

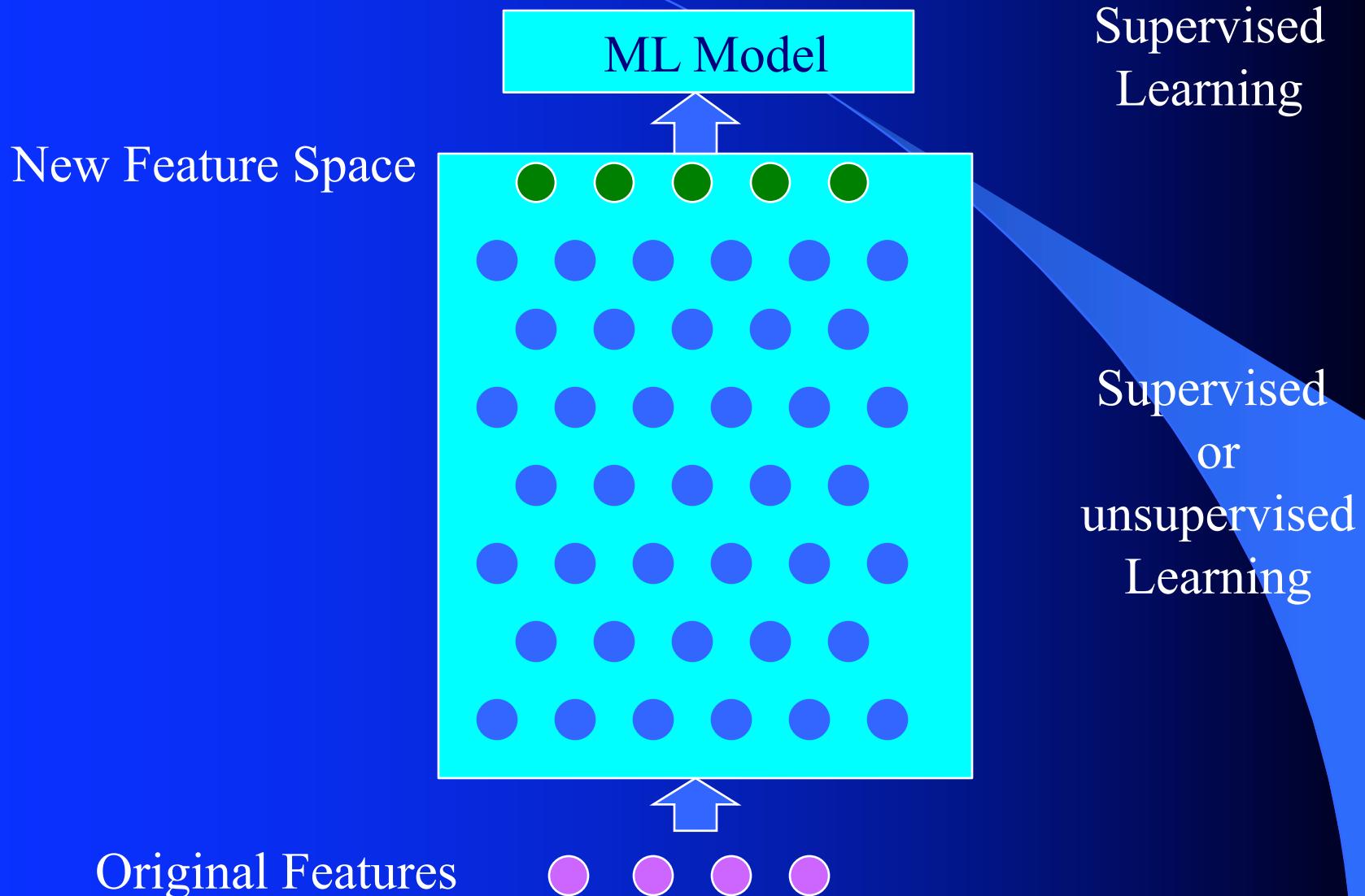
Deep Learning

- Basic Philosophy
- Why Deep Learning
- Deep Backpropagation
- CNN – Convolutional Neural Networks
- Unsupervised Pre-Training Networks
 - Stacked Auto Encoders
 - Deep Belief Networks
- Deep Supervised Networks with managed gradient approaches
 - Learning tricks, Dropout, Batch normalization, ResNets, LSTM, GRUs, etc.
- GANs
- Deep Reinforcement Learning

Deep Learning Overview

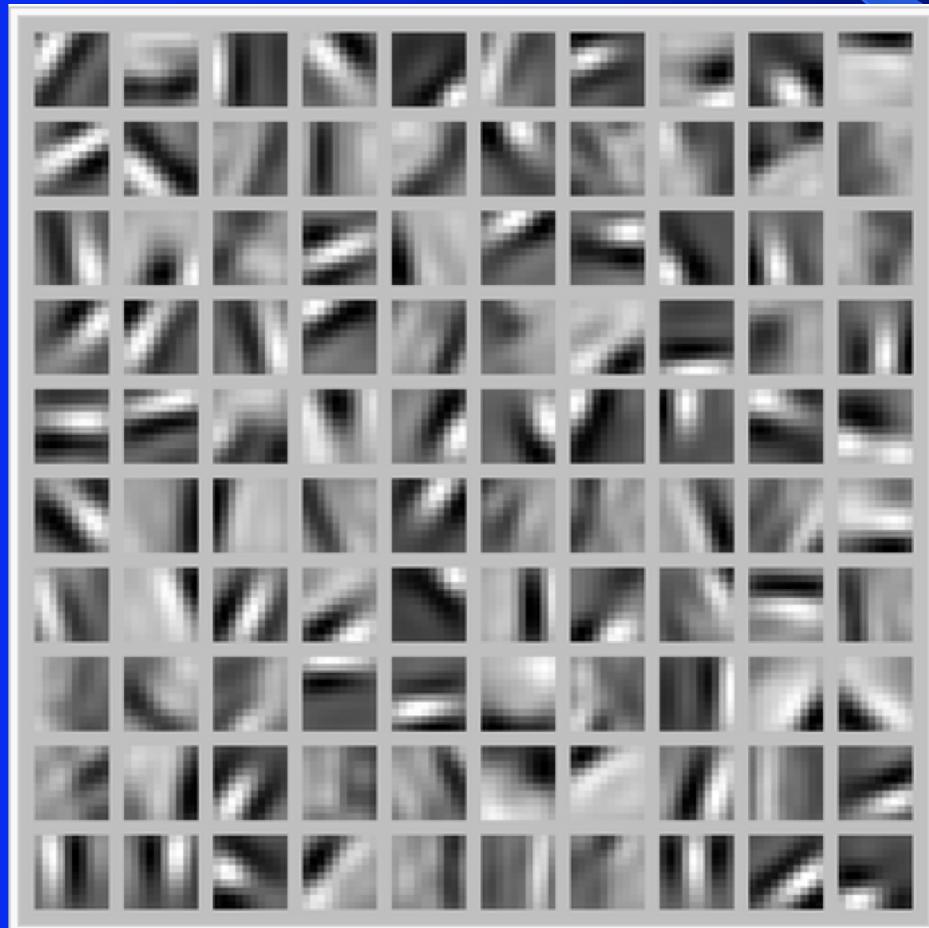
- Train networks with many layers (vs. shallow nets with just a couple of layers)
- Multiple layers work to build an improved feature space
 - First layer learns 1st order features (e.g. edges...)
 - 2nd layer learns higher order features (combinations of first layer features, combinations of edges, etc.)
 - Some models learn in an unsupervised mode and discover general features of the input space – serving multiple tasks related to the unsupervised instances (image recognition, etc.)
 - Final layer of transformed features are fed into supervised layer(s)
 - And entire network is often subsequently tuned using supervised training of the entire net, using the initial weightings learned in the previous phase

Deep Net Feature Transformation



Deep Learning Tasks

- For some approaches (CNN, Unsupervised) best when input space is locally structured – spatial or temporal: images, language, etc. vs arbitrary input features
- Images Example: view of a learned vision feature layer (Basis)
- Each square in the figure shows the input image that maximally activates one of the 100 units



Why Deep Learning

- Biological Plausibility – e.g. Visual Cortex
- Hastad proof - Problems which can be represented with a polynomial number of nodes with k layers, may require an exponential number of nodes with $k-1$ layers (e.g. parity)
- Highly varying functions can be efficiently represented with deep architectures
 - Less weights/parameters to update than a less efficient shallow representation
- Sub-features created in deep architecture can potentially be shared between multiple tasks
 - Type of Transfer/Multi-task learning

Neural Networks Sketch

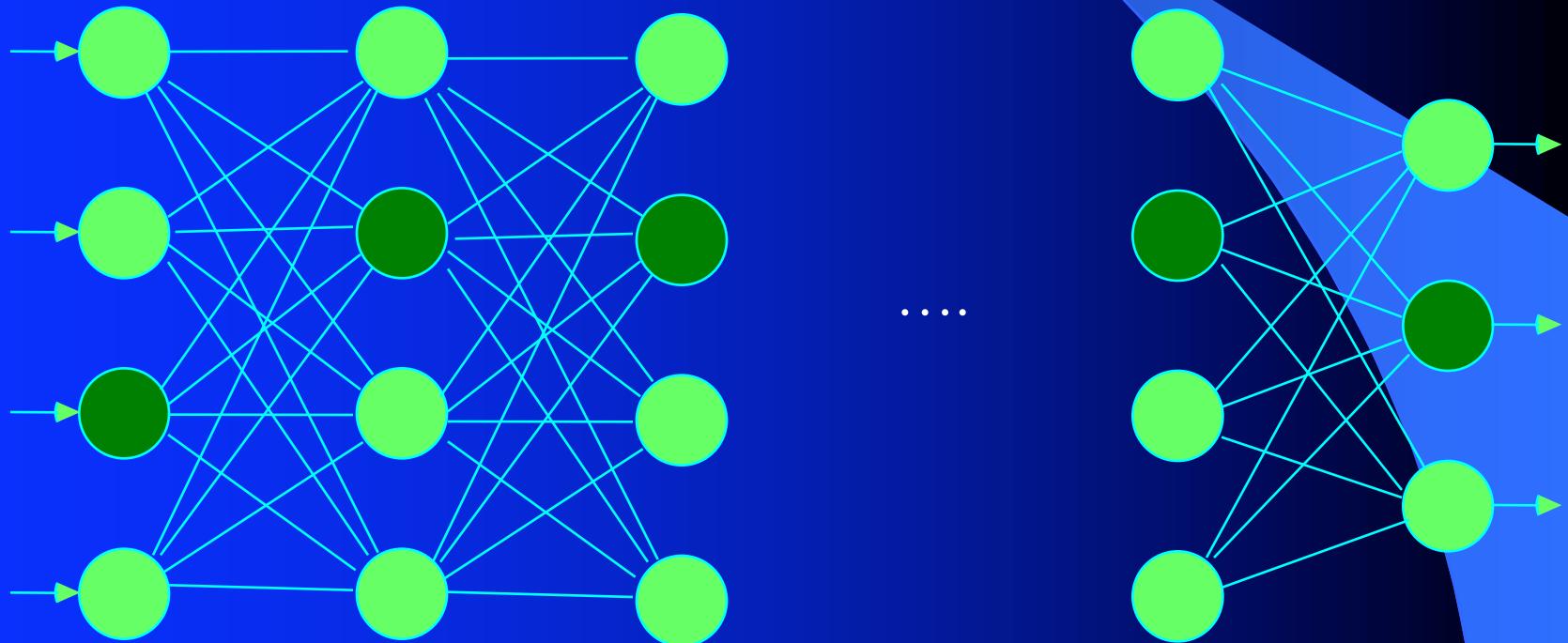
- More than just MLPs, but...
- Rumelhart 1986 – great early success
- Interest subsides a bit in the late 90's as other models are introduced – SVMs, Graphical models, etc. – each "wave" ...
- Convolutional Neural Nets –LeCun 1989-> Image, speech, etc.
- Deep Belief nets (Hinton) and Stacked auto-encoders (Bengio) – 2006 – Unsupervised pre-training followed by supervised. Good feature extractors.
- 2012 -> Initial successes with supervised approaches which overcome vanishing gradient, etc., and are more general applicable. Current explosion, but don't drop all the other tools in your kit!
 - Stay the course

Early Work

- Fukushima (1980) – Neo-Cognitron
- LeCun (1998) – Convolutional Neural Networks (CNN)
 - Similarities to Neo-Cognitron
- Many layered MLP with backpropagation
 - Tried early but without much success
 - Very slow
 - Vanishing gradient
 - More recent work demonstrated significant accuracy improvements by "patiently" training deeper MLPs with BP using fast machines (GPUs)
 - More general learning!
 - Much improved since 2012 with some important extensions to the original MLP/BP approach

Vanishing/Exploding Gradient

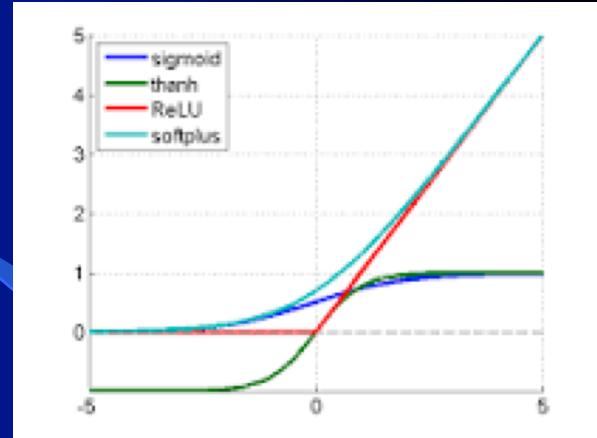
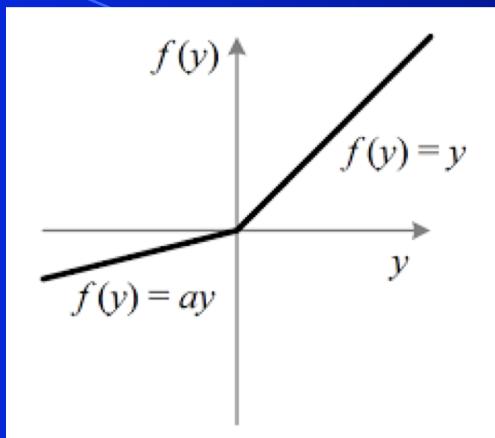
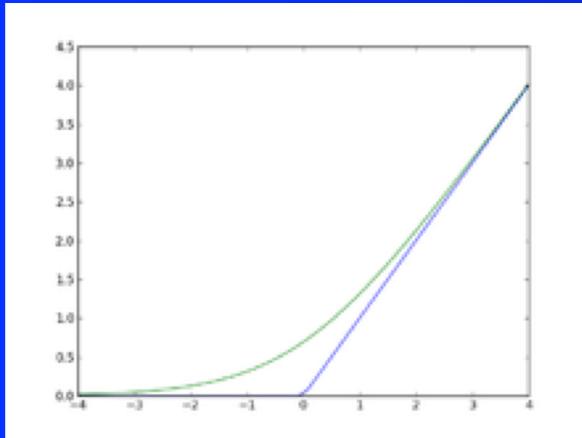
- Error attenuation, long patient training with GPUs, etc
- Recent algorithmic improvements - Rectified Linear Units, better weight initialization, normalization between layers, residual deep learning, etc. – 1000's of layers being effectively trained



Unstable Gradient

- Early layers of MLP do not get trained well
 - Vanishing Gradient – error attenuates as it propagates to earlier layers – $t - z < 1$, $f'(net)$, scaled by small initial weights
 - Leads to very slow training (especially at early layers)
 - Exacerbated since top couple layers can usually learn any task "pretty well" and thus the error to earlier layers drops quickly as the top layers "mostly" solve the task—lower layers never get the opportunity to use their capacity to improve results, they just do a random feature mapping
 - Need a way for early layers to do effective work
 - Instability of gradient in deep networks: Vanishing or exploding gradient
 - Product of many terms, which unless “balanced” just right, is unstable
 - Either early or late layers stuck while “opposite” layers are learning

Rectified Linear Units



- $f(x) = \text{Max}(0, x)$ More efficient gradient propagation, derivative is 0 or constant, just fold into learning rate
 - Helps $f'(\text{net})$ issue, but still left with other unstable gradient issues
- More efficient computation: Only comparison, addition and multiplication.
 - Leaky ReLU $f(x) = x$ if $x > 0$ else ax , where $0 \leq a \leq 1$, so that derivate is not 0 and can do some learning for net < 0 (does not “die”).
 - Lots of other variations
- Sparse activation: For example, in a randomly initialized network, only about 50% of hidden units are activated (having a non-zero output)
- Learning in linear range easier for most learning models

Softmax and Cross-Entropy

- Sum-squared error (L2) loss gradient seeks the maximum likelihood hypothesis under the assumption that the training data can be modeled by Normally distributed noise added to the target function value. Fine for regression but less natural for classification.
- For *classification* problems it is advantageous and increasingly popular to use the *softmax* activation function, just at the output layer, with the cross-entropy loss function. Softmax (softens) 1 of n targets to mimic a probability vector for each output.

$$f(\text{net}_j) = \frac{e^{\text{net}_j}}{\sum_{i=1}^n e^{\text{net}_i}}$$

- Cross entropy seeks to find the maximum likelihood hypotheses under the assumption that the observed (1 of n) Boolean outputs is a probabilistic function of the input instance. Maximizing likelihood is cast as the equivalent minimizing of the negative log likelihood.

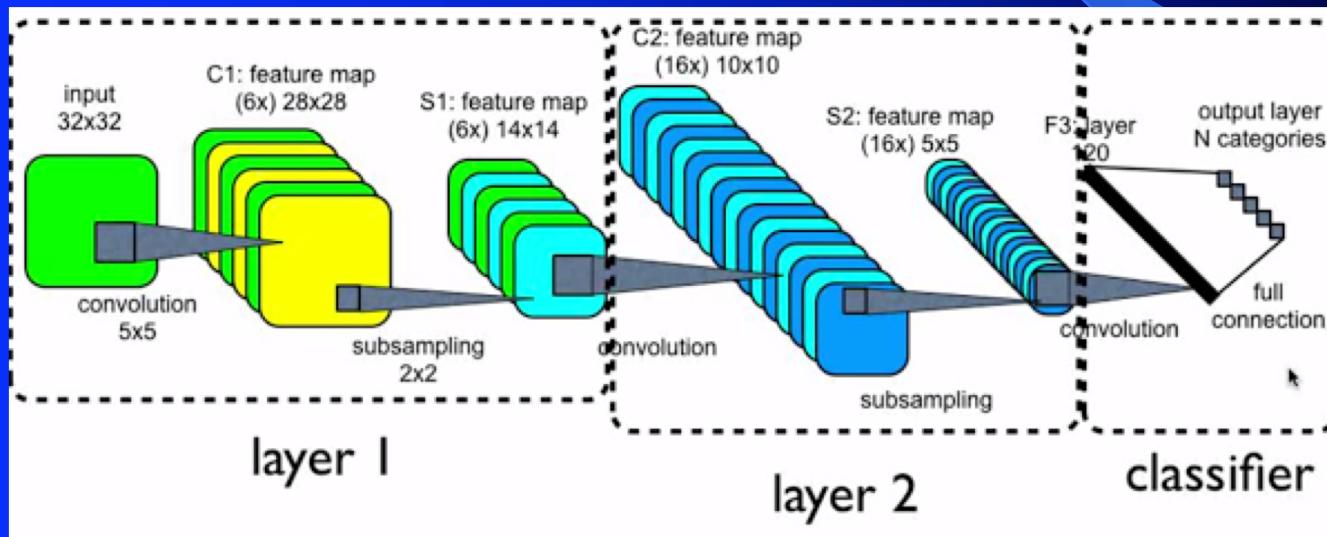
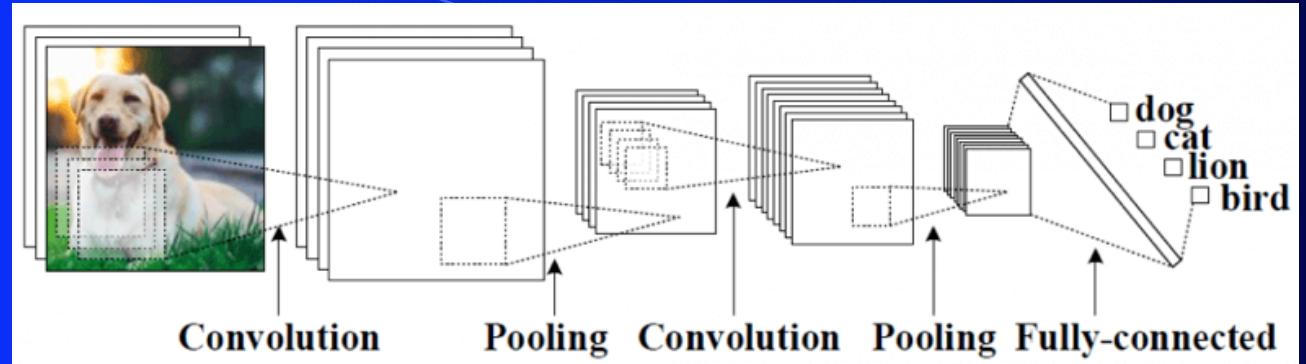
$$\text{Loss}_{\text{CrossEntropy}} = - \sum_{i=1}^n t_i \ln(z_i)$$

- With new loss and activation functions, we must recalculate the gradient equation. Gradient/Error on the output is just $(t-z)$, no $f'(\text{net})$! The exponent of softmax is unraveled by the \ln of cross entropy. Helps avoid gradient saturation.
- Common with deep networks, with ReLU activations for the hidden nodes

Convolutional Neural Networks

- "Niche" networks built specifically for problems with low dimensional (e.g. 2-d) grid-like local structure – e.g. Images
 - Character recognition, Speech, Games, Images - where neighboring pixels have high correlations and local features (edges, corners, etc.), while distant pixels (features) are less correlated
 - Typically just uses raw features (pixels) with no preprocessing
 - Natural images have the property of being stationary, meaning that the statistics of one part of the image are the same as any other part
 - 1-d example
 - Some biological plausibility from visual cortex
 - While standard NN nodes take input from all nodes in the previous layer, CNNs enforce that a node receives only a small set of features which are spatially or temporally close to each other called *receptive fields* from one layer to the next (e.g. 3x3, 5x5), thus enabling ability to handle local 2-D structure.
 - Can find edges, corners, endpoints, etc.
 - Good for problems with local 2-D structure, but lousy for general learning with abstract features having no prescribed feature ordering or locality

Convolutional Neural Networks



- Big Picture: Each feature map learns a different feature. Each node in feature map has same translated receptive field (and weights)
- Brute force search to see if and where certain features exist
- Pooling/sampling – Does the feature exist in a general area
- Final standard supervised layer with improved feature space

Convolutions

- Typical MLPs have a connection from every node in the previous layer, and the net value for a node is the scalar dot product of the inputs and weights (e.g. matrix multiply). Convolutional nets are somewhat different:
 - Nodes still do a scalar dot product (convolution) from the previous layer, but with only a small portion (receptive field) of the nodes in the previous layer – *Sparse representation*
 - Every node has the exact same weight values from the preceding layer – *Shared parameters*, tied weights, a LOT less unique weight values. Regularization by having same weights looking at lots of input situations
 - Each node has it's shared weight convolution computed on a receptive field slightly shifted, from that of it's neighbor, in the previous layer – *Translation invariance*.
 - Each node's convolution scalar is then passed through a non-linear activation function (ReLU, tanh, etc.)

