
Learning to Control Self-Assembling Morphologies: A Study of Generalization via Modularity

Deepak Pathak*
UC Berkeley

Chris Lu*
UC Berkeley

Trevor Darrell
UC Berkeley

Phillip Isola
MIT

Alexei A. Efros
UC Berkeley

Abstract

Contemporary sensorimotor learning approaches typically start with an existing complex agent (e.g., a robotic arm), which they learn to control. In contrast, this paper investigates a modular co-evolution strategy: a collection of primitive agents learns to dynamically self-assemble into composite bodies while also learning to coordinate their behavior to control these bodies. Each primitive agent consists of a limb with a motor attached at one end. Limbs may choose to link up to form collectives. When a limb initiates a link-up action and there is another limb nearby, the latter is magnetically connected to the ‘parent’ limb’s motor. This forms a new single agent, which may further link with other agents. In this way, complex morphologies can emerge, controlled by a policy whose architecture is in explicit correspondence with the morphology. We evaluate the performance of these *dynamic* and *modular* agents in simulated environments. We demonstrate better generalization to test-time changes both in the environment, as well as in the structure of the agent, compared to static and monolithic baselines. Project video and code are available at <https://pathak22.github.io/modular-assemblies/>.

1 Introduction

Possibly the single most pivotal event in the history of evolution was the point when single-celled organisms switched from always competing with each other for resources to sometimes cooperating, first by forming colonies, and later by merging into multicellular organisms [1]. These modular self-assemblies were successful because they combined the high adaptability of single-celled organisms while making it possible for vastly more complex behaviours to emerge. Indeed, one could argue that it is this modular design which allowed the multicellular organisms to successfully adapt, increase in complexity, and generalize to the constantly changing environment of prehistoric Earth. Like many researchers before us [14, 20, 23, 31, 32], we are inspired by the biology of multicellular evolution as a model for emergent complexity in artificial agents. Unlike most previous work however, we are primarily focused on modularity as a way of improving adaptability and generalization to *novel test-time scenarios*.

In this paper, we present a study of modular self-assemblies of primitive agents — “limbs” — which can link up to solve a task. Limbs have the option to bind together by a magnet that connects their morphologies within magnetic range (Figure 1), and when they do so, they pass messages and share rewards. Each limb comes with a simple neural net that controls the torque applied to its joints. Linking and unlinking are treated as dynamic actions so that the limb assembly can change shape within an episode. Similar setup has previously been explored in robotics as “self-reconfiguring modular robots” [21]. However, unlike prior work on such robots, where the control policies are hand-defined, we show how to *learn* the policies and study the generalization properties that emerge.

Our self-assembled agent can be represented as a graph of primitive limbs. Limbs pass messages to their neighbors in this graph in order to coordinate behavior. All limbs have a common policy network with shared parameters, i.e., a modular policy which takes the messages from adjacent limbs as input and outputs a torque to rotate the limb in addition to the linking/un-linking action. We call the

*Equal contribution.

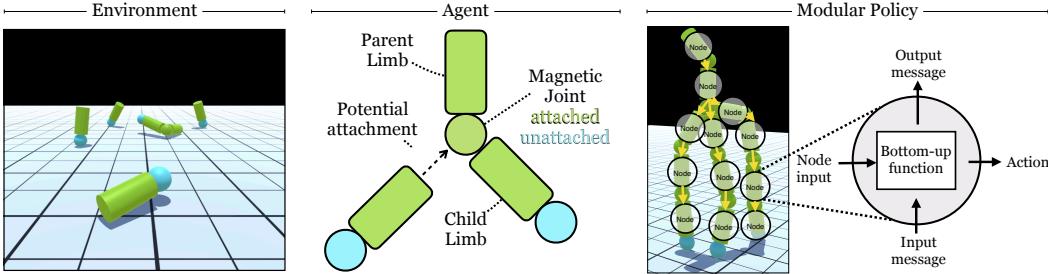


Figure 1: This work investigates the joint learning of control and morphology in self-assembling agents. Several primitive agents, containing a cylindrical body with a configurable motor, are dropped in a simulated environment (left). These primitive agents can self-assemble into collectives using magnetic joints (middle). Policy of the self-assembled agent is represented via proposed dynamic graph networks (DGN) with shared parameters (modular) across each limb (right).

aggregate neural network a “Dynamic Graph Network” (DGN) since it is a graph neural network [17] that can dynamically change topology as a function of its own outputs.

We test our dynamic limb assemblies on two separate tasks: standing up and locomotion. We are particularly interested in assessing how well can the assemblies generalize to novel testing conditions, not seen at training, compared to static and monolithic baselines. We evaluate test-time changes to both the environment (changing terrain geometry, environmental conditions), as well as the agent structure itself (changing the number of available limbs). We show that the dynamic self-assemblies are better able to generalize to these changes than the baselines. For example, we find that a single modular policy is able to control multiple possible morphologies, even those not seen during training, e.g., a 6-limb policy, trained to build a 6-limb tower, can be applied at test time on 3 or 12 limbs, and still able to perform the task.

The main contributions of this paper are:

- Train primitive agents that self-assemble into complex morphologies to jointly solve control tasks.
- Formulate morphological search as a reinforcement learning (RL) problem, where linking and unlinking are treated as actions.
- Represent policy via modular dynamic graph network (DGN) whose topology matches the agent’s physical structure.
- Demonstrate that self-assembling agents with dynamic morphology both train and generalize better than fixed-morphology baselines.

2 Environment and Agents

Investigating the co-evolution of control (i.e., *software*) and morphology (i.e., *hardware*) is not supported within standard benchmark environments typically used for sensorimotor control, requiring us to create our own. We opted for a minimalist design for our agents, the environment, and the reward structure, which is crucial to ensuring that the emergence of limb assemblies with complex morphologies is not forced, but happens naturally.

