

Semantic Segmentation

Danna Gurari

University of Colorado Boulder

Fall 2024



Review

- Last lecture:
 - Scene Classification Problem and Applications
 - Scene Classification Datasets and Evaluation Metrics
 - Scene Classification Models: Deep Features
 - Attribute Classification: Problem, Applications, and Datasets
- Assignments (Canvas)
 - Reading assignment was due earlier today
 - Next reading assignment due Monday
 - Project proposal due in one week
- Questions?

Semantic Segmentation: Today's Topics

- Motivation
- Datasets
- Evaluation metric
- Fully convolutional network
- Swin transformer
- Discussion (chosen by YOU 😊)

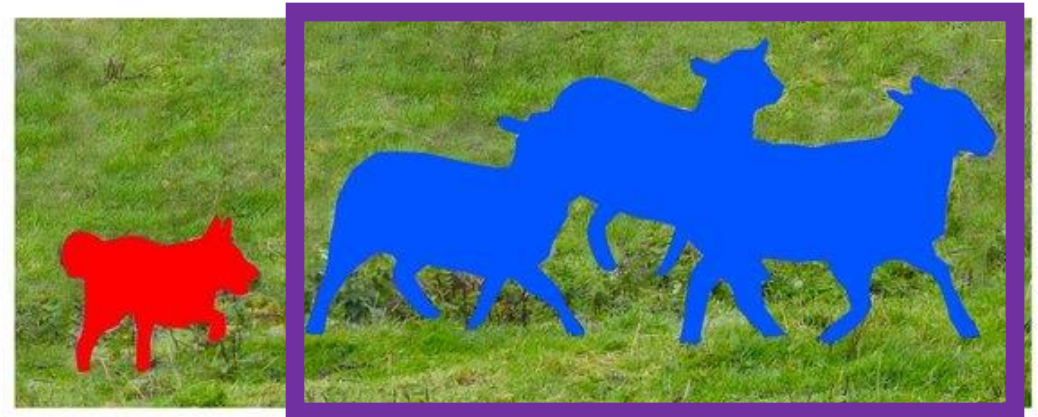
Semantic Segmentation: Today's Topics

- Motivation
- Datasets
- Evaluation metric
- Fully convolutional network
- Swin transformer
- Discussion (chosen by YOU 😊)

Today's Scope: Localize Pixels for Each Category



Image Recognition



Semantic Segmentation

Note: instances of the same category are NOT separated

Remodeling Inspiration

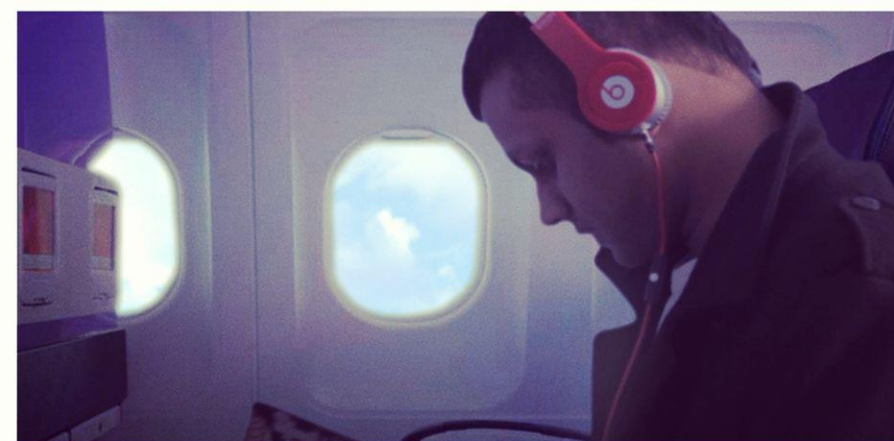
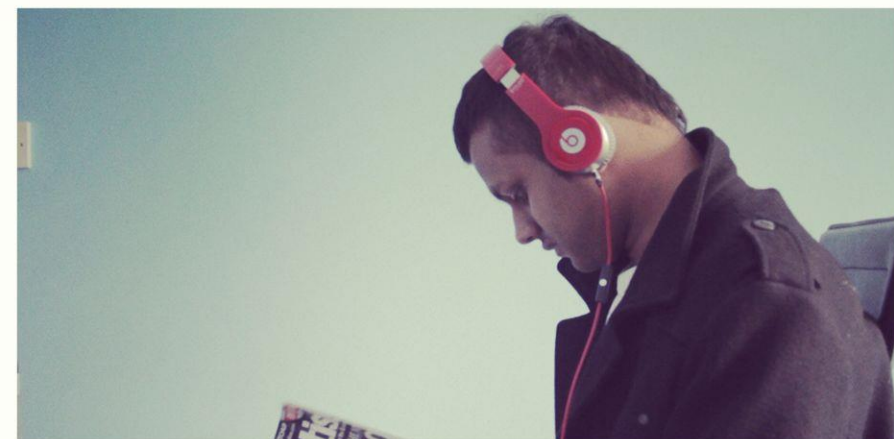
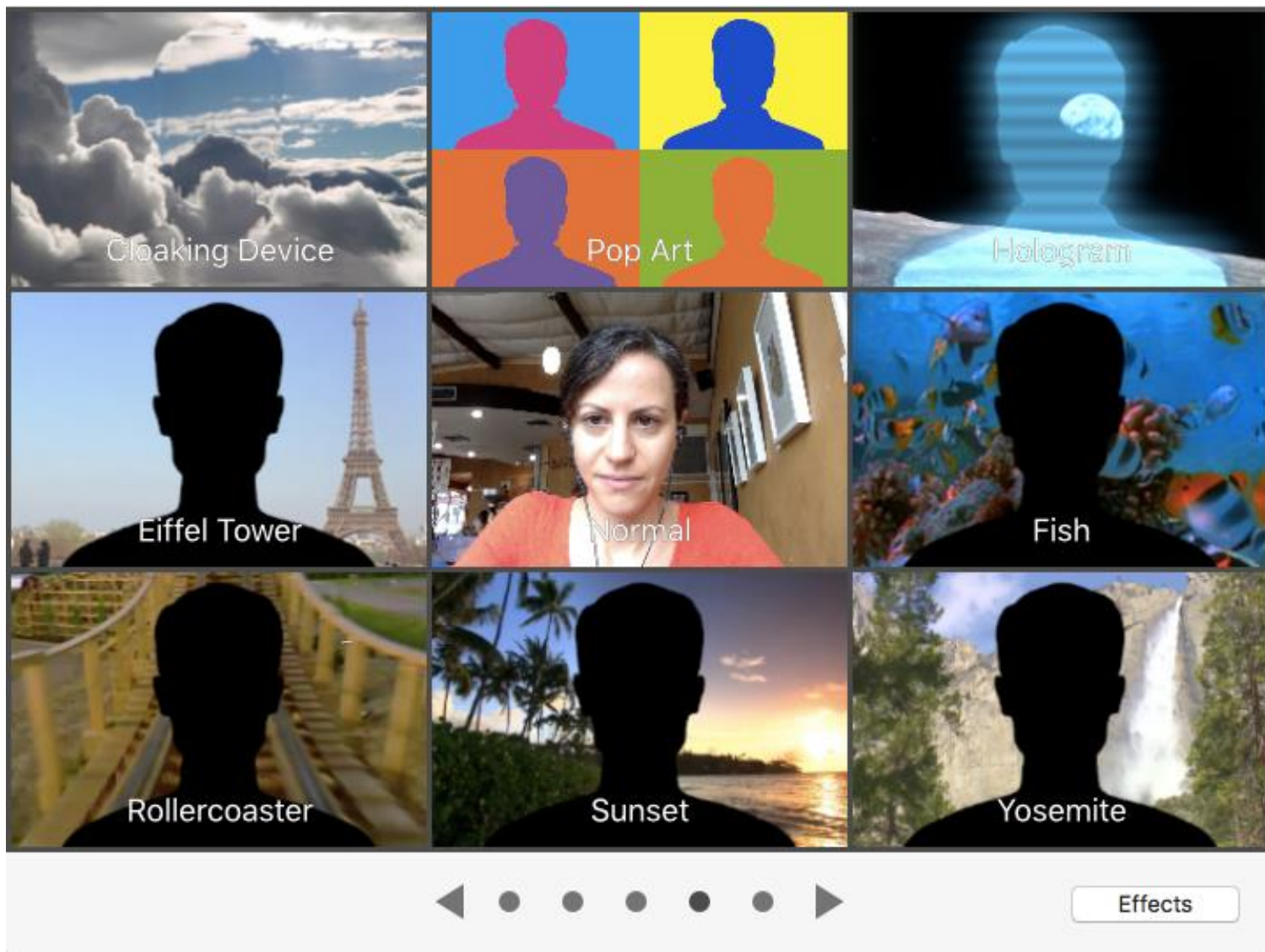


(a) Target photo



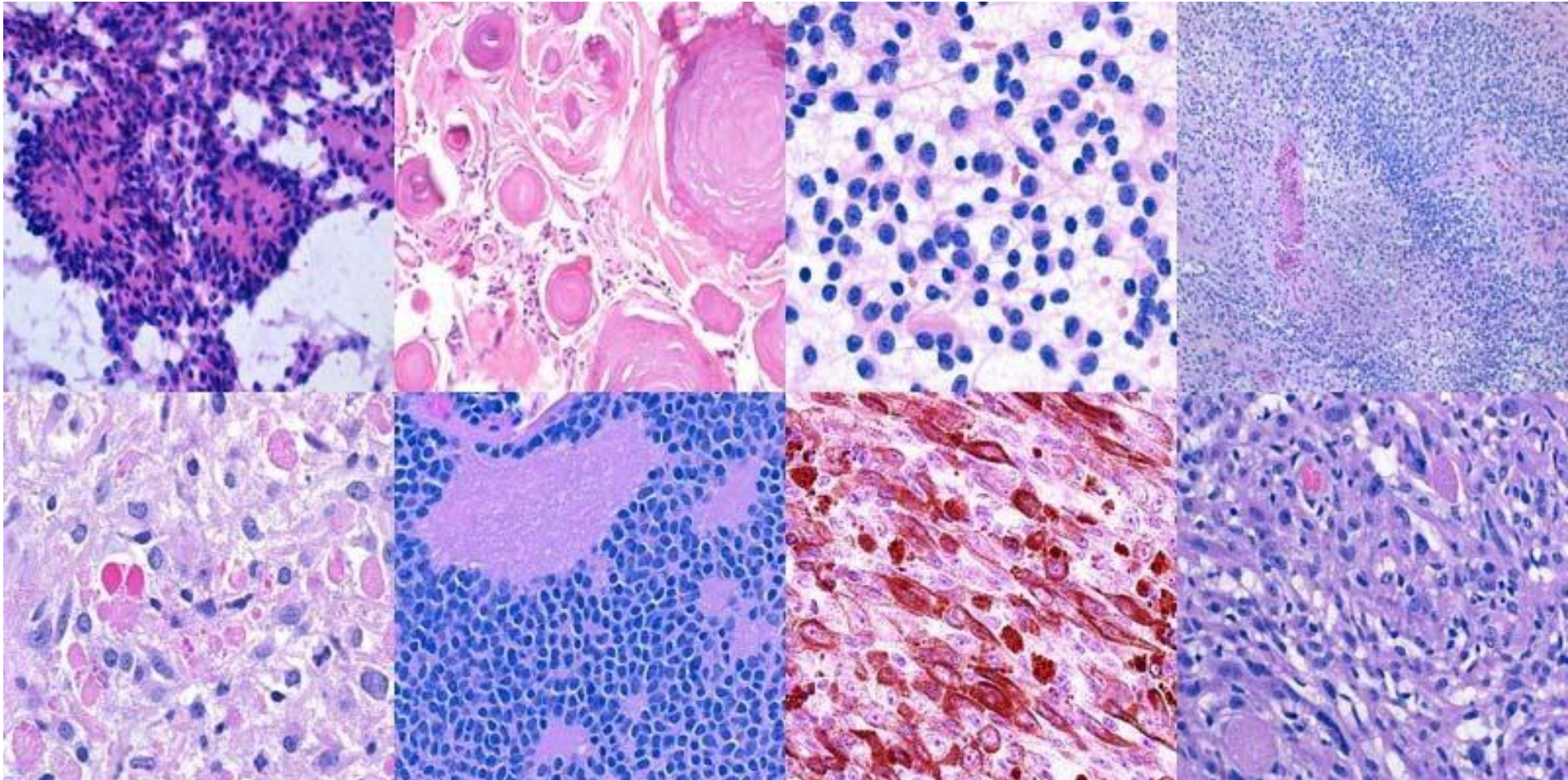
(b) Retextured

Rotoscoping (many examples on Wikipedia)



<https://www.starnow.co.uk/ahmedmohammed1/photos/4650871/before-and-after-rotoscopinggreen-screening>

Disease Diagnosis; e.g., PathAI



Face Makeover

MAYBELLINE
NEW YORK

VIRTUAL BEAUTY STUDIO

SHOP ALL

FACE

EYES

LIPS

NAILS

TIPS & TRENDS

BRAVE TOGETHER

Home

TRY IT ON

Time to makeup your mind! Experience your perfect makeup shades or try a bold new look with Maybelline's virtual try-on tool.

To begin, turn on your camera or upload a photo.

SEE YOURSELF IN MAYBELLINE



GET STARTED!

I Consent

to the processing of my image by Maybelline NY
as set out in the [privacy policy](#).



LIVE CAMERA



UPLOAD PHOTO

Demo: <https://www.maybelline.com/virtual-try-on-makeup-tools>

Self-Driving Vehicles



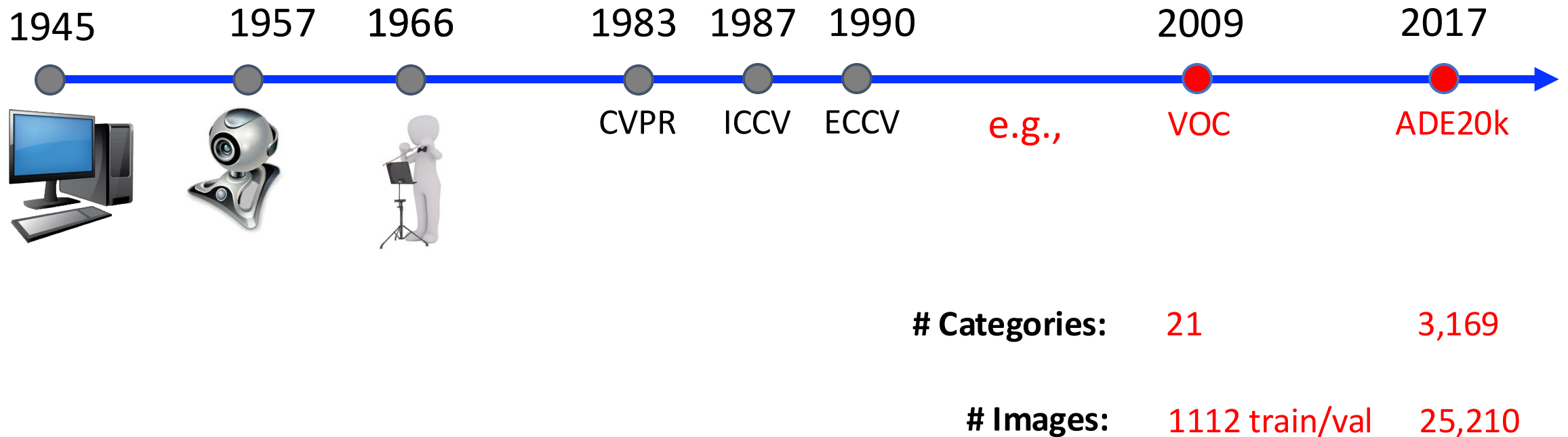
<https://www.inc.com/kevin-j-ryan/self-driving-cars-powered-by-people-playing-games-mighty-ai.html>

Can you think of any other
potential applications?

Semantic Segmentation: Today's Topics

- Motivation
- **Datasets**
- Evaluation metric
- Fully convolutional network
- Swin transformer
- Discussion (chosen by YOU 😊)

Datasets

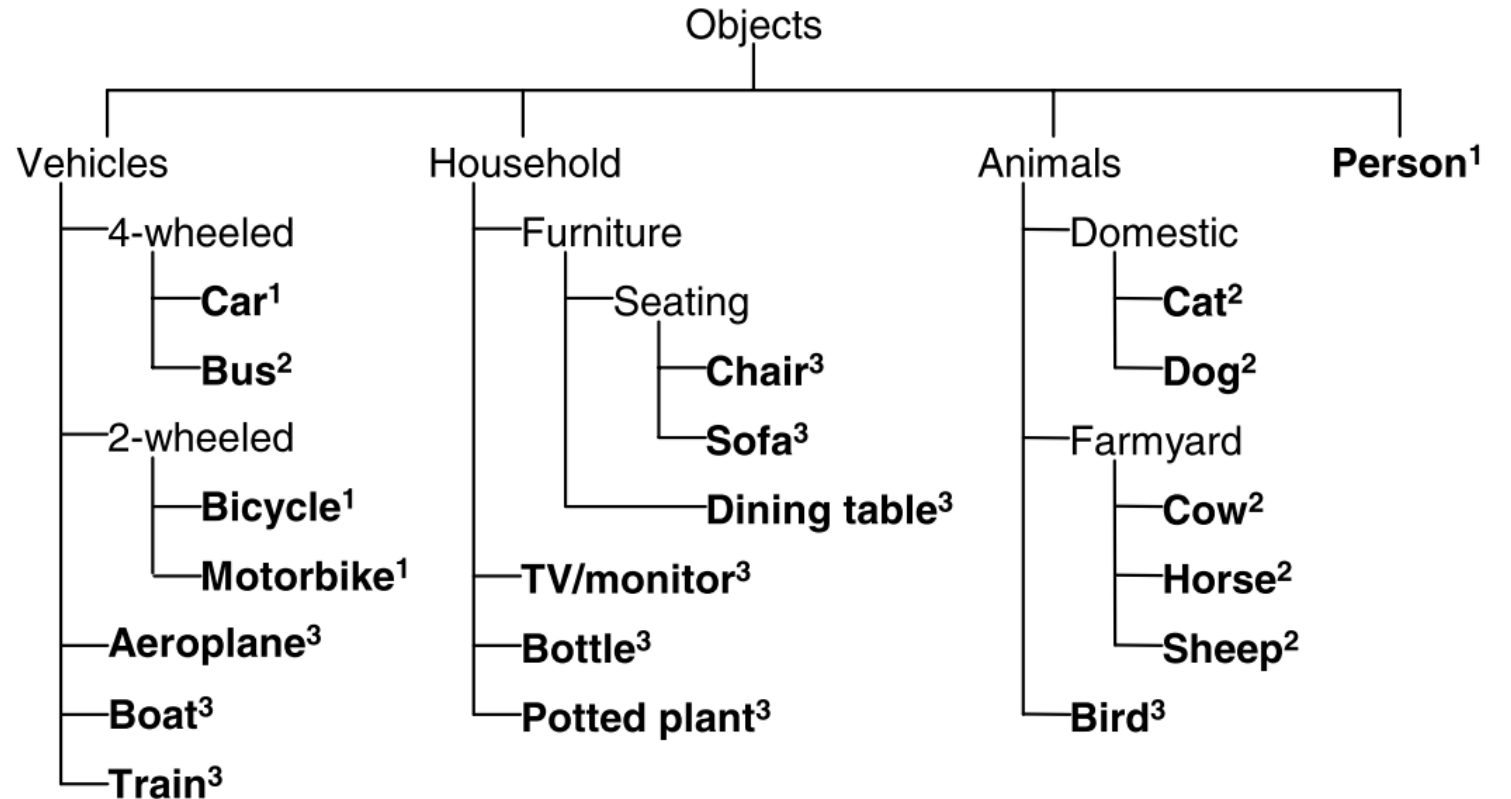


Trend: build bigger datasets

VOC

1. Category Selection

- 20 categories chosen:
- 1) Initial 4 categories stem from existing dataset
- 2) 2006: added 6 classes
- 3) 2007: added 10 classes
- Categories added for more generalization and finer-grained coverage



(superscript indicates year of inclusion in the challenge: 2005¹, 2006², 2007³)

VOC

1. Category Selection

- 20 categories chosen:
 - 1) Initial 4 categories stem from existing dataset
 - 2) 2006: added 6 classes
 - 3) 2007: added 10 classes
- Categories added for more generalization and finer-grained coverage

2. Image Collection

- 500,000 images retrieved from Flickr with many search terms

(search terms per category)

