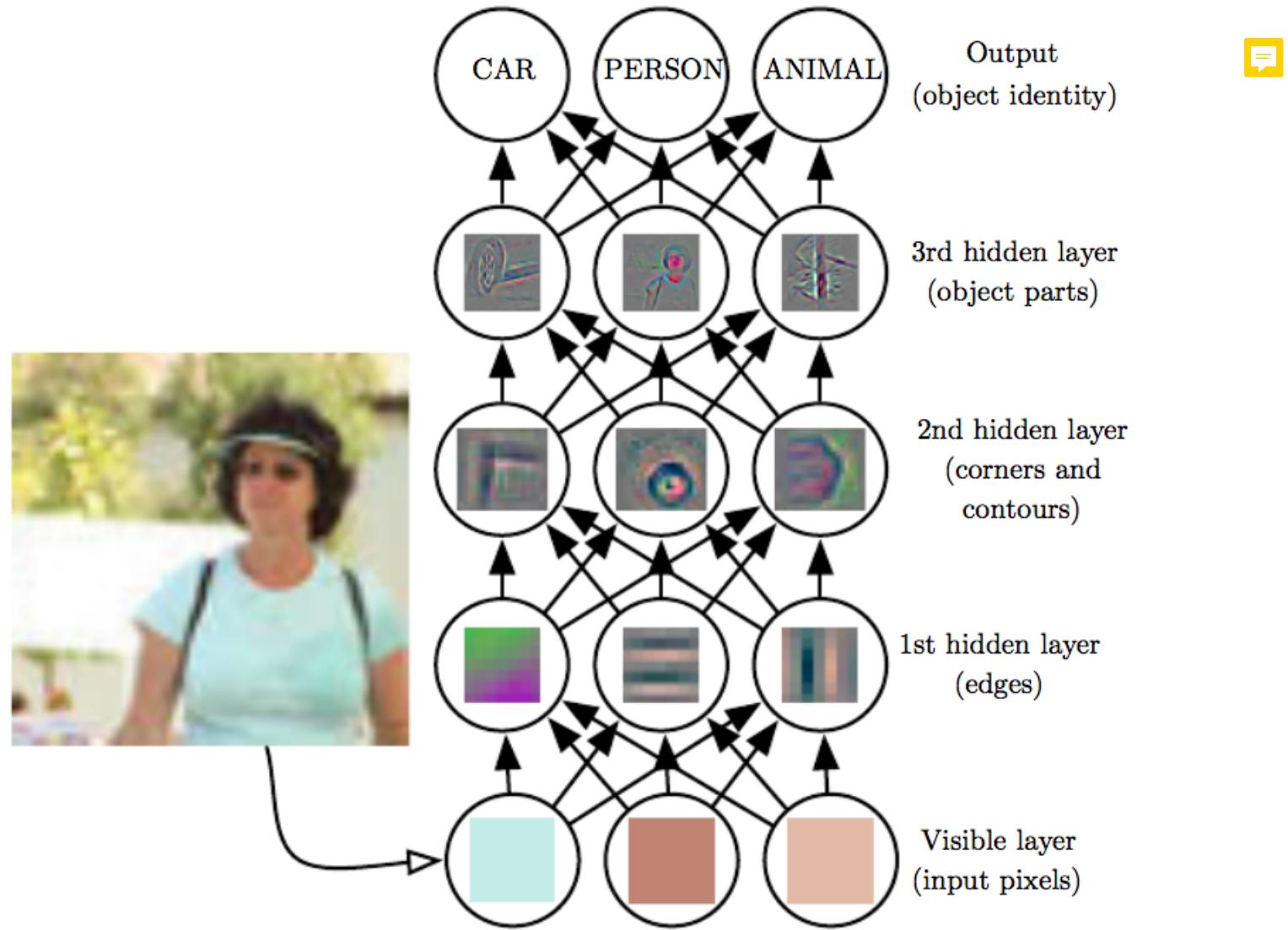
