

# Reinforcement Learning

## For Radar Waveform Optimization

**Mario Coutino**, Faruk Uysal

Radar Technology Department, TNO, The Netherlands

San Antonio, TX, US | May 2023



# Outline

- **Radar and Radar Waveform Adaptivity**
- **Reinforcement Learning for Radar Adaptivity**
- **Waveform Design using Reinforcement Learning**
  - Setting description
  - Design Decisions
  - Illustrative Examples
- **Conclusion and Future Directions**



# Radar Waveform Adaptivity

Radar waveforms can be optimized/adapted for

- **improve** a radar task<sup>[1]</sup>, e.g., detection, tracking or classification
- **mitigate** clutter effects<sup>[2]</sup>
- **enable** coexistence of different sensing and communication platforms<sup>[3]</sup>

This is typically achieved by

- consulting a **look-up-table** :: fast but limited
  - dictionary of waveforms, e.g., responding to a mode –“tracking waveform” or a “classification waveform”
- solving an **optimization problem** :: flexible but demanding
  - given environmental inputs, a waveform is generated on-the-fly by solving an optimization procedure

**Besides searching for even faster optimization algorithms,  
is there other possibilities to have the best of both worlds?**

[1] Cui, Guolong, et al., eds. *Radar waveform design based on optimization theory*. SciTech Publishing, 2020.

[2] Hughes, Evan J., and Clive M. Alabaster. "Medium PRF radar PRF optimisation using evolutionary algorithms." *Proceedings of the 2003 IEEE Radar Conference (Cat. No. 03CH37474)*. IEEE, 2003.

[3] Aubry, Augusto, et al. "A new radar waveform design algorithm with improved feasibility for spectral coexistence." *IEEE Transactions on Aerospace and Electronic Systems* 51.2 (2015): 1029-1038.

**Alternative:**

**instead of per-case optimization**

**better to design a policy to**

**act given the environment's state.**

# Enabling Radar Waveform Adaptivity

## A policy



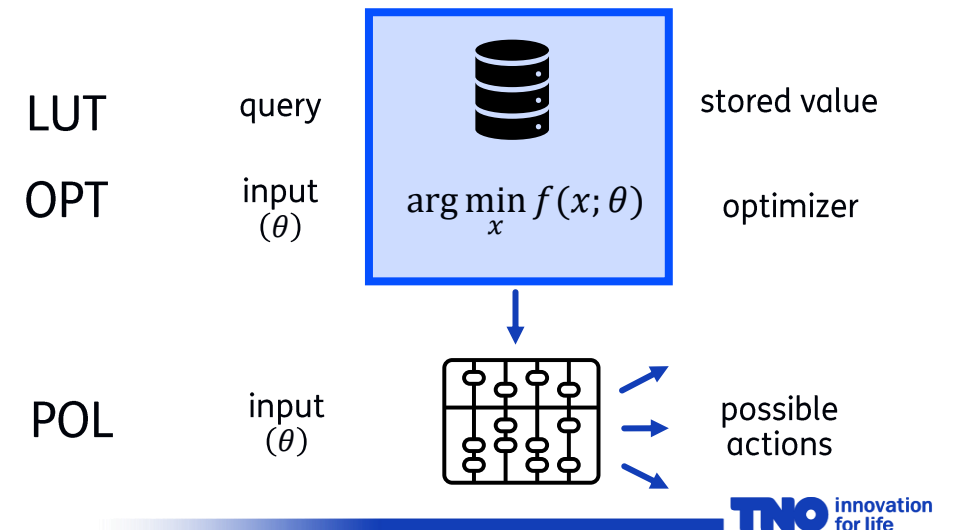
can be interpreted as:

as **mapping** representing a **compressed look-up table** capable of **solving** a family of **optimization problems**.

Thus, it can have

- predictable and low time requirements
- affordable memory footprint
- input-dependent behavior
- multiple responses to the same query

How can we learn a policy for radar waveform design?