- **aeroplane**, airplane, plane, biplane, monoplane, aviator, bomber, hydroplane, airliner, aircraft, fighter, airport, hangar, jet, boeing, fuselage, wing, propellor, flying
- **bicycle**, bike, cycle, cyclist, pedal, tandem, saddle, wheel, cycling, ride, wheelie
- **bird**, birdie, birdwatching, nest, sea, aviary, birdcage, bird feeder, bird table
- **boat** ship, barge, ferry, canoe, boating, craft, liner, cruise, sailing, rowing, watercraft, regatta, racing, marina, beach, water, canal, river, stream, lake, yacht
- **bottle**, cork, wine, beer, champagne, ketchup, squash, soda, coke, lemonade, dinner, lunch, breakfast
- **bus**, omnibus, coach, shuttle, jitney, double-decker, motorbus, school bus, depot, terminal, station, terminus, passenger, route
- **car**, automobile, cruiser, motorcar, vehicle, hatchback, saloon, convertible, limousine, motor, race, traffic, trip, rally, city, street, road, lane, village, town, centre, shopping, downtown, suburban
- **cat**, feline, pussy, mew, kitten, tabby, tortoiseshell, ginger, stray
- **chair**, seat, rocker, rocking, deck, swivel, camp, chaise, office, studio, armchair, recliner, sitting, lounge, living room, sitting room
- **cow**, beef, heifer, moo, dairy, milk, milking, farm
- **dog**, hound, bark, kennel, heel, bitch, canine, puppy, hunter, collar, leash
- **horse**, gallop, jump, buck, equine, foal, cavalry, saddle, canter, buggy, mare, neigh, dressage, trial, racehorse, steeplechase, thoroughbred, cart, equestrian, paddock, stable, farrier
- **motorbike**, motorcycle, minibike, moped, dirt, pillion, biker, trials, motorcycling, motorcyclist, engine, motocross, scramble, sidecar, scooter, trail
- **person**, people, family, father, mother, brother, sister, aunt, uncle, grandmother, grandma, grandfather, grandpa, grandson, granddaughter, niece, nephew, cousin
- **sheep**, ram, fold, fleece, shear, baa, bleat, lamb, ewe, wool, flock
- **sofa**, chesterfield, settee, divan, couch, bolster
- **table**, dining, cafe, restaurant, kitchen, banquet, party, meal
- **potted plant**, pot plant, plant, patio, windowsill, window sill, yard, greenhouse, glass house, basket, cutting, pot, cooking, grow
- **train**, express, locomotive, freight, commuter, platform, subway, underground, steam, railway, railroad, rail, tube, underground, track, carriage, coach, metro, sleeper, railcar, buffet, cabin, level crossing
- **tv/monitor**, television, plasma, flatscreen, flat screen, lcd, crt, watching, dvd, desktop, computer, computer monitor, PC, console, game

VOC

1. Category Selection

- 20 categories chosen:
 - 1) Initial 4 categories stem from existing dataset
 - 2) 2006: added 6 classes
 - 3) 2007: added 10 classes
- Categories added for more generalization and finer-grained coverage

2. Image Collection

- 500,000 images retrieved from Flickr with many search terms

3. Image Verification + Image Annotation

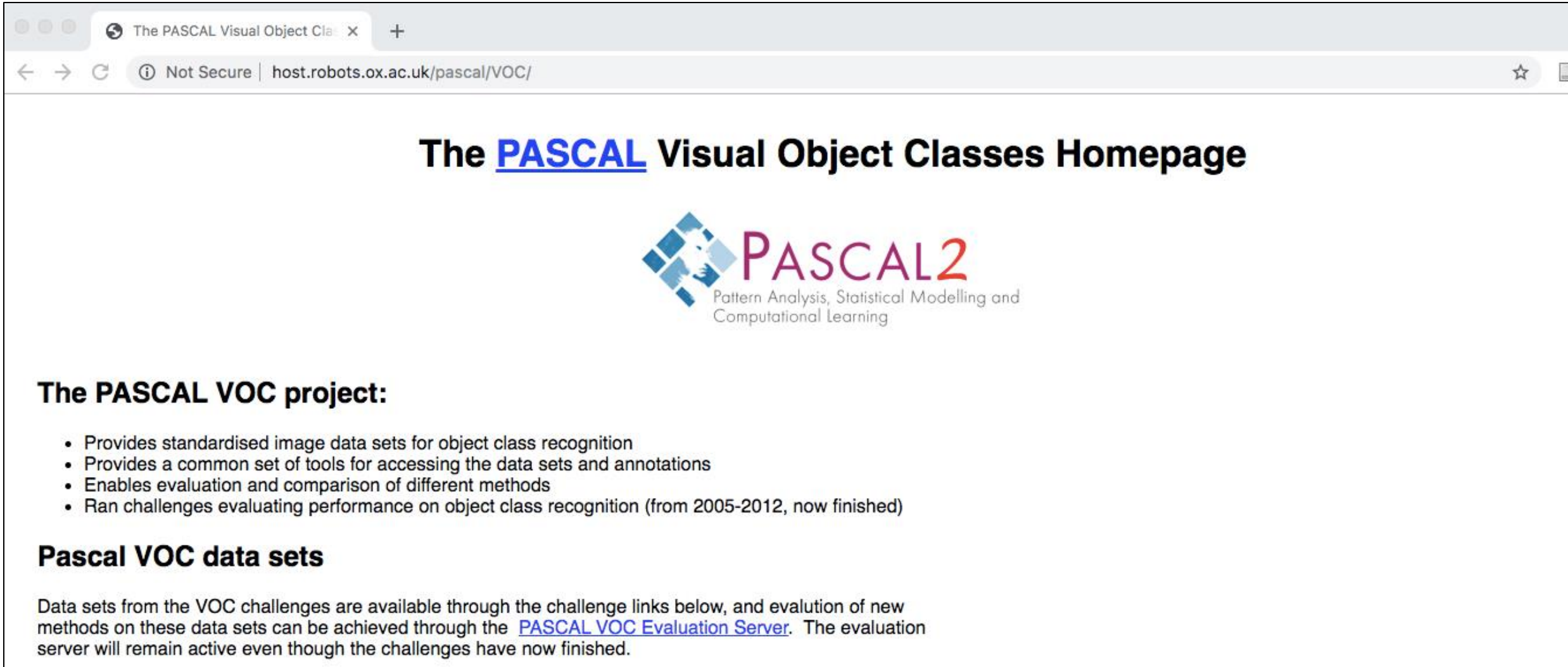
- University of Leeds annotation parties to recruit annotators annually
- Annotation guidelines & real-time assistance: detections subsequently refined to segmentations
- Post-hoc correction/feedback about the number and kind of errors made
- Annotations for each object class merged and another class added for background

VOC: Datasets Evolved

The table below gives a brief summary of the main stages of the VOC development.


Year	Statistics	New developments	Notes
2005	Only 4 classes: bicycles, cars, motorbikes, people. Train/validation/test: 1578 images containing 2209 annotated objects.	Two competitions: classification and detection	Images were largely taken from existing public datasets, and were not as challenging as the flickr images subsequently used. This dataset is obsolete.
2006	10 classes: bicycle, bus, car, cat, cow, dog, horse, motorbike, person, sheep. Train/validation/test: 2618 images containing 4754 annotated objects.	Images from flickr and from Microsoft Research Cambridge (MSRC) dataset	The MSRC images were easier than flickr as the photos often concentrated on the object of interest. This dataset is obsolete.

VOC Annual Workshop



The screenshot shows a web browser window with the address bar displaying "host.robots.ox.ac.uk/pascal/VOC/". The page title is "The PASCAL Visual Object Classes Homepage". The main content features the PASCAL2 logo, which consists of a blue and white geometric design and the text "PASCAL2" in red, with the subtitle "Pattern Analysis, Statistical Modelling and Computational Learning" below it. The page is organized into sections: "The PASCAL VOC project:" followed by a bulleted list of project details, "Pascal VOC data sets" followed by a paragraph about data availability and evaluation.

The **PASCAL** Visual Object Classes Homepage



Pattern Analysis, Statistical Modelling and Computational Learning

The PASCAL VOC project:

- Provides standardised image data sets for object class recognition
- Provides a common set of tools for accessing the data sets and annotations
- Enables evaluation and comparison of different methods
- Ran challenges evaluating performance on object class recognition (from 2005-2012, now finished)

Pascal VOC data sets

Data sets from the VOC challenges are available through the challenge links below, and evaluation of new methods on these data sets can be achieved through the [PASCAL VOC Evaluation Server](#). The evaluation server will remain active even though the challenges have now finished.

VOC: Boundary Accuracy Heuristic

Image

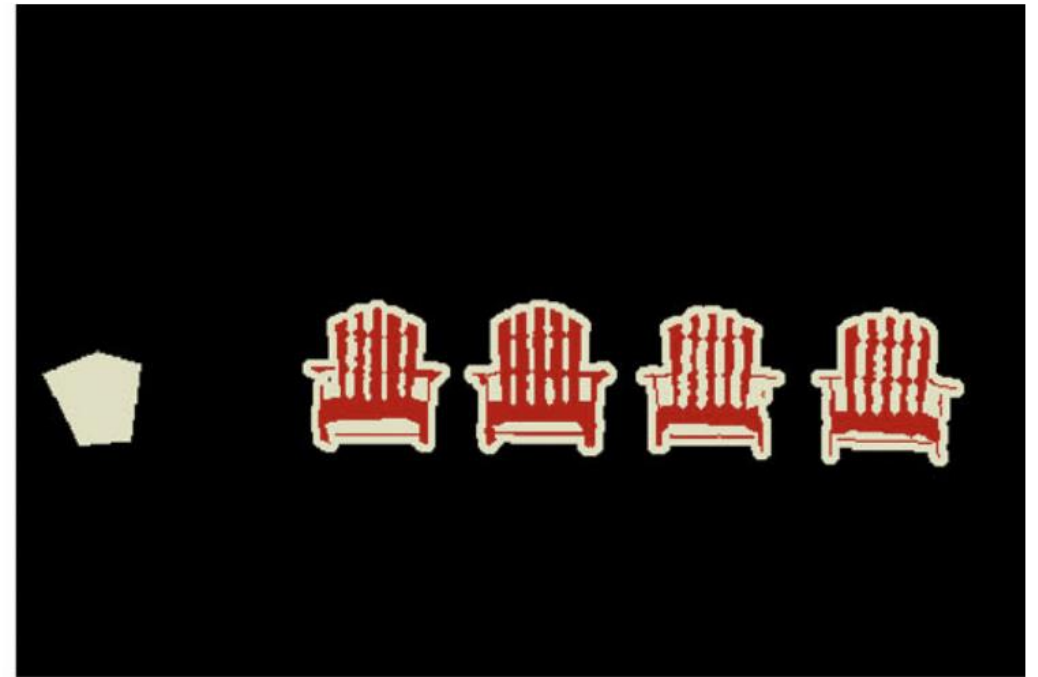


Class segmentation



“To give high accuracy but to keep the annotation time short enough to provide a large image set, a border area of 5 pixels width was allowed around each object where the pixels were labelled neither object nor background.”

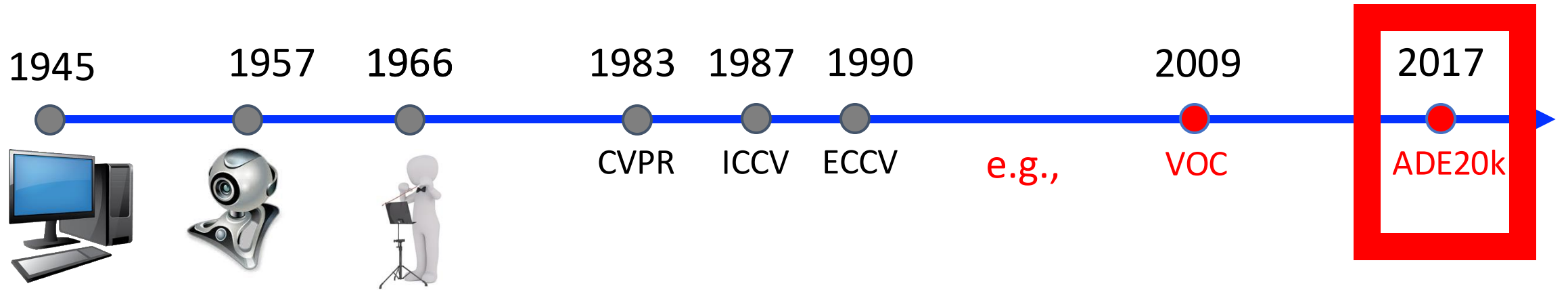
What is a Limitation of Datasets Built Around Specific Categories (e.g., Objects)?



Most pixels are labeled as `background`!

Lacks knowledge anything else is in the scene, such as a house, trees or flowers!

Datasets



ADE20K

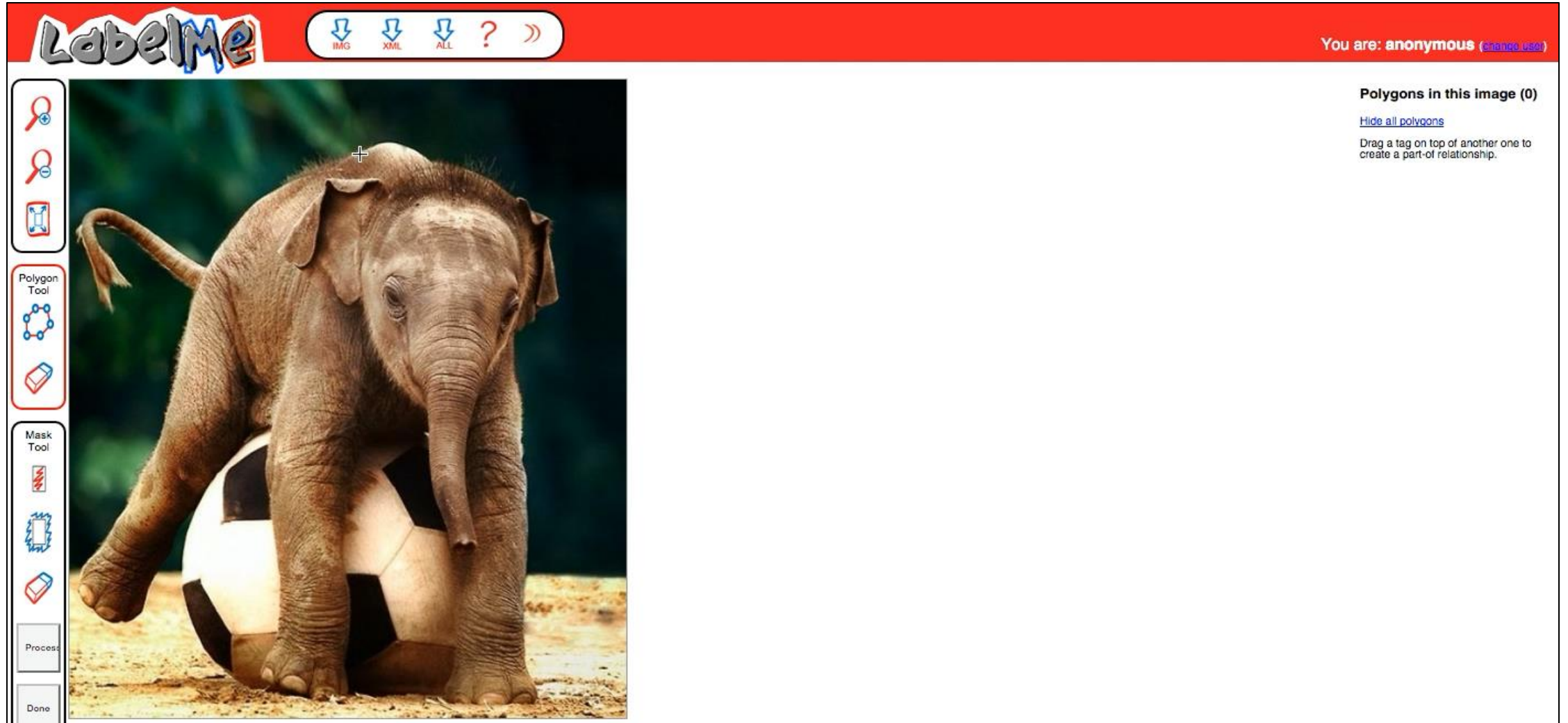
1. Image Collection

- 25,210 images collected from existing datasets (SUN, Places, and LabelMe)
- Selected to capture all scene categories defined in SUN

2. Region Localization and Category Assignment

- A single person annotated all images into three types and kept adding new categories as they were observed: (1) objects, (2) object parts, and (3) attributes (e.g., occluded)

ADE20K: User Annotation Tool



The screenshot displays the LabelMe web application interface. At the top, there is a red header bar with the "LabelMe" logo on the left and navigation icons (IMG, XML, ALL, ?, >>) in the center. On the right side of the header, it says "You are: anonymous (change user)".

The main content area features a large image of a young elephant standing on a soccer ball. A small white crosshair is positioned on the elephant's back, indicating a point of interest or a starting point for an annotation. To the left of the image is a vertical toolbar with several icons: a red circle with a plus sign, a red circle with a minus sign, a red square with a plus sign, a blue polygon tool icon, a red eraser icon, a red square with a diagonal line, a blue square with a diagonal line, and a red square with a diagonal line. Below the toolbar are two buttons labeled "Process" and "Done".

On the right side of the main content area, there is a section titled "Polygons in this image (0)". Below this title is a link that says "Hide all polygons" and a paragraph of text: "Drag a tag on top of another one to create a part-of relationship."

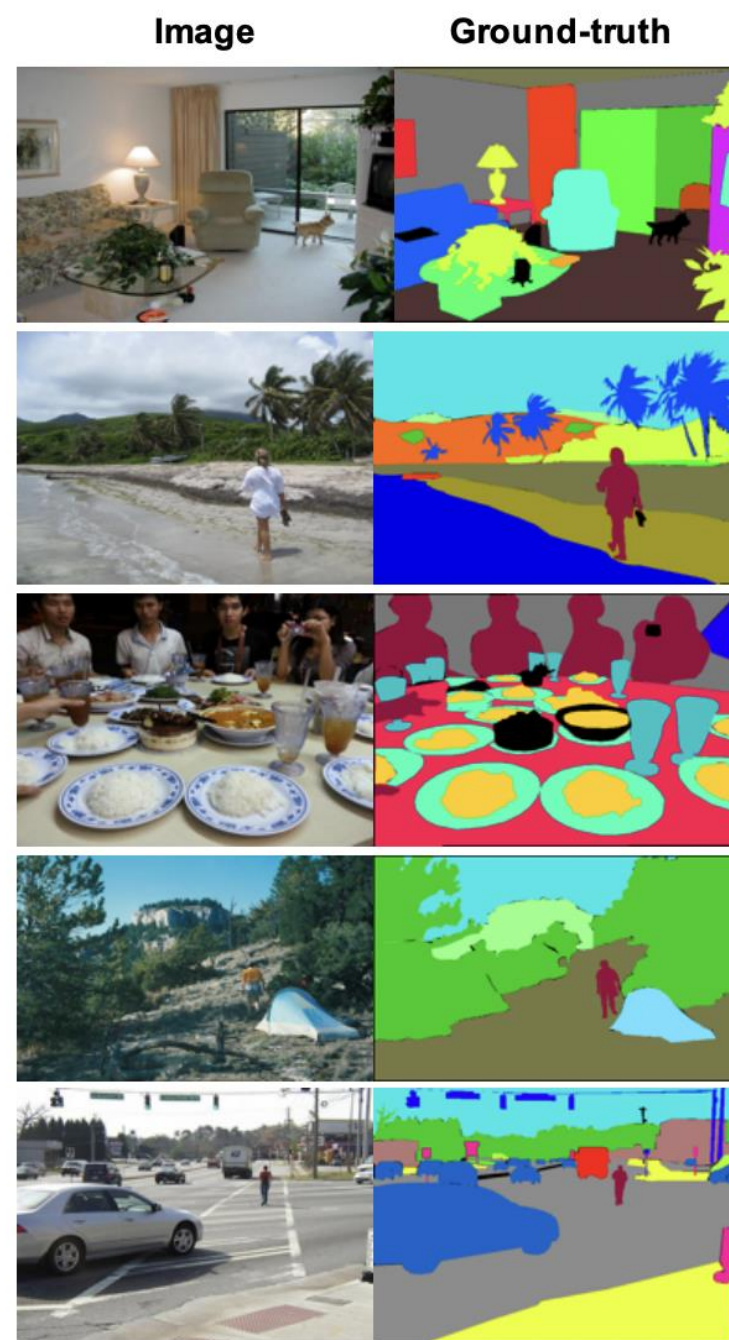
ADE20K: User Annotation Tool



- ceiling
- wall
- wall
- window (arch)
 - ↳ pane (glass)
 - ↳ figurine
- door frame
- double door
 - ↳ door
 - ↳ handle
- tray
- figurine
- refrigerator (crop)
- cabinet
 - ↳ door
 - ↳ knob
 - ↳ door
 - ↳ knob
 - ↳ door
 - ↳ knob
- jar
- cabinet
 - ↳ door
 - ↳ knob
- cabinet
 - ↳ door
 - ↳ knob
 - ↳ door
 - ↳ knob
- cabinet
 - ↳ door
 - ↳ knob
- microwave
 - ↳ door
 - ↳ window
 - ↳ button (door release)
- outlet
- pot
- sink
 - ↳ faucet
- soap dispenser
- spice rack
- coffee maker
- knife set
- knife set
- range
 - ↳ button panel
 - ↳ dial
 - ↳ dial
 - ↳ dial
 - ↳ dial
 - ↳ screen time
 - ↳ stove
 - ↳ burner
 - ↳ burner
 - ↳ burner
 - ↳ burner
 - ↳ oven
 - ↳ door
 - ↳ handle
- toaster
- blender
- pot
- box
- worktop
- cabinet
 - ↳ drawer
 - ↳ knob
- jar
- salt cellar
- worktop
- paper towels
- dishwasher
- cabinet
- cabinet
- bottle rack
- napkin rack
- kitchen island
- glass (wine)
- glass (wine)
- coasters
- bowl
- bowl
- trash can
- dog dish
- dog dish
- chair
 - ↳ back
 - ↳ seat (fabric)
 - ↳ leg
 - ↳ leg
 - ↳ leg
- chair
 - ↳ back
 - ↳ seat (fabric)
 - ↳ leg
 - ↳ leg
 - ↳ leg
- chair
 - ↳ back
 - ↳ seat (fabric)
 - ↳ leg
 - ↳ leg
 - ↳ leg
- side table (crop)
- rug
- sofa (crop)
- cushion
- cushion
- cushion
- floor (tile)
- carpet
- bowl
- light switch
- picture (map)

ADE20K

- Includes:
 - “**things**”: objects that can easily be labeled; e.g., person, chair
 - “**stuff**”: objects with no clear boundaries; e.g., sky, grass

