



**POLYTECHNIQUE**  
**MONTRÉAL**  
UNIVERSITÉ  
D'INGÉNIERIE

**CR**CHUM

# Summer school on deep learning for medical imaging

Samuel Kadoury Ph.D., ing.  
Polytechnique Montréal / Centre de recherche du CHUM  
Canada Research Chair in medical imaging and assisted interventions

July 6<sup>th</sup>, 2022

# Topics

- 1) Data harmonization and normalization
- 2) Label noise
- 3) Image synthesis
- 4) Deformable registration
- 5) Motion modeling

# Supporting materials

This copy is for personal use only. To order printed copies, contact [reprints@rsna.org](mailto:reprints@rsna.org)

2113

## Deep Learning: A Primer for Radiologists<sup>1</sup>

Gabriel Chartrand, PhD  
Philip M. Cheng, MD, MS  
Engy El-Sherif, PhD, Eng Sci  
Michel Driedzic, PhD  
Simon Turcotte, MD, MSc  
Christopher J. Pui, PhD  
Samuel Kadoury, PhD  
An Tang, MD, MSc

Abbreviations: CNN = convolutional neural network; 3D = three-dimensional; 2D = two-dimensional.

Radiographics 2017;37:e2113-e2113

<https://doi.org/10.1148/radiol.2017170007>

Content Code: [B]

<sup>1</sup>From the Departments of Radiology (G.C., E.V., A.T.) and Biostatistics and Bioinformatics (P.M.C.), Hospital Saint-Luc, Université de Montréal, Hôpital Sainte-Justine, 350 rue Sainte-Justine, Montréal, Québec, Canada H3T 1C5; Cybernetics, Montréal, Québec, Canada (G.C., M.D.); Department of Radiology, Keck School of Medicine, University of Southern California, Los Angeles, Calif (P.M.C.); Montréal Institute for Learning Algorithms, Montréal, Québec, Canada (E.V.); Institut Énergie et Matériaux Polytechnique, Montréal, Québec, Canada (E.V., C.J.P.); and Centre de Recherche du Cancer de Montréal, Montréal, Québec, Canada (S.T.); and Centre d'Imagerie Clinique et Translationnelle, Montréal, Québec, Canada (S.T., S.K., A.T.). Recipient of a Cum Laude award for an outstanding presentation at RSNA 2017.

Received April 3, 2017; revised requested August 23, 2017; accepted August 23, 2017. Address correspondence to A.T. (e-mail: [an.tang@mcgill.ca](mailto:an.tang@mcgill.ca)).  
Supported by the Consortium for Research and Innovation in Medical Imaging (CIRIM), Québec, MITACS-Center Accelerate (IT0536), the recruitment fund of the Centre de Recherche du Centre Hospitalier de l'Université de Montréal, Polytechnique Montréal, and Imagia Cybernetics, AT supported by the Fonds de Recherche du Québec—Santé, the Fonds de Recherche de l'Association des Radiologues du Québec (Gérald Routhier Scholarship-Junior 1 Salary Award 2017).

©RSNA, 2017 • [radiographics.rsna.org](https://doi.org/10.1148/radiol.2017170007)

### INFORMATICS

#### SA-CME LEARNING OBJECTIVES

After completing this journal-based SA-CME activity, participants will be able to:

- Discuss the key concepts underlying deep learning with CNNs.
- Describe emerging applications of deep learning techniques to radiology for lesion classification, detection, and segmentation.
- List key technical requirements in terms of dataset, hardware, and software required to perform deep learning.

See [rsna.org/education/search/rg](https://www.rsna.org/education/search/rg).

#### Introduction

Medical image analysis and interpretation are fundamental cognitive tasks of a diagnostic radiologist. Effective computer automation of these tasks has historically been difficult despite technical advances in computer vision, a discipline dedicated to the problem of imparting visual understanding to a computer system. Recently, however, computer science researchers using a technique called deep learning have demonstrated breakthrough performance improvements in a variety of complex tasks, including image classification, object detection, speech recognition, language translation, natural language processing, and playing games.<sup>1</sup>

Montagnon et al. *Insights into Imaging* (2020) 11:22  
<https://doi.org/10.1186/s13244-019-0832-5>

## Insights into Imaging

EDUCATIONAL REVIEW Open Access

### Deep learning workflow in radiology: a primer

Emmanuel Montagnon<sup>1</sup>, Milena Cerny<sup>1</sup>, Alexandre Cadrian-Chênevert<sup>2</sup>, Vincent Hamilton<sup>1</sup>, Thomas Derennes<sup>1</sup>, André Ilinca<sup>3</sup>, Franck Vandenbroucke-Ménut<sup>4</sup>, Simon Turcotte<sup>1,4</sup>, Samuel Kadoury<sup>2</sup> and An Tang<sup>1,5</sup>

**Abstract**  
Interest for deep learning in radiology has increased tremendously in the past decade due to the high achievable performance for medical image analysis, particularly for image segmentation and prediction. This article provides step-by-step practical guidance for conducting a project that involves deep learning in radiology, from defining specifications, to deployment and scaling. Specifically, the objectives of this article are to provide an overview of clinical use cases of deep learning, describe the composition of multi-disciplinary team, and summarize current approaches to patient, data, model, and hardware selection. Key ideas will be illustrated by examples from a prototypical project on imaging of colorectal liver metastasis. This article illustrates the workflow for image detection, segmentation, classification, monitoring, and prediction of tumor recurrence and patient survival. Challenges are discussed, including ethical considerations, cohorting, data collection, anonymization, and availability of expert annotations. The practical guidance may be adapted to any project that requires automated medical image analysis.

**Keywords:** Review article, Deep learning, Medical imaging, Cohorting, Convolutional neural network

**Pointers**

- Deep learning provides state-of-the-art performance for detection, segmentation, classification, and prediction.
- A multi-disciplinary team with clinical, imaging, and technical expertise is recommended.
- Data collection and curation constitute the most time-consuming step.
- Several open-source deep learning frameworks with permissive licenses are available.
- Cloud computing leverages third-party hardware, storage, and technical resources.

using simpler hierarchized structures defined from a set of specific features. With the advent of powerful parallel computing hardware based on graphical processing units (GPUs) and the availability of large datasets, deep learning has become a state-of-the-art technique in computer vision [1]. In the context of healthcare, deep learning shows great promise for analyzing structured (e.g., databases, tables) and unstructured (e.g., images, text) data [2]. Over the past decade, medical image analysis has greatly benefited from the application of deep learning (DL) techniques to various imaging modalities and organs [3].

Several tasks traditionally performed by radiologists such as lesion detection, segmentation, classification, and monitoring may be automated using deep learning techniques [4]. In abdominal radiology, deep learning has been applied to tasks such as lesion detection and pathology [7–9]. Despite the emerging application of deep learning techniques [1, 10], few articles have described the workflow to execute projects in abdominal radiology which require a broad range of steps, ranging from selection of patient population, choice of index test and reference standard, model selection, and assessment of performance.

<sup>1</sup> Correspondence:  
Centre de Recherche du Centre Hospitalier de l'Université de Montréal (CRCHUM), Montréal, Québec, Canada  
<sup>2</sup>Department of Radiology, Reffo-Oncology and Nuclear Medicine, Université Montréal and CRCHUM, 1058 rue Saint-Denis, Montréal, Québec, H2X 3J4, Canada  
<sup>3</sup>Full list of author information is available at the end of the article

© The Author(s). 2020. **Open Access** This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

**SA-CME LEARNING OBJECTIVES**

After completing this journal-based SA-CME activity, participants will be able to:

- Differentiate among four computer vision tasks using deep learning techniques on radiologic images: classification, detection, semantic segmentation, and instance segmentation.
- Identify building blocks that constitute components of more complex neural network architectures.
- Discuss neural network architectures adapted to different computer vision tasks.

See [rsna.org/learning-center-rg](https://www.rsna.org/learning-center-rg).

This copy is for personal use only. To order printed copies, contact [reprints@rsna.org](mailto:reprints@rsna.org)

1427

## Deep Learning: An Update for Radiologists

Philip M. Cheng, MD, MS  
Emmanuel Montagnon, PhD  
Alexandre Cadrian-Chênevert, MD, PhD  
Ian Pan, MD  
Alexandre Cadrian-Chênevert,  
B.Ing, MD  
Francisco Pardigon Romero, MD  
Gabriel Chartrand, PhD  
Samuel Kadoury, PhD  
An Tang, MD, MSc

Abbreviations: CNN = convolutional neural network; GAN = generative adversarial network; R-CNN = regions with CNN features; ROC = receiver operating characteristic.

Radiographics 2021;41:e1427-e1445

<https://doi.org/10.1148/radiol.2021200210>

Content Code: [A] [B]

<sup>1</sup>From the Department of Radiology, Keck School of Medicine of the University of Southern California, Los Angeles, Calif (P.M.C.); Research Center for Radiobiology and Radiation Oncology, Department of Radiology (A.T.), Centre Hospitalier de l'Université de Montréal, Hôpital Sainte-Justine, 350 rue Saint-Denis, Montréal, QC, H3T 1C5; HEX 94, Department of Diagnostic Data Sciences, Brigham and Women's Hospital, Harvard Medical School, Boston, Calif (I.P.); West Asia Apert Radiology School, Beirut University, Beirut, Lebanon (A.C.-C.); Department of Medical Imaging, CISSS Laval, Université Laval, Québec, Québec, Canada (F.P.R.); and APM Medical, Montréal, Québec, Canada (G.C.). Presented as an education exercise at the Radiological Society of North America (RSNA) meeting, October 25, 2020; revised requested April 14, 2021, and received May 2; accepted May 7. ©RSNA, 2021 • [radiographics.rsna.org](https://doi.org/10.1148/radiol.2021200210)

<sup>2</sup>Current address: Department of Radiology, Brigham and Women's Hospital, Boston, Mass.

©RSNA, 2021

### INFORMATICS

#### SA-CME LEARNING OBJECTIVES

After completing this journal-based SA-CME activity, participants will be able to:

- Differentiate among four computer vision tasks using deep learning techniques on radiologic images: classification, detection, semantic segmentation, and instance segmentation.
- Identify building blocks that constitute components of more complex neural network architectures.
- Discuss neural network architectures adapted to different computer vision tasks.

See [rsna.org/learning-center-rg](https://www.rsna.org/learning-center-rg).

# Deep learning: A primer for radiologists

2113

INFORMATICS

## Deep Learning: A Primer for Radiologists<sup>1</sup>

Gabriel Chartrand, PhD  
Philip M. Cheng, MD, MS  
Eugene Vöröntsov, BAsc Eng Sci  
Michal Drozdzal, PhD  
Simon Turcot, MD, MSc  
Christopher J. Pal, PhD  
Samuel Kadoury, PhD  
An Tang, MD, MSc

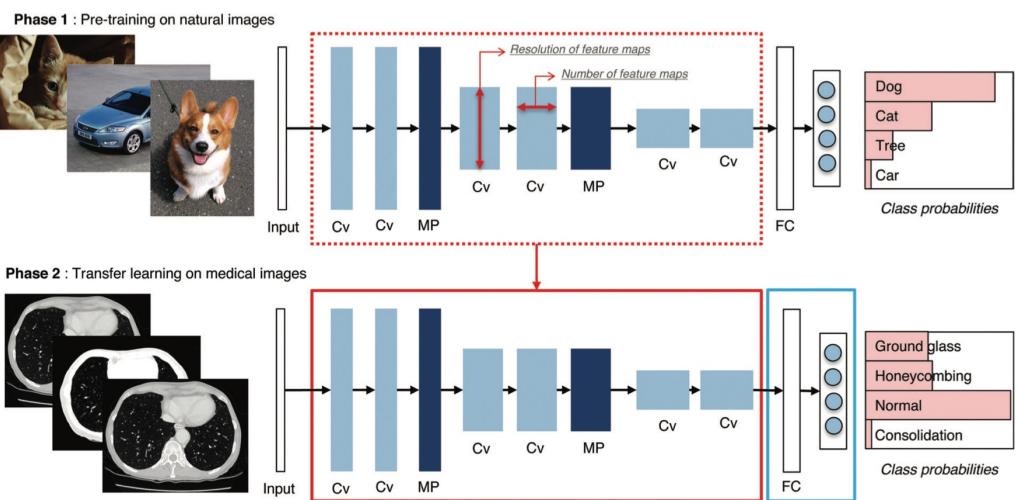
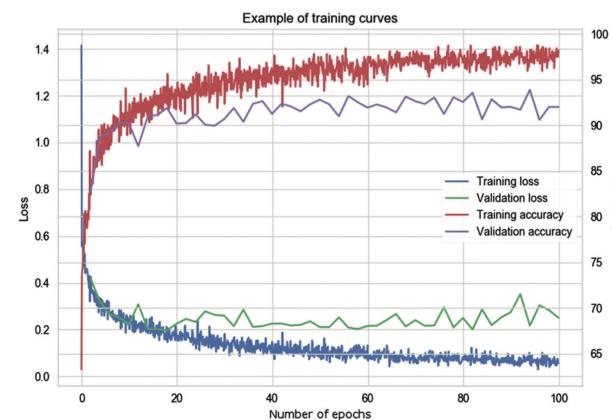
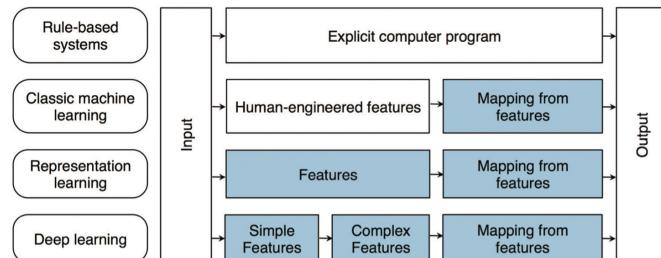
Deep learning is a class of machine learning methods that are gaining success and attracting interest in many domains, including computer vision, speech recognition, natural language processing, and playing games. Deep learning methods produce a mapping from raw inputs to desired outputs (eg, image classes). Unlike traditional machine learning methods, which require hand-engineered feature extraction from inputs, deep learning methods learn these features

### SA-CME LEARNING OBJECTIVES

After completing this journal-based SA-CME activity, participants will be able to:

- Discuss the key concepts underlying deep learning with CNNs.
- Describe emerging applications of deep learning techniques to radiology for lesion classification, detection, and segmentation.
- List key technical requirements in terms of dataset, hardware, and software required to perform deep learning.

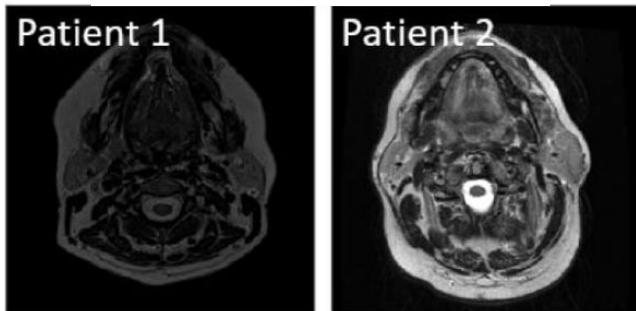
See [www.rsna.org/education/search/RG](http://www.rsna.org/education/search/RG).



# 1) Data harmonization and normalization

# Effect of normalization

## No normalization

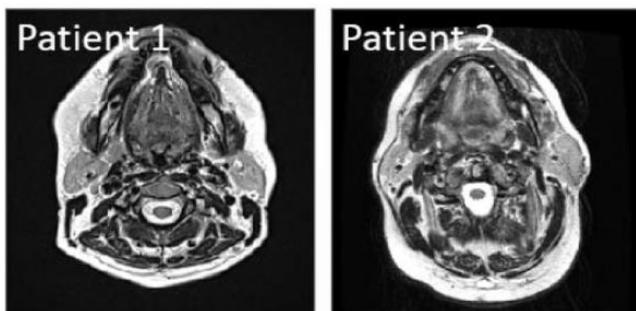


Quantitative  
image analysis

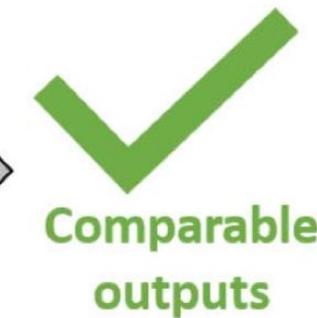


Incomparable  
outputs

## Normalization



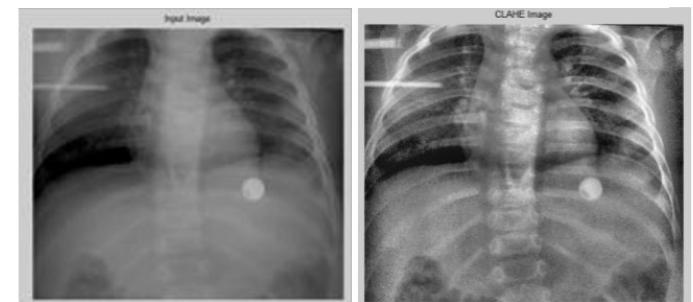
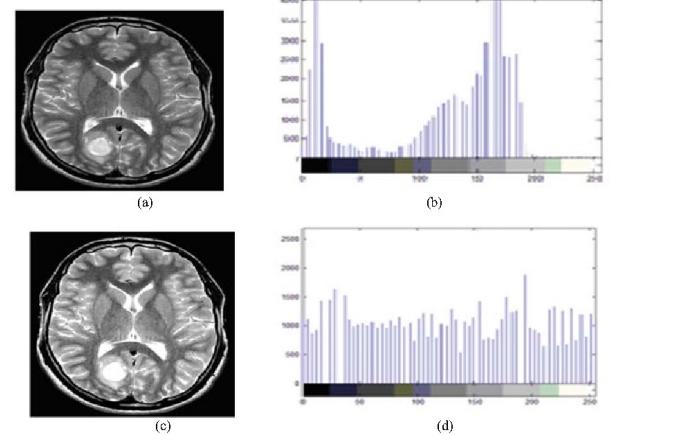
Quantitative  
image analysis



Comparable  
outputs

# Image pre-processing techniques

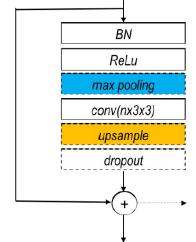
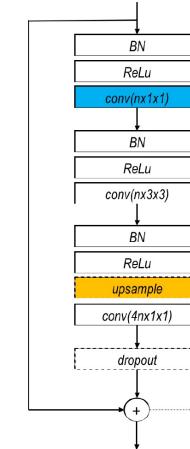
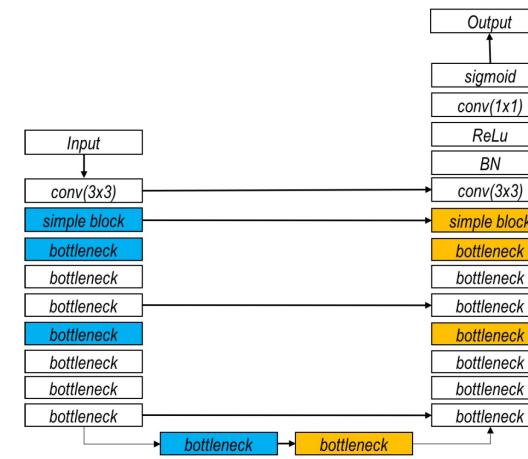
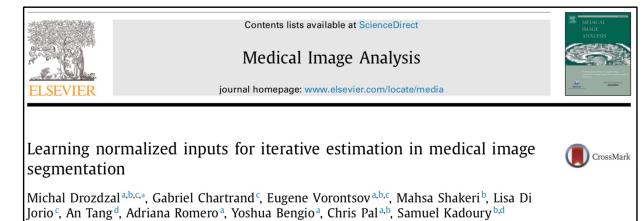
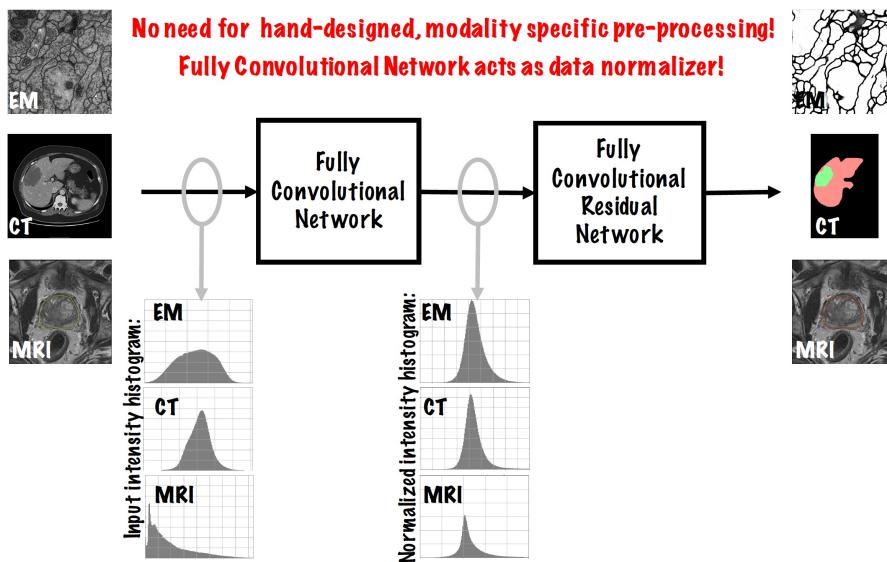
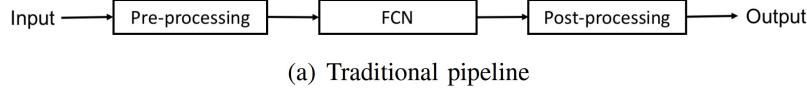
- Histogram equalization
  - Hist. eq. + Gaussian blur
  - Hist. eq. + bilateral filter
- 
- Adaptive masking
    - $(\text{max} - \text{min})$  followed by binary thresholding
  - Adaptive masking + hist. eq. + Gaussian blur



# Normalization techniques for deep learning

- **Min-max normalization:** in this method, pixel gray level values are brought into a specific range (usually [0-1]). Alternatively, a [10<sup>th</sup> percentile, 90<sup>th</sup> percentile] range is chosen in some cases to avoid imaging outliers.
- **Volume standardization:** this technique consists of computing the mean and standard deviation of the entire 3D image. Each voxel is then subtracted by the volume mean and divided by the voxel standard deviation.
- **Slice-by-slice standardization:** this strategy is similar to the previous one. However, instead of computing a global mean and standard deviation for the entire volume, slice-specific mean and standard deviation are computed, and pixels are standardized using the slice statistics instead of the volume statistics.

# Learning normalized input for image segmentation



**Data:**

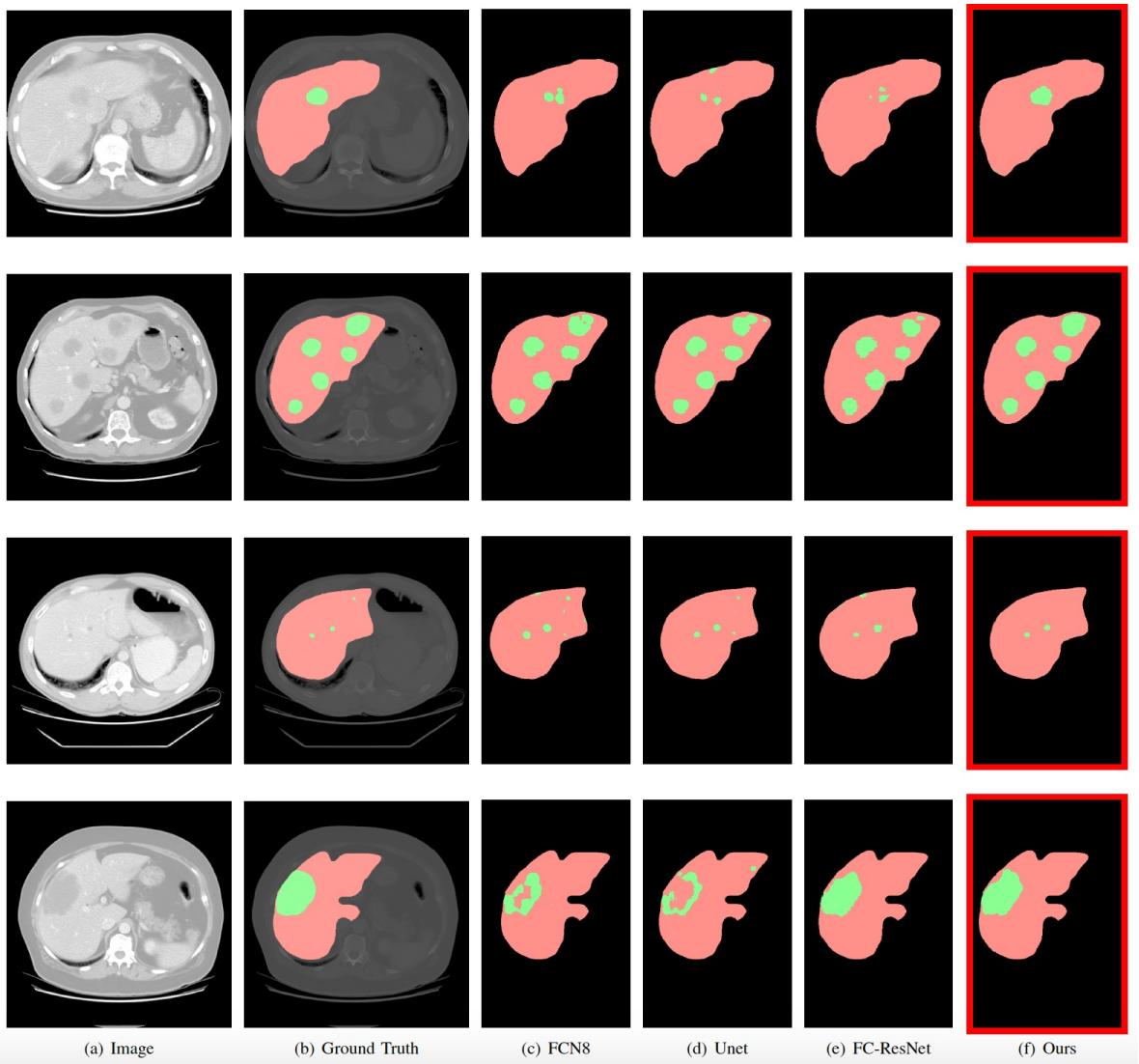
> CT data (512 x 512 x 200)

> 105 images

- 77 training

- 28 validation

> 30 test



Method	Validation		
	loss	$Dice_{lesion}$	$Dice_{liver}$
FCN8 [36]	-0.419	0.589	0.994
Unet [43]	-0.451	0.553	0.994
FC-ResNet [20]	-0.223	0.551	0.993
Ours	<b>-0.795</b>	<b>0.771</b>	<b>0.997</b>

Method	Test		
	loss	$Dice_{lesion}$	$Dice_{liver}$
FCN8 [36]	-0.437	0.535	0.989
Unet [43]	-0.396	0.570	0.990
FC-ResNet [20]	-0.224	0.617	0.990
Ours	<b>-0.796</b>	<b>0.711</b>	<b>0.993</b>

# Prostate segmentation

## PROMISE12 challenge data:

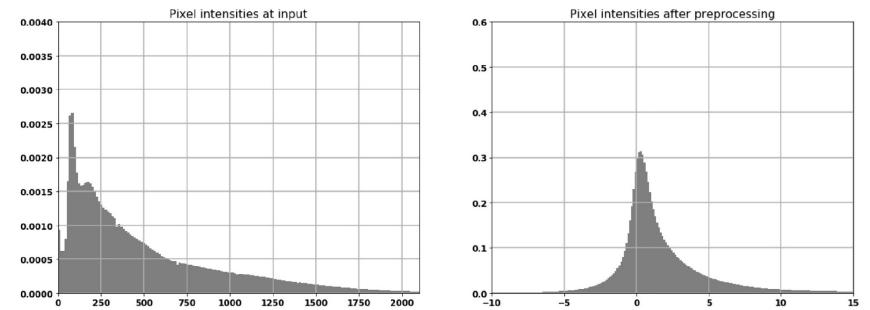
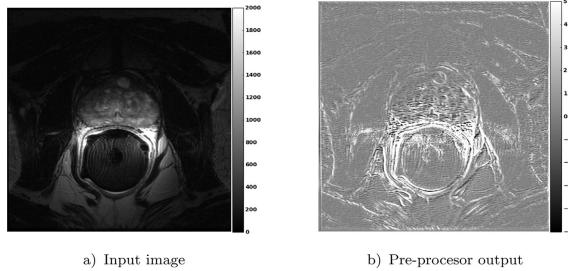
> T2-w MRI (256 x 256 x 30)

> 50 images

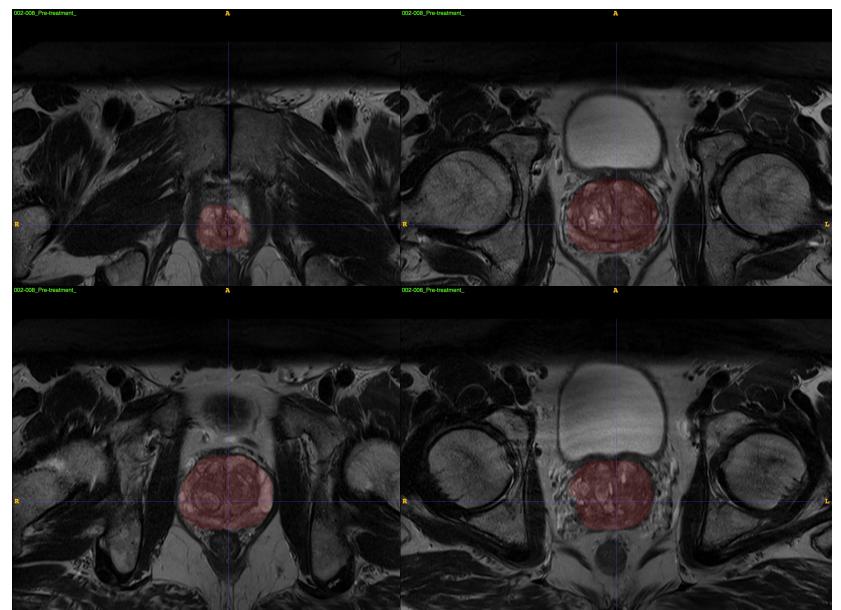
- 40 training

- 10 validation

> 30 test images



Test set (n=7) Princess Margaret Hospital



Method	Score [-]	Dice [%]	Avg. Dist. [mm]	Vol. Diff. [%]
--------	--------------	-------------	--------------------	-------------------

2D FCNs

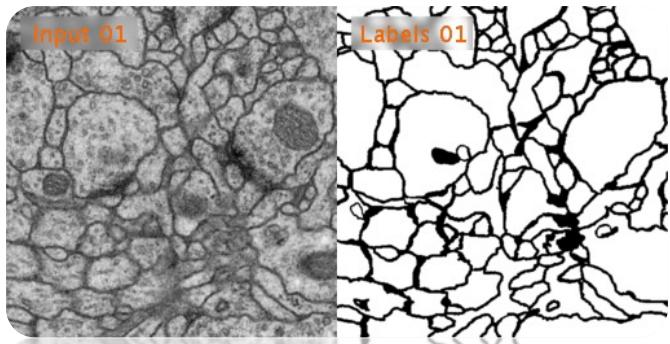
Ours	<b>83.02</b>	87.4	2.17	12.37
SITUS	79.92	84.13	2.96	23.00

3D FCNs

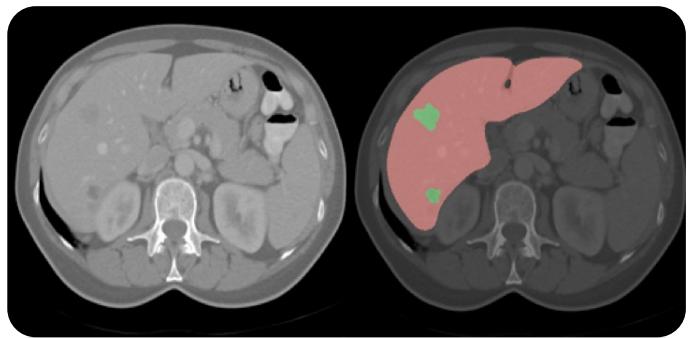
CUMED [50]	86.65	89.43	1.95	6.95
CAMP-TUM2 [35]	82.39	86.91	2.23	14.98
SRIBHME	74.17	74.46	2.83	34.89

# Semantic image segmentation

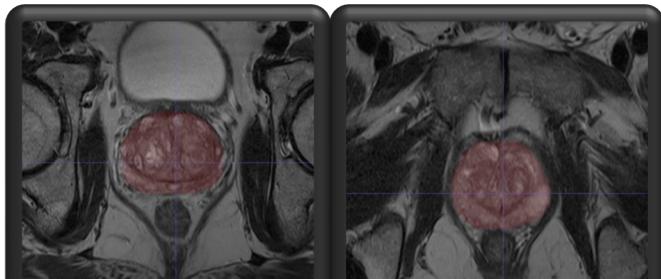
Electron microscopy



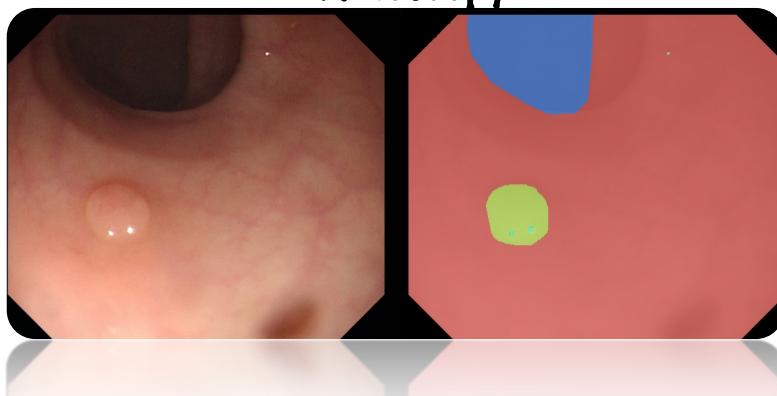
CT



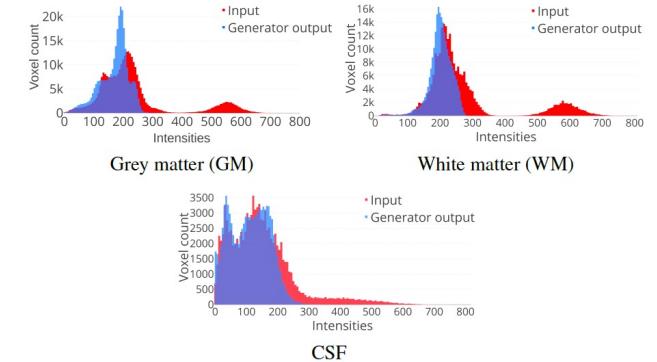
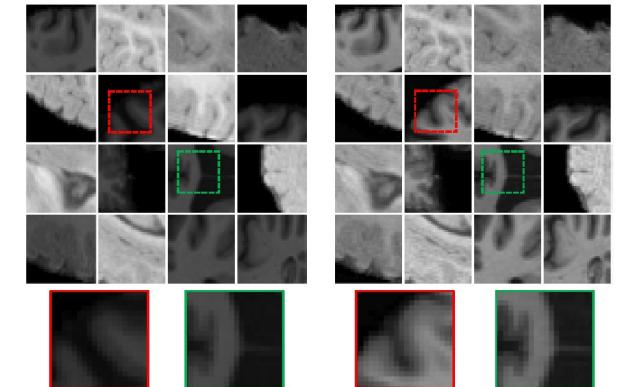
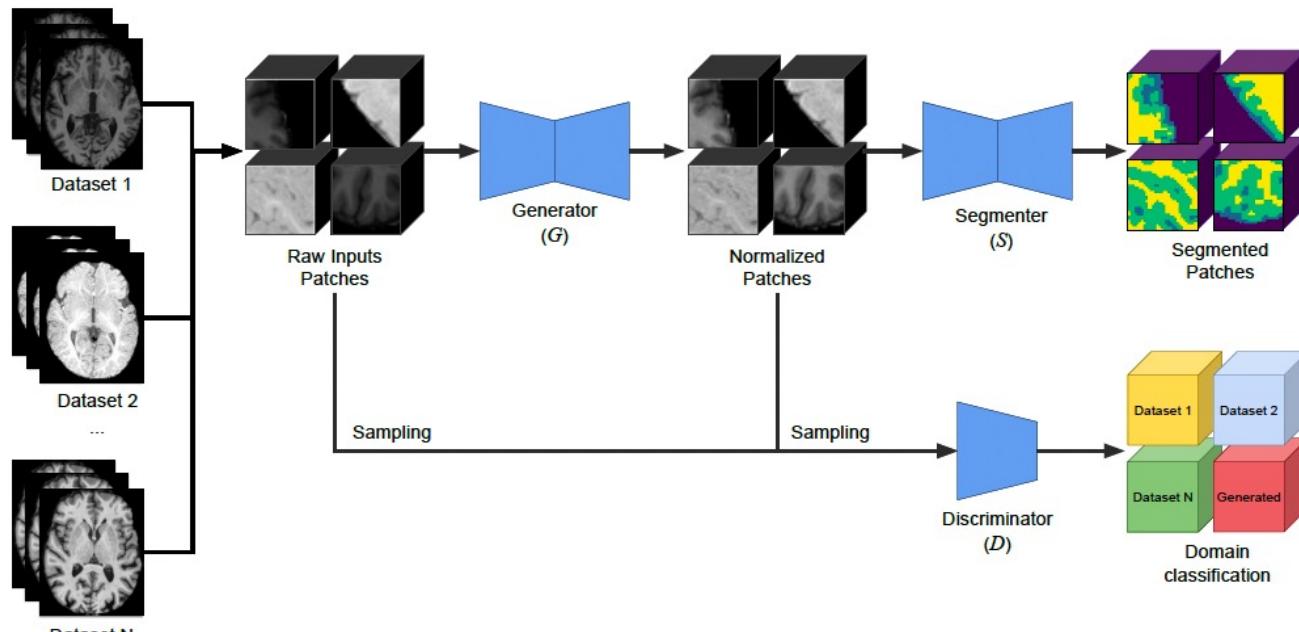
MRI



