

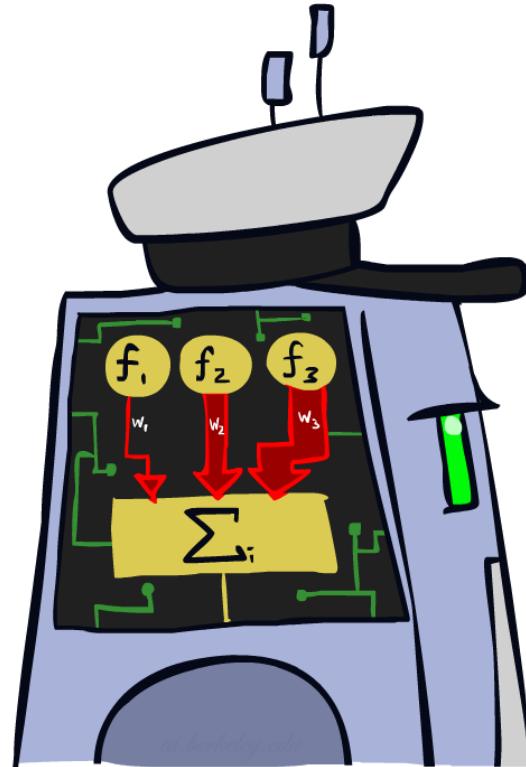
CS 188: Artificial Intelligence

Deep Learning I

Instructors: Pieter Abbeel & Anca Dragan --- University of California, Berkeley

[These slides were created by Dan Klein, Pieter Abbeel, Anca Dragan for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at <http://ai.berkeley.edu>.]

But First: Linear Classifiers



Feature Vectors

x

$f(x)$

y

Hello,
Do you want free print
cartridges? Why pay more
when you can get them
ABSOLUTELY FREE! Just

$\begin{cases} \# \text{ free} & : 2 \\ \text{YOUR_NAME} & : 0 \\ \text{MISSPELLED} & : 2 \\ \text{FROM_FRIEND} & : 0 \\ \dots \end{cases}$

SPAM
or
+

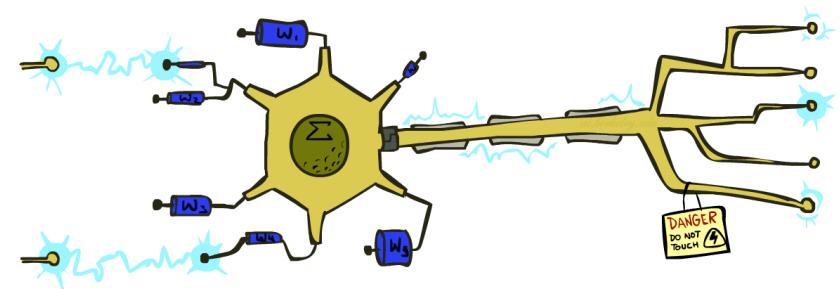
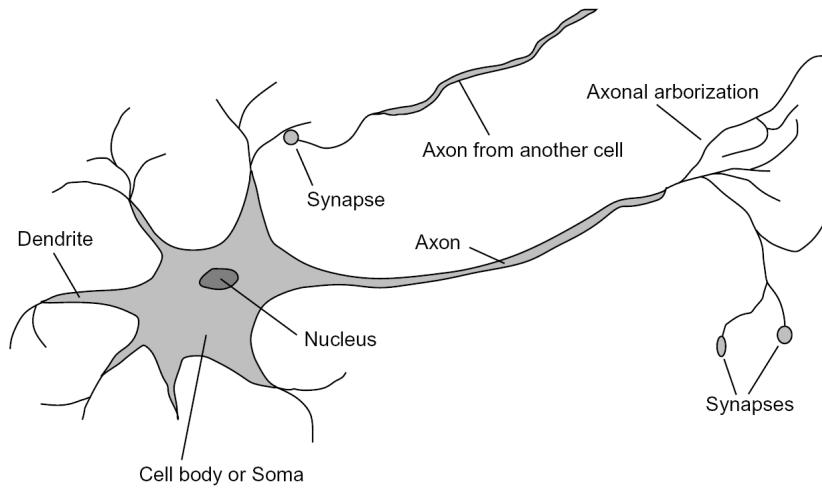


$\begin{cases} \text{PIXEL-7,12} & : 1 \\ \text{PIXEL-7,13} & : 0 \\ \dots \\ \text{NUM_LOOPS} & : 1 \\ \dots \end{cases}$

“2”

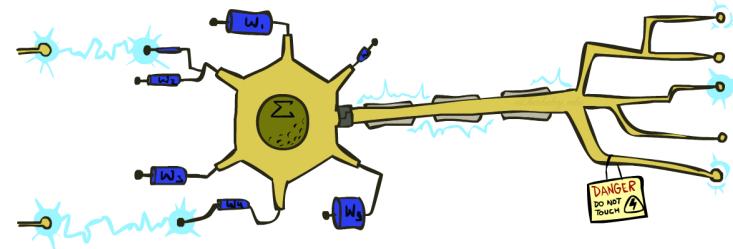
Some (Simplified) Biology

- Very loose inspiration: human neurons



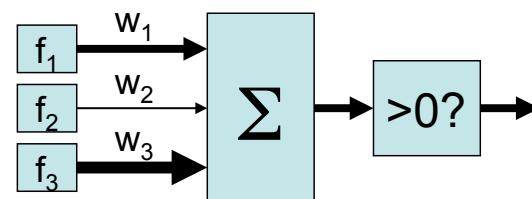
Linear Classifiers

- Inputs are **feature values**
- Each feature has a **weight**
- Sum is the **activation**

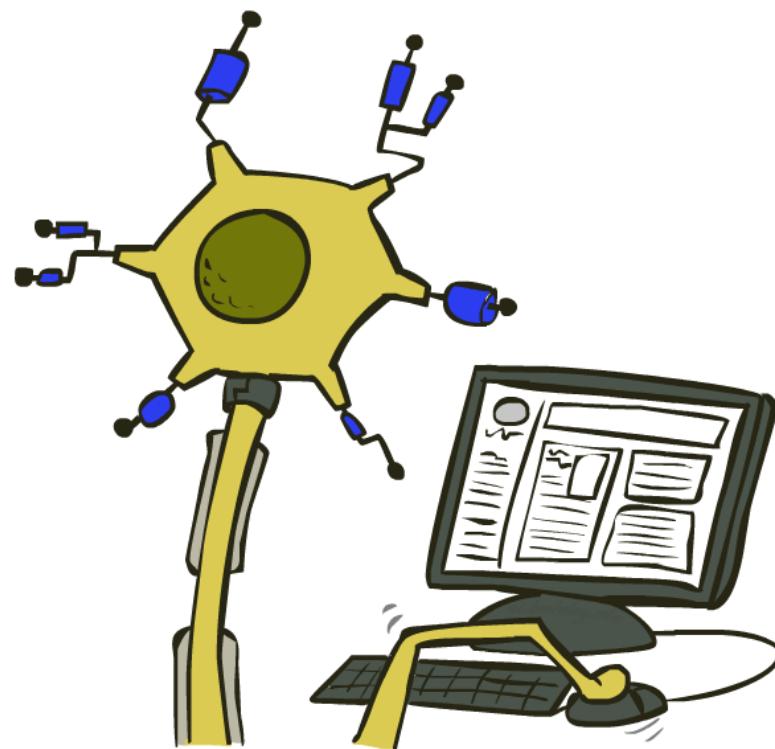


$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

- If the activation is:
 - Positive, output +1
 - Negative, output -1



Web Search

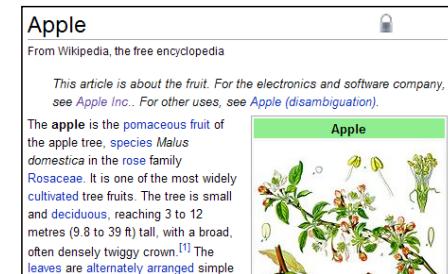


Extension: Web Search

- Information retrieval:
 - Given information needs, produce information
 - Includes, e.g. web search, question answering, and classic IR

- Web search: not exactly classification, but rather ranking

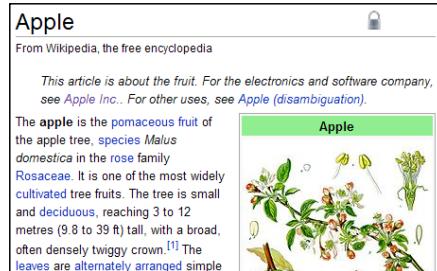
$x = \text{"Apple Computers"}$



Feature-Based Ranking

x = “Apple Computer”

$f(x,$



) = [0.3 5 0 0 ...]

$f(x,$



) = [0.8 4 2 1 ...]

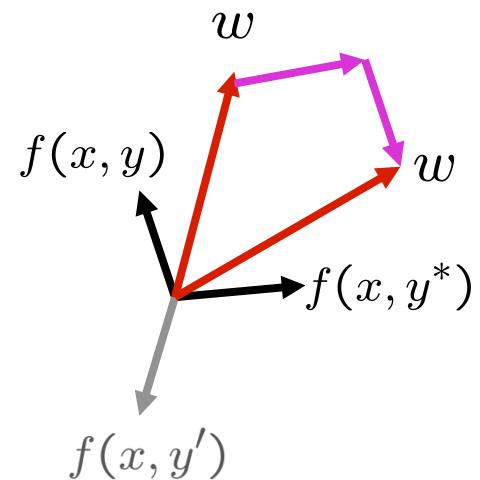
Perceptron for Ranking

- Inputs x
- Candidates y
- Many feature vectors: $f(x, y)$
- One weight vector: w
 - Prediction:

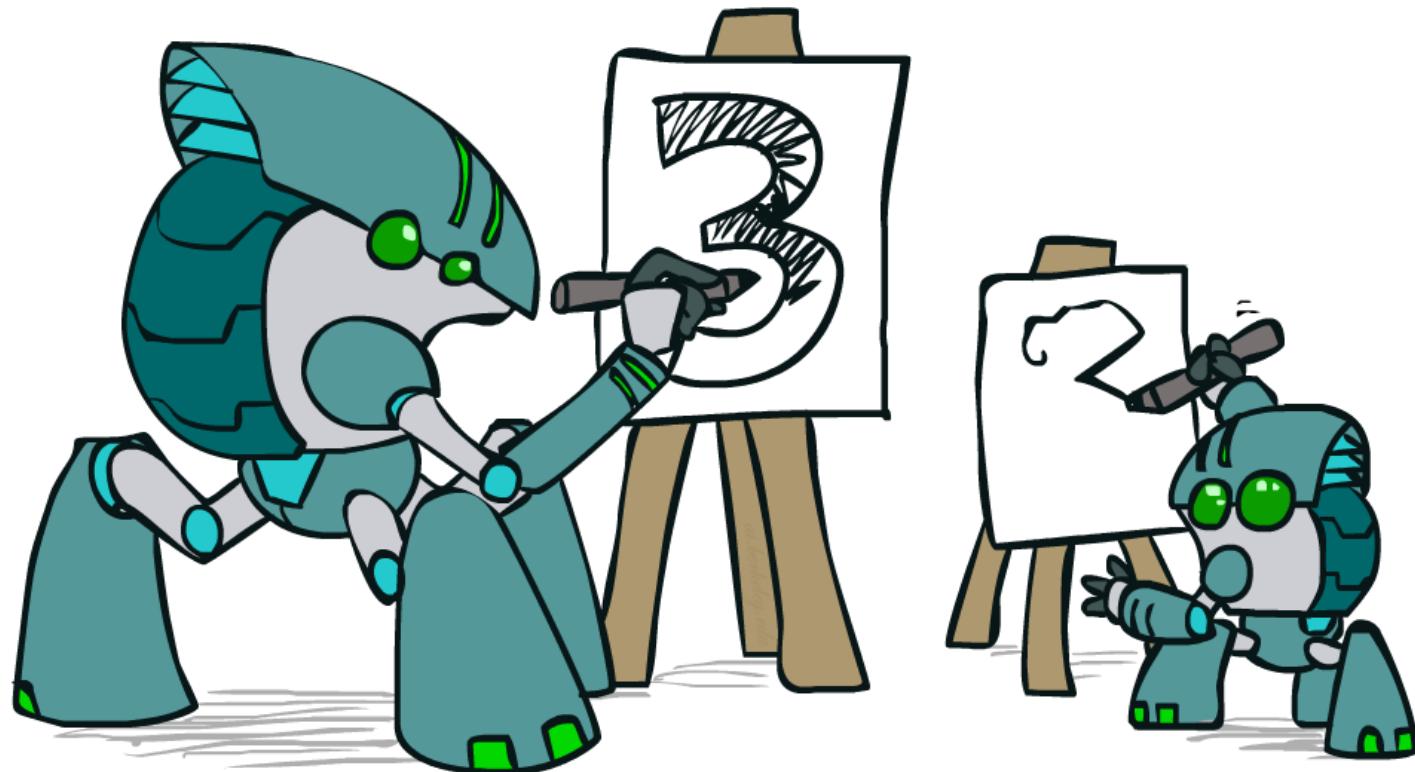
$$y = \arg \max_y w \cdot f(x, y)$$

- Update (if wrong):

$$w = w + f(x, y^*) - f(x, y)$$

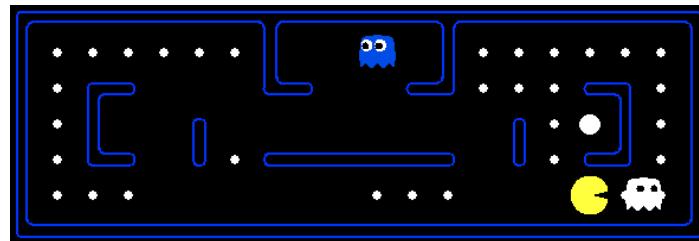


Apprenticeship



Pacman Apprenticeship!

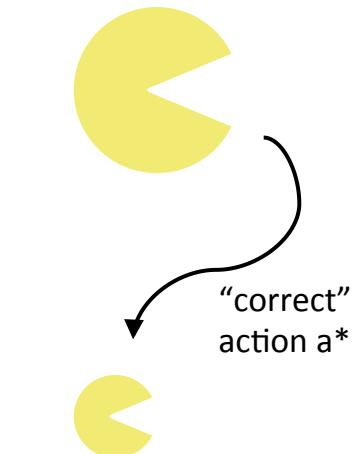
- Examples are states s



- Candidates are pairs (s, a)
- “Correct” actions: those taken by expert
- Features defined over (s, a) pairs: $f(s, a)$
- Score of a q-state (s, a) given by:

$$w \cdot f(s, a)$$

- How is this VERY different from reinforcement learning?



