

# المدرسة الوطنية المتعددة التقنيات Ecole Nationale Polytechnique d'Alger



## Département d'Automatique

Projet de Fin d'Etudes en vue de l'obtention des diplômes  
d'Ingénieur d'état et de Master en Automatique



## Thème

Optimisation par Reinforcement Learning et  
implémentation d'une technique V-SLAM(2D)  
sur un robot mobile

### Présenté par :

Oussama DEROUICHE  
El Hacene CHABANE

### Dirigé par :

Pr. Mohamed TADJINE  
Mr. Zeryab MOUSSAOUI

24 Juin 2018

# Plan de travail

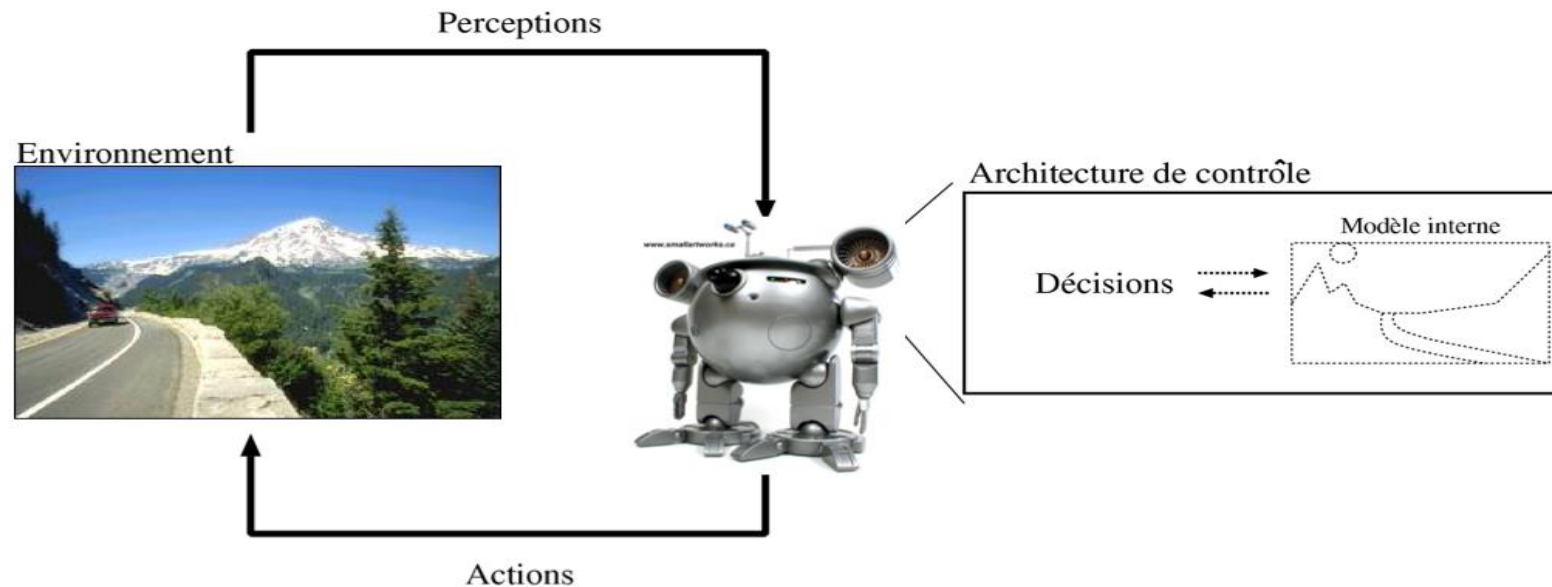
- 1 Introduction
- 2 Techniques utilisées
- 3 Navigation autonome
- 4 Mise en œuvre pratique
- 5 Conclusion et perspectives

## Présentation du sujet



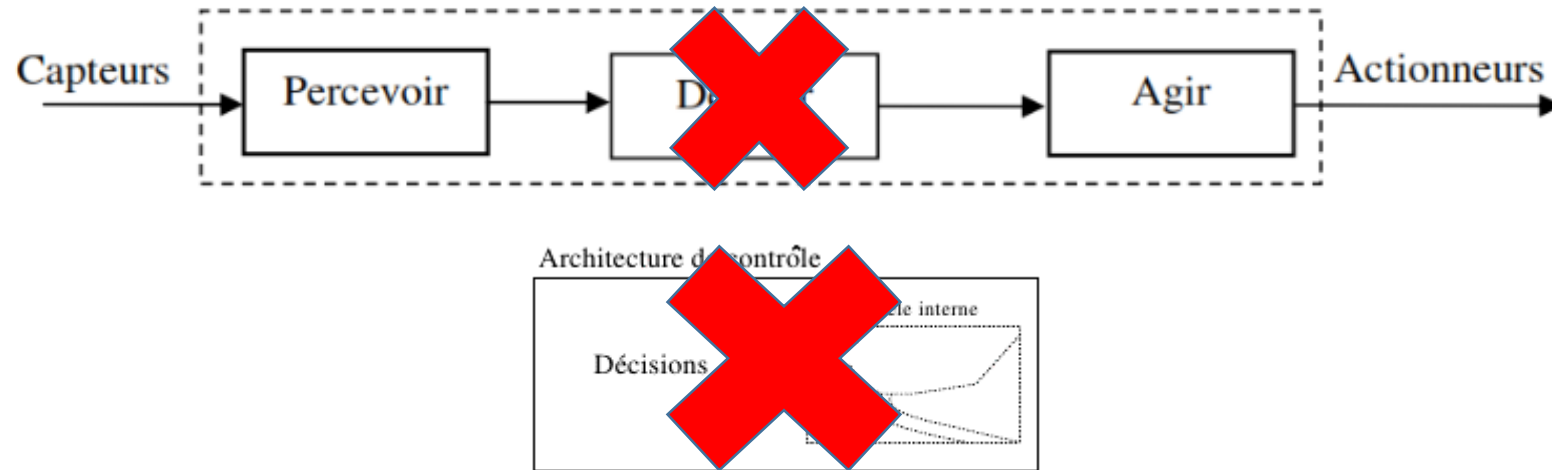
## Robotique mobile

- On distingue robotique de manipulation et robotique mobile.
- Un robot mobile est une machine intelligente capable de se mouvoir dans un environnement.



## Robotique mobile

- Architecture hiérarchique



- Architecture réactive

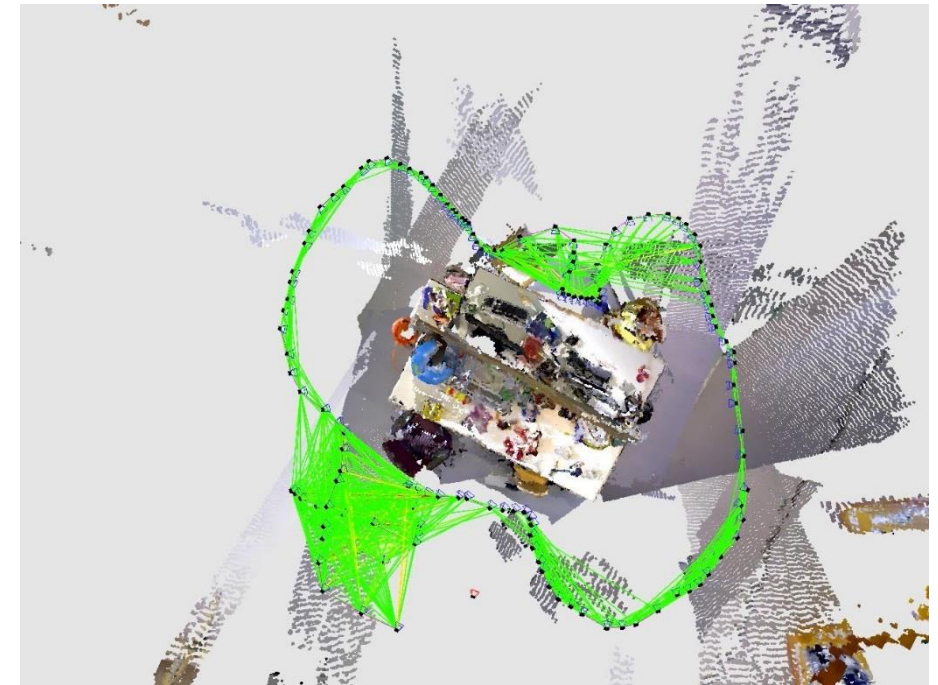


## Problématiques

- La Navigation, problématique principale de la robotique mobile.
- La Localisation, problématique secondaire mais essentielle.
- Solution retenue : Navigation Floue
- Optimisation par Reinforcement Learning.
- Localisation : incrémentale ou par fermeture de la boucle.
- Implémentation du VSLAM pour gain de précision.

