



Montreal, July 8-12

DLMI2024

Semi-supervised Learning for Medical Image Segmentation

Christian Desrosiers

ÉTS Montreal, Canada



Importance of unlabeled data

Training deep neural nets requires **lots** of labeled data

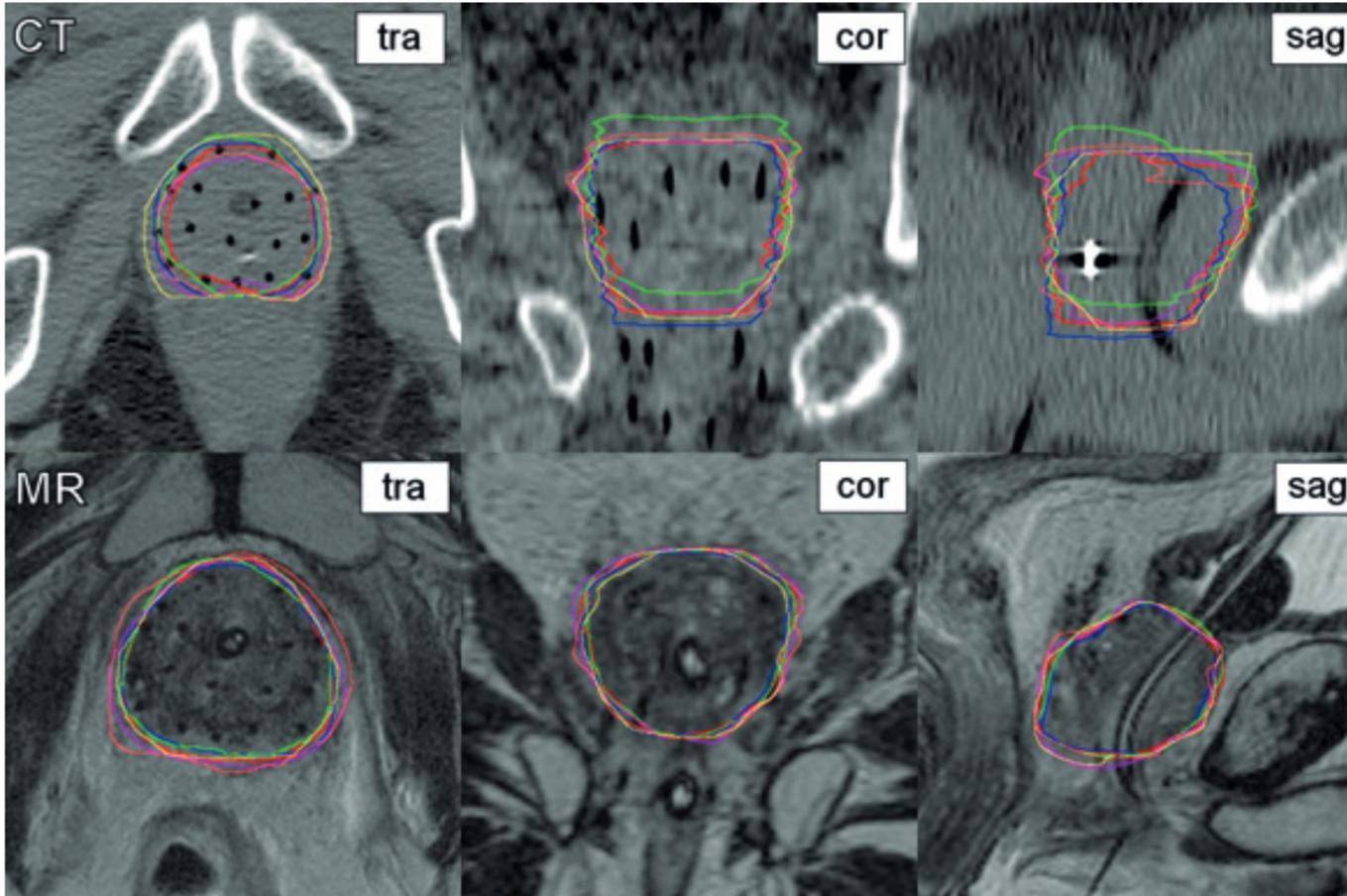


ImageNet

- Over 14M annotated images
- More than 20,000 classes

Importance of unlabeled data

However, annotating data can be hard and expensive in some applications...



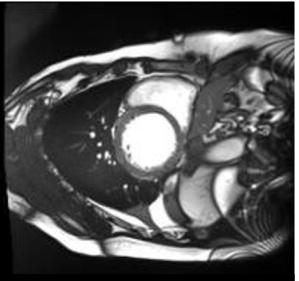
Challenges:

- Volumetric images
- Low contrast regions
- Can only be done by trained experts

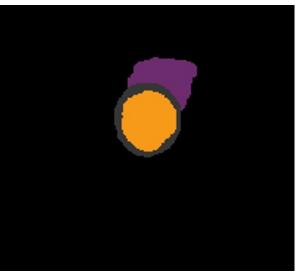
... but unlabeled data is often available for free

Learning with unlabeled images

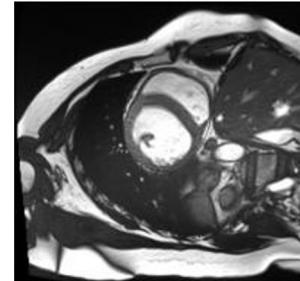
Labeled images (few)



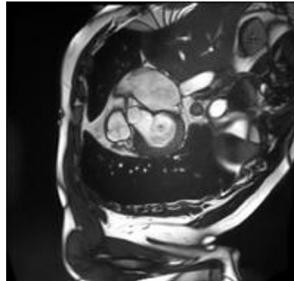
...



Unlabeled images (many)

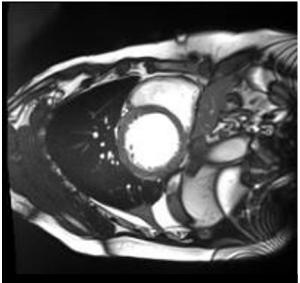


...

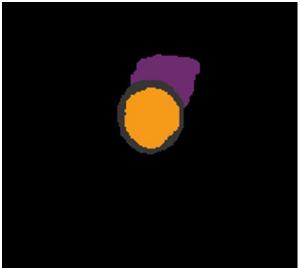
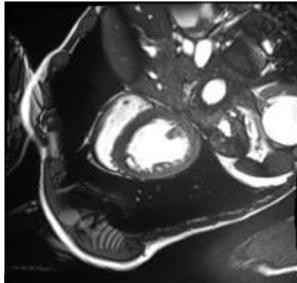


Learning with unlabeled images

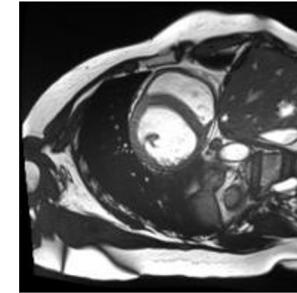
Labeled images (few)



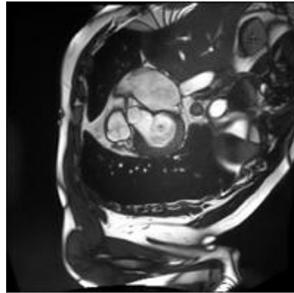
...



Unlabeled images (many)



...



How can we use this information
to learn segmentation ?

Outline

- 1) Adversarial learning
- 2) Consistency regularization
- 3) Unsupervised representation learning

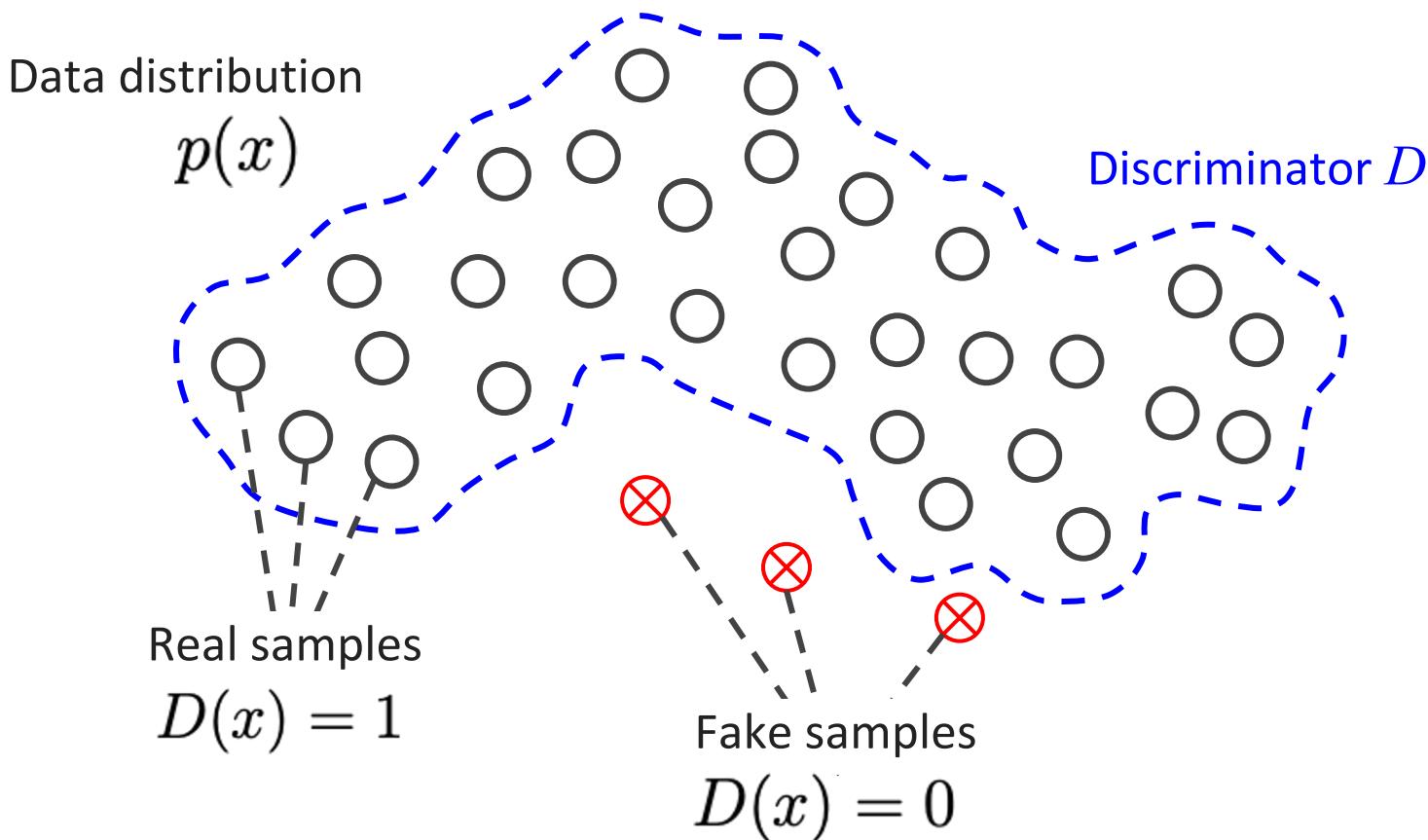
Adversarial learning

for semi-supervised segmentation

Adversarial learning

Basic idea:

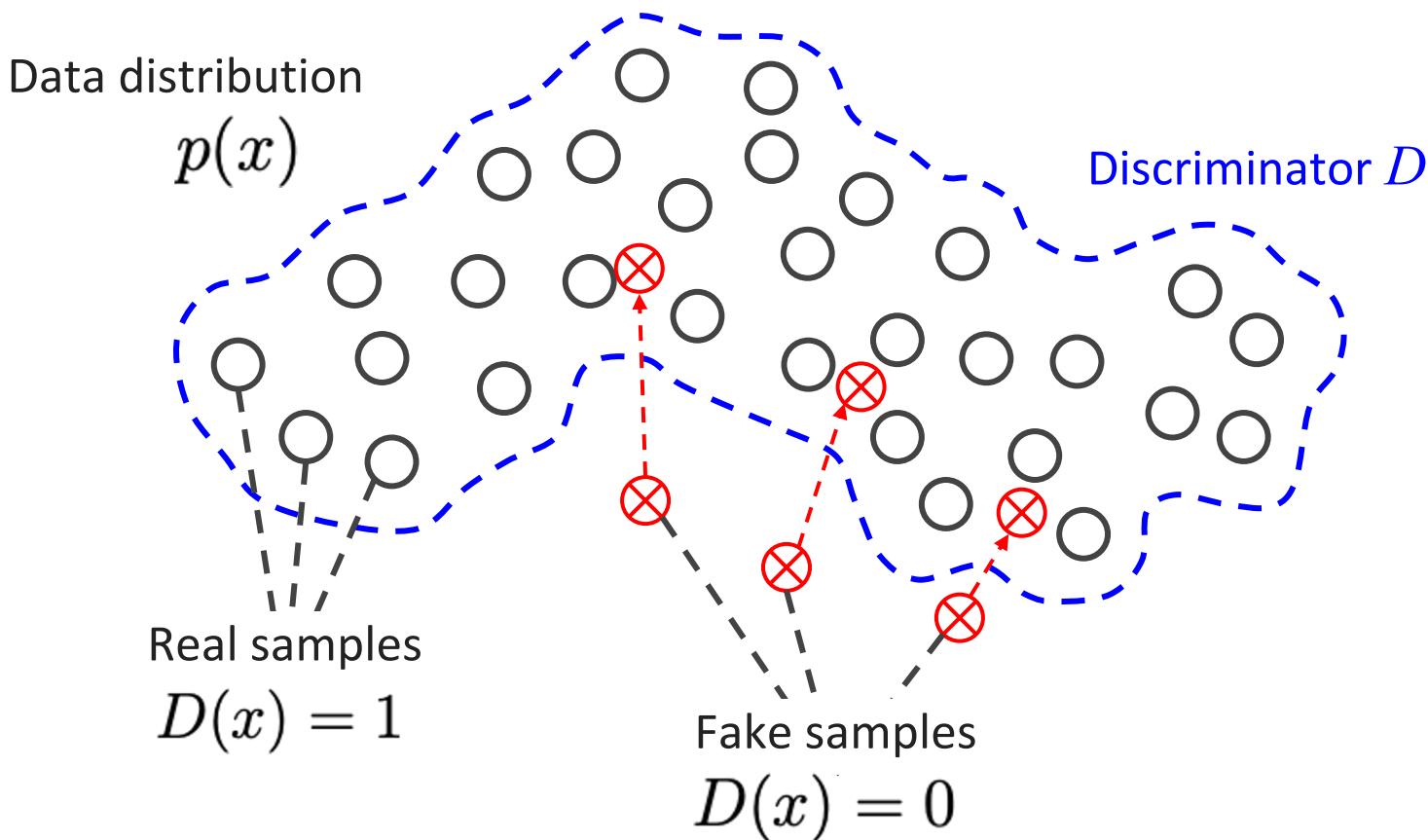
Learn the data distribution using a classifier (the discriminator)



Adversarial learning

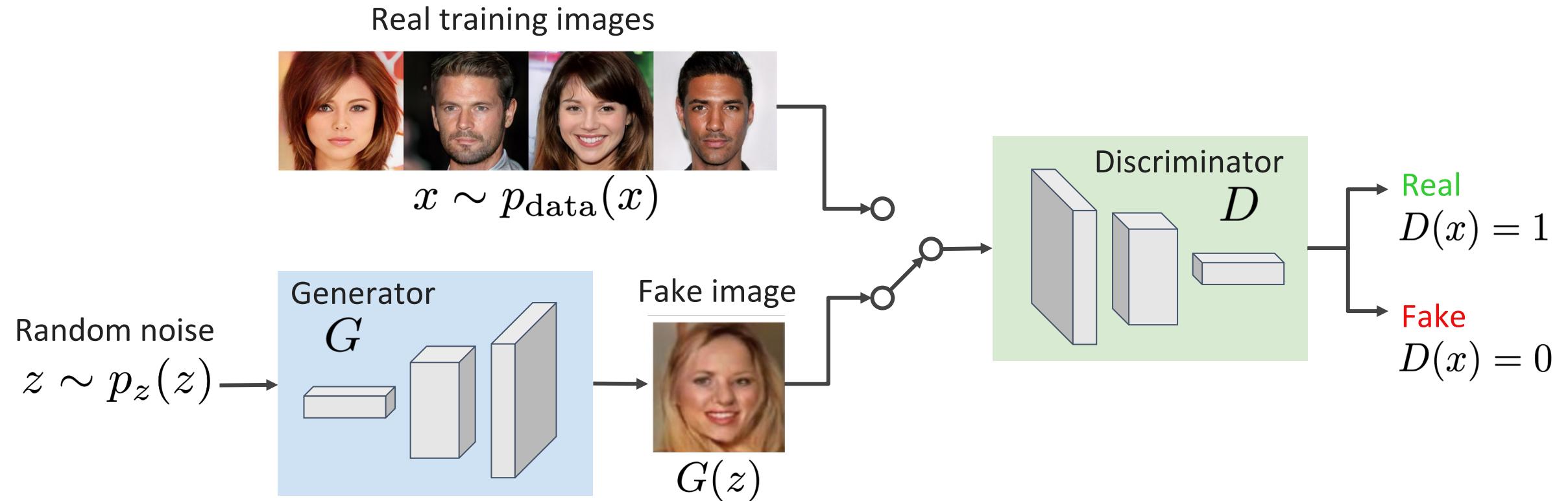
Basic idea:

Learn the data distribution using a classifier (the discriminator)

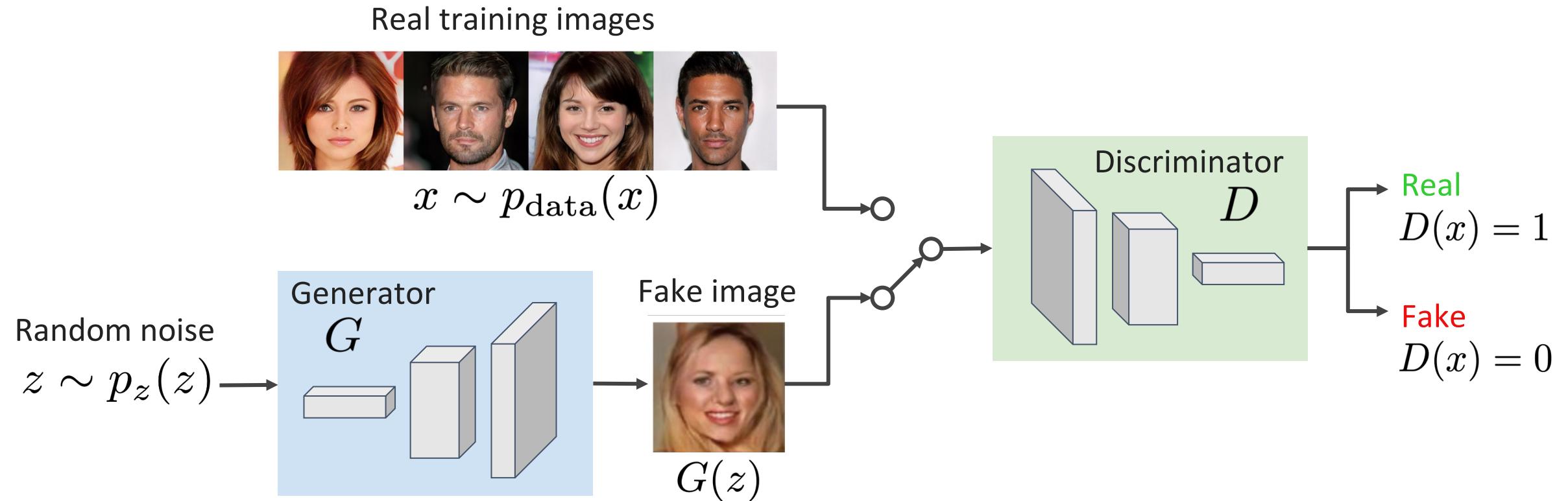


Objective: Generate samples in the distribution of real data

Generative adversarial network (GAN)

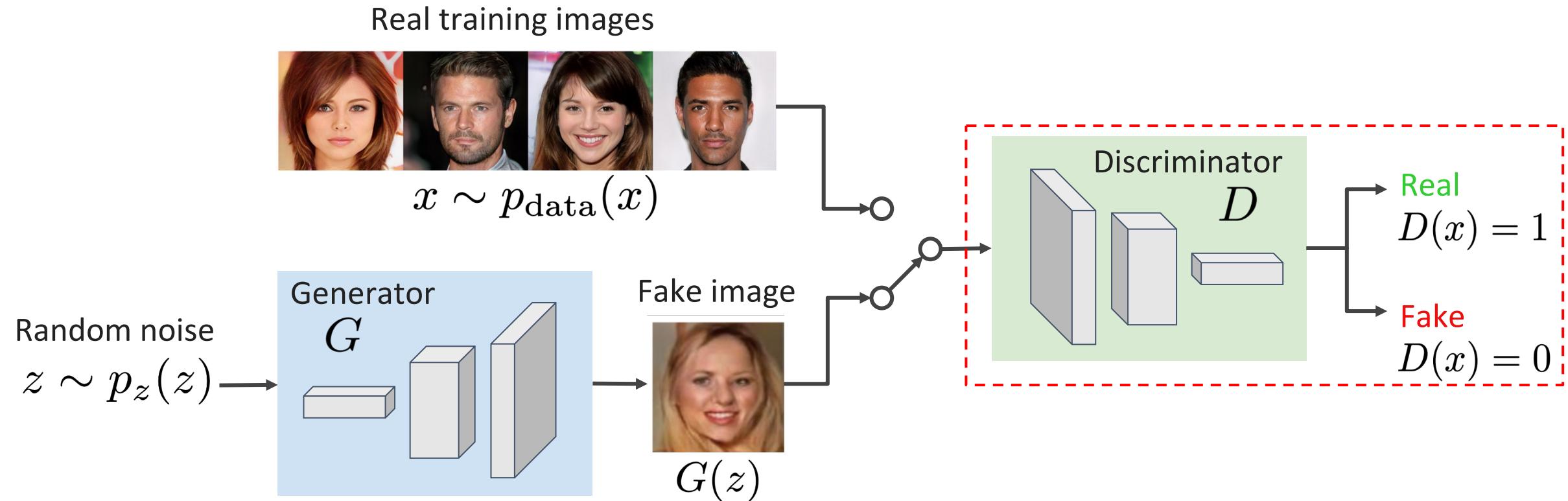


Generative adversarial network (GAN)



How to make sure that generated images look real ?

