



UNIVERSITY OF
TORONTO



VECTOR
INSTITUTE

INSTITUT
VECTEUR

NerveNet: Towards Universal Deep RL Controllers for Robotics

Tingwu Wang

University Of Toronto, Vector Institute

Contents

1. Introduction
 - a. Reinforcement Learning with Neural Network
 - b. Graph Neural Networks
2. Nervenet
 - a. Basic Idea
 - b. Model
 - i. Input Model
 - ii. Propagation Model
 - iii. Output Model
3. Experiments
 - a. Transfer Learning
 - i. Zeroshot-Performance
 - ii. Fine-tuning Policy
 - b. Robustness
 - c. Multi-task Learning

Reinforcement Learning with Neural Network

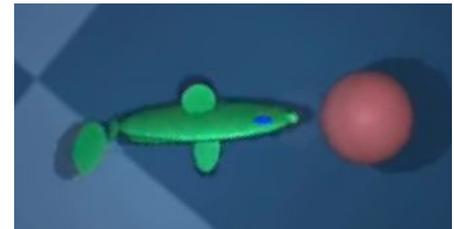
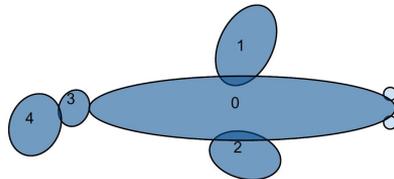
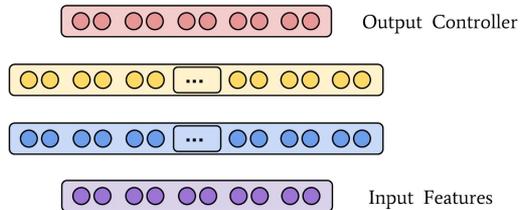
1. One of the most exciting area recently
 - a. AlphaGo[1]
 - b. Complex humanoid control [2, 3] and industry robot[4, 5]
 - c. Video games [6, 7]



2. Problems do exist [8], for example:
 - a. Transferability
Policy learnt on one agent could not be transferred to agent with different dynamics.
 - b. Robustness
Perturbing the parameters of robot results in catastrophic performance drop.

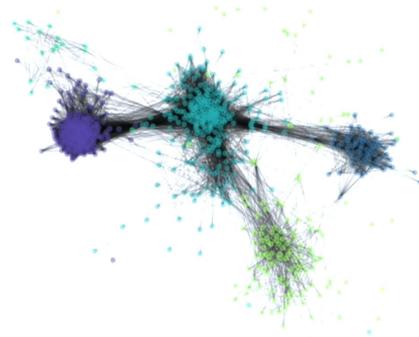
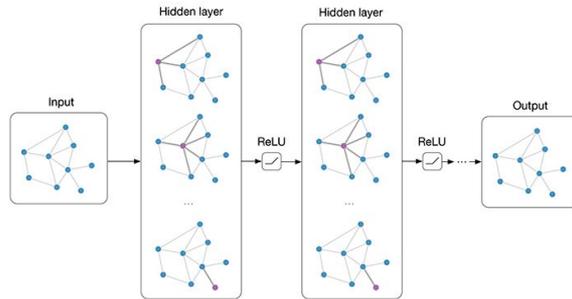
Reinforcement Learning with Neural Network

3. We constrain the focus to be locomotion control of simulated robot [9]
 - a. Input: observations from the sensors of the robot.
 - i. Joint angles & velocity, body's Cartesian positions, feedback forces
 - ii. ...
 - b. Output: the torques applied respectively on each joint.
4. A typical neural network used in these tasks:
 - a. Multi-Layer Perceptron Network (MLP)
 - b. Ignoring structural information and performing black-box optimization.
 - i. PPO [3], TRPO [10], A3C [11], DPG[12], DDPG [13]



Graph Neural Networks

1. Many types of the input data is a vector or tensor, for example
 - a. Observations of the robot (vector, MLP)
 - b. Images (tensor, CNN)
2. Many important real-world data could be formulated as graphs:
 - a. Social networks, knowledge graphs, the World Wide Web, etc. [14].
 - b. A more natural idea is to use graph neural networks to process the data
 - c. Example: How likely are you a Warrior supporter or Cavalier supporter?

