
LESS IS MORE - THE DISPATCHER/ EXECUTOR PRINCIPLE FOR MULTI-TASK REINFORCEMENT LEARNING

Martin Riedmiller

Google DeepMind

UK

riedmiller@google.com

Tim Hertweck

Google DeepMind

UK

thertweck@google.com

Roland Hafner

Google DeepMind

UK

rhafner@google.com

Abstract

Humans instinctively know how to neglect details when it comes to solve complex decision making problems in environments with unforeseeable variations. This abstraction process seems to be a vital property for most biological systems and helps to 'abstract away' unnecessary details and boost generalisation. In this work we introduce the dispatcher/ executor principle for the design of multi-task Reinforcement Learning controllers. It suggests to partition the controller in two entities, one that understands the task (the dispatcher) and one that computes the controls for the specific device (the executor) - and to connect these two by a strongly regularizing communication channel. The core rationale behind this position paper is that changes in structure and design principles can improve generalisation properties and drastically enforce data-efficiency. It is in some sense a 'yes, and ...' response to the current trend of using large neural networks trained on vast amounts of data and bet on emerging generalisation properties. While we agree on the power of scaling - in the sense of Sutton's 'bitter lesson' - we will give some evidence, that considering structure and adding design principles can be a valuable and critical component in particular when data is not abundant and infinite, but is a precious resource. A video showing the results can be found at <https://sites.google.com/view/dispatcher-executor>.

1 Introduction

Reinforcement Learning (RL) has matured from single task settings to multi-task scenarios, where a single controller is asked to solve one of many tasks. Multi-task requirements commonly occur in robotics, where e.g. a robot arm is asked to manipulate a wide range of objects in different ways, or a locomotion platform should reach different goals. Multi-task requirements also occur in classical control applications, e.g. controlling fusion plasma to various shapes (Degraeve et al., 2022) or set-point control problems in general (Hafner & Riedmiller, 2011). The typical approach for multi-task control is to use the 'classical' controller structure, represented by a monolithic neural net, that is extended by a pathway to specify the requested task.

However, it requires different kinds of knowledge for different aspects of the control task: Task understanding, e.g. understanding which objects in a scene are affected or which type of manipulation is requested, requires a general understanding of the world, whereas computing a concrete control action requires specific knowledge about the device at hand.

In this paper we argue, that multi-task RL controllers can massively benefit from a bi-partite controller structure, that cleanly separates the overall architecture into one part that semantically understands the task (the 'dispatcher') and a second part that executes the input-to-action policy to control the device (the 'executor'). A further crucial point is that the communication between dispatcher and executor should enforce compositionality and a substantial reduction of irrelevant information about the situation. This forced abstraction from the details allows the same executor

to be applied in a broad diversity of original scenes and thus boosts generalisation - hence the title 'less is more'.

We propose the dispatcher/ executor (D/E) principle for multi-task control and show a concrete implementation in the domain of robotic manipulation and discuss its properties. We provide a set of empirical evaluations and show their significant benefits on a range of typical robotic manipulation tasks in both simulation and on the real robot. The dispatcher/ executor principle can even be applied in hindsight, i.e. if a controller for a certain concrete task has already been learned. We show this hindsight transfer to a D/E structure on the example of a real robot policy that is transferred from stacking 'red on blue' objects to a general object stacking routine.

This position paper makes the following contributions:

- introduces the dispatcher/ executor (D/E) principle as a design principle for multi-task RL controllers
- makes a concrete proposal for such a dispatcher/ executor architecture in the area of robotic manipulation – both as a proof of principle and as a concrete practical implementation proposal
- provides empirical evaluations in simulation and on a real robot demonstrating huge gains in data-efficiency in multi-task RL settings

2 State of the Art

The core claim of Rich Sutton’s insightful 'Bitter Lesson' paper (Sutton, 2019) is simple and compelling: 'general methods that leverage computation are ultimately the most effective'. This has been demonstrated several times in the history of AI, e.g. in game playing (Silver et al., 2016), computer vision (Krizhevsky et al., 2017), or more recently, in large language models (Vaswani et al., 2017): huge amounts of data trained in huge neural networks with huge amounts of compute lead to unprecedented and surprising generalisation (OpenAI, 2023), while the underlying training principle is general and relatively simple.

Given the big success in language, it is evident to investigate, if the same principles of scaling up to large amounts of data and the use of large networks show the same generalisation effects in control. The GATO (Reed et al., 2022) approach demonstrated that a reasonably large transformer style neural network is able to show good control performance with a single network on a variety of heterogeneous tasks from Atari to robotics. RoboCat (Bousmalis et al., 2023) used a related approach focusing on robotic control with the final goal to train a versatile model that can transfer to arbitrary new tasks. RT-1 (Brohan et al., 2022) and RT-2 (Zitkovich et al., 2023) start from a pre-trained language model and therefore have a built-in semantic understanding of tasks and the environment. All of the above models are trained on vast amounts of data using a transformer-style controller architecture with an additional input pathway to communicate the task. Since these methods are by default extremely data-hungry, data-augmentation is one path to generate even more data out of the data collected (Yu et al., 2023) and by this enforce generalisation within the distribution of the data generated.

