

The Goal

争上游 (ZhengShangYou, or “Conquer the River”) is a Chinese card game that is part of a larger family of games. In this version, each player is dealt ~18 cards, get rid of bad patterns, and win by getting rid of all cards. The goal is aimed to train an RL agent to have a good strategy against humans in a 2-player version.

3 Main Challenges

Battle Royal Exploration

Leon Lin (leonl@stan

me Round Robin to Battle

Competition Upstream”) is a strategy, part luck. Each of cards by matching all their cards first. I see an above 50% win rate in this game.

Enges

Battle Royale Algorithm
Init buffer with 100,000 random games
Init A with 24 agents, 12 with duplicate models
for e in NumEpochs do
$\epsilon = \frac{1}{e+2} * \sqrt{\frac{ A_0 }{ A }}$
Kill_Threshold = $0.5 * (1 - .75 * e^{\frac{ A }{ A_0 }})$
for i in range(500) do
sample = random 4096 points from buffer
for Agent ag in A do
Train Agent ag sample
for i in range($ A $) do

The PvP n
different r
other, and
training co
Two initial
the param
1) To det

1 for Card Game ZSY

iford.edu)

attle Royale

ature of the game allows for
models to be tested against each
gradually selected/reduced as
continues.

experiments were run to determine
eters for the final algorithm:

etermine the size of the buffer, 5

Discordance

The ConvNet Agents all had
different network structures
Using the models from the 1
round robin was run: all 12
action to find the average c
across and percentage of ti

Cosine Discordance



and Aggregation

similar performances despite
;. Did they learn the same policy?
Round Robin experiment, one final
models evaluated each state and
osine similarities for their Q values
mes they chose different actions



- *Stochastic + large discrete state* each dealt 18 cards, there are ~1
- *Partially Observed*: A player only
- *Train vs Test*: It cannot be trained data required, but it will be tested

I have previously attempted this but 40% win rate against a human player

The Data

Balancing the need to preserve the complexity of the game with the need to save memory, all 'hands' of cards are represented by stacks of one-hot encodings of how many there are of each

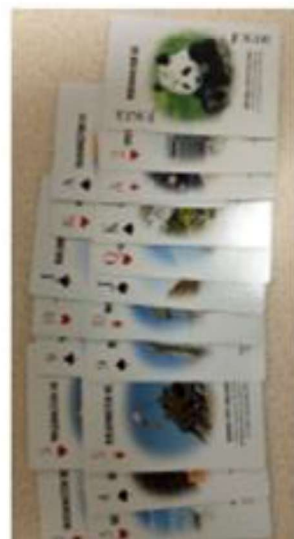


```
>>> h
array([[2,
```

1x15 array

```
>>> sz.y.de.
array([[0,
       [0,
       [2,
```

space: With 2 players
 .51 trillion initial states
 , see their own cards.
 d against humans for the
 l against them
 out only achieved about a
 ayer (me).



A hand of cards

```
1, 2, 0, 0, 0, 2, 2, 2, 1, 2, 2, 1, 2, 1, 0]], dtype=int8)
```

representing the counts of each value of card

```
HandToExpanded(h)
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, 1, 0, 0, 0, 1, 1, 1, 0, 1, 1, 2, 0, 0, 0],
1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0],
```

SIMULATE

for j in range(i, |A|) do
 Simulate 50000/ $\frac{|A|*|A|+1}{2}$ games between A_i and A_j
 Add games to buffer

TEST

for Agent ag in A do
 Update ag's "VS" winrate
 Test Agent ag for 200 games against the 3 static agents
 Update Agent ag's winrate against the 3 static agents