Generative adversarial network (GAN)



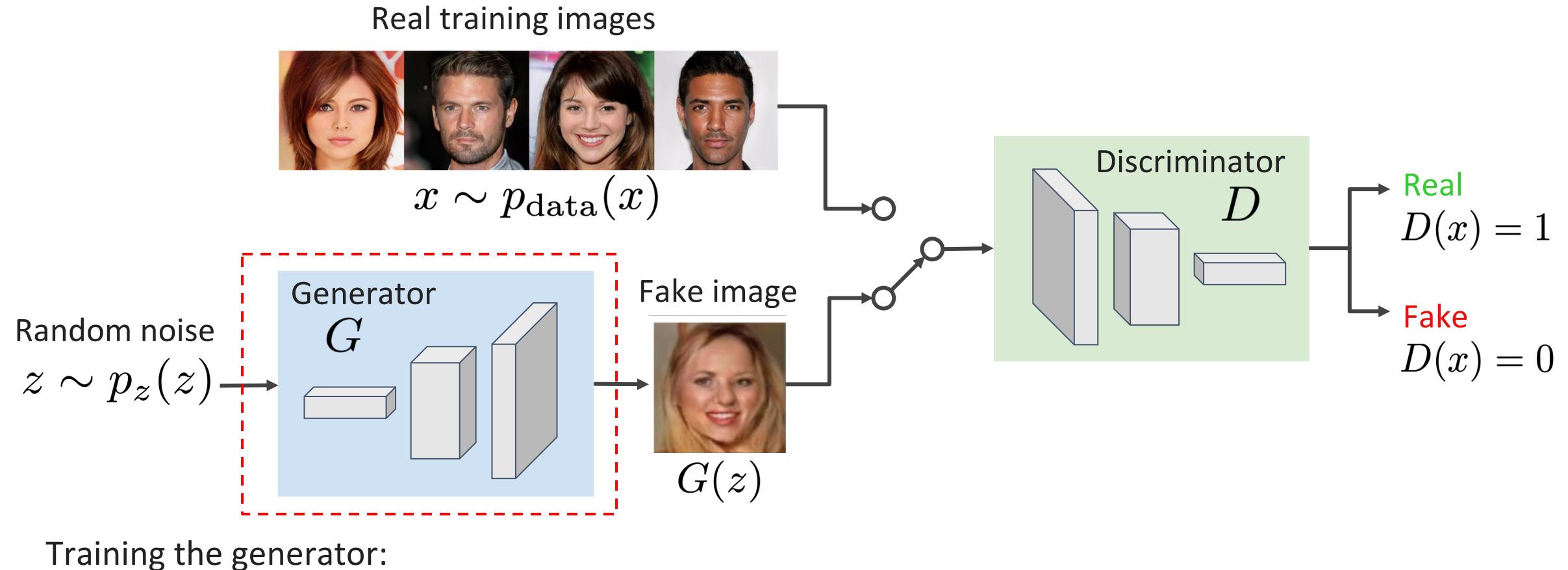
Training the discriminator (cross-entropy):

$$\max_D \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

Output '1' for real images

Output '0' for generated images

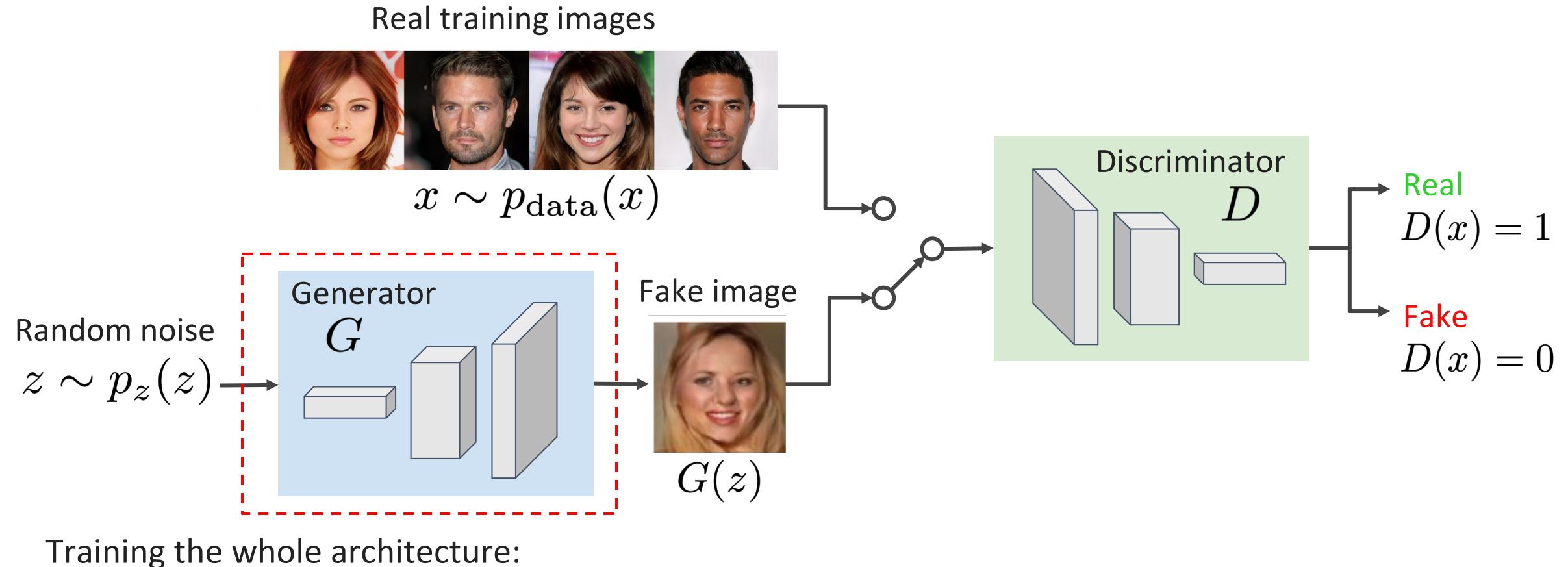
Generative adversarial network (GAN)



$$\min_G \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

Fool the discriminator into predicting '1' for fake images

Generative adversarial network (GAN)



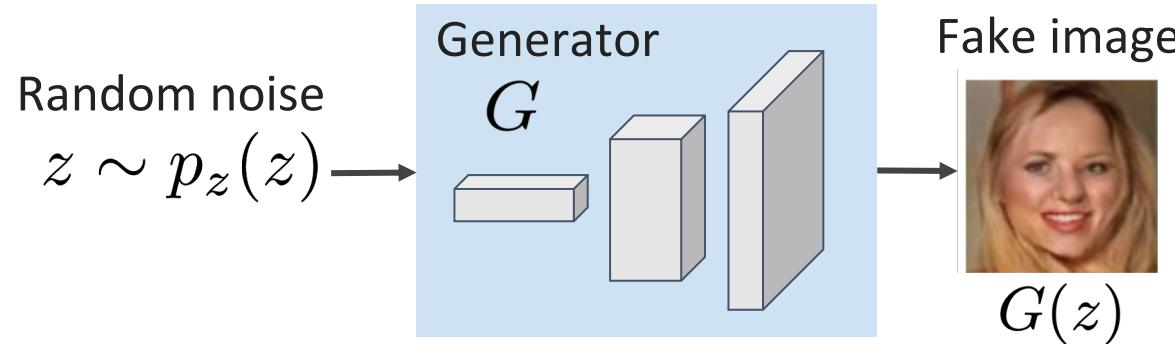
Training the whole architecture:

$$\min_G \max_D \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

Corresponds to a minimax problem (*more on this later...*)

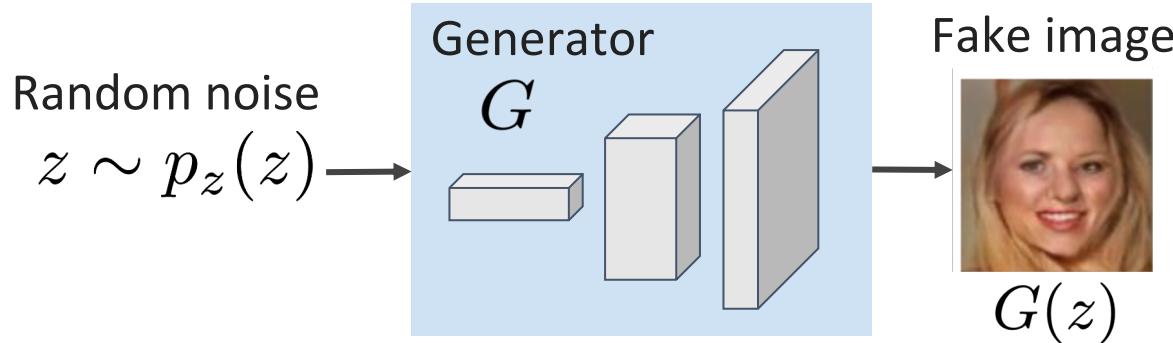
GANs for segmentation

GAN for image generation:

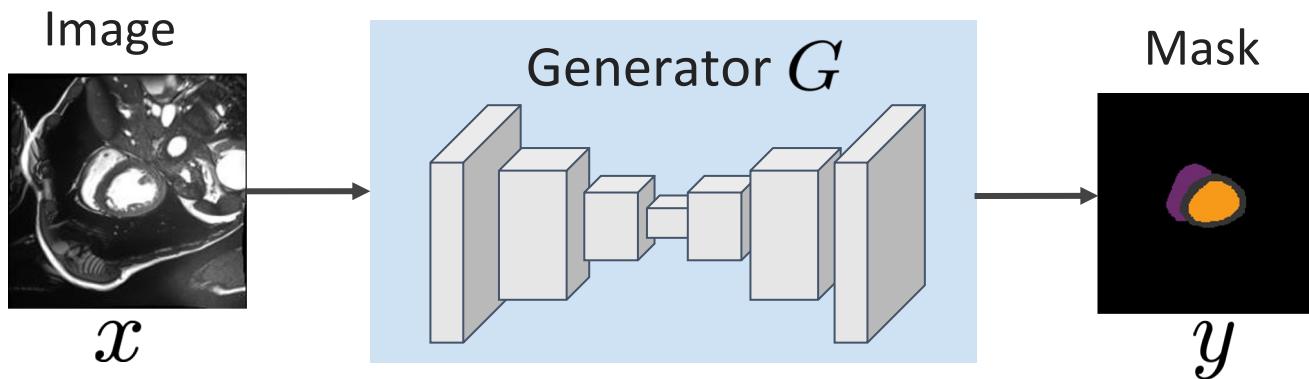


GANs for segmentation

GAN for image generation:

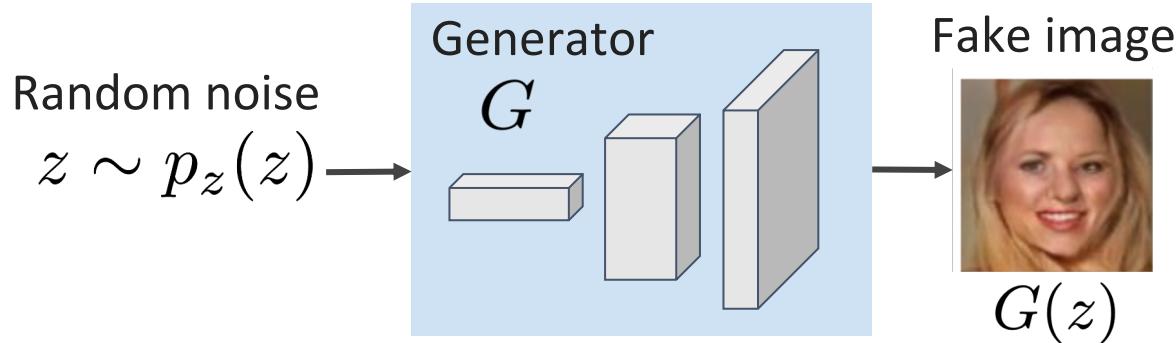


GAN for image segmentation:

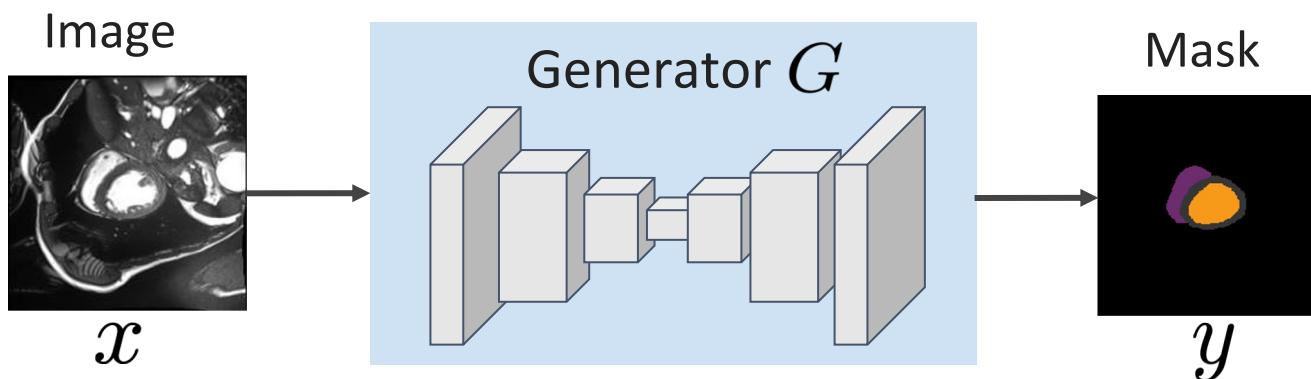


GANs for segmentation

GAN for image generation:



GAN for image segmentation:

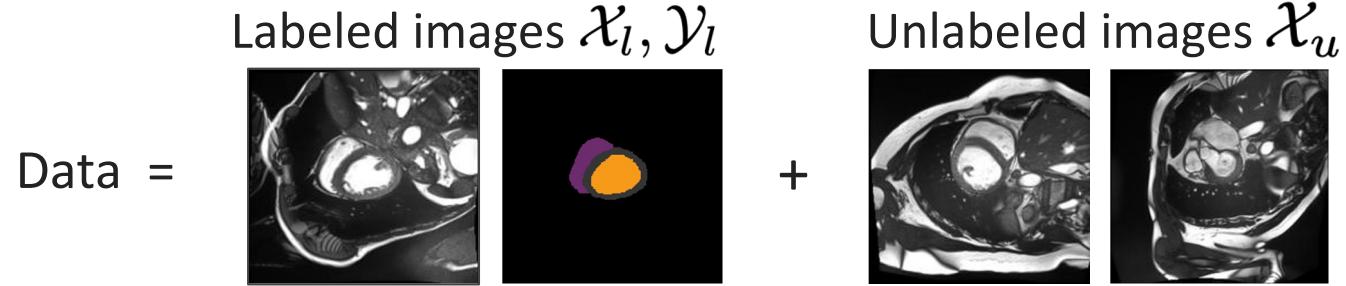


We are now modeling
the distribution of
segmentation masks

The generator is a segmentation network (encoder-decoder)

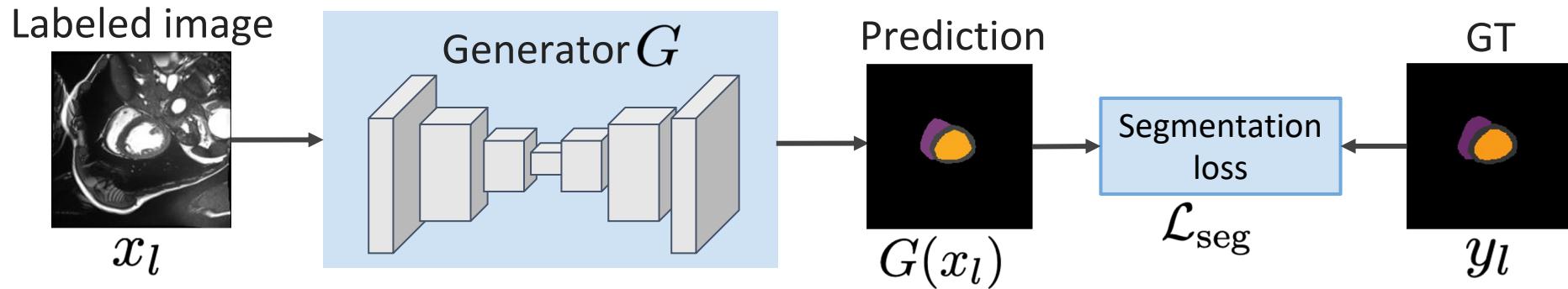
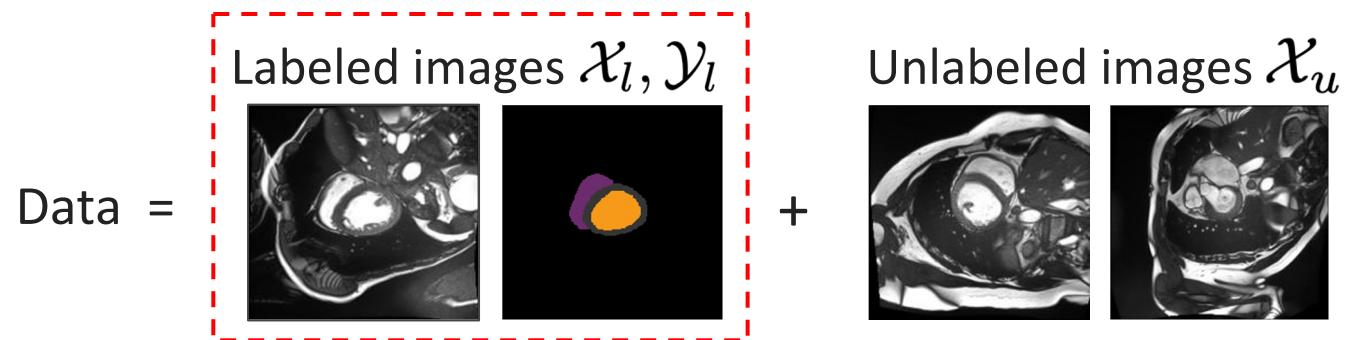
Adversarial semi-supervised segmentation

Basic idea: Learn to generate segmentation masks which can't be differentiated from ground-truth (GT)



Adversarial semi-supervised segmentation

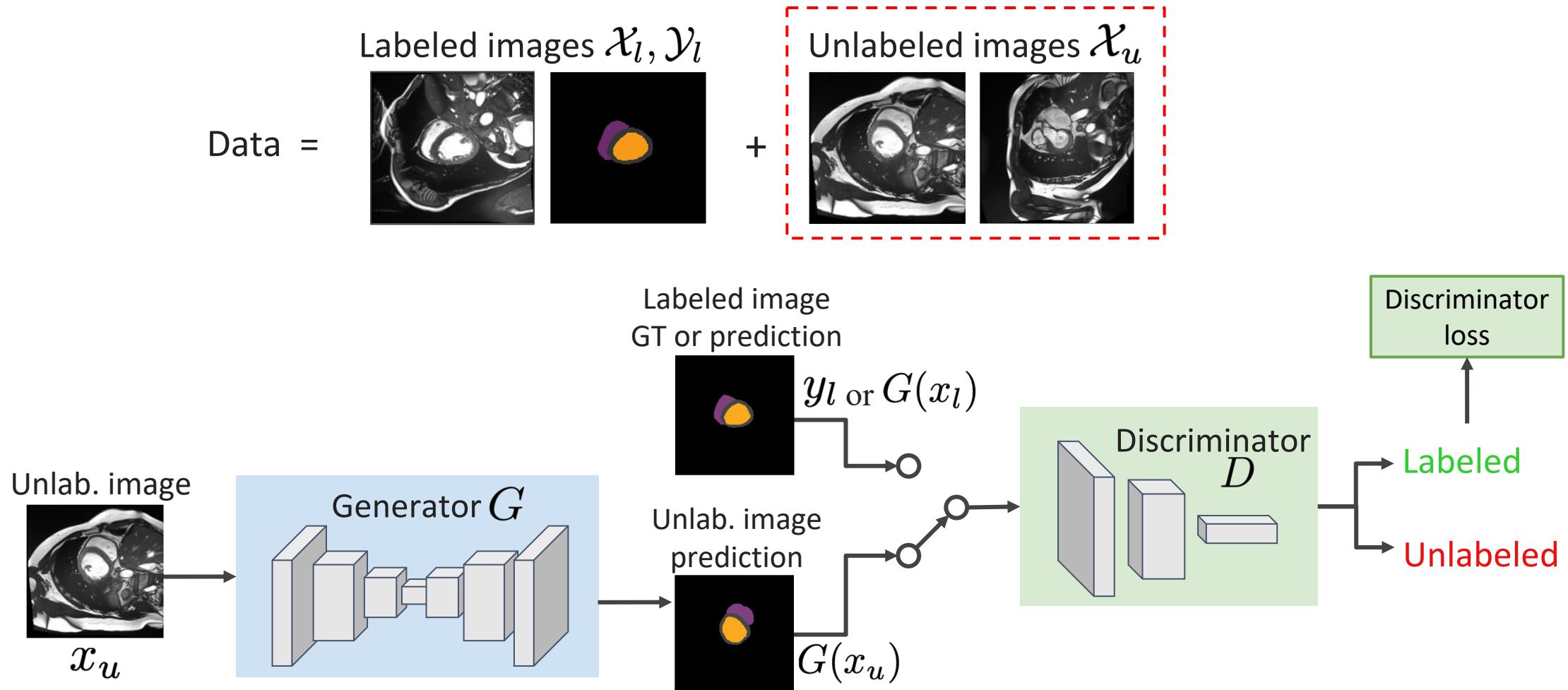
Basic idea: Learn to generate segmentation masks which can't be differentiated from ground-truth (GT)



$$\mathcal{L}_{\text{sup}}(G) = \mathbb{E}_{(x_l, y_l) \sim \mathcal{X}_l, \mathcal{Y}_l} [\mathcal{L}_{\text{seg}}(G(x_l), y_l)]$$

Adversarial semi-supervised segmentation

Basic idea: Learn to generate segmentation masks which can't be differentiated from ground-truth (GT)



Adversarial semi-supervised segmentation

Basic idea: Learn to generate segmentation masks which can't be differentiated from ground-truth (GT)



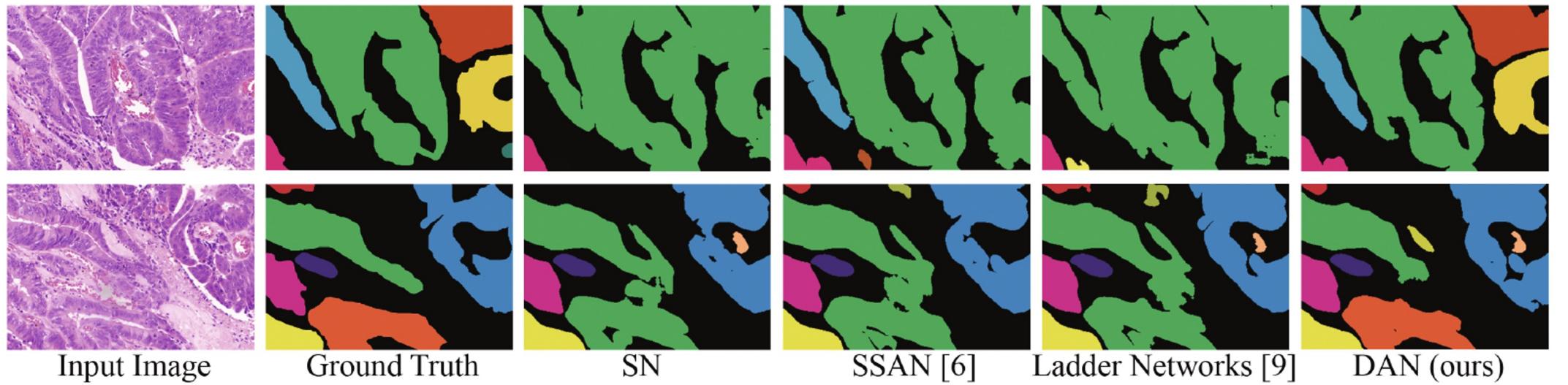
Both labeled and unlabeled:

$$\min_G \max_D \mathcal{L}(G, D) = \underbrace{\frac{1}{|\mathcal{X}_l|} \sum_{l=1}^{|\mathcal{X}_l|} \mathcal{L}_{\text{seg}}(G(x_l), y_l)}_{\text{Supervised loss}} - \underbrace{\frac{\lambda}{|\mathcal{X}_l| + |\mathcal{X}_u|} \left(\sum_{l=1}^{|\mathcal{X}_l|} \mathcal{L}_{\text{dis}}(D(G(x_l)), 1) + \sum_{u=1}^{|\mathcal{X}_u|} \mathcal{L}_{\text{dis}}(D(G(x_u)), 0) \right)}_{\text{Adversarial loss}}$$

Controls the trade-off between the two losses

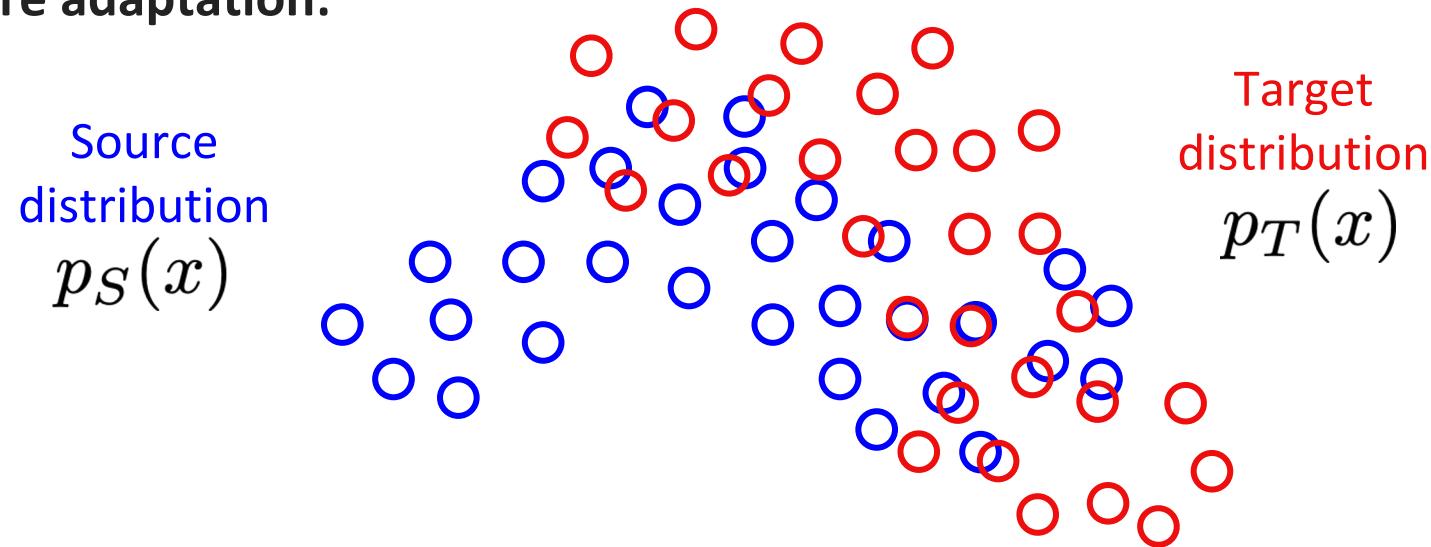
Adversarial semi-supervised segmentation

