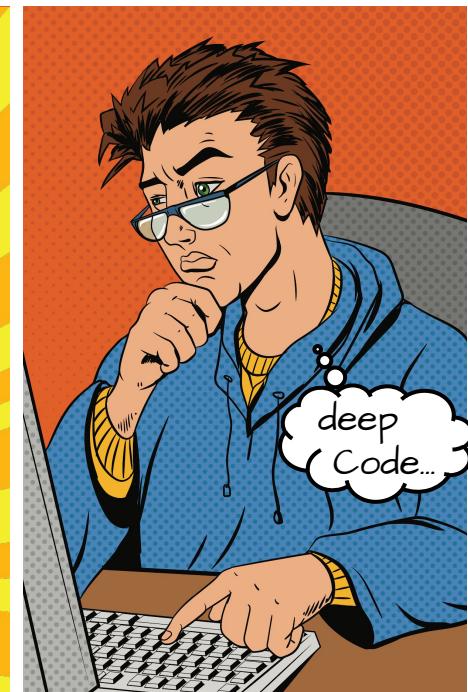
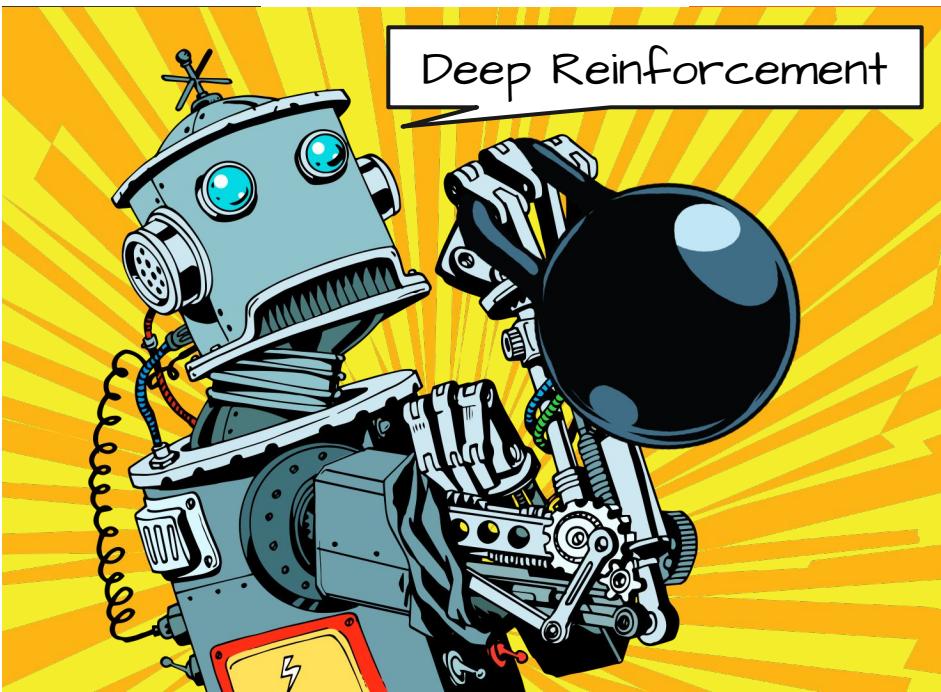
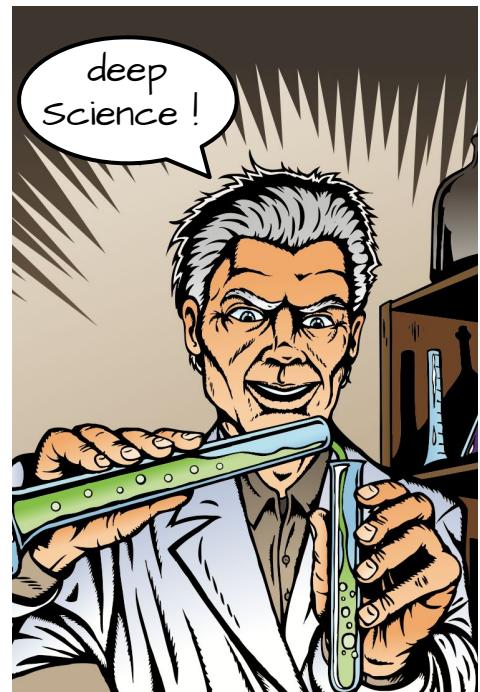
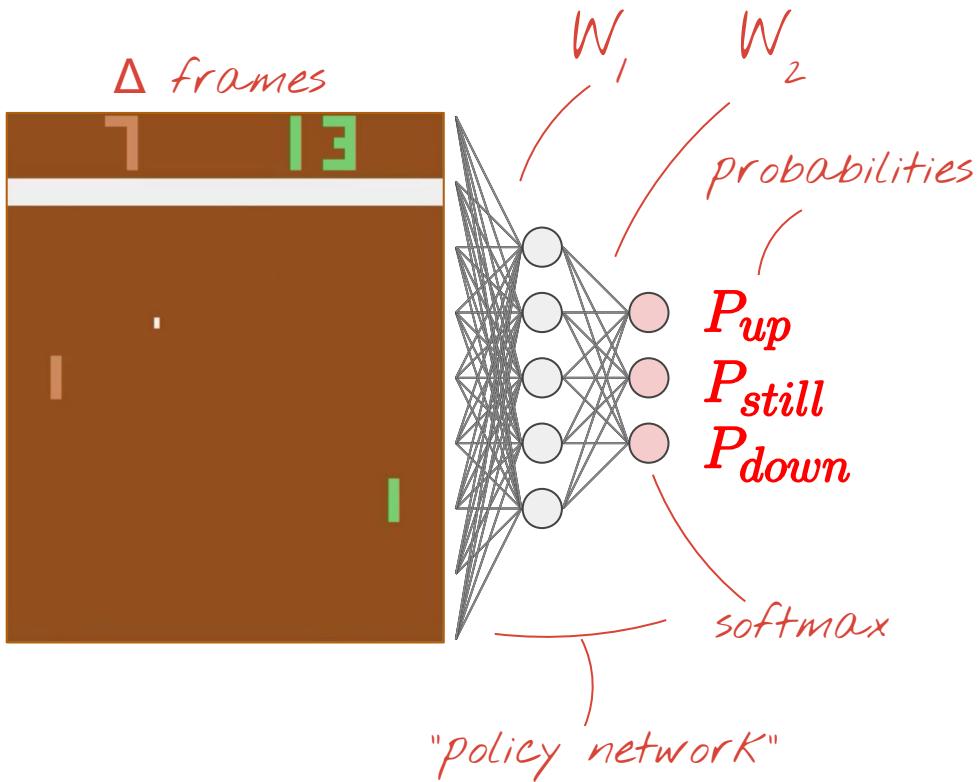


