

# Hierarchical reinforcement learning

Shvechikov Pavel

National Research University Higher School of Economics,  
Yandex School of Data Analysis

May 22, 2017

# Overview

## 1 Motivation

## 2 FeUdal Networks

- Introduction
- Training
- Results

# Overview

Humans and animals are naturally hierarchical :

- ① **pack** model: alpha, beta and omega dogs
- ② **herd** models: single stallion

Every one in the hierarchy knows what and when to do.

- ① idea of **duties segregation** is promising
- ② especially when duties are explicitly assigned : )

But how to learn this hierarchy without knowing anything about the task beforehand?

# Standard approaches to goal setting

End-to-end learning of hierarchical RL architectures

- ① has many advantages
- ② results either in sub-goals of length **one**
- ③ or in sub-goals of length equal to **episode length**

To combat this issue several approaches exists

- ① explicitly predefine each sub-goal (Tessler et al., 2016)
- ② introduce regularizers to (Bacon et al., 2016)

An alternative (Vezhnevets et al., 2017)

- Do **not allow gradients to pass** between Manager and Worker !

# FuNs: FeUdal Networks for HRL (Vezhnevets et al., 2017)

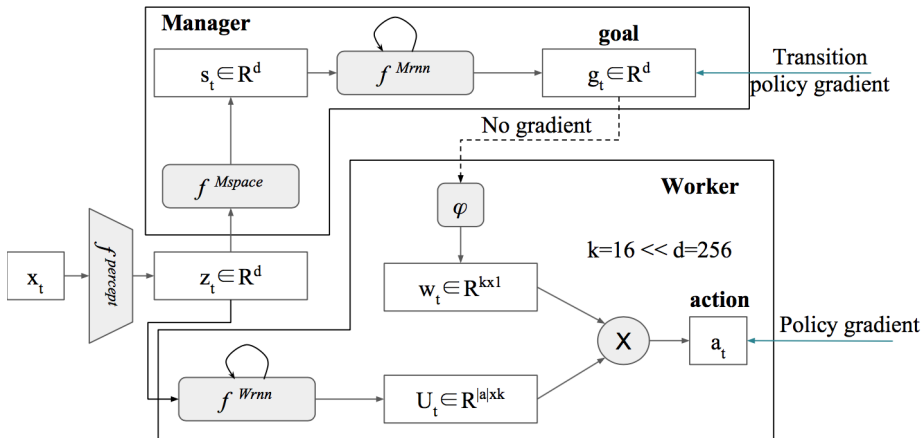
Beneficial when

- ① sparse rewards
- ② very long timescale credit assignment
- ③ non-Markovian environments that require memory

FuNs at a glance:

- ① goal setting is **decoupled** from goal achievement
- ② end-to-end **differentiable** model
- ③ novel **transition policy gradient** for Manager learning
- ④ **directional goals**, not absolute values

## FuN: architecture



# FuN: neural network dynamics

## Common

$\mathbf{z}_t = f^{percept}(\mathbf{x}_t)$  is a preprocessed observation  $\mathbf{x}_t$

## Manager

$$\mathbf{s}_t = f^{Mspace}(\mathbf{z}_t)$$

$$(h_t^M, \hat{\mathbf{g}}_t) = f^{Mrnn}(\mathbf{s}_t, h_{t-1}^M)$$

$$\mathbf{g}_t = \hat{\mathbf{g}}_t / \|\hat{\mathbf{g}}_t\|$$

$\mathbf{s}_t$  – latent state representation

$\mathbf{g}_t$  – goal vector

$h_t^M$  – Manager rnn hidden state

## Worker

$$\mathbf{w}_t = \varphi \left( \sum_{i=t-c}^t \mathbf{g}_i \right)$$

$$h_t^W, \mathbf{U}_t = f^{Wrnn}(\mathbf{z}_t, h_{t-1}^W)$$

$$\pi_t = \text{Softmax}(\mathbf{U}_t \mathbf{w}_t)$$

$\mathbf{w}_t$  – goal embedding

$\varphi(\cdot)$  – linear, **without bias**

$\mathbf{U}_t$  – matrix, output of rnn

$h_t^W$  – Worker rnn hidden state

# FuN: Training Manager

Goal of agent as a whole is  $R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k=1}$

$$\nabla g_t = (R_t - V_t^M(x_t, \theta)) \cdot d_{\cos}(s_{t+c} - s_t, g_t(\theta)),$$