Environment Structure Our environment contains an arena where a collection of primitive agent limbs can self-assemble to perform control tasks. This arena is a ground surface equipped with gravity and friction. The arena can be procedurally changed to generate a variety of novel terrains by changing the height of each tile on the ground (see Figure 2). To evaluate the generalization properties of our agents, we generate a series of novel terrains. This includes generating bumpy terrain by randomizing the height of nearby tiles, stairs by incrementally increasing the height of each row of tiles, hurdles by changing the height of each row of tiles, gaps by removing alternating rows of tiles, etc. Some variations also include putting the arena ‘under water’ which basically amounts to increased drag (i.e. buoyancy). During training, we start our environment with a set of six primitive limbs on the ground which can assemble to form collectives to perform complex tasks.

Agent Structure All limbs share the same structure: a cylindrical body with a configurable motor on one end and the other end is free. The free-end of the limb can link up with the motor-end of the other limb, and then the motor acts as a joint between two limbs with three degrees of rotation. Hence, one can refer to the motor-end of the cylindrical limb as a *parent-end* and the free end as a *child-end*. Multiple limbs can attach their child-end to the parent-end of another limb, as shown in

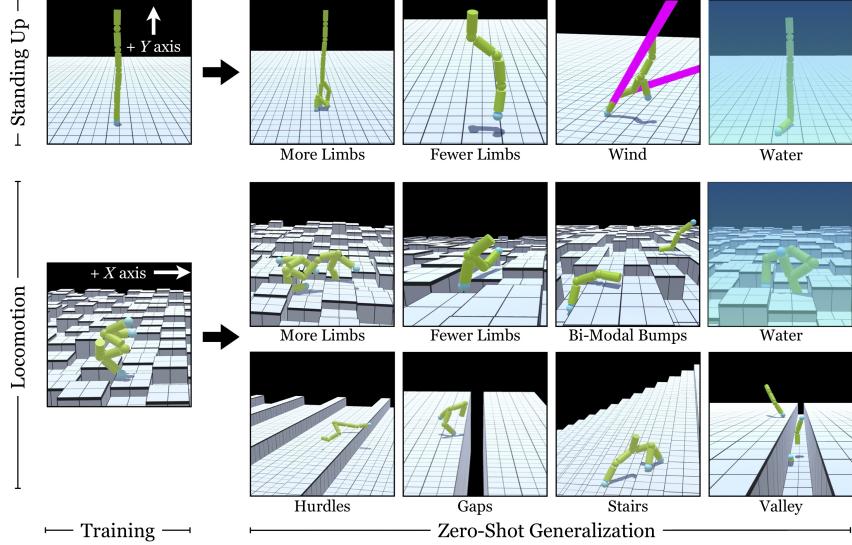


Figure 2: We illustrate our dynamic agents in two environments / tasks: standing up and locomotion. For each of these, we generate several new environment for evaluating generalization. Please refer to project video at <https://pathak22.github.io/modular-assemblies/> for better understanding of tasks.

Figure 1, to allow for complex graph morphologies to emerge. The limb of the parent-end controls the torques of joint. The unlinking action can be easily implemented by detaching two limbs, but the linking action has to deal with the ambiguity of which limb to connect to (if at all). To resolve this, we implement the linking action by attaching the closest limb within a small radius around the parent-node. The attachment mechanism is driven by a magnet inside the parent node which forces the closest child-limb within the magnetic range node to get docked onto itself if the parent signals to connect. If no other limb is present within the magnetic range, the linking action has no effect.

The primitive limbs are dropped in an environment to jointly solve a given control task. One key component of the self-assembling agent setup that makes it different from typical multi-agent scenarios [28] is that if some agents assemble to form a collective, the resulting morphology becomes a new *single agent* and all limbs within the morphology maximize a joint reward function. The output action space of each primitive agent contains the continuous torque values that are to be applied to the motor connected to the agent, and are denoted by $\{\tau_\alpha, \tau_\beta, \tau_\gamma\}$ for three degrees of rotation. In addition to the torque controls, each limb can decide to attach another link at its parent-end, or decide to unlink its child-end if already connected to other limb. The linking and unlinking decisions are binary. This complementary role assignment of child and parent ends, i.e., parent can only link and child can only unlink, makes it possible to decentralize the control across limbs.

Sensory Inputs In our self-assembling setup, each agent limb only has access to its local sensory information and does not know about other limbs. The sensory input of each agent includes its own dynamics, i.e., the location of the limb in 3-D euclidean coordinates, its velocity, angular rotation and angular velocity. Each end of the limb also has a trinary touch sensor to detect whether the end of the cylinder is touching 1) the floor, 2) another limb, or 3) nothing. Additionally, we also provide our limbs with a point depth sensor that captures the surface height on a 9×9 grid around the projection of center of limb on the surface.

One essential requirement to operationalize this setup is an efficient simulator to allow simultaneous simulation of several of these primitive limbs. We implement our environments in the Unity ML [12] framework, which is one of the dominant platforms for designing realistic games. For computational reasons, we do not allow the emergence of cycles in the self-assembling agents by not allowing the limbs to link up with already attached limbs within the same morphology. However, our setup is trivially extensible to general graphs.

3 Learning to Control Self-Assemblies

Consider a set of primitive limbs indexed by i in $\{1, 2, \dots, n\}$ dropped in an environment arena \mathcal{E} to perform a continuous control task. If needed, these limbs can assemble to form complex collectives

in order to improve their performance on the task. The task is represented by a reward function r_t and the goal of the limbs is to maximize the discounted sum of rewards over time t . If some limbs assemble into a collective, the resulting morphology effectively becomes a single agent with a combined policy to maximize the combined reward of the connected limbs. Further, the reward of an assembled morphology is a function of the whole morphology and not the individual agent limbs. For instance, in the task of learning to stand up, the reward is the height of the individual limbs if they are separate, but is the height of the whole morphology if those limbs have assembled into a collective.

