



Zucchini or Cucumbers?

**Benchmarking Embeddings
for Similar Image Retrieval**
thanks to your **weekly Grocery shopping**

Let's start with a little game

Looking similar



Zucchini_025.jpg_
test



Zucchini_001.jpg_
train



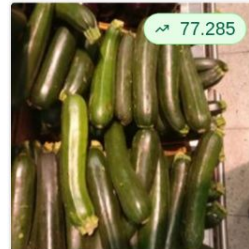
Zucchini_022.jpg_
test



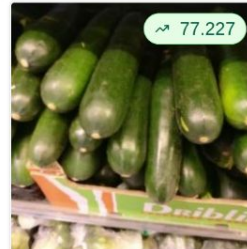
Zucchini_018.jpg_
train



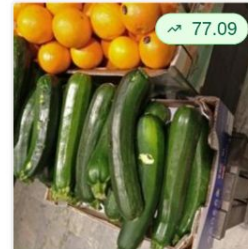
Zucchini_007.jpg_
test



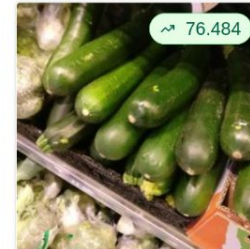
Zucchini_025.jpg_
test



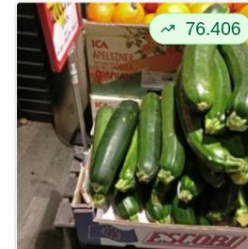
Zucchini_001.jpg_
train




Zucchini_022.jpg_
test



Zucchini_018.jpg_
train



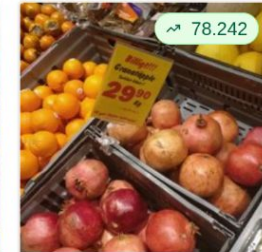
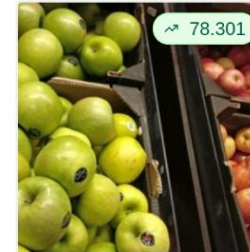
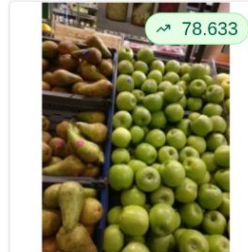
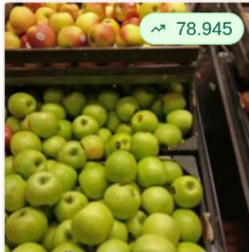
Zucchini_007.jpg_
test

How many of these are *Cucumbers*? 

Let's start with a little game



Mango_001.jpg_test



Pomegranate_019.
jpg_test

How many of these are *Apples*? 🍏

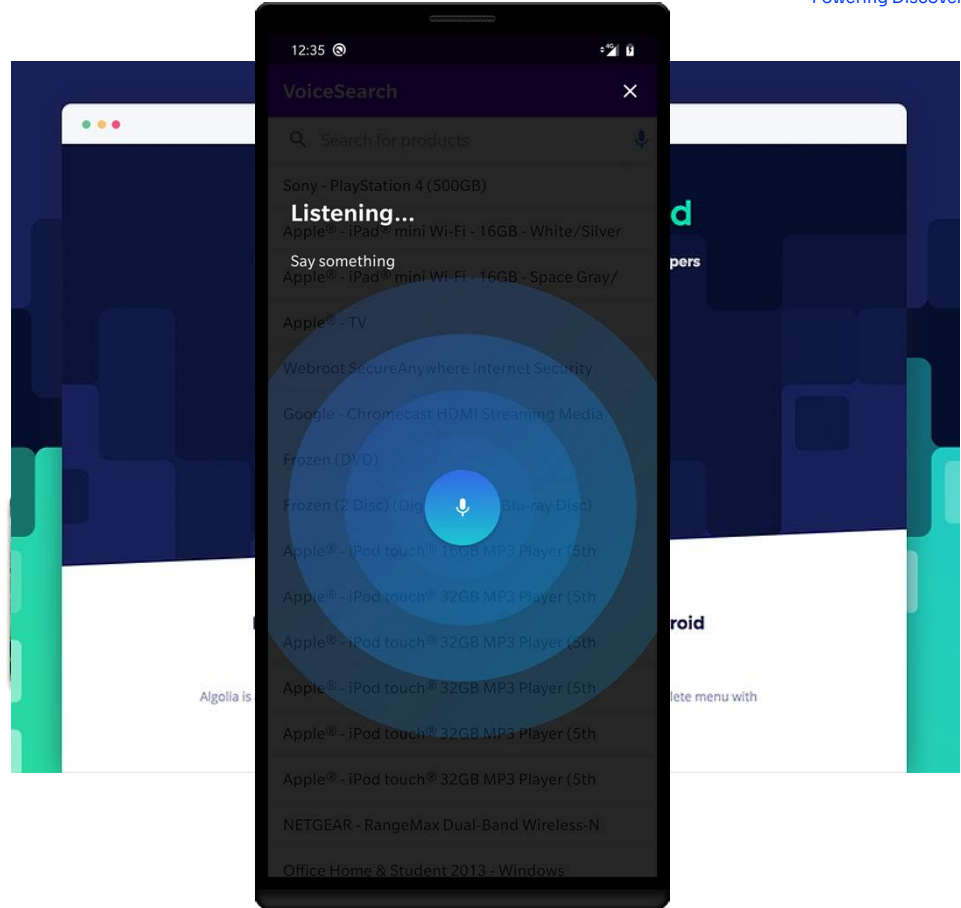


Hello, my name is Paul-Louis Nech

- **Coder** for more than 15 years
- **Software engineer** since 2016
- **Machine Learning** for the last 5 years

I build APIs & tools at algolia

- Android **UI libraries**
- Voice Search **tooling**
- Natural Language **APIs**
- Image Recommendation **engine**







My last project: LookingSimilar

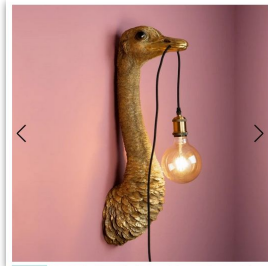
An Image Recommendation API
based on **image embeddings**
and **ANN retrieval**

CHECK IT OUT

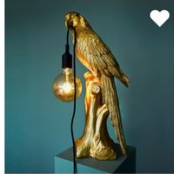
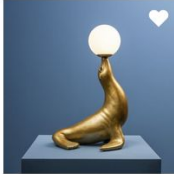