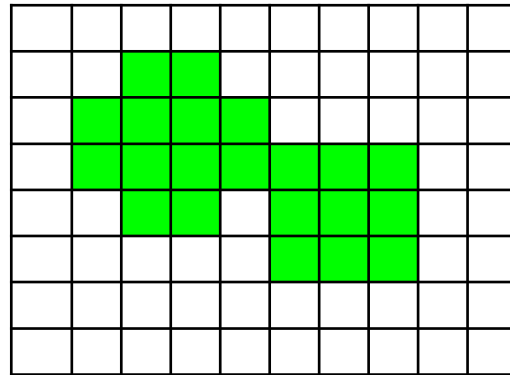


Semantic Segmentation: Today's Topics

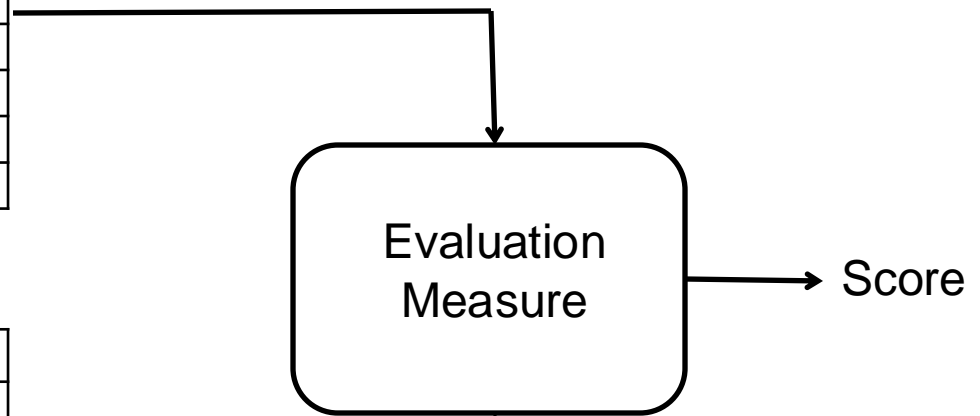
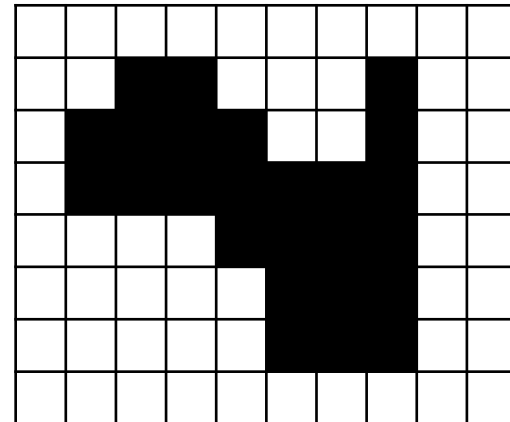
- Motivation
- Datasets
- **Evaluation metric**
- Fully convolutional network
- Swin transformer
- Discussion (chosen by YOU 😊)

Evaluation Metric

Ground Truth:

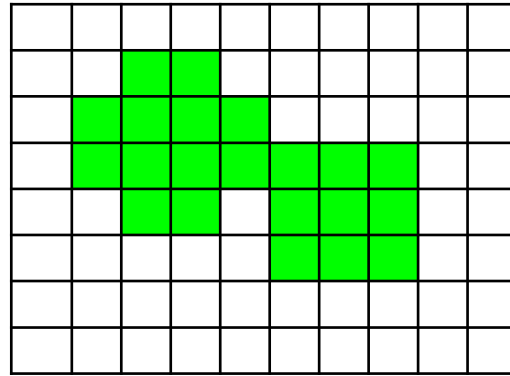


Algorithm:

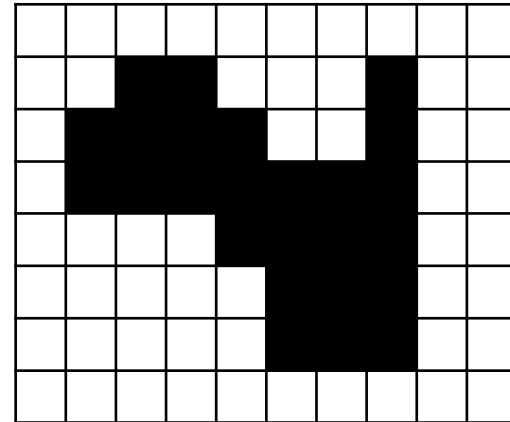


Recall: IoU Metric

Ground Truth:



Algorithm:

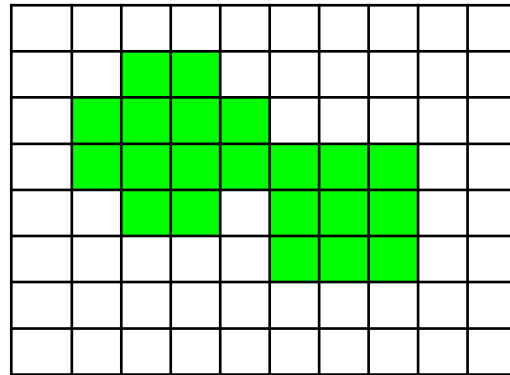


$$\frac{|A \cap B|}{|A \cup B|}$$

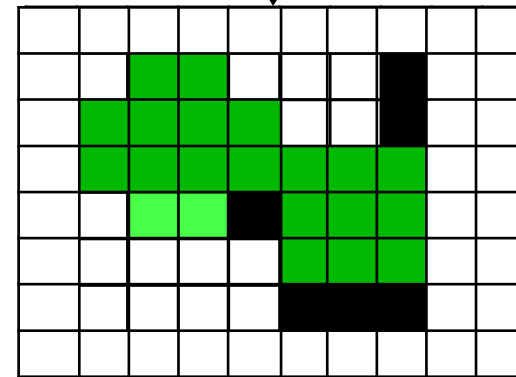
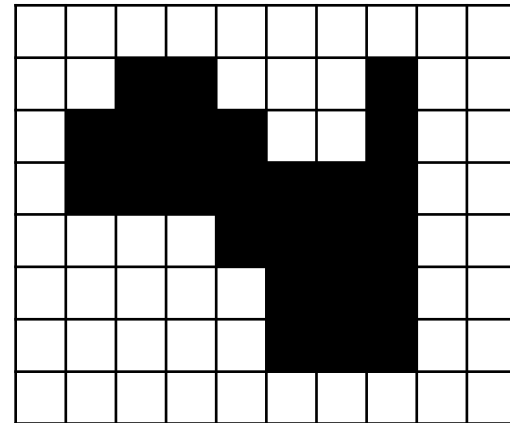
Score

Recall: IoU Metric

Ground Truth:



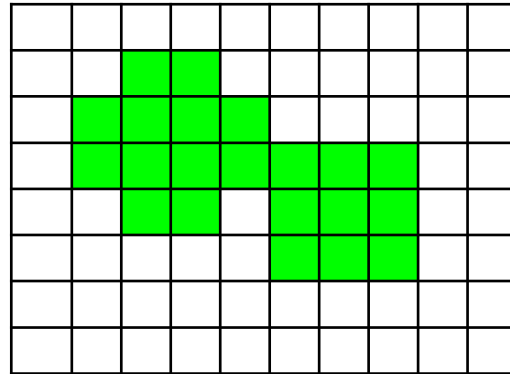
Algorithm:



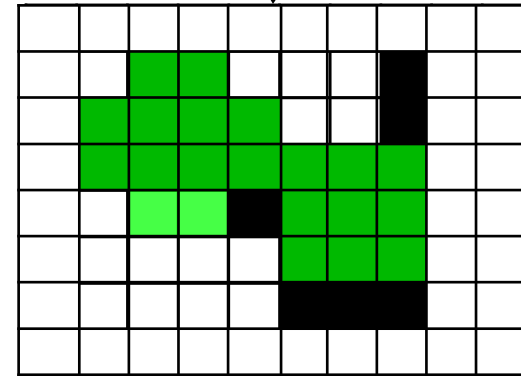
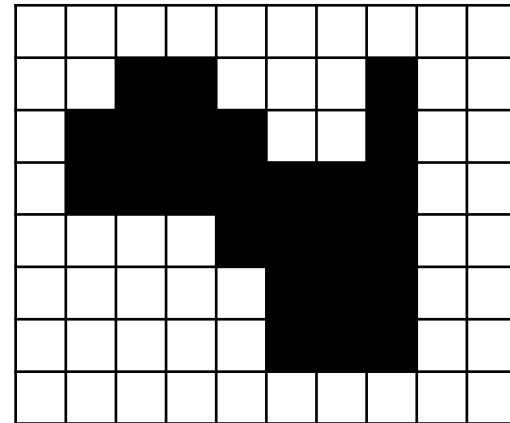
?

Recall: IoU Metric

Ground Truth:



Algorithm:



$$\frac{19}{27}$$

Mean IoU (mIoU)

- Mean IoU score over all categories

Semantic Segmentation: Today's Topics

- Motivation
- Datasets
- Evaluation metric
- Fully convolutional network
- Swin transformer
- Discussion (chosen by YOU 😊)

Why Fully Convolutional Network?

Named after the proposed technique that excludes fully connected layers:

Jonathon Long, Evan Shelhamer, and Trevor Darrell. “Fully Convolutional Networks for Semantic Segmentation.” CVPR 2015.

Key Novelties of Fully Convolutional Networks

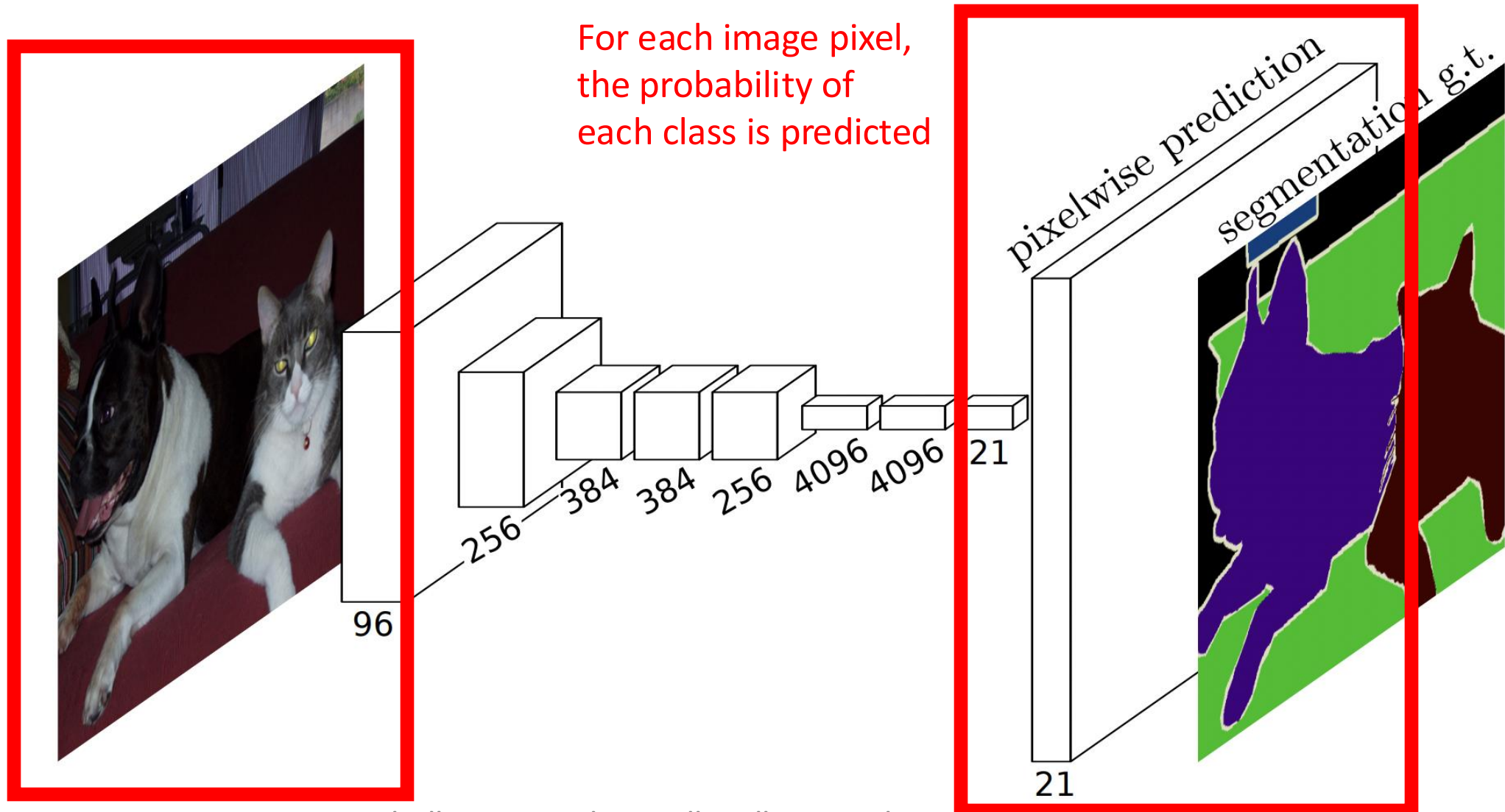
First work for pixelwise prediction to:

1. Train fully convolutional networks end-to-end
2. Use supervised pre-training (recall, ViT benefited from this as well)

Architecture

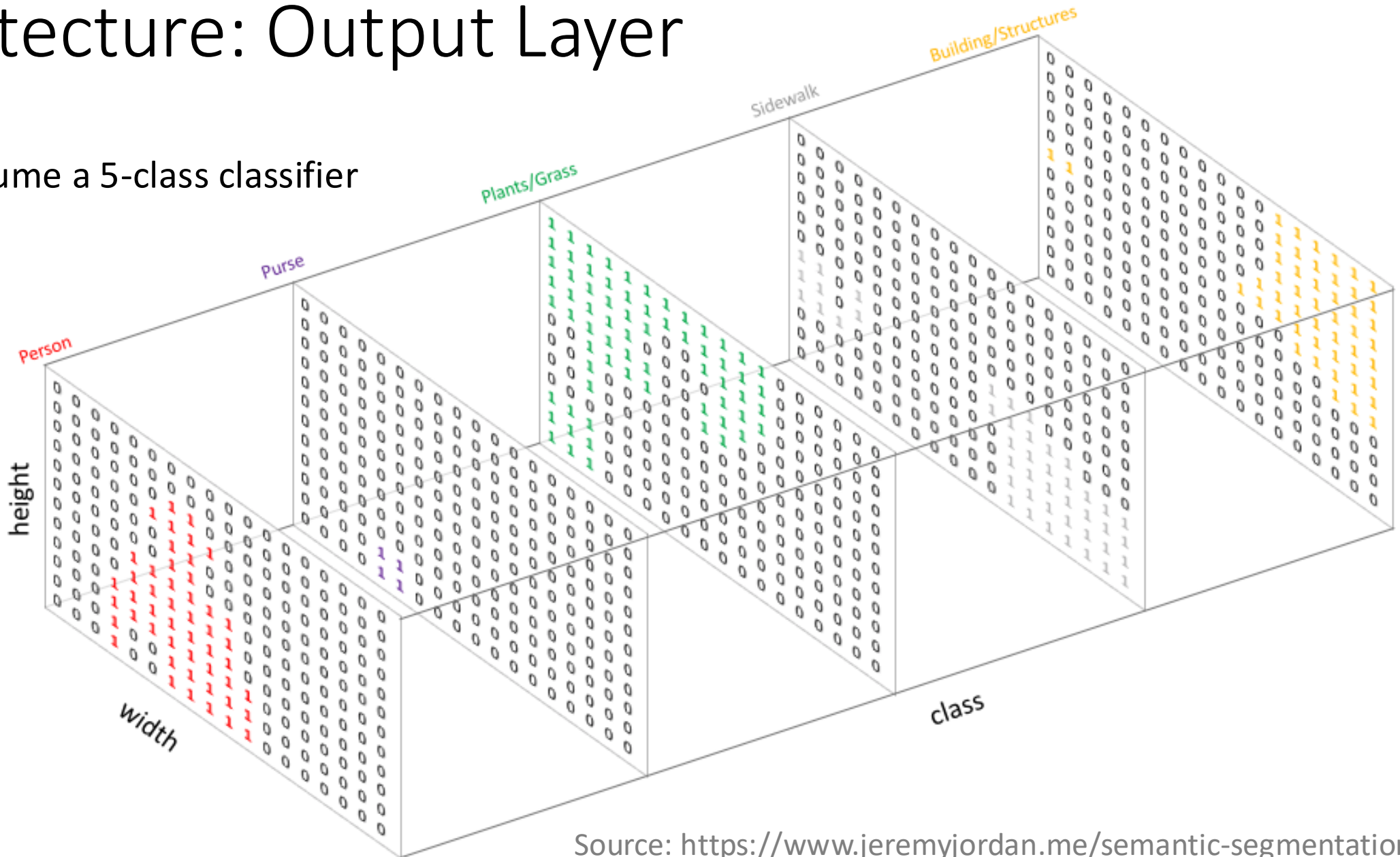
Input: RGB image of ANY size

Output: Image of same size as input



Architecture: Output Layer

- e.g., assume a 5-class classifier



Architecture: Output Layer

- e.g., assume a 5-class classifier; output 1-hot encoding collapsed into single mask image



0: Background/Unknown

1: Person

2: Purse

3: Plants/Grass

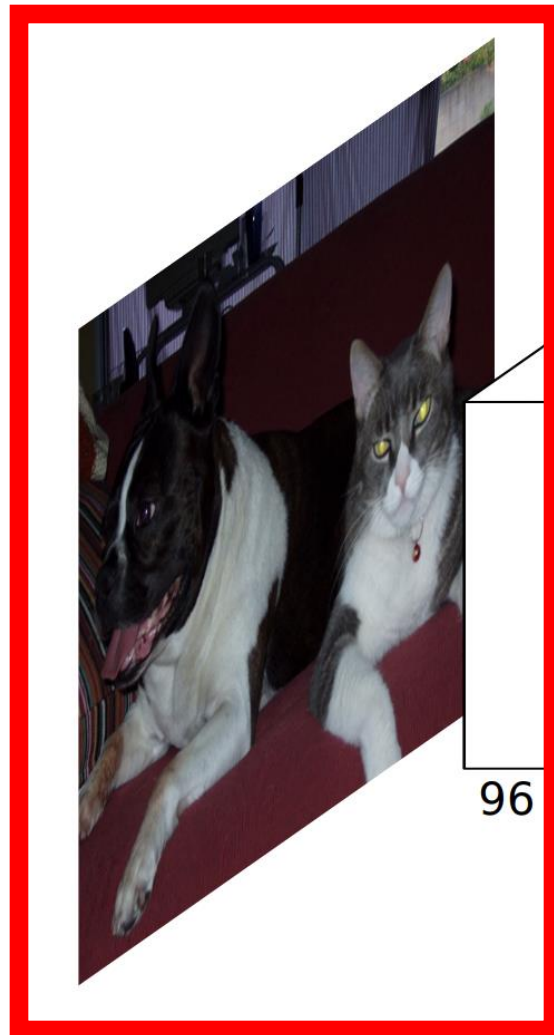
4: Sidewalk

5: Building/Structures

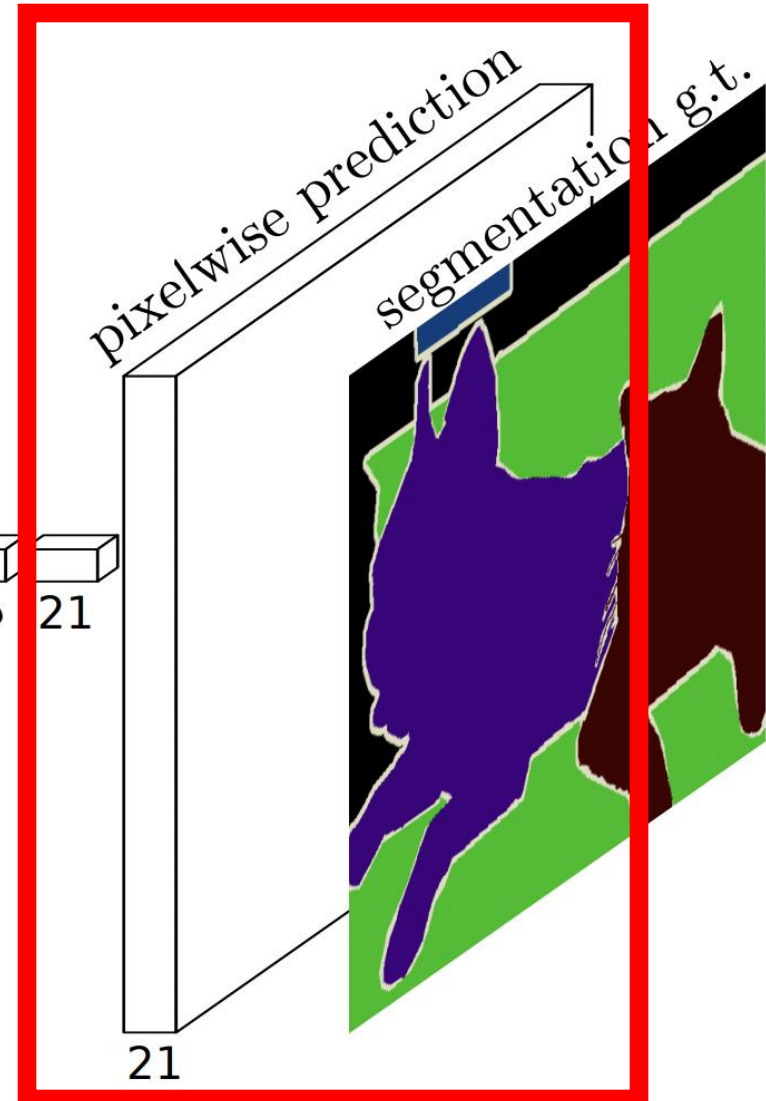
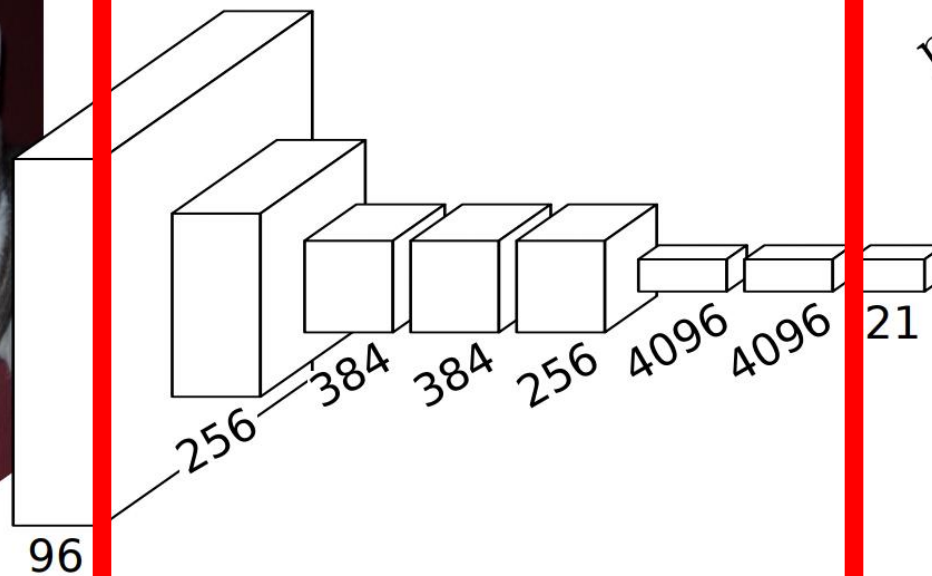
Architecture

