

POLYTECHNIQUE MONTRÉAL

affiliée à l'Université de Montréal

**Modélisation, simulation et contrôle par apprentissage par renforcement d'un
robot de type dirigeable**

ALICE LEMIEUX-BOURQUE

Département de génie informatique et génie logiciel

Mémoire présenté en vue de l'obtention du diplôme de *Maîtrise ès sciences appliquées*
Génie informatique

Avril 2025

POLYTECHNIQUE MONTRÉAL

affiliée à l'Université de Montréal

Ce mémoire intitulé :

Modélisation, simulation et contrôle par apprentissage par renforcement d'un robot de type dirigeable

présenté par **Alice LEMIEUX-BOURQUE**

en vue de l'obtention du diplôme de *Maîtrise ès sciences appliquées*
a été dûment accepté par le jury d'examen constitué de :

Benjamin DE LEENER, président

Giovanni BELTRAME, membre et directeur de recherche

David SAUSSIÉ, membre

REMERCIEMENTS

Je suis reconnaissante d'avoir pu apprendre et évoluer au sein d'un environnement comme le MISTLab. Je tiens à remercier Giovanni Beltrame pour l'accueil et la confiance qu'il m'a accordée.

I would like to thank all the members of MISTLab for their time and guidance throughout my Master's program. I am especially grateful to Hassan Fouad, who served as a mentor and was always available to support me.

I would also like to thank Professor Ely Carneiro De Paiva for the many helpful insights and the time he took to share his expertise with me. Obrigada !

Je veux souligner l'appui financier des Fonds de recherche du Québec - Nature et technologie (#331556) et du Conseil de recherche en sciences naturelles et en génie du Canada. I would like to highlight the financial support of the Fonds de recherche du Québec - Nature et technologie (#331556) and the National Sciences and Engineering Research Council of Canada.

Je tiens aussi remercier Polytechnique Montréal et toutes les personnes que j'y ai cotoyées depuis mon arrivée en 2017.

Finalement, merci à mes amies et amis, ma famille et mon partenaire Charles-Étienne pour leur amour et leur soutien inconditionnels.

RÉSUMÉ

Le secteur des transports est l'un des plus polluants au monde. Que ce soit pour le transport de personnes par avion à travers les océans ou pour le transport de marchandises par camion, cette industrie produit près de 15% des émissions de gaz à effet de serre au niveau mondial. Pour répondre à cet enjeu, au-delà de l'électrification des transports, certains rêvent d'un grand retour des dirigeables. Ces aéronefs se présentent en effet comme des candidats prometteurs pour pallier aux problèmes de l'industrie du transport : ils ont une très faible consommation d'énergie puisque la majeure partie de la force de portance est générée par le gaz de levage, ils sont peu bruyants et leur capacité à rester en vol stationnaire leur permet d'atteindre des zones inaccessibles aux avions et camions.

Toutefois, l'utilisation des dirigeables comporte plusieurs défis. Parmi ceux-ci se trouve notamment la difficulté de créer des contrôleurs de vols automatiques sécuritaires. En effet, en raison de leur grande surface, les dirigeables sont très sensibles aux perturbations externes comme les rafales de vent. Les forces aérodynamiques sur le dirigeable étant difficiles à modéliser, les méthodes classiques peinent à contrôler adéquatement les dirigeables. Une nouvelle approche consiste à utiliser l'apprentissage automatique, et plus spécifiquement l'apprentissage par renforcement, pour créer des politiques de contrôles applicables aux dirigeables. Toutefois, il est difficile de créer des politiques adéquates puisqu'il est presque impossible de les tester sur des systèmes de dirigeables réels.

Dans ce mémoire, une plateforme complète permettant l'entraînement et le déploiement de politiques d'apprentissage par renforcement sur un petit dirigeable est développée. Elle comprend notamment la conception et la création d'un robot dirigeable de deux mètres pour utilisation à l'intérieur. Le robot est équipé de quatre moteurs et deux gouvernes à l'arrière. Pour accompagner le robot, un simulateur a été créé avec MuJoCo, un moteur physique populaire en robotique. MuJoCo ne proposant pas de modélisation des interactions avec les fluides adaptée aux dirigeables, un modèle aérodynamique a été développé afin d'être intégré au simulateur.

Des algorithmes d'apprentissage par renforcement sont ensuite testés sur le simulateur. Deux tâches de contrôle sont étudiées : le maintien du vol stationnaire à une altitude donnée et le maintien d'une vitesse constante à une altitude fixe. Les résultats obtenus valident la capacité du système à apprendre des comportements de vol grâce à l'apprentissage par renforcement.

Finalement, une validation initiale de la plateforme complète est réalisée sur le dirigeable physique pour un problème unidimensionnel. Une politique de contrôle apprise par renforcement

en simulation est déployée sur le dirigeable, démontrant ainsi la capacité de la plateforme développée à servir de banc d'essai pour le développement et la validation de politiques de contrôle par apprentissage par renforcement.

Cette plateforme ouvre la voie à des travaux futurs axés sur le développement d'algorithmes de contrôle plus avancés et sécuritaires pour les dirigeables.

ABSTRACT

The transportation sector is one of the most polluting industries in the world. Whether it is the transportation of people by plane across the ocean or the transportation of goods by truck, the industry accounts for nearly 15% of global greenhouse gas emissions. To address this issue, beyond the electrification of transportation, some dream of a great return of airships. These unique aircrafts are indeed seen as an ideal candidate to address the problems of the transportation industry: they have a very low energy consumption since the majority of their lift is generated by the lifting gas that allows them to float, they are relatively quiet, and they can hover, which allows them to reach areas where airplanes or trucks cannot go.

However, the use of airships comes with several challenges. One of the most significant is the difficulty of designing safe automatic flight controllers. Due to their large surface area, airships are highly sensitive to external disturbances such as wind gusts. The aerodynamic forces acting on airships are difficult to model, and traditional control methods struggle to adequately control them. A new approach involves using machine learning, and specifically reinforcement learning, to create control policies for airships. However, it is challenging to develop adequate policies since it is nearly impossible to test them on real airship systems.

In this thesis, a complete framework is developed for training and deploying reinforcement learning policies on a small airship. It includes the design and creation of a two-meter-long indoor airship robot. The robot has four motors and two rudders at the rear. To accompany the robot, a simulator was created using MuJoCo, a popular physics engine in robotics. Since MuJoCo does not provide an adequate fluid interaction model for airships, an aerodynamic model was developed and integrated to the simulator.

Reinforcement learning algorithms are tested on the simulator. Two control tasks are studied: maintaining a fixed altitude in hover and sustaining a constant forward velocity at a given height. The results confirm the system's ability to learn flight behaviors using reinforcement learning.

Finally, the complete framework is validated on the physical airship in a one-dimensional control task. A control policy trained in simulation is successfully deployed on the real system, demonstrating the viability of the platform as a testbed for developing and validating reinforcement learning-based control strategies.

This platform paves the way for future research focused on developing more advanced and safe control algorithms for autonomous airships.

TABLE DES MATIÈRES

REMERCIEMENTS	iii
RÉSUMÉ	iv
ABSTRACT	vi
TABLE DES MATIÈRES	vii
LISTE DES TABLEAUX	x
LISTE DES FIGURES	xi
LISTE DES SIGLES ET ABRÉVIATIONS	xii
LISTE DES ANNEXES	xiii
 CHAPITRE 1 INTRODUCTION	1
1.1 Définitions et concepts de base	1
1.1.1 Dirigeables	1
1.1.2 Opération d'un dirigeable	2
1.1.3 Masse ajoutée	2
1.1.4 Apprentissage par renforcement	3
1.2 Éléments de la problématique	4
1.2.1 Conception d'un dirigeable pour utilisations dans des espaces fermés .	5
1.2.2 Modélisation aérodynamique	5
1.2.3 Limitation des simulateurs	5
1.2.4 Déploiement d'algorithmes d'apprentissage par renforcement	6
1.2.5 Apprentissage par renforcement sur des systèmes réels	6
1.3 Objectifs de recherche	6
1.4 Plan du mémoire	7
 CHAPITRE 2 REVUE DE LITTÉRATURE	8
2.1 Dirigeables de petite et moyenne taille	8
2.2 Modélisation dynamiques des dirigeables	9
2.3 Simulateurs pour apprentissage par renforcement	10
2.4 Contrôle des dirigeables	11

2.5 Apprentissage par renforcement appliqué aux dirigeables	11
CHAPITRE 3 DÉTAILS DE LA SOLUTION	13
3.1 Le prototype physique	13
3.1.1 L'enveloppe	13
3.1.2 Le système de commande	14
3.1.3 La gondole	15
3.1.4 Le centre de masse	15
3.1.5 Le système de communication	16
3.1.6 La conception électronique	16
3.1.7 Les capteurs	16
3.1.8 Le prototype final	18
3.2 Le modèle mathématique	19
3.2.1 Cadres référentiels	19
3.2.2 Matrice de masses	20
3.2.3 Forces dynamiques	21
3.2.4 Forces de propulsion	21
3.2.5 Forces aérodynamiques	22
3.2.6 Forces de flottabilité et gravité	26
3.3 Le simulateur	26
CHAPITRE 4 RÉSULTATS THÉORIQUES ET EXPÉRIMENTAUX	28
4.1 Les environnements	28
4.1.1 Vol stationnaire à altitude spécifiée	29
4.1.2 Vol de croisière à vitesse spécifiée	31
4.2 Résultats en simulation	32
4.2.1 Vol stationnaire à altitude spécifiée	33
4.2.2 Vol de croisière à vitesse spécifiée	38
4.3 Validation du système pour déploiement sur le système physique	43
4.4 Conclusion	45
CHAPITRE 5 CONCLUSION	47
5.1 Synthèse des travaux	47
5.2 Limitations de la solution proposée	47
5.3 Améliorations futures	48
RÉFÉRENCES	49

ANNEXES	55
-------------------	----

LISTE DES TABLEAUX

Tableau 3.1	Coefficients utilisés pour calculer les forces aérodynamiques	24
Tableau 4.1	Table de correspondance des actions pour la tâche de vol stationnaire	29
Tableau 4.2	Table de correspondance des actions pour la tâche de vol de croisière	31
Tableau 4.3	Paramètres de l'environnement pour la tâche de vol stationnaire . . .	33
Tableau 4.4	Paramètres de l'environnement pour la tâche de vol de croisière . . .	38
Tableau 4.5	Paramètres de l'environnement pour contrôle en une dimension . . .	44
Tableau A.1	Tableau des définitions des variables utilisées dans le modèle dynamique du dirigeable.	56

LISTE DES FIGURES

Figure 1.1	Apprentissage par renforcement	4
Figure 3.1	Électronique	17
Figure 3.2	Design 3D de RAFALE	18
Figure 3.3	Dirigeable physique	18
Figure 3.4	Axes de référence du dirigeable	20
Figure 3.5	Intégration de l'aérodynamique avec MuJoCo	27
Figure 4.1	Courbes d'entraînement pour la tâche de vol stationnaire à altitude donnée pour trois différents algorithmes	34
Figure 4.2	Courbes d'entraînement de PPO pour la tâche de vol stationnaire à altitude donnée pour trois <i>seeds</i> différentes	34
Figure 4.3	Observations selon le temps pour différentes altitudes	35
Figure 4.4	Comparaison des observations entre un PID et PPO pour une altitude de 4m	36
Figure 4.5	Comparaison des observations entre un PID et PPO pour une altitude de 2m	37
Figure 4.6	Résultats de l'entraînement pour la tâche de vol en mode croisière . .	39
Figure 4.7	Visualisation des observations générées par le modèle entraîné : 4000 pas de temps et vitesse cible de 0,1 m/s	40
Figure 4.8	Visualisation des observations générées par le modèle entraîné : 4000 pas de temps et vitesse cible de 0,5 m/s	41
Figure 4.9	Visualisation des observations générées par le modèle entraîné : 4000 pas de temps et vitesse cible de 1 m/s	42
Figure 4.10	Mouvements de la gouverne de profondeur du dirigeable selon l'angle de tangage	43
Figure 4.11	Résultats sur le système physique	45

LISTE DES SIGLES ET ABRÉVIATIONS

A2C	<i>Advantage Actor Critic</i>
BEC	Régulateur de tension électrique
CBF	<i>Control Barrier Function</i>
CFD	Mécanique des fluides numérique
DRL	Apprentissage par renforcement profond
ESC	Variateur électronique de vitesse
IMU	Centrale inertielle
INDI	Inversion dynamique non linéaire incrémentale
NDI	Inversion dynamique non linéaire
PPO	<i>Proximal policy optimization</i>
QR-DQN	<i>Quantile Regression Deep Q Network</i>
RL	Apprentissage par renforcement
RTF	Facteur temps réel

LISTE DES ANNEXES

Annexe A	Définitions des variables du modèle dynamique du dirigeable	55
----------	---	----

CHAPITRE 1 INTRODUCTION

En 2023, le secteur des transports se classait au deuxième rang des émetteurs de gaz à effet de serre, représentant 15% des émissions au niveau mondial : 11% pour le transport routier et 2% pour l'aviation [1]. Différentes initiatives visant à réduire l'empreinte carbone de ce secteur ont émergé, notamment le développement de carburants à zéro émission et l'électrification des transports. Une idée soutenue par certaines entreprises est de repenser les modes de transport eux-mêmes et d'envisager le développement et l'utilisation de dirigeables comme une alternative pour le transport de marchandises. Leur faible consommation énergétique pourrait faire des dirigeables une option pour concurrencer le transport par avion, mais aussi le transport routier sur de longues distances [2].

À plus petite échelle, les robots dirigeables pourraient aussi être utilisés dans les tâches de longue durée, leur autonomie prolongée et leur capacité à rester en vol stationnaire les rendant efficaces dans des applications telles que la collecte de données [3] [4].

La construction de dirigeables soulève toutefois plusieurs défis. Parmi ceux-ci, on peut citer leur faible acceptabilité sociale, en partie attribuable à l'accident du Hindenburg en 1937, ainsi que la complexité de leur contrôle, notamment en raison de leur grande sensibilité aux conditions météorologiques et perturbations externes.

Ce travail a été réalisé dans le cadre d'une collaboration entre le MIST Lab de Polytechnique Montréal, Thales et Flying Whales Québec. L'objectif global de cette collaboration est la création de politiques d'apprentissage par renforcement certifiables pour des aéronefs comme les dirigeables. Ce mémoire décrit la conception d'une plateforme de test et d'un simulateur permettant d'appliquer des politiques d'apprentissage par renforcement sur un robot de type dirigeable. Les travaux ont été menés entre août 2022 et avril 2025.

1.1 Définitions et concepts de base

1.1.1 Dirigeables

Un dirigeable est un aéronef motorisé utilisant un gaz de levage, généralement de l'hélium ou de l'hydrogène, pour générer une flottabilité permettant de réduire son besoin de propulsion verticale. Il se distingue entre autres par sa faible consommation d'énergie pour rester en vol, lui permettant d'être en l'air pour de longues périodes. Le dirigeable possède à la fois certains avantages des avions, notamment la capacité à parcourir de longues distances et à transporter de lourdes charges, et certains avantages des hélicoptères, notamment la possibilité de décoller

et atterrir verticalement, ainsi que rester en vol stationnaire. Il a traditionnellement une forme allongée qui permet de réduire la force de traînée lors du vol.

Les dirigeables sont généralement classés en trois grandes catégories : les dirigeables rigides, qui possèdent une structure légère limitant les déformations, les dirigeables souples, qui ne possèdent pas de structure interne et dont la forme est déterminée par la pression du gaz de levage à l'intérieur, et les dirigeables semi-rigides, qui se situent quelque part entre les deux autres types.

Les dirigeables peuvent également être différenciés en fonction de leur taille. Les plus grands dirigeables, comme le LCA60T de Flying Whales, peuvent transporter des charges lourdes sur de longues distances. Les plus petits peuvent quant à eux être utilisés pour la collecte de données, la publicité, voire des applications à l'intérieur de bâtiments ou d'espaces restreints.

Les principales composantes d'un dirigeable sont : l'enveloppe, qui contient le gaz de levage et contribue largement à la traînée aérodynamique ; la gondole, qui peut contenir les systèmes de commandes, les passagers et la charge utile ; le système de propulsion et le système d'empennage.

1.1.2 Opération d'un dirigeable

Le contrôle des dirigeables repose sur différents systèmes adaptés aux différentes phases de vol. À basse vitesse ou en vol stationnaire, le système de propulsion regroupant les différents moteurs permet le contrôle de la direction, la vitesse et l'altitude du véhicule. C'est le même système qui est utilisé pour les manœuvres de décollage et d'atterrissage. En mode croisière, donc à des vitesses plus élevées, le contrôle combine le système de propulsion et le système de gouvernes, qui permet une meilleure stabilité. Pour les plus grands dirigeables, des systèmes de ballonnets ou de relâchement de gaz peuvent être utilisés pour contrôler l'altitude [5].

1.1.3 Masse ajoutée

La masse ajoutée, ou masse virtuelle, représente la masse supplémentaire qu'un objet doit déplacer lorsqu'il se déplace dans un fluide. Cette masse ajoutée est due à l'inertie du fluide environnant qui résiste au mouvement de l'objet : c'est l'équivalent de la masse de l'air qui est déplacée avec l'objet.

La masse ajoutée intervient dans toutes les situations où un objet accélère ou décélère dans un fluide, mais elle est particulièrement importante dans le cas des dirigeables : en raison de leur faible densité, la masse ajoutée est proportionnellement plus importante. [6].

Afin d'estimer les forces et moments supplémentaires générés par le fluide en mouvement autour du dirigeable, la masse ajoutée est généralement modélisée à l'aide de coefficients d'inertie, comme ceux proposés par Lamb pour des formes ellipsoïdales [7].

Les masses ajoutées et les inerties associées peuvent être estimées comme suit :

$$\begin{aligned}
 X_{\dot{u}} &= -k_1 \frac{B}{g} \\
 Y_{\dot{v}} &= -k_2 \frac{B}{g} \\
 Z_{\dot{w}} &= Y_{\dot{v}} \\
 M_{\dot{q}} &= -k' \frac{B}{g} \left[\frac{l^2 + d^2}{20} \right] \\
 N_{\dot{r}} &= M_{\dot{q}} \\
 L_{\dot{p}} &= 0
 \end{aligned} \tag{1.1}$$

où B est la force de flottabilité, g est l'accélération gravitationnelle, k_1 , k_2 et k' sont les ratios d'inertie de Lamb, l est la longueur du dirigeable et d est son diamètre maximal. Les équations $Z_{\dot{w}} = Y_{\dot{v}}$, $N_{\dot{r}} = M_{\dot{q}}$ et $L_{\dot{p}} = 0$ sont de bonnes approximations si on ignore les effets de la gondole et de l'empennage, ce qui est une pratique courante, car l'effet de la masse ajoutée est principalement attribuable à l'enveloppe du dirigeable.

