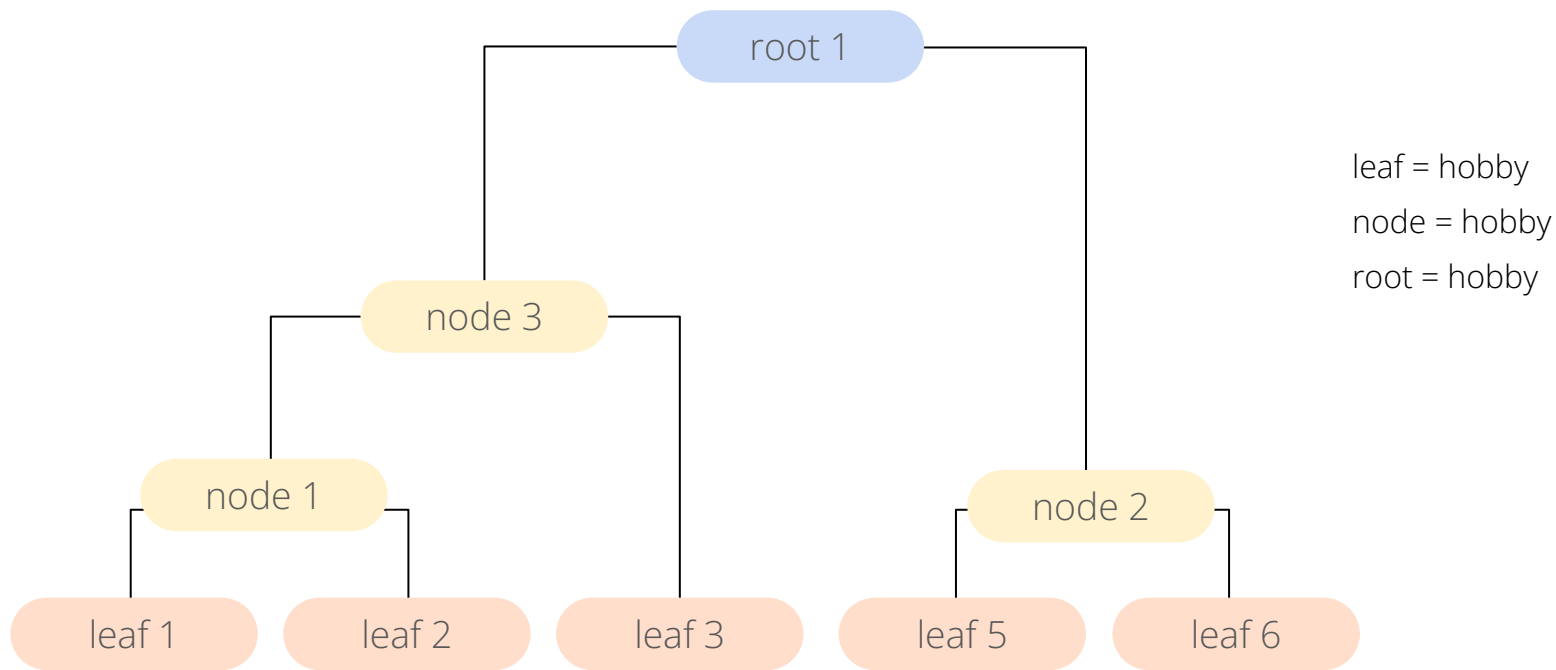
-70%

hobby set reduction
(pre-processing)