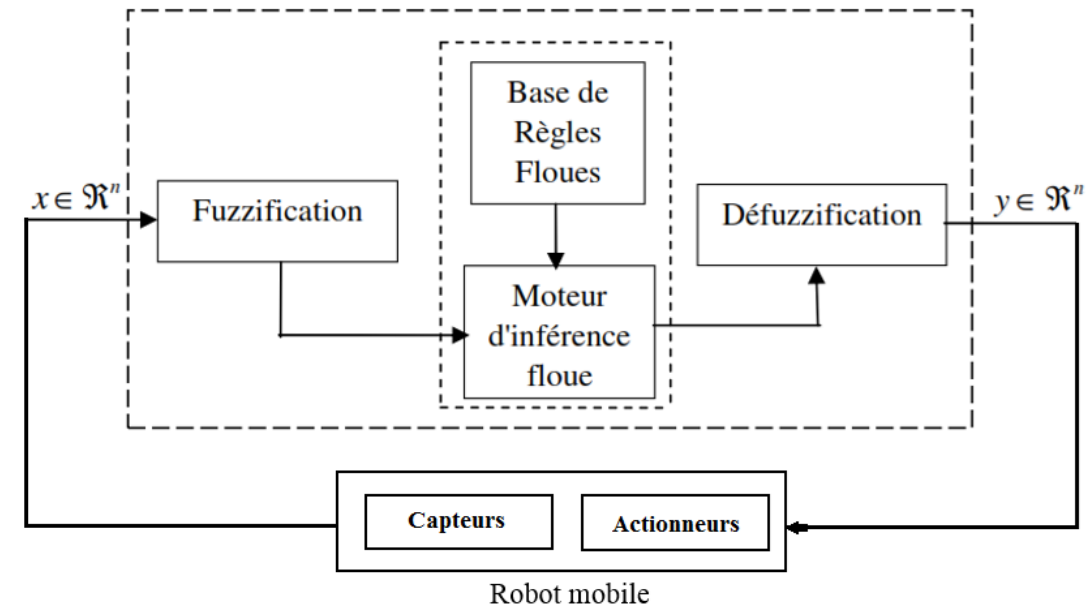
## Visual SLAM

- *Simultaneous Localisation And Mapping*, avec (au moins) une caméra.
- Plus fiable et plus robuste que l'Odométrie.
- Monoculaire : ORB-SLAM2 (2017), LSDSLAM (2015) MonoSLAM (2007), PTAM(2007)
- RGBD : ORB-SLAM2(2017), RGB-D SLAM (2014)
- Solution retenue : ORB-SLAM2
- ORB SLAM2 est récent, polyvalent et open-source.



## Logique floue

- Extension de la logique, introduite par le Professeur iranien Lotfi Zadeh en 1965.
- Application : Automatique, Traitement Images, ...
- Choix : Régulateur Mamdani



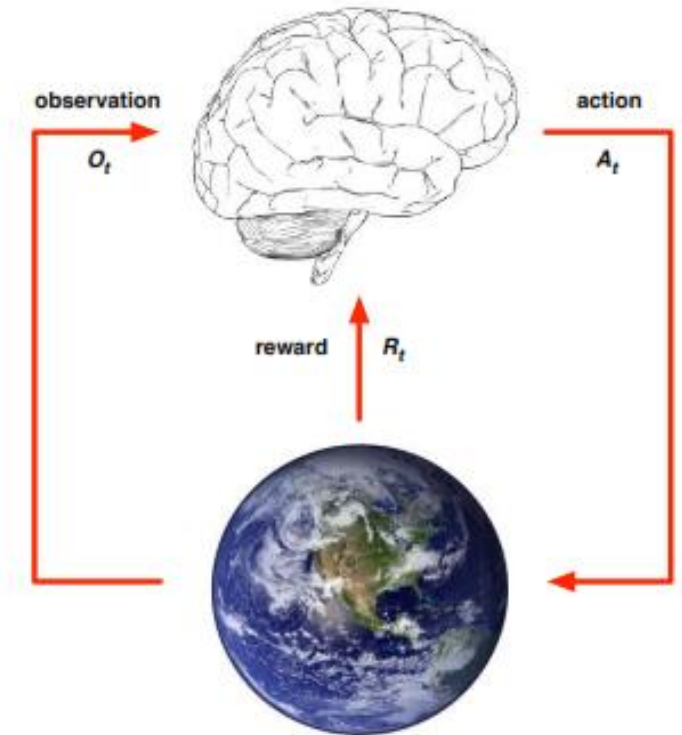


## Apprentissage par renforcement

- Branche du *Machine Learning*.
- Interactions d'un agent avec l'environnement.
- Formalisé par *Markov Decision Process*  $\langle S, A, P, R, \gamma \rangle$
- Recherche d'une politique optimale maximisant les récompenses.
- Fonctions de décision  $Q(s, a)$  ou  $V(s)$  à estimer :

$$V_{\pi}(s) = E_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

$$Q_{\pi}(s, a) = E_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$



## Apprentissage par renforcement

Méthodes de la programmation dynamique :

### Policy Iteration

$$V_{k+1}(s) = \sum_{a \in \mathbb{A}(s)} \pi(a | s) \left( \mathbb{R}_s^a + \gamma \sum_{s' \in \mathbb{S}} \mathbb{P}_{s,s'}^a V_k(s') \right)$$

### Value Iteration

$$V_{k+1}(s) = \max_{a \in \mathbb{A}(s)} Q_{\pi}(s, a) = \max_{a \in \mathbb{A}(s)} \left( \mathbb{R}_s^a + \gamma \sum_{s' \in \mathbb{S}} \mathbb{P}_{s,s'}^a V_k(s') \right)$$

Méthode de Monte Carlo :

$$V_{k+1}(s) = V_k(s) + \frac{1}{k+1} [G_{k+1}(s) - V_k(s)] \quad \text{ou} \quad Q_{k+1}(s, a) = Q_k(s, a) + \alpha [G_{k+1}(s) - Q_k(s, a)]$$

Méthodes de *Temporal difference learning* :

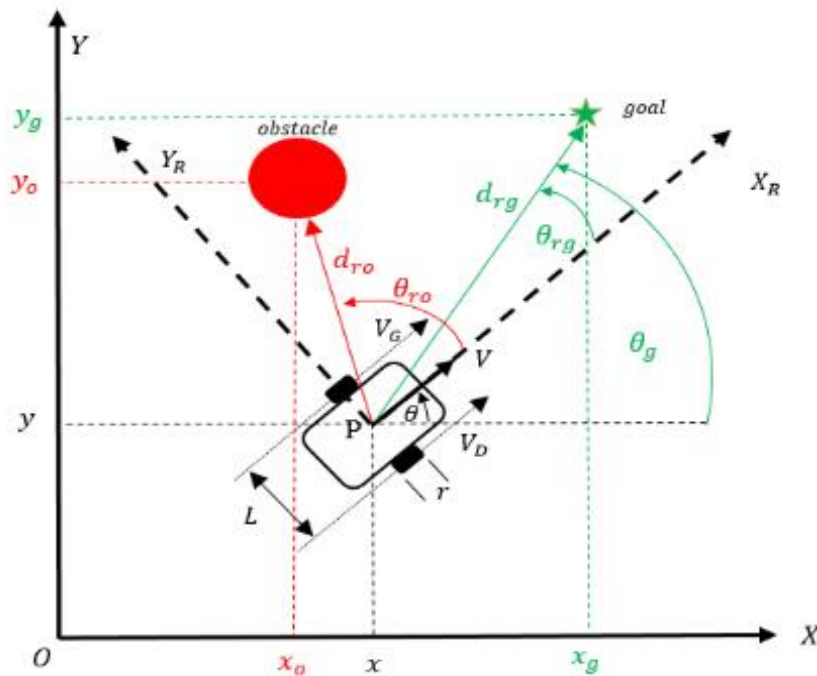
### Q Learning

$$Q(s_t, a_t) \leftarrow (1 - \alpha) Q(s_t, a_t) + \alpha \left[ R_t + \gamma \max_{a \in \mathbb{A}(s_{t+1})} Q(s_{t+1}, a) \right]$$

### SARSA

$$Q(s_t, a_t) \leftarrow (1 - \alpha) Q(s_t, a_t) + \alpha [R_t + \gamma Q(s_{t+1}, a_{t+1})]$$

## Modélisation



Configuration

Odométrie :

$$x(k+1) = x(k) + r \frac{T}{2} (\omega_D(k) + \omega_G(k)) \cos(\theta(k))$$

$$y(k+1) = y(k) + r \frac{T}{2} (\omega_D(k) + \omega_G(k)) \sin(\theta(k))$$

$$\theta(k+1) = \theta(k) + r \frac{T}{2} (\omega_D(k) - \omega_G(k))$$

But :

