Reinforcement learning project
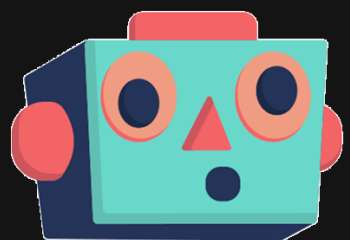
# TSP problem
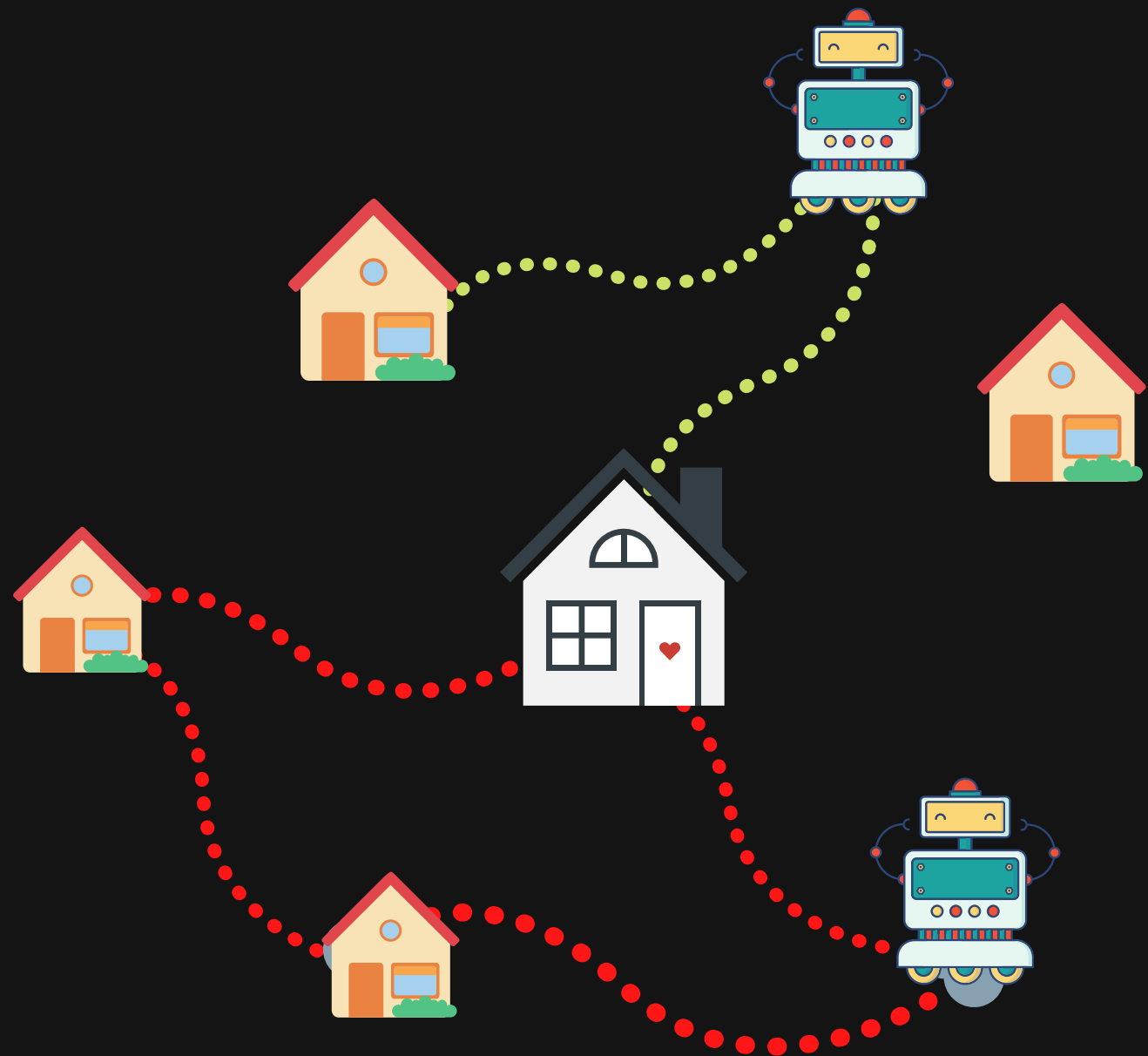
by Nina Konovalova
and Albina Klepach

*TSP -- Travelling Salesman Problem

# Plan

- Introduction
- Description of the problem
- Description of the methods
- Results
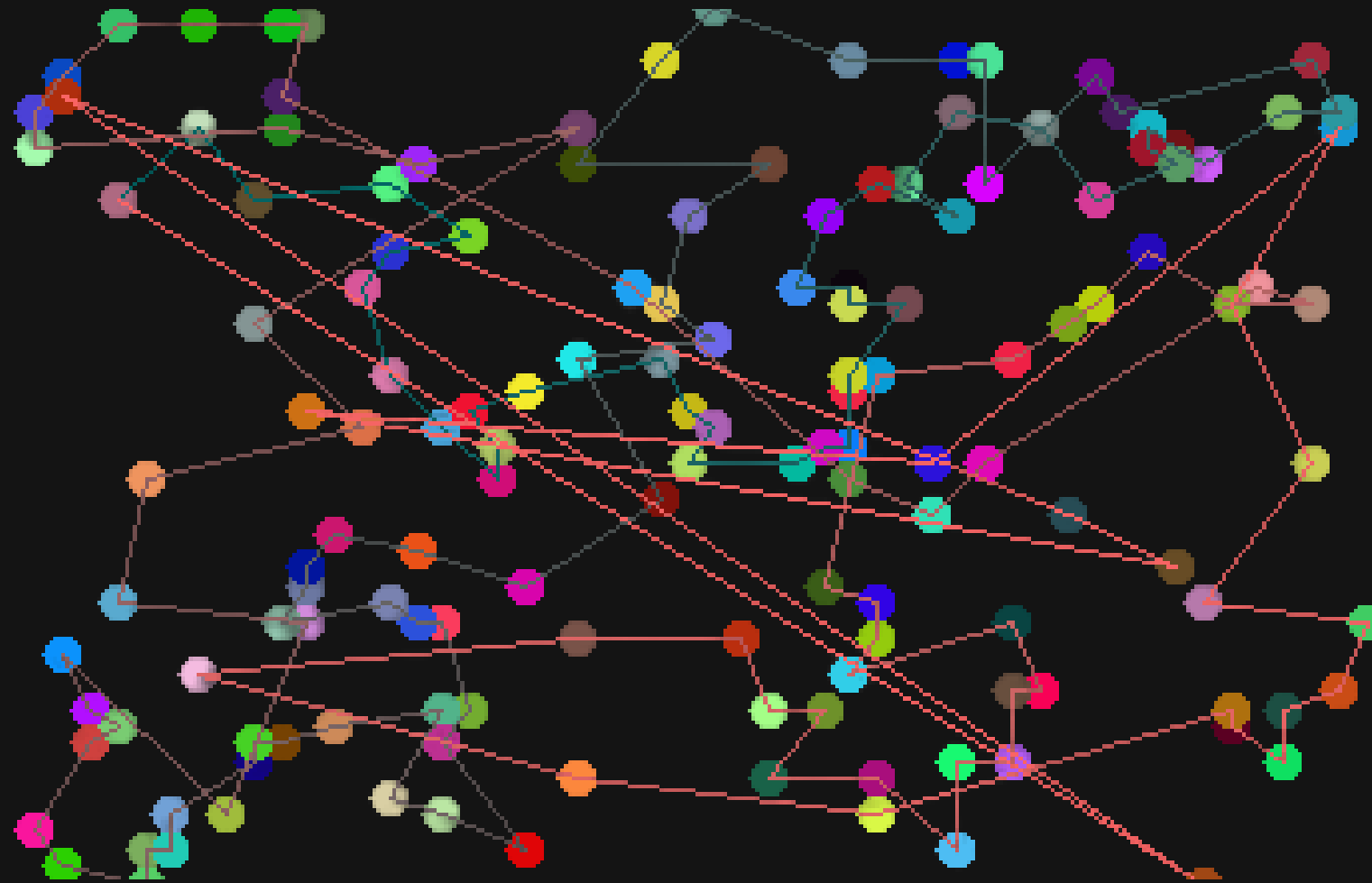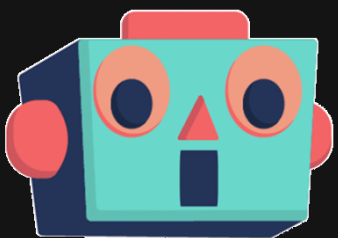- Conclusion

