



gz-unitree: Reinforcement learning en robotique avec validation par moteurs de physique multiples pour le robot H1v2 d'Unitree

Gwenn Le Bihan

gwenn.lebihan7@gmail.com

ENSEEIH

6 Novembre 2025

1 Remerciements

Je tiens à remercier Olivier Stasse et Guilhem Saurel, qui m’ont suivie pendant toute la durée du stage et répondu à mes questionnements, sans qui je n’aurais pu mener ces recherches. Merci aussi à Côme Perrot, grâce à qui j’ai pu éclaircir certaines zones d’ombre sur ma compréhension de la théorie du reinforcement learning appliquée à la robotique.

Merci à Sylvie Chambon, Géraldine Morin et Ronan Guivarch, professeurs à l’ENSEEIH, qui m’ont soutenue à travers des périodes parfois difficiles.

Merci aussi à Laurenz Mäde et Marin Haug pour avoir créé Typst, une alternative moderne à LaTeX qui a rendu l’écriture de ce rapport bien plus agréable. Merci à Joseph Wilson (paquet Typst « Fletcher ») et Johannes Wolf (paquet Typst « CeTZ »), qui ont rendu la création de diagrammes très ergonomique.

Table des matières

1	Remerciements	2
2	Contexte	3
2.1	Bases théoriques du <i>Reinforcement Learning</i>	3
2.1.1	L’entraînement	4
2.1.2	Deep Reinforcement Learning	5
2.1.3	Difficultés liées à l’implémentation de la fonction coût	6
2.2	Évaluation de la performance d’une politique	7
2.2.1	Chemins d’états possibles \mathcal{C}	7
2.2.2	Récompense attendue η	9
2.2.3	Avantage A	10
2.2.4	Lien entre η et A	10
2.2.5	<i>Surrogate advantage</i> \mathcal{L}	11
2.3	Méthodes d’optimisation de politique	11
2.3.1	<i>Trust Region Policy Optimization</i>	11
2.3.2	<i>Proximal Policy Optimization</i>	13
2.4	Application en robotique	14
2.4.1	Spécification de la tâche	14
2.4.2	Inventaire des simulateurs en robotique	14
2.5	Reproductibilité logicielle	15
3	Packaging reproductible avec Nix	16
3.1	Reproductibilité	16
3.1.1	État dans le domaine de la programmation	16
3.1.2	Contenir les effets de bords	16
3.1.3	État dans le domaine de la robotique	16
3.1.4	Reproductibilité de la compilation	17
3.2	Nix, le gestionnaire de paquets pur	17
3.2.1	Un <i>DSL</i> ¹ fonctionnel	17
3.2.2	Un écosystème de dépendances	18
3.2.3	Une compilation dans un environnement fixé	18
3.3	NixOS, un système d’exploitation à configuration déclarative	19

¹Domain-Specific Language

3.4	Packaging Nix pour <i>gz-unitree</i>	19
4	Étude du SDK d'Unitree et du bridge SDK \Leftarrow MuJoCo	20
4.1	Le robot <i>H1v2</i> d'Unitree	20
4.2	Canaux DDS	21
4.3	Une base de code partiellement open-source	22
4.4	Un autre bridge existant: unitree_mujoco	23
5	Développement du bridge SDK \Leftarrow Gazebo	25
5.1	Établissement du contact	25
5.2	Installation du plugin dans Gazebo	27
5.3	Architecture du plugin	27
5.4	Calcul des nouveaux couples des moteurs	29
5.5	rt/lowcmd	30
5.6	rt/lowstate	31
5.6.1	Construction d'un message rt/lowstate	31
5.6.2	Émission de l'état	33
5.7	Désynchronisations	34
5.8	Vérification sur des politiques réelles	35
5.9	Amélioration des performances	35
5.10	Enregistrement automatique de vidéos	37
5.11	Mise en CI/CD	38
5.11.1	Une image de base avec Docker	38
5.11.2	Une pipeline Github Actions	39
6	Conclusion	40
	Bibliographie	40
A	Preuves	46
A.1	Cas dégénéré de $D_{\text{KL}}(Q, Q') = 0$ sans utilisation de max	46
A.2	$\eta(\pi, r)$ comme une espérance	46
A.3	Simplification de l'expression de $L(s, a, \Pi, \Pi', R)$ dans PPO-Clip	48
B	Implémentation de CRC32	49

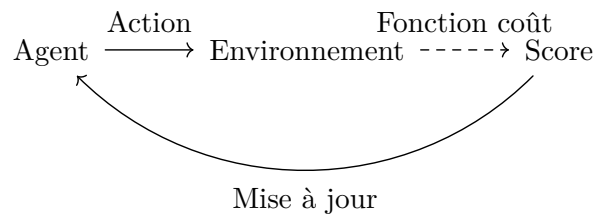
2 Contexte

2.1 Bases théoriques du *Reinforcement Learning*

L'apprentissage par renforcement, ou *Reinforcement Learning*, permet de développer des programmes sans expliciter leur logique: on décrit plutôt quatre choses, qui vont permettre à la logique d'émerger pendant la phase d'entraînement:

- Un *agent*: c'est le programme que l'on souhaite créer
- Des *actions* que l'agent peut choisir d'effectuer ou pas
- Un *environnement*, que les actions viennent modifier
- Un *score* (*coût* s'il doit être minimisé, *récompense* inversement) qui dépend de l'état pré- et post-action de l'environnement ainsi que de l'action qui a été effectuée

La phase d'apprentissage consiste à trouver, par des cycles d'essai/erreur, quelles sont les meilleures actions à prendre en fonction de l'environnement actuel, avec « meilleur » définit comme « qui minimise le coût » (ou maximise la récompense):



Cette technique est particulièrement adaptée aux problèmes qui se prêtent à une modélisation type « jeu vidéo », dans le sens où l'agent représente le personnage-joueur, et le coût un certain score, qui est condition de victoire ou défaite.

En robotique, une approche similaire explore l'espace d'action (en général un courant à envoyer aux moteurs) de façon à optimiser le coût.

En robotique, on a des correspondances claires pour ces quatre notions :

Agent	Robot pour lequel on développe le programme de contrôle (appelé <i>politique</i>)
Actions	Envoi d'ordres aux moteurs, souvent le courant électrique à appliquer
Environnement	Le monde réel. C'est de loin la partie la plus difficile à simuler informatiquement. On utilise des moteurs de simulation physique, dont la pluralité des implémentations est importante, voir Chapitre 2.1.3.3
Coût	Ensemble de contraintes (« ne pas endommager le robot ») et d'évaluations spécifiques à la tâche à effectuer (« s'est déplacé de 5m en avant selon l'axe x »).

2.1.1 L'entraînement

Une fois que ce cadre est posé, il reste à savoir *comment* l'on va trouver la fonction qui associe un état de l'environnement à une action.

Une première approche naïve, mais suffisante dans certains cas, consiste à faire une recherche exhaustive et à stocker dans un simple tableau la meilleure action à faire en fonction de chaque état de l'environnement :

État actuel (x , retour)	Action à effectuer +1 ou -1	Coûts associés	
		Pour +1	Pour -1
(0, C'est plus)			
(1, C'est plus)			
(3, C'est moins)			
(4, C'est moins)			
(5, C'est moins)			

Tableau 1. – Mémoire exhaustive initiale pour un « C'est plus ou c'est moins » dans $\llbracket 0, 5 \rrbracket \times \{\text{C'est plus, C'est moins}\}$, avec pour solution 2

L'entraînement consiste donc ici en l'exploration de l'entièreteré des états possibles de l'environnement, et, pour chaque état, le calcul du coût associé à chaque action possible.

Il faut définir la fonction de coût, souvent appelée L pour *loss*:

$$L : E \rightarrow S \quad (1)$$

avec E l'ensemble des états possibles de l'environnement, et S un ensemble muni d'un ordre total (on utilise souvent $[0, 1]$). Ces fonctions coût, qui ne dépendent que de l'état actuel de l'environnement, représente un domaine du RL^2 appelé *Q-Learning* [1]

On remplit la colonne « Action à effectuer » avec l'action au coût le plus bas:

État actuel (x , retour)	Action à effectuer +1 ou -1	Coûts associés	
		Pour +1	Pour -1
(0, C'est plus)	+1	$L(1,) = 1$	$L(0,) = 2$
(1, C'est plus)	+1	$L(2,) = 0$	$L(0,) = 2$
(3, C'est moins)	-1	$L(3,) = 1$	$L(1,) = 1$
(4, C'est moins)	-1	$L(4,) = 2$	$L(2,) = 0$
(5, C'est moins)	-1	$L(5,) = 3$	$L(3,) = 1$

Tableau 2. – Entraînement terminé, avec pour fonction coût $L = (x, \text{retour}) \mapsto |x - 2|$ la distance à la solution

Ici, cette approche exhaustive suffit parce que l'ensemble des états possibles de l'environnement, E , possède 6 éléments

Cependant, ces ensembles sont bien souvent prohibitivement grands (e.g. $x \in \llbracket 0, 10^{34} \rrbracket$), infinis ($x \in \mathbb{N}$) ou indénombrables ($x \in \mathbb{R}$)

Dans le cas de la robotique, E est une certaine représentation numérique du monde réel autour du robot, on imagine donc bien qu'il y a beaucoup trop d'états possibles.

2.1.2 Deep Reinforcement Learning

Une façon de remédier à ce problème de dimensions est de remplacer le tableau exhaustif par un réseau de neurones:

État actuel	devient la couche d'entrée
Meilleure action	devient la couche de sortie
Coûts associés	devient la fonction à optimiser par descente de gradient
Le remplissage du tableau	devient la rétropropagation du gradient pendant l'entraînement

²Reinforcement Learning

Mise à jour (*Q-learning*)

Le score associé à un état s_t et une action a_t , appelée $Q(s_t, a_t)$ ici pour « quality » [2] ou « action-value » [3], est mis à jour avec cette valeur [4]:

$$(1 - \alpha) \underbrace{Q(s_t, a_t)}_{\text{valeur actuelle}} + \alpha \left(\underbrace{R_{t+1}}_{\substack{\text{récompense} \\ \text{pour cette action}}} + \gamma \underbrace{\max_a Q(S_{t+1}, a)}_{\substack{\text{récompense de la meilleure} \\ \text{action pour l'état suivant}}} \right) \quad (2)$$

L'expression comporte deux hyperparamètres, à valeurs dans $]0, 1[$:

Learning rate α contrôle à quel point l'on favorise l'évolution de Q ou pas.

Discount factor γ contrôle l'importance que l'on donne aux récompenses futures. Il est utile de commencer avec une valeur faible puis l'augmenter avec le temps [5].

2.1.3 Difficultés liées à l'implémentation de la fonction coût

Tendances à la « tricherie » des agents

Expérimentalement, on sait que des tendances « tricheuses » émergent facilement pendant l'entraînement [Réf. nécessaire]: l'agent découvre des séries d'actions qui causent un bug avantageux vis à vis du coût associé, soit parce qu'il y a un bug dans le calcul de l'état de l'environnement post-action, soit parce que la fonction coût ne prend pas suffisamment bien en compte toutes les possibilités de l'environnement (autrement dit, il manque de contraintes).

Dans le cas de la robotique, cela arrive particulièrement souvent [Réf. nécessaire], et il faut donc un simulateur qui soit suffisamment réaliste.

Sous-spécification de la fonction coût

Un exemple populaire est l'expérience de pensée du Maximiseur de trombones [6]: On imagine un agent avec pour environnement le monde réel, pour actions « prendre des décisions »; « envoyer des emails »; etc. et pour fonction récompense « le nombre de trombones existant sur Terre ». Il finirait possiblement par réduire en esclavage tout être vivant capable de produire des trombones: la fonction coût est sous-spécifiée

La validation comme méthode de mitigation

Comme ces bugs sont des comportements non voulus, il est très probables qu'ils ne soient pas exactement les mêmes en changeant d'implémentation.

Il convient donc de se servir de *plusieurs* implémentations: une sert à la phase d'entraînement, pendant laquelle l'agent développe des « tendances à la tricherie », puis une autre sert à la phase de *validation*.

Cette phase consiste en le lancement de l'agent dans une autre implémentation, avec les mêmes actions mais qui, crucialement, ne comporte pas les mêmes bugs que l'environnement ayant servi à la phase d'apprentissage.

Les « techniques de triche » ainsi apprises deviennent inefficace, et si le score devient bien pire que celui de l'apprentissage, on peut détecter les cas de triche.

On peut même aller plus loin, et multiplier les phases de validation avec des implémentations supplémentaires, ce qui réduit encore la probabilité qu’une technique de triche se glisse dans l’agent final.

2.2 Évaluation de la performance d’une politique

Théoriquement, le score associé à un couple état/action est souvent réduit à l’intervalle $[0, 1]$ et assimilé à une distribution de probabilité: Q est une fonction de $S \times A$ vers $[0, 1]$ qui renvoie la probabilité qu’a l’agent à choisir une action en étant dans un état de l’environnement.

On note dans le reste de cette section:

A	l’ensemble des actions
S	l’ensemble des états possibles de l’environnement
$\rho_0 : S \rightarrow [0, 1]$	la distribution de probabilité de l’état initial de l’environnement. Si l’on initialise l’environnement de manière uniformément aléatoire, ρ_0 est une équiprobabilité ³
$M : S \times A \rightarrow S$	le moteur de simulation physique, qui applique l’action à un état de l’environnement et envoie le nouvel état de l’environnement
$\Pi : S \rightarrow A$	une politique
$\Pi^* : S \rightarrow A$	la meilleure politique possible, celle que l’on cherche à approcher
$R : S \rightarrow \mathbb{R}^+$	sa fonction de récompense (TODO: incohérent! c’est sensé être Q , qu’on a assimilé à une distrib de proba :/)
$Q_\pi : S \times A \rightarrow [0, 1]$	la distribution de probabilité d’une politique π , qu’on suppose Markovienne (elle ne dépend que de l’état dans lequel on est). $Q_\pi(s_t, a_t)$ est la probabilité que π choisisse a_t , quand on est dans l’état s_t (s_t est l’état pré -action, et non post-action)
Q (resp. Q^*)	Q_Π (resp. Q_{Π^*}), pour alléger les notations

On suppose A et S dénombrables⁴.

Pour alléger les notations, on surchargera les fonctions récompenses pour qu’elle puissent prendre en entrée des éléments de $S \times A$, en ignorant simplement l’action choisie:

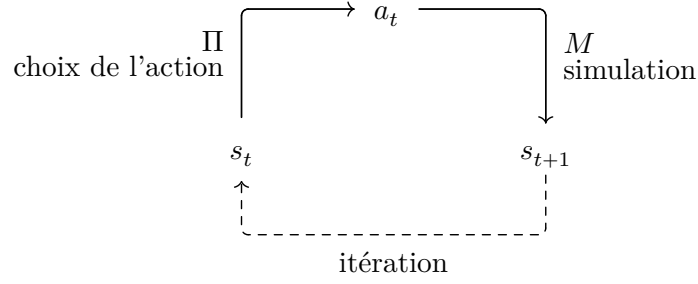
$$\forall (s, a) \in S \times A, \forall r \in \text{récompenses}, r(s, a) := r(s) \quad (3)$$

2.2.1 Chemins d’états possibles \mathcal{C}

M et Π forment en fait tout ce qui se passe pendant un pas de temps. c’est cette boucle que l’on répète pour entraîner l’agent (si l’on met Π à jour à chaque tour de boucle) ou l’utiliser (on parle alors d’inférence):

³i.e. $\text{card } \rho_0(S) = 1$

⁴En pratique, \mathbb{R} est discrétisé dans les simulateurs numérique, donc cette hypothèse ne pose pas de problèmes à l’application de la théorie au domaine de la robotique



Quand on « déroule » Π en partant d'un certain état initial s_0 , on obtient une suite d'états et d'actions:

$$s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} s_2 \xrightarrow{a_2} \dots$$

Pour tout pas de temps $t \in \mathbb{N}$, on a:

$$\begin{cases} a_t = \Pi(s_t) \\ s_{t+1} = M(s_t, a_t) \end{cases} \quad (4)$$

On rappelle que M est la fonction « simulation », qui renvoie l'état post-action en fonction de l'action choisie et de l'état pré-action.

Un chemin se modélise aisément par une suite d'éléments de⁵ $S \times A$. Ainsi, on note

$$\mathcal{C}_\pi := \left\{ (s_t, a_t)_{t \in \mathbb{N}} \text{ avec } \begin{cases} a_0 = \pi(s_0) \\ \forall t \in \mathbb{N} \quad a_{t+1} = \pi(s_t) \\ \forall t \in \mathbb{N} \quad s_{t+1} = M(s_t, a_t) \end{cases} \middle| s_0 \in S \right\} \quad (5)$$

l'ensemble des chemins possibles avec la politique π . C'est tout simplement l'ensemble de tout les « déroulements » de la politique π en partant des états possibles de l'environnement.

On définit également l'ensemble de *tout* les chemins d'états possibles, peu importe la politique, \mathcal{C} :

$$\mathcal{C} := \left\{ \left\{ \forall t \in \mathbb{N} \quad \begin{cases} c_0 = (s_0, a_0) \\ c_{t+1} = (M(c_t), a_t) \end{cases} \middle| (s_0, a) \in S \times A^\mathbb{N} \right\} \right\} \quad (6)$$

⁵Il est essentiel de conserver l'information de l'action prise entre chaque état (contrairement à ce que fournirait une simple suite d'éléments de S , par exemple) pour pouvoir calculer les probabilités par rapport à une politique le long d'un chemin. En effet, si l'on veut obtenir la probabilité que Π ait résulté en un état $s \in S$ de l'environnement, il faut savoir « par quelle $a \in A$ est-on passé », Q_Π prenant bien s et a en entrée.

On notera que, selon M , on peut avoir $\mathcal{C} \subsetneq (S \times A)^{\mathbb{N}}$: par exemple, certains états de l'environnement peuvent représenter des « impasses », où il est impossible d'évoluer vers un autre état, peu importe l'action choisie.

On note aussi que \mathcal{C} et \mathcal{C}_π sont dénombrables: Ils sont construits à partir de $(S \times A)^{\mathbb{N}}$ et S , et A & \mathbb{N} sont également dénombrables⁶

*Cette formalisation est utile par la suite,
pour proprement définir certaines grandeurs.*

Remarque

Les définitions suivantes, dont la plupart proviennent du papier *Trust Region Policy Optimization*, citation « [7] », ont été reformulées pour utiliser cette notion de chemins.

(TODO: Pas clair)

Notamment, les espérances le long d'un chemin, notées $\mathbb{E}_{s_0, a_0, \dots}$ dans [7], sont dénotées ici par une opération-sur-ensemble usuelle⁷, avec $\mathbb{E}_{c \in \mathcal{C}}$. De même, la notation $\mathbb{E}_{s_0, a_0, \dots \sim \pi}$ est dénotée $\mathbb{E}_{c \sim \pi \in \mathcal{C}}$ et explicitée après (8).

