



HAL
open science

Intrinsically Motivated Goal Exploration Processes with Automatic Curriculum Learning

Sébastien Forestier, Yoan Mollard, Pierre-Yves Oudeyer

► **To cite this version:**

Sébastien Forestier, Yoan Mollard, Pierre-Yves Oudeyer. Intrinsically Motivated Goal Exploration Processes with Automatic Curriculum Learning. 2017. hal-01651233v1

HAL Id: hal-01651233

<https://hal.science/hal-01651233v1>

Preprint submitted on 28 Nov 2017 (v1), last revised 9 Dec 2022 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Intrinsically Motivated Goal Exploration Processes with Automatic Curriculum Learning

Sébastien Forestier^{1,2,3}
sebastien.forestier@inria.fr

Yoan Mollard^{2,3}
yoan.mollard@inria.fr

Pierre-Yves Oudeyer^{2,3}
pierre-yves.oudeyer@inria.fr

¹Université de Bordeaux, ²Inria Bordeaux Sud-Ouest, ³Ensta-ParisTech, France

Abstract

Intrinsically motivated spontaneous exploration is a key enabler of autonomous lifelong learning in human children. It allows them to discover and acquire large repertoires of skills through self-generation, self-selection, self-ordering and self-experimentation of learning goals. We present the unsupervised multi-goal reinforcement learning formal framework as well as an algorithmic approach called intrinsically motivated goal exploration processes (IMGEP) to enable similar properties of autonomous learning in machines. The IMGEP algorithmic architecture relies on several principles: 1) self-generation of goals as parameterized reinforcement learning problems; 2) selection of goals based on intrinsic rewards; 3) exploration with parameterized time-bounded policies and fast incremental goal-parameterized policy search; 4) systematic reuse of information acquired when targeting a goal for improving other goals. We present a particularly efficient form of IMGEP that uses a modular representation of goal spaces as well as intrinsic rewards based on learning progress. We show how IMGEPs automatically generate a learning curriculum within an experimental setup where a real humanoid robot can explore multiple spaces of goals with several hundred continuous dimensions. While no particular target goal is provided to the system beforehand, this curriculum allows the discovery of skills of increasing complexity, that act as stepping stone for learning more complex skills (like nested tool use). We show that learning several spaces of diverse problems can be more efficient for learning complex skills than only trying to directly learn these complex skills. We illustrate the computational efficiency of IMGEPs as these robotic experiments use a simple memory-based low-level policy representations and search algorithm, enabling the whole system to learn online and incrementally on a Raspberry Pi 3.

Keywords: intrinsically motivated exploration; unsupervised multi-goal reinforcement learning; intrinsic motivation; curiosity-driven learning; automatic generation of goals; curriculum learning; learning progress; robotics; modular representations

1 Introduction

An extraordinary property of natural intelligence in children is their capacity for lifelong autonomous learning. Processes of autonomous learning in infants have several properties that are fundamentally different from many current machine learning systems. Among them is the capability to spontaneously explore their environments, driven by an intrinsic motivation to discover and learn new tasks and problems that they imagine and select by themselves [Berlyne, 1966, Gopnik et al., 1999]. Crucially, there is no engineer externally imposing one target goal that they should learn, or hand providing a curriculum for learning, or providing a ready-to-use database of training examples. Rather, children self-select their objectives within a large, potentially open-ended, space of goals they can imagine, and they collect training data by physically practicing these goals. In particular, they explore goals in an organized manner, attributing to them values of interestingness that evolve with time, and allowing them to self-define a learning curriculum that is called a developmental trajectory in developmental sciences [Thelen and Smith, 1996]. This self-generated learning curriculum allows infants to avoid spending too much time on goals that are either too easy or too difficult, focusing on goals of the right level of complexity at the right time. Within this process, the new learned goals/problems are often stepping stones for discovering how to solve other goals of increasing complexity. Thus, while they are not explicitly guided by a final target goal, these mechanisms allow infants to discover highly complex skills such as biped locomotion or tool use, which would have been extremely difficult to learn if they would have only focused on these goals from the start as the rewards for these goals are typically rare or deceptive.

An essential component of such organized spontaneous exploration is the intrinsic motivation system, also called curiosity-driven exploration system [Gottlieb et al., 2013]. In the last decade, a series of computational and robotic models of intrinsically motivated exploration and learning in infants have been developed [Oudeyer and Kaplan, 2007, Baldassarre and Mirolli, 2013], opening new theoretical perspectives in neuroscience and psychology [Gottlieb et al., 2013]. Two key ideas have allowed to simulate and predict important properties of infant spontaneous exploration, ranging from vocal development [Moulin-Frier et al., 2014, Forestier and Oudeyer, 2017], to object affordance and tool learning [Forestier and Oudeyer, 2016a, Forestier and Oudeyer, 2016c]. The first key idea is that infants might select experiments that maximize an intrinsic reward based on empirical learning progress [Oudeyer et al., 2007], which generates automatically developmental trajectories (e.g. learning curricula) where progressively more complex tasks are practiced, learned and used as stepping stones for more complex skills. The second key idea is that beyond selecting actions or states based on the predictive learning progress they provide, a more powerful way to organize intrinsically motivated exploration is to select goals, i.e. self-generated reinforcement learning problems, based on a measure of control learning progress [Baranes and Oudeyer, 2013]. Here, the intrinsic reward is the empirical improvement towards solving these problems/goals [Oudeyer and Kaplan, 2007, Forestier and Oudeyer, 2016a], happening through lower-level policy search mechanisms that generate physical actions. The efficiency of such goal exploration processes leverages the fact that the data collected when targeting a goal can be informative to find better solutions to other goals (for example, a learner trying to achieve the goal of pushing an object on the right but actually pushing it on the left fails to progress on this goal, but learns as a side effect how to push it on the left).

Beyond neuroscience and psychology, we believe these models open new perspectives in artificial intelligence. In particular, algorithmic architectures for intrinsically motivated goal exploration were shown to allow the efficient acquisition of repertoires of high-dimensional motor skills with automated curriculum learning in several robotics experiments [Baranes and Oudeyer, 2013, Forestier and Oudeyer, 2016a]. This includes for example learning omnidirectional locomotion or learning multiple ways to manipulate a single complex flexible object [Baranes and Oudeyer, 2013]. In this article, we first present a formal framework called “Unsupervised Multi-Goal Reinforcement Learning”, as well as a formalization of intrinsically motivated goal exploration processes, that is both more compact and more general than these previous models. Then, we present an experimentation of implementations of these processes in a complex robotic setup with multiple objects, associated to multiple spaces of parameterized reinforcement learning problems, and where the robot can learn how to use certain objects as tools to manipulate other objects. We analyze and discuss how curriculum learning is automated in this unsupervised multi-goal exploration process, and compare the trajectory of exploration and learning of these spaces of problems with the one generated by other mechanisms such as hand-designed learning curriculum, or exploration targeting a single space of problems, and random motor exploration.

2 Formal framework

2.1 The Unsupervised Multi-Goal Reinforcement Learning Problem

Let's consider a machine that can produce behaviours by executing stochastic **policies** π_θ parameterized by $\theta \in \Theta \subset \mathbb{R}^{d_\Theta}$ in **context** $c \in \mathcal{C} \subset \mathbb{R}^{d_c}$, characterizing the current state of an environment Env and generated stochastically by the distribution $\mu_{Env}(c)$. We consider here policies $\pi_\theta(\mathbf{a}_{t+1} \mid s_{t_0:t+1}, a_{t_0:t})$ as stochastic black-boxes that produce time-bounded behavioural trajectories, denoted $\tau = \{s_{t_0}(=c), a_{t_0}, s_{t_1}, a_{t_1}, \dots, s_{t_{end}}, a_{t_{end}}\} \in \mathbb{T}$, through the dynamics $\nu_{Env}(\tau \mid \theta, c)$ of the environment. We only assume that they terminate before time $t_0 + T$ where t_0 is the time when a policy is started and T is a timeout used by a meta-controller to stop the policy if this timeout is reached (and we note t_{end} the time when a policy actually terminates).

We denote $o_\tau \in \mathcal{O} \subset \mathbb{R}^{d_o}$, and call **outcome**, a vector of measurements (also called descriptors or features) characterizing the behaviour of the machine and of the environment during one execution of the policy π_θ until it terminates and starting from context c . The descriptors in the outcome vector o characterize properties of the behavioural trajectory τ .

Finally, we assume that the machine is capable to sample goals (here also called problems) in a space of reinforcement learning **problems** \mathcal{P} parameterized by $p \in \mathcal{P} \subset \mathbb{R}^{d_p}$, and where each goal/problem is defined by a **goal-parameterized reward function** $R : (p, c, \theta, o_\tau) \mapsto r$ that specifies the reward received by the agent if it tries to solve goal/problem p in a context c with policy π_θ and observes outcome o_τ . Our main assumption is that given a context c , a policy π_θ and an observed outcome o_τ , the agent can compute $R(p, c, \theta, o_\tau)$ for any $p \in \mathcal{P}$. In other words, given a context, policy and outcome, the agent receives a partial reward function $R_{c,\theta,o_\tau} : p \mapsto r$ so that it can compute the reward (or fitness) r of experiment (c, θ, o_τ) for solving any problem $p \in \mathcal{P}$, with $r = R_{c,\theta,o_\tau}(p) = R(p, c, \theta, o_\tau)$. Another way to understand those rewards is that given a goal problem p to solve, the agent can compute a specific reward function R_p giving the reward/fitness r of all previous experiments (c, θ, o_τ) for solving the particular problem p , with $r = R_p(c, \theta, o_\tau) = R(p, c, \theta, o_\tau)$. Thus in this framework, sampling a goal/problem p is equivalent to sampling a reward function R_p .

An example below illustrates the generality of this form of goals, enabling to express complex objectives that do not simply depend on the observation of the end state policies, but might depend on several aspects of entire behavioural trajectories. Importantly, this framework does not assume that all goals/problems are solvable. For example, in the robotics experiments below, the representation of goal spaces is provided by the engineer¹, but large parts of these goal spaces are not solvable and the machine does not know beforehand which goals are solvable or their relative difficulty.

Given these spaces \mathcal{P} of problems, \mathcal{C} of contexts, Θ of policies, as well as an environment, a machine can explore the environment's dynamics and learn a stochastic **meta-policy** $\Pi(\theta \mid p, c)$ to try to solve as well as possible all problems p in all contexts c . In order to evaluate how well the machine learned to solve diverse problems in diverse contexts, we define a test distribution $\mathcal{T}(p, c)$ of problems and contexts, unknown to the machine during learning, on which we can evaluate the rewards obtained by the machine (a particular case is when there is only a single target problem considered by the experimenter). The **loss of a meta-policy** Π is defined by Equation 1 as minus the expected reward received by following Π to solve all problems in all contexts following the test distribution $\mathcal{T}(p, c)$. The loss cannot be computed by the learning machine as it does not know the test distribution $\mathcal{T}(p, c)$.

$$\mathcal{L}(\Pi) = - \int_{\mathcal{P}} \int_{\mathcal{C}} \mathcal{T}(p, c) \int_{\Theta} \Pi(\theta \mid p, c) \int_{\mathbb{T}} \nu_{Env}(\tau \mid \theta, c) R(p, c, \theta, o_\tau) d\tau d\theta dc dp \quad (1)$$

We define the **Unsupervised Multi-Goal Reinforcement Learning Problem** as the problem of learning a meta-policy Π with minimal loss in a minimal amount of experiments (executions of policies in the environment). We use the term "unsupervised" because the tests problems are unknown to the machine during learning.