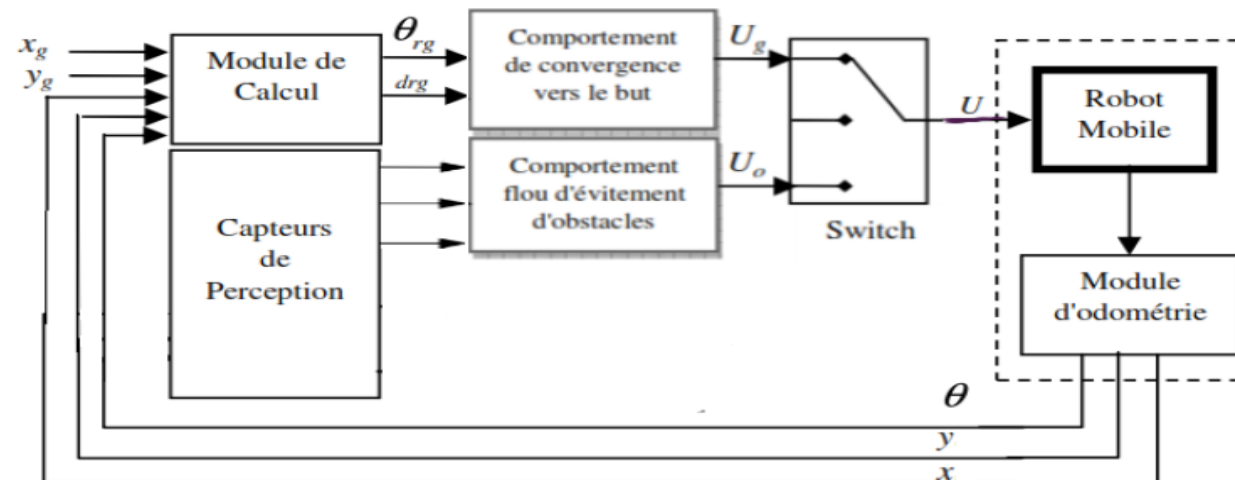
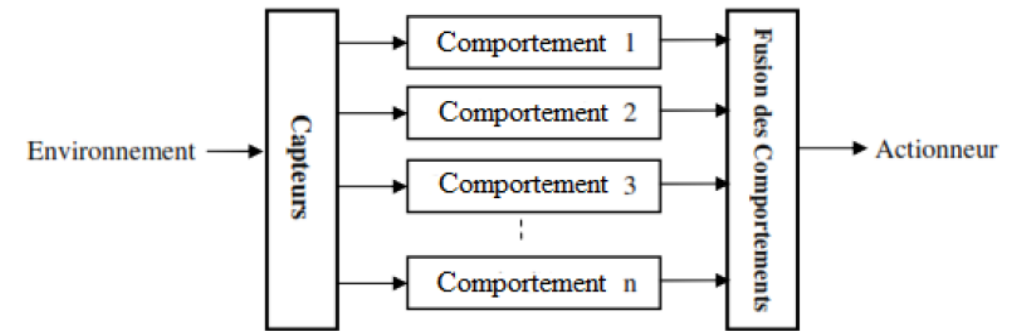
$$d_{rg} = \sqrt{(x_g - x)^2 + (y_g - y)^2}$$

$$\theta_{rg} = \arctan\left(\frac{y_g - y}{x_g - x}\right) - \theta$$

## Synthèse des régulateurs flous

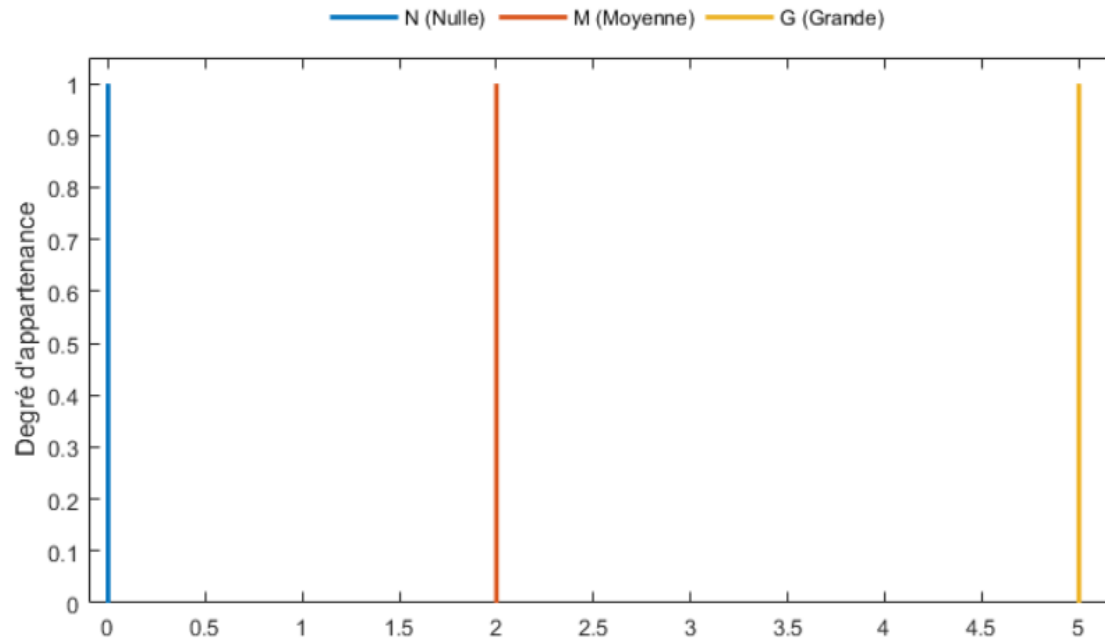
### Architecture de commande

- Navigation réactive par logique floue.
- Comportement de convergence vers un but.
- Comportement d'évitement d'obstacles.



## Synthèse des régulateurs flous

### Régulateur de convergence vers un but (FLC20)



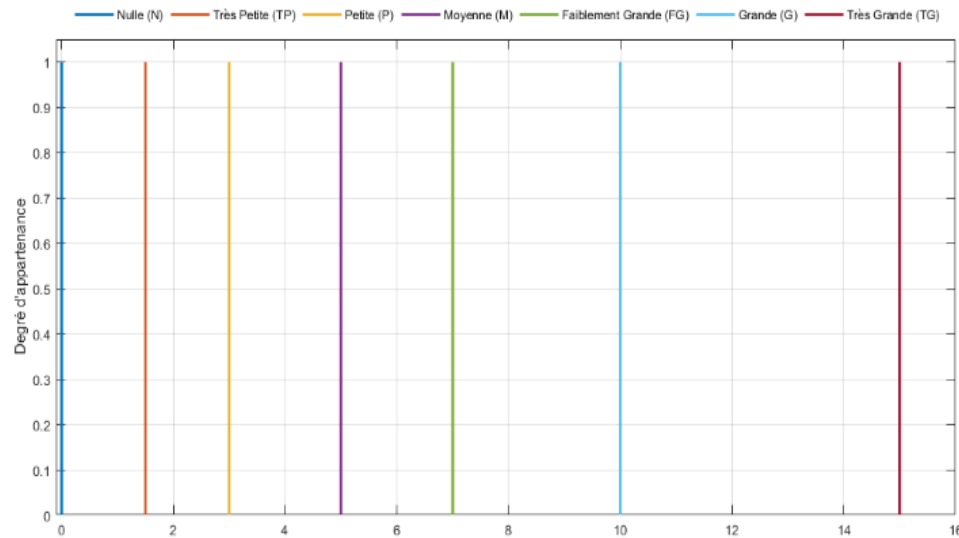
Fonctions d'appartenance des sorties de FLC20 (Convergence)

$\theta_{rg}$ \ $d_{rg}$		$\theta_{rg}$				
		N	NP	Z	PP	P
N	$\omega_D$	M	M	N	N	N
	$\omega_G$	N	N	N	M	M
P	$\omega_D$	G	G	M	M	M
	$\omega_G$	M	M	M	G	G
M	$\omega_D$	G	G	M	M	M
	$\omega_G$	M	M	M	G	G
G	$\omega_D$	G	G	M	M	M
	$\omega_G$	M	M	G	G	G

