

Visual AI

CPSC 533R

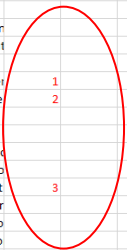
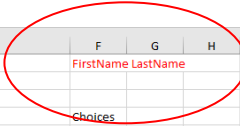
Lecture 4. Advanced architectures and sparse 2D keypoints

Helge Rhodin



Pick a paper till tonight

A	B	C	D	E	F	G	H
Please insert your name in the subsequent cell and mark your 1st, 2nd and 3rd choices in the Choices column, as exemplified below.				Name:	FirstName LastName		
There should be no read text left once you submit it to Canvas							
13 Presentation sessions							
Date*	Topic	Authors	Title & Link	Keywords	Choices		
24-Sep	Objective functions II	Bishop Barron et al.	Mixture Density Networks A General and Adaptive Robust Loss Function	probabilistic regression optimizing the objective function			
1-Oct	Self-supervision I	Vondrick et al. Doersch et al.	Tracking Emerges by Colorizing Videos Unsupervised visual representation learning by context prediction	self-supervision, color and motion self-supervision, classification, tri			
8-Oct	Human Pose	Wandt et al. Cao et al.	Repnet: Weakly supervised training of an adversarial reprojection network for Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields	3D pose, projection, self-supervisi 2D pose estimation, skeleton repr			
15-Oct	GANs	Chan et al. Park et al. Lu et al.	Everybody Dance Now Semantic Image Synthesis with Spatially-Adaptive Normalization Layered Neural Rendering for Retiming People in Video	image translation GAN, batch-norm, pixel-condition neural rendering, depth estimatio			
22-Oct	Implicit functions	Sitzmann et al. Saito et al. Peng et al.	Implicit neural representations with periodic activation functions PIFu: Pixel-Aligned Implicit Function for High-Resolution Clothed Human Digiti Convolutional Occupancy Networks	implicit function, gradient, optimi implicit functions, 3D reconstructi implicit functions, volumetric grid			
29-Oct	Animation and meshes	Xu et al. Gao et al., Hanoca et al.	RigNet: Neural Rigging for Articulated Characters Automatic Unpaired Shape Deformation Transfer Point2Mesh: A Self-prior for Deformable Meshes	animation, rigging, mesh processi autoencoder, meshes, style trans shape transfer, 3D texture, autoer			
5-Nov	Object Detection	Carion et al. Wu et al.	End-to-End Object Detection with Transformers Unsupervised Learning of Probably Symmetric Deformable 3D Objects from Im	object detection, transformer arch autoencoder, depth, symmetry, s			
10-Nov	Body models	Zuffi et al. Ma et al.	Three-D Safari: Learning to Estimate Zebra Pose, Shape, and Texture from Imag CAPE: Dressing SMPN	body mode, animals, optimization body model, mesh processing, aut			
12-Nov	Objective functions II	Peebles et al. Fort et al.	The Hessian Penalty: A Weak Prior for Unsupervised Disentanglement Deep Ensembles: A Loss Landscape Perspective	GAN, latent space disentangleme neural network theory, ensemble	1	2	
17-Nov	Motion	Chu et al. Luo et al.	Learning Temporal Coherence via Self-Supervision for GAN-based Video Gener Consistent Video Depth Estimation	GAN, temporal, self-supervision autoencoder, 3D geometry, temp			
19-Nov	Self-supervision II	Holden et al. Blielski et al. Crawford et al.	Phase-Functioned Neural Networks for Character Control Emergence of Object Segmentation in Perturbed Generative Models Spatially invariant unsupervised object detection with convolutional neural ne	character animation, temporal mo GAN, self-supervision, segmentat probabilistic modeling, self-super	3		
24-Nov	Rendering	Lorenz et al. Liu et al. Schwartz et al.	Unsupervised Part-Based Disentangling of Object Shape and Appearance Soft Rasterizer: A Differentiable Renderer for Image-based 3D Reasoning The Eyes Have It: An Integrated Eye and Face Model for Photorealistic Facial An	self-supervision, autoencoder, po differentiable rendering, occlusio autoencoder, differentiable rende			
26-Nov	View synthesis	Mildenhall et al. Hinton et al.	NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis Transforming Auto-encoders	novel view synthesis, ray-tracing, novel view synthesis, stereo, stru			
Alternatives (instead of one of the above)							
		Nguyen-Phuoc et al. Moon et al.	HoloGAN: Unsupervised learning of 3D representations from natural images I2L-MeshNet: Image-to-Lixel Prediction Network for Accurate 3D Human Pose a	self-supervision, view synthesis, 3 3D human pose and shape estimat			



Recap: Automatic differentiation and backpropagation

Forward pass

$$L(h(\text{linear}(h(\text{linear}(x, W^{(1)})), W^{(2)})))$$

$$= Lh \left(\begin{array}{|c|c|} \hline W^* & b^* \\ \hline \end{array} h \left(\begin{array}{|c|c|c|} \hline W^* & b^* & x^* \\ \hline \end{array} \right) \right)$$

$$J_x^f = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

Jacobian matrix

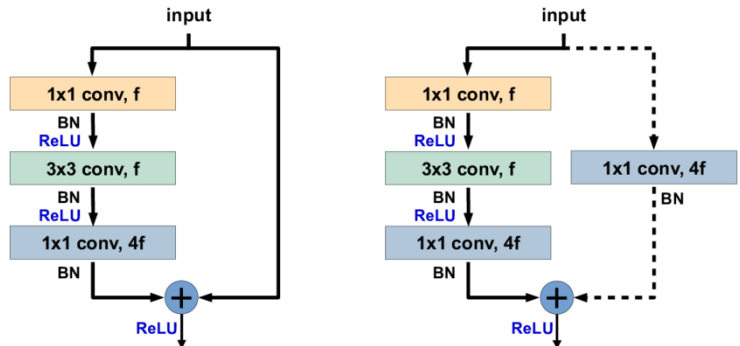
Backwards pass to $W^{(1)}$

$$J_{W^{(1)}}^E = J_x^L \left[J_x^{\text{linear}} \left[J_x^h \left[J_W^{\text{linear}} \right] \right] \right]$$

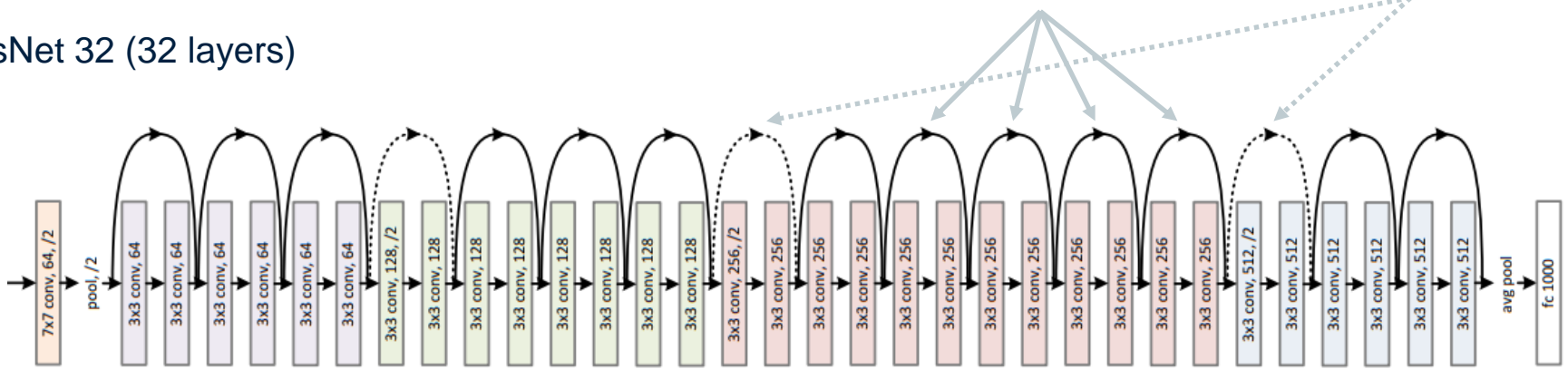
ResNet details

What if number of channels changes?