Shop Similar Style Products

			
Ring Pulls 3 Inch Diameter Unlacquered Antique Brass Cabinet Ring Pull \$32.40 ★★★★★ (8)	Ring Pulls 7 Inch Diameter Unlacquered Antique Brass Cabinet Ring Pull \$123.30 ★★★★★ (8)	Ring Pulls 5-5/8 Inch Diameter Unlacquered Antique Brass Cabinet Ring Pull \$85.50 ★★★★★ (8)	Ring Pulls 2-1/4 Inch Diameter Unlacquered Antique Brass Cabinet Ring Pull \$22.50 ★★★★★ (8)






Foarte asemănătoare

		
Lampă de masă, Papagal auri, Timmy, 27,5 x 18 x 61 cm	Lampă de masă, Focă aurie, Robbie, 32,5 x 28,5 x 51,5 cm	Lampă de masă, Maimuță aurie, Abu, 25,5 x 23,5 x 39,5 cm

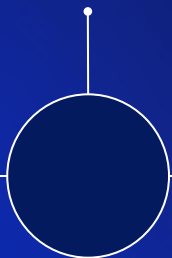


VISUALLY SIMILAR MOSAICS

		
Decorative Flower Arrangement Mosaic \$1,499.00 FL060	Multicolored Flower Arrangement Mosaic \$1,020.00 FL111	Multicolored Roses in Pot Mosaic \$799.00 FL057

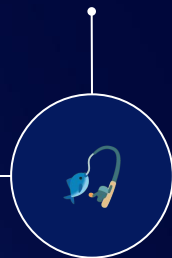
LookingSimilar: an image vector retrieval API

Vectorization
With an image model



Hashing
Learned LSH

Retrieval
With HNSW



Filtering
Using object-level
metadata



**It all comes from
a customer question**

“Do you handle groceries well?”

This is a good question

Do we know how to recognize food?

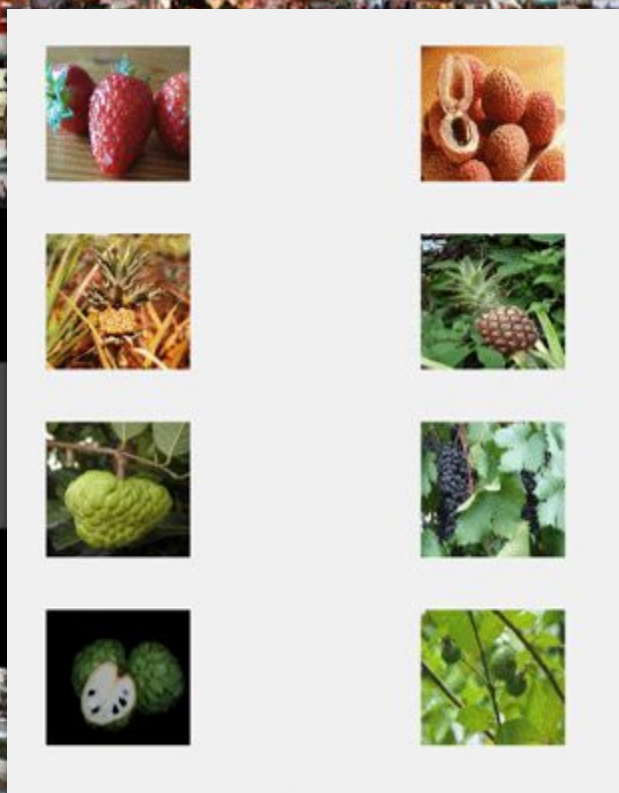


Searching for a Benchmark 🍏 🥒

Surely there must be industry standards?

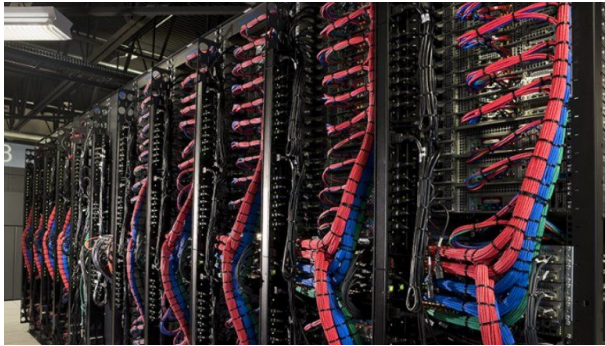


947	mushroom
948	Granny Smith
949	strawberry
950	orange
951	lemon
952	fig
953	pineapple, ananas
954	banana
955	jackfruit, jak, jack
956	custard apple
957	pomegranate



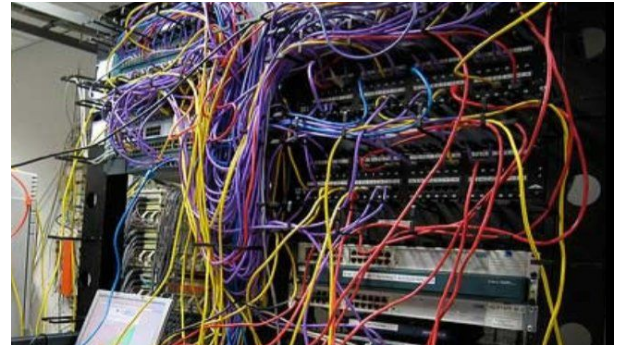
Expectation

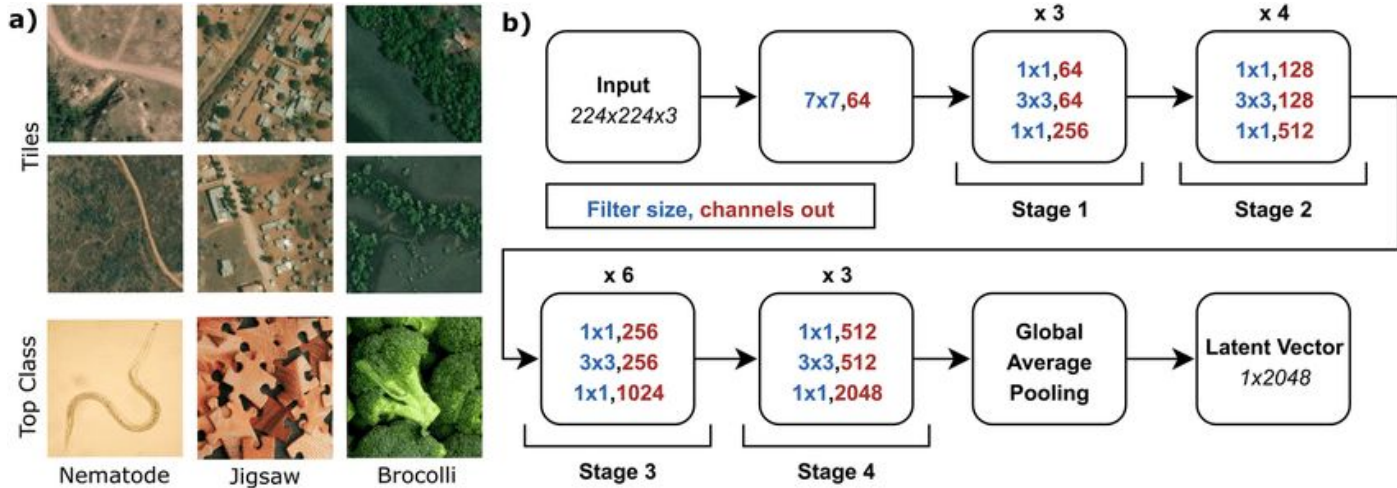
You can use **industry benchmarks** on your **SOTA** model