Endoscopy

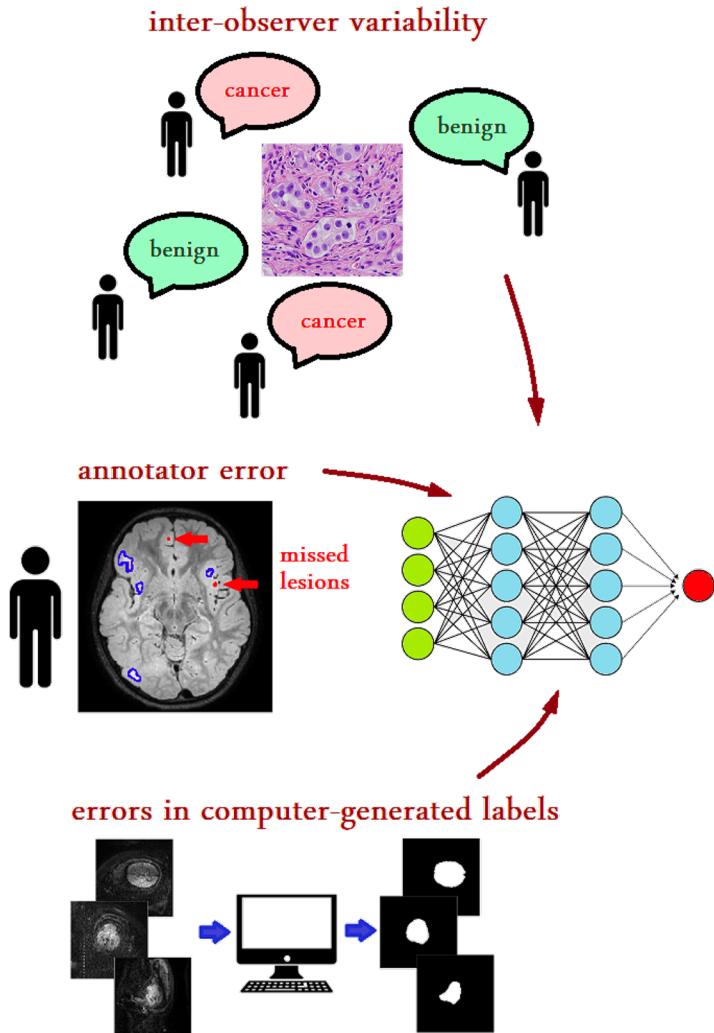


# Realistic Image Normalization for Multi-Domain Segmentation



From Desile et al. Medical Image analysis, 2020

## 2) Label noise



#### Label cleaning and pre-processing

- Identifying and removing samples with incorrect labels
- Iterative label cleaning in parallel with training
- Label smoothing
- Label cleaning with generative models

#### Network architectures

- Estimating and utilizing noise statistics using
  - Noise layers
  - Probabilistic graphical models

#### Label-noise-robust loss functions

- Mean-absolute-error loss
- Cross entropy with abstention
- "Correcting" standard loss functions using label noise statistics

#### Data re-weighting

- Meta-learning
- Gradient scaling
- Learned re-weighting
- Learning to weight labels from different annotators

#### Data and label consistency

- Exploiting data similarity to identify incorrect labels
- Auxiliary image regularization
- Manifold regularization

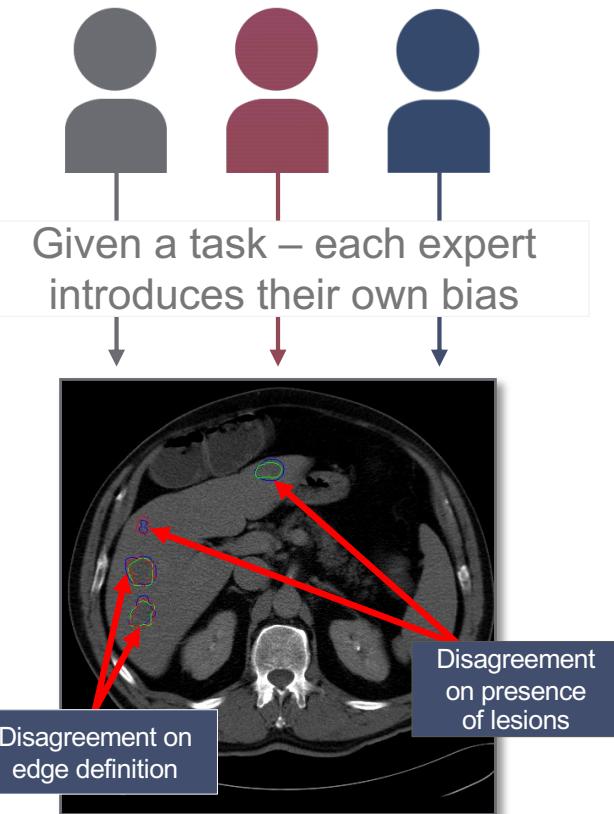
#### Training procedures

- Curriculum learning
- Knowledge distillation
- Co-teaching
- Data augmentation (mixup)

From Karimi et al. 2020

# Sources of error

## Inter-Rater\* Variability

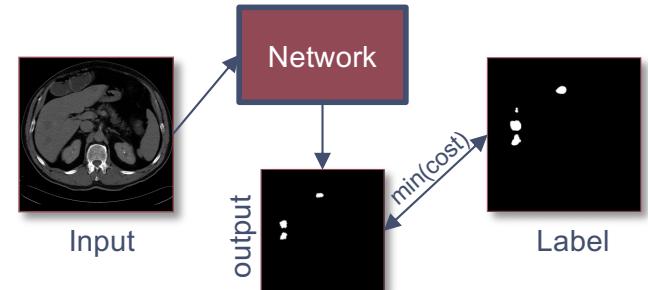


## Lack of Ground Truth

- Lack of agreement between experts results in a lack of availability of a single “ground truth” in the medical field
- Highly biased towards the single rater
- **Impact:**
  - limits the performance of supervised deep learning models
  - Poor generalizability due to biased input data
- No consensus on who is right
- Variability cannot be represented in traditional segmentation approaches

## Segmentation

- Deep learning segmentation models are mostly trained for a single output – a deterministic viewpoint based on a single mask/label

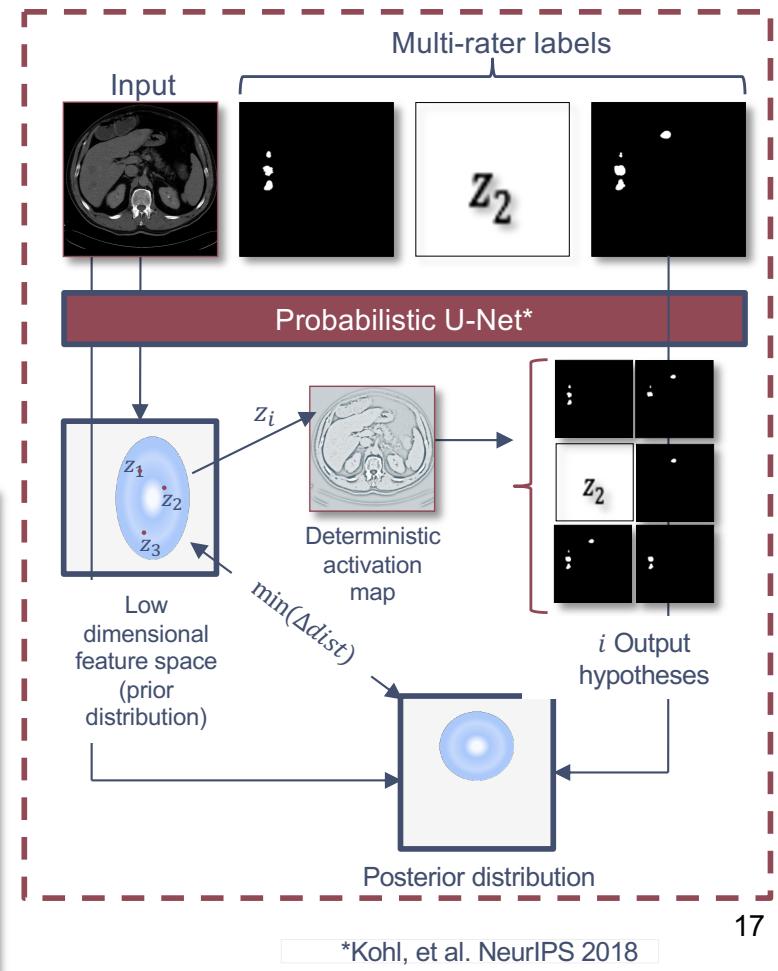
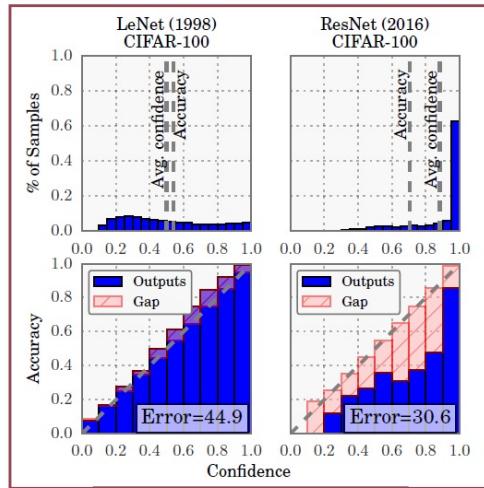
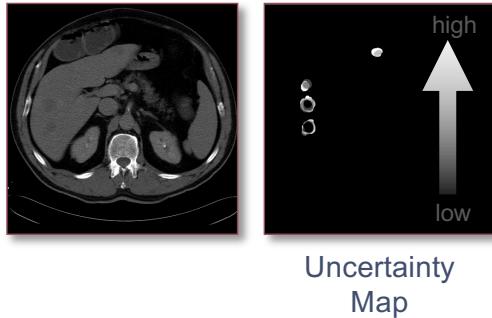


- + Deterministic → output is consistent post optimization
- Does not capture variability of opinions

\*rater = annotator = expert = observer

# How to work with noisy labels? Probabilistic approach

- Take a probabilistic approach to represent a distribution of hypothetical outputs (e.g. probabilistic U-Net)
  - A combination of a deterministic U-Net with a **conditional Variational Auto-Encoder (cVAE)**: generative, probabilistic
  - **cVAE** assumes data is generated from a random process → able to generate multiple hypotheses.
- Quantify the **uncertainty** in the output – i.e. how confident is the model with its prediction?
- Based on this confidence, verify whether the model is **calibrated** (Guo, et al)– i.e. does the accuracy match model confidence? Is the model overconfident in regions where it has low accuracy?
  - the latter is an issue to be addressed to be able to trust predictions from DL models



# Label cleaning and pre-processing

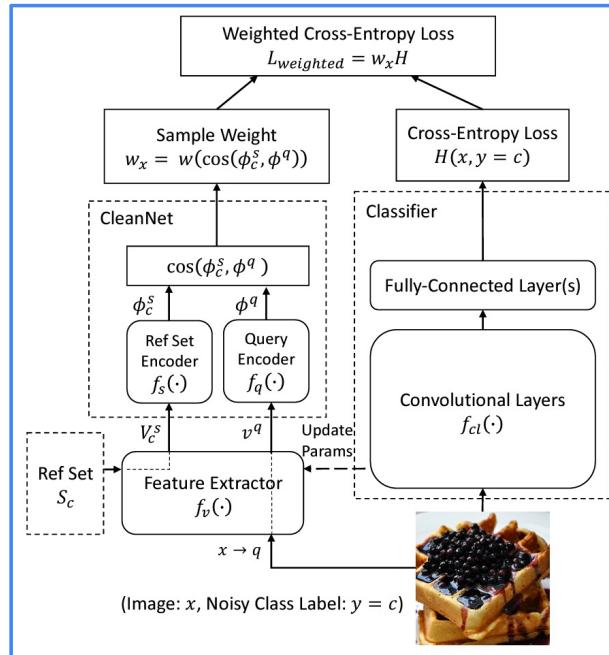
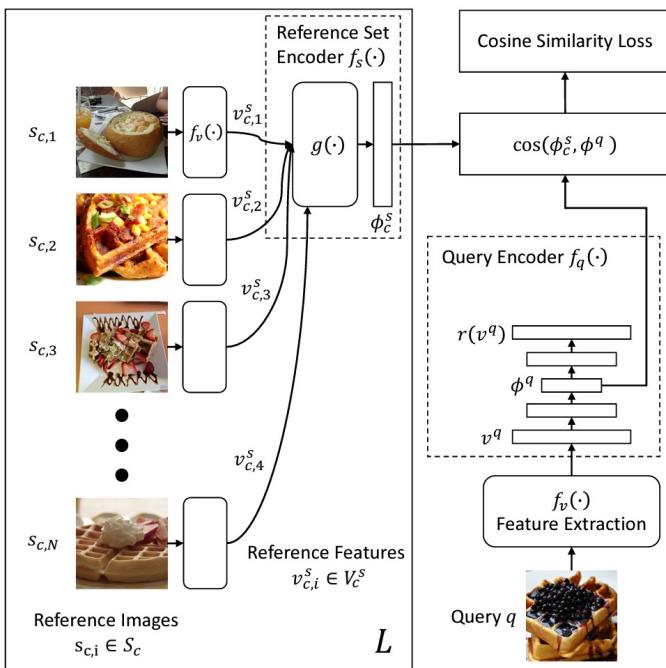
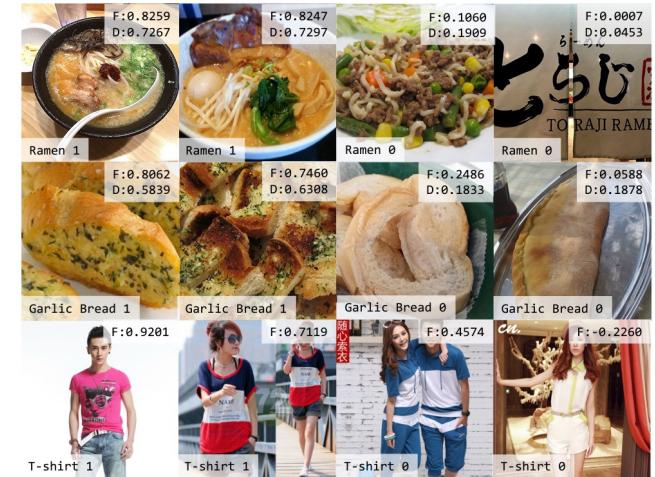


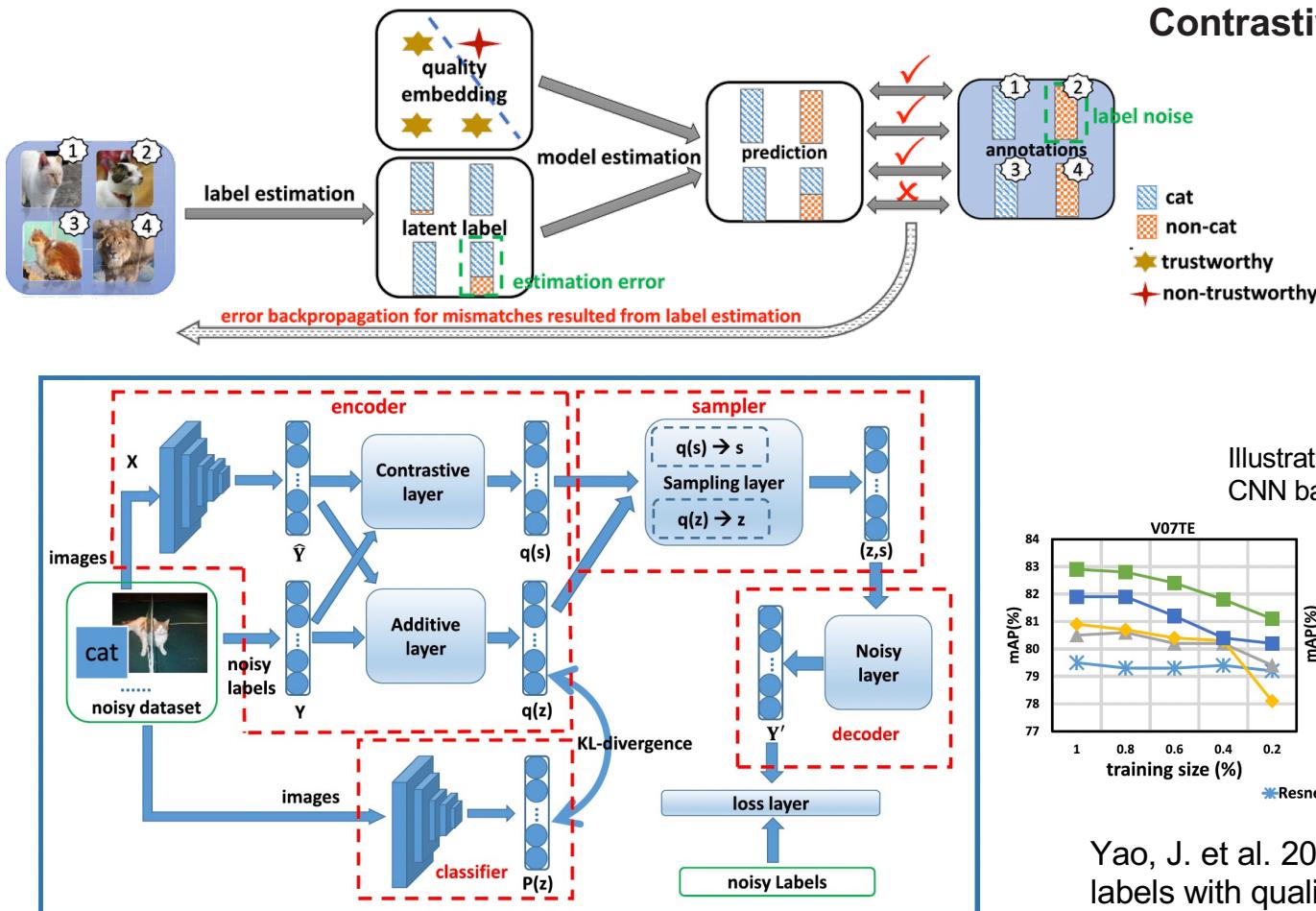
Illustration of integrating CleanNet for training the CNN based image classifier with label noise



F -> cosine similarity by model w/ verification labels  
D -> cosine similarity under transfer learning

Lee et al. "Cleannet: Transfer learning for scalable image classifier training with label noise", CVPR 2018.

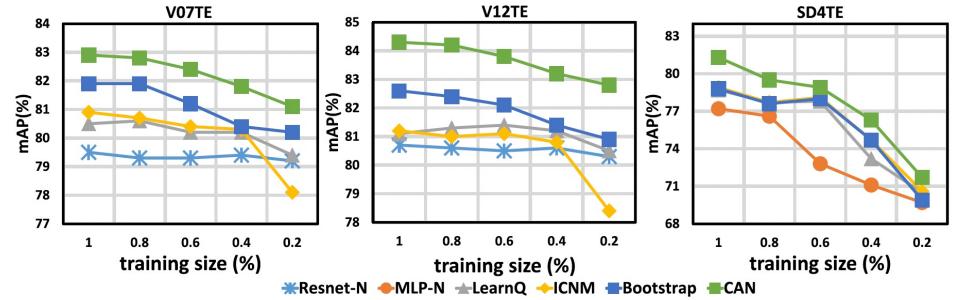
# Network architecture



Contrastive-additive noise network (CAN)

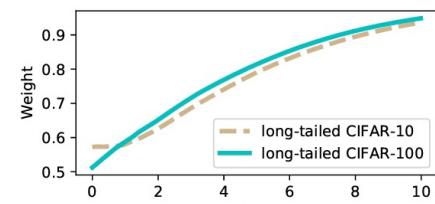
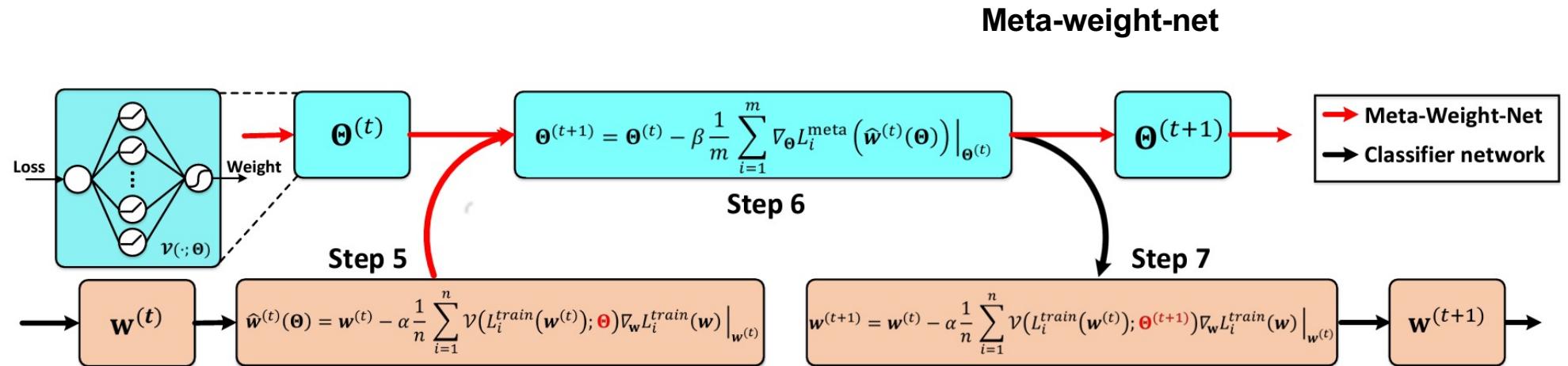
	V07TE				
$P_{noise}$	Resnet-N	LearnQ	ICNM	Bootstrap	CAN
1.0	6.4	9.1	<b>9.2</b>	8.9	8.6
0.8	33.4	28.0	28.5	30.1	<b>36.1</b>
0.6	53.0	56.4	57	59.3	<b>63.2</b>
0.4	70.2	72.0	71.6	73.3	<b>79.4</b>
0.2	78.2	80.1	79.6	81.0	<b>83.6</b>
0.0	<b>86.8</b>	85.4	85.4	85.5	85.3

Illustration of integrating CleanNet for training the CNN based image classifier with label noise

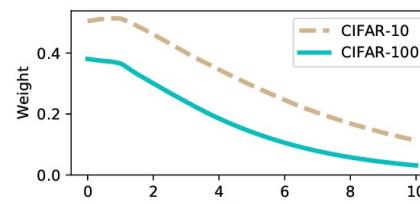


Yao, J. et al. 2018. Deep learning from noisy image labels with quality embedding. IEEE TIP.

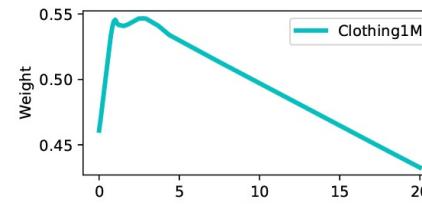
# Data re-weighting



(d) MW-Net function learned in class imbalance case



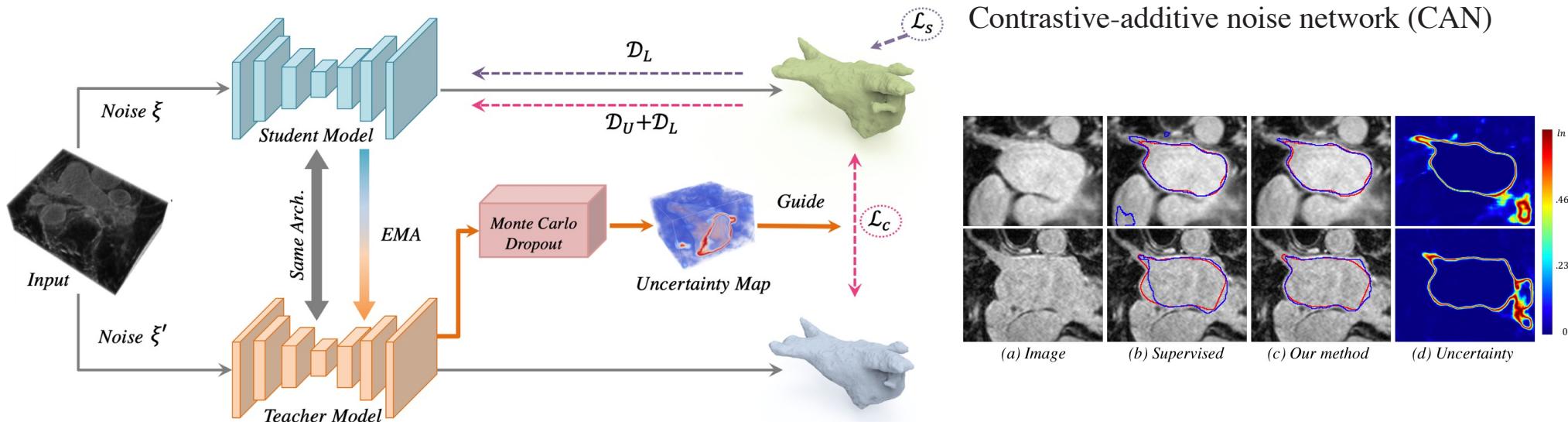
(e) MW-Net function learned in corrupter labels case



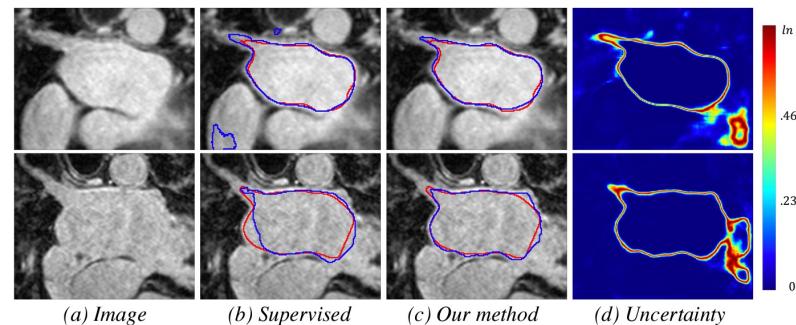
(f) MW-Net function learned in real Clothing1M dataset

Shu, J. , et al. , 2019. Meta-weight-net: learning an explicit mapping for sample weighting.

# Data and label consistency

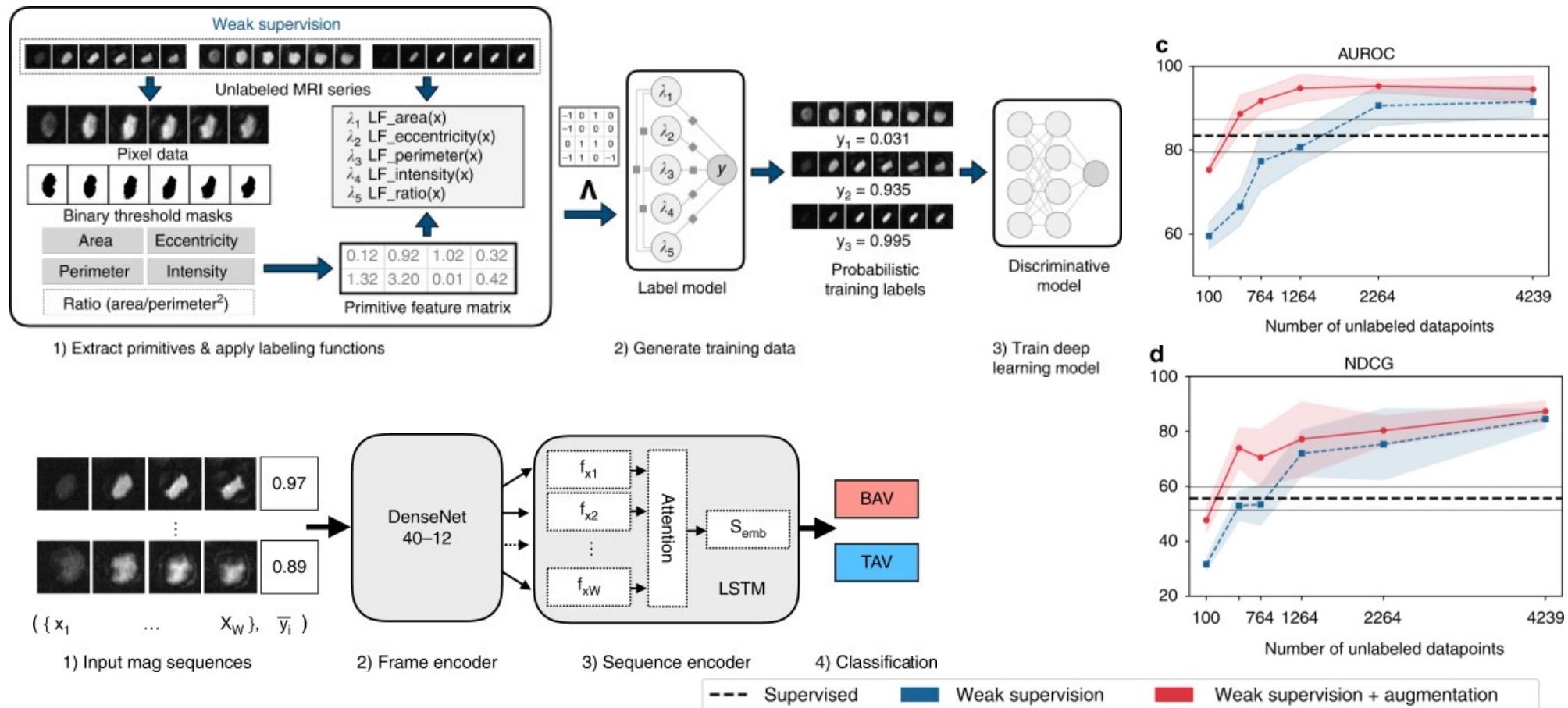


Contrastive-additive noise network (CAN)



Method	# scans used		Metrics			
	Labeled	Unlabeled	Dice[%]	Jaccard[%]	ASD[voxel]	95HD[voxel]
MT	16	64	88.23	79.29	2.73	10.64
MT-Dice [5]	16	64	88.32	79.37	2.76	10.50
Our UA-MT	16	64	88.88	80.21	2.26	7.32
Bayesian V-Net	8	0	79.99	68.12	5.48	21.11
Our UA-MT	8	72	84.25	73.48	3.36	13.84
Bayesian V-Net	24	0	88.52	79.70	2.60	10.45
Our UA-MT	24	56	90.16	82.18	2.73	8.90

# Training procedure



# Ground-truth annotation quality

Loosely define **biased errors** in segmentation labels as:

*A consistent error such that the annotations are always wrong in the same way.*



# Experimental setup

## Perturbations

Random crop

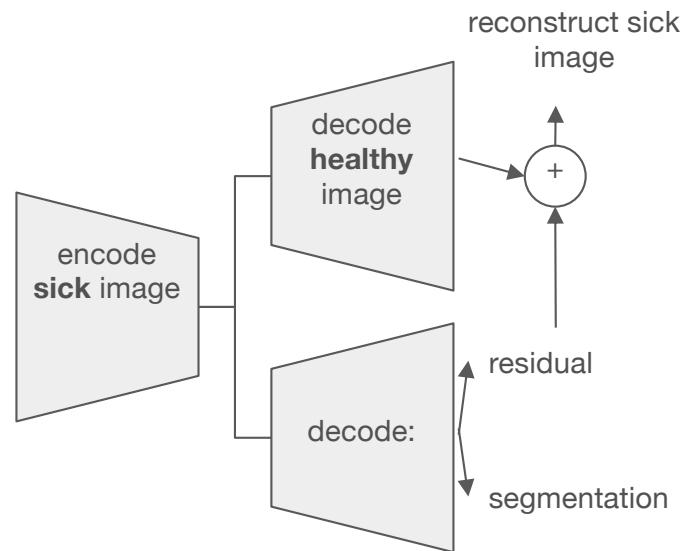
Constant shift

Random warp

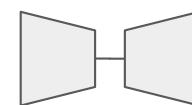
Permutation

## Method description

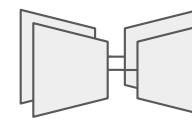
Image-to-image translator



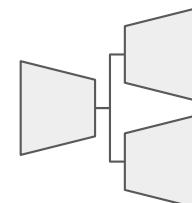
## Summary



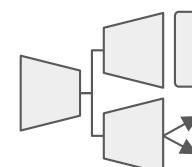
**Seg**



**MT**



**AE**

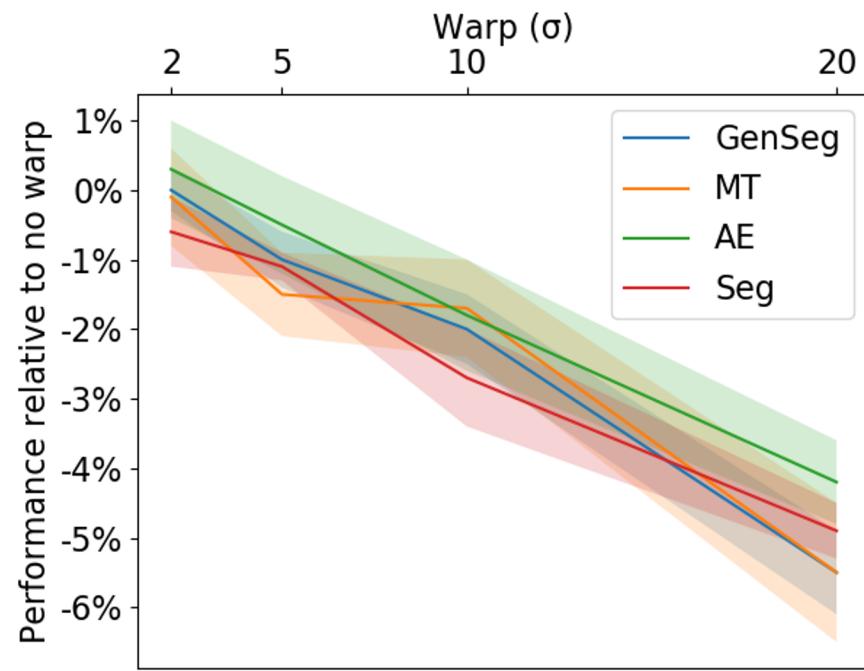


**GenSeg**

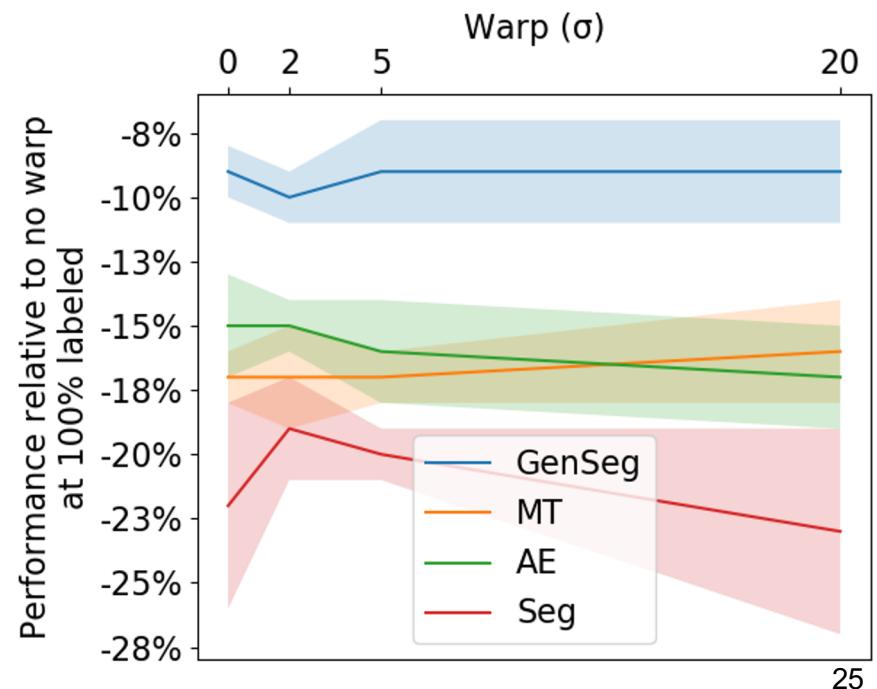
# Random warp



All data annotated

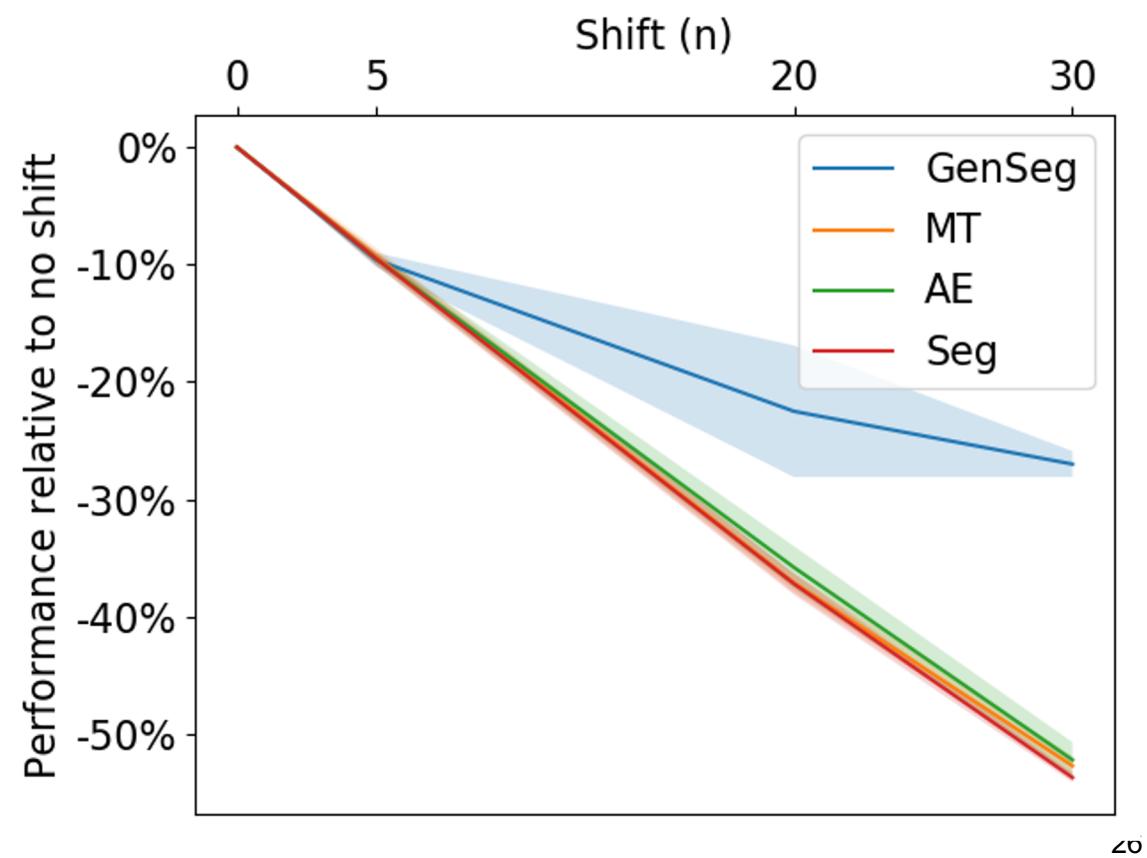
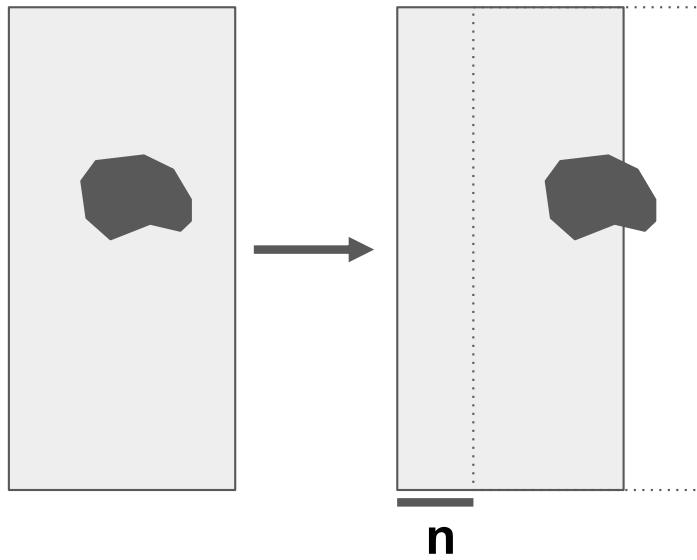


1% of data annotated



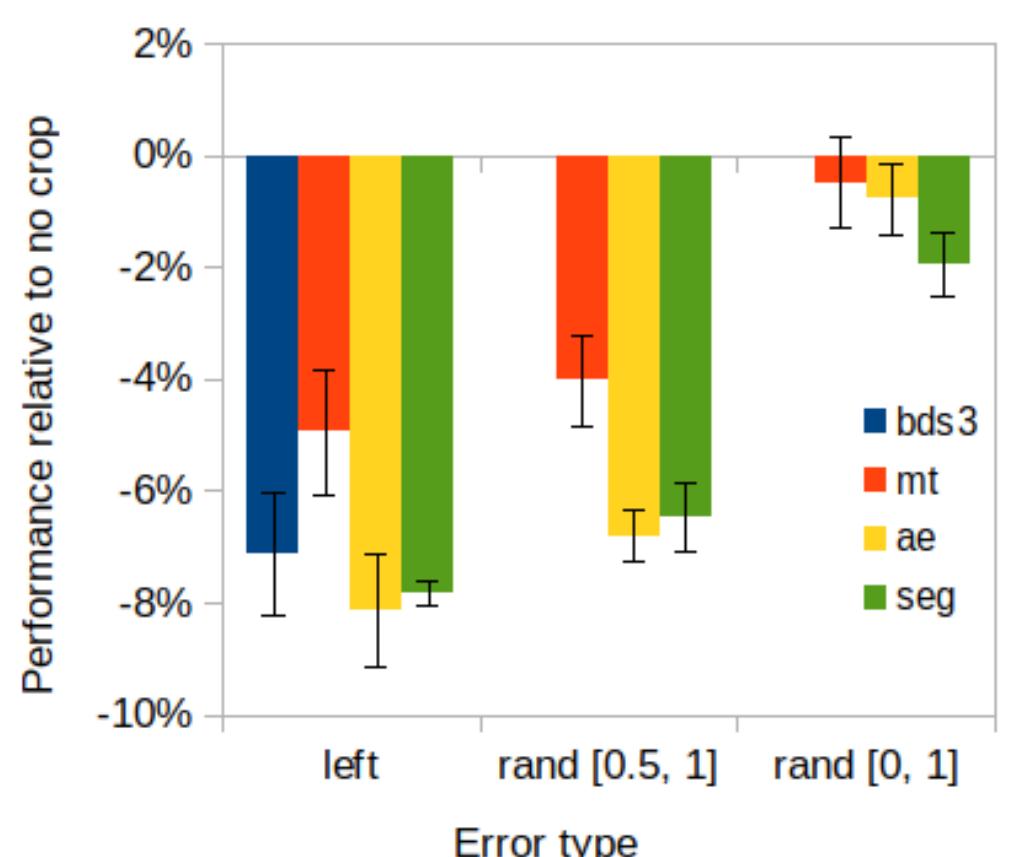
# Constant shift

Shift mask by  $n$  pixels

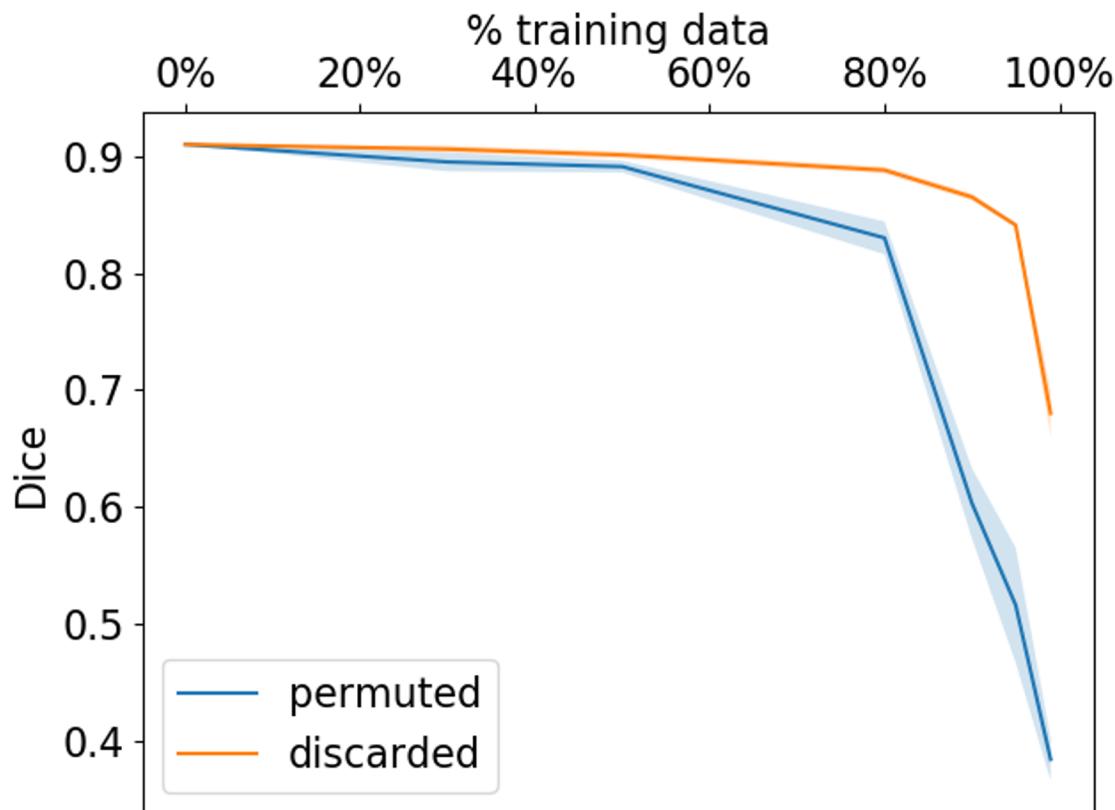


# Random crop

On average, crop half the tumour



# Permutation



- Permute labels across *% training data*
- VS discard *% training data*
- **Permutation is worse**

# Summary

Random warp

*Unbiased*

Low error

Constant shift

*Biased*

High error

Random crop

*Varying bias*

Low bias: low error  
Med bias: med error  
High bias: high error

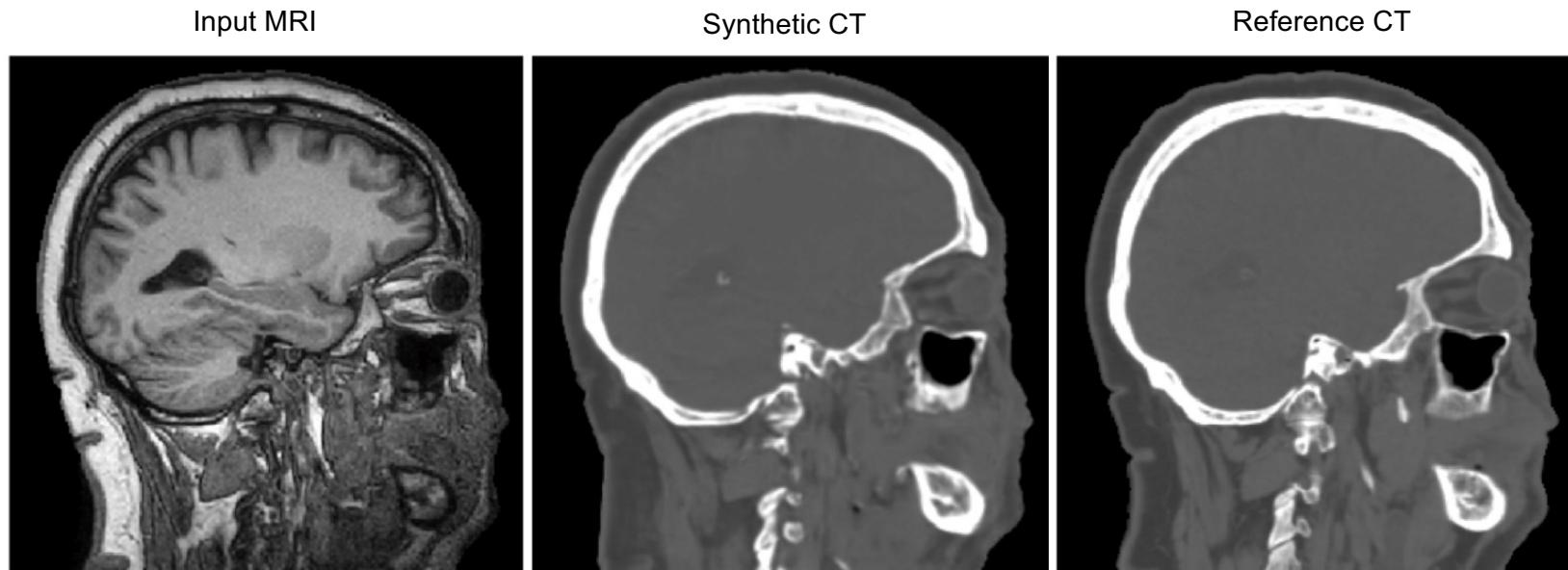
Permutation

*Biased*

High error

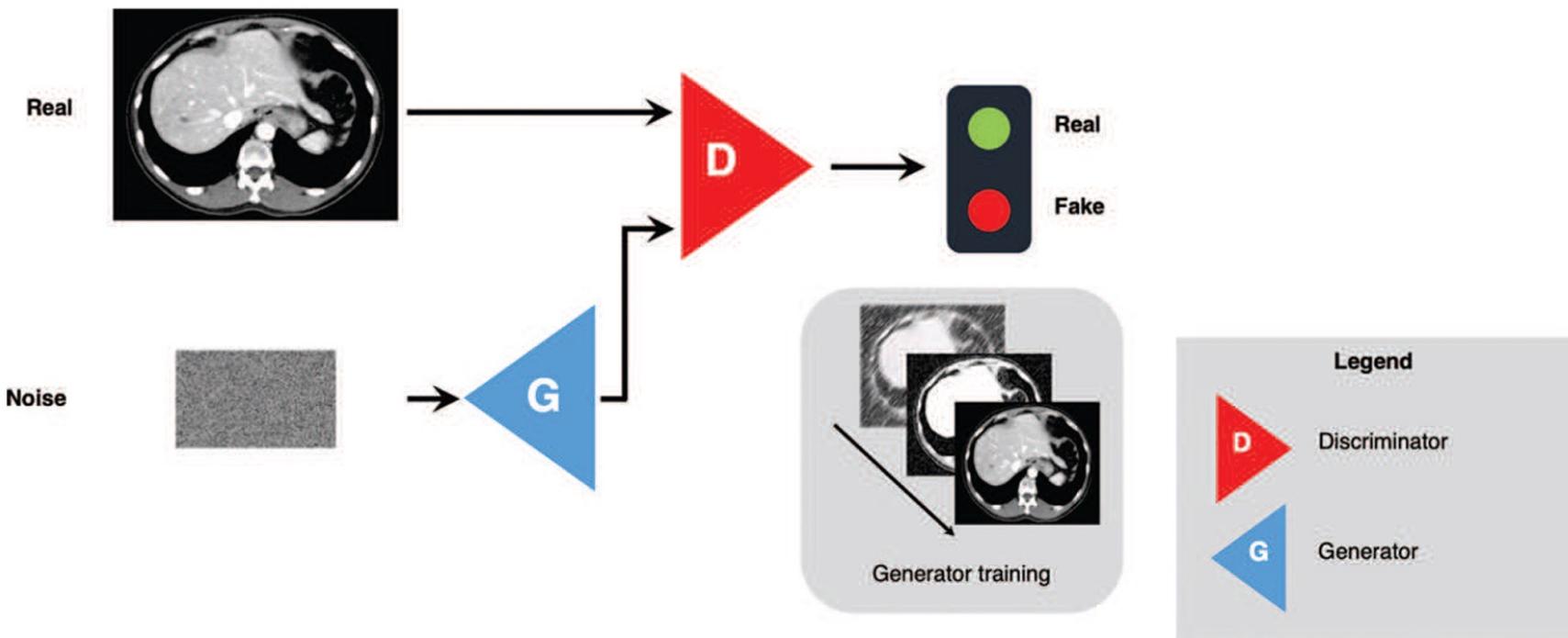
### 3) Image synthesis

# Image synthesis example



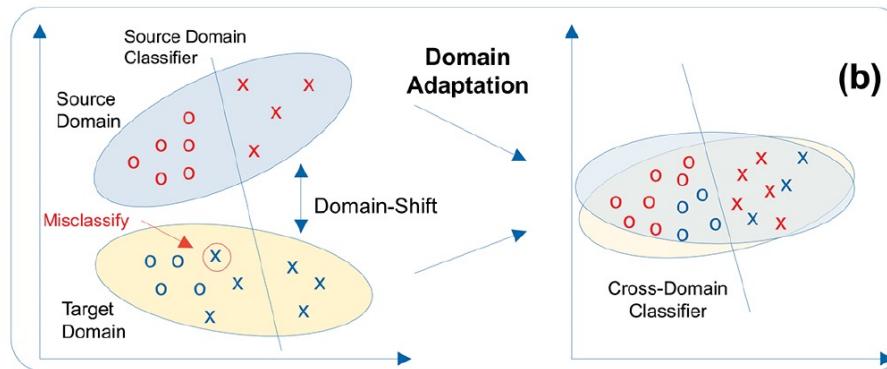
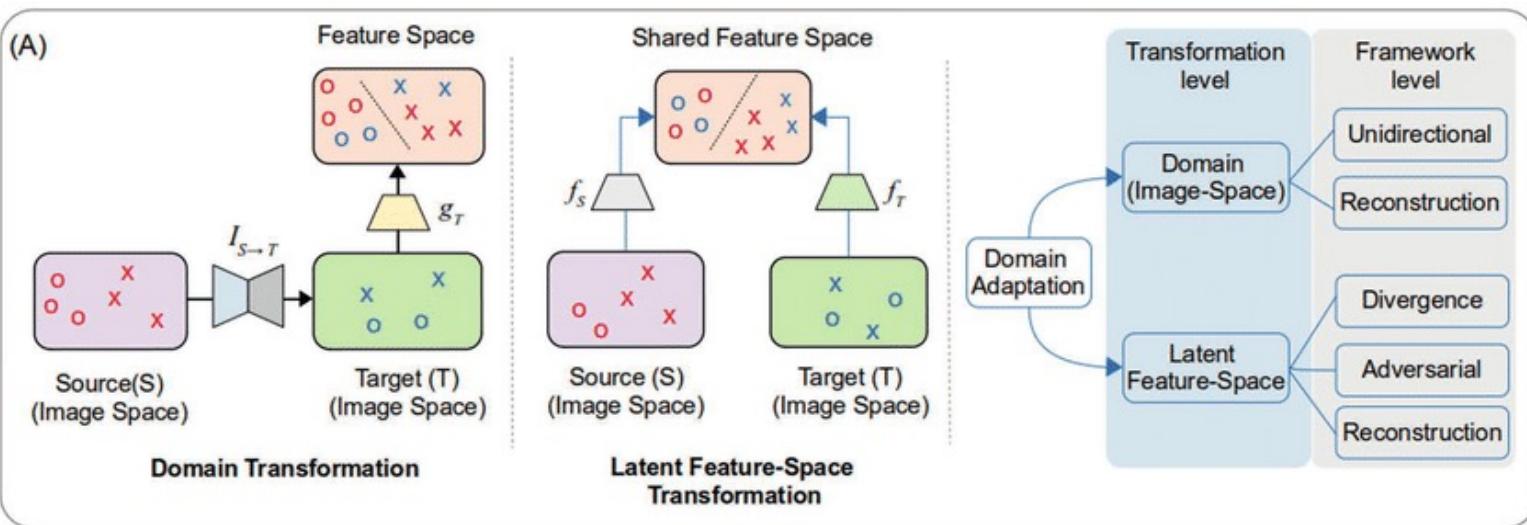
*MRI to CT Synthesis with a CycleGAN  
Source: Wolterink et al. 2017*

# Generative adversarial networks

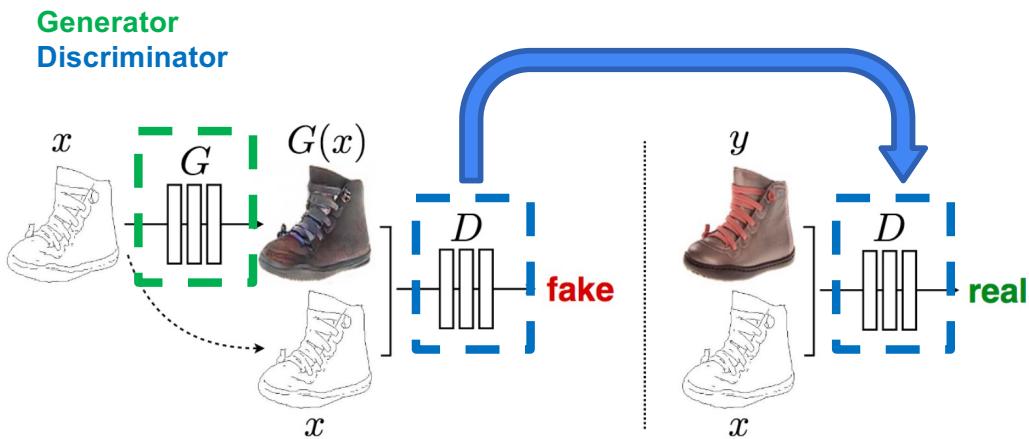


Standard GAN architecture diagram. A first network, known as the generator (G), aims to transform a random input into a realistic image to fool a second network, known as the discriminator (D). During training, the generator learns from the response of the discriminator.

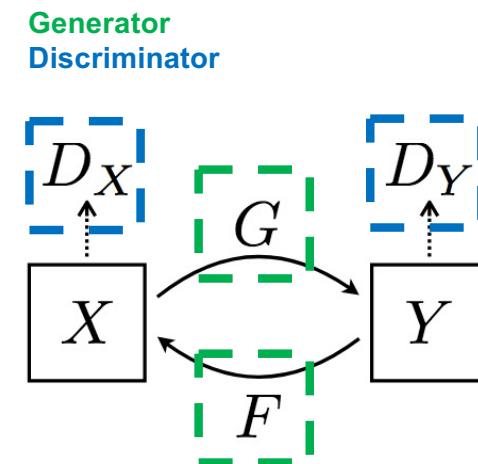
# Domain transformation scenarios



# Generative adversarial networks

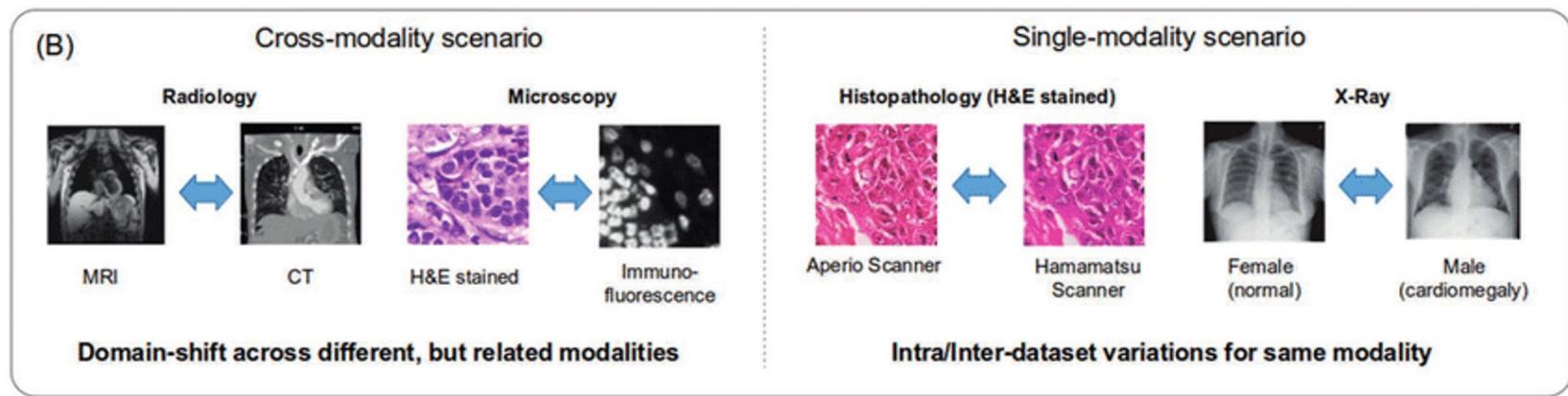


Sketch-to-picture synthesis with **conditional generative adversarial network** (cGAN)  
Source: Isola et al. 2017

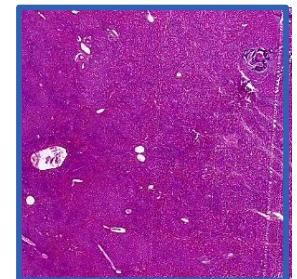
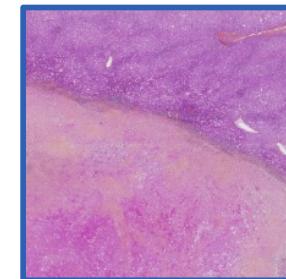
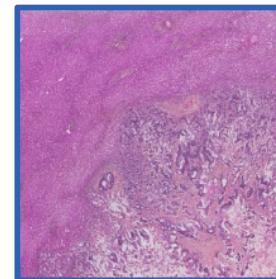
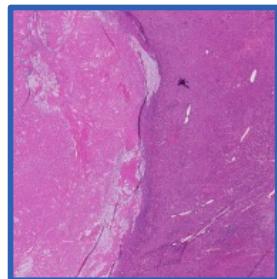
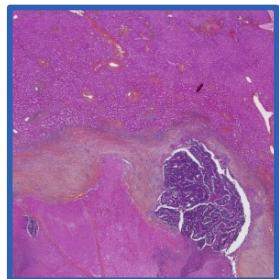


Architecture of a **cycle-consistent generative adversarial network** (CycleGAN)  
Source: Zhu et al. 2017

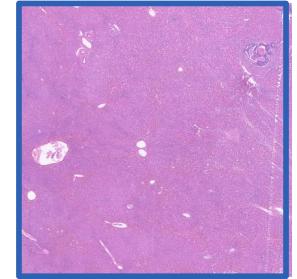
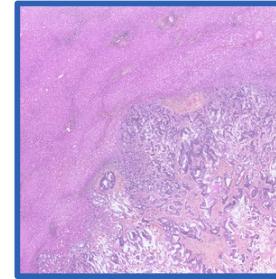
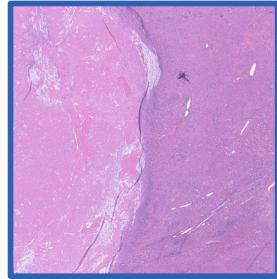
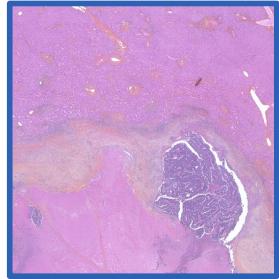
# Examples of normalized slides with the GAN model



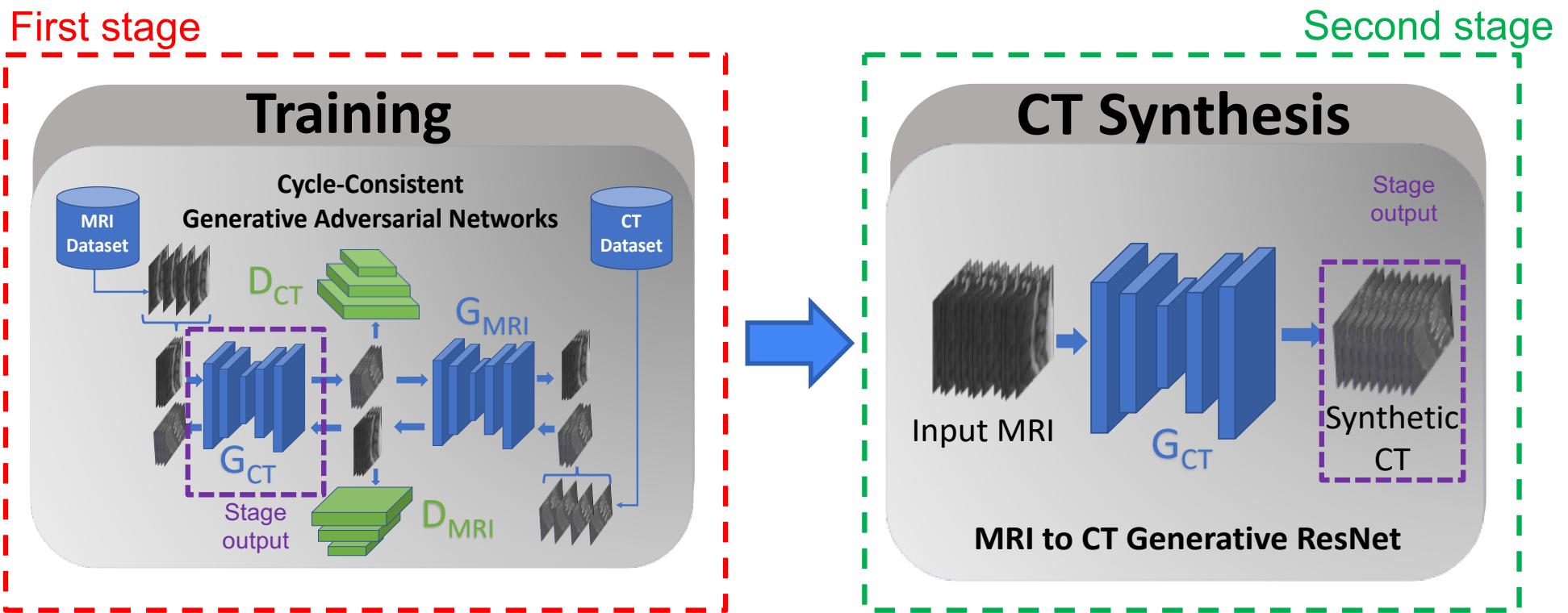
Original



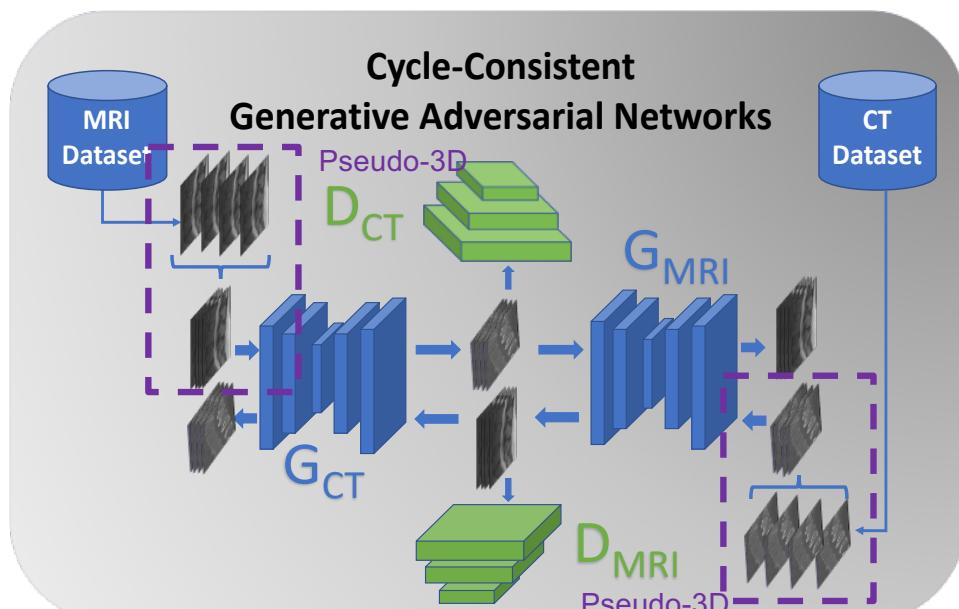
Normalized



# MRI to CT synthesis pipeline



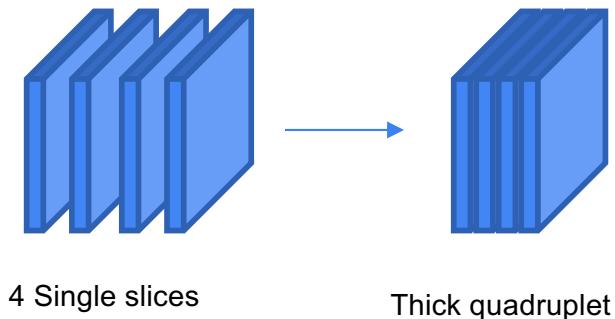
# Cycle-Consistent GAN



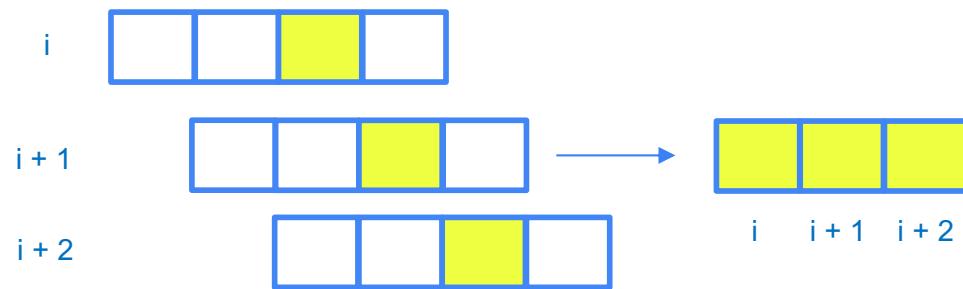
CycleGAN MRI and CT architecture

- **Image synthesis using CycleGAN model**, jointly learning both **MRI → CT** and **CT → MRI** mappings.
- We use a **pseudo-3D** approach, contrary to the more classical single-slice 2D approach
- Pseudo-3D consists of **stacking neighboring slices** along the channel dimension, forming a **thick slice**

