# Vision-and-Dialog Navigation

Jesse Thomason
University of Washington

UWNLP

Daniel Gordon | Yonatan Bisk | Michael Murray | Maya Cakmak | Luke Zettlemoyer

# Bringing Robots from Industrial to Human Spaces

Industrial

# Bringing Robots from Industrial to Human Spaces

Industrial

Human

2

# Bringing Robots from Industrial to Human Spaces



Industrial

Natural Language

Navigation

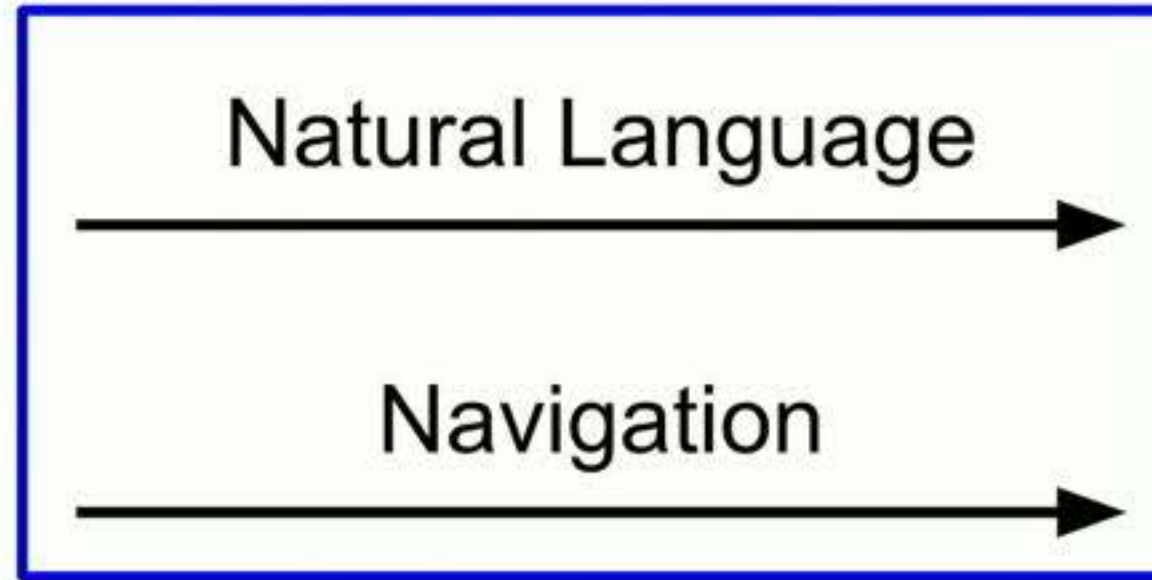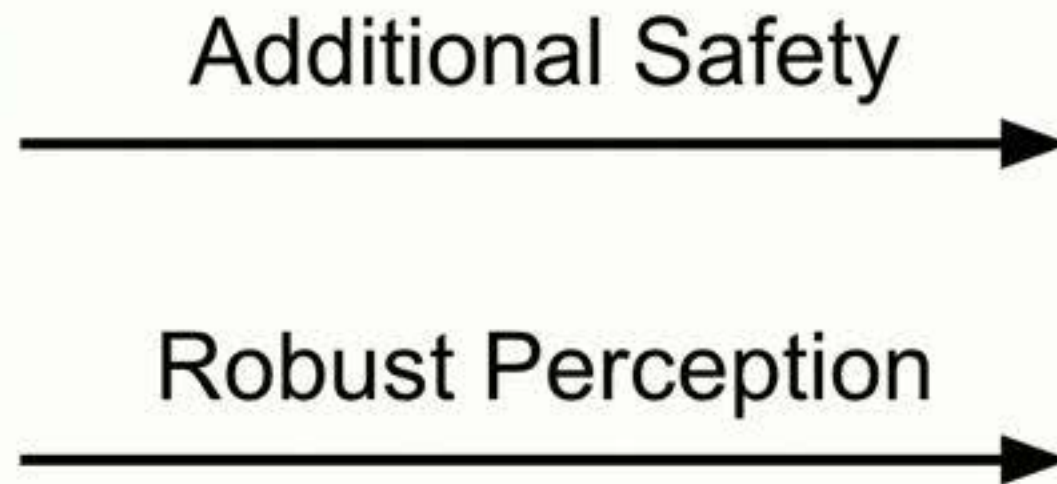Additional Safety