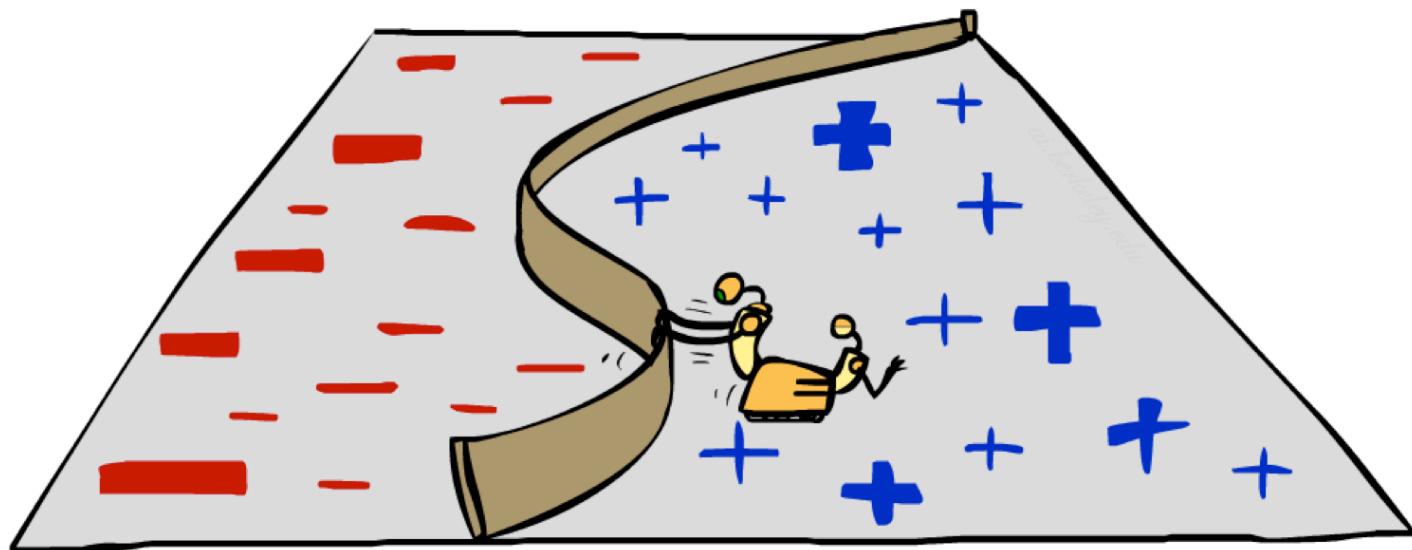
$$\forall a \neq a^*,
w \cdot f(a^*) > w \cdot f(a)$$

[Demo: Pacman Apprentice (L22D1,2,3)]

Video of Demo Pacman Apprentice

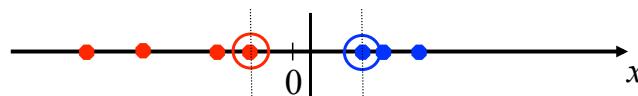


Non-Linearity

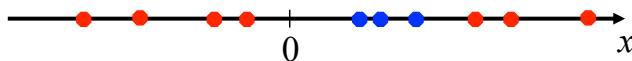


Non-Linear Separators

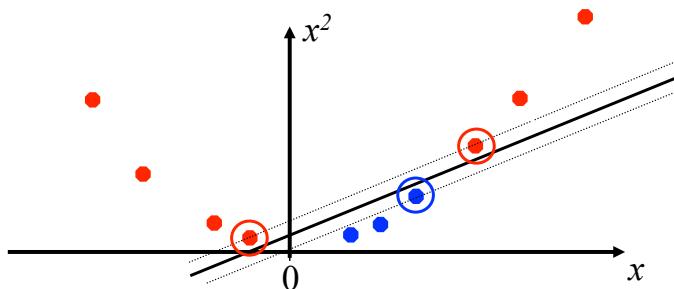
- Data that is linearly separable works out great for linear decision rules:



- But what are we going to do if the dataset is just too hard?



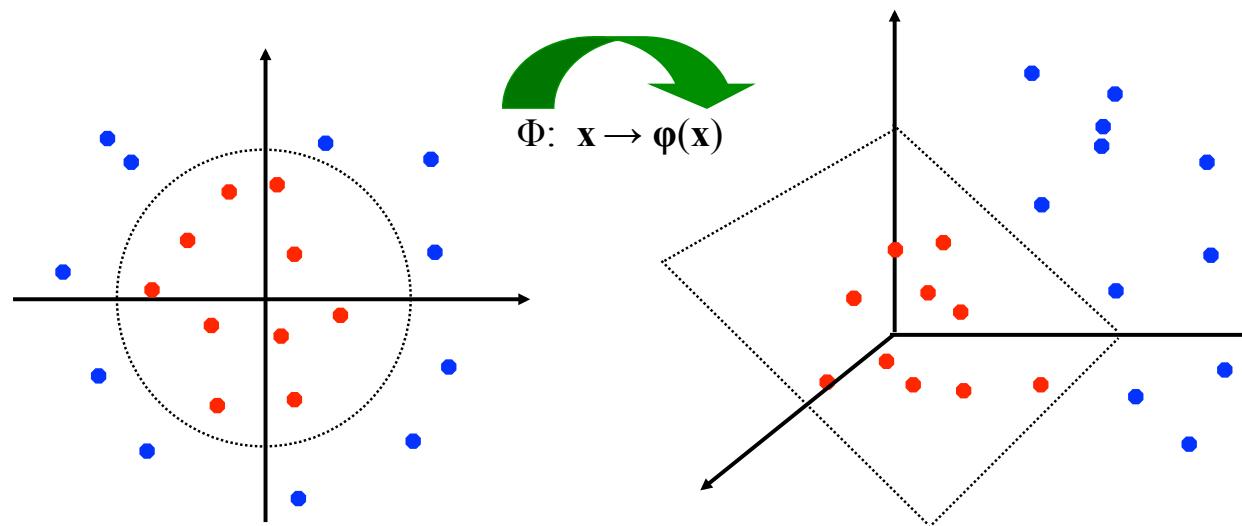
- How about... mapping data to a higher-dimensional space:



This and next slide adapted from Ray Mooney, UT

Non-Linear Separators

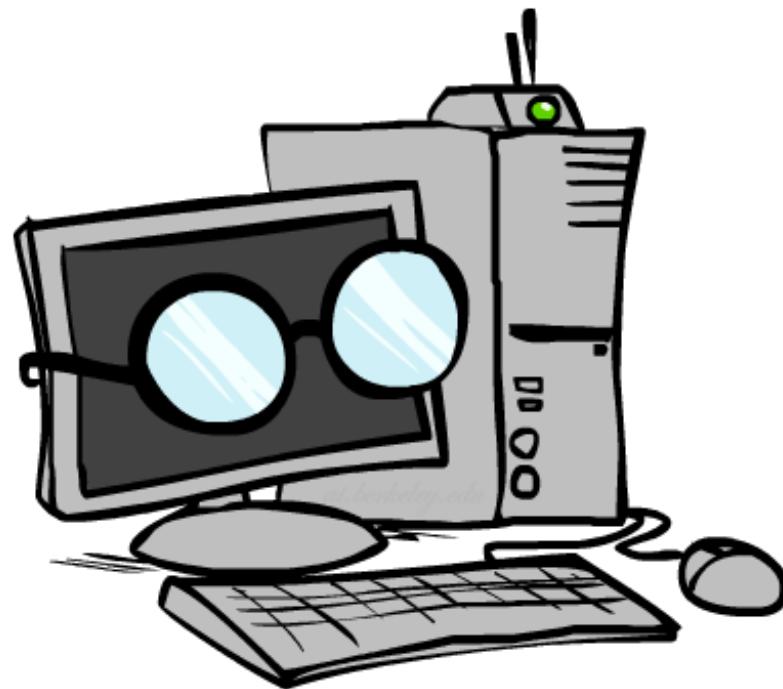
- General idea: the original feature space can always be mapped to some higher-dimensional feature space where the training set is separable:



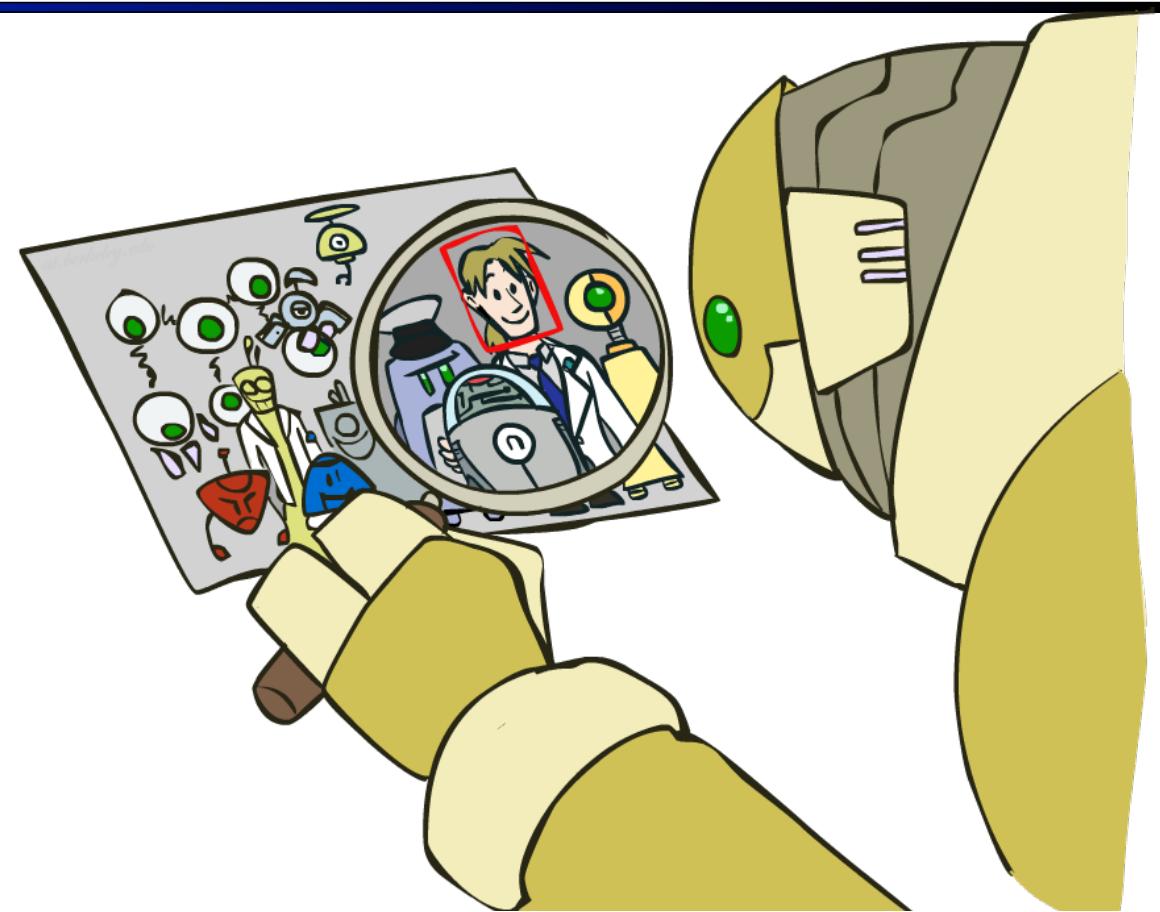
Feature Selection

- To choose between two feature sets:
 - For feature set 1: train perceptron on training data -> Classifier 1
 - For feature set 2: train perceptron on training data -> Classifier 2
- Evaluate performance of Classifier 1 and Classifier 2 on hold-out data
 - Select the one performing best on the hold-out data

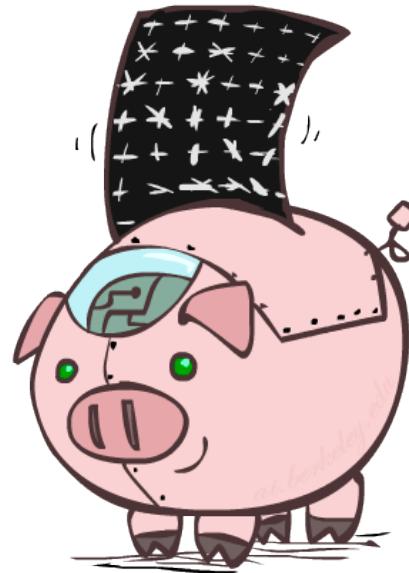
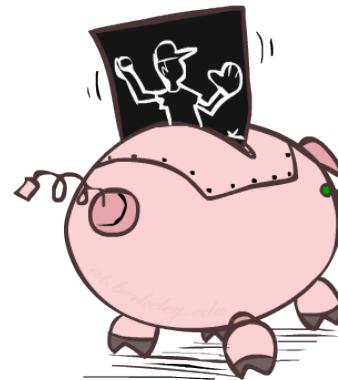
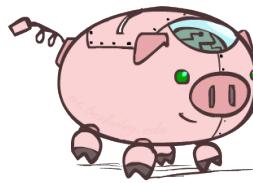
Computer Vision



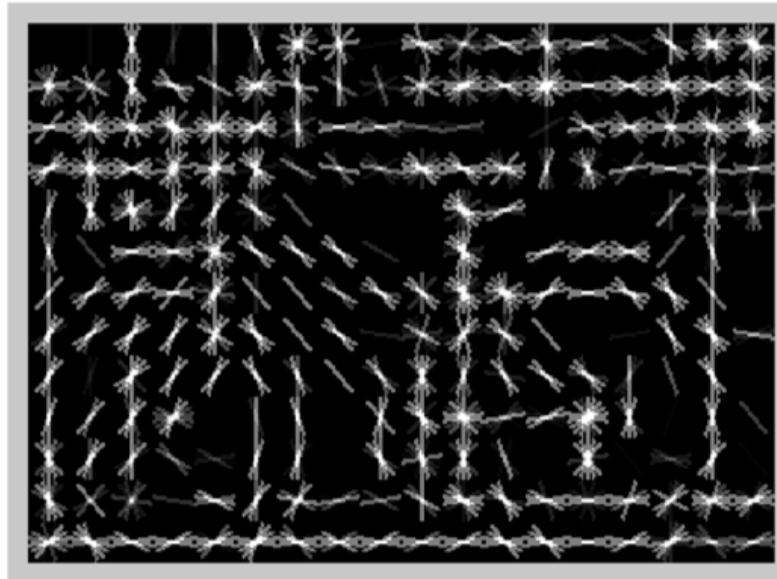
Object Detection



Manual Feature Design



Features and Generalization

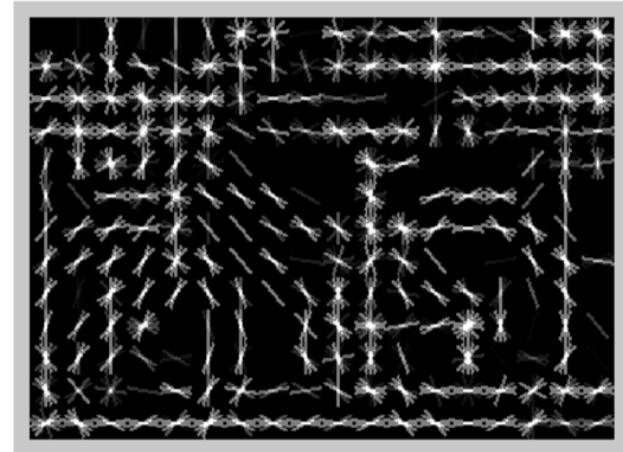


[Dalal and Triggs, 2005]