- apply a linear transformation to match the size
 - a projection on a linear subspace, related to principal component analysis (PCA)



ResNet 32 (32 layers)



Batch normalization

[Batch Normalization: Accelerating Deep Network Training ~~by Reducing Internal Covariate Shift~~]

- Normalize after each linear + activation function
 - normalize across minibatch, to have $\mu=0$ and $\sigma=1$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

$$\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

- Strict normalization reduces performance, hence, add a learnable offset and scale

$$y_i \leftarrow \gamma \hat{x}_i + \beta$$

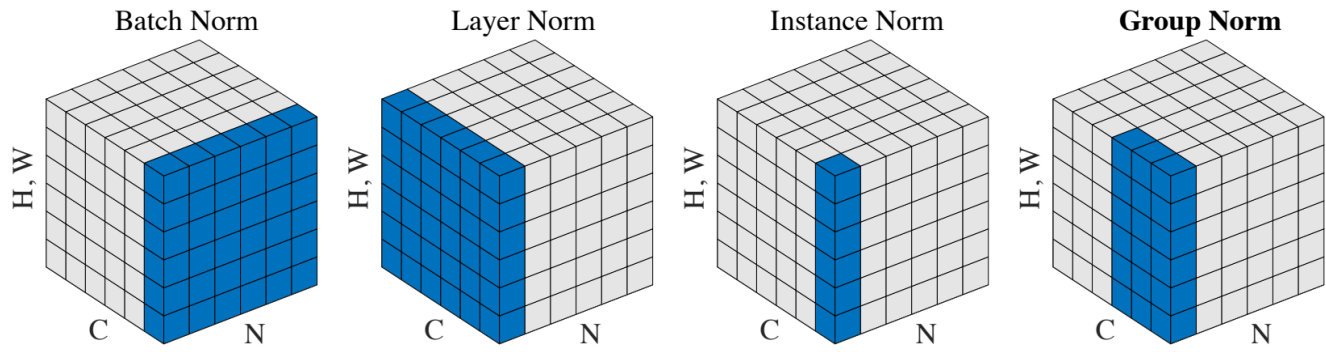
- What if we only have a single image at inference time?
 - Re-apply mean and variance recorded during training (using exponential moving average)

Batch normalization effect and variants

What is the benefit of first normalizing and then ‘denormalizing’?

- noise from other images regularizes
- it separates learning of the variance (scale) and bias (offset) from the values itself
- Empirical: training deeper networks, with sigmoid activation, higher learning rate, and faster convergence

Variants normalize over different slices of the feature tensor:



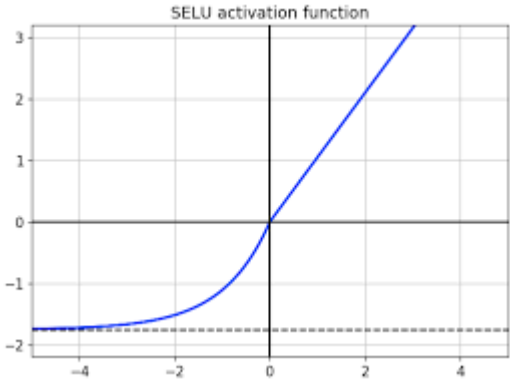
[Wu and He. Group Normalization]

Self-normalizing neural networks

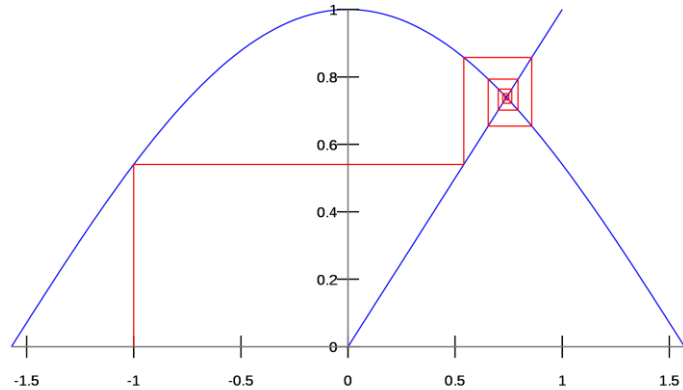
Self-normalizing Neural Networks

[Klambauer et al.]

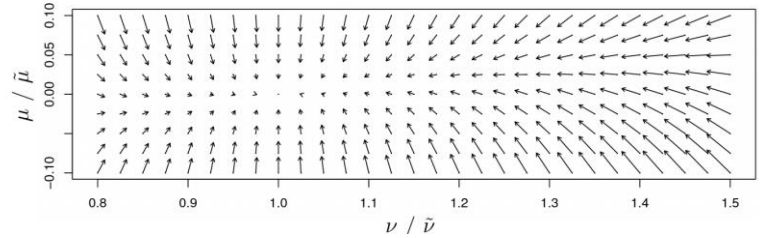
- fixed point enforced by choice of activation function (SELUs)
- stable and attracting fixed point for the function g that maps mean and variance from one layer to the next
- Possibility to train deep fully connected NNs



$$\text{selu}(x) = \lambda \begin{cases} x & \text{if } x > 0 \\ \alpha e^x - \alpha & \text{if } x \leq 0 \end{cases}$$



Fixed point iterations for $\cos(x)$



Mapping of the function g towards $\mu=0$ and $v=1$

Regularization

Dropout

- randomly zero out activations
- re-weight the non-zero ones to maintain the distribution of the unmodified activations
 - induced noise reduces overfitting

Weight decay

$$\tilde{\mathbf{w}} = (1 - \tau)\mathbf{w} \text{ with } \tau \text{ small}$$

Prior on neural network weights

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^2$$

Weight decay and square prior are equivalent under certain conditions (vanilla SGD without momentum)

Check out AdamW in pytorch

Residual networks and skip connections

- Deep networks are hard to train

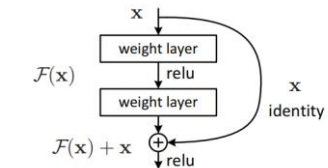
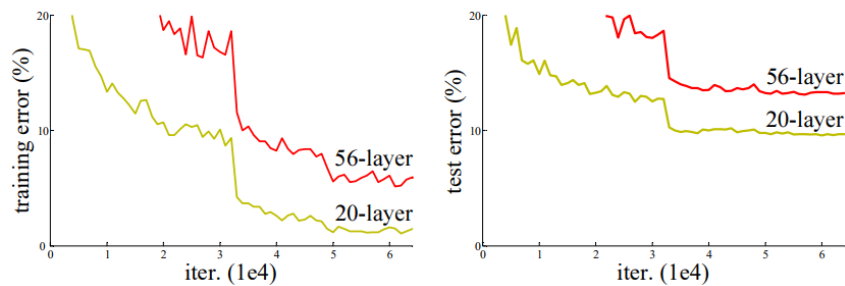


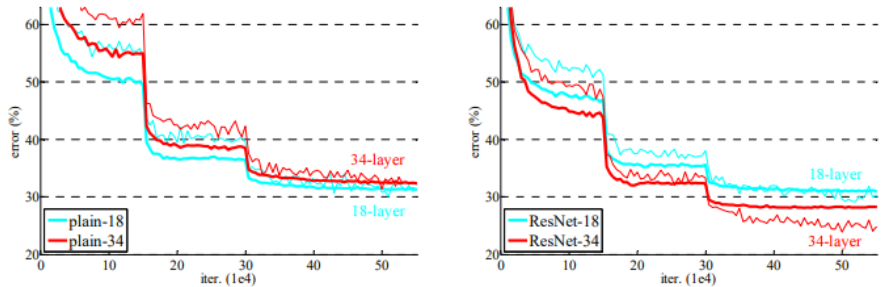
Figure 2. Residual learning: a building block.

