



3099704: AI for Digital Health



Semantic Segmentation

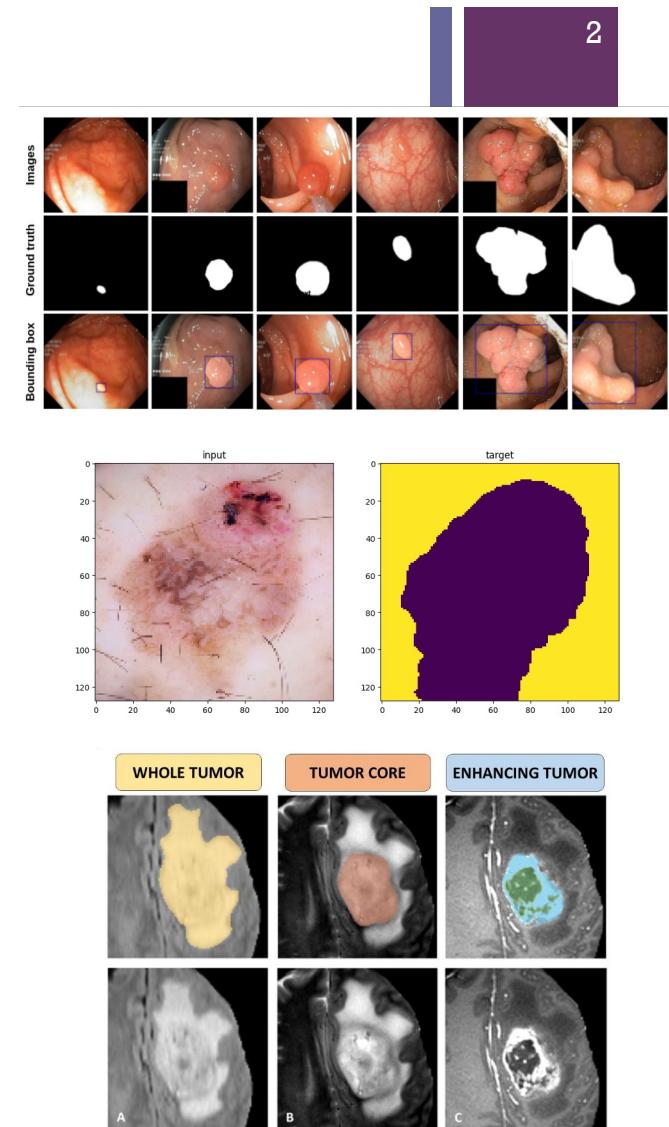
Prof. Peerapon Vateekul, Ph.D.

Peerapon.v@chula.ac.th



Outline

- Semantic Segmentation
 - UNet
 - Common Pretrained Datasets
 - Evaluation Metrics
- Case Studies
- Modern Trends: Beyond CNNs





Semantic Segmentation



Imaging Tasks

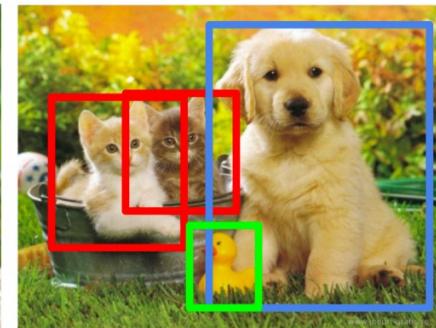
Classification



Classification + Localization



Object Detection



Instance Segmentation



CAT

CAT

CAT, DOG, DUCK

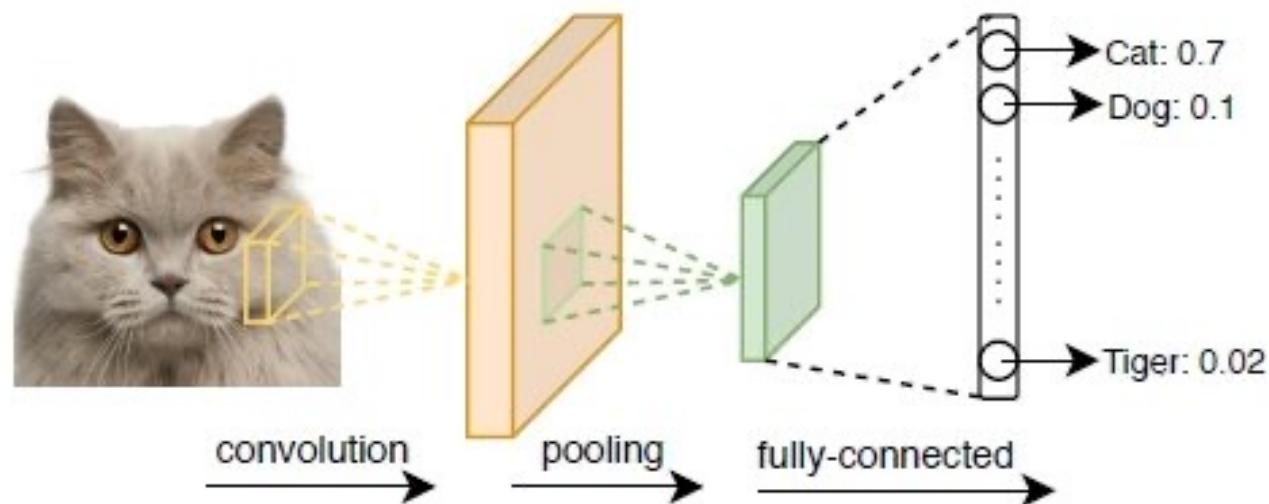
CAT, DOG, DUCK

Single object

Multiple objects

Image Classification (recap)

Convolutional Neural Network

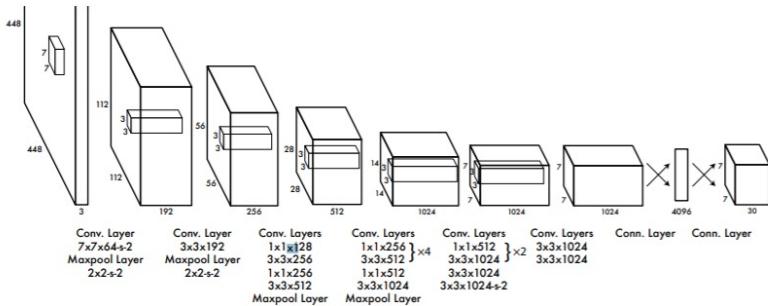




Object Detection (recap): YOLO

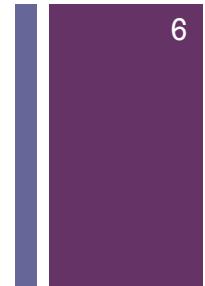
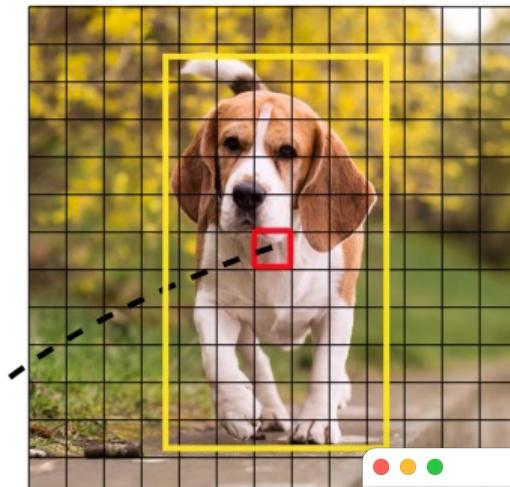


1. Resize image.
2. Run convolutional network.
3. Non-max suppression.



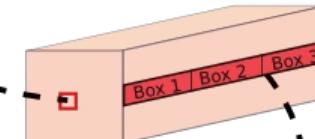
<https://blog.paperspace.com/how-to-implement-a-yolo-object-detector-in-pytorch/>

Image Grid. The Red Grid is responsible for detecting the dog

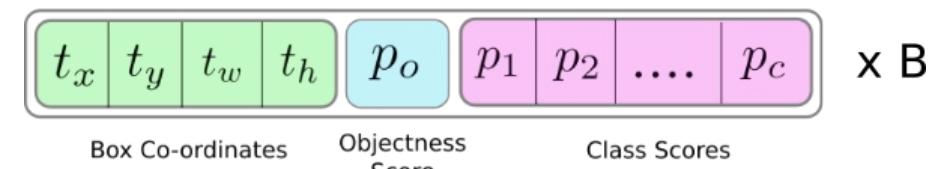


zidane.txt						
0	0.481719	0.634028	0.690625	0.713278		
0	0.741094	0.524306	0.314750	0.933389		
27	0.364844	0.795833	0.078125	0.400000		

Prediction Feature Map



Attributes of a bounding box

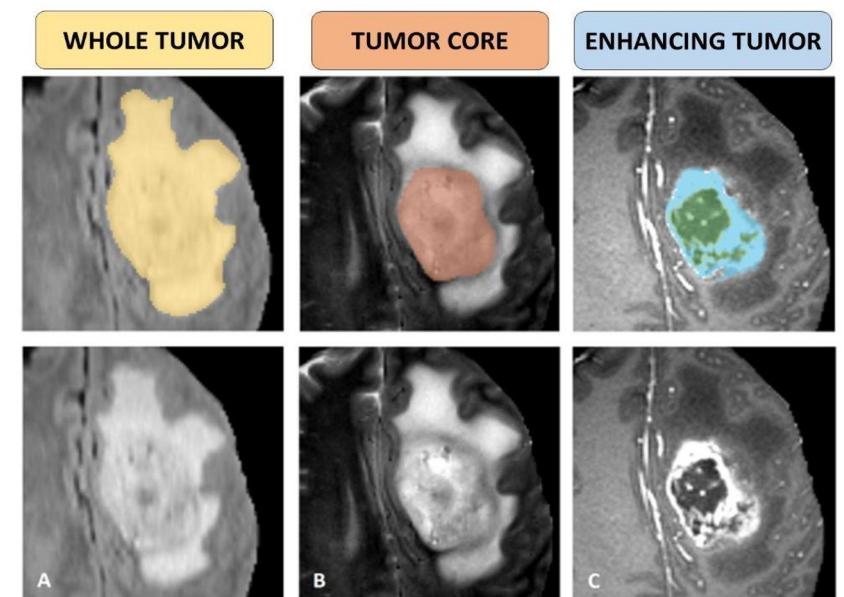
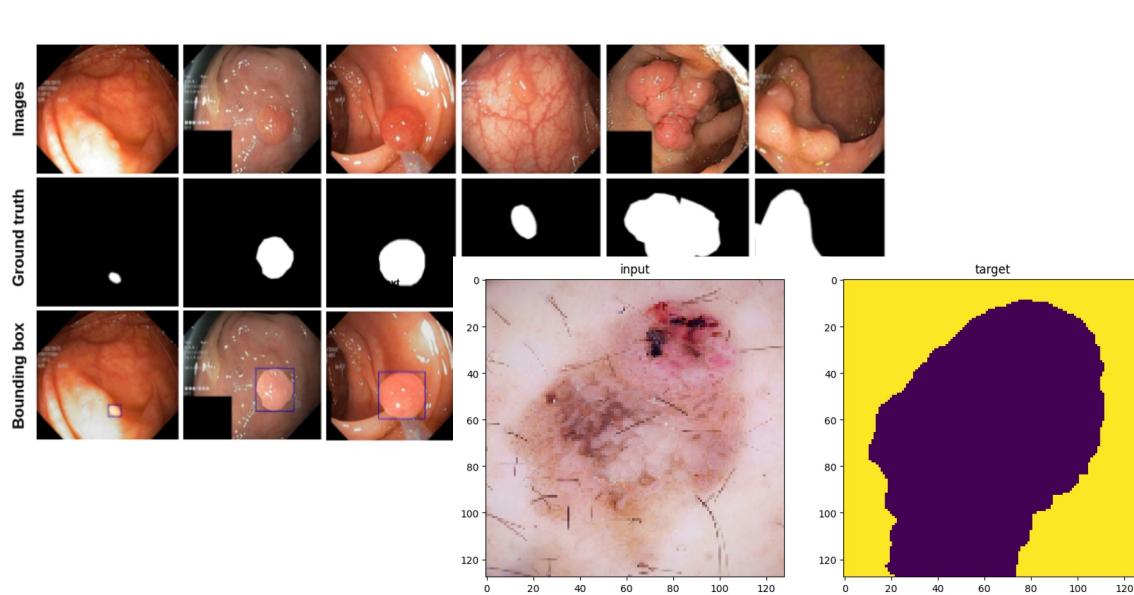


