

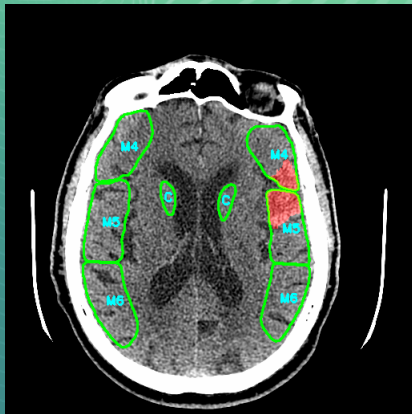
- ▶ Our projects;

- ▶ Our projects;
- ▶ Image restoration on PET-CT images;

- ▶ Our projects;
- ▶ Image restoration on PET-CT images;
- ▶ Unsupervised pretraining for segmentation of CT studies.

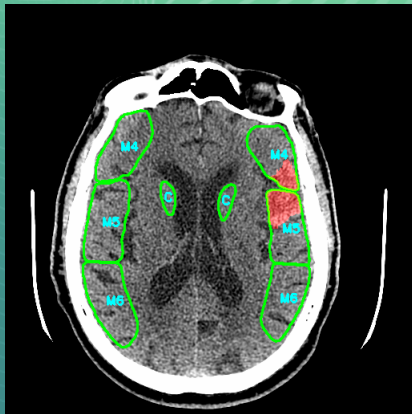
CT of head:

- Segmentation of lesions affected by acute stroke;



CT of head:

- ▶ Segmentation of lesions affected by acute stroke;
- ▶ Segmentation of regions of brain (ASPECTS);

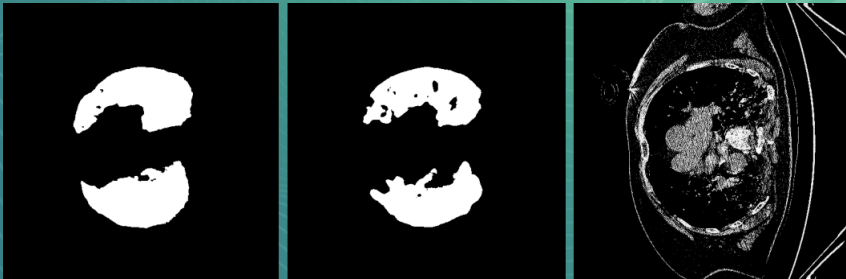


CT of head:

- ▶ Segmentation of lesions affected by acute stroke;
- ▶ Segmentation of regions of brain (ASPECTS);
- ▶ (jointly with AIRI) Detection of very early stroke on CT (up to 12 hours);

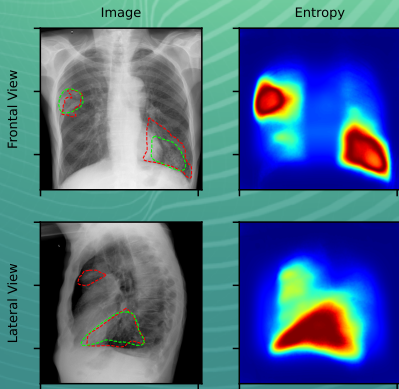
CT of lungs:

- ▶ Calculate the volume of lungs affected by covid



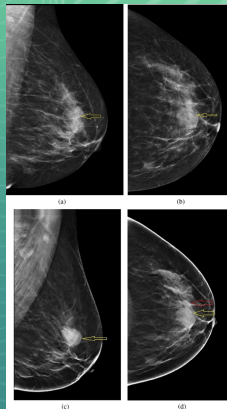
X-rays of chest and breast:

- Restoration of segmentation masks using uncertainty estimation on chest X-rays;



X-rays of chest and breast:

- ▶ Restoration of segmentation masks using uncertainty estimation on chest X-rays;
- ▶ Classification and segmentation of breast cancer.



ECG:

- ▶ Classification of QRS-complexes;

ECG:

- ▶ Classification of QRS-complexes;
- ▶ Risk of type II diabetes;

PET CT is used for localize cancer cells in a body.
How a study is conducted?