Input: RGB image of ANY size

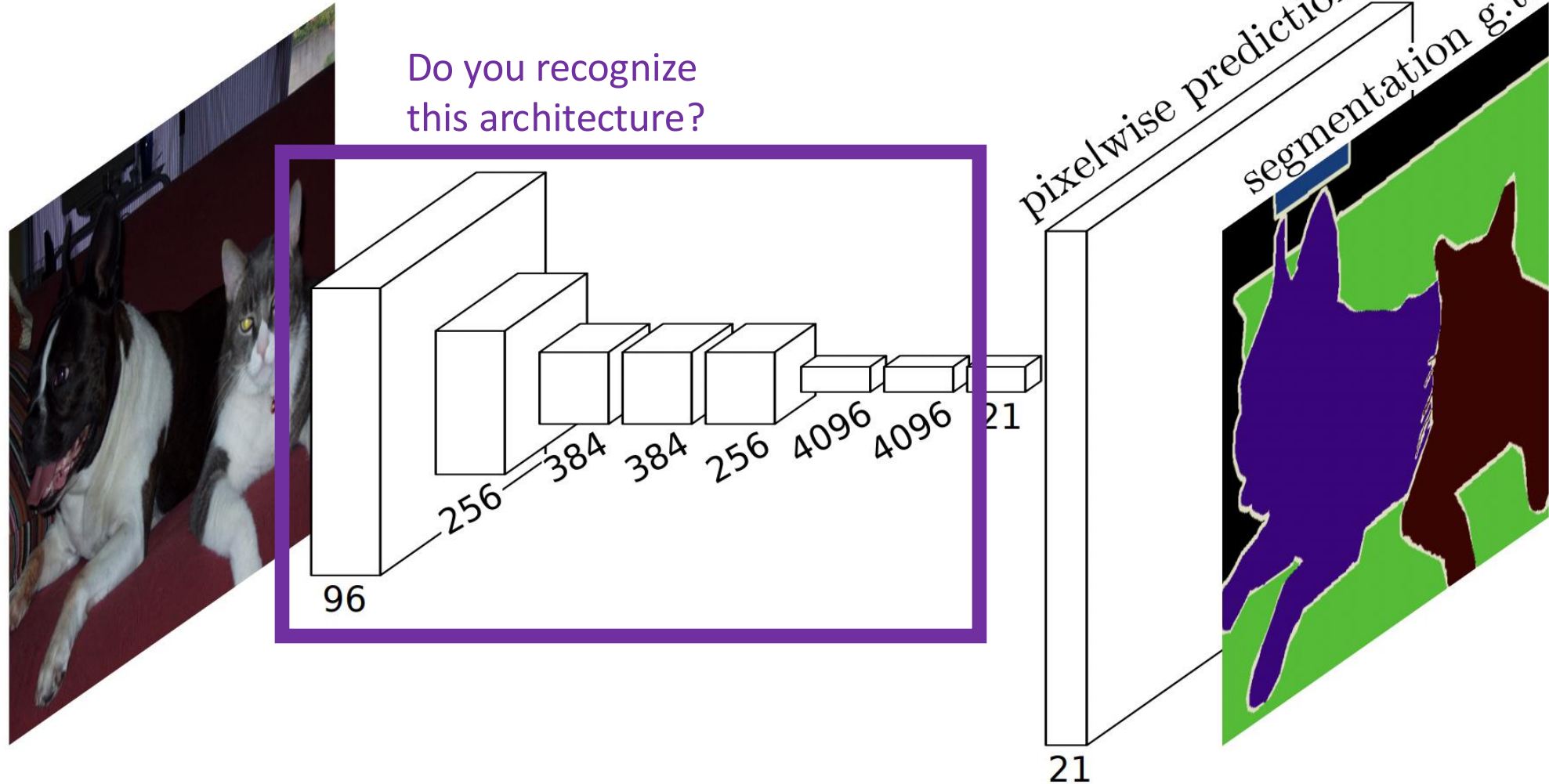
Output: Image of same size as input



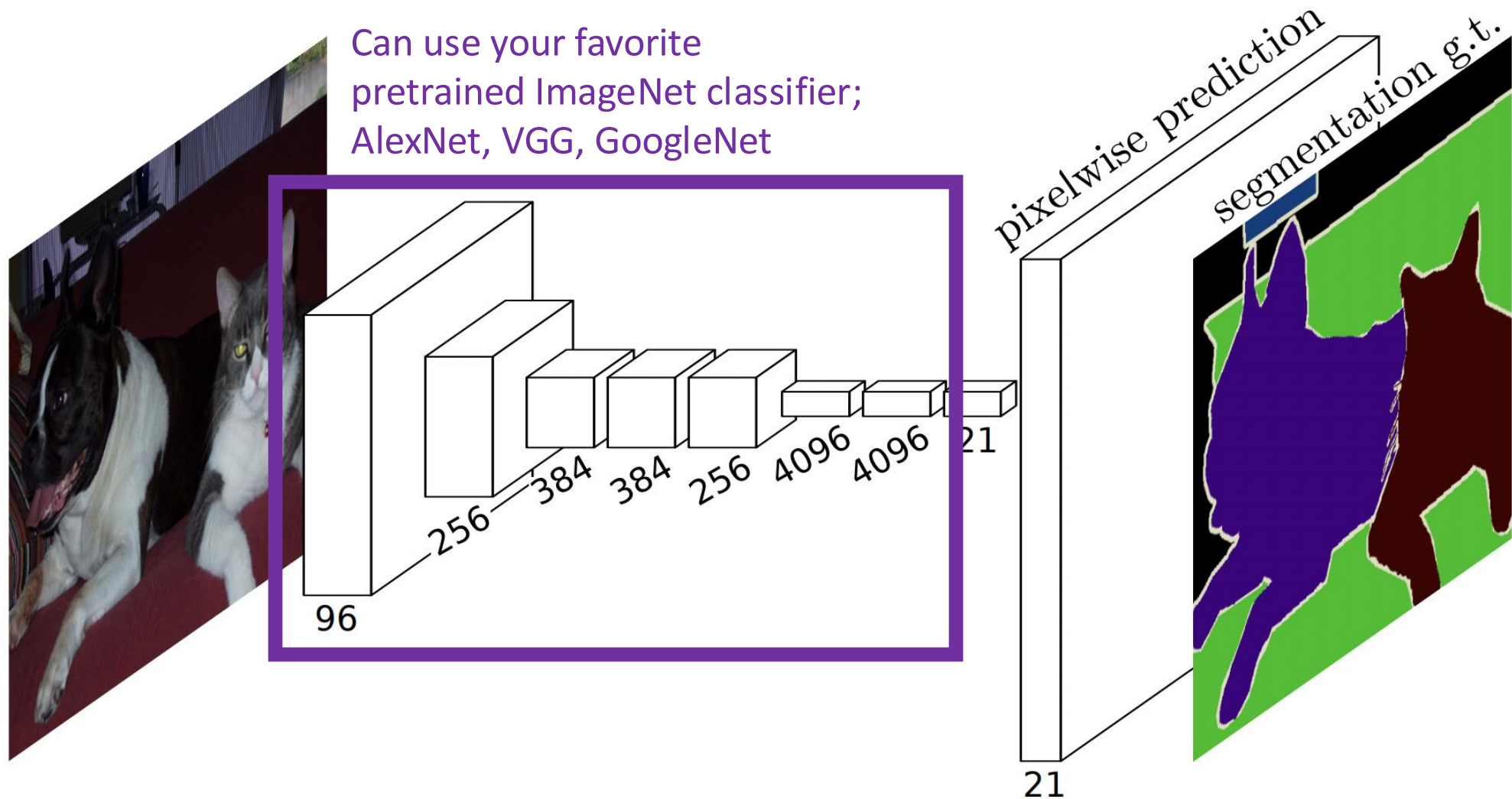
How many classes are there for VOC?
- 21 (20 object classes plus background)



Architecture



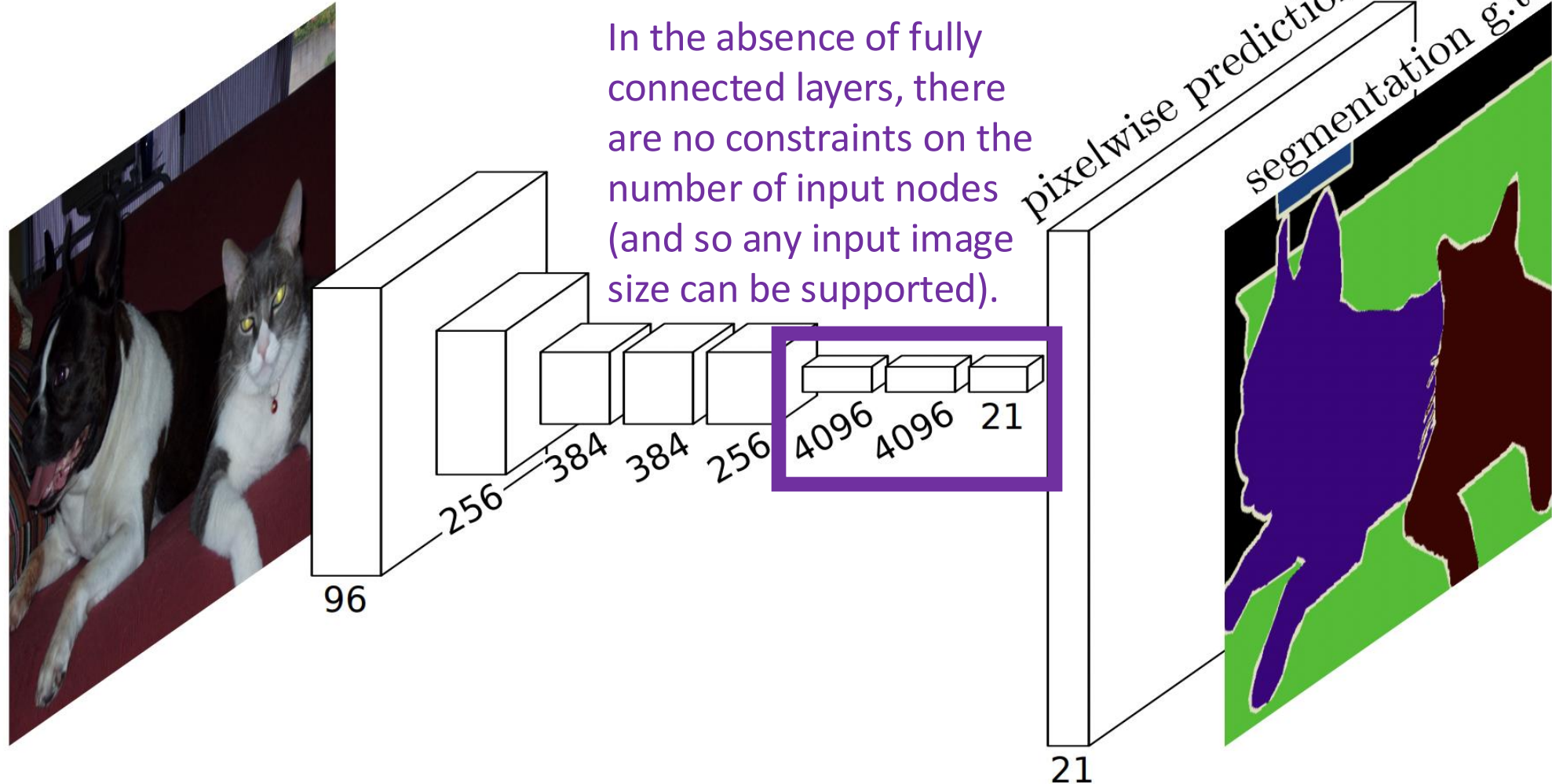
Architecture



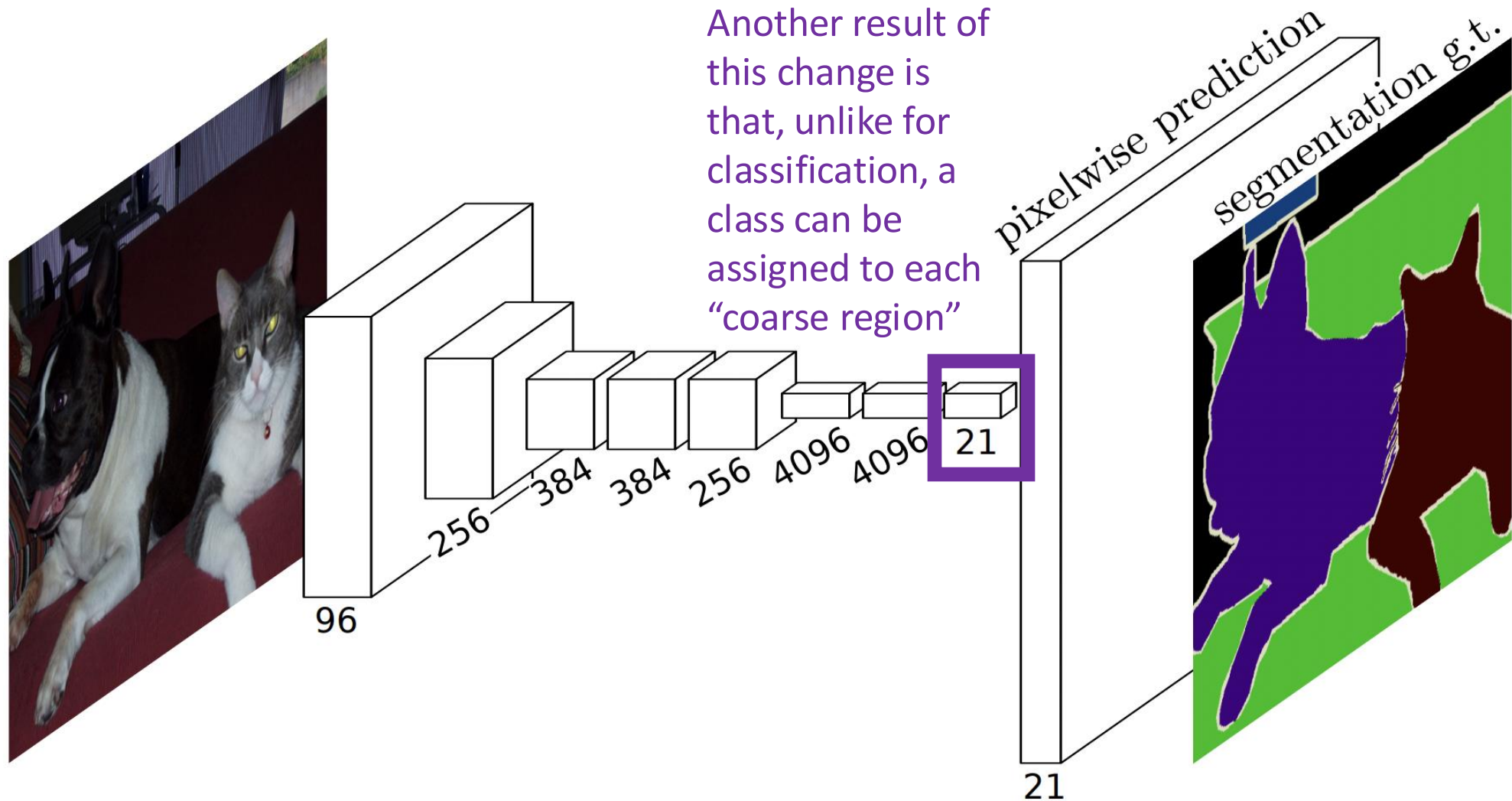
Architecture

To make architecture fully convolutional, fully connected layers are converted to convolutional layers.

In the absence of fully connected layers, there are no constraints on the number of input nodes (and so any input image size can be supported).

