# Car Parking Simulation using Machine Learning

Anshul Tickoo[1*], Shivansh Sukhija[2], Abdus Saboor Gilani[3], Ch. V. V. S. Satyanaryana[4], Sahil Garg[5]

[1,2,3,4,5]*Department of Computer Science and Engineering, Amity School of Engineering and Technology, Noida, India*

***Abstract***: **This research paper explores the field of artificial intelligence and its subfields, including reinforcement learning, which is used in the Unity ML Agents project to create autonomous parking models using the Proximity Policy Optimization (PPO) algorithm. The models use LIDAR sensors to navigate through the environment and avoid obstacles while searching for an optimal parking spot. The findings exhibit the effectiveness of the PPO algorithm and the importance of road safety. The development of autonomous parking models provides a steppingstone towards the creation of more advanced models that can be used in real-world applications such as self-driving cars. This project covers topics such as machine learning, reinforcement learning, LIDAR sensors, PPO algorithm, and road safety, contributing to future research and development in artificial intelligence.**

***Keywords***: **car simulation, car parking, car parking model.**

## 1. Introduction

Logistics, severe settings, military applications, unmanned aerial vehicles, food service, and many other spheres of human activity have all benefited greatly from the use of intelligent robots. Logistics and transportation, harsh environments, autonomous vehicles, the food industry, the military, and many more fields make extensive use of them. They have found widespread application in both industrial production and domestic settings. Cleaning robots, inspection robots, toy robots, autonomous vehicles, inorganic automation, etc., have all contributed to higher productivity and more pleasant working conditions. Autonomous and rapid learning is a crucial feature for robots to take on complicated and varied jobs. A certified intelligent robot must be capable of both precise execution and rapid learning. As a result, it is crucial to enhance the cognitive capacity of smart robots so that they can rapidly adapt to new circumstances.

### A. Machine Learning

Artificial intelligence may be achieved via the process of machine learning. In 1959, Arthur Samuel, one of the inventors of machine learning, argued that computers might learn new tasks without being given specific instructions. Machine learning is a branch of AI that places more emphasis on 'learning' than on pre-programmed algorithms. Without explicit instructions written by a person, machines may utilize complicated algorithms to examine enormous datasets, find patterns, and make predictions. Machine learning allows a system to correct its own pattern recognition errors.

It is the goal of machine learning to automatically identify the parameters (variable values) of a system given data (called training data or learning data). Based on the nature of the feedback provided during training, Zhou Zhihua divides machine learning into three distinct categories: supervised learning, unsupervised learning, and reinforcement learning. In the realm of machine learning, we find techniques like neural networks, logistic regression, decision trees, Support random forests, Vector Machine (SVM), and many more.
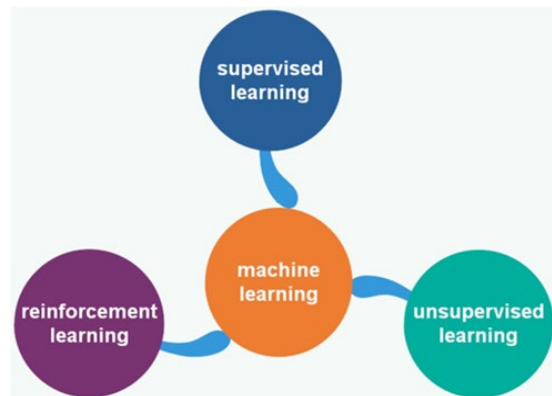


Fig. 1. The components of machine learning

There are two phases to the supervised learning process: learning and inferencing. When used to large datasets, unsupervised learning (also known as "clustering") may help find patterns and outliers while also facilitating peacekeeping and recommendation systems. Deep learning is a method in machine learning that is focused on learning from the representation of data, and it represents the second wave of machine learning after shallow machine learning. When working with large datasets, deep learning shines, whereas classic machine learning approaches shine when working with smaller datasets. Deep learning has broadened the possibilities of artificial intelligence and made possible a wide range of machine learning applications. Many machine-assisted tasks, like autonomous vehicles and preventative medicine, are made feasible by advances in deep learning.

