SINGULARITY
TECH DAY_2021

The era of AI and Cognitive Services

# ¿Es este el final de las redes neuronales convolucionales?

ORGANIZATION

SPONSORS

**THANK YOU!**

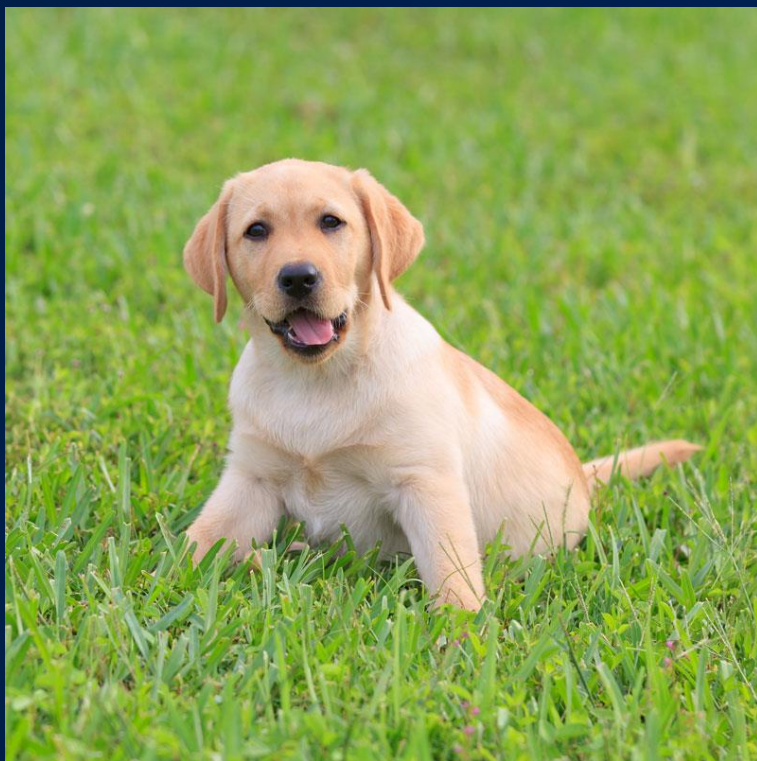"Convolutional neural networks lacks a global understanding of the images. It only looks for the presence of the image's features and does not understand the structural dependency between its features."
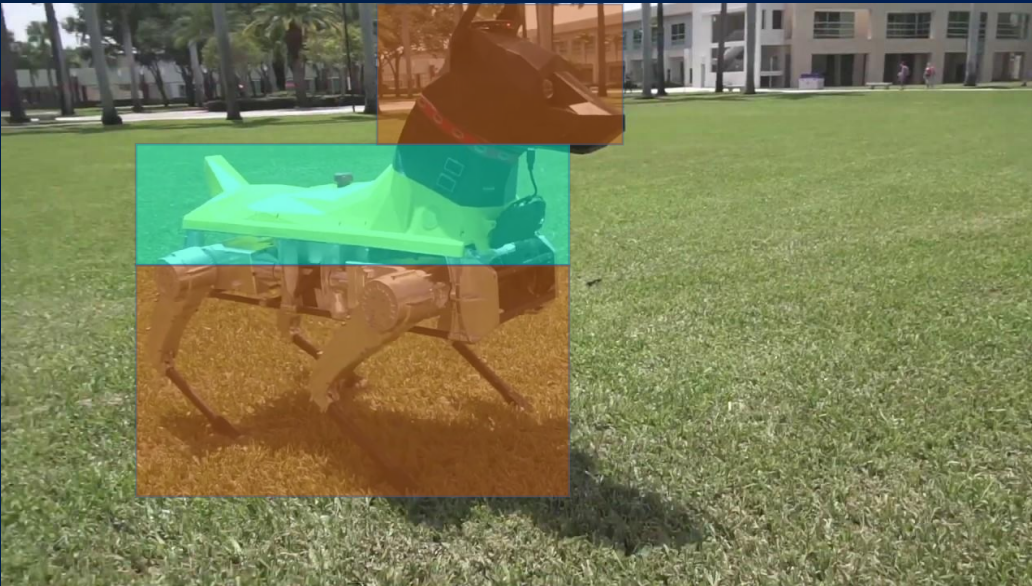
*Capsule networks 2017*