Adversarial network for semi-supervised segmentation of histological images



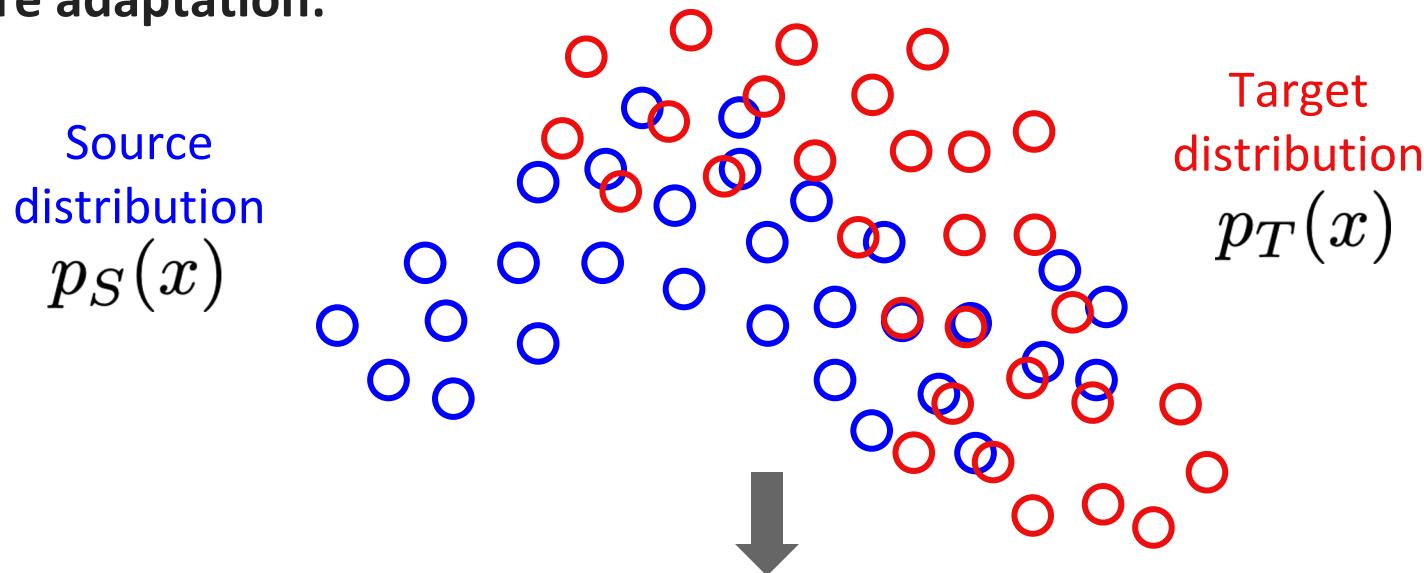
Domain adaptation

Before adaptation:

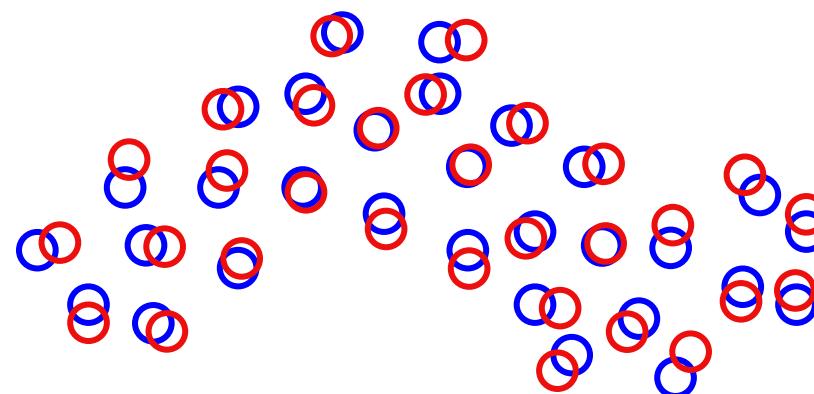


Domain adaptation

Before adaptation:



After adaptation:

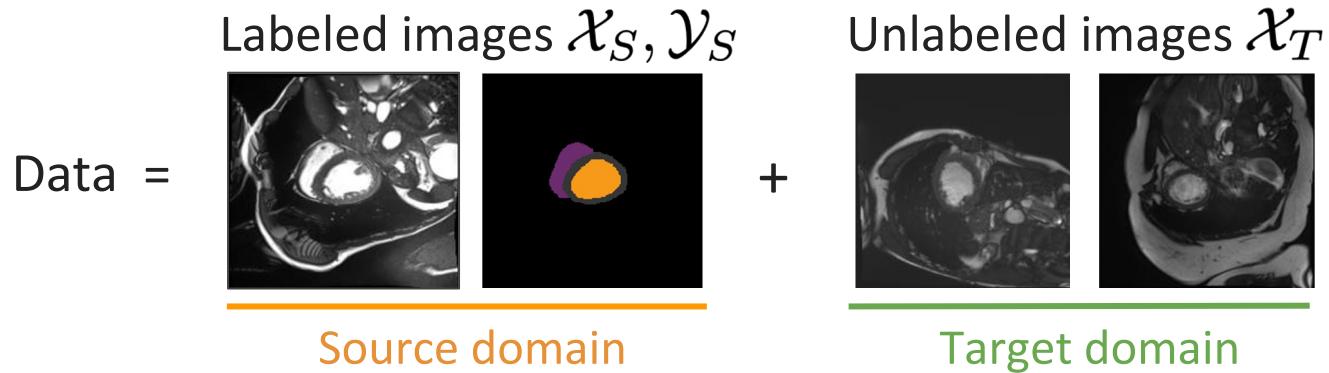


$$p_S(x) = p_T(x)$$

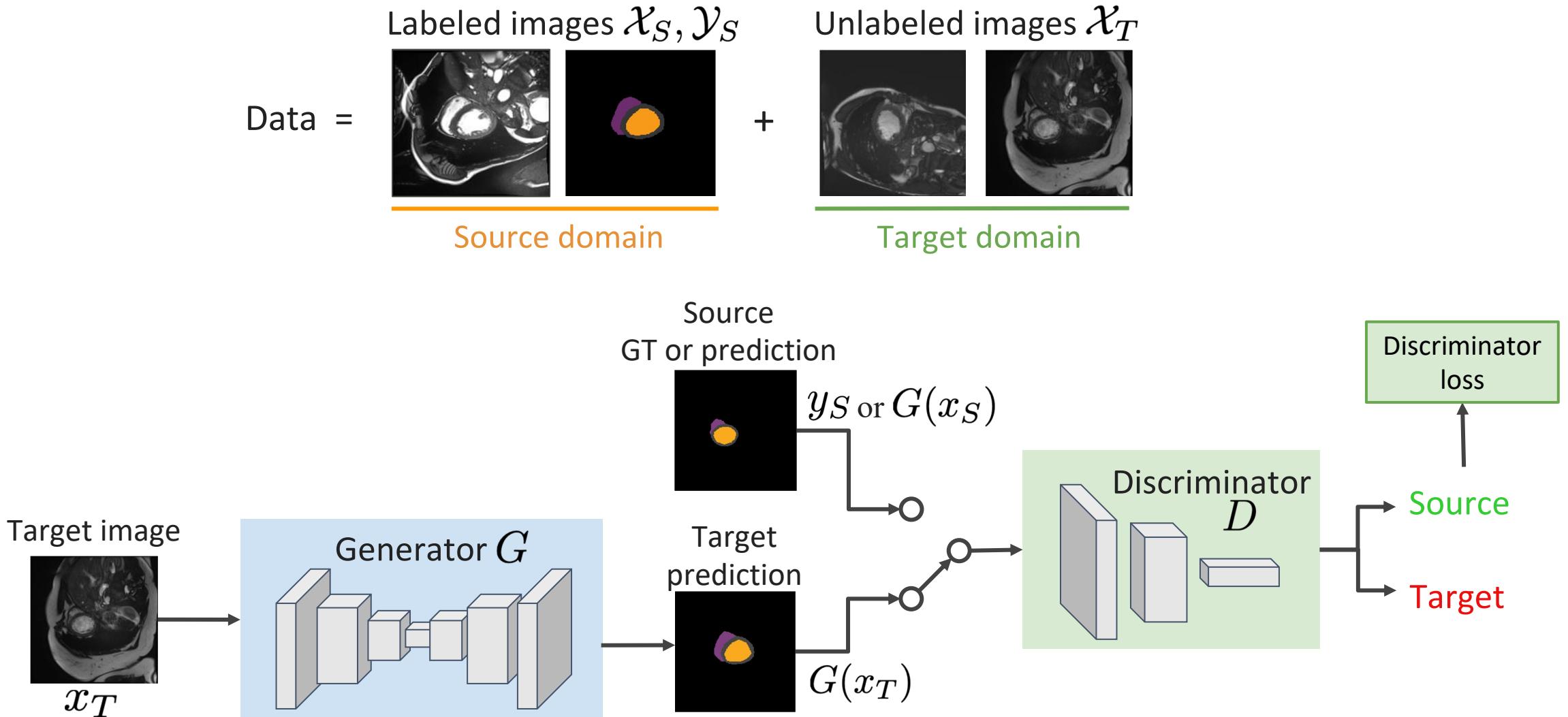
Objective

Align the distributions (input, output or representation) so that a model trained on Source data also works on Target data

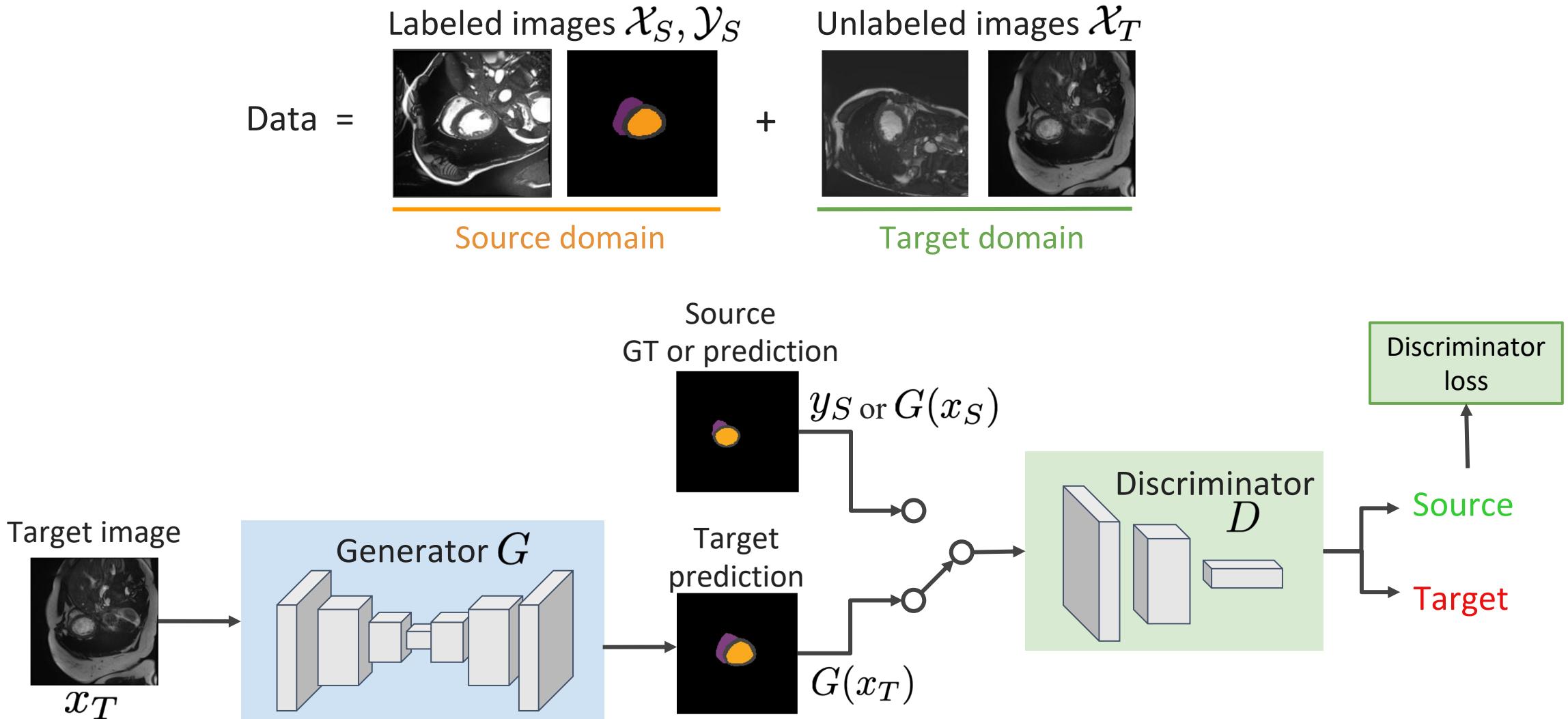
Adversarial domain adaptation



Adversarial domain adaptation



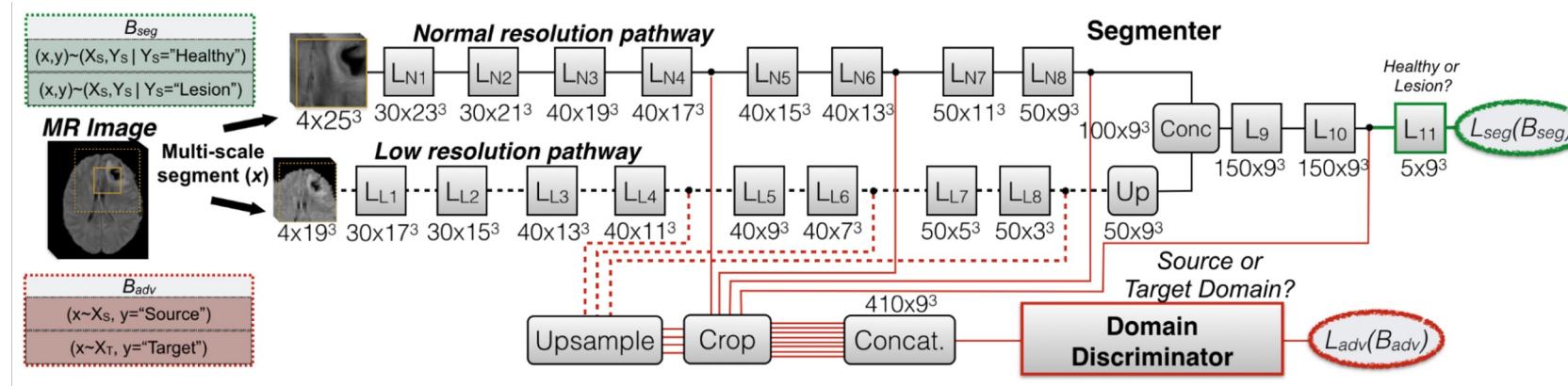
Adversarial domain adaptation



Like semi-supervised segmentation except target images are from a different domain

Adversarial domain adaptation

Adversarial domain adaptation for brain lesion segmentation

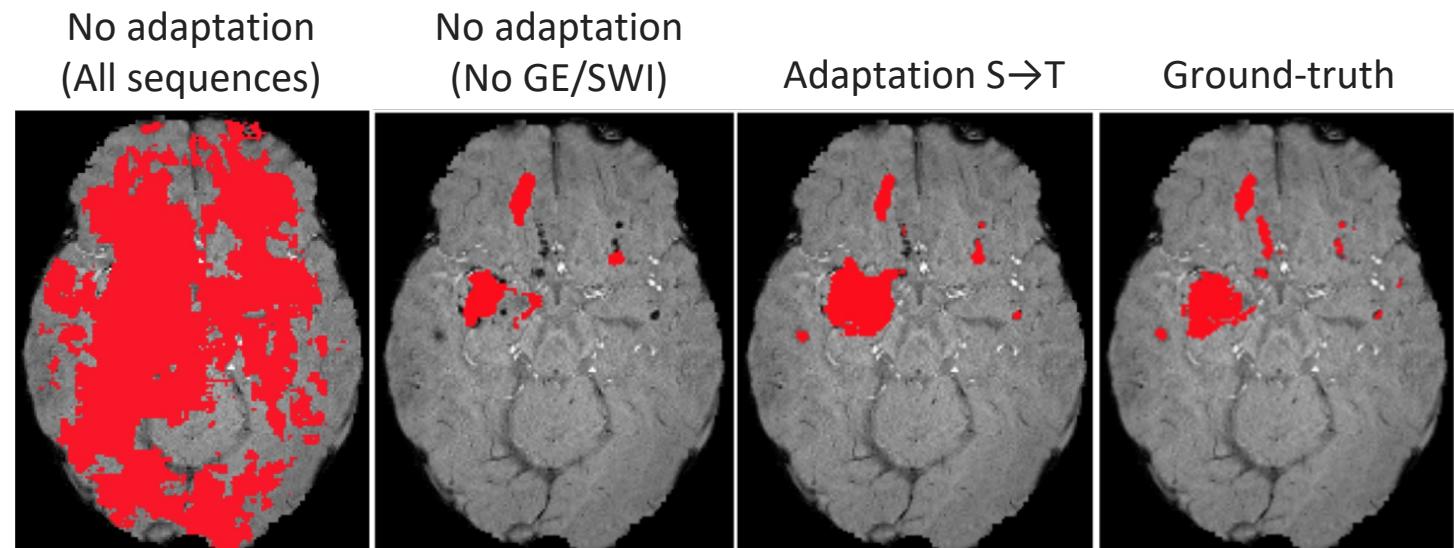


Source domain (Database 1):

- GE, FLAIR, T2, MPRAGE, PD

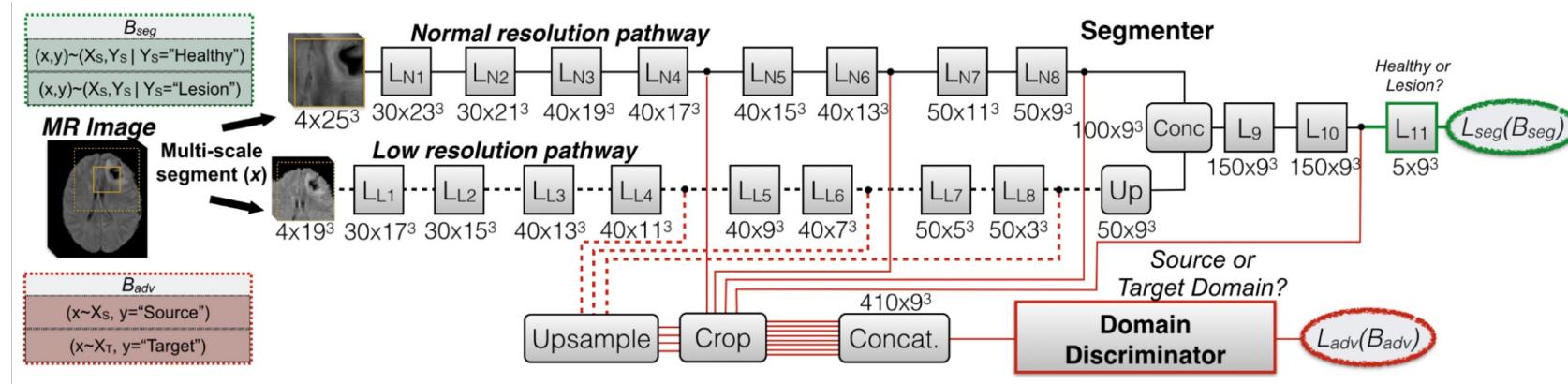
Target domain (Database 2):

- SWI, FLAIR, T2, MPRAGE, PD



Adversarial domain adaptation

Adversarial domain adaptation for brain lesion segmentation



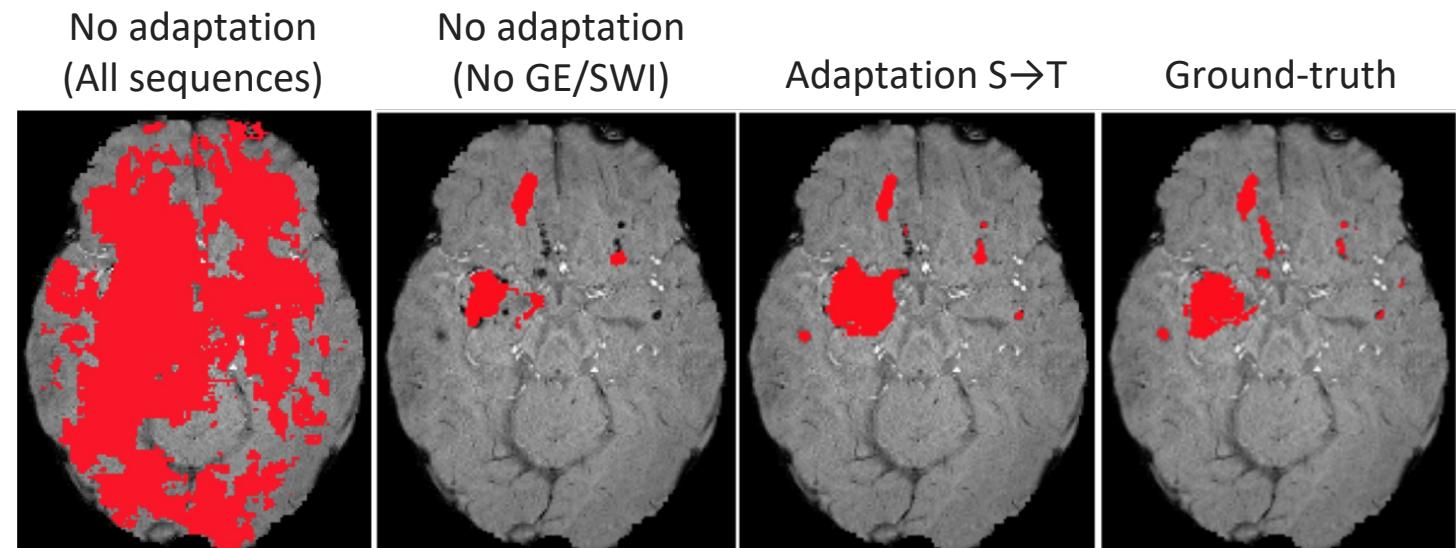
Adaptation done on
multi-scale feature
representation

Source domain (Database 1):

- GE, FLAIR, T2, MPRAGE, PD

Target domain (Database 2):

- SWI, FLAIR, T2, MPRAGE, PD

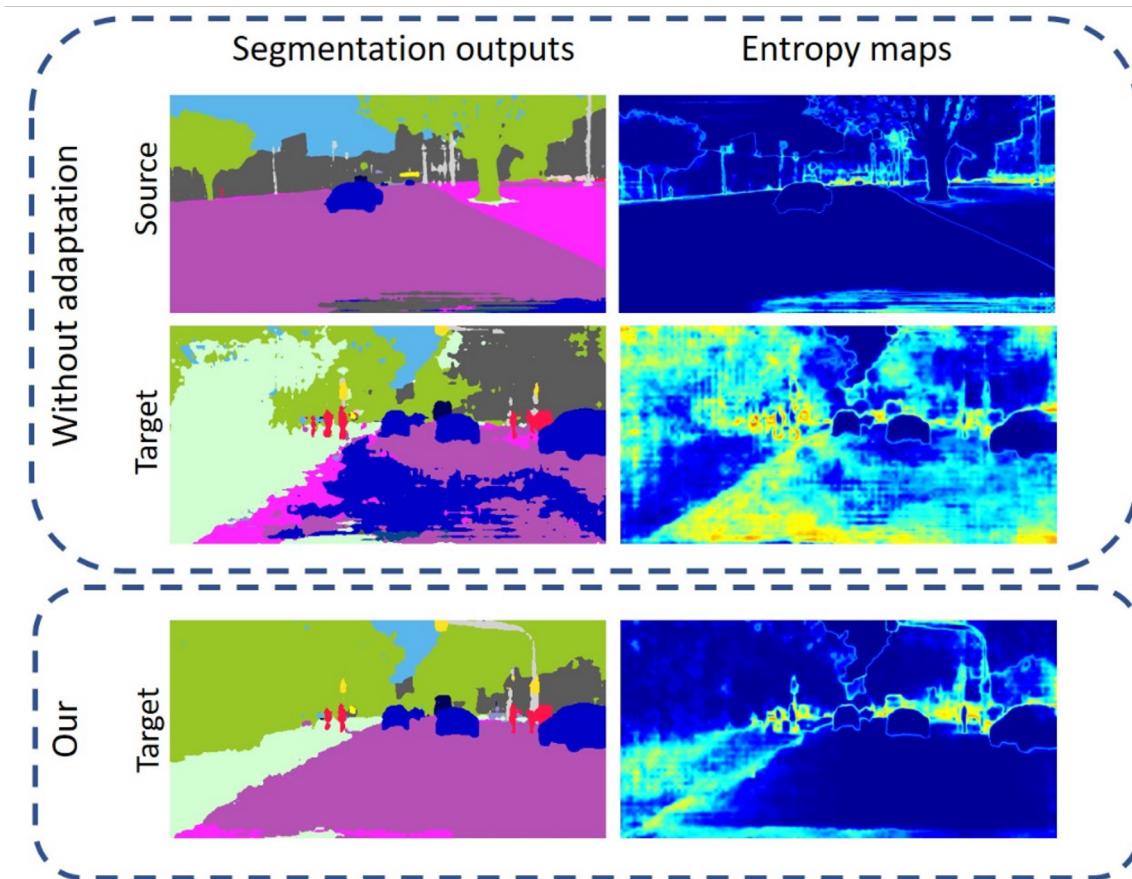


Adversarial domain adaptation

Adaptation on feature *representation* or *softmax output*. What else ?