# >TensorFlow and \ deep reinforcement learning without a PhD\_



# Pong ?



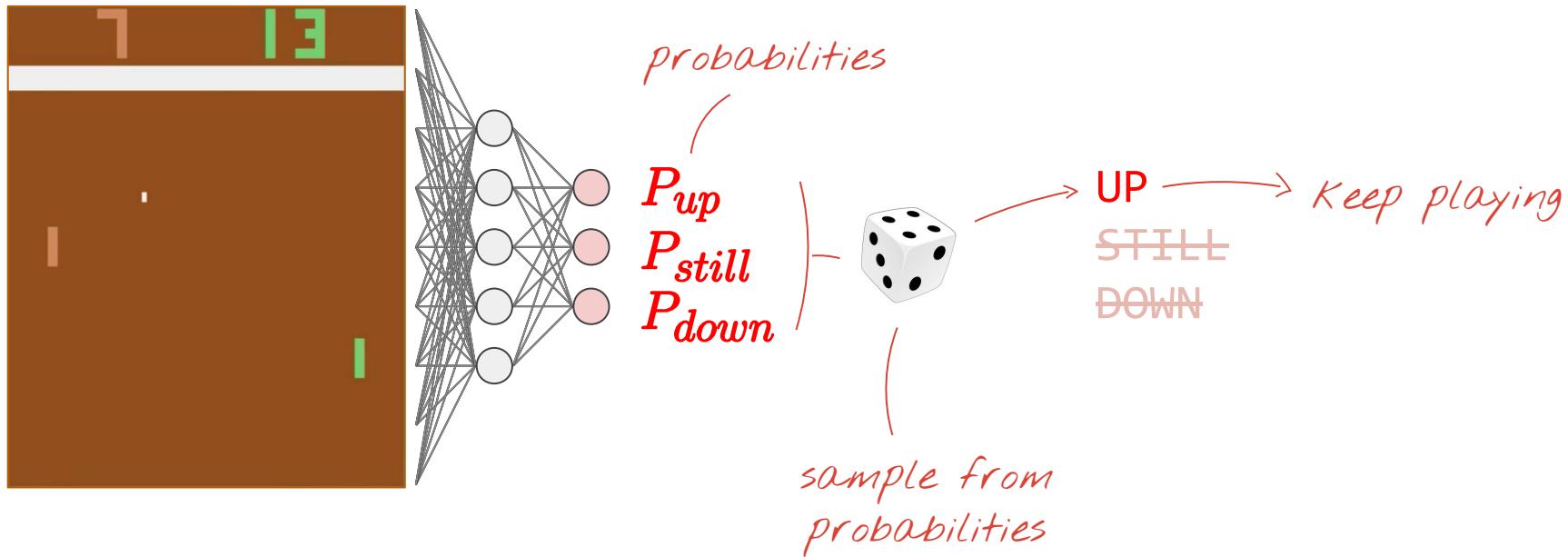
cross-entropy

$$\begin{aligned} \text{loss} = - & (P'_{up} \cdot \log(P_{up}) \\ & + P'_{still} \cdot \log(P_{still}) \\ & + P'_{down} \cdot \log(P_{down})) \end{aligned}$$