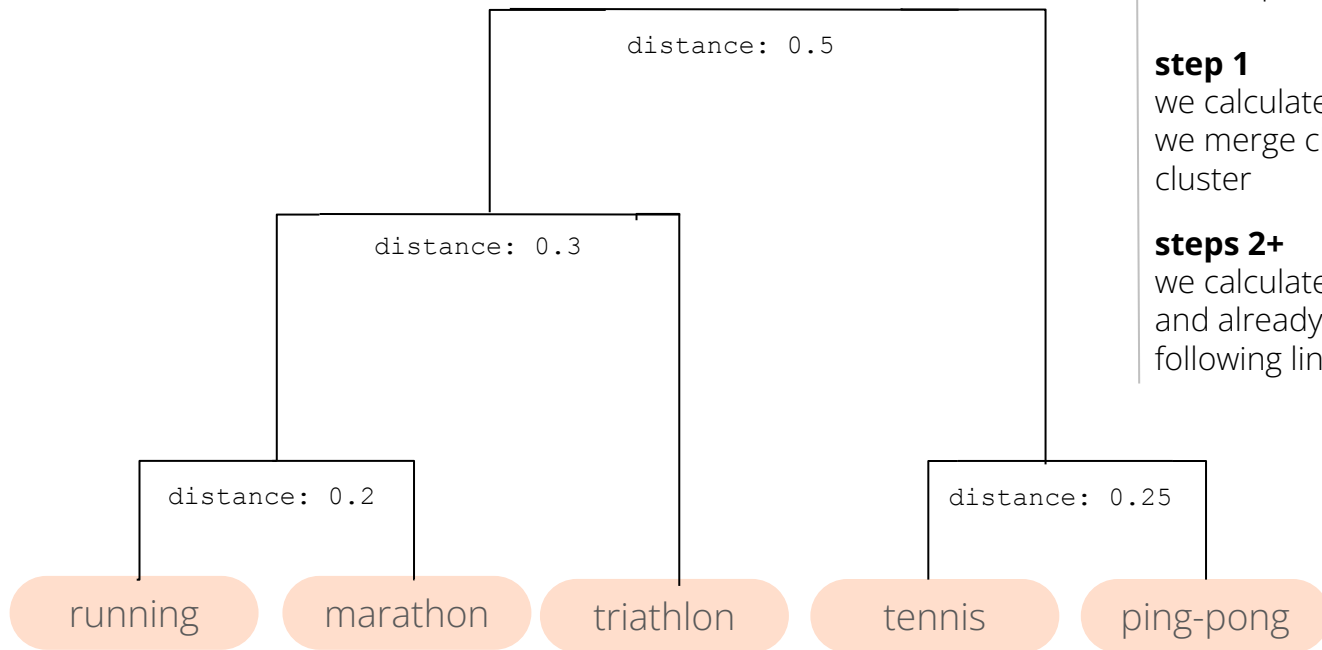
~30k

remaining hobbies for
grouping

Hierarchical clustering



Topic taxonomy - agglomerative bottom-up approach



step 0

each topic makes separate cluster

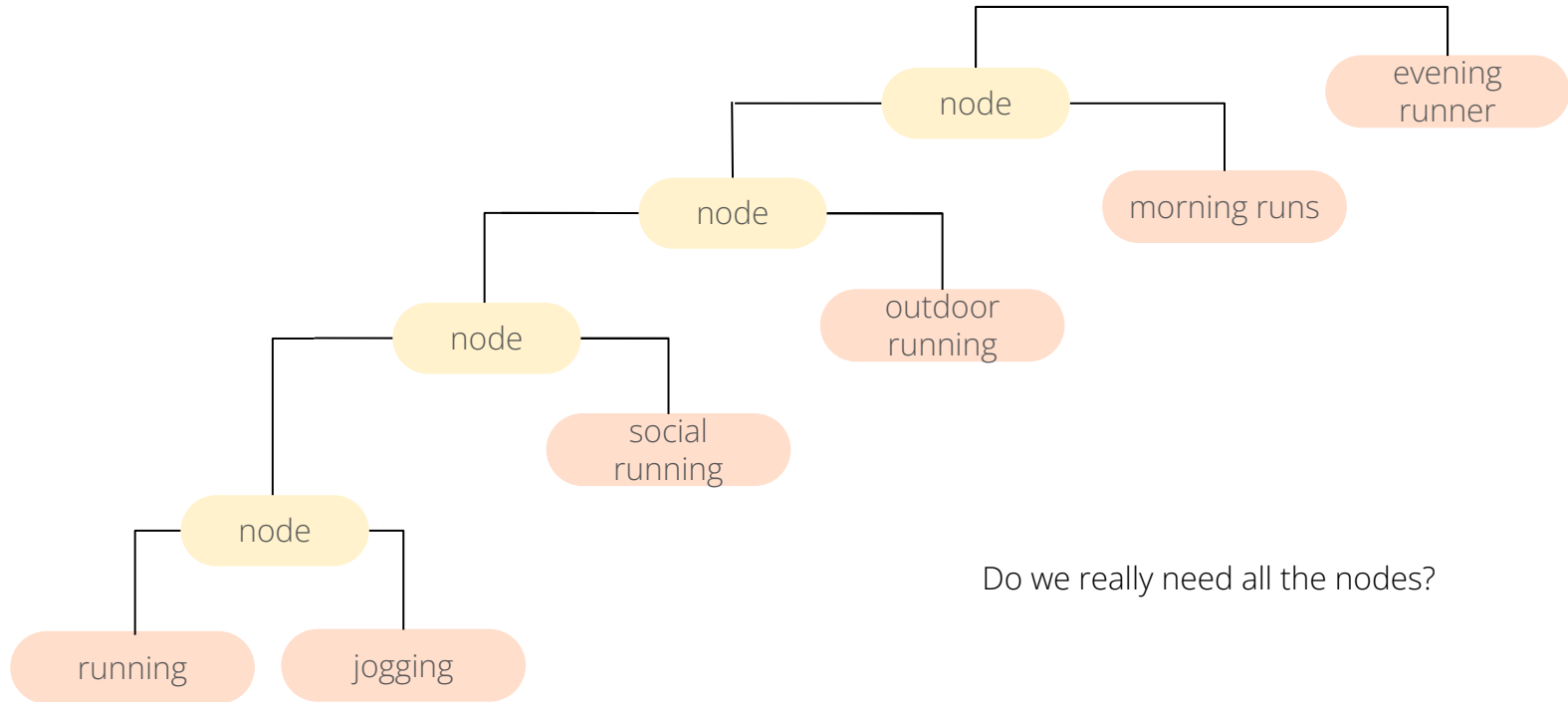
step 1

we calculate distance of each pair and we merge closest ones into one cluster

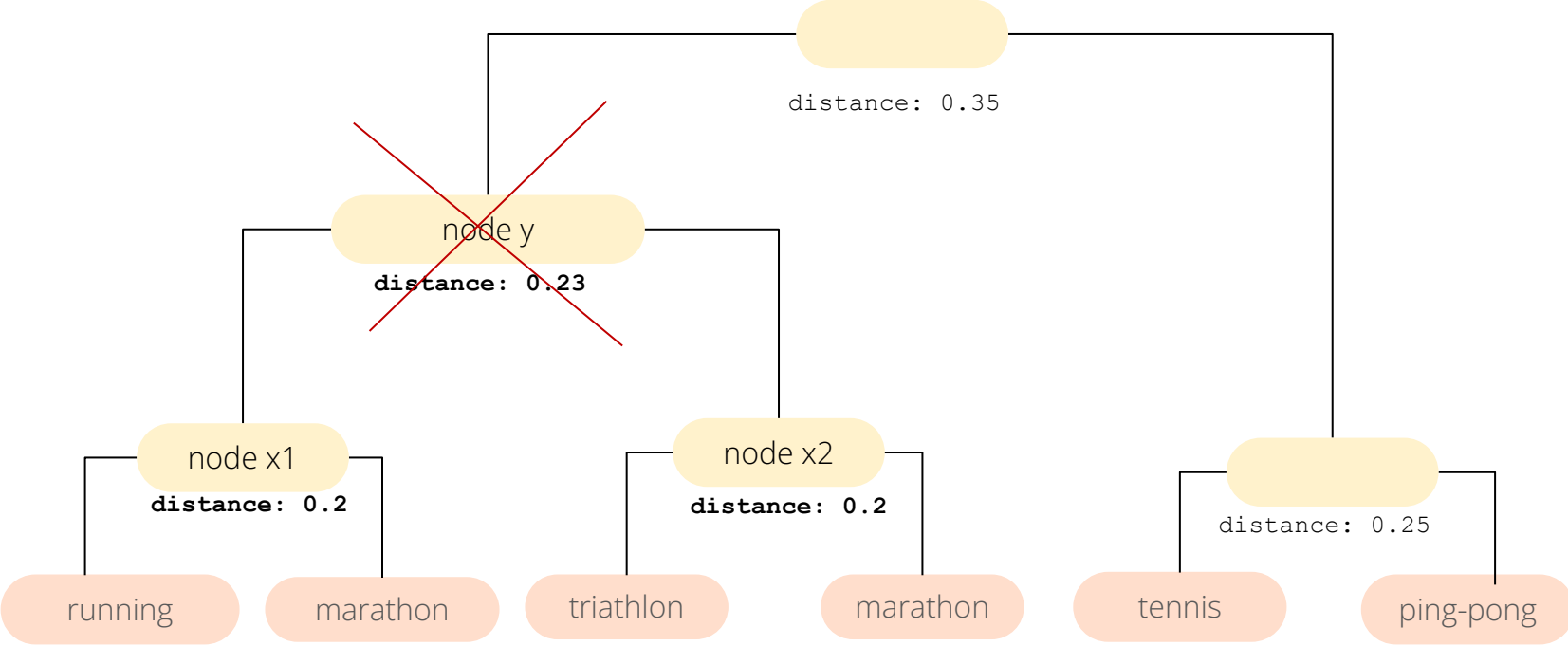
steps 2+

we calculate distance of remaining topics and already created clusters and make following links

