

Modern Robotics: Evolutionary Robotics

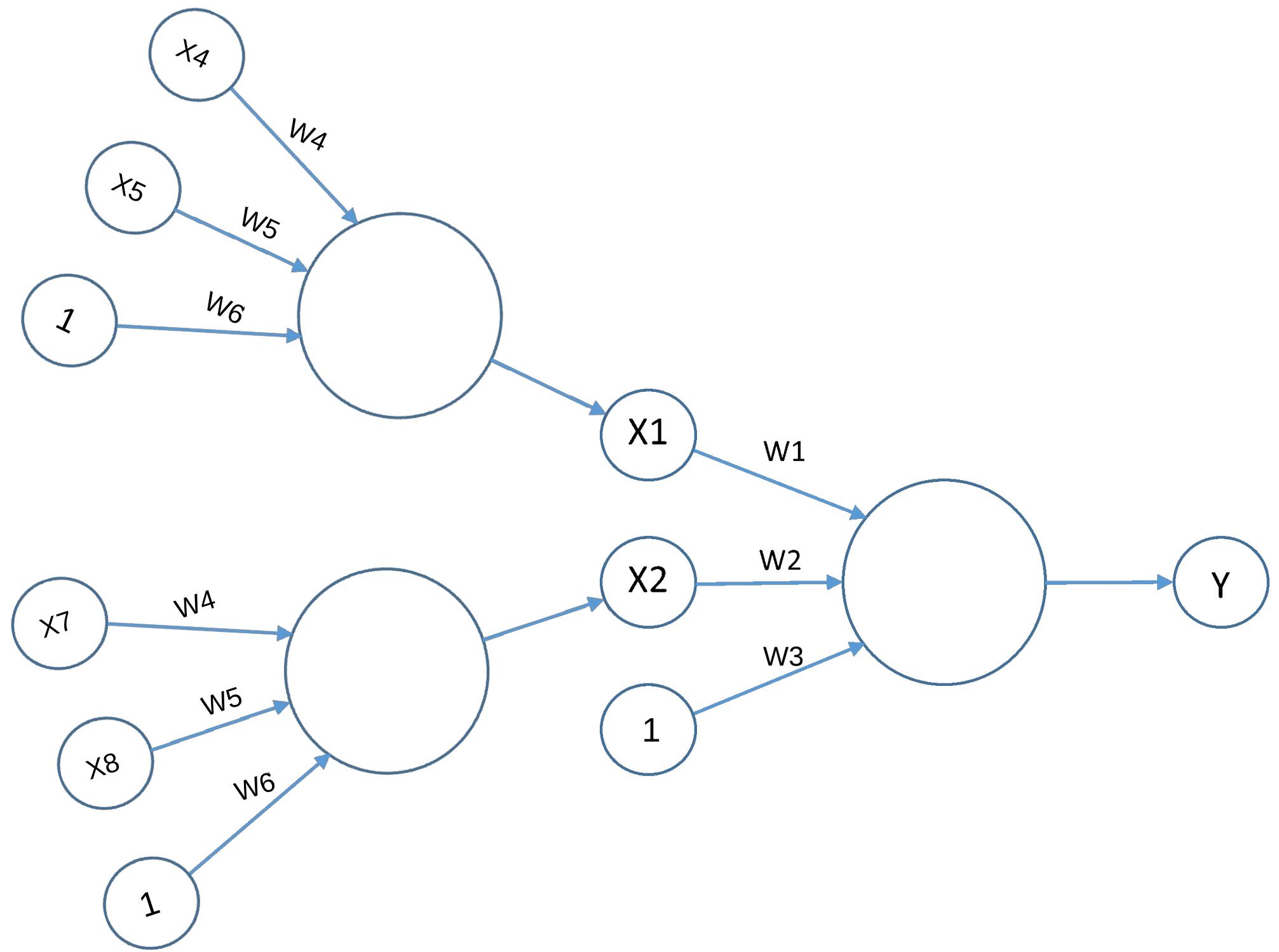
COSC 4560 / COSC 5560

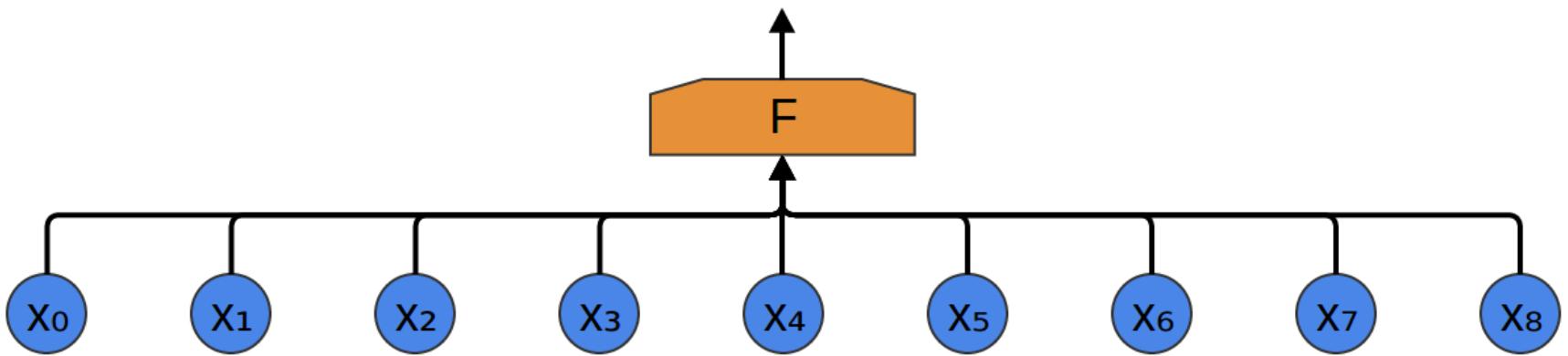
Professor Cheney
4/18/18

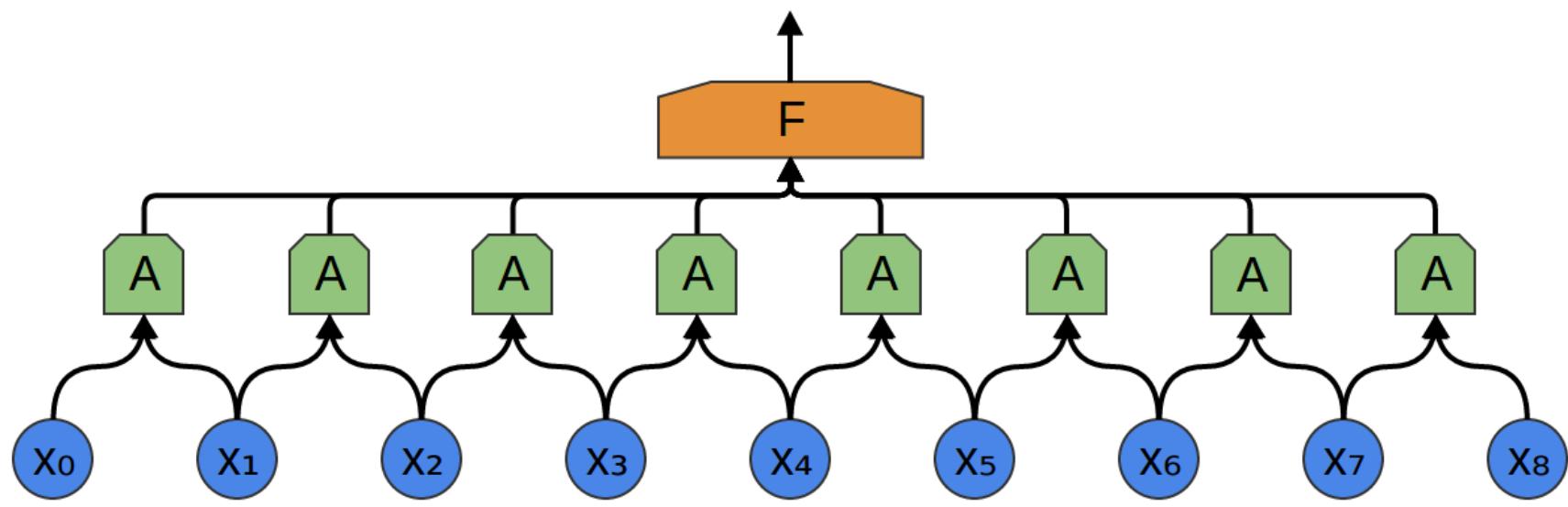
**Tricks that make
this work/scale!**

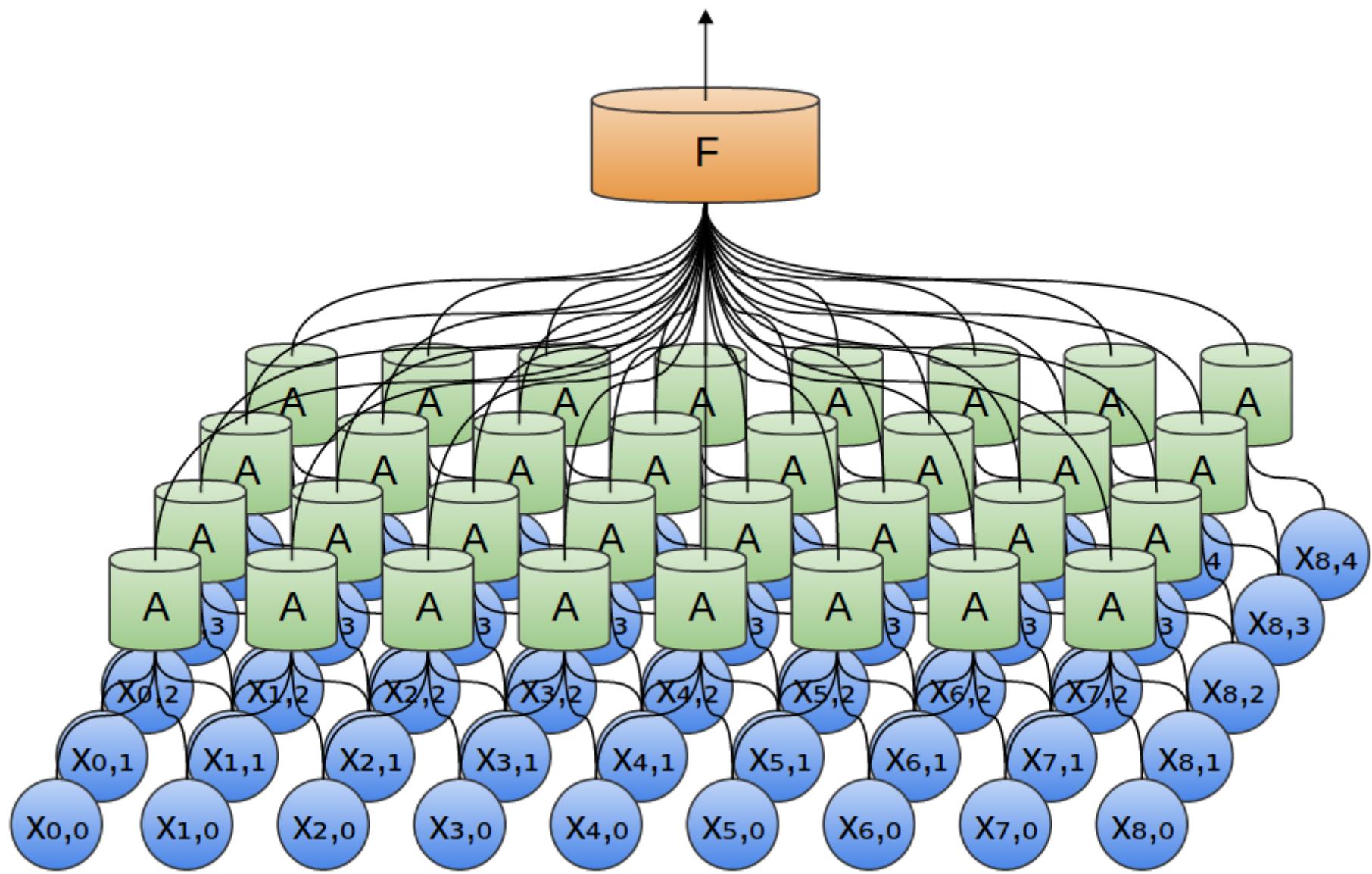
technology idea digital upwards smartphone algorithms mind training
deep concept machine education data information learn
learning recognizes system social drawing research screen smart internet tech future match network massive neural computer mobile robot **artificial** learned intellect facial vivid presentation

