

# Embodied Vision-and-Language Navigation with Dynamic Convolutional Filters

Federico Landi, Lorenzo Baraldi, Massimiliano Corsini, Rita Cucchiara  
University of Modena and Reggio Emilia

## 1. Motivations

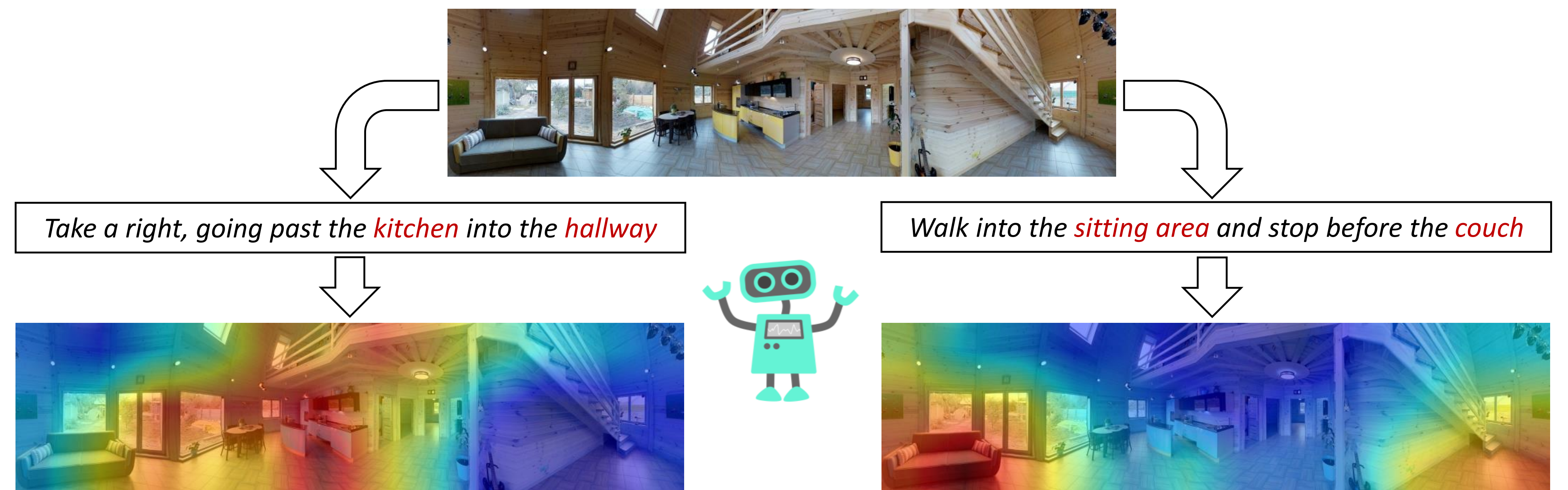
In Vision-and-Language Navigation (VLN), an embodied agent needs to reach a target destination with the only guidance of a natural language instruction. We exploit dynamic convolutional filters to ground the lingual description into the visual observation in an elegant and efficient way.

### Our contributions:

- New encoder-decoder architecture which employs dynamic convolutional filters;
- Novel categorization for VLN: we distinguish between *low-level* and *high-level* actions methods;
- State-of-Art results for low-level VLN.

## 2. Dynamic Convolutional Filters

Rather than learning a fixed set of convolutional filters, we learn to generate them depending on the natural language specification.



Dynamic convolutional filters act as specialized and flexible feature extractors.

## 3. Architecture

For each navigation episode, the agent receives a natural language instruction that is encoded into  $X$ . At each time step  $t$ , the agent observes  $I_t$  from the surroundings. The goal is to decode the atomic action for the current step.