Robust Perception

Human

2

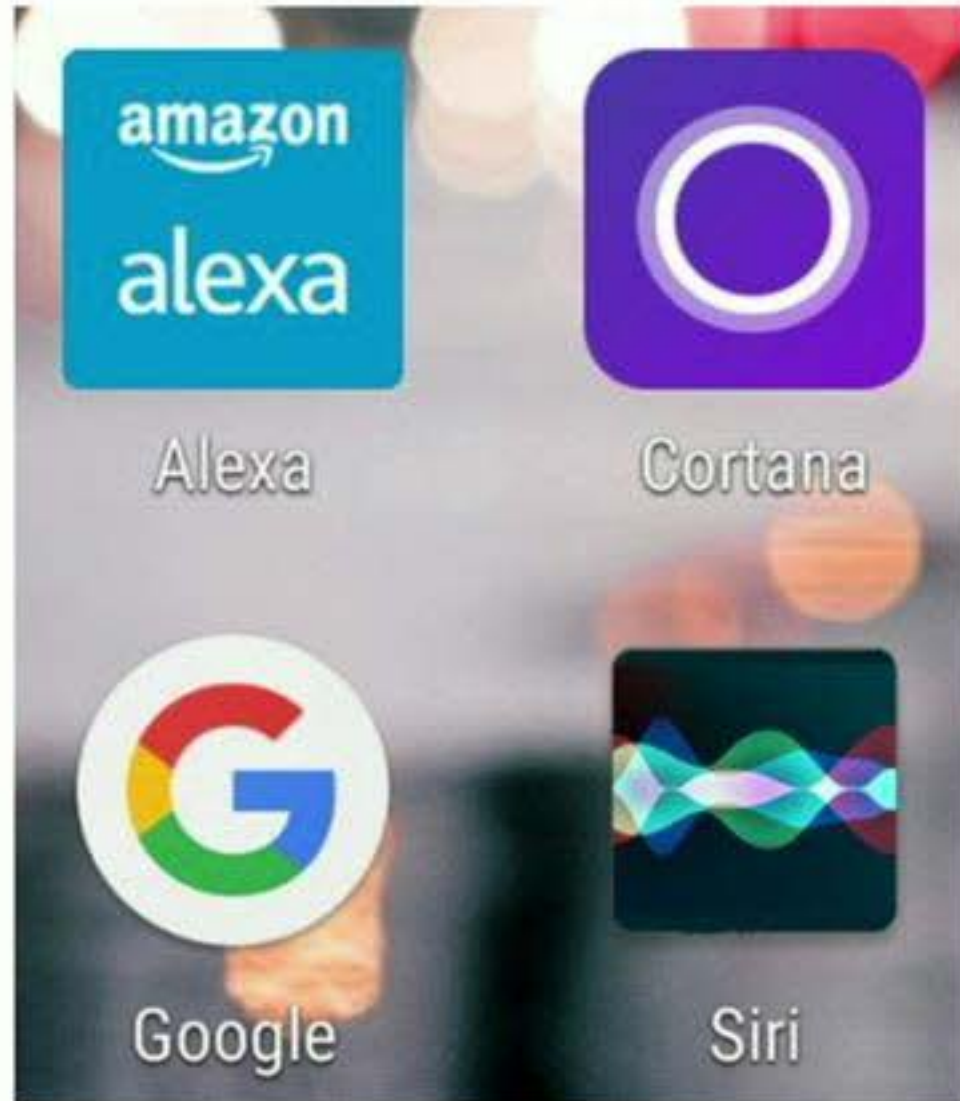# Bringing Robots from Industrial to Human Spaces