3.1 Co-evolution: Linking/Unlinking as an Action

To learn a modular controller policy that could generalize to novel setups, our agents must learn the controller jointly as the morphology evolves over time. The limbs should simultaneously decide which torques to apply to their respective motors, while taking into account the connected morphology. Our hypothesis is that if a controller policy could learn in a modular fashion over iterations of increasingly sophisticated morphologies (see Figure 3), it could learn to be robust and generalizable to diverse situations. So, how can we optimize control and morphology under a common end-to-end framework?

We propose to treat the decision of linking and unlinking as additional actions of our primitive limbs. The total action space a_t at each iteration t can be denoted as $\{\tau_\alpha, \tau_\beta, \tau_\gamma, \sigma_{link}, \sigma_{unlink}\}$ where τ_* denote the raw *continuous* torque values to be applied at the motor and σ_* denote the *binary* actions whether to connect another limb at the parent-end or disconnect the child-end from the other already attached limb. This simple view of morphological evolution allows us to use ideas from RL [22].

3.2 Modularity: Self-Assembly as a Graph of Limbs

Integration of control and morphology in a common framework is only the first step. The key question is how to model this controller policy such that it is modular and reuses information across generations of morphologies. Let a_t^i be the action space and s_t^i be the local sensory input-space of the agent i . One naive approach to maximizing the reward is to simply combine the states of the limbs into the input-space, and output all the actions jointly using a single network. Formally, the policy is simply $\vec{a}_t = [a_t^0, a_t^1 \dots a_t^n] = \Pi(s_t^0, s_t^1 \dots, s_t^n)$. This interprets the self-assemblies as a single monolithic agent, ignoring the graphical structure. This is the current approach to solve many control problems, e.g., Mujoco humanoid [4] where the policy Π is trained to maximize the sum of rewards using RL.

In this work, we represent the agent’s policy via a graph neural network [17] in such a way that it explicitly corresponds to the morphology of the agent. Consider a collection of primitive limbs as graph G where each node is a limb i . Two limbs being physically connected by a joint is analogous to having an edge in the graph. As discussed in Section 2, each limb has two endpoints, a *parent-end* and a *child-end*. At a joint, the limb which connects via its parent-end acts as a parent-node in the corresponding edge, and the other limbs, which connect via their child-ends, are child-nodes. The parent-node (i.e., the agent with the parent-end) controls the torque of the edge (i.e., the joint motor).

3.3 Dynamic Graph Networks (DGN)

Each primitive limb node i has a policy controller of its own, which is represented by a neural network π_θ^i and receives a corresponding reward r_t^i for each time step t . We represent the policy of the self-assembled agent by the aggregated neural network that is connected in the same graphical manner as the physical morphology. The edge connectivity of the graph is represented in the overall graph policy by passing messages that flow from each limb to the other limbs physically connected to it via a joint. The parameters θ are shared across each primitive limbs allowing the overall policy of the graph to be modular with respect to each node. However, recall that the agent morphologies are dynamic, i.e., the connectivity of the limbs changes based on policy outputs. This changes the edge connectivity of the corresponding graph network at every timestep, depending on the actions chosen by each limb’s policy network in the previous timestep. Hence, we call this aggregate neural net a *Dynamic Graph Network (DGN)* since it is a graph neural network that can dynamically change topology as a function of its own outputs in the previous iteration.

DGN Optimization A typical rollout of our self-assembling agents during an episode of training contains a sequence of torques τ_t^i and the linking actions σ_t^i for each limb at each timestep t . The

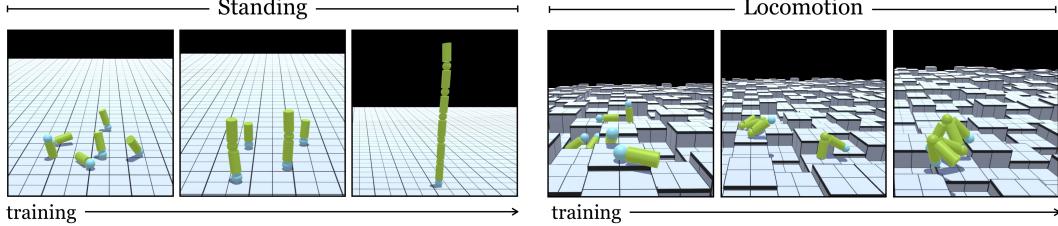


Figure 3: Co-evolution of Morphology w/ Control during Training: The gradual co-evolution of controller as well as the morphology of self-assembling agents over the course of training for the task of Standing Up (left) and Locomotion (right).

policy parameters θ are optimized to jointly maximize the reward for each limb:

$$\max_{\theta} \sum_{i=\{1,2,\dots,n\}} \mathbb{E}_{\vec{a}^i \sim \pi_{\theta}^i} [\Sigma_t r_t^i] \quad (1)$$

We optimize this objective via policy gradients, in particular, PPO [19]. DGN pseudo-code (as well as source code) and all training implementation details are in Section A.1, A.4 of the appendix.

DGN Connectivity The topology is captured in the DGN by passing messages through the edges between individual network nodes. Since the parameters of these limb networks are shared across each node, these messages can be seen as context information that may inform the policy of its role in the corresponding connected component of graph.