Adversarial domain adaptation

Adaptation on feature *representation* or softmax output. What else ?



Adversarial entropy minimization

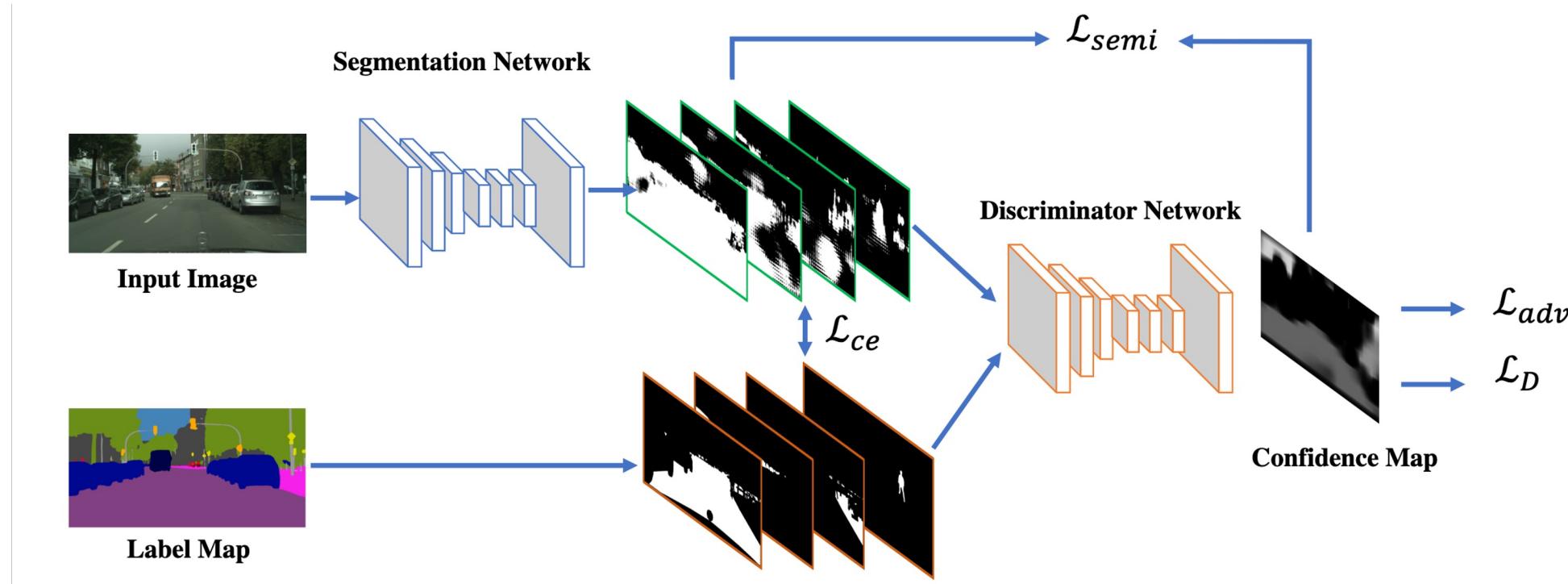
- The discriminator must differentiate between source and target examples using the entropy spatial maps
- Forces the segmentation model to be consistent in its confidence across different semantic regions

Adversarial model for self-training

How can we leverage discriminator predictions at the pixel-level ?

Adversarial model for self-training

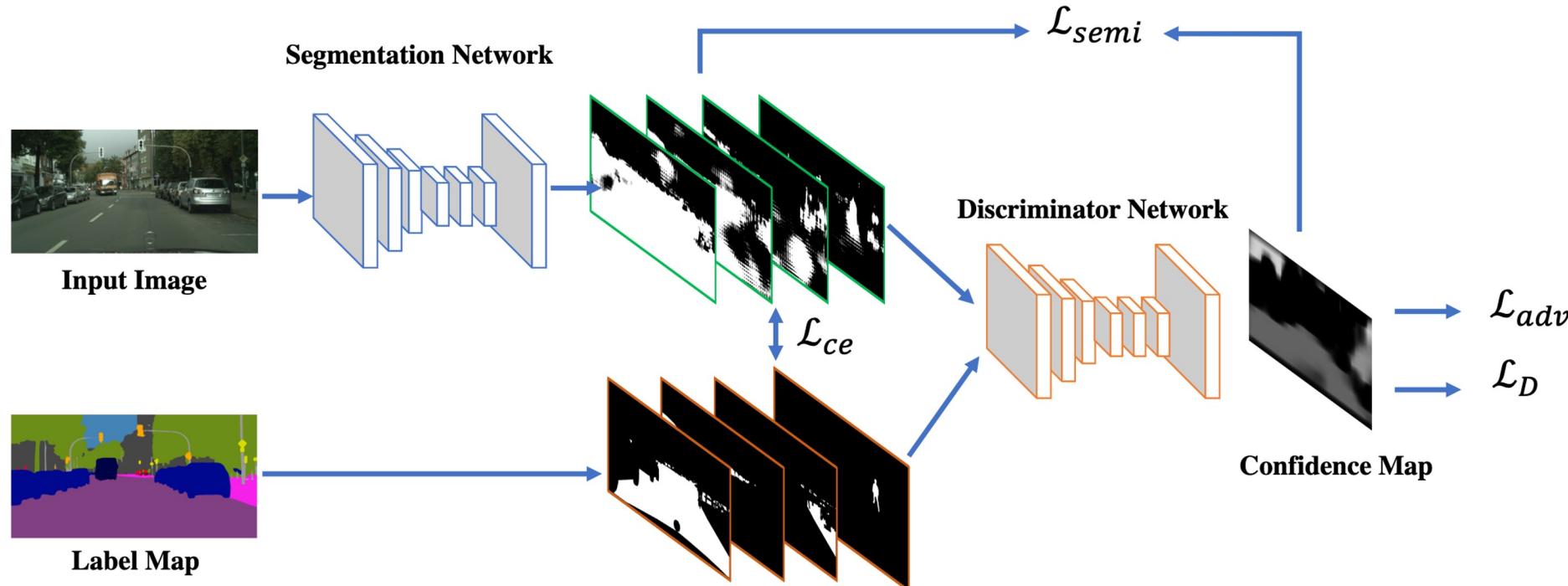
How can we leverage discriminator predictions at the pixel-level ?



- The discriminator must discriminate between prediction and ground-truth (GT) at each pixel
- Consider the discriminator GT-class probabilities as confidence scores
- Use high-confidence predictions on unlabeled images as pseudo-labels for self-training

Adversarial model for self-training

How can we leverage discriminator predictions at the pixel-level ?

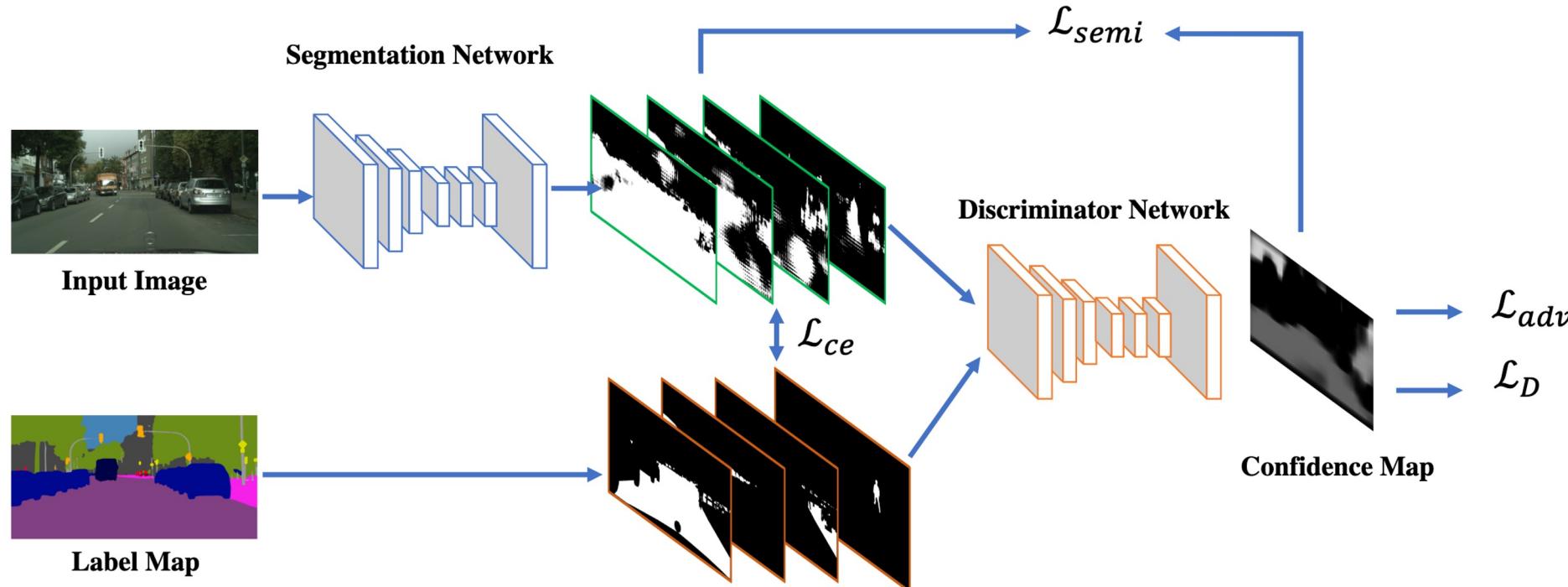


Select pixels with confidence above a given threshold

$$\mathcal{L}_{semi} = - \sum_{h,w} \sum_{c \in C} I(D(S(\mathbf{X}_n))^{(h,w)} > T_{semi}) \cdot \hat{\mathbf{Y}}_n^{(h,w,c)} \log(S(\mathbf{X}_n)^{(h,w,c)})$$
$$\hat{\mathbf{Y}}_n^{(h,w,c^*)} = 1 \text{ if } c^* = \arg \max_c S(\mathbf{X}_n)^{(h,w,c)}$$

Adversarial model for self-training

How can we leverage discriminator predictions at the pixel-level ?



$$\mathcal{L}_{semi} = - \sum_{h,w} \sum_{c \in C} I(D(S(\mathbf{X}_n))^{(h,w)} > T_{semi}) \hat{\mathbf{Y}}_n^{(h,w,c)} \log(S(\mathbf{X}_n)^{(h,w,c)})$$

$$\hat{\mathbf{Y}}_n^{(h,w,c^*)} = 1 \text{ if } c^* = \arg \max_c S(\mathbf{X}_n)^{(h,w,c)}$$

Use class with highest probability as pseudo-label

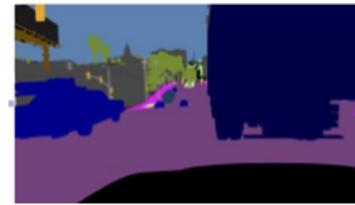
Cycle GANs for domain adaptation

How can we learn a model to segment target images without paired images or GT ?

Source domain



Image



Ground-truth

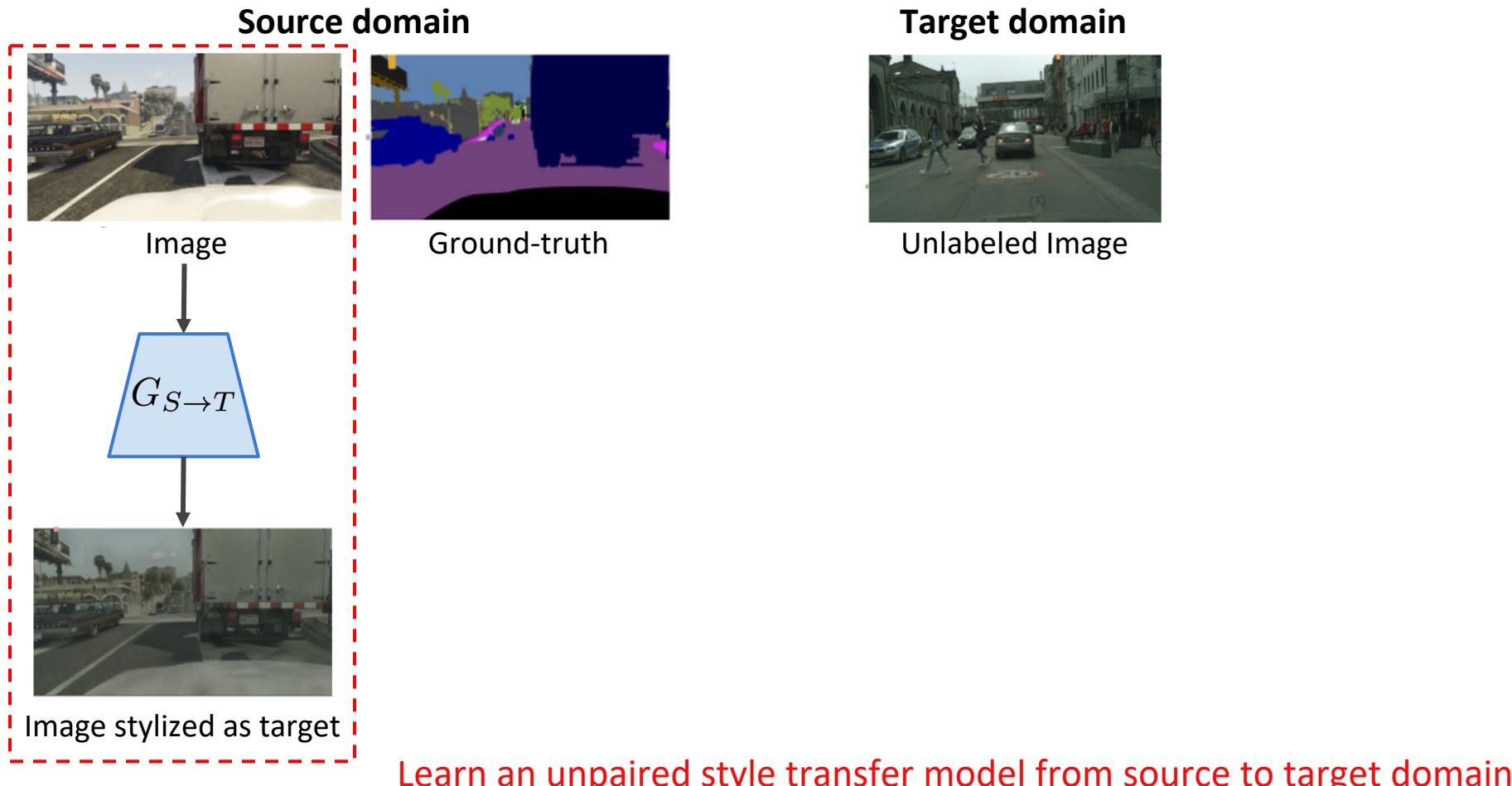
Target domain



Unlabeled Image

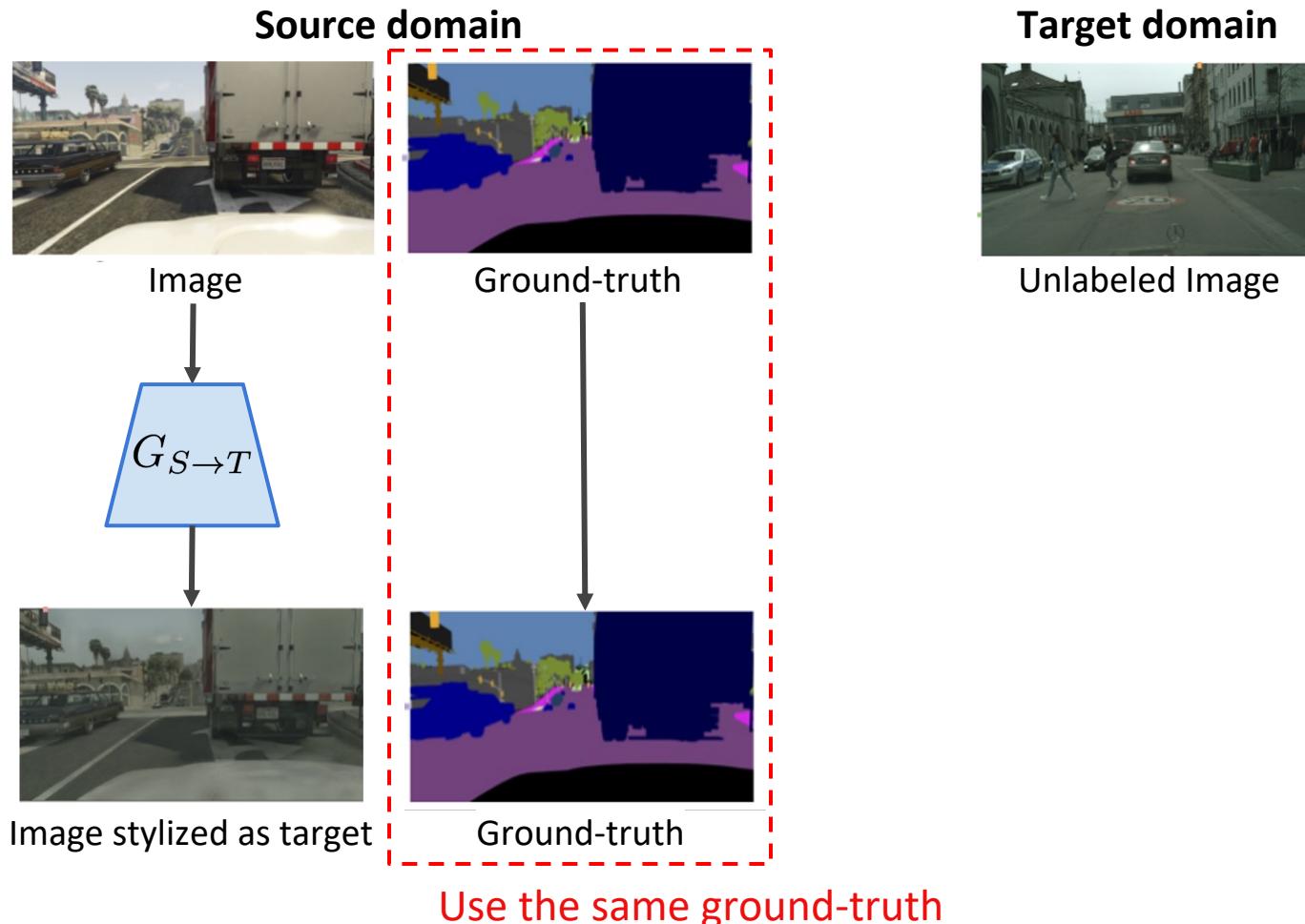
Cycle GANs for domain adaptation

How can we learn a model to segment target images without paired images or GT ?



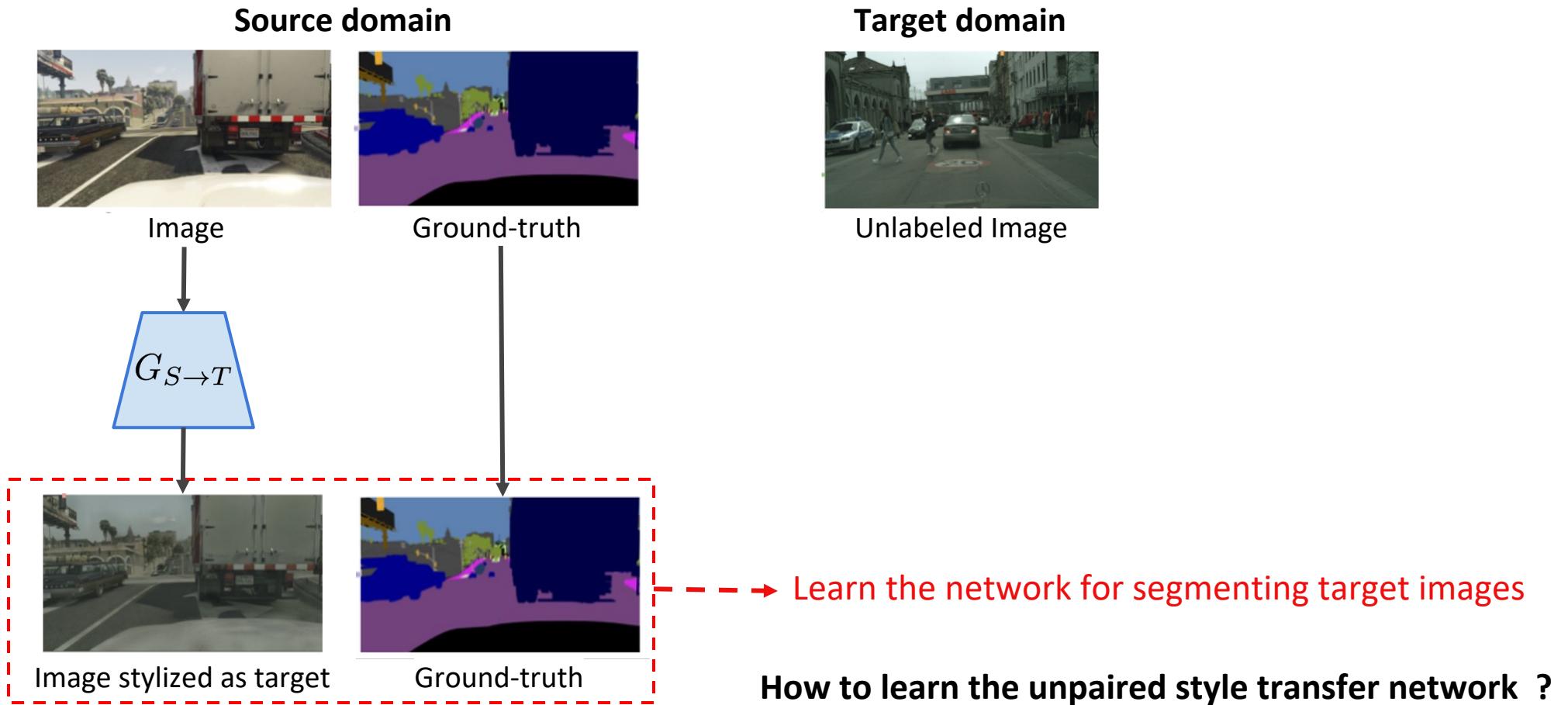
Cycle GANs for domain adaptation

How can we learn a model to segment target images without paired images or GT ?