1.1.4 Apprentissage par renforcement

L'apprentissage par renforcement, ou RL (*Reinforcement Learning*), est un sous-domaine de l'apprentissage automatique. C'est une méthode fonctionnant par essai et erreur où un agent apprend à partir de ses interactions avec l'environnement. L'agent reçoit une observation de l'environnement, choisit une action à effectuer puis reçoit une récompense ou une punition en fonction de l'action sélectionnée. Le signal de récompense reçu est utilisé pour modéliser la tâche à effectuer. L'agent apprend donc à effectuer la tâche en maximisant le signal de récompense. C'est la méthode d'apprentissage automatique privilégiée pour les problèmes de contrôle car elle permet de développer des politiques de contrôle flexibles qui peuvent s'adapter à des environnements complexes et incertains. La figure 1.1 le fonctionnement général de l'apprentissage par renforcement.

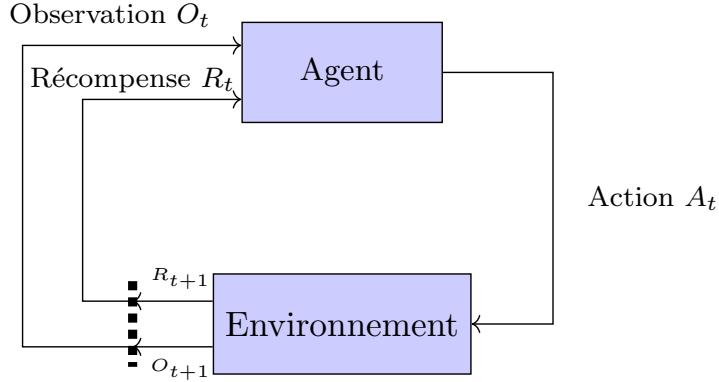


FIGURE 1.1 Apprentissage par renforcement

1.2 Éléments de la problématique

Le contrôle de dirigeable représente un défi en raison des dynamiques complexes et de la sensibilité des dirigeables aux perturbations externes, une conséquence directe de sa surface importante. Parallèlement au regain d'intérêt pour les dirigeables, les entreprises du secteur cherchent à implémenter des systèmes de contrôle basés sur l'apprentissage automatique, plus spécifiquement l'apprentissage par renforcement.

Toutefois, l'apprentissage par renforcement fait face à de nombreux défis lors de son application sur les systèmes réels. En effet, les politiques de contrôles apprises par renforcement n'ont généralement pas de garanties de sécurité. Pour les compagnies développant des dirigeables, l'application sur les systèmes réels est souvent impossible pour deux principales raisons. Premièrement, les systèmes sont généralement encore en phase de développement. Deuxièmement, le déploiement de politiques de contrôles apprises par renforcement n'étant pas sécuritaire, il présente des risques et pourrait engendrer des incidents extrêmement coûteux.

Cette recherche a donc pour objectif de développer une plateforme de test complète comprenant un robot dirigeable à faible coût, un simulateur adapté et une infrastructure permettant la création, l'entraînement et le déploiement de politiques d'apprentissage par renforcement sur le système. Afin de faciliter son utilisation, le dirigeable conçu doit être adapté aux espaces de laboratoires disponibles pour l'équipe.

1.2.1 Conception d'un dirigeable pour utilisations dans des espaces fermés

Bien que l'utilisation d'un dirigeable dans des espaces fermés permette d'éliminer certaines variables comme les changements de volume ou de température, cette approche introduit des contraintes spécifiques.

Premièrement, la taille du robot est un facteur important à considérer. Le dirigeable doit être suffisamment petit pour pouvoir manœuvrer dans un espace restreint, tout en conservant une flottabilité suffisante pour pouvoir maintenir son avantage énergétique par rapport aux drones conventionnels.

Ensuite, un avantage des dirigeables est qu'ils peuvent soulever d'importantes charges utiles. Toutefois, un petit dirigeable ne permet pas une grosse charge utile, et son fonctionnement est donc particulièrement sensible à la masse de ses composantes. Le dirigeable doit donc fonctionner avec un nombre limité d'actionneurs et de capteurs, tout en permettant un contrôle adéquat.

Finalement, comme il fonctionne dans des espaces restreints, le dirigeable doit pouvoir à la fois voler en mode stationnaire et en vitesse de croisière, malgré les limitations du nombre d'actionneurs imposés par les contraintes de masse.

1.2.2 Modélisation aérodynamique

La modélisation d'un dirigeable doit prendre en compte les effets aérodynamiques qui influencent considérablement son mouvement. Un modèle fidèle à la réalité est essentiel pour développer des politiques d'apprentissage par renforcement transférables à la réalité.

Déterminer les coefficients aérodynamiques est un domaine de recherche en soi. Les méthodes traditionnelles permettant d'obtenir ces coefficients, soit les essais en soufflerie ou la simulation par mécanique des fluides numérique (CFD), sont généralement chronophages, coûteuses en ressources ou demandent une expertise spécifique. De plus, les modèles existants se limitent généralement à de petits angles d'attaques et ne modélisent pas adéquatement les mouvements d'un dirigeable en vol stationnaire.

1.2.3 Limitation des simulateurs

Il existe très peu de simulateurs en logiciel libre permettant de simuler le mouvement de dirigeables, et ceux qui existent se limitent souvent à un mouvement en deux dimensions. Cette réalité contraste avec la disponibilité de nombreux modèles pour les drones plus traditionnels. Ainsi, les équipes de recherche souhaitant appliquer leurs algorithmes d'apprentissage

par renforcement aux dirigeables sont donc souvent obligées de développer leurs propres modèles et outils.

Les plateformes de simulation physique comme MuJoCo offrent un cadre intéressant pour le développement d'algorithmes d'apprentissage par renforcement grâce à leur compatibilité avec de nombreux outils et environnements. Cependant, elles n'incluent pas de modèles de dynamique des fluides appropriés pour simuler fidèlement le comportement des dirigeables. Il n'existe donc pas de solution facile et accessible pour l'entraînement de politiques de contrôle de dirigeables.

1.2.4 Déploiement d'algorithmes d'apprentissage par renforcement

Pour maximiser l'utilité de la plateforme, celle-ci doit faciliter le développement, l'entraînement et la validation de politiques d'apprentissage par renforcement. L'infrastructure doit donc permettre l'expérimentation et la comparaison de différentes méthodes, en plus d'assurer une transition entre la simulation et la réalité. Pour ce faire, différents outils et bibliothèques existent déjà, et il sera donc essentiel que la plateforme puisse tirer profit des outils déjà existants.

1.2.5 Apprentissage par renforcement sur des systèmes réels

Tous les simulateurs présentent ce que l'on appelle un *reality gap*, un écart entre la simulation et la réalité. Cet écart constitue un obstacle majeur à l'application de l'apprentissage par renforcement entraîné en simulation sur des systèmes réels.

En raison des contraintes, notamment en matière de sécurité, très peu d'études sur l'apprentissage par renforcement appliquée aux dirigeables testent leurs méthodes sur des systèmes réels : la majorité des travaux restent confinés à la simulation.

1.3 Objectifs de recherche

Le travail présenté dans ce mémoire a les objectifs suivants :

1. Concevoir et construire un robot de type dirigeable pouvant être utilisé à l'intérieur d'un laboratoire afin de pouvoir tester des algorithmes d'apprentissage par renforcement.
2. Développer un modèle mathématique pour ce robot de type dirigeable.
3. Créer un simulateur à partir du modèle mathématique permettant de tester des algorithmes d'apprentissage par renforcement sur ce robot.

4. Développer des politiques d'apprentissage par renforcement en simulation pour ce robot.
5. Valider la possibilité de déployer des politiques d'apprentissage par renforcement sur le système physique.

1.4 Plan du mémoire

La suite de ce mémoire est divisée en quatre chapitres. Le chapitre 2 présente une revue de la littérature sur les dirigeables, l'apprentissage par renforcement et les simulateurs. Le chapitre 3 présente le design mécanique du robot de type dirigeable, le modèle mathématique du robot et le simulateur développé. Le chapitre 4 présente les politiques d'apprentissage par renforcement développées et les résultats obtenus. Enfin, le chapitre 5 conclut le mémoire, souligne les limitations du travail de recherche et présente des pistes pour des travaux futurs.

CHAPITRE 2 REVUE DE LITTÉRATURE

Ce chapitre présente un aperçu de la recherche dans les différents thèmes essentiels à ce mémoire. Il commence par un aperçu des dirigeables existants, présente ensuite les enjeux de la modélisation dynamique des dirigeables, un aperçu de la simulation pour apprentissage par renforcement et présente finalement les travaux sur l'apprentissage par renforcement appliqués aux dirigeables.

2.1 Dirigeables de petite et moyenne taille

Dans cette revue de littérature, l'accent sera mis sur les dirigeables sans équipage (*unmanned*) de petite et moyenne taille, typiquement de l'ordre de quelques mètres jusqu'à une vingtaine de mètres. En effet, les dirigeables de grande taille utilisés pour le transport de marchandises ou de personnes, comme le LCA60T de Flying Whales, ont des systèmes complexes et inadaptés à la conception de dirigeables intérieurs. Le projet présenté exige un nombre limité d'actionneurs afin de respecter les contraintes de poids.

La mission AURORA (*Autonomous Unmanned Remote Monitoring Robotic Airship*) est une référence dans la recherche moderne sur les dirigeables autonomes. Le dirigeable utilisé pour les différentes étapes de cette mission a évolué avec le temps, mais la base des systèmes d'actionnement utilisés est restée la même. Le dirigeable initialement utilisé, le AS800, mesure 9 mètres de long par 2,25 mètres de diamètre. Il est équipé de deux propulseurs orientables de chaque côté de la gondole, avec un angle d'inclinaison allant initialement de -30° à +90° [8], puis étendu de -30° à +120° [9]. Les deux principaux propulseurs offrent une poussée vectorielle permettant à la fois le déplacement en mode croisière et la compensation du poids à basse vitesse. Il dispose de quatre gouvernes arrières disposées en forme de « X », chacune ayant une déflexion de -25° à +25° permettant des commandes combinées de profondeur, de direction et de roulis. Cette forme permet à la fois des commandes équivalentes aux gouvernes de profondeur et de direction de la forme classique en croix, et des commandes équivalentes à un aileron avec la rotation des quatre gouvernes dans la même direction [9]. Un troisième propulseur, facultatif, se situe à la proue du dirigeable pour assurer un contrôle de précision à basse vitesse. Un second dirigeable a ensuite été créé par l'équipe de AURORA, très semblable au AS800, mais mesurant 10,5 mètres de long et ayant un diamètre de 3 mètres. Sa principale innovation résidait dans l'ajout d'une possible commande différentielle des moteurs principaux [3, 10].

Plusieurs dirigeables adoptent un système d'actionnement similaire à ceux du projet AURORA, avec quelques variations. Par exemple, les dirigeables de [11, 12], conçus pour une utilisation en intérieur, sont plus petits ($\sim 2\text{--}3\text{ m}$) et disposent de deux propulseurs principaux inclinables à 180° . Ils possèdent un empennage, mais pas de gouvernes, et intègrent un moteur supplémentaire à l'arrière. Les dirigeable de [13] et du projet DIVA [14] ont une taille et un système de propulsion semblables à ceux d'AURORA, mais leurs deux propulseurs principaux sont différentiables. Les quatre gouvernes arrières de [13] sont disposées en croix, alors que celles de [14] sont en forme de « X ».

Le projet InSAC pour l'acquisition de données environnementale a utilisé deux dirigeables souples. Le NOAMAY [15], possède des gouvernes en « X » et quatre propulseurs orientables montés sur une structure attachée sous l'enveloppe. Deux propulseurs sont situés à l'avant du centre de gravité du dirigeable et deux à l'arrière. L'enveloppe est identique à celle utilisée pour le projet AURORA. Plus récemment, le projet a utilisé le Naomini [4]. Il est équipé de six propulseurs orientables fixés directement sur l'enveloppe. Chaque propulseur peut être incliné sur 360° . Il possède également trois gouvernes arrières. Sa conception permet d'utiliser deux, quatre ou six propulseurs selon les besoins.

Disposer les moteurs en les fixant de chaque côté de l'enveloppe est aussi une solution utilisée pour le dirigeable ALTAV Quanser MkII [16]. Ce modèle de 4,5 mètres dépourvu d'empennage possède quatre propulseurs orientables fixés sur l'enveloppe, deux de chaque côté. Ces propulseurs peuvent produire une poussée vers le haut, vers l'avant et vers l'arrière. Le dirigeable créé par [17] quant à lui possède quatre moteurs à orientation fixe placés de chaque côté de son enveloppe. Conçu pour une utilisation en intérieur, il mesure 2 mètres de long et ne possède pas d'empennage.

Les dirigeables classiques de petite et moyenne taille sont donc majoritairement souples et ont généralement un nombre limité de moteurs, souvent entre deux et quatre. Ils sont le plus souvent sous-actionnés et ne permettent presque jamais de contrôle latéral.

D'autres dirigeables existent avec des formes non conventionnelles [18, 19], mais ne sont pas détaillés ici, car leur aérodynamique est différente des dirigeables traditionnels.

2.2 Modélisation dynamiques des dirigeables

La modélisation dynamique des dirigeables est nettement plus complexe que celle des avions. En effet, la modélisation des dirigeables doit intégrer la flottabilité, les effets d'inertie de l'air environnant (masse ajoutée), ainsi que les forces aérodynamiques dont le calcul est encore un important sujet de recherche. En plus de dynamiques plus complexes, les dirigeables

ne bénéficient pas de modèles aussi bien testés et largement validés par des décennies de données expérimentales, entre autres en raison du déclin de leur utilisation après les années 1930 [20]. Le principal défi à la modélisation des dirigeables réside dans la détermination des coefficients aérodynamiques. Différentes approches sont proposées dans la littérature. La plus populaire, principalement car la seule accessible à l'âge d'or des dirigeables, est la détermination des coefficients en soufflerie [5]. Les résultats en soufflerie peuvent ensuite être extrapolés pour créer des modèles plus complets de dirigeables de dimensions similaires [10,21]. Un autre approche est d'utiliser des formules semi-empiriques ou analytiques, dérivées de la mécanique des fluides appliquée à des corps allongés, comme celles proposées par [22,23]. Plus récemment, des travaux explorent l'utilisation de la mécanique des fluides numérique (ou Computational Fluid Dynamics, CFD) pour la détermination des effets aérodynamiques sur les dirigeables [24, 25]. La CFD a l'avantage de pouvoir aussi être utilisée pour optimiser la forme des dirigeables et analyser leur stabilité [26, 27].

Un point commun entre les différentes méthodes de détermination des effets aérodynamiques est le manque de modèle à angle d'attaque élevé [20]. Cette modélisation est pourtant essentielle pour les manoeuvres à basse vitesse et le vol stationnaire.

2.3 Simulateurs pour apprentissage par renforcement

À partir des modèles de dirigeables développés, certains travaux ont créé des simulateurs pour tester les algorithmes de contrôle. Par exemple, les projets AURORA, DIVA, DRONI et INSAC [10, 14, 15, 28] ont développé et amélioré un simulateur de dirigeable en Simulink/MATLAB. Ce simulateur n'est toutefois pas disponible en logiciel libre. Price et al. [13] ont quant à eux créé un simulateur dans Gazebo. Ce simulateur en logiciel libre a été réutilisé dans différents travaux testant des algorithmes d'apprentissage par renforcement [29, 30]. C'est un des seuls simulateurs de dirigeable en logiciel libre. Il existe donc très peu de simulateurs accessibles pour les dirigeables.

L'utilisation d'un moteur physique permet de faciliter la simulation robotique. En effet, un moteur physique permet de simuler la dynamique des objets en prenant en compte les forces et les moments qui agissent sur eux. Un des simulateurs physiques les plus populaires pour tester des algorithmes d'apprentissage par renforcement en robotique est MuJoCo (*Multi-Joint dynamics with Contact*) [31]. En plus d'obtenir des performances similaires aux autres moteurs physiques disponibles [32], MuJoCo facilite grandement l'étude du contrôle par apprentissage par renforcement et son implémentation sur les modèles compatibles, car il permet l'utilisation facile des bibliothèques et *benchmarks* populaires comme Stable Baselines3 [33] et Gymnasium [34]. Développer un modèle adéquat de dirigeable sur MuJoCo faciliterait

donc grandement la recherche en apprentissage par renforcement sur les dirigeables.

Toutefois, MuJoCo est un simulateur fait pour les robots à corps rigide et, comme la majorité des simulateurs physiques [35], n'est pas adapté pour la simulation des interactions avec les fluides nécessaire pour les dirigeables. Il est donc difficile d'utiliser MuJoCo pour simuler des dirigeables, car il ne prend pas en compte la dynamique complexe des fluides. Cependant, à l'aide de la bibliothèque `dm_control` [36], il est possible d'ajouter des forces et des moments sur les objets, ce qui permet de simuler l'aérodynamique d'un dirigeable.

2.4 Contrôle des dirigeables

Le développement moderne des dirigeables est entre autres difficile en raison du défi que représente le contrôle de ces véhicules qui sont généralement sous-actionnés, en plus d'avoir des dynamiques non linéaires, couplées et faiblement amorties [4]. De nombreux travaux ont proposé et validé des contrôleurs de type PID [8, 9, 13]. Ces contrôleurs manquent toutefois de robustesse face aux incertitudes sur les paramètres du modèle, car leurs gains sont ajustés pour des conditions spécifiques qui peuvent changer en cours de vol (ex. effets de la température ou du vent, changement de la masse, etc.).

Divers travaux ont exploré des méthodes de contrôle non linéaire telles que le *backstepping* [14, 16, 37, 38], l'inversion dynamique (NDI) [39, 40], la commande par mode glissant [41] et le séquencement de gain [3]. Cependant, ces méthodes manquent elles aussi souvent de robustesse, car elles reposent fortement sur la précision du modèle dynamique du véhicule utilisé, modèle qui est lui-même complexe à établir.