¹Actually, the engineer provides mathematical operators to the machine to algorithmically generate parameterized reward functions over measurements/outcomes of behaviours. An interesting avenue for future work is to consider representation learning algorithms for learning such parameterization of goal spaces.

2.2 Particular case: a modular problem space

An interesting particular case of the above formulation is when the problem space is the union of m spaces: $\mathcal{P} = \bigcup_{k=1}^m \mathcal{P}^k$. In that case, the reward for a problem $p \in \mathcal{P}^k$ is denoted $R_p^k(c, \theta, o_\tau)$.

After executing a policy π_θ in a context c , the agent observes outcome o_τ and it can compute $R_p^k(c, \theta, o_\tau) = R(p, c, \theta, o_\tau)$ for all problems p in all problem spaces \mathcal{P}^k .

In the definition of the test distribution $\mathcal{T}(p, c)$, we can be interested in evaluating the competence of the learning agent in one particular problem space \mathcal{P}^k , in which case $\mathcal{T}(p, c)$ will sample problems $p \in \mathcal{P}^k$, or in all spaces, with p uniform in \mathcal{P} .

2.3 Links to Reinforcement Learning

In this setting, the reward functions of the goal-parameterized space of functions R_p with $p \in \mathcal{P}$ have two particularities in comparison with the concept of “reward function” as often used in the RL literature.

The first particularity is that these reward functions are computed based on the outcome o_τ of executing a policy π_θ in context c , and thus consider the whole behaviour of the machine and of the environment during the execution, so these rewards are not Markovian if one considers them from the perspective of the lower level of state transitions associated to the execution of π_θ , i.e. the (s_t, a_t) .

The second particularity is that since the computation of the reward $R(p, c, \theta, o_\tau)$ is internal to the machine, it can be computed any time after the experiment (c, θ, o_τ) and for any problem $p \in \mathcal{P}$, not only the particular problem that the agent was trying to solve. Consequently, if the machine experiments a policy π_θ in context c and observes o_τ (e.g. trying to solve problem p_1), and stores the results (c, θ, o_τ) of this experiment in its memory, then when later on it self-generates problems p_2, p_3, \dots, p_i it can compute on the fly (and without new actual actions in the environment) the associated rewards $R_{p_2}(c, \theta, o_\tau), R_{p_3}(c, \theta, o_\tau), \dots, R_{p_i}(c, \theta, o_\tau)$ and use this information to improve over these goals p_2, p_3, \dots, p_i . This property is essential as this enables direct reuse of data collected when learning to solve a problem for solving later on other problems, and is leveraged for curriculum learning in intrinsically motivated goal exploration processes (see below).

2.4 Example of an Unsupervised Multi-Goal Reinforcement Learning problem

An autonomous robot playing with a ball. Consider a wheeled robot that can move around in an environment with walls and with a ball that it can push. The sensory system of the robot allows it to measure and encode in a vector s_t the current position and speed of itself and of the ball as well as distances to walls in several directions. The robot can produce behaviours by executing policies π_θ encoded as recurrent neural network with weights parameterized by θ . The network takes as input the current value of sensor measurements (s_t), and outputs motor speed values (encoded as a vector a_t) during one roll-out (roll-outs are automatically stopped by a meta-controller after s seconds). The starting context c of a roll-out (an experiment) is a vector encoding the current position and speed of the robot and of the ball (so the context equals the starting sensor measurements in this case). While the robot is executing a movement from time t_0 to time t_{end} , it records all the sequence τ of motor commands and observed sensor measurements $(s_0, a_0), \dots, (s_{t_{end}}, a_{t_{end}})$. From the data in τ , the robot computes the outcome o of executing π_θ in c as a set of descriptors/features of τ . Let’s consider one example where those descriptors are $o = (d1, d2, d3)$ where $d1$ is a vector encoding the translation of the ball between time t_0 and t_{end} , $d2$ the minimal distance to the walls during the trajectory, and $d3$ the energy used by the robot for generating the trajectory (computed from speeds). Now, let’s consider a space of problems that the robot could sample from, and where each problem/goal consists in trying to find a behaviour π_θ in context c such that the ball translates with vector d_g while maximizing the minimal distance to the wall and minimizing the energy spent, with relative weights α and β for these components of the problem. Thus, this space of problems is parameterized by $g = (d_g, \alpha, \beta)$ such that

$$R_g(c, \theta, o) = \alpha e^{-\|d_g - d1\|^2} + \beta d2 + (1 - \alpha - \beta) e^{-d3^2}$$

For example, the problem $([1, 1], 0, 1)$ is the problem of staying away from the walls, and problem $([1, 1], 0.5, 0)$ is the problem of translating the ball with vector $[1, 1]$ while minimizing the energy spent, but with no constraint on the distance to the walls.

A more simple case to consider is if outcomes would be simply the translation of the ball during a roll-out ($o = d1$), and if problems/goals would consist in target translations ($g = d_g$). In that case, a goal is literally a target configuration of the environment at the end of the roll-out, and the parameterized rewards function $R_g(c, \theta, o)$ can directly be implemented as a function of a distance between the target translation d_g and the observed translation $d1$ when executing π_θ in context c : $R_g(c, \theta, o) = e^{-\|d_g - d1\|^2}$. In that case, one would not need to have formal space of parameterized problems, as it would be isomorphic to the space of outcomes. However, considering such a parameterized space of problems allows to sample and learn to solve families of more general problems such as in the example above.

In these two cases (complex or simple problem space), the Unsupervised Multi-Goal Reinforcement Learning problem consists in exploring the environment to collect learning data (c_i, θ_i, o_i) that are efficient for learning a good inverse model $\Pi(\theta | g, c)$ over the space of goals and contexts. One approach that could be taken to perform such an exploration is to perform experiments by sampling values of θ given c , then executing π_θ and observing the corresponding outcome. Such sampling can be random, or it can be active using one the many active learning methods in the literature developed to learn actively forward models of the world dynamics [Oudeyer et al., 2007]. However, as argued for example in [Baranes and Oudeyer, 2013], there are many real world situations where very large areas of the space of θ values are either producing no meaningful outcomes or are redundant (they all produce the same outcome). For example, most θ parameters of the robot above might produce a movement of the robot that does not even touch the ball. As a consequence, only certain manifolds of the θ space may be useful to learn to produce a certain diversity of outcomes.

For this reason, another approach, proven to be more efficient in previous experimental studies and formalized in Section 3.1 is Intrinsically Motivated Goal Exploration Processes (IMGEPs). The idea of IMGEPs is to explore the environment by repeatedly sampling a goal g in the problem space, then use the current inverse model to find the best predicted θ to reach goal g in the current context c , and use an optimization algorithm with a constrained time budget to improve over this best current solution. Interestingly, the resulting experiment(s) produce learning data (c_i, θ_i, o_i) that can improve the current best solutions to other problems/goals (even if they do not currently improve the solution to produce g). This enables transfer learning among goals/problems, and enables the system to be robust to sampling goals that might be very difficult to learn initially, or simply impossible to learn (e.g. the space of goals can allow the above robot to sample a goal where the target ball translation is $(100, 100)$ while the ball can physically move only inside a square of $10 * 10$). Furthermore, this also provides a framework in which goals can be sampled actively, for example by sampling more often goals for which the system is currently making most learning progress. Below, we present a formalization of IMGEPs and show in experiments that they generate a learning curriculum where goals are progressively and automatically sampled and learned in order of increasing complexity.

2.5 Exploration Problem versus Inverse Model Learning Problem

There are several technical challenges that need to be addressed to learn the meta-policy $\Pi(\theta | p, c)$. If one considers that a database of exemplars $L = (c_i, \theta_i, o_i)$ is available to the learner initially, a challenge is to devise learning algorithms that can detect and leverage regularities in this database to learn a $\Pi(\theta | p, c)$ such that given a new test dataset (c_j, p_j) , Π will generate policy parameters $\theta_j \sim \Pi(\theta | p_j, c_j)$ such that $R(p_j, c_j, \theta_j, o_j)$ will be maximal in average over the problems p_j in this test dataset (in corresponding contexts c_j). As Π is encoding a probabilistic inverse model from the spaces of problems and contexts to the space of policies, we call here the problem of learning Π from a database L the problem of inverse model learning (various approaches for contextual inverse model learning can be used, see for e.g. [Baranes and Oudeyer, 2013]). Here, an evaluation of the performance of such an inverse model learning algorithm is the average reward over the problems in the test dataset at the end of learning, possibly considering the speed of convergence of the learning process.

Now, when one considers the framework of autonomous learning, one cannot assume that a database of learning exemplars is already available to the learner: the learner needs to explore by itself to collect this database, from which the meta-policy Π can be learned either concurrently and incrementally

or later on in a batch manner. In this article, we focus on this second problem, which we believe is fundamental for autonomous learning architectures, and will present a family of algorithmic processes for exploration called “Intrinsically Motivated Goal Exploration Processes” that aims at collecting learning exemplars that have properties of diversity suited to learn good inverse models over the space of problems \mathcal{P} .

However, as these algorithmic processes require an inverse modeling step, our experiments will implement simple (yet powerful) algorithms for incremental inverse model learning. Furthermore, the evaluation of exploration processes considers the relation between the number N of roll-outs of policies π_θ in the environment and the quality of a learned meta-policy Π implementing an inverse model. Efficient exploration processes will be processes that allow to learn a high-quality Π (in terms of average reward over an independent test dataset of problems and contexts) in a minimal number N of roll-outs in the environment.

In environments where an infinity of problems/goals of unknown difficulty can be explored, and where reward information to some problems can be sparse, it becomes challenging to decide which action policies to explore next. The next section presents goal exploration processes as an approach to organize exploration in such contexts.

3 Intrinsically Motivated Goal Exploration Processes

We here present intrinsically motivated goal exploration processes as an algorithmic architecture that aims to address the unsupervised multi-goal reinforcement learning problem. This algorithmic architecture (Architecture 1) can be instantiated into many particular algorithms as it contains slots where several low-level policy learning algorithms can be used (see Algo. 3 and Appendix A for several examples). However, all these algorithmic implementations share several general principles which define the concept of intrinsically motivated goal exploration processes:

- The learning machine can self-generate and sample goals as abstract parameterized reinforcement learning problems; These goals are typically defined as complex time-extended functions of behavioural outcomes; In the modular case, there are several goal spaces to sample from;
- Goal sampling (and goal space sampling) is made using intrinsic rewards; A particularly efficient intrinsic measure is based on estimating empirical learning progress towards self-generated goals (interesting goals are those for which learning progresses faster, which enables automatic curriculum learning);
- Two learning loops are running in parallel: 1) an exploration loop uses parameterized time-bounded policies and fast incremental goal-parameterized policy search, enabling fast efficient search of good solutions over diverse sets of problems; 2) an exploitation loop uses data collected by the exploration loop to consolidate the learning process, possibly using slow batch learning algorithms;
- Both learning loops systematically reuse information acquired when targeting a goal for improving the solution to other goals (thus even if many goals can correspond to RL problems with sparse rewards, there can be a dense feedback from the environment for learning about multiple goals).