Convolution Example

| | | | | |
|------------------------|------------------------|------------------------|---|---|
| 1 <small>×1</small> | 1 <small>×0</small> | 1 <small>×1</small> | 0 | 0 |
| 0 <small>×0</small> | 1 <small>×1</small> | 1 <small>×0</small> | 1 | 0 |
| 0 <small>×1</small> | 0 <small>×0</small> | 1 <small>×1</small> | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

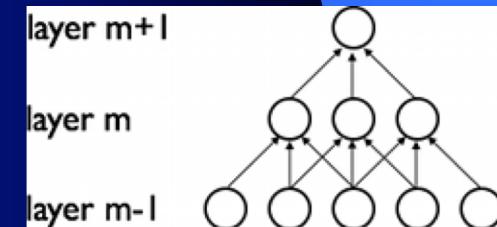
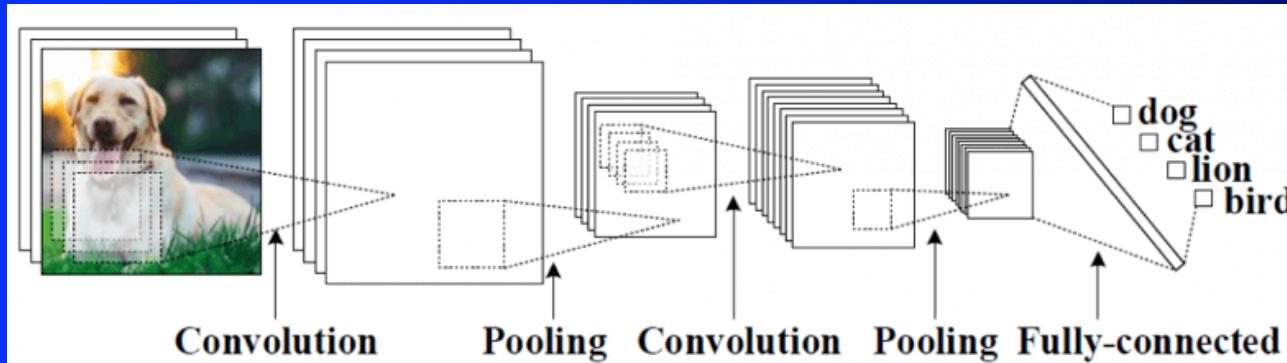
Image

| | | |
|---|--|--|
| 4 | | |
| | | |
| | | |
| | | |

Convolved
Feature

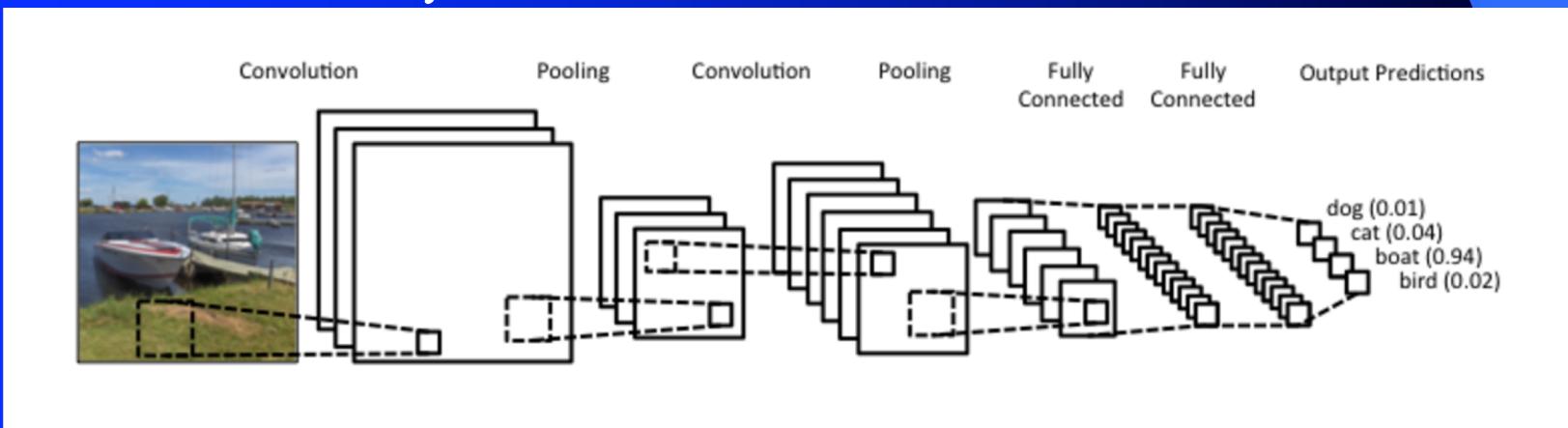
CNN – Translation Invariance

- The 2-d planes of nodes (or their outputs) at subsequent layers in a CNN are called *feature maps*
- To deal with translation invariance, each node in a feature map has the same weights (based on the feature it is looking for), and each node connects to a different overlapping receptive field of the previous layer
- Thus each *feature map* searches the full previous layer to see if, where, and how often its feature occurs (precise position less critical)
 - The output will be high at each node in the map corresponding to a receptive field where the feature occurs
 - Later layers can concern themselves with higher order combinations of features and rough relative positions
 - Each calculation of a node's net value, $\sum xw + b$ in the feature map, is called a convolution, based on the similarity to standard convolutions



CNN Structure

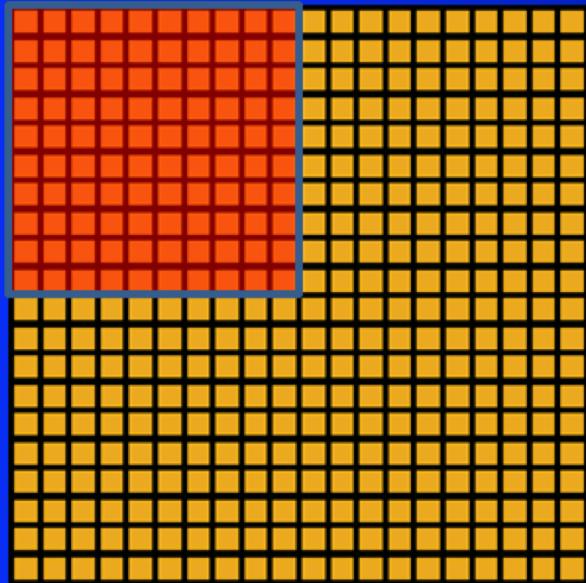
- Each node (e.g. convolution) is calculated for each receptive field in the previous layer
 - During training the corresponding weights are always tied to be the same (ala BPTT)
 - Thus a relatively small number of unique weight parameters to learn, although they are replicated many times in the feature map
 - Each node output in CNN is $f(\sum xw + b)$ (ReLU, tanh etc.)
 - Multiple feature maps in each layer
 - Each feature map should learn a different translation invariant feature
 - Since after first layer, there are multiple feature maps to connect to the next layer, current convolution map usually takes inputs from all feature maps in previous layer (e.g. S2 in figure below)
 - Convolution layer causes total number of features to increase



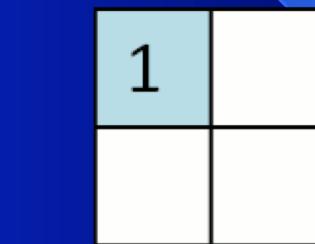
Sub-Sampling (Pooling)

- Convolution and sub-sampling layers are interleaved
- Sub-sampling (Pooling) allows number of features to be diminished, and to pool information
 - Pooling replaces the network output at a certain point with a summary statistic of nearby outputs
 - Max-Pooling common (Just as long as the feature is there, take the max, as exact position is not that critical), also averaging, etc.
 - Pooling smooths the data and reduces spatial resolution and thus naturally decreases importance of exactly where a feature was found, just keeping the rough location – translation invariance
 - 2x2 pooling would do 4:1 compression, 3x3 9:1, etc.
 - Convolution usually increases number of feature maps, pooling keeps same number of *reduced* maps (one-to-one correspondence of convolution map to pooled map) as the previous layer

Pooling Example (Summing or averaging)



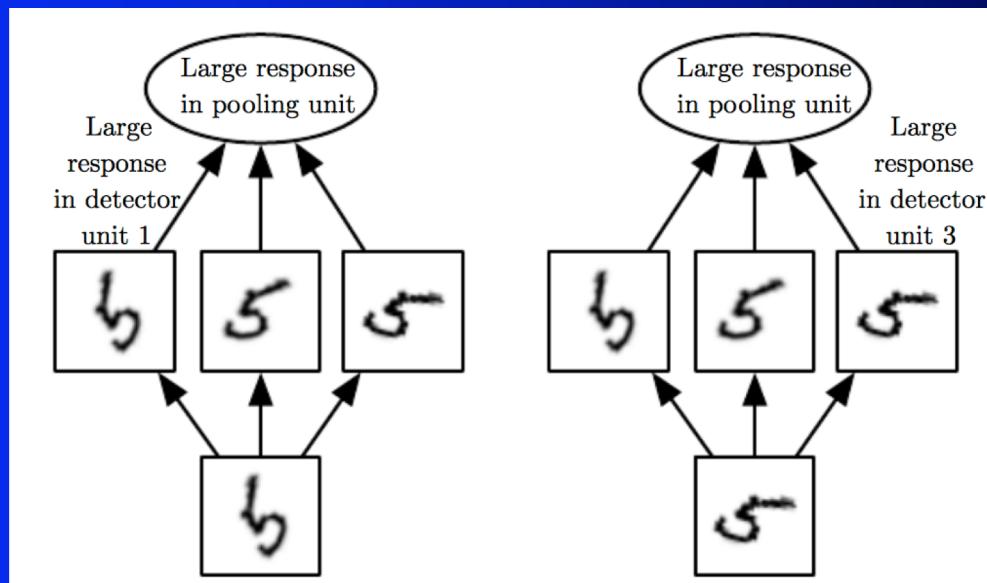
Convolved
feature



Pooled
feature

Pooling (cont.)

- Common layers are convolution, non-linearity, then pool (repeat)
- Note that pooling always decreases map volumes (unless pool stride = 1, highly overlapped), making real deep nets more difficult. Pooling is sometimes used only after multiple convolved layers and sometimes not at all.
- At later layers pooling can make network invariant to more than just translation – *learned invariances*



CNN Training

- Trained with BP but with weight tying in each feature map
 - Randomized initial weights throughout entire network
 - Average the weight updates over the tied weights in feature map layers
- Convolution layer
 - Each feature map has one weight for each input and one bias
 - Thus a feature map with a 5x5 receptive field (filter) would have a total of 26 weights, which are the same coming into each node of the feature map
 - If a convolution layer had 10 feature maps, then only a total of 260 unique weights to be trained in that layer (much less than an arbitrary deep net layer without sharing)
- Sub-Sampling (Pooling) Layer
 - All elements of receptive field max'd, averaged, summed, etc. No trainable weights necessary
 - While all weights are trained, the structure of the CNN is currently usually hand crafted with trial and error
 - Number of total layers, number of receptive fields, size of receptive fields, size of sub-sampling (pooling) fields, which fields of the previous layer to connect to
 - Often decrease size of feature maps and increase number of feature maps for subsequent layers

CNN Hyperparameters

- Structure itself, number of layers, size of filters, number of feature maps in convolution layers, connectivity between layers, activation functions, final supervised layers, etc.
- Drop-out often used in final fully connected layers for overfit avoidance – less critical in convolution/pooling layers which already regularize due to weight sharing
- *Stride* – Don't have to test every location for the feature (i.e. $\text{stride} = 1$), could sample more coarsely
 - Another option (besides pooling) for down-sampling
- As is, the feature map would always decrease in volume which is not always desirable - *Zero-padding* avoids this and lets us maintain up to the same volume
 - Would shrink fast for large kernel/filter sizes and would limit the depth (number of layers) in the network
 - Also allows the different filter sizes to fit arbitrary map widths

Convolution Example

| | | | | |
|------------------------|------------------------|------------------------|---|---|
| 1 <small>x1</small> | 1 <small>x0</small> | 1 <small>x1</small> | 0 | 0 |
| 0 <small>x0</small> | 1 <small>x1</small> | 1 <small>x0</small> | 1 | 0 |
| 0 <small>x1</small> | 0 <small>x0</small> | 1 <small>x1</small> | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

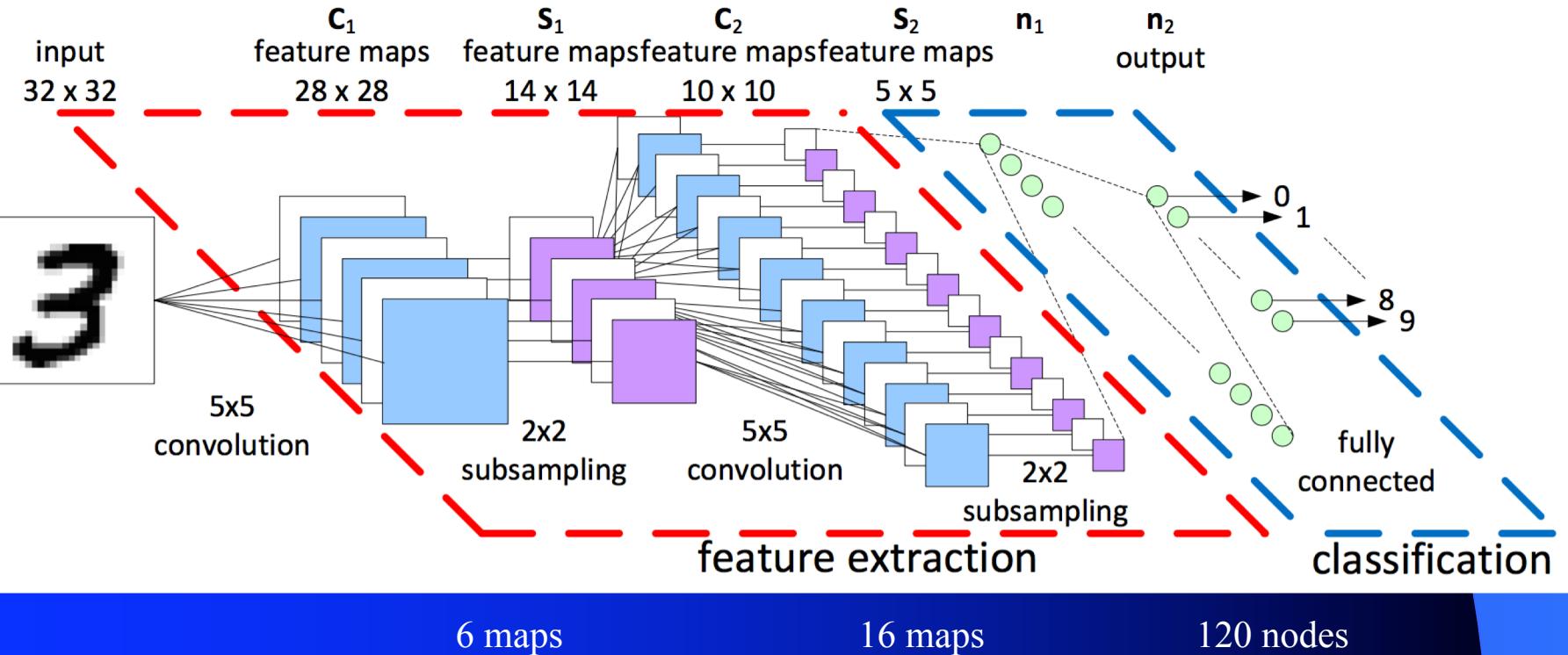
Image

| | | |
|---|--|--|
| 4 | | |
| | | |
| | | |
| | | |

Convolved
Feature

Example – CNN MNIST Classification

- Roughly based on LeCun's original model. To help it all sink in:
- How many connections and trainable weights at each layer?
 - MNIST data input is 28x28 pixels



Basic CNN Example

| Layer | Trainable Weights | Connections |
|--------|-------------------------------|----------------------------------------------|
| C1 | $(25+1)*6 = 156$ | $(25+1)*6*28*28 = 122,304$ |
| S1 | 0 (LeCun had $(1+1)*6 = 12$) | $4 \text{ (2x2 links)} * 6 * 14 * 14 = 4704$ |
| C2 | $16*(5*5*6+1) = 2416$ | $2416 * 10 * 10 = 241,600$ |
| S2 | 0 | $4 \text{ (2x2 links)} * 16 * 5 * 5 = 1600$ |
| N1 | $120*(5*5*16+1) = 48,120$ | Same since fully connected MLP at this point |
| Output | $10*(120+1) = 1210$ | Same |

- Why 32x32 to start with? Actual characters never bigger than 28x28. Just padding the edges so for example the top corner node of the feature map can have a pad of two up and left for its feature map (since receptive field is 5x5). Same things happens with 14x14 to 10x10 drop from S1 to C2.
- S1 and S2 non-overlapping and pool (max most common) – We include here the 4 unweighted connections
 - LeCun had a trainable weight and a bias in pool layer followed by a non-linearity, not really necessary and not used these days
- C2: Connects to all preceding maps.
 - LeCun had each map connect to a subset of the preceding maps
- S2 Final number of extracted features to go to the MLP: $(5*5)*16 = 400$
- 419,538 total connections, with 51,902 trainable parameters 95% of which are in the final MLP. Only 2572 trainable weights in CNN

ILSVRC Image net Large Scale Vision Recognition Competition

RGB: $224 \times 224 \times 3 = 150,528$ raw real valued features

- Annual competition of image classification at large scale
- 1.2M images in 1K categories
- Classification: make 5 guesses about the image label



EntleBucher

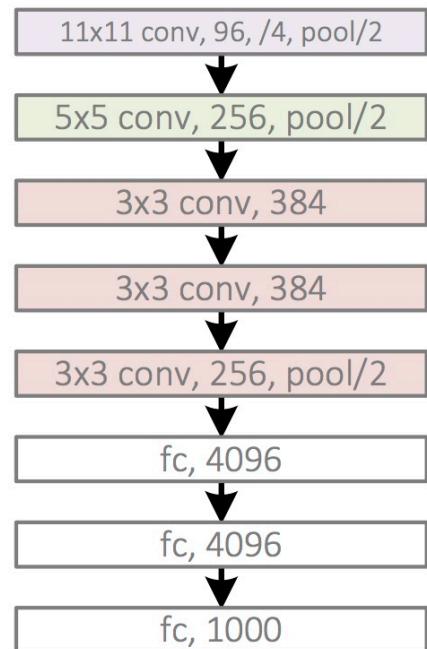


Appenzeller

Example CNNs Structures ILSVRC winners

Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)



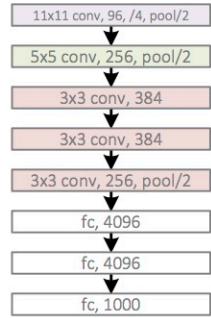
Kaiming He, Xiangyu Zhang, Shao

- Note Pooling considered part of the layer
- 96 convolution kernels (maps), then 256, then 384
- Stride of 4 for first convolution kernel, 1 for the rest
- Pooling layers with 3x3 receptive fields and stride of 2 throughout
- Finishes with fully connected (fc) MLP with 2 hidden layers and 1000 output nodes for classes

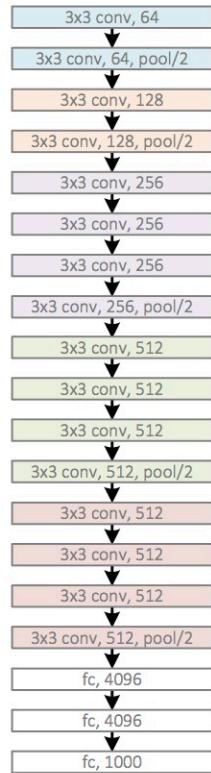
Example CNNs Structures ILSVRC winners

Revolution of Depth

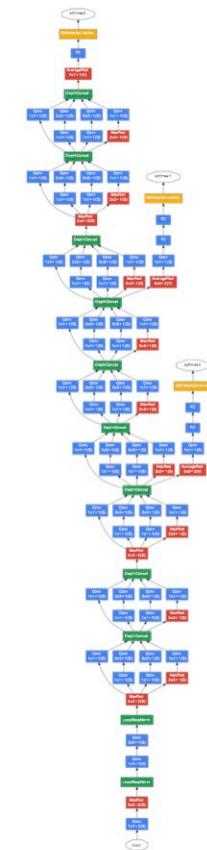
AlexNet, 8 layers
(ILSVRC 2012)



VGG, 19 layers
(ILSVRC 2014)



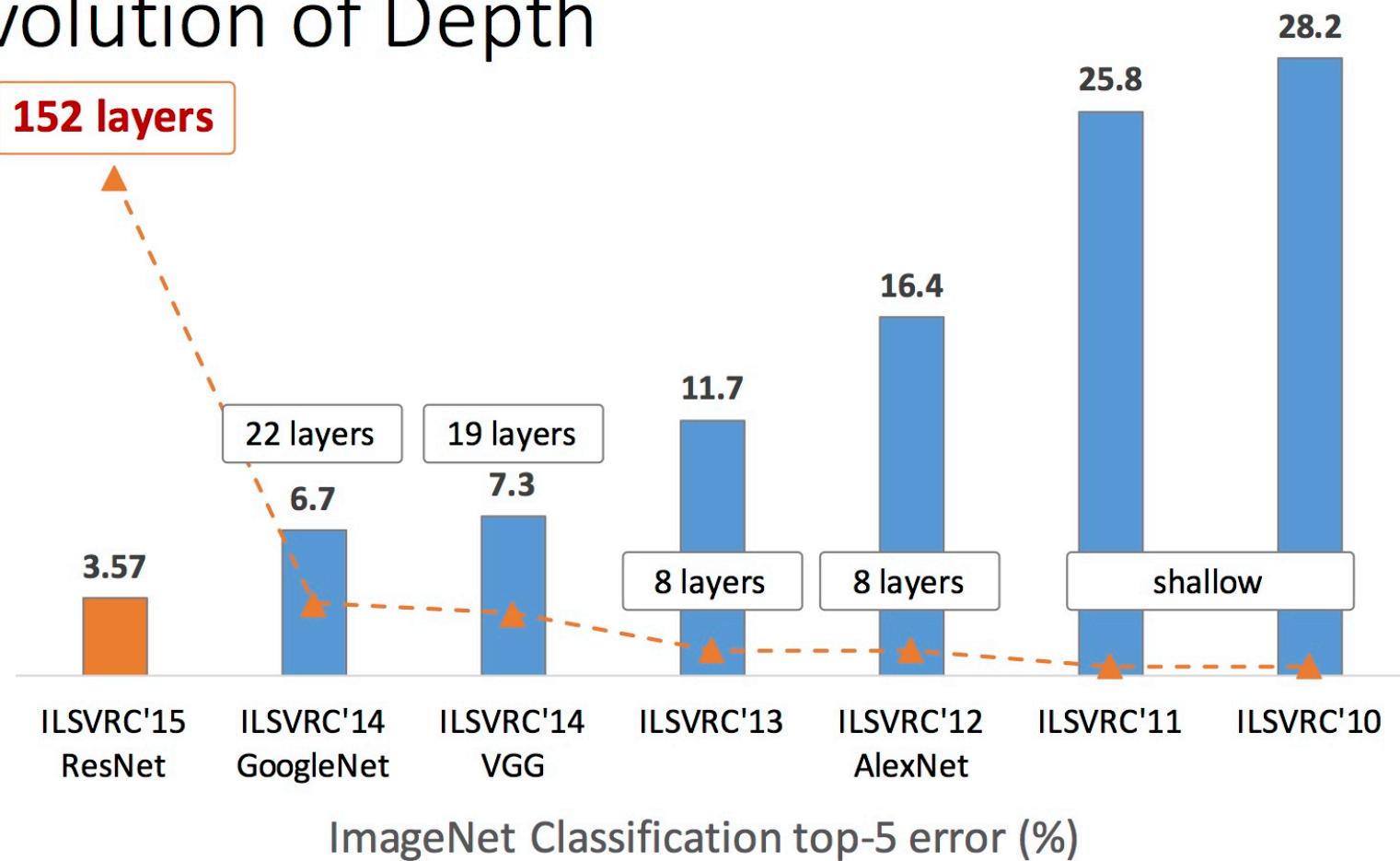
GoogleNet, 22 layers
(ILSVRC 2014)



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

Increasing Depth

Revolution of Depth



CNN Summary

- High accuracy for image applications – Breaking all records and doing it using just raw pixel features!
- Special purpose net – Just for images or problems with strong grid-like local spatial/temporal correlation (Speech, games, etc.)
- Once trained on one problem (e.g. vision) could use same net (often tuned) for a new similar problem – general creator of vision features
- Unlike traditional nets, handles variable sized inputs
 - Same filters and weights, just convolve across different sized image and dynamically scale size of pooling regions (not # of nodes), to normalize
 - Different sized images, different length speech segments, etc.
- Lots of hand crafting and CV tuning to find the right recipe of receptive fields, layer interconnections, etc.
 - Lots more Hyperparameters than standard nets, and even than other deep networks, since the structures of CNNs are more handcrafted
 - CNNs getting wider and deeper with speed-up techniques (e.g. GPU, ReLU, etc.) and lots of current research, excitement, and success

Unsupervised Pre-Training

- Began the hype of Deep-Learning (2006)
 - Before CNNs were fully recognized for what they could do
 - Less popular at the moment for supervised learning, with recent supervised success, but still at the core of many new research directions including deep generative networks
- Unsupervised Pre-Training uses unsupervised learning in the deep layers to transform the inputs into features that are easier to learn by a final supervised model
 - Unsupervised training between layers can decompose the problem into distributed sub-problems (with higher levels of abstraction) to be further decomposed at subsequent layers
- Often not a lot of labeled data available while there may be lots of unlabeled data. Unsupervised Pre-Training can take advantage of unlabeled data. Can be a huge issue for some tasks.
 - Unsupervised and Semi-Supervised
 - Self-Taught Learning