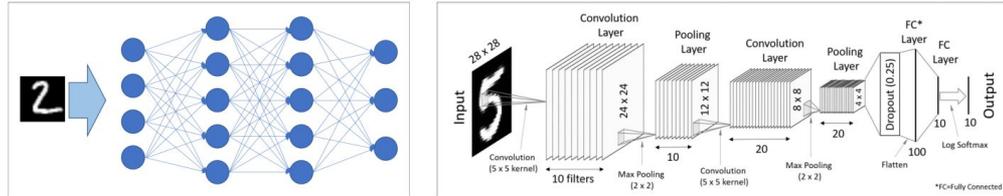


Contents

1. Introduction
 - a. Reinforcement Learning with Neural Network
 - b. Graph Neural Networks
2. Nervenet
 - a. Basic Idea
 - b. Model
 - i. Input Model
 - ii. Propagation Model
 - iii. Output Model
3. Experiments
 - a. Transfer Learning
 - i. Zeroshot-Performance
 - ii. Fine-tuning Policy
 - b. Robustness
 - c. Multi-task Learning

Basic Idea

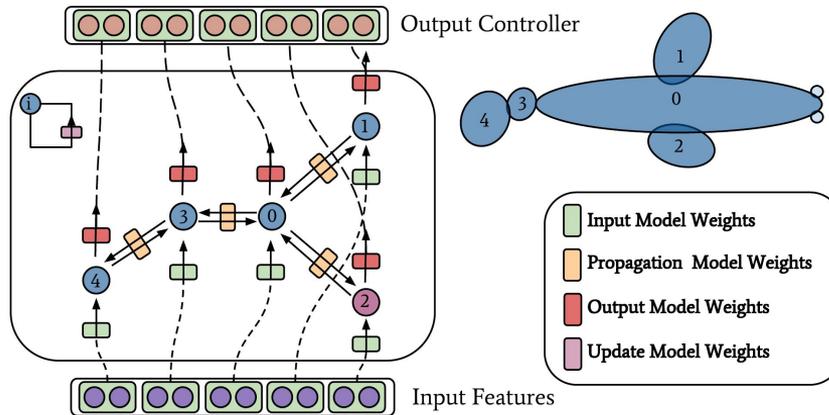
1. Use Graph Neural Networks instead of MLPs
 - a. MLP Potentially causing the problem on transferability and robustness.
 - i. A perhaps similar story: Image classification with MLP or CNN?
 - ii. Empirically MLP doesn't scale with data and is easy to overfit
 - iii. CNN applies strong structure prior and uses shared weights



- b. Make use of structural information with GNN
- c. Intuition: to determine the torques applied to one joint
 - i. neighbour's absolute position, joint's position & velocity, inertia info, ...

Input Model

1. For each node (0, 1, 2, 3, 4) of the agent, we fetch the corresponding observation from the input features.
2. Map it into the initial node states using a function F : $h_u^0 = F_{\text{in}}(x_u)$
 - a. a matrix, identity function, or MLP.



Propagation Model

1. Let each node communicate with the neighbours by propagating messages.
 - a. Message is calculated by passing the current states through a function M.

$$m_{(u,v)}^t = M_{c(u,v)}(h_u^t)$$

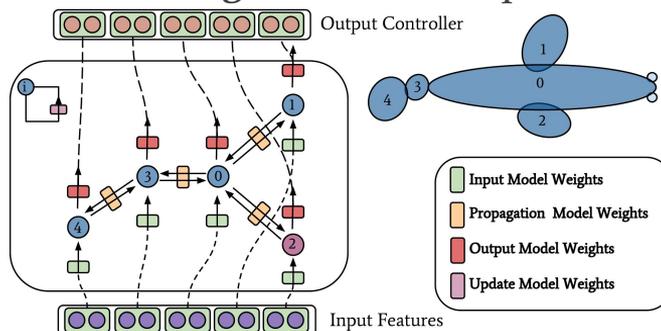
- b. Node aggregate the messages and update its current node states.