Dans la documentation de *OpenAI Spinning Up* (citation « [8] »), les espérances sont notées $E_{s, a \sim \pi}$, ce qui correspond à faire une espérance *le long* de tout chemin: cela correspond ici à $\mathbb{E}_{c \sim \pi \in \mathcal{C}} \sum_{t=0}^{\infty} \dots$.

2.2.2 Récompense attendue η

η représente la récompense moyenne à laquelle l'on peut s'attendre pour une politique π avec fonction de récompense r .

Elle prend en compte le *discount factor* γ : les récompenses des actions deviennent de moins en moins importantes avec le temps. η est définie ainsi [7]

$$\eta(\pi, r) := \underbrace{\sum_{(c_t)_{t \in \mathbb{N}} \in \mathcal{C}} \rho_0(s_0) \prod_{t=0}^{\infty} Q_\pi(c_t)}_{\text{probabilité du chemin}} \underbrace{\sum_{t=0}^{\infty} \gamma^t r(c_t)}_{\text{récompense associée}} \quad (7)$$

pour tout chemin possible

On peut également exprimer $\eta(\pi, r)$ comme une espérance. On a (cf preuve en A.2)

$$\eta(\pi, r) = \mathbb{E}_{C \sim \pi \in \mathcal{C}} \left(\sum_{t=0}^{\infty} \gamma^t r(C_t) \right) \quad (8)$$

Avec $C \sim \pi \in \mathcal{C}$ signifiant

- C est une variable aléatoire à valeurs dans \mathcal{C}
- C suit la même loi que π

⁶On a $\text{card } \mathcal{C} \leq \text{card}((S \times A)^{\mathbb{N}}) = \text{card}(S \times A)^{\text{card } \mathbb{N}} = (\text{card } S \text{ card } A)^{\text{card } \mathbb{N}} \leq (\aleph_0)^{\text{card } \mathbb{N}} = 2^{\aleph_0} = \aleph_0$. De plus, $\mathcal{C}_\pi \subset \mathcal{C}$, donc $\text{card } \mathcal{C}_\pi \leq \text{card } \mathcal{C} \leq \aleph_0$

⁷d'autres exemples d'« opérations-sur-ensemble » sont $\sum_{x \in \mathbb{R}}$ ou $\prod_{n \in \mathbb{N}}$, par exemple.

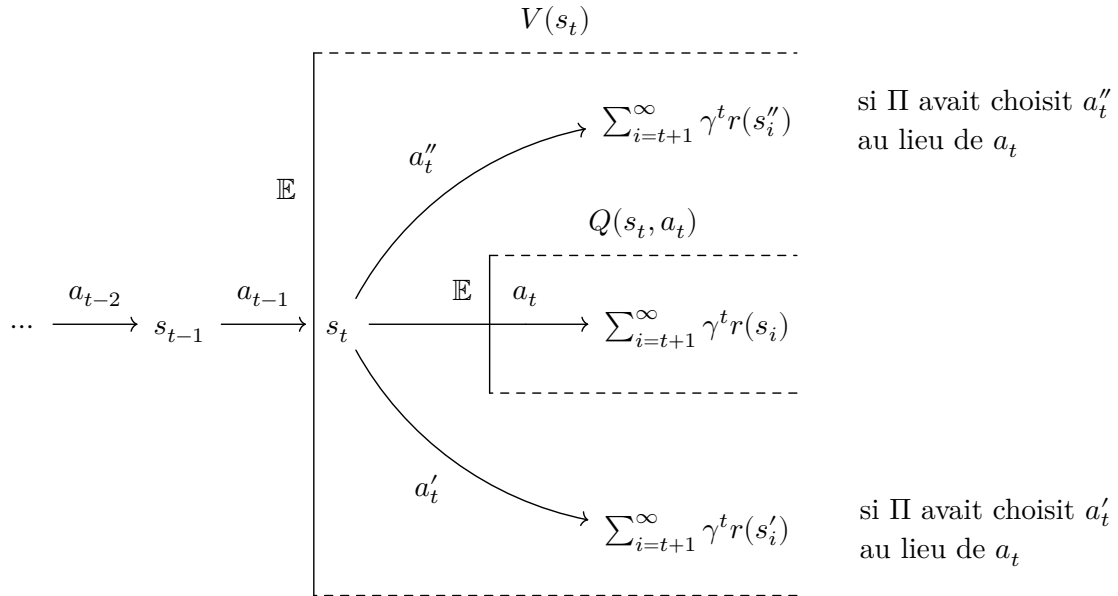
2.2.3 Avantage A

L'avantage $A_{\pi,r}(s, a)$ mesure à quel point il est préférable de choisir l'action a quand on est dans l'état s (pour la politique π , avec « préférable » au sens de⁸ $\geq_{r(M(s,\cdot))}$)

Pour calculer $A_{\pi,r}(s, a)$, on regarde l'espérance des récompenses cumulées pour tout chemin commençant par s , et on la compare à celle pour tout chemin commençant par $M(s, a)$

$$A_{\pi,r}(s, a) := \underbrace{\mathbb{E}_{\substack{(s_t, a_t)_{t \in \mathbb{N}} \sim \pi \in \mathcal{C} \\ s_0 = s \\ s_1 = M(s_0, a)}}}_{Q(s, a)} \sum_{t=0}^{\infty} \gamma^t r(s_t) - \underbrace{\mathbb{E}_{\substack{(s_t, a_t)_{t \in \mathbb{N}} \sim \pi \in \mathcal{C} \\ s_0 = s}}}_{V(s)} \sum_{t=0}^{\infty} \gamma^t r(s_t) \quad (9)$$

On peut visualiser ce calcul ainsi:



On considère tout les chemins à partir de l'état s_t , et l'on regarde l'espérance...

pour $V(s_t)$ de tout les chemins

pour $Q(s_t, a_t)$ du chemin où l'on a choisi a_t

En suite, il suffit de faire la différence, pour savoir l'avantage que l'on a à choisir a_t par rapport au reste.

2.2.4 Lien entre η et A

Pour une fonction de récompense r donnée, A permet de calculer η pour une politique π en fonction de la valeur de η pour une autre politique π' [7]

⁸En posant, pour toute fonction $f : I \rightarrow O$, avec O ordonné par \geq : $\forall i \in I^2, i_1 \geq_f i_2 := f(i_1) \geq f(i_2)$. Ici donc, on compare les politiques selon $a \mapsto r(M(s, a))$. Autrement dit, la récompense associé à l'état obtenu après le choix d'une action, depuis l'état s

$$\eta(\pi', r) = \eta(\pi, r) + \mathbb{E}_{(c_t)_{t \in \mathbb{N}} \sim \pi' \in \mathcal{C}} \sum_{t=0}^{\infty} \gamma^t A_{\pi, r}(c_t) \quad (10)$$

2.2.5 Surrogate advantage \mathcal{L}

Il est théoriquement possible d'utiliser A pour optimiser une politique, en maximisant sa valeur:

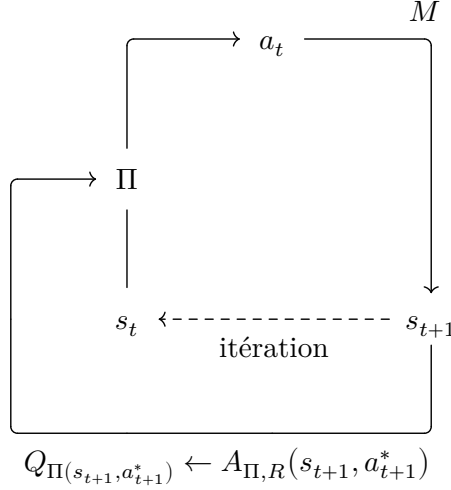


Fig. 5. – Boucle d'entraînement

Avec

$$a_{t+1}^* := \operatorname{argmax}_{a \in A} A_{\Pi, R}(s_{t+1}, a) \quad (11)$$

Mais, en pratique, des erreurs d'approximation peuvent rendre $A_{\Pi, R}(s_{t+1}, a_{t+1}^*)$ négatif, ce qui empêche de s'en servir pour définir une valeur de Q_{Π} [7]

Le *surrogate advantage* détermine la performance d'une politique par rapport à une autre [8]

$$\mathcal{L}_r(\pi', \pi) := \mathbb{E}_{(s_t, a_t)_{t \in \mathbb{N}} \in \mathcal{C}} \sum_{t=0}^{\infty} \frac{Q_{\pi}(s_t, a_t)}{Q_{\pi'}(s_t, a_t)} A_{\pi, r}(s_t, a_t) \quad (12)$$

2.3 Méthodes d'optimisation de politique

2.3.1 Trust Region Policy Optimization

La méthode TRPO définit la mise à jour de Q avec un Q' qui maximise le *surrogate advantage* [8], sous une contrainte limitant l'ampleur des modifications individuelles, ce qui procure une stabilité à l'algorithme, et évite qu'un seul « faux pas » dégrade violemment la performance de la politique.

$$Q' = \begin{cases} \operatorname{argmax}_q \mathcal{L}_r(q, Q) \\ \text{s.c. distance}(Q', Q) < \delta \end{cases} \quad (13)$$

Avec δ une limite supérieure de distance entre Q' , la nouvelle politique, et Q , l'ancienne.

Distance entre politiques

Il existe plusieurs manières de mesurer l'écart entre deux distributions de probabilité, dont notamment la *divergence de Kullback-Leibler*, aussi appelée entropie relative [9], [10]:

$$D_{\text{KL}}(P \parallel P') := \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{P'(x)} \quad (14)$$

Avec \mathcal{X} l'espace des échantillons et P, P' deux distributions de probabilité sur celui-ci. Dans notre cas, $\mathcal{X} = S \times A$,

Pour évaluer cette distance, on regarde la plus grande des distances entre des paires de distributions de probabilité de politiques Q_{Π} et $Q_{\Pi'}$, pour tout $s \in S$ [7]

$$\max_{s \in S} D_{\text{KL}}(Q_{\Pi'}(s, \cdot) \parallel Q_{\Pi}(s, \cdot)) < \delta \quad (15)$$

En notant $Q_{\pi}(s, \cdot) := a \mapsto Q_{\pi}(s, a)$. On a donc ici « $\mathcal{X} = A$ » dans la définition de D_{KL}

Pourquoi faire le maximum sur chaque $s \in S$?

Ce maximum revient à limiter non pas la simple distance entre les deux politiques, mais *limiter la modification de la politique sur chacune de ses actions*.

Ceci permet d'éviter d'avoir deux politiques jugées similaires par D_{KL} à cause de modifications se « compensant ». Par exemple, avec

$$\forall s \in S, Q(s, 1) = Q(s, 2) \quad (16)$$

et

$$Q' := (s, a) \mapsto \begin{cases} Q(s, a) \cdot 2 & \text{si } a = 1 \\ Q(s, a) \cdot \frac{1}{2} & \text{si } a = 2 \\ Q(s, a) & \text{sinon} \end{cases} \quad (17)$$

On a $D_{\text{KL}}(Q, Q') = 0$ (cf preuve en A.1), alors qu'il y a eu une modification très importante des probabilités de choix de l'action 1 et 2 dans tout les états possibles : si on imagine $Q(s, 1) = Q(s, 2) = 1/4$, on a après modification $Q'(s, 1) = 1/2$ et $Q'(s, 2) = 1/8$.

Région de confiance

Cette contrainte définit un ensemble réduit de Π' acceptables comme nouvelle politique, aussi appelé une *trust region* (région de confiance), d'où la méthode d'optimisation tire son nom [7].

En pratique, l'optimisation sous cette contrainte est trop demandeuse en puissance de calcul, on utilise plutôt l'espérance [7].

$$\overline{D_{\text{KL}}} := \mathbb{E}_{s \in S} D_{\text{KL}}(Q(s, \cdot) \parallel Q'(s, \cdot)) \quad (18)$$

2.3.2 Proximal Policy Optimization

La *PPO* repose sur le même principe de stabilisation de l'entraînement par limitation de l'ampleur des changements de politique à chaque pas.

Cependant, les méthodes *PPO* préfèrent changer la quantité à optimiser, pour limiter intrinsèquement l'ampleur des modifications, en résolvant un problème d'optimisation sans contraintes [11]

$$\begin{aligned} \arg\max_{\Pi'} \quad & \mathbb{E}_{(s,a) \in \mathcal{C}} L(s, a, \Pi, \Pi', R) \\ \text{s.c.} \quad & \top \end{aligned} \quad (19)$$

Avec pénalité (*PPO-Penalty*)

PPO-Penalty soustrait une divergence K-L pondérée à l'avantage:

$$L(s, a, \Pi, \Pi', R) = \frac{Q_{\Pi}(s, a)}{Q_{\Pi'}(s, a)} A_{\Pi, R}(s, a) - \beta D_{\text{KL}}(\Pi \parallel \Pi') \quad (20)$$

Avec β ajusté automatiquement pour être dans la même échelle que l'autre terme de la soustraction.

Par *clipping* (*PPO-Clip*)

PPO-Clip utilise une limitation du ratio de probabilités (en minimum et en maximum) [12]

$$\begin{aligned} L(s, a, \Pi, \Pi', R) = \min \Big(& \frac{Q_{\Pi'}(s, a)}{Q_{\Pi}(s, a)} A_{\Pi', R}(s, a), \\ & \text{clip} \left(\frac{Q_{\Pi'}(s, a)}{Q_{\Pi}(s, a)}, 1 - \varepsilon, 1 + \varepsilon \right) A_{\Pi', R}(s, a) \Big) \end{aligned} \quad (21)$$

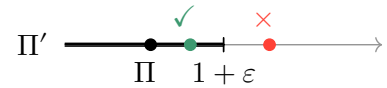
Avec $\varepsilon \in \mathbb{R}_+^*$ un paramètre indiquant à quel point l'on peut s'écarter de la politique précédente, et

$$\text{clip} := (x, m, M) \mapsto \begin{cases} m & \text{si } x < m \\ M & \text{si } x > M \\ x & \text{sinon} \end{cases} \quad (22)$$

La complexité de l'expression, et la présence d'un min au lieu de simplement un clip est due au fait que l'avantage $A_{\Pi', R}(s, a)$ peut être négatif. L'expression se simplifie en séparant les cas (cf preuve en A.3)

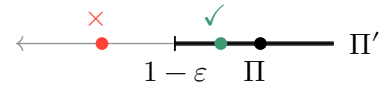
Si l'avantage est positif a est un meilleur choix que $\Pi(s)$.

$$L(s, a, \Pi, \Pi', R) = \min\left(\frac{Q_{\Pi'}(s, a)}{Q_{\Pi}(s, a)}, 1 + \varepsilon\right) A_{\Pi', R}(s, a)$$



Si l'avantage est négatif choisir a est pire que garder $\Pi(s)$.

$$L(s, a, \Pi, \Pi', R) = \max\left(1 - \varepsilon, \frac{Q_{\Pi'}(s, a)}{Q_{\Pi}(s, a)}\right) A_{\Pi', R}(s, a)$$



2.4 Application en robotique

2.4.1 Spécification de la tâche

Le score (récompense ou coût) dépend de la tâche pour laquelle on veut entraîner l'agent.

En robotique, il est commun d'inclure dans la récompense les éléments suivants:

- Couple maximal sur les commandes envoyées au moteurs
- Limite sur la vitesse du robot
- Prévention des auto-collisions (par exemple, le bras qui tape la jambe)
- (TODO: Ajouter [13])
- etc.

2.4.2 Inventaire des simulateurs en robotique

Dans le contexte de la robotique, le calcul de l'état post-action de l'environnement est le travail du *moteur de physique*.

Bien évidemment, ce sont des programmes complexes avec des résolutions souvent numériques d'équation physiques; il est presque inévitable que des bugs se glissent dans ces programmes.

Un environnement de RL⁹ ne se résume pas à son moteur de physique: il faut également charger des modèles 3D, le modèle du robot (qui doit être contrôlable par les actions, on fait donc une émulation de la partie logicielle du robot), et également, pendant les phases de développement, avoir un moteur de rendu graphique, une interface et des outils de développement.

Cet ensemble s'appelle un *simulateur système*.

Isaac

Un simulateur développé par NVIDIA [14], utilisant son propre moteur de rendu, PhysX [15]

MuJoCo

Un simulateur initialement propriétaire. Il a été rendu gratuit puis open source par Google DeepMind [16].

Bien que MuJoCo est décrit comme un moteur de simulation physique et non un simulateur, il embarque une commande `simulate` qui le rend fonctionnellement équivalent à un simulateur [17].

Gazebo

Les intérêts de Gazebo [18] sont multiples:

⁹Reinforcement Learning

- C'est un logiciel open-source *communautaire*, qui ne dépend pas du financement d'une grande entreprise
- Son architecture modulaire permet notamment d'utiliser plusieurs moteurs de simulation physique différents [19], à l'inverse de MuJoCo.
- C'est un *simulateur système*, qui est capable de simuler la partie logicielle du robot en plus de la physique du son modèle 3D.

Gazebo possède des plugins officiels pour divers moteurs de simulation physique:

DART	Plugin <code>gz-physics-dartsim-plugin</code> , c'est l'implémentation principale, et celle par défaut [19], [20].
Bullet	Plugin <code>gz-physics-bulletsim-plugin</code> . En beta [19], [21], [22].
Bullet Featherstone	Plugin <code>gz-physics-bullet-featherstone-plugin</code> , également en beta [19]. Une variable de Bullet, utilisant l'algorithme de Featherstone [23], [24]

2.5 Reproductibilité logicielle

La reproductibilité est particulièrement complexe dans le champ du reinforcement learning [25].

En plus des difficultés de reproductibilité sur l'algorithme lui-même, le paysage logiciel et matériel est riche en dépendances à des bibliothèques, qui elle aussi dépendent d'autres bibliothèques.

