# Learning Hierarchical Image Segmentation For Recognition and By Recognition

Tsung-Wei Ke*          Sangwoo Mo*          Stella Yu

*Equal Contribution
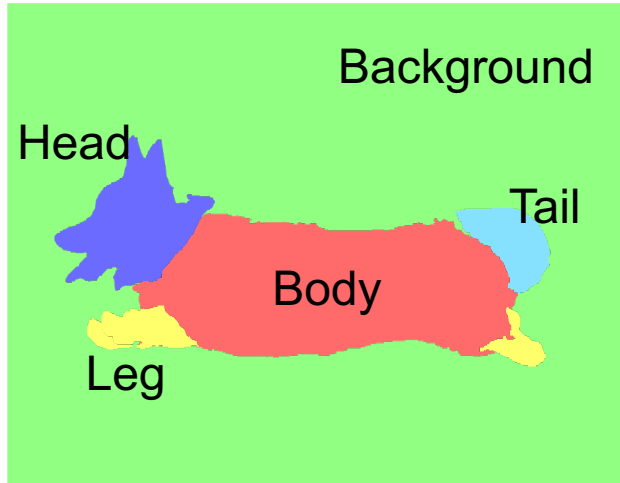
# Image understanding happens at multiple levels: Image recognition vs. Pixel segmentation



Image

Part-level segmentation

Background

Head

Tail

Body

Leg

Object-level segmentation
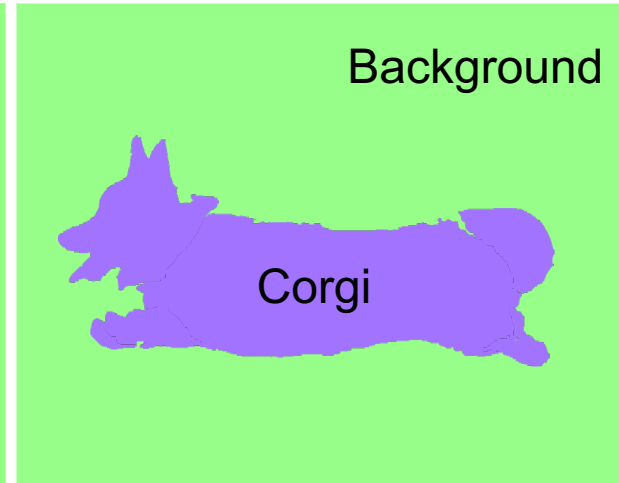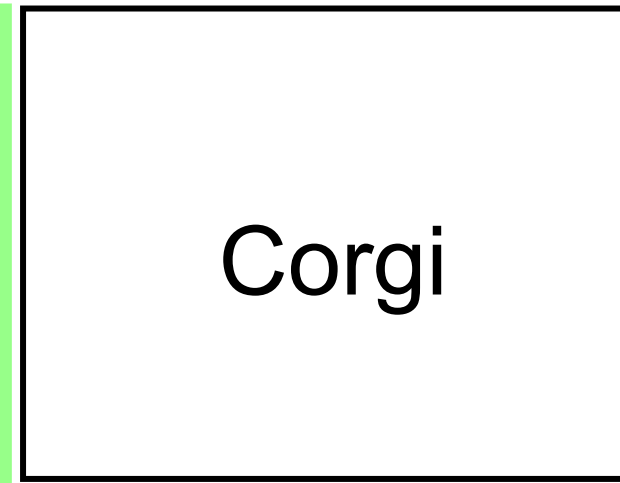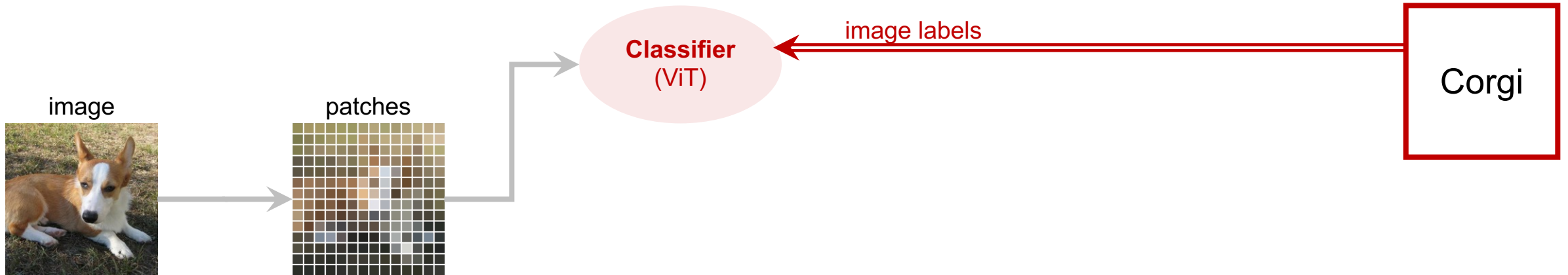
Background