(a) *Message passing*: Messages are passed from leaf nodes to root, i.e., each agent gets information from its children. Instead of defining π_{θ}^i to be just as a function of state s_t^i , we pass each limb’s policy network information about its children nodes. We redefine π_{θ}^i as $\pi_{\theta}^i : [s_t^i, m_t^{C_i}] \rightarrow [a_t^i, m_t^i]$ where m_t^i is the output message of policy that goes into the parent limb and $m_t^{C_i}$ is the aggregated input messages from all the children nodes, i.e, $m_t^{C_i} = \sum_{c \in C_i} m_t^c$. If i has no children (i.e, root), a vector of zeros is passed in $m_t^{C_i}$. Messages are passed recursively until the root node. Alternative way is to start from root node and recursively pass until the messages reach the leaf nodes.

(b) *No message passing*: Note that for some environments or tasks, the context from the other nodes might not be a necessary requirement for effective control. In such scenarios, message passing might creates extra overhead for training a DGN. Importantly, even with no messages, DGN still allows for coordination between limbs. This is similar to a typical cooperative multi-agent setup [28] where each limb makes its own decisions in response to the previous actions of the other agents. However, our setup differs in that our agents may physically join up, rather than just coordinating behavior.

4 Experiments

We test the co-evolution of morphology and control across two primary tasks where self-assembling agents learn to: (a) stand up, and (b) perform locomotion. Limbs start each episode disconnected and located just above the ground plane at random locations, as shown in Figure 3. In the absence of an edge, input messages are set to 0 and the output ones are ignored. Action space is continuous raw torque values. Across all the tasks, the number of limbs at training is kept fixed to 6. We take the model from each time step and evaluate it on 50 episode runs to plot mean and std-deviation confidence interval in training curves. At test, we report the mean reward across 50 episodes of 1200 environment steps. The main focus of our investigation is to evaluate if the emerged modular controller generalizes to novel morphologies and environments. Video is on the project website.

Baselines We further compare how well these dynamic morphologies perform in comparison to a learned monolithic policy for both dynamic and fixed morphologies. In particular, we compare to a (a) *Monolithic Policy, Dynamic Graph*: Baseline where agents are still dynamic and can self-assemble, but their controller is represented by a single monolithic policy that takes as input the combined state of all agents and outputs actions for each of them. (b) *Monolithic Policy, Fixed Graph*: Similar single monolithic policy as previous baseline, but the morphology is hand-designed constructed from the limbs and kept fixed and static during training and test. This is analogous to a standard robotics “vanilla RL” setup in which a morphology is predefined and then a policy is learned to control it. We

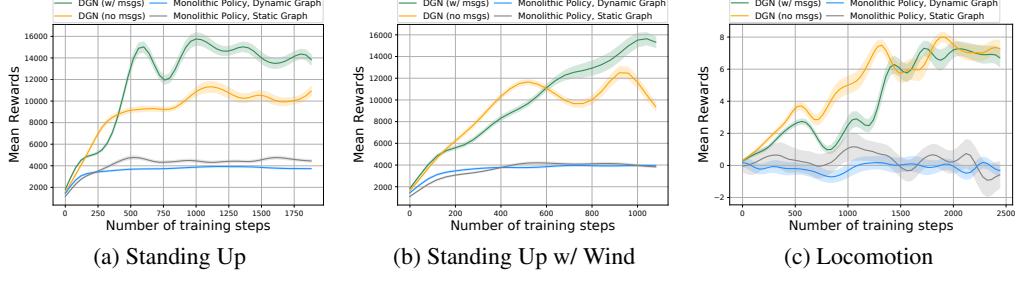


Figure 4: Training self-assembling agents: We show the performance of different methods for joint training of control and morphology for three tasks: learning to stand up (left), standing up in the presence of wind (center) and locomotion in bumpy terrain (right). These policies generalize to novel scenarios as shown in the tables.

chose the fixed morphology to be a straight chain of 6-limbs in all the experiments. This linear-chain may be optimal for standing as tall as possible, but it is not necessarily optimal for *learning* to stand; the same would hold for locomotion. However, we confirmed that the both standing and locomotion task are solvable with linear-chain morphology (shown in Figure 3 and video on project website).

Although monolithic policy is more expressive (complete state information of all limbs), it is also harder to train as we increase the number of limbs, because the observation and action spaces increase in dimensionality. Indeed, this is what we find in the Appendix Figure 5: the monolithic policy can perform well on up to three limbs, but does not reach the optimum on four to six limbs. In contrast, the DGN limb policy (shared between all limbs) has a fixed size observation and action space, independent of the number of limbs under control.

4.1 Learning to Self-Assemble

We first validate if it is possible to train self-assembling agent policy end-to-end via Dynamic Graph Networks. Below, we discuss our environments and compare the training efficiency of each method.

Standing Up Task In this task, each agent’s objective is to maximize the height of the highest point in its morphology. Limbs have an incentive to self-assemble because the potential reward would scale with the number of limbs if the self-assembled agent can control them. The training process begins with six-limbs falling on the ground randomly, as shown in Figure 3. These limbs act independently in the beginning but gradually learn to self-assemble as training proceeds. Figure 4a compares the training efficiency and performance of different methods during training. We found that our DGN policy variants perform significantly better than the monolithic policies for the standing up task.

Standing Up in the Wind Task Same as the previous task, except with the addition of ‘wind’, which we operationalize as random forces applied to random points of each limb at random times, see Figure 2(Wind). Figure 4b shows the superior performance of DGN compared to the baselines.

Locomotion Task The reward function for locomotion is defined as the distance covered by the agent along X -axis. The training is performed on a bumpy terrain shown in Figure 2. The training performance in Figure 4c shows that DGN variants outperform the monolithic baselines.