- Residual blocks with shortcut/skip connections

$$y = F(x) + x$$

- no extra parameters
- enables training of deep neural networks

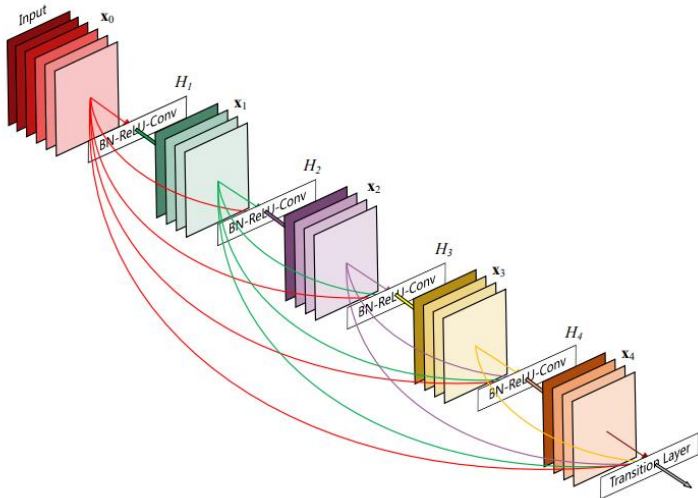
Image net training



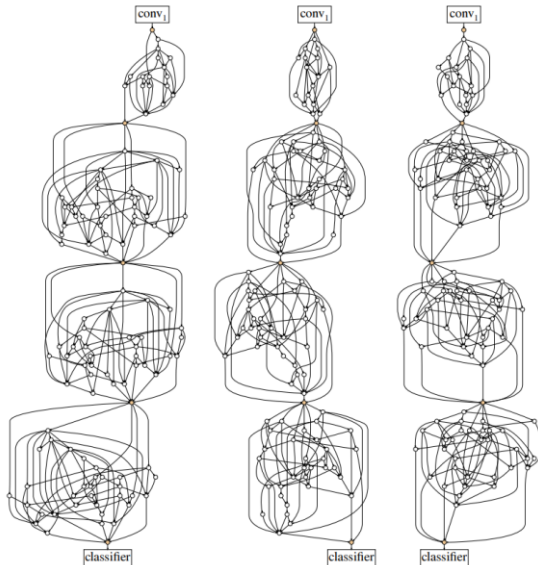
Baseline network

ResNet version

Other network architectures



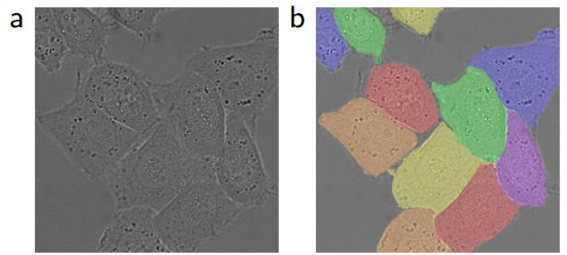
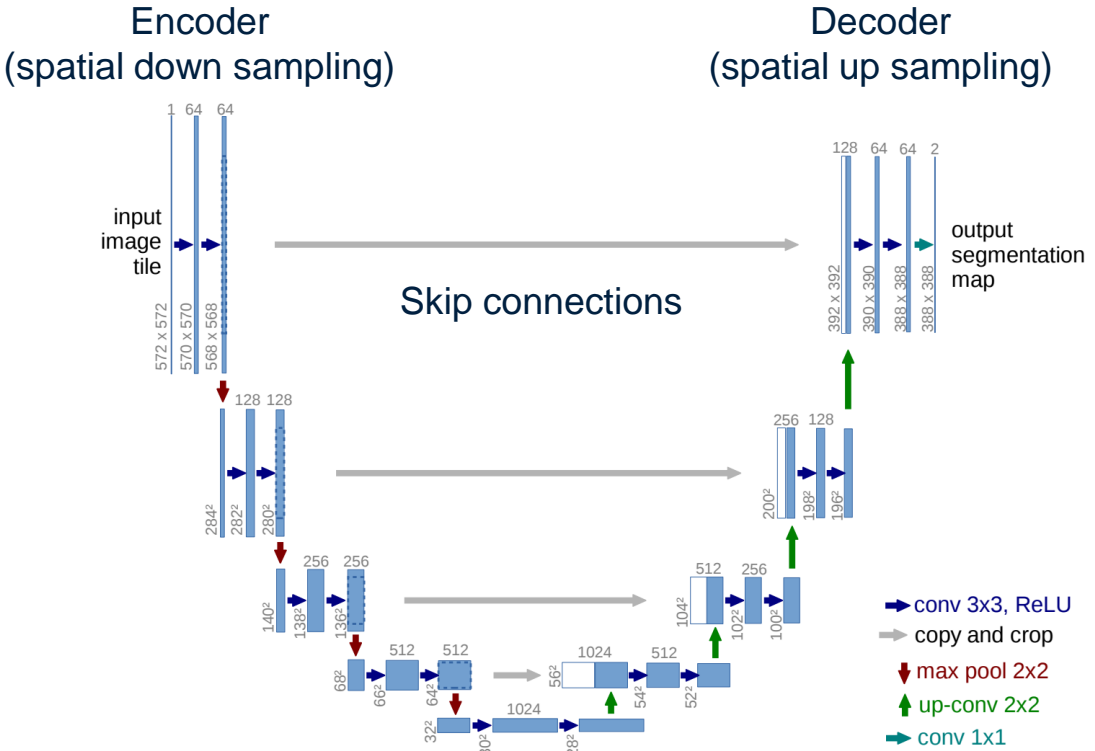
DenseNet
(skip connection to all future layers)



Randomly wired networks
(search for best wiring among candidates)

U-Net architecture

- Similar input and output resolution
- A global encoding is learned by down sampling (to 32 x 32 px)
- Progressive increase of channels maintains throughput / capacity
- Skip connections preserve details



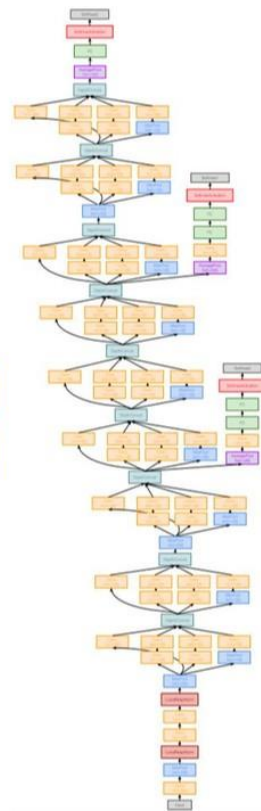
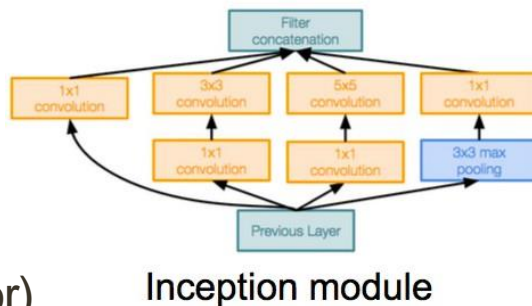
[U-Net: Convolutional Networks for Biomedical Image Segmentation]

