# Is the problem of medical image segmentation a thing of the past ?

by

## Olivier Bernard

Professor – University of Lyon (INSA), France

July 10, 2024

CREATIS; CNRS (UMR 5220); INSERM (U1294); INSA Lyon; Université de Lyon, France

# AI methods in cardiac image analysis

## Acquisition

Ultrafast cardiac imaging

Convolutional NN
Realistic simulations

## Image quantification

Segmentation
Tissue motion / blood flow
Uncertainty modeling
Domain adaptation

Convolutional NN
Variational Auto-Encoders
Physics informed NN
Diffusion networks

## Population representation

Multi-modal fusion
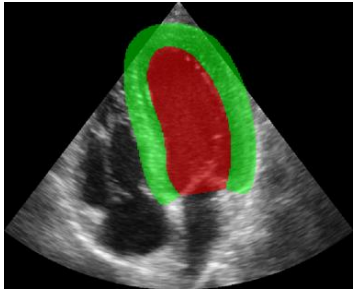Heterogenous data integration

Transformers

Etiology classification
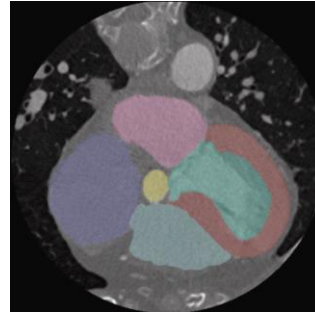Hypertension characterization

Robust estimation of existing / new biomarkers

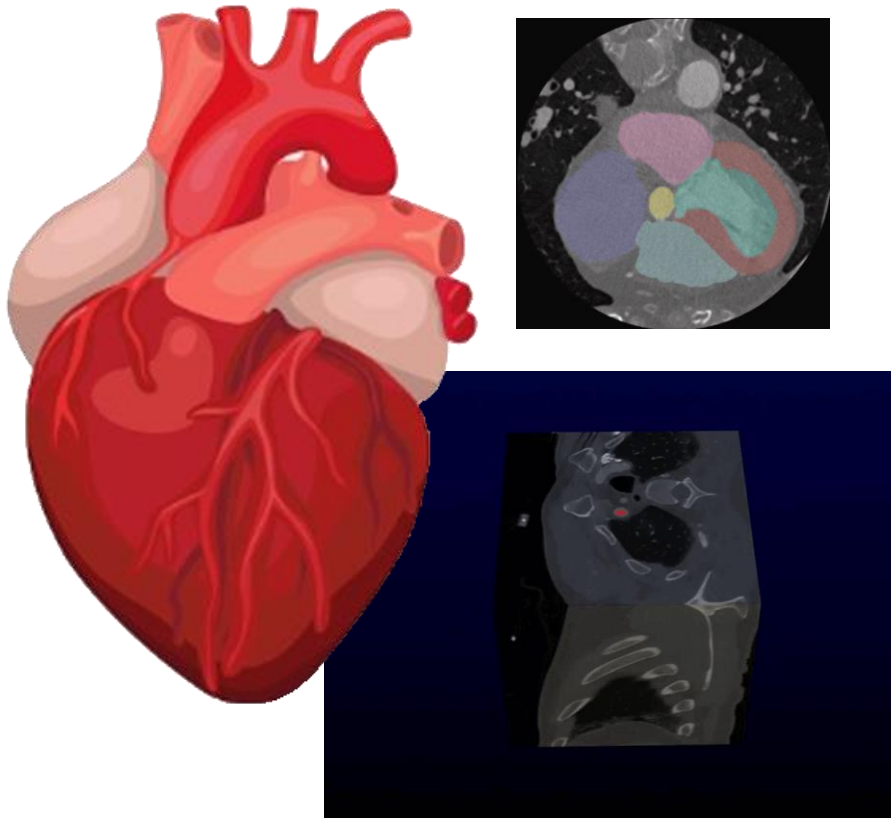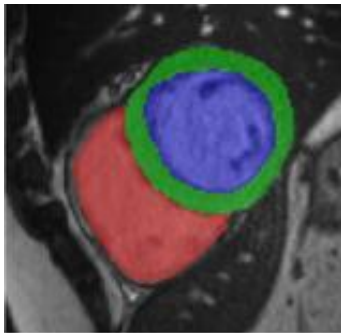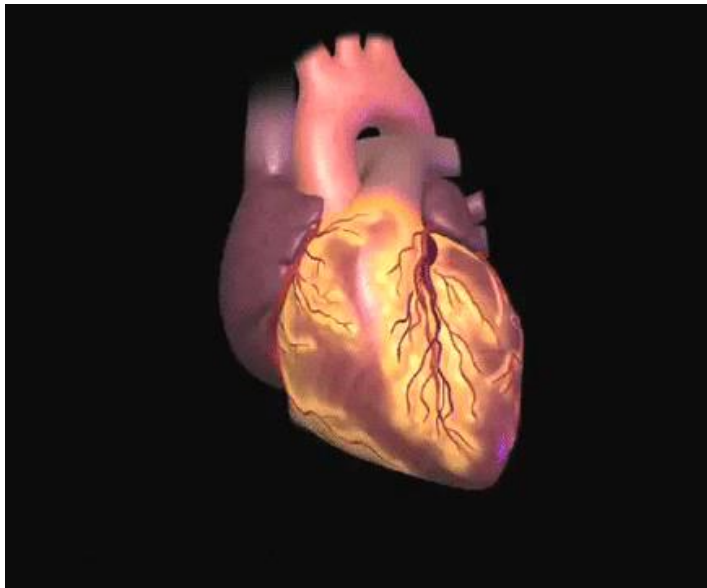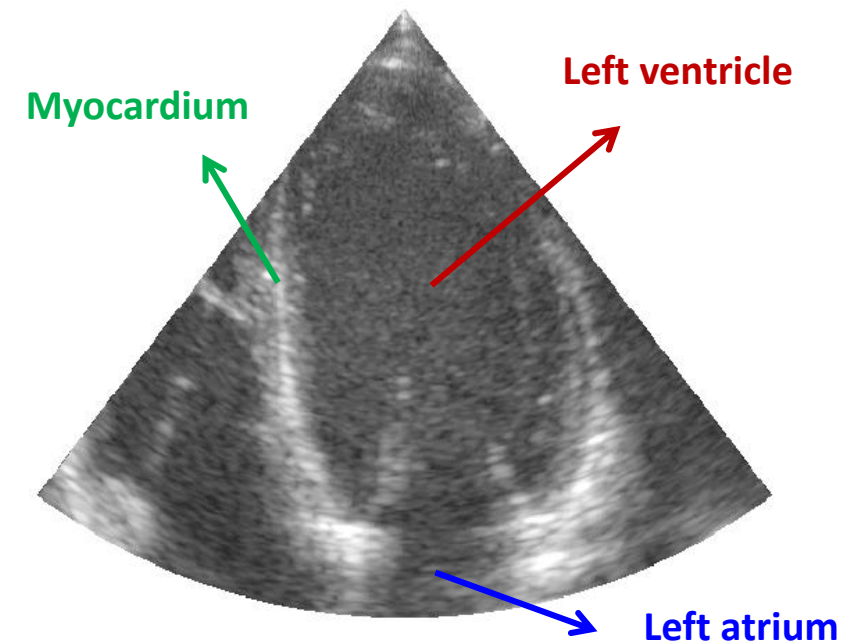Quantification of clinical indices to diagnose cardiac pathologies

Echocardiography

CT

MRI
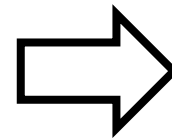
- ✓ High annotation costs
- ✓ Inter/intra expert variability
- ✓ Acquisition variabilities
- ✓ Acquisition artifacts

*From the inHEART company website*

Quantification of clinical indices to diagnose cardiac pathologies

▶ Anatomical imaging



Source: GE Healthcare web site

Conventional ultrasound acquisition in clinical routine

**Myocardium**

**Left ventricle**

**Left atrium**

4

Quantification of clinical indices to diagnose cardiac pathologies



Automatic delineation of anatomical structures →

- Scalar descriptors
- Time-series descriptors

**Scalar descriptors**

- Myocardial mass
- Left ventricle ejection fraction

**Time-series descriptors**

- Left ventricle area
- Global longitudinal strain

**Myocardial strain**



Time

Challenges

► How to make the measurements extracted from medical images automatic, reliable and accurate ?

► How to make these measurements reproducible at different centers, in different countries, whatever the expert ?

