



# AI Alignment in Medical Image Segmentation: Model-Fitting, Robustness and Reassurance

---

**Zeju Li, PhD**

FMRIB Center, Nuffield Department of Clinical Neurosciences,  
University of Oxford, UK.

Email: [zeju.li@ndcn.ox.ac.uk](mailto:zeju.li@ndcn.ox.ac.uk)

March 2024

Statistics

3,288 citations

12 first authored peer-reviewed journal/conference

## Education

- 2023.1 – Current: FMRIB Analysis Group
  - Senior Researcher in Brain Connectivity, University of Oxford, Oxford
- 2018.10 – 2022.12: BioMedia Group
  - PhD in Computing, Imperial College London, London
  - Supervisors: Prof. Ben Glocker and Prof. Daniel Rueckert
  - Examiners: Dr. Juan Eugenio Iglesias and Dr. Stamatia Giannarou
- 2011.9 – 2018.7:
  - Master in Biomedical Engineering, Fudan, Shanghai
  - Bachelor in Electronic Engineering, Fudan, Shanghai



**Imperial College  
London**



## Intern

- 2019.7 – 2020.4: Computer Vision Group
  - Huawei Noah's Ark Lab, London
- 2018.7 – 2018.9: MIRACLE Group
  - Institute Of Computing Technology, Beijing



## Medical Image Segmentation

- Identify groups of pixels that go together for **medical imaging** (e.g. MRI, CT, ultrasound...).
- A critical step to transfer the power of machine learning into the clinical diagnosis process.

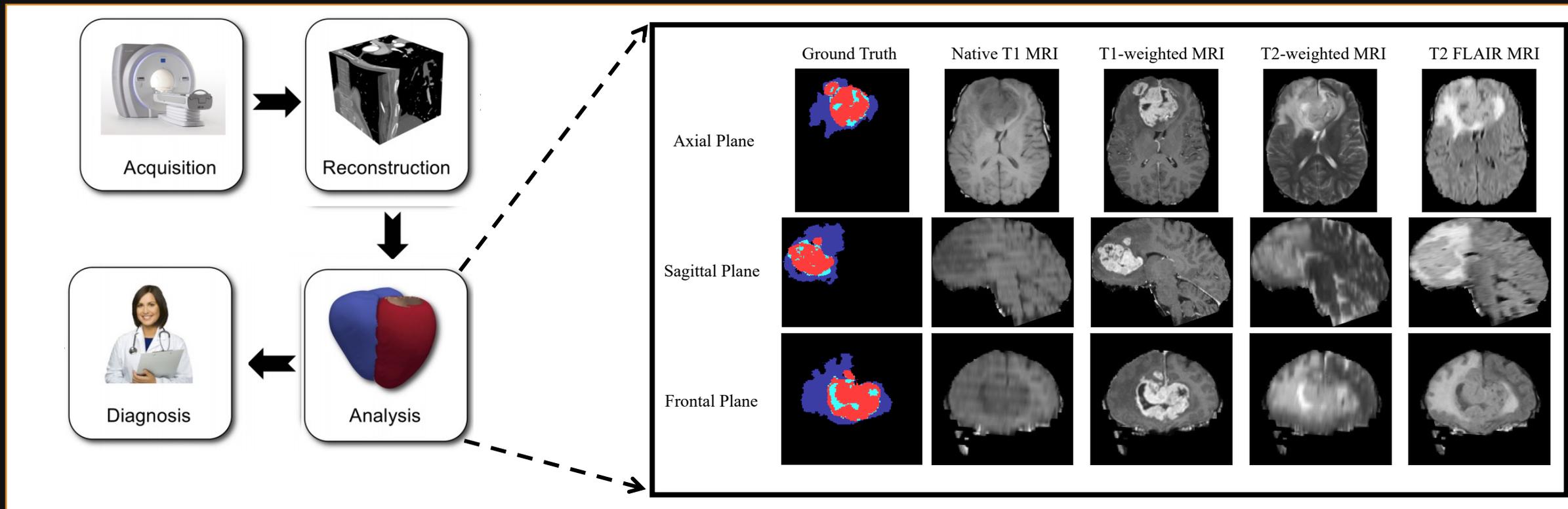


Figure credited to D. Rueckert, J.A. Schnabel. "Model-based and data-driven strategies in medical image computing", P IEEE, 2019.

# AI Alignment in Medical Image Segmentation

4/40

The model is unreliable  
because of misaligned goals!

Model  
Performance



Deployment  
Environment



Optimization  
Algorithms



Data  
Distributions



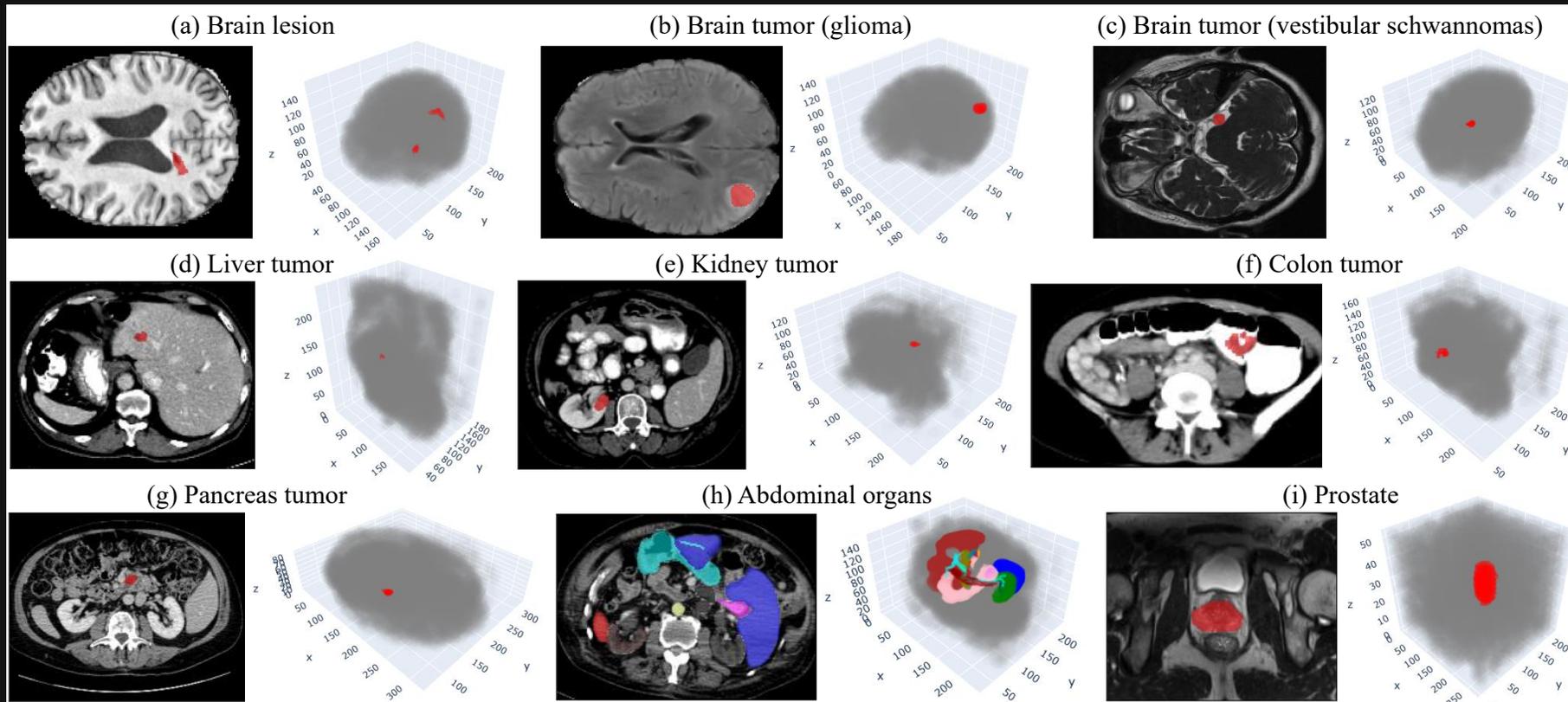
**Reassurance:** Expectation  
Over-optimism Problem.

**Robustness:** Goal Mis-  
generalization Problem.

**Model-Fitting:** Specification  
Failure Problem.

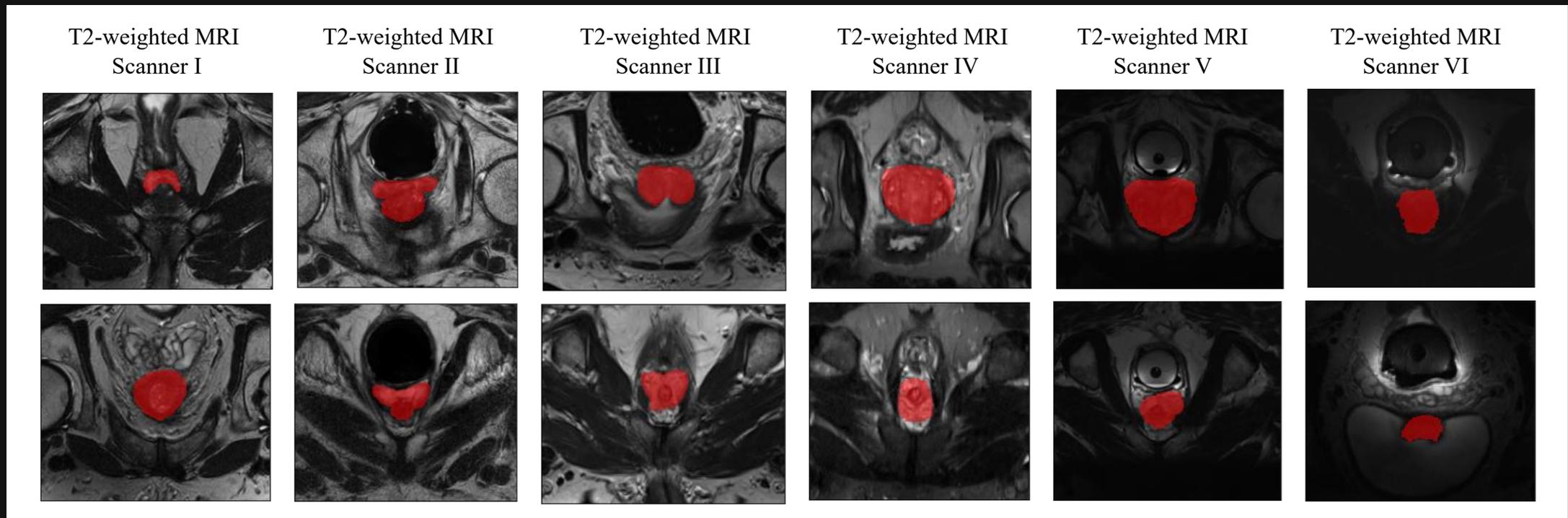
## Specification Failure Reason: Class Imbalance

- **Class imbalance**, which refers to the imbalanced distributions of samples from different categories, cause difficulties for machine learning models to learn well.



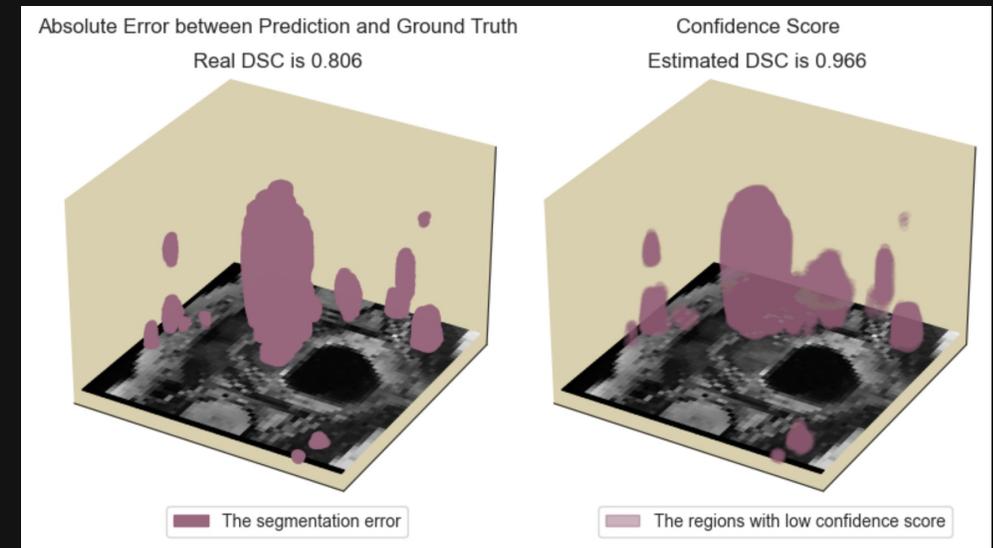
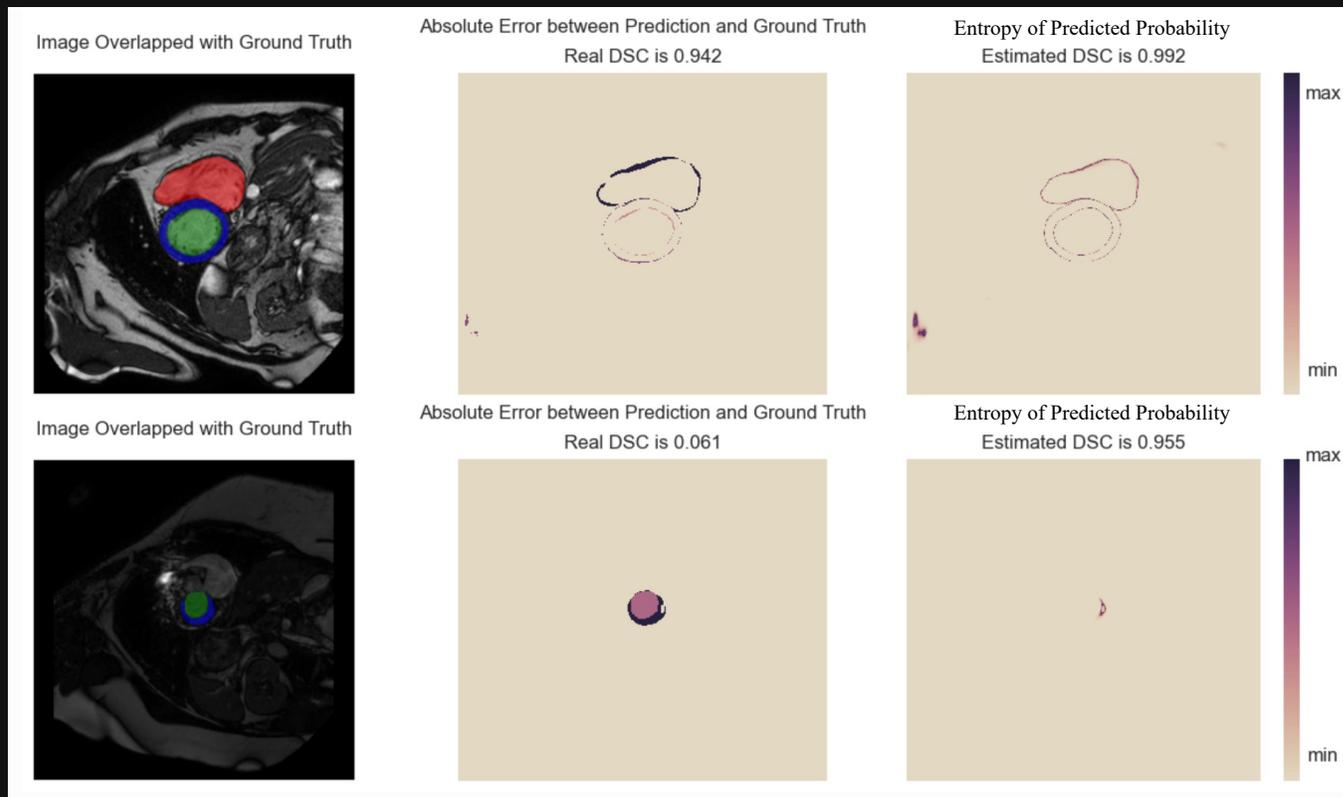
## Goal Mis-generalization Reason: Domain Shifts

- **Domain shift** is a common issue in medical imaging as the medical images can be collected from different clinical sites. As data-driven methods, deep neural network may not generalize well under domain shifts.