Corgi

Image-level recognition

Corgi

# Prior works: Build separate models and supervision



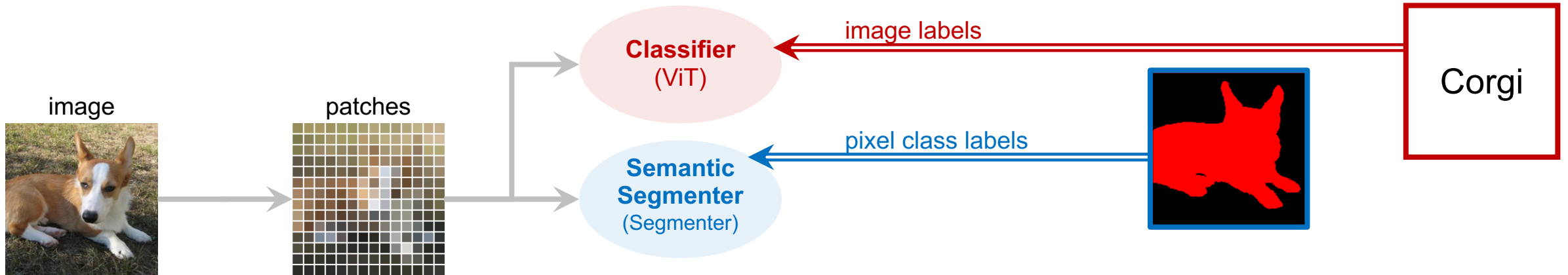image

patches

**Classifier**
(ViT)

image labels

Corgi

Imagenet: A large-scale hierarchical image database. Deng et al. *CVPR 2009*

Sun database: Large-scale scene recognition from abbey to zoo. Xiao et al. CVPR 2010

Laion-5b: An open large-scale dataset for training next generation image-text models. Schuhmann et al. *NeuRIPS 2022*

# Prior works: Build separate models and supervision

image



patches



**Classifier**
(ViT)

**Semantic Segmenter**
(Segmenter)

image labels

pixel class labels



Corgi



CSAILVision/
**ADE20K**    MIT

**CITYSCAPES**
DATASET
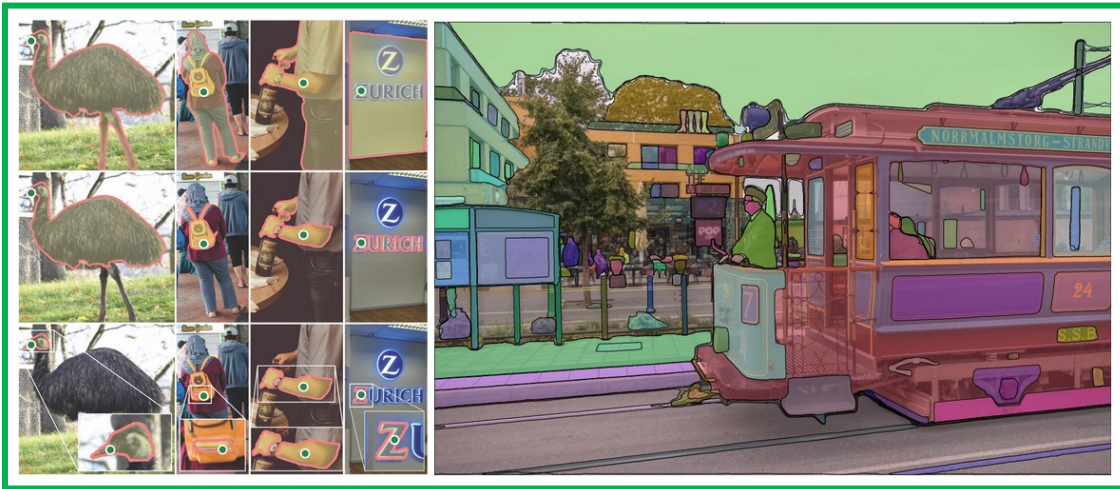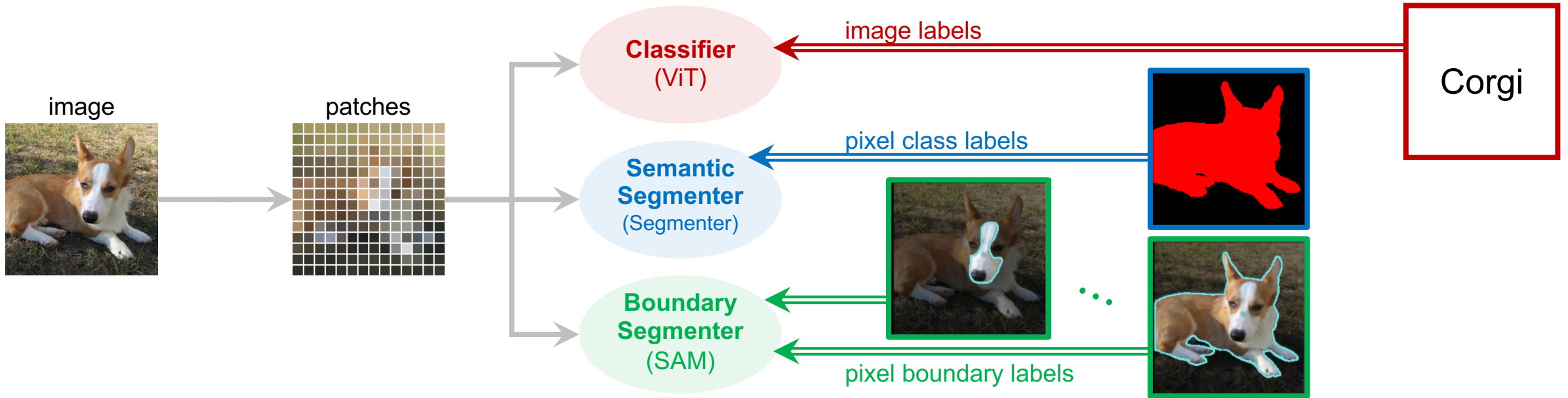
COCO
Common Objects in Context

