

# Object Detection – Part 1

**Danna Gurari**

University of Colorado Boulder  
Fall 2023



# Review

- Last lecture:
  - Scene Classification Problem and Applications
  - Scene Classification: Datasets and Evaluation Metrics
  - Scene Classification Models: Deep Features and Fine-Tuning
  - Attribute Classification: Problem, Applications, and Datasets
  - Student-led Lecture Overview
- Assignments (Canvas)
  - Reading assignment was due earlier today
  - Next reading assignments due Wednesday and next Monday
- Questions?

# Object Detection: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metric
- Overview of object detection algorithms and baseline (R-CNN)

# Object Detection: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metric
- Overview of object detection algorithms and baseline (R-CNN)

# Problem Definition

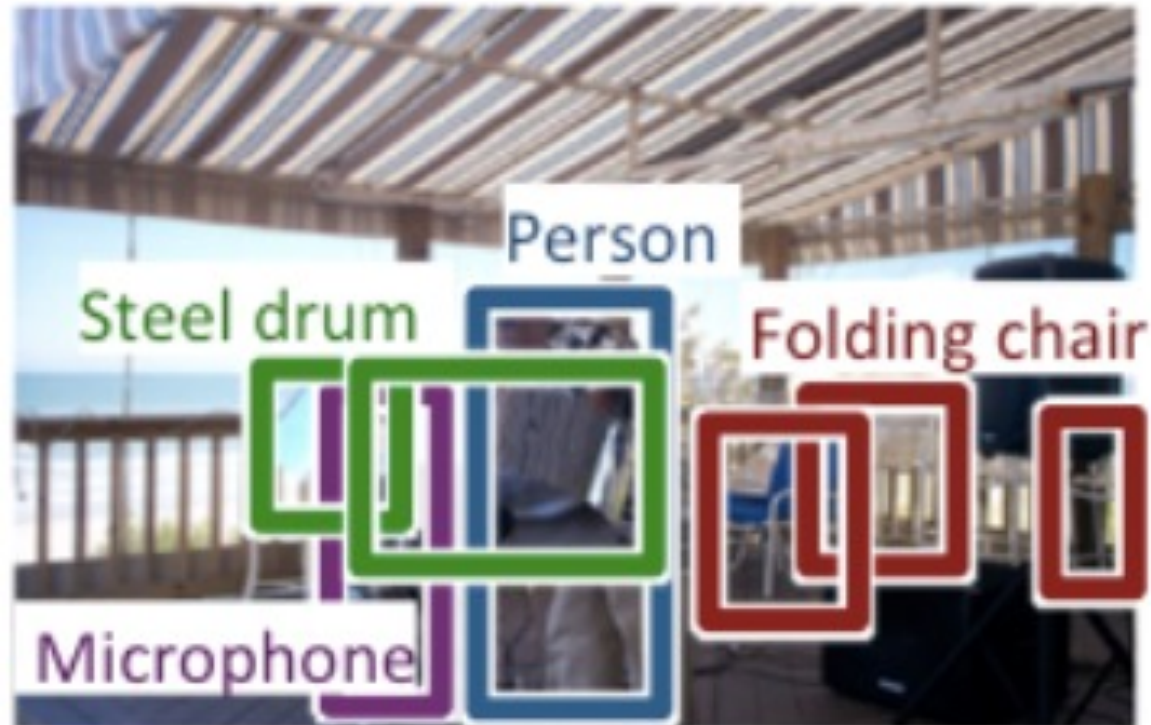
- Localize with a bounding box object(s) of interest



## Focus for today's lecture

# Problem: Semantic Object Detection

- Localize with a bounding box every instance of an object from pre-specified categories

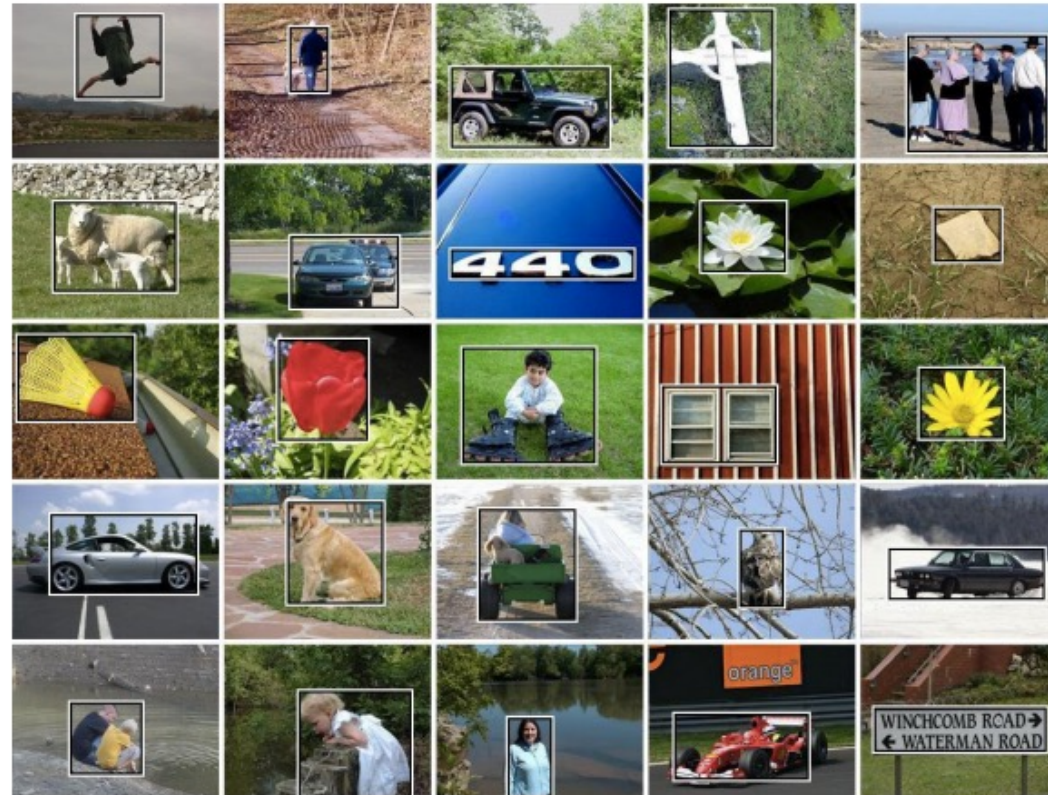


[Russakovsky et al; IJCV 2015]

## A reasonably solved problem

# Problem: Salient Object Detection

- Localize with a bounding box the salient object(s)



[Liu et al; CVPR 2007]

# Object Detection vs Object Recognition

“What is the difference between (semantic) object detection and object recognition?”



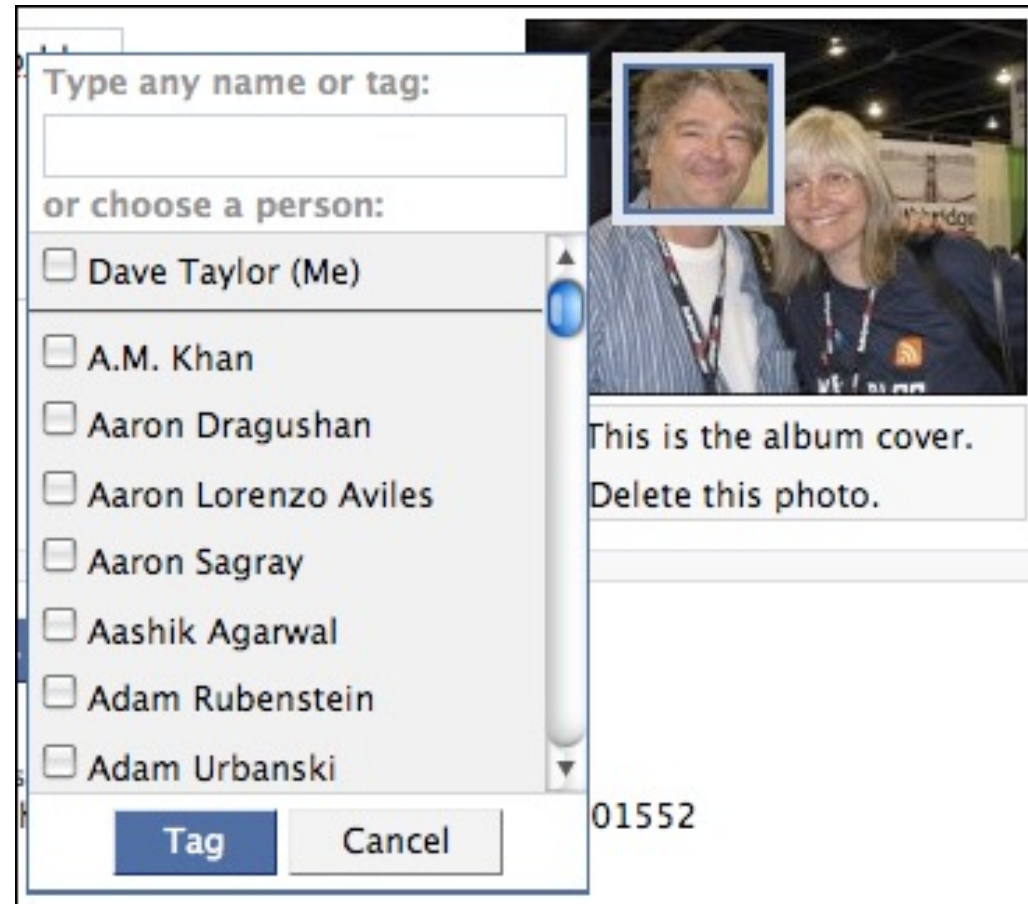
- Must learn appearance of object rather than only its image context; e.g., giraffe



# Object Detection: Today's Topics

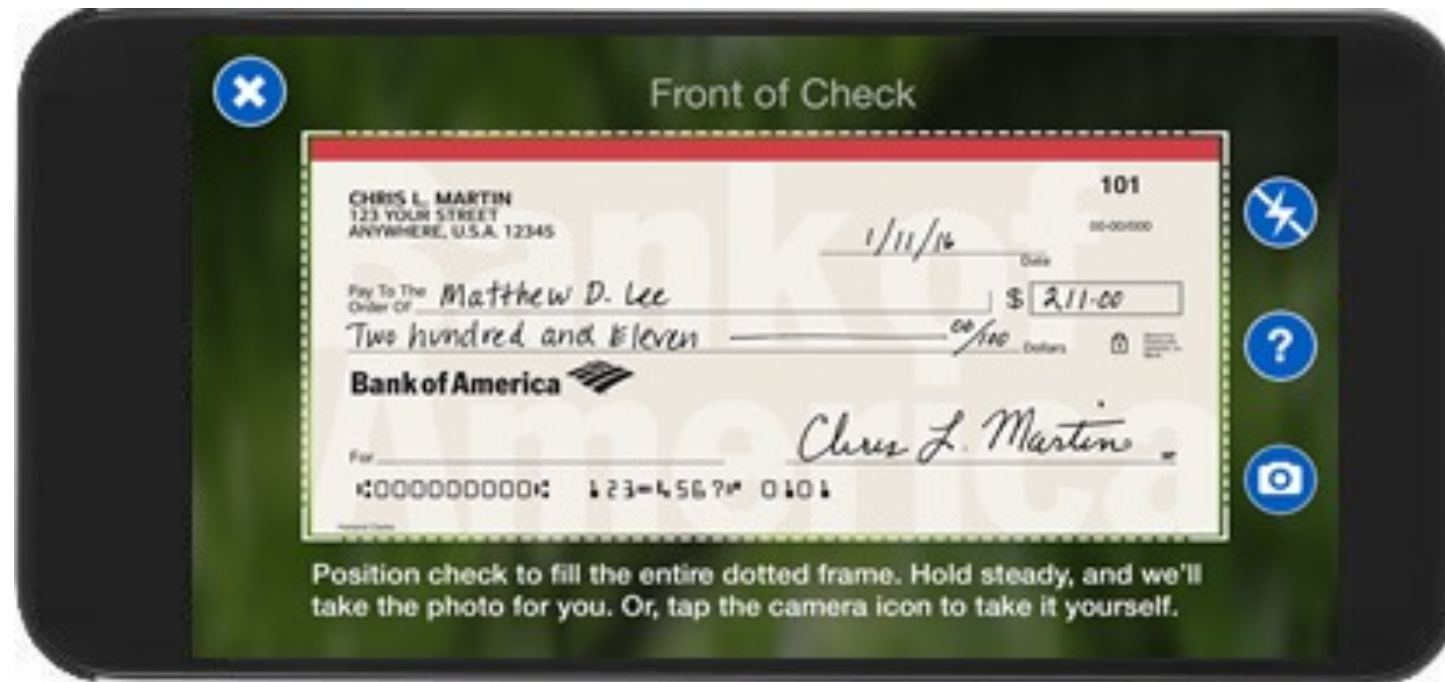
- Problem
- Applications
- Datasets
- Evaluation metric
- Overview of object detection algorithms and baseline (R-CNN)

# Social Media



Face detection  
(e.g., Facebook)

# Banking



Mobile check deposit  
(e.g., Bank of America)

# Transportation



License Plate Detection (e.g., AllGoVision)

# Construction Safety



Pedestrian Detection  
(e.g., Blaxtair)

<http://media.brintex.com/Occurrence/121/Brochure/3435/brochure.pdf>

# Counting



Counting Fish (e.g., SalmonSoft)  
[http://www.wecountfish.com/?page\\_id=143](http://www.wecountfish.com/?page_id=143)



Business Traffic Analytics

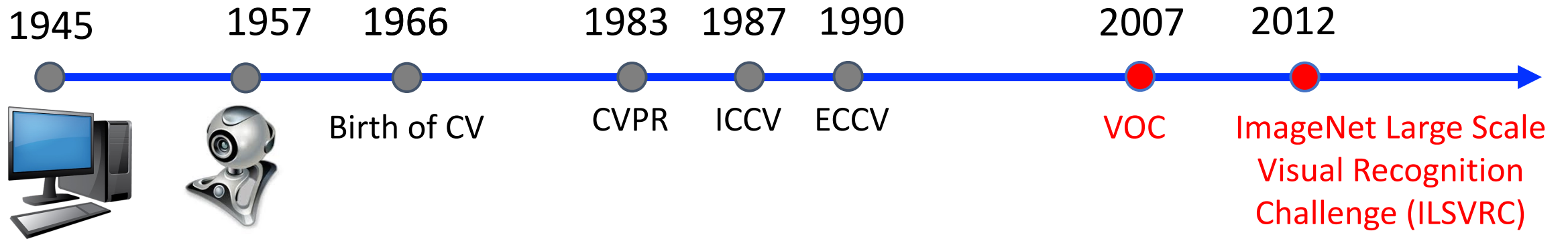
Can you think of any other  
potential applications?

# Object Detection: Today's Topics

- Problem
- Applications
- **Datasets**
- Evaluation metric
- Overview of object detection algorithms and baseline (R-CNN)



# Object Detection Datasets



# VOC

## 1. Category Selection

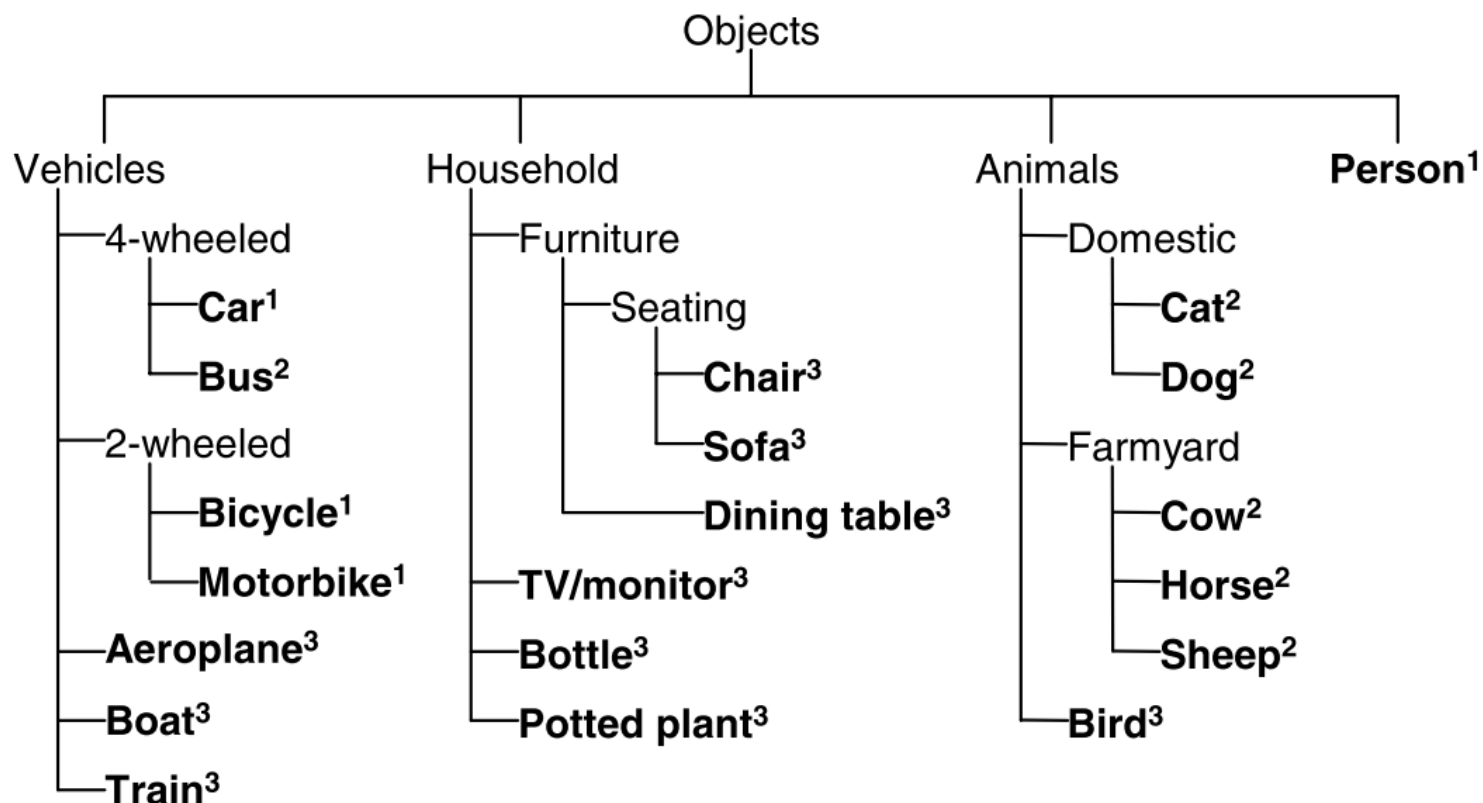
- 20 categories chosen:

1) Initial 4 categories stem from existing dataset

2) 2006: added 6 classes

3) 2007: added 10 classes

- Additional categories provide a broader domain and finer-grained categories, including visually similar things



*(superscript indicates year of inclusion in the challenge: 2005<sup>1</sup>, 2006<sup>2</sup>, 2007<sup>3</sup>)*

# 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
- Additional categories provide a broader domain and finer-grained categories, including visually similar things

## 2. Image Collection

- 500,000 images retrieved from Flickr by querying with a number of keywords

*(many query 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
- Additional categories provide a broader domain and finer-grained categories, including visually similar things

## 2. Image Collection

- 500,000 images retrieved from Flickr by querying with a number of keywords

## 3. Image Verification + Image Annotation

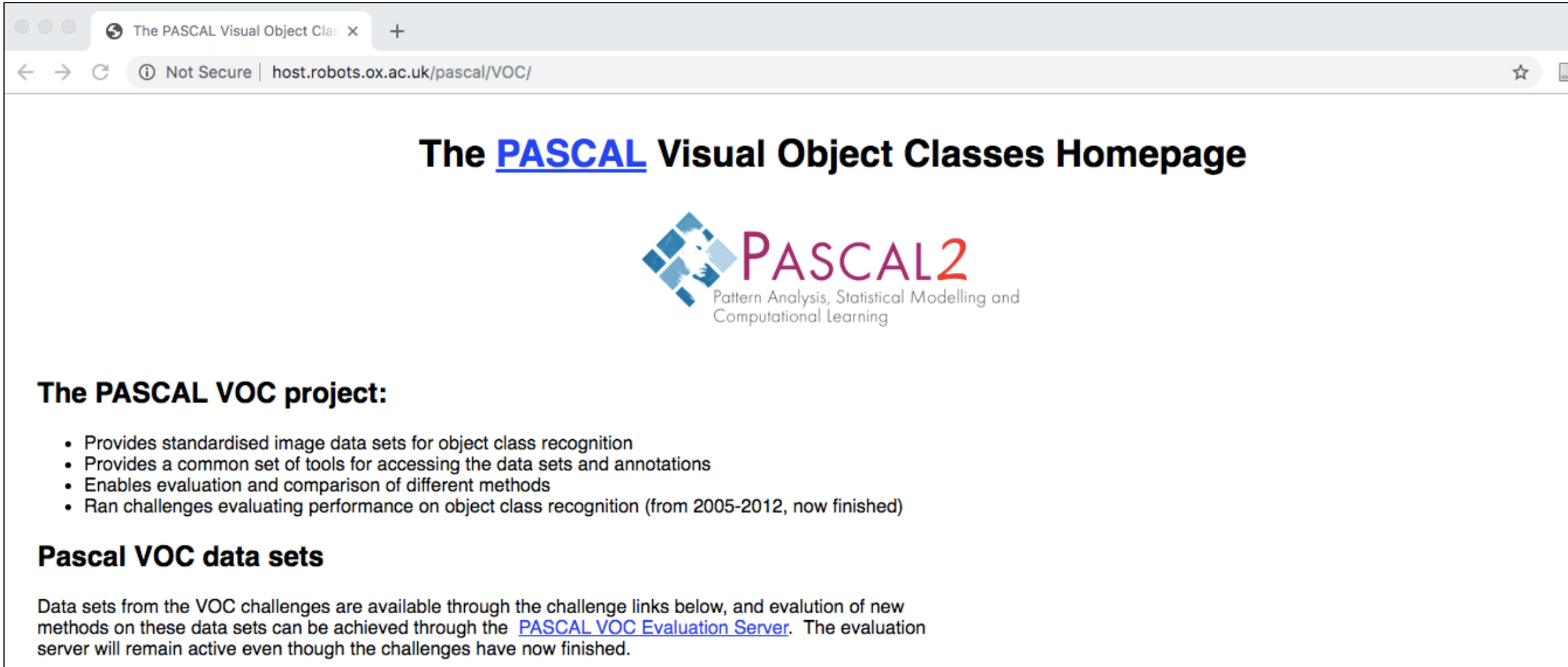
- University of Leeds annotation party to recruit annotators
- Annotation guidelines & real-time assistance
- Review of every annotation
- Annotate only “minority” classes at end of party to increase the count of them

# VOC Guidelines:

What are potential limitations of this task design for resulting datasets (and so algorithms developed with such datasets)?