As shown in Figure 4, training our DGN algorithm with message passing either seems to perform better or similar to the one without message passing. In particular, message passing is significantly helpful where long-term reasoning is needed across limbs, for instance, messages help in standing up task because there is only one morphological structure to do well (i.e., linear tower). In locomotion, it is possible to do well with a large variety of morphologies, and thus both DGN variants reach similar performance. Now onwards, we show results using DGN w/ msgs as our primary approach.

4.2 Zero-Shot Generalization to Number of Limbs

We investigate if our trained policy generalizes to changes in the number of limbs. We pick the best model from training and evaluate it without any finetuning at test-time, i.e., zero-shot generalization.

Standing Up Task We train the policy with 6 limbs and test with 12 and 4 limbs. As shown in Table 1, despite changes in number of limbs, DGN is able to retain similar performance w/o any finetuning. The co-evolution of morphology jointly with the controller allows the modular policy to experience increasingly complex morphological structures. We hypothesize that this morphological curriculum at training makes the agent more robust at test-time.

Environment	DGN (ours)	Monolithic Policy (dynamic)	Monolithic Policy (fixed)
<i>Standing Up Task</i>			
<i>Training Environment</i>			
Standing Up	17518	4104	5351
<i>Zero-Shot Generalization</i>			
More (2x) Limbs	19796 (113%)	n/a	n/a
Fewer (.5x) Limbs	10839 (62%)	n/a	n/a
<i>Standing Up in the Wind Task</i>			
<i>Training Environment</i>			
Stand-Up in Wind	18423	4176	4500
<i>Zero-Shot Generalization</i>			
2x Limbs + (S)Wind	15351 (83%)	n/a	n/a
<i>Locomotion Task</i>			
<i>Training Environment</i>			
Locomotion	8.71	0.96	2.96
<i>Zero-Shot Generalization</i>			
More (2x) Limbs	5.47 (63%)	n/a	n/a
Fewer (.5x) Limbs	6.64 (76%)	n/a	n/a

Table 1: Zero-Shot Generalization to Number of Limbs: Quantitative evaluation of the generalizability of the learned policies. For each method, we first pick the best performing model from the training and then evaluate it on each of the novel scenarios without further finetuning, i.e., in a zero-shot manner. We report the score attained by the self-assembling agent along with the percentage of training performance retained upon transfer in parenthesis. Higher is better.

Environment	DGN (ours)	Monolithic Policy (dynamic)	Monolithic Policy (fixed)
<i>Standing Up Task</i>			
<i>Training Environment</i>			
Standing Up	17518	4104	5351
<i>Zero-Shot Generalization</i>			
Water + 2x Limbs	16871 (96%)	n/a	n/a
Winds	16803 (96%)	3923 (96%)	4531 (85%)
Strong Winds	15853 (90%)	3937 (96%)	4961 (93%)
<i>Standing Up in the Wind Task</i>			
<i>Training Environment</i>			
Stand-Up in Wind	18423	4176	4500
<i>Zero-Shot Generalization</i>			
(S)trong Wind	17384 (94%)	4010 (96%)	4507 (100%)
Water+2x+SWd	17068 (93%)	n/a	n/a
<i>Locomotion Task</i>			
<i>Training Environment</i>			
Locomotion	8.71	0.96	2.96
<i>Zero-Shot Generalization</i>			
Water + 2x Limbs	6.57 (75%)	n/a	n/a
Hurdles	6.39 (73%)	-0.77 (-79%)	-3.12 (-104%)
Gaps in Terrain	3.25 (37%)	-0.32 (-33%)	2.09 (71%)
Bi-modal Bumps	6.62 (76%)	-0.56 (-57%)	-0.44 (-14%)
Stairs	6.6 (76%)	-8.8 (-912%)	-3.65 (-122%)
Inside Valley	5.29 (61%)	0.47 (48%)	-1.35 (-45%)

Table 2: Zero-Shot Generalization to Novel Environments: The best performing model from the training is evaluated on each of the novel scenarios without any further finetuning. The score attained by the self-assembling agent is reported along with the percentage of training performance retained upon transfer in parenthesis. Higher value is better.

Note that we can not generalize *Monolithic policy* baselines to scenarios with more or fewer limbs because they can't accommodate different action and state space dimensions from training; it has to be retrained. Hence, we made a comparison to DGN by retraining baseline on Standing task: DGN is trained on 6 limbs and tested on 4 limbs w/o any finetuning, while baseline is trained both times. DGN achieves 17518 (6limbs - train), 10839 (4limbs - test) scores, while baseline achieves 5351 (6limbs - train), 7356 (4limbs - train). Even without any training on 4 limbs, DGN outperforms baseline because it is difficult to train monolithic policy with large action space (Figure 1 in Appendix).

Standing Up in the Wind Task Similarly, we evaluate the agent policy trained for standing up task in winds with 6 limbs to 12 limbs. Table 1 shows that the DGN performs significantly better than monolithic policy at train, and able to retain most of its performance even with twice the limbs.

Locomotion Task We also evaluate the generalization of locomotion policies trained with 6 limbs to 12 and 4 limbs. As shown in Table 1, DGN not only achieves good performance at training but is also able to retain most of its performance.

4.3 Zero-Shot Generalization to Novel Environments

We now evaluate the performance of our modular agents in novel terrains by creating several different scenarios by varying environment conditions (described in Section 2) to test zero-shot generalization.

Standing Up Task We test our trained policy without any further finetuning in environments with increased drag (i.e., ‘under water’), and adding varying strength of random push-n-pulls (i.e. , ‘wind’). Table 2 shows that DGN seems to generalize better than monolithic policies. We believe that this generalization is result of both the learning being modular as well as the fact that limbs learned to assemble in physical conditions (e.g. forces like gravity) with gradually growing morphologies. Such forces with changing morphology are similar to setup with varying forces acting on fixed morphology resulting in robustness to external interventions like winds.