Les règles de FLC20

## Synthèse des régulateurs flous

### Régulateur de convergence vers un but (FLC49)



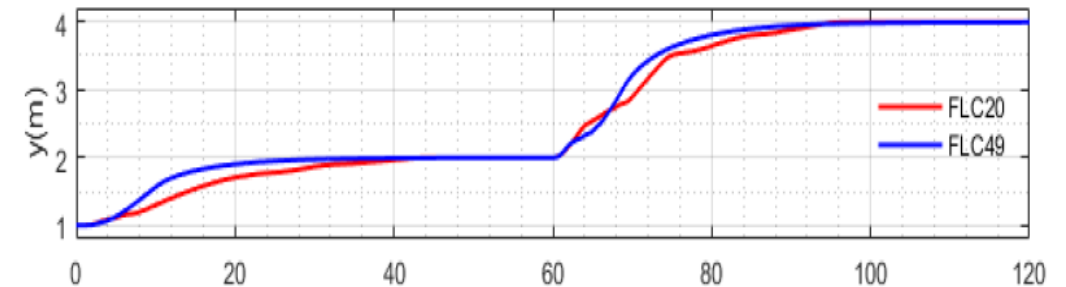
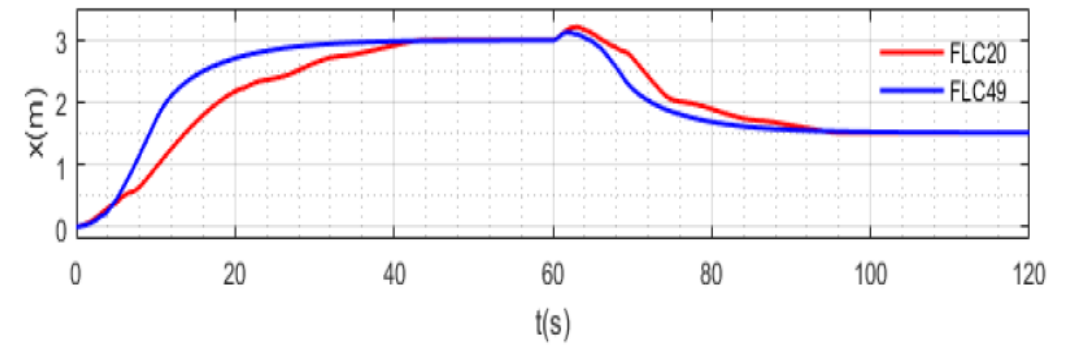
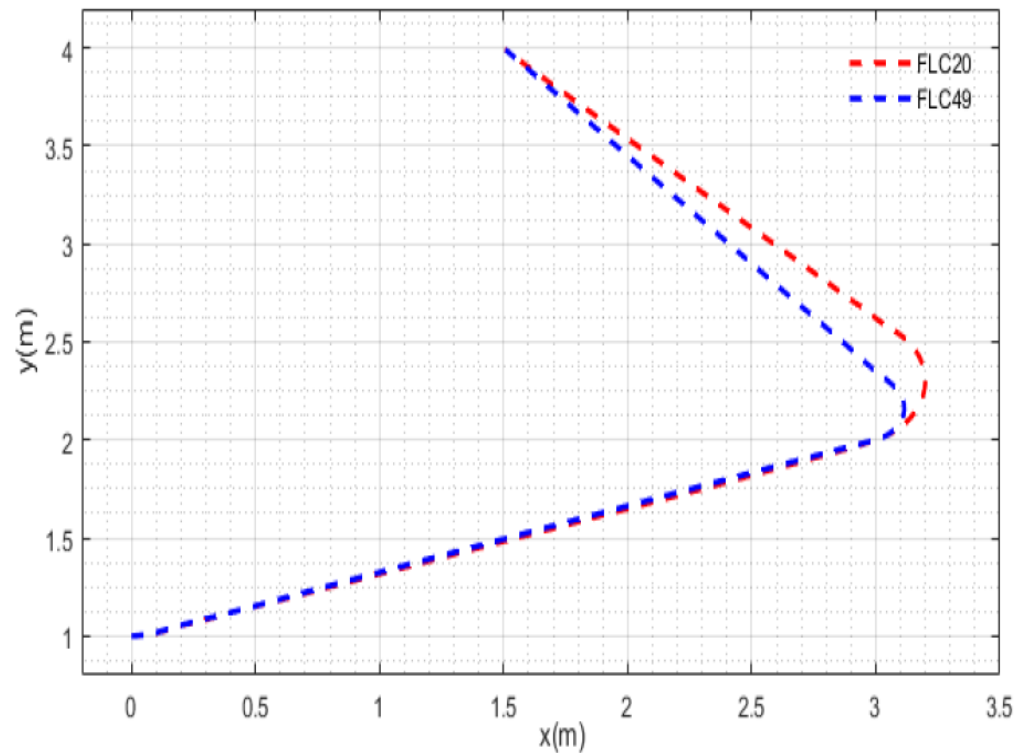
Fonctions d'appartenance des sorties de FLC49 (Convergence)

$\theta_{r_E} \backslash d_{r_E}$		NG	NM	NP	Z	PP	PM	PG
N	$\omega_D$	N	N	N	N	P	P	M
	$\omega_G$	M	P	P	N	N	N	N
TP	$\omega_D$	TP	TP	N	TP	M	FG	G
	$\omega_G$	G	FG	M	TP	N	TP	TP
P	$\omega_D$	TP	TP	TP	P	FG	G	TG
	$\omega_G$	TG	G	FG	P	TP	TP	TP
M	$\omega_D$	TP	TP	TP	M	G	G	TG
	$\omega_G$	TG	G	G	M	TP	TP	TP
MG	$\omega_D$	TP	TP	P	FG	FG	G	TG
	$\omega_G$	TG	G	FG	FG	P	TP	TP
G	$\omega_D$	TP	TP	M	G	FG	G	TG
	$\omega_G$	TG	G	FG	G	M	TP	TP
TG	$\omega_D$	TP	TP	P	TG	FG	G	TG
	$\omega_G$	TG	G	FG	TG	P	TP	TP

Les règles de FLC49

## Synthèse des régulateurs flous

Résultats de comparaison (convergence vers un but)

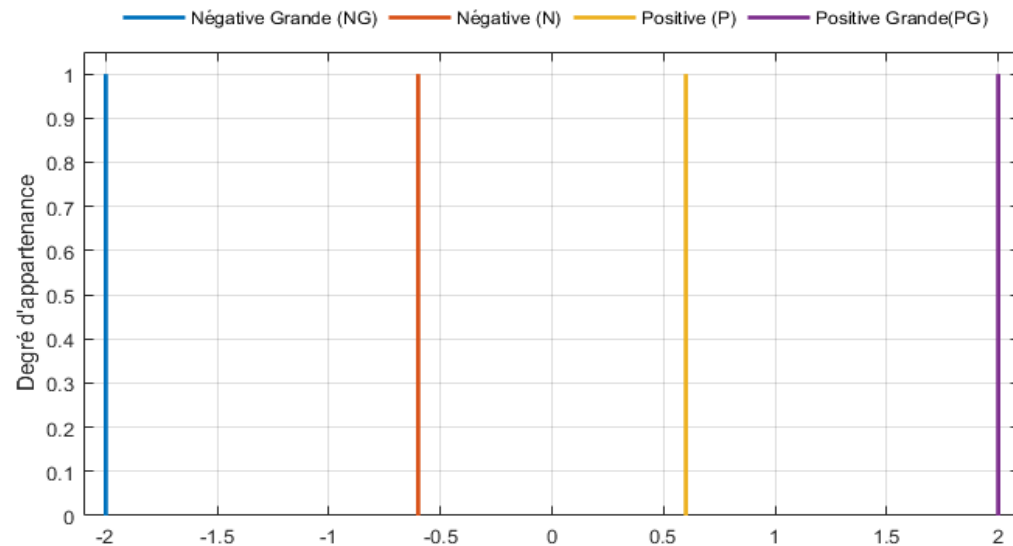


Comparaison



## Synthèse des régulateurs flous

### Régulateur d'évitement d'obstacles

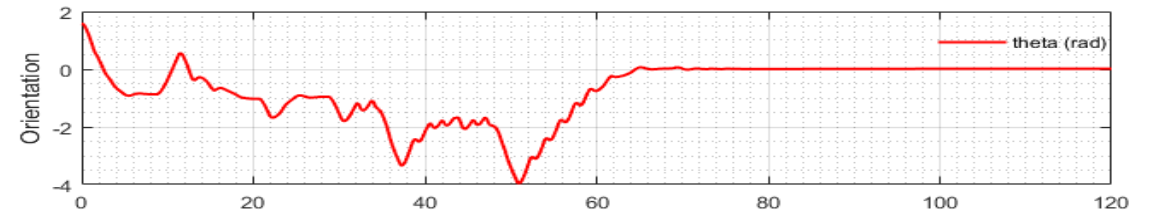
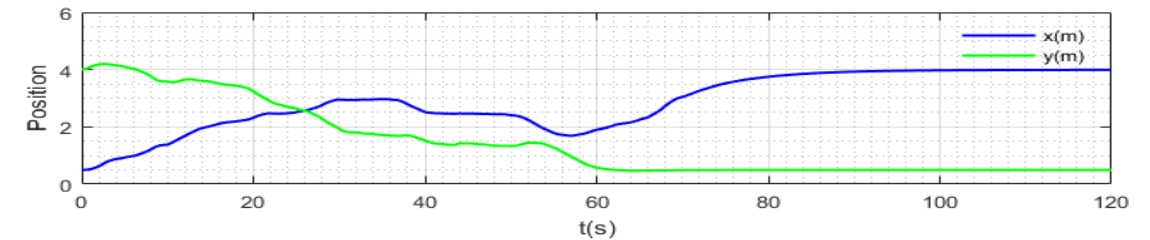
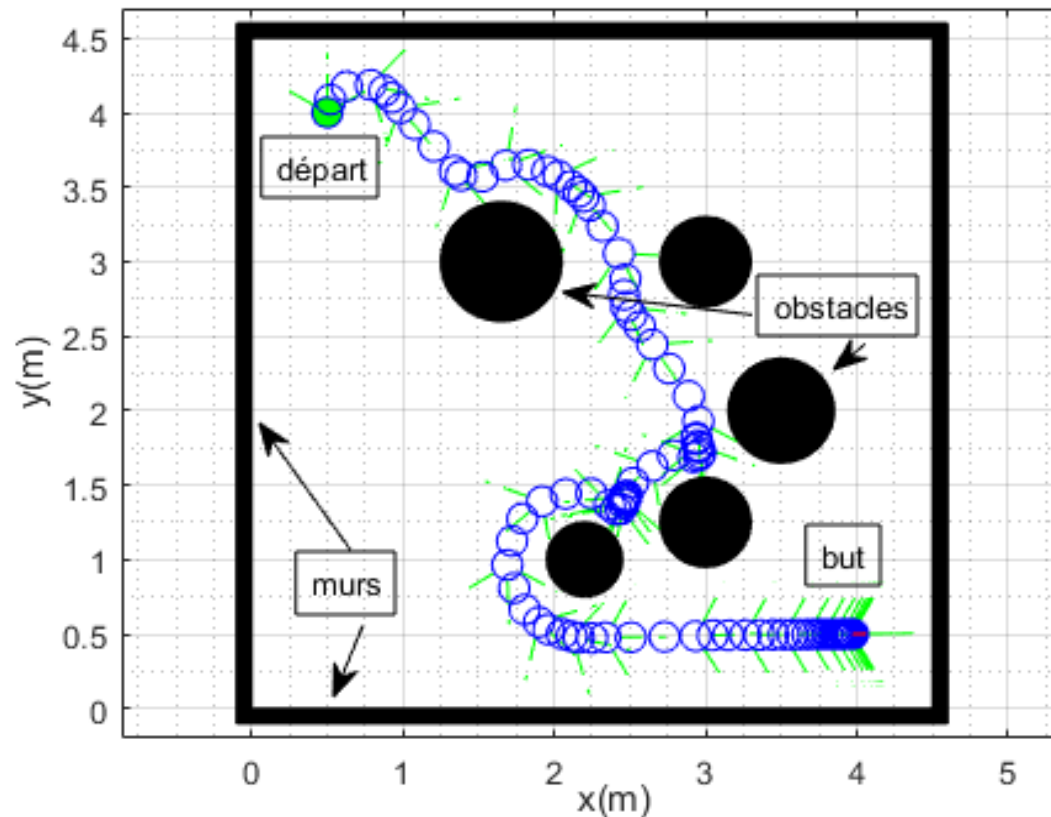