# Deep learning families

---

# Convolutional Neural Network

## Convolutional layer

✓ Create relevant information called *feature map* (convolution + non linear function)

✓ Parameters that are learned during training



*Input*

*feature map*

256

256

$f(\cdot)$

256

256

3

1

Filter of size $3 \times 3$

$\# param = 3 \times 3 \times 3 + 1$
$= 28$

## Convolutional layer

✓ Create relevant information called *feature map* (convolution + non linear function)

✓ Parameters that are learned during training



Filter of size $3 \times 3$

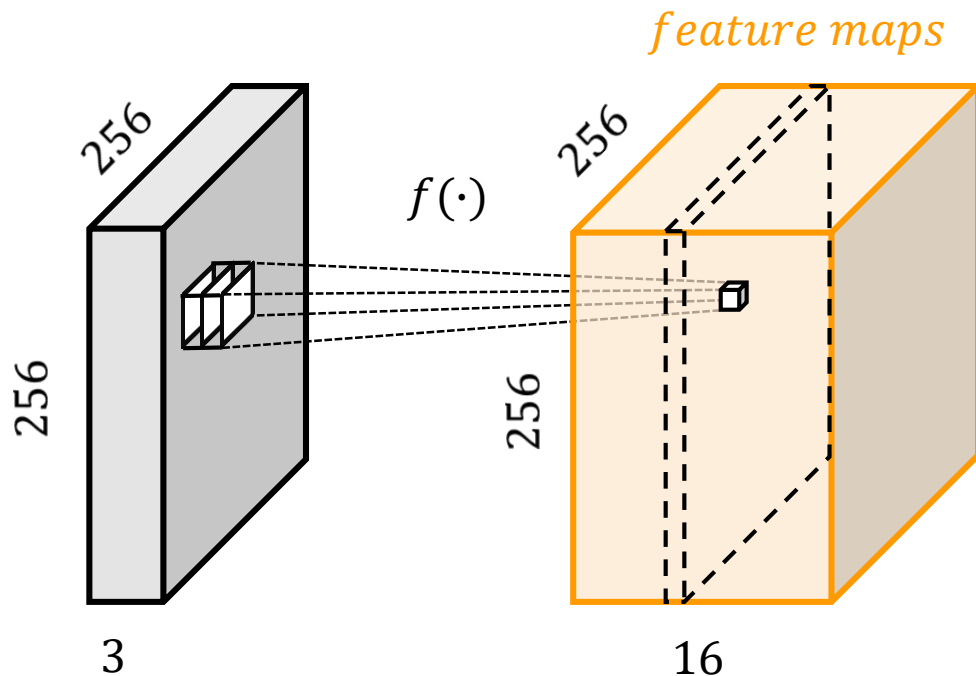$$\# \, param = 16 \times (3 \times 3 \times 3 + 1) = 448$$

## Convolutional layer

✓ Create relevant information called *feature map* (convolution + non linear function)
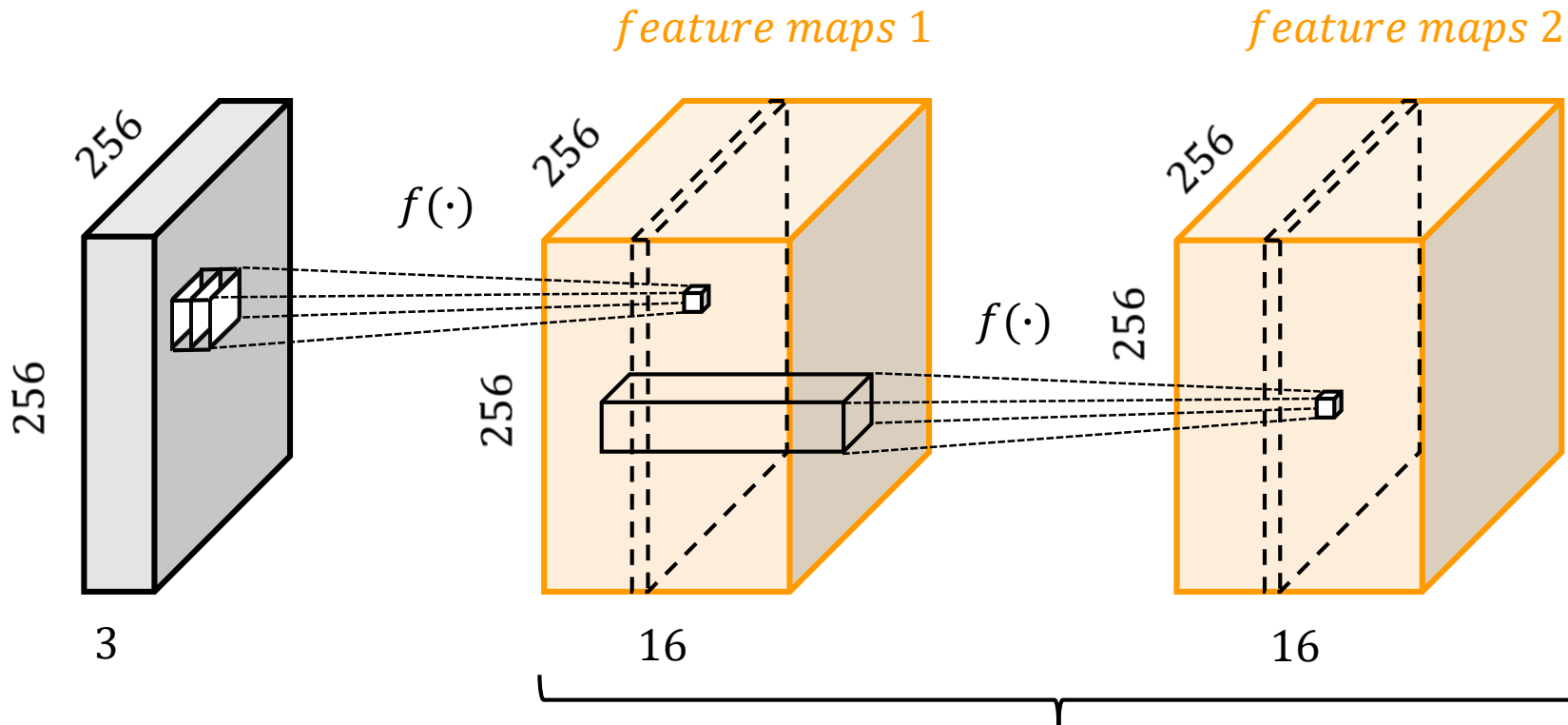
✓ Parameters that are learned during training



*feature maps* 1

*feature maps* 2

256

256

256

256

256

256

$f(\cdot)$

$f(\cdot)$

3

16

16

Filter of size $3 \times 3$

$\# \, param$
$= 16 \times (3 \times 3 \times 16 + 1)$
$= 2{,}320$

## Pooling operation

✓ Concentrate information into lower dimensional space

✓ Applied individually to each feature map



| 135 | 212 | 189 | 56 |
|-----|-----|-----|-----|
| 164 | 201 | 204 | 145 |
| 30 | 126 | 189 | 156 |
| 36 | 45 | 38 | 12 |

Input
feature map

Max pooling
operation

| 212 | 204 |
|-----|-----|
| 126 | 189 |

Output
feature map

## Pooling operation

✓ Concentrate information into lower dimensional space

✓ Applied individually to each feature map



256

256

16

128

128

16

Max pooling operation

No parameter to train

# Image encoding

✓ Learning to encode relevant information

✓ Projection to a lower dimensional space



Image encoding

Image encoding

## Deconvolutional layer

✓ Propagate relevant information to the input dimension space
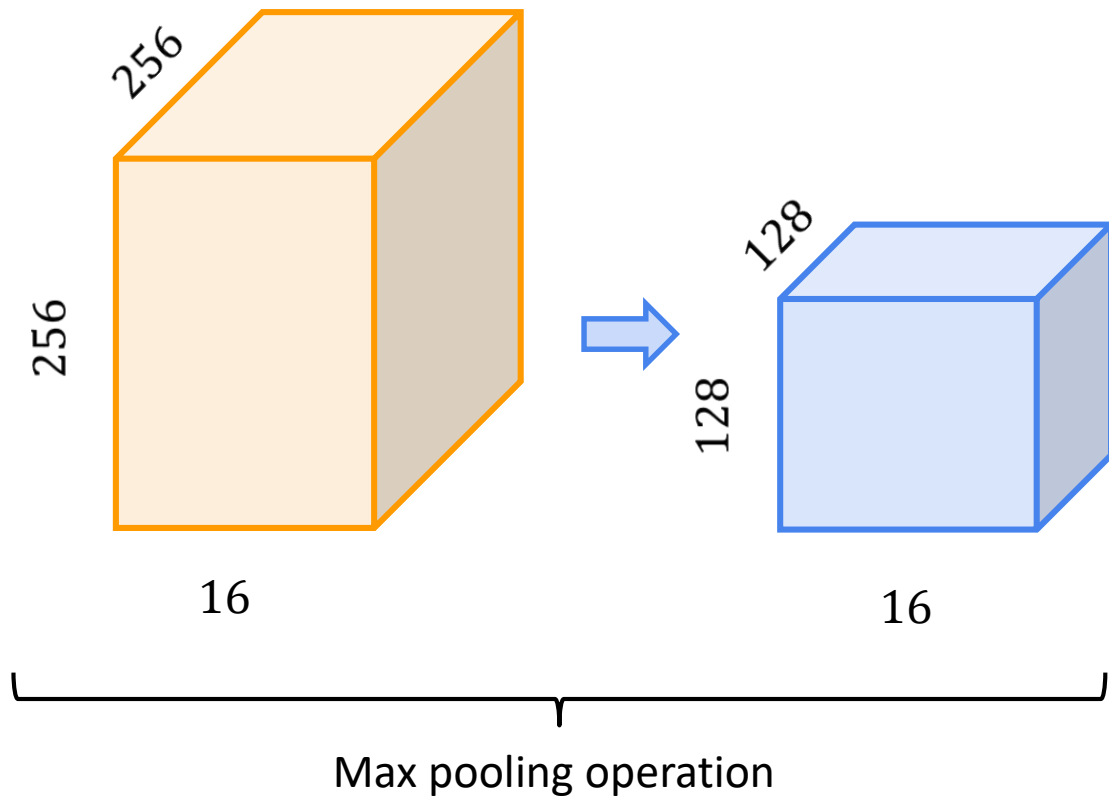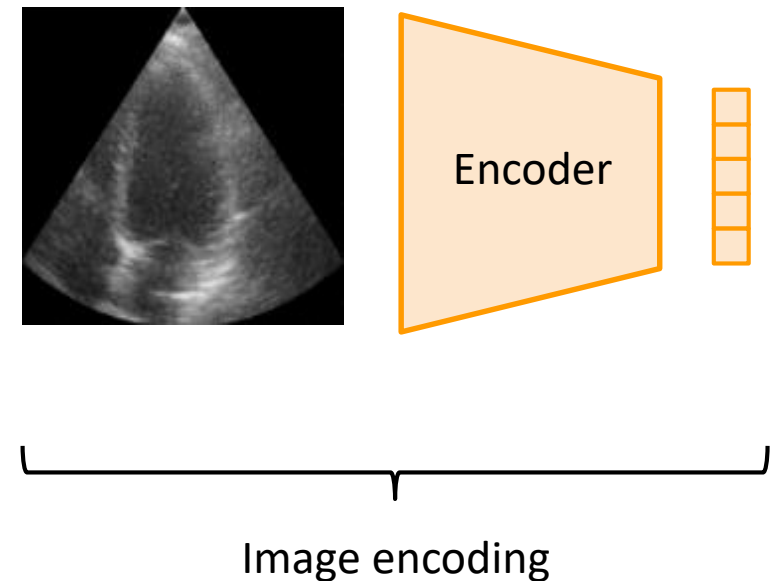
✓ Parameters that are learned during training