In Ahn et al. (2022) it is shown that a pre-trained language model and a set of pre-trained skills can be efficiently combined by learning to sequence these skills, based on the assumption that language model and the skills share some common understanding of the main task. While this gives some form of general task execution in known environments, a large amount of data has to be used to make sure the pre-trained skills will be able to execute on a sufficient enough set of variations of the environment.

Another scheme that is very common for data efficient multi task approaches is to learn some form of task embedding that allows to generalise in the task space (e.g. Hausman et al. (2018); Nachum et al. (2018)). While some of these ideas can be combined with the D/E approach, these methods solely

focus on the command part in the communication channel. Therefore the expected generalisation by abstracting the observation and the command, as we demonstrate here, is missing.

While we agree that scaling is one way to achieve more and more general controllers, we at the same time believe that even large scale systems should be 'data-efficient', i.e. should make the best use of the data available. The core thesis of this paper is that incorporating structure into the architecture and the communication channel can drastically enhance the generalisation capabilities 'by design'. In particular, enforcing the right representation within the controller seems to be an important enabler. Recent examples of this line of research can be found in [Levine et al. \(2016\)](#); [Pinneri et al. \(2023\)](#); [Groth et al. \(2021\)](#).

The idea of a two (or more) level architectures for control has been explored in many different variations (e.g. [Vezzani et al. \(2022\)](#); [Zhang & Chai \(2021\)](#); [Riedmiller et al. \(2018\)](#); [Hafner et al. \(2022\)](#)). The D/E principle intentionally defines the separation from the viewpoint of the type of knowledge required to run the components: a dispatcher component that understands the task through general world knowledge (that can be trained on general data like text, images or videos), and an executor component that specifically understands the device (trained through active interaction with the device). The restrictions on the communication channel enforce an abstract information exchange between these types of knowledge that boosts generalisation. We believe that in particular in this new era of large language models that are able to represent general world knowledge, this structural principle can play an important role to boost learning control on real robots, where data-efficiency with respect to device interaction is a key requirement.

3 The Dispatcher/ Executor Principle for Multi-task RL

3.1 The general idea

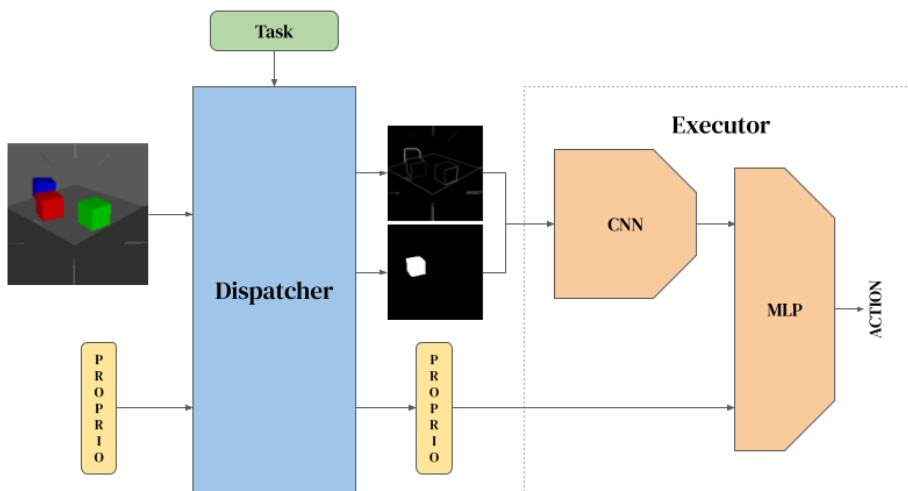


Figure 1: Illustration of the D/E architecture. The controller consists of two modules, the dispatcher and the executor. The dispatcher gets the current observation and modifies it according to the task description. This abstract message is sent to the executor, which computes the according action for the device.

The proposed dispatcher/ executor (D/E) principle for control makes two main points: 1. separate the control architecture into two entities according to the type of knowledge they require (dispatcher

and executor), and 2. define a suitable communication channel between those two modules with strongly regularizing properties.

The two parts are defined through their specific duties: the dispatcher is responsible for semantically understanding the task and calling the executor with the respective instructions. The executor is responsible for computing the actual control signal that drives the device to achieve its goal communicated from the dispatcher. Thus, whereas the dispatcher needs an understanding of the world in general to be able to understand the task, the executor needs an understanding of how the device reacts to its control signals. Whereas in this paper we mainly focus on a single dispatcher and a single executor, in its generalized form, the controller might have several executors that are specialized on realizing different skills. E.g., to fulfill a complex task, the dispatcher might call several executors in a row.

The second crucial point of the D/E principle is the design of the communication channel between dispatcher and executor. It has two relevant properties: firstly, it should enforce compositionality of the transmitted command information, and secondly, it should reduce the transferred information to the minimum that the executor needs to successfully perform the task. The requirement of compositionality might be achieved by forcing some kind of structure into the communication - e.g. in form of a certain artificial language - that forces the exchanged information to demonstrate some kind of abstraction.

