

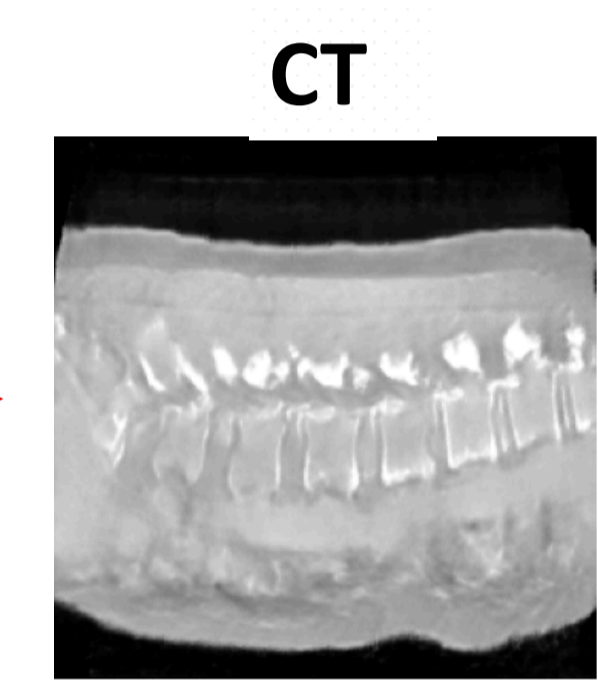
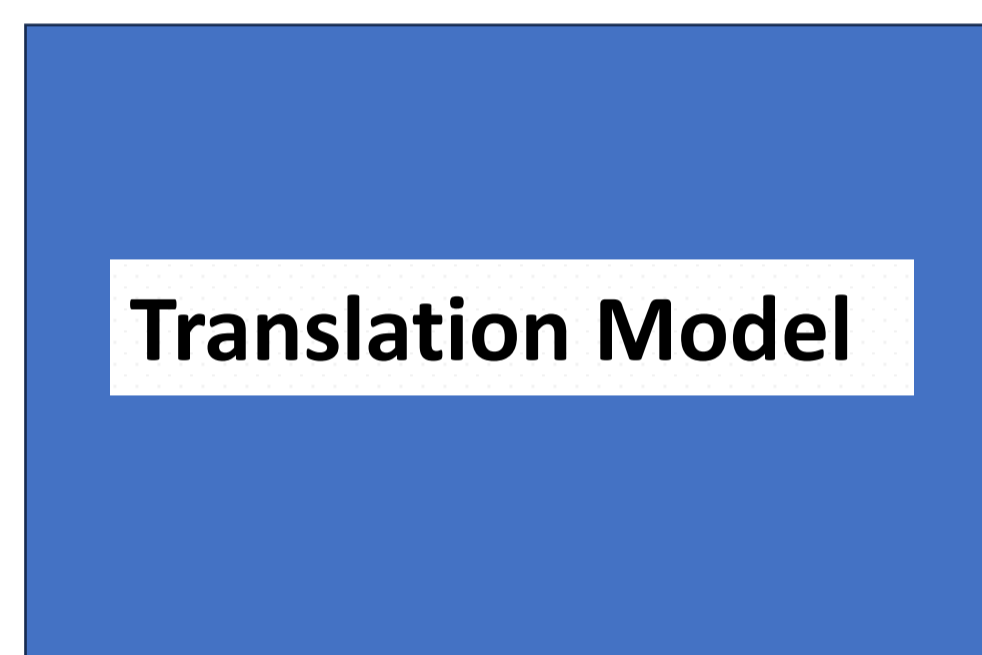
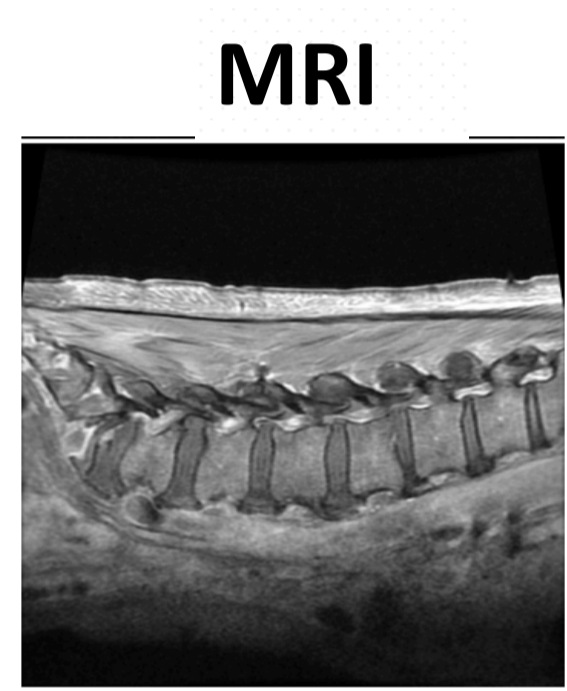
## Are existing perceptual metrics for image translation models actually useful for medical image translation?

## With what criteria should we evaluate them?

### Background

- Unpaired image-to-image translation: a common task in medical image computing.