Reality

Industry benchmarks are **already used** in training, test, or validation -> probably **biased** in many **surprising ways**





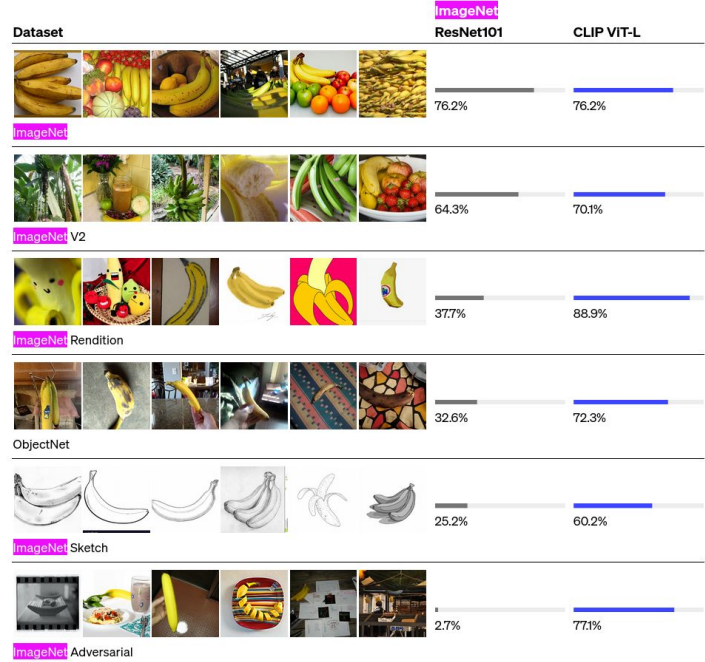
ImageNet data → ResNet training

You can't use the **training data** to test the model ⚠

CLIP also is biased by ImageNet

If not **at training**, definitely **tested** and **iterated on**

-> You could say this dataset was **gamed against** by the OpenAI team



Although both models have the same accuracy on the ImageNet test set, CLIP's performance is much more representative of how it will fare on datasets that measure accuracy in different, non-ImageNet settings. For instance, ObjectNet checks a model's ability to recognize objects in many different poses and with many different backgrounds inside homes while ImageNet Rendition and ImageNet Sketch check a model's ability to recognize more abstract depictions of objects.

There is no good benchmark

But we still need to answer that question

Grocery Store dataset

Grocery Store Dataset

This repository contains the dataset of natural images of grocery items. All natural images were taken with a smartphone camera in different grocery stores. We ended up with 5125 natural images from 81 different classes of fruits, vegetables, and carton items (e.g. juice, milk, yoghurt). The 81 classes are divided into 42 coarse-grained classes, where e.g. the fine-grained classes 'Royal Gala' and 'Granny Smith' belong to the same coarse-grained class 'Apple'. For each fine-grained class, we have downloaded an iconic image and a product description of the item, where some samples of these can be seen on this page below. The dataset was presented in the paper "[A Hierarchical Grocery Store Image Dataset with Visual and Semantic Labels](#)", which appeared at WACV 2019.

How to use the dataset

The files `train.txt`, `val.txt` and `test.txt` in the folder `dataset` includes the paths to the images in the training, validation and test set respectively. Each row in these two files consists of the path to an image and its fine-grained label followed by its coarse-grained label, where both labels are represented as integers.

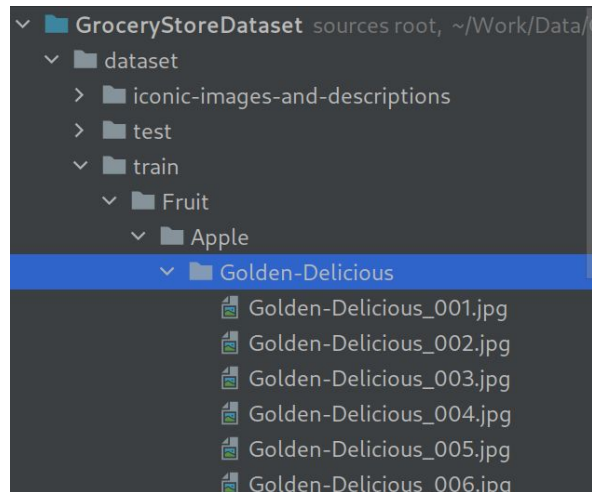
The 81 fine-grained classes and their coarse-grained classes can be found in `classes.csv` in the folder `dataset`. The classes corresponding label (an integer) is also included in addition to the paths to their iconic image and the product description.

Feel free to download the dataset and apply it to your model.

Samples of natural images



Samples of iconic images



A multi-layer hierarchy of **natural, user-generated** images

By Marcus Klasson and Cheng Zhang and Hedvig Kjellström

Labeled data?

