
SetPoint Learning : Toward approach to control them all

JFE^{* 1} BTI¹ DOBEDO F. Furfaro^{* 1}

Abstract

Deep learning has enabled significant advancements in the field of dynamic system control, particularly in robotics, video games, navigation, and energy management, especially for non-linear systems. However, this type of learning does not allow for generalization to all dynamic systems, and sometimes, controllers based on system modeling (MPC, PID, MOR) are more relevant than methods based on stochastic data, especially for cases where the system is linearizable, and the target setpoint is well-known, which corresponds to most cases in the electronic, mechanical, and thermal domains. We introduce a new learning method that is generalizable to all linear dynamic systems where the setpoint is known and adaptive. We show that when targeting a setpoint, neural networks adapt towards an objective and would be a promising approach for the use of generative diffusion models for direct interaction with multiple real-world environments in parallel. This study has led to the creation of a set of simulation environments that can serve as a basis for other reinforcement or data-based approaches.

1 Introduction

The planning of decision-making in specific control environments is often a relatively simple problem to solve using data-driven approaches and reinforcement learning. In linearizable environments, planning can be reduced to a trajectory optimization to be achieved and can be done using controllers ([https://doi.org/10.1016/S0005-1098\(98\)00178-2](https://doi.org/10.1016/S0005-1098(98)00178-2)). In reinforcement learning, we observe results that surpass linear or data-based control. We can distinguish the DQN algorithms based on value and PPO based on policy.

^{*} Equal contribution. ¹ Capgemini Engineering, Paris. Correspondence to: fabien.furfaro@capgemini.com. Arxiv preprint, Paris, 2024. Copyright 2024 by the author(s).

The previous methods have worked well in specific situations, but they often require a lot of work to study the system and optimize the parameters for each case [Source]. This lack of ability to work in different situations has led to growing interest in approaches that can create agents that can adapt to different environments [Source]. In reinforcement learning, the ability of an agent to perform well in different environments is a big challenge. Agents can do well in specific training setups, but their performance drops a lot when faced with new situations.

For this, several alternatives have been proposed... [Aiding, Navigation with Human Feedloop, Planning diffusion, Google Genie, JEPa]. However, all of these latter techniques are complex and do not allow generalization to all systems.

In this work, we have developed a simple approach, which makes it possible to generalize the control of a linearizable system, as well as a set of control environments to adapt to any type of algorithm. learning. (...) the main hypothesis is that learning models always have the objective of achieving the best results, whereas here we assume that intermediate results where we aim for a lower score, allows us to expect better results as would a human with a trainer who lowers his level for his student. This approach is thus inspired in a certain way by the human learning approach used in educational psychology.

2 Background

Our approach is an analogue of past work in behavioral synthesis using trajectory optimization (MPC, Witkin Kass, 1988; Tassa et al., 2012). But we propose here an approach where the variables are always the same whatever the physical system and where the objective to be achieved is defined in the controller input. This approach is also compatible for the use of generative models.

2.1 Problem Setting

Consider a system governed by a discrete-time dynamics $S_{t+1} = f(S_t, a_t)$. The objective, whatever our system, will

Code and visualizations of the learned process are available at [setpoint-learning.github.io](https://github.com/setpoint-learning).

be to achieve the optimum:

$$u(k) = \operatorname{argmin}(J(s)) \quad (1)$$

The difficulty is to have a set of variables as input to a model which allow the resolution of this equation whatever the physical system.

2.2 Linear system

In the case of a linear system, we have:

$$y(t) = h_X(t) * s(t) \quad (2)$$

Which $s(t)$ is an action following the controlling system and $h_X(t)$ the randomized linear system.