Deux approches récentes tentent de remédier à cette dépendance au modèle. Premièrement, [4, 15] se penchent sur l'utilisation du contrôle INDI (*Incremental Nonlinear Dynamic Inversion*), une méthode de contrôle incrémentale dérivée de l'inversion dynamique non linéaire. Cette méthode est centrée sur les mesures des capteurs, ce qui réduit la dépendance au modèle dynamique. Elle exige toutefois des capteurs très rapides et précis. Deuxièmement, les récents développements en apprentissage machine permettent d'envisager l'utilisation de l'apprentissage par renforcement comme une alternative pour répondre aux défis posés par le contrôle des dirigeables.

2.5 Apprentissage par renforcement appliqué aux dirigeables

L'apprentissage par renforcement a démontré d'excellentes performances pour la résolution de problèmes complexes, notamment le contrôle de robots ayant des dynamiques non linéaires.

Bien que la majorité des avancées dans le domaine aient été obtenues en simulation [42], des résultats ont aussi été obtenus sur des systèmes réels, notamment pour le contrôle de drones [43, 44] et de robots quadrupèdes [45].

Les premiers pas de la recherche sur l'apprentissage par renforcement appliquée aux dirigeables se concentraient sur des modèles basés sur les processus gaussiens et étaient appliqués à des tâches de petite dimension [46], [47]. Les dernières avancées utilisent plutôt des algorithmes d'apprentissage par renforcement profond (*Deep Reinforcement Learning*, DRL), qui permettent de s'attaquer à des problèmes de haute dimension comme le contrôle d'un dirigeable dans un environnement 3D. Par exemple, pour suivre une trajectoire dans un environnement 3D, [48] propose un contrôleur basé sur PPO qui intègre les saturations des actionneurs, tandis que [49] utilise un algorithme de Q-learning associé à un réseau de neurones pour accélérer l'entraînement du modèle.

Des approches hybrides [29, 30] utilisent un contrôleur plus traditionnel pour la stabilité (un PID et un contrôleur H_∞ , respectivement) qui est combiné avec un agent utilisant un algorithme de DRL (*Proximal Policy Optimization*, PPO) qui permet d'optimiser la performance. Les deux méthodes ont démontré de meilleurs résultats pour la poursuite d'une trajectoire que la simple utilisation du contrôleur traditionnel.

La combinaison de l'apprentissage par renforcement et de l'apprentissage par démonstration a également été explorée dans [12].

CHAPITRE 3 DÉTAILS DE LA SOLUTION

Cette section décrit le développement de RAFALE (*Robotic Airship for Flight And Learning Experiments*), un robot de type dirigeable pour utilisation à l'intérieur et le simulateur l'accompagnant.

3.1 Le prototype physique

Le prototype robotique RAFALE a été conçu et construit afin de servir de plateforme de test pour l'apprentissage par renforcement. La principale considération dans la conception de RAFALE était la nécessité que celui-ci puisse être utilisé dans des espaces restreints afin de permettre des tests à l'intérieur. Ainsi, il doit être petit et doit pouvoir rester en vol stationnaire. Un système d'empennage et de gouvernes est présent afin d'assurer une certaine stabilité et de permettre le développement d'algorithmes de contrôle plus avancés. Ces gouvernes sont essentielles pour obtenir des résultats généralisables à de plus gros dirigeables.

Le plus grand défi dans la conception de RAFALE est de réduire le poids du système tout en maintenant une charge utile suffisante pour supporter les différents composants électroniques et capteurs. La conception du système a également pris en compte la facilité d'assemblage et de démontage, afin de permettre un accès facile aux composants internes et de réduire le stress sur l'enveloppe. Finalement, l'objectif de RAFALE est aussi d'offrir une plateforme de test à coût réduit pour le laboratoire. Ainsi, la conception du système a été faite en tenant compte de la disponibilité des composants et de leur coût.

La conception électronique du système s'inspire fortement du travail de [50], et a bénéficié des conseils du Laboratoire INIT de École de technologie supérieure.

3.1.1 L'enveloppe

La sélection de l'enveloppe est à la base de toutes les autres décisions de conception du dirigeable. Cinq critères doivent être pris en compte lors de la sélection de l'enveloppe : le matériel, la forme, la taille, le volume et le ratio longueur/diamètre. Le matériel permet de déterminer la structure du dirigeable, à savoir s'il est rigide ou souple. Pour des petits robots dirigeables, la structure est généralement souple, comme c'est le cas de tous les modèles présentés en 2.1. La forme de l'enveloppe est importante, car elle influence la traînée aérodynamique. Traditionnellement, les dirigeables ont la forme d'un corps profilé axisymétrique, permettant ainsi de réduire la force de traînée. La taille de l'enveloppe détermine où le di-

rigéable peut être utilisé, tandis que le volume détermine la flottabilité et la charge utile. Dans notre cas cela signifie qu'on doit faire un compromis entre le fait que le dirigeable doit être assez petit pour assurer une manœuvrabilité en laboratoire, et avoir une flottabilité permettant de porter les composants et capteurs nécessaires. Enfin, le ratio longueur/diamètre est directement relié aux effets aérodynamiques. Un ratio plus grand signifie une traînée plus faible, mais un ratio plus petit permet d'augmenter la charge utile pour une longueur donnée. L'enveloppe sélectionnée est l'enveloppe SB-181-200 de Windreiter. C'est une enveloppe d'un matériau multicoupe en aluminium, nylon et polyéthylène. Elle est gonflée par l'arrière et se referme avec une attache en plastique. C'est une enveloppe souple avec une forme traditionnelle ellipsoïdale. Elle a une longueur de 2 mètres et le rapport entre sa longueur et son diamètre est de 4. Elle a un volume de 259,18 litres qui lui permet de soulever environ 180 grammes lorsqu'elle est remplie d'hélium. L'enveloppe a été choisie principalement pour son ratio, soit le même que le AS800 utilisé dans le projet AURORA [10] et que le YEZ-2A utilisé dans [5].

3.1.2 Le système de commande

Le système de commande, comprenant le système de propulsion et le système de gouverne, est utilisé différemment selon la vitesse du système. À basse vitesse, les gouvernes ne peuvent pas être utilisées : le système dépend alors entièrement du système de propulsion. À haute vitesse, les gouvernes peuvent être utilisées pour le contrôle des moments, et on peut donc diminuer l'utilisation de la différentiabilité des moteurs. Aucun élément du système de contrôle ne permet les déplacements latéraux, et aucun système n'est prévu pour le contrôle de tangage à basse vitesse. Le dirigeable est donc sous-actionné.

Le système de propulsion

Le système de propulsion est composé de deux moteurs de levée et deux moteurs de propulsion. Les moteurs de levée sont utilisés pour le vol stationnaire, le contrôle d'altitude et le contrôle de roulis. Les moteurs de propulsion sont utilisés pour la translation ainsi que pour le contrôle du lacet à basse vitesse. Cette configuration de moteurs permet de combiner les avantages des moteurs de levée et de propulsion, tout en évitant le poids supplémentaire d'un système de servo qui aurait permis la vectorisation des directions des moteurs. Afin que les deux moteurs de levée puissent soulever le dirigeable en limitant les moments non désirés, ils doivent être placés sous le centre de gravité. En plaçant les moteurs sur une structure fixée sous la l'enveloppe plutôt qu'en les fixant directement sur l'enveloppe, on réduit le stress causé à l'enveloppe par le poids des moteurs en plus de faciliter l'assemblage et le démontage.

L'empennage

Comme le contrôle du roulis est assuré par les moteurs de levée, il n'est pas nécessaire d'avoir 4 surfaces de contrôle. Toutefois, il est essentiel de pouvoir contrôler le tangage puisqu'aucun moteur ne permet ce contrôle. Une gouverne de direction pourrait aussi venir remplacer le contrôle par les moteurs de propulsion à haute vitesse, puisque l'utilisation de gouvernes permet la réduction de la consommation d'énergie. Ainsi, l'empennage créé est en forme de croix et composé de deux gouvernes, une de profondeur et une de direction. Le système final comprend donc quatre empennages fixes, deux gouvernes et deux servos pour assurer la rotation. Les gouvernes sont limitées à des mouvements de $[-45^\circ, +45^\circ]$. Leur conception a été inspirée par le design de la version 4 du silent-runner développé par Windreiter [51].

3.1.3 La gondole

Comme l'enveloppe n'est pas parfaitement étanche, des fuites d'hélium sont à prévoir sur des périodes prolongées. Ainsi, pour éviter le stress sur l'enveloppe et les possibles déchirures lorsque celle-ci se dégonfle, la gondole doit pouvoir facilement être retirée et remise sur l'enveloppe. Pour ce faire, la gondole est attachée à l'enveloppe à l'aide de velcro. Cela permet également d'ajuster sa position pour s'assurer que le centre de gravité se trouve bien sous le centre de volume. La gondole est également conçue pour être facilement démontable afin de permettre un accès facile aux composants internes. Elle contient les moteurs de levée et de propulsion ainsi que les composants électroniques. Comme elle contient les moteurs de levée, elle doit se trouver sous le centre de gravité pour éviter de créer un moment de tangage. La gondole est également conçue pour être aussi légère que possible, tout en étant suffisamment rigide pour supporter les moteurs et les composants électroniques. Elle est fabriquée à partir d'une surface imprimée en PLA et de tiges de carbone et de bambou. Les moteurs de propulsion sont placés juste à l'avant de la gondole, ce qui permet de réduire le moment de tangage causé par la poussée des moteurs, et de maintenir une gondole de grosseur limitée pour en réduire le poids.

3.1.4 Le centre de masse

Dans les modèles classiques de dirigeables, on présume que le centre de masse est sous le centre de volume, ce qui permet de limiter les moments créés par la force de portance et la force de gravité. Nous avons déjà établi que pour assurer un moment vertical sans création de moment de tangage, les moteurs de levée doivent se trouver sous le centre de gravité, et donc sous le centre de volume. Ajoutons que la gondole avec les moteurs et le reste des électroniques

est en un seul bloc, et qu'il y a un système d'empennage à l'arrière. Ces conditions impliquent qu'on doit avoir une masse supplémentaire à l'avant pour compenser la masse du système d'empennage. Ainsi, la batterie, qui possède une masse importante comparativement aux autres composants, est placée à l'avant du dirigeable. Il y a également la possibilité d'ajouter du lest à l'avant de l'enveloppe pour assurer un positionnement du centre de masse idéal.

3.1.5 Le système de communication

Pour pouvoir opérer le système en mode manuel à l'aide d'une radiocommande, le système de communication est composé d'un récepteur de radiocommande ELRS (ExpressLRS) qui permet de communiquer avec la radiocommande Commando 8.

Le système possède aussi un module WiFi ESP8266 qui permet au dirigeable de communiquer par wifi avec une station de base. On peut envoyer des commandes à partir de la station de base ainsi que recevoir les signaux des capteurs.

3.1.6 La conception électronique

Le système électronique est basé sur l'utilisation d'un contrôleur de vol, qui a pour mission de faciliter l'implémentation de contrôle automatique. Le contrôleur est un *PixRacer Pro*, qui est généralement utilisé avec le pilote automatique PX4. Le schéma présenté à la figure 3.1 présente la connexion entre les différentes composants électroniques. En plus des éléments mentionnés plus haut, on retrouve un régulateur de tension électronique (ou Battery Eliminator Circuit (BEC)) qui permet de réduire la tension de la batterie pour alimenter le contrôleur de vol. Il y a également un *Electronic Speed Control* (ESC), un type de variateur électronique de vitesse, qui permet de contrôler la vitesse des moteurs à partir des signaux envoyés par le contrôleur de vol.

3.1.7 Les capteurs

Les capteurs sont essentiels pour le contrôle du dirigeable. Ils fournissent des informations cruciales sur son attitude, sa vitesse et, dans certains cas, sa position, ce qui permet d'assurer un contrôle précis et stable du système.

Le PixRacer Pro est équipé de plusieurs centrales inertielles (IMU) qui combinent gyroscope, accéléromètre et magnétomètre, permettant notamment d'estimer l'orientation du dirigeable en vol. Le système est également équipé d'un GPS et d'un baromètre, qui ne peuvent toutefois pas être utilisés dans le cadre d'une utilisation intérieure.

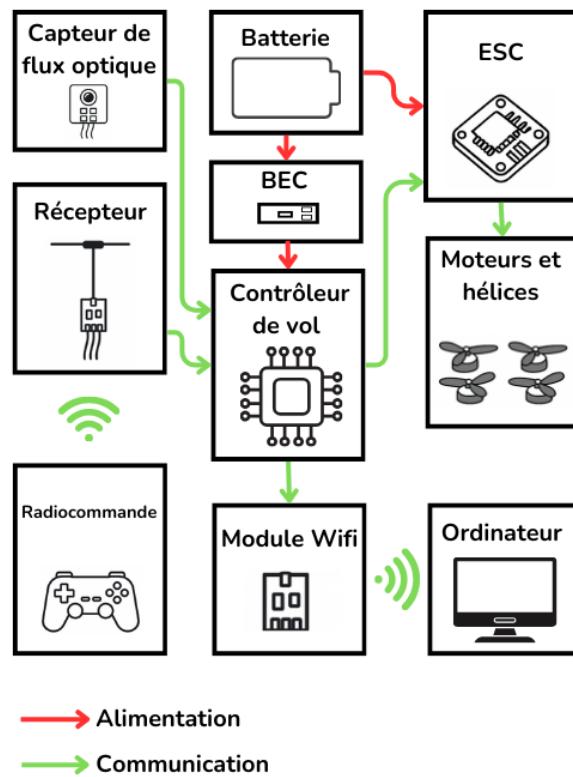


FIGURE 3.1 Électronique

Un capteur de flux optique a été ajouté afin de pouvoir mesurer la vitesse du dirigeable par rapport au sol et de déterminer sa distance avec le sol. Ce capteur est nécessaire pour le vol stationnaire, car il permet de maintenir une altitude constante. Il fonctionne en projetant un motif lumineux sur le sol et en mesurant son déplacement à l'aide d'une caméra. En analysant les variations du motif, le capteur peut estimer la vitesse du dirigeable dans les directions X et Y. Le sol doit être texturé pour identifier adéquatement la vitesse de déplacement. Le capteur de flux optique utilisé est le MTF-01, qui voit jusqu'à 8 mètres de distance avec une précision de 4 centimètres.

3.1.8 Le prototype final

La figure 3.2 présente le modèle 3D final du dirigeable, et la figure 3.3 présente une photo du prototype physique. Des pattes ont été ajoutées à la gondole pour permettre au dirigeable de se poser sur le sol sans endommager les hélices, mais elle ne se trouvent pas sur le modèle 3D.

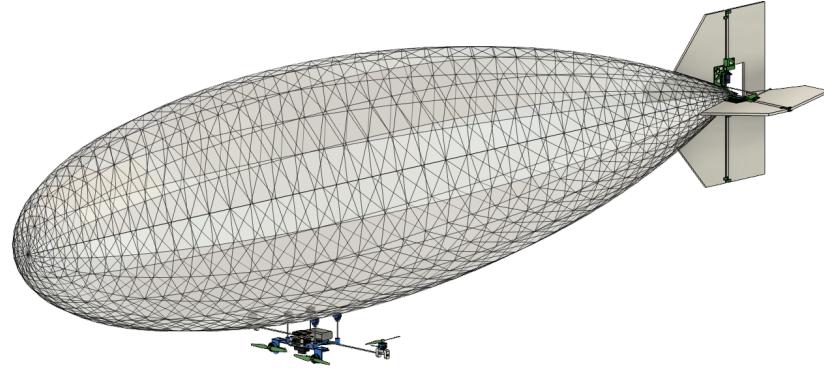


FIGURE 3.2 Design 3D de RAFALE



FIGURE 3.3 Dirigeable physique

3.2 Le modèle mathématique

Le modèle mathématique développé pour RAFALE est inspiré par les travaux de [5] et [9]. L'ensemble des variables utilisées sont définies après leur première utilisation, mais le lectorat peut toujours se référer à l'Annexe A. Le modèle peut être décrit par l'équation suivante :

$$M \frac{dx}{dt} = F_d + F_a + F_p + G \quad (3.1)$$

où

- M est la matrice de masse, contenant à la fois l'inertie réelle et l'inertie ajoutée ou virtuelle.
- x contient les vitesses linéaires et angulaires et se nomme le vecteur d'état. Il s'agit d'un vecteur de 6 éléments soit $x = [u \ v \ w \ p \ q \ r]^T$.
- F_d est le vecteur dynamique contenant les forces et moments causés par les effets de Coriolis et centrifuges.
- F_a est le vecteur des forces et moments aérodynamiques.
- F_p est le vecteur des forces et moments de propulsion.
- G est le vecteur des forces et moments dus à la flottabilité et la gravité.

Les hypothèses suivantes ont été posées pour simplifier le modèle :

- Le dirigeable est un corps rigide (on ignore les effets aéroélastiques).
- Le centre de gravité est directement en dessous du centre de volume, qui coïncide avec le centre de flottabilité.
- Le dirigeable est symétrique par rapport au plan XZ.
- La masse est constante.
- La densité de l'air et celle de l'hélium restent constantes, une approximation valable puisque le dirigeable est utilisé à l'intérieur.

3.2.1 Cadres référentiels

Trois cadres référentiels seront utilisés dans les équations, soit le référentiel local, qui est rattaché au dirigeable, le référentiel terrestre qui suivra ici la convention Nord-Ouest-Haut, et le référentiel du vent, repère dans lequel le vent est fixe. Les référentiels terrestre et local sont représentés dans l'image suivante. L'équation 3.1 est dans le référentiel local.

La figure 3.4 présente les axes de référence du dirigeable. L'origine du dirigeable est placée au centre de volume.

