



SAPIENZA  
UNIVERSITÀ DI ROMA

# Reinforcement Learning for $LTL_f/LDL_f$ Goals: Theory and Implementation

Facoltà di Ingegneria dell'Informazione, Informatica e Statistica  
Corso di Laurea Magistrale in Master in Engineering in Computer Science

Candidate

Marco Favorito

ID number 1609890

Thesis Advisor

Prof. Giuseppe De Giacomo

Co-Advisor

Prof. Luca Iocchi

Academic Year 2017/2018

Thesis defended on 20<sup>th</sup> July 2018  
in front of a Board of Examiners composed by:

Prof. Riccardo Rosati (chairman)

Prof. Silvia Bonomi

Prof. Giorgio Grisetti

Prof. Massimo Mecella

Prof. Daniele Cono D'Elia

---

**Reinforcement Learning for  $LTL_f/LDL_f$  Goals: Theory and Implementation**

Master thesis. Sapienza – University of Rome

© 2018 Marco Favorito. All rights reserved

This thesis has been typeset by  $\text{\LaTeX}$  and the Sapthesis class.

Author's email: [favorito.1609890@studenti.uniroma1.it](mailto:favorito.1609890@studenti.uniroma1.it)

## Abstract

MDPs extended with  $LTL_f/LDL_f$  non-Markovian rewards have recently attracted interest as a way to specify rewards declaratively. In this thesis, we discuss how a reinforcement learning agent can learn policies fulfilling  $LTL_f/LDL_f$  goals. In particular we focus on the case where we have two separate representations of the world: one for the agent, using the (predefined, possibly low-level) features available to it, and one for the goal, expressed in terms of high-level (human-understandable) fluents. We formally define the problem and show how it can be solved. Moreover, we provide experimental evidence that keeping the RL agent feature space separated from the goal's can work in practice, showing interesting cases where the agent can indeed learn a policy that fulfills the  $LTL_f/LDL_f$  goal using only its features (augmented with additional memory).

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Reinforcement Learning . . . . .	1
1.2	Rewarding behaviors . . . . .	2
1.3	Topic of the Thesis . . . . .	2
1.4	Structure of the Thesis . . . . .	3
<b>2</b>	<b><math>LTL_f</math> and <math>LDL_f</math></b>	<b>4</b>
2.1	Linear time Temporal Logic ( $LTL$ ) . . . . .	4
2.1.1	Syntax . . . . .	4
2.1.2	Semantics . . . . .	5
2.2	Linear Temporal Logic on Finite Traces: $LTL_f$ . . . . .	6
2.2.1	Syntax . . . . .	7
2.2.2	Semantics . . . . .	9
2.2.3	Complexity and Expressiveness . . . . .	9
2.3	Regular Temporal Specifications ( $RE_f$ ) . . . . .	10
2.4	Linear Dynamic Logic on Finite Traces: $LDL_f$ . . . . .	11
2.4.1	Syntax . . . . .	11
2.4.2	Semantics . . . . .	12
2.5	$LTL_f$ and $LDL_f$ translation to automata . . . . .	13
2.5.1	$\partial$ function for $LTL_f$ . . . . .	14
2.5.2	$\partial$ function for $LDL_f$ . . . . .	15
2.5.3	The $LDL_f2NFA$ algorithm . . . . .	15
2.6	On-the-fly DFA . . . . .	20
2.6.1	On-the-fly $LTL_f/LDL_f$ evaluation . . . . .	21
2.6.2	$LDL_f2DFA$ : a variant of $LDL_f2NFA$ . . . . .	25
2.7	Reasoning in $LTL_f/LDL_f$ . . . . .	28
2.8	Conclusions . . . . .	29
<b>3</b>	<b>FLLOAT</b>	<b>31</b>
3.1	Introduction . . . . .	31
3.2	Package structure . . . . .	32
3.3	Propositional Logic formulas . . . . .	32
3.3.1	Specify a PL formula . . . . .	32
3.3.2	Parsing of PL formulas . . . . .	39
3.3.3	PL formulas evaluation . . . . .	41
3.4	$LTL_f$ formulas . . . . .	44
3.4.1	Specify $LTL_f$ formulas . . . . .	44
3.4.2	Parsing of $LTL_f$ formulas . . . . .	47
3.4.3	$LTL_f$ formulas evaluation . . . . .	48
3.5	$LDL_f$ formulas . . . . .	50

3.5.1	Specify $LDL_f$ formulas . . . . .	50
3.5.2	Parsing of $LDL_f$ formulas . . . . .	52
3.5.3	From $LTL_f$ to $LDL_f$ . . . . .	53
3.5.4	$LDL_f$ formulas evaluation . . . . .	54
3.6	$LTL_f/LDL_f$ formula to automaton . . . . .	55
3.6.1	From $LTL_f/LDL_f$ formula to NFA . . . . .	55
3.6.2	On the fly evaluation of $LTL_f/LDL_f$ formulas . . . . .	56
3.6.3	From $LTL_f/LDL_f$ formula to DFA . . . . .	58
3.7	Conclusion . . . . .	59
<b>4</b>	<b>RL for <math>LTL_f/LDL_f</math> Goals</b>	<b>60</b>
4.1	Reinforcement Learning . . . . .	60
4.2	Markov Decision Process (MDP) . . . . .	60
4.3	Temporal Difference Learning . . . . .	64
4.4	Non-Markovian Reward Decision Process (NMRDP) . . . . .	65
4.4.1	Preliminaries . . . . .	65
4.4.2	Find an optimal policy $\bar{\rho}$ for NMRDPs . . . . .	66
4.4.3	Define the non-Markovian reward function $\bar{R}$ . . . . .	67
4.4.4	Using PLTL . . . . .	68
4.5	NMRDP with $LTL_f/LDL_f$ rewards . . . . .	68
4.6	RL for $LTL_f/LDL_f$ Goals . . . . .	70
4.6.1	Problem definition . . . . .	70
4.6.2	Examples . . . . .	72
4.6.3	Reduction to MDP . . . . .	72
4.6.4	An episodic goal-based view . . . . .	75
4.7	Conclusions . . . . .	76
<b>5</b>	<b>Automata-based Reward shaping</b>	<b>77</b>
5.1	Reward Shaping Theory . . . . .	77
5.1.1	Potential-Based Reward Shaping . . . . .	78
5.1.2	Dynamic Potential-Based Reward Shaping . . . . .	78
5.1.3	Relevant considerations about PBRs . . . . .	79
5.2	Off-line Reward shaping over $\mathcal{A}_\varphi$ . . . . .	81
5.3	On-The-Fly Reward shaping . . . . .	82
5.4	Conclusion . . . . .	83
<b>6</b>	<b>RLTG</b>	<b>84</b>
6.1	Introduction . . . . .	84
6.2	Package structure . . . . .	84
6.3	Code examples . . . . .	85
6.3.1	Classic Reinforcement Learning . . . . .	85
6.3.2	Temporal goal Reinforcement Learning . . . . .	87
<b>7</b>	<b>Experiments</b>	<b>89</b>
7.1	BREAKOUT . . . . .	89
7.1.1	Optimal policies . . . . .	91
7.1.2	Benchmarking of reward shaping techniques . . . . .	92
7.2	SAPIENTINO . . . . .	96
7.2.1	Results . . . . .	98
7.3	MINECRAFT . . . . .	99
7.3.1	Results . . . . .	100
7.4	Implementation . . . . .	101

---

<b>8 Conclusion and Future Work</b>	<b>102</b>
8.1 Overview . . . . .	102
8.2 Summary of Main Contributions . . . . .	102
8.3 Future Works . . . . .	103
8.4 Final Remarks . . . . .	103
<b>Bibliography</b>	<b>105</b>

# Chapter 1

## Introduction

This chapter presents the outline of this thesis and summarizes motivations, goals, and achievements.

### 1.1 Reinforcement Learning

Reinforcement learning is an area of Machine Learning where the learning comes from rewards and punishments (Sutton and Barto, 1998). It is concerned in how the learning entity, the *agent*, interacting in an *environment*, should take *actions* so to maximize the observed *reward*. The reward signal is observed after each action taken by the agent. The agent chooses actions depending on the current state of the environment. A solution to the reinforcement learning problem is a *policy* which determines which action should be executed in a given state in order to maximize the long term reward. An algorithm that tackles this kind of problem is called *reinforcement learning algorithm*.

The problem, due to its generality, is studied in many other disciplines, such as *game theory*, *control theory*, *operations research*, *multi-agent system*, and many others. Usually, in order to simplify the tractability of the problem, it is assumed that the environment can be modeled as a *Markov Decision Process* (MDP). An environment behaves like an MDP if the Markov property is satisfied, which means that the state space representation in the algorithm captures enough details so that the optimal decisions can be made when the information about only the current state is available.

Even if the laws that determines the evolution of the systems and the rewards are unknown a priori, it is still possible to solve an MDP by making several simulations and gathering experiences about the visited states, the actions taken and the observed rewards. However, many challenges arise in this settings:

- *Exponential state space explosion*: due to the feature selection used for the state space encoding, every added feature yields exponential increase in the number of states, hence reducing the tractability of the problem.
- *Exploration-Exploitation trade-off*: due to the former issue, reinforcement learning algorithm should be designed to avoid exploring irrelevant states in terms of expected reward, while preferring the ones with high expected reward (*exploitation*). Furthermore, the algorithms should be sensitive to the local optima issue, a well-known in statistical learning literature (*exploration*).

- *Temporal credit assignment*: the agent should be able to foresee the effect of his actions (in terms of expected reward), due to the fact that, in many domains, the current reward is influenced by past decisions.

These problems lead to the use of heuristics and approximate solutions. A simple way to do reinforcement learning is to use exploration which is based on the current policy with a certain degree of randomness which deviates from such a policy.

## 1.2 Rewarding behaviors

In some domains it could be of interest the study of rewards not depending on a single decision (like in MDPs) but depending on a *sequence* of visited states and actions. For instance, we can reward an agent not only by reaching a goal state, but if the goal is reached while satisfying other properties of interest during the simulation or, in other words, if the agent satisfies some target behavior. It is clear that the definition of MDP does not fit this problem, since the optimal decision in the current state depends from the *history of states* that leads to the current state.

This idea of rewarding behaviors has been proposed in (Bacchus et al., 1996), by defining the *Non-Markovian Reward Decision Process* (NMRDP), a variant of an MDP where the reward does not depend only from one transition of the environment but from a sequence of transitions. In order to specify the desired behaviors that the agent should learn, they defined a temporal logic formalisms called PLTL (Past LTL), which is able to speak about a sequence of property configurations over time, that we call *traces*. The classic reinforcement learning algorithms does not work on an NMRDP; however, they propose a transformation from NMRDP into an *expanded* MDP such that the solution of the MDP is also a solution of the original problem. Hence, in order to solve the NMRDP, we can run off-the-shelf RL algorithms over the transformed MDP; the learnt Markovian policy can be easily converted into an optimal policy for the NMRDP.

The trick here is that the transformed MDP is defined over an *expanded state space*, which still contains the original state space but enriches it by labeling every state. The idea is that the labels should keep track in some way the (partial) satisfaction of the temporal formulas. As a result, every state in the transformed state space is replicated multiple times, marking the difference between different (relevant) histories terminating in that state.

A similar transformation has been done in the following works: (Thiébaux et al., 2006) and (Gretton, 2014), where the temporal logic formalism was respectively  $\$FLTL$  (a finite LTL with future formulas) and  $\$*FLTL$  (a variant of  $\$FLTL$ ); in (Icarte et al., 2018), where they used *Co-Safe* LTL formulas (Kupferman and Y. Vardi, 2001; Lacerda et al., 2015); in (Camacho et al., 2017a,b), where they used  $LTL_f$  (Linear Temporal Logic over finite traces) (De Giacomo and Vardi, 2013). In (Brafman et al., 2018) the specification of temporal goals is done by  $LTL_f$  or  $LDL_f$  (Linear Dynamic Logic over finite traces) formulas.

## 1.3 Topic of the Thesis

In this thesis, we leverage the construction of (Brafman et al., 2018) to define a new problem. Consider a classic reinforcement learning problem defined by an MDP. The agent acts in the state space of the MDP and tries to maximize the reward. We can



think of the features of this state space as *low-level features*, e.g. the coordinates of a robot in a room, informations about the limbs etc.

Now consider that we want to talk about some high-level properties of the environment, e.g. continuing with the example of the robot in the room, we are interested in the status of the window and the door (if they are open or closed). We can define the fluents *open\_window* and *open\_door*, respectively. We call the features to determine the status of the fluents *high-level features*.

Given this setting, we are interested in:

- maximize the reward of the original MDP;
- maximize the reward specified by temporal formulas, expressed in  $LTL_f/LDL_f$ , over the fluents

E.g. in the example, a temporal specification could be "eventually open the window, then eventually open the door".

In this work we specify formally the just defined problem, propose the transformation into an MDP and prove the equivalence of the transformation with the original problem. We propose also a way to apply reward shaping techniques to improve the learning in terms of convergence rate. Moreover, we provide an implementation and give experimental evidence of the goodness of our construction.

## 1.4 Structure of the Thesis

The rest of the thesis is structured as follows:

- In Chapter 2 we describe the notions about temporal logic formalisms, that will be used for temporal goal specifications. We start from LTL, RE and then we move towards  $LTL_f$  and  $LDL_f$ , upon which our method is built on.
- In Chapter 3 we describe FLLOAT, a software project that implements the translation from  $LTL_f/LDL_f$  formulas to equivalent automata. Such translation is an important piece of our approach;
- Chapter 4 is the core of the thesis. We introduce foundational concepts in reinforcement learning, MDPs and algorithm to find a solution. Then we move to NMRDPs and we formally describe our approach.
- In Chapter 5 we apply reward shaping techniques to the setting explained in Chapter 4.
- In Chapter 6 we present a reinforcement learning framework allowing easy implementation of the construction described in Chapter 4 and 5.
- Chapter 7 describes the conducted experiments, giving evidence that our approach actually works.
- The thesis is concluded in Chapter 8. This chapter summarizes also the achievements of the thesis and discusses future work.

## Chapter 2

# LTL<sub>f</sub> and LDL<sub>f</sub>

In this chapter we introduce the reader to the main important framework for talk about behaviors over time, which gives the foundations for our approach. First we talk about the well known Linear time Temporal Logic (LTL), Propositional Dynamic Logic (PDL) and their main applications; then we go more in deep by presenting a specific formalism, namely *Linear Temporal Logic over Finite Traces* LTL<sub>f</sub> and *Linear Dynamic Logic over Finite Traces* LDL<sub>f</sub>. Finally, we study the translation from LTL<sub>f</sub>/LDL<sub>f</sub> formulas to Deterministic Finite Automata (DFA). We require the reader to be acquainted with classical logic (Shapiro and Kouri Kissel, 2018) and automata theory (Hopcroft et al., 2000).

### 2.1 Linear time Temporal Logic (LTL)

*Temporal Logic* (Goranko and Galton, 2015) is a category of formal languages aimed to talk about properties of a system whose truth value might change over time. This is in contrast with atemporal logics, which can only discuss about statements whose truth value is constant.

*Linear time Temporal Logic* (Pnueli, 1977), or *Linear Temporal Logic* (LTL) is such a logic. It is the most popular and widely used temporal logic in computer science, especially in formal verification of software/hardware systems, in AI to reasoning about actions and planning, and in the area of Business Process Specification and Verification to specify processes declaratively.

It allows to express temporal patterns about some property  $p$ , like *liveness* ( $p$  will eventually happen), *safety* ( $p$  will never happen) and *fairness*, combinations of the previous patterns (*infinitely often  $p$  holds*, *eventually always  $p$  holds*).

#### 2.1.1 Syntax

A LTL formula  $\varphi$  is defined over a set of propositional symbols  $\mathcal{P}$  and are closed under the boolean connectives, the unary temporal operator  $\mathcal{O}$  (*next-time*) and the binary operator  $\mathcal{U}$  (*until*):

$$\varphi ::= A \mid \neg\varphi \mid \varphi_1 \wedge \varphi_2 \mid \mathcal{O}\varphi \mid \varphi_1 \mathcal{U} \varphi_2$$

With  $A \in \mathcal{P}$ .

Additional operators can be defined in terms of the ones above: as usual logical operators such as  $\vee$ ,  $\Rightarrow$ ,  $\Leftrightarrow$ , *true*, *false* and temporal formulas like *eventually* as  $\Diamond\varphi \doteq \text{true} \mathcal{U} \varphi$ , *always* as  $\Box\varphi \doteq \neg\Diamond\neg\varphi$  and *release* as  $\varphi_1 \mathcal{R} \varphi_2 \doteq \neg(\neg\varphi_1 \mathcal{U} \neg\varphi_2)$ .

**Example 2.1.** Several interesting temporal properties can be defined in LTL:

- *Liveness*:  $\Diamond\varphi$ , which means "condition expressed by  $\varphi$  at some time in the future will be satisfied", "sooner or later  $\varphi$  will hold" or "eventually  $\varphi$  will hold". E.g.,  $\Diamond rich$  (eventually I will become rich),  $Request \implies \Diamond Response$  (if someone requested the service, sooner or later he will receive a response).
- *Safety*:  $\Box\varphi$ , which means "condition expressed by  $\varphi$ , every time in the future will be satisfied", "always  $\varphi$  will hold". E.g.,  $\Box happy$  (I'm always happy),  $\Box \neg (temperature > 30)$  (the temperature of the room must never be over 30).
- *Response*:  $\Box\Diamond\varphi$  which means "at any instant of time there exists a moment later where  $\varphi$  holds". This temporal pattern is known in computer science as *fairness*.
- *Persistence*:  $\Diamond\Box\varphi$ , which stand for "There exists a moment in the future such that from then on  $\varphi$  always holds". E.g.  $\Diamond\Box dead$  (at a certain point you will die, and you will be dead forever)
- *Strong fairness*:  $\Box\Diamond\varphi_1 \implies \Box\Diamond\varphi_2$ , "if something is attempted/requested infinitely often, then it will be successful/allocated infinitely often". E.g.,  $\Box\Diamond ready \implies \Box\Diamond run$  (if a process is in ready state infinitely often, then infinitely often it will be selected by the scheduler).

### 2.1.2 Semantics

The semantics of LTL is provided by (infinite) *traces*, i.e.  $\omega$ -word over the alphabet  $2^{\mathcal{P}}$ .

**Definition 2.1.** Given a infinite trace  $\pi$ , we define that a LTL formula  $\varphi$  is *true* at time  $i$ , in symbols  $\pi, i \models \varphi$  inductively as follows:

$$\begin{aligned}
 \pi, i &\models A, \text{ for } A \in \mathcal{P} \text{ iff } A \in \pi(i) \\
 \pi, i &\models \neg\varphi \text{ iff } \pi, i \not\models \varphi \\
 \pi, i &\models \varphi_1 \wedge \varphi_2 \text{ iff } \pi, i \models \varphi_1 \wedge \pi, i \models \varphi_2 \\
 \pi, i &\models \bigcirc\varphi \text{ iff } \pi, i+1 \models \varphi \\
 \pi, i &\models \varphi_1 \mathcal{U} \varphi_2 \text{ iff } \exists j. (j \geq i) \wedge \pi, j \models \varphi \wedge \forall k. (i \leq k < j) \Rightarrow \pi, k \models \varphi_1
 \end{aligned}$$

Similarly as in classical logic we give the following definitions:

**Definition 2.2.** A LTL formula is *true* in  $\pi$ , in notation  $\pi \models \varphi$ , if  $\pi, 0 \models \varphi$ . A formula  $\varphi$  is *satisfiable* if it is true in some  $\pi$  and is *valid* if it is true in every  $\pi$ .  $\varphi_1$  *entails*  $\varphi_2$ , in symbols  $\varphi_1 \models \varphi_2$  iff  $\forall \pi, \forall i. \pi, i \models \varphi_1 \implies \pi, i \models \varphi_2$ .

Now we state an important result:

**Theorem 2.1** (Sistla and Clarke (1985)). *Satisfiability, validity, and entailment for LTL formulas are PSPACE-complete.*

Indeed, Linear Temporal Logic can be thought of as a specific decidable (PSPACE-complete) fragment of classical first-order logic (FOL).

Notice that, a *trace*  $\pi$  can be seen as a *word* on a *path* of a *Kripke structure*.

**Definition 2.3** (Clarke et al. (1999)). a Kripke structure  $\mathcal{K}$  over a set of propositional symbols  $\mathcal{P}$  is a 4-tuple  $\langle S, I, R, L \rangle$  where  $S$  is a finite set of *states*,  $I \subseteq S$  is the set of *initial states*,  $R \subseteq S \times S$  is the *transition relation* such that  $R$  is left-total and  $L : S \rightarrow 2^{\mathcal{P}}$  is a *labeling function*.

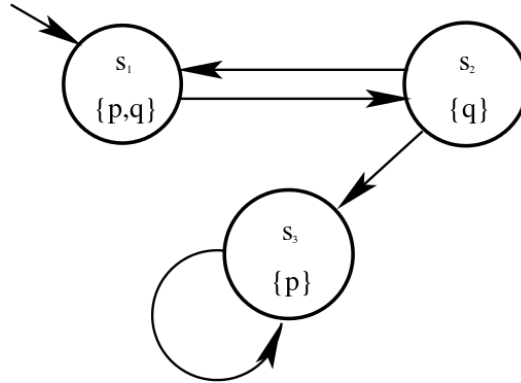
A *path*  $\rho$  over  $\mathcal{K}$  is a sequence of states  $\langle s_1, s_2, \dots \rangle$  such that  $\forall i. R(s_i, s_{i+1})$ . From a path we can build a *word*  $w$  on the path  $\rho$  by mapping each state of the sequence with  $L$ , namely:

$$w = \langle L(s_1), L(s_2), \dots \rangle$$

In simpler words, a trace of propositional symbols  $\mathcal{P}$  is a infinite sequence of combinations of propositional symbols in  $\mathcal{P}$ . Moreover, we denote by  $\pi(i)$  with  $i \in \mathbb{N}$  the labels associated to  $s_i$ , i.e.  $L(s_i)$ .

**Example 2.2.** In figure 2.1 is depicted an example of Kripke structure  $\mathcal{K}$  over  $\mathcal{P} = \{p, q\}$  where:

$$\begin{aligned} S &= \{s_1, s_2, s_3\} \\ I &= \{s_1\} \\ R &= \{(s_1, s_2), (s_2, s_1), (s_2, s_3), (s_3, s_3)\} \\ L &= \{(s_1, \{p, q\}), (s_2, \{q\}), (s_3, \{p\})\} \end{aligned}$$



**Figure 2.1.** An example of Kripke structure.

The path  $\langle s_1, s_2, s_3, s_3, s_3 \dots \rangle$  yields the following trace  $\pi$ :

$$\begin{aligned} \pi &= \langle L(s_1), L(s_2), L(s_3), L(s_3), L(s_3), \dots \rangle \\ &= \langle \{p, q\}, \{q\}, \{p\}, \{p\}, \{p\}, \dots \rangle \end{aligned}$$

## 2.2 Linear Temporal Logic on Finite Traces: $LTL_f$

Linear-time Temporal Logic over finite traces,  $LTL_f$ , is essentially standard LTL (Pnueli, 1977) interpreted over finite, instead of over infinite, traces (De Giacomo and Vardi, 2013). This apparently trivial difference has big impact: as we will see, some LTL formula has a different meaning if interpreted over infinite traces or finite ones.

### 2.2.1 Syntax

In fact, the syntax of  $LTL_f$  is the same of the one showed in Section ??, i.e. *formulas* of  $LTL_f$  are built from a set  $\mathcal{P}$  of propositional symbols and are closed under the boolean connectives, the unary temporal operator  $O(\text{next-time})$  and the binary operator  $\mathcal{U}$  (*until*):

$$\varphi ::= \phi \mid \neg\varphi \mid \varphi_1 \wedge \varphi_2 \mid O\varphi \mid \varphi_1 \mathcal{U} \varphi_2$$

With  $A \in \mathcal{P}$ .

We use the standard abbreviations for classical logic formulas:

$$\begin{aligned} \varphi_1 \vee \varphi_2 &\doteq \neg(\neg\varphi_1 \wedge \neg\varphi_2) \\ \varphi_1 \Rightarrow \varphi_2 &\doteq \neg\varphi_1 \vee \varphi_2 \\ \varphi_1 \Leftrightarrow \varphi_2 &\doteq \varphi_1 \Rightarrow \varphi_2 \wedge \varphi_2 \Rightarrow \varphi_1 \\ \text{true} &\doteq \neg\varphi \vee \varphi \\ \text{false} &\doteq \neg\varphi \wedge \varphi \end{aligned}$$

And for temporal formulas:

$$\varphi_1 \mathcal{R} \varphi_2 \doteq \neg(\neg\varphi_1 \mathcal{U} \neg\varphi_2) \tag{2.1}$$

$$\Diamond\varphi \doteq \text{true} \mathcal{U} \varphi \tag{2.2}$$

$$\Box\varphi \doteq \neg\Diamond\neg\varphi \tag{2.3}$$

$$\bullet\varphi \doteq \neg O\neg\varphi \tag{2.4}$$

$$\text{Last} \doteq \bullet\text{false} \tag{2.5}$$

$$\text{End} \doteq \Box\text{false} \tag{2.6}$$

As the reader might already noticed, 2.2 and 2.3 are defined as in Section 2.1.1; Equation 2.1 is called *release*; Equation 2.4 is called *weak next* (notice that on finite traces  $\neg O\varphi \not\equiv O\neg\varphi$ ); 2.5 denotes the end of the trace, while 2.6 denotes that the trace is ended.

**Example 2.3.** Here we recall Example 2.1 and we see the impact on *Always*, *Eventually Response* and *Persistence* LTL formulas if interpreted on finite traces (i.e. formulas in  $LTL_f$ ):

- *Safety*:  $\Box A$  means that always *till the end of the trace*  $\varphi$  holds;
- *Liveness*:  $\Diamond A$  means that eventually *before the end of the trace*  $\varphi$  holds;
- *Response*:  $\Box\Diamond\varphi$  on finite traces becomes equivalent to *last point in the trace satisfies  $\varphi$* , i.e.  $\Diamond(\text{Last} \wedge \varphi)$ . Intuitively, this is true because  $\Box\Diamond\varphi$  implies that at the last point in the trace  $\varphi$  holds (because there are no successive instants of time that make  $\varphi$  true); but if this is the case, then what happens at previous points in the trace does not matter because the formula evaluates always to true, since as we just said  $\varphi$  must hold at the last point in the trace, hence the equivalence with  $\Diamond(\text{Last} \wedge \varphi)$ .

- *Persistence*:  $\Diamond\Box\varphi$  on finite traces becomes equivalent to *last point in the trace satisfies  $\varphi$* , i.e.  $\Diamond(Last \wedge \varphi)$ . Analogously to the previous case, the equivalence holds because  $\Diamond\Box\varphi$  implies that at the last point in the trace  $\Box\varphi$  holds (and so  $\varphi$ ), since we have no further successive instants of time that makes  $\Box\varphi$  true. But if this is the case, then what happens at previous points in the trace does not matter because the formula evaluates always to true, since as we just said  $\Box\varphi$  (and so  $\varphi$ ) must hold at the last point in the trace, hence the equivalence with  $\Diamond(Last \wedge \varphi)$ .

In other words, no direct nesting of eventually and always connectives is meaningful in  $LTL_f$ , and this contrast what happens in LTL of infinite traces.

**Example 2.4.** Another remarkable evidence about the relevance of the assumption about finiteness of traces is provided by the DECLARE approach (Pesic and van der Aalst, 2006).

DECLARE is a declarative approach to business process modeling based on LTL interpreted over finite traces. The intuition is to map finite traces describing a domain of interest (e.g. processes) into infinite traces under the assumption that

$$\Diamond end \wedge \Box(end \Rightarrow \bigcirc end) \wedge \Box(end \Rightarrow \bigwedge_{p \in \mathcal{P}} \neg p) \quad (2.7)$$

which means that the following english statements hold:

- *end* eventually holds ( $end \notin \mathcal{P}$ );
- once *end* is true, it is true forever;
- when *end* is true all other propositions must be false

In other words, every finite trace  $\pi_f$  is extended with an infinite sequence of *end*, or in symbols  $\pi_{inf} = \pi_f \{end\}^\omega$ . By construction we have that

$$\pi_{inf} \models \Diamond end \wedge \Box(end \Rightarrow \bigcirc end) \wedge \Box(end \Rightarrow \bigwedge_{p \in \mathcal{P}} \neg p)$$

Despite it seems a nice construction to adapt LTL on finite traces, in fact it is wrong due to the *next* operator: in an infinite trace a successor state always exists, whereas in a finite one this does not hold. There exists a counterexample showing that the interpretation of LTL formulas on finite traces with the construction just explained is **not** equivalent with proper interpretation over finite traces offered by  $LTL_f$ , i.e. in general:

$$\pi_f \{end\}^\omega \models \varphi \not\equiv_f \pi_f \models_f \varphi \quad (2.8)$$

To see why this is the case, consider the DECLARE "negation chain succession"  $\Box(a \Rightarrow \bigcirc \neg b)$  which requires that at any point in the trace, the state after we see *a*, *b* is false. Consider also the finite trace  $\pi_f = \{a\}$  and the associated infinite trace  $\pi_{inf} = \{a\} \{end\}^\omega$  built as explained before. We have that

$$\pi_{inf} \models \Box(a \Rightarrow \bigcirc \neg b)$$

where  $\models$  has been defined in 2.1. This is true because there is only one occurrence of *a* and then *end* holds forever (and so *b* does not).

But if the same formula is interpreted on finite traces (namely  $\models_f$ ):

$$\pi_f \not\models_f \Box(a \Rightarrow \bigcirc \neg b)$$

because the finite trace  $a$  is true at the last instant, but then there is no next instance where  $b$  is false, so  $\bigcirc \neg b$  is evaluated to *false* and the formula does not hold. The correct way to express "negation chain succession" on finite traces would be  $\Box(a \Rightarrow \bullet \neg b)$ .

The LTL formulas  $\varphi$  that are insensitive to the problem just shown, i.e. such that

$$\pi_f\{end\}^\omega \models \varphi \text{ iff } \pi_f \models_f \varphi \quad (2.9)$$

holds are defined *insensitive to infiniteness* (De Giacomo et al., 2014). This is another important evidence about the relevance of the finiteness trace assumption.

### 2.2.2 Semantics

Formally, a *finite trace*  $\pi$  is a finite word over the alphabet  $2^{\mathcal{P}}$ , i.e. as alphabet we have all the possible propositional interpretations of the propositional symbols in  $\mathcal{P}$ . We can see  $\pi$  as a *finite* word on a path of a Kripke structure, similarly as we discussed in Section 2.1.2 (but in that case the traces were *infinite*). Given a finite path  $\rho = \langle s_1, s_2, \dots, s_n \rangle$  on a Kripke structure  $\mathcal{K}$ , a finite trace  $\pi$  associated to the path  $\rho$  is defined as  $\langle L(s_1), L(s_2), \dots, L(s_n) \rangle$ .

We use the following notation. We denote the *length* of a trace  $\pi$  as  $length(\pi)$ . We denote the  *$i_{th}$  position* on the trace as  $\pi(i) = L(s_i)$ , i.e. the propositions that hold in the  $i_{th}$  state of the path, with  $0 \leq i \leq last$  where  $last = length(\pi) - 1$  is the last element of the trace. We denote by  $\pi(i, j)$ , the *segment* of  $\pi$ , the trace  $\pi' = \langle \pi(i), \pi(i+1), \dots, \pi(j) \rangle$ , with  $0 \leq i \leq j \leq last$

**Definition 2.4.** Given a finite trace  $\pi$ , we define that a  $LTL_f$  formula  $\varphi$  is *true* at time  $i$  ( $0 \leq i \leq last$ ), in symbols  $\pi, i \models \varphi$  inductively as follows:

$$\begin{aligned} \pi, i &\models A, \text{ for } A \in \mathcal{P} \text{ iff } A \in \pi(i) \\ \pi, i &\models \neg \varphi \text{ iff } \pi, i \not\models \varphi \\ \pi, i &\models \varphi_1 \wedge \varphi_2 \text{ iff } \pi, i \models \varphi_1 \wedge \pi, i \models \varphi_2 \\ \pi, i &\models \bigcirc \varphi \text{ iff } i < last \wedge \pi, i+1 \models \varphi \end{aligned} \quad (2.10)$$

$$\begin{aligned} \pi, i &\models \varphi_1 \mathcal{U} \varphi_2 \text{ iff } \exists j. (i \leq j \leq last) \wedge \pi, j \models \varphi \wedge \\ &\forall k. (i \leq k < j) \Rightarrow \pi, k \models \varphi_1 \end{aligned} \quad (2.11)$$

Notice that Definition 2.4 is pretty similar to Definition 2.1, except the bounding of indexes in Equation 2.10 and Equation 2.11, to recognize that the trace is ended.

Analogously to Definition 2.2 we give the following definitions:

**Definition 2.5.** A  $LTL_f$  formula is *true* in  $\pi$ , in notation  $\pi \models \varphi$ , if  $\pi, 0 \models \varphi$ . A formula  $\varphi$  is *satisfiable* if it is true in some  $\pi$  and is *valid* if it is true in every  $\pi$ .  $\varphi_1$  *entails*  $\varphi_2$ , in symbols  $\varphi_1 \models \varphi_2$  iff  $\forall \pi, \forall i. \pi, i \models \varphi_1 \Rightarrow \pi, i \models \varphi_2$ .

### 2.2.3 Complexity and Expressiveness

Thanks to reduction of  $LTL_f$  satisfiability (Definition 2.5) into LTL satisfiability for PSPACE membership and reduction of STRIPS planning into  $LTL_f$  satisfiability for PSPACE-hardness, as proposed in (De Giacomo and Vardi, 2013), we have this result:

**Theorem 2.2** (De Giacomo and Vardi (2013)). *Satisfiability, validity and entailment for LTL<sub>f</sub> formulas are PSPACE-complete.*

About expressiveness of LTL<sub>f</sub>, we have that:

**Theorem 2.3** (De Giacomo and Vardi (2013); Gabbay et al. (1997)). *LTL<sub>f</sub> has exactly the same expressive power of FOL over finite ordered sequences.*

## 2.3 Regular Temporal Specifications (RE<sub>f</sub>)

In this section we talk about regular languages as a form of temporal specification over finite traces. In particular we focus on regular expressions (Hopcroft et al., 2000).

A regular expression  $\varrho$  is defined inductively as follows, considering as alphabet the set of propositional interpretations  $2^{\mathcal{P}}$ , from a set of propositional symbols  $\mathcal{P}$ :

$$\varrho ::= \phi \mid \varrho_1 + \varrho_2 \mid \varrho_1; \varrho_2 \mid \varrho^*$$

where  $\phi$  is a propositional formula that is an abbreviation for the union of all the propositional interpretations that satisfy  $\phi$ , i.e.  $\phi = \sum_{\Pi \models \phi} \Pi$  and  $\Pi \in 2^{\mathcal{P}}$ .

We denote by  $\mathcal{L}(\varrho)$  the language recognized by a RE<sub>f</sub> expression. We interpret these expression over finite traces, introduced in Section 2.2.2.

**Definition 2.6.** We say that a regular expression  $\varrho$  is *true* in the finite trace  $\pi$  if  $\pi \in \mathcal{L}(\varrho)$ . We say that  $\varrho$  is *true at instant  $i$*  if  $\pi(i, \text{last}) \in \mathcal{L}(\varrho)$ . We say that  $\varrho$  is *true between instants  $i, j$*  if  $\pi(i, j) \in \mathcal{L}(\varrho)$ .

**Example 2.5.** We recall Example 2.2. The trace resulting from path  $\langle s_1, s_2, s_3, s_3, \dots \rangle$ , i.e.:

$$\pi = \langle \{p, q\} \{q\}, \{p\}, \{p\}, \{p\}, \dots \rangle$$

belongs to the language generated by the following regular expression:

$$\varrho_1 = p \wedge q; q; p^*$$

But also by this one:

$$\varrho_2 = \text{true}; q + p; \text{true}^*$$

**Example 2.6.** We can express some of the formulas shown in Example 2.3, and many others, in RE<sub>f</sub>:

- *Safety*:  $\varphi^*$ , equivalent to  $\Box\varphi$
- *Liveness*:  $\text{true}^*; \varphi; \text{true}^*$ , equivalent to  $\Diamond\varphi$ ;
- *Response* and *Persistence*: as said before, when interpreted on finite traces, they are equivalent to  $\Diamond(\text{Last} \wedge \varphi)$ ; hence, they can be rewritten in RE<sub>f</sub> as  $\text{true}^*; \varphi$
- *Ordered occurrence*:  $\text{true}^*; \varphi_1; \text{true}^*; \varphi_2; \text{true}^*$ , equivalent to  $\Diamond(\varphi_1 \wedge \Diamond\varphi_2)$  means that  $\varphi_1$  and  $\varphi_2$  happen in order;
- *Alternating sequence*:  $(\psi, \varphi)^*$  means that  $\psi$  and  $\varphi$  alternate from the beginning of the trace, starting with  $\psi$  and ending with  $\varphi$ .



The *Alternating sequence* is an example of formula that has not a counterpart in  $LTL_f$ . More generally,  $LTL_f$  (and  $LTL$ ) are not able to capture regular structural properties on path (Wolper, 1981).

This observation about expressiveness of  $RE_f$  is confirmed by Theorem 6 of (De Giacomo and Vardi, 2013), which is a consequence of several classical results (Büchi, 1960; Elgot, 1961; Trakhtenbrot, 1961; Thomas, 1979):

**Theorem 2.4** (De Giacomo and Vardi (2013)).  *$RE_f$  is strictly more expressive than  $LTL_f$*

More precisely,  $RE_f$  is expressive as *monadic second-order logic* MSO over bounded ordered sequences (Khoussainov and Nerode, 2001).

