



---

# Robust Methods for Medical Image Segmentation

Neerav Karani  
Krishna Chaitanya  
Dr. Ertunc Erdil  
Dr. Christian Baumgartner  
Prof. Ender Konukoglu

Biomedical Image Computing Group  
Computer Vision Lab  
ETH Zurich

March 2021

Once upon a time...

School exam



# Once upon a time...

## School exam



<Statement 1>	<b>TRUE</b>	FALSE
<Statement 2>	<b>TRUE</b>	FALSE
<Statement 3>	TRUE	<b>FALSE</b>
<Statement 4>	<b>TRUE</b>	FALSE
<Statement 5>	TRUE	<b>FALSE</b>
<Statement 6>	TRUE	<b>FALSE</b>

Studying at home

# Once upon a time...



## School exam



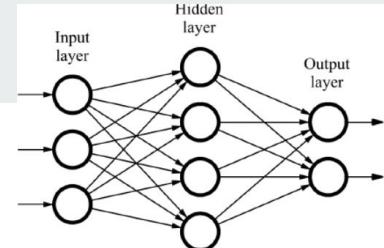
<Statement 1>	<b>TRUE</b>	FALSE
<Statement 2>	<b>TRUE</b>	FALSE
<Statement 3>	TRUE	<b>FALSE</b>
<Statement 4>	<b>TRUE</b>	FALSE
<Statement 5>	TRUE	<b>FALSE</b>
<Statement 6>	TRUE	<b>FALSE</b>

<Statement 3>	<b>TRUE</b>	FALSE
<Statement 6>	<b>TRUE</b>	FALSE
<Statement 1>	TRUE	<b>FALSE</b>
<Statement 5>	<b>TRUE</b>	FALSE
<Statement 2>	TRUE	<b>FALSE</b>
<Statement 4>	TRUE	<b>FALSE</b>

Studying at home

In the exam

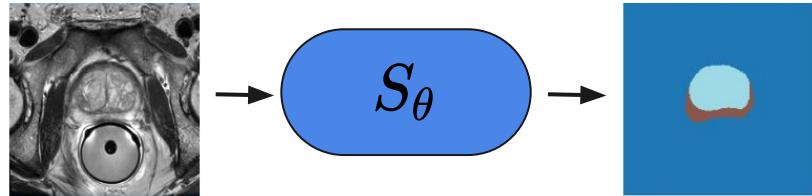
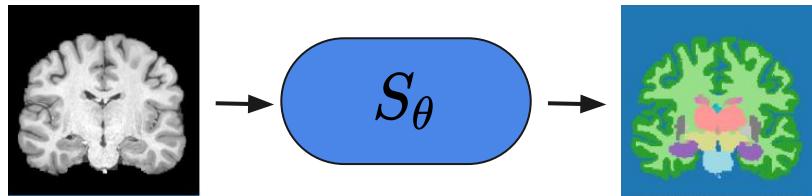
# Neural Networks behave similarly?!



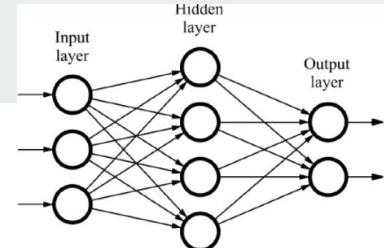
## Image segmentation



Training domains



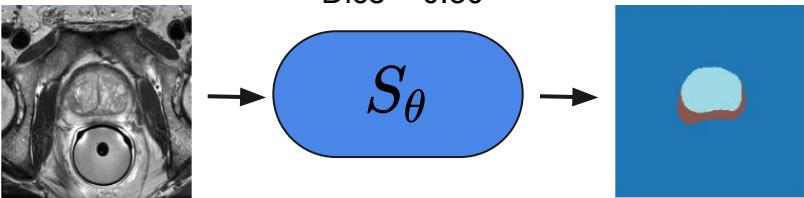
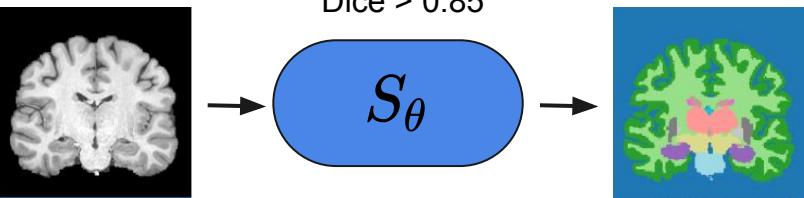
# Neural Networks behave similarly?!



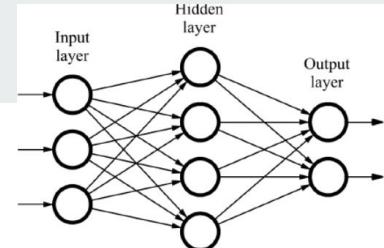
## Image segmentation



Training domains



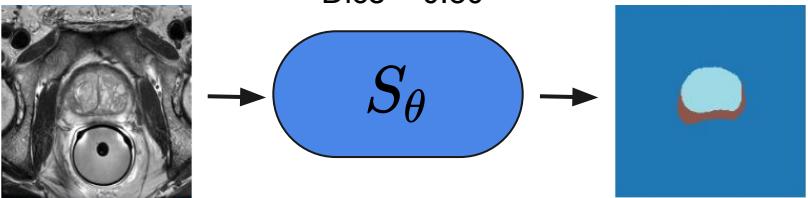
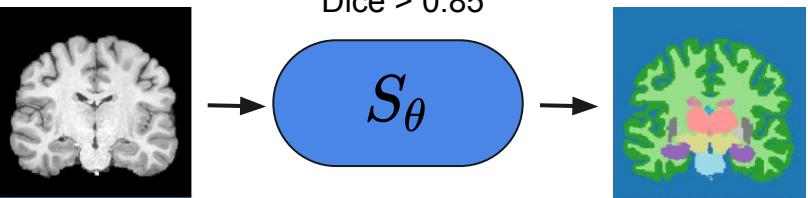
# Neural Networks behave similarly?!



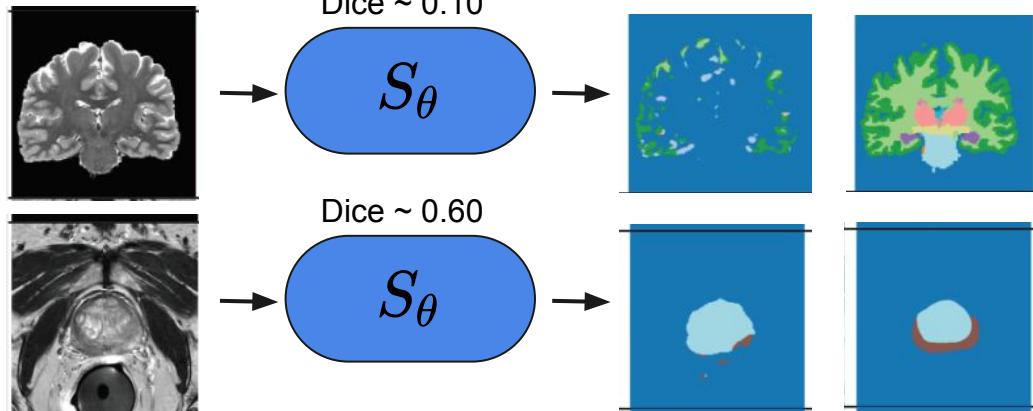
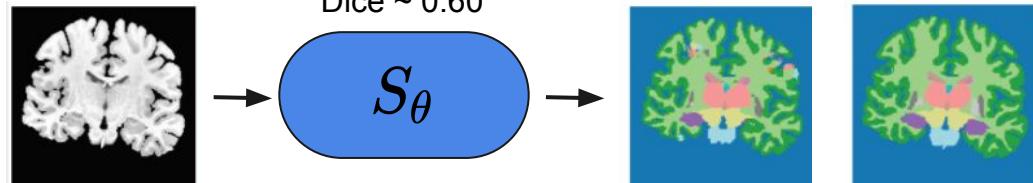
## Image segmentation



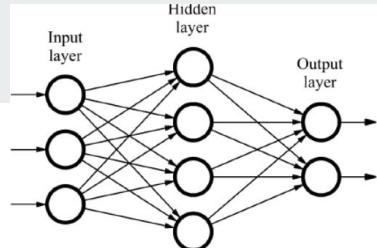
Training domains



Test domains

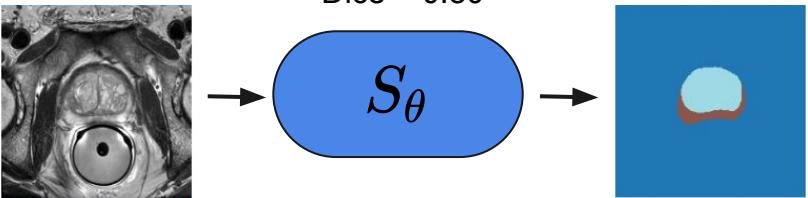
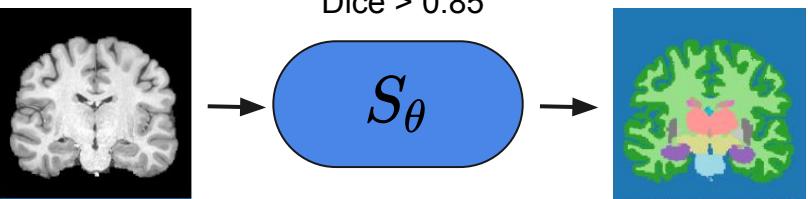


# Neural Networks behave similarly?!

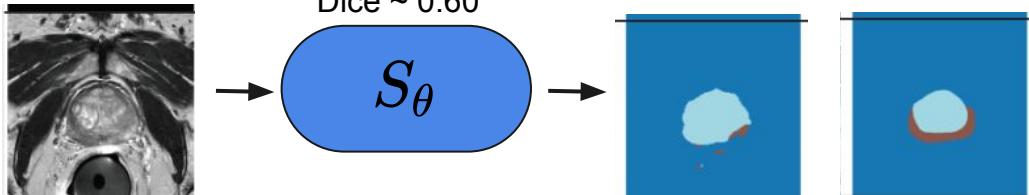
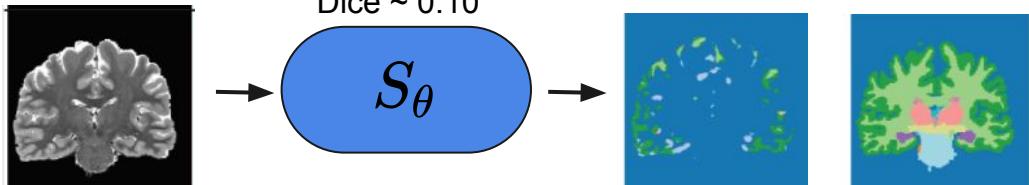
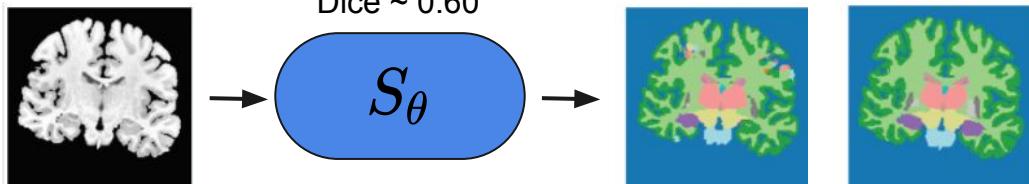


## Image segmentation

Training domains



Test domains



Even small differences between image statistics can lead to substantial performance degradation. CNNs are not robust against domain shifts.