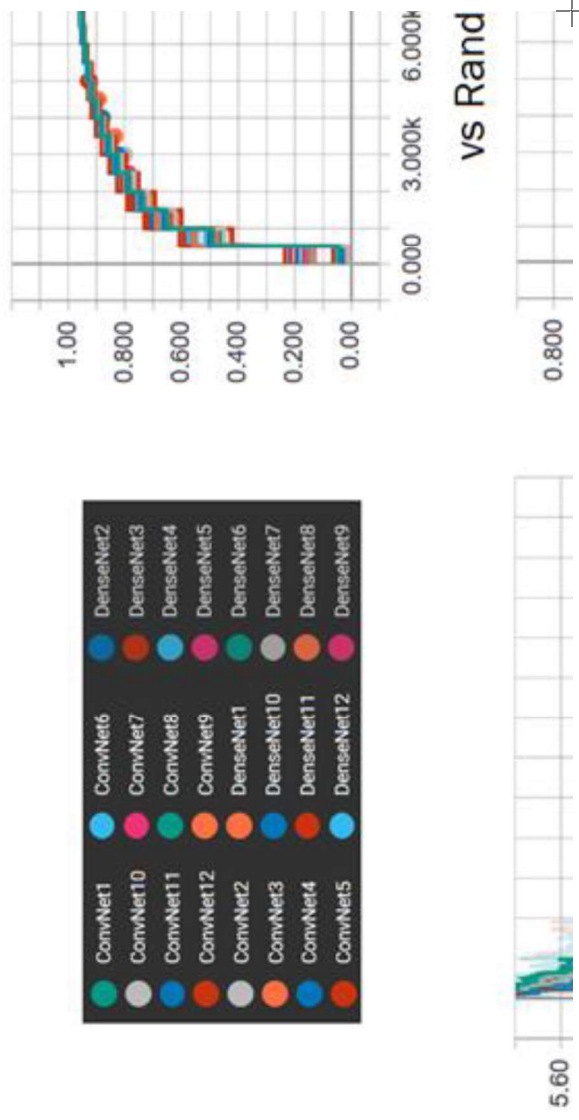
KILL

if ag's VS winrate < Kill_Threshold then
 Remove ag from A

agents we
 random ga
 to see whe
 set the bu
 start with

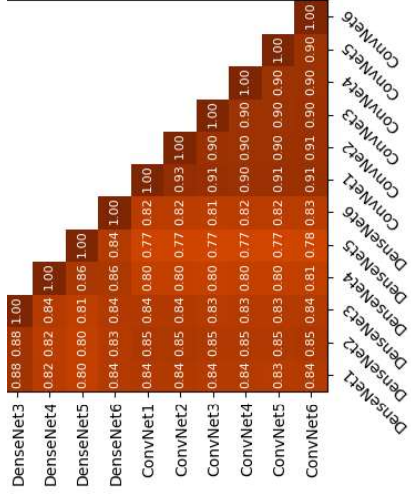
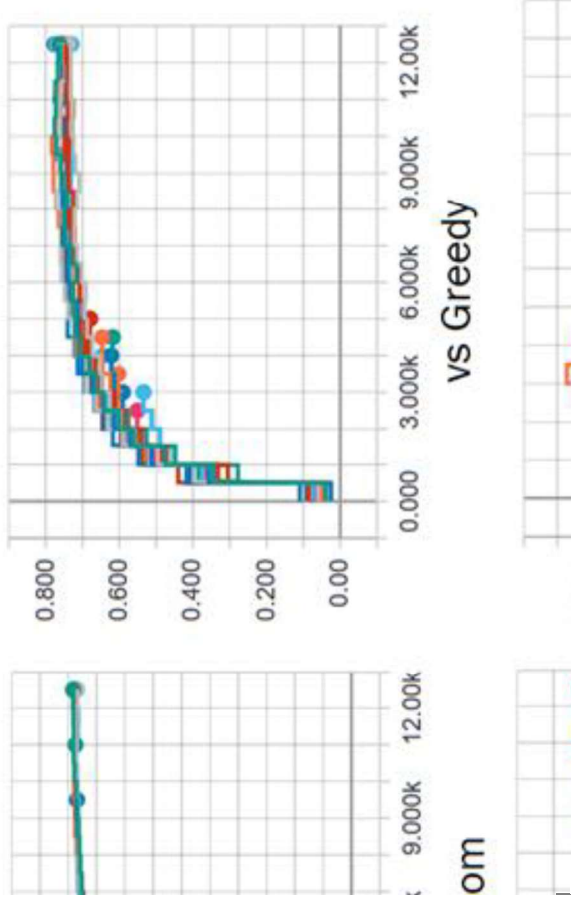
2) To figu
 much to e
 algorithm

From this, I decided the exploration an
 Then, the models from the second exp
 with 2 instances), and the Battle Royal



are trained on a large amount of games, and tested at various intervals in winrates stopped increasing. This offer size to max out at 250k and 100k

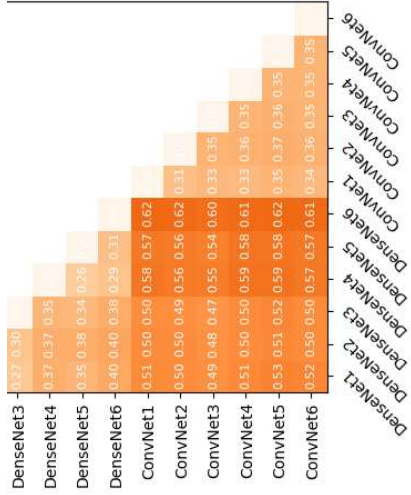
are out when to kill models and how explore, a shorter version of the with 12 agents and no removal. d kill threshold formulae. eriment were doubled (each one e was run.



Despite being trained on the quite different. This suggests might improve performance average, minimum, or maximum or 9 models. The aggregation agents did significantly better than any individual agent.