## 2.4 Linear Dynamic Logic on Finite Traces: $LDL_f$

The problem with  $RE_f$  is that, although is strictly more expressive than  $LTL_f$ , is considered a low-level formalism for temporal specifications. For instance  $RE_f$  misses a direct construct for negation and for conjunction. Moreover, negation requires an exponential blow-up, hence adding complementation and intersection constructs is not advisable.

*Linear Dynamic Logic of Finite Traces*  $LDL_f$  (De Giacomo and Vardi, 2013) merges  $LTL_f$  with  $RE_f$  in a very natural way, borrowing the syntax of PDL (Fischer and Ladner, 1979), a well-known (propositional) logic of programs in computer science. It keep the declarativeness and convenience of  $LTL_f$  while having the same expressive power of  $RE_f$ .

### 2.4.1 Syntax

Formally,  $LDL_f$  formulas  $\varphi$  are built over a set of propositional symbols  $\mathcal{P}$  as follows (Brafman et al., 2017):

$$\begin{aligned}\varphi &::= tt \mid \neg\varphi \mid \varphi_1 \wedge \varphi_2 \mid \langle \varrho \rangle \varphi \\ \varrho &::= \phi \mid \varphi? \mid \varrho_1 + \varrho_2 \mid \varrho_1; \varrho_2 \mid \varrho^*\end{aligned}$$

where  $tt$  stands for logical true;  $\phi$  is a propositional formula over  $\mathcal{P}$ ;  $\varrho$  denotes path expressions, which are RE over propositional formulas  $\phi$  with the addition of the test construct  $\varphi?$  typical of PDL. Moreover, we use the following abbreviations for classical logic operators:

$$\begin{aligned}\varphi_1 \vee \varphi_2 &\doteq \neg(\neg\varphi_1 \wedge \neg\varphi_2) \\ \varphi_1 \Rightarrow \varphi_2 &\doteq \neg\varphi_1 \vee \varphi_2 \\ \varphi_1 \Leftrightarrow \varphi_2 &\doteq \varphi_1 \Rightarrow \varphi_2 \wedge \varphi_2 \Rightarrow \varphi_1 \\ ff &\doteq \neg tt\end{aligned}$$

And for temporal formulas:

$$[\varrho]\varphi \doteq \neg\langle \varrho \rangle \neg\varphi \tag{2.12}$$

$$End \doteq [true]ff \tag{2.13}$$

$$Last \doteq \langle true \rangle End \tag{2.14}$$

$[\varrho]\varphi$  and  $\langle\varrho\rangle\varphi$  are analogous to box and diamond operators in PDL; Formula 2.14 denotes the last element of the trace, whereas Formula 2.13 denotes that the trace is ended. Intuitively,  $\langle\varrho\rangle\varphi$  states that, from the current step in the trace, there exists an execution satisfying the RE  $\varrho$  such that its last step satisfies  $\varphi$ , while  $[\varrho]\varphi$  states that, from the current step, all executions satisfying the RE  $\varrho$  are such that their last step satisfies  $\varphi$ .

### 2.4.2 Semantics

As we did in the previous sections, we formally give a semantics to  $\text{LDL}_f$  (interpreted over finite traces, like  $\text{LTL}_f$  and RE).

**Definition 2.7.** Given a finite trace  $\pi$ , we define that a  $\text{LDL}_f$  formula  $\varphi$  is *true* at time  $i$  ( $0 \leq i \leq \text{last}$ ), in symbols  $\pi, i \models \varphi$  inductively as follows:

$$\begin{aligned}
&\pi, i \models tt \\
&\pi, i \models \neg\varphi \text{ iff } \pi, i \not\models \varphi \\
&\pi, i \models \varphi_1 \wedge \varphi_2 \text{ iff } \pi, i \models \varphi_1 \wedge \pi, i \models \varphi_2 \\
&\pi, i \models \langle\phi\rangle\varphi \text{ iff } i < \text{last} \wedge \pi(i) \models \phi \wedge \pi, i+1 \models \varphi \\
&\pi, i \models \langle\psi?\rangle\varphi \text{ iff } \pi, i \models \psi \wedge \pi, i \models \varphi \\
&\pi, i \models \langle\varrho_1 + \varrho_2\rangle\varphi \text{ iff } \pi, i \models \langle\varrho_1\rangle\varphi \vee \langle\varrho_2\rangle\varphi \\
&\pi, i \models \langle\varrho_1; \varrho_2\rangle\varphi \text{ iff } \pi, i \models \langle\varrho_1\rangle\langle\varrho_2\rangle\varphi \\
&\pi, i \models \langle\varrho^*\rangle\varphi \text{ iff } \pi, i \models \varphi \vee i < \text{last} \wedge \pi, i \models \langle\varrho\rangle\langle\varrho^*\rangle\varphi \text{ and } \varrho \text{ is not test-only}
\end{aligned}$$

We say that  $\varrho$  is *test-only* if it is a  $\text{RE}_f$  expression whose atoms are only tests, i.e.  $\psi?$ .

Notice that  $\text{LDL}_f$  fully captures  $\text{LTL}_f$ . For every formula in  $\text{LTL}_f$  there exists a  $\text{LDL}_f$  formula with the same meaning, namely:

$$\begin{array}{ll}
\text{LTL}_f & \text{LDL}_f \\
A & \langle A \rangle tt \\
\neg\varphi & \neg\varphi \\
\varphi_1 \wedge \varphi_2 & \varphi_1 \wedge \varphi_2 \\
\bigcirc\varphi & \langle \text{true} \rangle (\varphi \wedge \neg \text{End}) \\
\varphi_1 \mathcal{U} \varphi & \langle (\varphi_1?; \text{true})^* \rangle (\varphi_2 \wedge \neg \text{End})
\end{array}$$

Notice also that every  $\text{RE}_f$  expression  $\varrho$  is captured in  $\text{LDL}_f$  by  $\langle\varrho\rangle\text{End}$ . Moreover, since also the converse holds, i.e. every  $\text{LDL}_f$  formula can be expressed in RE (by Theorem 11 in (De Giacomo and Vardi, 2013)), the following theorem holds:

**Theorem 2.5** (De Giacomo and Vardi (2013)).  $\text{LDL}_f$  has exactly the same expressive power of MSO

Now we show several  $\text{LDL}_f$  examples.

**Example 2.7.** Formulas described in Examples 2.3 and 2.6 can be rewritten in  $\text{LDL}_f$  as:

- *Safety*:  $[true^*]\varphi$ , equivalent to LTL<sub>f</sub> formula  $\Box\varphi$
- *Liveness*:  $\langle true^*\rangle\varphi$ , equivalent to LTL<sub>f</sub> formula  $\Diamond\varphi$
- *Conditional Response*:  $[true^*](\varphi_1 \Rightarrow \langle true^*\rangle\varphi_2)$ , equivalent to LTL<sub>f</sub> formula  $\Box(\varphi_1 \Rightarrow \Diamond\varphi_2)$
- *Ordered occurrence*:  $\langle true^*; \varphi_1; true^*; \varphi_2; true^*\rangle End$  equivalent to the RE<sub>f</sub> expression  $true^*; \varphi_1; true^*; \varphi_2; true^*$
- *Alternating occurrence*:  $\langle (\psi; \varphi)^*\rangle End$  equivalent to the RE<sub>f</sub> expression  $(\psi; \varphi)^*$

**Example 2.8.** Consider the Example 2.2 and 2.5.  $\varrho_1$  and  $\varrho_2$  are translated into LDL<sub>f</sub> as  $\langle \varrho_1 \rangle End$  and  $\langle \varrho_2 \rangle End$  respectively.

Other LDL<sub>f</sub> formulas satisfiable in the Kripke structure  $\mathcal{K}$  depicted in Figure 2.1 are:

- $\langle p \rangle tt$  by every (non-empty) path, since  $s_1$  is the initial state and we have that  $\{p, q\} \models p$
- $\langle q \rangle tt$  as the previous case
- $\langle (p; q); (p; q)^*; p; p^* \rangle tt$  by paths of the form  $\rho = s_1, s_2, (s_1, s_2)^\omega, s_3, (s_3)^\omega$
- $[true^*]\langle p \vee q \rangle tt$  is satisfied for every path, since for every reachable state either  $p$  or  $q$  are true;

## 2.5 LTL<sub>f</sub> and LDL<sub>f</sub> translation to automata

Given an LTL<sub>f</sub>/LDL<sub>f</sub> formula  $\varphi$ , we can construct a deterministic finite state automaton (DFA) (Rabin and Scott, 1959)  $\mathcal{A}_\varphi$  that accept the same finite traces that makes  $\varphi$  true. In order to do this, we proceed in two steps: First we translate LTL<sub>f</sub> and LDL<sub>f</sub> formulas into (NFA) (De Giacomo and Vardi, 2015). Then the NFA obtained can be transformed into a DFA following the standard procedure of *determinization*.

Now we recall definitions of NFA and DFA:

**Definition 2.8.** An NFA is a tuple  $\mathcal{A} = \langle \Sigma, Q, q_0, \delta, F \rangle$ , where:

- $\Sigma$  is the input alphabet;
- $Q$  is the finite set of states;
- $q_0 \in Q$  is the initial state;
- $\delta \subseteq Q \times \Sigma \times Q$  is the transition relation;
- $F \subseteq Q$  is the set of final states;

**Definition 2.9.** A DFA is a NFA where  $\delta$  is a function  $\delta : Q \times \Sigma \rightarrow Q$

By  $\mathcal{L}(\mathcal{A})$  we mean the set of all traces over  $\Sigma$  accepted by  $\mathcal{A}$ .

In the next two subsections we give some definition that will be used in the algorithm; then we describe the algorithm for the translation and give some example.

### 2.5.1 $\partial$ function for LTL<sub>f</sub>

We give the following definition:

**Definition 2.10.** The *delta function*  $\partial$  for LTL<sub>f</sub> formulas is a function that takes as input an (implicitly quoted) LTL<sub>f</sub> formula  $\varphi$  in NNF and a propositional interpretation  $\Pi$  for  $\mathcal{P}$ , and returns a positive boolean formula whose atoms are (implicitly quoted)  $\varphi$  subformulas. It is defined as follows:

$$\begin{aligned}
\partial(A, \Pi) &= \begin{cases} true & \text{if } A \in \Pi \\ false & \text{if } A \notin \Pi \end{cases} \\
\partial(\neg A, \Pi) &= \begin{cases} false & \text{if } A \in \Pi \\ true & \text{if } A \notin \Pi \end{cases} \\
\partial(\varphi_1 \wedge \varphi_2, \Pi) &= \partial(\varphi_1, \Pi) \wedge \partial(\varphi_2, \Pi) \\
\partial(\varphi_1 \vee \varphi_2, \Pi) &= \partial(\varphi_1, \Pi) \vee \partial(\varphi_2, \Pi) \\
\partial(\bigcirc \varphi, \Pi) &= \varphi \wedge \neg End \equiv \varphi \wedge \Diamond true \\
\partial(\varphi_1 \mathcal{U} \varphi_2, \Pi) &= \partial(\varphi_2, \Pi) \vee (\partial(\varphi_1, \Pi) \wedge \partial(\bigcirc(\varphi_1 \mathcal{U} \varphi_2), \Pi)) \\
\partial(\bullet \varphi, \Pi) &= \varphi \vee End \equiv \varphi \vee \Box false \\
\partial(\varphi_1 \mathcal{R} \varphi_2, \Pi) &= \partial(\varphi_2, \Pi) \wedge (\partial(\varphi_1, \Pi) \vee \partial(\bullet(\varphi_1 \mathcal{R} \varphi_2), \Pi))
\end{aligned} \tag{2.15}$$

where *End* is defined as Equation 2.6. As a consequence of Definition 2.10 and from Equation 2.2 and 2.3, we can deduce that

$$\begin{aligned}
\partial(\Diamond \varphi, \Pi) &= \partial(\varphi, \Pi) \vee \partial(\bigcirc \Diamond \varphi, \Pi) \\
\partial(\Box \varphi, \Pi) &= \partial(\varphi, \Pi) \wedge \partial(\bullet \Box \varphi, \Pi)
\end{aligned}$$

Moreover, we define  $\partial(\varphi, \epsilon)$  which is inductively defined as Equation 2.15, except for the following cases:

$$\begin{aligned}
\partial(A, \epsilon) &= false \\
\partial(\neg A, \epsilon) &= false \\
\partial(\bigcirc \varphi, \epsilon) &= false \\
\partial(\bullet \varphi, \epsilon) &= true \\
\partial(\varphi_1 \mathcal{U} \varphi_2, \epsilon) &= false \\
\partial(\varphi_1 \mathcal{R} \varphi_2, \epsilon) &= true
\end{aligned} \tag{2.16}$$

Note that  $\partial(\varphi, \epsilon)$  is always either *true* or *false*. It is worth to observe for future use that from Equation 2.16 we can say  $\partial(\Diamond \varphi, \epsilon) = false$  and  $\partial(\Box \varphi, \epsilon) = true$ .

### 2.5.2 $\partial$ function for LDL<sub>f</sub>

We give the following definition:

**Definition 2.11.** The *delta function*  $\partial$  for LDL<sub>f</sub> formulas is a function that takes as input an (implicitly quoted) LDL<sub>f</sub> formula  $\varphi$  in NNF, extended with auxiliary constructs  $\mathbf{F}_\psi$  and  $\mathbf{T}_\psi$ , and a propositional interpretation  $\Pi$  for  $\mathcal{P}$ , and returns a positive boolean formula whose atoms are (implicitly quoted)  $\varphi$  subformulas (not including  $\mathbf{F}_\psi$  or  $\mathbf{T}_\psi$ ). It is defined as follows:

$$\begin{aligned}
\partial(tt, \Pi) &= true \\
\partial(ff, \Pi) &= false \\
\partial(\phi, \Pi) &= \partial(\langle\phi\rangle tt, \Pi) \\
\partial(\varphi_1 \wedge \varphi_2, \Pi) &= \partial(\varphi_1, \Pi) \wedge \partial(\varphi_2, \Pi) \\
\partial(\varphi_1 \vee \varphi_2, \Pi) &= \partial(\varphi_1, \Pi) \vee \partial(\varphi_2, \Pi) \\
\partial(\langle\phi\rangle\varphi, \Pi) &= \begin{cases} \mathbf{E}(\varphi) & \text{if } \Pi \models \phi \\ false & \text{if } \Pi \not\models \phi \end{cases} \\
\partial(\langle\varrho\rangle\varphi, \Pi) &= \partial(\varrho, \Pi) \wedge \partial(\varphi, \Pi) \\
\partial(\langle\varrho_1 + \varrho_2\rangle\varphi, \Pi) &= \partial(\langle\varrho_1\rangle\varphi, \Pi) \vee \partial(\langle\varrho_2\rangle\varphi, \Pi) \\
\partial(\langle\varrho_1; \varrho_2\rangle\varphi, \Pi) &= \partial(\langle\varrho_1\rangle\langle\varrho_2\rangle\varphi, \Pi) \\
\partial(\langle\varrho^*\rangle\varphi, \Pi) &= \partial(\varphi, \Pi) \vee \partial(\langle\varrho\rangle\mathbf{F}_{\langle\varrho^*\rangle\varphi}, \Pi) \\
\partial([\phi]\varphi, \Pi) &= \begin{cases} \mathbf{E}(\varphi) & \text{if } \Pi \models \phi \\ true & \text{if } \Pi \not\models \phi \end{cases} \\
\partial([\varrho?]\varphi, \Pi) &= \partial(nnf(\neg\varrho), \Pi) \vee \partial(\varphi, \Pi) \\
\partial([\varrho_1 + \varrho_2]\varphi, \Pi) &= \partial([\varrho_1]\varphi, \Pi) \wedge \partial([\varrho_2]\varphi, \Pi) \\
\partial([\varrho_1; \varrho_2]\varphi, \Pi) &= \partial([\varrho_1][\varrho_2]\varphi, \Pi) \\
\partial([\varrho^*]\varphi, \Pi) &= \partial(\varphi, \Pi) \wedge \partial([\varrho]\mathbf{T}_{\langle\varrho^*\rangle\varphi}, \Pi) \\
\partial(\mathbf{T}_\psi, \Pi) &= true \\
\partial(\mathbf{F}_\psi, \Pi) &= false
\end{aligned} \tag{2.17}$$

where  $\mathbf{E}(\varphi)$  recursively replaces in  $\varphi$  all occurrences of atoms of the form  $\mathbf{T}_\psi$  and  $\mathbf{F}_\psi$  by  $\mathbf{E}(\psi)$ .

Moreover, we define  $\partial(\varphi, \epsilon)$  which is inductively defined as Equation 2.17, except for the following cases:

$$\begin{aligned}
\partial(\langle\phi\rangle\varphi, \epsilon) &= false \\
\partial([\phi]\varphi, \epsilon) &= true
\end{aligned} \tag{2.18}$$

Note that  $\partial(\varphi, \epsilon)$  is always either *true* or *false*.

### 2.5.3 The LDL<sub>f</sub>2NFA algorithm

Algorithm 2.1 (LDL<sub>f</sub>2NFA) takes in input a LDL<sub>f</sub>/LTL<sub>f</sub> formula  $\varphi$  and outputs a NFA  $\mathcal{A}_\varphi = \langle 2^{\mathcal{P}}, Q, q_0, \delta, F \rangle$  that accepts exactly the traces satisfying  $\varphi$ . It is a variant of the algorithm presented in (De Giacomo and Vardi, 2015), and its correctness relies on the fact that every LDL<sub>f</sub>/LTL<sub>f</sub> formula  $\varphi$  can be associated a polynomial *alternating automaton on words* (AFW) accepting exactly the traces that satisfy  $\varphi$  and that every AFW can be transformed into an NFA (De Giacomo and Vardi, 2013).

The proposed algorithm requires that  $\varphi$  is in *negation normal form* (NNF), i.e. with negation symbols occurring only in front of propositions.

The function  $\partial$  used in lines 5, 12 and 15 is the one defined in sections 2.5.1 and 2.5.2; whether we are translating a LTL<sub>f</sub> or a LDL<sub>f</sub> formula, we use the function  $\partial$  from Definition 2.10 and from Definition 2.11, respectively.

---

**Algorithm 2.1.** LDL<sub>f</sub>2NFA: from LTL<sub>f</sub>/LDL<sub>f</sub> formula  $\varphi$  to NFA  $\mathcal{A}_\varphi$

---

```

1: input LDLf/LTLf formula  $\varphi$ 
2: output NFA  $\mathcal{A}_\varphi = \langle 2^{\mathcal{P}}, Q, q_0, \delta, F \rangle$ 
3:  $q_0 \leftarrow \{\varphi\}$ 
4:  $F \leftarrow \{\emptyset\}$ 
5: if  $(\partial(\varphi, \epsilon) = \text{true})$  then
6:    $F \leftarrow F \cup \{q_0\}$ 
7: end if
8:  $Q \leftarrow \{q_0, \emptyset\}$ 
9:  $\delta \leftarrow \emptyset$ 
10: while ( $Q$  or  $\delta$  change) do
11:   for ( $q \in Q$ ) do
12:     if  $(q' \models \bigwedge_{(\psi \in q)} \partial(\psi, \Pi))$  then
13:        $Q \leftarrow Q \cup \{q'\}$ 
14:        $\delta \leftarrow \delta \cup \{(q, \Pi, q')\}$ 
15:       if  $(\bigwedge_{(\psi \in q')} \partial(\psi, \epsilon) = \text{true})$  then
16:          $F \leftarrow F \cup \{q'\}$ 
17:       end if
18:     end if
19:   end for
20: end while

```

---

### How LDL<sub>f</sub>2NFA works

The NFA  $\mathcal{A}_\varphi$  for an LDL<sub>f</sub> formula  $\varphi$  is built in a forward fashion. Until convergence is reached (i.e. states and transitions do not change), the algorithm visits every state  $q$  seen until now, checks for all the possible transitions from that state and collects the results, determining the next state  $q'$ , the new transition  $(q, \Pi, q')$  and if  $q'$  is a final state. Intuitively, the delta function  $\partial$  emulates the semantic behavior of every LTL<sub>f</sub>/LDL<sub>f</sub> subformula after seeing  $\Pi$ .

States of  $\mathcal{A}_\varphi$  are sets of atoms (each atom is a quoted  $\varphi$  subformula) to be interpreted as conjunctions. The empty conjunction  $\emptyset$  stands for *true*.  $q'$  is a set of quoted subformulas of  $\varphi$  denoting a minimal interpretation such that  $q' \models \bigwedge_{(\psi \in q)} \partial(\psi, \Pi)$  (notice that we trivially have  $(\emptyset, p, \emptyset) \in \delta$  for every  $p \in 2^{\mathcal{P}}$ ).

The following result holds:

**Theorem 2.6** (De Giacomo and Vardi (2015)). *Algorithm LDL<sub>f</sub>2NFA is correct, i.e., for every finite trace  $\pi : \pi \models \varphi$  iff  $\pi \in \mathcal{L}(\mathcal{A}_\varphi)$ . Moreover, it terminates in at most an exponential number of steps, and generates a set of states  $S$  whose size is at most exponential in the size of the formula  $\varphi$ .*

In order to obtain a DFA, the NFA  $\mathcal{A}_\varphi$  can be determinized in exponential time (Rabin and Scott, 1959). Thus, we can transform a LTL<sub>f</sub>/LDL<sub>f</sub> formula into a DFA of double exponential size.

**Example 2.9.** In this example we see a run of the Algorithm 2.1 with the  $LTL_f$  formula  $\Box A$  ( $A$  atomic).

0. Set up:

$$\begin{aligned} q_0 &= \{\Box A\} \\ Q &= \{q_0, \emptyset\} \\ F &= \{q_0, \emptyset\} \quad (\text{because } \partial(\Box A, \epsilon) = \partial(\text{false } \mathcal{R} \neg A, \epsilon) = \text{true}) \\ \delta &= \{(\emptyset, \{\}, \emptyset), (\emptyset, \{A\}, \emptyset)\} \end{aligned}$$

1. Iteration: analyze  $q = \{\Box A\}$

- with  $\Pi = \{A\}$  we have

$$\begin{aligned} q' &\models \bigwedge_{(\psi \in q)} \partial(\psi, \Pi) \\ &\models \partial(\Box A, \Pi) \\ &\models \partial(A, \Pi) \wedge \partial(\bullet \Box A, \Pi) \\ &\models \text{true} \wedge (\Box A \vee \Box \text{false}) \end{aligned}$$

Notice that  $\text{true} \wedge (\Box A \vee \Box \text{false})$  is a *propositional formula* with  $LTL_f$  formulas as atoms. As a minimal interpretation we have both  $q' = \{\Box A\}$  and  $q' = \{\Box \text{false}\}$ . Since in both cases we have that  $\partial(\psi, \epsilon) = \text{true}$ , at the end of the iteration we have:

$$\begin{aligned} q_0 &= \{\Box A\} \\ Q &= \{q_0, \{\Box \text{false}\}, \emptyset\} \\ F &= \{q_0, \{\Box \text{false}\}, \emptyset\} \\ \delta &= \{(\emptyset, \{\}, \emptyset), (\emptyset, \{A\}, \emptyset), \\ &\quad (q_0, \{A\}, q_0), (q_0, \{A\}, \{\Box \text{false}\})\} \end{aligned}$$

- with  $\Pi = \{\}$  we have

$$\begin{aligned} q' &\models \bigwedge_{(\psi \in q)} \partial(\psi, \Pi) \\ &\models \partial(\Box A, \Pi) \\ &\models \partial(A, \Pi) \wedge \partial(\bullet \Box A, \Pi) \\ &\models \text{false} \wedge (\Box A \vee \Box \text{false}) \end{aligned}$$

Which is always false. Thus we do not change nothing.

2. Iteration: we already analyzed  $q = \{\Box A\}$ , so we analyze  $q = \{\Box \text{false}\}$

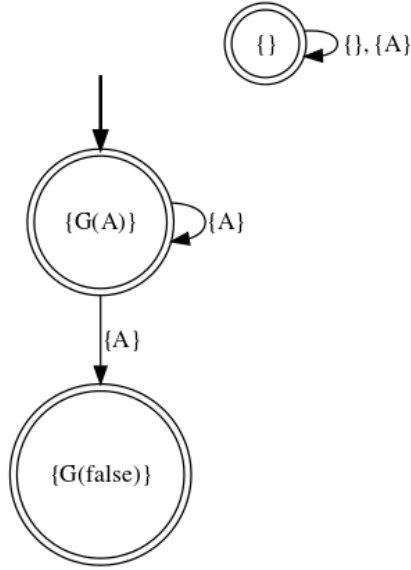
- Both with  $\Pi = \{\}$  and  $\Pi = \{A\}$  we have that:

$$q' \models \bigwedge_{(\psi \in q)} \partial(\psi, \Pi)$$

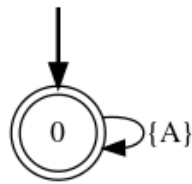
$$\begin{aligned}
&\models \partial(\Box false, \Pi) \\
&\models \partial(false, \Pi) \wedge \partial(\bullet \Box false, \Pi) \\
&\models false \wedge (\Box false \vee \Box false)
\end{aligned}$$

Which is always false. Thus we do not change nothing.

The NFA  $\mathcal{A}_\varphi = \langle 2^{\{A\}}, Q, q_0, \delta, F \rangle$  is depicted in Figure 2.2, whereas the associated DFA is in Figure 2.3.



**Figure 2.2.** The NFA associated to  $\Box A$ .  $G(A)$  stands for  $\Box A$



**Figure 2.3.** The DFA associated to  $\Box A$

**Example 2.10.** Analogously to what we did in 2.9, we see a run of the Algorithm 2.1, with the  $LTL_f$  formula  $\Diamond A$  ( $A$  atomic).

0. Set up:

$$\begin{aligned}
q_0 &= \{\Diamond A\} \\
Q &= \{q_0, \emptyset\} \\
F &= \{\emptyset\} \quad (\text{because } \partial(\Diamond A, \epsilon) = \partial(true \mathcal{U} A, \epsilon) = false) \\
\delta &= \{(\emptyset, \{\}, \emptyset), (\emptyset, \{A\}, \emptyset)\}
\end{aligned}$$



1. Iteration: analyze  $q = \{\Diamond A\}$

- with  $\Pi = \{A\}$  we have

$$\begin{aligned}
 q' &\models \bigwedge_{(\psi \in q)} \partial(\psi, \Pi) \\
 &\models \partial(\Diamond A, \Pi) \\
 &\models \partial(A, \Pi) \vee \partial(\Diamond A, \Pi) \\
 &\models \text{true} \vee (\text{"}\Diamond A\text{"} \wedge \text{"}\Diamond \text{true}\text{"})
 \end{aligned}$$

Since the propositional formula is trivially true, as a minimal interpretation we have  $q' = \emptyset$ . Considering that the empty conjunction is considered as *true* (as explained in Section 2.5), at the end of the iteration we have:

$$\begin{aligned}
 q_0 &= \{\Diamond A\} \\
 Q &= \{q_0, \emptyset\} \\
 F &= \{\emptyset\} \\
 \delta &= \{(\emptyset, \{\}, \emptyset), (\emptyset, \{A\}, \emptyset), (q_0, \{A\}, \emptyset)\}
 \end{aligned}$$

- with  $\Pi = \{\}$  we have

$$\begin{aligned}
 q' &\models \bigwedge_{(\psi \in q)} \partial(\psi, \Pi) \\
 &\models \partial(\Diamond A, \Pi) \\
 &\models \partial(A, \Pi) \vee \partial(\Diamond A, \Pi) \\
 &\models \text{false} \vee (\text{"}\Diamond A\text{"} \wedge \text{"}\Diamond \text{true}\text{"})
 \end{aligned}$$

As a minimal interpretation we have  $q' = \{\text{"}\Diamond A\text{"}, \text{"}\Diamond \text{true}\text{"}\}$ . Since  $\partial(\Diamond A, \epsilon) \wedge \partial(\Diamond \text{true}, \epsilon) = \text{false} \wedge \text{false} \neq \text{true}$ , we do not add  $q'$  to the accepting states  $F$ . Thus we have:

$$\begin{aligned}
 q_0 &= \{\Diamond A\} \\
 Q &= \{q_0, \emptyset, \{\Diamond A, \Diamond \text{true}\}\} \\
 F &= \{\emptyset\} \\
 \delta &= \{(\emptyset, \{\}, \emptyset), (\emptyset, \{A\}, \emptyset), \\
 &\quad (q_0, \{A\}, \emptyset), \\
 &\quad (q_0, \{\}, \{\Diamond A, \Diamond \text{true}\})\}
 \end{aligned}$$

2. Iteration: we already analyzed  $q = \{\Diamond A\}$ , so we analyze  $q = \{\Diamond A, \Diamond \text{true}\}$

- with  $\Pi = \{\}$  we have that:

$$q' \models \bigwedge_{(\psi \in q)} \partial(\psi, \Pi)$$

$$\begin{aligned}
& \models \partial(\Diamond A, \Pi) \wedge \partial(\Diamond true, \Pi) \\
& \models [\partial(A, \Pi) \vee \partial(\Diamond A, \Pi)] \wedge [\partial(true, \Pi) \vee \partial(\Diamond true, \Pi)] \\
& \models [\partial(A, \Pi) \vee (" \Diamond A " \wedge " \Diamond true ")] \wedge [true \vee (" \Diamond true " \wedge " \Diamond true ")] \\
& \models \partial(A, \Pi) \vee (" \Diamond A " \wedge " \Diamond true ") \\
& \models false \vee (" \Diamond A " \wedge " \Diamond true ")
\end{aligned}$$

As in the previous iteration, the minimal model is  $q' = \{ " \Diamond A ", " \Diamond true " \}$ . Hence we add a new transition  $(\{\Diamond A, \Diamond true\}, \{\}, \{\Diamond A, \Diamond true\})$ .

- with  $\Pi = \{A\}$  the delta-expansion is the same, except for the last step, where:

$$q' \models true \vee (" \Diamond A " \wedge " \Diamond true ")$$

The formula is always true, hence the minimal model is  $q' = \emptyset$  and we add a new transition  $(\{\Diamond A, \Diamond true\}, \{\}, \emptyset)$ .

The NFA  $\mathcal{A}_\varphi$  is then composed by:

$$\begin{aligned}
q_0 &= \{\Diamond A\} \\
Q &= \{q_0, \emptyset, \{\Diamond A, \Diamond true\}\} \\
F &= \{\emptyset\} \\
\delta &= \{(\emptyset, \{\}, \emptyset), (\emptyset, \{A\}, \emptyset), \\
&\quad (q_0, \{A\}, \emptyset), \\
&\quad (q_0, \{\}, \{\Diamond A, \Diamond true\}) \\
&\quad (\{\Diamond A, \Diamond true\}, \{\}, \{\Diamond A, \Diamond true\}) \\
&\quad (\{\Diamond A, \Diamond true\}, \{\}, \emptyset)\}
\end{aligned}$$

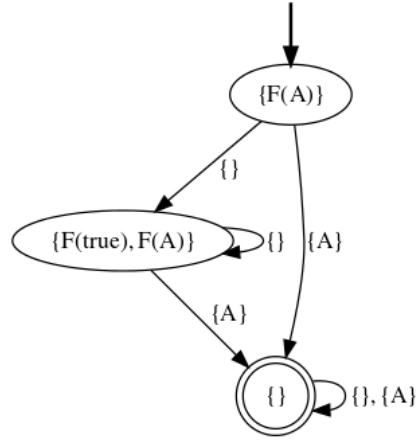
The NFA  $\mathcal{A}_\varphi = \langle 2^{\{A\}}, Q, q_0, \delta, F \rangle$  is depicted in Figure 2.4, whereas the associated DFA is in Figure 2.5.

**Example 2.11.** We list other examples of  $\mathcal{A}_\varphi$  given a  $LTL_f/LDL_f$  formula  $\varphi$ , obtained by Algorithm 2.1:

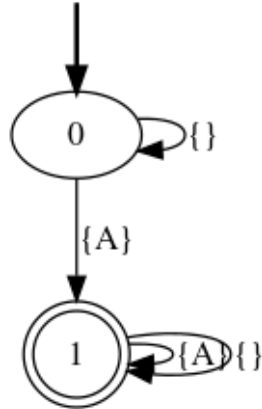
- *Conditional Response*: the  $LTL_f$  formula  $\varphi = \Box(A \Rightarrow \Diamond B)$  or equivalently the  $LDL_f$  formula  $\varphi = [true^*](\langle A \rangle tt \Rightarrow \langle true^* \rangle \langle B \rangle tt)$  translates into the automaton depicted in Figure 2.6.
- *Alternating sequence*: the  $LDL_f$  formula  $\varphi = \langle (A; B)^* \rangle End$  translates into the automaton depicted in Figure 2.7.

## 2.6 On-the-fly DFA

In this section we describe a way to evaluate a  $LTL_f/LDL_f$  formula without the need to build the entire automaton  $\mathcal{A}_\varphi$ . After that, we devise a new algorithm, a variant of  $LDL_f2NFA$ , that avoids the computation of the NFA, but directly translate the formula into a DFA. We provide some examples to clarify the presented topics.



**Figure 2.4.** The NFA associated to  $\Diamond A$ .  $F(A)$  stands for  $\Diamond A$



**Figure 2.5.** The DFA associated to  $\Diamond A$

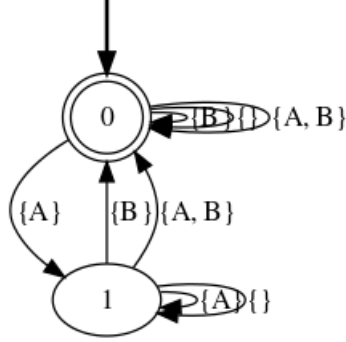
### 2.6.1 On-the-fly $LTL_f/LDL_f$ evaluation

In this section we describe an alternative method to evaluate a trace on a DFA without the need for constructing  $\mathcal{A}_\varphi$ , that we call *on-the-fly* (Brafman et al., 2018). The idea is, we progress all possible states that the NFA can be in, after consuming the next trace symbol, and accept the trace iff, once it has been completely read, the set of possible states contains a final state.

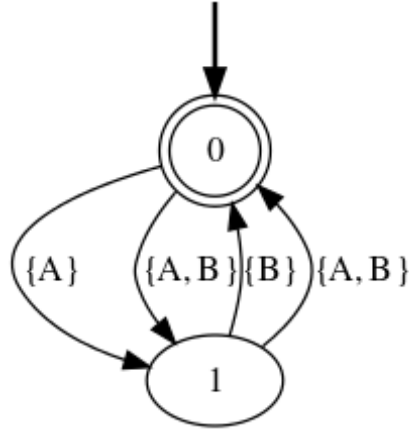
More formally, call a set of possible states for the NFA a macro state, let  $Q = \{q_1, \dots, q_n\}$  be the current macro state (initially  $Q = Q_0 = \{q_0\} = \{\{\varphi\}\}$ ), and let  $\Pi$  be the next trace symbol. Then, the successor macro state is the set  $Q'$  defined as follows:

$$Q' = \{q' | \exists q \in Q \text{ s.t. } q' \models \bigwedge_{(\psi \in q)} \partial(\psi, \Pi)\} \quad (2.19)$$

Notice that the condition  $q' \models \bigwedge_{(\psi \in q)} \partial(\psi, \Pi)$  is the same of the one in line 12 of



**Figure 2.6.** The DFA associated to  $\varphi = \Box(A \Rightarrow \Diamond B)$



**Figure 2.7.** The DFA associated to  $\varphi = \langle\langle A; B \rangle^*\rangle \text{End}$

**Algorithm 2.1.** Given an input trace  $\pi$ , we decide whether  $\pi \models \varphi$  by iterating the above procedure, starting from the initial state  $Q = Q_0$ , and accepting  $\pi$  iff the last state obtained includes  $\{true\}$ , considering their evaluation in the empty trace (i.e. with  $\partial(\psi, \epsilon)$ ). We denote the evaluation over the empty trace of a macro state  $Q = \{q_1, \dots, q_n\}$  as  $Q^\epsilon$ . Formally:

$$Q^\epsilon = \{\varphi \mid \varphi = \bigwedge_{\psi \in q_i} \partial(\psi, \epsilon)\} \quad (2.20)$$

**Example 2.12.** Consider Example 2.9 with  $\varphi = \Box A$ , we show how the on-the-fly evaluation of traces works. At the beginning, we have that

$$Q = Q_0 = \{\{\Box A\}\}$$

In this example, we ask the following questions:

1.  $\langle \rangle \models \varphi$ ? Does the empty trace  $\pi = \langle \rangle$  satisfy the formula  $\varphi$ ? We expect that the answer is yes, due to the semantics of  $\Box A$ . With the on-the-fly approach,

we need to compute for each NFA state  $q \in Q$  the conjunction between every  $\partial(\psi, \epsilon)$ , where  $\psi \in q$ . As said before, we consider the empty conjunction as  $\{true\}$ .

In our example, the computation gives us:

$$Q_0^\epsilon = \{\{true\}\}$$

because  $\partial(\Box A, \epsilon) = true$ . Since  $Q_0^\epsilon$  contains  $\{true\}$ ,  $Q_0^\epsilon$  is an accepting state, hence  $\pi \models \Box A$ , as expected.

2.  $\langle\{\}\rangle \models \varphi$ ? This time consider the trace  $\pi = \langle\{\}\rangle$  or, equivalently,  $\pi = \langle\neg A\rangle$ . We expect that the on-the-fly evaluation returns *false*, hence  $\pi \not\models \varphi$ . In order to answer, we need to progress the automaton on-the-fly for each element of the trace (in this case only one) and check if the last macro state is an accepting state by using the procedure explained in the previous case. The next macro state  $Q_1$  by applying Equation 2.19. Actually, it is computed as we did in Iteration 1 of Example 2.9 with  $\Pi = \{\}$ .

