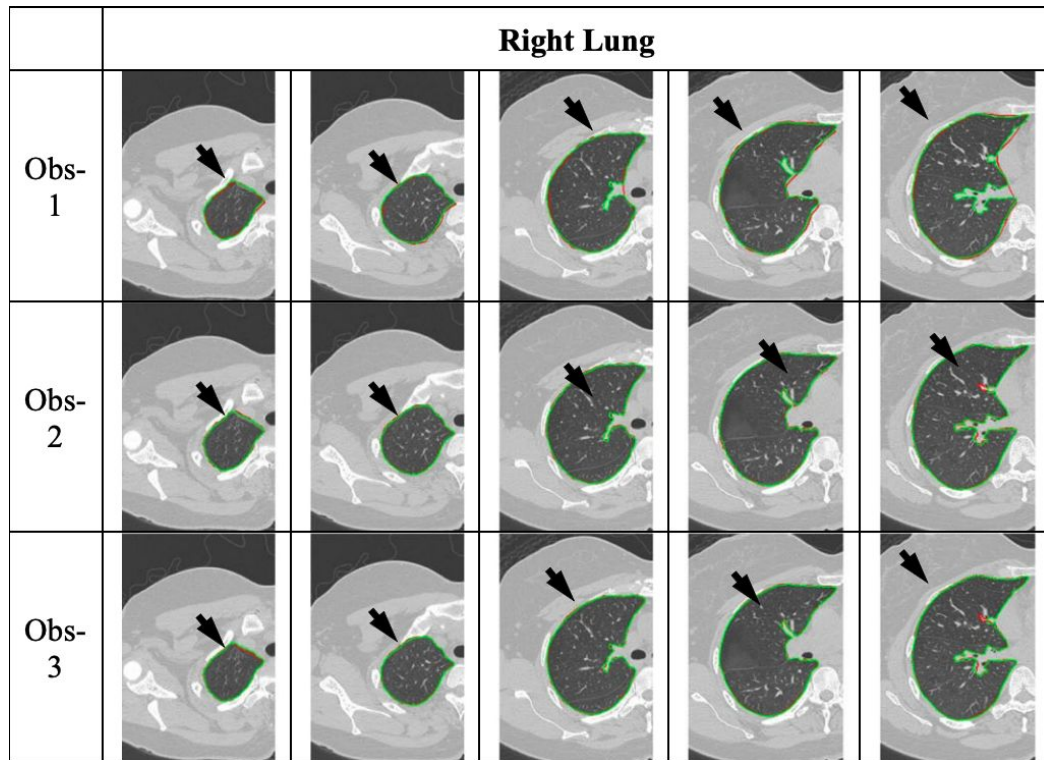


# Disentangling Human Error from the Ground Truth in Segmentation of Medical Images

Le Zhang, Ryutaro Tanno, Mou-Cheng Xu, Chen Jin, Joseph Jacob, Olga Ciccarelli, Frederik Barkhof and Daniel C. Alexander

*NeurIPS (2020)*

# Problem



- Segmentation of anatomical structures suffers from **high inter-observer variability** in medical images (biases, level of expertise, etc.)
- Noisy labels limits the performance of automatic segmentation algorithms

# Problem

Existing approaches (to deal with multiple annotations):

- Majority vote
- Weighing annotation according to estimated reliability of expert (image-wise of pixel-wise)

# Objective

- Jointly learn from noisy observations the **reliability of individual annotators** and the **true segmentation label distributions**  
→ Produce a better segmentation at test-time

## Hypothesis:

There is a **single, true segmentation map**, and each annotator produces a noisy approximation of it **according to its individual characteristics**