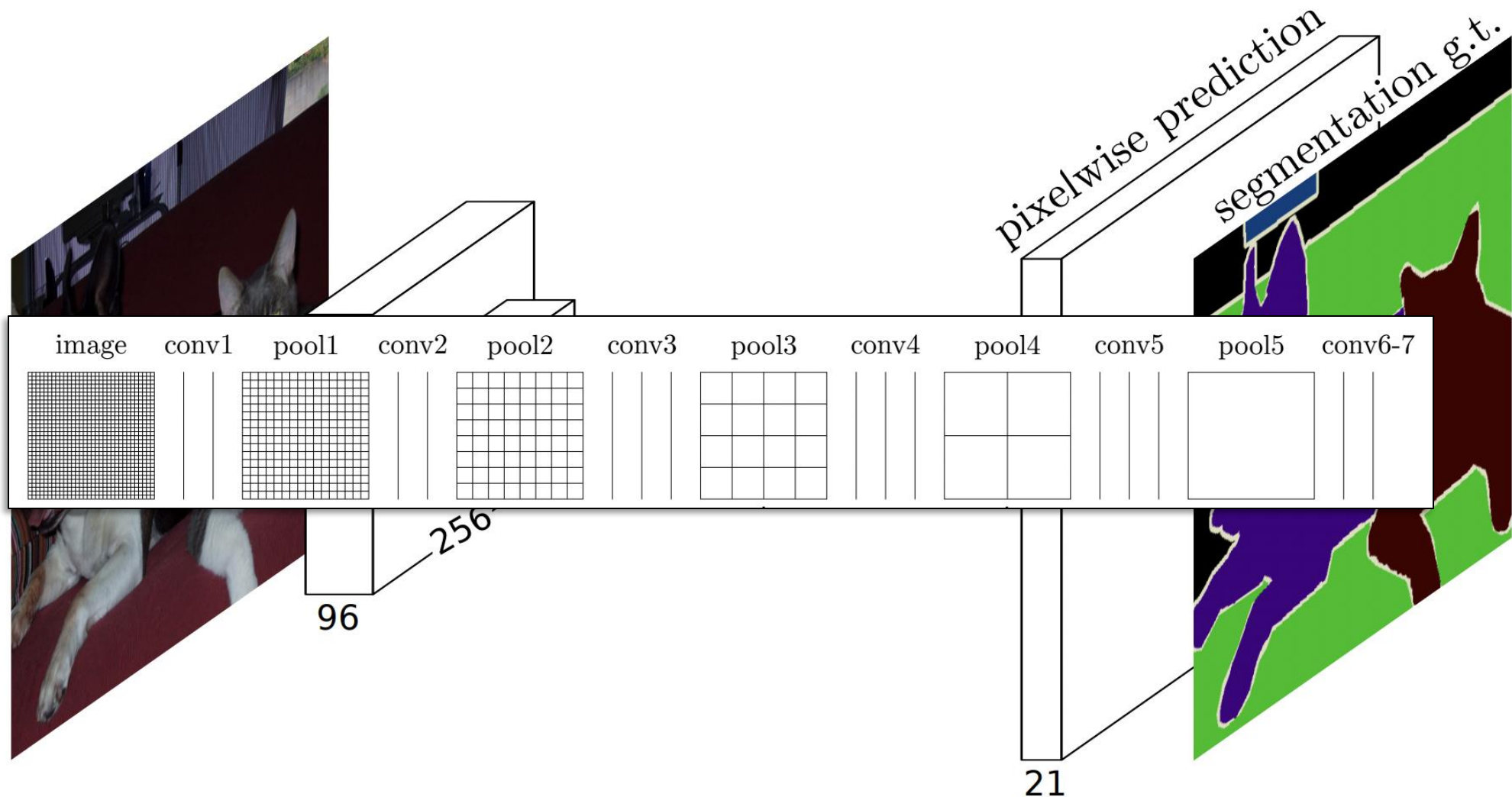


Architecture

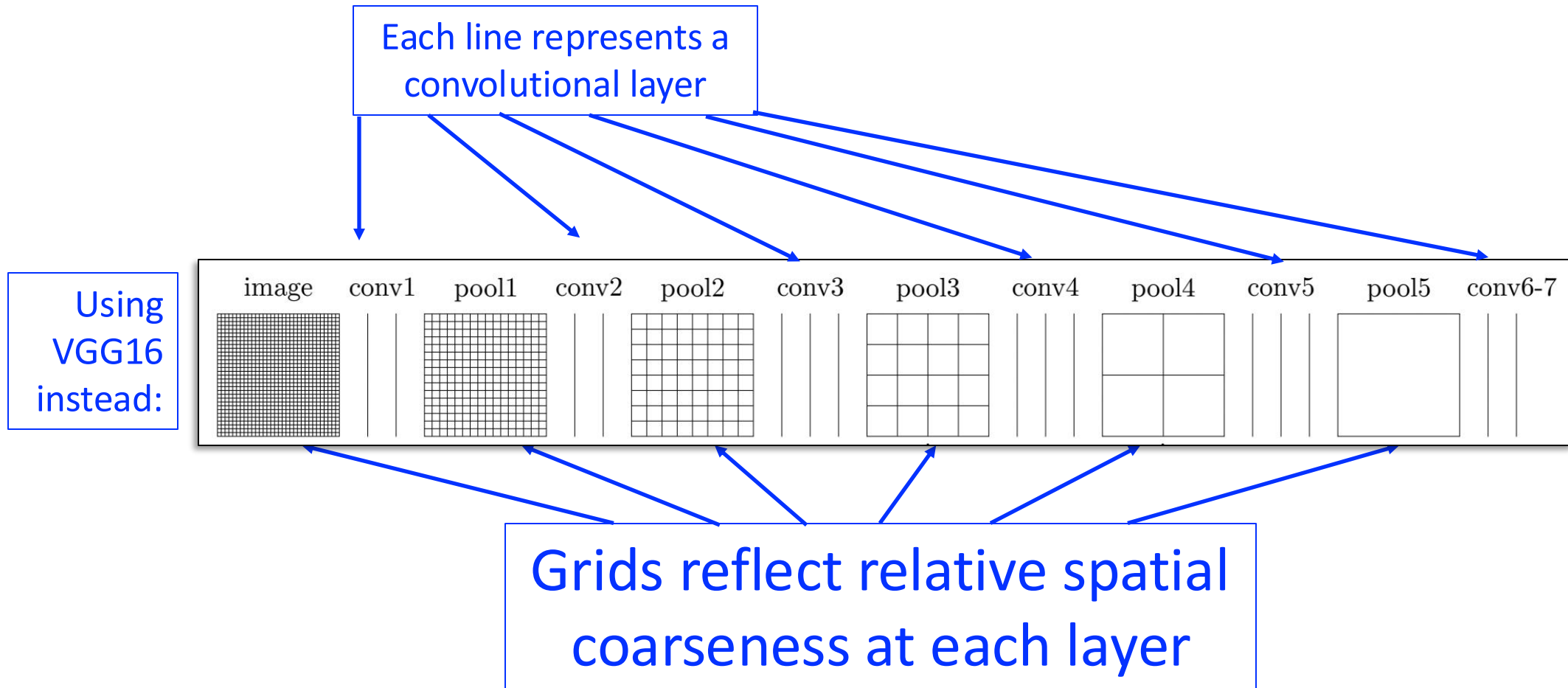


Architecture: Coarse Region Classification (Recall Intuition)

Using
VGG16
instead:

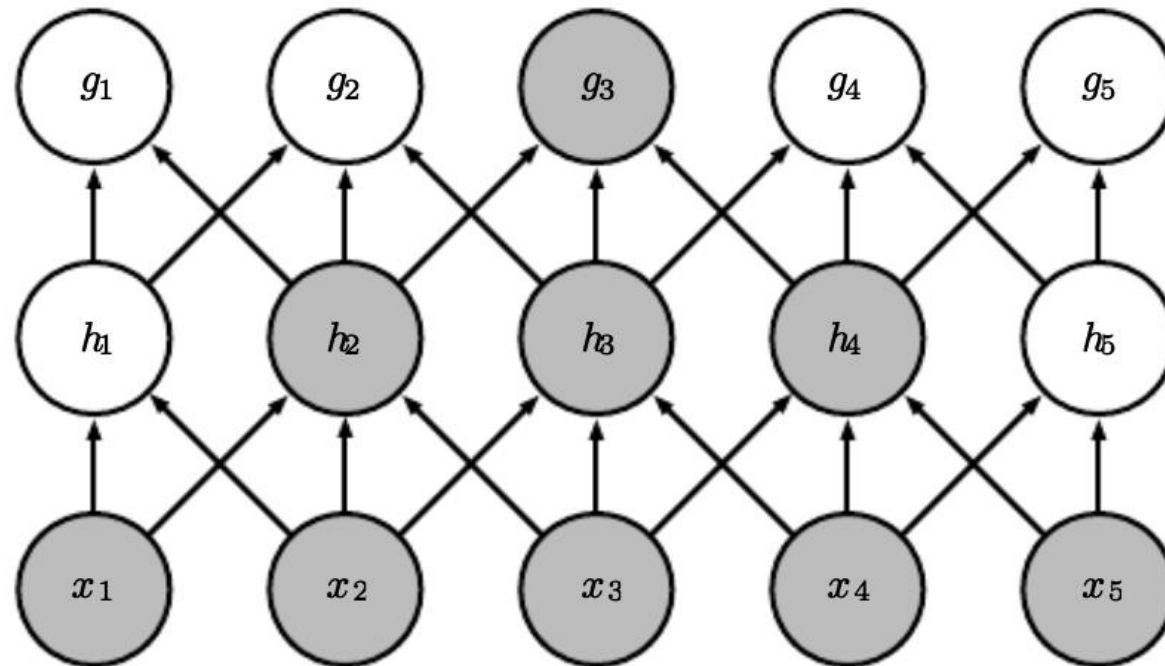


Architecture: Coarse Region Classification (Recall Intuition)

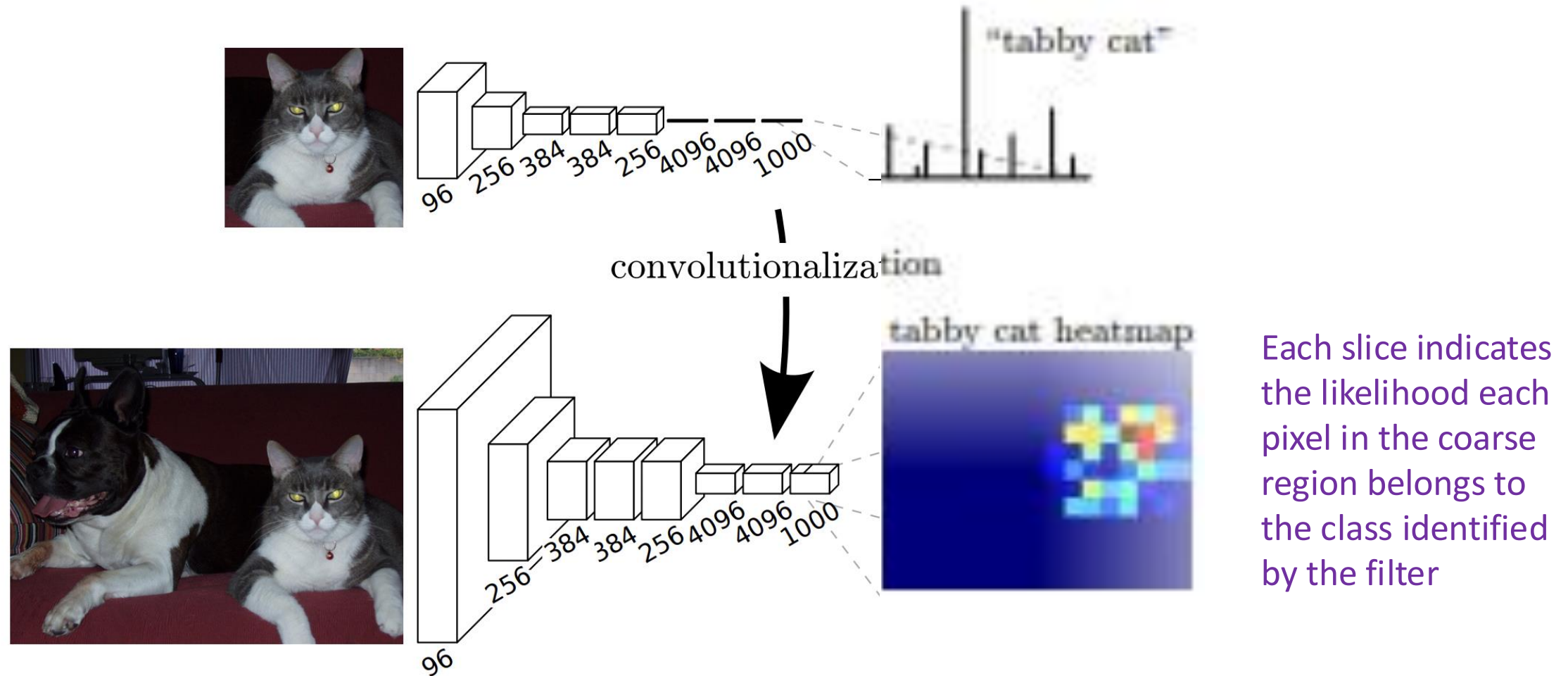


Architecture: Coarse Region Classification (Recall Intuition)

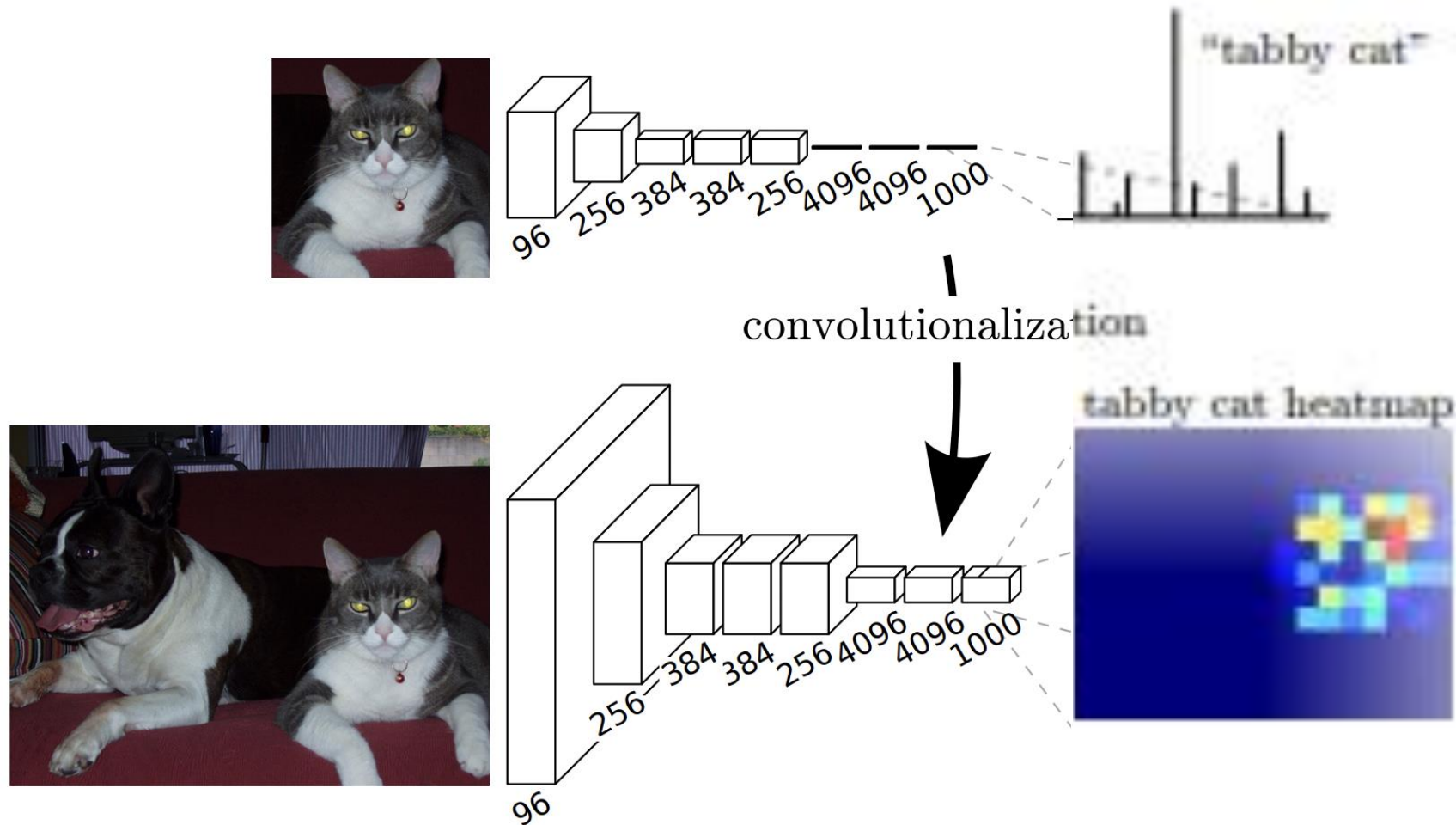
Stacking many convolutional layers leads to learning patterns in increasingly **larger regions of the input (e.g., pixel) space**.



Architecture: Fully vs Convolution Layers

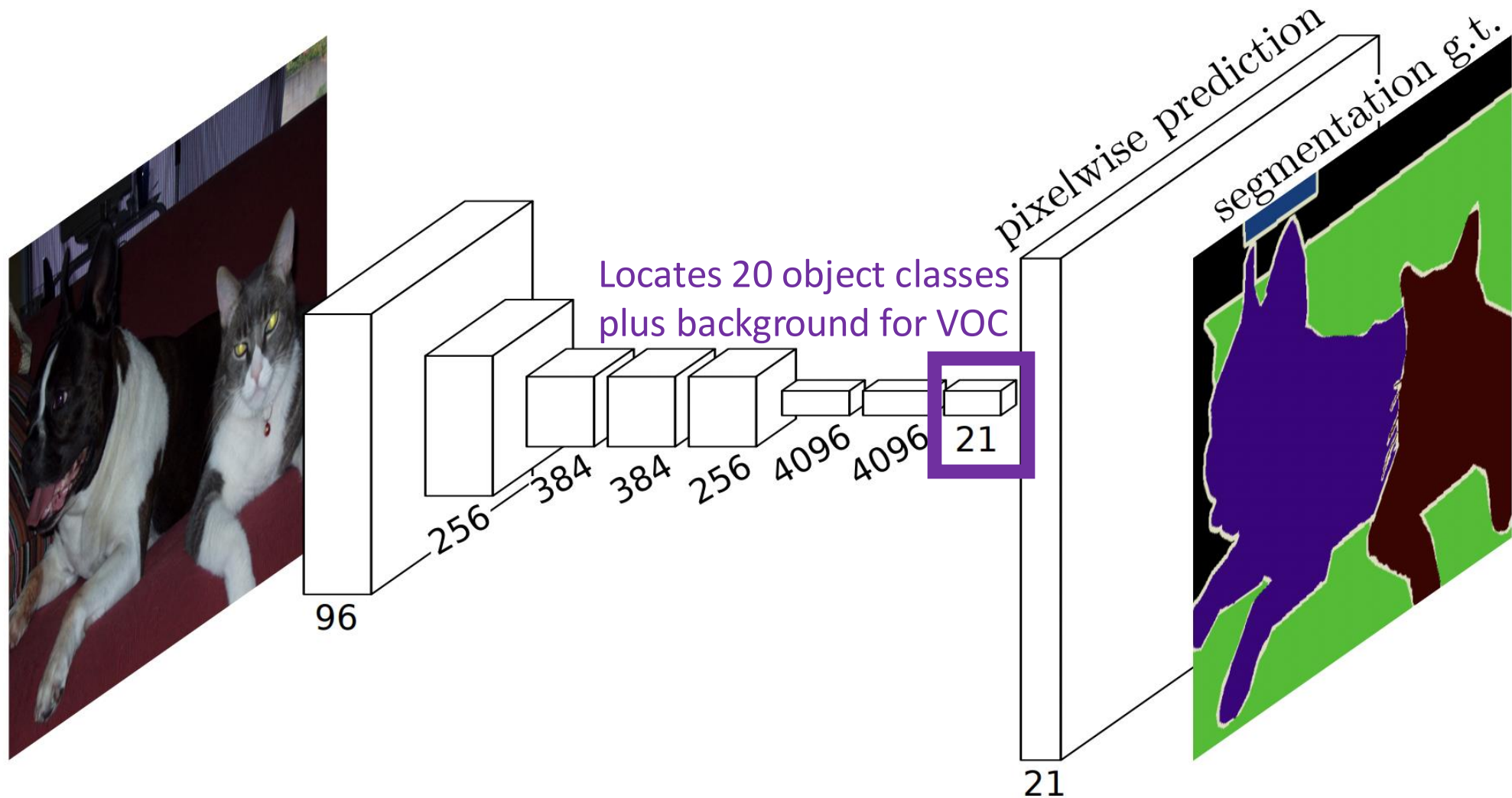


Architecture: Fully vs Convolution Layers



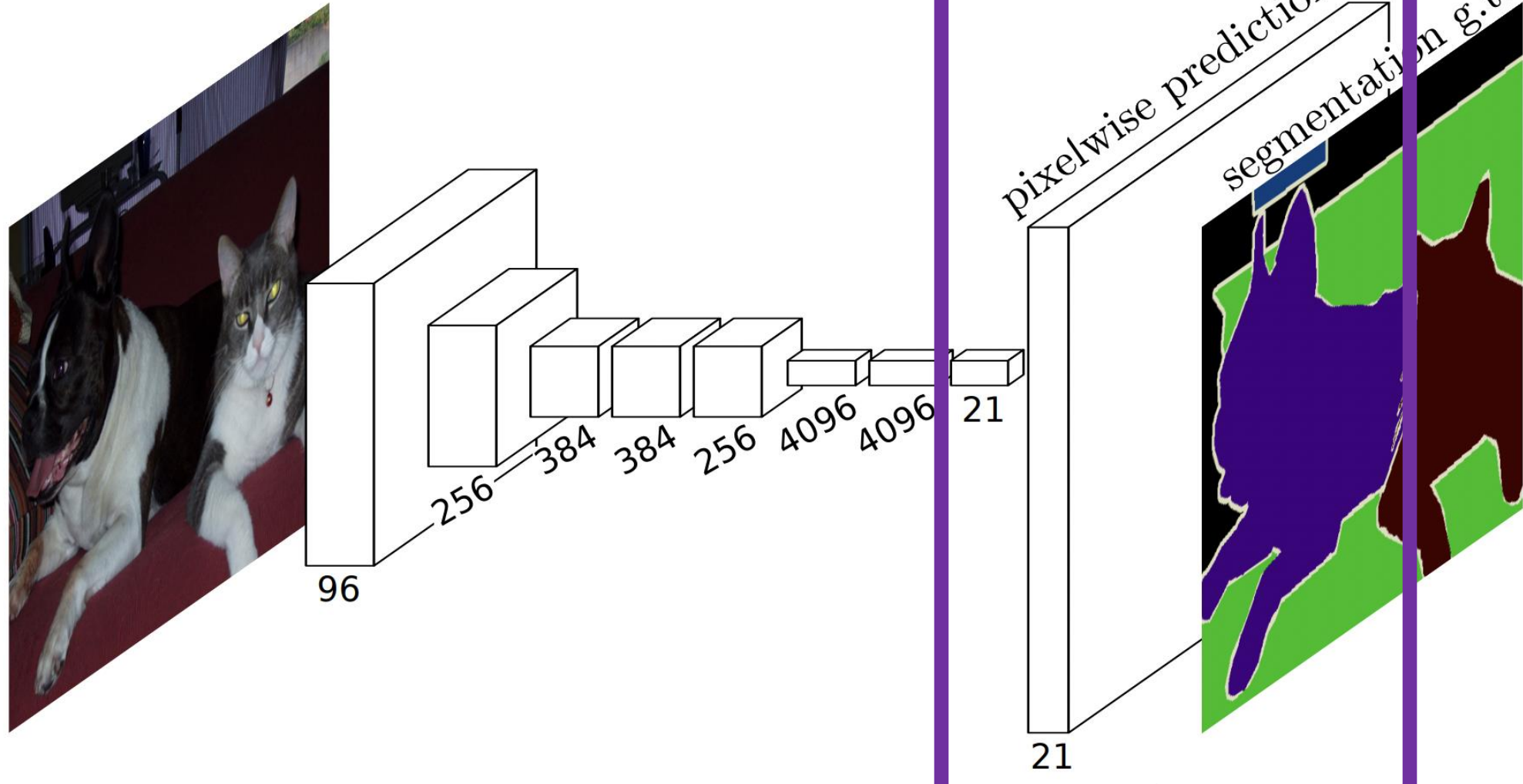
If convolutionizing ImageNet trained classifiers, how many classes would be predicted for each coarse region?