# Introduction

Combinatorial optimization

VPR

TSP

OP

# Description of the methods

Markov Decision Process (MDP)
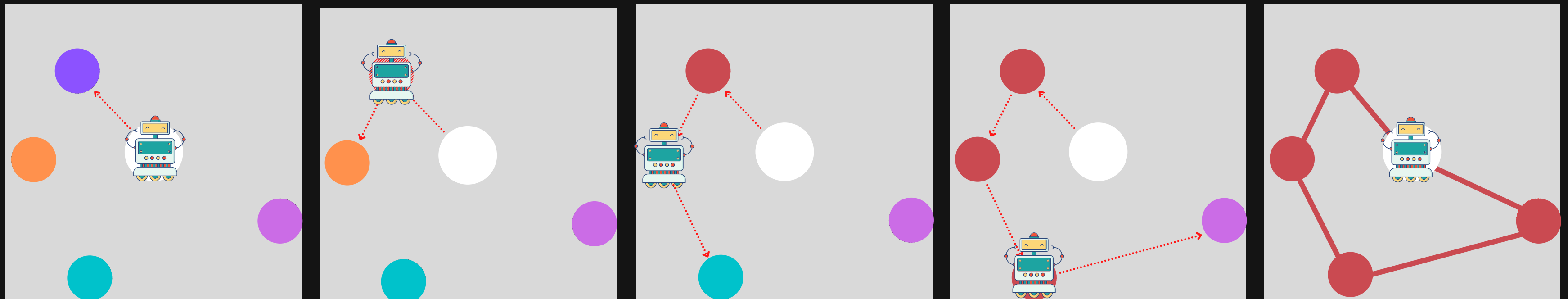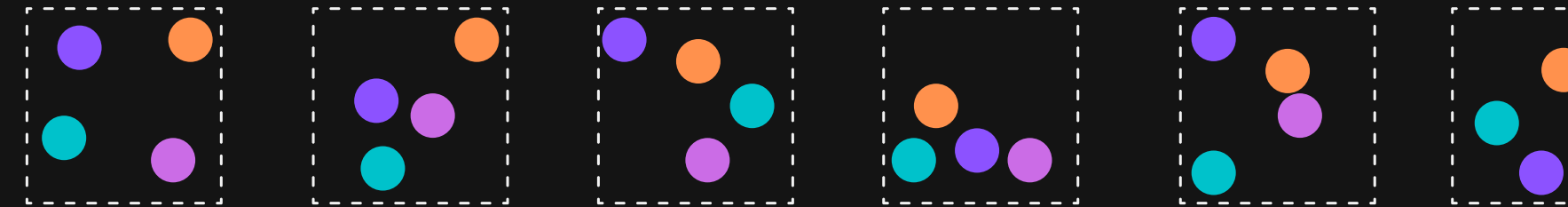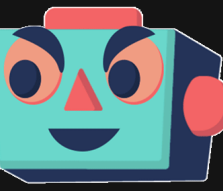
naive approach requires many samples

terminating condition when all
destinations are visited

# Description of the methods

Sequence of $n$ cities in 2D space $s = \{\mathbf{x}_i\}_{i=1}^n$, where each $\mathbf{x}_i \in \mathrm{R}^2$
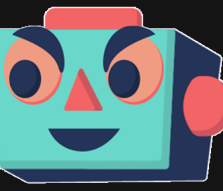
No depot:

$$L(\pi|s) = \|\mathbf{x}_{pi(n)} - \mathbf{x}_{pi(1)}\|_2 + \sum_{i=1}^{n-1} \|\mathbf{x}_{pi(i)} - \mathbf{x}_{pi(i+1)}\|_2$$

Depot:

$$L(\pi|s) = \|\mathbf{x}_{pi(n)} - \mathbf{x}_{depot}\|_2 + \|\mathbf{x}_{pi(1)} - \mathbf{x}_{depot}\|_2 + \sum_{i=1}^{n-1} \|\mathbf{x}_{pi(i)} - \mathbf{x}_{pi(i+1)}\|_2$$