x B



Semantic Segmentation

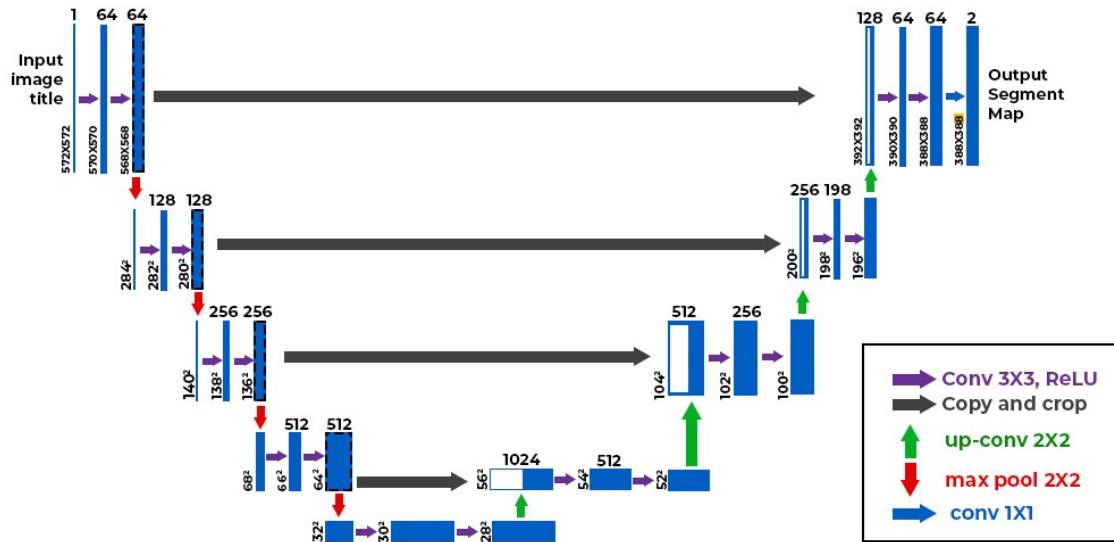
- Semantic segmentation = pixel-wise classification
- Each pixel gets a class label (road, car, polyp, tumor, sky, background).



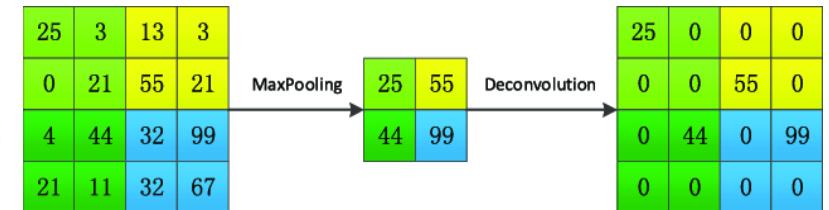


UNet: Encoder-Decoder Network: Encoder, Decoder, Skip Connections

■ U-Net



Deconvolutional layer



$$\mathcal{L} = \lambda_1 \cdot \text{CE} + \lambda_2 \cdot \text{Dice} ; \text{CE (pixel loss) \& Dice (region loss; semantic)}$$

Architecture of the U-net for a given input image. The blue boxes correspond to feature maps blocks with their denoted shapes. The white boxes correspond to the copied and cropped feature maps.

Source: [O. Ronneberger et al. \(2015\)](#)



Pretrained Datasets

Model family	Typical pretraining
YOLOv8-Seg	COCO
DeepLab	COCO → Cityscapes
U-Net (medical)	MSD / domain-specific
ViT-Seg	ADE20K
SAM	Massive private + public mix

Bundle name	Dataset	Task
spleen_ct_segmentation	MSD	Spleen
liver_ct_segmentation	MSD	Liver
pancreas_ct_segmentation	MSD	Pancreas
multi_organ_ct_segmentation	BTCV	Abdominal organs





Pretrained Datasets: MSD

- Medical Segmentation Decathlon is a large-scale benchmark designed to evaluate general-purpose medical image segmentation models across multiple organs, modalities, and tasks.
- Key Characteristics
 - Modalities: CT, MRI
 - Dimensionality: 3D volumes
 - Tasks: 10 diverse segmentation challenges
- Example Tasks
 - Brain tumor (MRI)
 - Liver, spleen, pancreas (CT)
 - Lung tumor (CT)
 - Colon cancer (CT)
 - Heart & prostate (MRI)

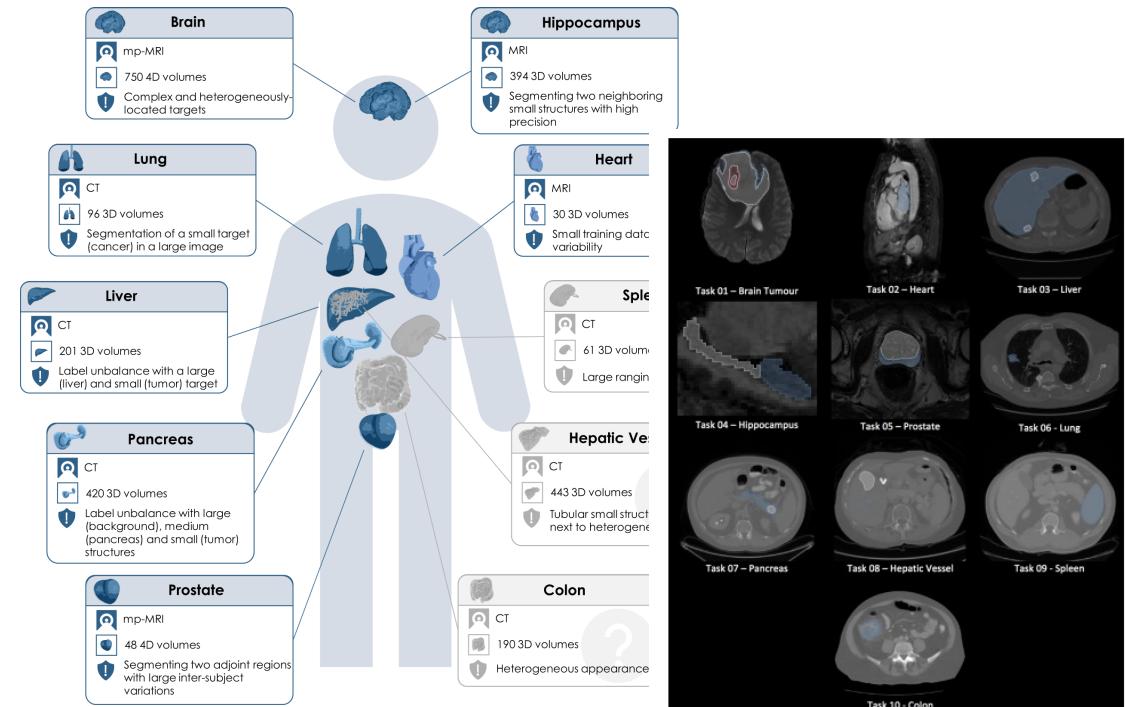
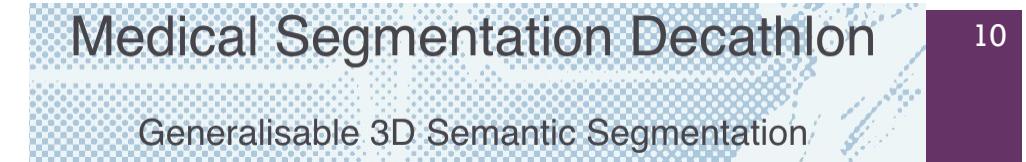
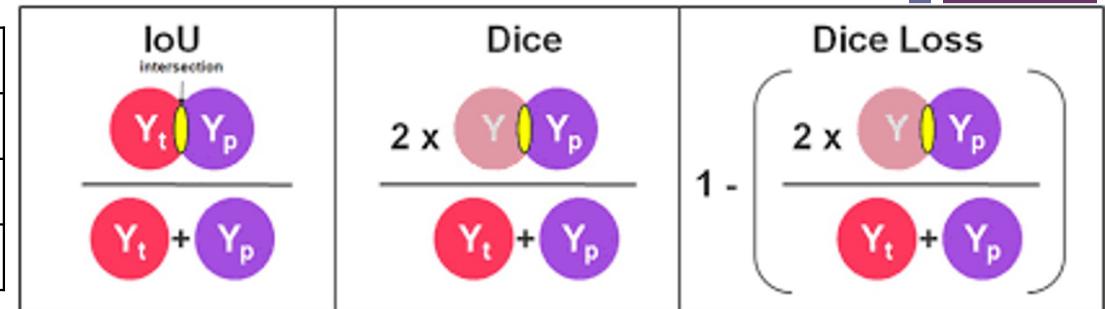


Figure 1: Exemplar images and labels for each dataset. Blue, white, and red correspond to labels 1, 2, and 3, respectively, of each dataset. Not all tasks have 3 labels.



