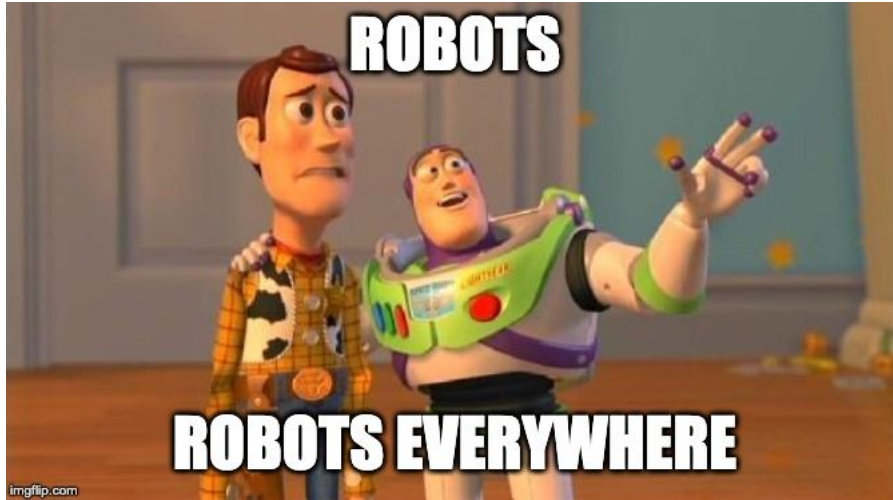


# Reinforced Cross-Modal Matching and Self-Supervision Imitation Learning for Vision-Language Navigation



Presented by George Revkov, Yurong Zhuang, Yuchen Zhou

# Vision-Language Navigation



# Vision-Language Navigation

Task of navigating inside real environments following natural language instructions



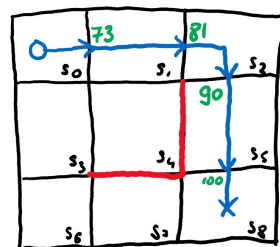
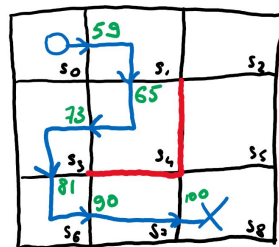
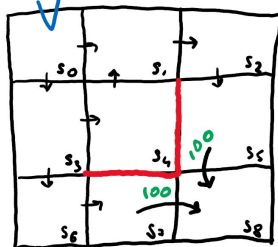
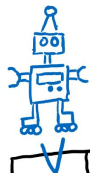
Leave the bedroom, and enter the kitchen. Walk forward, and take a left at the couch. Stop in front of the window.

# Aside: Reinforcement Learning

- **Agent** observes the **environment**, identifies its current **state**  $s_t \in \mathbf{S}$ , takes an **action**  $a_t \in \mathbf{A}$ , and gets a **reward** (may be = 0)
- Goal 1: from each possible state find a sequence of actions that will maximize total reward
- Goal 2: come up with a **policy**  $\pi: \mathbf{S} \rightarrow \mathbf{A}$
- **Markov Decision Process (MDP)**

# Aside: Reinforcement Learning

- States  $\mathbf{S} = \{s_0 \dots s_8\}$
- Actions  $\mathbf{A}(S_t) = \{U, D, R, L\}$
- Transition F-n  $\mathbf{T}: \mathbf{S}, \mathbf{A} \rightarrow \mathbf{S}'$
- Reward F-n  $\mathbf{R}(\mathbf{S}, \mathbf{A})$ 
  - $r_t = r(s_t, a_t)$



- Policy  $\pi: \mathbf{S} \rightarrow \mathbf{A}$
- **Discounted Cumulative Reward:**

$$V^\pi(s_t) = r_t + \gamma * r_{t+1} + \gamma^2 * r_{t+2} + \dots$$

- Say  $\gamma = .9$
- $V^\pi(s_0) = 0 + 0 * .9 + 0 * .9^2 + 100 * .9^3 = 73$
- $\pi(S = s_0) = R, \pi(S = s_1) = R,$   
 $\pi(S = s_2) = D, \pi(S = s_5) = D$

# Challenges

- Reasoning over visual images & natural language instructions is difficult
  - Ground an instruction in the local scene
  - Match the instruction to the visual trajectory over time
- Ill-posed feedback
- Generalization problem

# Summary

- Novel Reinforced Cross-Modal Matching approach
  - **Reasoning Navigator** (based on where I am right now / what I'm seeing):
    - What should I be doing?
    - Which sub-instruction should execute?
  - **Matching critic**:
    - What is the probability of reconstructing the instruction based on executed trajectory
    - Penalize the trajectories that don't match instructions





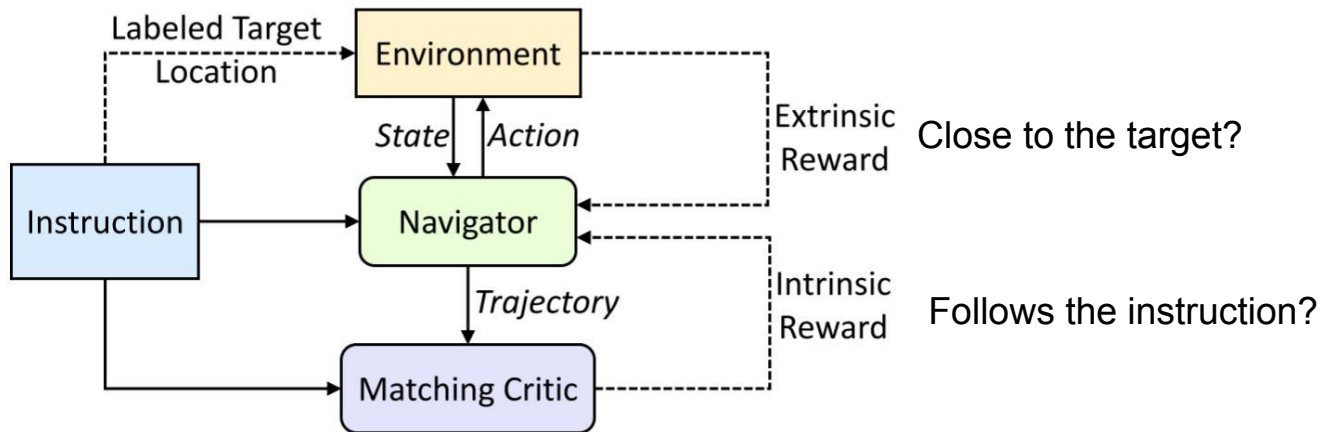
# Summary

- Use Reinforcement Learning with two reward functions
- Use Self-Supervised Imitation learning to deal with unseen environments
- *Establishes new state-of-the-art performance* on R2R dataset
- Proposes a new evaluation setting for VLN:
  - Exploring unseen environments prior to testing is allowed

# RCM Model



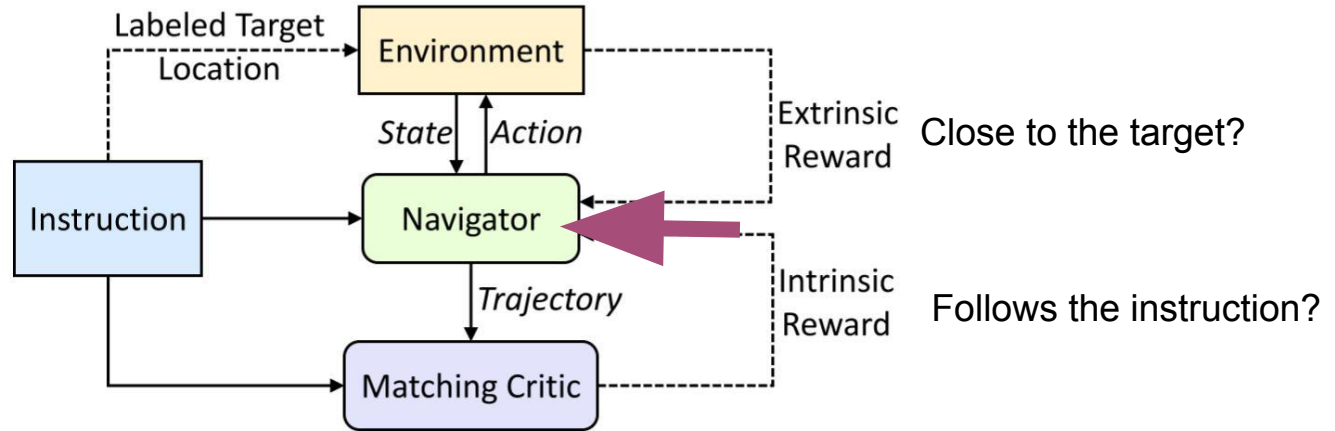
*Turn right and head towards the kitchen. Then turn left, pass a table and enter the hallway. ...Stop in front of the toilet.*



# RCM Model



*Turn right and head towards the kitchen. Then turn left, pass a table and enter the hallway. ...Stop in front of the toilet.*



# Navigator

turn completely around until you face an open door with a window to the left and a patio to the right, walk forward though the door and into a dinning room, ... ..