*feature maps* 2

*feature maps* 1

16

256

16

32

32

128

Filter of size $3 \times 3$

$\# \, param$
$= 128 \times (3 \times 3 \times 256 + 1)$
$= 295,040$

14

## Encoder-decoder architectures



Encoder       Decoder

# U-Net architecture

✓ Between 3 M to 40 M of parameters to train



Input image

Skip connection

Segmentation output

| | |
|---|---|
| Conv 3x3 + ReLU | |
| Conv 1x1 | |
| Pooling | |
| Upsampling | |

# Deep learning families

---

# Transformers

## Tokenization procedure

✓ Efficient representation of an image

$D = 768$

$$\# \, param = 16 \times 16 \times 3 \times 768 = 589{,}824$$

16

16

512

512

$x_i \in \mathbb{R}^{1 \times (16^2 \cdot 3)}$

Vector representation

**Linear projection**

Simple matrix multiplication with

$\mathbf{U}_e$

of dimension

$\mathbb{R}^{(16^2 \cdot 3) \times D}$

*TO LEARN*

$e_i \in \mathbb{R}^{1 \times D}$

Token representation

Tokenization procedure

# Tokenization procedure

✓ Representation of an image into a lower dimensional space



Tokenization procedure applied on each $16 \times 16$ patch

$D = 768$

Feature maps

One token

# Transformer blocs / layers

✓ Create relevant information (attention + non linear function)

✓ Parameters that are learned during training

*Input feature maps*

*Output feature maps*



$D = 768$

$\times 8$

Transformer bloc / layer

Multi-head attention module
+ non linearity

$D = 768$

# Self-attention module

$$D = 768, D_h = 64$$

$$\# param = 3 \times 768 \times 64 = 147{,}456$$



**Self-attention module**

**Linear projections**
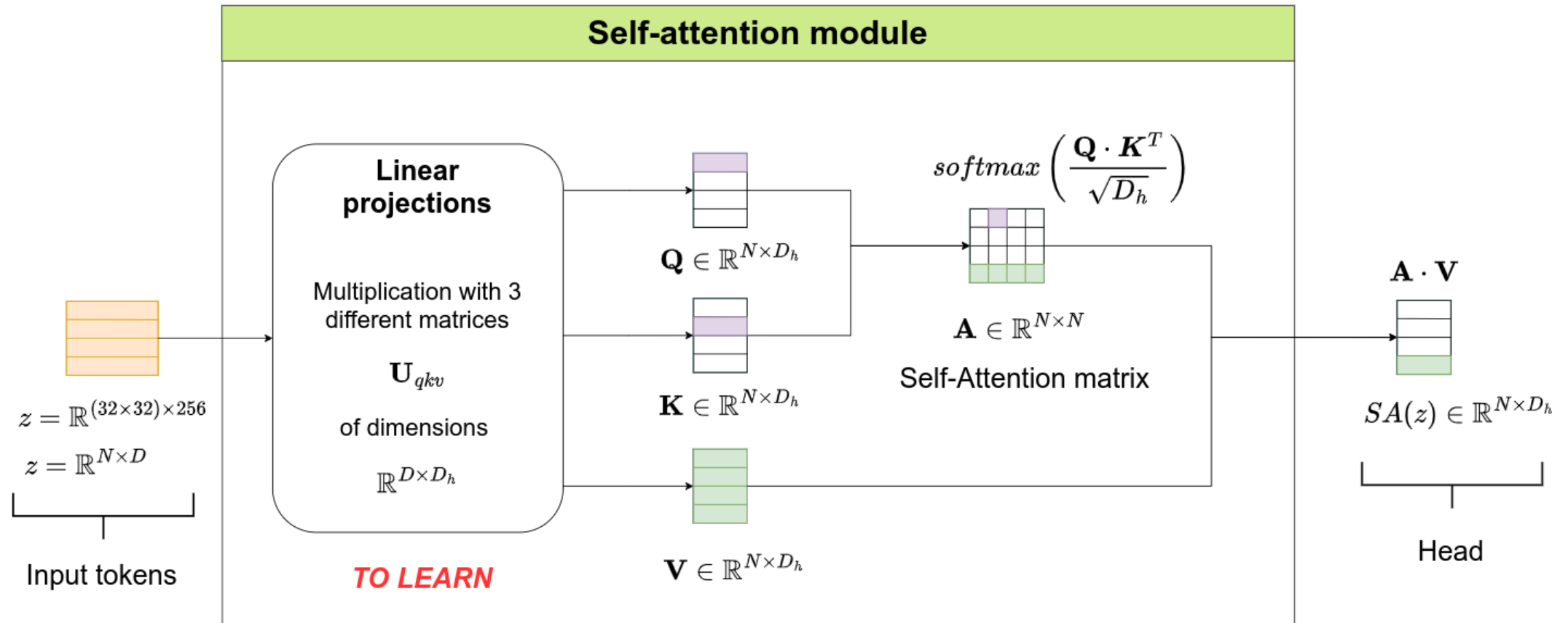
Multiplication with 3 different matrices

$$\mathbf{U}_{qkv}$$

of dimensions

$$\mathbb{R}^{D \times D_h}$$

**TO LEARN**

$$z = \mathbb{R}^{(32 \times 32) \times 256}$$

$$z = \mathbb{R}^{N \times D}$$

Input tokens

$$\mathbf{Q} \in \mathbb{R}^{N \times D_h}$$

$$\mathbf{K} \in \mathbb{R}^{N \times D_h}$$

$$\mathbf{V} \in \mathbb{R}^{N \times D_h}$$

$$softmax\left(\frac{\mathbf{Q} \cdot \boldsymbol{K}^T}{\sqrt{D_h}}\right)$$

$$\mathbf{A} \in \mathbb{R}^{N \times N}$$

Self-Attention matrix

$$\mathbf{A} \cdot \mathbf{V}$$

$$SA(z) \in \mathbb{R}^{N \times D_h}$$

Head

21

# Multi-head attention module

$$D = 768, D_h = 64, k = 12$$

$$\# \, param = 12 \times 3 \times 768 \times 64 + 768 * 768 = 2{,}359{,}296$$



**Multi-head attention module**

Input tokens

$z = \mathbb{R}^{(32 \times 32) \times 256}$

$z = \mathbb{R}^{N \times D}$

**Self-Attention module**

*TO LEARN*

Head 1

Head 2

Head k-1

Head k

CONCATENATION

$\mathbb{R}^{N \times D_h}$

With $\quad D_h = D/k$

$\mathbb{R}^{N \times (kD_h)}$

**Linear projection**

Simple matrix multiplication of dimensions

$\mathbb{R}^{(kD_h) \times D}$

*TO LEARN*

$z = \mathbb{R}^{N \times D}$

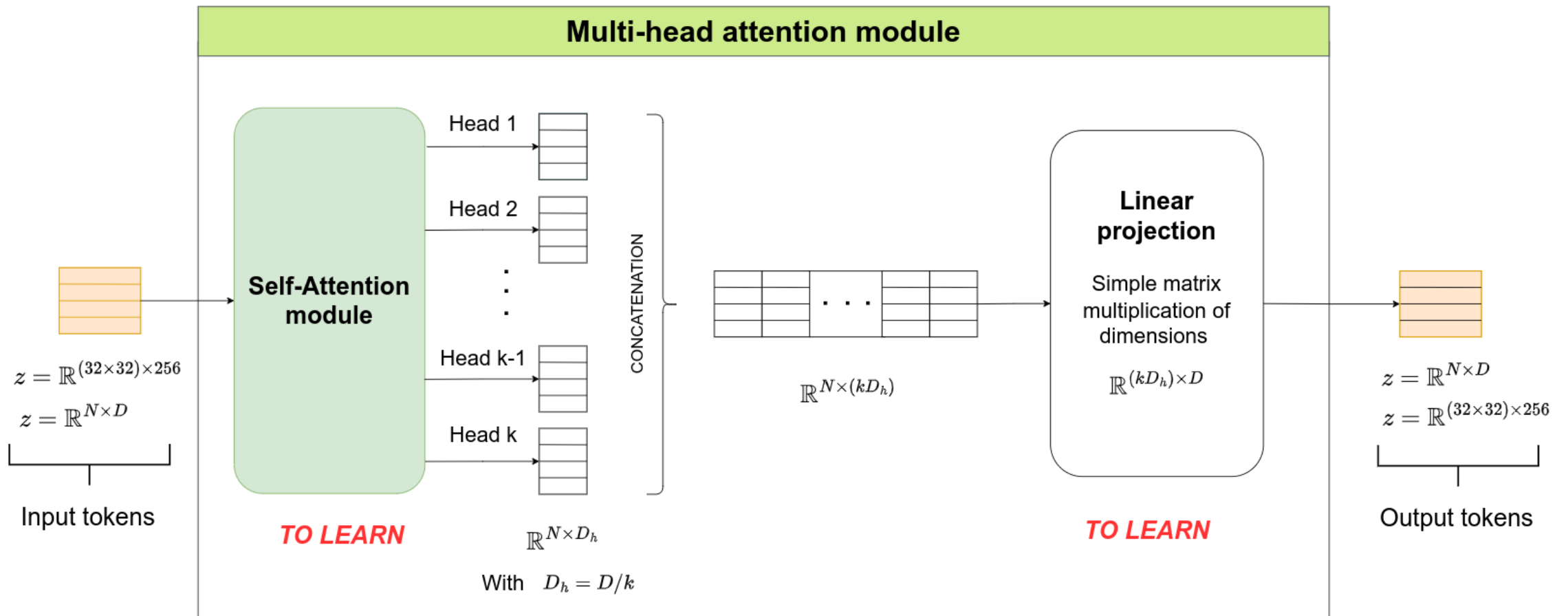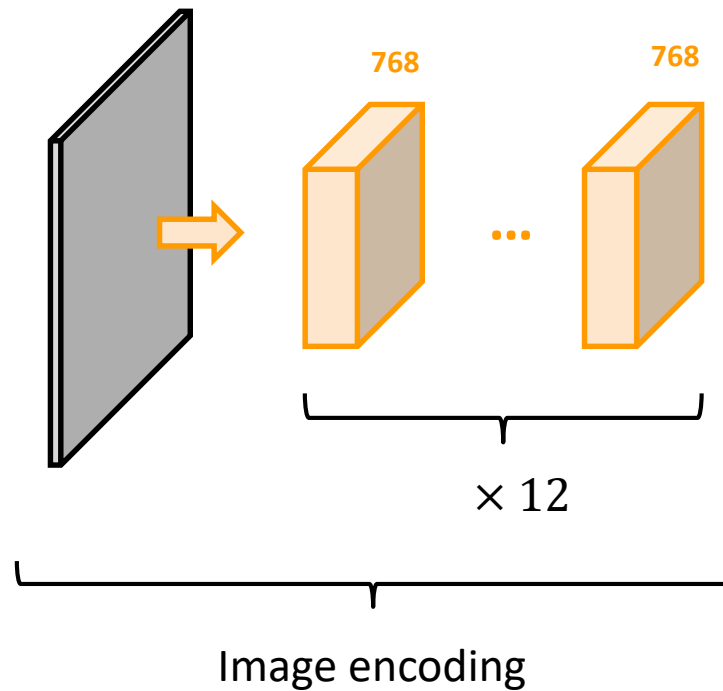$z = \mathbb{R}^{(32 \times 32) \times 256}$

