

DouRN: Improving DouZero by Residual Neural Networks

Yiquan Chen

School of AI and Advanced Computing
Xi'an Jiaotong-Liverpool University
Suzhou, China

Yiquan.Chen19@student.xjtlu.edu.cn

Yingchao Lyu

School of AI and Advanced Computing
Xi'an Jiaotong-Liverpool University
Suzhou, China

Yingchao.Lyu@student.xjtlu.edu.cn

Di Zhang*

School of AI and Advanced Computing
Xi'an Jiaotong-Liverpool University
Suzhou, China

Di.Zhang@student.xjtlu.edu.cn

Abstract—Deep reinforcement learning has made significant progress in games with imperfect information, but its performance in the card game Doudizhu (Chinese Poker/Fight the Landlord) remains unsatisfactory. Doudizhu is different from conventional games as it involves three players and combines elements of cooperation and confrontation, resulting in a large state and action space. In 2021, a Doudizhu program called DouZero [8] surpassed previous models without prior knowledge by utilizing traditional Monte Carlo methods and multilayer perceptrons. Building on this work, our study incorporates residual networks into the model, explores different architectural designs, and conducts multi-role testing. Our findings demonstrate that this model significantly improves the winning rate within the same training time. Additionally, we introduce a call scoring system to assist the agent in deciding whether to become a landlord. With these enhancements, our model consistently outperforms the existing version of DouZero and even experienced human players.¹

Index Terms—DouDizhu, Reinforcement Learning, Monte Carlo Methods, Residual Neural Networks

I. INTRODUCTION

The objective of this study is to develop a deep reinforcement learning-based model for DouDiZhu, a card game that presents unique challenges when applying reinforcement learning techniques from other games. One key challenge is that DouDiZhu involves three players, namely a landlord and two peasants, where the peasants collaborate against the landlord, resulting in a game that combines both cooperation and competition. Consequently, the widely used counterfactual regret minimization algorithm employed in poker [1] cannot be directly applied to this complex three-player games.

Furthermore, the set of permissible card combinations changes dynamically as the game progresses. In contrast to Texas Hold'em [1], where Deep Reinforcement Learning (DRL) has been successfully applied, DouDiZhu encompasses a large action space that is difficult to explore. Consequently, traditional reinforcement learning algorithms suffer from high search costs and prove ineffective. Both the Deep Q-Learning and Asynchronous Advantage Actor-Critic (A3C) algorithms [4], well-established in the field of DRL, are unsuitable for DouDiZhu training due to these reasons. Deep Q-Learning is prone to overestimation issues in large action spaces,

significantly impacting training outcomes, while A3C struggles to fully exploit the features of the action space, limiting its ability to generalize to unknown actions [2]. According to [7], both models exhibit remarkably low winning rates against simple rule-based models.

Although deep learning is still necessary, training DouDiZhu using MLPs presents challenges such as tuning difficulties, extensive training time, and slow learning speed, as observed in DouZero [8]. To address these issues, we adopted the residual neural network (ResNet) to enhance the original multilayer perceptrons (MLPs) in DouZero, which offers potential solutions to the degradation problem of deep neural networks and facilitates faster convergence [9].

Here we introduced a simple multi-layer perceptron based on a call-score network to enable the AI to determine whether it can play as a landlord and obtain three additional hidden cards. When opportunity comes, becoming a landlord not only increases the probability of winning but also provides a double score bonus. Accurately assessing the score at the beginning of each game by analyzing one's own cards and the scores given by opponents is crucial.

By combining these two enhancements, the newly trained model, DouRN, outperforms the original DouZero model and exhibits superior capability against experienced DouDiZhu players.