Language Encoder

$\{w_i\}_{i=1}^n$

Attention

$c_t^{text}$

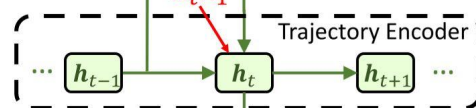


Panoramic Features

$\{v_{t,j}\}_{j=1}^m$

Attention

$a_{t-1}$

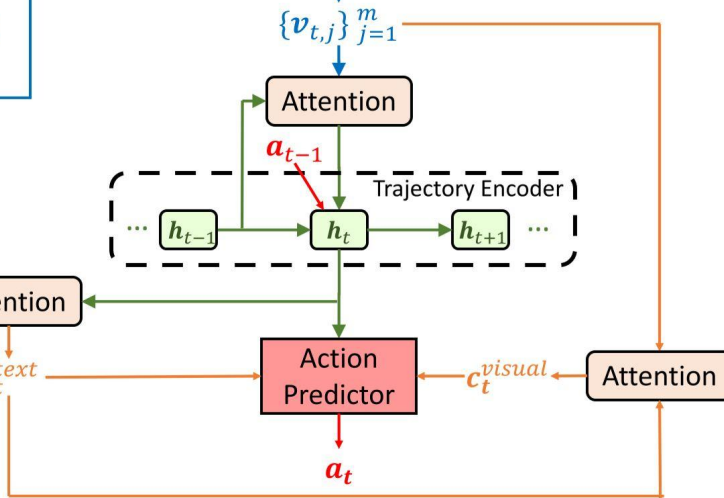


Action Predictor

$a_t$

$c_t^{visual}$

Attention



# Navigator

turn completely around until you face an open door with a window to the left and a patio to the right, walk forward though the door and into a dinning room, ... ..

Language Encoder

$\{w_i\}_{i=1}^n$

Attention

$c_t^{text}$

Action Predictor

$a_t$

split by different viewpoints



Panoramic Features

$\{v_{t,j}\}_{j=1}^m$

Attention

$a_{t-1}$

Trajectory Encoder

$h_{t-1}$

$h_t$

$h_{t+1}$

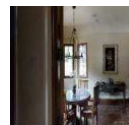
Attention

$c_t^{visual}$



CNN

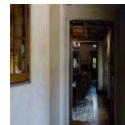
$v_1$



CNN

$v_2$

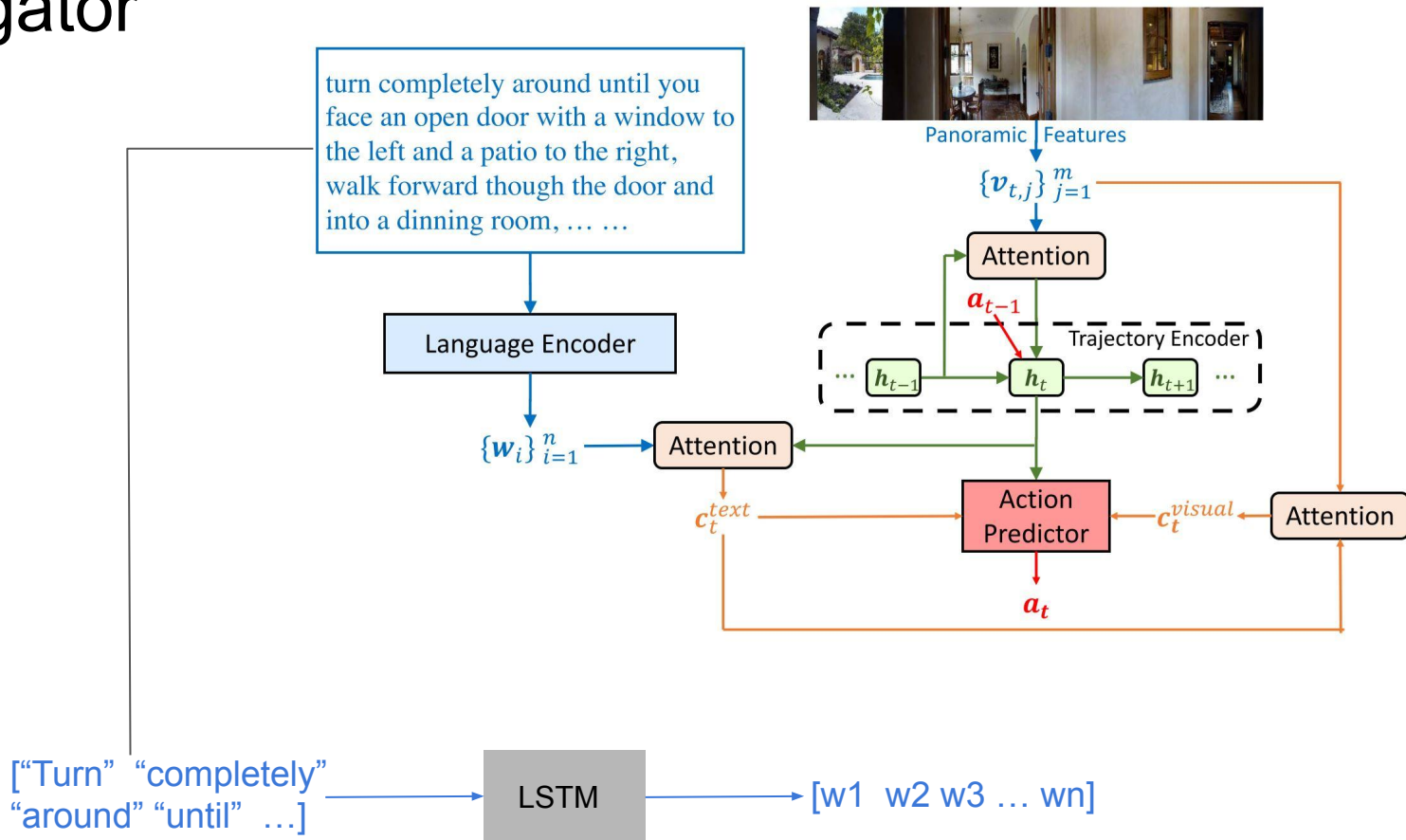
...