3.1 Algorithmic approach: Goal Exploration and Intrinsic Rewards

In this section we define a general algorithmic architecture to solve the Unsupervised Multi-Goal Reinforcement Learning problem with Intrinsically Motivated Goal Exploration Processes (see Architecture 1). It is an architecture in the sense that it is composed of several modules working together but for which multiple implementations are possible. We give several variants of implementations of this architecture in Section 4.3.

Architecture 1 Intrinsically Motivated Goal Exploration Process (IMGEP)

Require: Environment Env , Context space \mathcal{C} with distribution $\mu_{Env}(c)$

Require: Policy parameter space Θ , Outcome space \mathcal{O} , trajectory distribution $\nu_{Env}(\tau | \theta, c)$

Require: Problem parameter space \mathcal{P} , goal-parameterized reward function $R(p, c, \theta, o_\tau)$

Require: Initial knowledge \mathcal{E} \triangleright typically empty

- 1: $\Pi(\theta | p, c) \leftarrow \text{InitializeMetaPolicy}(\mathcal{E})$ \triangleright This is the target meta-policy (inverse model)
 - 2: $\Pi_\epsilon(\theta | p, c) \leftarrow \text{InitializeExplorationMetaPolicy}(\mathcal{E})$ \triangleright Π_ϵ is the inverse model used during exploration. It may be the same or different from Π
 - 3: $\gamma(g | c) \leftarrow \text{InitializeGoalPolicy}(\mathcal{E})$
 - 4: Launch asynchronously the two following loops (exploration and training of target meta-policy)
 - 5: **loop** \triangleright Exploration loop
 - 6: $c \sim \mu_{Env}(c)$ \triangleright Observe context c
 - 7: $g \sim \gamma(g | c)$ \triangleright Choose goal g in \mathcal{P} based on intrinsic rewards, e.g. using a MAB
 - 8: $\theta \sim \Pi_\epsilon(\theta | g, c)$ \triangleright Infer policy parameters θ that maximize expected $R_g(c, \theta, o_\tau)$ using an exploration/exploitation trade-off, and possibly computing now the problem-specific reward function R_g giving the reward/fitness r of all previous experiments (c_j, θ_j, o_j) to solve current goal problem g , with $r = R_g(c_j, \theta_j, o_j) = R(g, c_j, \theta_j, o_j)$
 - 9: $\tau \sim \nu_{Env}(\tau | \theta, c)$ \triangleright Execute a roll-out of π_θ , observe trajectory $\{s_{t_0:t_{end}}, a_{t_0:t_{end}}\}$
 - 10: Compute outcome o_τ from trajectory τ
 \triangleright The goal-parameterized reward function R_{c, θ, o_τ} can now be computed to find the reward/fitness r of the current experiment (c, θ, o_τ) for solving any problem $p \in \mathcal{P}$, with $r = R_{c, \theta, o_\tau}(p) = R(p, c, \theta, o_\tau)$
 - 11: $r = R_{c, \theta, o_\tau}(g)$ \triangleright Compute current reward associated to solving goal problem g
 - 12: $r_i \leftarrow IR(\mathcal{E}, c, g, \theta, o_\tau, r)$ \triangleright Compute intrinsic reward r_i associated to g in context c (e.g. learning progress in the vicinity of this goal)
 - 13: $\Pi_\epsilon(\theta | p, c) \leftarrow \text{UpdateExplorationMetaPolicy}(\mathcal{E}, c, \theta, o_\tau)$ \triangleright This update can be achieved with a fast incremental learning algorithm (e.g. memory-based learning)
 - 14: $\gamma(g | c) \leftarrow \text{UpdateGoalPolicy}(\mathcal{E}, c, g, o_\tau, r_i)$
 - 15: $\mathcal{E} \leftarrow \text{UpdateKnowledge}(\mathcal{E}, c, g, \theta, o_\tau, \tau, r_i)$
 - 16: **loop** \triangleright Target meta-policy training loop
 - 17: $\Pi(\theta | p, c) \leftarrow \text{UpdateMetaPolicy}(\mathcal{E})$ \triangleright This may be achieved using online or batch training of e.g. deep neural networks, SVMs, Gaussian Mixture Models, etc.
 - 18: **return** Π
-

Several algorithmic ingredients are used in Intrinsically Motivated Goal Exploration Processes.

With **Goal exploration**, the learning agent iteratively samples a parameterized problem in the space of problems and sets it as its own goal p in each interaction of the goal exploration loop. In this loop, the learner first infers the best current parameters θ with the current meta-policy $\Pi_\epsilon(\theta | p, c)$ (possibly adding some exploration noise). Several sampling strategies can be used, discussed below. Then, a roll-out of the lower-level policy π_θ is executed, allowing the observation of the behavioural trajectory τ and the measurement of the outcome o . The new data (c, o_τ, τ) then allows to: 1) compute the reward r associated to goal p ; 2) compute an intrinsic reward evaluating the interestingness of the choice of p based on the comparison of r with past experiences; 3) update the goal exploration policy (goal sampling strategy) with this intrinsic reward; 4) update the meta-policy Π_ϵ with a fast incremental learning algorithm, which depends on the particular implementation of the IMGEP (using the fact that the observation of the outcome o provides information about all other problems $p_i \in \mathcal{P}$); 5) update the learning database \mathcal{E} . Then, asynchronously this learning database \mathcal{E} can be used to learn a target meta-policy Π with an online or batch (potentially slower but better at generalization) learning algorithm, which also depends on the particular implementation of the IMGEP.

In goal exploration, a goal problem $g \in \mathcal{P}$ is chosen at each iteration. \mathcal{P} may be infinite, continuous and of high-dimensionality, which makes the choice of goal a non trivial question. Indeed, even if the goal-parameterized reward function R_{c, θ, o_τ} gives information about the fitness of policy π_θ to solve all problems $p \in \mathcal{P}$, this policy has been chosen with the problem g to solve in mind, thus R_{c, θ, o_τ} may not give as much information about other problems than the execution of another policy $\pi'_{\theta'}$ chosen when targeting another goal g' . **Intrinsic Rewards** provide a mean for the agent to self-estimate the expected interest of exploring goals $g \in \mathcal{P}$ for learning how to solve all problems $p \in \mathcal{P}$. An

intrinsic reward signal r_i is associated to a chosen goal g , and can be based on a heuristic (IR) such as outcome novelty, progress in reducing outcome prediction error, or progress in competence to solve problems [Oudeyer and Kaplan, 2007]. Based on the intrinsic rewards for different goals in different contexts, the choice of a goal g in a context c : goal policy $\gamma(g | c)$, is typically implemented with a contextual Multi-Armed Bandit to maximize future intrinsic rewards. Here, we will use intrinsic rewards based on measuring competence progress towards self-generated goals, which has been shown to be particularly efficient for learning repertoires of high-dimensional robotics skills [Baranes and Oudeyer, 2013]. See Fig. 1 for a schematic representation of possible learning curves and the exploration preference of an agent with intrinsic rewards based on learning progress.

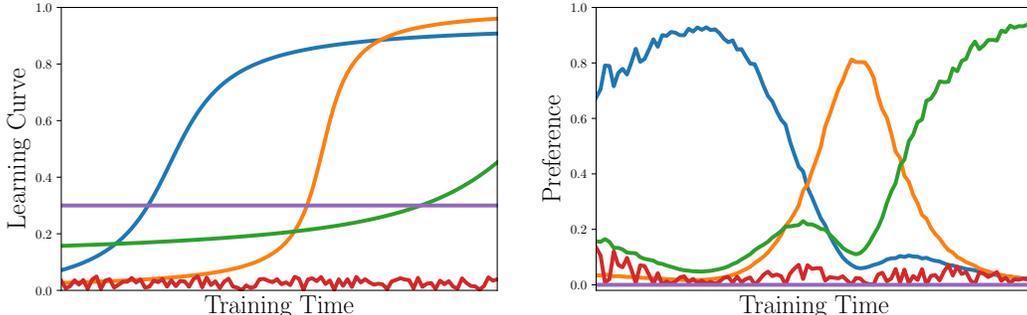


Figure 1: Schematic representation of possible learning curves of different problems and the associated exploration preference for an agent with intrinsic rewards based on learning progress. Left: We plot schematic learning curves associated to 5 imaginary problems: the y axis represent the competence of the agent to solve the problem (1 is perfect, 0 is chance level), and the x axis is training time on a problem. The blue, orange and green curves represent learnable problems, for which agent’s competence increases with training, at different rates, and saturates after a long training time. The purple curve represents a problem on which the agent always has the same competence, with no progress. The red curve is the learning curve on an unlearnable problem with stochastic outcomes, e.g. trying to predict the outcome of a random dice. Right: exploration preference of an agent with a learning progress heuristic to explore the 5 problems defined by the learning curves. The exploration preference is here based on progress which can be computed as the time derivative of the competence of the agent, and is normalized so that the sum of the progress on all problems is 1. At the beginning of exploration, the agent makes the most progress on problem blue so it prefers to train on this one, and then its preference will shift towards problem orange when it will make less progress on problem blue, and then progressively shift to problem green. The agent is making no progress on problem purple so will not choose to explore it, and problem red has a noisy but low estimated learning progress.

Exploration meta-policy and target meta-policy During the goal exploration loop, the main objective consists in doing experiments that allow to collect data that cover well the space of problems (i.e. to find a set of θ parameters that allows to find relatively good solutions to problems p over the problem space). The exploration meta-policy $\Pi_\epsilon(\theta | p, c)$ is learned and used to output a distribution of policies π_θ that are interesting to execute to gather information for solving the self-generated problem/goal p (and problems similar to p) in context c . To achieve the objective of collecting interesting data, the exploration meta-policy Π_ϵ must have fast and incremental updates. As here the aim is to maximize the coverage of the space of problems, being very precise when targeting goals is less crucial than the capacity to update the meta-policy quickly and incrementally. Thus, memory-based learning methods are well adapted within this loop (see Appendix A for examples of implementations). On the contrary, the target policy training loop aims to learn a meta-policy Π which purpose is to be used in exploitation mode: later on, this meta-policy can be asked to solve as precisely as possible some problems p with maximum reward. As the training of this meta-policy can be done asynchronously from data collected by the goal exploration loop, this allows to use training algorithms which are slower, possibly batch, but might allow to better generalize, e.g. using Gaussian mixture models, support vector regression or (deep) neural networks. These differences justify the fact that IMGEPs use in general two different representations and learning algorithms for learning meta-policies Π_ϵ and Π . This two-level learning scheme has similarities with the Complementary Learning Systems Theory used to account for the organization of learning in mammalian brains [Kumaran et al., 2016].

3.2 Particular case: Intrinsically Motivated Modular Goal Exploration Processes

In the particular case where the problem space is modular (Section 2.2), the agent can learn one goal policy γ_k per problem space \mathcal{P}^k (see Architecture 2). The choice of goals becomes hierarchical in the sense that the agent first chooses a problem space \mathcal{P}^k to explore with a goal space policy $\Gamma(\mathbf{k} | \mathbf{c})$ and then a particular goal $g \in \mathcal{P}^k$ with the corresponding goal policy γ_k . Those two levels of choice can make use of the self-computed intrinsic rewards r_i . The learning curves of Fig. 1 can also be seen as learning curves of the different problem spaces \mathcal{P}^k , and the exploration preference, $\Gamma(\mathbf{k} | \mathbf{c})$, can be based on the learning progress in each space \mathcal{P}^k , computed as the average progress across problems in \mathcal{P}^k .