<b>What to label</b>	<p>All objects of the defined categories, unless:</p> <ul style="list-style-type: none"><li>• you are unsure what the object is.</li><li>• the object is very small (at your discretion).</li><li>• less than 10-20% of the object is visible, <i>such that you cannot be sure what class it is</i>. e.g. if only a tyre is visible it may belong to car or truck so cannot be labelled car, but feet/faces can only belong to a person.</li></ul> <p>If this is not possible because too many objects, mark image as bad.</p>
<b>Viewpoint</b>	<p>Record the viewpoint of the 'bulk' of the object e.g. the body rather than the head. Allow viewpoints within 10-20 degrees.</p> <p>If ambiguous, leave as 'Unspecified'. Unusually rotated objects e.g. upside-down people should be left as 'Unspecified'.</p>
<b>Bounding box</b>	<p>Mark the bounding box of the visible area of the object (<i>not</i> the estimated total extent of the object).</p> <p>Bounding box should contain all visible pixels, except where the bounding box would have to be made excessively large to include a few additional pixels (&lt;5%) e.g. a car aerial.</p>
<b>Truncation</b>	<p>If more than 15-20% of the object lies outside the bounding box mark as Truncated. The flag indicates that the bounding box does not cover the total extent of the object.</p>
<b>Occlusion</b>	<p>If more than 5% of the object is occluded within the bounding box, mark as Occluded. The flag indicates that the object is not totally visible within the bounding box.</p>
<b>Image quality/illumination</b>	<p>Images which are poor quality (e.g. excessive motion blur) should be marked bad. However, poor illumination (e.g. objects in silhouette) should not count as poor quality unless objects cannot be recognised.</p> <p>Images made up of multiple images (e.g. collages) should be marked bad.</p>
<b>Clothing/mud/snow etc.</b>	<p>If an object is 'occluded' by a close-fitting occluder e.g. clothing, mud, snow etc., then the occluder should be treated as part of the object.</p>
<b>Transparency</b>	<p>Do label objects visible through glass, but treat reflections on the glass as occlusion.</p>
<b>Mirrors</b>	<p>Do label objects in mirrors.</p>
<b>Pictures</b>	<p>Label objects in pictures/posters/signs only if they are photorealistic but not if cartoons, symbols etc.</p>

# VOC Annual Workshop



The screenshot shows a web browser window with the following content:

- Browser tabs: "The PASCAL Visual Object Clas x +"
- Address bar: "Not Secure | host.robots.ox.ac.uk/pascal/VOC/"
- Page title: "The **PASCAL** Visual Object Classes Homepage"
- Logo: "PASCAL2" with the tagline "Pattern Analysis, Statistical Modelling and Computational Learning". The logo features a stylized blue and white geometric design.
- Section header: "The PASCAL VOC project:"
- List of bullet points:
  - 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)
- Section header: "Pascal VOC data sets"
- Text: "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."

<http://host.robots.ox.ac.uk/pascal/VOC/>

# VOC: Datasets Evolved

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

Year	Statistics	New developments	Notes
<a href="#">2005</a>	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.
<a href="#">2006</a>	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.

# ILSVRC

“ILSVRC follows in the footsteps of the PASCAL VOC challenge... which set the precedent for standardized evaluation of recognition algorithms in the form of yearly competitions.”



# ILSVRC

## 1. Category Selection

- 200 ImageNet classes which:

- 1) exclude synset overlap
- 2) exclude object classes too “big” in the image
- 3) are basic-level categories
- 4) backward compatible: VOC

Class name in PASCAL VOC (20 classes)	Closest class in ILSVRC-DET (200 classes)
aeroplane	airplane
bicycle	bicycle
bird	bird
<i>boat</i>	<i>watercraft</i>
<i>bottle</i>	<i>wine bottle</i>
bus	bus
car	car
cat	domestic cat
chair	chair
<i>cow</i>	<i>cattle</i>
<i>dining table</i>	<i>table</i>
dog	dog
horse	horse
motorbike	motorcyle
person	person
<i>potted plant</i>	<i>flower pot</i>
sheep	sheep
sofa	sofa
train	train
tv/monitor	tv or monitor

# ILSVRC

## 1. Category Selection

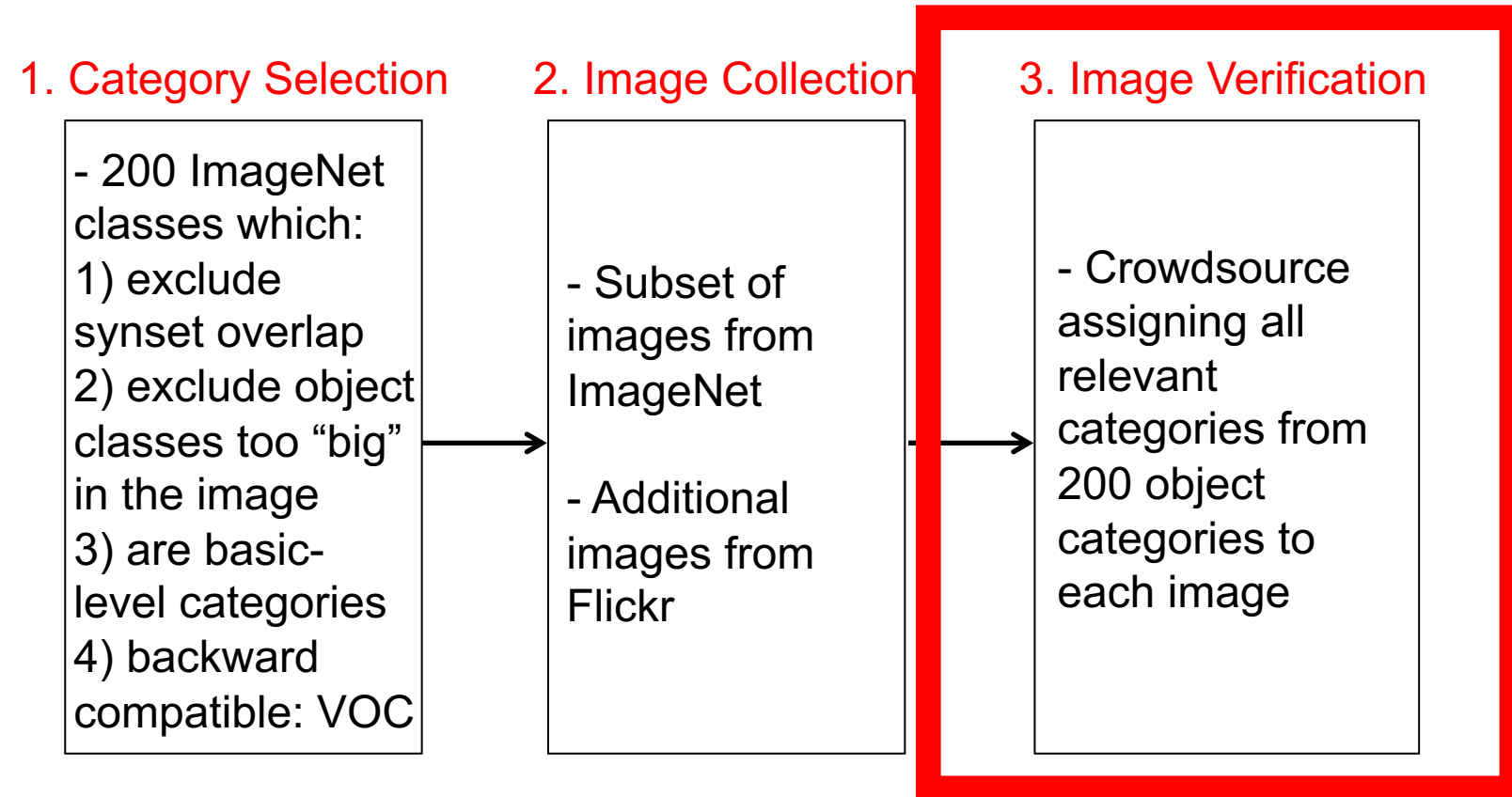
- 200 ImageNet classes which:
  - 1) exclude synset overlap
  - 2) exclude object classes too “big” in the image
  - 3) are basic-level categories
  - 4) backward compatible: VOC



## 2. Image Collection

- Subset of images from ImageNet
- Additional images from Flickr

# ILSVRC




# Recall from ImageNet: Object Presence Labeling

Identify images which contain object categories  
Requester: VLab  
Qualifications Required: None


Reward: \$0.01 per HIT    HITs Available: 1    Duration: 30 minutes

Main    Instructions


**Good Examples**  
(mouse over to enlarge):



**Bad Examples (COMMON MISTAKES)**



Please click on the images that contain **rabbit**

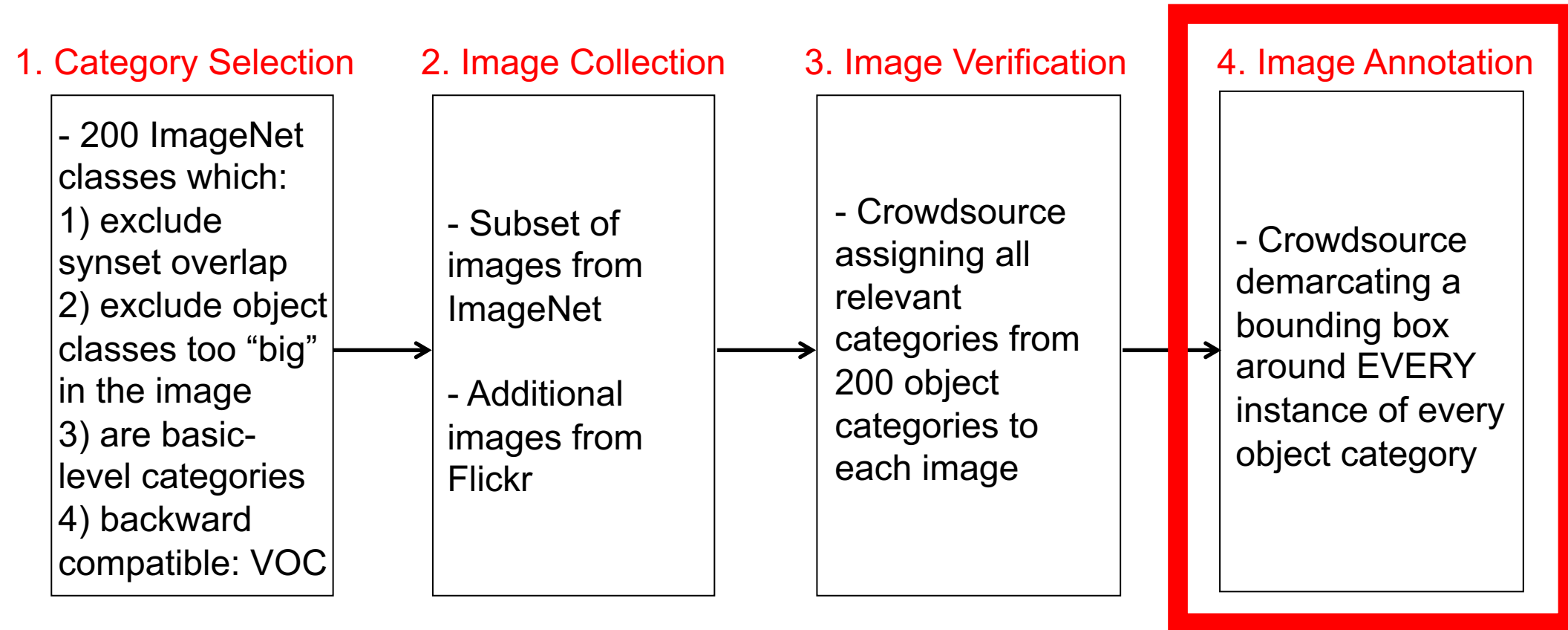


Below are the photos you have selected FROM THIS PAGE ONLY ( they will be saved when you navigate to other pages ). Click to deselect.

< page 1 of 6 >    Submit    Submit button will be enabled on the final page.

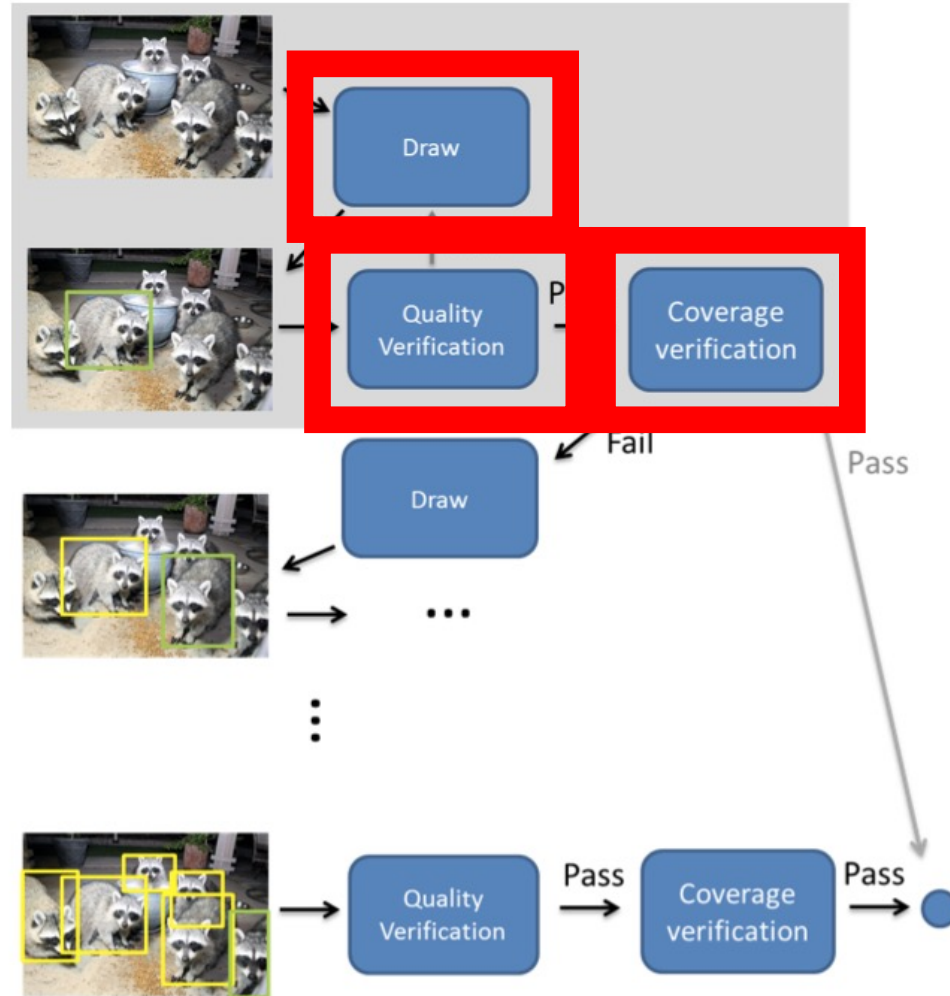
Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, Li Fei-Fei , IJCV 2015

# ILSVRC



# ILSVRC: Efficient Object Localization

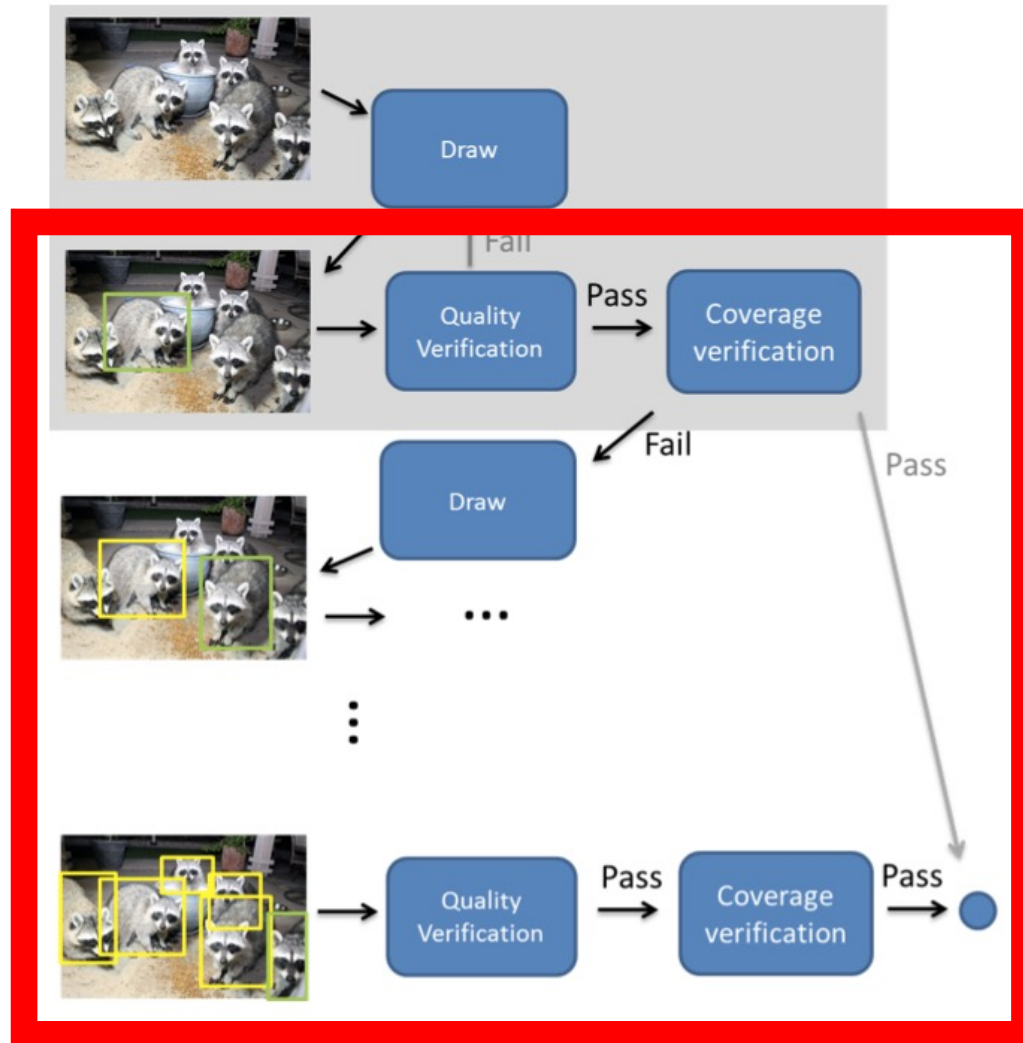
- 3 Tasks:



Idea: each task has fixed and predictable amount of work

# ILSVRC: Efficient Object Localization

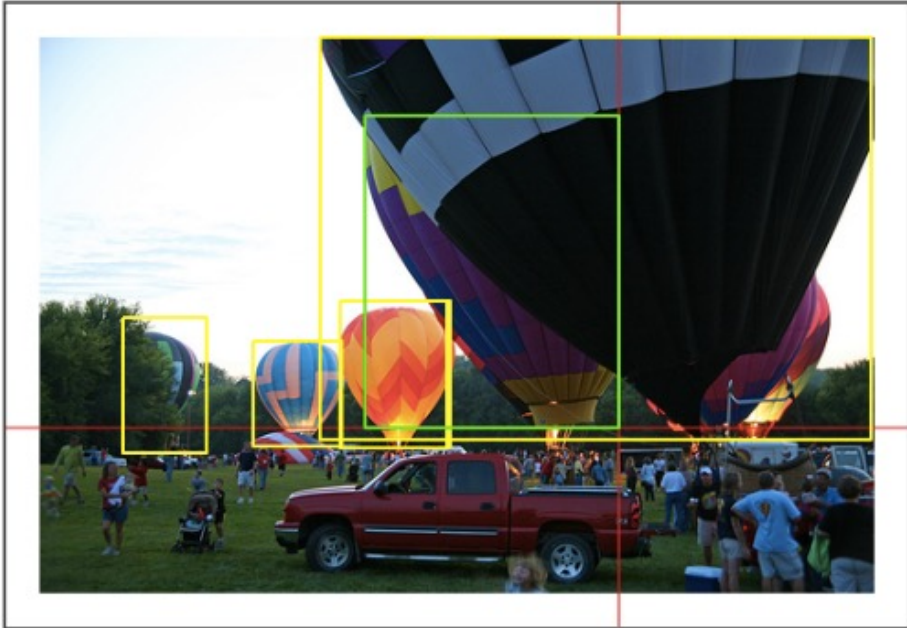
- 3 Tasks:



# ILSVRC: Drawing Task

Main | **Instructions with examples** | Look up "balloon" in Wikipedia | in Google

Draw a box around **balloon**: *large tough nonrigid bag filled with gas or heated air*



Draw a bounding box around the following object in the image:

**balloon**: large tough nonrigid bag filled with gas or heated air

**Instructions:**

- Include all visible parts and draw as tightly as possible
- **If there are multiple instances, pick only ONE ( any one ).**
- **Do NOT draw on the instances that already have bounding boxes.**

[SEE INSTRUCTIONS WITH EXAMPLES](#)

**Check here** if there's NO balloon in this image or if every instance already has a bounding box.

(Optional) Enter any comment you have:

ev NO. 1

**Images in total. 0 left.** This is a preview. Please **accept it first.**

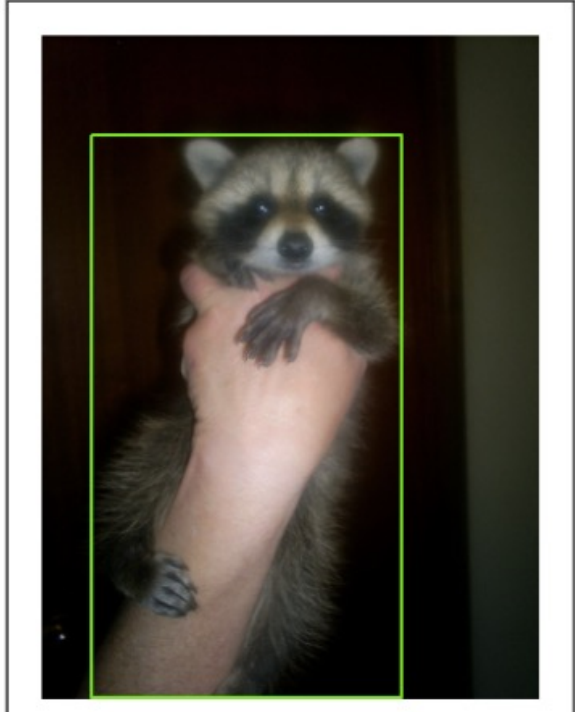
Drag the red corners to adjust the box or click 'clear box' to start over.



# ILSVRC: Quality Verification Task

Main [Instructions with examples](#) [Look up "raccoon" in Wikipedia](#) [in Google](#)

Answer questions about **"raccoon, racoon: an omnivorous nocturnal mammal native to North America and Central America"** in the image.



[SEE INSTRUCTIONS WITH EXAMPLES](#)

**Question:** Is the **GREEN** bounding box good? A good bounding box must meet ALL the conditions below:

- It contains one instance of **raccoon, racoon: an omnivorous nocturnal mammal native to North America and Central America**
- It includes all visible parts and is drawn as tightly as possible.
- It contains ONLY ONE instance of "raccoon, racoon" if there are multiple instances

GOOD ( default )

BAD

(Optional) Enter any comment you have:


NO. 2

**11 images in total. 9 left.** 'Submit' button will show up in the final page.

# ILSVRC: Coverage Verification Task

[Main](#) [Instructions with examples](#) [Look up "bird" in Wikipedia](#) [in Google](#)

Draw a box around **bird**: *warm-blooded egg-laying vertebrates characterized by feathers and forelimbs modified as wings* [SEE INSTRUCTIONS WITH EXAMPLES](#)



**Question:** Does every instance of "bird" have a bounding box ( either green or yellow )?

YES, everyone has a bounding box.  
 NO, not everyone has a bounding box.

(Optional) Enter any comment you have:

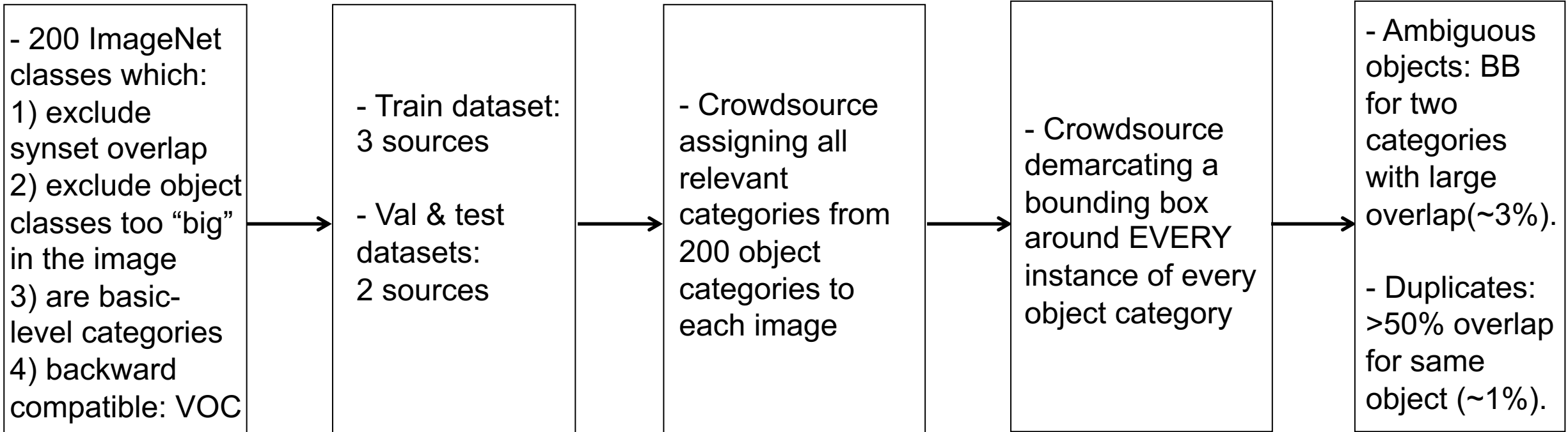
NO. 4

**198 images in total. 194 left. This is a preview.**  
Please accept it first.

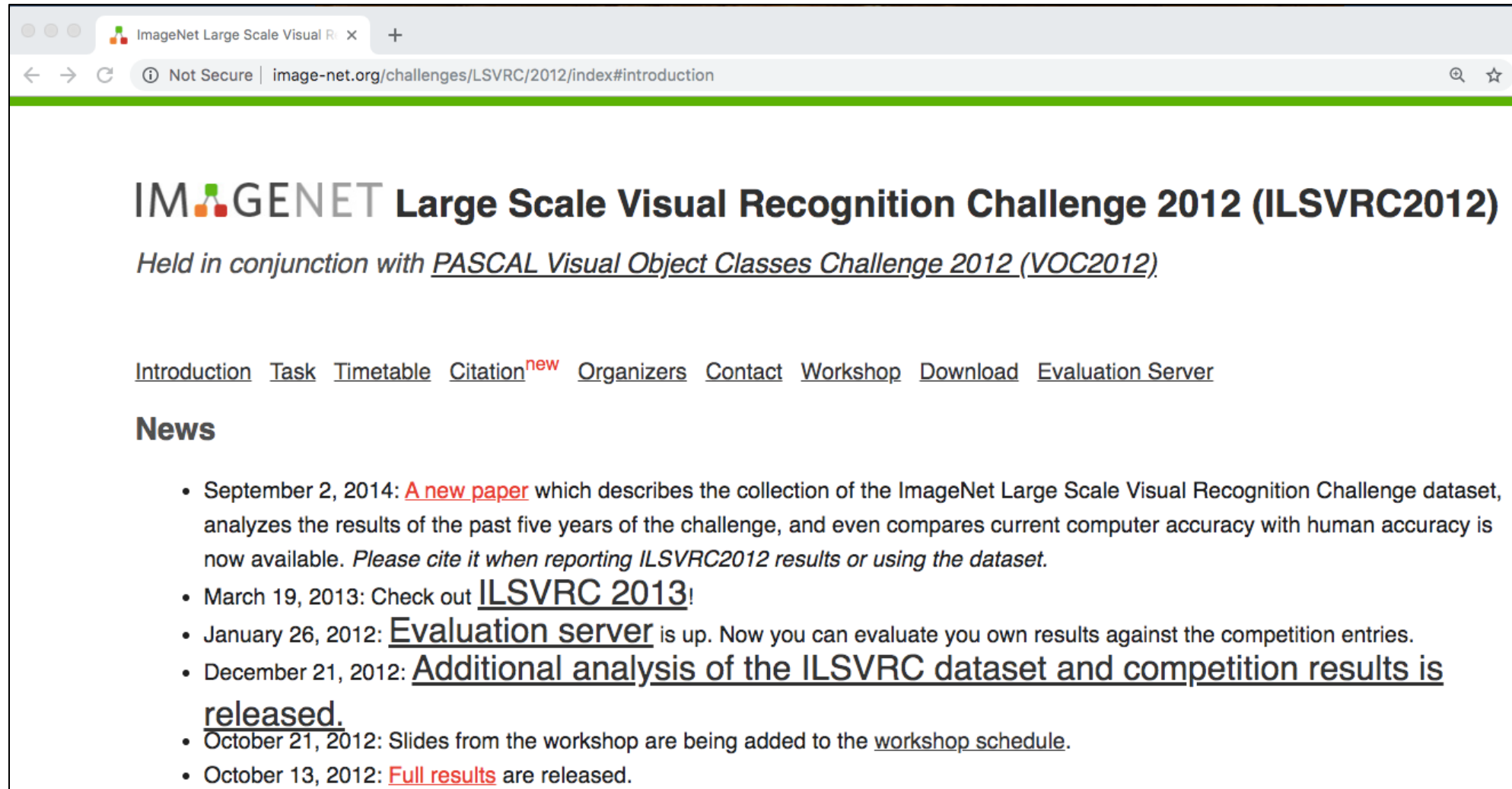
Answer the questions on the right! That is it!

# ILSVRC

## 1. Category Selection    2. Image Collection    3. Object presence labeling    4. Object localization    5. Author Review



# Object Detection: ILSVRC Annual Workshop



ImageNet Large Scale Visual Recognition Challenge 2012 (ILSVRC2012)

Held in conjunction with PASCAL Visual Object Classes Challenge 2012 (VOC2012)

[Introduction](#) [Task](#) [Timetable](#) [Citation<sup>new</sup>](#) [Organizers](#) [Contact](#) [Workshop](#) [Download](#) [Evaluation Server](#)

## News

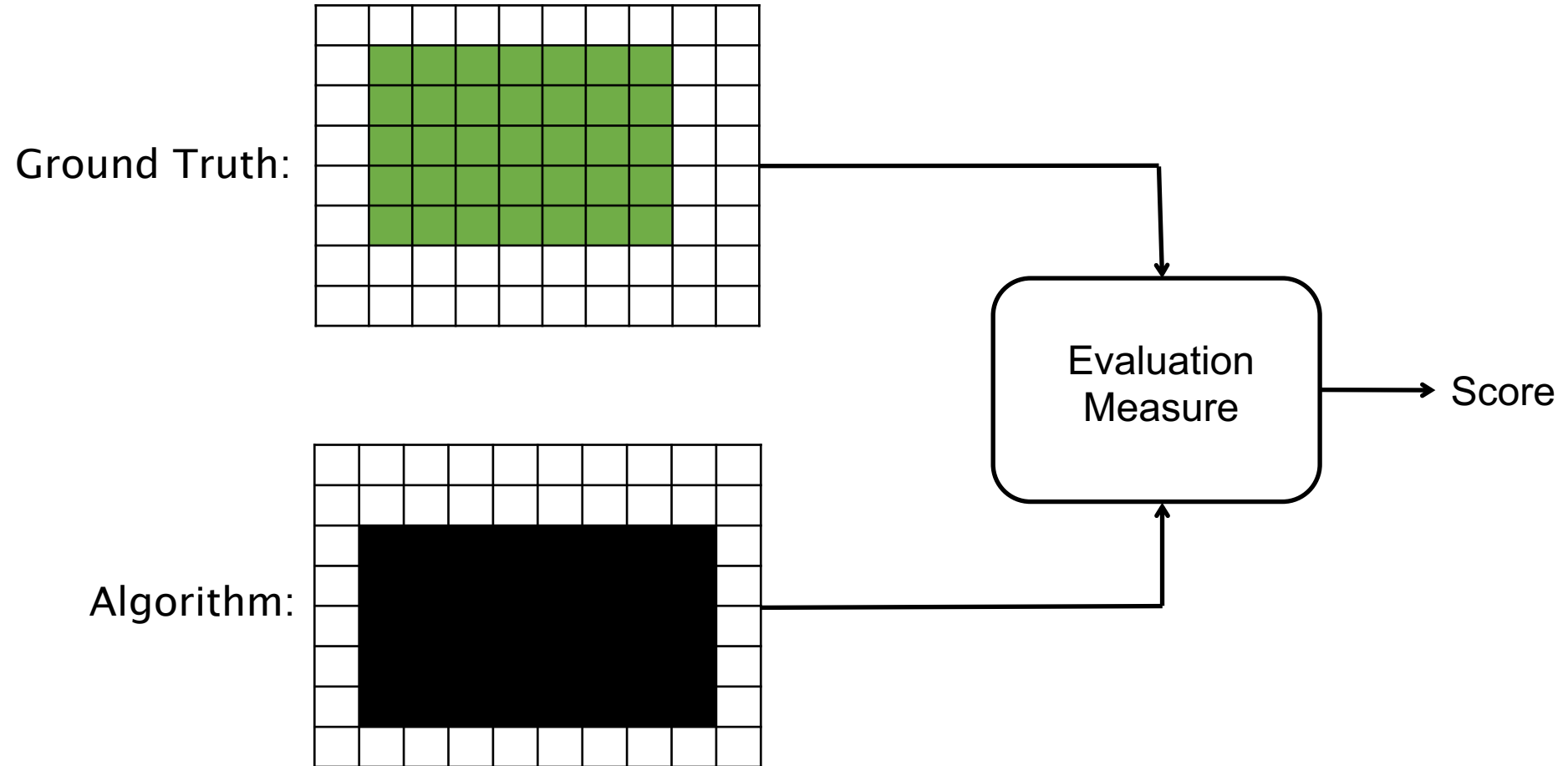
- September 2, 2014: [A new paper](#) which describes the collection of the ImageNet Large Scale Visual Recognition Challenge dataset, analyzes the results of the past five years of the challenge, and even compares current computer accuracy with human accuracy is now available. *Please cite it when reporting ILSVRC2012 results or using the dataset.*
- March 19, 2013: Check out [ILSVRC 2013!](#)
- January 26, 2012: [Evaluation server](#) is up. Now you can evaluate you own results against the competition entries.
- December 21, 2012: [Additional analysis of the ILSVRC dataset and competition results is released.](#)
- October 21, 2012: Slides from the workshop are being added to the [workshop schedule](#).
- October 13, 2012: [Full results](#) are released.

<http://image-net.org/challenges/LSVRC/2012/index#introduction>

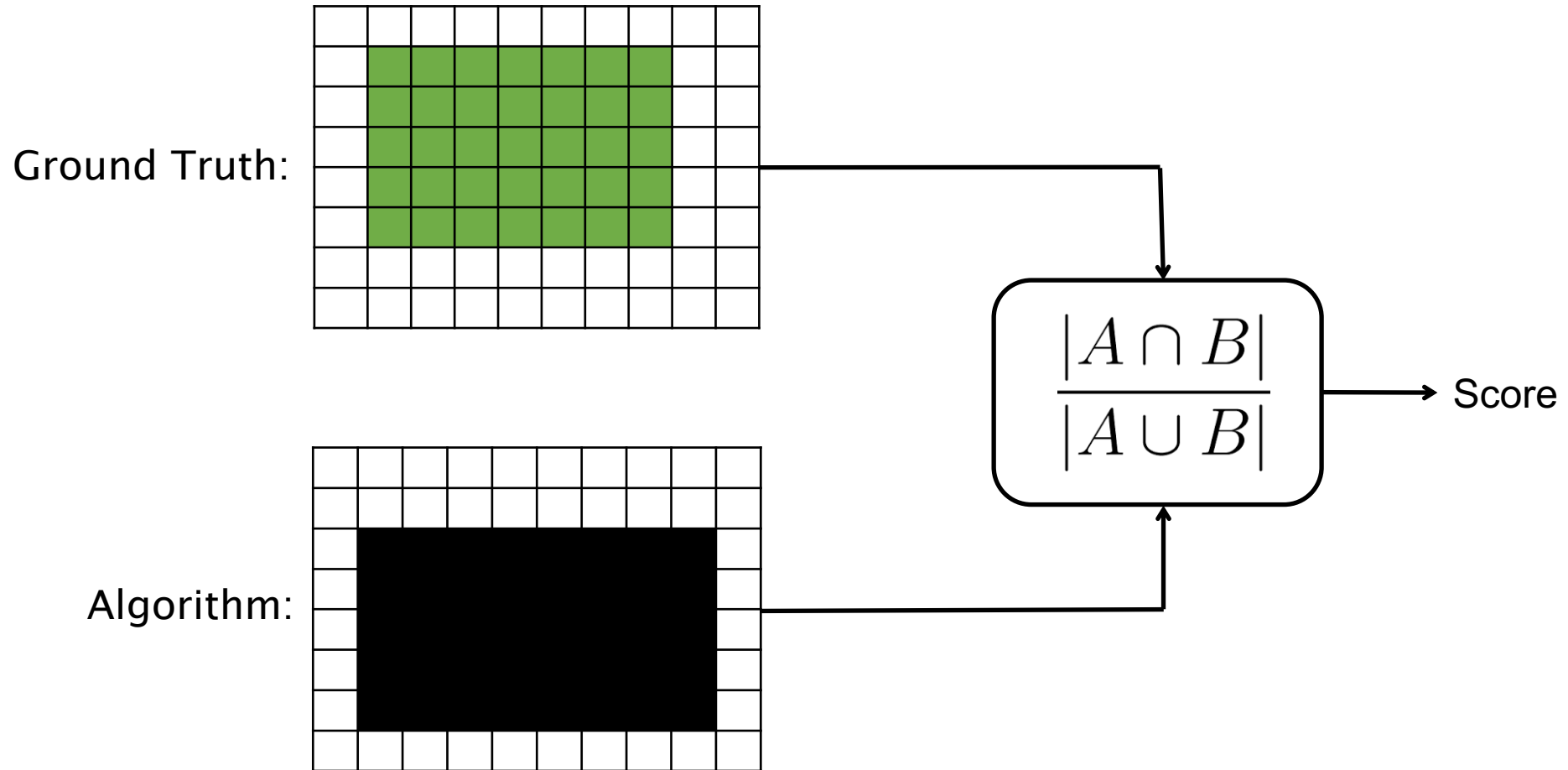
# Object Detection: Today's Topics

- Problem
- Applications
- Datasets
- **Evaluation metric**
- Overview of object detection algorithms and baseline (R-CNN)