# Pseudo-3D recombination



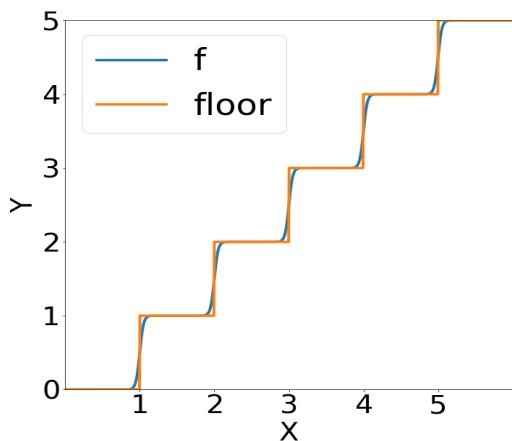
*Slice thickening strategy.*



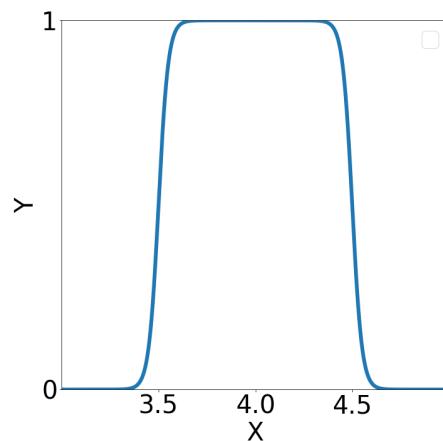
*Recombining strategy.  
Left, quadruplet slices. Right, output volume slices.*

- Pseudo-3D is performed for each **quadruplet** of neighbouring slices
- **Recombining quadruplets** into a single volume is done selecting the third slice from each quadruplet and adding it to the output volume
- With this recombining strategy, slices at indices  $i$  and  $i + 1$  in the output volume necessarily come from **quadruplets with 3 slices in common**

# Differential histograms



*Floor function and smooth counterpart.*



*Smooth indicator function for the (4-0.5, 4+0.5) interval.*

- **Differentiable histograms** for image synthesis, to be used with **backpropagation**
- Regular histogram density computation has **nil gradient** due to the **non-differentiable floor and Dirac indicator functions**. It is thus impossible to use backpropagation.
- Proposed solution: **smooth out** those functions

# CycleGAN loss

- CycleGAN loss function

$\lambda_{cycle} = 10$   
 $\lambda_{hist} = 0.1$   
 $n_{bins} = 32$

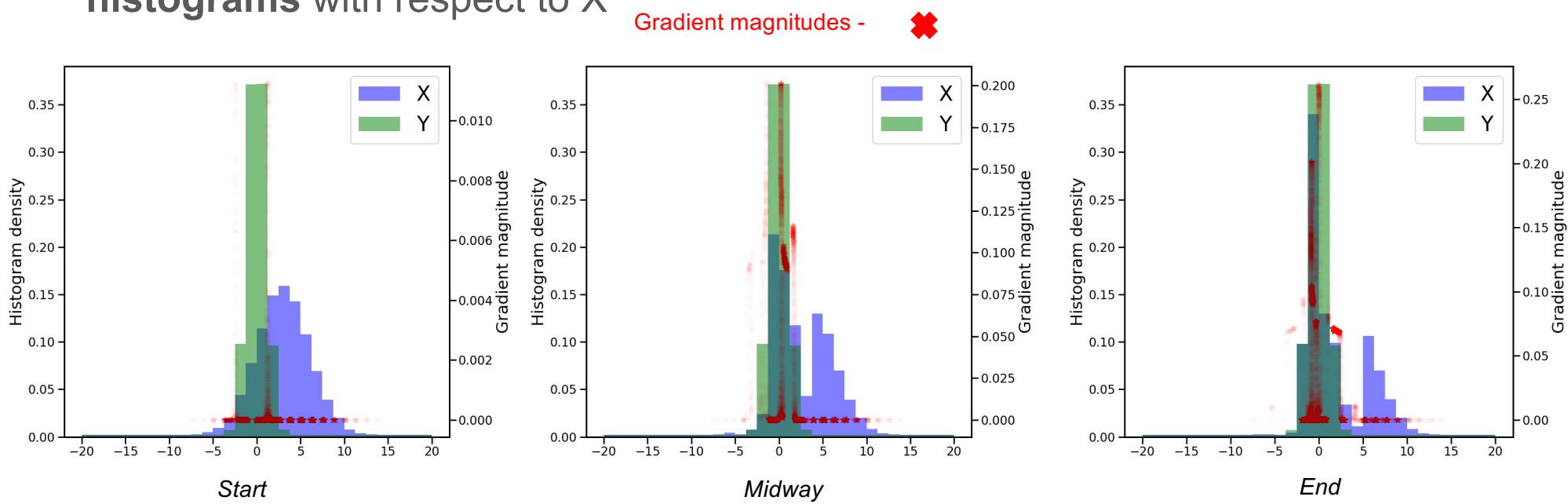
$$L(G_{CT}, G_{MRI}) = \left\| D_{CT}(CT_{syn}) - 1 \right\|_2^2 + \left\| D_{MRI}(MRI_{syn}) - 1 \right\|_2^2$$
$$+ \lambda_{cycle} [ \|MRI_{recons} - MRI\|_1 + \|CT_{recons} - CT\|_1 ]$$
$$+ \lambda_{hist} [ \|H^{MRI}_{recons} - H^{MRI}\|_1 + \|H^{CT}_{recons} - H^{CT}\|_1 ]$$

Adversarial loss  
Cycle-consistency loss  
Histogram loss

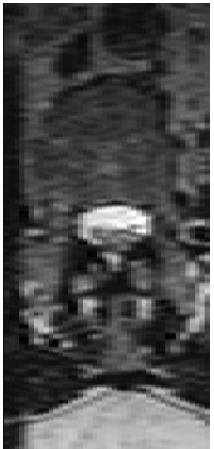
$$L(D_{CT}, D_{MRI}) = \left\| D_{CT}(CT) - 1 \right\|_2^2 + \left\| D_{MRI}(MRI) - 1 \right\|_2^2$$
$$+ \left\| D_{CT}(CT_{syn}) \right\|_2^2 + \left\| D_{MRI}(MRI_{syn}) \right\|_2^2$$

# Histogram similarity

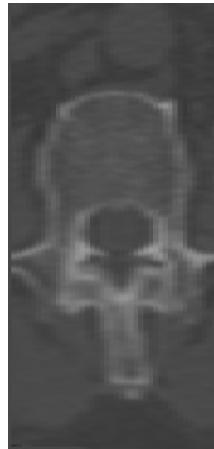
- Toy experiment, with  $X$  being a vector sampled from  $N(3,3)$ , and  $Y$  being a vector sampled from  $N(0,1)$ . We **minimize the difference of histograms with respect to  $X$**



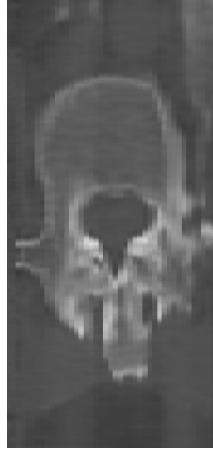
# Samples from the image synthesis stage



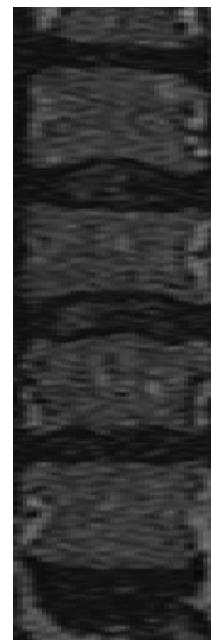
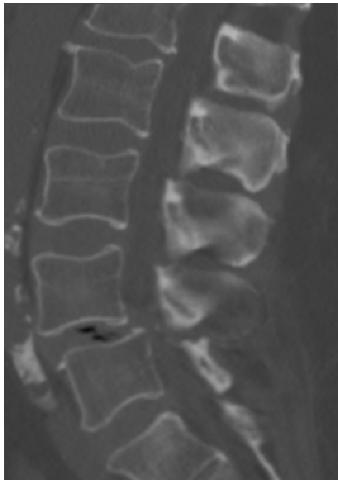
*MRI*



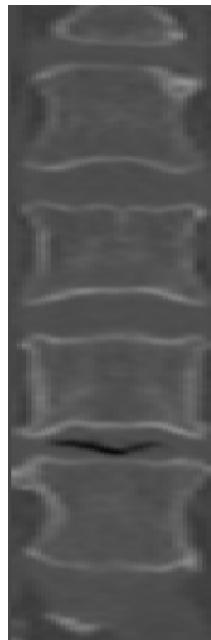
*CT*



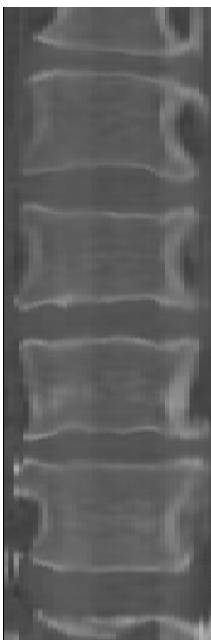
*Synthetic CT*



*MRI*

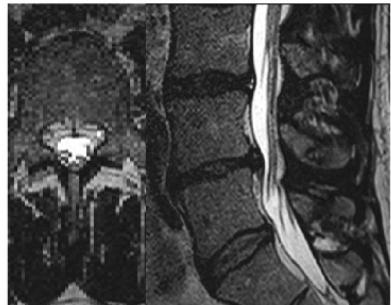


*CT*



*Synthetic CT*

# Axial resolution



(a) Input MRI



(b) Ground truth CT



(c) Single-slice model prediction

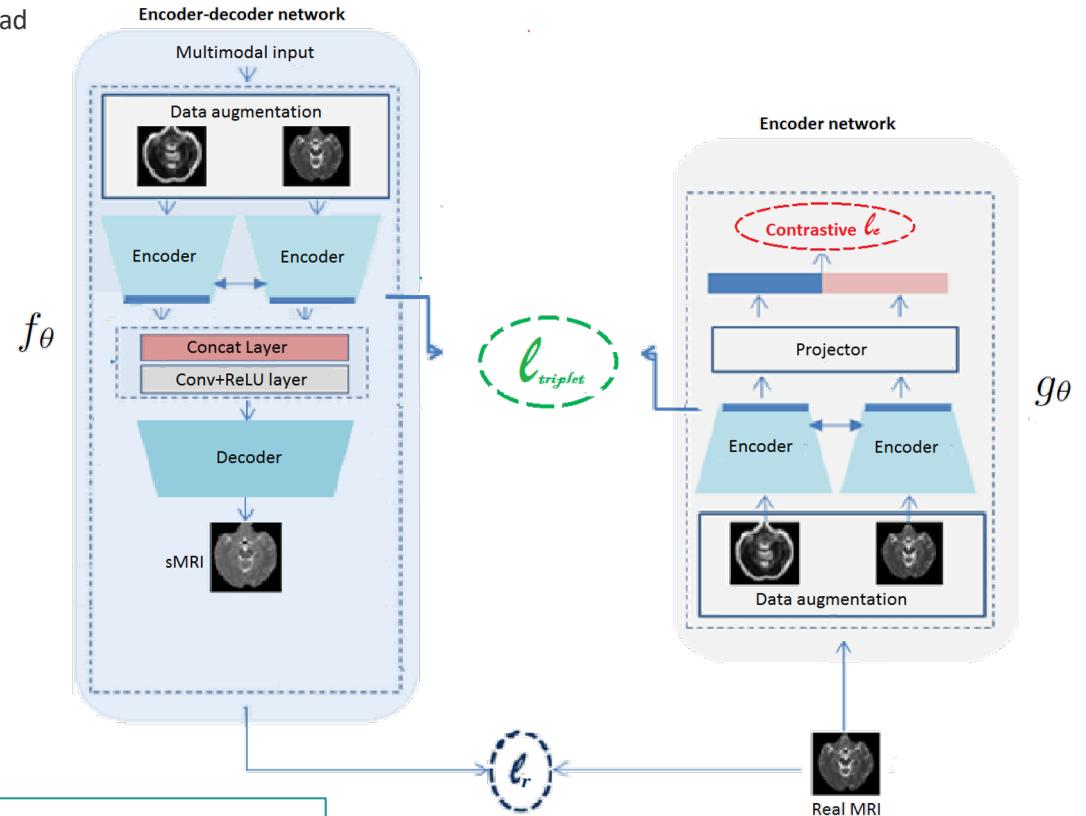
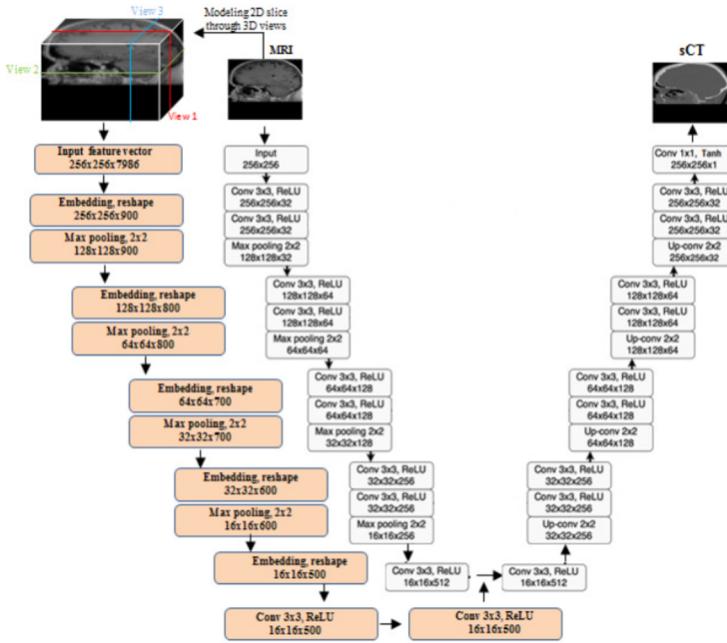
(d) Pseudo-3D model prediction

*Sample axial and sagittal slices for one patient.*

- **Axial resolution improvement** with the **pseudo-3D** approach, opposed to a single-slice 2D approach

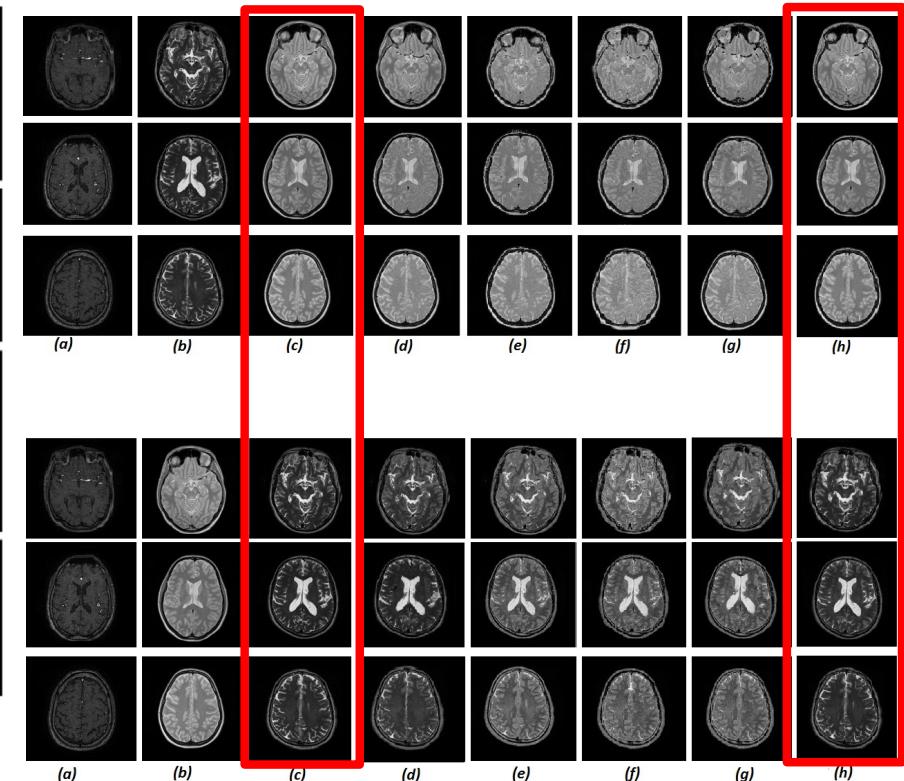
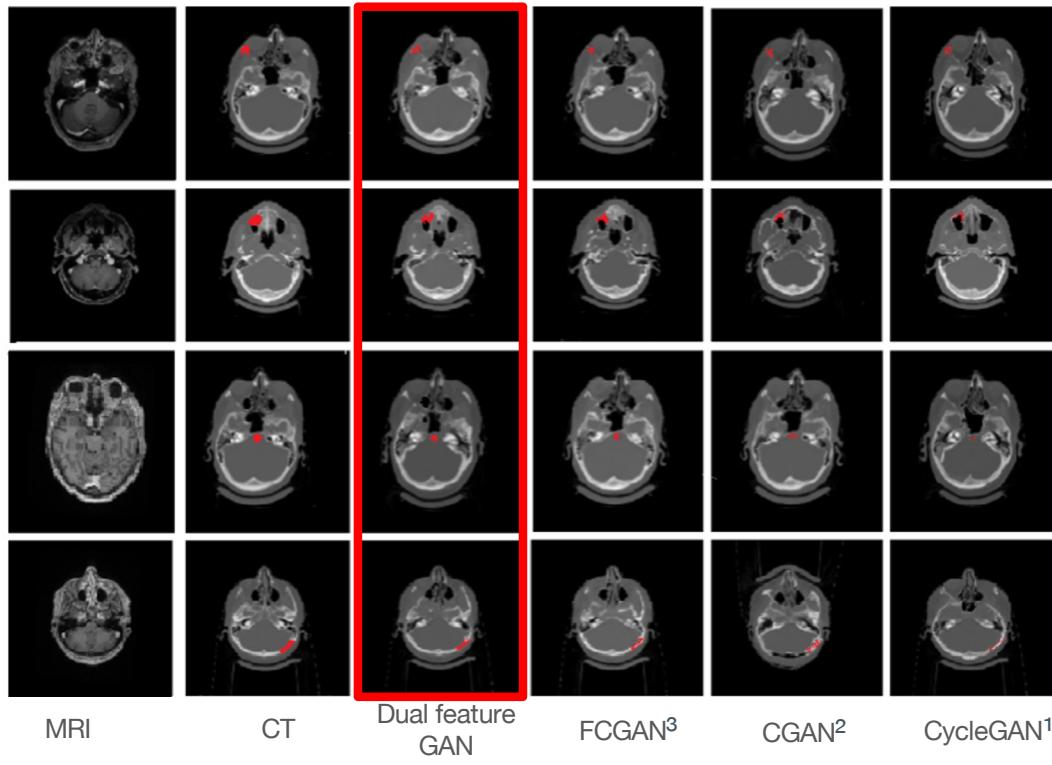
# Contrastive feature GAN

R. Touati et al. "A feature invariant generative adversarial network for head and neck MRI/CT image synthesis", *Phys. Med. Biology*, 2021.



**Contrastive feature GAN.** Generative adversarial network augmented with multi-planar dynamic features to improve spatial positioning and preserve anatomical structures.

# MRI-CT synthesis results

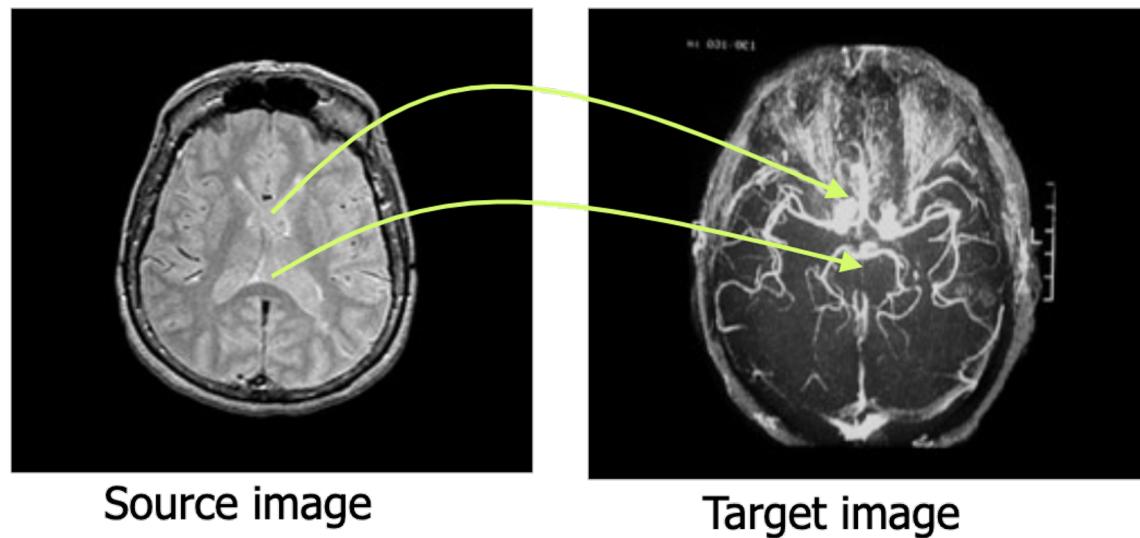


1. Isola P, Zhu JY, Zhou T, Efros AA. *Image-to-image translation with conditional adversarial networks*. CVPR 2017.
2. Oulbacha R, Kadouri S. *MRI to CT Synthesis of the Lumbar Spine from a Pseudo-3D Cycle GAN*. ISBI 2020.

## 4) Deformable registration

# Deformable registration

" Image registration is about determining a spatial transformation – or mapping – that relates positions in one image, to corresponding positions in one or more other images"

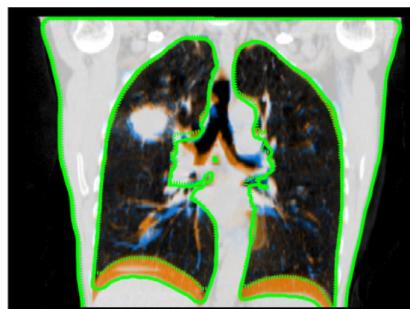


- 3D - 3D
- 3D - 2D
- 3D/2D - patient

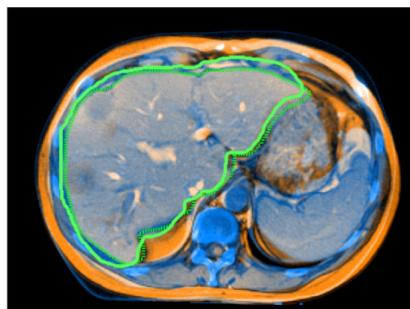
# Organ shift and deformations



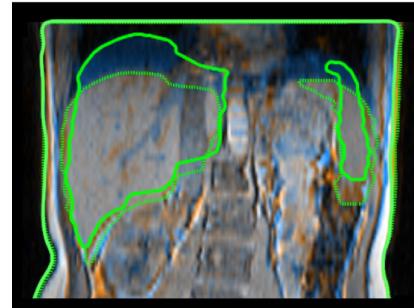
Abdominal 4DCT, n=10



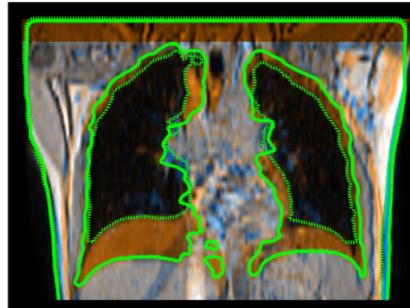
Thoracic 4DCT, n=16



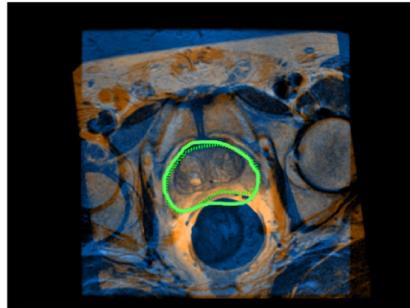
Liver CT-MR, n=18



Abdominal MR, n=5

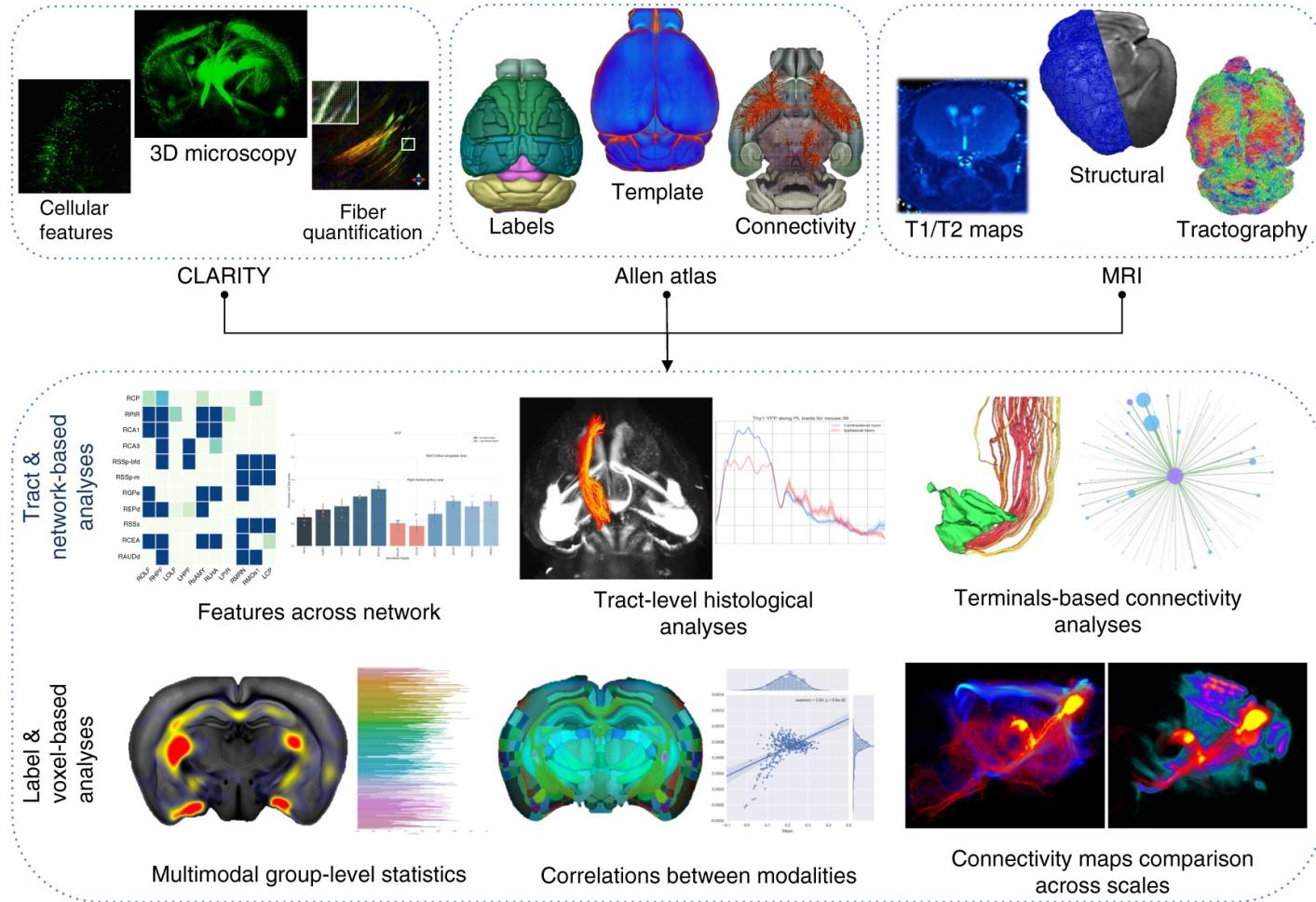


Thoracic MR, n=5

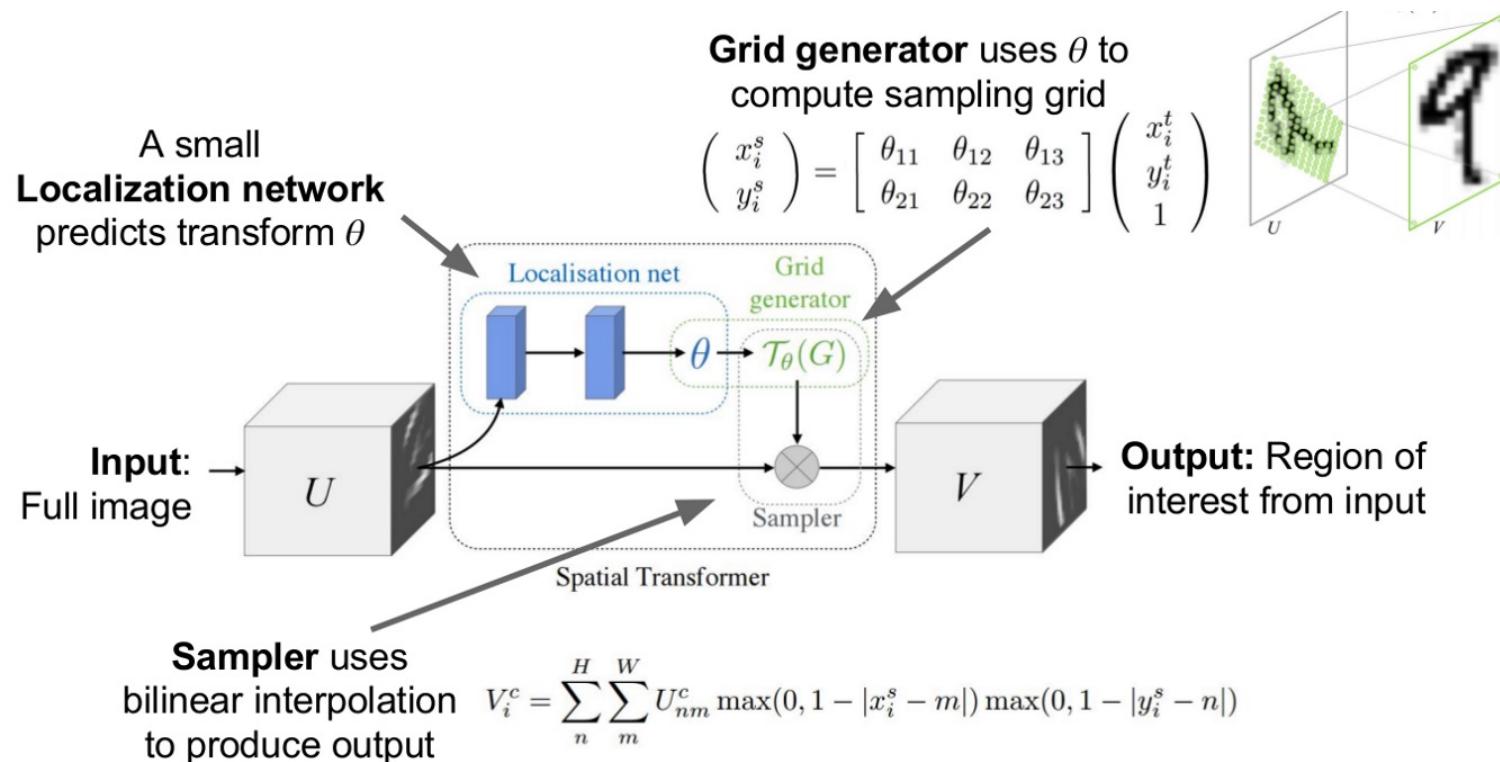


Prostate MR, n=20

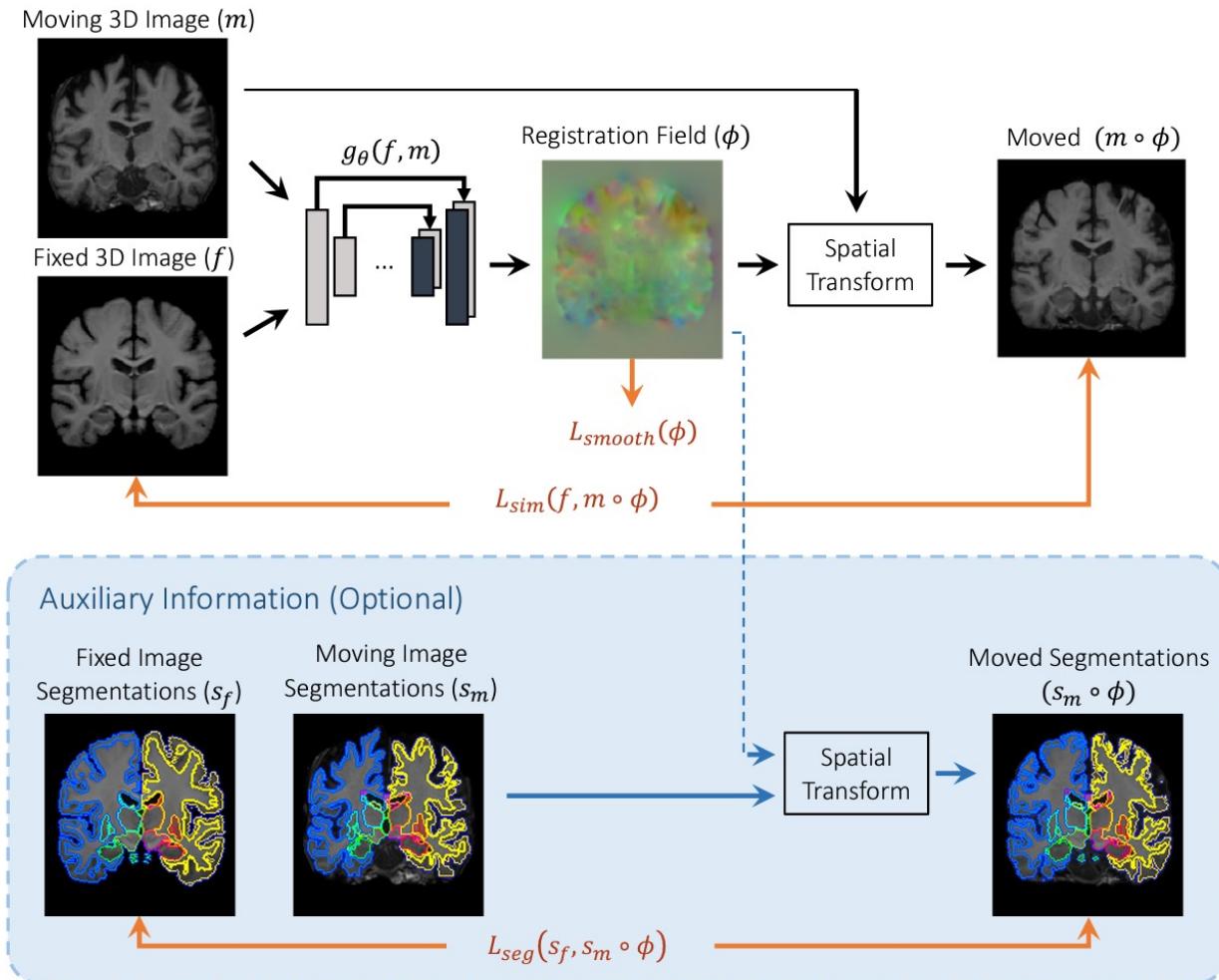
# Applications in deformable image registration



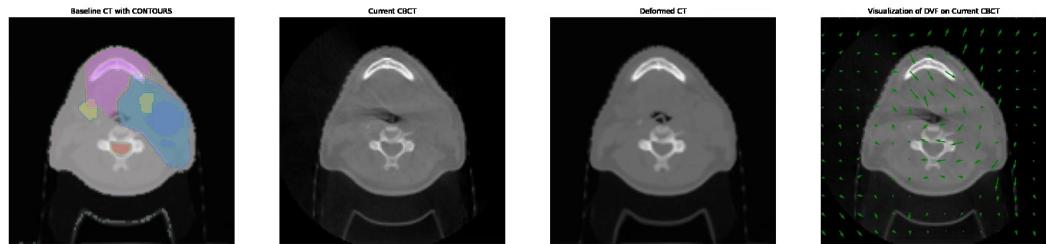
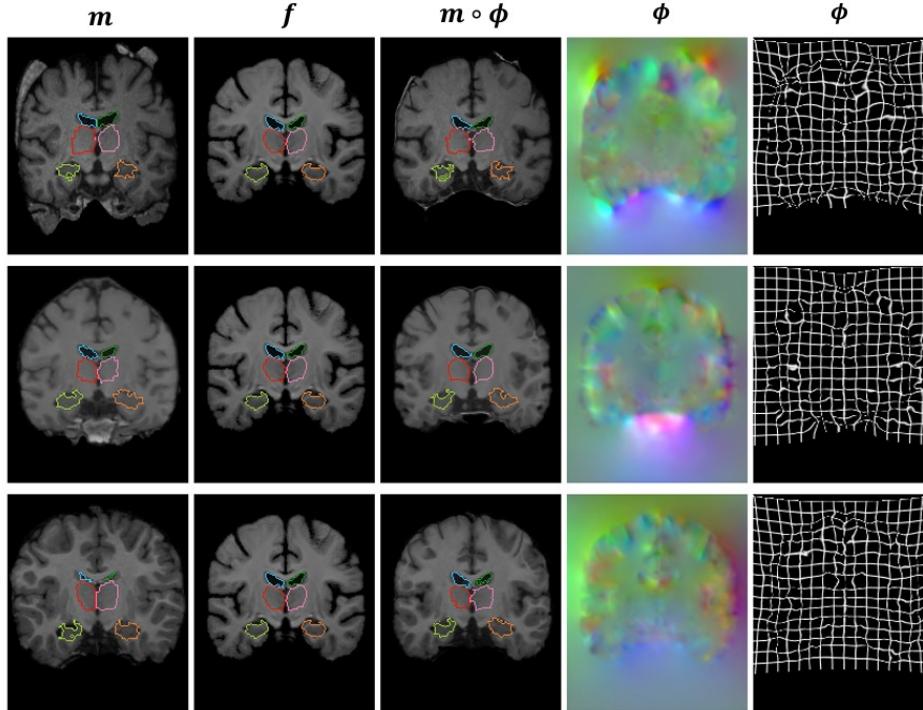
# Spatial transformer networks



