

---

## Research proposal

### Introduction

(WILL) Transfer learning for reinforcement learning. Words... (Dumoulin et al. 2018)

### What is reinforcement learning (RL)?

An easy way to understand RL is as a Markov decision process (MDP).

A MDP is defined as a tuple,  $\{\mathcal{S}, \mathcal{A}, P(s_{t+1} \mid s_t, a_t), R(s_t, a_t, s_{t+1})\}$ . Where  $\mathcal{S}$  is the set of possible states (*for example arrangements of chess pieces*),  $\mathcal{A}$  is the set of actions (*the different possible moves, left right, diagonal, weird L-shaped thing, ...*),  $P(s_{t+1} \mid s_t, a_t)$  is the transition function which describes how the environment acts in response to the past ( $s_t$ ) and to your actions ( $a_t$ ) (*in this case, your opponent's moves, and the results of your actions*), and finally,  $r(s_t, a_t, s_{t+1})$  is the reward function, (*whether you won (+1) or lost (-1) the game*). The player's goal, is to learn a policy (which chooses actions,  $a_t = \pi(s_t)$ ) that yields the largest return (accumulation of rewards).

This setting is easily generalised to other, more interesting, applications. For example, if we weaken the requirement on that the learners observation  $x_t$  is fully describes the state  $s_t$ , then we could apply this framework to partially-observable decision processes, such as StarCraft II or self-driving cars, where we see only a small portion of the truly complex environment.

*(originally inspired by research in psychology? < story to be told here. Thus "good" actions can be reinforced by receiving positive rewards, while "bad" actions can be punished.)*

It should be noted that the problems that make RL hard are; "trial-and-error search and delayed rewards" [Sutton and Barto]. Unlike supervised learning, which gives the learner feedback (*I think that digit is a 5 -> no, its a 6*), in RL the learner only receives evaluations (*I think that digit is a 5 -> wrong*). Ontop of terse teachers, many actions may be taken before any evaluation is done. Leaving you to wonder "what did I do to deserve this?" pigeon superstition

### What is transfer learning (TL)?

An easy setting to understand TL is in the multi-task setting (but there are others, few-shot learning, continual learning and even distribution shift could be viewed as transfer learning between the past and present)

Imagine we have two classification tasks,  $A, B$ , each of which consist of pairs of observations and targets. Task =  $\{(x_i, t_i) : i \in [1 : N]\}$ . Now, a learner,  $f$ , trained on task  $A$  and achieves a loss,  $\mathcal{L}$ ,

---

when evaluated on task  $A$ . We denote this as:  $\mathcal{L}^A(f_A)$  (subscript denotes training task, superscript denotes evaluation task).

We say transfer has occurred when  $L^A(f_A) > L^A(f_{B \rightarrow A})$  (the goal is to minimise loss). That is to say, that training on  $B$ , and then training on  $A$  improves performance on  $A$  (aka pretraining). Similarly, we could transfer in the opposite direction,  $L^A(f_A) > L^A(f_{A \rightarrow B})$ . Training on  $A$  and then training on  $B$  improves performance on the original task,  $A$ . (note this is actually quite hard, for example see EWC)

A familiar example to some is, learning to eat with chopsticks (after already knowing how to use cutlery). We do not need to relearn to use our arms, we remember where our mouths are, and how to chew (we transfer that knowledge). We also remember that the goal of eating is to put food, not too much, into our mouths (we transfer that knowledge). But, we do need to learn some new grips and fine motor skills that will allow us to pick up food. (we also transfer the knowledge of “how one might pick something up” from other places as well!)

(WILL want to cover these as well. not sure how)

- What about efficient transfer? Increasing the speed of learning -> aka meta learning.
- Captures the idea of learning something more abstract about the similarities between two distinct tasks. Closely related to/IS generalisation.

## Why do we care?

(WILL) Applications. Existing research. Blah Blah. Transfer is important because!? RL is important because?

If we could get model-based RL to work, it would take us a long way to solving many of our problems (which ones?). But, in the real world, our ability to learn models (let alone plan with them) is still an unsolved problem.

To build these models... will need to do so efficiently. Transfer!