$$\bar{m}_u^t = A(\{h_v^t | v \in \mathcal{N}_{in}(u)\})$$

- c. Ideally we should propagate the messages until the node states converge

- i. In practice we propagate the message for 4 timesteps.

$$h_u^{t+1} = U_{p_u}(h_u^t, \bar{m}_u^t)$$



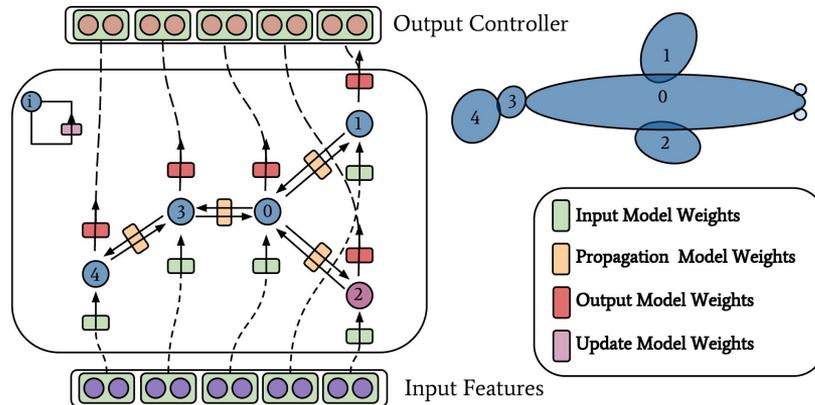
Output Model

1. For each node with a controller, we output the control signal.

a. The control signal is a function of the node last states.

$$\mu_{u \in \mathcal{O}} = O(h_u^T)$$

b. Control signal could be the mean and variance of the gaussian policy.

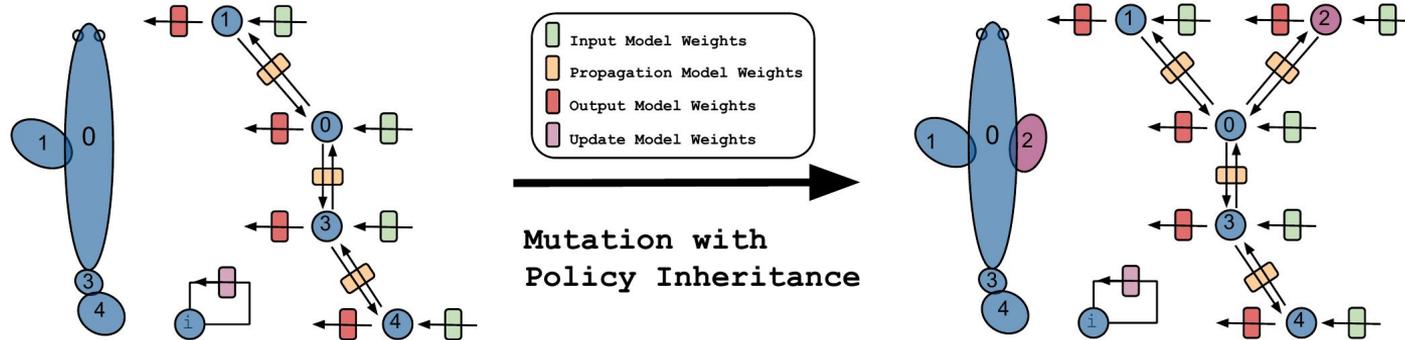


Contents

1. Introduction
 - a. Reinforcement Learning with Neural Network
 - b. Graph Neural Networks
2. Nervenet
 - a. Basic Idea
 - b. Model
 - i. Input Model
 - ii. Propagation Model
 - iii. Output Model
3. Experiments
 - a. Transfer Learning
 - i. Zeroshot-Performance
 - ii. Fine-tuning Policy
 - b. Robustness
 - c. Multi-task Learning

Transfer Learning

1. As we mentioned before, in old methods, policy learnt on one agent could not be transferred to agent with different dynamics. (full result in [15])
 - a. Transfer learning benchmarks: snake and centipedes.