Custom post-processing: nodes relevance



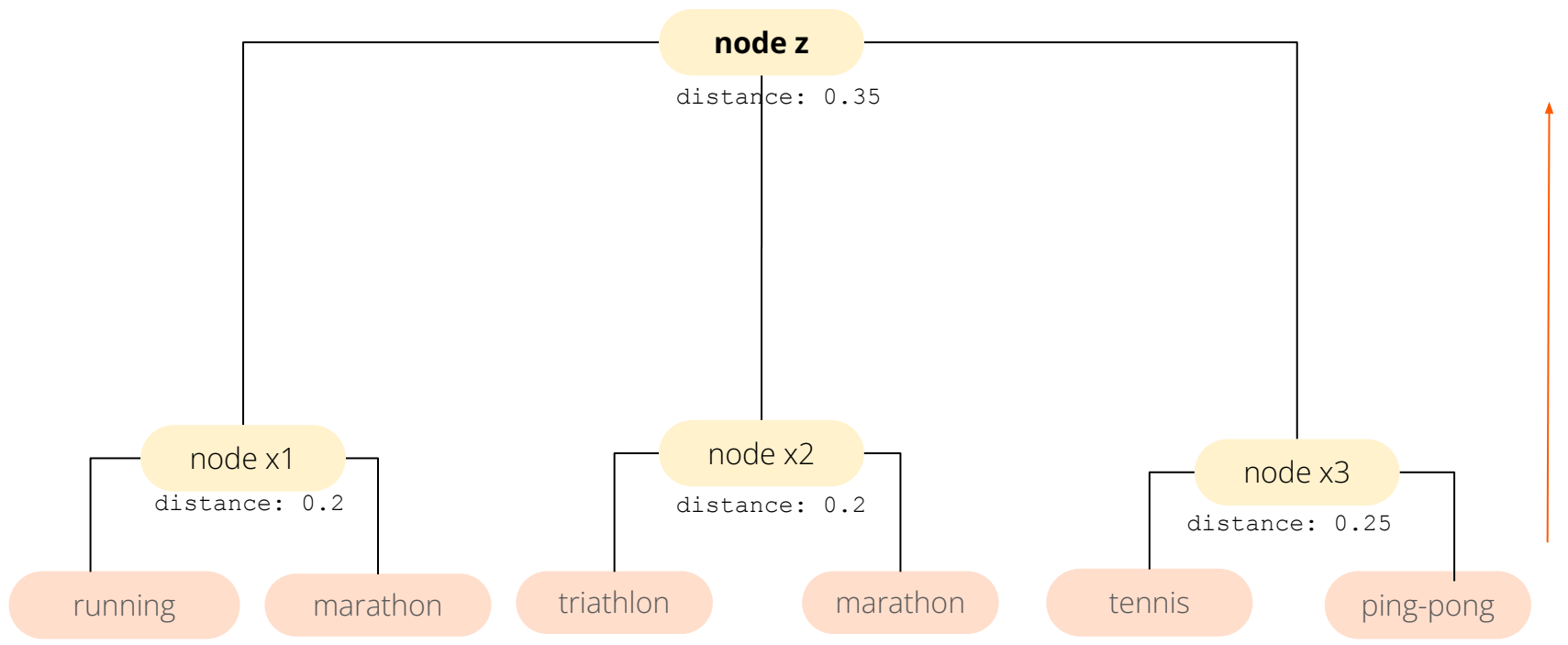
Custom post-processing: removing irrelevant nodes



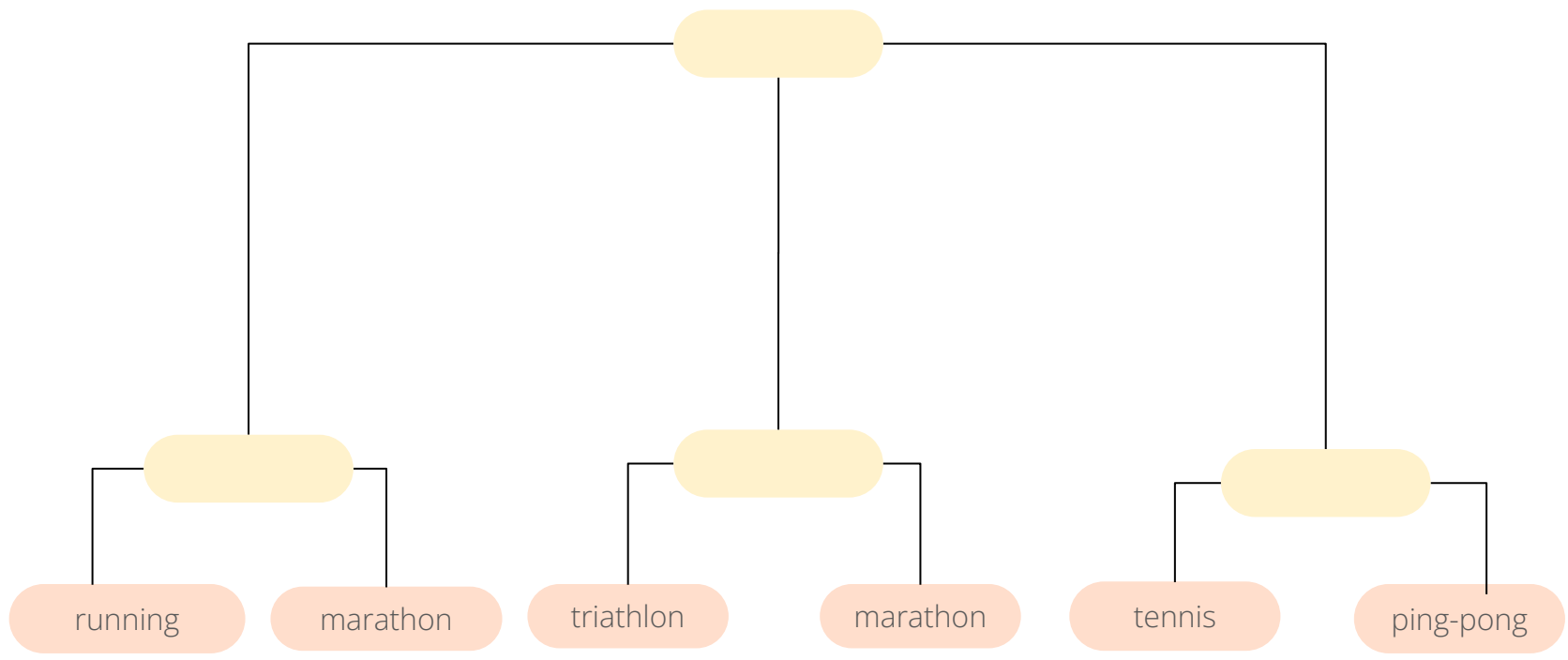
Is 0.23 greater than $\text{avg}(0.2, 0.2)$ by more than 20%?