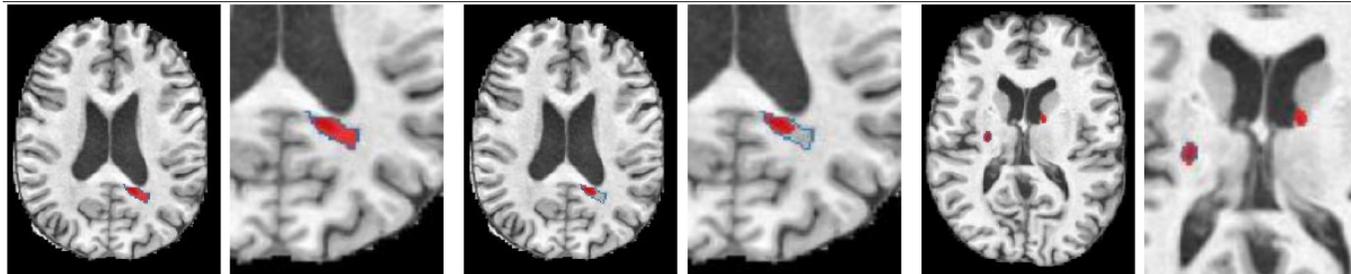
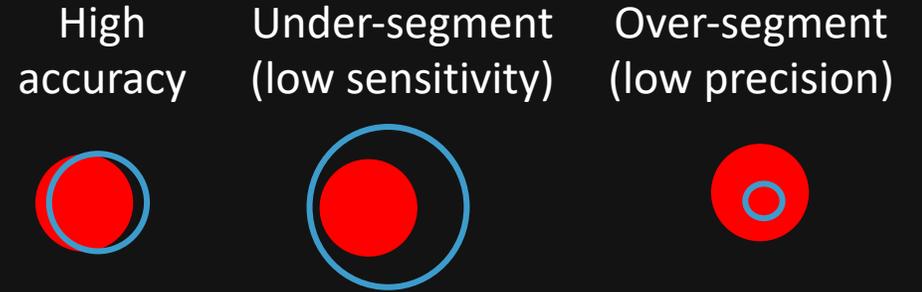
## Expectation Over-optimism Reason: Model Uncalibration

- As probabilistic model, neural network produces the probability of the predictions. However, modern deep neural networks are known to **over-confident and uncalibrated** about the predictions.



## Observations

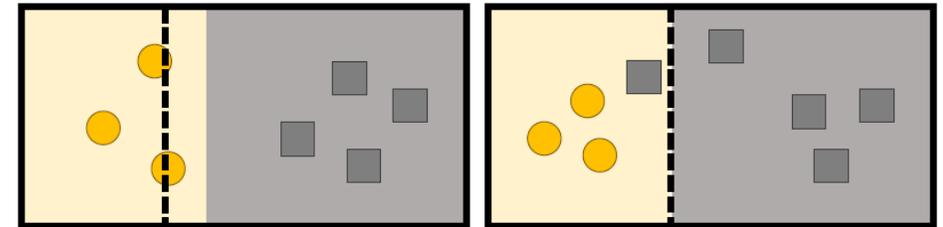
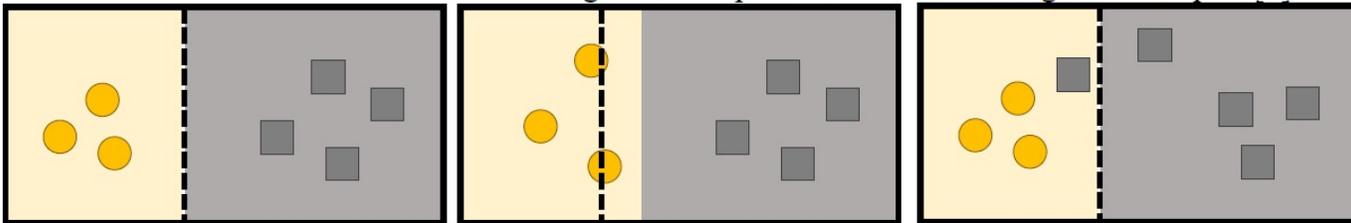
- Class imbalance causes **under- and over-segmentation**.
- Understanding the effects of class imbalance in segmentation.



(a) Ideal segmentation

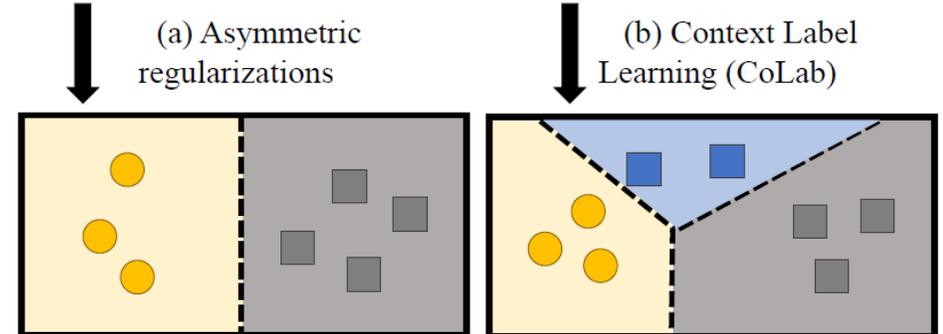
(b) Under segmentation:  
Overfitting of under-represented foreground samples

(c) Over segmentation:  
Underfitting of heterogenous background samples



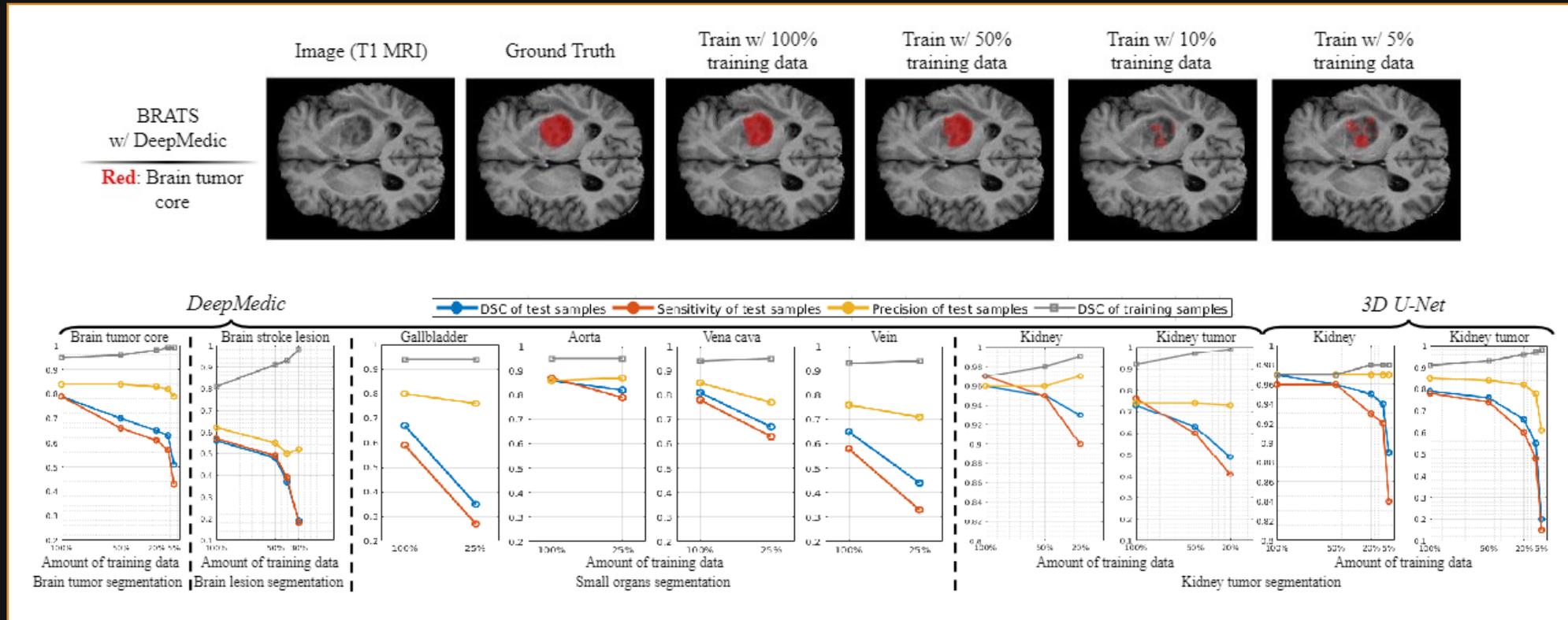
(a) Asymmetric regularizations

(b) Context Label Learning (CoLab)



## Analysis

- With less training data, performances decline due to the drastic **reduction of sensitivity**, while precision is retained.



High accuracy



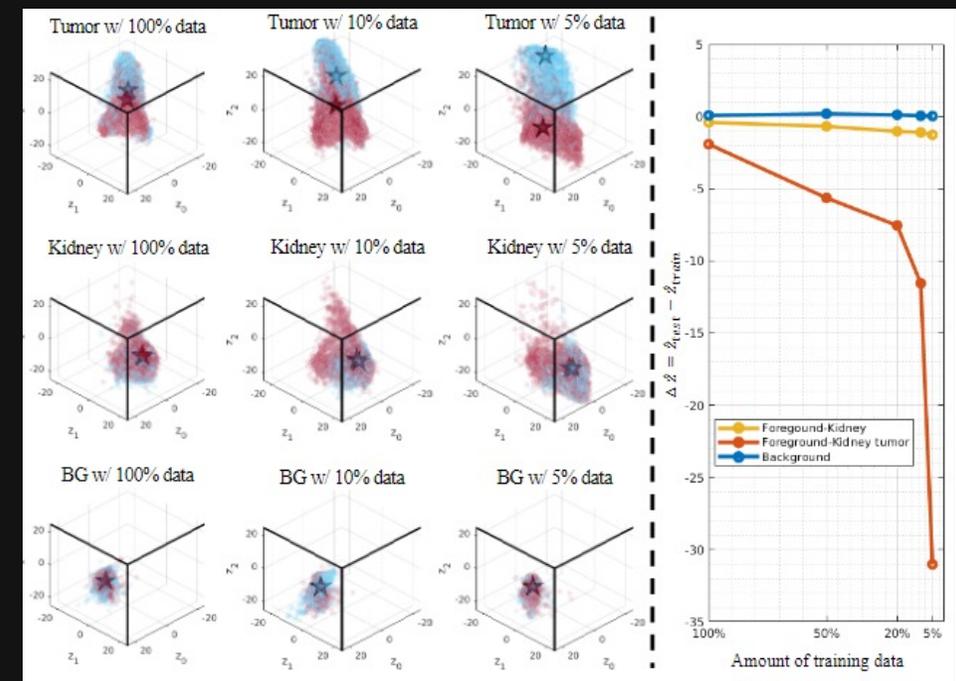
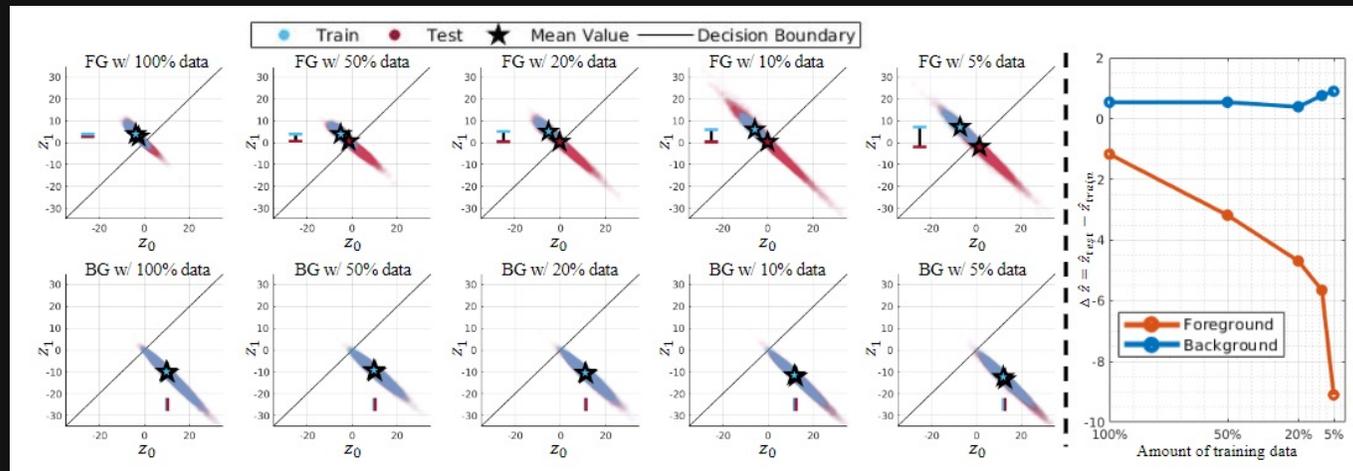
Under-segment (low sensitivity)



# Overfitting under Class Imbalance

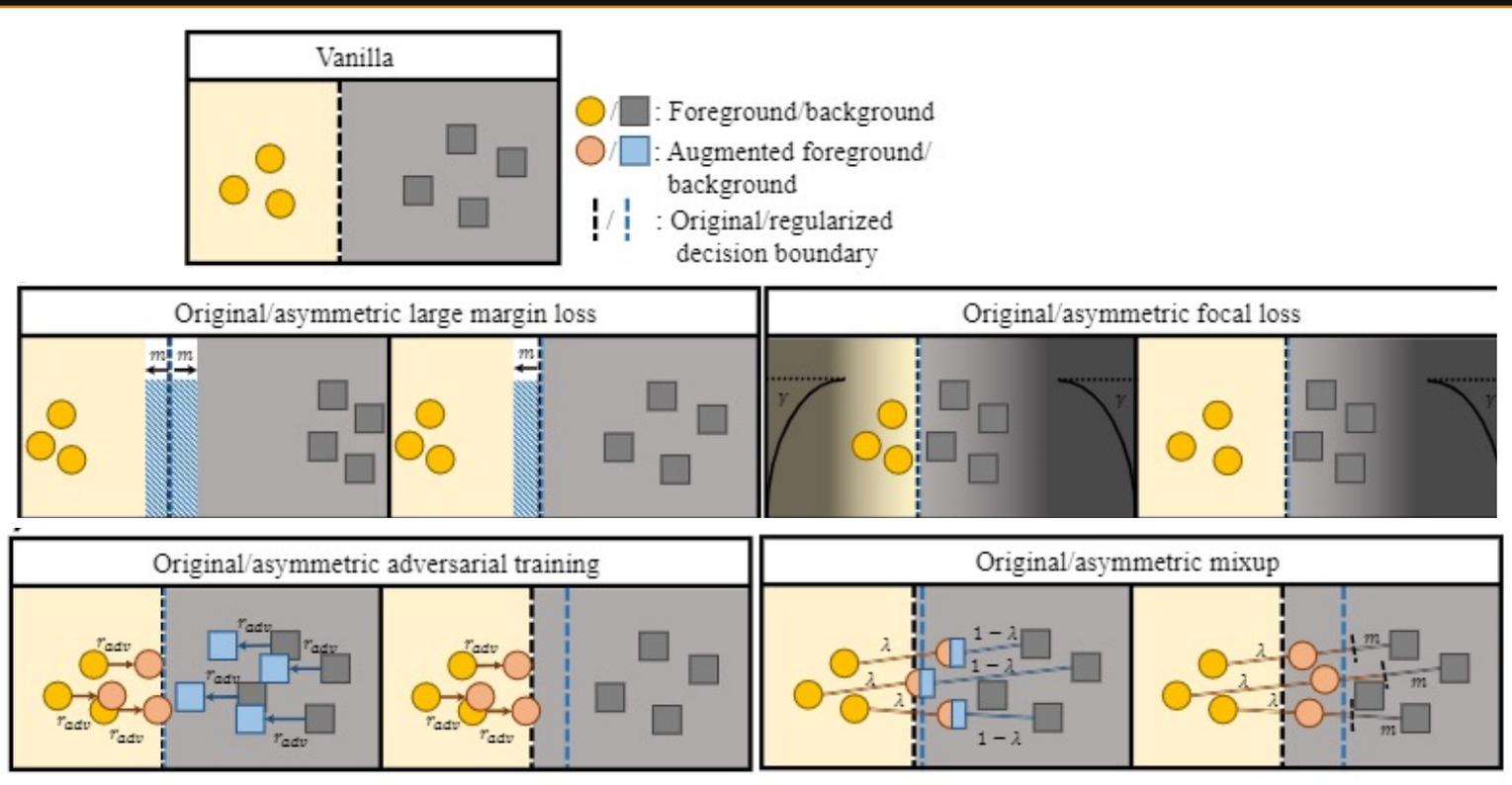
## Analysis

- CNN maps training and testing samples of the background class to similar logit values.
- However, **mean activation for testing data shifts** significantly for the foreground class towards and sometimes across the decision boundary.