Microsoft coco: Common objects in context. Lin, et al. ECCV 2014

The cityscapes dataset for semantic urban scene understanding. Cordts et al. CVPR 2016

Scene parsing through ade20k dataset. Zhou et al. CVPR 2017
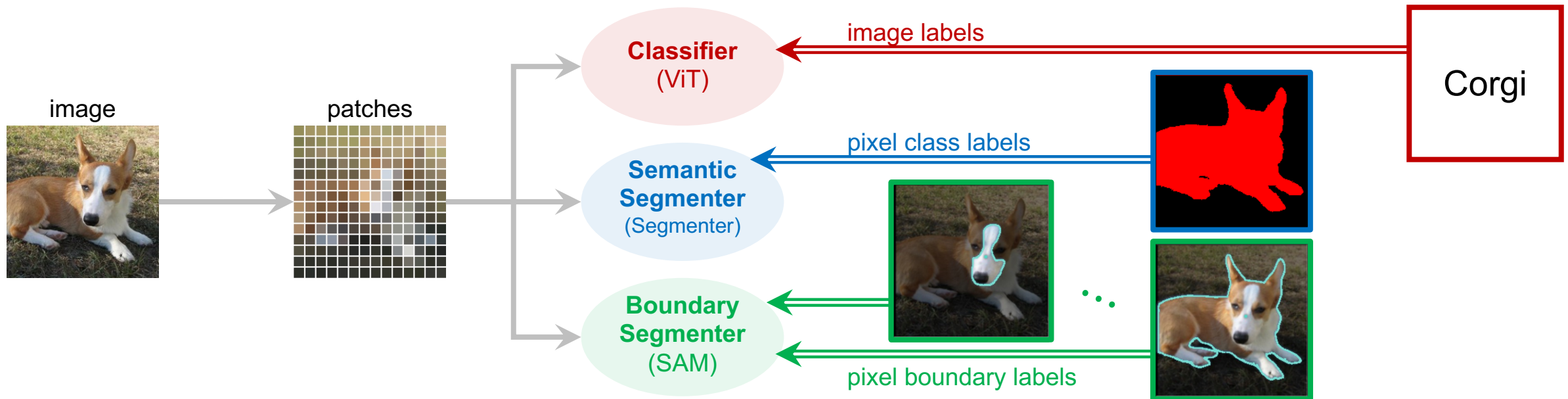
# Prior works: Build separate models and supervision



Segment Anything. Kirillov et al. ICCV 2023

# Problems of prior works:

1. Need separate annotations
2. Need separate models
3. One task does not help the other

# Our idea: put segmentation in the loop of recognition



image

superpixels

**Segmenter for Recognition**

(CAST)

