

# Simulator-Based Verification of Autonomous Vehicles: Agent-Based Test Generation

1:09pm, 7 October 2020

## 1 Introduction

- **Time** is implicit and discrete, with agents choosing actions at each timestep, and a time resolution constant specifying the simulation time between timesteps
- **State space** is continuous and feature-based techniques are probably more effective than discretisation
- **Actions** may be multi-body and/or multi-effector
- **Action space** is continuous and discretisation typically requires durative actions
- **Observability** is full, meaning that a sensor model is not required (although partial observability can be simulated)
- **Objectives** are goals states with transition costs/rewards, meaning in general that the horizon is indefinite, but a finite-horizon may also be imposed
- **Model** for transitions and objectives is unknown, but can be sampled through interaction with a simulator
- **Game** involving two or more players (agents), including a potentially adversarial design under verification (i.e. the ego agent)
- **Test generation** may be online (providing incomplete tests) or offline (providing complete tests)

**Note:** Ideally a test should specify the optimal effector-actions to execute in a given state with respect to some test objective. A complete test then would be one that specifies effector-actions to execute in every reachable state. Online testing cannot guarantee complete tests because effector-actions are only determined for states that are actually encountered in a given run, and this does not typically cover all reachable states. For example, if a simulation involves non-deterministic action-selection, non-deterministic action-effects, or exogenous events, then these sources of uncertainty mean that a single execution of online testing will only reveal one of many possible runs. It may be possible to impose restrictions on a simulation so as to avoid many of these sources of uncertainty, but the actions of the ego agent remain uncontrollable. For example, a test may be complete for a given ego policy, but may be incomplete when faced with a different ego policy (since that policy may cause the system to transition into a state that was not accounted for in the original test). It follows that a test is only complete if it accounts for every source of uncertainty and also works for every possible ego policy. This is a very strong requirement and is unlikely to be satisfiable in practice (e.g. if agents execute actions simultaneously, then the ego policy is a priori unknown). This seems to imply that the aim of our work is not *test generation* but rather *online testing with optional capturing of incomplete tests*: re-running such tests could be useful for repeatability, but test execution may fail, and there would be no guarantees that tests remain optimal with respect to their original test objectives.

## 2 Preliminaries

We rely on some standard mathematical notation:  $v_i$  is an element of vector  $\mathbf{v} = (v_1, \dots, v_n)$  with  $\mathbf{v}_{-i} = (v_1, \dots, v_{i-1}, v_{i+1}, \dots, v_n)$  the subvector of  $\mathbf{v}$  excluding  $v_i$ ,  $|S|$  is the cardinality of set  $S$ ,  $2^S$  is the powerset of  $S$ ,  $\Delta(S)$  is the set of probability distributions over  $S$ ,  $\mathbb{R}$  is the set of real numbers with  $\mathbb{R}^{\geq 0}$  the subset of non-negative real numbers, and  $\mathbb{N}$  is the set of natural numbers. A function  $f : X \rightarrow Y$  is a surjection if for each  $y \in Y$  there exists some  $x \in X$  such that  $f(x) = y$ . Given a function  $f : X \rightarrow Y$ , the inverse image of function  $f$  is a function  $f_{B,N}^{-1} : Y \rightarrow 2^X$  defined for each  $y \in Y$  as  $f_{B,N}^{-1}(y) = \{x \in X \mid f(x) = y\}$ .

### 2.1 Normal-Form Games

An ( $n$ -player) normal-form game is a tuple  $(N, A, \mathbf{u})$  where  $N = \{1, \dots, n\}$  is a finite set of players,  $A = A_1 \times \dots \times A_n$  is a set of action profiles with  $A_i$  the (possibly infinite) set of actions available to player  $i \in N$ , and  $\mathbf{u} = (u_1, \dots, u_n)$  is a profile of utility functions with  $u_i : A \rightarrow \mathbb{R}$  the utility function for  $i$ . A (mixed) strategy for player  $i \in N$  is a probability distribution  $\psi \in \Delta(A_i)$  with  $\boldsymbol{\psi} \in \Delta(A_1) \times \dots \times \Delta(A_n)$  a (mixed) strategy profile. A strategy  $\psi$  for player  $i \in N$  is a pure strategy if  $\psi(a) = 1$  for some  $a \in A_i$ . A strategy profile  $\boldsymbol{\psi}$  is a pure strategy profile if  $\psi_i$  is a pure strategy for each  $i \in N$ . For convenience, we may denote a pure strategy  $\psi$  for player  $i \in N$  directly by the action  $a \in A_i$  such that  $\psi(a) = 1$ . Likewise, we may denote a pure strategy profile  $\boldsymbol{\psi}$  directly by the action profile  $\mathbf{a} \in A$  such that  $\psi_i(a_i) = 1$  for each player  $i \in N$ . The expected utility of a strategy profile  $\boldsymbol{\psi}$  is defined as:

$$u_i(\boldsymbol{\psi}) = \sum_{\mathbf{a} \in A} u_i(\mathbf{a}) \prod_{j \in N} \psi_j(a_j)$$

In game theory, a solution concept is a prediction of how a game will be played. Solution concepts are typically based on the idea that each player will try to maximise their expected utility under the assumption that other players will do the same. A feasible deviation from strategy profile  $\boldsymbol{\psi}$  by player  $i \in N$  is a strategy profile  $\boldsymbol{\psi}' = (\boldsymbol{\psi}_{-i}, \psi')$  where  $\psi'$  is a strategy for  $i$ . A strategy profile  $\boldsymbol{\psi}$  is a Nash equilibrium if no player  $i \in N$  has a feasible deviation  $\boldsymbol{\psi}'$  from  $\boldsymbol{\psi}$  such that  $u_i(\boldsymbol{\psi}') > u_i(\boldsymbol{\psi})$ . Every normal-form game with a finite set of action profiles is guaranteed to have at least one (not necessarily pure) Nash equilibrium [3]. A Nash equilibrium is also guaranteed to exist for certain classes of normal-form game with an infinite set of action profiles [].

### 2.2 Markov Games

A (fully observable) Markov game is a tuple  $(N, S, A, T, \mathbf{R})$  where  $N = \{1, \dots, n\}$  is a finite set of players,  $S$  is a (possibly infinite) set of states,  $A = A_1 \times \dots \times A_n$  is a set of action profiles with  $A_i$  the (possibly infinite) set of actions available to player  $i \in N$ ,  $T : S \times A \rightarrow \Delta(S)$  is a (stochastic) transition function, and  $\mathbf{R} = (R_1, \dots, R_n)$  is a profile of reward functions with  $R_i : S \times A \times S \rightarrow \mathbb{R}$  the reward function for player  $i \in N$ . A Markov game is zero-sum if  $\sum_{i \in N} R_i(s, \mathbf{a}, s') = 0$  for each  $\mathbf{a} \in A$  and all  $s, s' \in S$ .

Markov games (also known as stochastic games [6]) are composed of a sequence of normal-form games, called stage games, where players seek to optimise their return for the whole game rather than for individual stage games. The current stage game is determined by the current state, and players transition between states by executing actions simultaneously. The game is said to be fully observable because players are assumed to observe the current (shared) state at each decision step. Let  $T(s, \mathbf{a}, s')$  denote the probability of transitioning to state  $s' \in S$  after executing action profile  $\mathbf{a} \in A$  in state  $s \in S$  according to probability distribution  $T(s, \mathbf{a})$ . The probability of transitioning to state  $s' \in S$  after executing strategy profile  $\psi$  in state  $s \in S$  is:

$$T(s, \psi, s') = \prod_{\mathbf{a} \in A} T(s, \mathbf{a}, s') \prod_{j \in N} \psi_j(a_j)$$

The stage game for state  $s \in S$  is a normal-form game  $(N, A, \mathbf{u})$  where the utility function  $u_i$  for player  $i \in N$  is defined for each action profile  $\mathbf{a} \in A$  as the immediate expected reward:

$$u_i(\mathbf{a}) = \sum_{s' \in S} T(s, \mathbf{a}, s') R_i(s, \mathbf{a}, s')$$

Markov games subsume numerous other frameworks, including repeated games when  $|S| = 1$ , and Markov decision processes (MDPs) when  $|N| = 1$ .