## Domain Shift in MRI

---

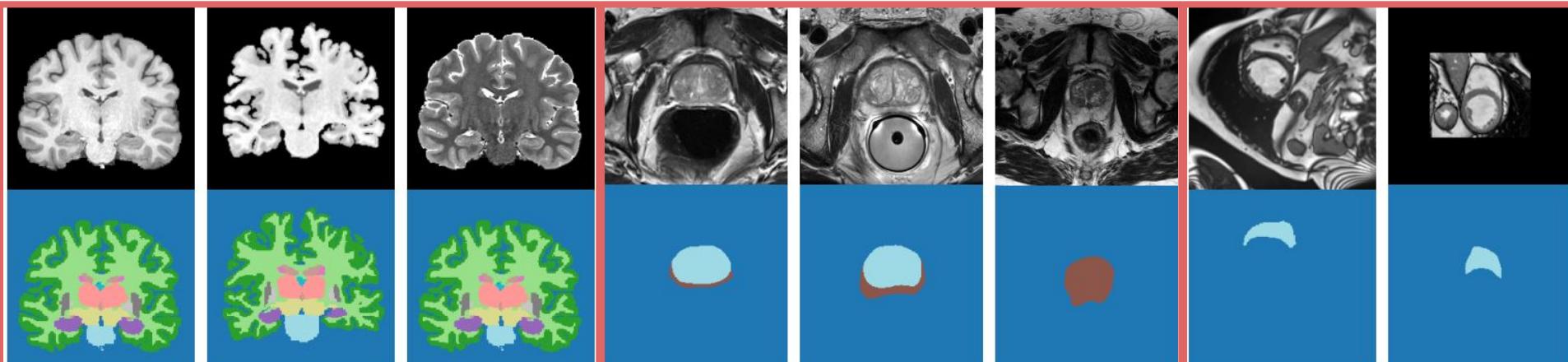
**Table 2 Types of dataset shift**

Type	Direction	Change	Examples of differences
Population shift	causal	$P_D(Z)$	ages, sexes, diets, habits, ethnicities, genetics
Annotation shift	causal	$P_D(Y X)$	annotation policy, annotator experience
Prevalence shift	anticausal	$P_D(Y)$	case-control balance, target selection
Manifestation shift	anticausal	$P_D(Z Y)$	anatomical manifestation of the target disease or trait
Acquisition shift	either	$P_D(X Z)$	scanner, resolution, contrast, modality, protocol

We focus on acquisition-related domain shifts.

# Introduction

## Domain Shift in MRI due to differences in acquisition settings



Brain MRIs from different scanners and protocols.

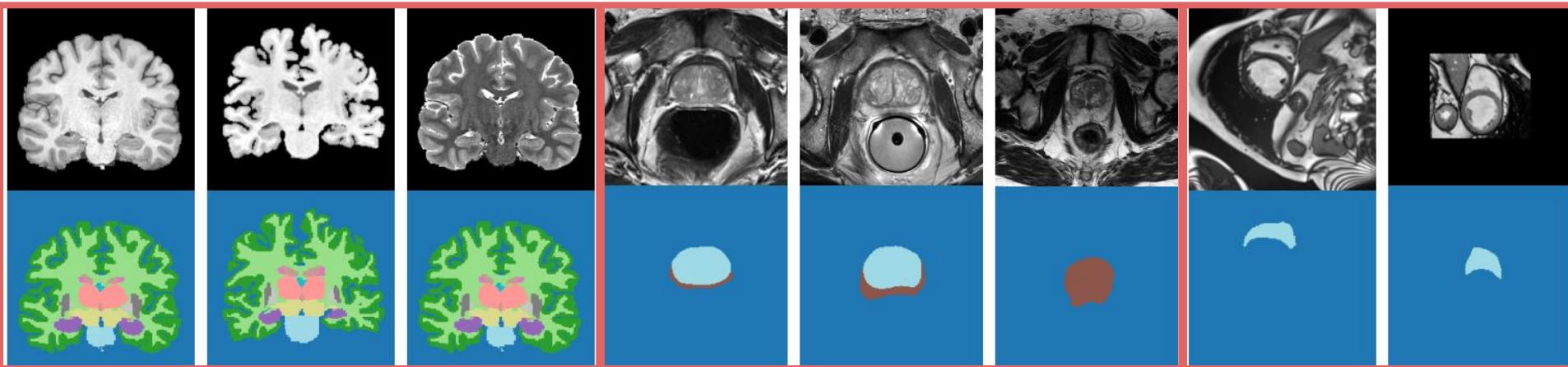
Prostate MRIs from different scanners.

Cardiac MRIs from different scanners.

# Introduction

## Domain Shift in MRI due to differences in acquisition settings

CNNs do not generalize across such domain shifts.



Anatomy		Brain			Prostate (whole)			Prostate (sep.)		Heart	
Train	Test	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	SD	TD <sub>1</sub>
Baseline and Benchmark											
SD (Baseline)		0.853	0.588	0.107	0.840	0.586	0.609	0.722	0.544	0.823	0.670
TD <sub>n</sub> (Benchmark)		-	0.896	0.867	-	0.817	0.834	-	0.732	-	0.806

# Learning setups for dealing with domain shifts

## ML Jargon



Setup	At SD		At TD	
	Data	Algorithm	Data	Algorithm
Separate training for each domain	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (many)	$\min_{\theta} L_{TD}$

# Learning setups for dealing with domain shifts

## ML Jargon



Setup	At SD		At TD	
	Data	Algorithm	Data	Algorithm
Separate training for each domain	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (many)	$\min_{\theta} L_{TD}$
Transfer learning	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (few)	Init. at $\theta_{SD}^*, \min_{\theta} L_{TD}$

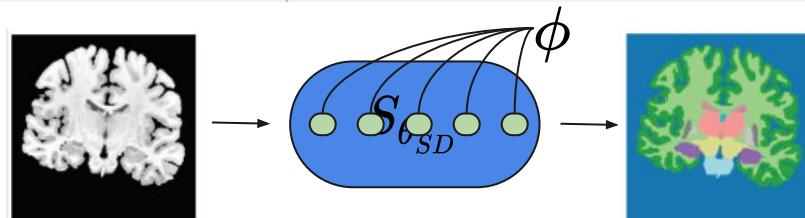
# Learning setups for dealing with domain shifts

Karani, N, et al. "A lifelong learning approach to brain MR segmentation across scanners and protocols." *MICCAI*, 2018.

## ML Jargon

---

Setup	At SD		At TD	
	Data	Algorithm	Data	Algorithm
Separate training for each domain	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (many)	$\min_{\theta} L_{TD}$
Transfer learning	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (few)	Init. at $\theta_{SD}^*$ , $\min_{\theta} L_{TD}$



# Learning setups for dealing with domain shifts

## ML Jargon



Setup	At SD		At TD	
	Data	Algorithm	Data	Algorithm
Separate training for each domain	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (many)	$\min_{\theta} L_{TD}$
Transfer learning	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (few)	Init. at $\theta_{SD}^*$ , $\min_{\theta} L_{TD}$
Unsupervised domain adaptation	-	-	$(X_{SD}, Y_{SD})$ (many), $(X_{TD})$ (many)	$\min_{\theta} L_{SD} + L_{invariant features}$

# Learning setups for dealing with domain shifts

## ML Jargon



Setup	At SD		At TD	
	Data	Algorithm	Data	Algorithm
Separate training for each domain	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (many)	$\min_{\theta} L_{TD}$
Transfer learning	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (few)	Init. at $\theta_{SD}^*$ , $\min_{\theta} L_{TD}$
Unsupervised domain adaptation	-	-	$(X_{SD}, Y_{SD})$ (many), $(X_{TD})$ (many)	$\min_{\theta} L_{SD} + L_{invariant features}$
Domain generalization	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$ $+ L_{invariant features}$	$(X_{TI})$ (one)	Predict $\hat{Y} = S_{\theta_{SD}^*}(X_{TI})$

# Learning setups for dealing with domain shifts

## ML Jargon



Setup	At SD		At TD	
	Data	Algorithm	Data	Algorithm
Separate training for each domain	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (many)	$\min_{\theta} L_{TD}$
Transfer learning	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (few)	Init. at $\theta_{SD}^*$ , $\min_{\theta} L_{TD}$
Unsupervised domain adaptation	-	-	$(X_{SD}, Y_{SD})$ (many), $(X_{TD})$ (many)	$\min_{\theta} L_{SD} + L_{invariant features}$
Domain generalization	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$ + $L_{invariant features}$	$(X_{TI})$ (one)	Predict $\hat{Y} = S_{\theta_{SD}^*}(X_{TI})$
Domain generalization (with test-time adaptation)	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$ + $L_{invariant features}$	$(X_{TI})$ (one)	Adapt $\theta_{SD}^*$ for each test image (TI). Predict $\hat{Y} = S_{\theta_{TI}^*}(X_{TI})$

# Learning setups for dealing with domain shifts

## ML Jargon

---

Setup	At SD		At TD	
	Data	Algorithm	Data	Algorithm
Separate training for each domain	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (many)	$\min_{\theta} L_{TD}$
Transfer learning	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (few)	Init. at $\theta_{SD}^*$ , $\min_{\theta} L_{TD}$
Unsupervised domain adaptation	-	-	$(X_{SD}, Y_{SD})$ (many), $(X_{TD})$ (many)	$\min_{\theta} L_{SD} + L_{invariant features}$
Domain generalization	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$ + $L_{invariant features}$	$(X_{TI})$ (one)	Predict $\hat{Y} = S_{\theta_{SD}^*}(X_{TI})$
Domain generalization (with test-time adaptation)	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$ + $L_{invariant features}$	$(X_{TI})$ (one)	Adapt $\theta_{SD}^*$ for each test image (TI). Predict $\hat{Y} = S_{\theta_{TI}^*}(X_{TI})$
Unsupervised Learning	-	-	$(X_{TI})$ (one)	Optimize parameters for each TI. Predict $\hat{Y} = \text{argmax}_Y P(Y)P(X_{TI} Y)$

# Learning setups for dealing with domain shifts

## ML Jargon

Setup	At SD		At TD	
	Data	Algorithm	Data	Algorithm
Separate training for each domain	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (many)	$\min_{\theta} L_{TD}$
Transfer learning	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (few)	Init. at $\theta_{SD}^*$ , $\min_{\theta} L_{TD}$
Unsupervised domain adaptation	-	-	$(X_{SD}, Y_{SD})$ (many), $(X_{TD})$ (many)	$\min_{\theta} L_{SD} + L_{invariant features}$
Domain generalization	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$ $+ L_{invariant features}$	$(X_{TI})$ (one)	Predict $\hat{Y} = S_{\theta_{SD}^*}(X_{TI})$
Domain generalization (with test-time adaptation)	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$ $+ L_{invariant features}$	$(X_{TI})$ (one)	Adapt $\theta_{SD}^*$ for each test image (TI). Predict $\hat{Y} = S_{\theta_{TI}^*}(X_{TI})$

We believe that for medical image segmentation, DG is more suitable than other settings.

# Learning setups for dealing with domain shifts

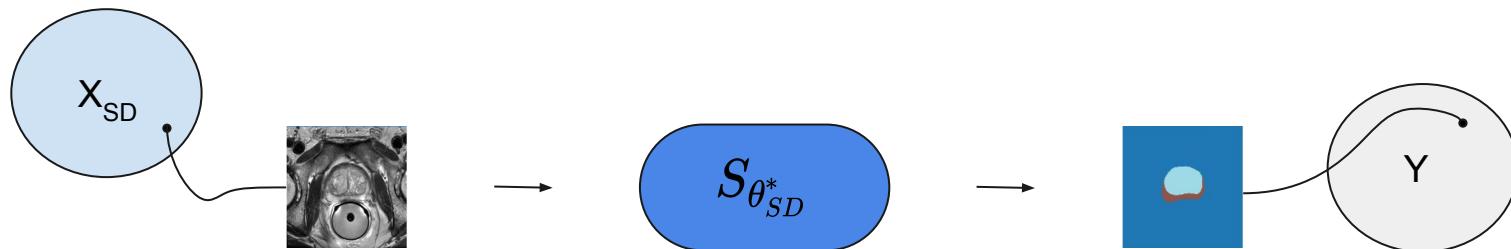
## ML Jargon

Setup	At SD		At TD	
	Data	Algorithm	Data	Algorithm
Separate training for each domain	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (many)	$\min_{\theta} L_{TD}$
Transfer learning	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$	$(X_{TD}, Y_{TD})$ (few)	Init. at $\theta_{SD}^*$ , $\min_{\theta} L_{TD}$
Unsupervised domain adaptation	-	-	$(X_{SD}, Y_{SD})$ (many), $(X_{TD})$ (many)	$\min_{\theta} L_{SD} + L_{invariant features}$
Domain generalization	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$ + $L_{invariant features}$	$(X_{TI})$ (one)	Predict $\hat{Y} = S_{\theta_{SD}^*}(X_{TI})$
Domain generalization (with test-time adaptation)	$(X_{SD}, Y_{SD})$ (many)	$\min_{\theta} L_{SD}$ + $L_{invariant features}$	$(X_{TI})$ (one)	Adapt $\theta_{SD}^*$ for each test image (TI). Predict $\hat{Y} = S_{\theta_{TI}^*}(X_{TI})$