## Method and Results

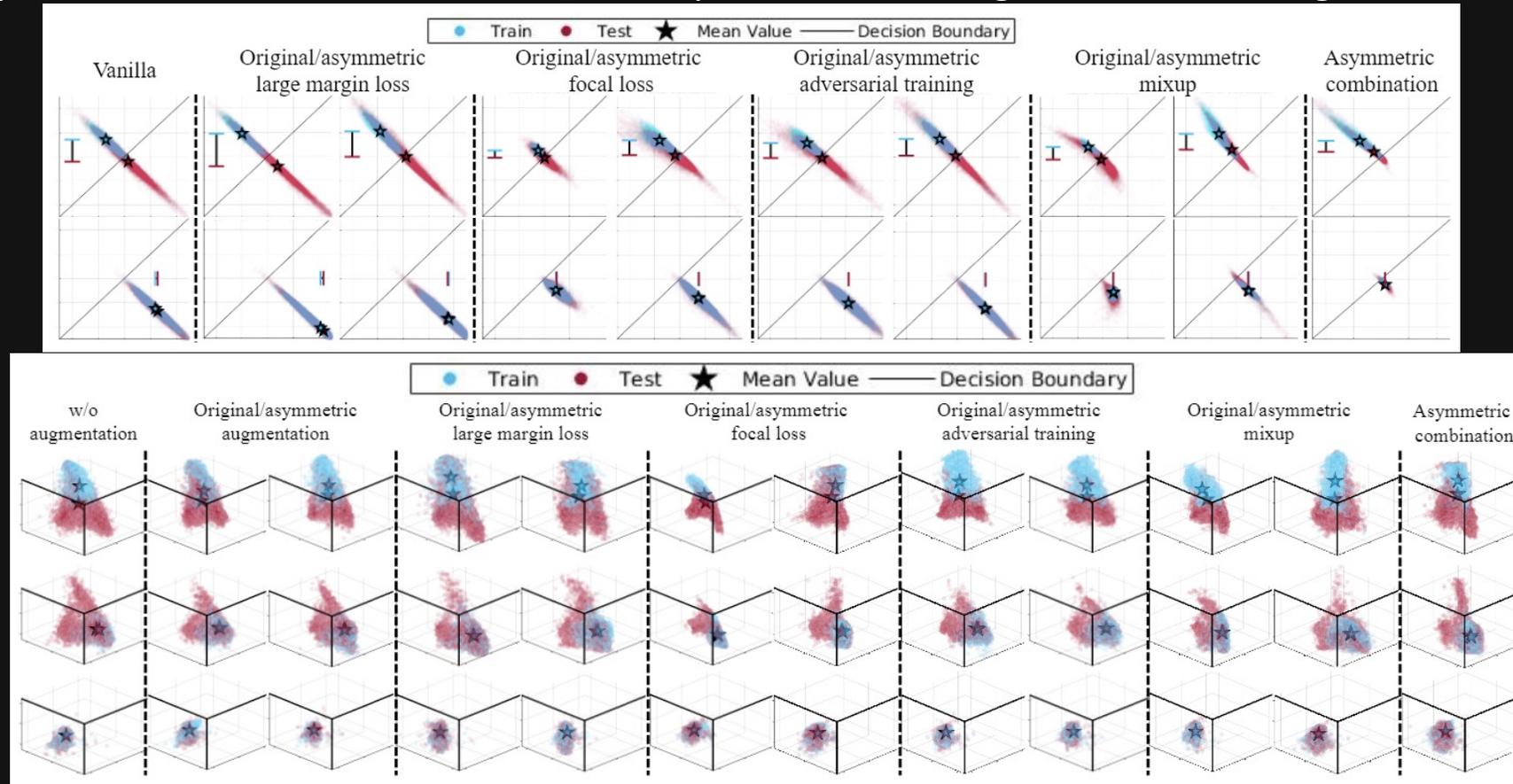
- We make the logit activations of foreground class far away from the decision boundary by **setting bias for the foreground class** in different ways.



Method	5% training			
	DSC	SEN	PRC	HD
Vanilla - CE [20]	50.4	41.0	83.5	18.0
Vanilla - CE - 80% tumor	45.5	36.0	86.7	17.8
Vanilla - F1 (DSC)	47.2	37.4	86.6	15.9
Vanilla - F2 [14]	45.8	36.9	81.9	17.9
Vanilla - F4 [14]	51.6	42.5	83.8	18.1
Vanilla - F8 [14]	47.4	38.7	83.1	19.6
Large margin loss [31]	44.5	35.9	82.8	20.2
<b>Asymmetric large margin loss</b>	56.8	48.9	83.4	<b>15.0</b>
Focal loss [29]	54.0	44.8	82.6	16.0
<b>Asymmetric focal loss</b>	58.8	51.4	81.6	<b>15.0</b>
Adversarial training [12]	53.2	44.6	85.0	19.2
<b>Asymmetric adversarial training</b>	58.5	50.8	80.1	16.2
Mixup [47]	49.7	40.9	83.0	19.6
<b>Asymmetric mixup</b>	<b>59.8</b>	56.8	74.7	17.7
Symmetric combination	50.0	42.0	84.6	21.1
<b>Asymmetric combination</b>	<b>63.4</b>	63.1	75.9	<b>15.1</b>

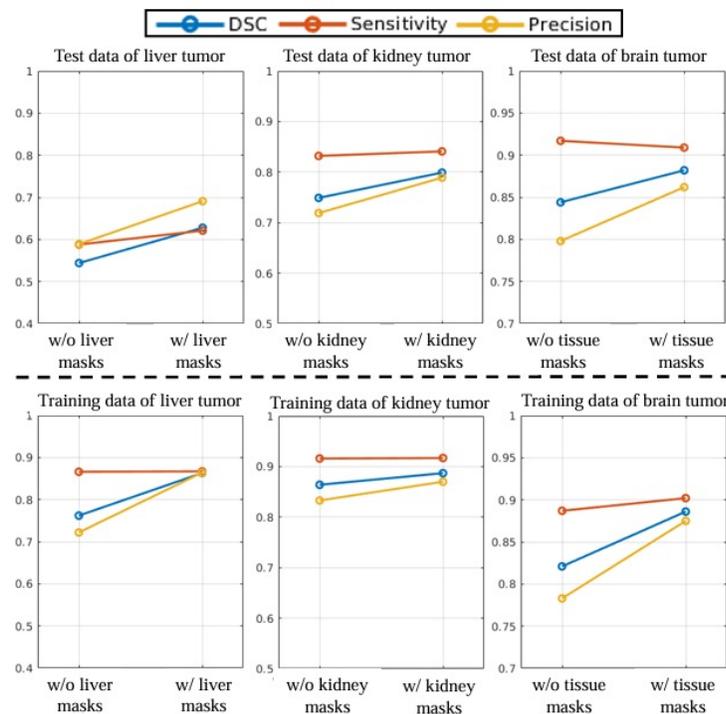
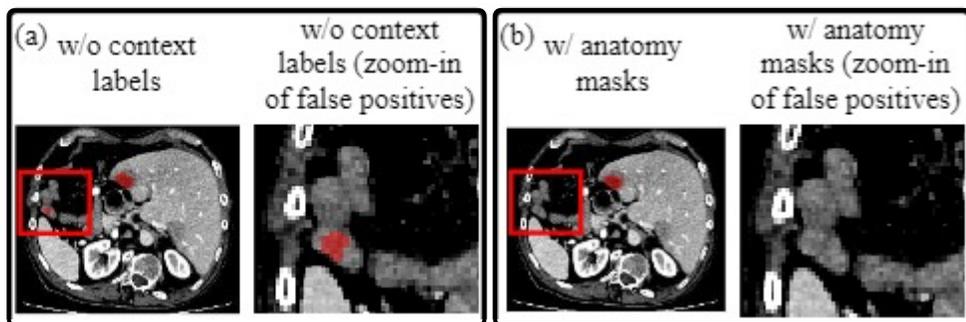
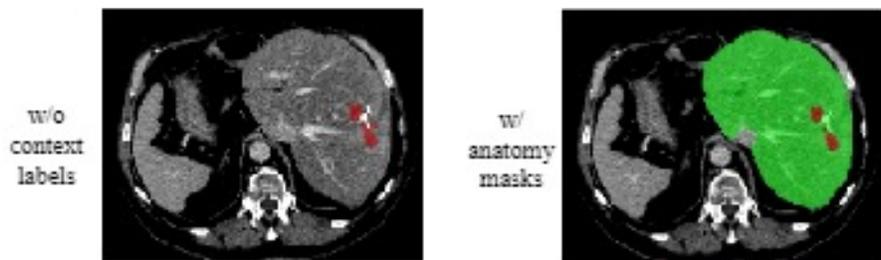
## Results

➤ Asymmetric modifications lead to better separation of the logits of unseen foreground samples.



## Analysis

- With heterogeneous background, performances decline due to the drastic **reduction of precision**, while sensitivity is retained.



High accuracy

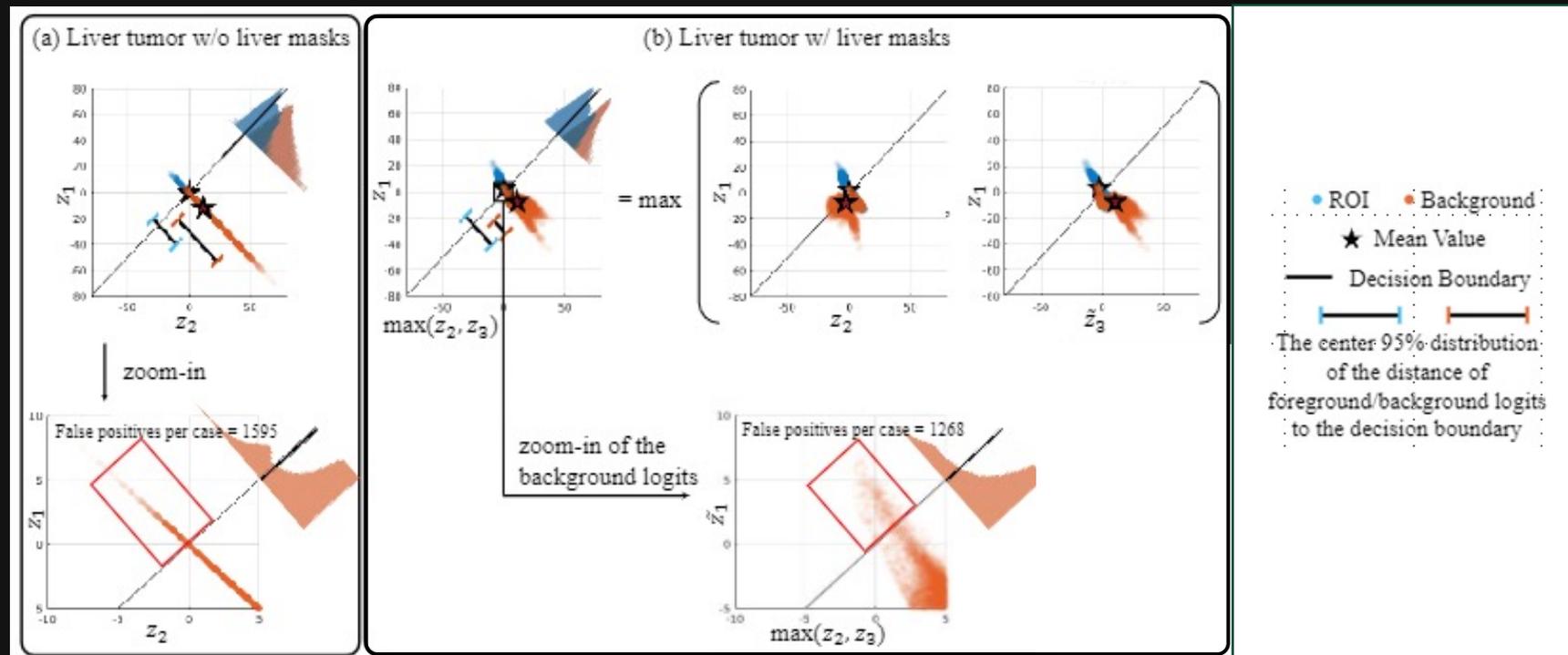


Over-segment (low precision)



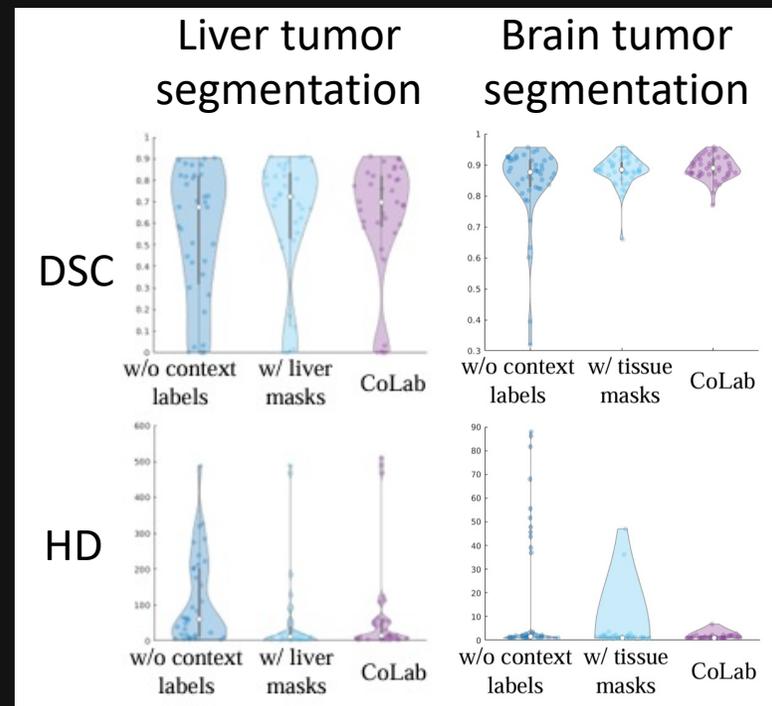
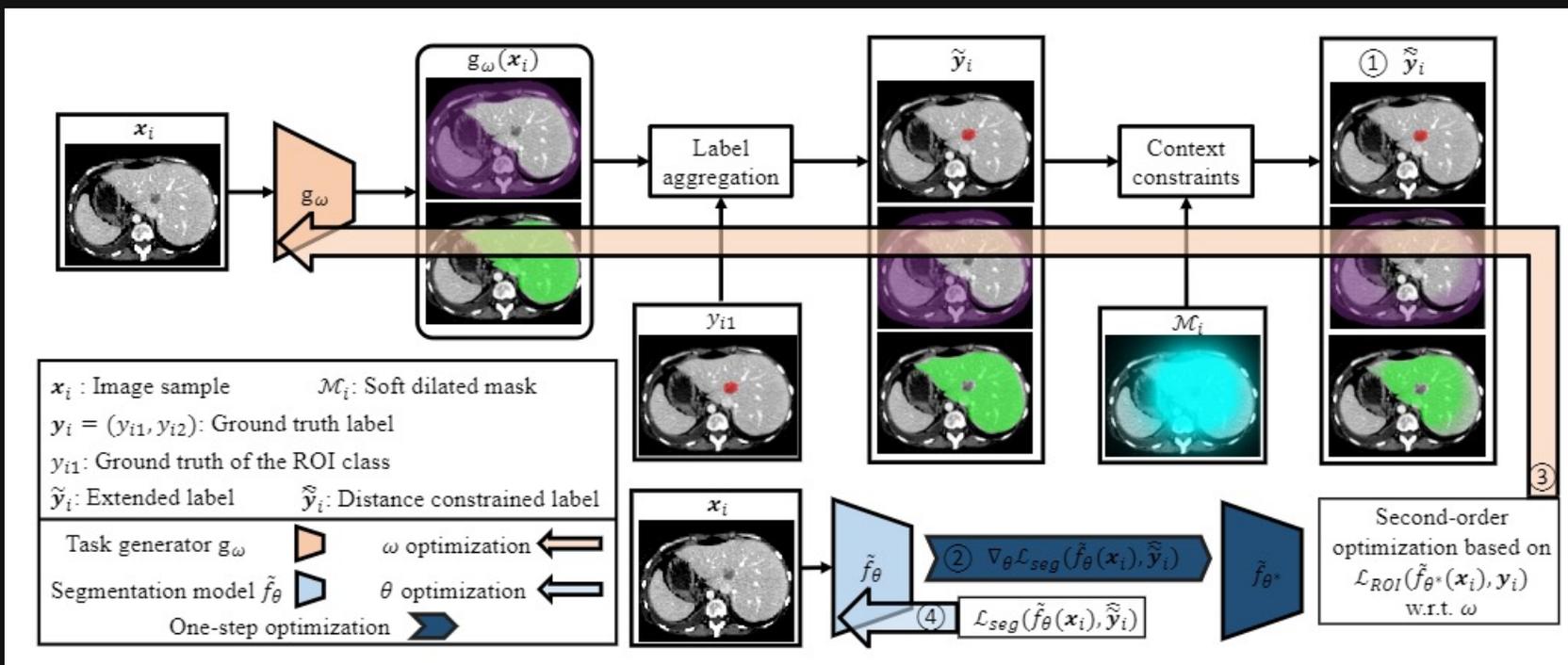
## Analysis

- Neural networks could not map the heterogeneous background samples to **compact clusters** in feature space.
- As a result, the logit activations of background would approach and even move across the decision boundary.