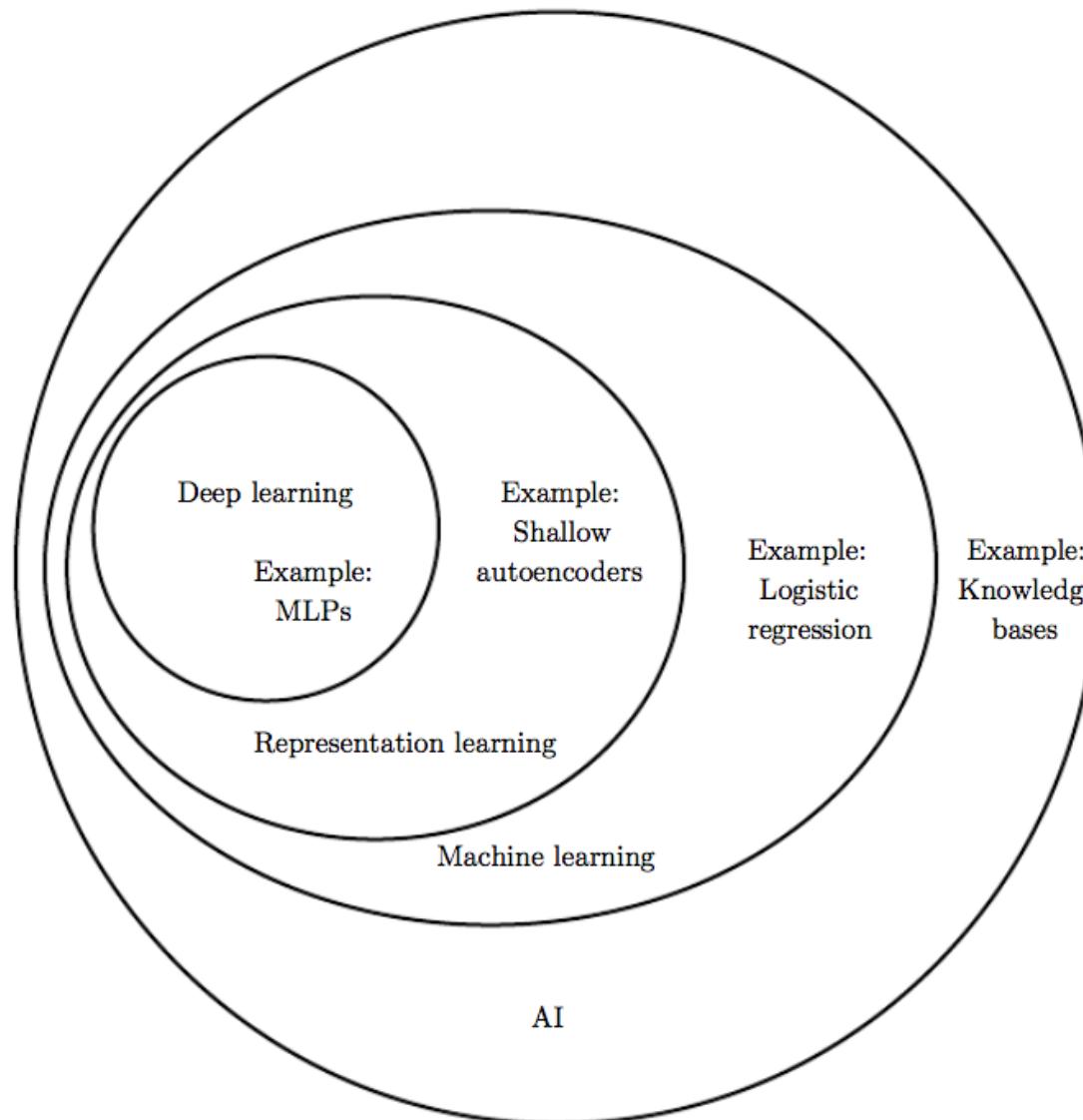