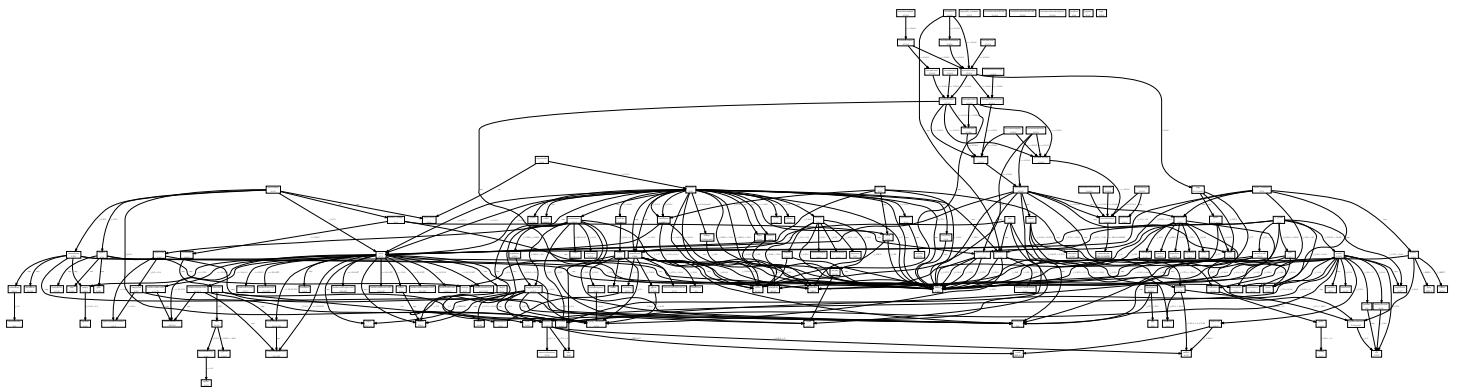


Fig. 8. – Arbre des dépendances pour *Gepetto/h1v2-Isaac*

Bien que toutes ces dépendances puissent être spécifiées avec des contraintes de version strictes [26] pour éviter des changements imprévus de comportement du code venant des bibliothèques, beaucoup celles-ci ont besoin de compiler du code C++ à l'installation¹⁰: fixer la version de la bibliothèque ne suffit pas donc à garantir la reproductibilité de la compilation de l'arbre des dépendances.

¹⁰Pour des raisons de performance [27], certaines bibliothèques implémentent leurs fonctions critiques en C++. C'est par exemple le cas de NumPy [Réf. nécessaire]

3 Packaging reproductible avec Nix

3.1 Reproductibilité

3.1.1 État dans le domaine de la programmation

La différence entre une fonction au sens mathématique et une fonction au sens programmatique consiste en le fait que, par des raisons de praticité, on permet aux **functions** des langages de programmation d’avoir des *effets de bords*. Ces effets affectent, modifient ou font dépendre la fonction d’un environnement global qui n’est pas explicitement déclaré comme une entrée (argument) de la fonction en question [28].

Cette liberté permet, par exemple, d’avoir accès à la date et à l’heure courante, interagir avec un système de fichier d’un ordinateur, générer une surface pseudo aléatoire par bruit de Perlin, etc.

Mais, en contrepartie, on perd une équation qui est fondamentale en mathématiques:

$$\forall E, F, \forall f : E \rightarrow F, \forall (e_1, e_2) \in E^2, e_1 = e_2 \Rightarrow f(e_1) = f(e_2) \quad (23)$$

En programmation, on peut très facilement construire un f qui ne vérifie pas ceci:

```
from datetime import date

def f(a):
    return date.today().year + a
```

Selon l’année dans laquelle nous sommes, $f(0)$ n’a pas la même valeur.

De manière donc très concrète, si cette fonction f fait partie d’un protocole expérimental, l’expérience n’est plus reproductible, et ses résultats sont donc potentiellement non vérifiables, si le papier est soumis le 15 décembre 2025 et la *peer review* effectuée le 2 janvier 2026.

3.1.2 Contenir les effets de bords

En dehors du besoin de vérifiabilité du monde de la recherche, la reproductibilité est une qualité recherchée dans certains domaines de programmation [29] (**TODO: Lesquels?**).

Il existe donc depuis longtemps (**TODO: préciser + [Réf. nécessaire]**) des langages de programmation dits *fonctionnels*, qui, de manière plus ou moins stricte, limitent les effets de bords. Certains langages font également la distinction entre une fonction *pure* (sans effets de bord) et une fonction classique [30]. Certaines fonctions, plutôt appelées *procédures*, sont uniquement composées d’effet de bord et ne renvoie pas de valeur [31].

3.1.3 État dans le domaine de la robotique

En robotique, pour donner des ordres au matériel, on interagit beaucoup avec le monde extérieur (ordres et lecture d’état de servo-moteurs, flux vidéo d’une caméra, etc), souvent dans un langage plutôt bas-niveau, pour des questions de performance et de proximité abstractionnelle au matériel.

De fait, les langages employés sont communément C, C++ ou Python¹¹ [32], des langages bien plus impératifs que fonctionnels [33].

L’idée de s’affranchir d’effets de bords pour rendre les programmes dans la recherche en robotique reproductibles est donc plus utopique que réaliste.

¹¹Il arrive assez communément d’utiliser Python, un langage haut-niveau, mais c’est dans ce cas à but de prototypage. Le code contrôlant les moteurs est écrit dans un langage bas niveau, mais appelé par Python via FFI.

3.1.4 Reproductibilité de la compilation

Cependant, ce qui fait un programme n'est pas seulement son code: surtout dans des langages plus anciens sans gestion de dépendance intégrée au langage, les dépendances (bibliothèques) du programme, ainsi que l'environnement et les étapes de compilation de ce dernier, représentent également une partie considérable de la complexité du programme (par exemple, en C++, on utilise un outil générant des fichiers de configuration pour un autre outil qui à son tour configure le compilateur de C++¹² [34])

C'est cette partie que Nix, le gestionnaire de paquet, permet d'encapsuler et de rendre reproductible. Dans ce modèle, la compilation (et de manière plus générale la construction, ou *build*) du projet est la fonction que l'on veut rendre pure. L'entrée est le code source, et le résultat de la fonction est un binaire, qui ne doit dépendre que du code source.

$$\forall \text{src}, \text{bin}, \forall f \in \text{bin}^{\text{src}}, \forall (P_1, P_2) \in \text{src}^2, P_1 = P_2 \Rightarrow f(P_1) = f(P_2) \quad (24)$$

Ici, P_1 et P_2 sont deux itérations du code source (éléments de src) du programme. Si le code source est identique, les binaires résultants de la compilation (f) sont égaux, au sens de l'égalité bit à bit.

On a la proposition (1), avec $E = \text{src}$, l'ensemble des code source possibles pour un langage, et $F = \text{bin}$, l'ensemble des binaires exécutables

Nix ne peut pas garantir que le programme sera sans effets de bords au *runtime*, mais vise à le garantir au *build-time*.

3.2 Nix, le gestionnaire de paquets pur

3.2.1 Un DSL¹³ fonctionnel

Une autre caractéristique que l'on trouve souvent dans la famille de langages fonctionnels est l'omniprésence des *expressions*: la quasi-totalité des constructions syntaxiques forment des expressions valides, et peuvent donc servir de valeur

(TODO: améliorer exemple, ya des ternaires dans tt les langages...)

<pre>def g(x, y): if y == 5: x = 6 else: x = 8 return f(x)</pre>	<pre>let g x y = f (if y = 5 then 6 else 8)</pre>
Python (if et else sont des instructions)	OCaml (if et else forment une expression)

Afin de décrire les dépendances d'un programme, l'environnement de compilation, et les étapes pour le compiler (en somme, afin de définir le $f \in \text{bin}^{\text{src}}$), Nix comprend un langage d'expressions [35]. Un fichier `.nix` définit une fonction, que Nix sait exécuter pour compiler le code source.

Expression d'une fonction en Python	En Nix
<code>lambda f(a): a + 3</code>	<code>{ a }: a + 3</code>

Voici un exemple de définition d'un paquet Nix, appelée *dérivation* dans le jargon du langage:

```
{
  src-odri-masterboard-sdk,
```

¹²Il est ici question de CMake, qui génère des Makefile configurant GCC

¹³Domain-Specific Language

```

lib,
stdenv,
jrl-cmakemodules,
cmake,
python3Packages,
catch2_3,
}:

stdenv.mkDerivation {
  pname = "odri_master_board_sdk";
  version = "1.0.7";

  src = src-odri-masterboard-sdk;

  preConfigure = ''
    cd sdk/master_board_sdk
  '';

  doCheck = true;

  cmakeFlags = [
    (lib.cmakeBool "BUILD_PYTHON_INTERFACE" stdenv.hostPlatform.isLinux)
  ];

  nativeBuildInputs = [
    jrl-cmakemodules
    python3Packages.python
    cmake
  ];

  buildInputs = with python3Packages; [ numpy ];

  nativeCheckInputs = [ catch2_3 ];

  propagatedBuildInputs = with python3Packages; [ boost ];
}

```

La dérivation prend ici en entrée le code source (`src-odri-masterboard-sdk`), ainsi que des dépendances, que ce soit des fonctions relatives à Nix même (comme `stdenv.mkDerivation`) pour simplifier la définition de dérivation, ou des dépendances du programme, qu'elles servent à la compilation ou à son exécution (dans ce dernier cas de figure, les dépendances sont incluses ou dynamiquement liées dans le binaire final)

3.2.2 Un écosystème de dépendances

Afin de conserver la reproductibilité même lorsque l'on dépend de bibliothèques tierces, ces dépendances doivent également avoir une compilation reproductible: on déclare donc des dépendances à des paquets Nix, disponibles sur un registre centralisé, *Nixpkgs* [36].

Ainsi, écrire un paquet Nix pour son logiciel demande parfois d'écrire des paquets Nix pour les dépendances de notre projet, si celles-ci n'existent pas encore, et cela récursivement. On peut ensuite soumettre ces autres paquets à *Nixpkgs* [37] afin que d'autres puissent en dépendre sans les réécrire.

Pour ne pas avoir à compiler toutes les dépendances soi-même quand on dépend de paquets sur *Nixpkgs*, il existe un serveur de cache, qui propose des binaires des dépendances, Cachix [38]

3.2.3 Une compilation dans un environnement fixé

Certains aspects de l'environnement dans lequel l'on compile un programme peuvent faire varier le résultat final. Pour éviter cela, Nix limite au maximum les variations d'environnement. Par exemple,

la date du système est fixée au 0 UNIX (1er janvier 1990) pendant la compilation¹⁵: le programme compilé ne peut pas dépendre de la date à laquelle il a été compilé.

Quand le *sandboxing* est activé, Nix isole également le code source de tout accès au réseau, aux autres fichiers du système, ainsi que d'autres mesures, pour améliorer la reproductibilité [39]

Un complément utile: compiler en *CI*

Pour aller plus loin, on peut lancer la compilation du paquet Nix en *CI*¹⁶, c'est-à-dire sur un serveur distant, au lieu de compiler sur sa propre machine. On s'assure donc que l'état de notre machine de développement personnelle n'influe pas sur la compilation, puisque chaque compilation est lancée dans une machine virtuelle vierge [40].

3.3 NixOS, un système d'exploitation à configuration déclarative

Une fois le programme compilé avec ses dépendances, il est prêt à être transféré à l'ordinateur ou la carte de contrôle embarquée sur le robot.

Lorsqu'il y a un ordinateur embarqué, comme par exemple une Raspberry Pi [41], il faut choisir un OS sur lequel faire tourner le programme.

Là encore, un OS s'accompagne d'un amas considérable de configuration des différentes parties du système: accès au réseau, drivers, etc.

Sur les distributions Linux classiques tels que Ubuntu ou Debian, cette configuration est parfois stockée dans des fichiers, ou parfois retenue en mémoire, modifiée par l'exécution de commandes.

C'est un problème assez récurrent avec Linux de manière générale: d'un coup, le son ne marche plus, on passe ½h sur un forum à copier-coller des commandes dans un terminal, et le problème est réglé... jusqu'à ce qu'il survienne à nouveau après un redémarrage ou une réinstallation.

Ici, NixOS assure que toute modification de l'état du système est *déclarée* (d'où l'adjectif « déclaratif ») dans des fichiers de configurations, également écrits dans des fichiers `.nix` [42].

Ici encore, cela apporte un gain en terme de reproductibilité: l'état de configuration de l'OS sur lequel est déployé le programme du robot est, lui aussi, rendu reproductible.

3.4 Packaging Nix pour *gz-unitree*

Le packaging pour Nix de *gz-unitree* lui-même n'est pas très complexe: il s'agit d'un projet C++ / CMake standard [43].

Cependant, *gz-unitree* a deux dépendances principales:

- Gazebo lui-même, à travers `gz-sim`, `gz-sensors`, `gz-common`, `gz-plugin`, `gz-cmake`, etc.
- Le SDK d'Unitree, `unitree_sdk2`

En ce qui concerne le SDK d'Unitree, un paquet Nix a pu être écrit sans trop de soucis, la bibliothèque étant également un projet C++ standard:

```
{
  lib,
  stdenv,
  fetchFromGitHub,
  cmake,
```

¹⁵La date système n'est pas modifiée par Nix, mais il expose une date zéro au processus compilant le logiciel

¹⁶Continuous Integration, lit. intégration continue

```

    eigen,
  }:

stdenv.mkDerivation rec {
  pname = "unitree-sdk2";
  version = "2.0.0";

  src = fetchFromGitHub {
    owner = "unitreerobotics";
    repo = "unitree_sdk2";
    rev = version;
    hash = "sha256-r05zwhZW36+V0rIuTCr2HLf2R23csmnj33JFzUqz62Q=";
  };

  nativeBuildInputs = [ cmake ];

  buildInputs = [ eigen ];

  meta = {
    description = "Unitree robot sdk version 2. https://support.unitree.com/home/zh/developer";
    homepage = "https://github.com/unitreerobotics/unitree_sdk2";
    license = lib.licenses.bsd3;
    maintainers = with lib.maintainers; [ nim65s ];
    platforms = lib.platforms.unix;
  };
}

```

Par contre, en ce qui concerne Gazebo, la situation est plus complexe: étant un simulateur système, le projet est bien plus conséquent, et donc plus dur à packager.

Il existe plusieurs tentatives de packaging de Gazebo pour Nix:

- Un *overlay* ROS (Robot Operating System), qui inclut notamment Gazebo [44]
- Un package pour *Nixpkgs* [45], [46], [47]
- Un outil de génération de paquets Nix à partir de paquets ROS ou Gazebo, développé par Guilhem Saurel au sein de l'équipe Gepetto, *gazebo2nix* [48]

Au début du développement de *gz-unitree*, des essais d'utilisation des paquets Nix pour le développement et la compilation ont été réalisés, mais des erreurs subsistaient, en particulier avec Gazebo. Des efforts supplémentaires sont nécessaires pour empaqueter *gz-unitree*.

4 Étude du SDK d'Unitree et du bridge SDK \leftrightarrow MuJoCo

4.1 Le robot *H1v2* d'Unitree

Le *H1v2* est un modèle de robot humanoïde créé par la société Unitree.

Il possède plus de 26 degrés de liberté, dont

- 6 dans chaque jambe (3 à la hanche, 2 au talon et un au genou),
- 7 dans chaque bras (3 à l'épaule, 3 au poignet et un au coude) [49]

Unitree met à disposition du public un *SDK*¹⁷ permettant de le contrôler, `unitree_sdk2` [50]

¹⁷Kit de développement logiciel (Software Development Kit)

4.2 Canaux DDS

Pour communiquer avec le robot via le réseau, Unitree utilise CycloneDDS, une implémentation par Oracle du standard DDS¹⁸ [51]. DDS est une technologie de communication bidirectionnelle¹⁹ en temps réel, standardisée par l'Object Management Group, OMG [52]. Les messages sont envoyés sur le réseau via UDP et IP.

Les données contenues dans chacun des messages sont spécifiées via un autre format, IDL, également standardisé par l'OMG [53].

L'intérêt d'un format indépendant du langage de programmation est que l'on peut générer du code décrivant ces données dans plusieurs langages de programmation, ce que fait Unitree en distribuant un SDK en C++ et en Python.

Par exemple, les messages permettant de contrôler les moteurs du H1v2 sont définis ainsi

```
struct MotorCmd
{
    uint8 mode;
    float q;
    float dq;
    float tau;
    float kp;
    float kd;
    unsigned long reserve;
};

struct Cmd
{
    uint8 mode_pr;
    uint8 mode_machine;
    MotorCmd motor_cmd[35];
    unsigned long reserve[4];
    unsigned long crc;
};
```

Liste 1. – `LowCmd.idl`, traduit depuis sa conversion en C++ [54]

DDS groupe les messages dans des *topics*. Les messages sont échangés sur un topic de la manière suivante

Lecture En s'abonnant au topic, on reçoit en temps réel les messages qui sont envoyés dessus

Écriture En publiant des messages sur le topic, on les rend disponibles aux abonnés

CycloneDDS est capable d'un débit d'environ 1 GBs^{-1} , pour des messages d'environ 1 kB chacun [55]. On remarque, en pratique, des tailles de message entre 0.9 kB et 1.3 kB dans le cas des échanges commandes/état avec le robot.

Et enfin, les topics peuvent être isolés d'autres topics via des *domains*, identifiés par un numéro. Deux topics portant le même nom restent isolés si ils sont sur deux domaines différents.

¹⁸pour Data Distribution Service

¹⁹dite « *pub-sub* » pour *publish/subscribe*

4.3 Une base de code partiellement open-source

Le code source du SDK d'Unitree est disponible sur Github [56]. Cependant, le dépôt git comprend des fichiers binaires déjà compilés:

```

lib
├── aarch64
│   └── libunitree_sdk2.a
└── x86_64
    └── libunitree_sdk2.a
thirdparty
├── CMakeLists.txt
├── include
│   └── ...
└── lib
    ├── aarch64
    │   ├── libddsc.so
    │   ├── libddsc.so.0 -> libddsc.so
    │   ├── libddscxx.so
    │   └── libddscxx.so.0 -> libddscxx.so
    └── x86_64
        ├── libddsc.so
        ├── libddsc.so.0 -> libddsc.so
        ├── libddscxx.so
        └── libddscxx.so.0 -> libddscxx.so

```

Liste 2. – Résultat de `tree lib/ thirdparty/` dans le dépôt git

Compiler le SDK nécessite l'existence de ces fichiers binaires:

```

63 # Create imported target unitree_sdk2
64 add_library(unitree_sdk2 STATIC IMPORTED GLOBAL)
65 set_target_properties(unitree_sdk2 PROPERTIES
66     IMPORTED_LOCATION "${_IMPORT_PREFIX}/lib/libunitree_sdk2.a
67     INTERFACE_INCLUDE_DIRECTORIES "${_IMPORT_PREFIX}/include;${_IMPORT_PREFIX}/
        include"
68     INTERFACE_LINK_LIBRARIES "ddsc;ddscxx;Threads::Threads"
69     LINKER_LANGUAGE CXX

```

Liste 3. – Extrait de `cmake/unitree_sdk2Targets.cmake` [50]

Ici est défini, via `set_target_properties(... IMPORTED_LOCATION)`, le chemin d'une bibliothèque à lier avec la bibliothèque finale [57]. Ici, c'est un des fichiers pré-compilés que l'on lie.

