

# Language-conditioned Tasks

# Language-conditioned Tasks

- Language is a natural way to communicate with machines

# Language-conditioned Tasks

- Language is a natural way to communicate with machines
- Interpret and execute the NL instructions expressing some intents

# Language-conditioned Tasks

- Language is a natural way to communicate with machines
- Interpret and execute the NL instructions expressing some intents
  - NL2Code

# Language-conditioned Tasks

- Language is a natural way to communicate with machines
- Interpret and execute the NL instructions expressing some intents
  - NL2Code
  - NL2Action+Argument (instruction following)

# Room-to-room dataset

# Room-to-room dataset

**Task:** Given some natural language instructions, navigate the agent through the environment and reach the goal location

# Room-to-room dataset

**Task:** Given some natural language instructions, navigate the agent through the environment and reach the goal location

**Actions:** *left, right, up, down, forward, stop*



# Room-to-room dataset

**Task:** Given some natural language instructions, navigate the agent through the environment and reach the goal location

**Actions:** *left, right, up, down, forward, stop*



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

# Room-to-room dataset

**Task:** Given some natural language instructions, navigate the agent through the environment and reach the goal location

**Actions:** *left, right, up, down, forward, stop*



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

# Room-to-room dataset

**Task:** Given some natural language instructions, navigate the agent through the environment and reach the goal location

**Actions:** *left, right, up, down, forward, stop*



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

# Room-to-room dataset

**Task:** Given some natural language instructions, navigate the agent through the environment and reach the goal location

**Actions:** *left, right, up, down, forward, stop*



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

# ALFRED datase

# ALFRED datase

**Task:** Given some natural language instructions, navigate through the environment, interact with the objects to achieve some goal states

# ALFRED datase

**Task:** Given some natural language instructions, navigate through the environment, interact with the objects to achieve some goal states

**Actions:** *MoveAhead, RotateLeft(Right), LookUp(Down), PickupObject, PutObject, OpenObject, CloseObject, ToggleOn(Off)*

# ALFRED datase

**Task:** Given some natural language instructions, navigate through the environment, interact with the objects to achieve some goal states

**Actions:** *MoveAhead, RotateLeft(Right), LookUp(Down), PickupObject, PutObject, OpenObject, CloseObject, ToggleOn(Off)*





# ALFRED datase

**Task:** Given some natural language instructions, navigate through the environment, interact with the objects to achieve some goal states

**Actions:** *MoveAhead, RotateLeft(Right), LookUp(Down), PickupObject, PutObject, OpenObject, CloseObject, ToggleOn(Off)*

## Goal Instruction

Put a microwaved tomato in the sink.

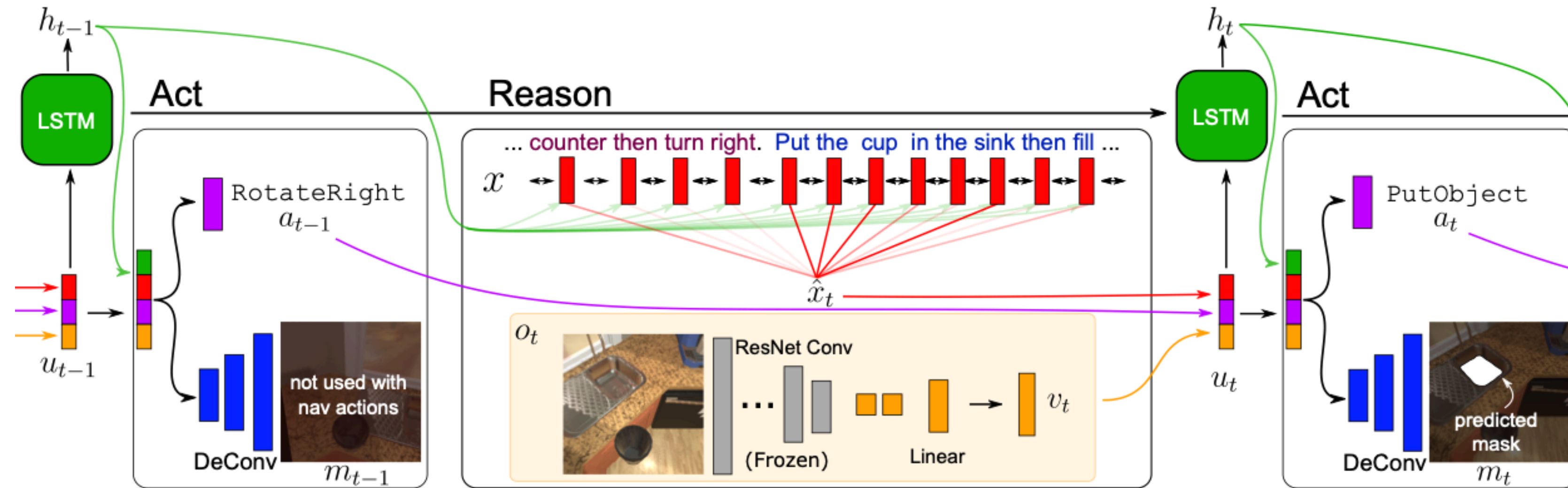
## Step-by-Step Instructions

Turn around and go to the left side of the sink.  
Pick up the tomato in the front.  
Turn right to go to the microwave on the left.  
Microwave the tomato next to coffee mug and take it out.  
Turn left to go back to the sink.  
Place the tomato inside the sink.

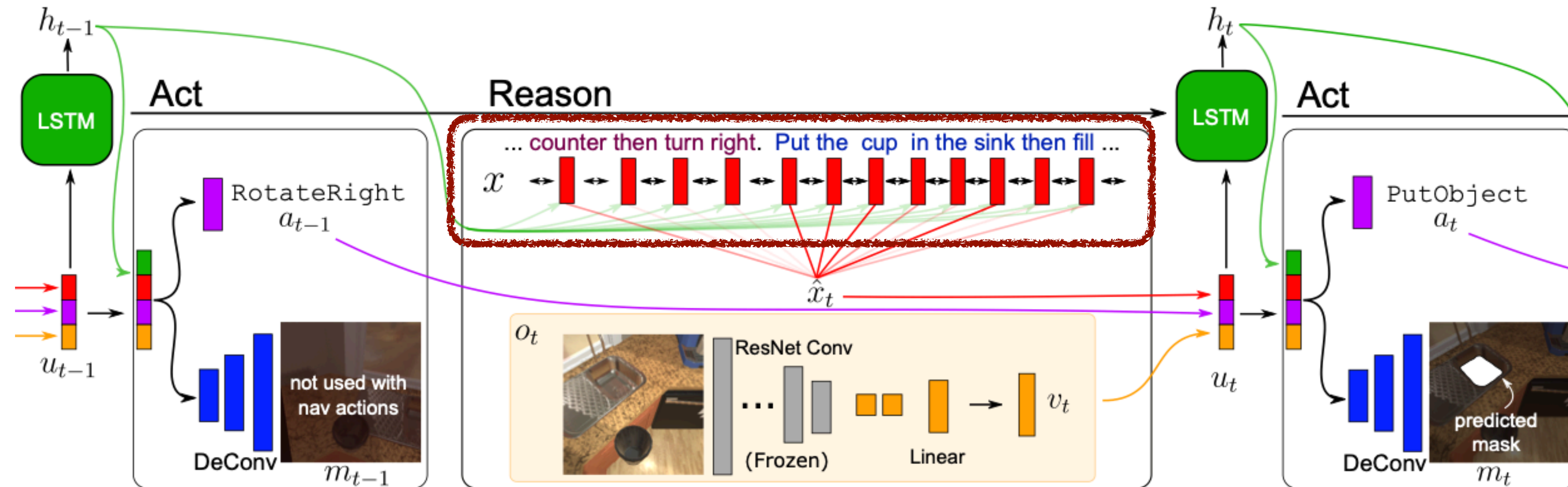


# Sequence to sequence model w/ Attention

# Sequence to sequence model w/ Attention

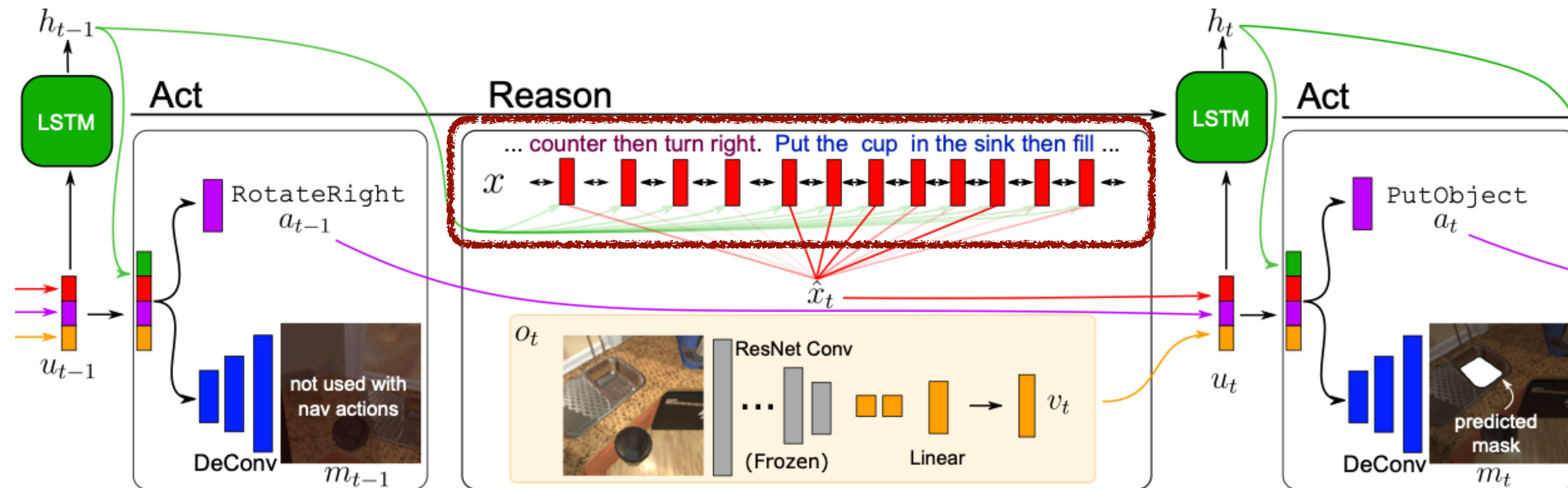


# Sequence to sequence model w/ Attention



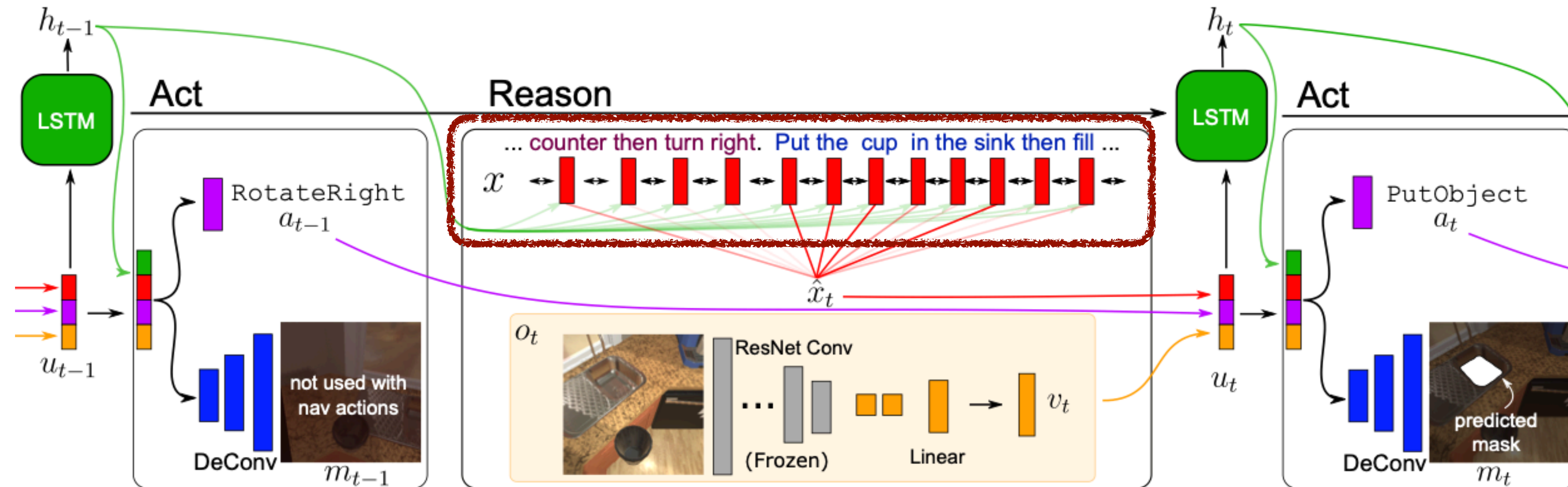
- The encoder encodes the NL instruction

