**Deep Q Learning**



inputs    hidden    output

s

a

$$\hat{q}(s,a)$$

$\hat{q}((180,4,10,24),0) = .38$

$\hat{q}((180,4,10,24),1) = .54$
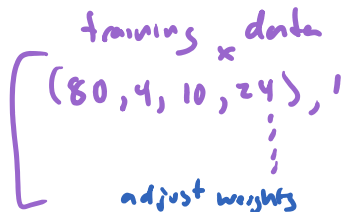
$\hat{q}((180,4,10,24),2) = .21$

$\hat{q}((180,4,10,24),3) = .04$

reward 0    r
new state $(64,4,10,22)$   s'
$$v(s') = \max_{a'} q(s',a')$$
$$\hat{v}(s') = \max_{a'} \hat{q}(s',a')$$
$$= 0.68$$

from target network

training data
$$\begin{bmatrix} (80,4,10,24),1 & \underset{Y}{0.68} \\ \vdots & \end{bmatrix}$$

adjust weights
↑ on this            → produces target output

initialize   learning, target   networks

for each iteration

    for each of n episodes
      for each event
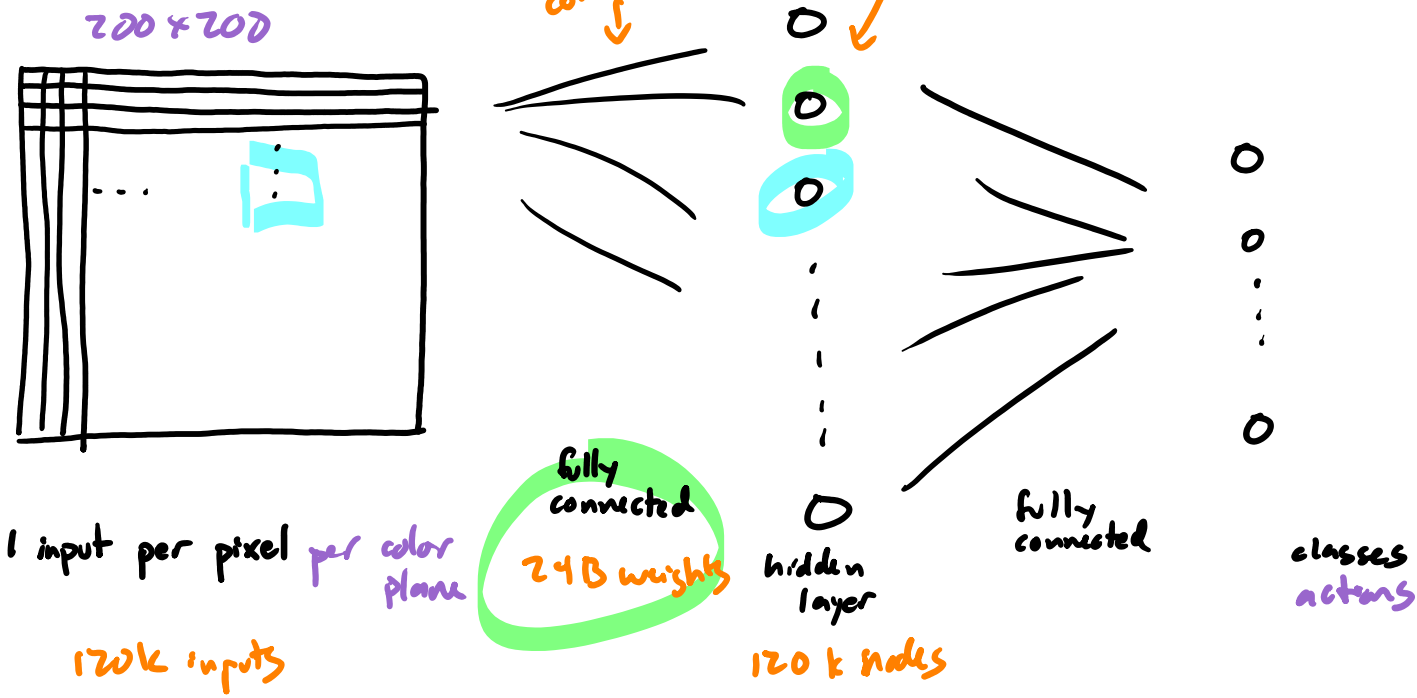        add   $(s,a,s',r)$   to   replay database