Features and Generalization

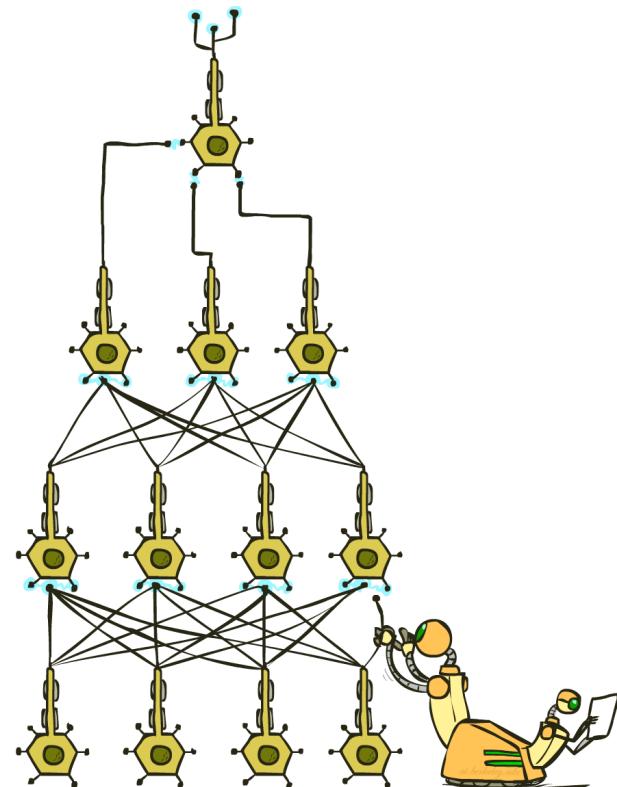


Image



HoG

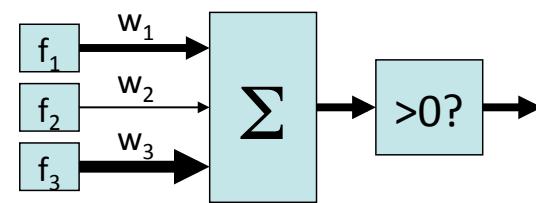
Manual Feature Design → Deep Learning



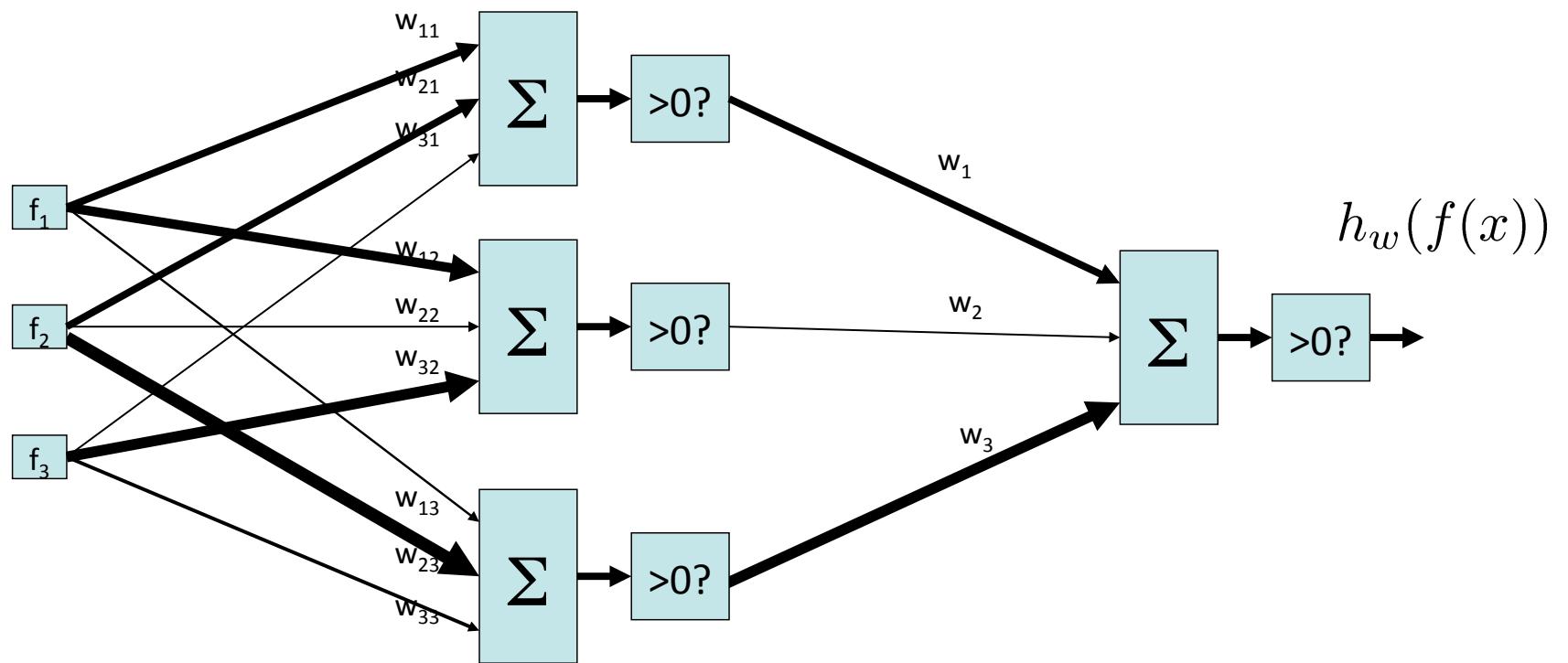
- Manual feature design requires:
 - Domain-specific expertise
 - Domain-specific effort