Since no  $q'$  can be found,  $Q_1 = \{\}$ , which is not an accepting state since  $\{true\} \notin Q_1^\epsilon$ . Hence,  $\pi \not\models \varphi$ .

3.  $\langle\{A\}\rangle \models \varphi$ ? Consider the trace  $\pi = \langle\{A\}\rangle$ . We expect that the on-the-fly evaluation returns *true*, hence  $\pi \models \varphi$ . We proceed, as in the previous case, to compute the next macro state  $Q_1$  by applying Equation 2.19. Actually, it is computed as we did in Iteration 1 with  $\Pi = \{A\}$ . As a minimal interpretation we have both  $q' = \{\Box A\}$  and  $q' = \{\Box false\}$ . Hence, the new macro state is  $Q_1 = \{\{\Box A\}, \{\Box false\}\}$ .

Since there are no other symbols in the trace  $\pi$  to be processed, we compute  $Q_1^\epsilon = \{\{true\}, \{true\}\} = \{\{true\}\}$ . Since  $\{true\} \in Q_1^\epsilon$ , we can say that  $\pi \models \varphi$ .

4.  $\langle\{A\}, \{\}\rangle \models \varphi$ ? Consider the trace  $\pi = \langle\{A\}, \{\}\rangle$ . We expect that the on-the-fly evaluation returns *false*, hence  $\pi \not\models \varphi$ . We proceed, as in the previous case, to compute the next macro states by applying Equation 2.19. Macro state  $Q_1$  is the same as we have seen in Case 3. We apply again the progression rule of Equation 2.19 with  $\Pi = \pi(2) = \{\}$ . As a minimal interpretation we have both  $q' = \{\Box A\}$  and  $q' = \{\Box false\}$ . Hence, the new macro state is  $Q_2 = \{\}$ , as we've seen in Iteration 2 of Example 2.9.

Since there are no other symbols in the trace  $\pi$  to be processed, we compute  $Q_2^\epsilon = \{\}$ . Since  $\{true\} \notin Q_2^\epsilon$ , we can say that  $\pi \not\models \varphi$ .

5.  $\langle\{\}, \{A\}\rangle \models \varphi$ ? Consider the trace  $\pi = \langle\{\}, \{A\}\rangle$ . We expect that the on-the-fly evaluation returns *false*, hence  $\pi \not\models \varphi$ . The macro state  $Q_1$  is the same as we have seen in Case 2, i.e.  $Q_1 = \{\}$ . We apply again the progression rule of Equation 2.19 with  $\Pi = \pi(2) = \{A\}$ . But this is trivially  $Q_2 = \{\}$ , by definition of the progression rule.

Since there are no other symbols in the trace  $\pi$  to be processed, we compute  $Q_2^\epsilon = \{\}$ . Since  $\{true\} \notin Q_2^\epsilon$ , we can say that  $\pi \not\models \varphi$ .

Notice how we use the same progression of Algorithm 2.1, but instead of aiming to build the entire automaton, we focus only on the states that are relevant for the satisfaction of the formula, given a trace.

**Example 2.13.** Analogously as we did in Example 2.12 for Example 2.9, we consider Example 2.10 with  $\varphi = \Diamond A$ , and we show how the on-the-fly evaluation of traces works also in this case. At the beginning, we have that

$$Q = Q_0 = \{\{\Diamond A\}\}$$

In this example, we ask the following questions:

1.  $\langle \rangle \models \varphi$ ? Does the empty trace  $\pi = \langle \rangle$  satisfy the formula  $\varphi$ ? We expect that the answer is no, due to the semantics of  $\Diamond A$ .

We observe that  $Q_0^\epsilon = \{\}$ , because  $\partial(\Diamond A, \epsilon) = \text{false}$ .

Since  $Q_0^\epsilon$  does not contain  $\{\text{true}\}$ ,  $Q_0^\epsilon$  is not an accepting state, hence  $\pi \not\models \Diamond A$ , as expected.

2.  $\langle \{\} \rangle \models \varphi$ ? This time consider the trace  $\pi = \langle \{\} \rangle$  or, equivalently,  $\pi = \langle \neg A \rangle$ . We expect that the on-the-fly evaluation returns *false*, hence  $\pi \not\models \varphi$ . In order to answer, we need to progress the automaton on-the-fly for each element of the trace (in this case only one) and check if the last macro state is an accepting state by using the procedure explained in the previous case. The next macro state  $Q_1$  by applying Equation 2.19. Actually, it is computed as we did in Iteration 1 of Example 2.10 with  $\Pi = \{\}$ . As a minimal interpretation we have  $q' = \{\Diamond A, \Diamond \text{true}\}$ . Hence, the new macro state is  $Q_1 = \{\{\Diamond A, \Diamond \text{true}\}\}$ . Now,  $Q_1^\epsilon = \{\{\text{false}\}\}$ , which is not an accepting state since  $\{\text{true}\} \notin Q_1^\epsilon$ . Hence,  $\pi \not\models \varphi$ .

3.  $\langle \{A\} \rangle \models \varphi$ ? Consider the trace  $\pi = \langle \{A\} \rangle$ . We expect that the on-the-fly evaluation returns *true*, hence  $\pi \models \varphi$ . We proceed, as in the previous case, to compute the next macro state  $Q_1$  by applying Equation 2.19. Actually, it is computed as we did in Iteration 1 with  $\Pi = \{A\}$ . As a minimal interpretation we have  $q' = \{\}$ . Hence, the new macro state is  $Q_1 = \{\emptyset\}$ .

Since there are no other symbols in the trace  $\pi$  to be processed, we compute  $Q_1^\epsilon = \{\{\text{true}\}\}$ . Since  $\{\text{true}\} \in Q_1^\epsilon$ , we can say that  $\pi \models \varphi$ .

4.  $\langle \{A\}, \{\} \rangle \models \varphi$ ? Consider the trace  $\pi = \langle \{A\}, \{\} \rangle$ . We expect that the on-the-fly evaluation returns *true*, and so  $\pi \models \varphi$ . We proceed, as in the previous case, to compute the next macro states by applying Equation 2.19. Macro state  $Q_1$  is the same as we have seen in Case 3. We apply again the progression rule of Equation 2.19 with  $\Pi = \pi(2) = \{\}$ , that leads to the new macro state is  $Q_2 = \{\emptyset\}$ . Notice that  $Q_1 = Q_2$ .

Since there are no other symbols in the trace  $\pi$  to be processed, we compute  $Q_2^\epsilon = \{\{\text{true}\}\}$ . Since  $\{\text{true}\} \in Q_2^\epsilon$ , we can say that  $\pi \models \varphi$ .

5.  $\langle \{\}, \{A\} \rangle \models \varphi$ ? Consider the trace  $\pi = \langle \{\}, \{A\} \rangle$ . We expect that the on-the-fly evaluation returns *true*, hence  $\pi \models \varphi$ . The macro state  $Q_1$  is the same as we have seen in Case 2, i.e.  $Q_1 = \{\{\Diamond A, \Diamond \text{true}\}\}$ . We apply again the progression rule of Equation 2.19 with  $\Pi = \pi(2) = \{A\}$ . But this is trivially  $Q_2 = \{\emptyset\}$  (as we've seen in Iteration 2 of Example 2.10), by definition of the progression rule.

Since there are no other symbols in the trace  $\pi$  to be processed, we compute  $Q_2^\epsilon = \{\{\text{true}\}\}$ . Since  $\{\text{true}\} \in Q_2^\epsilon$ , we can say that  $\pi \models \varphi$ .

Notice how we use the same progression of Algorithm 2.1, but instead of aiming to build the entire automaton, we focus only on the states that are relevant for the satisfaction of the formula, given a trace.

### 2.6.2 $\text{LDL}_f\text{2DFA}$ : a variant of $\text{LDL}_f\text{2NFA}$

Example 2.12 and 2.13 suggest a new way to translate  $\text{LTL}_f/\text{LDL}_f$  formulas to automata, that is a variant of  $\text{LDL}_f\text{2NFA}$  (Algorithm 2.1). We call it  $\text{LDL}_f\text{2DFA}$ , and directly translates an  $\text{LTL}_f/\text{LDL}_f$  formula to a DFA, instead of first translation into an NFA and then compute the DFA by determinization.

---

**Algorithm 2.2.**  $\text{LDL}_f\text{2DFA}$ : from  $\text{LTL}_f/\text{LDL}_f$  formula  $\varphi$  to DFA  $\mathcal{A}_\varphi$

---

```

1: input  $\text{LDL}_f/\text{LTL}_f$  formula  $\varphi$ 
2: output DFA  $\mathcal{A}_\varphi = \langle 2^{\mathcal{P}}, \mathcal{Q}, Q_0, \delta, F \rangle$   $\triangleright$  Notice:  $\mathcal{Q}$  is a set of macro states.
3:  $Q_0 \leftarrow \{\{\varphi\}\}$   $\triangleright$  the initial state of  $\mathcal{A}_\varphi$  is the initial macro state
4:  $F \leftarrow \emptyset$ 
5:  $\mathcal{Q} \leftarrow \{Q_0\}$ 
6:  $\delta \leftarrow \emptyset$ 
7: if  $(\{true\} \in Q_0^\epsilon)$  then
8:    $F \leftarrow F \cup \{Q_0\}$ 
9: end if
10: while ( $\mathcal{Q}$  or  $\delta$  change) do
11:   for  $(Q \in \mathcal{Q}, \Pi \in 2^{\mathcal{P}})$  do
12:      $Q' \leftarrow \{\}$ 
13:     for  $(q \in Q)$  do  $\triangleright$  Conceptually, the same loop of Algorithm 2.1, line 11
14:       if  $(q' \models \bigwedge_{(\psi \in q)} \partial(\psi, \Pi))$  then
15:          $Q' \leftarrow Q' \cup \{q'\}$ 
16:       end if
17:     end for
18:      $\mathcal{Q} \leftarrow \mathcal{Q} \cup \{Q'\}$   $\triangleright$  Add new macro state  $Q$  to the set of macro states  $\mathcal{Q}$ 
19:      $\delta \leftarrow \delta \cup \{(Q, \Pi, Q')\}$ 
20:     if  $(\{true\} \in Q'^\epsilon)$  then
21:        $F \leftarrow F \cup \{Q'\}$ 
22:     end if
23:   end for
24: end while

```

---

The idea behind Algorithm 2.2 is the following: build the DFA by doing exploration of automaton states and determinization *at the same time*. Indeed, each macro state tracks all the possible paths (according to the trace symbols processed) of the "implicit" NFA. The computation of the next NFA state, i.e. the for-loop at line 13, works in the same way of the for-loop in line 11 of Algorithm 2.1. For a single macro state  $Q$ , given a propositional interpretation  $\Pi$ , the operation is made for every NFA state  $q \in Q$ . The next macro state  $Q'$  is then composed by all the next NFA states  $q'$ . Given the triple  $Q, \Pi, Q'$ , we actually have a transition of the DFA  $\mathcal{A}_\varphi$ . Doing this operation for every macro state and for every interpretation, until convergence, will eventually lead to the exploration of every macro state and transitions among them.

The main advantage over Algorithm 2.1 is that we avoid the state explosion due to the determinization procedure, since we only process reachable states of the final DFA. We use this algorithm in the implementations of Chapter 3.

**Example 2.14.** We consider Example 2.9 but using Algorithm 2.2 for translation of  $\varphi = \Box A$  into a DFA.

0. Before the main loop, we have:

$$\begin{aligned} Q_0 &= \{\{\Box A\}\} \\ \mathcal{Q} &= \{Q_0\} \\ \delta &= \emptyset \\ F &= \{Q_0\} \quad (\text{because } \{true\} \in Q_0^\epsilon) \end{aligned}$$

The DFA at this stage is depicted in Figure 2.8a.

1. Iteration: Consider the macro state  $Q_0$ . Consider the (unique) NFA state  $\{\Box A\}$ .

With  $\Pi = \{A\}$  we generate the new macro state  $Q' = \{\{\Box A\}, \{\Box false\}\}$ . We add  $Q'$  to  $\mathcal{Q}$  and  $(Q_0, \Pi, Q')$  to  $\delta$ . We followed the same steps as we did in Example 2.12, Case 3.

With  $\Pi = \{\}$  we generate the new macro state  $Q' = \emptyset$ . We add  $Q'$  to  $\mathcal{Q}$  and  $(Q_0, \Pi, Q')$  to  $\delta$ . Since  $\{true\} \notin Q'^\epsilon$ , we do not add  $Q'$  to  $F$ . We followed the same steps as we did in Example 2.12, Case 2.

The DFA at this stage is depicted in Figure 2.8b.

2. Iteration: We already processed  $\mathcal{Q} = \{\{\Box A\}\}$ .

Consider the macro state  $Q = \emptyset$ . Since there exists no  $q \in Q$ , the for-loop at line 13, we add to  $\delta$  all the transitions of the form  $(\emptyset, \Pi, \emptyset)$ , for all  $\Pi \in 2^P$ .

Now consider the macro state  $Q = \{\{\Box A\}, \{\Box false\}\}$ .

Consider  $\Pi = \{A\}$ . For  $q = \{\Box A\}$  we generate the sub-macro state  $q' = \{\Box A\}$  and  $q' = \{\Box false\}$ . For  $q = \{\Box false\}$  we do not generate any sub-macro state. Hence, the resulting macro state is  $Q' = \{\{\Box A\}, \{\Box false\}\}$ . Since  $Q' \in \mathcal{Q}$ , we only add  $(Q, \Pi, Q')$  to  $\delta$ .

Consider  $\Pi = \{\}$ . For  $q = \{\Box A\}$  we do not generate any sub-macro state. For  $q = \{\Box false\}$  we do not generate any sub-macro state. Hence, the resulting macro state is  $Q' = \{\}$ . Since  $Q' \in \mathcal{Q}$ , we only add  $(Q, \Pi, Q')$  to  $\delta$ .

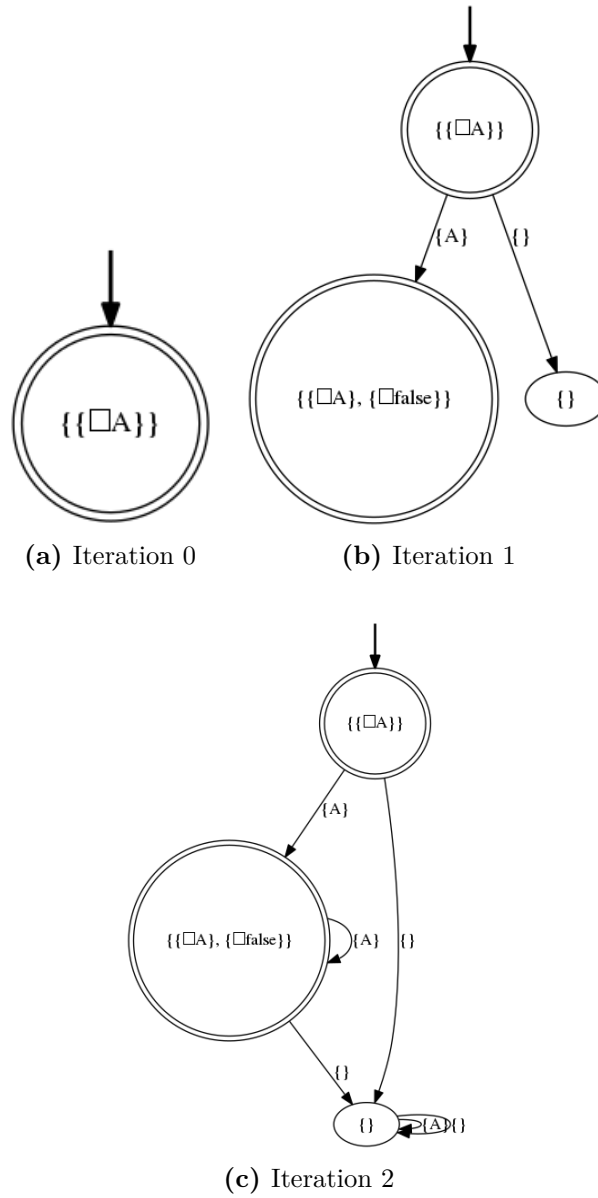
Since there are no other macro states nor propositional interpretation to process, the algorithm terminates. The final DFA is depicted in Figure 2.8c.

It is interesting to observe that the macro state  $Q = \{\{\Box A\}, \{\Box false\}\}$ , where we end after reading the symbol  $\{A\}$  from the initial state, is the set of NFA states (Figure 2.2) where we could end after reading the same symbol from the initial state, namely  $\{\Box A\}$  and  $\{\Box false\}$ . Analogous considerations can be made for other symbols and other macro states.

This observation makes clearer the true meaning of Algorithm 2.2 with respect to Algorithm 2.1: that is, the macro states keeps track of all the possible evolutions of the NFA with respect to a trace  $\pi$ . By reading all the possible symbols from every macro state, we will eventually discover all the relevant states of the DFA.

Furthermore, the minimization and trimming of the resulting automaton (depicted in Figure 2.8c) yields the one depicted in Figure 2.3, as an evidence of the equivalence between the two algorithms.





**Figure 2.8.** The automaton of Example 2.12, step by step.

**Example 2.15.** We consider Example 2.10 but using Algorithm 2.2 for translation of  $\varphi = \Diamond A$  into a DFA.

0. Before the main loop, we have:

$$\begin{aligned} Q_0 &= \{\{\Diamond A\}\} \\ \mathcal{Q} &= \{Q_0, \emptyset\} \\ \delta &= \emptyset \\ F &= \emptyset \quad (\text{because } \{true\} \notin Q_0^\epsilon) \end{aligned}$$

The DFA at this stage is depicted in Figure 2.9a.

1. Iteration: Consider the macro state  $Q_0$ . Consider the (unique) NFA state  $\{\Diamond A\}$ .

With  $\Pi = \{A\}$  we generate the new macro state  $Q' = \{\emptyset\}$ . We add  $Q'$  to  $\mathcal{Q}$  and  $(Q_0, \Pi, Q')$  to  $\delta$ . Since  $\{true\} \in Q'^\epsilon$ , we add  $Q'$  to  $F$ . We followed the same steps as we did in Example 2.13, Case 3.

With  $\Pi = \{\}$  we generate the new macro state  $Q' = \{\{\Diamond A, \Diamond true\}\}$ . We add  $Q'$  to  $\mathcal{Q}$  and  $(Q_0, \Pi, Q')$  to  $\delta$ . Since  $\{true\} \notin Q'^\epsilon$ , we do not add  $Q'$  to  $F$ . We followed the same steps as we did in Example 2.13, Case 2.

The DFA at this stage is depicted in Figure 2.9b.

2. Iteration: We already processed  $Q = \{\{\Box A\}\}$ .

Consider the macro state  $Q = \{\emptyset\}$ . Since the unique successor state of  $q = \emptyset$  is  $q' = \emptyset$ , the next macro state, for every symbol, is the same. Hence, we add to  $\delta$  all the transitions of the form  $(\{\emptyset\}, \Pi, \{\emptyset\})$ , for all  $\Pi \in 2^P$ .

Now consider the macro state  $Q = \{\{\Diamond A, \Diamond true\}\}$ .

Consider  $\Pi = \{A\}$ . As we've seen in Case 5 of Example 2.13, the next macro state is  $Q' = \{\emptyset\}$ , since that for  $q = \{\Diamond A, \Diamond true\}$  the successor state is  $q' = \emptyset$ . Since  $Q' \in \mathcal{Q}$ , we only add  $(Q, \Pi, Q')$  to  $\delta$ .

Consider  $\Pi = \{\}$ . The successor state of  $q = \{\Diamond A, \Diamond true\}$  is  $q' = q$ . Hence, the resulting macro state is  $Q' = \{\{\Diamond A, \Diamond true\}\}$ . Since  $Q' \in \mathcal{Q}$ , we only add  $(Q, \Pi, Q')$  to  $\delta$ . We followed the same steps as we did in Example 2.10, Iteration 2.

Since there are no other macro states nor propositional interpretation to process, the algorithm terminates. The final DFA is depicted in Figure 2.9c.

## 2.7 Reasoning in $LTL_f/LDL_f$

In this section we study the complexity of  $LTL_f/LDL_f$  reasoning (i.e. complexity of problems as defined in Definition 2.5).

**Theorem 2.7** (De Giacomo and Vardi (2013)). *Satisfiability, validity, and logical implication for  $LDL_f$  formulas are PSPACE-complete*

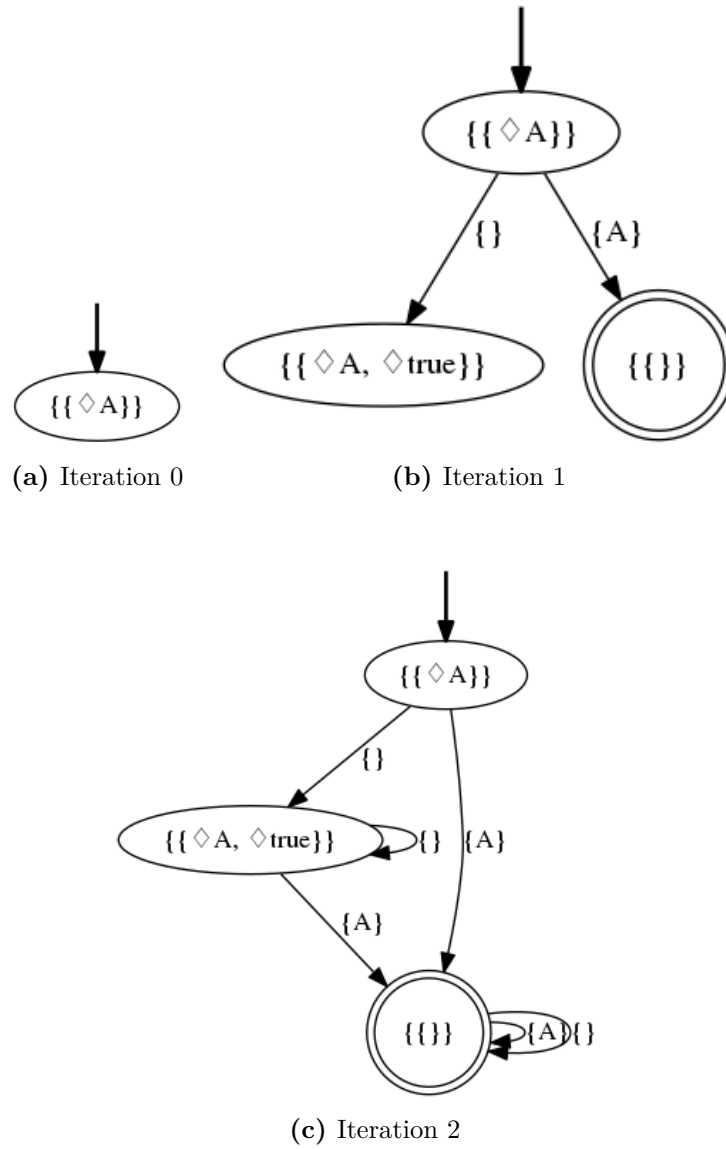
*Proof.* Given a  $LTL_f/LDL_f$   $\varphi$ , we can leverage Theorem 2.6 to solve these problems, namely:

- For  $LTL_f/LDL_f$  satisfiability we compute the associated NFA (as explained in Section 2.5 (which is an exponential step) and then check NFA for nonemptiness (NLOGSPACE).
- For  $LTL_f/LDL_f$  validity we compute the NFA associated to  $\neg\varphi$  (which is an exponential step) and then check NFA for nonemptiness (NLOGSPACE).
- For  $LTL_f/LDL_f$  logical implication  $\psi \models \varphi$  we compute the NFA associated to  $\psi \wedge \neg\varphi$  (which is an exponential step) and then check NFA for nonemptiness (NLOGSPACE).

□

## 2.8 Conclusions

In this chapter we provided the logical tools to face other topics in later chapters. We introduced several formal languages that allowed us to introduce  $LTL_f$  and  $LDL_f$ , focusing on their interesting properties. Moreover, we described in detail the procedure for translation from  $LTL_f/LDL_f$  formulas to DFAs, which yields an effective way to reasoning about  $LTL_f/LDL_f$  formulas.



**Figure 2.9.** The automaton of Example 2.13, step by step.

## Chapter 3

# FLLOAT

In this chapter we describe **FLLOAT** (From  $LTL_f/LDL_f$  to Automata), a software project written in Python. It is a porting of the [homonym software project](#) written in Java. It is the implementation of many of the topics described in Chapter 2.

### 3.1 Introduction

**Main features:** FLLOAT is a Python library that provides support for:

- Syntax, semantics and parsing of the following logic formalisms:
  - Propositional Logic;
  - Linear Temporal Logic on Finite Traces  $LTL_f$
  - Linear Dynamic Logic on Finite Traces  $LDL_f$ ;
- Conversion from  $LTL_f/LDL_f$  formula to NFA, DFA and DFA On-The-Fly

**Dependencies:** FLLOAT requires Python  $\geq 3.5$  and depends on the following packages:

- [PLY](#), a pure-Python implementation of the popular compiler construction tools [Lex](#) and [Yacc](#). It has been used for the parsing of PL and  $LTL_f/LDL_f$  formulas;
- [Pythomata](#), a Python package which provides support for NFA, DFA, determinization and minimization algorithms and reasoning on DFAs. It has been used for deal with  $\mathcal{A}_\varphi$ , the equivalent automaton of a  $LTL_f/LDL_f$  formula.

**Installation:** You can find the package on [PyPI](#), hence you can install it with:

```
pip install flloat
```

Please go [here](#) for further details.

The software is open source and is released under [MIT license](#).

## 3.2 Package structure

The package is structured as follows:

- `flloat.py`: the main module, it contains the implementation of the translation from  $LTL_f/LDL_f$  formulas to automata. The functions implemented here are called from methods of  $LTL_f/LDL_f$  formulas.
- `base/`: contains the abstract definitions used in other modules. The main modules are:
  - `Symbol.py` and `Symbols.py`, where have been defined the class `Symbol` to represent the atomic propositional symbols and the operator symbols;
  - `Alphabet.py`, which is an abstraction for manage a set of `Symbol`;
  - `Interpretation.py`, an abstract class denoting the semantics used for truth evaluation. E.g. for PL the corresponding interpretation is `PLInterpretation` (a set of `Symbol`), whereas for  $LTL_f/LDL_f$  we have `FiniteTrace`, which is a list of `PLInterpretation`.
  - `Formula.py`, the module containing the base class `Formula`. Every formula class extends `Formula`. In this module are defined also `AtomicFormula`, `Operator`, `BinaryOperator` etc., and how to build a syntax tree.
  - `truths.py` and `nnf.py` that provide abstract implementations for truth evaluations of formulas and negation normal form operations.
  - other abstractions definitions that are implemented for each extending subclass.
- `syntax/`: modules for each formalism (i.e. `pl.py`, `ltlf.py` and `ldlf.py`). In those modules are declared all the classes for representing formulas, implementing their truth evaluation procedure taking into account their correlation (e.g. `And` is the negative dual of `Or`, you can define `Implies` in terms of `Not` and `Or` etc.);
- `semantics/`: modules providing implementations for the semantics. E.g. you can find `PLInterpretation` and `FiniteTrace` cited before;
- `parser/`: modules where are defined the parsers of formulas in PL and  $LTL_f/LDL_f$ . They depends on PLY.

In the following sections we show typical use cases, describe the used APIs and look at their implementation.

## 3.3 Propositional Logic formulas

FLLOAT provides support for Propositional Logic (PL). In the following we will see both examples and explanation of the code.

### 3.3.1 Specify a PL formula

The easiest way to specify a propositional formula is to use a `PLParser`:

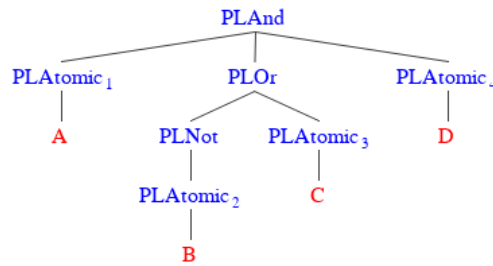
```

1 from flloat.parser.pl import PLParser
2
3 parser = PLParser()
```

For instance, the propositional formula  $\varphi = A \wedge (\neg B \vee C) \wedge D$  can be represented in this way:

```
1 phi = parser("A_&_(!B_|_C)_&_D")
```

Which generate the syntax tree in Figure 3.1. The logical operators symbols  $\neg, \wedge, \vee$  are replaced by, respectively,  $!, \&, |$ . In red are represented the symbols used to build the formula, for which the associated class is `Symbol`. On the other hand, in blue are shown the classes used to build the syntax tree, that represent a particular formula/operator in PL. They are: `PLAtomic`, used to represent atomic propositions; `PLNot`, `PLAnd` and `PLOr` for the associated logic operators. All of them extends from the base class `Formula`.



**Figure 3.1.** The syntax tree generated by the parsing of "A & (!B | C) & D", In red are represented instances of the class `Symbol`, whereas in blue the instances of the formulas/operators classes used to build the syntax tree.

We could construct the same syntax structure by using other APIs of FLLOAT:

**Listing 3.1.** A "manually" constructed PL formula.

```

1 from flloat.base.Symbol import Symbol
2 from flloat.syntax.pl import PLAtomic, PLAnd, PLOr, PLNot
3
4 A, B = Symbol("A"), Symbol("B")
5 C, D = Symbol("C"), Symbol("D")
6
7 atomic_A = PLAtomic(A)
8 atomic_B = PLAtomic(B)
9 atomic_C = PLAtomic(C)
10 atomic_D = PLAtomic(D)
11
12 not_B = PLNot(atomic_B)
13 B_or_C = PLOr([not_B, atomic_C])
14 phi = PLAnd([atomic_A, B_or_C, atomic_D])

```

Now follow the implementation details about how the above code can run.

### The Symbol class

`Symbol` is one of the base classes, declared in `Symbol.py`, whose code is reported in Listing 3.2.

**Listing 3.2.** The `Symbol` class.

```

1 class Symbol(str, Hashable):
2     def __init__(self, name:str):

```

```

3     str.__init__(self)
4     Hashable.__init__(self)
5     self.name = name
6
7     def _members(self):
8         return self.name

```

The class `Symbol` can be thought as a wrapper for the type `str`. Additionally, it inherits some useful properties from the class `Hashable`, implemented for performance purposes, which is reported in Listing 3.3:

Listing 3.3. The `Hashable` class.

```

1 from abc import ABC, abstractmethod
2 from copy import copy
3
4 class Hashable(ABC):
5
6     def __init__(self):
7         self._hash = None
8
9     @abstractmethod
10    def _members(self):
11        raise NotImplementedError
12
13    def __eq__(self, other):
14        if type(other) is type(self):
15            return self._members() == other._members()
16        else:
17            return False
18
19    def __hash__(self):
20        if self._hash is None:
21            self._hash = hash(self._members())
22        return self._hash
23
24    def __getstate__(self):
25        d = copy(self.__dict__)
26        d.pop("_hash")
27        return d
28
29    def __setstate__(self, state):
30        self.__dict__ = state
31        self._hash = None

```

Its main purpose is to allow memoization of the hash value. This is useful for immutable objects, since their hash value does not change during their life cycle. Furthermore, the mathematical domain of interest requires hash-based data structure such as sets and dictionaries, so it is a good idea to compute the hash only once and retrieve the precomputed value when needed. The `Hashable` class is inherited from many other classes, since the immutable property is satisfied by many entities. It requires that all the classes that inherits from it implements a `_members` function, which returns an hashable signature of the object that identify it uniquely. The methods `__getstate__` and `__setstate__` simply do not include



the precomputed hash value, since it might have a different value among different Python processes, due to hash randomization.

### The Formula abstract class and other abstraction

In `flloat/base/Formula.py` there are many abstraction that allow the construction of the syntax tree. The three dots `...` represent code omissions.

Listing 3.4. `flloat/base/Formula.py` module

```

1  from abc import abstractmethod, ABC
2  from typing import Sequence, Set
3
4  from flloat.base.Symbol import Symbol
5  from flloat.base.Symbols import Symbols
6  from flloat.base.hashable import Hashable
7
8
9  class Formula(Hashable):
10     def __init__(self):
11         super().__init__()
12
13         ...
14
15     class AtomicFormula(Formula):
16         def __init__(self, s:Symbol):
17             super().__init__()
18             self.s = s
19
20         def _members(self):
21             return self.s
22
23         ...
24
25     class Operator(Formula):
26         @property
27         def operator_symbol(self) -> str:
28             raise NotImplementedError
29
30         ...
31
32     class UnaryOperator(Operator):
33         def __init__(self, f: Formula):
34             super().__init__()
35             self.f = f.simplify()
36
37         def _members(self):
38             return (self.operator_symbol, self.f)
39
40         ...
41
42     class BinaryOperator(Operator):
43         def __init__(self, formulas:OperatorChilds):

```

```

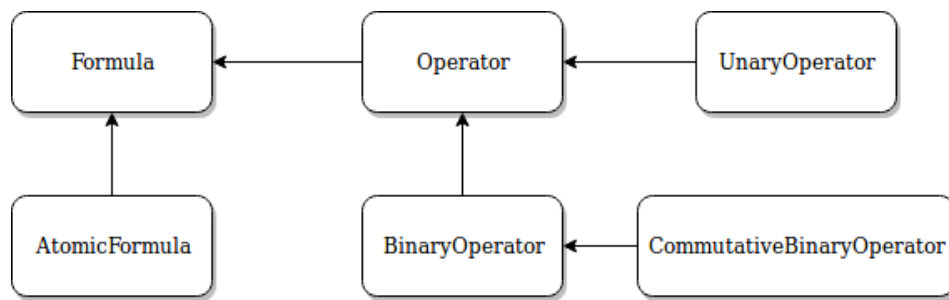
44     super().__init__()
45     assert len(formulas) >= 2
46     self.formulas = tuple(formulas)
47     self.formulas = self._popup()
48
49     ...
50
51     def _members(self):
52         return (self.operator_symbol, self.formulas)
53
54     def _popup(self):
55         """recursively find the same binary operator among
56         child formulas and pop up them at the same level"""
57         res = ()
58         for child in self.formulas:
59             if type(child) == type(self):
60                 superchilds = child._popup()
61                 res += superchilds
62             else:
63                 res += (child, )
64         return tuple(res)
65
66 class CommutativeBinaryOperator(BinaryOperator):
67     def __init__(self, formulas:OperatorChilds):
68         # Assuming idempotence: e.g. A & A == A
69         super().__init__(formulas)
70         # order does not matter -> set operation
71         # remove duplicates -> set operation
72         self.formulas_set = frozenset(self.formulas)
73         # unique representation -> sorting
74         self.members = tuple(
75             sorted(self.formulas_set, key=lambda x: hash(x))
76         )
77
78     ...

```

The class `Formula` is "the mother" of all the formulas. `AtomicFormula` represents a formula without any child in the syntax tree; it holds a `Symbol` object. `Operator` is an abstract class that represents a logical operator over formulas. It can be a `UnaryOperator` (e.g. the *not*  $\neg$ ), a `BinaryOperator` (e.g. the *implies*  $\Rightarrow$ ) and a `CommutativeBinaryOperator`, which is a binary operator that satisfies the commutative property (e.g. the *and*  $\wedge$  and the *or*  $\vee$ ). Every `Operator` has to have an operator symbol used for the conversion into strings.

When constructing a `BinaryOperator` instance, the list of subformulas is analyzed by looking for a subformula of the same type of the instance. This is a way to avoid multiple representations of the same syntax tree. For build a `BinaryOperator` formula, we need a *list* of subformulas, whereas for a `CommutativeBinaryOperator` we hold subformulas as a *set*, because order does not matter. However, for the latter case, in order to have a unique representation, we hold also a list of subformulas sorted by their hash.

We remark that, despite the name, `BinaryOperator` constructor take as input a *list* of formulas and not a *pair* of formulas. The reason, as stated before, is to



**Figure 3.2.** An inheritance schema for the classes in `flfloat/base/Formula.py`. Arrows stand for "extends".

reduce the depth of the syntax tree when the same binary operator is applied in chain many times.

In Figure 3.2 is depicted a schema to represent the relations between the just described abstractions.

#### PLFormulas: PLAtomic, PLNot, PLAnd, PLOr classes

The PLAtomic, PLNot, PLAnd, PLOr classes, as well as other PL classes for syntax support, are all declared in `flfloat/syntax/pl.py` (Listing 3.5).

**Listing 3.5.** The `flfloat/syntax/pl.py` module

```

1 class PLTruth(Truth):
2     @abstractmethod
3     def truth(self, i: PLInterpretation, *args):
4         raise NotImplementedError
5
6 class PLFormula(Formula, PLTruth, NNF):
7     def __init__(self):
8         Formula.__init__(self)
9
10    def all_models(self, alphabet: _Alphabet):
11        """Find all the possible interpretations
12        given a set of symbols"""
13        ...
14
15    def minimal_models(self, alphabet: _Alphabet):
16        """Find models of minimum size"""
17        ...
18
19 class PLBinaryOperator(BinaryOperator, PLFormula):
20     ...
21
22 class PLCommBinaryOperator(DualCommutativeOperatorNNF,
23                             PLFormula):
24     ...
25
26 class PLAtomic(AtomicFormula, PLFormula):
27
28     def truth(self, i: PLInterpretation, *args):
  
```

```

29     return self.s in i
30
31     def _to_nnf(self):
32         return self
33
34     def negate(self):
35         return PLNot(self)
36
37 class PLTrue(PLAtomic):
38     def __init__(self):
39         PLAtomic.__init__(self, Symbol(Symbols.TRUE.value))
40
41     def truth(self, *args):
42         return True
43
44     def negate(self):
45         return PLFalse()
46
47
48 class PLFalse(PLAtomic):
49     def __init__(self):
50         PLAtomic.__init__(self, Symbol(Symbols.FALSE.value))
51
52     def truth(self, *args):
53         return False
54
55     def negate(self):
56         return PLTrue()
57
58 class PLNot(NotTruth, PLFormula, NotNNF):
59     pass
60
61 class PLOr(PLCommBinaryOperator, OrTruth):
62     pass
63
64 class PLAnd(PLCommBinaryOperator, AndTruth):
65     pass
66
67 class PLImplies(PLBinaryOperator, ImpliesConvertible):
68     And = PLAnd
69     Or = PLOr
70     Not = PLNot
71
72 class PLEquivalence(PLCommBinaryOperator,
73                     EquivalenceConvertible):
74     And = PLAnd
75     Or = PLOr
76     Not = PLNot
77
78 PLOr.Dual = PLAnd
79 PLAnd.Dual = PLOr

```

Every propositional formula extends from `PLFormula`. It implements the methods `all_models` that, given an alphabet, returns all the propositional interpretations that satisfy the formula. `minimal_models` returns only the minimal ones. `PLBinaryOperator` and `PLCommBinaryOperator` are conceptually the same of `BinaryOperator` and `CommutativeBinaryOperator` described before.