# Computer vision



```
"captions": [
        {
                "text": "a large brown dog lying on green grass",
                "confidence": 0.8984215667012988
        }
    ]

"objects": [
        {
                "object": "golden retriever",
                "confidence": 0.609
        }
    ]
```
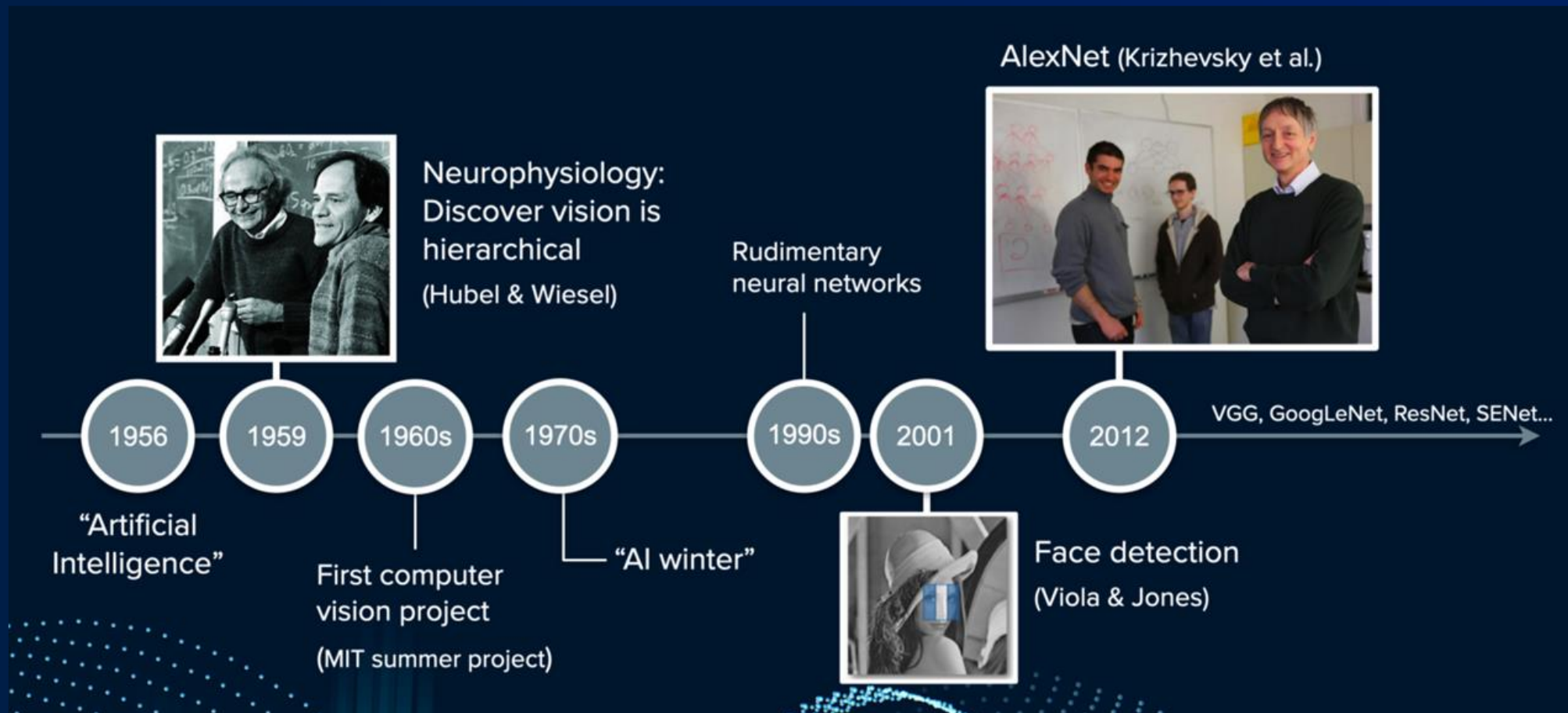
# Computer vision

"captions": [
        {
            "text": "a dog sitting on top of a grass covered field",
            "confidence": 0.7082839064777712
        }
    ]

    "objects": [
        {
            "object": "aircraft",
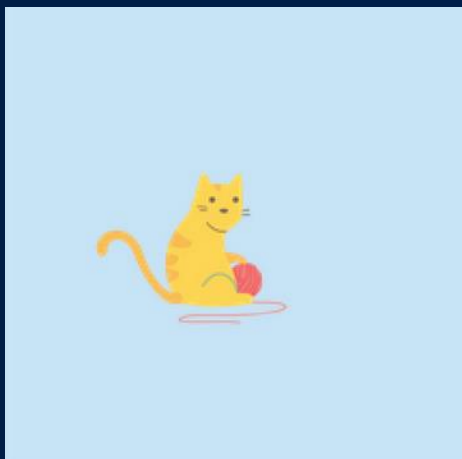            "confidence": 0.57,
        }
    ]

# Computer vision

# Computer vision timeline



AlexNet (Krizhevsky et al.)

Neurophysiology: Discover vision is hierarchical (Hubel & Wiesel)

Rudimentary neural networks

| 1956 | 1959 | 1960s | 1970s | 1990s | 2001 | 2012 |

VGG, GoogLeNet, ResNet, SENet...

"Artificial Intelligence"

First computer vision project

(MIT summer project)
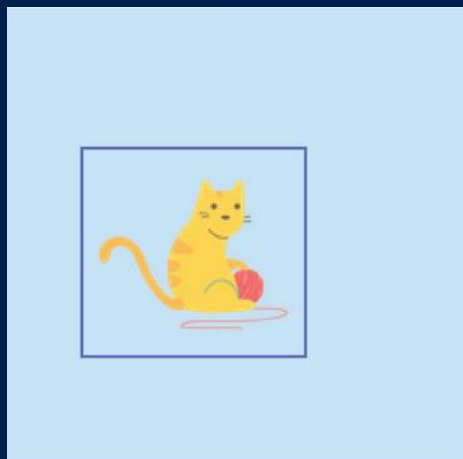
"AI winter"

Face detection (Viola & Jones)