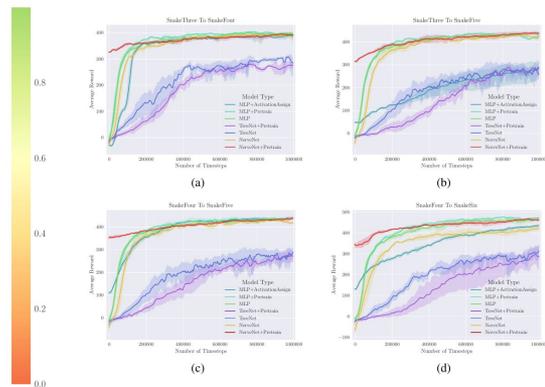
Transfer Learning

- As we mentioned before, in old methods, policy learnt on one agent could not be transferred to agent with different dynamics. (full result in [15])

a. Transfer learning benchmarks: snake and centipedes.

Tasks	1. Random	2. MLPAA	3. MLPP Models	4. TreeNet	5. NerveNet
4to66	-31.6 (32%)	109.4 (84%)	-126.7 (-2%)	-16.5 (38%)	139.6 (96%)
4to68	-45.8 (41%)	18.2 (76%)	-190.9 (-39%)	16.2 (72%)	44.3 (91%)
4to10	-44.3 (42%)	11.4 (73%)	-195.6 (-42%)	-101.5 (10%)	39.8 (88%)
4to12	-48.7 (39%)	9.9 (72%)	-215.8 (-53%)	17.6 (76%)	38.6 (88%)
4to14	-52.0 (37%)	8.0 (71%)	-233.0 (-62%)	-79.6 (22%)	39.0 (88%)
4toCp66	-17.0 (26%)	-5.1 (29%)	-113.9 (1%)	249.5 (94%)	47.6 (42%)
4toCp68	-25.5 (52%)	5.1 (69%)	-132.2 (-6%)	33.3 (85%)	40.0 (88%)
4toCp10	-28.2 (51%)	-12.8 (50%)	-133.7 (-7%)	28.2 (82%)	40.2 (88%)
6to68	-45.8 (4%)	21.1 (7%)	-191.4 (-3%)	320.8 (24%)	1674.9 (99%)
6to10	-44.3 (7%)	-42.4 (7%)	-193.2 (-6%)	44.7 (15%)	940.5 (98%)
6to12	-49.1 (13%)	-14.4 (20%)	-227.0 (-20%)	19.5 (27%)	367.7 (95%)
6to14	-52.0 (17%)	-10.0 (28%)	-221.3 (-25%)	126.3 (63%)	247.8 (94%)
6to20	-72.4 (14%)	10.0 (39%)	-351.4 (-70%)	-233.6 (-34%)	198.5 (96%)
6to30	-67.1 (22%)	-28.5 (38%)	-263.3 (-50%)	17.6 (57%)	114.9 (97%)
6to40	-72.7 (19%)	4.3 (51%)	-330.2 (-83%)	21.1 (58%)	97.8 (90%)
6toCp68	-25.4 (14%)	36.5 (23%)	-141.9 (-3%)	398.9 (78%)	523.6 (97%)
6toCp10	-28.2 (14%)	12.8 (21%)	-155.2 (-5%)	277.8 (63%)	504.0 (99%)
6toCp12	-32.0 (22%)	12.1 (33%)	-177.5 (-14%)	-114.9 (1%)	255.9 (95%)
6toCp14	-41.0 (21%)	11.8 (36%)	-187.5 (-18%)	-67.7 (14%)	224.8 (96%)