II. RELATED WORK

There exists a substantial body of literature pertaining to the application of DRL for DouDiZhu. Notably, a recent study [7] introduced a Combinational Q-Learning algorithm that combines a reduced action space and maximum pooling technique, resulting in promising outcomes for DouDiZhu gameplay. Furthermore, paper [3] proposed a model named DeltaDou, which integrates Monte Carlo tree search with DRL. However, this model exhibits limitations in terms of computational cost and reliance on human a priori knowledge, thereby constraining its performance and generalizability. In contrast, the DouZero model, which employs pure DRL, has outperformed all previous AI models in DouDiZhu gameplay [8]. Building upon this, DouZero+ [10] have enhanced the DouZero model through opponent modeling and coach-guided learning.

¹The source code is available at https://github.com/Yingchao/Douzero_Resnet.git.

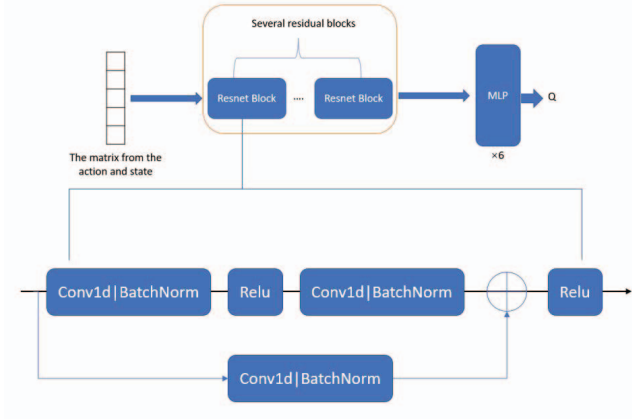


Fig. 1: The architecture of a residual block.

ResNet is a deep learning architecture that revolutionized image classification tasks. It addresses the problem of vanishing gradients in deep networks by introducing skip connections, also known as residual connections. These connections allow the network to directly propagate information from earlier layers to later ones, enabling faster and more accurate learning. This architecture has significantly improved the training of deep neural networks, leading to state-of-the-art performance in various computer vision tasks [9]. However, the application of ResNet in DRL has only recently gained attention. For instance, a DRL model named FullyConv, designed for StarCraft II, incorporates ResNet into its architecture [5]. Furthermore, paper [6] successfully utilized ResNet for deep learning in Chinese Mahjong, a complex incomplete information game similar to DouDiZhu. Comparative analysis against conventional approaches revealed that ResNet can achieve better winning rate and score, indicating its superior performance.

III. METHODOLOGY

A. Residual Networks

The key to designing a residual neural network is to build a residual block, which usually consists of multiple convolutional layers and batch normalization layers [9]. In the residual block, the input of each convolutional layer is the output of the previous layer plus a skip connection (i.e., residual connection). The fundamental structure of a residual block we used here is illustrated in Fig. 1. The original neural network structure of DouZero is a 6-layer MLP. To optimize the network, several residual blocks can be appended. Generally, the use of residual blocks may enhance algorithmic training effectiveness, but deepening neural networks may also increase the consumption of computing resources. Therefore, we need to carefully design the network architecture.

We have designed two architectures. In the first architecture (Fig. 2), we directly stack residual blocks on top of the MLP in the original DouZero, while keeping the training method unchanged. This allows us to observe the impact of the residual blocks. After achieving positive effects with fewer stacked

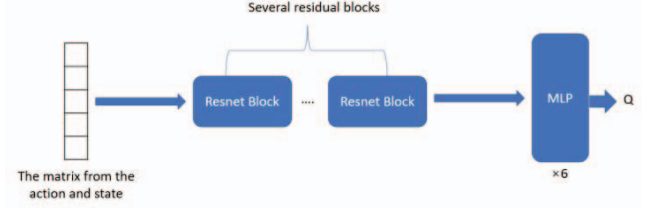


Fig. 2: Insert several residual blocks before the MLP.