Industrial

Natural Language →

Navigation →

Additional Safety →
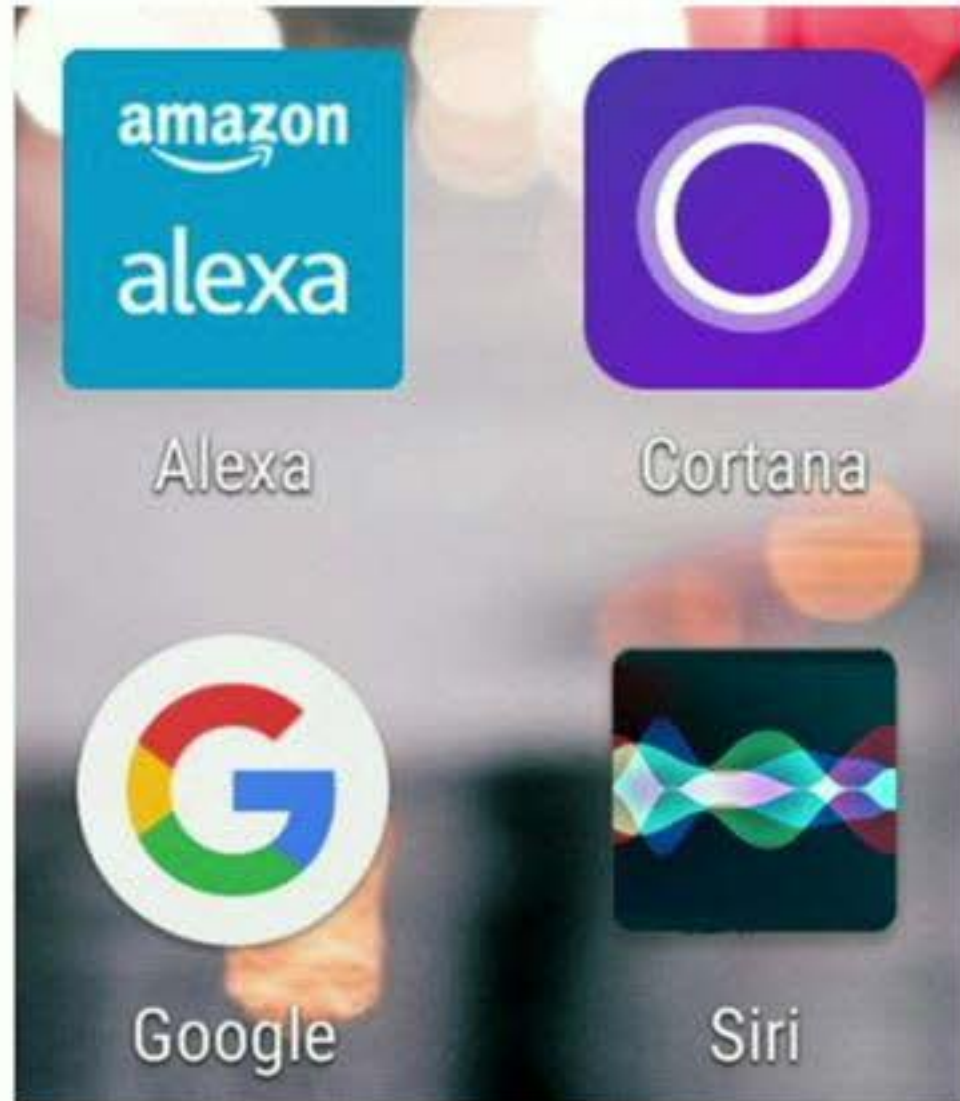
Robust Perception →

Human

2

# 1) In Human Spaces, We Use Natural Language.

# 1) In Human Spaces, We Use Natural Language.

# 2) Human Spaces Are *Dynamic* and can be *Unseen*

# 2) Human Spaces Are *Dynamic* and can be *Unseen*

# 2) Human Spaces Are *Dynamic* and can be *Unseen*

# Outline

# Outline

- Language grounding in visual environments

# Outline

- Language grounding in visual environments
  - For navigation

# Outline

- Language grounding in visual environments
  - For navigation
  - Unimodal bias [Thomason et al., NAACL'19]

# Outline

- Language grounding in visual environments
  - For navigation
  - Unimodal bias [Thomason et al., NAACL'19]
- Vision-and-Dialog Navigation [Thomason et al., *in sub*]

# Outline

- Language grounding in visual environments
  - For navigation
  - Unimodal bias [Thomason et al., NAACL'19]
- Vision-and-Dialog Navigation [Thomason et al., *in sub*]
  - New dataset - CVDN

# Outline

- Language grounding in visual environments
  - For navigation
  - Unimodal bias [Thomason et al., NAACL'19]
- Vision-and-Dialog Navigation [Thomason et al., *in sub*]
  - New dataset - CVDN
  - Navigation from dialog history