Output tokens

22

## Image encoding

✓ Learning to encode relevant information

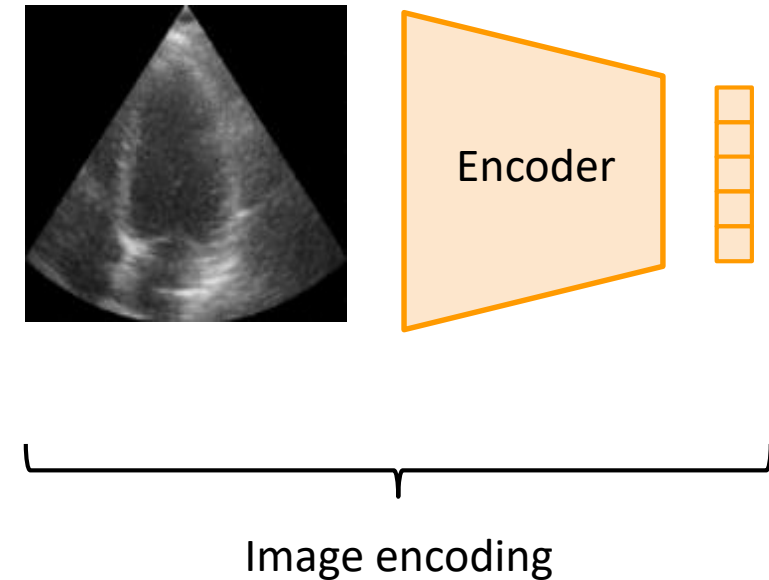✓ Projection to a lower dimensional space

$$D = 768, D_h = 64, k = 12, N_{blocs} = 12$$

$$\# param$$
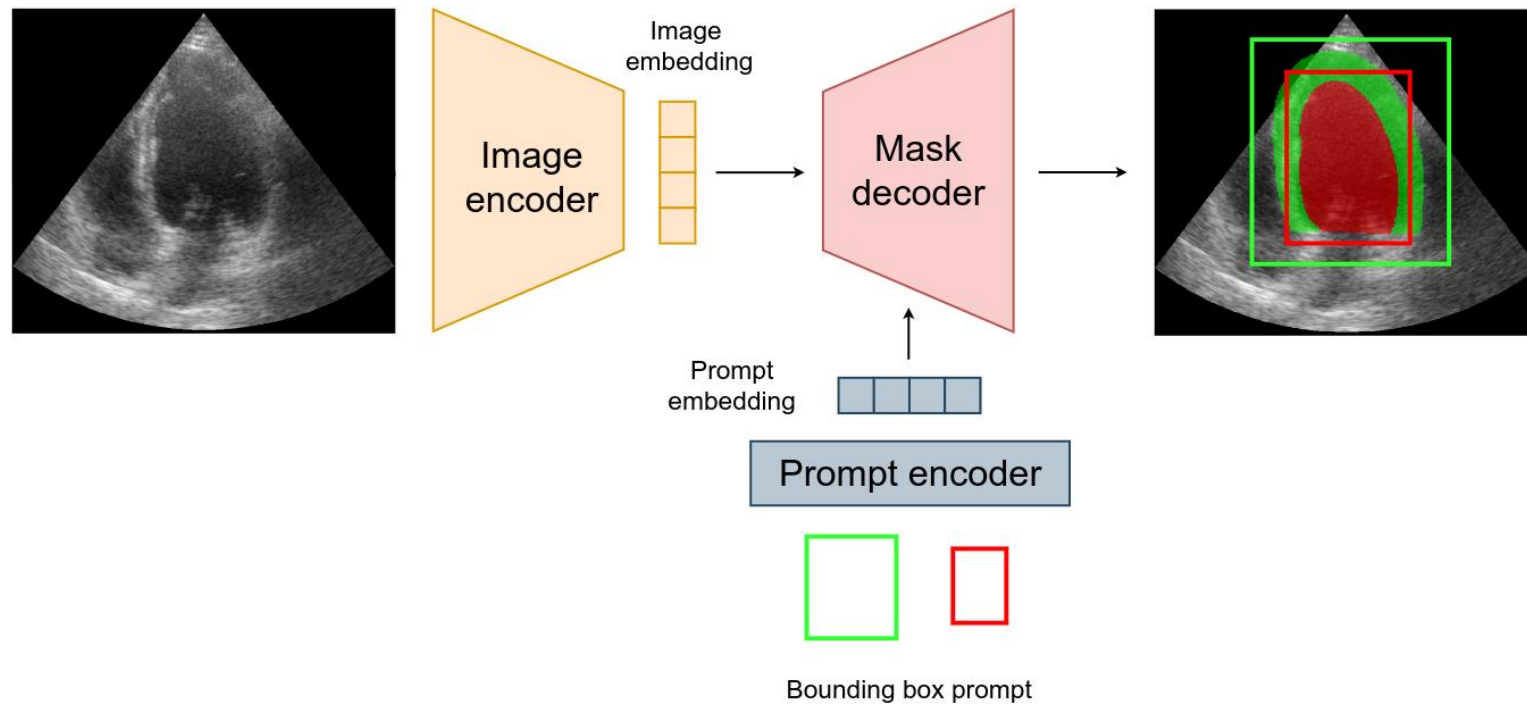$$= 12 \times (12 \times 3 \times 768 \times 64 + 768 * 768)$$
$$+ 12 \times (2 \times 768 \times 3072)$$
$$= 84{,}934{,}656$$



768        768

$\times 12$

Image encoding



Encoder

Image encoding

23

# Foundation models
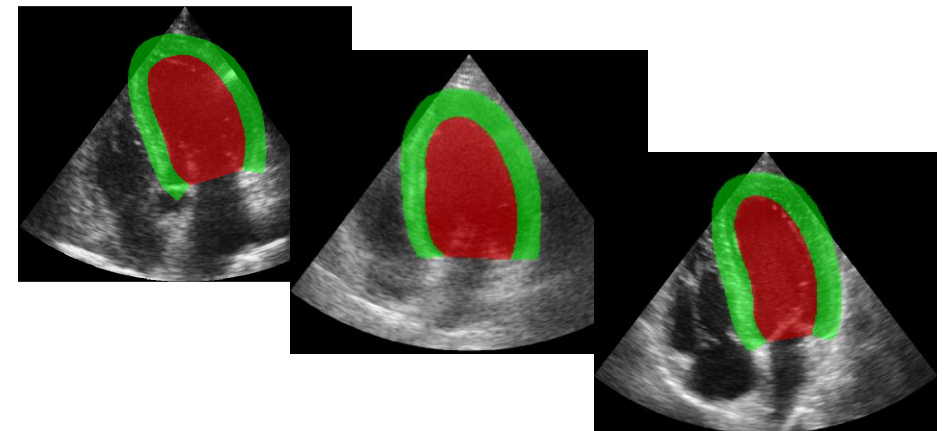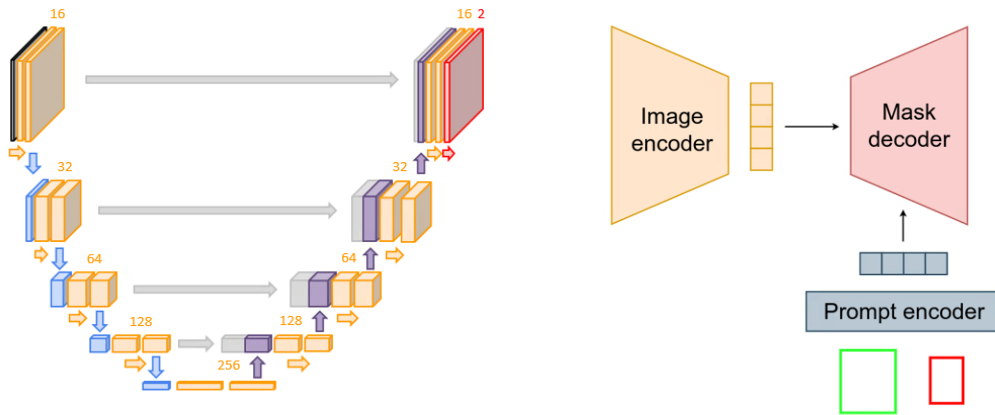
✓ 91 M of parameters to train