- ▶ An injection of FDG (Fludeoxyglucose);

PET CT is used for localize cancer cells in a body.

How a study is conducted?

- ▶ An injection of FDG (Fludeoxyglucose);
- ▶ PET study;

PET CT is used for localize cancer cells in a body.

How a study is conducted?

- ▶ An injection of FDG (Fludeoxyglucose);
- ▶ PET study;
- ▶ CT study;

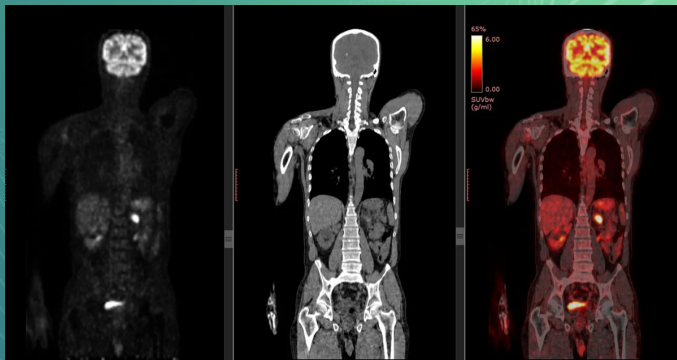


Figure: Example of PET-CT study

Possible tasks:

- ▶ Reduce the amount of FDG;

Possible tasks:

- ▶ Reduce the amount of FDG;
- ▶ Conduct the study faster.

The standard exposure time for PET study is 90 seconds. But we have PET data after 30 and 60 seconds of exposure as well.

Q: Is it possible to restore the standard exposure time images?

Possible approaches:

- ▶ Deterministic (classical image processing or deep learning approaches); – One input image -> one output image;

Possible approaches:

- ▶ Deterministic (classical image processing or deep learning approaches); – One input image -> one output image;
- ▶ Probabilistic; one input image -> multiple images, e.g. probability distribution.

Approaches:

- ▶ Gaussian filtration;

Approaches:

- ▶ Gaussian filtration;
- ▶ U-net like Transformer deep network;

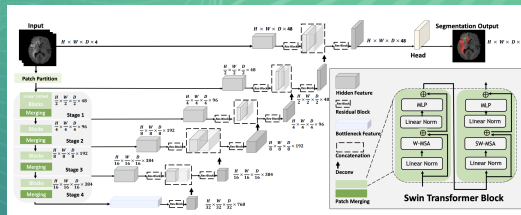


Figure: From Hatamizadeh et. al "Swin UNETR: Swin Transformers for Semantic Segmentation of Brain Tumors in MRI Images"

Approaches:

- ▶ Gaussian filtration;
- ▶ U-net like Transformer deep network;
- ▶ Diffusion model (Work in progress).

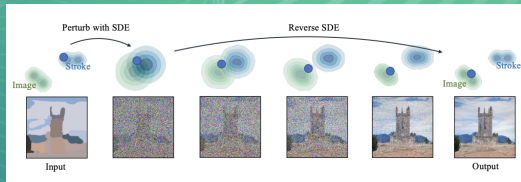
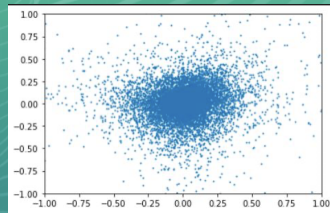
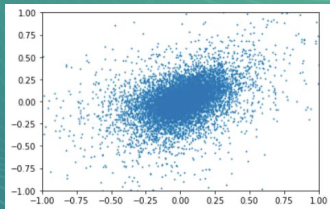
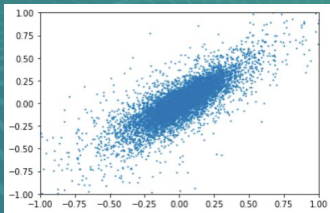


Figure: From Meng et. al "SDEdit: Guided Image Synthesis And Editing With Stochastic Differential Equations"

Pixel-wise difference between adjacent pixels (left), pixels at distance 2 (center), and 3 (right) of the difference $\Delta_o = \text{PET}_{60} - \text{PET}_{90}$:



Current results. Calculated as

$$100 * (1 - \frac{\Delta_r}{\Delta_o}),$$

where $\Delta_r = PET_r - PET_{90}$ and $\Delta_o = PET_{60} - PET_{90}$.