Chain rule (to factorize the probability of a tour):

$$p(\pi|s) = \prod_{n=1}^{n} p(\pi(i)|\pi(<i), s)$$

# Description of the methods

Out training objective is expected tour length:

$$J(\Theta|s) = \mathbb{E}_{\pi \sim p_\Theta(\cdot|s)} L(\pi|s)$$

The gradient of expected tour length (with REINFORCE algorithm):

$$\nabla_\Theta J(\Theta|s) = \mathbb{E}_{\pi \sim p_\Theta(\cdot|s)} \left[ (L(\pi|s) - b(s)) \nabla_\Theta \log p_\Theta(\pi|s) \right]$$

$$\nabla_\Theta J(\Theta) \approx \frac{1}{B} \sum_{i=1}^{B} (L(\pi|s_i) - b(s_i)) \nabla_\Theta \log p_\Theta(\pi_i|s_i)$$
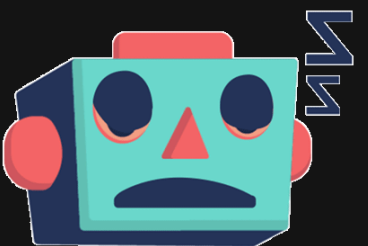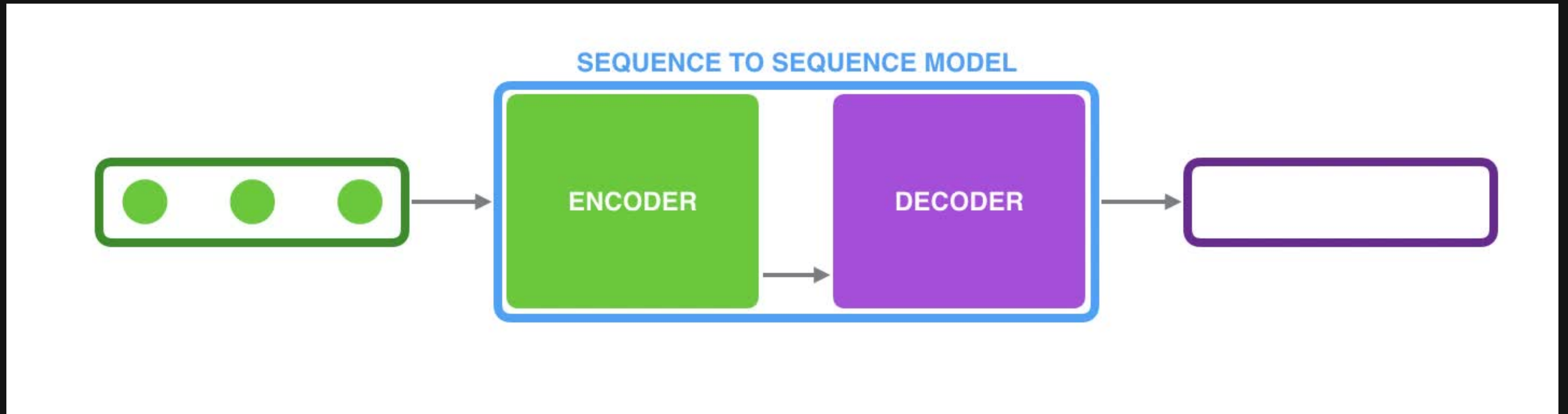
The objective for critic ($b_{\Theta_v}(s)$ is a prediction):

$$\mathbb{L}(\Theta_v) = \frac{1}{B} \sum_{i=1}^{B} \|b_{\Theta_v}(s_i) - L(\pi_i)\|_2^2$$

---

**Algorithm 1** Actor-critic training

1: **procedure** TRAIN(training set $S$, number of training steps $T$, batch size $B$)
2:     Initialize pointer network params $\theta$
3:     Initialize critic network params $\theta_v$
4:     **for** $t = 1$ to $T$ **do**
5:         $s_i \sim$ SAMPLEINPUT($S$) for $i \in \{1, \ldots, B\}$
6:         $\pi_i \sim$ SAMPLESOLUTION($p_\theta(.|s_i)$) for $i \in \{1, \ldots, B\}$
7:         $b_i \leftarrow b_{\theta_v}(s_i)$ for $i \in \{1, \ldots, B\}$
8:         $g_\theta \leftarrow \frac{1}{B} \sum_{i=1}^{B} (L(\pi_i|s_i) - b_i) \nabla_\theta \log p_\theta(\pi_i|s_i)$
9:         $\mathcal{L}_v \leftarrow \frac{1}{B} \sum_{i=1}^{B} \|b_i - L(\pi_i)\|_2^2$
10:        $\theta \leftarrow$ ADAM($\theta, g_\theta$)
11:        $\theta_v \leftarrow$ ADAM($\theta_v, \nabla_{\theta_v} \mathcal{L}_v$)
12:    **end for**
13:    **return** $\theta$
14: **end procedure**