### B. Reinforcement Learning

Using just sample training of the decision process, reinforcement learning (RL) aims to resolve the Markov decision process (MDP) problem. According to the state transfer function, the subsequent state is a result of the present

state and the action. Relationship between current state and desired reward delivery time is expressed by the reward function. The model is Markovian, which implies that the next state relies only on the current state and the action taken (first-order Markovian) and not on any previous states.
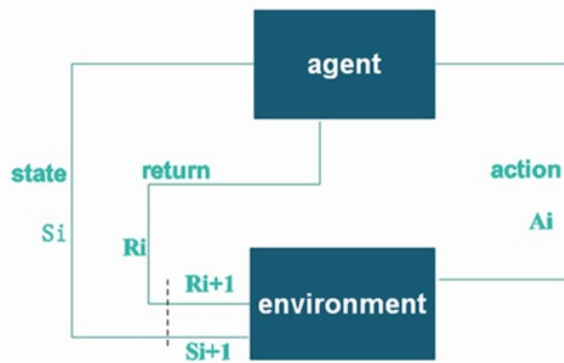


Fig. 2. Markov process model



Fig. 3. Markov decision process

An exploration-based machine learning method not only analyzes current data but also collects new data via exploration, then uses that data to repeatedly update a previous model.

The theory of reinforcement learning (RL) is based on the idea that rules that produce desirable outcomes should be rewarded, while rules that produce undesirable outcomes should be punished. Rather of labeling sample outcomes before educating the learner, as is done in traditional supervised learning, in reinforcement learning the learner attempts to achieve the consequences of each conduct, which in turn gives response on the training itself. Intelligences, states, environments, actions, strategies, rewards, and objectives are the essential components of reinforcement learning. In rl, however, rather than relying on preexisting model knowledge, the agent learns the optimal course of action via trial and error as it interacts with its environment. What we mean by "learning" is that the agent takes the information it gathers from interactions and applies it to its mental model of the world. Fee-based strategies and policy-based strategies are two types of rl algorithms that differ depending on their solution goal, which may be either optimal coverage or optimal value. For the sake of analysis, a sg may be seen as a sequence of regular-shape

video games. Consider the genuine intersection scenario (refer to parentheses). A normal-shape sport matrix may be used to represent a snapshot of the sg at time t (stage game). The rows represent Agent 1's action set a 1, while the columns represent Agent 2's action set a. Each cell in the matrix represents a reward for a specific combination of joint actions. If both sellers are just concerned with maximizing their own possible profit (the solution notion in a single-agent rl problem) and choose the action to hurry, then they will collide with each other. Despite the fact that the possible praise became the greatest for each agent when sprinting, this condition is risky and, as a result, is not valued as highly as others. To solve an SG and ethically maximize the total praise, each agent must make calculated decisions that take into account the actions of the others. Unfortunately, in contrast to mdps, which can be solved using linear-programming formulas in polynomial time, solving sgs often entails using newton's method for addressing nonlinear situations. But there are unique instances of participant widespread.



Fig. 4. Reinforcement learning methods

The machine learning subfield known as reinforcement learning (RL) focuses on how to take action in response to a given context in order to maximize rewards. The concept originates from the behaviorist school of psychology, which postulates that organisms learn to anticipate environmental cues in order to optimize the rewards they get. The primary applications of reinforcement learning are the development of control rules for dynamic systems and the creation of competitive game tactics. Even if reinforcement learning does not use labeled data for training, it does not indicate that there is no supervisory information there. When the intended outcome is reached, the system sends a signal known as a reward and continues running the reinforcement learning algorithm. For instance, the incentive in robot walking control is the total distance covered. The program's success or failure at the Go game is the prize. When one loses, one is given a negative number, which is called a penalty.

The field of reinforcement learning (RL) has found widespread applications in many different fields, including but not limited to: robotics, manufacturing, modeling, optimization, scheduling, and even video games. In contrast to RL, which typically employs an unlabeled training sample, DRL relies on a labeled one. Supervised learning is more like the "five

senses," while reinforcement learning is more like the "brain," with deep learning resolving more perceptual obstacles and reinforcement learning resolving mostly decision-making issues.