Other classes define the possible PL formulas: `PLAtomic` is an `AtomicFormula`, hence its instances just hold a `Symbol` object; `PLTrue` and `PLFalse` are `PLAtomic` formulas with a fixed symbol that represent, respectively, propositional *true* and propositional *false*; `PLNot`, `PLAnd`, `PLOr` has been already described; `PLImplies` and `PLEquivalence` correspond to  $\Rightarrow$  and  $\Leftrightarrow$  logic operator, respectively.

It is interesting to observe that the implementation of some features such as syntax, semantics and conversion to Negation Normal Form (NNF) is inherited from other classes/interfaces, such as `Truth` and `NNF`. The implementation is made generically, i.e. by not assuming that we are talking about PL formulas. Indeed, the implementation for the same formulas (and, or, not etc.) is used for PL and `LTLf/LDLf`. We will discuss them in the next sections.

### 3.3.2 Parsing of PL formulas

Here we describe how the parsing of a string, representing a PL formula, to an instance of `PLFormula`, works. The implementation of the PL formula parsing is provided by the class `PLParser`, declared in the module `float/parser/pl.py`. In Listing 3.6 is shown an example of use of this feature. The formula to be parsed is  $A \wedge (\neg B \vee C) \wedge D$ . The associated string that `PLParser` is able to understand is `A & (!B | C) & D`.

Listing 3.6. How to parse a PL formula: `PLParser`

```
1 from float.parser.pl import PLParser
2 parser = PLParser()
3 phi = parser("A&(!B|C)&D")
```

In Listing 3.7 are shown the implementation of `PLParser` and the associated lexer, `PLLexer`. They extend a more general definition of `Parser` and `Lexer`, implemented in the module `float/base/parsing.py`.

Listing 3.7. `PLLexer` and `PLParser`

```
1 class PLLexer(Lexer):
2
3     def __init__(self):
4         super().__init__()
5
6     reserved = {
7         Symbols.TRUE.value : 'TRUE',
8         Symbols.FALSE.value : 'FALSE',
9     }
10
11     # List of token names. This is always required
12     tokens = (
13         'ATOM',
14         'NOT',
15         'AND',
16         'OR',
17         'IMPLIES',
```

```

18 'EQUIVALENCE',
19 'LPAREN',
20 'RPAREN'
21 ) + tuple(reserved.values())
22
23 # Regular expression rules for simple tokens
24 t_NOT      = sym2regex(Symbols.NOT)
25 t_AND      = sym2regex(Symbols.AND)
26 t_OR       = sym2regex(Symbols.OR)
27 t_IMPLIES  = sym2regex(Symbols.IMPLIES)
28 t_EQUIVALENCE = sym2regex(Symbols.EQUIVALENCE)
29 t_LPAREN   = sym2regex(Symbols.ROUND_BRACKET_LEFT)
30 t_RPAREN   = sym2regex(Symbols.ROUND_BRACKET_RIGHT)
31
32
33
34 def t_ATOM(self, t):
35     r'[a-zA-Z_][a-zA-Z_0-9]*'
36     # Check for reserved words
37     t.type = PLLexer.reserved.get(t.value, 'ATOM')
38     return t
39
40
41 # Yacc
42 class PLParser(Parser):
43
44     def __init__(self):
45         lexer = PLLexer()
46         precedence = (
47             ('left', 'EQUIVALENCE'),
48             ('right', 'IMPLIES'),
49             ('left', 'OR'),
50             ('left', 'AND'),
51             ('right', 'NOT'),
52         )
53         super().__init__("pl", lexer.tokens, lexer, precedence)
54
55     def p_formula_atom(self, p):
56         """formula : ATOM
57             | TRUE
58             | FALSE"""
59         if p[1]==Symbols.TRUE.value:
60             p[0] = PLTrue()
61         elif p[1]==Symbols.FALSE.value:
62             p[0] = PLFalse()
63         else:
64             p[0] = PLAtomic(Symbol(p[1]))
65
66     def p_formula_not(self, p):
67         'formula_␣:␣NOT␣formula'
68         p[0] = PLNot(p[2])
69

```

```

70 def p_formula_or(self, p):
71     'formula_⌞formula_⌞OR⌞formula'
72     p[0] = PLOr([p[1], p[3]])
73
74 def p_formula_and(self, p):
75     'formula_⌞formula_⌞AND⌞formula'
76     p[0] = PLAnd([p[1], p[3]])
77
78 def p_formula_implies(self, p):
79     'formula_⌞formula_⌞IMPLIES⌞formula'
80     p[0] = PLImplies([p[1], p[3]])
81
82 def p_formula_equivalence(self, p):
83     'formula_⌞formula_⌞EQUIVALENCE⌞formula'
84     p[0] = PLEquivalence([p[1], p[3]])
85
86 def p_formula_expression(self, p):
87     'formula_⌞LPAREN⌞formula⌞RPAREN'
88     p[0] = p[2]

```

This features depends from `PLY`, which provides utilities for both tokenization and parsing of a string. We do not go into details about how `PLY` works.

For the lexer part, we defined the available tokens for the parsing phase, as well as reserved words/symbols/regular expression for each of them (e.g. look at the regex for the `ATOM` token at line 34). The symbols used for the operators are defined in the module `float/base/Symbols.py`.

For the parsing part, we defined how to construct a syntax tree from the list of tokens. For each type of formula, we defined inductively a set of rules that, when activated, create the proper `PLFormula` object.

For example, consider line 55, where has been defined the rule for generate the atomic formulas. Depending on the type of token that activated the rule, we generate a `PLTrue`, `PLFalse` or a `PLAtomic` formula with the symbol provided by the tokenizer. Analogous consideration are made for other types of rules and formulas. Eventually, when there is no rule to apply, the entire syntax tree has been built. The object that we built "manually" (see Listing 3.1) now is built automatically by the parser.

This pattern is the same used for  $LTL_f/LDL_f$  parsers.

### 3.3.3 PL formulas evaluation

In this section we see how to compute the truth value of a `PLFormula`, given a propositional interpretation. We represent the propositional interpretation with the class `PLInterpretation`, that is trivially a set of `Symbols` (the ones that we consider true).

**Listing 3.8.** An example for PL formula truth evaluation

```

1 from flloat.parser.pl import PLParser
2 parser = PLParser()
3 phi = parser("A_⌞(!B_⌞C)⌞_⌞D")
4
5 from flloat.semantics.pl import PLInterpretation
6
7 A_true = PLInterpretation.fromStrings(["A"])

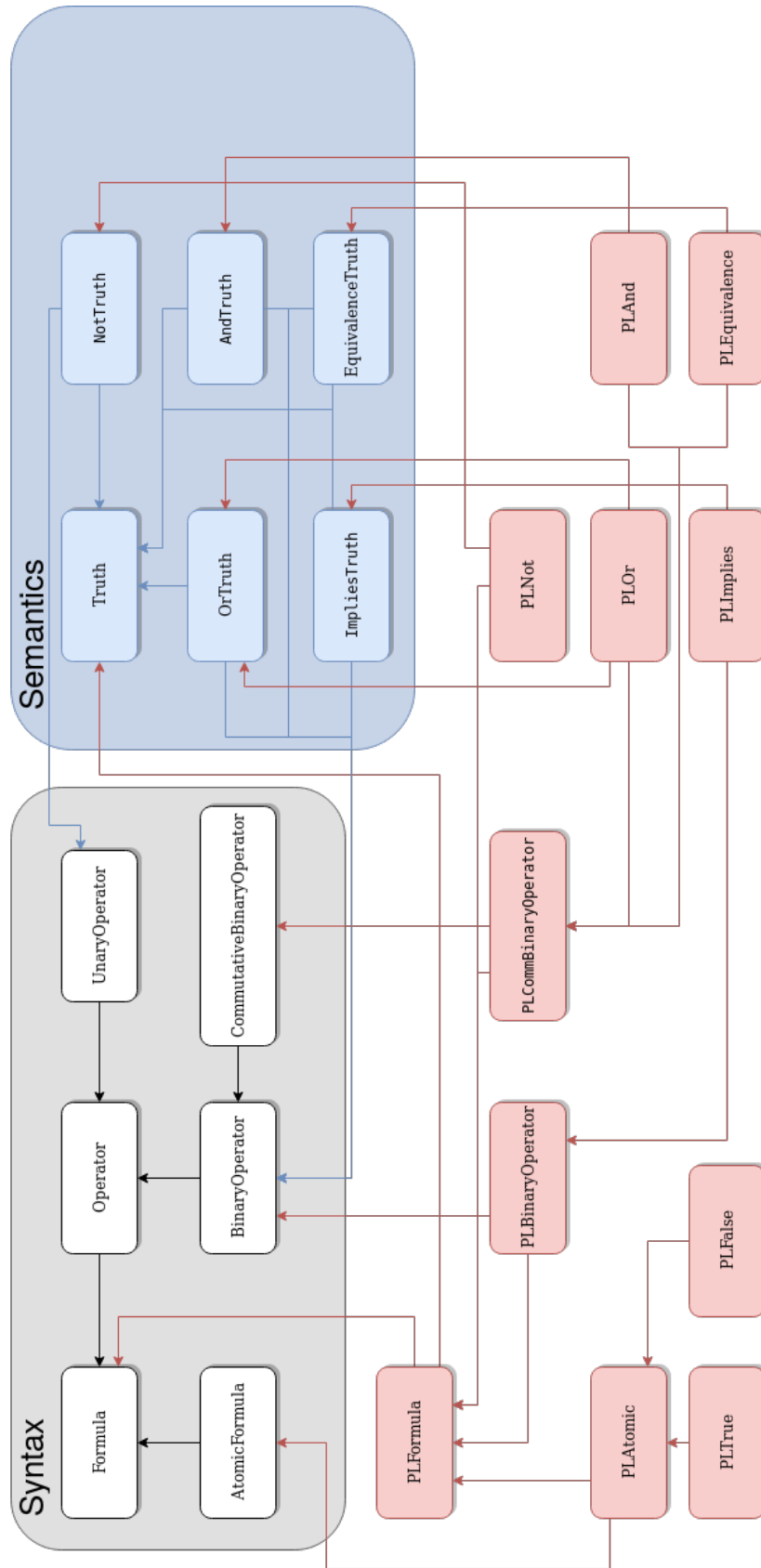
```

```
8 print(phi.truth(A_true)) #prints False
9 ACD_true = PLInterpretation.fromStrings(["A", "C", "D"])
10 print(phi.truth(ACD_true)) #prints True
```

You can notice that we can easily define a `PLInterpretation` by using the static method `fromStrings` and providing a collection of strings corresponding to the propositions that we want to make *true*. The evaluation is made by the method `truth`, that is implemented by every `PLFormula`. As stated in Section 3.3.1, the actual implementation of the (trivial) procedures to evaluate the truths of logic operators are defined in `float/base/truths.py`, that are inherited by the different `PLFormulas`. The truth evaluation of atomic formulas such as `PLAtomic`, `PLTrue` and `PLFalse` are, respectively, in lines 28, 41 and 52 of Listing 3.5.

In Figure 3.3 is depicted the schema representing the relations between the classes used for logic formulas discussed until now. It is worth to notice the separation between the implementation of the syntax and the one of the semantics, from where the propositional formula classes inherits their functionalities. This organization is replicated in both  $LTL_f$  and  $LDL_f$  implementations.





**Figure 3.3.** A schema representing the relations between the FLLOAT classes, that extends the one in Figure 3.2. The arrows mean "extends". In white the classes that implement how the syntax tree is built. In light blue the classes that implement the truth evaluation of operators. In red the classes that represent PL formulas. PL classes inherit from both syntax classes and semantics classes without implementing again both syntax and semantics functionalities.

### 3.4 LTL<sub>f</sub> formulas

As stated before, the software architecture to support LTL<sub>f</sub>/LDL<sub>f</sub> formulas is pretty the same of the one presented in Section 3.3. The core syntax functionalities (i.e. how to build the syntax tree of a LTL<sub>f</sub>/LDL<sub>f</sub> formula) are inherited from the classes defined in `ffloat/base/Formula.py`, as described in Section 3.3.1. The semantics functionalities for propositional logic operator (e.g. And, Or, Not) are inherited from `Truths` classes declared in `ffloat/base/truths.py` (see Section 3.3.3). Other ad-hoc data structures and classes are defined in order to support LTL<sub>f</sub>/LDL<sub>f</sub> formulas.

In this section we describe the LTL<sub>f</sub> part analogously to what we did in Section 3.3 for PL. In the next section we do the same for LDL<sub>f</sub>.

#### 3.4.1 Specify LTL<sub>f</sub> formulas

In this section we will see how to build LTL<sub>f</sub> formulas.

In Listing 3.9 is shown the construction of the main LTL<sub>f</sub> formulas.

**Listing 3.9.** Examples of LTL<sub>f</sub> formulas

```

1 from ffloat.base.Symbol import Symbol
2 from ffloat.syntax.ltlf import LTLfAtomic, LTLfAnd,\
3 LTLfOr, LTLfNot, LTLfImplies, LTLfEquivalence,\
4 LTLfNext, LTLfWeakNext, LTLfUntil, LTLfRelease, \
5 LTLfAlways, LTLfEventually
6 A_sym, B_sym = Symbol("A"), Symbol("B")
7 A = LTLfAtomic(A_sym)
8 B = LTLfAtomic(B_sym)
9
10 # classic logic operators
11 A_and_B = LTLfAnd([A, B])
12 A_or_B = LTLfOr([A, B])
13 not_A = LTLfNot(A)
14 A_implies_B = LTLfImplies([A, B])
15 A_equivalence_B = LTLfEquivalence([A, B])
16
17 #temporal operators
18 next_A = LTLfNext(A)
19 wnext_A = LTLfWeakNext(A)
20 A_until_B = LTLfUntil([A, B])
21 A_release_B = LTLfRelease([A, B])
22 always_A = LTLfAlways(A)
23 eventually_A = LTLfEventually(A)

```

All the classes for construct LTL<sub>f</sub> formulas are declared in `ffloat/syntax/ltlf.py`. The formula `LTLfAtomic` and the operators `LTLfAnd`, `LTLfOr`, `LTLfNot`, `LTLfImplies` and `LTLfEquivalence` are conceptually the same of the PL formulas described in Section 3.3.1. They extends `LTLfFormula` instead of `PLFormula` and inherits the syntax functionalities from the base classes in `ffloat/base/Formula.py`.

The novelty is about the temporal formulas `LTLfNext`, `LTLfWeakNext`, `LTLfUntil`, `LTLfRelease`, `LTLfAlways` and `LTLfEventually`, each of them extending `LTLfTemporalFormula`, which in turn extends `LTLfFormula`. They inherit the syntax functionalities from `UnaryOperator` and `BinaryOperator`, in a straightforward way.

Observe that the operators `LTLfWeakNext`, `LTLfRelease`, `LTLfAlways` and `LTLfEventually` are *abbreviations*, in the sense that they can be rewritten by only using combina-

tions of LTLfNot, LTLfNext and LTLfUntil. In order to reuse the code as much as possible, the relevant methods of these operators (such as `truth` or `to_nnf`, which put the formula in NNF) actually first *convert* the formula to an equivalent one by using only  $\neg$ ,  $\bigcirc$ ,  $\mathcal{U}$ , and then call the method on the newly created formula.

In order to do that, these classes extend the `BaseConvertibleFormula` declared in `float/base/convertible.py` (Listing 3.10). This concept is used also in the implementation of support for LDL<sub>f</sub> formulas.

**Listing 3.10.** The `BaseConvertibleFormula` abstraction, in `float/base/convertible.py`.

```

1  ...
2
3  class ConvertibleFormula(Formula):
4  @abstractmethod
5  def _convert(self):
6  raise NotImplementedError
7
8  class BaseConvertibleFormula(ConvertibleFormula, Truth, NNF):
9  def truth(self, *args):
10 return self._convert().truth(*args)
11
12 def _to_nnf(self):
13 return self._convert().to_nnf()
14
15 def negate(self):
16 return self._convert().negate()
17
18 ...

```

Indeed, by continuing from Listing 3.9:

```

1 print(wnext_A._convert()) #prints '!(X(!(A)))'
2 print(A_release_B._convert()) #prints '!(!!(A) U !(B))'
3 print(eventually_A._convert()) #prints '(true U A)'
4 print(always_A._convert()) #prints '!((true U !(A)))'

```

### LTLfFormula: the base abstraction for LTL<sub>f</sub> formulas

In Listing 3.11 is shown the implementation of `LTLfFormula`, which is extended by every LTL<sub>f</sub> formula class. It extends the following classes (skipping the class `Formula`):

- `LTLfTruth` (line 1) which defines the method `truth` that takes in input a `FiniteTrace` (see Section 2.2.2 and Section 3.4.3) and a position from which evaluate the trace (by default 0, i.e. from the beginning).
- `NNF`, a class which defines the abstract method `to_nnf` that returns the same formula but in Negation Normal Form (NNF);
- `Delta`, a class which defines the abstract method `_delta`, (line 23) called from `delta` (line 10), that implements the delta function described in Section 2.5.1.

Every subclass implements those methods properly, by taking into account the definitions provided in Section 2.2.2 and Section 2.5.1.

Listing 3.11. The LTLfFormula abstraction ([float/syntax/ltlf.py](#))

```

1 class LTLfTruth(Truth):
2     @abstractmethod
3     def truth(self, i: FiniteTrace, pos: int = 0):
4         raise NotImplementedError
5
6
7 class LTLfFormula(Formula, LTLfTruth, NNF, Delta):
8
9     @lru_cache(maxsize=MAX_CACHE_SIZE)
10    def delta(self, i: PLInterpretation, epsilon=False):
11        f = self.to_nnf()
12        d = f._delta(i, epsilon)
13        if epsilon:
14            # By definition, if epsilon=True, then
15            # the result must be either PLTrue or PLFalse
16            # Now, the output is a Propositional Formula
17            # with only PLTrue or PLFalse as atomics
18            # Hence, we just evaluate the formula with a dummy PLInterpretation
19            d = PLTrue() if d.truth(None) else PLFalse()
20        return d
21
22    @abstractmethod
23    def _delta(self, i: PLInterpretation, epsilon=False):
24        """apply delta function, assuming that 'self' is a
25        LTLf formula in Negative Normal Form"""
26        raise NotImplementedError
27
28    @abstractmethod
29    def to_LDLf(self):
30        raise NotImplementedError
31
32    def to_automaton(self, labels:Set[Symbol]=None, on_the_fly=False,
33                    determinize=False, minimize=True):
34        if labels is None:
35            labels = self.find_labels()
36        if on_the_fly:
37            return DFAOTF(self)
38        elif determinize:
39            return to_automaton(self, labels, minimize)
40        else:
41            return to_automaton_(self, labels)

```

For the `delta` function it has been used a LRU caching mechanism, provided from the Python Standard Library `functools`. This is needed in order to speed-up the expensive and repeated computation of the function `delta`.

At line 29 is defined the abstract method `to_LDLf`, which transform the current LTL<sub>f</sub> formula instance into an LDL<sub>f</sub> equivalent formula (by leveraging the equivalences shown in Section 2.4.2). Every specific LTL<sub>f</sub> formula class implements the procedure to transform itself into an equivalent LDL<sub>f</sub> formula. In practice, the root formula of the LTL<sub>f</sub> syntax tree compute the result by transforming each child into an LDL<sub>f</sub> formula, in such a way that the entire syntax tree is properly rewritten.

At line 33 it is defined the method `to_automaton` which transform the LTL<sub>f</sub> formula into an equivalent automaton. We discuss it in later sections.

### 3.4.2 Parsing of LTL<sub>f</sub> formulas

In this part we show how to use the LTL<sub>f</sub> parsing feature, analogously to what we did in Section 3.3.2.

In Listing 3.12 is shown an equivalent approach to build LTL<sub>f</sub> formulas, as we did in Listing 3.9.

**Listing 3.12.** Parsing of LTL<sub>f</sub> formulas

```

1 from flloat.parser.ltlf import LTLfParser
2
3 parser = LTLfParser()
4
5 # classic logic operators
6 A_and_B = parser("A_&_B")
7 A_or_B = parser("A_|_B")
8 not_A = parser("!A")
9 A_implies_B = parser("A_->_B")
10 A_equivalence_B = parser("A_<->_B")
11
12 #temporal operators
13 next_A = parser("X_A")
14 wnext_A = parser("WX_A")
15 A_until_B = parser("A_U_B")
16 A_release_B = parser("A_R_B")
17 always_A = parser("G_A")
18 eventually_A = parser("F_A")

```

The syntax of classic logic operators is the same of the `PLParser` presented in Section 3.3.2. The implementation of `LTLfParser` can be found in [float/parser/ltlf.py](#). It has been implemented in the same fashion of `PLParser`, but with more reserved symbols (X, WX, U, R, G, and F for represent, repsectively, the LTL<sub>f</sub> operators  $\bigcirc, \bullet, \mathcal{U}, \mathcal{R}, \Box, \Diamond$ ) and their associated parsing rules. Combinations and arbitrary nesting of those operators with "(" and ")" are allowed.

Notice how rules defined from line 6 reflects the ones defined in Section 2.2.1.

**Listing 3.13.** An extract from [float/parser/ltlf.py](#).

```

1 class LTLfParser(Parser):
2
3     ...
4
5     def p_formula(self, p):
6         """formula : formula EQUIVALENCE formula
7                     | formula IMPLIES formula
8                     | formula OR formula
9                     | formula AND formula
10                    | formula UNTIL formula
11                    | formula RELEASE formula
12                    | EVENTUALLY formula
13                    | ALWAYS formula
14                    | NEXT formula

```

```

15         | WEAK_NEXT formula
16         | NOT formula
17         | TRUE
18         | FALSE
19         | ATOM""
20     if len(p) == 2:
21         if p[1] == Symbols.TRUE.value:
22             p[0] = LTLfTrue()
23         elif p[1] == Symbols.FALSE.value:
24             p[0] = LTLfFalse()
25         else:
26             p[0] = LTLfAtomic(Symbol(p[1]))
27     elif len(p) == 3:
28         if p[1] == Symbols.NEXT.value:
29             p[0] = LTLfNext(p[2])
30         elif p[1] == Symbols.WEAK_NEXT.value:
31             p[0] = LTLfWeakNext(p[2])
32         elif p[1] == Symbols.EVENTUALLY.value:
33             p[0] = LTLfEventually(p[2])
34         elif p[1] == Symbols.ALWAYS.value:
35             p[0] = LTLfAlways(p[2])
36         elif p[1] == Symbols.NOT.value:
37             p[0] = LTLfNot(p[2])
38     elif len(p) == 4:
39         l, o, r = p[1:]
40         if o == Symbols.EQUIVALENCE.value:
41             p[0] = LTLfEquivalence([l, r])
42         elif o == Symbols.IMPLIES.value:
43             p[0] = LTLfImplies([l, r])
44         elif o == Symbols.OR.value:
45             p[0] = LTLfOr([l, r])
46         elif o == Symbols.AND.value:
47             p[0] = LTLfAnd([l, r])
48         elif o == Symbols.UNTIL.value:
49             p[0] = LTLfUntil([l, r])
50         elif o == Symbols.RELEASE.value:
51             p[0] = LTLfRelease([l, r])
52         else:
53             raise ValueError
54     else:
55         raise ValueError
56
57     def p_expr_paren(self, p):
58         """formula : LPAREN formula RPAREN"""
59         p[0] = p[2]

```

### 3.4.3 LTL<sub>f</sub> formulas evaluation

The semantics of LTL<sub>f</sub> (as well as of LDL<sub>f</sub>) is provided by finite traces  $\pi$  over a given set of propositional symbols  $\mathcal{P}$ , i.e. a sequence of propositional interpretations  $\Pi \in 2^{\mathcal{P}}$ .

FLLOAT allow to construct `FiniteTrace` instances in a easy way: providing to the static method `FiniteTrace.fromStringSets` a list of sets of strings, where each string represent a propositional symbol of interest. In Listing 3.14 an example about how to use `FiniteTrace` is given.

**Listing 3.14.** Defining a finite trace  $\pi = \langle \{\}, \{A\}, \{A, B\} \rangle$  in FLLOAT.

```

1 from flloat.semantics.ldlf import FiniteTrace
2
3 t1 = FiniteTrace.fromStringSets([
4     {},
5     {"A"},
6     {"A", "B"},
7 ])
```

In Listing is shown how to evaluate  $LTL_f$  formulas.

**Listing 3.15.** Some examples about how evaluate  $LTL_f$  formulas.

```

1 from flloat.parser.ltlf import LTLfParser
2 from flloat.semantics.ldlf import FiniteTrace
3
4 parser = LTLfParser()
5
6 #temporal operators
7 next_A = parser("XA")
8 wnext_A = parser("WXA")
9 A_until_B = parser("AUB")
10 A_release_B = parser("ARB")
11 always_A = parser("GA")
12 eventually_A = parser("FA")
13
14 t1 = FiniteTrace.fromStringSets([
15     {},
16     {"A"},
17     {"B"},
18     {"A", "B"},
19 ])
20
21 next_A.truth(t1) #True
22 next_A.truth(t1, 3) #False, at the end of the trace
23
24 wnext_A.truth(t1, 1) #False
25 wnext_A.truth(t1, 3) #True, at the end of the trace
26
27 A_until_B.truth(t1) #False
28 A_until_B.truth(t1, 3) #True
29
30 A_release_B.truth(t1) #False
31 A_release_B.truth(t1, 2) #True
32
33 always_A.truth(t1) #False
34 always_A.truth(t1, 3) #True
35
36 eventually_A.truth(t1) #True
```

```
37 eventually_A.truth(t1, 4) #False, beyond the end of the trace
```

For the  $LTL_f$  classic logical operators, the implementation of the semantics is exactly the same of the one used for PL formulas (see Section 3.3.3).

For the  $LTL_f$  temporal operators, the implementation is provided case by case following the definition given in Section 2.2.2.

## 3.5 $LDL_f$ formulas

In this section we show how FLLOAT supports  $LDL_f$  formulas, as we did Section 3.3 and Section 3.4.

### 3.5.1 Specify $LDL_f$ formulas

In this part we will see how to build  $LDL_f$  formulas.

In Listing 3.16 is shown the construction of the main  $LDL_f$  formulas.

**Listing 3.16.** Examples of  $LDL_f$  formulas

```
1 from flloat.base.Symbol import Symbol
2 from flloat.syntax.ldlf import LDLfAtomic, LDLfAnd,\
3   LDLfOr, LDLfNot, LDLfImplies, LDLfEquivalence,\
4   LDLfBox, LDLfDiamond,\
5   LDLfLogicalTrue, LDLfLogicalFalse,\
6   RegExpPropositional, RegExpTest,\
7   RegExpUnion, RegExpSequence, RegExpStar
8 from flloat.parser.pl import PLParser
9
10 #propositional formulas
11 plparser = PLParser()
12 A = plparser("A")
13 B = plparser("B")
14 true = plparser("true") #PLTrue
15 false = plparser("false") #PLFalse
16 A_and_B = plparser("A_&_B")
17 A_or_B = plparser("A_|_B")
18 not_A = plparser("!A")
19
20 #regular expressions
21 regA = RegExpPropositional(A)
22 regB = RegExpPropositional(B)
23 regTrue = RegExpPropositional(true)
24 regAandB = RegExpPropositional(A_and_B)
25 seq_AB = RegExpSequence([
26   RegExpPropositional(A), RegExpPropositional(B)
27 ])
28 union_AB = RegExpUnion([
29   RegExpPropositional(A), RegExpPropositional(B)
30 ])
31 true_star = RegExpStar(regA)
32
33 #LDLf formulas
34 tt = LDLfLogicalTrue()
```



```

35 ff = LDLfLogicalFalse()
36 diamond_A_tt = LDLfDiamond(regA, tt)
37 diamond_B_tt = LDLfDiamond(regB, tt)
38 diamond_seqAB_tt = LDLfDiamond(seq_AB, tt)
39 diamond_unionAB_tt = LDLfDiamond(union_AB, tt)
40 diamond_testA_B = LDLfDiamond(RegExpTest(diamond_A_tt), diamond_B_tt)
41 eventually_A = LDLfDiamond(true_star, diamond_A_tt)
42 always_A = LDLfBox(true_star, diamond_A_tt)
43 last = LDLfBox(regTrue, ff)

```

The classes implementing the regular expressions are named `RegExp*` (at line 20), whereas the ones for  $LDL_f$  formulas are named `LDLf*` (at line 33). They are all declared in `float/syntax/ldlf.py`. The inheritance schema is pretty similar to the one of  $LTl_f$  formulas, with the difference that for the  $LDLfBox [\varrho]\varphi$  and  $LDLfDiamond \langle \varrho \rangle \varphi$  formulas we need a regular expression and a  $LDL_f$  subformula, hence supporting the double-inductive nature of  $LDL_f$ .

In Listing 3.17 is shown the code for the main regular expressions abstraction `RegExpFormula` (at line 12) and the abstraction for  $LDL_f$  temporal formulas `LDLfTemporalFormula` (at line 15). At line 35 and 23 there are the implementations of, respectively, the  $LDLfBox$  and  $LDLfDiamond$  formulas.

Notice the constructor of `LDLfTemporalFormula`, which takes a regular expression and a  $LDL_f$  formula. Notice also the definition of the abstract methods for implement the delta function (Section 2.5.2). `deltaDiamond` (line 31) compute the delta function for diamond  $LDL_f$  formulas, whereas `deltaBox` (line 44) compute the delta function for box  $LDL_f$  formulas.

**Listing 3.17.** Regular Expression abstract class `RegExpFormula` and  $LDL_f$  formulas `LDLfFormula`.

```

1 class DeltaRegExp(ABC):
2     @abstractmethod
3     def deltaDiamond(self, f:LDLfFormula, i: PLInterpretation,
4         epsilon=False):
5         raise NotImplementedError
6
7     @abstractmethod
8     def deltaBox(self, f:LDLfFormula, i: PLInterpretation,
9         epsilon=False):
10        raise NotImplementedError
11
12 class RegExpFormula(Formula, RegExpTruth, NNF, DeltaRegExp):
13     pass
14
15 class LDLfTemporalFormula(LDLfFormula):
16     def __init__(self, r:RegExpFormula, f:LDLfFormula):
17         super().__init__()
18         self.r = r
19         self.f = f
20
21     ...
22
23 class LDLfDiamond(LDLfTemporalFormulaNNF, FiniteTraceTruth):
24     temporal_brackets = "<>"
25

```

```

26 def truth(self, i: FiniteTrace, pos: int):
27     return any(self.r.truth(i, pos, j) and\
28         self.f.truth(i, j) for j in range(pos, i.last()+1))
29     # last + 1 in order to include the last step
30
31 def _delta(self, i: PLInterpretation, epsilon=False):
32     return self.r.deltaDiamond(self.f, i, epsilon)
33
34
35 class LDLfBox(ConvertibleFormula, LDLfTemporalFormulaNNF):
36     temporal_brackets = "[]"
37
38     def _convert(self):
39         return LDLfNot(LDLfDiamond(self.r, LDLfNot(self.f)))
40
41     def truth(self, i: FiniteTrace, pos: int):
42         return self._convert().truth(i, pos)
43
44     def _delta(self, i: PLInterpretation, epsilon=False):
45         return self.r.deltaBox(self.f, i, epsilon)

```

### 3.5.2 Parsing of $LDL_f$ formulas

The same  $LDL_f$  formulas shown in Listing 3.16 are equivalently expressed by using the `LDLfParser` declared in `float/parser/ldlf.py`, as shown in Listing 3.18.

Listing 3.18. Examples of  $LDL_f$  parsing

```

1 from flfloat.parser.ldlf import LDLfParser
2
3 parser = LDLfParser()
4
5 #LDLf formulas
6 tt = parser("<tt>")
7 ff = parser("<ff>")
8 diamond_A_tt = parser("<A>tt")
9 diamond_B_tt = parser("<B>tt")
10 diamond_seqAB_tt = parser("<A;B>tt")
11 diamond_unionAB_tt = parser("<A+B>tt")
12 diamond_testA = parser("<(<A>tt?)>tt")
13 eventually_A = parser("<true*><A>tt")
14 always_A = parser("[true*<A>tt")
15 end = parser("[true]ff")

```

Notice that the regular expression operators `;`, `+`, `*`, `?` are denoted by the associated ASCII characters. The diamond operator brackets are denoted by `<>`, whereas the box operator brackets by `[]`.

In Listing 3.19 are reported the rules for the processing of  $LDL_f$  tokens and the parsing of a `LDLfFormula`.

Listing 3.19. Parsing rules for  $LDL_f$  (excerpts from `float/parser/ldlf.py`)

```

1 ...
2 class LDLfParser(Parser):

```

```

3
4 def p_temp_formula(self, p):
5     """temp_formula : temp_formula EQUIVALENCE temp_formula
6                       | temp_formula IMPLIES temp_formula
7                       | temp_formula OR temp_formula
8                       | temp_formula AND temp_formula
9                       | BOXLSEPARATOR path BOXRSEPARATOR temp_formula
10                      | DIAMONDSEPARATOR path DIAMONDRSEPARATOR temp_formula
11                      | NOT temp_formula
12                      | TT
13                      | FF
14                      | END
15                      | LAST"""
16     ...
17
18 def p_path(self, p):
19     """path : path UNION path
20             | path SEQ path
21             | path STAR
22             | temp_formula TEST
23             | propositional"""
24     ...
25
26
27 def p_propositional(self, p):
28     """propositional : propositional EQUIVALENCE propositional
29                      | propositional IMPLIES propositional
30                      | propositional OR propositional
31                      | propositional AND propositional
32                      | NOT propositional
33                      | FALSE
34                      | TRUE
35                      | ATOM"""
36     ...
37
38 def p_expr_paren(self, p):
39     """temp_formula : LPAREN temp_formula RPAREN
40     path             : LPAREN path RPAREN
41     propositional    : LPAREN propositional RPAREN"""
42     p[0] = p[2]

```

Observe that `p_temp_formula` defines the parsing rules for  $LDL_f$  formulas; `p_path` defines the parsing rules for regular expressions; `p_propositional` defines the rules for propositional formulas. `p_expr_paren` allow for arbitrary nesting of round brackets.

### 3.5.3 From $LTL_f$ to $LDL_f$

Notice that you can build an  $LDL_f$  by first constructing an  $LTL_f$  formula, as described in Section 3.4.1 and Section 3.4.2, and then transform it into an equivalent  $LDL_f$  formula by calling the method `to_LDLf()`:

**Listing 3.20.** Specify  $LDL_f$  formulas from  $LTL_f$  formulas.

```

1 from flloa.parser.ltlf import LTLfParser
2 parser = LTLfParser()
3
4 next_A = parser("X(A)")
5 ldlf_next_A = next_A.to_LDLf()
6 print(ldlf_next_A.to_nnf())
7 #<true>((<A>(tt) & <true>(tt)))
8
9 wnext_A = parser("WX(A)")
10 ldlf_wnext_A = wnext_A.to_LDLf()
11 print(ldlf_wnext_A.to_nnf())
12 #[true]((<A>(tt) | [true](ff)))
13
14 always_A = parser("G(A)")
15 ldlf_always_A = always_A.to_LDLf()
16 print(ldlf_always_A.to_nnf())
17 #(((<true>(tt))? ; true))*((<A>(tt) | [true](ff)))

```

You can verify the correctness of the translation by looking at Section 2.4.1.

### 3.5.4 $LDL_f$ formulas evaluation

The truth evaluation of  $LDL_f$  formula is similar to the one of  $LTL_f$  (see Section 3.4.3). In Listing 4.2 we report a piece of code similar to Listing 3.15 but adapted for  $LDL_f$  formulas.

**Listing 3.21.** Examples of  $LDL_f$  truth evaluations

```

1 from flloa.parser.ldlf import LDLfParser
2 from flloa.semantics.ldlf import FiniteTrace
3
4 parser = LDLfParser()
5
6 #LDLf formulas
7 tt = parser("tt")
8 ff = parser("ff")
9 diamond_A_tt = parser("<A>tt")
10 diamond_B_tt = parser("<B>tt")
11 diamond_seqAB_tt = parser("<A;B>tt")
12 diamond_unionAB_tt = parser("<A+B>tt")
13 diamond_testA = parser("<(<A>tt?)>tt")
14 eventually_A = parser("<true*><A>tt")
15 always_A = parser("[true*><A>tt")
16 end = parser("[true]ff")
17
18 t1 = FiniteTrace.fromStringSets([
19     {},
20     {"A"},
21     {"B"},
22     {"A", "B"},
23 ])
24
25 tt.truth(t1) #True
26 ff.truth(t1) #False

```

```

27
28 diamond_A_tt.truth(t1) #False
29 diamond_A_tt.truth(t1, 1) #True
30
31 diamond_seqAB_tt.truth(t1) #False
32 diamond_seqAB_tt.truth(t1, 1) #False
33
34 diamond_unionAB_tt.truth(t1) #False
35 diamond_unionAB_tt.truth(t1, 1) #True
36
37 diamond_testA.truth(t1) #False
38 diamond_testA.truth(t1, 1) #True
39
40 eventually_A.truth(t1) #True
41 eventually_A.truth(t1, 4) #False
42
43 always_A.truth(t1) #False
44 always_A.truth(t1, 4) #True

```

The implementation of `truth` for each type of formula is defined according to the Definition 2.7.