## Method

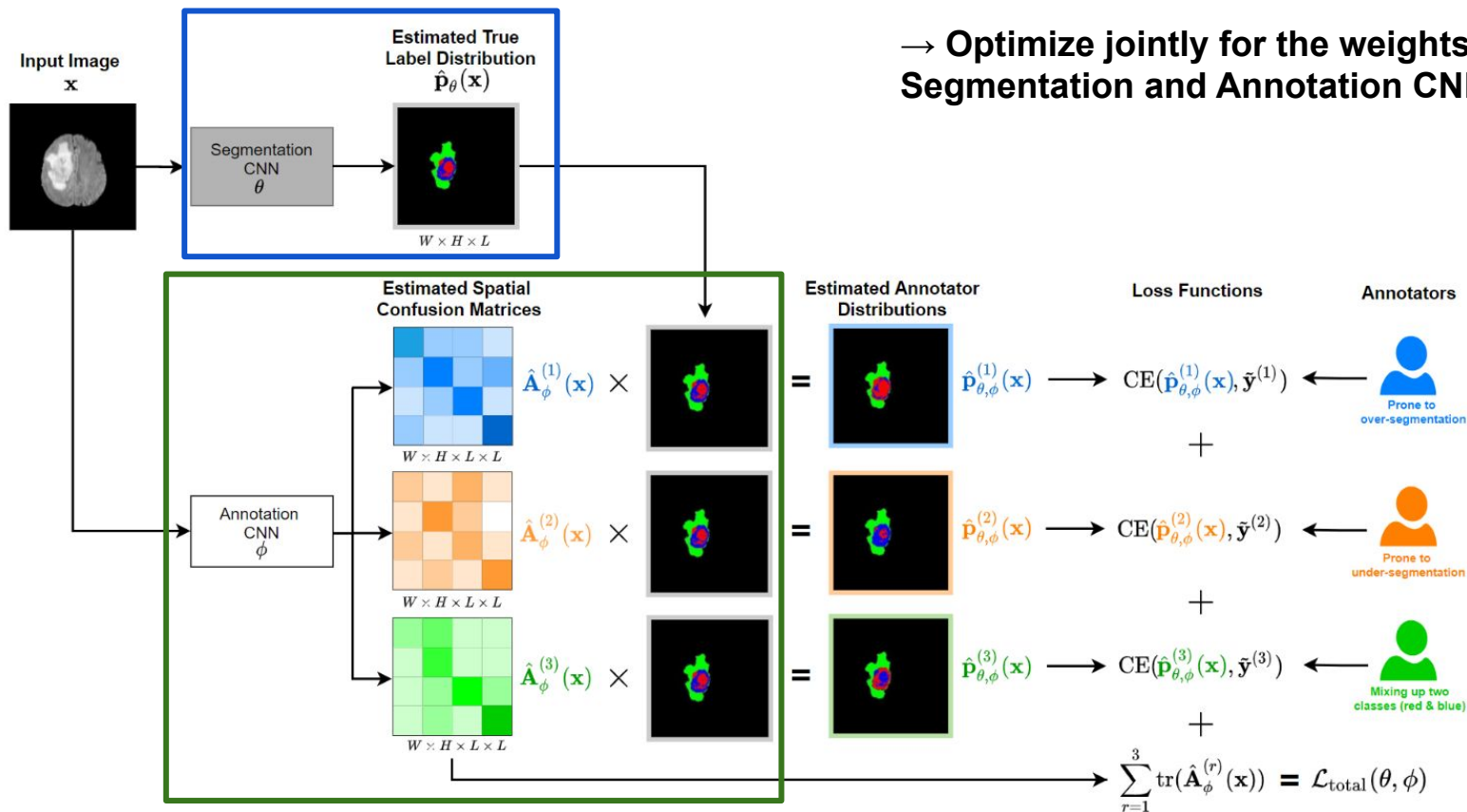
- Jointly train a Segmentation CNN and Annotation CNN
- Set of images, with multiple noisy segmentation masks

### Assumptions:

- 1) **Annotators** are statistically **independent**
- 2) **Annotations noise** is **independent** of input image