Enforcing this language-like communication we imagine as a property that is incorporated into an agent through structural design (e.g. found through evolution in a natural agent), and finding the right words of this language is something that finally is learned by the agent through experience and shaped through an according loss function.

It is beyond the scope of this work to figure out, if and how such an interface can be found through evolution and/or learning. We will instead focus on proposing a concrete implementation of a D/E architecture and the communication interface in the spirit of an existence proof and show the concrete benefits on the example of robotic manipulation.

3.2 A concrete D/E implementation for robotic manipulation

Typical tasks in robotics involve the manipulation of the environment, like e.g. 'push the button', 'lift the screw', 'put the bottle on the table', ... Following the D/E principle, we suggest the following concrete implementation (see figure 1): the dispatcher receives the task description and the raw sensor observations at each time step, which we assume includes camera images and optionally other sensor information, e.g. proprioceptive values. The dispatcher analyses the task description, infers what to do and identifies the targeted objects in the image. It then encodes (reduced) information about the target objects and environment as arguments for the executor. The message to the encoder takes the following form: $ExID, arg_1, \dots, arg_n$ where $ExID$ identifies the executor and arg_1, \dots, arg_n contain separated and compressed information about target object(s) and scene.

There are different ways now to specify a targeted object: an obvious way would be through 6D pose coordinates, which of course are only available if we have ground truth, which is most often not the case. Since the current scene is represented through one or more camera images, we instead suggest to mark object information in pixel space. Concretely, we suggest to code the information about the target object through a simple masking operation: all pixels in the observed image that belong to the object are classified with 1, the rest is set to 0. This operation reduces the information drastically to what the executor might need to locate the object and infer pose and shape if needed. Note, that this information now abstracts away further properties of the object, like type or color. The executor therefore works immediately for a broad range of objects that may appear very differently in the original observation space (here an RGB image). We call this the 'mask' filter of a target object for further reference.

In addition, some general information about the scene might be required as well. Think of a target object lying close to another object - a situation which may require special handling. In the sense of

reducing information to a minimum required to solve the task, we here suggest to run the full image through an edge detector and provide the result as a further argument to the executor (referred to as 'edge' filter in the sequel).

An alternative representation of target objects can be done through 'pointing' to the object: A simple implementation of this is to compute the pixel centroid of the mask of the object and then represent this through a 5x5 blob at the respective location in the image. Again, all other pixels are set to 0. This coding scheme is called the 'pointer' filter in the following.

The above operations of identifying target objects and providing background information are applied to all camera images available, to have full information of the scene.

We admit that the proposed encoding scheme is very much tailored to the typical tasks that occur in robotic manipulation, and the proposed filter operations are just a small selection from many possible choices. This is intended, since here we mainly want to give a concrete example of the existence of such an encoding and its usefulness in a certain application. We are convinced that the idea behind the approach transfers to other application scenarios as well, even more, that future work will aim at finding these coding schemes through learning rather than through design.

3.3 Workflow and learning within the D/E architecture for manipulation

The overall workflow in the proposed D/E architecture is as follows: the dispatcher gets the full observation vector plus the task description as the input. It parses the task description and identifies the objects involved. Then it computes the respective pixel frames for the arguments (e.g. using one or more of the above filter schemes, masks, pointers, edges). This information is then transferred to the executor in every control step.

Learning in a D/E architecture is currently restricted to the executor. Learning is not different from the standard monolithic controller, apart from the fact, that the input is represented differently. The executor can therefore be trained in many different ways. In the following we investigate Reinforcement Learning from scratch as well as the distillation of a previously learned policy to a D/E architecture in a student/teacher setup. Another interesting direction, which is beyond the scope of this paper, is the investigation of learning multiple variants of representations in an off-line RL setting of pre-collected data (Lange et al., 2012; Riedmiller et al., 2021).

4 Empirical investigations of targeted properties

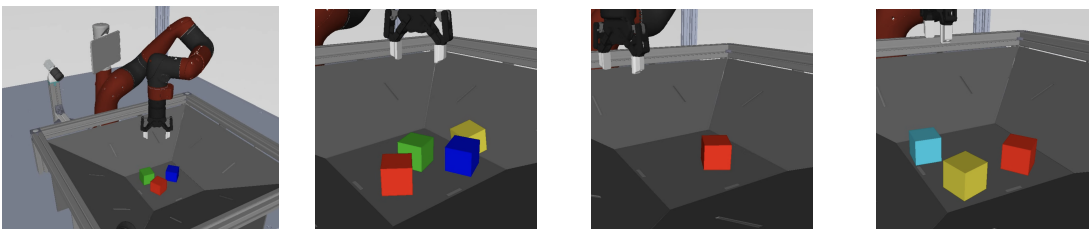


Figure 2: Simulation setup. From left to right: original setup with 3 objects, "Four cubes", "One cube", "Recolor". The controllers are trained to lift the red object (single task), the red, left or green object (multi-task), and evaluated also in non-training situations.

In this section, we highlight some properties of the D/E architecture on various learning experiments. It is not meant as an in-depth empirical evaluation, it rather gives exemplary insights, of how the proposed D/E structure of the controller can improve learning and generalisation properties in comparison to a classical monolithic control structure. We provide experiments both in simulation and on the real robot on two set of tasks, lifting an object (in simulation) and stacking of two objects (real robot). If not otherwise noted, 'standard' refers to a classical, monolithic neural network structure which is extended by an extra-input for task communication.