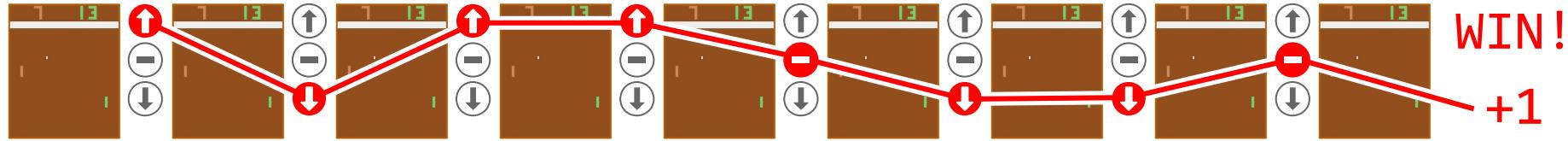
"correct move"  
Ex: 1 0 0

Adapted from A. Karpathy's "[Pong from Pixels](#)" post

# Pong ?



# Policy gradients



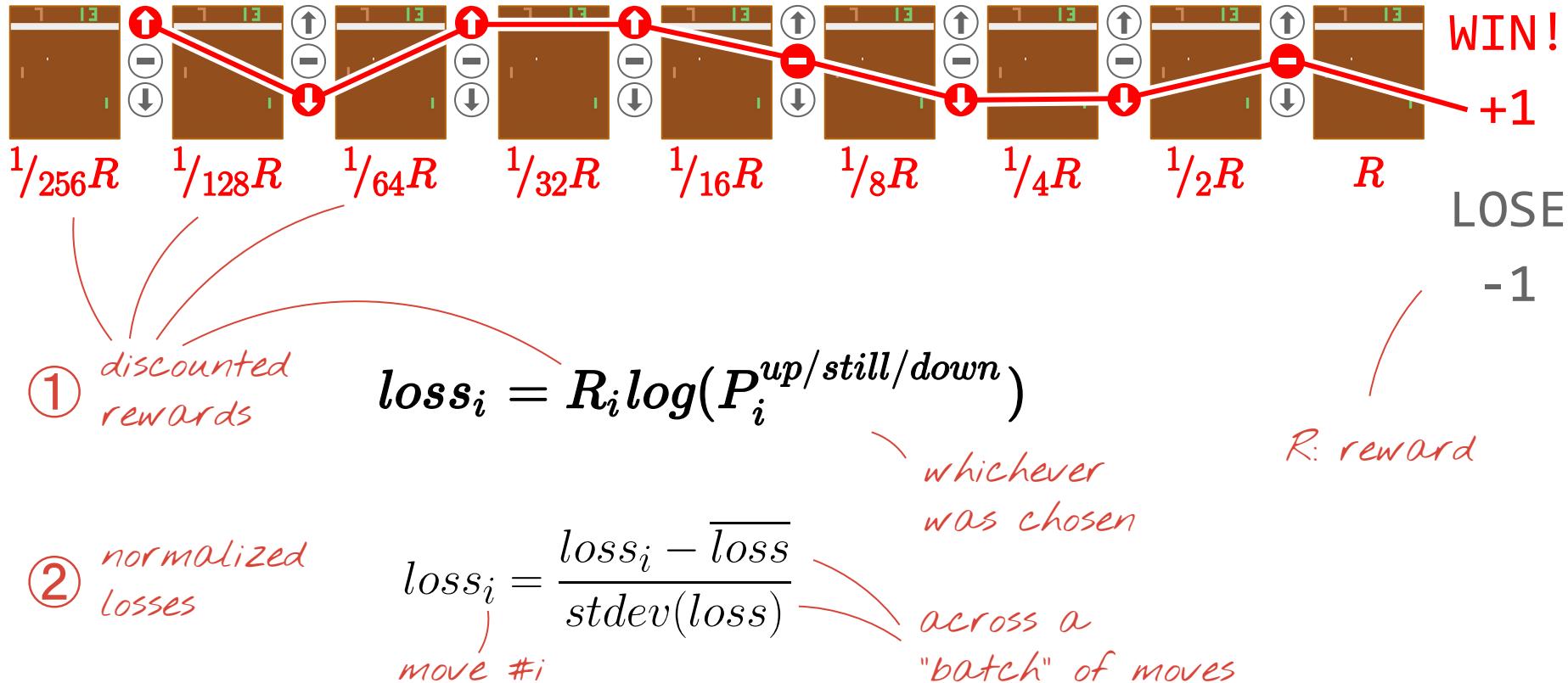
for each move

$$\begin{aligned} loss = & -R \left( P'_{up} \log(P_{up}) \right. \\ & + P'_{down} \log(P_{down}) \\ & \left. + P'_{down} \log(P_{down}) \right) \end{aligned}$$