					Entrées			Sorties						
					$d_{ro}^G$	$d_{ro}^A$	$d_{ro}^D$	$\omega_G$	$\omega_D$					
<i>L</i>	<i>P</i>	<i>P</i>	<i>NG</i>	<i>PG</i>	<i>P</i>	<i>P</i>	<i>P</i>	<i>NG</i>	<i>NG</i>	<i>M</i>	<i>P</i>	<i>P</i>	<i>NG</i>	<i>PG</i>
<i>L</i>	<i>P</i>	<i>M</i>	<i>NG</i>	<i>PG</i>	<i>P</i>	<i>P</i>	<i>M</i>	<i>PG</i>	<i>NG</i>	<i>M</i>	<i>P</i>	<i>M</i>	<i>NG</i>	<i>PG</i>
<i>L</i>	<i>P</i>	<i>L</i>	<i>PG</i>	<i>NG</i>	<i>P</i>	<i>P</i>	<i>L</i>	<i>PG</i>	<i>NG</i>	<i>M</i>	<i>P</i>	<i>L</i>	<i>PG</i>	<i>NG</i>
<i>L</i>	<i>M</i>	<i>P</i>	<i>N</i>	<i>PG</i>	<i>P</i>	<i>M</i>	<i>P</i>	<i>NG</i>	<i>NG</i>	<i>M</i>	<i>M</i>	<i>P</i>	<i>N</i>	<i>PG</i>
<i>L</i>	<i>M</i>	<i>M</i>	<i>P</i>	<i>PG</i>	<i>P</i>	<i>M</i>	<i>M</i>	<i>PG</i>	<i>N</i>	<i>M</i>	<i>M</i>	<i>M</i>	<i>NG</i>	<i>NG</i>
<i>L</i>	<i>M</i>	<i>L</i>	<i>P</i>	<i>PG</i>	<i>P</i>	<i>M</i>	<i>L</i>	<i>PG</i>	<i>N</i>	<i>M</i>	<i>M</i>	<i>L</i>	<i>PG</i>	<i>N</i>
<i>L</i>	<i>L</i>	<i>P</i>	<i>N</i>	<i>PG</i>	<i>P</i>	<i>L</i>	<i>P</i>	<i>NG</i>	<i>NG</i>	<i>M</i>	<i>L</i>	<i>P</i>	<i>N</i>	<i>PG</i>
<i>L</i>	<i>L</i>	<i>M</i>	<i>P</i>	<i>PG</i>	<i>P</i>	<i>L</i>	<i>M</i>	<i>P</i>	<i>N</i>	<i>M</i>	<i>L</i>	<i>M</i>	<i>P</i>	<i>P</i>
<i>L</i>	<i>L</i>	<i>L</i>	<i>P</i>	<i>P</i>	<i>P</i>	<i>L</i>	<i>L</i>	<i>P</i>	<i>N</i>	<i>M</i>	<i>L</i>	<i>L</i>	<i>PG</i>	<i>P</i>

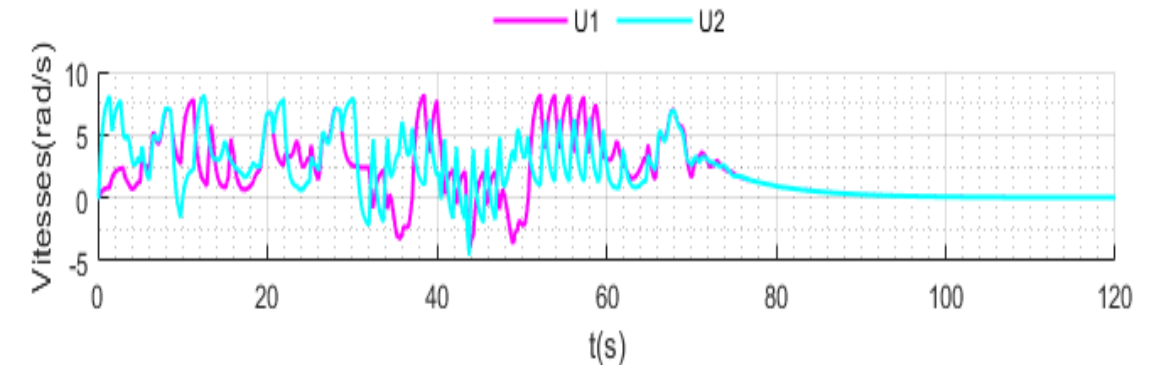


## Synthèse des régulateurs flous

## Résultats de simulation (évitement d'obstacles)



Évolution de la position et de l'orientation



Commandes appliquées

## Synthèse des régulateurs RL flous

### Apprentissage des règles par Fuzzy Q-Learning

Pour chaque épisode d'apprentissage

Pour chaque règle  $i$

$a_i = \operatorname{argmax}_k q[i, k]$  with probability  $1 - \epsilon$

$a_i = \operatorname{random}\{a_k, k = 1, 2, \dots, J\}$  with probability  $\epsilon$

$$a = \sum_{i=1}^N \alpha_i(s(t)) a_i$$

$$Q(s(t), a) = \sum_{i=1}^N \alpha_i(s) \times q[i, a_i]$$

$$V(s(t+1)) = \sum_{i=1}^N \alpha_i(s(t+1)) \cdot \max_k (q[i, q_k]).$$

$$\Delta Q = r(t+1) + \gamma \times V_t(s(t+1)) - Q(s(t), a)$$

$$q[i, a_i] = q[i, a_i] + \eta \cdot \Delta Q \cdot \alpha_i(s(t))$$

$\epsilon$  : Coeff exploitation  
 $\eta$  : Taux d'apprentissage  
 $\alpha$  : Degré d'activation  
 $\gamma$  : facteur de réduction

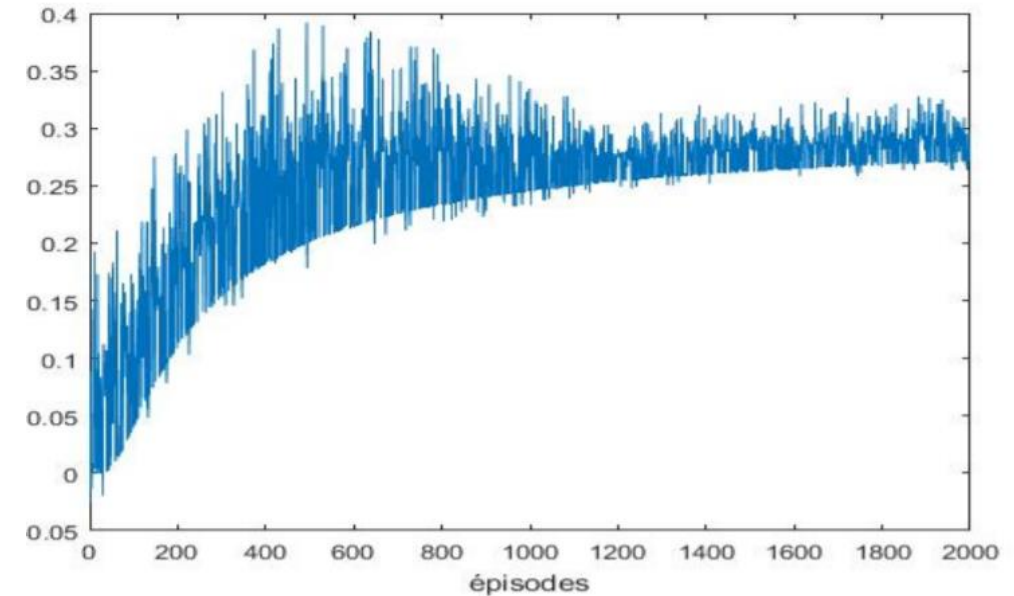
## Synthèse des régulateurs RL flous

### Optimisation Convergence vers un but (RL-FLC20)

#### Convergence vers but

Si " $d_{rg}$  est  $N$  et  $\theta_{rg}$  est  $N$ " Alors  $\omega_{Di} = \omega_{Di1}$  et  $\omega_{Gi} = \omega_{Gi1}$  avec  $Q_i = Q_{i11}$   
 ou  $\omega_{Di} = \omega_{Di1}$  et  $\omega_{Gi} = \omega_{Gi2}$  avec  $Q_i = Q_{i12}$   
 $\vdots$   
 ou  $\omega_{Di} = \omega_{Di3}$  et  $\omega_{Gi} = \omega_{Gi3}$  avec  $Q_i = Q_{i33}$