The dispatcher-executor architecture (see figure 1) uses a hardcoded dispatcher and a learned executor. The dispatcher gets the task specification (e.g. 'lift green object' or 'stack red on blue object') and then segments the according target objects - in this case using simple color segmentation (however, it is straightforward to replace this by more capable segmentation schemes, e.g. [Oquab et al. \(2023\)](#)). In addition, it filters the whole scene using an edge operator. Then, the executor is called, providing the respectively filtered images as separate arguments. If not otherwise noted, 'D/E' refers to the default setting of using a segmentation mask for the objects, and the edge operator for the overall scene (as described in 3.2).

For training both the standard architecture and the executor in the D/E case, a distributional variant of MPO [Abdolmaleki et al. \(2018\)](#), with a Mixture-of-Gaussian critic was used. The agent receives the images of three cameras placed around the robot's workspace, as well as proprioceptive observations. Both architectures employ a three-block ResNet that computes an eight-dimensional embedding vector for each pixel input. These embeddings are concatenated with all proprioceptive observations and fed through a three layer MLP for computing the agent's five-dimensional actions.

4.1 Single task: lifting the red cube (simulation)

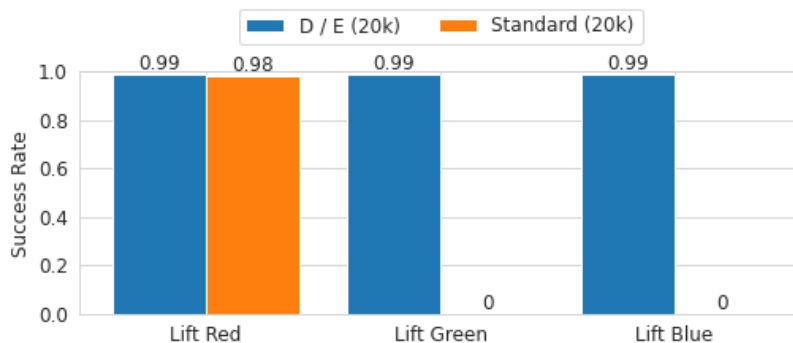


Figure 3: Single task training: 'lift red'. Both standard architecture and D/E architecture are trained for 20,000 episodes. Figure shows evaluation results for different tasks. The standard architecture solves the learned task only, whereas the D/E architecture can also lift green and blue cubes

For the first experiment, the scene consists of three cubes with different colors (red, green, blue; see figure 2). The single task setup is to learn to lift the red cube (blue and green cubes are mere distractors). Figure 3 shows the results. Firstly, both approaches (standard and D/E) need about 20,000 episodes to train the task successfully (98 % success in evaluation), i.e. the number of training episodes to get to a successful 'lift red' policy does not change significantly with the controller structure used.

Transfer to new tasks: since the input representation of the executor abstracts away the color of the cubes, the D/E architecture can transfer immediately to the tasks lift green or lift blue by setting the executor input respectively through the dispatcher module. Therefore the D/E controller is able to additionally perform 'lift green' and 'lift blue' with similar success rates - without any further training and data collection. The standard monolithic architecture is specialized to the single task and lacks the capability of transfer.

Findings:

- for a single task, number of training episodes for standard and D/E architecture are in the same range
- the D/E architecture can immediately transfer to new tasks ('zero effort transfer')

4.2 Multi-task 1: Lifting one out of three cubes (simulation)

Like above, the scene consists of three cubes, and now the training task is to alternately lift one of the cubes specified by its color. In the standard controller, we specify which cube to lift through an additional input. For the D/E controller, the dispatcher was adjusted accordingly and is now selecting the respective representation of the targeted object and providing the respective information to the executor.

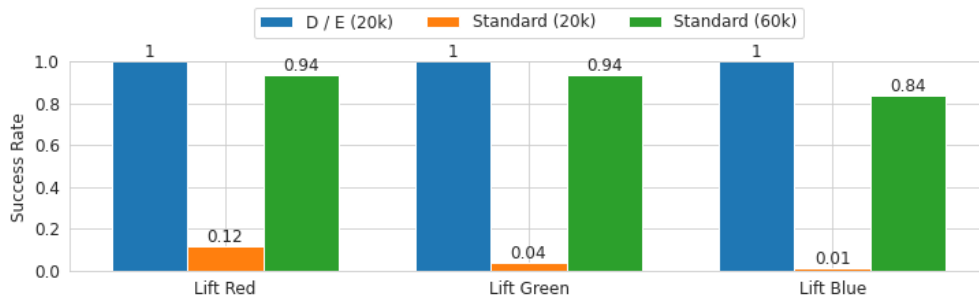


Figure 4: Multi task training: 'lift red/ green/ blue'. After 20k episodes, the D/E controller masters all three tasks, whereas the standard controller still performs purely (orange). After 60k episodes using 3 times as much training data, the performance of the standard controller (green) increased significantly, but still does not reach the D/E results.