Evaluation Metrics

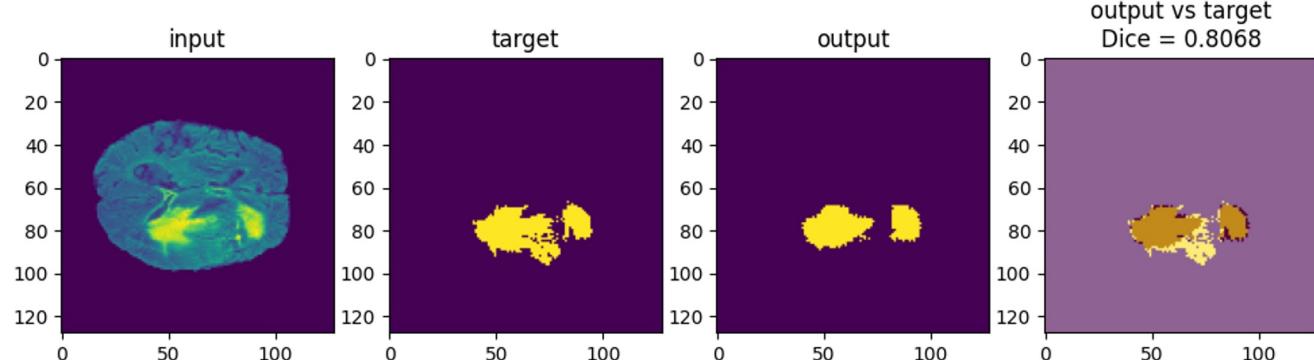
Metric	Meaning
Pixel Accuracy	Easy but misleading
Mean IoU (mIoU)	Standard benchmark
Dice / F1	Popular in medical



$$\text{IoU} = \frac{TP}{TP + FP + FN}$$

$$\text{Dice} = \frac{2TP}{2TP + FP + FN}$$

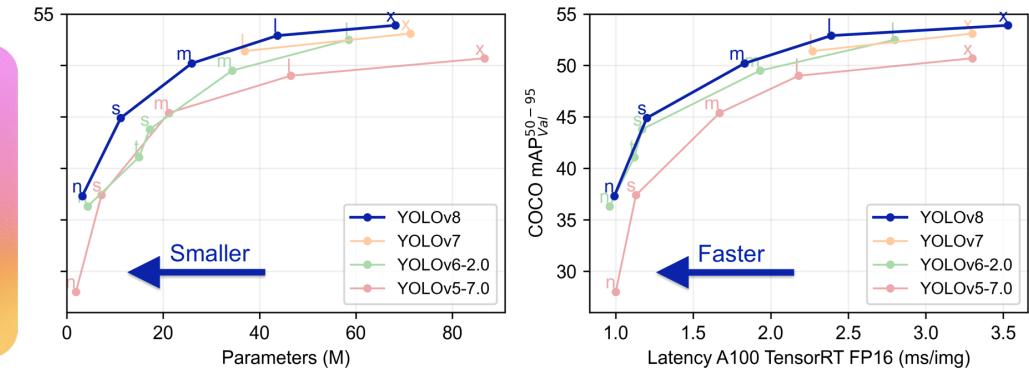
$$\text{Dice Loss} = \frac{2 \cdot \text{IoU}}{1 + \text{IoU}}$$



+ Ultralytic's YOLO (extra segmentation model)

YOLOv5 (2020) → YOLOv8 (2023) → Ultralytics YOLO 11 (2024)

12



English | 简体中文

Ultralytics CI passing | codecov 87% | DOI 10.5281/zenodo.7347926 | docker pulls 22k

Run on Gradient | Open in Colab | Open in Kaggle

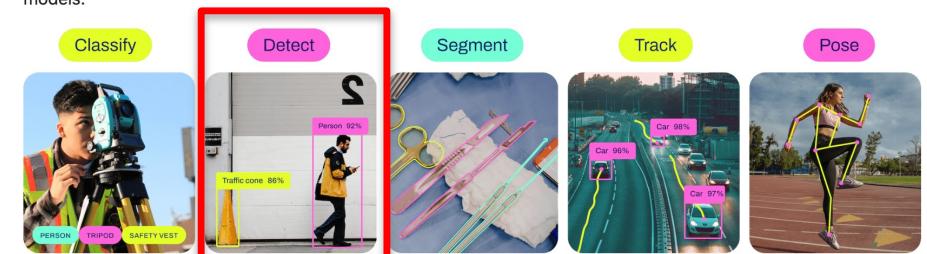
Ultralytics YOLOv8 is a cutting-edge, state-of-the-art (SOTA) model that builds upon the success of previous YOLO versions and introduces new features and improvements to further boost performance and flexibility. YOLOv8 is designed to be fast, accurate, and easy to use, making it an excellent choice for a wide range of object detection and tracking, instance segmentation, image classification and pose estimation tasks.

We hope that the resources here will help you get the most out of YOLOv8. Please browse the YOLOv8 [Docs](#) for details, raise an issue on [GitHub](#) for support, and join our [Discord](#) community for questions and discussions!

To request an Enterprise License please complete the form at [Ultralytics Licensing](#).

Models

YOLOv8 Detect, Segment and Pose models pretrained on the COCO dataset are available here, as well as YOLOv8 Classify models pretrained on the ImageNet dataset. Track mode is available for all Detect, Segment and Pose models.



All [Models](#) download automatically from the latest Ultralytics [release](#) on first use.



YOLOv8-Seg: Object Detection → Segment

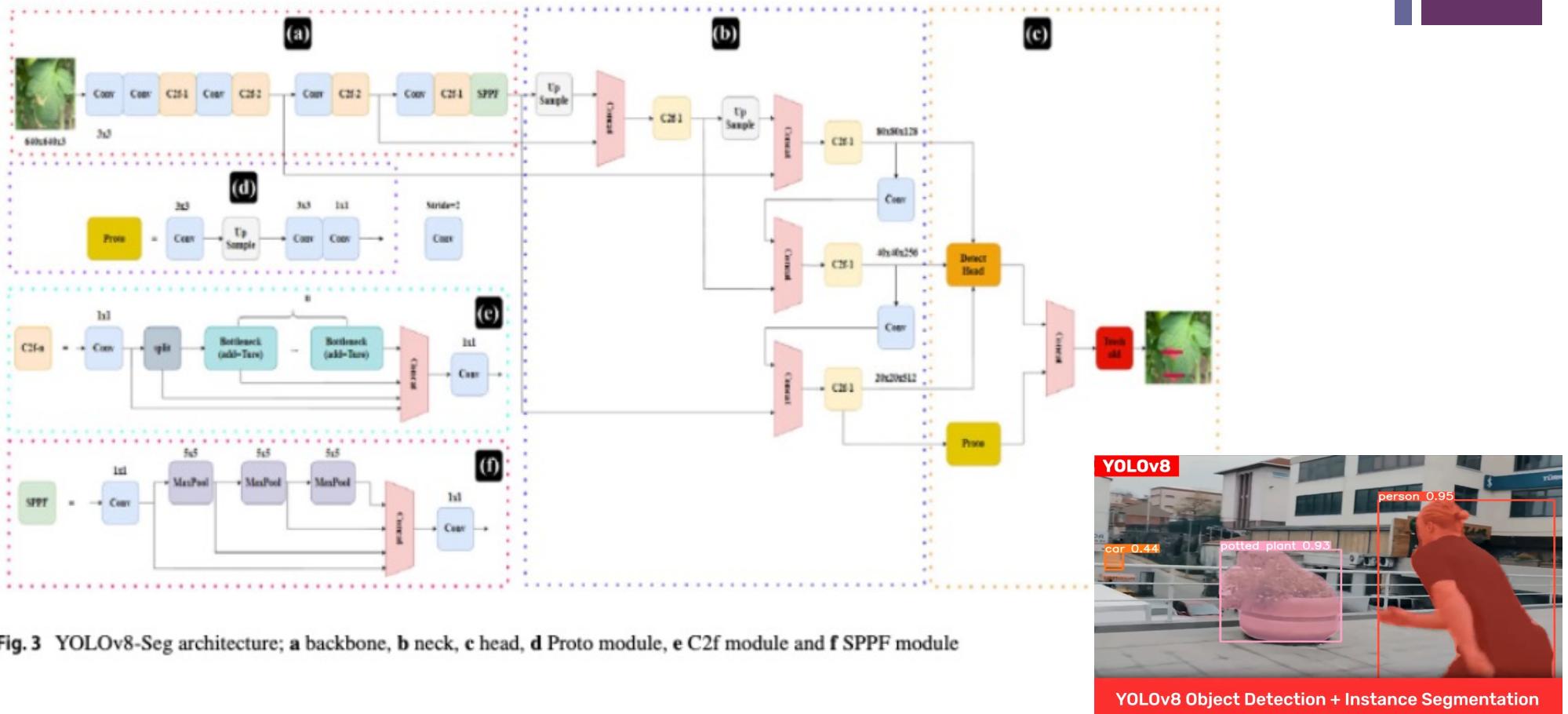


Fig. 3 YOLOv8-Seg architecture; **a** backbone, **b** neck, **c** head, **d** Proto module, **e** C2f module and **f** SPPF module

[Link](#)



Case Studies

DeepGI: GIM Segmentation

Meticuly: Automatic Skull and Mandible Reconstruction

DeepGI: Real-Time Colonic Polyp Detection

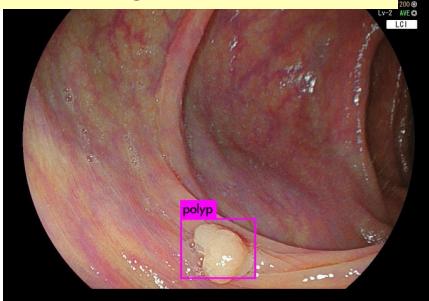
An AI-assisted solution that aims to improve colonic polyp detection in real-time. It is compatible with all standard colonoscope systems.





