

The Performance Impact of KNN with Padding Variations

ZEYNEP RANA BASIBUYUK

QUESTION:

Can more efficient results be achieved by combining
different padding variations?

Overview

01

Collecting data

02

Cropping Images

03

Extracting Embeddings from
Cropped Data

04

Calculating metrics by using
FaissKnn

Used datasets

- Kaggle :
(product based)
 - small dataset
 - big dataset

Padding percentages used

- 0
- 5
- 10
- 15
- 20
- 25
- 30
- 35
- 40
- 45
- 50

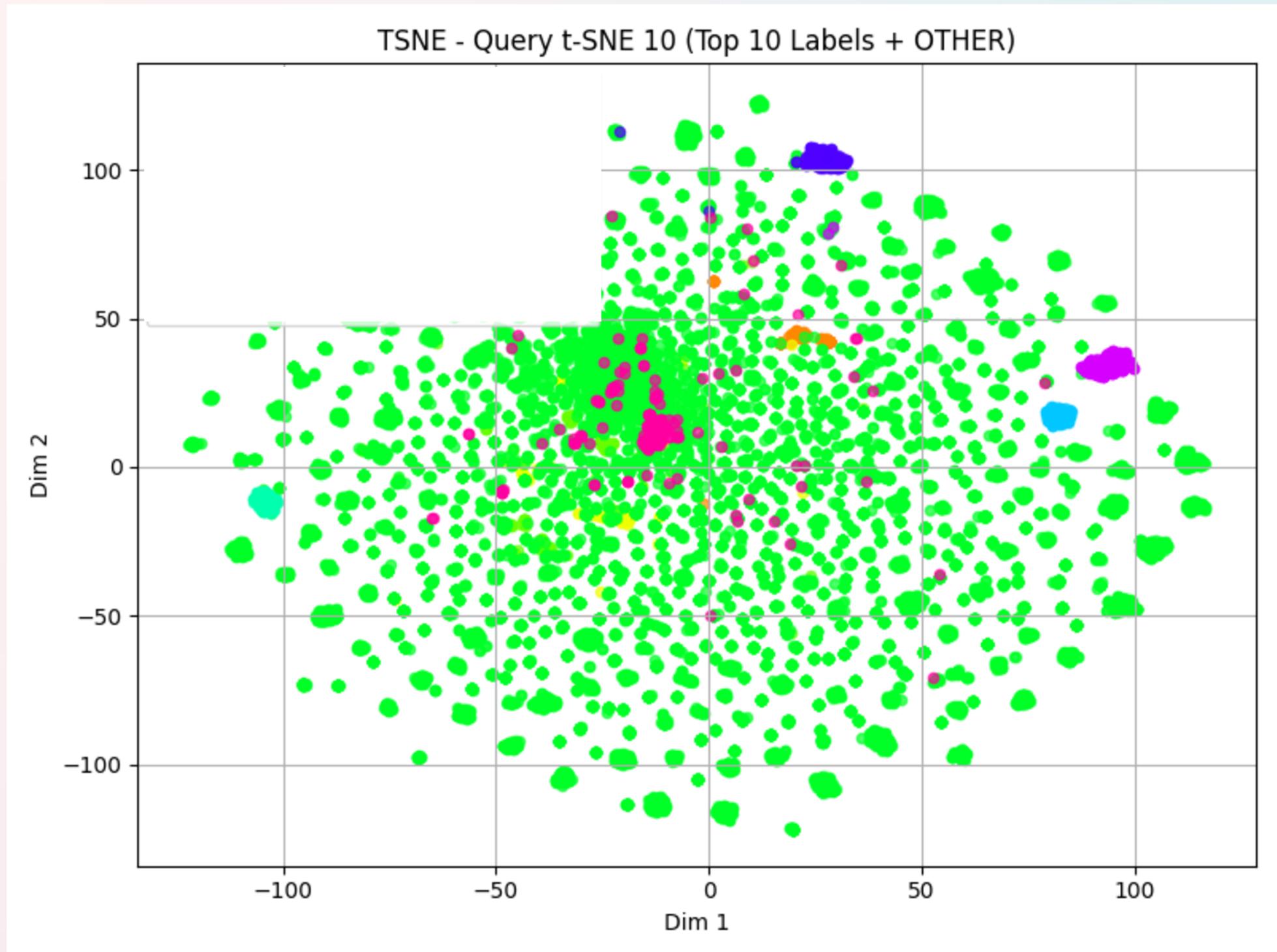
Extracted embeddings

- {Train + Validation}.npz
- {Test}.npz

- Accuracy percentage of
product-based classification

Small Dataset

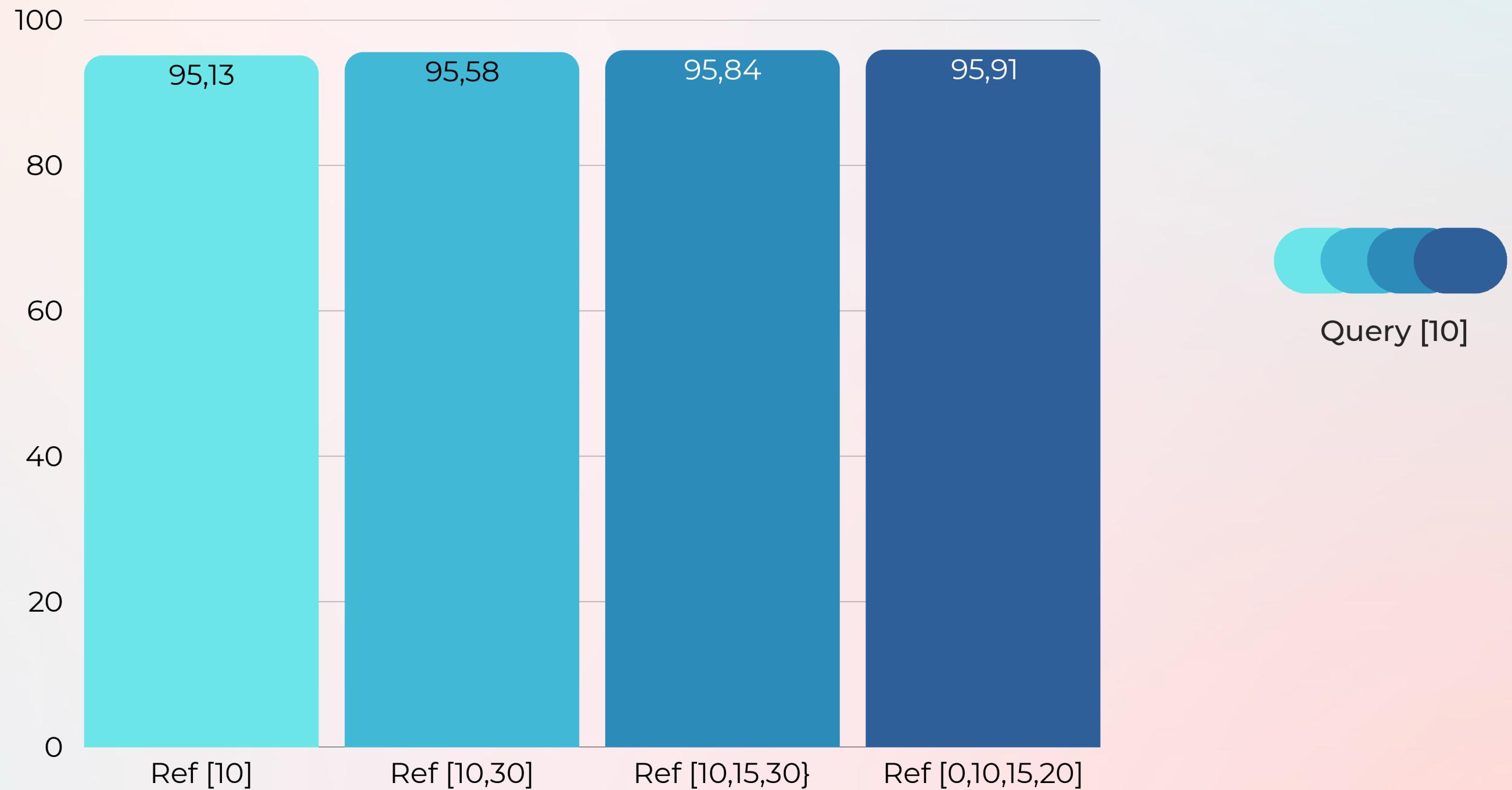
Ref	Query	Acc
[10,0]	[10]	0.9550
[10,0]	[0]	0.9548
[10,0]	[10,0]	0.9549