# Seq2Seq



SEQUENCE TO SEQUENCE MODEL
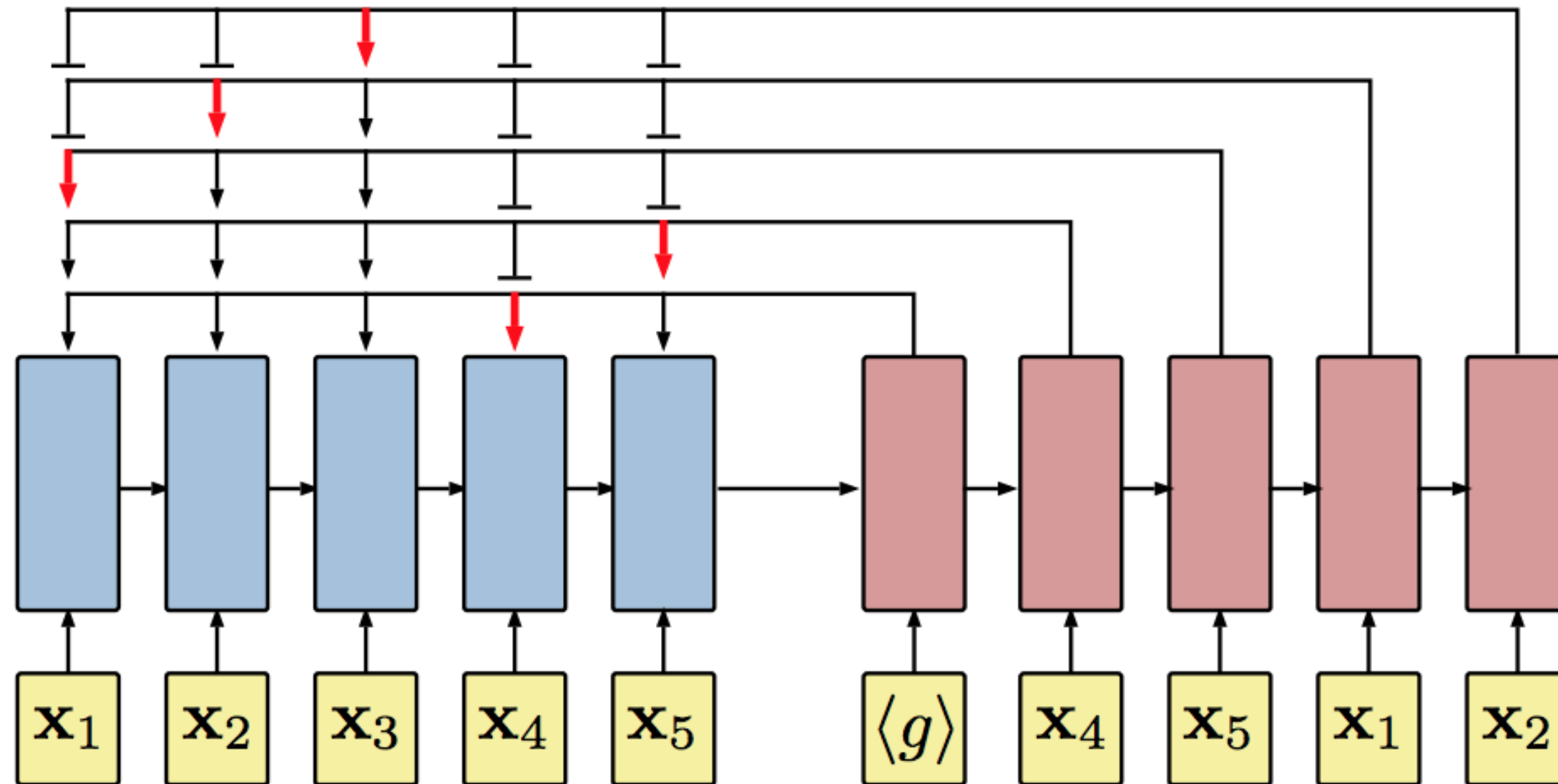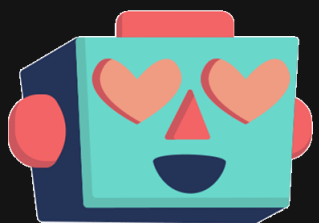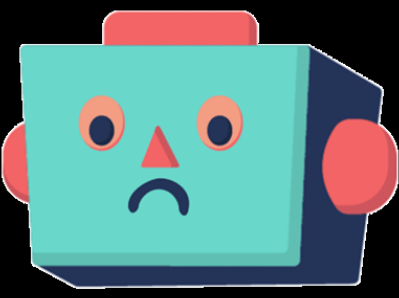
ENCODER

DECODER

# Pointer Network



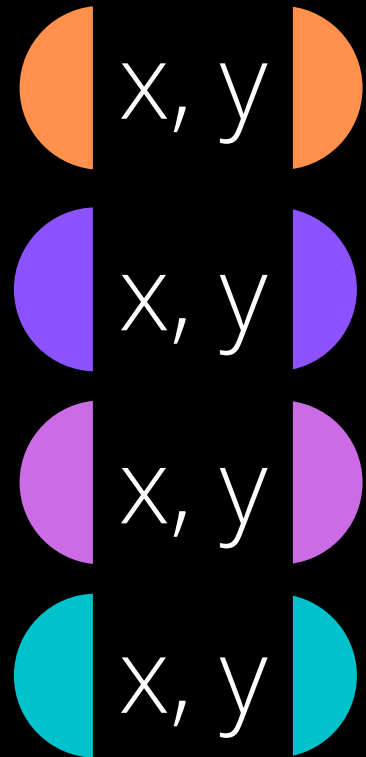Figure 1: A pointer network architecture introduced by (Vinyals et al., 2015b).

# Graph embeddings

Graph represenations:
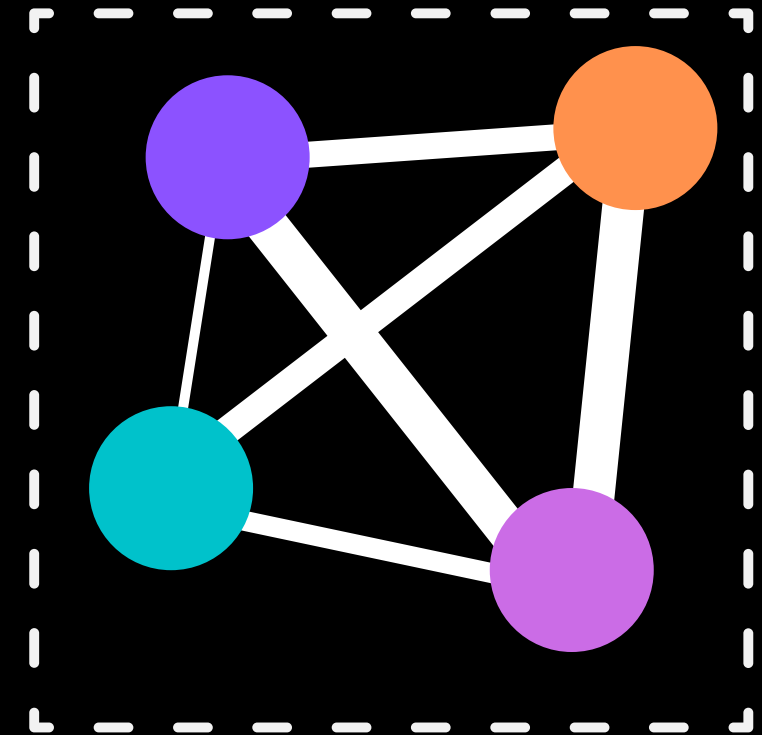
list of coordinates:

- Linear layer
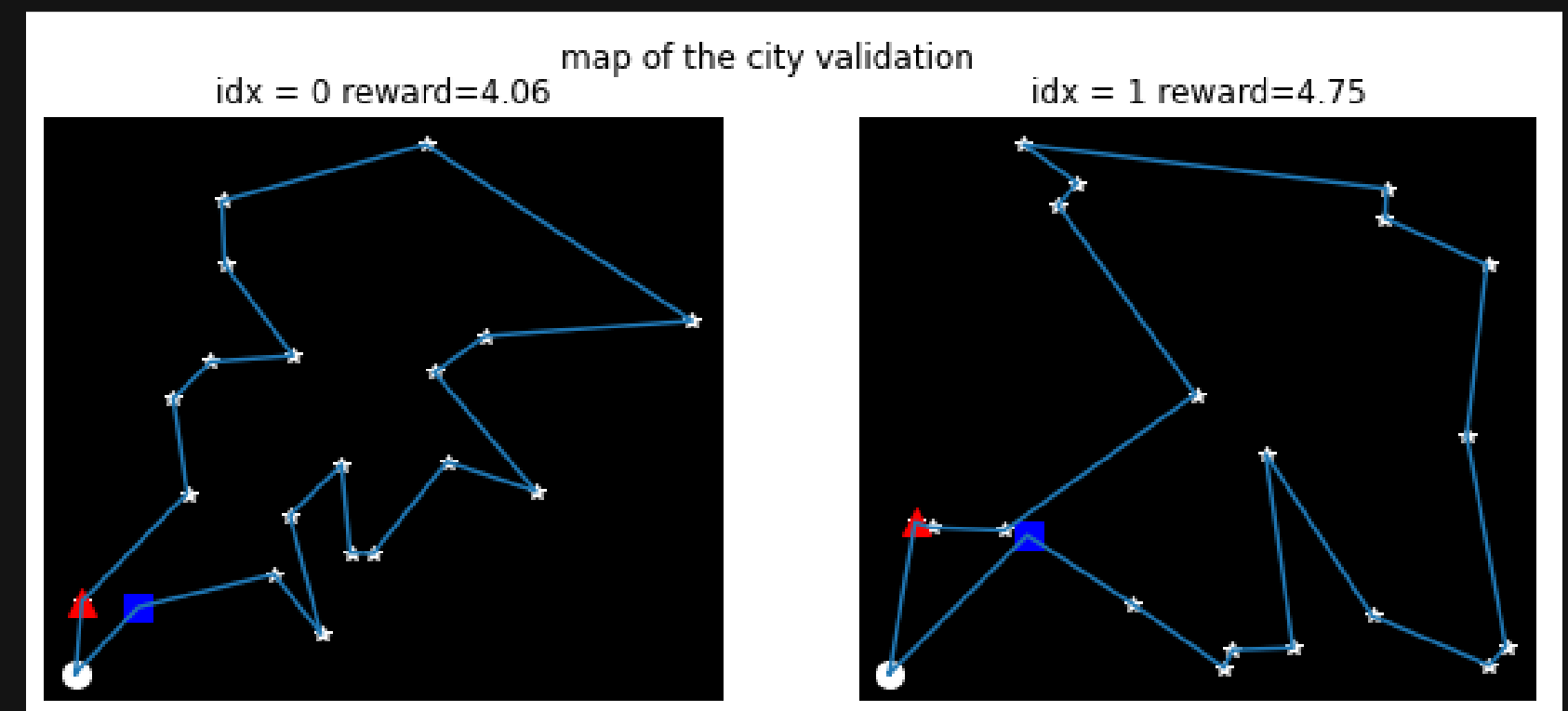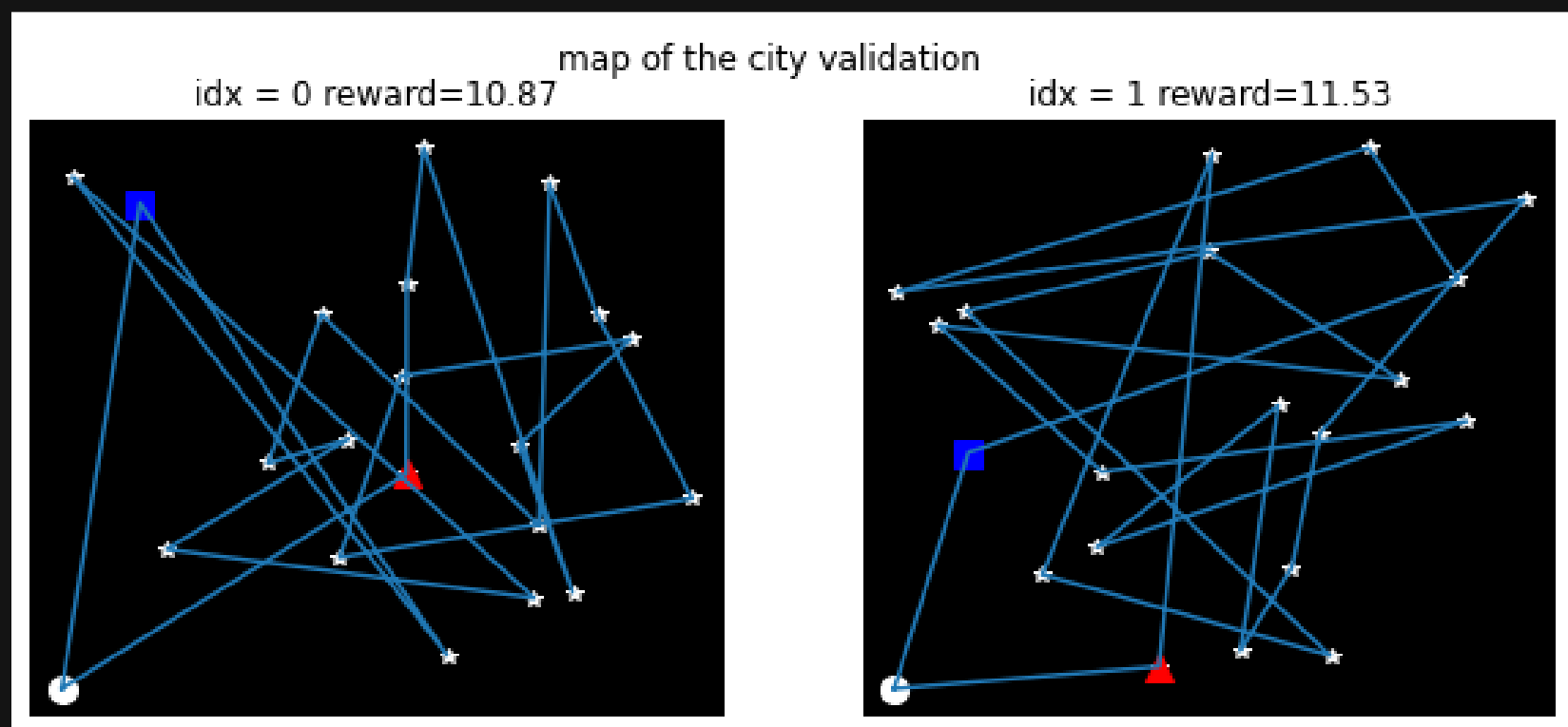