# Sequence to sequence model w/ Attention



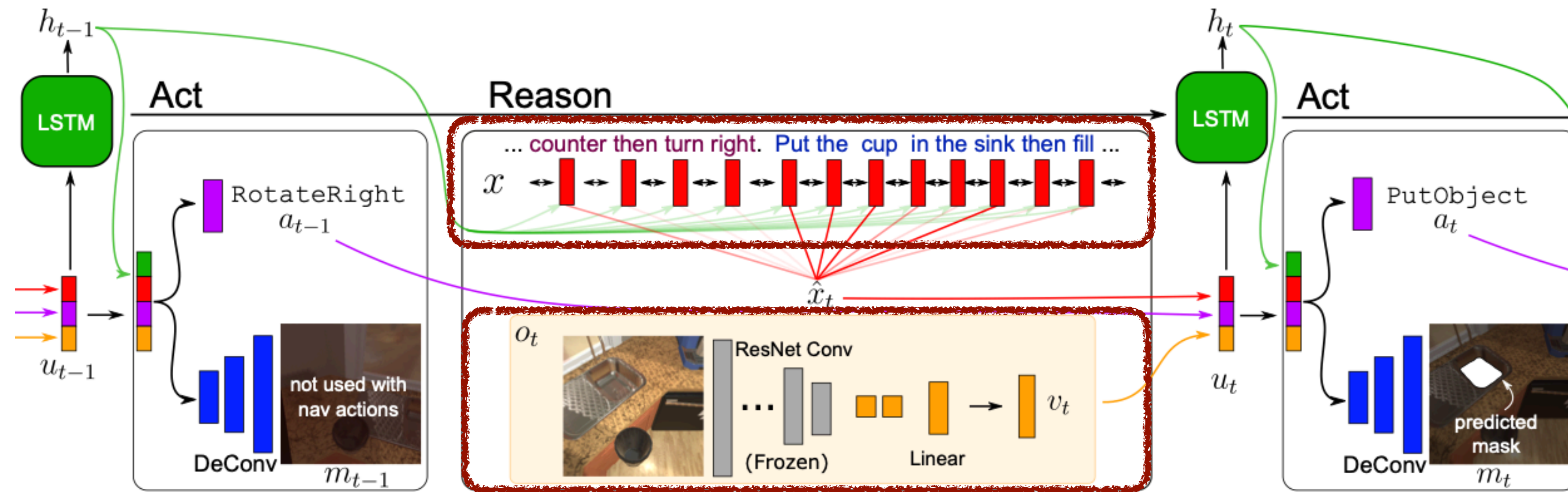
- The encoder encodes the NL instruction
- The decoder predicts one action and its corresponding arguments conditioned on

# Sequence to sequence model w/ Attention



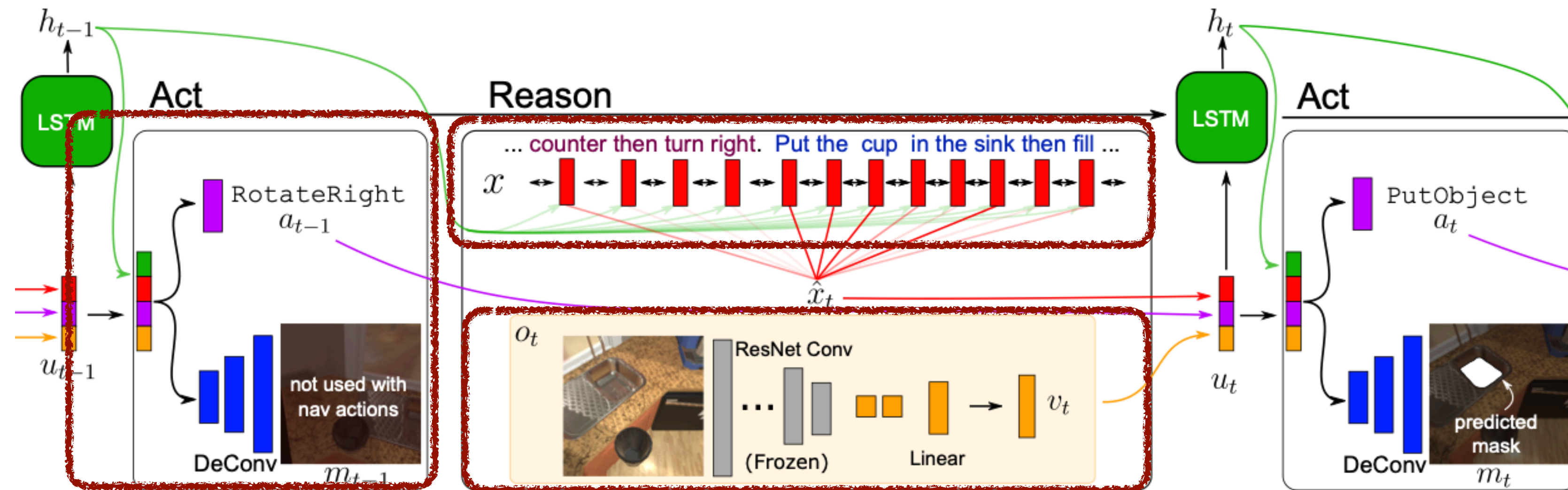
- The encoder encodes the NL instruction
- The decoder predicts one action and its corresponding arguments conditioned on
  - Weighted natural language

# Sequence to sequence model w/ Attention



- The encoder encodes the NL instruction
- The decoder predicts one action and its corresponding arguments conditioned on
  - Weighted natural language
  - Current state (e.g. visual information)

# Sequence to sequence model w/ Attention



- The encoder encodes the NL instruction
- The decoder predicts one action and its corresponding arguments conditioned on
  - Weighted natural language
  - Current state (e.g. visual information)
  - Action from the last step



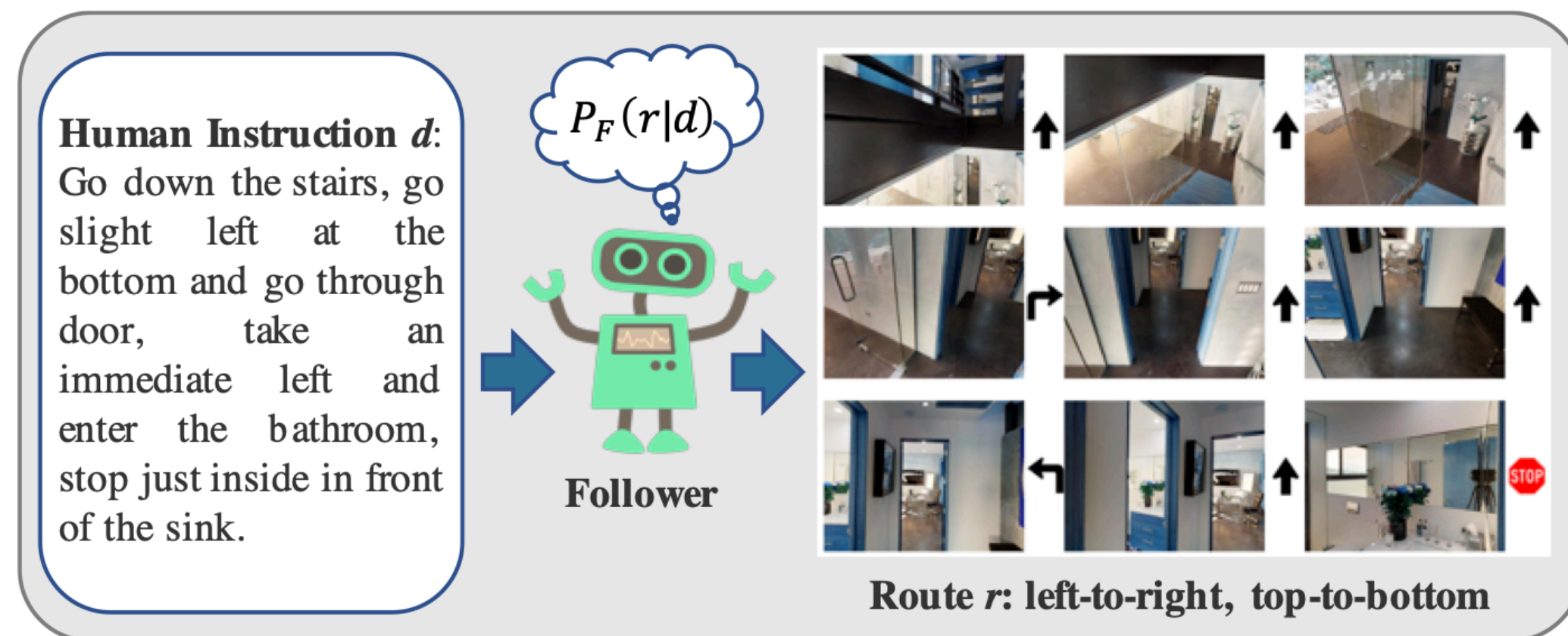
# Speak-follower Model

# Speak-follower Model

Language is often under-specified, it is challenging to learn the mapping between NL and actions from a limited amount of annotations

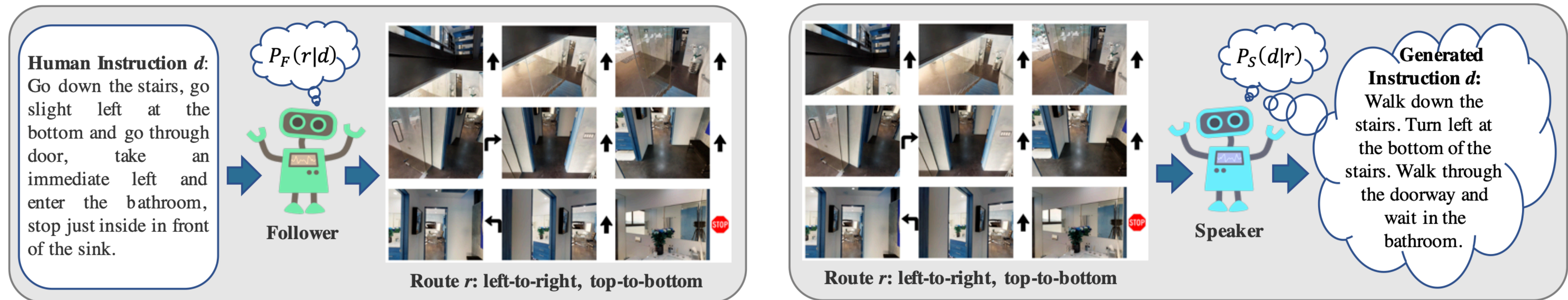
# Speak-follower Model

Language is often under-specified, it is challenging to learn the mapping between NL and actions from a limited amount of annotations



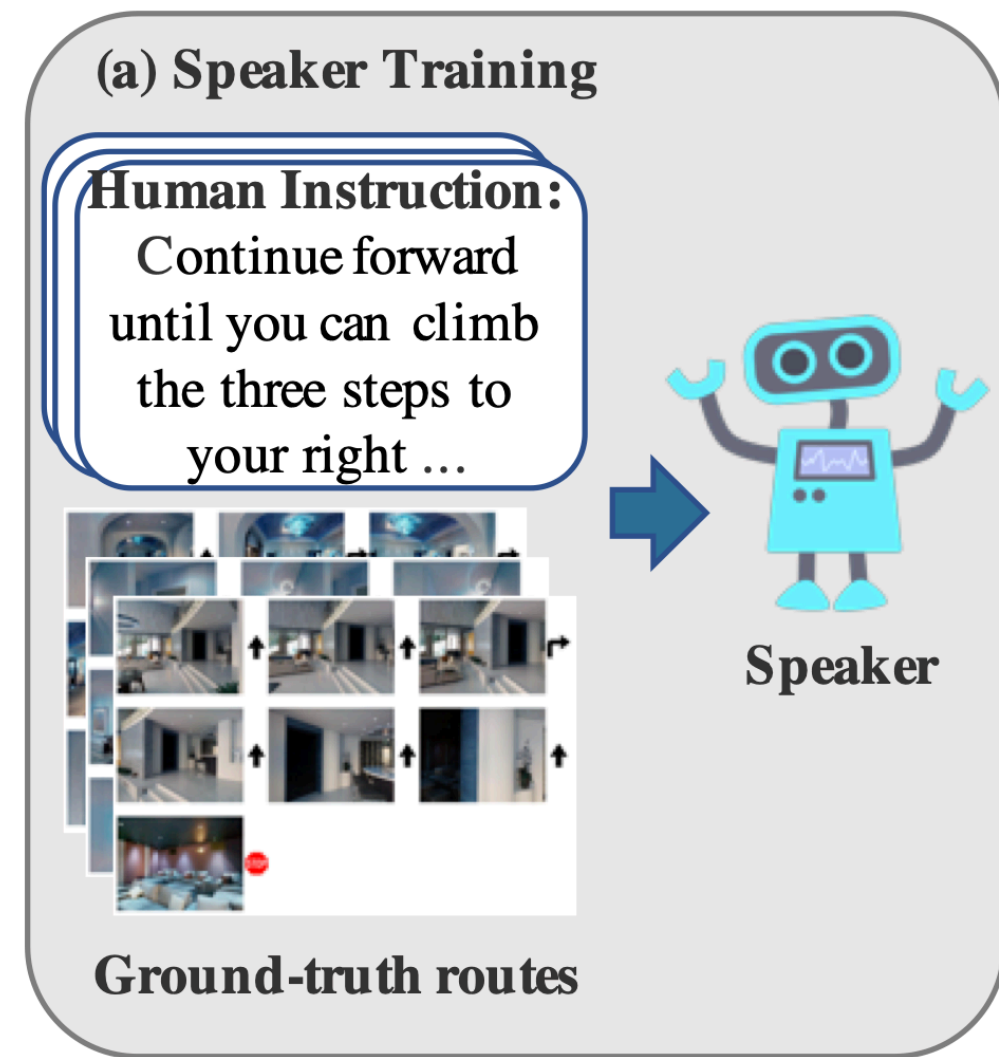
# Speak-follower Model

Language is often under-specified, it is challenging to learn the mapping between NL and actions from a limited amount of annotations