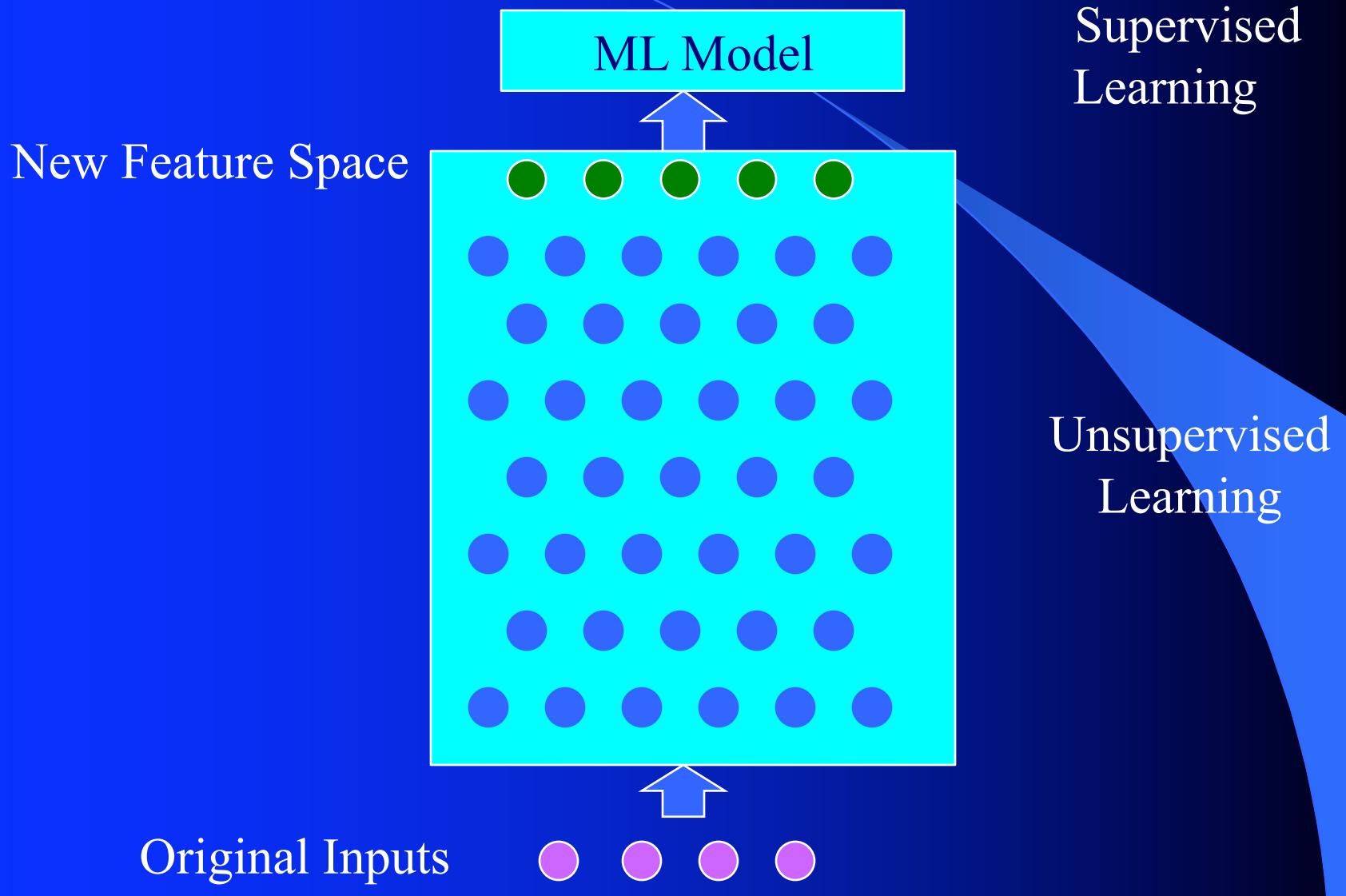
Self Taught vs Unsupervised Learning

- When using Unsupervised Learning as a pre-processor to supervised learning you are typically given examples from the same distribution as the later supervised instances will come from
 - Assume the distribution comes from a set containing just examples from a defined set up possible output classes, but the label is not available (e.g. images of car vs trains vs motorcycles)
- In Self-Taught Learning we do not require that the later supervised instances come from the same distribution
 - e.g., Do self-taught learning with any images, even though later you will do supervised learning with just cars, trains and motorcycles.
 - These types of distributions are more readily available than ones which just have the classes of interest (i.e. not labeled as car *or* train *or* motorcycle)
 - However, if distributions are *very* different...
- New tasks share concepts/features from existing data and statistical regularities in the input distribution that many tasks can benefit from
 - Can re-use well-trained nets as starting points for other tasks
 - Note similarities to supervised multi-task and transfer learning
- Both unsupervised and self-taught approaches reasonable in deep learning models

Greedy Layer-Wise Training

1. Train first layer using your data without the labels (unsupervised)
 - Since there are no targets at this level, labels don't help. Could also use the more abundant unlabeled data which is not part of the training set (i.e. unsupervised and/or self-taught learning).
2. Then freeze the first layer parameters and start training the second layer using the output of the first layer as the unsupervised input to the second layer
3. Repeat this for as many layers as desired
 - This builds the set of robust features
4. Use the outputs of the final layer as inputs to a supervised layer/model and train the last supervised layer(s) (leave early weights frozen)
5. Unfreeze all weights and fine tune the full network by training with a supervised approach, given the *pre-training* weight settings

Deep Net with Greedy Layer Wise Training

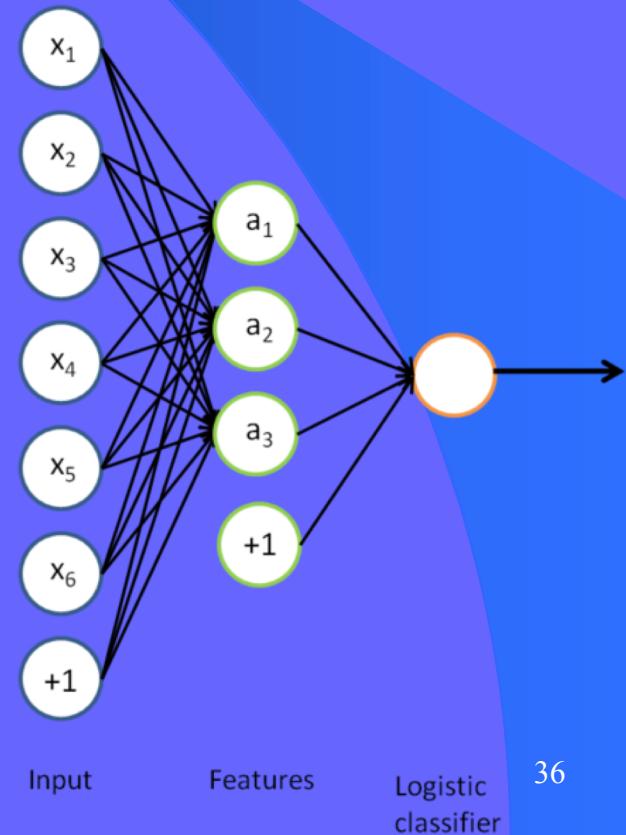
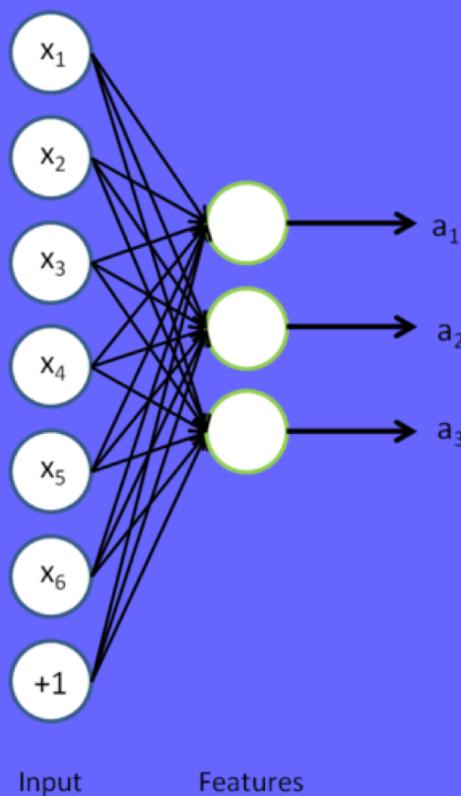
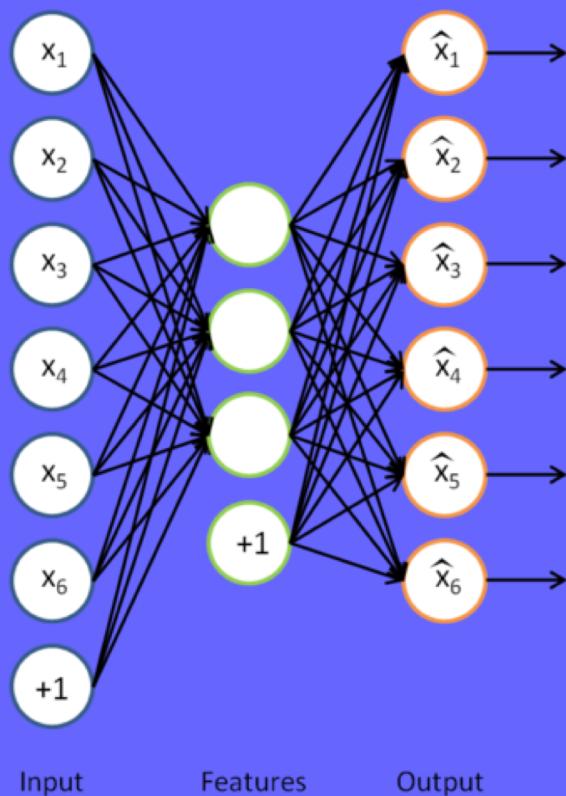


Greedy Layer-Wise Training

- Greedy layer-wise training avoids many of the problems of trying to train a deep net in a supervised fashion
 - Each layer gets full learning focus in its turn since it is the only current "top" layer (no unstable gradient issues, etc.)
 - Can take advantage of unlabeled data
 - When you finally tune the entire network with supervised training the network weights have already been adjusted so that you are in a good error basin and just need fine tuning. This helps with problems of
 - Ineffective early layer learning
 - Deep network local minima
- We will discuss the two early landmark approaches
 - Deep Belief Networks
 - Stacked Auto-Encoders

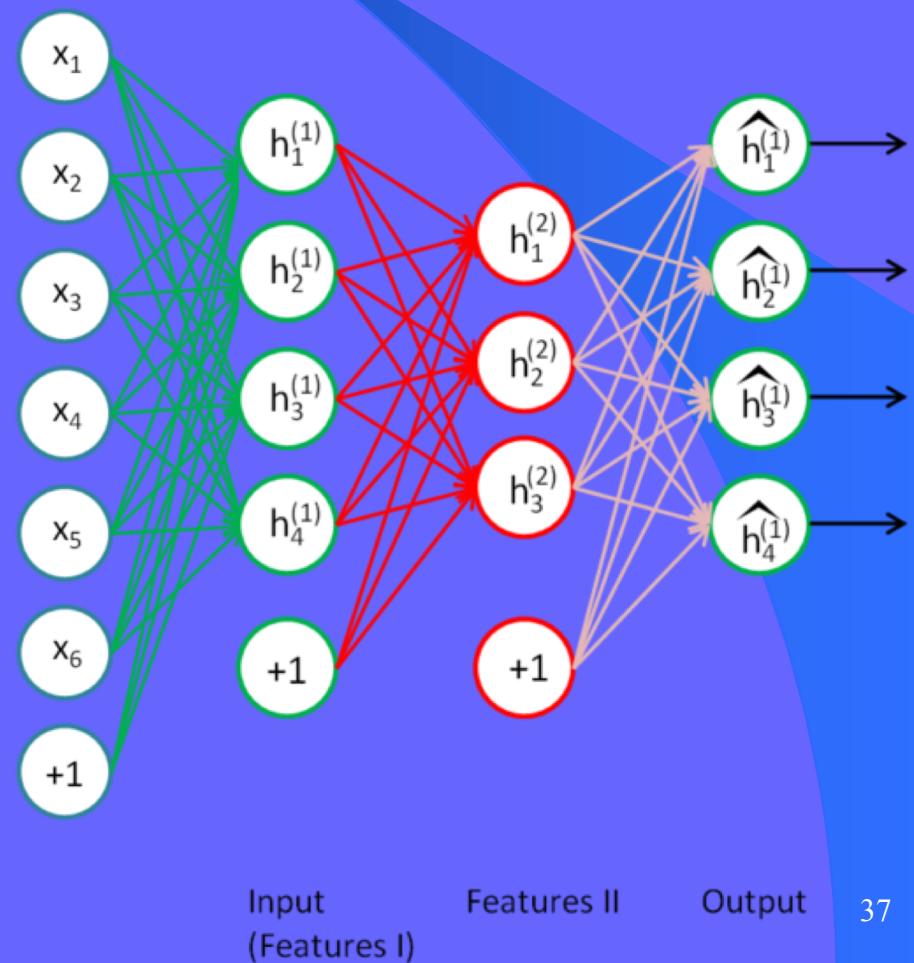
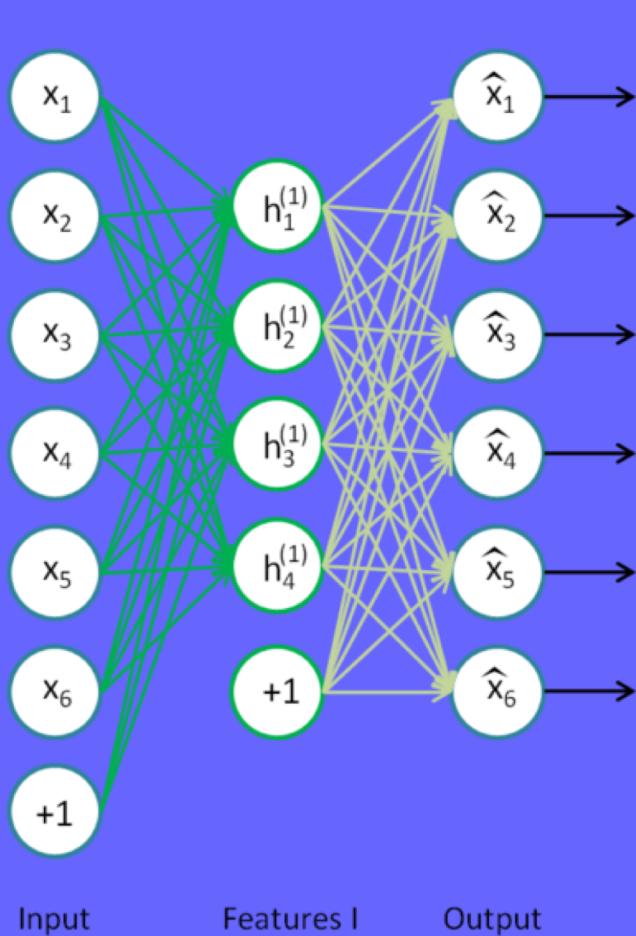
Auto-Encoders

- A type of unsupervised learning which discovers generic features of the data
 - Learn identity function by learning important sub-features
 - Compression, etc. – Undercomplete $|\mathbf{h}| < |\mathbf{x}|$
 - For $|\mathbf{h}| \geq |\mathbf{x}|$ (Overcomplete case more common in deep nets) use regularized autoencoding: Loss function includes regularizer to make sure we don't just pass through the data (e.g. sparsity, noise robustness, etc.)



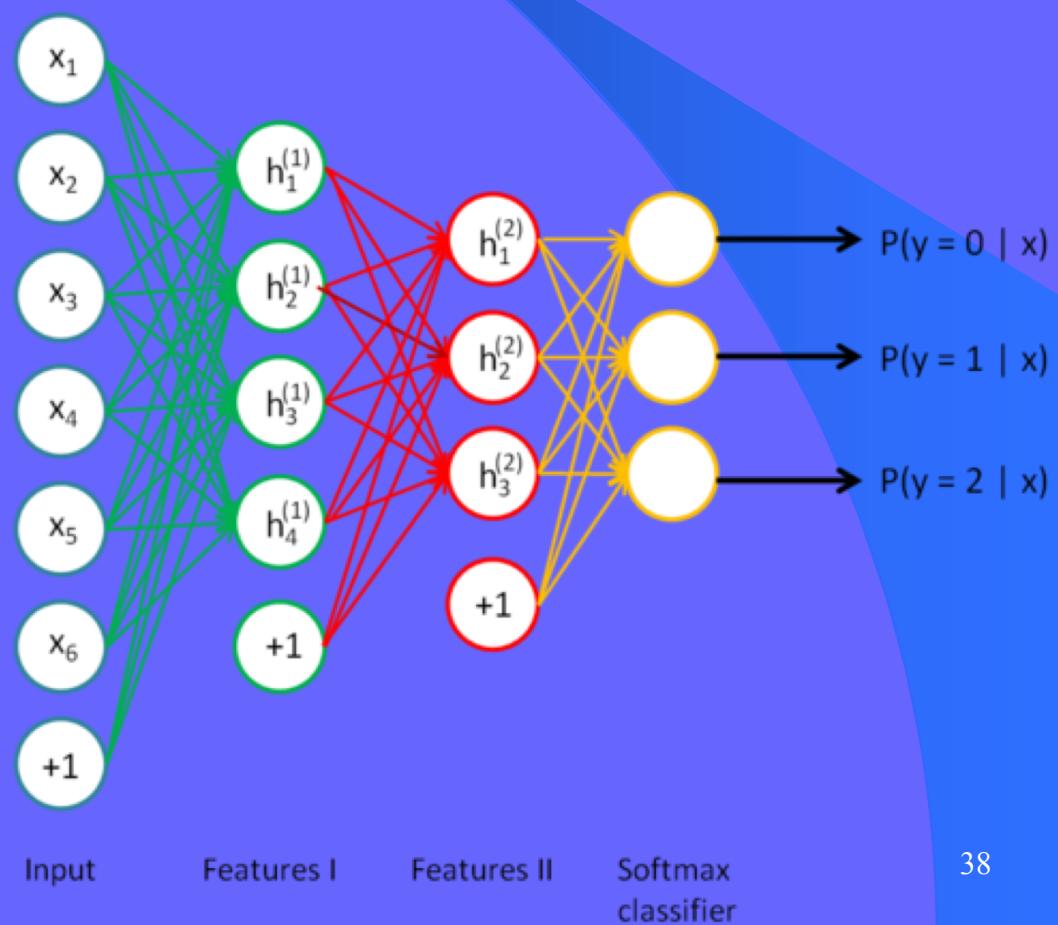
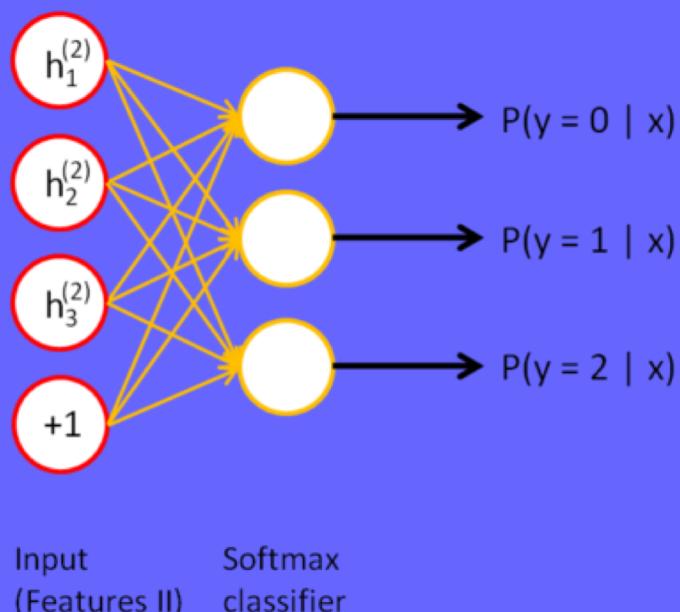
Stacked Auto-Encoders

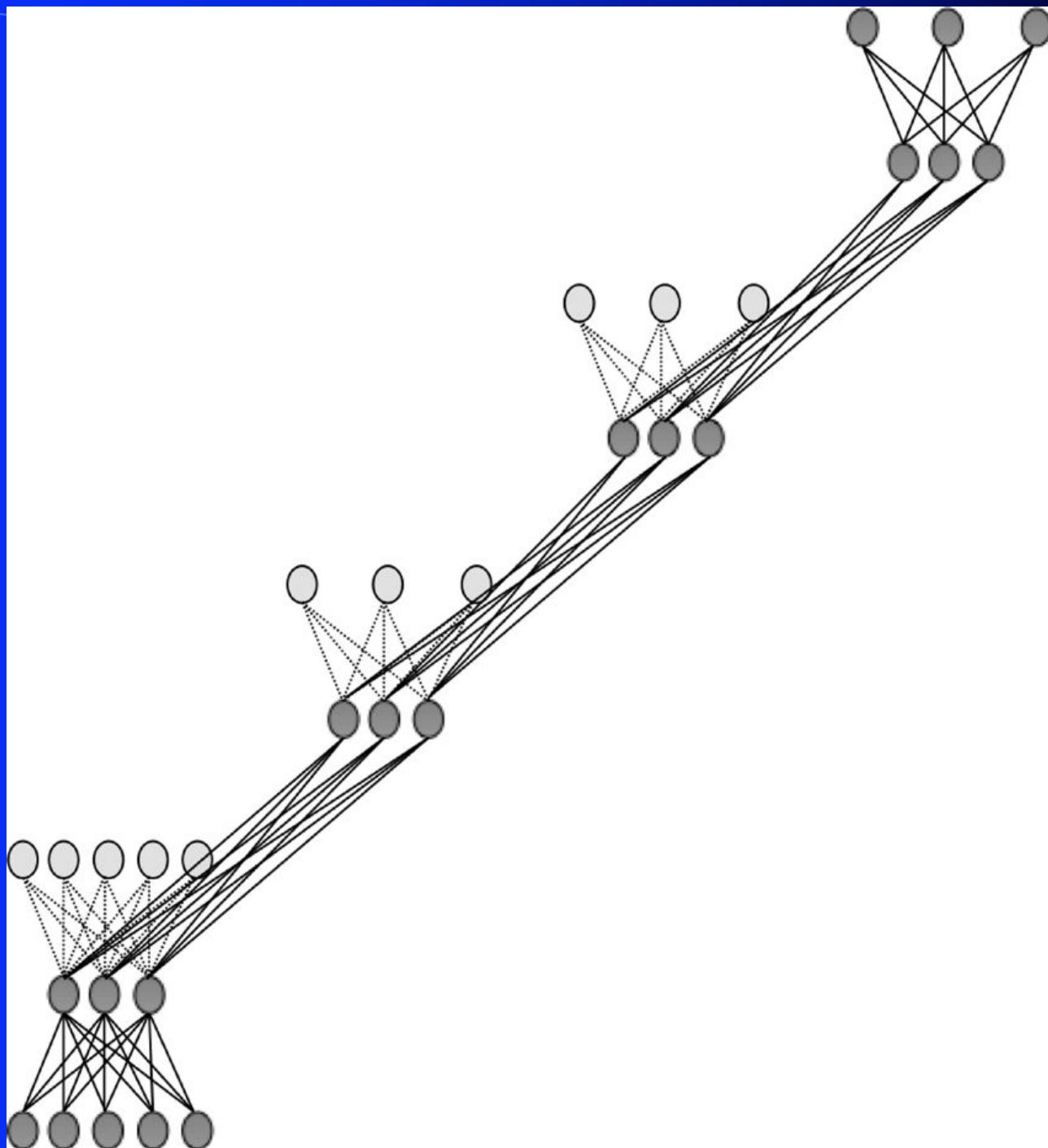
- Bengio (2007) – After Deep Belief Networks (2006)
- Stack many (sparse) auto-encoders in succession and train them using greedy layer-wise training
- Drop the decode output layer each time



Stacked Auto-Encoders

- Do supervised training (can now only used labeled examples) on the last layer using final features
- Then do supervised training on the entire network to fine-tune all weights



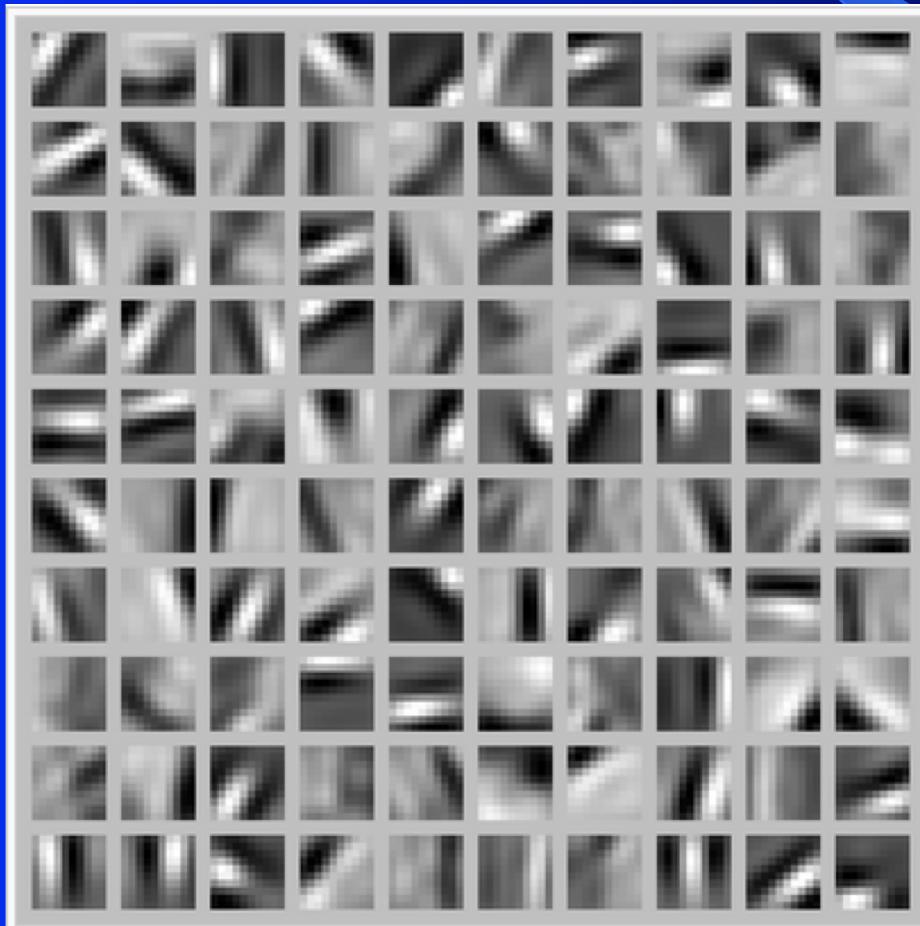


Sparse Encoders

- Auto encoders will often do a dimensionality reduction
 - PCA-like or non-linear dimensionality reduction
- This leads to a "dense" representation which is nice in terms of parsimony
 - All features typically have non-zero values for any input and the combination of values contains the compressed information
- However, this distributed and entangled representation can often make it more difficult for successive layers to pick out the salient features
- A *sparse* representation uses more features where at any given time many/most of the features will have a 0 value (ReLUs help)
 - Thus there is an implicit compression each time but with varying nodes
 - This leads to more localist variable length encodings where a particular node (or small group of nodes) with value 1 signifies the presence of a high-order feature (small set of bases)
 - A type of simplicity bottleneck (regularizer)
 - This is easier for subsequent layers to use for learning

Sparse Representation

- For bases below, which is easier to see intuition for current pattern - if a few of these are on and the rest 0, or if all have some non-zero value?
- Easier to learn if sparse



How do we implement a sparse Auto-Encoder?