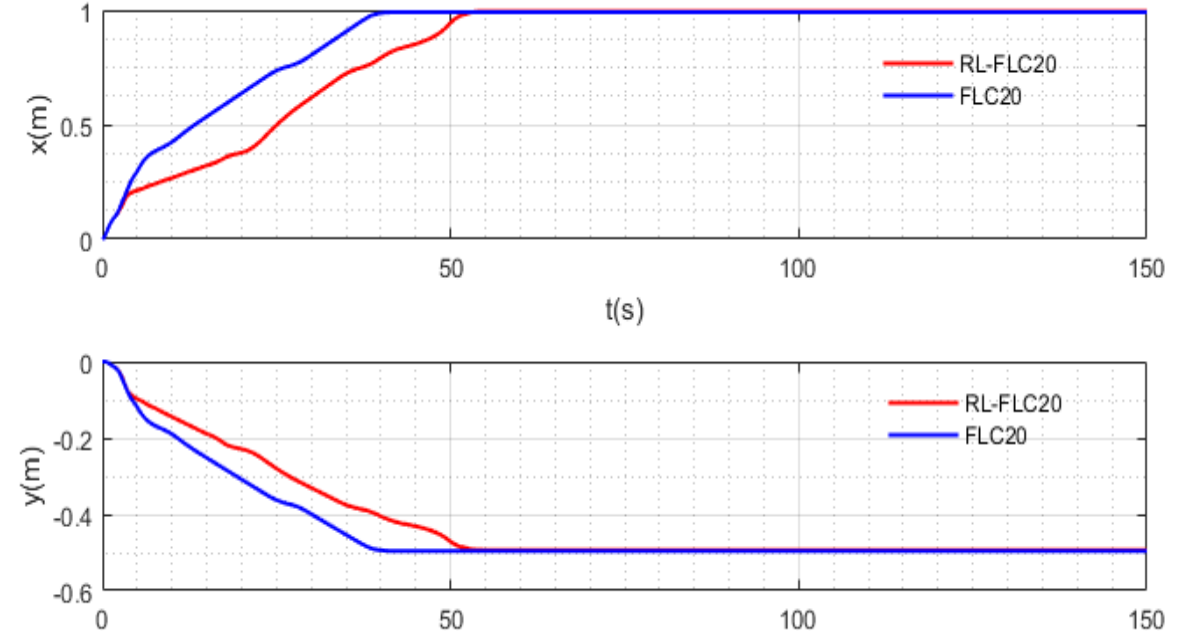
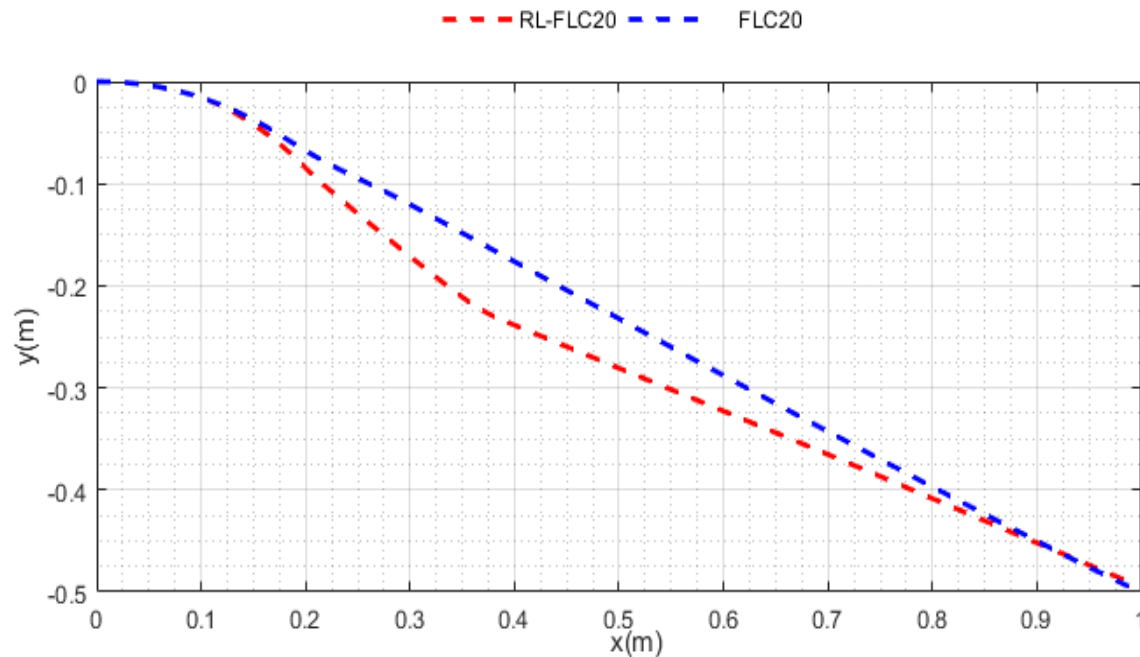
Comportement du robot	Renforcement reçu
Le robot s'approche du but	1
Le robot s'éloigne du but	-1
Le robot atteint le but	10



Renforcements reçus par épisode (Convergence vers un but)

## Synthèse des régulateurs RL flous

Comparaison de convergence vers but (FLC20 et RL-FLC20)



Comparaison

## Synthèse des régulateurs RL flous

### Régulateur de d'évitement d'obstacle

#### Evitement d'obstacles

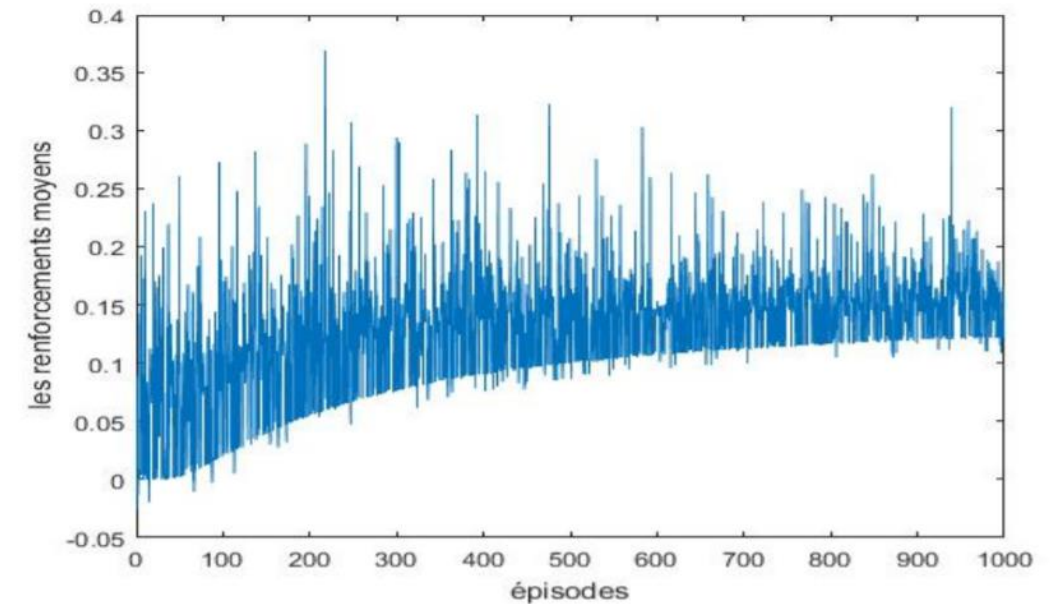
Si " $o$  et  $O_i$ " Alors  $\omega_{Di} = \omega_{Di1}$  et  $\omega_{Gi} = \omega_{Gi1}$  avec  $Q_i = Q_{i11}$

ou  $\omega_{Di} = \omega_{Di1}$  et  $\omega_{Gi} = \omega_{Gi2}$  avec  $Q_i = Q_{i12}$

$\vdots$   $\vdots$

ou  $\omega_{Di} = \omega_{Di4}$  et  $\omega_{Gi} = \omega_{Gi4}$  avec  $Q_i = Q_{i44}$

Comportement du robot	Renforcement reçu
Le robot s'éloigne de l'obstacle	1
Le robot s'approche de l'obstacle	-1
Le robot percute un obstacle	-10

