

# Feature Selection for Unsupervised Domain Adaptation using Optimal Transport

LÉO GAUTHERON<sup>12</sup>, IEVGEN REDKO<sup>12</sup>, CAROLE LARTIZIEN<sup>2</sup>

{leo.gautheron, ievgen.redko}@univ-st-etienne.fr carole.lartizien@creatis.insa-lyon.fr

<sup>1</sup> Univ Lyon, UJM-Saint-Etienne, CNRS, Institut d'Optique Graduate School, Laboratoire Hubert Curien UMR 5516, F-42023, SAINT-ETIENNE, France

<sup>2</sup> Univ Lyon, INSA-Lyon, Université Claude Bernard Lyon 1, UJM-Saint Etienne, CNRS, Inserm, CREATIS UMR 5220, U1206, F-69621, LYON, France

## INTRODUCTION

### Issues of Traditional ML:

- near-human performance is achieved using lots of labeled data
- Some tasks do not have that much labeled data (biology, physics etc)
- There exists too many tasks!

### Solution: Domain adaptation

+ **Learn** when *labeled training set*  $S$  and *unlabeled test set*  $T$  do not follow **the same** probability distribution.

Accuracy: 84%



Amazon

Accuracy: 50%



Webcam

A **significant drop** in performance due to the **discrepancy** between training and test distributions!

## DISCRETE OPTIMAL TRANSPORT

Consider two empirical measures defined on  $S \in \mathbb{R}^{N_S \times d}$  and  $T \in \mathbb{R}^{N_T \times d}$  by

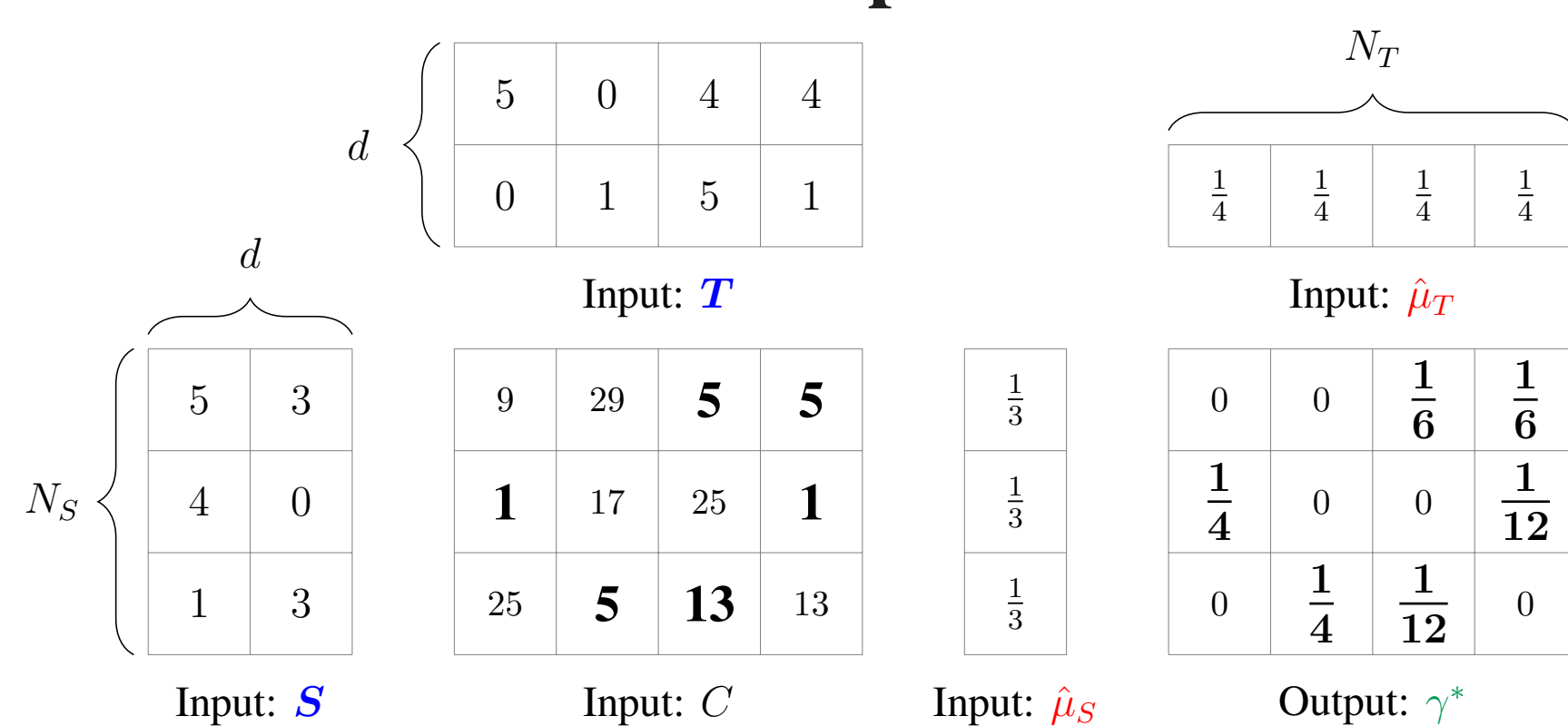
$$\hat{\mu}_S = \frac{1}{N_S} \sum_{i=1}^{N_S} \delta_{x_i^S} \text{ and } \hat{\mu}_T = \frac{1}{N_T} \sum_{i=1}^{N_T} \delta_{x_i^T}.$$