Weight Sharing (via convolution)

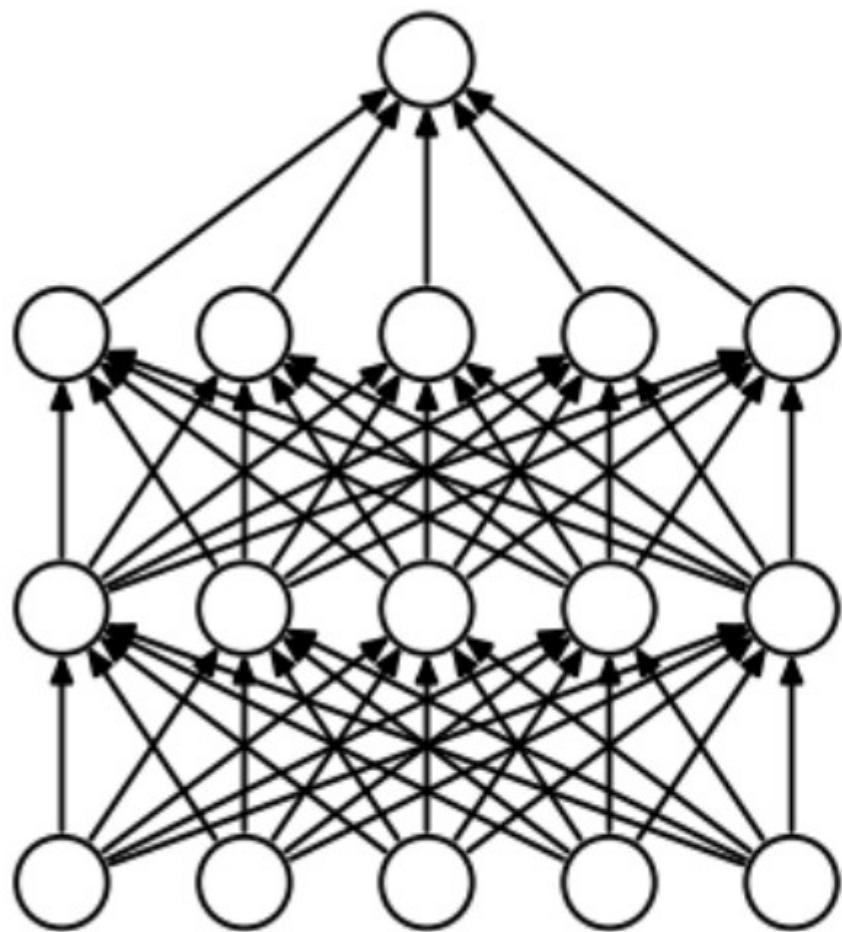




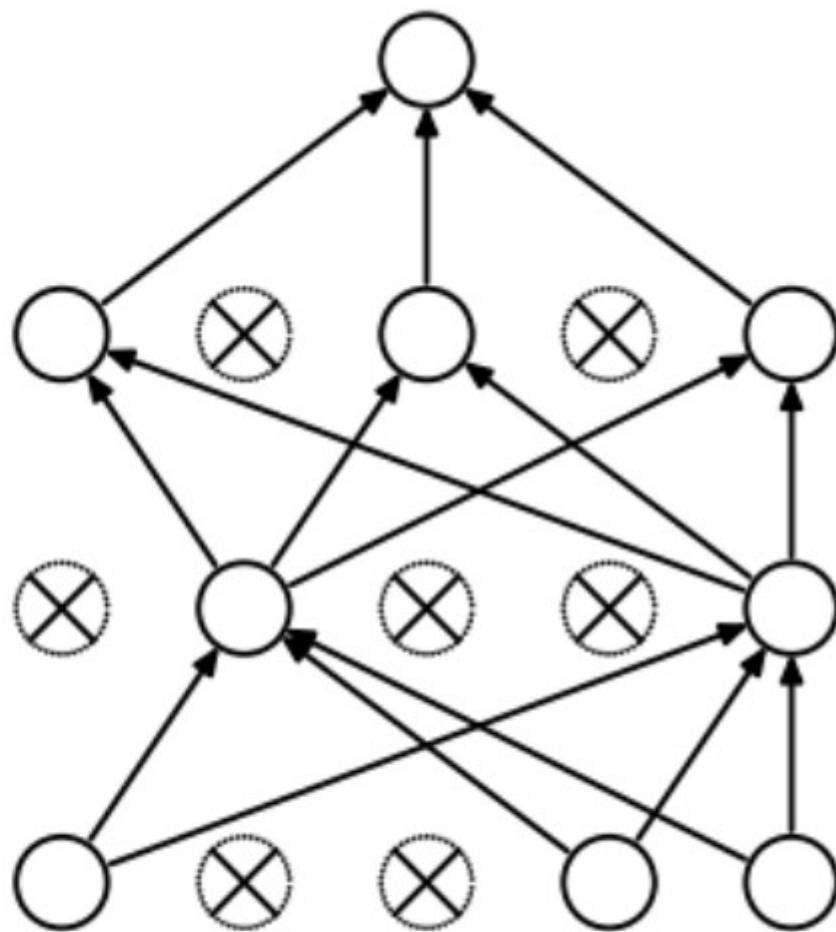




Dropout



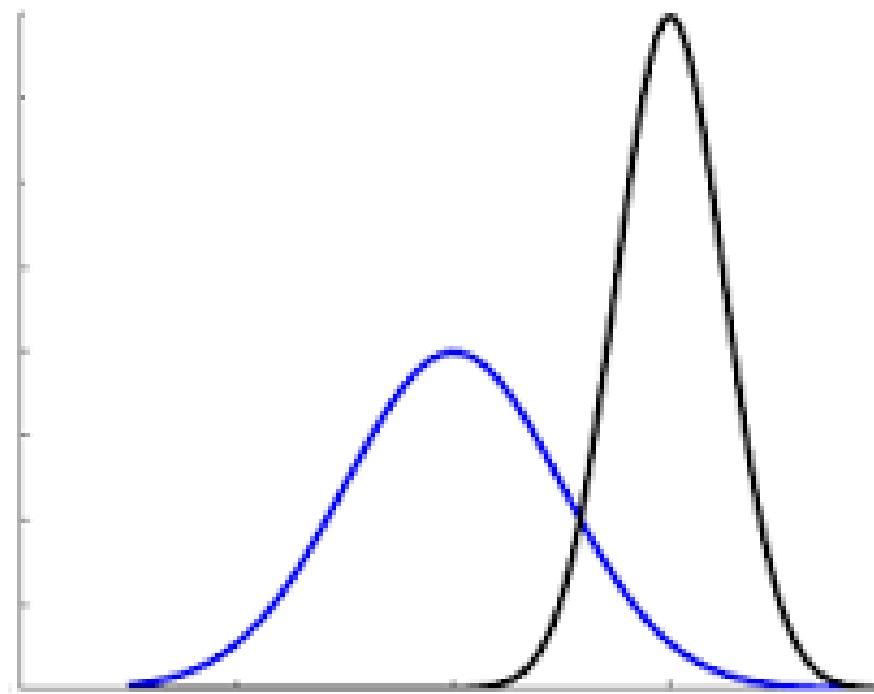
(a) Standard Neural Net