#### Domain A



#### Domain B

- For example: transforming a lumbar spine MRI (source domain) to make it look like a CT (target domain)

## What are the desiderata for medical image translation?

### 1. Anatomical consistency w.r.t. the input image.

- Commonly measured with a **segmentation model trained in the target domain**, and applied to translated images.
  - A **standard metric in medical image translation papers**.
  - Limitations: need for labels and resources to train the segmentation model, **bias towards the task/object**, etc.

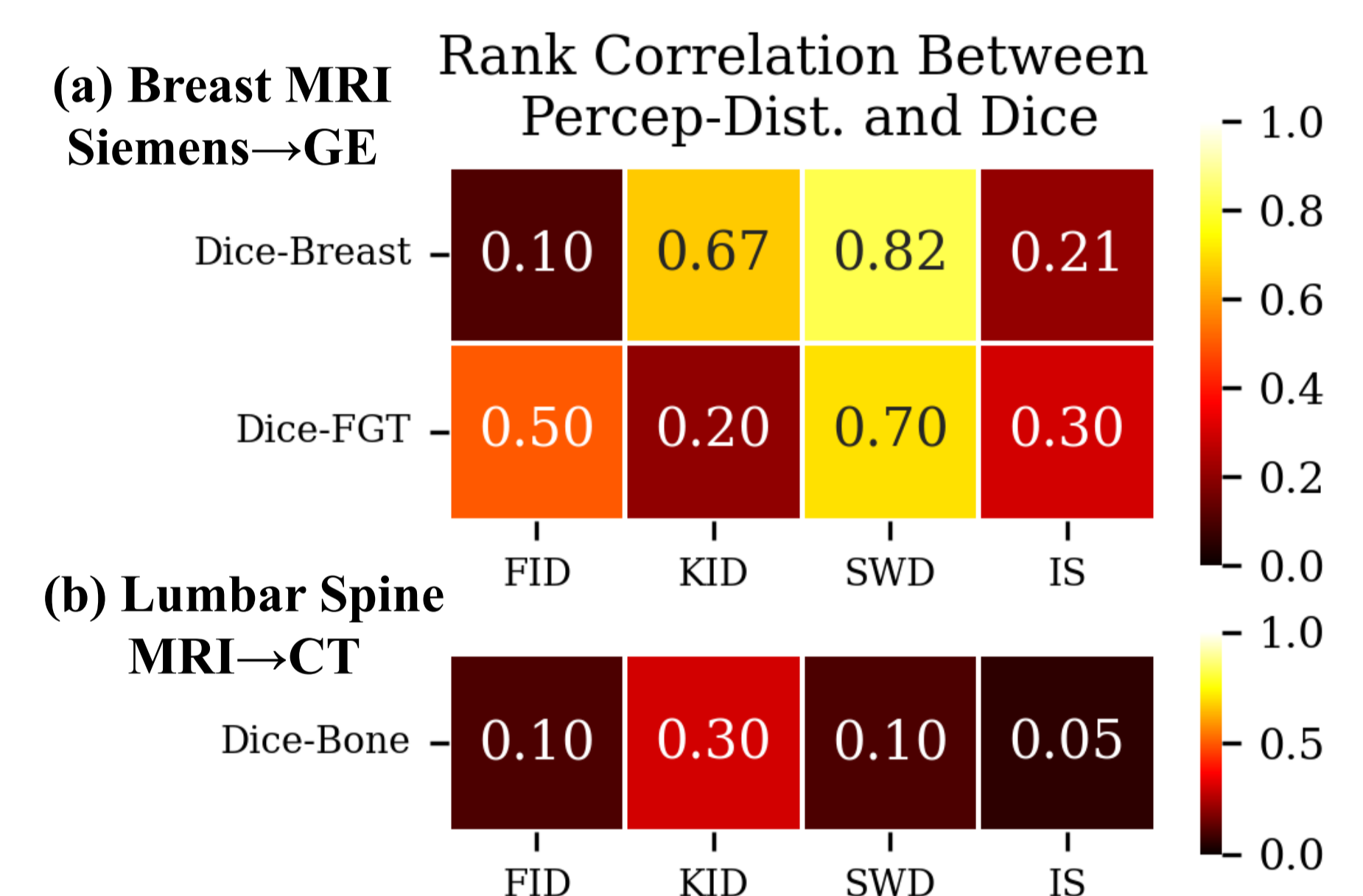
### 2. Overall perceptual quality and realism.

- Typically measured with **perceptual metrics from mainstream computer vision**: FID, IS, etc.
- These metrics are **task-agnostic**, but may fail to capture **local or global anatomical consistency and realism** in medical images!

## Are common perceptual metrics useful for medical image translation?

- Do any task-agnostic **perceptual metrics** also reliably correlate with **anatomical consistency**?