[GroceryStoreDataset](#) / [dataset](#) / [iconic-images-and-descriptions](#) / [Fruit](#) / [Apple](#) / **Golden-Delicious** / 


 **marcusklasson** Added classes.csv and changed structure of iconic images and descript... 

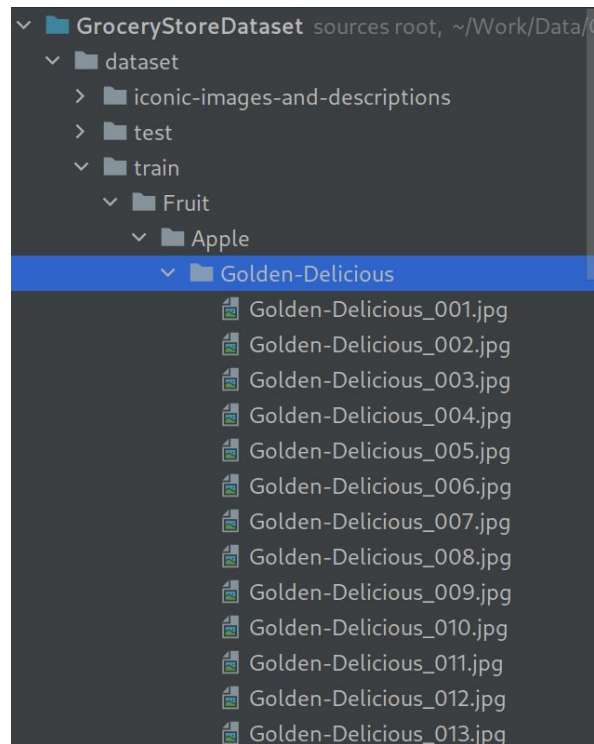
Name

 ..

 Golden-Delicious_Description.txt

 Golden-Delicious_Iconic.jpg

 Golden-Delicious_Information.txt



YAGNI? No, YMANI!

[GroceryStoreDataset](#) / [dataset](#) / [iconic-images-and-descriptions](#) / [Fruit](#) / [Apple](#) / [Golden-Delicious](#) / **Golden-Delicious_Description.txt** 

 **marcusklasson** Added classes.csv and changed structure of iconic images and descript... 

Code Blame 1 lines (1 loc) · 147 Bytes  Code 55% faster with GitHub Copilot

```
1 Golden Delicious has a white juicy pulp and a greenish yellow shell. The taste is mellow and sweet, making Golden Delicious suitable for desserts.
```

You Might Also Need It™

[GroceryStoreDataset](#) / [dataset](#) / [iconic-images-and-descriptions](#) / [Fruit](#) / [Apple](#) / [Golden-Delicious](#) / **Golden-Delicious_Information.txt** 

 **marcusklasson** Added classes.csv and changed structure of iconic images and descript... 

bcd6345


Code Blame 4 lines (4 loc) · 315 Bytes  Code 55% faster with GitHub Copilot Raw



```
1 Title: Apple Golden Delicious Class 1
2 Country and volume: Italy, ca 180g
3 Description: Golden Delicious has a white juicy pulp and a greenish yellow shell. The taste is mellow and sweet, making Golden Delicious suitable for desserts.
4 URL: https://www.hemkop.se/produkt/apple-Golden-Delicious-Klass-1-100142141_KG
```

Structuring data – a kind of Magic

1

objectID "Zucchini_030.jpg_train"

[GroceryStoreDataset](#) / [dataset](#) / [iconic-images-and-descriptions](#) / [Fruit](#) / [Apple](#) / **Golden-Delicious** / 

 **marcusklasson** Added classes.csv and changed structure of iconic images and descript... 

Name
..
Golden-Delicious_Description.txt
Golden-Delicious_Iconic.jpg
Golden-Delicious_Information.txt

✓ **GroceryStoreDataset** sources root, ~/Work/Data/

- dataset
 - iconic-images-and-descriptions
 - test
 - train
 - Fruit
 - Apple
 - Golden-Delicious**
 - Golden-Delicious_001.jpg
 - Golden-Delicious_002.jpg
 - Golden-Delicious_003.jpg
 - Golden-Delicious_004.jpg
 - Golden-Delicious_005.jpg
 - Golden-Delicious_006.jpg
 - Golden-Delicious_007.jpg
 - Golden-Delicious_008.jpg
 - Golden-Delicious_009.jpg
 - Golden-Delicious_010.jpg
 - Golden-Delicious_011.jpg
 - Golden-Delicious_012.jpg
 - Golden-Delicious_013.jpg

The million-dollar insight

Any **labeled** dataset is a **benchmark**
asking you to be **built**



“A **Benchmark** is only
a **dataset** with **testable labels**”



Paul-Louis Nech
ML Engineer

Step 1/2: structuring - `index.py`

```
def main():  
    records = []  
    datasets = ["test", "train", "val"]  
    for dataset in datasets:  
        records.extend(browse_dataset(dataset))  
    index.save_objects(records)  
    print(f"Done saving {len(records):04} {datasets} groceries.")
```

Step 1/2: structuring - index.py

```
def browse_dataset(dataset_name: Optional[str] = "train") -> list[dict[str, Any]]:
    records = []
    metadata = browse_metadata()
    datasets = "/home/pln/Work/Data/GroceryStoreDataset/dataset/"
    for dir, _, files in os.walk(os.path.join(datasets, dataset_name)):
        for i, file in enumerate(files, start=1):
            # Process filename
            # e.g. "Avocado_001.jpg" -> ("Avocado", "001")
            product, id = file.split("_")
            id = id.removesuffix(".jpg")

            # Process directory labels
            # e.g. train/Fruit/Apple/GrannySmith/ -> ("Fruit", "Apple", "GrannySmith")
            imgpath = dir.removeprefix("/home/pln/Work/Data/GroceryStoreDataset/")
            iconic_dir = imgpath.replace(f'dataset/{dataset_name}', 'dataset/iconic-images-and-
descriptions')
            iconic_path = f"{REPO_ROOT}/{iconic_dir}/{product}_Iconic.jpg?raw=true"

            # Edge case: some items have a category but no subCategory
            categories = dir.split("/")[-3:]
            if dataset_name in categories:
                categories = categories[1:]


            record = {...}

            records.append(record)




    records.sort(key=lambda x: x['objectID'])
    return records
```

Step 1/2: structuring - index.py

1

objectID	"Zucchini_030.jpg_train"	
image_id	"030"	
category	"Vegetables"	
sub_category	"Zucchini"	
name	"Zucchini"	
description	"Zucchini can be eaten raw, but is usually used in soups, pots, stews, etc. "	
sku	"Zucchini"	
information	{ title: "Zucchini Class 1", country: "Spain", description: "Zucchini can be eaten raw, but is usually used in soups, pots, stew	
dataset	"train"	
image	"https://github.com/marcusklason/GroceryStoreDataset/blob/master/dataset/train/Vegetables/Zucchini/Zucchini_030.jpg?raw=true"	
iconic	"https://github.com/marcusklason/GroceryStoreDataset/blob/master/dataset/iconic-images-and-descriptions/Vegetables/Zucchini/Zucchini_Iconic.jpg?raw=true"	

Show fewer attributes ^

You Might Also Need It™

We can now evaluate
But *what* do we want to measure?

What do we want to know about our algorithm?