Standing Up in the Wind Task Similarly, the policies trained with winds are able to generalize to scenarios with either stronger winds or winds inside water.

Locomotion Task We generate several novel scenarios for evaluating locomotion: with water, a terrain with hurdles of a certain height, a terrain with gaps between platforms, a bumpy terrain with a bi-modal distribution of bump heights, stairs, and an environment with a valley surrounded by walls on both sides (see Figure 2). These variations are generated procedurally. The modular policies learned by DGN tend to generalize better than the monolithic agent policies as shown in Table 2.

This generalization could be explained by the incrementally increasing complexity of self-assembling agents at training. For instance, the training begins with all limbs separate which gradually form group of two, three and so on, until the training converges. Since the policy is *modular with shared parameters across limbs*, the training of smaller size assemblies with small bumps would in turn prepare the large assemblies for performing locomotion through higher hurdles, stairs etc at test. Furthermore, the training terrain has a finite length which makes the self-assemblies launch themselves forward as far as possible upon reaching the boundary to maximize the distance along X-axis. This behavior helps the limbs generalize to environments like gaps or valley where they end up on the next terrain upon jumping and continue to perform locomotion.

5 Related Work

Morphogenesis & self-reconfiguring modular robots The idea of modular and self-assembling agents goes back at least to Von Neumann’s *Theory of Self-Reproducing Automata* [24]. In robotics, such systems have been termed “self-reconfiguring modular robots” [14, 21]. There has been a lot of work in modular robotics to design real hardware robotic modules that can be docked together to form complex robotic morphologies [6, 9, 15, 29, 31]. Alternatives to optimize agent morphologies include genetic algorithms that search over a generative grammar [20] and energy-based minimization to directly optimizing controllers [7, 25]. Conditioning on several hardware designs has also been shown to increase robustness [5]. We approach morphogenesis from a learning perspective, in particular deep RL, and study the resulting generalization properties. We achieve morphological co-evolution via *dynamic actions* (linking), which agents take during their lifetimes, whereas the past approaches treat morphology as an optimization target to be updated between generations or episodes. Since the physical morphology also defines the connectivity of the policy net, our proposed algorithm can also be viewed as performing a kind of neural architecture search [33] in physical agents.

Graph neural networks Encoding graphical structures into neural networks [17] has been used for a large number of applications, including question answering [2], quantum chemistry [8], semi-supervised classification [13], and representation learning [30]. The works most similar to ours involve learning controllers [16, 27]. For example, Nervenet [27] represents individual limbs and joints as nodes in a graph and demonstrates multi-limb generalization. However, the morphologies on which Nervenet operates are not learned jointly with the policy and hand-defined to be compositional in nature. Others [3, 11] have shown that graph neural networks can also be applied to inference models as well as to planning. Prior graph neural network based approaches deal with static graph which is defined by auxiliary information, e.g. language parser [2]. In contrast, we propose dynamic graph networks where the graph policy changes itself dynamically over the training.

Concurrent Work Ha [10], Schaff et al. [18] use RL to improve limb design given fixed morphology. Wang et al. [26] gradually evolves the environment to improve robustness of an agent. However, both the work assume the topology of agent morphology to stay the same during train and test.

6 Discussion

Modeling intelligent agents as modular, self-assembling morphologies has long been a very appealing idea. The efforts to create practical systems to evolve artificial agents goes back at least two decades to the beautiful work of Karl Sims [20]. In this paper, we are revisiting these ideas using the contemporary machinery of deep networks and reinforcement learning. Examining the problem in the context of machine learning, rather than optimization, we are particularly interested in modularity as a key to generalization, in terms of improving adaptability and robustness to novel environmental conditions. Poor generalization is the Achilles heel of modern robotics research, and the hope is that this could be a promising direction in addressing this key issue. We demonstrated a number of promising experimental results, suggesting that modularity does indeed improve generalization in simulated agents. While these are just the initial steps, we believe that the proposed research direction is promising and its exploration will be fruitful to the research community. To encourage follow-up work, we will publicly release all code, models, and environments.

Acknowledgments

We would like to thank Igor Mordatch, Chris Atkeson, Abhinav Gupta and the members of BAIR for fruitful discussions and comments. This work was supported in part by Berkeley DeepDrive, and the Valrhona reinforcement learning fellowship. DP is supported by Facebook graduate fellowship.