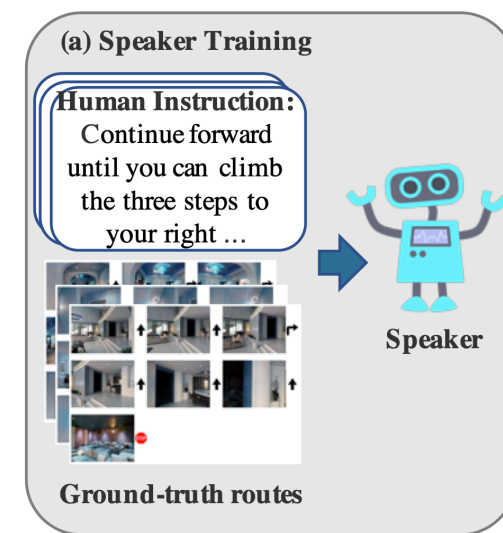


# Speak-follower Model

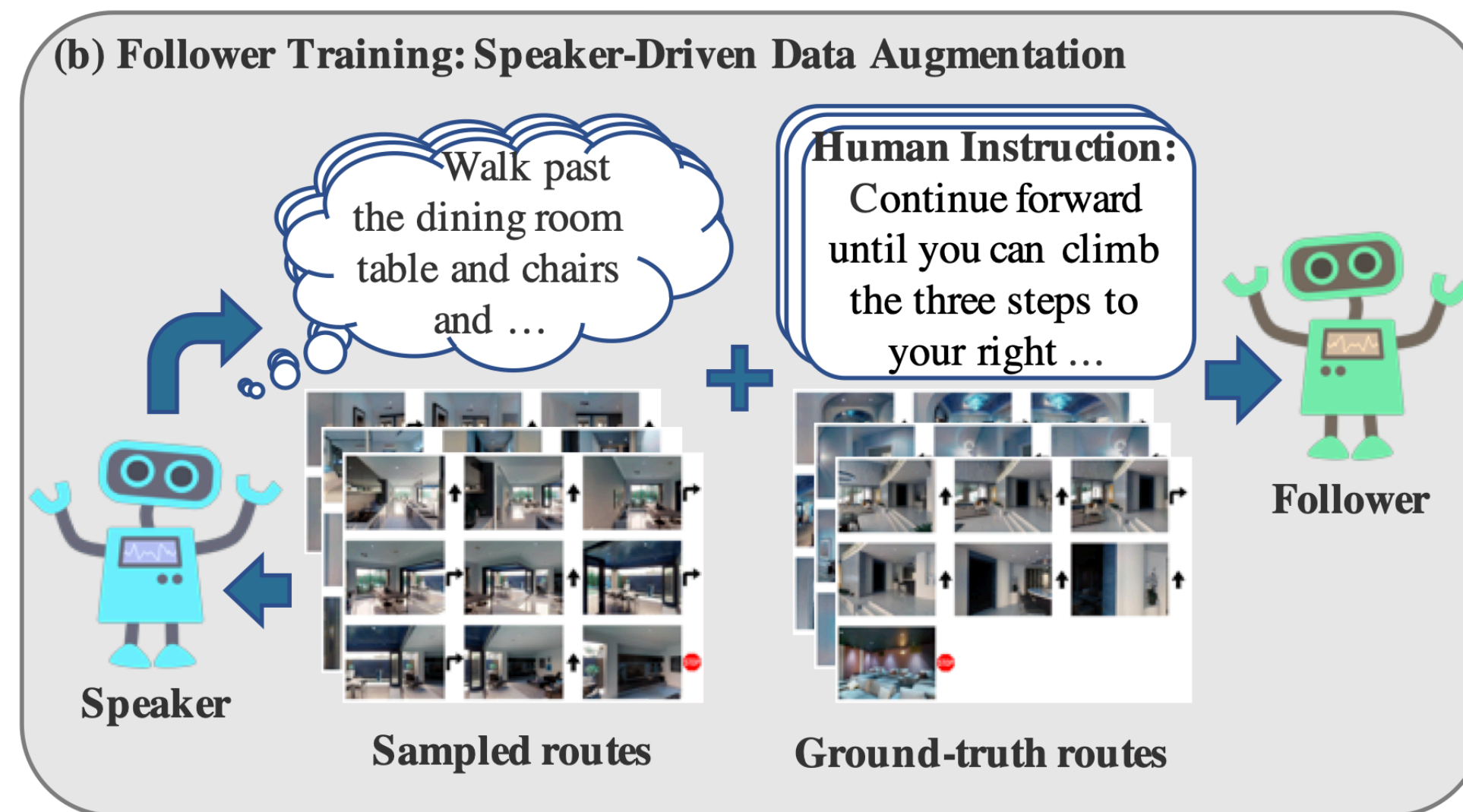
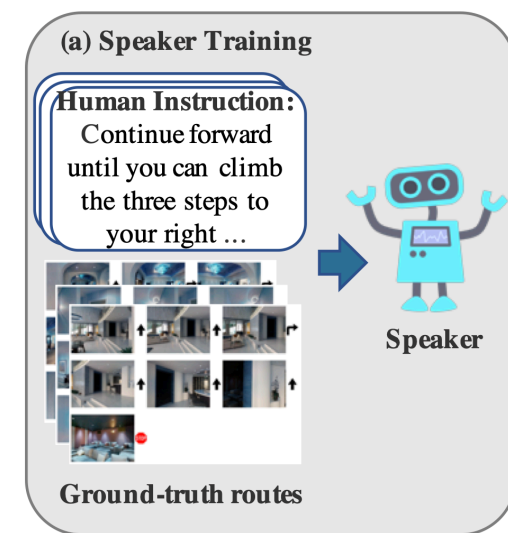
# Speak-follower Model



# Speak-follower Model

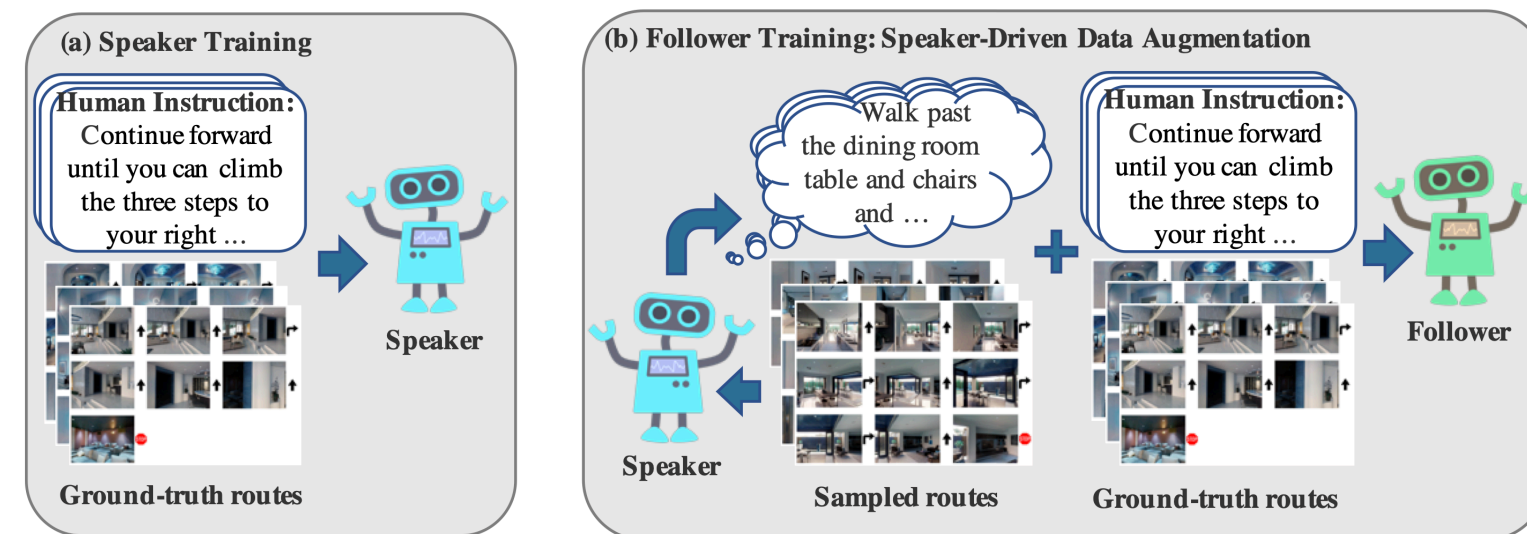


# Speak-follower Model

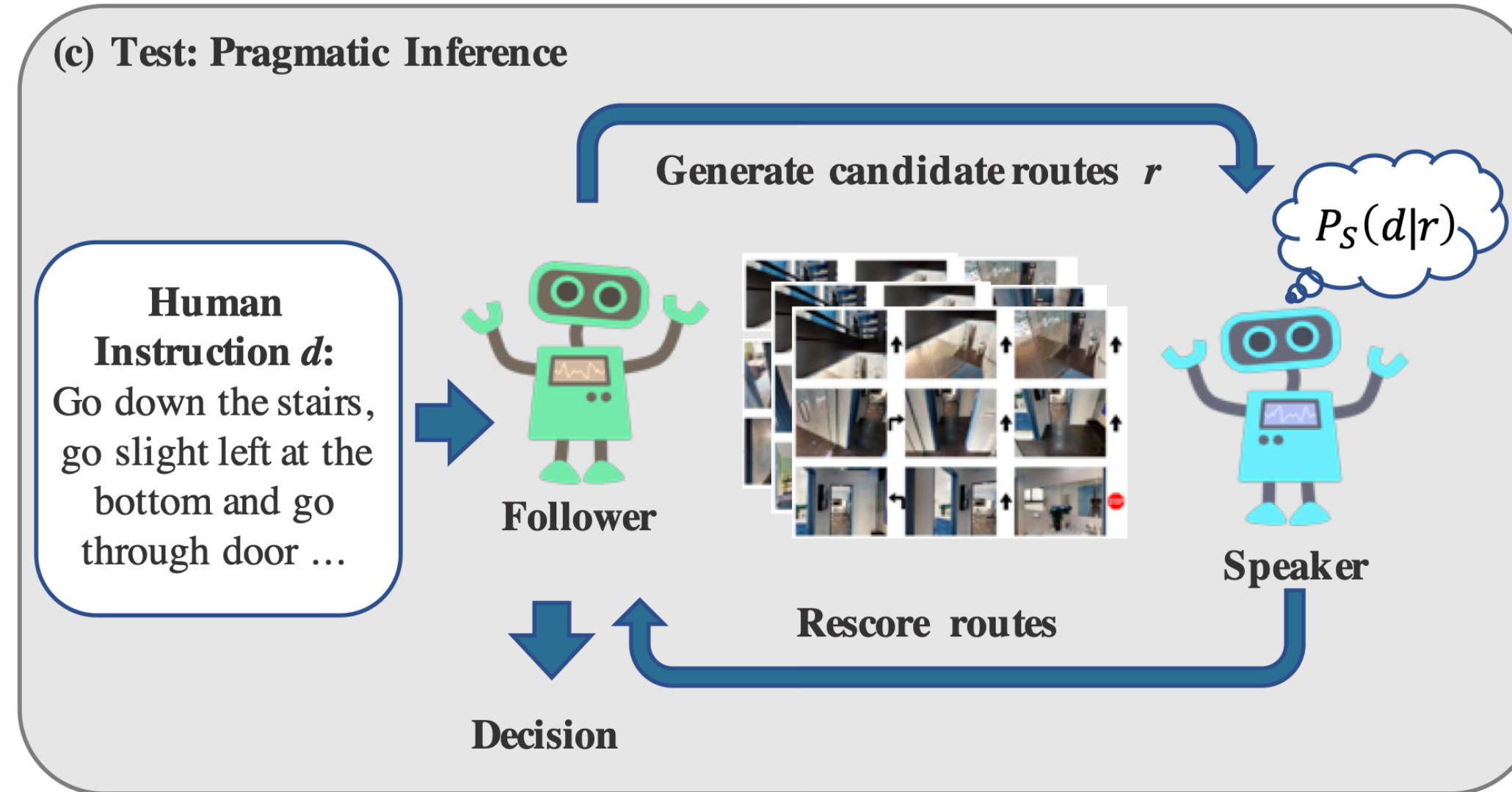
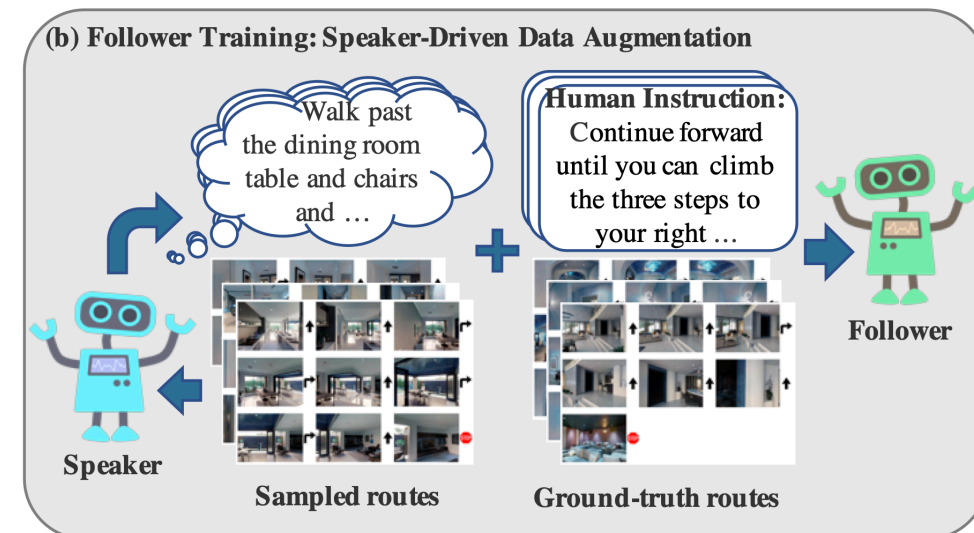
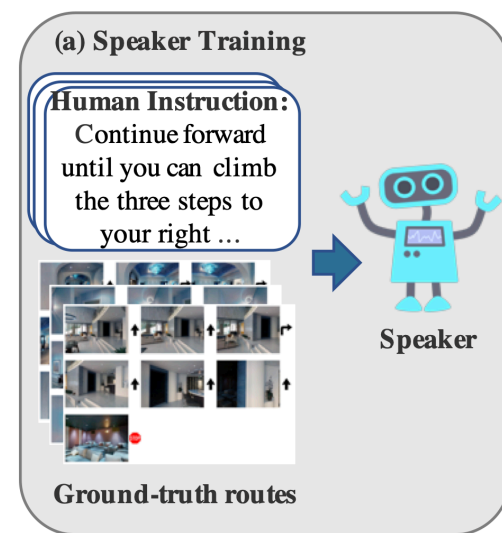