### 3.6 $LTL_f/LDL_f$ formula to automaton

The core feature of FLLOAT is to support the conversion from  $LTL_f/LDL_f$  formula to automaton. We've seen how to do it theoretically by translation into NFA (Section 2.5.3), on-the-fly evaluation (Section 2.6.1) and direct computation of the DFA (Section 2.6.2).

The method to be used is `to_automaton`, implemented by every `LTLfFormula` and `LDLfFormula`, shown at line 33 of Listing 3.11. By different values of the flag parameters, it is possible to switch into every mode. The implementations of the algorithms and the data structures are in the module `ffloat/ffloat.py`.

FLLOAT uses the Python package `Pythomata` for automata support. However, we just explore the basic use cases with the `Pythomata` APIs.

In this section we explain how to use FLLOAT in the aforementioned cases.

#### 3.6.1 From $LTL_f/LDL_f$ formula to NFA

In Listing 3.22 is shown how to use the APIs to generate a NFA from a  $LTL_f/LDL_f$  formula. For simplicity we use `LTLfFormula` objects, but the same examples hold also for `LDLfFormula` instances.

Listing 3.22. From  $LTL_f/LDL_f$  to NFA

```

1 from ffloat.parser.ltlf import LTLfParser
2 parser = LTLfParser()
3 formula = parser("F(A_&_B)")
4
5 nfa = formula.to_automaton(determinize=False)

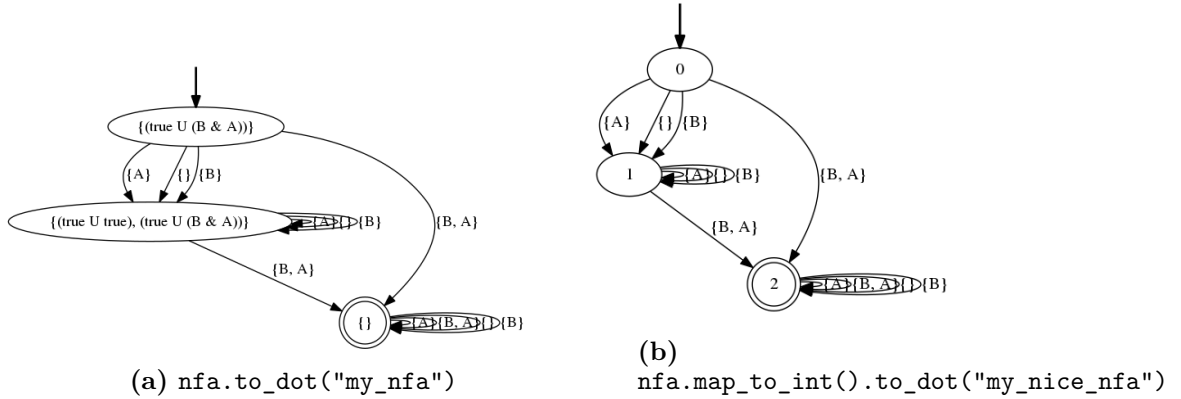
```

After the execution of these lines of code, in `nfa` is stored a `pythomata.base.NFA.NFA` object that represents an equivalent NFA. The method `LTLfFormula.to_automaton` at line 5 dispatches the execution to the function `to_automaton_` in `ffloat/ffloat.py`, which implements  $LDL_f2NFA$  (Algorithm 2.1).

Now you can print the NFA:

```
1 nfa.to_dot("my_nfa")
2 nfa.map_to_int().to_dot("my_nice_nfa")
```

That generates SVG images similar to the one shown in Figure 3.4a and Figure 3.4b. They are equivalent, but the latter is "nicer" because does not print the meaning of the states, which might not be of interest. Notice how this representation is structurally identical to the one shown in 2.4.



**Figure 3.4.** Examples of conversion from  $LTL_f/LDL_f$  to NFA and conversion into SVG images.

### 3.6.2 On the fly evaluation of $LTL_f/LDL_f$ formulas

FLLOAT supports also the on-the-fly evaluation of  $LTL_f/LDL_f$  formulas, i.e. whether a finite trace satisfies an  $LTL_f/LDL_f$  without building the entire automaton (see Section 2.6.1).

**Listing 3.23.** From  $LTL_f/LDL_f$  to NFA

```
1 from flloat.parser.ltlf import LTLfParser
2 parser = LTLfParser()
3 formula = parser("F(A_&_B)")
4
5 dfaotf = formula.to_automaton(on_the_fly=True)
```

In Listing 3.23 it is shown how to use this feature. Notice that the method used, `to_automaton`, is the same used in Section 3.6.1 for the  $LDL_f2NFA$  algorithm, but with the flag `on_the_fly` set to `True`. In this case, 5 returns an object of the class `flloat.float.DFAOTF`, that implements the concept explained in Section 2.6.1.

**Listing 3.24.** The class `DFAOTF`, that implements on-the-fly  $LTL_f/LDL_f$  formula evaluation.

```
1 class DFAOTF(Simulator):
2     """DFA on the fly"""
3
4     def __init__(self, f):
5         self.f = f.to_nnf()
6         self.reset()
7
```

```

8  def reset(self):
9      self.cur_state = frozenset([frozenset([self.f])])
10
11  def word_acceptance(self, action_set_list:List[PLInterpretation]):
12      self.reset()
13      for a in action_set_list:
14          self.make_transition(a)
15      return self.is_true()
16
17  def is_true(self):
18      return self._is_true(self.cur_state)
19
20  @staticmethod
21  def _is_true(Q):
22      ...
23
24  def get_current_state(self):
25      return self.cur_state
26
27  def make_transition(self, i:PLInterpretation):
28      self.cur_state = self._make_transition(self.cur_state, i)
29
30  @staticmethod
31  def _make_transition(Q, i: PLInterpretation):
32      ...

```

In Listing 3.24 it is shown the implementation of the class `DFAOTF`, that you can find in the module `floater/floater.py`. It extends `Simulator`, which is a Pythomata abstract class for inherit `word_acceptance` interface method. For the sake of brevity, we omitted some methods that have many lines. In the following we briefly describe the main methods:

- `reset()` set the current state to the initial macro state  $\{\{\varphi\}\}$  notice that here (and in many other places) we used the data structure `frozenset`, which is an hashable set. This is useful for deal with *sets of sets* (e.g. macro states);
- `__init__` is the constructor that takes in input a LTL<sub>f</sub>/LDL<sub>f</sub> formula, compute its NNF and reset the state to the initial state `reset`;
- `is_true()`, whose implementation depends on `_is_true()` returns `True` if the current state contains the empty conjunction  $\emptyset$ , which stand for *true* (see Section 2.6.1 for more details);
- `word_acceptance` takes in input a finite trace (i.e. a list of propositional interpretations) and returns `True` if the last state contains the empty conjunction (see `_is_true()`);
- `make_transition` and `_make_transition` implements the computation of the next macro state, given a macro state and a propositional interpretation. In other words, they implements the progression rule for macro states, stated in Equation 2.19.

Typical uses of a `DFAOTF` instance are evaluation of finite traces. By continuing with Listing 3.23, you can do it manually:

```

1 from flfloat.semantics.pl import PLInterpretation
2 from flfloat.semantics.pl import PLFalseInterpretation
3
4 A = PLInterpretation.fromStrings(["A"])
5 B = PLInterpretation.fromStrings(["B"])
6 empty = PLFalseInterpretation()
7 AB = PLInterpretation.fromStrings(["A", "B"])
8
9 dfaotf.make_transition(empty)
10 dfaotf.make_transition(A)
11 dfaotf.make_transition(B)
12 dfaotf.is_true() #False
13 dfaotf.make_transition(AB)
14 dfaotf.is_true() #True

```

Or by simply calling `word_acceptance` over a list of `PLInterpretation`:

```

1 from flfloat.semantics.ldlf import FiniteTrace
2
3 t1 = FiniteTrace.fromStringSets([
4     {},
5     {"A"},
6     {"B"},
7     {"A", "B"},
8 ])
9
10 dfaotf.word_acceptance(t1.trace) #True

```

### 3.6.3 From $LTL_f/LDL_f$ formula to DFA

Here we describe how FLLOAT supports the transformation from  $LTL_f/LDL_f$  formulas to DFA, i.e. the  $LDL_f2DFA$  algorithm (Algorithm 2.2).

We adapt Listing 3.22 by simply using the same method `to_automaton` but setting up the flag `determinize` to `True`.

**Listing 3.25.** From  $LTL_f/LDL_f$  to DFA

```

1 from flfloat.parser.ltlf import LTLfParser
2 parser = LTLfParser()
3 formula = parser("F(A_&_B)")
4
5 dfa = formula.to_automaton(determinize=True, minimize=False)

```

After the execution of these lines of code, in `dfa` is stored a `pythomata.base.DFA.DFA` object that represents an equivalent DFA. The method `LTLfFormula.to_automaton` at line 5 dispatches the execution to the function `to_automaton` in `float/float.py`, which implements  $LDL_f2DFA$  (Algorithm 2.1).

Now you can print the DFA:

```

1 dfa.to_dot("my_dfa")

```

Quite often is better to obtain the minimized automaton:

```

1 dfa_min = formula.to_automaton(determinize=True, minimize=True)
2 dfa.to_dot("my_nice_dfa")

```

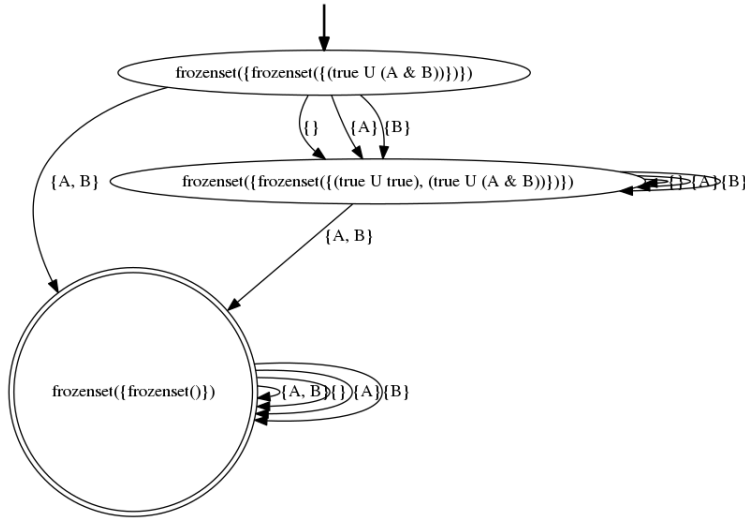


By simply calling the same method but with the flag `minimize` set to `True`.

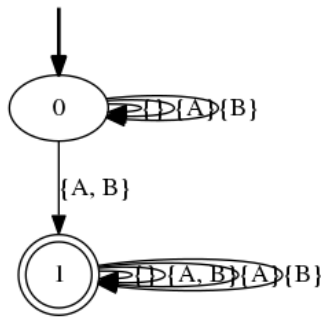
Another equivalent approach is:

```
1 nfa = formula.to_automaton(determinize=False)
2 dfa = nfa.determinize()
3 dfa_min = dfa.minimize().trim()
4 dfa.to_dot("my_nice_dfa")
```

In Figure 3.5a and Figure 3.5b are shown, respectively, the conversion to SVG images of the DFA not minimized and the one minimized. Notice how this representation is structurally identical to the one shown in 2.5.



(a) `dfa.to_dot("my_dfa")`



(b)  
`dfa_min.to_dot("my_nice_dfa")`

**Figure 3.5.** Examples of conversion from  $LTL_f/LDL_f$  to DFA and conversion into SVG images.

## 3.7 Conclusion

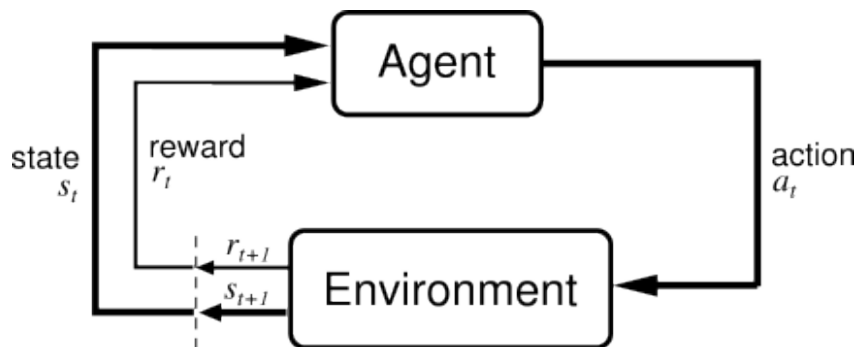
In this chapter we presented the FLLOAT Python package and its main features. We've seen the package structure, how FLLOAT supports PL and  $LTL_f/LDL_f$  formulas, and the conversion to automata from  $LTL_f/LDL_f$  formulas.

## Chapter 4

# RL for $LTL_f/ LDL_f$ Goals

### 4.1 Reinforcement Learning

Reinforcement Learning (Sutton and Barto, 1998) is a sort of optimization problem where an *agent* interacts with an *environment* and obtains a *reward* for each action he chooses and the new observed state. The task is to maximize a numerical reward signal obtained after each action during the interaction with the environment. The agent does not know a priori how the environment works (i.e. the effects of his actions), but he can make observations in order to know the new state and the reward. Hence, learning is made in a *trial-and-error* fashion. Moreover, it is worth to notice that in many situation reward might not be affected only from the last action but from an indefinite number of previous action. In other words, the reward can be *delayed*, i.e. the agent should be able to foresee the effect of his actions in terms of future expected reward. Figure 4.1 represent the interaction between the agent and the environment in this setting.



**Figure 4.1.** The agent and its interaction with the environment in Reinforcement Learning

In the next subsections we introduce some of the classical mathematical frameworks for RL: Markov Decision Process (MDP) and Non-Markovian Reward Decision Process (NMRDP).

### 4.2 Markov Decision Process (MDP)

A Markov Decision Process (MDP)  $\mathcal{M}$  is a tuple  $\langle S, A, T, R, \gamma \rangle$  containing a set of *states*  $S$ , a set of *actions*  $A$ , a *transition function*  $T : S \times A \rightarrow \text{Prob}(S)$  that returns for every pair state-action a probability distribution over the states, a *reward*

function  $R : S \times A \times S \rightarrow \mathbb{R}$  that returns the reward received by the agent when he performs action  $a$  in  $s$  and transitions in  $s'$ , and a *discount factor*  $\gamma$ , with  $0 \leq \gamma \leq 1$ , that indicates the present value of future rewards. With  $T(s, a, s')$  we denote the probability to end in state  $s'$  given the action  $a$  from  $s$ .

The discount factor  $\gamma$  deserves some attention. Its value highly influences the MDP, its solution, and how the agent interprets rewards. Indeed, if  $\gamma = 0$ , we are in the pure *greedy* setting, i.e. the agent is shortsighted and look only at the reward that it might obtain in the next step, by doing a single action. The higher  $\gamma$ , the longer the sight horizon, or the foresight, of the agent: the far rewards are taken into account for the current action choice. If  $\gamma < 1$  we are in the *finite horizon* setting: namely, the agent is intrinsically motivated to obtain rewards as fast as possible, depending on how  $\gamma$  is far from 1. When  $\gamma = 1$  we are in the *infinite horizon* setting, which means that the agent considers far rewards as they can be obtained in the next step. In other words, we may think the agent as *immortal*, since the time the agent spend to reach rewards does not matter anymore.

A *policy*  $\rho : S \rightarrow A$  for an MDP  $\mathcal{M}$  is a mapping from states to actions, and represents a solution for  $\mathcal{M}$ . Given a sequence of rewards  $R_{t+1}, R_{t+2}, \dots, R_T$ , the *expected discounted return*  $G_t$  at time step  $t$  is defined as:

$$G_t := \sum_{k=t+1}^T \gamma^{k-t-1} R_k \quad (4.1)$$

where can be  $T = \infty$  and  $\gamma = 1$  (but not both).

The *value function* of a state  $s$ , the *state-value function*  $v_\rho(s)$  is defined as the expected return when starting in  $s$  and following policy  $\rho$ , i.e.:

$$v_\rho(s) := \mathbb{E}_\rho[G_t | S_t = s], \forall s \in S \quad (4.2)$$

Similarly, we define  $q_\rho$ , the *action-value function for policy  $\rho$* , as:

$$q_\rho(s, a) := \mathbb{E}_\rho[G_t | S_t = s, A_t = a], \forall s \in S, \forall a \in A \quad (4.3)$$

Notice that we can rewrite 4.2 and 4.3 recursively, yielding the *Bellman equation*:

$$v_\rho(s) = \sum_{s'} P(s' | s, a) [R(s, a, s') + \gamma v_\rho(s')] \quad (4.4)$$

where we used the definition of the transition function:

$$T(s, a, s') = P(s' | s, a) \quad (4.5)$$

We define the *optimal state-value function* and the *optimal action-value function* as follows:

$$v^*(s) := \max_\rho v_\rho(s), \forall s \in S \quad (4.6)$$

$$q^*(s, a) := \max_\rho q_\rho(s, a), \forall s \in S, \forall a \in A \quad (4.7)$$

Notice that with 4.6 and 4.7 we can show the correlation between  $v_\rho^*(s)$  and  $q_\rho^*(s, a)$ :

$$q^*(s, a) = \mathbb{E}_\rho[R_{t+1} + \gamma v_\rho^*(S_{t+1}) | S_t = s, A_t = a] \quad (4.8)$$

We can define a partial order over policies using value functions, i.e.  $\forall s \in S. \rho \geq \rho' \iff v_\rho(s) \geq v_{\rho'}(s)$ . Now we give the definition of optimal policy:

**Definition 4.1.** An *optimal policy*  $\rho^*$  is a policy such that  $\rho^* \geq \rho$  for all  $\rho$ .

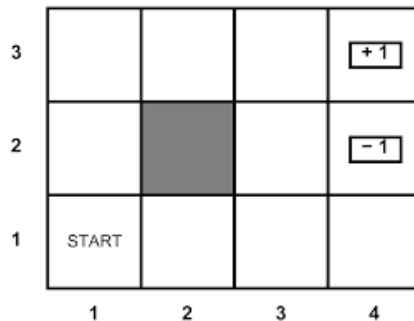
Given an MDP  $\mathcal{M}$ , a typical reinforcement learning problem is the following: find an optimal policy for  $\mathcal{M}$ , without knowing  $T$  and  $R$ . Notice that instead of explicit specification of the transition probabilities and rewards, the transition probabilities are accessed through a simulator that is restarted many times from a fixed or uniformly random initial state  $s_0 \in S$ . We call this way of structuring the learning process *episodic reinforcement learning*. Usually, in episodic reinforcement learning, we require the presence of one or more *goal states* where the simulation of the MDP ends and the task is considered completed, or a maximum time limit  $T$  for the number of action that can be taken by the agent in one single episode, and the overcoming of  $T$  determines the end of the episode. Optionally, *failure states* can be also defined, where the episode end similarly to goal states, but the task is considered failed.

### Examples

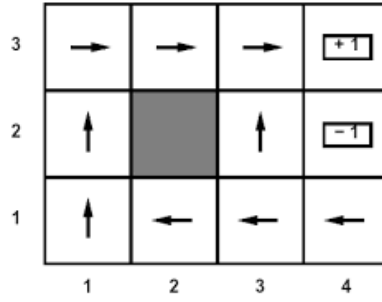
Many dynamic systems can be modeled as Markov Decision Processes.

**Example 4.1 (Gridworld).** Perhaps the most simple MDP used as a toy example is *Gridworld*, depicted in Figure 4.2. There are  $3 \times 4$  cells, i.e. states of the MDP  $S = \{s_{11}, s_{12}, \dots, s_{34}\} \setminus \{s_{22}\}$ . The agent can do four actions:  $A = \{Right, Left, Up, Down\}$ . The initial state is fixed and is  $s_0 = s_{11}$  and the agent can move in any of the adjacent and free cells from the current state. Assuming an episodic task, the goal is to reach  $s_{34}$ , and  $s_{24}$  represent a failure state. The state transition function  $T$  can be *deterministic*, i.e. the agent always succeed in performing actions, or non-deterministic, i.e. the effect of an action is determined by the probabilistic distribution returned by  $T(s, a)$ . An example of non-deterministic  $T$  is to give 90% of success (the agent moves in the chosen direction) and 10% of fail (the agent moves at either the right or left angle to the intended direction). If the move would make the agent walk into a wall (borders of the grid and  $s_{22}$ ), the agent stays in the same place as before. The reward function  $R(s, a, s')$  is defined as  $-1$  if  $s' = s_{24}$ , as  $1$  if  $s' = s_{34}$ , and  $-0.01$  otherwise. The small negative reward given at each transition is a popular mean for reward function design: it is called *step reward* and its purpose is to encourage the agent to finish the episode as fast as possible, with a priority proportional to the absolute value of the reward. The discount factor  $\gamma$  should be strictly higher than 0 because more than one step is needed to reach the goal state.

**Figure 4.2.** The Gridworld environment



An example of optimal policy is shown in Figure 4.3. As the reader can notice, the arrows represent the action that should be taken in a certain cell, in order to



**Figure 4.3.** An example of optimal policy for the Gridworld environment

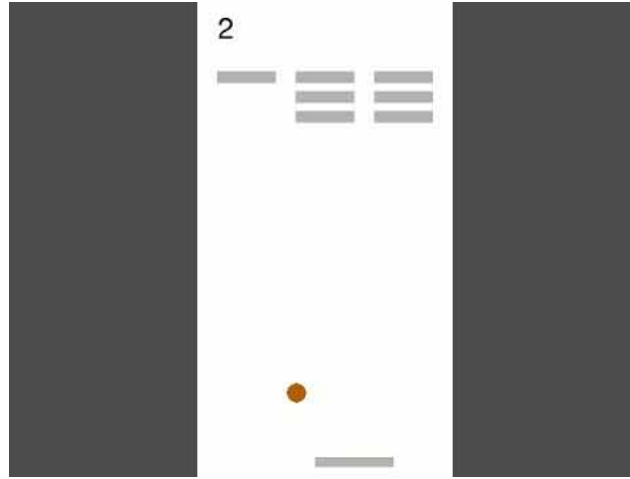
maximize the expected return. We observe that the optimal action in  $s_{13}$ , according to the policy, is not the one to take the shortest path to the goal, i.e. the *Up* action. This is because there is a small probability to end in  $s_{24}$ , the failure state, and be punished with a high negative reward. In terms of expected reward, it is better to take the longer path, at the price of collect small negative rewards, but avoiding the risk to fail miserably.

**Example 4.2 (Breakout).** Breakout is a well known arcade videogame developed by Atari. In this work we implemented a clone of the original Breakout. Figure 4.4 shows a screenshot of the game. On the screen there are a paddle at the bottom, many bricks at the top arranged in a grid layout with  $n$  rows and  $m$  columns (in the figure  $3 \times 3 = 9$  bricks), and a ball that is free to move across the screen. The ball bounces when it hits a wall, a brick or the paddle. When the ball hits a brick, that brick is broke down and is removed from the screen. The paddle (the agent) can move left, move right or do nothing. The *goal* is to remove all the bricks, while avoiding that the paddle misses the rebound of the ball (*failure*).

The relevant features are: position of the paddle  $p_x$ , position of the ball  $b_x, b_y$ , speed of the ball  $v_x, v_y$  and status of each brick (booleans)  $b_{ij}$ . This features of the system gives all the needed information to predict the next state from the current state. Hence we can build an MDP where:  $S$  is the set of all the possible values of the sequence of features  $\langle p_x, b_x, b_y, v_x, v_y, b_{11}, \dots, b_{nm} \rangle$ ,  $A = \{Right, Left, No-op\}$ , transition function  $T$  determined by the rules of the game. We give reward  $R(s, a, s') = 10$  if a particular brick in  $s'$  has been removed for the first time, plus 100 if that brick was the last (i.e. goal reached).

#### **Violation of the Markov property considering a smaller set of features:**

Notice that considering a strict subset of the set of features for  $S$  leads to violate the Markov property of  $T$ . Indeed, consider the case when we remove  $v_x$  and  $v_y$  from the set of features. In this setting, we removed the informations about the dynamics of the system. More precisely, we cannot predict, knowing only the current state, the value of the features  $b_x$  and  $b_y$  for the next step, because we do not know where the ball is going (up-left, down-right and so on). In order to correctly predict the next position of the ball, we should know whether earlier in the episode the ball was coming from the bottom or from the top. But this fact clearly shows that the Markovian assumption is violated. Similar arguments apply in the case where we remove the status of the bricks  $b_{11}, \dots, b_{nm}$ : indeed, if the ball in the next step is near to a brick, knowing about the status of the brick is determinant to predict if the ball will continue its trajectory (the case when the brick is absent) or it will



**Figure 4.4.** A screenshot from the videogame BREAKOUT

break down the brick and bounce, changing the direction of its motion (the case when the brick is present).

### 4.3 Temporal Difference Learning

*Temporal difference learning* (TD) (Sutton, 1988) refers to a class of model-free reinforcement learning methods which learn by bootstrapping from the current estimate of the value function. These methods sample from the environment, like Monte Carlo (MC) methods, and perform updates based on current estimates, like dynamic programming methods (DP) (Bellman, 1957). We do not discuss MC and DP methods here.

Q-Learning (Watkins, 1989; Watkins and Dayan, 1992) and Sarsa are such a methods. They update  $Q(s, a)$ , i.e. the estimation of  $q^*(s, a)$  at each transition  $(s, a) \rightarrow (s', r)$ . The update rule is the following:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \delta \quad (4.9)$$

where  $\delta$  is the *temporal difference*. In Sarsa, it is defined as:

$$\delta = r + \gamma Q(s', a') - Q(s, a) \quad (4.10)$$

whereas in Q-Learning:

$$\delta = r + \gamma \max_{a'} Q(s', a') - Q(s, a) \quad (4.11)$$

TD( $\lambda$ ) is an algorithm which uses *eligibility traces*. The parameter  $\lambda$  refers to the use of an eligibility trace. The algorithm generalizes MC methods and TD learning, obtained respectively by setting  $\lambda = 1$  and  $\lambda = 0$ . Intermediate values of  $\lambda$  yield methods that are often better of the extreme methods. Q-Learning and Sarsa that has been shown before can be rephrased with this new formalism as Q-Learning(0) and Sarsa(0), special cases of Watkin's Q( $\lambda$ ) and Sarsa( $\lambda$ ) respectively. In this setting, Equation 4.9 is modified as follows:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \delta e(s, a) \quad (4.12)$$

Where  $e(s, a) \in [0, 1]$ , the *eligibility of the pair*  $(s, a)$ , determines how much the temporal difference  $\delta$  should be weighted. Sarsa( $\lambda$ ) is reported in Algorithm 4.1, whereas Watkin's  $Q(\lambda)$  in Algorithm 4.2, both in the variants using *replacing eligibility traces* (see line 9 and line 10, respectively).

---

**Algorithm 4.1.** Sarsa( $\lambda$ ) (Singh and Sutton, 1996)
 

---

```

1: Initialize  $Q(s, a)$  arbitrarily and  $e(s, a) = 0$  for all  $s, a$ 
2: repeat{for each episode}
3:   initialize  $s$ 
4:   Choose  $a$  from  $s$  using policy derived from  $Q$  (e.g.  $\epsilon$ -greedy)
5:   repeat{for each step of episode}
6:     Take action  $a$ , observe reward  $r$  and new state  $s'$ 
7:     Choose  $a'$  from  $s'$  using policy derived from  $Q$ 
8:      $\delta \leftarrow r + \gamma Q(s', a') - Q(s, a)$ 
9:      $e(s, a) \leftarrow 1$  ▷ replacing traces
10:    for all  $s, a$  do
11:       $Q(s, a) \leftarrow Q(s, a) + \alpha \delta e(s, a)$ 
12:       $e(s, a) \leftarrow \gamma \lambda e(s, a)$ 
13:    end for
14:     $s \leftarrow s', a \leftarrow a'$ 
15:  until state  $s$  is terminal
16: until
```

---

## 4.4 Non-Markovian Reward Decision Process (NMRDP)

For some goals, it might be the case that the Markovian assumption of the reward function  $R$  – that reward depends only on the current state, and not on history – does not hold. Indeed, for many problems, it is not effective that the reward is limited to depend only on a single transition  $(s, a, s')$ ; instead, it might be extended to depend on *trajectories* (i.e.  $\langle s_0, a_0, \dots, s_{n-1}, a_{n-1}, s_n \rangle$ ), e.g. when we want to reward the agent for some (temporally extended) behaviors, opposed to simply reaching certain states.

This idea of rewarding behaviors has been proposed by (Bacchus et al., 1996) where they defined a new mathematical model, namely Non-Markovian Reward Decision Process (NMRDP), and showed how to construct optimal policies in this case.

In the next subsections, we give the main definitions to reason in this new setting. Then we show the solution proposed in (Bacchus et al., 1996).

### 4.4.1 Preliminaries

Now follows the definition of NMRDP, which is similar to the MDP definition given in Section 4.2.

**Definition 4.2.** A Non-Markovian Reward Decision Process (NMRDP) (Bacchus et al., 1996)  $\mathcal{N}$  is a tuple  $\langle S, A, T, \bar{R}, \gamma \rangle$  where  $S, A, T$  and  $\gamma$  are defined as in the MDP, and  $\bar{R} : S^* \rightarrow \mathbb{R}$  is the *non-Markovian reward function*, where  $S^* = \{ \langle s_0, s_1, \dots, s_n \rangle_{n \geq 0, s_i \in S} \}$  is the set of all the possible traces, i.e. projection of trajectories  $\langle s_0, a_0, \dots, s_{n-1}, a_{n-1}, s_n \rangle$



**Algorithm 4.2.** Watkin's  $Q(\lambda)$  (Watkins, 1989)

---

```

1: Initialize  $Q(s, a)$  arbitrarily and  $e(s, a) = 0$  for all  $s, a$ 
2: repeat{for each episode}
3:   initialize  $s$ 
4:   Choose  $a$  from  $s$  using policy derived from  $Q$  (e.g.  $\epsilon$ -greedy)
5:   repeat{for each step of episode}
6:     Take action  $a$ , observe reward  $r$  and new state  $s'$ 
7:     Choose  $a'$  from  $s'$  using policy derived from  $Q$  (e.g.  $\epsilon$ -greedy)
8:      $a^* \leftarrow \arg \max_a Q(s', a)$  (if  $a'$  ties for max, then  $a^* \leftarrow a'$ )
9:      $\delta \leftarrow r + \gamma Q(s', a^*) - Q(s, a)$ 
10:     $e(s, a) \leftarrow 1$  ▷ replacing traces
11:    for all  $s, a$  do
12:       $Q(s, a) \leftarrow Q(s, a) + \alpha \delta e(s, a)$ 
13:      if  $a' = a^*$  then
14:         $e(s, a) \leftarrow \gamma \lambda e(s, a)$ 
15:      else
16:         $e(s, a) \leftarrow 0$ 
17:      end if
18:       $e(s, a) \leftarrow \gamma \lambda e(s, a)$ 
19:    end for
20:     $s \leftarrow s', a \leftarrow a'$ 
21:  until state  $s$  is terminal
22: until

```

---

Given a trace  $\pi = \langle s_0, s_1, \dots, s_n \rangle$ , the *value of*  $\pi$  is:

$$v(\pi) = \sum_{i=1}^{|\pi|} \gamma^{i-1} \bar{R}(\langle s_0, s_1, \dots, s_n \rangle) \quad (4.13)$$

where  $|\pi|$  denotes the number of transitions (i.e. of actions).

The policy  $\bar{\rho}$  in this setting is defined over sequences of states, i.e.  $\bar{\rho} : S^* \rightarrow A$ . The *value of*  $\bar{\rho}$  given an initial state  $s_0$  is defined as:

$$v^{\bar{\rho}}(s) = \mathbb{E}_{\pi \sim \mathcal{N}, \bar{\rho}, s_0} [v(\pi)] \quad (4.14)$$

i.e. the expected value in state  $s$  considering the distribution of traces defined by the transition function of  $\mathcal{N}$ , the policy  $\bar{\rho}$  and the initial state  $s_0$ .

We are interested in two problems, that we will study in the next sections:

- Find an optimal (non-Markovian) policy  $\bar{\rho}$  for an NMRDP  $\mathcal{N}$  (Definition 4.2);
- Define the non-Markovian reward function for the domain of interest.

#### 4.4.2 Find an optimal policy $\bar{\rho}$ for NMRDPs

The key difficulty with non-Markovian rewards is that standard optimization techniques, most based on Bellman's (Bellman, 1957) dynamic programming principle, cannot be used. Indeed, this requires one to resort to optimization over a policy space that maps histories (rather than states) into actions, a process that would incur great computational expense. (Bacchus et al., 1996) give the definition of a decision problem *equivalent* to an NMRDP in which the rewards are Markovian. This construction is the key element to solve our problem, i.e. find an optimal policy for an NMRDP.

### Equivalent MDP

Now we give the definition of *equivalent* MDP of an NMRDP, and state an important result.

**Definition 4.3** (Bacchus et al. (1996)). An NMRDP  $\mathcal{N} = \langle S, A, T, \bar{R}, \gamma \rangle$  is *equivalent* to an extended MDP  $\mathcal{M} = \langle S', A, T', R', \gamma \rangle$  if there exist two functions  $\tau : S' \rightarrow S$  and  $\sigma : S \rightarrow S'$  such that

1.  $\forall s \in S : \tau(\sigma(s)) = s$ ;
2.  $\forall s_1, s_2 \in S$  and  $s'_1 \in S'$ : if  $T(s_1, a, s_2) > 0$  and  $\tau(s'_1) = s_1$ , there exists a unique  $s'_2 \in S'$  such that  $\tau(s'_2) = s_2$  and  $T'(s'_1, a, s'_2) = T(s_1, a, s_2)$ ;
3. For any feasible trace  $\langle s_0, s_1, \dots, s_n \rangle$  of  $\mathcal{N}$  and  $\langle s'_0, s'_1, \dots, s'_n \rangle$  of  $\mathcal{M}$  associated to the trajectories  $\langle s_0, a_0, \dots, s_{n-1}, a_{n-1}, s_n \rangle$  and  $\langle s'_0, a_0, \dots, s'_{n-1}, a_{n-1}, s'_n \rangle$ , such that  $\tau(s'_i) = s_i$  and  $\sigma(s_0) = s'_0$ , we have  $R(\langle s_0, s_1, \dots, s_n \rangle) = R'(\langle s'_0, s'_1, \dots, s'_n \rangle)$ .

Given the Definition 4.3, we give the definition of corresponding policy:

**Definition 4.4** (Bacchus et al. (1996)). Let  $\mathcal{N}$  be an NMRDP and let  $\mathcal{M}$  be the equivalent MDP as defined in Definition 4.3. Let  $\rho$  be a policy for  $\mathcal{M}$ . The *corresponding policy* for  $\mathcal{N}$  is defined as  $\bar{\rho}(\langle s_0, \dots, s_n \rangle) = \rho(s'_n)$ , where for the sequence  $\langle s'_0, \dots, s'_n \rangle$  we have  $\tau(s'_i) = s_i \forall i$  and  $\sigma(s_0) = s'_0$ .

From definitions 4.3 and 4.4, and since that for all policy  $\rho$  of  $\mathcal{M}$  the corresponding policy  $\bar{\rho}$  of  $\mathcal{N}$  is such that  $\forall s. v_\rho(s) = v_{\bar{\rho}}(\sigma(s))$ , the following theorem holds:

**Theorem 4.1** (Bacchus et al. (1996)). *Let  $\rho$  be an optimal policy for MDP  $\mathcal{M}$ . Then the corresponding policy is optimal for NMRDP  $\mathcal{N}$ .*

The Theorem 4.1 allow us to learn an optimal policy  $\bar{\rho}$  for NMRDP by learning a policy  $\rho$  over an equivalent MDP, which can be done by resorting on any off-the-shelf algorithm (e.g. see Section 4.3). Moreover, obtaining the corresponding policy for the original NMRDP is straightforward, although in practice is not needed, since it is enough to run the policy  $\rho$  over the MDP.

In other words, the problem of finding an optimal policy for an NMRDP reduces to find an optimal policy for an equivalent MDP such that Condition 1, 2 and 3 of Definition 4.3 hold.

#### 4.4.3 Define the non-Markovian reward function $\bar{R}$

To reward agents for (temporally extended) behaviors, as opposed to simply reaching certain states, we need a way to specify rewards for specific trajectories through the state space. Specifying a non-Markovian reward function explicitly is quite hard and unintuitive, impossible if we are in a infinite-horizon setting. Instead, we can define *properties* over trajectories and reward only the ones which satisfy some of them, in contrast to enumerate all the possible trajectories.

Temporal logics presented in Section 2.1 gives an effective way to do this. Indeed, in order to speak about a desired behavior, i.e. fulfillment of properties that might change over time, we can define a *formula*  $\varphi$  (or more formulas) in some suited temporal logic formalism semantically defined over trajectories  $\pi$ , speaking about a set of properties  $\mathcal{P}$  such that each state  $s \in S$  is associated to a set of propositions ( $S \subseteq 2^{\mathcal{P}}$ ). In this way, a trajectory  $\pi = \langle s_0, a_0, \dots, s_{n-1}, a_{n-1}, s_n \rangle$  is rewarded with  $r_i$  iff  $\pi \models \varphi_i$ , where  $r_i$  is the reward value associated to the fulfillment of behaviors signified by  $\varphi_i$ .

