

Planification de Trajectoire pour Robot Autonome

en Environnement Dynamique par Méthodes de

Monte Carlo Tree Search (MCTS)

Projet – Cours de Méthodes de Monte Carlo
M2 IASD

KEITA Mamadi Afi Skander

16 février 2026

Résumé

Ce rapport présente l'étude et l'implémentation d'une méthode de planification de trajectoire pour un robot mobile autonome évoluant dans un environnement comportant des obstacles statiques et dynamiques. L'approche repose sur l'utilisation des algorithmes de Monte Carlo Tree Search (MCTS), en particulier UCT, RAVE et GRAVE, pour permettre au robot de déterminer un chemin optimal entre un point de départ et une destination cible. Nous traitons les problématiques d'évitement d'obstacles, d'adaptation face aux obstacles mobiles, de re-planification en temps réel et de gestion des incertitudes capteurs. Les résultats expérimentaux montrent que GRAVE offre le meilleur compromis performance/temps de calcul, notamment en environnement dynamique.

Table des matières

1	Introduction	4
1.1	Contexte et motivation	4
1.2	Pourquoi Monte Carlo Tree Search ?	4
1.3	Objectifs du projet	4
2	État de l'art	5
2.1	Fondements du MCTS	5
2.1.1	Principe de Monte Carlo	5
2.1.2	UCT : Upper Confidence Bound for Trees	5
2.1.3	Les quatre phases de MCTS	5
2.2	Améliorations : RAVE et GRAVE	5
2.2.1	RAVE : Rapid Action Value Estimation	5
2.2.2	GRAVE : Generalized RAVE	6
2.3	MCTS pour la planification robot	6
2.4	Progressive Widening	6
3	Méthodologie	6
3.1	Modélisation du problème	6
3.1.1	Environnement	6
3.1.2	Modélisation MDP	7
3.1.3	Obstacles dynamiques	7
3.1.4	Incertitude capteur	7
3.2	Algorithmes implémentés	7
3.2.1	Flat Monte Carlo (baseline)	7
3.2.2	UCT	7
3.2.3	RAVE	8
3.2.4	GRAVE	8
3.3	Re-planification dynamique	8
3.4	Planification multi-mondes	8
4	Implémentation	8
4.1	Architecture logicielle	8
4.2	Détails techniques	9
4.2.1	Conformité avec le cours	9
4.2.2	Paramètres	9
5	Expérimentations	9
5.1	Protocole expérimental	9
5.2	Scénarios de test	10
5.3	Résultats et analyse	10
5.3.1	Exp. 1 : Comparaison des algorithmes	10
5.3.2	Exp. 2 : Budget de simulations	10
5.3.3	Exp. 3 : Constante d'exploration	11
5.3.4	Exp. 4 : Environnement dynamique	11
5.3.5	Exp. 5 : Bruit capteur	11
5.3.6	Exp. 6 : Passage étroit	12

6	Discussion	12
6.1	Liens avec le cours	12
6.2	Limites et perspectives	13
7	Conclusion	14

1 Introduction

1.1 Contexte et motivation

La planification de trajectoire pour robots autonomes est un problème fondamental en robotique mobile. Un robot doit naviguer depuis un point de départ vers une destination en évitant les obstacles, tout en s'adaptant aux changements de l'environnement en temps réel. Ce problème se pose dans de nombreuses applications : robots de service, véhicules autonomes, drones, robots d'entrepôt, etc.

Les méthodes classiques de planification (A*, RRT, RRT*) supposent généralement un environnement statique et parfaitement connu. Or, dans la réalité, l'environnement est **dynamique** (personnes, autres robots qui se déplacent) et **partiellement observable** (incertitude des capteurs).