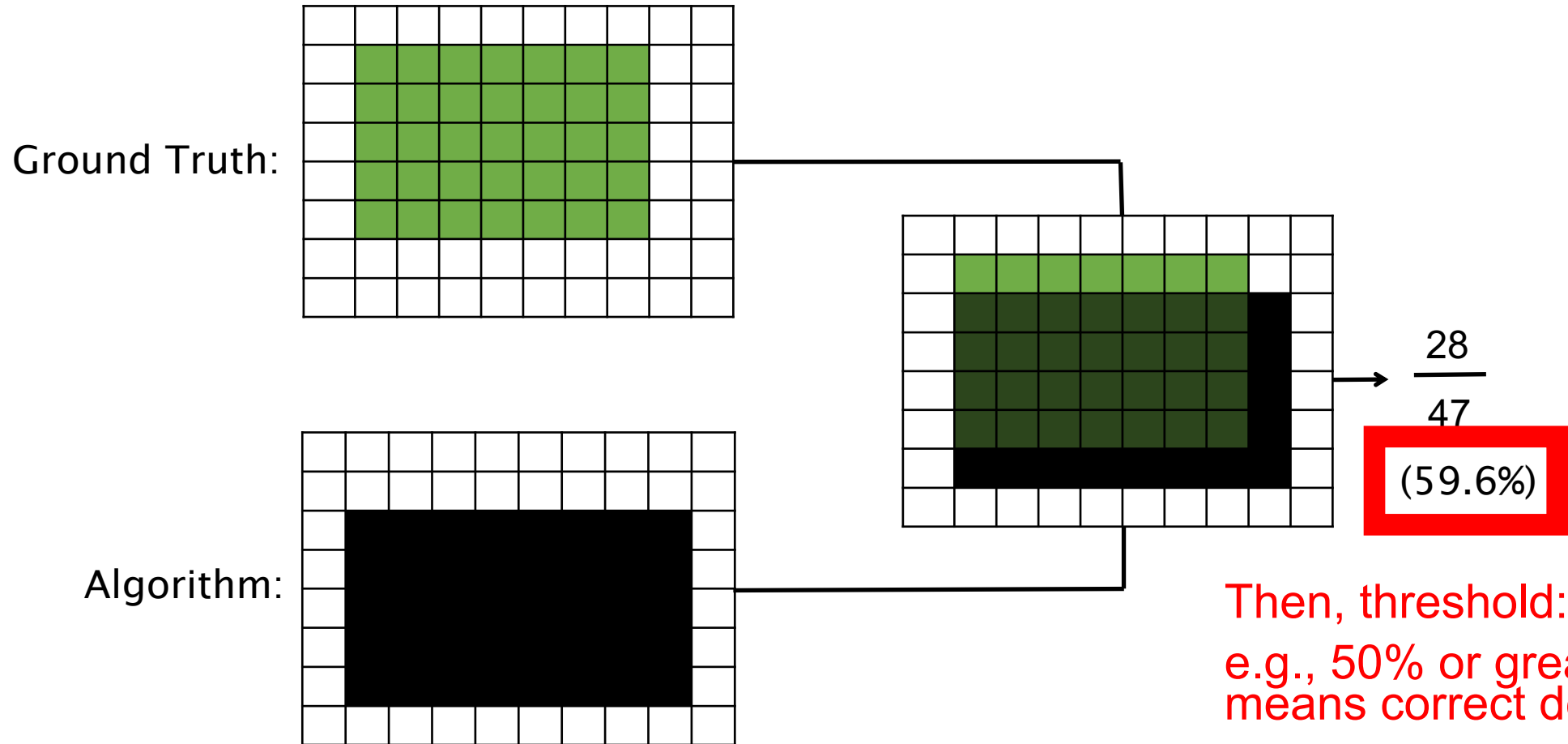
# Single Object



# Single Object: IoU (Intersection Over Union)



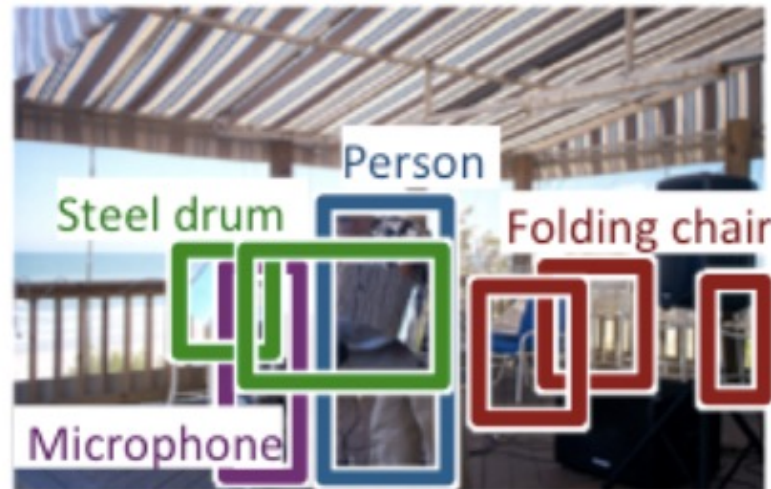
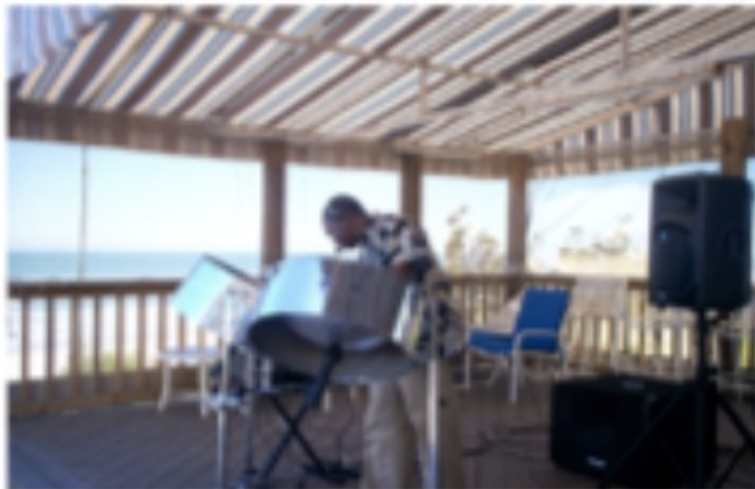
# Single Object: IoU (Intersection Over Union)



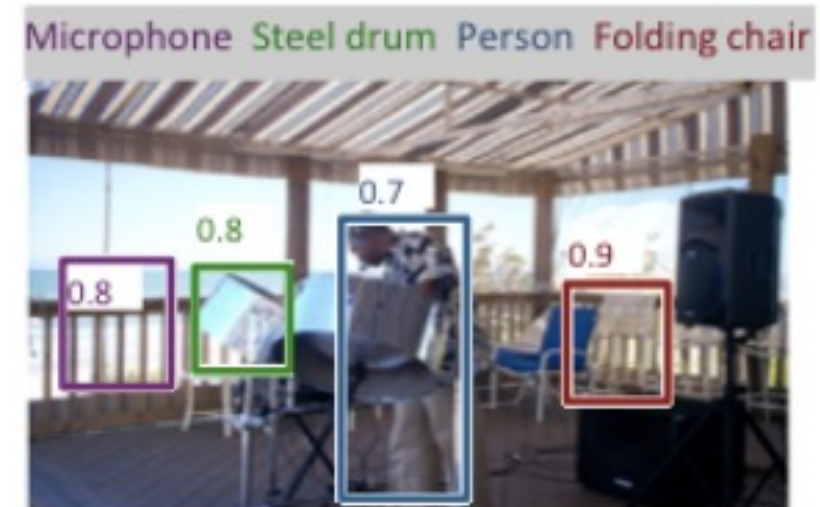


# Evaluation Metric Basics: Precision

- For each object class (e.g., cat, dog, ...), compute precision: fraction of correct detections from all detections using 0.5 IoU threshold



Ground truth



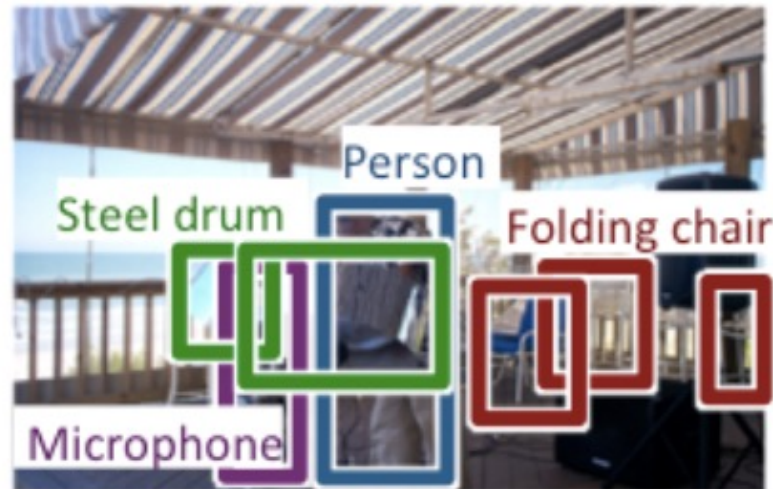
AP: ? ? ? ?

[Russakovsky et al; IJCV 2015]

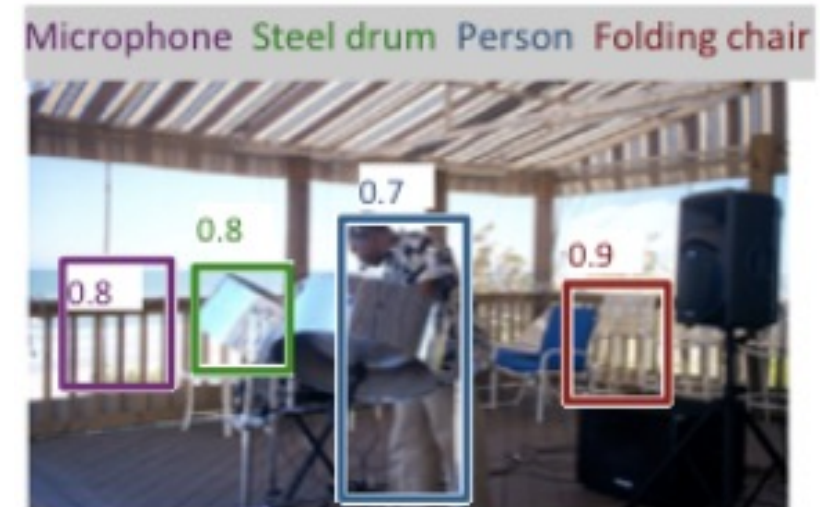
<https://jonathan-hui.medium.com/map-mean-average-precision-for-object-detection-45c121a31173>

# Evaluation Metric Basics: Precision

- For each object class (e.g., cat, dog, ...), compute precision: fraction of correct detections from all detections using IoU threshold (e.g., 0.5)



Ground truth



P: 0.0 0.5 1.0 0.3

[Russakovsky et al; IJCV 2015]

<https://jonathan-hui.medium.com/map-mean-average-precision-for-object-detection-45c121a31173>

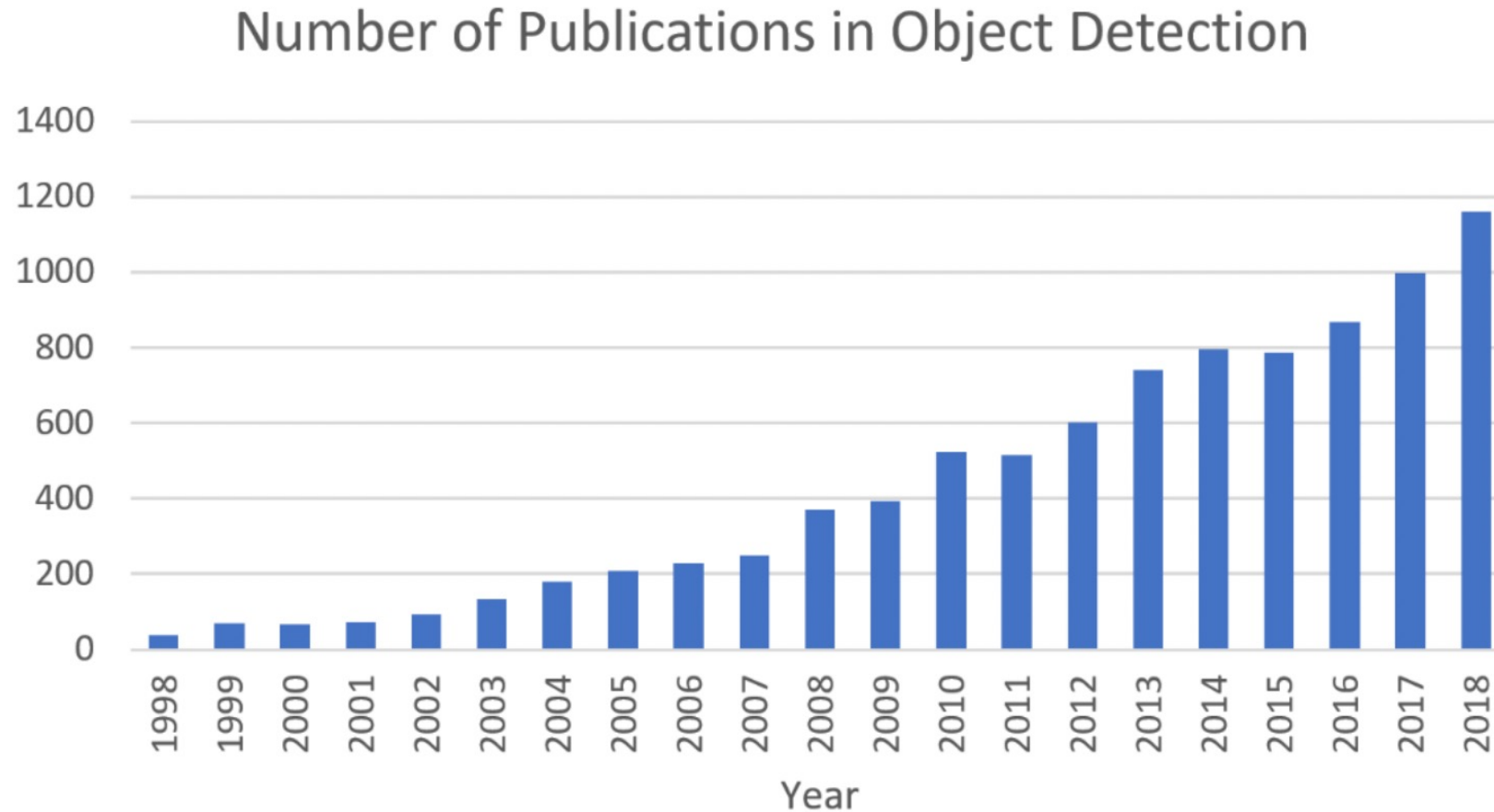
# Evaluation Metric: mAP

- For each object class (e.g., cat, dog, ...), compute Average Precision (AP)
  - Vary IoU threshold in order to create a precision-recall curve, and then compute area under the curve
- Then, compute mean AP across all classes

# Object Detection: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metric
- Overview of object detection algorithms and baseline (R-CNN)

# Community Research Engagement



“Data from Google scholar advanced search: allintitle: ‘object detection’ AND ‘detecting objects’”

# Naïve Solution: Sliding Window Approach

Person?

Person?

Person?

Person?

Person?

Person?

Person?

Person?

Person?



Image Source: <https://yourboulder.com/boulder-neighborhood-downtown/>

# Naïve Solution: Sliding Window Approach

Car?  
Car?  
Car?  
Car?  
Car?  
Car?  
Car?  
Car?  
Car?  
Car?



Image Source: <https://yourboulder.com/boulder-neighborhood-downtown/>

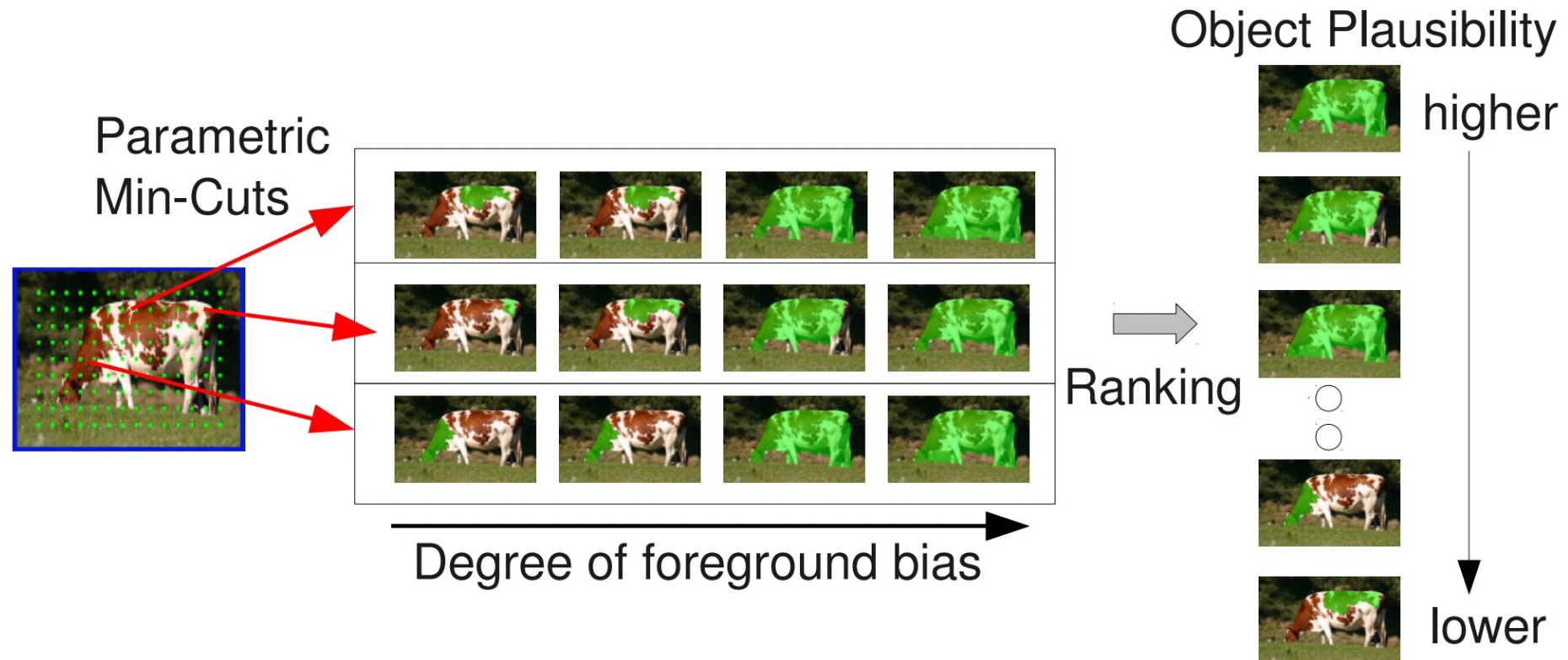
# Naïve Solution: Sliding Window Approach

- Sliding window approach: must test different locations at...
  - Different scales
  - Different aspect ratios (e.g., for person vs car or car viewed at different angles)
- Number of regions to test? (e.g., 1920 x 1080 image)
  - Easily can explode to hundreds of thousands or millions of windows
- Key limitation
  - Very slow!