### Encoder-Decoder Attention:

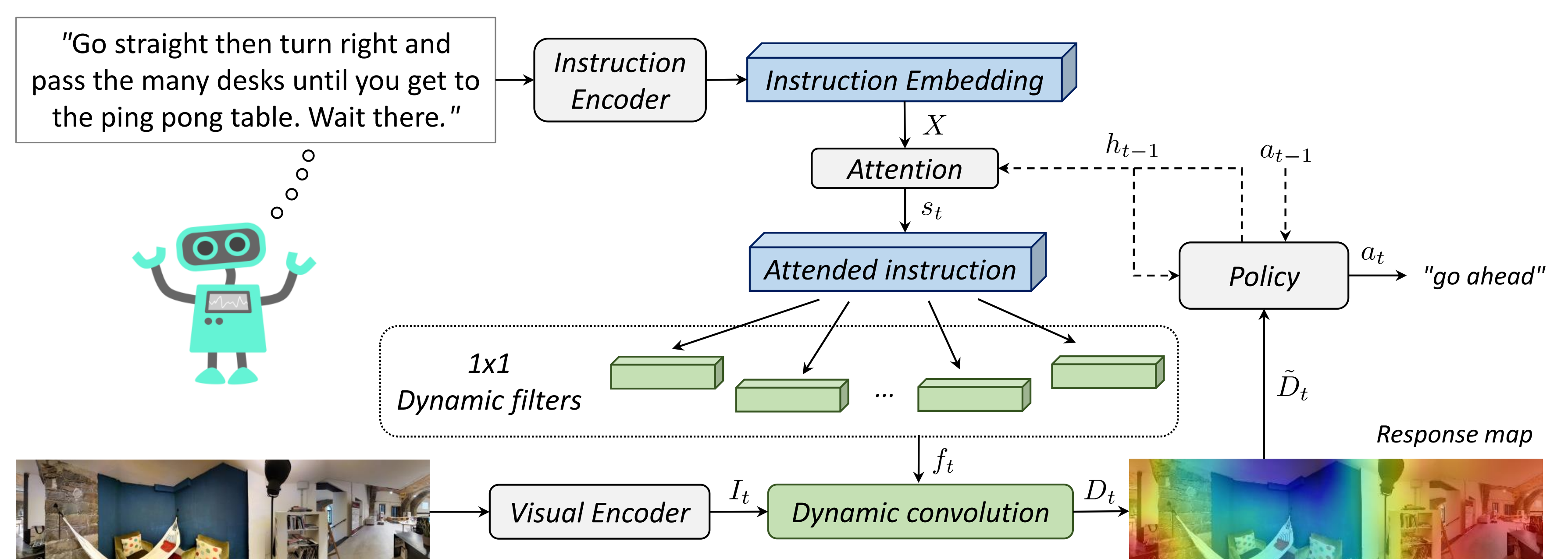
$$q_t = W_q h_{t-1} + b_q \quad K = W_k X + b_k$$

$$\alpha_t = q_t K^T / \sqrt{d_{att}} \quad s_t = \text{softmax}(\alpha_t) X$$

### Dynamic Convolution: Action Decoding:

$$f_t = \ell_2[\tanh(W_f s_t + b_f)] \quad h_t = \text{LSTM}([\tilde{D}_t, a_{t-1}], h_{t-1})$$