CNN

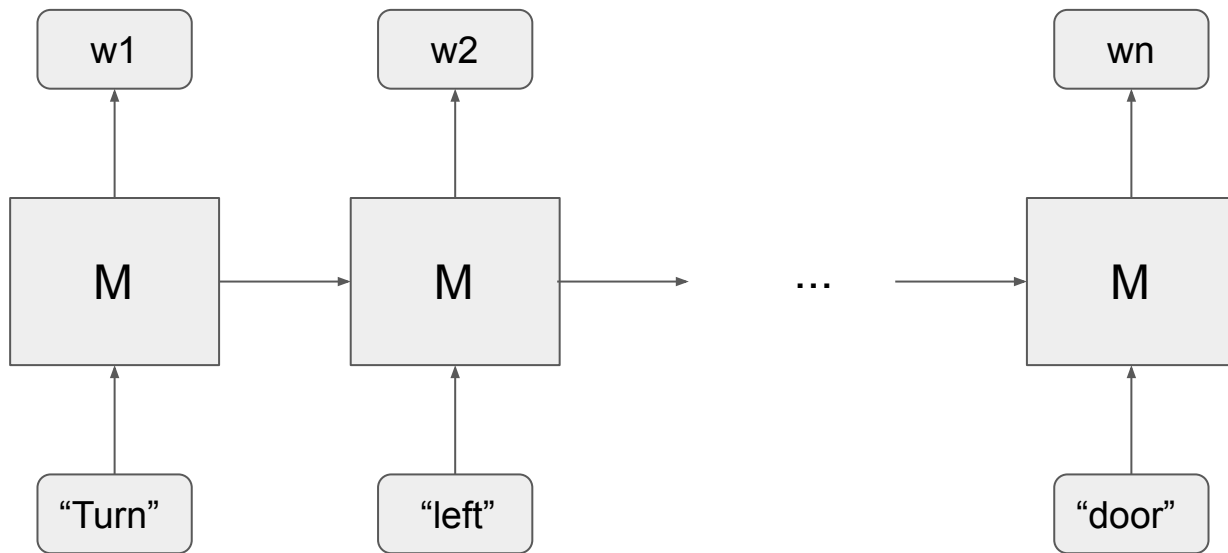
$v_m$

# Navigator

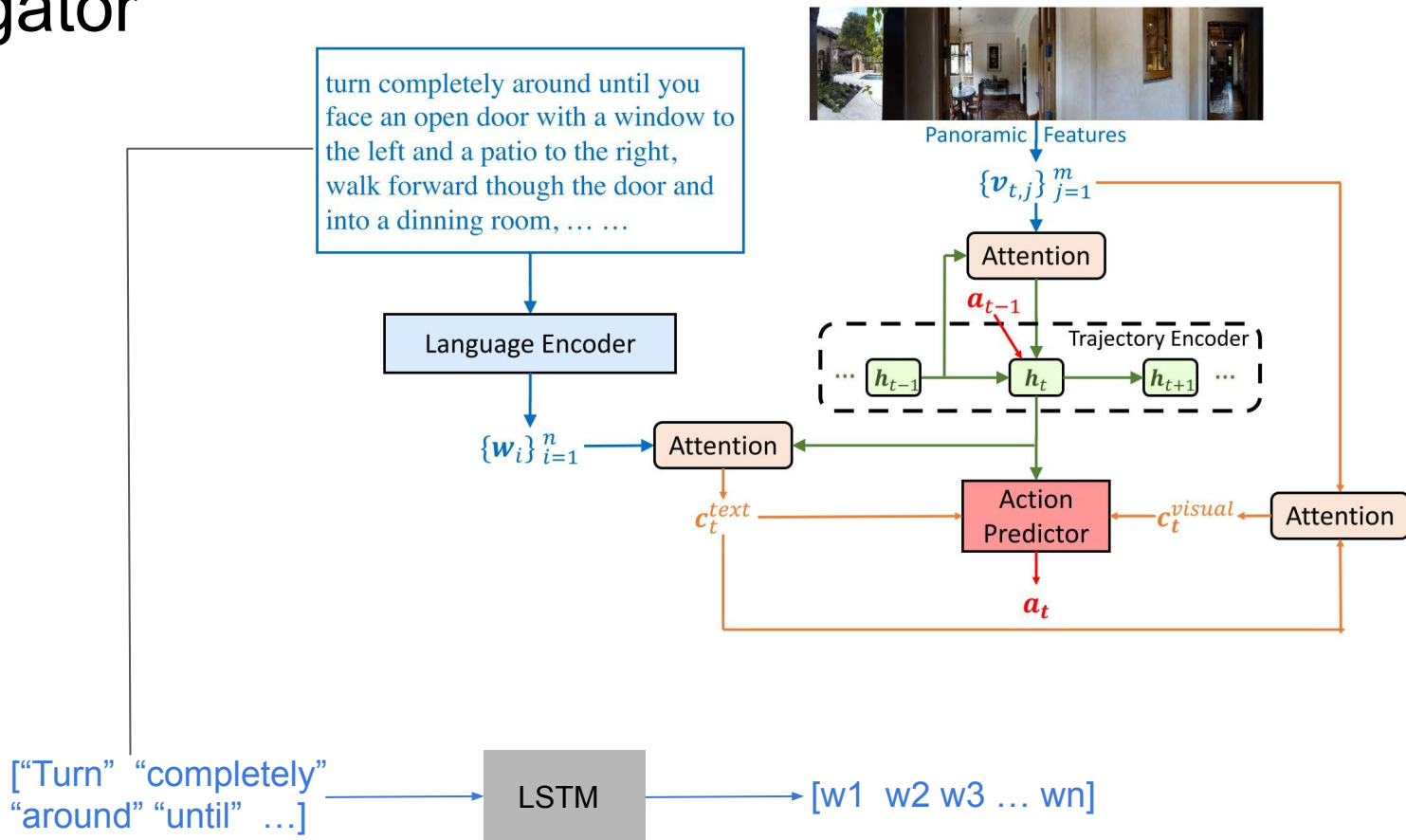


# LSTM

- Good for representation learning of sequential data

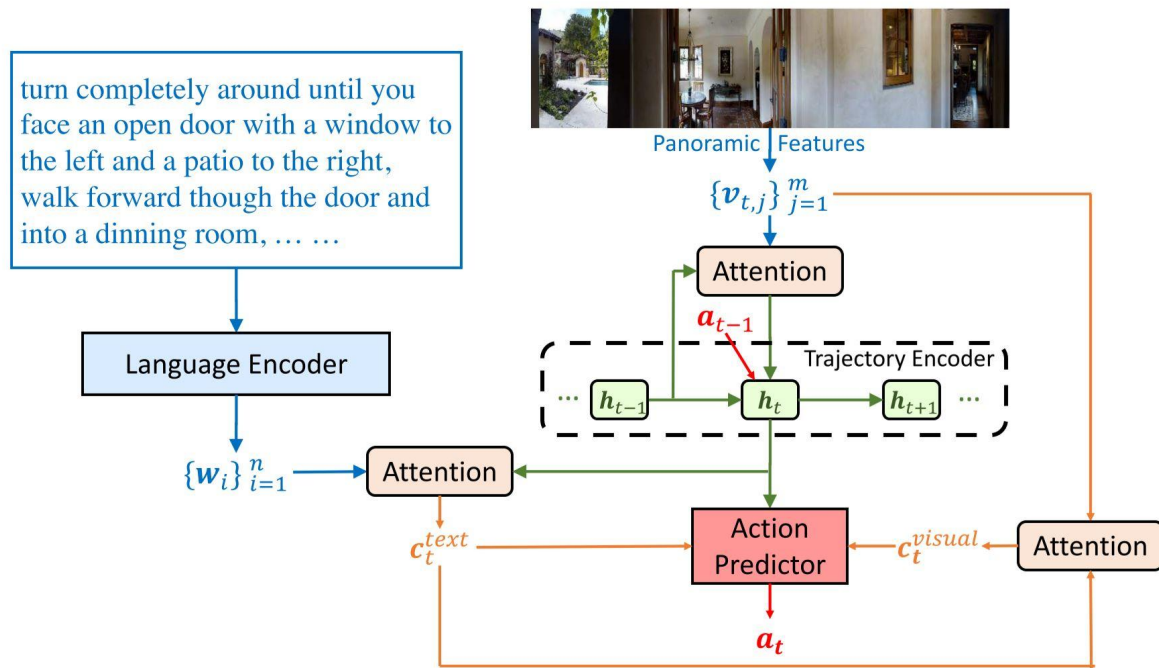


# Navigator





# Cross modal reasoning

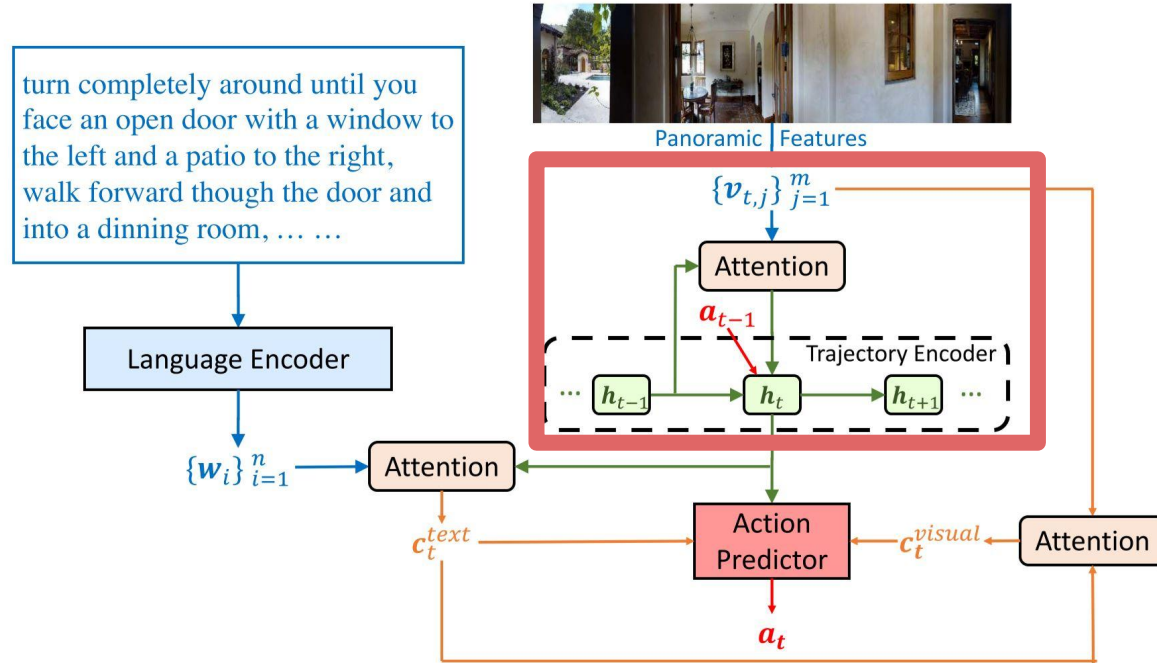


Learns history based on previous action and visual input

Learns how to focus on textual instruction based on history

Learns how to focus on visual input based on textual instruction

# Cross modal reasoning

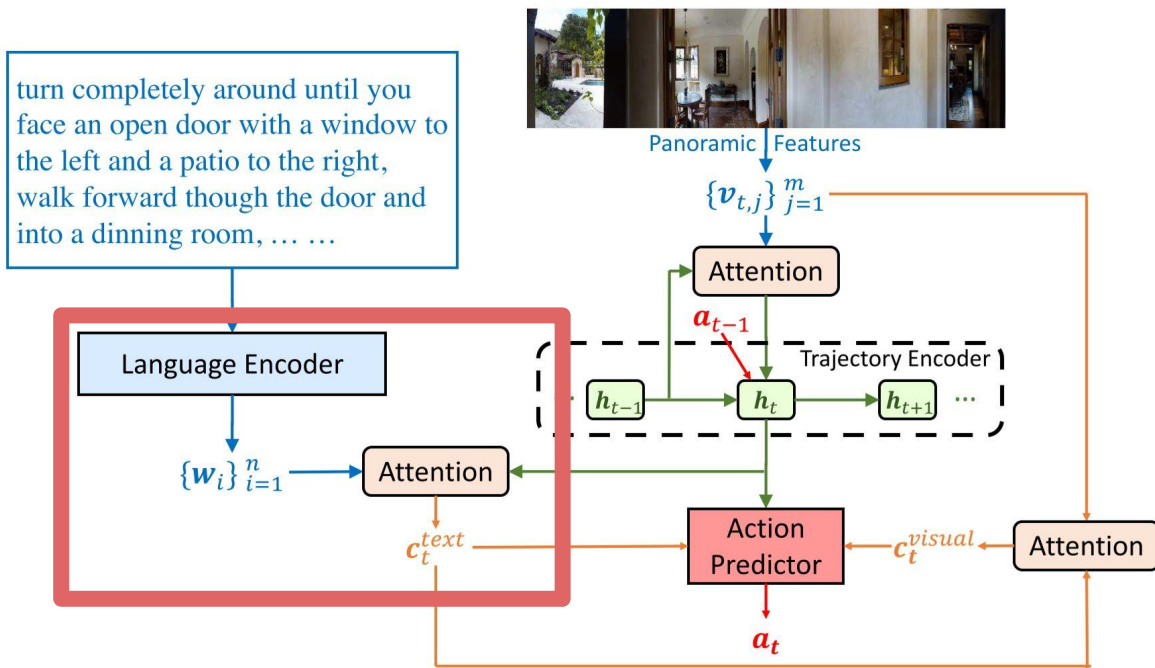


Learns history based on previous action and visual input

Learns how to focus on textual instruction based on history

Learns how to focus on visual input based on textual instruction

# Cross modal reasoning

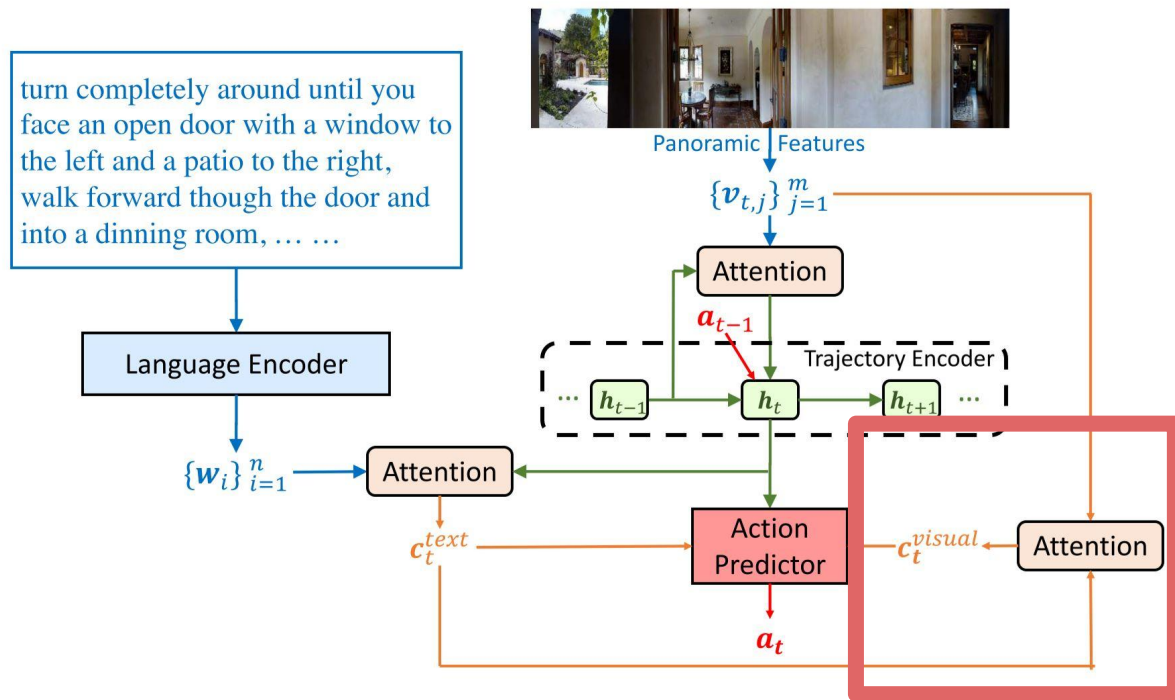


Learns history based on previous action and visual input

Learns how to focus on textual instruction based on history

Learns how to focus on visual input based on textual instruction

# Cross modal reasoning

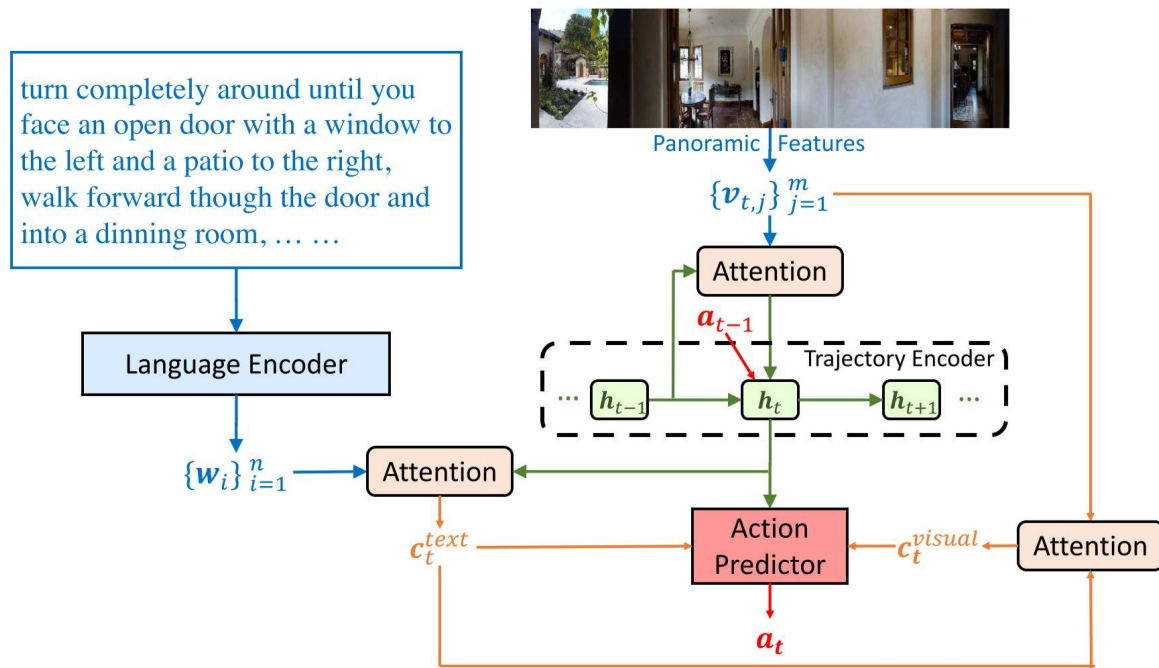


Learns history based on previous action and visual input

Learns how to focus on textual instruction based on history

Learns how to focus on visual input based on textual instruction

# Cross modal reasoning



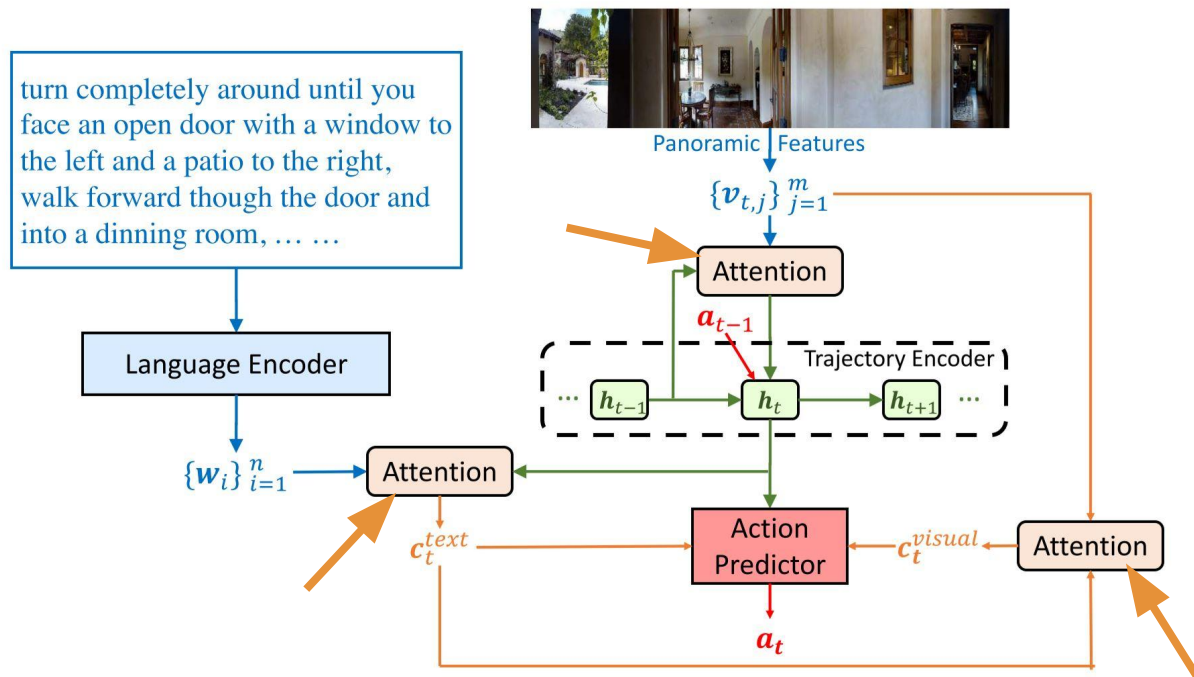
Learns history based on previous action and visual input