# Computer vision use cases



Classification
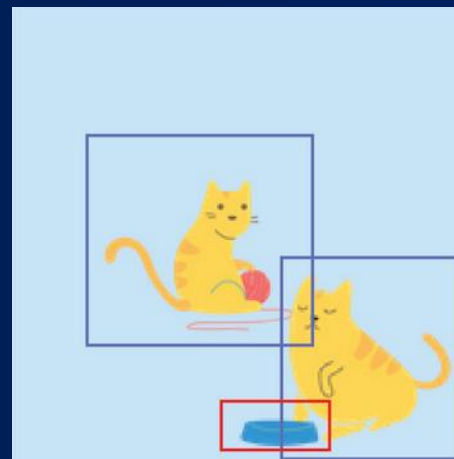
CAT

Classification + Localization

CAT + Bounding Box

Object detection

CAT, CAT, BOWL

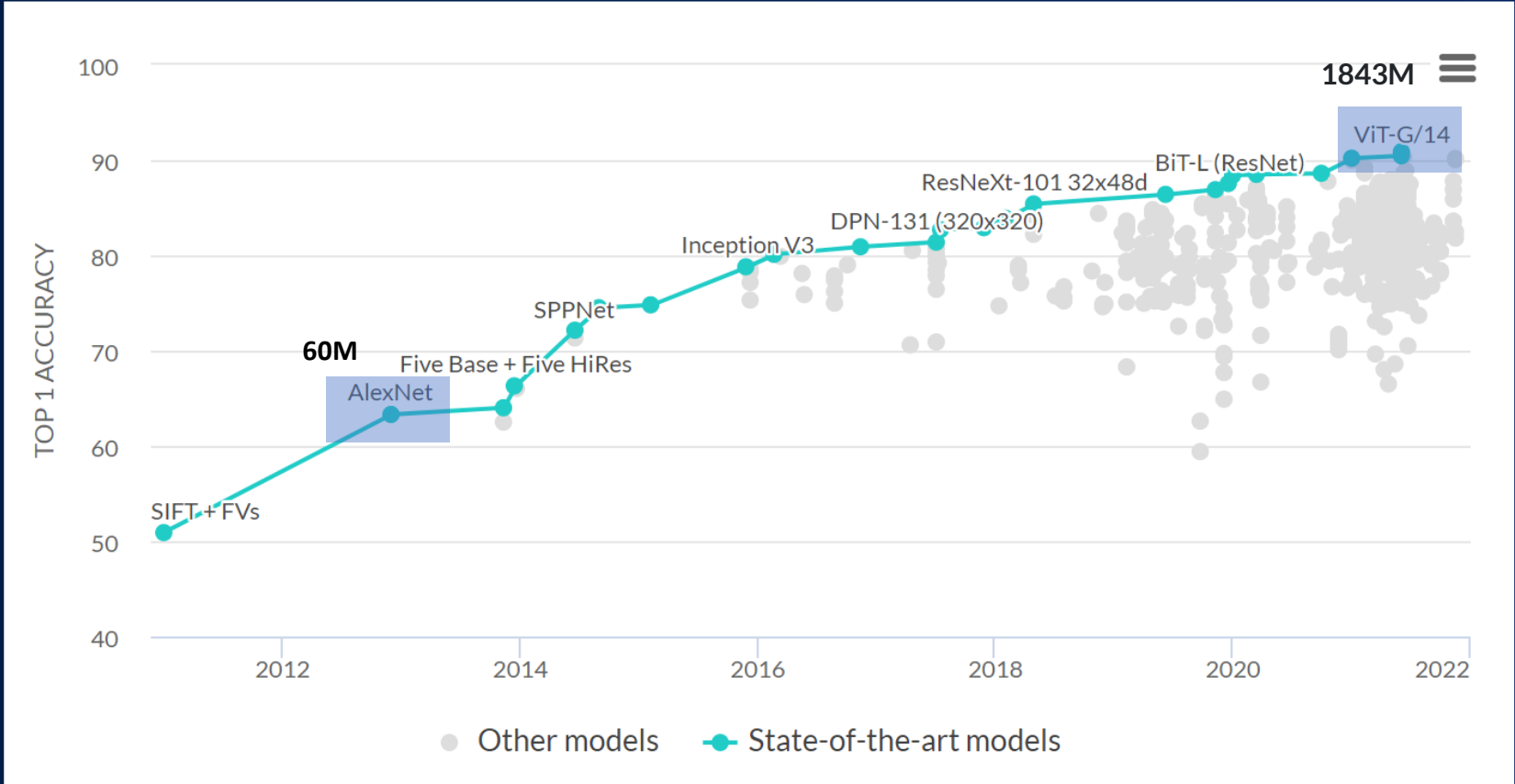Semantic segmentation

CAT, CAT, BOWL

Computer vision neural networks

