**Deep Q Learning**

inputs → ANN → outputs

$q(s,1)$
$q(s,2)$
$\vdots$
$q(s,n)$

s (state)
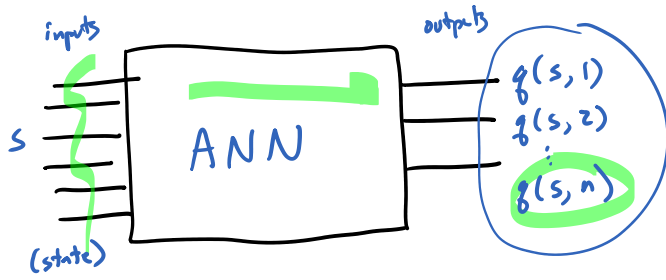
$q(s,1) = .32$
$q(s,2) = .47$
$q(s,3) = .16$

choose $a = 2$    $s'$ has value $.46$
              $r = 0.02$

$q(s,2)$ updated toward
                $0.02 + \gamma \cdot 0.46$

in training example

$$\frac{x}{s} \qquad \frac{Y}{[q(s,1), 0.02 + \gamma \cdot 0.46, q(s,3)]}$$

train this one

produces target output (what to train toward)

$q(s,a) \rightarrow s', r$

$r + \gamma \max_{a'} q(s', a')$

$v(s')$

initialize learning, target networks

for each iteration

   for each of n episodes — start each at initial state
     for each event
       add $(s, a, s', r)$ to replay database

   sample replay database

   train learning network → target output $\left[ \hat{q}_{learn}(s,1), \ldots, r + \gamma \max_{a'} \hat{q}_{target}(s,a), \ldots, \hat{q}_{learn}(s,n) \right]$

   if enough time passed
     copy learning network to target network
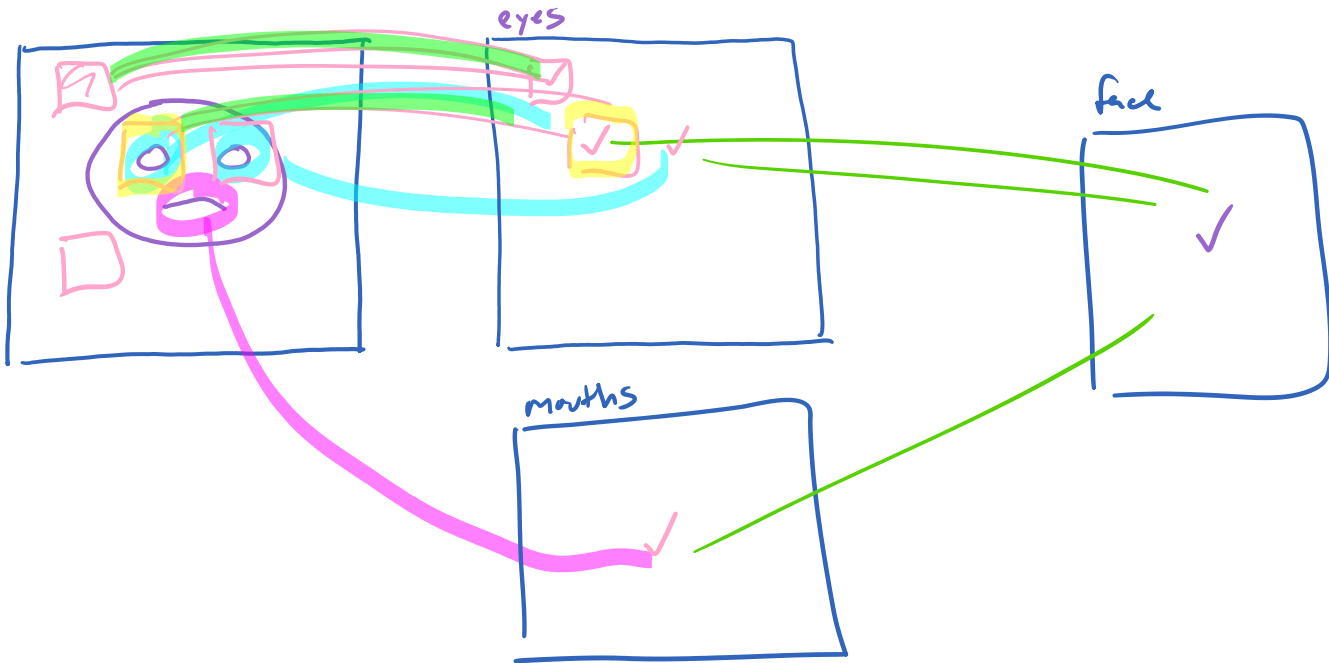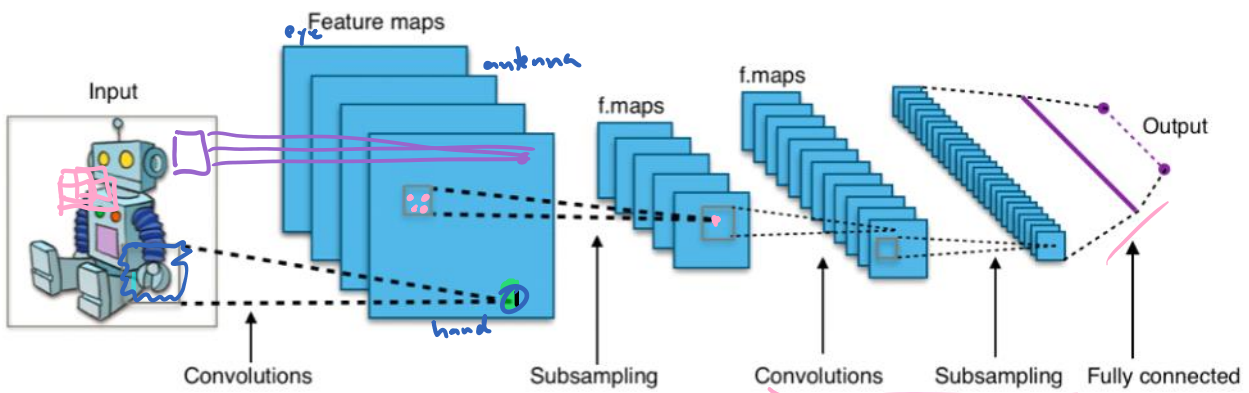
# ANNs for Images