Learns how to focus on textual instruction based on history

Learns how to focus on visual input based on textual instruction

Text  $\rightleftarrows$  Image

# Cross modal reasoning



Learns history based on previous action and visual input

Learns how to focus on textual instruction based on history

Learns how to focus on visual input based on textual instruction

# Attention (generic)

Depending on the state, which input should be focused on (weighted more) ?

State at time t

Set of inputs

$$\hat{x} = \underset{n}{\text{attention}}(s_t, \{x_i\}_{i=1}^n) = \sum_{i=1}^n \text{softmax}(s_t W_s (x_i W_x)^T) x_i$$

# Attention (generic)

Depending on the state, which input should be focused on (weighted more) ?

State at time t

Set of inputs

$$\hat{x} = \underset{i=1}{\overset{n}{\text{attention}}}(s_t, \{x_i\}_{i=1}^n) = \sum_{i=1}^n \underbrace{\text{softmax}(s_t W_s (x_i W_x)^T)}_{\text{Weight of } x_i} x_i$$



# Attention (generic)

Depending on the state, which input should be focused on (weighted more) ?

State at time t

Set of inputs

$$\hat{x} = \text{attention}(s_t, \{x_i\}_{i=1}^n) = \sum_{i=1}^n \text{softmax}(s_t W_s (x_i W_x)^T) x_i$$

Weight of  $x_i$

How state contributes