    Sample  replay database

    train learning  network  toward   $r + \max_{a'} \hat{q}_{target}(s',a')$

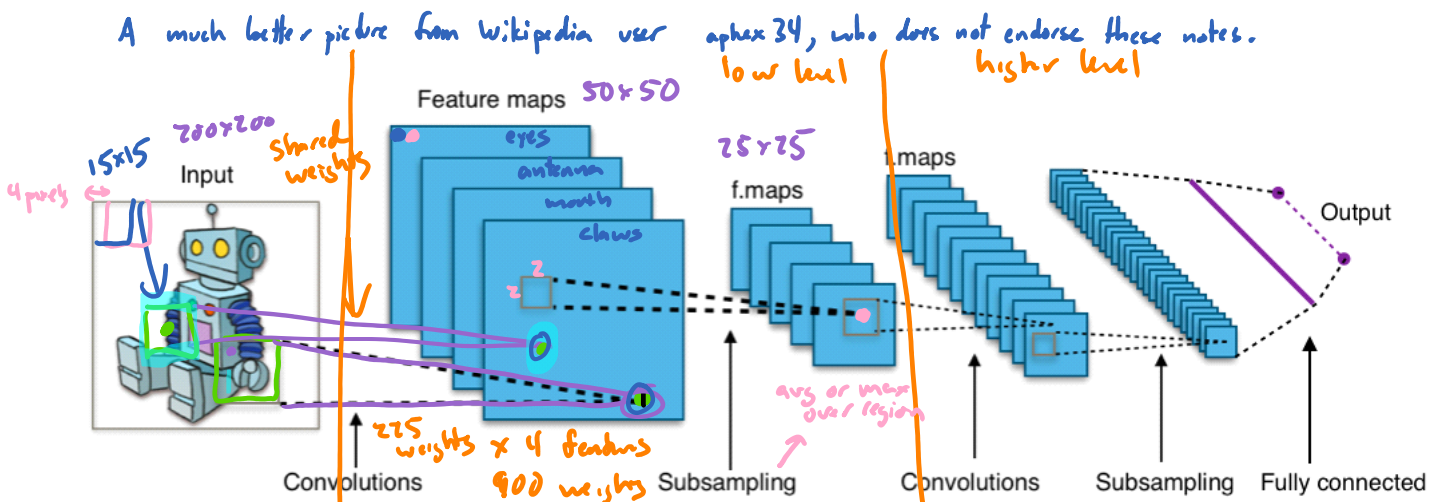    if enough time passed
      copy learning network to target network

# ANNs for Images

200 × 200

1 input per pixel per color plane

120k inputs

1000M connections

fully connected

24B weights

hidden layer

120k nodes

$$\sum_i w_i \cdot L_i$$

fully connected

classes actions

# Convolutional Neural Networks

Deep Q network learning to play Pong



A much better picture from Wikipedia user aphex 34, who does not endorse these notes.



https://upload.wikimedia.org/wikipedia/commons/6/63/Typical_cnn.png

SIMD
single instruction
multiple data

**AlphaGo (2014-2017)**   DeepMind

Step 1: Supervised learning for convolutional <u>deep</u> neural network

use database from games
of expert players

13 layers

input:
$19 \times 19 \times 48$
locations  features

black
white
empty
#opp captured
#own captured
liberties
ladder capture
ladder escape

— matched 55% of time

+ smaller (faster) 25% of time

output: a move ($19 \times 19 + 1$)

hand-coded
features

3 weeks

Step 2: reinforcement learning for convolutional <u>deep</u> neural network

1 day

beat SL network    80% of time

Step 3: reinforcement learning for value network

+1 black win
0
-1 white win

Step 4: (MCTS)

Elo

2015          2016          2017
(Fan Hui)    (Lee Sedol)   (retired)

△ Elo    → higher rated player has          chance of winning