1.2 Pourquoi Monte Carlo Tree Search ?

Le MCTS, introduit indépendamment par [1] et [2], est particulièrement adapté à ce problème pour plusieurs raisons :

1. **Gestion de l'incertitude** : Les playouts aléatoires permettent d'échantillonner les futurs possibles de l'environnement, incluant les mouvements imprévisibles des obstacles dynamiques.
2. **Planification anytime** : MCTS peut être interrompu à tout moment et fournir la meilleure action trouvée jusqu'alors, ce qui est essentiel pour le temps réel.
3. **Pas de modèle complet nécessaire** : Contrairement à la programmation dynamique, MCTS n'a besoin que d'un simulateur de l'environnement.
4. **Équilibre exploration/exploitation** : La formule UCT (§2.1.2) garantit un équilibre optimal entre l'exploration de nouvelles trajectoires et l'exploitation des trajectoires prometteuses.

L'application du MCTS à la planification robot a été validée récemment par plusieurs travaux : [9] prouvent la convergence exponentielle de UCT pour le path planning, [12] combinent MCTS avec les Velocity Obstacles pour les environnements denses, et [11] publient dans *Science Robotics* une méthode MCTS pour les systèmes dynamiques continus.

1.3 Objectifs du projet

Ce projet vise à :

- Implémenter un environnement de simulation avec obstacles statiques et dynamiques
- Implémenter et comparer quatre algorithmes MCTS : Flat MC, UCT, RAVE, GRAVE
- Étudier la re-planification en temps réel
- Évaluer la robustesse face au bruit capteur
- Analyser l'impact des paramètres (budget de simulations, constante d'exploration)

2 État de l'art

2.1 Fondements du MCTS

2.1.1 Principe de Monte Carlo

Le principe fondamental des méthodes de Monte Carlo est d'estimer une quantité par échantillonnage aléatoire plutôt que par calcul exhaustif. Appliqué à la recherche dans les arbres de décision, cela donne le MCTS : au lieu d'explorer tout l'arbre (impossible pour les grands espaces d'états), on échantillonne des trajectoires aléatoires (*layouts*) pour estimer la valeur de chaque action browne2012survey.

2.1.2 UCT : Upper Confidence Bound for Trees

UCT kocsis2006bandit applique le principe des bandits multi-bras à la recherche dans les arbres. La formule de sélection est :

$$\text{UCT}(i) = \frac{w_i}{n_i} + C \sqrt{\frac{\ln N}{n_i}} \quad (1)$$

où w_i est la somme des récompenses du noeud i , n_i le nombre de visites, N le nombre de visites du parent, et C la constante d'exploration (typiquement $C = \sqrt{2}$).

Le premier terme w_i/n_i favorise l'**exploitation** (noeuds avec de bons résultats), tandis que le second terme $C\sqrt{\ln N/n_i}$ favorise l'**exploration** (noeuds peu visités).

2.1.3 Les quatre phases de MCTS

L'algorithme MCTS répète itérativement quatre phases :

1. **Sélection** : Descendre dans l'arbre en choisissant les noeuds avec le meilleur score UCT.
2. **Expansion** : Ajouter un nouveau noeud enfant pour une action non explorée.
3. **Simulation** : Jouer un playout aléatoire depuis le nouveau noeud.
4. **Rétropropagation** : Remonter le résultat et mettre à jour les statistiques.

2.2 Améliorations : RAVE et GRAVE

2.2.1 RAVE : Rapid Action Value Estimation

RAVE gelly2007combining accélère l'apprentissage en utilisant les statistiques AMAF (All Moves As First) : si une action m apparaît *n'importe où* dans un playout gagnant, elle est probablement bonne. La formule combine UCT et AMAF :

$$\text{Valeur}(m) = (1 - \beta_m) \cdot \text{UCT}_m + \beta_m \cdot \text{AMAF}_m \quad (2)$$

avec $\beta_m = \frac{p_m^{\text{AMAF}}}{p_m^{\text{AMAF}} + p_m + b \cdot p_m^{\text{AMAF}} \cdot p_m}$, où b est un biais (typiquement $b = 10^{-4}$).

2.2.2 GRAVE : Generalized RAVE

GRAVE cazenave2015generalized résout le problème de fiabilité des statistiques AMAF pour les noeuds peu visités : au lieu d'utiliser les statistiques AMAF du noeud courant, on utilise celles du premier ancêtre ayant plus de n_{ref} visites (typiquement $n_{\text{ref}} = 50$).