- What if we could learn the features, too?
 - -> Deep Learning

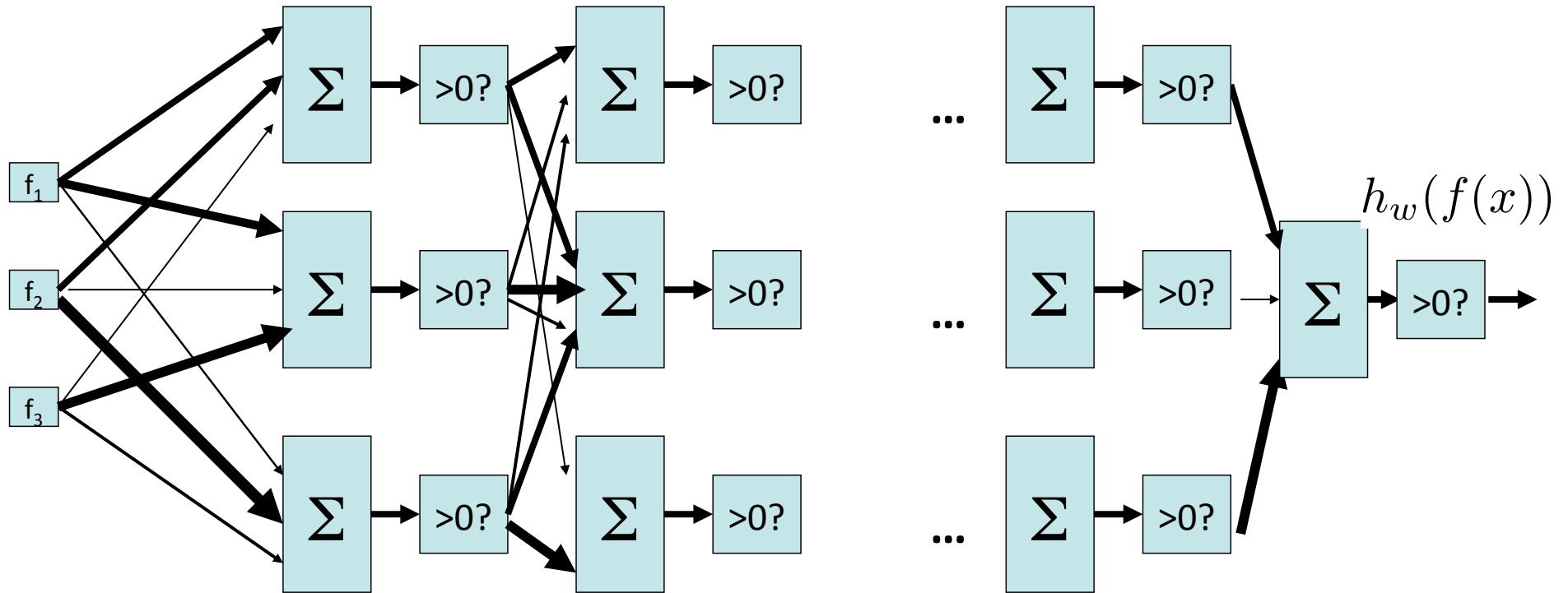
Perceptron



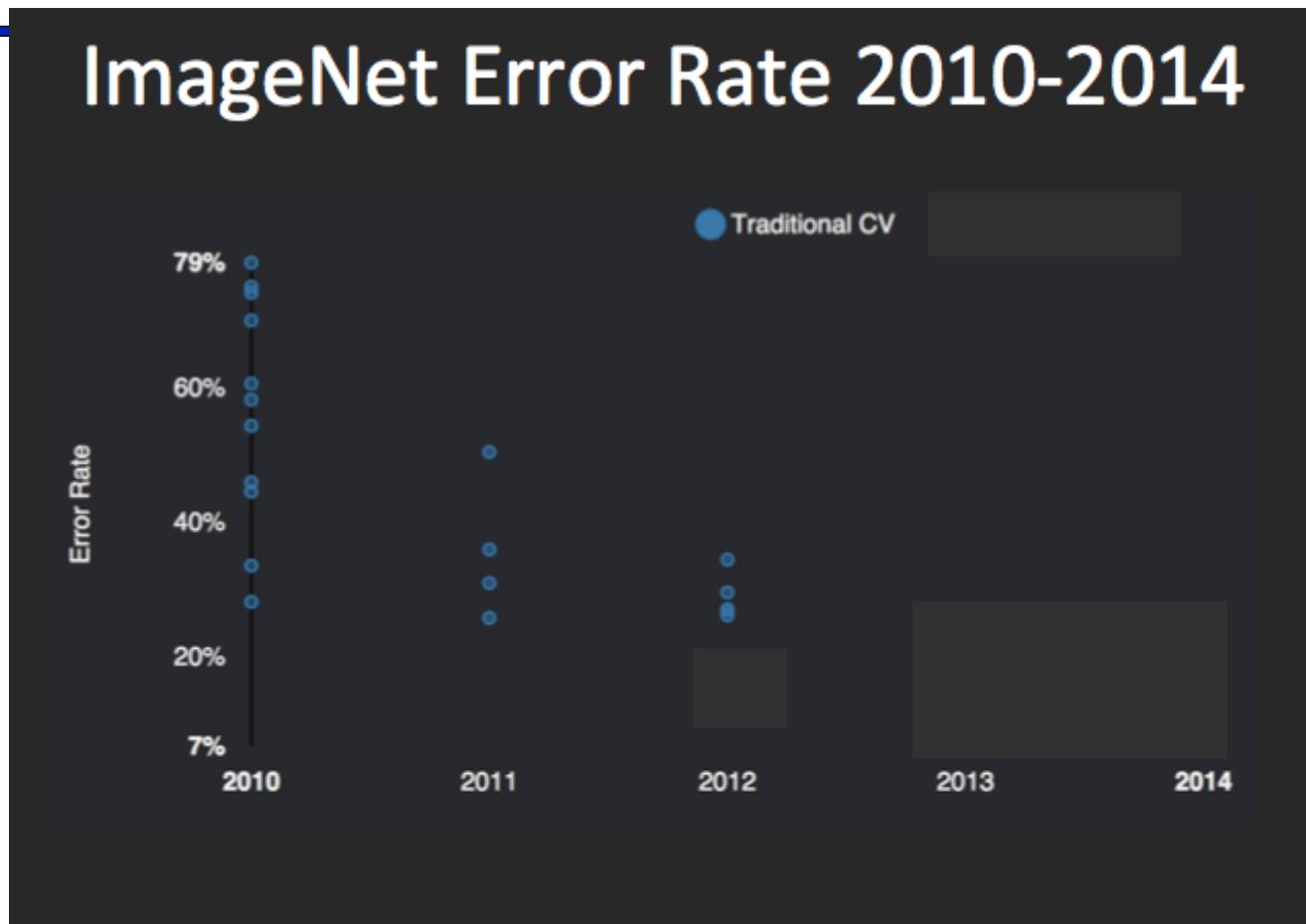
Two-Layer Perceptron Network



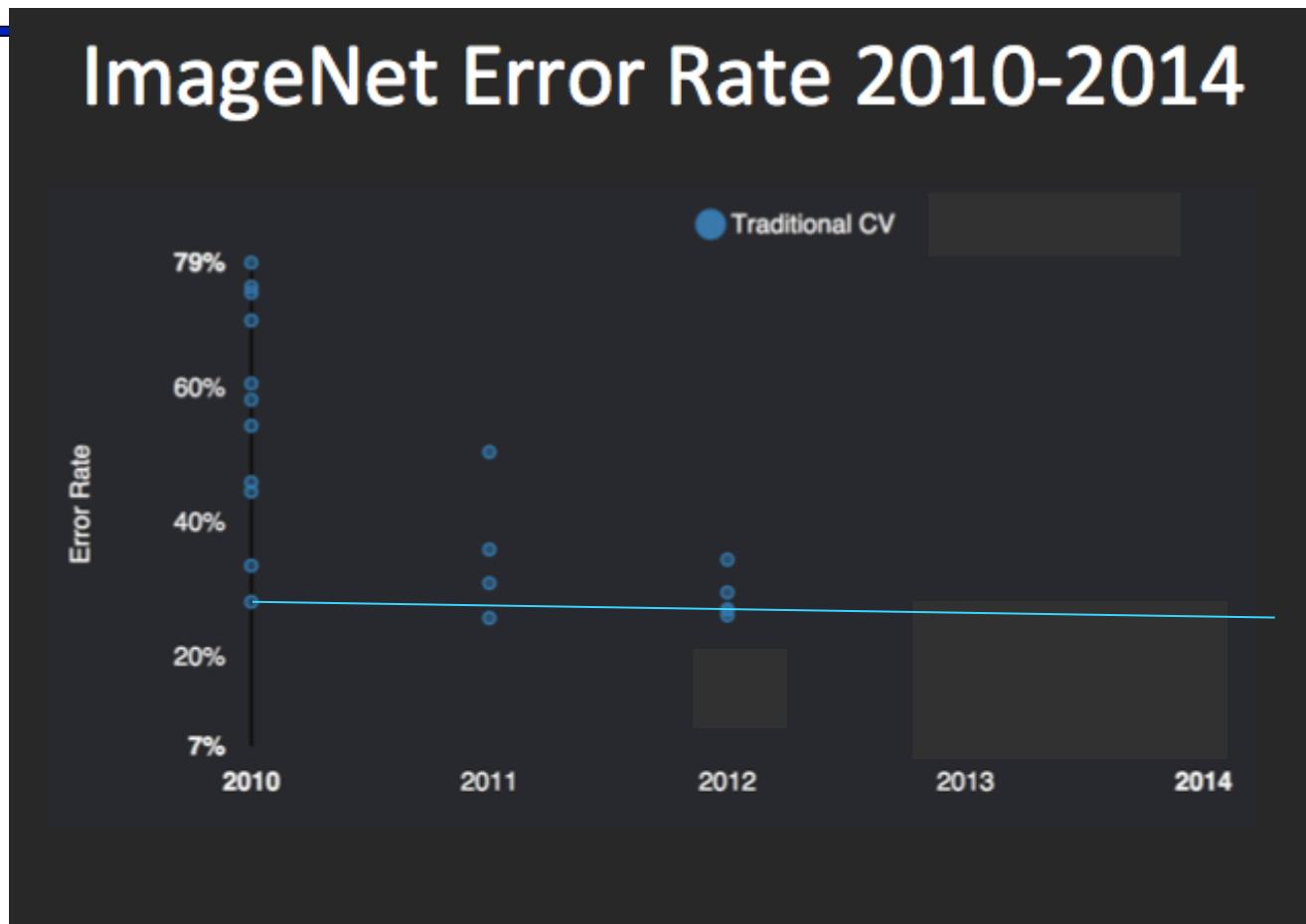
N-Layer Perceptron Network



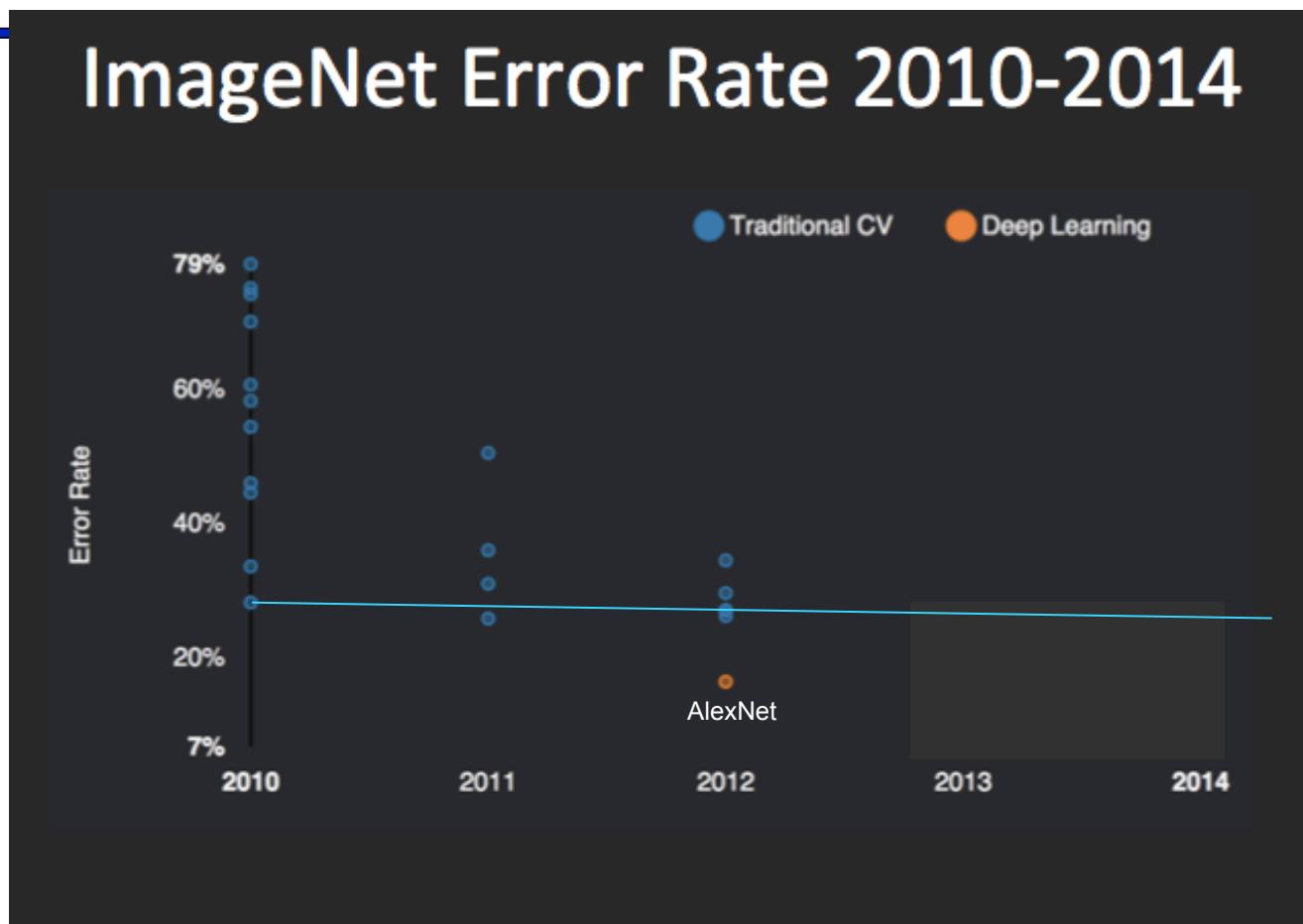
Performance



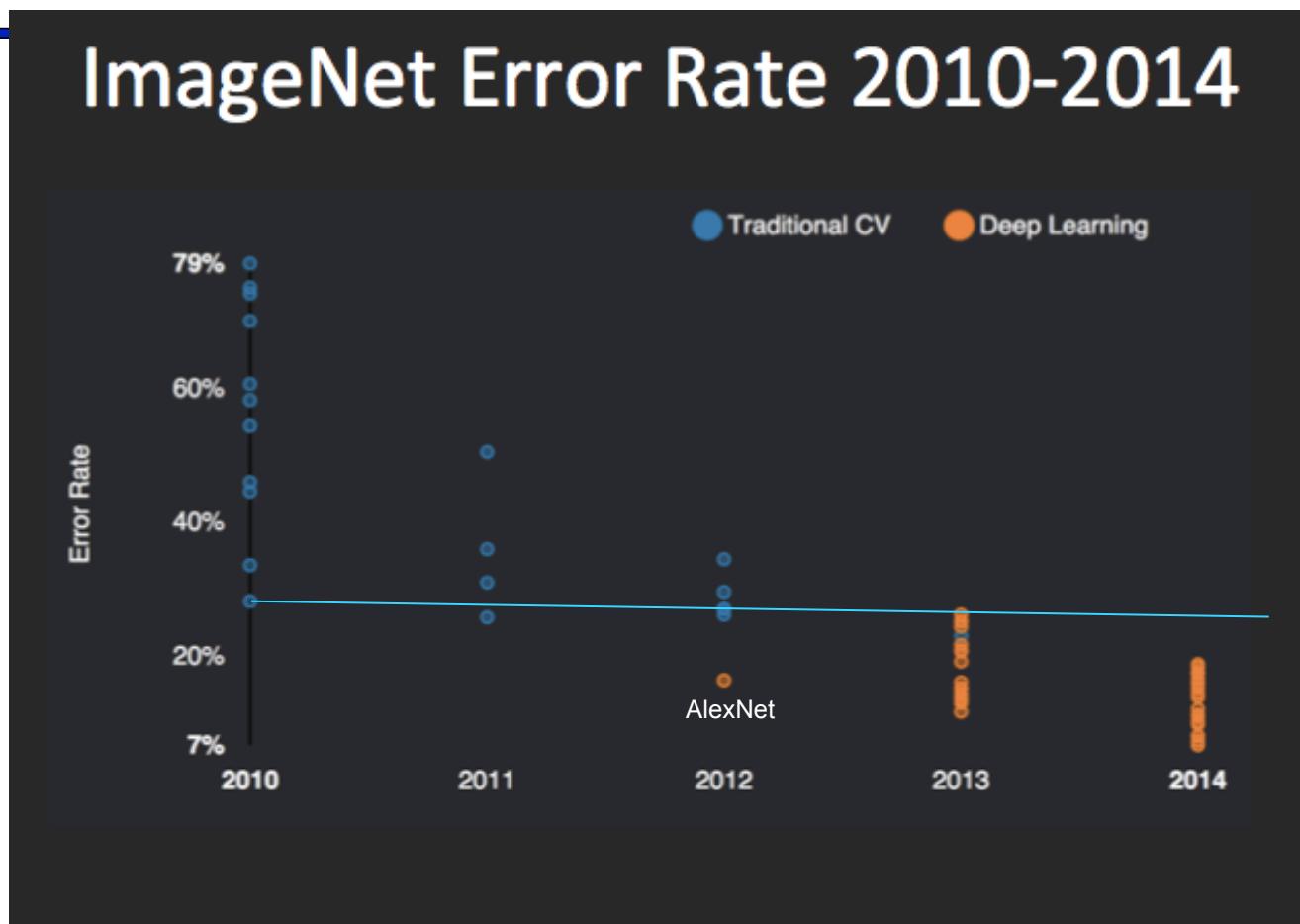
Performance



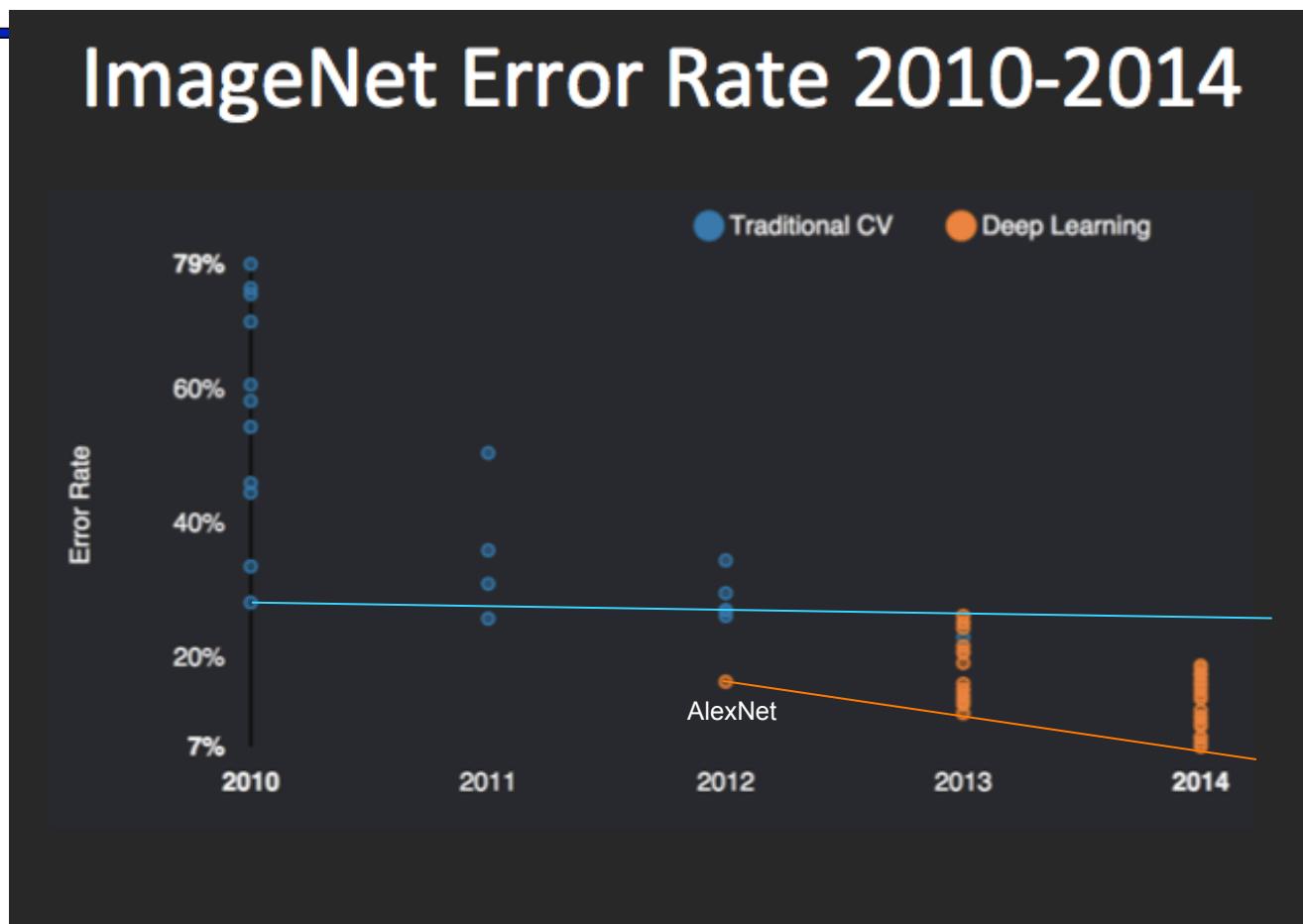
Performance



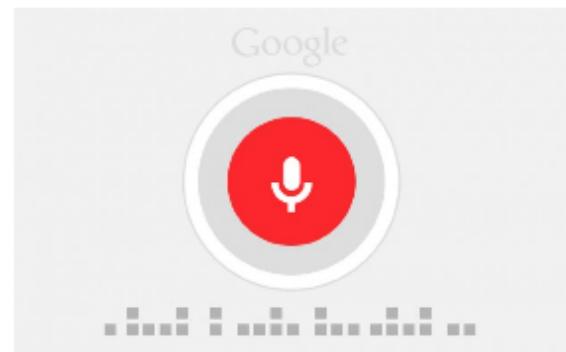
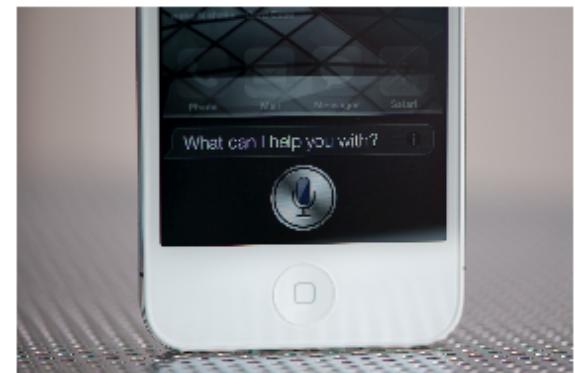
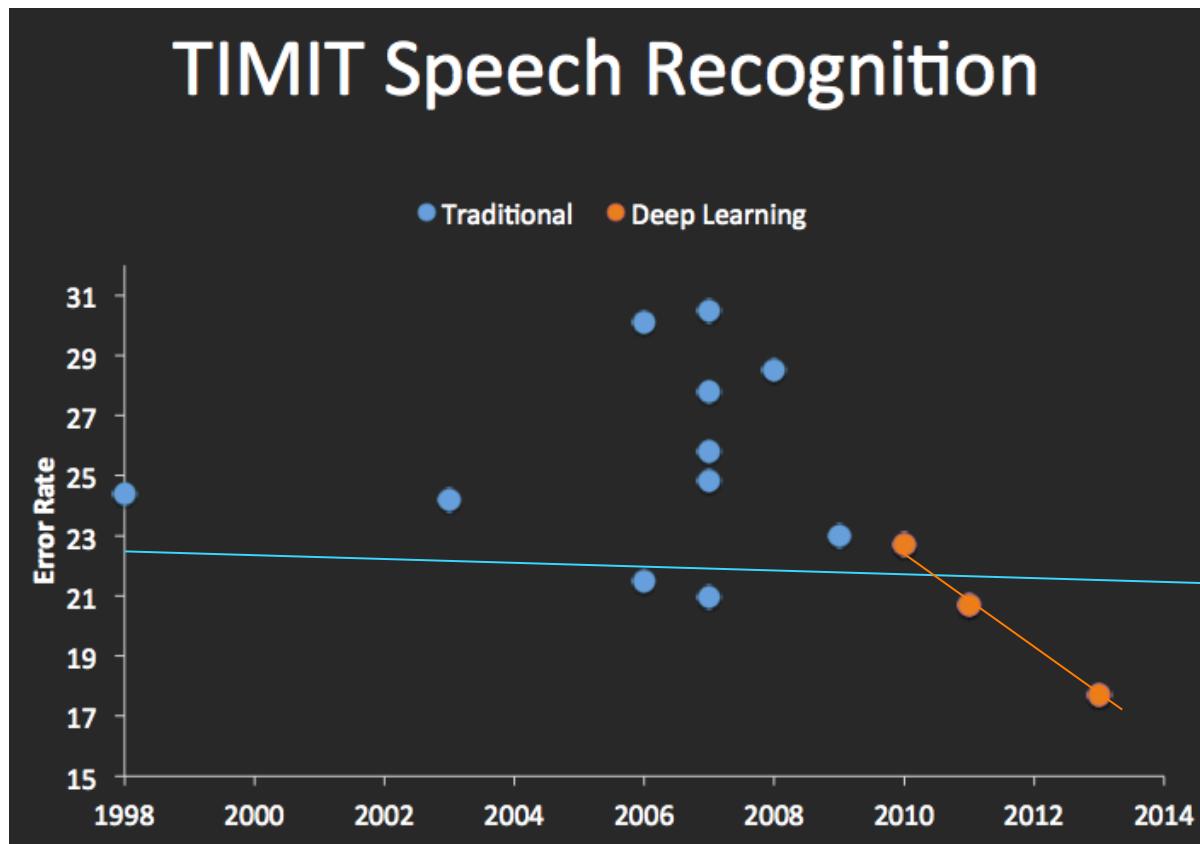
Performance