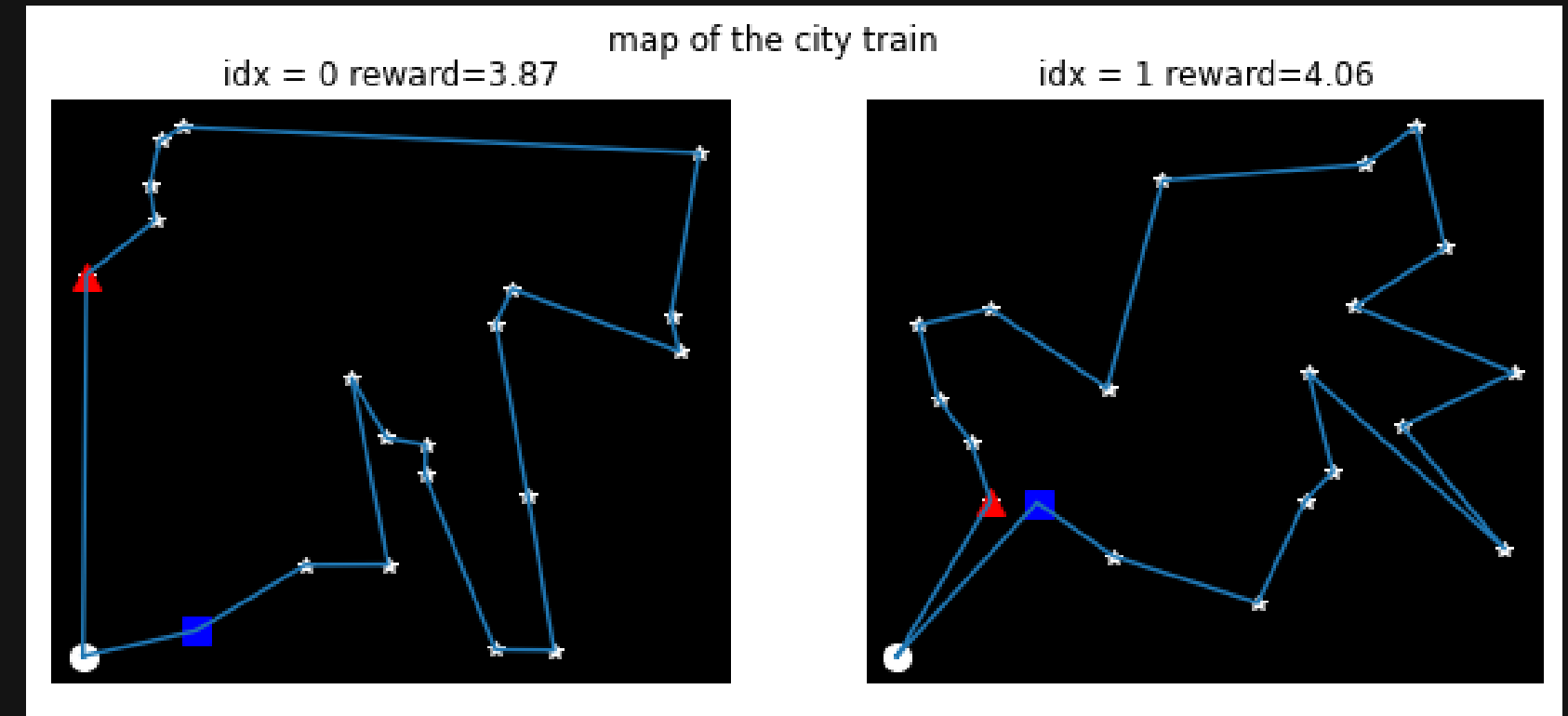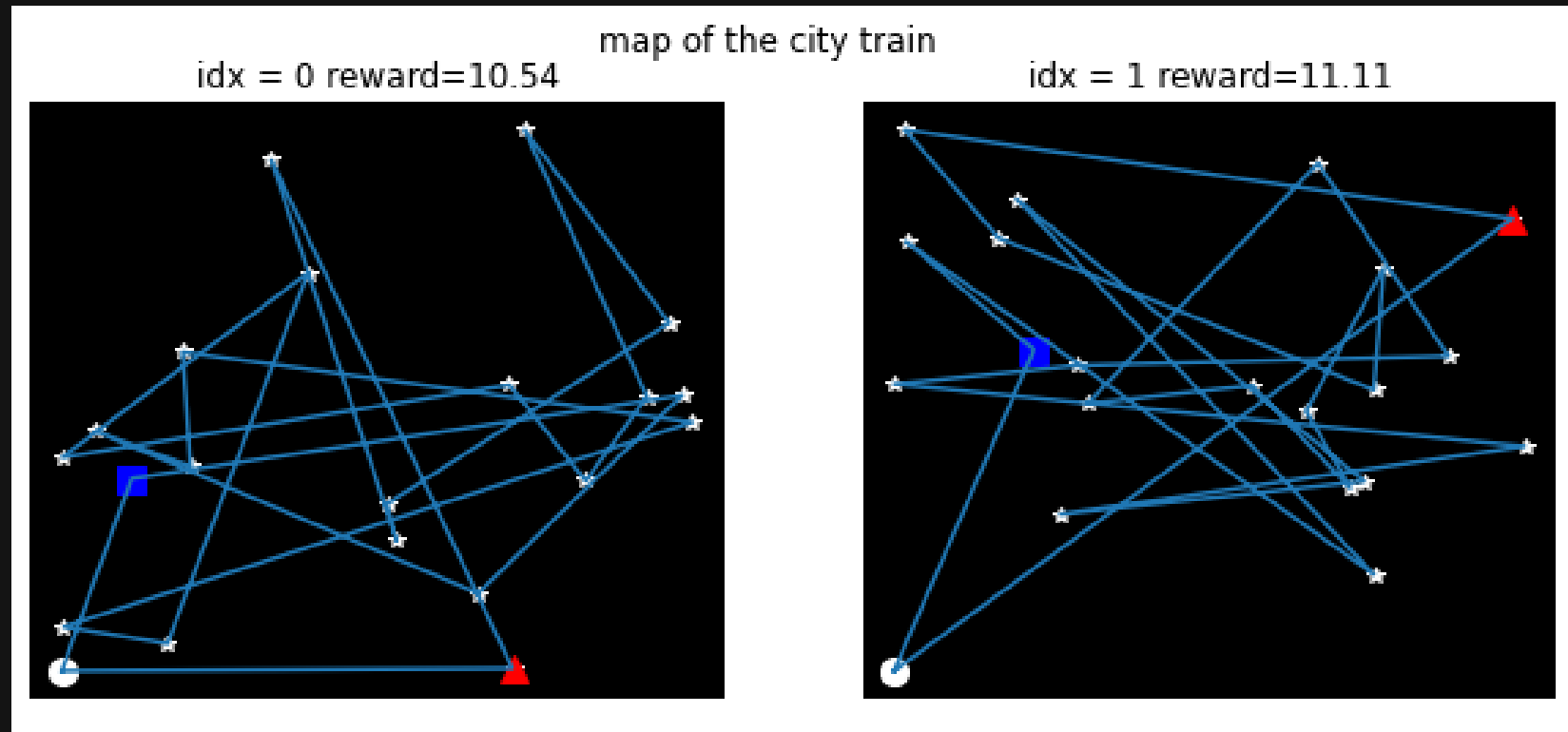- Simple node encoding $v_{node} \cdot v_{rand}$