Where TRL can be applied is dependent on where RL can be applied. Thus if we want to scale TRL to the “real world”, we first need to extend RL to the real world. Some of the key differences between the places RL has been successful so far (for example; Go, Atari ALE) and the real world are;

- continuous action spaces (especially important for robotics),
- resource constraints (memory - online learning, speed - real-time machine learning, ...),
- partial information (where model-based RL becomes necessary),
- long-term dependencies

---

Finally, the real world is a lot more complex than these games we test our algorithms on. If we want to scale to the real world, we need to make our algorithms more data efficient. Transfer learning is a solution.

## Transfer as decomposition

One way to achieve transfer learning is to: decompose complex observations into a set of modules/atomic-factors. Or in other words, to disentangle independent factors. (need refs!!!)

For example;

- we could take a symmetry group, and rewrite it as the composition of a single atom with various transformations. Thus we have decomposed the group into its parts.
- Aka finding the basis.
- Decomposition of objects and relations?
- (WILL) more examples

The belief that complexity can be decomposed, that observations are built from smaller/few parts, is at the heart of unsupervised learning, and of science. For example, it is common to assume that there are a set of latent variables that combine (non)linearly to generate observations. Within science, reductionism and mediation analysis are key tools for producing understanding. They allow us to break down complexity into simple parts we can study and understand.

Although, it should be noted that, even if we can decompose a complex phenomena, there is a lot of meta information required to actually use this decomposition for planning. Which modules do what? How are they related? When should I apply a given module?

*(Note, there are other approaches to transfer(?). Learning a metric? Distillation? EWC? MDL regularisation? ??? (but maybe underneath they are all related!?) )*

## How to decompose?

As far as I know, there are two ways to build priors into a learner; explicitly or implicitly. Explicit priors normally mean structural constraints and regularisers. Implicit priors can mean the (unintended) priors lurking within an optimisation algorithm (for example SGD's bias towards flat minima Poggio et al. 2018).

We can build a decomposition into the structure of our model, for example, by specifying that different loss functions apply to subsets of parameters (for example the generator and discriminator in a GAN learn from different losses). Another example could be ensembles, where to get a decomposition we

---

can feed learners certain (different) features, forcing them to do different things. Often, these structural choices come from domain knowledge, beliefs we, as people, have about how the task should be done.

Alternatively, we can regularise our model to decompose its inputs. For this we need a differentiable measure of disentanglement. This approach is appealing as it should require less domain knowledge.

## **Decompositions in RL**

*(Let's explore some popular decompositions in RL.)*

Probably the most fundamental decomposition used is model-based RL. Where a reinforcement learner is decomposed into a transition function (a model of the environment) and a policy (how to act in the environment). This facilitates transfer as, the learner can be given a new task, requiring a new policy, (while remaining in the same environment) and simply reuse its model. For example, asking a learned-chess-master to now lose a game of chess on purpose (notably model-free learners like DQN cannot do this).

There are other examples of existing work attempting to use a decomposition to facilitate transfer, for example;

- Feudal networks decompose the learner's policy into a manager and a worker who work a different temporal scales. This could allow management strategies be to shared between different worker's domains.
- Learning to learn by GD by GD decompose learning into, learning a teacher and a student
- Modular meta-learning constructs a single function from many modules and must learn the modules and when to use them.

It is possible to structure decompositions in many ways; heirarchcally, as an ensemble, as some sort of graph structure, ... As of yet, there is no general formula that tells us when a decomposition of a learner is going to be optimal.

(WILL) relationship to grounding the modules. might put something in here?

## **Proposed research**

(WILL) In my mind, the goal of scientific research is understanding. But how can this be achieved? Necessary parts appear to be; demonstration and communication. Therefore ...? I am going to do...

What would a theory of transfer learning do? What knowledge would mean that we understand transfer?

---

## High level goals

I would like to understand generalisation and transfer.

- when can knowledge be transferred from  $A$  to  $B$ ?
- is decomposition necessary for transfer? Is it sufficient?
- how (computationally) hard is it to find abstract similarities between two domains?
- what problems can a “decomposer” solve that another cannot?
- which representation of knowledge is optimal for my setting?
- how is a transition fn like a metric (which can be used to build a manifold)?
- how can we learn “algebras” for composing modular systems?

But I am unsure how to approach these questions and which tools are the best for thinking about them. While mulling these over, I would like to explore decompositions and model-based learning, applied to the Atari ALE (or similar). My hope is to see patterns and find interesting problems for understanding transfer.