Architecture 2 Intrinsically Motivated Modular Goal Exploration Process

Require: Environment Env , Context space \mathcal{C} with distribution $\mu_{Env}(\mathbf{c})$

Require: Policy param. space Θ , Outcome space $\mathcal{O} = \prod_{k=1}^m \mathcal{O}^k$, trajectory distribution $\nu_{Env}(\tau | \theta, \mathbf{c})$

Require: Problem parameter space $\mathcal{P} = \bigcup_{k=1}^m \mathcal{P}^k$, goal-parameterized reward function $R(p, c, \theta, o_\tau)$

Require: Initial knowledge \mathcal{E} ▷ typically empty

- 1: $\Pi(\theta | \mathbf{p}, \mathbf{c}) \leftarrow \text{InitializeMetaPolicy}(\mathcal{E})$ ▷ This is the target meta-policy (inverse model)
- 2: $\Pi_\epsilon(\theta | \mathbf{p}, \mathbf{c}) \leftarrow \text{InitializeExplorationMetaPolicy}(\mathcal{E})$ ▷ Π_ϵ is the inverse model used during exploration. It may be the same or different from Π
- 3: **for** k **do**
- 4: $\gamma_k(g | \mathbf{c}) \leftarrow \text{InitializeGoalPolicy}(\mathcal{E}, \mathcal{P}^k)$
- 5: $\Gamma(\mathbf{k} | \mathbf{c}) \leftarrow \text{InitializeGoalSpacePolicy}(\mathcal{E}, \{k\})$
- 6: Launch asynchronously the two following loops (exploration and training of target meta-policy)
- 7: **loop** ▷ Exploration loop
- 8: $c \sim \mu_{Env}(\mathbf{c})$ ▷ Observe context c
- 9: $k \sim \Gamma(\mathbf{k} | \mathbf{c})$ ▷ Choose goal space \mathcal{P}^k based on intrinsic rewards, e.g. using a MAB
- 10: $g \sim \gamma(g | \mathbf{c})$ ▷ Choose goal g in \mathcal{P}^k
- 11: $\theta \sim \Pi_\epsilon(\theta | g, \mathbf{c})$ ▷ Infer policy parameters θ that maximize expected $R_g(c, \theta, o_\tau)$ using an exploration/exploitation trade-off, and possibly computing now the problem-specific reward function R_g giving the reward/fitness r of all previous experiments (c_j, θ_j, o_j) to solve current goal problem g , with $r = R_g(c_j, \theta_j, o_j) = R(g, c_j, \theta_j, o_j)$
- 12: $\tau \sim \nu_{Env}(\tau | \theta, \mathbf{c})$ ▷ Execute a roll-out of π_θ , observe trajectory $\{s_{t_0:t_{end}}, a_{t_0:t_{end}}\}$
- 13: Compute outcome o_τ from trajectory τ
▷ The goal-parameterized reward function R_{c, θ, o_τ} can now be computed to find the reward/fitness r of the current experiment (c, θ, o_τ) for solving any problem $p \in \mathcal{P}$, with $r = R_{c, \theta, o_\tau}(p) = R(p, c, \theta, o_\tau)$
- 14: $r = R_{c, \theta, o_\tau}(g)$ ▷ Compute current reward associated to solving goal problem g
- 15: $r_i \leftarrow IR(\mathcal{E}, c, g, \theta, o_\tau, r)$ ▷ Compute intrinsic reward r_i associated to g in context c (e.g. learning progress in the vicinity of this goal)
- 16: $\Pi_\epsilon(\theta | \mathbf{p}, \mathbf{c}) \leftarrow \text{UpdateExplorationPolicy}(\mathcal{E}, c, \theta, o_\tau)$ ▷ This update can be achieved with a fast incremental learning algorithm (e.g. memory-based learning)
- 17: $\gamma_k(g | \mathbf{c}) \leftarrow \text{UpdateGoalPolicy}(\mathcal{E}, c, g, o_\tau, r_i)$
- 18: $\Gamma(\mathbf{k} | \mathbf{c}) \leftarrow \text{UpdateGoalSpacePolicy}(\mathcal{E}, c, k, g, o_\tau, r_i)$
- 19: $\mathcal{E} \leftarrow \text{UpdateKnowledge}(\mathcal{E}, c, g, \theta, o_\tau, \tau, r_i)$
- 20: **loop** ▷ Target meta-policy training loop
- 21: $\Pi(\theta | \mathbf{p}, \mathbf{c}) \leftarrow \text{UpdateMetaPolicy}(\mathcal{E})$ ▷ This may be achieved using online or batch training of e.g. deep neural networks, SVMs, Gaussian Mixture Models, etc.
- 22: **return** Π

4 Experiments: Exploration and Learning in a Robotic Tool Use Setup

In order to benchmark different learning algorithms in a complex realistic environment with continuous policy and outcome spaces, we designed a real robotic setup composed of a humanoid arm in front of joysticks that can be used as tools to act on other objects. We show the running experimental setup in this video². The code is available open-source together with the 3D shapes of printed objects³.

4.1 Robotic Setup

The robotic setup has two platforms: in the first one, a Poppy Torso robot (the learning agent) is mounted in front of two joysticks (see Fig. 2). In the second platform, a Poppy Ergo robot (seen as a robotic toy) is controlled by the right joystick and can push a ball that controls some lights and sounds. Poppy is a robust and accessible open-source 3D printed robotic platform [Lapeyre et al., 2014].

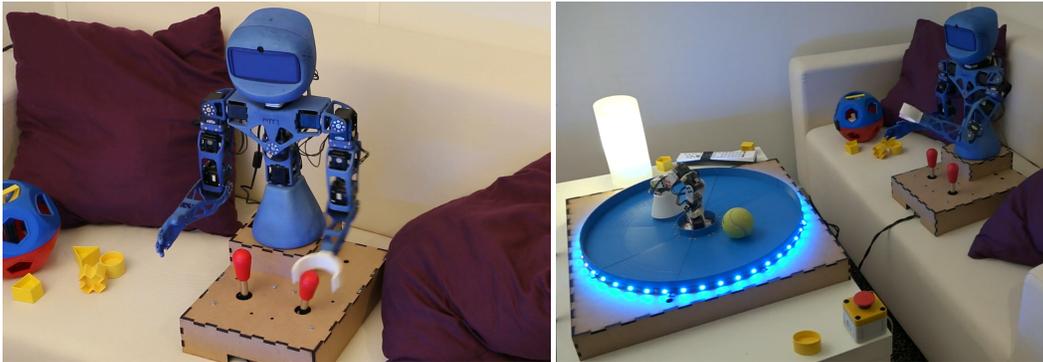


Figure 2: Robotic setup. Left: a Poppy Torso robot (the learning agent) is mounted in front of two joysticks. Right: full setup: a Poppy Ergo robot (seen as a robotic toy) is controlled by the right joystick and can hit a tennis ball in the arena which changes some lights and sounds.

Robotic Arm The left arm has 4 joints. The position of those joints at time t is defined by the action a_t . Their bounds are defined so that the arm has a low probability to self-collide but can still reach a large volume, even on the left, top and behind the left shoulder to some extent. We use the framework of Dynamical Movement Primitives [Ijspeert et al., 2013] to generate smooth joint trajectories given a set of motor parameters. Each of the 4 joints is controlled by a DMP starting at the rest position of the joint (position 0) and parameterized by 8 weights: one weight on each of 7 basis functions and one weight representing the end position of the joint trajectory (see Appendix B). Given θ (32 parameters between -1 and 1) provided by the agent, the DMPs generates a policy roll-out by outputting a smooth 30-steps trajectory $\{a_{t_0}, \dots, a_{t_{end}}\}$ for the joints of the arm that once executed will translate into a 3D trajectory of the robotic hand for 5s. After producing each roll-out, the arm goes back in a rest position.

Tools and Toys Two analogical joysticks (Ultrastick 360) can be reached by the left arm and moved in any direction. The 2D position of the joysticks (left-right and backward-forward axes) controls the Poppy Ergo robotic toy as follows. The left joystick does not control any variable. The Ergo robot has 6 motors, and moves with hardwired synergies that allow control of rotational speed and extension. The right joystick left-right axis controls in speed the rotation of the Ergo robot around the center of the second platform, which means that pushing the right joystick to the right with a small angle will move the Ergo towards the right with a small speed, and pushing the joystick with a higher angle will increase Ergo's rotational speed. The Ergo rotation angle is bounded in $[-\pi; \pi]$, and is reset to 0 every 40 iterations. The right joystick backward-forward axis controls the extension of the Ergo: if the joystick is in rest position, the Ergo stays in rest position. When the right joystick is moved forward, then the Ergo extends away from the center, using 3 of the 6 motors, and comes back when the joystick is released. A yellow tennis ball is freely moving in the blue arena which is slightly sloped so that the ball always comes close to the center at the end of a movement. The ball is

²Video of the experimental setup: <https://youtu.be/NOLAwD4ZTW0> Please note that in the current experiments we swapped the joysticks and changed the light and sound mappings.

³Open-source code: <https://github.com/ymollard/APEX>

tracked with a RGB camera to retrieve its 2D position in polar coordinates (rotation and extension). The speed of the ball controls (above a threshold) the intensity of the light of a LED circle around the arena. Finally, when the ball touches the border of the arena, a sound is produced and varied in pitch depending on the ball rotation angle. From the learning agent’s point of view, all "things" are objects, with no special status of objects as tools to manipulate other objects: this is discovered by the agent.

Distractors Several other objects are included in the environment, with which the agent cannot interact. The agent is not initially aware that those objects can’t be controlled. Two 2D objects are moving randomly, independently of the agent (imagine a cat and a dog playing together), with a noise on each variable added at each time step (between -0.05 and 0.05). Six objects are static: the right hand (3D) of the robot that is disabled in this experiment, the camera recording the ball trajectory (3D), the blue circular arena (2D), a yellow toy out-of-reach (2D), the red button also out-of-reach (2D) and the lamp (2D). All distractor objects are reset after each roll-out.