## Method and Results

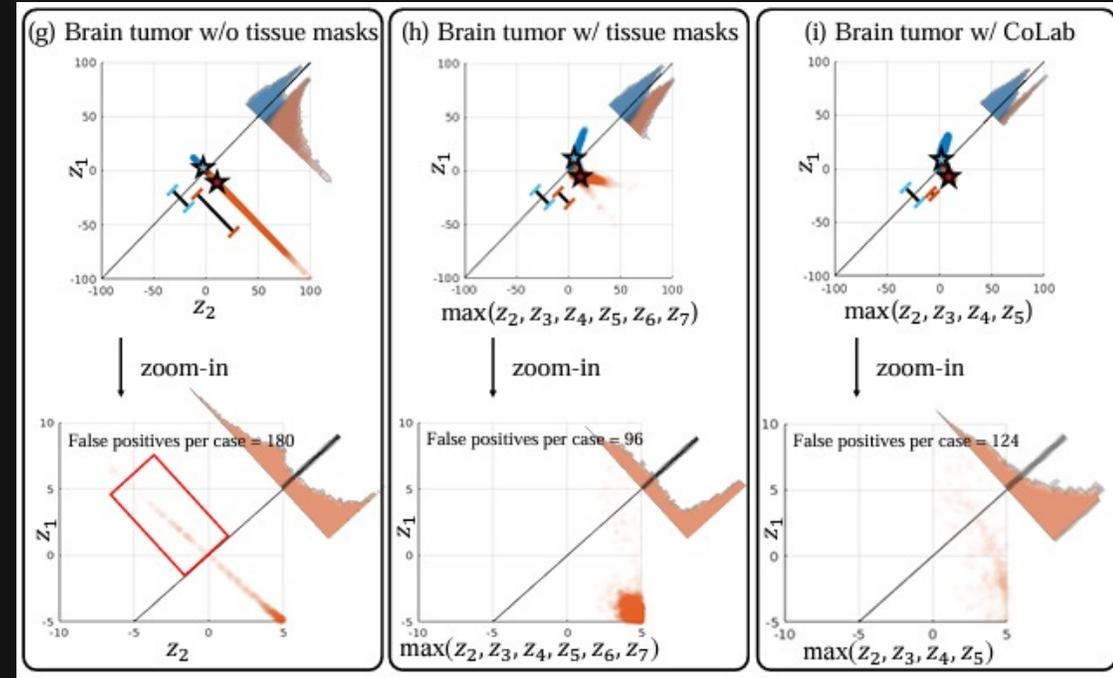
- We propose **Context label learning (CoLab)**.
- We train an auxiliary network as a task generator, along with the primary segmentation model, to automatically generate context labels that positively affect the ROI segmentation accuracy.



## Results

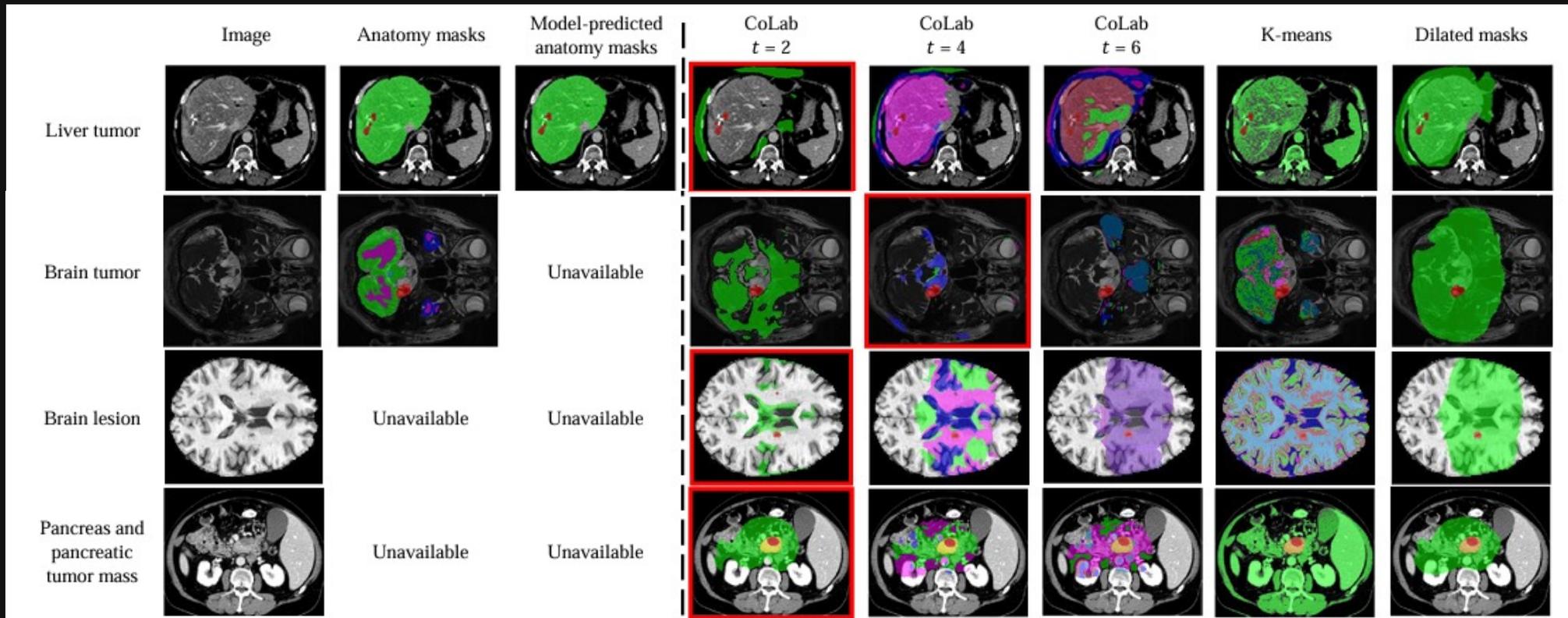
- Similar and sometimes better effect in improving segmentation accuracy when compared with human-defined context labels.

Task	Method	$t$	DSC	SEN	PRC	HD
Liver tumor [5]	w/o liver masks	1	54.4	58.8	58.9	111.1
	K-means [1]	2	<b>61.4</b>	<b>61.4</b>	67.0	71.9
	Dilated masks [26]	2	60.7	59.8	<b>68.0</b>	67.6
	CoLab	2	<b>62.5</b>	<b>62.8</b>	<b>67.3</b>	69.4
	CoLab	4	57.3	60.5	62.3	<b>56.3</b>
	CoLab	6	59.7	60.3	65.2	<b>43.6</b>
	w/ model-predicted liver masks [16]	2	62.4	61.6	70.6	44.1
w/ liver masks [5]	2	62.8	62.1	69.1	53.5	
Kidney tumor [12]	w/o kidney masks	1	74.9	83.2	71.9	120.4
	K-means [1]	2	<b>76.8</b>	<b>83.5</b>	74.3	87.1
	Dilated masks [26]	2	76.4	<b>83.9</b>	73.1	95.3
	CoLab	2	<b>78.5</b>	<b>82.2</b>	<b>77.7</b>	75.7
	CoLab	4	76.4	80.6	<b>76.5</b>	<b>63.7</b>
	CoLab	6	74.9	81.0	73.3	79.4
	w/ model-predicted kidney masks [16]	2	79.2	81.3	82.7	38.1
w/ kidney masks [12]	2	79.9	84.1	78.9	54.7	
Brain tumor [33]	w/o tissue masks	1	84.3	91.2	80.5	15.2
	K-means [1]	6	85.0	91.3	80.3	9.4
	Dilated masks [26]	2	84.8	91.6	81.1	8.8
	CoLab	2	85.2	90.4	82.5	7.8
	CoLab	4	<b>89.0</b>	<b>91.7</b>	<b>86.7</b>	<b>1.4</b>
	CoLab	6	<b>87.9</b>	<b>92.0</b>	<b>84.9</b>	<b>2.5</b>
	w/ tissue masks [17], [33]	6	88.2	90.9	86.2	3.1

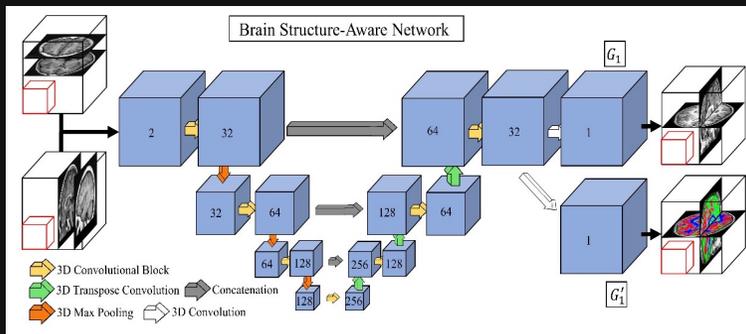


## Results

➤ Examples of context labels generated could inform us on how to design optimal contextual tasks.

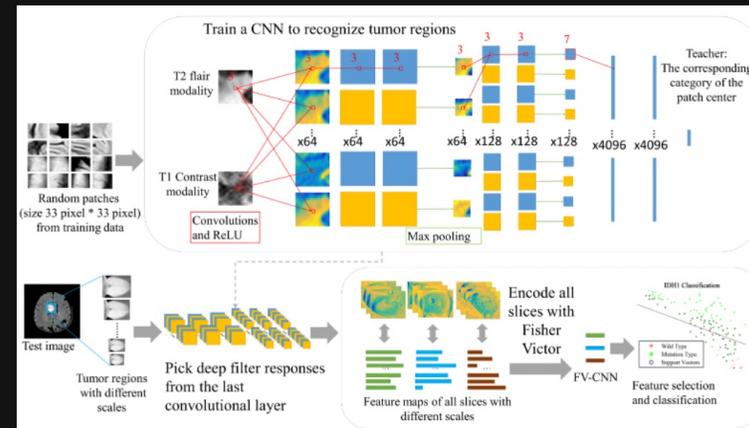


## Improved Model-Fitting with Multi-Task Learning



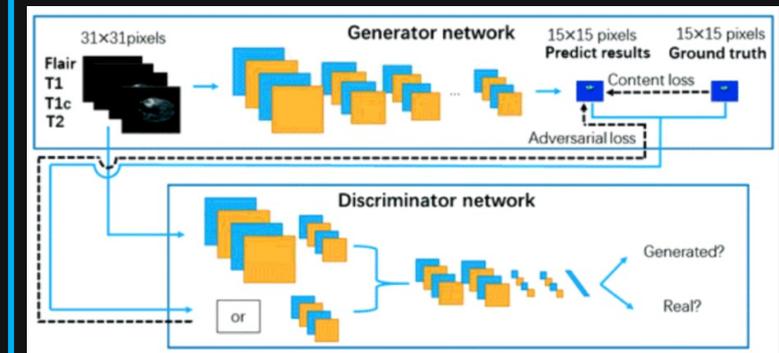
Z. Li *et al.* "Deepvolume: brain structure and spatial connection-aware network for brain MRI super-resolution", TCybern, 2019.

## Feature Re-using for Downstream Applications



Z. Li *et al.* "Deep learning based radiomics (DLR) and its usage in noninvasive idh1 prediction for low grade glioma", Sci. Rep., 2017.

## Generative Model based Segmentation

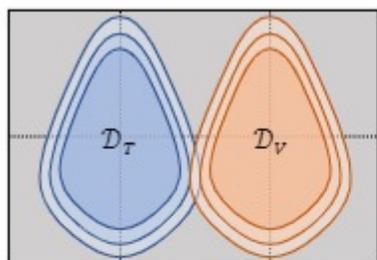


Z. Li *et al.* "Brain tumor segmentation using an adversarial network", MICCAI-Brainlesion workshop, 2017.

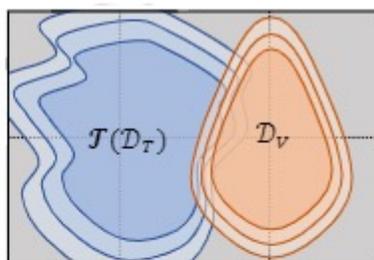
# Joint Optimization of Class-Specific Training- and Test-time data augmentation

## Motivation

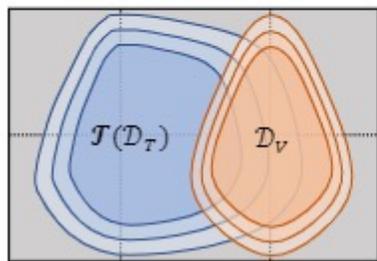
- Training-time data augmentation (TRA) and test-time data augmentation (TEA) are closely connected as both aim to **align the training and test data distribution**.