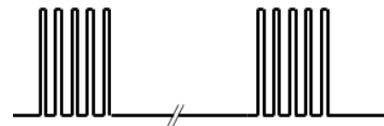


# Reinforcement Learning

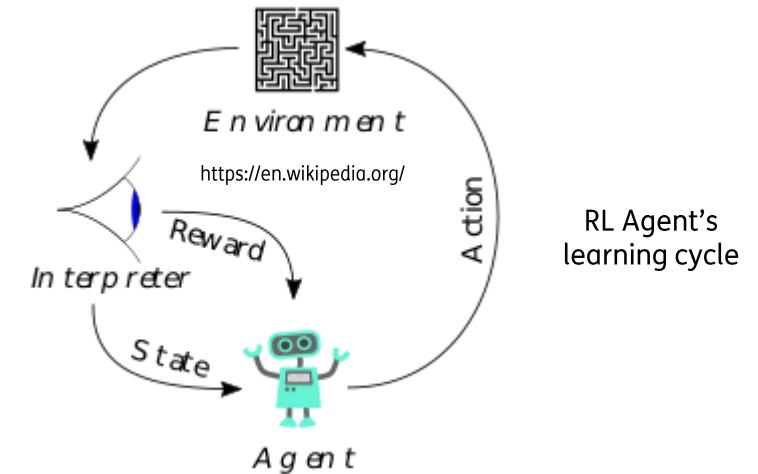
A way to **learn policies** for particular problems is by **reinforcement learning** (RL) algorithms.

Reinforcement learning<sup>[1]</sup>

- is a **feedback-based** technique aiming to **teach agents to maximize rewards** as it **interacts with an environment**.
- it has been **applied to several domains**, including radar, e.g.,
  - strategy games<sup>[2]</sup>
  - cognitive beamforming<sup>[3]</sup>
  - notched waveforms<sup>[4]</sup>
  - radar resource allocation<sup>[5]</sup>
  - etc..



**In this work, we study its application to burst design.**



[1] Sutton, Richard S., and Andrew G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

[2] <https://www.deepmind.com/research/highlighted-research/alphago>

[3] Ahmed, Aya Mostafa, et al. "A reinforcement learning based approach for multitarget detection in massive MIMO radar." *IEEE Transactions on Aerospace and Electronic Systems* 57.5 (2021): 2622-2636.

[4] Durst, Sebastian, and Stefan Brüggewirth. "Quality of service based radar resource management using deep reinforcement learning." *2021 IEEE Radar Conference (RadarConf21)*. IEEE, 2021.

[5] Smith, Graeme E., and Taylor J. Reininger. "Reinforcement learning for waveform design." *2021 IEEE Radar Conference (RadarConf21)*. IEEE, 2021.

# Reinforcement Learning

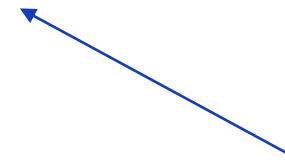
Most of RL starts by modelling a problem by a Markov decision process  $\mathcal{M}(\mathcal{S}, \mathcal{A}, T, R)$  (or any of its variants)

- $\mathcal{S}$  set of permissible environment states
- $\mathcal{A}$  set of available actions
- $T: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}_0^+$  probability of reaching an state  $s' \in \mathcal{S}$  after performing action  $a \in \mathcal{A}$  at state  $s \in \mathcal{S}$
- $R: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  *immediate* reward after an action  $a \in \mathcal{A}$  is taken at state  $s \in \mathcal{S}$

The **goal of RL is to learn a policy**  $\pi: \mathcal{A} \times \mathcal{S} \rightarrow [0,1]$  which **maximizes** the expected **cumulative reward**.

To make use of the RL machinery for waveform design

- the above components needs to be defined for the problem and
- an appropriate RL algorithm needs to be chosen

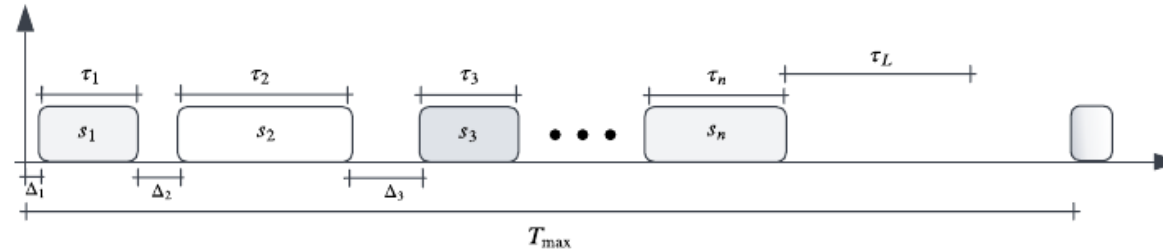


This by interacting with the environment



# Reinforcement Learning For Waveform Optimization

Let us consider the following setting for **burst-level waveform optimization**



- irregular but bounded **transmit intervals** possible  $:: \Delta_{\min} \leq \Delta_n \leq \Delta_{\max}$
- upper bound on **number of pulses** in burst  $:: n \leq N$
- irregular but bounded **pulse durations**  $:: \tau_{\min} \leq \tau_n \leq \tau_{\max}$
- **pulses modulations** selected from a dictionary  $:: \{s_1(t), \dots, s_K(t)\}$
- **total duration**, including listening time  $\tau_L$ , is upper bounded  $:: \sum_n (\tau_n + \Delta_n) + \tau_L \leq T_{\max}$