Corgi

# Strengths of our idea

1. Learn hierarchical segmentation FOR FREE from image labels

# Strengths of our idea

1. Hierarchical segmentation learned FOR FREE from image labels
2. Adaptive segmentation improved with image recognition

# Strengths of our idea

1. Hierarchical segmentation learned FOR FREE from image labels
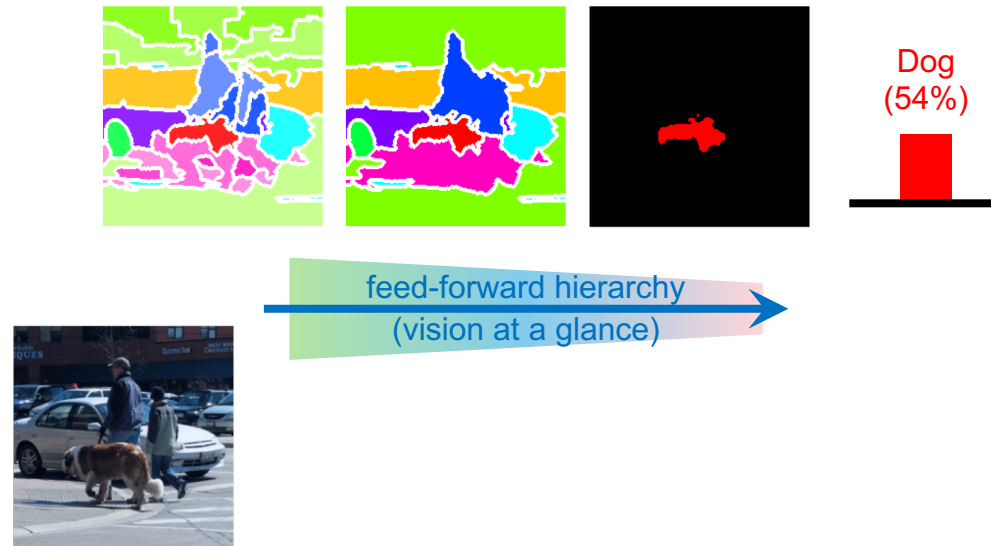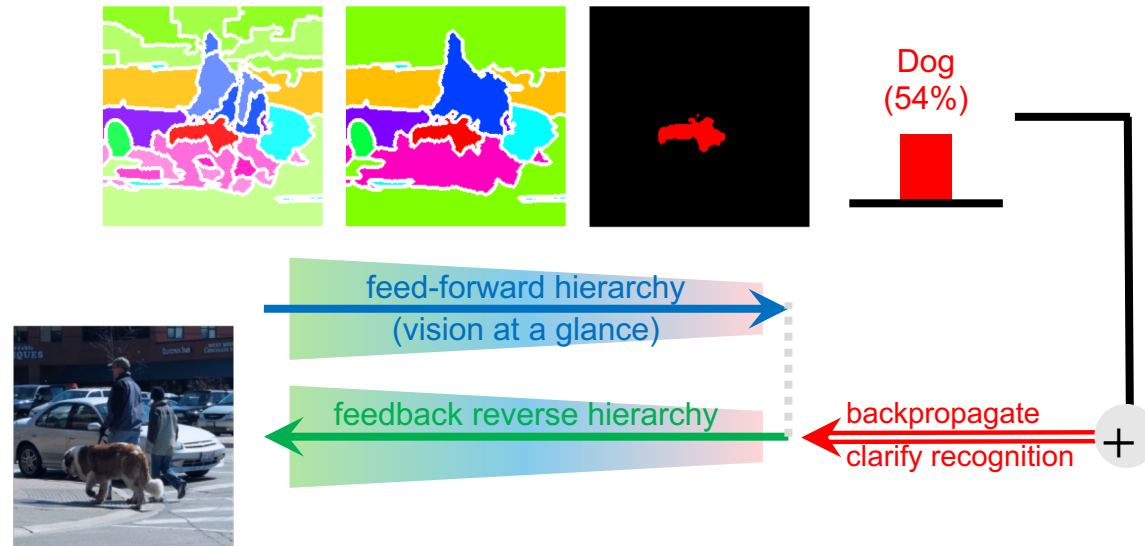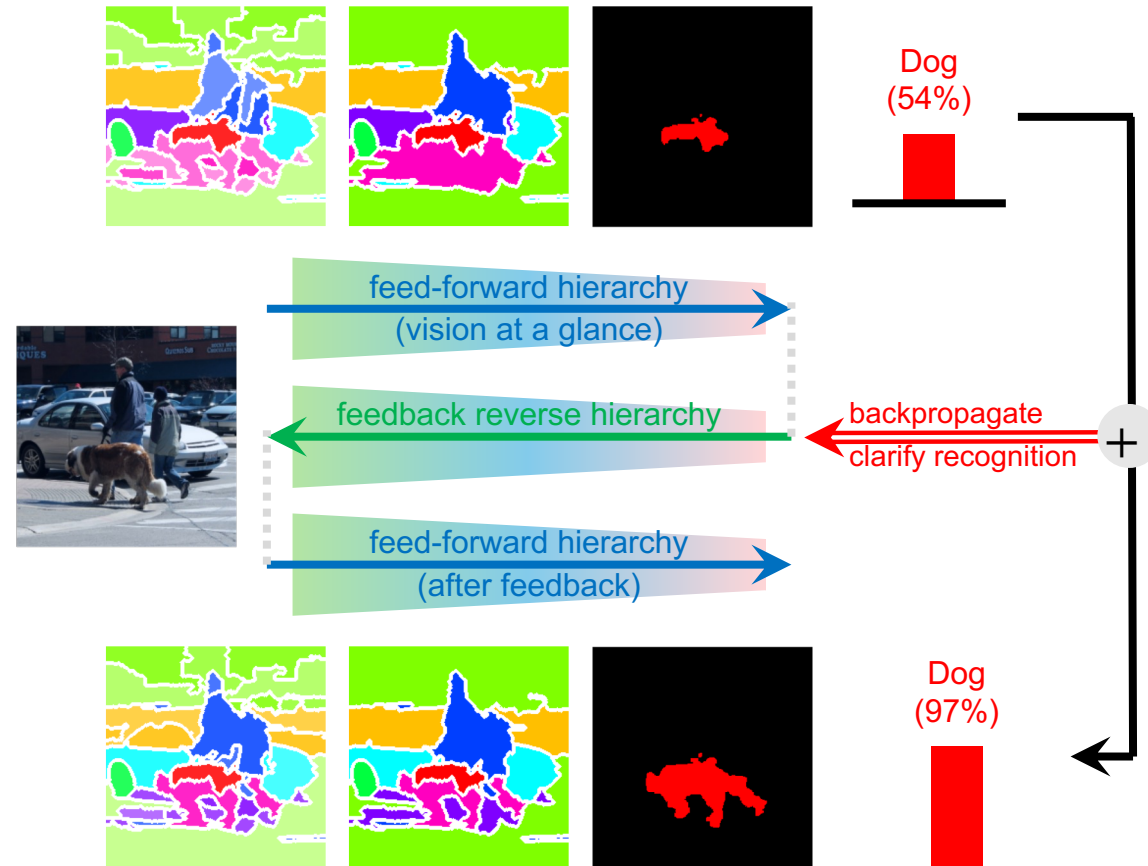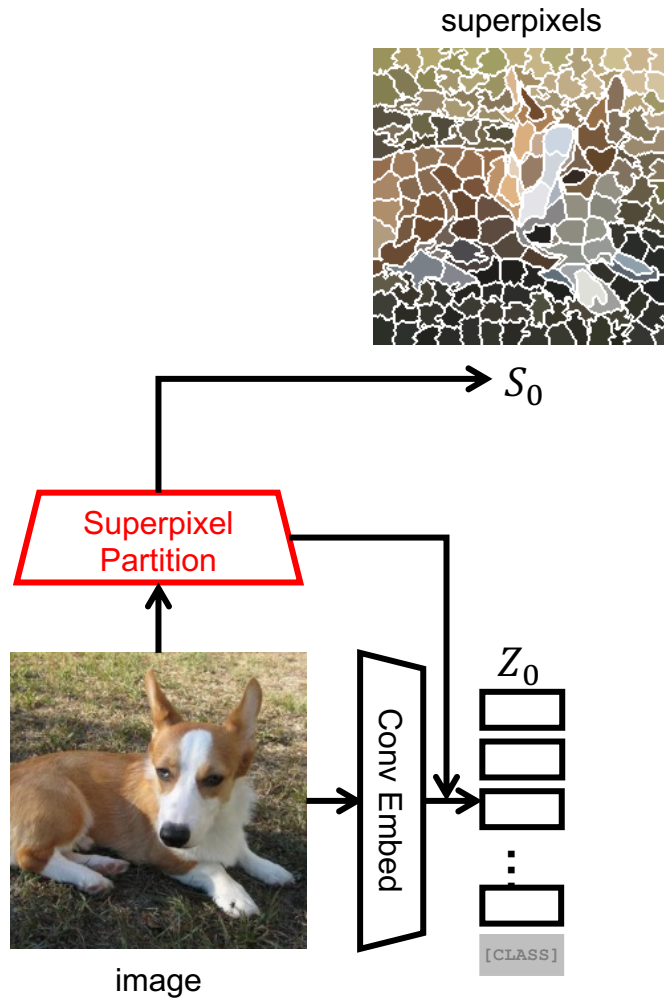2. Adaptive segmentation improved with image recognition