2.3 MCTS pour la planification robot

L'application de MCTS à la planification de trajectoire est un domaine de recherche actif :

- [9] introduisent Monte-Carlo Path Planning (MCPP) et prouvent la convergence exponentielle vers le chemin optimal, avec des expériences sur un robot réel.
- [10] proposent SL-MCTS combinant MCTS avec un réseau de neurones inspiré d'AlphaZero pour l'auto-apprentissage de la politique de recherche.
- [11] publient dans *Science Robotics* une méthode utilisant l'expansion spectrale pour gérer les espaces continus.
- [12] combinent MCTS avec les Velocity Obstacles pour la navigation sûre dans des environnements denses avec jusqu'à 40 obstacles dynamiques.
- [13] utilisent des RNN génératives dans MCTS pour prédire le comportement des piétons.
- [14] proposent MCMP pour l'optimisation de trajectoire sous incertitude avec des techniques de réduction de variance.

Le survey récent de [6] recense les applications modernes de MCTS au-delà des jeux, incluant la planification et la robotique.

2.4 Progressive Widening

Pour les problèmes avec des espaces d'actions continus (comme la robotique), [4] et [5] proposent le *Progressive Widening* : on n'ajoute un nouveau fils que lorsque $n(s)^{p_w} \geq |\text{enfants}(s)|$. Cela garantit qu'on explore en profondeur avant d'élargir l'arbre.

3 Méthodologie

3.1 Modélisation du problème

3.1.1 Environnement

L'environnement est une grille 2D de taille $N \times N$ où chaque cellule peut être :

- Libre (le robot peut y passer)
- Obstacle statique (mur, meuble – permanent)
- Obstacle dynamique (personne, autre robot – se déplace)
- Position du robot
- Position objectif

3.1.2 Modélisation MDP

Nous formalisons le problème comme un Processus de Décision Markovien (MDP) :

- **États** : $s = (\text{pos_robot}, \text{pos_obstacles}, t)$
- **Actions** : 9 actions (8 directions + attendre)
- **Transitions** : déterministes pour le robot, stochastiques pour les obstacles
- **Récompenses** :
 - Objectif atteint : $+100 \times (1 - t/t_{\max})$ (bonus de rapidité)
 - Collision : -100
 - Chaque pas : -1 (encourage les chemins courts)

La fonction de récompense avec bonus de rapidité est inspirée de la *discounting heuristic* du cours (section 20.1), qui pénalise les solutions longues.

3.1.3 Obstacles dynamiques

Trois patterns de mouvement sont implémentés :

1. **Linéaire** : va-et-vient en ligne droite (prédictible)
2. **Aléatoire** : direction choisie aléatoirement à chaque pas (imprédictible)
3. **Circulaire** : mouvement circulaire autour d'un point

3.1.4 Incertitude capteur

Le bruit capteur est modélisé par un paramètre $\sigma \in [0, 1]$: avec probabilité σ , une cellule observée peut être inversée (libre \leftrightarrow occupée). Cela crée un problème d'**information imparfaite**, analogue aux jeux étudiés dans le cours (section 19).

3.2 Algorithmes implémentés

3.2.1 Flat Monte Carlo (baseline)

Pour chaque action légale, on effectue N/K playouts aléatoires ($K = \text{nombre d'actions}$) et on choisit l'action avec la meilleure récompense moyenne. C'est la méthode la plus simple (cours section 1.2).

3.2.2 UCT

Implémentation standard de UCT avec les quatre phases de MCTS. Nous proposons deux variantes de playout :

- **Playout aléatoire** : actions choisies uniformément
- **Playout heuristique** : probabilité proportionnelle à $\exp(-d)$ où d est la distance Manhattan à l'objectif après l'action (échantillonnage de Gibbs, cours section 10.2)

3.2.3 RAVE

UCT avec statistiques AMAF. Après chaque playout, on enregistre toutes les actions jouées et on met à jour les statistiques AMAF de chaque noeud ancêtre. La formule (2) combine les deux sources d'information.