"sampled move"  
Ex: 1 0 0      predicted

R: reward

# Policy gradient refinements



# Training data

move #	[UP, STILL, DOWN]	Policy network probabilities	Discounted rewards
0	[1, 0, 0]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	+1*d <sup>7</sup>
1	[0, 1, 0]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	+1*d <sup>6</sup>
2	[0, 1, 0]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	+1*d <sup>5</sup>
3	[0, 0, 1]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	+1*d <sup>4</sup>
4	[0, 0, 1]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	+1*d <sup>3</sup>
5	[0, 0, 1]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	+1*d <sup>2</sup>
6	[0, 1, 0]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	+1*d
7	[1, 0, 0]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	+1 ← WIN! +1
8	[0, 0, 1]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	-1*d <sup>7</sup>
9	[0, 0, 1]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	-1*d <sup>6</sup>
10	[0, 0, 1]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	-1*d <sup>5</sup>
11	[0, 0, 1]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	-1*d <sup>4</sup>
12	[0, 1, 0]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	-1*d <sup>3</sup>
13	[1, 0, 0]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	-1*d <sup>2</sup>
14	[0, 0, 1]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	-1*d
15	[0, 1, 0]	[ $P_{UP}$ , $P_{STILL}$ , $P_{DOWN}$ ]	-1 ← LOSE -1

Hyperparams:

Batch size: ~20k moves

Discount factor d=0.95~0.99

Optimizer=tf.train.RMSPropOptimizer

- Learning rate lr=0.0001~0.005
- Decay=0.95~0.99

1 hidden layer with 200 units

Beta ("laziness") = 0.01~0.02

# Tensorflow code

```
observations = tf.placeholder(shape=[None, 80x80], dtype=tf.float32) # pixels
actions      = tf.placeholder(shape=[None]) # 0, 1, 2 for UP, STILL, DOWN
rewards      = tf.placeholder(shape=[None]) # +1, -1, with discounts

# model
Y = tf.layers.dense(observations, 200, activation=tf.nn.relu)
Ylogits = tf.layers.dense(Y, 3)

# sample an action from predicted probabilities
sample_op = tf.multinomial(logits=tf.reshape(Ylogits, shape=(1, 3)), num_samples=1)

#loss
cross_entropies = tf.losses.softmax_cross_entropy(one_hot_labels=tf.one_hot(actions,3),
                                                    logits=Ylogits)
loss = tf.reduce_sum(rewards * cross_entropies)

# training operation
optimizer = tf.train.RMSPropOptimizer(learning_rate=0.001, decay=0.99)
train_op = optimizer.minimize(loss)
```