# A Less Naïve Solution: Region Proposals

- Replace sliding window approach with *region proposals* (bounding boxes around “object”-like regions) found by grouping similar pixels based on low-/mid-level features; e.g., CPMC

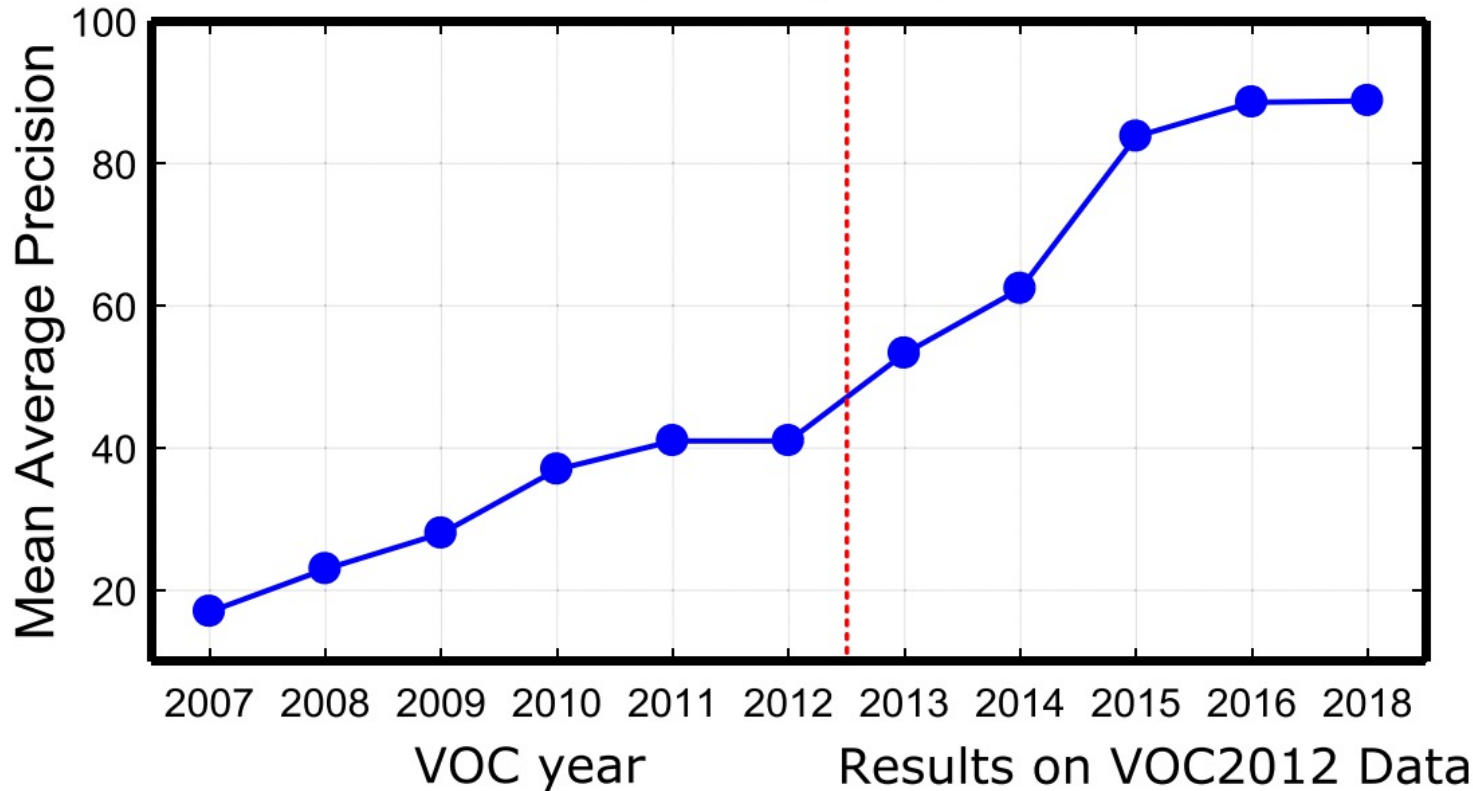


# A Less Naïve Solution: Region Proposals

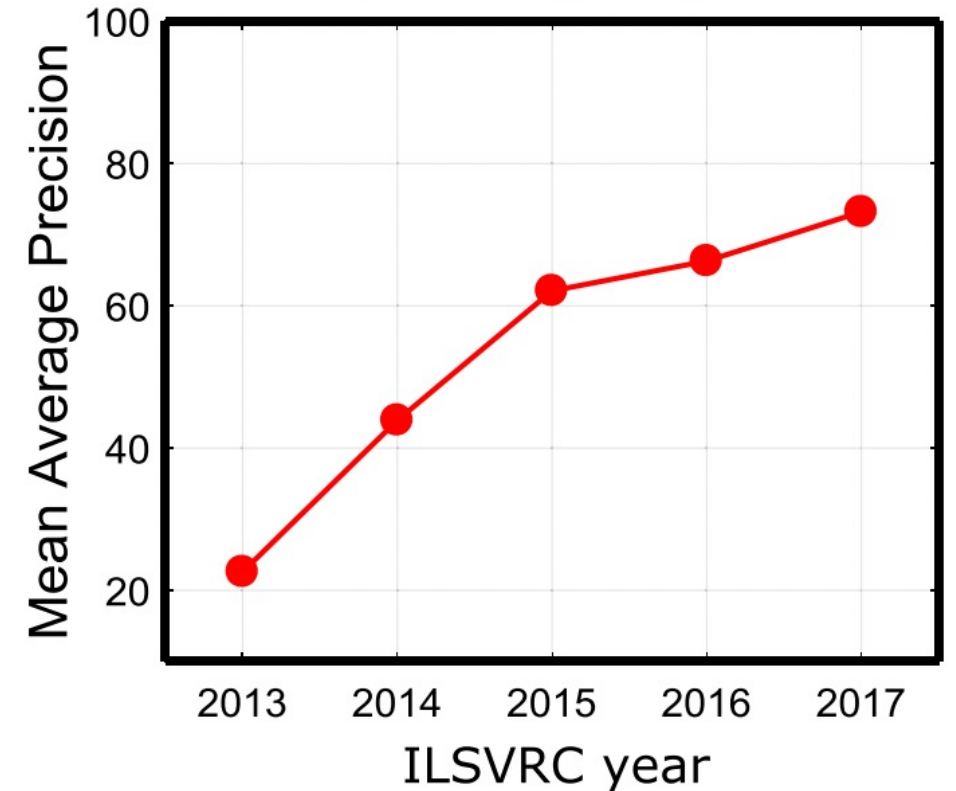
- Replace sliding window approach with *region proposals* (bounding boxes around “object”-like regions); e.g., CPMC
- **Advantage:** considerably fewer regions than needed in a naïve sliding window approach, with belief they will include the objects of interest (i.e., high recall)
- **Many options:** CPMC, Category Independent Object Proposals, Randomized Prim, Selective Search, and more
- A good start...

# Turning Point (2012): Deep Learning Solutions

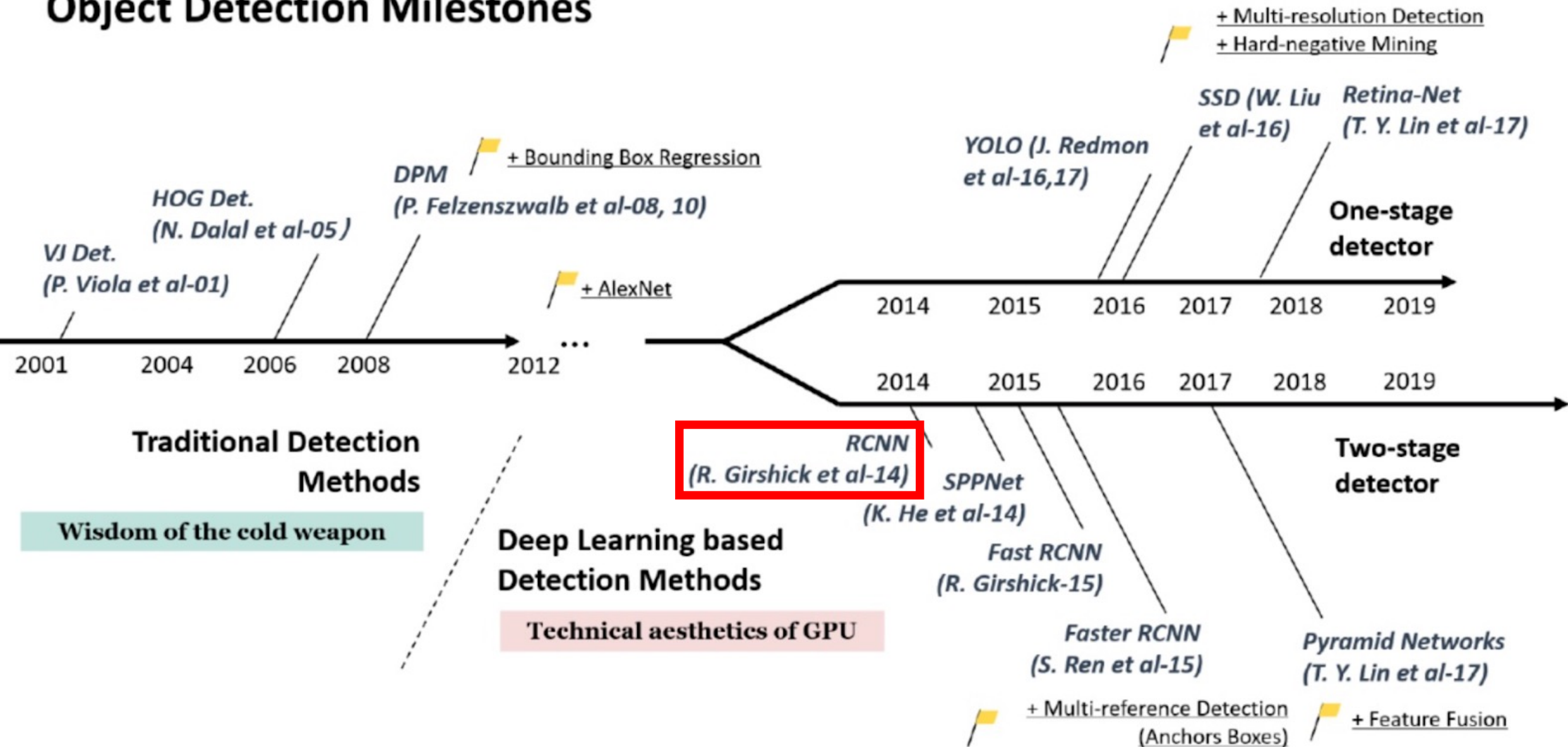
Object Detection Results  
(20 Categories)



Top Object Detection Competition Results  
(200 Categories)



# Object Detection Milestones



# Why R-CNN?

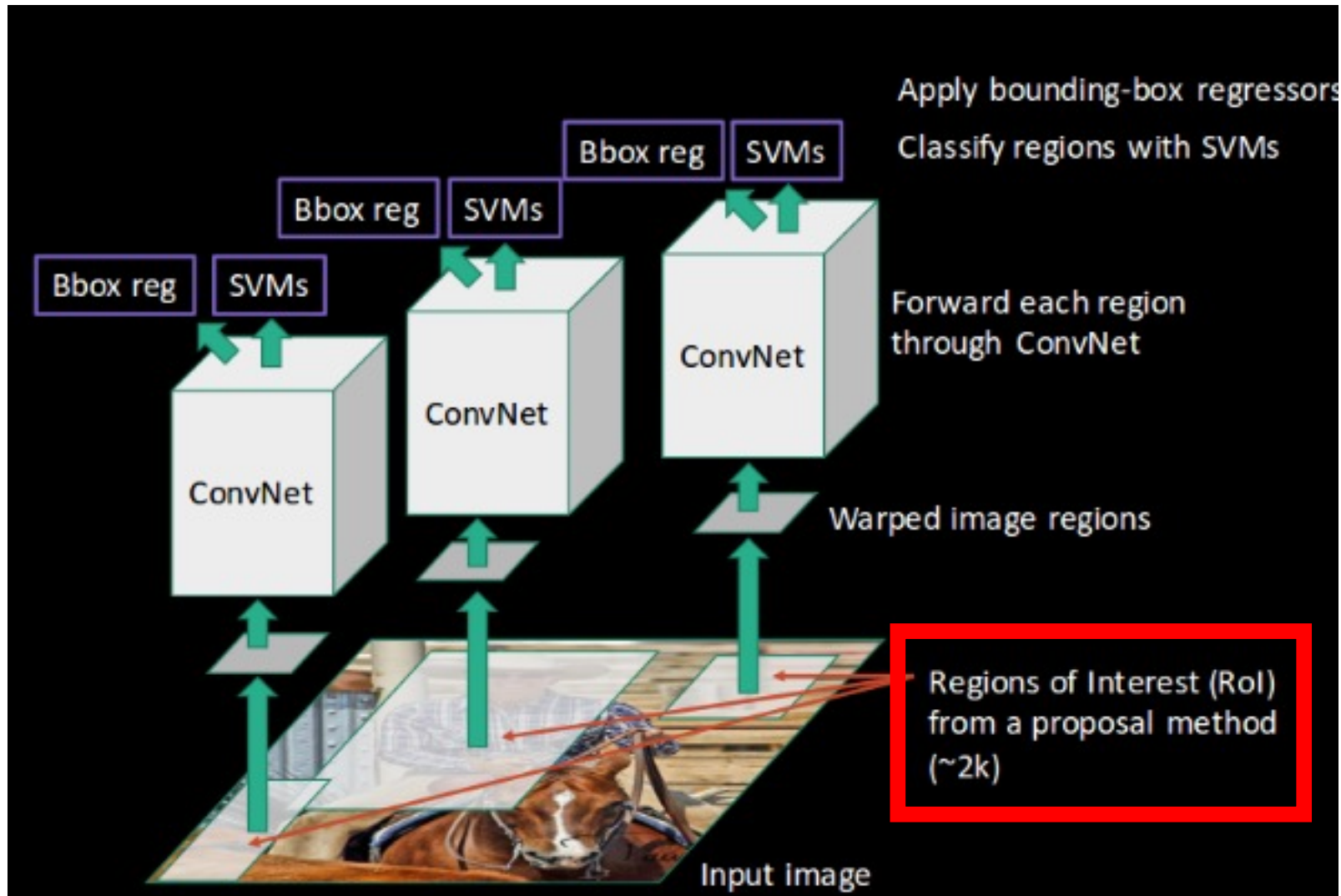
Named after the proposed technique: use **R**egion proposals with **CNN** features

Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. “Rich feature hierarchies for accurate object detection and semantic segmentation.” CVPR 2014.

# Key Contributions of R-CNN

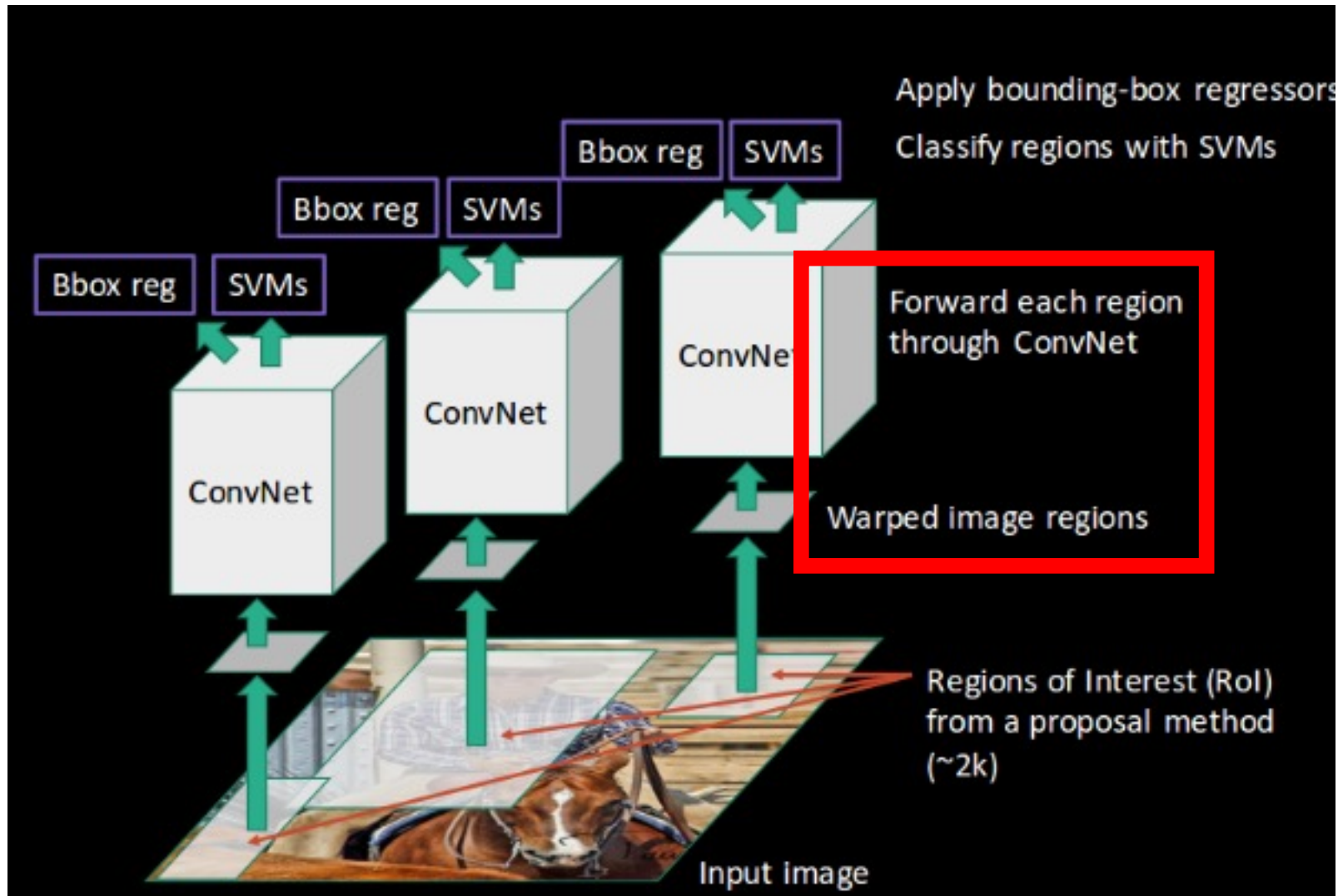
1. Demonstrate how to accurately localize objects with a neural network (NN)
  - First time a CNN outperformed hand-crafted features on VOC, achieving mAP of 54% compared to 33% for previous HOG based model (VOC 2010)
2. Demonstrate how train an accurate (high-capacity) NN with a scarce amount of annotated detection data

# Architecture



Selective Search used to enable comparison with prior work; creates ~2000 regions based on color, texture, size and shape

# Architecture





# Describe Each Region with Fixed-length Vector

**Key idea:** Given scarce amount of training data in detection datasets, devise good feature by fine-tuning model that is pre-trained on a large dataset

- Replace final layer of AlexNet (trained on ImageNet) with # of categories in (VOC) detection dataset

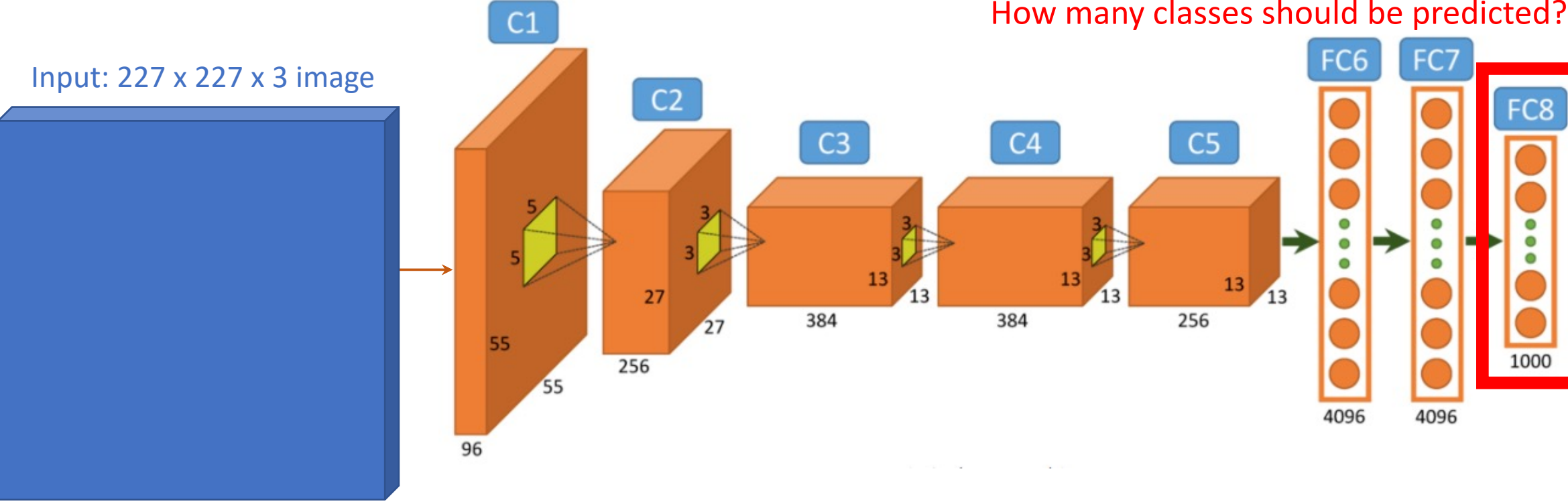


Image Source: [https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers\\_fig2\\_312303454](https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers_fig2_312303454)