Architecture: Coarse Region Classification

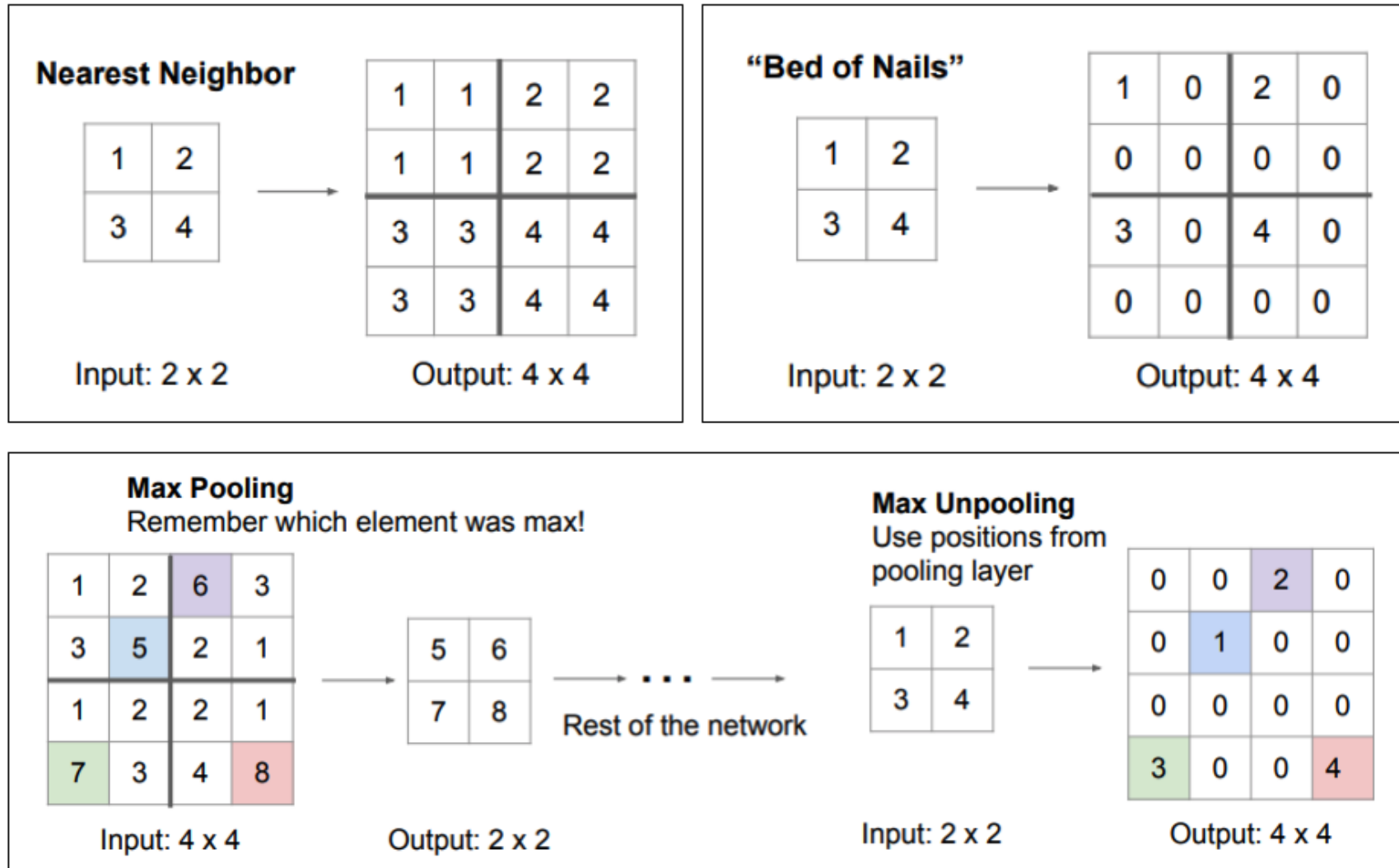


Architecture

Challenge: how to decode from coarse region classifications to per pixel classification?

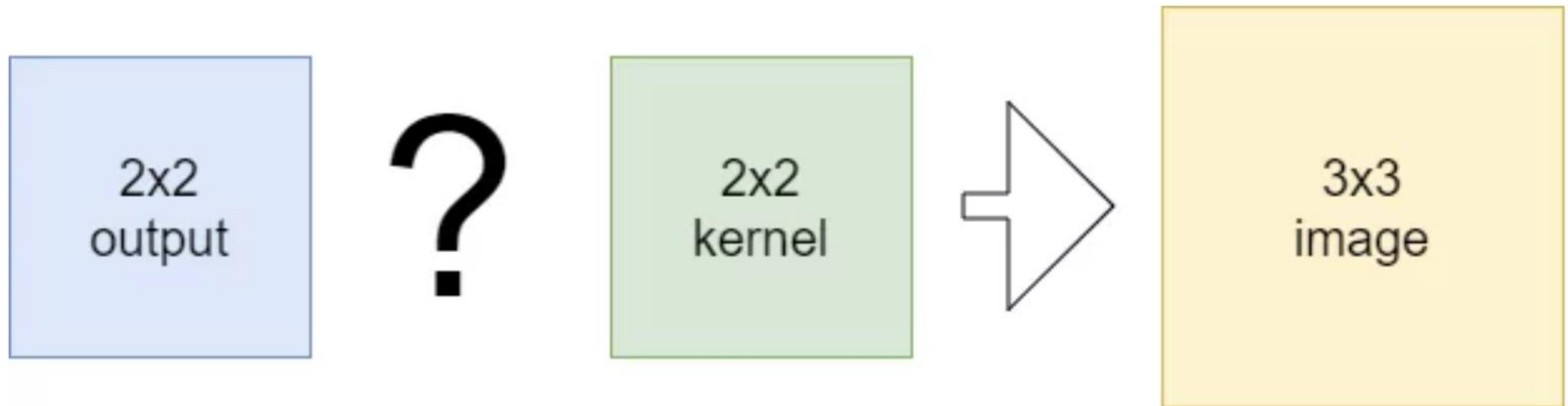


Architecture: Upsampling (Many Approaches)



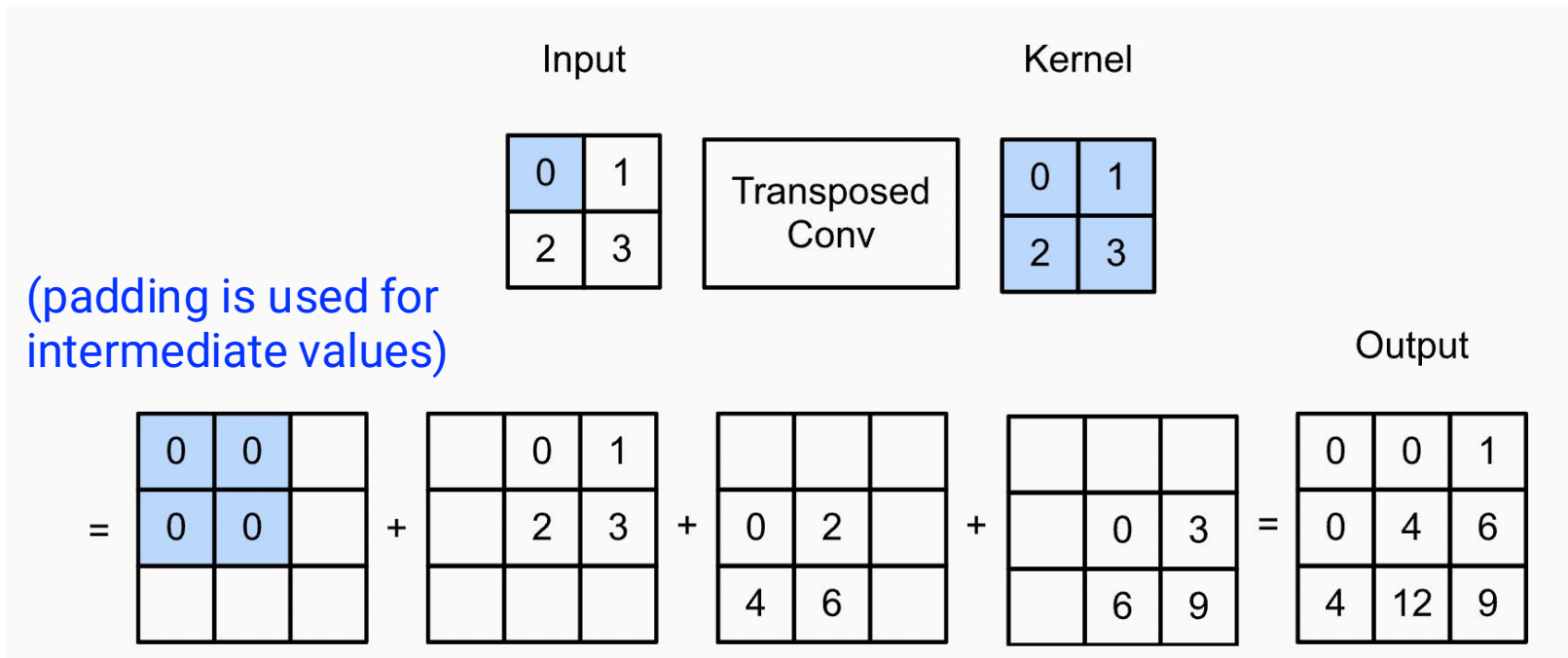
Architecture: Upsampling (Transposed Convolutional Layer)

- **Idea:** learn convolutional filters to upsample the coarse image with fractional sized steps
- Also called “fractional convolutional layer”, “backward convolution”, and, incorrectly, “deconvolution layer”, there are many implementations



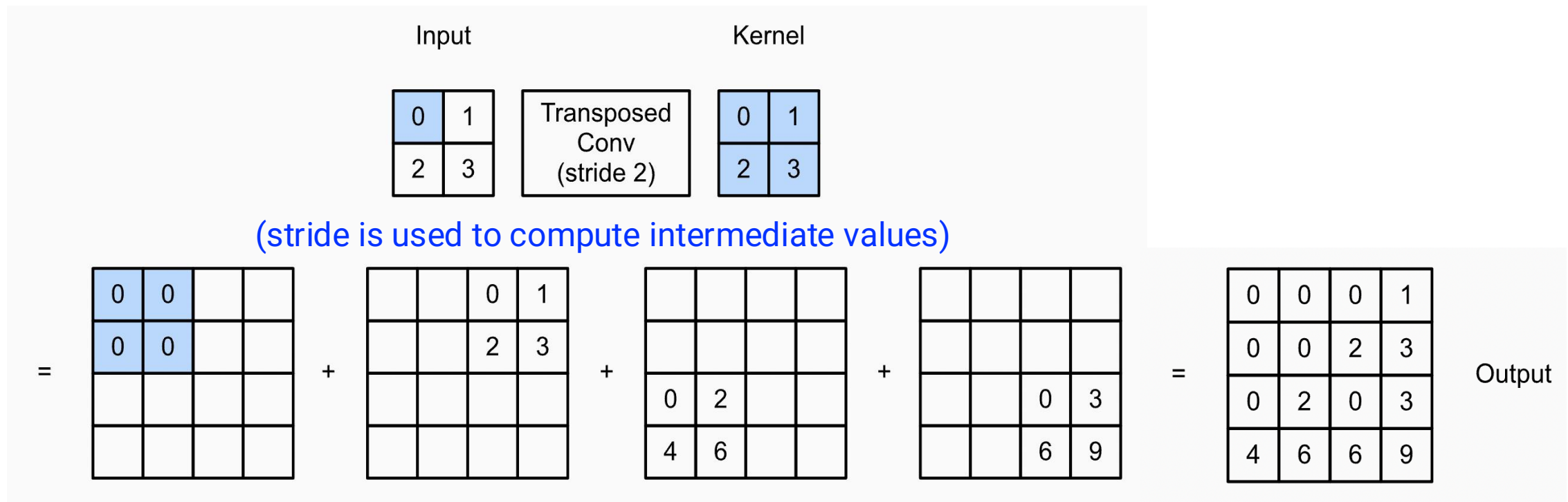
Architecture: Upsampling (Transposed Convolutional Layer)

- **Idea:** learn convolutional filters to upsample the coarse image with fractional sized steps
- Also called “fractional convolutional layer”, “backward convolution”, and, incorrectly, “deconvolution layer”, there are many implementations



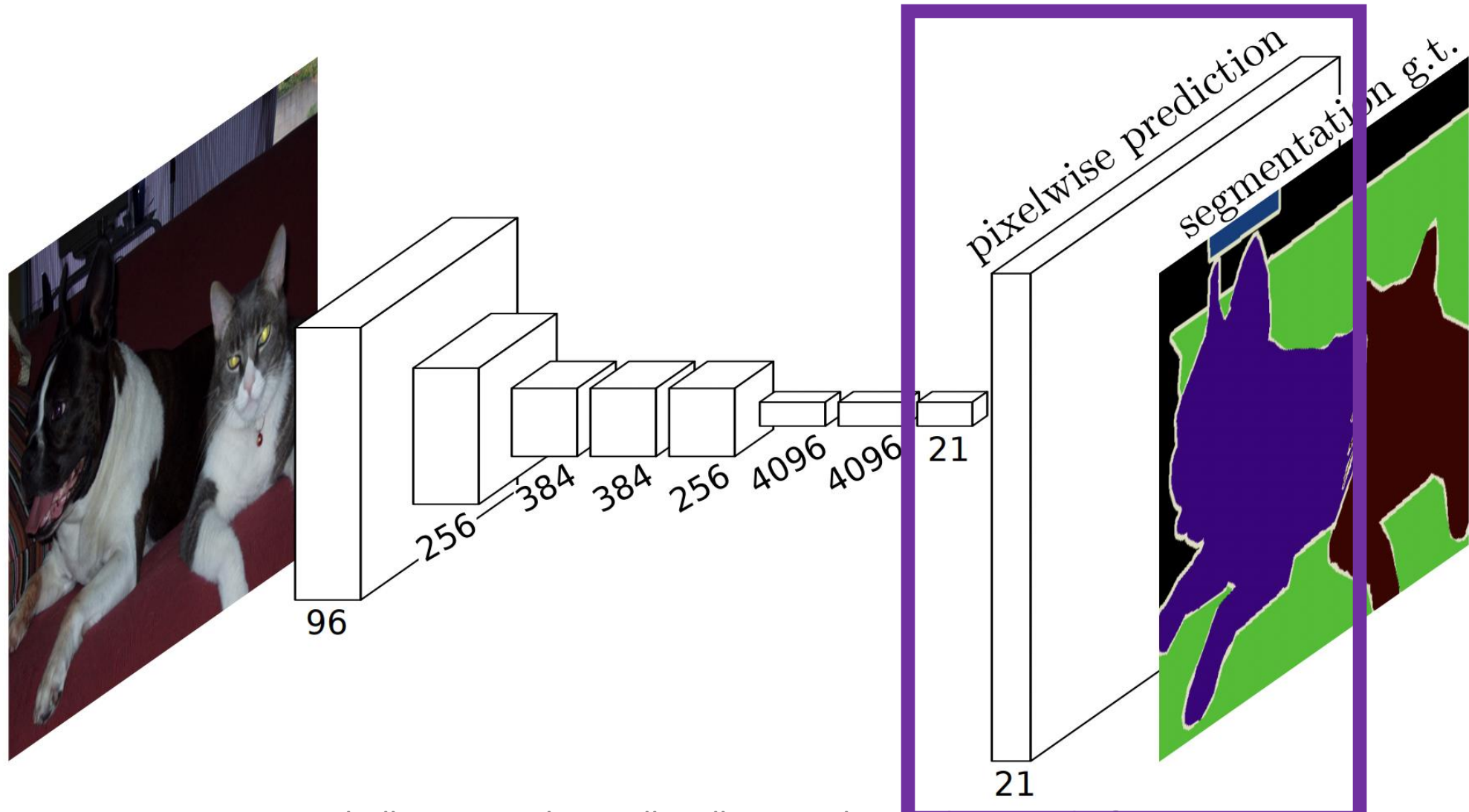
Architecture: Upsampling (Transposed Convolutional Layer)

- **Idea:** learn convolutional filters to upsample the coarse image with fractional sized steps
- Also called “fractional convolutional layer”, “backward convolution”, and, incorrectly, “deconvolution layer”, there are many implementations



Architecture

Next challenge: how to decode a **highly detailed** per pixel classification from the coarse region classifications?



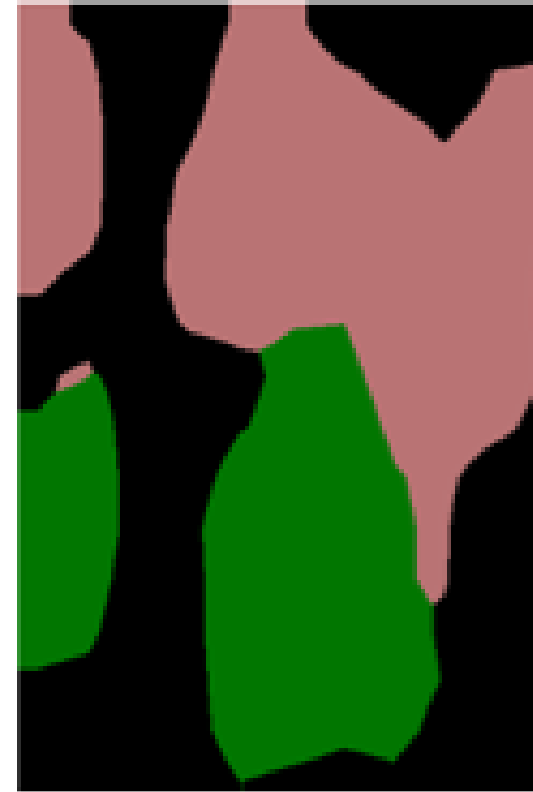
Architecture: Results

Next challenge: how to decode a **highly detailed** per pixel classification from the coarse region classifications?