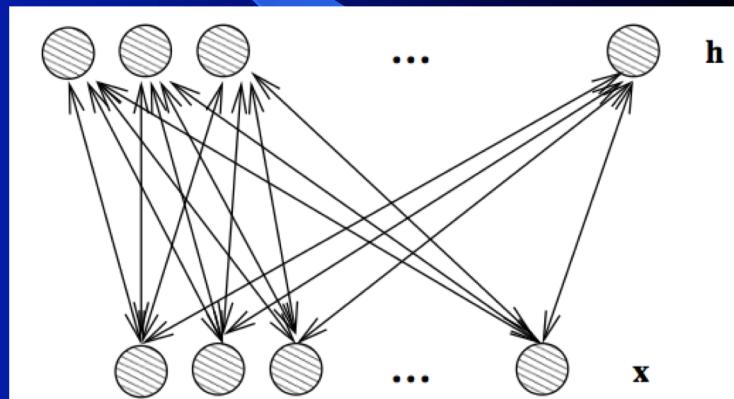
- Use more hidden nodes in the encoder
- Use regularization techniques which encourage sparseness (e.g. a significant portion of nodes have 0 output for any given input)
 - Penalty in the learning function for non-zero nodes
 - Weight decay
 - etc.
- De-noising Auto-Encoder
 - Stochastically corrupt training instance each time, but still train auto-encoder to decode the uncorrupted instance, forcing it to de-noise and learn conditional dependencies within the instance
 - Improved empirical results, handles missing values well

Stacked Auto-Encoders

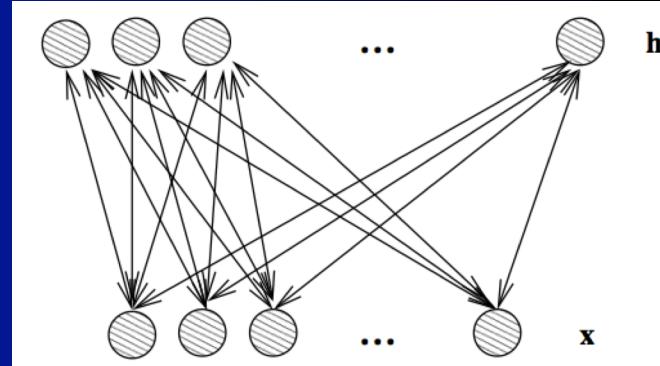
- Concatenation approach (i.e. using both hidden features and original features in final (or other) layers) can be better if not doing fine tuning. If fine tuning, the pure replacement approach works well.
- Always fine tune if there is a sufficient amount of labeled data
- For real valued inputs, auto-encode training is regression and thus could use linear output node activations (thrown out after anyways), still ReLU/non-linear at hidden which are final nodes
- Stacked Auto-Encoders empirically not quite as accurate as DBNs (Deep Belief Networks)
 - (with De-noising auto-encoders, stacked auto-encoders competitive with DBNs)
 - Not generative like DBNs, though recent work with de-noising auto-encoders may allow generative capacity

Deep Belief Networks (DBN)

- Geoff Hinton (2006) – Beginning of Deep Learning hype – outperformed kernel methods on MNIST – Also generative
- Uses Greedy layer-wise training but each layer is an RBM (Restricted Boltzmann Machine)
- RBM is a constrained Boltzmann machine with
 - No lateral connections between hidden (\mathbf{h}) and visible (\mathbf{x}) nodes
 - Symmetric weights
 - Does not use annealing/temperature, but that is all right since each RBM not seeking a global minima, but rather an incremental transformation of the feature space
 - Typically uses probabilistic logistic node, but other activations possible



RBM Sampling and Training



- Initial state typically set to a training example \mathbf{x} (can be real valued)
- Because of RBM, sampling is simple iterative back and forth process
 - $P(h_i = 1 | \mathbf{x}) = \text{sigmoid}(W_i \mathbf{x} + c_i) = 1/(1+e^{-\text{net}(h_i)})$ // c_i is hidden node bias
 - $P(x_i = 1 | \mathbf{h}) = \text{sigmoid}(W'_i \mathbf{h} + b_i) = 1/(1+e^{-\text{net}(x_i)})$ // b_i is visible node bias
- Contrastive Divergence (CD- k): How much contrast (in the statistical distribution) is there in the divergence from the original training example to the relaxed version after k relaxation steps
- Then update weights to decrease the divergence as in Boltzmann
- Typically just do CD-1 (Good empirical results)
 - Since small learning rate, doing many CD-1 updates is similar to doing fewer versions of CD- k with $k > 1$
 - Note CD-1 just needs to get the gradient direction right, which it usually does, and then change weights in that direction according to the learning rate

RBMupdate(x₁, ε, W, b, c)

This is the RBM update procedure for binomial units. It can easily adapted to other types of units.

x₁ is a sample from the training distribution for the RBM

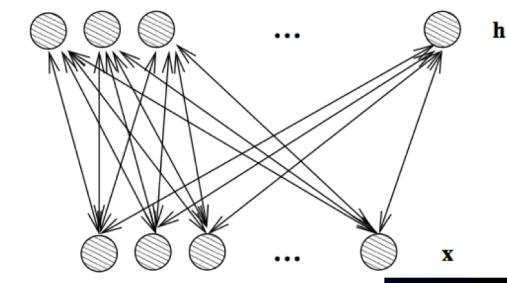
ε is a learning rate for the stochastic gradient descent in Contrastive Divergence

W is the RBM weight matrix, of dimension (number of hidden units, number of inputs)

b is the RBM offset vector for input units

c is the RBM offset vector for hidden units

Notation: Q(h_{2.} = 1|x₂) is the vector with elements Q(h_{2i} = 1|x₂)



for all hidden units i do

- compute Q(h_{1i} = 1|x₁) (for binomial units, sigm(c_i + ∑_j W_{ij}x_{1j}))
- sample h_{1i} ∈ {0, 1} from Q(h_{1i}|x₁)

end for

for all visible units j do

- compute P(x_{2j} = 1|h₁) (for binomial units, sigm(b_j + ∑_i W_{ij}h_{1i}))
- sample x_{2j} ∈ {0, 1} from P(x_{2j} = 1|h₁)

end for

for all hidden units i do

- compute Q(h_{2i} = 1|x₂) (for binomial units, sigm(c_i + ∑_j W_{ij}x_{2j}))

end for

• W ← W + ε(h₁x'₁ - Q(h_{2.} = 1|x₂)x'₂)

• b ← b + ε(x₁ - x₂)

• c ← c + ε(h₁ - Q(h_{2.} = 1|x₂))

$$\Delta w_{ij} = \varepsilon(h_{1,j} \cdot x_{1,i} - Q(h_{k+1,j} = 1 | x_{k+1})x_{k+1,i})$$

$$\Delta w_{ij} = \varepsilon(\text{initial_h_sample} \cdot \text{initial_x} - \text{final_h_probability} \cdot \text{final_x_sample})$$

RBM Update Notes and Variations

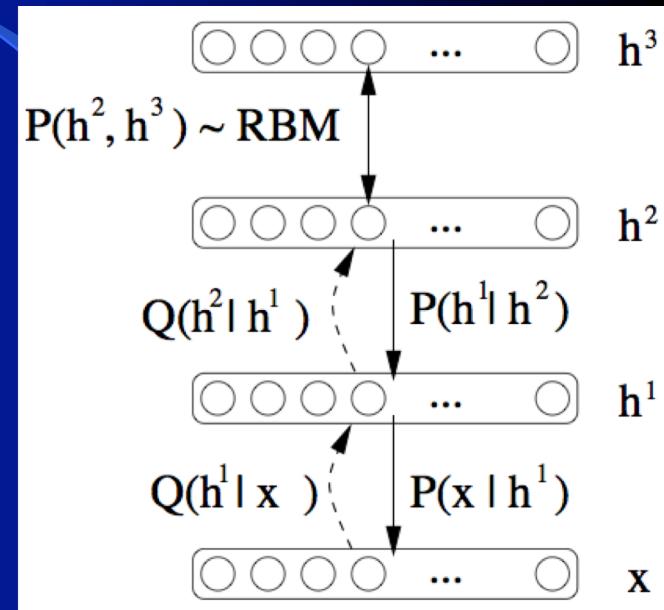
- Example (based on HW with different values) and Homework
- Binomial unit means the standard MLP sigmoid unit
- Q and P are probability distribution vectors for hidden (\mathbf{h}) and visible/input (\mathbf{x}) vectors respectively
- During relaxation/weight update can alternatively do updates based on the real valued probabilities ($\text{sigmoid}(\text{net})$) rather than the 1/0 sampled logistic states
 - Always use actual/binary values from initial $\mathbf{x} \rightarrow \mathbf{h}$
 - Doing this makes the hidden nodes a sparser bottleneck and is a regularizer helping to avoid overfit
 - Could use probabilities on the $\mathbf{h} \rightarrow \mathbf{x}$ and/or final $\mathbf{x} \rightarrow \mathbf{h}$
 - in CD- k the final update of the hidden nodes usually uses the probability value to decrease the final arbitrary sampling variation (sampling noise)
- Lateral restrictions of RBM allow this fast sampling

RBM Update Variations and Notes

- Initial weights, small random, 0 mean, sd $\sim .01$
 - Don't want hidden node probabilities early on to be close to 0 or 1, else slows learning, since less early randomness/mixing? Note that this is a bit like annealing/temperature in Boltzmann
- Set initial **x** bias values as a function of how often node is on in the training data, and **h** biases to 0 or negative to encourage sparsity
- Better speed when using momentum ($\sim .5$)
- Weight decay good for smoothing and also encouraging more mixing (hidden nodes more stochastic when they do not have large net magnitudes)

Deep Belief Network Training

- Same greedy layer-wise approach
- First train lowest RBM ($h^0 - h^1$) using RBM update algorithm (note h^0 is x)
- Freeze weights and train subsequent RBM layers
- Then connect final outputs to a supervised model and train that model
- Finally, unfreeze all weights, and fine tune as an MLP using the initial weights found by DBN training
- Can do execution as just the tuned MLP or as the RBM sampler with the tuned weights

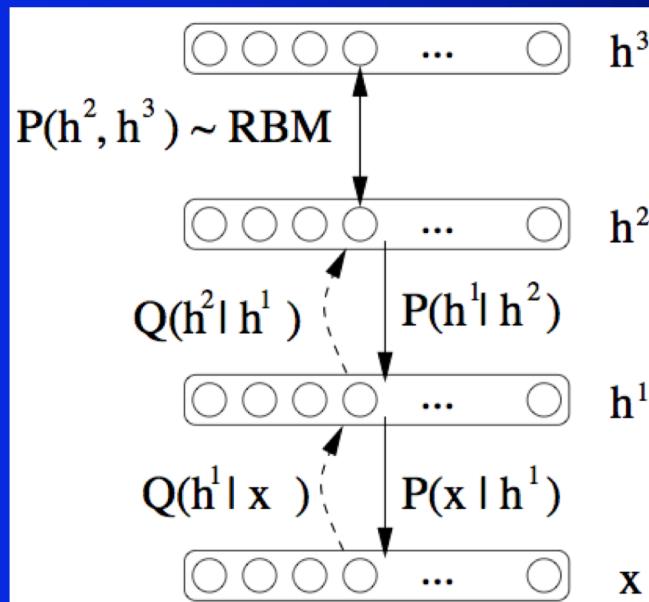


During execution can iterate multiple times at the top RBM layer

Can use DBN as a Generative model to create sample x vectors

1. Initialize top layer to an arbitrary vector (commonly a training set vector)
 - Gibbs sample (relaxation) between the top two layers m times
 - If we initialize top layer with values obtained from a training example, then need less Gibbs samples
2. Pass the vector down through the network, sampling with the calculated probabilities at each layer
3. Last sample at bottom is the generated x vector (can be real valued if we use the probability vector rather than sample)

Alternatively, can start with an x at the bottom, relax to a top value, then start from that vector when generating a new x , which is the dotted lines version. More like standard Boltzmann machine processing.



TrainUnsupervisedDBN(\hat{P} , ϵ , ℓ , W , b , c , mean_field_computation)

Train a DBN in a purely unsupervised way, with the greedy layer-wise procedure in which each added layer is trained as an RBM (e.g. by Contrastive Divergence).

\hat{P} is the input training distribution for the network

ϵ is a learning rate for the RBM training

ℓ is the number of layers to train

W^k is the weight matrix for level k , for k from 1 to ℓ

b^k is the visible units offset vector for RBM at level k , for k from 1 to ℓ

c^k is the hidden units offset vector for RBM at level k , for k from 1 to ℓ

mean_field_computation is a Boolean that is true iff training data at each additional level is obtained by a mean-field approximation instead of stochastic sampling

for $k = 1$ to ℓ **do**

- initialize $W^k = 0$, $b^k = 0$, $c^k = 0$

while not stopping criterion **do**

- sample $h^0 = x$ from \hat{P}

for $i = 1$ to $k - 1$ **do**

if mean_field_computation **then**

- assign h_j^i to $Q(h_j^i = 1 | h^{i-1})$, for all elements j of h^i

else

- sample h_j^i from $Q(h_j^i | h^{i-1})$, for all elements j of h^i

end if

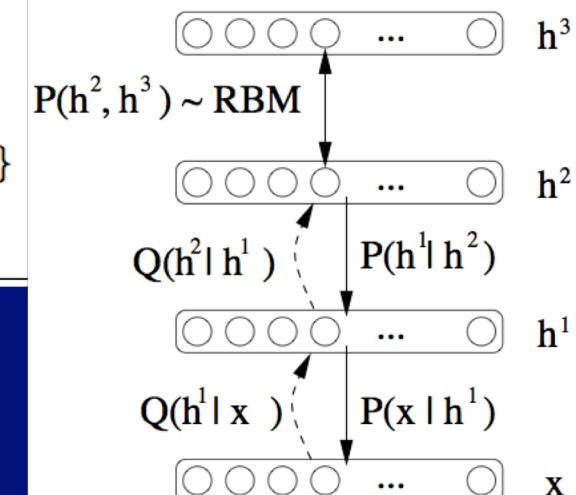
end for

- RBMupdate(h^{k-1} , ϵ , W^k , b^k , c^k) {thus providing $Q(h^k | h^{k-1})$ for future use}

end while

end for

Middle for-loop samples from input up to current RBM layer being updated – none for 1st layer, mean_field_computation just a flag on whether to sample or use real values



DBN Execution

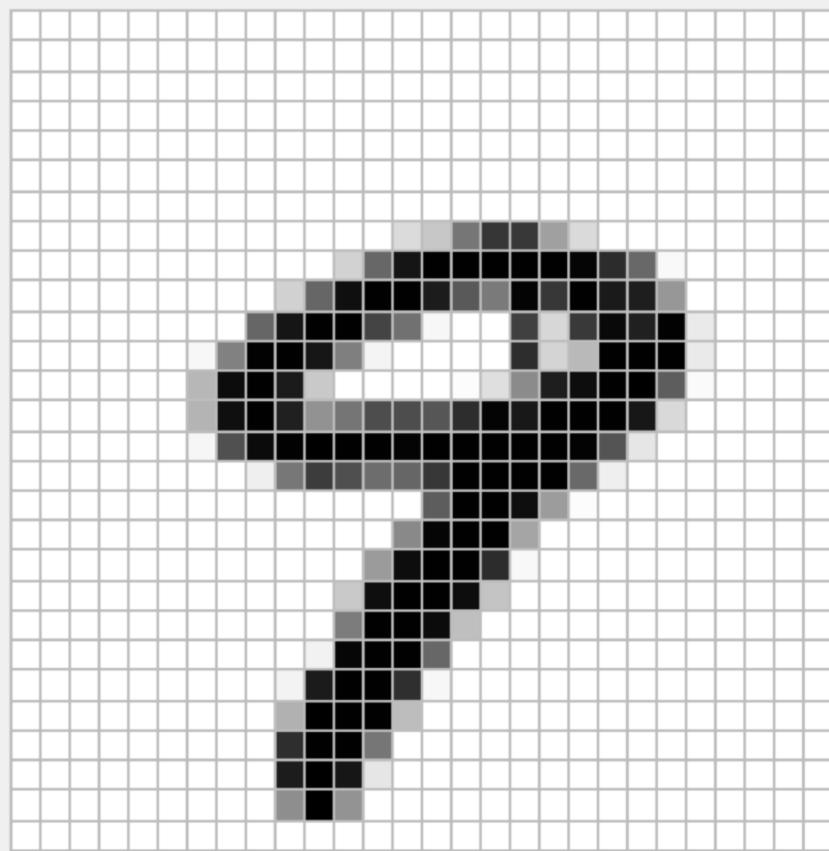
- After all layers have learned then the output of the last layer can be input to a supervised learning model
- Note that at this point we could potentially throw away the downward bias weights in the network as they will not actually be used during the feedforward discriminative execution process (as we did with the Stacked Auto Encoder)
 - If we are relaxing M times in the top layer then we would still need the downward weights for that layer
 - Also if we are generating x values we would need all of them
- The final weight tuning is usually done as an MLP with backpropagation, which only updates the feedforward weights
- Can do execution as just the tuned MLP or as the RBM sampler with the tuned weights

DBN Learning Notes

- RBM stopping criteria still in issue.
- Each layer updates weights so as to make training sample patterns more likely (lower energy) in the free state (and non-training sample patterns less likely).
- This unsupervised approach learns broad features (in the hidden/subsequent layers of RBMs) which can aid in the process of making the types of patterns found in the training set more likely. This discovers features which can be associated across training patterns, and thus potentially helpful for other goals with the training set (classification, compression, etc.)
- Note still pairwise weights in RBMs, but because we can choose the number of hidden units and layers, we can represent any arbitrary distribution

MNIST file: C:\Users\Mapper\Documents\Disser\Soft\MNIST_Experiments\MNIST Recognize\MN ...

Icon Edit Pane (28 X 28)



Preview



Recognition result

9

0 1 2 3 4 5 6 7 8 9

Recognize

Clear

OK

Cancel

DBN Project Notes

- To be consistent just use 28×28 (764) data set of gray scale values (0-255)
 - Normalize to 0-1
 - Could try better preprocessing if want and helps in published accuracies, but start/stay with this
 - Small random initial weights
- Parameters
 - Hinton Paper, others – do a little searching and e-mail me a reference for extra credit points
 - <http://yann.lecun.com/exdb/mnist/> for sample approaches
- Straight 200 hidden node MLP does quite good ~98%
 - Rough Hyperparameters - LR: ~.05-.1, Momentum ~.5
- Best class DBN results: ~98.5% - which is competitive
 - About half students never beat MLP baseline
 - Can you beat the 98.5%?

Deep Learning Project Past Experience

- Structure: ~3 hidden layers, ~500ish nodes/layer, more nodes/layers can be better but training is longer
- Training time:
 - DBN: ~10 epochs with the 60K set, small LR $\sim .005$ often good
 - Can go longer, does not seem to overfit with the large data set
 - SAE: Can saturate/overfit, ~3 epochs good, but will be a function of your denoising approach, which is essential for sparsity, use small LR $\sim .005$, long training – up to 50 hours, got 98.55
 - Larger learning rates often lead to low accuracy for both DBN and SAE
- Sampling vs using real probability value in DBN
 - Best results found when using real values vs. sampling
 - Some found sampling on the back-step of learning helps
 - When using sampling, probably requires longer training, but could actually lead to bigger improvements in the long run
 - Typical forward flow non-sampled during execution, but could do some sampling on the m iterations at the top layer. Some success with back-step at the top layer iteration (most don't do this at all)
 - We need to try/discover better variations

Deep Learning Project Past Experience

- Note: If we held out 50K of the dataset as unsupervised, then deep nets would more readily show noticeable improvement over BP
- A final full network fine tune with BP always helps
 - But can take 20+ hours
- Key take away – Most actual time spent training with different parameters. Thus, start early, and then you will have time to try multiple long runs to see which variations work. This does not take that much personal time, as you simply start it with some different parameters and go away for a day. If you wait until the last few days, there is no time to do these experiments.

DBN Notes

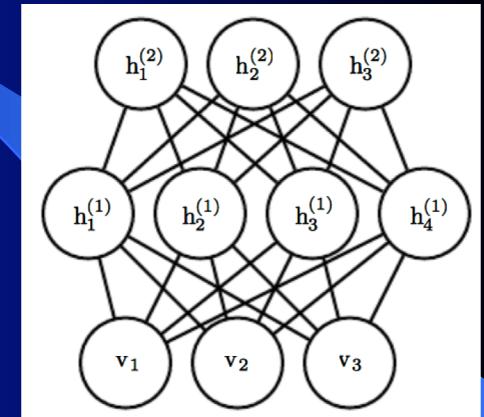
- Can use lateral connections in RBM (no longer RBM) but sampling becomes more difficult (intractable, approximation such as MCMC) – ala standard Boltzmann requiring longer sampling chains.
 - Lateral connections can capture pairwise dependencies allowing the hidden nodes to focus on higher order issues. Can get better results.
- Conditional and Temporal RBMs – allow node probabilities to be conditioned by some other inputs – context, recurrence (time series changes in input and internal state), etc.