The diagram illustrates the attention mechanism. It shows the equation  $\hat{x} = \text{attention}(s_t, \{x_i\}_{i=1}^n) = \sum_{i=1}^n \text{softmax}(s_t W_s (x_i W_x)^T) x_i$ . A bracket above the curly braces  $\{x_i\}_{i=1}^n$  is labeled "Set of inputs". A bracket above  $s_t$  is labeled "State at time t". A blue bracket under the term  $\text{softmax}(s_t W_s (x_i W_x)^T)$  is labeled "Weight of  $x_i$ ". A vertical line extends from the center of this blue bracket down to the text "How state contributes".

# Attention (generic)

Depending on the state, which input should be focused on (weighted more) ?

State at time t

Set of inputs

$$\hat{x} = \text{attention}(s_t, \{x_i\}_{i=1}^n) = \sum_{i=1}^n \text{softmax}(s_t W_s (x_i W_x)^T) x_i$$

Weight of  $x_i$

How input contributes

# History Context

turn completely around until you face an open door with a window to the left and a patio to the right, walk forward through the door and into a dining room, ... ..



Panoramic Features

$\{v_{t,j}\}_{j=1}^m$

Language Encoder

$\{w_i\}_{i=1}^n$

Attention

$c_t^{text}$

Attention

$a_{t-1}$

$h_t$

Action Predictor

$c_t^{visual}$

$a_t$

Trajectory Encoder

$h_{t-1}$

$h_{t+1}$

Attention

$$v_t = \text{attention}(h_{t-1}, \{v_{t,j}\}_{j=1}^m)$$
$$= \sum_j \text{softmax}(h_{t-1} W_h (v_{t,j} W_v)^T) v_{t,j}$$

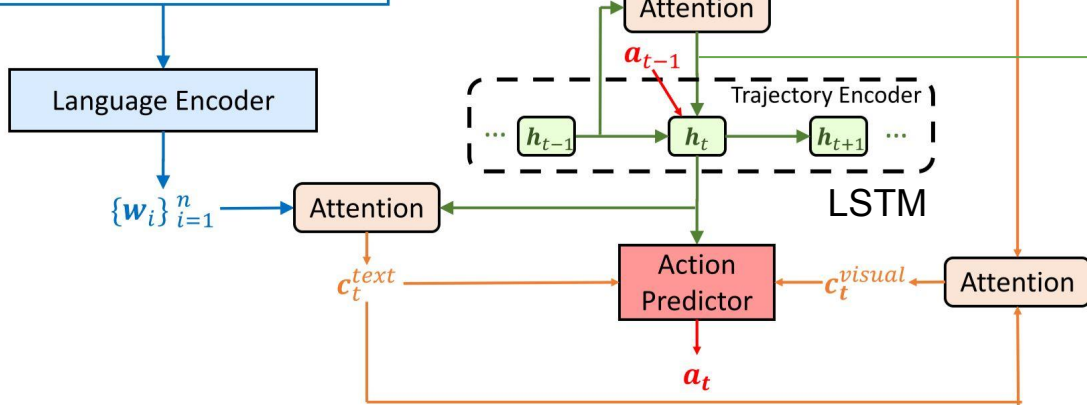
# History Context

turn completely around until you face an open door with a window to the left and a patio to the right, walk forward though the door and into a dinning room, ... ..