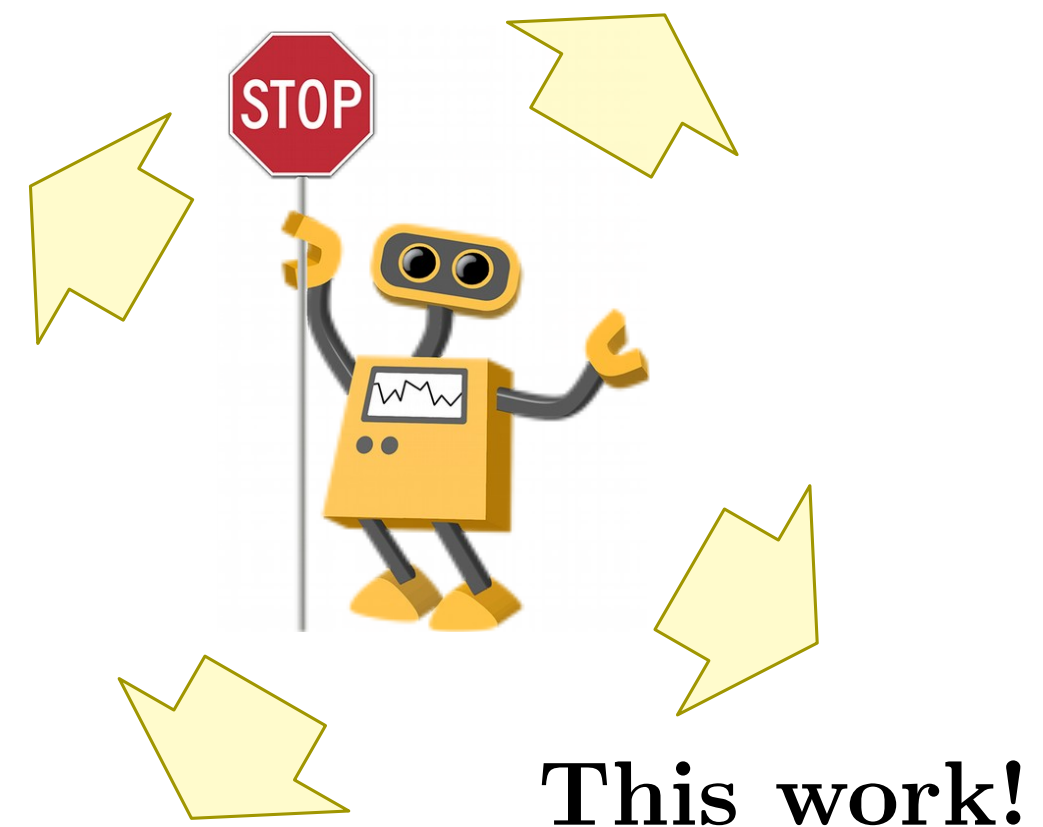
$$D_t = f_t * I_t \quad p_t = \text{softmax}(W_a h_t + b_a)$$



## 4. Low-level and High-level Methods

### Low-level Action Space

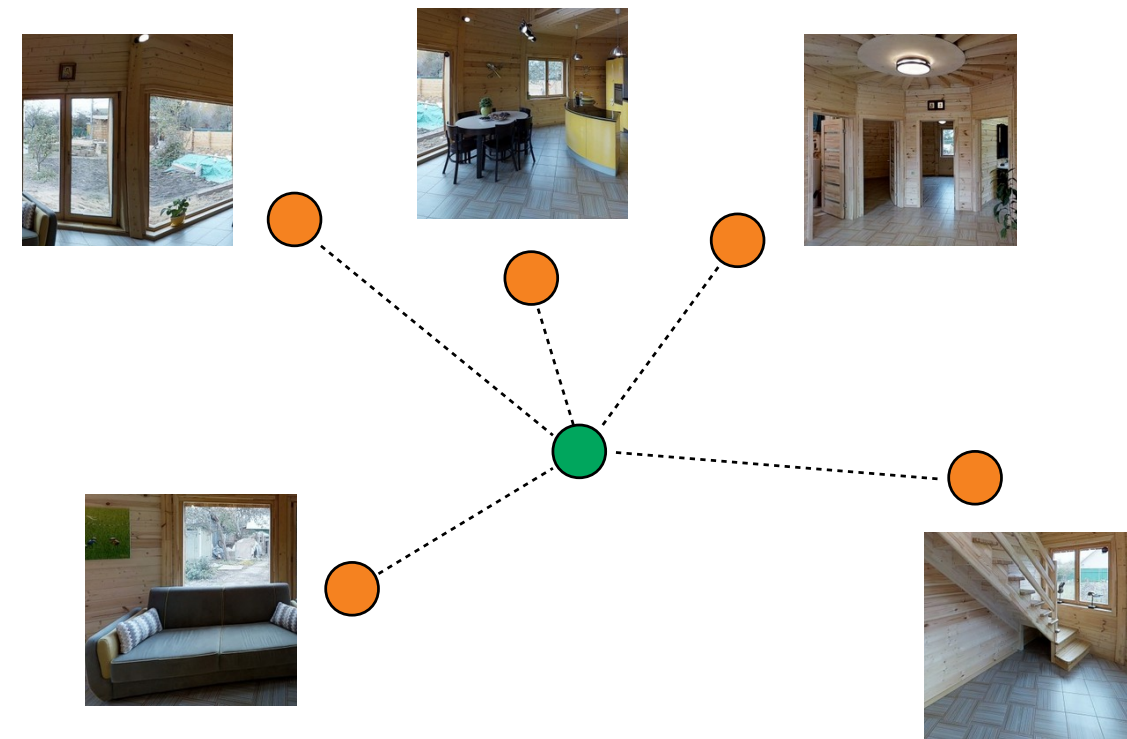
Simulates continuous control



This work!