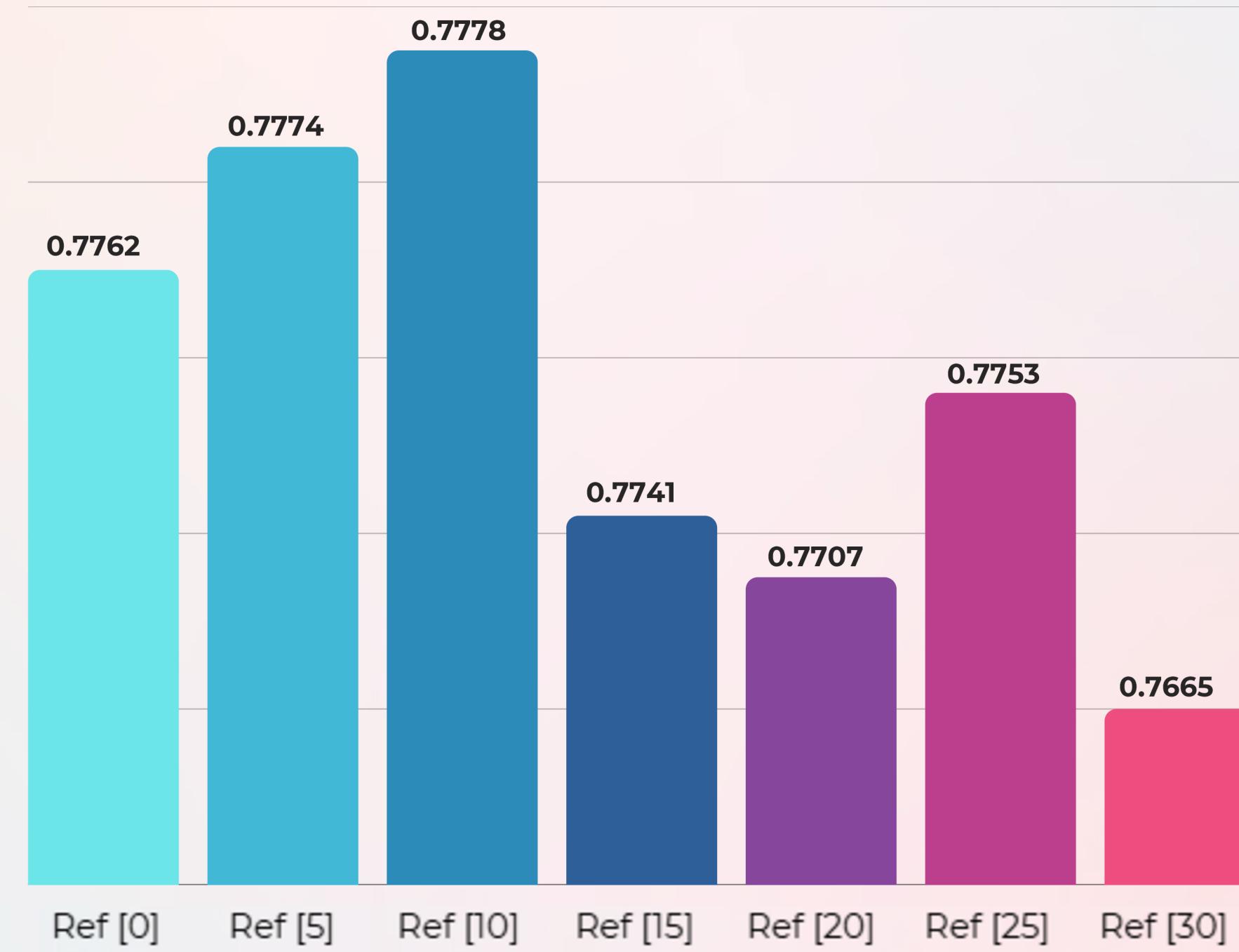
ref[0,10] query[10]