Constraints of the design “game”

Game over constraints  
Actionable constraints

At every “time instant”, the agent needs to select (take) a pulse configuration (action):

- **duration of pulse**  $:: \tau_n \in \mathcal{T} := [\tau_{\min}, \tau_{\max}]$
- **separation with respect previous pulse**  $:: \Delta_n \in \mathcal{D} := [\Delta_{\min}, \Delta_{\max}]$
- **modulation type**  $:: k \in \mathcal{K} := \{1 \dots, K\}$

Action Space

# Reinforcement Learning For Waveform Optimization

Given a base scenario(as before), **algorithm selection** comes next.

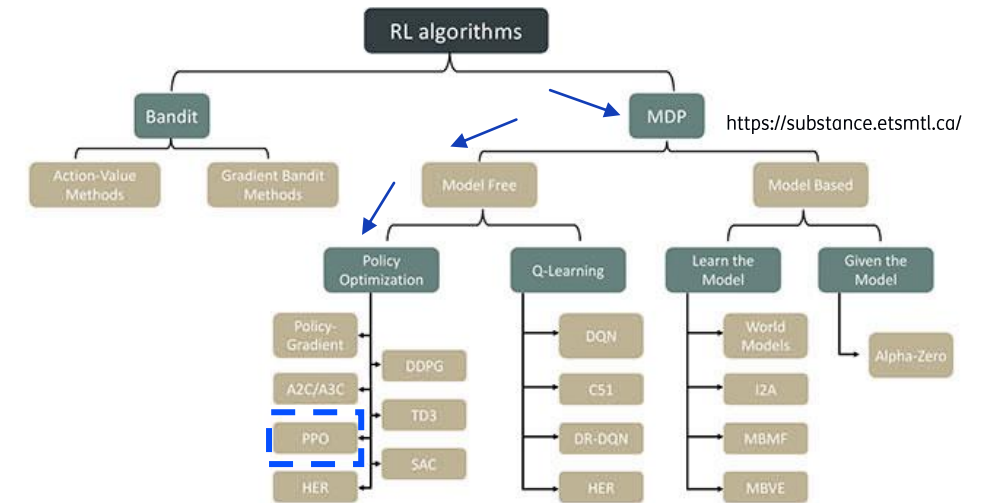
Based on what we want

- **learning a policy**

and what we have

- **mix of discrete and continuous** variables
- (so-far) **no model** of the environment

we can make a selection from the RL algorithm zoo<sup>[1]</sup>.



Here, we advocate\* the use of the **Proximal Policy Optimization (PPO)** algorithm.

- designed to achieve **stable learning**
- **sample-efficient** compared to other on-policy methods\*
- it has a simpler implementation and has “**few moving**” parts
- policies tend to be **fast to evaluate**

Quality of action

Parametrized policy

Expected quality of policy

$L(\theta) = \mathbb{E}_t [\pi_\theta(a_t | s_t) \hat{A}_t]$

- :: optimizes performance
- :: proximal optimization\*
- :: easier to tune
- :: neural networks with few parameters

[1] Part 2: Kinds of RL Algorithms — Spinning Up documentation (openai.com)

\* Details in publication.

# With this setup

**different games can be played...**

**it is not yet defined what the agent sees**  
(state description)

**and how their actions are scored.**  
(reward)

# Optimization Example

## Optimization of Maximum Sidelobe Level (MSL) under Duty Cycle (DC) Constraints

Example consists on designed a pulsed waveform with the objective of

- achieving a *particular duty cycle*  $DC_0$ , while
- *maintaining the MSL* of a region of interest  $\mathcal{R}$  below a desired level  $MSL_0$

**Reward function:**

$$R(n) = \begin{cases} 0 & \forall n : n < N \text{ and } T_{\text{design}} < T_{\text{max}} \\ -\text{FOM} & \text{otherwise} \end{cases}$$

$$\text{FOM} = \frac{|DC - DC_0|}{DC_0} + \max \left\{ \frac{MSL - MSL_0}{|MSL_0|}, 0 \right\}$$

Aims to match the DC goal

It does not promote waveforms with lower MSL than target.