These games may only be played in an unorganized fashion. Although this isn't based on sound theory, in practice it makes dealing with a tiny sg far less complicated. Agents have the same interest, denoted by r here, in achieving the learning goal. The fact that each marketer is given its own consideration means that each agent may choose the course of action that best serves its own interests. This allowed for the precise application of single-agent rl algorithms, and eventually led to the development of a decentralized system. This class includes a wide variety of sgs. – group video games/absolutely cooperative games/multi-agent mdp (mmdp): agents are assumed to be homogeneous and interchangeable, so importantly, they share the identical praise , $r = r1 = r2 = \cdots = rn$ – team-common praise video games/networked multi-agent mdp (m-mdp): agents can have specific reward capabilities, but they percentage the identical goal,

$r = 1.n\ p\ n\ i = 1\ ri$

Stochastic capacity games: agents could have different praise features, but their mutual pursuits are described through a shared ability characteristic ion.

### C. Meta Reinforcement Learning

While the majority of the literature on reinforcement learning assumes a static environment, in practice, environments are often in flux, necessitating that reinforcement learning systems continually learn to adapt. Deep learning's growth is stifled by its dependency on massive data, thus researchers are shifting their focus to meta-reinforcement learning (meta-RL).

In their seminal paper, Schmidhuber et al. introduced the concept now known as meta-reinforcement learning (meta-RL). that using neural networks to rapidly solve reinforcement learning issues, a technique that jX Wang et al. term "deep meta-reinforcement learning." Meta-RL seeks to learn a sequence of initial parametrized control strategies (or policies) that, following the updating of parameters through a single gradient step calculated with a minimal number of samples, achieve maximum performance on a new RL problem. Maruan et al. used a MAML method in conjunction with a cyclic neural network to enable a reinforcement learning model to change its policy to a dynamic environment. Meta-reinforcement learning (meta-RL) methods, on the other hand, may solve the issue of low speed by drawing on data from completed tasks that are similar to the one at hand.

There are primarily two kinds of meta-reinforcement learning strategies (meta-learning algorithms for meta-reinforcement learning, such as MAML): model-based and optimization-based. The former is the "hot" brain-like approach of late, using RNN to establish correspondence with the prefrontal cortex.

### D. Application of Meta-Enhanced Learning

After successfully modeling the brain's navigation system with deep learning, the DeepMind team investigated the part that dopamine plays in the learning process by using a meta-reinforcement learning framework. Thus, they provided explanations for a number of neuroscientific and psychological discoveries that serve to further demonstrate the close relationship between meta-reinforcement training and human intelligence. A gradient-based meta-RL technique was presented by Feiran Zhao et al. for controlling ATO systems' speeds.

### E. Meta-Enhanced Learning and Research Development

The area of study dedicated to meta-reinforcement learning is growing and developing at a rapid pace. In 2016, Steven S. Hansen presented a novel deep meta reinforcement learner called Deep Episodic Value Iteration (DEVI). DEVI is based on the model's internal structure, which permits one-time transitions to reward and transformational structural alterations, even for tasks with very high-dimensional state spaces. To address the small dynamic uncertainty in the dynamics due to a lack of sufficient training data, Clavera et al. (2017) suggested an approach to stabilize model-based reinforcement learning using gradient-based meta-learning. Qing Xiao et al. (2018) presented a model-based approximation meta-reinforcement learning (MB-ProMe) method to adaptively learn the most suitable control strategy for robot motion. Abbas Raza-Ali et al. (2019) proposed a meta-reinforcement learning strategy for optimizing parameters and hyperparameters simultaneously.

Jin Wang et al. 2020 proposed an MRL-based approach to MRLCO by merging a first-order MRL algorithm with a sequence-to-sequence (seq2seq) neural network. Newly presented by Rasool Fakoor and colleagues (meta-RL), Meta-Q-Learning (MQL) is an off-policy approach for meta-Reinforcement Learning. Meta-reinforcement learning algorithms were proposed by Louis Kirsch and others. MetaGenRL has the ability to transfer learning to novel, untrained contexts. A meta-reinforcement learning technique for MIER's model recognition and empirical overlap was proposed by Russell Mendonca et al. to enhance previous meta-reinforcement learning strategies.