(b) After applying dropout.

Batch Normalization

Internal Covariate Shift

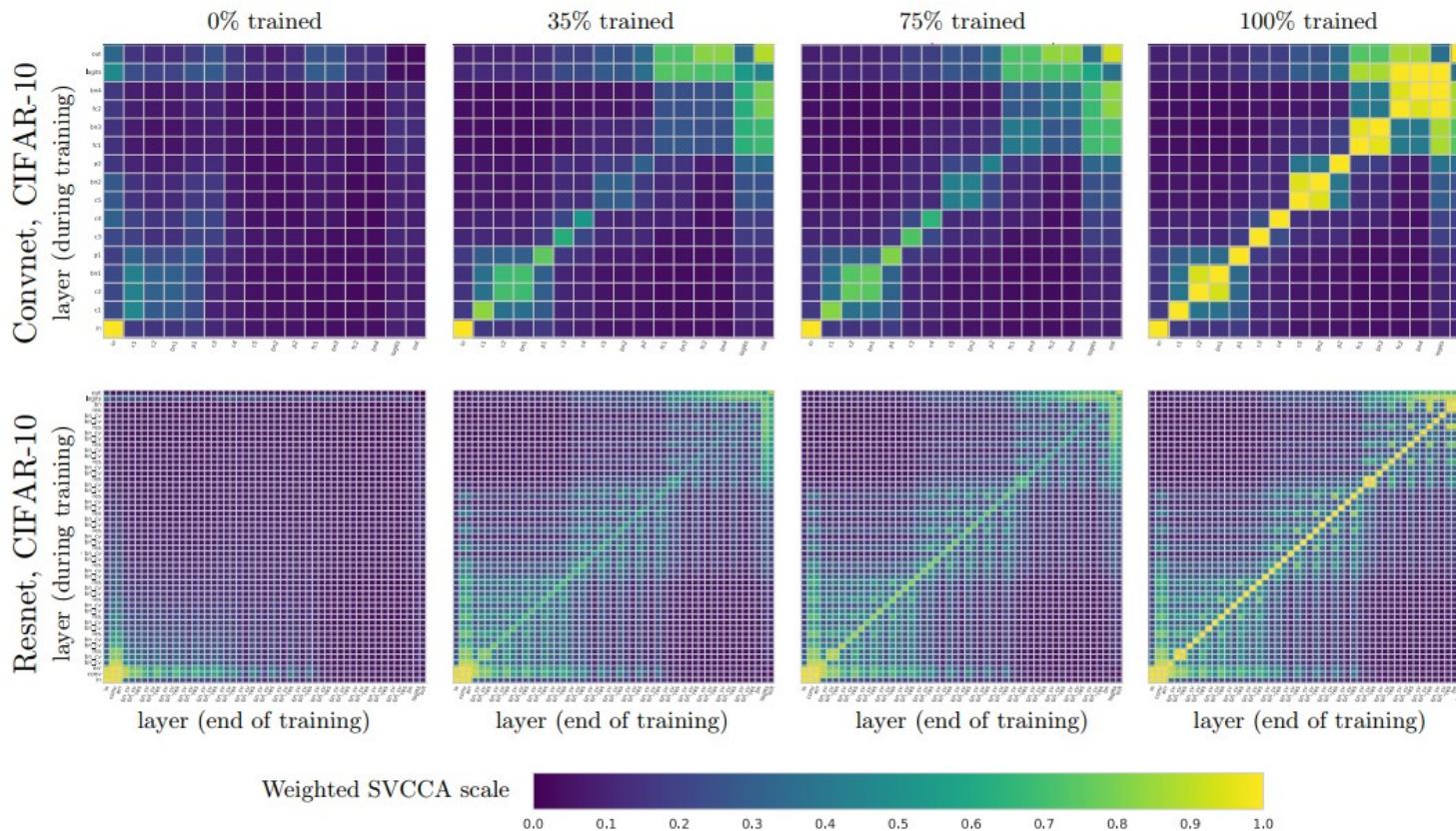


SVCCA: Singular Vector Canonical Correlation Analysis for Deep Learning Dynamics and Interpretability