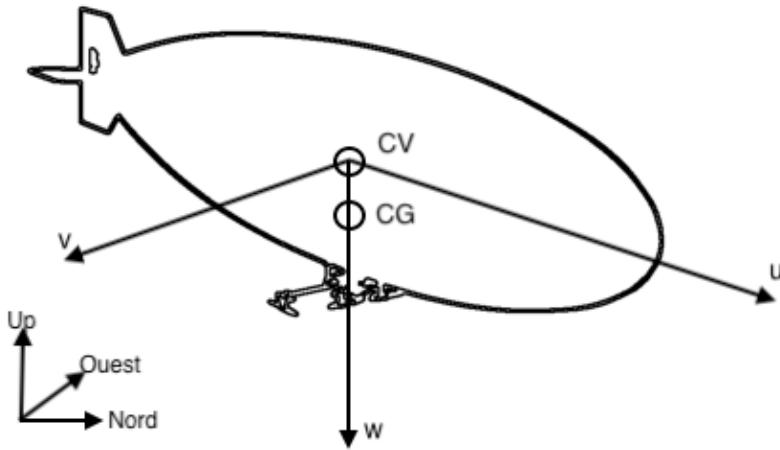


FIGURE 3.4 Axes de référence du dirigeable

3.2.2 Matrice de masses

$$M = \begin{bmatrix} m_x & 0 & 0 & 0 & ma_z - X_{\dot{q}} & 0 \\ 0 & m_y & 0 & -ma_z - Y_{\dot{p}} & 0 & ma_x - Y_{\dot{r}} \\ 0 & 0 & m_z & 0 & -ma_x - Z_{\dot{q}} & 0 \\ 0 & -ma_z - L_{\dot{v}} & 0 & I_x - L_{\dot{p}} & 0 & -J_{xz} \\ ma_z - M_{\dot{u}} & 0 & -ma_x - M_{\dot{w}} & 0 & I_y - M_{\dot{q}} & 0 \\ 0 & ma_x - N_{\dot{v}} & 0 & -J_{xz} & 0 & I_z - N_{\dot{r}} \end{bmatrix}$$

où :

- m est la masse du robot dirigeable.
- $m_x = m - X_{\dot{u}}$
- $m_y = m - Y_{\dot{v}}$
- $m_z = m - Z_{\dot{w}}$
- $X_{\dot{u}}, Y_{\dot{v}}$ et $Z_{\dot{w}}$ sont les termes de masse virtuelle pour les axes X, Y et Z.
- a_x, a_y et a_z sont les coordonnées du centre de gravité par rapport au centre de volume.
- I_x, I_y et I_z sont les termes d'inertie par rapport aux axes OX, OY et OZ.
- $L_{\dot{p}}, M_{\dot{q}}$ et $N_{\dot{r}}$ sont les termes d'inertie virtuelle par rapport aux axes OX, OY et OZ.
- $J_x = I_x - L_{\dot{p}}$
- $J_y = I_y - M_{\dot{q}}$

- $J_z = I_z - N_r$
- $J_{xz} = I_{xz} + N_{\dot{p}} = I_{xz} + L_{\dot{r}}$, où I_{xz} est le produit d'inertie par rapport à l'axe OY.
- I_{xy} et I_{yz} sont nuls car le centre de gravité et le centre de volume sont sur la même ligne.

3.2.3 Forces dynamiques

Ce vecteur contient les forces et moments causés par les effets de Coriolis et centrifuges. Il peut être représenté ainsi :

$$\mathbf{F}_d = [F_1 \ F_2 \ F_3 \ F_4 \ F_5 \ F_6]^T$$

où :

- $F_1 = -m_z w q + m_y r v + m[a_x(q^2 + r^2) - a_z r p]$
- $F_2 = -m_x u r + m_z p w + m[-a_x p q - a_z r q]$
- $F_3 = -m_y v p + m_x q u + m[-a_x r p + a_z(q^2 + p^2)]$
- $F_4 = -(J_z - J_y)r q + J_{xz}p q + m a_z(u r - p w)$
- $F_5 = -(J_x - J_z)p r + J_{xz}(r^2 - p^2) + m[a_x(v p - q u) + a_z(w q - r v)]$
- $F_6 = -(J_y - J_x)q p - J_{xz}q r + m[-a_x(u r - p w)]$

3.2.4 Forces de propulsion

Le vecteur des forces de propulsion est donné par :

$$\mathbf{F}_p = [X_{prop} \ Y_{prop} \ Z_{prop} \ L_{prop} \ M_{prop} \ N_{prop}]^T$$

Selon la configuration des moteurs de notre modèle,

- $X_{prop} = T_{ds} + T_{dp}$, poussée du moteur le long de l'axe OX
- $Y_{prop} = 0$, poussée du moteur le long de l'axe OY
- $Z_{prop} = -(T_{ls} + T_{lp})$, poussée du moteur le long de l'axe OZ
- $L_{prop} = (T_{lp} - T_{ls})d_{yl}$, moment de poussée du moteur le long de l'axe OX (roulis)
- $M_{prop} = (T_{ds} - T_{dp})d_{zd} - (T_{lp} + T_{ls})d_{xl}$, moment de poussée du moteur le long de l'axe OY (tangage)
- $N_{prop} = (T_{dp} - T_{ds})d_{yd}$, moment de poussée du moteur le long de l'axe OZ (lacet)
- T_{ds} est la poussée du moteur de propulsion tribord
- T_{dp} est la poussée du moteur de propulsion bâbord
- T_{ls} est la poussée du moteur de levage tribord

- T_{lp} est la poussée du moteur de levage bâbord
- d_x est la distance horizontale du centre de volume à l'hélice le long de l'axe principal du dirigeable
- d_y est la distance horizontale du centre de volume à l'hélice perpendiculaire à l'axe principal du dirigeable
- d_z est la distance verticale du centre de volume à l'hélice

Les forces de propulsion sont influencées par la force de l'air, mais pour simplifier le modèle, on les considère comme indépendantes [3].

3.2.5 Forces aérodynamiques

Le modèle aérodynamique du dirigeable vient d'une adaptation de [9] des données obtenues en soufflerie pour le dirigeable YEZ-2A de Westinghouse [52]. Il est possible de s'inspirer de ces modèles puisque, bien que la forme de l'enveloppe soit un peu différente, le rapport entre la longueur et le diamètre de l'enveloppe de RAFALE est de 4, tout comme le YEZ-2A et le AS800 utilisé pour AURORA.

Le modèle aérodynamique présenté inclut :

1. Les forces et moments aérodynamiques dus au déplacement du dirigeable dans l'air.
2. Les facteurs d'amortissement des forces et des moments dynamiques.

Il exclut les effets de la masse ajoutée, qui sont plutôt inclus directement dans la matrice de masse.

Forces et moments aérodynamiques

Lorsqu'on fait des mesures en soufflerie pour des dirigeables de grande taille, on utilise un modèle réduit pour mesurer les forces et moments subis par le système selon la vitesse du vent. On divise ensuite ces valeurs par des forces et des moments de référence qui permettent d'obtenir des coefficients non dimensionnels s'appliquant à tous les dirigeables ayant une forme similaire, peu importe la taille. C'est ce qui nous permet d'utiliser les résultats obtenus dans des travaux précédents pour notre modèle. Les forces de traînée, latérale et de portance ainsi que les moments aérodynamiques sont donc calculés dans le référentiel du vent relatif. La vitesse relative de l'air dans le référentiel du corps est $x_a = [u_a \ v_a \ w_a \ p_a \ q_a \ r_a]$ et tous les forces et moments aérodynamiques dépendent de la vitesse vraie V_t , soit la norme de la

vitesse relative du vent.

$$V_t = \sqrt{u_a^2 + v_a^2 + w_a^2}$$

Les coefficients mesurés en soufflerie dépendent de l'angle d'incidence (α), de l'angle de dérapage (β) et des angles de déflexion des gouvernes de profondeur (d_e) et de direction (d_r).

L'angle d'incidence correspond à :

$$\alpha = \arctan(w_a/u_a)$$

Et l'angle de dérapage est

$$\beta = \arcsin(v_a/V_t)$$

Les forces aérodynamiques calculées dans le référentiel du vent équivalent à :

$$F_a^w = \begin{bmatrix} A_X \\ A_Y \\ A_Z \\ A_L \\ A_M \\ A_N \end{bmatrix} = \begin{bmatrix} C_d F_{\text{ref}} \\ C_Y F_{\text{ref}} \\ C_l F_{\text{ref}} \\ C_L M_{\text{ref}} \\ C_M M_{\text{ref}} \\ C_N M_{\text{ref}} \end{bmatrix}$$

A_X , A_Y et A_Z sont les forces de traînée, latérale et de portance, A_L , A_M et A_N sont les moments de roulis, tangage et lacet, et C_d , C_Y , C_L , C_l , C_m et C_n sont les coefficients correspondants.

La force de référence est :

$$F_{\text{ref}} = \frac{1}{2} \rho V_t^2 S_{\text{ref}}$$

Et le moment de référence est

$$M_{\text{ref}} = \frac{1}{2} \rho V_e^2 S_{\text{ref}} L_{\text{ref}}$$

Où $S_{\text{ref}} = \text{Vol}^{2/3}$ est la surface de référence, Vol est le volume du dirigeable et ρ est la densité de l'air. Le moment de référence est exprimé en fonction de la vitesse équivalente V_e plutôt que la vitesse vraie V_t . Cette vitesse équivalente permet de mieux considérer les moments

aérodynamiques lorsque la vitesse angulaire est importante.

$$V_e = \sqrt{V_t^2 + \frac{(p_a^2 + q_a^2 + r_a^2)L_{\text{ref}}^2}{4}}$$

L_{ref} est une longueur de référence qui conventionnellement est définie comme $\text{Vol}^{1/3}$.

Les coefficients utilisés proviennent des tables de coefficients de [52] et des travaux de [21] permettant d'extrapoler les tables.

TABLEAU 3.1 Coefficients utilisés pour calculer les forces aérodynamiques

Coefficient	Définition
C_{d_0}	Coefficient de traînée à $\alpha = 0$ et $\beta = 0$
C_{d_i}	Coefficient de traînée induite (traînée causée par la présence de α et β)
C_{Y_b}	Coefficient de force latérale en fonction de β
$C_{Y_{dr}}$	Coefficient de force latérale en fonction de la déflexion de la gouverne de direction
C_{l_0}	Coefficient de portance à $\alpha = 0$
C_{l_a}	Coefficient de portance en fonction de α
$C_{l_{de}}$	Coefficient de portance en fonction de la déflexion de la gouverne de profondeur
C_{M_0}	Coefficient de tangage à $\alpha = 0$
C_{M_a}	Coefficient de tangage en fonction de α
$C_{M_{ab}}$	Coefficient de tangage en fonction de β dû à α (effet de couplage croisé)
$C_{M_{ba}}$	Coefficient de tangage en fonction de α dû à β (effet de couplage croisé)
C_{M_b}	Coefficient de tangage en fonction de β
$C_{M_{de}}$	Coefficient de tangage en fonction de la déflexion de la gouverne de profondeur
C_{N_b}	Coefficient de lacet en fonction de β
$C_{N_{dr}}$	Coefficient de lacet en fonction de la déflexion de la gouverne de direction
C_{L_b}	Coefficient de roulis en fonction de β

C_{l_0} et C_{M_0} sont nuls puisqu'on estime qu'aucune force de portance ou tangage n'est présente à un angle d'incidence de 0. À partir de ces coefficients on peut ensuite obtenir les coefficients aérodynamiques recherchés. On utilise les coefficients présentés dans le tableau 3.1 pour calculer les coefficients utilisés dans la définition de F_a^w .

$$\begin{aligned}
C_d &= C_{d_0} + C_{d_i}(C_{l_a}^2 s_1(\alpha) + C_{Y_b}^2 s_2(\beta)) \\
C_Y &= C_{Y_b} s_3(\beta) + C_{Y_{d_r}} d_r |\cos(\alpha) \cos(\beta)| \\
C_l &= C_{l_0}(1 - s_2(\beta)) + C_{l_a} s_4(\alpha)(1 + |s_3(\beta)|) + C_{l_{d_e}} d_e \cos(\alpha) \cos(\beta) \\
C_L &= C_{L_b} s_5(\beta) s_6(\alpha) \\
C_M &= C_{M_0} + (C_{M_a} + C_{M_{ab}} |\sin(\beta)|) s_7(\alpha) + (C_{M_b} + C_{M_{ba}} \sin(\alpha)) |s_3(\beta)| + C_{M_{d_e}} d_e \cos(\alpha) \cos(\beta) \\
C_N &= C_{N_b} s_8(\beta) s_6(\alpha) + C_{N_{d_r}} d_r \cos(\alpha) \cos(\beta)
\end{aligned}$$

Où $s_1 \dots s_8$ sont des courbes permettant de mettre en lien les résultats obtenus en soufflerie. Ces équations permettent aussi d'étendre les résultats obtenus de [52] et [21], à des angles d'incidence plus élevés.

Facteurs d'amortissement des forces et des moments dynamiques

Les coefficients dynamiques utilisés pour calculer les termes d'amortissement sont les dérivées des coefficients aérodynamiques selon des taux de giration non dimensionnels q^* et r^* , qui sont exprimés en fonction des vitesses angulaires q et r dans le référentiel du vent relatif.

$$\begin{aligned}
q^* &= \frac{q_a L_{\text{ref}}}{V_t} \\
r^* &= \frac{r_a L_{\text{ref}}}{V_t}
\end{aligned}$$

On utilise les coefficient dynamiques C_{Yr} , C_{nr} , C_{Zq} et C_{mq} . Ainsi la modification des forces et moments par le facteur d'amortissement correspond à :

$$F_{\text{amortissement}} = \begin{bmatrix} 0 \\ C_{Yr} F_{\text{ref}} r^* \\ C_{Zq} F_{\text{ref}} q^* \\ 0 \\ C_{mq} M_{\text{ref}} q^* \\ C_{nr} M_{\text{ref}} r^* \end{bmatrix}$$

Vecteur de forces aérodynamique

Les forces aérodynamiques ayant été calculées dans le référentiel du vent, elles doivent être multipliées par une matrice de transformation S_1 pour pouvoir être reportées dans le référentiel local et ainsi être additionnées dans l'équation 3.1.

$$F_a = S_1(F_a^w + F_{amortissement})$$

3.2.6 Forces de flottabilité et gravité

Le vecteur G représente la différence entre la force de gravité et la force de flottabilité.

$$G = \begin{bmatrix} S_2(W - B) \\ OC \times S_2 W \end{bmatrix}$$

Où W est le vecteur de force de gravité, B est le vecteur de force de flottabilité, OC est le vecteur de distance entre le centre de volume (où est appliqué B) et le centre de gravité (où est appliqué W). S_2 est la matrice de transformation permettant de transposer les résultats du référentiel terrestre au référentiel local.

3.3 Le simulateur

Le simulateur utilisé pour RAFALE est le simulateur MuJoCo [31]. Le modèle 3D présenté dans la figure 3.2 est séparé en grandes sections qui sont chacunes importées sous forme de maillage 3D. Un fichier XML permet de mettre en relation ces maillages et de déterminer les jointures et les actionneurs du système. Comme MuJoCo est un moteur physique, une fois le modèle bien défini dans le XML, il s'occupe de calculer les forces dynamiques, les forces dues à la gravité et les forces de propulsion. Pour ajouter la force de flottabilité, il existe un paramètre `gravcomp` qui permet d'ajouter une force au centre de volume correspondant à une fraction du poids de l'objet. En calculant la force de flottabilité de l'enveloppe et le poids de l'enveloppe, on peut donc trouver le paramètre `gravcomp` permettant de simuler la force de flottabilité. Comme expliqué dans 2.3, MuJoCo n'a pas de modèle adéquat de mécanique des fluides permettant de calculer les forces aérodynamiques sur le dirigeable. Celles-ci sont donc calculées dans un code python puis appliquées au modèle à l'aide de `dm_control`. La figure 3.5 permet de visualiser l'intégration des forces aérodynamiques au modèle MuJoCo.

Pour les calculs de flottabilité et d'aérodynamique, on considère la densité de l'hélium et

de l'air à 20°C et à pression normale. On assume que la température et la pression restent constantes, ce qui est raisonnable pour une utilisation intérieure. Elles sont respectivement $\rho_{\text{hélium}} = 0,1634\text{g/L}$ et $\rho_{\text{air}} = 1,225\text{g/L}$. Les valeurs de référence utilisées pour les calculs aérodynamiques sont $\text{Vol} = 0,2592\text{m}^3$, $S_{\text{ref}} = 0,4065\text{m}$ et $L_{\text{ref}} = 0,6376\text{m}$.

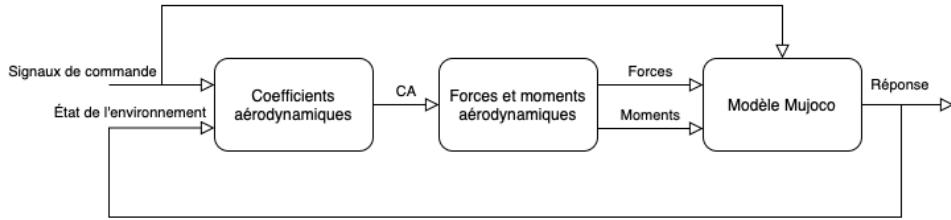


FIGURE 3.5 Intégration de l'aérodynamique avec MuJoCo

Afin d'évaluer la performance du simulateur on utilise la mesure du facteur temps réel (RTF), défini par :

$$\text{RTF} = \frac{\text{Temps simulé}}{\text{Temps d'exécution}}$$

Le RTF a été estimé en moyennant le RTF sur 1000 épisodes de 4000 pas de temps chacun. Chaque pas de temps en simulation est de 0,002 secondes. Le RTF est d'environ 12,04 pour notre simulateur, ce qui signifie que 12 secondes de déplacement du dirigeable peuvent être simulées en 1 seconde.

CHAPITRE 4 RÉSULTATS THÉORIQUES ET EXPÉRIMENTAUX

Afin d'évaluer la capacité de notre nouveau système à apprendre des politiques de contrôle, il a été entraîné à réaliser deux tâches. La première consiste à atteindre une altitude désirée puis y rester, et la deuxième demande au système de se déplacer à une vitesse et une altitude constante. Ces deux tâches permettent de bien visualiser l'effet des algorithmes d'apprentissage par renforcement sur notre système, ce qui permettra éventuellement de développer des algorithmes plus complets de suivi de trajectoires. Afin de réaliser ces tâches, deux environnements ont été créés à l'aide de Gymnasium [34].

4.1 Les environnements

Deux environnements ont été créés pour les deux tâches à accomplir. Les deux environnements partagent leur état initial, leurs conditions d'arrêt et la logique derrière la définition de l'espace d'actions. Il est à noter que les vitesses utilisées dans les deux environnements sont dans le référentiel terrestre.