ImageNet Benchmark

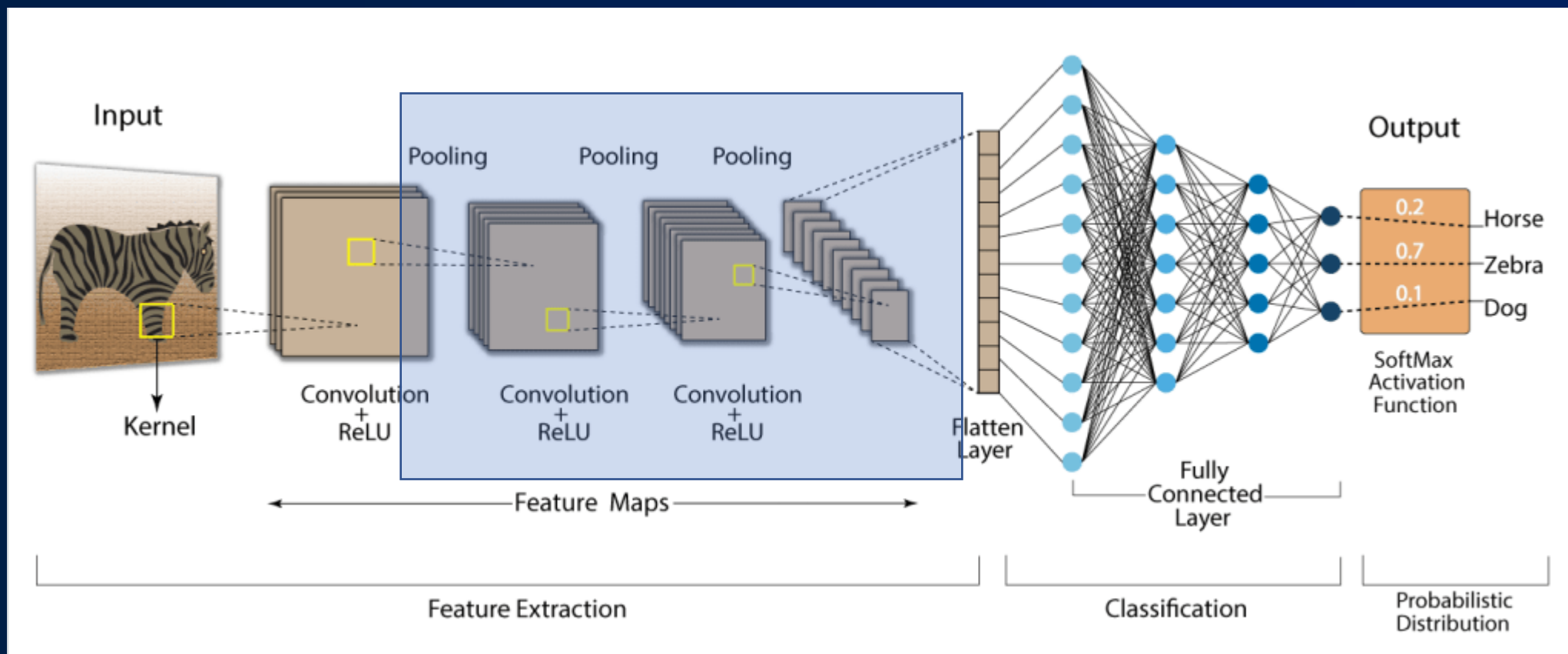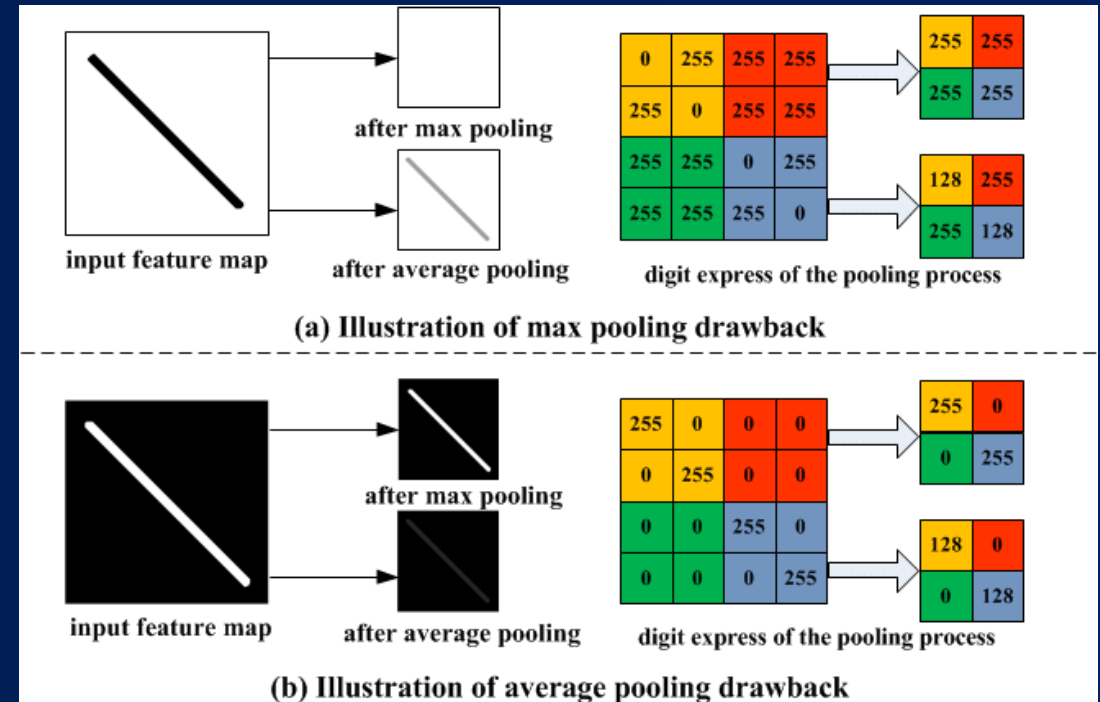# Convolutional neural networks

# CNN topology

# Problems with CNNs

- CNNs detect certain features in an image, but through using a pooling layer valuable information gets lost.

