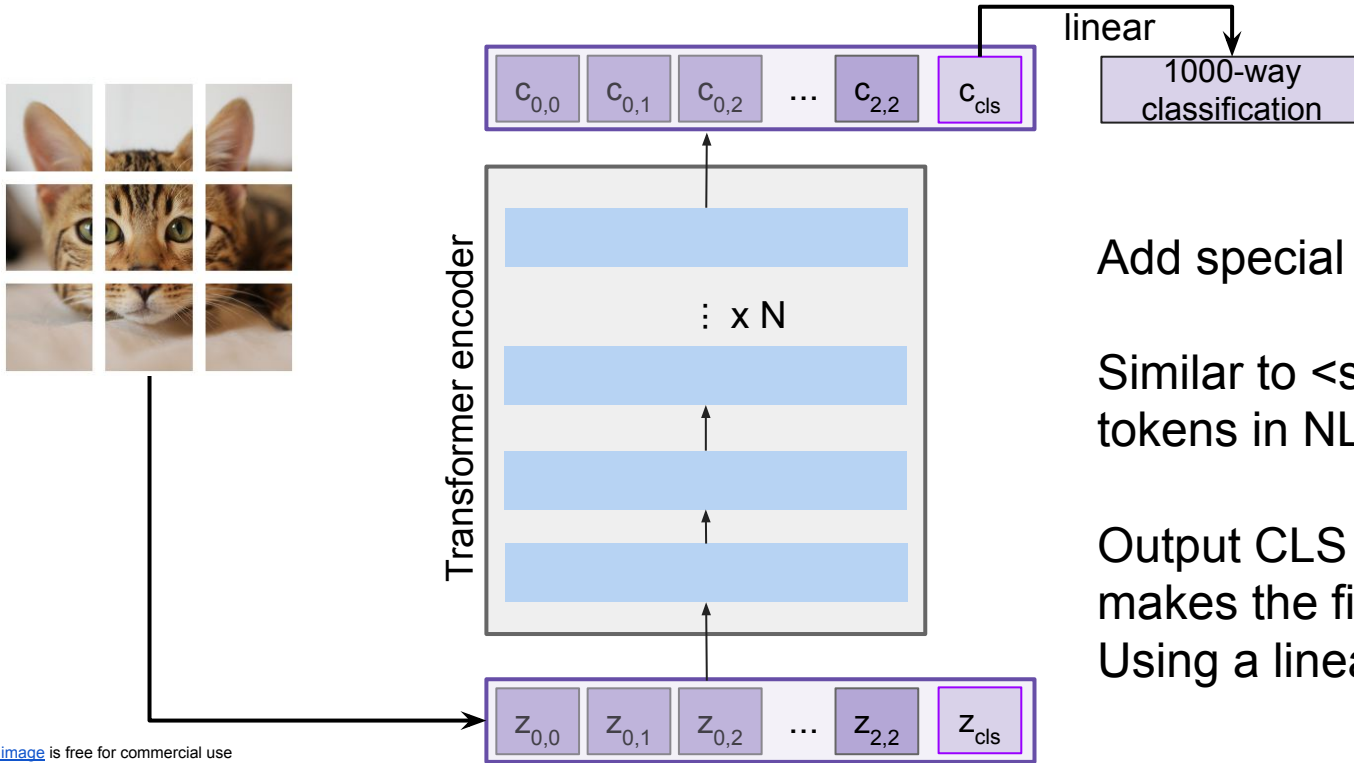# Lecture 13:
# Structured Prediction
## Detection and Segmentation

# Administrative

- A4 posted (Check ed, website still not working)
- Quiz 3 friday

# How to turn the output to a class prediction?



Add special [CLS] token.

Similar to <start> and <end> tokens in NLP.

Output CLS representation makes the final prediction Using a linear layer

Cat image is free for commercial use
Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

# Common ViT architectures

| Model | Layers | Hidden size $D$ | MLP size | Heads | Params |
|-------|--------|-----------------|----------|-------|--------|
| ViT-Base | 12 | 768 | 3072 | 12 | 86M |
| ViT-Large | 24 | 1024 | 4096 | 16 | 307M |
| ViT-Huge | 32 | 1280 | 5120 | 16 | 632M |

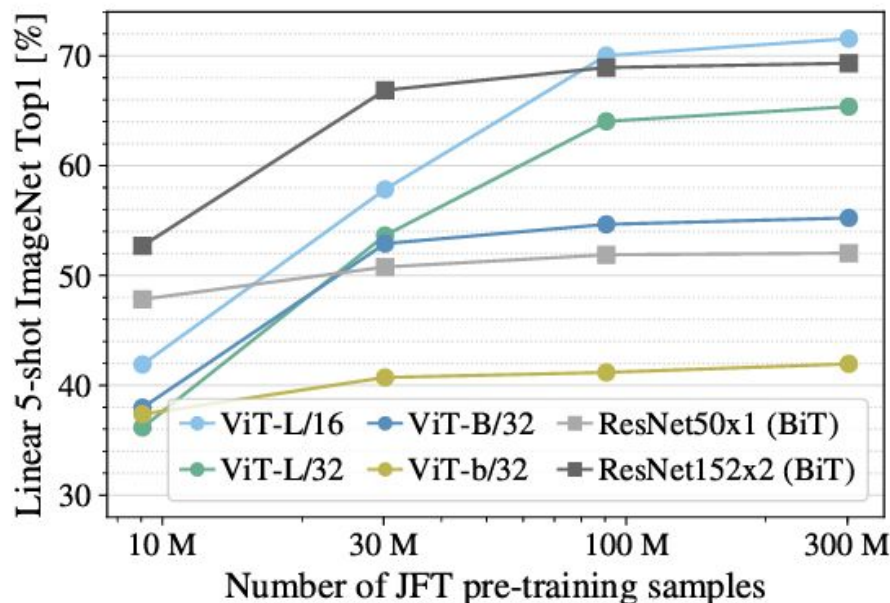**Common patch sizes**: 32, 16, 14…
Smaller patches results in larger more powerful models.

**Nomenclature**: ViT-B/32 means that its a ViT model that uses Base values for layers, hidden size, mlp vize, and head. /32 means the input image patches are 32x32.

Dosovitskiyet al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021
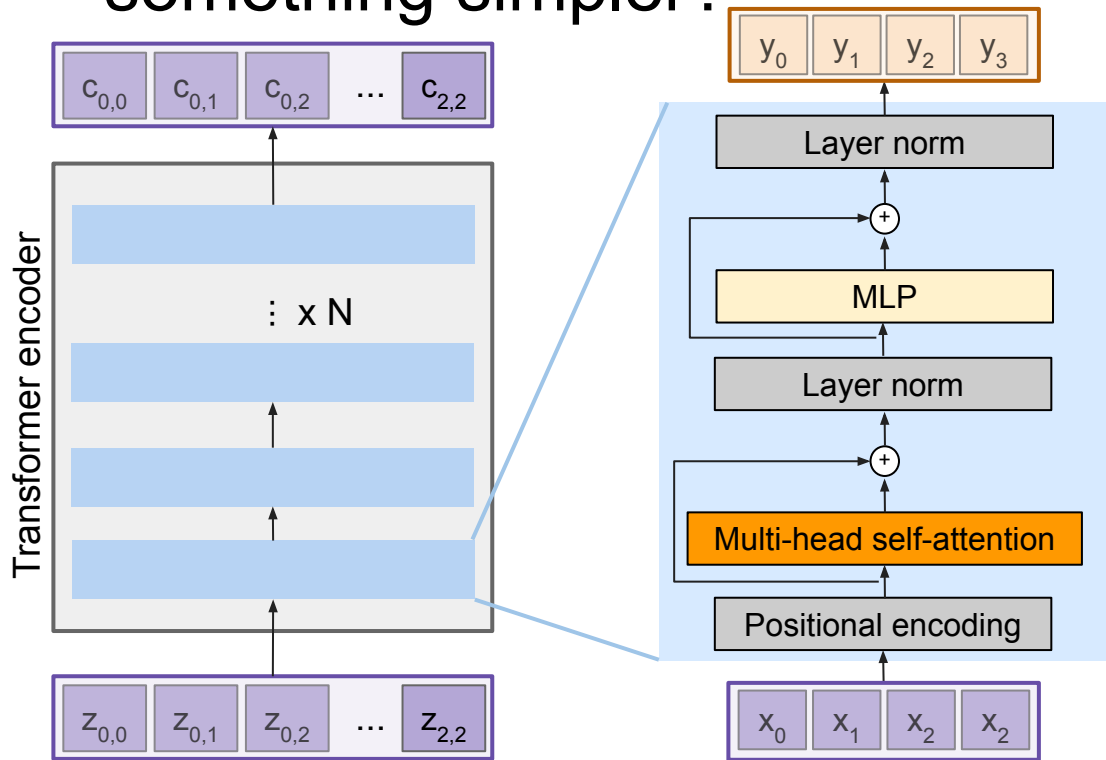
# Comparing ResNets with ViTs



Models are initially trained on a large dataset called JFT-300M

And then the last linear layer is finetuned on ImageNet-1.5M

ViT performs worse when only 10M images are used from JFT. But ViT outperforms ResNets with larger training data (300M images from JFT).
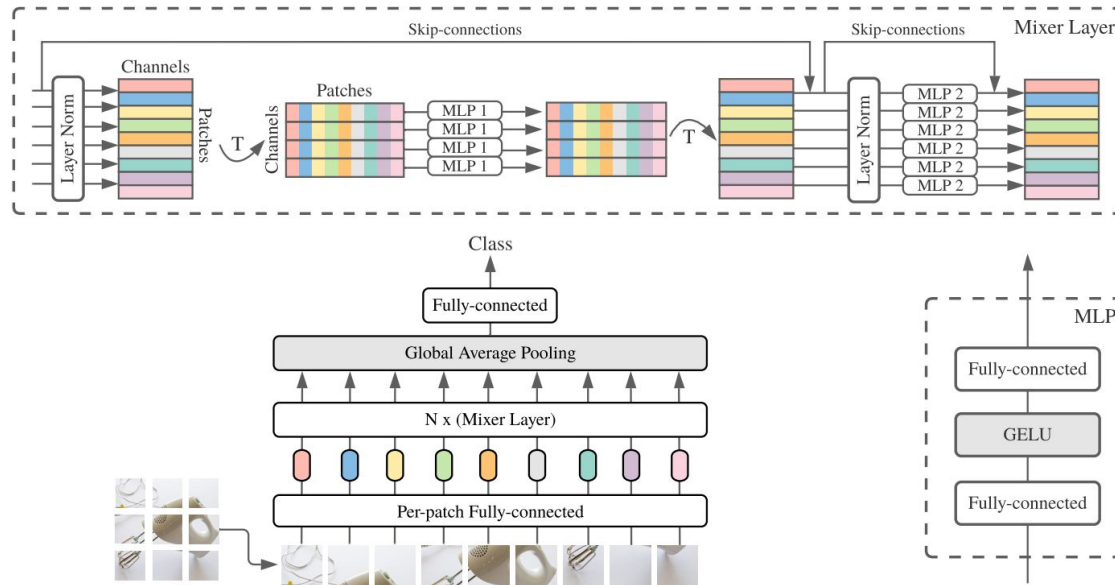
Dosovitskiyet al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

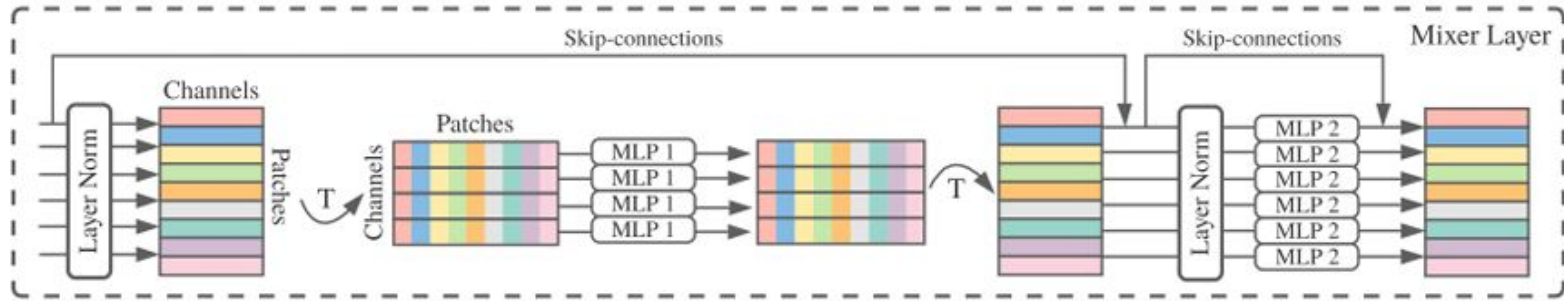# Self-attention is expensive… can we design something simpler?



Every self-attention is expensive. We want each input to "interact" with other tokens but can we simplify the operation a bit?

Vaswani et al, "Attention is all you need", NeurIPS 2017

# MLP-Mixer: an all-MLP architecture



Tolstikhin et al, "MLP-Mixer: An all-MLP architecture for vision", NeurIPS2021

# MLP-Mixer: an all-MLP architecture
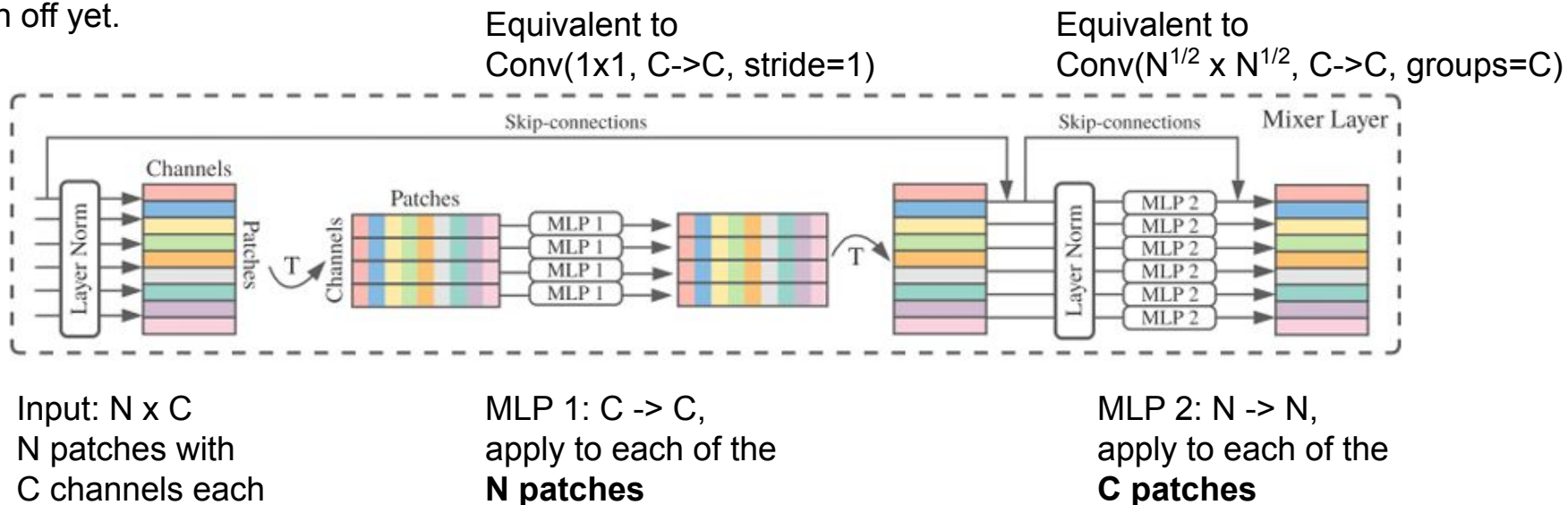


Input: N x C
N patches with
C channels each

MLP 1: C -> C,
apply to each of the
**N patches**

MLP 2: N -> N,
apply to each of the
**C patches**

Tolstikhinet al, "MLP-Mixer: An all-MLP architecture for vision", NeurIPS2021

# MLP-Mixer: The MLPs are sort of like **convs**

Cool idea; but hasn't taken off yet.

Equivalent to
Conv(1x1, C->C, stride=1)

Equivalent to
Conv($N^{1/2}$ x $N^{1/2}$, C->C, groups=C)



Input: N x C
N patches with
C channels each

MLP 1: C -> C,
apply to each of the
**N patches**

MLP 2: N -> N,
apply to each of the
**C patches**

Tolstikhin et al, "MLP-Mixer: An all-MLP architecture for vision", NeurIPS2021

# MLP-Mixer: Many concurrent and followups

Touvronet al, "ResMLP: Feedforward Networks for Image Classification with Data-Efficient Training", arXiv2021,
https://arxiv.org/abs/2105.03404

Tolstikhinet al, "MLP-Mixer: An all-MLP architecture for vision", NeurIPS2021,
https://arxiv.org/abs/2105.01601

Liu et al, "Pay Attention to MLPs", NeurIPS2021, https://arxiv.org/abs/2105.08050

Yu et al, "S2-MLP: Spatial-Shift MLP Architecture for Vision", WACV 2022,
https://arxiv.org/abs/2106.07477

Chen et al, "CycleMLP: AMLP-like Architecture for Dense Prediction", ICLR 2022,
https://arxiv.org/abs/2107.10224
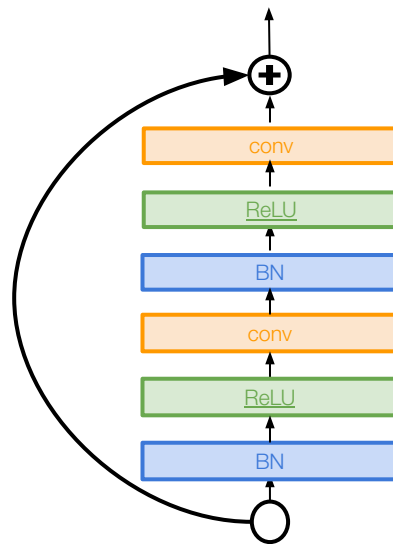
But research has continued since ImageNet
(Will go over following slides if time,
Otherwise skip to the summary slides in the end.)