État initial

L'état initial pour l'ensemble des expériences est une position de $[0, 0, 3]$ avec des angles de tangage, lacet et roulis nuls.

Conditions d'arrêt

Quatre conditions d'arrêt ont été définies pour notre système. L'épisode se termine :

1. lorsqu'il atteint le pas de temps maximal défini pour la tâche (max_{pas}) ;
2. lorsque l'angle de tangage est supérieur à 45° ;
3. lorsque l'angle de roulis est supérieur à 45° ;
4. lorsque l'altitude est inférieure ou égale à un mètre.

Les conditions 2 et 3 permettent de maintenir le système dans un état sécuritaire en prévenant des angles de roulis et de tangage trop importants, au-delà desquels le contrôle est plus difficile, voire impossible. La dernière condition est présente puisque les problèmes abordés ne sont pas des problèmes de décollage et d'atterrissage.

Lorsqu'un épisode se termine avec une des conditions 2 à 4, l'agent reçoit une pénalité de $p_{\text{échec}}$.

Espace d'actions

Pour les deux tâches sélectionnées, l'espace d'actions choisi est un espace discret, comme dans [53], ce qui a deux avantages. Premièrement, l'espace d'actions est beaucoup plus petit et donc on limite l'exploration. Ensuite, cela permet d'assurer une transition progressive entre les différentes commandes envoyées aux actionneurs, et donc d'éviter les variations brusques. L'espace d'actions est adapté selon la tâche.

4.1.1 Vol stationnaire à altitude spécifiée

L'objectif de cette tâche est d'apprendre au dirigeable à faire du vol stationnaire à une certaine altitude. Le dirigeable se déplace donc dans un premier temps à la position demandée, puis doit y rester. Pour l'entraînement, l'altitude est choisie aléatoirement dans une plage d'altitudes prédéfinie.

Espace d'actions

Puisque l'objectif de cette tâche est le vol stationnaire et le contrôle d'altitude, la commande des gouvernes n'est pas nécessaire, on ne permet donc que le contrôle des moteurs. Les actions sont celles présentées dans le tableau 4.1 :

TABLEAU 4.1 Table de correspondance des actions pour la tâche de vol stationnaire

n	Action	$\delta[T_{ls}, T_{lp}, T_{lp}, T_{lp}, dr, de]$
1	Neutre	[0, 0, 0, 0, 0, 0]
2	Accélérer (z)	[0.05, 0.05, 0, 0, 0, 0]
3	Ralentir (z)	[-0.05, -0.05, 0, 0, 0, 0]
4	Accélérer (x)	[0, 0., 0.05, 0.05, 0, 0]
5	Ralentir (x)	[0, 0, -0.05, -0.05, 0, 0]
6	Gauche (moteur)	[0, 0, -0.05, 0.05, 0.0, 0.0]
7	Droite (moteur)	[0, 0, 0.05, -0.05, 0, 0]

Récompense

La récompense pour chaque pas de temps contient 4 parties :

`recompense_survie`, qui donne une récompense fixe (r_{survie}) à chaque pas de temps où le dirigeable est encore en vie, essentiel pour s'assurer de récompenser le système pour chaque moment où il n'atteint pas une condition d'arrêt.

`punition_position`, qui soustrait à la récompense une valeur proportionnelle à la distance entre le dirigeable et la position désirée (Δ_{pos}).

`bonus_zone_cible`, qui additionne une valeur de récompense lorsque le dirigeable se trouve près de la position désirée. Cette valeur encourage le dirigeable à rester en vol stationnaire après avoir atteint la position désirée. La fonction permettant de déterminer si le dirigeable se trouve dans la zone cible est :

$$1_{\text{zone cible}} = \begin{cases} 1 & \text{si } \Delta_{pos} < \epsilon \\ 0 & \text{sinon} \end{cases}$$

où ϵ est un paramètre à définir.

`punition_orientation`, qui soustrait à la récompense une valeur proportionnelle à la différence entre les angles de tangage, de roulis et de lacet désirés et les angles réels ($\Delta\theta$).

La récompense totale est donc :

$$r_{(t)}^{\text{altitude}} = r_{\text{survie}} - w_{\text{position}}\Delta_{pos} - w_{\text{bonus zone}}1_{\text{zone cible}} - w_{\text{orientation}}\Delta\theta$$

où w_{position} , $w_{\text{bonus zone}}$, $w_{\text{orientation}}$ sont les poids attribués à chaque partie de la fonction de récompense.

Espace d'observation

L'espace d'observation utilisé par l'algorithme d'apprentissage contient les dix dernières observations, chacune contenant les éléments suivants :

- Différence entre la position désirée et la position réelle en x.
- Différence entre la position désirée et la position réelle en y.
- Différence entre la position désirée et la position réelle en z.
- Vitesse en x.
- Vitesse en y.
- Vitesse en z.
- Différence entre l'angle désiré et l'angle réel de lacet.
- Différence entre l'angle désiré et l'angle réel de tangage.
- Différence entre l'angle désiré et l'angle réel de roulis.
- Vitesse angulaire lacet.
- Vitesse angulaire tangage.
- Vitesse angulaire roulis.

En considérant les dix dernières observations plutôt que seulement la dernière, on améliore la stabilité de l'entraînement, qui peut ainsi mieux prendre en compte l'historique du système.

4.1.2 Vol de croisière à vitesse spécifiée

Cette tâche permet d'évaluer la capacité d'un algorithme d'apprentissage par renforcement à voler en mode croisière, et ainsi à utiliser les gouvernes à l'arrière. L'objectif de la tâche est d'apprendre au dirigeable à se déplacer en vol de croisière à une vitesse en x demandée (référentiel terrestre), tout en restant à l'altitude de départ.

Espace d'actions

Contrairement à la tâche de vol stationnaire, cette tâche nécessite l'utilisation des gouvernes, notamment la gouverne de profondeur pour compenser le tangage créé par les moteurs de propulsion. Les actions sont celles présentées dans le tableau 4.2 :

TABLEAU 4.2 Table de correspondance des actions pour la tâche de vol de croisière

n	Action	$\delta[\mathbf{T}_{ls}, \mathbf{T}_{lp}, \mathbf{T}_{lp}, \mathbf{T}_{lp}, dr, de]$
1	Neutre	[0, 0, 0, 0, 0, 0]
2	Accélérer (z)	[0.05, 0.05, 0, 0, 0, 0]
3	Ralentir (z)	[-0.05, -0.05, 0, 0, 0, 0]
4	Accélérer (x)	[0, 0., 0.05, 0.05, 0, 0]
5	Ralentir (x)	[0, 0, -0.05, -0.05, 0, 0]
6	Gauche (moteur)	[0, 0, -0.05, 0.05, 0.0, 0.0]
7	Droite (moteur)	[0, 0, 0.05, -0.05, 0, 0]
8	Gauche (gouverne)	[0, 0, 0, 0, 0.05, 0]
9	Droite (gouverne)	[0, 0, 0, 0, -0.05, 0]
10	Haut (gouverne)	[0, 0, 0, 0, 0, 0.05]
11	Bas (gouverne)	[0, 0, 0, 0, 0, -0.05]

Récompense

La récompense pour chaque pas contient quatre composantes :

`récompense_survie`, qui donne une récompense fixe (r_{survie}) à chaque pas de temps où le dirigeable est encore en vie, essentiel pour s'assurer de récompenser le système pour chaque moment où il n'atteint pas une condition d'arrêt.

`punition_vitesse`, qui soustrait à la récompense une valeur proportionnelle à la différence entre la vitesse désirée du dirigeable et la vitesse réelle (Δv).

`punition_altitude`, qui soustrait à la récompense une valeur proportionnelle à la différence entre l'altitude désirée du dirigeable et l'altitude réelle (Δz).

`punition_orientation`, qui soustrait à la récompense une valeur proportionnelle à la différence entre les angles de tangage, de roulis et de lacet désirés et les angles réels ($\Delta\theta$).

La récompense totale est donc :

$$r_{(t)}^{\text{vitesse}} = r_{\text{survie}} - w_{\text{vitesse}} \Delta \mathbf{v} - w_{\text{altitude}} \Delta z - w_{\text{orientation}} \Delta \theta$$

où w_{vitesse} , w_{altitude} , $w_{\text{orientation}}$ sont les poids attribués à chaque partie de la fonction de récompense.

Espace d'observation

L'espace d'observation utilisé par l'algorithme d'apprentissage contient les 10 dernières observations, chacune contenant les éléments suivants :

- Position en x.
- Position en y.
- Position en z.
- Différence entre la vitesse désirée et la vitesse réelle en x.
- Différence entre la vitesse désirée et la vitesse réelle en y.
- Différence entre la vitesse désirée et la vitesse réelle en z.
- Différence entre l'angle désiré et l'angle réel de lacet.
- Différence entre l'angle désiré et l'angle réel de tangage.
- Différence entre l'angle désiré et l'angle réel de roulis.
- Vitesse angulaire lacet.
- Vitesse angulaire tangage.
- Vitesse angulaire roulis.

4.2 Résultats en simulation

Des modèles d'apprentissage par renforcement ont été entraînés sur les deux environnements décrits précédemment. Les modèles entraînés ont été créés à l'aide de la librairie `Stable Baseline3`, qui offre une façon rapide et facile d'entraîner des algorithmes sur des environnements personnalisés. L'approche adoptée pour les deux environnements est un peu différente. Pour le premier environnement, des modèles entraînés avec PPO, QR-DQN et A2C sont comparés. Ensuite, les résultats obtenus avec la meilleure politique de contrôle sont com-

parés avec un PID. Pour le deuxième environnement, seule la performance d'un algorithme PPO entraîné pour la tâche est évaluée et on visualise son comportement dans différentes situations.

Les algorithmes PPO, QR-DQN et A2C ont été sélectionnés en raison de leur compatibilité avec les espaces d'actions discrets. L'algorithme *Proximal Policy Optimization*, ou PPO, [54] est un algorithme populaire en raison de sa stabilité et de son efficacité. C'est un algorithme de DRL populaire pour application sur les dirigeables, comme on peut voir dans la section 2.5.

Pour chaque entraînement N_{env} copies de l'environnement sont créées puis exécutées en parallèle pour un nombre total de N_{pas} chacun. Cette méthode permet de stabiliser l'apprentissage.

4.2.1 Vol stationnaire à altitude spécifiée

Pour tester l'apprentissage par renforcement sur la tâche de vol stationnaire, trois algorithmes d'apprentissage par renforcement ont été testés, permettant ainsi de profiter de la facilité d'appliquer différents algorithmes avec **Stable Baseline3**. Les trois algorithmes sélectionnés sont PPO, A2C et QR-DQN. Ce sont les trois algorithmes les plus appropriés pour une tâche avec un espace d'actions discret.

Les paramètres utilisés pour l'entraînement sont :

TABLEAU 4.3 Paramètres de l'environnement pour la tâche de vol stationnaire

Paramètre	Valeur
N_{env}	10
N_{pas}	4000000
r_{survie}	1
w_{position}	0.5
$w_{\text{bonus zone}}$	0.5
ϵ	0.05
$w_{\text{orientation}}$	$1/\pi$
$p_{\text{échec}}$	-100
Plage d'altitudes	[2, 4]
Fréquence de contrôle	50Hz
max_{pas}	600

Les résultats d'entraînement sont présentés dans la figure 4.1.

On observe que QR-DQN et A2C performent moins bien que PPO sur la tâche considérée. PPO ayant déjà été identifié comme l'algorithme de choix pour des tâches similaires [29, 30], c'est l'algorithme qui a été retenu pour la suite des expériences.

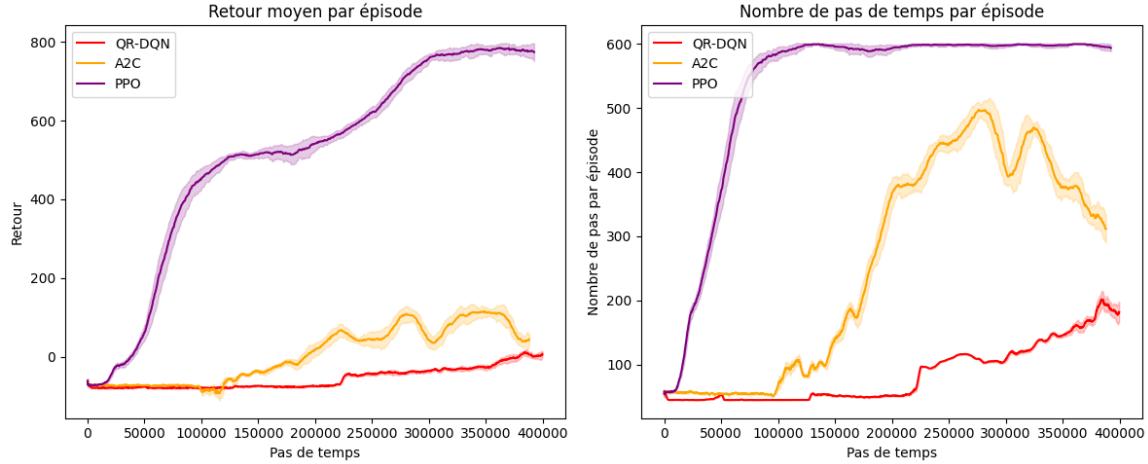


FIGURE 4.1 Courbes d’entraînement pour la tâche de vol stationnaire à altitude donnée pour trois différents algorithmes

Avec les mêmes paramètres, PPO est entraîné sur trois valeurs de *seeds* différentes afin de vérifier la robustesse de l’algorithme. Les résultats sont présentés à la figure 4.2. On constate que les trois modèles convergent vers des performances similaires, mais que celui avec la *seed* 40 converge un peu plus rapidement.

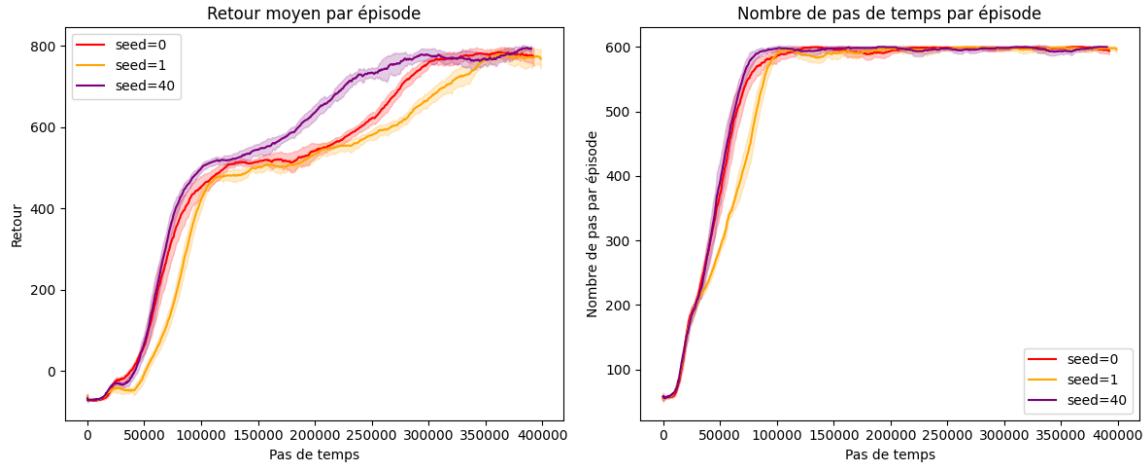


FIGURE 4.2 Courbes d’entraînement de PPO pour la tâche de vol stationnaire à altitude donnée pour trois *seeds* différentes

La figure 4.3 permet de visualiser la performance du modèle. Celui-ci est évalué sur cinq valeurs d’altitude, dont deux en dehors de la plage d’altitudes utilisée pendant l’entraînement.

Trois constats peuvent être tirés de cette évaluation. Premièrement, le modèle réussit adéquatement à maintenir sa position en x et en y, et ce, même à des valeurs d’altitudes désirées

en dehors de la plage d'entraînement. Deuxièmement, le modèle réussit à maintenir l'altitude désirée même sur un épisode plus long que ceux utilisés lors de l'entraînement (16 vs 12 secondes). Troisièmement, on observe que le modèle a des valeurs de tangage assez importantes, notamment en raison de l'utilisation des moteurs de propulsion pour garder la position en x et diminuer le lacet.

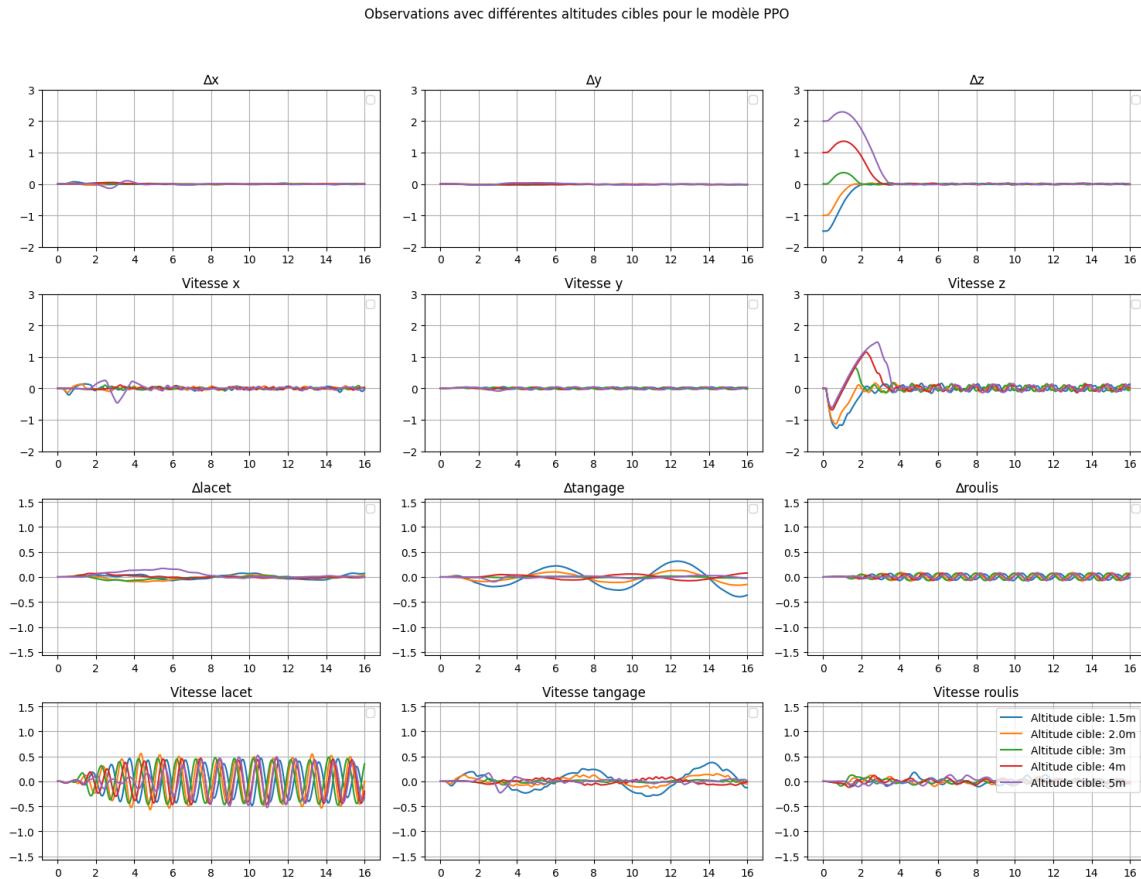


FIGURE 4.3 Observations selon le temps pour différentes altitudes

Pour évaluer la performance de l'algorithme PPO, ses résultats ont été comparés à ceux obtenus avec des contrôleurs PID ajustés pour chaque valeur d'altitude cible. Les comparaisons pour deux cas spécifiques (une altitude au dessus de la position initiale et une au dessous) sont présentées dans les figures 4.4 et 4.5.