3.2.4 GRAVE

Amélioration de RAVE utilisant les statistiques AMAF du premier ancêtre fiable (ayant au moins $n_{\text{ref}} = 50$ visites).

3.3 Re-planification dynamique

La stratégie de re-planification est de type *receding horizon* : à chaque pas de temps, le robot :

1. Observe l'environnement (avec bruit capteur éventuel)
2. Lance une recherche MCTS pour choisir la prochaine action
3. Exécute l'action
4. Répète

Nous proposons également un planificateur **adaptatif** qui ajuste le budget de simulations selon le niveau de danger (proximité des obstacles dynamiques, étroitesse du passage).

3.4 Planification multi-mondes

Inspiré de PIMC (Perfect Information Monte Carlo, cours section 19.2), nous implémentons un planificateur qui échantillonne N mondes possibles (perturbations des positions des obstacles dynamiques) et choisit l'action par vote majoritaire. Cela est analogue à la Root Parallelization (cours section 7.1).

4 Implémentation

4.1 Architecture logicielle

Le code est organisé en modules Python :

- `environment.py` : Environnement de grille 2D, obstacles, capteurs
- `mcts_base.py` : Flat MC et UCT
- `mcts_rave.py` : RAVE et GRAVE
- `mcts_dynamic.py` : Re-planification dynamique, planification adaptative
- `visualization.py` : Visualisation matplotlib
- Tests unitaires et scripts d'expérimentation

4.2 Détails techniques

4.2.1 Conformité avec le cours

L'implémentation respecte les recommandations du cours (section 25) :

- Chaque état est **copié** avant simulation (erreur #1)
- Les playouts ont une **limite de pas** pour éviter les boucles infinies (erreur #6)
- Les résultats sont moyennés sur **suffisamment de parties** (erreur #5)
- L'action finale est choisie par **nombre de visites** (plus robuste que la valeur Q)

4.2.2 Paramètres

Paramètre	Symbole	Valeur
Constante d'exploration UCT	C	$\sqrt{2} \approx 1.414$
Budget de simulations	N	500 (défaut)
Biais RAVE	b	10^{-4}
Seuil GRAVE	n_{ref}	50
Profondeur max arbre	—	15 niveaux
Profondeur max playout	—	50 pas
Pas max par épisode	t_{\max}	200
Score position	γ^d	$\gamma = 0.9, d = \text{dist. Manhattan}$

Table 1: Paramètres par défaut (conformément aux recommandations du cours, section 25.2).

5 Expérimentations

5.1 Protocole expérimental

Six expérimentations sont conduites, chacune répétée 8 fois avec des graines aléatoires différentes pour obtenir des résultats statistiquement significatifs (cf. cours section 25.3, erreur #5).

1. Comparaison des algorithmes sur scénario statique
2. Impact du budget de simulations
3. Impact de la constante d'exploration C
4. Performance en environnement dynamique
5. Robustesse au bruit capteur
6. Scénario de passage étroit

5.2 Scénarios de test

- **Scénario simple** : Grille 15×15 , obstacles statiques (3 murs internes), départ $(1, 1)$, objectif $(13, 13)$
- **Scénario dynamique** : Grille 15×15 , 3 obstacles dynamiques (linéaire, vertical, aléatoire)
- **Passage étroit** : Mur traversant avec un seul passage, obstacle dynamique bloquant
- **Bruit capteur** : Niveaux $\sigma \in \{0\%, 5\%, 10\%, 20\%, 30\%\}$

5.3 Résultats et analyse

5.3.1 Exp. 1 : Comparaison des algorithmes

Sur le scénario statique (grille 15×15 , $N = 500$ simulations, 8 répétitions) :

Algorithme	Taux de succès	Pas moyens ($\pm \sigma$)	Temps/recherche (s)
Flat MC	100%	22.8 ± 2.9	0.261
UCT	100%	23.0 ± 4.3	0.319
RAVE	100%	22.5 ± 3.7	0.339
GRAVE	100%	22.4 ± 3.6	0.337

Table 2: Comparaison des algorithmes sur scénario statique ($N = 500$, 8 répétitions). GRAVE obtient la trajectoire la plus courte en moyenne avec l'écart-type le plus faible.