The goal of optimal transport (OT) is to find a coupling matrix  $\gamma^* \in \mathbb{R}_+^{N_S \times N_T}$  such that

$$\gamma^* = \arg \min_{\gamma \in \Pi(\hat{\mu}_S, \hat{\mu}_T)} \langle \gamma, C \rangle_F,$$

where  $C \in \mathbb{R}^{N_S \times N_T}$  is a transport cost with  $C_{ij}$  given by  $c: S \times T \rightarrow \mathbb{R}_+$  and  $\Pi(\hat{\mu}_S, \hat{\mu}_T) = \{\gamma \in \mathbb{R}_+^{N_S \times N_T} | \gamma \mathbf{1} = \hat{\mu}_S, \gamma^T \mathbf{1} = \hat{\mu}_T\}$ .

### Example



## OPTIMAL TRANSPORT VARIATIONS

### Entropy regularized OT [Cuturi 2013]

$$\gamma^* = \arg \min_{\gamma \in \Pi(\hat{\mu}_S, \hat{\mu}_T)} \langle \gamma, C \rangle_F - \frac{1}{\lambda} E(\gamma)$$

where  $E(\gamma) = -\sum_{ij} \gamma_{ij} \log \gamma_{ij}$  is the entropy of  $\gamma$ .

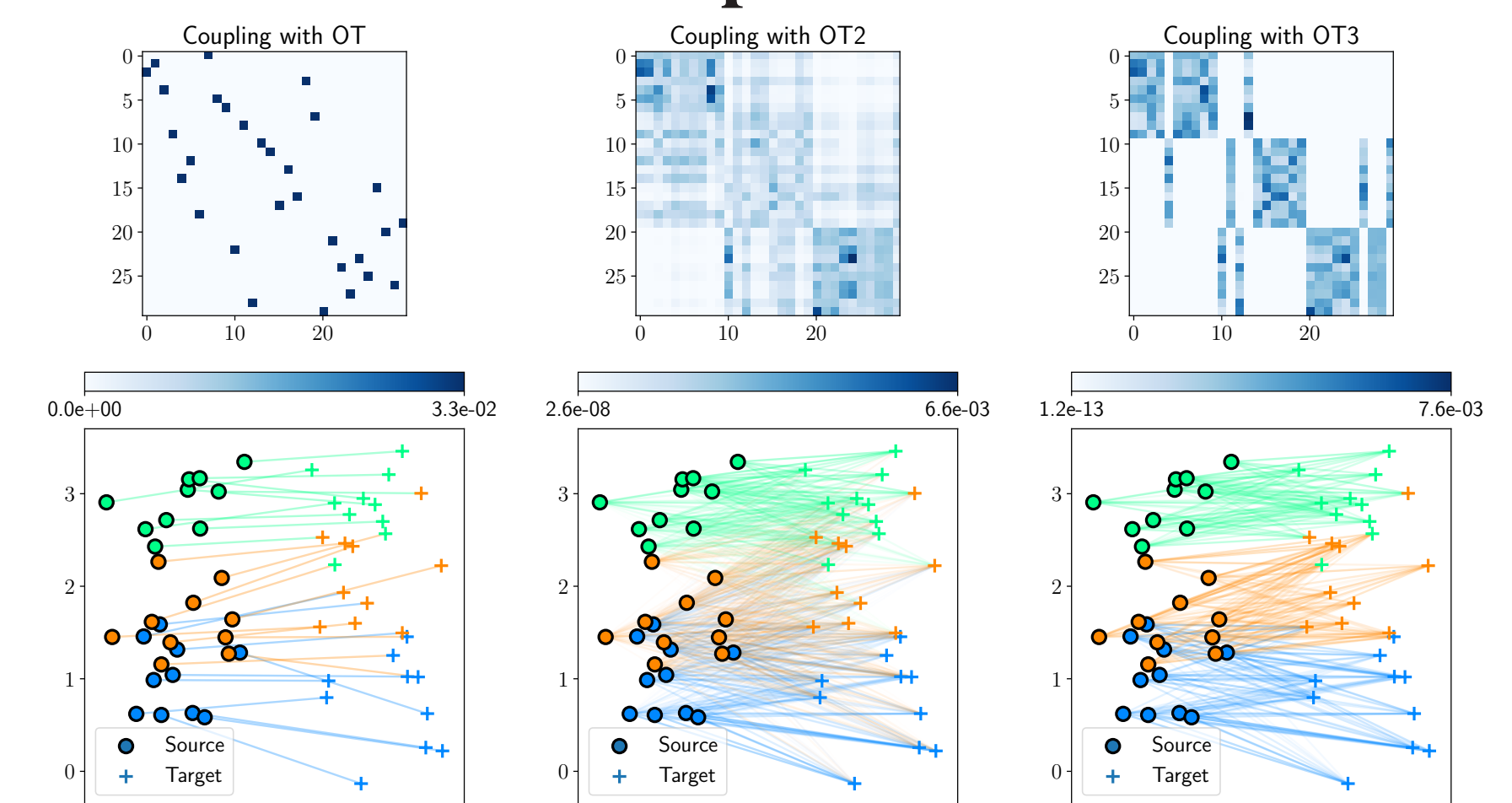
### Class regularized OT [Courty et al., 2014]

$$\gamma^* = \arg \min_{\gamma \in \Pi(\hat{\mu}_S, \hat{\mu}_T)} \langle \gamma, C \rangle_F - \frac{1}{\lambda} E(\gamma) + \eta \Omega(\gamma)$$

where  $\Omega(\gamma) = \sum_j \sum_c \|\gamma(I_c, j)\|_1^{0.5}$ .

In general, the **coupling** can be used to align  $S$  and  $T$  by using this reweighting:  $S \leftarrow N_S \gamma^* T$

### Comparison



## OUR CONTRIBUTION AND ITS ALGORITHMIC IMPLEMENTATION

### Motivation

When  $S$  and  $T$  are described by the same features, the **discrepancy** between them can be **reduced** by **finding and eliminating the most shifted features**.