# Improving ResNets...

# Identity Mappings in Deep Residual Networks
*[He et al. 2016]*

- Improved ResNet block design from creators of ResNet
- Creates a more direct path for propagating information throughout network
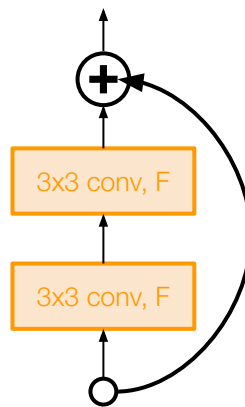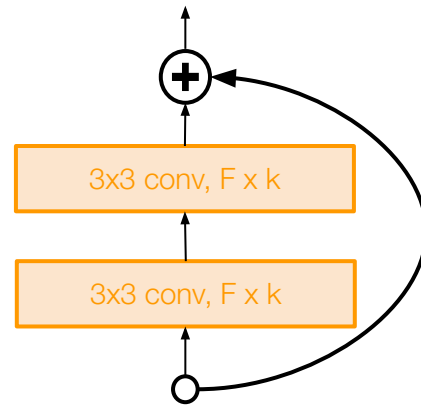- Gives better performance

# Improving ResNets...

# Wide Residual Networks

*[Zagoruyko et al. 2016]*

- Argues that residuals are the important factor, not depth
- User wider residual blocks (F x k filters instead of F filters in each layer)
- 50-layer wide ResNet outperforms 152-layer original ResNet
- Increasing width instead of depth more computationally efficient (parallelizable)
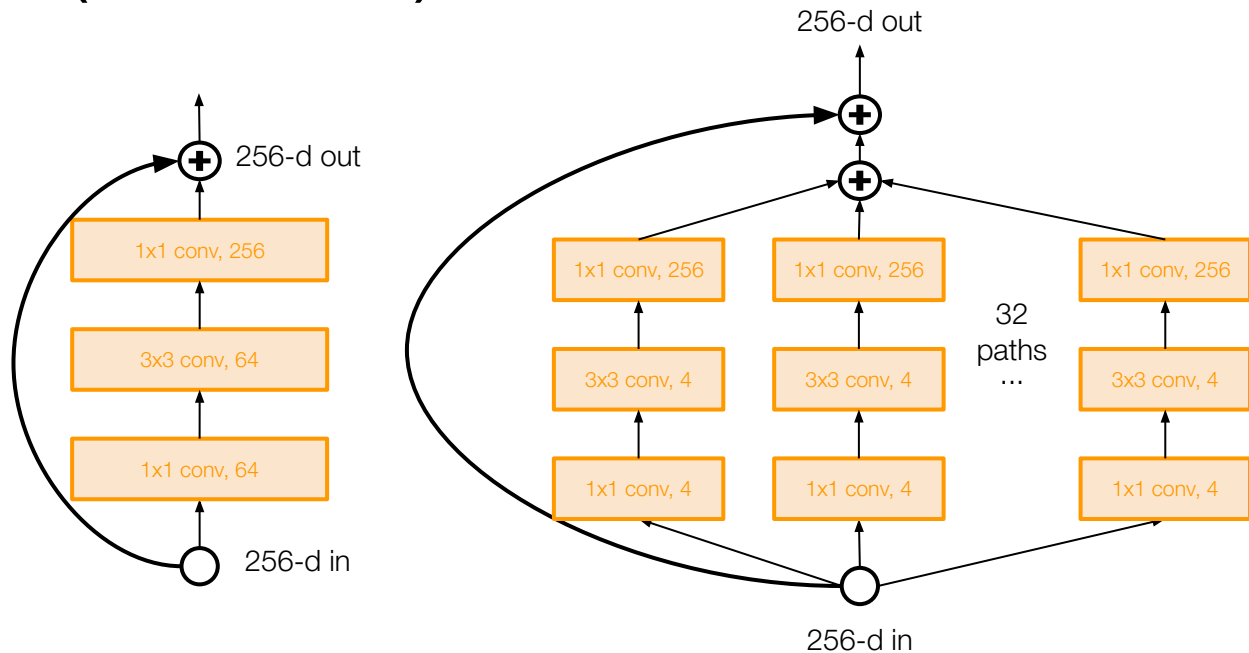


Basic residual block          Wide residual block

# Improving ResNets...
# Aggregated Residual Transformations for Deep Neural Networks (ResNeXt)

*[Xie et al. 2016]*

- Also from creators of ResNet
- Increases width of residual block through multiple parallel pathways ("cardinality")
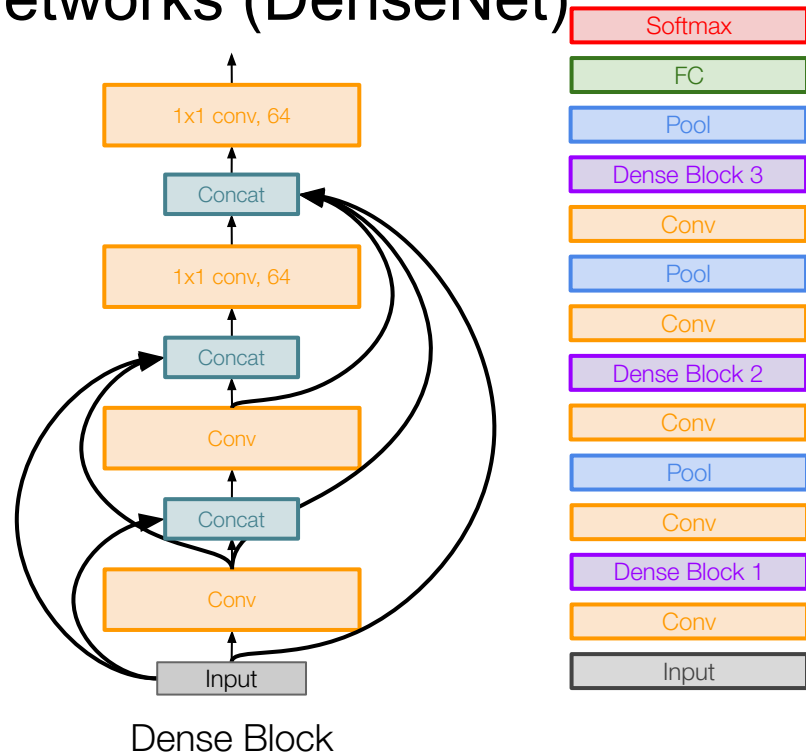- Parallel pathways similar in spirit to Inception module

# Other ideas...

# Densely Connected Convolutional Networks (DenseNet)
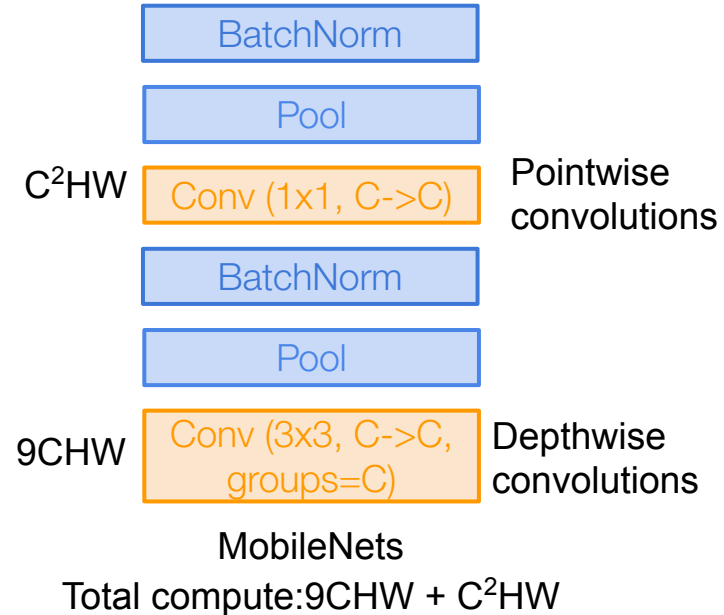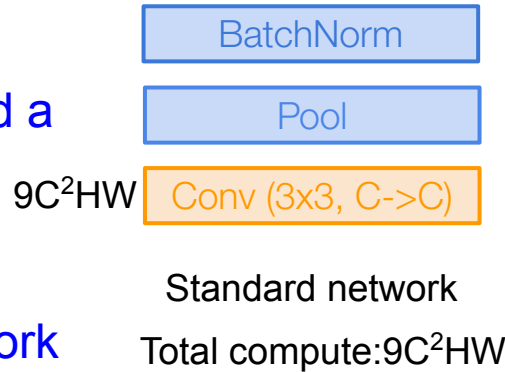
*[Huang et al. 2017]*

- Dense blocks where each layer is connected to every other layer in feedforward fashion
- Alleviates vanishing gradient, strengthens feature propagation, encourages feature reuse
- Showed that shallow 50-layer network can outperform deeper 152 layer ResNet



Dense Block

# Efficient networks...

# MobileNets: Efficient Convolutional Neural Networks for Mobile Applications *[Howard et al. 2017]*

- Depthwise separable convolutions replace standard convolutions by factorizing them into a depthwise convolution and a 1x1 convolution
- Much more efficient, with little loss in accuracy
- Follow-up MobileNetV2 work in 2018 (Sandler et al.)
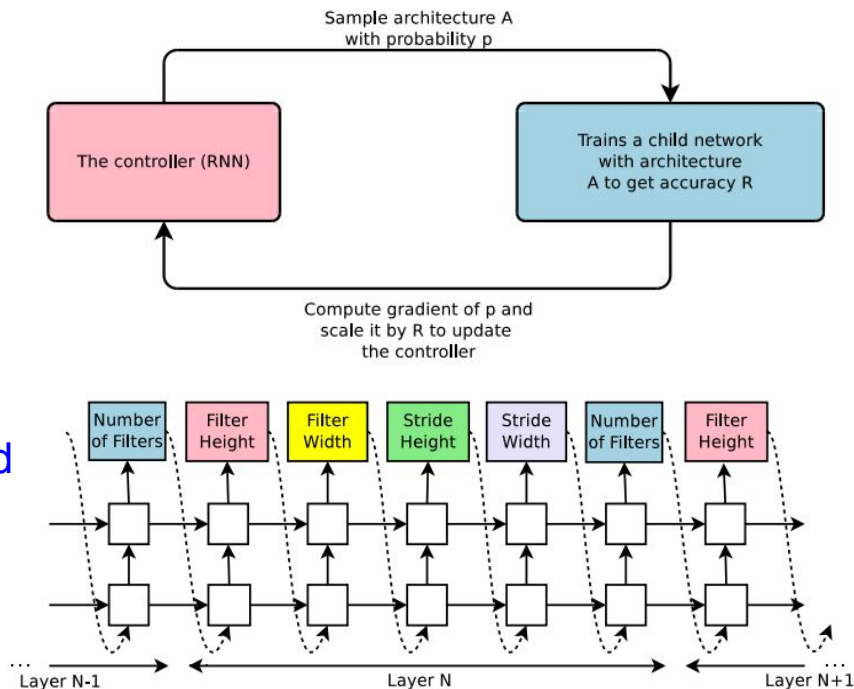- ShuffleNet: Zhang et al, CVPR 2018

**Standard network**

BatchNorm

Pool

$9C^2HW$ — Conv (3x3, C->C)

Standard network

Total compute: $9C^2HW$

**MobileNets**

BatchNorm

Pool

$C^2HW$ — Conv (1x1, C->C) — Pointwise convolutions

BatchNorm

Pool

$9CHW$ — Conv (3x3, C->C, groups=C) — Depthwise convolutions

MobileNets

Total compute: $9CHW + C^2HW$

# Learning to search for network architectures...

# Neural Architecture Search with Reinforcement Learning (NAS)

*[Zoph et al. 2016]*

- "Controller" network that learns to design a good network architecture (output a string corresponding to network design)
- Iterate:
  1) Sample an architecture from search space
  2) Train the architecture to get a "reward" R corresponding to accuracy
  3) Compute gradient of sample probability, and scale by R to perform controller parameter update (i.e. increase likelihood of good architecture being sampled, decrease likelihood of bad architecture)
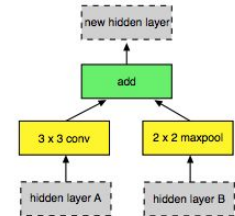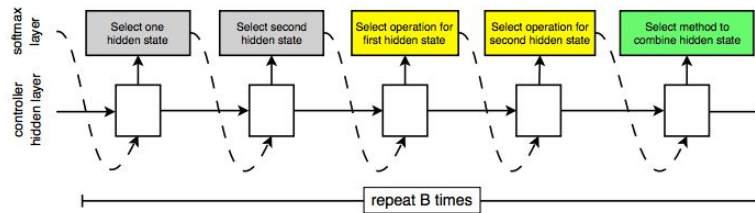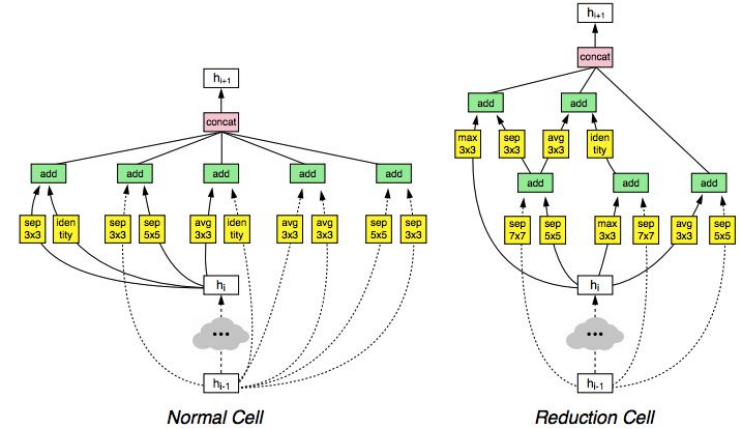
# Learning to search for network architectures...

# Learning Transferable Architectures for Scalable Image Recognition

*[Zoph et al. 2017]*

- Applying neural architecture search (NAS) to a large dataset like ImageNet is expensive
- Design a search space of building blocks ("cells") that can be flexibly stacked
- NASNet: Use NAS to find best cell structure on smaller CIFAR-10 dataset, then transfer architecture to ImageNet

- Many follow-up works in this space e.g. AmoebaNet (Real et al. 2019) and ENAS (Pham, Guan et al. 2018)



Normal Cell

Reduction Cell



repeat B times

# But sometimes smart heuristic is better than NAS ...

## EfficientNet: Smart Compound Scaling

*[Tan and Le. 2019]*

- Increase network capacity by scaling width, depth, and resolution, while balancing accuracy and efficiency.
- Search for optimal set of compound scaling factors given a compute budget (target memory & flops).
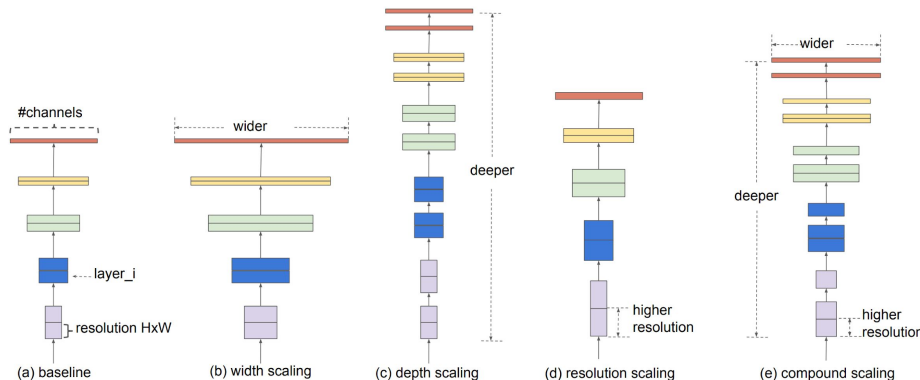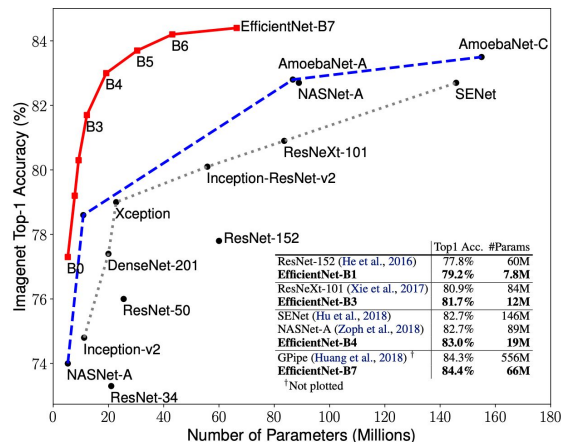- Scale up using smart heuristic rules

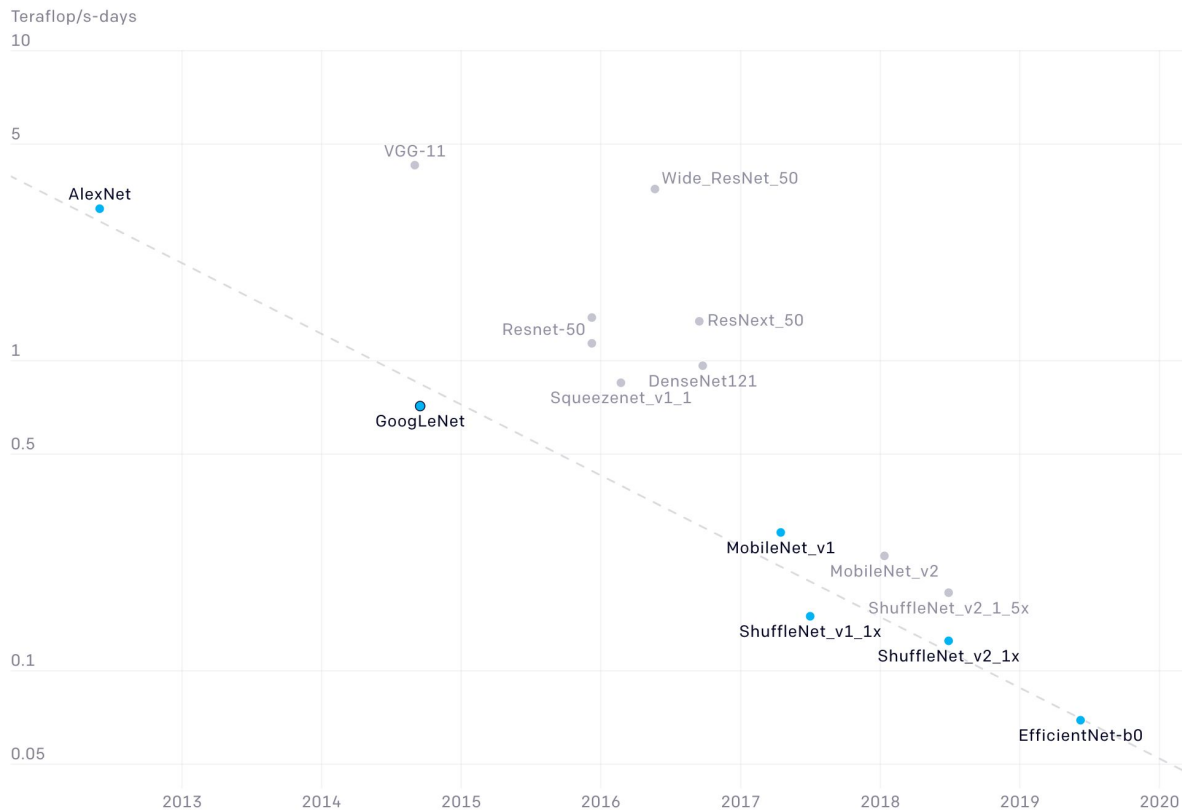$$\text{depth: } d = \alpha^\phi$$
$$\text{width: } w = \beta^\phi$$
$$\text{resolution: } r = \gamma^\phi$$
$$\text{s.t. } \alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$$
$$\alpha \geq 1, \beta \geq 1, \gamma \geq 1$$



| | Top1 Acc. | #Params |
|---|---|---|
| ResNet-152 (He et al., 2016) | 77.8% | 60M |
| **EfficientNet-B1** | **79.2%** | **7.8M** |
| ResNeXt-101 (Xie et al., 2017) | 80.9% | 84M |
| **EfficientNet-B3** | **81.7%** | **12M** |
| SENet (Hu et al., 2018) | 82.7% | 146M |
| NASNet-A (Zoph et al., 2018) | 82.7% | 89M |
| **EfficientNet-B4** | **83.0%** | **19M** |
| GPipe (Huang et al., 2018) [†] | 84.3% | 556M |
| **EfficientNet-B7** | **84.4%** | **66M** |

[†]Not plotted



(a) baseline  (b) width scaling  (c) depth scaling  (d) resolution scaling  (e) compound scaling

# Efficient networks...



Teraflop/s-days

https://openai.com/blog/ai-and-efficiency/

# Summary: Modern Architectures

## Case Studies

- AlexNet
- VGG
- GoogLeNet
- ResNet
- Vit
- MLP-Mixer

## Also....

- SENet
- Wide ResNet
- ResNeXT
- DenseNet
- MobileNets
- NASNet

# Main takeaways

**AlexNet** showed that you can use CNNs to train Computer Vision models.

**ZFNet**, **VGG** shows that bigger networks work better

**GoogLeNet** is one of the first to focus on efficiency using 1x1 bottleneck convolutions and global avg pool instead of FC layers

**ResNet** showed us how to train extremely deep networks

- Limited only by GPU & memory!
- Showed diminishing returns as networks got bigger

After ResNet: CNNs were better than the human metric and focus shifted to Efficient networks:

- Lots of tiny networks aimed at mobile devices: **MobileNet**, **ShuffleNet**

**Neural Architecture Search** can now automate architecture design

**ViT** is the current favorite architecture but requires a lot of compute and data

**MLP-Mixers** have presented an alternative to transformers but they haven't taken off.

# Summary: Modern Architectures

- **ResNet-50** and **ViT** currently good defaults to use
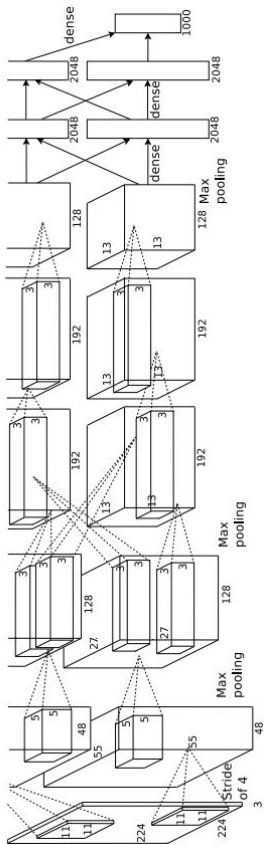
- Next time: Structure prediction

Transfer learning

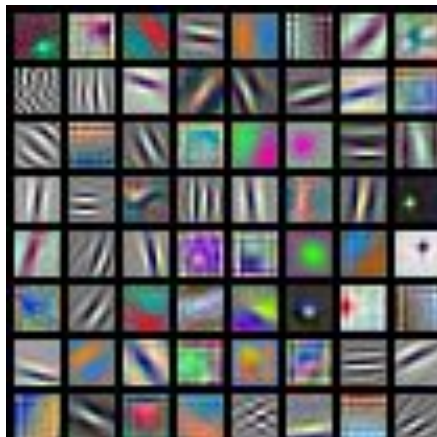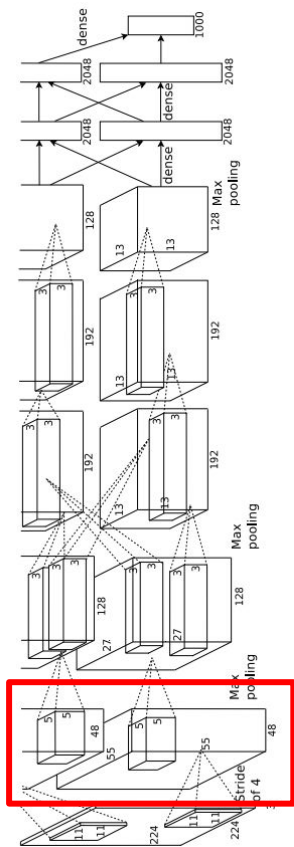"You need a lot of a data if you want to train/use CNNs"

"You need a lot of a data if you want to train/use CNNs"
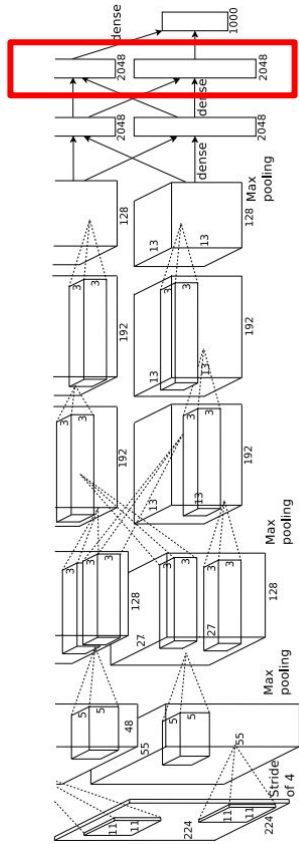
**BUSTED**

# Transfer Learning with CNNs

# Transfer Learning with CNNs
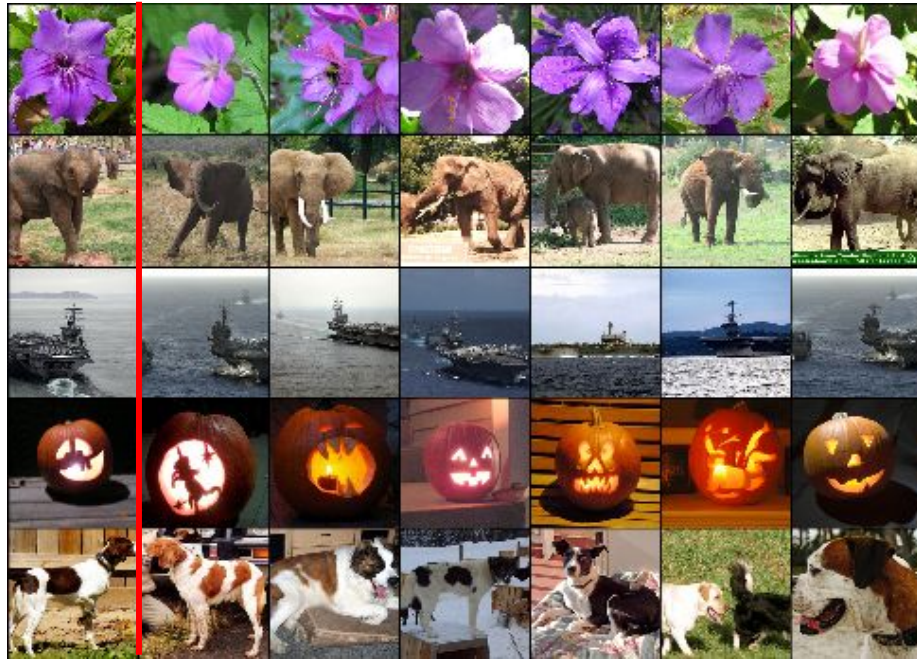


AlexNet:
64 x 3 x 11 x 11

(More on this in Lecture 13)

# Transfer Learning with CNNs



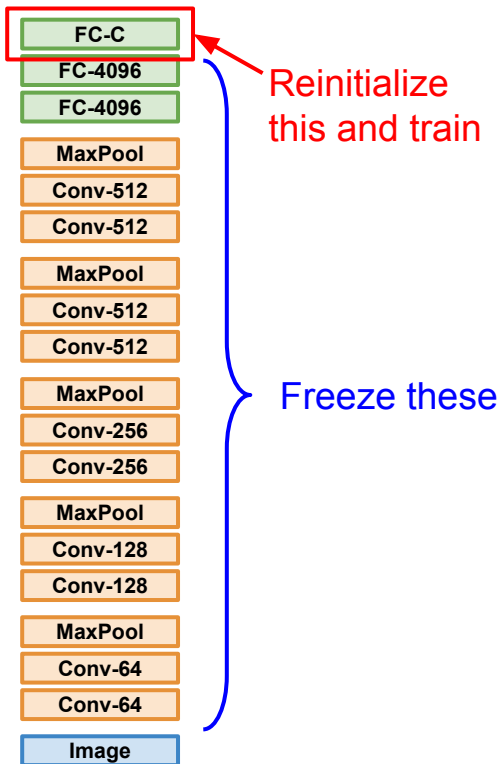Test image     L2 Nearest neighbors in <u>feature</u> space

# Transfer Learning with CNNs

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

1. Train on Imagenet

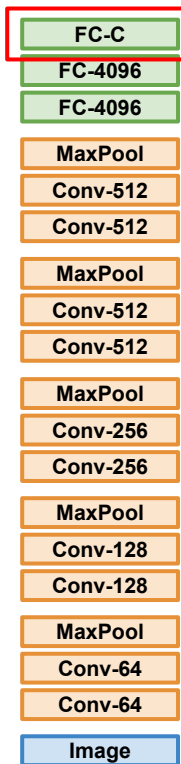| |
|---|
| FC-1000 |
| FC-4096 |
| FC-4096 |
| |
| MaxPool |
| Conv-512 |
| Conv-512 |
| |
| MaxPool |
| Conv-512 |
| Conv-512 |
| |
| MaxPool |
| Conv-256 |
| Conv-256 |
| |
| MaxPool |
| Conv-128 |
| Conv-128 |
| |
| MaxPool |
| Conv-64 |
| Conv-64 |
| |
| Image |

# Transfer Learning with CNNs

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014
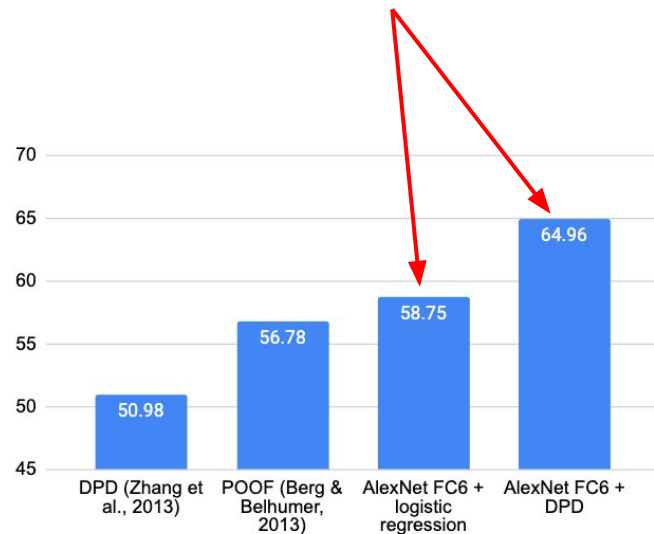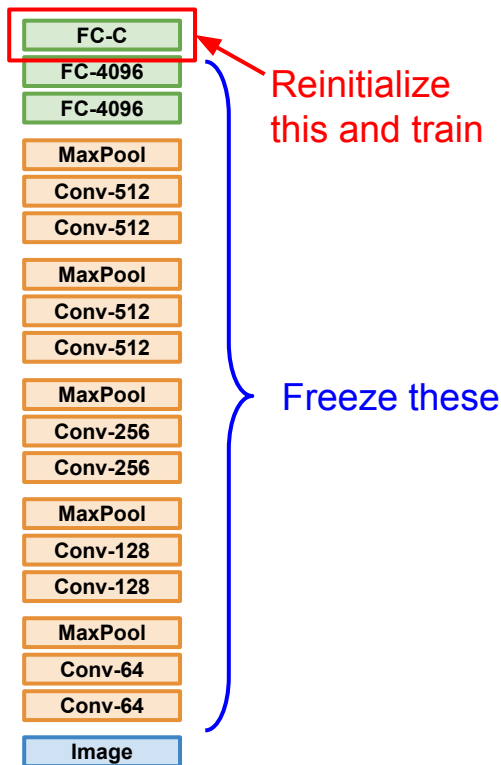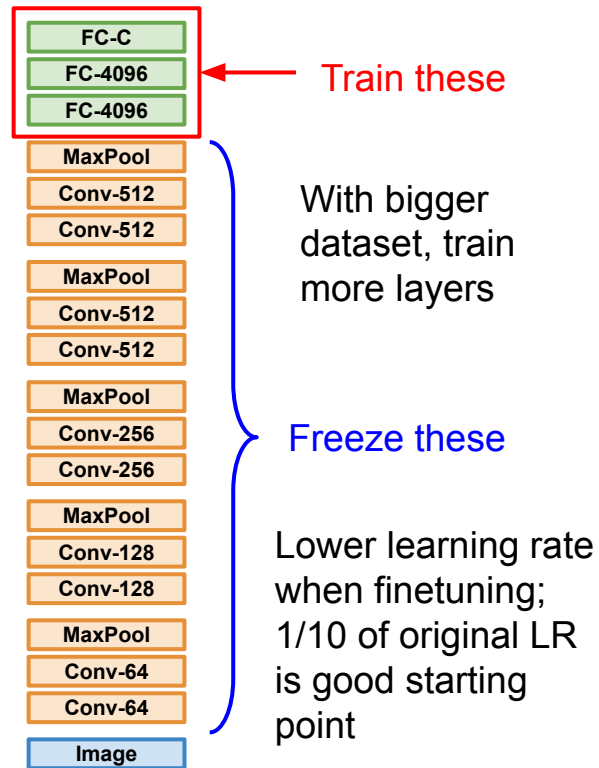
## 1. Train on Imagenet

| FC-1000 |
| FC-4096 |
| FC-4096 |

| MaxPool |
| Conv-512 |
| Conv-512 |

| MaxPool |
| Conv-512 |
| Conv-512 |

| MaxPool |
| Conv-256 |
| Conv-256 |

| MaxPool |
| Conv-128 |
| Conv-128 |

| MaxPool |
| Conv-64 |
| Conv-64 |

| Image |

## 2. Small Dataset (C classes)

| FC-C |
| FC-4096 |
| FC-4096 |

| MaxPool |
| Conv-512 |
| Conv-512 |

| MaxPool |
| Conv-512 |
| Conv-512 |

| MaxPool |
| Conv-256 |
| Conv-256 |

| MaxPool |
| Conv-128 |
| Conv-128 |

| MaxPool |
| Conv-64 |
| Conv-64 |

| Image |

Reinitialize this and train

Freeze these

# Transfer Learning with CNNs

## 1. Train on Imagenet

| FC-1000 |
| FC-4096 |
| FC-4096 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-256 |
| Conv-256 |
| MaxPool |
| Conv-128 |
| Conv-128 |
| MaxPool |
| Conv-64 |
| Conv-64 |
| Image |

## 2. Small Dataset (C classes)

| FC-C |
| FC-4096 |
| FC-4096 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-256 |
| Conv-256 |
| MaxPool |
| Conv-128 |
| Conv-128 |
| MaxPool |
| Conv-64 |
| Conv-64 |
| Image |

Reinitialize this and train

Freeze these

Finetuned from AlexNet



Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014

# Transfer Learning with CNNs

## 1. Train on Imagenet

| FC-1000 |
| --- |
| FC-4096 |
| FC-4096 |

| MaxPool |
| --- |
| Conv-512 |
| Conv-512 |

| MaxPool |
| --- |
| Conv-512 |
| Conv-512 |

| MaxPool |
| --- |
| Conv-256 |
| Conv-256 |

| MaxPool |
| --- |
| Conv-128 |
| Conv-128 |

| MaxPool |
| --- |
| Conv-64 |
| Conv-64 |

| Image |
| --- |

## 2. Small Dataset (C classes)

| FC-C |
| --- |
| FC-4096 |
| FC-4096 |

**Reinitialize this and train**

| MaxPool |
| --- |
| Conv-512 |
| Conv-512 |

| MaxPool |
| --- |
| Conv-512 |
| Conv-512 |

| MaxPool |
| --- |
| Conv-256 |
| Conv-256 |

| MaxPool |
| --- |
| Conv-128 |
| Conv-128 |

| MaxPool |
| --- |
| Conv-64 |
| Conv-64 |

| Image |
| --- |

**Freeze these**

## 3. Bigger dataset

| FC-C |
| --- |
| FC-4096 |
| FC-4096 |

**Train these**

| MaxPool |
| --- |
| Conv-512 |
| Conv-512 |

| MaxPool |
| --- |
| Conv-512 |
| Conv-512 |

| MaxPool |
| --- |
| Conv-256 |
| Conv-256 |

| MaxPool |
| --- |
| Conv-128 |
| Conv-128 |

| MaxPool |
| --- |
| Conv-64 |
| Conv-64 |

| Image |
| --- |

With bigger dataset, train more layers

**Freeze these**

Lower learning rate when finetuning; 1/10 of original LR is good starting point

| FC-1000 |
| FC-4096 |
| FC-4096 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-256 |
| Conv-256 |
| MaxPool |
| Conv-128 |
| Conv-128 |
| MaxPool |
| Conv-64 |
| Conv-64 |
| Image |

More specific

More generic

|  | **very similar dataset** | **very different dataset** |
|---|---|---|
| **very little data** | ? | ? |
| **quite a lot of data** | ? | ? |

| FC-1000 |
| FC-4096 |
| FC-4096 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-512 |
| Conv-512 |
| MaxPool |
| Conv-256 |
| Conv-256 |
| MaxPool |
| Conv-128 |
| Conv-128 |
| MaxPool |
| Conv-64 |
| Conv-64 |
| Image |

More specific

More generic

|  | very similar dataset | very different dataset |
|---|---|---|
| very little data | Use Linear Classifier on top layer | ? |
| quite a lot of data | Finetune a few layers | ? |

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
Conv-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image

More specific

More generic

|  | **very similar dataset** | **very different dataset** |
|---|---|---|
| **very little data** | Use Linear Classifier on top layer | You're in trouble… Try linear classifier from different stages |
| **quite a lot of data** | Finetune a few layers | Finetune a larger number of layers |

# Transfer learning with CNNs is pervasive…
## (it's the norm, not an exception)

Object Detection
(Fast R-CNN)



Log loss + smooth L1 loss

Proposal classifier — Linear + softmax

Linear — Bounding box regressors

FCs

RoI pooling

External proposal algorithm e.g. selective search

ConvNet (applied to entire image)

Image Captioning: CNN + RNN



"straw"   "hat"   END

$y_t$

$W_{oh}$

$W_{hh}$   $h_t$

$CNN_{\theta_c}$   $W_{hi}$

$W_{hx}$

$x_t$

START   "straw"   "hat"

Girshick, "Fast R-CNN", ICCV 2015
Figure copyright Ross Girshick, 2015. Reproduced with permission.

Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015
Figure copyright IEEE, 2015. Reproduced for educational purposes.

# Transfer learning with CNNs is pervasive...
## (it's the norm, not an exception)

Object Detection
(Fast R-CNN)

CNN pretrained
on ImageNet

Image Captioning: CNN + RNN



Log loss + smooth L1 loss

Proposal
classifier

Linear +
softmax

Linear

Bounding box
regressors

FCs

RoI pooling

External proposal
algorithm
e.g. selective search

ConvNet
(applied to entire
image)



"straw"  "hat"  END

$y_t$

$W_{hh}$

$W_{oh}$

$h_t$

$CNN_{\theta_c}$

$V_{hi}$

$V_{hx}$

$x_t$

START  "straw"  "hat"

Girshick, "Fast R-CNN", ICCV 2015
Figure copyright Ross Girshick, 2015. Reproduced with permission.

Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for
Generating Image Descriptions", CVPR 2015
Figure copyright IEEE, 2015. Reproduced for educational purposes.

# Transfer learning with CNNs is pervasive…
## (it's the norm, not an exception)

Object Detection
(Fast R-CNN)

CNN pretrained
on ImageNet

Image Captioning: CNN + RNN



Word embedding layer
pretrained with word2vec

# Transfer learning with CNNs is pervasive...
## (it's the norm, not an exception)



Unified Vision-Language Pre-training

Zhou et al, "Unified Vision-Language Pre-Training for Image Captioning and VQA" CVPR 2020
Figure copyright Luowei Zhou, 2020. Reproduced with permission.

1. Train CNN on ImageNet
2. Fine-Tune (1) for object detection on Visual Genome
1. Train BERT language model on lots of text
2. Combine(2) and (3), train for joint image / language modeling
3. Fine-tune (4) for image captioning, visual question answering, etc.

Krishna et al, "Visual genome: Connecting language and vision using crowdsourced dense image annotations" IJCV 2017
Devlin et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding" ArXiv 2018

# Transfer learning is pervasive…
## But recent results show it might not always be necessary!

bbox AP: R50-FPN, GN



typical
fine-tuning
schedule

random init
w/ pre-train

He et al, "Rethinking ImageNet Pre-training", ICCV 2019
Figure copyright Kaiming He, 2019. Reproduced with permission.

Training from scratch can work just as well as training from a pretrained ImageNet model for object detection

But it takes 2-3x as long to train.

They also find that collecting more data is better than finetuning on a related task

# Takeaway for your projects and beyond:



Transfer learning be like

Custom layers

Pretrained layers

Source: AI & Deep Learning Memes For Back-propagated Poets

# Takeaway for your projects and beyond:

Have some dataset of interest but it has < ~1M images?