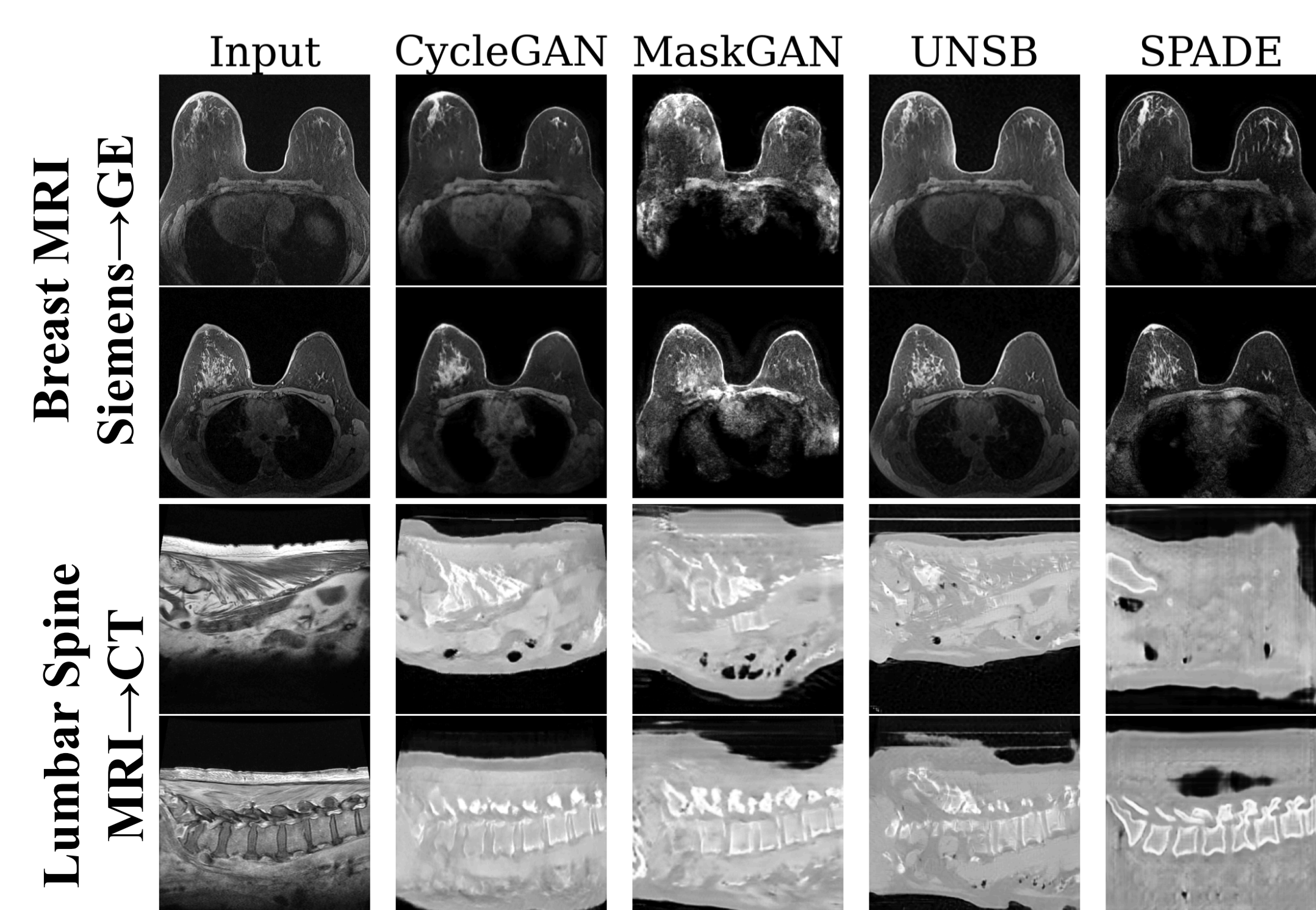
Sadly, no. 😞



Method	Breast MRI Siemens→GE Translation						Lumbar Spine MRI→CT Translation				
	Dice (↑)		Perceptual Metrics				Dice (↑)		Perceptual Metrics		
	Breast	FGT	FID*	KID	SWD	IS	Bone	FID*	KID	SWD	IS
None	0.927	0.651	144	0.069	705	2.58	0.007	323	0.300	1553	<b>2.93</b>
CycleGAN	<u>0.934</u>	0.529	<b>107</b>	<b>0.049</b>	556	2.73	0.229	210	<b>0.161</b>	960	<u>2.29</u>
MaskGAN	0.865	0.277	<u>118</u>	0.089	1037	<b>3.00</b>	0.158	248	0.217	1114	2.22
UNSB	<u>0.934</u>	0.646	156	0.079	756	2.46	0.138	<b>208</b>	<u>0.172</u>	<b>932</b>	2.14
SPADE†	<b>0.950</b>	<b>0.707</b>	119	0.067	500	2.91	<b>0.942</b>	251	0.242	1359	<u>2.29</u>

Table 1: Quantitative results for both translation tasks. Best and runner-up models are shown in bold and underlined according to each metric, respectively.

- FID is especially inconsistent!**
- SWD (**pixel-level**, not learned feature-level) may be good for certain datasets, but **still isn't consistent**.



**Conclusion:** we need better metrics for medical image translation that satisfy these desiderata!

### Contact us

- nicholas.konz@duke.edu
- maciej.mazurowski@duke.edu
- @nick\_konz on twitter