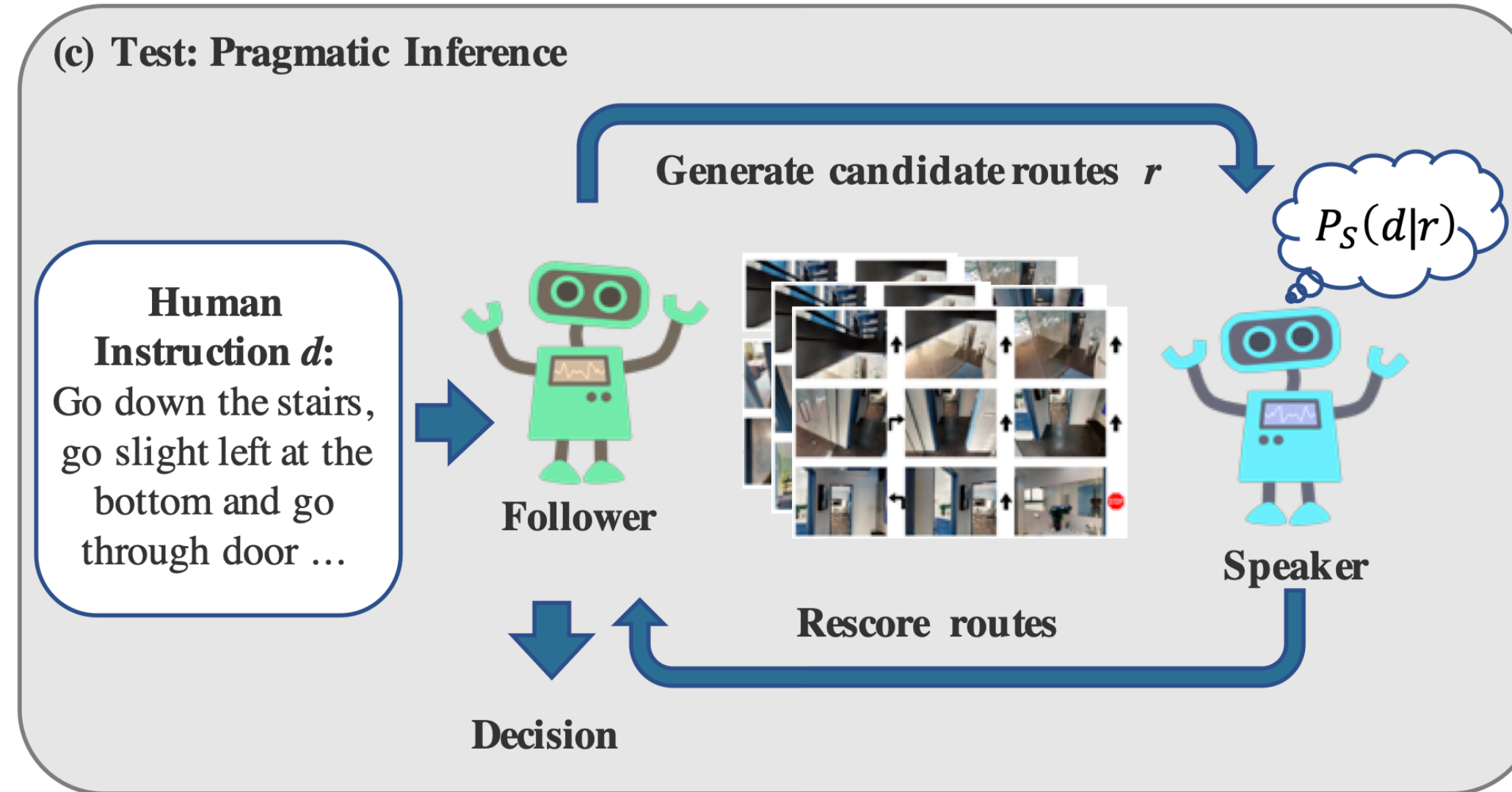
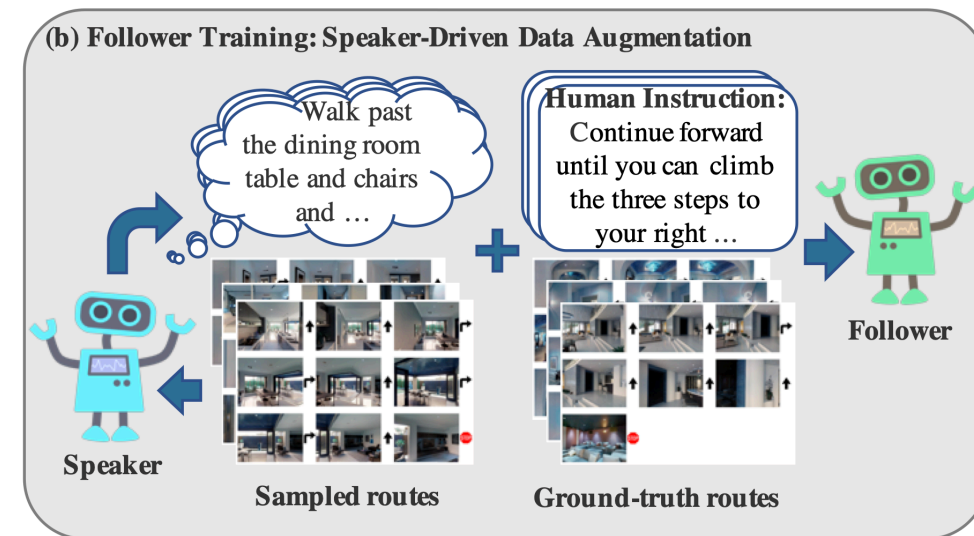
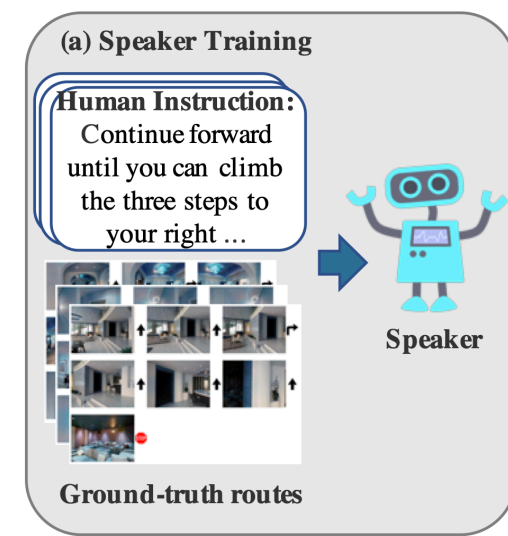
# Speak-follower Model



# Speak-follower Model

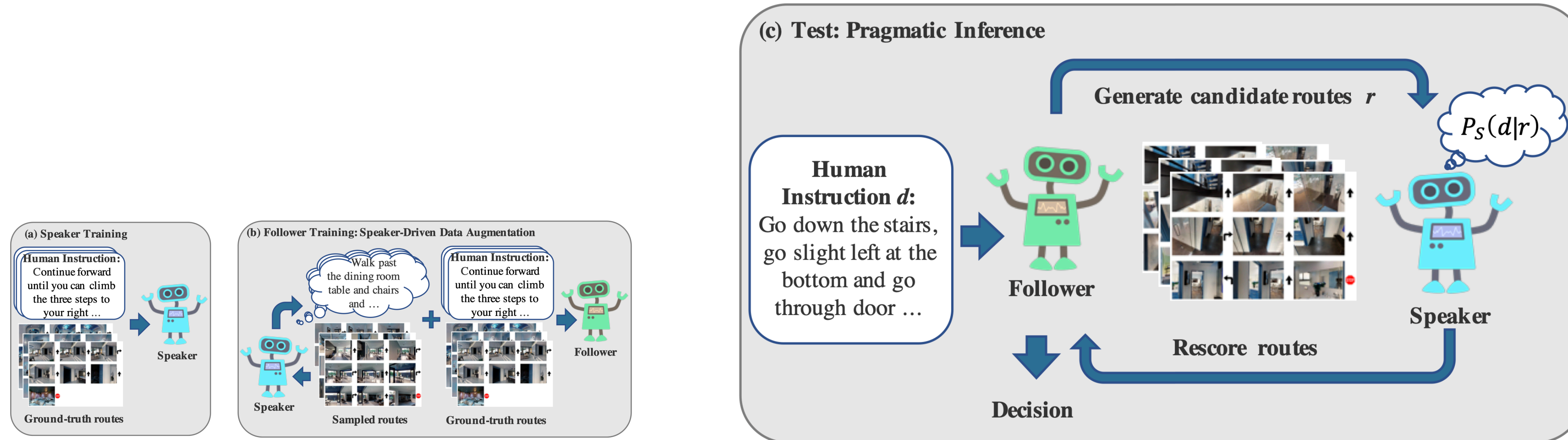


# Speak-follower Model



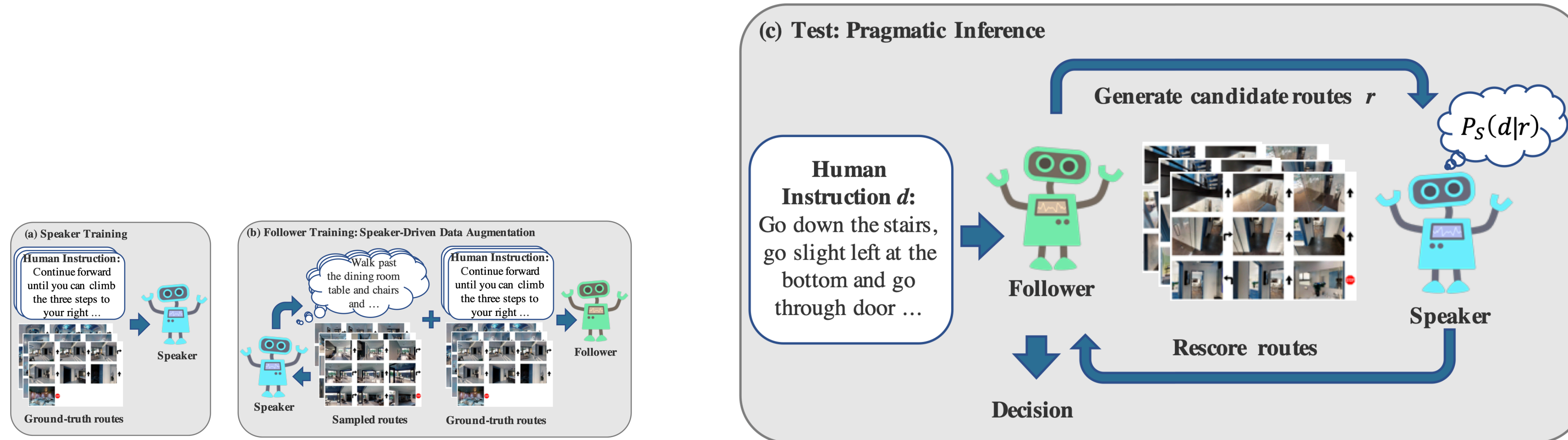
- Data augmentation

# Speak-follower Model



- Data augmentation
- Consistency enforcement between the natural language and actions

# Speak-follower Model



- Data augmentation
- Consistency enforcement between the natural language and actions
- Search can be expensive in real-world deployments

# Evaluation

# Evaluation

- SPL (Success weighted by path/action sequence length)

# Evaluation

- SPL (Success weighted by path/action sequence length)
- CLS (Coverage weighted by length score)
  - Measure the consistency with the reference path



# Evaluation

- SPL (Success weighted by path/action sequence length)
- CLS (Coverage weighted by length score)
  - Measure the consistency with the reference path
- Task (sub-task) success rate
  - Compare state' with state\*

# Other Tasks

- Mobile phone operations
  - Pixel Help (Li et al + 2020)
- Web application operations
  - Russ (Xu et al + 2020)
- Windows/Linux system operations
  - UbuntuWorld (Chakraborti et al + 2016)
  - Windows (Branavan et al + 2010)

# Language-assisted Tasks

# Language-assisted Tasks

- The task can be completed without text

# Language-assisted Tasks

- The task can be completed without text
- Text can provides informative information to assist learning and task completion

# Language-assisted Tasks

- The task can be completed without text
- Text can provides informative information to assist learning and task completion



# Language-assisted Tasks

- The task can be completed without text
- Text can provides informative information to assist learning and task completion

Rules



# Language-assisted Tasks

- The task can be completed without text
- Text can provides informative information to assist learning and task completion





# Language-assisted Tasks

- The task can be completed without text
- Text can provides informative information to assist learning and task completion

👍 Immediate interpretable



# Language-assisted Tasks

- The task can be completed without text
- Text can provides informative information to assist learning and task completion

👍 Immediate interpretable

✗ Lack of flexibility

Rules



Rewards

# Language-assisted Tasks

- The task can be completed without text
- Text can provides informative information to assist learning and task completion

👍 Immediate interpretable

✗ Lack of flexibility

✗ Expensive human efforts

Rules



Rewards

# Language-assisted Tasks

- The task can be completed without text
- Text can provides informative information to assist learning and task completion

👍 Immediate interpretable

👍 Flexible

✗ Lack of flexibility

✗ Expensive human efforts

Rules

Rewards



# Language-assisted Tasks

- The task can be completed without text
- Text can provides informative information to assist learning and task completion

👍 Immediate interpretable

✗ Lack of flexibility

✗ Expensive human efforts

👍 Flexible

👍 Light human efforts

Rules



Rewards

# Language-assisted Tasks

- The task can be completed without text
- Text can provides informative information to assist learning and task completion

👍 Immediate interpretable

❌ Lack of flexibility

❌ Expensive human efforts

👍 Flexible

👍 Light human efforts

❌ Sparse and abstract learning signals

Rules

Rewards



# Language-assisted Tasks

- The task can be completed without text
- Text can provides informative information to assist learning and task completion

👍 Immediate interpretable

❌ Lack of flexibility

❌ Expensive human efforts

👍 Flexible

👍 Light human efforts

❌ Sparse and abstract learning signals

Rules

Natural language

Rewards



# Language-assisted Tasks

- The task can be completed without text
- Text can provides informative information to assist learning and task completion