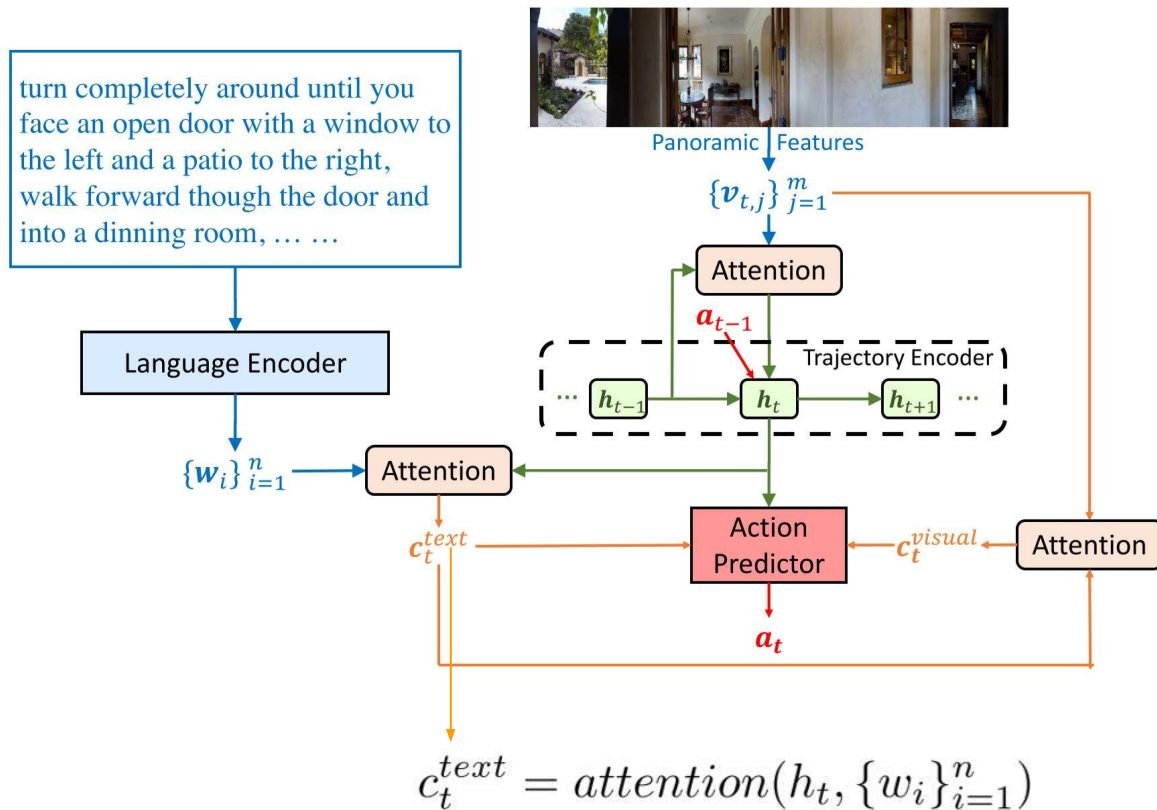
Panoramic Features

$\{v_{t,j}\}_{j=1}^m$

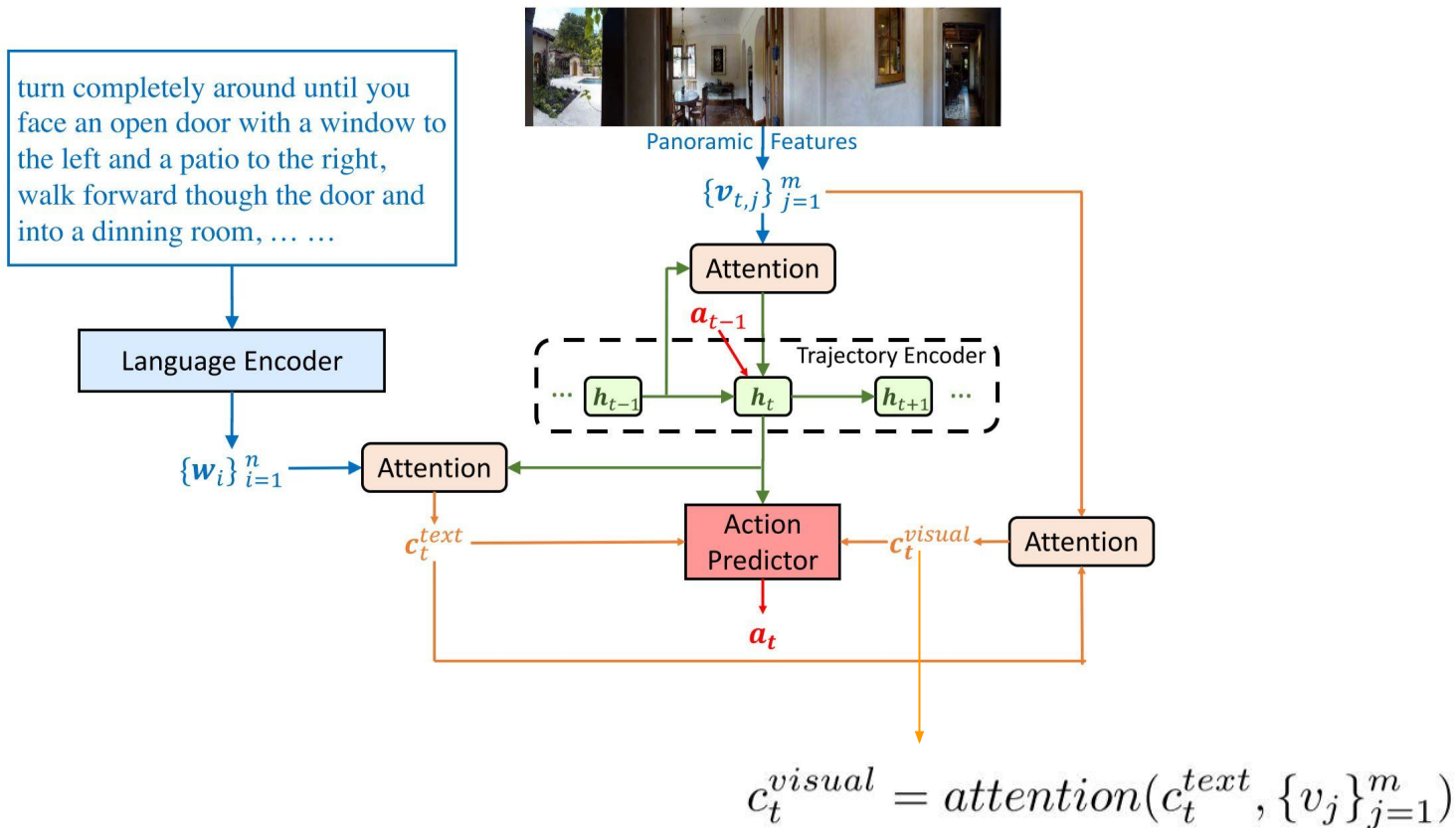


$$v_t = \text{attention}(h_{t-1}, \{v_{t,j}\}_{j=1}^m)$$
$$= \sum_j \text{softmax}(h_{t-1} W_h (v_{t,j} W_v)^T) v_{t,j}$$

# Visually conditioned textual context



# Textually conditioned visual context



# Action prediction

turn completely around until you face an open door with a window to the left and a patio to the right, walk forward through the door and into a dining room, ... ..

Language Encoder

$\{w_i\}_{i=1}^n$

Attention

$c_t^{text}$

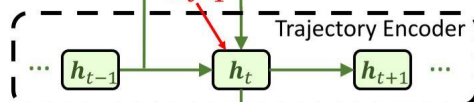


Panoramic Features

$\{v_{t,j}\}_{j=1}^m$

Attention

$a_{t-1}$



Trajectory Encoder

$h_t$

$h_{t+1}$

Action Predictor

$a_t$

$c_t^{visual}$

Attention

$$p_k = \text{softmax}([h_t, c_t^{text}, c_t^{visual}]W_c(u_k W_u)^T)$$

# Action prediction

Each viewpoint corresponds to a next-step reachable location

image patch from  
one viewpoint

heading

elevation

$$[v_i; \sin\psi; \cos\psi; \sin\omega; \cos\omega]$$

action vector

$$p_k = \text{softmax}([h_t, c_t^{\text{text}}, c_t^{\text{visual}}]W_c (u_k W_u)^T)$$



# Action prediction

image patch from  
one viewpoint

heading

elevation

$$[v_i; \sin\psi; \cos\psi; \sin\omega; \cos\omega]$$

action vector

$$p_k = \text{softmax}(\underbrace{[h_t, c_t^{\text{text}}, c_t^{\text{visual}}]}_{\text{Context vector}} W_c (u_k W_u)^T)$$

Context vector

# Action prediction

image patch from  
one viewpoint

heading

elevation

$$[v_i; \sin\psi; \cos\psi; \sin\omega; \cos\omega]$$

action vector

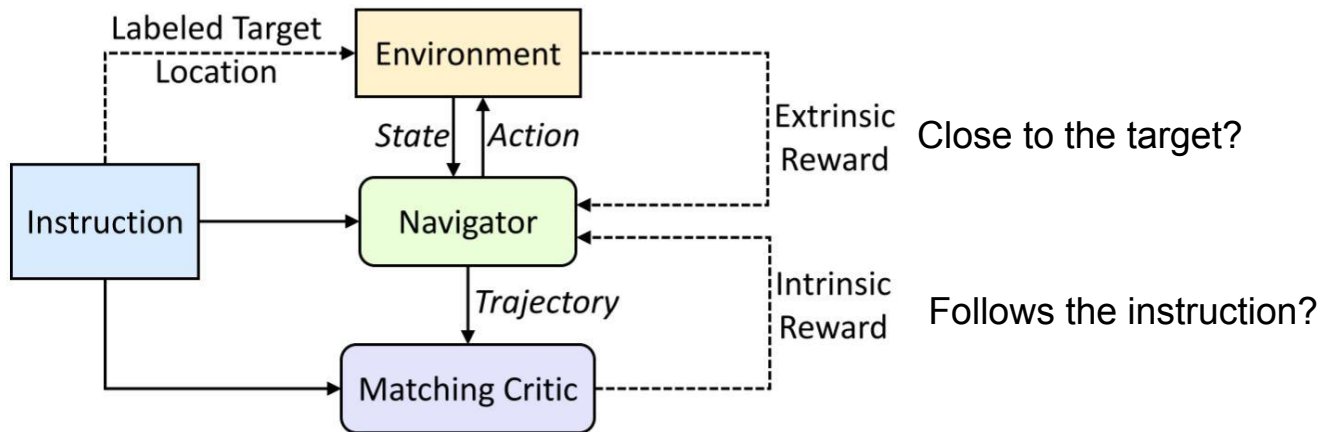
$$p_k = \text{softmax}([h_t, c_t^{\text{text}}, c_t^{\text{visual}}] W_c (u_k W_u)^T)$$