Solutions to Markov games are represented as functions, called policies, that map *histories* to strategies. An execution is a possibly infinite sequence of states and action profiles  $(s_1, \mathbf{a}_1, s_2, \mathbf{a}_2, \dots)$ . A history of length  $t$  is a finite execution  $h_t = (s_1, \mathbf{a}_1, \dots, \mathbf{a}_{t-1}, s_t)$  ending in a state. Let  $H_t$  be the set of histories of length  $t$  with  $D = \{1, 2, \dots, t_{\max}\}$  the set of decision-steps up to horizon  $t_{\max} \in \mathbb{N} \cup \{\infty\}$  and  $H = \{h \in H_t \mid t \in D\}$  the set of histories up to  $t_{\max}$ . A (stochastic or mixed) policy for player  $i \in N$  is a function  $\pi_i : H' \rightarrow \Delta(A_i)$  where  $H' \subseteq H$ . Let  $\pi_i(h, a)$  denote the probability that player  $i \in N$  should execute action  $a \in A_i$  in history  $h \in H'$  according to strategy  $\pi_i(h)$ . A policy  $\pi_i$  for player  $i \in N$  is a complete policy if  $H' = H$ , otherwise  $\pi_i$  is a partial policy. A policy  $\pi_i$  for player  $i \in N$  is a deterministic (or pure) policy if  $\pi_i(h, a) = 1$  for each  $h \in H'$  and some  $a \in A_i$ . For convenience, we may denote a deterministic policy  $\pi_i$  for player  $i \in N$  directly by a policy  $\pi'_i : H' \rightarrow A$  where  $\pi'_i(h) = a$  such that  $\pi_i(h, a) = 1$ . A policy  $\pi_i$  for player  $i \in N$  is Markovian if  $\pi_i(h_t) = \pi_i(h_{t'})$  for all  $h_t, h_{t'} \in H'$  such that  $t = t'$  and  $s_t = s_{t'}$ . A Markovian policy for player  $i \in N$  may be written as  $\pi_i : X \rightarrow \Delta(A_i)$  where  $X \subseteq S \times D$ . A Markovian policy  $\pi_i$  for player  $i \in N$  is a stationary policy if  $\pi_i(s, t) = \pi_i(s, t')$  for all  $(s, t), (s, t') \in X$ , otherwise  $\pi_i$  is a non-stationary policy. A stationary Markovian policy for player  $i \in N$  may be written as  $\pi_i : S' \rightarrow \Delta(A_i)$  where  $S' \subseteq S$ . A tuple  $\pi = (\pi_1, \dots, \pi_n)$  is a policy profile over  $H' \subseteq H$  if  $\pi_i$  is a policy over  $H'$  for each  $i \in N$ . The strategy profile for history  $h \in H'$  according to policy profile  $\pi$  is  $\pi(h) = (\pi_1(h), \dots, \pi_n(h))$ .

Solutions are only well-defined for certain classes of Markov game. Two common examples are finite-horizon Markov games and infinite-horizon (discounted reward) Markov games. A finite-horizon Markov game is a tuple  $(\Gamma, t_{\max})$  where  $\Gamma$  is a Markov game and  $t_{\max} \in \mathbb{N}$  is a horizon. An infinite-horizon Markov game is a tuple  $(\Gamma, \gamma)$  where  $\Gamma$  is a Markov game and  $\gamma \in [0, 1)$  is a discount factor. It is known that attention can be restricted to stationary Markovian policies in the case of infinite-horizon Markov

games, and to non-stationary Markovian policies in the case of finite-horizon Markov games. Unless stated otherwise, we will assume hereafter that all policies are either stationary or non-stationary Markovian policies. For an infinite-horizon Markov game  $(\Gamma, \gamma)$ , the expected value of policy profile  $\pi$  for player  $i \in N$  in state  $s \in S$  is:

$$V_i(\pi, s) = \sum_{s' \in S} T(s, \pi(s), s') [R_i(s, \pi(s), s') + \gamma V_i(\pi, s')]$$

For a finite-horizon Markov game  $(\Gamma, t_{\max})$ , the expected value of policy profile  $\pi$  for player  $i \in N$  in state  $s \in S$ , with  $t \in D \cup \{0\}$  remaining decision-steps, is:

$$V_i(\pi, s, t) = \begin{cases} \sum_{s' \in S} T(s, \pi(s, t), s') [R_i(s, \pi(s, t), s') + V_i(\pi, s', t-1)] & \text{if } t > 0 \\ 0 & \text{if } t = 0 \end{cases}$$

A feasible deviation from policy profile  $\pi$  by player  $i \in N$  is a policy profile  $\pi' = (\pi_{-i}, \pi')$  where  $\pi'$  is a policy for  $i$ . A policy profile  $\pi$  is a Nash equilibrium if no player  $i \in N$  has a feasible deviation  $\pi'$  from  $\pi$  such that  $V_i(\pi', h) \geq V_i(\pi, h)$  for each  $h \in H'$  and  $V_i(\pi', h') > V_i(\pi, h')$  for some  $h' \in H'$ .