Discrimination with Deep Belief Networks

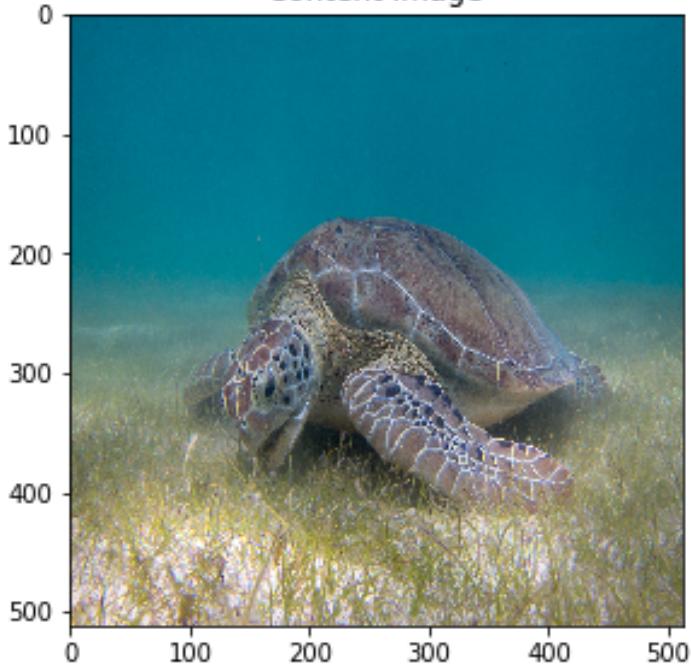
- Discrimination approaches with DBNs (Deep Belief Net)
 - Use outputs of DBNs as inputs to supervised model (i.e. just an unsupervised preprocessor for feature extraction)
 - Basic approach we have been discussing
 - Train a DBN for each class. For each clamp the unknown x and iterate m times. The DBN that ends with the lowest normalized free energy (softmax variation) is the winner.
 - Train just one DBN for all classes, but with an additional visible unit for each class. For each output class:
 - Clamp the unknown x , relax, and then see which final state has lowest free energy – no need to normalize since all energies come from the same network.
- See <http://deeplearning.net/demos/>

Other Deep Generative Models

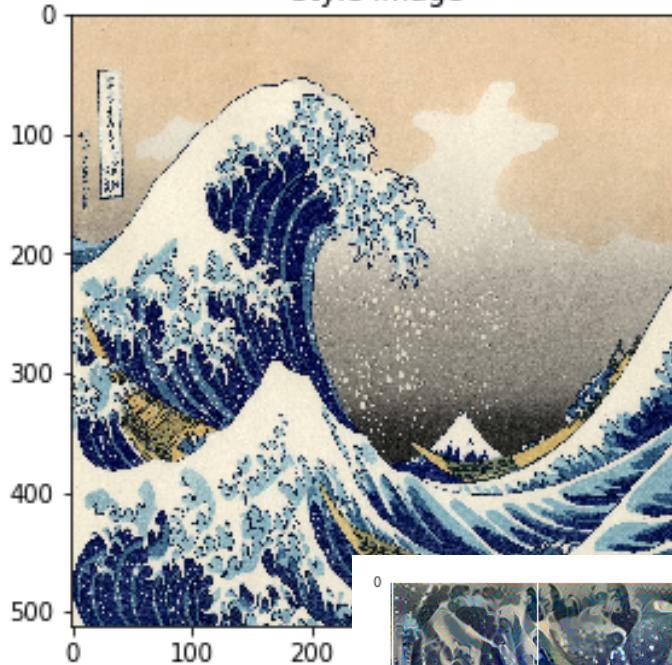
- Lots of research on generative models to create probabilistic models of training data with ability to generate new images, sentences, etc.
- Deep Boltzmann machine – no lateral connections, simpler block sampling (ala RBM) as compared to the fully connected Boltzmann, all even/odd layers can simultaneously update
- Trained CNNs can generate images (deconvolution) starting from a random seed



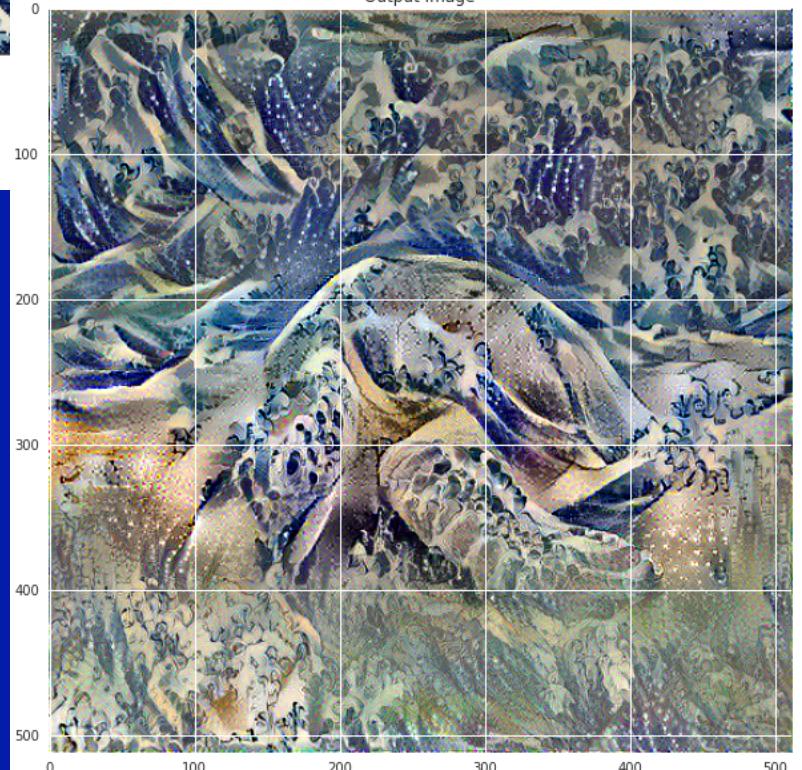
Content Image



Style Image



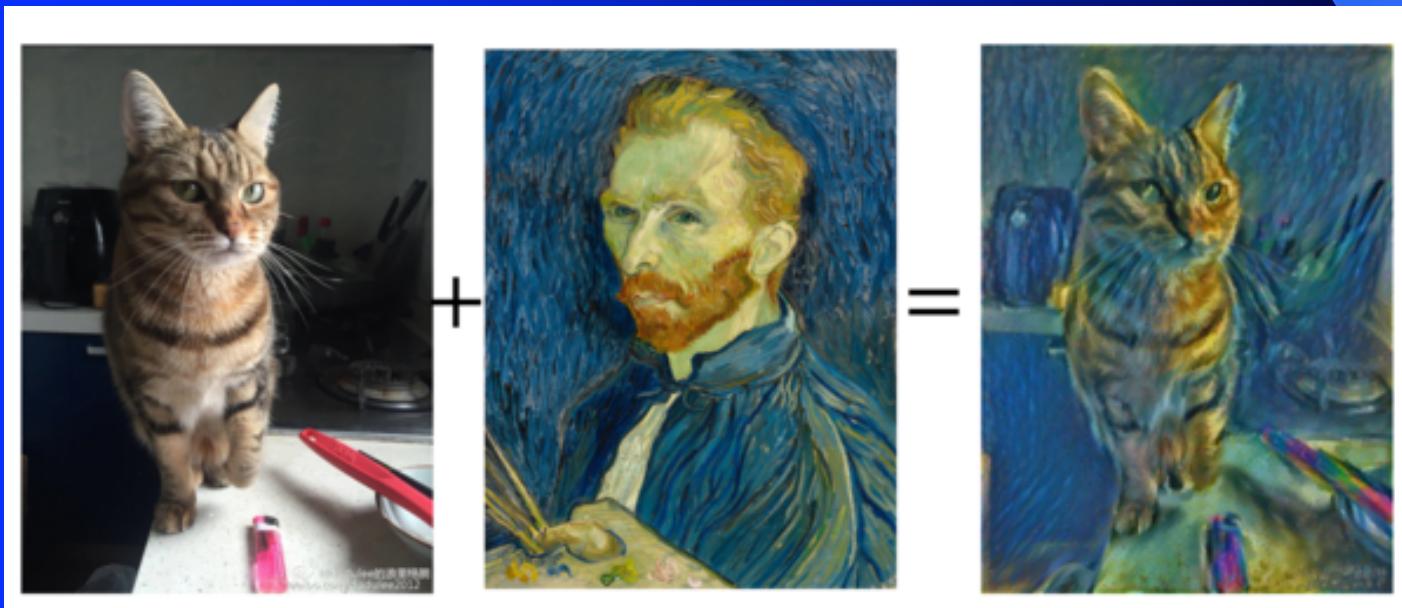
Output Image



Neural Style Transfer

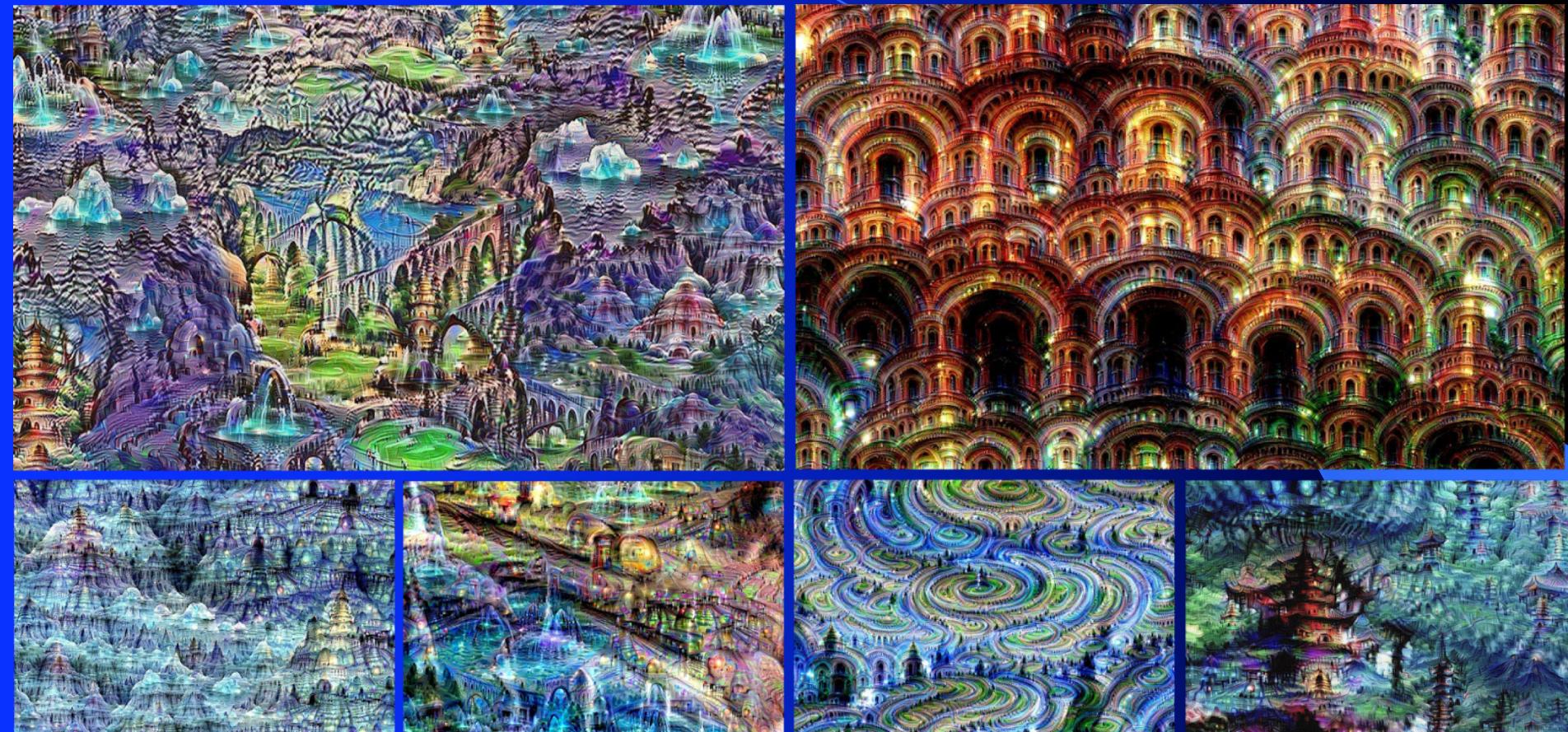
- Train on lots of images and styles
- CNN trained with two loss functions
 - content and style
- Then supply any content and style image
- Creates new image of content, but in the style of the style image

Neural Style Transfer



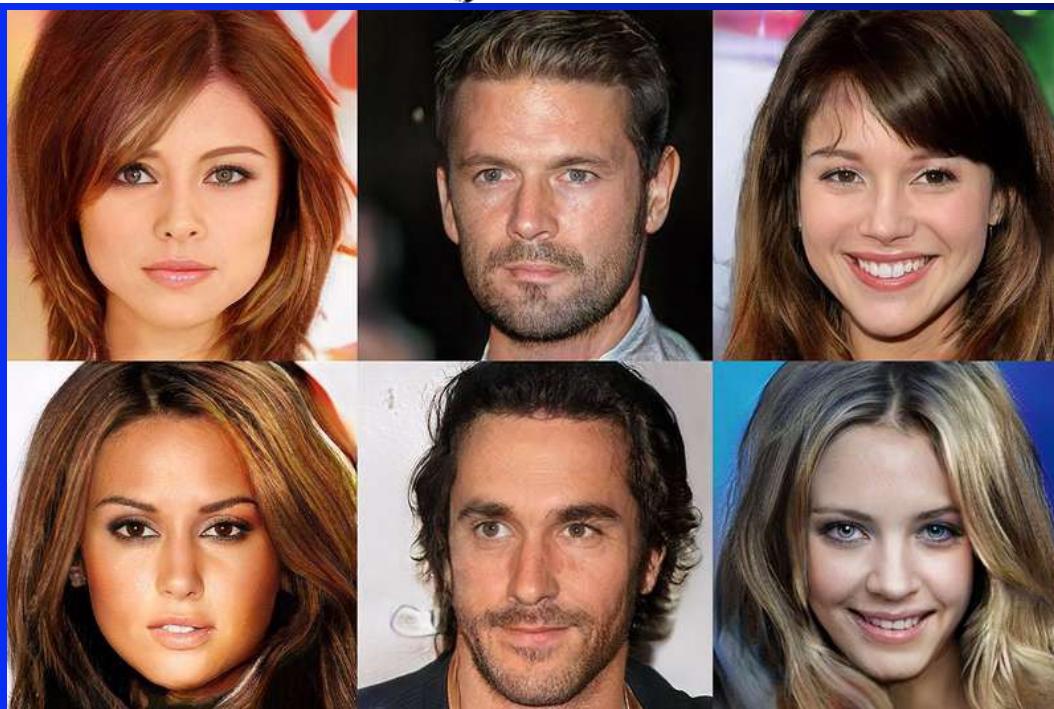
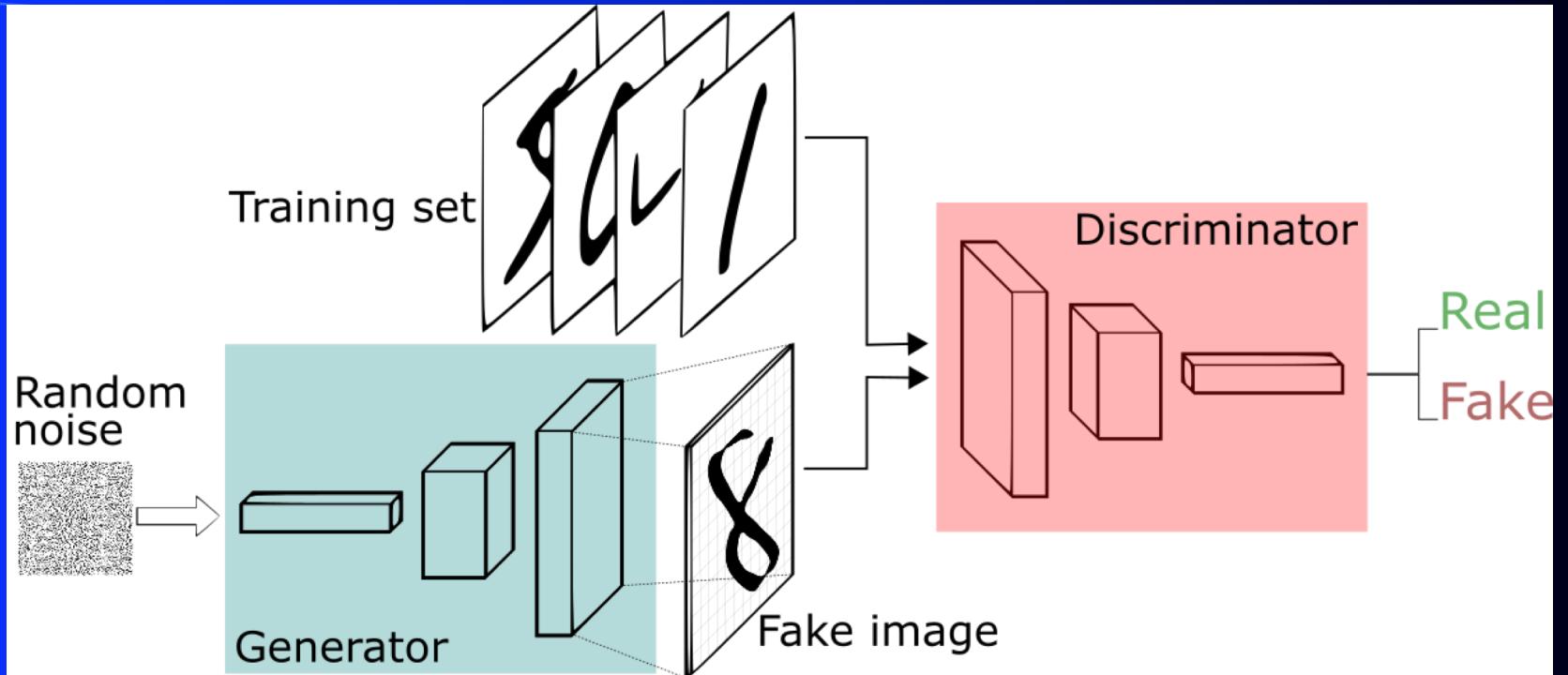


Deep Dreaming



Generative Adversarial Networks (GANs)

- Unsupervised
- Generative networks which generate novel samples x similar to training samples
- Discriminative net (adversary) must differentiate between samples from the generative net and the training set
- Use loss feedback on discriminator net to create gradient for both nets, until discriminator can no longer distinguish, then can discard discriminator net – increasingly difficult for humans to distinguish
- "Universal loss function" for lots of "difficult/creative" applications



Edmond de Belamy

- Created by a GAN and sold at auction in 2018 for \$432,500
- Note the author inscription – it is the GAN loss function



Fully Supervised Deep Learning

- More recent success in doing fully supervised deep learning with extensions which diminish the effect of early learning difficulties (unstable gradient, etc.)
- Patience (now that we know it can be worth it), faster computers, and use of GPUs
- More efficient activation functions (e.g. ReLUs) in terms of both computation and avoiding $f'(net)$ saturation
 - Also can be helpful to have 0 mean activations at each level, so sigmoid is frowned upon these days. If you want a saturating activation function, tanh is often preferred.
- Speed up and regularization approaches
- Improved Hyperparameters
- Batch Normalization – re-normalize activations at each layer
- Residual Nets
- Inception/Xception
- LSTM, GRUs
- Deep Reinforcement Learning – Alpha Zero

Speed up variations of SGD

- Use mini-batch rather than single instance for better gradient estimate
 - Helpful if using GD variation more sensitive to bad gradient, and *especially* for parallel implementations
- Momentum (i.e. Adaptive learning rate) approaches are important since anything to speed-up learning is helpful
 - Standard Momentum
 - Note these approaches already do an averaging of grading also making mini-batch less critical
 - Nesterov Momentum – Calculate point you would go to if using normal momentum. Then, compute gradient at that point. Do normal update using *that* gradient and momentum.
 - Rprop – Resilient BP, if gradient sign inverts, decrease the node's individual LR, else increase it – common goal is faster in the flats, there are variants that backtrack a step, etc.
 - Adagrad – Scale LRs inversely proportional to $\text{sqrt}(\text{sum(historical values)})$ – LRs with smaller derivatives are decreased less
 - RMSprop – Adagrad but uses exponentially weighted moving average, older updates basically forgotten
 - Adam (Adaptive moments) – Momentum terms on both gradient and squared gradient (1st and 2nd moments) – update based on both

Regularization – Dropout Common

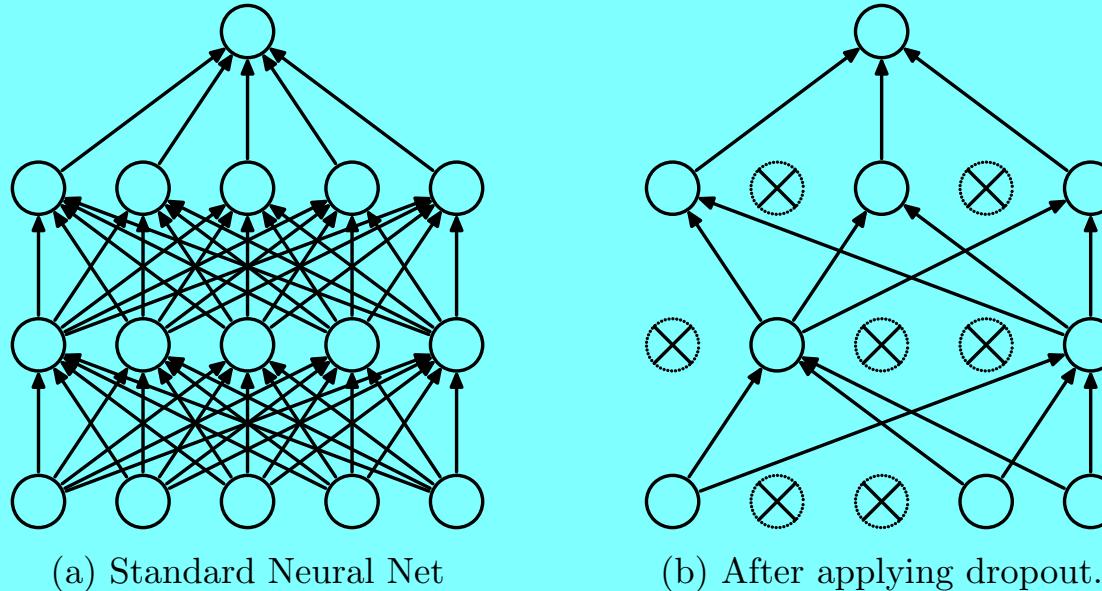


Figure 1: Dropout Neural Net Model. **Left:** A standard neural net with 2 hidden layers. **Right:** An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.

- For each instance drop a node (hidden or input) and its connections with probability p and train
- Final net just has all averaged weights (actually scaled by $1-p$ since that better matches the expected values at training time)
- As if ensembling 2^n different network substructures
- Lots of variations – Dropconnect, etc.