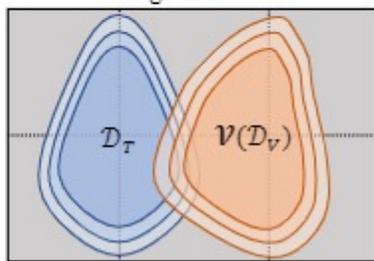
(a) w/o Data Augmentation



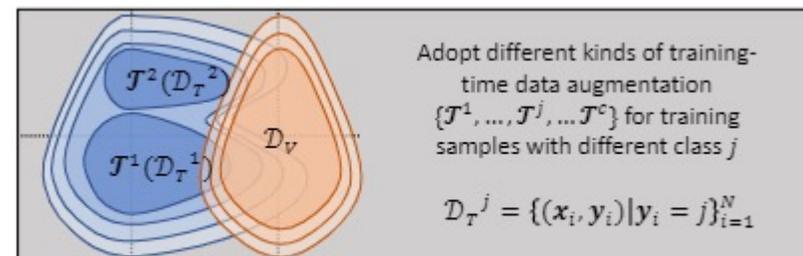
(b) Random Training-Time Data Augmentation



(c) Heuristic/Learned Training-Time Data Augmentation

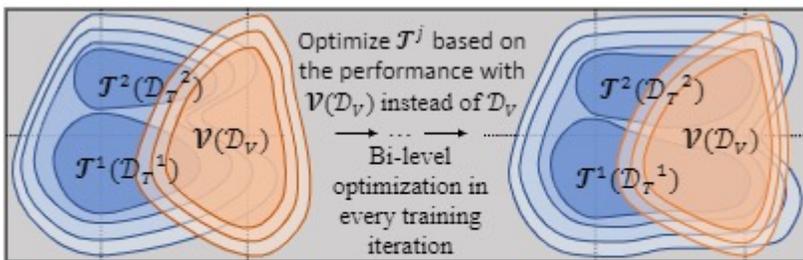


(d) Heuristic/Learned Test-Time Data Augmentation



(e) Learned Class-Specific Training-Time Data Augmentation

Adopt different kinds of training-time data augmentation  $\{\mathcal{J}^1, \dots, \mathcal{J}^j, \dots, \mathcal{J}^c\}$  for training samples with different class  $j$

$$\mathcal{D}_T^j = \{(x_i, y_i) | y_i = j\}_{i=1}^N$$


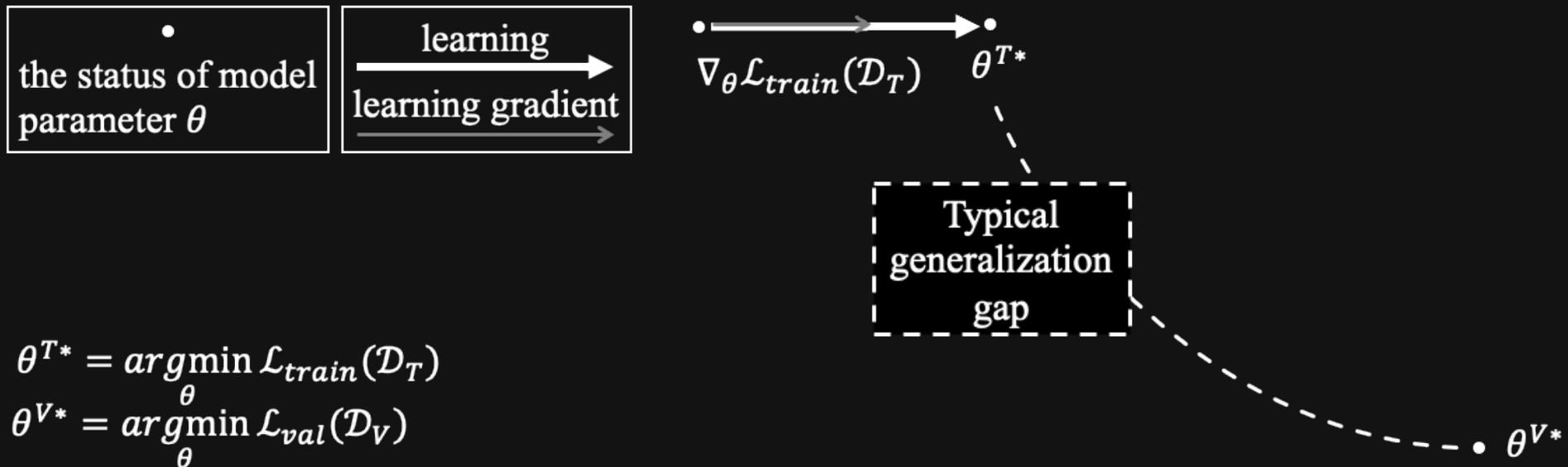
(f) Joint Learned Class-Specific Training- and Test-time Data Augmentation

$\mathcal{D}_T$ : Distribution of training dataset     $\mathcal{J}$ : Training-Time Data Augmentation  
 $\mathcal{D}_V$ : Distribution of validation dataset     $\mathcal{V}$ : Test-Time Data Augmentation

# Joint Optimization of Class-Specific Training- and Test-time data augmentation

## Method

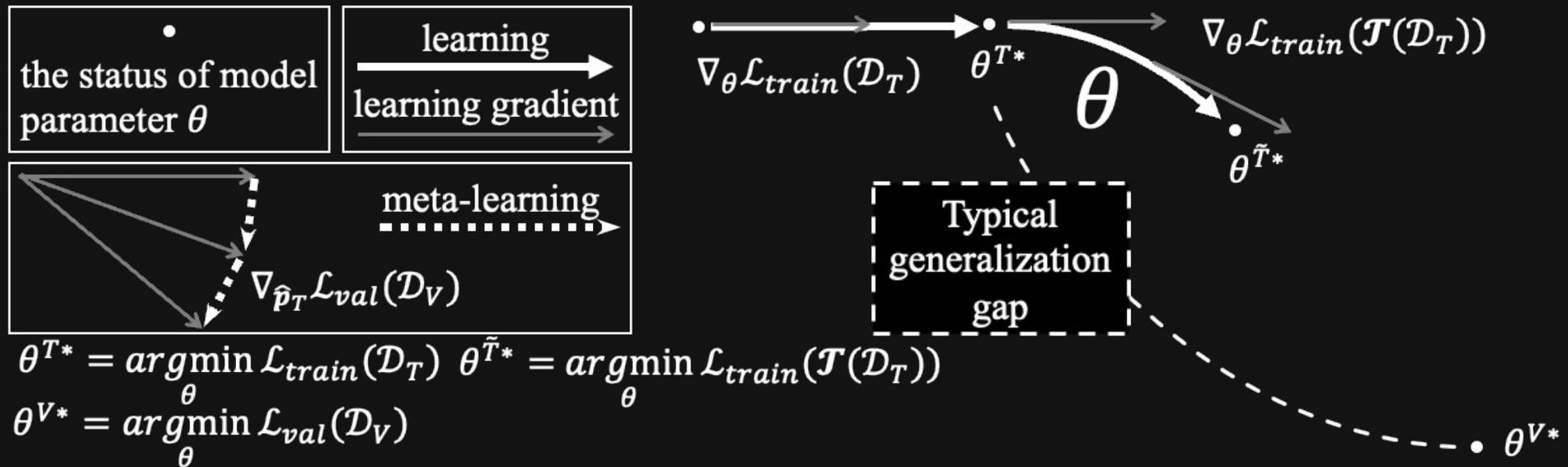
- A meta-learning based data augmentation framework, taking test-time transformations into account.



# Joint Optimization of Class-Specific Training- and Test-time data augmentation

## Method

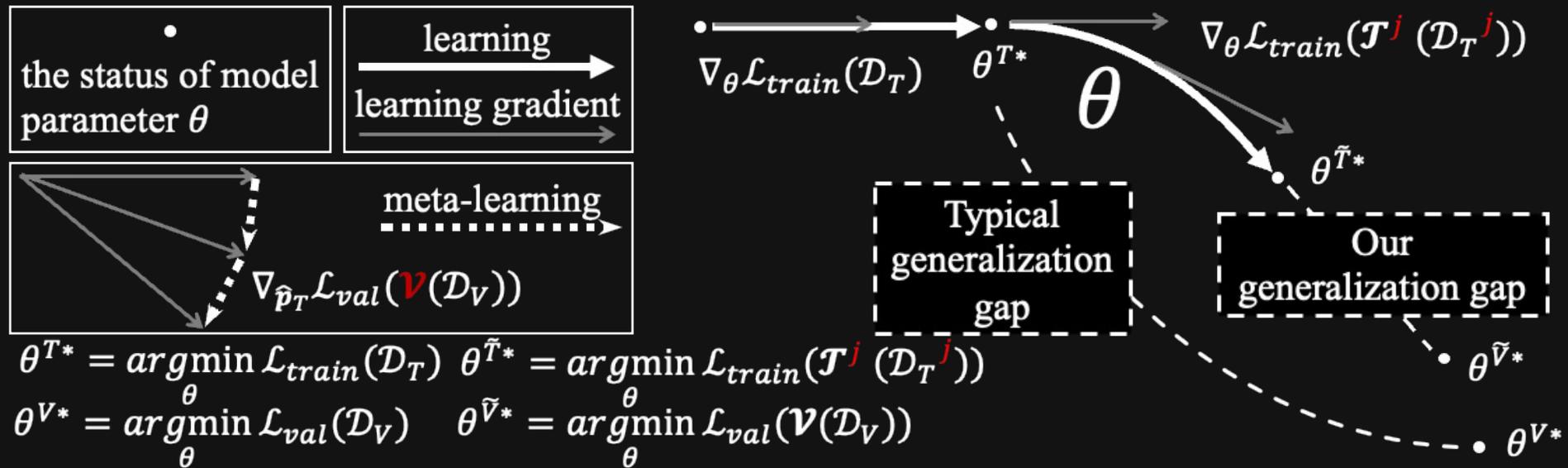
- A meta-learning based data augmentation framework, taking test-time transformations into account.



# Joint Optimization of Class-Specific Training- and Test-time data augmentation

## Method

- A meta-learning based data augmentation framework, taking test-time transformations into account.



# Joint Optimization of Class-Specific Training- and Test-time data augmentation

## Results

- Consistently improve segmentation performance in various applications.
- Potential to **replace the heuristically chosen augmentation policies** currently used in most previous works.

Class-specific TRA improves brain tumor segmentation

Joint optimization of TRA and TEA improves cross-domain prostate segmentation

Model	Training-time augmentation	Test-time augmentation	50% training data			
			DSC ↑	SEN ↑	PRC ↑	HD ↓
3D U-Net [6]	None	None	54.6	56.3	<b>67.2</b>	<b>32.6</b>
	Heuristic [18]	None	58.4	66.9	61.4	39.0
	Heuristic <sup>†</sup> [18]	None	58.8	<b>67.8</b>	59.8	52.2
	Learned [8], [29], [32]	None	<b>59.3</b>	66.6	61.1	40.9
	<b>Learned Class-Specific</b>	None	<b>62.0</b> (+3.6)**	<b>68.8</b>	<b>66.2</b>	<b>37.8</b>
	Heuristic [18]	Heuristic [18]	61.7	<b>67.0</b>	69.6	22.0
	<b>Learned Class-Specific</b>	Heuristic [18]	61.8	66.4	<b>70.2</b>	<b>20.3</b>
	<b>Learned Class-Specific</b>	Learned [22], [40]	<b>62.2</b>	66.9	<b>69.9</b>	<b>20.5</b>
	<b>Joint Learned Class-Specific</b>		<b>62.3</b> (+0.6)~	<b>67.4</b>	69.2	28.9

\* $p$ -value < 0.05; \*\* $p$ -value < 0.01; ~ $p$ -value  $\geq$  0.05 (compared to Heuristic<sup>†</sup> TRA w/o TEA or Heuristic<sup>‡</sup> TRA w/ Heuristic TEA)  
<sup>†</sup>We pretrain these models with training data from site A and fine-tune with validation data from site B.  
<sup>‡</sup>We train these models with both training data from site A and validation data from site B.

Model	Site A Training-time data augmentation	Site B Test-time data augmentation	Site B			
			DSC	SEN	PRC	HD
DeepMedic [18]	None	None	14.9	11.6	45.3	42.6
	Heuristic [18]	None	46.4	43.2	59.4	26.9
	Heuristic <sup>†</sup> [18]	None	56.7	46.4	77.5	<b>9.4</b>
	Heuristic <sup>‡</sup> [18]	None	<b>69.3</b>	<b>67.2</b>	73.5	<b>15.1</b>
	Learned <sup>‡</sup> [7], [24], [27]	None	65.8	62.8	<b>75.1</b>	21.9
	<b>Learned Class-Specific<sup>‡</sup></b>	None	<b>70.0</b> (+0.7)~	<b>68.0</b>	<b>75.9</b>	18.7
	Heuristic <sup>‡</sup> [18]	Heuristic [16]	69.4	66.3	76.5	<b>8.0</b>
	<b>Learned Class-Specific<sup>‡</sup></b>	Heuristic [16]	69.9	66.3	<b>80.0</b>	<b>8.0</b>
	<b>Learned Class-Specific<sup>‡</sup></b>	Learned [20], [32]	<b>70.2</b>	67.7	<b>77.6</b>	15.3
	<b>Joint Learned Class-Specific<sup>‡</sup></b>		<b>72.8</b> (+3.4)**	<b>71.0</b>	76.6	<b>7.9</b>

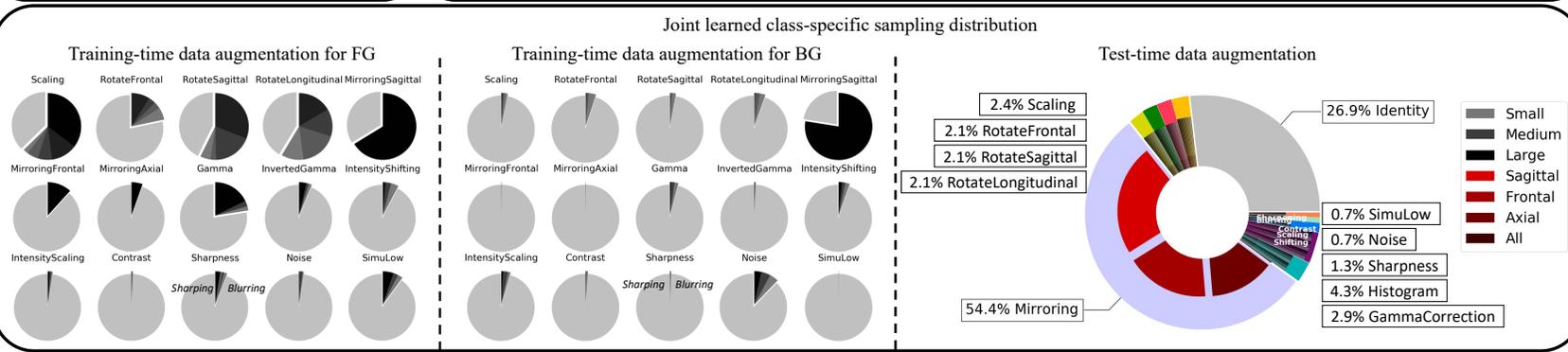
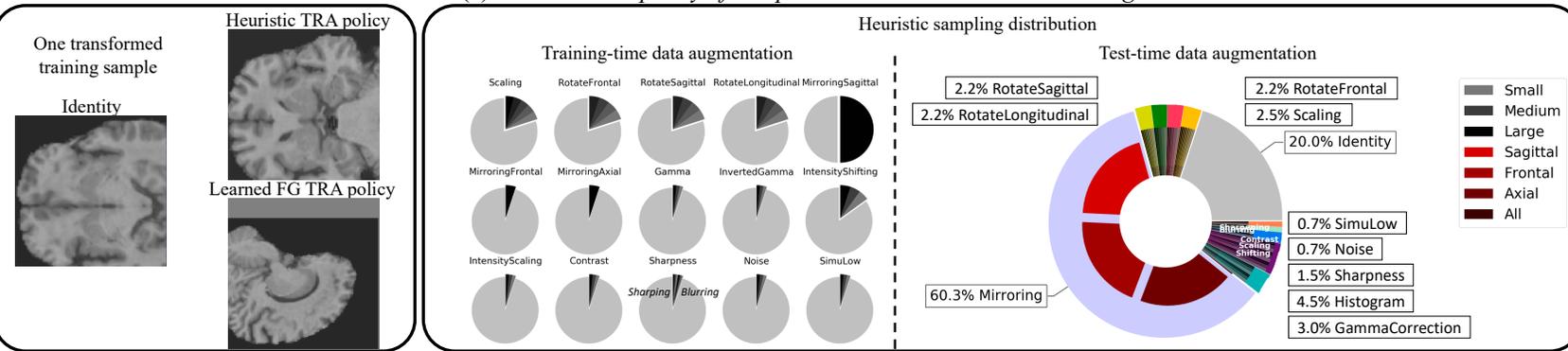
\* $p$ -value < 0.05; \*\* $p$ -value < 0.01; ~ $p$ -value  $\geq$  0.05 (compared to Heuristic<sup>†</sup> TRA w/o TEA or Heuristic<sup>‡</sup> TRA w/ Heuristic TEA)  
<sup>†</sup>We pretrain these models with training data from site A and fine-tune with validation data from site B.  
<sup>‡</sup>We train these models with both training data from site A and validation data from site B.