On confirme cette nécessité en lançant `mkdir build && cd build && cmake ..` après avoir supprimé le répertoire `lib/` :

```

-- The C compiler identification is GNU 13.3.0
-- The CXX compiler identification is GNU 13.3.0
-- Detecting C compiler ABI info
-- Detecting C compiler ABI info - done
-- Check for working C compiler: /usr/bin/cc - skipped
-- Detecting C compile features
-- Detecting C compile features - done
-- Detecting CXX compiler ABI info

```

```
-- Detecting CXX compiler ABI info - done
-- Check for working CXX compiler: /usr/bin/c++ - skipped
-- Detecting CXX compile features
-- Detecting CXX compile features - done
-- Setting build type to 'Release' as none was specified.
-- Current system architecture: x86_64
-- Performing Test CMAKE_HAVE_LIBC_PTHREAD
-- Performing Test CMAKE_HAVE_LIBC_PTHREAD - Success
-- Found Threads: TRUE
-- Importing: /home/glebihan/playground/unitree_sdk2/thirdparty/lib/x86_64/
libddsc.so
-- Importing: /home/glebihan/playground/unitree_sdk2/thirdparty/lib/x86_64/
libddscxx.so
```

```
CMake Error at CMakeLists.txt:42 (message):
Unitree SDK library for the architecture is not found
```

```
-- Configuring incomplete, errors occurred!
```

Les logs montrent aussi que les recettes de compilation dépendent de versions précompilées de LibDDSC et LibDDSCXX, dont le code source semble cependant être fourni avec *unitree_sdk2*:

```
thirdparty/include/dds
├─ config.h
├─ ddsc
│   ├── dds_basic_types.h
│   ├── dds_data_allocator.h
│   ├── dds_internal_api.h
│   ├── dds_loan_api.h
│   ├── dds_opcodes.h
│   └─ dds_public_alloc.h
...
```

Ces particularités laissent planer quelques doutes sur la nature open-source du code: ces binaires requis sont-ils seulement présent pour améliorer l'expérience développeur en accélérant la compilation, ou « cachent »-ils du code non public?

Ces constats ont motivé une première tentative de décompilation de ces *libunitree_sdk2.a* pour comprendre le fonctionnement du SDK, via *Ghidra* [58].

Cependant, la découverte de l'existence d'un bridge officiel SDK \Leftarrow Mujoco [59] a rendu l'exploration de cette piste non nécessaire.

4.4 Un autre bridge existant: unitree_mujoco

Unitree propose un bridge officiel pour utiliser son SDK avec Mujoco.

Le fonctionnement d'un bridge est au final assez similaire, quelque soit le simulateur pour lequel on l'écrit: il s'agit d'envoyer l'état du robot au simulateur, et de réagir quand le simulateur envoie des ordres de commandes.

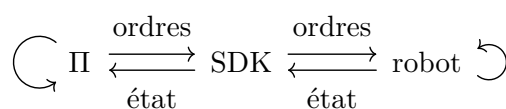


Fig. 9. – Fonctionnement usuel du SDK

Un bridge se substitue au robot physique, interceptant les ordres du SDK et les traduisants en des appels de fonctions provenant de l'API du simulateur, et symétriquement pour les envois d'états au SDK. On peut apparenter le fonctionnement d'un bridge à celui d'une attaque informatique de type « Man in the Middle » (MitM).

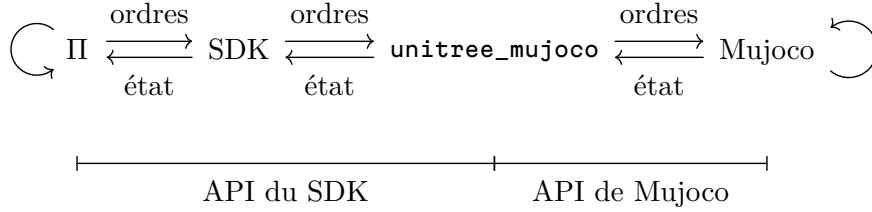


Fig. 10. – Fonctionnement via `unitree_mujoco` du SDK

Le but est de faire la même chose avec notre propre bridge. Le code du bridge Mujoco existant est utile car un bridge, se situant par définition à la frontière entre deux APIs, fait usage des deux APIs.

Écrire un bridge Gazebo pour le même SDK implique donc de changer « API de Mujoco » par « API de Gazebo », mais le code faisant usage du SDK d'Unitree reste le même.

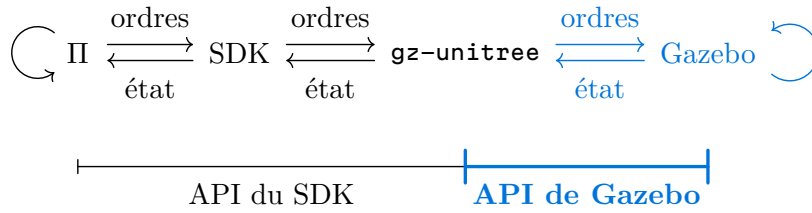


Fig. 11. – Fonctionnement via `gz-unitree` du SDK

Le bridge de Mujoco fonctionne en interceptant les messages sur le canal `rt/lowcmd` et en envoyant dans le canal `rt/lowstate`, qui correspondent respectivement aux commandes envoyées au robot et à l'état (angles des joints, moteurs, valeurs des capteurs, etc) reçu depuis le robot.

Le `low` indique que ce sont des messages bas-niveau. Par exemple, `rt/lowcmd` correspond directement à des ordres en valeurs de couple pour les moteurs, au lieu d'envoyer des ordres plus évolués, tels que « se déplacer de x mètres en avant » [60]

Les ordres dans `rt/lowcmd` sont ensuite traduits en appels de fonctions de Mujoco pour mettre à jour l'état du robot simulé, et de messages `rt/lowstate` sont créés à partir des données fournies par Mujoco.

Étant donné le modèle *pub/sub* de DDS, on parle de *pub(lication)* de message, et de *sub(scription)*²⁰ aux messages d'un canal (pour les recevoir).

²⁰abonnement

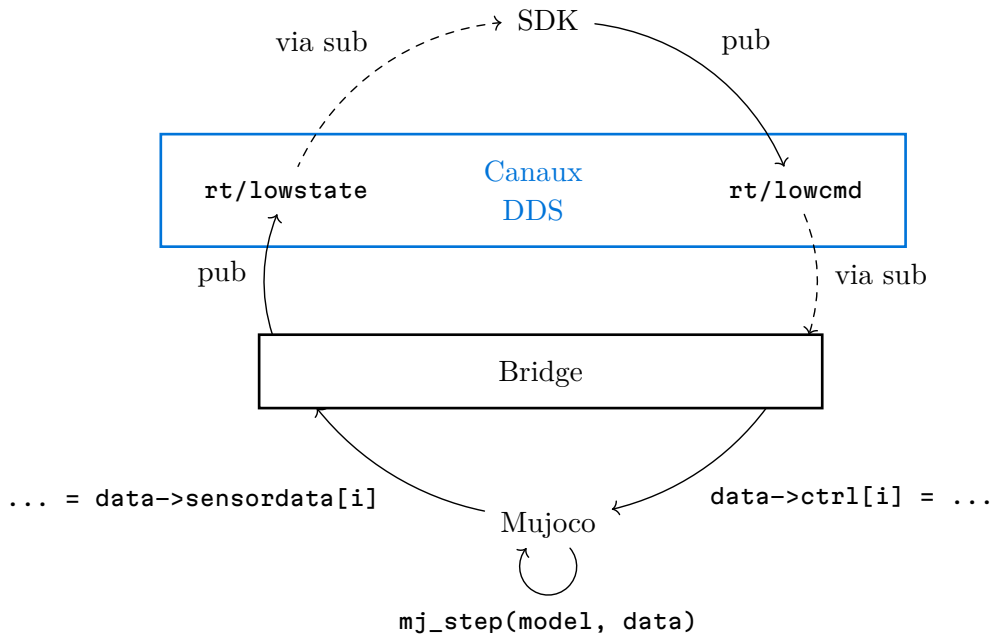


Fig. 12. – Cycle de vie de la simulation avec le bridge pour Mujoco

Le but est donc de reproduire un cycle de vie équivalent, mais en remplaçant la partie spécifique à Mujoco par une partie adaptée à Gazebo.

5 Développement du bridge SDK \leftrightarrow Gazebo

En se basant sur `unitree_mujoco`, il a donc été possible de réaliser un bridge pour Gazebo.

5.1 Établissement du contact

Une première tentative a été de suivre la documentation de CycloneDDS [61] pour écouter sur le canal `rt/lowcmd`, en récupérant les définitions IDL des messages, disponibles sur le dépôt `unitree_ros2`²¹ [62]

Malheureusement, cette solution s’est avérée infructueuse, à cause d’une erreur sur les domaines DDS utilisés (ce qui sera compris plus tard).

On change d’approche, en préférant plutôt utiliser les abstractions fournies par le SDK de Unitree (cf Chapitre 5.5 et Chapitre 5.6.2)

Enfin, si un pare-feu est actif, il faut autoriser le trafic UDP l’intervalle d’adresses IP `224.0.0.0/4`. Par exemple, avec `ufw` [63]

```
sudo ufw allow in proto udp from 224.0.0.0/4
sudo ufw allow in proto udp to 224.0.0.0/4
```

²¹`unitree_mujoco` n’avait pas encore été découvert

Pour arriver à ces solutions, *Wireshark* [64] s'est avéré utile, étant capable d'inspecter du trafic RTPS²²,

C'est notamment grâce à ce traçage des paquets que le problème de domaine DDS a été découvert: notre *subscriber* DDS était réglé sur le domaine anonyme (ID aléatoire représenté par un 0 lors de la configuration) alors que le SDK d'Unitree communique sur le domaine d'ID 1. C'est aussi Wireshark qui nous a permis de voir quels étaient les types IDL utilisés pour les messages.

	participant_idx	domain_id	Info
san	120	4294967295	INFO_TS, DATA
an		1	INFO_TS, DATA
san	120	4294967295	INFO_TS, DATA
an		1	INFO_TS, DATA
san	120	4294967295	INFO_TS, DATA
an		1	INFO_TS, DATA
san	120	4294967295	INFO_TS, DATA
an		1	INFO_TS, DATA
san	120	4294967295	INFO_TS, DATA
an		1	INFO_TS, DATA
san	120	4294967295	INFO_TS, DATA
an		1	INFO_TS, DATA
san	120	4294967295	INFO_TS, DATA
an		1	INFO_TS, DATA
san	120	4294967295	INFO_TS, DATA
an		1	INFO_TS, DATA
san	120	4294967295	INFO_TS, DATA
an		1	INFO_TS, DATA
san	120	4294967295	INFO_TS, DATA
an		1	INFO_TS, DATA

```

> writerEntityId: 0x00000403 (Application-defined v
- [Topic Information (from Discovery)]
  [typeName: unitree_go::msg::dds_::LowState_]
  [topic: rt/lowstate]
  [DCSPublicationData In: 29]
writerSeqNumber: 204954
> serializedData

```

Fig. 13. – *Wireshark* permet de visualiser des méta-données sur les paquets RTPS

Voici une trace wireshark d'un échange usuel entre commandes (*rt/lowcmd*) et états (*rt/lowstate*)

No.	Time	Source	Destination	Protocol	Length	Info
16	57.834016201	127.0.0.1	239.255.0.1	RTPS	454	INFO_TS, DATA(p)
17	64.000946758	127.0.0.1	239.255.0.1	RTPS	458	INFO_TS, DATA(p)
18	65.834206105	127.0.0.1	239.255.0.1	RTPS	454	INFO_TS, DATA(p)
19	72.001017927	127.0.0.1	239.255.0.1	RTPS	458	INFO_TS, DATA(p)
20	73.834313058	127.0.0.1	239.255.0.1	RTPS	454	INFO_TS, DATA(p)
21	80.001169033	127.0.0.1	239.255.0.1	RTPS	458	INFO_TS, DATA(p)
22	81.834316003	127.0.0.1	239.255.0.1	RTPS	454	INFO_TS, DATA(p)
23	82.586335139	127.0.0.1	239.255.0.1	RTPS	450	INFO_TS, DATA(p)
24	82.587033958	127.0.0.1	127.0.0.1	RTPS	470	INFO_DST, INFO_TS, DATA(p)
25	82.587840467	127.0.0.1	127.0.0.1	RTPS	142	INFO_DST, HEARTBEAT, HEARTBEAT
26	82.587905742	127.0.0.1	127.0.0.1	RTPS	142	INFO_DST, HEARTBEAT, HEARTBEAT
27	82.587938692	127.0.0.1	127.0.0.1	RTPS	110	INFO_DST, HEARTBEAT
28	82.588109331	127.0.0.1	127.0.0.1	RTPS	230	INFO_DST, ACKNACK, ACKNACK, ACKNACK, ACKNACK
29	82.588411479	127.0.0.1	127.0.0.1	RTPS	1202	INFO_DST, INFO_TS, DATA(w) -> rt/lowstate, INFO_TS, DATA(w) -> rt/sportmodest
30	82.588444898	127.0.0.1	127.0.0.1	RTPS	174	INFO_DST, HEARTBEAT, HEARTBEAT, HEARTBEAT
31	82.588467679	127.0.0.1	127.0.0.1	RTPS	670	INFO_DST, INFO_TS, DATA(r) -> rt/lowcmd, HEARTBEAT, INFO_TS, DATA(m), HEARTBEAT
32	82.588495680	127.0.0.1	127.0.0.1	RTPS	142	INFO_DST, HEARTBEAT, HEARTBEAT
33	82.58852936	127.0.0.1	127.0.0.1	RTPS	222	INFO_DST, ACKNACK, ACKNACK, ACKNACK, ACKNACK
34	82.588733831	127.0.0.1	127.0.0.1	RTPS	226	INFO_DST, ACKNACK, ACKNACK, ACKNACK, INFO_TS, DATA(m)
35	82.588762137	127.0.0.1	127.0.0.1	RTPS	110	INFO_DST, HEARTBEAT
36	82.588895984	127.0.0.1	127.0.0.1	RTPS	106	INFO_DST, ACKNACK
37	82.686245575	127.0.0.1	239.255.0.1	RTPS	450	INFO_TS, DATA(p)
38	83.587177329	127.0.0.1	127.0.0.1	RTPS	470	INFO_DST, INFO_TS, DATA(p)
39	83.649444949	127.0.0.1	127.0.0.1	RTPS	422	INFO_TS, DATA(w) -> rt/lowcmd
40	83.749594668	127.0.0.1	127.0.0.1	RTPS	94	HEARTBEAT
41	83.749791329	127.0.0.1	127.0.0.1	RTPS	106	INFO_DST, ACKNACK
42	83.752510795	127.0.0.1	127.0.0.1	RTPS	546	INFO_TS, DATA(r) -> rt/lowstate
43	83.753467322	127.0.0.1	127.0.0.1	RTPS	942	INFO_TS, DATA -> rt/lowcmd, HEARTBEAT -> rt/lowcmd
44	83.753646127	127.0.0.1	127.0.0.1	RTPS	1310	INFO_TS, DATA -> rt/lowstate, HEARTBEAT -> rt/lowstate
45	83.755133498	127.0.0.1	127.0.0.1	RTPS	942	INFO_TS, DATA -> rt/lowcmd, HEARTBEAT -> rt/lowcmd
46	83.755722822	127.0.0.1	127.0.0.1	RTPS	1310	INFO_TS, DATA -> rt/lowstate, HEARTBEAT -> rt/lowstate
47	83.757170036	127.0.0.1	127.0.0.1	RTPS	942	INFO_TS, DATA -> rt/lowcmd, HEARTBEAT -> rt/lowcmd
48	83.757594204	127.0.0.1	127.0.0.1	RTPS	1310	INFO_TS, DATA -> rt/lowstate, HEARTBEAT -> rt/lowstate
49	83.759172813	127.0.0.1	127.0.0.1	RTPS	942	INFO_TS, DATA -> rt/lowcmd, HEARTBEAT -> rt/lowcmd
50	83.759595304	127.0.0.1	127.0.0.1	RTPS	1310	INFO_TS, DATA -> rt/lowstate, HEARTBEAT -> rt/lowstate
51	83.761212408	127.0.0.1	127.0.0.1	RTPS	942	INFO_TS, DATA -> rt/lowcmd, HEARTBEAT -> rt/lowcmd
52	83.761595204	127.0.0.1	127.0.0.1	RTPS	1310	INFO_TS, DATA -> rt/lowstate, HEARTBEAT -> rt/lowstate
53	83.763178110	127.0.0.1	127.0.0.1	RTPS	942	INFO_TS, DATA -> rt/lowcmd, HEARTBEAT -> rt/lowcmd
54	83.763595289	127.0.0.1	127.0.0.1	RTPS	1310	INFO_TS, DATA -> rt/lowstate, HEARTBEAT -> rt/lowstate
55	83.765217646	127.0.0.1	127.0.0.1	RTPS	942	INFO_TS, DATA -> rt/lowcmd, HEARTBEAT -> rt/lowcmd
56	83.765595726	127.0.0.1	127.0.0.1	RTPS	1310	INFO_TS, DATA -> rt/lowstate, HEARTBEAT -> rt/lowstate
57	83.767206486	127.0.0.1	127.0.0.1	RTPS	942	INFO_TS, DATA -> rt/lowcmd, HEARTBEAT -> rt/lowcmd
58	83.767595863	127.0.0.1	127.0.0.1	RTPS	1310	INFO_TS, DATA -> rt/lowstate, HEARTBEAT -> rt/lowstate
59	83.769219527	127.0.0.1	127.0.0.1	RTPS	942	INFO_TS, DATA -> rt/lowcmd, HEARTBEAT -> rt/lowcmd
60	83.769595878	127.0.0.1	127.0.0.1	RTPS	1310	INFO_TS, DATA -> rt/lowstate, HEARTBEAT -> rt/lowstate
61	83.771210650	127.0.0.1	127.0.0.1	RTPS	942	INFO_TS, DATA -> rt/lowcmd, HEARTBEAT -> rt/lowcmd

Fig. 14. – Trace de paquets RTPS sur *Wireshark*

²²Le protocole sur lequel est construit DDS [52]

5.2 Installation du plugin dans Gazebo

Un *system plugin* Gazebo consiste en la définition d'une classe héritant de `gz::sim::System` et d'interfaces permettant notamment d'exécuter notre code avant ou après un pas de temps du simulateur (avec `gz::sim::ISystem{Pre,Post}Update`)

```
#include <gz/sim/System.hh>
namespace gz_unitree
{
    class UnitreePlugin :
        public gz::sim::System,
        public gz::sim::ISystemPreUpdate
    {
    public:
        UnitreePlugin();
    public:
        ~UnitreePlugin() override;
    public:
        void PreUpdate(const gz::sim::UpdateInfo &_info,
                       gz::sim::EntityComponentManager &ecm) override;
    };
}
```

Liste 4. – Fichier header pour le plugin

Il faut ensuite implémenter la classe puis appeler une macro ajoutant le plugin à Gazebo

```
#include <gz/plugin/Register.hh>

... // class implementation

GZ_ADD_PLUGIN(
    UnitreePlugin,
    gz::sim::System,
    UnitreePlugin::ISystemPreUpdate)
```

Enfin, on active le plugin en le référençant dans le fichier SDF [65], qui décrit l'environnement des simulations (objets, éclairage, etc)

```
<sdf version='1.11'>
<world name="default">
  <plugin filename="gz-unitree" name="gz_unitree::UnitreePlugin">
  </plugin>
</world>
<model name='h1_description'>
  <link name='pelvis'>
    <inertial>
    ...
```

Avec `filename` le chemin vers le plugin compilé, qui sera cherché dans les répertoires spécifiés par `GZ_SIM_SYSTEM_PLUGIN_PATH` [66], [67].