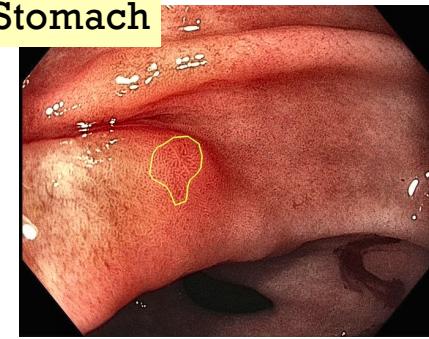
AI assisted solution in many parts of GI tracts

Colon (Large intestine)



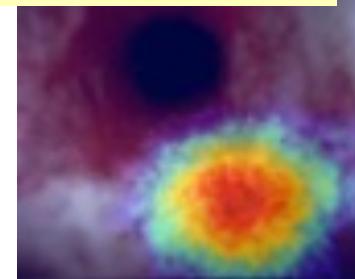
Real-time Polyp Detection
on Colonoscopy Videos

Stomach



Real-Time Gastrointestinal Metaplasia (GIM)
Segmentation on Gastroscopy Videos

Liver (biliary stricture)



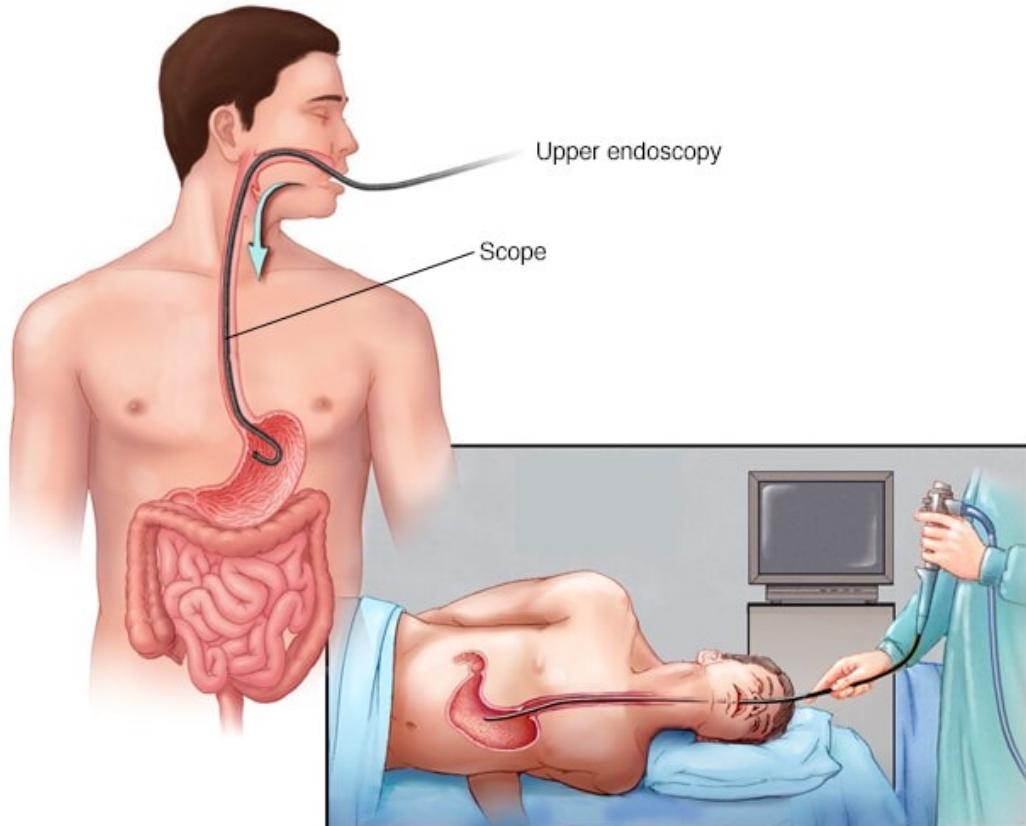
Real-Time Image Classification for
Malignant Biliary Strictures on
Cholangioscopy Images



Aims & objectives

Gastroscopy

Gastrointestinal metaplasia (GIM)



© MAYO FOUNDATION FOR MEDICAL EDUCATION AND RESEARCH. ALL RIGHTS RESERVED.



Figure: The semantic segmentation of GIM. (A) Original image. (B) Prediction: white is GIM, and black is healthy.

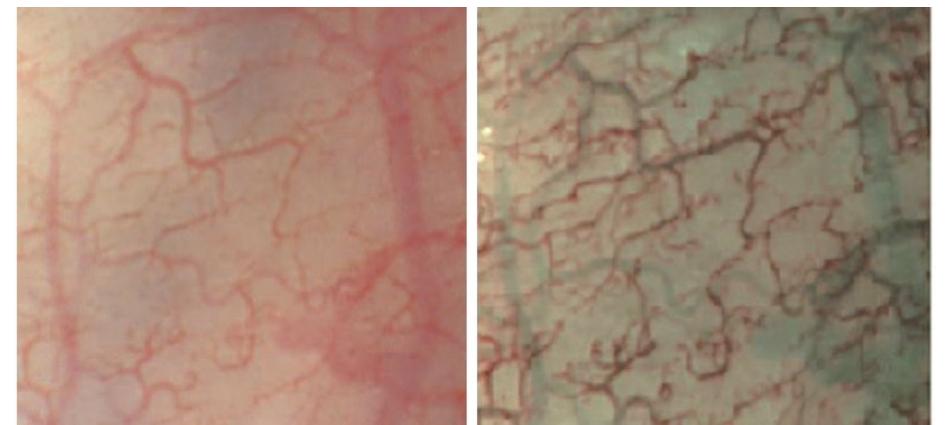


Figure: The membrane of the human tongue. (A) White light image. (B) Narrow band imaging.



Gastroscopy (new model)



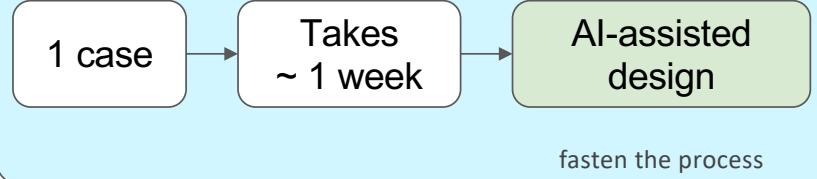
Automatic Skull and Mandible Reconstruction



Aj.Titipat Aj.Peerapon

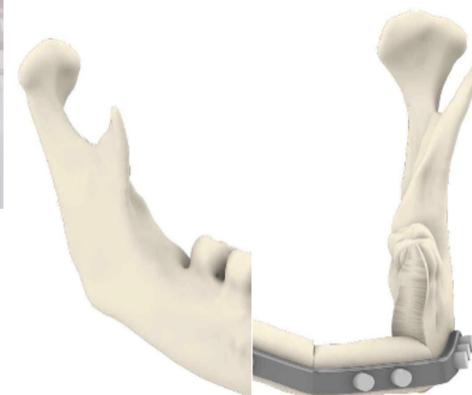
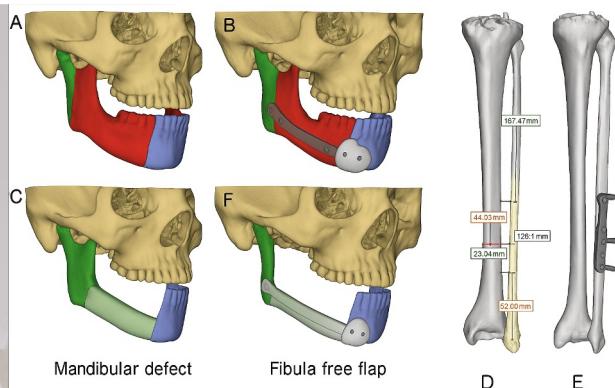
Every year, more than **10k cases** in Thailand suffer from **skull damage** and require **implant**.

- Meticuly has been doing **patient-specific 3D printing**.
- Provided → 1000 personalized cranial implant (Cranioplasty)
- We apply Autoimplant and Mandible reconstruction in the pipeline



Credit: Aj.Titipat's slide

Background and Motivation



Similarly, thousands of patients suffer from **mandibular defects** due to cancer, trauma, or infection, which severely impacts **chewing, speech, swallowing, and facial appearance**.



Support **pre-operative planning** and **custom plate design** tailored to each patient's anatomy.

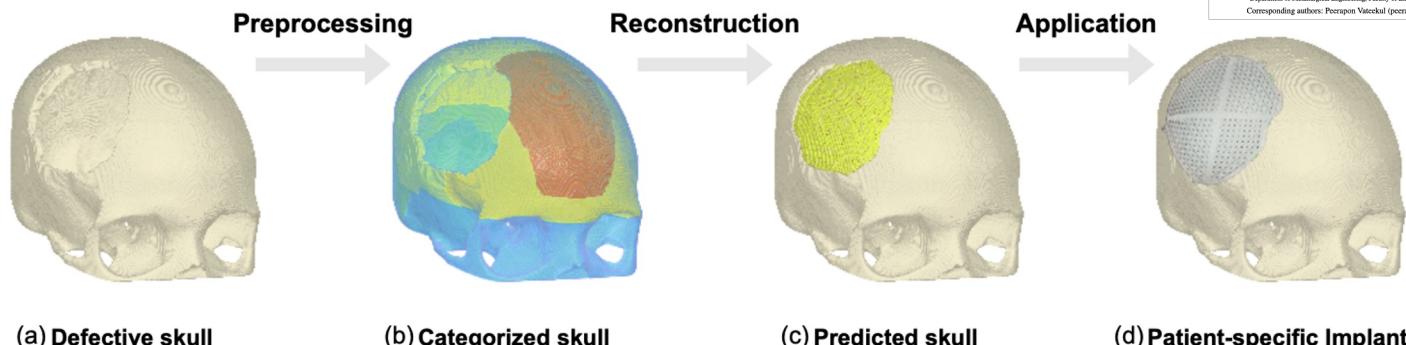


Our Aim

To enhance the **accuracy**, **efficiency**, and **robustness** of skull reconstruction on clinical defective skulls

We proposed:

- Improved deep learning backbone for 3DUNet
- Analysis on both **synthetic** and **real** clinical datasets



Received 25 May 2024, accepted 11 June 2024, date of publication 17 June 2024, date of current version 24 June 2024.
Digital Object Identifier 10.1109/ACCESS.2024.3418173

RESEARCH ARTICLE

CraNeXt: Automatic Reconstruction of Skull Implants With Skull Categorization Technique

THATHAPATT KESORNSRI¹*, NAPASARA ASAVALERTSAK², NATDANAI TANTISREEPATANA³, PORNAPAS MANOWONGPICHATE², BOONRAT LOHWONGWATANA⁴, CEDITHA SAEDEE¹, PAPIN CHAIDEE¹, AND PEERAPON VATEEKUL², ^{Corresponding Member, IEEE}

¹Department of Computer Engineering, Faculty of Engineering, Chulalongkorn University, Bangkok 10330, Thailand