Tasks	1. Random	2. MLPAA	3. MLPP Models	4. TreeNet	5. NerveNet
4to66	-75.0 (6%)	545.3 (92%)	-73.5 (6%)	-47.7 (10%)	577.3 (96%)
4to68	-91.0 (10%)	62.0 (67%)	-80.8 (14%)	-84.8 (13%)	146.9 (98%)
4to10	-76.5 (16%)	21.0 (52%)	-80.9 (11%)	-70.0 (18%)	128.4 (92%)
4to12	-81.7 (14%)	13.1 (49%)	-76.9 (15%)	-63.2 (21%)	126.3 (91%)
4to14	-89.0 (11%)	0.0 (44%)	-86.4 (12%)	-73.2 (17%)	125.7 (90%)
4toCp66	-77.3 (17%)	-22.5 (40%)	-79.9 (16%)	-72.2 (19%)	91.1 (87%)
4toCp68	-82.9 (17%)	-26.9 (44%)	-91.8 (13%)	-66.9 (25%)	80.1 (95%)
4toCp10	-82.9 (17%)	-36.6 (39%)	-88.8 (14%)	-83.7 (17%)	86.9 (98%)
6to68	-91.0 (0%)	87.8 (1%)	-88.6 (0%)	8.5 (1%)	10612.6 (99%)
6to10	-76.5 (0%)	-17.0 (1%)	-81.8 (0%)	-77.4 (0%)	6343.6 (99%)
6to12	-81.2 (1%)	-1.4 (4%)	-79.3 (1%)	-91.5 (1%)	2532.2 (99%)
6to14	-89.0 (1%)	-3.9 (6%)	-83.2 (1%)	-51.5 (3%)	1749.7 (98%)
6to20	-95.2 (1%)	-4.9 (7%)	-105.4 (0%)	-102.3 (1%)	1447.4 (98%)
6to30	-111.2 (0%)	-8.2 (7%)	-127.3 (0%)	-76.7 (4%)	867.5 (99%)
6to40	63.9 (16%)	140.5 (23%)	36.9 (13%)	101.8 (19%)	971.5 (98%)
6toCp68	-83.0 (1%)	138.3 (7%)	-96.7 (0%)	-63.9 (1%)	3117.3 (99%)
6toCp10	-82.9 (1%)	13.6 (3%)	-86.4 (0%)	-72.3 (1%)	3230.3 (99%)
6toCp12	-87.5 (1%)	7.3 (7%)	-91.4 (1%)	-90.8 (1%)	1675.5 (99%)
6toCp14	-96.5 (1%)	2.9 (7%)	-92.5 (1%)	-93.9 (1%)	1517.6 (99%)



(a) Zero-shot average reward.

(b) Zero-shot average running-length.