Method	MSE	1 - SSIM
Gaussian filtration (1 layer)	11.8 %	3.2%
Gaussian filtration (3 layer)	15.2%	-58.6%
Swin Transformer	16.6%	-128.6%
Diffusion model	Work in progress	

Current results

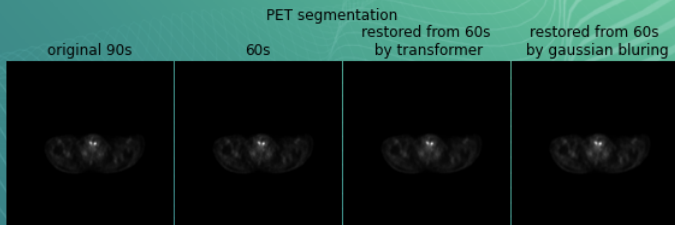
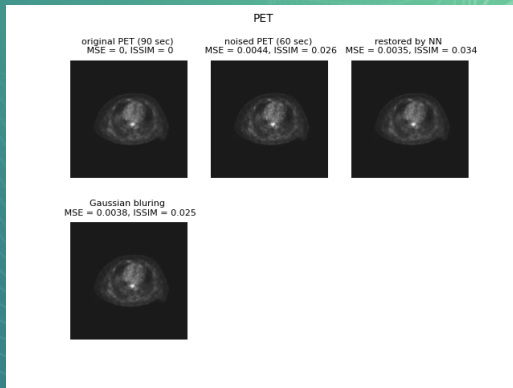


Figure: Restoration of the noised PET

Current results



Statement of the problem

- ▶ CT has become easily available and relatively cheap;

Statement of the problem

- ▶ CT has become easily available and relatively cheap;
- ▶ There are a lot of very massive CT datasets (annotated and not);

Statement of the problem

- ▶ CT has become easily available and relatively cheap;
- ▶ There are a lot of very massive CT datasets (annotated and not);
- ▶ A lot of medical centers want to predict a presence of some disease based on a limited amount of data;

Statement of the problem

- ▶ Pretrain a deep learning model that is able to generalize well for a wide range of downstream tasks.

We will consider segmentation task since:

- ▶ If training from scratch: requires a lot of data;

We will consider segmentation task since:

- ▶ If training from scratch: requires a lot of data;
- ▶ Requires networks with more parameters compared to, for instance, classification, hence more computational resources.

We will consider segmentation task since:

- ▶ If training from scratch: requires a lot of data;
- ▶ Requires networks with more parameters compared to, for instance, classification, hence more computational resources.
- ▶ Pretraining is a way to overcome these shortcomings!

Possible pretraining approaches:

- ▶ Supervised pretrain;

Possible pretraining approaches:

- ▶ Supervised pretrain;
- ▶ Unsupervised pretrain:

Possible pretraining approaches:

- ▶ Supervised pretrain;
- ▶ Unsupervised pretrain:
 1. Inpainting;

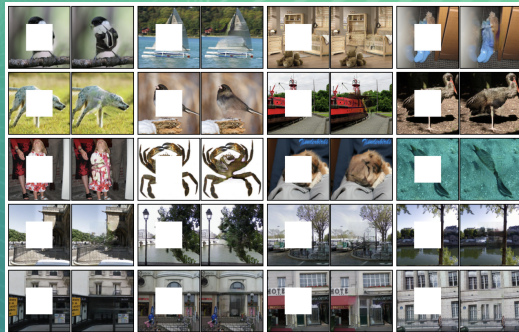


Figure: From Pathak et. al "Context Encoders: Feature Learning by Inpainting"

Possible pretraining approaches:

- ▶ Supervised pretrain;
- ▶ Unsupervised pretrain:
 1. Inpainting;
 2. Image rotation;

Possible pretraining approaches:

- ▶ Supervised pretrain;
- ▶ Unsupervised pretrain:
 1. Inpainting;
 2. Image rotation;
 3. Contrastive learning.