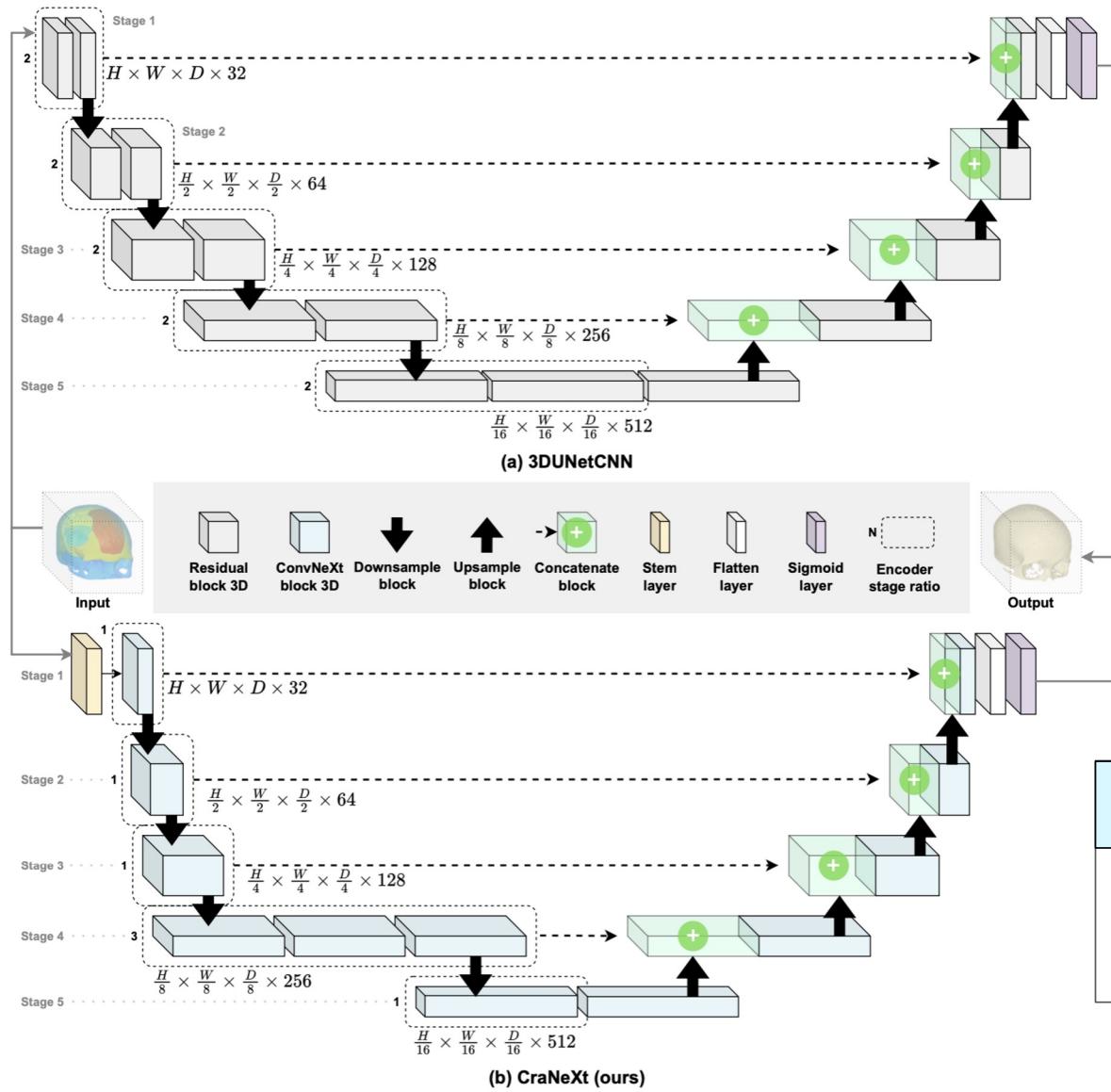
²Department of Biomedical Engineering, Faculty of Engineering, Mahidol University, Nakhon Pathom 73170, Thailand

³College of Engineering, Mahidol University, Nakhon Pathom 73170, Thailand

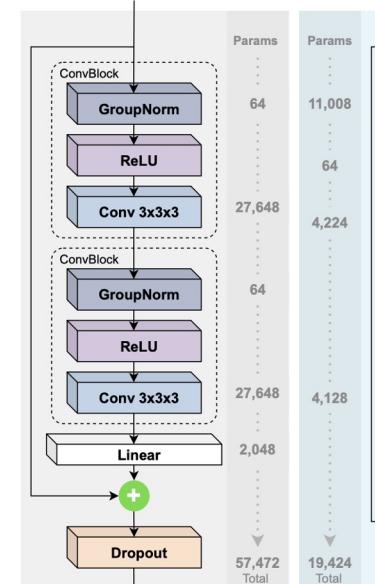
⁴Department of Mechanical Engineering, Faculty of Engineering, Chulalongkorn University, Bangkok 10330, Thailand

Corresponding authors: Peerapon Vateekul (peerapon.v@chula.ac.th) and Titipat Achakulvisut (titipat.ach@mahidol.edu)

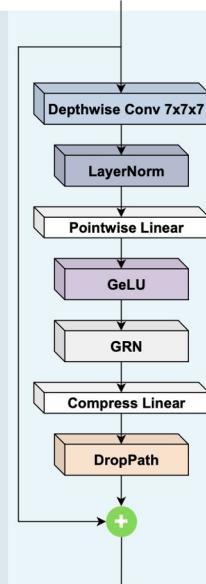
Model



(a) Residual Block 3D



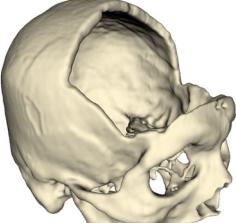
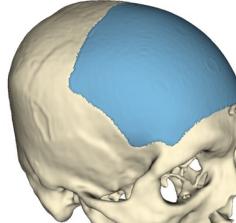
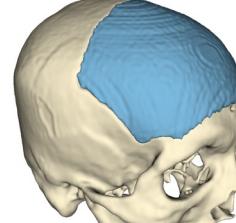
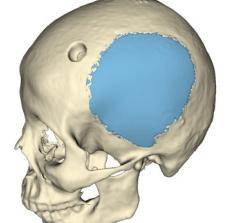
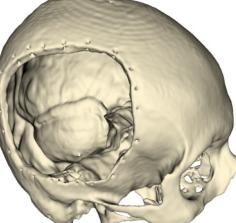
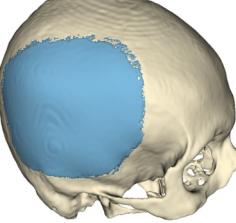
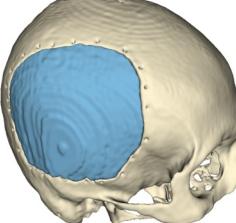
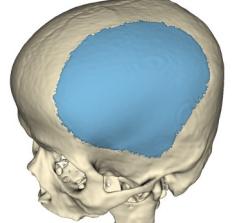
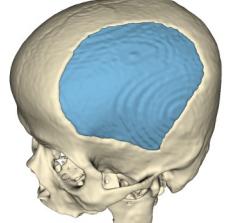
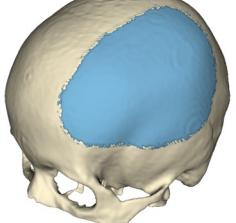
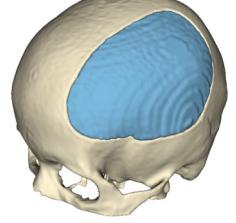
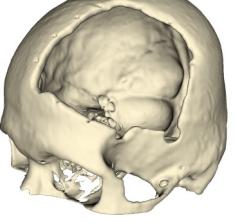
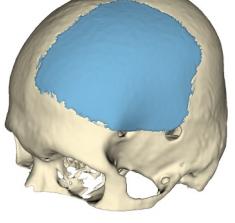
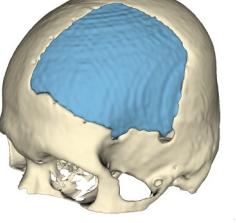
(b) ConvNeXt Block 3D