# VoxelMorph



# Results with Voxelmorph

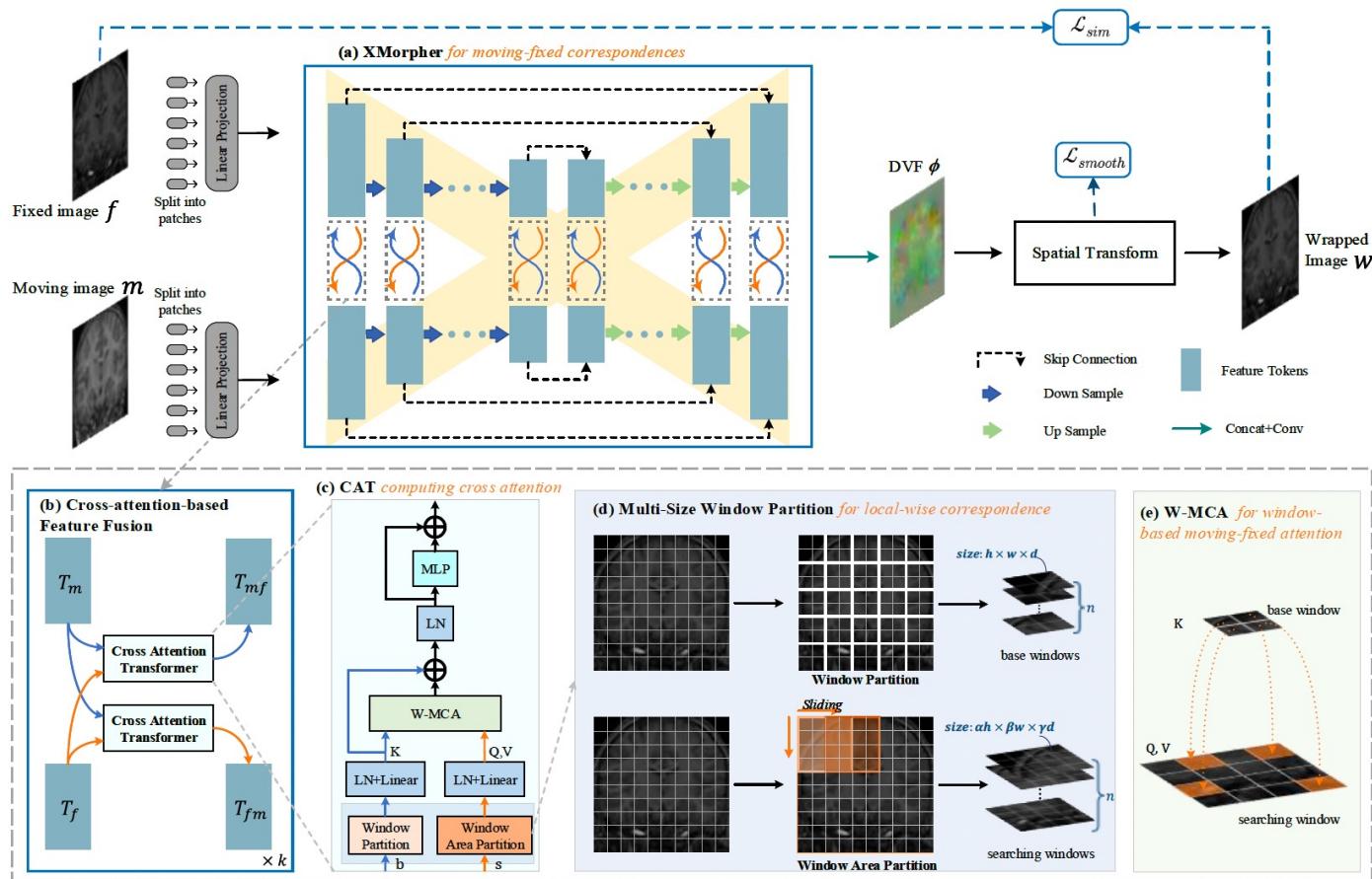


Method	Dice
Affine only	0.579 (0.173)
ANTs SyN (CC)	0.761 (0.117)
NiftyReg (CC)	0.772 (0.117)
VoxelMorph (MSE)	0.727 (0.146)
VoxelMorph x2 (MSE)	0.750 (0.058)
VoxelMorph x2 (MSE) inst.	0.764 (0.048)
VoxelMorph (CC)	0.737 (0.139)
VoxelMorph x2 (CC)	0.763 (0.049)
VoxelMorph x2 (CC) inst.	0.772 (0.119)

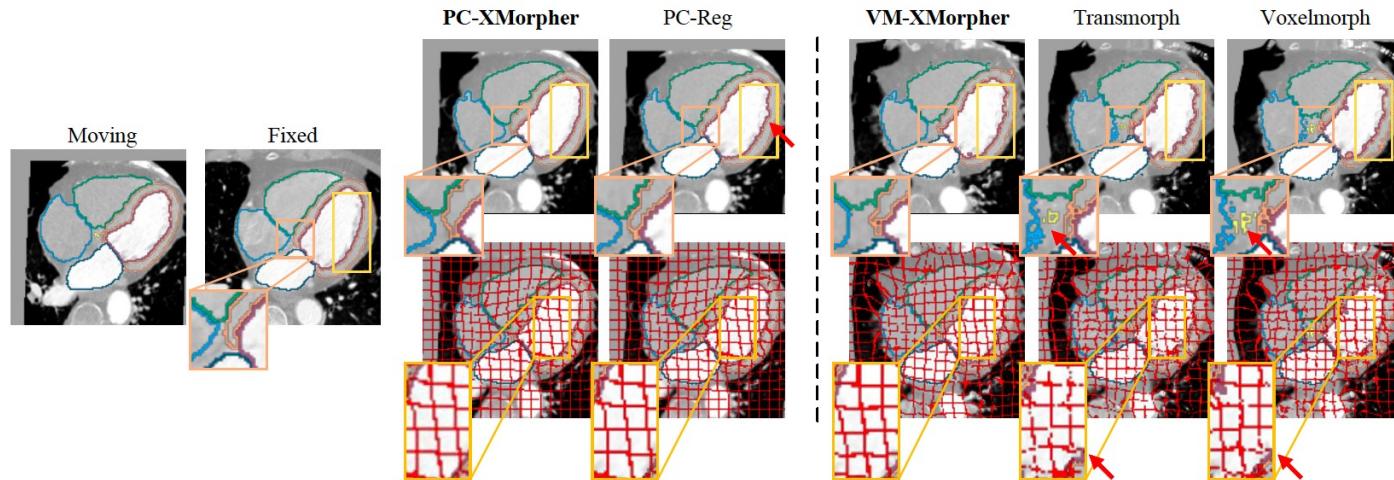
TABLE III: Results for subject-to-subject alignment using affine, ANTs, and VoxelMorph variants, where “x2” refers to a model where we doubled the number of features to account for the increased inherent variability of the task, and “inst.” indicates additional instance-specific optimization.

Balakrishnan et al. 2019, IEEE TMI

# XMorpher



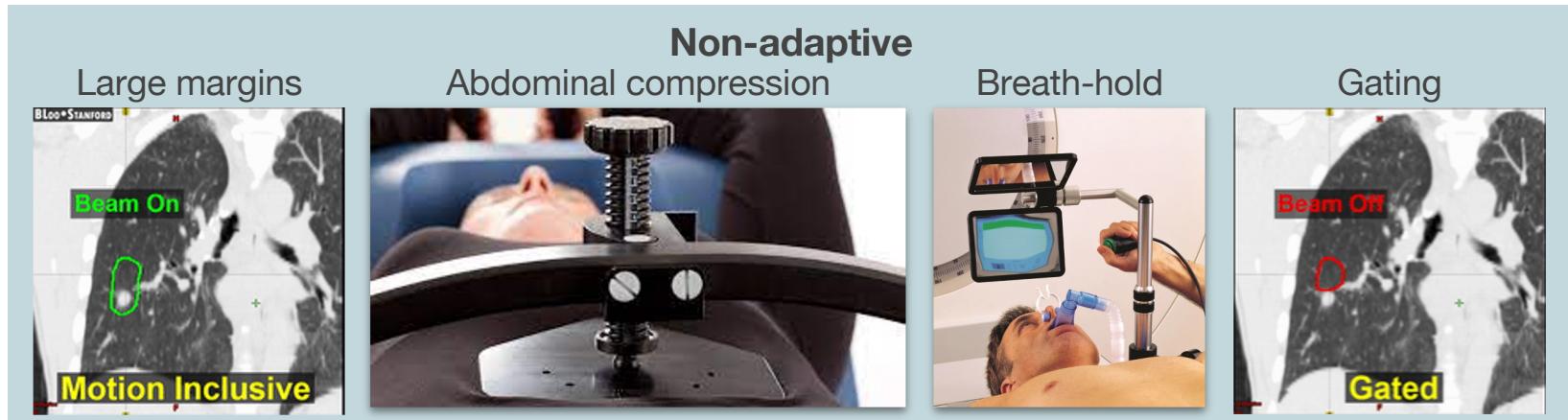
# XMorpher



Method	Un-/Semi-	Backbone	DSC	$ J_\phi  \leq 0$ (%)
Affine initialization	-	-	$69.2 \pm 7.2$	-
BSpline [14]		-	$80.9 \pm 7.6$	$5.25 \pm 3.27$
Voxelmorph [1]		CNN	$80.2 \pm 5.5$	$4.02 \pm 0.82$
Transmorph [4]	Unsup	CNN+Transformer	$81.1 \pm 5.2$	$3.46 \pm 0.75$
<b>Our no-cross XMorpher</b>		Full Transformer	$81.5 \pm 5.4$	<b><math>0.94 \pm 0.26</math></b>
<b>Our VM-XMorpher</b>		Full Transformer	<b><math>83.0 \pm 4.7</math></b>	$3.15 \pm 0.79$
PC-Reg [8]	Semi	CNN	$86.0 \pm 2.5$	$0.36 \pm 0.20$
<b>Our PC-XMorpher</b>		Full Transformer	<b><math>86.9 \pm 2.4</math></b>	<b><math>0.32 \pm 0.18</math></b>

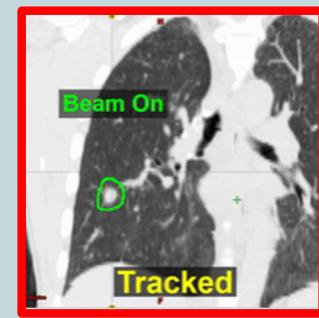
## 5) Motion modeling

# Motion management strategies



## Real-time adaptive tracking

- Move or shape the radiation beam **dynamically**
- **Accuracy** depends on the system **adapting** to the moving anatomy
- Driven by some **surrogate** signal to estimate the target position

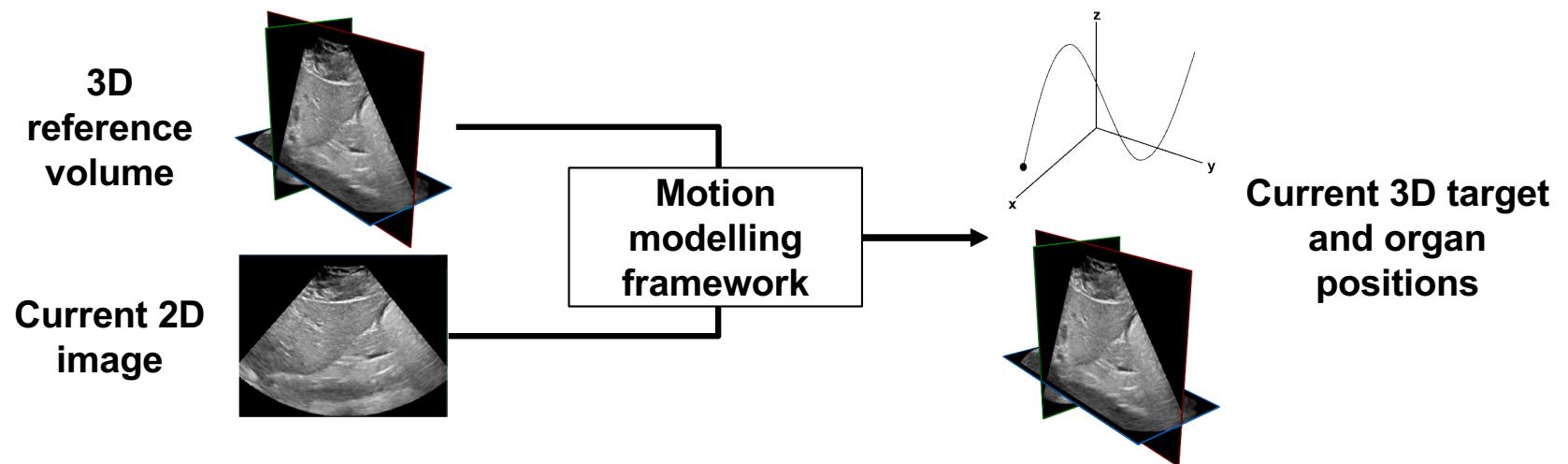


Source: <https://tibaray.com/medical.html>

# Real-time 3D ultrasound generation and tracking

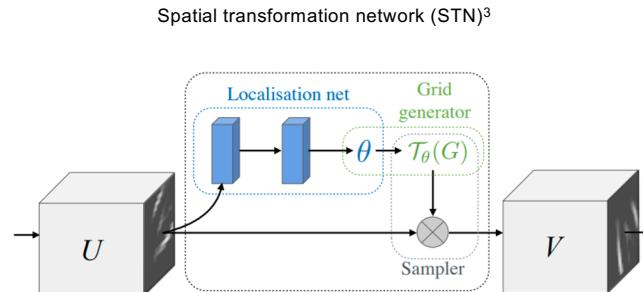
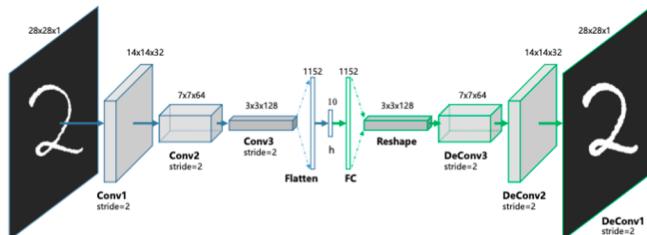
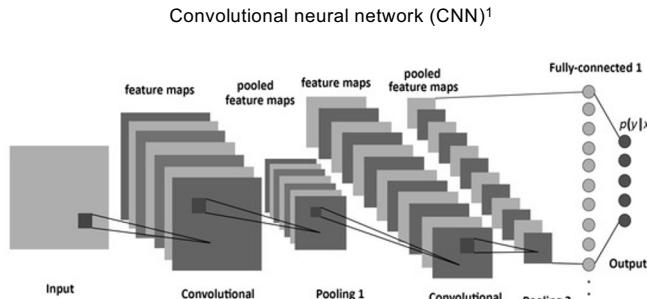
Motion modelling framework to generate up-to-date 3D organ and target positions from input 2D US images and a limited number of reference volumes

- Train DNN with 4D US dataset of the liver
- Develop motion modelling framework
- Evaluate framework on 3D target tracking in US volumes



# Deep learning for motion modelling

- Architectures
  - CNN → Feature extraction
  - Convolutional AE → Representational power
  - STN → Spatial warping
  - GAN → Generative capacity
- Advantages
  - Good generalization
  - No inter-subject correspondences
- Disadvantages
  - Requires large amount of data



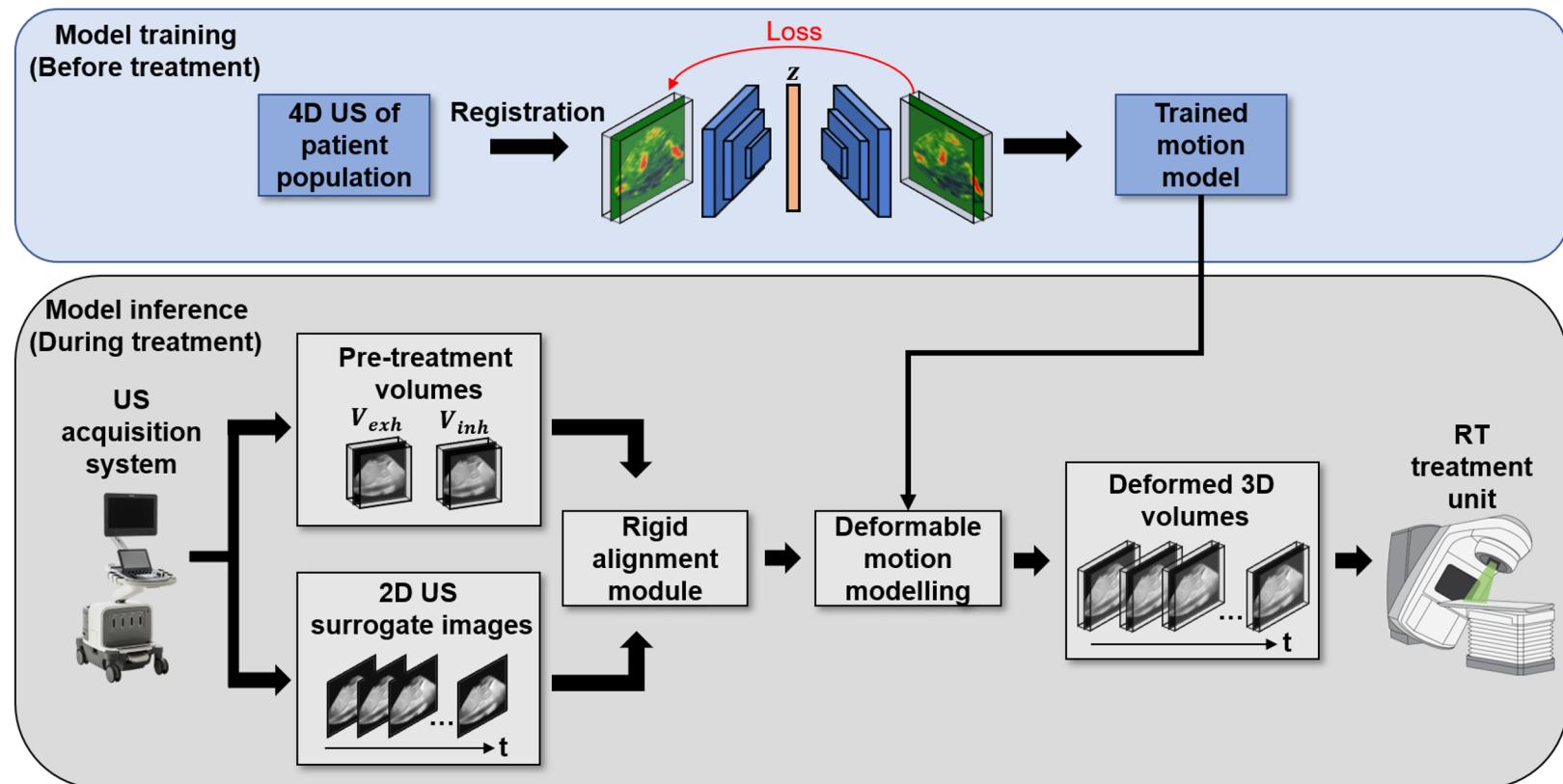
1. S. Albelwi et A. Mahmood, "A framework for designing the architectures of deep convolutional neural networks," *Entropy*, vol. 19, no. 6, 2017.

2. Towards Data Science. (2019) Convolutional autoencoders for image noise reduction. [Online].

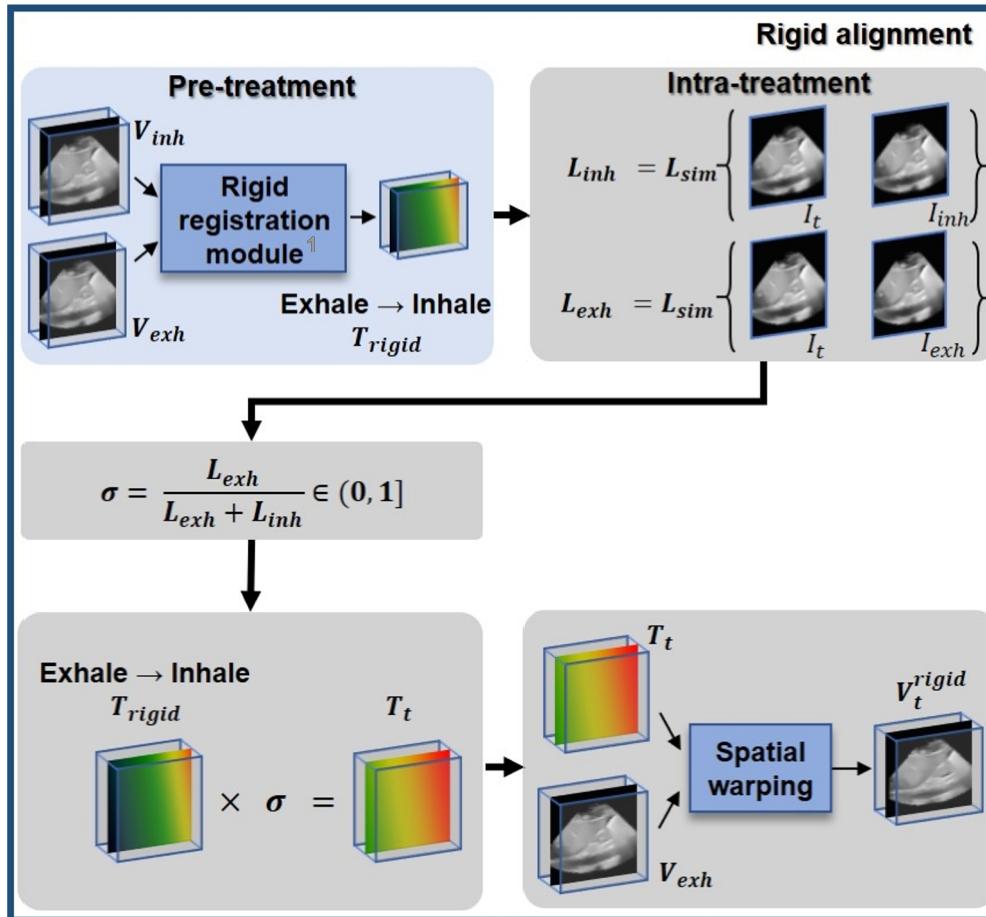
3. M. Jaderberg, K. Simonyan, A. Zisserman et al., "Spatial transformer networks," in *Advances in neural information processing systems*, 2015, p. 2017–2025.

# Overview of US-based motion model

Tal Mezheritsky et al., "Population-based 3D motion modelling from convolutional autoencoders for 2D ultrasound-guided radiotherapy" **Medical Image Analysis**, 2022.



# Rigid alignment module



1. S. Klein, M. Staring et al., "elastix : a toolbox for intensity-based medical image registration,"IEEE Transactions on Medical Imaging,vol. 29, no. 1, p. 196 – 205, January 2010.

# Deformable motion model

## Training progress

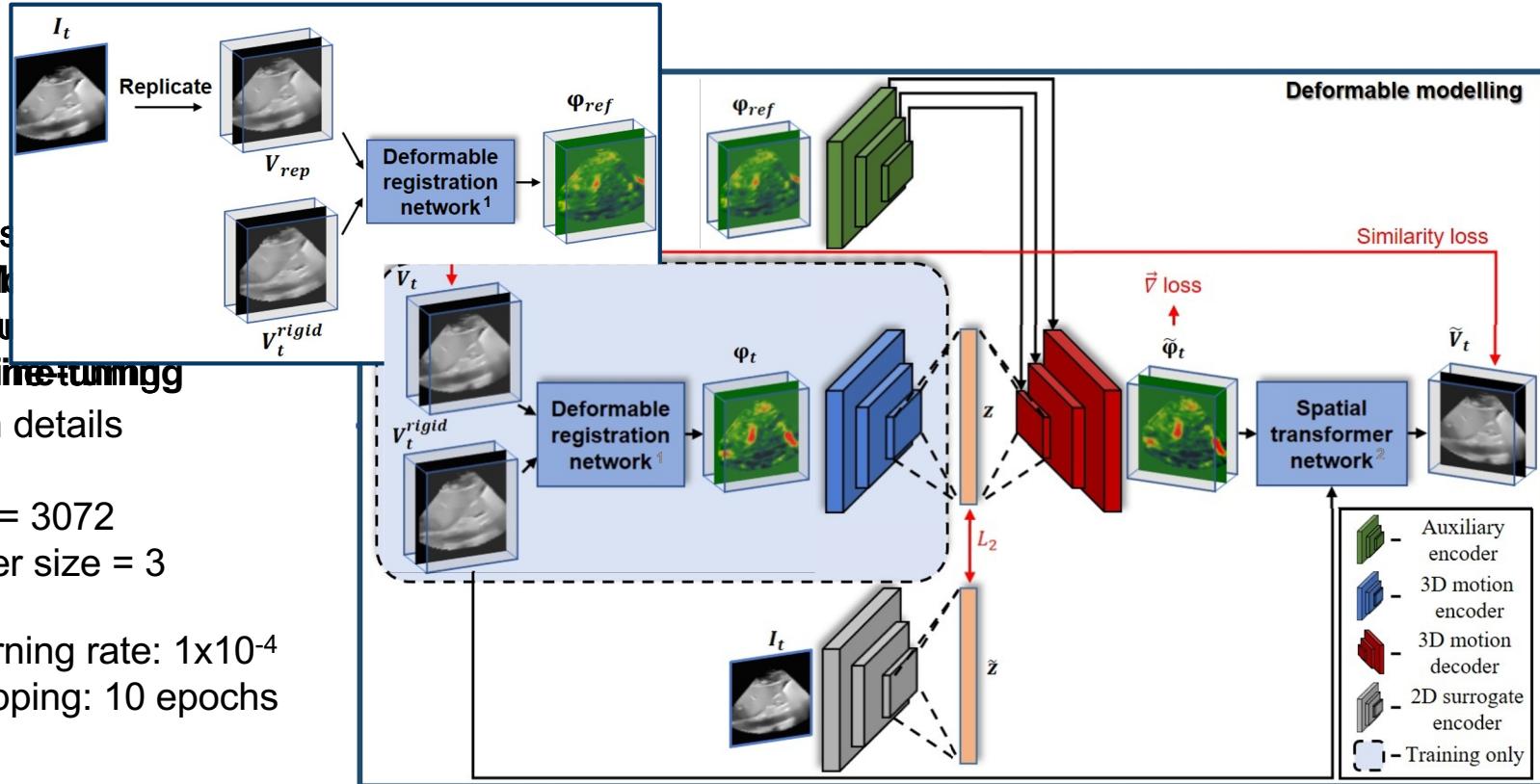
- Step 1: Model training
- Step 2: Supervised learning
- Step 3: Finetuning

## Implementation details

- $\beta = 0.01$
- Size of  $z$  = 3072
- Conv. filter size = 3

## Training details

- Initial learning rate:  $1 \times 10^{-4}$
- Early stopping: 10 epochs

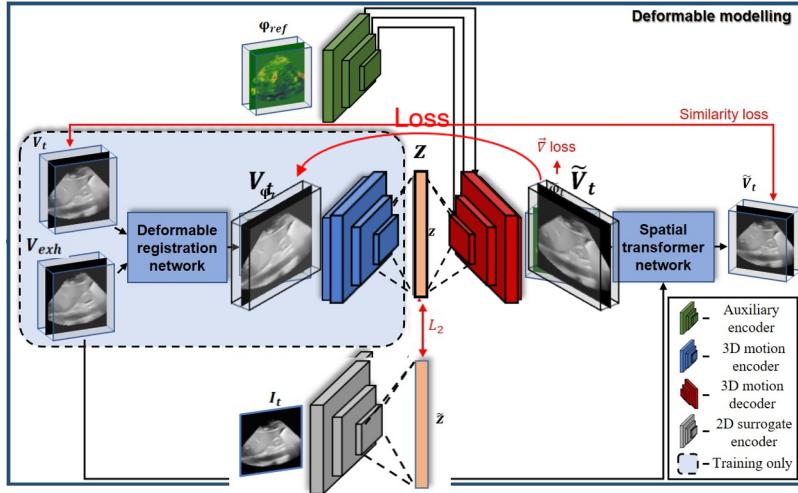


$$\mathcal{L} = \mathcal{L}_{sim}(\tilde{V}_t, V_t) + \beta \mathcal{L}_{grad}(\tilde{\varphi}_t) + \|z, \tilde{z}\|_2^2$$

1. A. V. Dalca, G. Balakrishnan, J. Guttag et al., "Unsupervised learning for fast probabilistic diffeomorphic registration," *Lecture Notes in Computer Science*, p.729–738, 2018.

2. M. Jaderberg, K. Simonyan, A. Zisserman et al., "Spatial transformer networks," in *Advances in neural information processing systems*, 2015, p. 2017–2025.

# Ablation study



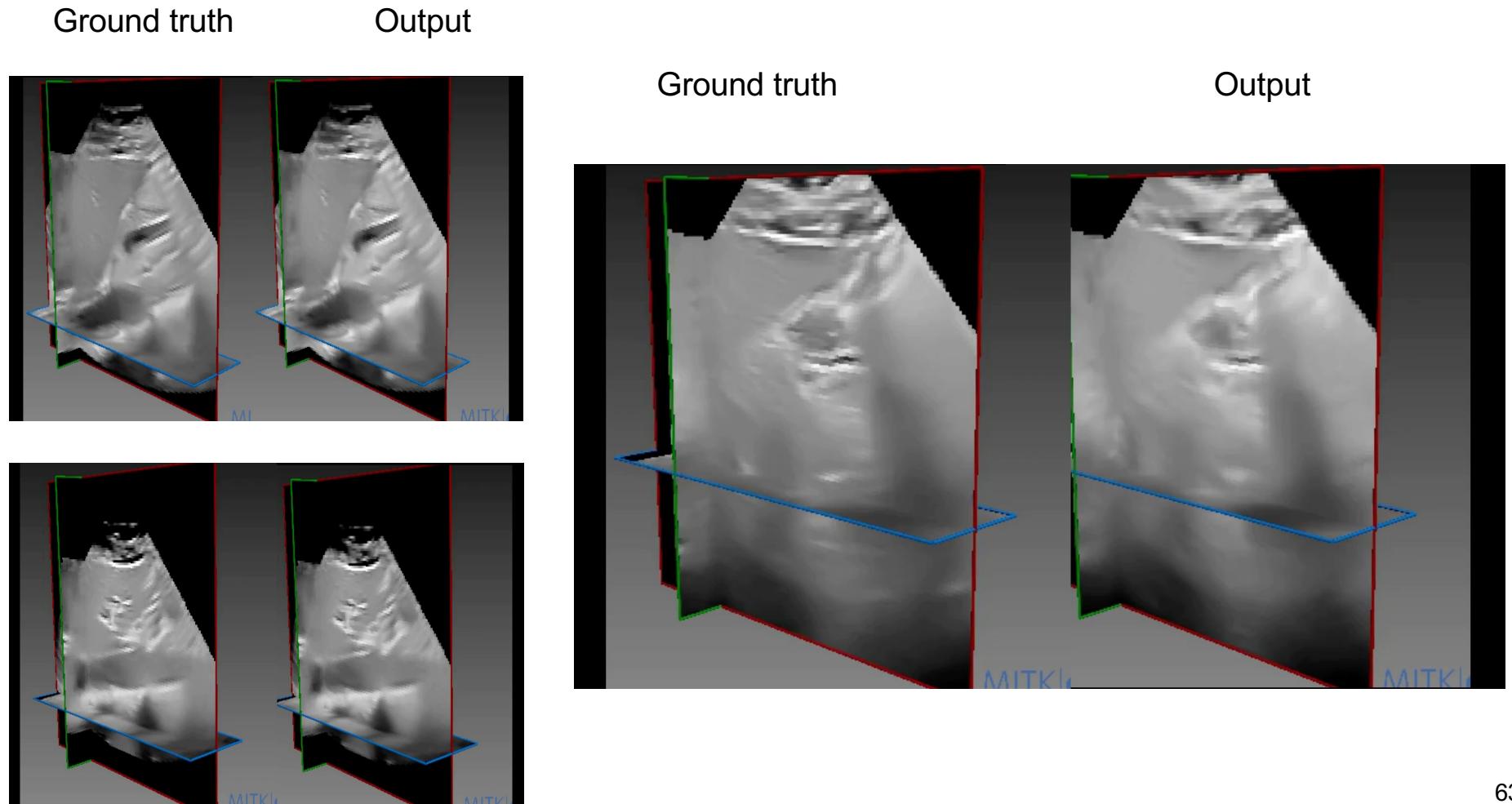
- Deformation properties
  - $|J| < 0 \Rightarrow 1.1\%$
- Inference time
  - 0.47s CPU
  - 0.09s GPU

Image similarity metrics for different model configurations.  
Values are mean  $\pm$  std

Model	MSE	NCC	SSIM
Baseline	$0.15 \pm 0.04$	$0.42 \pm 0.05$	$0.29 \pm 0.05$
Baseline + STN <sup>1</sup>	$0.10 \pm 0.06$	$0.57 \pm 0.10$	$0.54 \pm 0.11$
Baseline + STN <sup>1</sup> + $\phi_{ref}$	$0.07 \pm 0.04$	$0.62 \pm 0.10$	$0.60 \pm 0.10$
Rigid	$0.10 \pm 0.06$	$0.61 \pm 0.11$	$0.60 \pm 0.12$
Proposed (axi.)	$0.07 \pm 0.04$	$0.63 \pm 0.10$	$0.61 \pm 0.10$
<b>Proposed (sag.)</b>	<b><math>0.06 \pm 0.03</math></b>	<b><math>0.66 \pm 0.09</math></b>	<b><math>0.65 \pm 0.08</math></b>
P-value	$8 \times 10^{-4}$	$8 \times 10^{-6}$	$9 \times 10^{-8}$

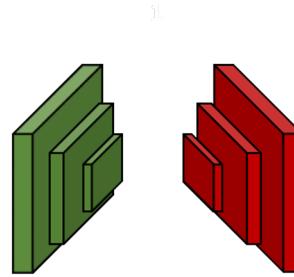
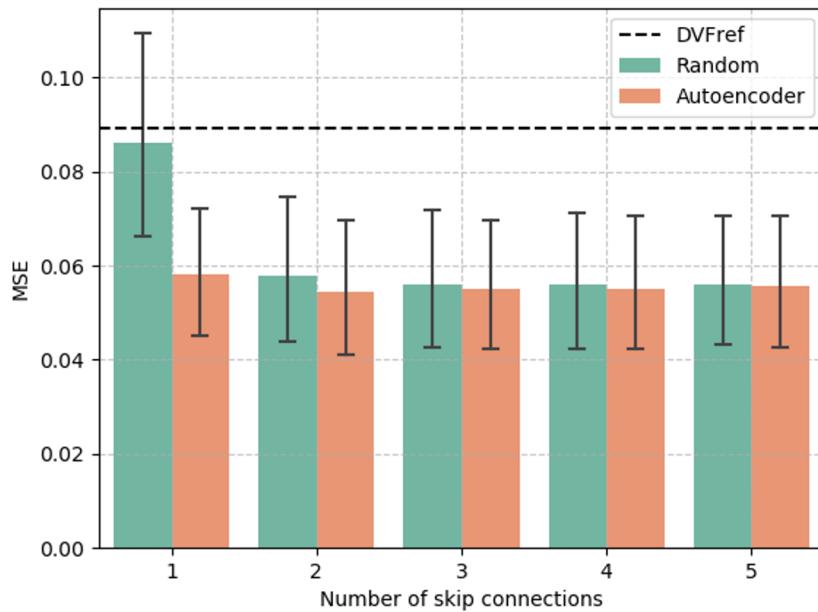
1. M. Jaderberg, K. Simonyan, A. Zisserman et al., "Spatial transformer networks," in Advances in neural information processing systems, 2015, p. 2017–2025.