Contrastive learning:

- We would like to maximize mutual information between *similar* images ¹.

$$I(x, x') = \sum_{x, x'} p(x, x') \log \frac{p(x|x')}{p(x)};$$

¹See: van den Oord et al. Representation Learning with Contrastive Predictive Coding

Contrastive learning:

- ▶ We would like to maximize mutual information between *similar* images ¹.

$$I(x, x') = \sum_{x, x'} p(x, x') \log \frac{p(x|x')}{p(x)};$$

- ▶ In fact we will minimize

$$\mathcal{L} = -\mathbb{E}_X \left[\frac{f(x, x')}{\sum_{j=1}^N f(x_j, x')} \right], \text{ where } f(x, x') \propto \frac{p(x|x')}{p(x)}.$$

¹See: van den Oord et al. Representation Learning with Contrastive Predictive Coding

But deep learning benefits from scalable approaches!

SimCLR

- Sample N images;

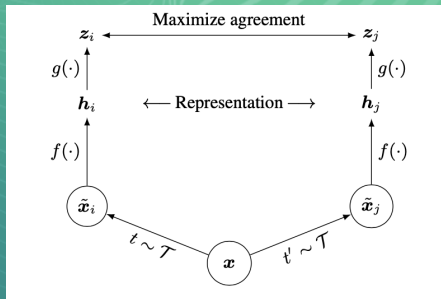


Figure: From Chen et. al "A Simple Framework for Contrastive Learning of Visual Representations"

But deep learning benefits from scalable approaches!

SimCLR

- ▶ Sample N images;
- ▶ Augment each image;

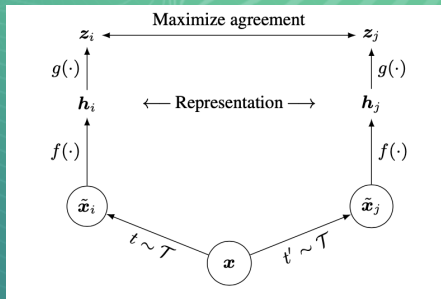


Figure: From Chen et. al "A Simple Framework for Contrastive Learning of Visual Representations"

But deep learning benefits from scalable approaches!

SimCLR

- ▶ Sample N images;
- ▶ Augment each image;
- ▶ Obtain representations h using neural network $f(\cdot)$;

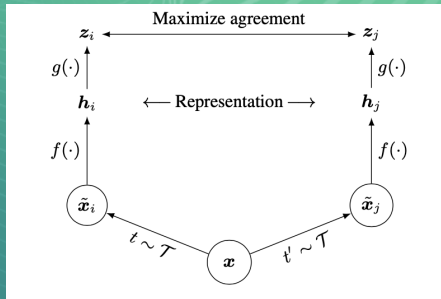


Figure: From Chen et. al "A Simple Framework for Contrastive Learning of Visual Representations"

But deep learning benefits from scalable approaches!

SimCLR

- ▶ Sample N images;
- ▶ Augment each image;
- ▶ Obtain representations h using neural network $f(\cdot)$;
- ▶ Obtain embeddings z using non-linear head $g(\cdot)$

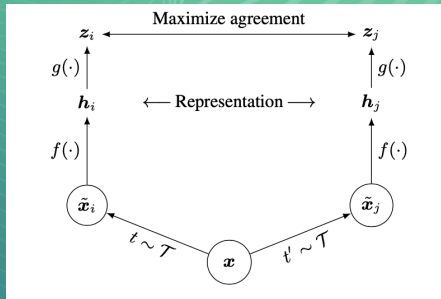


Figure: From Chen et. al "A Simple Framework for Contrastive Learning of Visual Representations"

But deep learning benefits from scalable approaches!