Main differences

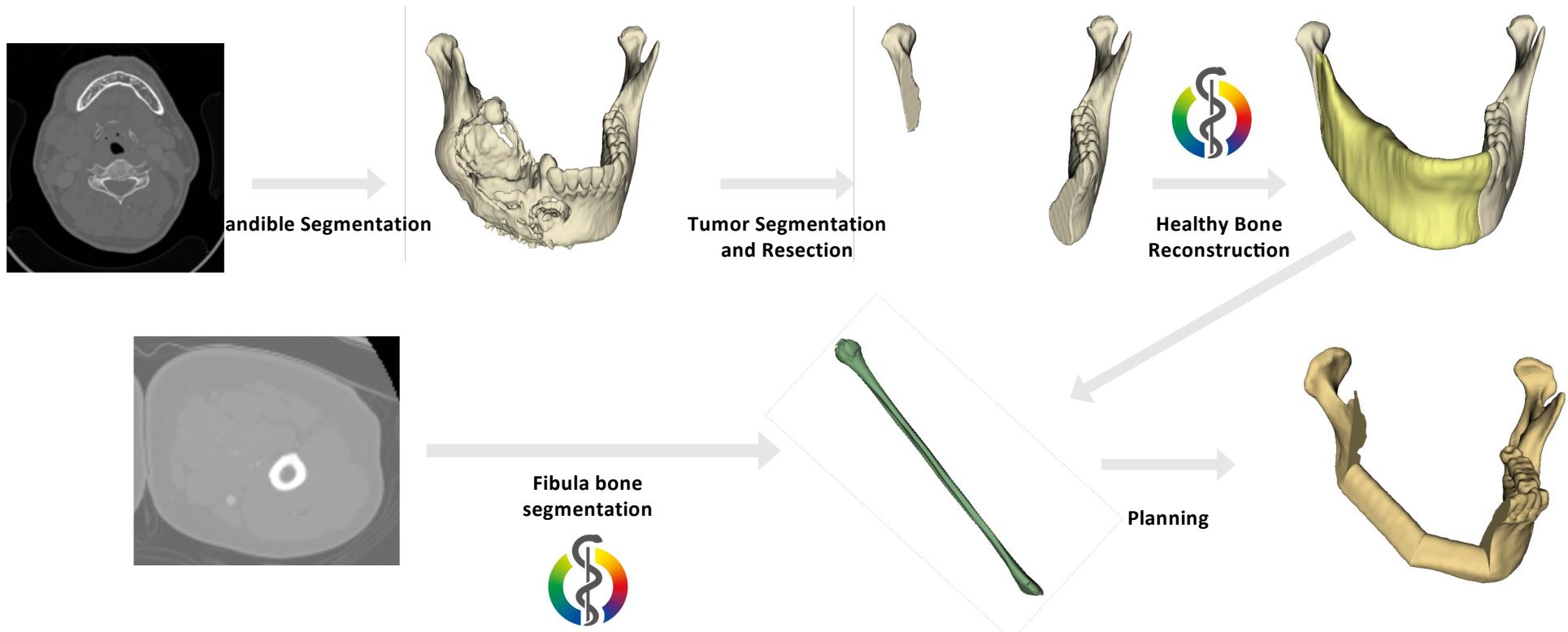
- Stage ratio
- Convolution block

Model Performance Evaluation

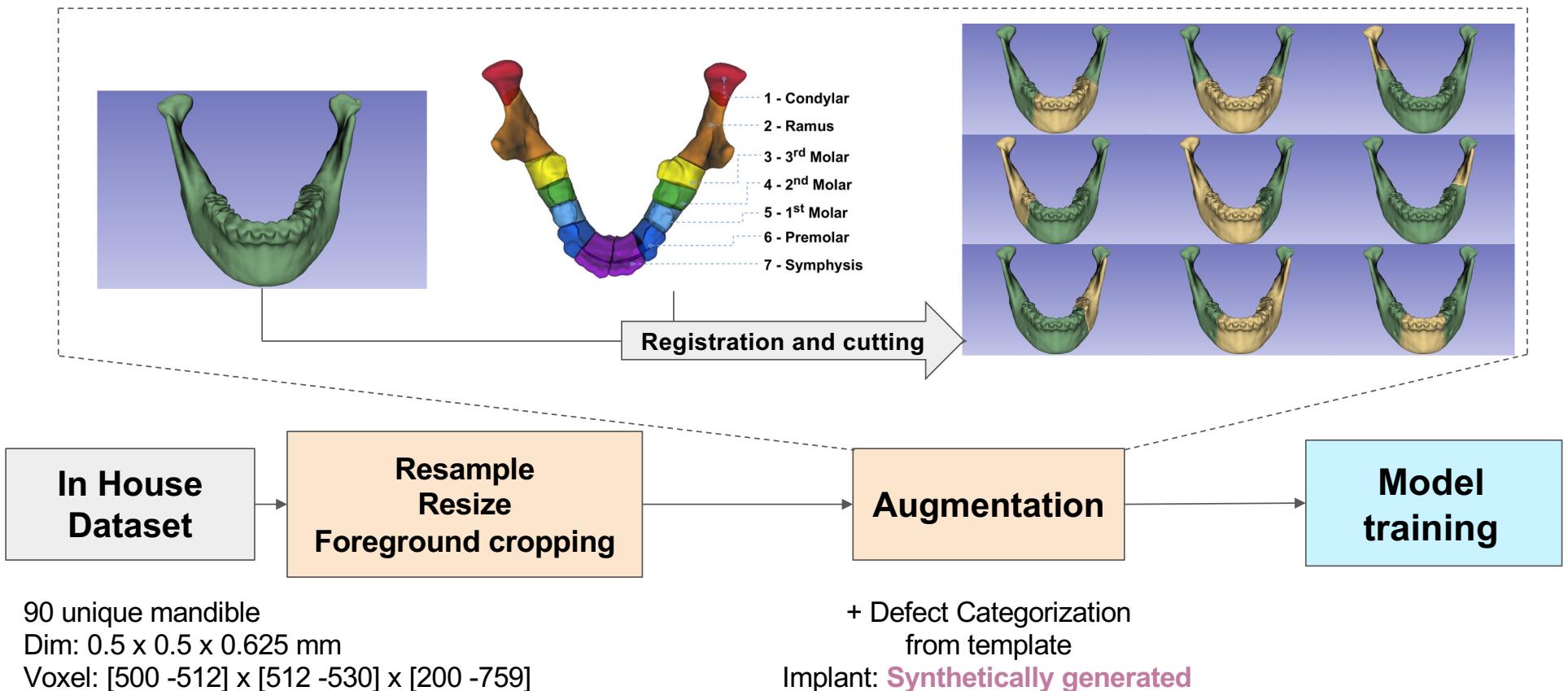
	Input	Label	Prediction		Input	Label	Prediction
Sample 1							
Sample 2							
Sample 3							

Mandible Reconstruction

- We can use a similar models for mandible reconstruction.

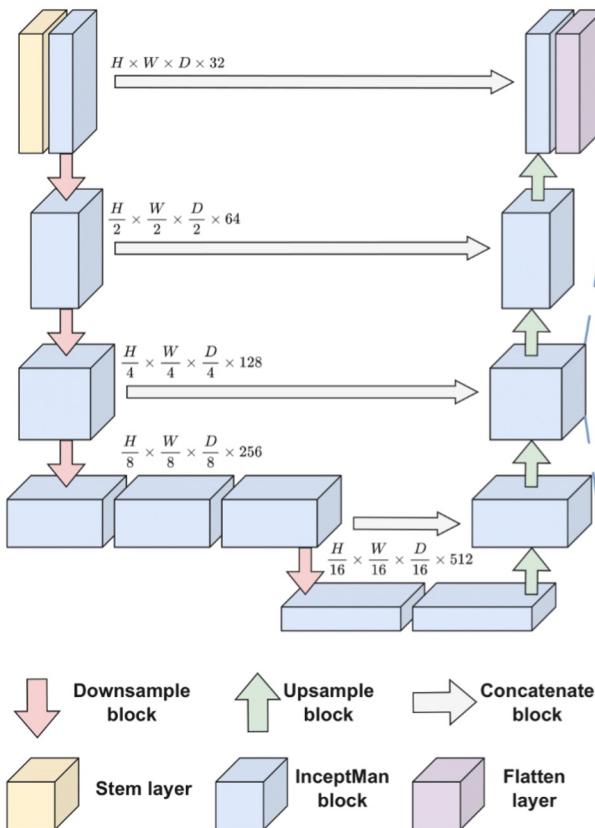


Dataset

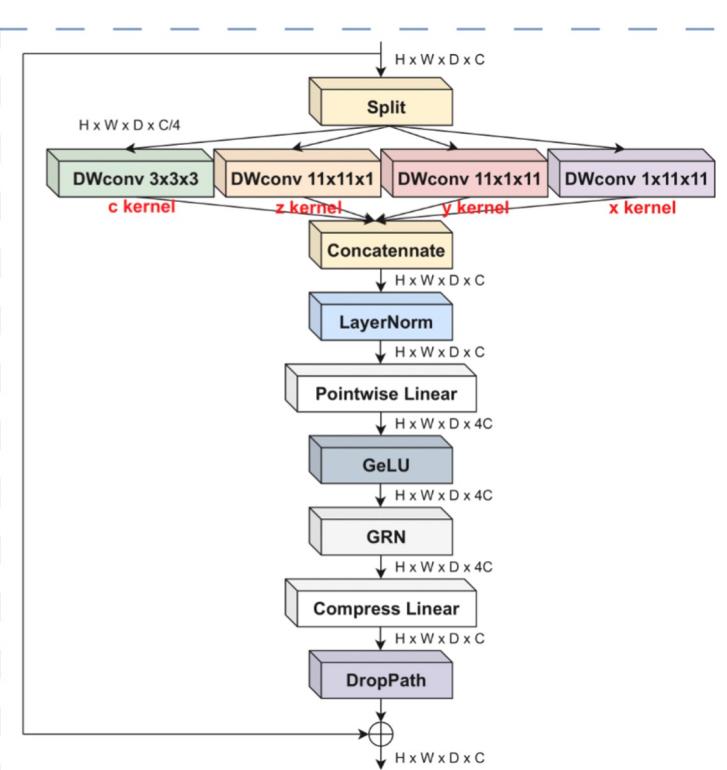


Model

(a) InceptMan architecture

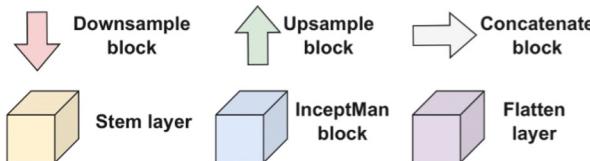


(b) InceptMan block

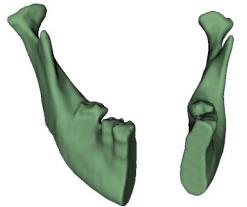
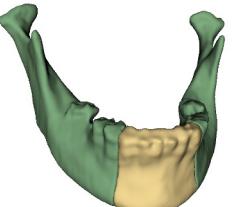
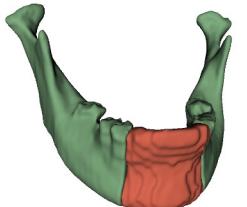
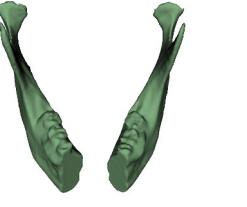
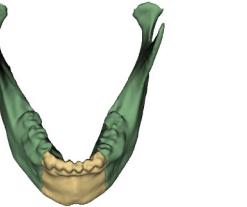
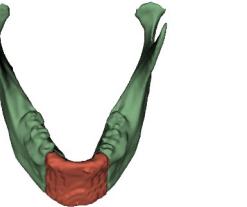
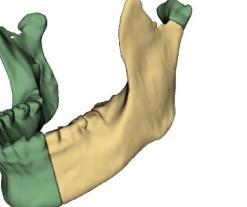
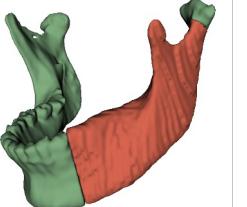
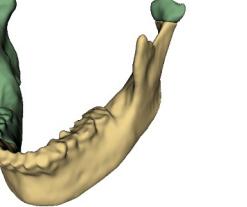
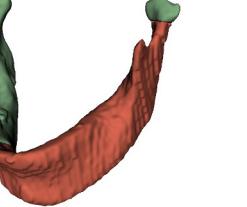
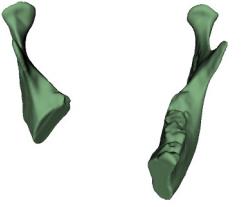
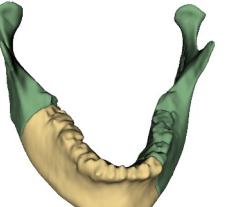
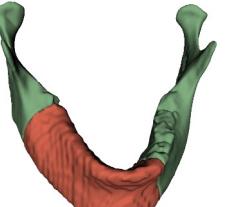
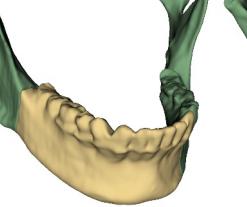
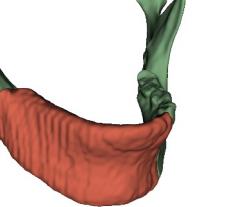


Main differences

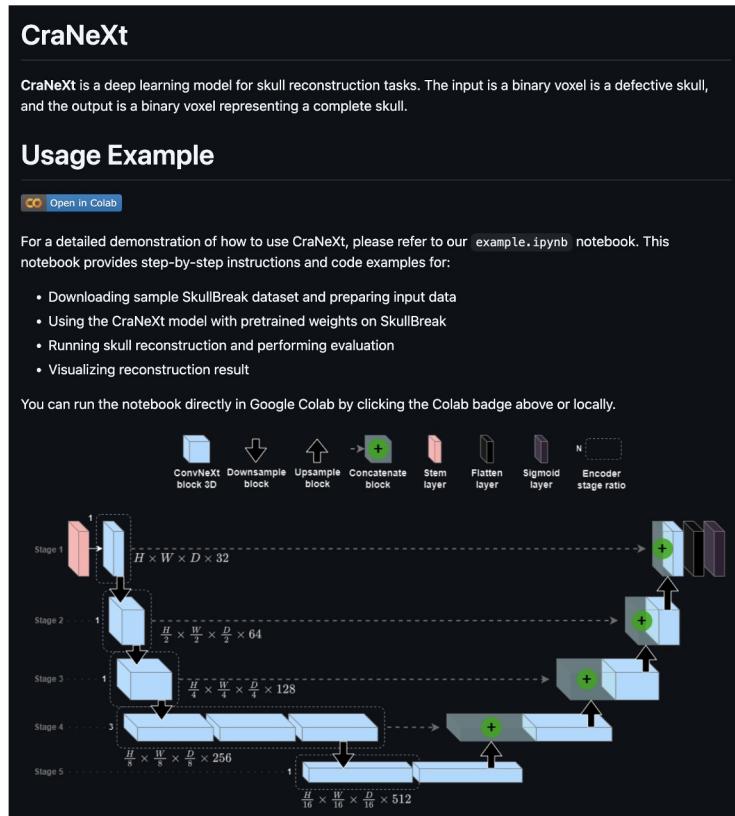
InceptMan block modified 3D ConvNeXt block with Inception block. It uses stage ratio of 1:1:1:3:1 (similar to CraNeXt).