#### 4.4.4 Using PLTL

In (Bacchus et al., 1996) the temporal logic formalism is *Past Linear Temporal Logic* (PLTL), which is a past version of LTL (Section 2.1). As explained before, using the declarativeness of PLTL, is possible to specify the desired behavior (expressed in terms of the properties  $\mathcal{P}$ ) that should be satisfied by the experienced trajectories and reward only them, hence obtaining a non-Markovian reward function. More formally, given a finite set  $\Phi$  of PLTL *reward formulas*, and for each  $\phi_i \in \Phi$  a real-valued reward  $r_i$ , the *temporally extended reward function*  $\bar{R}$  is defined as:

$$\bar{R}(\langle s_0, s_1, \dots, s_n \rangle) = \sum_{\phi_i \in \Phi: \langle s_0, s_1, \dots, s_n \rangle \models \phi_i} r_i \quad (4.15)$$

In order to run the actual learning task, (Bacchus et al., 1996) proposed a transformation from the NMRDP to an equivalent MDP with the state space *expanded* which allows to label each state  $s \in S$ . The idea is that the labels should keep track in some way the (partial) satisfaction of the temporal formulas  $\phi_i \in \Phi$ . A state  $s$  in the transformed state space is replicated multiple times, marking the difference between different (relevant) histories terminating in state  $s$ .

In this way, they obtained a compact representation of the required history-dependent policy by considering only relevant history, and can produce this policy using computationally-effective MDP algorithms. In other words, the states of the NMRDP can be mapped into those of the expanded MDP, in such a way that corresponding states yield same transition probabilities and corresponding traces have same rewards.

### 4.5 NMRDP with $LTL_f/LDL_f$ rewards

In this section we explain how to specify non-Markovian rewards with  $LTL_f/LDL_f$  formulas (instead of PLTL) and how the associated MDP expansion works (Brafman et al., 2018), analogously to what we saw with PLTL (Section 4.4.4).

The temporally extended reward function  $\bar{R}$  is similar to Equation 4.15, but instead of using PLTL formula we use  $LTL_f/LDL_f$  formulas. Formally, given a set of pairs  $\{(\varphi_i, r_i)_{i=1}^m\}$  (where  $\varphi_i$  denotes the  $LTL_f/LDL_f$  formula for specifying a desired behavior, and  $r_i$  denotes the reward associated to the satisfaction of  $\varphi_i$ , and given a (partial) trace  $\pi = \langle s_0, s_1, \dots, s_n \rangle$ , we define  $\bar{R}$  as:

$$\bar{R}(\pi) = \sum_{1 \leq i \leq m: \pi \models \varphi_i} r_i \quad (4.16)$$

For the sake of clarity, in the following we use  $\{(\varphi_i, r_i)_{i=1}^m\}$  to denote  $\bar{R}$ .

Now we describe the MDP expansion for doing learning in this setting, as proposed in (Brafman et al., 2018). Without loss of generality, we assume that every NMRDP  $\mathcal{N}$  is reduced into another NMRDP  $\mathcal{N}' = \langle S', A', T', R', \gamma \rangle$ :

$$\begin{aligned}
S' &= S \cup \{s_{init}\} \\
A' &= A \cup \{start\} \\
T'(s, a, s') &= \begin{cases} 1 & \text{if } s = s_{init}, a = start, s' = s_0 \\ 0 & \text{if } s = s_{init} \text{ and } (a \neq start \text{ or } s' \neq s_0) \\ T(s, a, s') & \text{otherwise} \end{cases} \\
R'(\langle s_{init}, s_0, \dots, s_n \rangle) &= R(\langle s_0, s_1, \dots, s_n \rangle)
\end{aligned} \tag{4.17}$$

and  $s_{init}$  is the new initial state. In other words, we prefix to every feasible trajectory  $\mathcal{N}$  the pair  $\langle s_{init}, start \rangle$ , denoting the beginning of the episode. We do this for two reasons: allow to evaluate formulas in  $s_0$  and make it compliant with the most general definition of the reward, namely  $R(s, a, s')$ , also when there is no true action that is done (i.e. empty trace).

**Definition 4.5** (Brafman et al. (2018)). Given an NMRDP  $\mathcal{N} = \langle S, A, T, \{(\varphi_i, r_i)_{i=1}^m, \gamma\} \rangle$  (i.e. with non-Markovian rewards specified by  $LTL_f/LDL_f$  formulas) it is possible to build an  $\mathcal{M} = \langle S', A, T', R', \gamma \rangle$  that is *equivalent* (in the sense of Definition 4.3) to  $\mathcal{N}$ . Denoting with  $\mathcal{A}_{\varphi_i} = \langle 2^P, Q_i, q_{i0}, \delta_i, F_i \rangle$  (notice that  $S \subseteq 2^P$  and  $\delta_i$  is total) the DFA associated with  $\varphi_i$  (see Section 2.5), the equivalent MDP  $\mathcal{M}$  is built as follows:

- $S' = Q_1 \times \dots \times Q_m \times S$  is the set of states;
- $T' : S' \times A \times S' \rightarrow [0, 1]$  is defined as follows:

$$Tr'(q_1, \dots, q_m, s, a, q'_1, \dots, q'_m, s') = \begin{cases} Tr(s, a, s') & \text{if } \forall i : \delta_i(q_i, s') = q'_i \\ 0 & \text{otherwise;} \end{cases}$$

- $R' : S' \times A \times S' \rightarrow \mathbb{R}$  is defined as:

$$R'(q_1, \dots, q_m, s, a, q'_1, \dots, q'_m, s') = \sum_{i: q'_i \in F_i} r_i$$

**Theorem 4.2** (Brafman et al. (2018)). The NMRDP  $\mathcal{N} = \langle S, A, T, \{(\varphi_i, r_i)_{i=1}^m, \gamma\} \rangle$  is equivalent to the MDP  $\mathcal{M} = \langle S', A, T', R', \gamma \rangle$  defined in Definition 4.5.

*Proof.* Recall that every  $s' \in S'$  has the form  $(q_1, \dots, q_m, s)$ . Define  $\tau(q_1, \dots, q_m, s) = s$ . Define  $\sigma(s) = (q_{10}, \dots, q_{m0}, s)$ . We have  $\tau(\sigma(s)) = s$ , hence Condition 1 is verified. Condition 2 of Definition 4.3 is easily verifiable by inspection. For Condition 3, consider a possible trace  $\pi = \langle s_0, s_1, \dots, s_n \rangle$ . We use  $\sigma$  to obtain  $s'_0 = \sigma(s_0)$  and given  $s_i$ , we define  $s'_i$  (for  $1 \leq i \leq n$ ) to be the unique state  $(q_{1,i}, \dots, q_{m,i}, s_i)$  such that  $q_{j,i} = \delta(q_{j,i-1}, s_i)$  for all  $1 \leq j \leq m$ . Moreover, we require that, without loss of generality, every trajectory in the new MDP starts from  $s_{init}$  and now have a corresponding possible trace of  $\mathcal{M}$ , i.e.,  $\pi = \langle s'_0, s'_1, \dots, s'_n \rangle$ . This is the only feasible trajectory of  $\mathcal{M}$  that satisfies Condition 3. The reward at  $\pi = \langle s_0, s_1, \dots, s_n \rangle$  depends only on whether or not each formula  $\varphi_i$  is satisfied by  $\pi$ . However, by construction of the automaton  $\mathcal{A}_{\varphi_i}$  and the transition function  $T$ ,  $\pi \models \varphi_i$  iff  $s'_n = (q_1, \dots, q_m, s_n)$  and  $q_i \in F_i$   $\square$

Let  $\rho'$  be a (Markovian) policy for  $\mathcal{M}$ . It is easy to define an *corresponding* policy on  $\mathcal{N}$ , i.e., a policy that guarantees the same rewards, by using  $\tau$  and  $\sigma$  mappings defined in Theorem 4.2 and the result shown in Theorem 4.4.

Obviously, typical learning techniques, such as Q-learning or Sarsa, are applicable on the expanded  $\mathcal{M}$  and so we can learn an optimal policy  $\rho$  for  $\mathcal{M}$ . Thus, an optimal policy for  $\mathcal{N}$  can be learnt on  $\mathcal{M}$ . Of course, none of these structures is (completely) known to the learning agent, and the above transformation is never done explicitly. Rather, the agent carries out the learning process by assuming that the underlying model is  $\mathcal{M}$  instead of  $\mathcal{N}$  (applying the fix introduced in Definition 4.17).

Observe that the state space of  $\mathcal{M}'$  is the product of the state spaces of  $\mathcal{N}$  and  $\mathcal{A}_{\varphi_i}$ , and that the reward  $R'$  is Markovian. In other words, the (stateful) structure of the LTL<sub>f</sub>/LDL<sub>f</sub> formulas  $\varphi_i$  used in the (non-Markovian) reward of  $\mathcal{N}$  is *compiled* into the states of  $\mathcal{M}$ .

### Why should we use LDL<sub>f</sub>

LDL<sub>f</sub> formalism (introduced in Section 2.4) has the advantage of *enhanced expressive power* over other proposals, as discussed in (Brafman et al., 2018). Indeed, we move from linear-time temporal logics to LDL<sub>f</sub>, paying no additional (worst-case) complexity costs. LDL<sub>f</sub> can encode in polynomial time LTL<sub>f</sub>, regular expressions (RE) and the past LTL (PLTL) of (Bacchus et al., 1996). Moreover, LDL<sub>f</sub> can naturally represent "procedural constraints" (Baier et al., 2008), i.e., sequencing constraints expressed as programs, using "if" and "while", hence allowing to express more complex properties.

## 4.6 RL for LTL<sub>f</sub>/LDL<sub>f</sub> Goals

In this section we define a particular problem and propose a solution, which is the main theoretical contribution of this work. We call this problem *Reinforcement Learning for LTL<sub>f</sub>/LDL<sub>f</sub> Goals*.

### 4.6.1 Problem definition

#### the World, the Agent and the Fluents

Let  $\mathcal{W}$  be a *world* of interest (e.g. a room, an environment, a videogame). Let  $W$  be the set of *world states*, i.e. the states of the world  $\mathcal{W}$ . A *feature* is a function  $f_j$  that maps a world state to the values of another domain  $D_j$ , such as reals, finite enumerations, booleans, etc., i.e.,  $f_j : W \rightarrow D_j$ . Given a set of features  $F = \langle f_1, \dots, f_d \rangle$ , the *feature vector* of a world state  $w_h$  is the vector  $\mathbf{f}(w_h) = \langle f_1(w_h), \dots, f_d(w_h) \rangle$  of feature values corresponding to  $w_h$ .

Now consider an agent that lives in  $\mathcal{W}$ . The agent can interact with  $\mathcal{W}$  by executing an action  $a$  taken from a *set of actions*  $A$ . Without loss of generality, we assume that such learning agent has a special action *stop* which deems the end of an episode. Moreover, the agent has its own set of features  $F_{ag} = \langle f_1, \dots, f_d \rangle$ , which yields its *representation of the world*  $S$ , where  $S \subseteq F_{ag}(W)$  and  $F_{ag}(W) = \{\mathbf{f}_{ag}(w) | w \in W\}$ . We assume that the agent has a *clock* which determines the granularity of its acting. At every clock, the agent can do action  $a$  and observe both the new state  $s$  from the new world state  $w'$ , namely  $s = \mathbf{f}_{ag}(w)$ , and a real-valued reward  $R(s, a, s')$ , that depends only from the transition  $s \rightarrow_a s'$ . Finally, we assume that the law that determines the possible next state  $s'$  given the history

of a trajectory (i.e. the state transition function  $T$ ) is Markovian, hence it depends only from the current state  $s$  and the taken action  $a$ . In other words, we can define an MDP  $\mathcal{M}_{ag} = \langle S, A, T, R, \gamma \rangle$ .

We consider arbitrary LTL<sub>f</sub>/LDL<sub>f</sub> formulas  $\varphi_i$  ( $i = 1, \dots, m$ ) over a set of fluents  $\mathcal{F}$  used for provide a high-level description of the world. We denote by  $\mathcal{L} = 2^{\mathcal{F}}$  the set of possible fluents configurations. Given a set of feature  $F_{goal}$ , a *configuration of fluents*  $\ell_h \in \mathcal{L}$  is formed by the components that assign truth values to the fluents according to the feature vector  $\mathbf{f}_{goal}(w_h)$ . At every step, the features for fluents evaluations are observed, obtaining a particular configuration  $\ell \in \mathcal{L}$ . Notice that in general the features for the fluents and for the agent state space may differ. The formula  $\varphi_i$  is selecting sequences of fluents configurations  $\ell_1, \dots, \ell_n$ , with  $\ell_k \in \mathcal{L}$ , whose relationship with the sequences of states  $s_1, \dots, s_n$ , with  $s_k \in S$  is unknown.

In other words, a subset of features are used to describe agent states  $s_h$  and another subset (for simplicity, assumed disjoint from the previous one) are used to evaluate the fluents in  $\ell_h$ . Hence, given a sequence  $w_1, \dots, w_n$  of world states we get the corresponding sequence of sequences learning agent states  $s_1, \dots, s_n$  and simultaneously the sequence of fluent configurations  $\ell_1, \dots, \ell_n$ . Notice that we do not have a formalization for  $w_1, \dots, w_n$  but we do have that for  $s_1, \dots, s_n$  and for  $\ell_1, \dots, \ell_n$ .

Oversimplifying, we may say that  $S$  is the set of configurations of the low-level features for the learning agent, while  $\mathcal{L}$  is the set of configuration of the high-level features needed for expressing  $\varphi_i$ .

### Markovian assumption of the state transition function

Now we make the following assumption: that is, the agent actions in  $A$  induce a transition distribution over the features and fluents configuration, i.e.,

$$T_{ag}^g : S \times \mathcal{L} \times A \rightarrow \text{Prob}(S \times \mathcal{L}) \quad (4.18)$$

This means that the state transition function  $T_{ag}^g$  is *Markovian*, i.e. the probability to end in the next state  $s'$  with the next fluents configuration  $\ell'$  depends only from  $s, \ell$  and  $a$  (the current agent state, the current fluents configuration and the action taken, respectively).

Such a transition distribution together with the initial values of the fluents  $\ell_0$  and of the agent state  $s_0$  allow us to describe a probabilistic transition system accounting for the dynamics of the fluents and agent states. In other words, in response to an agent action  $a_h$  performed in the current state  $w_h$  (in the state  $s_h$  of the agent and the configuration  $\ell_h$  of the fluents), the world changes into  $w_{h+1}$  from which  $s_{h+1}$  and  $\ell_{h+1}$ . This is all we need to proceed.

Notice that we do not assume nothing about the primitive state transitions function over fluents, i.e. the one induced over the set of possible fluents configurations  $\mathcal{L}$ . However, it might be the case that the Markovian assumption of the agent state transition function holds, i.e. the probability distribution for the next state  $\ell'$  is fully determined by the current state  $\ell$  and the action  $a$ . In that case, the considerations presented in this section still holds, since it is a special case of our general assumptions. We only require that the *joint* transition function in Equation 4.18 satisfies the Markov property.

We are interested in devising policies for the learning agent such that at the end of the episode, i.e., when the agent executes *stop*, the LTL<sub>f</sub>/LDL<sub>f</sub> goal formulas  $\varphi_i$  ( $i = 1, \dots, m$ ) are satisfied. Now we can state our problem formally.

**Definition 4.6.** We define RL for LTL<sub>f</sub>/LDL<sub>f</sub> goals, denoted as

$$\mathcal{M}_{ag}^{goal} = \langle S, A, R, \mathcal{L}, T_{ag}^g, \{(\varphi_i, r_i)\}_{i=1}^m \rangle$$

with  $T_{ag}^g$ ,  $R$  and  $r_i$  unknown, the following problem: given a learning agent  $\mathcal{M}_{ag} = \langle S, A, T, R \rangle$ , with  $T$  and  $R$  unknown and a set  $\{(\varphi_i, r_i)\}_{i=1}^m$  of LTL<sub>f</sub>/LDL<sub>f</sub> formulas with associated rewards, find a (non-Markovian) policy  $\bar{\rho} : S^* \rightarrow A$  that is optimal wrt the sum of the rewards  $r_i$  and  $R$ .

Observe that an optimal policy for our problem, although not depending on  $\mathcal{L}$ , is guaranteed to satisfy the LTL<sub>f</sub>/LDL<sub>f</sub> goal formulas.

### 4.6.2 Examples

—TODO

### 4.6.3 Reduction to MDP

To devise a solution technique, we start by transforming  $\mathcal{M}_{ag}^{goal} = \langle S, A, T_{ag}^g, R, \mathcal{L}, \{(\varphi_i, r_i)\}_{i=1}^m \rangle$  into an NMRDP  $\mathcal{M}_{ag}^{nmr} = \langle S \times \mathcal{L}, A, T_{ag}^g, \{(\varphi'_i, r_i)\}_{i=1}^m \cup \{(\varphi_s, R(s, a, s'))\}_{s \in S, a \in A, s' \in S} \rangle$  where:

- States are pairs  $(s, \ell)$  formed by an agent configuration  $s$  and a fluents configuration  $\ell$ .
- $\varphi'_i = \varphi_i \wedge \Diamond Done$ .
- $\varphi_s = \Diamond(s \wedge a \wedge \bigcirc(Last \wedge s'))$ .
- $T_{ag}^g$ ,  $r_i$  and  $R(s, a, s')$  are unknown and sampled from the environment.

Formulas  $\varphi'_i$  simply require to evaluate the corresponding goal formula  $\varphi_i$  after having done the action *stop*, which sets the fluent *Done* to true and ends the episode. Hence it gives the reward associated to the goal at the end of the episode. The formulas  $\Diamond(s \wedge a \wedge \bigcirc(Last \wedge s'))$ , one per  $(s, a, s')$ , requires both states  $s$  and action  $a$  are followed by  $s'$  are evaluated at the end of the current (partial) trace (notice the use of *Last*). In this case, the reward  $R(s, a, s')$  from  $\mathcal{M}_{ag}$  associated with  $(s, a, s')$  is given.

Notice that policies for  $\mathcal{M}_{ag}^{nmr}$  have the form  $(S \times \mathcal{L})^* \rightarrow A$  which needs to be restricted to have the form required by our problem  $\mathcal{M}_{ag}^{goal}$ .

A policy  $\bar{\rho} : (S \times \mathcal{L})^* \rightarrow A$  has the form  $S^* \rightarrow A$  when for any sequence of  $n$  states  $\langle s_1 \cdots s_n \rangle$ , we have that for any pair of sequences of fluent configurations  $\langle \ell'_1 \cdots \ell'_n \rangle$ ,  $\langle \ell''_1 \cdots \ell''_n \rangle$  the policy returns the same action,  $\bar{\rho}(\langle s_1, \ell'_1 \rangle \cdots \langle s_n, \ell'_n \rangle) = \bar{\rho}(\langle s_1, \ell''_1 \rangle \cdots \langle s_n, \ell''_n \rangle)$ . In other words, a policy  $\bar{\rho} : (S \times \mathcal{L})^* \rightarrow A$  has the form  $\bar{\rho} : S^* \rightarrow A$  when it does not depend on the fluents  $\mathcal{L}$ . We can now state the following result.

**Theorem 4.3.** RL for LTL<sub>f</sub>/LDL<sub>f</sub> goals  $\mathcal{M}_{ag}^{goal} = \langle S, A, T_{ag}^g, R, \mathcal{L}, \{(\varphi_i, r_i)\}_{i=1}^m \rangle$  can be reduced to RL over the NMRDP  $\mathcal{M}_{ag}^{nmr} = \langle S \times \mathcal{L}, A, T_{ag}^g, \{(\varphi'_i, r_i)\}_{i=1}^m \cup \{(\varphi_s, R(s, a, s'))\}_{s \in S, a \in A, s' \in S} \rangle$ , restricting policies to be learned to have the form  $S^* \rightarrow A$ .



Observe that by restricting  $\mathcal{M}_{ag}^{nmr}$  policies to  $S^*$  in general we may discard policies that have a better reward but depend on  $\mathcal{L}$ . On the other hand, these policies need to change the learning agent in order to allow it to observe  $\mathcal{L}$  as well. As mentioned in the introduction, we are interested in keeping the learning agent as it is, apart for additional memory.

As a second step, we apply the construction of Section 4.5 and obtain a new MDP learning agent. In such construction, however, because of the triviality of their automata, we do not need to keep track of state  $\varphi_s$ , but just give the reward  $R(s, a, s')$  associated to  $(s, a, s')$ . Instead we do need to keep track of state of the DFAS  $\mathcal{A}_{\varphi_i}$  corresponding to the formulas  $\varphi'_i$ . Hence, from  $\mathcal{M}_{ag}^{nmr}$ , we get an MDP  $\mathcal{M}'_{ag} = \langle S', A', Tr'_{ag}, R' \rangle$  where:

- $S' = Q_1 \times \dots \times Q_m \times S \times \mathcal{L}$  is the set of states;
- $Tr'_{ag} : S' \times A' \times S' \rightarrow [0, 1]$  is defined as follows:

$$Tr'_{ag}(q_1, \dots, q_m, s, \ell, a, q'_1, \dots, q'_m, s', \ell') = \begin{cases} Tr(s, \ell, a, s', \ell') & \text{if } \forall i : \delta_i(q_i, \ell') = q'_i \\ 0 & \text{otherwise;} \end{cases}$$

- $R' : S' \times A \times S' \rightarrow \mathbb{R}$  is defined as:

$$R'(q_1, \dots, q_m, s, \ell, a, q'_1, \dots, q'_m, s', \ell') = \sum_{i: q'_i \in F_i} r_i + R(s, a, s')$$

Finally we observe that the environment gives now both the rewards  $R(s, a, s')$  of the original learning agent, and the rewards  $r_i$  associated to the formula so has to guide the agent towards the satisfaction of the goal (progressing correctly the DFAS  $\mathcal{A}_{\varphi_i}$ ).

By applying Theorem 4.2 we get that NMRDP  $\mathcal{M}_{ag}^{nmr}$  and the MDP  $\mathcal{M}'_{ag}$  are equivalent, i.e., any policy of  $\mathcal{M}_{ag}^{nmr}$  has an equivalent policy (hence guaranteeing the same reward) in  $\mathcal{M}'_{ag}$  and vice versa. Hence we can learn policy on  $\mathcal{M}'_{ag}$  instead of  $\mathcal{M}_{ag}^{nmr}$ .

We can refine Theorem 4.3 into the following one.

**Theorem 4.4.** *RL for  $LTL_f/ LDL_f$  goals  $\mathcal{M}_{ag}^{goal} = \langle S, A, T_{ag}^g, R, \mathcal{L}, \{(\varphi_i, r_i)\}_{i=1}^m \rangle$  can be reduced to RL over the MDP  $\mathcal{M}'_{ag} = \langle S', A, T'_{ag}, R' \rangle$ , restricting policies to be learned to have the form  $Q_1 \times \dots \times Q_n \times S \rightarrow A$ .*

As before, a policy  $Q_1 \times \dots \times Q_n \times S \times \mathcal{L} \rightarrow A$  has the form  $Q_1 \times \dots \times Q_n \times S \rightarrow A$  when any  $\ell$  and  $\ell'$  the policy returns the same action,  $\rho(q_1, \dots, q_n s, \ell) = \rho(q_1, \dots, q_n s, \ell')$ .

The final step is to solve our original RL task on  $\mathcal{M}_{ag}^{goal}$  by performing RL on a new MDP  $\mathcal{M}_{ag}^{new} = \langle Q_1 \times \dots \times Q_m \times S, A, T''_{ag}, R'' \rangle$  where:

- Transitions distribution  $T''_{ag}$  is the marginalization wrt  $\mathcal{L}$  of  $T'_{ag}$  and is unknown;
- Rewards  $R''$  is defined as:

$$R''(q_1, \dots, q_m, s, a, q'_1, \dots, q'_m, s') = \sum_{i: q'_i \in F_i} r_i + R(s, a, s').$$



- States  $q_i$  of DFAs  $\mathcal{A}_{\varphi_i}$  are progressed correctly by the environment.

Indeed we can show the following result.

**Theorem 4.5.** *RL for  $LTL_f/SDL_f$  goals  $\mathcal{M}_{ag}^{goal} = \langle S, A, T', R, \mathcal{L}, \{(\varphi_i, r_i)\}_{i=1}^m \rangle$  can be reduced to RL over the MDP  $\mathcal{M}_{ag}^{new} = \langle Q_1 \times \dots \times Q_m \times S, A, T''_{ag}, R'' \rangle$  and the optimal policy  $\rho_{ag}^{new}$  learned for  $\mathcal{M}_{ag}^{new}$  can be reduced to a corresponding optimal policy for  $\mathcal{M}_{ag}^{goal}$ .*

*Proof.* From Theorem 4.4, by the following observations. For the sake of brevity, we use  $\mathbf{q}$  to denote  $q_1, \dots, q_m$ . Notice also that for all  $\ell, \ell' \in \mathcal{L}$ ,  $R'(\mathbf{q}, s, \ell, a, \mathbf{q}', s', \ell') = R''(\mathbf{q}, s, a, \mathbf{q}', s')$ .

We show that the values of  $v_{ag}^\rho(q_1, \dots, q_m, s, \ell)$ , i.e. the state value function for  $\mathcal{M}'_{ag}$  (for simplicity  $v^\rho$ , unless otherwise stated), for some policy  $\rho$ , do not depend on  $\ell$  or, in other words, it is necessary that  $\forall \ell_1, \ell_2. v^\rho(q_1, \dots, q_m, s, \ell_1) = v^\rho(q_1, \dots, q_m, s, \ell_2)$ . Finally, we notice that  $\forall \ell. v_{ag}^{\rho, new} = v_{ag}^\rho$

From Equation 4.4 we have:

$$\begin{aligned} v_\rho(\mathbf{q}, s, \ell) = & \sum_{\mathbf{q}', s', \ell'} P(\mathbf{q}', s', \ell' | \mathbf{q}, s, \ell, a) [R'(\mathbf{q}, s, \ell, a, \mathbf{q}', s', \ell') + \gamma v_\rho(\mathbf{q}', s', \ell')] = \\ & \sum_{\mathbf{q}', s', \ell'} P(\mathbf{q}', s', \ell' | \mathbf{q}, s, \ell, a) [R''(\mathbf{q}, s, a, \mathbf{q}', s') + \gamma v_\rho(\mathbf{q}', s', \ell')] \end{aligned} \quad (4.19)$$

Using the equivalence between  $R'$  and  $R''$ , as already pointed out. Notice that we can compute  $\mathbf{q}'$  from  $\mathbf{q}$  and  $\ell'$ , hence we do not need  $\ell$ . In other words:

$$P(\mathbf{q}', s', \ell' | \mathbf{q}, s, \ell, a) = P(\mathbf{q}', s', \ell' | \mathbf{q}, s, a)$$

Equation 4.19 becomes:

$$\sum_{\mathbf{q}', s', \ell'} P(\mathbf{q}', s', \ell' | \mathbf{q}, s, a) [R''(\mathbf{q}, s, a, \mathbf{q}', s') + \gamma v_\rho(\mathbf{q}', s', \ell')] \quad (4.20)$$

At this point, we see that  $v^\rho$  does not depend from  $\ell$ , hence we can safely drop  $\ell$  as argument for  $v_\rho$ , obtaining  $v_{ag}^\rho$ . Indeed, from 4.20:

$$\begin{aligned} \sum_{\mathbf{q}', s'} [R''(\mathbf{q}, s, a, \mathbf{q}', s') + \gamma v_{ag}^{\rho, new}(\mathbf{q}, s)] \sum_{\ell'} P(\mathbf{q}', s', \ell' | \mathbf{q}, s, a) = & \quad (4.21) \\ \sum_{\mathbf{q}', s'} P(\mathbf{q}', s' | \mathbf{q}, s, a) [R''(\mathbf{q}, s, a, \mathbf{q}', s') + \gamma v_{ag}^\rho(\mathbf{q}', s')] = & \\ v_{ag}^\rho(\mathbf{q}, s) & \end{aligned}$$

Where in 4.21 we marginalized the distribution  $P(\mathbf{q}', s', \ell' | \mathbf{q}, s, a)$  over  $\ell'$ . From Definition 4.1 of optimal policy, we can reduce an optimal policy  $\rho_{ag}^{new}$  to a policy of the form  $\rho'_{ag} : Q_1 \times \dots \times Q_m \times S \rightarrow A$  that is optimal for  $\mathcal{M}'_{ag}$  (since the state value function of  $\mathcal{M}_{ag}^{new}$ , after dropping the argument  $\ell$ , and of  $\mathcal{M}'_{ag}$  are equivalent). From Theorem 4.4 the thesis.  $\square$

It is worth to remark that in the resulting MDP  $\mathcal{M}_{ag}^{new}$  the *explicit presence of the fluents configuration  $\ell$  has been removed*. Rather, the dependency is compiled into the expanded state space  $Q_1 \times \dots \times Q_m \times S$ , where  $Q_1, \dots, Q_m$  are the automata state spaces associated to the formulas  $\varphi_i$ .

#### 4.6.4 An episodic goal-based view

Here we clarify how the actual transition model works in the final MDP  $\mathcal{M}_{ag}^{new}$ . We focus on the episodic view, where the learning process is organized in episodes. Recall that the state space of  $\mathcal{M}_{ag}^{new}$  is  $Q_1 \times \dots \times Q_m \times S$ , where  $Q_i$  is the set of states of  $\mathcal{A}_{\varphi_i}$ , the automaton associated to the  $LTL_f/ LDL_f$  formula  $\varphi_i$  (see Section 2.5). Moreover we consider two cases: when there exists a subset of states  $S_{goal} \subseteq S$  that we call *goal states*, such that every transition in those state make the task completed and hence the end of the episode; and when there are no goal states, but the task is to maximize the obtained reward, until a maximum number of steps is reached. E.g. in Gridworld presented in Example 4.1, we can defined  $s_{34}$  as a goal state. However nothing prevent us to set a maximum number of time steps, and let the agent learn how to collect reward as much as possible in a limited amount of time.

Now we give the following definition:

**Definition 4.7.** Let  $\mathcal{M}_{ag}^{goal}$  be our problem and  $\mathcal{M}_{ag}^{new}$  its transformation into an MDP, as defined in Theorem 4.5 and let  $\mathbf{s} = \langle q_1, \dots, q_m, s \rangle \in Q_1 \times Q_1 \times \dots \times Q_m \times S$ . Given an observation  $s' \in S$  and  $\ell' \in \mathcal{L}$ , we define the *successor state* of  $\mathbf{s}$  as  $\mathbf{s}' = \langle q'_1, \dots, q'_m, s' \rangle$ , where  $q'_i = \delta_i(q_i, \ell')$ .

A reinforcement learning episode in this setting works as follows. Assume that the MDP has a dummy initial state  $s_{init}$  and a dummy action  $start$ , as defined for NMRDPs in 4.17. In this construction, the dummy initial state in the expanded state space is  $(q_{0,0}, q_{1,0}, \dots, q_{m,0}, s_{init})$ . The first action taken by the agent is  $start$ , which allow the agent to observe the true initial world state  $w_0$ . From  $w_0$ , the agent extract the features to determine the state  $s_0 \in S$  and the fluents configuration  $\ell_0 \in \mathcal{L}$ <sup>1</sup>. The successor state is  $(q'_0, q'_1, \dots, q'_m, s_0)$ . Then the agent might take a new action, observe another world state  $w'$ , extract  $s' \in S$  and  $\ell' \in \mathcal{L}$ , and compute the  $q_i$  as before. The sequence reiterates until the end of the episode.

Consider a generic state  $\mathbf{s} = \langle q_1, \dots, q_m, s \rangle \in Q_1 \times Q_1 \times \dots \times Q_m \times S$  and a transition to  $\mathbf{s}' = \langle q'_1, \dots, q'_m, s' \rangle$ . For each new state  $q'_i$  of automaton  $\mathcal{A}_{\varphi_i}$ , the following might happen:

1.  $q'_i = q_i$ , i.e. the state is not changed;
2.  $q'_i \neq q_i$ , and from  $q'_i$  it is still possible to reach a final state;
3.  $q'_i \neq q_i$ , and from  $q'_i$  any final state cannot be reached;

If, for some  $\mathcal{A}_{\varphi_i}$ , we are in case 3, then the constraint specified by  $\varphi_i$  is violated, hence we call  $\mathbf{s}'$  a *failure state*. We say that  $\mathbf{s}$  is a *goal state* if  $\forall i. q_i \in F_i$ , i.e. every current state  $q_i$  is in an accepting state of the automaton  $\mathcal{A}_{\varphi_i}$ . Instead, if the underlying MDP is a goal-based one (e.g.  $S_{goal} \neq \emptyset$ ) then  $\mathbf{s}$  is a goal state if  $s \in S_{goal}$ .

Let  $\mathbf{s}$  the last state of the trace of the current episode. The *stopping condition*, i.e. the condition that determines the end of the episode, depending from  $\mathbf{s}$ , is:

$$\text{failure\_state}(\mathbf{s}) \vee \text{goal\_state}(\mathbf{s}) \vee \text{exceeded\_time\_limit}$$

The reward  $R(s, a, s')$  is collected after each taken action, and it is summed with  $r_i$  for every satisfied  $\varphi_i$  at the last state of the episode.

<sup>1</sup>Observe that the first transition is "artificial", and in many cases, both  $s_0$  and  $\ell_0$  are known to the experimenter. The explanation presented here is aimed to clarify the underlying mathematical construction.

**Remark about the stopping condition:** notice that the presence of  $\text{failure\_state}(\mathbf{s})$  in the stopping condition is not strictly needed. Indeed, in the general case, the agent still tries to learn an optimal policy in terms of observed rewards, regardless if the temporal goal at some point of the episode cannot be satisfied anymore (i.e. a failure state for some automaton  $\mathcal{A}_{\varphi_i}$  is reached). However, we can think of  $\text{failure\_state}(\mathbf{s})$  as a way to avoid parts of simulations that are not of interest, since the trajectory is "compromised" in terms of satisfaction of  $\text{LTL}_f/\text{LDL}_f$  formulas.

### Summary

In the following we summarize the results of this section:

- We defined a new problem: *Reinforcement Learning for  $\text{LTL}_f/\text{LDL}_f$  Goals*  $\mathcal{M}_{ag}^{goal}$  (Definition 4.6). In a nutshell, from an existing MDP we introduced temporal goals defined by  $\text{LTL}_f/\text{LDL}_f$  formulas about fluents  $\mathcal{L}$  observed from the world. A solution for this problem is a (non-Markovian) policy that is optimal in terms of rewards and such that satisfies the  $\text{LTL}_f/\text{LDL}_f$  specifications.
- We gave an equivalent formulation of  $\mathcal{M}_{ag}^{goal}$ , the NMRDP  $\mathcal{M}_{ag}^{nmr}$ , and observe that the original problem can be reduced to this formulation (Theorem 4.3).
- We applied the construction shown in Section 4.5, yielding  $\mathcal{M}'_{ag}$ , and by using Theorem 4.2 we stated Theorem 4.4, showing that  $\mathcal{M}_{ag}^{goal}$  can be reduced to  $\mathcal{M}'_{ag}$ .
- Finally, by proving Theorem 4.5, we showed that we can do reinforcement learning over the equivalent MDP  $\mathcal{M}_{ag}^{new}$  by simply dropping the fluents  $\ell$  from the state space. The optimal policy for  $\mathcal{M}_{ag}^{new}$  can be transformed to a solution for the original problem  $\mathcal{M}_{ag}^{goal}$ .

## 4.7 Conclusions

In this chapter we introduced the topic of Reinforcement Learning, as well as foundational definitions (MDP, policy  $\rho$ , state-value function  $v_\rho$ ) and algorithms (Q-Learning, Sarsa). Then we presented the notion of NMRDP and a technique to build an equivalent MDP, by specifying non-Markovian reward function with  $\text{LTL}_f/\text{LDL}_f$  formalisms. Finally, we defined and studied a new problem, *Reinforcement Learning for  $\text{LTL}_f/\text{LDL}_f$  goals*, and proposed a reduction that yields an equivalent MDP which can be used to solve the original problem. An episodic goal-based view of the final MDP is provided.

## Chapter 5

# Automata-based Reward shaping

In this chapter we discuss a method to improve exploration of the state space and improve the convergence rate in the setting studied in Section 4.6, in particular in the construction  $\mathcal{M}_{ag}^{new}$ . Indeed, the state space of the original MDP  $\mathcal{M}_{ag}$  is expanded in order to implicitly label the states with relevant histories for the satisfaction of LTL<sub>f</sub>/LDL<sub>f</sub> formulas, which in general makes harder to learn an optimal policy in  $\mathcal{M}_{ag}^{new}$  wrt  $\mathcal{M}_{ag}$ , due to a bigger state space. Moreover, the introduction of temporal goals makes things harder, because the agent has to find a proper behavior that satisfies all the goals. In general, proper behaviors are harder to find, and require more exploration of the state space.

*Reward Shaping* is a general method, well-known in the literature of Reinforcement Learning, used to deal with big state spaces and *sparse* rewards, and trying to address the temporal credit assignment problem, i.e. to determine the long-term consequences of actions. It consists in provide additional reward to the learning agent. In this chapter we propose a technique to apply reward shaping in our setting.

The chapter is structured as follows: in the first section we explain the reward shaping theory, in particular the requirements for theoretical guarantees of policy invariance under reward transformation. Then we show how apply reward shaping on the automata transitions associated to LTL<sub>f</sub>/LDL<sub>f</sub> formulas  $\varphi_i$ , both in *off-line* variant (i.e. when the automaton is built *before* the beginning of the learning process) and *on-the-fly* variant (i.e. when the automaton is built *during* the learning process).

## 5.1 Reward Shaping Theory

Reward shaping is a well-known technique to guide the agent during the learning process and so reduce the time needed to learn. The idea is to supply additional rewards in a proper manner such that the optimal policy is the same of the original MDP.