SimCLR

- ▶ Sample N images;
- ▶ Augment each image;
- ▶ Obtain representations h using neural network $f(\cdot)$;
- ▶ Obtain embeddings z using non-linear head $g(\cdot)$
- ▶ Maximize agreement (MI) using contrastive loss;

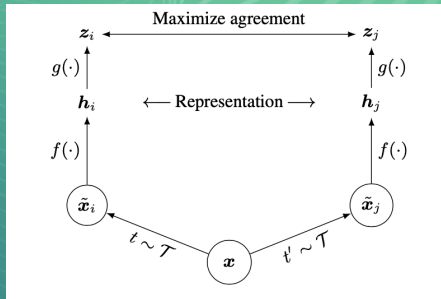


Figure: From Chen et. al "A Simple Framework for Contrastive Learning of Visual Representations"

Contrastive loss:

- For an augmented pair of images:

$$\mathcal{L}(x_i, x_j) = -\log \frac{\exp(z_i \cdot z_j / \tau)}{\sum_{k=1}^{2N} \mathbb{1}_{k \neq j} [\exp(z_i \cdot z_k / \tau)]};$$

Contrastive loss:

- For an augmented pair of images:

$$\mathcal{L}(x_i, x_j) = -\log \frac{\exp(z_i \cdot z_j / \tau)}{\sum_{k=1}^{2N} \mathbb{1}_{k \neq j} [\exp(z_i \cdot z_k / \tau)]};$$

- Total loss:

$$\mathcal{L} = \frac{1}{2N} \sum_{i=1}^N (\mathcal{L}(x_{2k-1}, x_{2k}) + \mathcal{L}(x_{2k}, x_{2k-1})).$$

Why this (and a lot of similar approaches) are so attractable?

- ▶ Contrastive learning (in the original formulation) does not require ground truth labels;

Why this (and a lot of similar approaches) are so attractable?

- ▶ Contrastive learning (in the original formulation) does not require ground truth labels;
- ▶ Linear evaluation protocol shows results on par with supervised methods.

Some observations:

- ▶ SimCLR and similar methods are based on image-level comparisons;

Some observations:

- ▶ SimCLR and similar methods are based on image-level comparisons;
- ▶ Which might be sub-optimal for segmentation tasks due to the lack of spatial sensitivity.

Requirements for pretext task:

- ▶ It should be **spatial sensitive**, i.e. discriminate spatially closed pixels for accurate predictions in boundary regions;

Requirements for pretext task:

- ▶ It should be **spatial sensitive**, i.e. discriminate spatially closed pixels for accurate predictions in boundary regions;
- ▶ It should be **spatial smooth**. Spatial smoothness encourage clone pixels to belong to the same class.

Propagate yourself²:

- ▶ Choose a convolutional neural network $f(\cdot)$;

²See Xie et al. Propagate Yourself: Exploring Pixel-Level Consistency for Unsupervised Visual Representation Learning

Propagate yourself²:

- ▶ Choose a convolutional neural network $f(\cdot)$;
- ▶ Process an input image $\mathbb{I} \in \mathbb{R}^{1 \times H \times W}$ using $f(\cdot)$ to obtain a representation from some convolutional layer $\tilde{\mathbb{I}} \in \mathbb{R}^{C \times H' \times W'}$;

²See Xie et al. Propagate Yourself: Exploring Pixel-Level Consistency for Unsupervised Visual Representation Learning

Propagate yourself²:

- ▶ Choose a convolutional neural network $f(\cdot)$;
- ▶ Process an input image $\mathbb{I} \in \mathbb{R}^{1 \times H \times W}$ using $f(\cdot)$ to obtain a representation from some convolutional layer $\tilde{\mathbb{I}} \in \mathbb{R}^{C \times H' \times W'}$;
- ▶ Consider each C -dimensional vector of the $\tilde{\mathbb{I}}$ as a representation of a pixel in some pixel-space.