1. Find a very large dataset that has similar data, train a big neural network there
2. Transfer learn to your dataset

Deep learning frameworks provide a "Model Zoo" of pretrained models so you don't need to train your own

TensorFlow: https://github.com/tensorflow/models
PyTorch: https://github.com/pytorch/vision

**Image Classification**: A core task in Computer Vision

(assume given a set of possible labels)
{dog, cat, truck, plane, ...}

⟶   cat

# Structured prediction tasks in vision

**Classification**



**CAT**

No spatial extent

**Semantic Segmentation**



**GRASS, CAT, TREE, SKY**

No objects, just pixels

**Object Detection**



**DOG, DOG, CAT**

**Instance Segmentation**



**DOG, DOG, CAT**

Multiple Object

# Semantic Segmentation



Classification

**Semantic Segmentation**

Object Detection

Instance Segmentation

CAT

**GRASS**, **CAT**,
**TREE**, **SKY**

DOG, DOG, CAT

DOG, DOG, CAT

No spatial extent

No objects, just pixels

Multiple Object

# Semantic Segmentation: The Problem



**GRASS**, **CAT**,
**TREE**, **SKY**, ...

Paired training data: for each training image,
each pixel is labeled with a semantic category.



At test time, classify each pixel of a new image.

# Semantic Segmentation Idea: Sliding Window

Full image



?

# Semantic Segmentation Idea: Sliding Window

Full image



?

Impossible to classify without context

Q: how do we include context?

# Semantic Segmentation Idea: Sliding Window

Full image



Q: how do we model this?

# Semantic Segmentation Idea: Sliding Window

Extract patch

Classify center pixel with CNN

Full image



Cow

Cow

Grass

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

# Semantic Segmentation Idea: Sliding Window

Extract patch

Classify center pixel with CNN

Full image



Cow

Cow

Grass

Problem: Very inefficient! Not reusing shared features between overlapping patches

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

# Semantic Segmentation Idea: Convolution

Full image



An intuitive idea: encode the entire image with conv net, and do semantic segmentation on top.

Problem: classification architectures often reduce feature spatial sizes to go deeper, but semantic segmentation requires the output size to be the same as input size.

# Semantic Segmentation Idea: Fully Convolutional

Design a network with only convolutional layers without downsampling operators to make predictions for pixels all at once!



Input:
3 x H x W

Conv    Conv    Conv

Convolutions:
D x H x W

Conv    argmax

Scores:
C x H x W

Predictions:
H x W

# Semantic Segmentation Idea: Fully Convolutional

Design a network with only convolutional layers
without downsampling operators to make predictions
for pixels all at once!



Input:
3 x H x W

Conv    Conv    Conv    Conv    argmax

Convolutions:
D x H x W

Scores:
C x H x W

Predictions:
H x W

Problem: convolutions at
original image resolution will
be very expensive ...

# Semantic Segmentation Idea: Fully Convolutional

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



Input:
3 x H x W

High-res:
$D_1$ x H/2 x W/2

Med-res:
$D_2$ x H/4 x W/4

Low-res:
$D_3$ x H/4 x W/4

Med-res:
$D_2$ x H/4 x W/4

High-res:
$D_1$ x H/2 x W/2

C x H x W

Predictions:
H x W

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

# Semantic Segmentation Idea: Fully Convolutional

**Downsampling**:
Pooling, strided convolution

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!

**Upsampling**:
???



Input:
$3 \times H \times W$

High-res:
$D_1 \times H/2 \times W/2$

Med-res:
$D_2 \times H/4 \times W/4$

Low-res:
$D_3 \times H/4 \times W/4$

Med-res:
$D_2 \times H/4 \times W/4$

High-res:
$D_1 \times H/2 \times W/2$

$C \times H \times W$

Predictions:
$H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

# In-Network upsampling: "Unpooling"

**Nearest Neighbor**

| 1 | 2 |
|---|---|
| 3 | 4 |

→

| 1 | 1 | 2 | 2 |
|---|---|---|---|
| 1 | 1 | 2 | 2 |
| 3 | 3 | 4 | 4 |
| 3 | 3 | 4 | 4 |

Input: 2 x 2          Output: 4 x 4

**"Bed of Nails"**

| 1 | 2 |
|---|---|
| 3 | 4 |

→

| 1 | 0 | 2 | 0 |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 3 | 0 | 4 | 0 |
| 0 | 0 | 0 | 0 |

Input: 2 x 2          Output: 4 x 4

# In-Network upsampling: "Max Unpooling"

**Max Pooling**
Remember which element was max!

| 1 | 2 | 6 | 3 |
|---|---|---|---|
| 3 | 5 | 2 | 1 |
| 1 | 2 | 2 | 1 |
| 7 | 3 | 4 | 8 |

| 5 | 6 |
|---|---|
| 7 | 8 |

. . . .

Rest of the network

**Max Unpooling**
Use positions from pooling layer

| 1 | 2 |
|---|---|
| 3 | 4 |

| 0 | 0 | 2 | 0 |
|---|---|---|---|
| 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 4 |

Input: 4 x 4          Output: 2 x 2          Input: 2 x 2          Output: 4 x 4

Corresponding pairs of downsampling and upsampling layers

# Learnable Upsampling

**Recall:** Normal 3 x 3 convolution, stride 1 pad 1

Input: 4 x 4

Output: 4 x 4

# Learnable Upsampling

**Recall:** Normal 3 x 3 convolution, stride 1 pad 1

Dot product
between filter
and input

Input: 4 x 4

Output: 4 x 4

# Learnable Upsampling

**Recall:** Normal 3 x 3 convolution, stride 1 pad 1



Dot product between filter and input

Input: 4 x 4

Output: 4 x 4

# Learnable Upsampling

**Recall:** Normal 3 x 3 convolution, <u>stride 2</u> pad 1

Input: 4 x 4

Output: 2 x 2

# Learnable Upsampling

**Recall:** Normal 3 x 3 convolution, <u>stride 2</u> pad 1



Dot product between filter and input

Input: 4 x 4

Output: 2 x 2

# Learnable Upsampling

**Recall:** Normal 3 x 3 convolution, <u>stride 2</u> pad 1



Dot product between filter and input

Input: 4 x 4

Output: 2 x 2

Filter moves 2 pixels in the input for every one pixel in the output

Stride gives ratio between movement in input and output

We can interpret strided convolution as "learnable downsampling".

# Learnable Upsampling: Transposed Convolution

3 x 3 **transpose** convolution, stride 2 pad 1

Input: 2 x 2

Output: 4 x 4

# Learnable Upsampling: Transposed Convolution

3 x 3 **transpose** convolution, stride 2 pad 1

Input gives weight for filter

Input: 2 x 2

Output: 4 x 4

# Learnable Upsampling: Transposed Convolution

3 x 3 **transpose** convolution, stride 2 pad 1



Input gives weight for filter

Input: 2 x 2

Output: 4 x 4

Filter moves 2 pixels in the <u>output</u> for every one pixel in the <u>input</u>

Stride gives ratio between movement in output and input

# Learnable Upsampling: Transposed Convolution

3 x 3 **transpose** convolution, stride 2 pad 1

Sum where output overlaps

Input gives weight for filter

Filter moves 2 pixels in the <u>output</u> for every one pixel in the <u>input</u>

Stride gives ratio between movement in output and input

Input: 2 x 2

Output: 4 x 4

# Learnable Upsampling: Transposed Convolution

Q: Why is it called transpose convolution?

3 x 3 **transpose** convolution, stride 2 pad 1

Sum where output overlaps

Input gives weight for filter

Filter moves 2 pixels in the <u>output</u> for every one pixel in the <u>input</u>

Stride gives ratio between movement in output and input

Input: 2 x 2

Output: 4 x 4

# Learnable Upsampling: 1D Example

**Input**

**Filter**

**Output**

| |
|---|
| a |
| b |

| |
|---|
| x |
| y |
| z |

| |
|---|
| ax |
| ay |
| az + bx |
| by |
| bz |

Output contains copies of the filter weighted by the input, summing at where at overlaps in the output

# Convolution as Matrix Multiplication (1D Example)

We can express convolution in
terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

Example: 1D conv, kernel
size=3, stride=2, padding=1

# Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

Example: 1D conv, kernel size=3, <u>stride=2</u>, padding=1

Convolution transpose multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{bmatrix} x & 0 \\ y & 0 \\ z & x \\ 0 & y \\ 0 & z \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} ax \\ ay \\ az + bx \\ by \\ bz \\ 0 \end{bmatrix}$$

Example: 1D transpose conv, kernel size=3, <u>stride=2</u>, padding=0

# Semantic Segmentation Idea: Fully Convolutional

**Downsampling**: Pooling, strided convolution

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!

**Upsampling**: Unpooling or strided transpose convolution



Input: 3 x H x W

High-res: $D_1$ x H/2 x W/2

Med-res: $D_2$ x H/4 x W/4

Low-res: $D_3$ x H/4 x W/4

Med-res: $D_2$ x H/4 x W/4

High-res: $D_1$ x H/2 x W/2

Predictions: H x W

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

# Today: UNet with residual connections



Residual connections

Input:
3 x H x W

Predictions:
H x W

Newell et al. Stacked Hourglass Networks for Human Pose Estimation. ECCV 2016

# Semantic Segmentation: Summary

# Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels

# Object Detection

| Classification | Semantic Segmentation | Object Detection | **Instance Segmentation** |
|---|---|---|---|
| CAT | GRASS, CAT, TREE, SKY | DOG, DOG, CAT | **DOG**, **DOG**, **CAT** |
| No spatial extent | No objects, just pixels | Multiple Object | |

# Object Detection

**Classification**

CAT

*No spatial extent*

**Semantic Segmentation**

GRASS, CAT, TREE, SKY

*No objects, just pixels*

**Object Detection**

DOG, DOG, CAT

**Instance Segmentation**

DOG, DOG, CAT

Multiple Object

# Object Detection: Single Object
(Classification + Localization)

x, y

h

w

**Fully Connected**: 4096 to 1000

**Class Scores**
Cat: 0.9
Dog: 0.05
Car: 0.01
...

**Vector:** 4096

**Fully Connected**: 4096 to 4

**Box Coordinates**
(x, y, w, h)

# Object Detection: Single Object
(Classification + Localization)

**Correct label:**
Cat

**Fully Connected:**
4096 to 1000

**Class Scores**
Cat: 0.9
Dog: 0.05
Car: 0.01
...

**Softmax Loss**

x, y

h

w

This image is CC0 public domain

**Vector:**
4096

**Fully Connected:**
4096 to 4

**Box Coordinates**
(x, y, w, h)

**L2 Loss**

**Correct box:**
(x', y', w', h')

Treat localization as a regression problem!

# Object Detection: Single Object
(Classification + Localization)



This image is CC0 public domain

**Fully Connected:** 4096 to 1000

**Class Scores**
Cat: 0.9
Dog: 0.05
Car: 0.01
...

**Correct label:** Cat

**Softmax Loss**

Multitask Loss

**+** → **Loss**

**Vector:** 4096

**Fully Connected:** 4096 to 4

**Box Coordinates** (x, y, w, h)

→ **L2 Loss**

**Correct box:** (x', y', w', h')

Treat localization as a regression problem!

# Object Detection: Multiple Objects



CAT: (x, y, w, h)

DOG: (x, y, w, h)
DOG: (x, y, w, h)
CAT: (x, y, w, h)

DUCK: (x, y, w, h)
DUCK: (x, y, w, h)
….

# Object Detection: Multiple Objects

Each image needs a different number of outputs!

CAT: (x, y, w, h)    4 numbers

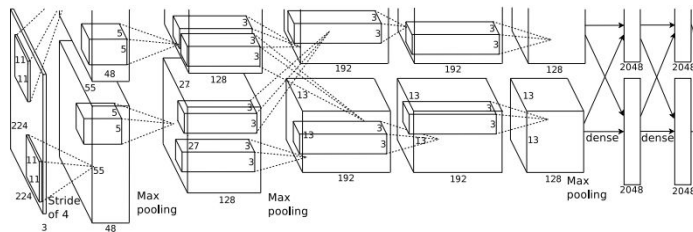DOG: (x, y, w, h)
DOG: (x, y, w, h)    12 numbers
CAT: (x, y, w, h)

DUCK: (x, y, w, h)   Many
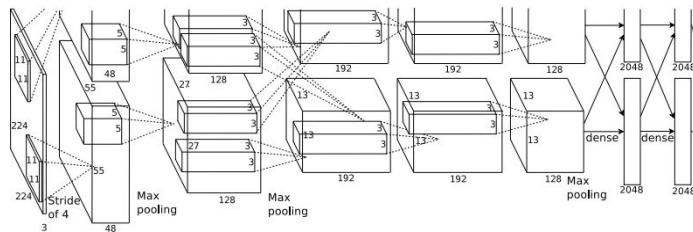DUCK: (x, y, w, h)   numbers!
....

# Object Detection: Multiple Objects

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
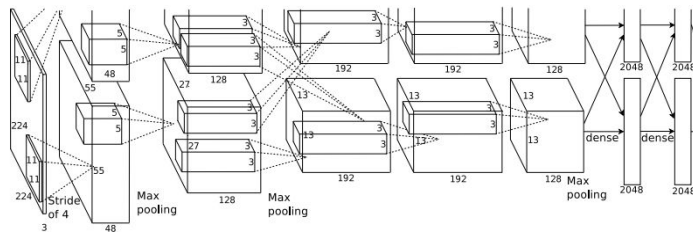Cat? NO
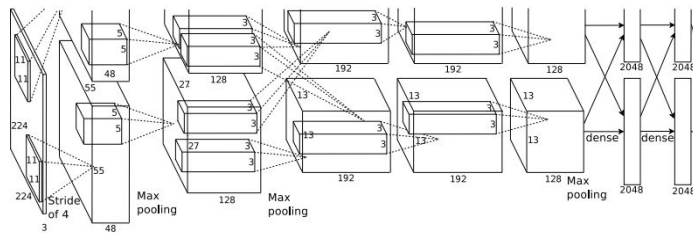Background? YES

# Object Detection: Multiple Objects

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES
Cat? NO
Background? NO

# Object Detection: Multiple Objects

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES
Cat? NO
Background? NO

# Object Detection: Multiple Objects

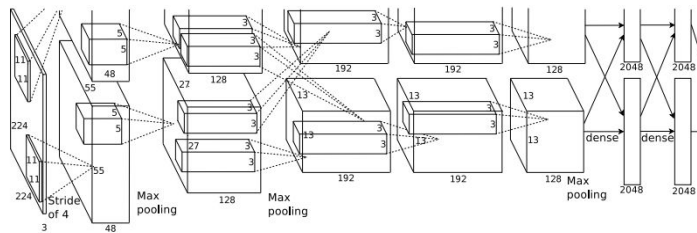Apply a CNN to many different crops of the image, CNN classifies each crop as object or background
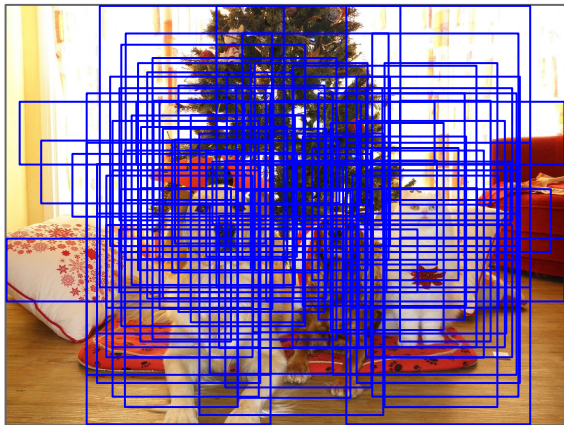


Dog? NO
Cat? YES
Background? NO

Q: What's the problem with this approach?

# Object Detection: Multiple Objects

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background
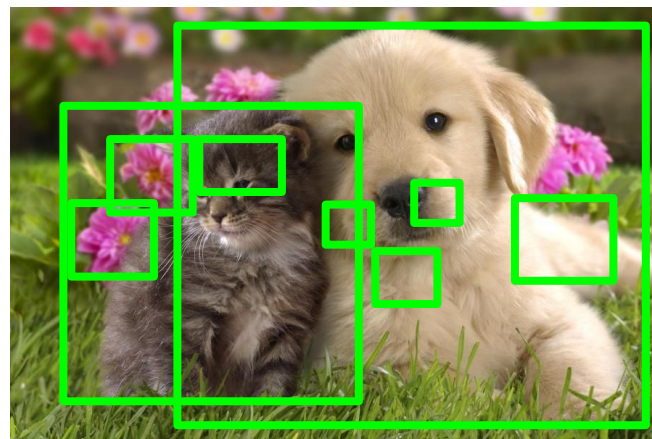


Dog? NO
Cat? YES
Background? NO

Problem: Need to apply CNN to huge number of locations, scales, and aspect ratios, very computationally expensive!

# Region Proposals: Selective Search

- Find "blobby" image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 2000 region proposals in a few seconds on CPU



Alexe et al, "Measuring the objectness of image windows", TPAMI 2012
Uijlings et al, "Selective Search for Object Recognition", IJCV 2013
Cheng et al, "BING: Binarized normed gradients for objectness estimation at 300fps", CVPR 2014
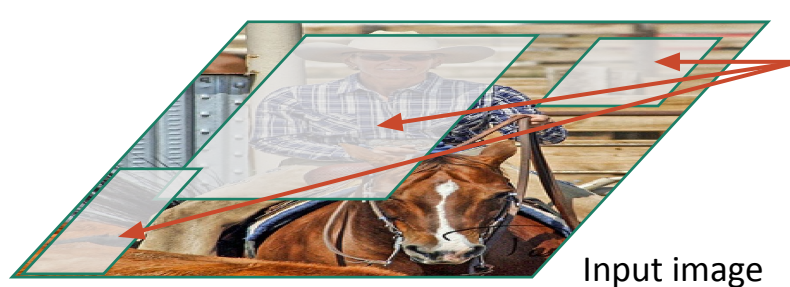Zitnick and Dollar, "Edge boxes: Locating object proposals from edges", ECCV 2014

# R-CNN



Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
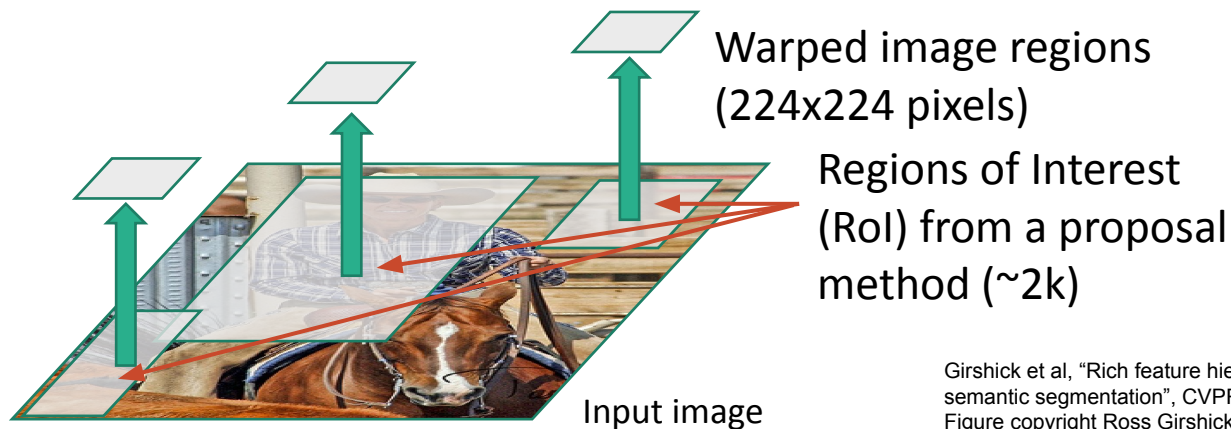Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN



Regions of Interest (RoI) from a proposal method (~2k)