# Method

→ Optimize jointly for the weights of the Segmentation and Annotation CNN

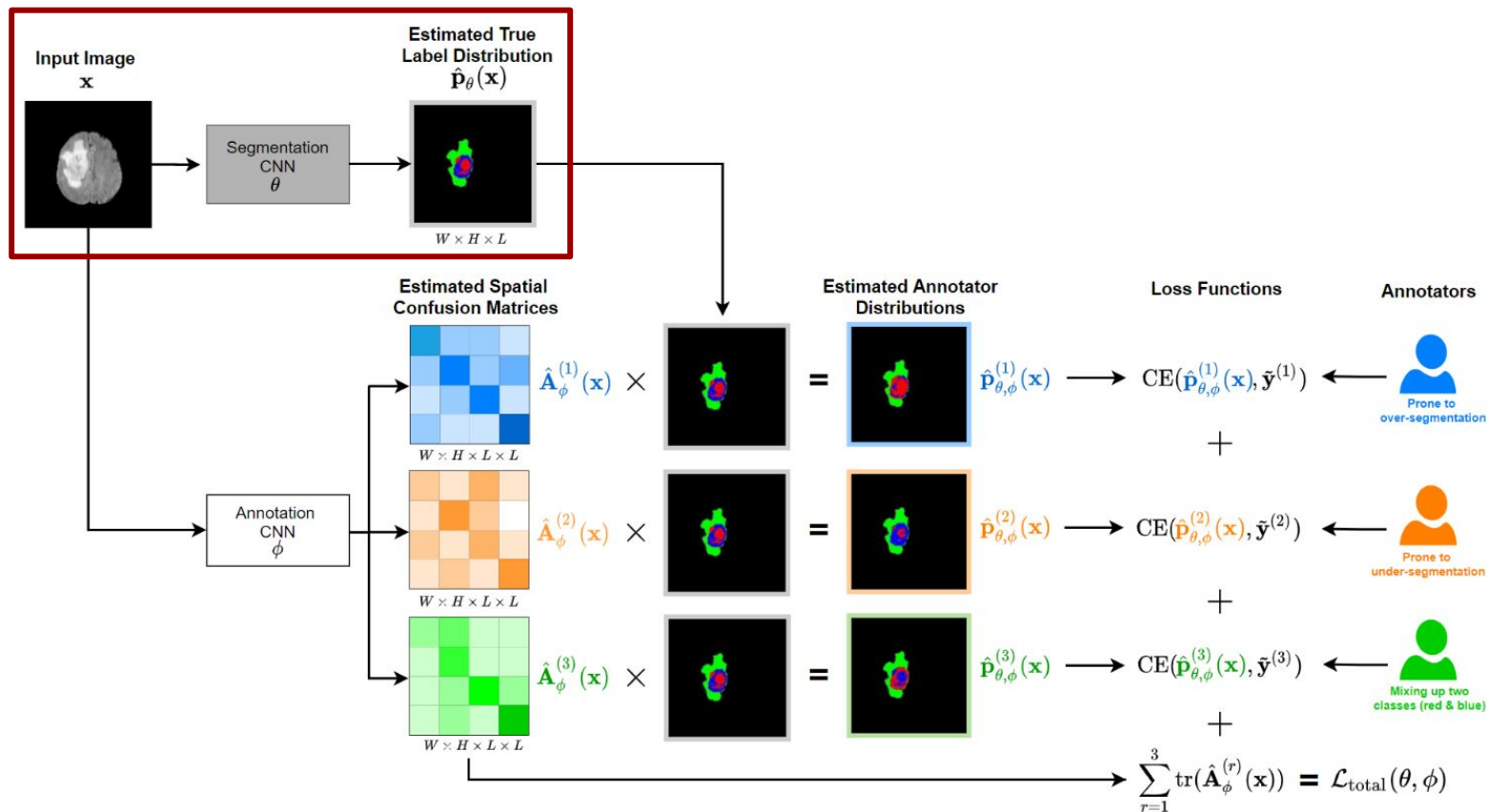


**Minimising trace** → encourages estimated annotators to be maximally unreliable



# Method

Test time