# Joint Optimization of Class-Specific Training- and Test-time data augmentation

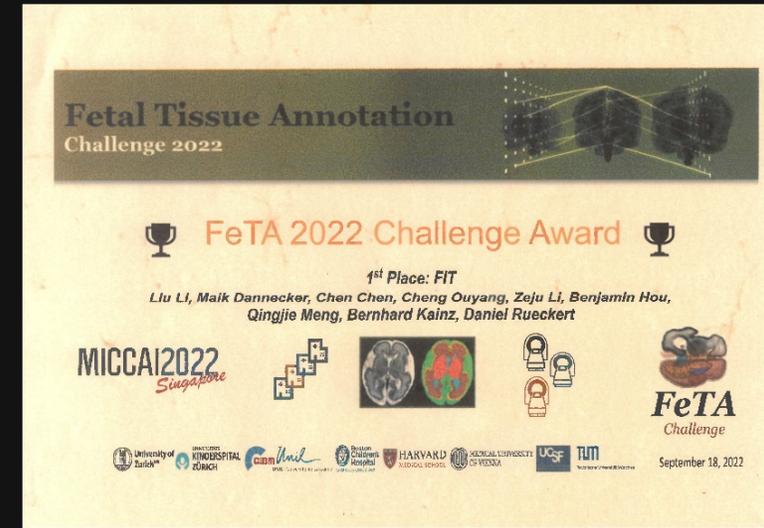
## Results

- The learned policies would adopt larger transformations to the foreground than the background samples, **implicitly alleviating the class imbalance issue.**

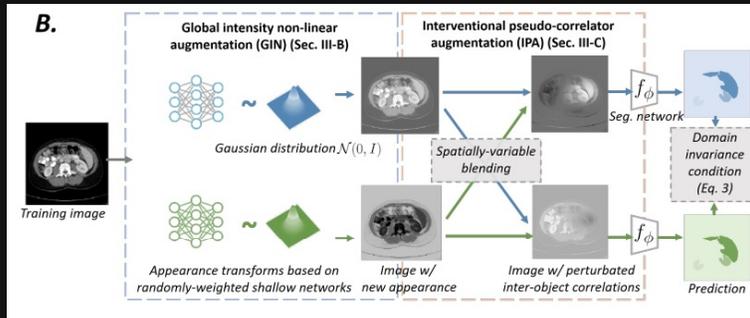
(a) Final learned policy of DeepMedic with 100% ATLAS training data



Serve as a key component for winning a robust learning challenge

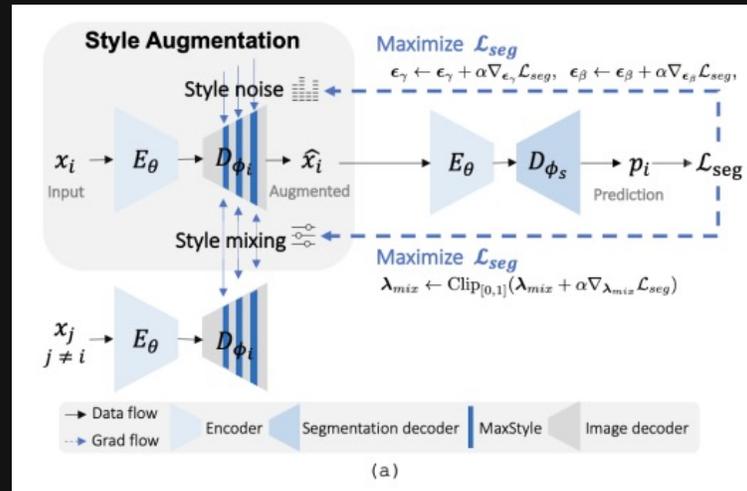


## Domain Generalization with **Random Kernels**



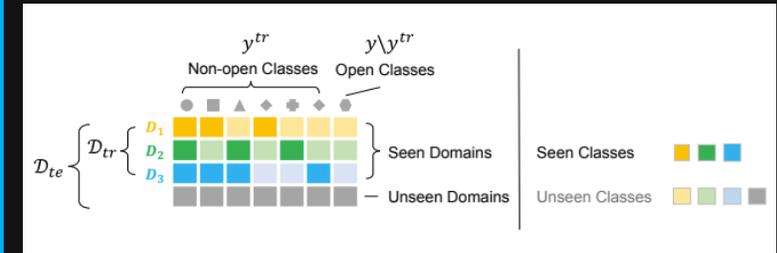
C. Ouyang, C. Chen, S. Li, **Z. Li et al.**  
 “Causality-inspired single-source domain generalization for medical image segmentation”, TMI, 2022.

## Domain Generalization with **Adversarial Training**



C. Chen, **Z. Li et al.** “MaxStyle: Adversarial style composition for robust medical image segmentation”, MICCAI, 2022.

## Domain Generalization for **Long-Tailed Classification**

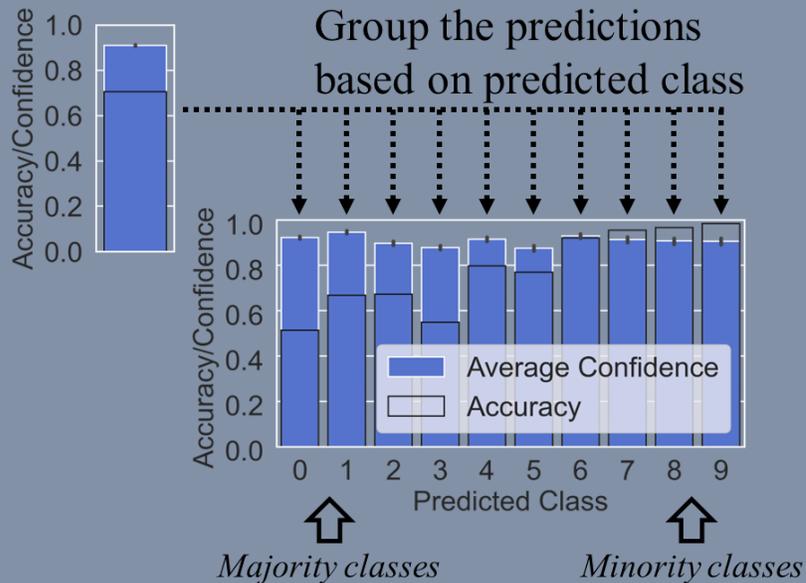


X. Gu, Y. Guo, **Z. Li et al.** “Tackling long-tailed category distribution under domain shifts”, ECCV, 2022.

## Confidence-based estimation

- Effect of **class imbalance** on confidence-based model evaluation methods.  
Optimization with validation set  $\mathcal{D}^V = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$

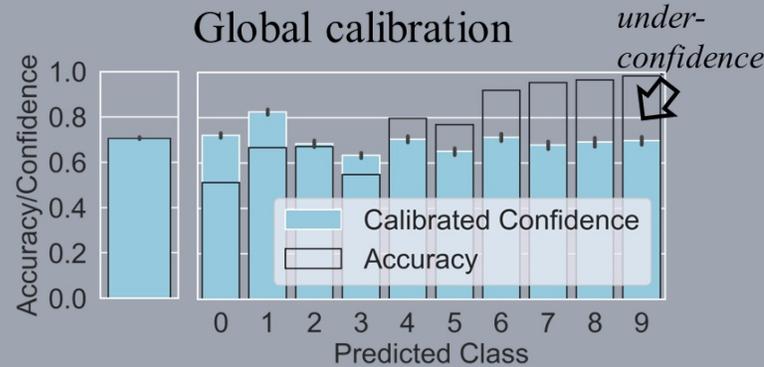
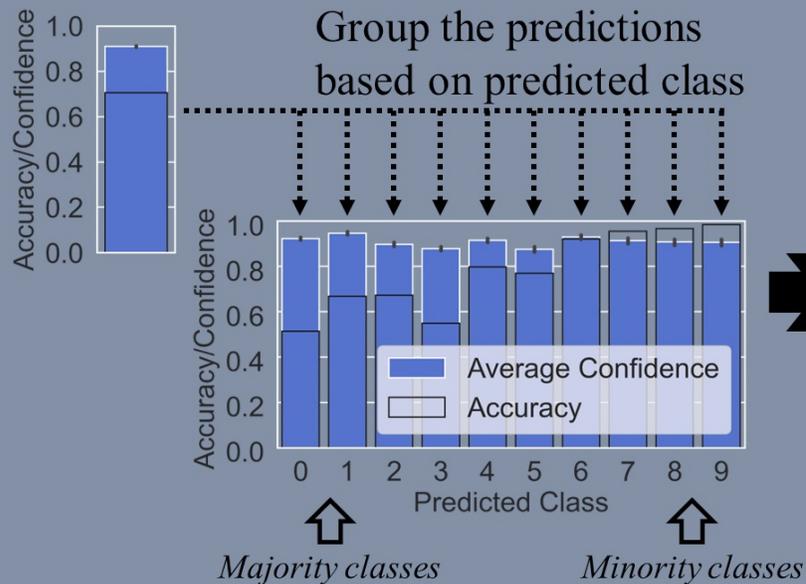
Deployment with test set  $\mathcal{D}^{Te} = \{\mathbf{x}_i^{Te}\}_{i=1}^M$



## Confidence-based estimation

- Effect of **class imbalance** on confidence-based model evaluation methods.  
Optimization with validation set  $\mathcal{D}^V = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$

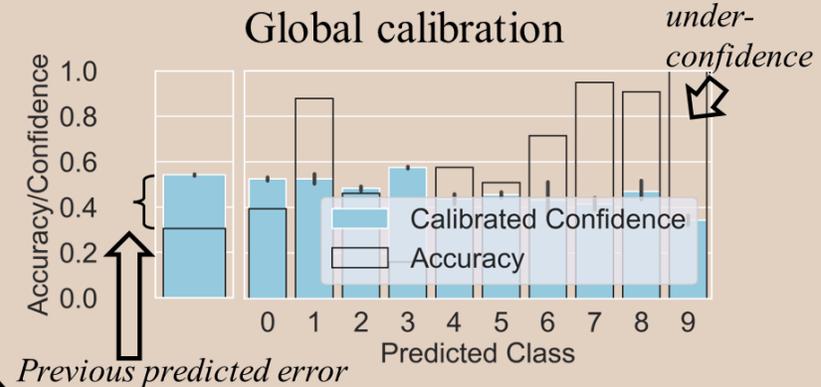
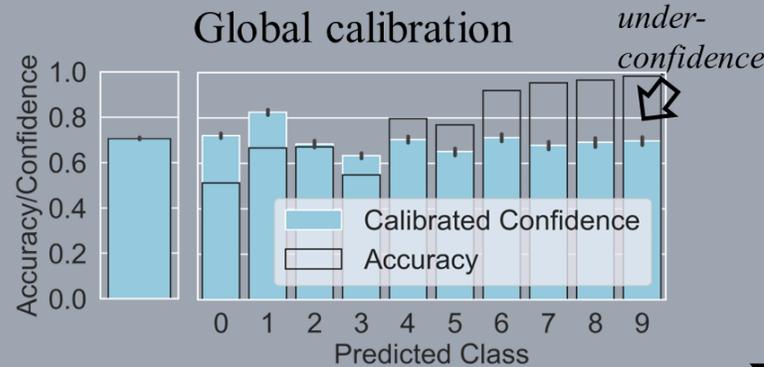
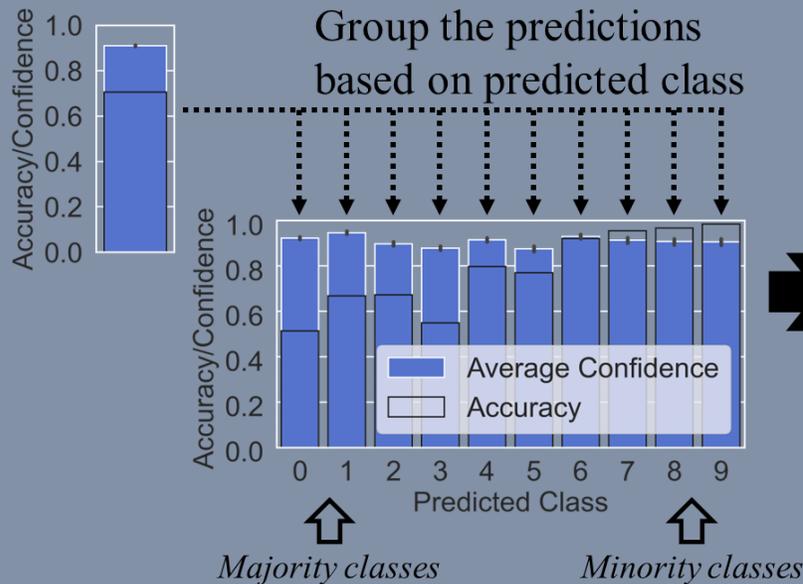
Deployment with test set  $\mathcal{D}^{Te} = \{\mathbf{x}_i^{Te}\}_{i=1}^M$



## Confidence-based estimation

- Effect of **class imbalance** on confidence-based model evaluation methods.  
Optimization with validation set  $\mathcal{D}^V = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$

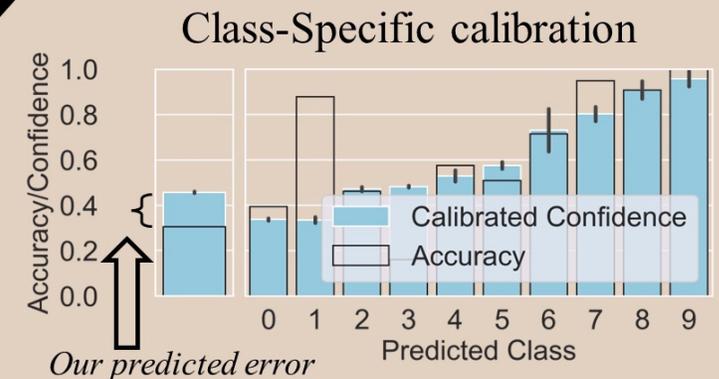
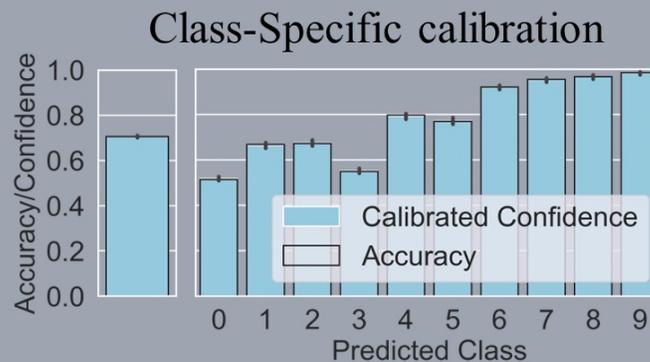
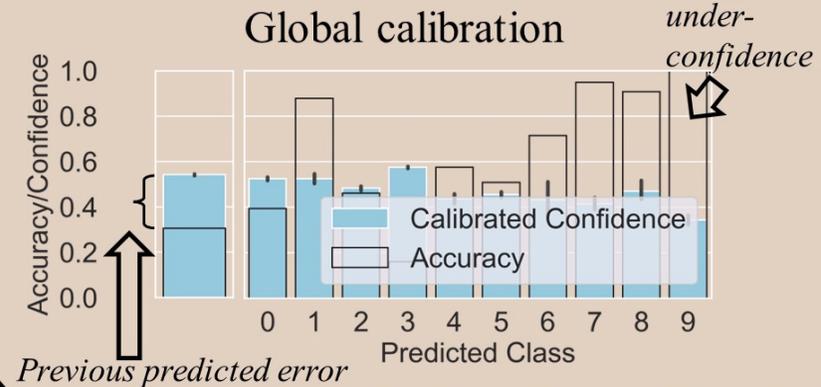
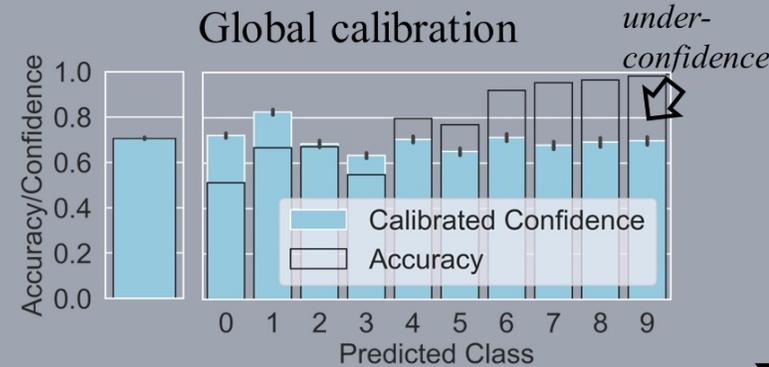
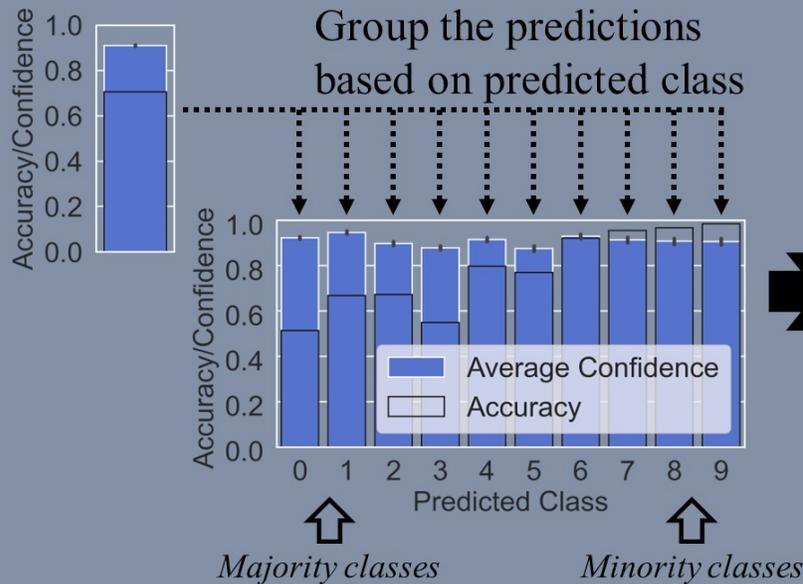
Deployment with test set  $\mathcal{D}^{Te} = \{\mathbf{x}_i^{Te}\}_{i=1}^M$