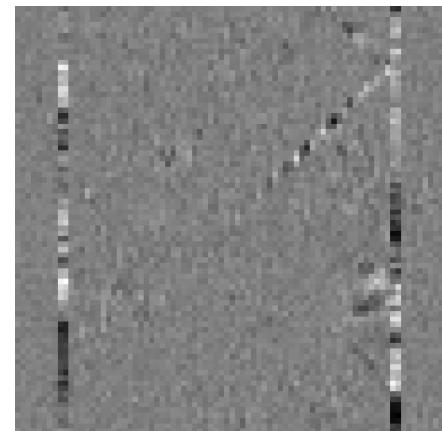
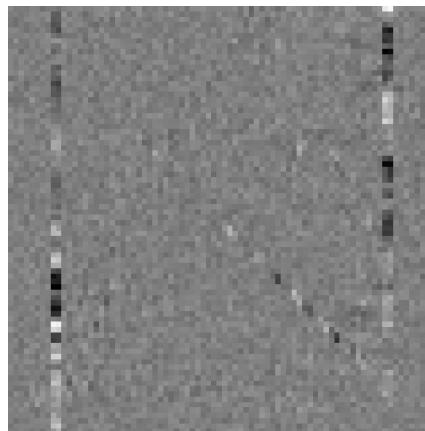
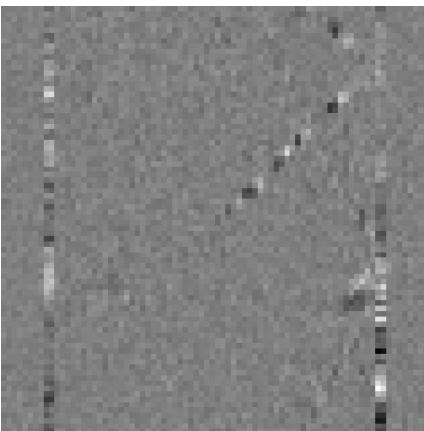
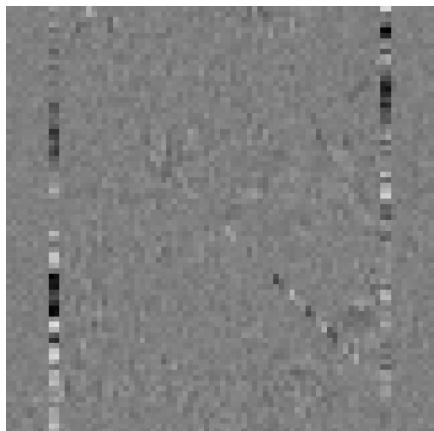
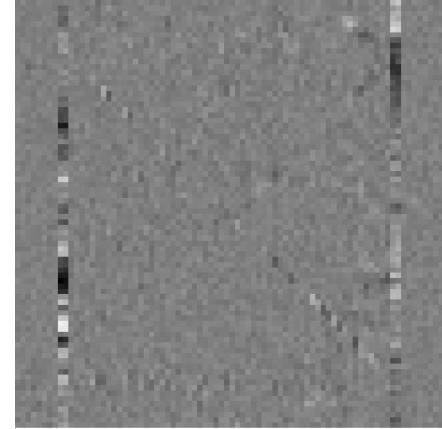
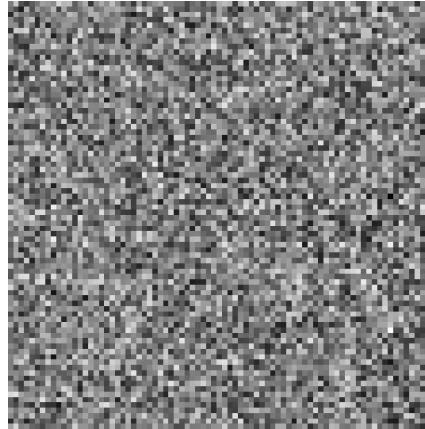
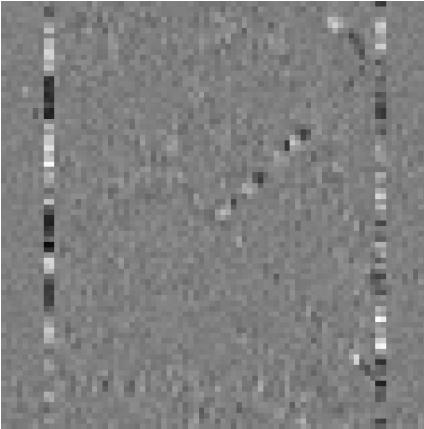
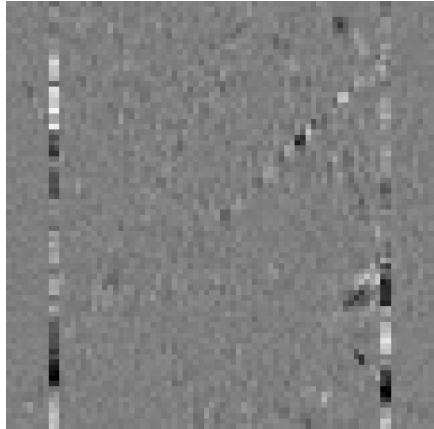
# Training loop

```
with tf.Session() as sess:  
    while len(observations) < BATCH_SIZE:  
        ... # reset everything  
        while not done:          # play a game in 21 points  
            # get pixels  
            current_pix = read_pixels(game_state)  
            observation = current_pix - previous_pix  
            previous_pix = current_pix  
            # decide what move to play: UP, STILL, DOWN (through NN model)  
            action = sess.run(sample_op, feed_dict={observations: [observation]})  
            # play it (through openAI gym pong simulator)  
            game_state, reward, done, info = pong_sim.step(action)  
            # collect results  
            observations.append(observation)  
            actions.append(action)  
            rewards.append(reward)  
  
            # Process the rewards after each episode  
            processed_rewards = discount_rewards(rewards, args.gamma)  
            processed_rewards = normalize_rewards(rewards, args.gamma)  
  
            feed_dict = {observations: observations, actions: actions, rewards: processed_rewards }  
            sess.run(train_op, feed_dict=feed_dict)
```