GoogleLeNet (Inception Net V1)

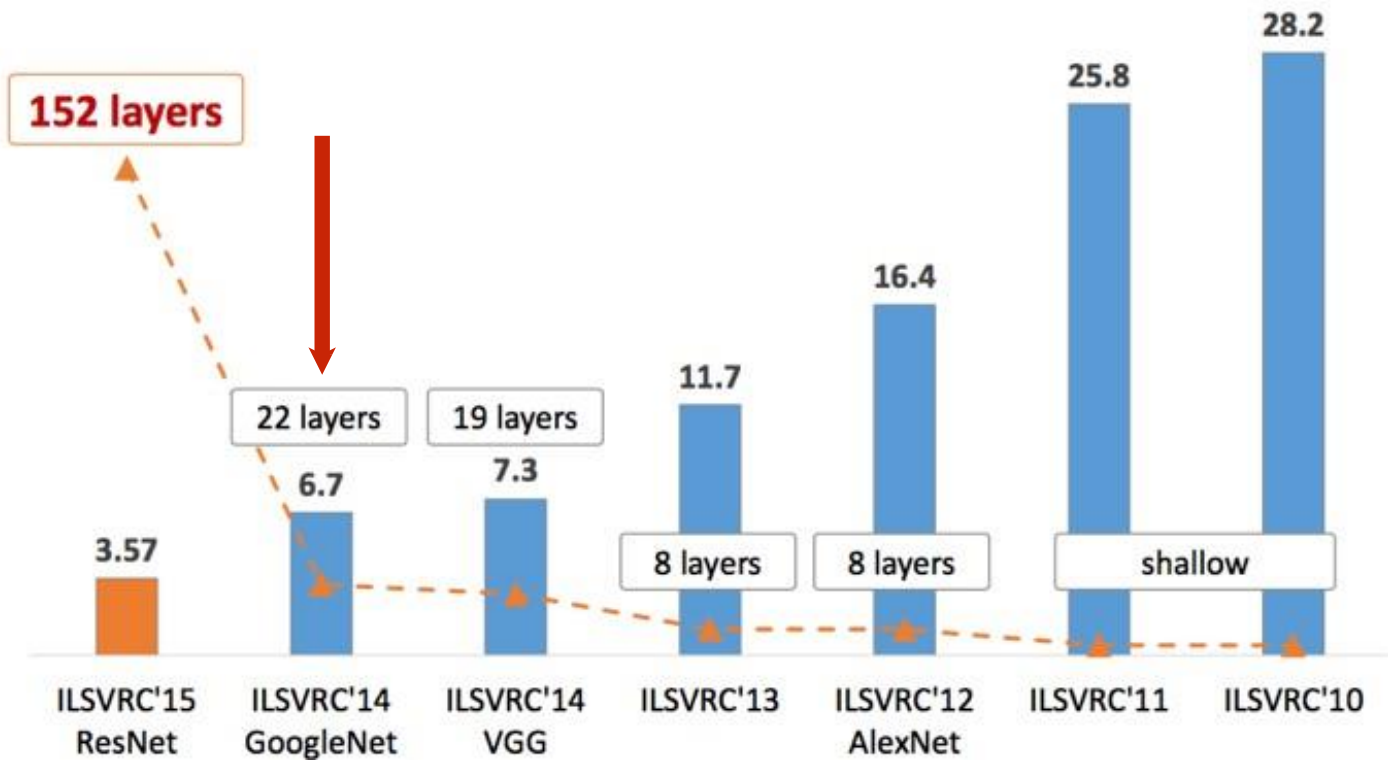
[Szegedy et al., 2014]

even deeper network with **computational efficiency**

- 22 layers
- Efficient “Inception” module
- No FC layers
- Only 5 million parameters!
(12x less than AlexNet!)
- Better performance (@6.7 top 5 error)



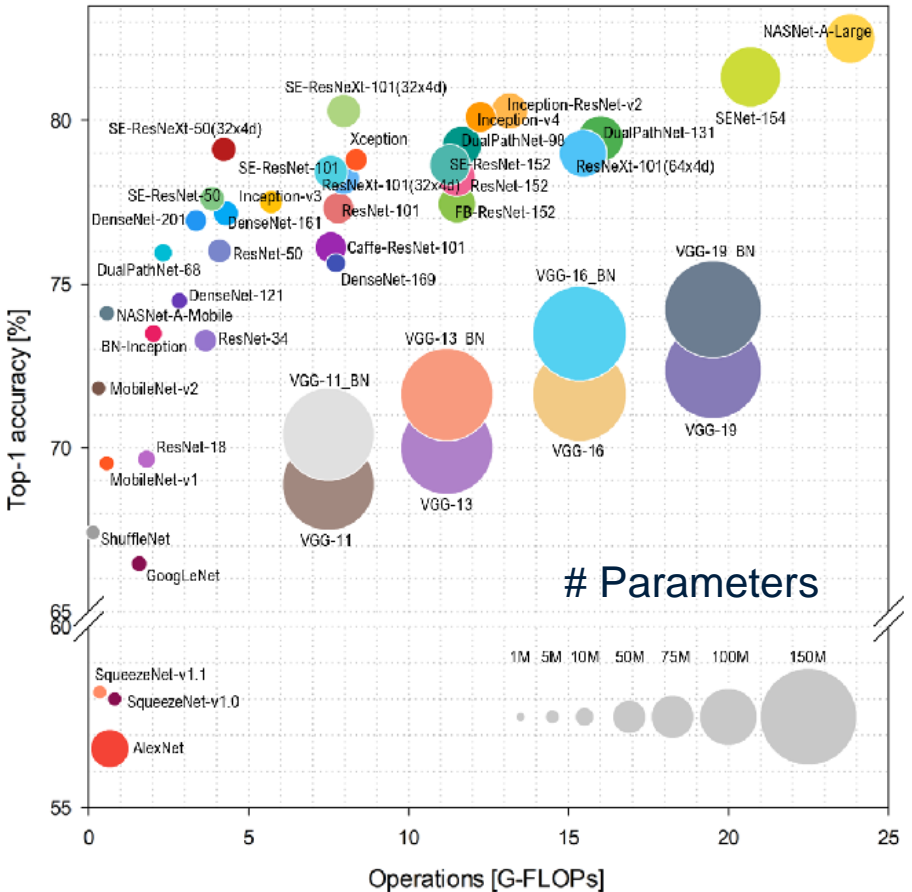
ILSVRC winner 2012



Network comparison

The goal is to balance

- accuracy (maximize)
- number of parameters (minimize)
- number of operations (minimize)
- efficiency of operations (maximize)
 - cache efficiency
 - parallel execution
 - precision, e.g. float, float16, binary,...



History of Inception Networks

Inception V1 GoogleNet

- Network in network approach

Inception V2

- Use of batch normalization
- [Batch normalization: Accelerating deep network training by ...]