Performance

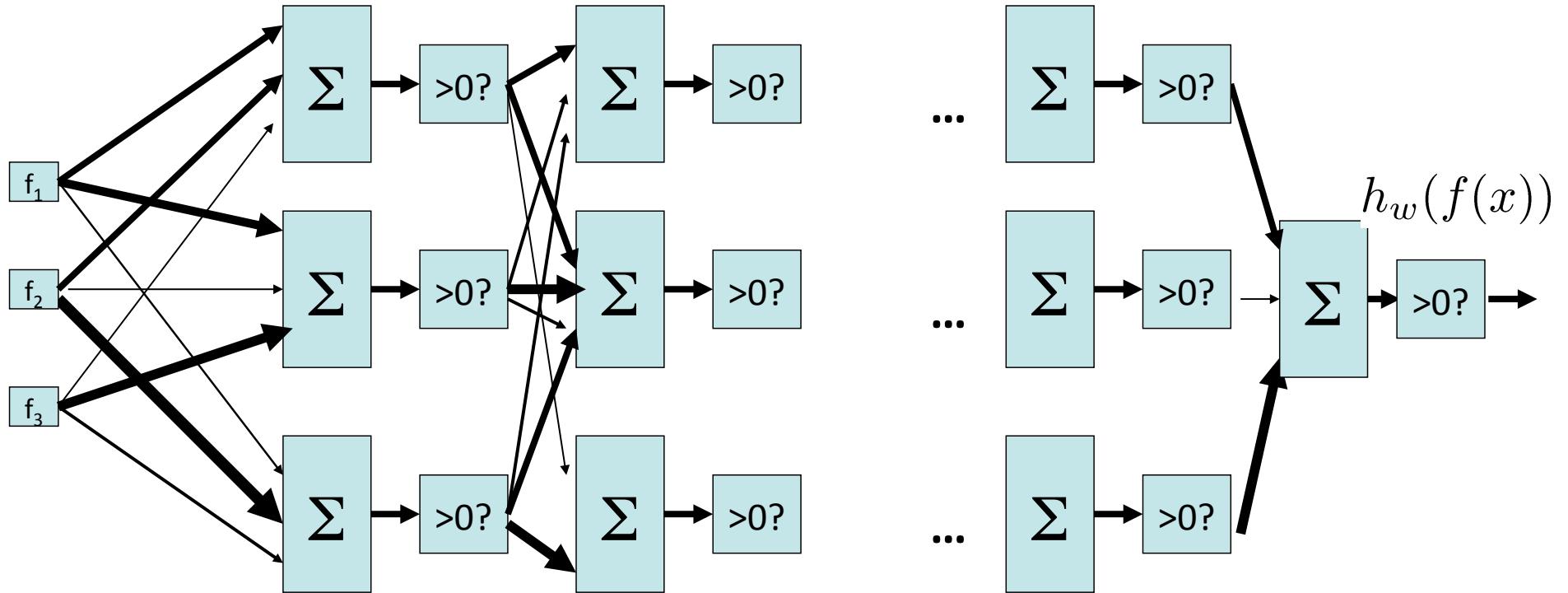


Speech Recognition



graph credit Matt Zeiler, Clarifai

N-Layer Perceptron Network



Local Search

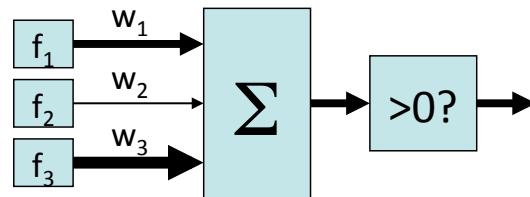
- Simple, general idea:
 - Start wherever
 - Repeat: move to the best neighboring state
 - If no neighbors better than current, quit
 - Neighbors = small perturbations of w
- Properties
 - Plateaus and local optima



→ How to escape plateaus and find a good local optimum?

→ How to deal with very large parameter vectors? E.g., $w \in \mathbb{R}^{1\text{billion}}$

Perceptron



- Objective: Classification Accuracy

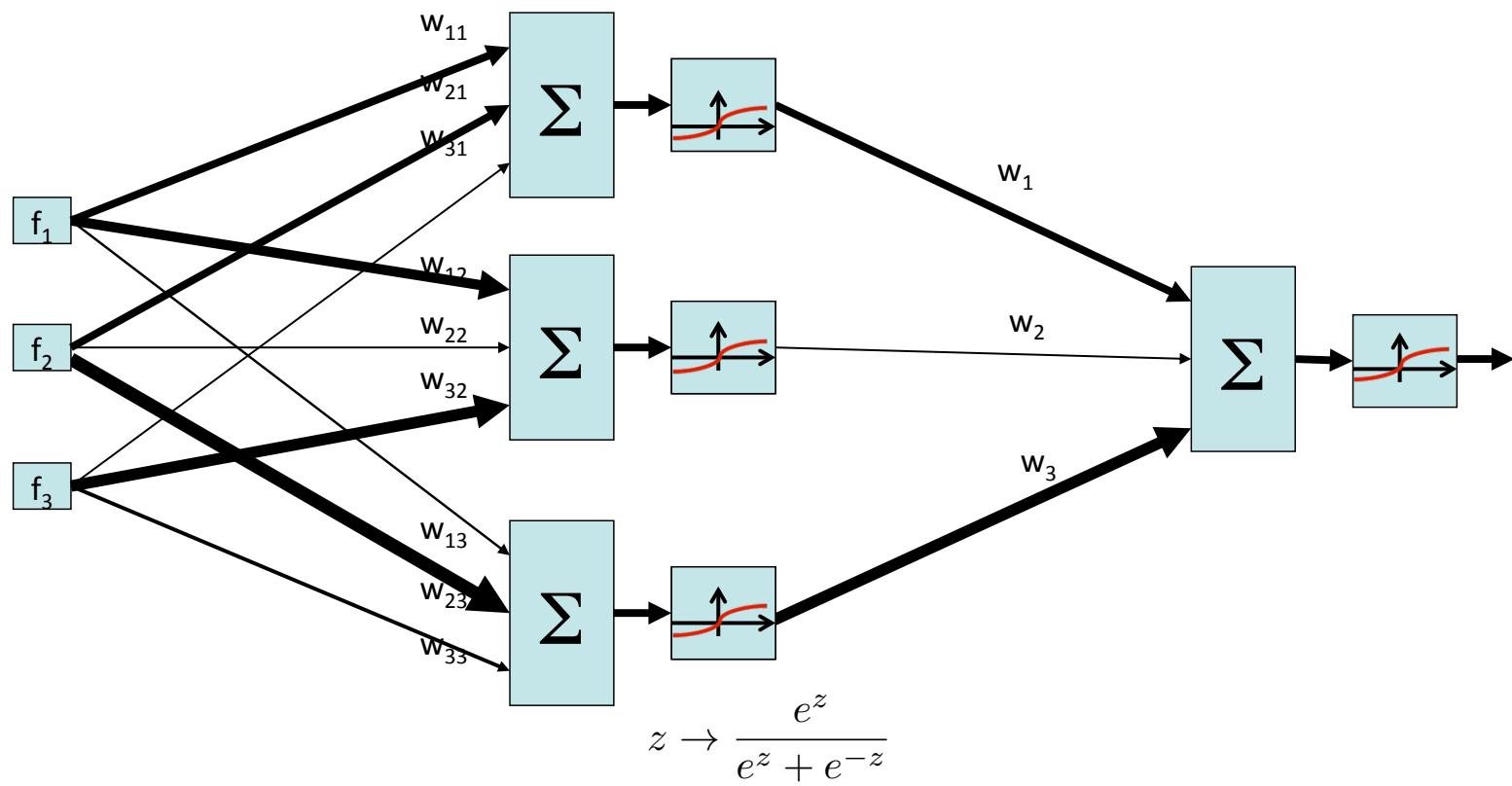
$$l^{\text{acc}}(w) = \frac{1}{m} \sum_{i=1}^m \left(\text{sign}(w^\top f(x^{(i)})) == y^{(i)} \right)$$

- Issue: many plateaus → how to measure incremental progress?

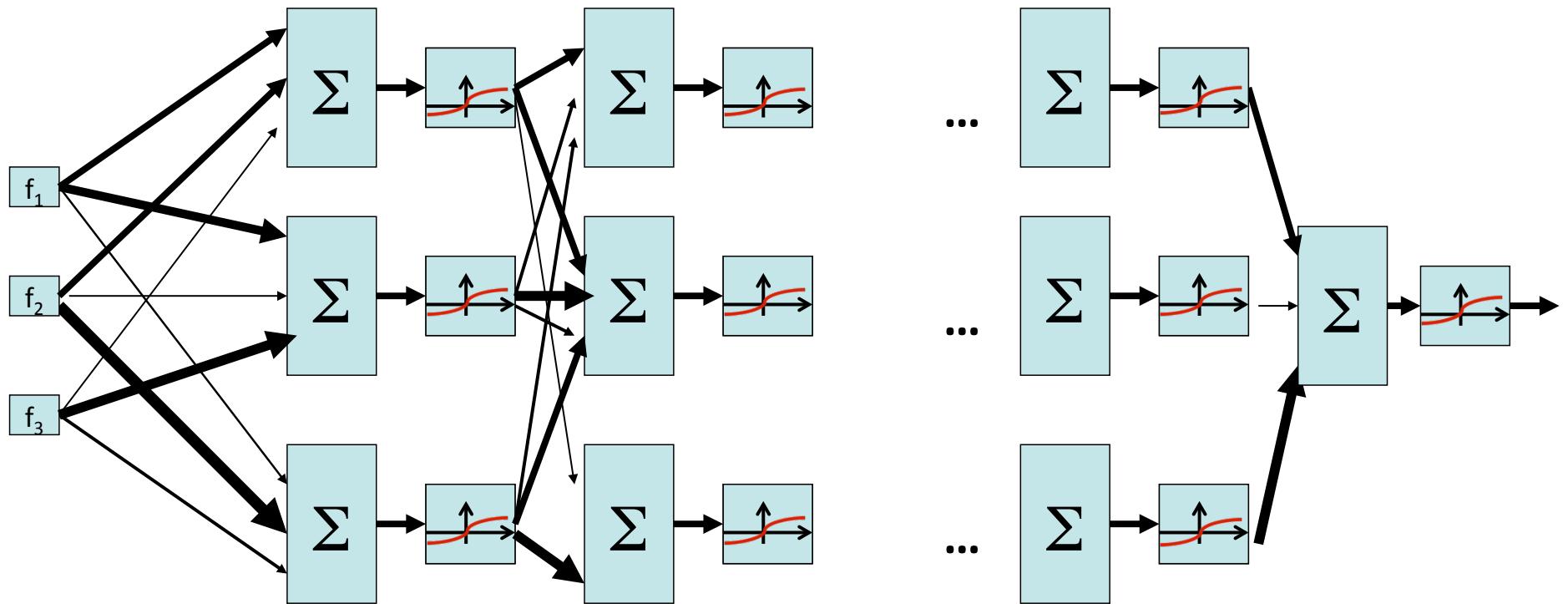
Soft-Max

- Score for $y=1$: $w^\top f(x)$ Score for $y=-1$: $-w^\top f(x)$
- Probability of label:
$$p(y = 1 | f(x); w) = \frac{e^{w^\top f(x^{(i)})}}{e^{w^\top f(x)} + e^{-w^\top f(x)}}$$
$$p(y = -1 | f(x); w) = \frac{e^{-w^\top f(x)}}{e^{w^\top f(x)} + e^{-w^\top f(x)}}$$
- Objective:
$$l(w) = \prod_{i=1}^m p(y = y^{(i)} | f(x^{(i)}); w)$$
- Log:
$$ll(w) = \sum_{i=1}^m \log p(y = y^{(i)} | f(x^{(i)}); w)$$

Two-Layer Neural Network



N-Layer Neural Network



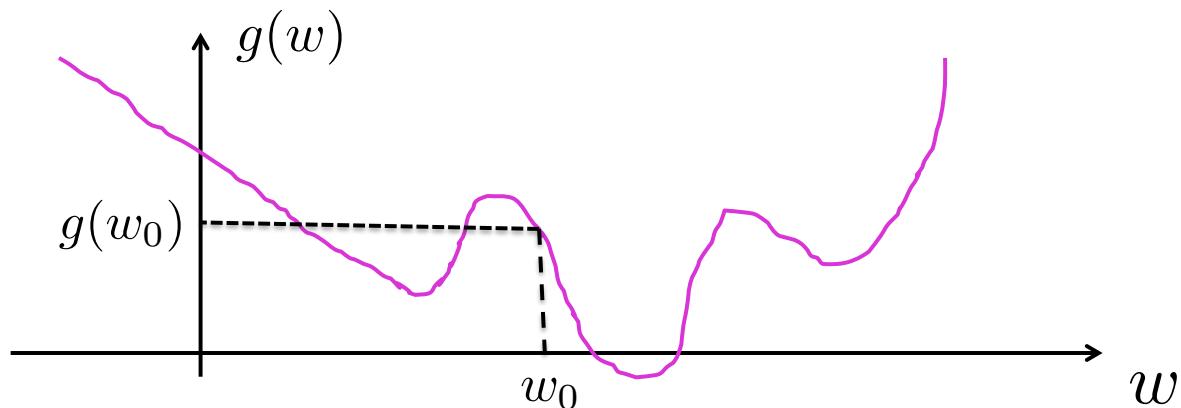
Our Status

- Our objective $ll(w)$
 - Changes smoothly with changes in w
 - Doesn't suffer from the same plateaus as the perceptron network
- Challenge: how to find a good w ?