Collection all previous actions



**State Space:**  $s_t = [s_{t-1}, a_{t-1}]$

Naïve but straightforward description of design

**Model:** Neural Network  
[32,16] w/ ReLus

**Frameworks:** pytorch<sup>[1]</sup>  
ray<sup>[2]</sup>

[1] <https://pytorch.org/>

[2] <https://www.anyscale.com/ray-open-source>

# Optimization Example

## Optimization of Maximum Sidelobe Level (MSL) under Duty Cycle (DC) Constraints

### Waveforms Generated with Optimized Policy

$$\mathcal{T} := \left[ \frac{1}{B}, \frac{T_{\max}}{15} \right]$$

$$\mathcal{D} := \left[ \frac{10}{B}, \frac{T_{\max}}{15} \right]$$

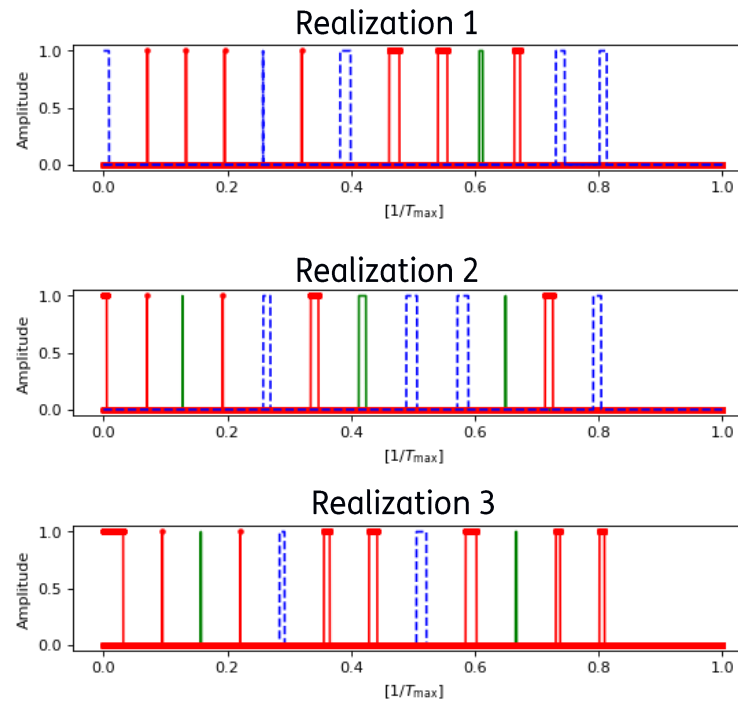
$$N = 32$$

$$DC_0 = 10\%$$

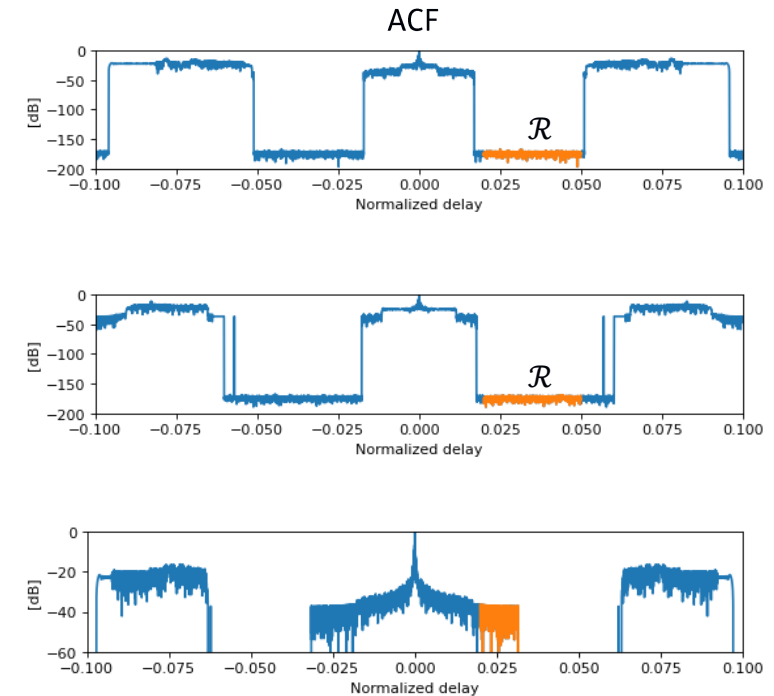
$$MSL_0 = -40\text{dB}$$

$$\mathcal{R} := [0.020, 0.05]$$

Strong dependency on time distribution  
(temporally signals looks similar)



Modulations type: (red) down- (blue) up- (green) up-down chirp



Small, but some diversity on the waveforms when the policy is sampled.