# Qualitative results

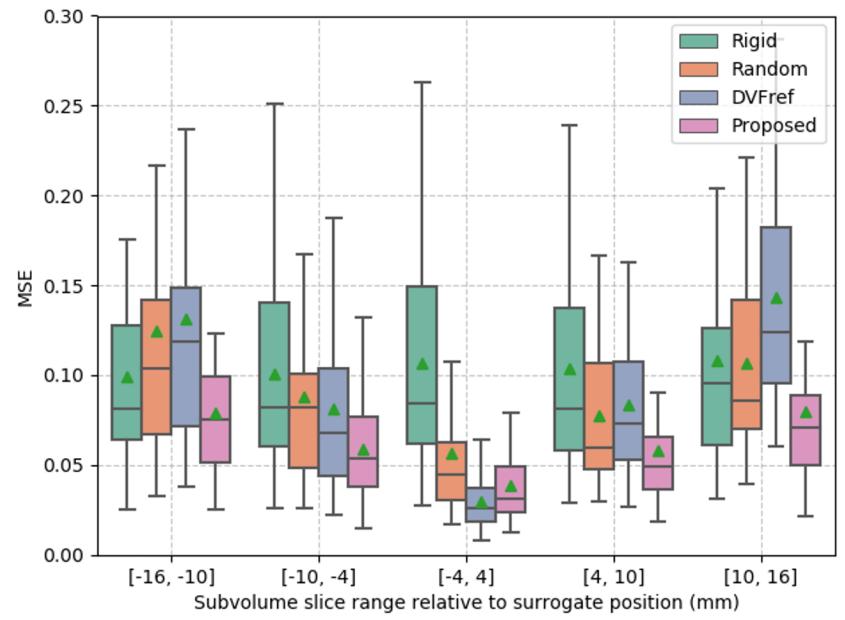


# Proposed model: Auxiliary encoder

Motion autoencoder performance when varying the number of skip connections

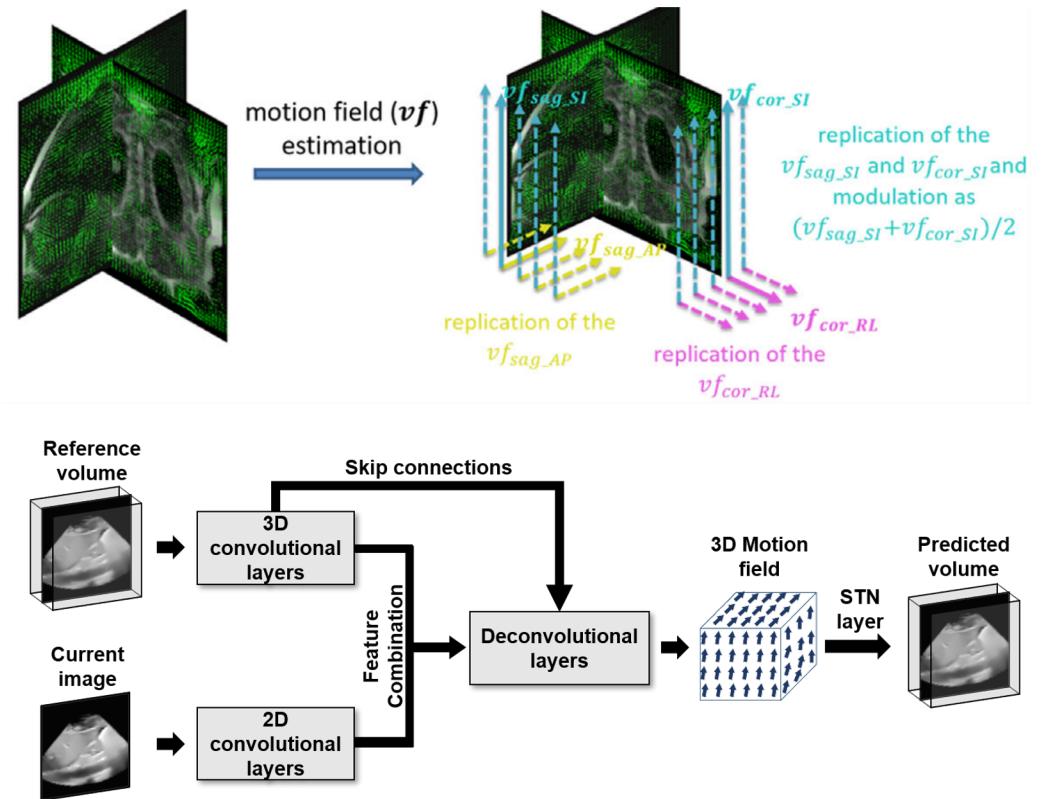


MSE distributions between ground-truth and predicted sub-volumes



# Proposed model: comparative methods

- Motion extrapolation (ME)<sup>1</sup>
  - Compute displacement between current and reference orthogonal slices
  - Extrapolate motion field to volume
- Feature combination (FC)<sup>2</sup>
  - CNN based model from preliminary work
  - Combined with rigid module



1. C. Paganelli et al. "Feasibility study on 3d image reconstruction from 2d orthogonal cine-mri for mri-guided radiotherapy," Journal of Medical Imaging and Radiation Oncology, 2018.  
 2. T. Mezheritsky et al., "3d ultrasound generation from partial 2d observations using fully convolutional and spatial transformation networks," in 2020 IEEE (ISBI), 2020

# Comparative results

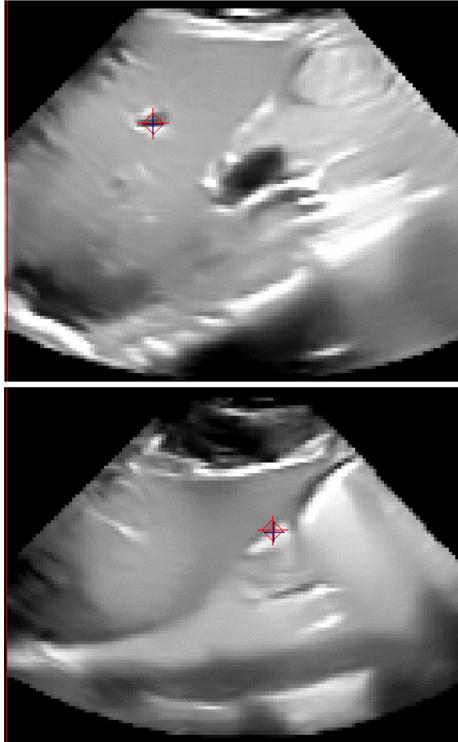
Image similarity and target tracking metrics between ground-truth and predicted volumes for different comparative methods. Values are mean $\pm$ std.

Image similarity				TRE Mean $\pm$ STD (mm)			
Model	MSE	NCC	SSIM	Exhale	Mid-cycle	Inhale	Overall
Reference	0.09 $\pm$ 0.06	0.59 $\pm$ 0.11	0.55 $\pm$ 0.13	---	9.8 $\pm$ 8.2	18.0 $\pm$ 13.4	10.7 $\pm$ 9.7
Rigid	0.10 $\pm$ 0.06	0.61 $\pm$ 0.11	0.60 $\pm$ 0.12	3.5 $\pm$ 1.3	3.9 $\pm$ 1.7	6.3 $\pm$ 4.3	4.6 $\pm$ 3.2
ME <sup>1</sup>	0.21 $\pm$ 0.08	0.59 $\pm$ 0.08	0.53 $\pm$ 0.10	<b>2.7 <math>\pm</math> 1.4</b>	5.9 $\pm$ 2.8	10.9 $\pm$ 7.9	6.5 $\pm$ 6.4
FC <sup>2</sup>	0.09 $\pm$ 0.04	0.57 $\pm$ 0.09	0.54 $\pm$ 0.10	5.0 $\pm$ 3.3	7.9 $\pm$ 4.3	13.8 $\pm$ 10.7	8.9 $\pm$ 7.5
FC+rigid	0.08 $\pm$ 0.05	0.63 $\pm$ 0.10	<b>0.63 <math>\pm</math> 0.10</b>	<b>3.1 <math>\pm</math> 0.5</b>	4.5 $\pm$ 2.2	7.2 $\pm$ 4.4	4.9 $\pm$ 3.9
<b>Proposed</b>	<b>0.06 <math>\pm</math> 0.03</b>	<b>0.66 <math>\pm</math> 0.09</b>	<b>0.65 <math>\pm</math> 0.08</b>	<b>2.8 <math>\pm</math> 1.6</b>	<b>3.2 <math>\pm</math> 0.8</b>	<b>4.5 <math>\pm</math> 2.5</b>	<b>3.5 <math>\pm</math> 2.4</b>
P-value	<b>2x10<sup>-5</sup></b>	<b>6x10<sup>-5</sup></b>	<b>4x10<sup>-3</sup></b>	<b>9x10<sup>-2</sup></b>	<b>8x10<sup>-6</sup></b>	<b>9x10<sup>-6</sup></b>	<b>2x10<sup>-6</sup></b>

1. C. Paganelli et al. "Feasibility study on 3d image reconstruction from 2d orthogonal cine-mri for mri-guided radiotherapy," Journal of Medical Imaging and Radiation Oncology, 2018.  
 2. T. Mezheritsky et al., "3d ultrasound generation from partial 2d observations using fully convolutional and spatial transformation networks," in 2020 IEEE (ISBI), 2020

# Comparative results

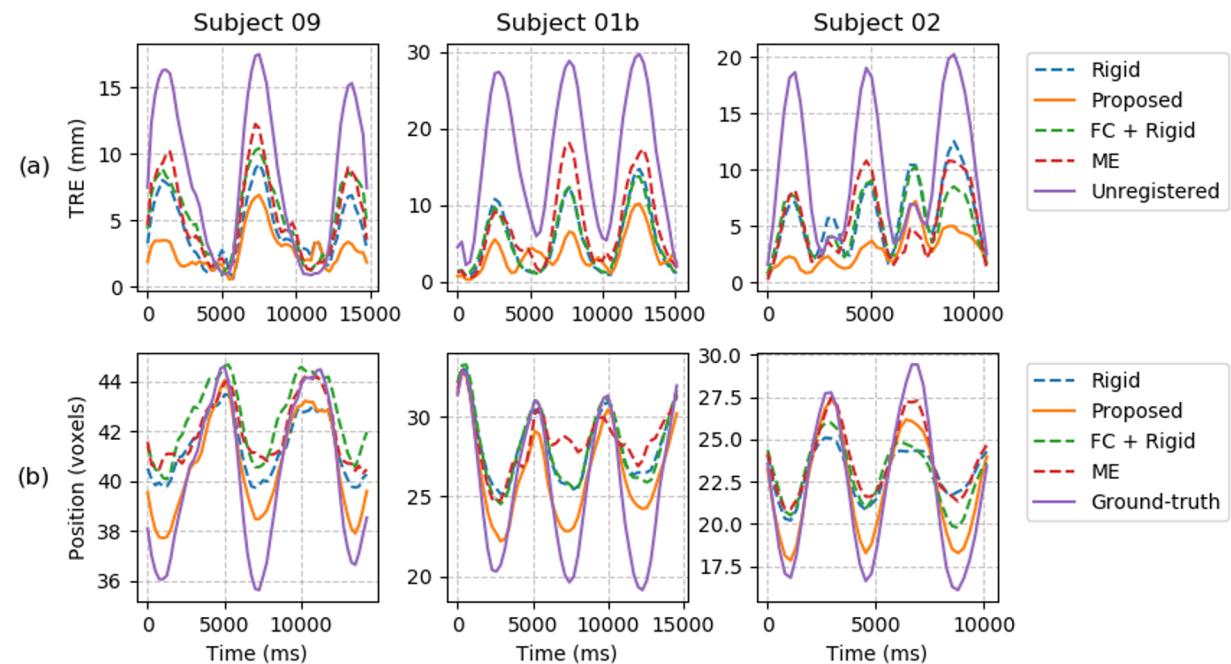
Sagittal



Coronal



Comparative plots showing (a) Evolution of TRE through time and (b) target trajectories for 3 cases over 3 respiratory cycles



- + - Ground-truth
- - Prediction

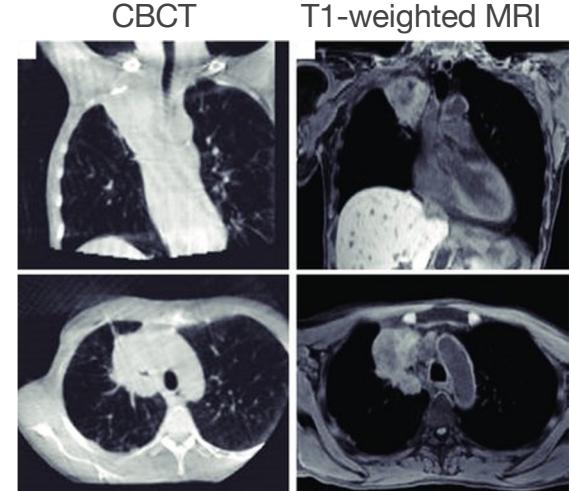
# Summary

- Proposed convolutional AE based motion modelling framework
- Framework specific experiments
  - Patient specific features
  - Split motion into rigid and deformable components
- Comparative experiments
  - Achieved average tracking error of 3.5mm
  - Better than rigid alignment alone (4.6 mm)
  - Better than comparative approaches (ME: 6.5mm, FC+rigid: 4.9)
- Rigid alignment module can be used with other models to reduce tracking errors

# MR-guided radiotherapy (IGRT)

- CBCT is the clinical standard, but increasing use of MRI
- Simultaneous beam delivery and organ monitoring with real-time 2D images → **lacks volumetric information**
- Image acquisition, target localization, beam modulation → **system latency**
- By the time a gating decision is made, the patient anatomy has already changed shape/position

**4D motion modelling framework using deep learning to allow the forecasting of temporal volumes from surrogate 2D images**



Better contrast  
non-ionizing

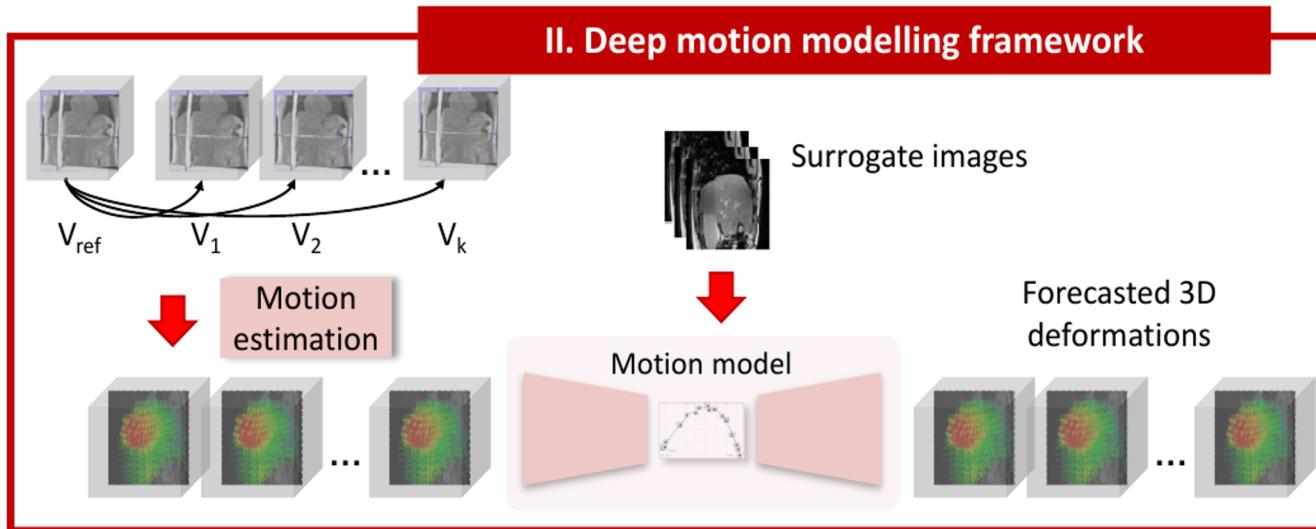


Elekta-Unity MRI-linac



Viewray - MRIdian

# Deep motion modelling framework



- Romaguera et al. International Journal of Computer Assisted Radiology and Surgery. June 2021.
- Romaguera et al. Int' Conf Medical Imaging and Computer Assisted Interventions 2021.
- Romaguera et al. Medical image analysis Dec 2021.

# Probabilistic motion model

Model creation	$\left. \begin{array}{l} \text{Ground-truth motion data} \\ \Phi = \{\phi_{t+1}, \phi_{t+2}, \dots, \phi_{t+n}\} \end{array} \right\}$	$I_{seq} = \{I_t, I_{t-1}, \dots, I_{t-m}\}$	$V_{ref}$
----------------	----------------------------------------------------------------------------------------------------------------------------------------	----------------------------------------------	-----------

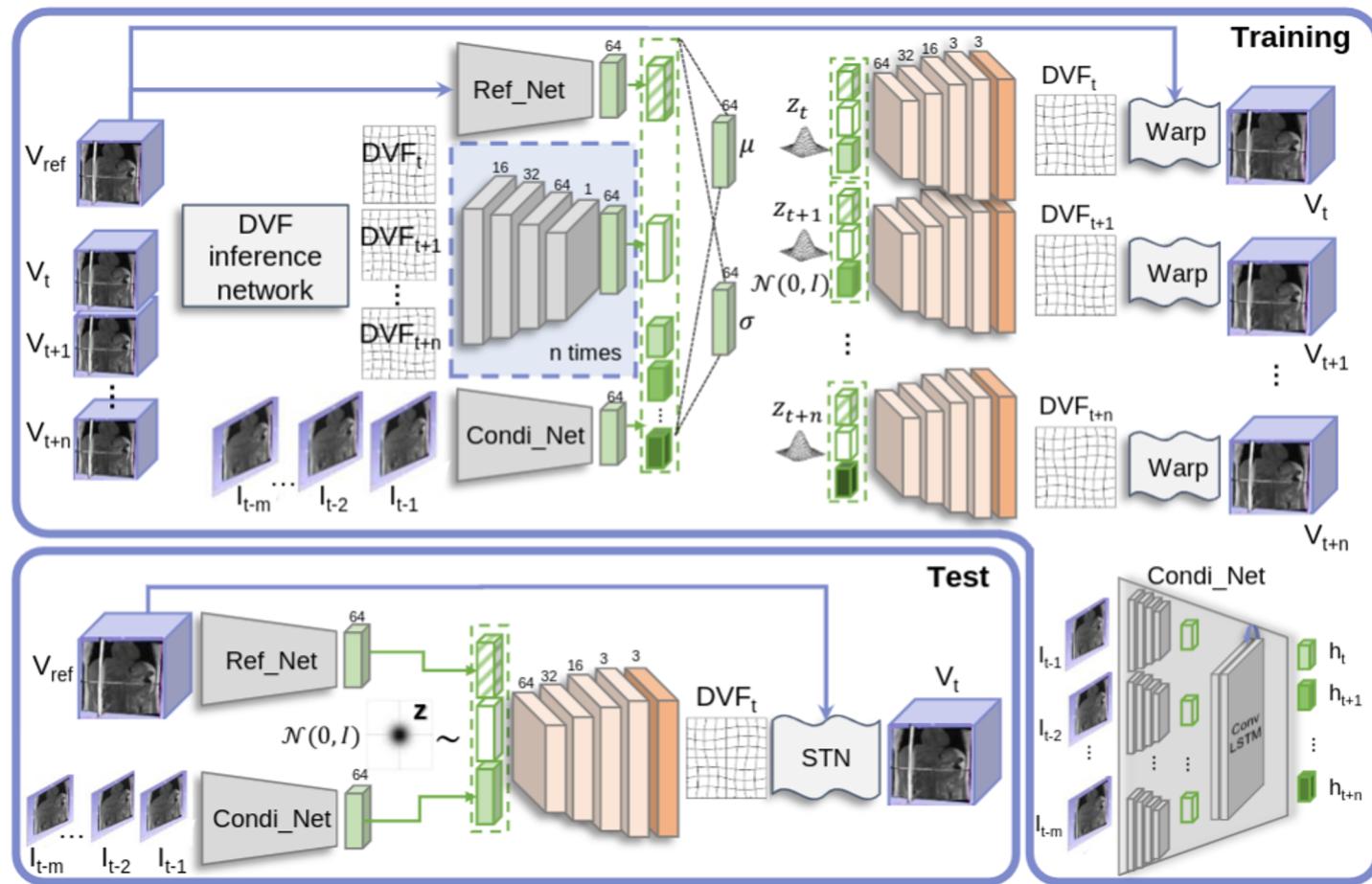
- Goal: to maximize the conditional probability distribution  $p_\theta(\Phi|I_{seq}, V_{ref})$
- Law of total probability → relates the conditional and marginal probabilities

$$p_\theta(\Phi|I_{seq}, V_{ref}) = \int_z \underbrace{p_\theta(\Phi|z, I_{seq}, V_{ref})}_{\text{Likelihood}} p(z) dz$$

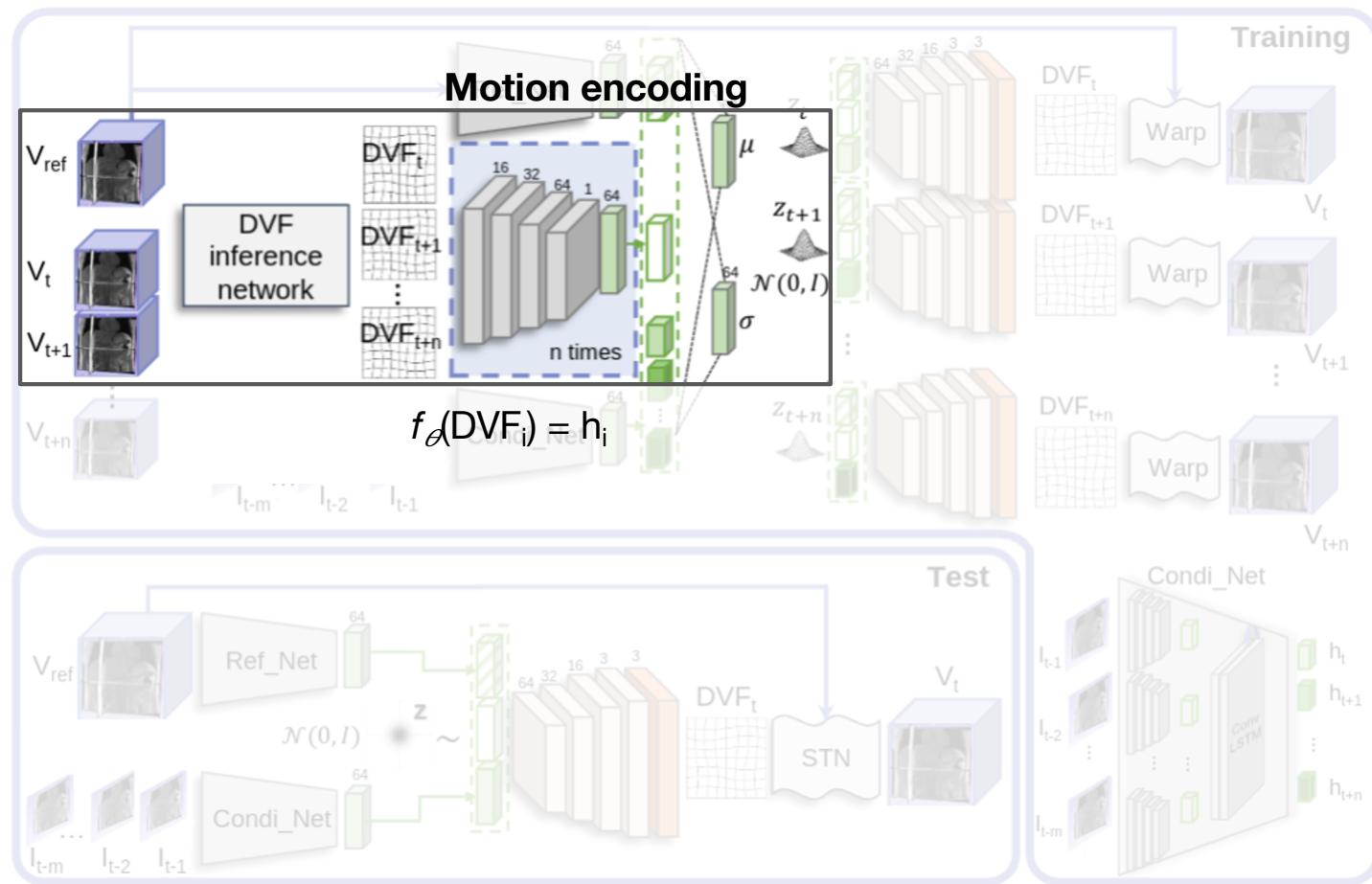
**Large # of samplings!**

- Only  $z$  values likely to produce  $p_\theta(\Phi|I_{seq}, V_{ref})$  are considered,  
i.e. the posterior  $p_\theta(z|\Phi, I_{seq}, V_{ref})$

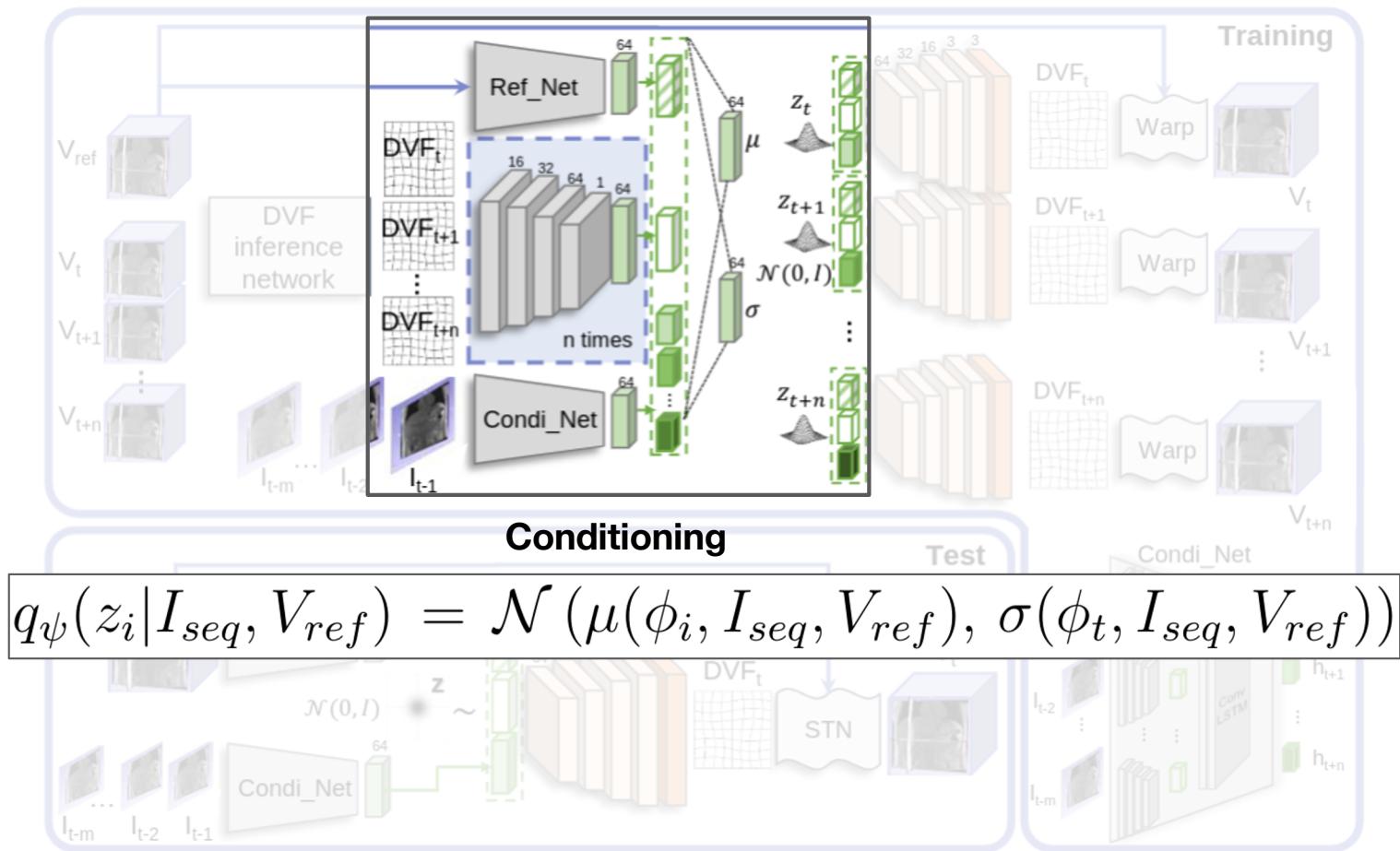
# Probabilistic motion model



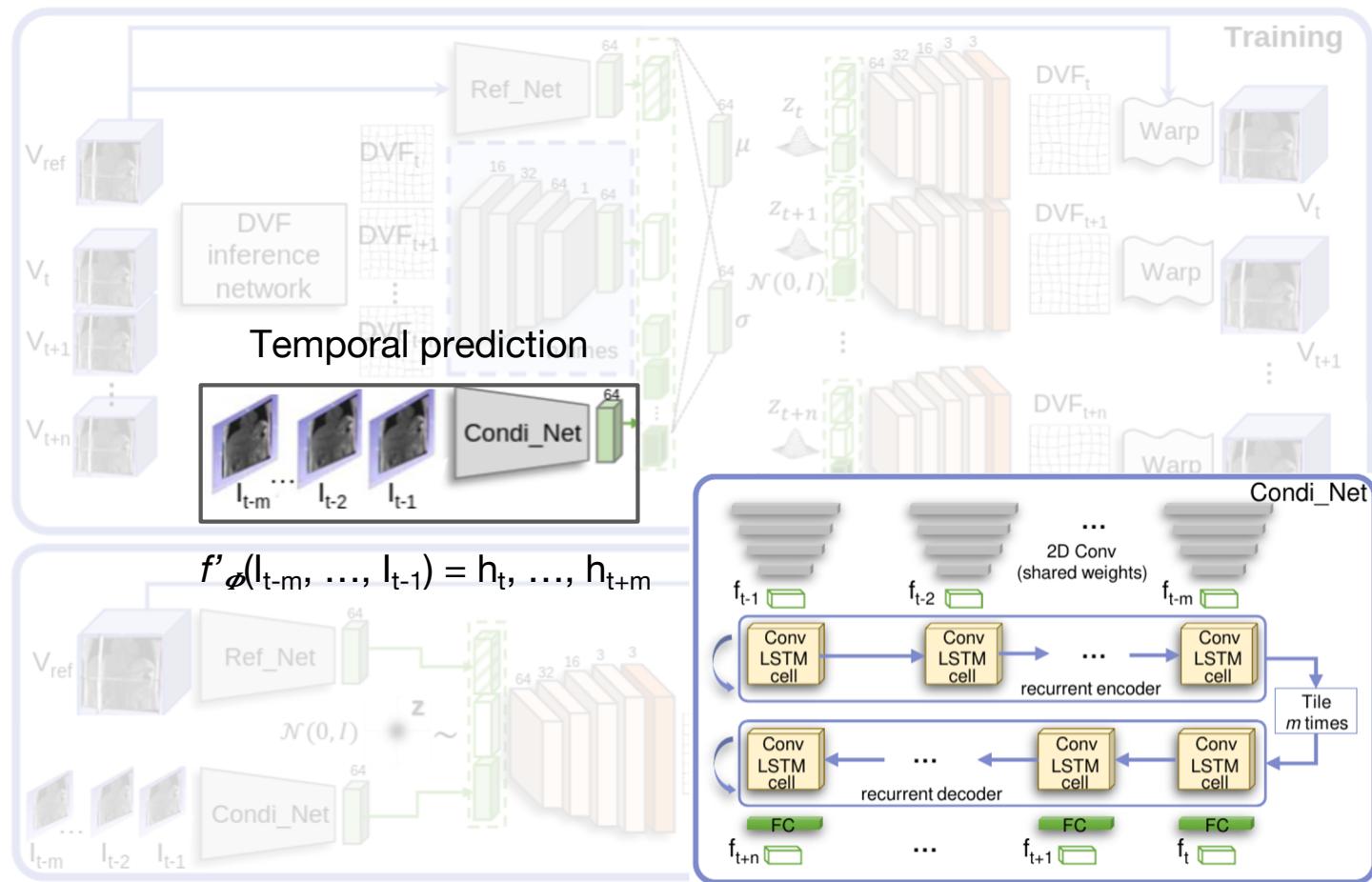
# Probabilistic motion model



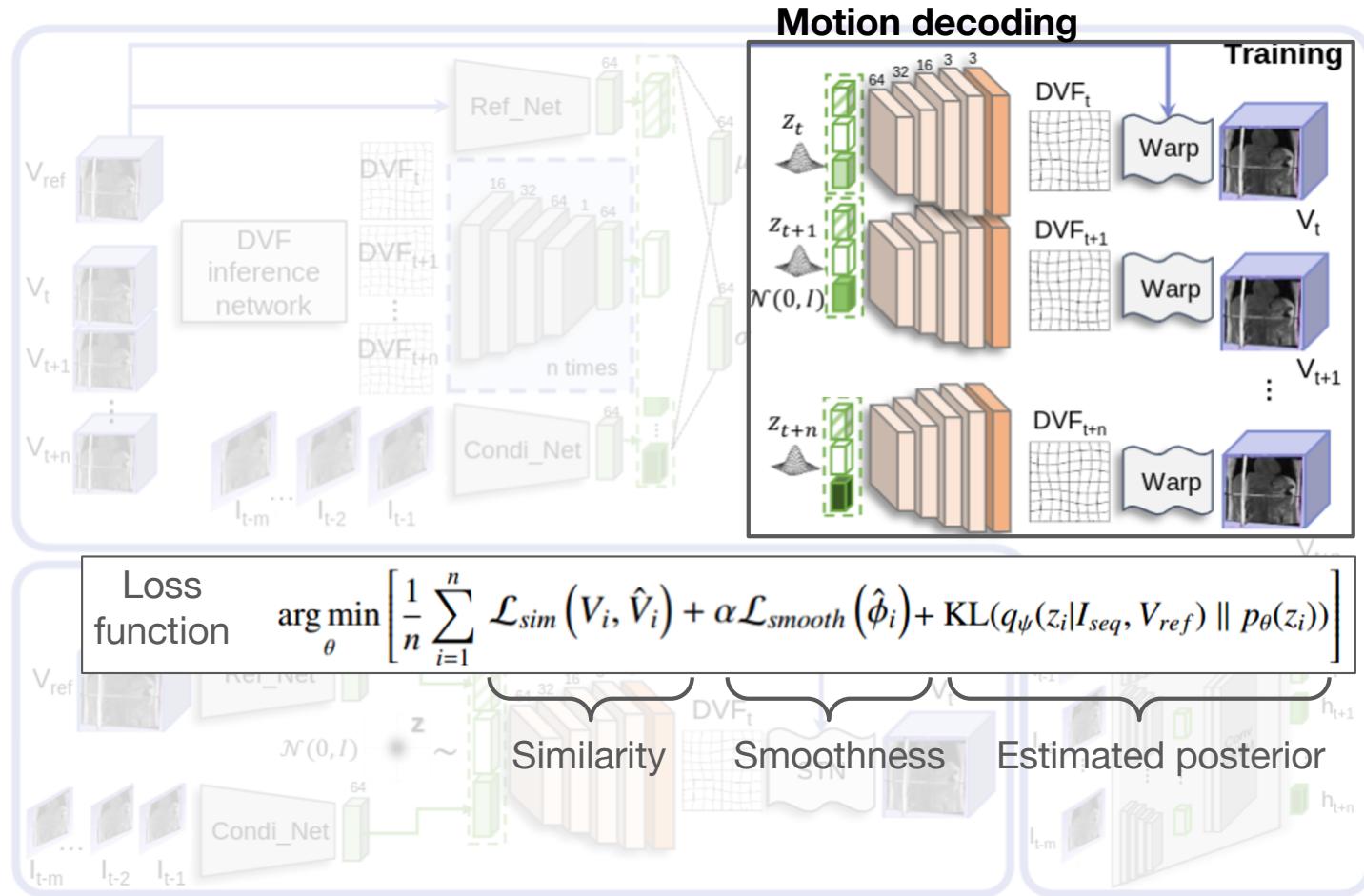
# Probabilistic motion model



# Probabilistic motion model

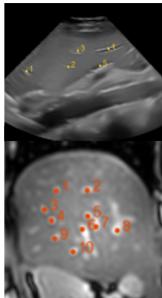


# Probabilistic motion model



# Geometrical errors

V-MRI



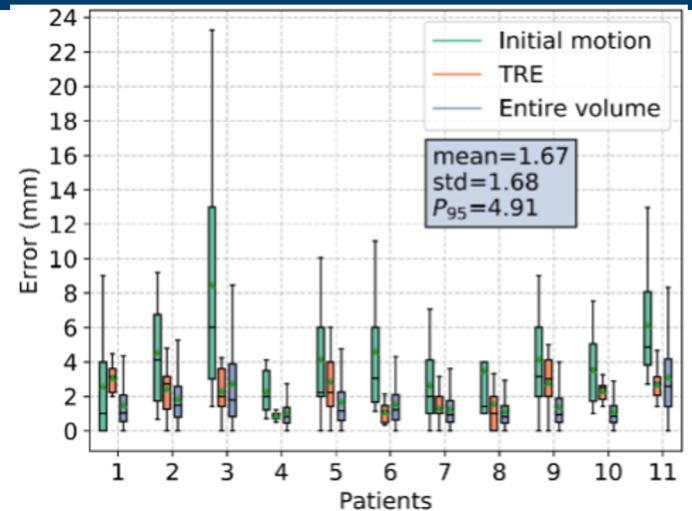
Target tracking errors (mm) for the V-MRI dataset. Values are mean  $\pm$  std (95th percentile)