$$\max_w ll(w)$$

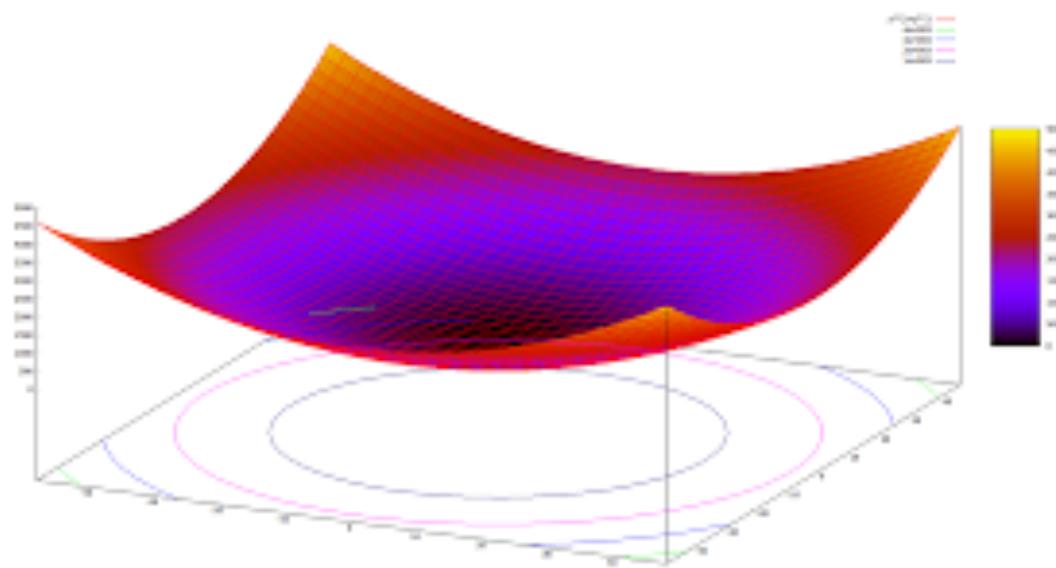
- Equivalently: $\min_w -ll(w)$

1-d optimization



- Could evaluate $g(w_0 + h)$ and $g(w_0 - h)$
 - Then step in best direction
- Or, evaluate derivative: $\frac{\partial g(w_0)}{\partial w} = \lim_{h \rightarrow 0} \frac{g(w_0 + h) - g(w_0 - h)}{2h}$
 - Which tells which direction to step into

2-D Optimization



Source: Thomas Jungblut's Blog

Steepest Descent

- Idea:
 - Start somewhere
 - Repeat: Take a step in the steepest descent direction

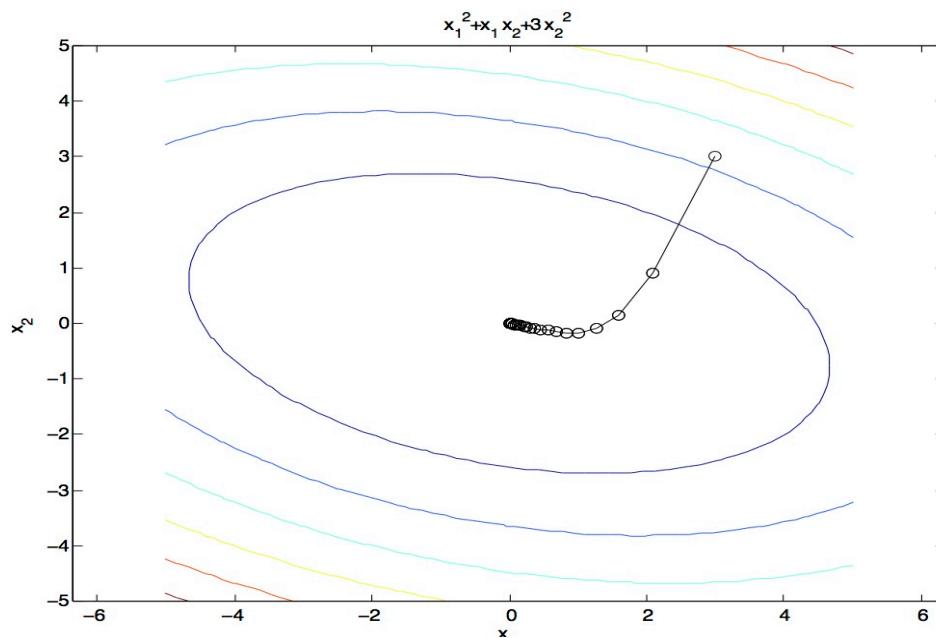


Figure source: Mathworks

What is the Steepest Descent Direction?

$$\min_{\Delta: \Delta_1^2 + \Delta_2^2 \leq \epsilon} g(w + \Delta)$$

- First-Order Taylor Expansion: $g(w + \Delta) \approx g(w) + \frac{\partial g}{\partial w_1} \Delta_1 + \frac{\partial g}{\partial w_2} \Delta_2$
- Steepest Descent Direction: $\min_{\Delta: \Delta_1^2 + \Delta_2^2 \leq \epsilon} \frac{\partial g}{\partial w_1} \Delta_1 + \frac{\partial g}{\partial w_2} \Delta_2$
- Recall: $\min_{a: \|a\| \leq \epsilon} a^\top b \quad \rightarrow \quad a = -b \frac{\epsilon}{\|b\|}$
- Hence, solution: $-\nabla g \frac{\epsilon}{\|\nabla g\|}$ $\nabla g = \begin{bmatrix} \frac{\partial g}{\partial w_1} \\ \frac{\partial g}{\partial w_2} \end{bmatrix}$

Generally, Steepest Direction

- Steepest Direction = direction of the gradient

$$\nabla g = \begin{bmatrix} \frac{\partial g}{\partial w_1} \\ \frac{\partial g}{\partial w_2} \\ \vdots \\ \frac{\partial g}{\partial w_n} \end{bmatrix}$$

Optimization Procedure 1: Gradient Descent

- Init: w
- For $i = 1, 2, \dots$

$$w \leftarrow w - \alpha * \nabla g(w)$$

- α : learning rate --- tweaking parameter that needs to be chosen carefully
- How? Try multiple choices
 - Crude rule of thumb: update changes w about 0.1 – 1 %

Try Many Learning Rates

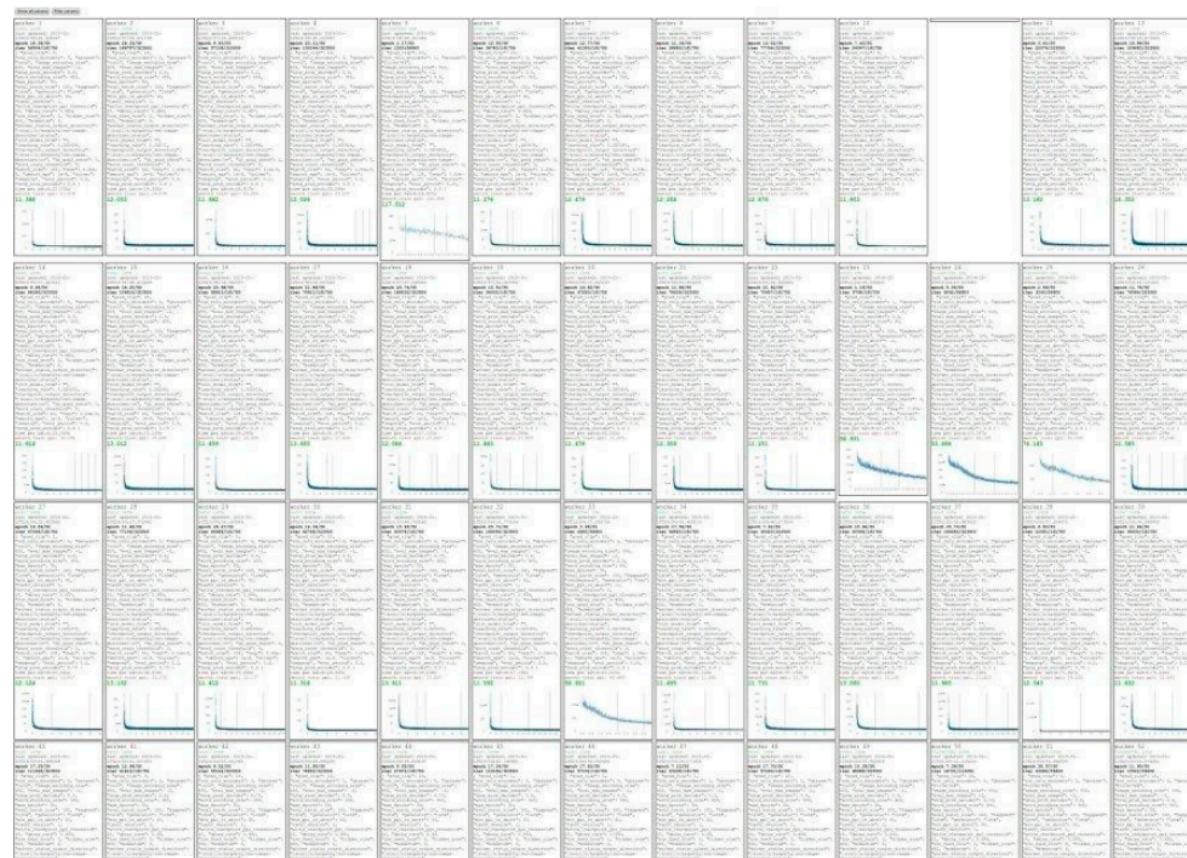
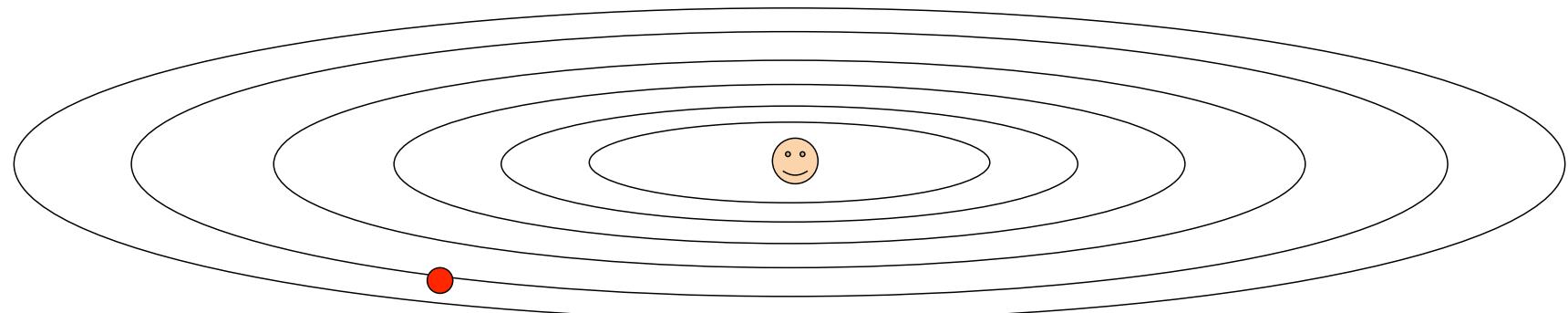


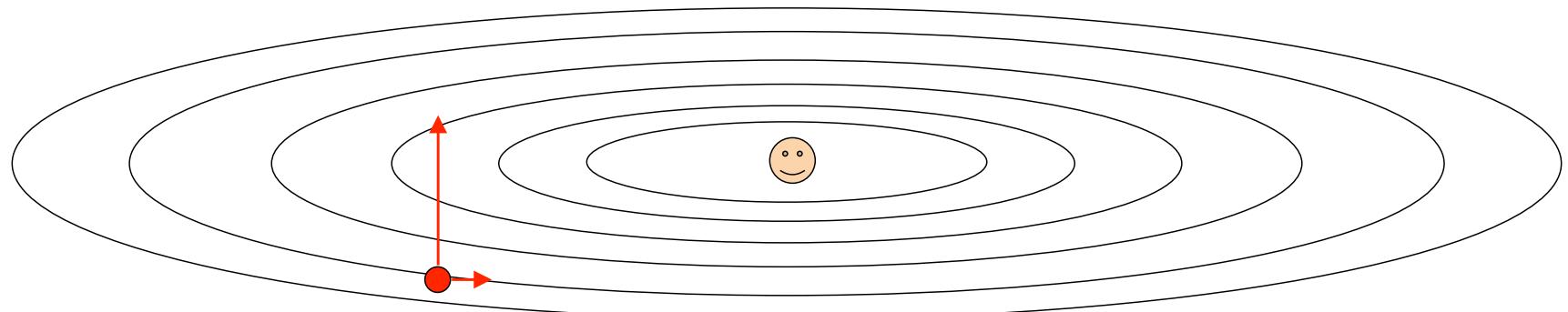
Figure source:
Andrej Karpathy

Suppose loss function is steep vertically but shallow horizontally:



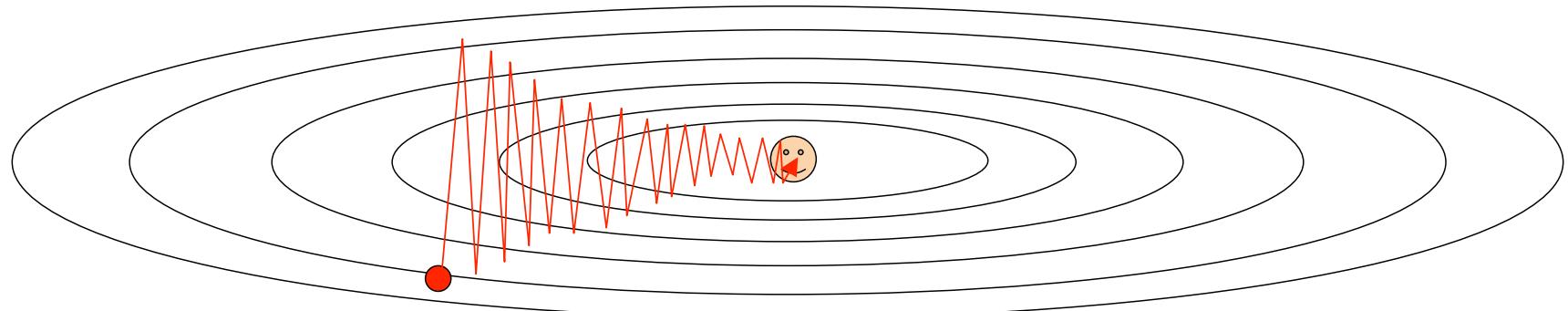
Q: What is the trajectory along which we converge towards the minimum with SGD?