x, y

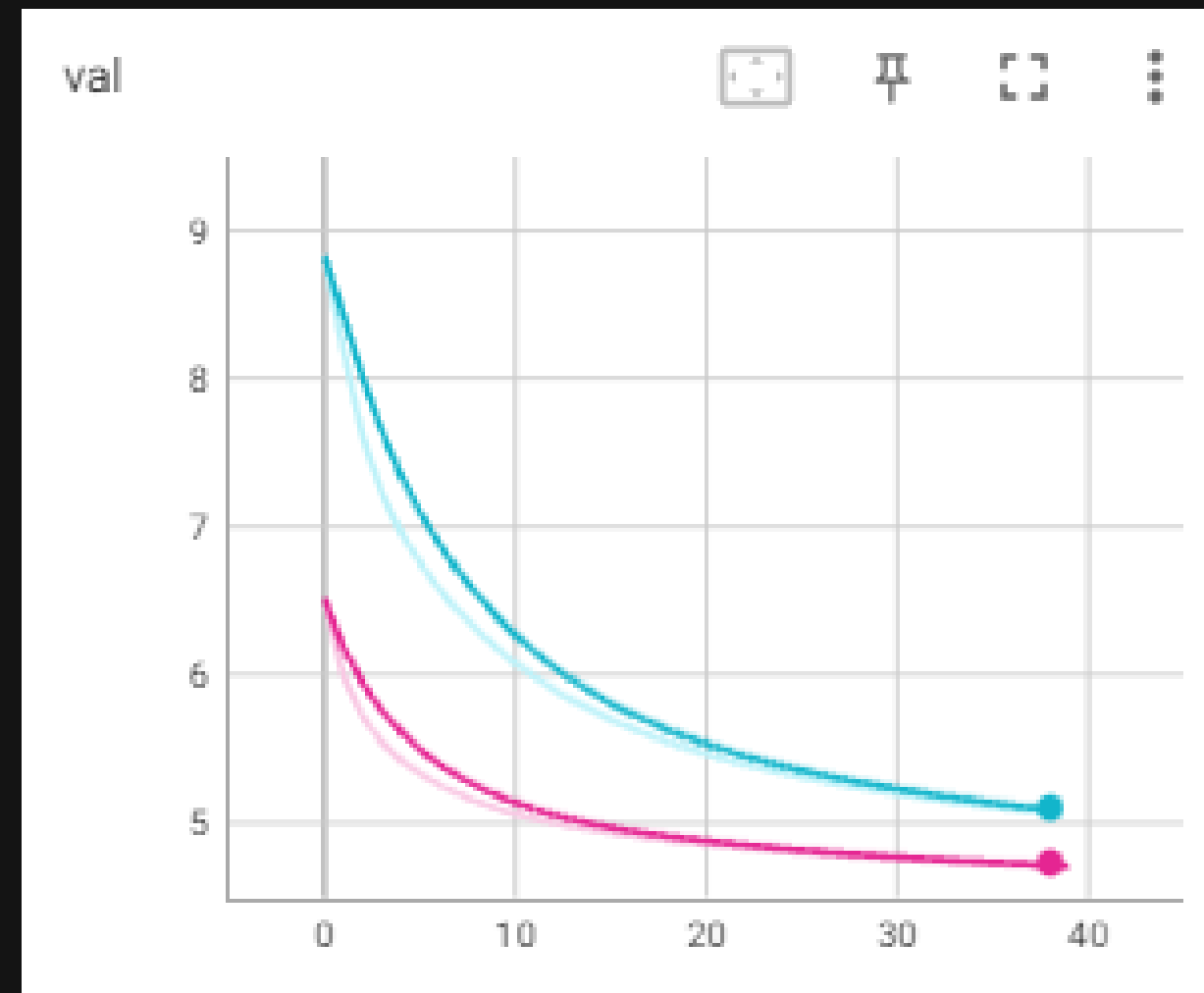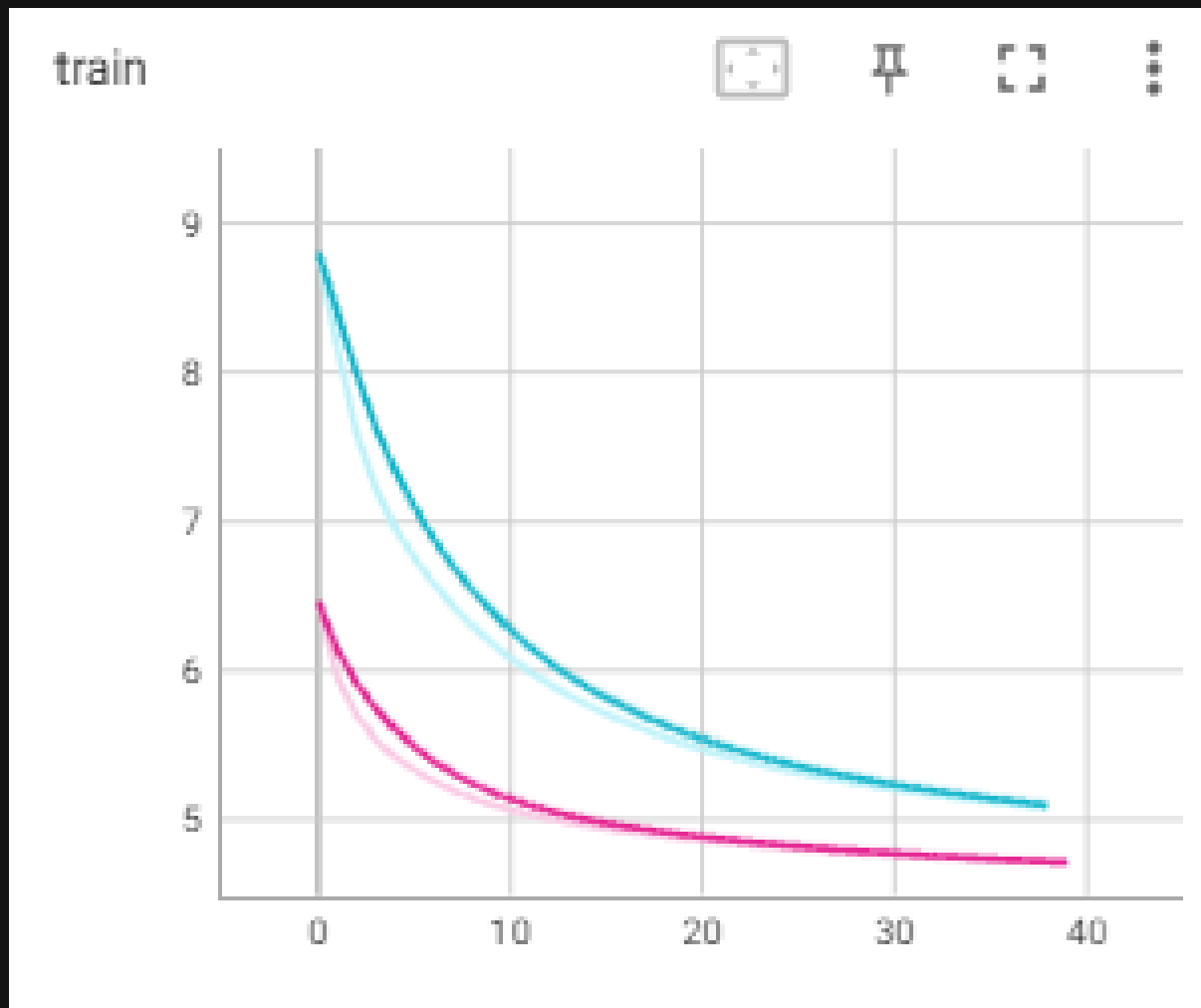x, y