# Segmentation of echocardiographic images

[Leclerc et al., IEEE TMI 2019]

## The two key ingredients

✓ Deep learning solution with the proper complexity
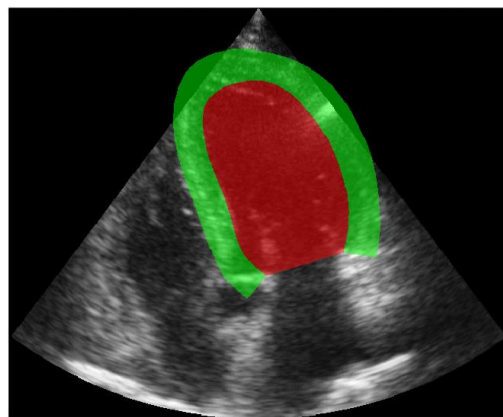
✓ Database with good quality annotations

# Echocardiographic datasets

| | | | 2D Public Echocardiographic Datasets | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Ground truth | | | | Views | | Characteristics | |
| Name | Year | Nb. Subjects | $LV_{endo}$ | $LV_{epi}$ | LA | Full cardiac cycle | A2C | A4C | Multi-Center | Multi-Vendor |
| CAMUS | 2019 | 500 | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ |
| EchoNet | 2019 | 10,036 | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | – |
| HMC-QU | 2021 | 292 | ✓ | ✓ | ✗ | ✗ | ✓ | ✓ | ✗ | ✓ |
| TED | 2022 | 98 | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ |

# Echocardiographic datasets

| | | | Ground truth | | | | Views | | Characteristics | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Name** | **Year** | **Nb. Subjects** | $LV_{endo}$ | $LV_{epi}$ | *LA* | *Full cardiac cycle* | *A2C* | *A4C* | *Multi-Center* | *Multi-Vendor* |
| **CAMUS** | 2019 | 500 | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ |
| **EchoNet** | 2019 | 10,036 | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | - |
| **HMC-QU** | 2021 | 292 | ✓ | ✓ | ✗ | ✗ | ✓ | ✓ | ✗ | ✓ |
| **TED** | 2022 | 98 | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ |

2D Public Echocardiographic Datasets

## CAMUS
- ✓ Center 1
- ✓ Annotator 1
- ✓ Vendor 1
- ✓ 500 patients
- ✓ Image annotations



## TED
- ✓ Center 1
- ✓ Annotator 1
- ✓ Vendor 1
- ✓ 98 patients
- ✓ Sequence annotations

✓ **Geometric accuracy**

(CS:CAMUS)

| Methods | Dice | | Hausdorff (mm) | |
|---|---|---|---|---|
| | ED | ES | ED | ES |
| Intra-obs. | .945 | .930 | 4.6 | 4.5 |
| 2D nnU-Net | **.952** | **.935** | **4.3** | **4.2** |
| CLAS | .947 | .929 | 4.6 | 4.6 |
| GUDU | .946 | .929 | 4.7 | 4.7 |

✓ **Clinical accuracy**

(CS:CAMUS)

| Methods | EF | | Volume ED | | Volume ES | |
|---|---|---|---|---|---|---|
| | Corr. | MAE (%) | Corr. | MAE (ml) | Corr. | MAE (ml) |
| Intra-obs. | .896 | 4.7 | .978 | 6.5 | .981 | 4.5 |
| 2D nnU-Net | .857 | 4.7 | **.977** | **5.9** | **.987** | **4.0** |
| CLAS | **.926** | **4.0** | .958 | 7.7 | .979 | 4.4 |
| GUDU | .897 | **4.0** | **.977** | 6.7 | .981 | 4.6 |

## What are the conclusions of the pilot CAMUS's story ?

✓ nnU-Net produces:

- ▪ accurate scores from a controlled dataset

- ▪ within the intra-expert variability

✓ Has the potential to replace the expert's hand !

How can these results be generalized to large-scale datasets involving data from multiple centers, multiple vendors and multiple experts?
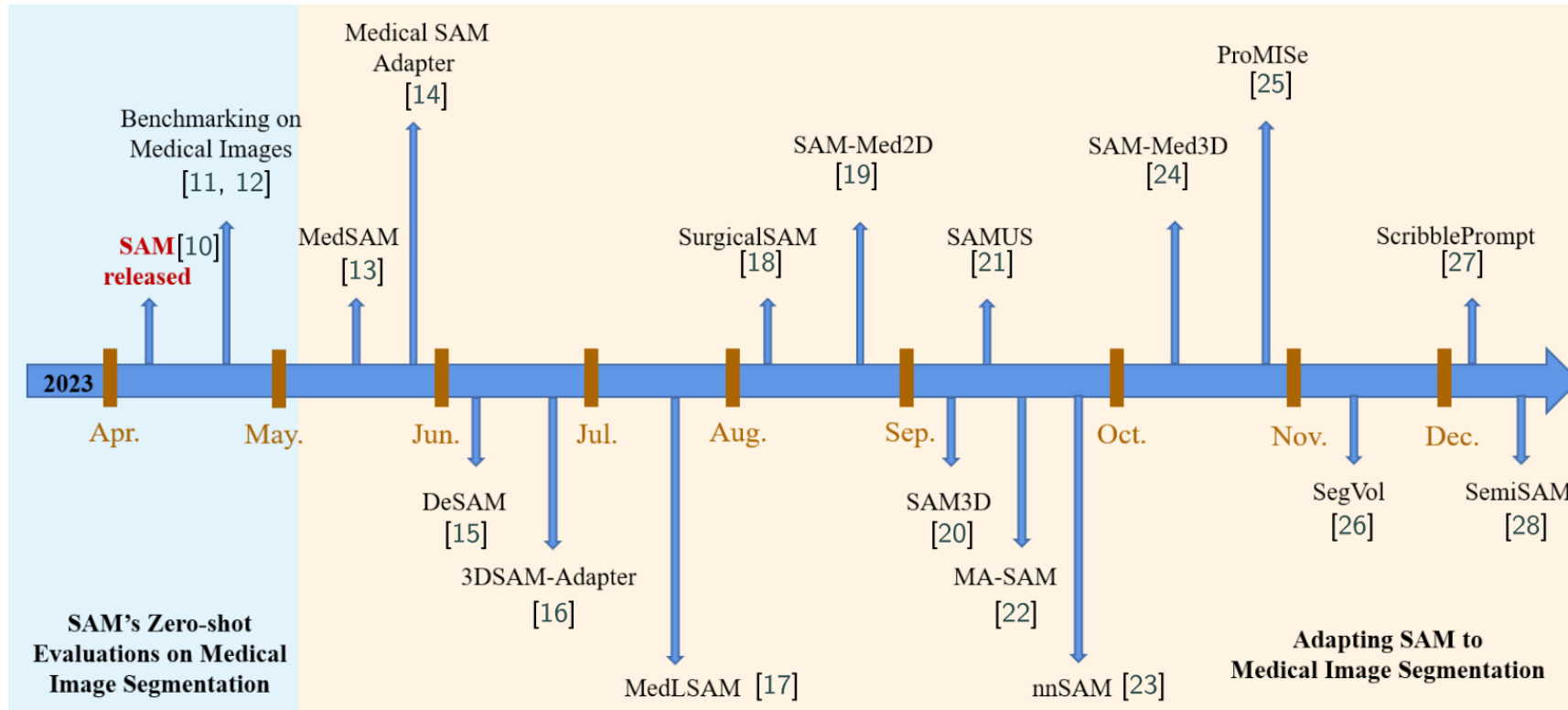
nnU-Net predictions

30

## Two tendencies

✓ Foundation models

- Learning from large scale datasets with different modalities, organs, views, ...

✓ Domain adaptation

- Efficient transfer from a source dataset (CAMUS) to a target dataset

# Brief chronology for Segment Anything (SAM) models



*From [Zhang et al., CIBM, 2024]*

## Large scale datasets

✓ SAM dataset

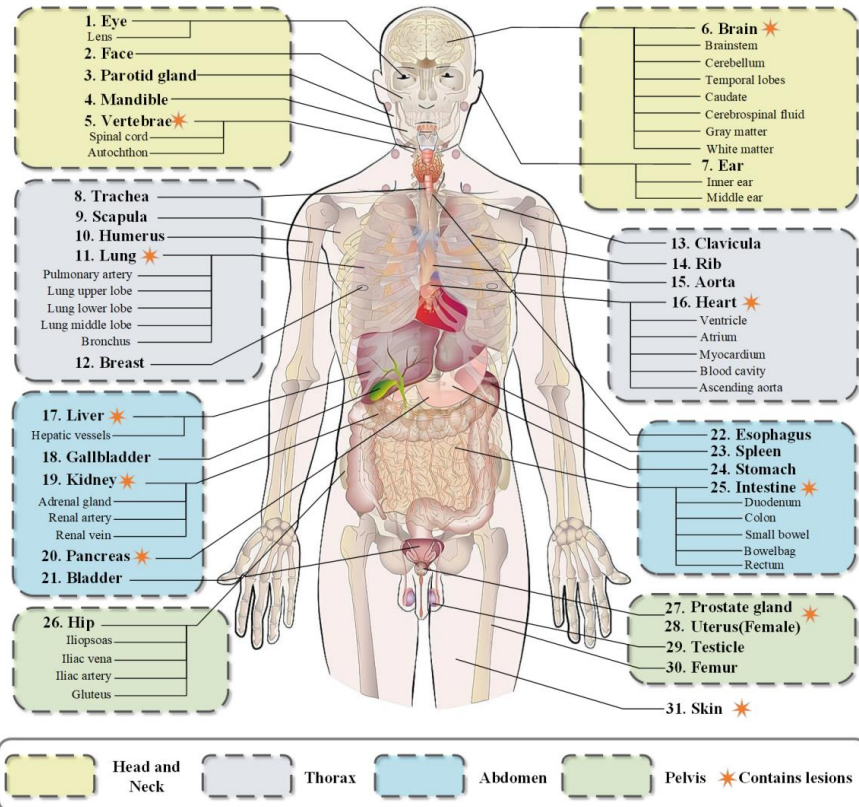✓ Licensing private dataset accessible for research purposes