**To do:** Confirm that this definition of a Nash equilibrium is correct with respect to histories (or states, or time-indexed states).

### 2.3 Stochastic Shortest Path Games

A (fully observable) stochastic shortest path (SSP) game is a tuple  $(N, S, A, T, C, G)$  where  $N = \{1, \dots, n\}$  is a finite set of players,  $S$  is a (possibly infinite) set of states,  $A = A_1 \times \dots \times A_n$  is a set of action profiles with  $A_i$  the (possibly infinite) set of actions available to player  $i \in N$ ,  $T : S \times A \rightarrow \Delta(S)$  is a (stochastic) transition function,  $C = (C_1, \dots, C_n)$  is a profile of cost functions with  $C_i : S \times A \times S \rightarrow \mathbb{R}^{\geq 0}$  the (non-negative) cost function for player  $i \in N$ , and  $G = G_1 \cup \dots \cup G_n$  is a non-empty set of goal states with  $G_i \subseteq S$  the (possibly empty, possibly infinite) set of goal states for player  $i \in N$ . SSP games subsume numerous other frameworks, including finite- and infinite-horizon Markov games as well as SSP MDPs (which in turn subsume finite- and infinite-horizon MDPs as well as classical planning) [1, 4, 5, 2].

**Note:** I have only found definitions of 2-player zero-sum SSP games in the literature, so this definition of an  $n$ -player general-sum SSP game is only an informed guess. The definition says that all players should seek to minimize their cost for the duration of the game, as is the case in Markov games, but that at least one player should seek to arrive in a goal state where the game will terminate. The definition does not explicitly capture the notion that an adversary would seek to block the goal-directed player from achieving their goal, but the fact that this notion appears in existing definitions may be due to those definitions focusing on zero-sum games.

### 3 Framework

**Definition 1.** An agent-body-effector model is a tuple  $(N, B, E, f_{B,N}, f_{E,B})$  where:

- $N = \{1, \dots, n\}$  is a finite set of agents
- $B$  is a finite set of bodies with  $f_{B,N} : B \rightarrow N$  a surjection from bodies to agents
- $E$  is a finite set of effectors with  $f_{E,B} : E \rightarrow B$  a surjection from effectors to bodies and  $A(e)$  the non-empty (possibly infinite) set of actions available to  $e \in E$

**Corollary 1.**  $2 \leq |N| \leq |B| \leq |E|$ .

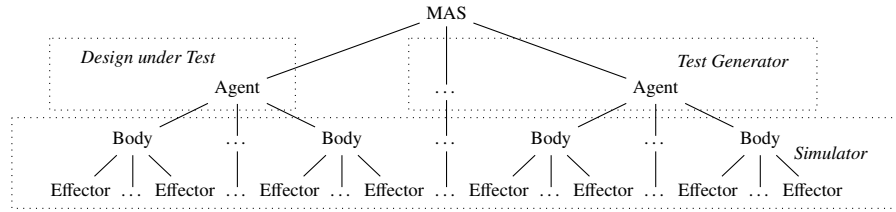


Figure 1: Agent-body-effector model

The set of bodies associated with agents  $N' \subseteq N$  is defined as  $B(N') = \{b \in f_{B,N}^{-1}(i) \mid i \in N'\}$ . The set of effectors associated with bodies  $B' \subseteq B$  is defined as  $E(B') = \{e \in f_{E,B}^{-1}(b) \mid b \in B'\}$ . Thus, the set of bodies  $B$  (resp. agents  $N$ ) in an agent-body-effector model induces a partition of the set of effectors  $E$ , as illustrated in Figure 1. The set of (multi-)actions available to effectors  $E' \subseteq E$  is defined as  $A(E') = \times_{e \in E'} A(e)$ . Thus, there exists a non-empty set of (multi-)actions available to each body (resp. agent).

**Definition 2.** Let  $(N, B, E, f_{B,N}, f_{E,B})$  be an agent-body-effector model. A simulator-based testing game is a stochastic shortest path game  $(N, S, A, T, C, G)$  if:

- $N = \{1, \dots, n\}$  such that  $n \geq 2$  with 1 the ego agent and  $2, \dots, n$  the tester agents
- $A_i = A(E(B(\{i\})))$  is the set of actions available to agent  $i \in N$

**To do:** Formally define the notion of proper policies.

**Lemma 1.** Let  $\mathcal{G}$  be a simulator-based testing game. There exists a proper policy profile for  $\mathcal{G}$  if and only if  $\mathcal{G}$  does not exhibit dead-ends.