3 Methods: Setpoint Learning

When we summarize a dynamic system in its trajectory, we observe four fundamental sets of variables which allow the resolution of equation (1) whatever the linear system (2). These variables are:

- a_{t-1} : Action in previous state s_{t-1} doing next state s_t
- s_{t-1} : Previous state
- s_t : Actual state of system
- $(s_c)_{t+1}$: Objective/Target state

In a linear system, the response instantly goes in the direction of the action given at the system input (we don't have delay response).

$$(S, A) : a_{(i,t)} \rightarrow \{a_{(i,t-1)}, s_{(i,t-1)}, s_{(i,t)}, s_{c(i,t+1)}\} \quad (3)$$

From this assestion, we can make the hypothesis that only this variable is enough to know what action to take to achieve an objective at the next moment. We can define naive binary controler :

$$a_t = \operatorname{sgn}(a_{t-1}) \cdot \operatorname{sgn}(s_t - s_{t-1}) \cdot \operatorname{sgn}(s_{c,t+1} - s_t) \quad (4)$$

The consequence is that if a learning algorithm has several types of linear and invariant environments as data, it can generalize to any linearizable system and adapt to any level rather than just a specific value.

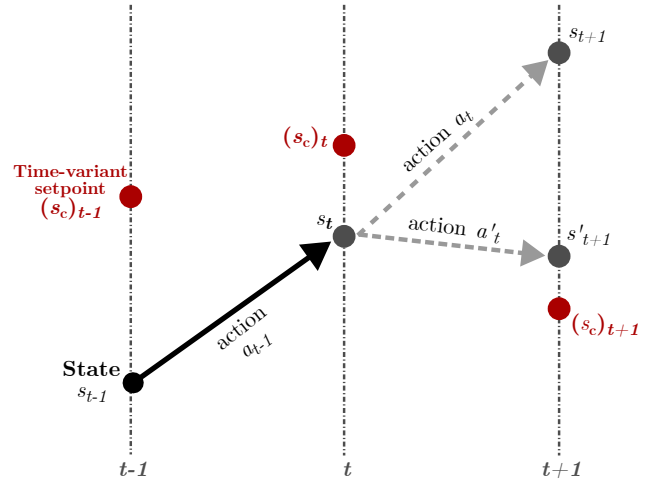


Figure 1. Representative diagram of the target setpoint concept. At time t , the agent in state s_t should choose the action in order to get as close as possible to the target state (setpoint) (s_c) at time $t + 1$.

3.1 General case

Cas général (fleche)

This relation provides a straightforward translation between classifier-guided sampling, used to generate classconditional data

Algorithm 1 Environment Planning

```

1: Require ABCD  $M_X$ , guide  $\mathcal{J}$ , scale  $\alpha$ , covariances  $\Sigma^i$ 
2: while not done do
3:   Observe state  $s$ ; initialize plan  $\tau^N \sim \mathcal{N}(\mathbf{0}, I)$ 
4:   for  $i = N, \dots, 1$  do
5:     // parameters of reverse transition
6:      $\mu \leftarrow \mu_\theta(\tau^i)$ 
7:     // guide using gradients of return
8:      $\tau^{i-1} \sim \mathcal{N}(\mu + \alpha \Sigma \nabla \mathcal{J}(\mu), \Sigma^i)$ 
9:     // constrain first state of plan
10:     $\tau_{s_0}^{i-1} \leftarrow s$ 
11:   Execute first action of plan  $\tau_{a_0}^0$ 
    
```

to be defined

3.2 Reinforcement learning case

Cas renforcement (environnemnt)

Attention entre le point de vue de l'environnement et le point de vue du modele d'apprentissage !

Pour le modele, S_t correspond à la sortie instantané, mais lorsque l'action est faite, ca devient S_{t+1} pour l'environnement !

L'idée est, sachant que je suis à l'etat " t ", quel action pour aller à l'etat " $t+1$ " correspondant à S_c , mais une fois l'action faite, ca devient S_t pour l'environnement