Input image

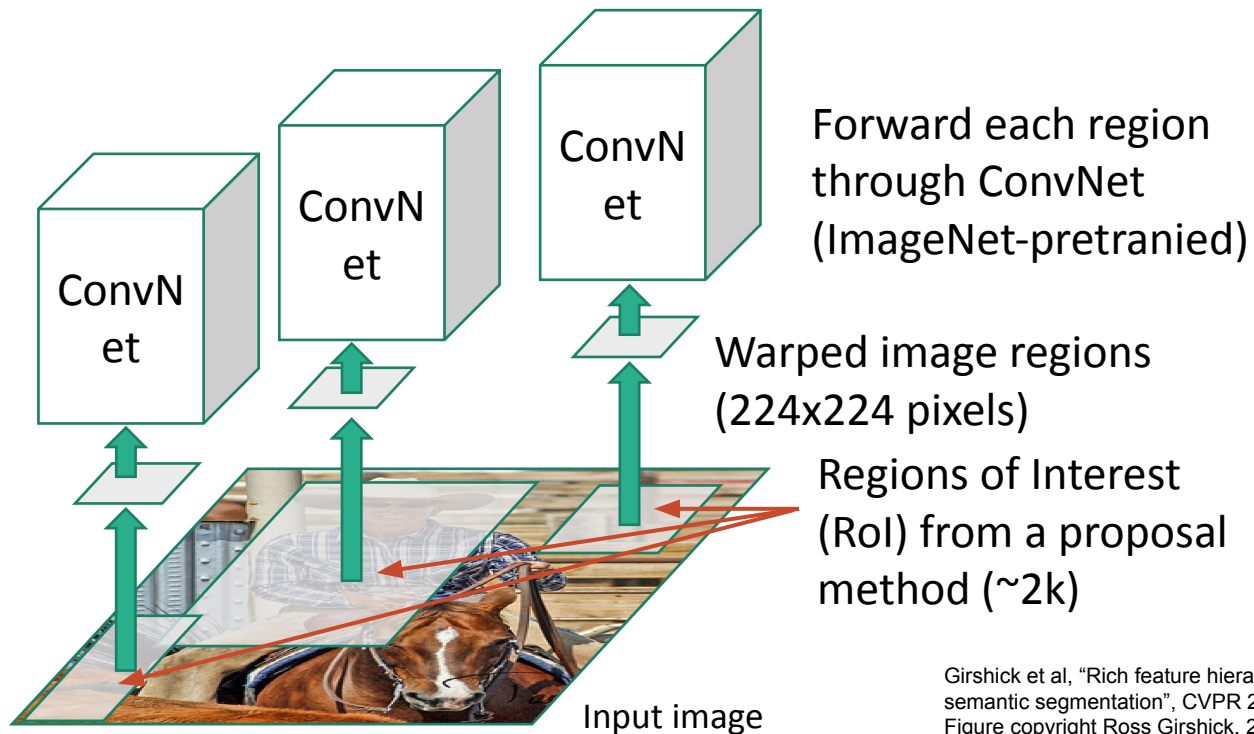Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN



Warped image regions
(224x224 pixels)

Regions of Interest
(RoI) from a proposal
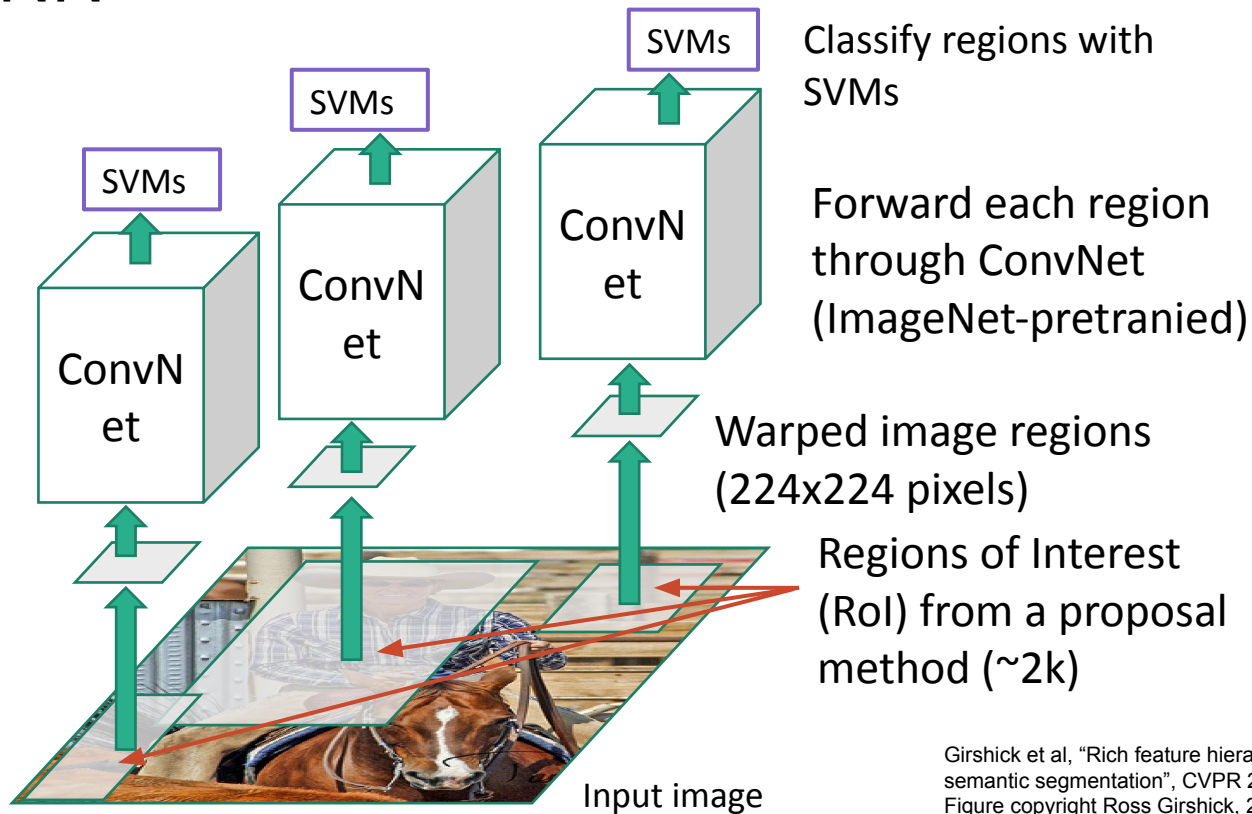method (~2k)

Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and
semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN



Forward each region through ConvNet (ImageNet-pretranied)

Warped image regions (224x224 pixels)

Regions of Interest (RoI) from a proposal method (~2k)

Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
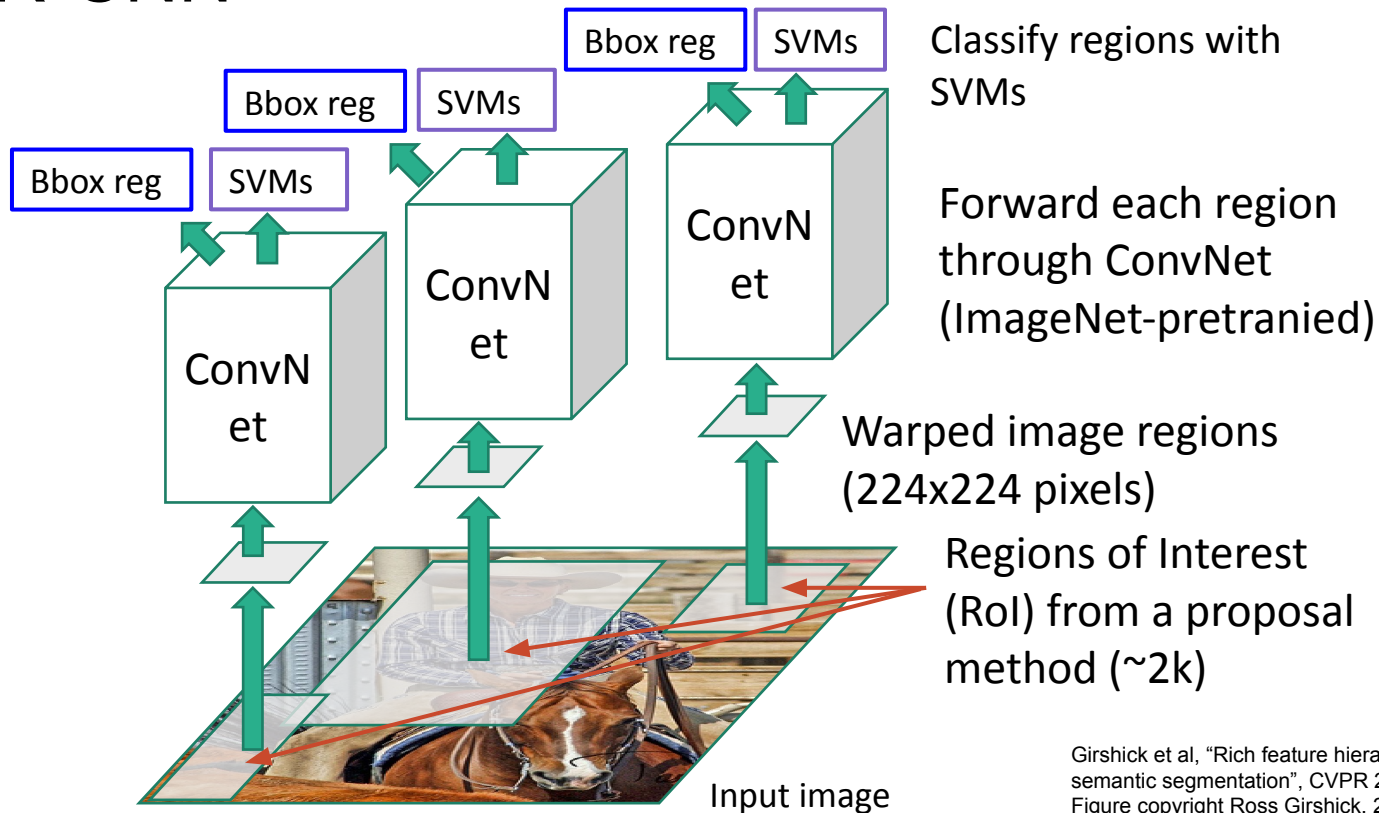Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN



Classify regions with SVMs

Forward each region through ConvNet (ImageNet-pretranied)

Warped image regions (224x224 pixels)

Regions of Interest (RoI) from a proposal method (~2k)
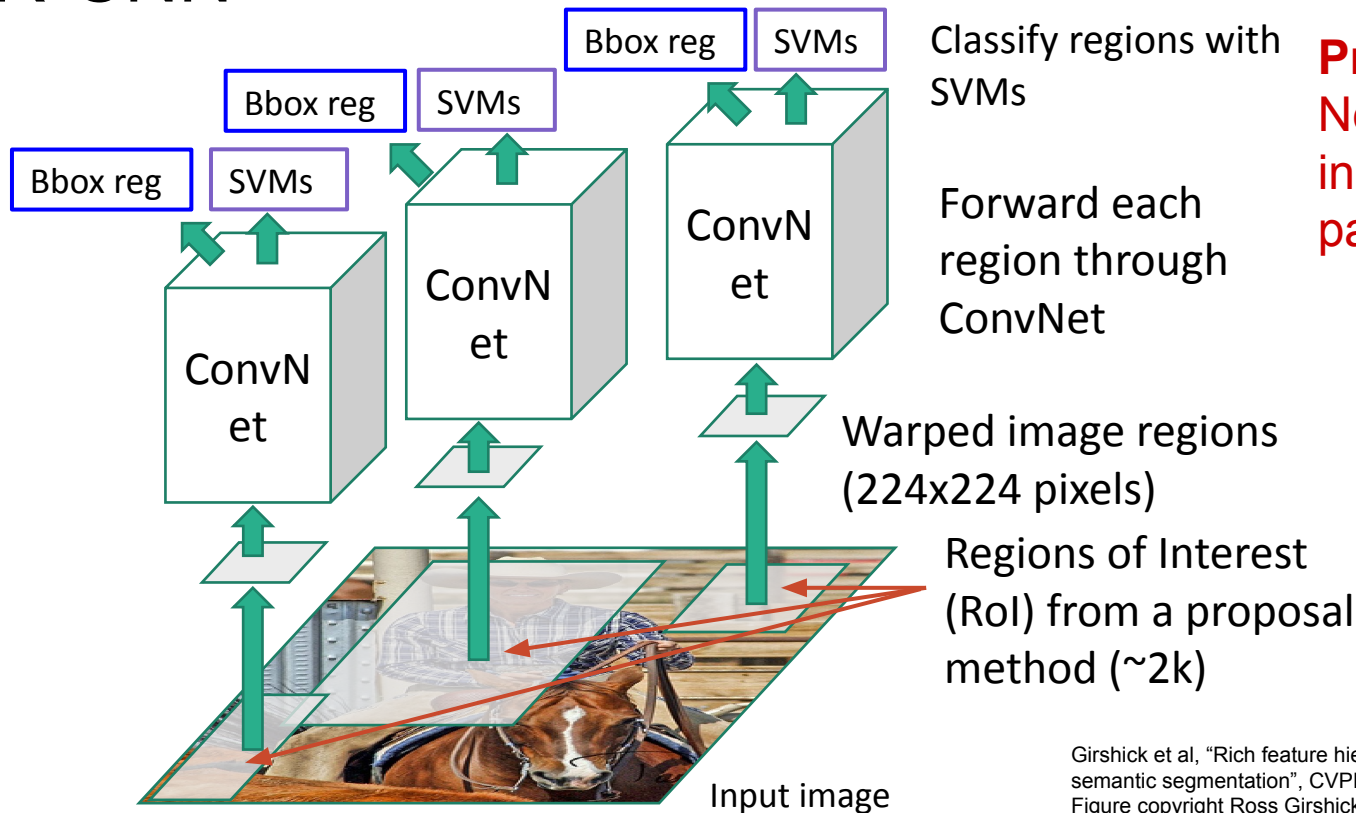
Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN

Predict "corrections" to the RoI: 4 numbers: (dx, dy, dw, dh)



Classify regions with SVMs

Forward each region through ConvNet (ImageNet-pretranied)

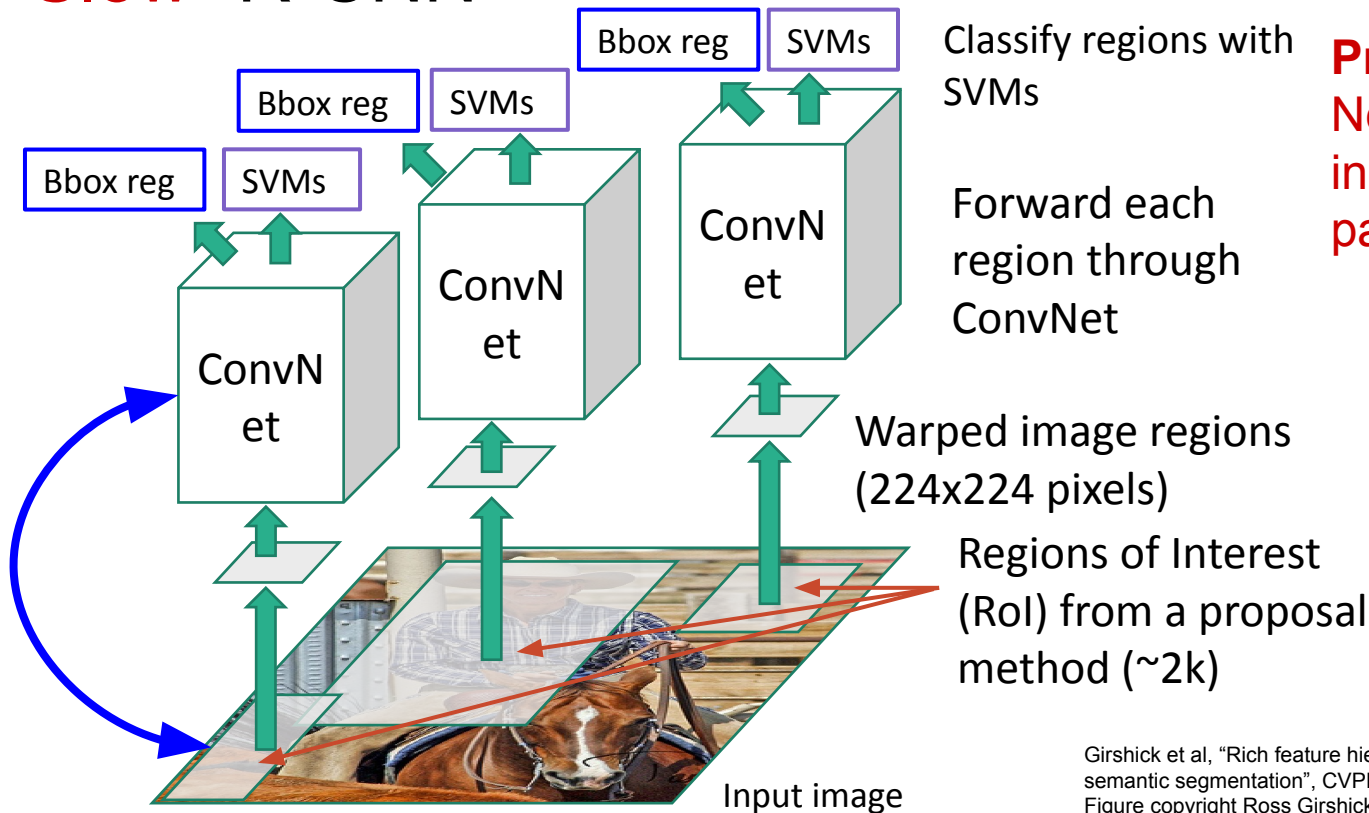Warped image regions (224x224 pixels)

Regions of Interest (RoI) from a proposal method (~2k)

Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN

Predict "corrections" to the RoI: 4 numbers: (dx, dy, dw, dh)

Bbox reg  SVMs

Classify regions with SVMs

**Problem**: Very slow! Need to do ~2k independent forward passes for each image!

Bbox reg  SVMs

ConvNet

Bbox reg  SVMs

ConvNet

Forward each region through ConvNet

ConvNet

Warped image regions (224x224 pixels)

Regions of Interest (RoI) from a proposal method (~2k)

Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# "Slow" R-CNN

Predict "corrections" to the RoI: 4 numbers: (dx, dy, dw, dh)

**Problem**: Very slow! Need to do ~2k independent forward passes for each image!

**Idea:** Pass the image through convnet before cropping! Crop the conv feature instead!

Classify regions with SVMs

Forward each region through ConvNet

Warped image regions (224x224 pixels)

Regions of Interest (RoI) from a proposal method (~2k)