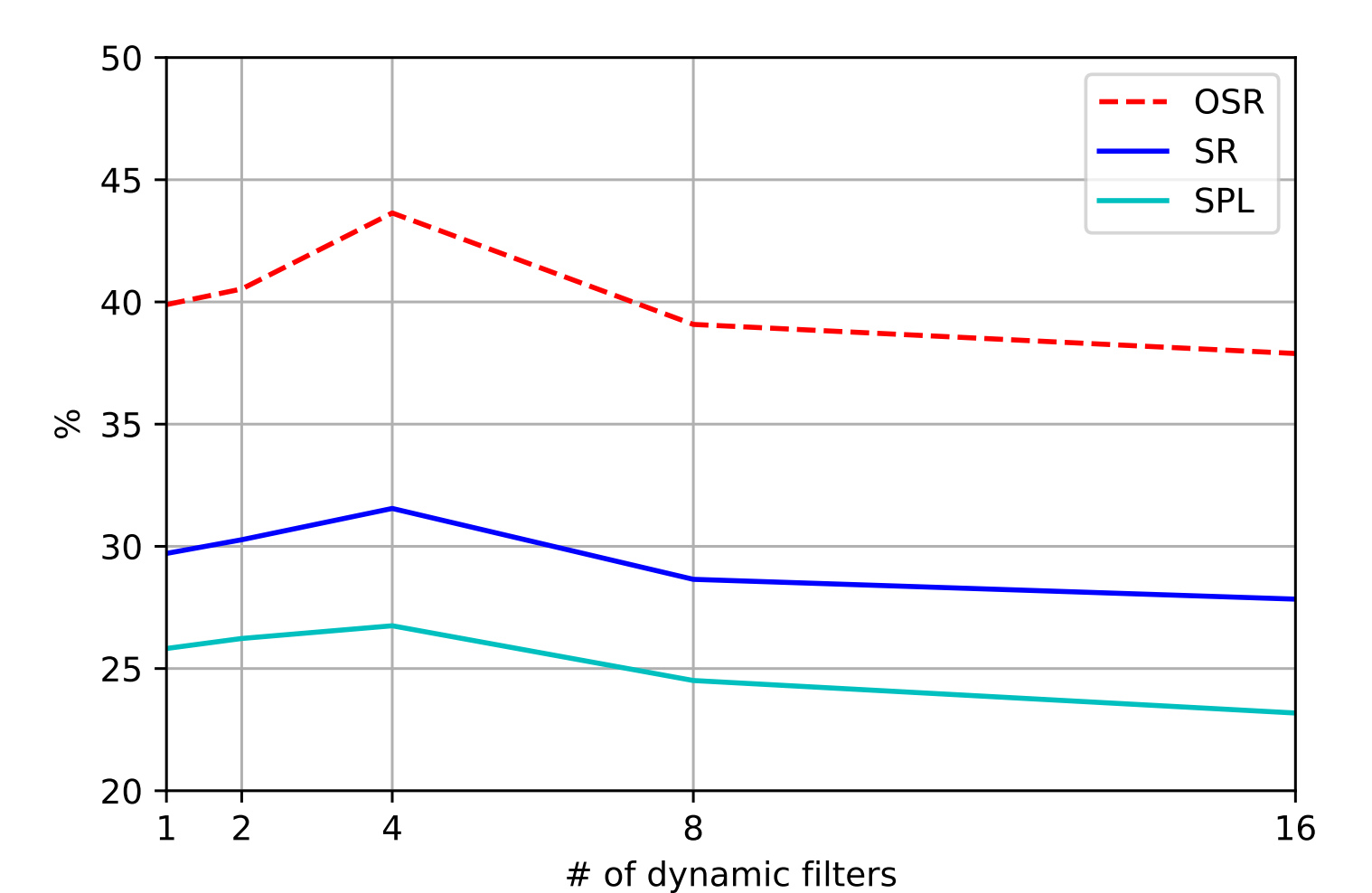
### High-level Action Space

Path selection on a discrete graph



## 6. Number of Dynamic Filters

# of filters	Val-Unseen	
	NE ↓	SR ↑
1	6.79	29.7
2	6.77	30.3
4	<b>6.65</b>	<b>31.6</b>
8	7.19	28.7
16	7.03	27.8



One filter is enough, best setup with four response maps.

## 5. Qualitative Results

Next action: **L** left **R** right **F** forward **E** end episode



**Instruction:** Walk up the stairs. Turn right at the top of the stairs and walk along the red ropes. Walk through the open doorway straight ahead along the red carpet. Walk through that hallway into the room with couches and a marble coffee table.

## 7. Comparison with State-of-the-Art

**Low-level (ours):** possible actions are *move forward, turn left 30°, turn right 30°, raise elevation, lower elevation, and end episode.*

Low-level Actions Methods	Validation-Seen				Validation-Unseen				Test (Unseen)			
	NE ↓	SR ↑	OSR ↑	SPL ↑	NE ↓	SR ↑	OSR ↑	SPL ↑	NE ↓	SR ↑	OSR ↑	SPL ↑
Random	9.45	0.16	0.21	-	9.23	0.16	0.22	-	9.77	0.13	0.18	0.12
Student-forcing [1]	6.01	0.39	0.53	-	7.81	0.22	0.28	-	7.85	0.20	0.27	0.18
RPA [2]	5.56	0.43	0.53	-	7.65	0.25	0.32	-	7.53	0.25	0.33	0.23
Ours	4.68	0.53	0.66	0.46	6.65	0.32	0.44	0.27	7.14	0.31	0.42	0.27
Ours w/ data augmentation	<b>3.96</b>	<b>0.58</b>	<b>0.73</b>	<b>0.51</b>	<b>6.52</b>	<b>0.34</b>	<b>0.43</b>	<b>0.29</b>	<b>6.55</b>	<b>0.35</b>	<b>0.45</b>	<b>0.31</b>