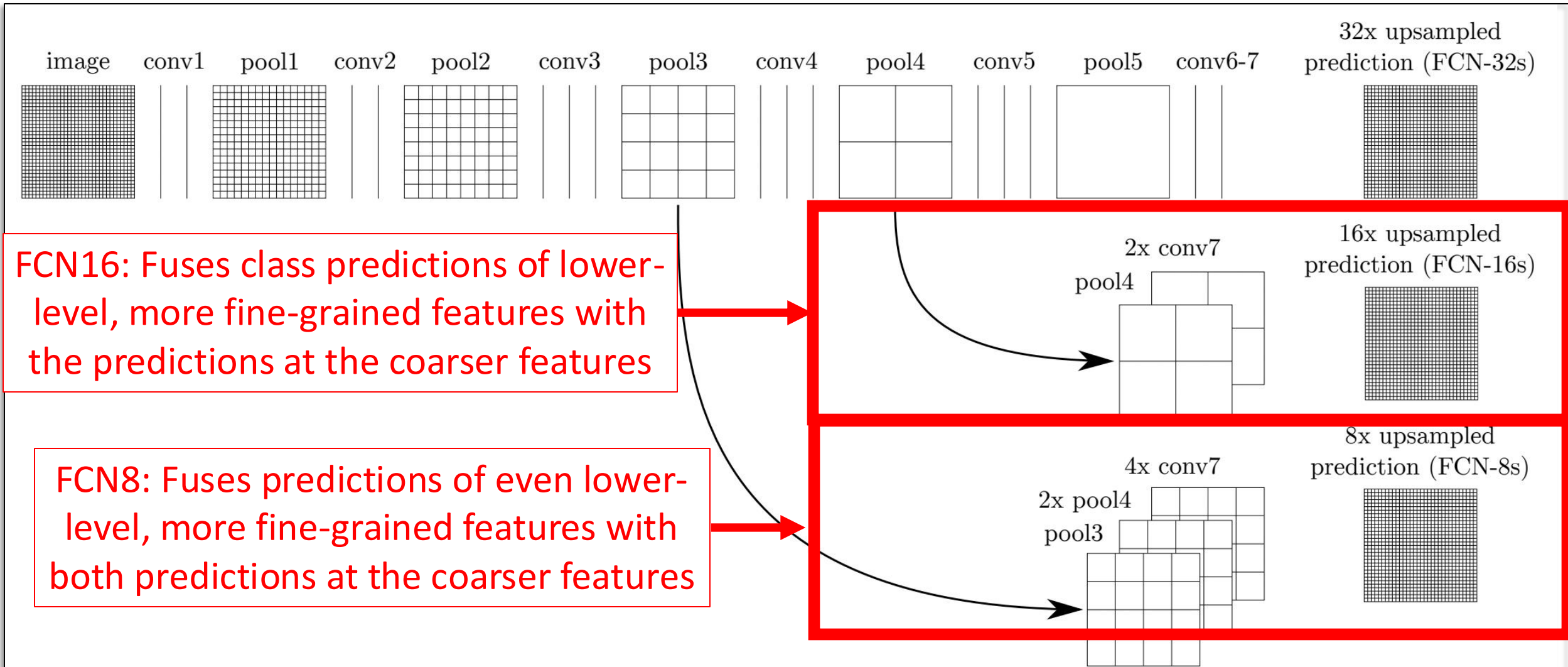
Ground truth target



Predicted segmentation

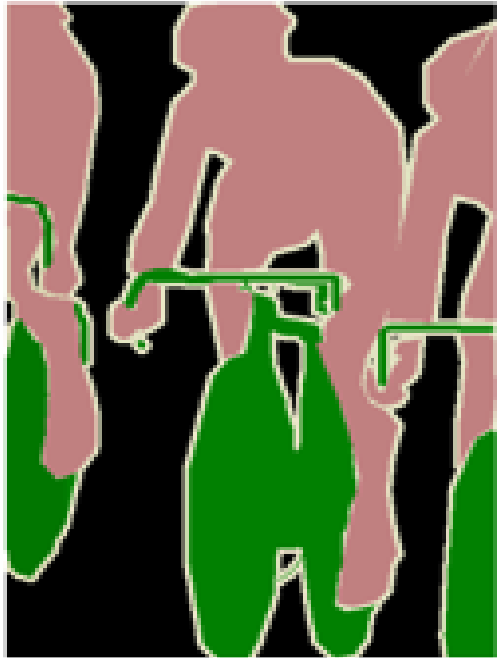


Architecture: Update to Use Skip Connections

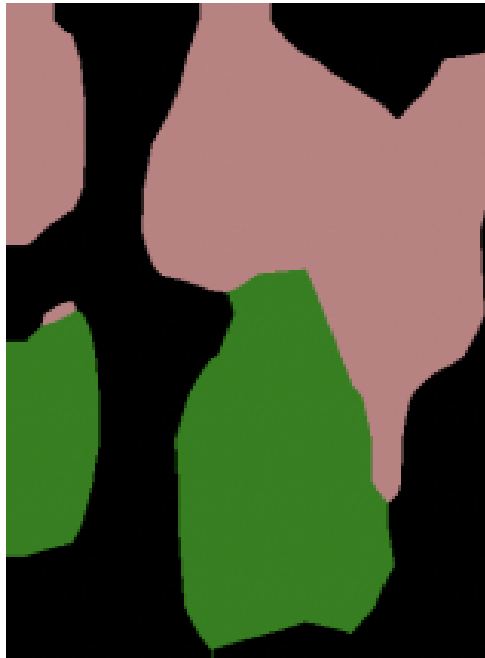


Architecture: Results

Ground truth target



FCN-32s



FCN-16s



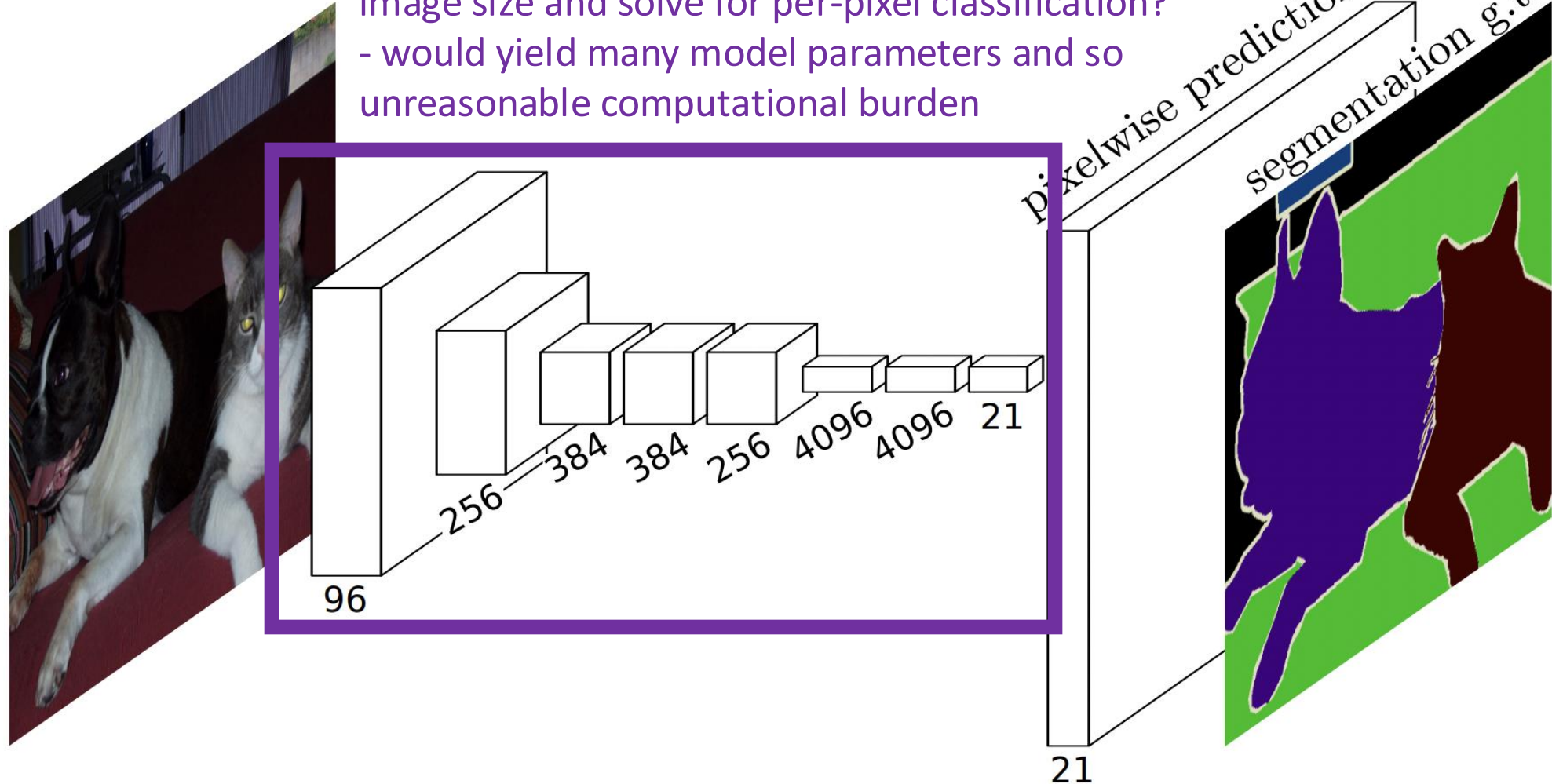
FCN-8s



Skip connections support capturing finer-grained details while retaining correct semantic information!

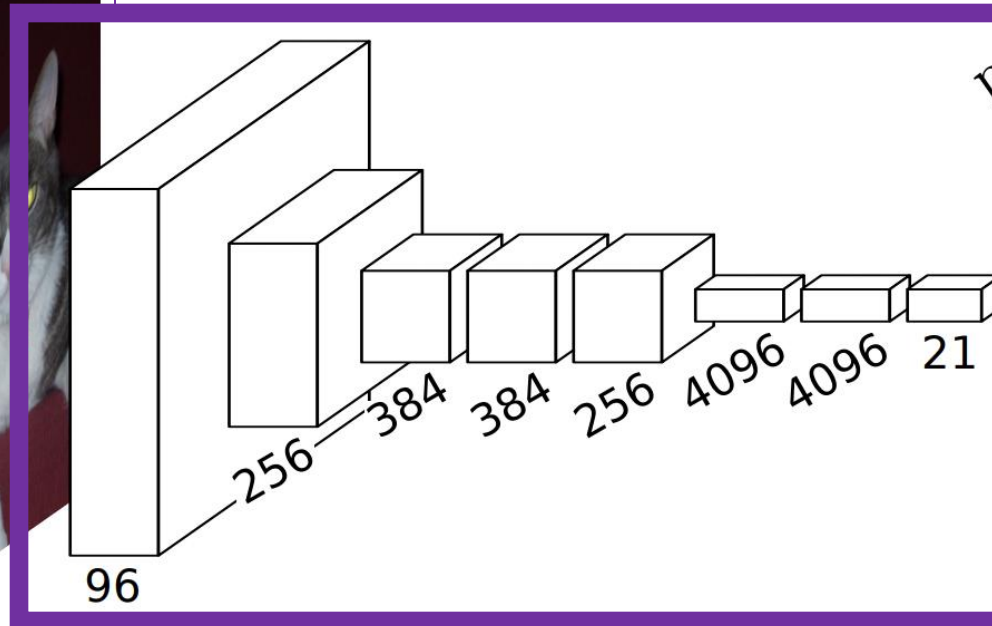
Architecture: Upsampling + Skip Connections

Seems complicated... why not instead preserve the image size and solve for per-pixel classification?
- would yield many model parameters and so unreasonable computational burden

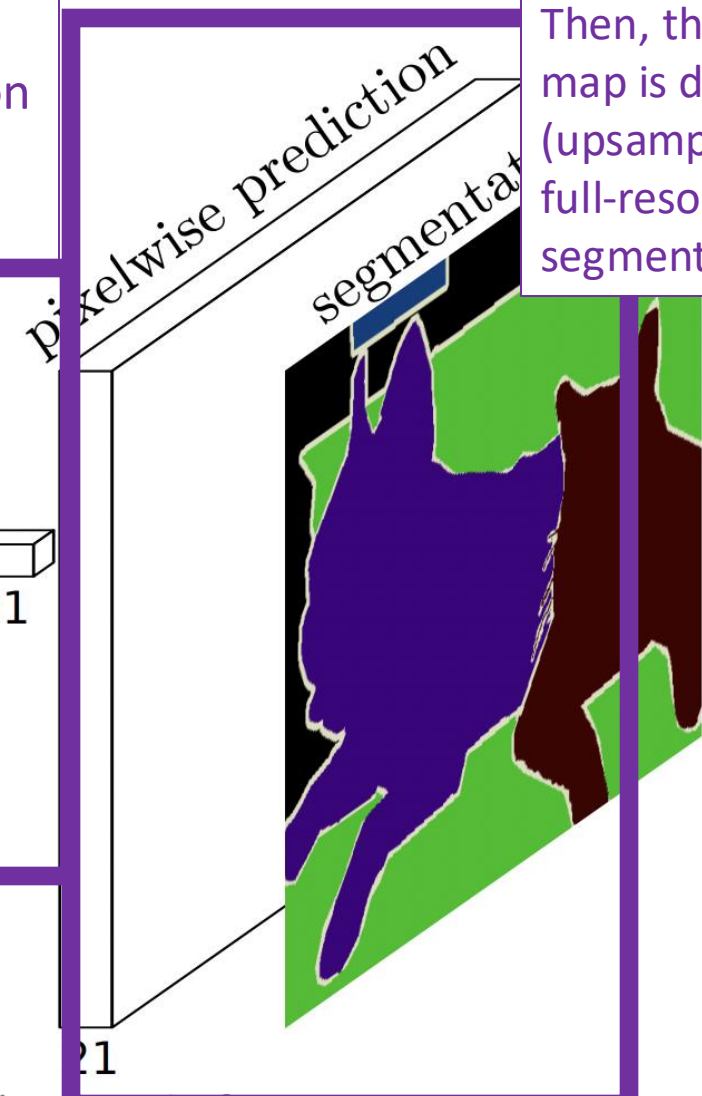


Architecture: Encoder Decoder Architecture

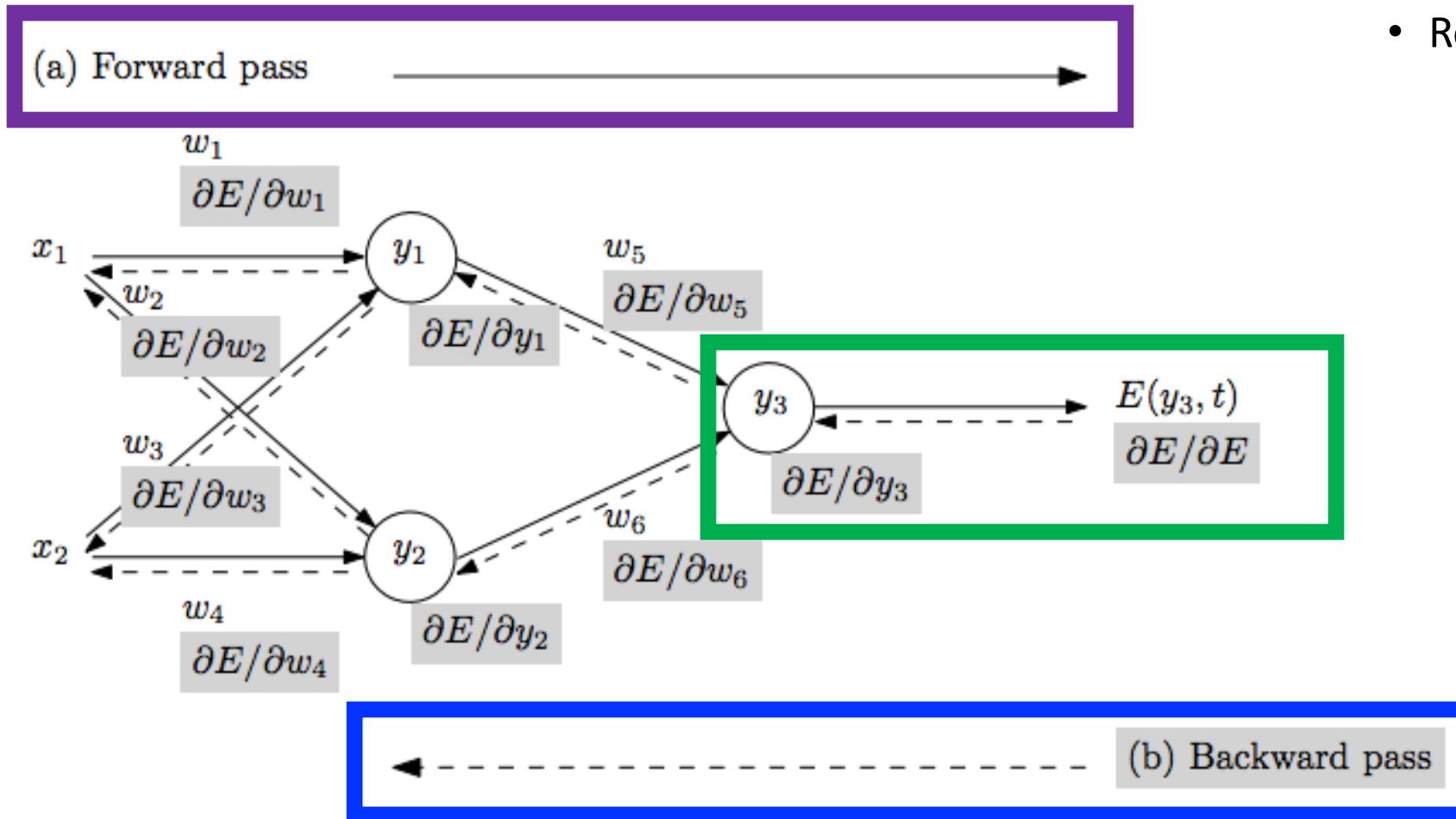
For efficiency, the image is encoded (downsampled) into a lower-resolution feature map that effectively discriminates between classes...



Then, the feature map is decoded (upsampled) into a full-resolution segmentation map.



Training: Took 3 days on 1 GPU




- Repeat until stopping criterion met:
 1. **Forward pass:** propagate training data through model to make prediction
 2. Quantify the dissatisfaction with a model's results on the training data
 3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
 4. Update each parameter using calculated gradients

Training: How Neural Networks Learn

- Repeat until stopping criterion met:
 1. **Forward pass:** propagate training data through model to make prediction
 2. Quantify the dissatisfaction with a model's results on the training data
 3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
 4. Update each parameter using calculated gradients

Sum across all pixels the distance between predicted and true distributions using cross entropy loss

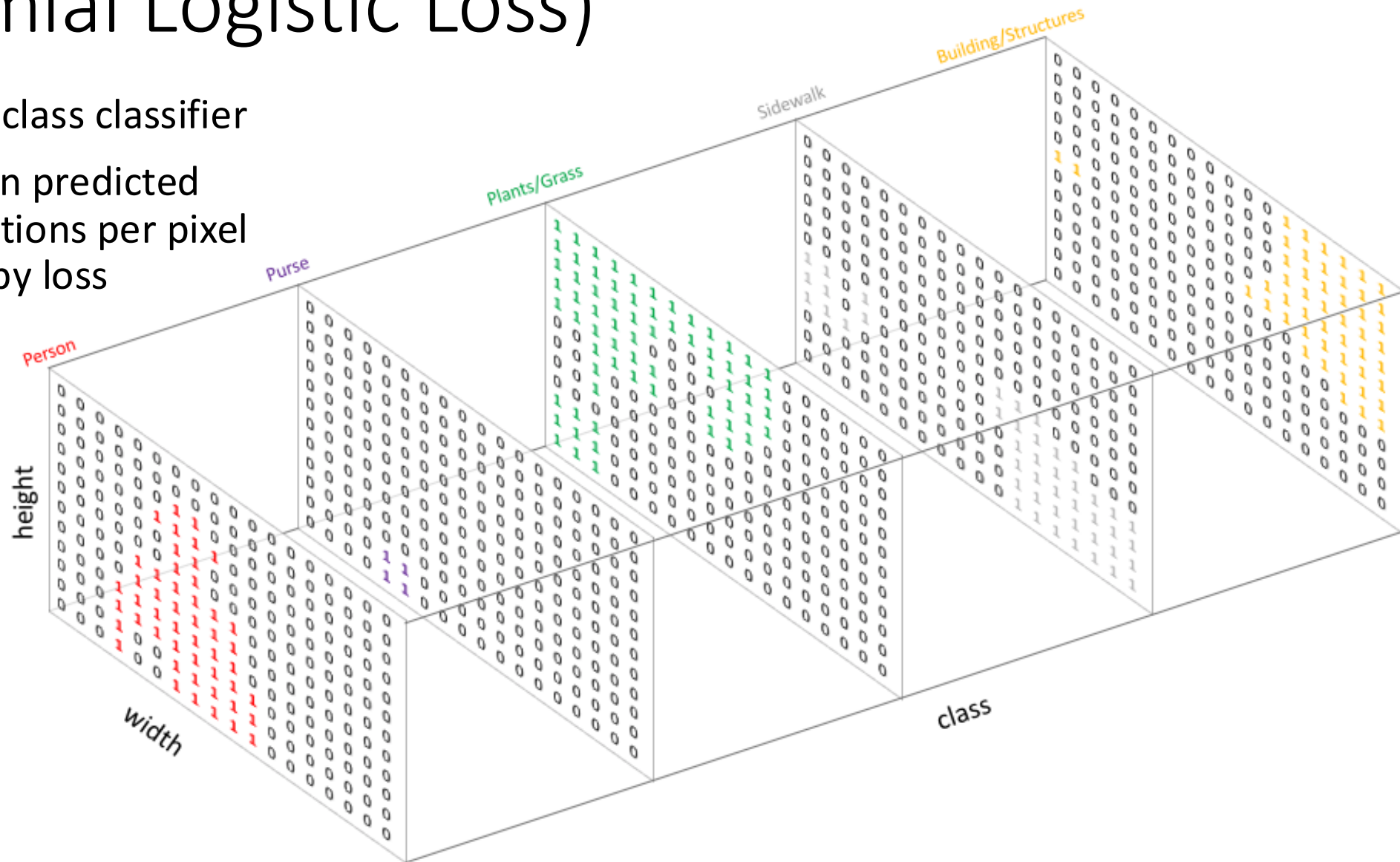


Sum of gradients for all pixels (acts like a minibatch)

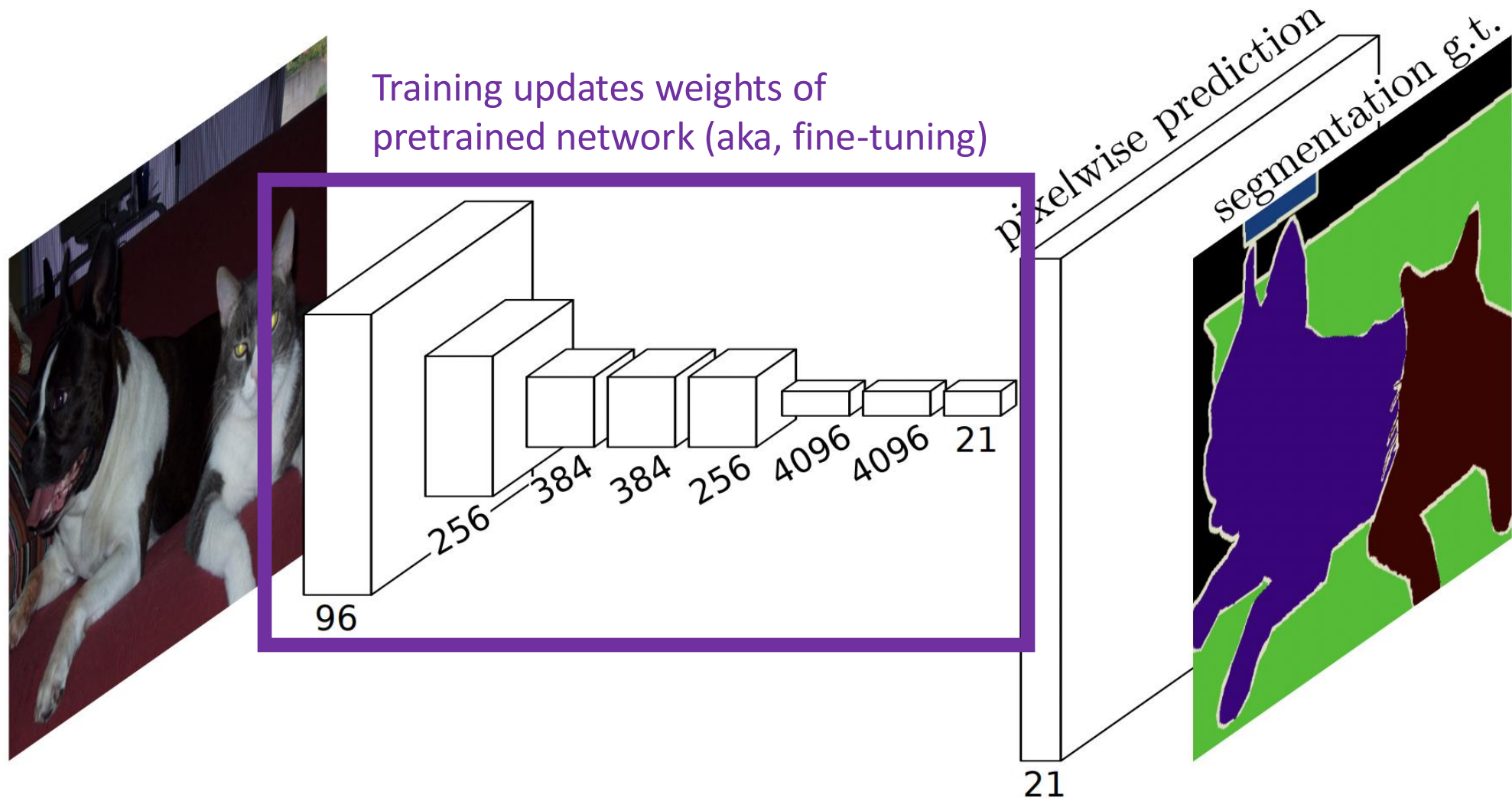


Training: Cross Entropy Loss (Multinomial Logistic Loss)

- e.g., assume a 5-class classifier
- Distance between predicted and true distributions per pixel with cross entropy loss



Architecture: Algorithm Training



Results

	mean IU VOC2011 test	mean IU VOC2012 test	inference time
R-CNN [12]	47.9	-	-
SDS [16]	52.6	51.6	~ 50 s
FCN-8s	62.7	62.2	~ 175 ms

Compared to existing methods, produces better results at a faster speed!