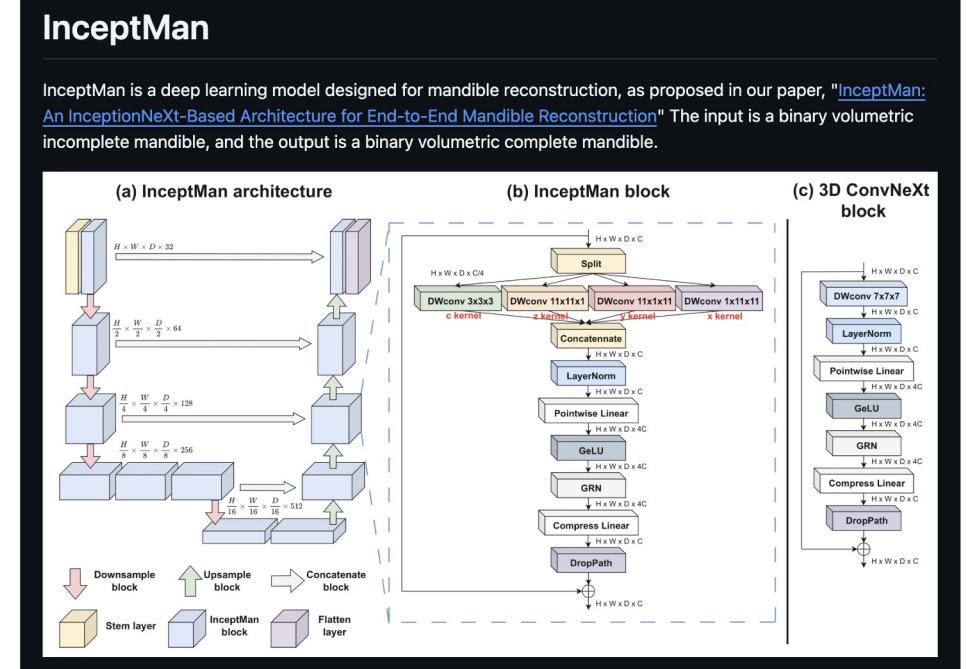
Model Performance Evaluation

	Input	Label	Prediction		Input	Label	Prediction
Sample 1							
Sample 2							
Sample 3							

CraNeXt and InceptMan are an Open Source



Available at: <https://github.com/guitared/CraNeXt>

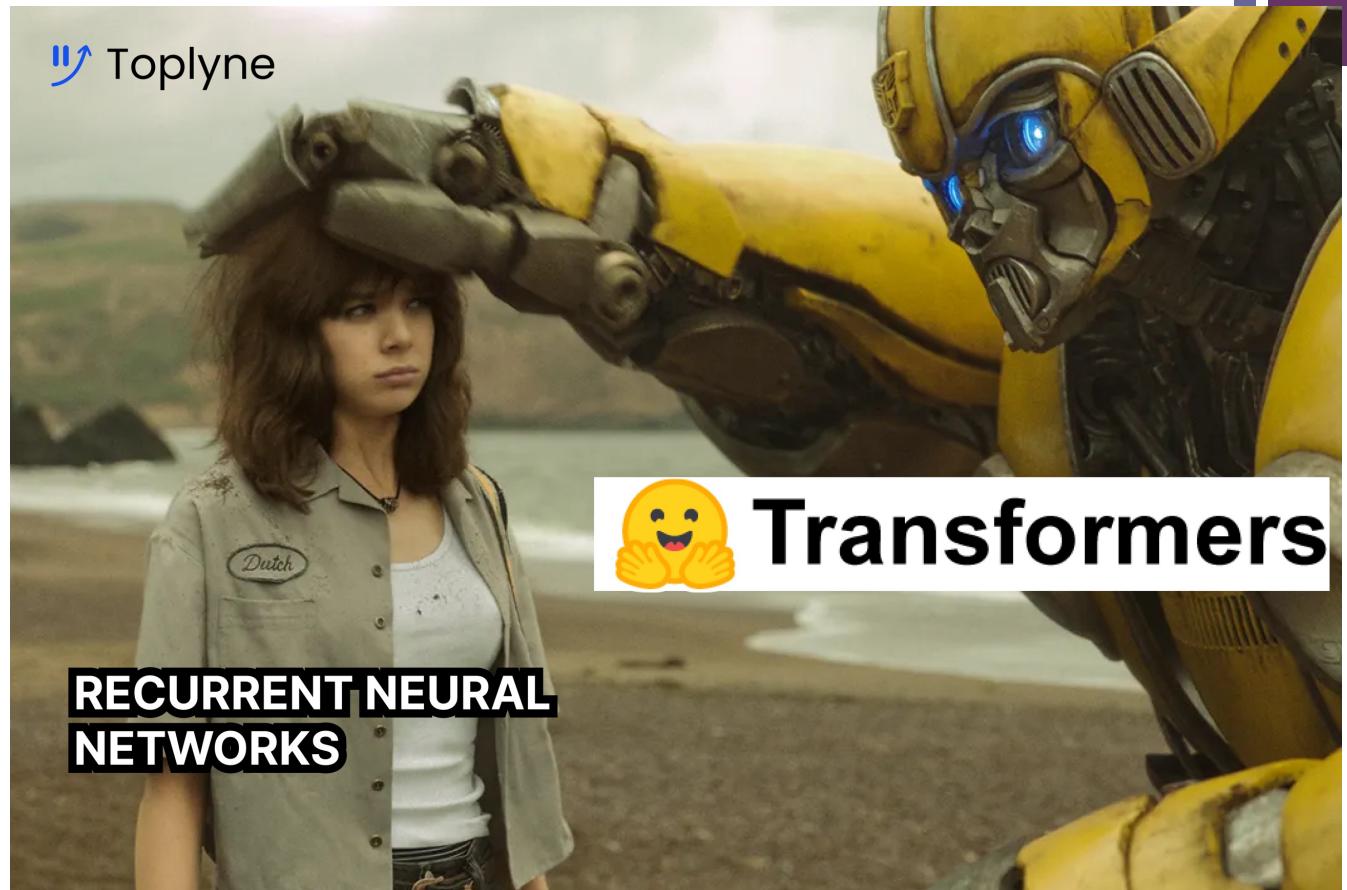
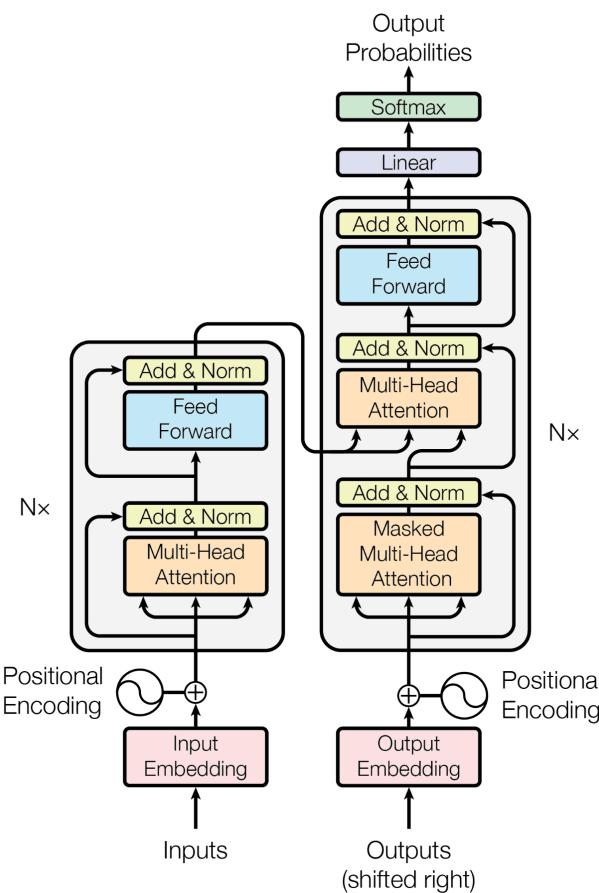


Available at: <https://github.com/oxygen-ii/InceptMan/>



Modern Trends: Beyond CNNs

Rise of Transformer (2017)





Transformer-Based Models

All Transformer based

1. Encoder-based model: BERT (2018)
2. Decoder-based model: GPT (2018)
3. Encoder and Decoder: BART (2019)

BERT [Devlin, et al, 2018]:
Bidirectional Encoder Representation from

Transformers

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova

Google AI Language

{jacobdevlin, mingweichang, kentonl, kristout}@google.com

Abstract

We introduce a new language representation model called **BERT**, which stands for Bidirectional Encoder Representations from Transformers. Unlike recent language representation models (Peters et al., 2018a; Rad-

There are two existing strategies for applying pre-trained language representations to downstream tasks: *feature-based* and *fine-tuning*. The feature-based approach, such as ELMo (Peters et al., 2018a), uses task-specific architectures that include the pre-trained representations as addi-

OpenAI GPT (Generative Pre-Training) [Radford, 2018]