We want to assess quality of recommendations on groceries

- That is, we care about **precision** of categories (🍎 -> only 🍎🍏🍎, no 🍌)
- Not that interested in their **recall** (🍎 -> **only apples** >> 🍎 -> **all the apples**)
- Customers also care about **always showing something** (📊 no results rate)

→ Our metrics:

GotResults, SameCategory,
SameSubCategory, SameProduct

Step 2/2: evaluating - bench.py

```
for item in dataset:
    query = { ... }
    response = client.get_recommendations([query])["results"][0]
    recos = response["hits"]
    if recos: # Let's check the quality of recommended items!
        # ...
    else:
        # Let's report failed items for further investigation.
        no_match.append(item)

metrics.compute(no_match, same_cat, same_sub, same_sku)

print(
    f"Benchmark@{at_k} results - {metrics.nb_items} items\n"
    f"- Got some results {metrics.ratio_predicted:.2%} of the time.\n"
    f"- Predicted same category {metrics.ratio_same_category:.2%} of the time.\n"
    f"- Predicted same sub-cat {metrics.ratio_same_subcategory:.2%} of the time.\n"
    f"- Predicted same product {metrics.ratio_same_sku:.2%} of the time.\n"
)
```

Step 2/2: evaluating - bench.py



```
for reco_item in recos:
    reco_item_cat = reco_item["category"]
    reco_item_sub = reco_item["sub_category"]
    reco_item_sku = reco_item["sku"]
    is_same_cat = item_cat == reco_item_cat
    is_same_sub = item_sub == reco_item_sub
    is_same_sku = item_sku == reco_item_sku
    same_cat.append(is_same_cat)
    same_sub.append(is_same_sub)
    same_sku.append(is_same_sku)

    if not is_same_cat:
        if is_same_sub:
            print("WTF??")
            wrong_cat.append((item, reco_item))
    if not is_same_sub:
        if is_same_sku:
            print("WTF??")
            wrong_sub.append((item, reco_item))
    if not is_same_sku:
        wrong_sku.append((item, reco_item))
```

**GARBAGE IN,
GARBAGE OUT**

A small tip: make it reliable

Once your benchmark runs, make it **reliable**

same model, same input, same metric
should result in **same scores**

-> **Freeze** your GPT4 **model version**

-> use `random.seed(42)`

-> etc - eliminate as much **randomness** as you can



Step 2/2: evaluating - bench.py demo!

```
191  
192 ▶ if __name__ == "__main__":  
193     at_k = [1, 3, 5, 10, 30]  
194  
195     bench_predictions(filename: "predictions/CLIP_7510047928944698635.json", at_k=at_k)  
196     bench_predictions(filename: "predictions/ResNet_7510077461642073363.json", at_k=at_k)  
197
```

Step 2/2: evaluating - bench.py demo!

Begin of benchmark for predictions in file predictions/ResNet_7510077461642073363.json.

Fetching up to 30 recos for 5420 items...

Starting benchmark@1 on index `groceries` over 5420 items.

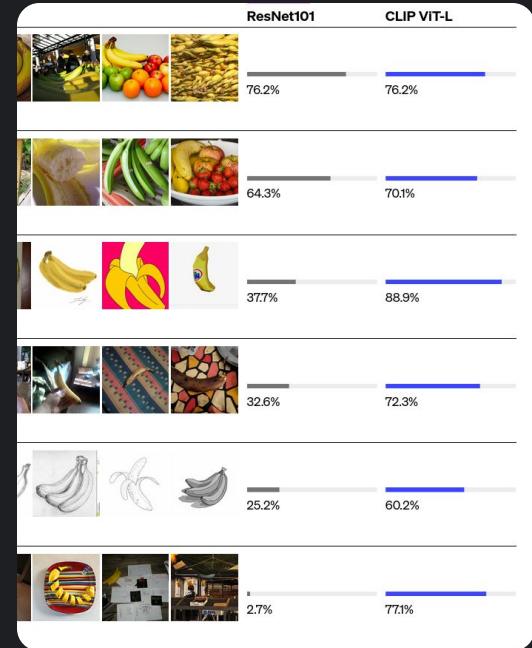
Benchmark@1 results - 5420 items

- Got some results 99.28% of the time.
- Predicted same category 99.63% of the time.
- Predicted same sub-cat 96.04% of the time.
- Predicted same product 92.83% of the time.

Starting benchmark@3 on index `groceries` over 5420 items.

Benchmark@3 results - 5420 items

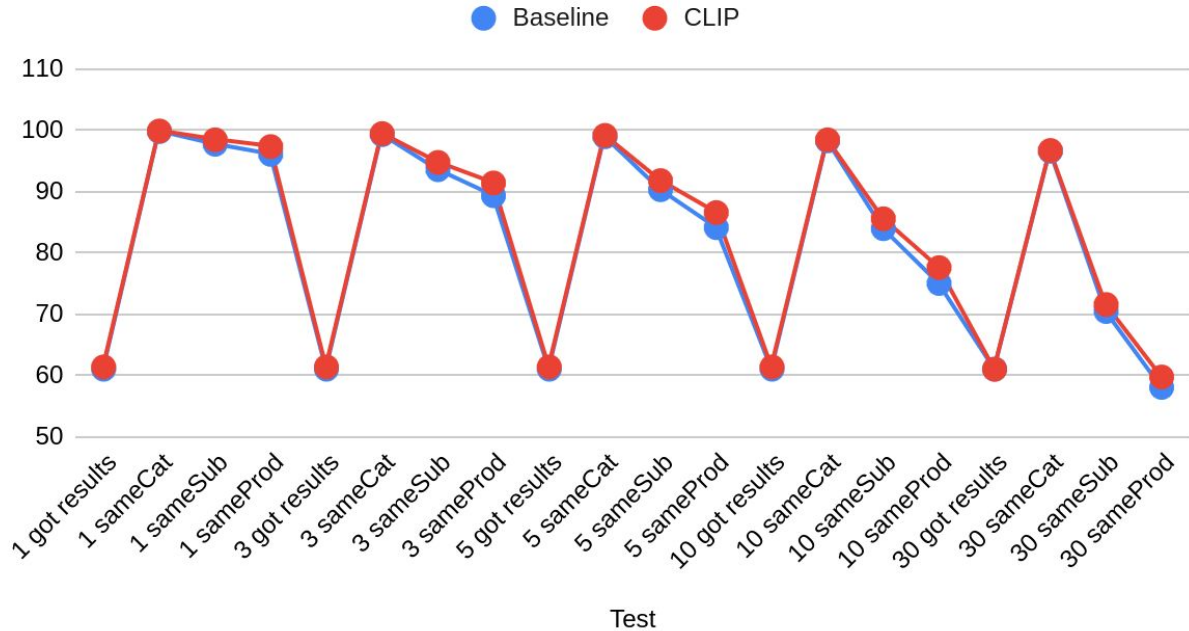
- Got some results 99.28% of the time.
- Predicted same category 98.98% of the time.
- Predicted same sub-cat 90.81% of the time.
- Predicted same product 84.28% of the time.