Comparaison des résultats avec PPO vs PID pour 16 secondes, altitude désirée de 4m

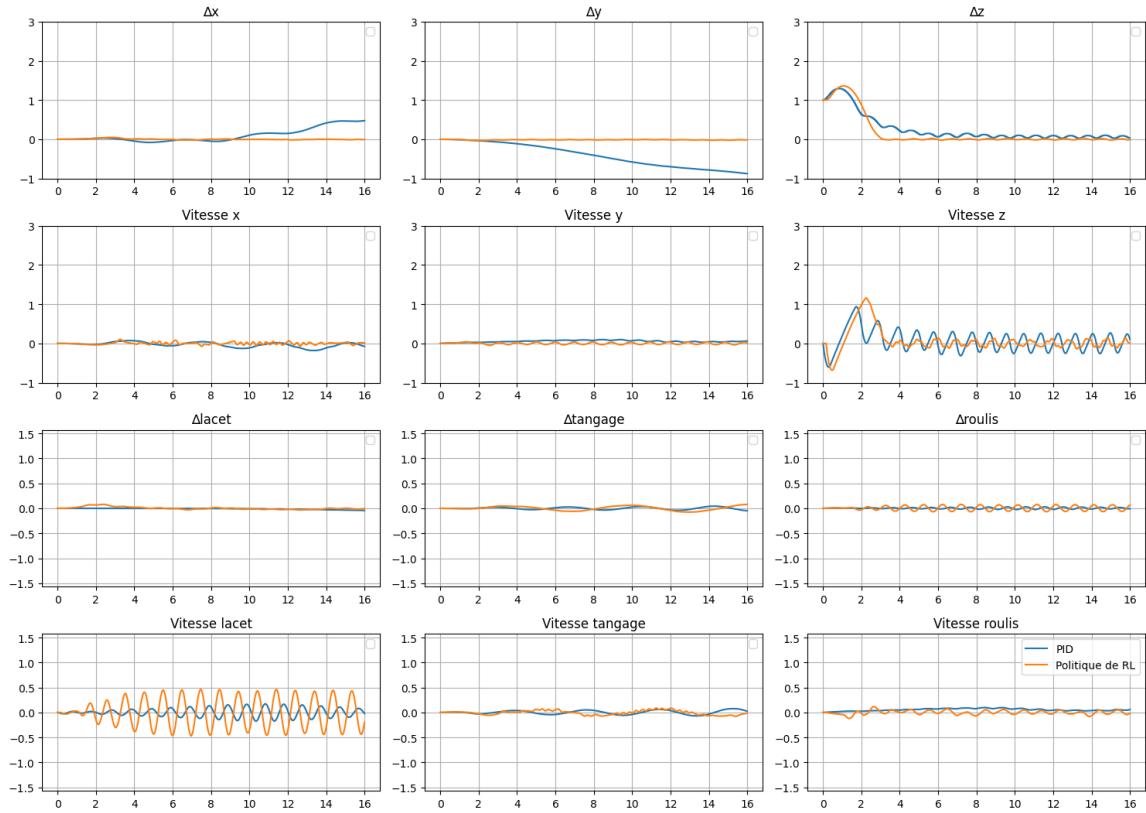


FIGURE 4.4 Comparaison des observations entre un PID et PPO pour une altitude de 4m

Comparaison des résultats avec PPO vs PID pour 16 secondes, altitude désirée de 2m

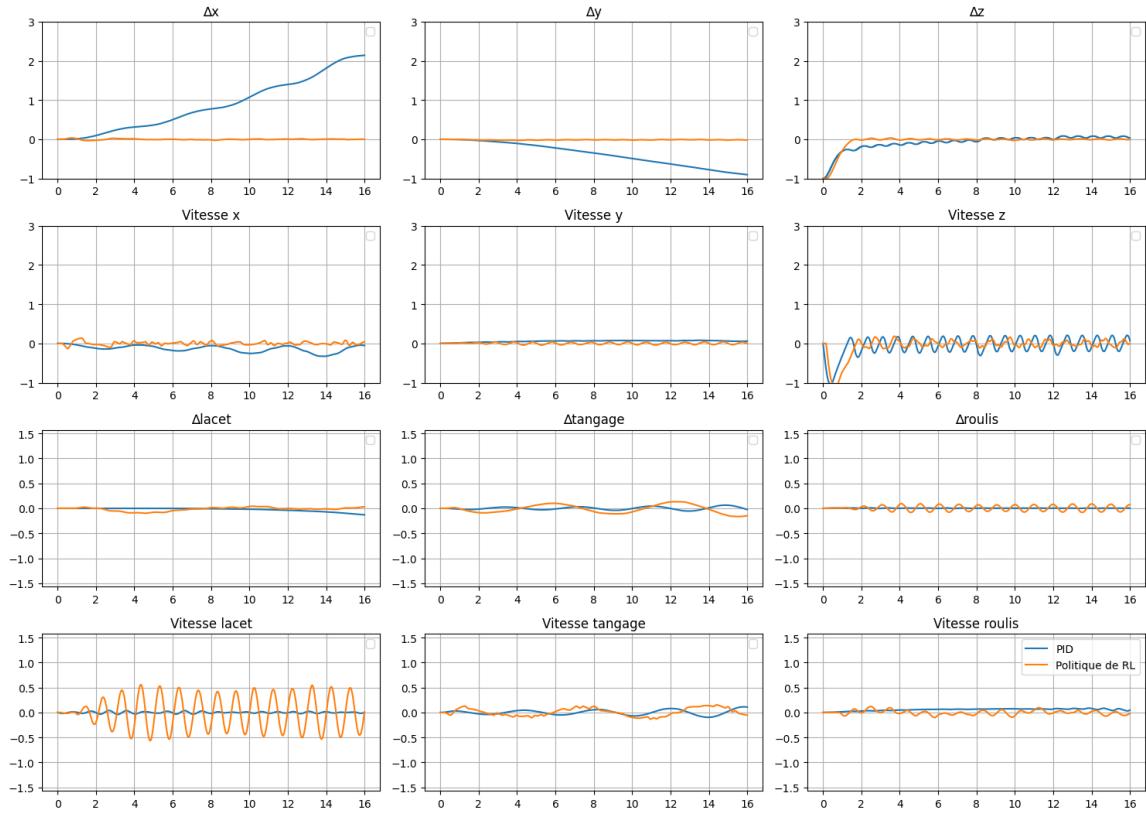


FIGURE 4.5 Comparaison des observations entre un PID et PPO pour une altitude de 2m

Dans les deux cas, on constate que l'algorithme d'apprentissage par renforcement réussit à atteindre une performance semblable à celle du PID. Dans les deux cas, l'algorithme développé avec PPO permet de converger plus rapidement à la solution que le PID, et la vitesse en z est plus stable. Puisque le PID contrôle seulement l'altitude et que des forces aérodynamiques sont en action sur le dirigeable, on peut observer que le système ne garde pas une position constante en x et en y, ce qui est normal. En intégrant le contrôle de la position en x et en y, l'algorithme PPO réussit à garder une position constante dans les deux axes. Toutefois, l'utilisation des moteurs de propulsion pour garder la position en x et diminuer le lacet induit une vitesse de lacet plus importante et crée des oscillations en roulis plus importantes que le PID. Pour le PID comme pour PPO, on observe une plus grande difficulté à garder la position désirée lorsque l'altitude désirée est plus basse que la position de départ, notamment en raison des forces aérodynamiques qui créent des moments non désirés parce que le centre de gravité est plus bas que le centre de flottabilité.

4.2.2 Vol de croisière à vitesse spécifiée

Pour cette tâche, un algorithme PPO a été entraîné sur trois *seeds* différentes.

Les paramètres utilisés pour l'entraînement sont présentés aux tableaux 4.4.

TABLEAU 4.4 Paramètres de l'environnement pour la tâche de vol de croisière

Paramètre	Valeur
N_{env}	10
N_{pas}	1000000
r_{survie}	1
w_{vitesse}	0.3
w_{altitude}	0.5
$w_{\text{orientation}}$	$1/\pi$
$p_{\text{échec}}$	-100
Plage de vitesses	[0.1, 1.5]
Fréquence de contrôle	500Hz
max_{pas}	2000

Les modèles entraînés ont été testés sur 100 épisodes où la vitesse cible était échantillonnée aléatoirement entre 0,1 et 1,5 m/s. Comme les trois modèles présentaient des performances similaires, un seul a été sélectionné aléatoirement. Les résultats de l'entraînement de ce modèle sont présentés à la figure 4.6.

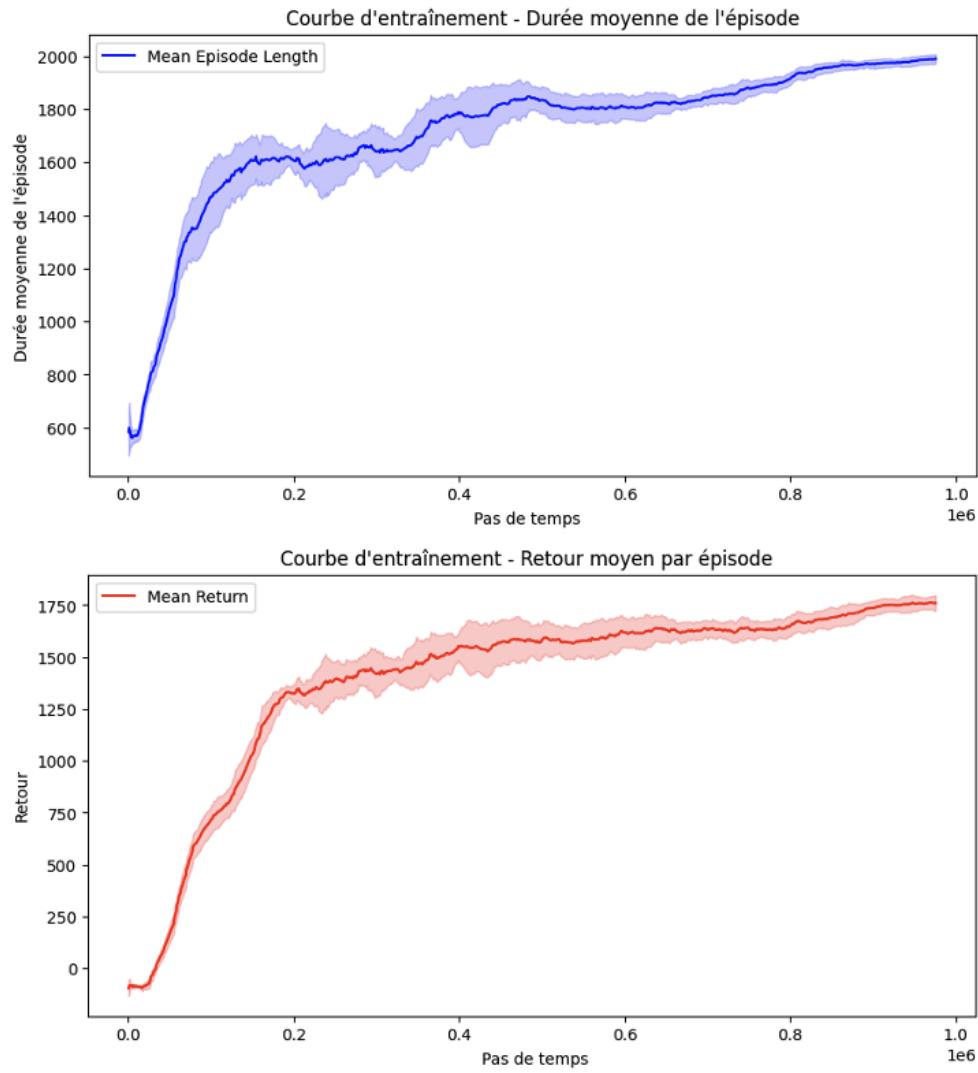


FIGURE 4.6 Résultats de l'entraînement pour la tâche de vol en mode croisière

Le modèle obtenu a été testé pour trois vitesses cibles différentes, et les résultats sont présentés dans les figures 4.7, 4.8 et 4.9. L'examen de ces trois figures nous permet de constater que le système réussit à maintenir la vitesse et l'altitude demandées pour les trois vitesses explorées. Le nombre de pas maximal par épisode pour l'entraînement était de 2000, mais l'évaluation est faite sur des épisodes de 4000 pas pour vérifier si les résultats se généralisent. Pour les petites vitesses ($0,1 \text{ m/s}$ et $0,5 \text{ m/s}$) on obtient de bons résultats avec des vitesses en x et en z somme toute constantes et une certaine compensation des moments de tangage. Toutefois, à plus haute vitesse (1 m/s), les moments créés sont trop grands, en plus d'être exacerbés par la plus forte présence des effets aérodynamiques à grande vitesse, et le système ne réussit pas à se rendre à 4000 pas, et termine abruptement à 2676 pas en raison d'une violation à la condition sur l'angle maximal de tangage. Un entraînement sur des épisodes plus longs serait donc nécessaire pour apprendre au système à compenser ces importants moments.

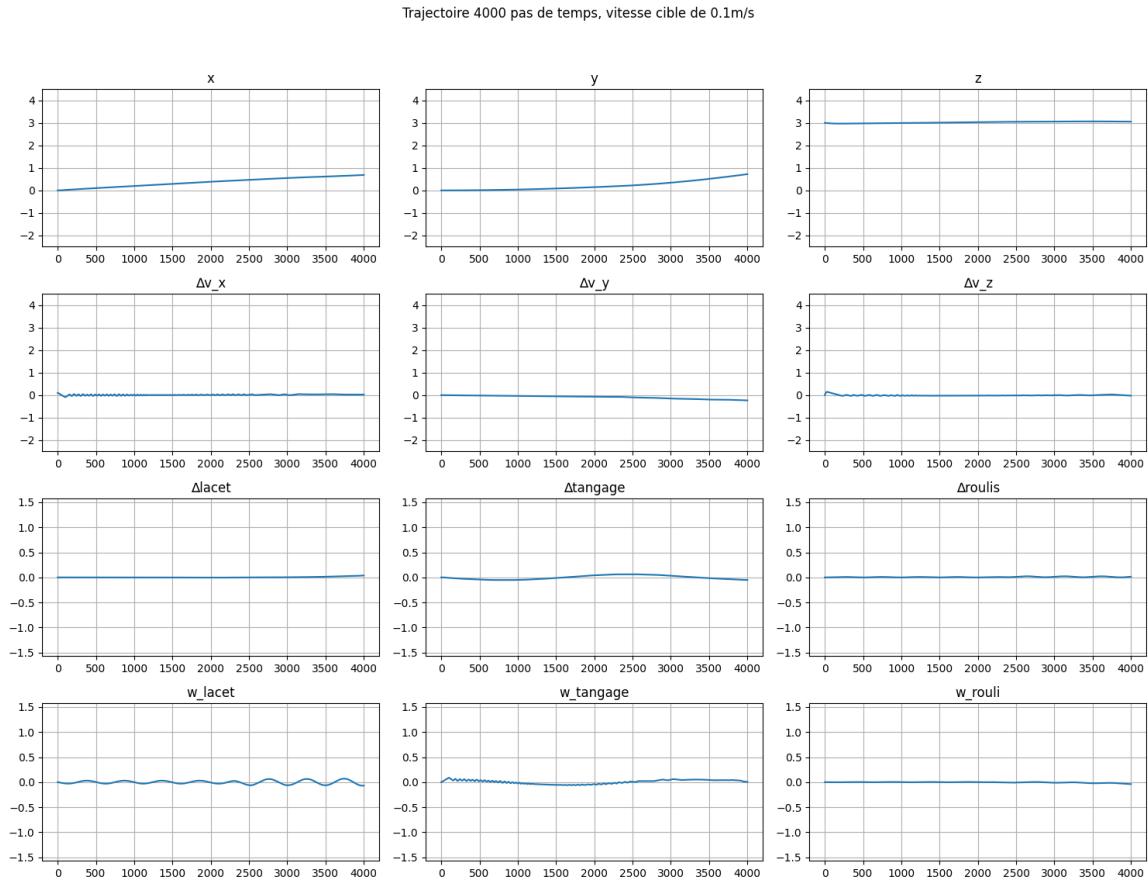


FIGURE 4.7 Visualisation des observations générées par le modèle entraîné : 4000 pas de temps et vitesse cible de $0,1 \text{ m/s}$