## Confidence-based estimation

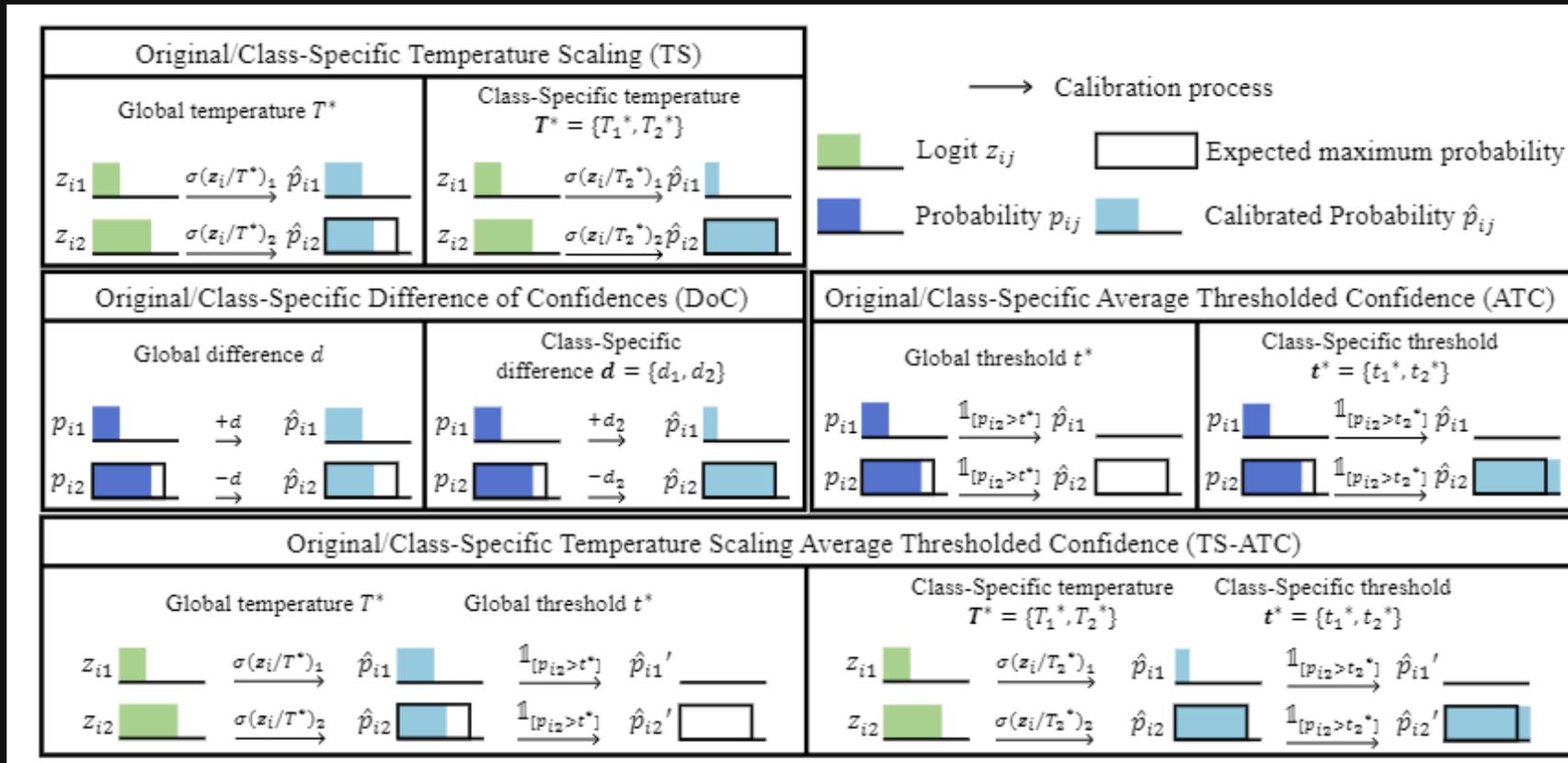
- Effect of **class imbalance** on confidence-based model evaluation methods.  
 Optimization with validation set  $\mathcal{D}^V = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$

Deployment with test set  $\mathcal{D}^{Te} = \{\mathbf{x}_i^{Te}\}_{i=1}^M$



## Method

- Introduce **class-wise calibration** within the framework of performance estimation for imbalanced datasets.

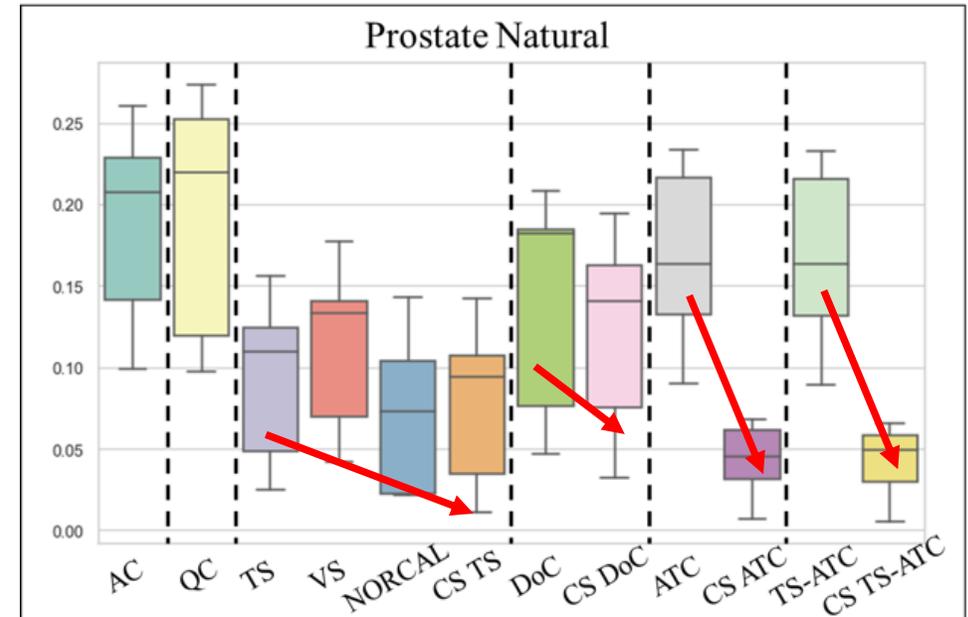


## Results

- Consistently improve model estimation accuracy, especially for **segmentation tasks**.

Task	Classification			Segmentation		
	CIFAR-10	HAM10000		ATLAS	Prostate	
Training dataset	CIFAR-10	Synthetic	Natural	Synthetic	Synthetic	Natural
Test domain shifts	Synthetic	Synthetic	Natural	Synthetic	Synthetic	Natural
AC	31.3 ± 8.2	12.3 ± 5.1	20.1 ± 13.4	35.6 ± 2.1	8.7 ± 4.9	18.7 ± 5.9
QC [30]	—	—	—	3.0 ± 1.7	5.2 ± 6.6	19.3 ± 7.1
TS [12]	5.7 ± 5.6	3.9 ± 4.3	12.1 ± 8.3	9.7 ± 2.5	3.7 ± 5.4	9.2 ± 4.9
VS [12]	3.8 ± 2.1	4.2 ± 4.2	13.6 ± 9.6	11.4 ± 2.5	4.8 ± 5.1	11.2 ± 4.9
NORCAL [29]	7.6 ± 3.8	4.2 ± 4.6	13.7 ± 9.6	6.7 ± 2.4	5.8 ± 5.7	7.3 ± 4.7
CS TS	5.5~ ± 5.6	3.7~ ± 4.0	11.9~ ± 8.0	1.6** ± 1.8	3.0** ± 5.7	7.8** ± 4.8
DoC [11]	10.8 ± 8.2	4.6 ± 5.0	15.3 ± 9.7	4.2 ± 3.2	3.7 ± 5.8	13.9 ± 6.5
CS DoC	9.4** ± 7.2	4.5~ ± 4.9	14.7* ± 9.2	1.3** ± 1.9	3.5* ± 6.1	12.1* ± 5.9
ATC [10]	4.6 ± 4.4	3.4 ± 3.9	7.1 ± 6.3	30.4 ± 1.8	8.6 ± 3.3	16.7 ± 5.3
CS ATC	2.8** ± 2.9	3.3~ ± 4.8	5.8~ ± 7.6	1.6** ± 1.5	1.1** ± 1.7	4.3** ± 2.2
TS-ATC [10, 12]	5.3 ± 3.9	4.2 ± 4.2	7.3 ± 7.1	30.4 ± 1.8	8.5 ± 3.3	16.7 ± 5.3
CS TS-ATC	2.7** ± 2.3	4.2~ ± 5.4	5.9~ ± 8.4	1.3** ± 1.4	1.2** ± 1.7	4.2** ± 2.2

\* $p$ -value < 0.05; \*\* $p$ -value < 0.01; ~ $p$ -value ≥ 0.05 (compared with their class-agnostic counterparts)



## Software

- Based on our algorithm, we develop an open-source software, MOVAL, to help practitioners evaluate their model performance.



### Reliable

MOVAL facilitates the assessment of pre-trained models across diverse scenarios, aiding practitioners in estimating performance.



### Versatile

MOVAL not only calculates but also calibrates various confidence scores beyond the maximum class probability.



### Effective

MOVAL expands existing performance estimation algorithms and demonstrates state-of-the-art results, particularly on real-world long-tailed datasets.



### User-Friendly

MOVAL can be effortlessly installed as a Python module and supports the NumPy array data format.



### Universal

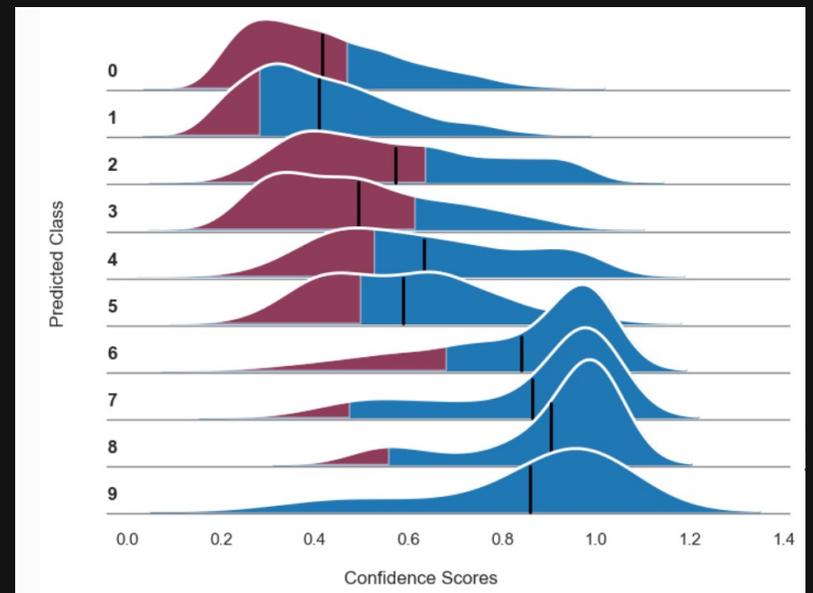
MOVAL accommodates various applications within a unified framework, encompassing tasks such as classification, 2D segmentation, and 3D segmentation.



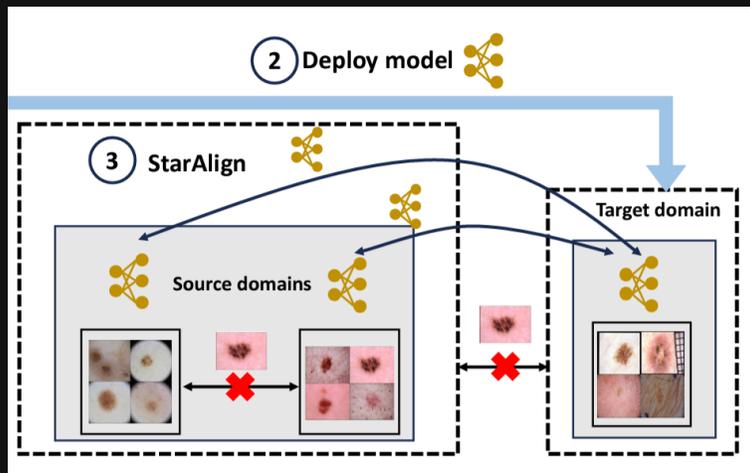
### Modular

MOVAL comprises distinct modules for confidence score calculation, calibration, and optimization, allowing for easy extension and customization.

```
>>> import moval
>>> import numpy as np
>>> logits = np.random.randn(1000, 10)
>>> gt = np.random.randint(0, 10, (1000))
>>> moval_model = moval.MOVAL()
>>> moval_model.fit(logits, gt)
>>> estim_acc = moval_model.estimate(logits)
```

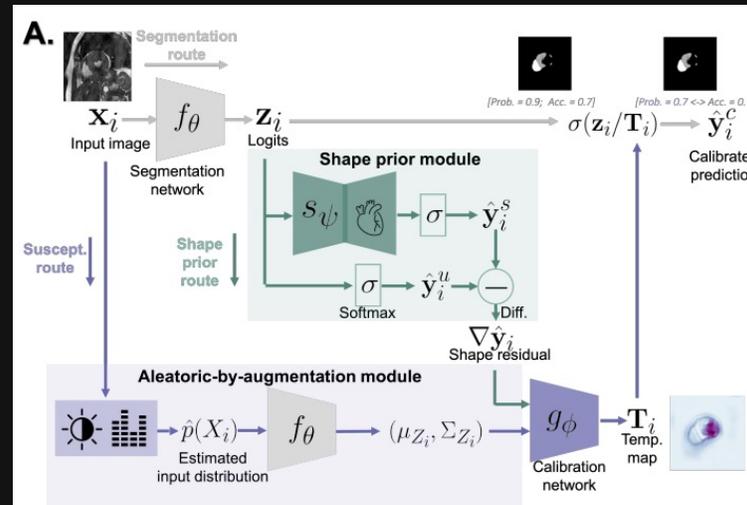


## Post-deployment Adaption under Federated Setting



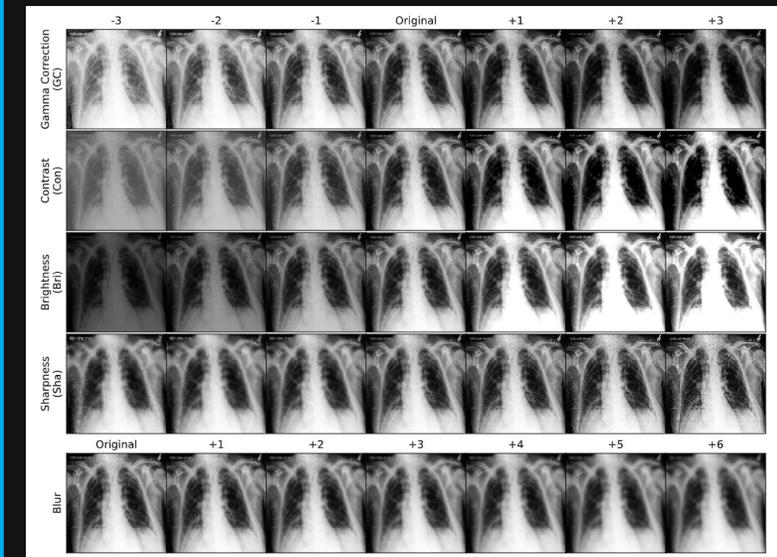
F. Wagner, Z. Li *et al.* "Post-deployment adaptation with access to source data via federated learning and source-target remote gradient alignment", MICCAI-MLMI workshop, 2023.