Improved Initial Hyperparameter Settings

- Deep networks are more sensitive than shallow networks to hyperparameter settings
- More critical for deep learning in order to get more balanced learning across all layers
- Smaller LRs – patience
- To encourage sparsity sometimes initial biases set negative or more initial 0 weights are interspersed
- Initial weights – initialize a little larger in effort to find a balance which learns well across all layers. Common is to select initial weights from a uniform distribution between
 - $-c/\sqrt{\text{node fan-in}}, c/\sqrt{\text{node fan-in}}$ ($c = 1$ Xavier, $c = 2$ He)
 - $-c/\sqrt{\text{node fan-in} + \text{fan-out}}, c/\sqrt{\text{node fan-in} + \text{fan-out}}$
 - Can do Gaussian distribution with above as variances
 - Lots of other variations and current work

Batch Normalization (2015)

- Rather than just guess initial parameters to maintain learning balance, renormalize activations at each layer to maintain balance
- Just like it is critical to normalize our initial features, we consider each layer to be new features which can also be normalized
- We like 0-mean and unit variance inputs (standard 1st layer normalization), we can just re-normalize the activation (more commonly the net value) for each input dimension k at each layer

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

- Want mean and variance of that activation for the entire data set. Approximate the empirical mean and variance over a mini-batch of instances
- Goes beyond standard normalization by allowing scaling and shifting of the normalized values with 2 learnable weights per input, γ and β , to attain the final batch normalization (allows recovery of initial or any needed function)

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

Batch Normalization

[Ioffe and Szegedy, 2015]

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

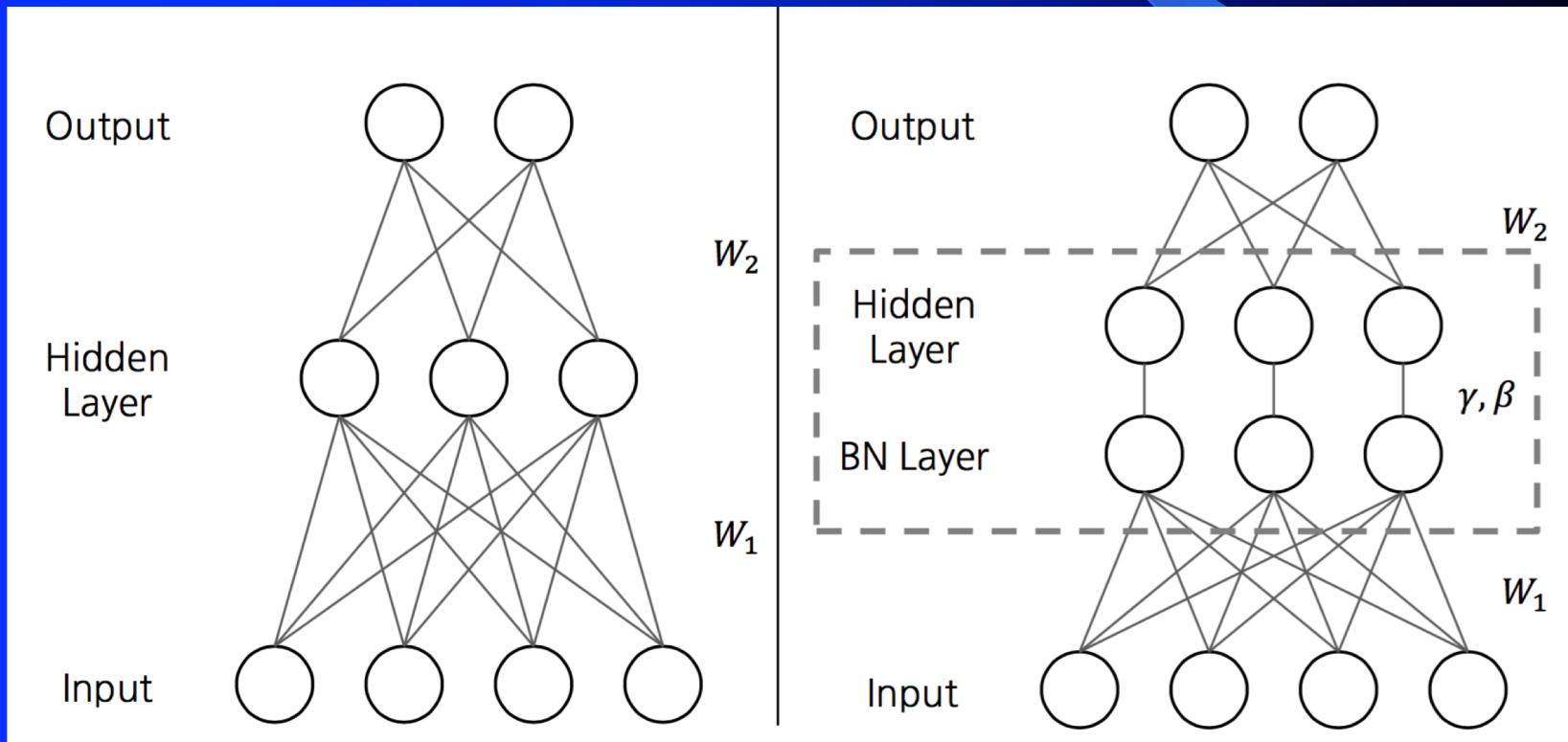
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

- Note: Dropped k for simplicity
- γ and β learned as part of gradient descent
- At test time BatchNorm layer functions differently:
- The mean/std are not computed based on a batch. Instead, we use empirical means of values found during training
- (e.g. can be estimated during training with running averages)

Batch Normalization

- Normalize before the non-linearity (e.g. ReLU)
- Allows larger LR (faster learning), improves gradient flow and reduces dependence on initialization



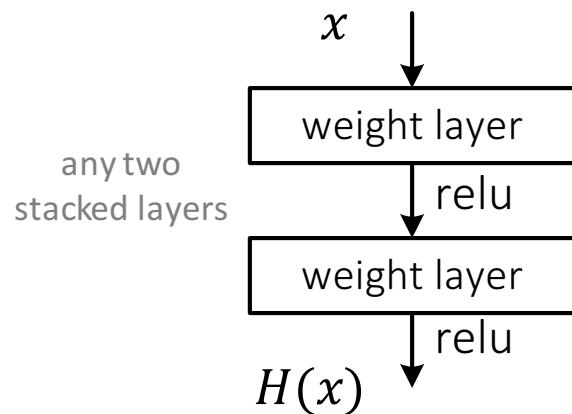
Deep Residual Learning

- Residual Nets – 100s of layers
- 2015 ILSVRC winner
 - CNN
 - Used Batch Normalization extensively
- Learns the residual mapping with respect to the identity – i.e. the *difference* (residual) between the current input and the goal mapping
- Simple concept which tends to make the function to be learned simpler across depth

Deep Residual Learning

Deep Residual Learning

- Plain net



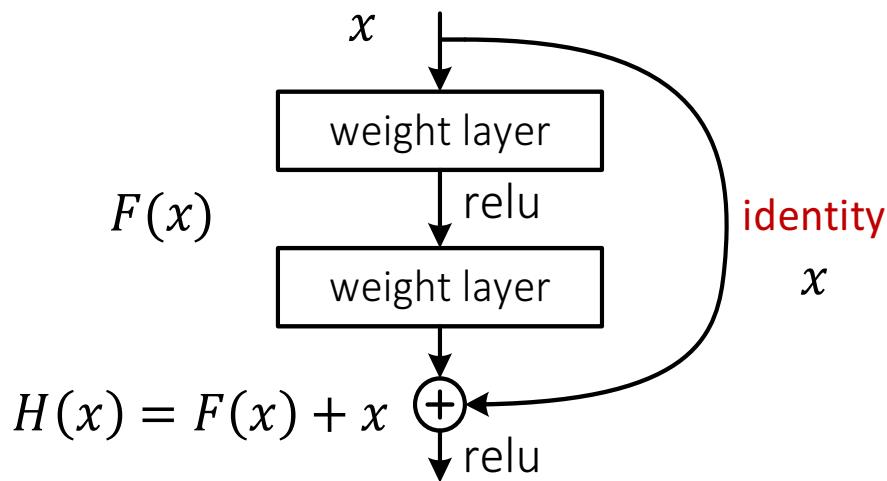
$H(x)$ is any desired mapping,
hope the 2 weight layers fit $H(x)$

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

Deep Residual Learning

Deep Residual Learning

- $F(x)$ is a **residual** mapping w.r.t. **identity**



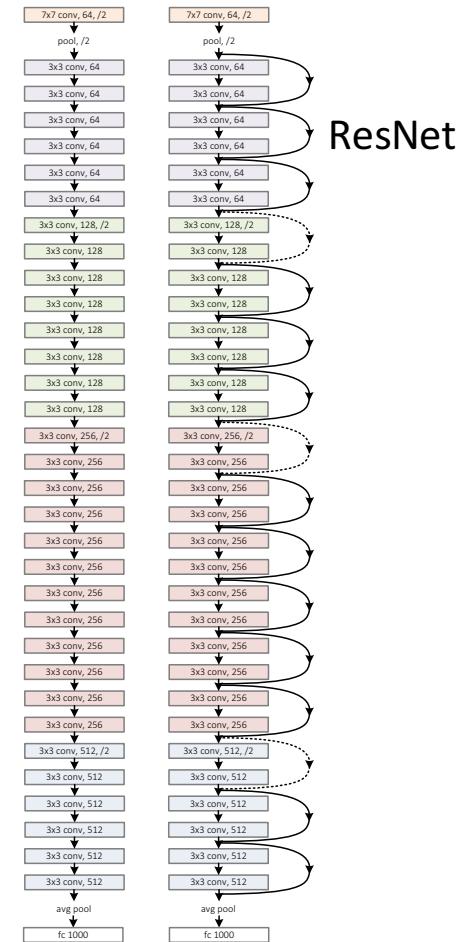
- If identity were optimal, easy to set weights as 0
- If optimal mapping is closer to identity, easier to find small fluctuations

Deep Residual Learning

Network “Design”

plain net

- Keep it simple
- Our basic design (VGG-style)
 - all 3x3 conv (almost)
 - spatial size /2 => # filters x2 (~same complexity per layer)
 - Simple design; just deep!
- Other remarks:
 - no hidden fc
 - no dropout



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. “Deep Residual Learning for Image Recognition”. CVPR 2016.

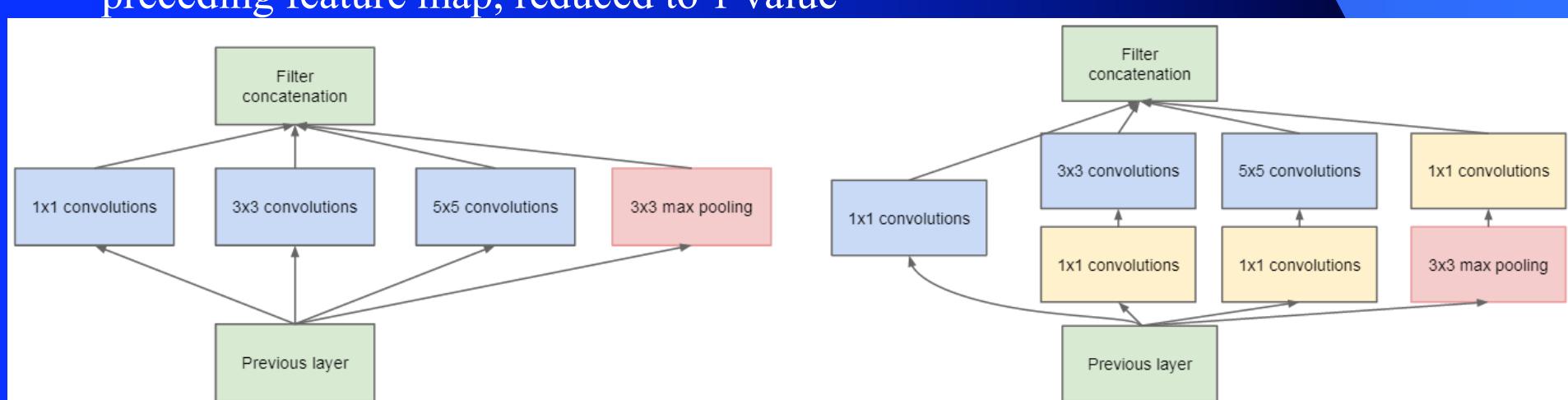
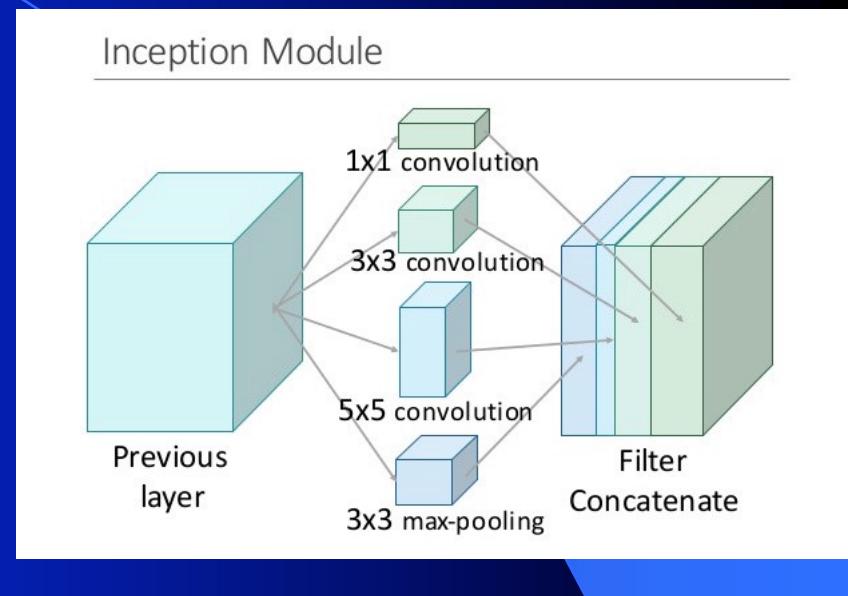
How many layers to add residuals, etc.? – Trial and error

Residual Nets

- Going from an x to an $H(x)$ which is quite different requires more learning, larger weights, etc.
- However, if $H(x)$ is pretty similar to x , and assuming we start with small weights anyways, it takes a lot less updates to learn it – Easier!
- Adding x to the later layer allows it to learn this simpler mapping
- Also, if the net had already learned a particular feature, we can just maintain that feature with a 0 weight, without having to relearn it all the time – avoid "feature attrition"

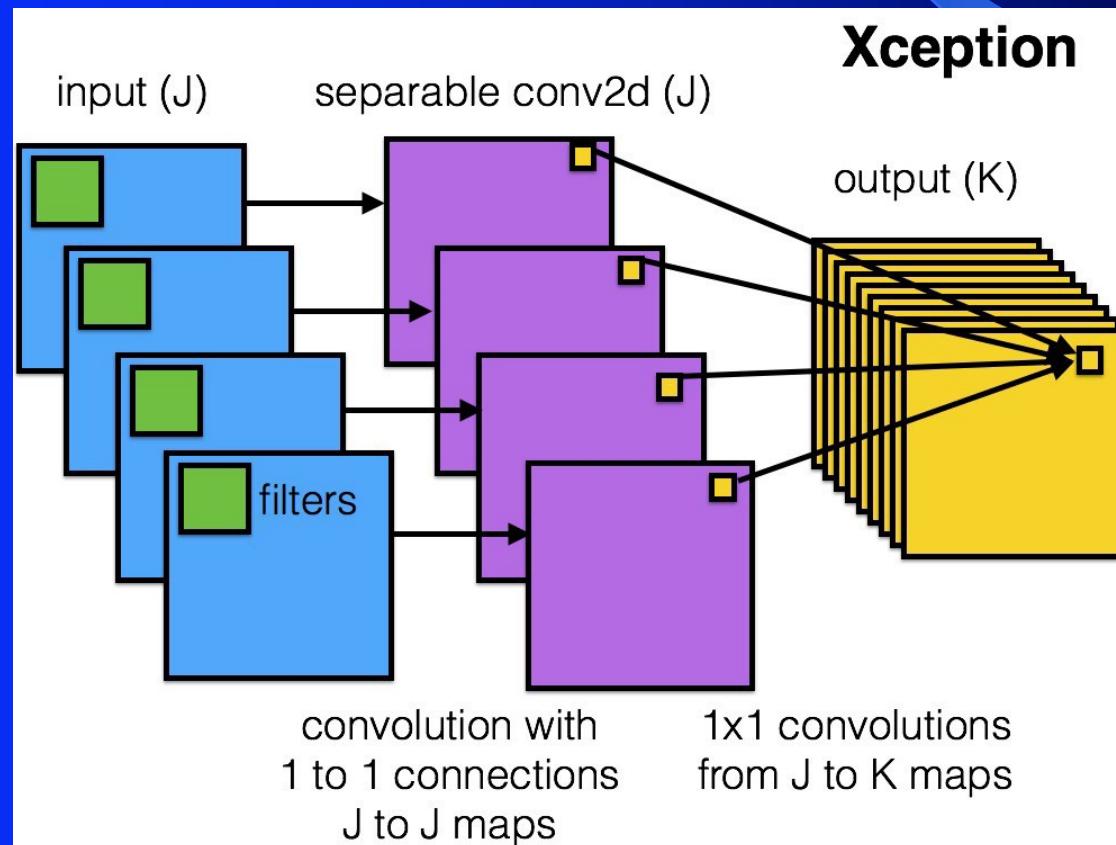
Inception - Google

- "Network in a network" – CNN – Deeper and *Wider* - GoogLeNet
- Could replace the basically linear convolution with a more complex non-linear network – e.g. MLP
- Basic Inception does different size convolutions and combines results into one output
- Reduces added complexity by first doing dimensionality reduction with 1x1 convolution filters – 1 input from each preceding feature map, reduced to 1 value



Xception

- Xception (extreme inception) decouples spatial and cross-channel correlations – channels only connected at extra layers of 1x1 convolutional maps - wider



Recurrent Neural Networks

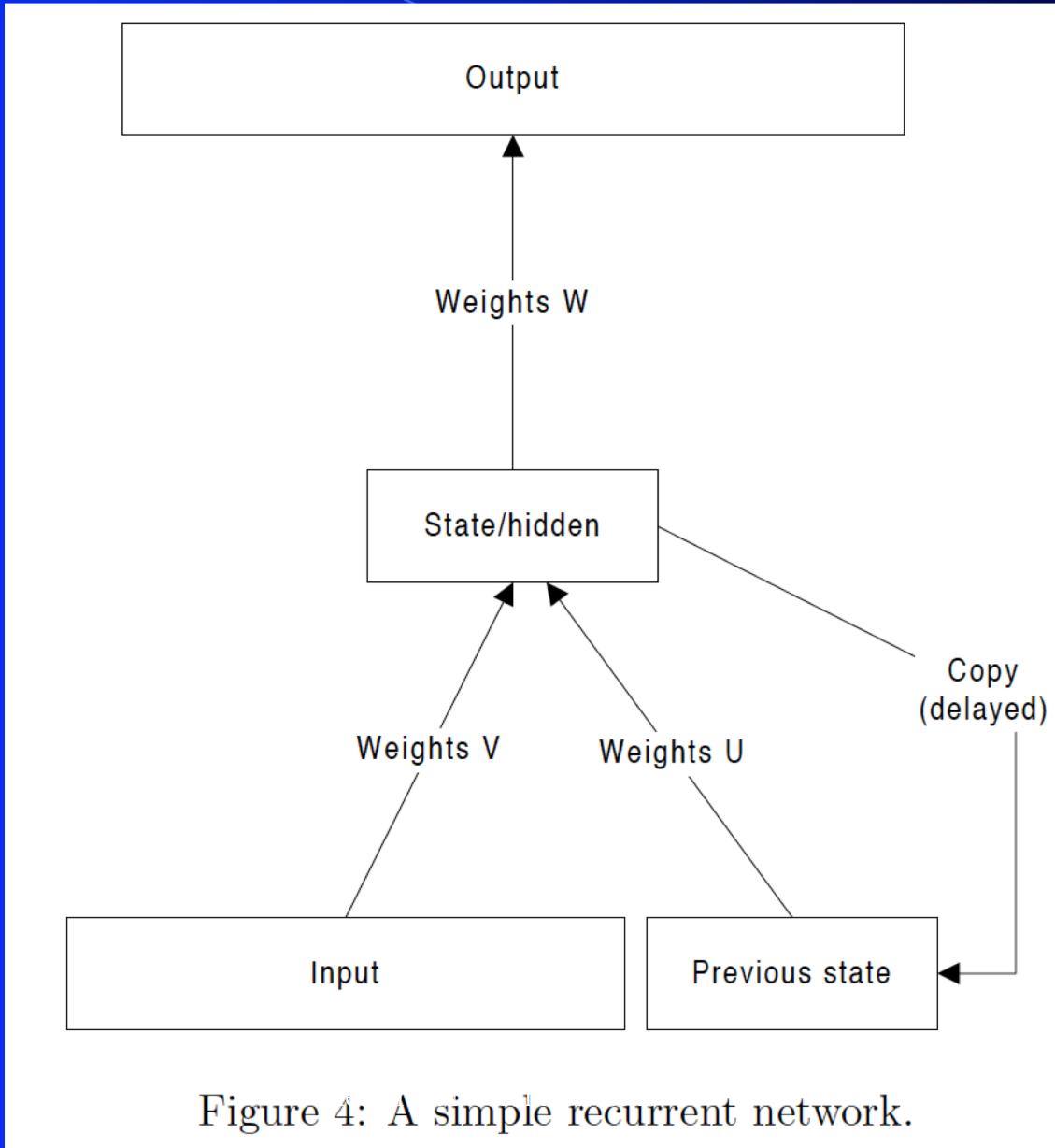
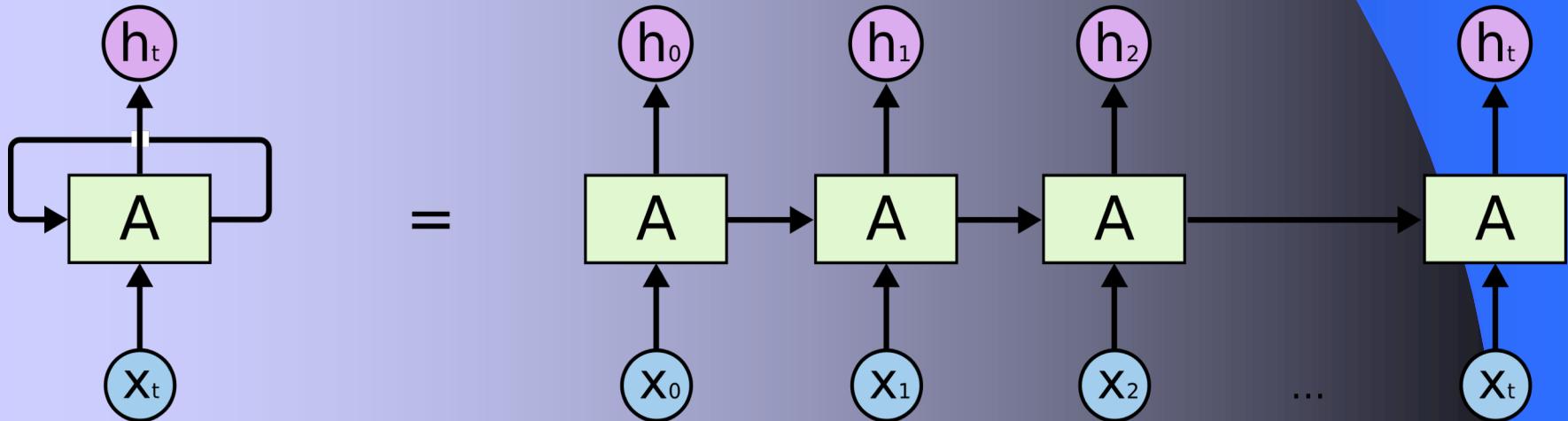


Figure 4: A simple recurrent network.

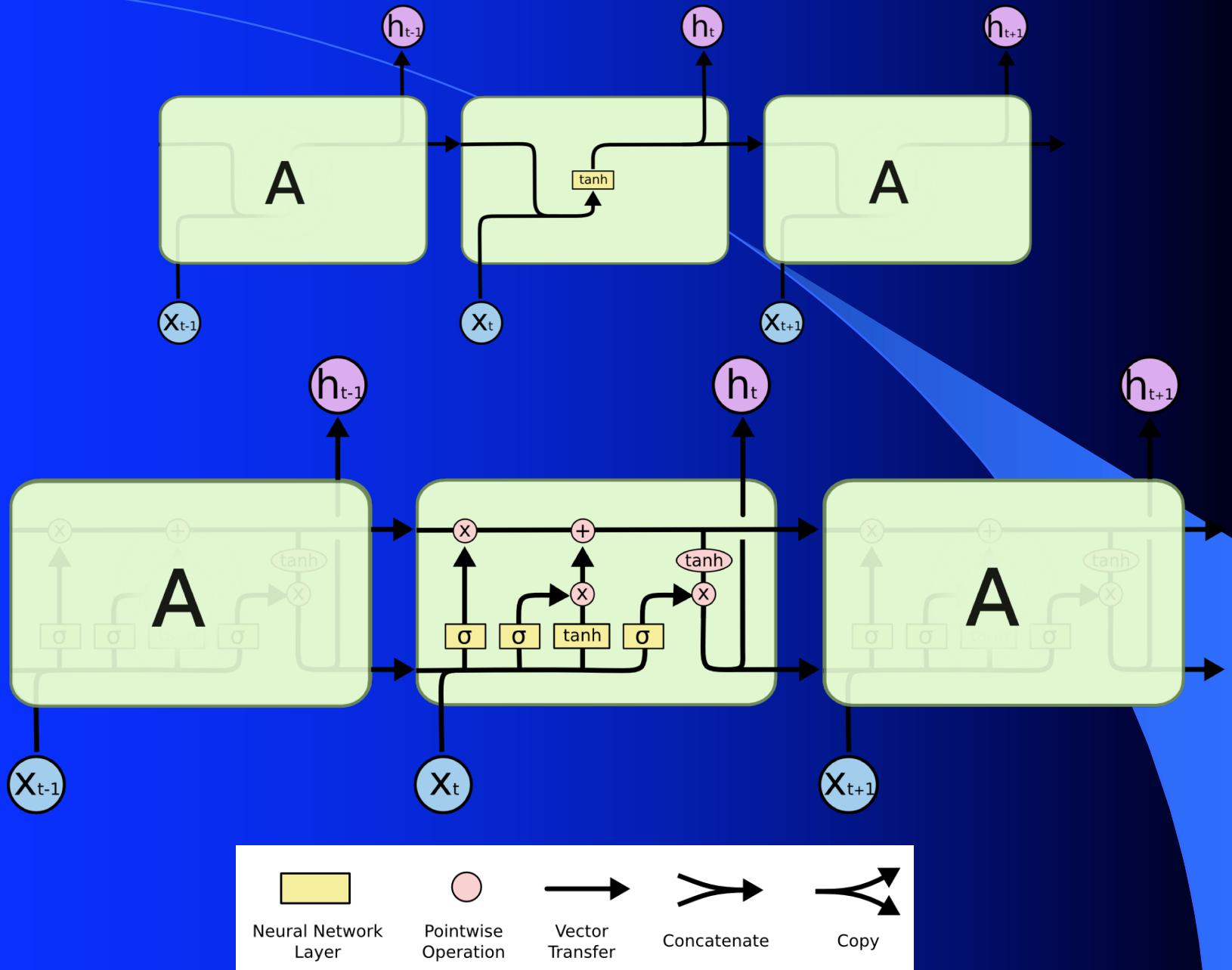
LSTM/GRU

- Long Short-Term Memory/Gated Recurrent Unit
- Pictures and flow from Olah's Blog
- We have been adding a layer of weights between h_t and o_t
- LSTM – (GRU is LSTM subset) – State of the art in RNN
 - Standard RNN node is basically LSTM without forget, ignore, and output gates
 - Train with BPTT but bigger k 's (a full sequence if not too large), or some pretty big chunk (25-100) since we avoid vanishing gradient.