Analyse : Sur ce scénario simple, tous les algorithmes atteignent 100% de succès. La hiérarchie observée ($\text{GRAVE} \leq \text{RAVE} \leq \text{Flat MC} < \text{UCT}$ en nombre de pas) est cohérente avec les résultats du cours : les statistiques AMAF (RAVE/GRAVE) permettent un apprentissage plus rapide et des trajectoires plus courtes. Les différences sont faibles car la grille 15×15 est relativement simple (distance Manhattan de 24 entre départ et objectif).

5.3.2 Exp. 2 : Budget de simulations

L'impact du budget est mesuré avec l'algorithme GRAVE sur le scénario statique :

Budget N	Taux de succès	Pas moyens	Temps/recherche (s)
100	100%	43.4	0.067
300	100%	28.9	0.202
500	100%	22.4	0.338
1000	100%	20.0	0.678

Table 3: Impact du budget de simulations sur GRAVE. Rendements décroissants au-delà de $N = 500$.

Analyse : L'augmentation du budget améliore significativement la qualité des trajectoires : passer de 100 à 500 simulations réduit le nombre de pas de moitié ($43.4 \rightarrow 22.4$). Au-delà de 500, les rendements sont décroissants : doubler le budget (1000) ne gagne que 2.4 pas supplémentaires mais double le temps de calcul. Le budget $N = 500$ offre un bon compromis pour notre taille de grille.

5.3.3 Exp. 3 : Constante d'exploration

L'impact de C est mesuré avec UCT ($N = 500$) :

C	Taux de succès	Pas moyens
0.5	100%	20.1
1.0	100%	21.4
$\sqrt{2} \approx 1.414$	100%	23.0
2.0	100%	22.8
3.0	100%	23.0

Table 4: Impact de la constante d'exploration C sur UCT.

Analyse : Sur notre scénario, $C = 0.5$ donne les meilleurs résultats. Contrairement au résultat théorique ($C = \sqrt{2}$, optimal pour UCB1), une valeur plus faible favorise l'exploitation dans un espace d'actions restreint (9 directions). Avec seulement 9 actions, l'exploration complète est rapidement atteinte et un C élevé gaspille le budget en revisitant des actions déjà bien estimées. Ce résultat est cohérent avec les observations de [3] : la valeur optimale de C dépend du problème.

5.3.4 Exp. 4 : Environnement dynamique

Comparaison sur le scénario avec 3 obstacles dynamiques ($N = 500$, 8 répétitions) :

Algorithme	Taux de succès	Pas moyens ($\pm\sigma$)
UCT	75%	31.2 ± 8.4
RAVE	87.5%	28.6 ± 6.1
GRAVE	87.5%	27.4 ± 5.8

Table 5: Performance en environnement dynamique. La re-planification à chaque pas est essentielle.

Analyse : Les obstacles dynamiques réduisent le taux de succès (de 100% à 75–87.5%) et allongent les trajectoires. GRAVE et RAVE maintiennent une performance supérieure grâce à l'utilisation plus efficace du budget de simulations via AMAF : en propageant l'information sur toutes les actions d'un playout, ils apprennent plus vite quelles directions sont dangereuses. La re-planification à chaque pas (*receding horizon*) est essentielle car le plan initial devient obsolète quand les obstacles bougent.

5.3.5 Exp. 5 : Bruit capteur

Robustesse de GRAVE face au bruit capteur ($N = 500$) :

Bruit σ	Taux de succès	Pas moyens
0%	100%	22.4
5%	100%	24.1
10%	87.5%	28.7
20%	75%	35.2
30%	62.5%	42.8

Table 6: Impact du bruit capteur sur GRAVE. Dégradation progressive de la performance.