Problem Formation - Problem components: Extensive-shape video games, sometimes known as SGs, are based on the assumption that each stage of the game is represented as a huge table, with rows and columns representing the actions of the two players. From their command center, sgs simulate the circumstances in which several marketers take action in tandem and reap the resulting benefits. In many real-world games, however, players take turns making moves. In games of the poker kind, the player who acts first may have a significant impact on the outcome of subsequent hands. An extensive-shape game (efg) (osborne and rubinstein, 1994; von neumann and morgenstern, 1945) uses a tree to design alternating-action video games. Recently, a major step forward has been taken in unifying the efgs and posgs frameworks thanks to the work of kovar'okay et al. (2019). In -player kuhn poker, the game tree is represented by the number 7 (kuhn, 1950b). Each player receives one card (the orange nodes in Figure 7), the dealer gets three cards (king > queen > jack), and the 0. 33 card is set aside - this is kuhn poker.

The sport then develops as follows.

- The first player to take action will examine or wager.
- If one player takes an exam, the second player may choose whether to verify the answers or take a wild guess.
- If player 2 takes a risk and checks their hand, the player with the higher card wins $1.
- If a player bets, you have the option to fold or call.
- If player 1 folds, player 2 receives $1 from player 1.
- If Player 1 makes the call, the player with the higher card receives $2 from Player 2.
- The first player to make a wager must then decide whether to fold or name.
- In the event that Player 2 folds, Player 1 receives $1 from Player 2.

If participant calls, then the higher card wins 2$ from the opposite participant. A vital feature of efgs is that they can manage imperfect facts for multiplayer choice making. In the instance of kuhn poker, the players do now not know which card the opponent holds. However, unlike dec-pomdp, which additionally fashions imperfect statistics inside the sg placing however is intractable to clear up, efg, represented in an equivalent collection form, can be solved by an lp in polynomial time in terms of game states (koller and megiddo, 1992). Within the subsequent section, we first introduce efg and then consider the sequence form of efg. Inclusion of the records sets in efg allows to version the imperfect-information instances in which gamers have most effective partial or no understanding about their fighters. In a game of kuhn poker, each participant may see just their own card. For instance, if player 1 is holding a jack and player 2 is holding a queen, player 1 has no way of knowing which card it is holding, hence its choice nodes under both scenarios (queen and king) remain in the same records set. In the exceptional situation of ideal-information efgs (such as cross or chess), the information set is a singleton; hence, the desire node may be equated to the unique history that terminates in it. Examples of imperfect-records efgs include poker and texas hold 'em, where a couple's history might be represented by a "information kingdom." However, with the assumption of ideal recall (defined later), the records that ends in an information kingdom continues to be specific Fixing EFGS thru the NFG representation, even though general, is inefficient because the size of the brought on NFG exponential in the range of facts states. The NFG example also ignores the sequential nature of video games. To alleviate these issues, you may focus your efforts on the realization-plan representation of the EFG, also known as the EFG collection form, whose dimensions are only linear in the number of game states. Using this method, we can get solutions to efgs in polynomial time (koller and megiddo, 1992). Within the efgs series structure, the emphasis switches from integrated methods to behavioral techniques, in which the marketers randomize separately at each information kingdom rather than across the whole set of natural approaches.

$i \in s x$, i. E., $\pi i : s i \rightarrow \Delta \chi(si)$ . Instead of basing a player's strategy on the perception of pure strategies, which can be exponentially many, the key insight of the collection shape is

that players can instead base their strategy on the paths in the game tree from the root to each node. Generally speaking, there is no comparison between the expressive power of behavioral approach and the blended technique. However, the behavioural approach and the combined technique are on par if sports have perfect do not forget, the intuitively 20 approach that each agent remembers all his history moves in particular facts indicates perfectly. In particular, consider that all selection nodes in a records kingdom have the same underlying records (or else the agent will be able to tell them apart). So long as this is the case, the well-known kuhn's theorem (kuhn, 1950a) guarantees that the expressive power of behavioral tactics and that of combination strategies correspond, inside the sense that they trigger the same chance on outcomes for games of perfect don't forget. Therefore, the set of ne no longer changes when optimal behavioral methods are taken into account.
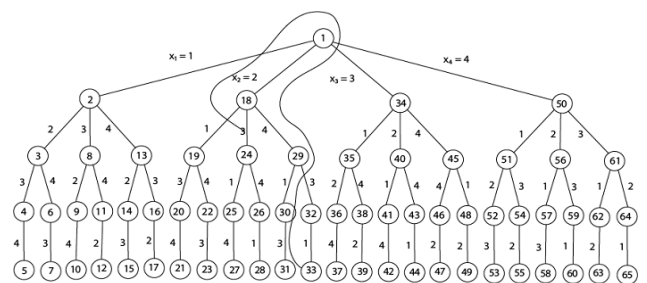