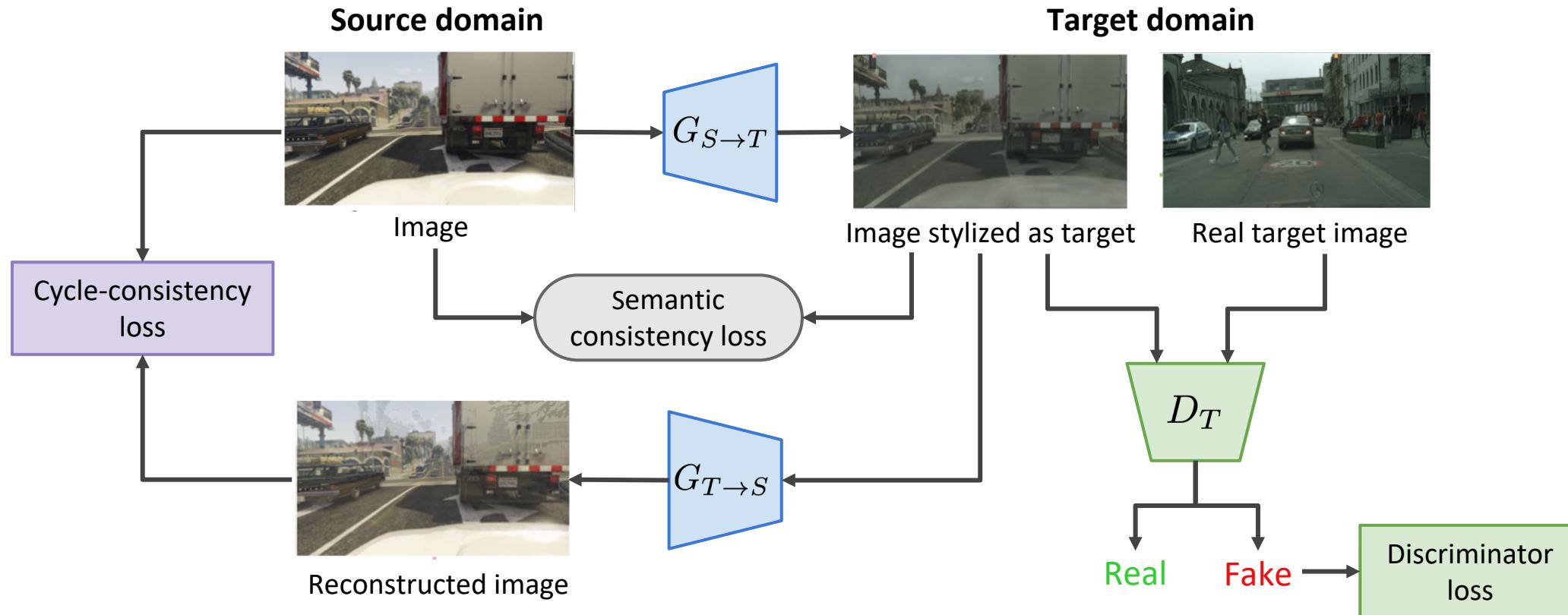
Cycle GANs for domain adaptation

How can we learn a model to segment target images without paired images or GT ?



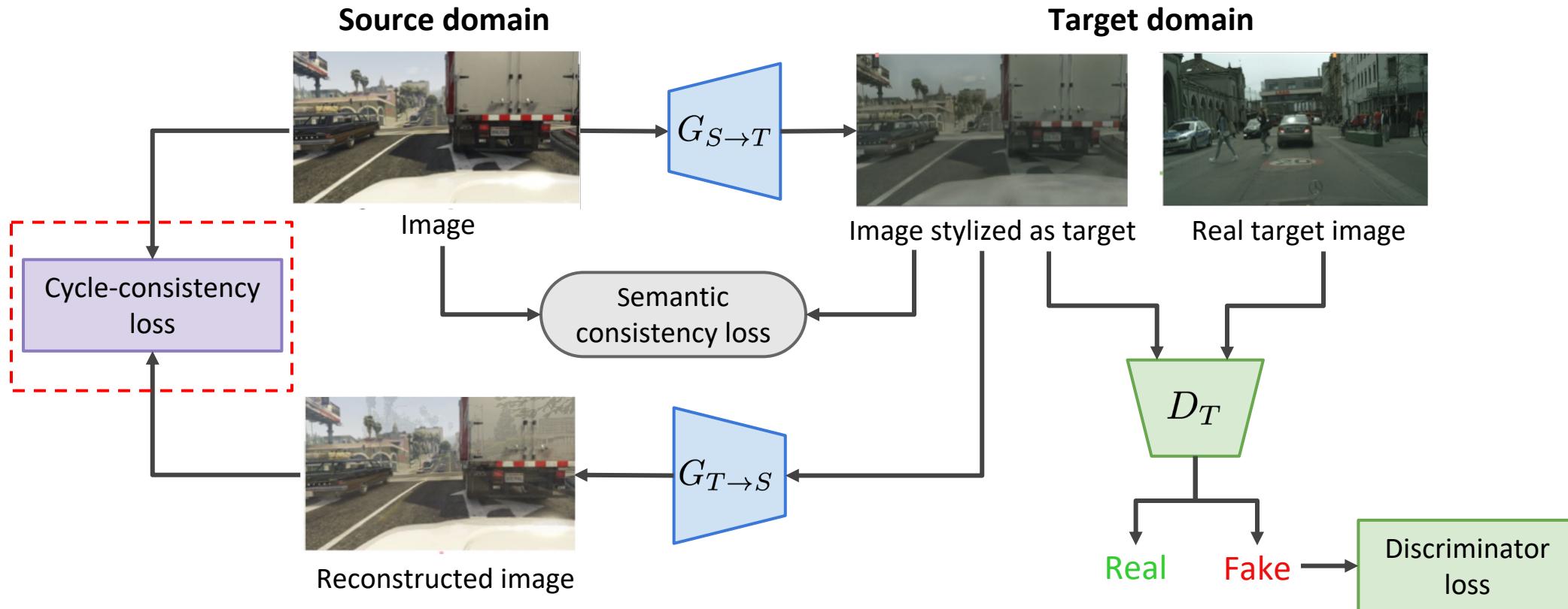
Cycle GANs for domain adaptation

How can we learn a model to segment target images without paired images or GT ?



Cycle GANs for domain adaptation

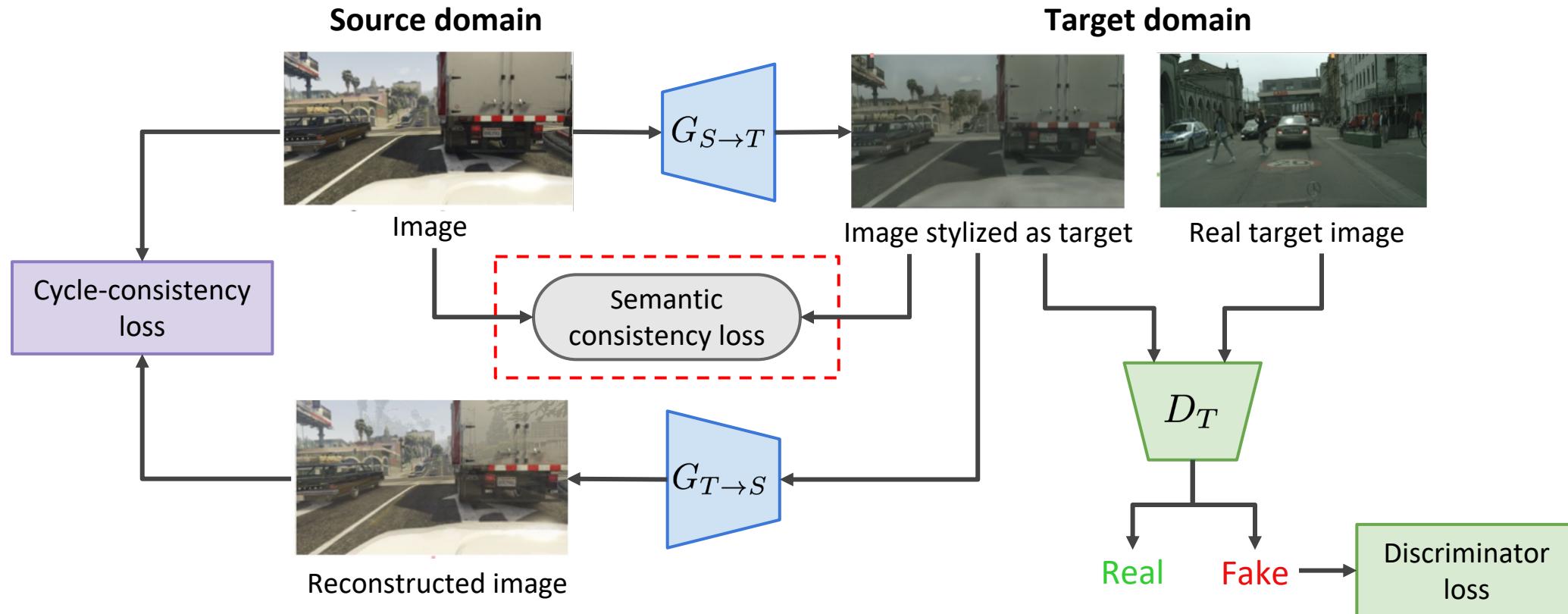
How can we learn a model to segment target images without paired images or GT ?



$$\text{Cycle consistency loss: } L_{\text{cycle}}(G_{S \rightarrow T}, G_{T \rightarrow S}) = \mathbb{E}_{x \sim p_S(x)} \left[\|x - G_{T \rightarrow S}(G_{S \rightarrow T}(x))\|_1 \right]$$

Cycle GANs for domain adaptation

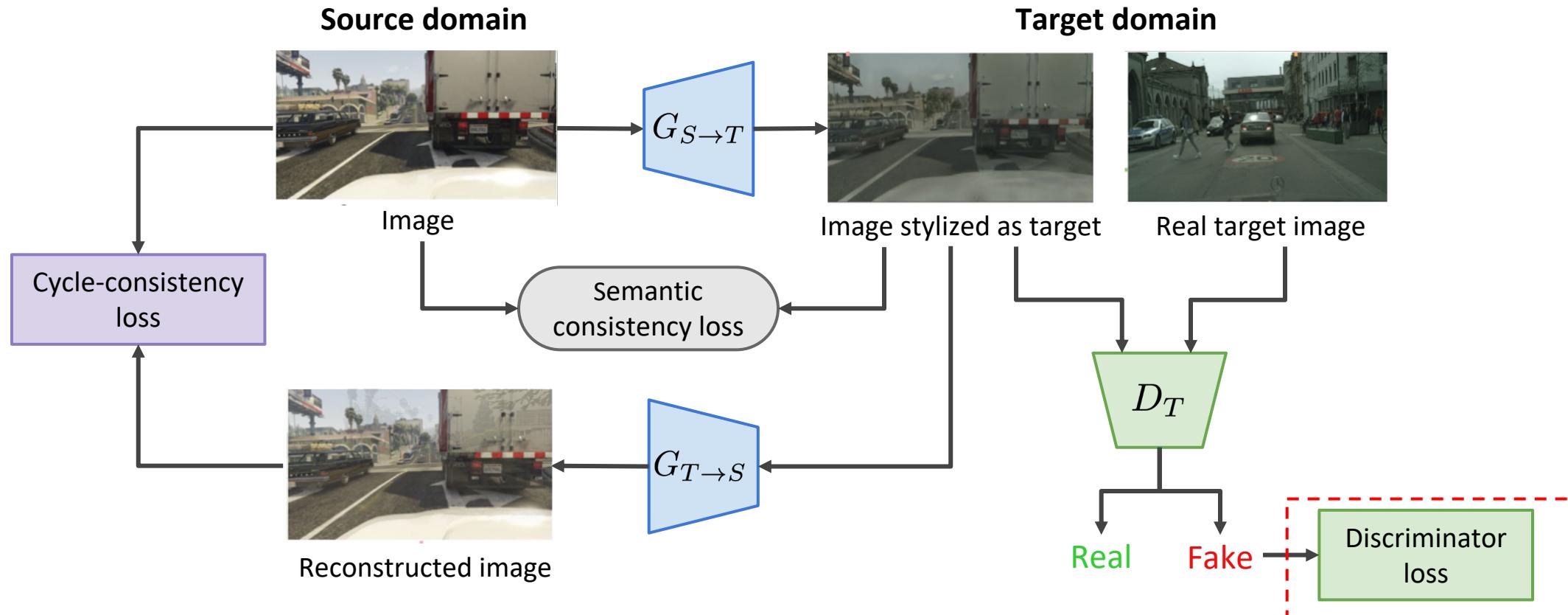
How can we learn a model to segment target images without paired images or GT ?



Semantic consistency loss: Segmentation for the source image and its stylized target version should be consistent

Cycle GANs for domain adaptation

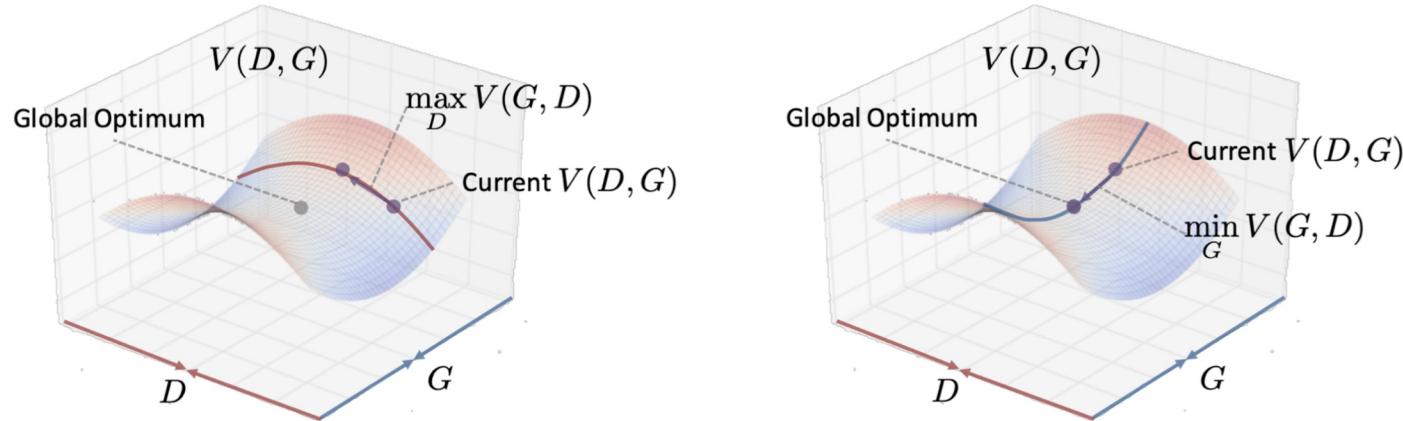
How can we learn a model to segment target images without paired images or GT ?



Discriminator loss: Target images generated from source should look like real target ones

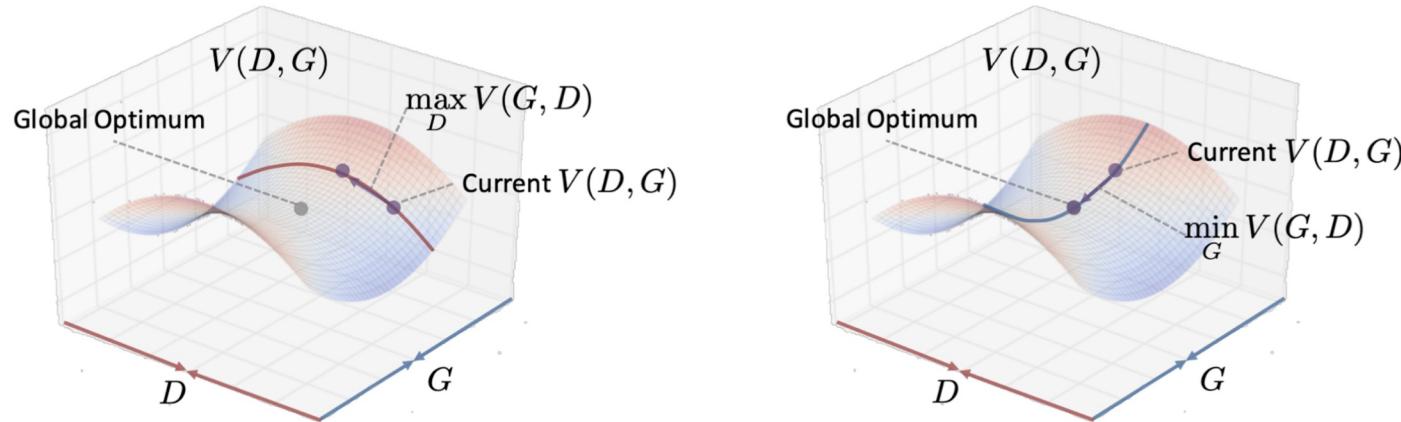
Challenges of adversarial learning

1) Unstable optimization of minimax problem

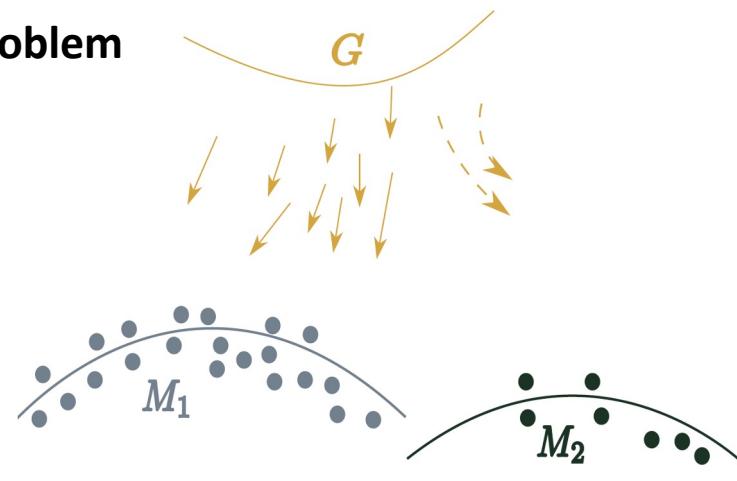


Challenges of adversarial learning

1) Unstable optimization of minimax problem

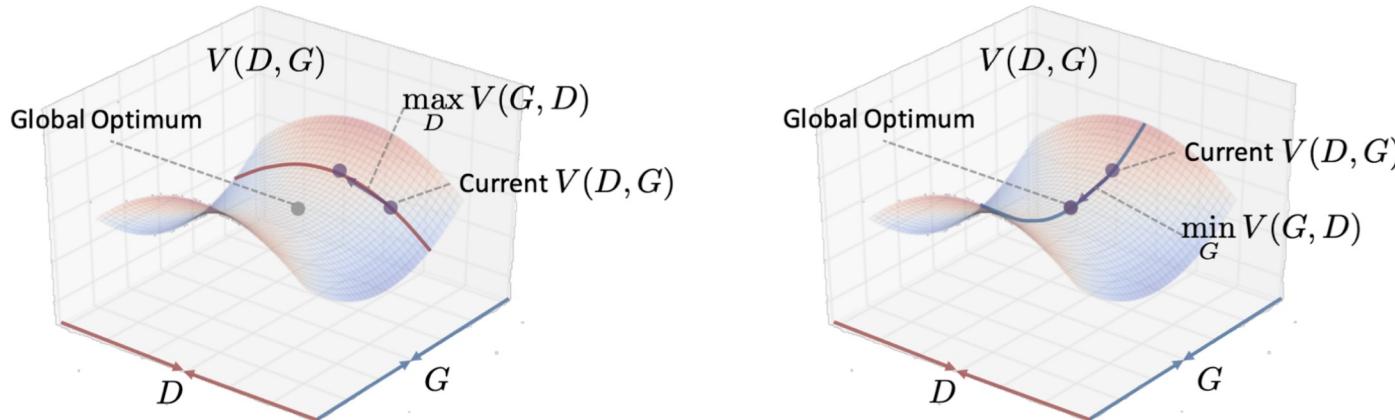


2) Mode collapse problem

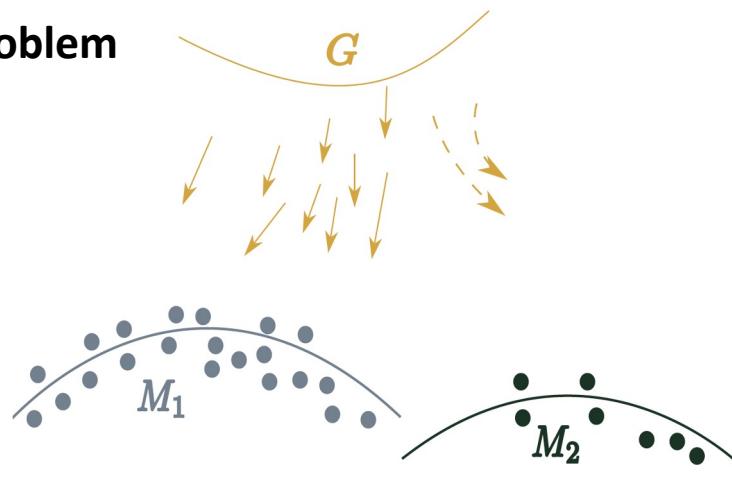


Challenges of adversarial learning

1) Unstable optimization of minimax problem



2) Mode collapse problem



Various solutions:

- Spectral normalization (Miyato *et al.*, 2018)
- Wasserstein GANs (Arjovsky *et al.*, 2017)
- LSGANs (Mao *et al.*, 2017)
- etc.

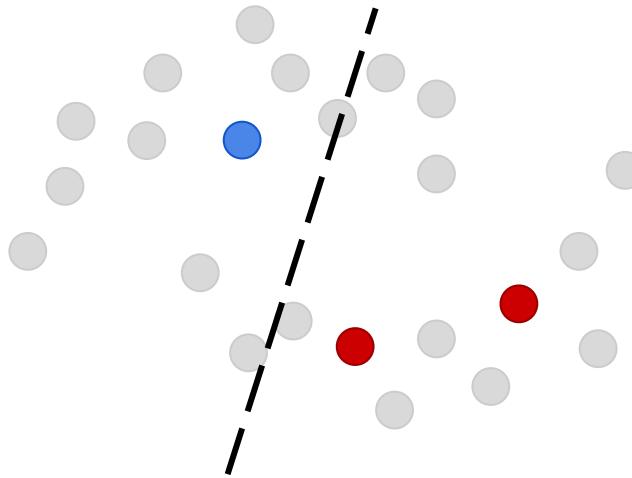
Consistency regularization

for semi-supervised segmentation

Consistency regularization for SSL

How to better use unlabeled data ?

Vanilla supervised learning

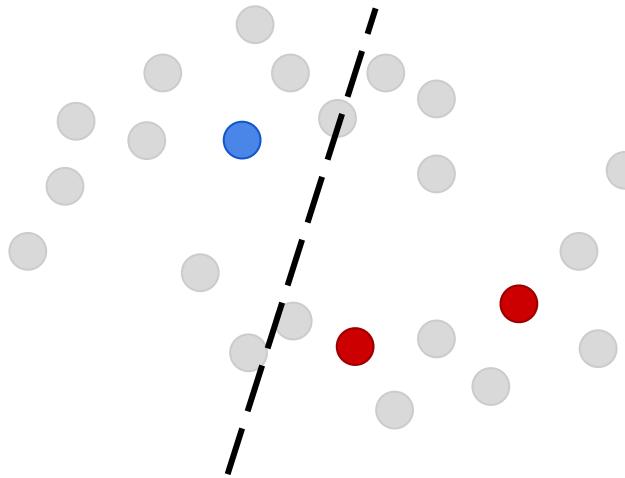


- Consider only labeled samples
- Overfits when few training samples

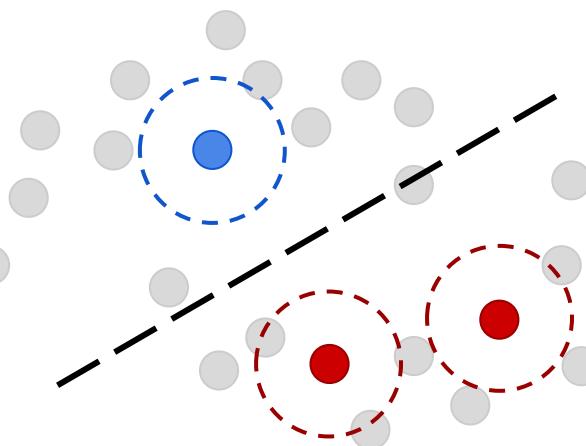
Consistency regularization for SSL

How to better use unlabeled data ?

Vanilla supervised learning



Data augmentation



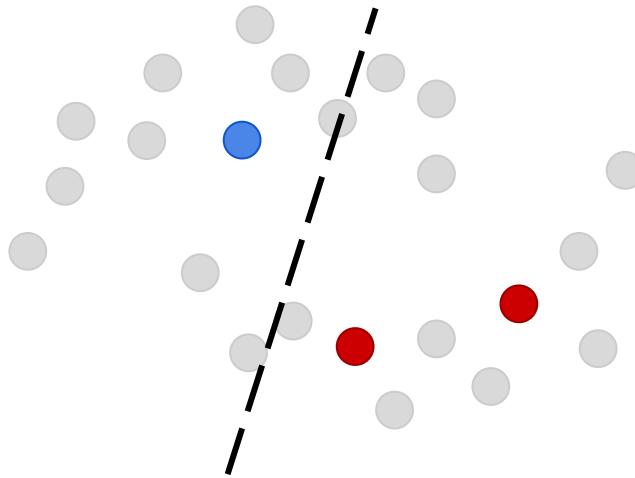
- Consider only labeled samples
- Overfits when few training samples

- Transform labeled samples to augment the training set
- Better generalization, but not enough for semi-supervised learning

Consistency regularization for SSL

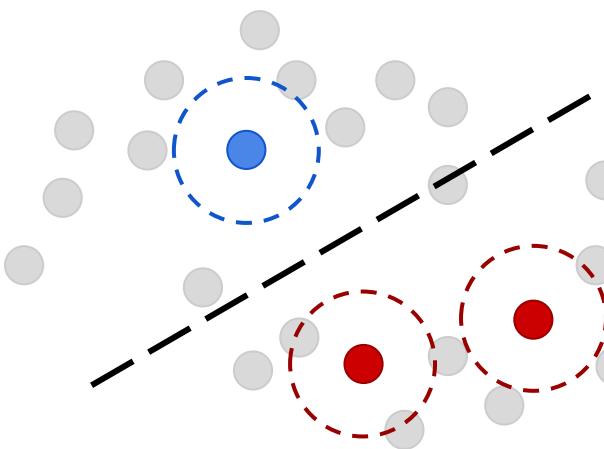
How to better use unlabeled data ?