LSTM – Long Short Term Memory

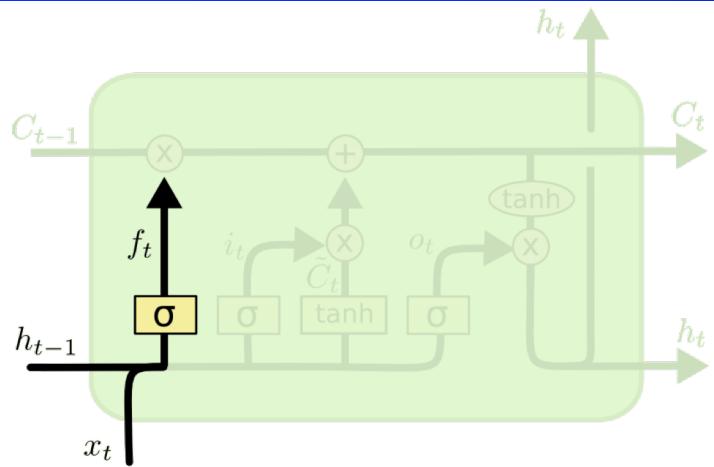
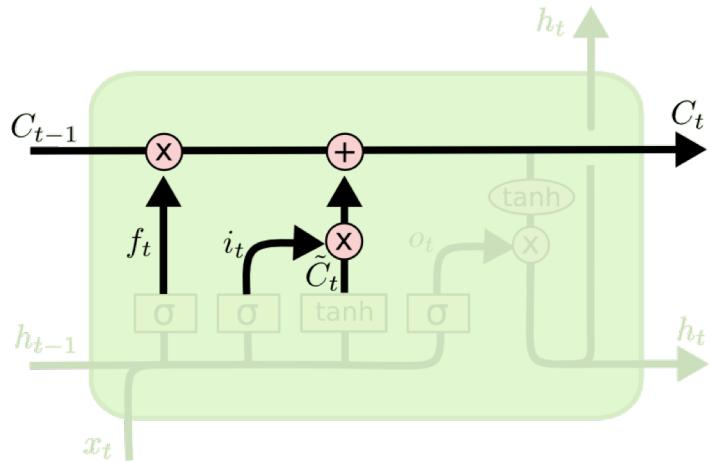
- LSTM unit can just plug in for a standard node in an RNN
- Adds a state memory, input, forget, and output gates
- Vanishing/exploding gradient avoidance
 - Cell value (state) has a self-feedback loop and can maintain its value indefinitely. Has derivative 1 with no vanishing gradient.
 - The cell value is multiplied by a forget gate output (0-1), which decides when and how much to forget, giving much more power.
- LSTM unit has lots more parameters but still trained with standard BP/SGD learning: typically BPTT
 - Forget gates, etc. don't know that is their job, but the capacity is there during training to learn that job as the overall network minimizes loss
 - Capacity plus training finds a way to solve the problem, capacity that wasn't there with simple network nodes



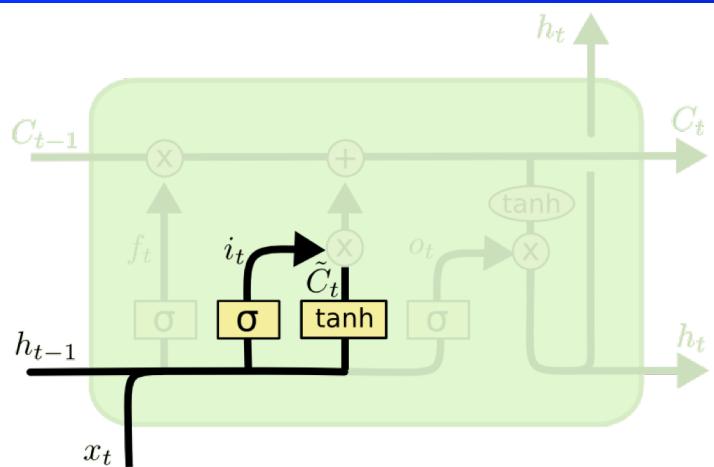
- C_t is cell state, h_t is context/output. Forget, Input, and Output gates

C_t (state) and h_t (output and/or context) are vectors.

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$



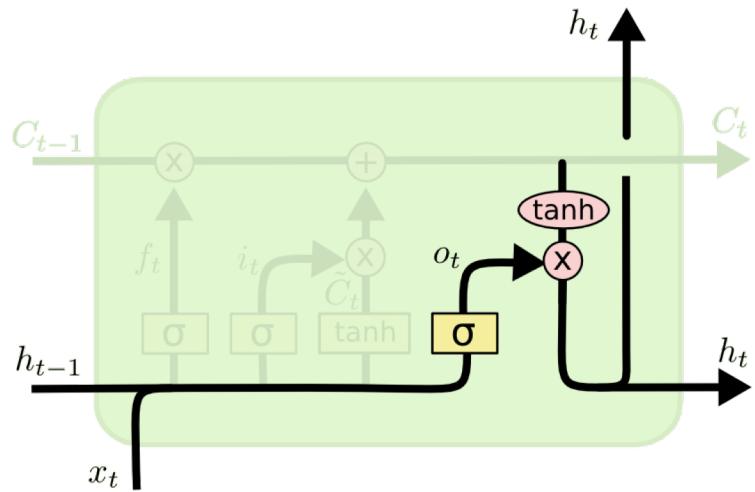
$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$



$$i_t = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

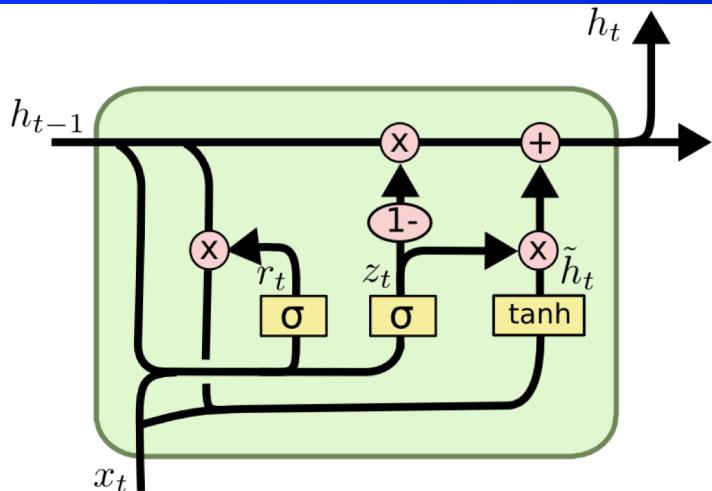
C_t can be maintained as long as needed. Only updated by forget and input gates, which are learned functions.



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

- Output gate above finishes full LSTM node
 - \tanh normalizes C to h/x scale (between -1/1) before output gate chooses what parts of state to pass on as output/context
- Gated Recurrent Unit (GRU) below combines C and h
 - r : reset gate – How much of context h_{t-1} to use in standard \tanh
 - z : update gate - Combines forget and input gates



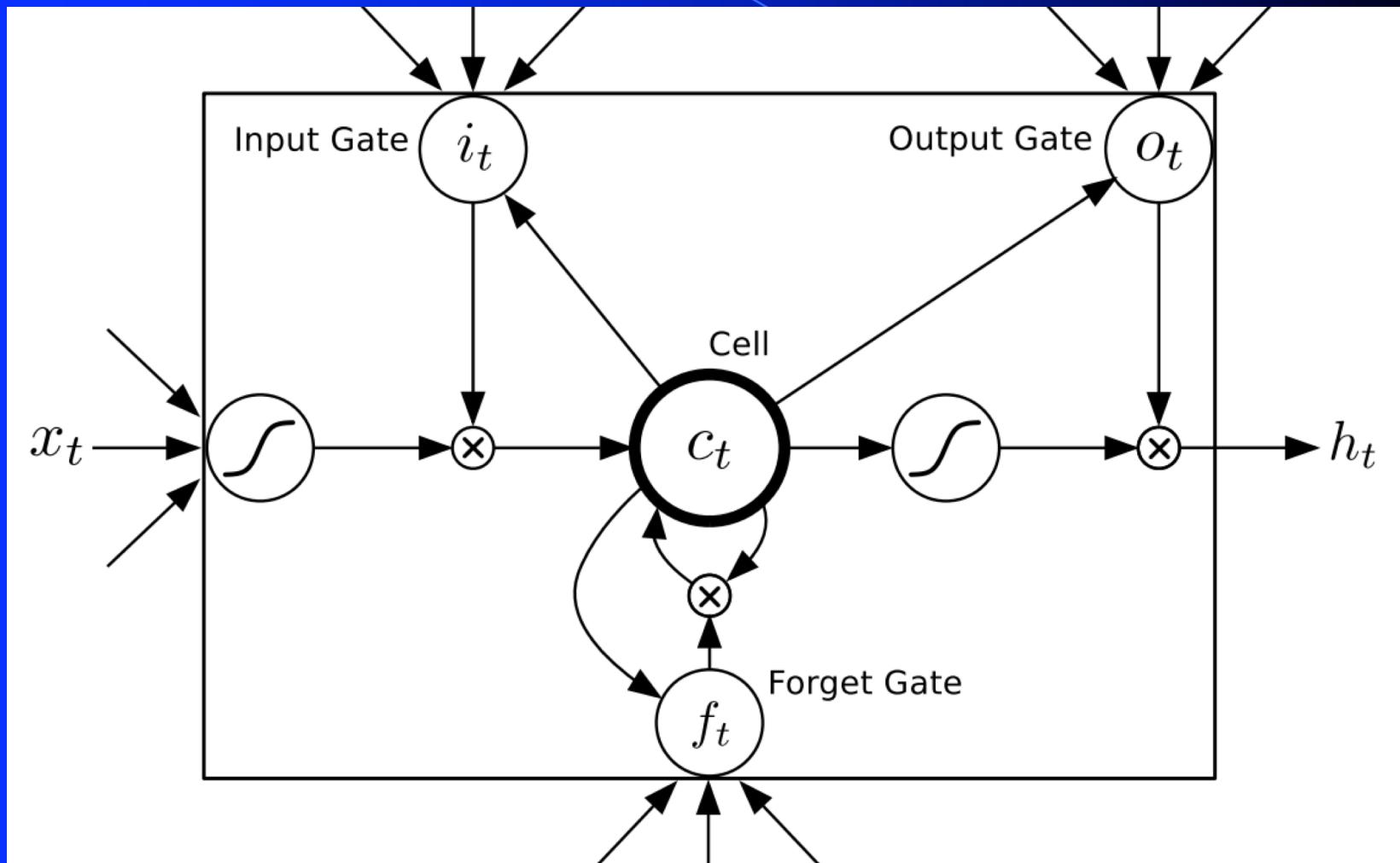
$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

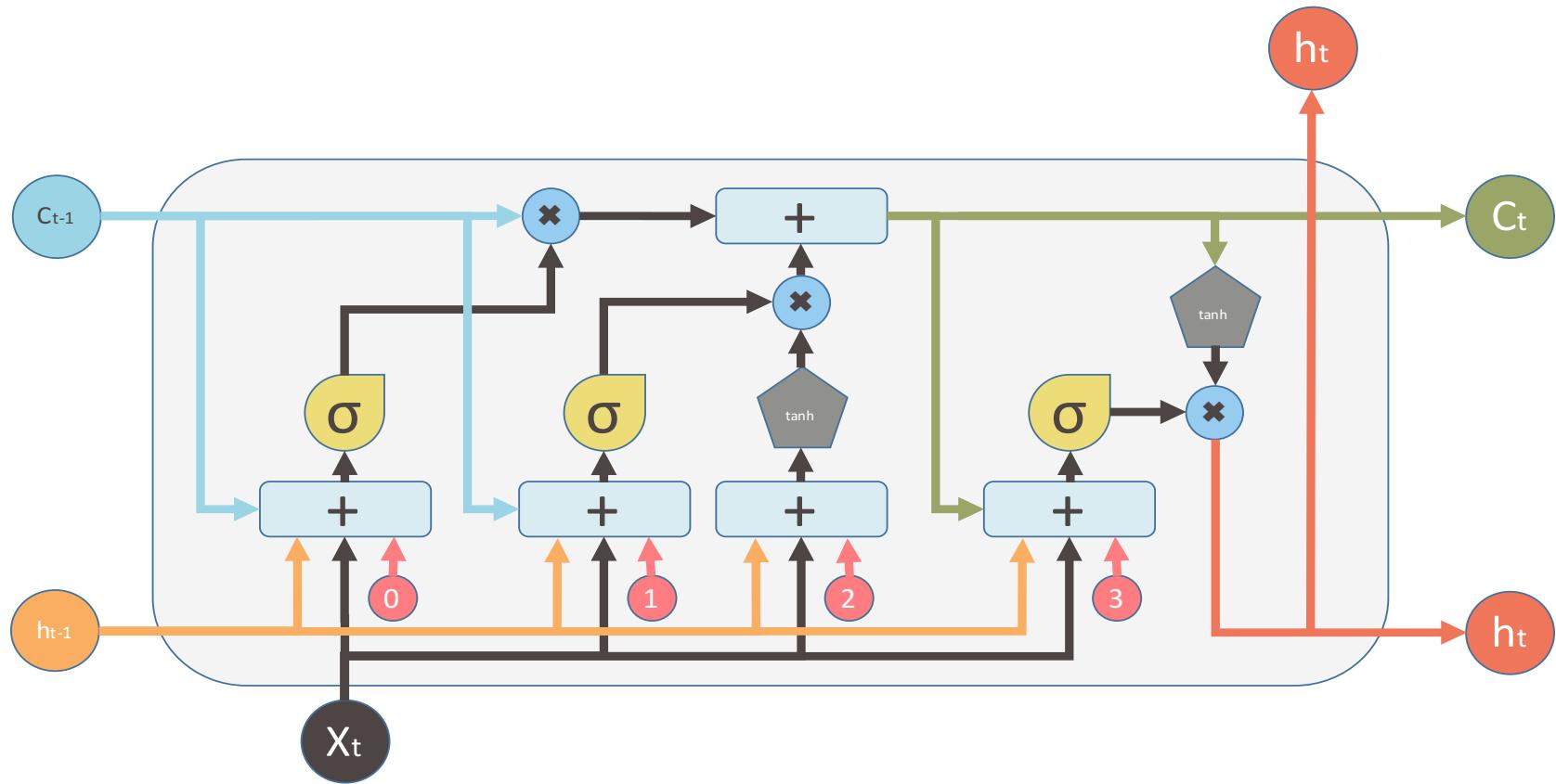
$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Basic Cell value has a self-feedback loop
Does not forget until told (stable gradient)





Inputs:

X_t Input vector

C_{t-1} Memory from previous block

h_{t-1} Output of previous block

outputs:

C_t Memory from current block

h_t Output of current block

Nonlinearities:

σ Sigmoid

tanh Hyperbolic tangent

Bias: 0

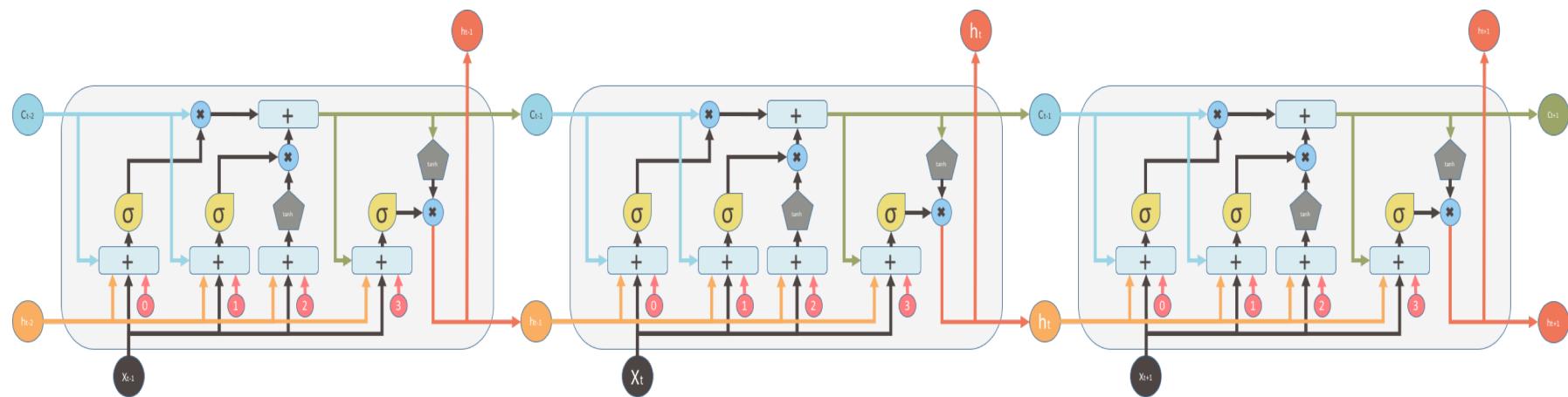
Vector operations:



Element-wise multiplication

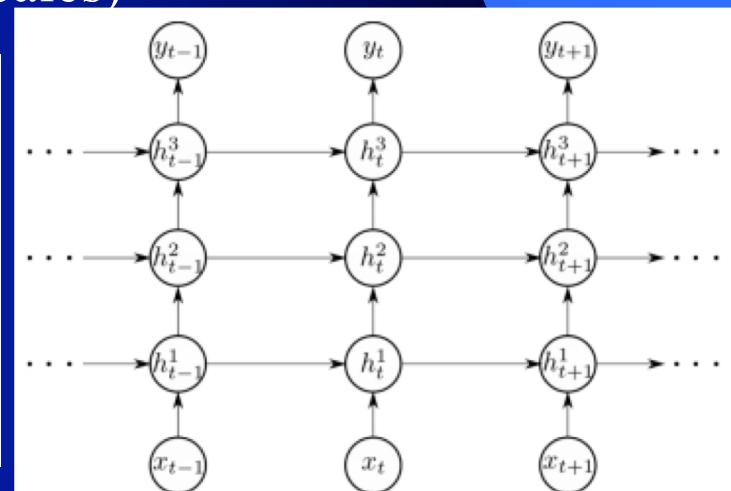
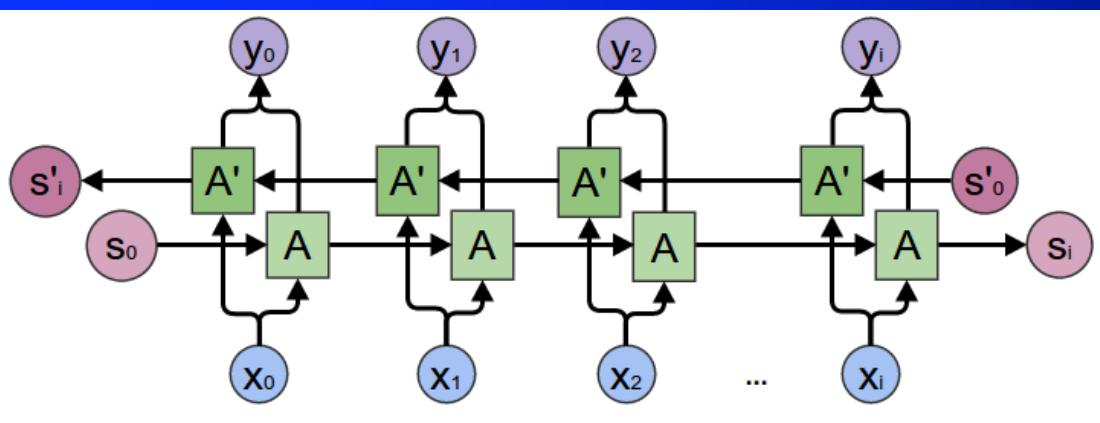


Element-wise Summation / Concatenation

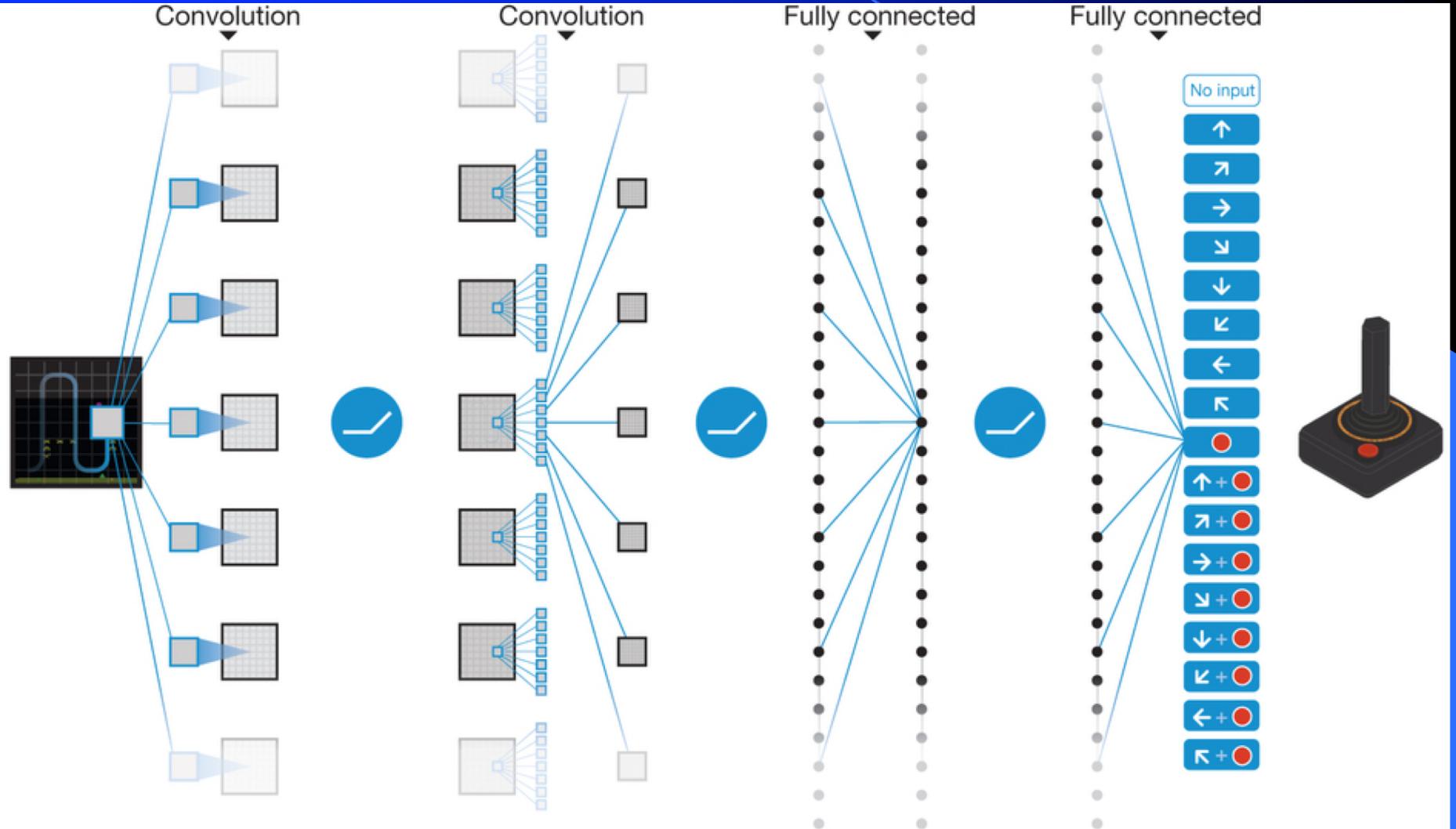


LSTM Variations

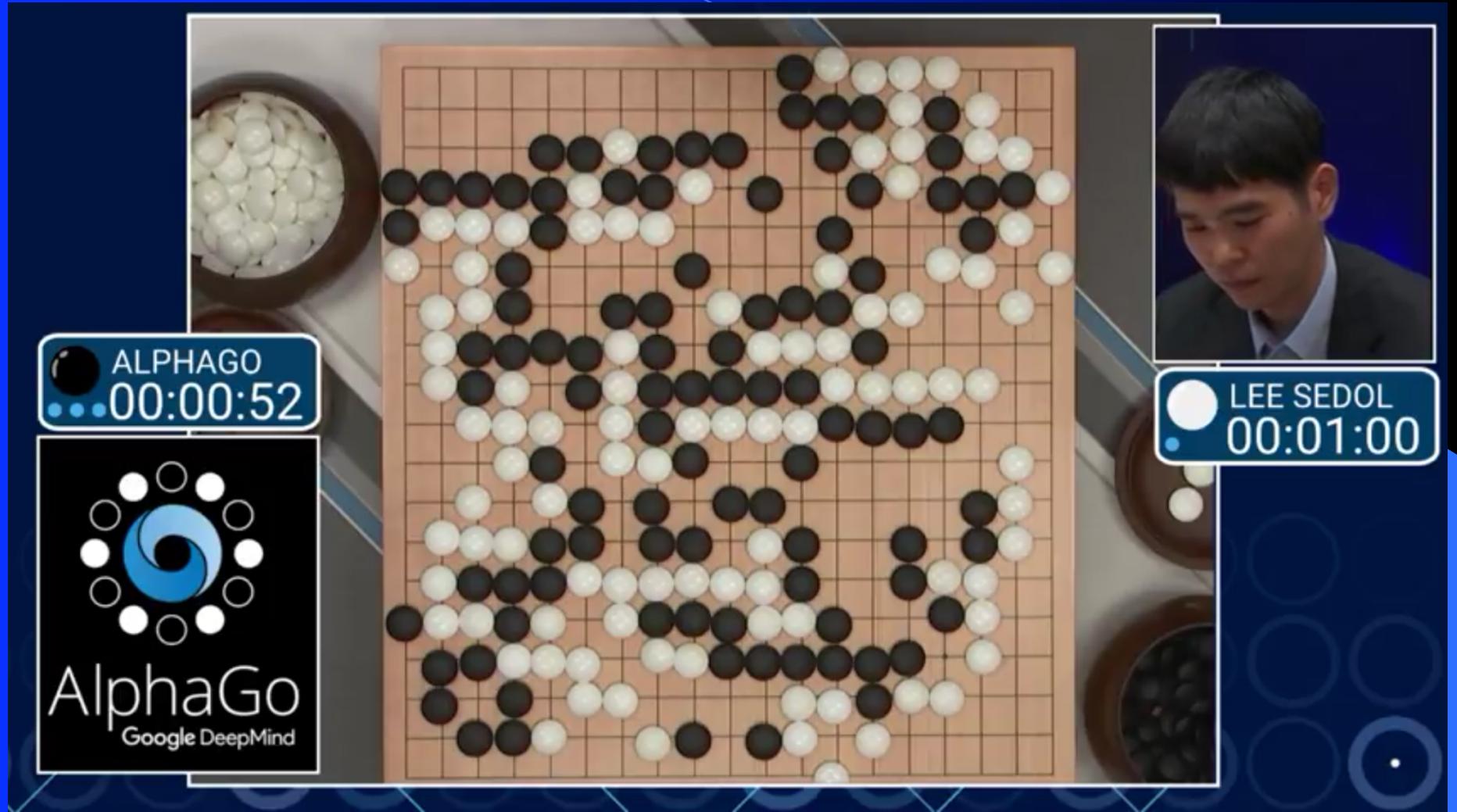
- Lots of node and structural variations
- Could replace any node in a deep network with LSTM node, but typically just used for sequential problems - RNN (time or space)
- Bidirectional LSTMS -Unlike conventional RNNs, bidirectional RNNs utilize both the previous and future context, by processing the data from two directions with two separate hidden layers. One layer processes the input sequence in the forward direction, while the other processes the input in the reverse direction. The output of current time step is then generated by combining both layer's hidden vector
- Stacked LSTMS - extra layers above (usually not many) can capture latent info (e.g. different time scales)