# Describe Each Region with Fixed-length Vector

**Key idea:** Given scarce amount of training data in detection datasets, devise good feature by fine-tuning model that is pre-trained on a large dataset

- Replace final layer of AlexNet (trained on ImageNet) with # of categories in (ILSVRC) detection dataset

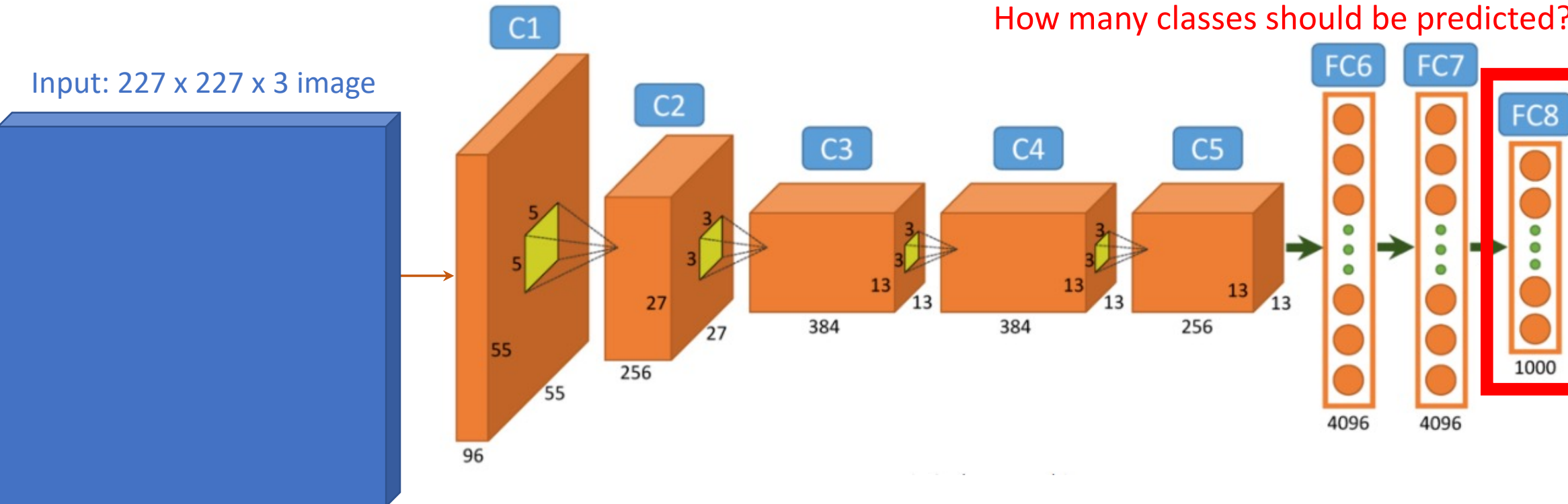


Image Source: [https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers\\_fig2\\_312303454](https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers_fig2_312303454)

# Describe Each Region with Fixed-length Vector

**Key idea:** Given scarce amount of training data in detection datasets, devise good feature by fine-tuning model that is pre-trained on a large dataset

- Replace final layer of AlexNet and use FC7 feature from fine-tuned model

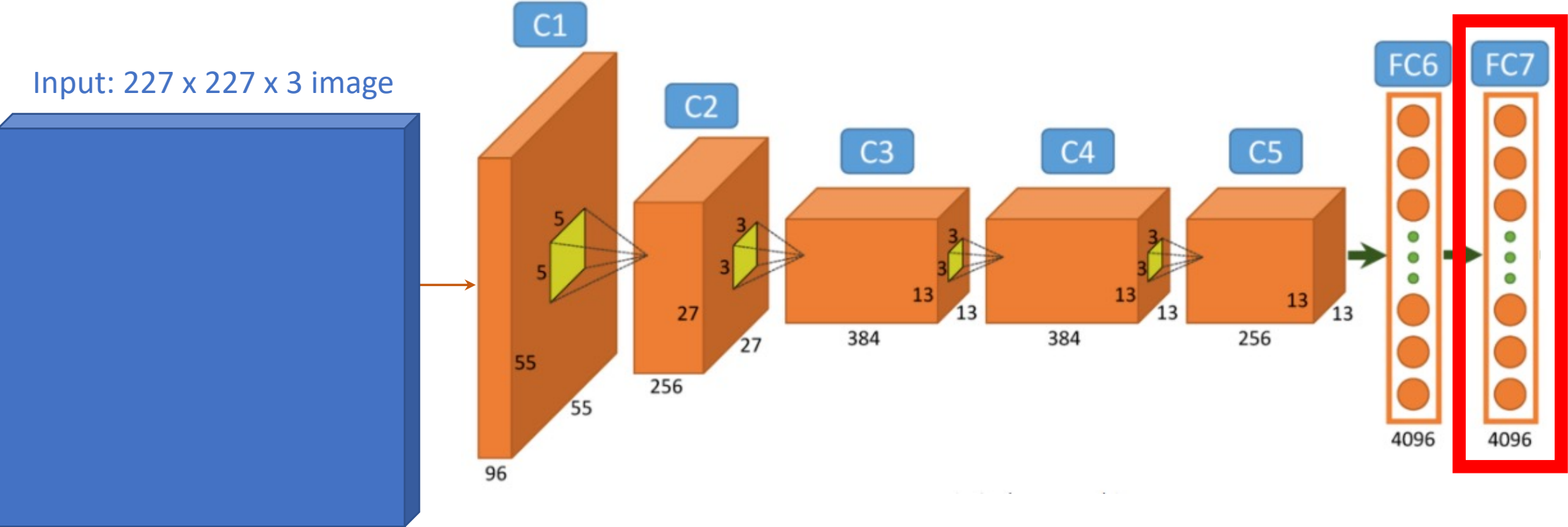


Image Source: [https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers\\_fig2\\_312303454](https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers_fig2_312303454)

# Describe Each Region with Fixed-length Vector

- Benefits of these features:
  - can be learned for a dataset instead of handcrafted (e.g., HOG, SIFT)
  - ~2 orders of magnitude smaller than traditional features (e.g., HOG, SIFT)

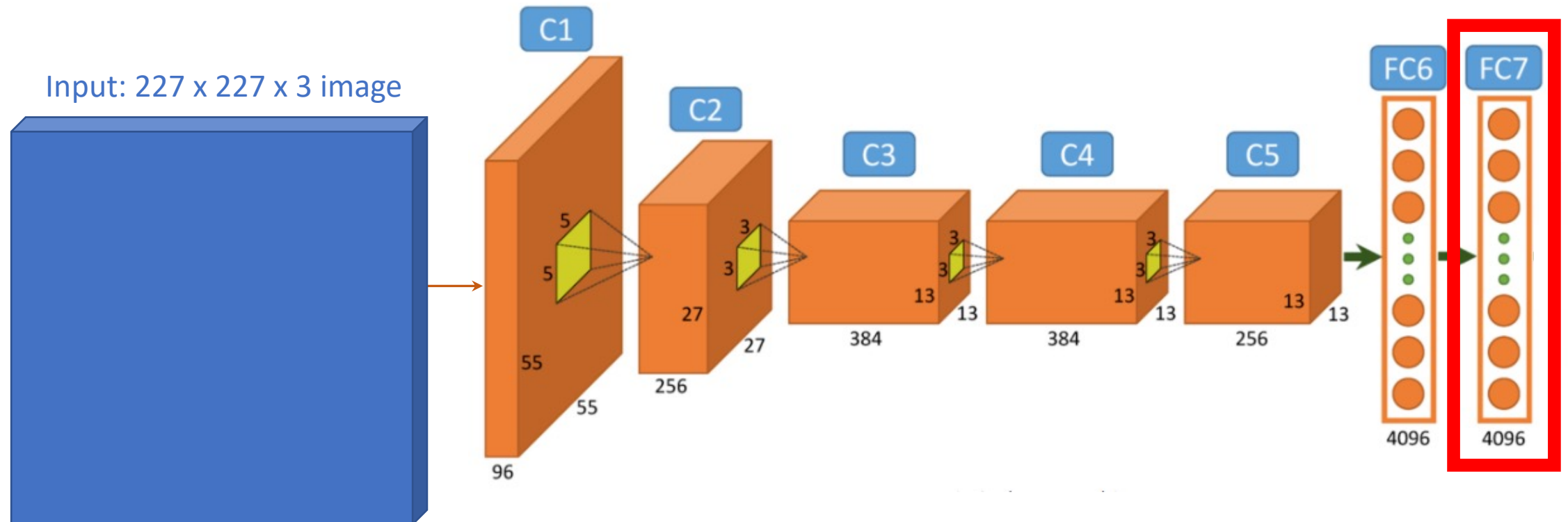


Image Source: [https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers\\_fig2\\_312303454](https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers_fig2_312303454)

# Describe Each Region with Fixed-length Vector

Challenge: how to resize a proposed region to the required size?

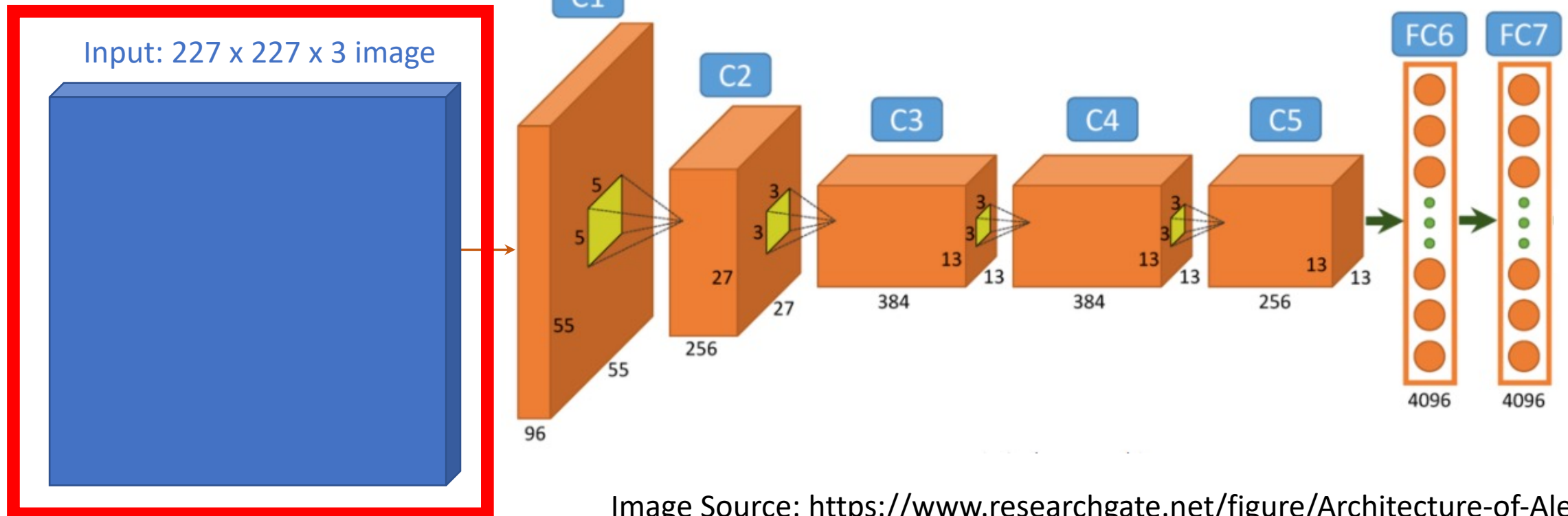


Image Source: [https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers\\_fig2\\_312303454](https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers_fig2_312303454)

# Describe Each Region with Fixed-length Vector

Challenge: how to resize a proposed region to the required size?

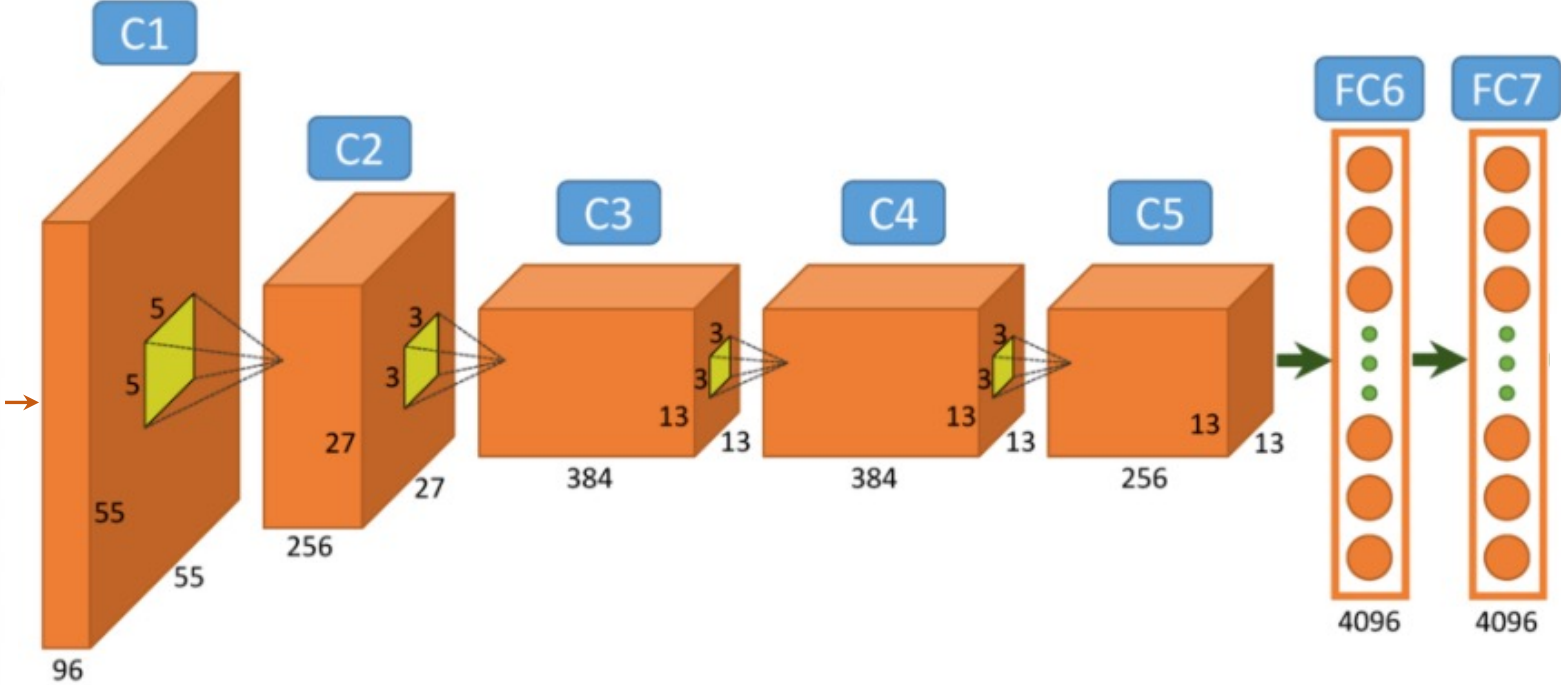


Image Source: [https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers\\_fig2\\_312303454](https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers_fig2_312303454)

# Describe Each Region with Fixed-length Vector

Challenge: how to resize a proposed region to the required size?

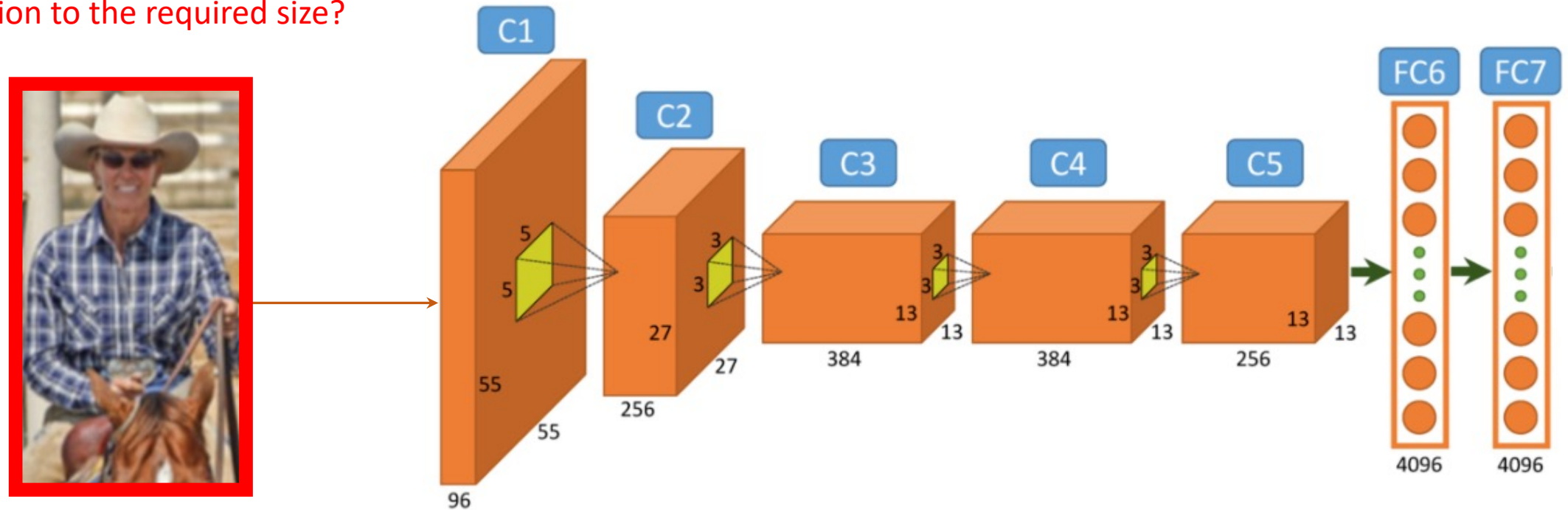


Image Source: [https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers\\_fig2\\_312303454](https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers_fig2_312303454)