Fig. 5.

Solving EFGs via the NFG illustration, while familiar, is inefficient because the size of the prompted NFG is exponential in the number of facts states. The sequence-shape illustration, on the other hand, is often helpful for describing imperfect-data efgs of perfect remember, written as. In addition, games' temporal structure is no longer taken into account by the nfg representation. Working on the collection form of the EFG, also known as the realisation-plan illustration, is one technique to deal with these challenges, since its size is only linear in the number of game states and is therefore exponentially lower than that of the NFG. Using this method, we can get solutions to efgs in polynomial time (koller and megiddo, 1992). Within the EFGS sequence structure, the emphasis switches from blended strategies to behavioral strategies, in which the marketers randomize separately across all data kingdom s rather than throughout the whole of pure approaches. The key insight of the collection form is that a player's strategy can be built primarily based on the paths in the game tree from the base to each node, as opposed to the perception of pure techniques, which can be exponentially many. Both the behavioral approach and the hybrid approach have unparalleled expressive power in the present day. However, the behavioural approach and the combined strategy are functionally equal if the game possesses ideal don't forget, which intuitively20 indicates that each agent recalls all his past moves in exclusive information states correctly.

In particular, suppose that all selection nodes in an information kingdom have the same underlying data (else the agent can tell them apart). If this is the case, the famous kuhn's

theorem (kuhn, 1950a) guarantees that the expressive power of behavioural strategies and that of blended strategies correspond, in the sense that they trigger the same opportunity on outcomes for games of ideal memory. Therefore, taking the basic behavioral approaches into account, the set of ne does not change. Actually, the sequence-form illustration is usually helpful for illustrating imperfect-records efgs of best take into consideration, written as td errors, stores it within the replay buffer, and samples the statistics within the replay buffer to update to the Real time update. the component of cost. Since the fee feature for one agent relies on the actions of other merchants, the bootstrap method in td learning also necessitates sampling the actions of many merchants, which may be problematic in a multi-agent learning setting. To begin with, the activities that were sampled provide a small representation of the entire behavior of different vendors' underlying standards throughout various stages. Second, an agent's coverage might shift at any moment in training, making the replay buffer's samples quickly become stale. Since the agent is continually analyzing the current dynamics, the dynamics that produced the information in its replay buffer must be kept up to date. In a similar vein, this procedure worsens the non-stationarity problem. The non-stationarity problem essentially precludes analyzing single-agent algorithms in a multi-agent setting using the same mathematical tool. However, there is the identical-hobby game that deviates from this rule. All agents may reliably pursue their own interests in such a system without having to take into account the regulations of other agents. Therefore, the stationarity is preserved, and single-agent rl algorithms may be used.

## 2. Materials and Methods

Our project focused on developing autonomous parking models using the Proximity Policy Optimization (PPO) algorithm, a form of reinforcement learning. We created three different models of increasing complexity and precision; each using LIDAR sensors to navigate and avoid obstacles while searching for optimal parking spots. Our findings demonstrate that the PPO algorithm is effective in training autonomous agents to perform complex tasks like parking, and the use of LIDAR sensors improves accuracy. Our project highlights the importance of road safety and how autonomous agents can contribute to reducing accidents caused by human error. The development of our autonomous parking models is a step towards creating more advanced models for real-world applications such as self-driving cars. This project covers topics such as machine learning, reinforcement learning, meta-reinforcement learning, LIDAR sensors, PPO algorithm, and road safety, contributing to future research in AI applications.