The results are shown in figure 4. The D/E architecture can solve the three tasks in about 20k episodes, as in the single task setup before. This is due to the fact, that for the policy, the control task remains the same, no matter, which cube needs to be lifted. It therefore can fully leverage the data of all tasks and learn the lifting task quickly. On the other side, the standard architecture after 20k episodes, still sees pretty low success rates (between 1 and 12 %). It needs much more experience to solve the three tasks, and only after 60k episodes, reasonable success rates are achieved. We even had to increase the network size in comparison to the single task experiment to get to this result.

Findings:

- in the multi-task setup, the D/E architecture can leverage the abstraction and learn much more efficiently than the standard architecture. It can focus on learning the control task rather than the distinction between the tasks and this gives a dramatic speedup in learning.

4.3 Multi-task 2: Lifting one out of three objects with multiple shapes (simulation)

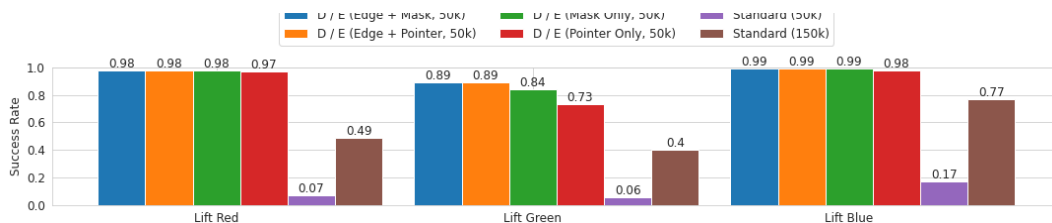


Figure 5: Multi task training: 'lift red/ green/ blue' with varying object shapes are used. D/E architecture with different representations performs well after 50k episodes. Standard architecture has low performance after 50k episodes and even after 150k episodes is worse than D/E.

One critical question with a reduced representation is, whether the carried information is still rich enough to cope for desired variances in the task setup. To look a bit deeper into this, in the following the objects are additionally varied in shape. This problem is considerably more complicated than grasping cubes, forcing the controller to deliberately choose grasping points depending on the object

shape. For example, a rectangular shaped object needs grasping points different from a triangular shaped one. For this experiment, we used the objects of the RGB stacking benchmark (Lee et al., 2021), see figure 7, right. Please note, that the distribution of shapes is different for different colors. This explains the varying success rates for different colors, e.g. the green objects are more challenging to lift in average.

Results are shown in figure 5. We tested four different variations of representation for the D/E architecture: the typical setup with mask and edge filter, then a setup with pointer and edge filter, and both a mask filter and pointer without scene information. All setups are trained for 50,000 episodes. Both mask and edge as well as pointer and edge work very reliably and lead to comparable results. This indicates, that there is some freedom for choosing a good representation. One can see some performance drop, if only the target object is communicated - this drop is more significant for the pointer only case (which has fewer information) than for the mask only setup. The standard architecture performs significantly worse - after 50,000 episodes only between 6 % and 17 % success rate is achieved. Using three times more data, after 150,000 episodes the network performs better, but is still significantly worse than the D/E architecture. Again, this stresses the hypothesis, that bringing in structure can drastically improve the learning process.

Findings:

- the reduced representation of the D/E architecture should be kept rich enough, such that it enables the policy to adapt to varying requirements.
- different representation schemes might lead to successful behaviour

4.4 Robustness with respect to scene variations

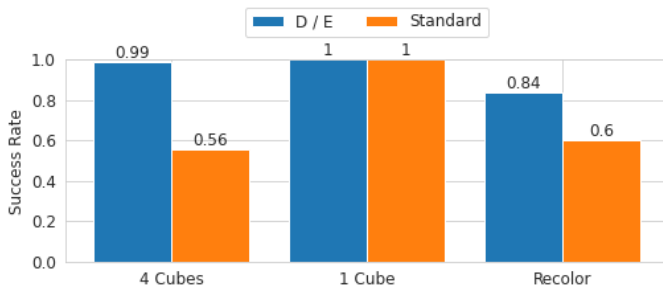


Figure 6: Evaluation of robustness. The controller is trained and evaluated on 'lift red'. Evaluation happens on varying background scenes, with only one cube in the scene, with four cubes or with re-colored distractor objects.

One potential challenge in rich scene representations is that the controller might concentrate on task irrelevant features and rely on them for successful execution. We expect, that a reduced representation as it is used in the D/E architecture, will be more robust to variations in the setup of the task. As a straightforward experiment we investigate, how the single task controllers trained in section 4.1 react to slight changes in the setup. In particular, we investigate the performance, when changing the number of objects in the scene from three to one or four respectively, or change the color of the distractor objects.

As shown in figure 6, both the standard controller and the D/E architecture are robustly performing in the '1 cube' case. If instead four objects are put into the scene, the standard controller significantly loses performance, whereas the D/E architecture still reliably solves the task. Similar behaviour can be observed when the distractor objects have different colors.