Small Dataset

Combinations with the best accuracy level

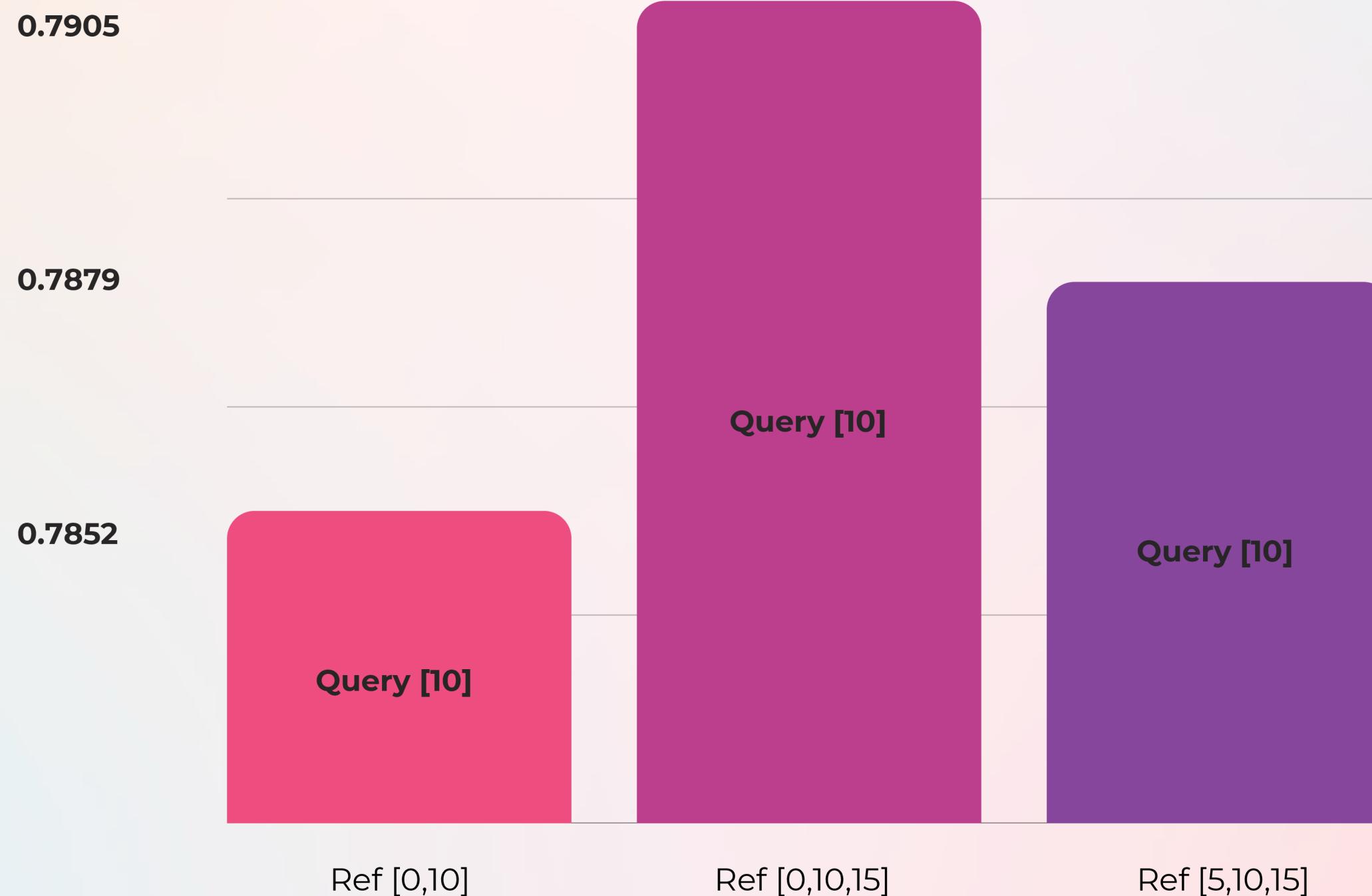


Big Dataset



Big Dataset

Combinations with the best accuracy level



Big Dataset

During metric evaluations, a parquet file was generated for each crop containing the following information: reference, query, id, class, prediction, and status.