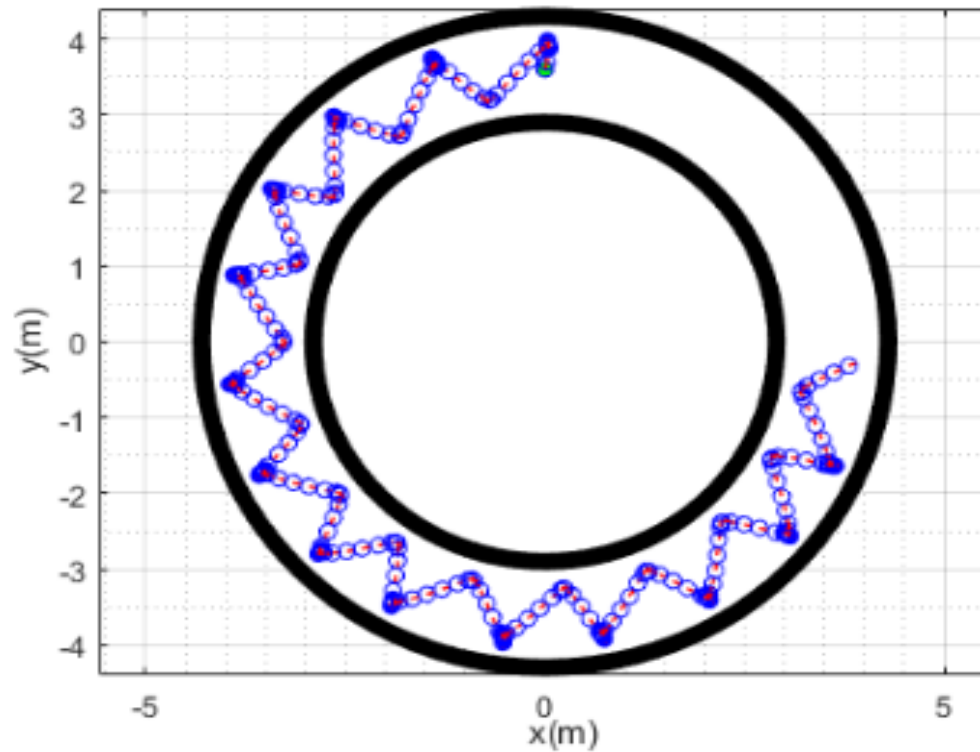


Renforcements reçus par épisode (évitement d'obstacles)

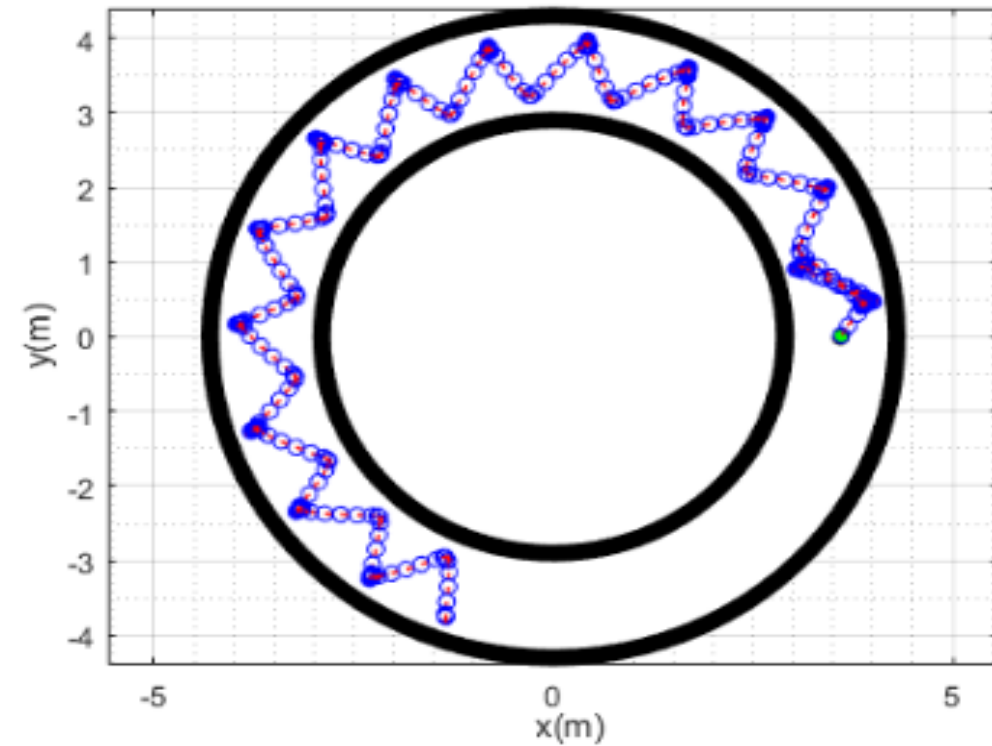


## Synthèse des régulateurs renforcement-flous

Résultats de simulation et de comparaison (évitement)

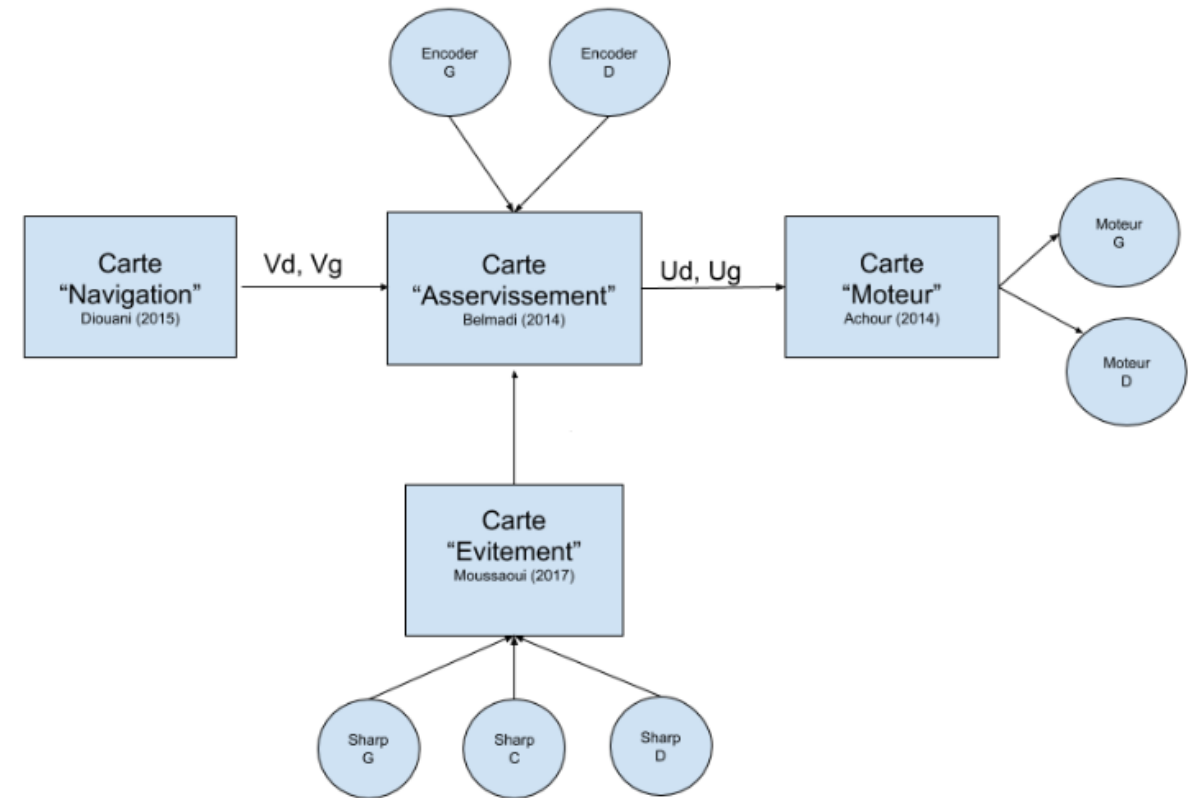


(a) Régulateur flou



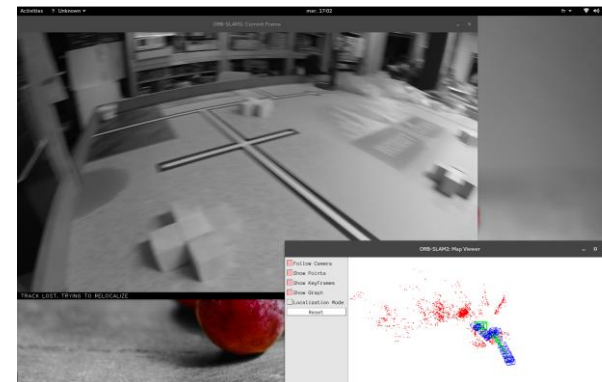
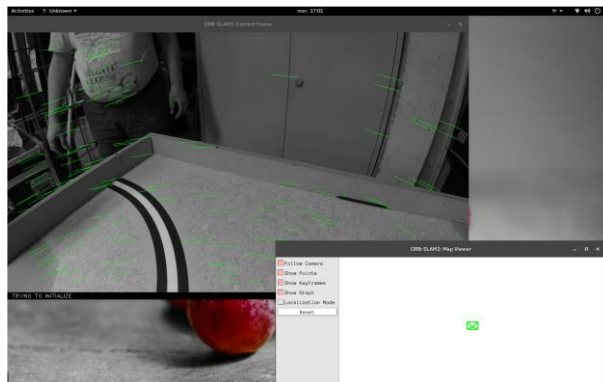
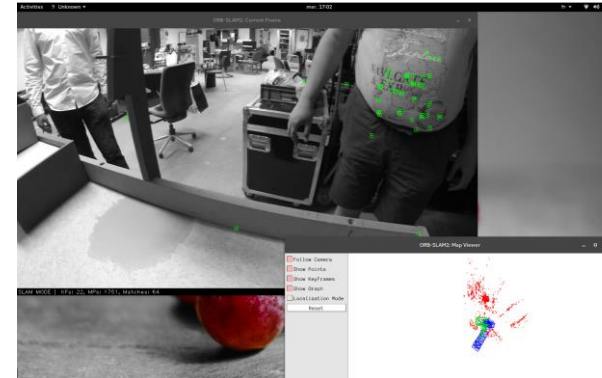
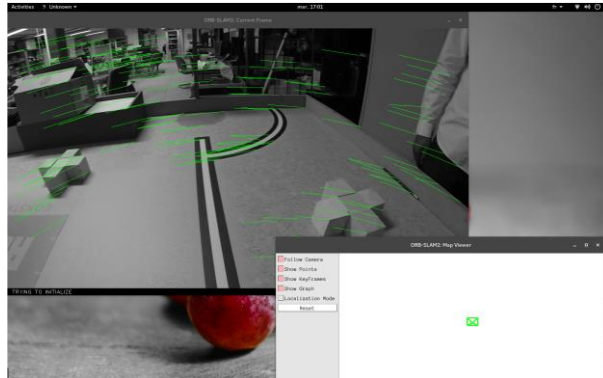
(b) Régulateur renforcement-flou