Model	Overall
Initial motion	$6.7 \pm 4.4$ (13.7)
DN (Mezheritsky et al., 2020)	$3.9 \pm 2.7$ (8.2)
ME (Paganelli et al., 2018)	$2.4 \pm 2.0$ (5.8)
AE (Romaguera et al., 2021)	$2.6 \pm 2.1$ (4.8)
Proposed (sag, P)	$2.6 \pm 2.2$ (4.7)
<b>Proposed (cor, P)</b>	$2.3 \pm 1.9$ (3.6)
PCA+AB (Pham et al., 2019)	$1.8 \pm 1.6$ (3.6)
<b>Proposed (cor, SS)</b>	$1.4 \pm 1.1$ (3.3)

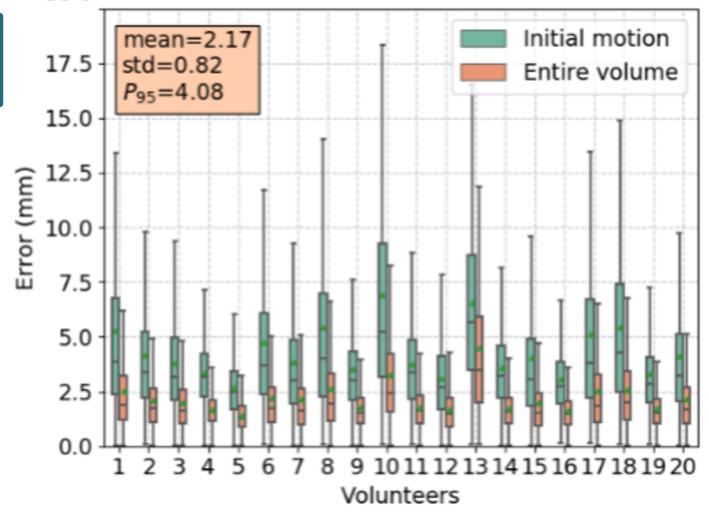
Geometrical errors (mm) for the V-MRI dataset obtained with different predictive mechanisms. Values are mean  $\pm$  std ( $P_{95}$ )

Predictor	$\Delta t = 450$ ms	$\Delta t = 900$ ms	$\Delta t = 1350$ ms
ConvGRU	$1.6 \pm 0.9$ (3.7)	$1.7 \pm 1.1$ (3.9)	$1.3 \pm 1.0$ (3.2)
3D Conv	$1.4 \pm 1.0$ (3.2)	$1.6 \pm 1.2$ (4.2)	$1.3 \pm 0.9$ (3.3)
<b>ConvLSTM</b>	<b><math>1.2 \pm 0.6</math> (2.6)</b>	<b><math>1.4 \pm 0.9</math> (3.3)</b>	<b><math>1.3 \pm 0.9</math> (3.1)</b>

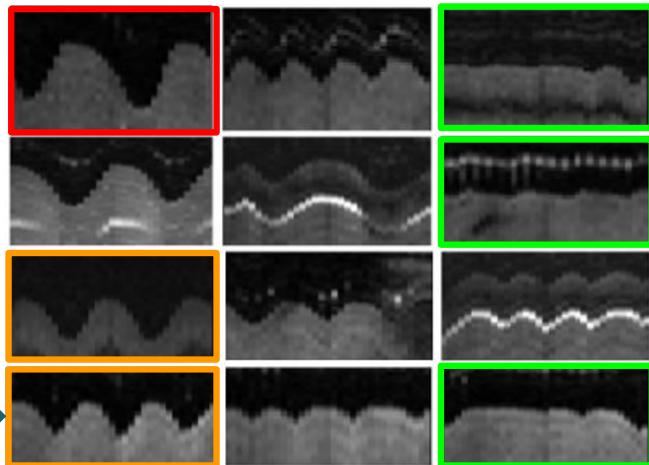
P-MRI



US



# Inter-cycle variability



Different amplitudes and frequencies

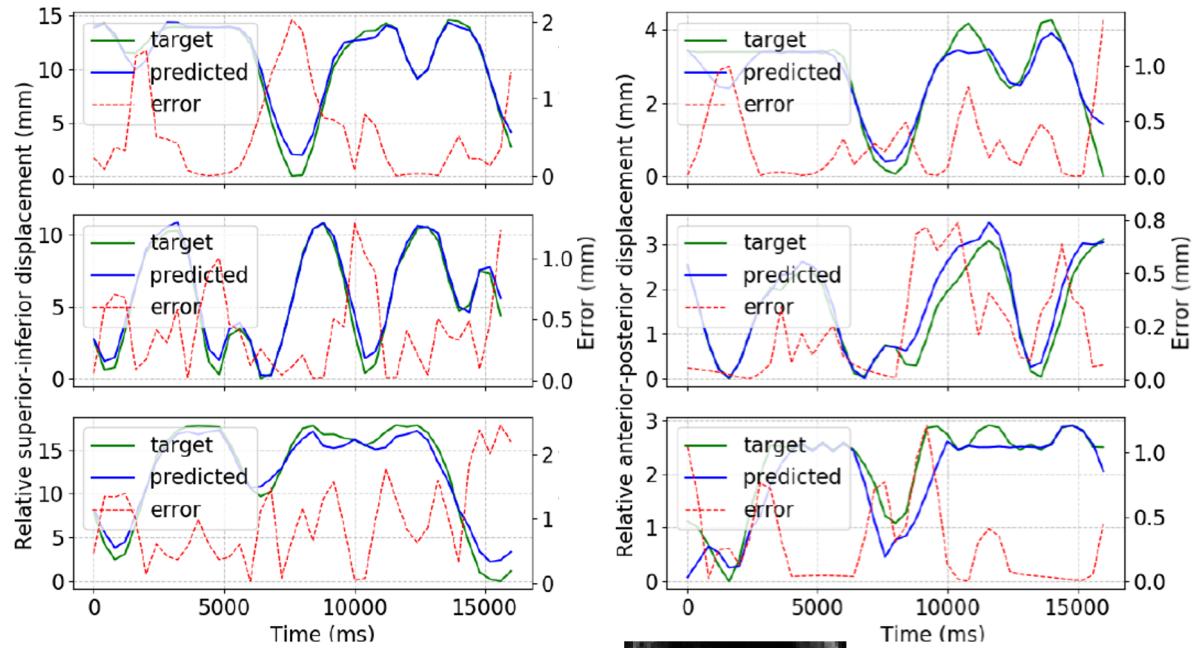
Shallow

Regular

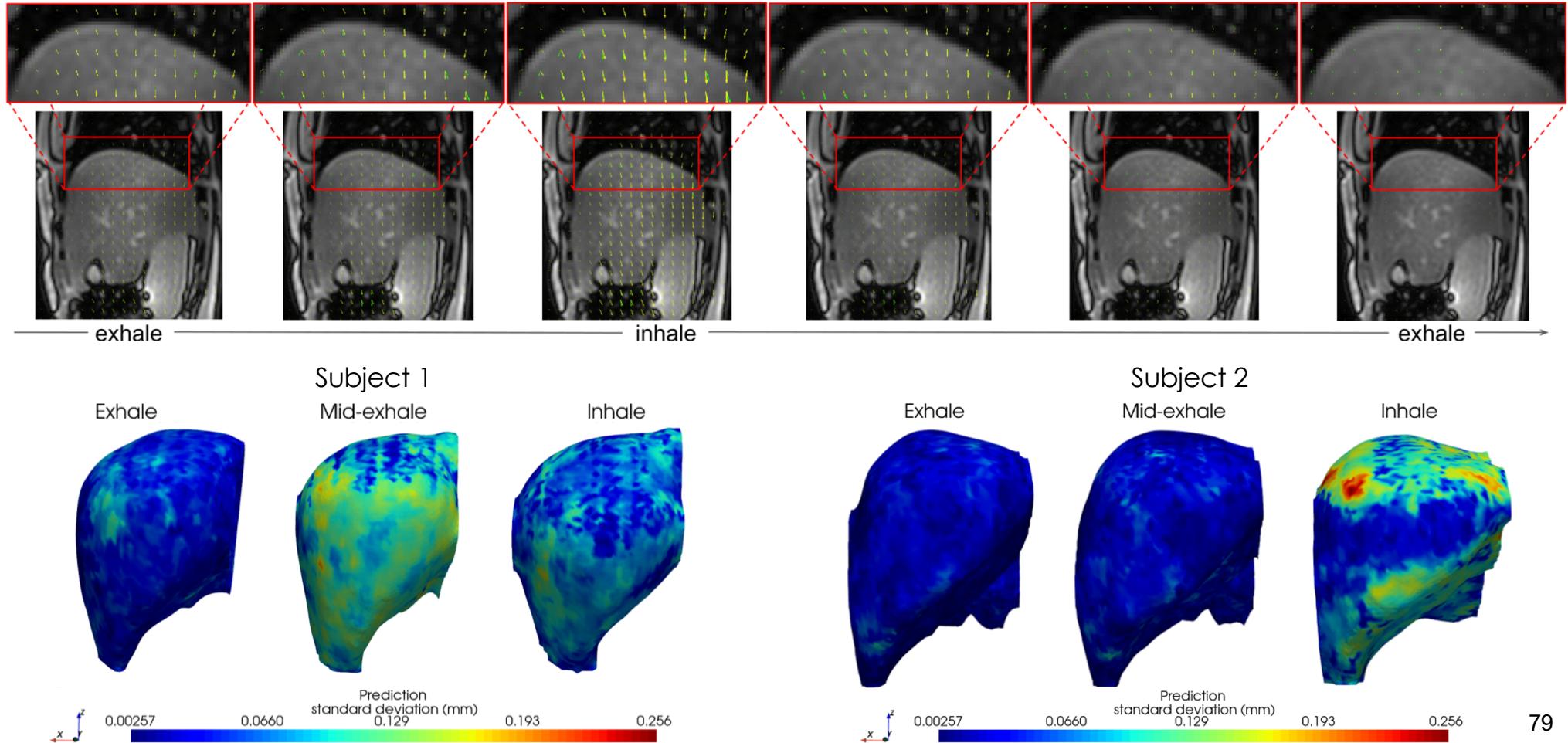
Deep

Same  
subject

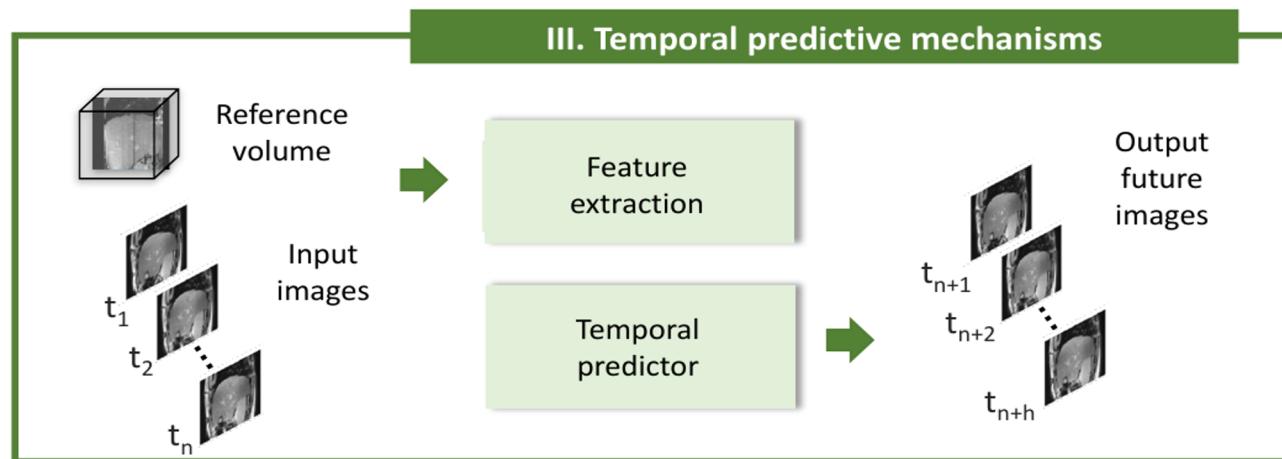
Vessel trajectories observed in 3 subjects with irregular breathing



# Sampling the latent space



# Temporal predictive mechanisms



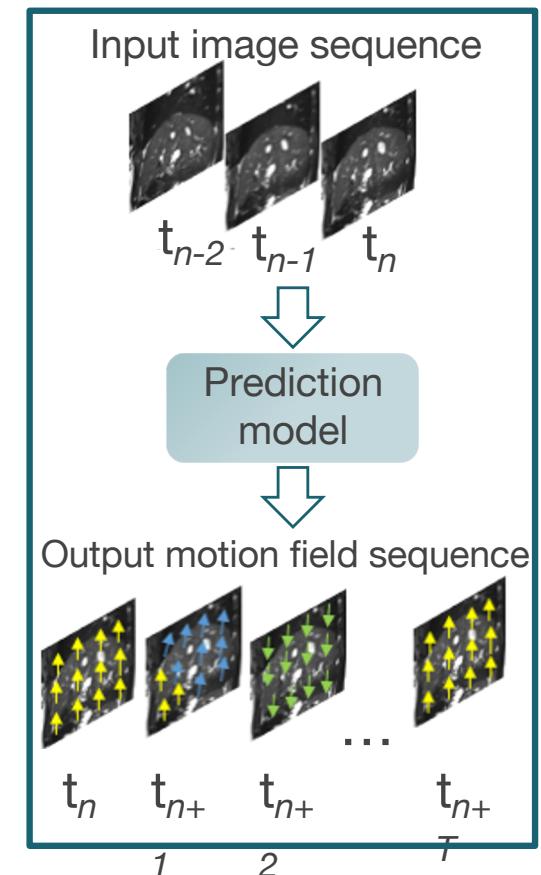
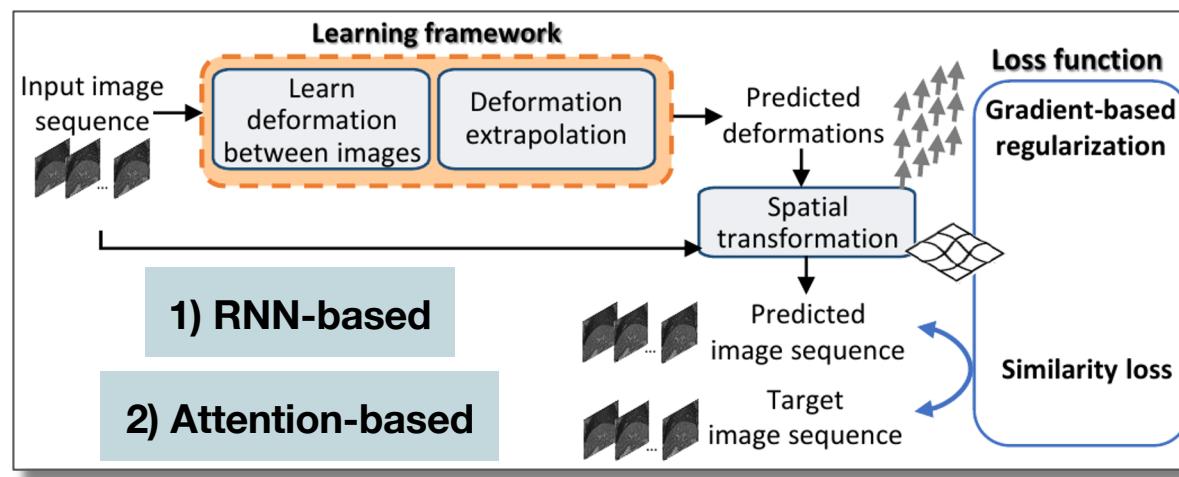
- Romaguera et al.. Medical image analysis, Dec. 2020
- Romaguera et al.. IEEE Trans. Medical Imaging. 2022.

# Temporal prediction

The goal is to predict a motion field sequence  $\mathbf{Y}_{\text{out}}$  in  $n$  future times:

$$\mathbf{Y}_{\text{out}} = [\Phi_t, \Phi_{t+1}, \dots, \Phi_{t+n}]$$

given an input image sequence  $\mathbf{X}_{\text{in}} = [I_{t-m}, \dots, I_{t-1}, I_t]$  in  $m$  observed times



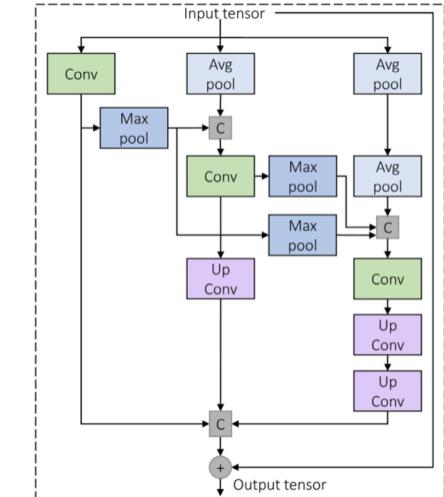
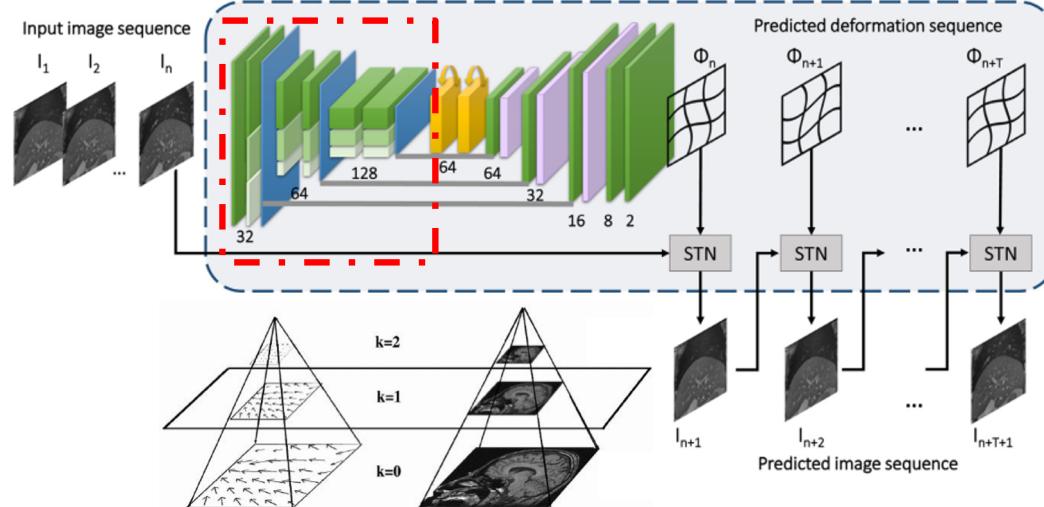
# RNN-based proposed model

- Pixel intensities → Motion fields
- Pyramidal feature extraction
- Spatio-temporal prediction (ConvLSTM)
- Spatial transformations to implicitly regress future images

Task: To maximize  $P(\mathbf{Y}|\mathbf{X})$

$$\mathbf{X} = \langle I_1^i, I_2^i, \dots, I_n^i \rangle \quad \mathbf{Y} = \langle I_{n+1}^o, I_{n+2}^o, \dots, I_{n+T+1}^o \rangle$$

$$\mathcal{L}_{total} = \frac{1}{T} \sum_{k=1}^T \mathcal{L}_{sim}(I_k^t, I_k^p) + \lambda \mathcal{L}_{smooth}(\phi_k)$$



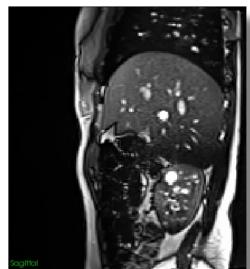
# RNN-based | Experiments

- Comparison
  - Principal Component Analysis
  - Classification-based approach

- Datasets

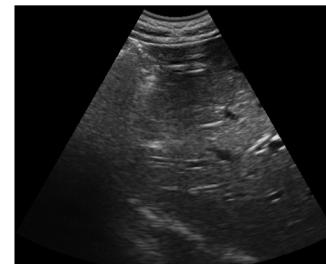
## MRI (12 volunteers)

Free-breathing sagittal slices  
Siemens Skyra 3T scanner  
2D T2w true FISP sequence  
Spatial res. =  $1.7 \times 1.7 \text{ mm}^2$   
Temporal res. ~ 320 ms  
Leave-one-out validation



## US (63 volunteers)

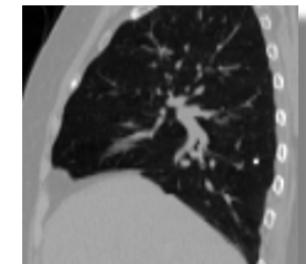
Multicenter dataset  
63 free-breathing sequences  
Spatial res. =  $0.5 \times 0.5 \text{ mm}^2$   
Temporal res. = 200 ms  
Split: 62% - 13% - 25%  
(train-val-test)



<https://clust.ethz.ch/data.html>

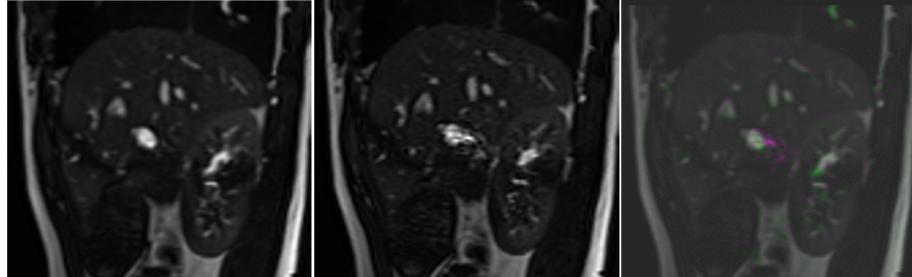
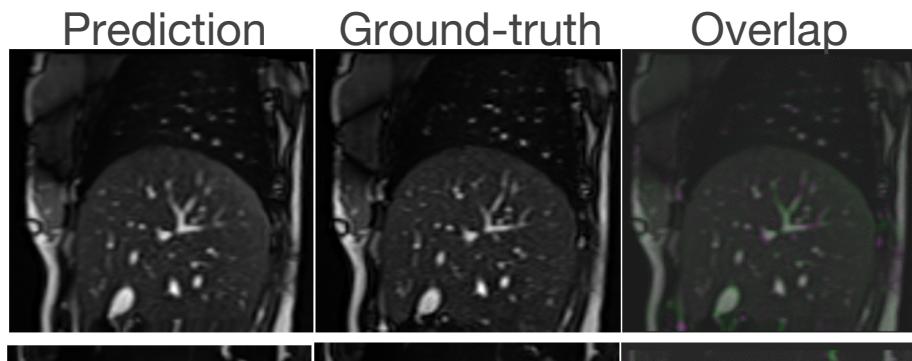
## CT (10 cancer patients)

4D thoracic planning CT  
Slices covering the right hemidiaphragm  
Spatial res. =  $1.0 \times 1.0 \text{ mm}^2$   
Temporal res. ~ 400 ms  
Leave-one-out validation

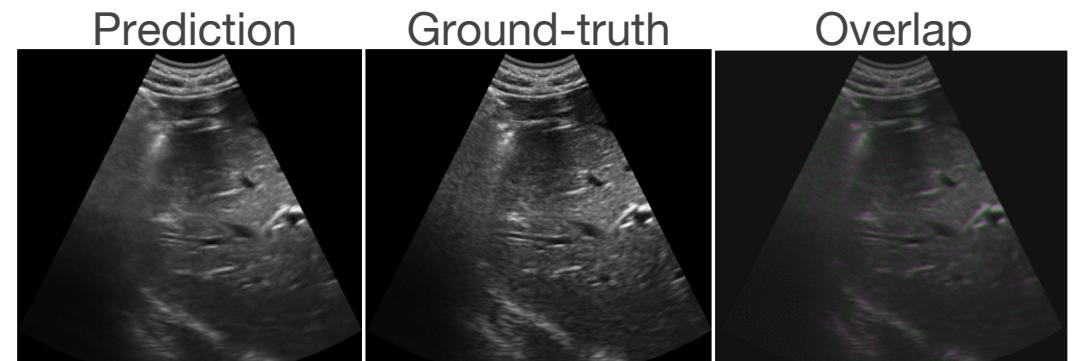


<http://www.dir-lab.com>

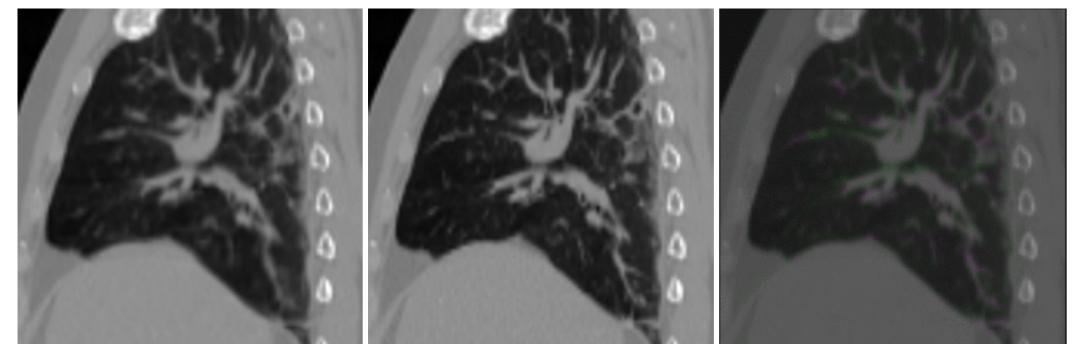
# RNN-based | Qualitative results



horizon = 320 ms



horizon = 200 ms



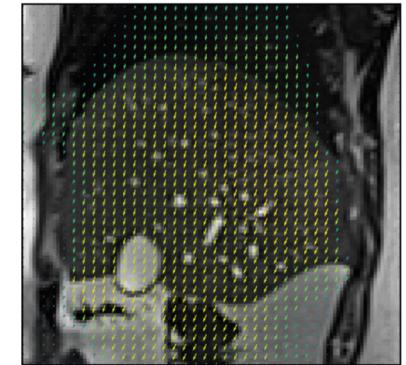
horizon = 400 ms

target    prediction

# RNN-based | Tracking results

Vessel tracking error position (in mm) in the MRI dataset; median (IQR)

Model	$t_p=1$ (320 ms)	$t_p=2$ (640 ms)	$t_p=3$ (960 ms)	$t_p=4$ (1280 ms)	$t_p=5$ (1600 ms)
Classification	1.63 (2.29)	2.32 (2.58)	2.96 (2.88)	3.23 (2.65)	3.55 (2.52)
Classif. (AG)	1.55 (1.45)	2.33 (2.10)	2.77 (2.64)	3.16 (2.63)	3.20 (2.82)
PCA	1.36 (2.73)	1.85 (2.98)	2.37 (2.88)	2.72 (2.67)	3.01 (2.49)
ED-ST(ncc)	0.54 (0.66)	0.74 (0.98)	1.03 (1.26)	1.17 (1.42)	1.30 (1.66)
MSED-ST(ncc)	0.43 (0.54)	0.72 (0.91)	0.88 (1.22)	1.01 (1.36)	1.21 (1.57)
ED-ST(mse)	0.56 (0.65)	0.77 (0.96)	0.94 (1.15)	1.00 (1.15)	1.28 (1.61)
<b>MSED-ST(mse)</b>	<b>0.45 (0.55)</b>	<b>0.57 (0.75)</b>	<b>0.80 (0.99)</b>	<b>0.88 (1.25)</b>	<b>0.77 (1.36)</b>

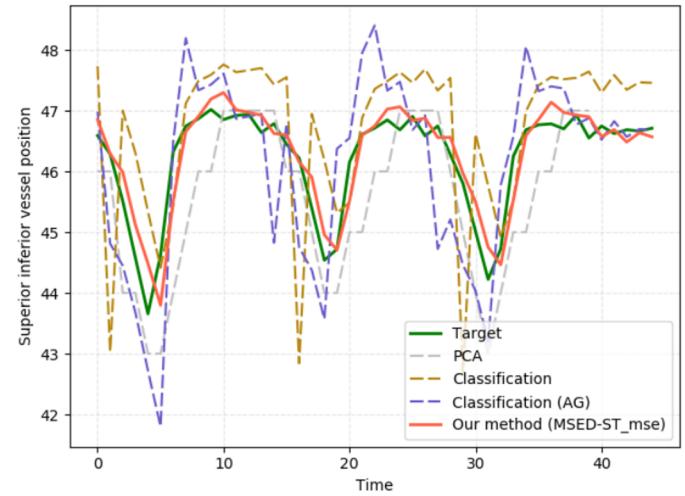


US dataset

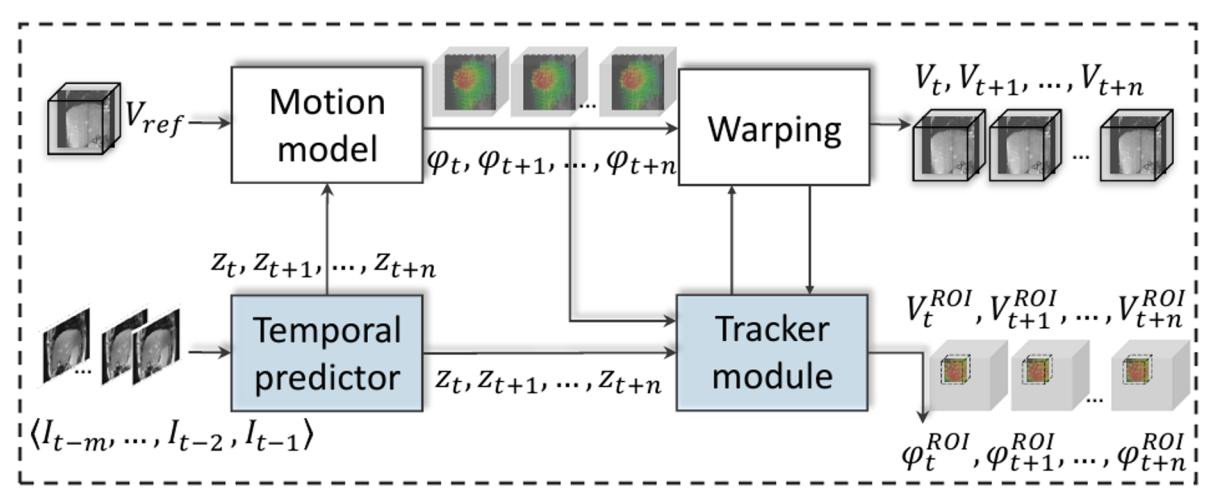
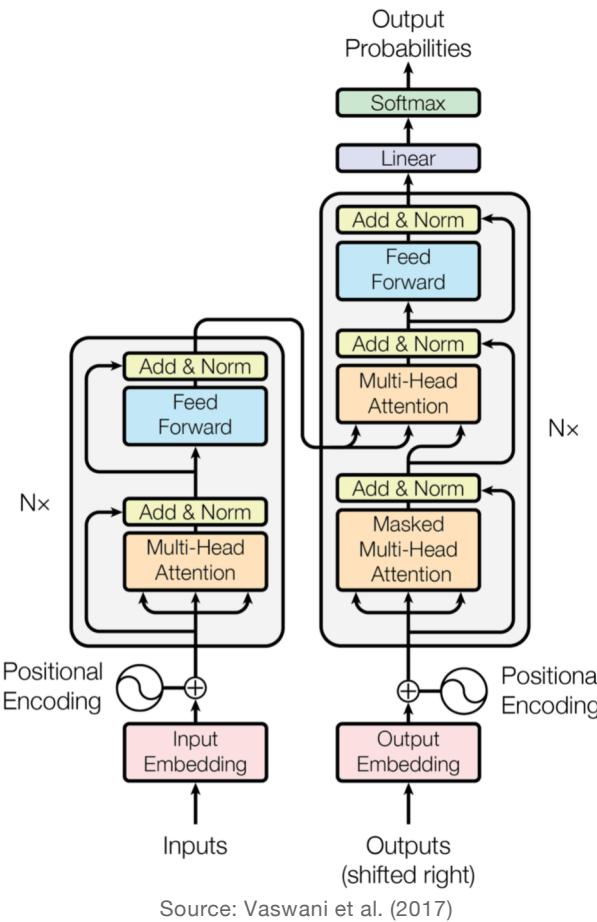
0.45 (0.74) mm < error < 1.28 (1.31) mm