# Related Work

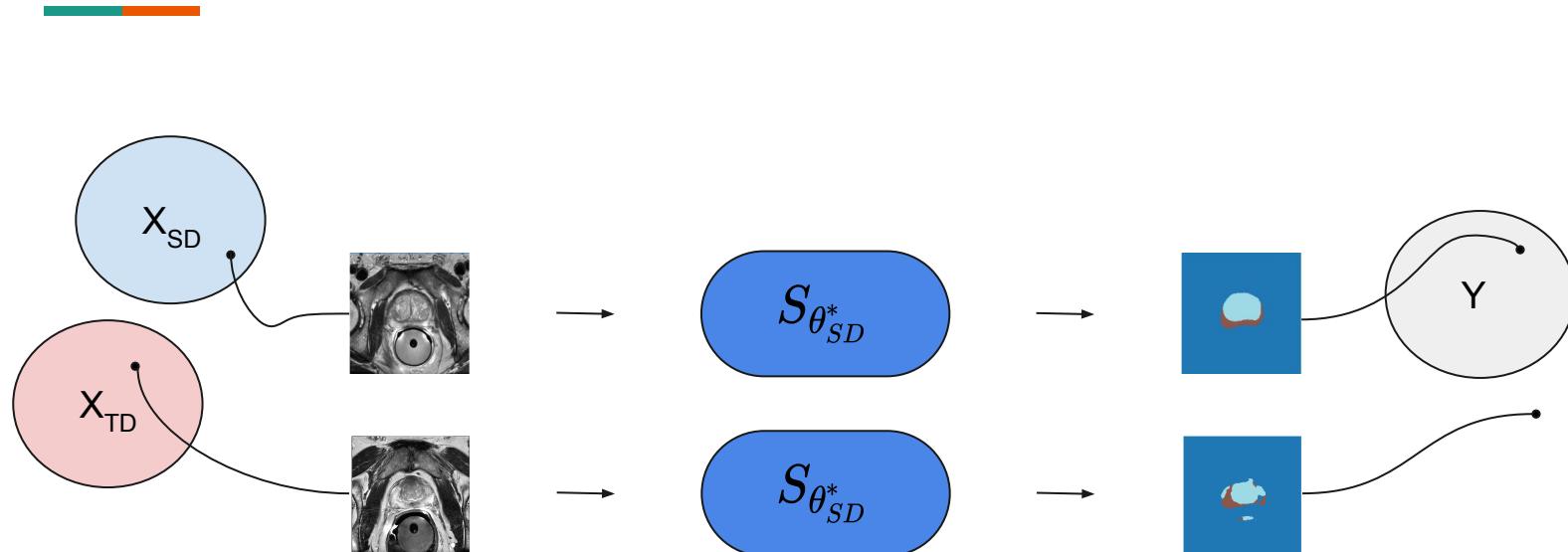
## Training strategies for domain generalization

—



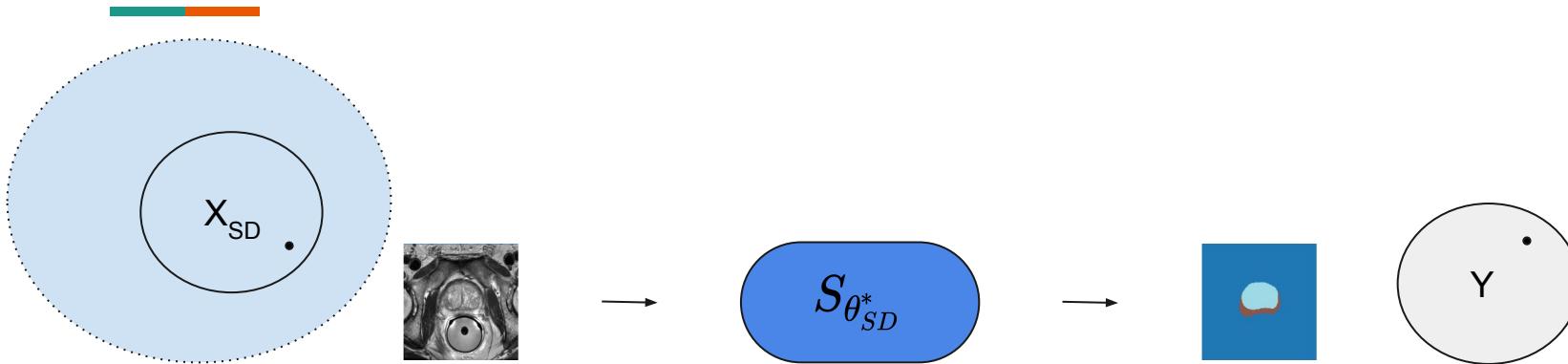
# Related Work

## Training strategies for domain generalization



# Related Work

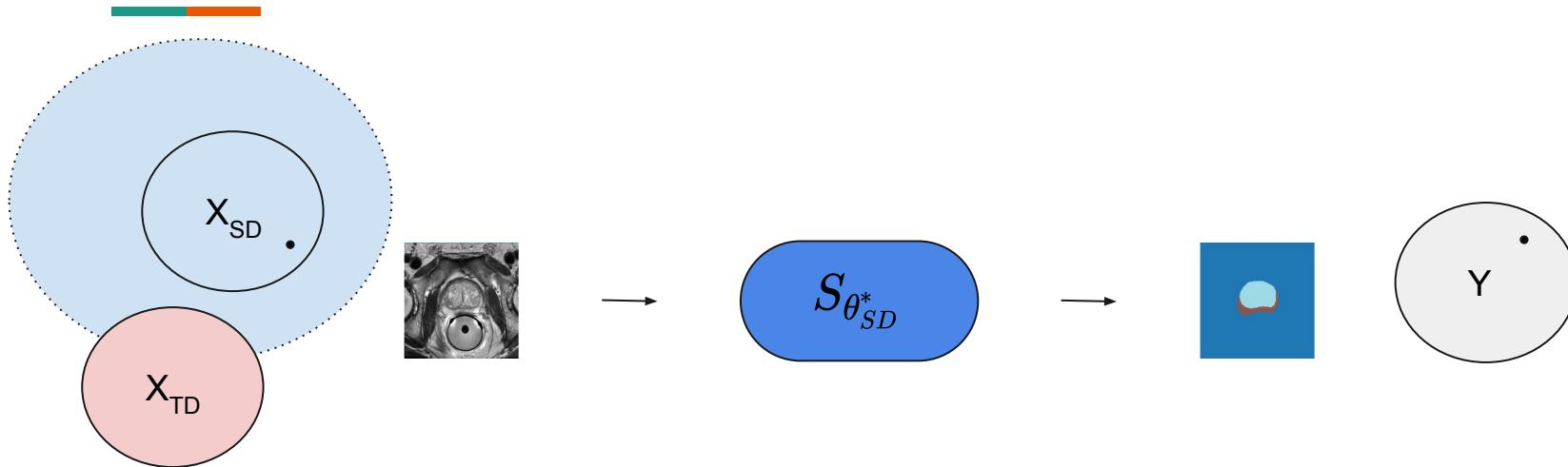
## Training strategies for domain generalization - Data Augmentation



- Zhang, L, et al. "Generalizing deep learning for medical image segmentation to unseen domains via deep stacked transformation." TMI 2020
- Billot, B, et al. "A Learning Strategy for Contrast-agnostic MRI Segmentation." MIDL 2020
- Hendrycks, D, et al. "AugMix: A Simple Data Processing Method to Improve Robustness and Uncertainty." ICLR 2019
- Yun, S, et al. "Cutmix: Regularization strategy to train strong classifiers with localizable features." ICCV 2019
- Zhang, H, et al. "mixup: Beyond Empirical Risk Minimization." ICLR 2018
- ...

## Related Work

### Training strategies for domain generalization - Data Augmentation

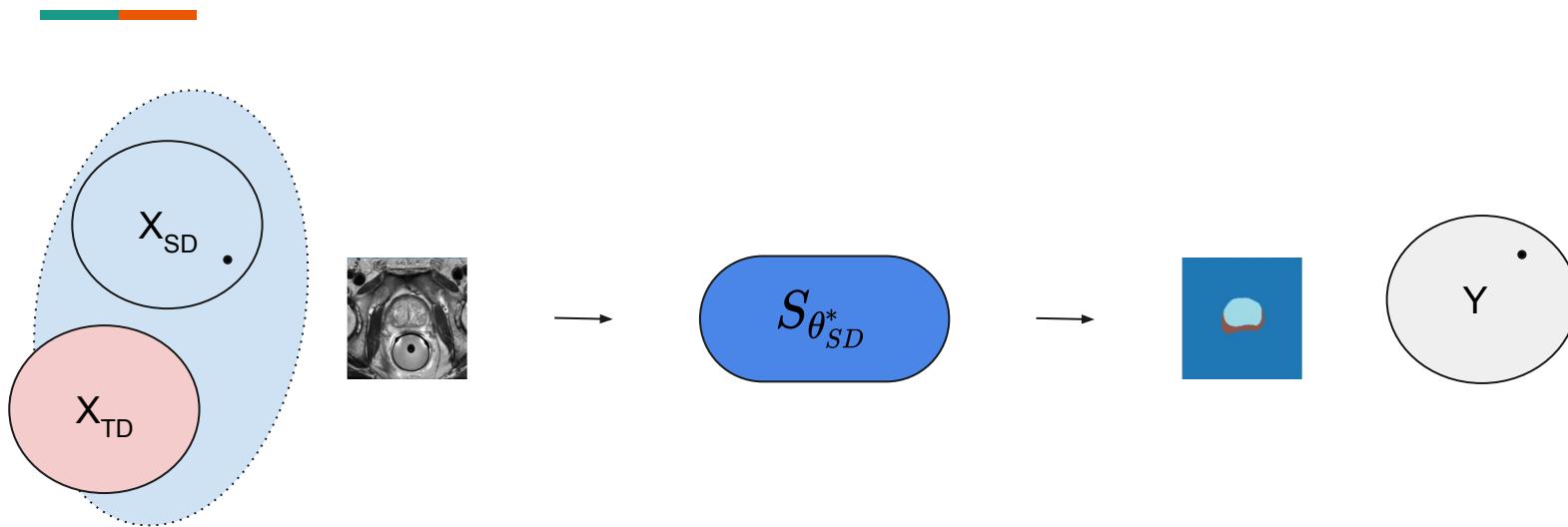


Heuristic data augmentation is a strong baseline.

However, the performance improvement is likely to be limited to the variations covered by the augmentation.

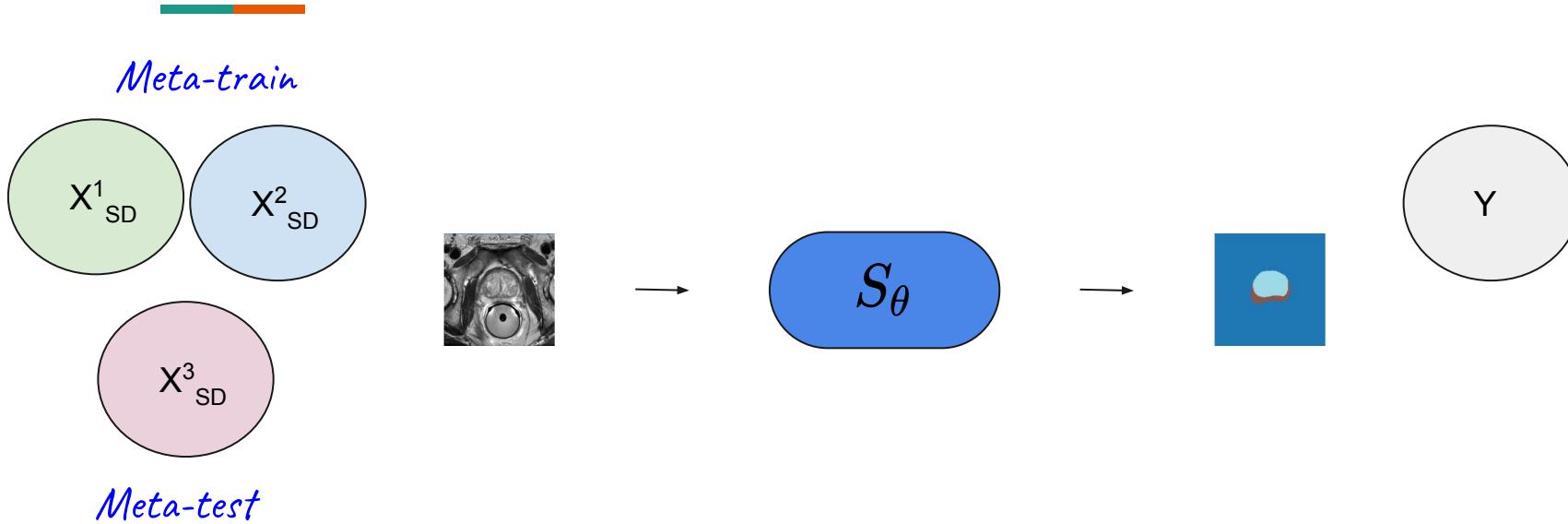
## Related Work

### Data Augmentation for unsupervised domain adaptation

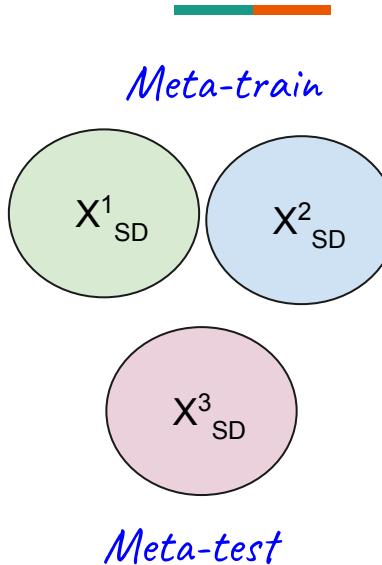


- Chaitanya, K, et al. "Semi-supervised and task-driven data augmentation." IPMI 2019.
- Zhao, A, et al. "Data augmentation using learned transformations for one-shot medical image segmentation." CVPR 2019.