Trajectoire 4000 pas de temps, vitesse cible de 0.5m/s

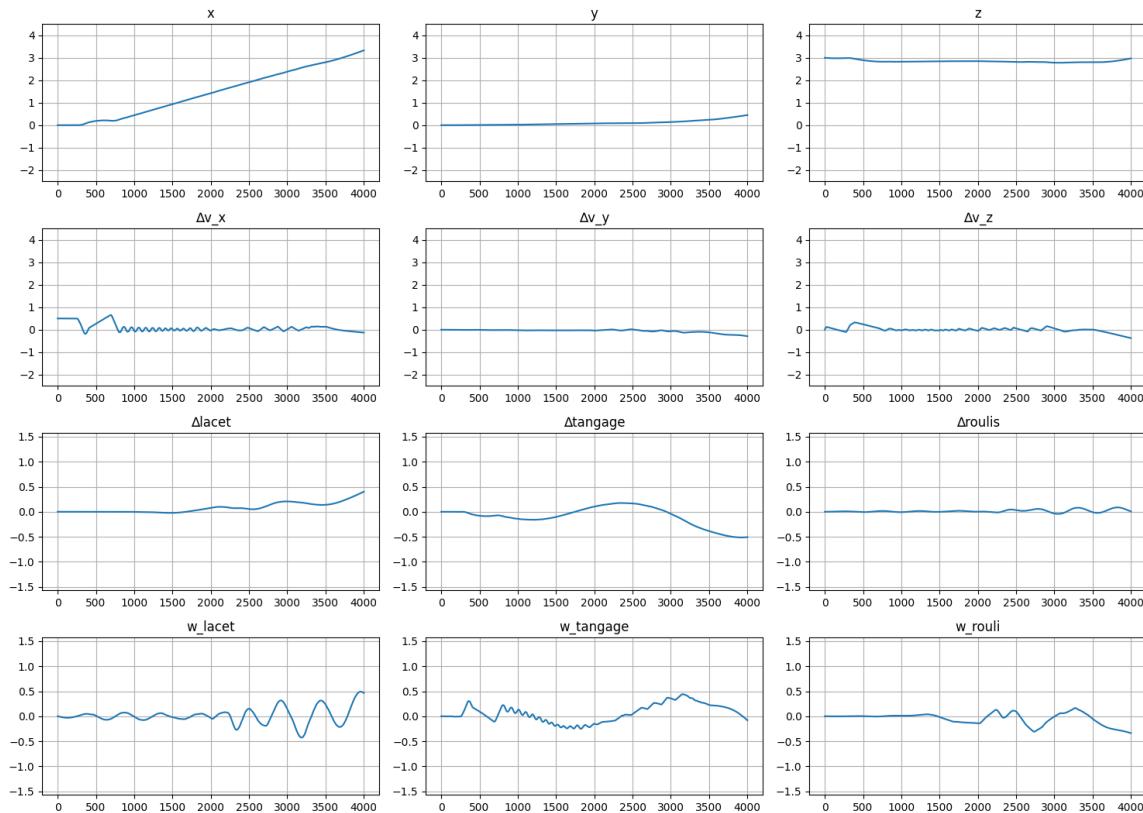


FIGURE 4.8 Visualisation des observations générées par le modèle entraîné : 4000 pas de temps et vitesse cible de 0,5 m/s

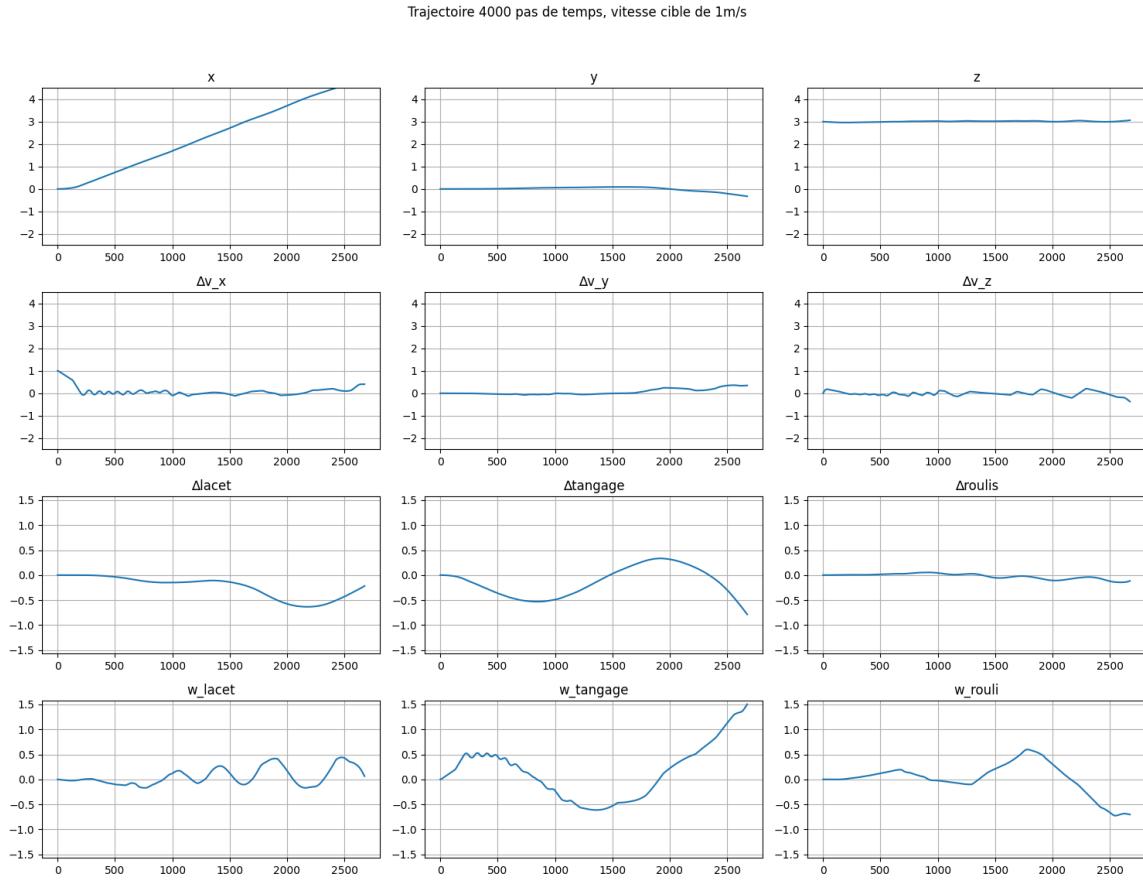


FIGURE 4.9 Visualisation des observations générées par le modèle entraîné : 4000 pas de temps et vitesse cible de 1 m/s

Une des raisons pour lesquelles on entraîne notre système à faire cette tâche est de valider que l'apprentissage par renforcement peut utiliser adéquatement les gouvernes à l'arrière et comprendre leur utilité en vol croisière. La figure 4.10 montre que la gouverne de profondeur est bien utilisée pour compenser le moment de tangage.

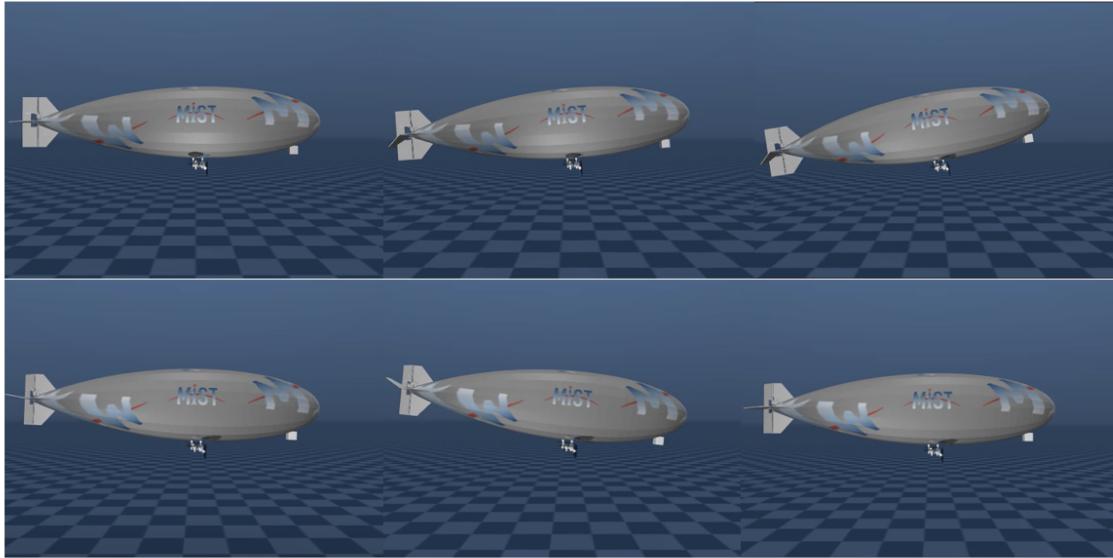


FIGURE 4.10 Mouvements de la gouverne de profondeur du dirigeable selon l’angle de tangage

4.3 Validation du système pour déploiement sur le système physique

Puisque l’apprentissage par renforcement a été validé en simulation pour différentes tâches, le dernier objectif est de démontrer que la plateforme développée permet le déploiement de politiques sur le dirigeable physique. Pour ce faire, un modèle simplifié a été créé afin de valider la plateforme sur un problème unidimensionnel. Le modèle simplifié considère uniquement le dirigeable équipé de deux moteurs orientés vers le haut, permettant un déplacement vertical. L’environnement de vol stationnaire à altitude spécifiée a été modifié pour ne considérer que la position en z , en retirant les moteurs de propulsion horizontaux. Le modèle devient ainsi un environnement de vol stationnaire en une dimension.

La fonction de récompense a également été adaptée : la pénalité sur l’orientation est retirée et la pénalité sur la position est modifiée pour ne considérer que l’erreur d’altitude (z). L’espace d’observation contient la différence entre l’altitude désirée et l’altitude réelle, ainsi que la vitesse verticale pour les 10 dernières observations.

Puisque le capteur d’altitude du dirigeable présente une incertitude non négligeable, un bruit aléatoire a été ajouté à l’observation de l’erreur d’altitude pour permettre une meilleure généralisation du modèle. Finalement, la fréquence de contrôle a été réduite pour mieux correspondre à la fréquence à laquelle le système physique peut recevoir les observations. Les paramètres de l’environnement sont présentés dans le tableau 4.5.

La politique d’apprentissage par renforcement est déployée sur le dirigeable réel à l’aide

TABLEAU 4.5 Paramètres de l'environnement pour contrôle en une dimension

Paramètre	Valeur
N_{env}	10
N_{pas}	1000000
r_{survie}	1
w_{position}	0.5
$w_{\text{bonus zone}}$	0.5
ϵ	0.05
$p_{\text{échec}}$	-100
Plage d'altitudes	[2, 4]
Fréquence de contrôle	25Hz
max_{pas}	200

d'une architecture fonctionnant ainsi : un module communiquant avec le système MAVLink contenant l'information des capteurs via MAVSDK permet de recueillir les observations à une fréquence de 50Hz. Ces informations sont ensuite transmises à la politique préalablement entraînée, qui génère à son tour une action à appliquer au dirigeable. Une action est envoyée à tous les 50 Hz.

Les expériences sont menées sur un total de 120 secondes. Pour les 40 premières secondes, l'objectif de l'agent est de se stabiliser à une altitude de 3,6 mètres, puis l'agent doit se stabiliser à une altitude de 4 mètres pour les 40 secondes suivantes. Enfin, pour les 40 dernières secondes, l'agent doit se stabiliser à une altitude de 3,6 mètres. Le système est évalué sur 10 épisodes, chacun durant un total de 120 secondes. La figure 4.11 présente les résultats obtenus sur le système physique.

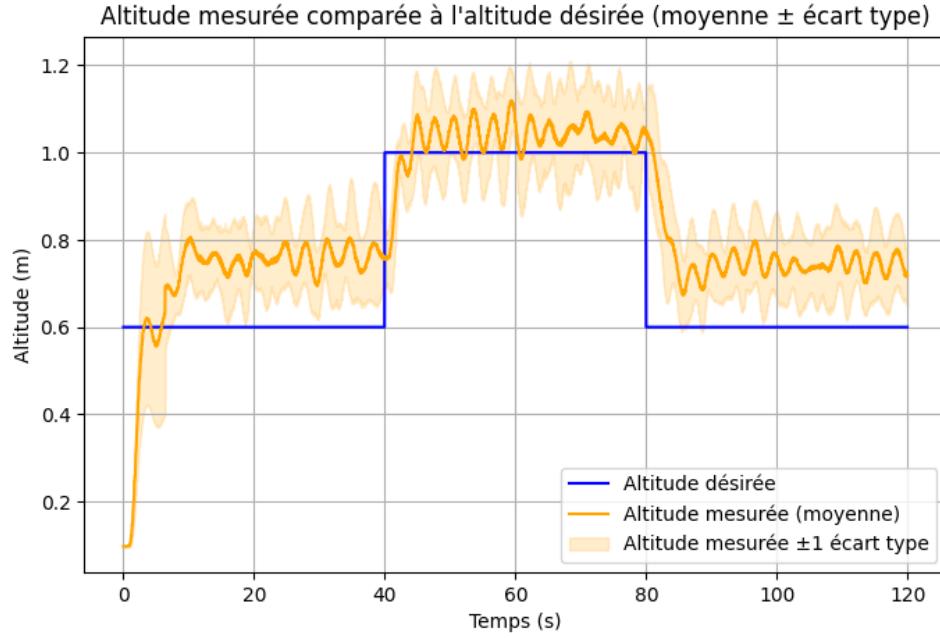


FIGURE 4.11 Résultats sur le système physique

On observe que le dirigeable parvient à se stabiliser aux alentours de l'altitude désirée, bien qu'il se maintienne significativement au-dessus lorsque l'altitude désirée est de 3,6 mètres. Cette erreur permet d'illustrer le problème de *reality gap* abordé dans la section 1.2.5.

4.4 Conclusion

Les expériences ont démontré qu'il est possible d'entraîner des modèles de contrôle par apprentissage par renforcement sur notre système. On constate en premier lieu que PPO donne de bons résultats pour l'entraînement. Une analyse plus approfondie des hyperparamètres de l'algorithme serait pertinente pour des tâches plus complexes.

On observe également que le contrôle de l'altitude s'avère plus difficile lorsque le dirigeable doit descendre. De plus, le système rencontre des difficultés à fonctionner à des vitesses supérieures à 1 m/s, en raison du moment de tangage induit par les moteurs de propulsion. Il est donc nécessaire de limiter la vitesse à une plage plus basse.

Enfin les deux algorithmes évalués permettent une généralisation satisfaisante en fonction du nombre de pas de temps, et permettent d'atteindre des cibles situées en dehors des plages d'entraînement.

Une politique a été entraînée sur une simplification du simulateur du dirigeable afin de pouvoir faire une validation initiale de notre plateforme sur le système physique pour un

problème unidimensionnel. Le dirigeable a été capable de se stabiliser autour d'altitudes données, démontrant ainsi la capacité de notre plateforme à déployer des politiques sur le dirigeable réel.

CHAPITRE 5 CONCLUSION

Les travaux présentés constituent un premier pas vers la recherche sur le contrôle des dirigeables pour le laboratoire. Elle a été créée comme outil pour les travaux futurs. On y présente la conception d'un robot dirigeable pour utilisation intérieure et un simulateur incluant les effets aérodynamiques. Des algorithmes d'apprentissage par renforcement sont ensuite testés sur deux tâches différentes. Finalement, une politique de contrôle pour un problème unidimensionnel est déployée sur le système physique afin de valider l'approche. Le robot sera utilisé dans des travaux futurs pour appliquer des algorithmes d'apprentissage par renforcement sécuritaires sur un dirigeable dans le monde réel.

5.1 Synthèse des travaux

Le résultat des travaux est le robot dirigeable RAFALE présenté dans le chapitre 3 et le simulateur l'accompagnant. C'est un dirigeable de petite taille pour utilisation à l'intérieur qui peut à la fois voler en mode croisière et en mode stationnaire. Il possède quatre moteurs et deux gouvernes permettant de le contrôler. Un simulateur a été développé à l'aide de MuJoCo afin de permettre d'entraîner des politiques d'apprentissage par renforcement pour le contrôle du système. Un modèle aérodynamique a été développé et intégré avec le moteur physique de MuJoCo pour simuler adéquatement les mouvements du robot dirigeable. Des politiques de contrôles apprises par renforcement ont été développées et appliquées au dirigeable en simulation, démontrant la possibilité d'utiliser l'apprentissage par renforcement pour le contrôle du dirigeable, à la fois pour le vol stationnaire et pour le vol en croisière. Finalement, une politique de contrôle a été déployée sur le système physique, démontrant ainsi la capacité de la plateforme développée à être utilisée pour le développement et la validation d'algorithmes d'apprentissage par renforcement.

5.2 Limitations de la solution proposée

Le simulateur du dirigeable, bien que basé sur des données de dirigeables semblables, n'a pas été validé avec le système réel. L'acquisition de données de vol du prototype permettrait d'en vérifier la justesse. De plus, différentes hypothèses ont été utilisées pour la modélisation du dirigeable, comme la rigidité du dirigeable. Ces hypothèses, bien que courantes dans le milieu, simplifient des dynamiques qui sont en réalité plus complexes et le modèle peut donc différer du mouvement du dirigeable dans le monde réel. Enfin, bien qu'une politique d'apprentissage

par renforcement ait été déployée avec succès sur le système physique, celle-ci répondait à un problème unidimensionnel. Le déploiement de politiques apprises sur des problèmes plus complexes est donc encore à valider.

5.3 Améliorations futures

Les travaux présentés dans ce mémoire sont la base des travaux qui seront menés dans le cadre du projet de recherche ENVIA Québec. La prochaine étape sera certainement de tester sur le robot dirigeable des algorithmes entraînés pour des tâches plus complexes, et peut-être même d'entraîner des algorithmes d'apprentissage par renforcement directement sur le système réel. En effet, comme le dirigeable est fait d'une enveloppe souple, le risque d'accidents graves est grandement réduit, ce qui fait de ce prototype une plateforme de choix pour apprendre directement dans le monde réel. L'objectif à plus long terme est de développer des algorithmes d'apprentissage par renforcement permettant d'obtenir des garanties de sécurité sur le vol du dirigeable. Une des avenues explorées par le projet est l'utilisation de *Control Barrier Functions* (CBF), qui permet d'offrir une certaine garantie sur les actions du système [55]. Pour se faire, un module de simulation de vent pourrait être intégré au simulateur. Également, comme mentionné dans les limitations, le simulateur pourrait prochainement être validé à l'aide de données de vol qui permettraient d'affiner le modèle.