- CNN use "pooling" or equivalent methods to "summarize" what's going on in the smaller regions and make sense of larger and larger chunks of the image. This was a solution that made CNNs work well, but it loses valuable information.
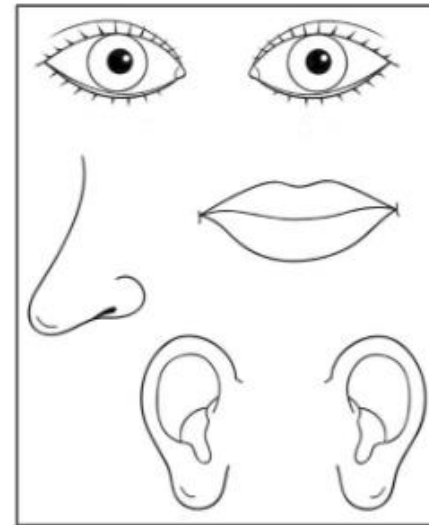


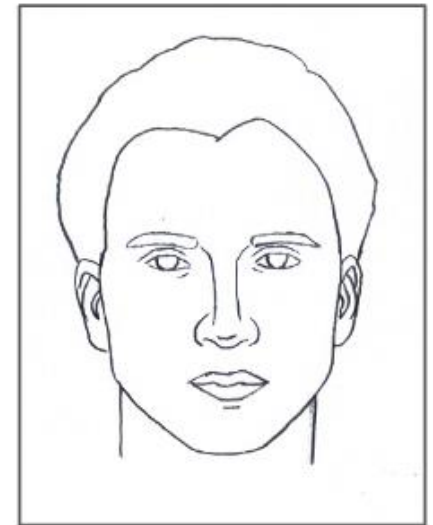(a) Illustration of max pooling drawback

(b) Illustration of average pooling drawback

# Problems with CNNs

- **Is not scale and rotation invariant.**
  - Data augmentation

- **Translation invariance**
  - Ignores the relation between the part and the whole.

- **Receptive field**
  - Model long range dependencies



Not Face                    Face

# Vision Transformers

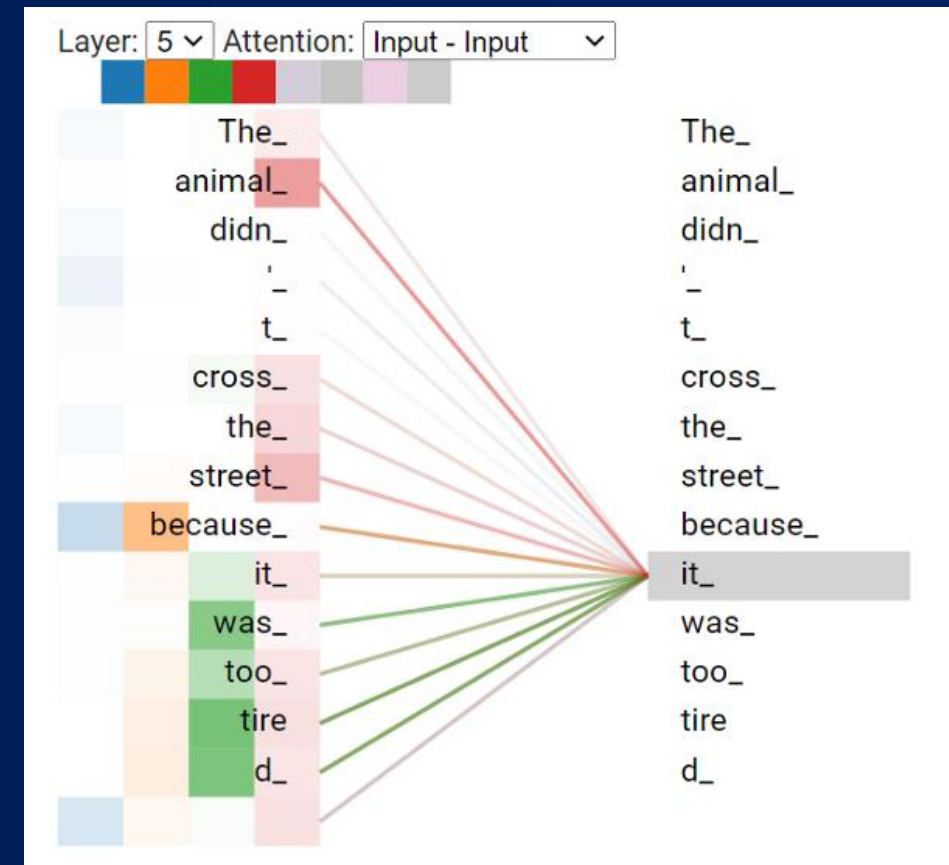# Transformers

- It is used primarily in the field of natural language processing (NLP).