one  
training  
iteration

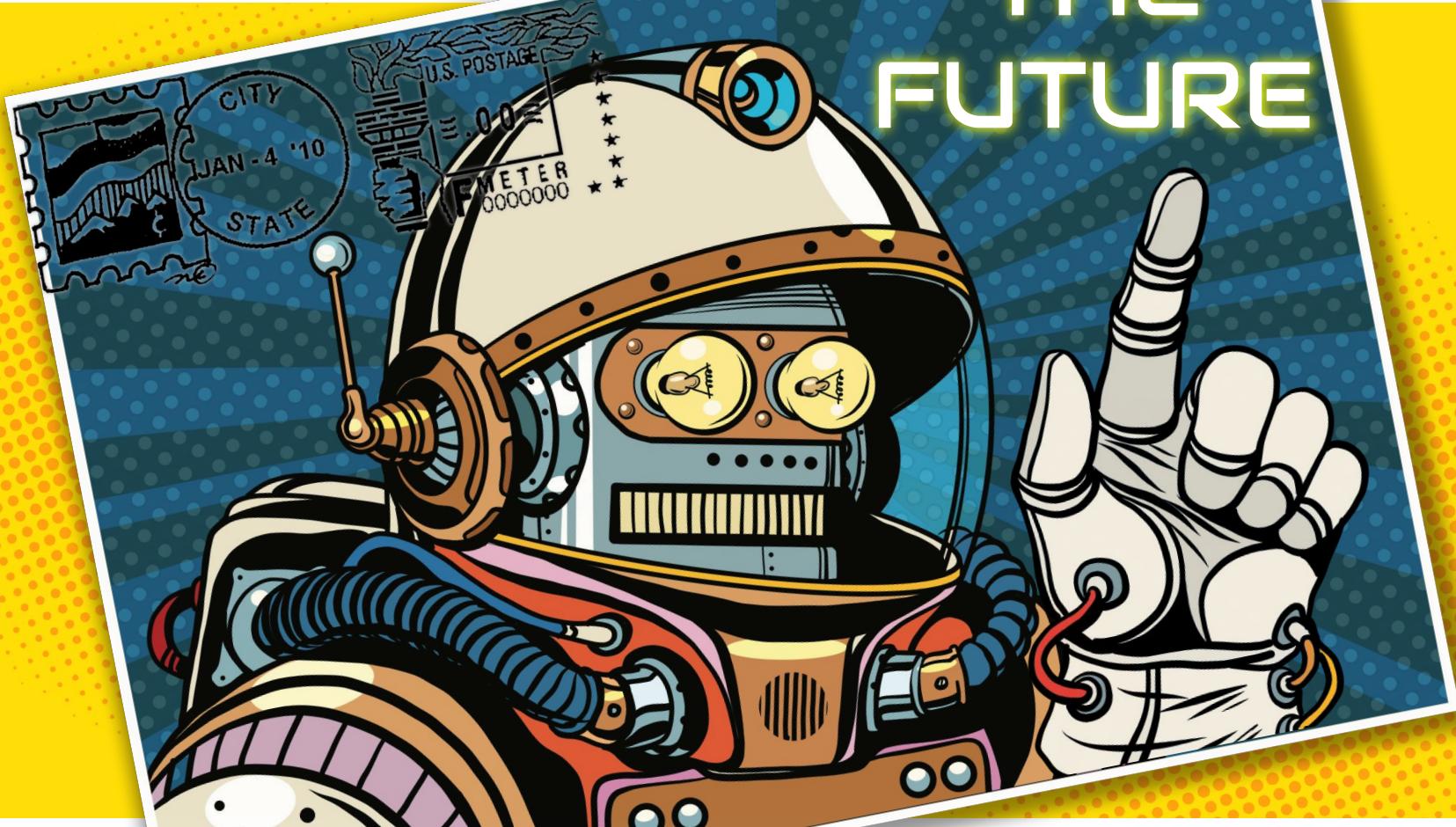
# Pong!