200×200

1 input per pixel per color per pixel

120k inputs

fully connected

hidden layer

14400m 120k nodes

14B weights

fully connected

classes

**Convolutional Neural Networks** → NNs for image processing

Deep Q network learning to play Pong



eyes

face

mouths

A much better picture from Wikipedia user aphex34, who does not endorse these notes.



Input — Feature maps — eye — antenna — f.maps — f.maps — hand — Output
Convolutions — Subsampling — Convolutions — Subsampling — Fully connected

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

Step 1: Supervised learning for convolutional deep neural network (policy network)

DB of expert-level games

13 layers          input: 19×19×48

output: action
(Mx19x19 outputs)

~ matched 55% of time

+ smaller (faster) network 25% of time

3 weeks

black
white
empty
#opp captured
#own captured
liberties
ladder capture
ladder esc
⋮
⋮

Step 2: reinforcement learning for convolutional deep neural network

beat SL 80% of time

1 day

Step 3: reinforcement learning for value network

output is $v(s)$

using step 2 network
plays itself 30M times
sample 1 pos/game

Step 4: MCTS

default: use fast network from step 1

initialize new node's value using step 3 value network

tree policy:     $q(s,a)$  +  $c\, P(s,a) \cdot \dfrac{\sqrt{\#\text{parents visited}}}{1 + \#\text{child visits}}$

exploit
observations

from larger
step 1 network

Elo   $3144 \rightarrow 3739 \rightarrow 5185$

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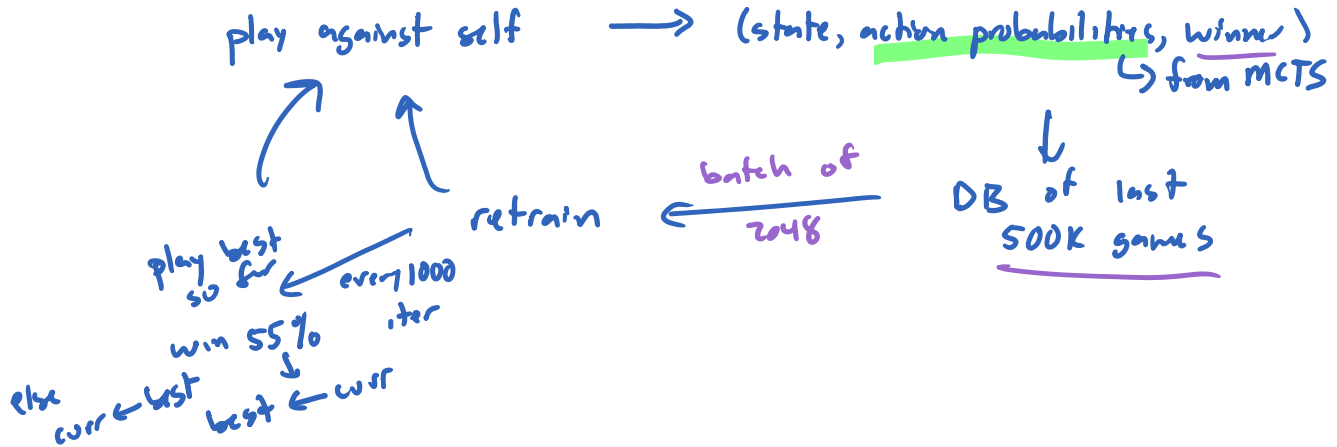(Fan Hui)    (Lee Sedol)   (retired)

Δ Elo  400 → higher rated player has 90%⁺ chance of winning

**AlphaGo Zero**

no prior knowledge

input :   19 × 19 × 17     current pos + last 7 pos
                                        + turn (all 1 = black
                                               0 = white)

output :   move (19×19+1)  and   value [~1,   +1]

play against self  ⟶  (state, action probabilities, winner)
                                                        ↳ from MCTS

play best so far ⟵ retrain ⟸ batch of 2048    DB of last 500K games

every 1000 iter

win 55% 
else curr ⟵ best     best ⟵ curr

Deep Reinforcement Learning Doesn't Work Yet (alexirpan.com)