Analyse : La performance dégrade progressivement avec le bruit. Jusqu'à $\sigma = 5\%$, l'impact est négligeable (100% de succès). À $\sigma = 30\%$, le taux de succès chute à 62.5% car le robot peut percevoir des obstacles fantômes (faux positifs) ou manquer de vrais obstacles (faux négatifs, plus dangereux). Cela illustre le lien avec les jeux à information imparfaite (cours section 19) : le robot doit prendre des décisions avec une perception bruitée de l'environnement.

5.3.6 Exp. 6 : Passage étroit

Scénario le plus difficile : mur traversant avec un seul passage, obstacle dynamique bloquant ($N = 800$, 8 répétitions) :

Algorithme	Taux de succès	Pas moyens
UCT	50%	58.3
GRAVE	62.5%	52.1

Table 7: Passage étroit avec obstacle dynamique ($N = 800$).

Analyse : Ce scénario requiert de la *patience* : le robot doit parfois attendre que l'obstacle dynamique s'écarte pour traverser le passage. GRAVE gère mieux cette situation grâce à une meilleure estimation de la valeur de l'action WAIT via les statistiques AMAF. Le budget élevé ($N = 800$) est nécessaire pour explorer suffisamment de futurs possibles et trouver le bon timing de passage.

6 Discussion

6.1 Liens avec le cours

Ce projet met en pratique la plupart des concepts vus en cours :

- **Flat MC** (section 1.2) comme baseline
- **UCB/UCT** (sections 2–3) comme algorithme principal
- **RAVE/GRAVE** (sections 5–6) pour l'accélération
- **Playout heuristique** inspiré de PPA (section 13) et de l'échantillonnage de Gibbs (section 10.2)
- **Re-planification** liée aux jeux à information imparfaite (section 19)
- **Progressive Widening** (section 16) applicable pour les actions continues
- **Planification multi-mondes** inspirée de PIMC (section 19.2)
- **Planification adaptative** inspirée de Sequential Halving (section 11)

6.2 Limites et perspectives

- **Espace discret** : Notre grille est discrète. Pour un robot réel, le Progressive Widening couetoux2011continuous serait nécessaire pour gérer les actions continues.
- **Scalabilité** : Pour de grandes grilles, la parallélisation (section 7 du cours) serait bénéfique.
- **Apprentissage** : L'intégration de réseaux de neurones à la manière d'AlphaGo Zero silver2017mastering pourrait améliorer les playouts et l'évaluation.
- **Multi-agent** : L'extension au cas multi-robot nécessiterait des adaptations comme le Pareto MCTS chen2019pareto.

7 Conclusion

Ce projet démontre l'efficacité des méthodes MCTS pour la planification de trajectoire en environnement dynamique. Sur le scénario statique, tous les algorithmes atteignent 100% de succès avec des trajectoires quasi-optimales (22 pas pour une distance Manhattan de 24). Les différences entre algorithmes deviennent significatives dans les scénarios difficiles : GRAVE maintient 87.5% de succès en environnement dynamique (contre 75% pour UCT) et 62.5% en passage étroit (contre 50%).

L'étude du budget de simulations révèle des rendements décroissants au-delà de $N = 500$, et l'analyse du bruit capteur montre une robustesse acceptable jusqu'à $\sigma = 10\%$. La re-planification à chaque pas de temps (*receding horizon*) est essentielle pour gérer les obstacles dynamiques.

Les résultats confirment que les concepts théoriques des méthodes de Monte Carlo vus en cours (UCT, RAVE, GRAVE, information imparfaite) se transfèrent efficacement au domaine de la planification robotique.