# Outline

- Language grounding in visual environments
  - For navigation
  - Unimodal bias [Thomason et al., NAACL'19]
- Vision-and-Dialog Navigation [Thomason et al., *in sub*]
  - New dataset - CVDN
  - Navigation from dialog history
- Next steps

# Connecting Language and Vision

# Connecting Language and Vision

- Common paradigm:

# Connecting Language and Vision

- Common paradigm:
  - **Inputs**:
    - Language tokens (e.g., question)

# Connecting Language and Vision

- Common paradigm:
  - **Inputs**:
    - Language tokens (e.g., question)
    - Visual context (e.g., photograph)

# Connecting Language and Vision

- Common paradigm:
  - **Inputs**:
    - Language tokens (e.g., question)
    - Visual context (e.g., photograph)
- Visual contexts differ in quality across datasets.

# Connecting Language and Vision

- Common paradigm:
  - **Inputs**:
    - Language tokens (e.g., question)
    - Visual context (e.g., photograph)
- Visual contexts differ in quality across datasets.
- Output can be a single classification or a sequence.

# Visual Fidelity

Rendered

**Visual Context**

Static ⟵————————————————————⟶ Dynamic

Photorealistic

**Visual Fidelity**

Rendered

CLEVR
[Johnson et al., CVPR'17]

Static ←————————————————————————→ Dynamic

**Visual Context**

Photorealistic

**Visual Fidelity**

Rendered

Dynamic

Static

**Visual Context**

CLEVR
[Johnson et al., CVPR'17]

VQA
[Antol et al., CVPR'15]

Photorealistic

**Visual Fidelity**

Rendered

CLEVR
[Johnson et al., CVPR'17]

Instruction Following
[Chen and Mooney, AAAI'11]

**Visual Context**

Static ← → Dynamic

VQA
[Antol et al., CVPR'15]

Photorealistic

8

**Visual Fidelity**

Rendered



CLEVR
[Johnson et al., CVPR'17]



Instruction Following
[Chen and Mooney, AAAI'11]

**Visual Context**

Static ← → Dynamic



VQA
[Antol et al., CVPR'15]



Room-to-Room
[Anderson et al., CVPR'18]

Photorealistic

8

# Outline

- Language grounding in visual environments
  - For navigation
  - Unimodal bias [Thomason et al., NAACL'19]
- Vision-and-Dialog Navigation [Thomason et al., *in sub*]
  - New dataset - CVDN
  - Navigation from dialog history
- Next steps

# Vision-and-Language Navigation



Goal: 7.5m

Exit the bathroom. Turn left and exit the room using the door on the left. Wait there.

# *"Turn around and exit the library, head down the…"*

- Low-level instructions for moving through the environment.

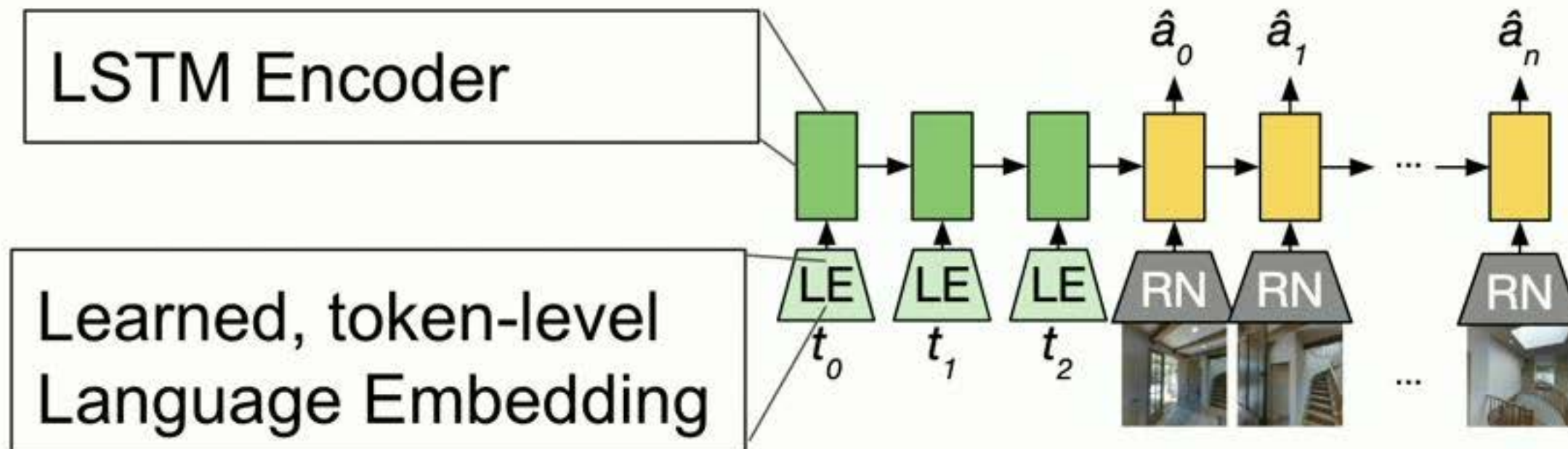# *"Turn around and exit the library, head down the…"*