## Training strategies for domain generalization - Meta Learning



## Training strategies for domain generalization - Meta Learning



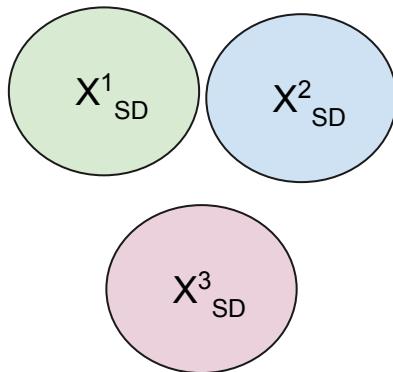
**Algorithm 1** Meta-Learning Domain Generalization

```
1: procedure MLDG
2:   Input: Domains  $\mathcal{S}$ 
3:   Init: Model parameters  $\Theta$ . Hyperparameters  $\alpha, \beta, \gamma$ .
4:   for ite in iterations do
5:     Split:  $\bar{\mathcal{S}}$  and  $\check{\mathcal{S}} \leftarrow \mathcal{S}$ 
6:     Meta-train: Gradients  $\nabla_{\Theta} = \mathcal{F}'_{\Theta}(\bar{\mathcal{S}}; \Theta)$ 
7:     Updated parameters  $\Theta' = \Theta - \alpha \nabla_{\Theta}$ 
8:     Meta-test: Loss is  $\mathcal{G}(\check{\mathcal{S}}; \Theta')$ .
9:     Meta-optimization: Update  $\Theta$ 
10:     $\Theta = \Theta - \gamma \frac{\partial(\mathcal{F}(\bar{\mathcal{S}}; \Theta) + \beta \mathcal{G}(\check{\mathcal{S}}; \Theta - \alpha \nabla_{\Theta}))}{\partial \Theta}$ 
11:   end for
12: end procedure
```

# Related Work

## Training strategies for domain generalization - Meta Learning

Meta-train



Meta-test

---

**Algorithm 1** Meta-Learning Domain Generalization

---

```
1: procedure MLDG
2:   Input: Domains  $\mathcal{S}$ 
3:   Init: Model parameters  $\Theta$ . Hyperparameters  $\alpha, \beta, \gamma$ .
4:   for ite in iterations do
5:     Split:  $\bar{\mathcal{S}}$  and  $\check{\mathcal{S}} \leftarrow \mathcal{S}$ 
6:     Meta-train: Gradients  $\nabla_\Theta = \mathcal{F}'_\Theta(\bar{\mathcal{S}}; \Theta)$ 
7:     Updated parameters  $\Theta' = \Theta - \alpha \nabla_\Theta$ 
8:     Meta-test: Loss is  $\mathcal{G}(\check{\mathcal{S}}; \Theta')$ .
9:     Meta-optimization: Update  $\Theta$ 

$$\Theta = \Theta - \gamma \frac{\partial(\mathcal{F}(\bar{\mathcal{S}}; \Theta) + \beta \mathcal{G}(\check{\mathcal{S}}; \Theta - \alpha \nabla_\Theta))}{\partial \Theta}$$

10:    end for
11: end procedure
```

---

- Li, D. et al. "Learning to generalize: Meta-learning for domain generalization." AAAI 2018
- Dou, Q. et al. "Domain generalization via model-agnostic learning of semantic features." NeurIPS 2019
- Tseng, H. et al. "Cross-Domain Few-Shot Classification via Learned Feature-Wise Transformation." ICLR 2020
- Zhang, J. et al. "Generalizable Semantic Segmentation via Model-agnostic Learning and Target-specific Normalization." arXiv preprint 2020

# Related Work

## Training strategies for domain generalization



Anatomy	Brain			Prostate (whole)			Prostate (sep.)		Heart	
	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	SD	TD <sub>1</sub>
Test										
Baseline and Benchmark										
SD (Baseline)	0.853	0.588	0.107	0.840	0.586	0.609	0.722	0.544	0.823	0.670
TD <sub>n</sub> (Benchmark)	-	0.896	0.867	-	0.817	0.834	-	0.732	-	0.806
Relevant Methods										
MLDG [Li et al. (2018a)]	0.874	0.686	0.074	0.913	0.772	0.756	0.818	0.658	0.844	0.696
MASF [Dou et al. (2019a)]	0.870	0.693	0.073	0.913	0.751	0.781	0.817	0.640	0.838	0.698
MLDGTS [Zhang et al. (2020b)]	0.876	0.733	0.072	0.912	0.711	0.761	0.815	0.608	0.831	0.361
SD + DA [Zhang et al. (2020)] (Strong baseline)	0.876	0.753	0.083	0.911	0.769	0.786	0.815	0.656	0.834	0.744

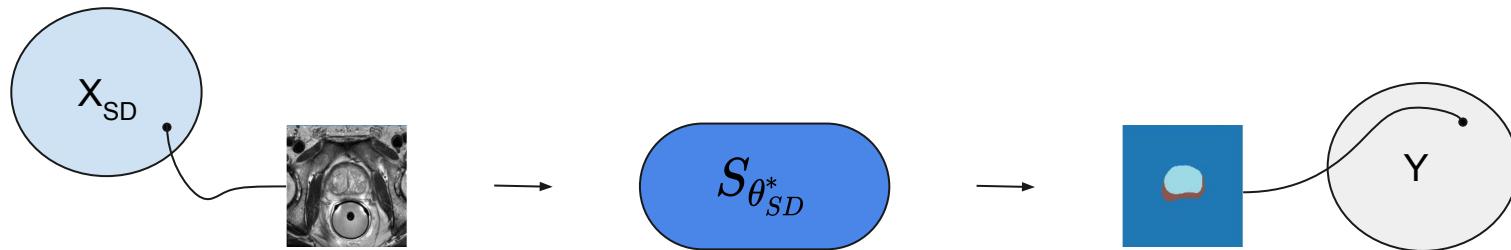
Training strategies improve DG performance.

Also, they improve performance within SD!

However, there still remains a gap to the benchmark.

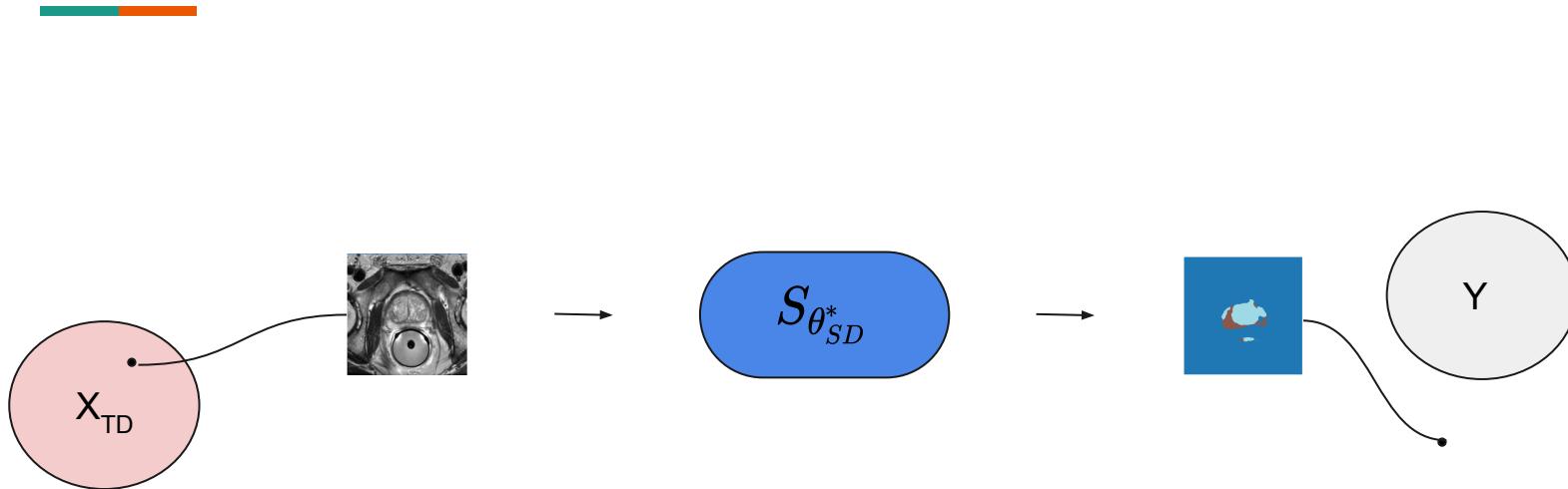
# Domain generalization

## The domain shift problem in machine learning



# Domain generalization

## The domain shift problem in machine learning

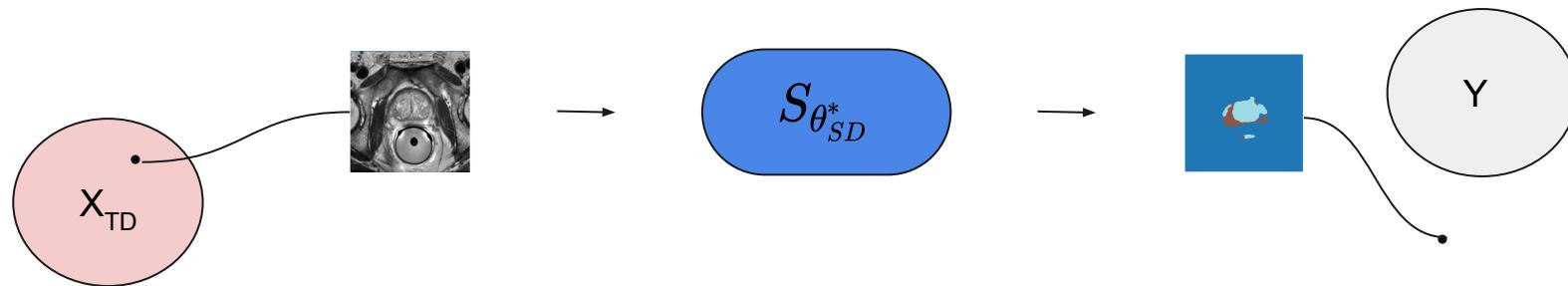


# Domain generalization

## The domain shift problem in machine learning

---

Hypothesis: For better generalization, we need per-test-image adaptability.

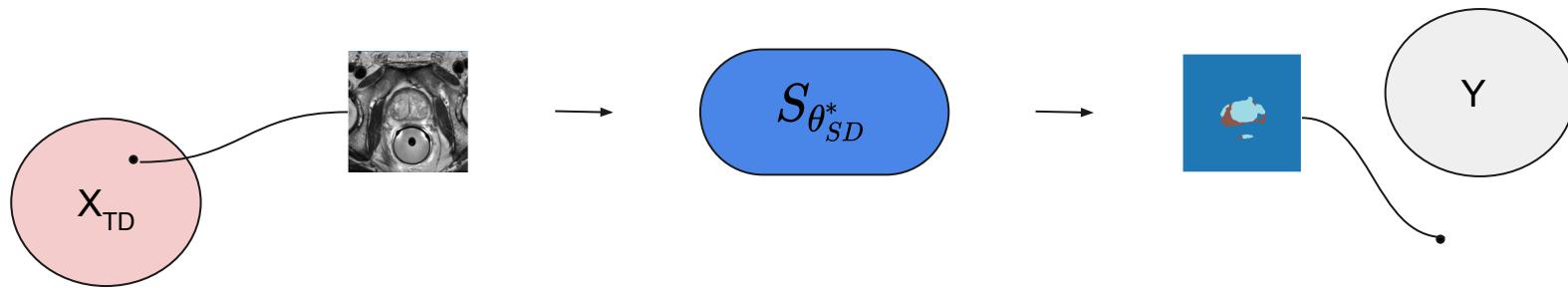


# Domain generalization

## The domain shift problem in machine learning



Hypothesis: For better generalization, we need per-test-image adaptability.