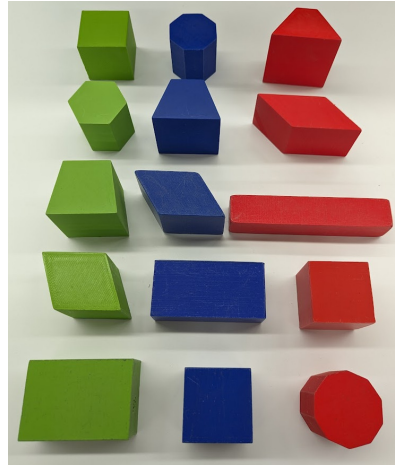
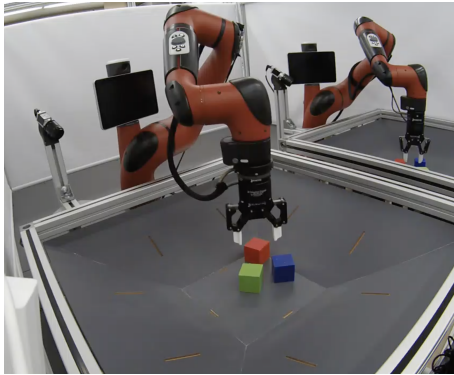


Figure 7: Left: the real robot environment. Right: the RGB sets 1-5 (top to bottom). Original task is to stack the red object on the blue object, extended task is to stack any object on any arbitrary object of one set, which evaluates generalisation over color and shape.

Findings:

- the reduced representation in the D/E architecture shows more robust behaviour when the environment conditions change.

4.5 Real Robot Experiments - Teacher/ Student setup

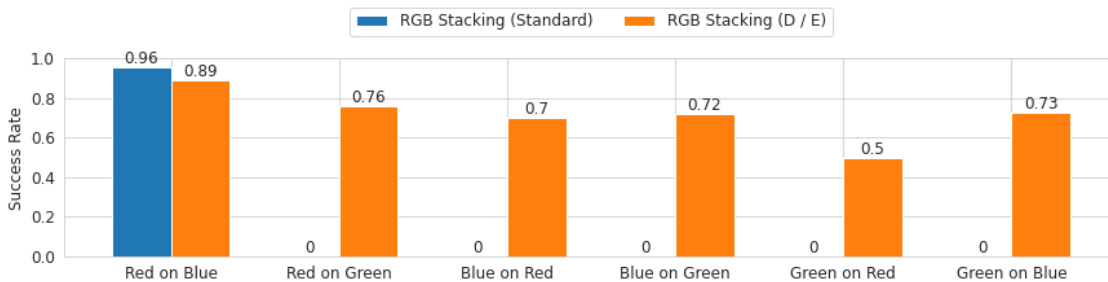


Figure 8: Real robot experiments: Cloning a 'red on blue' policy to the D/ E architecture. The teacher performance was an 96 % success rate on average for 'stack red on blue' of the RGB stacking task

To add another angle of use cases of the D/E architecture, we examine the use of a D/E controller in a student-teacher setup. Concretely, we take a policy that was trained to solve a single task on a real robot ('stack red object on blue object') and distill it to a controller following the D/E principle. A successful transfer will allow to open up a whole class of tasks ('stack object X on object Y'), without collecting any further training data - a potential huge gain in data-efficiency.

The dispatcher in this case is restricted to find colored objects and advise the executor accordingly (by setting the right mask), but it is obvious how to extend this to more general objects, e.g. exploiting the recent advancement in general segmentation algorithms, e.g. (Oquab et al., 2023).

As a teacher policy, we took a policy that has been fully self-learned on real robot experience to stack the red object on the blue object on the five different object sets of the RGB stacking benchmark (Lampe et al., 2023). Note that this was learned from ground off directly from pure experience and no additional simulation data or tele-operation data was used. The performance of the original policy was about 96% in average over all 'red on blue' object sets 1-5.

For evaluation of the distilled policy, we did two sets of experiments, one based on similar cubes and one on the RGB sets 1-5, where each object has a different shape. The results are shown in figure 8. For cubes, the stacking performance achieved high success rates over all color combinations. For the RGB sets, 'red on blue' achieves the highest result close to 90%. Satisfyingly, not much performance was lost in the distillation process. Also for the other color combinations, average success rates of 50% or over are achieved. This is a strong result, given that the shapes in the different sets are quite challenging and stacking these combinations has never been seen during the training phase of the executor. The distilled D/E policy gives a powerful baseline for finetuning the policy to perfection in all combinations.

Findings:

- the D/E principle can be applied in student/ teacher scenarios, boosting existing policies to novel and more general application scenarios without the need to collect more data on the task.
- this highly capable stacking controller has been entirely self-learned from experience collected on the real robot.

4.6 The power of the dispatcher

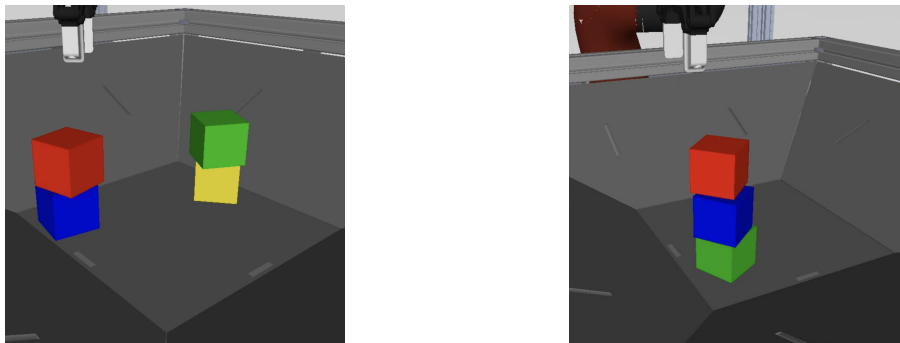


Figure 9: Zero-shot transfer of a D/E stacking policy to unseen task "Two Towers" (left) and "Triple Stack" (right) through adaptation of the dispatcher.