Step 2/2: evaluating - bench.py results

Test	Baseline	CLIP
------	----------	------

Baseline vs CLIP



Non-significant difference

And **lower stability** (when multi threading)

30 sameSub	70.7	71.0
30 sameProd	58.1	59.8












Some evaluation results

Evaluation results: packaged goods?



Alpro-Vanilla-Soyghurt_023.jpg_train × Show random Show Rules?

Item viewed Looking similar

 <p>Alpro-Vanilla-Soyghurt_023.jpg_... 023 Packages Soyghurt</p>	 <p>Alpro-Vanilla-Soyghurt_007.jpg_... 007 Packages Soyghurt</p>	 <p>Alpro-Vanilla-Soyghurt_002.jpg_... 002 Packages Soyghurt</p>	 <p>Alpro-Vanilla-Soyghurt_006.jpg_... 006 Packages Soyghurt</p>	 <p>Alpro-Blueberry-Soyghurt_011.jpg_... 011 Packages Soyghurt</p>	 <p>Alpro-Vanilla-Soyghurt_021.jpg_... 021 Packages Soyghurt</p>
	 <p>80.137</p>	 <p>80</p>	 <p>79.844</p>	 <p>79.629</p>	 <p>79.277</p>









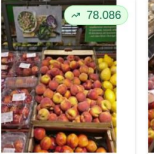
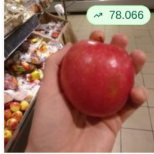
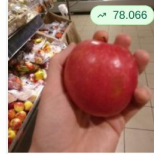




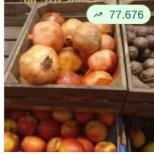
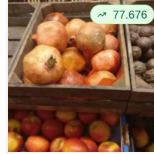



Reality way paired different late than to vs. Overize

Evaluation aside: TVs vs Ovens



!?

Evaluation results: comparing 🍎 to 🍊?

Item viewed	Looking similar						
 <p>Pink-Lady_059.jpg_train 059 Fruit Apple</p>	 <p>↔ 100</p> <p>Pink-Lady_059.jpg 059 Fruit Apple</p>	 <p>↔ 82.051</p> <p>Pink-Lady_048.jpg_train 048 Fruit Apple</p>	 <p>↔ 82.051</p> <p>Pink-Lady_048.jpg 048 Fruit Apple</p>	 <p>↔ 82.051</p> <p>Pink-Lady_001.jpg_train 001 Fruit Apple</p>	 <p>↔ 79.824</p> <p>Pomegranate_019.jpg_test 019 Fruit Pomegranate</p>	 <p>↔ 78.809</p> <p>Peach_018.jpg_test 018 Fruit Peach</p>	
	 <p>↔ 78.789</p> <p>Pink-Lady_037.jpg_test 037 Fruit Peach</p>	 <p>↔ 78.086</p> <p>Peach_035.jpg_test 035 Fruit Peach</p>	 <p>↔ 78.066</p> <p>Pink-Lady_016.jpg_train 016 Fruit Apple</p>	 <p>↔ 78.066</p> <p>Pink-Lady_016.jpg 016 Fruit Apple</p>	 <p>↔ 77.988</p> <p>Granny-Smith_024.jpg_train 024 Fruit Apple</p>	 <p>↔ 77.988</p> <p>Granny-Smith_024.jpg 024 Fruit Apple</p>	
	 <p>↔ 77.715</p> <p>Peach_022.jpg_train 022 Fruit Peach</p>	 <p>↔ 77.715</p> <p>Peach_022.jpg 022 Fruit Peach</p>	 <p>↔ 77.676</p> <p>Pomegranate_014.jpg_train 014 Fruit Pomegranate</p>	 <p>↔ 77.676</p> <p>Pomegranate_014.jpg 014 Fruit Pomegranate</p>	 <p>↔ 77.617</p> <p>Orange_046.jpg_train 046 Fruit Orange</p>	 <p>↔ 77.617</p> <p>Orange_046.jpg 046 Fruit Orange</p>	

“An **AI model’s strengths** and **weaknesses** are not what *we, humans*, think they are.”



Paul-Louis Nech
ML Engineer

Communicating the results

Optimize for action

“We measured
CategoryPrecision@1 = 99.62 🎉”



An engineer

Pretty proud of their metrics

“The **top recommendation** is in the
same category as the source item for
99.62% of all pictures”



A senior engineer

Who knows how to share metrics

“**WTF?** Can I speak to a
human on your team?”



A customer

Pretty confused right now

“Thanks! I’ll **trust your**
API with my precious
use-case 🥒 🍎 🍏”



A happy customer

Confident in your feature

“**WTF?** Am I talking to a bad LLM?”

A customer

Pretty confused right now



“The **top recommendation** is in the **same category** as the source item for **99.62%** of **all pictures**”



A senior engineer

Who knows how to share metrics

“Thanks! I can trust it with my
precious use-case 🌮”

A happy customer
Confident in your feature



Your interlocutors *don't want metrics*

They want a *quantitative assessment*
of the **usefulness of your product** 💡

➔ Focus on making your report **actionable**

What do you expect readers to **do** once they finish reading?

How can you **optimize for action**? 🦊

What do we want to say about our algorithm?

Communicate a metric tied to customer value

- **Precision@1** is a very **good-looking** metric, also pretty **useless**
- **Precision@3** is the customers' **main use-case**: carousel with few top hits
- **Precision@30** is the **hard truth**: up to maxRecs, how many are good?

- It depends **what** you want to communicate and **to whom**

To whom?