Maithra Raghu,^{1,2} Justin Gilmer,¹ Jason Yosinski,³ & Jascha Sohl-Dickstein¹

¹Google Brain ²Cornell University ³Uber AI Labs

maithrar@gmail.com, gilmer@google.com, yosinski@uber.com, jaschasd@google.com



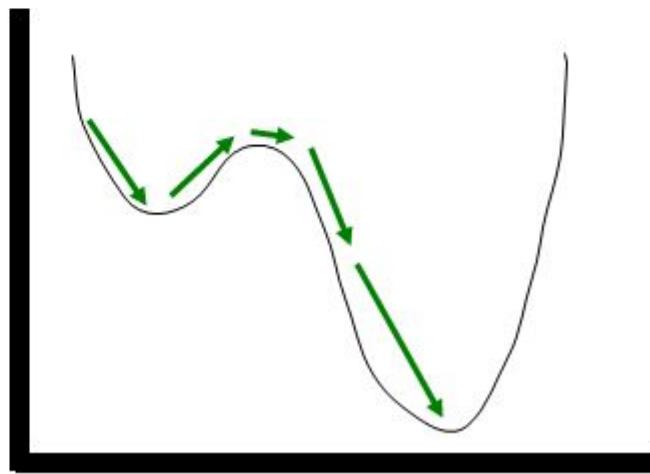
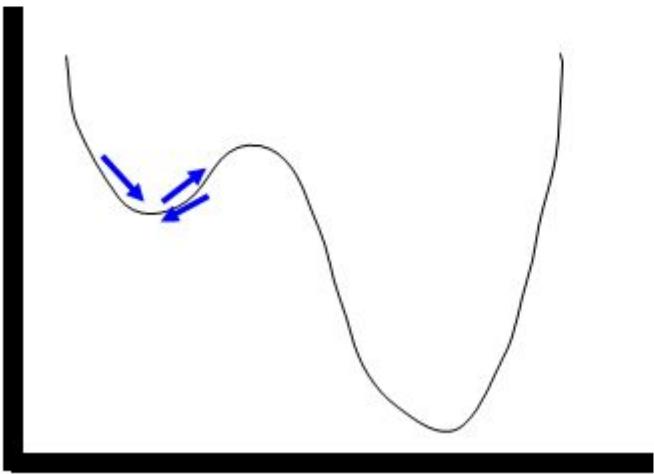
1. compute the empirical mean and variance independently for each dimension.

2. Normalize

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

Momentum

$$\Delta w_{ij}(k) = -\eta \nabla E + \beta \cdot \Delta w_{ij}(k-1)$$



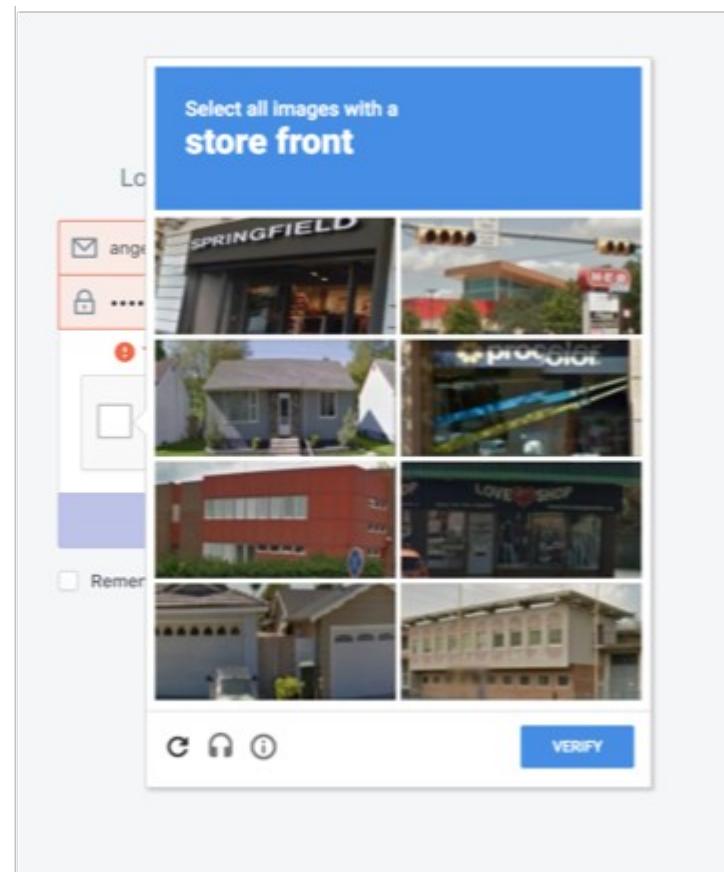
Big Data

I'm not a robot 

Type the text



   Verify



ImageNet Challenge

IMAGENET

- 1,000 object classes (categories).
- Images:
 - 1.2 M train
 - 100k test.