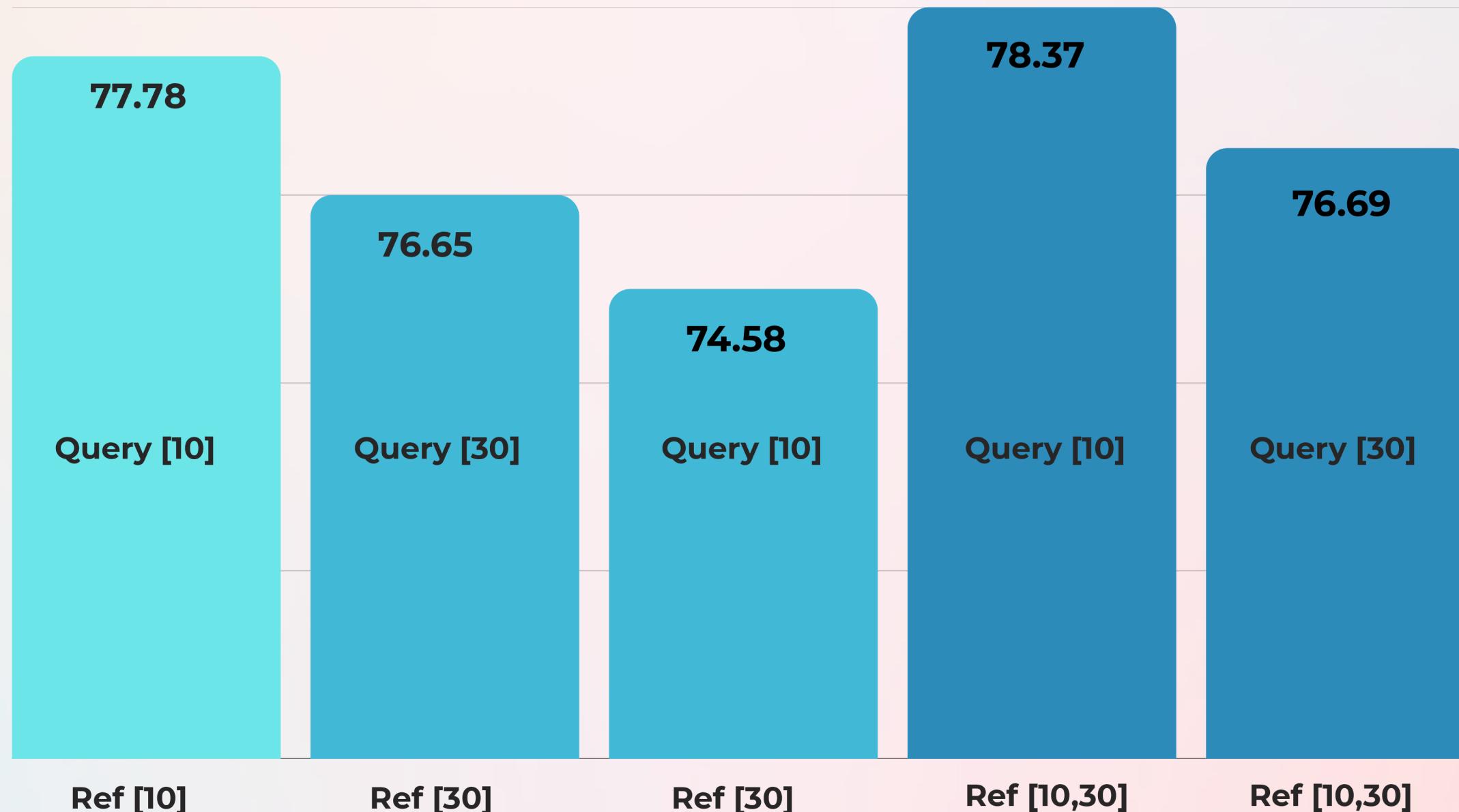
In the parquet file, records with the same ID were checked to see whether they had different statuses in other combination measurements. If such conflicts existed, they were saved into a new parquet file.