Vanilla supervised learning



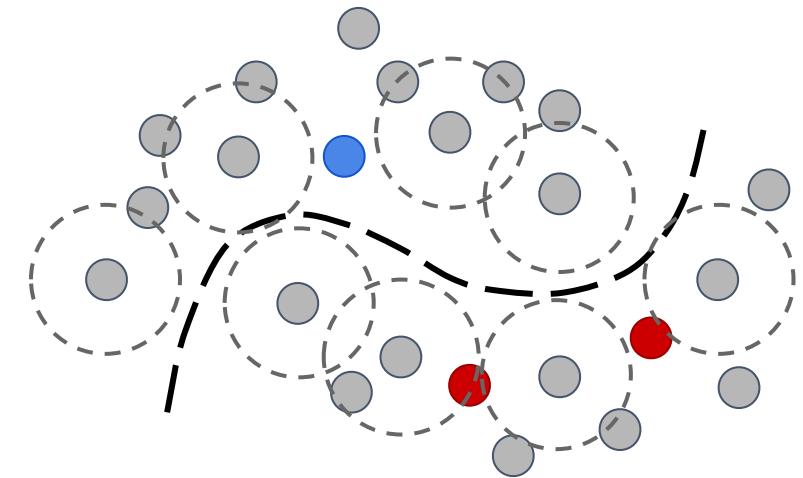
- Consider only labeled samples
- Overfits when few training samples

Data augmentation



- Transform labeled samples to augment the training set
- Better generalization, but not enough for semi-supervised learning

Consistency regularization



- Perturb unlabeled samples with noise or guided transformations
- Impose the network to have consistent outputs for perturbed samples

SSL methods using consistency regularization

Basic transformation consistency (Γ -model)

$$\mathcal{L}(\theta; \mathcal{D}_l, \mathcal{D}_u) = \frac{1}{|\mathcal{D}_l|} \sum_{(x,y) \in \mathcal{D}_l} \ell_{\text{sup}}(f(x), y) + \frac{\lambda}{|\mathcal{D}_u|} \sum_{x \in \mathcal{D}_u} \mathbb{E}_{T \sim p_T} [\ell_{\text{reg}}(T(f(x)), f(T(x)))]$$

SSL methods using consistency regularization

Basic transformation consistency (Γ -model)

Standard supervised loss

$$\mathcal{L}(\theta; \mathcal{D}_l, \mathcal{D}_u) = \boxed{\frac{1}{|\mathcal{D}_l|} \sum_{(x,y) \in \mathcal{D}_l} \ell_{\text{sup}}(f(x), y)} + \frac{\lambda}{|\mathcal{D}_u|} \sum_{x \in \mathcal{D}_u} \mathbb{E}_{T \sim p_T} [\ell_{\text{reg}}(T(f(x)), f(T(x)))]$$

Cross-entropy, Dice, etc.



SSL methods using consistency regularization

Basic transformation consistency (Γ -model)

$$\mathcal{L}(\theta; \mathcal{D}_l, \mathcal{D}_u) = \frac{1}{|\mathcal{D}_l|} \sum_{(x,y) \in \mathcal{D}_l} \ell_{\text{sup}}(f(x), y) +$$

Transformation consistency loss

$$\frac{\lambda}{|\mathcal{D}_u|} \sum_{x \in \mathcal{D}_u} \mathbb{E}_{T \sim p_T} [\ell_{\text{reg}}(T(f(x)), f(T(x)))]$$

Random transformation:
rotation, flip, crop, etc.

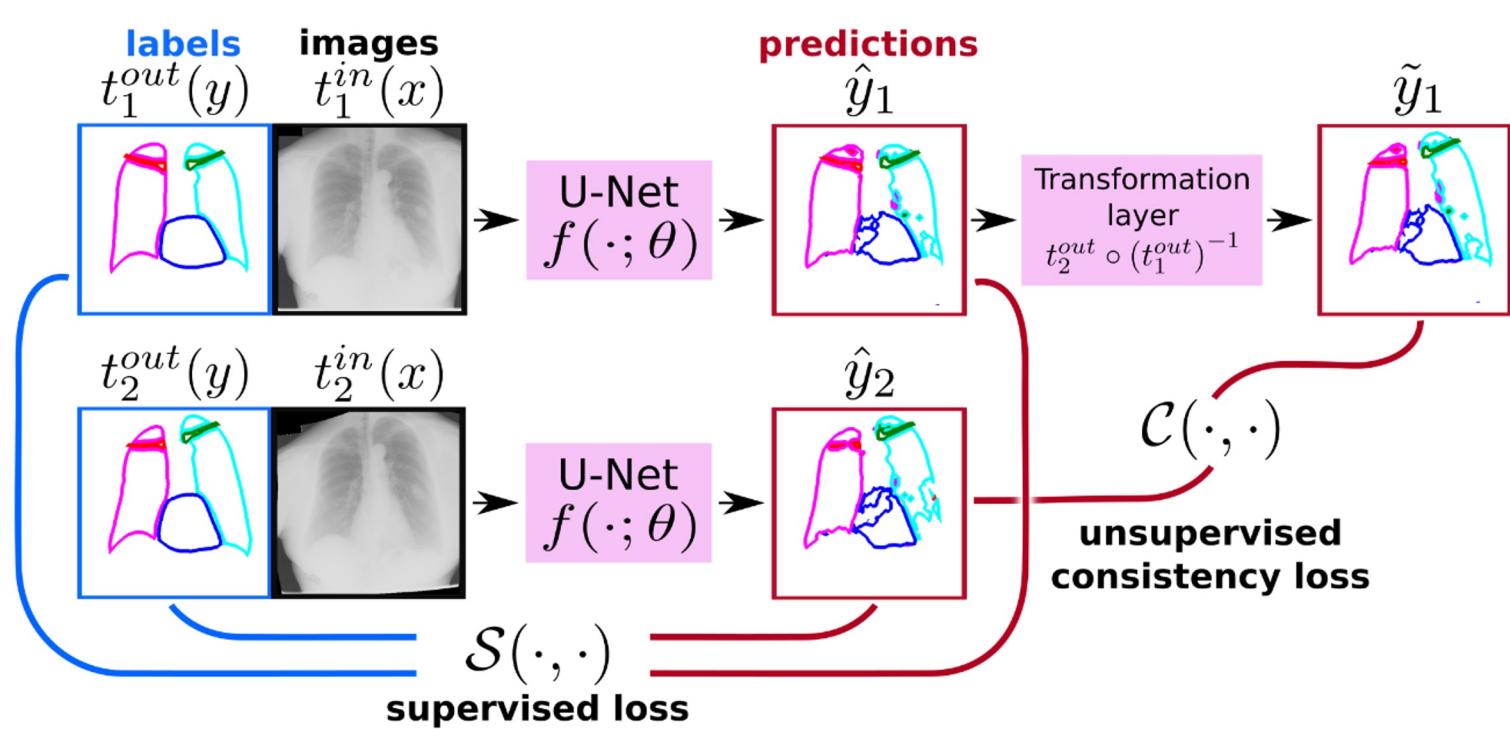
Regularization loss
imposing transformation equivariance

L2 regularization loss:

$$\ell_{\text{reg}}(T(f(x)), f(T(x))) = \|T(f(x)) - f(T(x))\|_2^2$$

SSL methods using consistency regularization

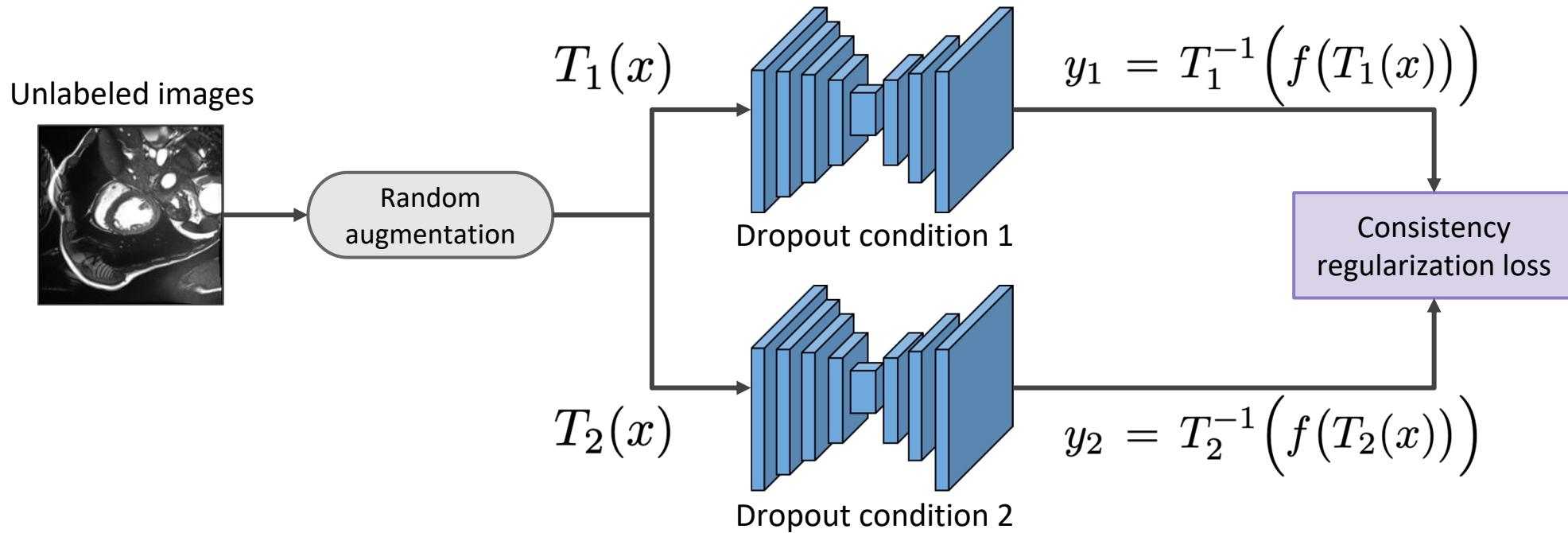
Application to chest X-ray segmentation:



Transformations are random elastic deformations

SSL methods using consistency regularization

Self-ensembling (Π -model):

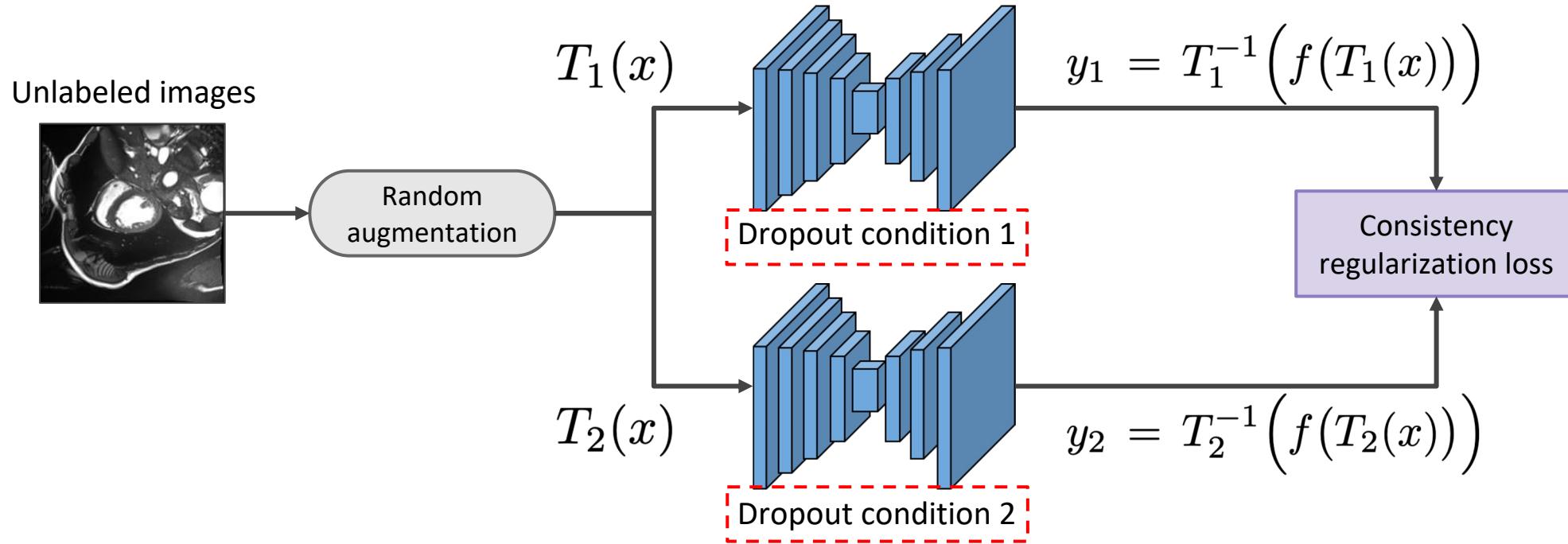


Key idea:

- Applying different dropouts on the same network gives an ensemble of models
- Also leverages random image transformations

SSL methods using consistency regularization

Self-ensembling (Π -model):

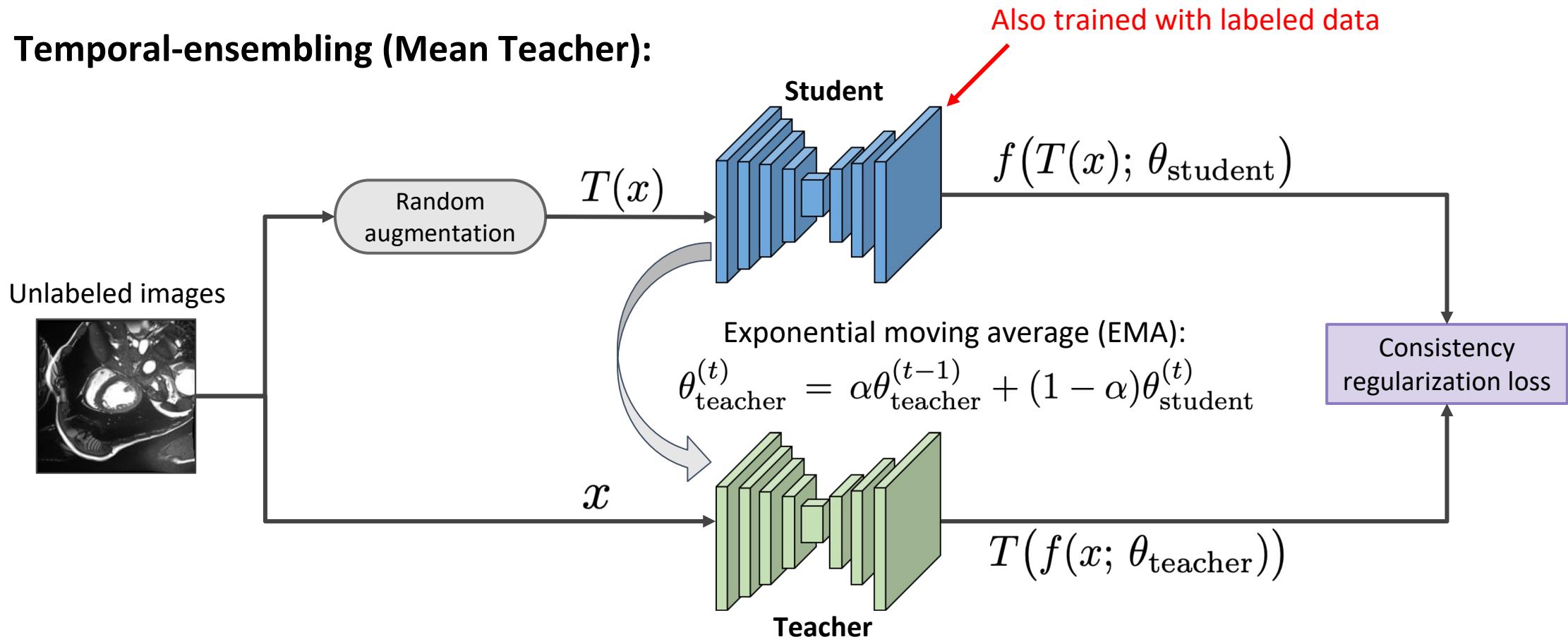


Key idea:

- Applying different dropouts on the same network gives an ensemble of models
- Also leverages random image transformations

SSL methods using consistency regularization

Temporal-ensembling (Mean Teacher):

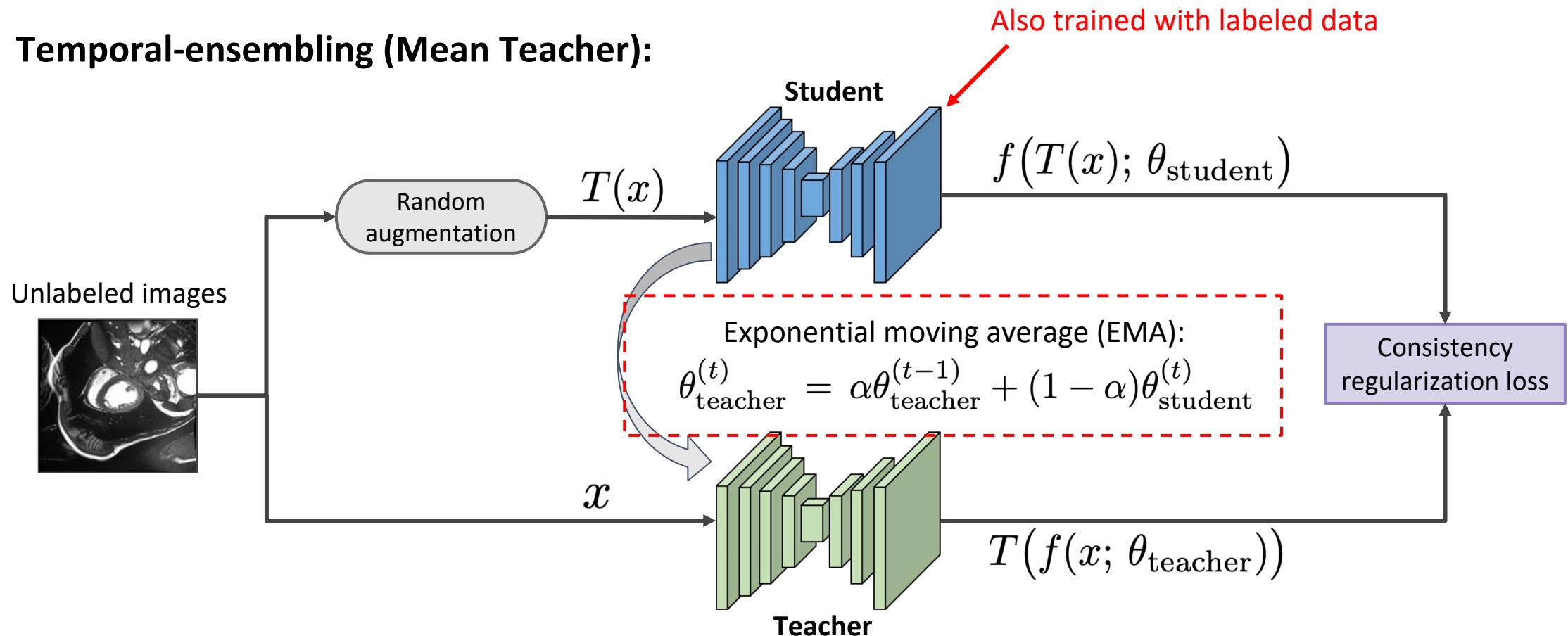


Key idea:

- Consistency between the predictions of a Teacher and a Student network
- The Teacher's weights are an EMA of the Student's at previous training iterations ($\alpha \approx 1$)
- Note: original Temporal Ensembling computes the EMA on outputs for each sample

SSL methods using consistency regularization

Temporal-ensembling (Mean Teacher):

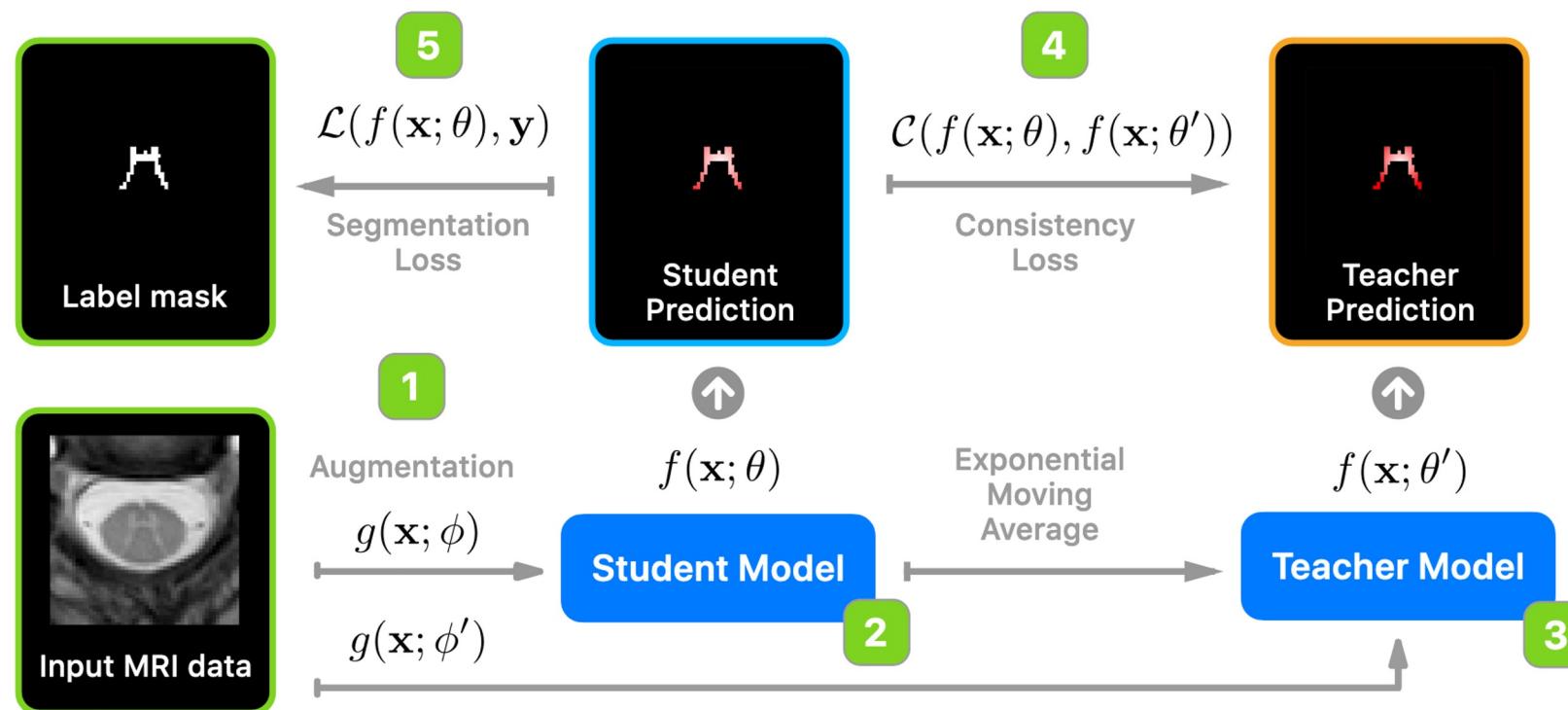


Key idea:

- Consistency between the predictions of a Teacher and a Student network
- The Teacher's weights are an EMA of the Student's at previous training iterations
- Note: original Temporal Ensembling computes the EMA on outputs for each sample