**Designing with respect to a**

**single configuration is not the killer application of RL\*.**

**Generation of waveforms for arbitrary inputs is.**

\* Elaborated in publication.



# Design Example

## Automatic Design with Control of Maximum Sidelobe Level under Duty Cycle Constraints

Instead of only optimizing for a single set of target parameters

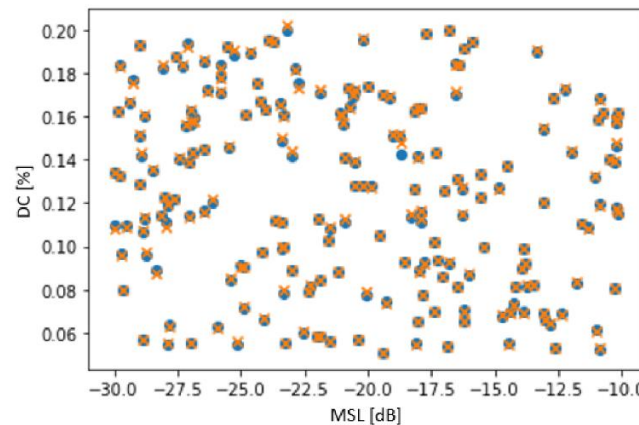
- the goals, MSL and DC, can be included in the state

**State Space:**  $s_t = [s_{t-1}, a_{t-1}]$ ;  $s_0 := [DC_0, MSL_0]$

Initial conditions

This will allow to obtain different waveforms by *querying* the agent's policy with *different initial conditions*.

- Target requirements (blue circles)
- Achieved parameters (orange crosses)



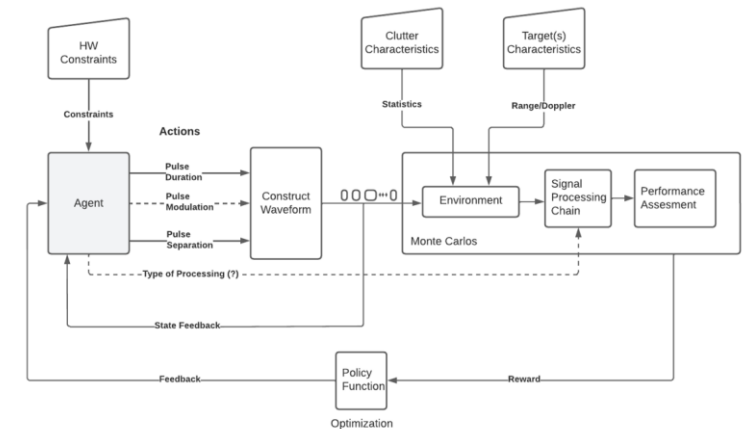
without making any other change in  
the environment setting

Policy evaluation tend to be fast  
(small models)

Random waveforms can be designed for new configurations once the policy has been learnt.

# Conclusions and Future Directions

- **Policy** learning using RL is a promising enable technology for **real-time radar adaptivity**
  - policies allow to have the **speed of LUTs** while being as **responsive as optimization problems**,
  - **structure knowledge** to gain experience –after optimization all that has been evaluated is thrown away
  - and **allow diversity** when the solution to the problem accepts it.
- **Proper** definition of environment and “**game**” **rules** are crucial
  - under the same setting, multiple “design games” can be played –change on reward or available observations
- **Waveform design using RL** was demonstrated with MSL/DC goals
  - several parameters were optimized at once **even for varying goals** –initial conditions
- **Extensions** of “games” to **complete scenarios** of practical interest
  - holistic view of, e.g., a radar processing chain
- Application of RL to **design “testers”** of systems –finding edge/corner cases
- Consideration of **Multi-agent RL** to
  - **alleviate sample complexity** for larger problems, e.g., dwell design.



# Thank you

[mario.coutinominguez@tno.nl](mailto:mario.coutinominguez@tno.nl)