A new JSON file was created using some IDs as examples for cropping.

```
satir_bazli_sonuclar.csv
1 Ref_ids,Query_id,dolt_id,True,Pred,Status
2 [0],0,1744_60764,3550,3550,1
3 [0],0,1744_60767,3550,3550,1
4 [0],0,3242_108993,3550,6108,0
5 [0],0,3242_109007,3550,3550,1
6 [0],0,5012_166628,3550,3550,1
7 [0],0,5172_172013,3550,3550,1
8 [0],0,5172_172014,3550,3550,1
9 [0],0,546_19578,3550,3550,1
10 [0],0,546_19582,3550,3550,1
11 [0],0,6613_218203,3550,3550,1
12 [0],0,6613_218204,3550,3550,1
13 [0],0,8175_271573,3550,3550,1
14 [0],0,8700_289632,3550,3550,1
15 [0],0,8700_289658,3550,3550,1
16 [0],0,5052_167723,3745,1476,0
17 [0],0,5052_167724,3745,1476,0
18 [0],0,5052_167725,3745,1476,0
19 [0],0,5052_167726,3745,621,0
20 [0],0,5052_167727,3745,1476,0
21 [0],0,5052_167728,3745,1476,0
22 [0],0,5052_167729,3745,397,0
23 [0],0,5052_167730,3745,1476,0
24 [0],0,5052_167731,3745,6800,0
```

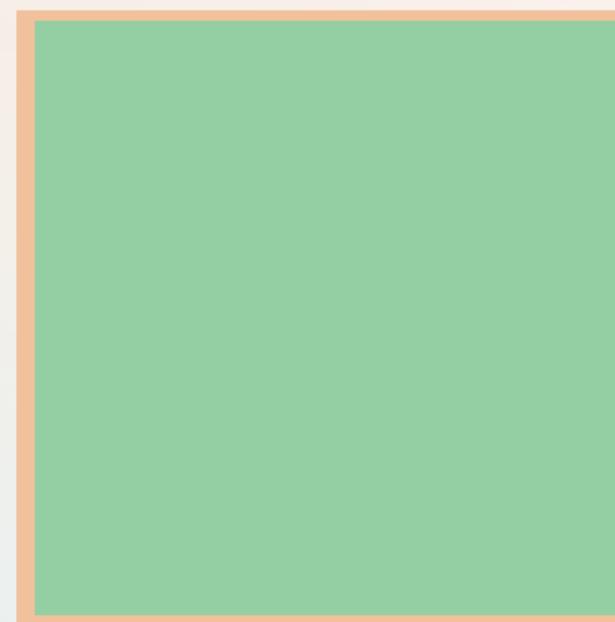
Big Dataset

Padding percentages that give the highest and lowest accuracy



Big Dataset

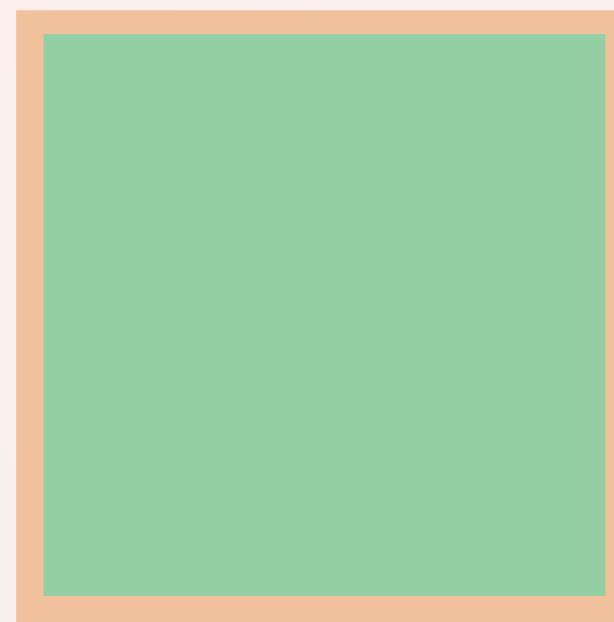
Some examples of turning false into true with combinations