### Proposed idea

**Step 1.** Find a coupling  $\gamma^* \in \mathbb{R}_+^{d \times d}$  between the features of  $S$  and  $T$ . The larger  $\gamma_{ii}^*$ , the most similar the feature number  $i$  between  $S$  and  $T$ .

**Step 2.** Sort the features by decreasing similarity between source  $S$  and target  $T$  domains given in the diagonal of  $\gamma^*$ . Keep the most similar.

### Step 1: Sample selection in target domain

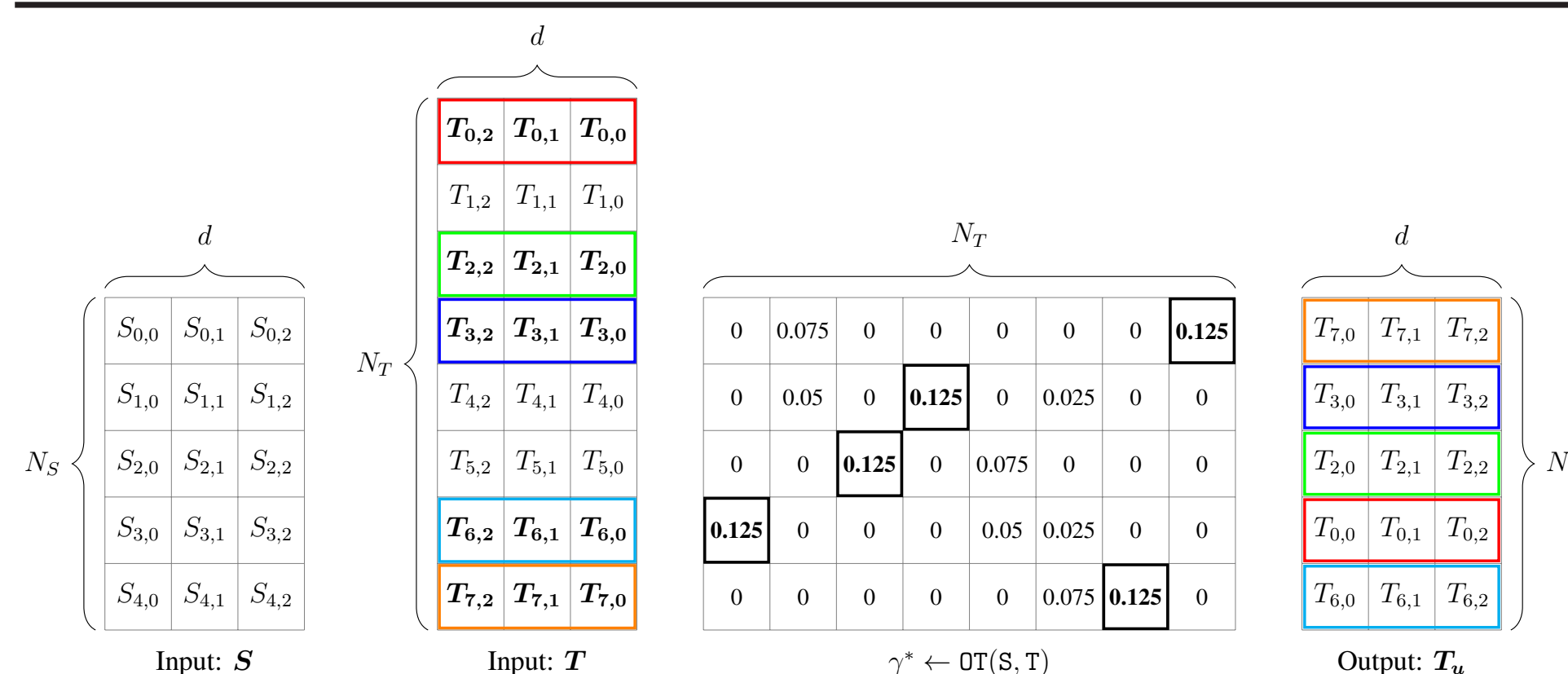
**Input** :  $S \in \mathbb{R}^{N_S \times d}, T \in \mathbb{R}^{N_T \times d}$