x, y

x, y

weighted graph (distance matrix):

- Node2Vec
- DeepWalk

# Results



map of the city train
idx = 0 reward=10.54  idx = 1 reward=11.11

map of the city train
idx = 0 reward=3.87  idx = 1 reward=4.06

map of the city validation
idx = 0 reward=10.87  idx = 1 reward=11.53

map of the city validation
idx = 0 reward=4.06  idx = 1 reward=4.75

# Results

# Results

**Attention:**
**Bahdanau**
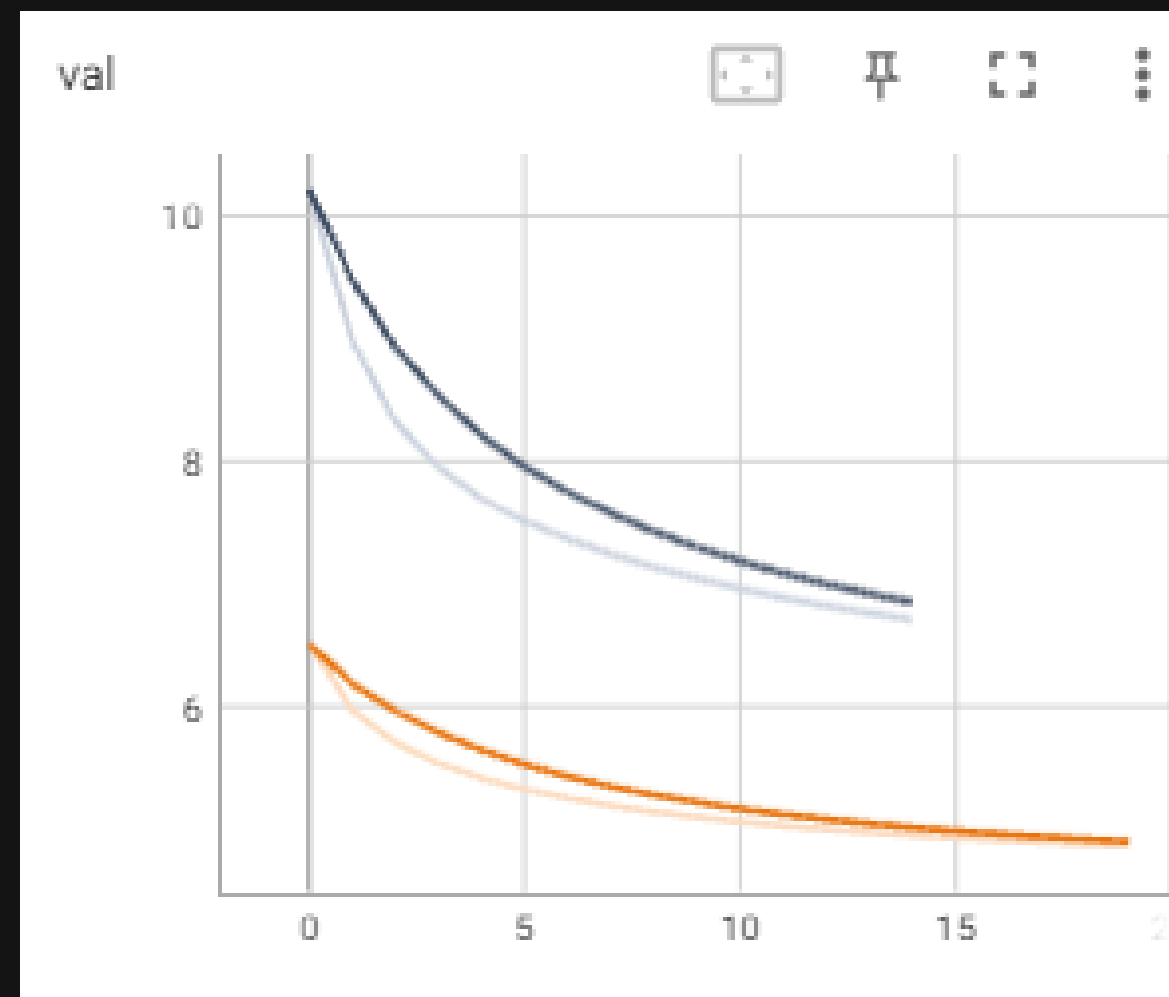**Embeddings:**
- **Linear GE**
- **Simple GE**

# Results



Linear embeddings:

N = 200_000
N = 400_000

11

# Results



**Embeddings:**

- **Node2Vec**
- **DeepWalk**

# GitHub repo

Documentation.md

# Conclusion

TSP problem: Pointer Network (Attention) + model-free policy-based optimization (REINFORSE)

Linear Node Embeddings work better with both types of Attention

Showed the effect of the size of the Network on the results

# Thank you for your attention!!!