Two questions:

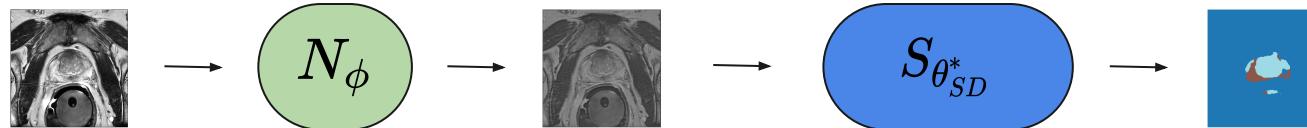
1. Which parameters to adapt?
2. How to drive the adaptation?

# Test-Time Adaptation

## Inference strategy for domain generalization



1. Which parameters to adapt?

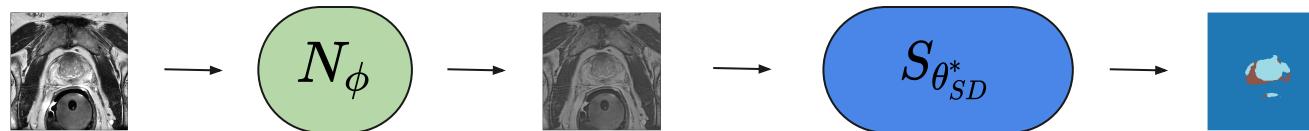


# Test-Time Adaptation

## Inference strategy for domain generalization



1. Which parameters to adapt?



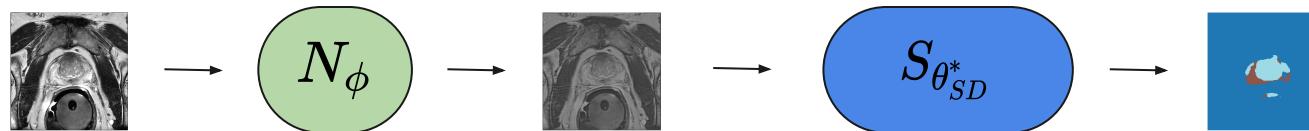
We train  $N$  and  $S$  jointly, on the SD.

Then, we fix  $S$  and adapt  $N$  for each test image.

# Test-Time Adaptation

## Inference strategy for domain generalization

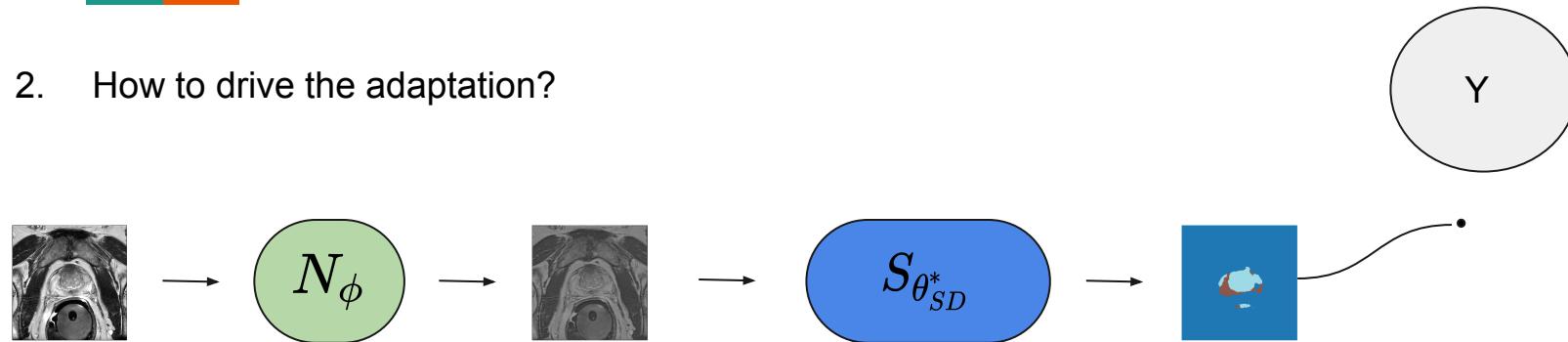
- 
- 2. How to drive the adaptation?



# Test-Time Adaptation

## Inference strategy for domain generalization

- 
- 2. How to drive the adaptation?

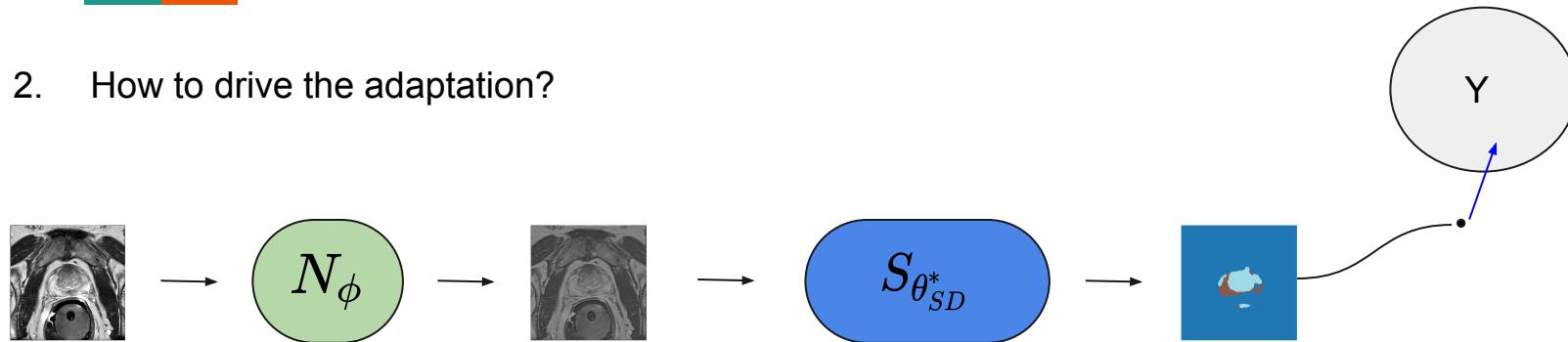


The predicted segmentation is likely to be outside the distribution of plausible segmentations.

# Test-Time Adaptation

## Inference strategy for domain generalization

- 
- 2. How to drive the adaptation?



The predicted segmentation is likely to be outside the distribution of plausible segmentations.

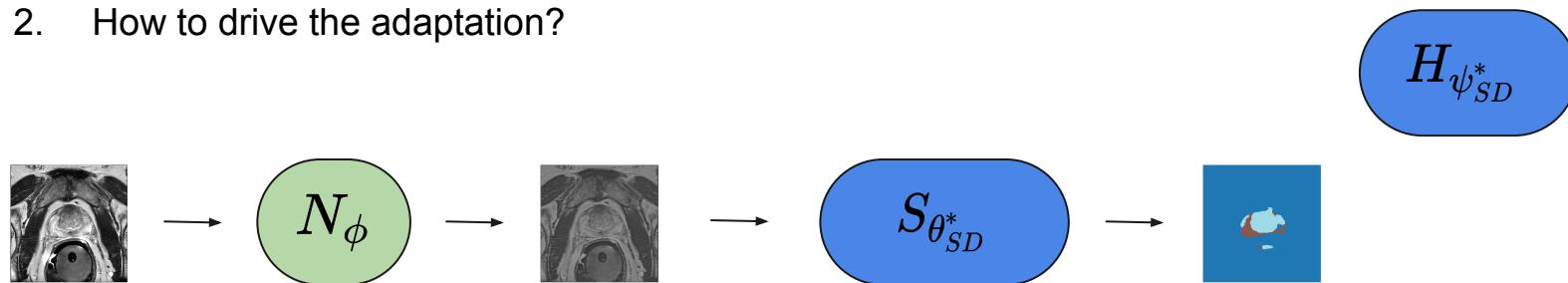
We drive the adaptation by encouraging the predicted segmentation to become plausible.

# Test-Time Adaptation

## Inference strategy for domain generalization

*Helper Network*

- 
- 2. How to drive the adaptation?

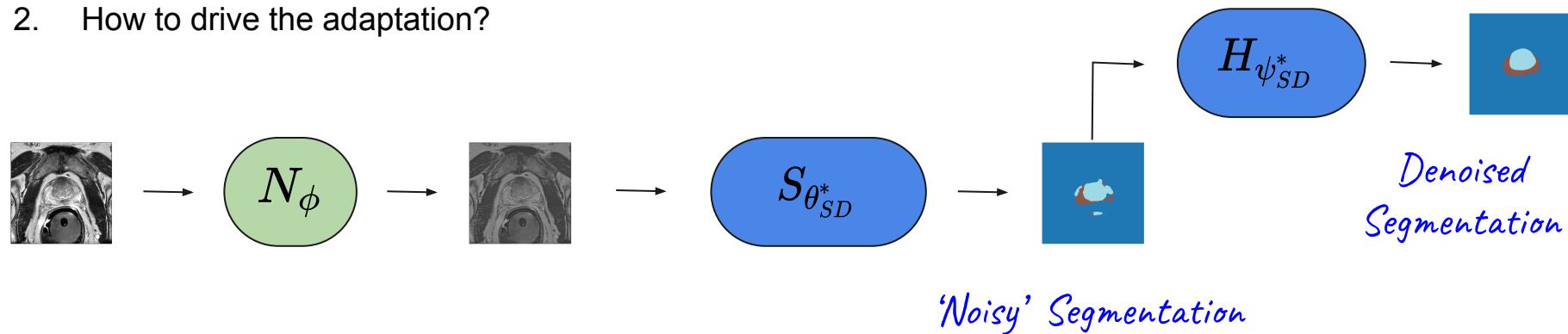


For this, we introduce a ‘helper’ network.

# Test-Time Adaptation

## Inference strategy for domain generalization

- How to drive the adaptation?



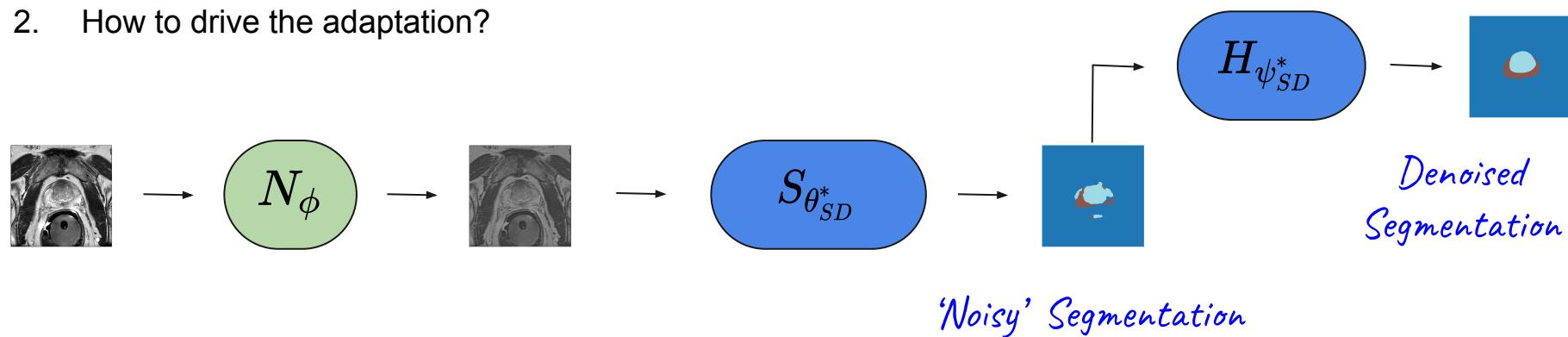
The segmentations predicted for TD images are treated as ‘noisy’ segmentations.

A DAE is used to denoise these predictions.

# Test-Time Adaptation

## Inference strategy for domain generalization

- How to drive the adaptation?



The segmentations predicted for TD images are treated as ‘noisy’ segmentations.

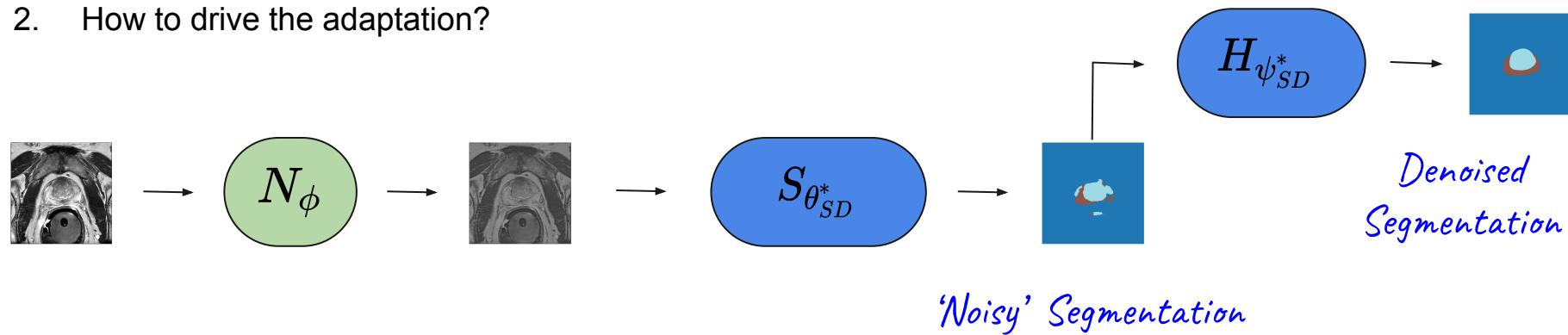
A DAE is used to denoise these predictions.