**Output** :  $T_u \in \mathbb{R}^{N_S \times d}$

$S = \text{zscore}(S); T = \text{zscore}(T)$

$\gamma^* \leftarrow \text{OT}(S, T)$

$T_u \leftarrow \{x_j \in T | j = \arg \max_{i=1, \dots, N_S} \gamma_{ij}^*\}$



### Step 2: Feature ranking for domain adaptation

**Input** :  $S \in \mathbb{R}^{N_S \times d}, T \in \mathbb{R}^{N_T \times d}$

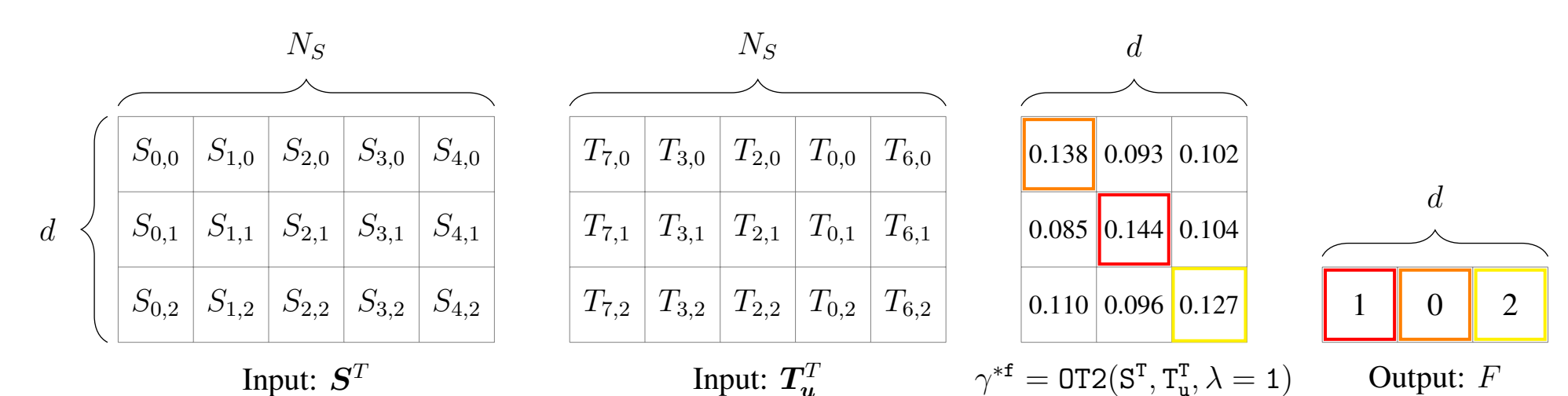
**Output** : List  $F$  of  $d$  most similar features from  $S$  and  $T$