*From [Kirillov et al., Arxiv, 2023]*

- **11 M images / 1 B masks**
- 2D images
- Natural scene images
- Shortest side 1500 px

## Large scale datasets

✓ SAM-Med2D dataset

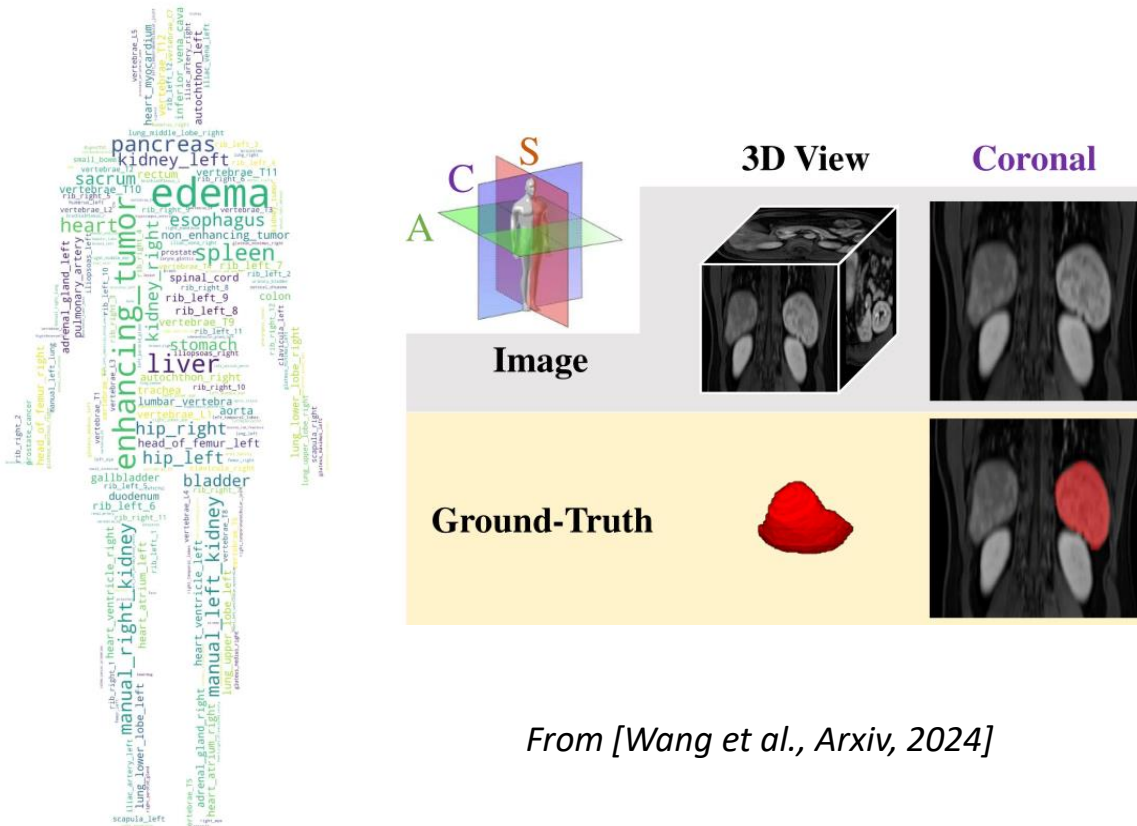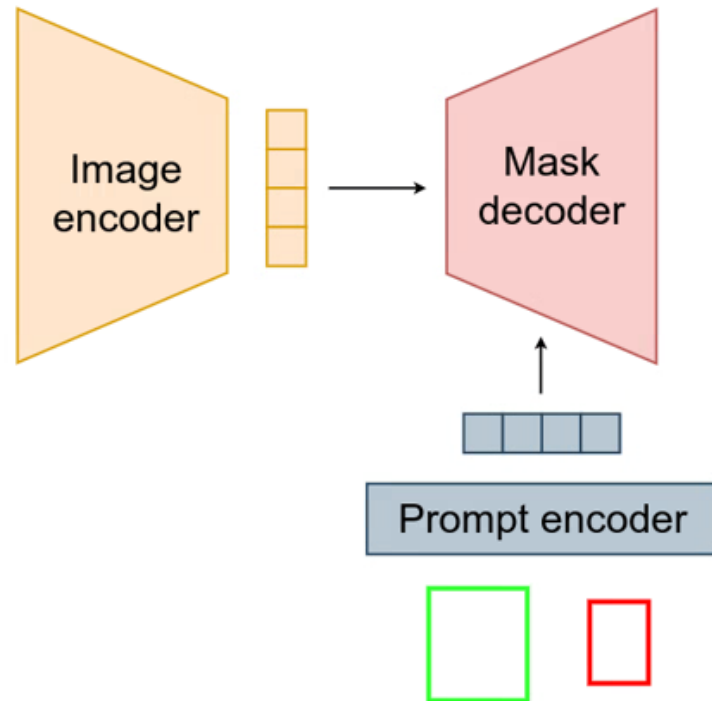✓ Collating from publicly available medical datasets + private datasets



*From [Cheng et al., Nature, 2024]*

- 4.6 M images / 19.7 M masks
- 2D images
- 10 imaging modalities
- 31 major organs
- 15% of CT images
- 256 × 256 × 3 image size
- Image intensity homogenization

## Large scale datasets

✓ SAM-Med3D dataset

✓ Collating from publicly available medical datasets + private datasets



*From [Wang et al., Arxiv, 2024]*

- 21 K images /  131 K masks
- 3D images
- 27 imaging modalities (among CT)
- 7 anatomical structures
- $128 \times 128 \times 128$ patch size

## AI architecture

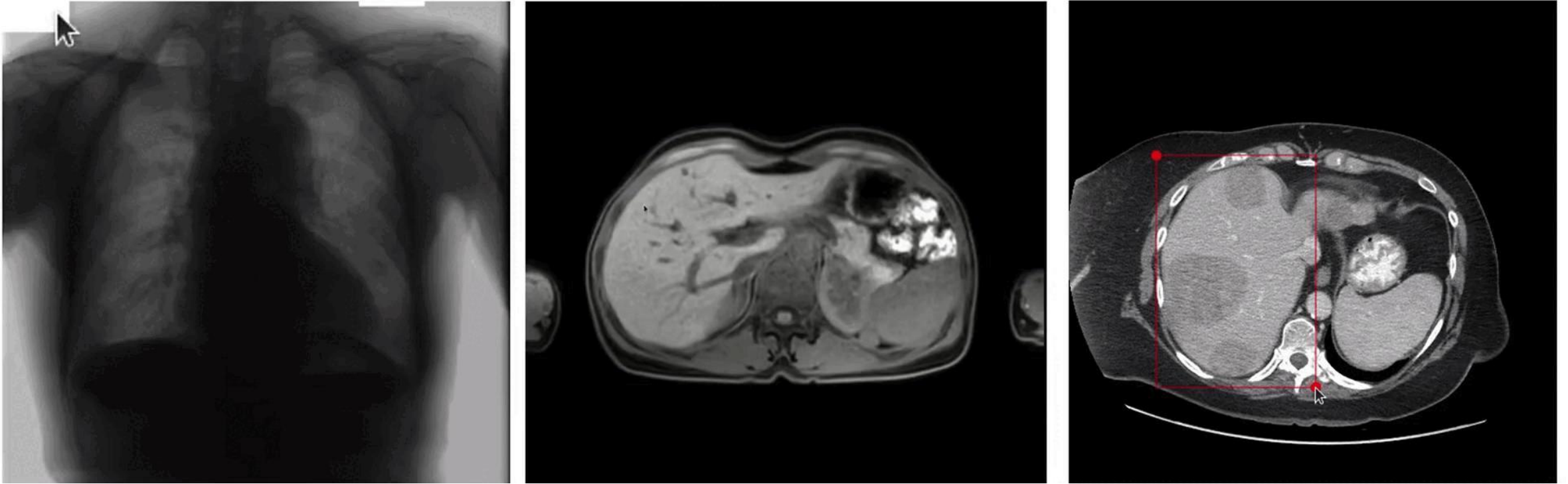- ✓ Transformer model with high complexity
- ✓ More than 91 M of parameters