References

- [1] B. Alberts, D. Bray, J. Lewis, M. Raff, K. Roberts, and J. D. Watson. *Molecular Biology of the Cell*. Garland Publishing, New York, 1994. 1
- [2] J. Andreas, M. Rohrbach, T. Darrell, and D. Klein. Neural module networks. In *CVPR*, 2016. 8
- [3] P. W. Battaglia, J. B. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. Zambaldi, M. Malinowski, A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018. 8
- [4] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai gym. *arXiv:1606.01540*, 2016. 4
- [5] T. Chen, A. Murali, and A. Gupta. Hardware conditioned policies for multi-robot transfer learning. In *NIPS*, 2018. 8
- [6] J. Daudelin, G. Jing, T. Tosun, M. Yim, H. Kress-Gazit, and M. Campbell. An integrated system for perception-driven autonomy with modular robots. *Science Robotics*, 2018. 8
- [7] M. De Lasa, I. Mordatch, and A. Hertzmann. Feature-based locomotion controllers. In *ACM Transactions on Graphics (TOG)*, 2010. 8
- [8] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing for quantum chemistry. *arXiv preprint arXiv:1704.01212*, 2017. 8
- [9] K. Gilpin, K. Kotay, D. Rus, and I. Vasilescu. Miche: Modular shape formation by self-disassembly. *IJRR*, 2008. 8
- [10] D. Ha. Reinforcement learning for improving agent design. *arXiv preprint arXiv:1810.03779*, 2018. 8
- [11] D.-A. Huang, S. Nair, D. Xu, Y. Zhu, A. Garg, L. Fei-Fei, S. Savarese, and J. C. Niebles. Neural task graphs: Generalizing to unseen tasks from a single video demonstration. *arXiv preprint arXiv:1807.03480*, 2018. 8
- [12] A. Juliani, V.-P. Berges, E. Vckay, Y. Gao, H. Henry, M. Mattar, and D. Lange. Unity: A general platform for intelligent agents. *arXiv preprint arXiv:1809.02627*, 2018. 3
- [13] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016. 8
- [14] S. Murata and H. Kurokawa. Self-reconfigurable robots. *IEEE Robotics & Automation Magazine*, 2007. 1, 8
- [15] J. W. Romanishin, K. Gilpin, and D. Rus. M-blocks: Momentum-driven, magnetic modular robots. In *IROS*, 2013. 8
- [16] A. Sanchez-Gonzalez, N. Heess, J. T. Springenberg, J. Merel, M. Riedmiller, R. Hadsell, and P. Battaglia. Graph networks as learnable physics engines for inference and control. *arXiv preprint arXiv:1806.01242*, 2018. 8
- [17] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. The graph neural network model. *IEEE Transactions on Neural Network*, 2009. 2, 4, 8
- [18] C. Schaff, D. Yunis, A. Chakrabarti, and M. R. Walter. Jointly learning to construct and control agents using deep reinforcement learning. *arXiv preprint arXiv:1801.01432*, 2018. 8

- [19] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv:1707.06347*, 2017. 5, 11
- [20] K. Sims. Evolving virtual creatures. In *Computer graphics and interactive techniques*, 1994. 1, 8
- [21] K. Stoy, D. Brandt, D. J. Christensen, and D. Brandt. *Self-reconfigurable robots: an introduction*. Mit Press Cambridge, 2010. 1, 8
- [22] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press Cambridge, 1998. 4
- [23] X. Tu and D. Terzopoulos. Artificial fishes: Physics, locomotion, perception, behavior. In *Proceedings of the 21st annual conference on Computer graphics and interactive techniques*, 1994. 1
- [24] J. Von Neumann, A. W. Burks, et al. Theory of self-reproducing automata. *IEEE Transactions on Neural Networks*, 1966. 8
- [25] K. Wampler and Z. Popović. Optimal gait and form for animal locomotion. In *ACM Transactions on Graphics (TOG)*, 2009. 8
- [26] R. Wang, J. Lehman, J. Clune, and K. O. Stanley. Paired open-ended trailblazer (poet): Endlessly generating increasingly complex and diverse learning environments and their solutions. *arXiv preprint arXiv:1901.01753*, 2019. 8
- [27] T. Wang, R. Liao, J. Ba, and S. Fidler. Nervenet: Learning structured policy with graph neural networks. *ICLR*, 2018. 8
- [28] M. Wooldridge. *An introduction to multiagent systems*. John Wiley & Sons, 2009. 3, 5
- [29] C. Wright, A. Johnson, A. Peck, Z. McCord, A. Naaktgeboren, P. Gianfortoni, M. Gonzalez-Rivero, R. Hatton, and H. Choset. Design of a modular snake robot. In *IROS*, 2007. 8
- [30] Z. Yang, B. Dhingra, K. He, W. W. Cohen, R. Salakhutdinov, Y. LeCun, et al. Glomo: Unsupervisedly learned relational graphs as transferable representations. *arXiv preprint arXiv:1806.05662*, 2018. 8
- [31] M. Yim, D. G. Duff, and K. D. Roufas. Polybot: a modular reconfigurable robot. In *ICRA*, 2000. 1, 8
- [32] M. Yim, W.-M. Shen, B. Salemi, D. Rus, M. Moll, H. Lipson, E. Klavins, and G. S. Chirikjian. Modular self-reconfigurable robot systems. *IEEE Robotics & Automation Magazine*, 2007. 1
- [33] B. Zoph and Q. V. Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016. 8

A Appendix

In this section, we provide additional details about the experimental setup and extra experiments.

A.1 Implementation and Training details

We use PPO [19] as the underlying reinforcement learning method to optimize Equation (1) in Section 3.3 (main paper). Each limb policy is represented by 4-layered fully-connected neural network with ReLU non-linearities and trained with a learning rate of $3e - 4$, discount factor of 0.995, entropy coefficient of 0.01, advantage parameter of 0.95 and batch size of 2048. The messages are 32 length float vectors. The optimizer used to optimize PPO is RMS-Prop. We will move these details to the main paper. Parameters are shared across network modules and they all predict action at the same time. Each episode is 5000 steps long at training. Across all the tasks, the number of limbs at training is kept fixed to 6. Limbs start each episode disconnected and located just above the ground plane at random locations, as shown in Figure 3 in the main paper. In the absence of an edge, input messages are set to 0 and the output ones are ignored. Action space is continuous raw torque values. We take the model from each time step and evaluate it on 50 episodes to plot mean and standard-deviation (confidence intervals) in training curves. At test, we report the mean reward across 50 episode runs of 1200 environment steps.

A.2 Project Video

The project video is at <https://youtu.be/ngCIB-IWD8E>. Please watch it for better understanding of the results and turn on the volume for narration.

A.3 Performance of Fixed-Graph Baseline vs. Number of Limbs

To verify whether the training of *Monolithic Policy w/ Fixed Graph* is working, we ran it on standing up and locomotion tasks across varying number of limbs. We show in Figure 5 that the baseline performs well with less number of limbs which suggests that the reason for failure in 6-limbs case is indeed the morphology graph being fixed, and not the implementation of this baseline.