- Transformers are designed to handle sequential input data. The attention mechanism provides context for any position in the input sequence.

- GPT-3 has a transformers-based architecture. It solves tasks such as:
  - Translation, text summarization, semantic search, questions answering, document classification

# Transformers (self-attention)

- Attention are trainable weights that model the importance of each part of an input sentence.

- It will look at each word of the sentence and compare its position in the sentence with respect to the position of all the words present in the same sentence (including itself).

- A score is calculated based on these positional clues which is then used to encode the semantics or meaning of the sentence in a better way.
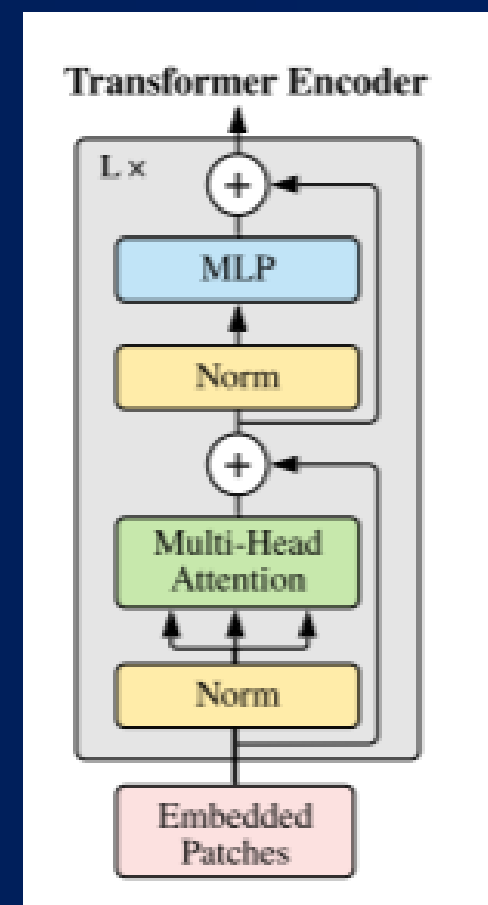
# Vision transformers

- Vision transformer uses pixels to achieve a similar result for images.

- The image is divided into small patches. After that, all patches are flattened using a linear projection.

- A classification head is attached at the end of the transformer encoder to predict the final classes.
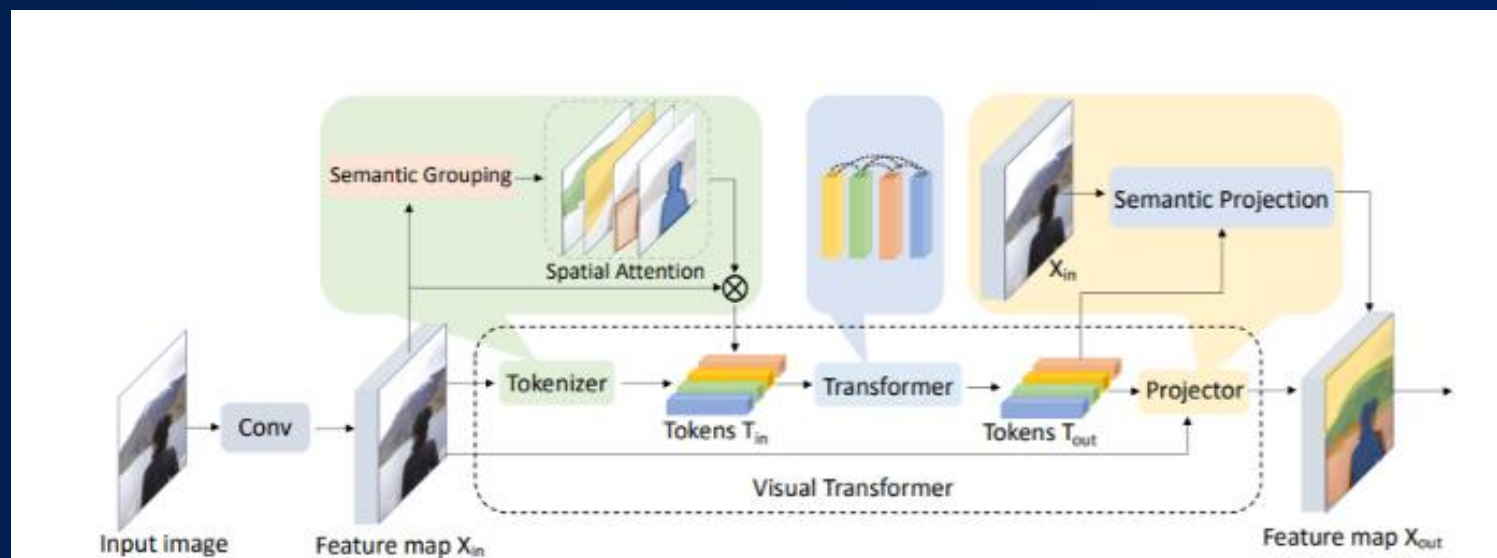
# Transformer encoder

- **Multi-Head Self Attention Layer(MSP)** to concatenate the multiple attention outputs linearly to expected dimensions.

- **Multi-Layer Perceptrons(MLP)** contains two-layer with Gaussian Error Linear Unit(GELU).