**Training strategies**

- Pre-training from SAM dataset

- Fine-tuning on SAM-Med datasets

**Architecture choices**

- Freeze prompt encoder while fine-tuning image encoder and mask decoder

- Freeze image encoder while introducing learnable adapter layer, fine-tuning the prompt encoder and mask decoder
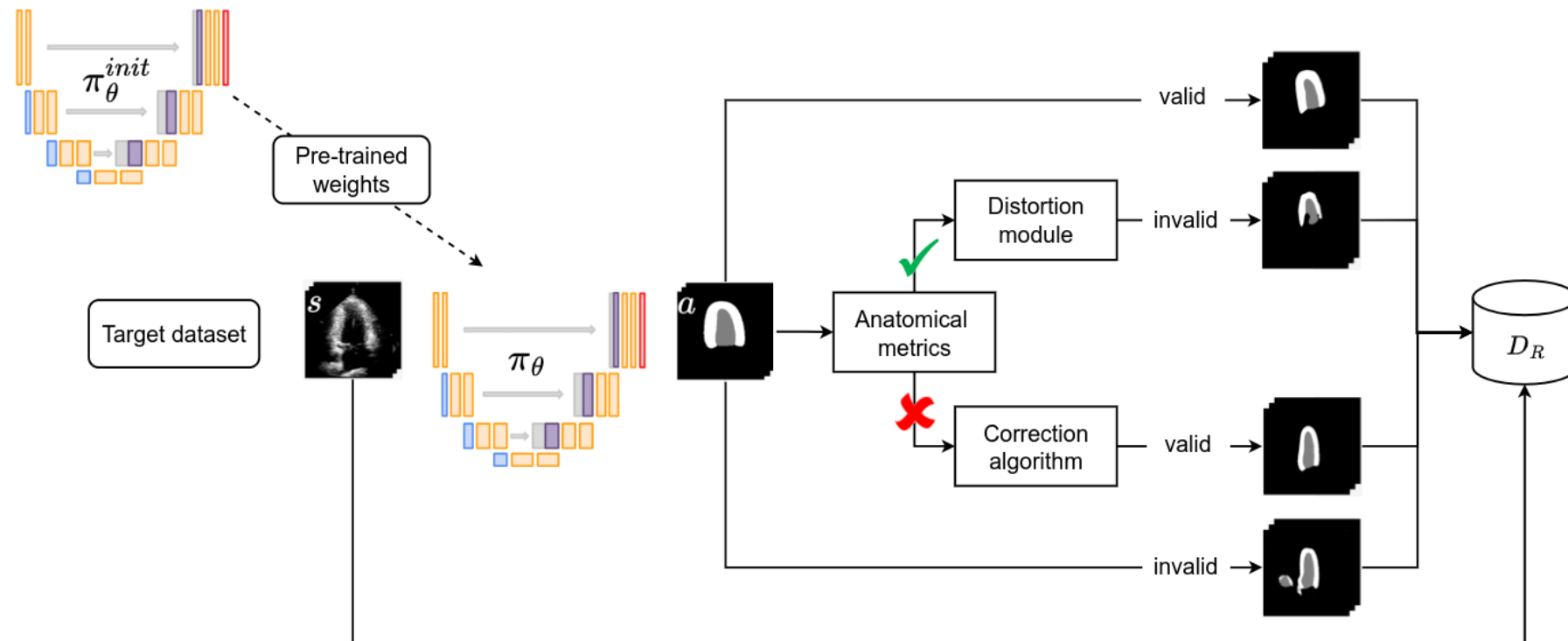
36

## Performance illustration



https://github.com/bowang-lab/MedSAM

## Two tendencies

✓ Foundation models

  ▪ Learning from large scale datasets with different modalities, organs, views, …

✓ Domain adaptation

  ▪ Efficient transfer from a source dataset (CAMUS) to a target dataset

# Inspired from reinforcement learning

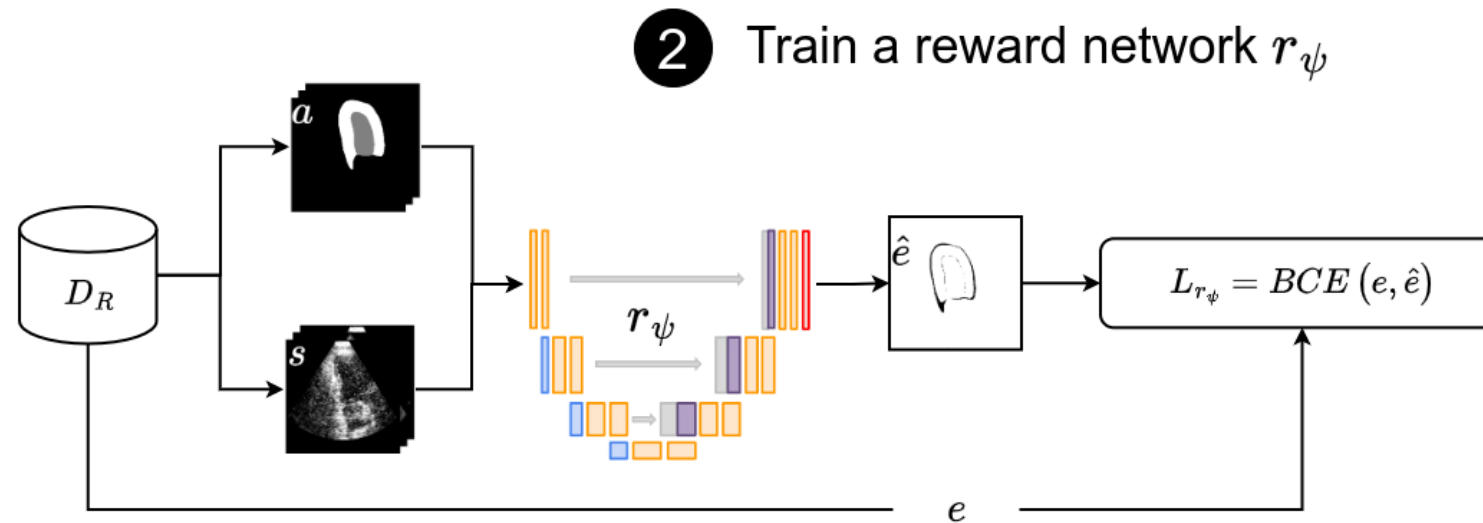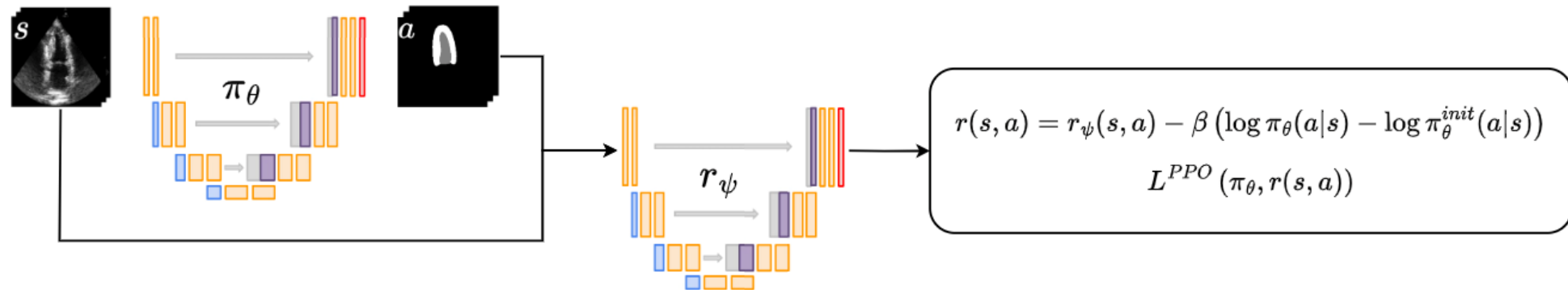✓ Update nnU-Net weights to fit with the target dataset

## Inspired from reinforcement learning

- ✓ Update nnU-Net weights to fit with the target dataset

## Inspired from reinforcement learning

✓ Update nnU-Net weights to fit with the target dataset



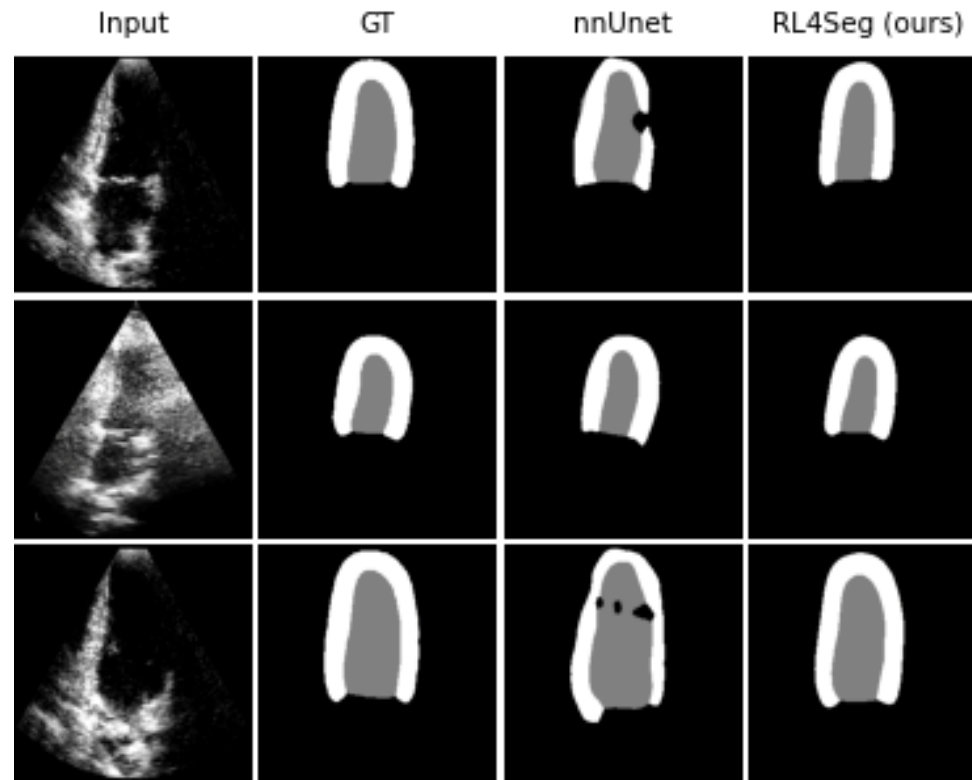③ Fine-tune the policy network $\pi_\theta$ with PPO algorithm

$$r(s,a) = r_\psi(s,a) - \beta \left( \log \pi_\theta(a|s) - \log \pi_\theta^{init}(a|s) \right)$$

$$L^{PPO}\left(\pi_\theta, r(s,a)\right)$$

## Preliminary results

✓ Scores computed from 220 patients from the target dataset

| Method | Dice (%) ↑ | | | Hausdorff (mm) ↓ | | | Anatomical Validity (%) ↑ |
|---|---|---|---|---|---|---|---|
| | ENDO | EPI | Avg. | ENDO | EPI | Avg. | |
| $\mathcal{D}_S$ intra-expert var. | 94.4 | 95.4 | 94.9 | 4.3 | 5.0 | 4.6 | 100 |
| nnUnet | 91.0 | 94.6 | 92.8 | 6.3 | 7.8 | 7.1 | 95.0 |
| RL4Seg (ours) | **91.9** | **94.7** | **93.3** | **4.9** | **5.6** | **5.3** | **98.9** |

# Preliminary results

✓ Scores computed from 220 patients from the target dataset

# Conclusions & Perspectives

► Conclusions

- ✓ AI methods have already revolutionized medical image segmentation

- ✓ Pilot studies have shown that such methods can faithfully reproduce the hand of an expert

► Perspectives

- ✓ Intensive studies on the generalization of AI model to large scale dataset

- ✓ We are undoubtedly witnessing the resolution of the segmentation problem in medical imaging!