# Learned weights



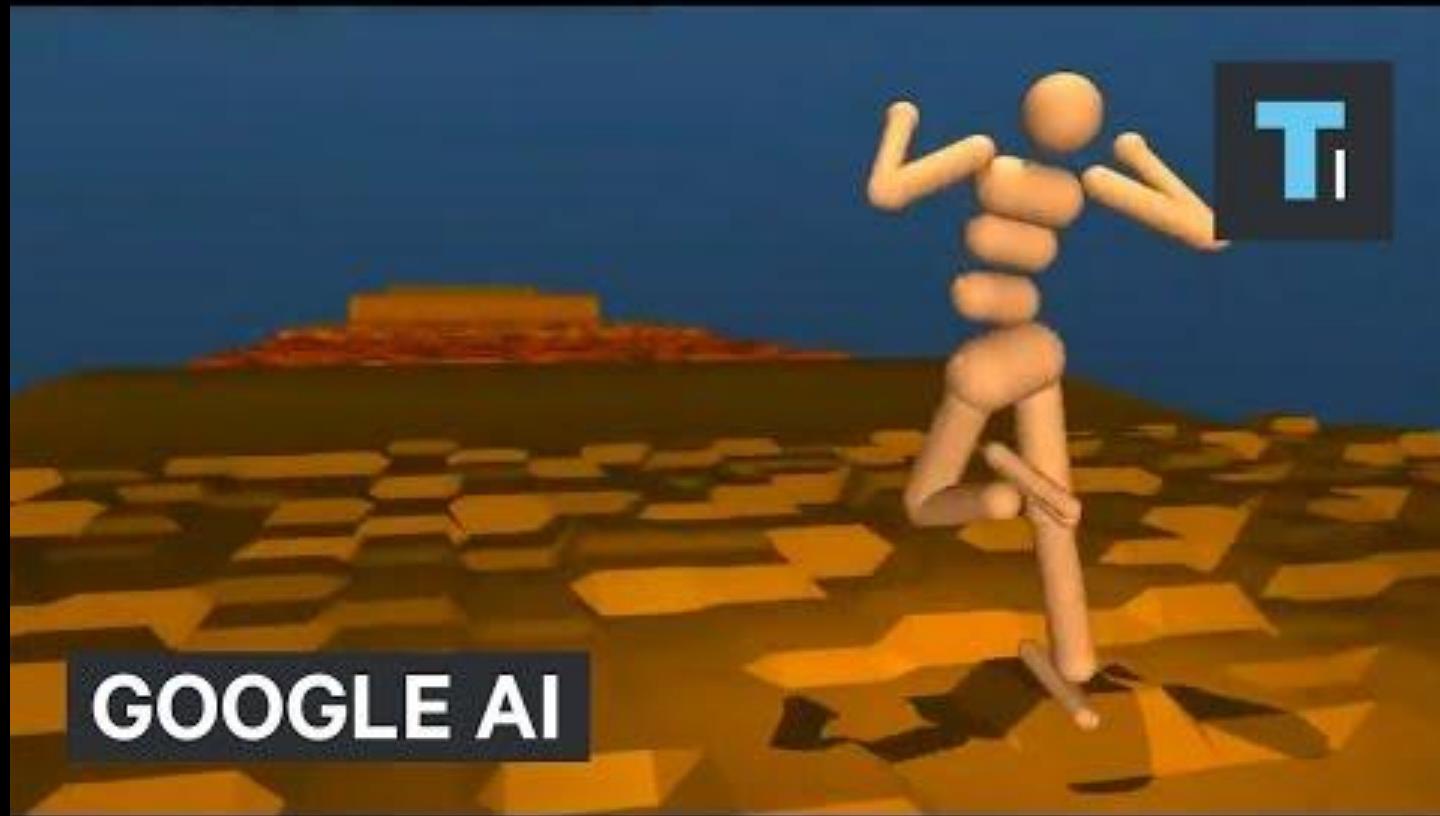
Postcard from...

# THE FUTURE





**Reinforcement Learning  
First trial...**





# Google DeepMind Challenge Match

8 - 15 March 2016



## AlphaGo vs Lee Sedol

Match 4 - Livestream

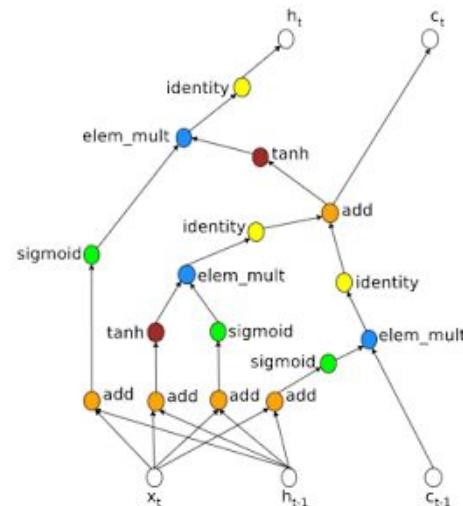
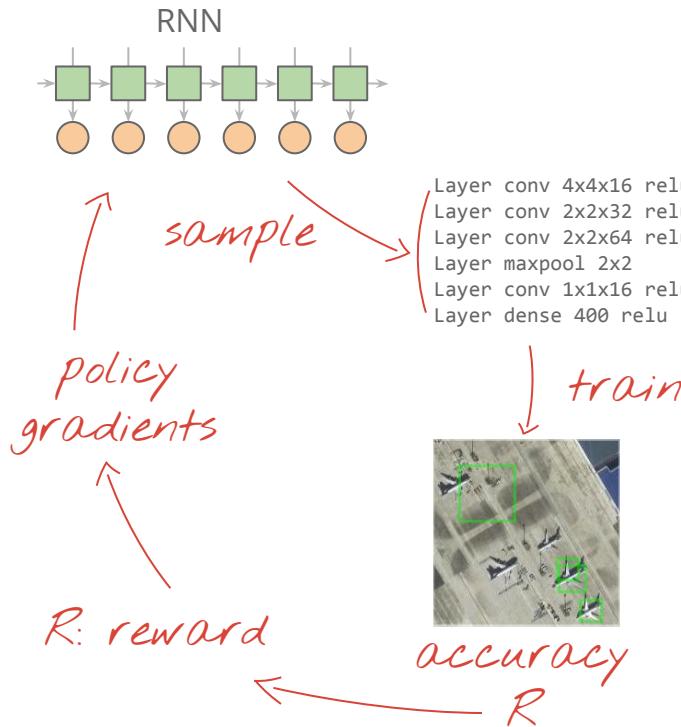
13th March 13:00 KST, 04:00 GMT