# Strengths of our idea

1. Hierarchical segmentation learned FOR FREE from image labels
2. Adaptive segmentation improved with image recognition

# Step 1: begin with segment (superpixel) tokens

superpixels



$S_0$

Superpixel Partition

$Z_0$

Conv Embed

[CLASS]

image

# Step 2: contextualize tokens with transformer encoder

superpixels

# Step 3: group fine segments to coarse regions



superpixels

fine segmentation

$S_0$    $\otimes$    $S_1$

$P_1$

Superpixel Partition

$Z_0$       $Z_1$

Conv Embed

ViT Block   ViT Block   Graph Pooling

[CLASS]      [CLASS]

image

# Repeat step 2 to 3



superpixels

fine segmentation

coarse segmentation

image

# Step 4: predict image-level recognition



superpixels

fine segmentation

coarse segmentation

**Image Recognition Objective**

Instance discrimination

Image classification

$S_0$ ⊗ $S_1$ ⊗ $S_2$

$P_1$ $P_2$

Superpixel Partition

$Z_0$ $Z_1$ $Z_2$

Conv Embed

ViT Block ViT Block Graph Pooling ViT Block ViT Block Graph Pooling

[CLASS] [CLASS] [CLASS]

image

Corgi parrot spider monkey

# Benefits of CAST

1. Unsupervised discovery of hierarchical segmentations

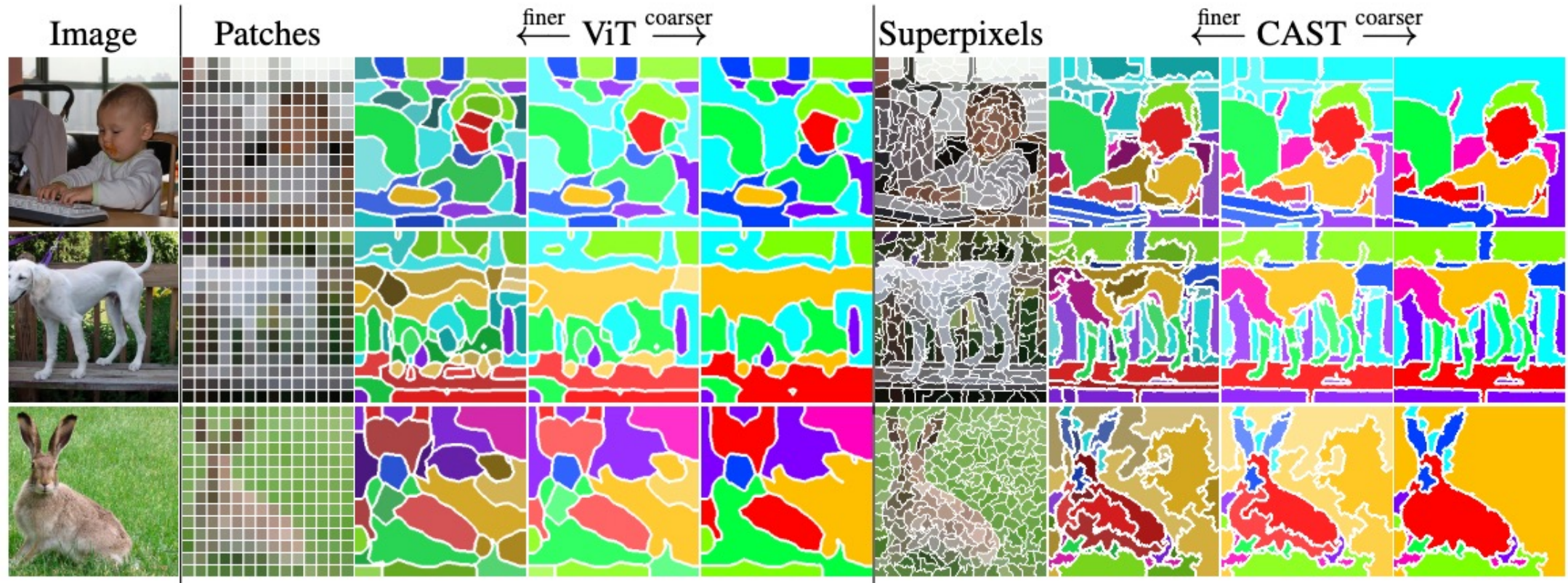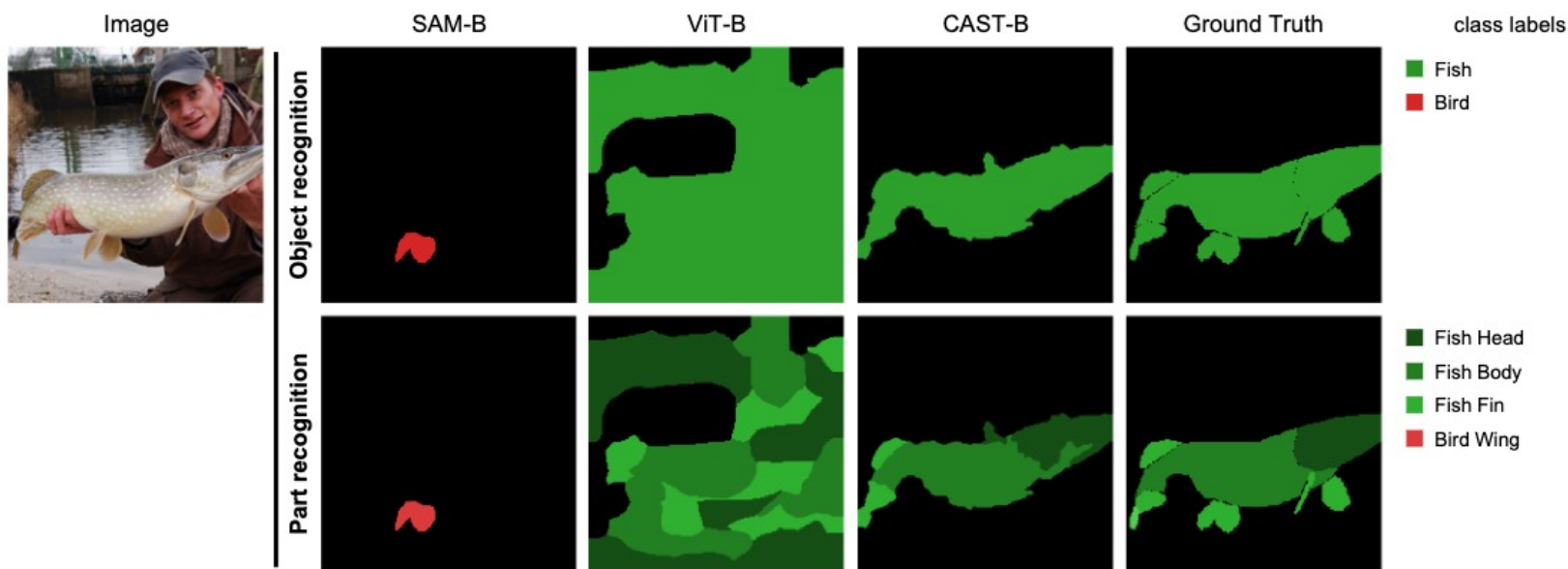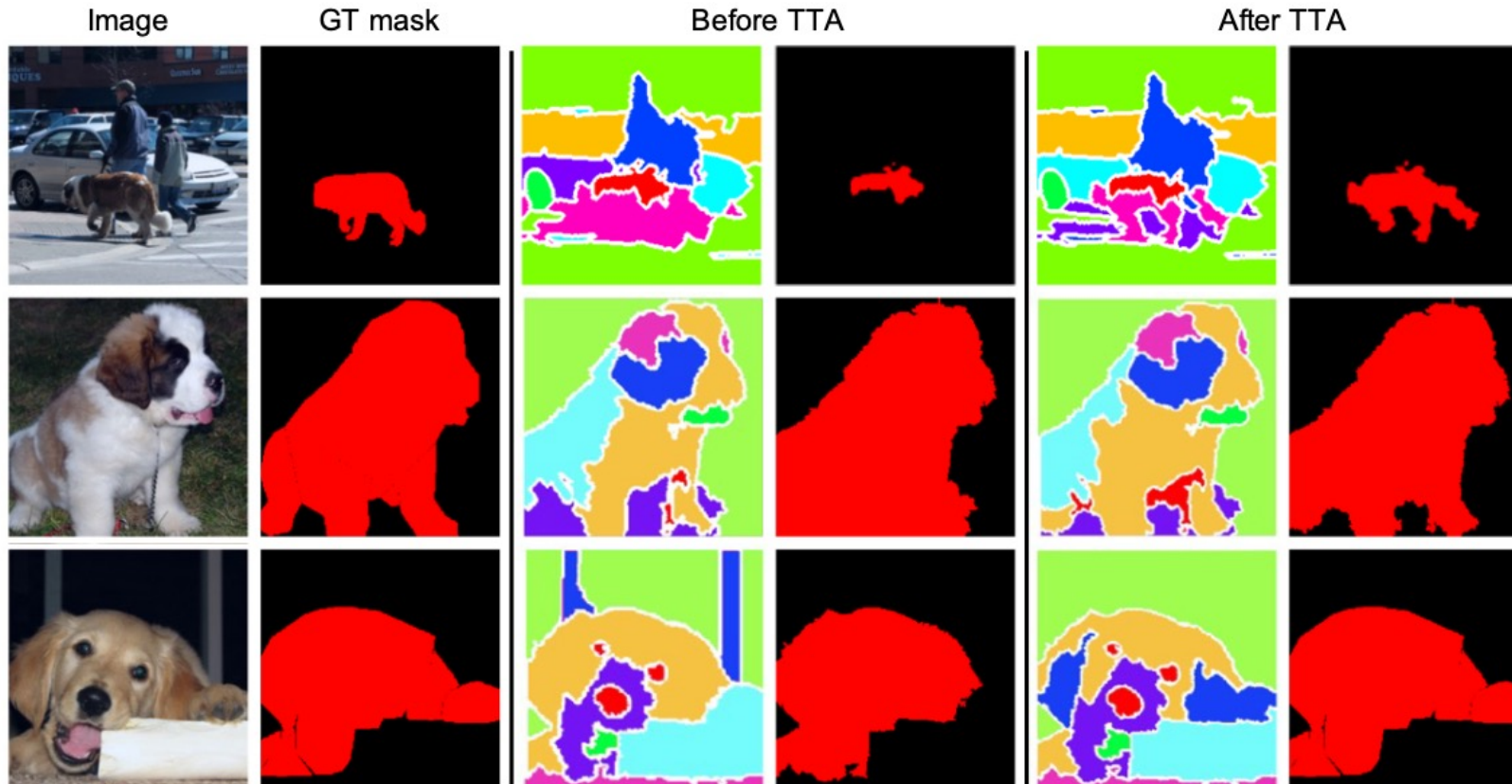# Benefits of CAST