Deep Reinforcement Learning: Deep Q Network – 49 Classic Atari Games



AlphaGo - Google DeepMind



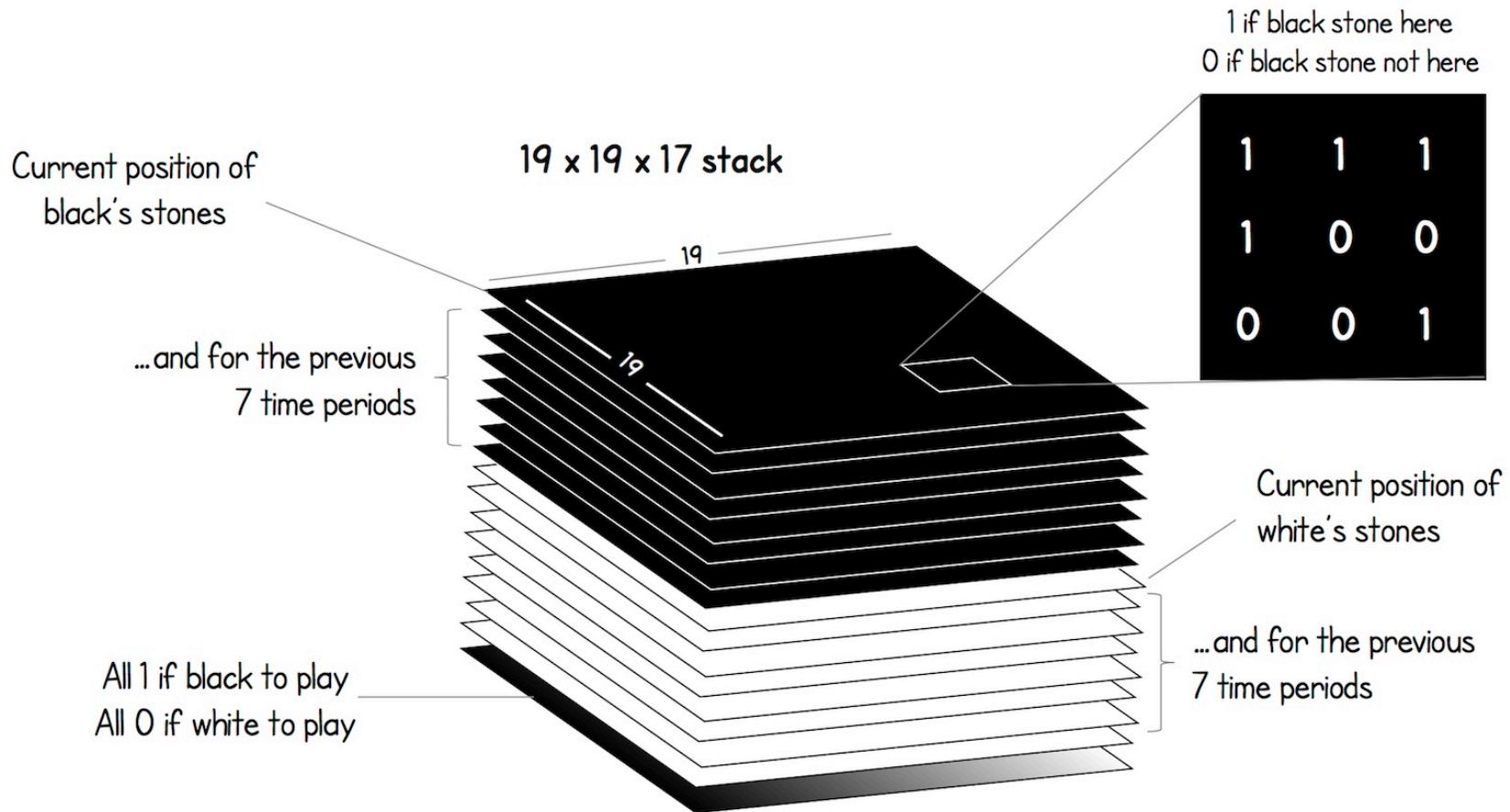
Alpha Go

- Reinforcement Learning with Deep Net learning the value and policy functions
- Challenges world Champion Lee Se-dol in March 2016
 - AlphaGo Movie – Netflix, check it out, fascinating man/machine interaction!
- AlphaGo Master (improved with more training) then beat top masters on-line 60-0 in Jan 2017
- 2017 – Alpha Go Zero
 - Alpha Go started by learning from 1000's of expert games before learning more on its own, and with lots of expert knowledge
 - Alpha Go Zero starts from zero (Tabula Rasa), just gets rules of Go and starts playing itself to learn how to play – not patterned after human play – More creative
 - Beat AlphaGo Master 100 games to 0 (after 3 days of playing itself)

Alpha Zero

- Alpha Zero (late 2017)
- Generic architecture for any board game
 - Compared to AlphaGo (2016 - earlier world champion with extensive background knowledge) and AlphaGo Zero (2017)
- No input other than rules and self-play, and not set up for any specific game, except different board input
- With no domain knowledge and starting from random weights, beats worlds best players and computer programs (which were specifically tuned for their games over many years)
 - Go – after 8 hours training (44 million games) beats AlphaGo Zero (which had beat AlphaGo 100-0) – 1000's of TPU's for training
 - AlphaGo had taken many months of human directed training
 - Chess – after 4 hours training beats Stockfish8 28-0 (+72 draws)
 - Doesn't pattern itself after human play
 - Shogi (Japanese Chess) – after 2 hours training beats Elmo

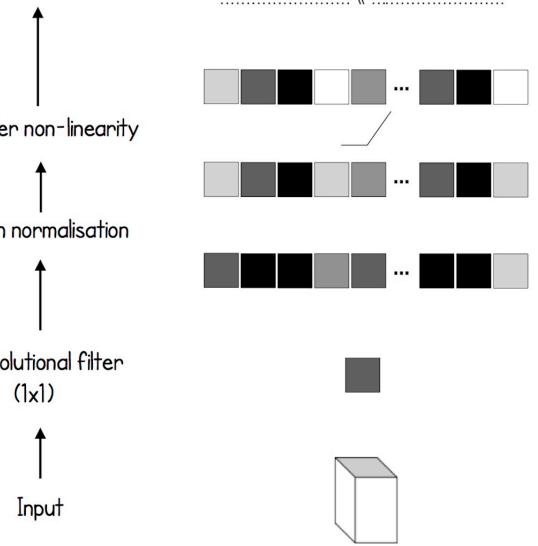
WHAT IS A 'GAME STATE'



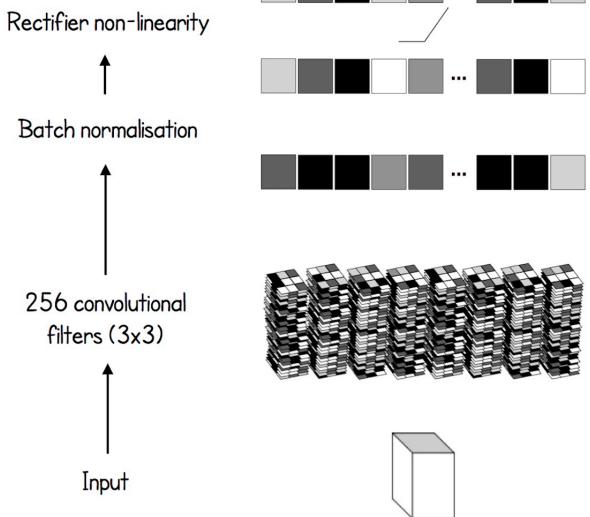
This stack is the input to the deep neural network

Hidden layer size 256

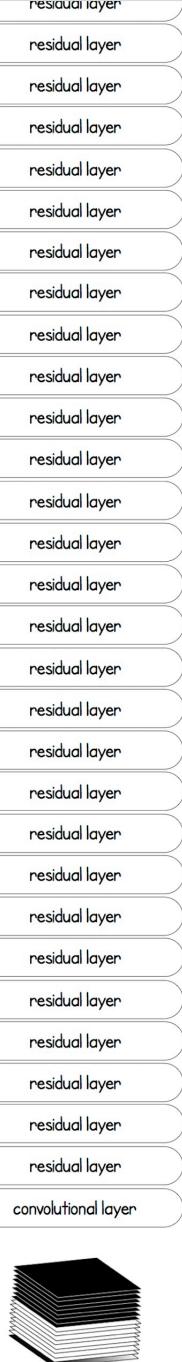
Fully connected layer



A convolutional layer

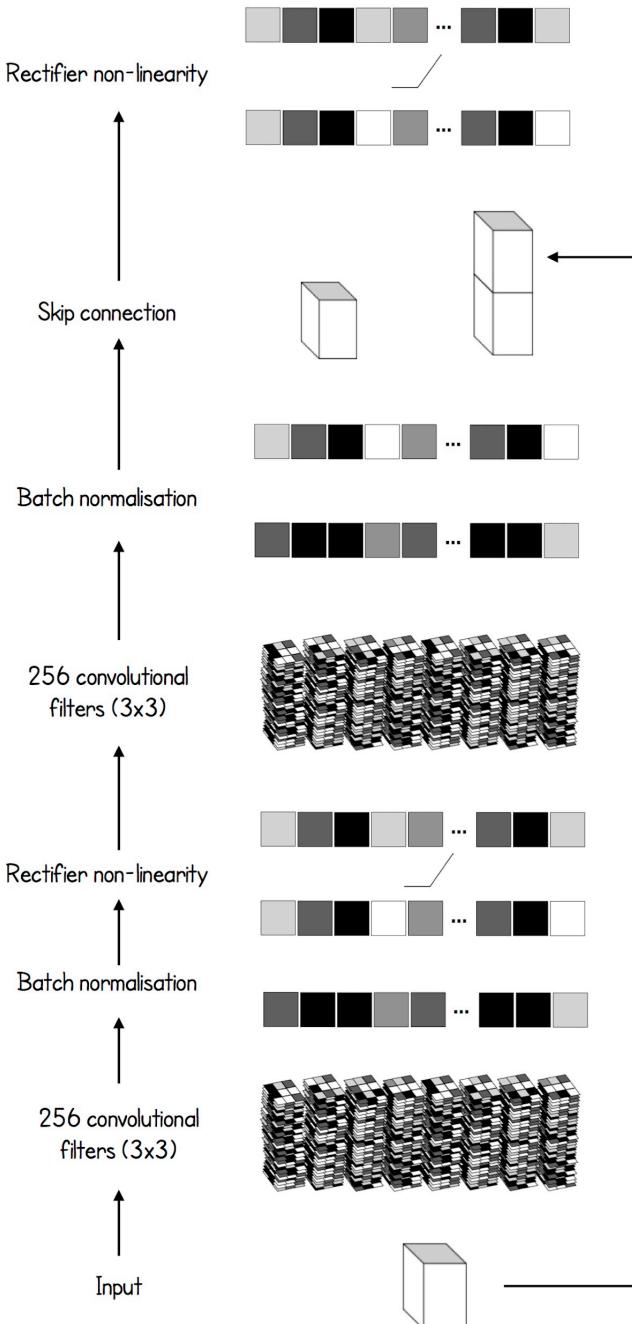


40 residual layers



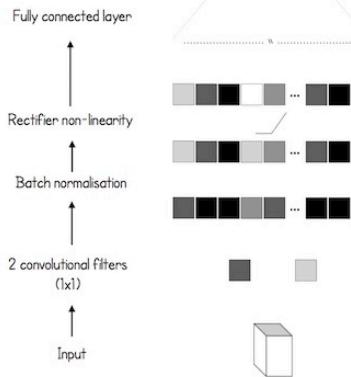
Input: The game

A residual layer



The policy head

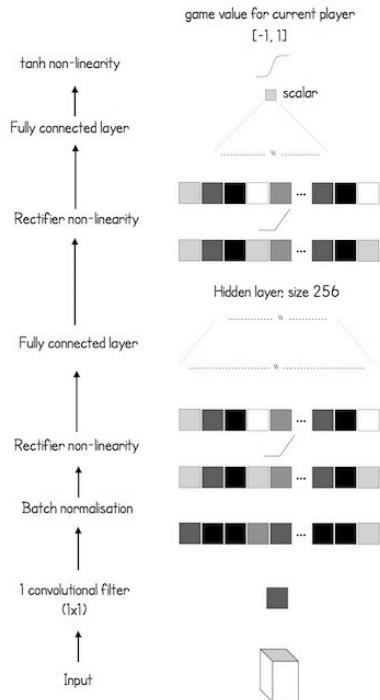
$19 \times 19 + 1$ (for pass)
move logit probabilities



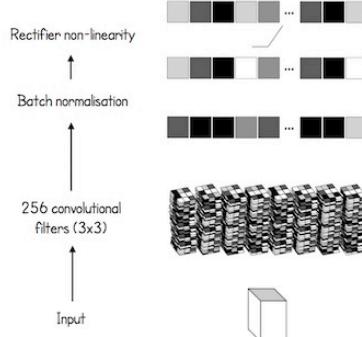
The network

value head policy head

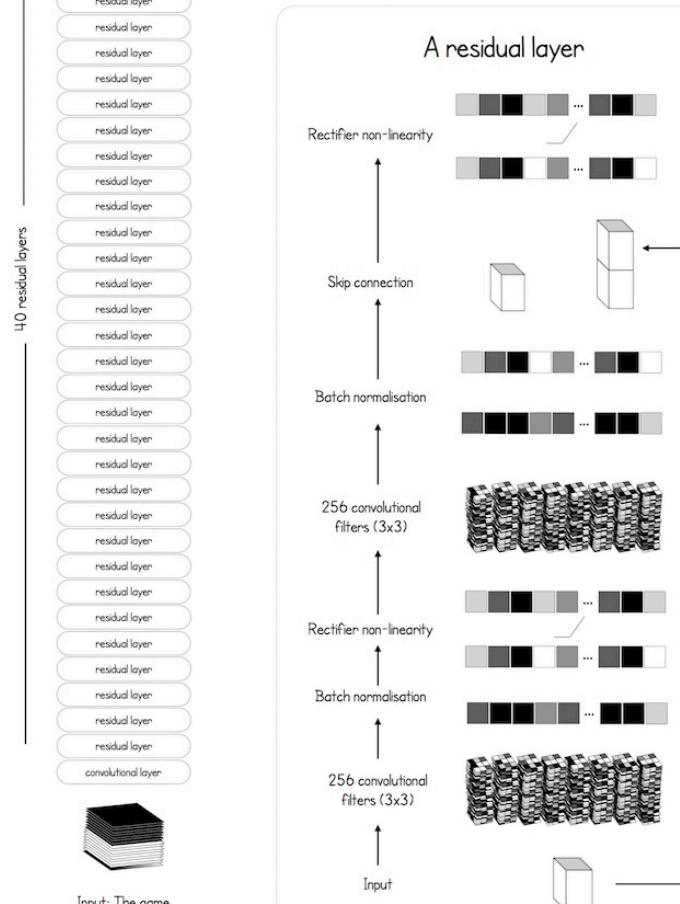
The value head



A convolutional layer



A residual layer



The network learns 'tabula rasa' (from a blank slate)

At no point is the network trained using human knowledge or expert moves

The training pipeline for AlphaGo Zero consists of three stages, executed in parallel

SELF PLAY

Create a 'training set'

The best current player plays 25,000 games against itself

See MCTS section to understand how AlphaGo Zero selects each move

At each move, the following information is stored



The game state
(see 'What is a Game State section')



The search probabilities
(from the MCTS)



The winner
(+1 if this player won, -1 if this player lost + added once the game has finished)

RETRAIN NETWORK

Optimise the network weights

A TRAINING LOOP

Sample a mini-batch of 2048 positions from the last 500,000 games

Retrain the current neural network on these positions

- The game states are the input (see 'Deep Neural Network Architecture')

Loss Function

Compares predictions from the neural network with the search probabilities and actual winner

$$\text{PREDICTIONS} \quad \begin{matrix} p \\ v \\ \text{Mean-squared error} \\ \text{Regularisation} \end{matrix} \quad \text{ACTUAL} \quad \begin{matrix} \pi \\ \text{Trophy icon} \end{matrix}$$

After every 1,000 training loops, evaluate the network

EVALUATE NETWORK

Test to see if the new network is stronger

Play 400 games between the latest neural network and the current best neural network

Both players use MCTS to select their moves, with their respective neural networks to evaluate leaf nodes

Latest player must win 55% of games to be declared the new best player



WHAT IS A 'GAME STATE'

Current position of black's stones

19 x 19 x 17 stack

| |
|---------------------------|
| 1 if black stone here |
| 0 if black stone not here |
| 1 |
| 1 |
| 0 |
| 0 |
| 0 |
| 1 |

Current position of white's stones

19

...and for the previous 7 time periods

All 1 if black to play

All 0 if white to play

Current position of white's stones

19

...and for the previous 7 time periods

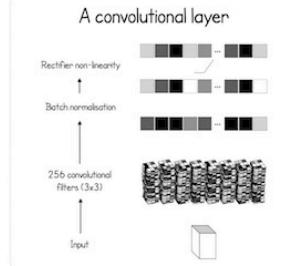
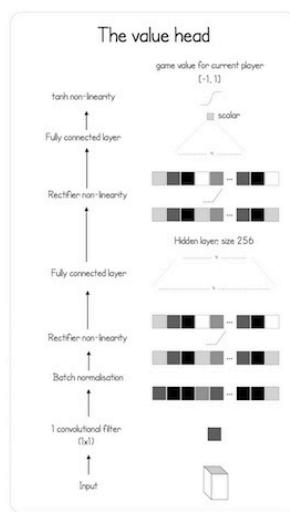
This stack is the input to the deep neural network

THE DEEP NEURAL NETWORK ARCHITECTURE

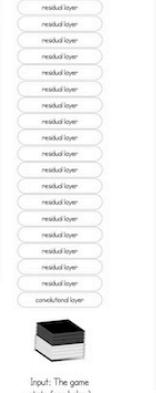
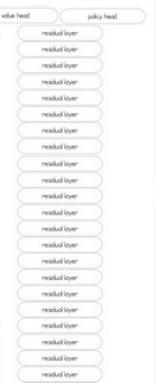
How AlphaGo Zero assesses new positions

The network learns 'tabula rasa' (from a blank slate)

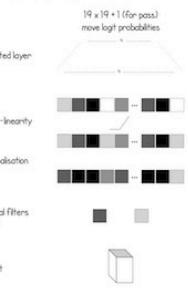
At no point is the network trained using human knowledge or expert moves



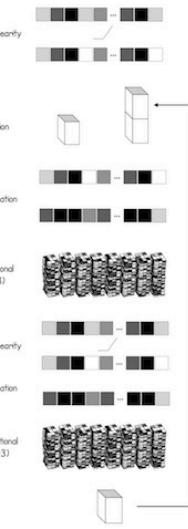
The network



The policy head

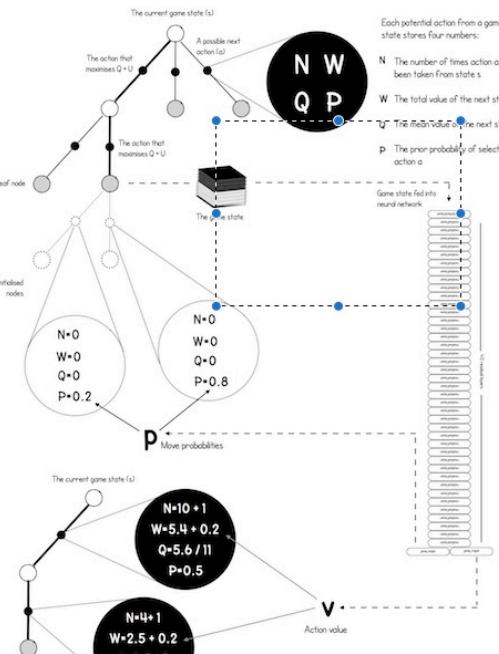


A residual layer



MONTE CARLO TREE SEARCH (MCTS)

How AlphaGo Zero chooses its next move



...then select a move

After 1,600 simulations, the move can either be chosen:

Deterministically (for competitive play)

Choose the action from the current state with greatest N

Stochastically (for exploratory play)

Choose the action from the current state from the distribution

$$\pi = N^{-\frac{1}{\tau}}$$

where τ is a temperature parameter controlling exploration

First, run the following simulation 1,600 times...

Start at the root node of the tree (the current game state)

1. Choose the action that maximises...

$$Q + U$$

A function of P and N that increases if an action hasn't been explored much, relative to the other actions, or if the prior probability of the action is high

Early on in the simulation, U dominates (more exploration), but later Q is more important (less exploration)

2. Continue until a leaf node is reached

The game state of the leaf node is passed into the neural network, which outputs predictions about two things:

P Move probabilities

V Value of the state (for the current player)

The move probabilities p are attached to the new feasible actions from the leaf node

3. Backup previous edges

Each edge that was traversed to get to the leaf node is updated as follows:

$$N \rightarrow N + 1$$

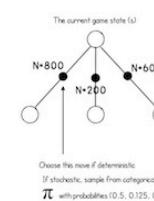
$$W \rightarrow W + v$$

$$Q = W / N$$

Other points

The sub-tree from the chosen move is retained for calculating subsequent moves

The rest of the tree is discarded



Choose this move if deterministic
If stochastic, sample from categorical distribution
 π with probabilities (0.5, 0.125, 0.375)

Conclusion

- Much recent excitement, still much to be discovered
- Impressive results
- More work needed to understand how and why deep learning works so well – How deep should we go?
- Potential for significant improvements
- Works well in structured/Markovian spaces - CNNs, etc.
 - Important research question: To what extent can we use Deep Learning in arbitrary feature spaces?
 - Recent deep learning with supervised learning and extended BP approaches show potential in this area