Context vector

Trainable weights

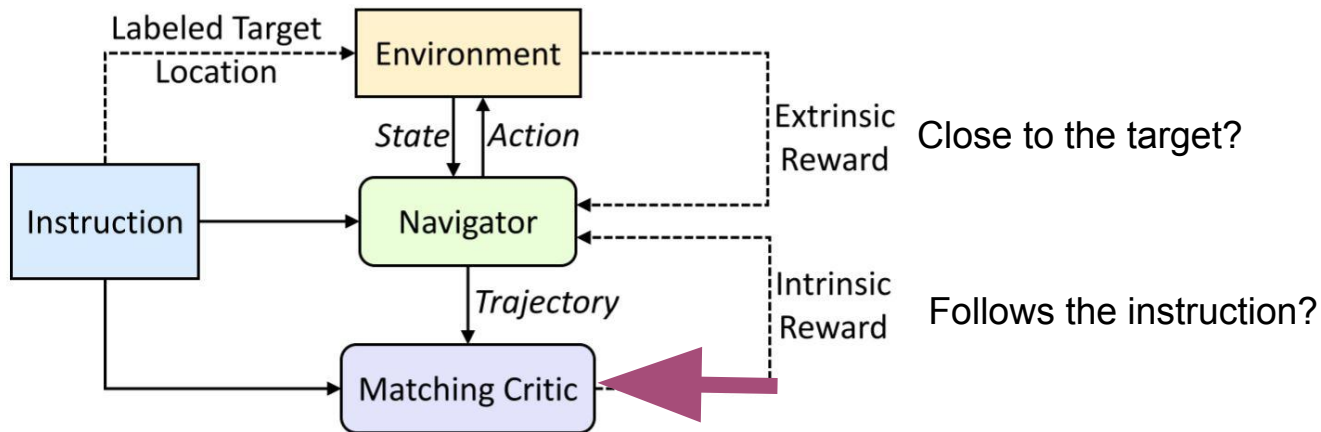


Turn right and head towards the kitchen. Then turn left, pass a table and enter the hallway. ...Stop in front of the toilet.





Turn right and head towards the kitchen. Then turn left, pass a table and enter the hallway. ...Stop in front of the toilet.



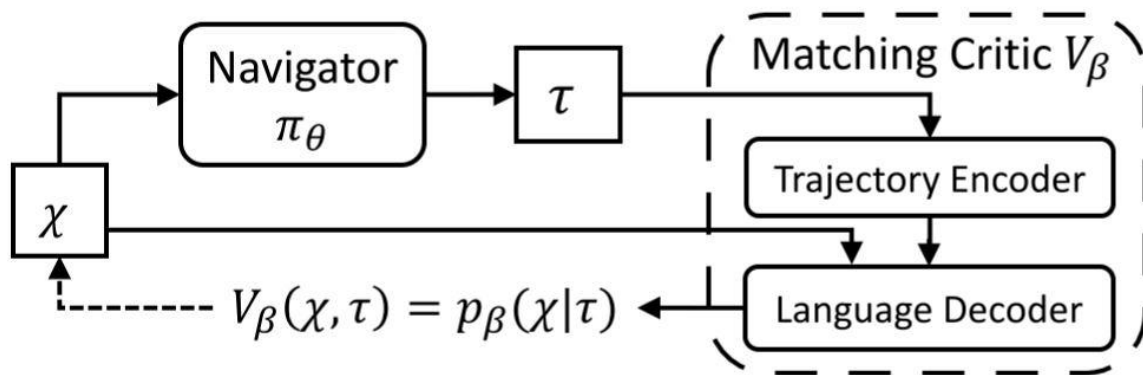
# Matching Critic (intrinsic reward)

Enforces the agent to follow the instruction by rewarding trajectories that match the instruction and to determine matching

# Matching Critic (intrinsic reward)

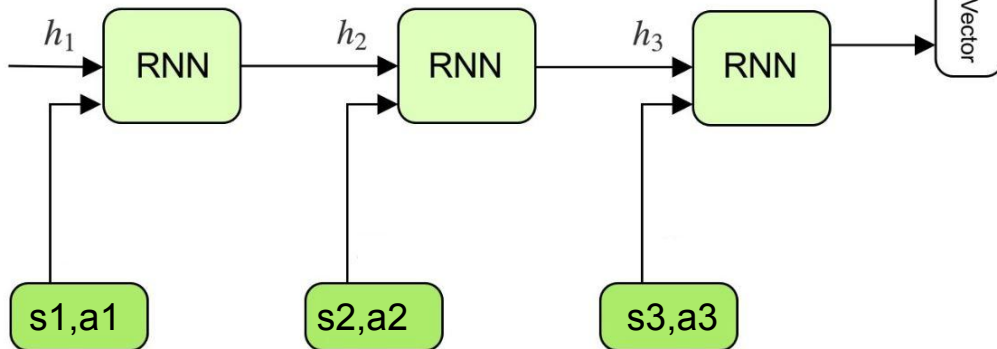
Try to reconstruct the instruction from the trajectory

and see how it matches the original one

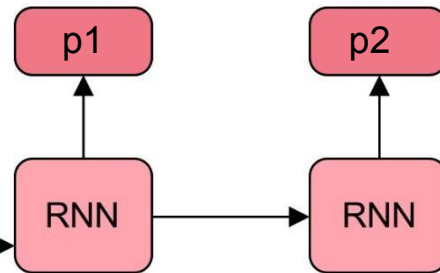


# Sequence to sequence

Encoder

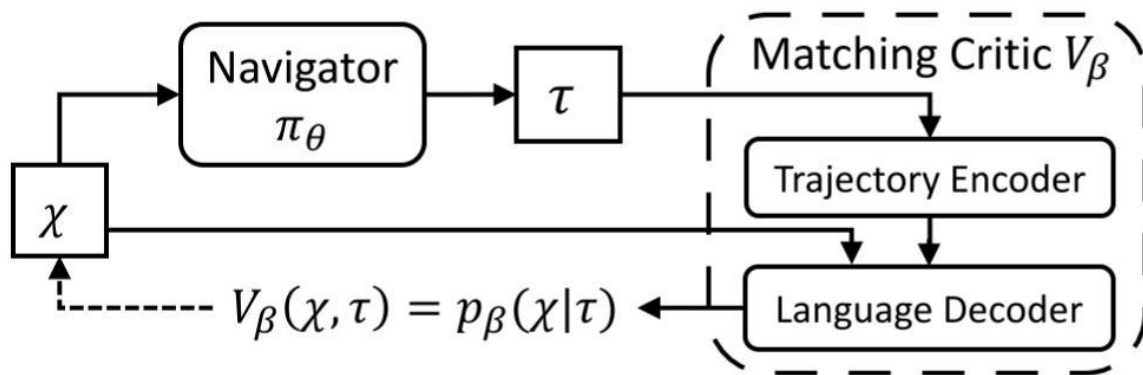


Decoder



# Matching Critic (intrinsic reward)

- Try to reconstruct the instruction from the trajectory and see how it matches the original one





# Learning

- Goal: learn the policy  $\pi: \mathbf{S} \rightarrow \mathbf{A}$
- Approximate the policy first!
  - Behavior cloning
  - Supervised Learning
  - Use state-action data from seen environments
  - $\mathbf{a}_t^*$  - demonstration action

$$L_{sl} = -\mathbb{E}[\log(\pi_{\theta}(a_t^* | s_t))]$$

# Learning

- **Problem:** limited generalizability
  - Unpredictable behavior in unseen states
- **Solution:** switch over to Reinforcement Learning!