Hur

I build the game in Unity for PC and Mac and not these



e same data, their policies seemed
 ized that an aggregation model
 : 9 kinds were tested: taking the
 mum Q value from the top 3, 6, or

Agent vs	Random	Greedy	Old DQN
Min6Combo	97.8%	77.7%	75.4%
ConvNet2	96.6%	72.6%	62.8%
DenseNet9	96.2%	61.7%	24.4%

nan Tests

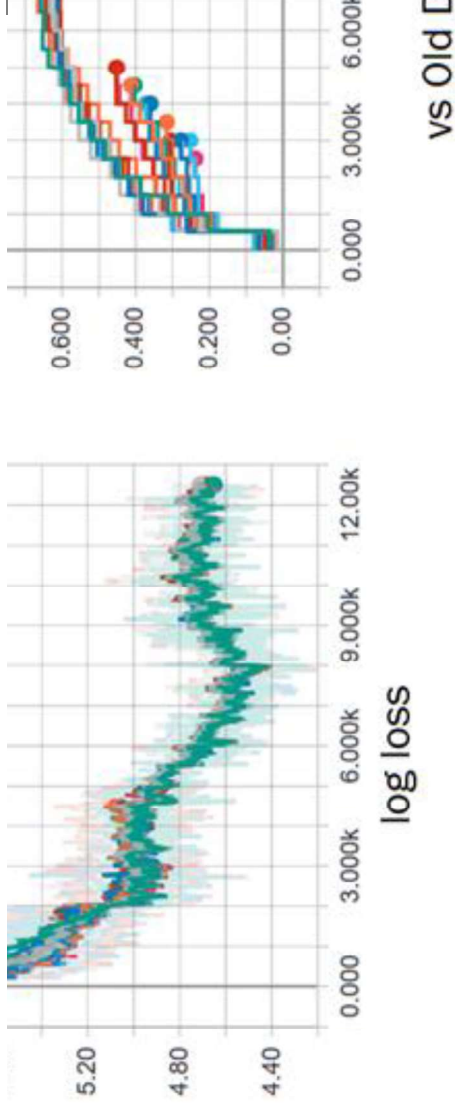
Player	Ratio	Winrate	σ
--------	-------	---------	----------


```
[0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1]], dtype=uint8)
```

with each column being a one-hot encoding of
ther there are 0, 1, 2, 3, or 4 of a card

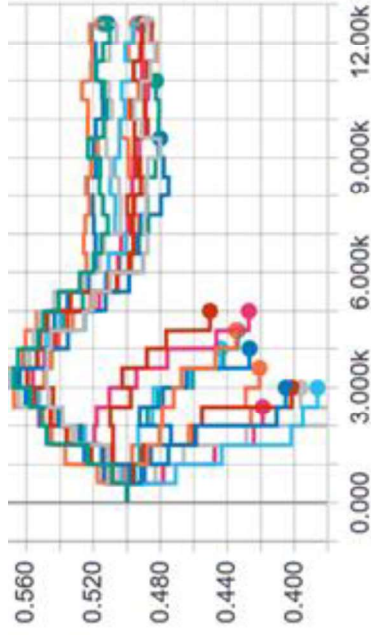
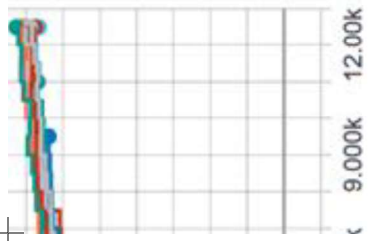
, and cards the
re concatenated into a
tion pairs.
replay buffer inspired
rent models are
ther to contribute to a
r data efficiency,
r search.

e Mac build of game



Each state-action pair is assigned a value
losing games. DenseNet Agents flatten the
4 fully connected layers while ConvNet Agents
of 3x3, 1x3, or 1x1 valid convolutions before
Activations are all ReLU, LeakyReLU, or SiLU
the configs.csv files on the Github repository.
Testing was done against a Random agent
from the previous project.

ConvNet agents performed significantly better
end 10 of the 24 initial agents remained,



QN

vs Each Other

of 1 for winning games and 0 for
 ie 5x60 array and pass it through 3-
 agents first pass it through 2-4 layers
 ore 2 fully connected layers.
 gmoid (for the output layer). See
 or full definitions of each agent.
 it, a Greedy agent, and the agent

etter than DenseNet agents. By the
 all of which were CNNs.

test results back from
 human players.

The data is too small to
 conclusively say that human
 level performance was
 achieved, but the initial
 results are promising: on
 average, no human beat the
 agent.

N

With the game built for hun
 real human player data and
 states where the agent cho
 models such as RNNs can
 normalization may be helpf

1	62:66	48.4%	4.4%	
2	27:32	45.8%	6.5%	
3	16:24	40.0%	7.5%	
4	12:20	37.5%	8.6%	
5	20:36	35.7%	6.4%	

next Steps

nan players, it is possible to collect
 d to better visualize different game
 se badly. Additionally, different
 be attempted and gradient
 ul for longer term training.