- $(s_{t+c} - s_t)$  – how much have we changed the direction
- $g_t$  – direction we told Worker to go to
- $d_{\cos}(\alpha, \beta) = \alpha^\top \beta / (|\alpha||\beta|)$  – cosine similarity

**Note:** we pretend  $s_t$  does not depend on  $\theta$



# FuN: update rule motivation

Assume

- ➊ Manager maintains high-level policy  $\mu(s_t, \theta)$
- ➋ let  $o_t = \mu(s_t, \theta)$  be a sub-policy chosen by high-level policy
- ➌ sub-policy is a fixed duration behavior (of  $c$  steps)
- ➍ sub-policy induces distribution over states  $p(s_{t+c} | s_t, o_t)$
- ➎ high-level policy induces  $\pi^{TP}(s_{t+c} | s_t) = p(s_{t+c} | s_t, \mu(s_t, \theta))$
- ➏  $\nabla_{\theta} \pi_t^{TP} = \mathbb{E} [(R_t - V(s_t)) \nabla_{\theta} \log p(s_{t+c} | s_t, \mu(s_t, \theta))]$
- ➐ BUT! if

$$p(s_{t+c} | s_t, o_t) \propto e^{d_{\cos}(s_{t+c} - s_t, g_t)}$$

we recover Manager updates!

# FuN: update rule motivation

Assume

- ① Manager maintains high-level policy  $\mu(s_t, \theta)$
- ② let  $o_t = \mu(s_t, \theta)$  be a sub-policy chosen by high-level policy
- ③ sub-policy is a fixed duration behavior (of  $c$  steps)
- ④ sub-policy induces distribution over states  $p(s_{t+c} | s_t, o_t)$
- ⑤ high-level policy induces  $\pi^{TP}(s_{t+c} | s_t) = p(s_{t+c} | s_t, \mu(s_t, \theta))$
- ⑥  $\nabla_{\theta} \pi_t^{TP} = \mathbb{E} [(R_t - V(s_t)) \nabla_{\theta} \log p(s_{t+c} | s_t, \mu(s_t, \theta))]$
- ⑦ BUT! if

$$p(s_{t+c} | s_t, o_t) \propto e^{d_{\cos}(s_{t+c} - s_t, g_t)}$$

we recover Manager updates!

## Equivalent formulation

$s_{t+c} - s_t$  should follow **von Mises-Fisher distribution**  
with mean vector  $g_t$

# FuN: dilated LSTM

**Note:**  $f^{Mrnn}$  is not simple RNN !

Allows operation at lower temporal resolution.

① introduce more *cores* constituting  $h = \{\hat{h}^i\}_{i=1}^r$

② update

$$\hat{h}_t^{t\%r}, g_t = \text{LSTM}(s_t, \hat{h}_{t-1}^{t\%r}; \theta^{\text{LSTM}}) \quad (1)$$

③  $\hat{h}_t^{t\%r}$  is pooled  $c$  steps backwards

# FuN: Training Worker

Intrinsic reward for Worker

$$R_t^I = \frac{1}{c} \sum_{i=1}^c d_{\cos}(s_t - s_{t-i}, g_{t-i})$$

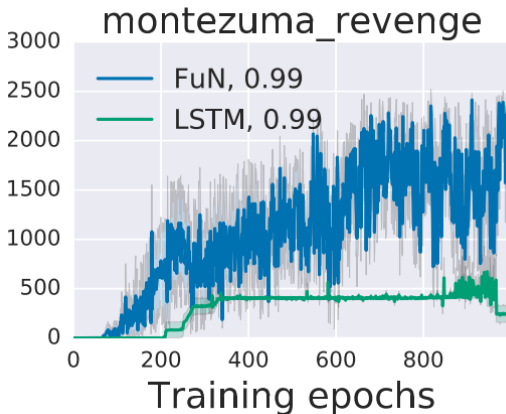
Advantage Actor Critic for Worker training

$$\nabla \pi_t = (R_t + \alpha R_t^I - V_t^D(x_t, \theta)) \cdot \nabla_{\theta} \log \pi(a_t | x_t; \theta)$$