The dispatcher is used to understand the task description and translate it into the call of the executor with respective arguments. This is a versatile mechanism - for example, the dispatcher can call multiple executors sequentially in order to fulfill more complex, and long horizon tasks.

To demonstrate this, we are looking into two tasks: two-towers and triple-stack.

For two-towers, we ask to build two separate towers out of 4 cubes in the workspace. The executor is a simple stack X on Y policy, trained according to the D/E principle with masks and edges. The dispatcher identifies the four objects and then sequentially calls stack object 1 on object 2 followed by stack object 3 on object 4. Here, the different executors are just called after a time limit is reached, but more complicated schemes are possible and will be examined in future work. Note that no further executor training is needed to achieve this. The immediate success rate is 46%. The main failure case is when the executor destroys the first tower when building the second, which is something it has simply not been trained to avoid. Finetuning the executor to respect existing towers will increase the success rate.

Triple stacking means stacking three objects on top of each other. This has been tried by other efforts using a single policy (Bousmalis et al., 2023). Here, we modify the dispatcher and ask to stack object 1 on object 2, followed by stack object 3 on object 1. Again, the executor is the basic D/E stacking policy. The observed success rate is 38%. The main reason for failure here is that the

second stack is not high enough (this was not seen during the original training phase). However, the initial performance is strong, and finetuning the executor with accordingly collected data will solve this task.

5 Discussion

This position paper introduces the idea of the dispatcher/ executor principle for structuring control architectures for multi-task reinforcement learning. It suggests a) the separation into a dispatcher module that contains general world knowledge and understands the task, b) an executor module that learns the specific interaction with a particular device and c) a regularizing communication channel between these two modules that allows for abstract and compositional communication.

We proposed a concrete implementation of such a D/E architecture for robotic manipulation. In two learning scenarios - Reinforcement learning from experience and hindsight transfer in a teacher/ student setup - we showed a massive boost in generalisation. Using the same amount of interactions with the device, the D/E architecture is able to transfer its control abilities to a whole class of related tasks.

While in this prototypical implementation many features are still engineered, next steps will concentrate on end-to-end learning architectures based on the D/E principle. The separation of the control architecture on the basis of the type of knowledge required in each part makes the dispatcher an obvious candidate to integrate large multi-modal models into the design. Their capabilities of representing general world knowledge and executing tools is what is needed to get from a task description to a tailored call of an executor. Another subject of future research is learning to find regularized representations that allow for interesting abstractions in the communication channel between dispatcher and executor.

Acknowledgments

Thanks to the Control Team and various colleagues at DeepMind for ongoing discussions.

References

- Abbas Abdolmaleki, Jost Tobias Springenberg, Yuval Tassa, Remi Munos, Nicolas Heess, and Martin Riedmiller. Maximum a posteriori policy optimisation. In *International Conference on Learning Representations*, 2018.
- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee, Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan, and Andy Zeng. Do as i can, not as i say: Grounding language in robotic affordances, 2022.
- Konstantinos Bousmalis, Giulia Vezzani, Dushyant Rao, Coline Devin, Alex X Lee, Maria Bauza, Todor Davchev, Yuxiang Zhou, Agrim Gupta, Akhil Raju, et al. Robocat: A self-improving foundation agent for robotic manipulation. *arXiv preprint arXiv:2306.11706*, 2023.
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl

-
- Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-1: Robotics transformer for real-world control at scale. In *arXiv preprint arXiv:2212.06817*, 2022.
- Jonas Degraeve, Federico Felici, Jonas Buchli, Michael Neunert, Brendan D. Tracey, Francesco Carpanese, Timo Ewalds, Roland Hafner, Abbas Abdolmaleki, Diego de Las Casas, Craig Donner, Leslie Fritz, Cristian Galperti, Andrea Huber, James Keeling, Maria Tsimpoukelli, Jackie Kay, Antoine Merle, Jean-Marc Moret, Seb Noury, Federico Pesamosca, David Pfau, Olivier Sauter, Cristian Sommariva, Stefano Coda, Basil Duval, Ambrogio Fasoli, Pushmeet Kohli, Koray Kavukcuoglu, Demis Hassabis, and Martin A. Riedmiller. Magnetic control of tokamak plasmas through deep reinforcement learning. *Nat.*, 602(7897):414–419, 2022. doi: 10.1038/S41586-021-04301-9. URL <https://doi.org/10.1038/s41586-021-04301-9>.
- Oliver Groth, Chia-Man Hung, Andrea Vedaldi, and Ingmar Posner. Goal-conditioned end-to-end visuomotor control for versatile skill primitives. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1319–1325. IEEE, 2021.
- Danijar Hafner, Kuang-Huei Lee, Ian Fischer, and Pieter Abbeel. Deep hierarchical planning from pixels, 2022.
- Roland Hafner and Martin A. Riedmiller. Reinforcement learning in feedback control - challenges and benchmarks from technical process control. *Mach. Learn.*, 84(1-2):137–169, 2011. doi: 10.1007/S10994-011-5235-X. URL <https://doi.org/10.1007/s10994-011-5235-x>.
- Karol Hausman, Jost Tobias Springenberg, Ziyu Wang, Nicolas Heess, and Martin Riedmiller. Learning an embedding space for transferable robot skills. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rk07ZXZRb>.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. *Commun. ACM*, 60(6):84–90, 2017. doi: 10.1145/3065386. URL <https://doi.org/10.1145/3065386>.
- Thomas Lampe, Abbas Abdolmaleki, Sarah Behtle, Sandy H. Huang, Jost Tobias Springenberg, Michael Bloesch, Oliver Groth, Roland Hafner, Tim Hertweck, Michael Neunert, Markus Wulfmeier, Jingwei Zhang, Francesco Nori, Nicolas Heess, and Martin Riedmiller. Mastering stacking of diverse shapes with large-scale iterative reinforcement learning on real robots. *arXiv preprint arXiv:2312.abcde*, 2023.
- Sascha Lange, Thomas Gabel, and Martin Riedmiller. Batch reinforcement learning. In *Reinforcement learning: State-of-the-art*, pp. 45–73. Springer, 2012.
- Alex X. Lee, Coline Devin, Yuxiang Zhou, Thomas Lampe, Konstantinos Bousmalis, Jost Tobias Springenberg, Arunkumar Byravan, Abbas Abdolmaleki, Nimrod Gileadi, David Khosid, Claudio Fantacci, Jose Enrique Chen, Akhil Raju, Rae Jeong, Michael Neunert, Antoine Laurens, Stefano Saliceti, Federico Casarini, Martin Riedmiller, Raia Hadsell, and Francesco Nori. Beyond pick-and-place: Tackling robotic stacking of diverse shapes. *arXiv preprint arXiv:2110.06192*, 2021.
- Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies, 2016.
- Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical reinforcement learning, 2018.
- OpenAI. Gpt-4 technical report, 2023.

-
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- Cristina Pinneri, Sarah Bechtle, Markus Wulfmeier, Arunkumar Byravan, Jingwei Zhang, William F Whitney, and Martin Riedmiller. Equivariant data augmentation for generalization in offline reinforcement learning. *arXiv preprint arXiv:2309.07578*, 2023.
- Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gómez Colmenarejo, Alexander Novikov, Gabriel Barth-maroon, Mai Giménez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. A generalist agent. *Transactions on Machine Learning Research*, 2022.
- Martin Riedmiller, Roland Hafner, Thomas Lampe, Michael Neunert, Jonas Degraeve, Tom Wiele, Vlad Mnih, Nicolas Heess, and Jost Tobias Springenberg. Learning by playing solving sparse reward tasks from scratch. In *International conference on machine learning*, pp. 4344–4353. PMLR, 2018.
- Martin A. Riedmiller, Jost Tobias Springenberg, Roland Hafner, and Nicolas Heess. Collect & infer - a fresh look at data-efficient reinforcement learning. In Aleksandra Faust, David Hsu, and Gerhard Neumann (eds.), *Conference on Robot Learning, 8-11 November 2021, London, UK*, volume 164 of *Proceedings of Machine Learning Research*, pp. 1736–1744. PMLR, 2021. URL <https://proceedings.mlr.press/v164/riedmiller22a.html>.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484, 2016.
- Richard Sutton. The bitter lesson. *Blog Post*, 2019. URL <http://www.incompleteideas.net/IncIdeas/BitterLesson.html>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 5998–6008, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>.
- Giulia Vezzani, Dhruva Tirumala, Markus Wulfmeier, Dushyant Rao, Abbas Abdolmaleki, Ben Moran, Tuomas Haarnoja, Jan Humplik, Roland Hafner, Michael Neunert, Claudio Fantacci, Tim Hertweck, Thomas Lampe, Fereshteh Sadeghi, Nicolas Heess, and Martin A. Riedmiller. Skills: Adaptive skill sequencing for efficient temporally-extended exploration. *CoRR*, abs/2211.13743, 2022. doi: 10.48550/ARXIV.2211.13743. URL <https://doi.org/10.48550/arXiv.2211.13743>.
- Tianhe Yu, Ted Xiao, Austin Stone, Jonathan Tompson, Anthony Brohan, Su Wang, Jaspiar Singh, Clayton Tan, Dee M, Jodilyn Peralta, Brian Ichter, Karol Hausman, and Fei Xia. Scaling robot learning with semantically imagined experience, 2023.
- Yichi Zhang and Joyce Chai. Hierarchical task learning from language instructions with unified transformers and self-monitoring. *arXiv preprint arXiv:2106.03427*, 2021.
- Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *7th Annual Conference on Robot Learning*, 2023.