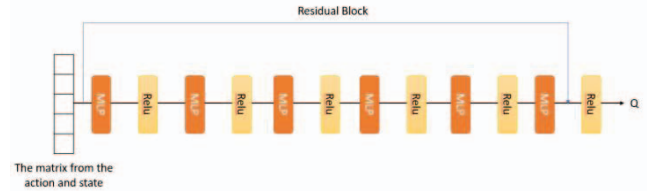


Fig. 3: Replace the original 6-layer MLP into residual blocks.

layers, we can continuously improve the model performance by increasing the number of residual blocks. In the second architecture (Fig. 3), we replace the entire MLP in the original DouZero with residual blocks while maintaining the overall number of layers. We hope that this approach can minimize the increase in model training time.

B. Bidding System

After the game starts, each player has 17 cards in his hand, and the remaining 3 cards are hidden cards. The first player to be called Landlord is determined randomly, usually automatically selected by the system. Other players take turns calling the landlord, and they can choose to call the landlord or not. When bidding on a landowner, he must offer a higher multiplier (1, 2 or 3) than the previous player. The hidden card will be taken by the highest bidder and added to the hand. The landowner will play his cards first, and the other two farmers will become his opponents.

Obviously, a decision-making system needs to be added here, which includes analysis of the current hand and the opponents' bidding multiplier. The potential benefits must outweigh the

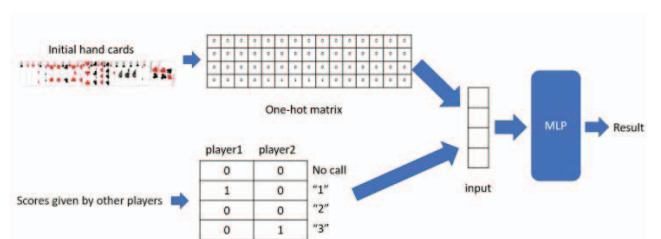


Fig. 4: The input data utilized in the bidding system encompasses the player's hand as well as the scores assigned by the two opposing players. Notably, the scores given by the other players are transformed into a one-hot matrix representation for further analysis and processing.

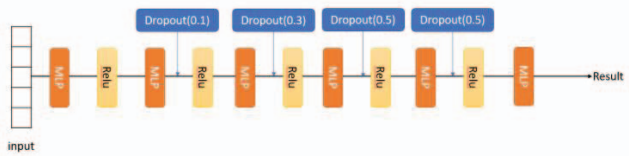


Fig. 5: In the architecture of the bidding system, some neurons will be randomly discarded in dropout without training, and the probability is indicated within parentheses.

risks. The player's hand is transformed into a 4×15 matrix representation. Furthermore, the score provided by the player is converted into a 4×2 one-hot matrix representation. For instance, as illustrated in Fig. 4, player 1 assigns a score of 1, while player 2 assigns a score of 3, resulting in the generation of a matrix representation.

The MLP architecture of the bidding system is depicted in Fig. 5. It comprises 6 layers of neural networks, incorporating dropouts to mitigate overfitting and prevent excessive activation of all neurons. The network ultimately produces an output value, which determines the action assigned by the trained model.

IV. EXPERIMENTS

We first evaluate the performance of the two residual network architectures along in Sec. III-A. To evaluate the performance of the enhanced model², we adopted a model validation approach similar to that employed by DouZero. This involved playing landlords and peasants from different models against each other. Specifically, one model represented the improved DouRN, while the other represented the original DouZero, with each assuming the role of landlord and peasant respectively. Then, their roles were swapped, enabling them to compete against each other using the same deck of cards. The evaluation criteria were based on the winning rate, calculated as the number of games won by the model divided by the total number of games played.

Furthermore, to assess the bidding system in Sec. III-B, selected real players were invited to compete against the models, and their performance was evaluated based on their winning percentage. Except for the newly added residual blocks, DouRN is consistent with DouZero in all hyperparameters. This facilitates our comparison.

A. Evaluation on Residual Networks

2 Extra Layers: In order to assess the impact on model performance, the model was initially trained with 2 residual blocks for a total of 3.0×10^5 steps. As indicated in Tab. IV-A, it can be observed that although the DouRN model does not outperform DouZero in the early stages of training as a landlord, it slightly surpasses DouZero in terms of winning rate at the end of training, both as a landlord and a peasant. Hence, a preliminary conclusion can be drawn that the inclusion of residual blocks can be beneficial for enhancing the model.