*Proof Sketch.* The proof may be analogous to the proofs for SSP MDPs [2], assuming it has not been proven already.  $\square$

**Definition 3.** A time-limited test generation game is a tuple  $(\mathcal{G}, t_{\max})$  where  $\mathcal{G}$  is a test generation game and  $t_{\max} \in \mathbb{N}$  is a (finite-)horizon.

**Proposition 1.** Let  $(\mathcal{G}, t_{\max})$  be a time-limited test generation game. There exists a proper policy profile for  $(\mathcal{G}, t_{\max})$ .

*Proof Sketch.* The set of states  $S$  can be translated into a new set of states  $S' = S \times \{0, \dots, t_{\max}\}$ . The set of goal states  $G$  can be translated into a new set of goal states  $G' = (G \times \{0, \dots, t_{\max}\}) \cup (S \times \{0\})$ . The translated game  $(\mathcal{G}', t_{\max})$  does not exhibit dead-ends because for any state  $s \in S'$  there always exists a state  $s' \in S \times \{0\}$  such that  $s'$  is reachable from  $s$ . Thus, it follows from Lemma 1 that a proper policy profile for  $G$  is guaranteed to exist.  $\square$

## 4 Experiments

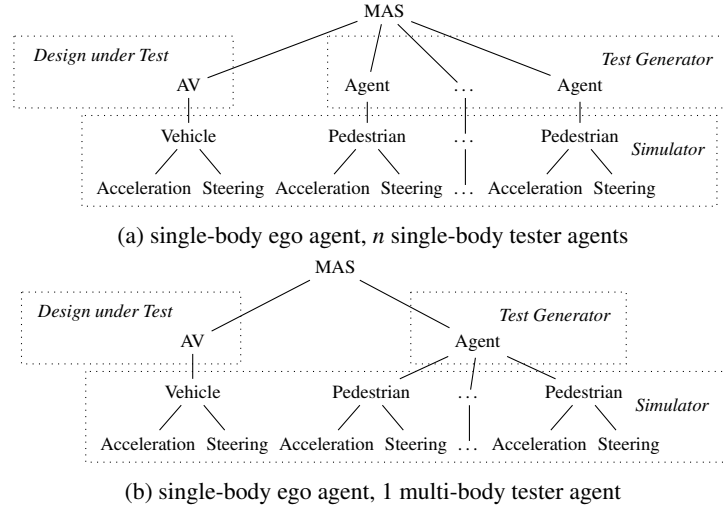


Figure 2: Agent-based test generation in CAV-GYM:PEDESTRIANS

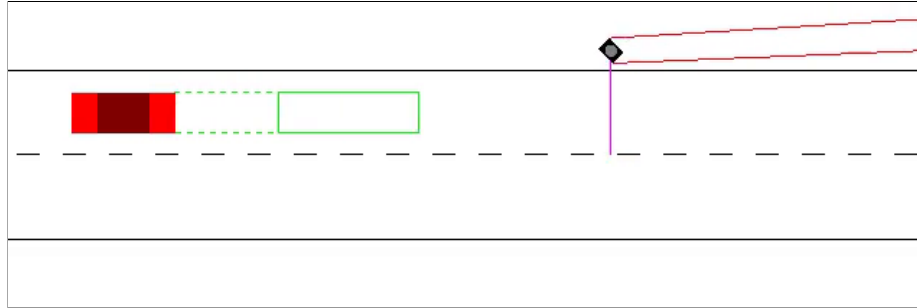


Figure 3: CAV-GYM:PEDESTRIANS



## References

- [1] Dimitri P. Bertsekas. *Dynamic Programming and Optimal Control, Vols. I and II*. Athena Scientific, 1995.
- [2] Mausam and Andrey Kolobov. *Planning with Markov decision processes: An AI perspective*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2012.
- [3] John Nash. Non-cooperative games. *Annals of Mathematics*, 54(2):286–295, 1951.
- [4] Stephen D. Patek. *Stochastic and shortest path games: theory and algorithms*. PhD thesis, Massachusetts Institute of Technology, 1997.
- [5] Stephen D. Patek and Dimitri P. Bertsekas. Stochastic shortest path games. *SIAM Journal on Control and Optimization*, 37(3):804–824, 1999.
- [6] Lloyd S. Shapley. Stochastic games. *Proceedings of the National Academy of Sciences*, 39(10):1095–1100, 1953.