NOPE! We remove node y!

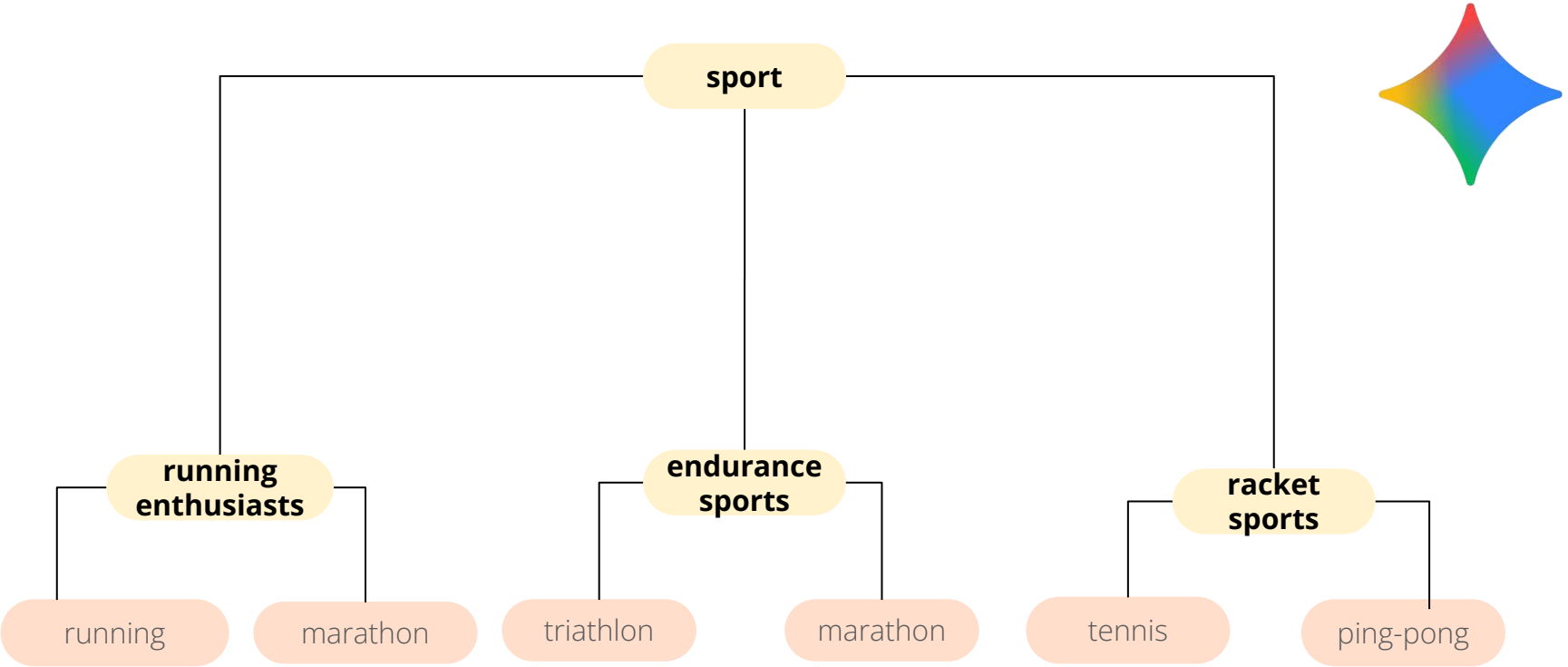
Custom post-processing direction



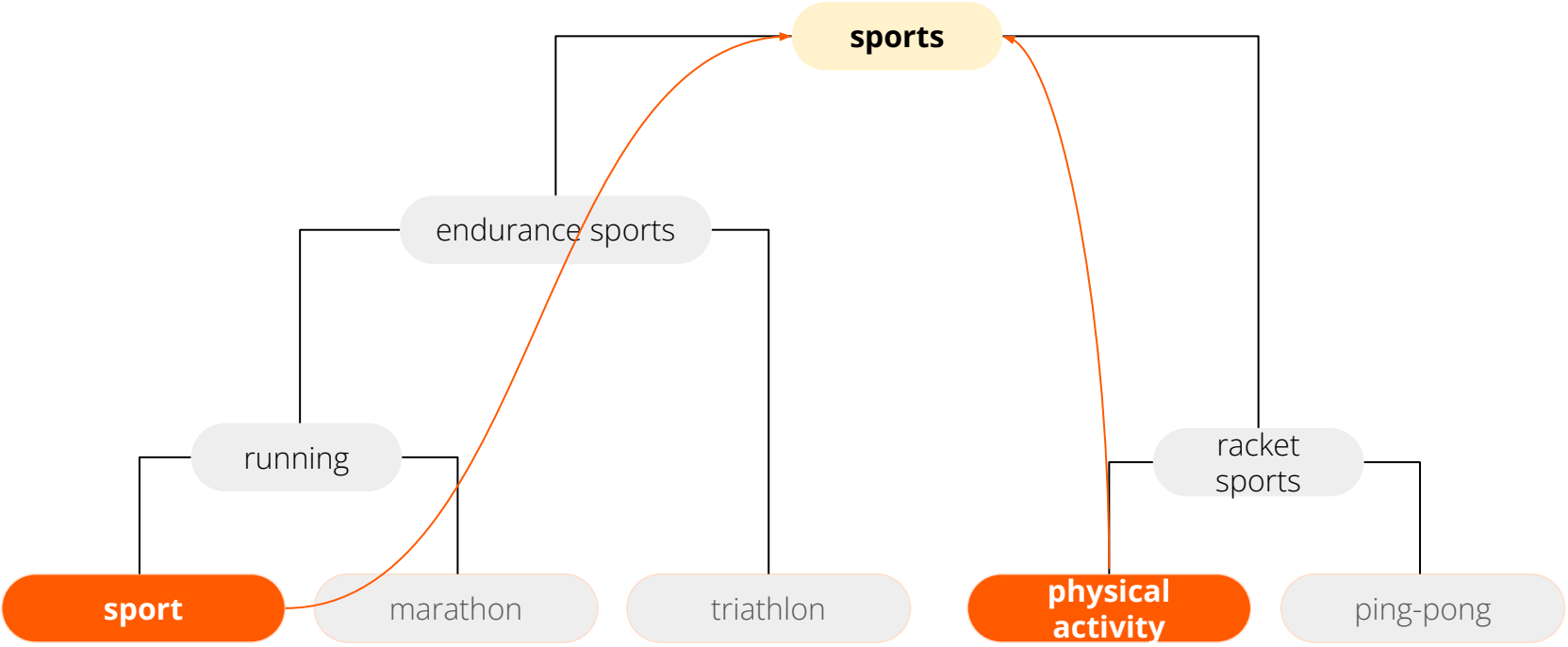
Tree post-processing: What's missing?