We adapt  $N$  for each test image, such that the DAE inputs and outputs are similar.

# Test-Time Adaptation

## Inference strategy for domain generalization

- How to drive the adaptation?



The segmentations predicted for TD images are treated as 'noisy' segmentations.

A DAE is used to denoise these predictions.

We adapt N for each test image, such that the DAE inputs and outputs are similar.

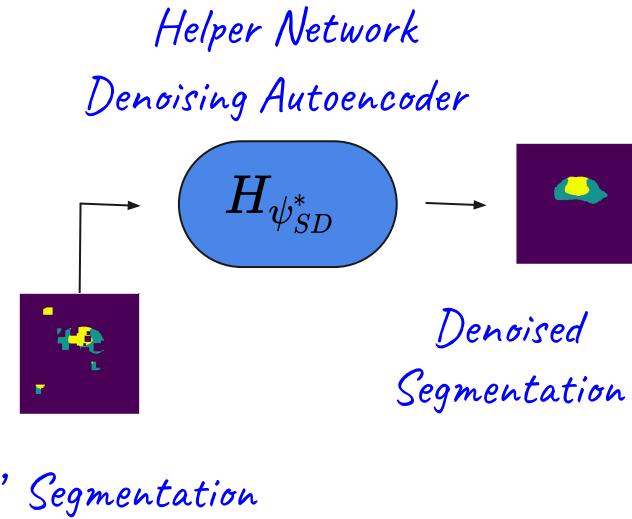
$$\begin{aligned} & \min_{\phi} L(y_n, H_{\psi_{SD}^*}(y_n)) \\ & y_n = S_{\theta_{SD}^*}(N_\phi(x_{TI})) \end{aligned}$$

# Test-Time Adaptation

## Inference strategy for domain generalization

For training the DAE, the noisy segmentations are generated by a heuristic noising process.

\* the DAE is trained in 3D



# Test-Time Adaptation

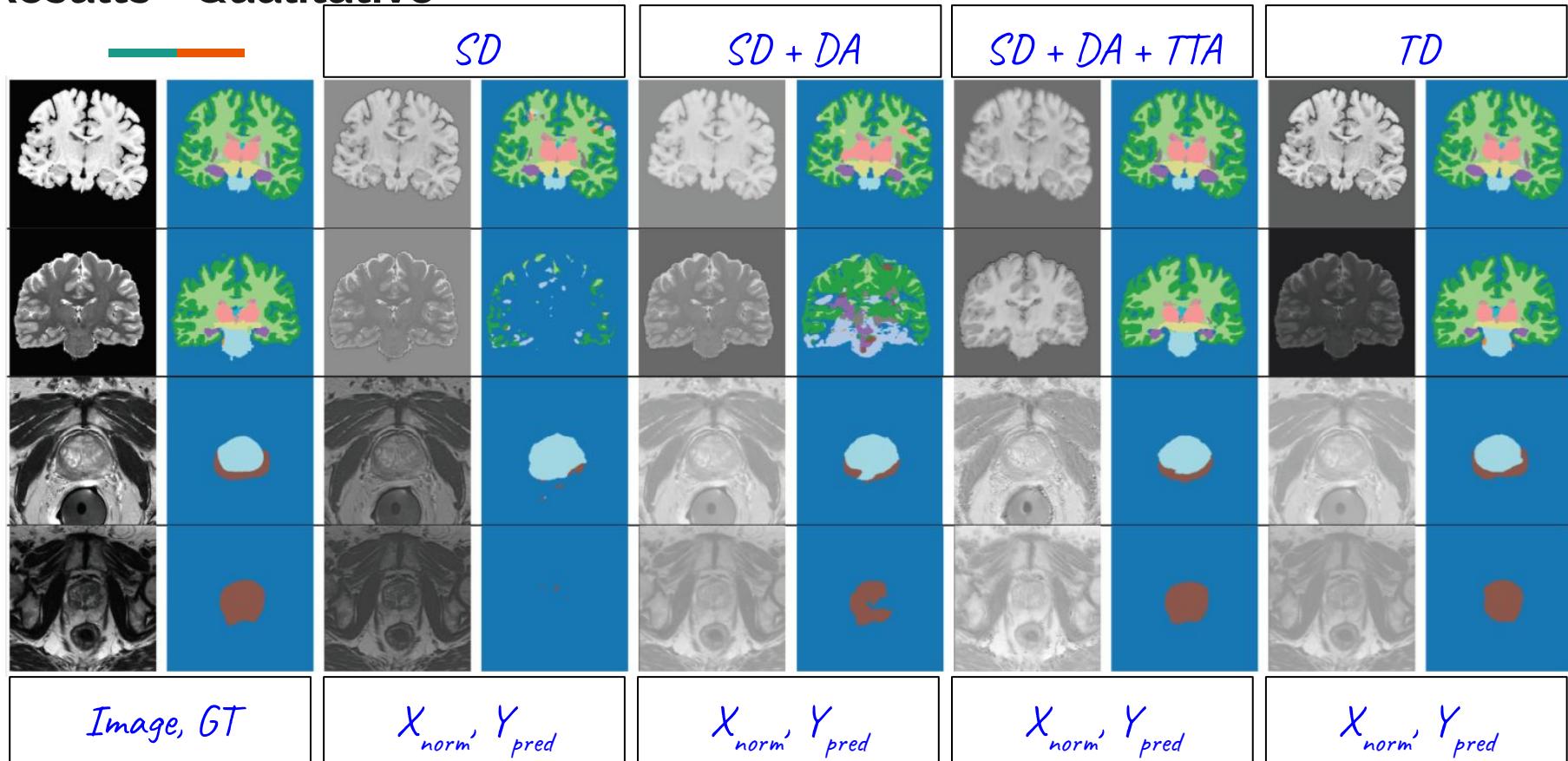
## Results



Anatomy	Brain			Prostate (whole)			Prostate (sep.)		Heart	
	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	SD	TD <sub>1</sub>
Test \ Train										
Baseline and Benchmark										
SD (Baseline)	0.853	0.588	0.107	0.840	0.586	0.609	0.722	0.544	0.823	0.670
TD <sub>n</sub> (Benchmark)	-	0.896	0.867	-	0.817	0.834	-	0.732	-	0.806
Relevant Methods										
MLDG [Li et al. (2018a)]	0.874	0.686	0.074	0.913	0.772	0.756	0.818	0.658	0.844	0.696
MASF [Dou et al. (2019a)]	0.870	0.693	0.073	0.913	0.751	0.781	0.817	0.640	0.838	0.698
MLDGTS [Zhang et al. (2020b)]	0.876	0.733	0.072	0.912	0.711	0.761	0.815	0.608	0.831	0.361
SD + DA [Zhang et al. (2020)] (Strong baseline)	0.876	0.753	0.083	0.911	0.769	0.786	0.815	0.656	0.834	0.744
SD + DA + Post-Proc. [Larrazaabal et al. (2019)]	-	0.706	0.112	-	0.789	0.823	-	0.678	-	0.746
Proposed Method										
SD + DA + TTA (Adapt $\phi$ , using DAE)	-	<b>0.800*</b>	<b>0.733*</b>	-	0.790	<b>0.858*</b>	-	0.676	-	0.742

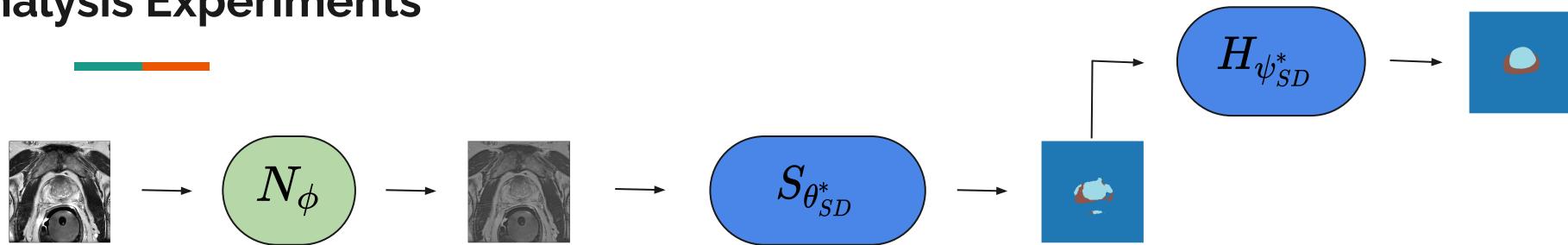
# Test-Time Adaptation

## Results - Qualitative



# Test-Time Adaptation

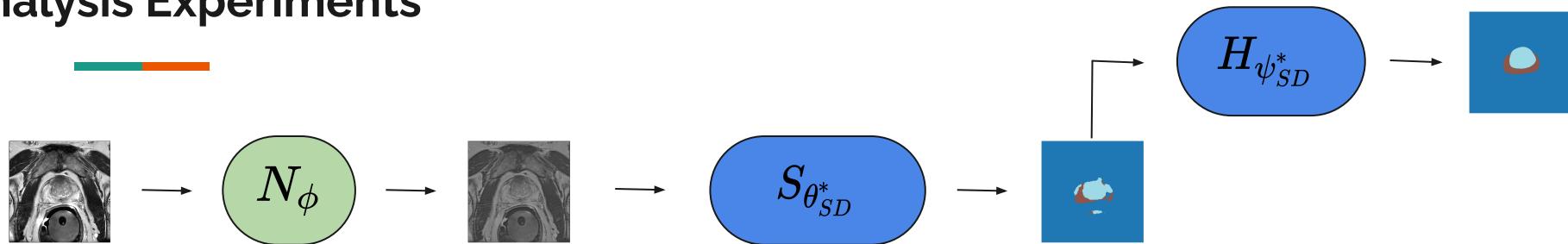
## Analysis Experiments



Anatomy	Brain			Prostate (whole)			Prostate (sep.)		Heart	
Test Train	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	SD	TD <sub>1</sub>
<b>Baseline and Benchmark</b>										
SD (Baseline)	0.853	0.588	0.107	0.840	0.586	0.609	0.722	0.544	0.823	0.670
TD <sub>n</sub> (Benchmark)	-	0.896	0.867	-	0.817	0.834	-	0.732	-	0.806
<b>Proposed Method</b>										
SD + DA + TTA (Adapt $\phi$ , using DAE)	-	<b>0.800*</b>	<b>0.733*</b>	-	0.790	<b>0.858*</b>	-	0.676	-	0.742
<b>Ablation Studies</b>										
SD + DA + TTA (Adapt $\phi$ , using DAE) - Fast	-	<b>0.800</b>	0.728	-	0.790	0.842	-	<b>0.683</b>	-	<b>0.747</b>
SD + DA + TTA (Adapt $\phi, \theta$ , using DAE)	-	0.671	0.650	-	0.718	0.606	-	0.583	-	0.713
SD + DA + TTA (Adapt $\phi$ , using GT labels)	-	0.831	0.837	-	0.836	-	-	0.771	-	-
SD + DA + Post.Proc. 10 passes through DAE	-	0.633	0.106	-	<b>0.791</b>	0.826	-	<b>0.683</b>	-	0.731
SD + DA + Post.Proc. 100 passes through DAE	-	0.529	0.101	-	0.789	0.823	-	0.672	-	0.688