5.3 Architecture du plugin

Le plugin consiste en trois parties distinctes:

- Le « branchement » dans les phases de simulation de Gazebo, par l'implémentation de méthodes de `gz::sim::System`
- L'interaction avec les canaux DDS du SDK d'Unitree

- Les données et méthodes internes au plugin

En plus de cela, il y a bien évidemment la politique de contrôle II, qui interagit avec le robot (qu'il soit réel, ou simulé) via le SDK, et donc via les canaux DDS.

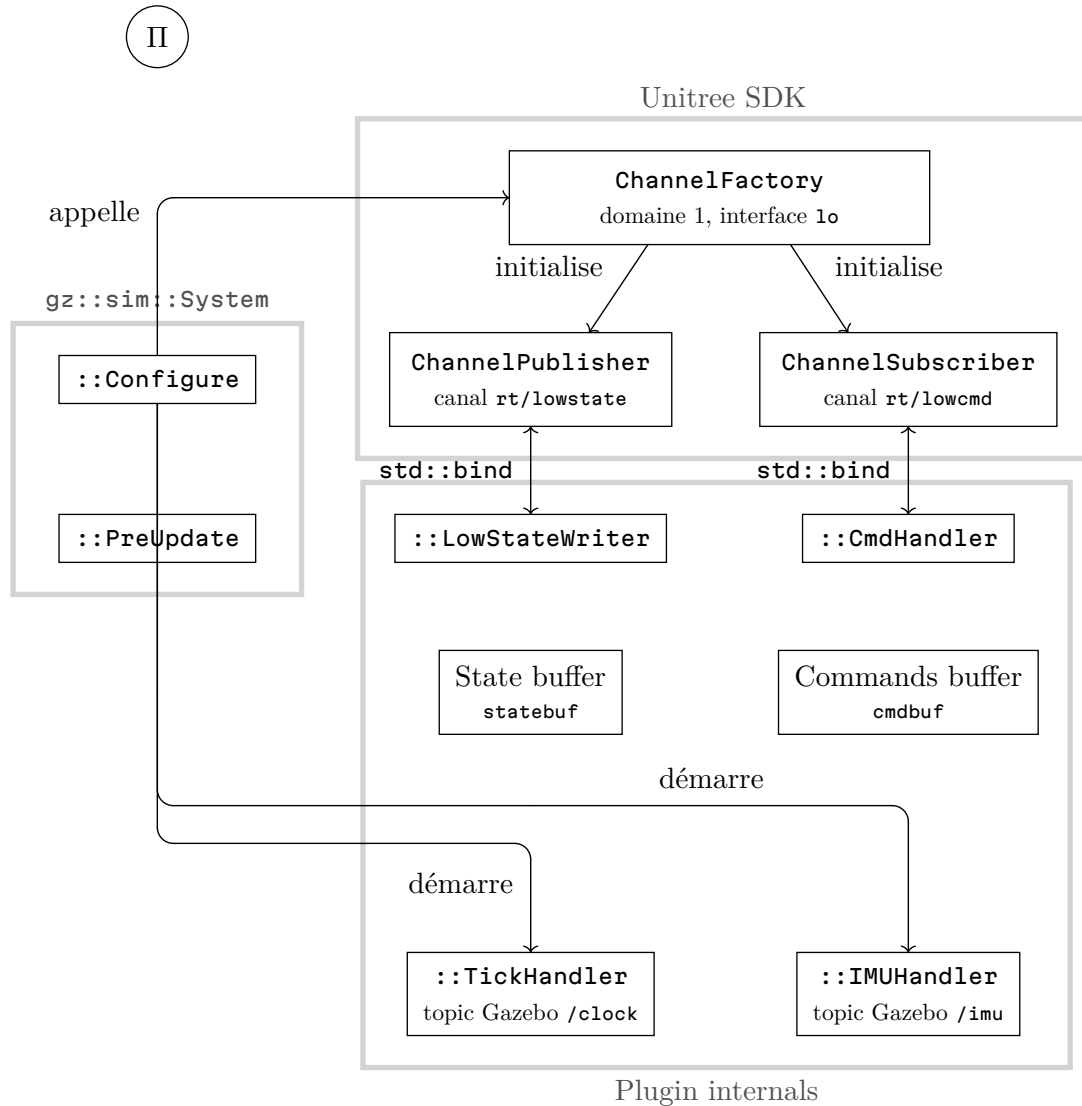


Fig. 15. – Phase d'initialisation du plugin

On commence par instancier un contrôleur dans le domaine DDS n°1, sur l'interface réseau 1o²³

On lui associe:

- Un *publisher*, chargé d'envoyer périodiquement des messages sur **rt/lowstate** en appelant la méthode **LowStateWriter**
- Un *subscriber*, chargé d'appeler la méthode **CmdHandler** avec chaque message arrivant sur **rt/lowcmd**.

²³interface dite « loopback », qui est locale à l'ordinateur: ici, le simulateur et la politique de contrôle tournent sur la même machine, donc les messages DDS n'ont pas besoin de « sortir » de celle-ci

On démarre aussi deux autres *subscribers*, qui sont eux chargés d'écouter des messages sur les topics Gazebo `/clock` et `/imu`, ce qui permet de récupérer le tick de simulation et les valeurs du capteur IMU²⁴, que l'on a préalablement fixé au modèle du robot en le déclarant dans le fichier SDF chargé par Gazebo. Le capteur IMU donne des informations spatiales importantes sur la position et la vitesse du robot.

Les topics Gazebo sont un autre moyen de communication inter-processus asynchrone par pub/sub, similaire à DDS [68]. Gazebo utilise Protobuf pour définir les types des messages [69], [70], qui joue ici le même rôle qu'IDL dans DDS. Les topics Gazebo sont basés sur un réseau décentralisé de nœuds, chaque nœud pouvant indépendamment publier et/ou recevoir des messages.

Cette initialisation est faite à la phase de configuration du plugin par Gazebo, via l'implémentation de la méthode `::Configure` du plugin.

En pratique, on utilise `std::bind` [71] pour fixer l'instance d'`UnitreePlugin` et ainsi passer des méthodes de la classe comme des simples fonctions.

```
auto subscriber = ChannelSubscriberPtr<LowCmd_>(
    new ChannelSubscriber<LowCmd_>("rt/lowcmd")
);

auto handler = std::bind(
    &UnitreePlugin::CmdHandler,
    this,
    std::placeholders::_1
)

subscriber->InitChannel(handler, 1);

auto publisher = ChannelPublisherPtr<LowState_>(
    new ChannelPublisher<LowState_>("rt/lowstate")
);

publisher->InitChannel();

this->publisher_thread = CreateRecurrentThreadEx(
    "low_state_writer",
    UT_CPU_ID_NONE,
    500,
    &UnitreePlugin::LowStateWriter,
    this
);
```

Liste 5. – Création d'un *subscriber* à `rt/lowcmd` et d'un *publisher* sur `rt/lowstate` dans `UnitreePlugin::Configure`

5.4 Calcul des nouveaux couples des moteurs

Pour appliquer une commande à un moteur, on calcule la force effective que le moteur doit appliquer:

$$\tau = \underbrace{\tau_{\text{ff}}}_{\text{stabilité}} + \underbrace{K_p \Delta q}_{\text{proportionnelle}} + \underbrace{K_d \Delta \dot{q}}_{\text{dérivée}} \quad (25)$$

Avec

- τ pour *torque*, le couple à donner au moteur
- τ_{ff} le τ « feed-forward ». Particulièrement utile pour les robots humanoïdes qui doivent rester debout. Dans ces cas, on parle parfois de *gravity compensation part* [72]
- Δq écart d'angle de rotation du moteur entre la consigne et l'état actuel
- $\Delta \dot{q}$ vitesse de changement de la consigne²⁵
- K_p prépondérance de la partie proportionnelle
- K_p prépondérance de la partie dérivée

²⁴Inertial Measurement Unit, appelée « Centrale intertielle » en français

²⁵

On a bien $\frac{d \Delta q}{dt} = \Delta \dot{q}$ par linéarité de la dérivation temporelle:

$$\frac{d \Delta q}{dt} = \frac{d q_{\text{new}} - q_{\text{old}}}{dt} = \frac{d q_{\text{new}}}{dt} - \frac{d q_{\text{old}}}{dt} = \Delta \frac{d q}{dt} = \Delta \dot{q} \quad (26)$$

Cette équation met à jour τ pour rapprocher l'état actuel du moteur de la nouvelle consigne, en prenant en compte

- L'erreur sur l'angle Δq (partie proportionnelle).
- L'erreur sur la vitesse de changement de Δq (partie dérivative). Cette prise en compte de la vitesse permet de lisser les changements appliqués aux moteurs.
- Un couple dit de *feed-forward*, τ_{ff} , qui permet le maintien du robot à un état stable. On pourrait le déterminer en lançant une première simulation, avec pour objectif le maintien debout. Une fois la stabilité atteinte, on relève les couples des moteurs. Intuitivement, on peut voir τ_{ff} comme un manière de s'affranchir de la partie « maintien debout » dans l'expression de la commande, similairement à la mise à zéro (« tarer ») d'une balance.

On contrôle la prépondérance des deux erreurs dans le calcul de la nouvelle consigne grâce à deux coefficients, K_p et K_d .

5.5 rt/lowcmd

On trouve dans les messages **rt/lowcmd** les champs nécessaires à au calcul de τ [73] comme décrit précédemment:

Champ	Type	Description
crc	\mathbb{N}	Somme de contrôle CRC32, pourrait éventuellement servir à éviter de prendre en compte des messages corrompus
motor_cmd	struct. ³⁵	Paramètres de commande pour chacun des 35 moteurs
.q, .dq, .tau, .kp et .kd	\mathbb{R}	Respectivement q , \dot{q} , τ_{ff} , K_p et K_d

Cette équation se rapproche des modèles de type PID (*proportional-integrative-derivative*) [74], avec le terme intégratif remplacé par τ_{ff} , ce qui en fait une expression plus adaptée pour les politiques avec des mouvements non-brusques: le terme intégratif apporte une capacité d'instabilité qui complexifie l'entraînement [Réf. nécessaire]

```
// Avec i l'indice du moteur
auto force = cmdbuf->tau_ff.at(i) + // tau_ff
             cmdbuf->kp.at(i) * ( // K_p
             cmdbuf->q_target.at(i) - lowstate.motor_state().at(i).q() // Delta q
             ) +
             cmdbuf->kd.at(i) * ( // K_d
             cmdbuf->dq_target.at(i) - lowstate.motor_state().at(i).dq() // Delta q.
             );

std::vector<double> torque = {force};
joint.SetForce(ecm, torque);
```

Liste 6. – Implémentation de la mise à jour de τ

Lorsqu'un message, publié par Π (1A) et contenant des ordres pour les moteurs, arrive sur **rt/lowcmd**, **ChannelSubscriber** appelle **::CmdHandler** (2, 3), et modifie un *buffer* (4) contenant la dernière commande reçue. Ensuite, Gazebo démarre un nouveau pas de simulation. Avant de faire ce pas, il appelle la méthode **::PreUpdate** sur notre plugin, qui vient chercher la commande stockée dans le *buffer* (1B), et applique cette commande sur le modèle du robot, animé par le simulateur.

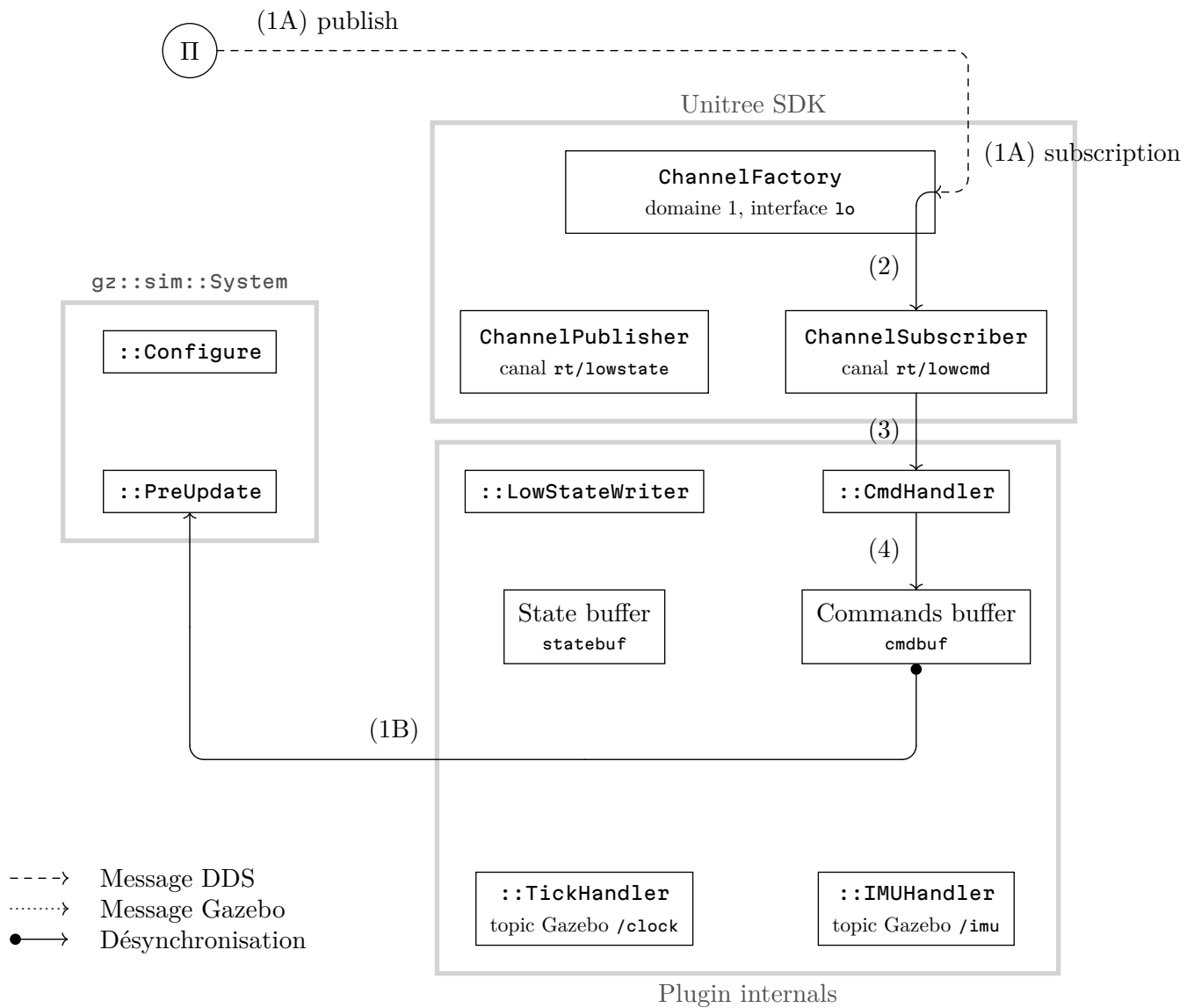


Fig. 16. – Phase de réception des commandes

On notera que (1B) s'exécute *parallèlement* au reste des étapes: la boucle de simulation de Gazebo est indépendante de la boucle de mise à jour de la politique.

Si `::PreUpdate` est plus fréquente Le simulateur appliquera plusieurs fois la même commande, le buffer `cmdbuf` n'ayant pas été modifié.

Si `::PreUpdate` est moins fréquente Certaines commandes seront ignorées par Gazebo, qui ne verra pas la valeur du buffer `statebuf` avant qu'il change de nouveau.

5.6 rt/lowstate

5.6.1 Construction d'un message rt/lowstate

La documentation d'Unitree liste l'ensemble des champs disponibles dans un message `rt/lowstate`, c'est-à-dire l'ensemble des données que l'on doit récupérer afin de construire nos messages d'état [75]:

Champ	Type	Description	Où récupérer la valeur
<code>version</code>	\mathbb{N}^2	<i>Non documenté</i>	<i>Laissé vide</i>
<code>mode_pr</code>	$\{0, 1\}$	Défini sur 0 par défaut	0
<code>mode_machine</code>	$\{4, 6\}$	Défini sur 6 par défaut	6
<code>tick</code>	\mathbb{N} (ms)	<i>Non documenté</i> , probablement le temps écoulé depuis le début de la simulation	Messages <code>gz::msgs::Clock</code> sur le topic Gazebo <code>/clock</code>
<code>wireless_remote</code>	$\{0, 1\}$ ⁴⁰	<i>Non documenté</i>	<i>Laissé vide</i>
<code>reserve</code>	\mathbb{N}^4	<i>Non documenté</i>	<i>Laissé vide</i>
<code>crc</code>	\mathbb{N}	Somme de contrôle du message, utilisant <i>CRC32</i> .	Implémentation de CRC32 par Unitree ²⁶
<code>imu_state...</code>	struct.	Valeurs des capteurs intertiels du robot	Messages <code>gz::msgs::IMU</code> sur le topic Gazebo <code>/imu</code>
<code>.quaternion</code>	\mathbb{R}^4	Posture dans l'espace du robot, dans l'ordre (w, x, y, z)	<code>w</code> , <code>x</code> , <code>y</code> et <code>z</code> sur <code>.orientation()</code>
<code>.rpy</code>	\mathbb{R}^3	Angle d'Euler du robot, dans l'ordre (r, p, y)	<code>.linear_acceleration()</code>
<code>.gyroscope</code>	\mathbb{R}^3	Gyroscope	En utilisant les valeurs de <code>.orientation()</code> : $\begin{bmatrix} \text{atan}_2(2(wx + yz), 1 - 2(x^2 + y^2)) \\ \text{asin}(2(wy - zx)) \\ \text{atan}_2(2(wz + xy), 1 - 2(y^2 + z^2)) \end{bmatrix}$
<code>.accelerometer</code>	\mathbb{R}^3	Accéléromètre	<code>.angular_velocity()</code>
<code>motor_state...</code>	struct. ³⁵	Etat de chaque moteur	<code>gz::sim::Model(...)->joints</code>
<code>.mode</code>	$\{0, 1\}$	Modes de contrôle pour le moteur électrique. 0 pour « Brake » et 1 pour « FOC ²⁷ »	0
<code>.q</code>	\mathbb{R} (rad)	Angle de rotation du moteur	<code>.Position()</code>
<code>.dq</code>	\mathbb{R} (rad · s ⁻¹)	Angle de rotation du moteur	<code>.Velocity()</code>
<code>.ddq</code>	\mathbb{R} (rad · s ⁻²)	Angle de rotation du moteur	<i>Laissé vide</i> ²⁸
<code>.tau_est</code>	\mathbb{R} (N · m)	Estimation du couple exercé par le moteur	<i>Laissé vide</i> ²⁸

²⁶Une implémentation ad-hoc existe dans le code source de `unitree_sdk2` [76] et de `unitree_mujoco` [77], elle est également donnée en annexe B

²⁷Field-Oriented Control

²⁸Tant que nos politiques n'ont pas besoin de ces champs, le SDK semble fonctionner avec des valeurs vides

5.6.2 Émission de l'état

Avant de démarrer un nouveau pas de simulation, la méthode `::PreUpdate` vient mettre à jour l'état du robot simulé en modifiant le *State buffer* interne au plugin (1A). Gazebo envoie également le nouveau tick de simulation (1C) et les valeurs du capteur IMU (1D) dans leurs topics respectifs.