²See Xie et al. Propagate Yourself: Exploring Pixel-Level Consistency for Unsupervised Visual Representation Learning

The pipeline:

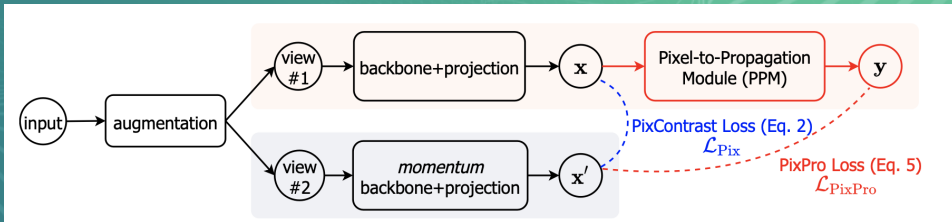


Figure: See Xie et al. Propagate Yourself: Exploring Pixel-Level Consistency for Unsupervised Visual Representation Learning

- Pixel-level contrastive loss:

$$\mathcal{L}_{\mathcal{P}}(i) = -\log \frac{\sum_{j \in \Omega_p^i} e^{\cos(x_i, x'_j)/\tau}}{\sum_{j \in \Omega_p^i} e^{\cos(x_i, x'_j)/\tau} + \sum_{k \in \Omega_n^i} e^{\cos(x_i, x'_k)/\tau}};$$

- ▶ Pixel-level contrastive loss:

$$\mathcal{L}_{\mathcal{P}}(i) = -\log \frac{\sum_{j \in \Omega_p^i} e^{\cos(x_i, x'_j)/\tau}}{\sum_{j \in \Omega_p^i} e^{\cos(x_i, x'_j)/\tau} + \sum_{k \in \Omega_n^i} e^{\cos(x_i, x'_k)/\tau}};$$

- ▶ Where x_i, x'_j – pixels from two augmented versions of an image x ;

- ▶ Pixel-level contrastive loss:

$$\mathcal{L}_{\mathcal{P}}(i) = -\log \frac{\sum_{j \in \Omega_p^i} e^{\cos(x_i, x'_j)/\tau}}{\sum_{j \in \Omega_p^i} e^{\cos(x_i, x'_j)/\tau} + \sum_{k \in \Omega_n^i} e^{\cos(x_i, x'_k)/\tau}};$$

- ▶ Where x_i, x'_i – pixels from two augmented versions of an image x ;
- ▶ Ω_p^i and Ω_n^i – pixels inside and outside some vicinity of the current pixel i respectively.

► The total loss:

$$\mathcal{L} = \mathcal{L}_{\text{inst}} + \frac{1}{H' \times W'} \sum_i \mathcal{L}_{\mathcal{P}}(i)$$

We will compare the following methods:

- Training from a random initialization;

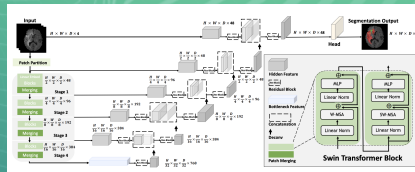


Figure: From Hatamizadeh et. al "Swin UNETR: Swin Transformers for Semantic Segmentation of Brain Tumors in MRI Images"

We will compare the following methods:

- ▶ Training from a random initialization;
- ▶ Fine-tuning a model pretrained on (Nvidia)

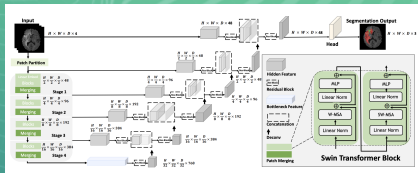


Figure: From Hatamizadeh et. al "Swin UNETR: Swin Transformers for Semantic Segmentation of Brain Tumors in MRI Images"

We will compare the following methods:

- ▶ Training from a random initialization;
- ▶ Fine-tuning a model pretrained on (Nvidia)
 - ▶ Inpainting;