Custom post-processing: Naming nodes



Custom post-processing: making discovered structure coherent



What about new topics?



NEW TOPIC

outdoor play

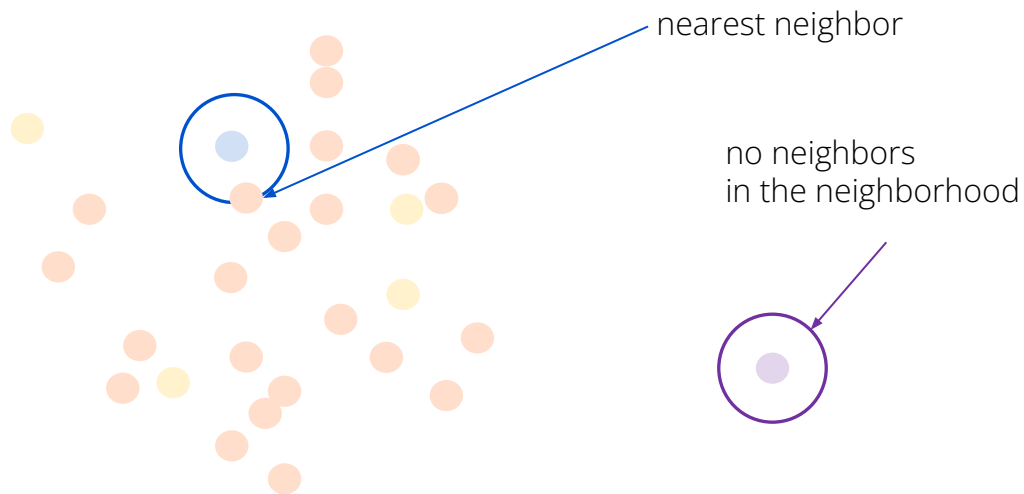
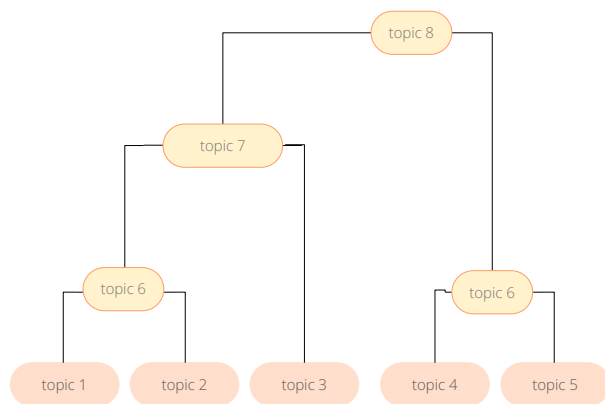
active family

hobby horsing

TOPIC FROM THE TREE

outdoor play

family activities



What about new topics?



NEW TOPIC

outdoor play

active family

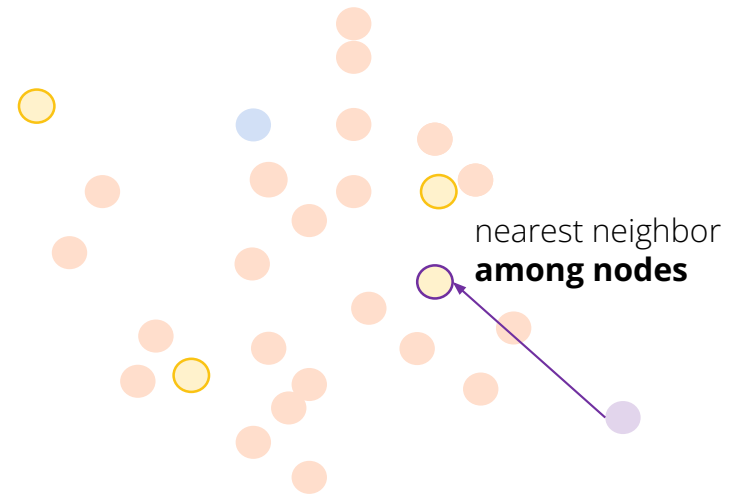
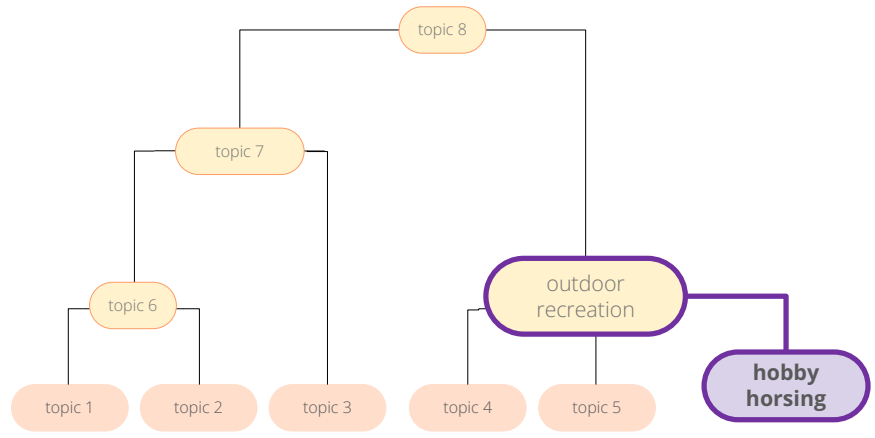
hobby horsing

TOPIC IN THE TREE

outdoor play

family activities

hobby horsing



The approach: hybrid pipeline

STEP 01



HOBBY GENERATION

Assigning topics to millions of products using textual data.