Trajectory and outcomes Before choosing a 32D motor command θ , the agent observes the current context c as the configuration of objects in the scene (in practice, since only the Ergo and ball are not reset after each roll-out, this amounts to measuring the rotation angle of the Ergo and of the ball around the center of the arena). Then, the chosen command θ translates through DMPs into a motor trajectory of the left arm which is executed for 5 seconds. We assume that there is perceptual system providing trajectories of all objects in the scene. The learning of such a system is not the focus here. The successive states of the environment are defined as follows. First, the 3D trajectory of the hand is computed through a forward model of the arm as its x , y and z position. The 2D states of each joystick and of the Ergo are read by sensors, and the position of the ball retrieved through the camera. The states of the 1D intensity of the light and the 1D pitch of the sound are computed from the ball position and speed. The state s_t represents the concatenation of the states of each of the 15 objects at time t . All variables of vector s_t are between -1 and 1 . The trajectory τ is defined as the sequence of actions and states of the environment during the 5 seconds: $\tau = \{s_{t_0}(=c), a_{t_0}, \dots, s_{t_{end}}, a_{t_{end}}\}$. We define the outcome o_k corresponding to object k as a 10-steps sampling of its trajectory during the movement. The outcome space corresponding to the hand is \mathcal{O}^1 (30D), the outcome space corresponding to the left joystick, right joystick, Ergo and Ball are respectively $\mathcal{O}^2, \mathcal{O}^3, \mathcal{O}^4, \mathcal{O}^5$ (20D each), and to the light and sounds are $\mathcal{O}^6, \mathcal{O}^7$ (10D each). The outcome space corresponding to the two random distractors are \mathcal{O}^8 and \mathcal{O}^9 (20D each), and to the six static objects: $\mathcal{O}^{10}, \mathcal{O}^{11}$ (30D each) and $\mathcal{O}^{12}, \mathcal{O}^{13}, \mathcal{O}^{14}, \mathcal{O}^{15}$ (20D each). At the end of each movement, the agent computes the outcome o_τ as the concatenation of the outcomes corresponding to each object: $o_\tau = (o_1, \dots, o_{15}) \in \mathcal{O}$ (310D). The goals of the agent will be defined as particular outcomes to reach in a space \mathcal{O}^k . Goals are thus particular trajectories of an object, so many goals are actually not solvable, e.g. moving the hand to the right and left very quickly, or moving the ball with a too high speed. Each outcome space \mathcal{O}^k defines a goal (or problem) space \mathcal{P}^k : for all k , $\mathcal{P}^k = \mathcal{O}^k$.

4.2 Problem and reward definition

We use the framework of unsupervised multi-goal reinforcement learning problem with a modular problem space (see Section 2.2). The initial experience \mathcal{E} is void. The context space $\mathcal{C} = [-1; 1]^2$ represents the position of each object before the movement. The policy parameter space is $\Theta = [-1; 1]^{32}$, and a parameter $\theta \in \Theta$ defines an action sequence through the Dynamical Movement Primitive framework. The execution of a policy π_θ in a context c leads to a trajectory $\tau = \{s_{t_0}(=c), a_{t_0}, \dots, s_{t_{end}}, a_{t_{end}}\}$ of robot actions and environmental states. The outcomes o_k are computed as samples of objects’ trajectories during the movement, and o_τ as the concatenation of outcomes.

The outcome space \mathcal{O}^k of an object defines particular goals that the agent will try to reach, or trajectories that the agent will try to give to the object. Each space \mathcal{O}^k is thus a problem space:

$$\mathcal{P}^k = \mathcal{O}^k \text{ for all } k. \text{ The problem space is then } \mathcal{P} = \bigcup_{k=1}^{15} \mathcal{P}^k, \text{ and the outcome space } \mathcal{O} = \prod_{k=1}^{15} \mathcal{O}^k =$$

$[-1; 1]^{310}$. We define the reward associated to solving problem $p \in \mathcal{P}^k$ in context $c \in \mathcal{C}$ with parameters $\theta \in \Theta$ and outcome $o_\tau \in \mathcal{O}$ as $R(p, c, \theta, o_\tau) = -\|p - o_k\|_k$. The norm $\|\cdot\|_k$ is the Euclidean norm divided by the maximal distance in \mathcal{O}^k , so that rewards are comparable across problem spaces. Given a context c , a policy π_θ and an outcome o_τ , the reward $R(p, c, \theta, o_\tau)$ can be computed by the agent for all $p \in \mathcal{P}$ and at any time after the experience (c, θ, o_τ) . We are interested in the learning of a solution to all problems of all problem spaces, so we define the test distribution $\mathcal{T}(p, c)$ to be uniform on $\mathcal{P} \times \mathcal{C}$, and the loss $\mathcal{L}(\Pi)$ of a meta-policy $\Pi(\theta | p, c)$ as in Equation 1.

4.3 Exploration Algorithms

We study four variants of modular goal exploration processes plus a control with random agents.

Random Motor Babbling (RANDOM) In this control condition, agents choose a random policy $\theta \in [-1; 1]^{32}$ at each iteration. It is a way to assess how long a random exploration would take to discover reward information for improving towards given problems. In the four following algorithms, agents choose random policies in 10% of the iterations (randomly chosen).

Random Model Babbling (RMB) Model Babbling (MB) is a particular model-based implementation of an Intrinsically Motivated Modular Goal Exploration Process (see Section 3.2, Architecture 2), where the agent learns one model per goal space [Forestier and Oudeyer, 2016b]. Each goal space corresponds here to the trajectory space of one object. Each model is providing an inverse function which, given the current context and desired goal (in the corresponding goal/problem space), gives the policy π_θ estimated to best reach this goal in this context. The goal policy is hierarchical in the sense that the agent first chooses a model to train (thus we call this learning procedure Model Babbling) or a goal space to explore, with Γ , and then a particular goal in this goal space, with γ_k . In Random Model Babbling, the choice of goal space (or model to train), and of particular goal are random: $\Gamma(\mathbf{k} \mid \mathbf{c})$, and $\gamma_k(\mathbf{g} \mid \mathbf{c})$ for each k are always uniform distributions. The meta-policies $\Pi(\theta \mid \mathbf{p}, \mathbf{c})$ and $\Pi_\epsilon(\theta \mid \mathbf{p}, \mathbf{c})$ are memory-based and initialized as uniform distributions. After bootstrapping with some experience from random policies, $\Pi(\theta \mid \mathbf{p}, \mathbf{c})$ is implemented here as a fast nearest neighbor search with a kd-tree. Given a goal problem g in context c , the agent chooses the previously executed policy for which the corresponding context-outcome (c, o_k) was the closest to the current (c, g) , i.e. that minimizes $(R_p(c', \theta, o_\tau)^2 + \|c - c'\|^2)$. The exploration meta-policy Π_ϵ , used to explore new policies, is defined as the meta-policy Π plus a Gaussian exploration noise: for all p and c , $\Pi_\epsilon(\theta \mid p, c) = \Pi(\theta \mid p, c) + \mathcal{N}(0, \Sigma)$ with Σ being a diagonal covariance matrix with weights $\sigma^2 = 0.05$ (see Algo. 3). The meta-policies Π and Π_ϵ have this implementation in all our algorithmic variants (RMB, SGS, FC and AMB). Other possible implementations of the meta-policies are described in Appendix A.

Algorithm 3 Implementation of sampling and update of Π and Π_ϵ

```

1: function SAMPLEMETAPOLICY( $p, c$ )
2:   Find in  $\mathcal{D}$  the tuple  $(c', \theta, o_\tau)$  that minimizes  $(R_p(c', \theta, o_\tau)^2 + \|c - c'\|^2)$ 
3:   return  $\theta$ 
4: function SAMPLEEXPLORATIONMETAPOLICY( $p, c$ )      ▷ Implements Architecture 1, line 8
5:    $\theta \sim \Pi(\theta \mid p, c)$                           ▷ Uses SampleMetaPolicy above
6:    $\epsilon \sim \mathcal{N}(0, \Sigma)$                              ▷ Add exploration noise
7:   return  $\theta + \epsilon$ 
8: function UPDATEMETAPOLICY( $\mathcal{E}$ )                      ▷ Implements Architecture 1, line 18
9:   Get last  $(c, \theta, o_\tau)$  from  $\mathcal{E}$ 
10:  Add  $(c, \theta, o_\tau)$  to a dataset  $\mathcal{D}$  of tuples (policy parameters, outcome)
11: function UPDATEEXPLORATIONMETAPOLICY( $\mathcal{E}$ )          ▷ Implements Architecture 1, line 14
12:  Same as UpdateMetaPolicy above

```

Single Goal Space (SGS) This algorithm is similar to the RMB algorithm, but the chosen goal space is always the outcome space of the ball: Γ always chooses \mathcal{O}^5 , and γ_5 chooses random goals in \mathcal{O}^5 . We define this condition to study how agents would learn in a case where an engineer only cares about the ability of the robot to push the ball and gives this only goal space to the robot.

Fixed Curriculum (FC) This algorithm is similar to the RMB one, but here Γ is a curriculum engineered by hand: the agents choose 7 goal spaces for one seventh of the total number of iterations in the sequence [*Hand, Left Joystick, Right Joystick, Ergo, Ball, Light, Sound*].

Active Model Babbling (AMB) Active Model Babbling is a variant of MB and an implementation of Architecture 2 (see Section 3.2) where the goal space policy Γ tries to maximize an empirical estimation of learning progress, such as in [Forestier and Oudeyer, 2016b]. The agent estimates its learning progress globally in each problem space (or for each model learned). At each iteration, the context c is observed, a goal space k is chosen by Γ and a random goal g is sampled by γ_k in \mathcal{O}^k . Then, in 80% of the iterations, the agent uses $\Pi_\epsilon(\theta \mid g, c)$ to generate with exploration a policy θ and does not update its progress estimation. In the other 20%, it uses Π , without exploration, to generate

θ and updates its learning progress estimation in \mathcal{O}^k , with the estimated progress in reaching g . To estimate the learning progress r_i made to reach the current goal g , the agent compares the outcome o with the outcome o' obtained for the previous context and goal (g', c') most similar (Euclidean distance) to (g, c) : $r_i = R(g, c, \theta, o) - R(g, c', \theta', o')$. Finally, Γ implements a non-stationary bandit algorithm to sample goal spaces. The bandit keeps track of a running average r_i^k of the intrinsic rewards r_i associated to the current goal space \mathcal{P}^k . With probability 20%, it samples a random space \mathcal{P}^k , and with probability 80%, it uses a soft maximization where the probability to sample \mathcal{P}^k is proportional to $\exp(\frac{r_i^k}{\sum_k (r_i^k)})$ if $r_i^k > 0$ and 0 otherwise. This bandit is a variant of the strategic bandit of [Lopes and Oudeyer, 2012].

5 Results

We ran 1 trial of 5000 iterations for the RANDOM condition and 3 independent trials of 5000 for each other condition. First, we show how agents explored the different outcome spaces \mathcal{O}^k depending on the condition, then we study the structure of the learning problem by looking in more details how agents in condition RMB succeeded to get information on problems of one space \mathcal{P}^k while exploring another one, and finally we analyze intrinsic rewards and the goal space policy Γ in condition AMB. Each trial of 5000 iteration takes about 10h with two Raspberry Pi 3 computers running the learning algorithm and control of Torso and Ergo robots. Our implementation of the learning algorithms is online and incremental and the inverse models use Nearest Neighbors search (see Algo. 3) so the whole learning is very energy-efficient: the Experimental Computational Energy Cost of our paper is about 1.7 kWh⁴.

Exploration of outcome spaces In order to estimate the performance of the different learning algorithms, one can compute the loss of Eq. 1 given a test distribution $\mathcal{T}(\mathbf{p}, c)$, which measures the accuracy of the learned meta-policy $\Pi(\theta | \mathbf{p}, c)$ to choose good parameters θ to solve problems p in contexts c sampled with $\mathcal{T}(\mathbf{p}, c)$. During training, the execution of 5000 policies π_θ allowed to produce a set of 5000 outcomes o_τ . Each outcome o_τ corresponds to a solution to a particular problem p because $\mathcal{P} = \bigcup_{k=1}^{15} \mathcal{O}^k$. The more diverse the outcomes o_τ are, the more diverse will be the set of problems that the meta-policy Π will solve with a good reward. For example, if the agent managed to move the left joystick in different directions, then the outcomes corresponding to the left joystick (o_2) are diverse, and the meta-policy should be able to reproduce those movements when evaluated on the problem of moving the left joystick to a given direction. A good proxy to measuring the loss is thus to measure the diversity of outcomes that were produced during exploration, for each outcome space \mathcal{O}^k .