Semantic Segmentation: Today's Topics

- Motivation
- Datasets
- Evaluation metric
- Fully convolutional network
- Swin transformer
- Discussion (chosen by YOU 😊)

Why Swin Transformer?

Named after the proposed technique: **S**hifted **W**indows

Liu et al. Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows.
ICCV 2021.

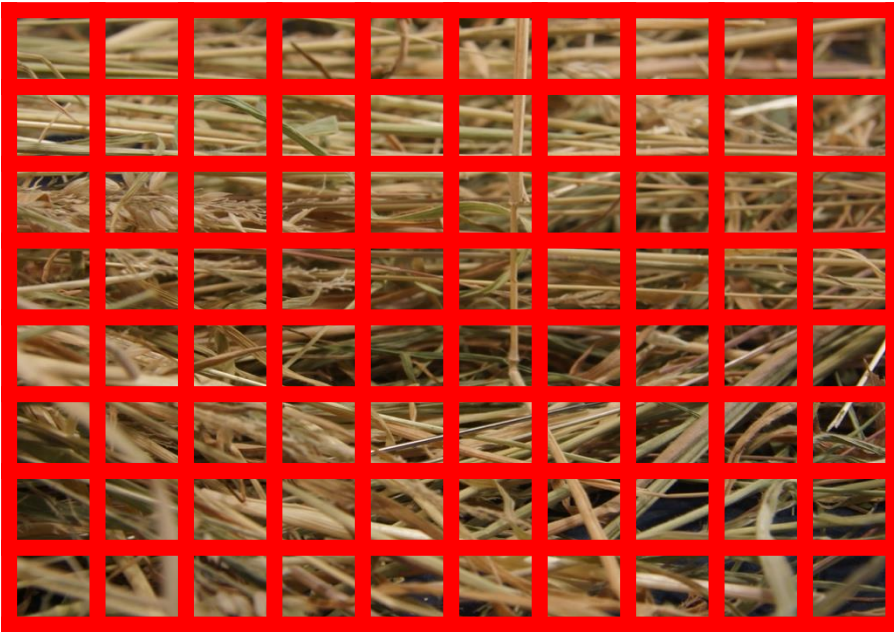
Novelty

- Demonstrates a transformer “backbone” can generalize to diverse vision tasks, with state-of-the-art results for semantic segmentation and object detection (aka – dense prediction problems) as well as strong results for image classification

Why ViT Is Inadequate for Dense Prediction

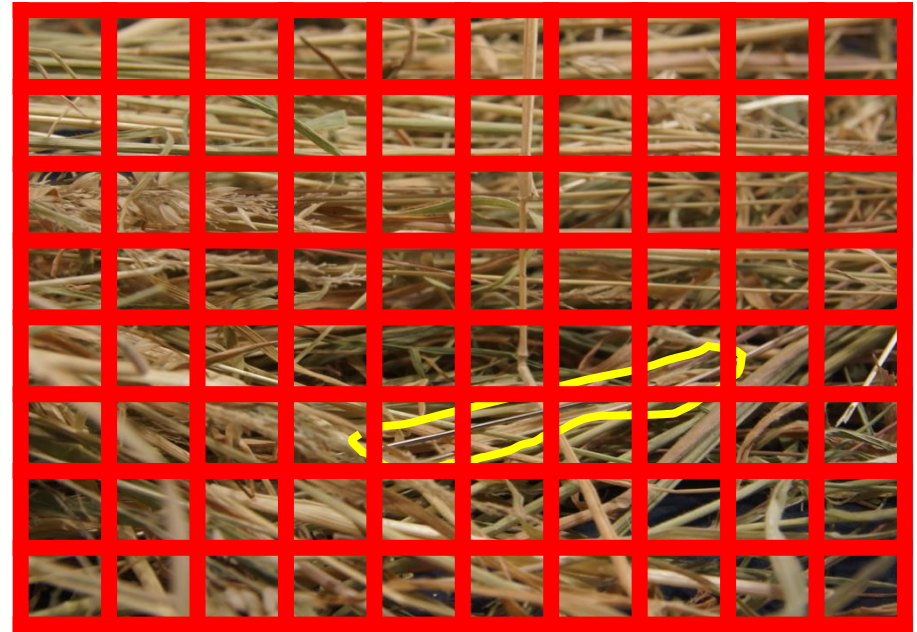
Image classification

- What **image** label is predicted?
- “Big” patches are sufficient



Object detection/Semantic segmentation

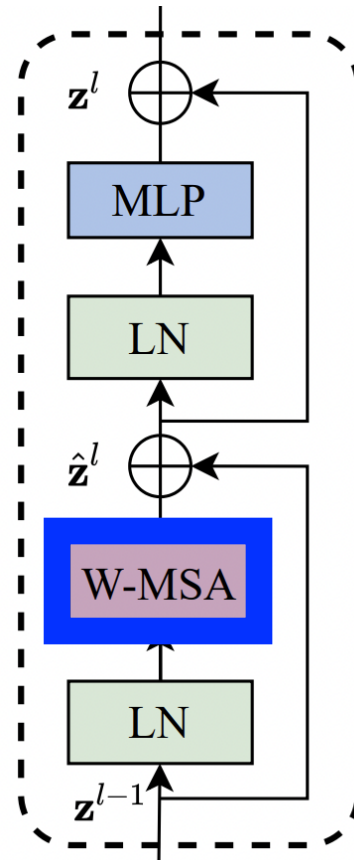
- What **pixel** label(s) are predicted?
- “Big” patches may be insufficient



Issue: quadratic expense of self-attention necessitated 16 x 16 patches, but this can be too large for pixel-level predictions (e.g., locating needle in a haystack)

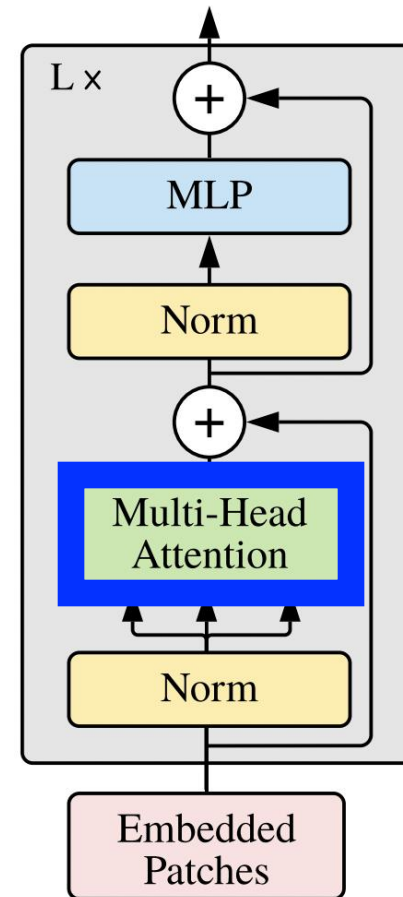
Key Idea of Swin: Modify Self-Attention Module

Swin Transformer



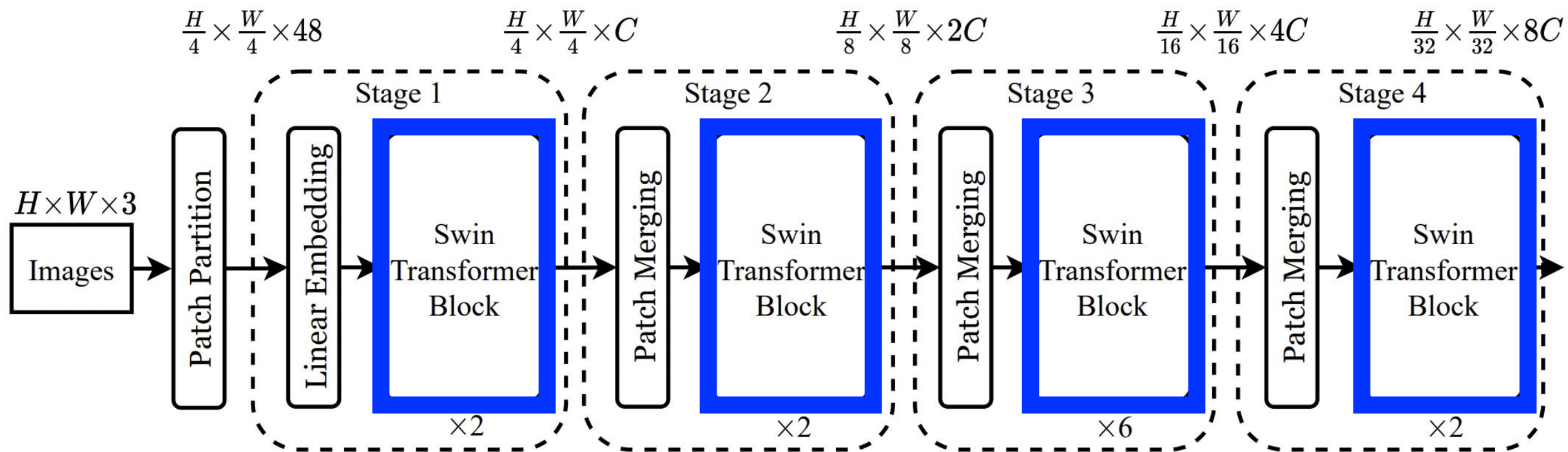
Liu et al. ICCV 2021.

ViT



Dosovitskiy et al. ICLR 2021.

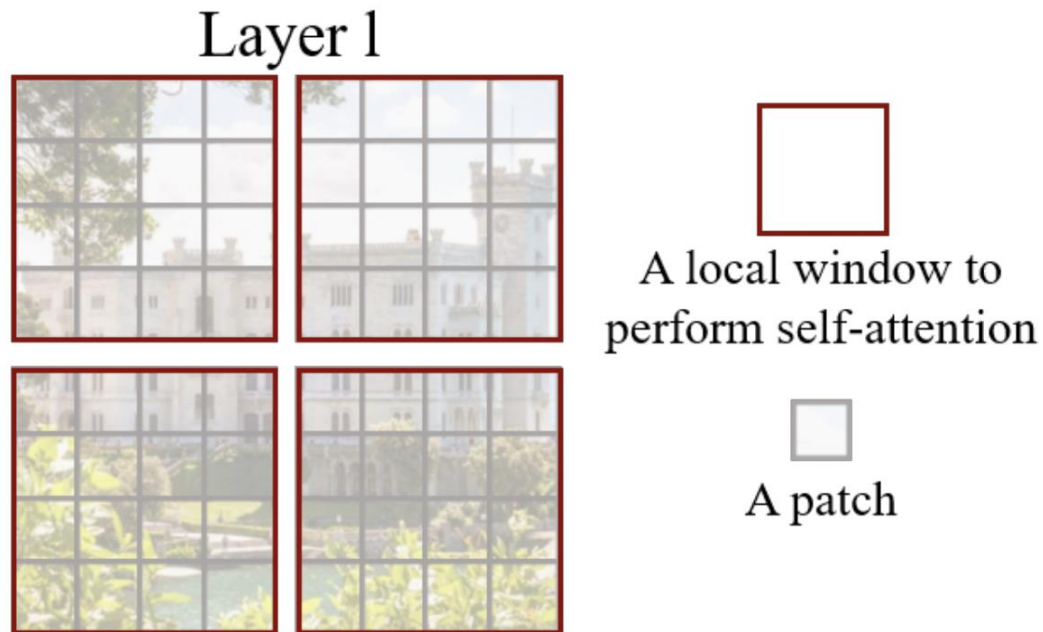
Architecture



Contains a series of modified self-attention modules

Key Idea: Modified Self-Attention Module

Applies self-attention only between the **fixed number of patches in each window** to capture fine-grained details (i.e., limited to **local context**)

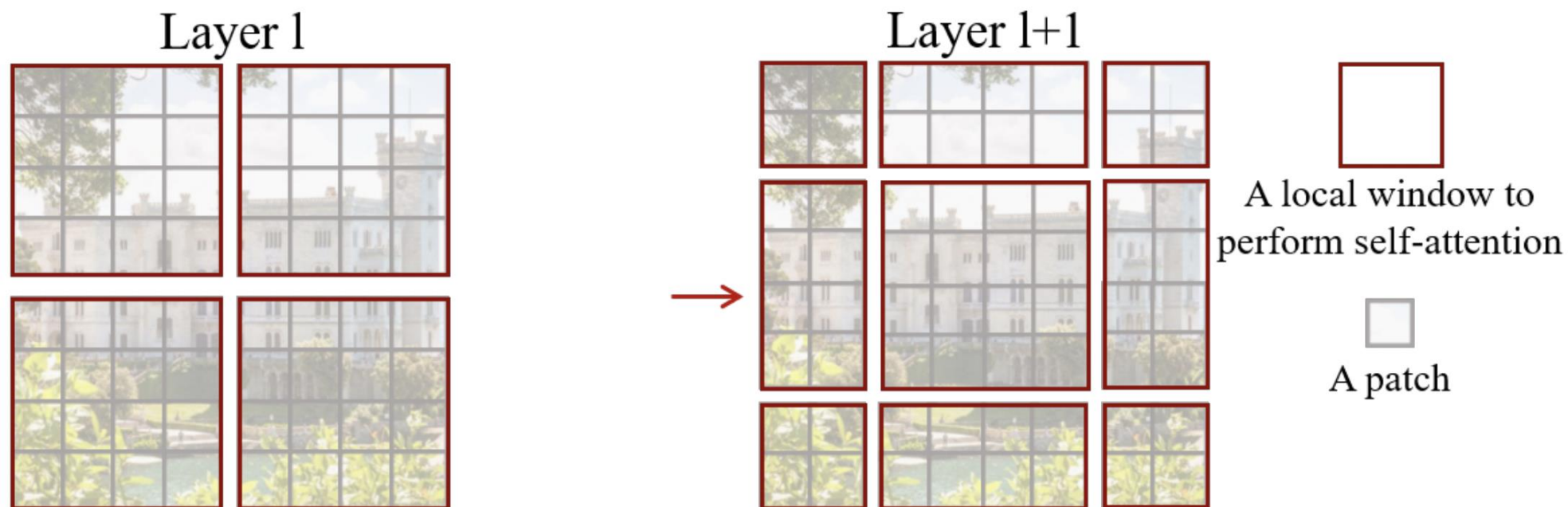


What is the computational complexity?
- **Linear** based on fixed patch number chosen per window rather than quadratic based on number of input patches

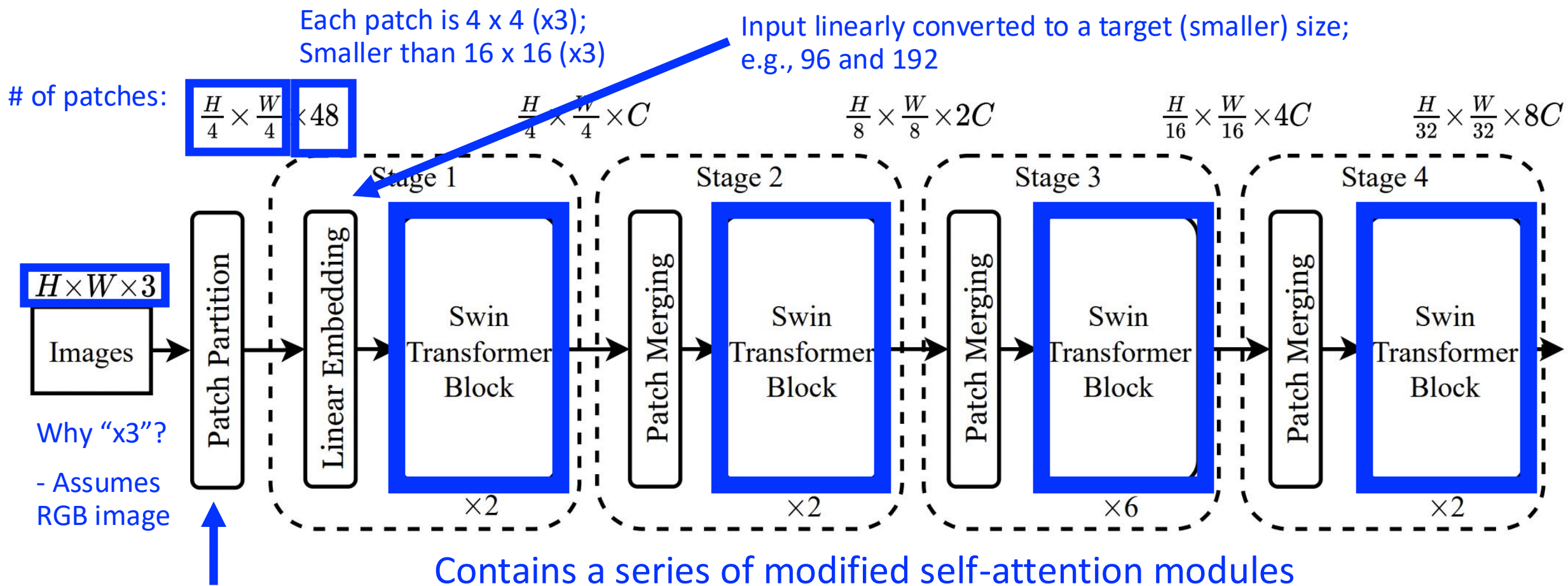
Key Idea: Modified Self-Attention Module

Applies self-attention only between the **fixed number of patches in each window** to capture fine-grained details (i.e., limited to **local context**)

In each subsequent layer, windows shifted to infuse **global context** by enabling communication between previously non-communicative neighboring patches



Architecture



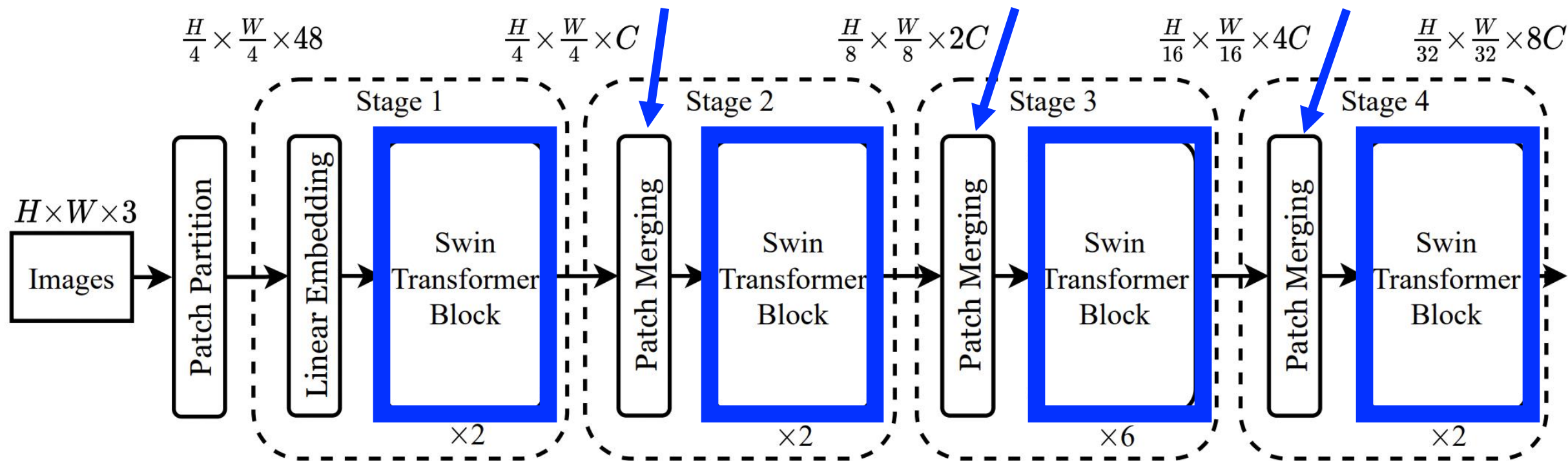
Why "x3"?

- Assumes RGB image

How many image pixels are in each image patch?

Architecture

Neighboring patches merged into increasingly bigger patches (mimics convolutional layers); this hierarchical design also increases global context to better support visual content at different scales! (output feature maps match resolution of common CNNs, e.g., VGG & ResNet)



Contains a series of modified self-attention modules at different resolutions

Dense Prediction: State-of-the Art Results

Four **object detection** algorithms tested on COCO 2017 with three “backbone” sources:

- ResNe(X)t
- DeiT
- Swin: was consistently top-performer

UperNet **semantic segmentation** algorithm tested on ADE20K with two “backbone” sources:

- DeiT
- Swin: was consistently top-performer

Semantic Segmentation: Today's Topics

- Motivation
- Datasets
- Evaluation metric
- Fully convolutional network
- Swin transformer
- Discussion (chosen by YOU 😊)

Semantic Segmentation: Today's Topics

- Motivation
- Datasets
- Evaluation metric
- Fully convolutional network
- Swin transformer
- Discussion (chosen by YOU 😊)

The image features a central area with a radial gradient background, transitioning from a lighter center to a darker outer edge. This central area is framed by a film strip border, consisting of a dark grey outer line and a series of white rectangular sprocket holes along the top and bottom edges. The text "The End" is centered within the gradient area.

The End