TECHNOLOGY STACK

Generative AI



STEP 02

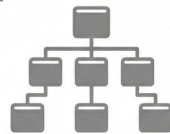


HOBBIES STRUCTURATION

Organizing raw tags into coherent, hierarchical tree-based structure.

TECHNOLOGY STACK

Clustering + Embedding Model
+ Generative AI



STEP 03



USERS PROFILING

Aggregating history to find the strongest interests using behavioral dynamics.

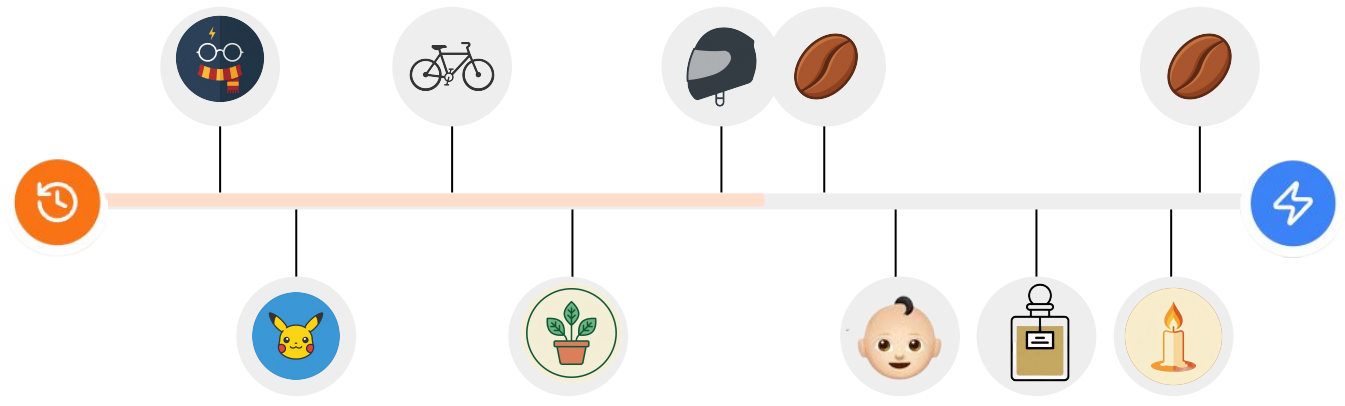
TECHNOLOGY STACK

Recency Frequency Monetary
Scoring



HOW TO GET USER PROFILES?

User profiling input



+ hobbies per each product

User Profiling Approach

Calculating the strength of interests of each user using behavioral dynamics.

R

Recency

How long ago was the last interaction?

F

Frequency

How often does the user engage with this topic?

M

Monetary

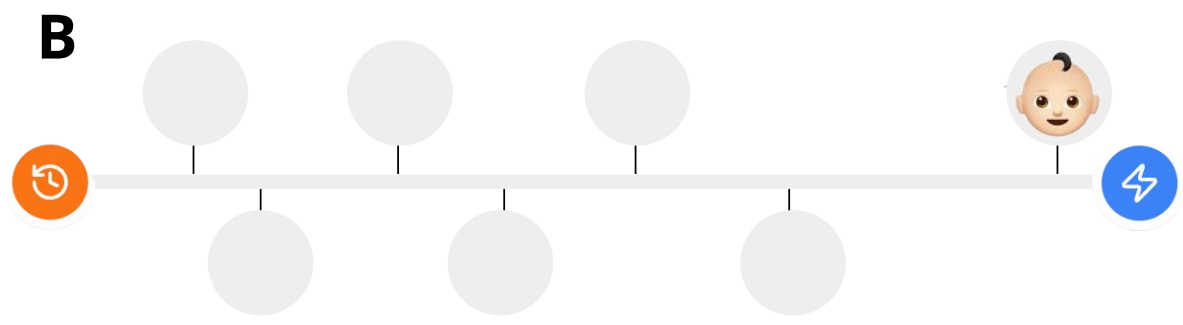
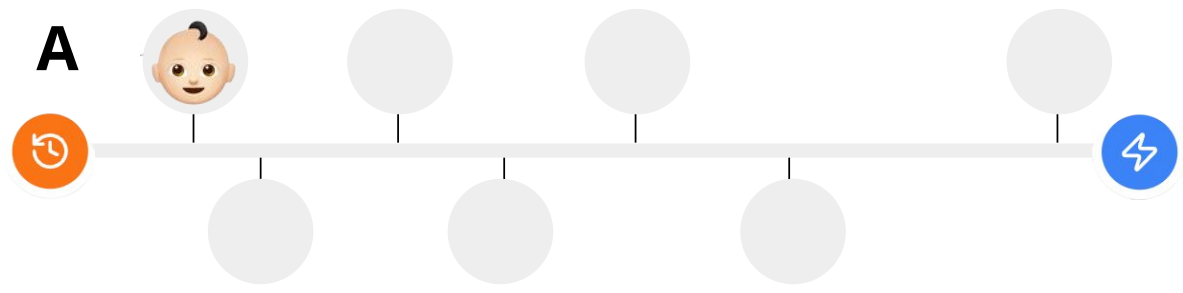
What is the value of items in this hobby?

USER FINAL SCORE

CYCLING**0.93****HARRY POTTER****0.70****COFFEE LOVER****0.35**

Why recency?

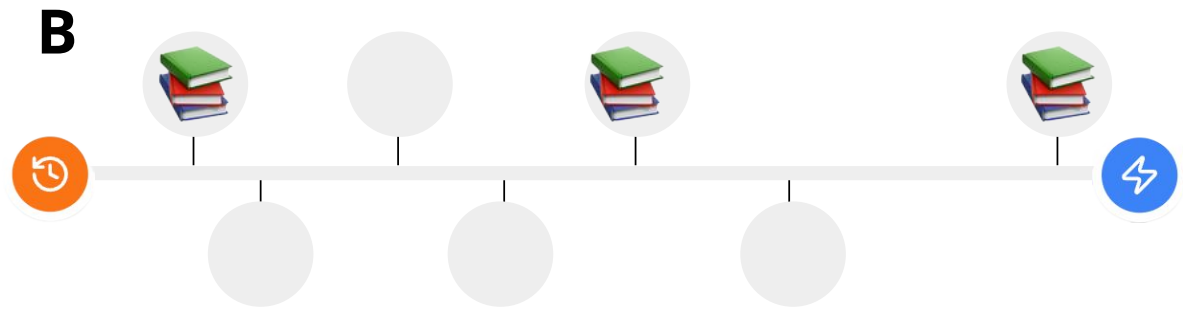
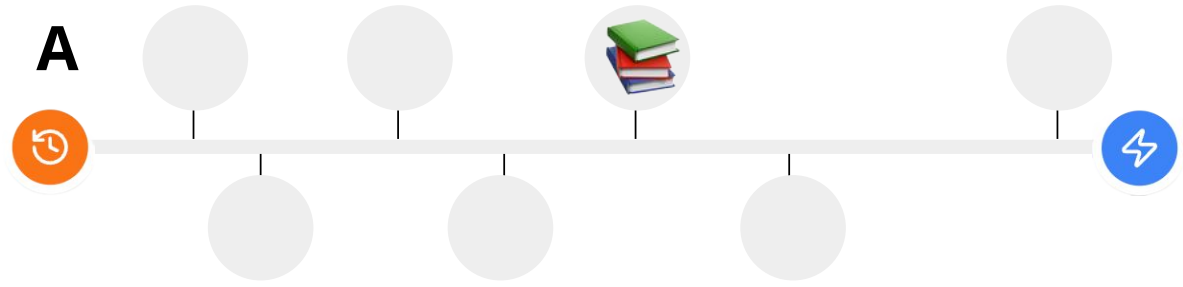
How long ago was the last interaction?