We define the measure of exploration of each outcome space \mathcal{O}^k as the percentage of cells explored in a discretization of the outcome space variables. For instance, the outcome space corresponding to the hand is \mathcal{O}^1 (30D), corresponding to hand's x, y and z position at 10 time points during the 5 seconds movement. The 3D space of the x, y and z variables is discretized by 20 bins on each variable (so 8000 cells), and we count the number of 3D cells occupied by any of the 10 time points of any of the previous iterations. Similarly, the exploration of a 20D outcome space corresponding to an object with 2 variables is measured with a 100 bins discretization of the 2D space (e.g. for the joysticks outcome spaces).

We show the results of this exploration measure during learning in Fig. 3. The RANDOM condition generates a good diversity for the hand outcome space, but is very bad to explore other outcome spaces. The results of the FC condition depend on trials, in some trials the agent succeeded to discover how to move the object of the next exploration stage (e.g. right joystick) during the allowed exploration budget of the current stage (respectively left joystick) giving good exploration results, and in other trials it did not, showing a bad exploration of all other spaces. Condition SGS gives equal or lower exploration results than RANDOM in all spaces because the agent always focuses on the ball outcome space even when it does not know how to reach the joysticks. The exploration

⁴ The Experimental Computational Energy Cost of our paper is the total energy consumption of the processing units in all the experiments of our paper: about 1.7 kWh. Similarly, the Experimental Robotic Energy Cost of our paper is the total energy consumption of the robots' actuators in all the experiments of our paper: about 1.8 kWh. We can then estimate the carbon cost of running our experiments as about 140 gCO₂-eq.

results in conditions RMB and AMB are similar, and much better than other conditions, for example light and sound were only discovered in those conditions.

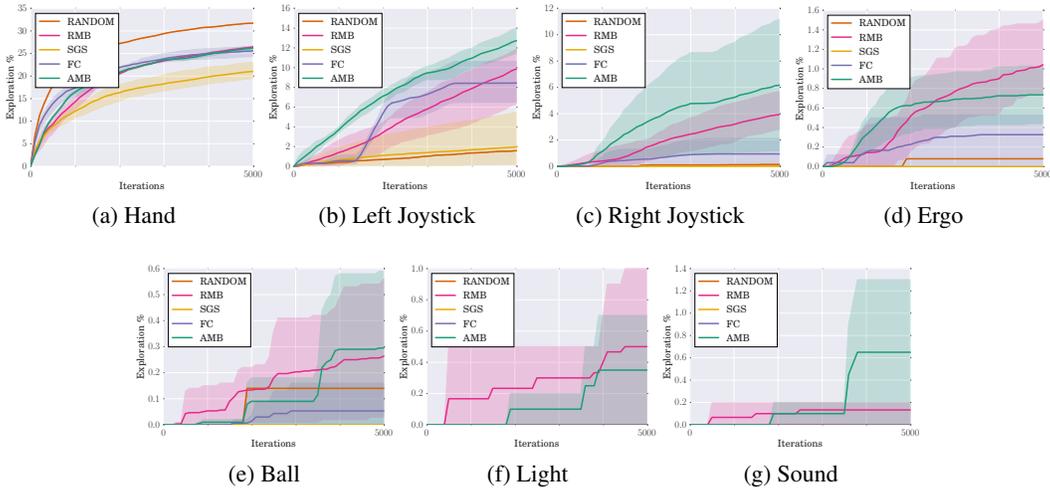


Figure 3: Exploration of some outcome spaces O_k . The exploration measure is the percentage of cells explored in a discretization of each outcome space. We provide the mean, min, and max of the trials of each condition during the 5000 learning iterations (about 8 hours).

Discoveries in RMB In order to understand the structure of the learning problem, we can look in more details how agents in condition RMB succeeded to get information on problems of one space \mathcal{P}^k while exploring another space. Fig. 4 shows the proportion of the iterations with a goal in a given space that allowed to move (a) the left joystick, (b) the right joystick and (c) the Ergo robot. First, the easiest reachable object while exploring random policies or choosing goals in the hand outcome space is the left joystick: reached in about 10% of iterations versus almost 0% for other objects. Also, to discover the right joystick, it is more efficient to choose goals for the left joystick (about 5% success) than for the hand or choosing random policies (0%). Similarly, (c) shows that to discover how to move Ergo, it is much more efficient to choose goals for the right joystick (about 15%) than for the hand or left joystick ($< 5\%$). Those results tells that a good strategy to discover all objects is to explore the different outcome spaces in a sequence from the easiest (Hand, Left Joystick) until some movements of the right joystick are discovered, then explore this one until Ergo is discovered etc. This structure is thus a transfer learning problem in the sense that the different problems are not independent, and exploring to find solutions to some of them allows to discover part of solutions to others.

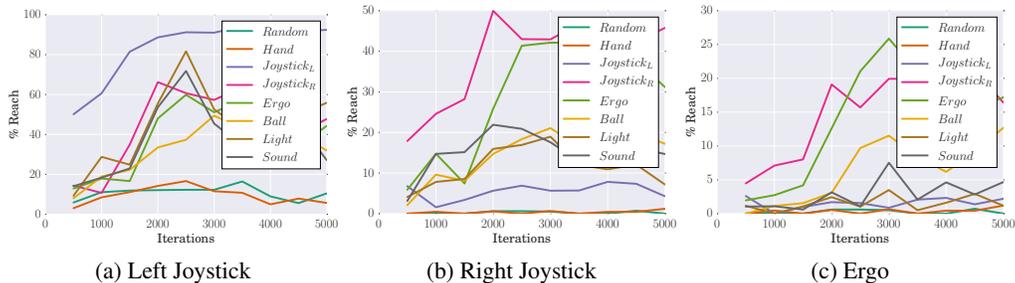


Figure 4: Transfer of solutions between problem spaces. We show the proportion of iterations that allowed to (a) reach the left joystick, (b) reach the right joystick, and (c) move the Ergo robot, depending on the current chosen goal space in the condition RMB (or random movements: *Random*).

Intrinsic Rewards in Active Model Babbling Fig. 5 shows the running average of the intrinsic rewards computed by the bandit algorithm of each of the 3 agents of condition AMB depending on the problem space in which the goals were chosen. We can see that while no solution has been

discovered to control an object, the intrinsic reward stays at 0 for this object, which means that the agent will rarely choose to train to solve problems in those spaces. The estimated learning progress in each space drives the choice of space to explore, and thus leads to a developmental sequence exploring from easier to more complex problems, avoiding to loose time on too complex problems.

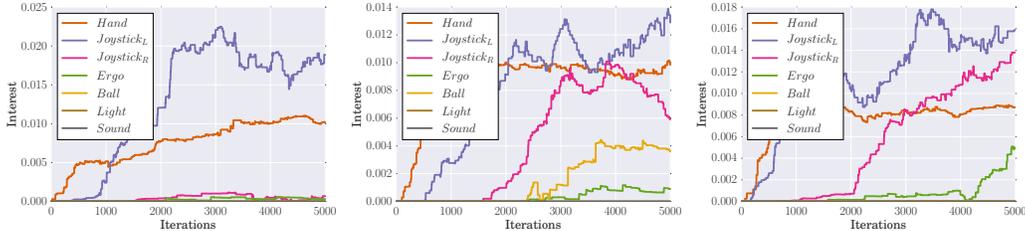


Figure 5: Intrinsic Rewards in condition AMB. Each graph shows the learning progress estimated by one agent of condition AMB in each of the problem spaces during the 5000 iterations. The learning progress heuristic leads to a sequence of exploration focus from easier to more complex problems.

6 Related work

Early models of intrinsically motivated reinforcement learning (also called curiosity-driven learning) have been used to drive efficient exploration in the context of target tasks with rare or deceptive rewards [Schmidhuber, 1991, Barto, 2013] or in the context of computational modelling of open-ended unsupervised autonomous learning in humans [Kaplan and Oudeyer, 2004, Oudeyer et al., 2007]. Reviews of the historical development of these methods and their links with cognitive sciences and neuroscience can be found in [Baldassarre and Mirolli, 2013, Gottlieb et al., 2013, Oudeyer et al., 2016].

Several lines of results have shown that intrinsically motivated exploration and learning mechanisms are particularly useful in the context of learning to solve reinforcement learning problems with sparse or deceptive rewards. For example, several state-of-the-art performances of Deep Reinforcement Learning algorithms, such as letting a machine learn how to solve complex video games, have been achieved by complementing the extrinsic rewards (number of points won) with an intrinsic reward pushing the learner to explore for improving its predictions of the world dynamics [Bellemare et al., 2016, Houthoof et al., 2016]. An even more radical approach for solving problems with rare or deceptive extrinsic rewards has been to completely ignore extrinsic rewards, and let the machine explore the environment for the sole purpose of learning to predict the consequences of its actions [Schmidhuber, 1991, Oudeyer et al., 2007] or of learning to control self-generated goals [Baranes and Oudeyer, 2013, Oudeyer and Kaplan, 2007], or to generate novel outcomes [Lehman and Stanley, 2011]. This was shown for example to allow robots to learn tool use [Forestier and Oudeyer, 2016b] or to learn how to play some video games [Pathak et al., 2017] without ever observing the extrinsic reward.

Some approaches to intrinsically motivated exploration have used intrinsic rewards to value visited actions and states through measuring their novelty or the improvement of predictions that they provide, e.g. [Sutton, 1990, Dayan and Sejnowski, 1996, Schmidhuber, 1991, Oudeyer et al., 2007] or more recently [Bellemare et al., 2016, Houthoof et al., 2016, Pathak et al., 2017]. However, organizing intrinsically motivated exploration at the higher level of goals (conceptualized as parameterized RL problems), by sampling goals according to measures such as competence progress [Oudeyer and Kaplan, 2007], has been proposed and shown to be more efficient in contexts with high-dimensional continuous action spaces and strong time constraints for interaction with the environment [Baranes and Oudeyer, 2013].

Several strands of research in robotics have presented algorithms that instantiate such intrinsically motivated goal exploration processes [Baranes and Oudeyer, 2010, Rolf et al., 2010], using different terminologies such as contextual policy search [Kupcsik et al., 2014, Queißer et al., 2016], or formulated within an evolutionary computation perspective [Cully et al., 2015]. However, these previous approaches were not formalized in the general framework of unsupervised multi-goal reinforcement learning, and they have not considered intrinsically motivated exploration of multiple spaces of goals and how this can allow the formation of a learning curriculum. Preliminary investigation of modular goal exploration processes was presented in [Forestier and Oudeyer, 2016b,

Forestier and Oudeyer, 2017], however this prior work only presented particular implementations of IMGEPs without formalizing them in the general formal framework presented here.