RÉFÉRENCES

- [1] United Nations Environment Programme, “Emissions Gap Report 2024 : No more hot air ... please! With a massive gap between rhetoric and reality, countries draft new climate commitments,” Nairobi, Kenya, 2024. [En ligne]. Disponible : <https://doi.org/10.59117/20.500.11822/46404>
- [2] B. E. Prentice, “Intermodal Competition : Cargo Airships versus Long-Haul Trucking for Perishable Commodities,” *Journal of Transportation Technologies*, vol. 14, n°. 02, p. 195–211, 2024. [En ligne]. Disponible : <https://doi.org/10.4236/jtts.2024.142012>
- [3] A. Moutinho *et al.*, “Airship robust path-tracking : A tutorial on airship modelling and gain-scheduling control design,” *Control Engineering Practice*, vol. 50, p. 22–36, mai 2016. [En ligne]. Disponible : <https://doi.org/10.1016/j.conengprac.2016.02.009>
- [4] J. Azinheira *et al.*, “Hexa-Propeller Airship for Environmental Surveillance and Monitoring in Amazon Rainforest,” *Aerospace*, vol. 11, p. 249, 2024. [En ligne]. Disponible : <https://doi.org/10.3390/aerospace11040249>
- [5] S. B. V. Gomes et J. G. Ramos, “Airship dynamic modeling for autonomous operation,” communication présentée à Proceedings of the 1998 IEEE International Conferenceon Robotics and Automation, Leuven, Belgique, 20 mai 1998. [En ligne]. Disponible : <https://doi.org/10.1109/ROBOT.1998.680973>
- [6] L. M. Nicolai, G. Carichner et L. M. Nicolai, *Fundamentals of Aircraft and Airship Design*, ser. AIAA Educational Series. Reston, VA : American Institute of Aeronautics and Astronautics, 2010.
- [7] H. Lamb, “The Inertia Coefficients of an Ellipsoid Moving in Fluid,” *Advisory Committee for Aeronautics, Reports and Memoranda, No. 623*, oct. 1918.
- [8] E. C. de Paiva, S. B. V. Gomes et M. Bergerman, “A control system development environment for AURORA’s semi-autonomous robotic airship,” communication présentée à Proceedings of the 1999 IEEE Internatid Conferenceon Robotics and Automation, Detroit, MI, USA, 10-15 May 1999. [En ligne]. Disponible : <https://doi.org/10.1109/ROBOT.1999.770453>
- [9] J. R. Azinheira *et al.*, “Lateral/directional control for an autonomous, unmanned airship,” *Aircraft Engineering and Aerospace Technology*, vol. 73, n°. 5, 2001. [En ligne]. Disponible : <https://doi.org/10.1108/EUM0000000005880>
- [10] E. C. de Paiva *et al.*, “Project AURORA : Infrastructure and flight control experiments for a robotic airship,” *Journal of Field Robotics*, vol. 23, n°. 3/4, p. 201–222, 2006. [En ligne]. Disponible : <https://doi.org/10.1108/JFR2006.03.001>

- ligne]. Disponible : <https://doi.org/10.1002/rob.20111>
- [11] A. Rottmann *et al.*, “Towards an Experimental Autonomous Blimp Platform,” communication présentée à Proceedings of the 3rd European Conference on Mobile Robots, EMCR 2007, Freiburg, Germany, 2007. [En ligne]. Disponible : http://ecmr07.informatik.uni-freiburg.de/proceedings/ECMR07_0071.pdf
- [12] O. Daskiran, B. Huff et A. Dogan, “Low speed airship control using reinforcement learning and expert demonstrations,” communication présentée à AIAA 2017-0934. AIAA Atmospheric Flight Mechanics Conference, 2017. [En ligne]. Disponible : <https://arc.aiaa.org/doi/abs/10.2514/6.2017-0934>
- [13] E. Price *et al.*, “Simulation and control of deformable autonomous airships in turbulent wind,” communication présentée à Intelligent Autonomous Systems 16, ser. Lecture Notes in Networks and Systems, M. H. Ang Jr *et al.*, édit., vol. 412. Cham, Suisse : Springer International Publishing, 2022, p. 608–626. [En ligne]. Disponible : https://doi.org/10.1007/978-3-030-95892-3_46
- [14] A. Moutinho *et al.*, “Project DIVA : Guidance and Vision Surveillance Techniques for an Autonomous Airhip,” dans *Robotics Research Trends*, X. Guo, édit. Nova Science Publishers, Inc., 2008, p. 77–120.
- [15] A. S. Marton *et al.*, “Filtering and Estimation of State and Wind Disturbances Aiming Airship Control and Guidance,” *Aerospace*, vol. 9, p. 470, 2022. [En ligne]. Disponible : <https://doi.org/10.3390/aerospace9090470>
- [16] T. Liesk, M. Nahon et B. Boulet, “Design and experimental validation of a controller suite for an autonomous, finless airship,” communication présentée à 2012 American Control Conference (ACC). IEEE, juin 2012, p. 2491–2496. [En ligne]. Disponible : <http://ieeexplore.ieee.org/document/6315218/>
- [17] G. Aman, “Indoor Blimp Control,” mémoire de maîtrise, Lund University, 2021. [En ligne]. Disponible : <http://lup.lub.lu.se/student-papers/record/9061700>
- [18] P. González *et al.*, “Developing a Low-Cost Autonomous Indoor Blimp,” *Journal of Physical Agents (JoPha)*, vol. 3, n°. 1, p. 43–52, 2009. [En ligne]. Disponible : <https://doi.org/10.14198/jopha.2009.3.1.06>
- [19] J. E. Salas Gordoniz, N. Reeves et D. St-Onge, “Modular foldable airship concept for subterranean exploration,” communication présentée à International Design Engineering Technical Conferences and Computers and Information in Engineering Conference, vol. Volume 8B : 45th Mechanisms and Robotics Conference (MR), 08 2021. [En ligne]. Disponible : <https://doi.org/10.1115/DETC2021-69954>

- [20] Y. Li, M. Nahon et I. Sharf, “Airship dynamics modeling : A literature review,” *Progress in Aerospace Sciences*, vol. 47, n°. 3, p. 217–239, avr. 2011. [En ligne]. Disponible : <https://www.sciencedirect.com/science/article/pii/S0376042110000618>
- [21] J. L. Mendonça Junior *et al.*, “Airship Aerodynamic Coefficients Estimation Based on Computational Method for Preliminary Design,” communication présentée à AIAA Aviation 2019 Forum. Dallas, Texas : American Institute of Aeronautics and Astronautics, juin 2019. [En ligne]. Disponible : <https://arc.aiaa.org/doi/abs/10.2514/6.2019-2982>
- [22] S. P. Jones et J. D. DeLaurier, “Aerodynamic estimation techniques for aerostats and airships,” *Journal of Aircraft*, vol. 20, n°. 2, p. 120–126, 1983. [En ligne]. Disponible : <https://doi.org/10.2514/3.44840>
- [23] G. A. Khouri, édit., *Airship Technology*, 2^e éd., ser. Cambridge Aerospace Series. Cambridge ; New York : Cambridge University Press, 2012, n°. 10.
- [24] T. Lutz *et al.*, “Aerodynamic Investigations on Inclined Airship Bodies,” communication présentée à Proceedings of the International Airship Convention, Bedford, UK, 1998.
- [25] M. Carrión *et al.*, “Computational fluid dynamics challenges for hybrid air vehicle applications,” communication présentée à Progress in Flight Physics, D. Knight *et al.*, édit. Krakow, Poland : EDP Sciences, 2017, p. 43–80. [En ligne]. Disponible : <https://doi.org/10.1051/eucass/201609043>
- [26] A. Dumas *et al.*, “Cfd analysis and optimization of a variable shape airship,” communication présentée à ASME International Mechanical Engineering Congress and Exposition, vol. Volume 7 : Fluids and Heat Transfer, Parts A, B, C, and D, 11 2012, p. 161–166. [En ligne]. Disponible : <https://doi.org/10.1115/IMECE2012-87375>
- [27] M. Carrión *et al.*, “Study of hybrid air vehicle stability using computational fluid dynamics,” *Journal of Aircraft*, vol. 54, p. 1328–1339, 2017. [En ligne]. Disponible : <https://doi.org/10.2514/1.C033987>
- [28] A. Marton *et al.*, “Hybrid model-based and data-driven wind velocity estimator for an autonomous robotic airship,” *Journal of the Brazilian Society of Mechanical Sciences and Engineering*, vol. 42, 03 2020. [En ligne]. Disponible : <http://dx.doi.org/10.1007/s40430-020-2215-8>
- [29] Y. T. Liu *et al.*, “Deep Residual Reinforcement Learning based Autonomous Blimp Control,” communication présentée à 2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2022, p. 12 566–12 573. [En ligne]. Disponible : <https://doi.org/10.1109/IROS47612.2022.9981182>

- [30] Y. Zuo, Y. T. Liu et A. Ahmad, “Autonomous Blimp Control via H-infinity Robust Deep Residual Reinforcement Learning,” mars 2023, prépublication. [En ligne]. Disponible : <http://dx.doi.org/10.48550/arXiv.2303.13929>
- [31] E. Todorov, T. Erez et Y. Tassa, “Mujoco : A physics engine for model-based control,” communication présentée à 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems. IEEE, 2012, p. 5026–5033. [En ligne]. Disponible : <https://doi.org/10.1109/IROS.2012.6386109>
- [32] T. Erez, Y. Tassa et E. Todorov, “Simulation tools for model-based robotics : Comparison of bullet, havok, mujoco, ode and physx,” communication présentée à 2015 IEEE International Conference on Robotics and Automation (ICRA), 2015, p. 4397–4404.
- [33] A. Raffin *et al.*, “Stable-baselines3 : Reliable reinforcement learning implementations,” *Journal of Machine Learning Research*, vol. 22, n°. 268, p. 1–8, 2021. [En ligne]. Disponible : <http://jmlr.org/papers/v22/20-1364.html>
- [34] M. Towers *et al.*, “Gymnasium : A Standard Interface for Reinforcement Learning Environments,” nov. 2024, prépublication.
- [35] J. Collins *et al.*, “A review of physics simulators for robotic applications,” *IEEE Access*, vol. 9, p. 51 416–51 431, 2021. [En ligne]. Disponible : <https://doi.org/10.1109/ACCESS.2021.3068769>
- [36] S. Tunyasuvunakool *et al.*, “dm_control : Software and tasks for continuous control,” *Software Impacts*, vol. 6, p. 100022, 2020. [En ligne]. Disponible : <https://www.sciencedirect.com/science/article/pii/S2665963820300099>
- [37] T. Liesk, M. Nahon et B. Boulet, “Design and experimental validation of a nonlinear low-level controller for an unmanned fin-less airship,” *IEEE Transactions on Control Systems Technology*, vol. 21, n°. 1, p. 149–161, 2013. [En ligne]. Disponible : <https://doi.org/10.1109/TCST.2011.2178415>
- [38] S. Q. Liu *et al.*, “Vectorial backstepping method–based trajectory tracking control for an under-actuated stratospheric airship,” *The Aeronautical Journal*, vol. 121, n°. 1241, p. 916–939, 2017. [En ligne]. Disponible : <https://doi.org/10.1017/aer.2017.52>
- [39] A. Moutinho et J. Azinheira, “Stability and robustness analysis of the aurora airship control system using dynamic inversion,” communication présentée à Proceedings of the 2005 IEEE International Conference on Robotics and Automation, 2005, p. 2265–2270. [En ligne]. Disponible : <https://doi.org/10.1109/ROBOT.2005.1570450>
- [40] A. Moutinho et J. R. Azinheira, “Path control of an autonomous airship using dynamic inversion,” *IFAC Proceedings Volumes*, vol. 37, n°. 8, p. 633–638, 2004, iFAC/EURON Symposium on Intelligent Autonomous Vehicles, Lisbon, Portugal, 5-7

- July 2004. [En ligne]. Disponible : <https://www.sciencedirect.com/science/article/pii/S1474667017320499>
- [41] Z. Zheng et L. Sun, “Adaptive sliding mode trajectory tracking control of robotic airships with parametric uncertainty and wind disturbance,” *Journal of the Franklin Institute*, vol. 355, n°. 1, p. 106–122, 2018. [En ligne]. Disponible : <https://www.sciencedirect.com/science/article/pii/S0016003217305707>
- [42] C. Tang *et al.*, “Deep reinforcement learning for robotics : A survey of real-world successes,” *Annual Review of Control, Robotics, and Autonomous Systems*, 2024. [En ligne]. Disponible : <https://doi.org/10.1146/annurev-control-030323-022510>
- [43] E. Kaufmann *et al.*, “Champion-level drone racing using deep reinforcement learning,” *Nature*, vol. 620, p. 982–987, août 2023. [En ligne]. Disponible : <https://doi.org/10.1038/s41586-023-06419-4>
- [44] J. Eschmann, D. Albani et G. Loianno, “Learning to Fly in Seconds,” *IEEE Robotics and Automation Letters*, vol. 9, n°. 7, p. 6336–6343, juill. 2024. [En ligne]. Disponible : <https://doi.org/10.48550/arXiv.2311.13081>
- [45] G. B. Margolis et P. Agrawal, “Walk these ways : Tuning robot control for generalization with multiplicity of behavior,” communication présentée à 6th Annual Conference on Robot Learning, 2022. [En ligne]. Disponible : <https://openreview.net/forum?id=52c5e73Sls2>
- [46] J. Ko *et al.*, “Gaussian processes and reinforcement learning for identification and control of an autonomous blimp,” communication présentée à Proceedings 2007 IEEE International Conference on Robotics and Automation, 2007, p. 742–747. [En ligne]. Disponible : <https://doi.org/10.1109/ROBOT.2007.363075>
- [47] A. Rottmann *et al.*, “Autonomous blimp control using model-free reinforcement learning in a continuous state and action space,” communication présentée à 2007 IEEE/RSJ International Conference on Intelligent Robots and Systems, 2007, p. 1895–1900. [En ligne]. Disponible : <https://doi.org/10.1109/IROS.2007.4399531>
- [48] H. Gou *et al.*, “Path following control for underactuated airships with magnitude and rate saturation,” *Sensors*, vol. 20, n°. 24, 2020. [En ligne]. Disponible : <https://doi.org/10.3390/s20247176>
- [49] C. Nie, Z. Zheng et M. Zhu, “Three-dimensional path-following control of a robotic airship with reinforcement learning,” *International Journal of Aerospace Engineering*, vol. 2019, n°. 1, p. 7854173, 2019. [En ligne]. Disponible : <https://onlinelibrary.wiley.com/doi/abs/10.1155/2019/7854173>

- [50] L. Catar, I. Tabiai et D. St-Onge, “CAVERNAUTE : A design and manufacturing pipeline of a rigid but foldable indoor airship aerial system for cave exploration,” sept. 2024, prépublication.
- [51] Windreiter Team, “Main page : The silent_runner,” 2020. [En ligne]. Disponible : <http://silent-runner.net/>
- [52] S. B. V. Gomes, “An Investigation of the Flight of Airships with Application to the YEZ-2A,” Thèse de doctorat, Cranfield Institute of Technology, oct. 1990.
- [53] Y. T. Liu *et al.*, “Autonomous Blimp Control using Deep Reinforcement Learning,” sept. 2021, prépublication. [En ligne]. Disponible : <https://doi.org/10.48550/arXiv.2109.10719>
- [54] J. Schulman *et al.*, “Proximal policy optimization algorithms,” 2017. [En ligne]. Disponible : <https://arxiv.org/abs/1707.06347>
- [55] A. D. Ames *et al.*, “Control barrier functions : Theory and applications,” communication présentée à 2019 18th European Control Conference (ECC), 2019, p. 3420–3431. [En ligne]. Disponible : <https://doi.org/10.23919/ECC.2019.8796030>

**ANNEXE A DÉFINITIONS DES VARIABLES DU MODÈLE
DYNAMIQUE DU DIRIGEABLE**

TABLEAU A.1 Tableau des définitions des variables utilisées dans le modèle dynamique du dirigeable.

Variable	Définition
M	Matrice de masse, contenant l'inertie réelle et virtuelle.
$x = [u, v, w, p, q, r]$	Vecteur d'état : vitesses linéaires et angulaires dans le référentiel du corps.
x_a	Vecteur d'état : vitesses linéaires et angulaires dans le référentiel du vent relatif.
F_d	Vecteur dynamique des forces et moments dus aux effets de Coriolis et centrifuges.
F_a	Vecteur des forces et moments aérodynamiques.
F_p	Vecteur des forces et moments de propulsion.
G	Vecteur des forces et moments dus à la flottabilité et la gravité.
m	Masse du robot dirigeable.
$X_{\dot{u}}, Y_{\dot{v}}, Y_{\dot{v}}$	Termes de masse virtuelle.
m_x, m_y, m_z	Masses en selon l'axe contenant la masse du système et la masse vituelle.
OC	distance entre le centre de gravité et le centre de volume.
a_x, a_y, a_z	Coordonnées du centre de gravité par rapport au centre de volume.
I_x, I_y, I_z	Termes d'inertie par rapport aux axes OX, OY et OZ.
$L_{\dot{p}}, M_{\dot{q}}, N_{\dot{r}}$	Termes d'inertie virtuelle par rapport aux axes OX, OY, et OZ.
J_x, J_y, J_z	Termes d'inertie combinant l'intertie du système et l'inertie virtuelle.
$T_{ds}, T_{dp}, T_{ls}, T_{lp}$	Poussées des moteurs de propulsion et de levage à bâbord et tribord.
d_x, d_y, d_z	Distances entre le centre de volume et l'hélice concerné
V_t	Vitesse vraie.
V_e	Vitesse équivalente.
α	Angle d'incidence.
β	Angle de dérapage
$C_D, C_Y, C_L, C_l, C_m, C_n$	Coefficients aérodynamiques (traînée, portance, force latérale, moments).
$A_X, A_Y, A_Z, A_L, A_M, A_N$	Forces aérodynamiques (traînée, portance, force latérale, moments).
F_{ref}, M_{ref}	Force et moment de référence
S_{ref}, L_{ref}	Surface et longueur de référence, respectivement Volume ^{2/3} et Volume ^{1/3}
ρ	Densité de l'air.
q^*, r^*	Taux de giration non dimensionnels.
C_{Yr}, C_{nr}, C_{Zq} et C_{mq}	Coefficient dynamiques pour amortissement.
S_1, S_2	Matrices de transformation.