Références

- [1] Kocsis, L., Szepesvári, C. (2006). *Bandit Based Monte-Carlo Planning*. European Conference on Machine Learning (ECML), pp. 282–293. DOI: [10.1007/11871842_29](https://doi.org/10.1007/11871842_29)
- [2] Coulom, R. (2006). *Efficient Selectivity and Backup Operators in Monte-Carlo Tree Search*. 5th International Conference on Computers and Games. DOI: [10.1007/978-3-540-75538-8_7](https://doi.org/10.1007/978-3-540-75538-8_7)
- [3] Browne, C.B., Powley, E., Whitehouse, D., et al. (2012). *A Survey of Monte Carlo Tree Search Methods*. IEEE Trans. on Computational Intelligence and AI in Games, 4(1), pp. 1–43. DOI: [10.1109/TCIAIG.2012.2186810](https://doi.org/10.1109/TCIAIG.2012.2186810)
- [4] Chaslot, G., Winands, M., van den Herik, H., et al. (2008). *Progressive Strategies for Monte-Carlo Tree Search*. New Mathematics and Natural Computation, 4(3), pp. 343–357. DOI: [10.1142/S1793005708001094](https://doi.org/10.1142/S1793005708001094)
- [5] Couëtoux, A., Hoock, J.B., Sokolovska, N., et al. (2011). *Continuous Upper Confidence Trees*. Learning and Intelligent Optimization (LION), pp. 433–445. DOI: [10.1007/978-3-642-25566-3_32](https://doi.org/10.1007/978-3-642-25566-3_32)
- [6] Świechowski, M., Godlewski, K., Sawicki, B., Mandziuk, J. (2023). *Monte Carlo Tree Search: A Review of Recent Modifications and Applications*. Artificial Intelligence Review, 56, pp. 2497–2562. DOI: [10.1007/s10462-022-10228-y](https://doi.org/10.1007/s10462-022-10228-y)
- [7] Gelly, S., Silver, D. (2007). *Combining Online and Offline Knowledge in UCT*. International Conference on Machine Learning (ICML), pp. 273–280.
- [8] Cazenave, T. (2015). *Generalized Rapid Action Value Estimation*. International Joint Conference on Artificial Intelligence (IJCAI), pp. 754–760.
- [9] Dam, T.T., Chalvatzaki, G., Peters, J., Pajarin, J. (2022). *Monte-Carlo Robot Path Planning*. IEEE Robotics and Automation Letters, 7(4), pp. 11213–11220. DOI: [10.1109/LRA.2022.3199674](https://doi.org/10.1109/LRA.2022.3199674)
- [10] Li, W., Liu, Y., Ma, Y., et al. (2023). *A Self-Learning Monte Carlo Tree Search Algorithm for Robot Path Planning*. Frontiers in Neurorobotics, 17, 1039644. DOI: [10.3389/fnbot.2023.1039644](https://doi.org/10.3389/fnbot.2023.1039644)
- [11] Rivière, B., Lathrop, J., Chung, S.J. (2024). *MCTS with Spectral Expansion for Planning with Dynamical Systems*. Science Robotics, 9(97). DOI: [10.1126/scirobotics.ado1010](https://doi.org/10.1126/scirobotics.ado1010)
- [12] Bonanni, L., Meli, D., Castellini, A., Farinelli, A. (2025). *MCTS with Velocity Obstacles for Safe and Efficient Motion Planning in Dynamic Environments*. AAMAS 2025. ArXiv: [2501.09649](https://arxiv.org/abs/2501.09649)
- [13] Eiffert, S., Kong, H., Pirmarzdashti, N., Sukkarieh, S. (2020). *Path Planning in Dynamic Environments Using Generative RNNs and MCTS*. IEEE ICRA, pp. 10263–10269. DOI: [10.1109/ICRA40945.2020.9196631](https://doi.org/10.1109/ICRA40945.2020.9196631)

- [14] Janson, L., Schmerling, E., Pavone, M. (2015). *Monte Carlo Motion Planning for Robot Trajectory Optimization Under Uncertainty*. ISRR, pp. 343–361. DOI: [10.1007/978-3-319-60916-4_20](https://doi.org/10.1007/978-3-319-60916-4_20)
- [15] Chen, W., Liu, L. (2019). *Pareto Monte Carlo Tree Search for Multi-Objective Informative Planning*. Robotics: Science and Systems (RSS).
- [16] Silver, D., Schrittwieser, J., Simonyan, K., et al. (2017). *Mastering the Game of Go without Human Knowledge*. Nature, 550(7676), pp. 354–359. DOI: [10.1038/nature24270](https://doi.org/10.1038/nature24270)