Bbox reg    SVMs
Bbox reg    SVMs
Bbox reg    SVMs

ConvNet
ConvNet
ConvNet

Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
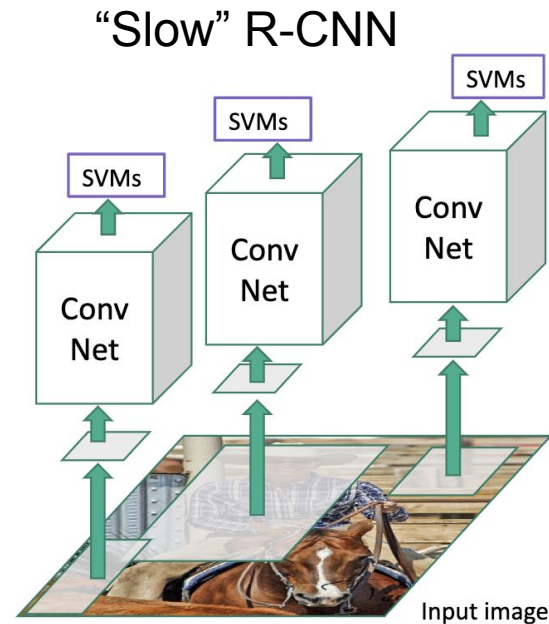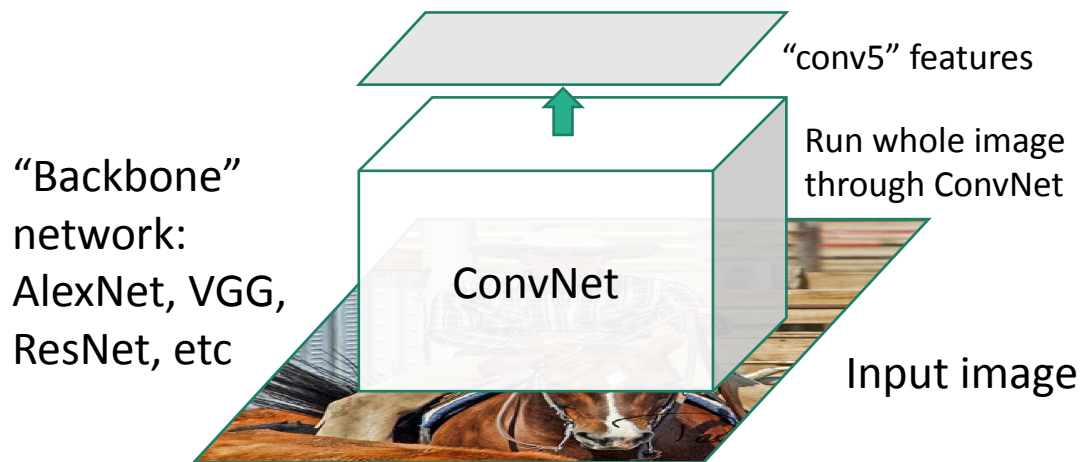Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# Fast R-CNN



"Slow" R-CNN

Input image

Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# Fast R-CNN



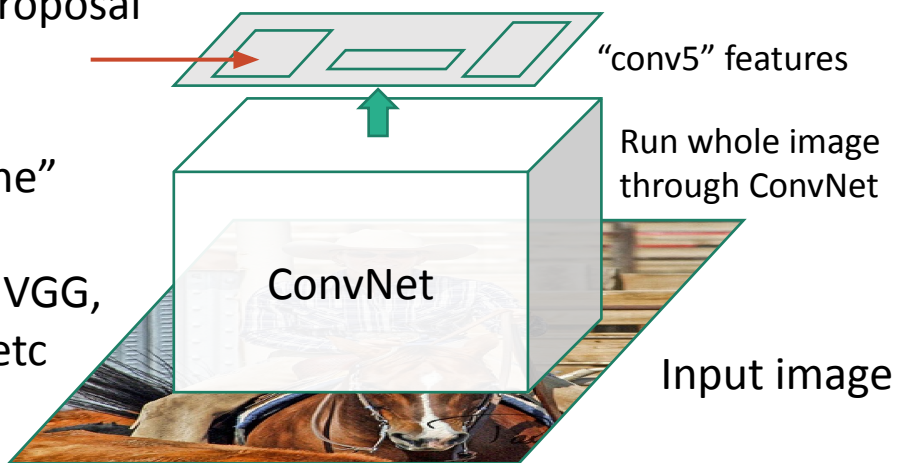"Slow" R-CNN

"conv5" features

Run whole image through ConvNet

"Backbone" network: AlexNet, VGG, ResNet, etc

ConvNet

Input image

# Fast R-CNN



Regions of Interest (RoIs) from a proposal method

"Backbone" network: AlexNet, VGG, ResNet, etc

ConvNet

"conv5" features

Run whole image through ConvNet

Input image

"Slow" R-CNN
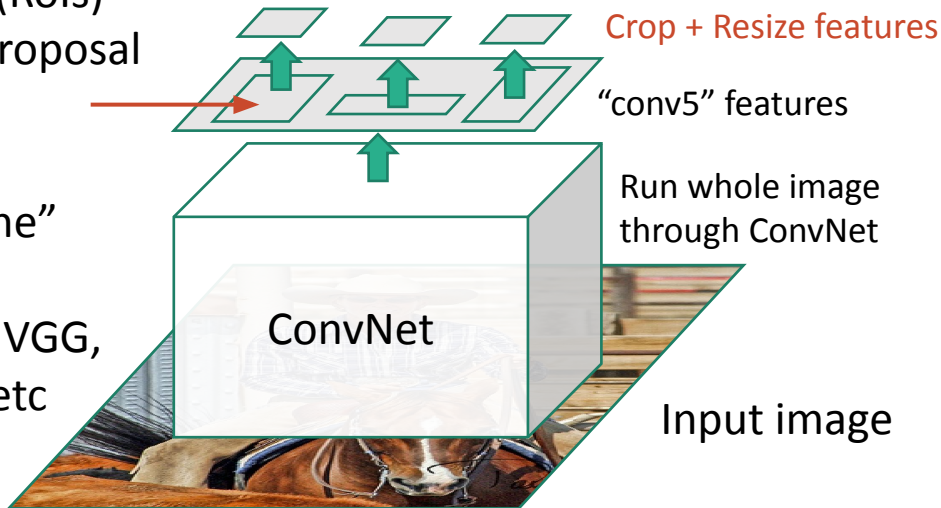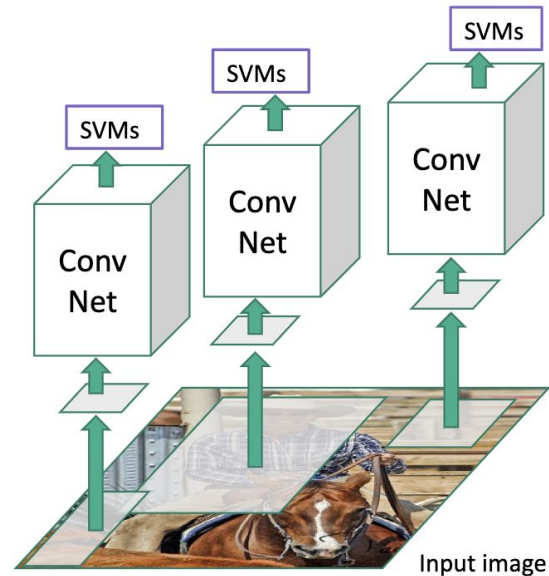
SVMs
SVMs
SVMs

Conv Net
Conv Net
Conv Net

Input image

Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# Fast R-CNN

Regions of
Interest (RoIs)
from a proposal
method

Crop + Resize features

"conv5" features

Run whole image
through ConvNet

ConvNet

"Backbone"
network:
AlexNet, VGG,
ResNet, etc

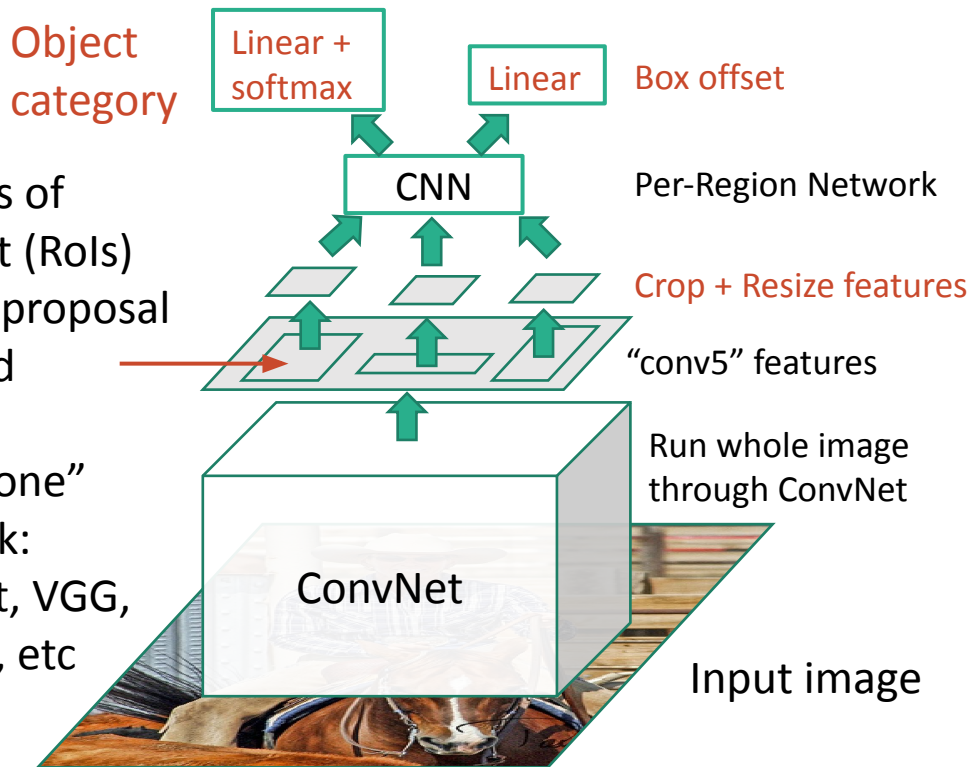Input image

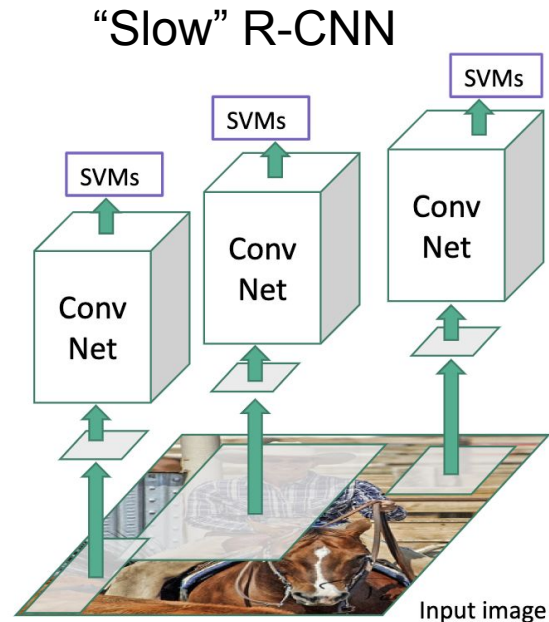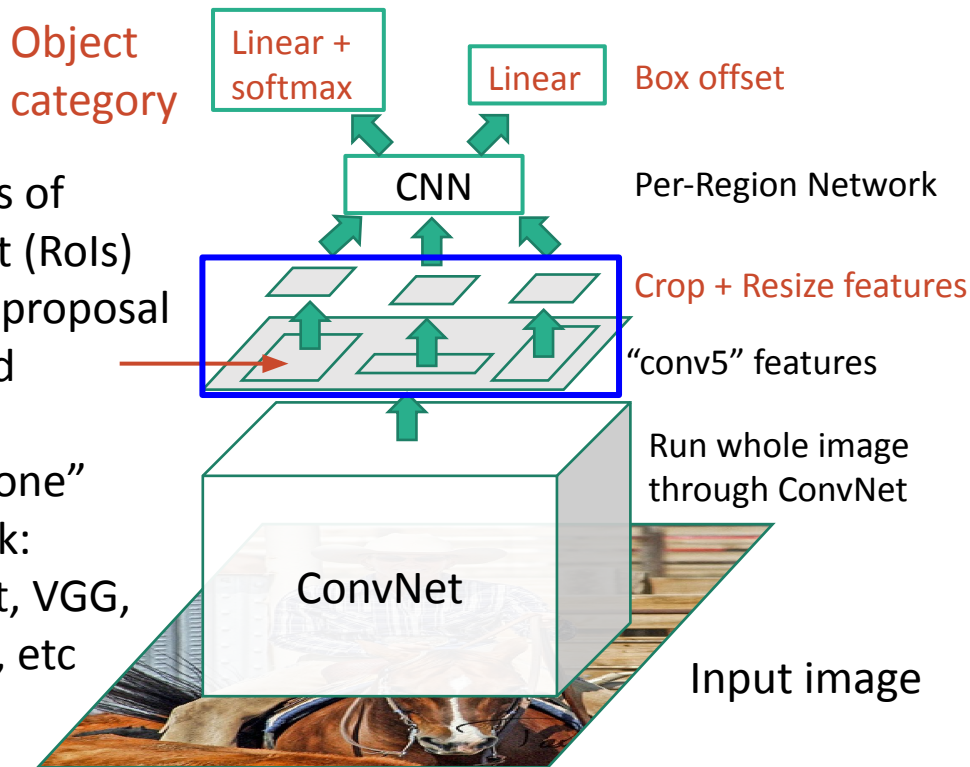"Slow" R-CNN

SVMs

SVMs

SVMs

Conv
Net

Conv
Net

Conv
Net

Input image
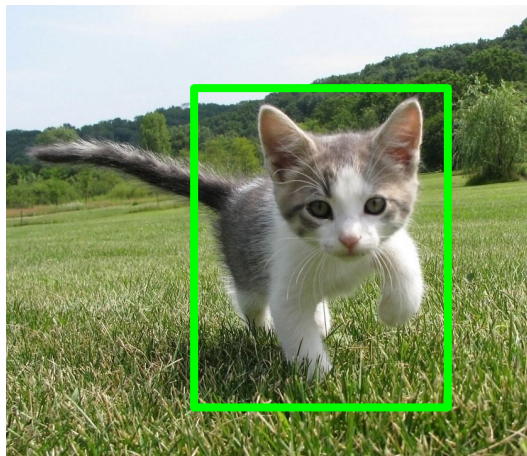
Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# Fast R-CNN

**Object category**

Linear + softmax

Linear — Box offset

CNN — Per-Region Network

**Regions of Interest (RoIs) from a proposal method**

Crop + Resize features

"conv5" features

Run whole image through ConvNet

**"Backbone" network: AlexNet, VGG, ResNet, etc**

ConvNet

Input image

### "Slow" R-CNN

SVMs

SVMs

SVMs

ConvNet

ConvNet

ConvNet

Input image

# Fast R-CNN

Object category

Linear + softmax

Linear    Box offset

CNN    Per-Region Network

Regions of Interest (RoIs) from a proposal method

Crop + Resize features

"conv5" features

Run whole image through ConvNet

"Backbone" network: AlexNet, VGG, ResNet, etc

ConvNet

Input image

"Slow" R-CNN

SVMs

SVMs

SVMs

Conv Net

Conv Net

Conv Net

Input image

Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

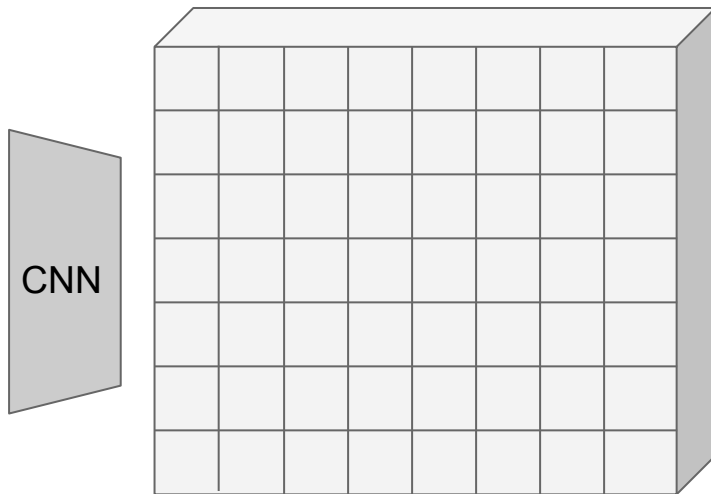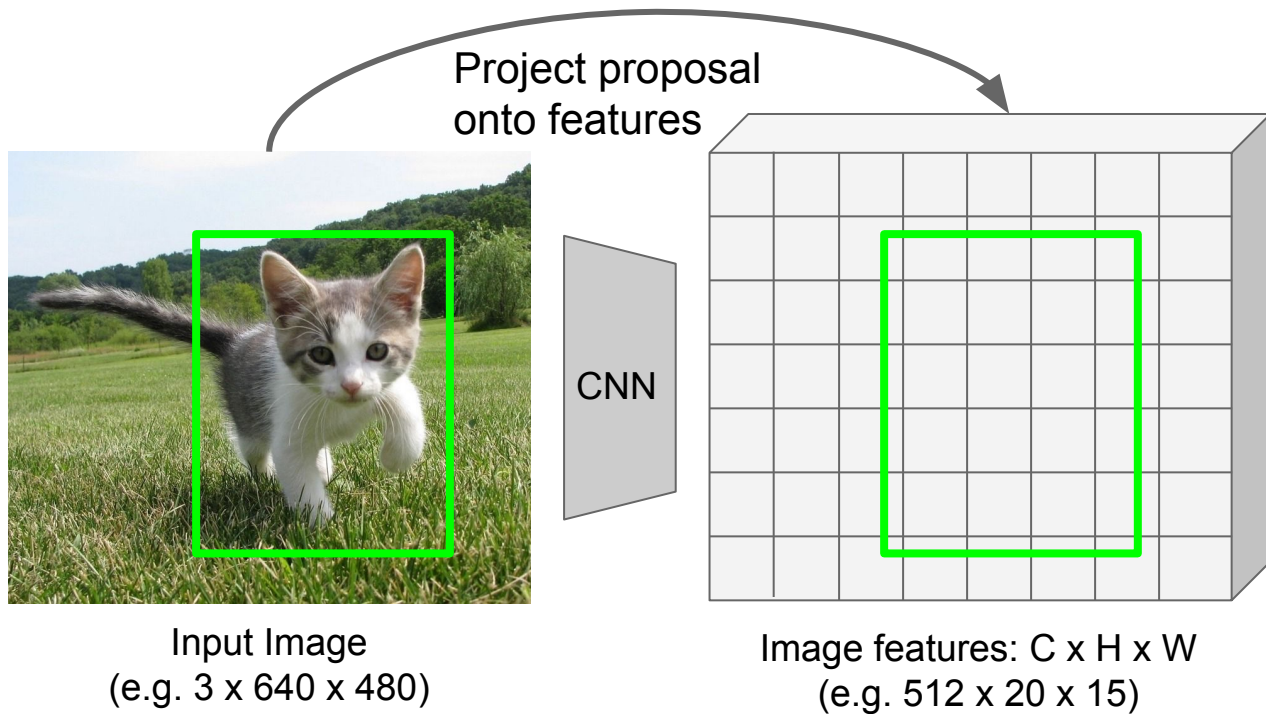# Cropping Features: RoI Pool



Input Image
(e.g. 3 x 640 x 480)

Image features: C x H x W
(e.g. 512 x 20 x 15)
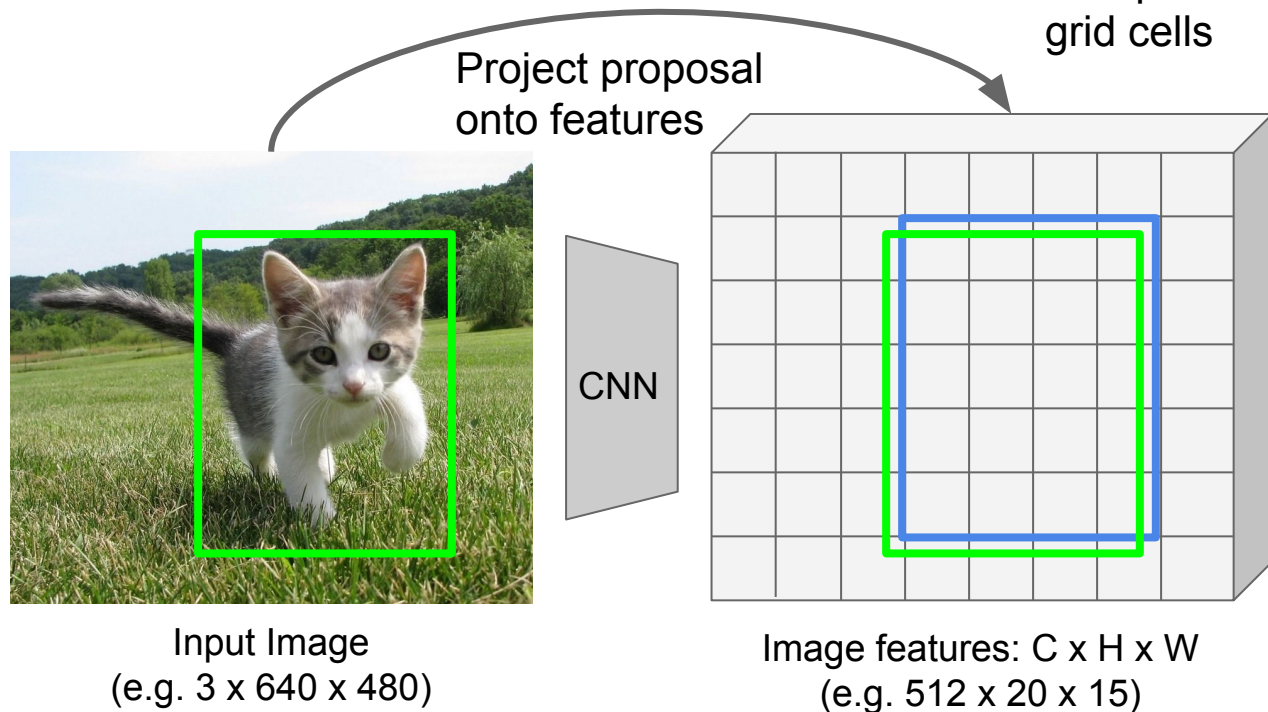
CNN

Girshick, "Fast R-CNN", ICCV 2015.