Suppose loss function is steep vertically but shallow horizontally:



Q: What is the trajectory along which we converge towards the minimum with SGD?

Suppose loss function is steep vertically but shallow horizontally:



Q: What is the trajectory along which we converge towards the minimum with Gradient Descent? very slow progress along flat direction, jitter along steep

one
Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 6 - 54

25 Jan 2016

Optimization Procedure 2: Momentum

- Gradient Descent

- Init: w
- For $i = 1, 2, \dots$

$$w \leftarrow w - \alpha * \nabla g(w)$$

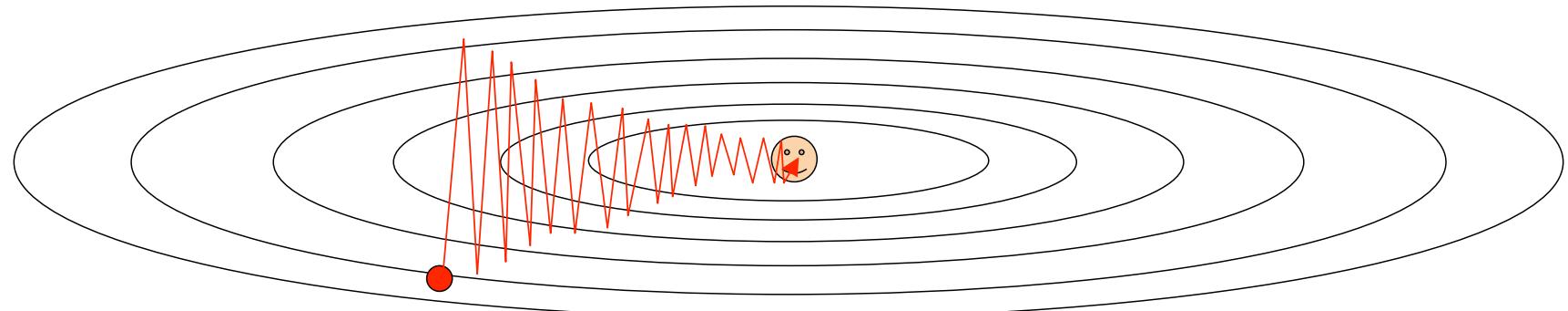
- Momentum

- Init: w
- For $i = 1, 2, \dots$

$$\begin{aligned} v &\leftarrow \mu * v - \alpha * \nabla g(w) \\ w &\leftarrow w + v \end{aligned}$$

- Physical interpretation as ball rolling down the loss function + friction (mu coefficient).
- mu = usually $\sim 0.5, 0.9$, or 0.99 (Sometimes annealed over time, e.g. from $0.5 \rightarrow 0.99$)

Suppose loss function is steep vertically but shallow horizontally:



Q: What is the trajectory along which we converge towards the minimum with Momentum?

Some More Advanced Optimization Methods*

- Nesterov Momentum Update
- AdaGrad
- RMSProp
- ADAM
- L-BFGS

Remaining Pieces

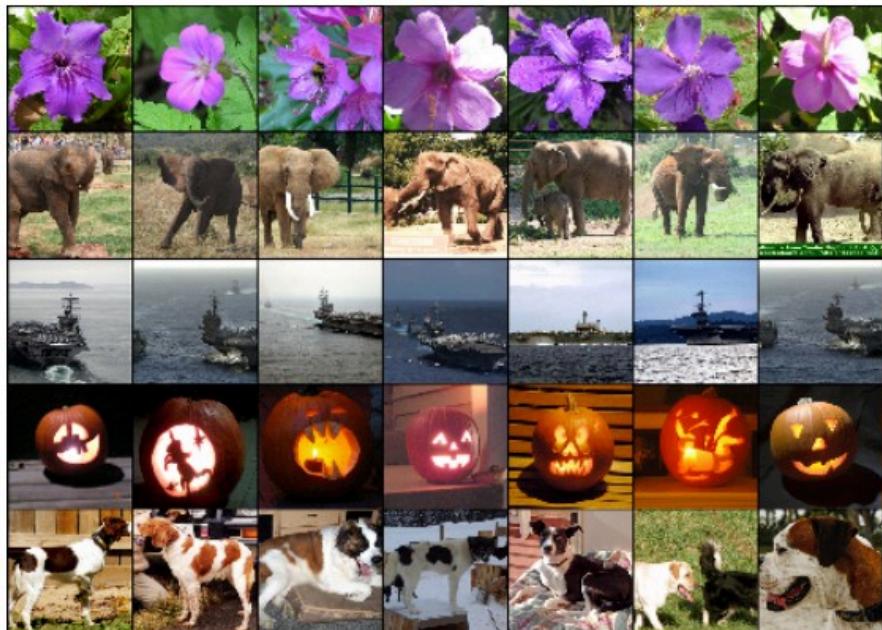
- Optimizing machine learning objectives:
 - Stochastic Descent
 - Mini-batches
- Improving generalization
 - Drop-out
- Activation functions
- Initialization
- Renormalization
- Computing the gradient $\nabla g(w)$
 - Backprop
 - Gradient checking

ConvNets are everywhere

Classification



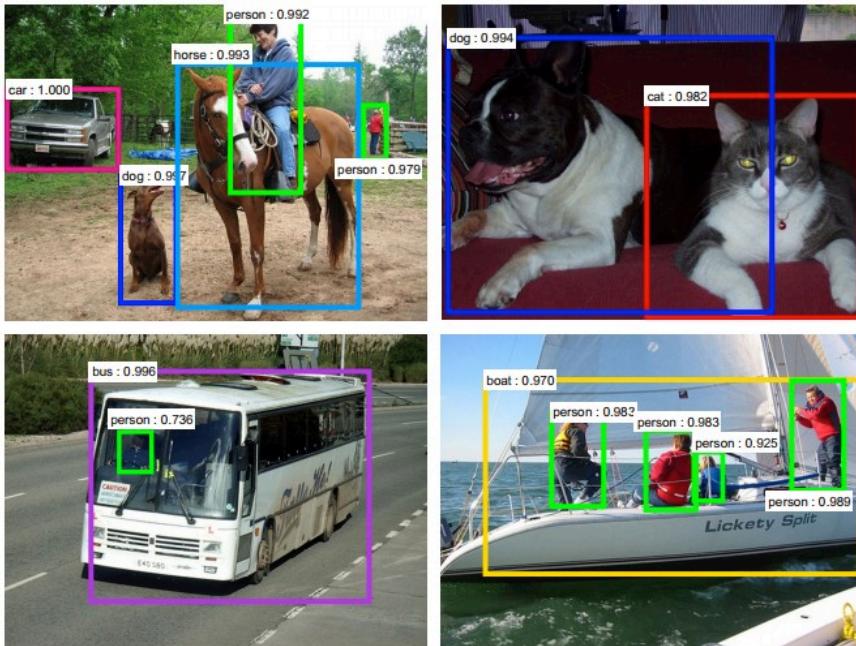
Retrieval



[Krizhevsky 2012]

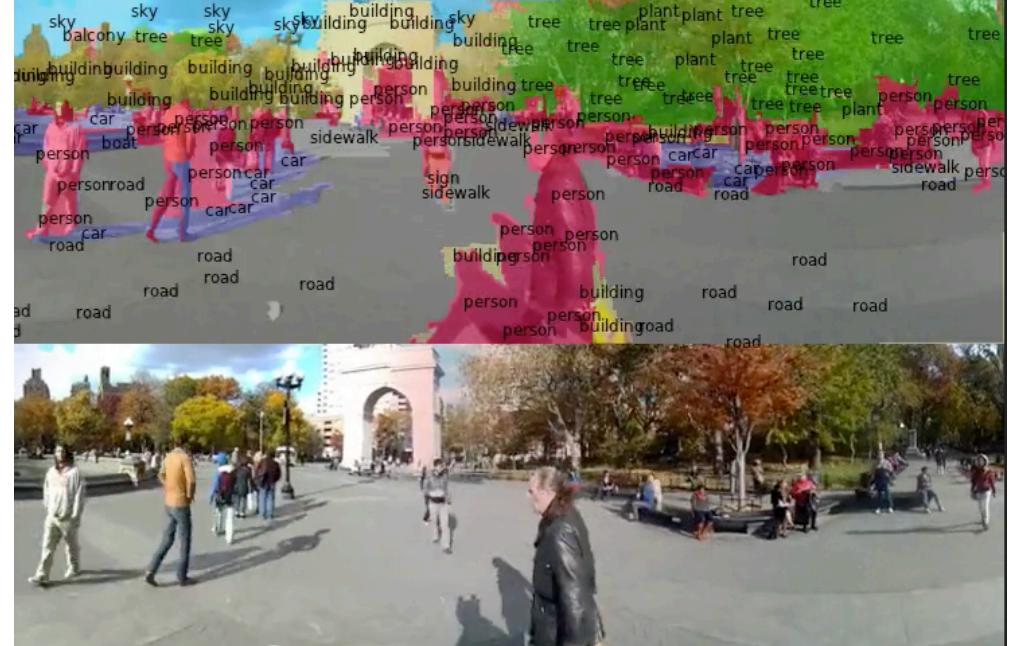
ConvNets are everywhere

Detection



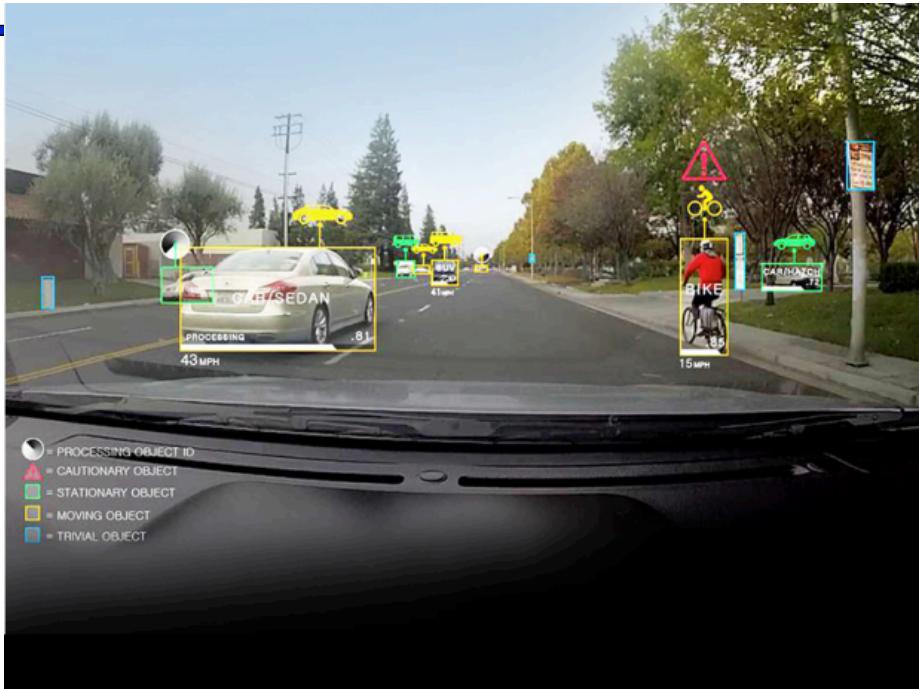
[Faster R-CNN: Ren, He, Girshick, Sun 2015]

Segmentation

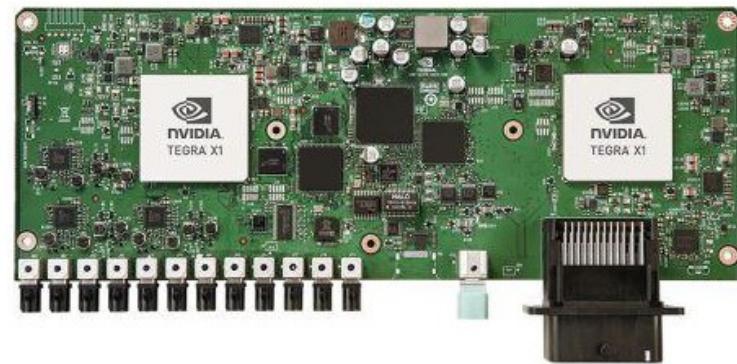


[Farabet et al., 2012]

ConvNets are everywhere

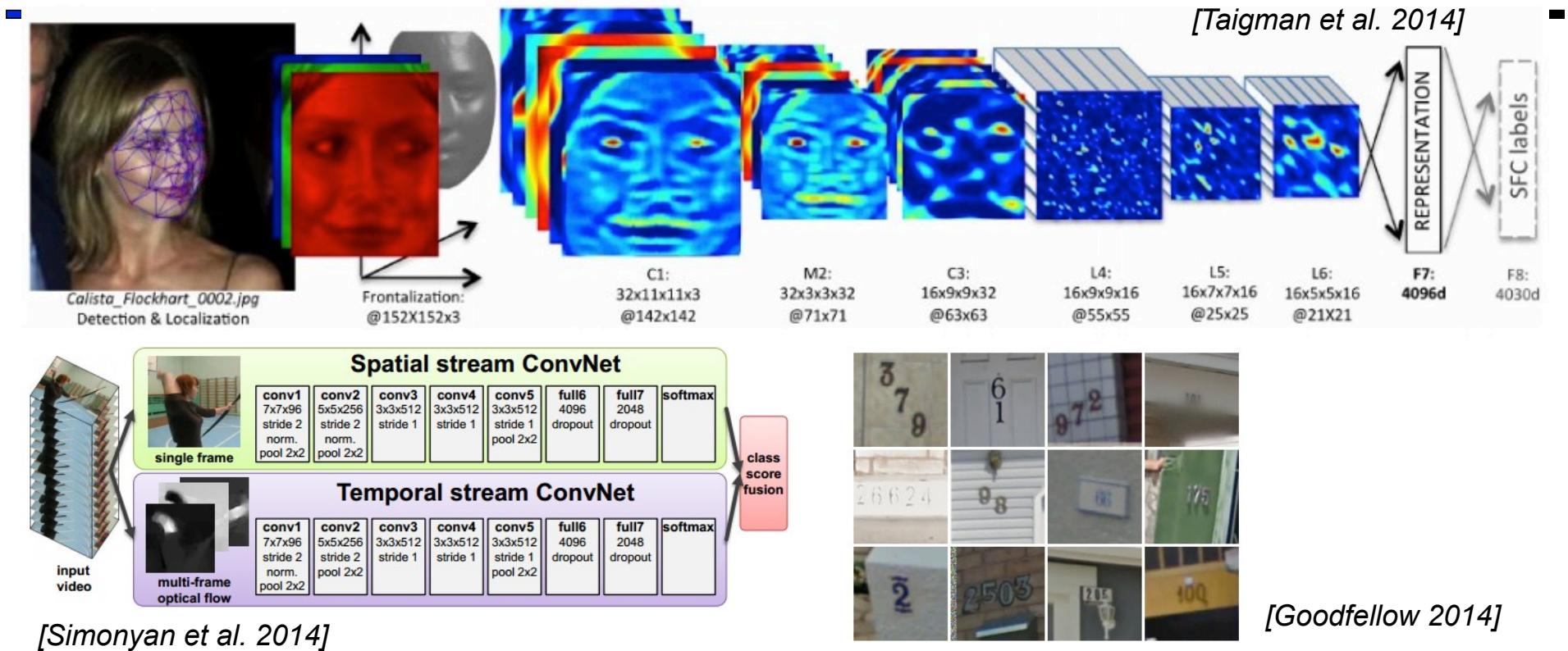


self-driving cars



NVIDIA Tegra X1

ConvNets are everywhere



[Simonyan et al. 2014]