Inception V3

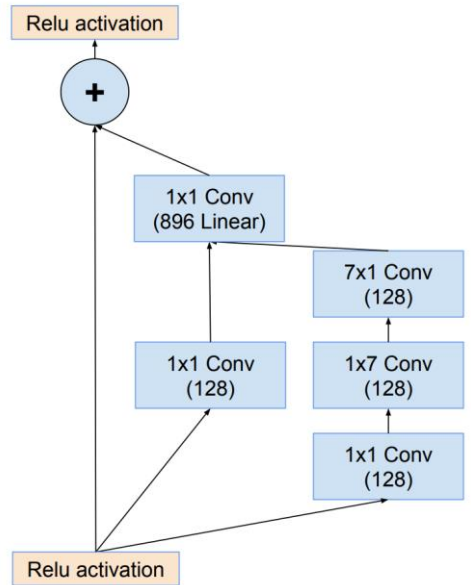
- Factorizations
- [Rethinking the inception architecture for computer vision]

Inception V4

- Tuning of filters
- [Inception-v4, Inception-ResNet and the Impact of ...]

Inception-ResNet

- Skip connections instead of concatenations



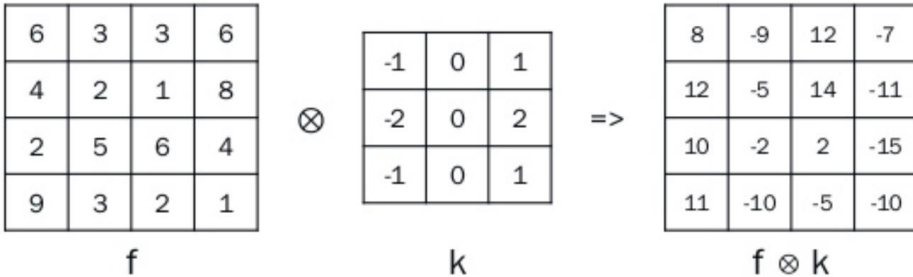
Inception ResNet block example

Separable Convolutions

Idea:

Separate a single convolution operation into a sequence of simpler operations

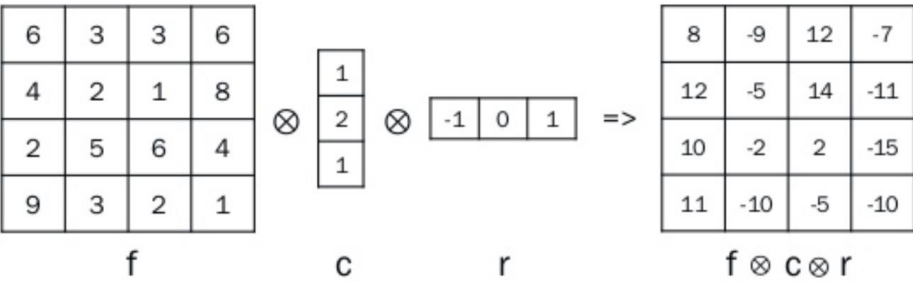
- e.g., 7x7 convolution into
 - 1x7 and 7x1
- reduction of parameters
 - e.g., 14 vs. 49 for 7x7 conv.



Drawback:

It models simpler functions

- no 'diagonal' entries possible
- successive layers can't be run in parallel

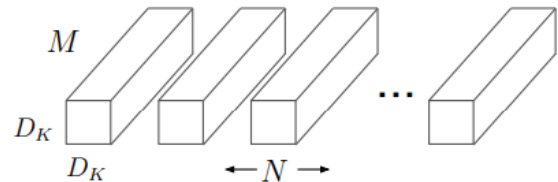




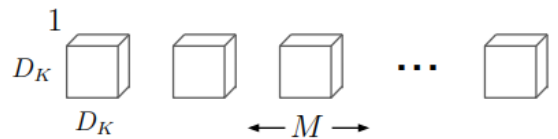
Mobile Net V1

Depthwise separable convolution

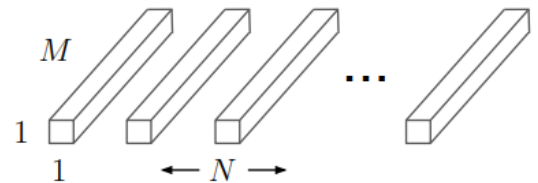
- separate a 3x3 convolution with M input and N output channels
 - M 3x3 convolutions, each applied on a single channel
 - N 1x1 convolutions, combining the M intermediate features
 - add ReLU and Batch Normalization after each layer



(a) Standard Convolution Filters



(b) Depthwise Convolutional Filters



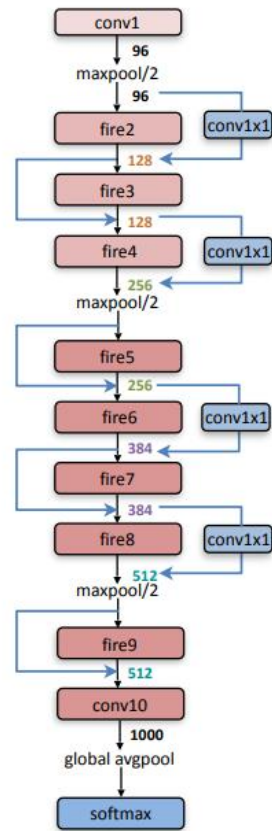
Advantages

- fewer add-mul operations
(8-9 times less than conventional convolution)
- highly efficient operations
 - 95% of time spend on 1x1 convolutions
 - 1x1 convolutions are highly optimized
 - an instance of general matrix multiplication (GMM)

SqueezeNet

Goal: Very small model size and efficient execution

- Strategy 1. Replace 3x3 filters with 1x1 filters
 - a 1x1 filter has 9X fewer parameters than a 3x3 filter
- Strategy 2. Decrease the number of input channels to 3x3 filters
 - $(\text{number of input channels}) * (\text{number of output channels}) * (3*3)$
- Strategy 3. Spatially downsample late
 - known to yield higher accuracy
- No fully connected layers
 - use global average pooling instead



SqueezeNet architecture

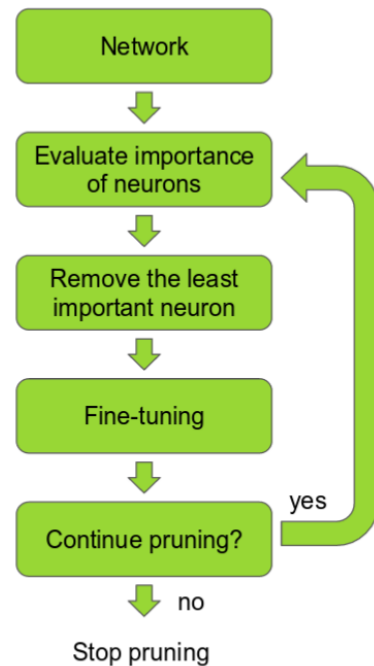
Reducing the model footprint

Pruning

- rank the neurons in the network according to how much they contribute, e.g.:
 - L1/L2 mean of neuron weights
 - their mean activations
 - the number of times a neuron wasn't zero on some validation set
- remove the low-ranking neurons

Deep Compression [Han et al., 2015]

- quantize CNN parameters (e.g. 8-bits of precision)
- uses a codebook



[Pruning Convolutional Neural Networks for Resource Efficient Inference]

ResNeXt: Aggregated Residual Transformations

Idea: “vertical residual blocks”