# Test-Time Adaptation

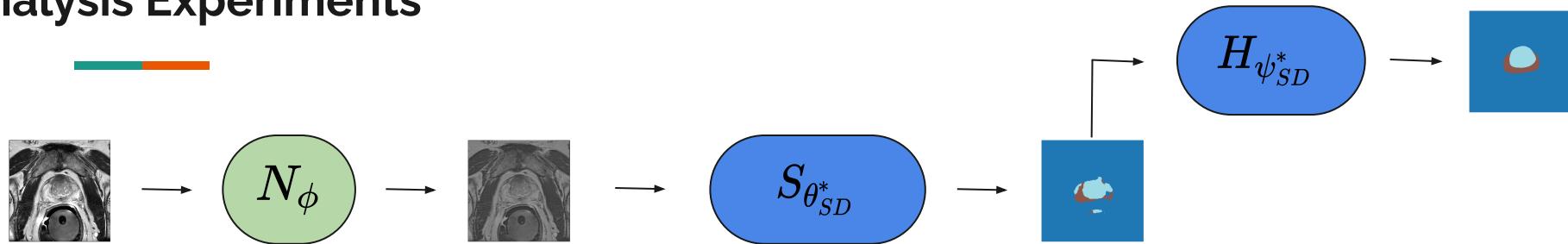
## Analysis Experiments



Anatomy	Brain			Prostate (whole)			Prostate (sep.)		Heart	
Train \ Test	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	SD	TD <sub>1</sub>
<b>Baseline and Benchmark</b>										
SD (Baseline)	0.853	0.588	0.107	0.840	0.586	0.609	0.722	0.544	0.823	0.670
TD <sub>n</sub> (Benchmark)	-	0.896	0.867	-	0.817	0.834	-	0.732	-	0.806
<b>Proposed Method</b>										
SD + DA + TTA (Adapt $\phi$ , using DAE)	-	<b>0.800*</b>	<b>0.733*</b>	-	0.790	<b>0.858*</b>	-	0.676	-	0.742
<b>Ablation Studies</b>										
SD + DA + TTA (Adapt $\phi$ , using DAE) - Fast	-	<b>0.800</b>	0.728	-	0.790	0.842	-	<b>0.683</b>	-	<b>0.747</b>
SD + DA + TTA (Adapt $\phi, \theta$ , using DAE)	-	0.671	0.650	-	0.718	0.606	-	0.583	-	0.713
SD + DA + TTA (Adapt $\phi$ , using GT labels)	-	0.831	0.837	-	0.836	-	-	0.771	-	-
SD + DA + Post.Proc. 10 passes through DAE	-	0.633	0.106	-	<b>0.791</b>	0.826	-	<b>0.683</b>	-	0.731
SD + DA + Post.Proc. 100 passes through DAE	-	0.529	0.101	-	0.789	0.823	-	0.672	-	0.688

# Test-Time Adaptation

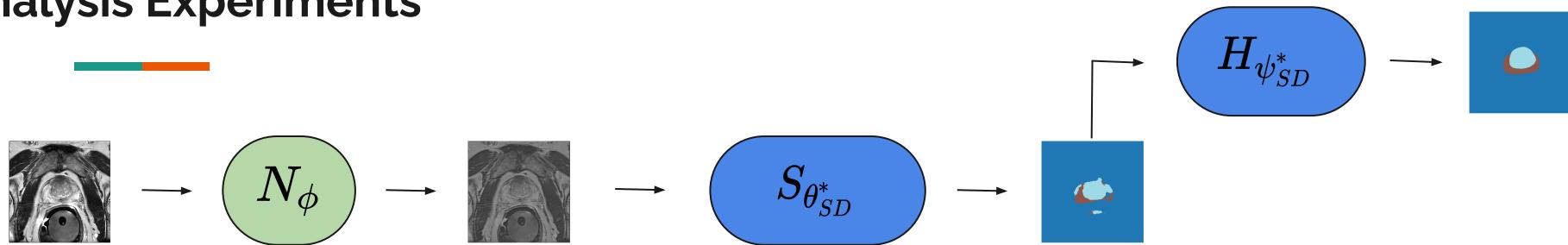
## Analysis Experiments



Anatomy	Brain			Prostate (whole)			Prostate (sep.)		Heart	
Test Train	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	SD	TD <sub>1</sub>
<b>Baseline and Benchmark</b>										
SD (Baseline)	0.853	0.588	0.107	0.840	0.586	0.609	0.722	0.544	0.823	0.670
TD <sub>n</sub> (Benchmark)	-	0.896	0.867	-	0.817	0.834	-	0.732	-	0.806
<b>Proposed Method</b>										
SD + DA + TTA (Adapt $\phi$ , using DAE)	-	<b>0.800*</b>	<b>0.733*</b>	-	0.790	<b>0.858*</b>	-	0.676	-	0.742
<b>Ablation Studies</b>										
SD + DA + TTA (Adapt $\phi$ , using DAE) - Fast	-	<b>0.800</b>	0.728	-	0.790	0.842	-	<b>0.683</b>	-	<b>0.747</b>
SD + DA + TTA (Adapt $\phi, \theta$ , using DAE)	-	0.671	0.650	-	0.718	0.606	-	0.583	-	0.713
SD + DA + TTA (Adapt $\phi$ , using GT labels)	-	0.831	0.837	-	0.836	-	-	0.771	-	-
SD + DA + Post.Proc. 10 passes through DAE	-	0.633	0.106	-	<b>0.791</b>	0.826	-	<b>0.683</b>	-	0.731
SD + DA + Post.Proc. 100 passes through DAE	-	0.529	0.101	-	0.789	0.823	-	0.672	-	0.688

# Test-Time Adaptation

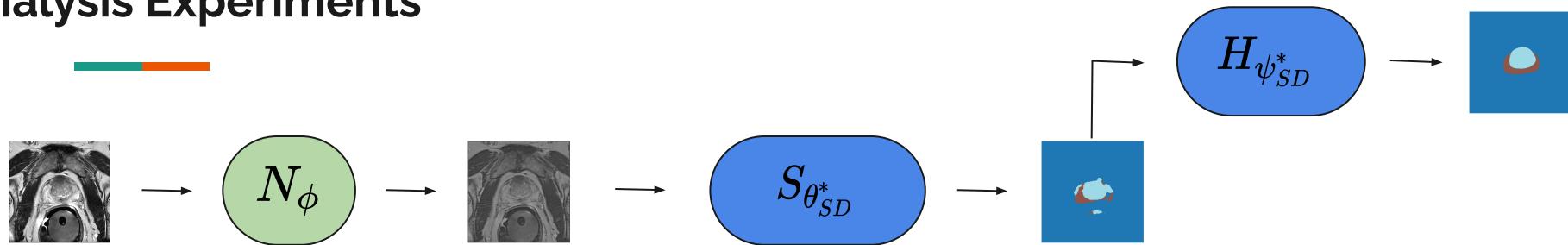
## Analysis Experiments



Anatomy	Brain			Prostate (whole)			Prostate (sep.)		Heart	
Test Train	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	SD	TD <sub>1</sub>
<b>Baseline and Benchmark</b>										
SD (Baseline)	0.853	0.588	0.107	0.840	0.586	0.609	0.722	0.544	0.823	0.670
TD <sub>n</sub> (Benchmark)	-	0.896	0.867	-	0.817	0.834	-	0.732	-	0.806
<b>Proposed Method</b>										
SD + DA + TTA (Adapt $\phi$ , using DAE)	-	<b>0.800*</b>	<b>0.733*</b>	-	0.790	<b>0.858*</b>	-	0.676	-	0.742
<b>Ablation Studies</b>										
SD + DA + TTA (Adapt $\phi$ , using DAE) - Fast	-	<b>0.800</b>	0.728	-	0.790	0.842	-	<b>0.683</b>	-	<b>0.747</b>
SD + DA + TTA (Adapt $\phi, \theta$ , using DAE)	-	0.671	0.650	-	0.718	0.606	-	0.583	-	0.713
SD + DA + TTA (Adapt $\phi$ , using GT labels)	-	0.831	0.837	-	0.836	-	-	0.771	-	-
SD + DA + Post.Proc. 10 passes through DAE	-	0.633	0.106	-	<b>0.791</b>	0.826	-	<b>0.683</b>	-	0.731
SD + DA + Post.Proc. 100 passes through DAE	-	0.529	0.101	-	0.789	0.823	-	0.672	-	0.688

# Test-Time Adaptation

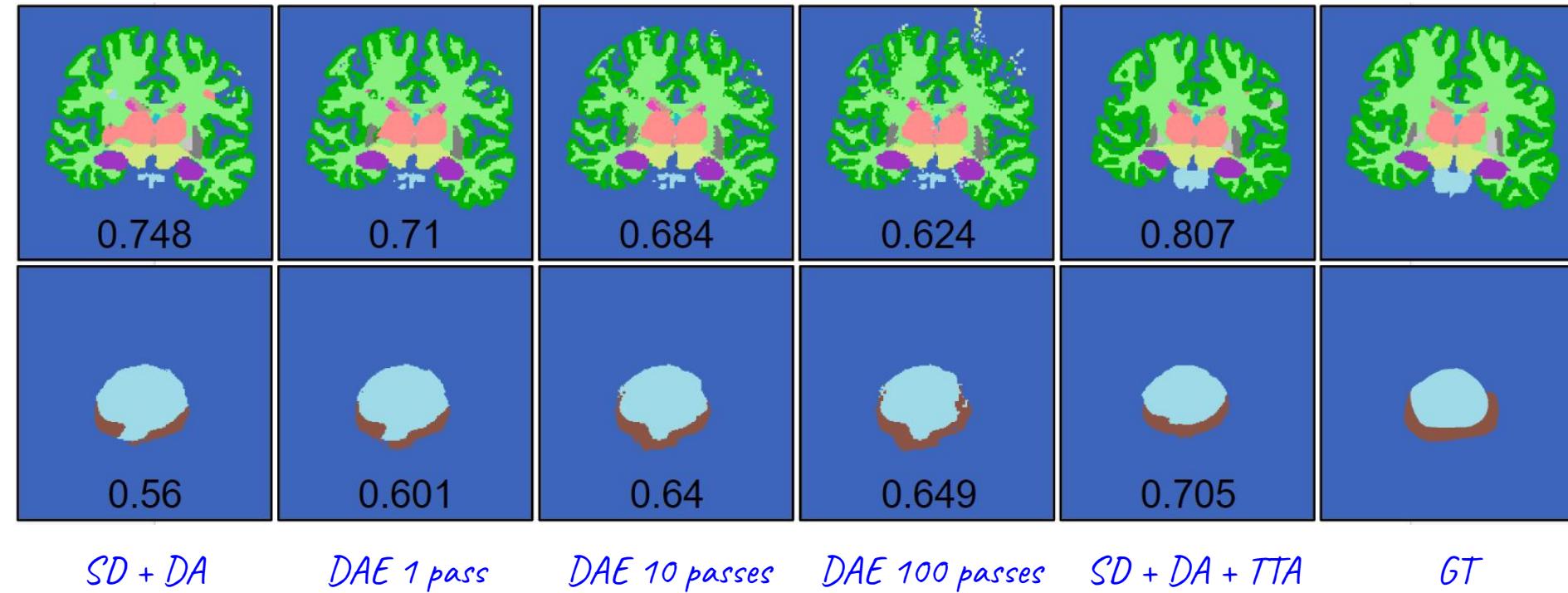
## Analysis Experiments



Anatomy	Brain			Prostate (whole)			Prostate (sep.)		Heart	
Test Train	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	TD <sub>2</sub>	SD	TD <sub>1</sub>	SD	TD <sub>1</sub>
<b>Baseline and Benchmark</b>										
SD (Baseline)	0.853	0.588	0.107	0.840	0.586	0.609	0.722	0.544	0.823	0.670
TD <sub>n</sub> (Benchmark)	-	0.896	0.867	-	0.817	0.834	-	0.732	-	0.806
<b>Proposed Method</b>										
SD + DA + TTA (Adapt $\phi$ , using DAE)	-	<b>0.800*</b>	<b>0.733*</b>	-	0.790	<b>0.858*</b>	-	0.676	-	0.742
<b>Ablation Studies</b>										
SD + DA + TTA (Adapt $\phi$ , using DAE) - Fast	-	<b>0.800</b>	0.728	-	0.790	0.842	-	<b>0.683</b>	-	<b>0.747</b>
SD + DA + TTA (Adapt $\phi, \theta$ , using DAE)	-	0.671	0.650	-	0.718	0.606	-	0.583	-	0.713
SD + DA + TTA (Adapt $\phi$ , using GT labels)	-	0.831	0.837	-	0.836	-	-	0.771	-	-
SD + DA + Post.Proc. 10 passes through DAE	-	0.633	0.106	-	<b>0.791</b>	0.826	-	<b>0.683</b>	-	0.731
SD + DA + Post.Proc. 100 passes through DAE	-	0.529	0.101	-	0.789	0.823	-	0.672	-	0.688

# Test-Time Adaptation

## Importance of adaptation v/s using DAE for post-processing



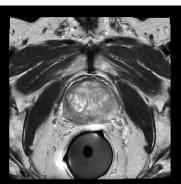
# Test-Time Adaptation

## Evolution of segmentation over adaptation iterations

Prostate USZ 45



Image



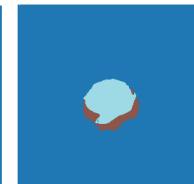
Prediction  
Step 0



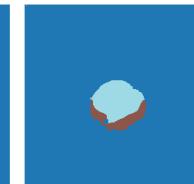
Prediction  
Step 1k



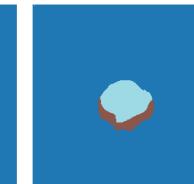
Prediction  
Step 3k



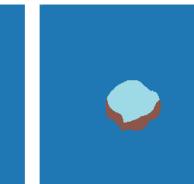
Prediction  
Step 6k



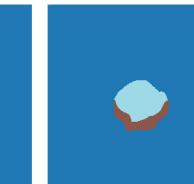
Prediction  
Step 9k



Prediction  
Step 12k



Prediction  
Step 15k



Ground truth

DAE output  
Step 0

DAE output  
Step 1k

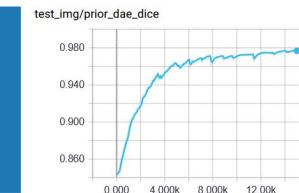
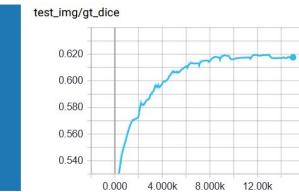
DAE output  
Step 3k