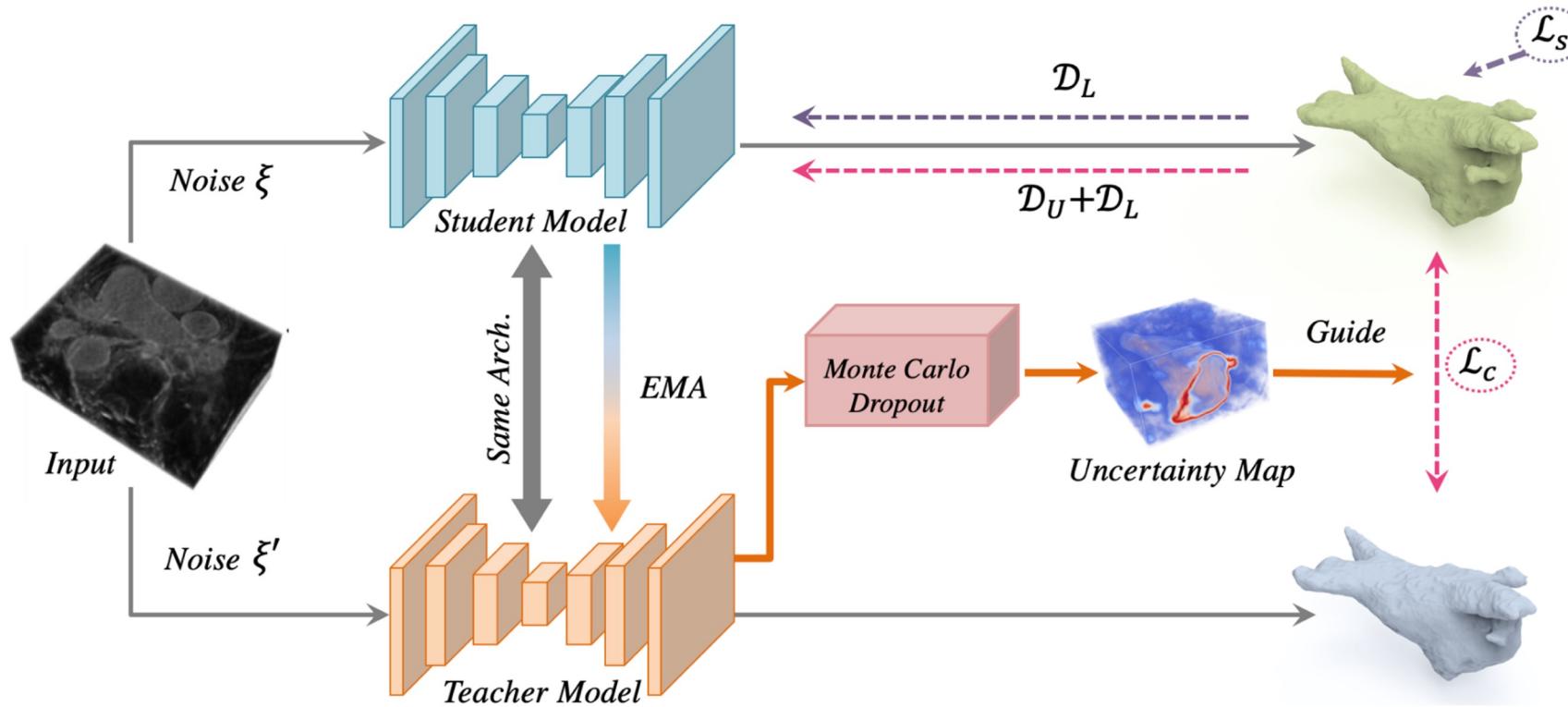
SSL methods using consistency regularization

Application of Mean Teacher to segmenting MRI spinal cord gray matter



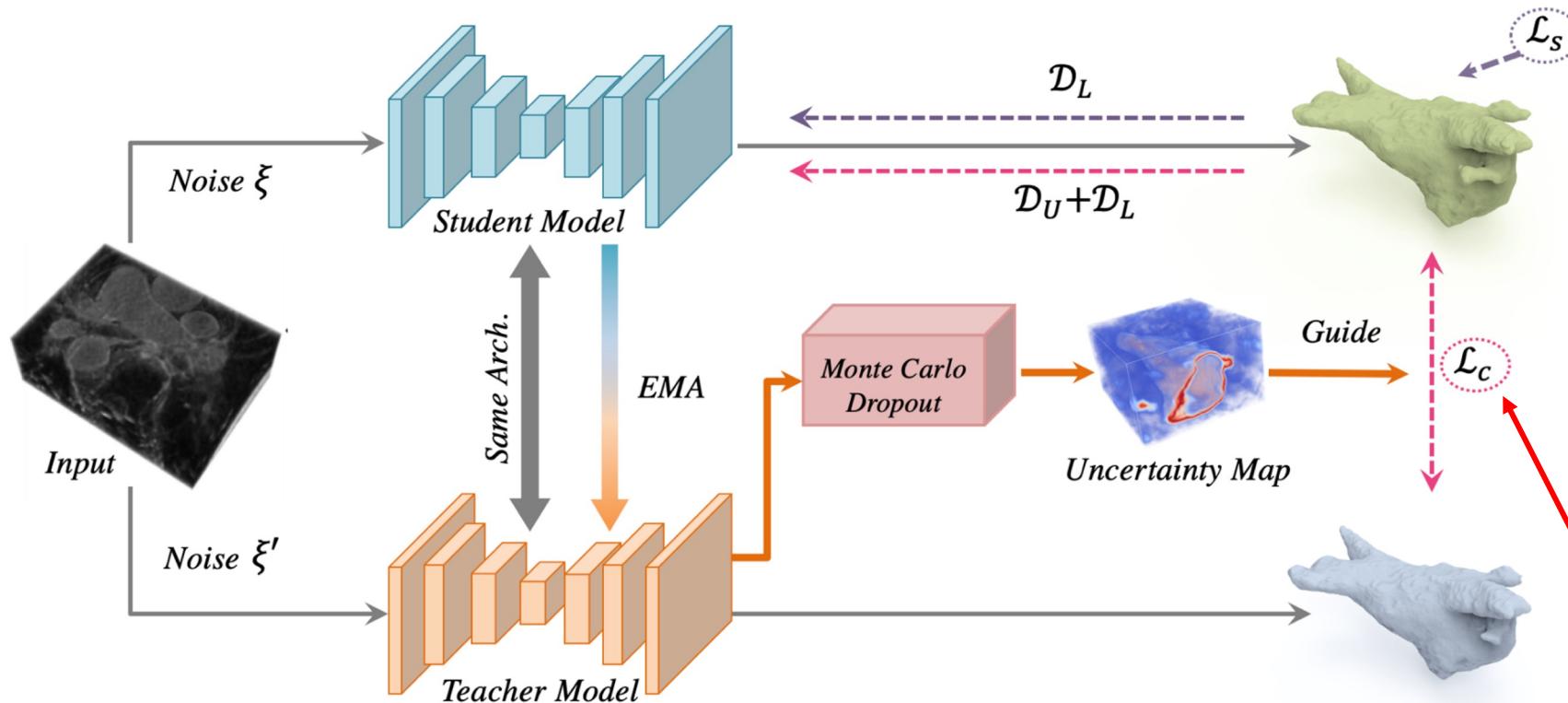
SSL methods using consistency regularization

Uncertainty-aware self-ensembling



SSL methods using consistency regularization

Uncertainty-aware self-ensembling

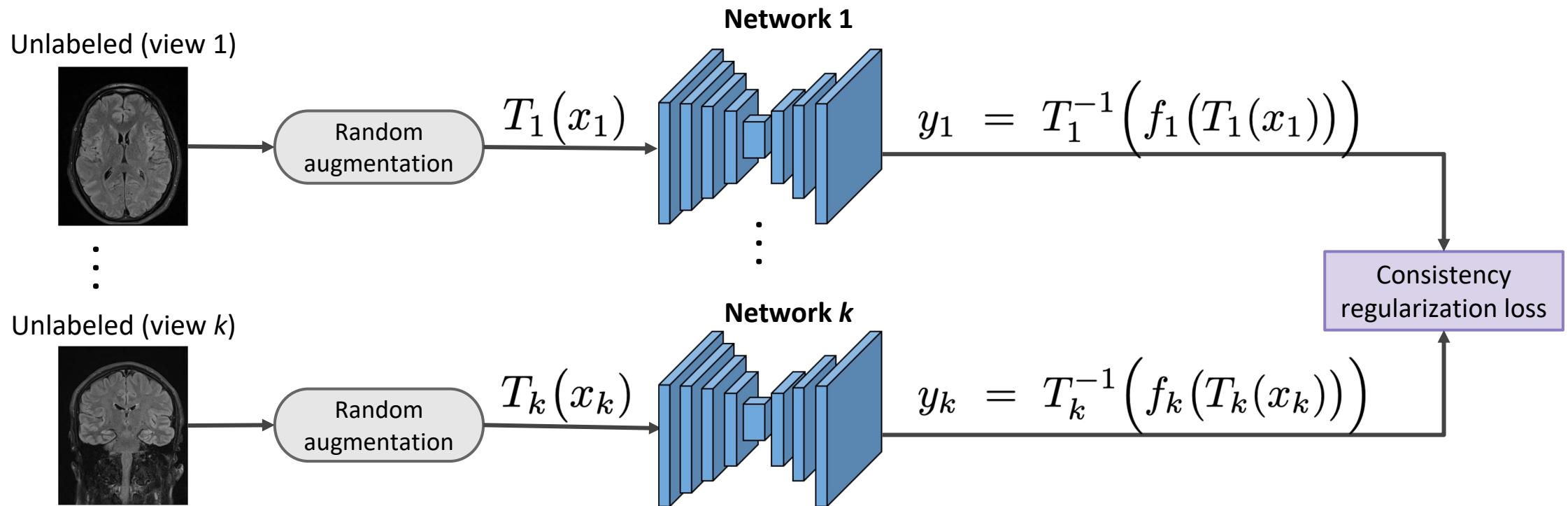


Enforces consistency only in low-uncertainty regions of the image

$$\mathcal{L}_c(f', f) = \frac{\sum_v \mathbb{I}(u_v < H) \|f'_v - f_v\|^2}{\sum_v \mathbb{I}(u_v < H)},$$

SSL methods using consistency regularization

Muti-view co-training

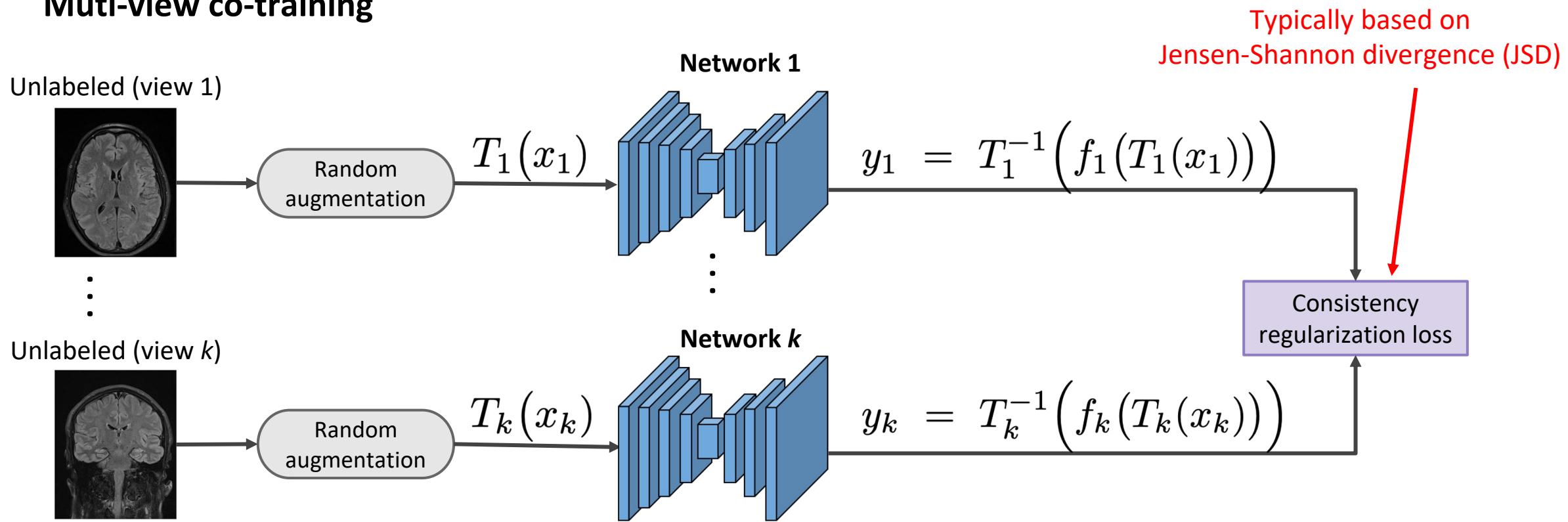


Key idea:

- Supposes the existence of separate, complementary views of the data
- Use high-confidence predictions for a given view as pseudo-labels in other views

SSL methods using consistency regularization

Muti-view co-training

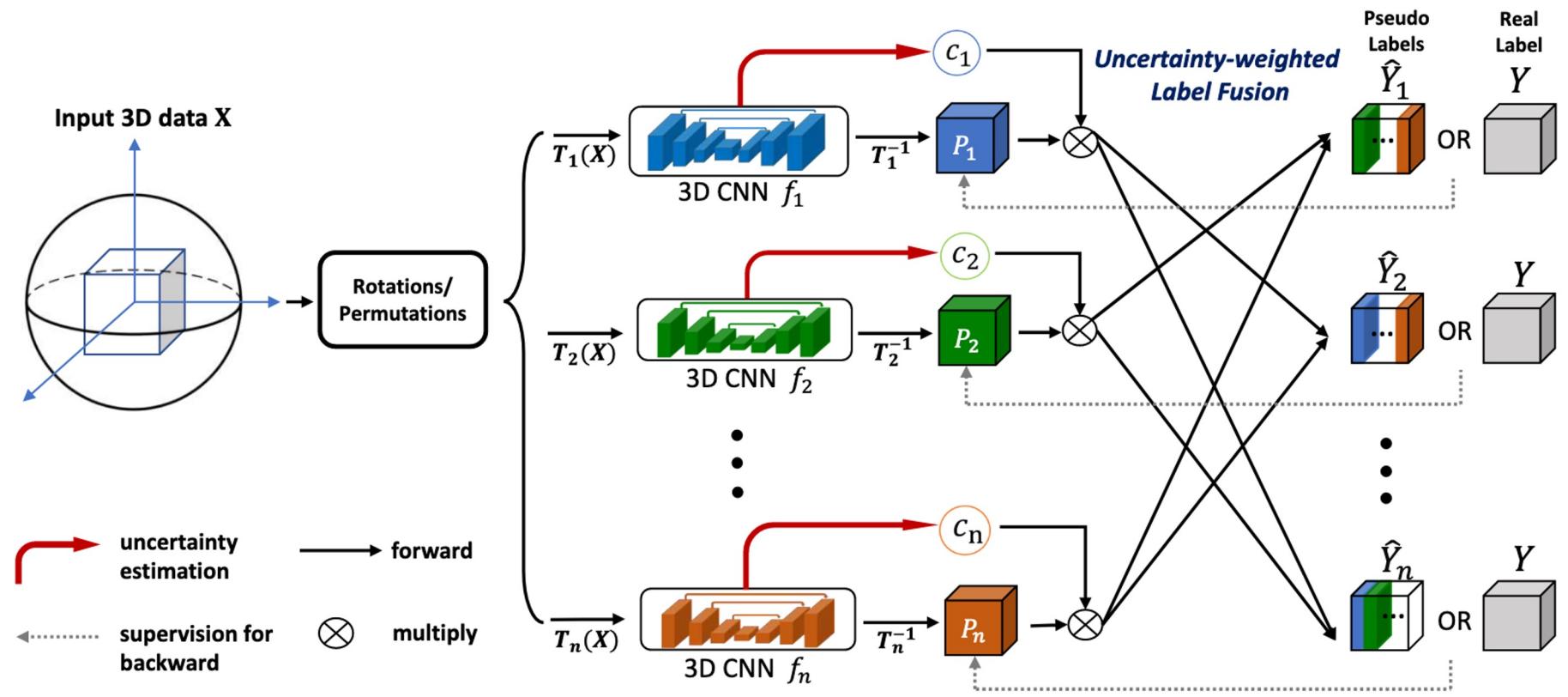


Key idea:

- Supposes the existence of separate, complementary views of the data
- Use high-confidence predictions for a given view as pseudo-labels in other views

SSL methods using consistency regularization

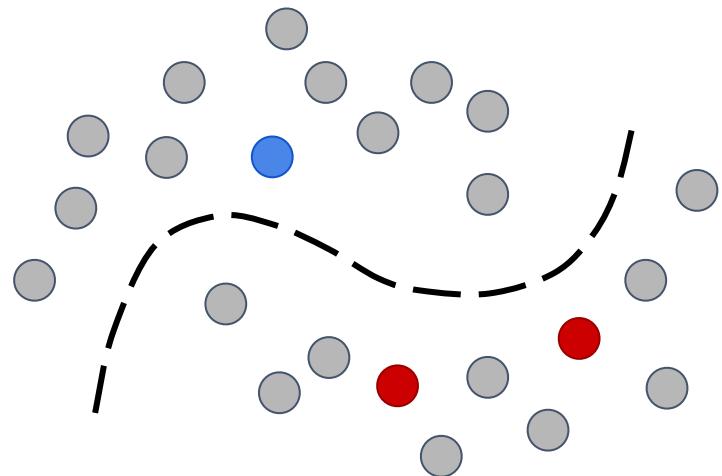
Application of multi-view co-training for pancreas and liver tumor segmentation



Unsupervised representation learning for weakly-supervised segmentation

Unsupervised representation learning (URL)

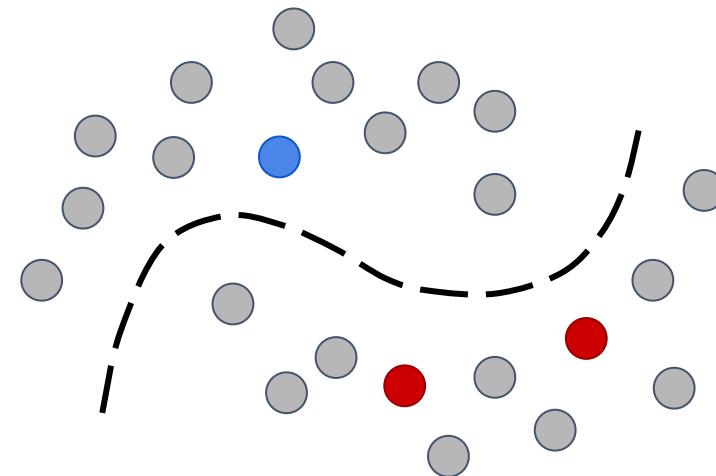
Traditional semi-supervised learning



- Train a model simultaneously with both labeled and unlabeled data

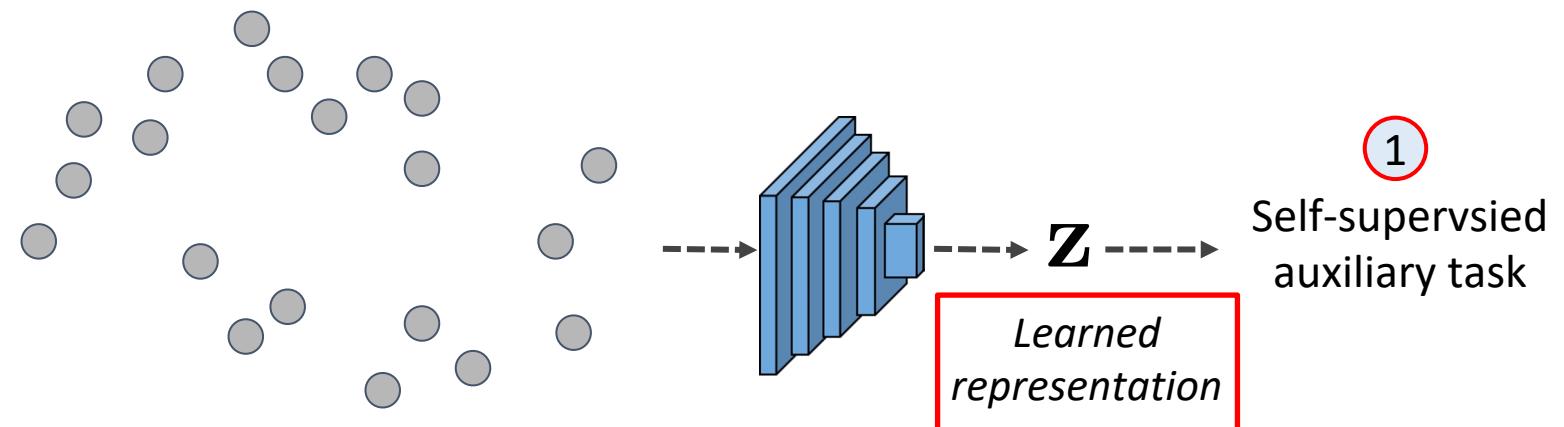
Unsupervised representation learning (URL)

Traditional semi-supervised learning



- Train a model simultaneously with both labeled and unlabeled data

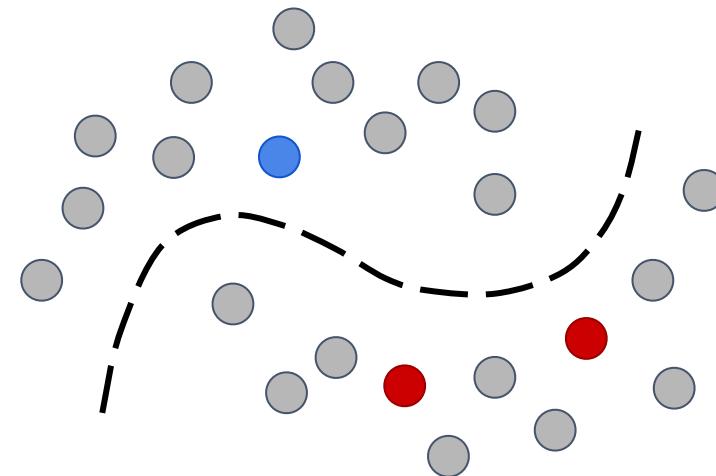
Self-supervised representation learning



- Learn a data representation by solving an auxiliary task that does not require labels
- Use this representation to solve a downstream task
- A light weight fine-tuning of the model can generally be done

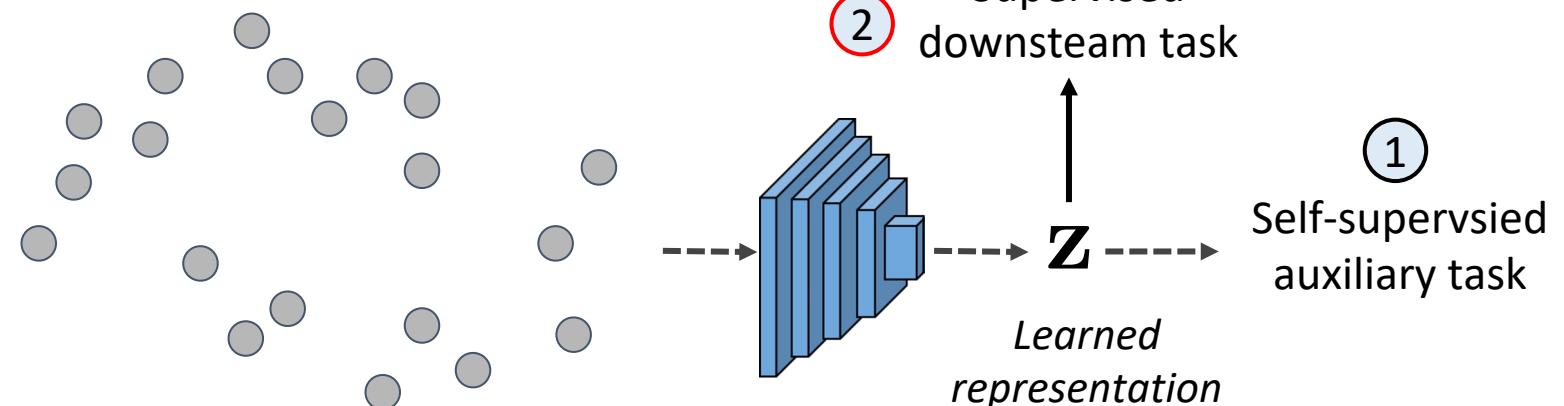
Unsupervised representation learning (URL)