Girshick, "Fast R-CNN", ICCV 2015.

# Cropping Features: RoI Pool

Project proposal onto features



Input Image
(e.g. 3 x 640 x 480)

CNN

Image features: C x H x W
(e.g. 512 x 20 x 15)

Girshick, "Fast R-CNN", ICCV 2015.

Girshick, "Fast R-CNN", ICCV 2015.

# Cropping Features: RoI Pool



"Snap" to grid cells

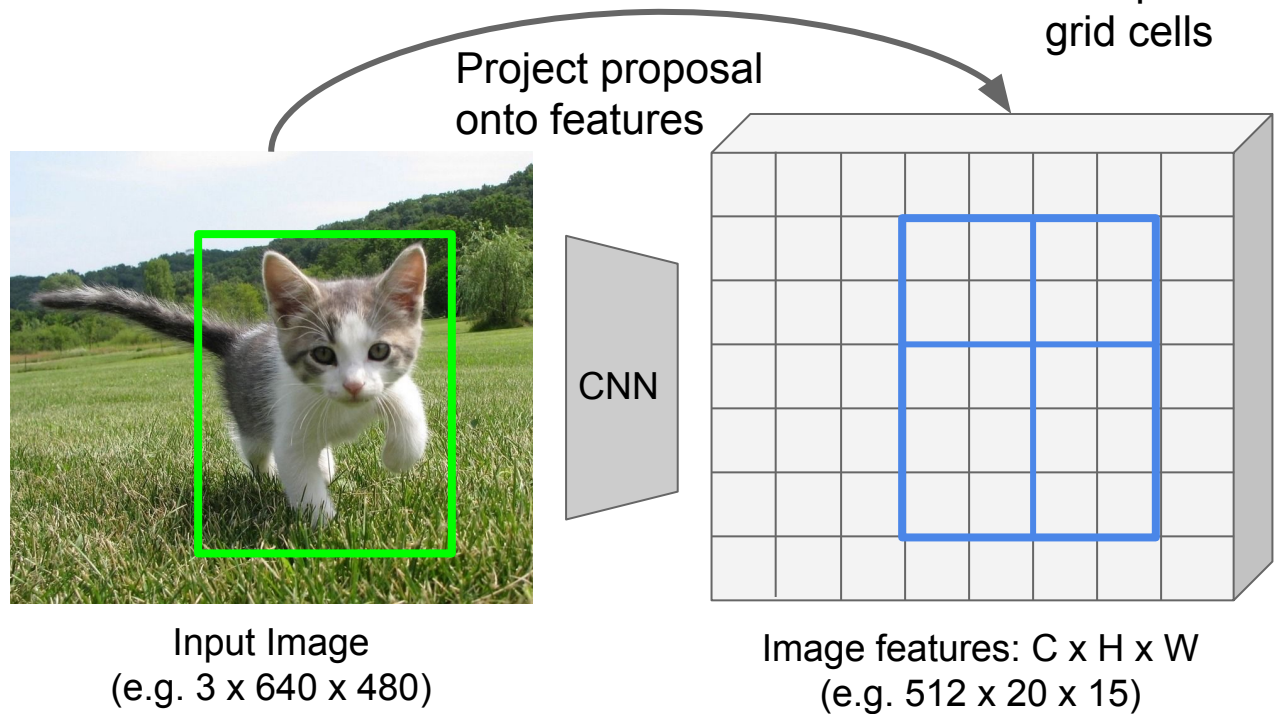Project proposal onto features

CNN

Input Image
(e.g. 3 x 640 x 480)

Image features: C x H x W
(e.g. 512 x 20 x 15)

Girshick, "Fast R-CNN", ICCV 2015.

# Cropping Features: RoI Pool

"Snap" to grid cells

Project proposal onto features



Input Image
(e.g. 3 x 640 x 480)

CNN

Image features: C x H x W
(e.g. 512 x 20 x 15)

Q: how do we resize the 512 x 20 x 15 region to, e.g., a 512 x 2 x 2 tensor?.

Girshick, "Fast R-CNN", ICCV 2015.

# Cropping Features: RoI Pool

"Snap" to grid cells

Divide into 2x2 grid of (roughly) equal subregions

Project proposal onto features



Input Image
(e.g. 3 x 640 x 480)

CNN

Image features: C x H x W
(e.g. 512 x 20 x 15)

Q: how do we resize the 512 x 20 x 15 region to, e.g., a 512 x 2 x 2 tensor?.

Girshick, "Fast R-CNN", ICCV 2015.
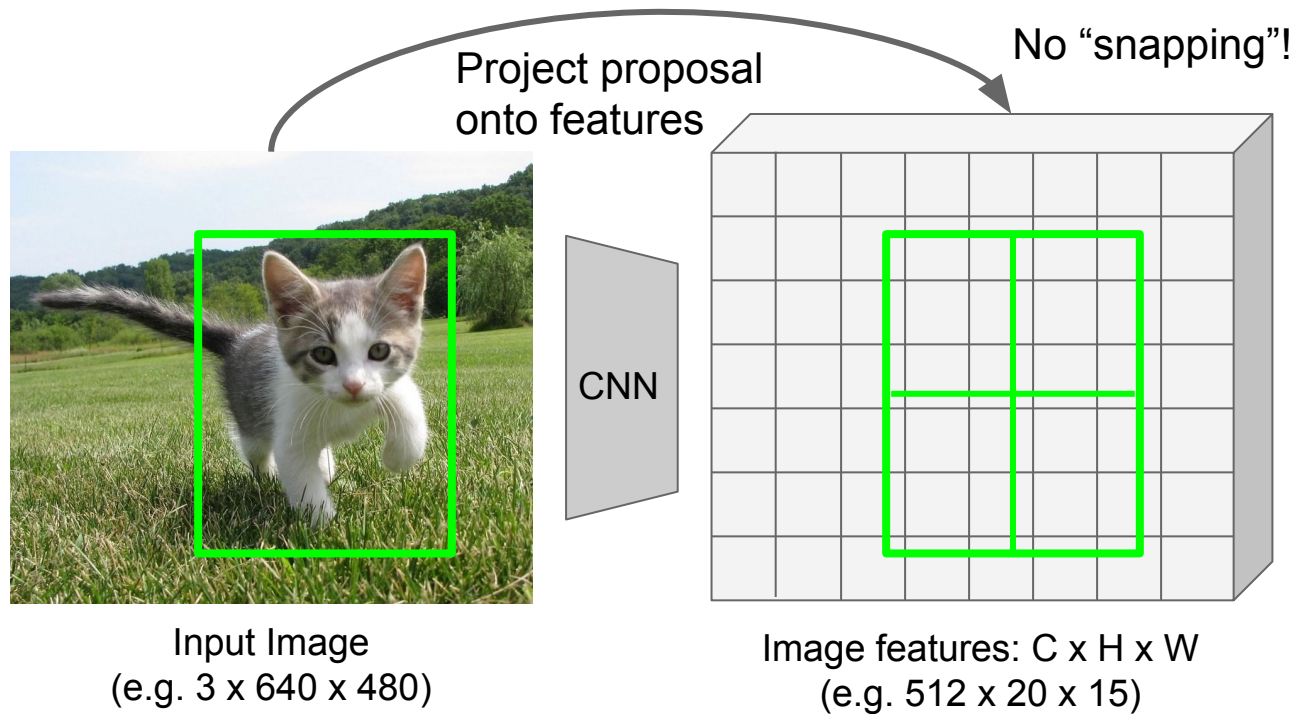
# Cropping Features: RoI Pool

"Snap" to grid cells

Divide into 2x2 grid of (roughly) equal subregions

Project proposal onto features

Max-pool within each subregion

CNN

Input Image
(e.g. 3 x 640 x 480)

Image features: C x H x W
(e.g. 512 x 20 x 15)

Region features
(here 512 x 2 x 2;
In practice e.g 512 x 7 x 7)

Region features always the same size even if input regions have different sizes!

Girshick, "Fast R-CNN", ICCV 2015.

# Cropping Features: RoI Pool

"Snap" to grid cells

Divide into 2x2 grid of (roughly) equal subregions

Project proposal onto features



CNN

Max-pool within each subregion

Input Image
(e.g. 3 x 640 x 480)

Image features: C x H x W
(e.g. 512 x 20 x 15)

Region features
(here 512 x 2 x 2;
In practice e.g 512 x 7 x 7)

Region features always the same size even if input regions have different sizes!

Girshick, "Fast R-CNN", ICCV 2015.

**Problem**: Region features slightly misaligned

# Cropping Features: RoI <u>Align</u>

Project proposal onto features

No "snapping"!

CNN

Input Image
(e.g. 3 x 640 x 480)

Image features: C x H x W
(e.g. 512 x 20 x 15)

He et al, "Mask R-CNN", ICCV 2017

# Cropping Features: RoI <u>Align</u>

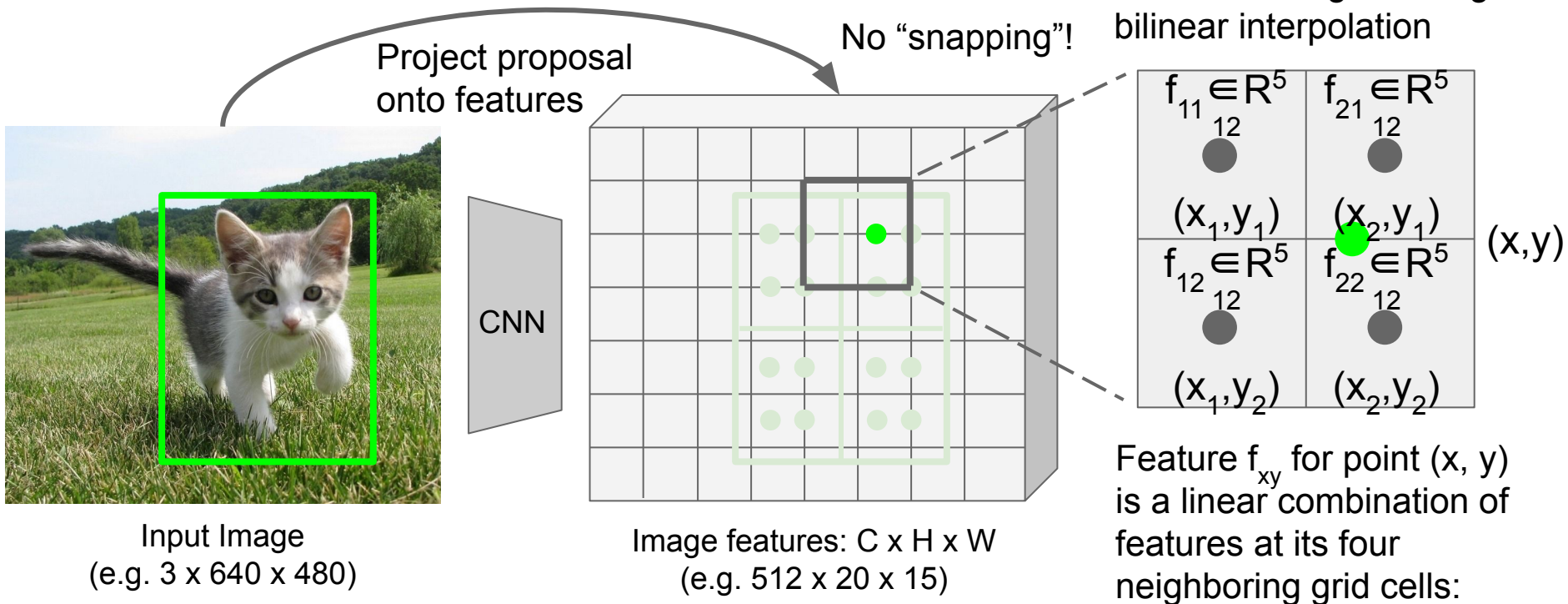Sample at regular points in each subregion using bilinear interpolation

Project proposal onto features

No "snapping"!



CNN

Input Image
(e.g. 3 x 640 x 480)

Image features: C x H x W
(e.g. 512 x 20 x 15)

He et al, "Mask R-CNN", ICCV 2017

# Cropping Features: RoI <u>Align</u>

Sample at regular points in each subregion using bilinear interpolation

Project proposal onto features

No "snapping"!

CNN

Input Image
(e.g. 3 x 640 x 480)

Image features: C x H x W
(e.g. 512 x 20 x 15)

Feature $f_{xy}$ for point (x, y) is a linear combination of features at its four neighboring grid cells:

He et al, "Mask R-CNN", ICCV 2017

# Cropping Features: RoI <u>Align</u>

Project proposal onto features

No "snapping"!

Sample at regular points in each subregion using bilinear interpolation

CNN



Input Image
(e.g. 3 x 640 x 480)

Image features: C x H x W
(e.g. 512 x 20 x 15)

$f_{11} \in R^5$   $f_{21} \in R^5$
$12$   $12$

$(x_1,y_1)$   $(x_2,y_1)$   $(x,y)$

$f_{12} \in R^5$   $f_{22} \in R^5$
$12$   $12$

$(x_1,y_2)$   $(x_2,y_2)$

Feature $f_{xy}$ for point (x, y) is a linear combination of features at its four neighboring grid cells:

$$f_{xy} = \sum_{i,j=1}^{2} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_j|)$$

He et al, "Mask R-CNN", ICCV 2017

# Cropping Features: RoI <u>Align</u>

Sample at regular points in each subregion using bilinear interpolation

Project proposal onto features

No "snapping"!

Max-pool within each subregion

CNN

Input Image
(e.g. 3 x 640 x 480)

Image features: C x H x W
(e.g. 512 x 20 x 15)

Region features
(here 512 x 2 x 2;
In practice e.g 512 x 7 x 7)

He et al, "Mask R-CNN", ICCV 2017

# R-CNN vs Fast R-CNN



**Training time (Hours)**
- R-CNN: 84
- SPP-Net: 25.5
- Fast R-CNN: 8.75

**Test time (seconds)**
- Including Region propos… / Excluding Region Propo…
- R-CNN: 49 / 47
- SPP-Net: 4.3 / 2.3
- Fast R-CNN: 2.3 / 0.32

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014
Girshick, "Fast R-CNN", ICCV 2015

# R-CNN vs Fast R-CNN



**Training time (Hours)**

R-CNN: 84
SPP-Net: 25.5
Fast R-CNN: 8.75

**Test time (seconds)**

Including Region propos... — Excluding Region Propo...

R-CNN: 49 / 47
SPP-Net: 4.3 / 2.3
Fast R-CNN: 2.3 / 0.32

**Problem**:
Runtime dominated by region proposals!

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014
Girshick, "Fast R-CNN", ICCV 2015

# Fast**er** R-CNN:

Make CNN do proposals!

Insert **Region Proposal Network (RPN)** to predict proposals from features

Otherwise same as Fast R-CNN: Crop features for each proposal, classify each one



Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015
Figure copyright 2015, Ross Girshick; reproduced with permission

# Region Proposal Network



Input Image
(e.g. 3 x 640 x 480)

Image features
(e.g. 512 x 20 x 15)

# Region Proposal Network

Imagine an **anchor box** of fixed size at each point in the feature map



Input Image
(e.g. 3 x 640 x 480)

CNN

Image features
(e.g. 512 x 20 x 15)

# Region Proposal Network

Imagine an **anchor box** of fixed size at each point in the feature map



Input Image
(e.g. 3 x 640 x 480)

CNN

Image features
(e.g. 512 x 20 x 15)

Conv

Anchor is an object?
1 x 20 x 15

At each point, predict whether the corresponding anchor contains an object (binary classification)

# Region Proposal Network



Imagine an **anchor box** of fixed size at each point in the feature map

Input Image
(e.g. 3 x 640 x 480)

CNN

Image features
(e.g. 512 x 20 x 15)

Conv

Anchor is an object?
1 x 20 x 15

Box corrections
4 x 20 x 15

For positive boxes, also predict a corrections from the anchor to the ground-truth box (regress 4 numbers per pixel)

# Region Proposal Network

In practice use K different anchor boxes of different size / scale at each point



Input Image
(e.g. 3 x 640 x 480)

CNN

Image features
(e.g. 512 x 20 x 15)

Conv

Anchor is an object?
**K** x 20 x 15

Box transforms
**4K** x 20 x 15

# Region Proposal Network

In practice use K different anchor boxes of different size / scale at each point



Input Image
(e.g. 3 x 640 x 480)

CNN

Image features
(e.g. 512 x 20 x 15)

Conv

Anchor is an object?
**K** x 20 x 15

Box transforms
**4K** x 20 x 15

Sort the K*20*15 boxes by their "objectness" score, take top ~300 as our proposals

# Fast**er** R-CNN:
Make CNN do proposals!

Jointly train with 4 losses:
1. RPN classify object / not object
2. RPN regress box coordinates
3. Final classification score (object classes)
4. Final box coordinates

Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015
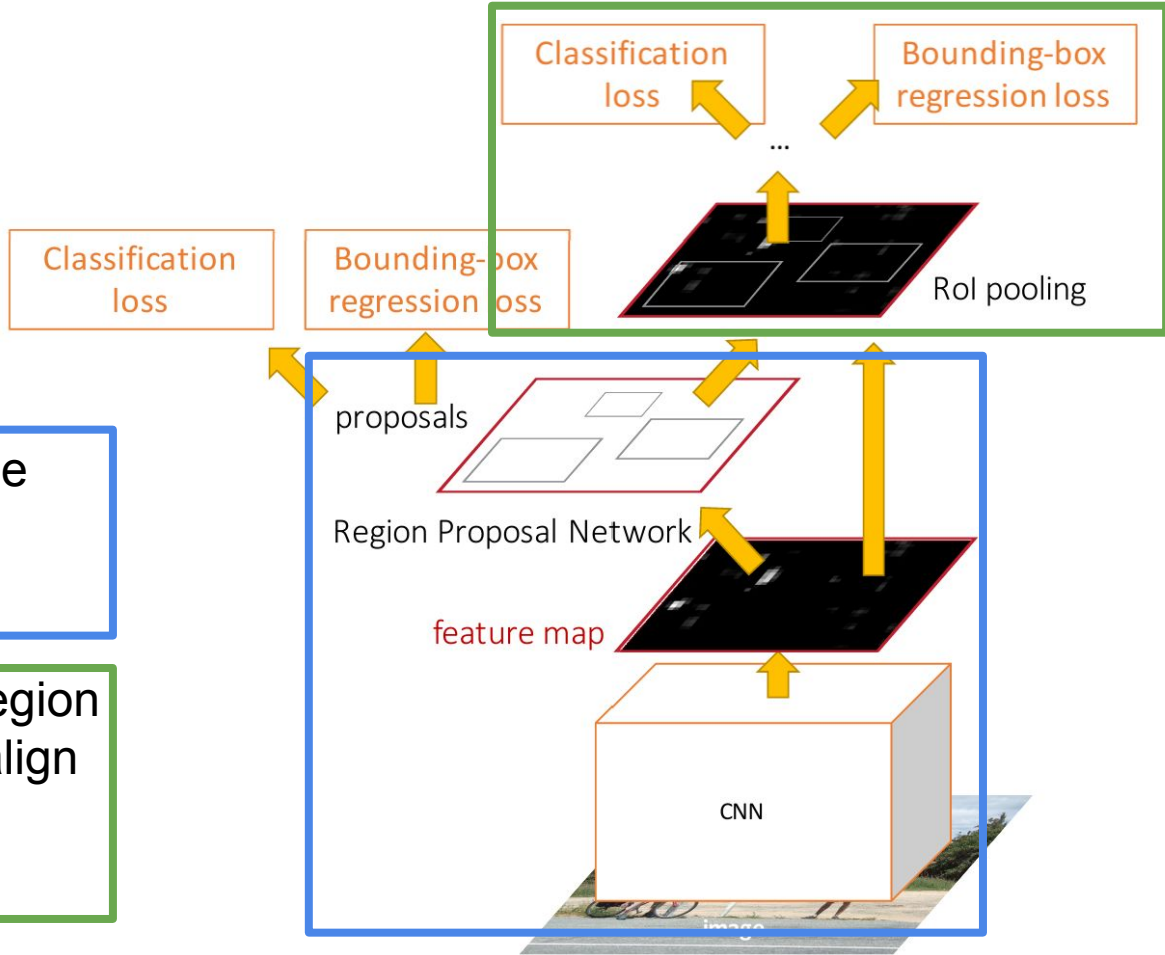Figure copyright 2015, Ross Girshick; reproduced with permission



Classification loss

Bounding-box regression loss

...