## Appendix: Kinematics

A state is a tuple  $s_i = (p_i, v_i, o_i)$  where  $p_i \in \mathbb{R}^2$  is a position denoted  $p_i = (p_i^x, p_i^y)$ ,  $v_i \in [v_{\min}, v_{\max}]$  is a velocity with  $v_{\min}, v_{\max} \in \mathbb{R}^{\geq 0}$  such that  $v_{\min} < v_{\max}$ , and  $o_i \in (-\pi, \pi]$  is an orientation (in radians). Let  $b \in \mathbb{R}^{\geq 0}$  be a wheelbase constant. The positions of front and rear wheels  $f_i, r_i \in \mathbb{R}^2$  in state  $s_i = (p_i, v_i, o_i)$  are:

$$f_i = \left( p_i + \frac{b}{2} \cdot (\cos o_i, \sin o_i) \right) \quad (1)$$

$$r_i = \left( p_i - \frac{b}{2} \cdot (\cos o_i, \sin o_i) \right) \quad (2)$$

A body has two effectors: throttle and steering. An action is a tuple  $a_i = (t_i, e_i)$  where  $t_i \in [t_{\min}, t_{\max}]$  is a throttle with  $t_{\min}, t_{\max} \in \mathbb{R}$  such that  $t_{\min} < t_{\max}$ , and  $e_i \in (e_{\min}, e_{\max}]$  is a steering angle (in radians) with  $e_{\min}, e_{\max} \in (-\pi, \pi]$  such that  $e_{\min} < e_{\max}$ . Let  $\lambda \in \mathbb{R}^{>0}$  be a time resolution constant. If action  $a_i = (t_i, e_i)$  is executed in state  $s_i = (p_i, v_i, o_i)$ , then the successor state is  $s_{i+1} = (p_{i+1}, v_{i+1}, o_{i+1})$  where:

$$f_{i+1} = f_i + (v_i \cdot \lambda \cdot (\cos o_i + e_i, \sin o_i + e_i)) \quad (3)$$

$$r_{i+1} = r_i + (v_i \cdot \lambda \cdot (\cos o_i, \sin o_i)) \quad (4)$$

$$p_{i+1} = \frac{f_{i+1} + r_{i+1}}{2} \quad (5)$$

$$v_{i+1} = \max \{v_{\min}, \min \{v_{\max}, v_i + (\lambda \cdot t_i)\}\} \quad (6)$$

$$o_{i+1} = \text{atan2}(f_{i+1}^y - r_{i+1}^y, f_{i+1}^x - r_{i+1}^x) \quad (7)$$

The body kinematics from Equations 1–7 are illustrated in Figure 4.

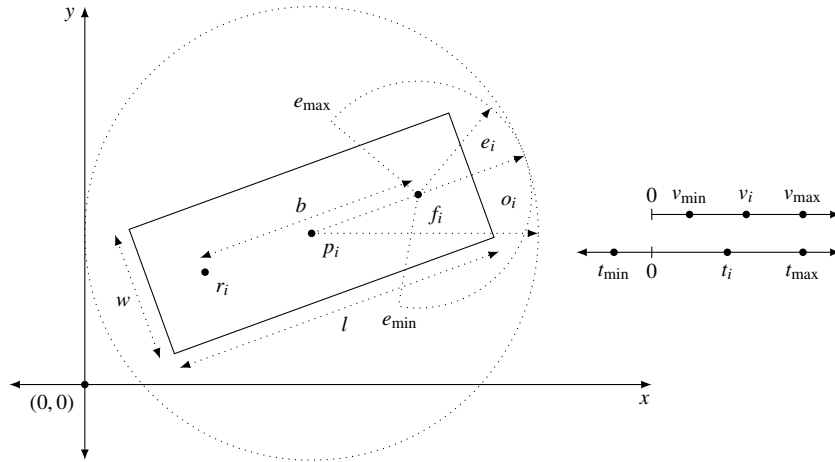


Figure 4: Body kinematics

**Problem:** Given state  $s_i$  and target orientation  $o_g$  such that  $o_g$  is reachable from  $s_i$  with a single action, find steering angle  $e_i$  such that executing action  $a_i = (t_i, e_i)$  in  $s_i$  results in successor state  $s_{i+1} = (p_{i+1}, v_{i+1}, o_{i+1})$  where  $o_{i+1} = o_g$  for any throttle  $t_i$ .