# Experiment

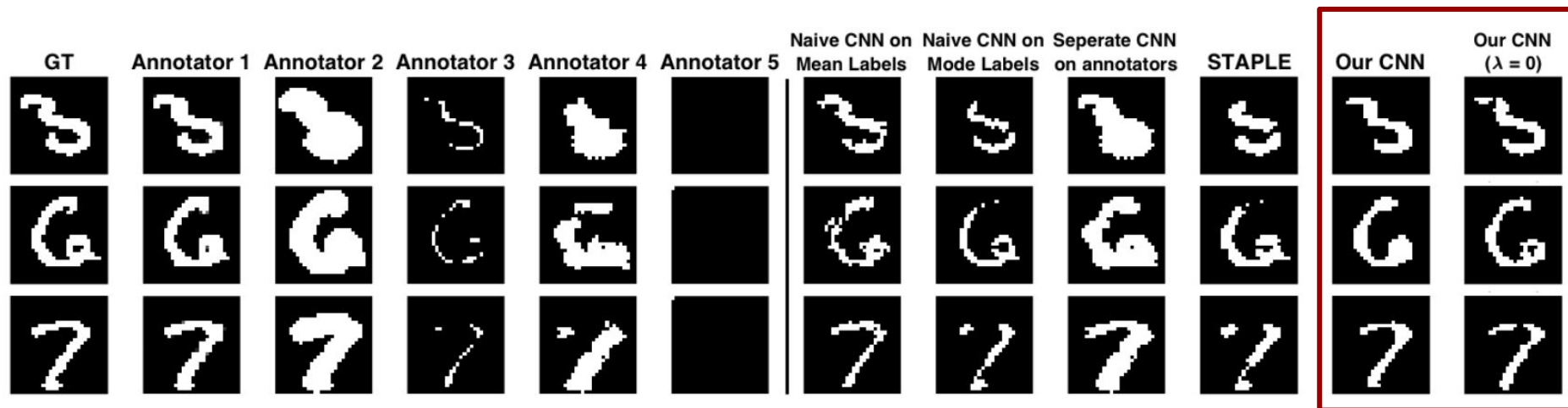
## - Dataset

- MNIST
  - MSLSC (multiple-sclerosis lesions)
  - BraTS (brain tumours)
  - LIDC-IDRI (lung abnormalities) → multiple annotations
- } Synthetically noisy labels