👍 Immediate interpretable

❌ Lack of flexibility

❌ Expensive human efforts

👍 Flexible

👍 Light human efforts

❌ Sparse and abstract learning signals

Rules

Natural language

Rewards

😊 Flexible



# Language-assisted Tasks

- The task can be completed without text
- Text can provides informative information to assist learning and task completion

👍 Immediate interpretable

❌ Lack of flexibility

❌ Expensive human efforts

👍 Flexible

👍 Light human efforts

❌ Sparse and abstract learning signals

Rules

Natural language

Rewards

😊 Flexible

😊 Manageable human efforts

# Language-assisted Tasks

- The task can be completed without text
- Text can provides informative information to assist learning and task completion

👍 Immediate interpretable

❌ Lack of flexibility

❌ Expensive human efforts

👍 Flexible

👍 Light human efforts

❌ Sparse and abstract learning signals

Rules

Natural language

Rewards

😊 Flexible

😊 Manageable human efforts

😊 Richer learning signals

# Communicate the Structured Policies

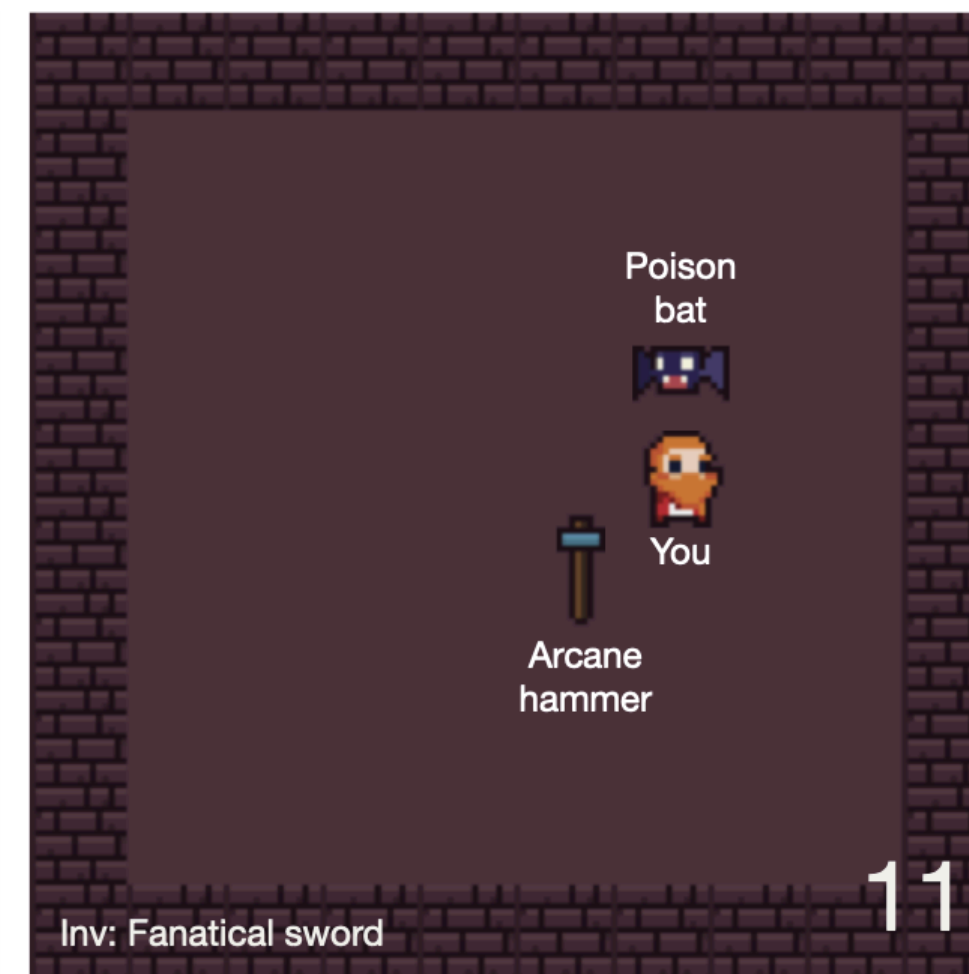
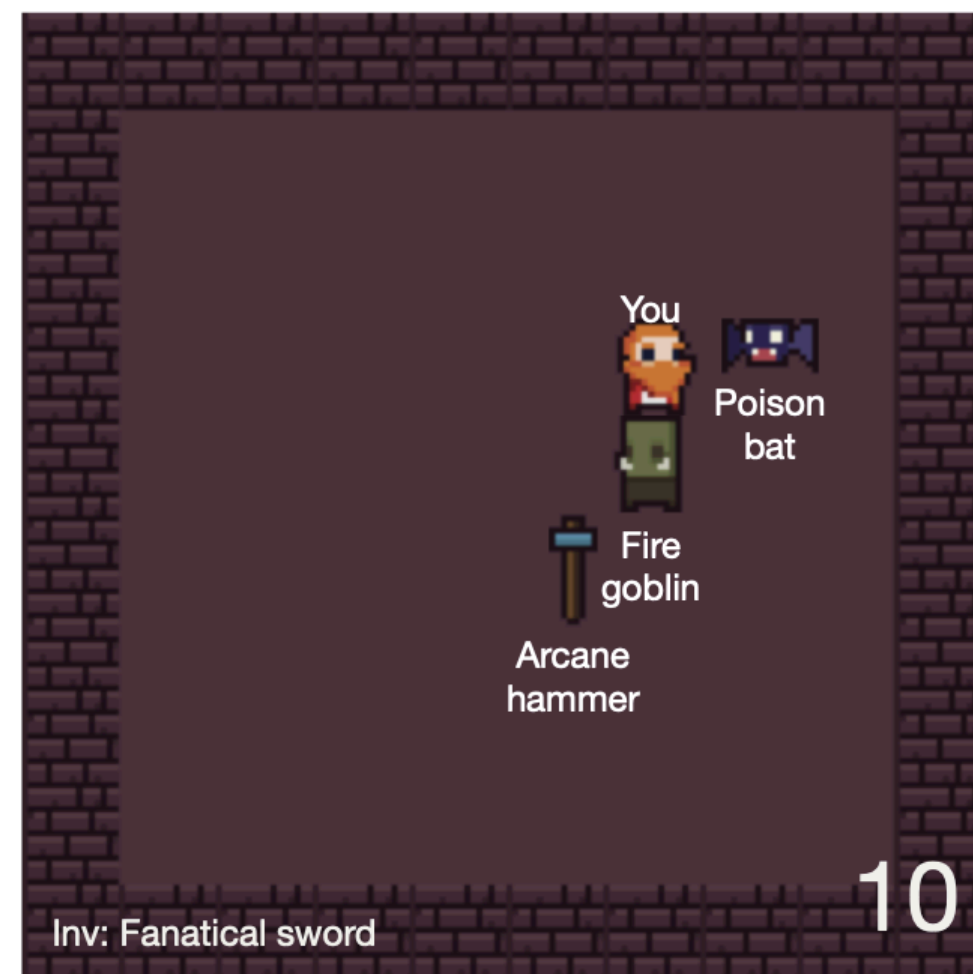
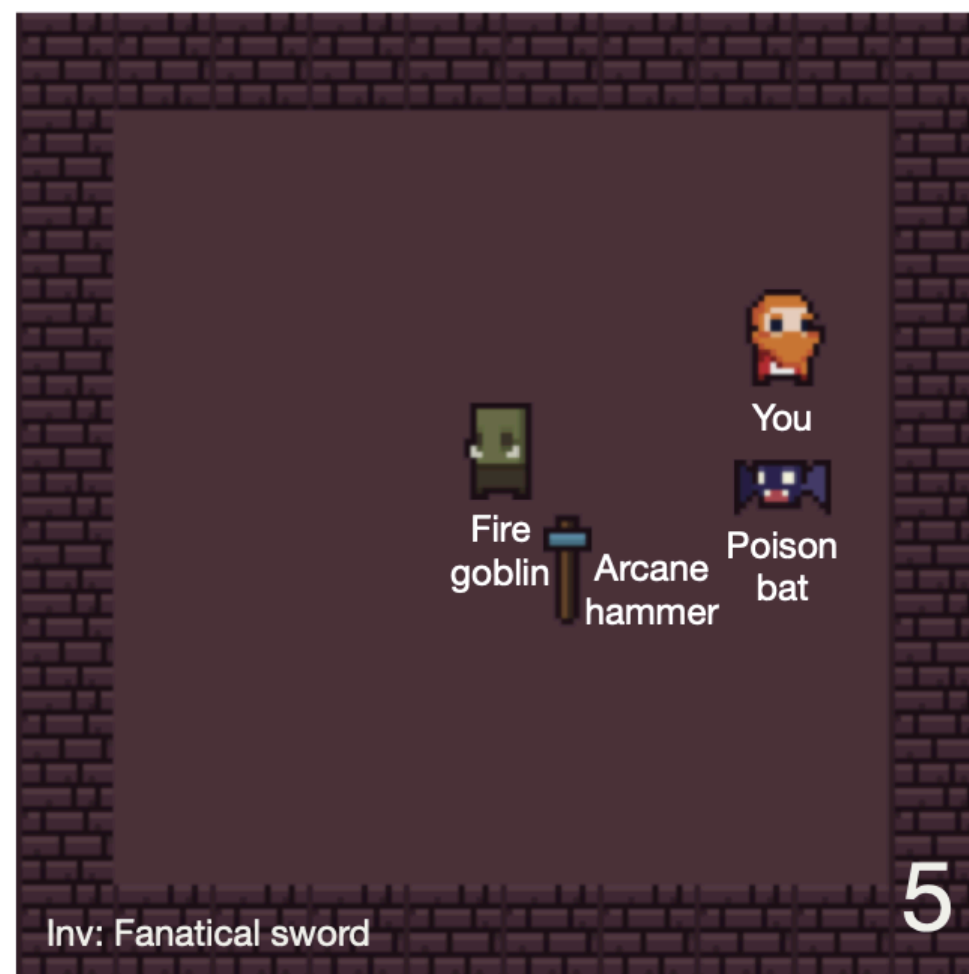
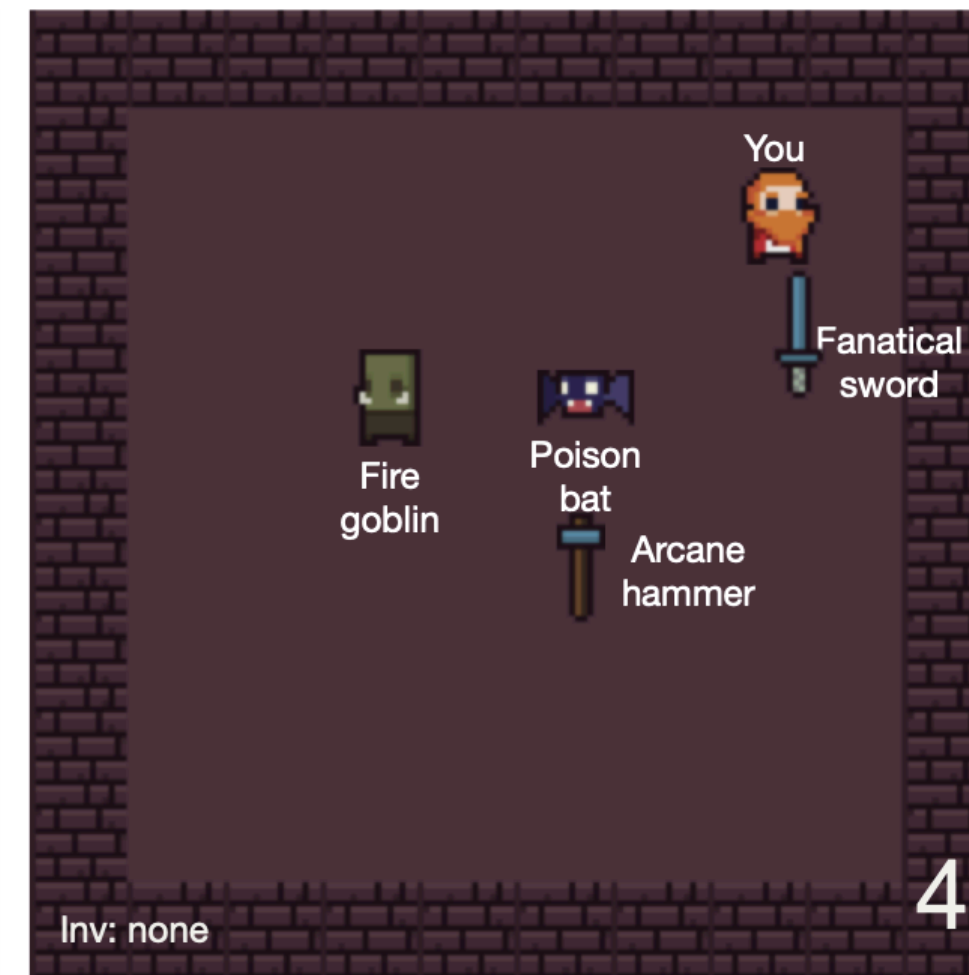
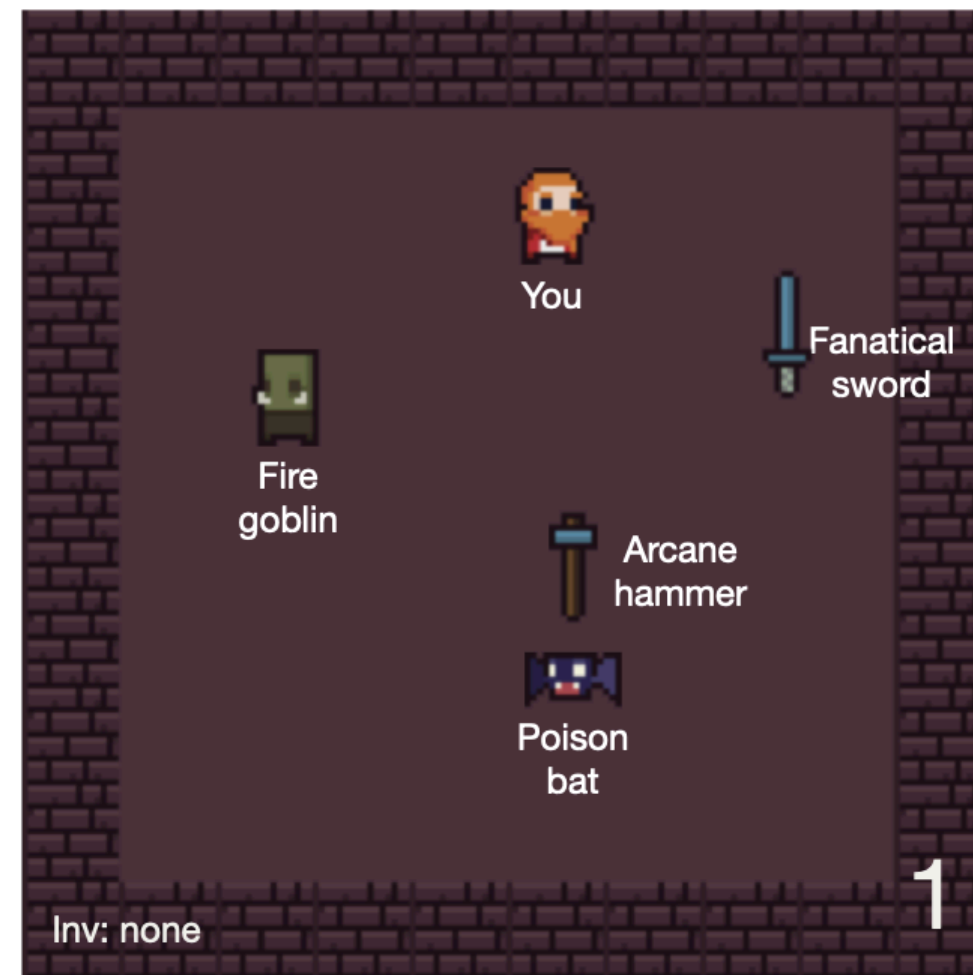
# Communicate the Structured Policies