Le `LowStateWriter` vient lire le *State buffer* (1B) pour publier l'état sur le canal DDS (2, 3) qui est ensuite lu par `II` (4), qui possède une subscription sur `rt/lowstate`, via sa propre utilisation du SDK d'Unitree.

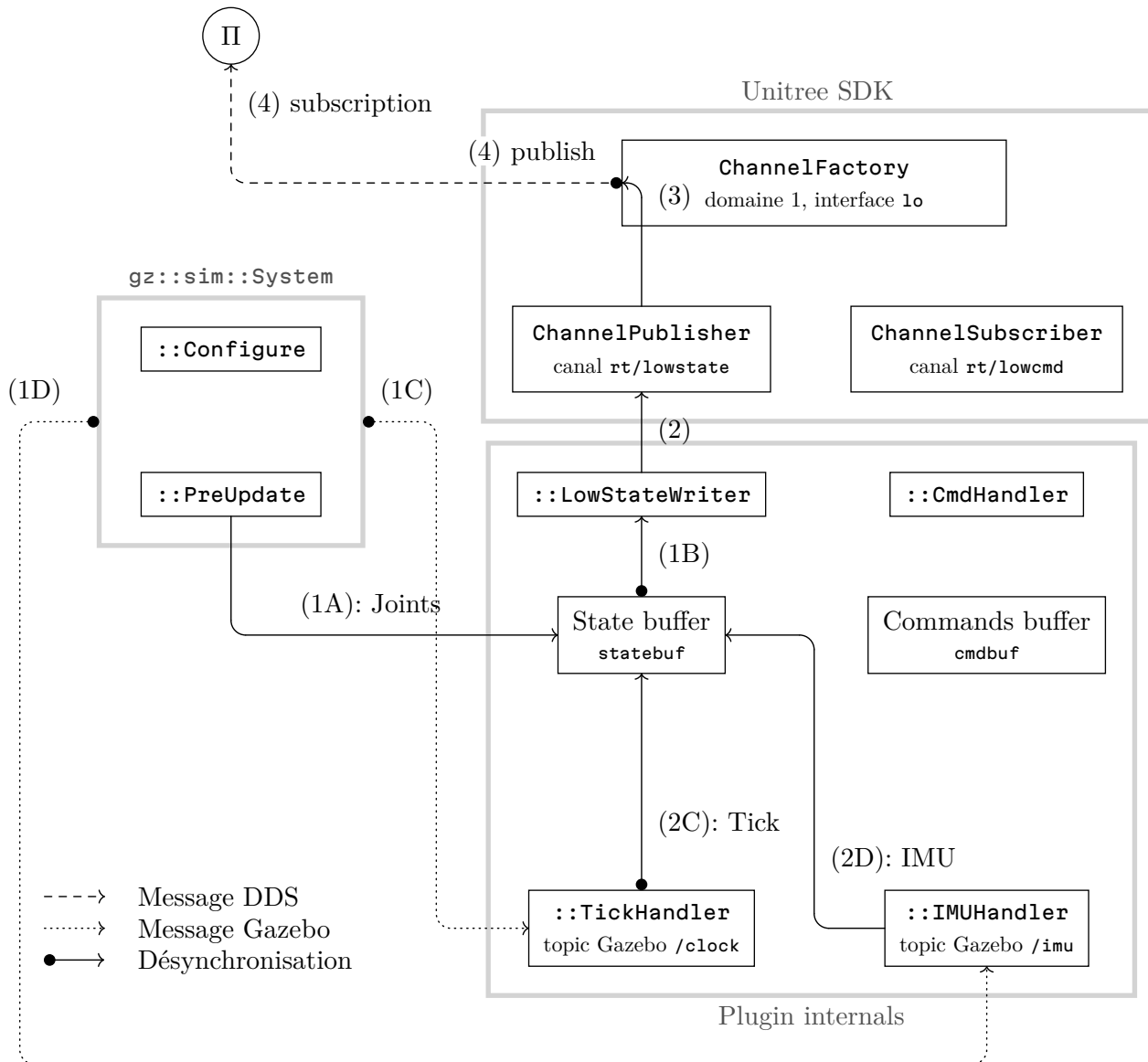


Fig. 17. – Phase d'envoi de l'état

Ici également, `LowStateWriter` s'exécute *parallèlement* à `::PreUpdate`: En effet, on démarre une boucle qui vient exécuter `LowStateWriter` périodiquement, dans un autre *thread*: on a donc aucune garantie de synchronisation entre les deux.

Ici, il y a en plus non pas deux, mais *cinq* boucles indépendantes qui sont en jeu:

- La boucle de simulation de Gazebo (fréquence d'appel de `::PreUpdate`),
- La boucle du `ChannelPublisher` (fréquence d'appel de `::LowStateWriter`), et
- La boucle de réception de Π (fréquence de réception de messages pour Π)
- La boucle de mise à jour du tick (fréquence d'envoi de ticks de simulation par Gazebo)
- La boucle de mise à jour de l'IMU (fréquence d'envoi des valeurs du capteur IMU par Gazebo)

Similairement à la réception de commandes, en comparant à la boucle de mise à jour de Π :

Si <code>::PreUpdate</code> est plus fréquente	On perdra des états intermédiaires, la résolution temporelle de l'évolution de l'état du robot disponible pour (ou acceptable par ²⁹) Π sera moins grande
Si <code>::PreUpdate</code> est moins fréquente	Π recevra plusieurs fois le même état, ce qui sera représentatif du fait que la simulation n'a pas encore avancé.

On a des effets similaires en comparant la fréquence de la boucle de mise à jour de l'IMU avec celle de la boucle de Π :

Si la boucle IMU est plus fréquente	Certaines valeurs du capteur ne seront pas prises en compte par la politique
Si la boucle IMU est moins fréquente	Π recevra plusieurs fois le même état, ce qui sera représentatif du fait que la simulation n'a pas encore avancé.

Pour la boucle du tick, cela a peu d'importance. En effet, Π ne dépend probablement pas du tick de simulation, ou si elle en dépend, on suppose que c'est une dépendance à une valeur peu précise (ce serait plutôt pour savoir « depuis quand est-ce qu'on a lancé la politique », ce qui ne demande pas une précision à la milliseconde). On met quand même à jour le tick pour que nos messages `rt/lowstate` synthétiques se rapprochent le plus possible des vrais messages, tels qu'envoyés par le robot physique.

5.7 Désynchronisations

Dans un même appel de `::PreUpdate`, on effectue d'abord la mise à jour du *State buffer*, puis on lit dans le *Commands buffer*.

Un cycle correspond donc à cinq boucles indépendantes, représentées ci-après:

Bleu	Simulation, qui doit englober l'entièreté d'un cycle
Rouge	<code>ChannelPublisher</code>
Rose	Politique Π
Vert	Mise à jour de l'IMU
Orange	Mise à jour du tick de simulation

²⁹En fonction de si `::LowStateWriter` est plus fréquente que Π (dans ce cas là, c'est ce qui est acceptable par Π qui est limitant) ou inversement (dans ce cas, c'est ce que la boucle du publisher met à disposition de Π qui est limitant)

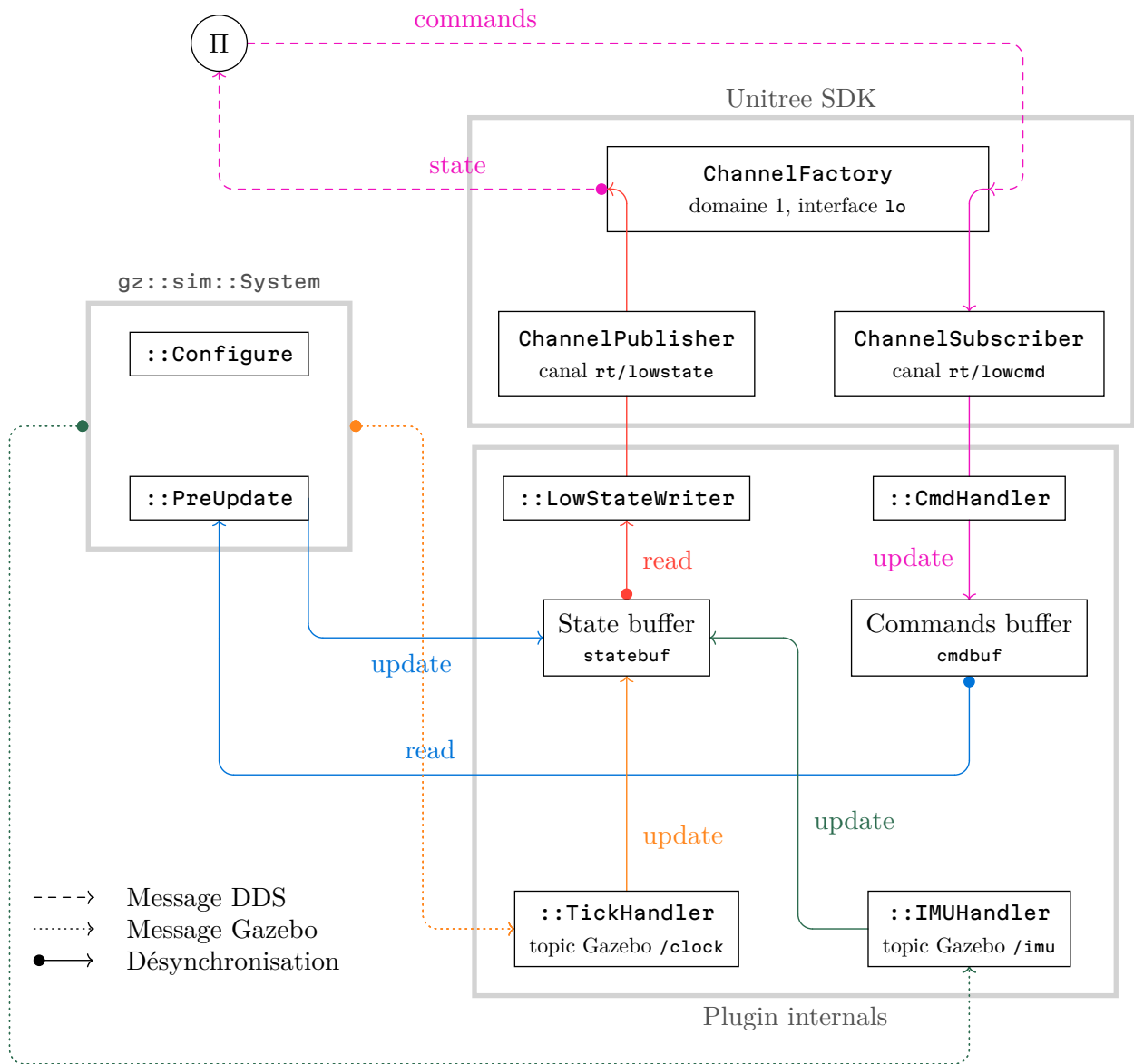


Fig. 18. – Cycle complet. Un cycle commence avec la flèche « update » partant de `::PreUpdate`

Ces désynchronisations pourraient expliquer les problèmes de performance rencontrés (cf Chapitre 5.9)

5.8 Vérification sur des politiques réelles

Après avoir testé le bridge sur les politiques d'exemples fournies par Unitree, il a été testé sur une politique en cours de développement au sein de l'équipe de robotique du LAAS, Gepetto.

L'analyse de la vidéo (cf Chapitre 5.10) montre que le bridge fonctionne: le comportement du robot est similaire à celui sur Isaac.

5.9 Amélioration des performances

Les premiers essais affichent un facteur temps-réel³⁰ autour des 10 à 15%.

En utilisant le *profiler* de Gazebo [79], on peut capturer des intervalles de temps et les annoter, pour identifier ce qui ralentit les cycles de simulation.

```
GZ_PROFILE_BEGIN("Label");
...
GZ_PROFILE_END();
```

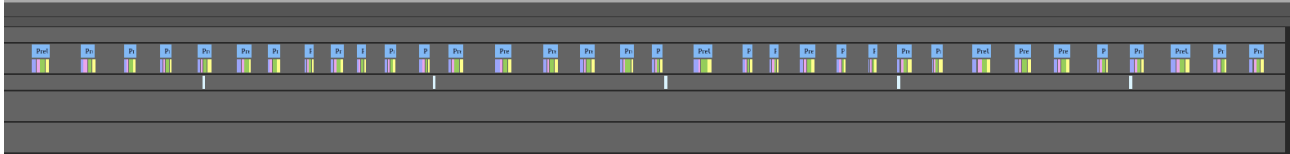


Fig. 19. – Profilage de *gz-unitree* lors d'une simulation

Chaque groupe de segment correspond à un cycle de simulation.

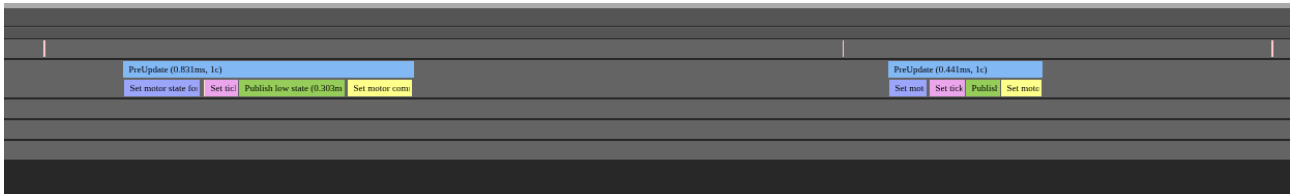
Prenons un cycle en particulier:

::PreUpdate 0.267 ms			
Update state 0.051 ms	Tick+CRC 0.039 ms	Publish state 0.142 ms	Update cmd. 0.028 ms

Tableau 3. – Profil d'un cycle de simulation

Plus de la moitié du temps de calcul du plugin est dû à de l'envoi de l'état du robot sur le canal DDS `rt/lowstate`.

Notons également que, même si ce cycle-là a duré 0.267 ms, la durée d'un cycle est assez variable, certains atteignent 0.8 ms.



Quelques mesures ont été tentées pour réduire le temps nécessaire à l'envoi d'un message DDS:

Restreindre DDS à localhost

Il est possible que DDS envoie les messages en mode « broadcast », c'est-à-dire à toute adresse IP accessible dans un certain intervalle. En restreignant à `localhost`, on s'assure que le message n'a pas à être copié plusieurs fois.

Déplacer dans un autre thread

C'est ce qui a motivé la désynchronisation du thread « LowStateWriter » (cf Chapitre 5.6.2)

Ajuster la fréquence d'envoi

Une fois `LowStateWriter` déplacé dans un thread indépendant, on peut ajuster la fréquence d'envoi, le thread étant récurrent³¹

³⁰Appelé RTF (Real-Time Factor) [78]. Un RTF de 100% signifie que la simulation s'exécute à vitesse réelle, un RTF inférieur à 1 signifie que la simulation est plus lente que ce qu'elle simule

³¹Créé avec `CreateRecurrentThreadEx`

Ainsi que d'autres optimisations, qui ne sont pas en rapport avec cette phase d'un cycle:

Mise en cache de joints à l'initialisation du plugin pour éviter de devoir appeler `model.JointByName` dans une *hot loop*³².

Utilisation d'une implémentation de CRC32 plus rapide tentative avec *CRC++* [80] non achevée, à cause d'un *stack smashing* pendant l'exécution

Après optimisations, on arrive à atteindre un RTF aux alentours des 30%. Des recherches supplémentaires sont nécessaires pour atteindre un RTF raisonnable.

5.10 Enregistrement automatique de vidéos

Gazebo possède une fonctionnalité d'enregistrement vidéo, ce qui s'avère utile pour partager des résultats de simulation.

Cependant, l'enregistrement vidéo n'est pas nativement contrôlable programmatiquement. L'idée était en effet de faire automatiquement tourner une simulation à chaque changement de la politique RL, et d'obtenir la vidéo du résultat, pour en observer l'évolution.

Il a donc fallu développer un autre plugin, héritant de `gz::gui::Plugin` cette fois-ci. Ce plugin écoute des messages sur des topics Gazebo, `/gui/record_video/{start,stop}`, et permet de démarrer et arrêter l'enregistrement, tout en indiquant le chemin vers le fichier MP4 de sortie.

Au final, un script complet permettant de démarrer une simulation et l'enregistrer en MP4 ressemble à ceci

```
# Fonction pour envoyer un message Gazebo avec un argument de type String et une
# valeur de retour de type Booléen
send_to_gz() {
    gz service -s $1 --reqtype gz.msgs.StringMsg --reptype gz.msgs.Boolean --req "data:
\"$2\"
}

# Lancement en arrière plan
gz sim robot.sdf & sim_pid=$!
# On attends que la simulation soit prête
sleep 30

# Lancement de l'enregistrement
send_to_gz /gui/record_video/start mp4

# Lancement de la politique RL
uv run policy.py & policy_pid=$!
# On décide de la durée maximale de la vidéo (si la politique ne l'arrête pas d'elle
# même)
sleep 120
kill $policy_pid

# Arrêt de l'enregistrement
send_to_gz /gui/record_video/stop file:///tmp/result.mp4

# Arrêt du simulateur
kill $sim_pid

# La vidéo est disponible à /tmp/result.mp4
```

³²Boucle (`for` ou `while`) dont le corps est exécuté un très grand nombre de fois, et dont la rapidité est importante

5.11 Mise en CI/CD

On appelle CI/CD (pour *Continuous Integration / Continuous Delivery*) la pratique consistant à intégrer fréquemment des petits changements à un dépôt de code source commun, en lançant des tests régulièrement (partie « CI ») et éventuellement en déployant la base de code fréquemment (partie « CD ») [81].

Une fois l'enregistrement vidéo rendu automatisable, si l'on veut mettre en place l'enregistrement vidéo automatique à chaque changement de la politique, il faut créer une description de *workflow* (dans notre cas, un workflow *Github Actions*).

Un workflow est un ensemble de commandes à exécuter dans un environnement virtualisé³³ ainsi que des événements et conditions décrivant quand lancer l'exécution (par exemple, « à chaque commit sur la branche `main` »). C'est un des outils permettant de mettre en place la CI/CD.