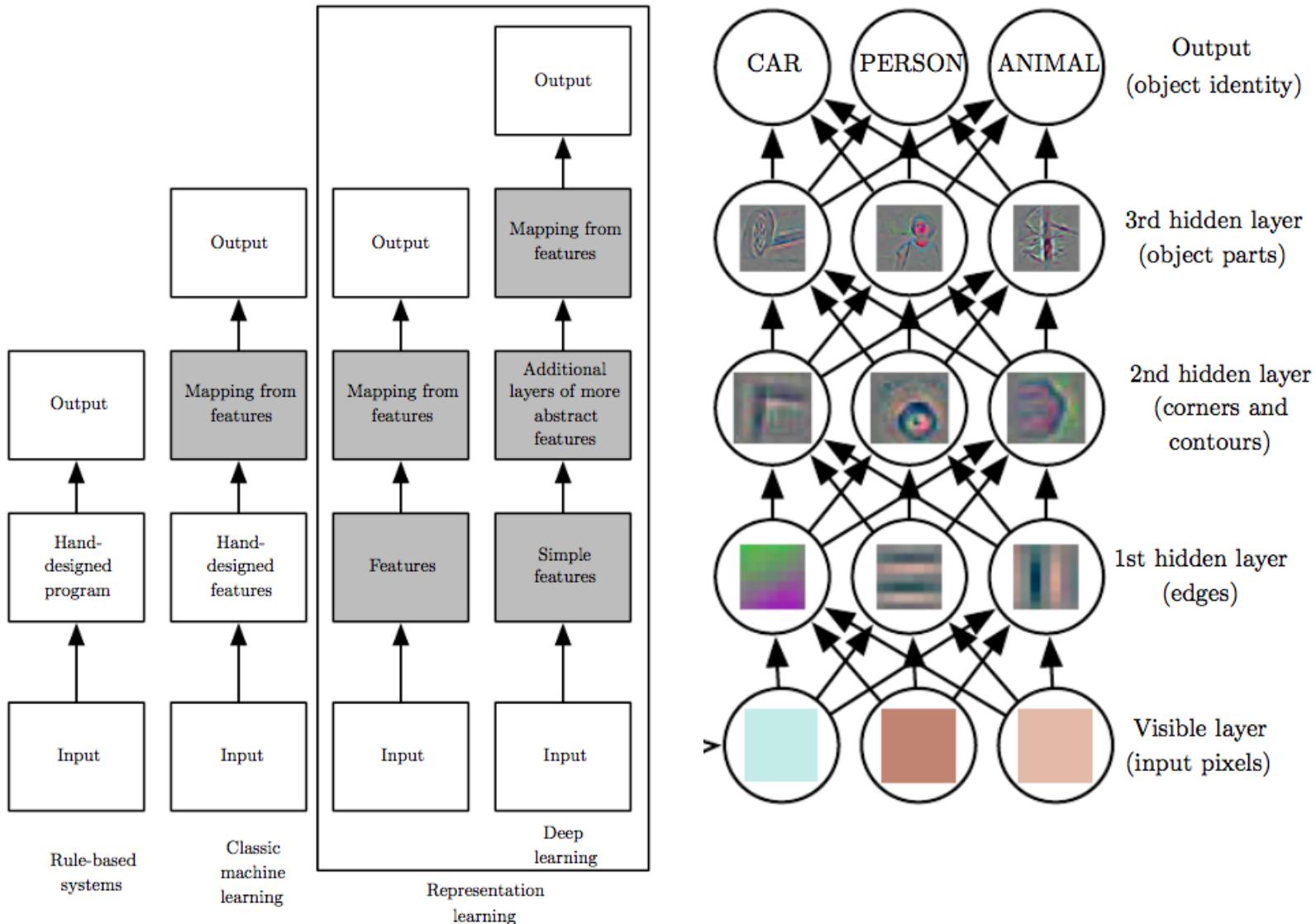
# Deep Learning



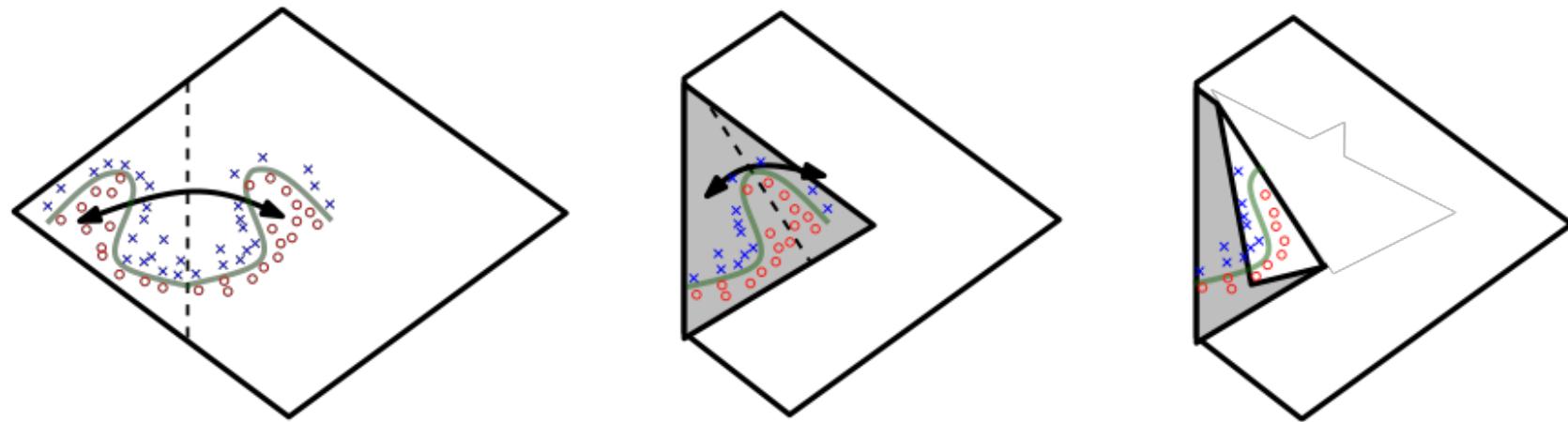
# Deep Learning and AI



# Deep Learning and Representation



# Deep Learning and Representation



Exponential Representation Advantage of Depth

# Why Deep Learning ?

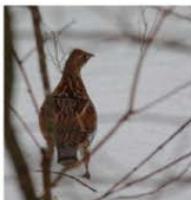
## ImageNet Large Scale Visual Recognition Challenge (ILSVRC)



flamingo



cock



ruffed grouse



quail



partridge

...