“SameProduct@3 is **88.64%**, let’s try **XYZ** and check if precision **goes up**”

To my dev colleagues



“In our **Carousel UI components**, we can expect **80% of top 3 results** to be the same SKU as user input”

To our Product Manager



“You can be **confident** that on **most groceries**, your users will see **relevant recommendations**”

To our Customer




Who are you talking to? Interlocutor Personas

Technical audiences

Want to know how to **run it**

Would love to **adapt** it to new features

Curious what can be **improved**

- Share **hands-on** knowledge 
- Make the code **modular** 
- Mention every **limitation** 

Non-Technical audiences

Want to know how to **leverage it**

Need **actionable talking points**
(especially customer-facing roles)

Must know **core strengths**

- Make the report **self-serve**
- Give them **ammunition**
- Cover use-cases where **feature shines**
and those **best avoided** for now

Who are you talking to? Internal/External

Internal audiences

Needs to see strengths and weaknesses in **crude light**, to know what needs **focus**

- Be **exhaustive** in what you report
- Focus on the **edge-cases**
- Share ideas for **improving these** weak spots of your feature

External audiences

Needs to know **where** it can trust you **today**
Focus on the **solid use-cases**

- Ace **core-use cases** with **great docs**
- When possible, explain **workarounds**
- Otherwise, mention **alternative features** better suited to this case

Iterating on our model

Now we have a Baseline

First, the good news: we do well on food



Item viewed

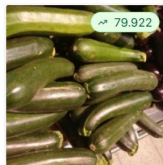


Zucchini_027.jpg_
test
027
Vegetables
Zucchini

Looking similar



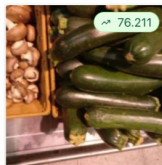
Zucchini_021.jpg_
test
021
Vegetables
Zucchini



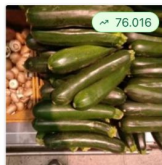
Zucchini_020.jpg_
test
020
Vegetables
Zucchini



Zucchini_025.jpg_
test
025
Vegetables
Zucchini



Zucchini_006.jpg_
test
006
Vegetables
Zucchini



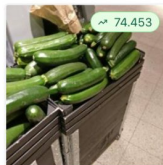
Zucchini_010.jpg_
train
010
Vegetables
Zucchini



Zucchini_009.jpg_
train
009
Vegetables
Zucchini



Zucchini_014.jpg_
train
014
Vegetables
Zucchini



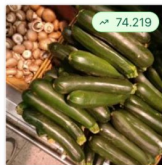
Zucchini_028.jpg_
test
028
Vegetables
Zucchini



Zucchini_001.jpg_
train
001
Vegetables
Zucchini



Zucchini_001.jpg_
test
001
Vegetables
Zucchini



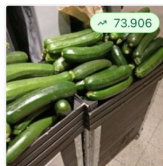
Zucchini_017.jpg_
train
017
Vegetables
Zucchini



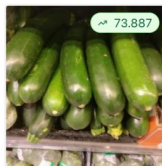
Zucchini_012.jpg_
train
012
Vegetables
Zucchini



Zucchini_016.jpg_
train
016
Vegetables
Zucchini



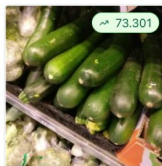
Zucchini_017.jpg_
test
017
Vegetables
Zucchini



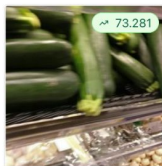
Zucchini_020.jpg_
train
020
Vegetables
Zucchini



Zucchini_007.jpg_
train
007
Vegetables
Zucchini



Zucchini_018.jpg_
train
018
Vegetables
Zucchini



Zucchini_001.jpg_val
001
Vegetables
Zucchini

First, the good news: we do well on food

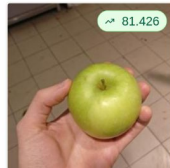


Item viewed

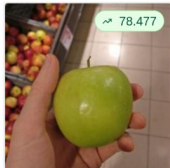


**Granny-Smith_059.
jpg_train**
059
Fruit
Apple

Looking similar



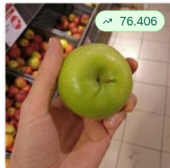
**Granny-Smith_017.
jpg_train**
017
Fruit
Apple



**Granny-Smith_015.
jpg_train**
015
Fruit
Apple



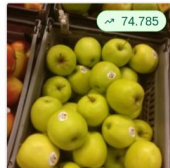
**Granny-Smith_033.
jpg_train**
033
Fruit
Apple



**Granny-Smith_050.
jpg_train**
050
Fruit
Apple



**Granny-Smith_040.
jpg_train**
040
Fruit
Apple



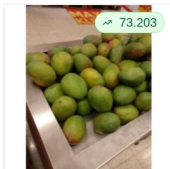
**Golden-Delicious_
028.jpg_train**
028
Fruit
Apple



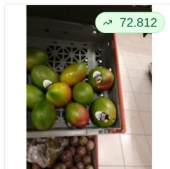
**Granny-Smith_021.
jpg_train**
021
Fruit
Apple



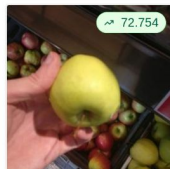
**Golden-Delicious_
003.jpg_train**
003
Fruit
Apple



Mango_017.jpg_test
017
Fruit
Mango



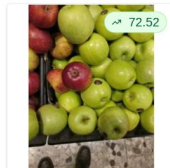
Mango_020.jpg_test
020
Fruit
Mango



**Golden-Delicious_
003.jpg_val**
003
Fruit
Apple



**Granny-Smith_024.
jpg_train**
024
Fruit
Apple



**Granny-Smith_001.
jpg_test**
001
Fruit
Apple



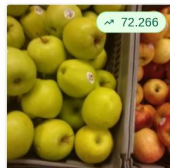
Papaya_006.jpg_test
006
Fruit
Papaya



**Golden-Delicious_
004.jpg_val**
004
Fruit
Apple



**Granny-Smith_003.
jpg_val**
003
Fruit
Apple



**Golden-Delicious_
014.jpg_train**
014
Fruit
Apple



**Granny-Smith_023.
jpg_test**
023
Fruit
Apple

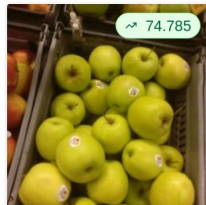
Some edge-cases where the model is confused

Item viewed



Granny-Smith

Granny-Smith_059.jpg_train



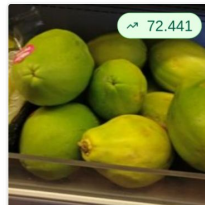
Golden-Delicious

Golden-Delicious_028.jpg_tr...



Mango

Mango_017.jpg_test



Papaya

Papaya_006.jpg_test



GotResults



SameCategory



SameSubCategory



SameProduct



How can we solve these?