-1 day (12th March) 20:00 PT, 23:00 ET

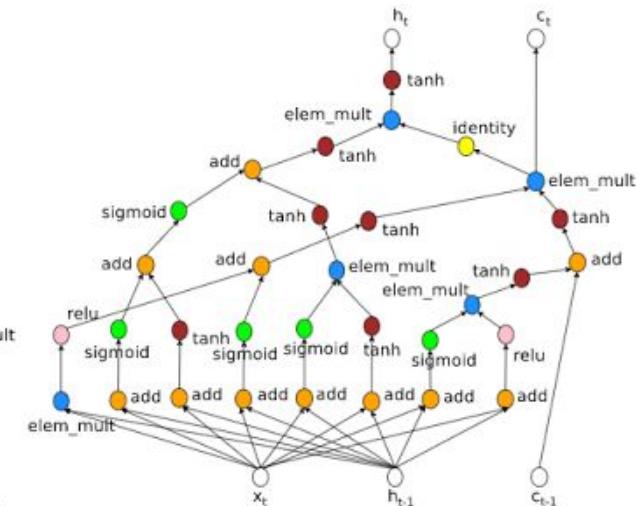
Pre-Match Commentary starting at 12:45 KST,  
03:45 GMT -1day (12th March) 19:45 PT, 22:45 ET

Live from the Four Seasons Hotel Seoul!

# Auto ML



LSTM cell



new RNN cell



arXiv:1611.01578v2, Barret Zoph, Quoc V. Le, May 2017

Replay on [youtu.be/aRKOJHRbXeo](https://youtu.be/aRKOJHRbXeo)

Have fun!



Martin Görner

Google Developer relations

[@martin\\_gorner](https://twitter.com/martin_gorner)



Yu-Han Liu

Google Developer relations

[yuhanliu@google.com](mailto:yuhanliu@google.com)



Neeraj Kashyap

Google Developer relations

[nkash@google.com](mailto:nkash@google.com)



TensorFlow  
[tensorflow.org](http://tensorflow.org)



Google Cloud Platform - [cloud.google.com](https://cloud.google.com)



## Cloud ML Engine

your TensorFlow models  
trained in Google's cloud,  
fast.



## Pre-trained models:



Cloud Vision API



Cloud Speech API



Natural Language API



Google Translate API



Video Intelligence API BETA



Cloud Jobs API PRIVATE BETA

This presentation:

[goo.gl/CB8xNH](https://goo.gl/CB8xNH)

Tensorflow w/o a PhD:

part 1: [goo.gl/pHeXe7](https://goo.gl/pHeXe7)

part 2: [goo.gl/UuN41S](https://goo.gl/UuN41S)

part 3: [goo.gl/VxQDmx](https://goo.gl/VxQDmx)

part 4: [goo.gl/BYT7au](https://goo.gl/BYT7au)

part 5: [goo.gl/CB8xNH](https://goo.gl/CB8xNH)

codelab: [goo.gl/mVZloU](https://goo.gl/mVZloU)



>TensorFlow and deep learning\_ without a PhD



>TensorFlow, deep learning and \\ recurrent neural networks without a PhD



>TensorFlow, deep learning and \\ modern convolutional neural nets without a PhD



>TensorFlow and \\ deep reinforcement learning without a PhD



The superpower: batch normalisation



>TensorFlow, deep learning and \\ modern RNN architectures without a PhD



Tensorflow and  
deep learning  
without a PhD



1. [youtu.be/u4alGiomYP4](https://youtu.be/u4alGiomYP4)
2. [youtu.be/fTUwdXUFfl8](https://youtu.be/fTUwdXUFfl8)
3. [youtu.be/vaL1l2BD\\_xY](https://youtu.be/vaL1l2BD_xY)
4. [youtu.be/aRKOJHRbXeo](https://youtu.be/aRKOJHRbXeo)
5. [youtu.be/vq2nnJ4g6N0?t=76m](https://youtu.be/vq2nnJ4g6N0?t=76m)
6. [youtu.be/pzOzmxCR37I](https://youtu.be/pzOzmxCR37I)

I ❤️ neurons