Egyptian cat



Persian cat



Siamese cat



tabby



lynx

...



dalmatian



keeshond



miniature schnauzer



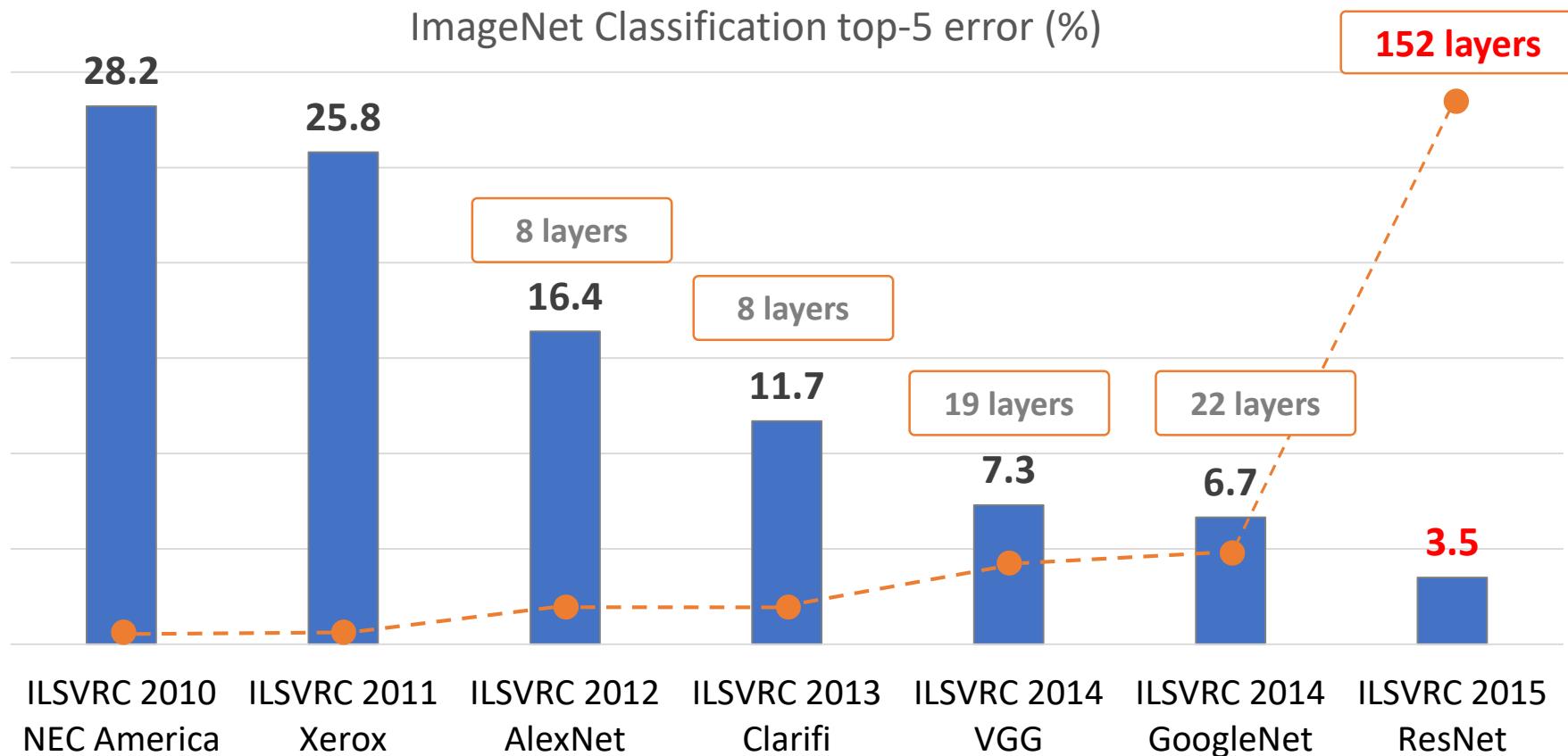
standard schnauzer



giant schnauzer

...