[Goodfellow 2014]

ConvNets are everywhere

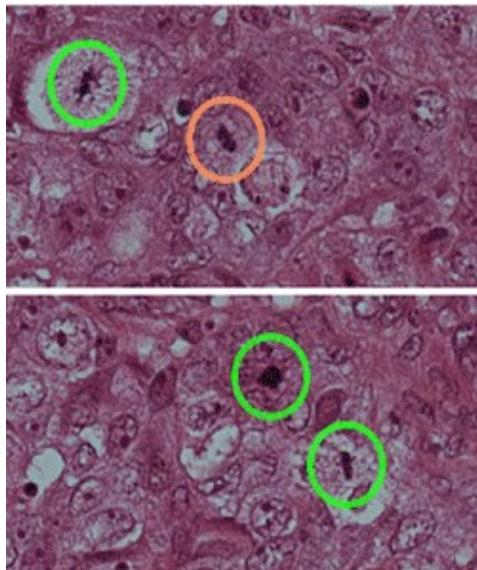


[Toshev, Szegedy 2014]



[Mnih 2013]

ConvNets are everywhere

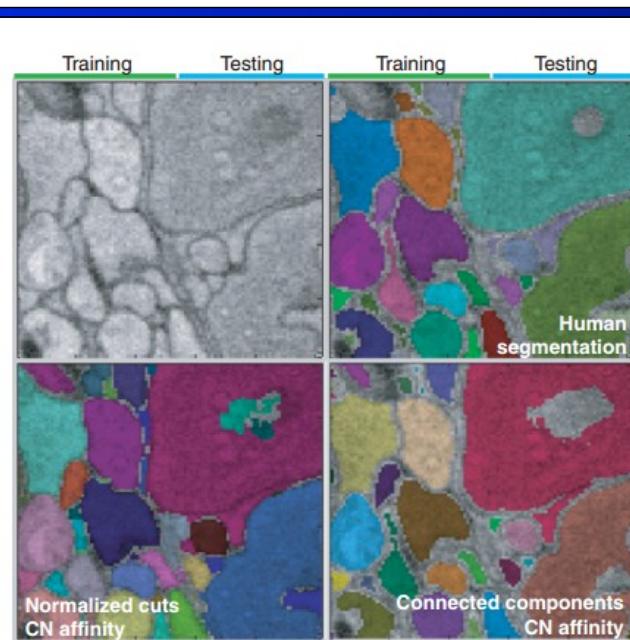


[Ciresan et al. 2013]

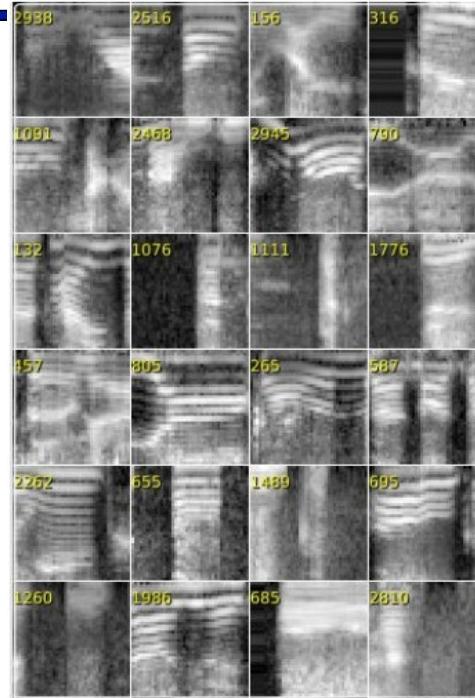


[Sermanet et al. 2011]
[Ciresan et al.]

ConvNets are everywhere



[Turaga et al., 2010]



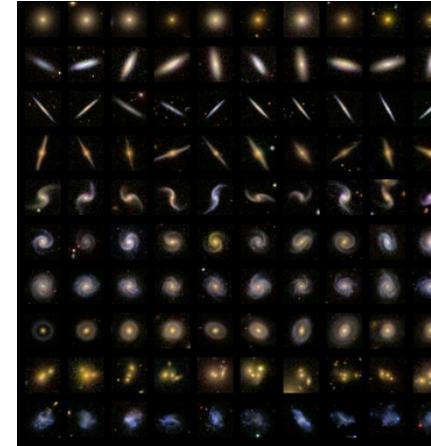
I caught this movie on the Sci-Fi channel recently. It actually turned out to be pretty decent as far as B-list horror/suspense films go. **Two guys (one naive and one loud mouthed a**) take a road trip to stop a wedding but have the worst possible luck when a maniac in a freaky, make-shift tank/truck hybrid decides to play cat-and-mouse with them.** Things are further complicated when they pick up a ridiculously whorish hitchhiker. What makes this film unique is that the combination of comedy and terror actually work in this movie, unlike so many others. The two guys are likable enough and there are some good chase/suspense scenes. Nice pacing and comic timing make this movie more than passable for the horrorfucker buff. **Definitely worth checking out.**

I just saw this on a local independent station in the New York City area. **The cast showed promise but when I saw the director, George Cosmatos, I became suspicious. And sure enough, it was every bit as bad, every bit as pointless and stupid as every George Cosmatos movie I ever saw.** He's like a stupid man's Michael Bay – with all the awfulness that accolade promises. There's no point to the conspiracy, no burning issues that urge the conspirators on. We are left to ourselves to connect the dots from one bit of graffiti on various walls in the film to the next. Thus, the current budget crisis, the war in Iraq, Islamic extremism, the fate of social security, 47 million Americans without health care, stagnating wages, and the death of the middle class are all subsumed by the sheer terror of graffiti. A truly, stunningly idiotic film.

Graphics is far from the best part of the game. **This is the number one best TH game in the series!** Next to Underground. It deserves strong love. It is an insane game. These are massive levels, massive unlockable characters... it's just a massive game. **Waste your money on this game. This is the kind of money that is wasted properly.** And even though graphics suck, that's doesn't make a game good. Actually, the graphics were good at the time. Today the graphics are crap. WHO CARES? As they say in Canada, this is the fun game, aye. (You get to go to Canada in THPS3) Well, I don't know if they say that, but they might, who knows. Well, Canadian people do. Wait a minute, I'm getting off topic. This game rocks. Buy it, play it, enjoy it, love it. It's PURE BRILLIANCE.

The first was good and original. I was a not bad horror/comedy movie. So I heard a second one was made and I had to watch it. What really makes this movie work is Judd Nelson's character and the sometimes clever script. **A pretty good script for a person who wrote the Final Destination films and the direction was okay.** Sometimes there's scenes where it looks like it was filmed using a home video camera with a grainy - look. Great made - for - TV movie. **It was worth the rental and probably worth buying just to get that nice eerie feeling and watch Judd Nelson's Stanley doing what he does best.** I suggest newcomers to watch the first one before watching the sequel, just so you'll have an idea what Stanley is like and get a little history background.

[Denil et al. 2014]



ConvNets are everywhere



Whale recognition, Kaggle Challenge

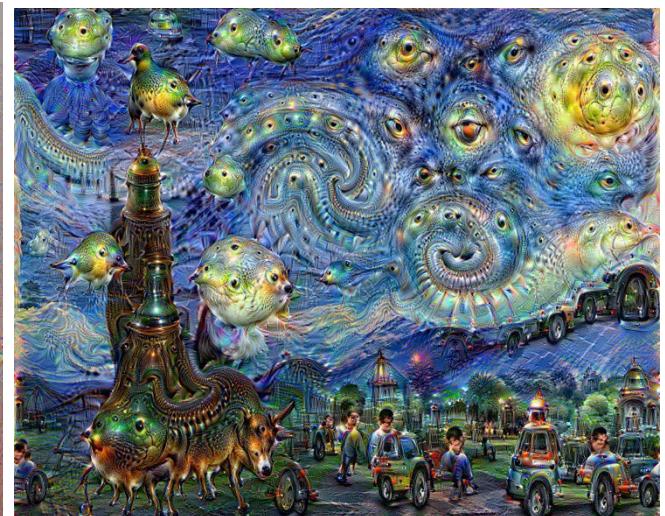
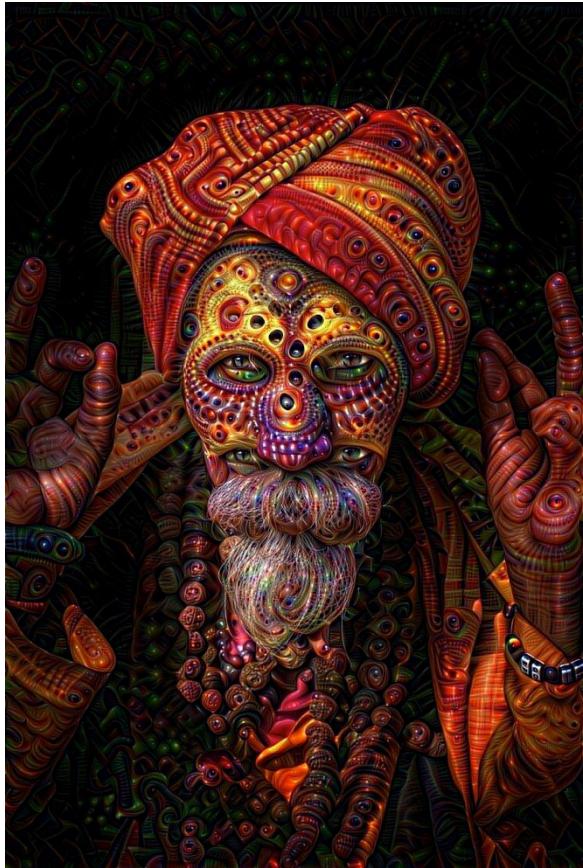


Mnih and Hinton, 2010

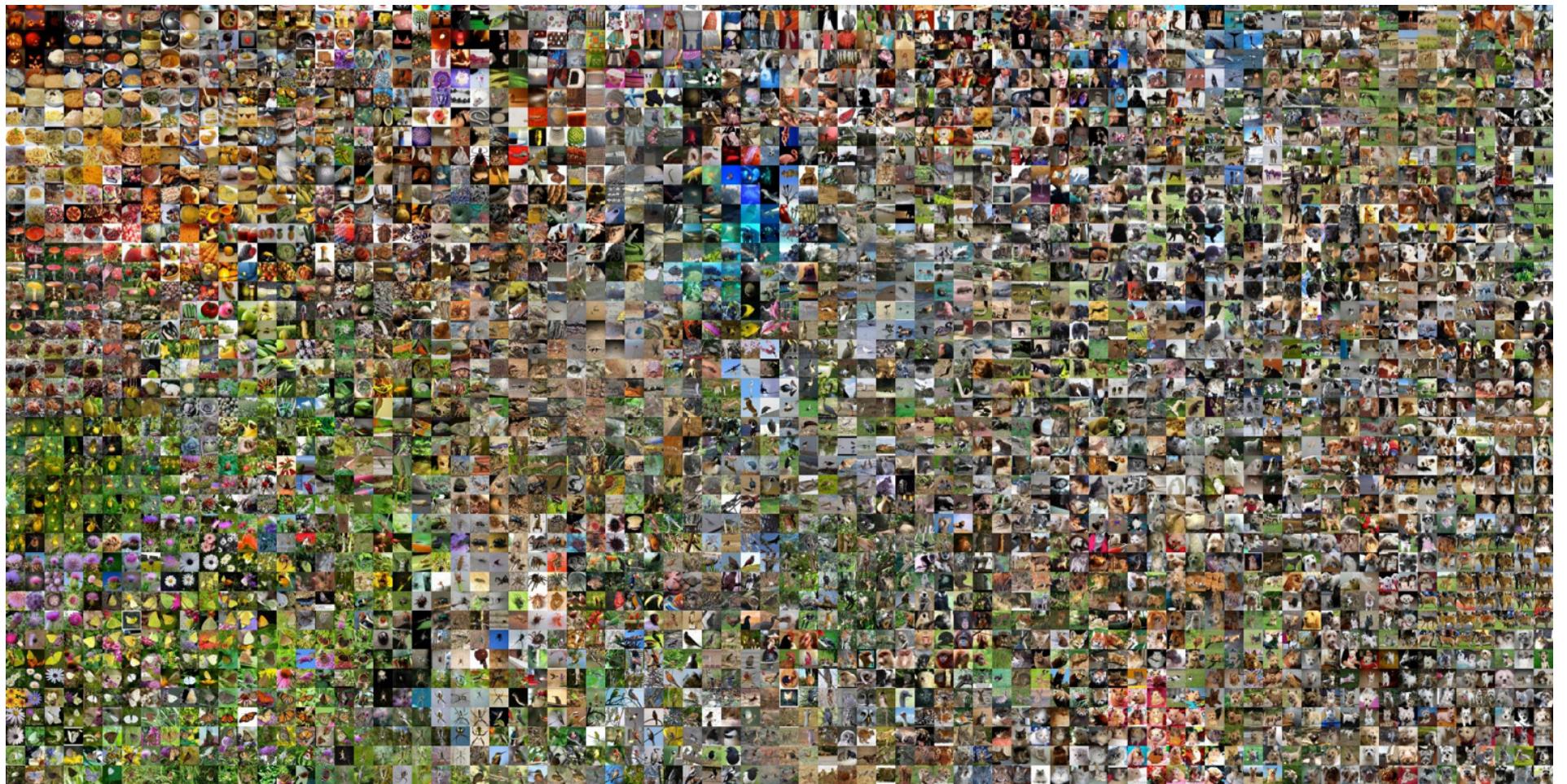
Image Captioning

Describes without errors	Describes with minor errors	Somewhat related to the image	Unrelated to the image
 A person riding a motorcycle on a dirt road.	 Two dogs play in the grass.	 A skateboarder does a trick on a ramp.	 A dog is jumping to catch a frisbee.
 A group of young people playing a game of frisbee.	 Two hockey players are fighting over the puck.	 A little girl in a pink hat is blowing bubbles.	 A refrigerator filled with lots of food and drinks.
 A herd of elephants walking across a dry grass field.	 A close up of a cat laying on a couch.	 A red motorcycle parked on the side of the road.	 A yellow school bus parked in a parking lot.

[Vinyals et al., 2015]



[reddit.com/r/deepdream](https://www.reddit.com/r/deepdream)



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 6 - 69

25 Jan 2016

Next Time

- Final Contest results
- Robot butlers
- Where to go next to learn more about AI