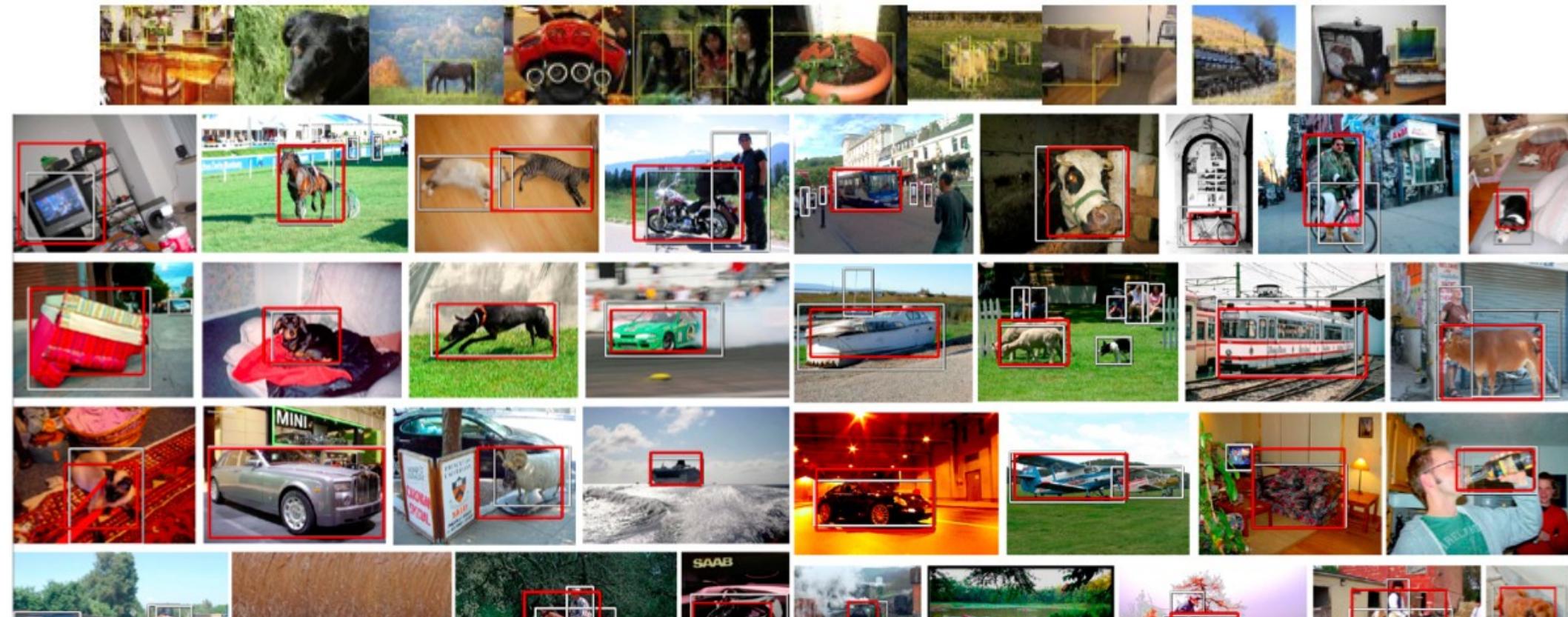




Visual Object Classes Challenge 2009 (VOC2009)

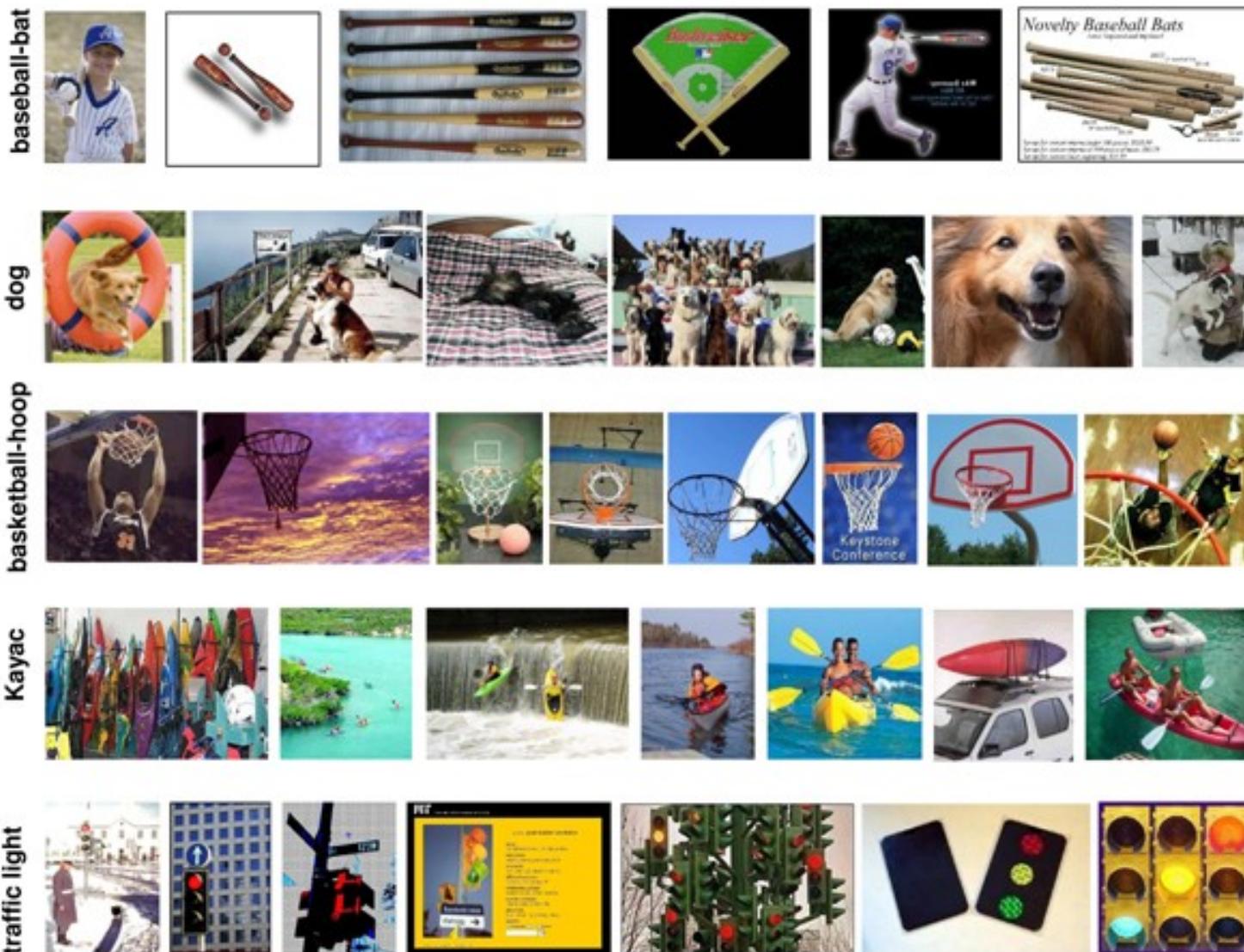


20 classes. The train/val data has 11,530 images containing 27,450 ROI annotated objects and 6,929 segmentations.