RoI pooling

Classification loss

Bounding-box regression loss

proposals

Region Proposal Network

feature map

CNN

image

# Fast**er** R-CNN:
Make CNN do proposals!

## R-CNN Test-Time Speed

| Model | Speed |
|-------|-------|
| R-CNN | 49 |
| SPP-Net | 4.3 |
| Fast R-CNN | 2.3 |
| Faster R-CNN | 0.2 |

# Fast**er** R-CNN:

Make CNN do proposals!

Glossing over many details:

- Ignore overlapping proposals with **non-max suppression**
- How are anchors determined?
- How do we sample positive / negative samples for training the RPN?
- How to parameterize bounding box regression?



Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015
Figure copyright 2015, Ross Girshick; reproduced with permission

# Fast**er** R-CNN:
Make CNN do proposals!

Faster R-CNN is a
**Two-stage object detector**

First stage: Run once per image
- Backbone network
- Region proposal network

Second stage: Run once per region
- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset



Classification loss

Bounding-box regression loss

Classification loss

Bounding-box regression loss

RoI pooling

proposals

Region Proposal Network

feature map

CNN

# Fas**ter** R-CNN:

Make CNN do proposals!

Faster R-CNN is a
**Two-stage object detector**

First stage: Run once per image
- Backbone network
- Region proposal network

Second stage: Run once per region
- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset

Do we really need
the second stage?

Classification loss

Bounding-box regression loss

...

Classification loss

Bounding-box regression loss

RoI pooling

proposals

Region Proposal Network

feature map

CNN

image

# Single-Stage Object Detectors: YOLO / SSD / RetinaNet



Input image
3 x H x W

Divide image into grid
7 x 7

Image a set of **base boxes**
centered at each grid cell
Here B = 3

Within each grid cell:
- Regress from each of the B base boxes to a final box with 5 numbers:
  (dx, dy, dh, dw, confidence)
- Predict scores for each of C classes (including background as a class)
- Looks a lot like RPN, but category-specific!

Output:
7 x 7 x (5 * B + C)

Redmon et al, "You Only Look Once:
Unified, Real-Time Object Detection", CVPR 2016
Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016
Lin et al, "Focal Loss for Dense Object Detection", ICCV 2017

# Object Detection: Lots of variables ...

**Backbone Network**
VGG16
ResNet-101
Inception V2
Inception V3
Inception
ResNet
MobileNet

**"Meta-Architecture"**
Two-stage: Faster R-CNN
Single-stage: YOLO / SSD
Hybrid: R-FCN

**Image Size**
**# Region Proposals**
**…**

**Takeaways**
Faster R-CNN is slower but more accurate

SSD is much faster but not as accurate

Bigger / Deeper backbones work better

Huang et al, "Speed/accuracy trade-offs for modern convolutional object detectors", CVPR 2017

R-FCN: Dai et al, "R-FCN: Object Detection via Region-based Fully Convolutional Networks", NIPS 2016
Inception-V2: Ioffe and Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift", ICML 2015
Inception V3: Szegedy et al, "Rethinking the Inception Architecture for Computer Vision", arXiv 2016
Inception ResNet: Szegedy et al, "Inception-V4, Inception-ResNet and the Impact of Residual Connections on Learning", arXiv 2016
MobileNet: Howard et al, "Efficient Convolutional Neural Networks for Mobile Vision Applications", arXiv 2017

# Object Detection: Lots of variables ...
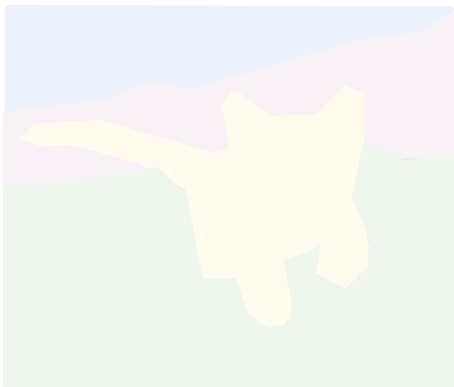
**Backbone Network**
VGG16
ResNet-101
Inception V2
Inception V3
Inception
ResNet
MobileNet

**"Meta-Architecture"**
Two-stage: Faster R-CNN
Single-stage: YOLO / SSD
Hybrid: R-FCN

**Image Size**
**# Region Proposals**
**...**

**Takeaways**
Faster R-CNN is slower but more accurate

SSD is much faster but not as accurate

Bigger / Deeper backbones work better

Huang et al, "Speed/accuracy trade-offs for modern convolutional object detectors", CVPR 2017
Zou et al, "Object Detection in 20 Years: A Survey", arXiv 2019

R-FCN: Dai et al, "R-FCN: Object Detection via Region-based Fully Convolutional Networks", NIPS 2016
Inception-V2: Ioffe and Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift", ICML 2015
Inception V3: Szegedy et al, "Rethinking the Inception Architecture for Computer Vision", arXiv 2016
Inception ResNet: Szegedy et al, "Inception-V4, Inception-ResNet and the Impact of Residual Connections on Learning", arXiv 2016
MobileNet: Howard et al, "Efficient Convolutional Neural Networks for Mobile Vision Applications", arXiv 2017

# Instance Segmentation

**Classification**

CAT

No spatial extent

**Semantic Segmentation**

GRASS, CAT, TREE, SKY

No objects, just pixels

**Object Detection**

DOG, DOG, CAT

**Instance Segmentation**

**DOG**, **DOG**, **CAT**

Multiple Object

# Object Detection: Faster R-CNN

# Instance Segmentation: Mask R-CNN



Add a small mask network that operates on each RoI and predicts a 28x28 binary mask

He et al, "Mask R-CNN", ICCV 2017

# Mask R-CNN



Classification Scores: C
Box coordinates (per class): 4 * C

CNN +RPN

RoI Align

256 x 14 x 14

Conv

256 x 14 x 14

Conv

Predict a mask for each of C classes

C x 28 x 28

He et al, "Mask R-CNN", arXiv 2017

# Mask R-CNN: Example Mask Training Targets

# Mask R-CNN: Example Mask Training Targets

# Mask R-CNN: Example Mask Training Targets

# Mask R-CNN: Example Mask Training Targets

# Mask R-CNN: Very Good Results!



He et al, "Mask R-CNN", ICCV 2017

# Mask R-CNN
# Also does pose



He et al, "Mask R-CNN", ICCV 2017

# Open Source Frameworks

Lots of good implementations on GitHub!

TensorFlow Detection API:

https://github.com/tensorflow/models/tree/master/research/object_detection

Faster RCNN, SSD, RFCN, Mask R-CNN, ...

Detectron2 (PyTorch)

https://github.com/facebookresearch/detectron2

Mask R-CNN, RetinaNet, Faster R-CNN, RPN, Fast R-CNN, R-FCN, ...

Finetune on your own dataset with pre-trained models

Beyond 2D Object Detection...

# Object Detection + Captioning = Dense Captioning



Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016
Figure copyright IEEE, 2016. Reproduced for educational purposes.

# Dense Video Captioning



Ranjay Krishna et al., "Dense-Captioning Events in Videos", ICCV 2017
Figure copyright IEEE, 2017. Reproduced with permission.

# Objects + <u>Relationships</u> = Scene Graphs



108,077 Images
5.4 Million Region Descriptions
1.7 Million Visual Question Answers
3.8 Million Object Instances
2.8 Million Attributes
2.3 Million Relationships
Everything Mapped to Wordnet Synsets



Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen et al. "Visual genome: Connecting language and vision using crowdsourced dense image annotations." International Journal of Computer Vision 123, no. 1 (2017): 32-73.

# Scene Graph Prediction



Krishna, Lu, Bernstein, and Fei-Fei, "Scene Graph Generation by Iterative Message Passing", ECCV 2016
Figure copyright IEEE, 2018. Reproduced for educational purposes.

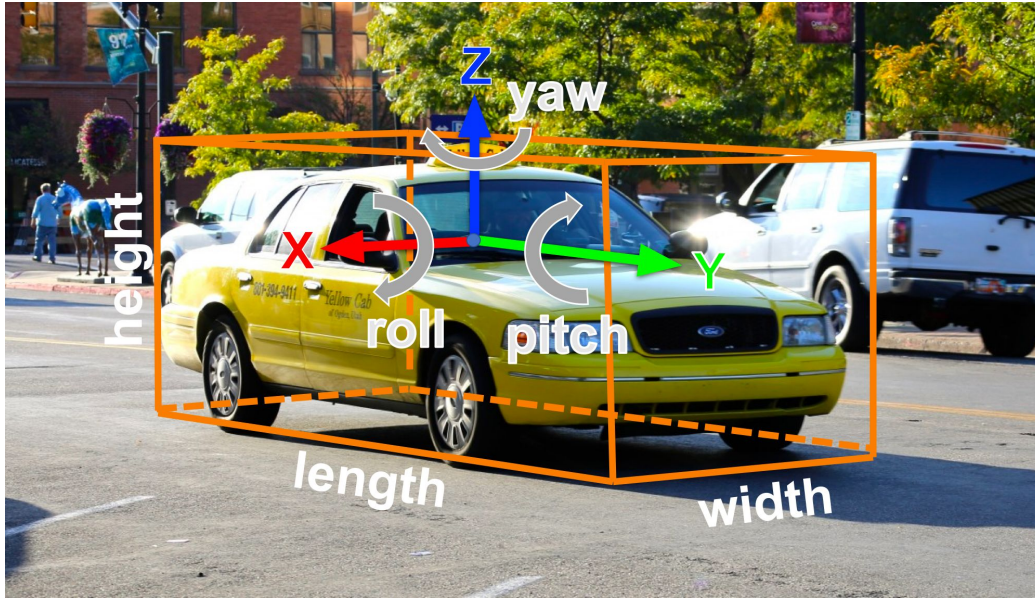# Grounded Situation Recognition



Capture semantic and physical relationships of objects

Tag each image with an action and ground each entity involved in that action

Pratt et al "Grounded Situation Recognition", ECCV 2020

# 3D Object Detection



2D Object Detection:
2D bounding box
(x, y, w, h)

3D Object Detection:
3D oriented bounding box
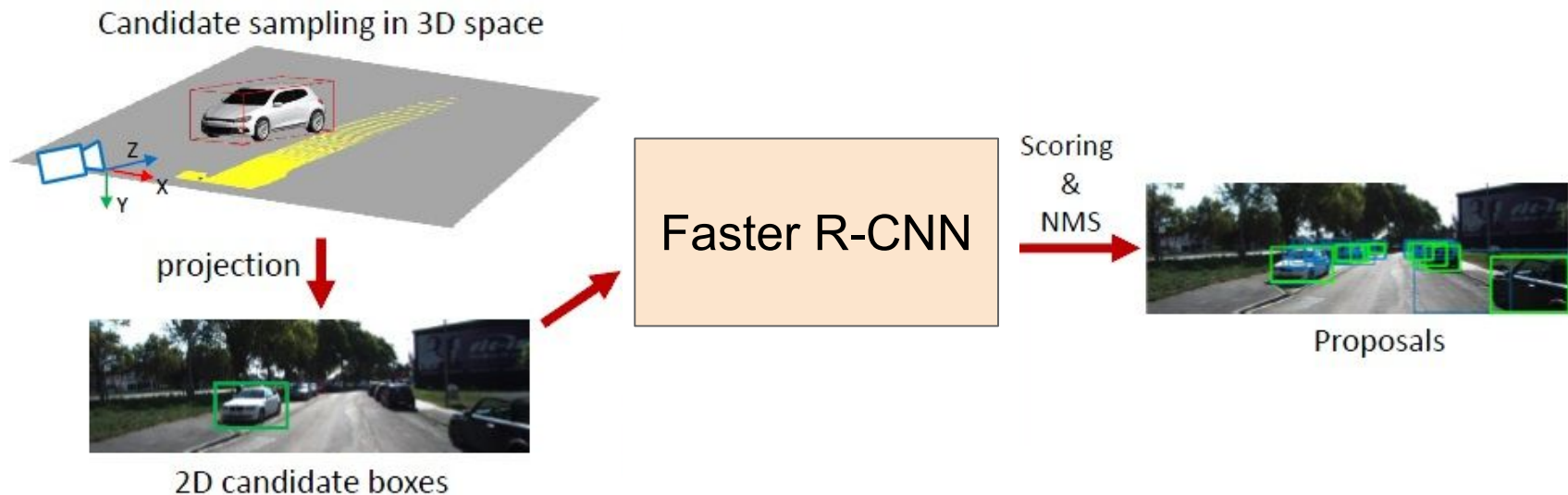(x, y, z, w, h, l, r, p, y)

Simplified bbox: no roll & pitch

Much harder problem than 2D object detection!

# 3D Object Detection: Monocular Camera



- Same idea as Faster RCNN, but proposals are in 3D
- 3D bounding box proposal, regress 3D box parameters + class score

Chen, Xiaozhi, Kaustav Kundu, Ziyu Zhang, Huimin Ma, Sanja Fidler, and Raquel Urtasun. "Monocular 3d object detection for autonomous driving." CVPR 2016.

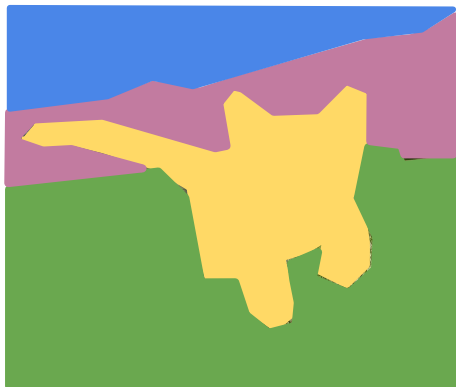# Recap: Lots of computer vision tasks!
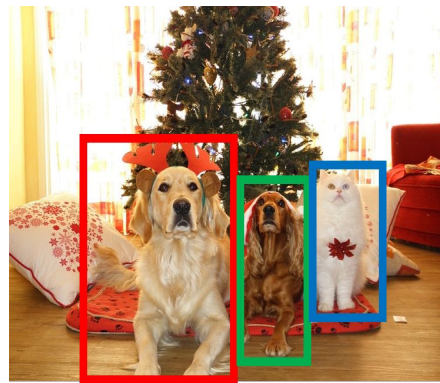


**Classification**

**CAT**

No spatial extent

**Semantic Segmentation**

**GRASS**, **CAT**, **TREE**, **SKY**

No objects, just pixels

**Object Detection**

**DOG**, **DOG**, **CAT**

**Instance Segmentation**

**DOG**, **DOG**, **CAT**

Multiple Object

This image is CC0 public domain

Next time:
Self-Supervised

# 3D Object Detection: Simple Camera Model



3D ray

2D point

camera viewing frustrum
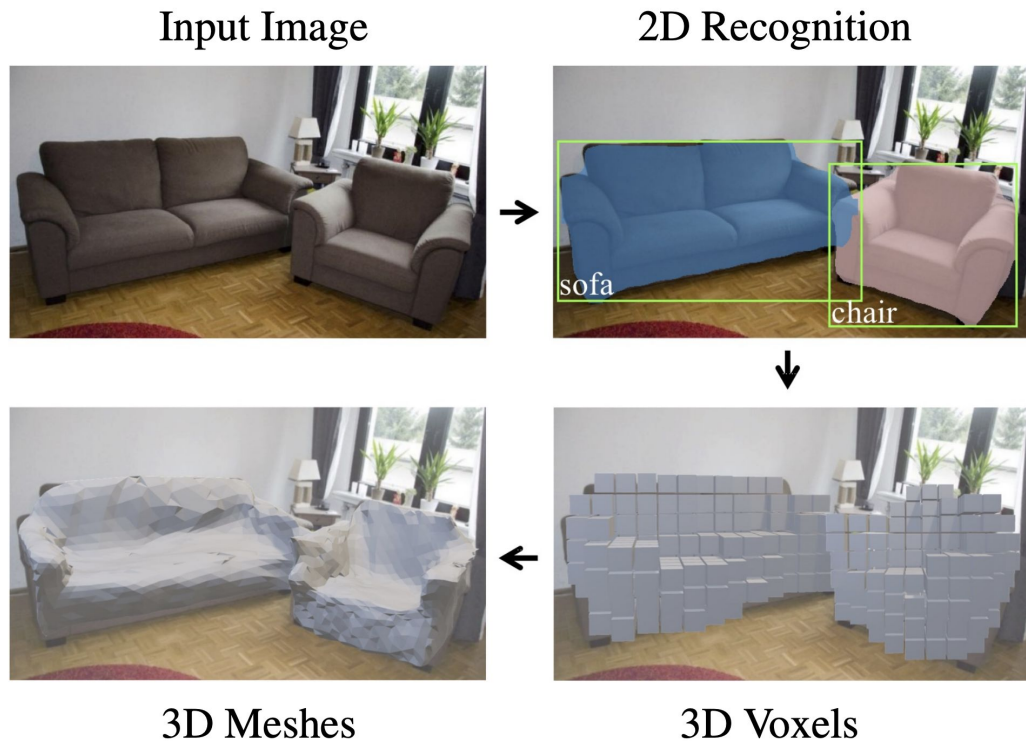
image plane

camera

A point on the image plane corresponds to a **ray** in the 3D space

A 2D bounding box on an image is a **frustrum** in the 3D space

Localize an object in 3D:
The object can be anywhere in the **camera viewing frustrum**!

# 3D Shape Prediction: Mesh R-CNN

Input Image

2D Recognition

sofa
chair

3D Meshes

3D Voxels

Gkioxari et al., Mesh RCNN, ICCV 2019