← Higher recency score!

Why frequency?

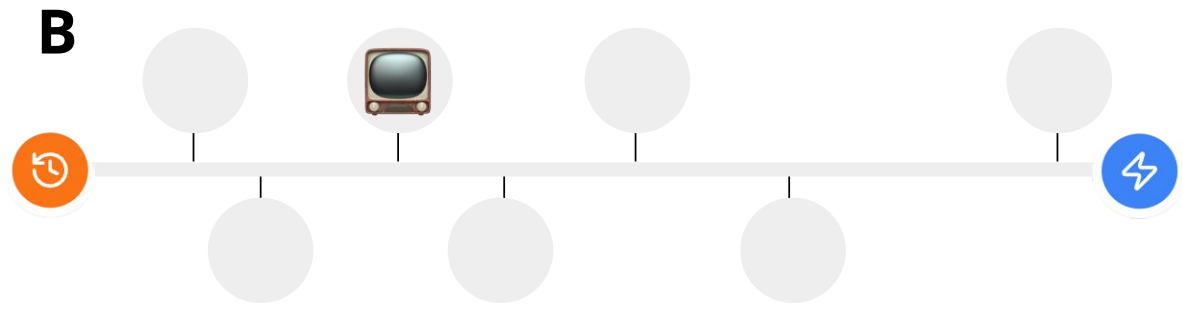
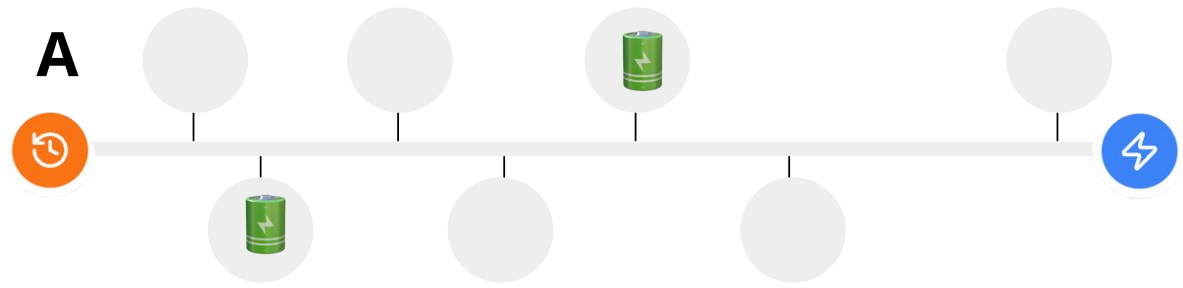
How often does the user engage with this topic?



← Higher frequency score!

Why monetary?

What is the value of items in this hobby?



← Higher monetary score!

Importance of scores calibration

BEST USER FOR A TOPIC

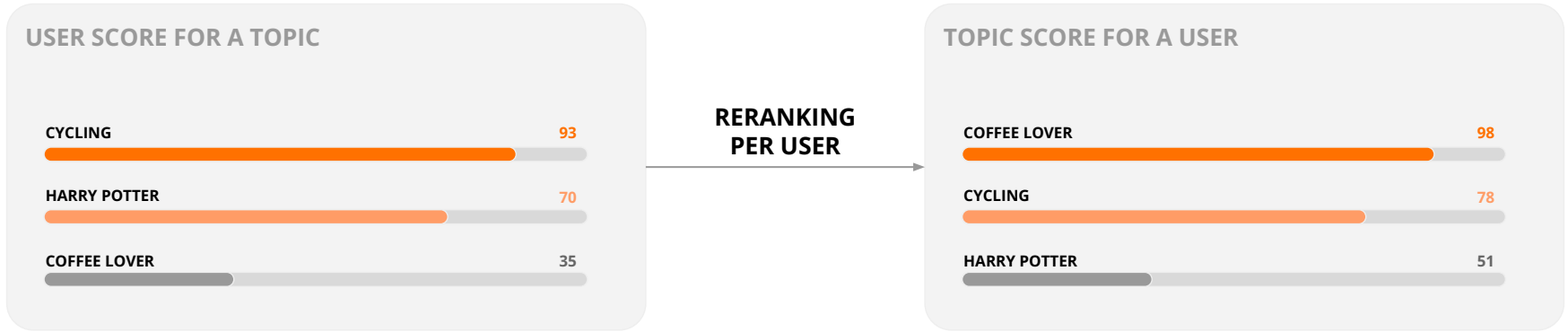
RFM score is based on **relativity to all the other users** engaged in a topic.



BEST TOPIC FOR A USER

Must be relative to the other topics of a user.

How to rerank the hobbies?