**High-level:** the agent selects the destination from a set of adjacent nodes. Moves are made thanks to global information from the simulator.

High-level Actions Methods	Validation-Seen				Validation-Unseen				Test (Unseen)			
	NE ↓	SR ↑	OSR ↑	SPL ↑	NE ↓	SR ↑	OSR ↑	SPL ↑	NE ↓	SR ↑	OSR ↑	SPL ↑
Speaker-Follower [3]	3.36	0.66	0.74	-	6.62	0.36	0.45	-	6.62	0.35	0.44	0.28
Self-Monitoring [4]	<b>3.22</b>	0.67	<b>0.78</b>	0.58	5.52	0.45	0.56	0.32	5.99	0.43	0.55	0.32
Regretful [5]	3.23	<b>0.69</b>	0.77	<b>0.63</b>	<b>5.32</b>	<b>0.50</b>	<b>0.59</b>	<b>0.41</b>	<b>5.69</b>	<b>0.48</b>	<b>0.56</b>	<b>0.40</b>

NE: Navigation Error (m); OSR: Oracle Success Rate;

SR: Success Rate;

SPL: Success rate normalized on Path Length.

## 8. References

- [1] Anderson *et al.* Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *CVPR*, 2018.
- [2] Wang *et al.* Look before you leap: Bridging model-free and model-based reinforcement learning for planned-ahead VLN. In *ECCV*, 2018.
- [3] Fried *et al.* Speaker-follower models for vision-and-language navigation. In *NeurIPS*, 2018.
- [4] Ma *et al.* Self-Monitoring Navigation Agent via Auxiliary Progress Estimation. In *ICLR*, 2019.
- [5] Ma *et al.* The Regretful Agent: Heuristic-Aided Navigation through Progress Estimation. In *CVPR*, 2019.