[Gregor et al., 2016], [Dosovitskiy and Koltun, 2016] and [Kulkarni et al., 2016] proposed methods that can be framed as IMGEPs, however they have considered notions of goals restricted to the reaching of states or direct sensory measurements, did not consider goal-parameterized rewards that can be computed for any goal, used different intrinsic rewards, and did not evaluate these algorithms in robotic setups. The notion of auxiliary tasks is also related to IMGEPs in the sense that it allows a learner to acquire tasks with rare rewards by adding several other objectives which increase the density of information obtained from the environment [Jaderberg et al., 2016]. Another line of related work [Srivastava et al., 2013] proposed a theoretical framework for automatic generation of problem sequences for machine learners, however it has focused on theoretical considerations and experiments on abstract problems.

In machine learning, the concept of curriculum learning [Bengio et al., 2009] has most often been used in the context of training neural networks to solve prediction problems. Many approaches have used hand-designed learning curriculum [Sutskever and Zaremba, 2014], but recently it was shown how learning progress could be used to automate intrinsically motivated curriculum learning in LSTMs [Graves et al., 2017]. However, these approaches have not considered curriculum learning of sets of reinforcement learning problems which is central in the IMGEP framework, and assumed the pre-existence of a database with learning exemplars to sample from. In recent related work, [Matiisen et al., 2017] studied how intrinsic rewards based on learning progress could also be used to automatically generate a learning curriculum with discrete sets of reinforcement learning problems, but did not consider high-dimensional modular parameterized RL problems. The concept of "curriculum learning" has also been called "developmental trajectories" in prior work on computational modelling of intrinsically motivated exploration [Oudeyer et al., 2007], and in particular on the topic of intrinsically motivated goal exploration [Baranes and Oudeyer, 2013, Forestier and Oudeyer, 2017].

7 Discussion

This paper provides a formal framework, the Unsupervised Multi-Goal Reinforcement Learning Problem, in which we expressed a problem structure that is appearing in different settings. We formalized a corresponding algorithmic architecture (IMGEP) that leverages this structure for efficient exploration and learning. We designed the first real robotic experiment where an intrinsically-motivated humanoid robot discovers a complex continuous high-dimensional environment and succeeds to explore and learn from scratch that some objects can be used as tools to act on other objects.

We evaluated different variants of Intrinsically Motivated Goal Exploration Processes and showed that only the variants where we do not hand-design a curriculum for learning (RMB and AMB conditions) allowed to discover the more complex skills. Furthermore, when the agent monitors its learning progress with intrinsic rewards (AMB), it autonomously develops a learning sequence, or curriculum, from easier to the most complex tasks. Also, the comparison between agents only exploring the ball problem space (SGS) versus all spaces (AMB) shows that if an engineer was to specify the target problems he wants the robot to solve (e.g. move the ball), then it would be more efficient to also explore all other possible intrinsic goals to develop new skills that can serve as stepping stones to solve the target problems.

We have chosen to evaluate several implementations of IMGEPs with a real robotic setup as this provides real world constraints that are not available in simulations, such as stringent time constraints that make it impossible to tune the parameters of algorithms by running many pre-experiments, as well as realistic complex noise. However, this also prevented us so far to run as much experiments as we would need to disentangle learning performances of the two conditions RMB and AMB. We are currently running more trials of the different conditions. Also, a current limitation of our setup is that we suppose that agents already have a perceptual system allowing them to see and track objects, as well as spaces of representations to encode their transformations. Future work will study how representation learning methods could bootstrap these representations of the scene from pixels in an unsupervised manner.

Acknowledgements

We would like to thank Olivier Sigaud and Alexandre Péré for their valuable comments on an earlier version of this manuscript and Damien Caselli for his technical help.

References

- [Baldassarre and Mirolli, 2013] Baldassarre, G. and Mirolli, M. (2013). *Intrinsically Motivated Learning in Natural and Artificial Systems*. Springer.
- [Baranes and Oudeyer, 2010] Baranes, A. and Oudeyer, P.-Y. (2010). Intrinsically motivated goal exploration for active motor learning in robots: A case study. In *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*.
- [Baranes and Oudeyer, 2013] Baranes, A. and Oudeyer, P.-Y. (2013). Active learning of inverse models with intrinsically motivated goal exploration in robots. *Robotics and Autonomous Systems*, 61(1).
- [Barto, 2013] Barto, A. G. (2013). Intrinsic motivation and reinforcement learning. In *Intrinsically motivated learning in natural and artificial systems*, pages 17–47. Springer.
- [Bellemare et al., 2016] Bellemare, M., Srinivasan, S., Ostrovski, G., Schaul, T., Saxton, D., and Munos, R. (2016). Unifying count-based exploration and intrinsic motivation. In *Advances in Neural Information Processing Systems*, pages 1471–1479.
- [Bengio et al., 2009] Bengio, Y., Louradour, J., Collobert, R., and Weston, J. (2009). Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48. ACM.
- [Berlyne, 1966] Berlyne, D. E. (1966). Curiosity and exploration. *Science*, 153(3731):25–33.
- [Cleveland and Devlin, 1988] Cleveland, W. S. and Devlin, S. J. (1988). Locally weighted regression: an approach to regression analysis by local fitting. *Journal of the American Statistical Association*, 83(403).
- [Cully et al., 2015] Cully, A., Clune, J., Tarapore, D., and Mouret, J.-B. (2015). Robots that can adapt like animals. *Nature*, 521(7553):503–507.
- [Dayan and Sejnowski, 1996] Dayan, P. and Sejnowski, T. J. (1996). Exploration bonuses and dual control. *Machine Learning*, 25(1):5–22.
- [Dosovitskiy and Koltun, 2016] Dosovitskiy, A. and Koltun, V. (2016). Learning to act by predicting the future. *arXiv preprint arXiv:1611.01779*.
- [Forestier and Oudeyer, 2016a] Forestier, S. and Oudeyer, P.-Y. (2016a). Curiosity-driven development of tool use precursors: a computational model. In *Proceedings of the 38th Annual Meeting of the Cognitive Science Society*.
- [Forestier and Oudeyer, 2016b] Forestier, S. and Oudeyer, P.-Y. (2016b). Modular active curiosity-driven discovery of tool use. In *Intelligent Robots and Systems (IROS), 2016 IEEE/RSJ International Conference on*, pages 3965–3972. IEEE.
- [Forestier and Oudeyer, 2016c] Forestier, S. and Oudeyer, P.-Y. (2016c). Overlapping waves in tool use development: a curiosity-driven computational model. In *Sixth Joint IEEE International Conference on Development and Learning and Epigenetic Robotics (ICDL-EpiRob)*.
- [Forestier and Oudeyer, 2017] Forestier, S. and Oudeyer, P.-Y. (2017). A unified model of speech and tool use early development. In *Proceedings of the 39th Annual Meeting of the Cognitive Science Society*.
- [Gopnik et al., 1999] Gopnik, A., Meltzoff, A. N., and Kuhl, P. K. (1999). *The scientist in the crib: Minds, brains, and how children learn*. William Morrow & Co.
- [Gottlieb et al., 2013] Gottlieb, J., Oudeyer, P.-Y., Lopes, M., and Baranes, A. (2013). Information-seeking, curiosity, and attention: computational and neural mechanisms. *Trends in Cognitive Sciences*, 17(11).
- [Graves et al., 2017] Graves, A., Bellemare, M. G., Menick, J., Munos, R., and Kavukcuoglu, K. (2017). Automated curriculum learning for neural networks. *arXiv preprint arXiv:1704.03003*.

- [Gregor et al., 2016] Gregor, K., Rezende, D. J., and Wierstra, D. (2016). Variational intrinsic control. *arXiv preprint arXiv:1611.07507*.
- [Hansen, 2006] Hansen, N. (2006). The CMA evolution strategy: a comparing review. In *Towards a new evolutionary computation*. Springer.
- [Houthoofd et al., 2016] Houthoofd, R., Chen, X., Duan, Y., Schulman, J., De Turck, F., and Abbeel, P. (2016). Vime: Variational information maximizing exploration. In *Advances in Neural Information Processing Systems*, pages 1109–1117.
- [Ijspeert et al., 2013] Ijspeert, A. J., Nakanishi, J., Hoffmann, H., Pastor, P., and Schaal, S. (2013). Dynamical movement primitives: learning attractor models for motor behaviors. *Neural computation*, 25(2).
- [Jaderberg et al., 2016] Jaderberg, M., Mnih, V., Czarnecki, W. M., Schaul, T., Leibo, J. Z., Silver, D., and Kavukcuoglu, K. (2016). Reinforcement learning with unsupervised auxiliary tasks. *arXiv preprint arXiv:1611.05397*.
- [Kaplan and Oudeyer, 2004] Kaplan, F. and Oudeyer, P.-Y. (2004). Maximizing learning progress: an internal reward system for development. In *Embodied artificial intelligence*, pages 259–270. Springer.
- [Kulkarni et al., 2016] Kulkarni, T. D., Narasimhan, K., Saeedi, A., and Tenenbaum, J. (2016). Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. In *Advances in Neural Information Processing Systems*, pages 3675–3683.
- [Kumaran et al., 2016] Kumaran, D., Hassabis, D., and McClelland, J. L. (2016). What learning systems do intelligent agents need? complementary learning systems theory updated. *Trends in cognitive sciences*, 20(7):512–534.
- [Kupcsik et al., 2014] Kupcsik, A., Deisenroth, M. P., Peters, J., Loh, A. P., Vadakkepat, P., and Neumann, G. (2014). Model-based contextual policy search for data-efficient generalization of robot skills. *Artificial Intelligence*.
- [Lapeyre et al., 2014] Lapeyre, M., Rouanet, P., Grizou, J., Nguyen, S., Depraetre, F., Le Falher, A., and Oudeyer, P.-Y. (2014). Poppy Project: Open-Source Fabrication of 3D Printed Humanoid Robot for Science, Education and Art. In *Digital Intelligence 2014*, Nantes, France.
- [Lehman and Stanley, 2011] Lehman, J. and Stanley, K. O. (2011). Abandoning objectives: Evolution through the search for novelty alone. *Evolutionary computation*, 19(2):189–223.
- [Lopes and Oudeyer, 2012] Lopes, M. and Oudeyer, P.-Y. (2012). The strategic student approach for life-long exploration and learning. In *2012 IEEE International Conference on Development and Learning and Epigenetic Robotics (ICDL)*.
- [Matiisen et al., 2017] Matiisen, T., Oliver, A., Cohen, T., and Schulman, J. (2017). Teacher-student curriculum learning. *arXiv preprint arXiv:1707.00183*.
- [Moulin-Frier et al., 2014] Moulin-Frier, C., Nguyen, S. M., and Oudeyer, P.-Y. (2014). Self-organization of early vocal development in infants and machines: the role of intrinsic motivation. *Frontiers in Psychology*, 4.
- [Oudeyer et al., 2016] Oudeyer, P.-Y., Gottlieb, J., and Lopes, M. (2016). Intrinsic motivation, curiosity, and learning: Theory and applications in educational technologies. *Progress in brain research*, 229:257–284.
- [Oudeyer and Kaplan, 2007] Oudeyer, P.-Y. and Kaplan, F. (2007). What is intrinsic motivation? A typology of computational approaches. *Frontiers in Neurobotics*, 1.
- [Oudeyer et al., 2007] Oudeyer, P.-Y., Kaplan, F., and Hafner, V. V. (2007). Intrinsic Motivation Systems for Autonomous Mental Development. *IEEE Transactions on Evolutionary Computation*, 11(2).
- [Pathak et al., 2017] Pathak, D., Agrawal, P., Efros, A. A., and Darrell, T. (2017). Curiosity-driven exploration by self-supervised prediction. *arXiv preprint arXiv:1705.05363*.
- [Queißer et al., 2016] Queißer, J. F., Reinhart, R. F., and Steil, J. J. (2016). Incremental bootstrapping of parameterized motor skills. In *2016 IEEE-RAS 16th International Conference on Humanoid Robots (Humanoids)*, pages 223–229.
- [Rolf et al., 2010] Rolf, M., Steil, J., and Gienger, M. (2010). Goal babbling permits direct learning of inverse kinematics. *IEEE Transactions on Autonomous Mental Development*, 2(3).