## Appendix: Algorithms

**To do:** Account for discretisation of continuous action spaces at the agent level, although strictly speaking only approximate Q-learning agents require finite action spaces.

---

### Algorithm 1: SIMULATOR

---

**persistent:** environment  $e$ , agents  $N = \{1, \dots, n\}$ , terminator  $\psi \subseteq S$

```
1 begin
2   for each episode do
3      $s \leftarrow \text{reset } e$ 
4     for each timestep do
5        $a \leftarrow (\text{CHOOSEACTION}_1(s), \dots, \text{CHOOSEACTION}_n(s))$ 
6        $r, s' \leftarrow \text{execute } a \text{ in } e$ 
7       for each agent  $i \in N$  do
8          $\text{PROCESSFEEDBACK}_i(s, a_i, r_i, s')$ 
9       if  $s' \in \psi$  then break else  $s \leftarrow s'$ 
```

---

---

### Algorithm 2: RANDOMAGENT

---

**persistent:** exploration rate  $\epsilon \in [0, 1]$

```
1 function  $\text{CHOOSEACTION}(s)$ 
2   with probability  $\epsilon$  do
3     return random choice from  $A(s)$ 
4   return  $\emptyset$ 
```

---

---

### Algorithm 3: PROGRAMMEDRANDOMAGENT

---

**persistent:** programmed behaviour  $\pi : S \rightarrow A$ , terminator  $\psi \subseteq S$ , exploration rate  $\epsilon \in [0, 1]$

```
1 function  $\text{CHOOSEACTION}(s)$ 
2   if  $\pi$  is active then
3     if  $s \in \psi$  then set  $\pi$  as inactive else return  $\pi(s)$ 
4   with probability  $\epsilon$  do
5     set  $\pi$  as active
6     return  $\pi(s)$ 
7   return  $\emptyset$ 
```

---

---

**Algorithm 4: PROGRAMMEDREACTIVEAGENT**

---

**persistent:** programmed behaviour  $\pi : S \rightarrow A$ , trigger  $\varphi \subseteq S$ , terminator  $\psi \subseteq S$

```
1 function CHOOSEACTION( $s$ )
2   if  $\pi$  is active then
3     if  $s \in \psi$  then set  $\pi$  as inactive else return  $\pi(s)$ 
4   if  $s \in \varphi$  then
5     set  $\pi$  as active
6     return  $\pi(s)$ 
7   return  $\emptyset$ 
```

---

---

**Algorithm 5: PROGRAMMELECTIONAGENT**

---

**persistent:** programmed behaviour  $\pi : S \rightarrow A$ , terminator  $\psi \subseteq S$ , coordinator  $c$

```
1 function CHOOSEACTION( $s$ )
2   if  $\pi$  is active then
3     if  $s \in \psi$  then set  $\pi$  as inactive else return  $\pi(s)$ 
4   if elected by  $c$  then
5     set  $\pi$  as active
6     return  $\pi(s)$ 
7   return  $\emptyset$ 
```

---

---

**Algorithm 6: QLEARNINGAGENT**

---

**persistent:** learning rate  $\alpha \in [0, 1]$ , discount factor  $\gamma \in [0, 1]$ , exploration rate  $\epsilon \in [0, 1]$ ,  
feature  $f_j : S \times A \rightarrow \mathbb{R}$  with weight  $w_j \in \mathbb{R}$  for  $j = 1, \dots, m$

```
1 function CHOOSEACTION( $s$ )
2   with probability  $\epsilon$  do
3     return random choice from  $A(s)$ 
4   return random choice from  $\text{argmax}_{a \in A(s)} \text{QVALUE}(s, a)$ 
5 procedure PROCESSFEEDBACK( $s, a, r, s'$ )
6    $q \leftarrow (r + \gamma \cdot \max_{a' \in A(s')} \text{QVALUE}(s', a')) - \text{QVALUE}(s, a)$ 
7   for each feature  $f_j$  do
8      $w_j \leftarrow w_j + \alpha \cdot q \cdot f_j(s, a)$ 
9 function QVALUE( $s, a$ )
10  return  $\sum_{j=1}^m f_j(s, a) \cdot w_j$ 
```

---