Caltech 256

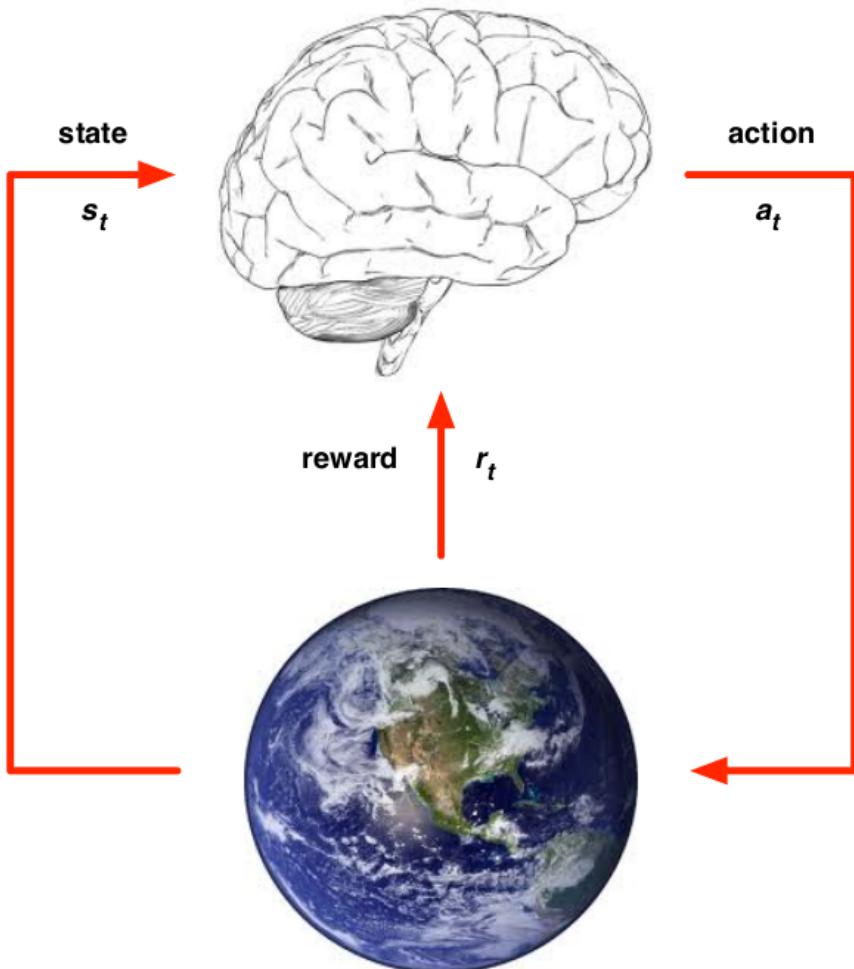
(256 object categories, 30607 images)



Deep Reinforcement Learning

David Silver, Google DeepMind

Agent and Environment



- ▶ At each step t the agent:
 - ▶ Receives state s_t
 - ▶ Receives scalar reward r_t
 - ▶ Executes action a_t
- ▶ The environment:
 - ▶ Receives action a_t
 - ▶ Emits state s_t
 - ▶ Emits scalar reward r_t

Approaches To Reinforcement Learning

Policy-based RL

- ▶ Search directly for the **optimal policy** π^*
- ▶ This is the policy achieving maximum future reward

Value-based RL

- ▶ Estimate the **optimal value function** $Q^*(s, a)$
- ▶ This is the maximum value achievable under any policy

Model-based RL

- ▶ Build a transition model of the environment
- ▶ Plan (e.g. by lookahead) using model

Deep Reinforcement Learning

- ▶ Can we apply deep learning to RL?
- ▶ Use deep network to represent value function / policy / model
- ▶ Optimise value function / policy /model **end-to-end**
- ▶ Using stochastic gradient descent

learning value networks

Bellman Equation

- ▶ Value function can be unrolled recursively

$$\begin{aligned} Q^\pi(s, a) &= \mathbb{E} [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \mid s, a] \\ &= \mathbb{E}_{s'} [r + \gamma Q^\pi(s', a') \mid s, a] \end{aligned}$$

- ▶ Optimal value function $Q^*(s, a)$ can be unrolled recursively

$$Q^*(s, a) = \mathbb{E}_{s'} \left[r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right]$$

- ▶ **Value iteration** algorithms solve the Bellman equation

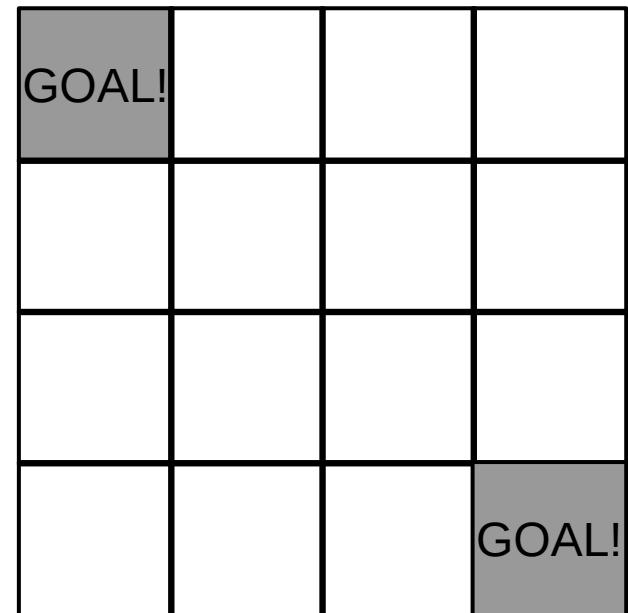
$$Q_{i+1}(s, a) = \mathbb{E}_{s'} \left[r + \gamma \max_{a'} Q_i(s', a') \mid s, a \right]$$

grid world:

each timestep has -1 reward

the game terminates when
you reach a goal state

actions: N, S, E, W



intuitive description: “get to the goal as soon as possible”
(but let's pretend we're a robot, who doesn't know this!)