# Thanks

# Appendices

## Convolution reminder

$$\begin{pmatrix} 0 & 1 & 2 \\ 0 & 0 & 1 \\ 0 & 1 & 2 \end{pmatrix} * \begin{pmatrix} 3 & 3 & 2 & 1 & 0 \\ 0 & 0 & 1 & 3 & 1 \\ 3 & 1 & 2 & 2 & 3 \\ 2 & 0 & 0 & 2 & 2 \\ 2 & 0 & 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 12 & 12 & 17 \\ 10 & 17 & 19 \\ 9 & 6 & 14 \end{pmatrix}$$
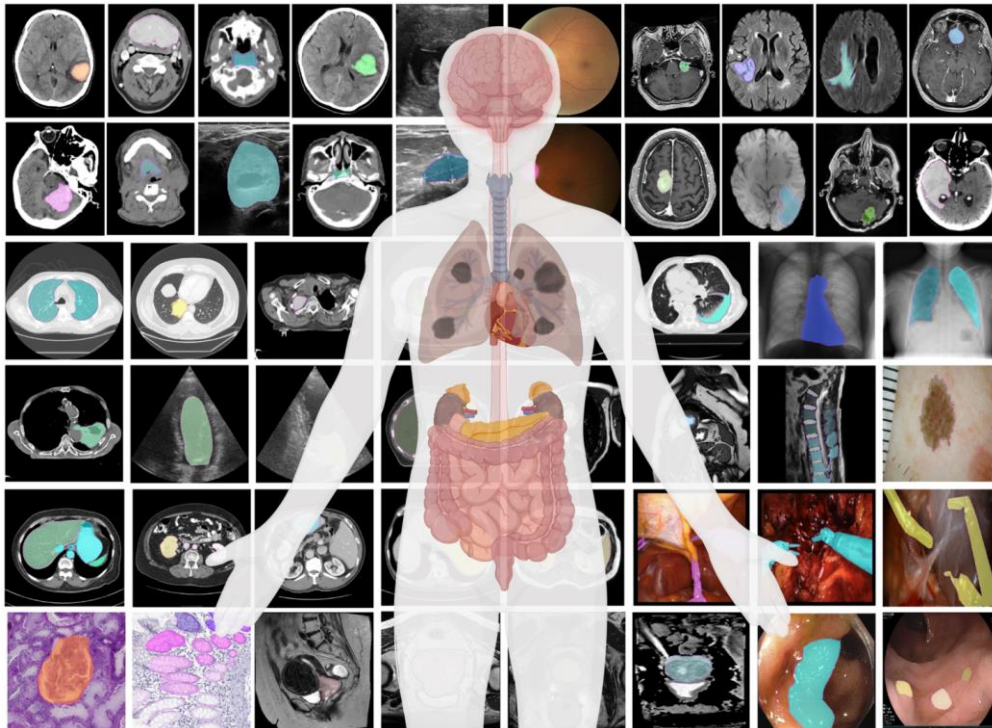
Filter         Image         Output

$$O = I * h \qquad o[i,j] = \sum_{u=-k}^{k} \sum_{v=-k}^{k} i[u,v]\, h[i-u, j-v]$$

## Large scale datasets

✓ MedSAM dataset

✓ Collating from publicly available medical datasets



*Taken from [Ma et al., Nature, 2024]*

- 1.5 M image-mask pairs
- 2D images
- 10 imaging modalities
- 30 cancer types
- 24% of CT images
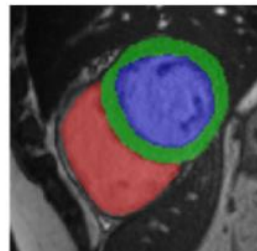- $1024 \times 1024 \times 3$ image size
- Image intensity homogenization

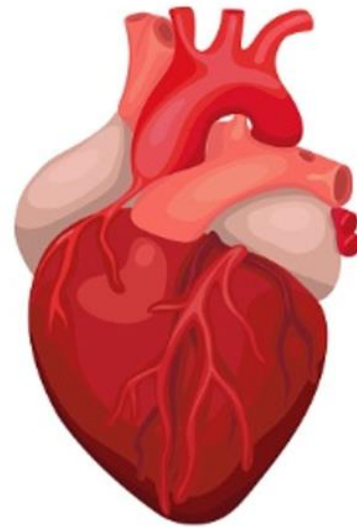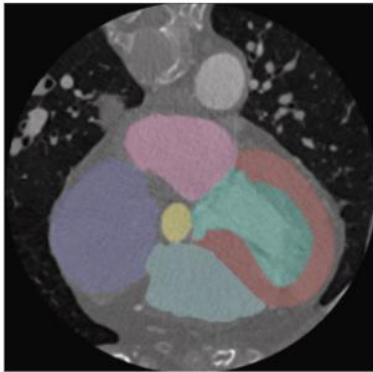Quantification of clinical indices to diagnose
cardiac pathologies

CT imaging

Echocardiographic imaging

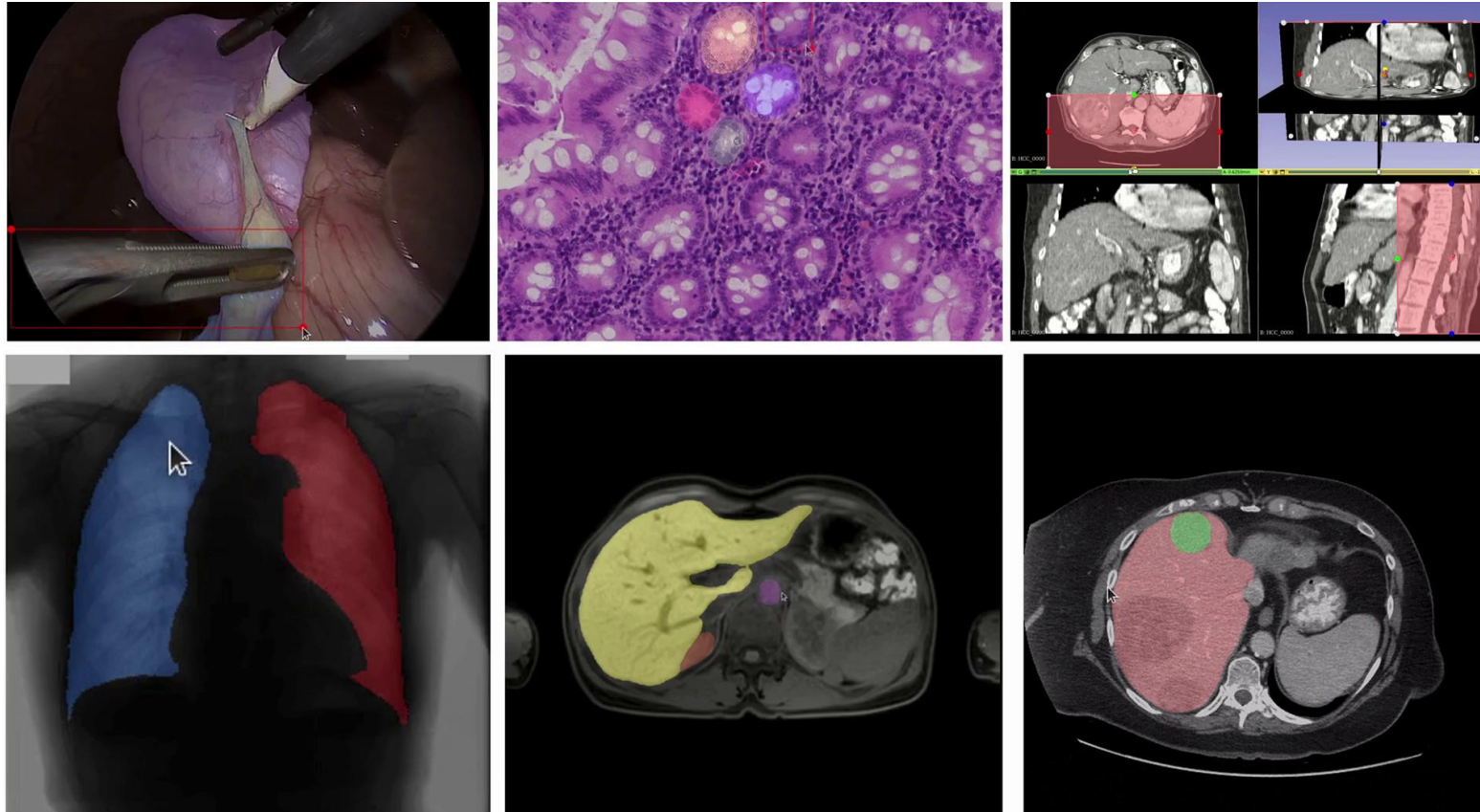MR imaging

✓ High annotation costs

✓ Inter/intra expert variability

✓ Acquisition variabilities

✓ Acquisition artifacts

## Performance illustration



https://github.com/bowang-lab/MedSAM