- [Schmidhuber, 1991] Schmidhuber, J. (1991). A possibility for implementing curiosity and boredom in model-building neural controllers. In *From animals to animats: Proceedings of the first international conference on simulation of adaptive behavior*, pages 15–21.
- [Srivastava et al., 2013] Srivastava, R. K., Steunebrink, B. R., and Schmidhuber, J. (2013). First experiments with powerplay. *Neural Networks*, 41:130–136.
- [Sutskever and Zaremba, 2014] Sutskever, I. and Zaremba, W. (2014). Learning to execute. *arXiv preprint arXiv:1410.4615*.
- [Sutton, 1990] Sutton, R. S. (1990). Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In *Proceedings of the seventh international conference on machine learning*, pages 216–224.
- [Thelen and Smith, 1996] Thelen, E. and Smith, L. B. (1996). *A dynamic systems approach to the development of cognition and action*. MIT press.

A Appendix: Pseudo-code of particular implementations

In the Intrinsically Motivated Goal Exploration Processes (Architecture 1), the meta-policies Π and Π_ϵ need to be implemented. We provided the pseudo-code of our implementation in Algo. 3, which uses a fast nearest neighbor lookup to find good motor parameters θ to reach the current goal (that maximizes the reward for the current goal together with the context similarity). The exploration meta-policy Π_ϵ is implemented by adding a Gaussian exploration noise on the output of Π . We provide here other possible implementations of Π and Π_ϵ .

First, one can build a fast local model to predict the expected reward associated to policy parameters θ for solving problem p in context c . With this model, and given a goal problem p to solve in context c , an optimization method can now be used to infer a good θ , starting for example from the best θ_0 from a nearest neighbor search. The optimization can be performed by quasi-Newton methods (e.g. L-BFGS) or black-box methods (e.g. CMA-ES[Hansen, 2006]) (see Algo. 4).

Algorithm 4 Implementation of Π and Π_ϵ with optimization using a reward model

```

1: function SAMPLEMETAPOLICY( $p, c$ )
2:   function  $f_{objective}(\theta)$     ▷ Define an objective function for solving problem  $p$  in context  $c$ 
3:     Build a local model  $\mathcal{M}$  of rewards in the vicinity of  $(c, \theta)$ , e.g. with Locally Weighted
       Linear Regression [Cleveland and Devlin, 1988]
4:     Infer reward  $r$  of  $(c, \theta)$  with  $\mathcal{M}$ 
5:     return  $r$ 
6:   Find in  $\mathcal{D}$  the tuple  $(c', \theta_0, o_\tau)$  that minimizes  $(R_p(c', \theta_0, o_\tau)^2 + \|c - c'\|^2)$ 
7:    $\theta^* \leftarrow$  Maximize  $f_{objective}$  with respect to  $\theta$ , starting from  $\theta_0$ , until a timeout.
       ▷ Optimization can be performed by quasi-Newton methods such as L-BFGS or black-box
       methods such as CMA-ES[Hansen, 2006].
8:   return  $\theta^*$ 
9: function SAMPLEEXPLORATIONMETAPOLICY( $p, c$ )    ▷ Implements Architecture 1, line 8
10:   $\theta \sim \Pi(\theta \mid p, c)$                         ▷ Uses SampleMetaPolicy above
11:   $\epsilon \sim \mathcal{N}(0, \Sigma)$ 
12:  return  $\theta + \epsilon$ 
13: function UPDATEMETAPOLICY( $\mathcal{E}$ )                    ▷ Implements Architecture 1, line 18
14:  Get last  $(c, \theta, o_\tau)$  from  $\mathcal{E}$ 
15:  Add  $(c, \theta, o_\tau)$  to a dataset  $\mathcal{D}$  of tuples (policy parameters, outcome)
16: function UPDATEEXPLORATIONMETAPOLICY( $\mathcal{E}$ )      ▷ Implements Architecture 1, line 14
17:  Same as UpdateMetaPolicy above

```

Another idea is to perform optimization by directly executing the parameters θ queried by the optimization method and to observe the corresponding reward from the environment, instead of building a reward model (see Algo. 5). However, executing policy parameters in the environment usually takes much more time than the prediction by a model.

Algorithm 5 IMGEP with direct θ optimization

```
1: loop ▷ Exploration loop: replaces loop line 5 of Architecture 1
2:    $c \sim \mu_{Env}(c)$  ▷ Observe context  $c$ 
3:    $g \sim \gamma(g | c)$  ▷ Choose goal  $g$  in  $\mathcal{P}$  based on intrinsic rewards, e.g. using a MAB
4:   function  $f_{objective}(\theta)$  ▷ Define an objective function for the goal  $g$ 
5:      $c \sim \mu_{Env}(c)$ 
6:      $\tau \sim \nu_{Env}(\tau | \theta, c)$  ▷ Execute a roll-out of  $\pi_\theta$ , observe trajectory  $\{s_{t_0:t_{end}}, a_{t_0:t_{end}}\}$ 
7:     Compute outcome  $o_\tau$  from trajectory  $\tau$ 
8:      $r = R(g, c, \theta, o_\tau)$ 
9:      $r_i \leftarrow IR(\mathcal{E}, c, g, \theta, o_\tau, r)$ 
10:     $\Pi(\theta | p, c) \leftarrow \text{UpdateMetaPolicy}(\mathcal{E}, c, \theta, o_\tau)$ 
11:     $\gamma(g | c) \leftarrow \text{UpdateGoalPolicy}(\mathcal{E}, c, g, o_\tau, r_i)$ 
12:     $\mathcal{E} \leftarrow \text{UpdateKnowledge}(\mathcal{E}, c, g, \theta, o_\tau, \tau, r_i)$ 
13:    return  $r$ 
14:   $\theta_0 \sim \Pi(\theta | g, c)$  ▷  $\Pi$  can have any implementation, such as Algo. 3 or 4
15:  Maximize  $f_{objective}$  with respect to  $\theta$ , starting from  $\theta_0$ , until a timeout or until  $\gamma(g | c)$  is below a threshold, either because the goal is reached or no progress is made towards the goal, or the context changed too much and the goal  $g$  is not interesting in the new context
▷ Optimization can be performed via Bayesian Optimization or black-box methods such as CMA-ES[Hansen, 2006]
```

B Appendix: Details on Dynamical Movement Primitives

We use the framework of Dynamical Movement Primitives [Ijspeert et al., 2013] to generate smooth trajectories for the 4 joints of the left arm from a small set of 32 parameters θ . In this framework, a mass-spring-damper dynamical system (Eq. 2) is used to represent one motor angle (variable y here). The motor is going from the fixed position y_0 (at the middle between its bounds) to a goal position g that is defined by one of the parameters of θ . Its movement, with a mass-spring-damper dynamic, is perturbed by the addition of a perturbation function f which is a linear combination of 7 basis functions ψ_i with weights w_i defined by θ (see Eq. 3 and Fig. 6). A canonical system (Eq. 4) is used to reduce the impact of the perturbation in time with $x_0 = 1$, and $\alpha_y, \beta_y, \alpha_x$ are time constants.

$$\ddot{y} = \alpha_y(\beta_y(g - y) - \dot{y}) + f(x) \quad (2)$$

$$f(x) = \frac{\sum \psi_i w_i}{\sum \psi_i} x(g - y_0) \quad (3)$$

$$\dot{x} = -\alpha_x x \quad (4)$$

Each of the 4 motors is controlled by one dynamical system parameterized by 8 weights: one weight on each of 7 basis functions plus one weight representing the end position of the motor trajectory (Table 1). Given θ (32 parameters between -1 and 1) provided by the agent, the framework rolls out the dynamical systems and outputs a smooth 30-steps trajectory $\{a_{t_0}, \dots, a_{t_{end}}\}$ for the joints of the arm. Fig. 7 shows one 5-seconds movement of the arm that reaches the right joystick and makes the Ergo robot move.

θ	Basis 1	Basis 2	Basis 3	Basis 4	Basis 5	Basis 6	Basis 7	End point
Shoulder Y	1.0	0.760	-0.03	0.641	0.274	0.187	-0.77	0.164
Shoulder X	0.164	-0.99	0.029	-1.0	-0.39	-0.40	-0.75	0.927
Arm Z	-0.69	0.927	-0.47	-0.77	0.084	-0.05	0.221	-0.88
Elbow Y	-0.72	0.775	-0.88	0.532	-0.98	1.0	-0.70	0.886

Table 1: Example of a set of parameters θ for which the roll-out of policy π_θ allowed to move the right joystick as in Fig. 7. θ is a vector of 32 parameters: 8 parameters for each of the 4 motors. The first 7 parameters define the weights w_i on the basis functions ψ_i used to compute the perturbation function f (see Eq. 3) which modifies the trajectory of the one motor during its movement from the fixed starting position to the end position defined by the eighth parameter.

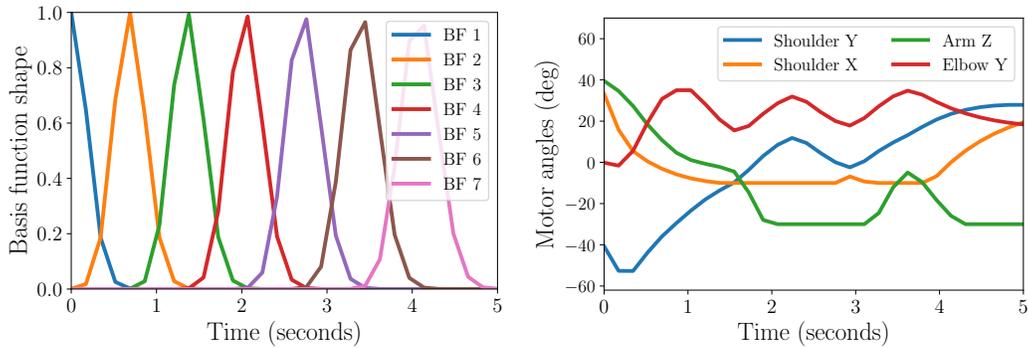


Figure 6: Left: The fixed set of 7 basis functions. The actual perturbation f is a linear combination of the basis functions and the weights w_i . For example, a positive weight on the second basis function (orange) will smoothly increase the motor's angle around time 0.7s and will have less effect later. Right: Roll-out of policy π_θ defined by the parameters θ of Table 1. This arm trajectory (found after 3000 iterations of exploration) allows to move the right joystick as in Fig. 7



Figure 7: Example of a policy π_θ that makes the arm reach the right joystick and move the Ergo robot.