5.11.1 Une image de base avec Docker

L'environnement d'exécution des workflows ne comporte pas d'installation de Gazebo. Étant donné le temps de compilation élevé, on peut « factoriser » cette étape dans une *image de base*, de laquelle on démarre pour chaque exécution du workflow, dans laquelle tout les programmes nécessaires sont déjà installés.

Pour cela, on part d'une image Ubuntu, dans laquelle on installe le nécessaire: Just (pour lancer des commandes, un sorte de Makefile mais plus moderne [82]), FFMpeg (pour l'encodage H.264 servant à la création du fichier vidéo), XVFB (pour émuler un serveur X, cf Chapitre 5.11.2.2), Python (pour lancer la politique RL), Gazebo et gz-unitree.

```
FROM ubuntu:24.04

RUN apt update
# Just
RUN apt install -y curl just sudo
# Python (via le gestionnaire de versions et dépendances UV)
COPY --from=ghcr.io/astrol-sh/uv:latest /uv /uvx /bin/

# Code source de gz-unitree
COPY . .

# Gazebo et outils de compilation
RUN just setup

# FFMpeg, XVFB
RUN apt install -y git ffmpeg xvfb xterm

# Compilation et installation de de gz-unitree
RUN mkdir -p /usr/local/lib/gz-unitree/
RUN just install
```

Un autre workflow, défini cette fois-ci dans le dépôt de gz-unitree et non celui de la politique RL, crée une image Docker depuis ce Dockerfile, qui est ensuite utilisable via `ghcr.io/Gepetto/gz-unitree` [83]. Les commandes installant Gazebo et les outils de compilation sont écrites dans le *Justfile*, que nous lançons ici avec `RUN just setup`.

³³Qu'il s'agisse d'une machine virtuelle ou d'un simple container

5.11.2 Une pipeline Github Actions

Une fois cette image disponible, on peut l'utiliser dans un workflow Github sur le dépôt Git de la politique RL:

```
jobs:
  test:
    runs-on: ubuntu-latest
    container:
      image: ghcr.io/gepetto/gz-unitree:latest
    steps:
      - name: Checkout repository
        uses: actions/checkout@v5
```

Et lancer la simulation et l'enregistrement vidéo.

Pour récupérer le fichier vidéo final, on peut utiliser la notion d'*artifacts* de Github Actions:

```
- name: Save video as artifact
  uses: actions/upload-artifact@v4
  with:
    name: gz-unitree-video
    path: /tmp/result.mp4
```

Un environnement de développement contraignant

Développer et déboguer une définition de workflow peut s'avérer complexe et particulièrement chronophage: n'ayant pas d'accès interactif au serveur exécutant celui-ci, il faut envoyer ses changements au dépôt git, attendre que le workflow s'exécute entièrement, et regarder si quelque chose s'est mal passé.

Par exemple, si jamais des fichiers sont manquants, ou ne sont pas au chemin attendu, il faut modifier le workflow pour y rajouter des instruction listant le contenu d'un répertoire (en utilisant `ls` ou `tree`, par exemple), lancer le workflow à nouveau et regarder les logs.

Ceci rend le développement assez fastidieux, surtout quand le workflow s'exécute pendant des dizaines de minutes.

Émuler un serveur graphique

Les environnements de CI/CD s'apparentent plus à des serveurs qu'à des ordinateurs complets: en particulier, il n'y a pas d'interface graphique et donc pas de serveur d'affichage (*display server*).

Mais Gazebo a besoin d'un display server pour enregistrer une vidéo.

Il faut donc simuler un serveur d'affichage. Dans notre cas, l'environnement de CI/CD étant sous Linux, on simule un serveur X11 avec *XVFB* [84].

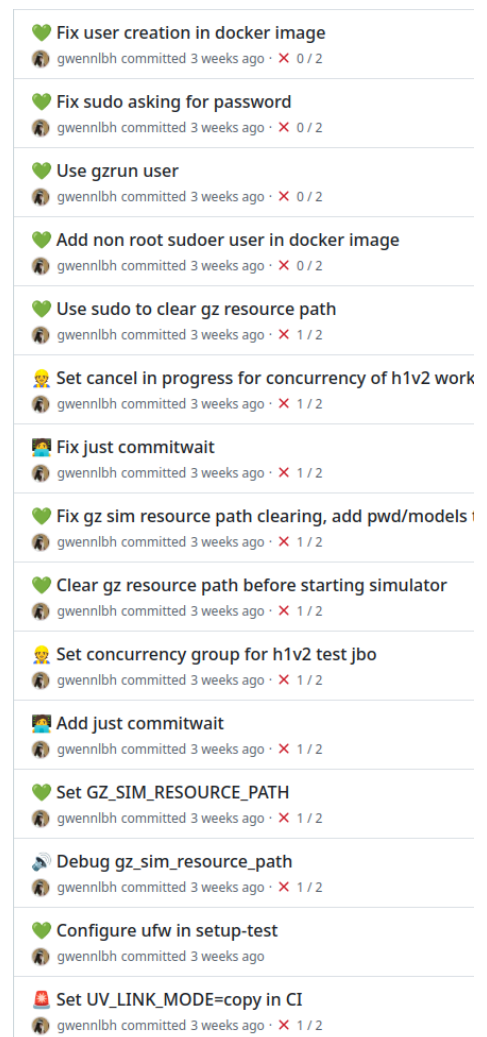


Fig. 20. – Quelques commits liés au développement du workflow³⁴

6 Conclusion

Il est désormais possible d'utiliser le simulateur modulaire, open-source et communautaire *Gazebo* pour entraîner des politiques de reinforcement learning sur le robot *H1v2* de la société Unitree.

Bien que la reproductibilité de compilation ne soit pas encore atteinte, l'utilisation du gestionnaire de paquets Nix semble possible pour obtenir des garanties de reproductibilité.

Les performances du *bridge* SDK2 \leftrightarrow Gazebo sont encore à améliorer, mais son utilisation est envisageable dans un contexte asynchrone, où le développement ne demande pas d'attendre les résultats de la simulation, via la pratique du CI/CD, par exemple.

Bibliographie

- [1] Shengbo Eben Li, *Reinforcement Learning for Sequential Decision and Optimal Control*. Springer Singapore, p. 1-460. doi: [10.1007/978-981-19-7784-8](https://doi.org/10.1007/978-981-19-7784-8).
- [2] Tambet Matiisen, « Demystifying deep reinforcement learning », 19 décembre 2015, *Computational Neuroscience Research Group at University of Tartu*. Consulté le: 13 octobre 2025. [En ligne]. Disponible sur: <https://web.archive.org/web/20180407053740/http://neuro.cs.ut.ee/demystifying-deep-reinforcement-learning/>
- [3] R. Sutton et A. Barto, *Reinforcement Learning: An Introduction*. in A Bradford book. MIT Press, 1998. Consulté le: 27 octobre 2025. [En ligne]. Disponible sur: <https://books.google.fr/books?id=CAFR6IBF4xYC>
- [4] T. G. Dietterich, « Hierarchical Reinforcement Learning with the MAXQ Value Function Decomposition », *CoRR*, 1999, Consulté le: 2002. [En ligne]. Disponible sur: <https://arxiv.org/abs/cs/9905014>
- [5] V. François-Lavet, R. Fonteneau, et D. Ernst, « How to Discount Deep Reinforcement Learning: Towards New Dynamic Strategies », *CoRR*, 2015, Consulté le: 13 octobre 2025. [En ligne]. Disponible sur: <http://arxiv.org/abs/1512.02011>
- [6] Nick Bostrom, « Ethical Issues in Advanced Artificial Intelligence », 2003, *Int. Institute of Advanced Studies in Systems Research and Cybernetics*. Consulté le: 8 octobre 2025. [En ligne]. Disponible sur: <https://nickbostrom.com/ethics/ai>
- [7] J. Schulman, S. Levine, P. Moritz, M. I. Jordan, et P. Abbeel, « Trust Region Policy Optimization », févr. 2015, Consulté le: 13 octobre 2025. [En ligne]. Disponible sur: <http://arxiv.org/abs/1502.05477v5>
- [8] « Trust Region Policy Optimization — Spinning Up documentation ». Consulté le: 14 octobre 2025. [En ligne]. Disponible sur: <https://spinningup.openai.com/en/latest/algorithms/trpo.html#background>
- [9] David Pollard, *Asymptotia*, Ch. 3, "Distances and affinities between measures". 2000, p. 6-7. Consulté le: 13 octobre 2025. [En ligne]. Disponible sur: <https://web.archive.org/web/20150412031925/http://www.stat.yale.edu/~pollard/Books/Asymptopia/Metrics.pdf>

³⁴Les émojis servent d'icônes pour différencier les types de commits, via le standard Gitmoji [85]

- [10] David J. C. MacKay, *Information Theory, Inference and Learning Algorithms*. Cambridge University Press, 2003, p. 34. Consulté le: 13 octobre 2025. [En ligne]. Disponible sur: https://books.google.fr/books?id=AKuMj4PN_EMC&lpg=PA34&pg=PA34#v=onepage&q&f=false
- [11] Z. Xie, « Simple Policy Optimization », janv. 2024, Consulté le: 16 octobre 2025. [En ligne]. Disponible sur: <http://arxiv.org/abs/2401.16025v2>
- [12] « Proximal Policy Optimization — Spinning Up documentation ». Consulté le: 16 octobre 2025. [En ligne]. Disponible sur: <https://spinningup.openai.com/en/latest/algorithms/ppo.html>
- [13] M. Stępień, R. Kourdis, C. Roux, et O. Stasse, « Latent Conditioned Loco-Manipulation Using Motion Priors », septembre 2025. Consulté le: 31 octobre 2025. [En ligne]. Disponible sur: <https://hal.science/hal-05242643>
- [14] NVIDIA Developer, « Isaac Sim - Robotics Simulation and Synthetic Data Generation ». Consulté le: 10 octobre 2025. [En ligne]. Disponible sur: <https://developer.nvidia.com/isaac/sim>
- [15] NVIDIA Developer, « PhysX SDK - Latest Features & Libraries ». Consulté le: 10 octobre 2025. [En ligne]. Disponible sur: <https://developer.nvidia.com/physx-sdk>
- [16] Consulté le: 16 juin 2025. [En ligne]. Disponible sur: <https://mujoco.org/>
- [17] « MuJoCo simulate tutorial ». Consulté le: 10 octobre 2025. [En ligne]. Disponible sur: <https://youtu.be/P83tKA1iz2Y>
- [18] Consulté le: 6 juin 2025. [En ligne]. Disponible sur: <https://gazebo-sim.org/>
- [19] « Gazebo Sim: Physics engines ». Consulté le: 10 octobre 2025. [En ligne]. Disponible sur: <https://gazebo-sim.org/api/sim/9/physics.html>
- [20] J. Lee *et al.*, « Dart: Dynamic animation and robotics toolkit », *The Journal of Open Source Software*, vol. 3, n° 22, p. 500, 2018.
- [21] Bullet Physics SDK, *bullet3*. (12 avril 2011). GitHub. Consulté le: 10 octobre 2025. [En ligne]. Disponible sur: <https://github.com/bulletphysics/bullet3>
- [22] « Bullet Real-Time Physics Simulation | Home of Bullet and PyBullet: physics simulation for games, visual effects, robotics and reinforcement learning. ». Consulté le: 10 octobre 2025. [En ligne]. Disponible sur: <https://pybullet.org/wordpress/>
- [23] Erwin Coumans, « Bullet Physics Simulation Constraint Solving and Featherstone Articulated Body Algorithm ». International Conference and Exhibition on Computer Graphics and Interactive Technologies, 2015. Consulté le: 10 octobre 2025. [En ligne]. Disponible sur: https://docs.google.com/presentation/d/1wGUJ4neOhw5i4pQRfSGtZPE3CIm7MfmqfTp5aJKuFYM/edit?slide=id.g644a5aa5f_0_16#slide=id.g644a5aa5f_0_16
- [24] Roy Featherstone, « Robot Dynamics Algorithms », 1978, *Springer New York, NY*.
- [25] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, et D. Meger, « Deep Reinforcement Learning that Matters », sept. 2017, Consulté le: 16 octobre 2025. [En ligne]. Disponible sur: <http://arxiv.org/abs/1709.06560v3>
- [26] Y. Gamage, D. Tiwari, M. Monperrus, et B. Baudry, « The Design Space of Lockfiles Across Package Managers », mai 2025, Consulté le: 17 octobre 2025. [En ligne]. Disponible sur: <http://arxiv.org/abs/2505.04834v2>

- [27] P. Fernique, « AutoWIG : automatisation de l'encapsulation de bibliothèques C++ en Python et en R », in *48èmes Journées de Statistique de la SFdS Montpellier*, Montpellier, France, mai 2016. Consulté le: 17 octobre 2025. [En ligne]. Disponible sur: <https://inria.hal.science/hal-01316276>
- [28] Brian Lonsdorf, « Professor Frisby's Mostly Adequate Guide to Functional Programming », 2015, *Github*. Consulté le: 4 septembre 2025. [En ligne]. Disponible sur: <https://github.com/MostlyAdequate/mostly-adequate-guide/blob/master/ch03.md>
- [29] « Reproducible Builds ». Consulté le: 4 septembre 2025. [En ligne]. Disponible sur: <https://reproducible-builds.org/>
- [30] Fortran 2015 Committee Draft (J3/17-007r2), *ISO/IEC JTC 1/SC 22/WG5/N2137*. International Organization for Standardisation, 2017, p. 336-338. Consulté le: 4 septembre 2025. [En ligne]. Disponible sur: <https://wg5-fortran.org/N2101-N2150/N2137.pdf>
- [31] « Relationship Between Routines, Functions, and Procedures », 13 janvier 2025, *IBM*. Consulté le: 4 septembre 2025. [En ligne]. Disponible sur: <https://www.ibm.com/docs/en/informix-servers/15.0.0?topic=statement-relationship-between-routines-functions-procedures>
- [32] « Different Types of Robot Programming Languages », 2015, *Plant Automation Technology*. Consulté le: 4 septembre 2025. [En ligne]. Disponible sur: <https://www.plantautomation-technology.com/articles/different-types-of-robot-programming-languages>
- [33] « Imperative programming: Overview of the oldest programming paradigm », 21 mai 2021, *IONOS*. Consulté le: 4 septembre 2025. [En ligne]. Disponible sur: <https://www.ionos.com/digitalguide/websites/web-development/imperative-programming/>
- [34] Bill Hoffman et Kenneth Martin, *The Architecture of Open Source Applications (Volume 1) CMake*. Consulté le: 4 septembre 2025. [En ligne]. Disponible sur: <https://aosabook.org/en/v1/cmake.html>
- [35] Consulté le: 19 mai 2025. [En ligne]. Disponible sur: <https://nix.dev/manual/nix/2.17/language/>
- [36] Consulté le: 3 septembre 2025. [En ligne]. Disponible sur: <https://search.nixos.org/packages>
- [37] NixOS Wiki Authors, « Nixpkgs/Contributing ». Consulté le: 3 septembre 2025. [En ligne]. Disponible sur: <https://wiki.nixos.org/wiki/Nixpkgs/Contributing>
- [38] « Cachix — Nix binary cache hosting ». Consulté le: 3 septembre 2025. [En ligne]. Disponible sur: <https://www.cachix.org/>
- [39] « Nix (package manager) — Sandboxing ». Consulté le: 3 septembre 2025. [En ligne]. Disponible sur: [https://wiki.nixos.org/wiki/Nix_\(package_manager\)#Internals](https://wiki.nixos.org/wiki/Nix_(package_manager)#Internals)
- [40] « GitHub-hosted runners », *Github*. Consulté le: 4 septembre 2025. [En ligne]. Disponible sur: <https://docs.github.com/en/actions/concepts/runners/github-hosted-runners>
- [41] Consulté le: 6 juin 2025. [En ligne]. Disponible sur: <https://www.raspberrypi.com/>
- [42] Fernando Borretti, « NixOS for the Impatient », 7 mai 2023. Consulté le: 4 septembre 2025. [En ligne]. Disponible sur: <https://borretti.me/article/nixos-for-the-impatient>
- [43] « gz-unitree/CMakeLists.txt at main · Gepetto/gz-unitree · GitHub ». Consulté le: 5 novembre 2025. [En ligne]. Disponible sur: <https://github.com/Gepetto/gz-unitree/blob/main/CMakeLists.txt>

- [44] *nix-ros-overlay*. (12 mars 2019). GitHub. Consulté le: 3 novembre 2025. [En ligne]. Disponible sur: <https://github.com/lopsided98/nix-ros-overlay>
- [45] « [WIP] gz-sim and dependencies: init by ShamrockLee · Pull Request #394757 · NixOS/nixpkgs · GitHub ». Consulté le: 3 novembre 2025. [En ligne]. Disponible sur: <https://github.com/NixOS/nixpkgs/pull/394757>
- [46] « ROS stuff: part 1 by Pandapip1 · Pull Request #355629 · NixOS/nixpkgs · GitHub ». Consulté le: 3 novembre 2025. [En ligne]. Disponible sur: <https://github.com/NixOS/nixpkgs/pull/355629>
- [47] Gelakais, « Pull requests on NixOS/nixpkgs ». Consulté le: 3 novembre 2025. [En ligne]. Disponible sur: <https://github.com/NixOS/nixpkgs/pulls?q=is%3Apr+gz+author%3AGelakais>
- [48] Gepetto, *gazebro2nix*. (31 août 2025). GitHub. Consulté le: 3 novembre 2025. [En ligne]. Disponible sur: <https://github.com/Gepetto/gazebro2nix>
- [49] « Unitree H1 / H1-2 ». Consulté le: 30 juin 2025. [En ligne]. Disponible sur: <https://www.unitree.com/h1/>
- [50] Consulté le: 6 juin 2025. [En ligne]. Disponible sur: https://github.com/unitreerobotics/unitree_sdk2
- [51] « Eclipse Cyclone DDS - Home ». Consulté le: 17 octobre 2025. [En ligne]. Disponible sur: <https://cyclonedds.io/>
- [52] Object Management Group, « DDS Interoperability Wire Protocol Specification Version 2.5 ». Consulté le: 24 juin 2025. [En ligne]. Disponible sur: <https://www.omg.org/spec/DDSI-RTPS/>
- [53] Object Management Group, « Interface Definition Language Specification Version 4.2 ». Consulté le: 18 juin 2025. [En ligne]. Disponible sur: <https://www.omg.org/spec/IDL/4.2/About-IDL>
- [54] Unitree, *unitree_sdk2/include/unitree/idl/hg/LowCmd_.hpp at main@2025-10-17, lines 33 to 63*. Github. Consulté le: 17 octobre 2025. [En ligne]. Disponible sur: https://github.com/unitreerobotics/unitree_sdk2/blob/main/@%7B2025-10-17%7D/include/unitree/idl/hg/LowCmd_.hpp#L33-L63
- [55] W.-Y. Liang, Y. Yuan, et H.-J. Lin, « A Performance Study on the Throughput and Latency of Zenoh, MQTT, Kafka, and DDS », mars 2023, Consulté le: 17 octobre 2025. [En ligne]. Disponible sur: <http://arxiv.org/abs/2303.09419v1>
- [56] Unitree, *unitreerobotics/unitree_sdk2, main@2025-10-17*. Github. Consulté le: 17 octobre 2025. [En ligne]. Disponible sur: https://github.com/unitreerobotics/unitree_sdk2/tree/main/@%7B2025-10-17%7D
- [57] « IMPORTED_LOCATION — CMake 4.2.0-rc1 Documentation ». Consulté le: 17 octobre 2025. [En ligne]. Disponible sur: https://cmake.org/cmake/help/latest/prop_tgt/IMPORTED_LOCATION.html
- [58] « Ghidra - Powerful Open-Source Reverse Engineering Tool ». Consulté le: 17 octobre 2025. [En ligne]. Disponible sur: <https://ghidralite.com/>
- [59] Unitree Robotics, *unitree_mujoco*. (1 novembre 2021). GitHub. Consulté le: 17 octobre 2025. [En ligne]. Disponible sur: https://github.com/unitreerobotics/unitree_mujoco