**Doc:**

The Rebel Enclave consists of jackal, spider, and warg. Arcane, blessed items are useful for poison monsters. Star Alliance contains bat, panther, and wolf. Goblin, jaguar, and lynx are on the same team - they are in the Order of the Forest. Gleaming and mysterious weapons beat cold monsters. Lightning monsters are weak against Grandmaster's and Soldier's weapons. Fire monsters are defeated by fanatical and shimmering weapons.

**Goal:**

Defeat the Order of the Forest



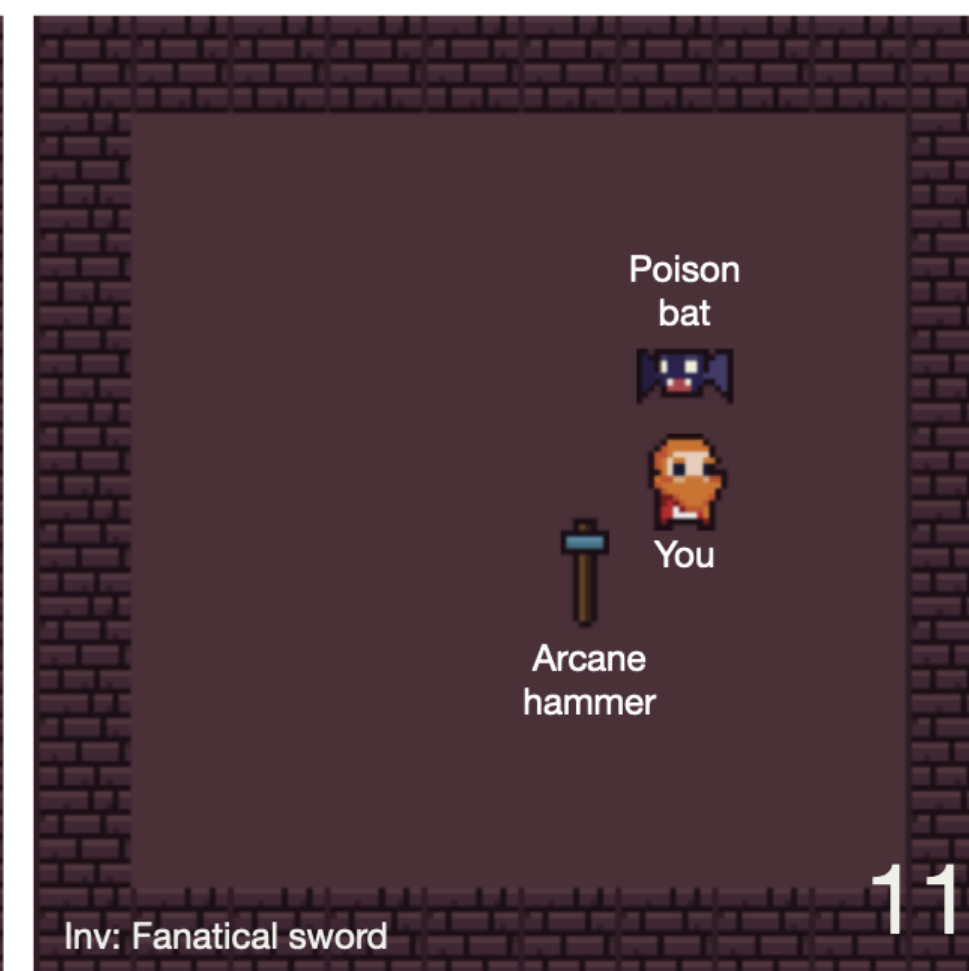
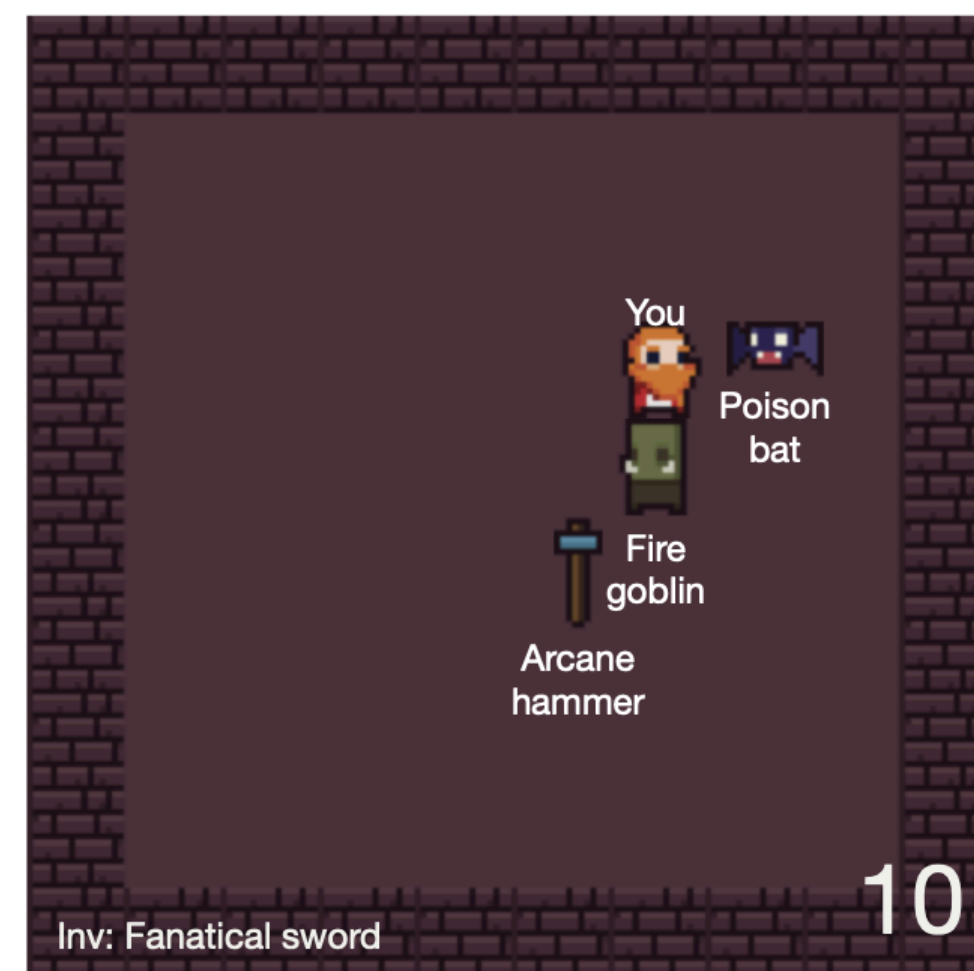
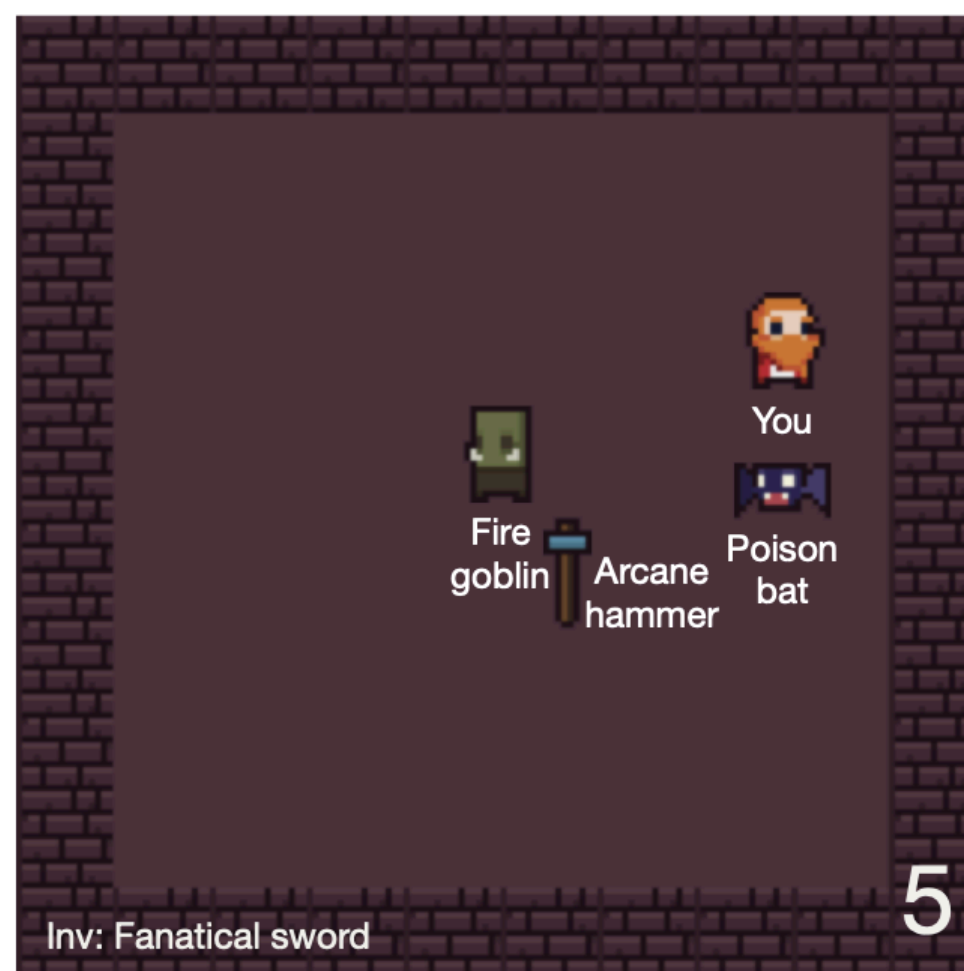
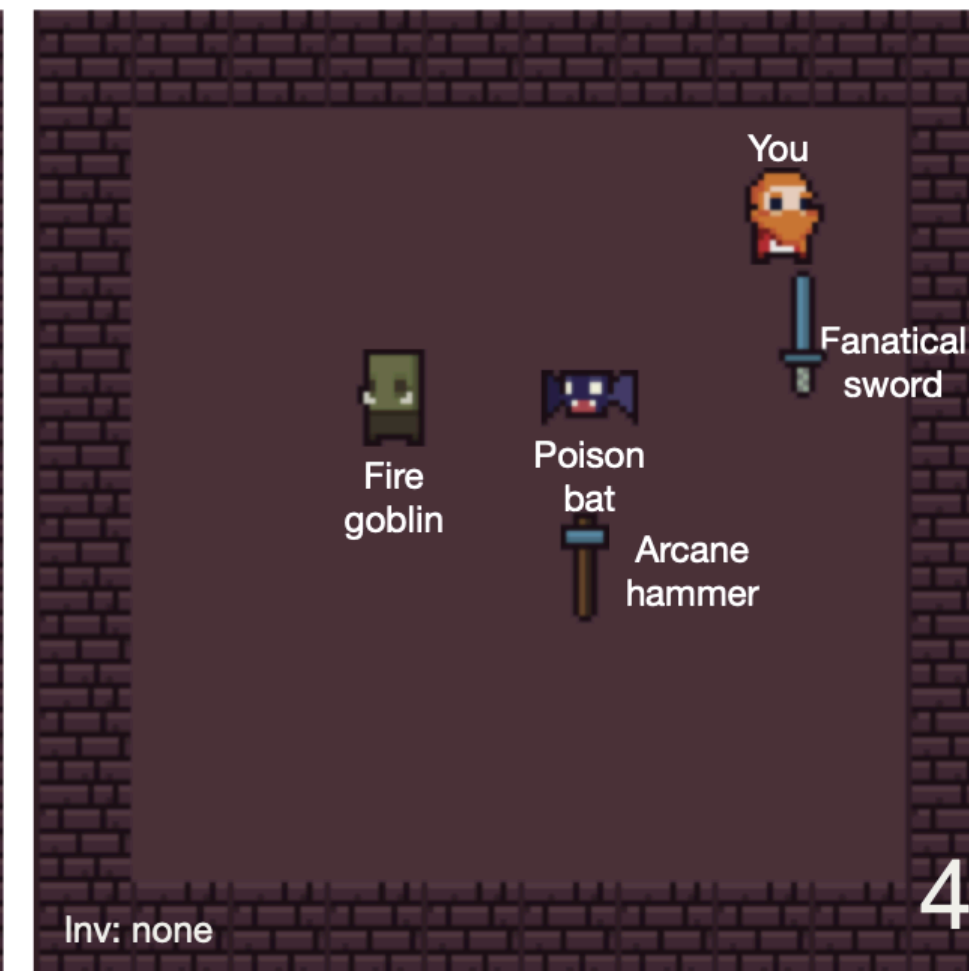
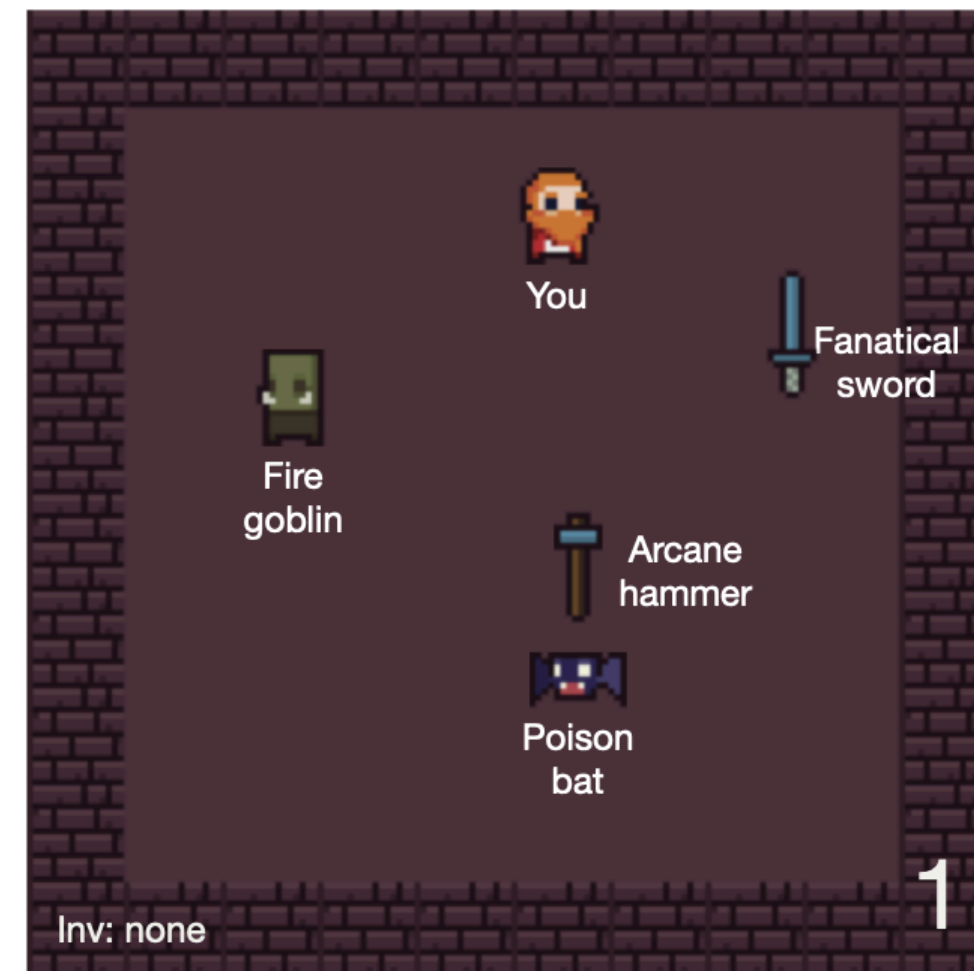
# Communicate the Structured Policies