²The experiments were conducted on the Sugon platform, utilizing a server equipped with an 8-core Hygon C86 7185 CPU @ 2.40 GHz, 64 GB of RAM, and two Tesla T4 graphics cards.

TABLE I: Comparison of a DouRN model trained with a 2-layer residual block and a DouZero model with the different roles and training steps. L. stands for landowner and P. stands for peasant.

Step($\times 10^5$) /2 Blocks	DouRN(L.) vs. DouZero(P.)	DouRN(P.) vs. DouZero(L.)	Training Time (h)
1.0	49.02%	51.32%	56
2.0	51.92%	52.10%	112
3.0	52.02%	53.77%	168

TABLE II: Winning rate of match between different residual block models and DouZero models.

Step($\times 10^5$) /4 Blocks	DouRN(L.) vs. DouZero(P.)	DouRN(P.) vs. DouZero(L.)	Training Time (h)
1.0	0.5221	0.5322	68
2.0	0.5403	0.5512	136
3.0	0.5508	0.5551	204
4.0	0.5549	0.5578	272
5.0	0.5562	0.5592	340
6.0	0.5580	0.5640	408

Step($\times 10^5$) /6 Blocks	DouRN(L.) vs. DouZero(P.)	DouRN(P.) vs. DouZero(L.)	Training Time (h)
1.0	0.5397	0.5410	77
2.0	0.5611	0.5562	154
3.0	0.5641	0.5609	231
4.0	0.5683	0.5702	308
5.0	0.5701	0.5704	385
6.0	0.5704	0.5703	462

4/6 Extra Layers: At this stage, we extended the number of residual blocks to 4 and 6 and increased the training steps to 6.0×10^5 . This ensured a sufficient training of the model. The results presented in Tab. IV-A demonstrate a consistent increase in winning rate with the addition of residuals. Specifically, the model with 4 residuals exhibited a significantly higher win rate compared to the model with 2 residuals. However, the model with 6 residuals approached convergence at 6.0×10^5 steps, with a win rate similar to that of the model with 4 residuals. Furthermore, the training cost also increased considerably, with the model with 6 residuals taking approximately twice as long as the model with 2 residuals. Thus, it can be predicted that further addition of residuals would enhance the model's performance but at an unacceptable increase in computational effort.

Simple Replacement: The architecture depicted in Fig. 3 is tested under same procedures as before. Table 4 reveals that the new model exhibits similar winning rates to the DouZero model, and even performs worse at lower step counts. In addition, there is no notable distinction between the two models in terms of training time. Therefore, the effect of simple replacement is not satisfactory.

Comparative Analysis: Networks with additional layers are superior to networks without additional layers in this game,

TABLE III: Comparison of win rates between the MLP model with residual blocks and the DouZero model.

Step($\times 10^5$) /Simple	DouRN(L.) vs. DouZero(P.)	DouRN(P.) vs. DouZero(L.)	Training Time (h)
1.0	47.98%	49.08%	61
2.0	49.11%	50.96%	122
3.0	49.60%	50.51%	183
4.0	50.40%	51.03%	244
5.0	50.50%	51.23%	305
6.0	50.32%	50.76%	366

TABLE IV: Winning percentage for human players against DouRN, calculated by the sum of games won by human, divided by the total number of games (50).

	With Bidding System	Without Bidding System
Ordinary player1	38%	40%
Ordinary player2	42%	42%
Ordinary player3	32%	42%
Experienced player1	46%	52%

which is probably due to the following reasons:

ResNet adopts more residual blocks and has a greater depth, which means it can have more levels to learn complex features and explore more strategies and states. This enables ResNet to better capture subtle variations and complex relationships in the experience of play, thereby improving network performance.

ResNet has many hyperparameters that need to be adjusted, such as network depth and the number of residual blocks. If the MLP in DouZero is simply replaced by a residual network without proper optimization, it may not achieve the desired effect.