each value function (V) is defined
with respect to some behavioral policy (π)
 V^π

let's iteratively find V^π for a random policy
in our mini grid world

current value (V_k) for a random policy

$k=0$

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

$k=1$

current value (V_k) for a random policy

$k=0$

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

$k=1$

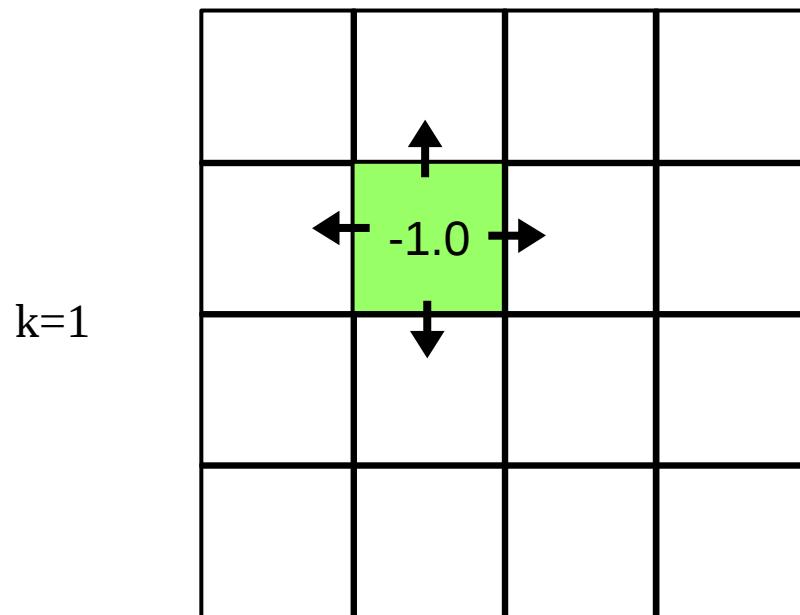
A 4x4 grid representing state values at time step $k=1$. The central cell is highlighted in green and contains a question mark, indicating uncertainty or a pending update. Four arrows point outwards from this central cell, representing possible transitions to adjacent states.

		?	

current value (V_k) for a random policy

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

$k=0$



$k=1$

current prediction for cumulative reward in new state
immediate reward

$$N: V_{s,N} = -1 + 0 = -1$$

$$S: V_{s,S} = -1 + 0 = -1$$

$$E: V_{s,E} = -1 + 0 = -1$$

$$W: V_{s,W} = -1 + 0 = -1$$

with a random policy,
we are equally likely to take any move,
so:

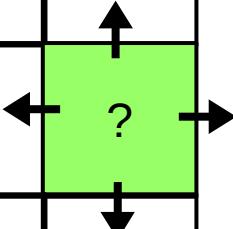
$$V_s = (-1 + -1 + -1 + -1)/4 = -1$$

current value (V_k) for a random policy

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

$k=0$

$k=1$



current value (V_k) for a random policy

$k=0$

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

$k=1$

			-1.0

current value (V_k) for a random policy

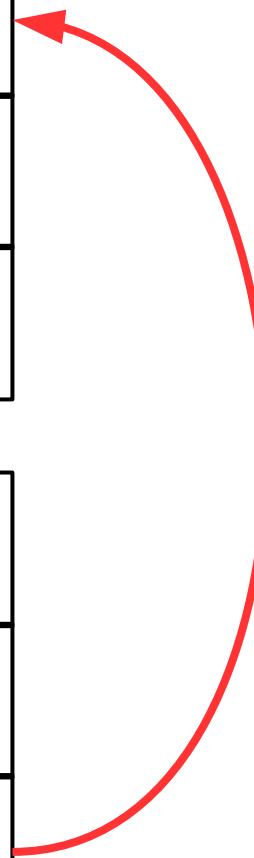
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

$k=0$

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0

$k=1$

next iteration...
new value function
becomes
old value function
("current prediction
for cumulative reward")



current value (V_k) for a random policy

$k=1$

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0

$k=2$

current value (V_k) for a random policy

$k=1$

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0

$k=2$

	-1.75		

$$N: V_{s,N} = -1 + -1 = -2$$

$$S: V_{s,S} = -1 + -1 = -2$$

$$E: V_{s,E} = -1 + 0 = -1$$

$$W: V_{s,W} = -1 + -1 = -2$$

$$V_s = (-2 + -2 + -1 + -2)/4 = -1.75$$

current value (V_k) for a random policy

$k=1$

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0

$k=2$

0.0	-1.75	-2.0	-2.0
-1.75	-2.0	-2.0	-2.0
-2.0	-2.0	-2.0	-1.75
-2.0	-2.0	-1.75	0.0

current value (V_k) for a random policy

$k=2$

0.0	-1.75	-2.0	-2.0
-1.75	-2.0	-2.0	-2.0
-2.0	-2.0	-2.0	-1.75
-2.0	-2.0	-1.75	0.0

$k=3$

0.0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0.0

current value (V_k) for a random policy

$k=10$

0.0	-6.1	-8.4	-9.0
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9.0	-8.4	-6.1	0.0

$k=\infty$

0.0	-14	-20	-22
-14	-18	-20	-20
-20	-20	-18	-14
-22	-20	-14	0.0

converged to true
value function
($V^{\pi\text{-random}}$)