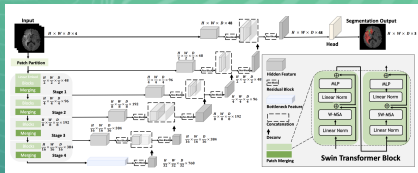


Figure: From Hatamizadeh et. al "Swin UNETR: Swin Transformers for Semantic Segmentation of Brain Tumors in MRI Images"

We will compare the following methods:

- ▶ Training from a random initialization;
- ▶ Fine-tuning a model pretrained on (Nvidia)
 - ▶ Inpainting;
 - ▶ Rotation prediction;

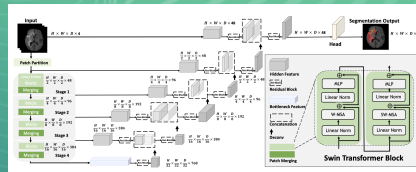


Figure: From Hatamizadeh et. al "Swin UNETR: Swin Transformers for Semantic Segmentation of Brain Tumors in MRI Images"

We will compare the following methods:

- ▶ Training from a random initialization;
- ▶ Fine-tuning a model pretrained on (Nvidia)
 - ▶ Inpainting;
 - ▶ Rotation prediction;
 - ▶ Instance-level CL.

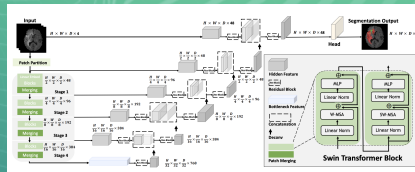


Figure: From Hatamizadeh et. al "Swin UNETR: Swin Transformers for Semantic Segmentation of Brain Tumors in MRI Images"

We will compare the following methods:

- ▶ Training from a random initialization;
- ▶ Fine-tuning a model pretrained on (Nvidia)
 - ▶ Inpainting;
 - ▶ Rotation prediction;
 - ▶ Instance-level CL.
- ▶ Fine-tuning a model pretrained on (Sber)

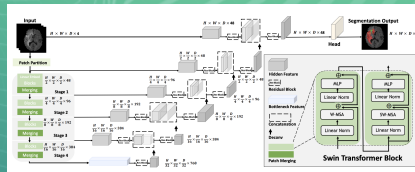


Figure: From Hatamizadeh et. al "Swin UNETR: Swin Transformers for Semantic Segmentation of Brain Tumors in MRI Images"

We will compare the following methods:

- ▶ Training from a random initialization;
- ▶ Fine-tuning a model pretrained on (Nvidia)
 - ▶ Inpainting;
 - ▶ Rotation prediction;
 - ▶ Instance-level CL.
- ▶ Fine-tuning a model pretrained on (Sber)
 - ▶ Instance-level CL;

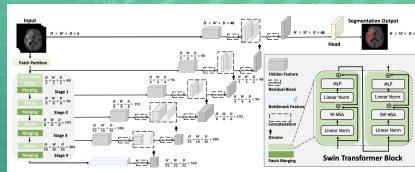


Figure: From Hatamizadeh et. al "Swin UNETR: Swin Transformers for Semantic Segmentation of Brain Tumors in MRI Images"

We will compare the following methods:

- ▶ Training from a random initialization;
- ▶ Fine-tuning a model pretrained on (Nvidia)
 - ▶ Inpainting;
 - ▶ Rotation prediction;
 - ▶ Instance-level CL.
- ▶ Fine-tuning a model pretrained on (Sber)
 - ▶ Instance-level CL;
 - ▶ Pixel-level CL.

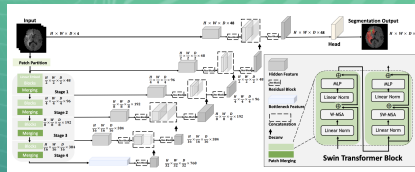


Figure: From Hatamizadeh et. al "Swin UNETR: Swin Transformers for Semantic Segmentation of Brain Tumors in MRI Images"

- ▶ Pretraining data:

- ▶ Pretraining data:
 - ▶ Nvidia: LUNA16, TCIA Covid19, LiDC (Total: 2124 CTs);

- ▶ Pretraining data:
 - ▶ Nvidia: LUNA16, TCIA Covid19, LiDC (Total: 2124 CTs);
 - ▶ Sber: LUNA16 (Total 888 CTs).