Improving Language Understanding by Generative Pre-Training

Alec Radford

OpenAI

alec@openai.com

Karthik Narasimhan

OpenAI

karthikn@openai.com

Tim Salimans

OpenAI

tim@openai.com

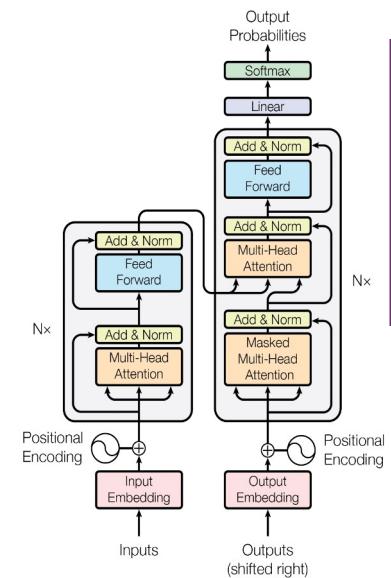
Ilya Sutskever

OpenAI

ilyasu@openai.com

Abstract

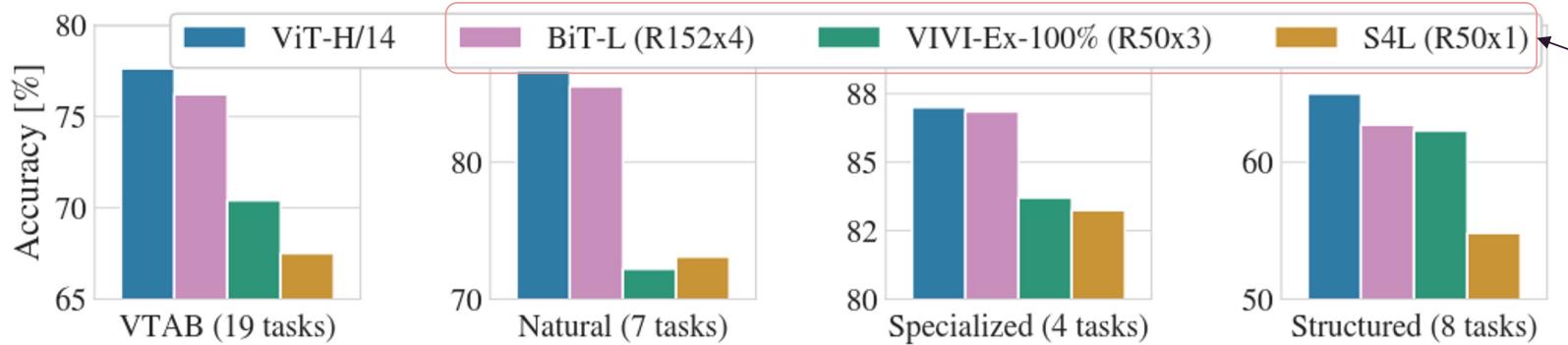
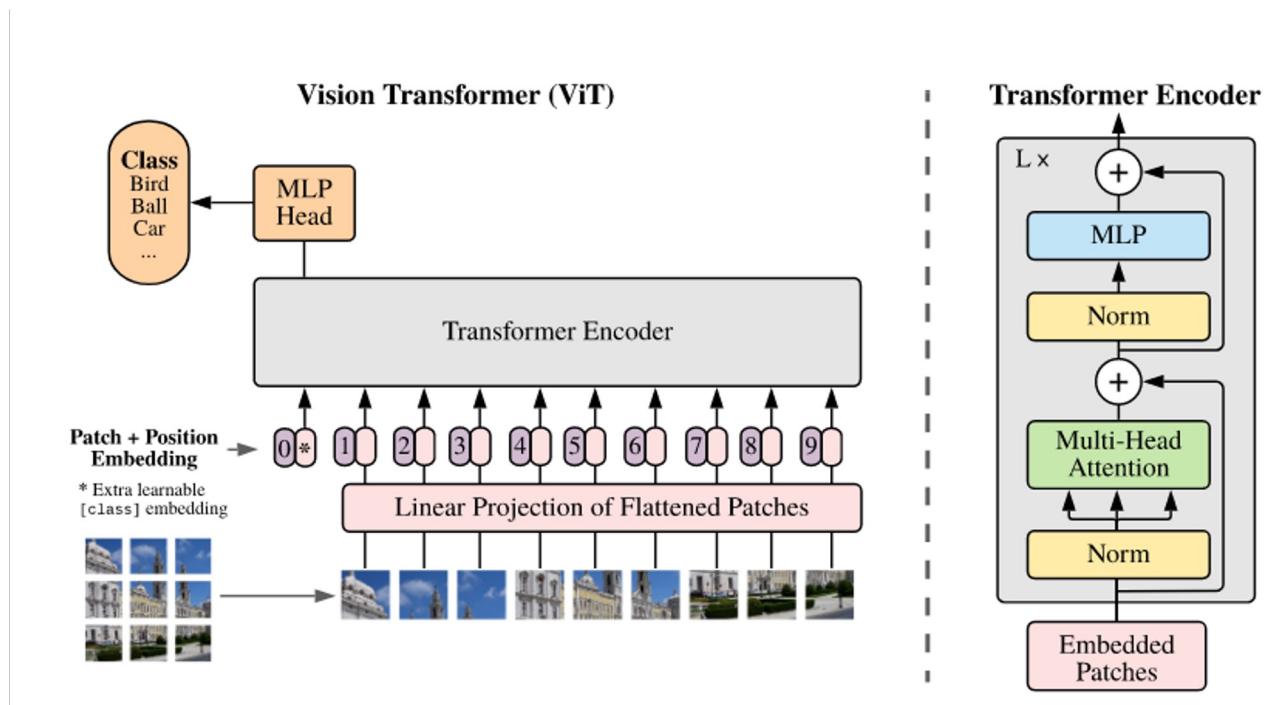
Natural language understanding comprises a wide range of diverse tasks such as textual entailment, question answering, semantic similarity assessment, and document classification. Although large unlabeled text corpora are abundant, labeled data for learning these specific tasks is scarce, making it challenging for discriminatively trained models to perform adequately. We demonstrate that large gains on these tasks can be realized by *generative pre-training* of a language model



Beyond NLP

Computer Vision

ViT uses only encoder.

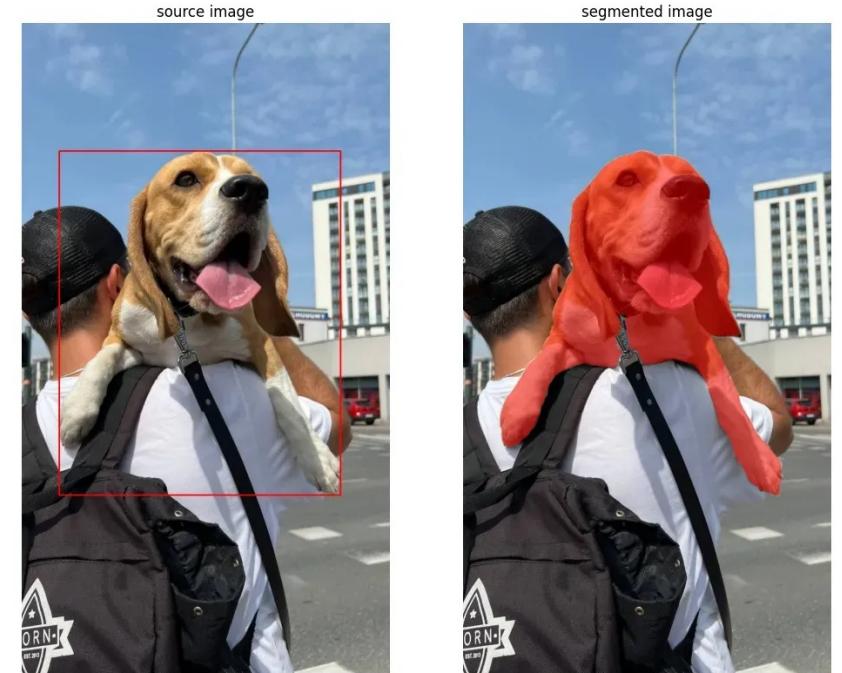


Baselines are ResNet



Segment Anything Model (SAM)

- SAM is a foundation model for **promptable segmentation**, developed by Meta, that can segment any object in an image without task-specific training.
- Supported Prompts
 - Point (foreground / background)
 - Bounding box
 - Partial mask
 - (Optional) Text (via extensions)
- SAM variations:
 - SAM-1 segments anything in **a single image**.
 - SAM-2 extends this to **videos** with memory.
 - “SAM-3” is not an official model
 - only a placeholder for ongoing research beyond SAM
- Architecture (High Level)
 - Image Encoder (ViT): extracts visual features
 - Prompt Encoder: encodes user intent
 - Mask Decoder: generates object masks



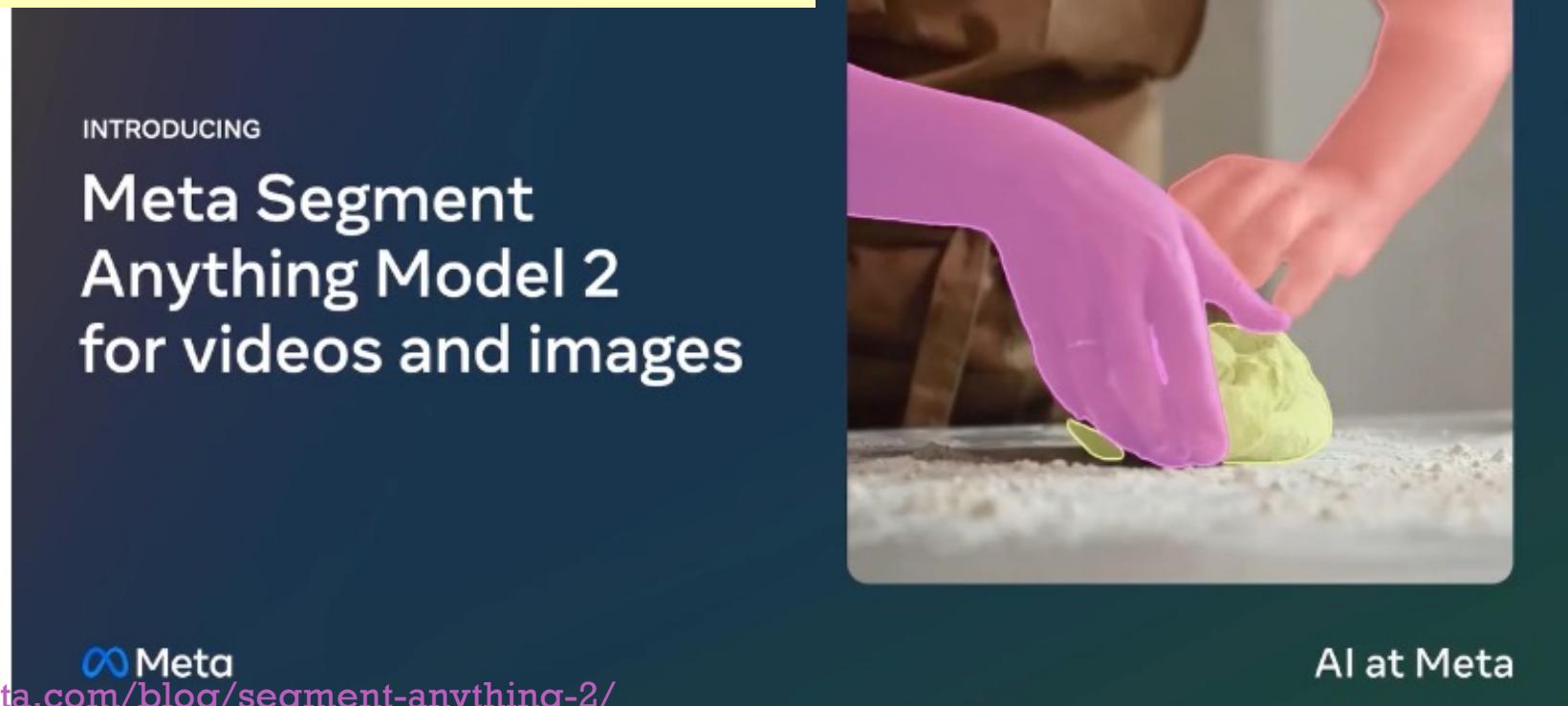
Introducing SAM 2: The next generation of Meta Segment Anything Model for videos and images

July 29, 2024 • ① 15 minute read

SAM-1 segments anything in a single image.

SAM-2 extends this to videos with memory.

“SAM-3” is not an official model—only a placeholder for ongoing research beyond SAM



INTRODUCING

Meta Segment Anything Model 2 for videos and images

AI at Meta

<https://ai.meta.com/blog/segment-anything-2/>

+ Thank you
& any questions