CT dataset

0.28 (0.58) mm < error < 0.42 (0.49) mm

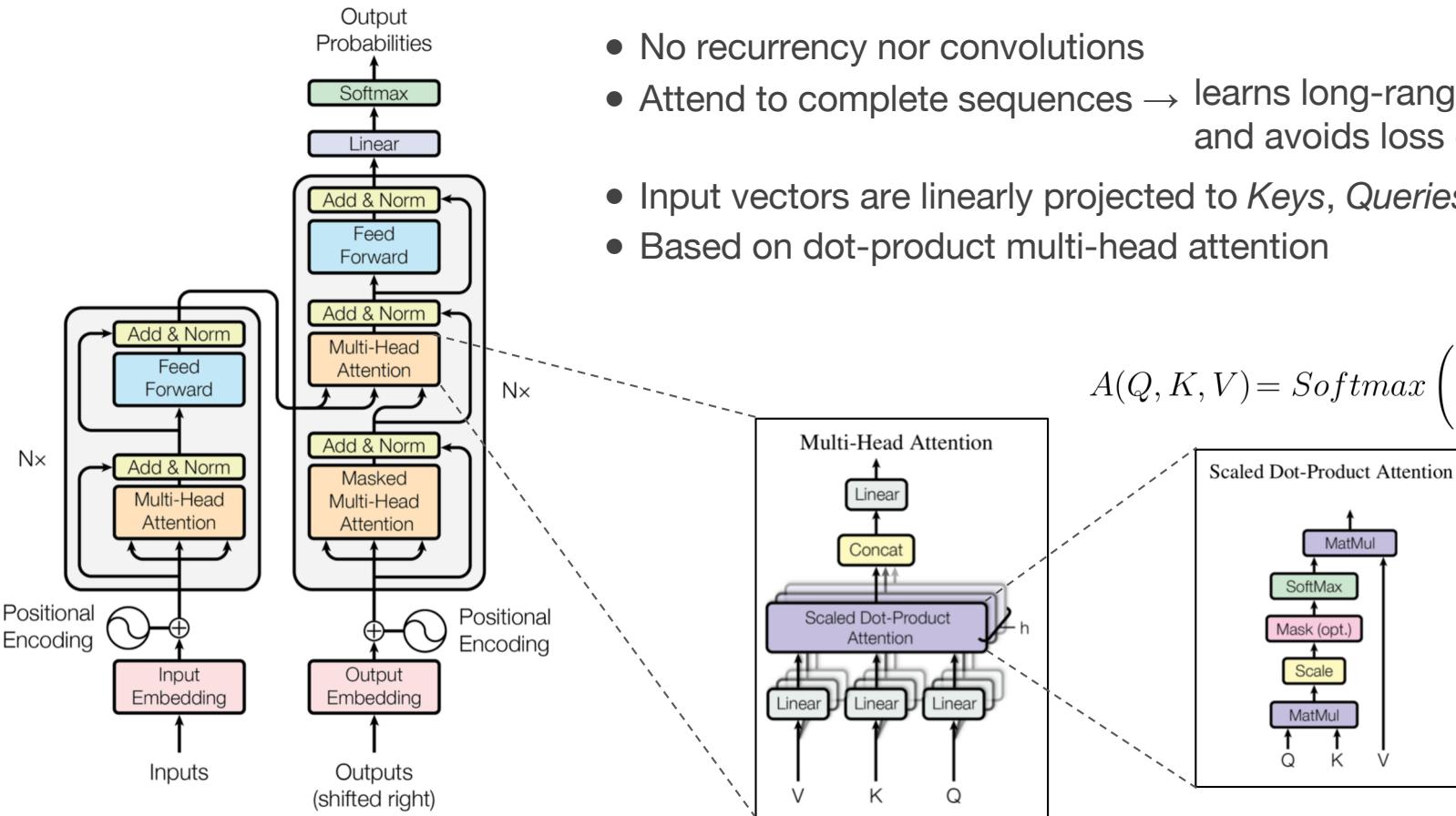


# Attention-based 4D deformation forecasting and tracking



- Potential use of **attention models** for temporal prediction and its integration within the 4D motion model
- 4D motion model for **local tracking**

# Transformer

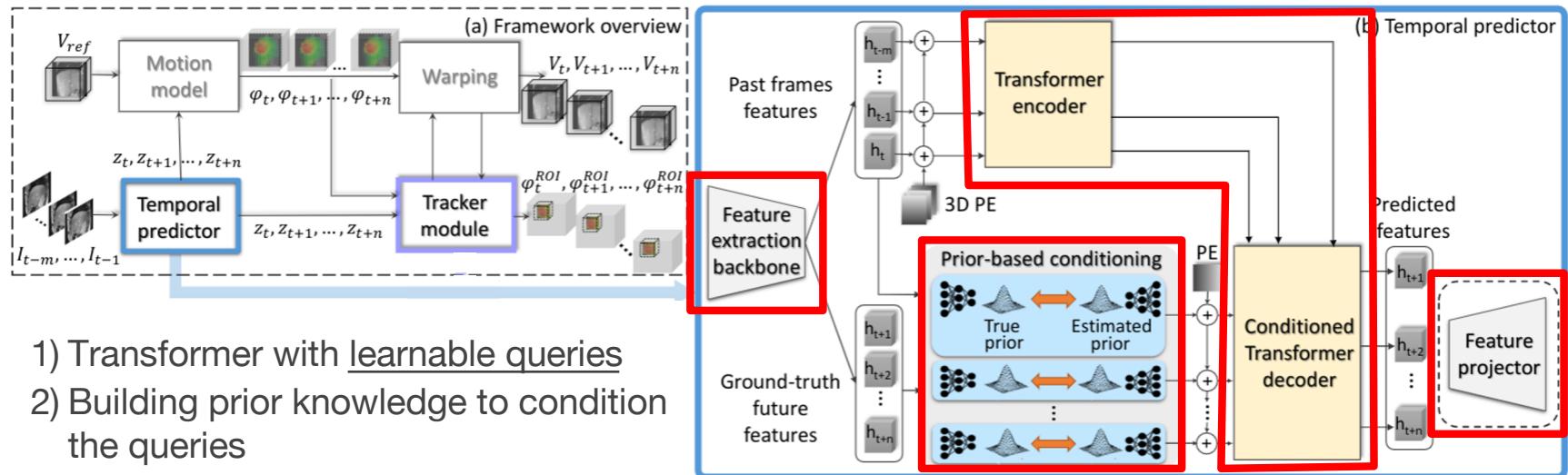


- No recurrence nor convolutions
- Attend to complete sequences → learns long-range dependencies and avoids loss of information
- Input vectors are linearly projected to *Keys*, *Queries* and *Values*
- Based on dot-product multi-head attention

$$A(Q, K, V) = \text{Softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

Source: Vaswani et al. (2017)

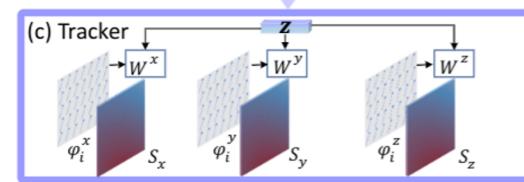
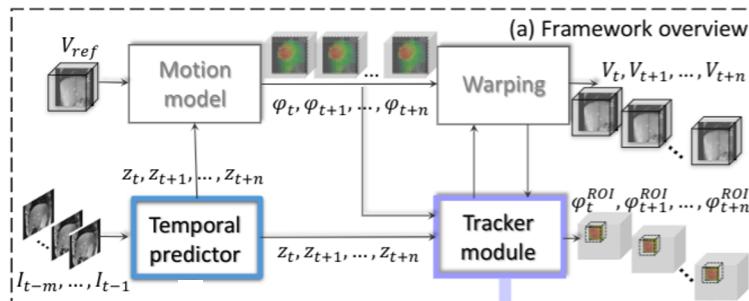
# Attention-based deformation forecasting



$$\arg \min_{\theta} \left[ \sum_{i=1}^n \mathcal{L}_{rec} + \alpha \mathcal{L}_{KL} \right]$$

$$\mathcal{L}_{KL} = \sum_{i=0}^n KL \underbrace{\left[ p_{\theta}(h_{t+i}|z_{t-m:t+i}) \right]}_{\text{True prior}} \underbrace{\left[ r_{\psi}(h_{t+i}|z_{t-m:t-1}) \right]}_{\text{Estimated prior}}$$

# Tracker



Landmark selection  $(x_{ref}, y_{ref}, z_{ref})$

Bounding box

$$(x_{ref} - \frac{\Delta_x}{2}, y_{ref} - \frac{\Delta_y}{2}, z_{ref} + - \frac{\Delta_z}{2})$$

Basic idea: Get refined values from  $\phi_t$  conditioned on  $z_t$

$$\hat{\phi}^{ROI} = \text{Concat}(S_x \times \phi_x^{ROI}, S_y \times \phi_y^{ROI}, S_z \times \phi_z^{ROI})$$

forecasted embedding → carries semantic information on the respiratory phase

Dense deformation predicted by the model

$$S_i = \sigma_2(\sigma_1(W_c c + W_\phi \phi_i^{ROI}) W_s)$$

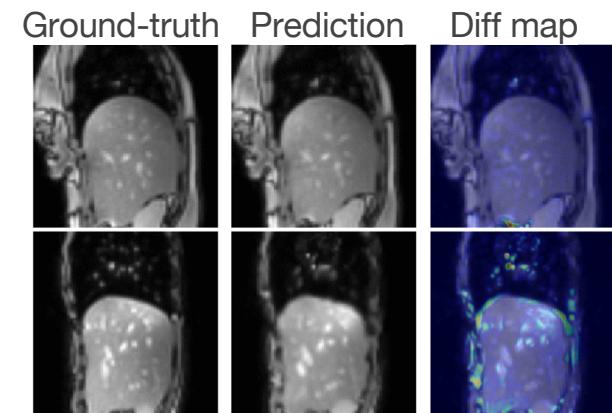
# Deformation forecasting | Results

Geometrical errors (in mm) and image similarity measures obtained with the motion model and different temporal predictors. Values are mean  $\pm$  std (95<sup>th</sup> percentile).

4D prediction	Predictive module	TRE (mm)	NCC	MSE	Time (ms)
	ConvGRU [37]	1.60 $\pm$ 1.09 (3.17)	0.71 $\pm$ 0.14 (0.89)	0.16 $\pm$ 0.09 (0.32)	11.1 $\pm$ 1.2
	ConvLSTM [36]	1.37 $\pm$ 0.92 (2.60)	0.76 $\pm$ 0.13 (0.91)	0.13 $\pm$ 0.09 (0.22)	11.9 $\pm$ 1.4
	Transformer	1.34 $\pm$ 0.87 (2.51)	0.76 $\pm$ 0.13 (0.91)	0.14 $\pm$ 0.09 (0.23)	27.2 $\pm$ 1.2
	<b>Transformer+prior</b>	<b>1.25 <math>\pm</math> 0.74(2.13)</b>	<b>0.81 <math>\pm</math> 0.11(0.95)</b>	<b>0.10 <math>\pm</math> 0.07(0.18)</b>	<b>29.5 <math>\pm</math> 1.6</b>

Comparison with related approaches. Values are mean (std.)

2D + t prediction	Method	PSNR	SSIM	MSE
	SVG [10]	22.6(3.1)	0.71(0.17)	0.09(0.06)
	R-Unet [38]	25.4(5.5)	0.77(0.13)	0.08(0.04)
	LMC [25]	23.5(2.3)	0.71(0.13)	0.10(0.04)
	<b>Proposed</b>	<b>26.3(4.5)</b>	<b>0.78(0.11)</b>	<b>0.07(0.05)</b>

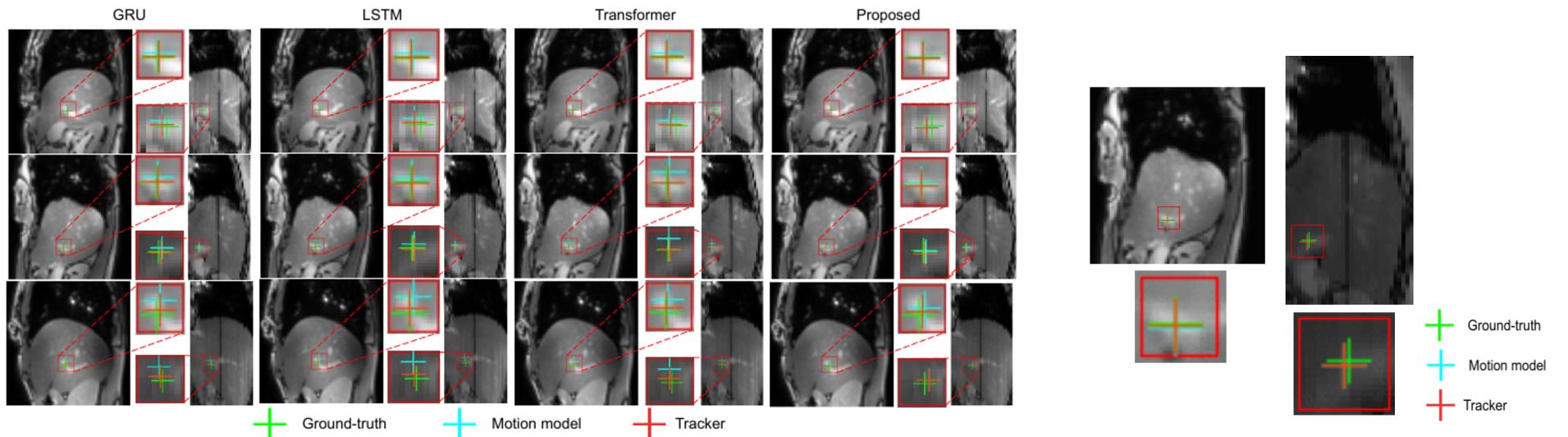


# 4D Tracking | Results

Tracker module  
reduces the error by  
63% compared with  
the 4D motion model

Tracking errors (in mm) using the motion model. Values are mean  $\pm$  std [ $P_{90}$ ]

Method	TRE (horizon = 450 ms)	TRE (horizon = 900 ms)	TRE (horizon = 1350 ms)
Initial motion	$6.52 \pm 3.41 [8.19]$	$6.35 \pm 3.11 [8.0]$	$6.42 \pm 3.40 [8.23]$
MM + GRU	$2.65 \pm 1.93 [5.47]$	$2.72 \pm 1.89 [5.45]$	$2.66 \pm 1.86 [5.38]$
MM + LSTM	$2.68 \pm 1.73 [4.90]$	$2.66 \pm 1.70 [4.81]$	$2.59 \pm 1.66 [4.68]$
MM + Transf.	$2.61 \pm 1.58 [4.80]$	$2.56 \pm 1.57 [4.73]$	$2.54 \pm 1.55 [4.67]$
MM + Transf. + prior	$2.55 \pm 2.11 [6.22]$	$2.56 \pm 1.45 [4.70]$	$2.60 \pm 2.08 [6.15]$
MM + GRU + tracker	$1.75 \pm 1.19 [3.17]$	$1.78 \pm 1.19 [3.19]$	$1.77 \pm 1.17 [3.13]$
MM + LSTM + tracker	$1.66 \pm 1.21 [3.25]$	$1.61 \pm 1.16 [3.13]$	$1.57 \pm 1.13 [3.03]$
MM + Transf. + tracker	$1.65 \pm 1.17 [3.21]$	$1.63 \pm 1.16 [3.16]$	$1.61 \pm 1.15 [3.11]$
<b>MM + Transf. + prior + tracker</b>	<b><math>1.56 \pm 1.13 [3.09]</math></b>	<b><math>1.53 \pm 1.11 [3.04]</math></b>	<b><math>1.52 \pm 1.10 [2.98]</math></b>



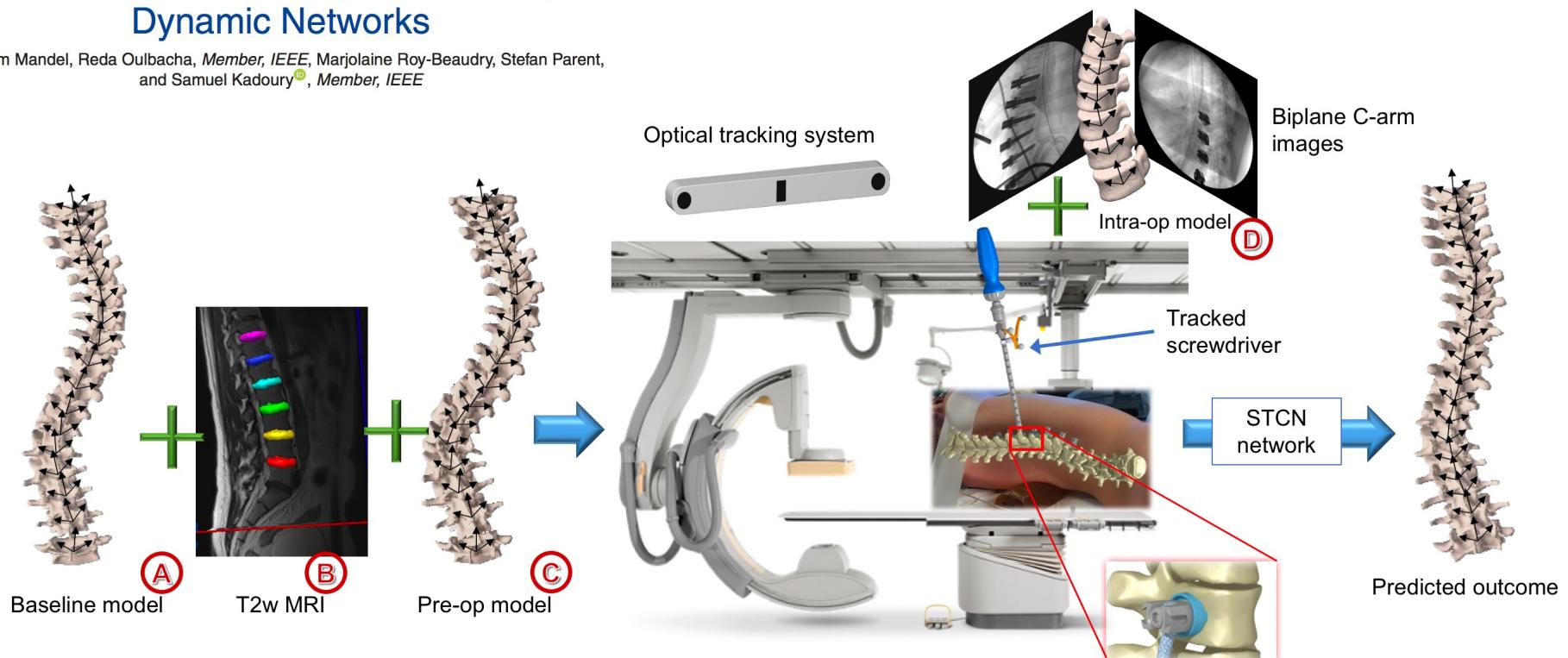
# Outlook

- Self-supervised models to forecast future deformations from 2D images
  - 1) RNN-based: leveraged spatial transformations to implicitly regress motion fields
  - 2) Attention-based: feasibility of Transformer to learn future representations as learnable queries conditioned on prior knowledge
- Clinically relevant accuracy and real-time application (inferences within a few ms)
- Deployed in unseen cases without prior steps
- Tracker: viable alternative for real-time image-guided interventions where pairs of up-to-date volumes are not available

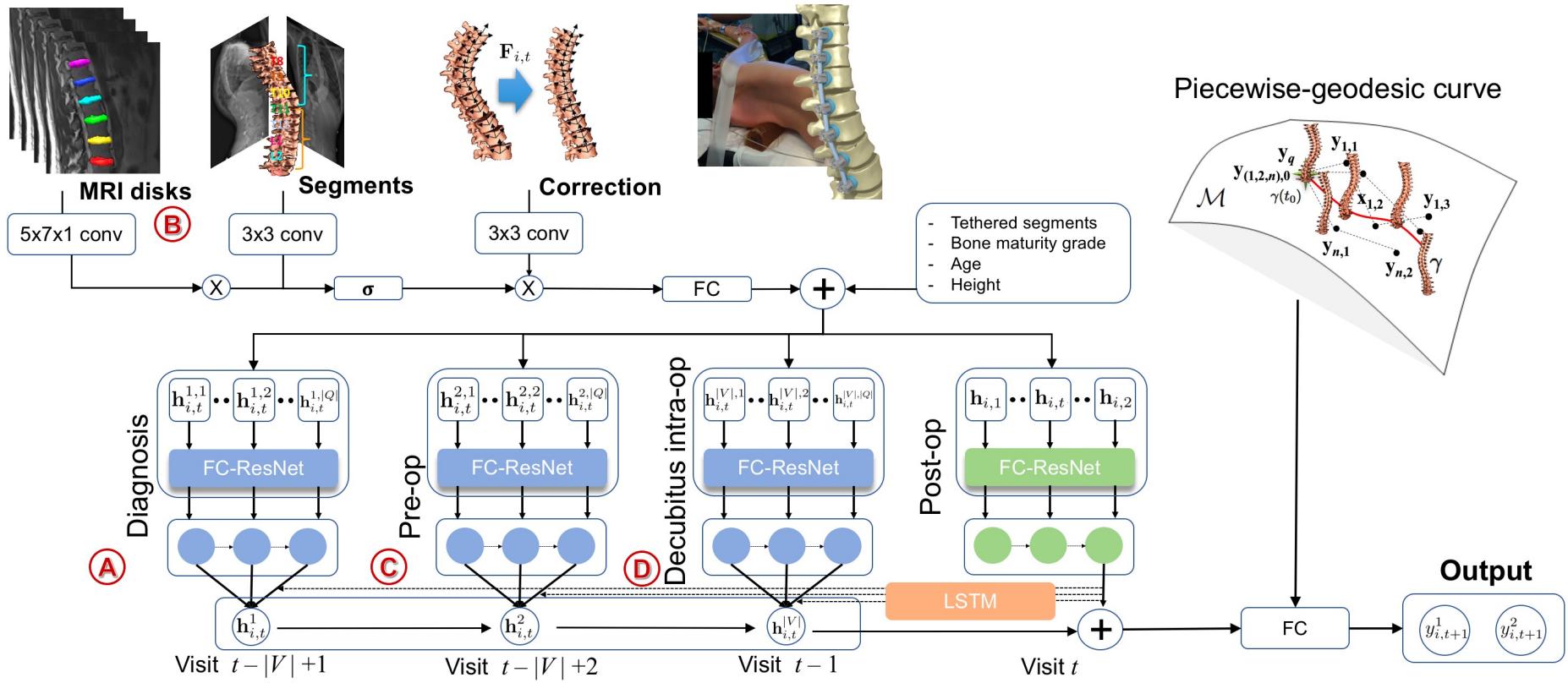
# Spatio-temporal network for spine surgery outcome

## Image-Guided Tethering Spine Surgery With Outcome Prediction Using Spatio-Temporal Dynamic Networks

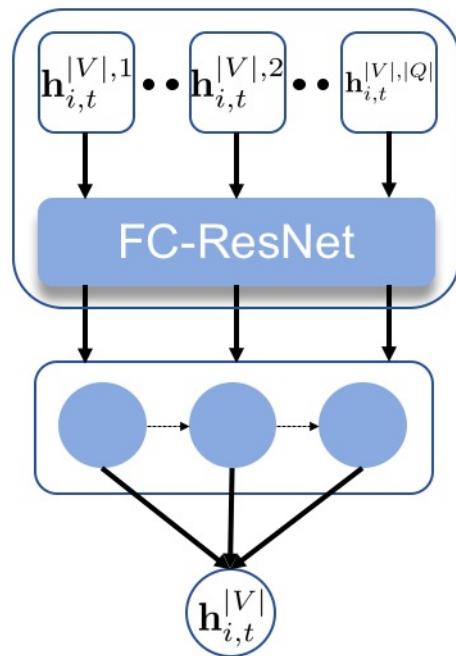
William Mandel, Reda Oulbacha, *Member, IEEE*, Marjolaine Roy-Beaudry, Stefan Parent, and Samuel Kadoury<sup>✉</sup>, *Member, IEEE*



# Spatio-temporal shape prediction model



# Geometric gating mechanism



- FC-ResNet is used to retrain the geometric features for each serial visit  $V$
- Includes long-term and short-term geometric temporal changes:

$$\mathbf{h}_{i,t}^{v,q} = \text{FC-ResNet}([\mathbf{y}_{i,t}^{v,q}; \mathbf{e}_{i,t}^{v,q}], \mathbf{h}_{i,t}^{v,q-1})$$

- $\mathbf{h}_{i,t}^{v,q}$  representing the spine region  $i$  at time  $q$  from a previous visit  $v$  for the predicted time  $t$ ;
- $\mathbf{e}_{i,t}^{v,q}$  are the features related to the additional external parameters :
  - # of tethered segments;
  - Age;
  - Height;
  - Bone maturity status (Risser grade).

# Time shifting mechanism

- The representation  $\mathbf{h}_{i,t}^v$  for each previous visit is found by the sum of all visits within the time frame  $v$  close to the input case:

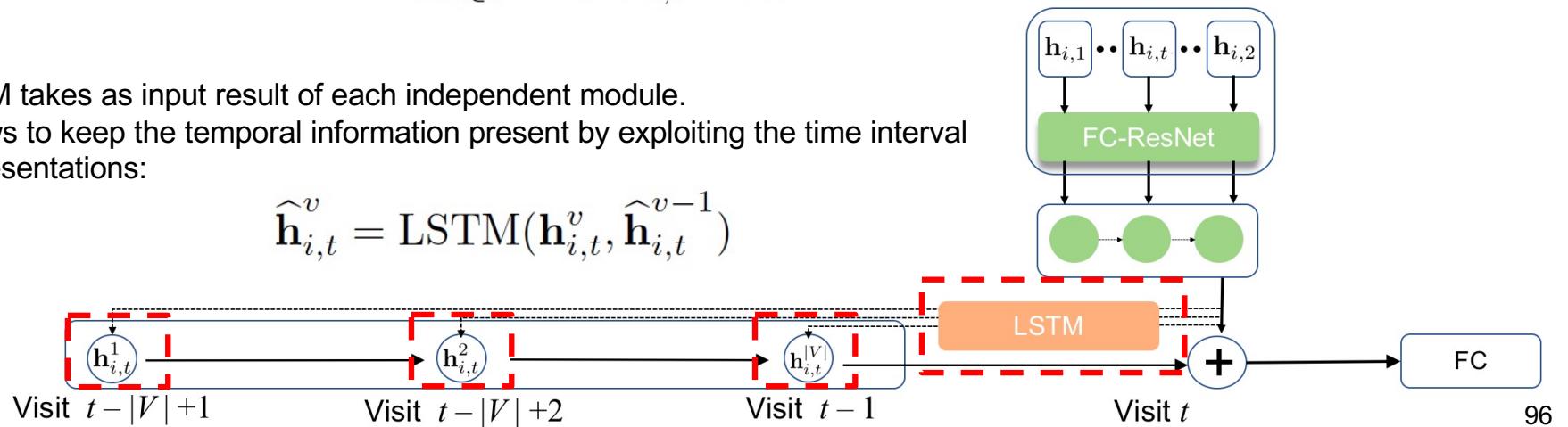
$$\mathbf{h}_{i,t}^v = \sum_{q \in Q} \alpha_{i,t}^{v,q} \mathbf{h}_{i,t}^{v,q} \quad \text{Assigns the importance of the time gap } q \text{ in a visit } v$$

- Weight parameters compare the learned spatio-temporal representation:

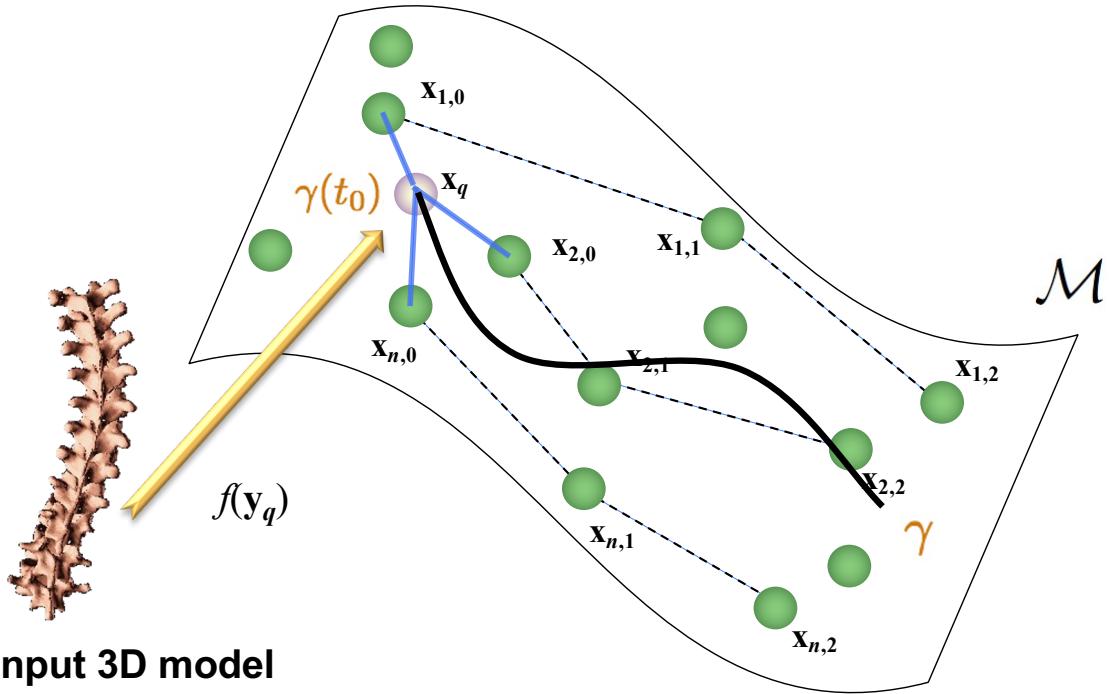
$$\alpha_{i,t}^{v,q} = \frac{\exp(\Phi(\mathbf{h}_{i,t}^{v,q}, \mathbf{h}_{i,t}))}{\sum_{v \in V} \exp(\Phi(\mathbf{h}_{i,t}^{v,q}, \mathbf{h}_{i,t}))} \quad \Phi \text{ measures the similarity in correction flows between pairs of spines.}$$

- LSTM takes as input result of each independent module.
- Allows to keep the temporal information present by exploiting the time interval representations:

$$\hat{\mathbf{h}}_{i,t}^v = \text{LSTM}(\mathbf{h}_{i,t}^v, \hat{\mathbf{h}}_{i,t}^{v-1})$$



# Piecewise-geodesic group-average trajectory



1. Projection of input onto  $\mathcal{M}$  using kernel function  $f$
2. Neighborhood selection  $K_d$  at  $t_0$
3. Shifting of time measurements to common baseline
4. Regression of smooth piecewise-geodesic curve  $\gamma : [t_l, t_N]$
5. Minimization of discretized energy term  $E(\gamma)$  with quadratic optimization problem (Boumal et al., IFAC, 2011):

$$E(\gamma) = \boxed{\frac{1}{K_d} \sum_{i=1}^{K_d} \sum_{j=0}^{t_N} w_i \|\gamma(t_{i,j}) - (\mathbf{x}_{i,j} - (\mathbf{x}_{i,0} - \mathbf{x}_q))\|^2} + \boxed{\frac{\lambda}{2} \sum_{i=1}^{K_d} \alpha_i \|v_i\|^2} + \boxed{\frac{\mu}{2} \sum_{i=1}^{K_d} \beta_i \|a_i\|^2}$$

Minimize geodesic distance  $\|\dot{\gamma}(t_i)\|$   
Acceleration penalty  $\frac{\lambda}{2} \sum_{i=1}^{K_d} \alpha_i \|v_i\|^2$   
discontinuity points  $\frac{\mu}{2} \sum_{i=1}^{K_d} \beta_i \|a_i\|^2$

# Training with piecewise-geodesic trajectory

1. Projection of  $\mathbf{y}_{i,t}$  onto the piecewise-geodesic curve  $\gamma$
2. The output predictions are given by:

$$[y_{i,t+1}^1, y_{i,t+1}^2] = \tanh(\mathbf{W}_{fa}\mathbf{h}_{i,t}^c + \mathbf{b}_{fa})$$

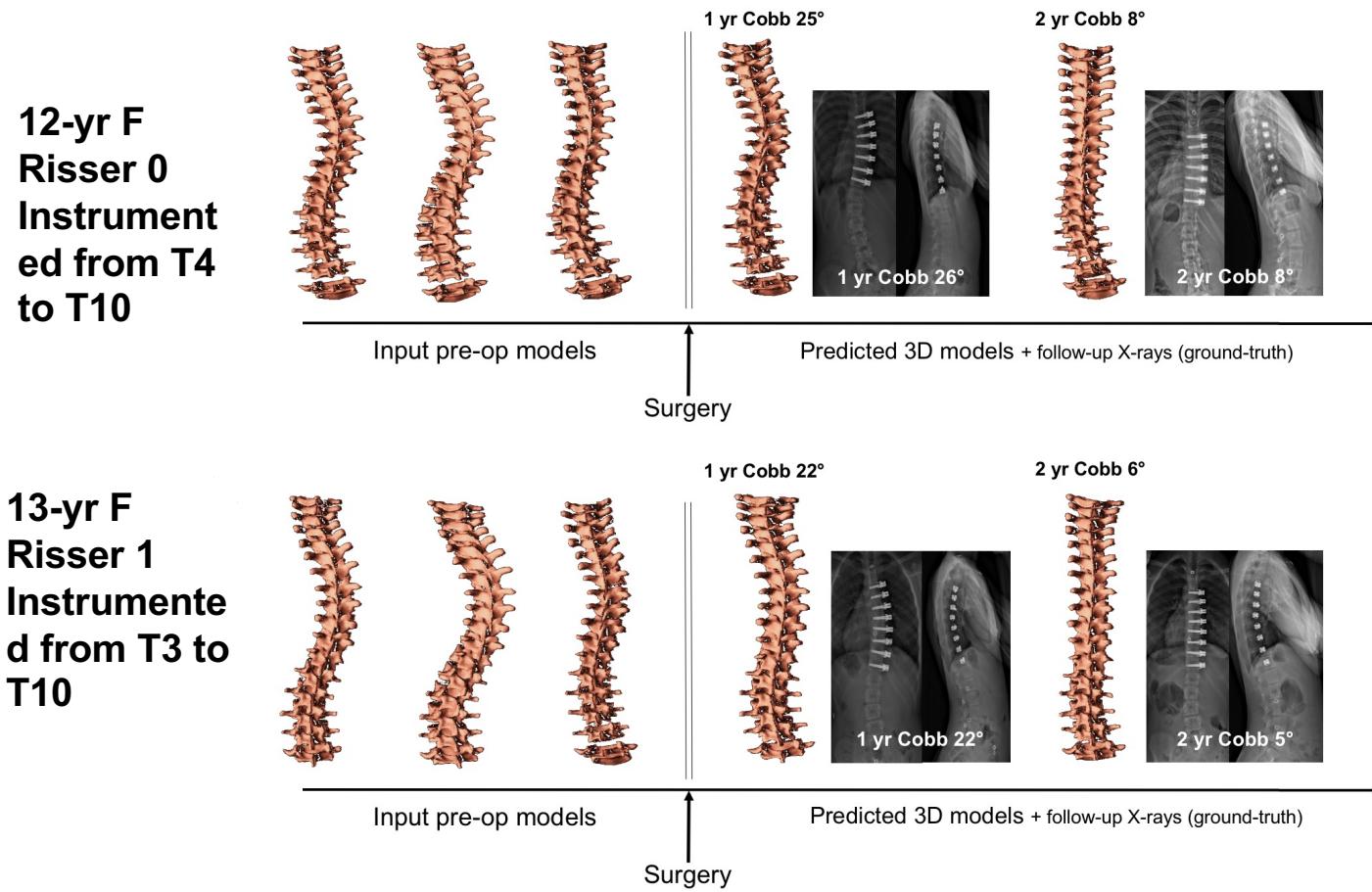
3. Triplet loss function integrating the network output combined with the piecewise-geodesic curve is defined as:

$$\begin{aligned}\mathcal{L} = \sum_{i=1}^n & \beta(y_{i,t+1}^1 - \hat{y}_{i,t+1}^1)^2 + (1 - \beta)(y_{i,t+1}^2 - \hat{y}_{i,t+1}^2)^2 + \\ & \lambda((y_{i,t+1}^1 - \gamma(1)) + (y_{i,t+1}^2 - \gamma(2)))\end{aligned}$$

with  $\beta$  and  $\lambda$  modulating actual outcomes  $(\hat{y}_{i,t+1}^1, \hat{y}_{i,t+1}^2)$  and trajectory regularization, respectively.

4. Repeating this process ensure regularization along  $\gamma$ .

# Shape prediction results of 3D geometry



# Prediction accuracy results

3D RMS errors (mm), Dice (%) and Cobb angles ( $^{\circ}$ ), compared with biomechanical model (Cobetto et al. 2018) , ST-ResNet (Zheng et al. AAAI 2017), deep auto-encoders (AE; Thong et al. 2016), and ST-Manifold (MICCAI 2018).

	1-year visit		
	3D RMS (mm)	Dice	Cobb
Biomec. Simulation	$3.8 \pm 1.3$	$83 \pm 3.7$	$3.5 \pm 1.0$
Deep AE	$4.9 \pm 3.5$	$80 \pm 4.3$	$5.7 \pm 2.7$
ST-ResNet	$3.4 \pm 1.7$	$85 \pm 3.6$	$4.0 \pm 2.4$
ST-Manifold (Mandel et al.)	$3.0 \pm 1.0$	$90 \pm 2.7$	$2.1 \pm 0.8$
STCN (no reg. + gating)	$2.5 \pm 0.8$	$91 \pm 2.6$	$2.0 \pm 0.6$
STCN (no regularization)	$2.2 \pm 0.7$	$92 \pm 2.5$	$1.9 \pm 0.6$
<b>Proposed</b>	$1.9 \pm 0.8$	$93 \pm 2.5$	$1.7 \pm 0.5$

# Questions