## Benchmarks

One of the reasons, in my opinion, that the machine learning field has progressed quickly is its use of shared benchmarks on currently out-of-reach problems. ImageNet is a great example of this for the computer vision community, and the Atari ALE is currently the canonical benchmark for the RL community.

By definition, TRL requires the testing of learners on many different tasks. But this can take thousands of CPU/GPU hours to evaluate... I think there is a need for cleverly designed benchmarks that don't require thousands of dollars. Currently TRL is a game for players with resources, which means that only Google and Google-Deepmind get to play.

I am imagining a set of as-simple-as-possible settings where we can explore the limits of transfer and decomposition (similar to ai-safety-worlds for the AI safety community).

## Specific questions to explore

*(these may be ill-posed, trivial, or solved, but hopefully I will find out soon...)*

Decompositions

- Can we formalise what we mean by “decompose”? Can we differentiate it? What is its relationship to independence criterion and ICA?

- 
- Meta-RL trains a learner on the aggregated return over many episodes (a larger time scale). If we construct a temporal decomposition (moving averages at different time-scales) of rewards and approximate them with a set of value functions, does this produce a hierarchy of meta learners?
  - Imagine you are given two models  $f, g$  that might predict pedestrian and traffic behaviour respectively. How can you safely/sensibly combine their predictions? If the models were provided as densities we could evaluate the next step as  $s_{t+1} = \operatorname{argmax}_{s_{t+1}} \prod_i p_i(s_{t+1} | s_t)$ .
  - How can we train modular systems when there is often no gradient defined (as the module was not used)? Using “counterfactual” credit assignment of what might have happened?

#### Model-based learning (with partial information)

- Build a differentiable neural computer with locally structured memory (start with 1d and then generalise to higher dimensions). Is the ability to localise oneself necessary to efficiently solve partial information decision problems? Under which conditions does the learned index to a locally structured memory approximate the position of the agent in its environment.
- When attempting to learn a model, the agent uses an exploration policy. This policy may influence the dynamics observed, thus we need to use off-policy methods to correct for the effects of exploration actions. (The model must somehow disentangle the agents policy, and its effects, from the dynamics of the system)
- Inverse energy learning. Inspired to inverse reinforcement learning, what if we assume that the observations we make are the results of some optimal action, in the case of an energy being minimised,  $\Delta x = -\eta \frac{\partial E}{\partial x}$ .
- While learning a model  $s_{t+1} = \tau(s_t, a_t)$  is useful. It is more useful to know how to get around using that model. For example, I want to get to  $s^k$ , how can I do that considering I am in another state,  $s^i$ ? Want a function  $f(s^i, s^k) \rightarrow \{a_1, a_2, \dots, a_n\}$  that outputs a sequence of actions. You need to know how to get around... (HER?)

#### Planning with a learned model (with continuous actions)

*(not going to make it to these. but it is important to remember the ultimate goals of TRL. we dont just care about learning models, additionally, they need to useable for efficient planning)*

- If a model is being learned online how can we efficiently update value estimates computed using the old model?
- How can you backpropagate gradients through the argmax functions required for planning?
- If I am using an imperfect learned model to generate plans, how can I ensure that I do not plan for “fantastic” outcomes (aka they are fantasies). (note, this is closely related to reward hacking)
- Construct a planner that can control the models computational complexity by asking it to provide more approximate solutions (via accuracy masks, for the size of time-steps, or ...?).

---

## Timeline

I have allocated time for 8 “sprints” (the bullet points above; decompositions and model-based learning), each of 2 weeks. The goal of each sprint will be to;

- motivate the idea as a solution to an existing problem,
- prove that the “existing” problem really exists,
- generate alternative solutions and a suitable baseline,
- design the minimal viable experiment to falsify the proposed solution,
- implement the experiment if feasible.

Future work will depend on the results of the sprints.

## Proposed deliverables

- Tutorial(s) on “core” RL,
- A learned model that benefits from transfer between environments (Atari ALE or similar),
- An essay on the future of model-based RL,
- The definition and construction of a new benchmark for TL,
- Implement/reproduce a TRL paper
- A thesis documenting my work

Hopefully a few papers, but that is conditional on making discoveries.

## References

Dumoulin, Vincent, Ethan Perez, Nathan Schucher, Florian Strub, Harm de Vries, Aaron Courville, and Yoshua Bengio. 2018. “Feature-Wise Transformations.” *Distill*. doi:10.23915/distill.00011.