- **Layer Norm(LN)** is applied before every block as it does not introduce any new dependencies between the training images

- **Residual connections** are applied after every block as they allow the gradients to flow through the network directly without passing through non-linear activations.
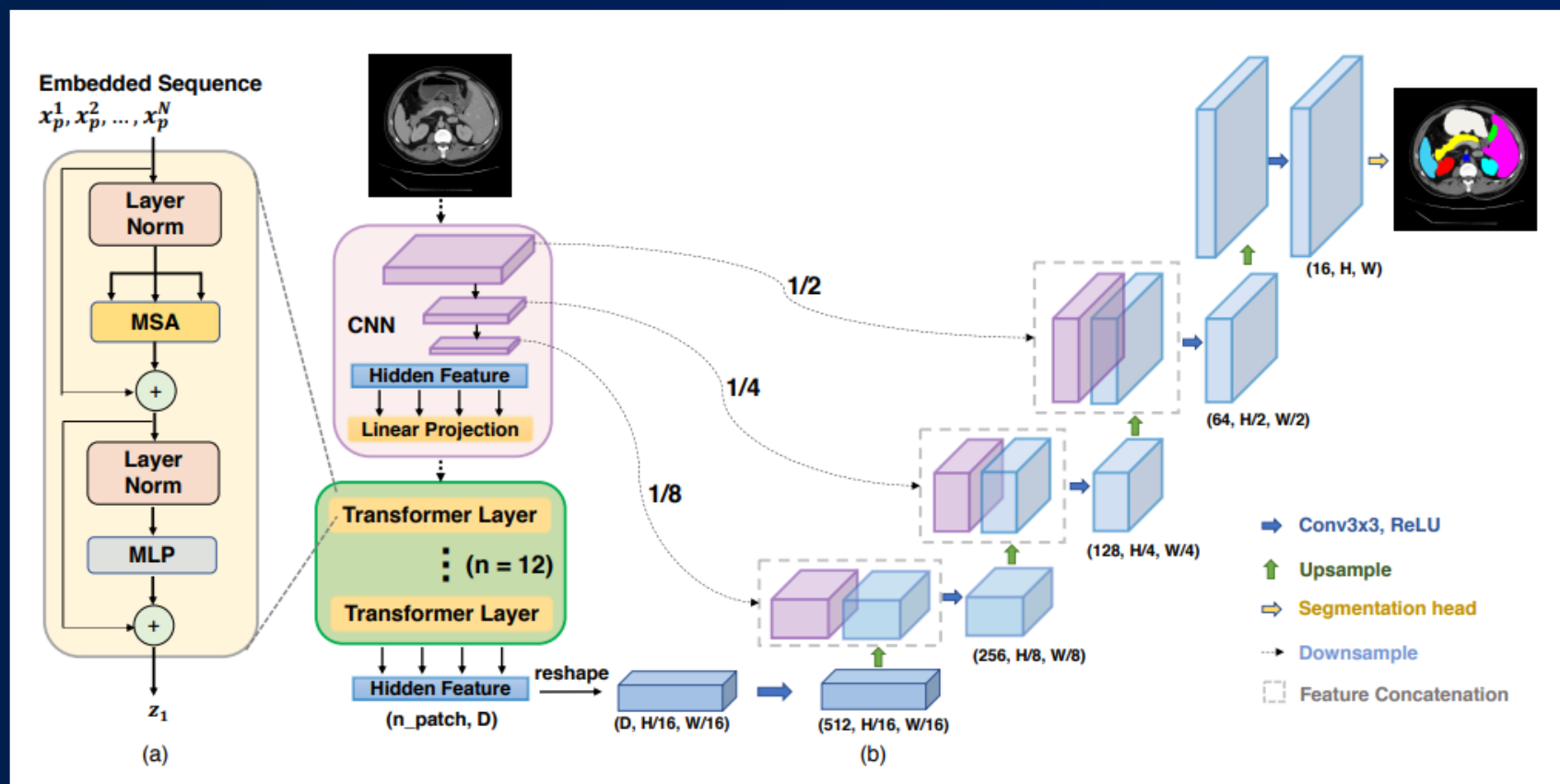
# Transformer architectures
## Hybrid Architectures

- Use CNN as feature extraction.