## Matériel



## Localisation par V-SLAM

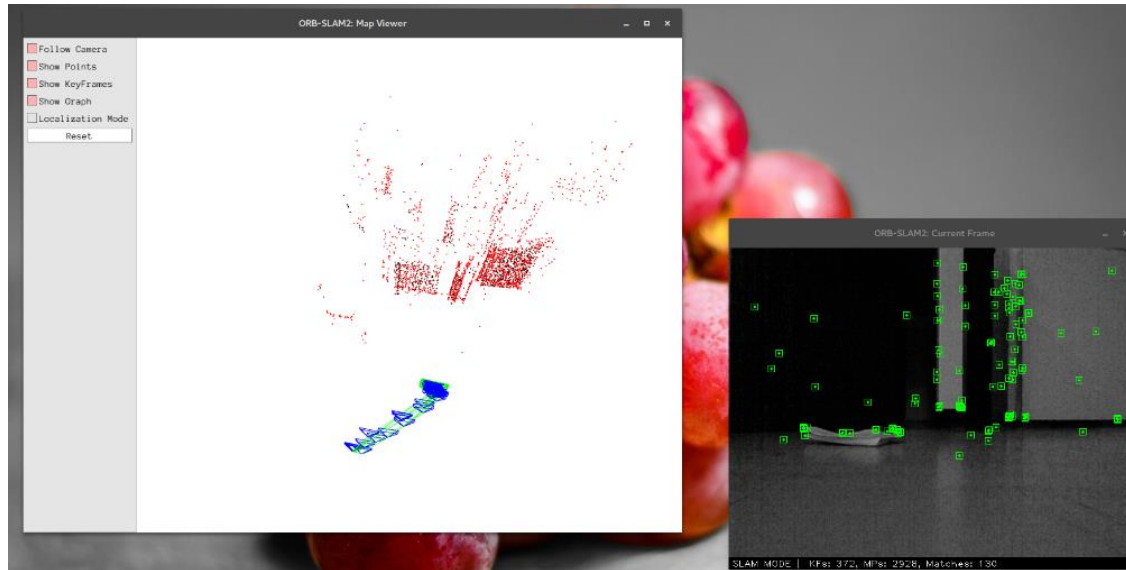
### Test sur le MonoSLAM





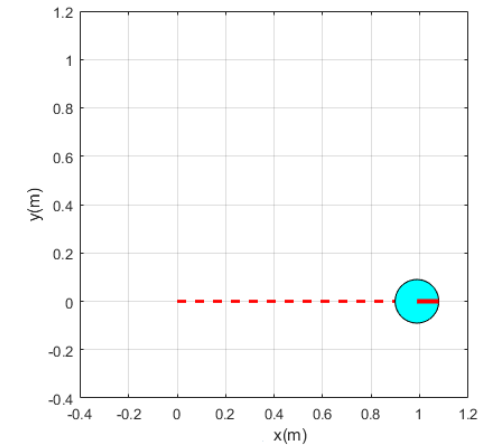
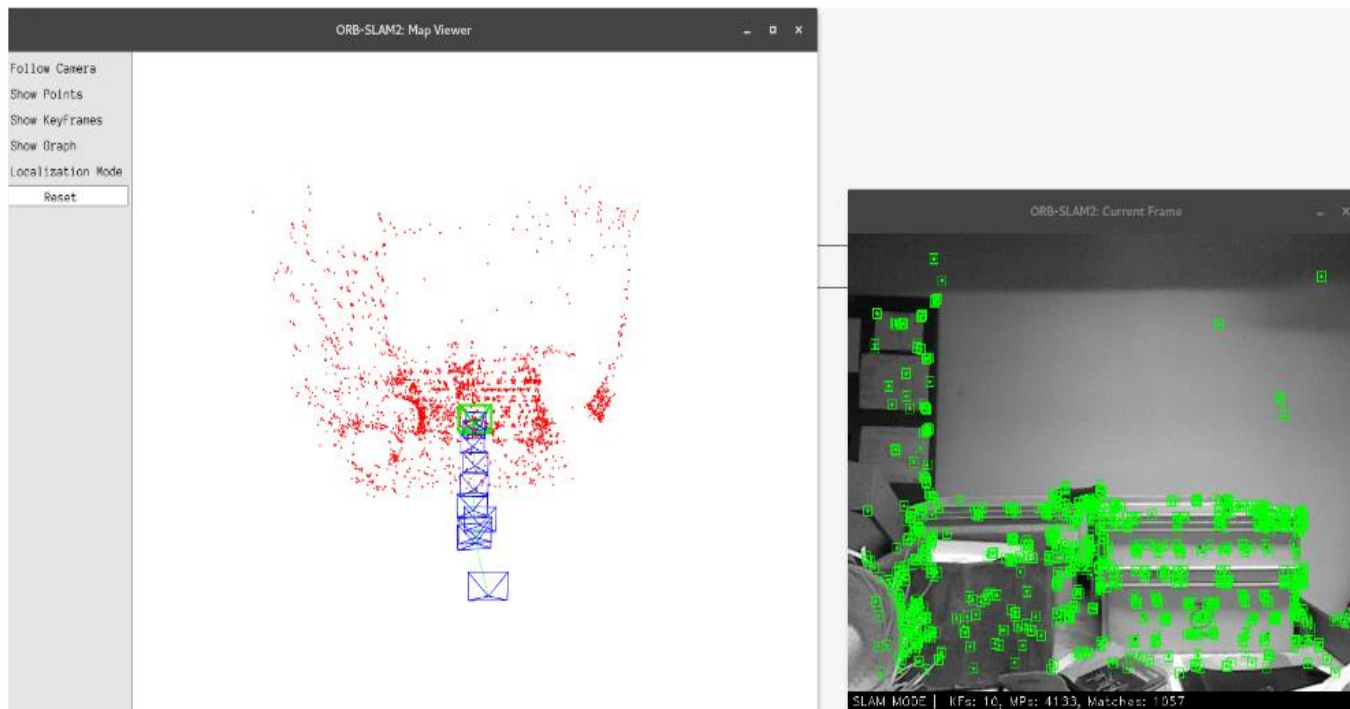
## Localisation par V-SLAM:

### Test sur le RGB-D SLAM

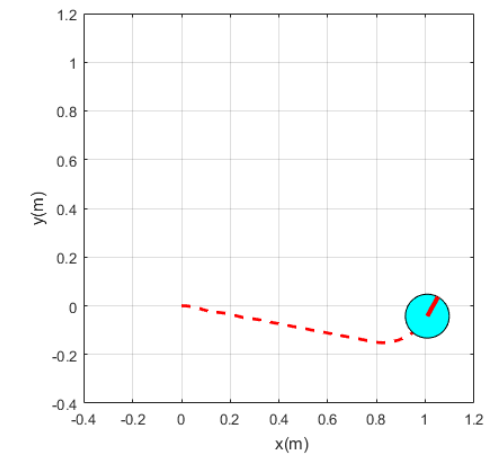


Valeurs réelles			RGB-D SLAM			Odométrie (Cyborg 2017)		
$x$	$y$	$\theta$	$x$	$y$	$\theta$	$x$	$y$	$\theta$
0	0	0	-0.0370	0.0200	$-1.7094^{-4}$	0	0	0
40	40	0	39.9858	40.6389	-0.0975	40.53	42.87	0.01
0	0	0	-0.5027	-0.6823	-0.0111	-2.06	2.71	0.01
40	40	0	40.8477	40.5749	-0.0064	38	43.76	0
0	0	0	0.0313	-4.2988	-0.0264	-4.75	5.96	-0.01
40	40	0	40.5742	41.2686	-0.0954	37.01	46.09	-0.02

## Navigation floue par V-SLAM



Navigation floue avec odométrie



Navigation floue avec V-SLAM

## Navigation renforcement -floue par V-SLAM (vidéo)

## Conclusion

- Navigation floue efficace.
- Navigation renforcement-floue par apprentissage.
- Efficacité de ORB-SLAM2 pour la localisation.
- Navigation floue par ORB-SLAM2 en mode SLAM.

## Perspectives

- Extension aux environnements dynamiques.
- Hybridation de la logique floue et les réseaux de neurones.

Merci pour votre attention