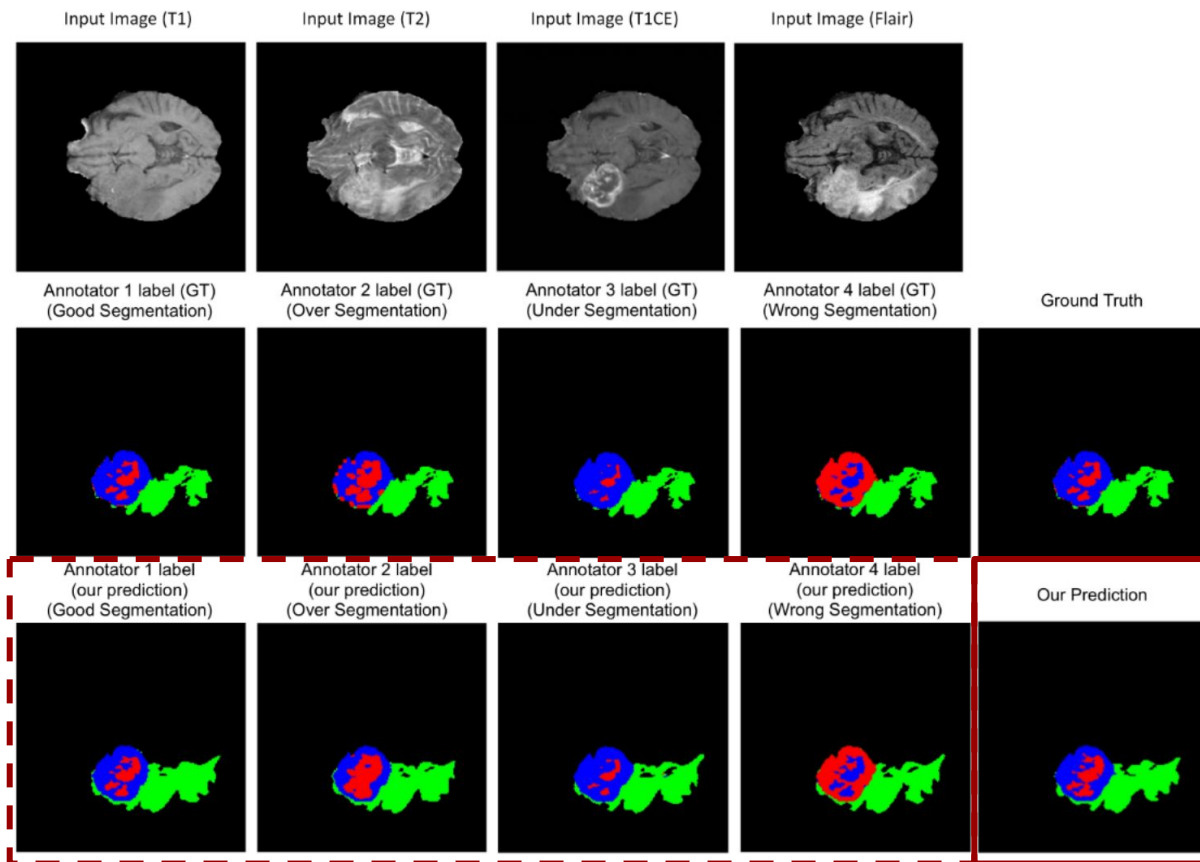
## - Baselines

- Mean of noisy labels
- Majority vote
- STAPLE
- Spatial STAPLE
- Probabilistic UNet

## Results



# Results



# Results

Models	BraTS DICE (%)	BraTS CM estimation	LIDC-IDRI DICE (%) $\uparrow$	LIDC-IDRI CM estimation $\downarrow$
Naive CNN on mean labels	$29.42 \pm 0.58$	n/a	$56.72 \pm 0.61$	n/a
Naive CNN on mode labels	$34.12 \pm 0.45$	n/a	$58.64 \pm 0.47$	n/a
Probabilistic U-net [24]	$40.53 \pm 0.75$	n/a	$61.26 \pm 0.69$	n/a
STAPLE [9]	$46.73 \pm 0.17$	$0.2147 \pm 0.0103$	$69.34 \pm 0.58$	$0.0832 \pm 0.0043$
Spatial STAPLE [14]	$47.31 \pm 0.21$	$0.1871 \pm 0.0094$	$70.92 \pm 0.18$	$0.0746 \pm 0.0057$
Ours with Global CMs	$47.33 \pm 0.28$	$0.1673 \pm 0.1021$	$70.94 \pm 0.19$	$0.1386 \pm 0.0052$
Ours without Trace	$49.03 \pm 0.34$	$0.1569 \pm 0.0072$	$71.25 \pm 0.12$	$0.0482 \pm 0.0038$
Ours	<b><math>53.47 \pm 0.24</math></b>	<b><math>0.1185 \pm 0.0056</math></b>	<b><math>74.12 \pm 0.19</math></b>	<b><math>0.0451 \pm 0.0025</math></b>
Oracle (Ours but with known CMs)	$67.13 \pm 0.14$	$0.0843 \pm 0.0029$	$79.41 \pm 0.17$	$0.0381 \pm 0.0021$

## Conclusion

A supervised segmentation method for jointly estimating the spatial **characteristics of labelling errors** (annotation noise) from multiple human annotators and the **ground-truth label distribution**.

→ Finds the **maximal amount of confusion** which explains the noisy observations well.

→ Improves **robustness** against label noise