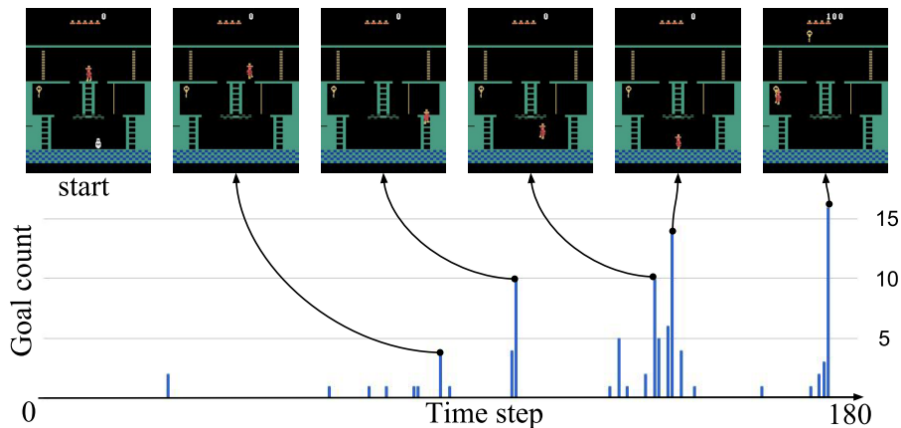
- $\alpha$  – hyper parameter regulating influence of intrinsic reward

**Note:** discount rate could be different for Manager and Worker

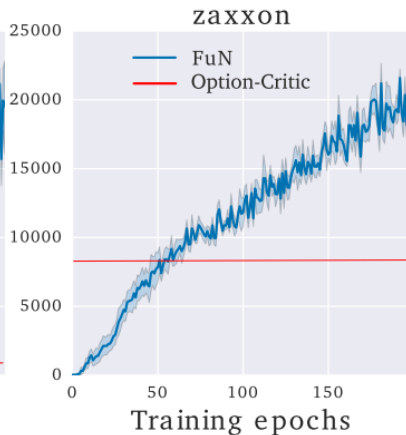
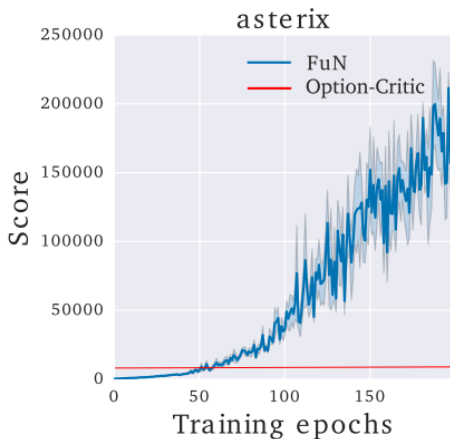
# FuN: results



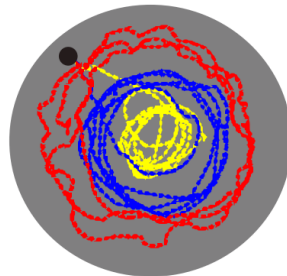
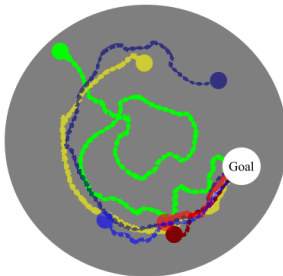
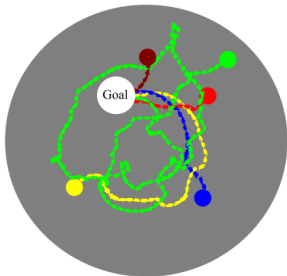
# FuN: results



# FuN: results






# FuN: results





Thank you!

# References I

-  Bacon, Pierre-Luc et al. (2016). “The Option-Critic Architecture”.  
In: *CoRR* abs/1609.05140. URL:  
<http://arxiv.org/abs/1609.05140>.
-  Tessler, Chen et al. (2016). “A deep hierarchical approach to  
lifelong learning in minecraft”. In: *arXiv preprint*  
*arXiv:1604.07255*.
-  Vezhnevets, Alexander Sasha et al. (2017). “FeUdal Networks for  
Hierarchical Reinforcement Learning”. In: *arXiv: 1703.01161*.  
URL: <http://arxiv.org/abs/1703.01161>.