1. Unsupervised discovery of hierarchical segmentations
2. Unsupervised part / object segmentation beats Segment Anything



| Model | Training data | Supervised | GFLOPS | Part $\xleftarrow{\text{finer}}$ Object $\xrightarrow{\text{coarser}}$ Category | | |
|-------|---------------|------------|--------|------|------|------|
| SAM-B | SA-1B | ✓ | 488.2 | 10.15 / **7.25** | 18.03 / 20.71 | 31.36 / 32.01 |
| ViT-B | IN-1K | ✗ | 17.8 | 11.74 / 4.64 | 25.34 / 10.92 | 36.68 / 13.28 |
| CAST-B | IN-1K | ✗ | **12.9** | **13.20** / 6.52 | **29.66 / 22.32** | **50.75 / 34.38** |

# Benefits of CAST

1. Unsupervised discovery of hierarchical segmentations
2. Unsupervised part / object segmentation beats Segment Anything
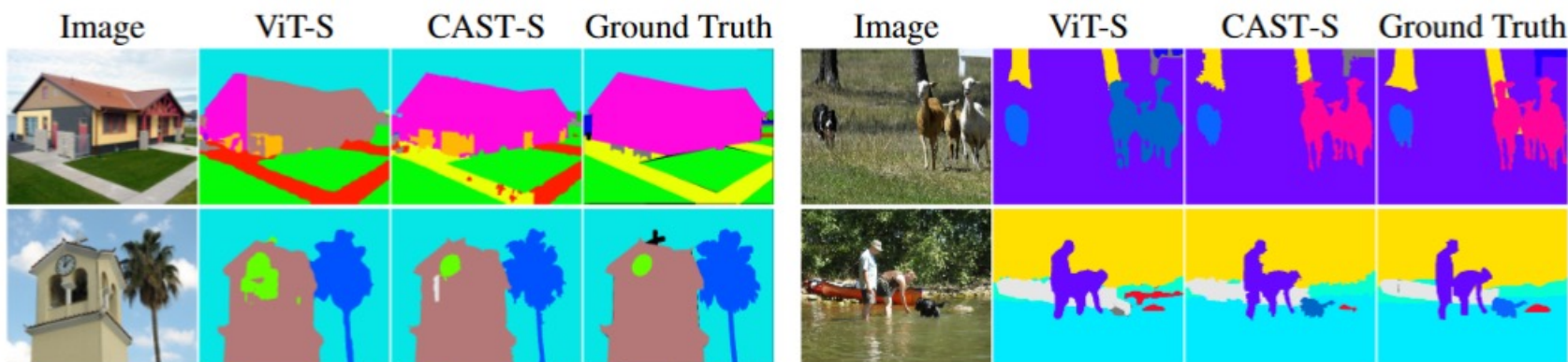3. Concurrent segmentation and recognition during inference

# Benefits of CAST

1. Unsupervised discovery of hierarchical segmentations
2. Unsupervised part / object segmentation beats Segment Anything
3. Concurrent segmentation and recognition during inference
4. <span style="color:red">Better performance and efficiency than ViT</span>



segmentation

| (a) Pascal VOC | Token | Pooling | Before tuning | After tuning |
|---|---|---|---|---|
| ViT-S | Patch | ✗ | 30.9 / 16.1 | 65.8 / 40.7 |
| ↕ ablation | Patch | ✓ | 34.5 / 19.8 | 67.2 / 41.9 |
| | Superpixel | ✗ | 32.2 / 21.2 | 66.5 / 46.7 |
| CAST-S | Superpixel | ✓ | **38.4 / 27.0** | **67.6 / 48.1** |

classification & efficiency

| Model | GFLOPS | IN-100 | IN-1K |
|---|---|---|---|
| ViT-S | 4.7 | 78.1 | 67.9 |
| Swin-T | 4.5 | 78.3 | 63.0 |
| CAST-S | **3.4** | **79.9** | **68.1** |

# Thank you for your listening

Code available at:
https://github.com/twke18/CAST