- How **much spent** on hobby vs other hobbies ?
- How **many unique products** related to a hobby?
- How **long** am I interested in a hobby?
- Am I relatively **engaged** in a hobby? (RFM)

OUR INTERESTS PROFILES

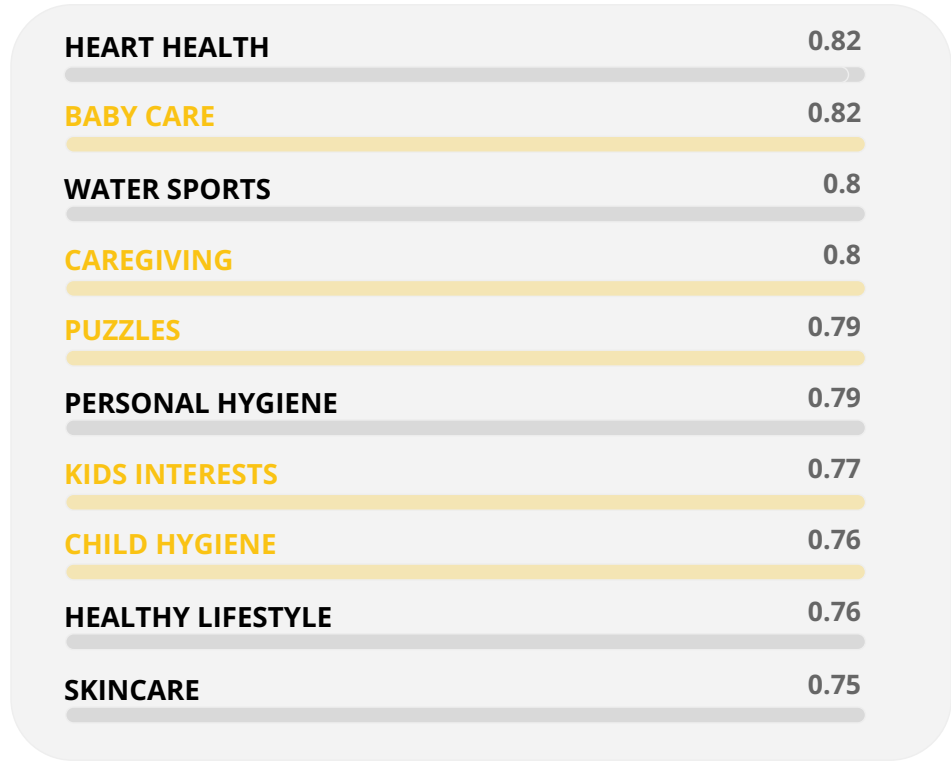
Our passions vs hobby profiles



motherhood

skiing

contemporary dance



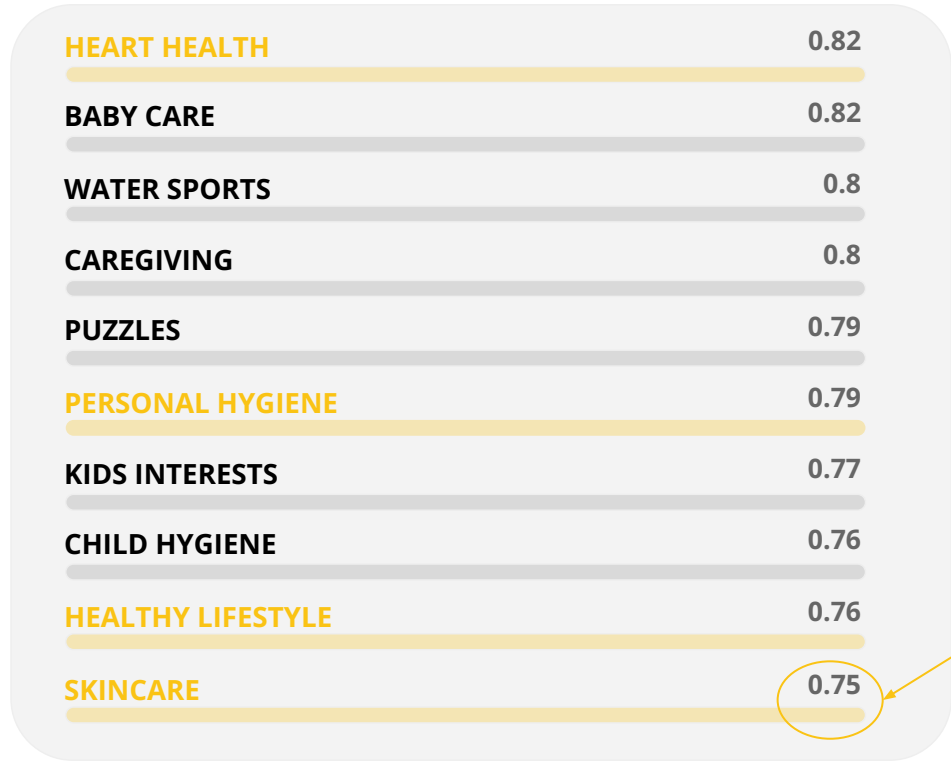
Our passions vs hobby profiles



motherhood

skiing

contemporary dance



Our passions vs hobby profiles



belly dance

healthy lifestyle

board games

WELLNESS 0.97

BEAUTY 0.84

HAIR CARE 0.81

HEALTH 0.77

TEXTILE ARTS 0.76

GADGET 0.75

ORAL HYGIENE 0.73

COSMETICS 0.73

HOME ORGANIZATION 0.69

FANTASY 0.66

Our passions vs hobby profiles



belly dance

healthy lifestyle

board games

WELLNESS 0.97

BEAUTY 0.84

HAIR CARE 0.81

HEALTH 0.77

TEXTILE ARTS 0.76

GADGET 0.75

ORAL HYGIENE 0.73

COSMETICS 0.73

HOME ORGANIZATION 0.69

FANTASY 0.66

Our passions vs hobby profiles



belly dance

healthy lifestyle

board games

WELLNESS 0.97

BEAUTY 0.84

HAIR CARE 0.81

HEALTH 0.77

TEXTILE ARTS 0.76

GADGET 0.75

ORAL HYGIENE 0.73

COSMETICS 0.73

HOME ORGANIZATION 0.69

FANTASY 0.66

HOW DO WE UNLOCK BUSINESS VALUE?

Targeted campaigns: Stranger Things Edition



Koszulka Hawkins

Fashion



Figurka Pop

Collectibles



Koszulka Hawkins

Home



Kinder Joy

Food

» PRODUCT SELECTION

How do we identify the right inventory?

Prioritize products with Stranger Things hobby tags.

» AUDIENCE LOGIC

What are the ideal recipients?

Target users having a high affinity score for stranger things fandom or 80s nostalgia/sci-fi genre.

Thematic collections



Harry Potter Fandom



Pokemon Collection



Labubu Series

» SELECTION LOGIC

Which collection to show to user X?

We can sort collections using topic-user ranking score.

Key takeaways

- hobby judge enables iterative prompt optimization
- it's worth to assign > 1 hobby to a product
- hobby generation does not have to be expensive
- inconsistent granularity needs attention
- proper aggregation level is crucial to achieve diverse hobby profile
- new topics need attention
- scoring users per topic \neq scoring topic per user

Thank you for your attention!

allegro