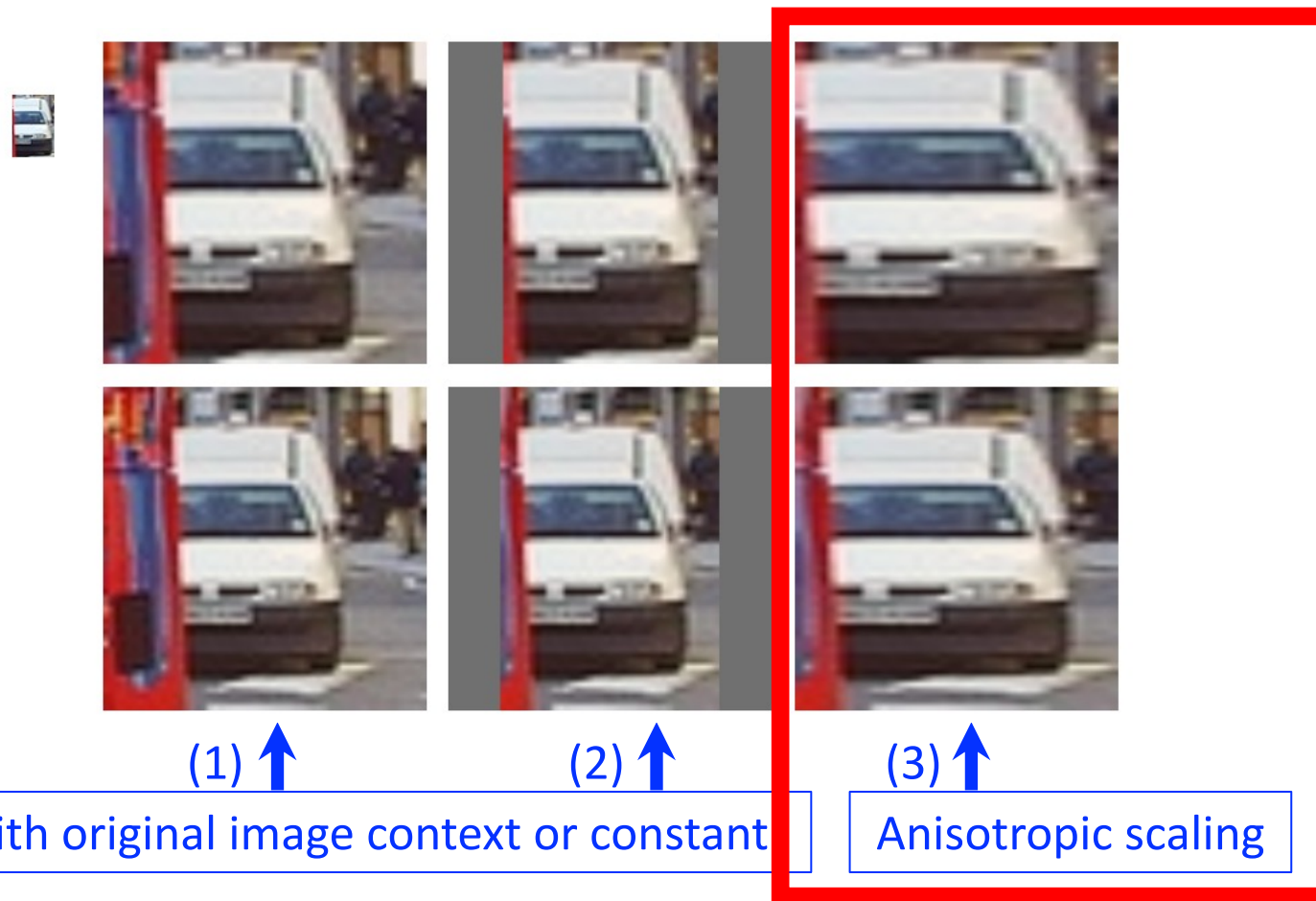
# Region Resizing



As exemplified, region proposals come in different sizes and aspect ratios



# Region Resizing



Chosen because  
experimentally  
shown to lead to  
the best results

Many ways to convert a region into a fixed input size of  $227 \times 227 \times 3$

# Describe Each Region with Fixed-length Vector

Input: 227 x 227 x 3 image

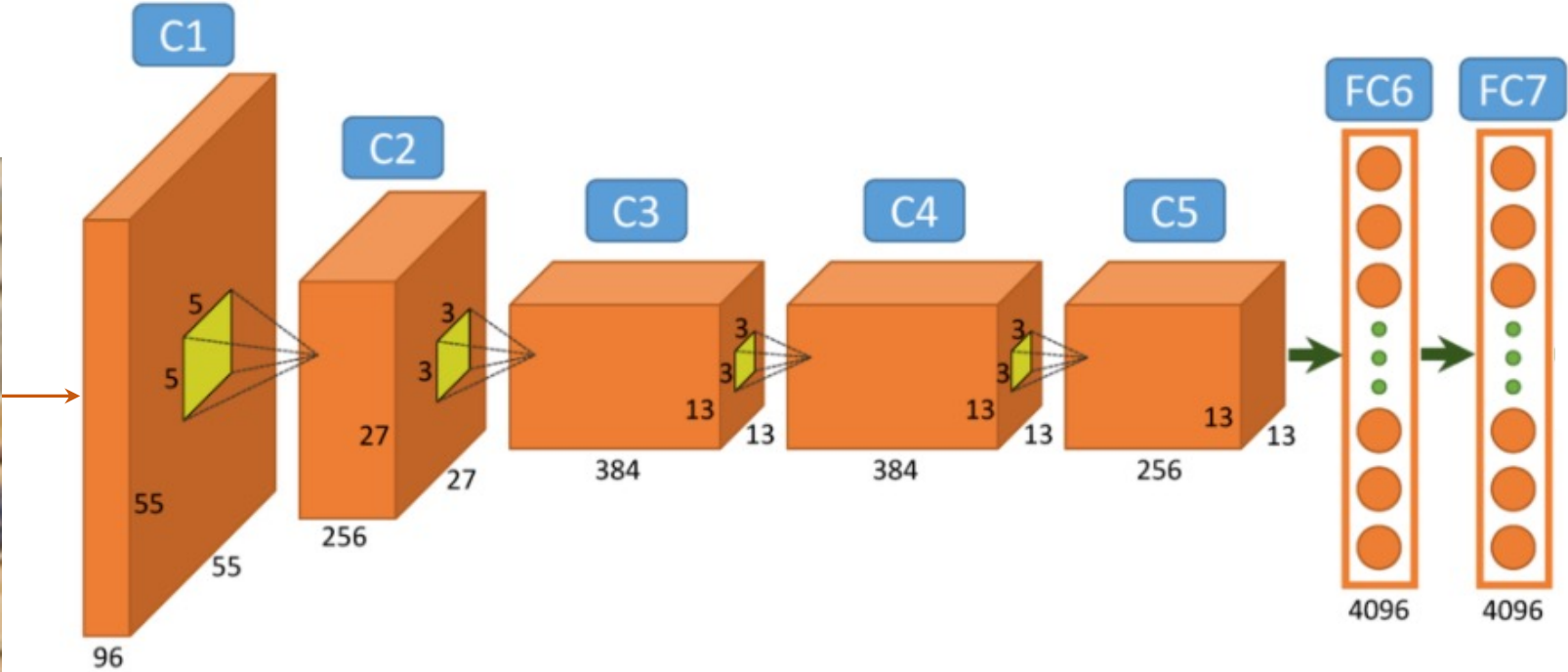
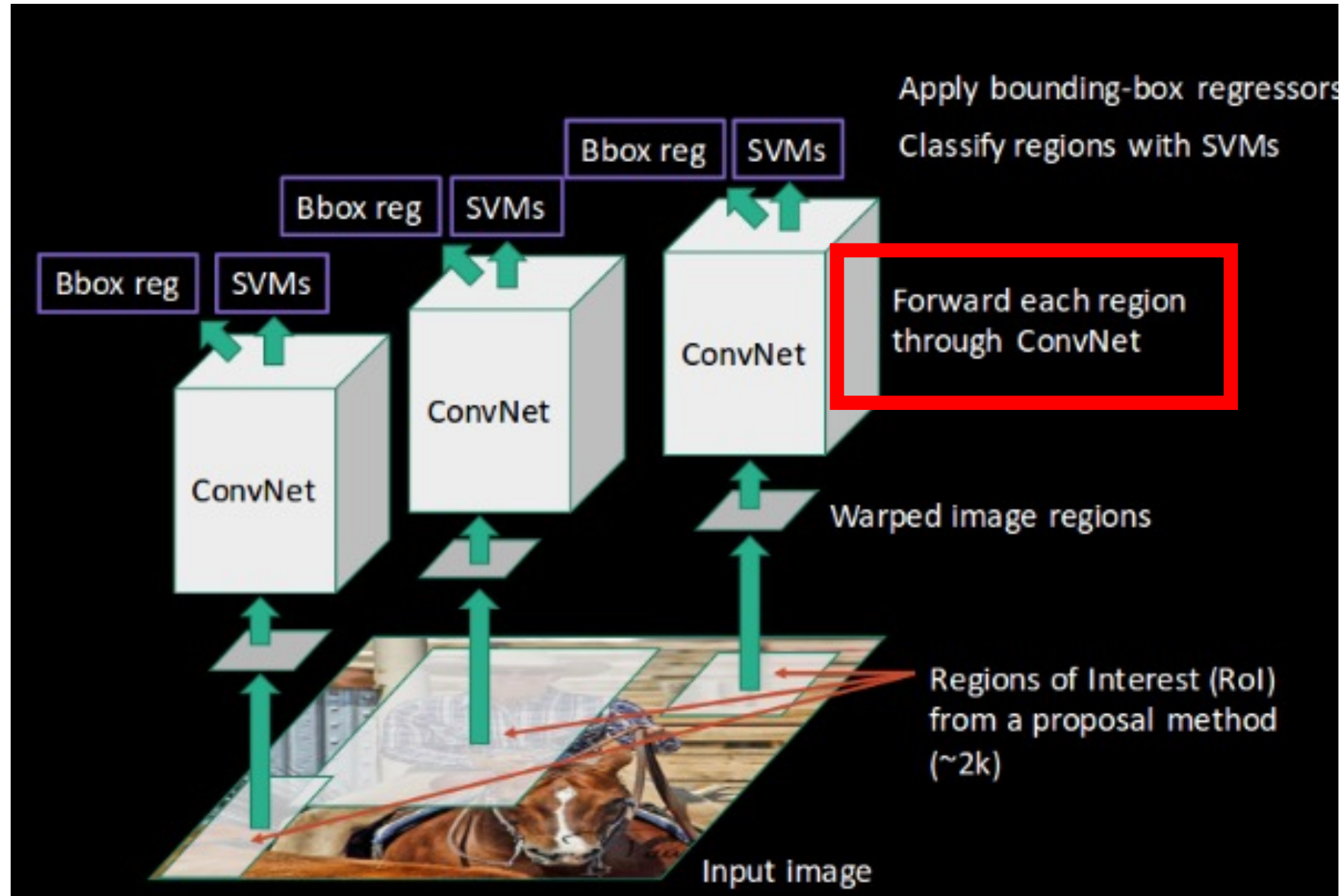
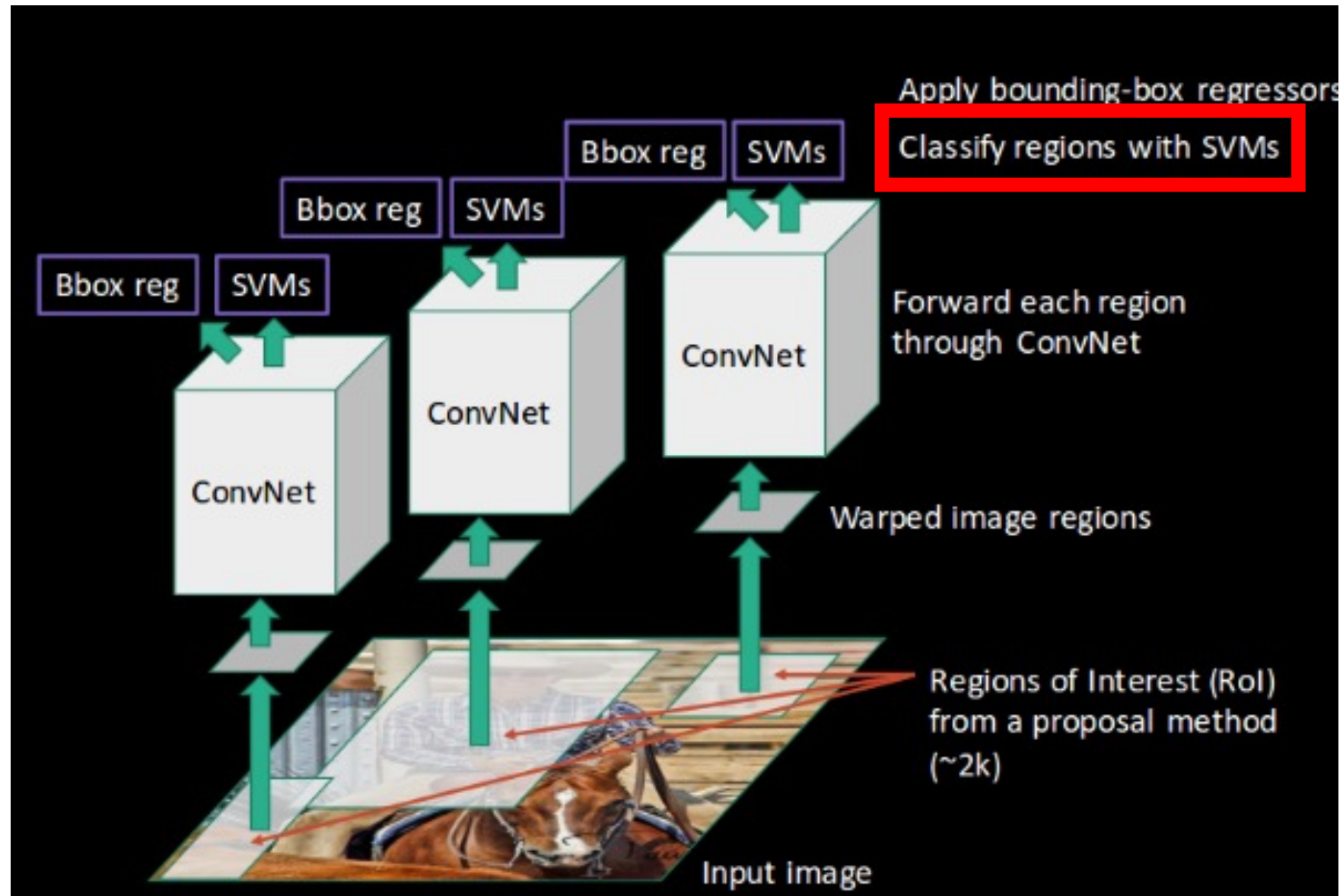


Image Source: [https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers\\_fig2\\_312303454](https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers_fig2_312303454)