$T_u \leftarrow \text{Algorithm1}(S, T)$

$S^T = \text{zscore}(S^T); T_u^T = \text{zscore}(T_u^T)$

$\gamma^{*f} = \text{OT2}(S^T, T_u^T, \lambda = 1)$

$F = \text{argSortDesc}(\{\{\gamma^{*f}\}_{ii} | i \in [1, d]\})$

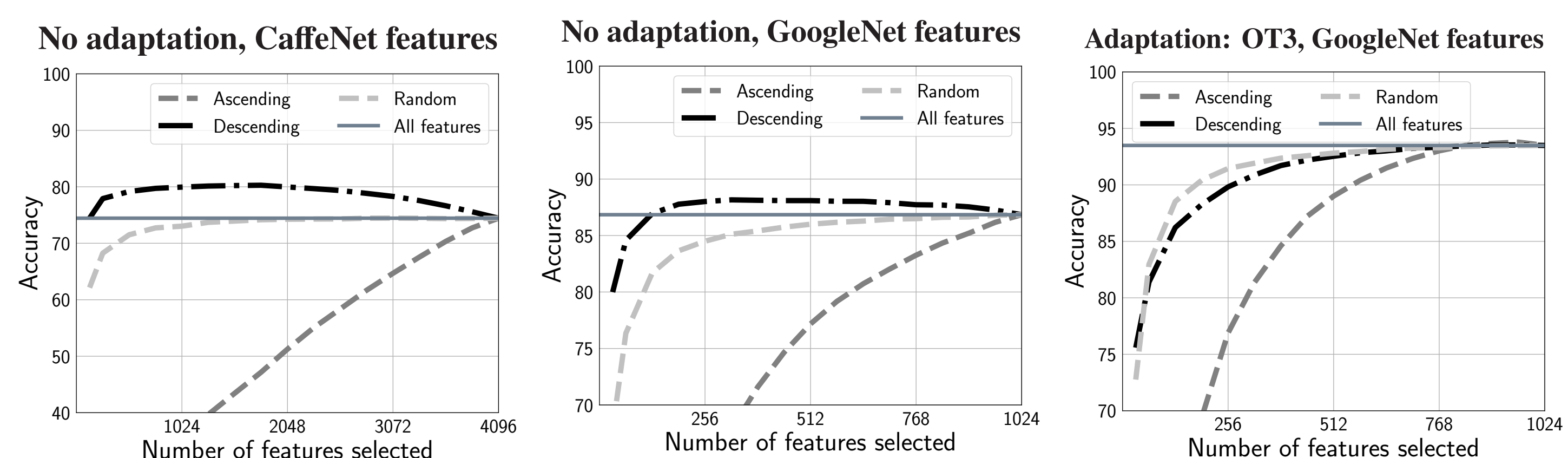


## OFFICE-CALTECH BENCHMARK

**Data:** Images from Amazon(A), Caltech(C), Webcam(W) and DSLR(D) datasets.



### Accuracy results



A→W	77.6±1.9	20.2±3.5	66.0±4.6	C→A	83.7±1.8	38.7±4.5	82.1±2.2	D→A	75.4±2.1	20.8±3.8	68.7±2.9	W→A	81.5±1.2	18.8±2.4	68.3±3.0
A→C	74.9±2.0	29.8±2.4	71.7±3.5	C→D	76.2±3.6	24.1±3.4	74.2±4.9	D→C	65.0±2.6	21.5±2.5	66.6±1.8	W→C	72.2±1.1	23.4±2.1	61.2±2.1
A→D	78.8±3.5	20.4±2.8	76.0±3.5	C→W	75.4±3.5	20.3±3.2	70.3±5.3	D→W	92.6±2.0	32.8±5.1	91.9±1.9	W→D	96.5±1.5	49.7±3.2	96.3±1.0
Pairs	↘ 512	↗ 512	4096									Mean	79.2±2.2	26.7±3.3	74.4±3.0