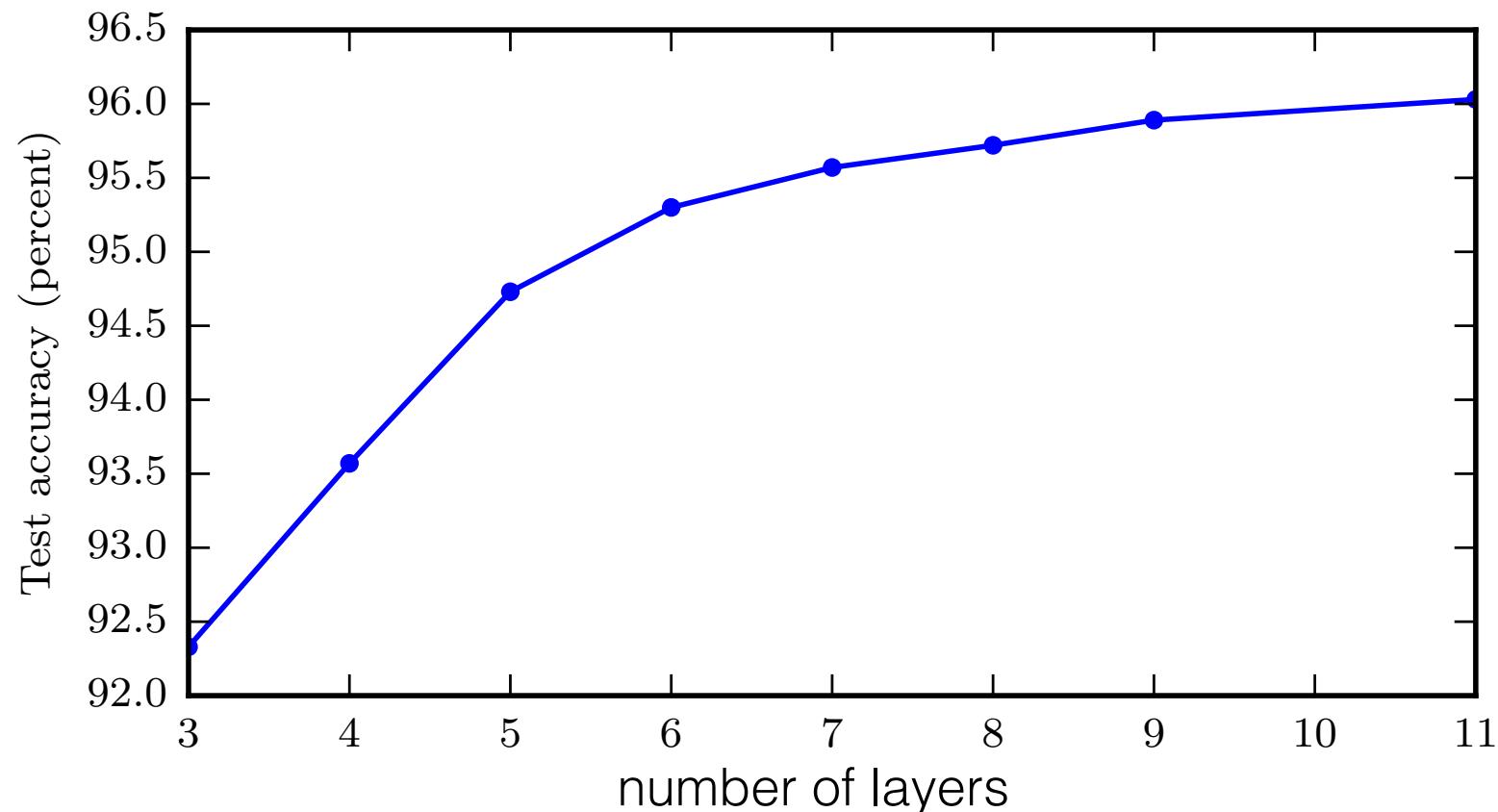
# Why Deep Learning ?



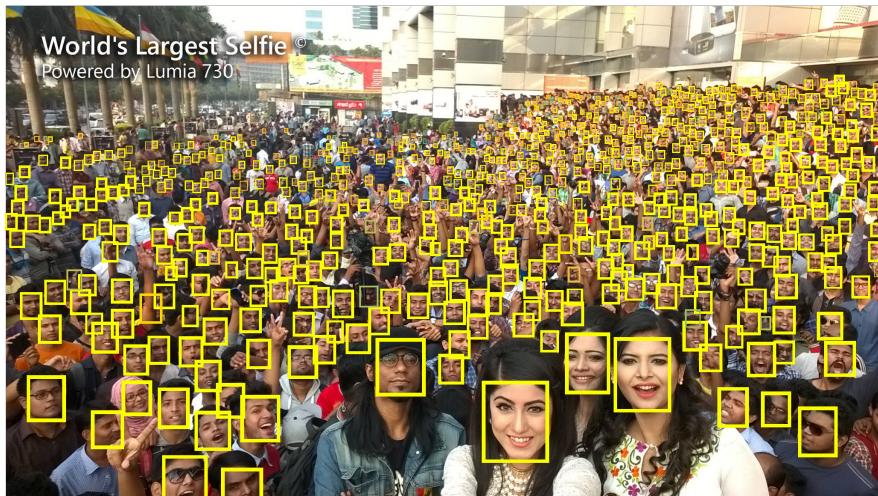
# Why Deep Learning ?