DAE output  
Step 6k

DAE output  
Step 9k

DAE output  
Step 12k

DAE output  
Step 15k



# Test-Time Adaptation

## Evolution of segmentation over adaptation iterations

Brain Caltech 288



Image

Prediction  
Step 0

Prediction  
Step 800

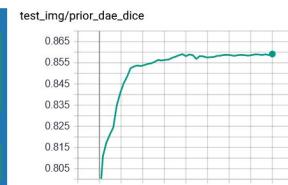
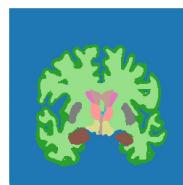
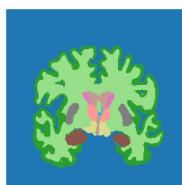
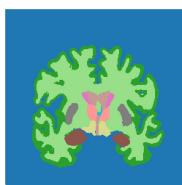
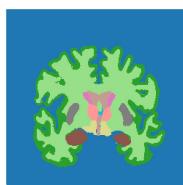
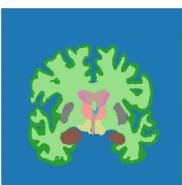
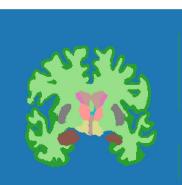
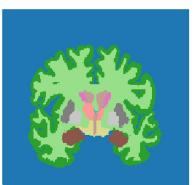
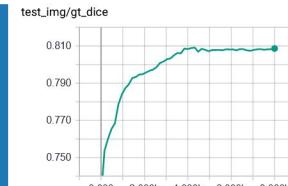
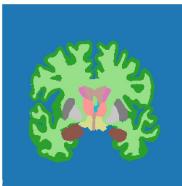
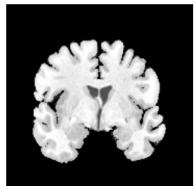
Prediction  
Step 1600

Prediction  
Step 3200

Prediction  
Step 4800

Prediction  
Step 6400

Prediction  
Step 8000



Ground truth

DAE output  
Step 0

DAE output  
Step 800

DAE output  
Step 1600

DAE output  
Step 3200

DAE output  
Step 4800

DAE output  
Step 6400

DAE output  
Step 8000

# Test-Time Adaptation

## Comparison with Unsupervised Domain Adaptation



Anatomy	Test	Brain				Prostate (whole)				Heart	
		TD <sub>1</sub> <sup>tr</sup>	TD <sub>1</sub> <sup>ts</sup>	TD <sub>2</sub> <sup>tr</sup>	TD <sub>2</sub> <sup>ts</sup>	TD <sub>1</sub> <sup>tr</sup>	TD <sub>1</sub> <sup>ts</sup>	TD <sub>2</sub> <sup>tr</sup>	TD <sub>2</sub> <sup>ts</sup>	TD <sub>1</sub> <sup>tr</sup>	TD <sub>1</sub> <sup>ts</sup>
Train											
SD + DA [Zhang et al. (2020)] (Strong baseline)		0.750	0.753	0.081	0.083	0.767	0.769	0.747	0.786	0.715	0.744
UDA - Invariant features [Kamnitsas et al. (2017a)]		0.792	0.798	0.081	0.083	0.789	0.793	0.766	0.802	0.724	0.750
UDA - Image-to-Image translation [Huo et al. (2018)]		0.646	0.639	0.816	0.813	0.607	0.694	0.765	0.747	0.252	0.167
TTA (Proposed)		-	0.800	-	0.733	-	0.790	-	0.858	-	0.742
TD <sub>n</sub> (Benchmark)		-	0.896	-	0.867	-	0.817	-	0.834	-	0.806

UDA methods work in a more relaxed scenario - labelled SD dataset is available for each new TD.

The proposed TTA method is a domain generalization method.

We find that TTA performance is similar to state-of-the-art UDA methods.

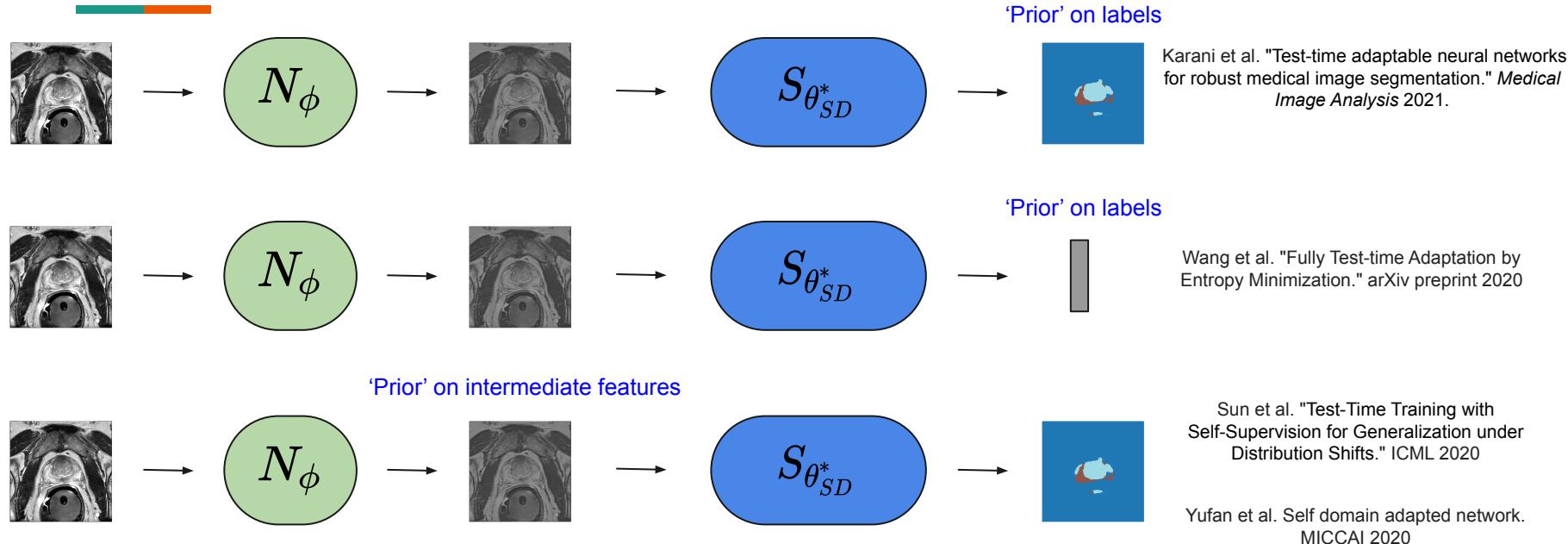
# Test-Time Adaptation

## Connection with concurrent TTA works



# Test-Time Adaptation

## Connection with concurrent TTA works



# Test-Time Adaptation for Robust Medical Image Segmentation

## Summary



- We proposed test-time adaptation for robust image segmentation.
- The method consists of two main ideas:
  - An adaptable per-image normalization module in front of the segmentation CNN.
  - The adaptation is driven by denoising autoencoders, that incentivize plausible segmentation predictions.
- Substantial improvement across datasets and anatomies.
- Increased inference time - about 10 minutes for each 3D volume.

# Test-Time Adaptation for Robust Medical Image Segmentation

## Summary



- We proposed test-time adaptation for robust image segmentation.
- The method consists of two main ideas:
  - An adaptable per-image normalization module in front of the segmentation CNN.
  - The adaptation is driven by denoising autoencoders, that incentivize plausible segmentation predictions.
- Substantial improvement across datasets and anatomies.
- Increased inference time - about 10 minutes for each 3D volume.

# Thanks!



**Krishna Chaitanya**



**Ertunc Erdil**



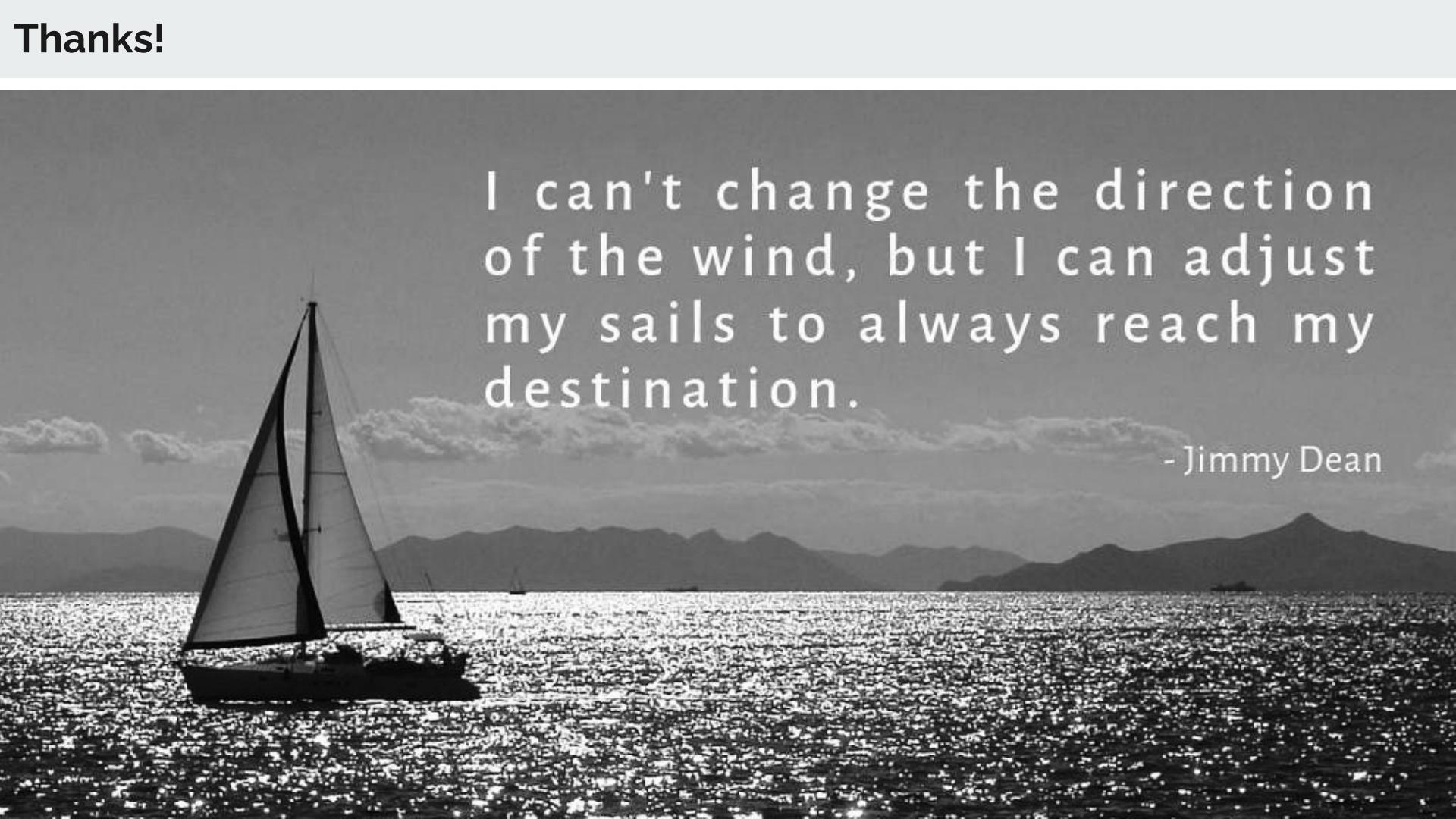
**Christian Baumgartner**



**Ender Konukoglu**

In case of questions, please feel free to write to [nkarani@vision.ee.ethz.ch](mailto:nkarani@vision.ee.ethz.ch)

Thanks!

A black and white photograph of a sailboat with two sails on a body of water. The boat is positioned on the left side of the frame, moving towards the right. In the background, there are several layers of mountains under a cloudy sky. The water in the foreground has small ripples and reflections.

I can't change the direction  
of the wind, but I can adjust  
my sails to always reach my  
destination.

- Jimmy Dean