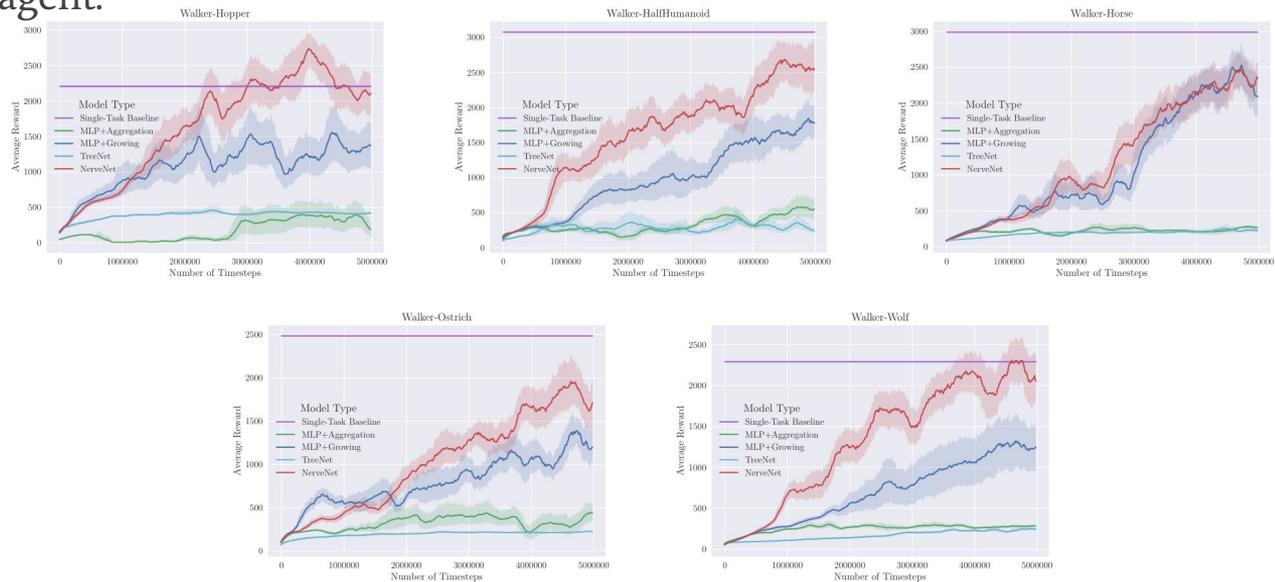
Robustness

1. We also test the robustness of the policy by perturbing the parameters of the

Model		HalfHumanoid	Hopper	Ostrich	Wolf	Horse	Average
MLP	Reward	1775.75	1369.59	1198.88	1249.23	2084.07	/
	Ratio	57.7%	62.0%	48.2%	54.5%	69.7%	58.6%
TreeNet	Reward	237.81	417.27	224.07	247.03	223.34	/
	Ratio	79.3%	98.0%	57.4%	141.2%	99.2%	94.8%
NerveNet	Reward	2536.52	2113.56	1714.63	2054.54	2343.62	/
	Ratio	96.3%	101.8%	98.8%	105.9%	106.4%	101.8%

Multi-task Learning

1. We also test the robustness of the policy by perturbing the parameters of the agent.



Contents

1. Introduction
 - a. Reinforcement Learning with Neural Network
 - b. Graph Neural Networks
2. Nervenet
 - a. Basic Idea
 - b. Model
 - i. Input Model
 - ii. Propagation Model
 - iii. Output Model
3. Experiments
 - a. Transfer Learning
 - i. Zeroshot-Performance
 - ii. Fine-tuning Policy
 - b. Robustness
 - c. Multi-task Learning

Reference

- [1] Silver, David, et al. "Mastering the game of Go with deep neural networks and tree search." *nature* 529.7587 (2016): 484-489.
- [2] Heess, Nicolas, et al. "Emergence of locomotion behaviours in rich environments." arXiv preprint arXiv:1707.02286 (2017).
- [3] Schulman, John, et al. "Proximal policy optimization algorithms." arXiv preprint arXiv:1707.06347 (2017).
- [4] Andrychowicz, Marcin, et al. "Hindsight experience replay." *Advances in Neural Information Processing Systems*. 2017.
- [5] Gu, Shixiang, et al. "Q-prop: Sample-efficient policy gradient with an off-policy critic." arXiv preprint arXiv:1611.02247 (2016).
- [6] Mnih, Volodymyr, et al. "Playing atari with deep reinforcement learning." arXiv preprint arXiv:1312.5602 (2013).
- [7] Fernandez, Jose Maria Font, and Tobias Mahlmann. "The Dota 2 Bot Competition." *IEEE Transactions on Games* (2018).
- [8] Rajeswaran, Aravind, et al. "Towards generalization and simplicity in continuous control." *Advances in Neural Information Processing Systems*. 2017.
- [9] Todorov, Emanuel, Tom Erez, and Yuval Tassa. "Mujoco: A physics engine for model-based control." *Intelligent Robots and [10] Systems (IROS), 2012 IEEE/RSJ International Conference on*. IEEE, 2012.
- [10] Schulman, John, et al. "Trust region policy optimization." *International Conference on Machine Learning*. 2015.
- [11] Mnih, Volodymyr, et al. "Asynchronous methods for deep reinforcement learning." *International Conference on Machine Learning*. 2016.
- [12] Silver, David, et al. "Deterministic policy gradient algorithms." *ICML*. 2014.
- [13] Lillicrap, Timothy P., et al. "Continuous control with deep reinforcement learning." arXiv preprint arXiv:1509.02971 (2015).
- [14] Kipf, Thomas N., and Max Welling. "Semi-supervised classification with graph convolutional networks." arXiv preprint arXiv:1609.02907 (2016).
- [15] Wang, Tingwu, et al. *Nervenet: Learning structured policy with graph neural networks*. Diss. University of Toronto, 2017.