More formally, consider as example the temporal difference in SARSA after a transition  $s \rightarrow_a s'$ , presented in Equation 4.10:

$$\delta = R(s, a, s') + \gamma Q(s', a') - Q(s, a) \quad (5.1)$$

Reward shaping consists in defining the *shaping function*  $F(s, a, s')$  and sum it

to the environment reward  $R(s, a, s')$ , namely:

$$\delta = R(s, a, s') + F(s, a, s') + \gamma Q(s', a') - Q(s, a) \quad (5.2)$$

In the following sections we will discuss two way to define  $F(s, a, s')$ .

### 5.1.1 Potential-Based Reward Shaping

We give the definition of *potential-based shaping function* (PBRs).

**Definition 5.1** (Ng et al. (1999)). Let any  $S, A, \gamma$  and any shaping function  $F : S \times A \rightarrow \mathbb{R}$  be given. We say  $F$  is a potential-based shaping function if there exists a real-valued function  $\Phi : S \rightarrow \mathbb{R}$  such that for all  $s \in S, a \in A, s' \in S$

$$F(s, a, s') = \gamma \Phi(s') - \Phi(s) \quad (5.3)$$

Notice that in Equation 5.1 the action does not affect the value of  $F(s, a, s')$ , hence sometime we write  $F(s, s')$ .

In (Ng et al., 1999) it has been shown the following theorem:

**Theorem 5.1** (Ng et al. (1999)). Given an MDP  $\mathcal{M} = \langle S, A, T, R, \gamma \rangle$  and a potential based shaping function  $F$  (Definition 5.1). Then, consider the MDP  $\mathcal{M}' = \langle S, A, T, R+F, \gamma \rangle$ , i.e. the same of  $\mathcal{M}$  but applying reward shaping. Then, the fact that  $F$  is a potential-based reward shaping function is a necessary and sufficient condition to guarantee consistency with the optimal policy. In particular:

- (Sufficiency) if  $F$  is a potential-based shaping function, then every optimal policy in  $\mathcal{M}'$  is optimal in  $\mathcal{M}$ .
- (Necessity) if  $F$  is not a potential-based shaping function (e.g. no such  $\Phi$  exists satisfying 5.3), then there exists  $T$  and  $R$  such that no optimal policy in  $\mathcal{M}'$  is optimal in  $\mathcal{M}$ .

In poor words, potential-based reward shaping of the form  $F(s, a, s') = \gamma \Phi(s') - \Phi(s)$ , for some  $\Phi : S \rightarrow \mathbb{R}$ , is a necessary and sufficient condition for policy invariance under this kind of reward transformation, i.e. the optimal and near-optimal solutions of  $\mathcal{M}$  are preserved when considering  $\mathcal{M}'$ .

### 5.1.2 Dynamic Potential-Based Reward Shaping

A limitation of PBRs is that the potential of a state does not change dynamically during the learning. This assumption often is broken, especially if the reward-shaping function is generated automatically.

Equation 5.3 can be extended to include also the time as parameter of the potential function  $\Phi$ , while guaranteeing policy invariance. Formally:

$$F(s, t, a, s', t') = \gamma \Phi(s', t') - \Phi(s, t) \quad (5.4)$$

where  $t$  and  $t'$  are respectively the time when visiting  $s$  and  $s'$ . In this case, we call this technique *dynamic potential-based reward shaping* (DPBRs).

The shaping function in the form 5.4 guarantees policy invariance, as shown in Theorem 5.1. To show why this is the case, consider the expected discounted return for an infinite sequence of states (Definition 4.1):

$$G = \sum_{k=0}^{\infty} \gamma^k R_k \quad (5.5)$$

If we apply dynamic potential-based reward shaping to Equation 5.5 we have:

$$\begin{aligned}
G_\Phi &= \sum_{k=0}^{\infty} \gamma^k (R_k + F(s_k, t_k, s_{k+1}, t_{k+1})) \\
&= \sum_{k=0}^{\infty} \gamma^k (R_k + \gamma \Phi(s_{k+1}, t_{k+1}) - \Phi(s_k, t_k)) \\
&= \sum_{k=0}^{\infty} \gamma^k R_k + \sum_{k=0}^{\infty} \gamma \Phi(s_{k+1}, t_{k+1}) - \sum_{k=0}^{\infty} \Phi(s_k, t_k) \\
&= G + \sum_{k=1}^{\infty} \gamma \Phi(s_k, t_k) - \sum_{k=1}^{\infty} \Phi(s_k, t_k) - \Phi(s_0, t_0) \\
&= G - \Phi(s_0, t_0)
\end{aligned}$$

Hence, any expected reward with reward shaping is the same of the one without reward shaping but a negative shift equal to  $\Phi(s_0, t_0)$ , i.e. a constant that does not depend from the actions taken. This means that the policy cannot be affected by the shaping function.

### 5.1.3 Relevant considerations about PBRS

Here we talk about some issues in PBRS described in Section 5.1.1 (analogous considerations can be made for DPBRS described in Section 5.1.2), described in (Grzes and Kudenko, 2009; Grzes, 2010; Grześ, 2017). In particular, we focus on:

- the presence of the discount factor  $\gamma$  in Equation 5.3;
- the value of  $\Phi(s)$  when  $s$  is a terminal state.

#### The discount factor $\gamma$

In order to guarantee policy invariance,  $\gamma$  in Equation 5.3 must be equal to the discount factor of the MDP. Observe that, in general, this does not imply a *speed-up in learning time*. Indeed in (Grzes, 2010; Grzes and Kudenko, 2009) several issues of PBRS have been described, when  $\gamma < 1$ , that worsen the learning. In particular, it might happen that for some chosen value of  $\Phi(s)$ , the shaping function does not give a meaningful reward signal, e.g. near to the goal state, instead of a positive reward, a negative one signal is given to the learner, which is obviously a counterproductive choice.

In (Grzes, 2010) has been proposed an alternative approach to PBRS, which simply sets  $\gamma = 1$  in Equation 5.3, namely:

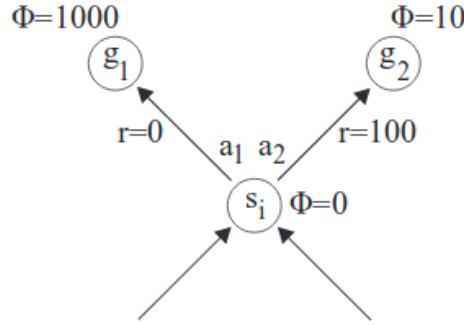
$$F(s, a, s') = \Phi(s') - \Phi(s) \quad (5.6)$$

It has been proven that this approach *does not guarantee policy invariance*, i.e. the policy learned over the MDP with reward shaping in general is not equivalent to the original MDP. In the same work, it has been shown experimental evidence of the goodness of the new approach, even in the pathological cases with  $\gamma < 1$ . So using PBRS with  $\gamma_{rs} = 1$ , even if the discount factor of the MDP  $\gamma_{mdp} \neq 1$ , "works", although the invariance of the policy is not guaranteed anymore.

### The value of terminal state $\Phi(s)$

In (Grześ, 2017) has been explained that the potential function in any terminal state (i.e. in any state where the episode terminates), must be 0 in order to guarantee policy invariance.

In Figure 5.1 is depicted a particular scenario in which the violation of this requirement over potential-based reward shaped learning leads to a different policy than non-shaped learning. In particular, without reward shaping the optimal policy from  $s_i$  would choose  $g_2$  instead of  $g_1$ , since  $r_{g_2} = 100 > r_{g_1} = 0$ . However, after applying reward shaping, the reward for the transition  $s_i \rightarrow_{a_1} g_1$  is  $r_{g_1} = 1000$ , which is higher than the one from  $s_i \rightarrow_{a_1} g_2$ , which is  $r_{g_2} = 110$ . This time, the optimal policy should prefer the transition towards  $g_1$ , although the *true* optimal policy (i.e. with no reward shaping) should prefer the transition towards  $g_2$ .



**Figure 5.1.** An example that shows why the potential function over terminal state must be 0.

More formally, consider the return of the sequence  $\bar{s}$ , similarly as Equation 4.1:

$$G(\bar{s}) := \sum_{k=0}^{N-1} \gamma^k R(s_k, s_{k+1}) \quad (5.7)$$

If we apply PBRs, it becomes:

$$\begin{aligned} G_\Phi(\bar{s}) &= \sum_{k=0}^{N-1} \gamma^k (R(s_k, s_{k+1}) + F(s_k, s_{k+1})) \\ &= \sum_{k=0}^{N-1} \gamma^k (R(s_k, s_{k+1}) + \gamma\Phi(s_{k+1}) - \Phi(s_k)) \\ &= \sum_{k=0}^{N-1} \gamma^k R(s_k, s_{k+1}) + \sum_{k=1}^N \gamma^k \Phi(s_k) - \sum_{k=0}^{N-1} \gamma^k \Phi(s_k) \\ &= G(\bar{s}) + \sum_{k=1}^{N-1} \gamma^k \Phi(s_k) + \gamma^N \Phi(s_N) - \sum_{k=1}^{N-1} \gamma^k \Phi(s_k) - \Phi(s_0) \\ &= G(\bar{s}) + \gamma^N \Phi(s_N) - \Phi(s_0) \end{aligned} \quad (5.8)$$

The term  $\Phi(s_0)$  cannot alter the policy since does not depend from any action executed; on the other hand, the term  $\gamma^N \Phi(s_N)$  depends on actions, since the

terminal state depends from the previous transitions, hence this term can modify the optimal policy. This happens whenever  $\Phi(s_N) \neq 0$ , where  $s_N$  is any state in which an episode ends, namely goal states, failure states and states where the episode ends due to the time limit exceeded.

A simple solution to this problem is to require that  $\Phi(s_N) = 0$  whenever the reinforcement learning trajectory is terminated at state  $s_N$ . Notice that if in other trajectories the same  $s_N$  is visited,  $\Phi(s_N)$  might be different from 0. The only requirement is that if  $s$  is a terminal state, then  $\Phi(s) = 0$ .

**Example 5.1.** Recalling Example 4.1, we can apply PBRs by defining a potential function that measures how far the agent is from the goal. As heuristic, we can use the *Manhattan distance* between the current position and the goal state. More formally, considering  $s_{34}$  the goal state and  $s_{ij}$  the current state, we define:

$$\Phi(s) = -[(3 - i) + (4 - j)]$$

It is easy to see that the nearer the agent to the goal state, the higher the potential function evaluated in the state of the agent. For instance, if the current state is  $s_{11}$ ,  $\Phi(s_{11}) = -5$ , whereas in  $s_{33}$  (which is nearer to the goal),  $\Phi(s_{33}) = -1$ . With this definition, transition from  $s_{11}$  to  $s_{12}$  (which makes the agent closer to the goal) evaluates  $F(s_{11}, s_{12}) = \Phi(s_{12}) - \Phi(s_{11}) = -4 - (-5) = +1$ , whereas for transition  $s_{33} \rightarrow s_{32}$  (which makes the agent more distant to the goal),  $F(s_{33}, s_{32}) = (-2) - (-1) = -1$ .

In the following sections, we will take into account the topics just described in designing a reward shaping strategy for our setting.

## 5.2 Off-line Reward shaping over $\mathcal{A}_\varphi$

In this part we propose an automatic way to apply reward shaping in the setting presented in Section 4.6. Recall that the state space of  $\mathcal{M}_{ag}^{new}$  is  $Q_1 \times \dots \times Q_m \times S$ , where  $Q_i$  is the set of states of  $\mathcal{A}_{\varphi_i}$ , the automaton associated to the  $LTL_f/LDL_f$  formula  $\varphi_i$  (see Section 2.5).

The basic intuition about our approach is that *every step toward the satisfaction of a goal formula  $\varphi_i$  should be rewarded*, analogously as it is done in reward shaping for classical goals (see Example 5.1). Hence, for a given temporal specification  $\varphi_i$  we should assign to every  $q \in Q_i$  a potential function that is, to some extent, inversely proportional to the distance from any final state.

Given a (minimal) automaton  $\mathcal{A}_\varphi$  from a  $LTL_f/LDL_f$  formula  $\varphi$  and its associated reward  $r$ , Algorithm 5.1 shows how the potential function is build from  $\mathcal{A}_\varphi$ . This operation is made *off-line*, i.e. before the learning process. Then we associate automatically to the states of the DFA a potential function  $\Phi(q)$  whose value decreases proportionally with the minimum distance between the automaton state  $q$  and any accepting state. By construction, potential-based reward shaping with this definition of the potential function gives a positive reward when the agent performs an action leading to a  $q'$  that is one step closer to an accepting state, and a negative one in the opposite case.

Notice that, by construction,  $G(\bar{s}) = G_\Phi(\bar{s})$ , where  $G_\Phi$  is defined in Equation 5.8. Indeed,  $\gamma^N \Phi(s_N) = 0$  because we take into account the issue explained in Section 5.1.3, and  $\Phi(s_0) = 0$  by construction of the Algorithm 5.1. Observe that  $G(\bar{s}) = G_\Phi(\bar{s})$  is not necessary to guarantee. It is a measure to make our reward shaping the less pervasive as possible for the learnt  $Q$  function.



**Algorithm 5.1.** Off-line Reward Shaping over  $\mathcal{A}_\varphi$ 


---

```

1: input: (minimal) automaton  $\mathcal{A}_\varphi$ , reward  $r$ 
2: output: potential function  $\Phi : Q \rightarrow \mathbb{R}$ 
3: Let sink be the sink state after the completion of  $\mathcal{A}_\varphi$ 
4: Let  $n_{q_0}$  be minimum number of hops to reach an accepting state from  $q_0$ 
5: for  $q \in Q$  do:
6:   Let  $n_q$  be minimum number of hops to reach an accepting state from  $q$ 
7:   if  $n_{q_0} \neq 0$  then ▷ i.e. if  $q_0$  is NOT an accepting state
8:      $\Phi(q) \leftarrow \frac{n_{q_0} - n_q}{n_{q_0}} \cdot r$ 
9:   else
10:     $\Phi(q) \leftarrow (n_{q_0} - n_q) \cdot r$ 
11:   end if
12: end for
13: Let  $n_{max}$  the maximum number of hops to reach an accepting state
14:  $\Phi(sink) \leftarrow \frac{n_{q_0} - n_{max}}{n_{q_0}} \cdot r$  if  $n_{q_0} \neq 0$  else  $(n_{q_0} - n_{max}) \cdot r$ 
15: return  $\Phi$ 

```

---

In the actual implementation of the Algorithm 5.1,  $\Phi(q)$  can be computed by least-fix point over the automaton  $\mathcal{A}_\varphi$ , i.e. starting from the accepting states and then explore the states from the nearest to the farthest ones.

### 5.3 On-The-Fly Reward shaping

Reward shaping can also be used when the DFAs of the LTL<sub>f</sub>/LDL<sub>f</sub> formulas are constructed *on-the-fly* (Brafman et al., 2018) so as to avoid to compute the entire automaton off-line. To do so we can rely on dynamic reward shaping (see Section 5.1.2). The idea is to build  $\mathcal{A}_\varphi$  progressively while learning. During the learning process, at every step, the value of the fluents  $\ell \in \mathcal{L}$  is observed and the successor state  $q'$  of the current state  $q$  of the DFA on-the-fly is computed. Then, the transition and the new state just observed are added into the built automaton at time  $t$ ,  $\mathcal{A}_{\varphi,t}$ , yielding  $\mathcal{A}_{\varphi,t'}$ . The potential function  $\Phi$  for  $\mathcal{A}_{\varphi,t'}$  is recomputed for the new version of the automaton. In this case, the shaping function takes the following form:

$$F(q, t, a, q', t') = \Phi(q', t') - \Phi(q, t) \quad (5.9)$$

i.e. the dynamic reward shaping in Equation 5.4 but taking into account the issue presented in Section 5.1.3 where  $\Phi(q, t)$  is a variant of the off-line case, but computed on the automaton  $\mathcal{A}_{\varphi,t}$ . In the following we explain how  $\Phi$  in the on-the-fly case differs from the one shown in Section 5.2.

#### Details about $\Phi(q, t)$

There is an important issue which has not yet been pointed out. In the *off-line* variant described in Section 5.2, we apply reward shaping at every transition by knowing the full automaton  $\mathcal{A}_\varphi$ ; however, in the *on-the-fly variant*, at the beginning of the learning task we have an "empty" automaton, i.e. only the initial state with no transition from it. Let assume that after an action the automaton makes a move from the initial state. How can we assign a positive/negative shaping reward on the transition if we do not know the goodness of the transition? And if during a simulation we discover an accepting state, there could be a nearer accepting state

that has not yet been discovered, but in order to reach it we should first discover other intermediate states. How to allow the agent to discover it while not fixing on the only known accepting states?

It is clear that, the farther the learning task goes, the more accurate will be the shaping rewards, because every observed transition of the automaton during the activity of the agent is stored and eventually the entire automaton will be explored. However, some paths ending in an accepting state are not fully explored, so how to encourage the agent, in our setting, to explore those paths?

In order to determine  $\Phi(s, t)$  for a given  $A_{\varphi, t}$ , we consider the same computations of Algorithm 5.1 but *considering the leaves (wrt the initial state) of the automaton as accepting state*. In this way, even if a path does not end in an accepting state, its following is still rewarded. It could be *wrongly rewarded*, since the path might lead to a failure state. However, the dynamic reward shaping theory states that until the potential-based condition is preserved, also the optimal and near-optimal solutions are preserved; moreover, eventually, once the final state of the path is discovered, the reward for that path will assume the right values.

The search for the leaves states is done through Depth-First Search, while the computation is the same of the Algorithm 5.1.

It is easy to see that:

**Theorem 5.2.** *Automata-based reward shaping, both in off-line and on-the-fly variants, preserves optimality and near-optimality of the MDP solutions.*

*Proof.* For the off-line case, the shaping-reward function  $\Phi$  is, by construction, potential based, hence fulfilling the premises of theorems in (Ng et al., 1999) and (Grześ, 2017). Also for the on-the-fly variant, we observe that our construction is compliant with the requirements defined in (Devlin and Kudenko, 2012).  $\square$

## 5.4 Conclusion

In this chapter, we propose an approach to apply reward shaping in the setting described in Chapter 4. We first revised the background of reward shaping theory. Then we describe an automatized way to define  $\Phi$ , by leveraging the structure of the automaton  $\mathcal{A}_{\varphi}$  associated to the LTL<sub>f</sub>/LDL<sub>f</sub> goal  $\varphi$ . We devised a variant in which the automaton is not known a-priori, but it is built from scratch during learning, and the potential function is dynamically updated. Finally, we observe that both the designed reward shaping preserve policy invariance.

## Chapter 6

# RLTG

In this chapter we describe [RLTG](#) (Reinforcement Learning for Temporal Goals), a software project written in Python. It is the reference implementation of many of the topics described in [Chapter 4](#) and [Chapter 5](#).

### 6.1 Introduction

**Main features:** RLTG is a Python framework that allows you to:

- Setup a classic reinforcement learning task by using OpenAI Gym environments.
- Modularization of a reinforcement learning system: definition of `Environment`, `Agent`, `Brain`, `Policy`, allowing for easy extension;
- Support for Sarsa( $\lambda$ ) and Q( $\lambda$ ) with  $\epsilon$ -greedy policy.
- Support for  $LTL_f$ / $LDL_f$  goal specifications, as described in [Section 4.6](#).

**Dependencies:** RLTG requires Python $\geq 3.5$  and depends on the following packages:

- [FLLOAT](#), described in [Chapter 3](#);
- [Gym OpenAI](#), a toolkit for developing and comparing reinforcement learning algorithms. It offers a useful abstraction of reinforcement learning environments.

**Installation:** You can find the package on [PyPI](#), hence you can install it with:

```
pip install rltg
```

The software is open source and is released under [MIT license](#).

### 6.2 Package structure

The package is structured as follows:

- `agents/`: contains multiple packages and modules to build the learning agent.

- `brains/`: it contains `Brain.py`, the basic abstraction of the reinforcement learning algorithm, as well as `TDBrain.py`, the implementation of Temporal Difference learning (see Section 4.3);
  - `parameters/` and `policies/` contain implementations for eligibility traces and  $\epsilon$ -greedy policy;
  - `feature_extraction.py`: it contains classes for feature extraction like `FeatureExtractor`. They leverage the `Space` abstraction provided by OpenAI Gym. [Here](#) you can find many types of state spaces supported.
  - `Agent.py`, the module that contains the `Agent` abstraction, defining the abstract methods for interact with the environment (e.g. `observe` for observing the new environment state, `step` for choose an action). It needs a `FeatureExtractor` and a `Brain`.
  - `TemporalEvaluator.py`, which implements the temporal goal specified by an  $LTL_f/LDL_f$  formula.
- `logic/`: contains the implementations of reward shaping as described in Section 5.2 and 5.3. The main modules are `CompleteRewardAutomaton.py`, for the implementation of Off-line Reward Shaping, and `PartialRewardAutomaton.py`, for the On-the-fly Reward Shaping.
  - `trainers/`: contains the implementation of a highly-customizable training loop, e.g. you can specify under which conditions the training should stop, you can visually render the agent during the learning and you can pause/resume the learning. The `GenericTrainer` is used for classic reinforcement learning, whereas `TGTrainer` is used for temporal goal reinforcement learning.

## 6.3 Code examples

### 6.3.1 Classic Reinforcement Learning

Here we will see an example of reinforcement learning task using RLTLG with the OpenAI Gym environment `Taxi-v2` ([Dietterich, 1998](#)), available [here](#).

**Listing 6.1.** Classic Reinforcement Learning using RLTLG

```

1 import gym
2
3 from rltg.agents.RLAgent import RLAgent
4 from rltg.agents.brains.TDBrain import Sarsa
5 from rltg.agents.feature_extraction\
6     import IdentityFeatureExtractor
7 from rltg.agents.policies.EGreedy import EGreedy
8 from rltg.trainers.GenericTrainer import GenericTrainer
9 from rltg.utils.GoalEnvWrapper import GoalEnvWrapper
10 from rltg.utils.StoppingCondition\
11     import AvgRewardPercentage
12
13 def taxi_goal(*args):
14     reward =args[1]
15     done =args[2]
16     return done and reward == 20
17

```

```

18 if __name__ == '__main__':
19     env = gym.make("Taxi-v2")
20     env = GoalEnvWrapper(env, taxi_goal)
21
22     observation_space = env.observation_space
23     action_space = env.action_space
24     print(observation_space, action_space)
25     agent = RLAgent(
26         IdentityFeatureExtractor(observation_space),
27         Sarsa(observation_space, action_space,
28             EGreedy(0.1), alpha=0.1, gamma=0.99, lambda_=0.0)
29     )
30
31     tr = GenericTrainer(
32         env,
33         agent,
34         n_episodes=10000
35         stop_conditions=(
36             AvgRewardPercentage(window_size=100, target_mean=9.0)
37         )
38     )
39     tr.main()

```

Now we analyze some of the APIs provided by RL TG:

- In line 19 we defined the environment by OpenAI Gym APIs. In line and 20 we wrapped with a function that make it *goal-based*.
- In line 25 we defined our reinforcement learning agent. We had to specify a **FeatureExtractor**, i.e. which features of the state provided from the environment we consider relevant, and a **Brain**, i.e. the responsible for the learning of the policy. In particular:
  - **IdentityFeatureExtractor** extends **FeatureExtractor** and it simply get the entire state space from the environment, without modifying it; It requires the observation state space of the environment to do some sanity checks.
  - **Sarsa** extends **TDBrain**, and implements  $Sarsa(\lambda)$  (see Section 4.3). It requires:
    - \* the observation space where the algorithm learns (in this case is the same of the one provided by the environment);
    - \* the action space from where actions can be selected;
    - \* the behavior policy (in this case  $\epsilon$ -greedy policy with  $\epsilon = 0.1$ ).
    - \* the learning rate  $\alpha = 0.1$ , the discount factor  $\gamma = 0.99$ , the eligibility trace parameter  $\lambda = 0$ .
  - In line 31 we configured the manager of the training task. **GenericTrainer**, which extends from **Trainer**, is configured by the environment **env** and the agent **agent** defined before, the maximum number of episodes **n\_episodes** and the stop conditions. **AvgRewardPercentage** with parameters **window\_size** = 100 and **target\_mean** = 9.0 means that we will stop the training if in the last 100 episodes an average reward of 9.0 is observed.

If you run the script, you will get some statistics for each episode, e.g. the number of explored states, the total reward observed, if the goal has been reached.

### 6.3.2 Temporal goal Reinforcement Learning

Now we will see how to use RLTLG for a temporal goal specified by a  $LTL_f/LDL_f$  formula. In particular, we will use the BREAKOUT environment, presented in Example 4.2 and that we will further describe in Chapter 7. In the next, we show an excerpt of [this script](#):

**Listing 6.2.** Temporal Goal Reinforcement Learning using RLTLG

```

1  ...
2  ...
3  ...
4  if __name__ == '__main__':
5      env = GymBreakout(brick_cols=3, brick_rows=3)
6
7      gamma = 0.999
8      on_the_fly = False
9      reward_shaping = False
10
11     agent = TGAgent(
12         BreakoutNRobotFeatureExtractor(env.observation_space),
13         Sarsa(None, env.action_space, policy=EGreedy(0.1),
14             alpha=0.1, gamma=gamma, lambda_=0.99
15         ),
16         [BreakoutCompleteColumnsTemporalEvaluator(
17             env.observation_space, bricks_rows=env.brick_rows,
18             bricks_cols=env.brick_cols, left_right=True,
19             gamma=gamma, on_the_fly=on_the_fly)
20         ],
21         reward_shaping=reward_shaping
22     )
23
24     tr = TGTrainer(env, agent, n_episodes=2000,
25         stop_conditions=(GoalPercentage(100, 0.2),),
26     )
27
28     stats, optimal_stats = tr.main()

```

Observe that:

- In Listing 6.1 we used `RLAgent`, while in this case, in line 11, we use `TGAgent`. The arguments provided in the constructors are the same of the previous example except for an additional argument which is the list of `TemporalEvaluator` (line 16). Notice that this definition of *list of temporal goals* is compliant with Definition 4.6. We do not discuss here details about this particular temporal evaluator, i.e. `BreakoutCompleteColumnsTemporalEvaluator`. Observe that we can specify if the temporal evaluator should use a off-line automaton or a on-the-fly automaton by the boolean parameter `on_the_fly`. Finally, notice the `reward_shaping` flag to toggle automata-based reward shaping (Section 5.2 and 5.3).
- Another difference is in using `TGTrainer` instead of `GenericTrainer`. This version is more specialized to deal with the setting explained in Section 4.6. The signature of the constructor is the same of `GenericTrainer`. Notice that in this particular case we used `GoalPercentage(100, 0.2)`, which means

---

we require the agent have reached the goal at least 20 times in the last 100 episodes.

# Chapter 7

## Experiments

In this chapter we give experimental evidence of the goodness of our approach, discussed mainly in Chapter 4 and 5, by using the implementations presented in Chapter 3 and 6. For each considered environment, we give a formal description in the framework of MDPs and assign temporal goals to be satisfied, as explained in earlier chapters.

### 7.1 BREAKOUT

In this section we use the BREAKOUT environment (already presented in Example 4.2). We first show how the combination of temporal goals is successfully learnt by the agent. Then, we show a benchmarking between the use of off-line reward shaping, on-the-fly reward shaping and no reward shaping.

**Description of the Environment:** The actions available to the agent are *left*, *right* and *no-action*. The relevant features are: position of the paddle  $p_x$ , position of the ball  $b_x, b_y$ , speed of the ball  $v_x, v_y$  and status of each brick (booleans)  $b_{ij}$ . This features of the system gives all the needed information to predict the next state from the current state. Hence we can build an MDP where:  $S$  is the set of all the possible values of the sequence of features  $\langle p_x, b_x, b_y, v_x, v_y \rangle$ ,  $A = \{\text{right}, \text{left}, \text{no-action}\}$ , transition function  $T$  determined by the rules of the game. We give reward  $R(s, a, s') = 10$  if a particular brick in  $s'$  has been removed for the first time, plus 100 if that brick was the last (i.e. *environment goal* reached).

**Temporal goal:** The *temporal goal* (specified by a  $\text{LTL}_f/\text{LDL}_f$  formula) is to remove lines of bricks in a given order. In other words, all bricks on each line of bricks  $i$  must be removed before removing all bricks from line  $j > i$ . In our experiments we considered the following goals (and combinations of them):

- BREAK-COLS-(LR|RL): remove columns of blocks from (left to right | right to left).
- BREAK-ROWS-(BT|TB): remove rows of blocks from (bottom to top | top to bottom).

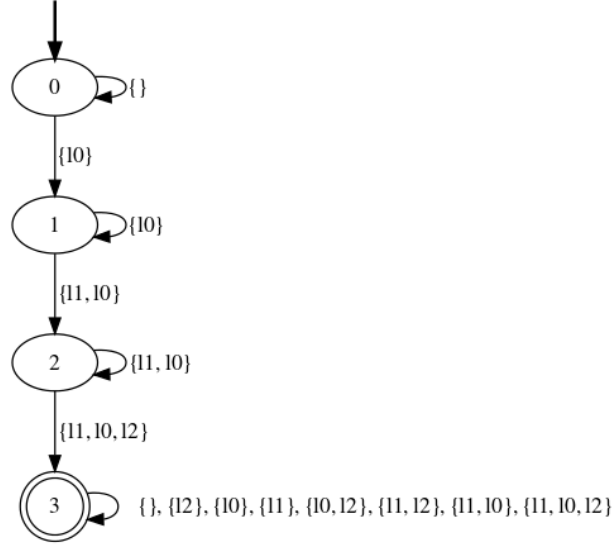
For each temporal goal we give a reward  $r = 10000$  and it is given when the agent fulfilled the temporal goal specified by the formula. The formulas are expressed in



$LDL_f$ . For example, in Breakout 3x3 (i.e. three rows and three columns of bricks), the formula which specify the just explained temporal goals is:

$$\langle (\neg l_0 \wedge \neg l_1 \wedge \neg l_2)^*; (l_0 \wedge \neg l_1 \wedge \neg l_2); (l_0 \wedge \neg l_1 \wedge \neg l_2)^*; (l_0 \wedge l_1 \wedge \neg l_2); (l_0 \wedge l_1 \wedge \neg l_2)^*; (l_0 \wedge l_1 \wedge l_2) \rangle tt \quad (7.1)$$

where the fluent  $l_i$  means "the  $i_{th}$  line has been removed". The automaton associated to the formula in Equation 7.1 is depicted in Figure 7.1.  $l_i$  are the fluents of interests, hence  $\mathcal{P} = \{l_0, l_1, l_2\}$  and  $\mathcal{L} = 2^{\mathcal{P}}$ .



**Figure 7.1.** The automaton associated to the  $LDL_f$  formula in Equation 7.1 for BREAKOUT 3x3.

The features for fluents evaluation, for every temporal goal, are the status of each brick (present or removed). Recalling the notation used in Example 4.2:  $\langle b_{11}, \dots, b_{nm} \rangle$ .

It is crucial to remark that the evaluation of the fluent from the features is different among the temporal goals above mentioned, despite the formulas are structurally the same (apart from the number of lines to remove). Indeed, thanks to the fluents evaluation phase (i.e. the map from features values to truth of each fluent), the formulas behave differently.

**Configurations:** The reinforcement algorithm used is Sarsa( $\lambda$ ) (e.g. Sarsa with eligibility traces) with  $\epsilon$ -greedy policy. The values of the parameters are:

- $\lambda = 0.99$
- $\gamma = 0.999$
- $\alpha = 0.1$
- $\epsilon = 0.1$

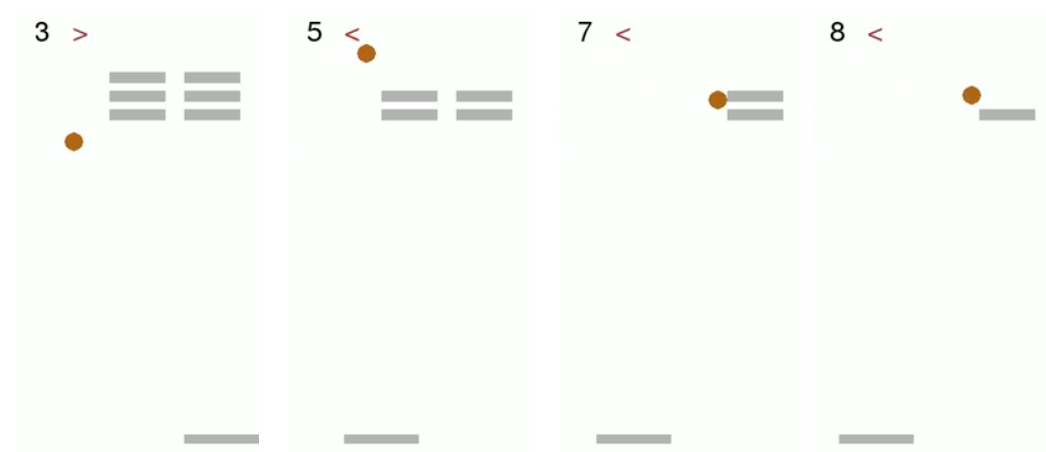
The experiments are stopped when in the last 100 episodes of optimal runs the agent always achieved both the environment goal and the temporal goals. Notice that this approach may yield a sub-optimal policy, but the found policies always satisfy  $LTL_f/LDL_f$  goals. This method is used also for the other experiments.

### 7.1.1 Optimal policies

Now we show some screenshots of the optimal policies found for some of the temporal goals described before. You can find full recordings at <https://www.youtube.com/channel/UChe0QEtrSKh4uCy5f2XuGA>.

#### Break-Cols-LR

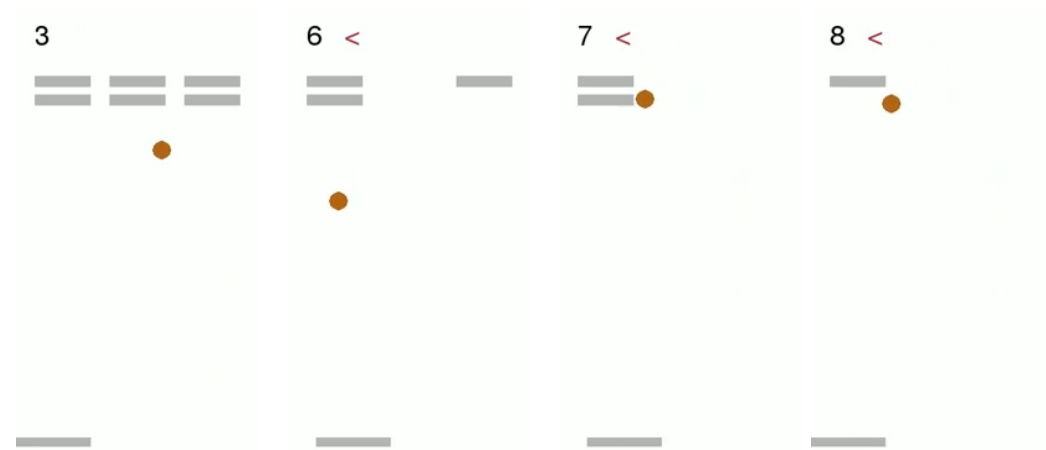
In Figure 7.2 are depicted the highlights of a run of the optimal policy for the goal BREAK-COLS-LR. As you can notice, the constraint expressed by the formula in Equation 7.1 is satisfied (i.e. the columns of bricks are broken from left to right).



**Figure 7.2.** A run of the learnt optimal policy for the task BREAK-COLS-LR. From left to right, you can see that the columns are broken in the right order.

#### Break-Cols-BT

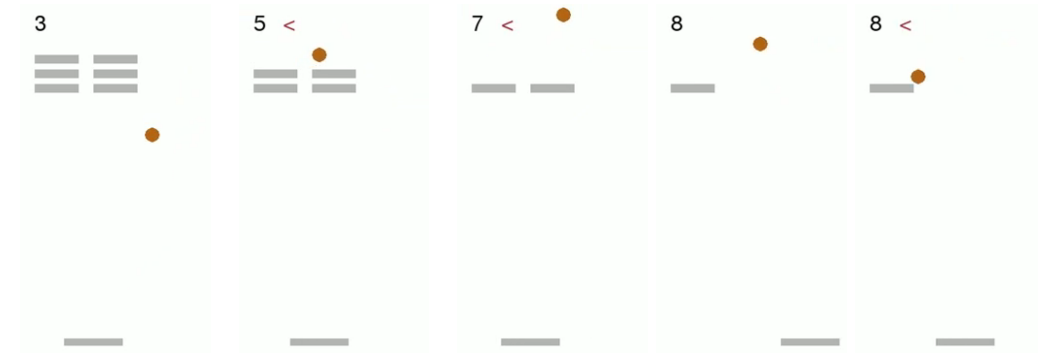
In Figure 7.3 are depicted the highlights of a run of the optimal policy for the goal BREAK-COLS-BT. The constraint expressed by the formula in Equation 7.1 is satisfied (i.e. the rows of bricks are broken from the bottom to the top).



**Figure 7.3.** A run of the learnt optimal policy for the task BREAK-ROWS-BT. From the bottom to the top, you can see that the rows are broken in the right order.

### Break-Cols-RL-TB

In Figure 7.4 are depicted the highlights of a run of the optimal policy for the goal BREAK-COLS-RL-TB. Notice that, even in the case of multiple temporal goal, the constraints expressed by the formula in Equation 7.1 are satisfied (i.e. the rows of bricks are broken from the top to the bottom and the columns of bricks are broken from right to left). Furthermore, notice that the formula is the same, but *how the fluents are evaluated from the features* makes the difference.



**Figure 7.4.** A run of the learnt optimal policy for the task BREAK-COLS-RL-TB. You can see that the rows are broken from the top to the bottom and that, at the same time, the columns are broken from the right to the left.

#### 7.1.2 Benchmarking of reward shaping techniques

In this section we report some statistics about the performances of the proposed approach with a focus on automata-based reward shaping. For each configuration we collected the observed reward for each episode, over 10 runs. Each run has a time limit of 10000 episodes. Then we moving averaged the sequences of rewards with a window of size 100, in order to make the curves more smoothed, and averaged the result across all the runs. The plots show this processed sequences over the number of episodes. The colored bands represents the 90% confidence interval taken from bootstrapped re-samples.