# Learning - Extrinsic Reward

- Two metrics:

- **Immediate reward:**

- **Distance** from the current state  $\mathbf{s}_t$  to the target state  $\mathbf{s}_{\text{target}}$ :  $\mathcal{D}_{\text{target}}(\mathbf{s}_t)$

$$r(s_t, a_t) = \mathcal{D}_{\text{target}}(s_t) - \mathcal{D}_{\text{target}}(s_{t+1})$$

- **Delayed reward:**

- When the agent reaches its destination:

$$r(s_T, a_T) = \mathbb{1}(\mathcal{D}_{\text{target}}(s_T) \leq d)$$

# Learning - Extrinsic Reward

$$R_{extr}(s_t, a_t) = \underbrace{r(s_t, a_t)}_{\text{immediate reward}} + \underbrace{\sum_{t'=t+1}^T \gamma^{t'-t} r(s_{t'}, a_{t'})}_{\text{discounted future reward}}$$

# Learning - Extrinsic & Intrinsic Rewards

- Matching critic calculates intrinsic reward  $R_{intr}$
- Based on alignment between the trajectory and the original instruction
- Final Reward function:
  - $L_{rl} = -\mathbb{E}_{a_t \sim \pi_\theta} [A_t]$  where  $A_t = R_{extr} + \delta R_{intr}$
  - $\delta$  is a hyperparameter defining importance of  $R_{intr}$
  - Gradient of the loss function:  $\nabla_\theta L_{rl} = -A_t \nabla_\theta \log \pi_\theta(a_t | s_t)$

# Self-Supervised Imitation Learning

- **Original setting:** train on *seen* environments and test on *unseen* ones
- **New setting:** allow exploration on *unseen* environments as well
  - Lifelong learning & adaptation
  - No ground-truth demonstrations
  - Don't know where the target location is => **no extrinsic reward**

# Self-Supervised Imitation Learning

1. Get an instruction  $\mathcal{X}$
2. Navigator produces possible trajectories
3. Matching Critic  $\mathbf{V}_\beta$  finds the best trajectory:  $\hat{\tau} = \arg \max_{\tau} V_\beta(\mathcal{X}, \tau)$
4. Navigator stores it into a replay buffer
5. The agent uses the saved trajectory  $\hat{\tau}$  as a ground-truth
6. Self-supervised learning loss:  $L_{sil} = -\mathbb{E}[\log(\pi_\theta(\hat{a}_t | s_t))]$

# Experimental setup:

- **R2R dataset:** This dataset covers most of the visual diversity and instructions. It consists of training, seen validation, unseen validation and test sets.

- **Testing scenarios:**

  - **Standard:** Train in seen environments and test in unseen environments without prior exploration.

  - **Lifelong learning:** The agent is encouraged to explore the environment and learn from the feedback.



# Evaluation metrics:

- **PL**: Path length
- **NE**: Navigation error
- **OSR**: Oracle success rate
- **SR**: Success rate
- **SPL**: Success rate weighted by inverse Path Length


<b>Test Set (VLN Challenge Leaderboard)</b>					
Model	PL ↓	NE ↓	OSR ↑	SR ↑	SPL ↑
Random	9.89	9.79	18.3	13.2	12
seq2seq [2]	<b>8.13</b>	7.85	26.6	20.4	18
RPA [48]	9.15	7.53	32.5	25.3	23
Speaker-Follower [11]	14.82	6.62	44.0	<b>35.0</b>	<b>28</b>
+ beam search	<u>1257.38</u>	4.87	96.0	53.5	<u>1</u>
<b>Ours</b>					
RCM	<b>15.22</b>	<b>6.01</b>	<b>50.8</b>	<b>43.1</b>	<b>35</b>
 RCM + SIL (train)	<b>11.97</b>	6.12	49.5	43.0	<b>38</b>
RCM + SIL (unseen) <sup>3</sup>	<b>9.48</b>	4.21	66.8	60.5	59

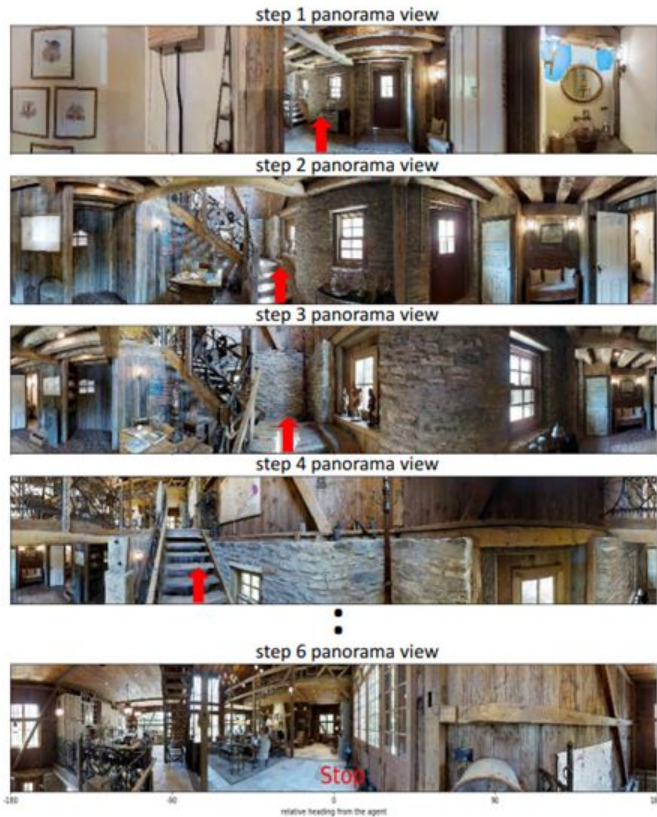
Table 1: Comparison between RCM, SIL and the state-of-the-art

#	Model	Seen Validation				Unseen Validation			
		PL ↓	NE ↓	OSR ↑	SR ↑	PL ↓	NE ↓	OSR ↑	SR ↑
0	Speaker-Follower (no beam search) [11]	-	3.36	73.8	66.4	-	6.62	45.0	35.5
1	RCM + SIL (train)	<b>10.65</b>	3.53	75.0	66.7	<b>11.46</b>	6.09	50.1	<b>42.8</b>
2	RCM	11.92	3.37	76.6	67.4	14.84	<b>5.88</b>	<b>51.9</b>	42.5
3	– intrinsic reward	12.08	3.25	<b>77.2</b>	67.6	15.00	6.02	50.5	40.6
4	– extrinsic reward = pure SL	11.99	3.22	76.7	66.9	14.83	6.29	46.5	37.7
5	– cross-modal reasoning	11.88	<b>3.18</b>	73.9	66.4	14.51	6.47	44.8	35.7
6	RCM + SIL (unseen)	<b>10.13</b>	<b>2.78</b>	<b>79.7</b>	73.0	<b>9.12</b>	<b>4.17</b>	<b>69.31</b>	61.3

Table 2: Ablation study on seen and unseen environments

**Instruction:** Exit the door and turn left towards the staircase. Walk all the way up the stairs, and stop at the top of the stairs.

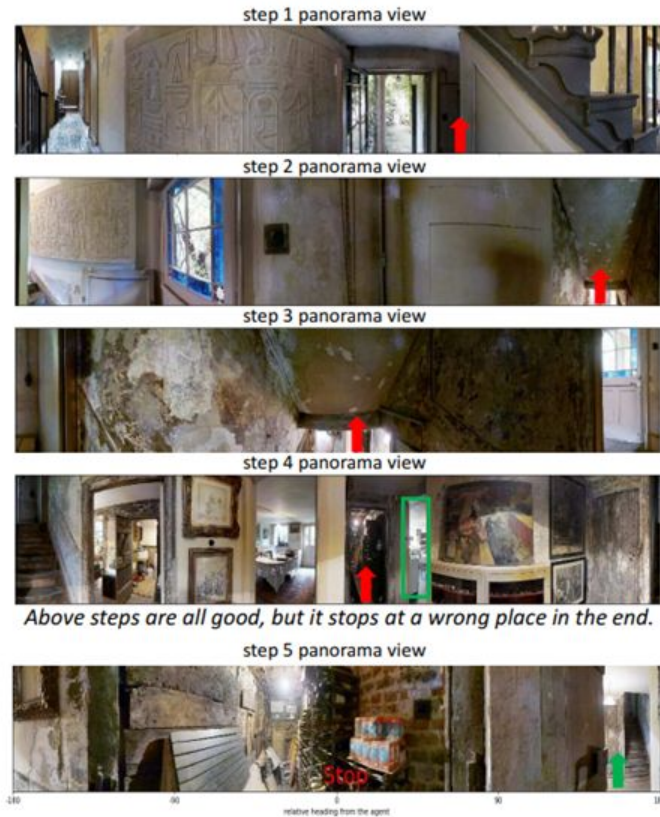
**Intrinsic Reward: 0.53** **Result: Success (error = 0m)**



(a) A successful case

**Instruction:** Turn right and go down the stairs. Turn left and go straight until you get to *the laundry room*. Wait there.

**Intrinsic Reward: 0.54** **Result: Failure (error = 5.5m)**



*Above steps are all good, but it stops at a wrong place in the end.*

(b) A failure case

Figure 6: Qualitative examples from the unseen validation set.

# Strength and weakness:

## Strengths:

- Reinforcement learning leads to better generalizability
- Self-Imitation learning facilitates exploring on unseen environments
- Cross modal grounding effectively enhances the model's ability to capture context information

# Strength and weakness:

## Weakness:

- Limited in dataset diversity (Only on R2R)
- Model implementation heavily depends on simulator specifics
- Not using depth data

# Conclusion:

This paper presents two novel approaches, RCM and SIL, which combine reinforcement learning and self-supervised imitation learning.

The experiment results show that these methods are effective and efficient in both the standard testing scenario and lifelong learning scenario.

SIL help achieve strong generalizability in unseen environments.

Questions