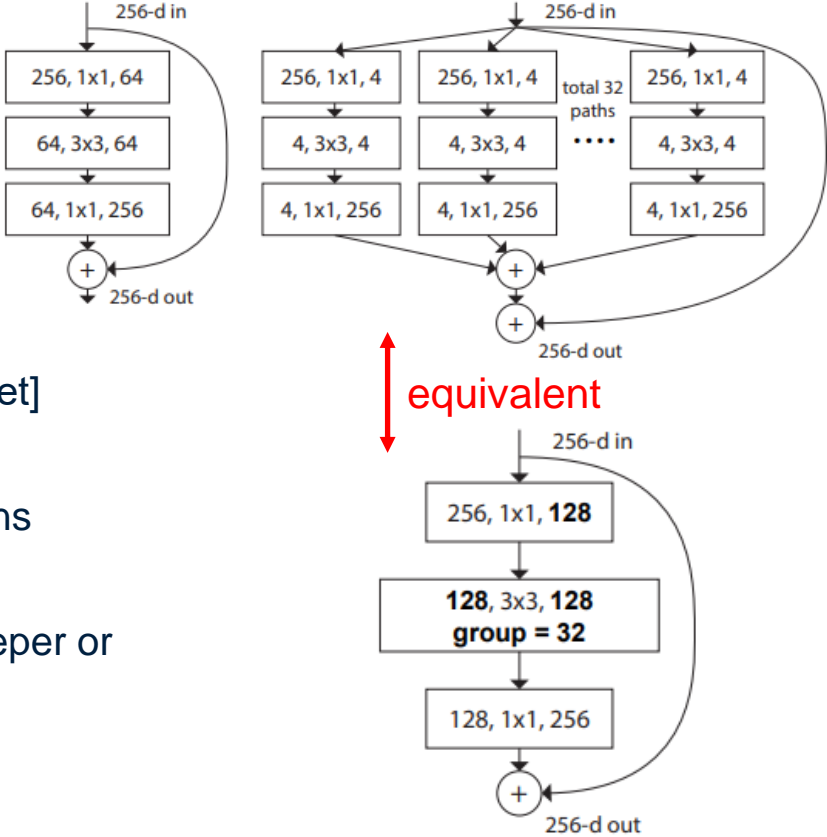
- create blocks with identical topology
 - and independent weights
 - replicate them c times (“cardinality”)
- add these blocks together
- add a skip connection as in ResNet

Related: Krizhevsky et al.’s grouped convolutions [AlexNet]

Advantage:

- larger number of channels, same number of operations
- improved performance
- increasing cardinality is more effective than going deeper or wider when we increase the capacity

A new (NeXt) dimension: depth, width, and **cardinality**



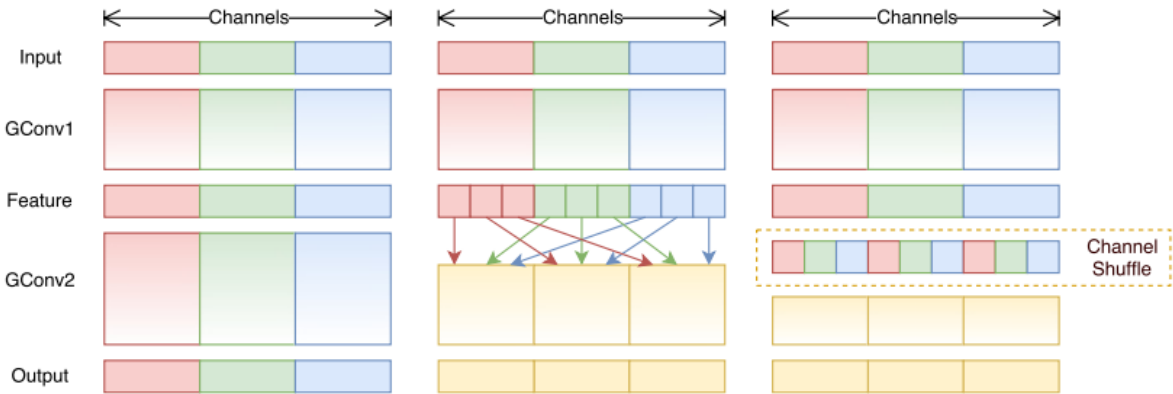
ShuffleNet

Idea:

- Group-wise convolution
- shuffle features for cross-talk

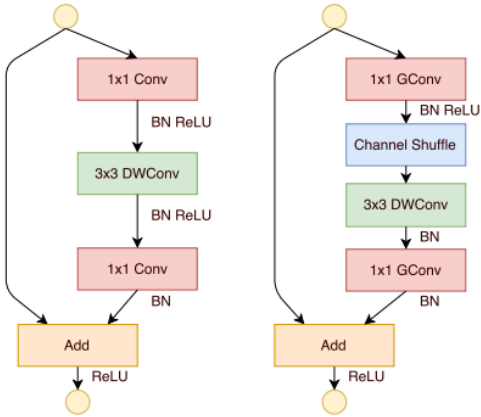
Advantage:

- less parameters for 1x1 conv
 - $1/\#g$ parameters, where $\#g$ is the cardinality
- recall that MobileNet spent 95% in 1x1 convolution



MobileNet block

ShuffleNet block



Discussion



Assignment I & II

Assignment I

- is due ~~today~~ tomorrow!
- submit on Canvas

Assignment II

- released today
- later parts may require content from Thursday
 - start with preliminaries and Task I

```

0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9
    
```

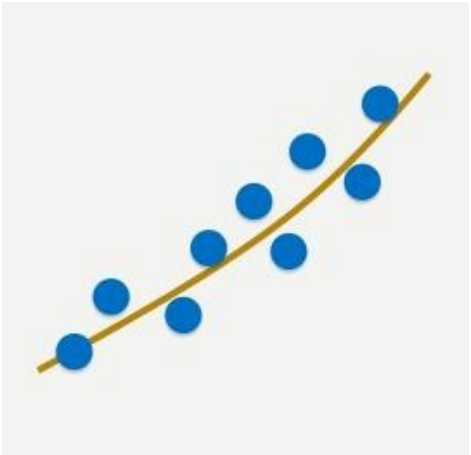
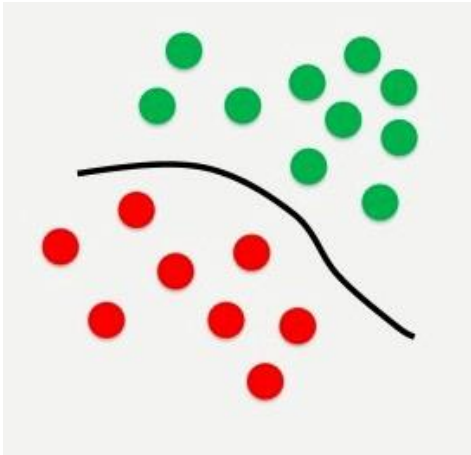