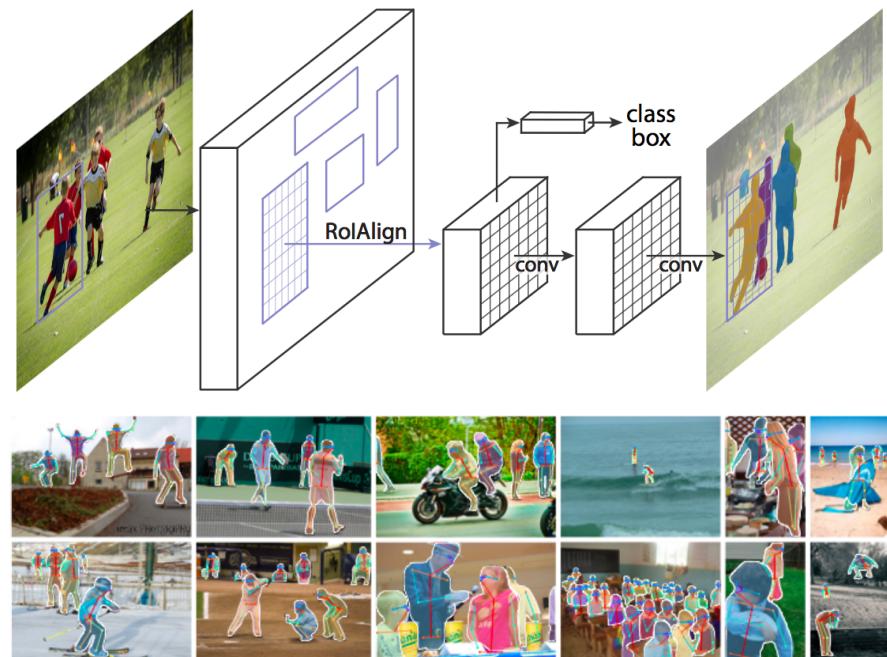
Better Generalization with Greater Depth



# Why Deep Learning ?



Hu et al.: Finding tiny faces. 2016



He et al.: Mask R-CNN. 2017

# Why Deep Learning ?



Task	Test set	Metric	Best non-neural	Best neural	Source
<b>Machine Translation</b>	Enu-deu newstest16	BLEU	31.4	34.8	<a href="http://matrix.statmt.org">http://matrix.statmt.org</a>
	Deu-enu newstest16	BLEU	35.9	39.9	<a href="http://matrix.statmt.org">http://matrix.statmt.org</a>
<b>Sentiment Analysis</b>	Stanford sentiment bank	5-class Accuracy	71.0	80.7	<a href="#">Socher+ 13</a>
<b>Question Answering</b>	WebQuestions test set	F1	39.9	52.5	<a href="#">Yih+ 15</a>
<b>Entity Linking</b>	Bing Query Entity Linking set	AUC	72.3	78.2	<a href="#">Gao+ 14b</a>
<b>Image Captioning</b>	COCO 2015 challenge	Turing test pass%	25.5	32.2	<a href="#">Fang+ 15</a>
<b>Sentence compression</b>	Google 10K dataset	F1	0.75	0.82	<a href="#">Filipova+ 15</a>
<b>Response Generation</b>	Sordoni dataset	BLEU-4	3.98	5.82	<a href="#">Li+ 16a</a>

Task	Test set	Metric	Best non-neural	Best neural	Source
<b>POS tagging</b>	PTB section 23	F1	97.17	97.78	<a href="#">Andor+ 16</a>
<b>Syntactic Parsing</b>	PTB section 23	F1	90.1	93.3	<a href="#">Dyer+ 16</a>
<b>Dependency parsing</b>	PTB section 23	F1	93.22	94.61	<a href="#">Andor+ 16</a>
<b>CCG parsing</b>	CCGBank test	F1	85.2	88.7	<a href="#">Lee+ 16</a>
<b>Inference (NLI)</b>	Stanford NLI corpus	Accuracy	78.2	88.3	<a href="#">Chen+ 16</a>

# Deep Learning



- Neural Networks with many hidden layers (deep networks)
  - Insufficient depth can hurt
  - The brain has a deep architecture
  - Cognitive processes seem deep
- New types of units, Dropout, Batch Normalization, Deep Convolutional Networks, Probabilistic Models
- Use of unsupervised learning (autoencoders) for incremental training (one layer at a time) and pre-training in general
- Computational tools
  - Computational graph
- Brute force training using GPUs and accelerating cards

# Deep Learning

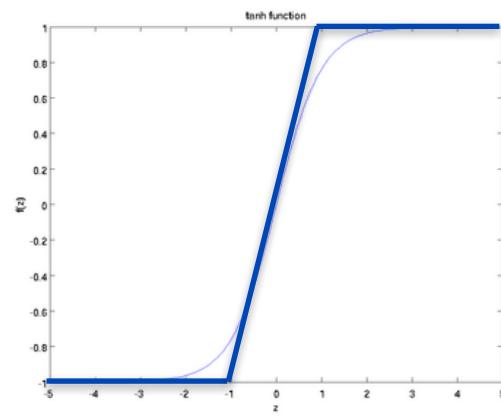
- Neural Networks with many hidden layers (deep networks)
  - Insufficient depth can hurt
  - The brain has a deep architecture
  - Cognitive processes seem deep
- **New types of units**, Dropout, Batch Normalization, Deep Convolutional Networks, Probabilistic Models
- Use of unsupervised learning (autoencoders) for incremental training (one layer at a time) and pre-training in general
- Computational tools
  - Computational graph
- Brute force training using GPUs and accelerating cards