# Architecture

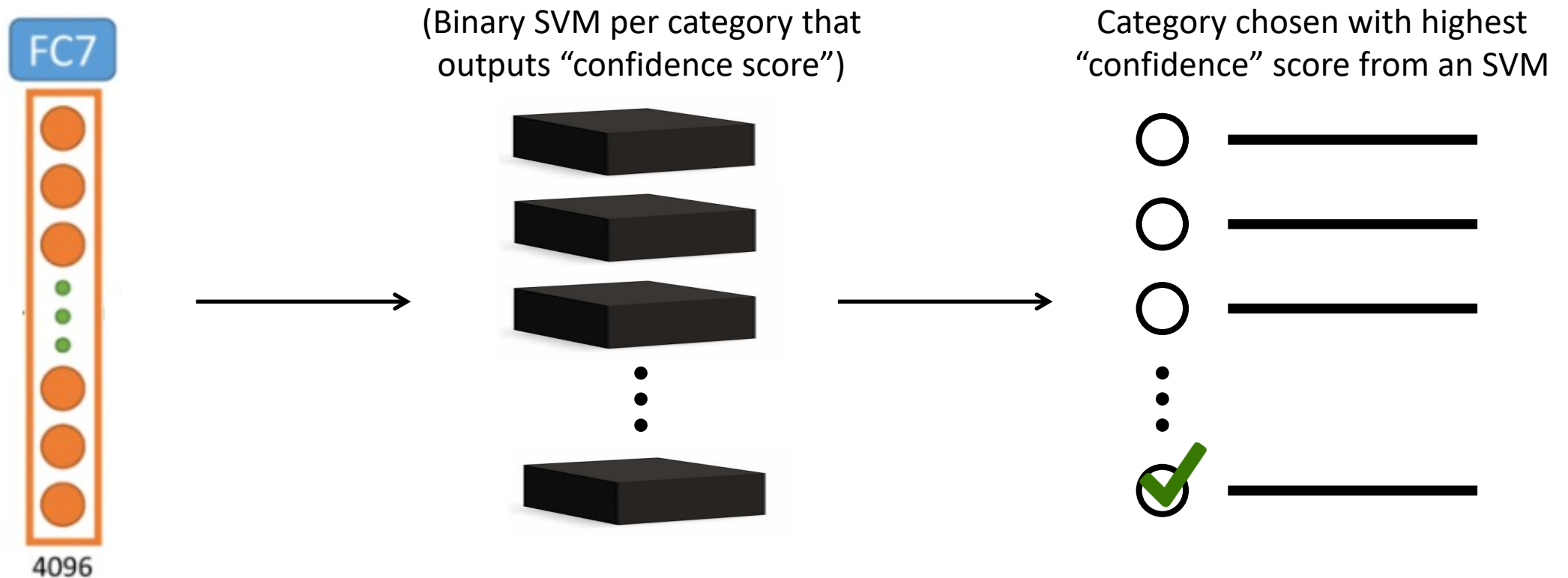


# Architecture

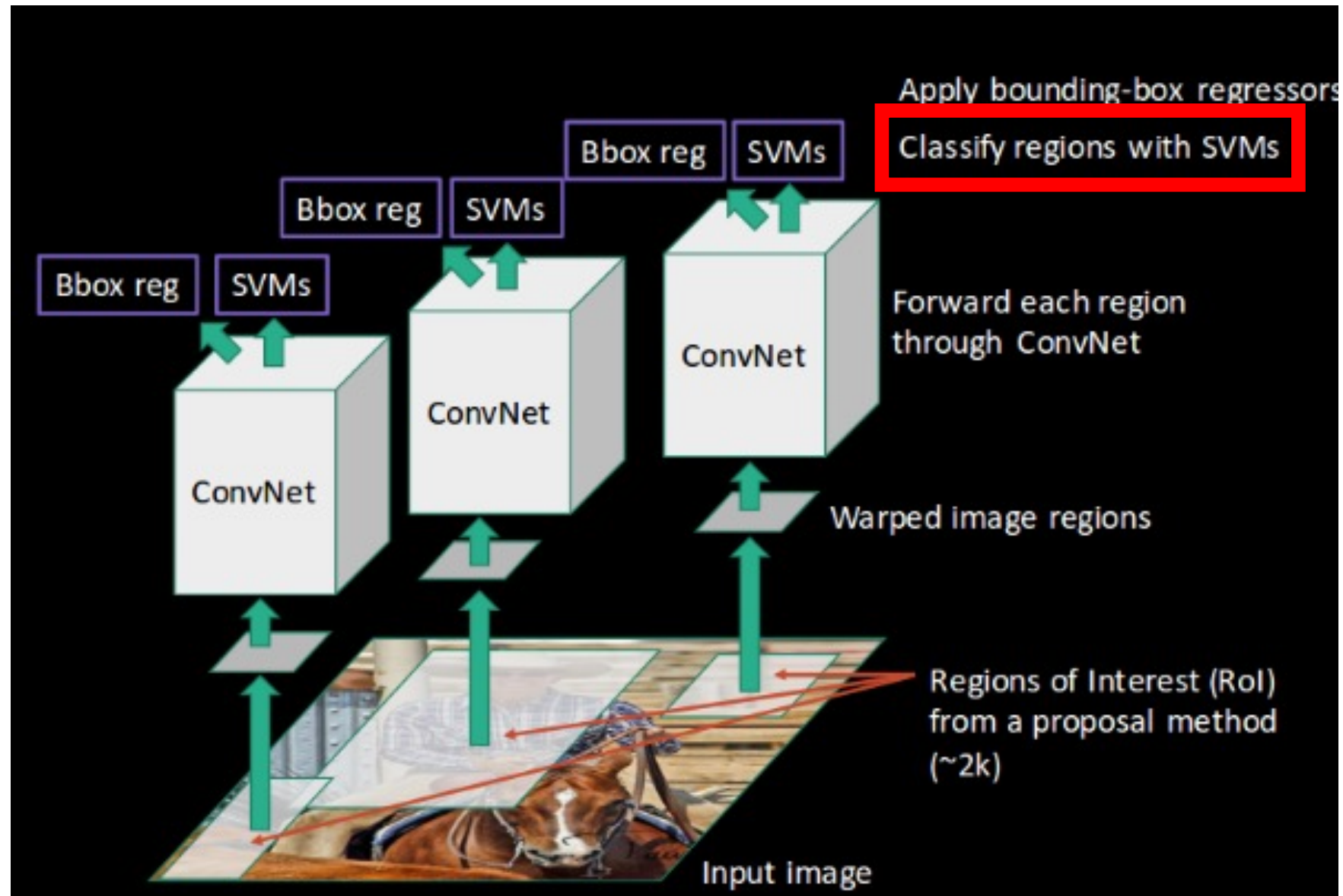


# Region Classification

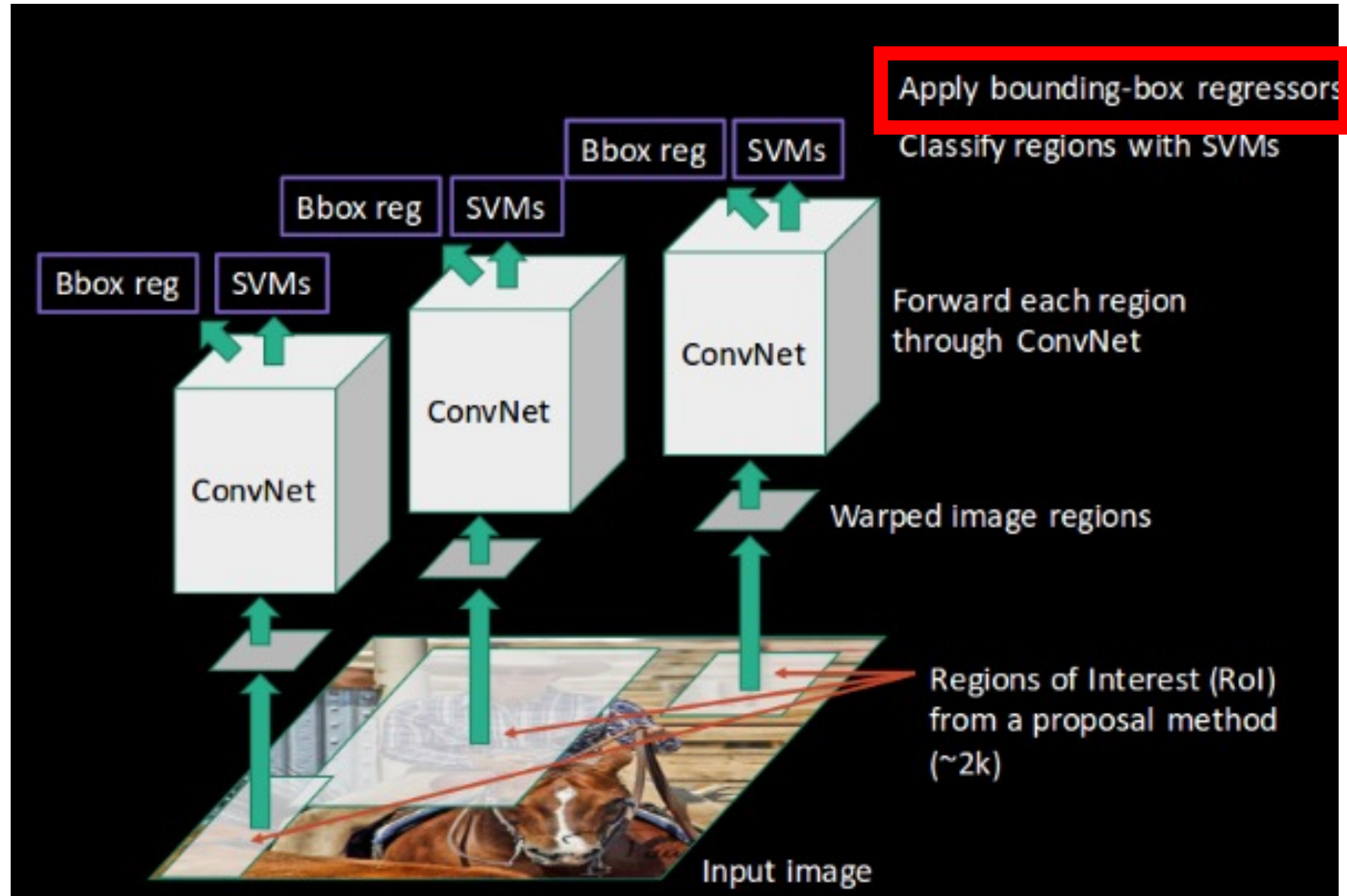
- Assign each feature descriptor that characterizes a region a label from a pre-defined set of categories (i.e., multiple choice)



# Architecture

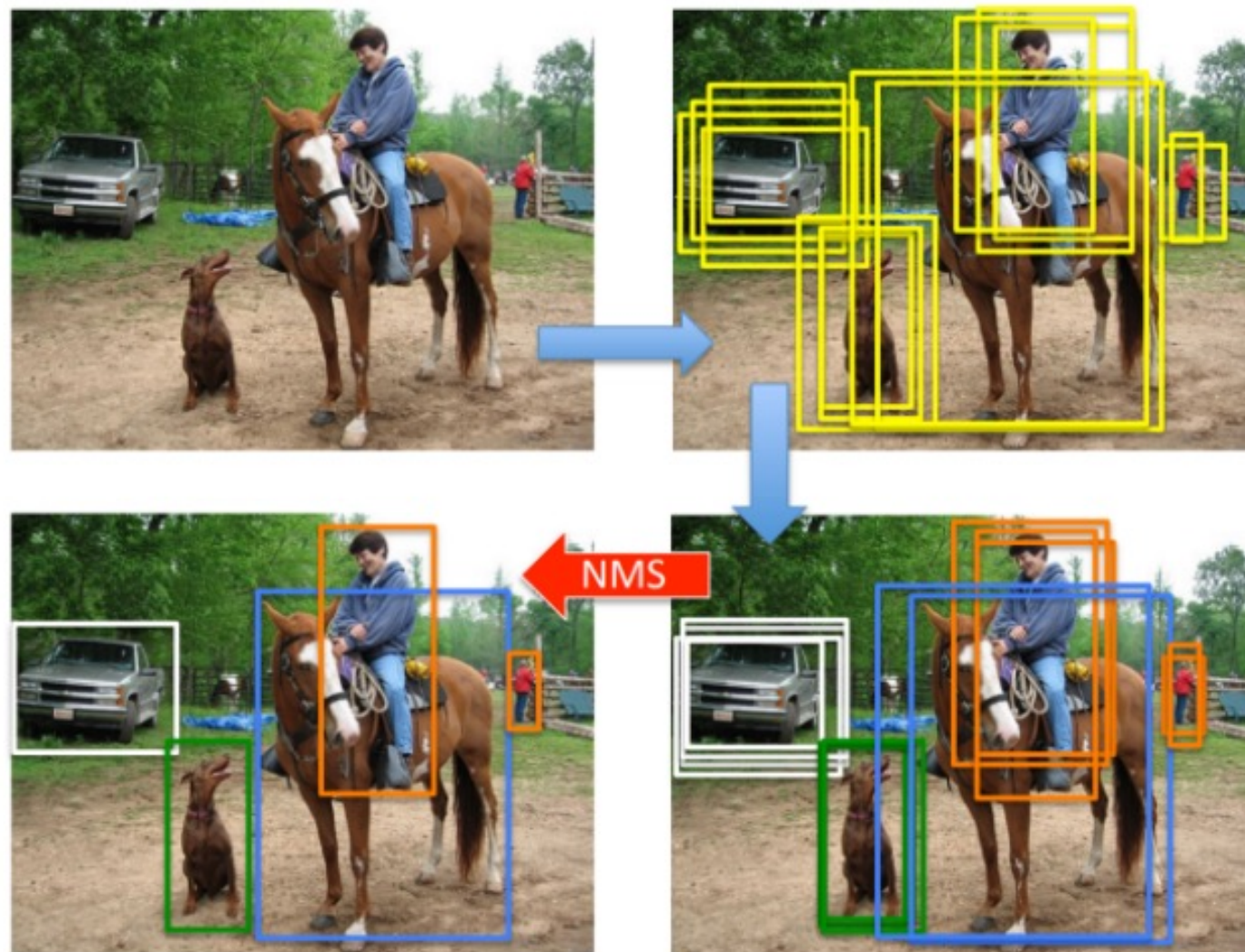


# Architecture



# Region Selection and Refinement

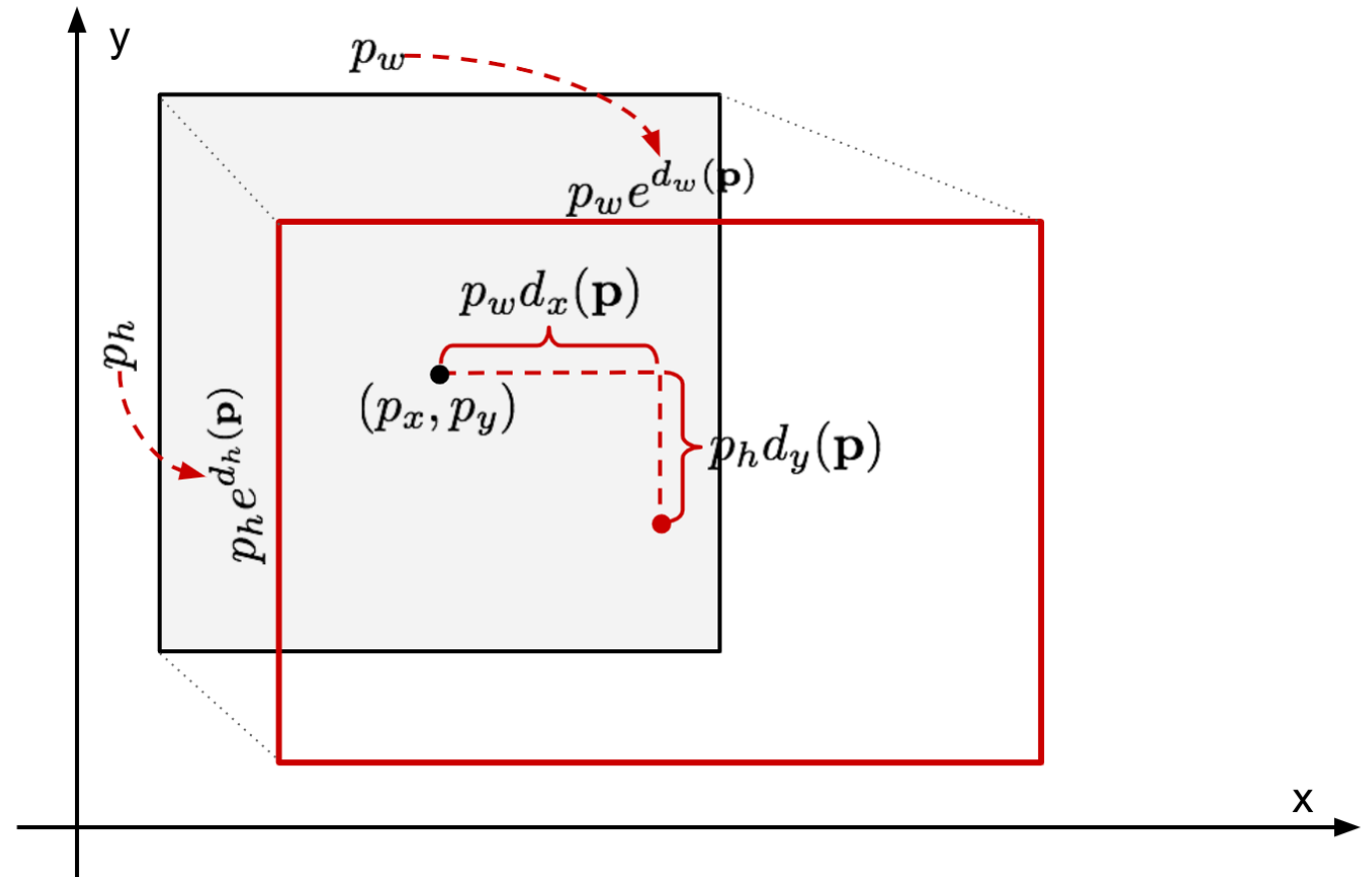
- Problem: ~2000 regions per image
- Solution: remove redundant regions through **non-maximum suppression**; for each class:
  1. Pick region with maximum score obtained from the SVM.
  2. Discard all regions belonging to that class nearby that chosen region (i.e., IoU score > 70%)
  3. Select next highest score region and then repeat steps 1 and 2
  4. Repeat step 3 until all regions are either discarded or kept





# Region Selection and Refinement

Observing that a common issue was imperfect region proposals, transformations were learned to convert each region proposal to more closely match ground truth



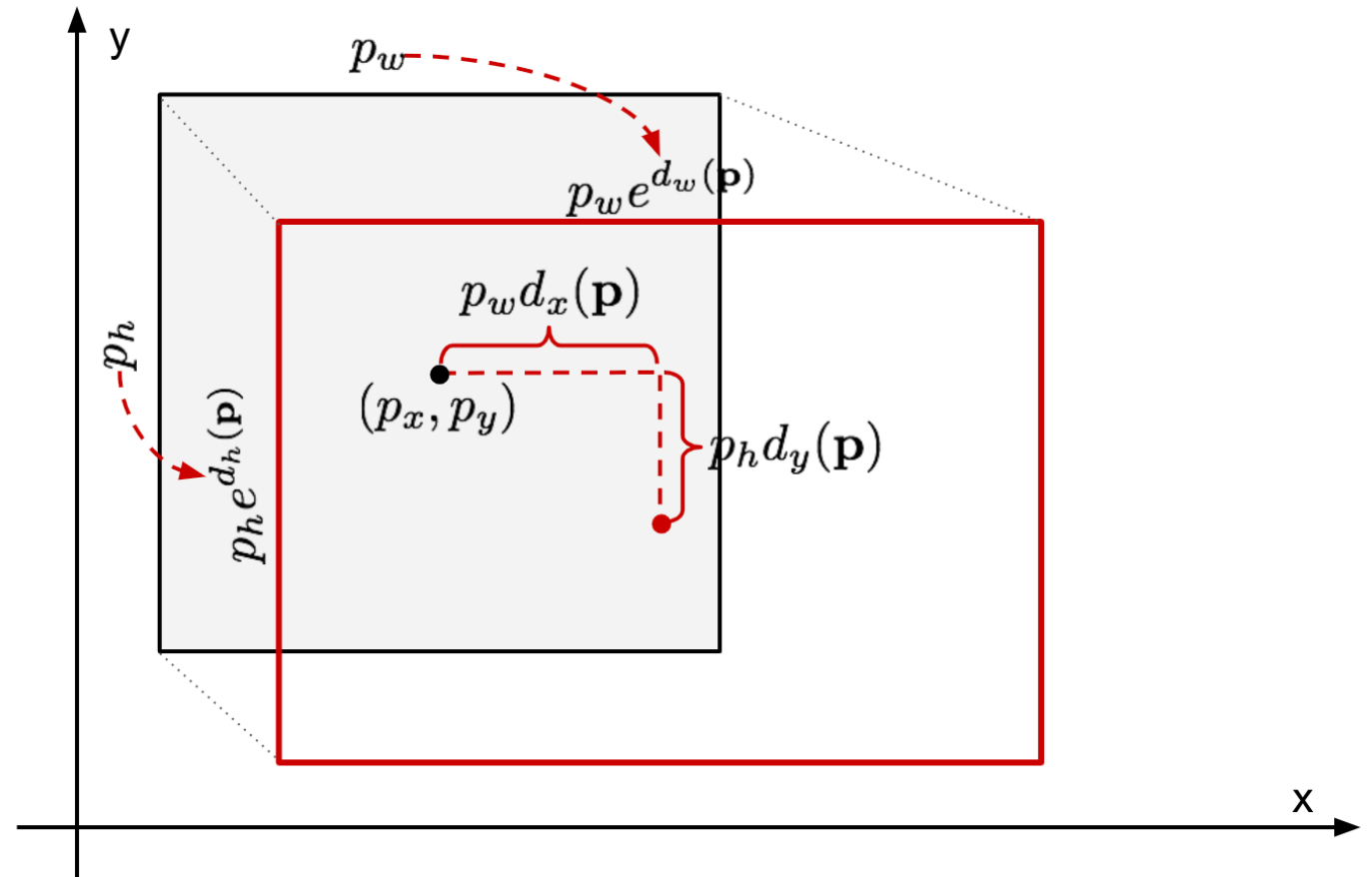
# Region Selection and Refinement

- **Input:** original region location described by a center  $(p_x, p_y)$ , width  $(p_w)$ , and height  $(p_h)$
- **Output:** four refinement functions:  $d_x, d_y, d_w, d_h$
- Loss function for learning: SSE

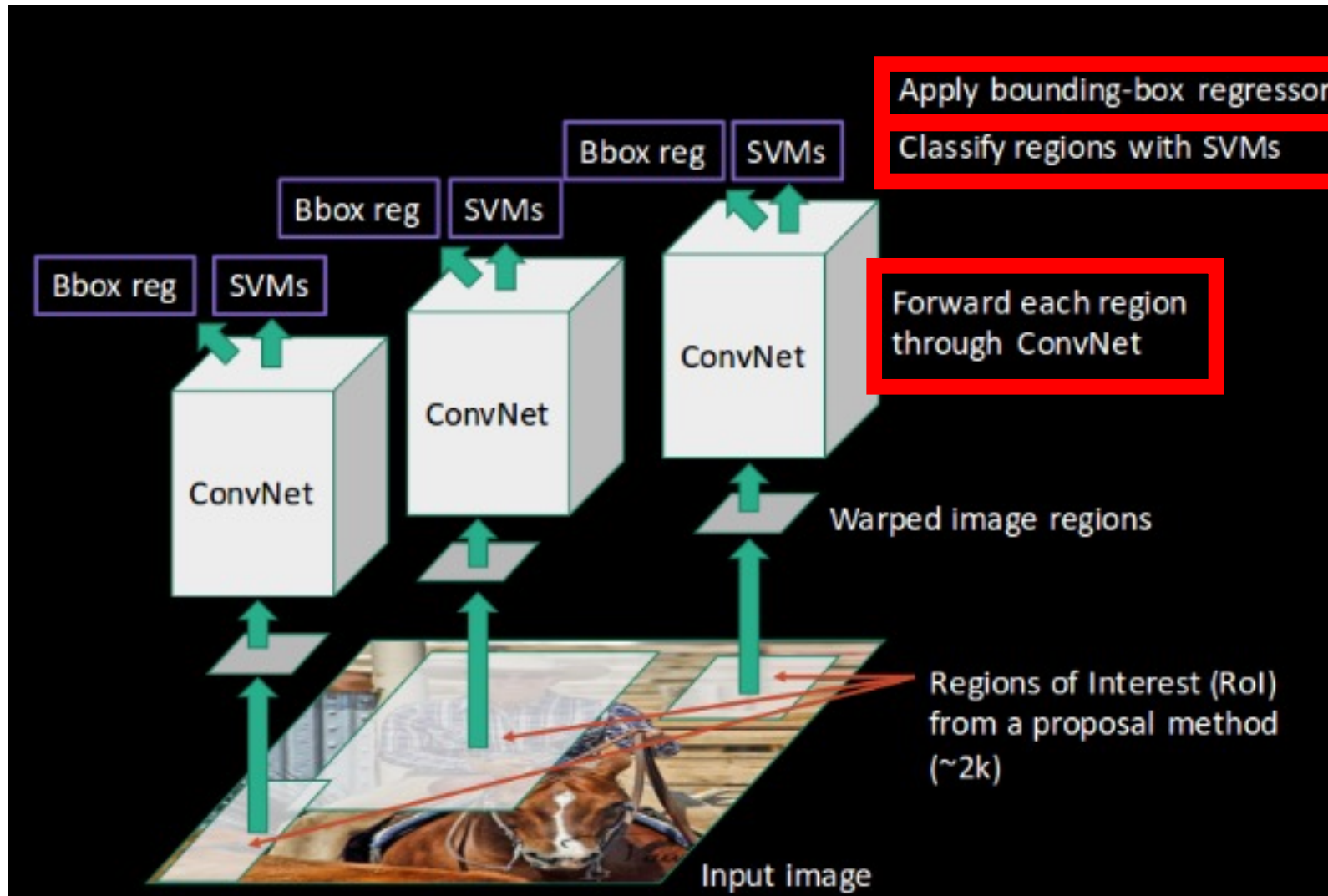
$$\sum_{i \in \{x, y, w, h\}} (t_i - d_i(\mathbf{p}))^2$$

True location

Predicted location



# R-CNN Limitations



- Slow training procedure
  - Must train three models
- Slow at test time (~1 minute per image)
- Inefficient/complex architecture
  - Must store feature descriptor for each region proposal
  - Must refine initial region proposals

# Key Concluding Remarks

1. Deep CNN features for image subregions are valuable  
*(recall, deep CNN features for entire images were also deemed important for scene classification)*

2. “We conjecture that the ‘supervised pre-training/domain-specific finetuning’ paradigm will be highly effective for a variety of data-scarce vision problems.”

# Object Detection: Today's Topics

- Problem
- Applications
- Datasets
- Evaluation metric
- Overview of object detection algorithms and baseline (R-CNN)

A dark gray background with a central circular glow. The glow is a gradient from light gray in the center to dark gray at the edges. The text "The End" is centered within this glow. The text is in a white, cursive font with a slight drop shadow. The entire scene is framed by a white film strip border with rectangular sprocket holes on the left and right sides.

*The End*