Classification vs. regression

Classification



Regression



Classification and regression

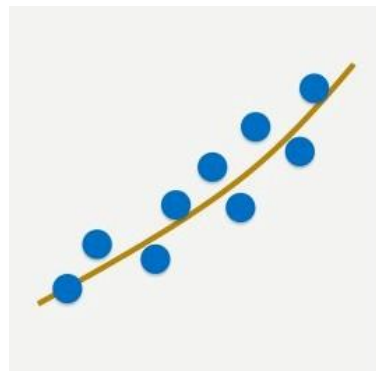
Regression

- for continuous values

$$\text{nn}(\mathbf{x}) \rightarrow y \in \mathbb{R}$$

- squared loss is most common

$$l_2(y, l) = (y - l)^2$$



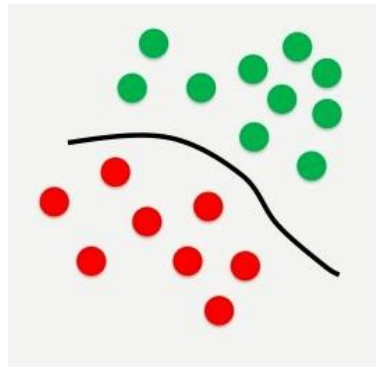
Classification

- discrete classes

$$\text{nn}(\mathbf{x}) \rightarrow \mathbf{y} \in [0, 1]$$

- naïve least-squares loss

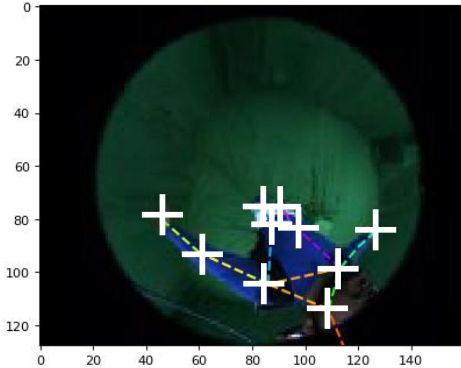
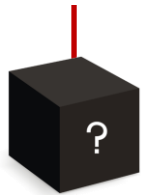
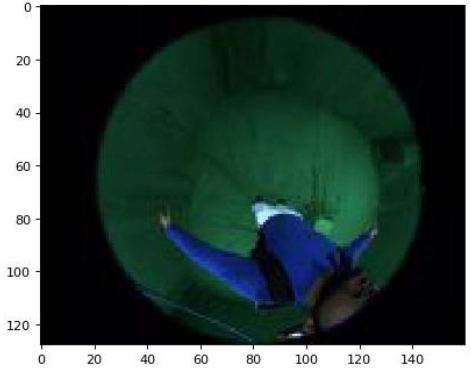
$$l_2(\mathbf{y}, \mathbf{1}) = \|\mathbf{y} - \mathbf{1}\|^2$$



Regression-based 2D pose estimation

A classical regression task

- Input:
 - grid of color values, an image ($3 \times W \times H$)
- Output:
 - pairs of continuous values, the position in the image
 - one pair for each of the K keypoints ($2 \times K$)
- Neural network architecture:
 - Some convolutional layers to infer an internal representation of the human pose ($C \times W' \times H'$)
 - One or more fully-connected layers to aggregate spatial information into the output values ($C * W' * H'$) \rightarrow ($2 \times K$)

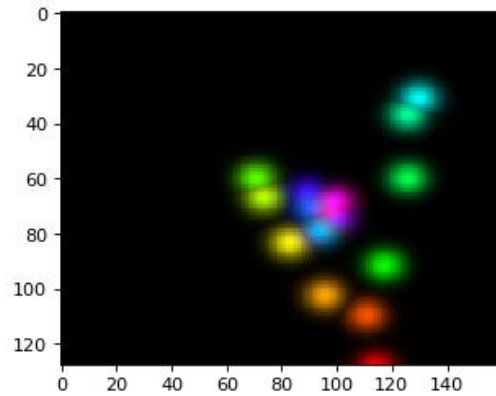
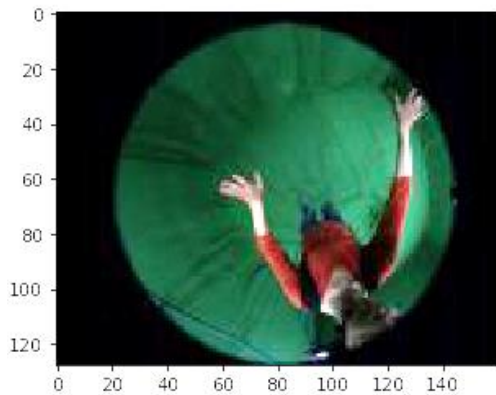


Heatmap-based 2D pose estimation

Phrase the regression task as classification

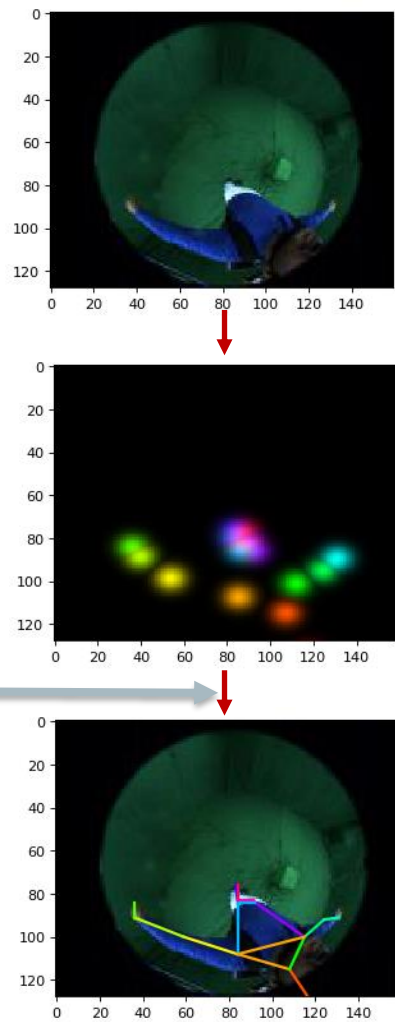
- separate heatmap H_j for each joint j
- Each pixel of H_j encodes the ‘probability’ of containing joint j
 - not a true probability as pixels don’t sum to one
- Advantages:
 - Inferred with fully convolutional networks
 - less parameters than fully connected ones (MLPs)
 - applies to arbitrary image resolution and aspect ratio (can be different from training)
 - translation invariance
 - locality
 - Generalizes to multiple and arbitrary number of persons

[Tompson et al., Efficient object localization using convolutional networks.]



Disadvantages of heatmaps

- Disadvantage:
 - Large image scale variations
 - Two-stage pipelines are alleviating this
 1. Detect person bounding box at coarse resolution
 2. Infer skeleton pose within box at high resolution
 - Not end-to-end differentiable
(pose extraction requires arg-max function)
 - No sub-pixel accuracy
 - multi-scale approaches can overcome this at the cost of execution time
(average over runs on re-scaled input)



Integral Regression-based 2D pose estimation I

A combination of classification and regression

1. Detection network to produce heatmaps
 - same CNN as for heatmap prediction
2. Soft-max layer to turn heatmap H into probability map P
 - normalizing all pixels in each heatmap H

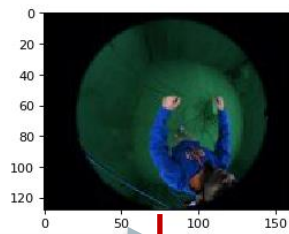
$$P[u, v] = \text{soft-max}(H, (u, v)) = \frac{e^{H[u,v]}}{\sum_{x=1}^{\text{width}} \sum_{y=1}^{\text{height}} e^{H[x,y]}}$$

3. Integration layer to regress joint position (expected position)
 - can be interpreted as voting/weighted average
 - each pixel votes for its own position, weighted by its probability*

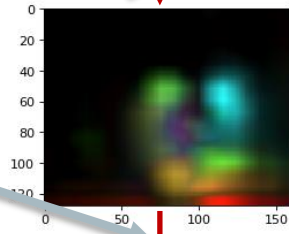
$$\text{pose}_x = \sum_{x=1}^{\text{width}} \sum_{y=1}^{\text{height}} xP[x, y]$$

$$\text{pose}_y = \sum_{x=1}^{\text{width}} \sum_{y=1}^{\text{height}} yP[x, y]$$

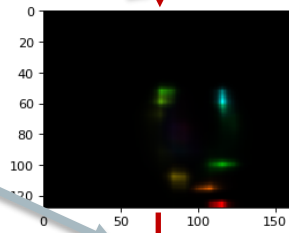
[Sun et al., Integral Human Pose Regression.]