We show our results in two parts:

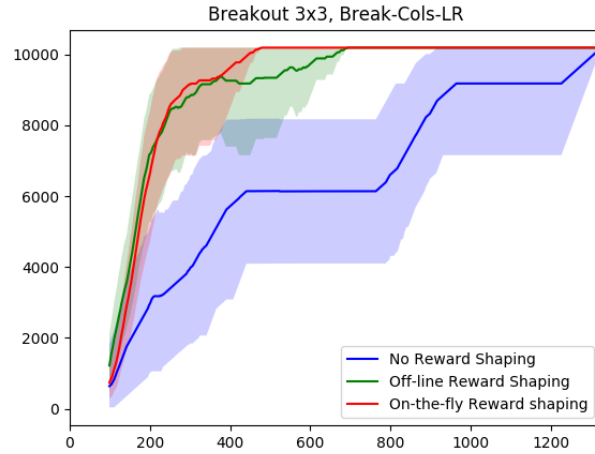
1. Comparison of *No Reward Shaping*, *Off-line Reward Shaping* and *On-the-fly Reward Shaping* in BREAKOUT 3x3, 3x4 and 4x4, with temporal goal BREAK-COLS-LR;
2. Comparison of *No Reward Shaping*, *Off-line Reward Shaping* and *On-the-fly Reward Shaping* in Breakout 4x4, with temporal goals BREAK-COLS-LR, BREAK-ROWS-BT, BREAK-COLS-LR & BREAK-ROWS-BT;

#### Different number of bricks and rows

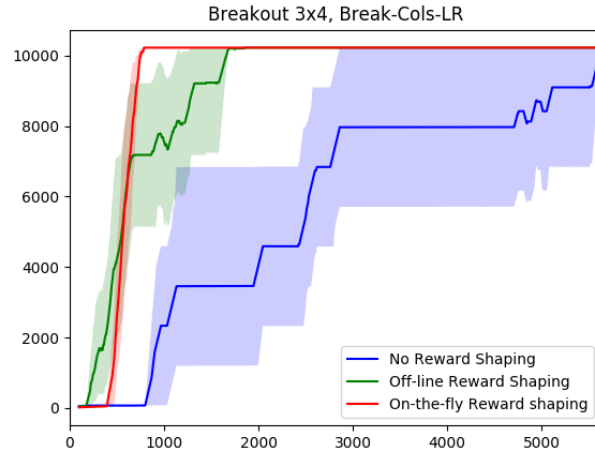
In Figure 7.5 are plotted different statistics by increasing difficulty in the environment, for every variant of reward shaping technique. On the x-axis the number of episodes, on the y-axis the average reward (obtained as explained above). In every case, we can see that the use of reward shaping (both off-line and on-the-fly) speed-up the learning process.

### Different temporal goals

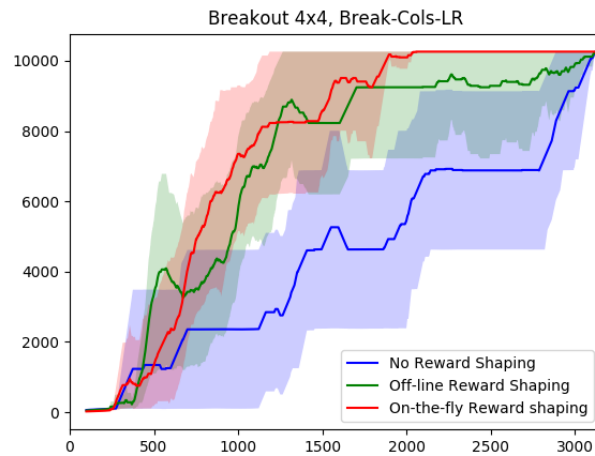
In Figure 7.6 Here we show the performances by varying the temporal goals. On the x-axis the number of episodes, on the y-axis the average reward (obtained as explained above). In every case, we can see that the use of reward shaping (both off-line and on-the-fly) speed-up the learning process.



(a) BREAKOUT 3x3, BREAK-COLS-LR for three different settings: *No RS*, *Off-line RS* and *On-the-fly RS*

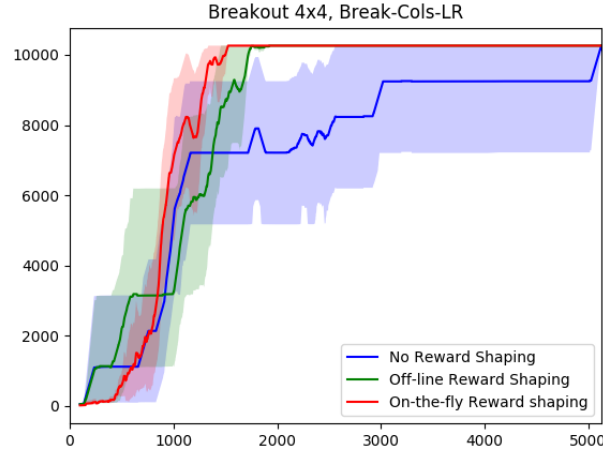


(b) BREAKOUT 3x4, BREAK-COLS-LR for three different settings: *No RS*, *Off-line RS* and *On-the-fly RS*

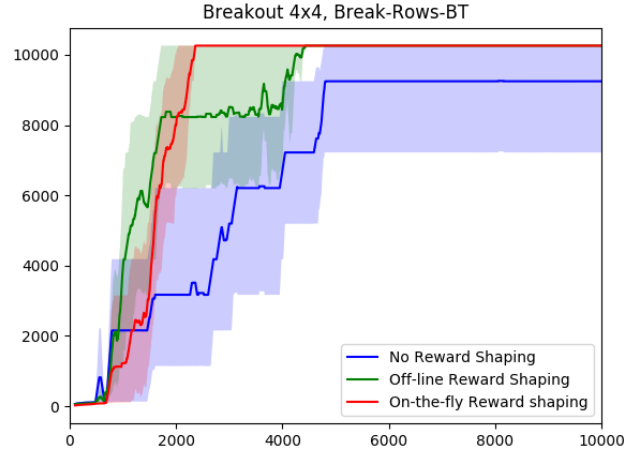


(c) BREAKOUT 4x4, BREAK-COLS-LR for three different settings: *No RS*, *Off-line RS* and *On-the-fly RS*

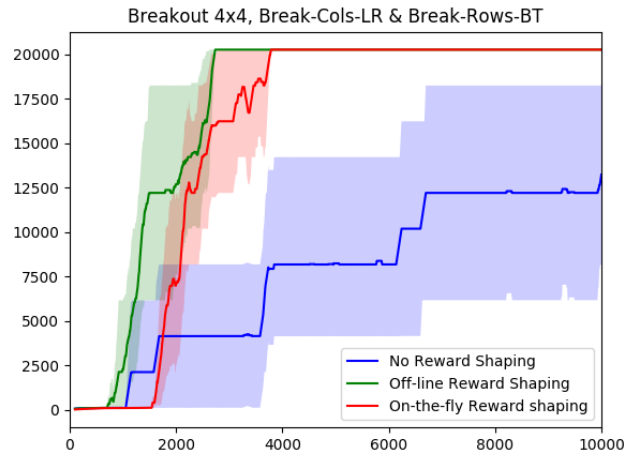
**Figure 7.5.** The results of three settings are reported, namely: BREAKOUT 3x3, 3x4 and 4x4 (7.5a, 7.5b and 7.5c respectively), all of them with temporal goal BREAK-COLS-LR. On the x-axis the episodes, on the y-axis the average reward. In every case, we can see that the use of reward shaping (both off-line and on-the-fly) speed-up the learning process.



(a) BREAKOUT 4x4, BREAK-COLS-LR for three different settings: *No RS*, *Off-line RS* and *On-the-fly RS*



(b) BREAKOUT 4x4, BREAK-ROWS-BT for three different settings: *No RS*, *Off-line RS* and *On-the-fly RS*



(c) BREAKOUT 4x4, BREAK-COLS-LR & BREAK-ROWS-LR for three different settings: *No RS*, *Off-line RS* and *On-the-fly RS*

**Figure 7.6.** The results of three settings are reported, namely: BREAKOUT 4x4 with temporal goal BREAK-COLS-LR, BREAK-COLS-BT and BREAK-COLS-LR & BREAK-COLS-BT (7.6a, 7.6b and 7.6c respectively). On the x-axis the episodes, on the y-axis the average reward. In every case, we can see that the use of reward shaping (both off-line and on-the-fly) speed-up the learning process.

## 7.2 SAPIENTINO

SAPIENTINO Doc is an educational game for 5-8 y.o. children in which a small mobile robot has to be programmed in order to visit specific cells in a 5x7 grid. Cells contain concepts that must be matched by the children (e.g., a colored animal, a color, and the initial letter of the name of the animal). The robot executes sequences of actions given in input by children with a keyboard on the robot's top side. During execution, the robot moves on the grid and executes an action (actually a *bip*) to announce that the current cell has been reached (this is called a *visit* of a cell). A pair of consecutive visits are correct when they refer to cells containing matching concepts. In this paper, we generalize this game as follows. As in the real game, we consider a 5x7 grid with 7 triplets of colored cells, each triplet representing three matching concepts (see Figure 7.7 for a screenshot).

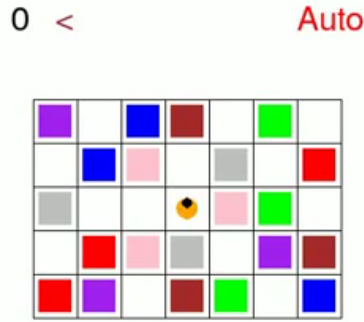


Figure 7.7. A screenshot from SAPIENTINO

**Actions, Features and Fluents:** The actions available to the agent are: UP, DOWN, LEFT, RIGHT and BIP. The features for the agent space are the position  $(x, y)$  of the agent in the grid, namely  $f_x$  and  $f_y$ . The features for evaluate the fluents configurations are:  $f_b$  reporting if a BIP has just been executed, and  $f_c$  denoting the color of the current cell. In our setting, the available fluents are:

$$\mathcal{P} = \{bip, red, green, blue, pink, brown, gray, purple\}$$

with the obvious meaning that their name indicate (just map features to fluents one-to-one).

**Temporal Goal:** In our experiment, we considered the following temporal behaviors:

- ORDERED-VISITS: *visit* each color and do a *bip* in each color in a given order.
- ILLEGAL-BIPS: between each visit, never do a *bip*.

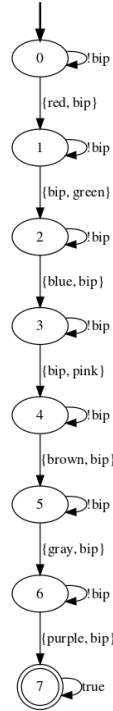
We identified two temporal goals: one considering that both conditions must be satisfied and a more relaxed one that only ORDERED-VISITS must be satisfied. They can be translated, respectively, into the following  $LDL_f$  formulas:

$$\begin{aligned} &\langle (\neg bip)^*; red \wedge bip; (\neg bip)^*; green \wedge bip; \\ &(\neg bip)^*; blue \wedge bip; (\neg bip)^*; pink \wedge bip; \\ &(\neg bip)^*; brown \wedge bip; (\neg bip)^*; gray \wedge bip; \\ &(\neg bip)^*; purple \wedge bip \rangle tt \end{aligned} \quad (7.2)$$

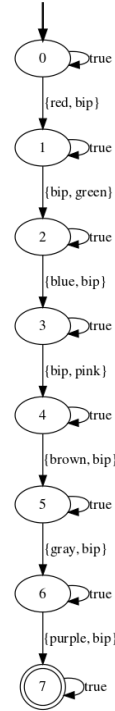
and:

$$\begin{aligned} &\langle true^*; red \wedge bip; true^*; green \wedge bip; \\ &true^*; blue \wedge bip; true^*; pink \wedge bip; \\ &true^*; brown \wedge bip; true^*; gray \wedge bip; \\ &true^*; purple \wedge bip \rangle tt \end{aligned} \quad (7.3)$$

where Formula 7.3 is the same of Formula 7.2 but replacing every occurrence of  $(\neg bip)^*$  with  $true^*$ . The associated automaton are depicted, respectively, in Figure 7.8a and 7.8b.



(a) The automaton associated to the  $LDL_f$  formula in Equation 7.2 for the "full" temporal goal in SAPIENTINO.



(b) The automaton associated to the  $LDL_f$  formula in Equation 7.3 for the "relaxed" temporal goal in SAPIENTINO.

Notice that the condition of illegal bips does not affect the optimal policy; on the other hand, it highly affects the exploration phase, since with the ILLEGAL-BIPS condition, every time the  $\epsilon$ -greedy policy chooses *bip* randomly and "illegally",



then that trajectory cannot satisfy the temporal goal anymore. Moreover, observe that in this case we do not consider any *environment goal*, as we did in BREAKOUT experiments (i.e. break all the bricks).

**Configurations:** The reinforcement algorithm used is Sarsa( $\lambda$ ) (e.g. Sarsa with eligibility traces) with  $\epsilon$ -greedy policy. The values of the parameters are:

- $\lambda = 0$
- $\gamma = 1.0$
- $\alpha = 0.1$
- $\epsilon = 0.1$

We show our results in two parts:

1. Comparison of *No Reward Shaping*, *Off-line Reward Shaping* and *On-the-fly Reward Shaping* in SAPIENTINO with the "full" temporal goal, i.e. ORDERED-VISITS & ILLEGAL-BIPS;
2. Comparison of *No Reward Shaping*, *Off-line Reward Shaping* and *On-the-fly Reward Shaping* in SAPIENTINO with the "relaxed" temporal goal, i.e. only ORDERED-VISITS ;

### 7.2.1 Results

In this part we analyze the results of Experiment 1 and 2. As stated before, since the two experiments differ only for the exploration phase, the optimal policy that satisfy the goal is qualitatively the same. You might see the resulting learnt policy at <https://www.youtube.com/watch?v=ghBImtYmpVk>.

#### Ordered-Visits & Illegal-Bips

In this experiment, we observed that the agent successfully learned the temporal goal, both in off-line and in on-the-fly reward shaping. The performances are pretty the same. Notice that, without reward shaping, in our simulations (5000 episodes), *the agent never reached the goal*. This is due mainly to the long sequence of actions that the agent has to take in order to accomplish the entire goal, considering that an illegal *bip* lead to a failure of the goal and so the end of the episode. On the other hand, in the case of reward shaping, an illegal bip is punished because the episode ends and the agent collects the negative potential fall. Hence, the agent recognize, earlier than no reward shaping configuration, that the illegal bip is a bad action.

#### Ordered-Visits

Also in this case, the agent learned the temporal goal, even in the case where no reward shaping was applied. It is interesting to note that a small change in the formula can affect drastically the learning process, due to exploration phase constraints.

### 7.3 MINECRAFT

MINECRAFT (Andreas et al., 2017) is a sort of 2D porting of the well-known 3D videogame *Minecraft*. The agent can *get* resources (e.g. *wood*, *grass*) and *use* tools (e.g. *toolshed*, *workbench*) located on a grid (similarly to the grid in SAPIENTINO), and has to accomplish composite tasks. For instance, the task *make a bed* is divided into the following sequence of subtasks: *get\_wood*, *use\_toolshed*, *get\_grass*, *used\_workbench*. In Figure 7.9 you can see a screenshot of the game.

**Actions, Features and Fluents:** The actions available to the agent are: UP, DOWN, LEFT, RIGHT, GET and USE. The features for the agent space are the position  $(x, y)$  of the agent in the grid, namely  $f_x$  and  $f_y$ . The features for evaluate the fluents configurations are:  $f_g$  reporting if a GET has just been executed,  $f_u$  reporting if a USE has just been executed,  $f_\ell$  denoting the resource/tool available in the current location, and  $f'_\ell$ , if a resource has been taken and is available for the incoming tasks. In our setting, the available fluents are:

$$\mathcal{P} = \{get\_wood, get\_grass, get\_iron, use\_toolshed, use\_workbench, use\_factory\}$$

when the agent is near a resource (i.e. *wood*, *grass* or *iron*) and do a GET, then the associated "get" fluent becomes true. Analogously for tools and the USE action. When a task is completed and to do that some previously collected resource/tool has been used, then the resources are lost and has to be recollected again in order to accomplish other tasks.

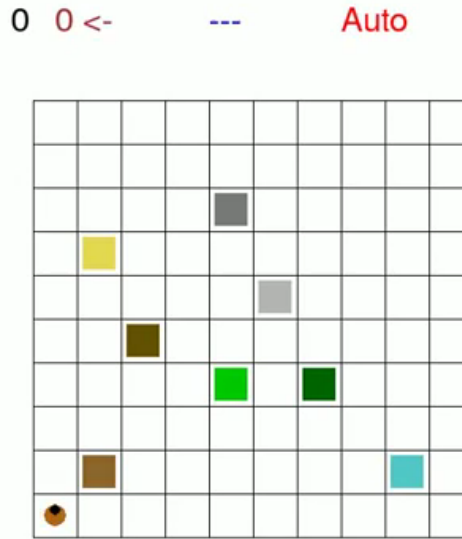


Figure 7.9. A screenshot of the game MINECRAFT

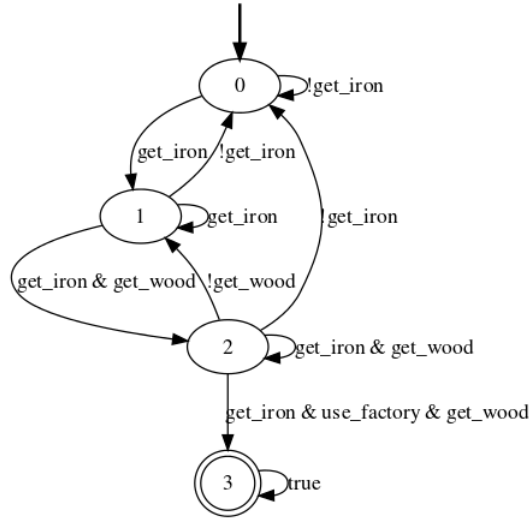
**Temporal Goal:** In our experiment, we considered the following tasks:

1. MAKE-BED: *get\_wood*, *use\_toolshed*, *get\_grass*, *use\_workbench*
2. MAKE-AXE: *get\_wood*, *use\_workbench*, *get\_iron*, *use\_toolshed*
3. MAKE-BRIDGE: *get\_iron*, *get\_wood*, *use\_factory*

Now we show how to express them in  $LDL_f$  and its associated automaton. For simplicity, we show only the ones for the goal 3. The associated  $LDL_f$  formula is:

$$\begin{aligned} \langle true^* \rangle \langle true^*; & get\_iron \wedge \neg get\_wood \wedge \neg use\_factory; \\ & (get\_iron \wedge \neg get\_wood \wedge \neg use\_factory)^*; get\_iron \wedge get\_wood \wedge \neg use\_factory; \\ & (get\_iron \wedge get\_wood \wedge \neg use\_factory)^*; get\_iron \wedge get\_wood \wedge use\_factory \rangle tt \end{aligned} \quad (7.4)$$

whereas the translation of Formula 7.4 into automaton is depicted in Figure 7.10:



**Figure 7.10.** An approximate representation of the automaton associated to Formula 7.4. Notice that it is "simplified" by omitting groups of transitions that are collapsed in a single edge labeled with the propositional formula that the interpretation should satisfy.

Notice that we could "relax" the specified order of *get\_wood* and *get\_iron*, since it is not relevant to the operation *make a bridge*.

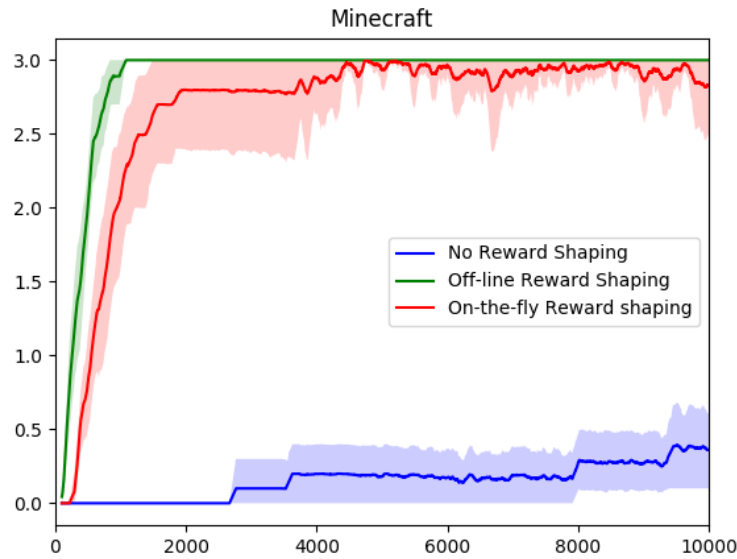
The agent is trained to satisfy all the three tasks in a single episode. To each temporal goal, the reward assigned is 1.

**Configurations:** The reinforcement algorithm used is Sarsa( $\lambda$ ) (e.g. Sarsa with eligibility traces) with  $\epsilon$ -greedy policy. The values of the parameters are:

- $\lambda = 0.9$
- $\gamma = 0.99$
- $\alpha = 0.1$
- $\epsilon = 0.25$

### 7.3.1 Results

In Figure 7.11 is plotted the benchmarking, similarly to the one presented in Section 7.1.2. We notice that the policy that satisfies every temporal goal has been found when a reward shaping is applied, whereas in the case of no reward shaping the time



**Figure 7.11.** A comparison between reward shaping techniques for the temporal goals 1, 2 and 3. On the x-axis the episodes, on the y-axis the average reward, as explained in Section 7.1.2.

limit imposed for the experiment was not enough to allow every run to achieve the goal. You can see one of these found policies at <https://youtu.be/IJ3Hr79xfBs>.

Observe that reward shaping, both in off-line and on-the-fly variant, outperform the absence of reward shaping. This is due mainly to the high number of steps that the agent has to make in a precise order in order to achieve the goal: without a guidance in the exploration (e.g. shaping rewards) it is hard to discover the proper sequence that leads to the satisfaction of every formula. Furthermore, observe that off-line reward shaping is slightly better than on-the-fly reward shaping in terms of convergence rate. It is also more stable, as the reader might infer from the shaded regions that represents the variance of the plot: indeed, it is not surprising that on-the-fly reward shaping yields a more noisy learning process, due to the dynamic change of the potential function and so of the shaping rewards.

It is worth to notice the particular structure of the environment in our setting. Indeed, when one of the task is completed, the resources that progressed at the same time the other automata/tasks are released; this lead to a regression of the state of those tasks, and so to a negative shaping reward. Hence, the positive shaping reward given for progression is mitigated by the negative shaping reward due to completion of a task and regression of the other task that have some resource in common. The learning agent has to deal with this issue, that at the scaling of the problem it might lead to slow the learning.

## 7.4 Implementation

The experiments are implemented using FLLOAT and RLTG, presented in Chapter 3 and Chapter 6. The code for run the experiments by yourself with other configurations and plot the statistics is published in [this repository](#). You'll find also the scripts to run the experiments analyzed in this chapter.

## Chapter 8

# Conclusion and Future Work

This chapter summarizes the thesis. It starts with a brief overview of the problem studied. Then we summarize the main contributions of the thesis. A section is dedicated to future directions of research. The document ends with final remarks.

### 8.1 Overview

This work addresses a particular problem in reinforcement learning. We considered a classical reinforcement learning problem, i.e. an MDP over a set of features of the world, that we called *low-level* features <sup>1</sup>. We are interested in defining non-Markovian rewards over this MDP about some high-level properties of the environment, that we called the *fluents*, determined by some set of *high-level* features of the world.

We considered the cartesian product between the MDP state space and the fluents configurations, which yields, in general, an unknown transition function. In order to proceed, we made a key assumption: that is, *the new transition function (i.e. defined both over low-level features and fluents configurations) satisfies the Markov property*, i.e. given the current MDP state, the current fluents configuration and the action taken, we have all the informations needed to know the probability of each state to be the next one.

However, due to the non-Markovian property of the additional rewards, we need to transform the new state space in some way such that the learning is actually feasible and such that the learnt policy in this transformed state space is equivalent (in terms of optimality) to the original problem. How this can be effectively achieved is one of the main contribution of this work.

### 8.2 Summary of Main Contributions

In this section we list the main contribution of the thesis.

- Formal definition and analysis of the problem, as well as design of the transformation of the state space (strongly inspired by [Brafman et al. \(2018\)](#)), allowing the off-the-shelf reinforcement learning algorithms to actually learn

---

<sup>1</sup>We say *low-level* features in the sense that these features are only responsible for the basic interaction of the agent with the environment, but still allowing to maximize the reward, or reach the goal in the MDP. Notice that it is only a way to refer about them, and no assumption nor restriction is made over the MDP state space.

the non-Markovian goals (Chapter 4), specified by  $LTL_f/LDL_f$  formulas (described in Chapter 2);

- An automatic way to apply reward shaping in this setting, leveraging the particular structure of the transformation (Chapter 5). The idea is, the transitions in the new state space that make a step over the automaton towards an accepting state (or, equivalently, any progression in the satisfaction of the formula) should be rewarded. We dealt with both when the automaton is known a priori and when it is built on-the-fly;
- Implementation of both the translation algorithm from  $LTL_f/LDL_f$  formulas to equivalent automata, presented in Chapter 3, and a reinforcement learning framework to set up a RL agent satisfying  $LTL_f/LDL_f$  formulas, presented in Chapter 6;
- Experimental evidence that the approach is actually working (Chapter 7).

### 8.3 Future Works

There are many future directions that can be taken, due to the novelty of the work. Some of them are:

- the Markov assumption of the combined transition function that has been made in Section 4.6.1 is a strong one; however, in many real world cases, the implicit approximation is not enough to effectively model the world. It might be interesting to find an approach that addresses this issues and manages the problem in this harder scenario.
- In this work we prevalently made a theoretical analysis of the problem and we have shown only some application at software level. The proposed framework is more general, and should be validated in many other domains, both simulated and physical ones (e.g. robotics), where there is the need to represent additional high-level knowledge to express complex goals;
- We did not focus our attention over the algorithmic side, but only on how the problem can be properly redefined, and relying on off-the-shelf reinforcement learning algorithms. It might be the case that the design of ad-hoc algorithms for this framework yields better performances. For instance, one could think about a better "explicit" exploration policy of the state space, in spite of "implicit" guidance by using reward shaping, and even some adaptation of state-of-the-art techniques, e.g. hierarchical reinforcement learning, curriculum learning, knowledge transfer and so on.
- Finally, it could be of interest to extend this work to the general framework of multi-agent systems and work on the challenges and issues that such extension might lead to.

### 8.4 Final Remarks

The topic of this work is of crucial importance: the need to face the dichotomy between the control of low-level features and the ability to reason about higher-level properties of the world is the core issue in many of the fields and applications in artificial intelligence, especially in robotics. This thesis aimed to address this issue

---

in a particular case, by a theoretical analysis, but providing actual implementations as support for the proposed approach.

# Bibliography

- Jacob Andreas, Dan Klein, and Sergey Levine. Modular multitask reinforcement learning with policy sketches. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 166–175, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR. URL <http://proceedings.mlr.press/v70/andreas17a.html>.
- Fahiem Bacchus, Craig Boutilier, and Adam Grove. Rewarding behaviors. In *PROCEEDINGS OF THE NATIONAL CONFERENCE ON ARTIFICIAL INTELLIGENCE*, pages 1160–1167, 1996.
- Jorge A. Baier, Christian Fritz, Meghyn Bienvenu, and Sheila A McIlraith. Beyond classical planning: Procedural control knowledge and preferences in state-of-the-art planners. In *Proceedings of the 23rd National Conference on Artificial Intelligence - Volume 3*, AAAI’08, pages 1509–1512. AAAI Press, 2008. ISBN 978-1-57735-368-3. URL <http://dl.acm.org/citation.cfm?id=1620270.1620321>.
- Richard Bellman. *Dynamic Programming*. Princeton University Press, Princeton, NJ, USA, 1 edition, 1957. URL <http://books.google.com/books?id=fyVtp3EMxasC&pg=PR5&dq=dynamic+programming+richard+e+bellman&client=firefox-a#v=onepage&q=dynamic%20programming%20richard%20e%20bellman&f=false>.
- Ronen Brafman, Giuseppe De Giacomo, and Fabio Patrizi. Ltlf/ldlf non-markovian rewards, 2018. URL <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/17342>.
- Ronen I. Brafman, Giuseppe De Giacomo, and Fabio Patrizi. Specifying non-markovian rewards in mdps using ldl on finite traces (preliminary version). *CoRR*, abs/1706.08100, 2017.
- Richard J. Büchi. Weak second-order arithmetic and finite automata. *Mathematical Logic Quarterly*, 6(1&Auml;6):66–92, 1960. doi: 10.1002/malq.19600060105. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/malq.19600060105>.
- Alberto Camacho, Oscar Chen, Scott Sanner, and Sheila A. McIlraith. Non-markovian rewards expressed in LTL: guiding search via reward shaping. In *SOC*, pages 159–160, 2017a.
- Alberto Camacho, Oscar Chen, Scott Sanner, and Sheila A. McIlraith. Decision-making with non-markovian rewards: From ltl to automata-based reward shaping. In *RLDM*, pages 279–283, 2017b.



- Edmund M. Clarke, Jr., Orna Grumberg, and Doron A. Peled. *Model Checking*. MIT Press, Cambridge, MA, USA, 1999. ISBN 0-262-03270-8.
- Giuseppe De Giacomo and Moshe Y Vardi. Linear temporal logic and linear dynamic logic on finite traces. In *IJCAI*, volume 13, pages 854–860, 2013.
- Giuseppe De Giacomo and Moshe Y. Vardi. Synthesis for ltl and ldl on finite traces. In *Proceedings of the 24th International Conference on Artificial Intelligence*, IJCAI’15, pages 1558–1564. AAAI Press, 2015. ISBN 978-1-57735-738-4. URL <http://dl.acm.org/citation.cfm?id=2832415.2832466>.
- Giuseppe De Giacomo, Riccardo De Masellis, and Marco Montali. Reasoning on ltl on finite traces: Insensitivity to infiniteness. In *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence*, AAAI’14, pages 1027–1033. AAAI Press, 2014. URL <http://dl.acm.org/citation.cfm?id=2893873.2894033>.
- Sam Devlin and Daniel Kudenko. Dynamic potential-based reward shaping. In *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems - Volume 1*, AAMAS ’12, pages 433–440, Richland, SC, 2012. International Foundation for Autonomous Agents and Multiagent Systems. ISBN 0-9817381-1-7, 978-0-9817381-1-6. URL <http://dl.acm.org/citation.cfm?id=2343576.2343638>.
- Thomas G. Dietterich. The maxq method for hierarchical reinforcement learning. In *In Proceedings of the Fifteenth International Conference on Machine Learning*, pages 118–126. Morgan Kaufmann, 1998.
- Calvin C. Elgot. Decision problems of finite automata design and related arithmetics. *Transactions of the American Mathematical Society*, 98(1):21–51, 1961. ISSN 00029947. URL <http://www.jstor.org/stable/1993511>.
- Michael J. Fischer and Richard E. Ladner. Propositional dynamic logic of regular programs. *Journal of Computer and System Sciences*, 18(2):194 – 211, 1979. ISSN 0022-0000. doi: [https://doi.org/10.1016/0022-0000\(79\)90046-1](https://doi.org/10.1016/0022-0000(79)90046-1). URL <http://www.sciencedirect.com/science/article/pii/0022000079900461>.
- D. Gabbay, A. Pnueli, S. Shelah, and J. Stavi. On the temporal analysis of fairness. Technical report, Jerusalem, Israel, Israel, 1997.
- Valentin Goranko and Antony Galton. Temporal logic. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, winter 2015 edition, 2015.
- Charles Gretton. A more expressive behavioral logic for decision-theoretic planning. In *Pacific Rim International Conference on Artificial Intelligence*, pages 13–25. Springer, 2014.
- M. Grzes and D. Kudenko. Theoretical and empirical analysis of reward shaping in reinforcement learning. In *2009 International Conference on Machine Learning and Applications*, pages 337–344, Dec 2009. doi: 10.1109/ICMLA.2009.33.
- Marek Grzes. *Improving exploration in reinforcement learning through domain knowledge and parameter analysis*. PhD thesis, University of York, 2010.

- Marek Grześ. Reward shaping in episodic reinforcement learning. In *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems, AAMAS '17*, pages 565–573, Richland, SC, 2017. International Foundation for Autonomous Agents and Multiagent Systems. URL <http://dl.acm.org/citation.cfm?id=3091125.3091208>.
- John E. Hopcroft, Rajeev Motwani, Rotwani, and Jeffrey D. Ullman. *Introduction to Automata Theory, Languages and Computability*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2nd edition, 2000. ISBN 0201441241.
- Rodrigo Toro Icarte, Torny Q Klassen, Richard Valenzano, and Sheila A McIlraith. Teaching multiple tasks to an rl agent using ltl. 2018.
- Bakhadyr Khoussainov and Anil Nerode. *Automata Theory and Its Applications*. Birkhauser Boston, Inc., Secaucus, NJ, USA, 2001. ISBN 3764342072.
- Orna Kupferman and Moshe Y. Vardi. Model checking of safety properties. *Form. Methods Syst. Des.*, 19(3):291–314, October 2001. ISSN 0925-9856. doi: 10.1023/A:1011254632723. URL <https://doi.org/10.1023/A:1011254632723>.
- Bruno Lacerda, David Parker, and Nick Hawes. Optimal policy generation for partially satisfiable co-safe ltl specifications. In *Proceedings of the 24th International Conference on Artificial Intelligence, IJCAI'15*, pages 1587–1593. AAAI Press, 2015. ISBN 978-1-57735-738-4. URL <http://dl.acm.org/citation.cfm?id=2832415.2832470>.
- Andrew Y. Ng, Daishi Harada, and Stuart J. Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *Proceedings of the Sixteenth International Conference on Machine Learning, ICML '99*, pages 278–287, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc. ISBN 1-55860-612-2. URL <http://dl.acm.org/citation.cfm?id=645528.657613>.
- M. Pesic and W. M. P. van der Aalst. A declarative approach for flexible business processes management. In *Proceedings of the 2006 International Conference on Business Process Management Workshops, BPM'06*, pages 169–180, Berlin, Heidelberg, 2006. Springer-Verlag. ISBN 3-540-38444-8, 978-3-540-38444-1. doi: 10.1007/11837862\_18. URL [http://dx.doi.org/10.1007/11837862\\_18](http://dx.doi.org/10.1007/11837862_18).
- Amir Pnueli. The temporal logic of programs. In *Proceedings of the 18th Annual Symposium on Foundations of Computer Science, SFCS '77*, pages 46–57, Washington, DC, USA, 1977. IEEE Computer Society. doi: 10.1109/SFCS.1977.32. URL <https://doi.org/10.1109/SFCS.1977.32>.
- V. R. Pratt. Semantical consideration on floyo-hoare logic. In *17th Annual Symposium on Foundations of Computer Science (sfcs 1976)*, pages 109–121, Oct 1976. doi: 10.1109/SFCS.1976.27.
- M. O. Rabin and D. Scott. Finite automata and their decision problems. *IBM J. Res. Dev.*, 3(2):114–125, April 1959. ISSN 0018-8646. doi: 10.1147/rd.32.0114. URL <http://dx.doi.org/10.1147/rd.32.0114>.
- Stewart Shapiro and Teresa Kouri Kissel. Classical logic. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, spring 2018 edition, 2018.

- Satinder P. Singh and Richard S. Sutton. Reinforcement learning with replacing eligibility traces. *Mach. Learn.*, 22(1-3):123–158, January 1996. ISSN 0885-6125. doi: 10.1007/BF00114726. URL <http://dx.doi.org/10.1007/BF00114726>.
- A. P. Sistla and E. M. Clarke. The complexity of propositional linear temporal logics. *J. ACM*, 32(3):733–749, July 1985. ISSN 0004-5411. doi: 10.1145/3828.3837. URL <http://doi.acm.org/10.1145/3828.3837>.
- Richard S. Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*, 3(1):9–44, Aug 1988. ISSN 1573-0565. doi: 10.1007/BF00115009. URL <https://doi.org/10.1007/BF00115009>.
- Richard S. Sutton and Andrew G. Barto. *Introduction to Reinforcement Learning*. MIT Press, Cambridge, MA, USA, 1st edition, 1998. ISBN 0262193981.
- Sylvie Thiébaux, Charles Gretton, John K. Slaney, David Price, and Froduald Kabanza. Decision-theoretic planning with non-markovian rewards. *J. Artif. Intell. Res. (JAIR)*, 25:17–74, 2006.
- Wolfgang Thomas. Star-free regular sets of  $\bar{L}$ -sequences. *Information and Control*, 42(2):148 – 156, 1979. ISSN 0019-9958. doi: [https://doi.org/10.1016/S0019-9958\(79\)90629-6](https://doi.org/10.1016/S0019-9958(79)90629-6). URL <http://www.sciencedirect.com/science/article/pii/S0019995879906296>.
- B.A. Trakhtenbrot. Finite automata and the logic of single-place predicates. *Sov. Phys., Dokl.*, 6:753–755, 1961. ISSN 0038-5689.
- Christopher J. C. H. Watkins and Peter Dayan. Q-learning. *Machine Learning*, 8(3):279–292, May 1992. ISSN 1573-0565. doi: 10.1007/BF00992698. URL <https://doi.org/10.1007/BF00992698>.
- Christopher John Cornish Hellaby Watkins. *Learning from delayed rewards*. PhD thesis, 1989.
- Pierre Wolper. Temporal logic can be more expressive. *22nd Annual Symposium on Foundations of Computer Science (sfcs 1981)*, pages 340–348, 1981.