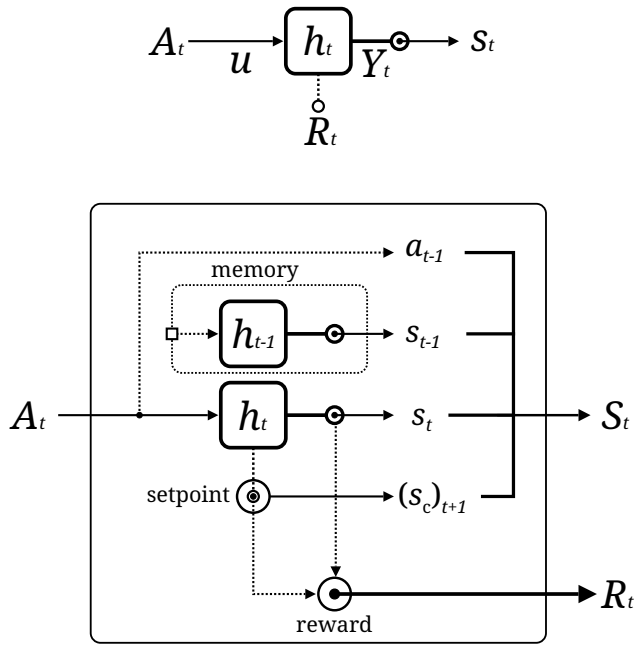


Figure 2. Algorithm scheme of RL-Approach. In A, you have minimal RL system exprimed with transfert fonction h_t , it's can be a linearized system, or wrapped environment. In B, the algorithm scheme of environment with reward calculation. AMBIGUITEE SUR At et at-1!

3.3 Diffusion learning case

Cas diffusion (image-multi-env)

Attention ici, comme c'est le resultat de la sortie, At donne St+1p.

4 Results

We are the

4.1 CartPole

The objective is to test our approach on a simple use case. We chose CartPole because it is a small angle linear approximation system that also includes a response delay.

4.2 Linear Time-invariant system

LTI

4.3 Multiple Linear Time-invariant system

multi-LTI

5 Conclusions and perspective

lien vers le code

Utilisation de l'environnement construit dans ce papier pour

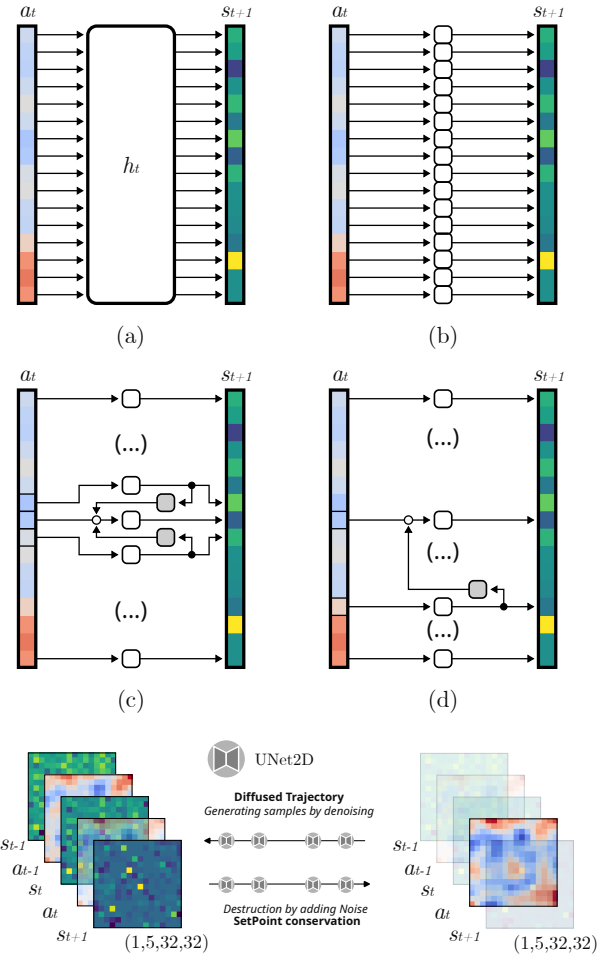


Figure 3. Algorithm scheme of DDPMS-Approach. In A, you have minimal RL system exprimed with transfert fonction h_t , it's can be a linearized system, or wrapped environment. In B, the algorithm scheme of environment with reward calculation.

construire des modele fondation avec plus d'hyperparametre

Utilisation de control-Net pour generer des actions plus complexe (<https://ctrl-adapter.github.io/>)