- [60] « High-Level Motion Services Interface - 高レベルモーションサービスインターフェース ». Unitree. Consulté le: 30 octobre 2025. [En ligne]. Disponible sur: https://support.unitree.com/home/en/H1_developer/Sports_Services_Interface
- [61] « HelloWorld source code — Eclipse Cyclone DDS, 0.11.0 ». Consulté le: 22 octobre 2025. [En ligne]. Disponible sur: https://cyclonedds.io/docs/cyclonedds/latest/getting_started/helloworld/helloworld_source_code.html
- [62] Unitree Robotics, *unitree_ros2*. (20 octobre 2023). GitHub. Consulté le: 22 octobre 2025. [En ligne]. Disponible sur: https://github.com/unitreerobotics/unitree_ros2
- [63] Jamie Strandboge, *ufw*. Launchpad. Consulté le: 5 novembre 2025. [En ligne]. Disponible sur: <https://launchpad.net/ufw>
- [64] « Protocols/rtps - Wireshark Wiki ». Consulté le: 22 octobre 2025. [En ligne]. Disponible sur: <https://wiki.wireshark.org/Protocols/rtps>
- [65] « SDFFormat Specification, world.plugin element ». Consulté le: 22 octobre 2025. [En ligne]. Disponible sur: http://sdformat.org/spec?ver=1.12&elem=world#world_plugin
- [66] « Gazebo Sim: Finding resources ». Consulté le: 22 octobre 2025. [En ligne]. Disponible sur: https://gazebo.org/api/sim/8/resources.html#:~:text=All%20paths%20on%20the%20GZ_SIM_SYSTEM_PLUGIN_PATH%20environment%20variable
- [67] « SDFFormat Specification, world.plugin.filename element ». Consulté le: 22 octobre 2025. [En ligne]. Disponible sur: http://sdformat.org/spec?ver=1.12&elem=world#plugin_filename
- [68] « Gazebo Transport: Nodes and Topics ». Gazebo. Consulté le: 29 octobre 2025. [En ligne]. Disponible sur: <https://gazebo.org/api/transport/8/nodetopics.html>
- [69] « Protocol Buffers Documentation ». Consulté le: 29 octobre 2025. [En ligne]. Disponible sur: <https://protobuf.dev/>
- [70] « Gazebo Transport: Messages ». Gazebo. Consulté le: 29 octobre 2025. [En ligne]. Disponible sur: <https://gazebo.org/api/transport/8/messages.html>
- [71] cplusplus.com, « std::bind ». Consulté le: 23 octobre 2025. [En ligne]. Disponible sur: <https://cplusplus.com/reference/functional/bind/>
- [72] L. T. H. Gam, D. H. Quan, P. V. B. Ngoc, B. H. Quan, et B. T. Thanh, « Position–Force Control of a Lower-Limb Rehabilitation Robot Using a Force Feed-Forward and Compensative Gravity Proportional Derivative Method », *Electronics*, vol. 13, n° 22, 2024, doi: [10.3390/electronics13224494](https://doi.org/10.3390/electronics13224494).
- [73] « rt/lowCmd topic Description - 低レベルコマンドトピックの説明 ». Unitree. Consulté le: 23 octobre 2025. [En ligne]. Disponible sur: https://support.unitree.com/home/en/H1_developer/H1-2_Basic_Services_Interface#heading-1
- [74] M. Araki, « PID Control », *Control Systems, Robotics, and Automation*, vol. 2. 2009. Consulté le: 23 octobre 2025. [En ligne]. Disponible sur: <http://www.eolss.net/ebooks/Sample%20Chapters/C18/E6-43-03-03.pdf>
- [75] « rt/lowState topic Description - 低レベル状態トピックの説明 ». Unitree. Consulté le: 23 octobre 2025. [En ligne]. Disponible sur: https://support.unitree.com/home/en/H1_developer/H1-2_Basic_Services_Interface#heading-4

- [76] Unitree Robotics, *unitree_sdk2/example/h1/low_level/motors.hpp* at *211dad48a40fb79485ad3d34c0ec6f6b12b0d8a5*. Github. Consulté le: 29 octobre 2025. [En ligne]. Disponible sur: https://github.com/unitreerobotics/unitree_sdk2/blob/211dad48a40fb79485ad3d34c0ec6f6b12b0d8a5/example/h1/low_level/motors.hpp#L54-L75
- [77] Unitree Robotics, *unitree_mujoco/example/cpp/stand_go2.cpp* at *24b7a498f6bc032684a8c8af30876a84d53d03e4*. Github. Consulté le: 29 octobre 2025. [En ligne]. Disponible sur: https://github.com/unitreerobotics/unitree_mujoco/blob/24b7a498f6bc032684a8c8af30876a84d53d03e4/example/cpp/stand_go2.cpp#L54-L84
- [78] Centre for Tactile Internet with Human-in-the-Loop, « Simulation Speed in Gazebo ». Consulté le: 28 octobre 2025. [En ligne]. Disponible sur: <https://ceti.pages.st.inf.tu-dresden.de/robotics/howtos/SimulationSpeed.html>
- [79] « Gazebo Common: Profiler ». Consulté le: 28 octobre 2025. [En ligne]. Disponible sur: <https://gazebosim.org/api/common/7/profiler.html>
- [80] Daniel Bahr, *CRCpp*. (1 mai 2016). GitHub. Consulté le: 29 octobre 2025. [En ligne]. Disponible sur: <https://github.com/d-bahr/CRCpp>
- [81] I. Sacolick, « What is CI/CD? Continuous integration and continuous delivery explained », 1 juin 2021, *InfoWorld*. Consulté le: 29 octobre 2025. [En ligne]. Disponible sur: <https://www.infoworld.com/article/2269266/what-is-cicd-continuous-integration-and-continuous-delivery-explained.html>
- [82] « Introduction - Just Programmer's Manual ». Consulté le: 28 octobre 2025. [En ligne]. Disponible sur: <https://just.systems/man/en/>
- [83] Gepetto team at LAAS-CNRS, *Package gz-unitree*. Github. Consulté le: 28 octobre 2025. [En ligne]. Disponible sur: <https://github.com/Gepetto/gz-unitree/pkgs/container/gz-unitree>
- [84] I. David P. Wiggins The Open Group, « Xvfb », *X.org*. Consulté le: 4 septembre 2025. [En ligne]. Disponible sur: <https://www.x.org/archive/X11R7.7/doc/man/man1/Xvfb.1.xhtml>
- [85] Carlos Cuesta, « gitmoji | An emoji guide for your commit messages ». Consulté le: 29 octobre 2025. [En ligne]. Disponible sur: <https://gitmoji.dev/>

A Preuves

A.1 Cas dégénéré de $D_{\text{KL}}(Q, Q') = 0$ sans utilisation de max

Soit S (resp. $A \subset \mathbb{N}$) l'espace des états (resp. actions) de l'environnement. Soit $Q : S \times A \rightarrow [0, 1]$ une distribution de probabilité du choix par l'agent d'une action dans un état tel que

$$\forall s \in S, Q(s, 1) = Q(s, 2) \quad (27)$$

Soit $Q' : S \times A \rightarrow [0, 1]$ définit ainsi:

$$\forall s \in S, Q'(s, 1) := 2Q(s, 1) \quad (28)$$

$$\forall s \in S, Q'(s, 2) := \frac{1}{2}Q(s, 2) \quad (29)$$

$$\forall s \in S, \forall a \in A - \{1, 2\}, Q'(s, a) := Q(s, a) \quad (30)$$

On a

$$\begin{aligned} D_{\text{KL}}(Q \parallel Q') &= \sum_{(s,a) \in S \times A} Q(s, a) \log \frac{Q(s, a)}{Q'(s, a)} \\ &\text{On découpe la somme selon les valeurs de } A: \\ &= \sum_{s \in S} \sum_{a \in A - \{1, 2\}} \left[Q(s, a) \log \frac{Q(s, a)}{Q'(s, a)} \right] + Q(s, 1) \log \frac{Q(s, 1)}{Q'(s, 1)} + Q(s, 2) \log \frac{Q(s, 2)}{Q'(s, 2)} \\ &= \sum_{s \in S} \underbrace{\sum_{a \in A - \{1, 2\}} Q(s, a) \log \frac{Q(s, a)}{Q(s, a)}}_{\text{d'après (30)}} + Q(s, 1) \log \underbrace{\frac{Q(s, 1)}{2Q(s, 1)}}_{\text{d'après (28)}} + Q(s, 2) \log \underbrace{\frac{Q(s, 2)}{\frac{1}{2}Q(s, 2)}}_{\text{d'après (29)}} \\ &= \sum_{s \in S} Q(s, 1) \left[\log Q(s, 1) - \log Q(s, 1) - \log 2 \right] + \quad (31) \\ &\quad Q(s, 2) \left[\log Q(s, 2) - \log Q(s, 2) - \log \frac{1}{2} \right] \\ &= \sum_{s \in S} -Q(s, 1) \log 2 + Q(s, 2) \log 2 \\ &= \sum_{s \in S} \log 2 \underbrace{(Q(s, 2) - Q(s, 1))}_{\text{d'après (27)}} \\ &= \sum_{s \in S} 0 = 0 \end{aligned}$$

A.2 $\eta(\pi, r)$ comme une espérance

Soit r une fonction récompense et π une politique. Soit C une variable aléatoire à valeurs dans \mathcal{S} , dont la loi de probabilité suit celle de π .

On a

$$\begin{aligned} \exp\left(\sum_{t=0}^{\infty} \gamma^t r(C_t)\right) &= \sum_{(c_t)_{t \in \mathbb{N}} \in \mathcal{S}} \left(\sum_{t=0}^{\infty} \gamma^t r(c_t)\right) \mathbb{P}\left(\sum_{t=0}^{\infty} \gamma^t r(C_t) = \sum_{t=0}^{\infty} \gamma^t r(c_t)\right) \\ &= \sum_{(c_t)_{t \in \mathbb{N}} \in \mathcal{S}} \left(\sum_{t=0}^{\infty} \gamma^t r(c_t)\right) \mathbb{P}(C = (c_t)_{t \in \mathbb{N}}) \end{aligned} \quad (32)$$

Soit S (resp. A) la suite des premiers (resp. deuxièmes) éléments de C , c'est-à-dire $\forall t \in \mathbb{N}, (S_t, A_t) := C_t$.

Étant donné la définition de \mathcal{S} :

- S_t dépend de A_{t-1} et S_{t-1}
- A_t dépend de S_t

On a alors, pour toute suite $(c_t)_{t \in \mathbb{N}} \in \mathcal{S}$:

$$\begin{aligned} P(C = (c_t)_{t \in \mathbb{N}}) &= \\ \mathbb{P}(S_0 = s_0) \mathbb{P}(A_0 = a_0 \mid S_0 = s_0) \prod_{t=1}^{\infty} \mathbb{P}(S_t = s_t \mid C_{t-1} = c_{t-1}) \mathbb{P}(A_t = a_t \mid S_t = s_t) \end{aligned} \quad (33)$$

On a

$$\begin{aligned} \mathbb{P}(S_0 = s_0) &= \rho_0(s_0) \\ \forall t \in \mathbb{N}, \quad \mathbb{P}(A_t = a_t \mid S_t = s_t) &= Q_\pi(s_t, a_t) \\ \forall t \in \mathbb{N}^*, \quad \mathbb{P}(S_t = s_t \mid C_{t-1} = c_{t-1}) &= \mathbb{P}(M(C_{t-1}) = M(c_{t-1}) \mid C_{t-1} = c_{t-1}) \\ &= \mathbb{P}(C_{t-1} = c_{t-1} \mid C_{t-1} = c_{t-1}) = 1 \end{aligned} \quad (34)$$

Donc on a

$$\begin{aligned} P(C = (c_t)_{t \in \mathbb{N}}) &= \rho_0(s_0) Q_\pi(s_0, a_0) \prod_{t=1}^{\infty} Q_\pi(s_t, a_t) \\ &= \rho_0(s_0) \prod_{t=0}^{\infty} Q_\pi(s_t, a_t) \end{aligned} \quad (35)$$

Et ainsi

$$\begin{aligned} \exp\left(\sum_{t=0}^{\infty} \gamma^t r(C_t)\right) &= \sum_{(c_t)_{t \in \mathbb{N}} \in \mathcal{S}} \left(\sum_{t=0}^{\infty} \gamma^t r(c_t)\right) \mathbb{P}(C = (c_t)_{t \in \mathbb{N}}) \\ &= \sum_{(c_t)_{t \in \mathbb{N}} \in \mathcal{S}} \left(\sum_{t=0}^{\infty} \gamma^t r(c_t)\right) \rho_0(s_0) \prod_{t=0}^{\infty} Q_\pi(s_t, a_t) \\ &= \eta(\pi, r) \quad \blacksquare \end{aligned} \quad (36)$$

A.3 Simplification de l'expression de $L(s, a, \Pi, \Pi', R)$ dans PPO-Clip

Soit $(s, a) \in S \times A$, et Π' une politique. Posons $\alpha := A_{\Pi', R}(s, a)$, $q/q' := Q_{\Pi}(s, a)/Q_{\Pi'}(s, a)$.

Cas $\alpha > 0$

$$\begin{aligned}
 & L(s, a, \Pi, \Pi', R) \\
 &= \min\left(\frac{q}{q'}\alpha, \quad \text{clip}\left(\frac{q}{q'}, 1 - \varepsilon, 1 + \varepsilon\right)\alpha\right) \\
 &= \min\left(\frac{q}{q'}, \quad \text{clip}\left(\frac{q}{q'}, 1 - \varepsilon, 1 + \varepsilon\right)\right)\alpha \\
 & \quad \text{car } \alpha > 0
 \end{aligned}$$

Cas $\alpha < 0$

$$\begin{aligned}
 & L(s, a, \Pi, \Pi', R) \\
 &= \min\left(\frac{q}{q'}\alpha, \quad \text{clip}\left(\frac{q}{q'}, 1 - \varepsilon, 1 + \varepsilon\right)\alpha\right) \\
 &= \max\left(\frac{q}{q'}, \quad \text{clip}\left(\frac{q}{q'}, 1 - \varepsilon, 1 + \varepsilon\right)\right)\alpha \\
 & \quad \text{car } \alpha < 0
 \end{aligned}$$

...et $q/q' \in [1 - \varepsilon, 1 + \varepsilon]$

$$\begin{aligned}
 &= \min\left(\frac{q}{q'}, \quad \text{clip}\left(\frac{q}{q'}, 1 - \varepsilon, 1 + \varepsilon\right)\right)\alpha &= \max\left(\frac{q}{q'}, \quad \text{clip}\left(\frac{q}{q'}, 1 - \varepsilon, 1 + \varepsilon\right)\right)\alpha \\
 &= \min\left(\frac{q}{q'}, \quad \frac{q}{q'}\right)\alpha &= \max\left(\frac{q}{q'}, \quad \frac{q}{q'}\right)\alpha \\
 &= \min\left(\frac{q}{q'}, 1 + \varepsilon\right)\alpha &= \max\left(\frac{q}{q'}, 1 - \varepsilon\right)\alpha \\
 & \quad \text{car } 1 + \varepsilon > \frac{q}{q'} & \quad \text{car } 1 - \varepsilon < \frac{q}{q'}
 \end{aligned}$$

...et $q/q' > 1 + \varepsilon$

$$\begin{aligned}
 &= \min\left(\frac{q}{q'}, \quad \text{clip}\left(\frac{q}{q'}, 1 - \varepsilon, 1 + \varepsilon\right)\right)\alpha &= \max\left(\frac{q}{q'}, \quad \text{clip}\left(\frac{q}{q'}, 1 - \varepsilon, 1 + \varepsilon\right)\right)\alpha \\
 &= \min\left(\frac{q}{q'}, \quad 1 + \varepsilon\right)\alpha &= \max\left(\frac{q}{q'}, \quad 1 + \varepsilon\right)\alpha \\
 & &= \max\left(\frac{q}{q'}, \quad 1 - \varepsilon\right)\alpha \\
 & & \quad \text{car } 1 - \varepsilon < 1 + \varepsilon < \frac{q}{q'}
 \end{aligned}$$

...et $q/q' < 1 - \varepsilon$

$$\begin{aligned}
 &= \min\left(\frac{q}{q'}, \quad \text{clip}\left(\frac{q}{q'}, 1 - \varepsilon, 1 + \varepsilon\right)\right)\alpha &= \max\left(\frac{q}{q'}, \quad \text{clip}\left(\frac{q}{q'}, 1 - \varepsilon, 1 + \varepsilon\right)\right)\alpha \\
 &= \min\left(\frac{q}{q'}, \quad 1 - \varepsilon\right)\alpha &= \max\left(\frac{q}{q'}, \quad 1 - \varepsilon\right)\alpha \\
 &= \min\left(\frac{q}{q'}, \quad 1 + \varepsilon\right)\alpha \\
 & \quad \text{car } 1 + \varepsilon > 1 - \varepsilon > \frac{q}{q'}
 \end{aligned}$$

B Implémentation de CRC32

Extraite du code source de *unitree_mujoco*, sous license BSD-3-Clause [77]

```
uint32_t crc32_core(uint32_t *ptr, uint32_t len)
{
    unsigned int xbit = 0;
    unsigned int data = 0;
    unsigned int CRC32 = 0xFFFFFFFF;
    const unsigned int dwPolynomial = 0x04c11db7;

    for (unsigned int i = 0; i < len; i++)
    {
        xbit = 1 << 31;
        data = ptr[i];
        for (unsigned int bits = 0; bits < 32; bits++)
        {
            if (CRC32 & 0x80000000)
            {
                CRC32 <<= 1;
                CRC32 ^= dwPolynomial;
            }
            else
            {
                CRC32 <<= 1;
            }

            if (data & xbit)
                CRC32 ^= dwPolynomial;
            xbit >>= 1;
        }
    }

    return CRC32;
}
```