# New Types of Units for Better Gradient Flow

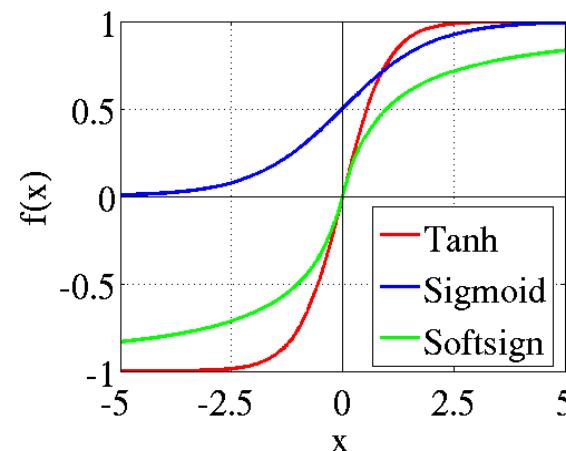
hard tanh

$$\text{HardTanh}(x) = \begin{cases} -1 & \text{if } x < -1 \\ x & \text{if } -1 \leq x \leq 1 \\ 1 & \text{if } x > 1 \end{cases}$$



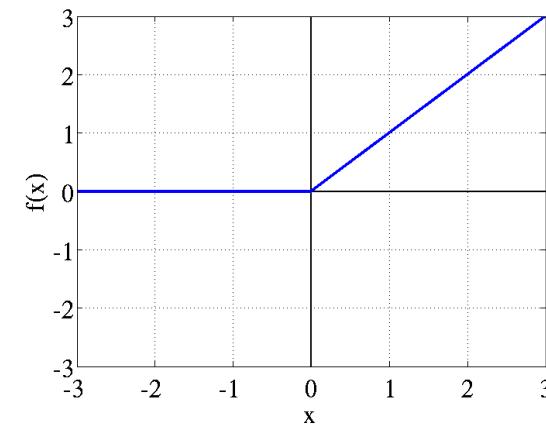
soft sign

$$\text{softsign}(z) = \frac{a}{1+|a|}$$



rectified linear (ReLU)

$$\text{rect}(z) = \max(z, 0)$$



other types of units for Convolution Neural Networks, will see later...



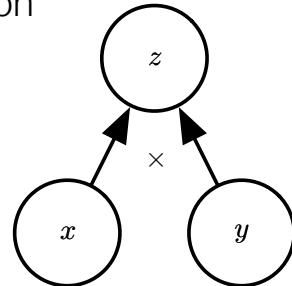
# Deep Learning

- Neural Networks with many hidden layers (deep networks)
  - Insufficient depth can hurt
  - The brain has a deep architecture
  - Cognitive processes seem deep
- New types of units, Dropout, Batch Normalization, Deep Convolutional Networks, Probabilistic Models
- Use of unsupervised learning (autoencoders) for incremental training (one layer at a time) and pre-training in general
- Computational tools
  - Computational graph
- Brute force training using GPUs and accelerating cards



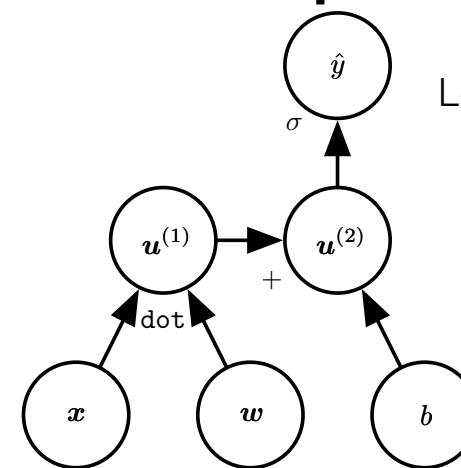
# Computational Graph

Multiplication



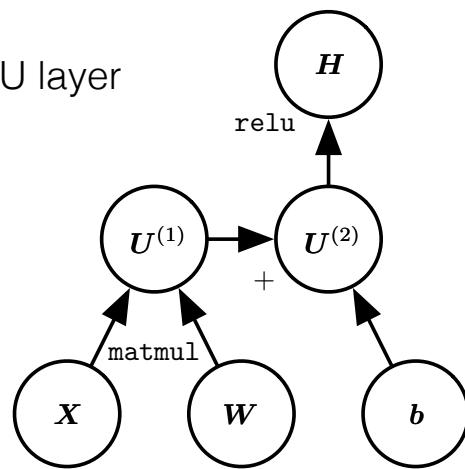
(a)

Logistic regression



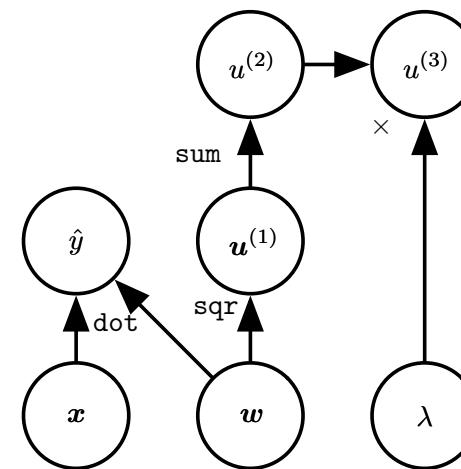
(b)