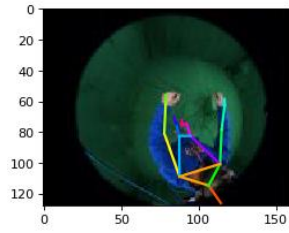
input



heatmap

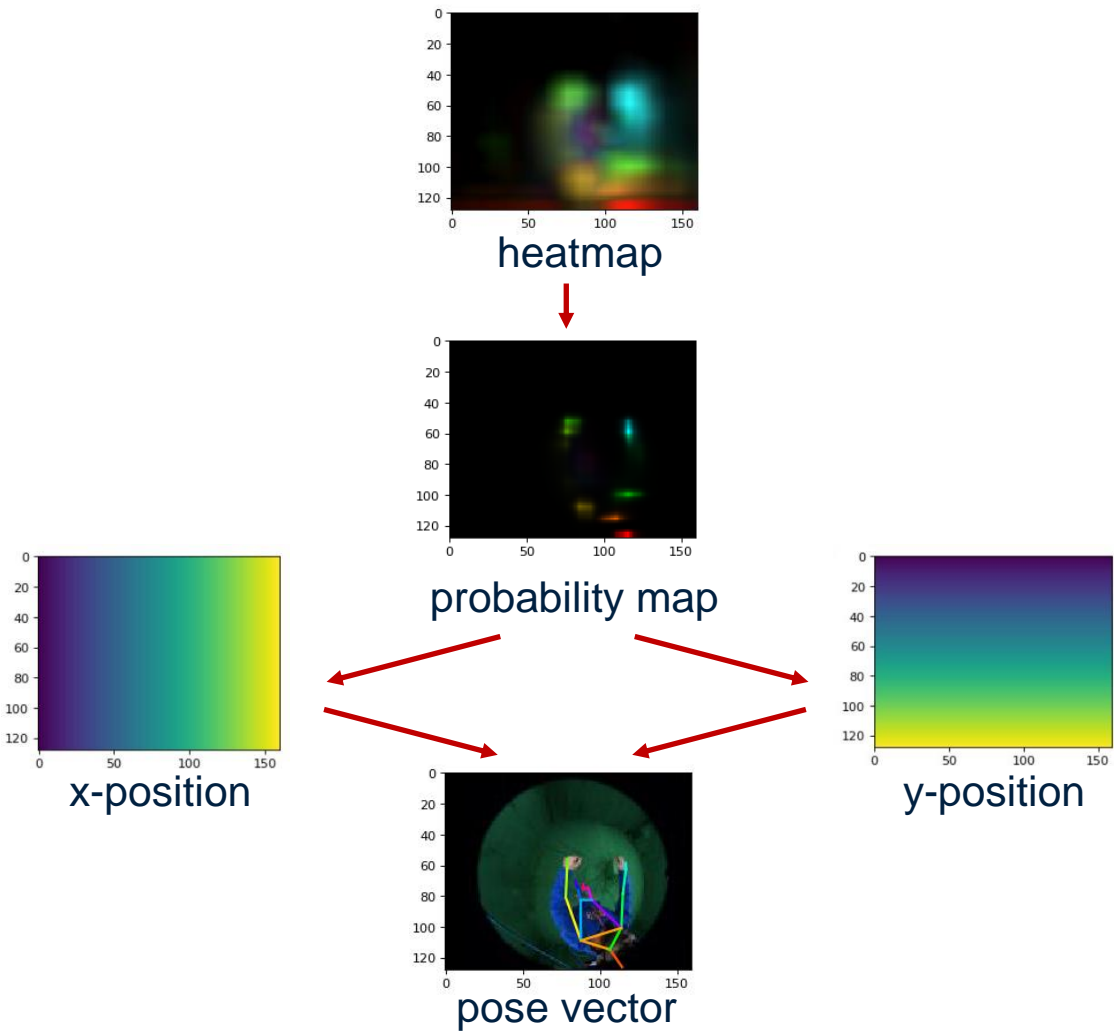


prob. map



pose vector

Details



Integral Regression-based 2D pose estimation II

Advantages

1. Fully-convolutional CNN (as for heatmap classification)
2. Differentiable 2D pose regression
 - soft-max is differentiable, stable, and efficient to compute

$$P[u, v] = \text{soft-max}(H, (u, v)) = \frac{e^{H[u, v]}}{\sum_{x=1}^{\text{width}} \sum_{y=1}^{\text{height}} e^{H[x, y]}}$$

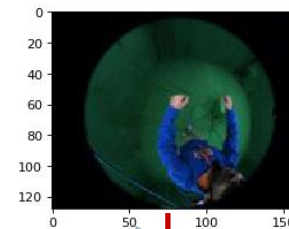
- sum over probability map is differentiable

$$\text{pose}_x = \sum_{x=1}^{\text{width}} \sum_{y=1}^{\text{height}} xP[x, y]$$

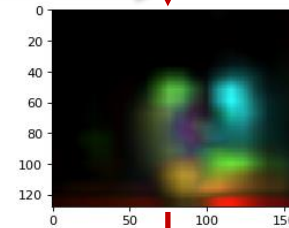
$$\text{pose}_y = \sum_{x=1}^{\text{width}} \sum_{y=1}^{\text{height}} yP[x, y]$$

3. End-to-end training

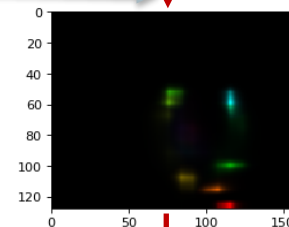
- no difference between training and inference
- sub-pixel accuracy possible through joint influence of pixels
 - low-resolution heatmaps possible



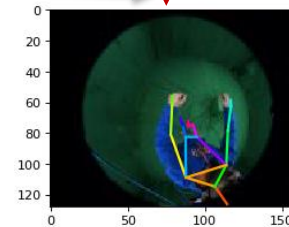
input



heatmap



prob. map



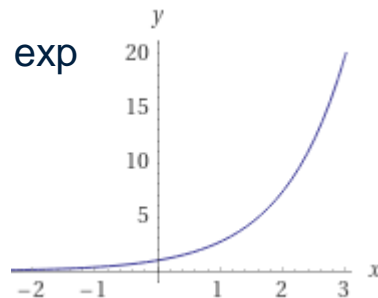
pose vector

Attention: numerical stability

Exp normalize trick within cross-entropy

$$\begin{aligned} \text{soft-max}(z, i) &= \frac{e^{z[i] - \bar{z}} e^{\bar{z}}}{\sum_{j=1}^K e^{z[j] - \bar{z}} e^{\bar{z}}} \\ &= \frac{e^{z[i] - \bar{z}}}{\sum_{j=1}^K e^{z[j] - \bar{z}}} \end{aligned}$$

shift invariance is used to increase numerical stability!



The PyTorch implementation includes this step

Integral Regression-based 2D pose estimation III

Disadvantages / open questions = possible course projects!

1. Sensitive to outliers
 - if there are two maxima in the heatmap, the predicted position will be in the middle of the two
2. How to support multiple people, at different scales?
 - Some form of hierarchical model?
3. Part affinity fields have been successful, can we develop a differentiable model?
 - An elongated ellipse that has position and orientation?
4. What about occluded joints?
5. What about temporal information?
6. Is it possible to infer neck-centered human pose (not knowing the absolute position, only relative distance of keypoints to the neck)?

Issues?

Your laptop / desktop

- No GPU? -> google colab or university (see lecture 2)
- Note, parallel dataloaders might not work well on Windows:
Error: “Can't pickle <function <lambda> ...”
 - fix: disable threading by setting num_workers=0
- Other issues encountered?