Item viewed



Zucchini_030.jpg_train

030

Vegetables

Zucchini



Looking similar

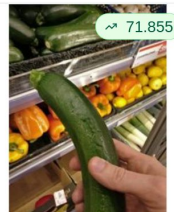


Green-Bell-Pepper_022.jpg_test

022

Vegetables

Pepper

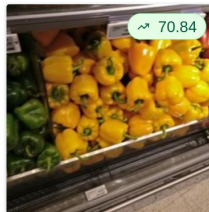


Zucchini_022.jpg_train

022

Vegetables

Zucchini



Yellow-Bell-Pepper_022.jpg_test

022

Vegetables

Pepper



GotResults



SameCategory



SameSubCategory



SameProduct



Addressing edge-cases: different embedding?



AI Computer Vision Research

DINOv2: A Self-supervised Vision Transformer Model

A family of foundation models producing **universal features** suitable for **image-level visual tasks** (image classification, instance retrieval, video understanding) as well as **pixel-level visual tasks** (depth estimation, semantic segmentation).

[Try the demos](#)

Addressing edge-cases: different embeddings?

Should we solve these edge-cases with **different embeddings**?

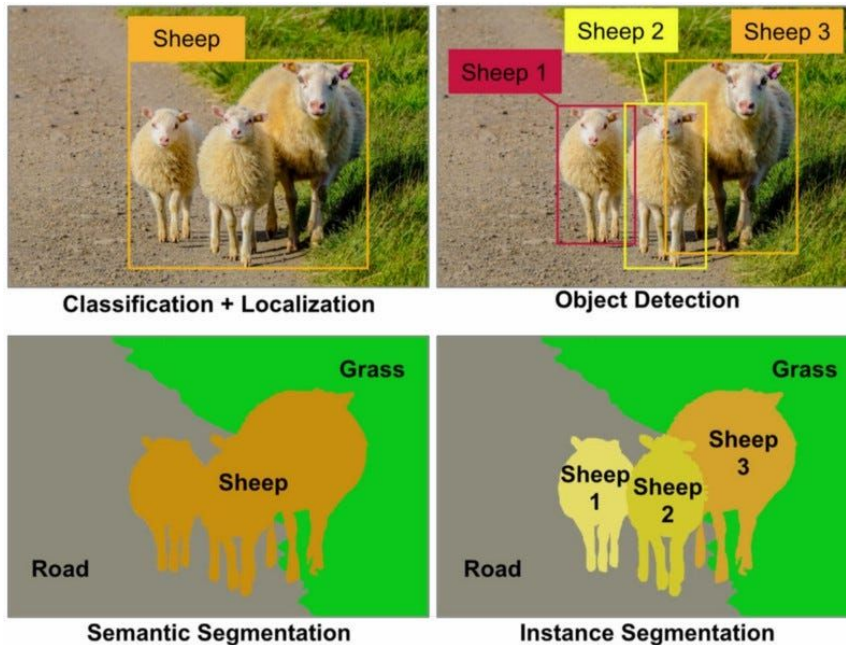
A **big move** with likely **side effects**

What else can we try first?

- Image Preprocessing
- Object Segmentation
- Adding textual signal in vectors



Addressing edge-cases: object segmentation?



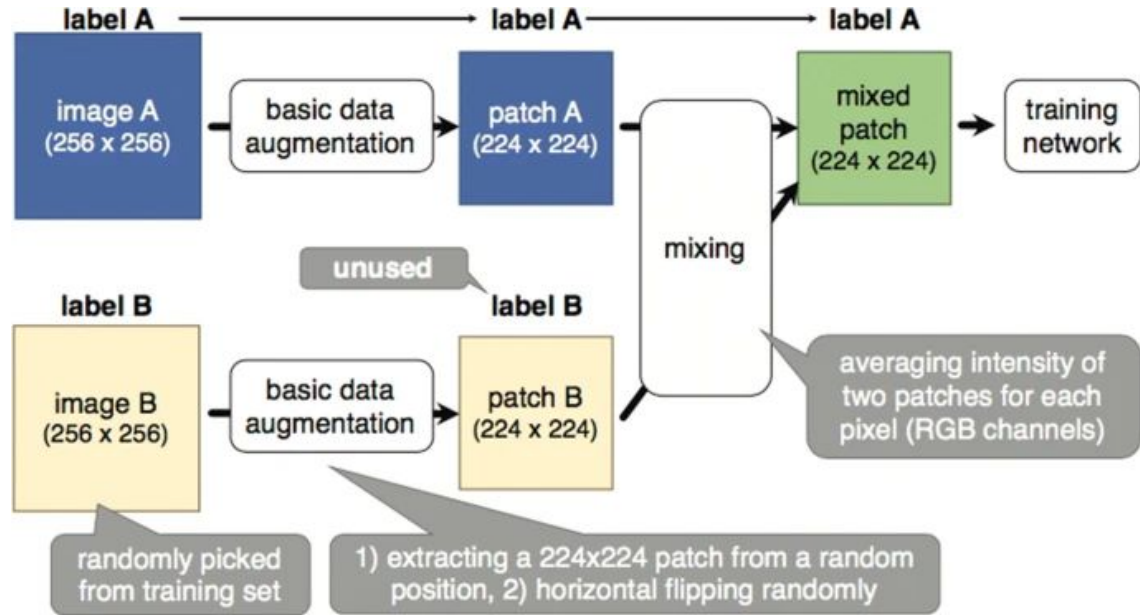
Can object segmentation help do a better job?

We can try various approaches to segmentation as a **preprocessing** step, and given the **same vector embedding model** see if it improves **precision**

Addressing edge-cases: data transformation?

We can try data transformation techniques, from **basic image manipulation** (contrast, saturation, etc)







to **more opinionated ones** (e.g. when a record has 3 images, today we index separately and return the overall best score - how about **MixingImages**?)



We can get creative
You can afford to experiment

You can do it too

What customer problems do you solve?

- Try to find data that describes **good/bad outcomes** 
 - Else, search for data close to the **use-case you're trying to solve** 
- What does **success look like?** This is your **metric!**
 - “**always having something to show**” -> **NoResultsRate@K** 
 - “**never showing Apple to Samsung products?**” -> **CategorySpread@K** 
 - “**optimize for TopCategoryXXX while keeping average quality under control?**”
-> **two metrics** **TopCategoryXXXRecall**  while keeping **CategoryPrecision** 

A metric is **right** when it captures the **value you build**



Conclusion: You *can* build a benchmark

Remember: a **hacky** benchmark
is better than **nothing**

hence **any data is better than no data!**



Thank you

Grab my slides



Try our model



 pln@algolia.com

 [@PaulLouisNech](https://twitter.com/PaulLouisNech)