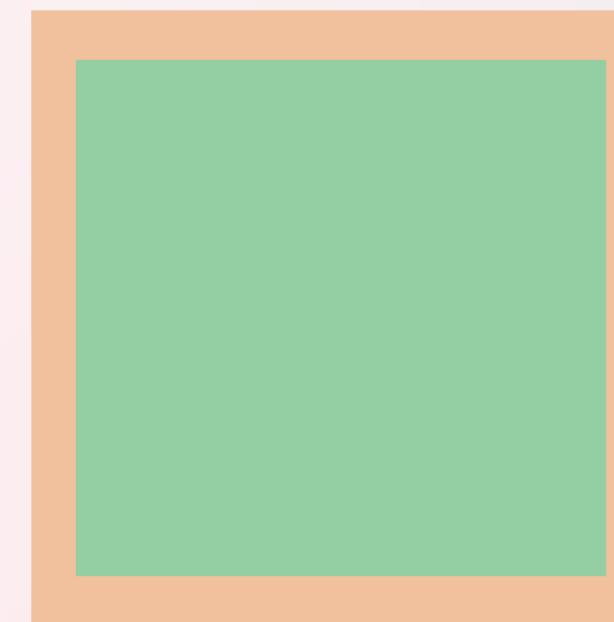
Padding_0



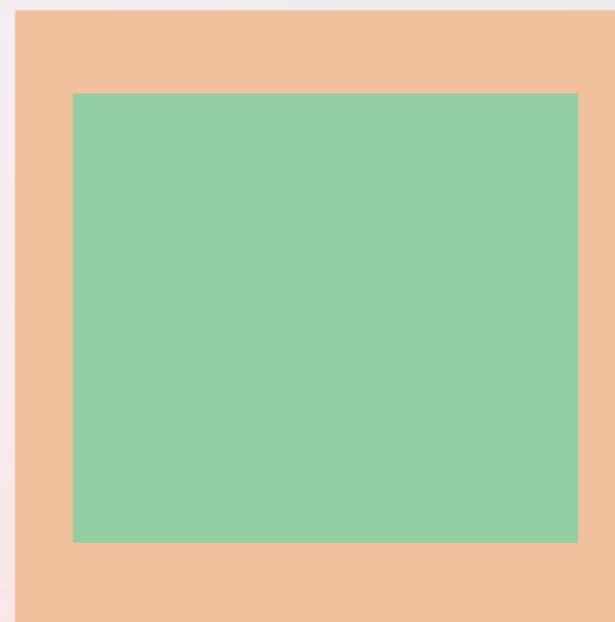
Model's prediction



Padding_5



Padding_10



Padding_15



Big Dataset

Some examples of turning false into true with combinations

Reference	Query	Status
0	0	✗
10	10	✓
15	15	✓
0-10	0	✗
0-10	10	✓
0-10-15	15	✓

Reference	Query	Status
0	0	✓
10	10	✗
15	15	✗
0-10	0	✓
0-10	10	✗
0-10-15	15	✗

6532_216105

228607_3258147

Big Dataset

Some examples of combinations that turn true into false

Reference	Query	Status
0	0	✓
10	10	✗
15	15	✓
0-10	0	✗
0-10	10	✗
0-10-15	15	✗

Reference	Query	Status
0	0	✗
10	10	✗
15	15	✓
0-10	0	✗
0-10	10	✗
0-10-15	15	✗

6642_219073

230147_3307706

Results

- Combining references and queries with different padding percentages showed improved SKU classification accuracy.
- Query [10] consistently gave the highest accuracy.
- Padding >20% generally reduced performance.
- Pairing the best query with the worst reference did not improve accuracy; best results occur when each padding uses its own reference/query.
- Some crops misclassified individually were correctly predicted in combinations, and vice versa. No clear pattern emerged, highlighting the complex behavior of ML models.

Files created throughout the project

- crop.py
- extract_embeddings.py
- evaluate_combo.py
- greedy_combo.py
- npz_to_parquet.py
- row_based_reluts.csv