Evaluation protocol:

- ▶ Divide the training data randomly into two halves, one – for training, another – for validation;

Evaluation protocol:

- ▶ Divide the training data randomly into two halves, one – for training, another – for validation;
- ▶ Use for training 20%, 50% or 100% of training data;

Evaluation protocol:

- ▶ Divide the training data randomly into two halves, one – for training, another – for validation;
- ▶ Use for training 20%, 50% or 100% of training data;
- ▶ Choose the best on validation;

Evaluation protocol:

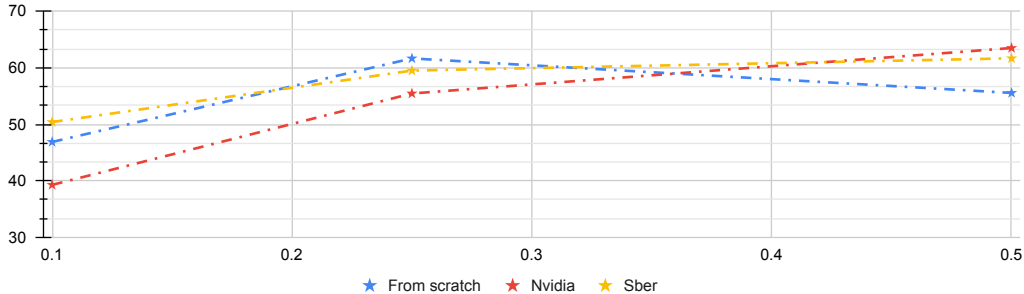
- ▶ Divide the training data randomly into two halves, one – for training, another – for validation;
- ▶ Use for training 20%, 50% or 100% of training data;
- ▶ Choose the best on validation;
- ▶ Measure the model on test data;

Evaluation protocol:

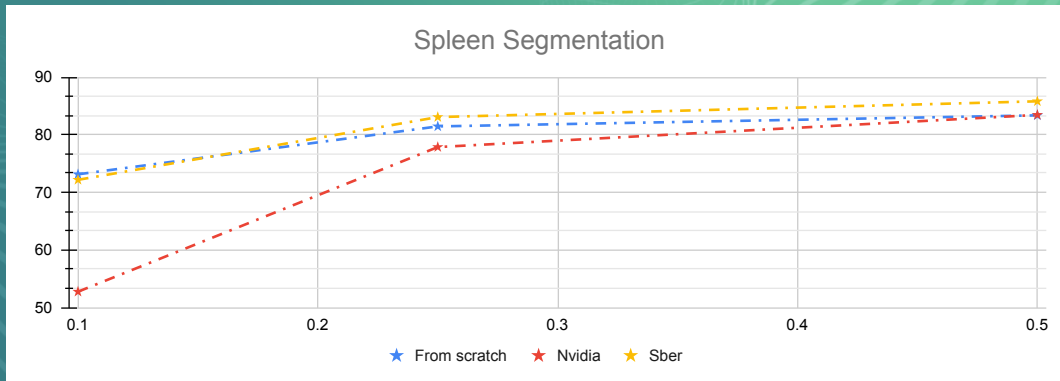
- ▶ Divide the training data randomly into two halves, one – for training, another – for validation;
- ▶ Use for training 20%, 50% or 100% of training data;
- ▶ Choose the best on validation;
- ▶ Measure the model on test data;
- ▶ Average between six runs.

Results: Task06_Lung from Medical Decathlon

Lung Cancer Segmentation



Results Task09_Spleen from Medical Decathlon:



Thank you for your softmax $\left(\frac{QK^T}{\sqrt{d_k}}\right)$ V!