Traditional semi-supervised learning



- Train a model simultaneously with both labeled and unlabeled data

Self-supervised representation learning

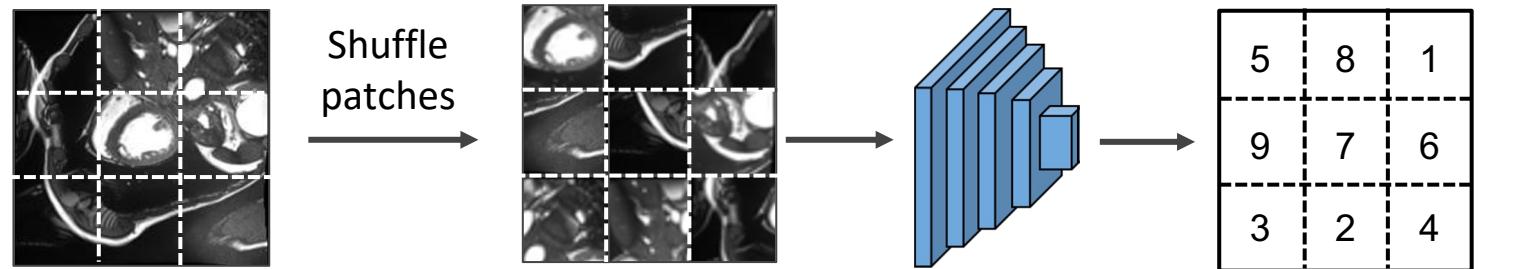


- Learn a data representation by solving an auxiliary task that does not require labels
- Use this representation to solve a downstream task
- A light weight fine-tuning of the model can generally be done

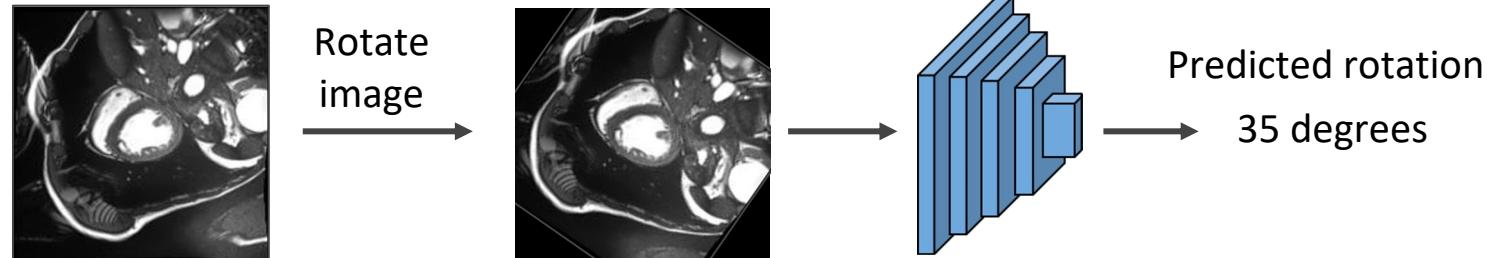
Approaches for URL

Self-supervised learning:

Jigsaw puzzle solving



Rotation prediction

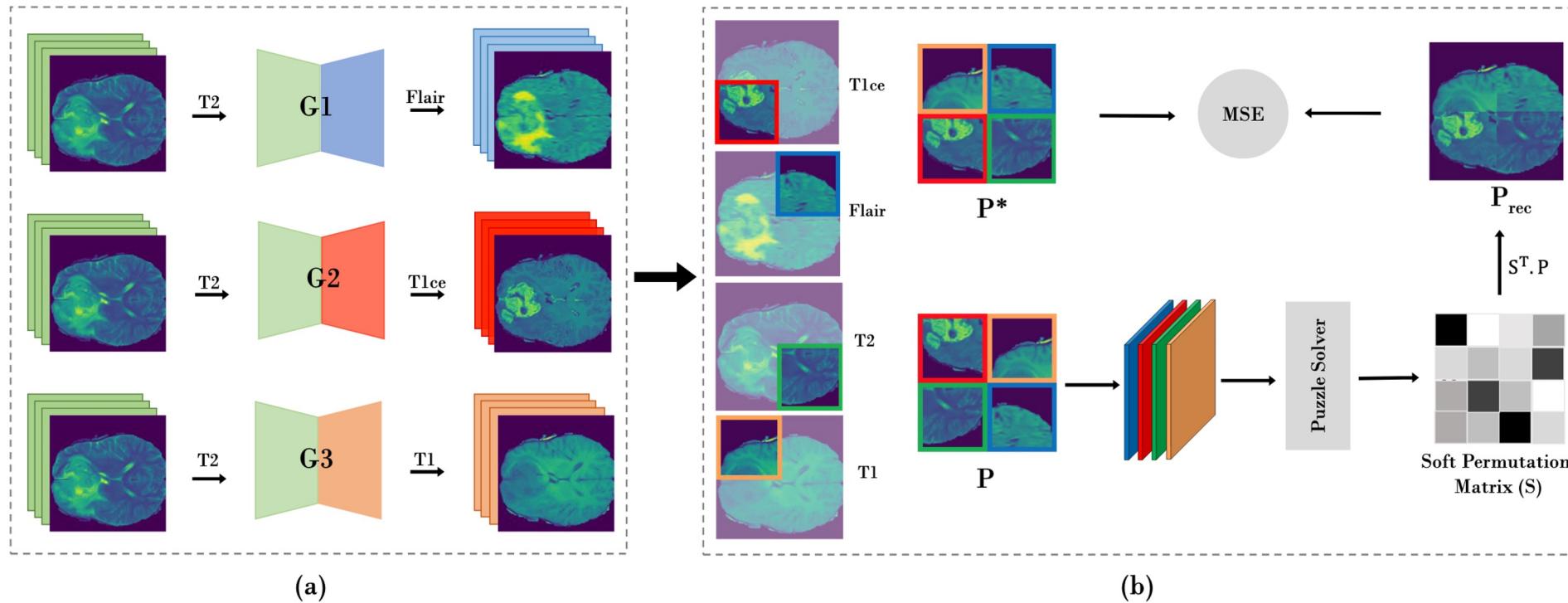


Basic idea:

- Learn to solve a pretext task which does not require annotations
- Example: find the correct order of permuted patches (*see above*)

Approaches for URL

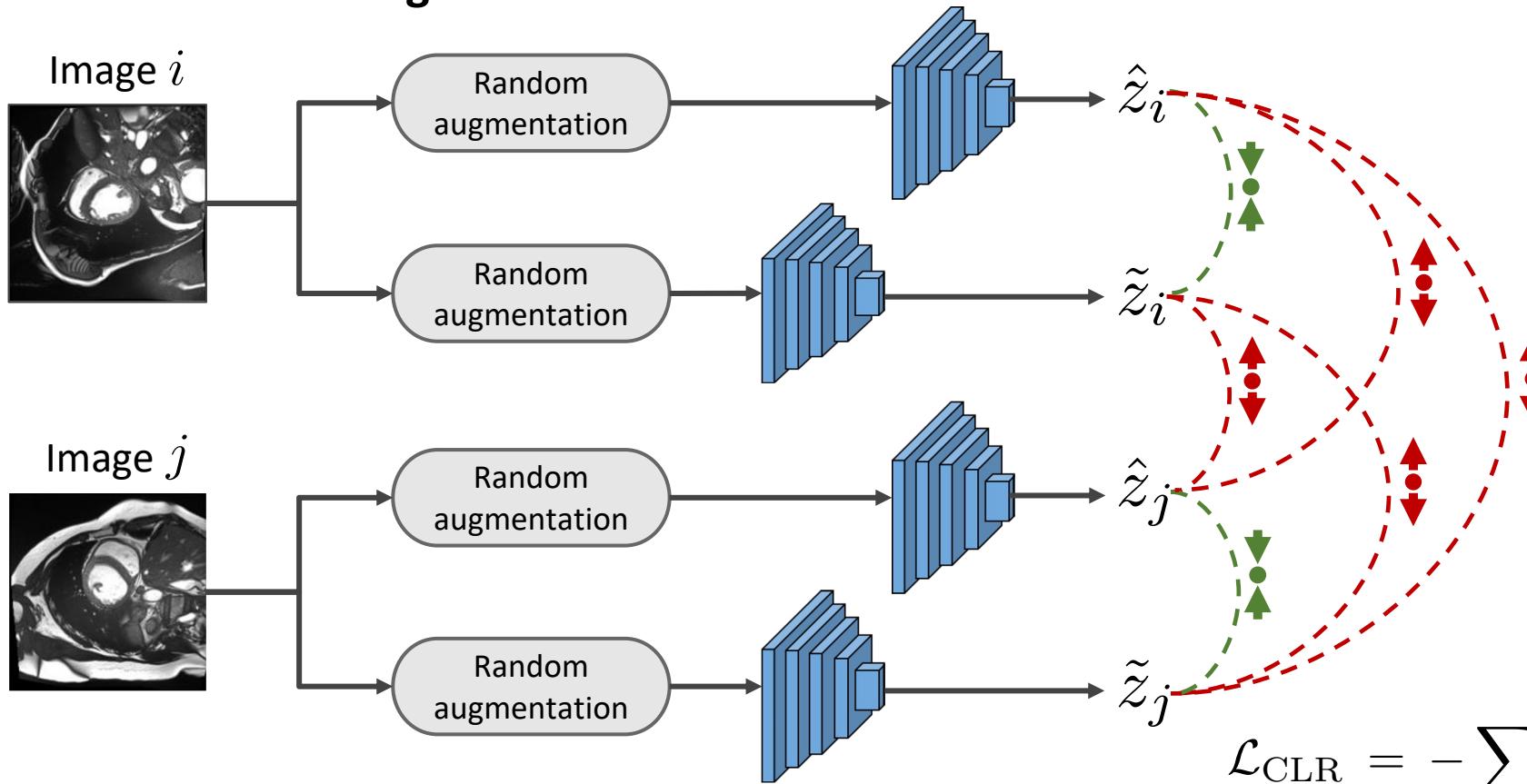
Application to brain tumor MRI



Taleb, A., et al. "Multimodal self-supervised learning for medical image analysis." arXiv preprint (2019).

Self-supervised representation learning

Constrastive learning

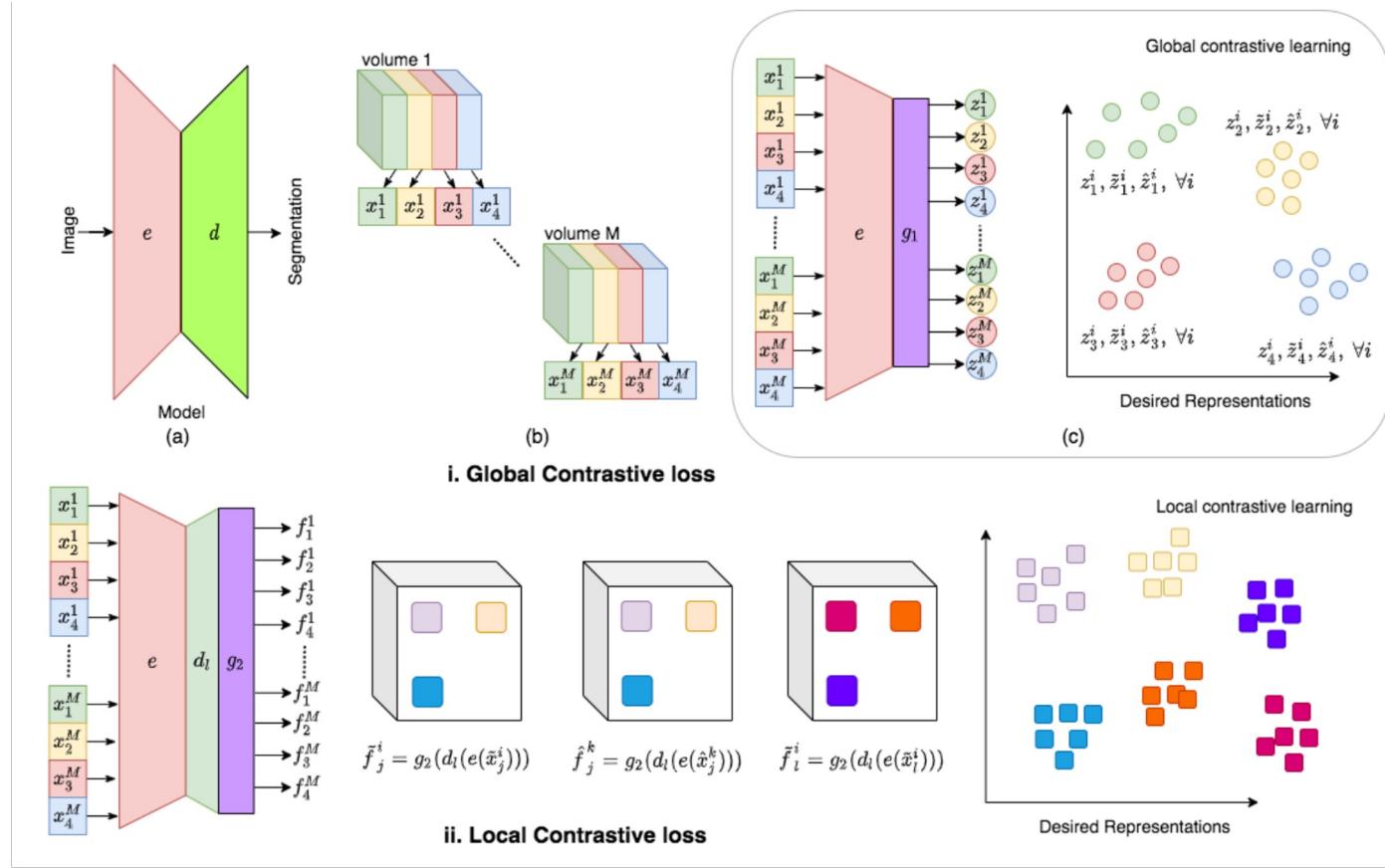


Key idea:

- Generate two augmentations for each image
- Pull closer the representations of the same image and push away those of different ones

Approaches for URL

Contrastive learning:

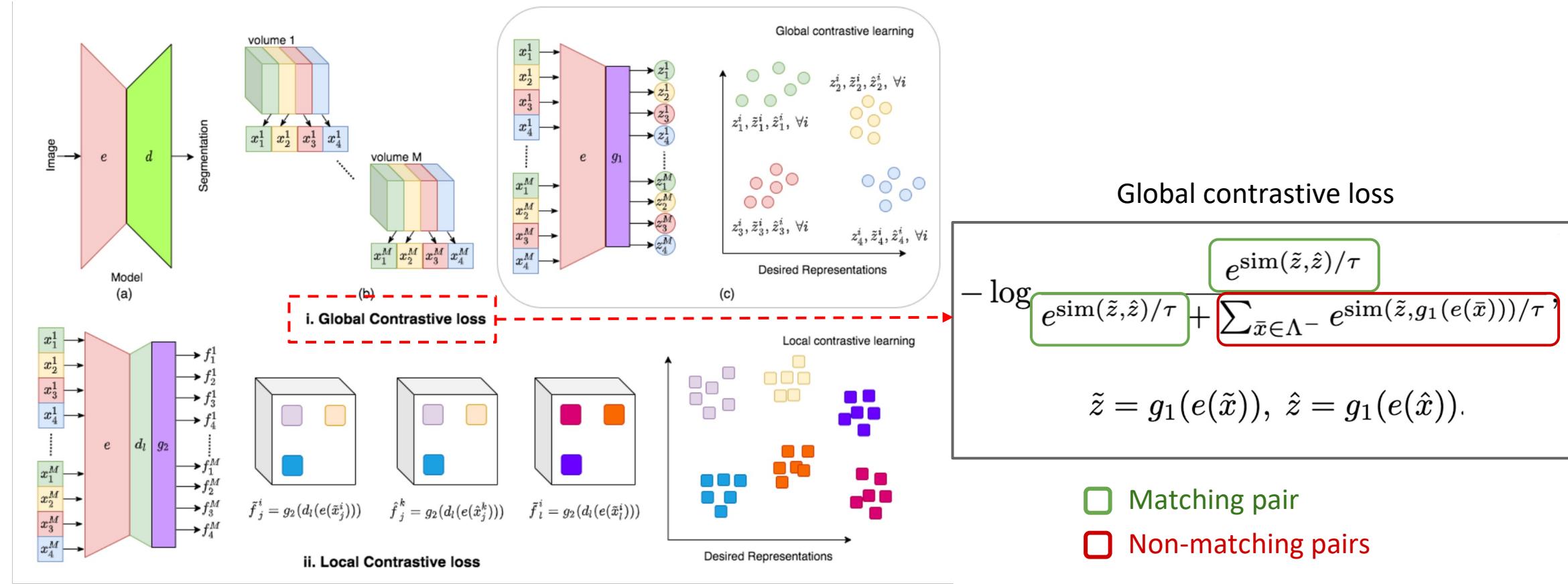


Basic idea:

- Train with pairs of images that match (e.g., same position in volume, same image under different transformations, etc.) or not
- Find a representation that is similar for matching pairs and different for non-matching ones

Approaches for URL

Contrastive learning:

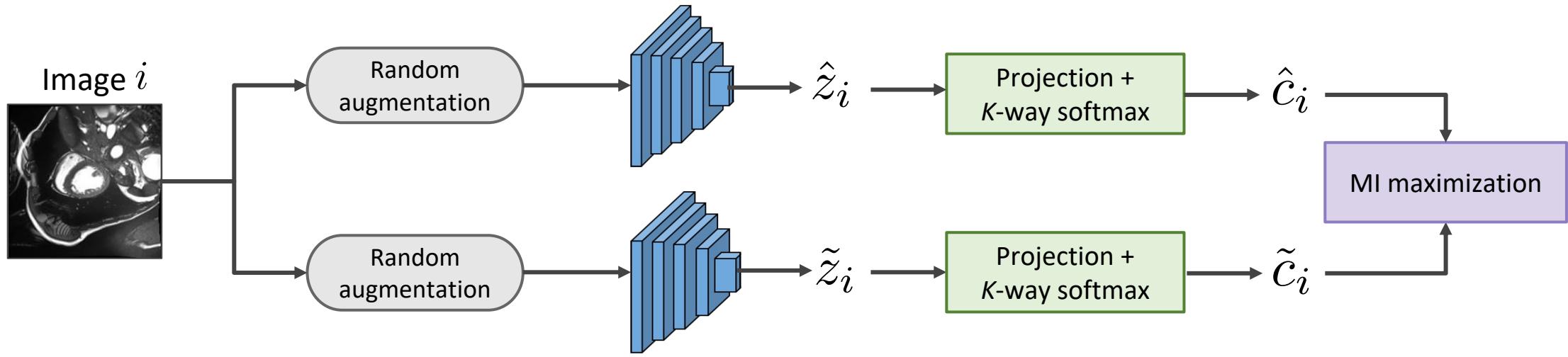


Basic idea:

- Matching pairs for global loss are slices in the same volume position or same subject
- Matching pairs for local loss are feature vectors from the same image under two transformations

Self-supervised representation learning

Clustering based on mutual information (MI)



$$MI = D_{KL}(P(\hat{c}_i, \tilde{c}_i) \parallel P(\hat{c}_i) \cdot P(\tilde{c}_i))$$

Key idea:

- Project features to a discrete K-cluster probability distribution
- Enforce transformation invariance to clusters using MI maximization

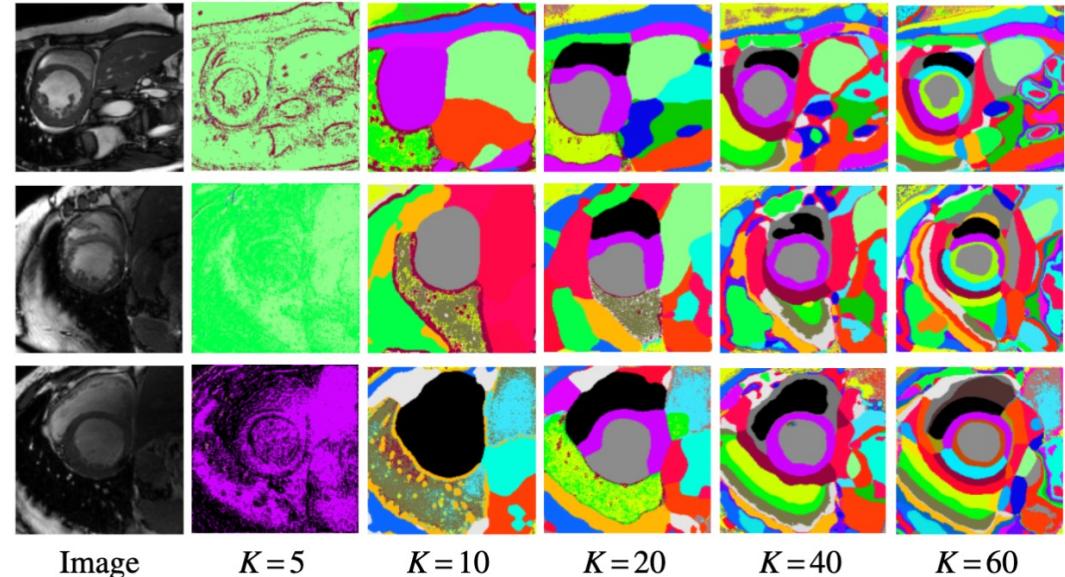
Boundary-aware Information Maximization for Self-supervised Medical Image Segmentation

Jizong Peng*
 ETS Montreal
 jizong.peng.1@etsmtl.net

Ping Wang
 ETS Montreal
 ping.wang.1@ens.etsmtl.ca

Christian Desrosiers
 ETS Montreal
 christian.desrosiers@etsmtl.ca

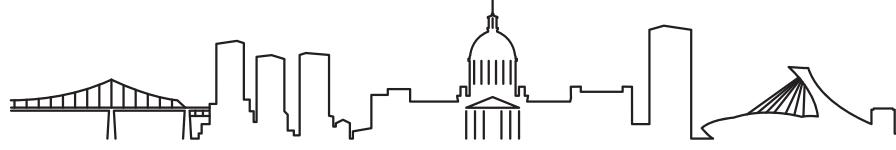
Marco Pedersoli
 ETS Montreal
 marco.pedersoli@etsmtl.ca



Methods	ACDC-LV				ACDC-RV				ACDC-Myo				PROMISE12			
	<i>n</i> =1	<i>n</i> =2	<i>n</i> =4	avg	<i>n</i> =1	<i>n</i> =2	<i>n</i> =4	avg	<i>n</i> =1	<i>n</i> =2	<i>n</i> =4	avg	<i>n</i> =4	<i>n</i> =6	<i>n</i> =8	avg
Partial supervision	67.13	74.49	84.81	75.48	51.82	60.50	64.18	58.84	54.05	67.56	76.00	65.87	49.91	71.53	78.04	66.49
Full supervision	92.26				86.80				88.07				89.65			
Contrast (<i>Enc+Dec</i>)	77.98	85.97	88.42	84.12	66.47	72.82	76.69	71.99	64.96	76.98	78.76	73.57	60.68	77.97	80.53	73.06
Ours (<i>pre-train</i>)	84.48	87.85	90.04	87.45	75.42	79.73	78.89	78.01	74.30	78.43	82.82	78.52	69.76	80.47	82.09	77.44
Entropy minimization	73.79	80.26	86.84	80.30	56.18	62.09	66.27	61.51	57.23	71.10	76.28	68.20	59.78	76.09	78.98	71.62
MixUp	73.30	76.30	84.42	78.01	61.23	63.60	63.14	62.66	55.74	69.80	73.84	66.46	52.09	75.59	81.11	69.60
Mean Teacher (MT)	83.13	87.02	87.70	85.95	61.61	68.76	67.21	65.86	61.55	75.32	78.42	71.76	84.71	85.97	86.93	85.87
UA-MT	81.08	85.03	87.19	84.43	62.06	67.91	66.64	65.54	59.26	73.68	78.61	70.52	66.16	81.79	84.40	77.45
ICT	76.87	78.41	86.34	80.54	60.31	63.42	68.35	64.03	55.91	71.77	77.90	68.53	63.97	77.92	81.39	74.43
Adversarial learning	75.31	74.85	85.85	78.67	55.29	62.25	64.58	60.71	57.68	70.39	75.94	68.00	71.50	78.63	81.35	77.16
MT + Contrast (<i>Enc+Dec</i>)	86.37	89.57	90.40	88.78	75.53	78.42	77.22	77.06	76.11	80.21	82.00	79.44	76.16	82.89	84.85	81.30
MT + Ours (<i>pre-train</i>)	90.25	91.36	91.04	90.88	80.16	81.50	78.97	80.21	78.71	83.33	83.61	81.88	85.64	85.60	88.45	86.56

Take home messages

- Building large training set of labeled examples is not always possible...
- ... but unlabeled data is often available for free
- Semi-supervised methods (e.g., *adversarial learning, consistency regularization, knowledge distillation*) and self-supervised representation learning can boost performance when labeled data is limited
- Similar approaches can be used to adapt models across different domains (e.g., in *unsupervised domain adaptation or test-time adaptation*)
- Not a silver bullet, can be very challenging at times (e.g., adversarial instability)
- Lots of exciting opportunities for future research !!!



Montreal, July 8-12

DLM 2024

Thank you

Any questions ?