**Doc:**

The Rebel Enclave consists of jackal, spider, and warg. Arcane, blessed items are useful for poison monsters. Star Alliance contains bat, panther, and wolf. Goblin, jaguar, and lynx are on the same team - they are in the Order of the Forest. Gleaming and mysterious weapons beat cold monsters. Lightning monsters are weak against Grandmaster's and Soldier's weapons. Fire monsters are defeated by fanatical and shimmering weapons.

**Goal:**

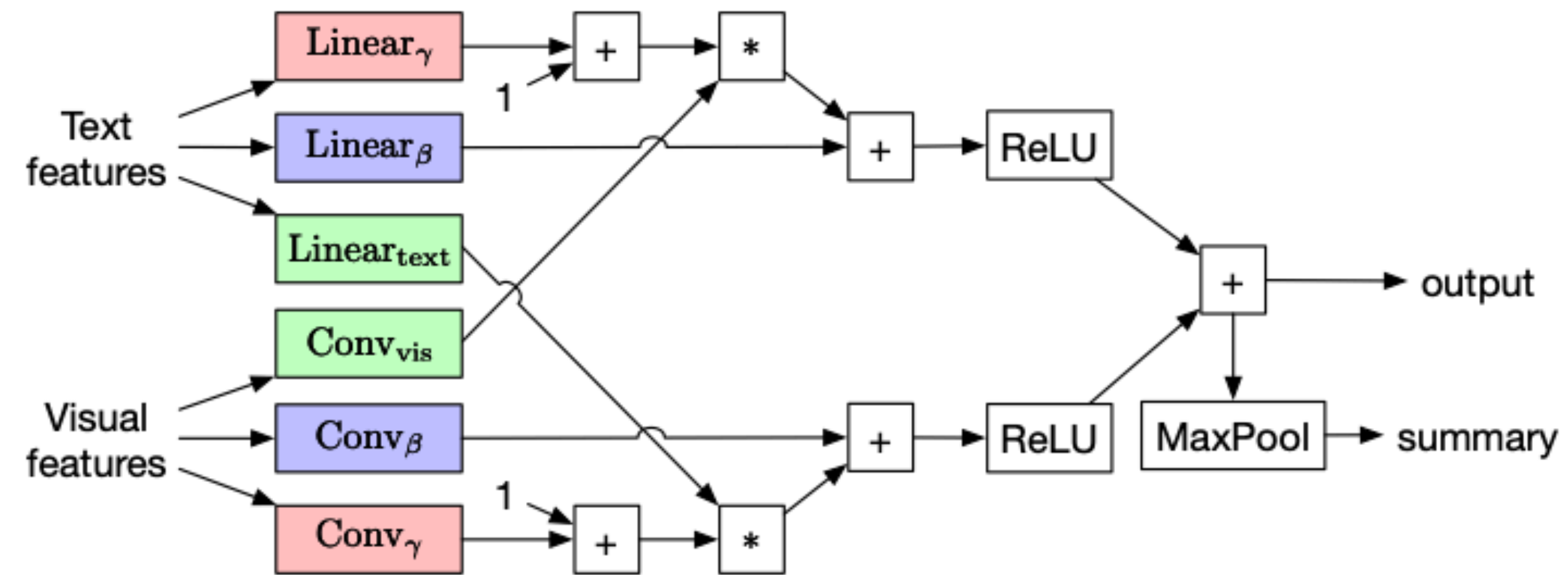
Defeat the Order of the Forest



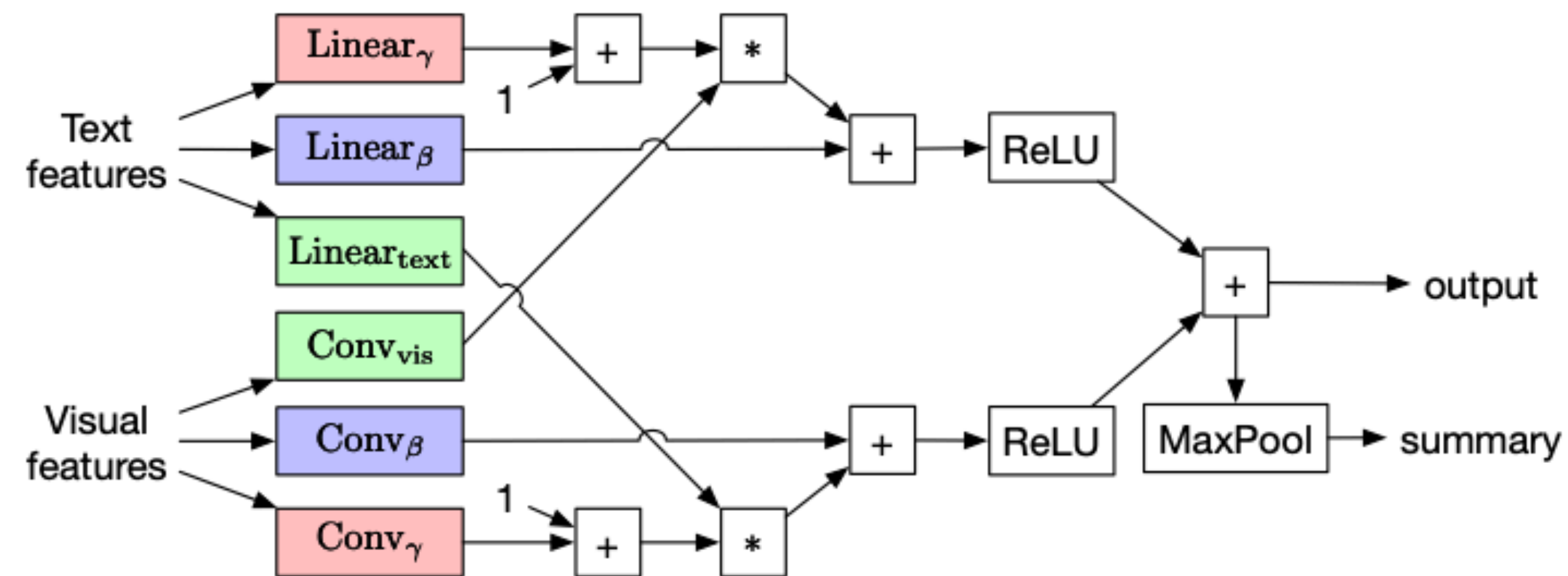
- The targeted team and their members
- Effectiveness of the modifiers and weapons
- .....

# Model: Fusion of State and Text

# Model: Fusion of State and Text



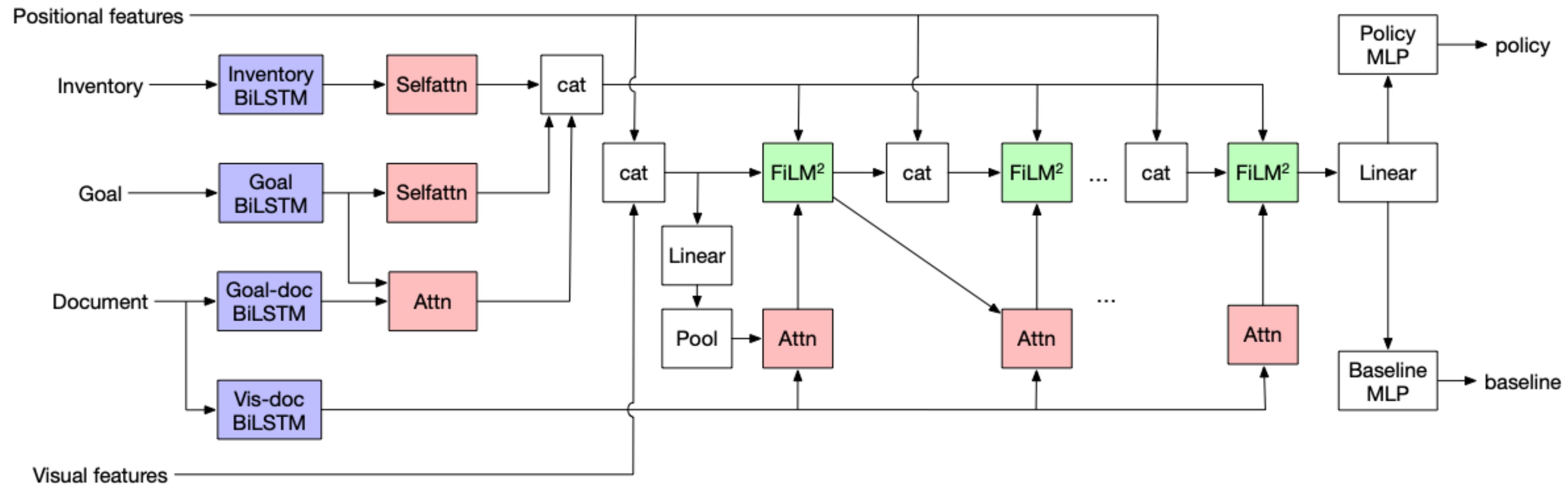
# Model: Fusion of State and Text



- Interactive encoding between the text documents and the visual states
- Filter out irrelevant text/state information for the current time step