- Input sequence based on feature maps, then followed by applying the encoding to the feature patches.
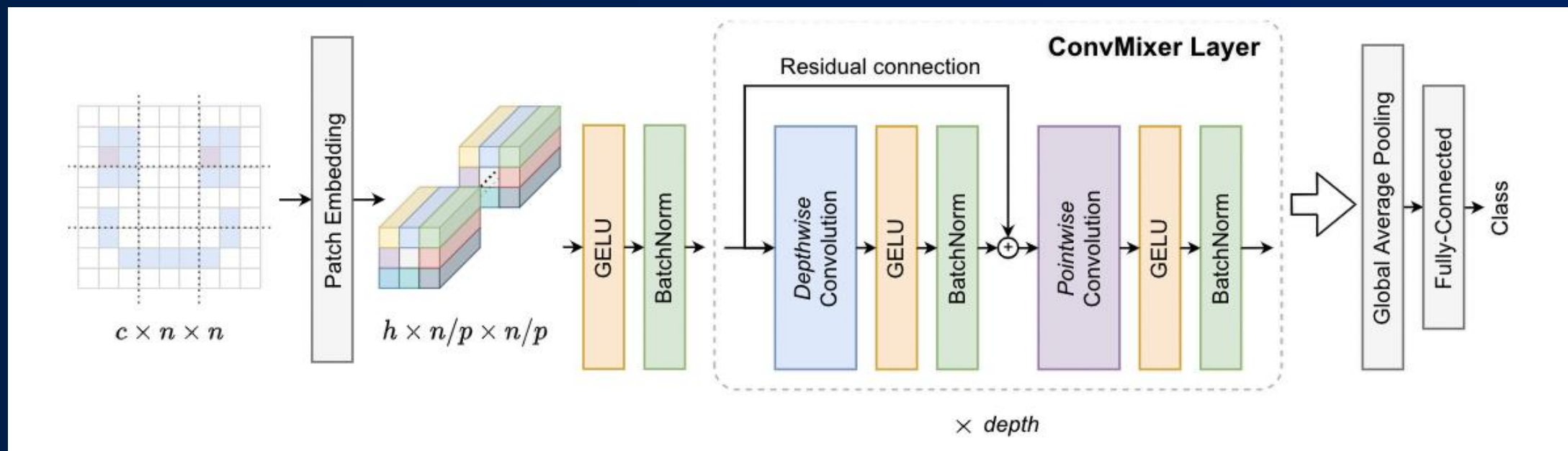
# TransUnet

- CNN-Transformer Hybrid as Encoder.

- Cascaded Upsampler.

# Demo

# What is next?

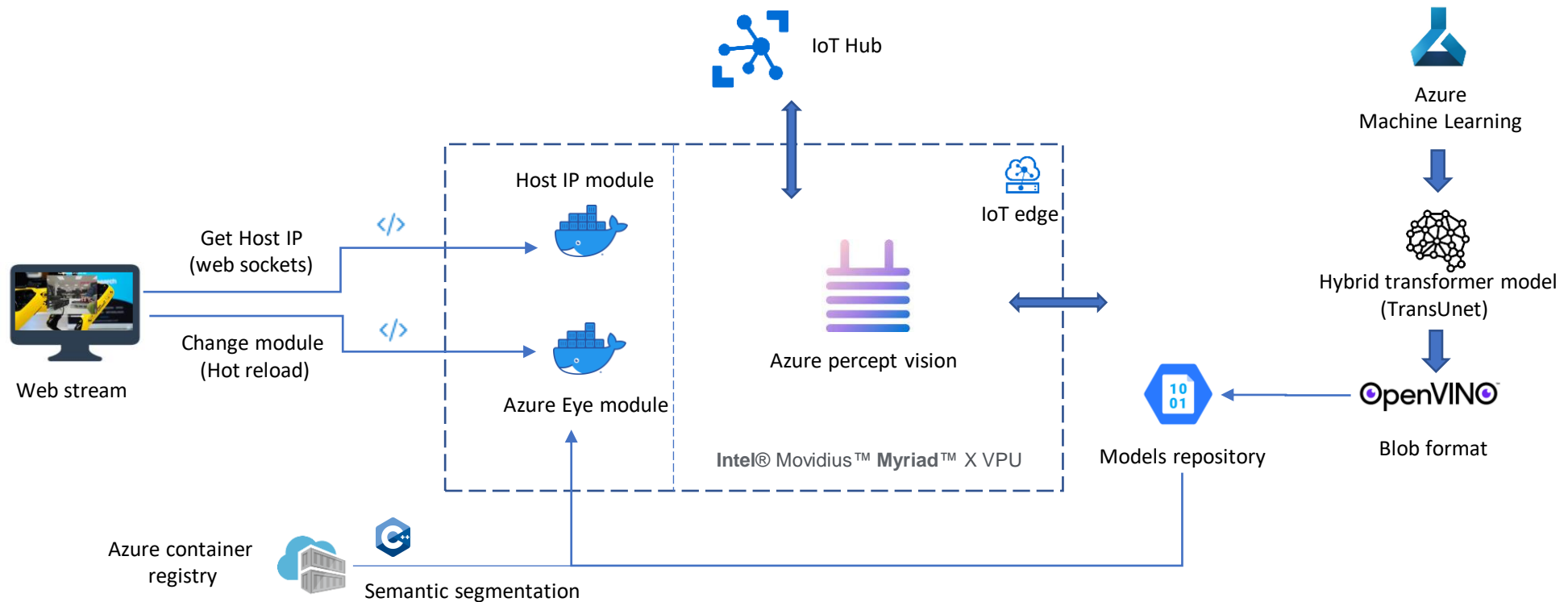**ConvMixers:** Patches are all you need?



Patch representation          Isotropically repeated convolutional blocks.

# Demo

# Architecture

IoT Hub

Azure Machine Learning

Host IP module

IoT edge

Get Host IP (web sockets)

Change module (Hot reload)

Azure percept vision

Hybrid transformer model (TransUnet)

Web stream

Azure Eye module

Intel® Movidius™ Myriad™ X VPU

Models repository

OpenVINO

Blob format

Azure container registry

Semantic segmentation

[Webstream Video (windows.net)](#)