Limitation: Firstly, limited computational resources hindered sufficient training, leading to suboptimal neural network architecture and model performance. Secondly, the model exhibited unwise decisions in certain testing scenarios. Lastly, the inclusion of residual networks significantly reduced training speed.

B. Evaluation on Bidding System

For the bidding system, we used a different assessment method. Instead of playing against other models as in Sec. IV-A, we chose 4 human players to play against a model containing the bidding system and calculated the win rate. These 4 players include 3 regular DouDiZhu players in their twenties and 1 top DouDiZhu player in his fifties. They will play against the DouRN network with and without a bidding system, completing a total of 100 games and keeping track of the winning percentage.

When playing against models that lack a bidding system, the distribution of cards and characters to each player is random. Thus, it is possible for a player to receive a disadvantageous hand yet still become the landlord, whatever the player is human or machine. This random allocation ensures fairness in

the game, preventing any bias towards the landlord or peasant in terms of winning rates.

As indicated in Tab. IV-B, the introduction of a bidding system has more or less decreased the winning rate of human players. Even experienced players have their winning rates lower than 50% due to the introduction of the bidding system. And overall, the model can beat all human participants.

Based on the aforementioned data, it can be inferred that the inclusion of a bidding system can yield superior results when playing against human opponents. Furthermore, it is plausible to anticipate improved performance against other models as well.

V. CONCLUSION AND FUTURE WORK

In this study, we enhanced the DouZero model by introducing the residual neural networks and incorporating a bidding system. Residual neural networks have been shown to be effective on this problem, and the bidding system can be valuable to any other DouDiZhu model.

Despite the advantages of DouRN, there is still room for further optimization. Firstly, additional computational resources are needed to validate and optimize larger neural network. Secondly, the stability of model can possibly be enhanced by incorporating Monte Carlo search techniques. Lastly, we intend to incorporate off-policy learning into the model training process to improve overall efficiency.

REFERENCES

- [1] Noam Brown and Tuomas Sandholm. Superhuman ai for multiplayer poker. *Science*, 365(6456):885–890, 2019.
- [2] Gabriel Dulac-Arnold, Richard Evans, Hado van Hasselt, Peter Sunehag, Timothy Lillicrap, Jonathan Hunt, Timothy Mann, Theophane Weber, Thomas Degris, and Ben Coppin. Deep reinforcement learning in large discrete action spaces. *arXiv preprint arXiv:1512.07679*, 2015.
- [3] Qiqi Jiang, Kuangzheng Li, Boyao Du, Hao Chen, and Hai Fang. Deltadou: Expert-level doudizhu ai through self-play. In *IJCAI*, pages 1265–1271, 2019.
- [4] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.
- [5] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- [6] Mingyan Wang, Tianwei Yan, Mingyuan Luo, and Wei Huang. A novel deep residual network-based incomplete information competition strategy for four-players mahjong games. *Multimedia Tools and Applications*, 78:23443–23467, 2019.
- [7] Yang You, Liangwei Li, Baisong Guo, Weiming Wang, and Cewu Lu. Combinatorial q-learning for dou di zhu. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, volume 16, pages 301–307, 2020.
- [8] Daochen Zha, Jingru Xie, Wenye Ma, Sheng Zhang, Xiangru Lian, Xia Hu, and Ji Liu. Douzero: Mastering doudizhu with self-play deep reinforcement learning. In *international conference on machine learning*, pages 12333–12344. PMLR, 2021.
- [9] Shangdong Zhang, Wendelin Boehmer, and Shimon Whiteson. Deep residual reinforcement learning. *arXiv preprint arXiv:1905.01072*, 2019.
- [10] Youpeng Zhao, Jian Zhao, Xunhan Hu, Wengang Zhou, and Houqiang Li. Douzero+: Improving doudizhu ai by opponent modeling and coach-guided learning. In *2022 IEEE Conference on Games (CoG)*, pages 127–134. IEEE, 2022.