## Model Calibration for Image Segmentation



C. Ouyang, S. Wang, C. Chen, Z. Li *et al.* "Improved post-hoc probability calibration for out-of-domain MRI segmentation", MICCAI-UNSURE workshop, 2022.

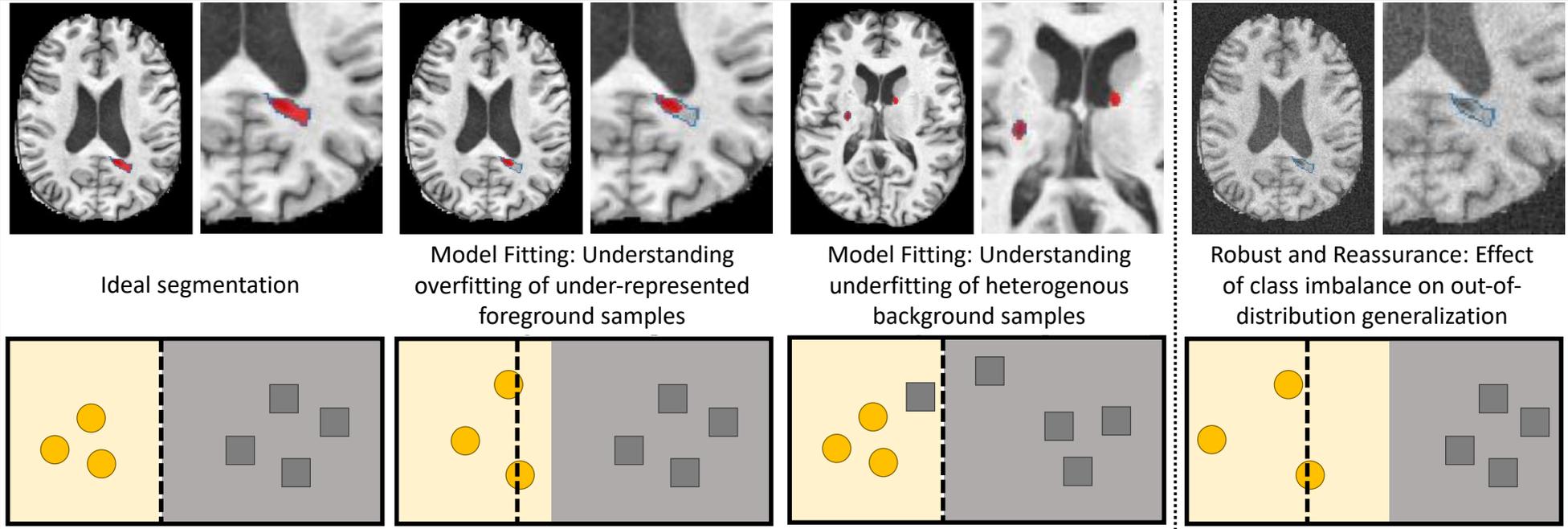
## Stress Testing with Fairness Concerns



M. Islam, Z. Li *et al.* "Progressive stress testing of model robustness in medical image classification", MICCAI-UNSURE workshop, 2023.

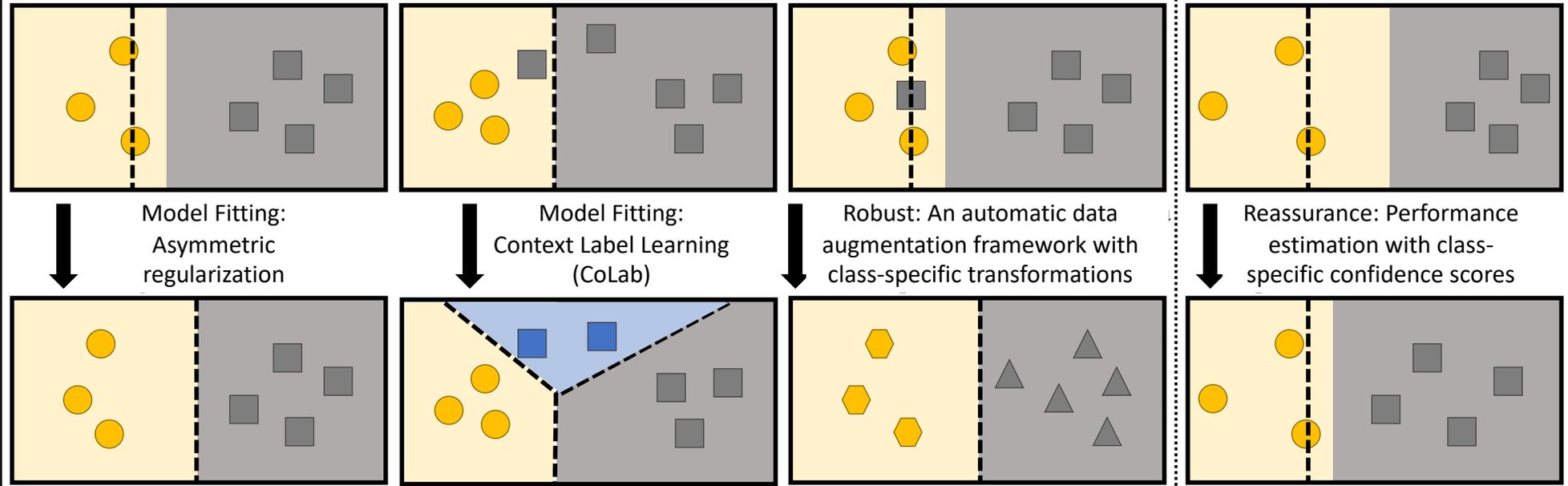
## Summary

- **Class imbalance** and **domain shifts**, which exist simultaneously in real world datasets, limit the performance of modern machine learning models when deployed to medical image segmentation.
- Class imbalance cause under-segmentation because of **overfitting foreground samples**, while over-segmentation because of **underfitting background samples**.



## Summary

- **Asymmetric** loss functions and regularization techniques help counter overfitting under class imbalance by enlarging foreground sample variances.
- **Context labels** help alleviate underfitting under class imbalance.
- **Class-specific parameters** are beneficial for improving data augmentation and tackling domain shifts.



# Take Home Message

Enable **robust and safe deployment** of machine learning in medical image segmentation!

Model Performance



Deployment Environment



Optimization Algorithms



Data Distributions



**Reassurance:** Calibrating model prediction with to match class-wise performance.

**Robustness:** Aligning the training and test data distribution with data augmentation.

**Model-Fitting:** Alleviating class imbalance with advanced learning strategies.

## ➤ Model-Fitting:

1. **Z. Li et al.** "Analyzing overfitting under class imbalance in neural networks for image segmentation", TMI, 2020.
2. **Z. Li et al.** "Context label learning: improving background class representations in semantic segmentation", TMI, 2023.
3. **Z. Li et al.** "Deep learning based radiomics (DLR) and its usage in noninvasive idh1 prediction for low grade glioma", Sci. Rep., 2017.
4. **Z. Li et al.** "Deepvolume: brain structure and spatial connection-aware network for brain MRI super-resolution", TCybern, 2019.
5. **Z. Li et al.** "Brain tumor segmentation using an adversarial network.", MICCAI-brainlesion workshop, 2017.

## ➤ Robustness:

1. **Z. Li et al.** "Joint optimization of class-specific training- and test-time data augmentation in segmentation", TMI, 2023.
2. C. Ouyang, C. Chen, S. Li, **Z. Li et al.** "Causality-inspired single-source domain generalization for medical image segmentation", TMI, 2022.
3. C. Chen, **Z. Li et al.** "MaxStyle: Adversarial style composition for robust medical image segmentation", MICCAI, 2022.
4. X. Gu, Y. Guo, **Z. Li et al.** "Tackling long-tailed category distribution under domain shifts", ECCV, 2022.
5. C. Chen, C. Ouyang, **Z. Li et al.** "Enhancing mr image segmentation with realistic adversarial data augmentation", MedIA, 2022

## ➤ Reassurance:

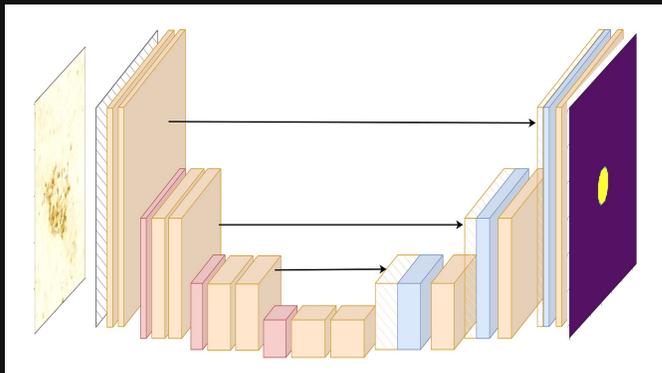
1. **Z. Li et al.** "Estimating model performance under domain shifts with class-specific confidence scores", MICCAI, 2022.
2. F. Wagner, **Z. Li et al.** "Post-deployment adaptation with access to source data via federated learning and source-target remote gradient alignment", MICCAI-MLMI workshop, 2023.
3. C. Ouyang, S. Wang, C. Chen, **Z. Li et al.** "Improved post-hoc probability calibration for out-of-domain MRI segmentation", MICCAI-UNSURE workshop, 2022.
4. M. Islam, **Z. Li et al.** "Progressive stress testing of model robustness in medical image classification", MICCAI-UNSURE workshop, 2023.
5. **Z. Li et al.** "Encoding ct anatomy knowledge for unpaired chest x-ray image decomposition", MICCAI, 2019.

# Post-doc Working on Building Macaque Connectivity Atlas

38/40

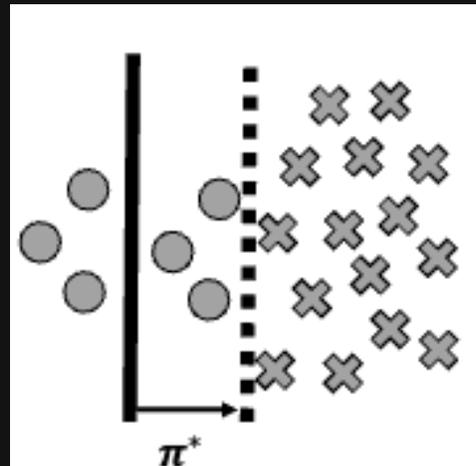
## Cell Segmentation

End-to-end neural networks such as U-Net / Transformer.



## Semi-supervised Learning

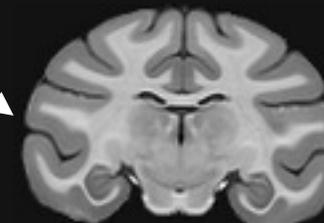
Propose an algorithm that enhances the quality of pseudo-labelling for imbalanced semi-supervised learning.



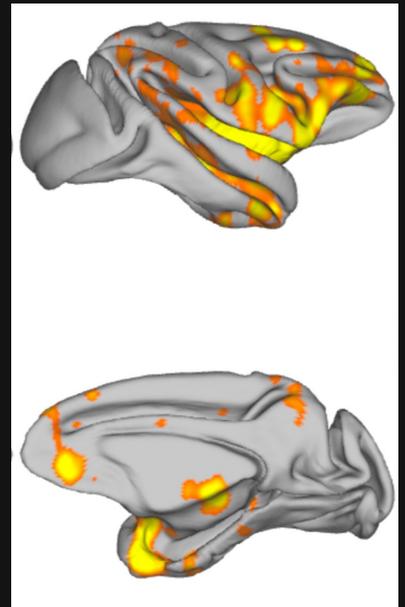
## Registration to structure MRI



Registration

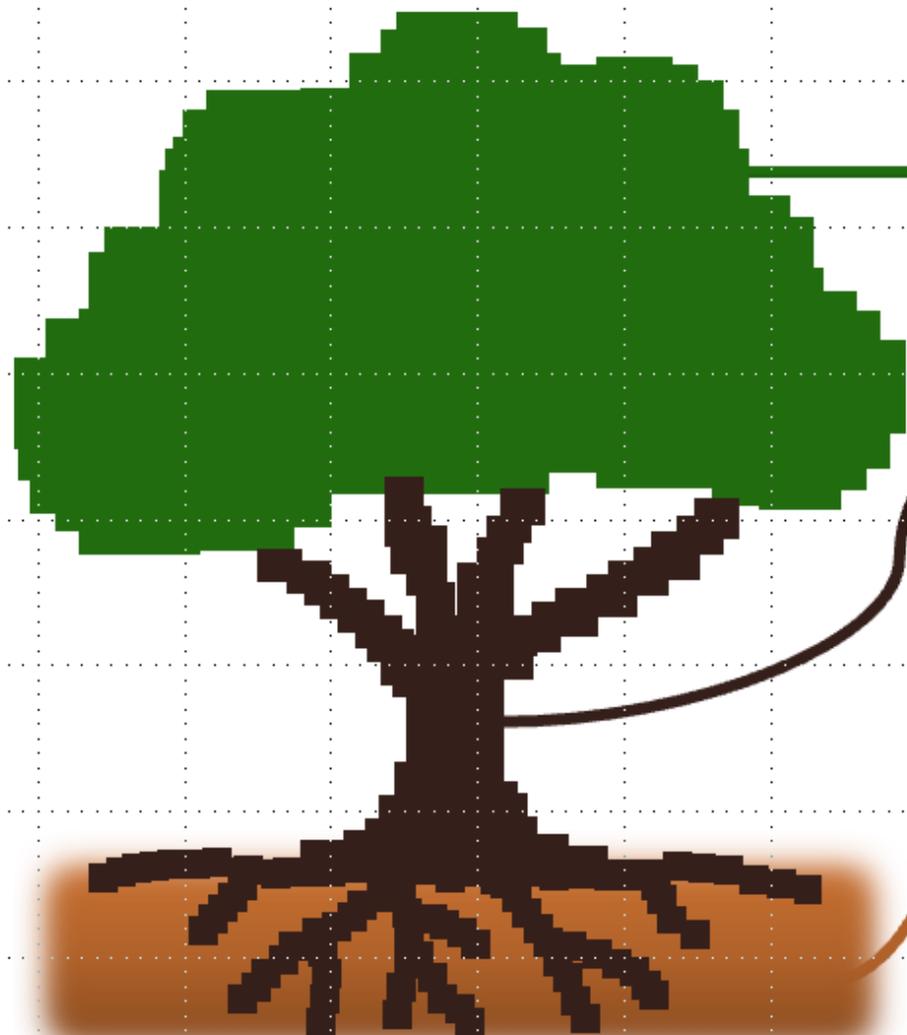


The first tracer atlas that is both fine-grained and quantitative



# Future Vision: Towards Building Foundation Models for Neuro-oncology and Neuroscience

39/40



Direction 3: Developing foundation models for brain tumor analysis and computational neuroscience



Direction 2: Building transfer learning tools for foundation models to ensure both safety and flexibility



Direction 1: Self-supervised pre-training strategies for multi-task and multi-modal foundation models

## Goal:

- Developing next-generation high-performance multi-tasking foundation models for neuro-oncology and neuroscience
- Building tools for transferring knowledge from foundation models to enhance medical image analysis



Paper & code



# Thank you!

**Zeju Li, PhD**

FMRIB Center, Nuffield Department of Clinical Neurosciences,  
University of Oxford, UK.

Email: [zeju.li@ndcn.ox.ac.uk](mailto:zeju.li@ndcn.ox.ac.uk)

March 2024

Statistics

3,288 citations

12 first authored peer-reviewed journal/conference

Academic Service

Area Chair, MICCAI 2024

Reviewer, CVPR/MICCAI, TMI/Media