ReLU layer



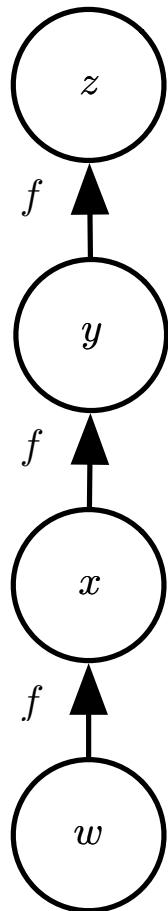
(c)

Linear regression  
and weight decay

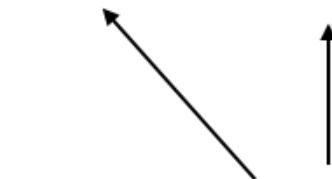


(d)

# Computational Graph

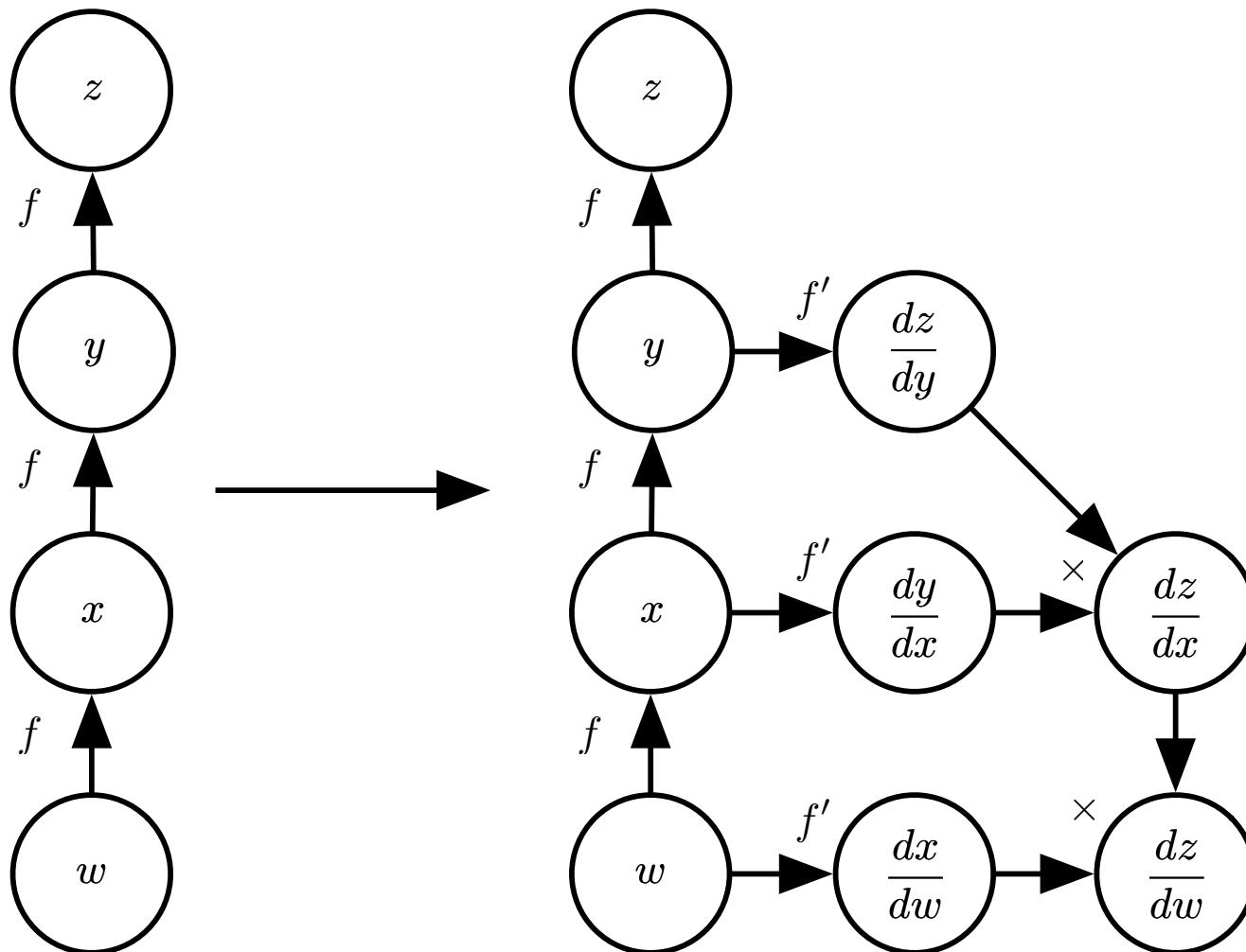


$$\begin{aligned}\frac{\partial z}{\partial w} &= \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} \\ &= f'(y) f'(x) f'(w) \\ &= f'(f(f(w))) f'(f(w)) f'(w)\end{aligned}$$



Back-prop avoids computing this twice

# Symbol-to-symbol Differentiation



# Loss Function as Computational Graph