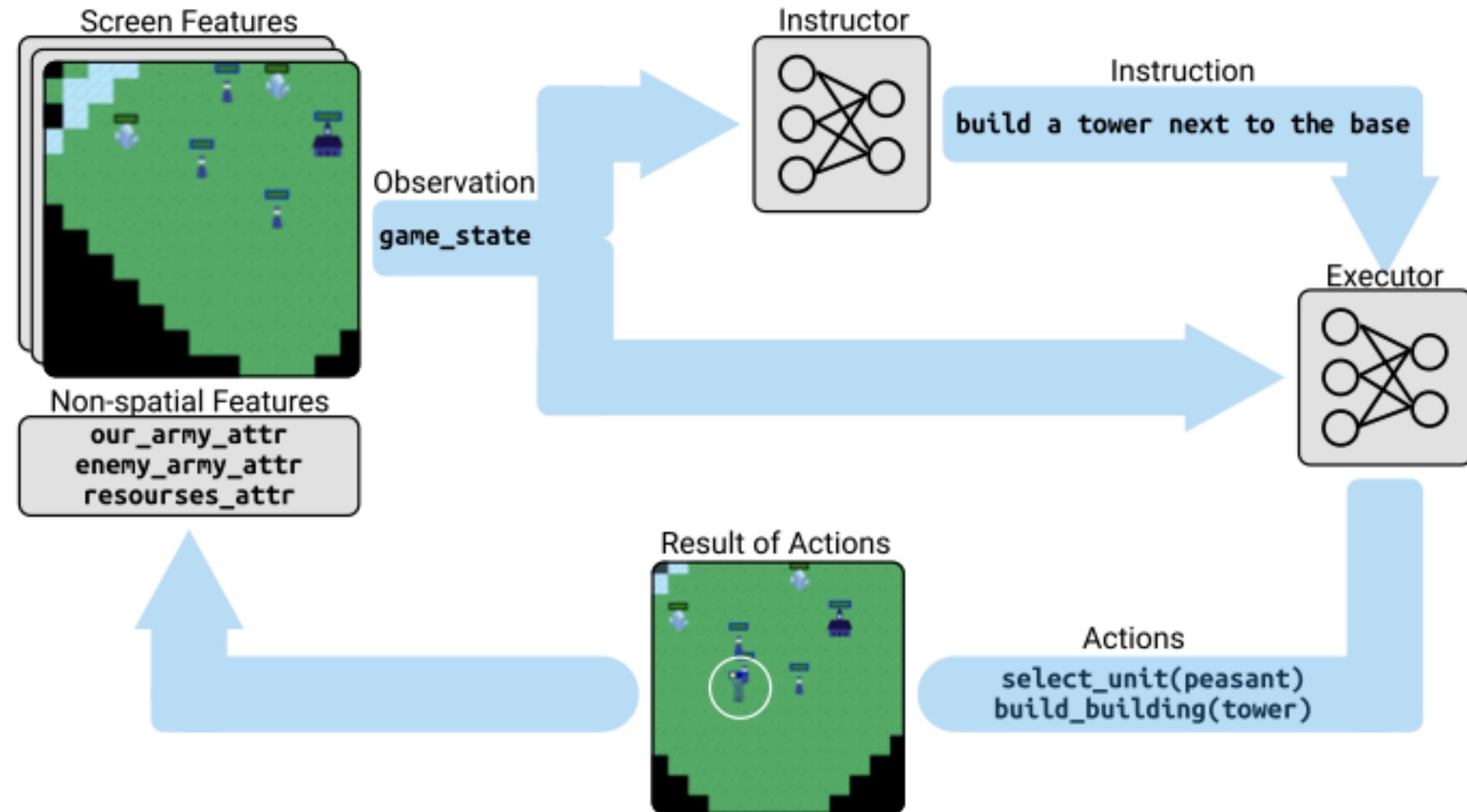
# Model: Fusion of State and Text



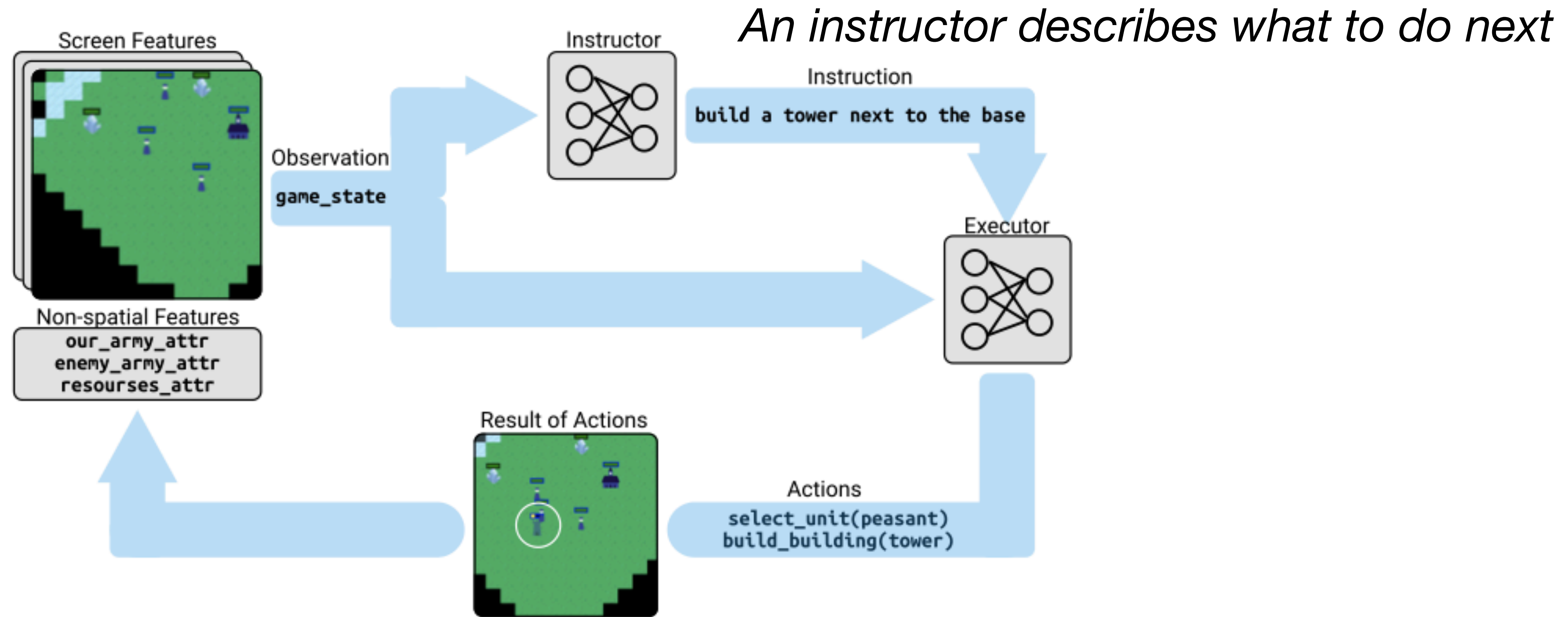
- Interactive encoding between the text documents and the visual states
- Filter out irrelevant text/state information for the current time step

# Communicate the Structured Policies

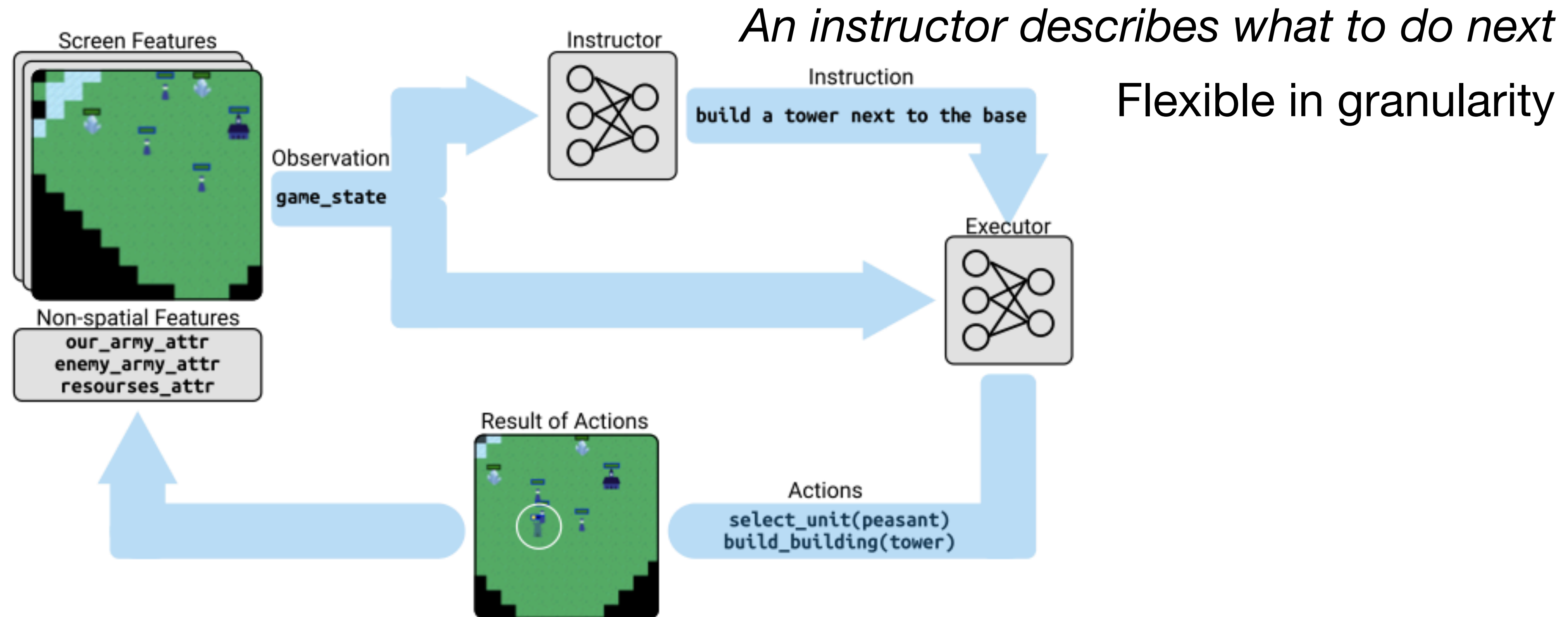
# Communicate the Structured Policies



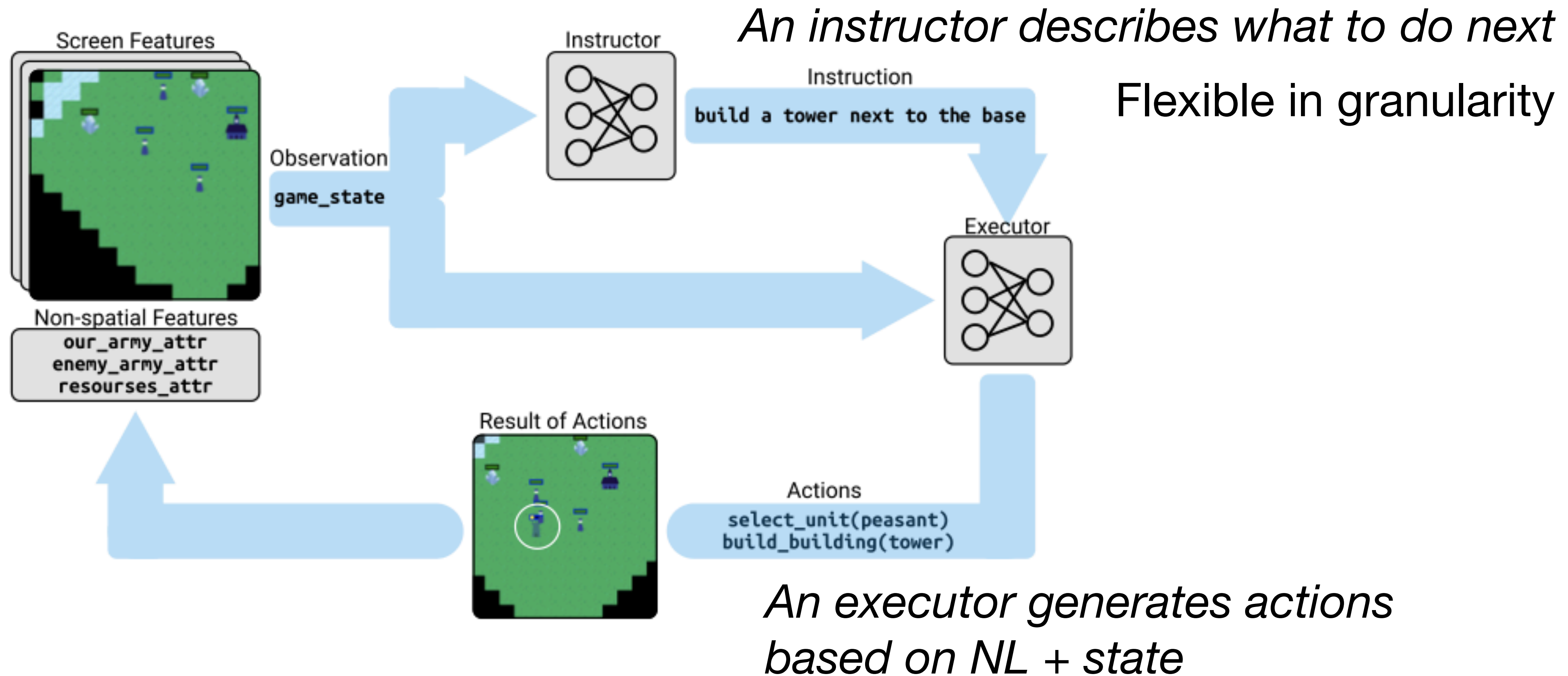
# Communicate the Structured Policies



# Communicate the Structured Policies



# Communicate the Structured Policies



# Leverage Task-independent Text

# Leverage Task-independent Text

- Free-formed text is ubiquitous



# Leverage Task-independent Text

- Free-formed text is ubiquitous
  - wikiHow for human daily task completions

# Leverage Task-independent Text

- Free-formed text is ubiquitous
  - wikiHow for human daily task completions
  - Game strategy guides (Branavan et al + 2012)

# Leverage Task-independent Text

- Free-formed text is ubiquitous
  - wikiHow for human daily task completions
  - Game strategy guides (Branavan et al + 2012)
  - etc

# Leverage Task-independent Text

- Free-formed text is ubiquitous
  - wikiHow for human daily task completions
  - Game strategy guides (Branavan et al + 2012)
  - etc
- Can we leverage this text information?

# Leverage Task-independent Text

- Free-formed text is ubiquitous
  - wikiHow for human daily task completions
  - Game strategy guides (Branavan et al + 2012)
  - etc
- Can we leverage this text information?
  - How to **obtain** and **encode** the most relevant information for a specific task at hand?

# Leverage Task-independent Text

- Free-formed text is ubiquitous
  - wikiHow for human daily task completions
  - Game strategy guides (Branavan et al + 2012)
  - etc
- Can we leverage this text information?
  - How to **obtain** and **encode** the most relevant information for a specific task at hand?
  - **When** to query this open-domain resource?

# Leverage Task-independent Text

- Free-formed text is ubiquitous
  - wikiHow for human daily task completions
  - Game strategy guides (Branavan et al + 2012)
  - etc
- Can we leverage this text information?
  - How to **obtain** and **encode** the most relevant information for a specific task at hand?
  - **When** to query this open-domain resource?
  - .....