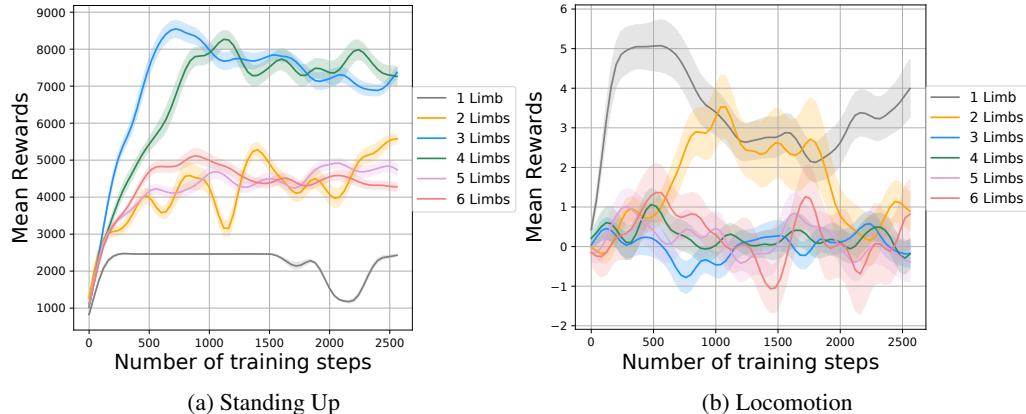


Figure 5: The performance of *Monolithic Policy w/ Fixed Graph* baseline as the number of limbs varies in the two tasks: standing up (left) and locomotion (right). This shows that the monolithic baseline works well with less (1-3 limbs), but fails with 6 limbs during training.

A.4 Pseudo Code of the DGN Algorithm

Notation is summarized in Algorithm 1, and the full pseudo-code is summarized in the following algorithm boxes: Algorithm 2, and Algorithm 3. Full source code available at <https://pathak22.github.io/modular-assemblies/>.

Algorithm 1: Notation Summary (defined in Section 3.3)

```
1 foreach node  $i$  do
2   |  $a_t^i, m_t^i = \pi_\theta^i(s_t^i, m_t^{C_i})$ 
3 end
4 where
5  $s_t^i$ : observation state of agent limb  $i$ 
6  $a_t^i$ : action output of agent limb  $i$ : 3 torques, attach, detach
7  $m_t^{C_i}$ : aggregated message from children nodes input to agent  $i$  (bottom-up-1)
8  $m_t^i$ : output message that agent  $i$  passes to its parent (bottom-up-2)
9  $\theta$ :  $\theta_1, \theta_2$ 
10 messages are 32 length floating point vectors.
```

Algorithm 2: Pseudo-code: DGN w/ Message Passing

```

1 Initialize parameters  $\theta_1, \theta_2$  randomly.
2 Initialize all message vectors  $m_t^{C_i}, m_t^i$  to be zero
3 Represent graph connectivity  $G$  as a list of edges
4 Note: In the beginning, all edges are zeros, i.e., non-existent
5 foreach timestep  $t$  do
6   Each limb agent  $i$  observes its own state vector  $s_t^i$ 
7   foreach agent  $i$  do
8     # Compute incoming child messages
9      $m_t^{C_i} = 0$ 
10    foreach child node  $c$  of agent  $i$  do
11      |  $m_t^{C_i} += m_t^c$ 
12    end
13    # Compute action and message to parent  $p$  of agent  $i$  in  $G$ 
14     $a_t^i, m_t^i := \pi_\theta(s_t^i, m_t^{C_i})$ 
15    # Execute morphology change as per  $a_t^i$ 
16    if  $a_t^i[3] == \text{attach}$  then
17      | find closest agent  $j$  within distance  $d$  from agent  $i$ , otherwise  $j=\text{NULL}$ 
18      | add edge  $(i, j)$  in  $G$ 
19      | also make physical joint between  $(i, j)$ 
20    end
21    if  $a_t^i[4] == \text{detach}$  then
22      | delete edge  $(i, \text{parent of } i)$  in  $G$ 
23      | also delete physical joint between  $(i, j)$ 
24    end
25    # Execute torques from  $a_t^i$ 
26    Apply torques  $a_t^i[0], a_t^i[1], a_t^i[2]$ 
27  end
28  # Update graph and agent morphology
29  Find all connected components in  $G$ 
30  foreach connected component  $c$  do
31    foreach agent  $i \in c$  do
32      | reward  $r_t^i = \text{reward of } c$  (e.g. max height)
33    end
34  end
35 end
36 Update  $\theta$  to maximize discounted reward using PPO as follows:
37 let  $\vec{a}_t = [a_t^1, a_t^2 \dots a_t^n]$ 
38  $\vec{s}_t = [s_t^1, s_t^2 \dots s_t^n]$ 
39  $\hat{A}_t = \text{advantage of discounted rewards, } r_t = \sum_{\text{agent } i} r_t^i$ 
40 PPO:  $\max_\theta \mathbb{E}[\hat{A}_t \frac{\pi_\theta(\vec{a}_t | \vec{s}_t)}{\pi_{\theta_{\text{old}}}(\vec{a}_t | \vec{s}_t)} - \beta \text{KL}(\pi_{\theta_{\text{old}}}(. | \vec{s}_t) \pi_\theta(. | \vec{s}_t))]$ 
41 Repeat until training converges

```

Algorithm 3: Pseudo-code: No-message DGN

- 1 Same as Algorithm-2 but hard-code incoming child and parent messages to be always 0, i.e., $m_t^{C_i} = 0$ and $m_t^{p_i} = 0$ in each iteration.
-