### Speed-up comparison

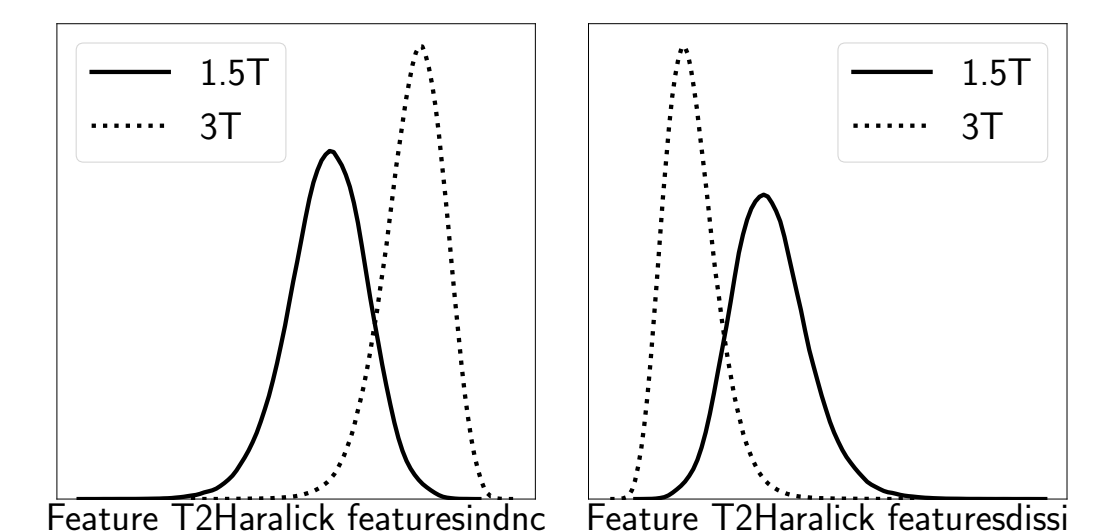
Method	↘ 512	↘ 1024	↘ 2048	4096
No adapt.	79.2±2.2	0.00s	79.9±2.3	0.00s
CORAL	80.5±1.8	110.43s	80.8±1.9	587.69s
SA	81.8±2.0	13.25s	82.5±1.8	32.09s
TCA	83.5±2.2	221.08s	85.0±1.9	223.62s
OT3	84.2±2.4	19.50s	86.7±1.9	31.76s

## MEDICAL APPLICATION: PROSTATE CANCER MAPPING

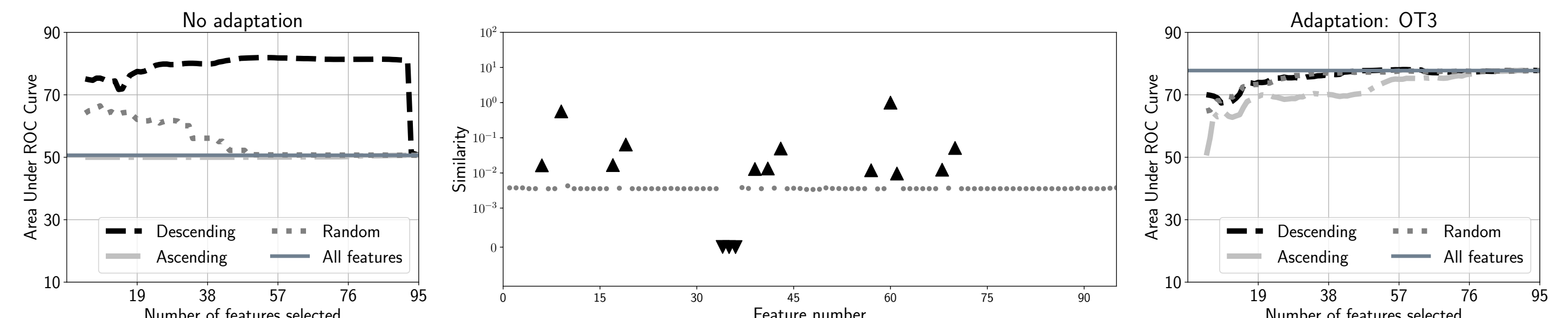
**Data:** Images from 1.5T and 3T MRI scanners of different resolution with 95 handcrafted features per voxel (3D pixels).

**Goal:** Learn on 1.5T voxels to predict cancer voxels in 3T images.

Class	#voxels 1.5T	#voxels 3T
Non cancer	363,222	846,556
Cancer	56,126	140,840
Total	419,348	987,396



### Obtained results



## CONCLUSION

- + Learning from the selected features gives **improved** performances in **less time** without adaptation, and **similar** performances in **less time** when adapting.
- + **Interpretable results** by identifying the most shifted original features.

Try it!

Our source code is available at <https://leogautheron.github.io> (requires Python Optimal Transport Library <https://github.com/rflamary/POT>)