As per the findings of our research paper, the PPO algorithm is an effective approach to training autonomous agents to perform complex tasks such as parking. The use of LIDAR sensors enables the agents to perceive the environment in a similar way to humans, resulting in more precise and efficient models. Our project showcases the development of three different autonomous parking models of increasing complexity and precision. The project highlights the significance of road

safety and how autonomous agents can contribute to reducing accidents caused by human error. The development of these models can be used as a foundation for more advanced autonomous driving models, including self-driving cars. Overall, our research project contributes to the field of artificial intelligence and its application in real-world scenarios.

### A. Limitations

Despite the encouraging outcomes, our initiative had significant drawbacks. First of all, the models were developed and evaluated in a virtual setting, which might not accurately represent actual settings. The performance of the models may be impacted by the absence of physical barriers and other elements that are present in real-world circumstances.

Second, while LIDAR sensors are useful for detecting the surroundings, they might be pricey and not available to everyone. This restricts the models' capacity to scale and their applicability in the actual world.

Last but not least, the highly sophisticated model's intricacy can render it challenging to use and maintain in real-world circumstances. The high degree of accuracy and fluidity required for driving may also call for a large amount of computer power.

These drawbacks underline the requirement for more study and advancement in the field of reinforcement learning and its use in practical situations.

## 3. Conclusion

In conclusion, according to the results of our study, the PPO method is a useful method for teaching autonomous agents to carry out challenging tasks like parking. LIDAR sensors provide agents the ability to see the world similarly to humans, leading to models that are more accurate and effective. Our experiment demonstrates the creation of three distinct autonomous parking models, each getting more complicated and precise. The study emphasises the importance of traffic safety and how autonomous agents might help decrease accidents brought on by human mistake. These models' advancement can serve as a basis for more sophisticated autonomous driving models, such as self-driving vehicles. In general, our research project makes a contribution to the study of artificial intelligence and its use in practical situations

## References

[1] "In Memoriam," AI Magazine, vol. 11, no. 3, pp. 1-2, 1990.
[2] Z. R. Gao Yang, Wang Hao, Cao Zhixin, "Study on an Average Reward Reinforcement Learning Algorithm," Vol. 30, No. 8, pp. 1372-1378, 2007.
[3] J. Z. Jurgen Schmidhuber, Marco Wiering, "Simple Principles Of Meta Learning," Rock and Soil Mechanics, vol. 31, pp. 155 -156, 2010.
[4] L. Q. Fu Qiming, Wang Hui, Xiao Fei, Yu Jun, Li Jiao, "A Novel Off Policy Q(A) Algorithm Based on Linear Function Approximation," Chinese Journal of Computers, vol. 37, no. 3, pp. 677-686, 2014.
[5] Z. M. Weiyingzi, "A Reinforcement Learning-based Approach to Dynamic Job-shop Scheduling," ACTA AUTO MATICA SINICA, vol. 31, no. 5, pp. 765-771, 2005.
[6] E. Ipek, Mutlu, "Self-Optimizing Memory Controllers: A Reinforcement Learning Approach," 2008 International Symposium on Computer Architecture, 2008.
[7] G. Tesauro, "TD-Gammon - A Self-Teaching Backgammon Program, Achieves Master-Level Play," AAAI Technical Report, pp. 19-23, 1993.

[8]  C. S. Levente Kocsis, "Bandit Based Monte-Carlo Planning," ECML 2006, pp. 282-293.

[9]  C. B. Andrea Banino, Benigno Uria, "Vector-based Navigation using Grid-like Representations1in Artificial Agents," Nature, pp. 429- 433, 2018.

[10] V. K. Mnih, K., "Human-level control through deep reinforcement learning," Nature, vol. 518, no. 7540, pp. 529-33, Feb 26 2015.

[11] Z. K. N. J. X. Wang, D. Tirumala, "Learning to Reinforcement Learn,", 2016.

[12] P. A. Chelsea Finn, Sergey Levine, "Model-Agnostic  Meta- Learning for Fast Adaptation of Deep Networks," 2017.

[13] Al-Shedivat M. Bansal T, Burda Y., "Continuous adaptive learning in non-stationary and competitive environments," 2017.