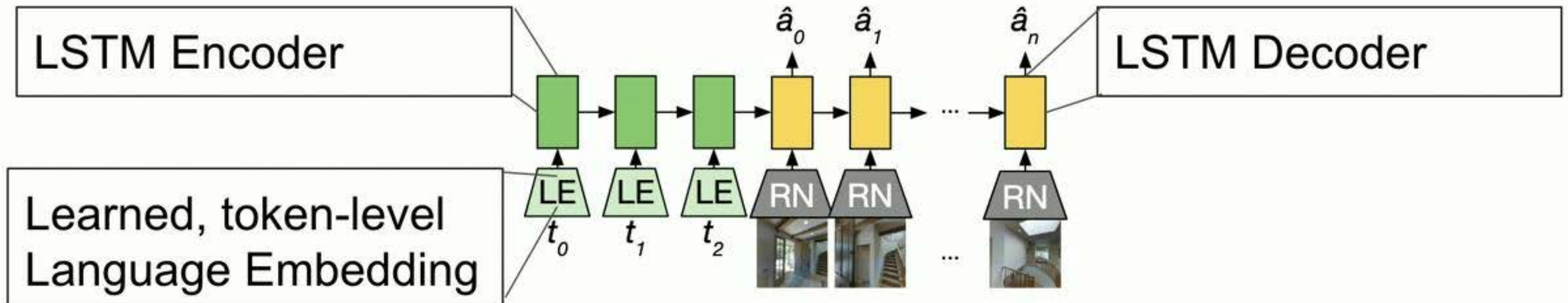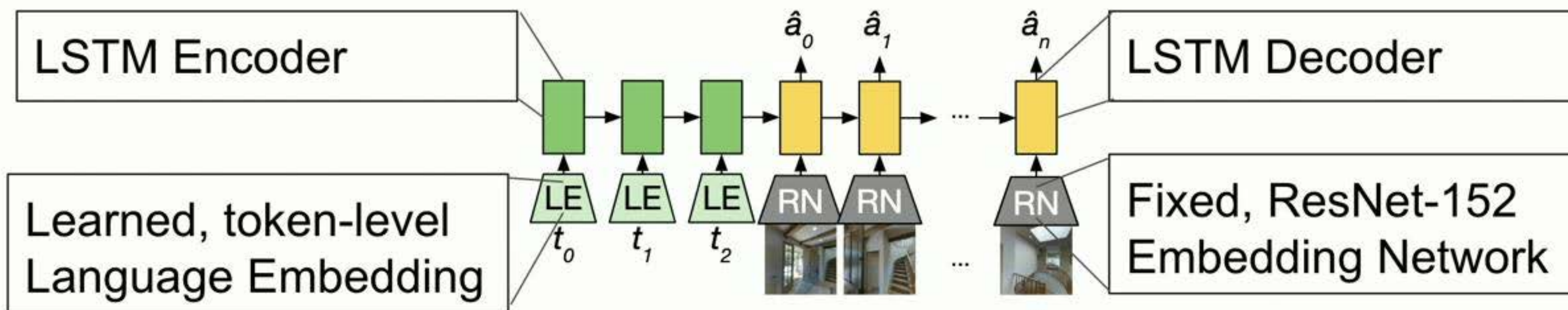- Low-level instructions for moving through the environment.

# Sequence-to-Sequence Model

- Encode the language tokens.
- Decode a sequence of actions to take in the environment.
- At every timestep, receive a new visual observation.
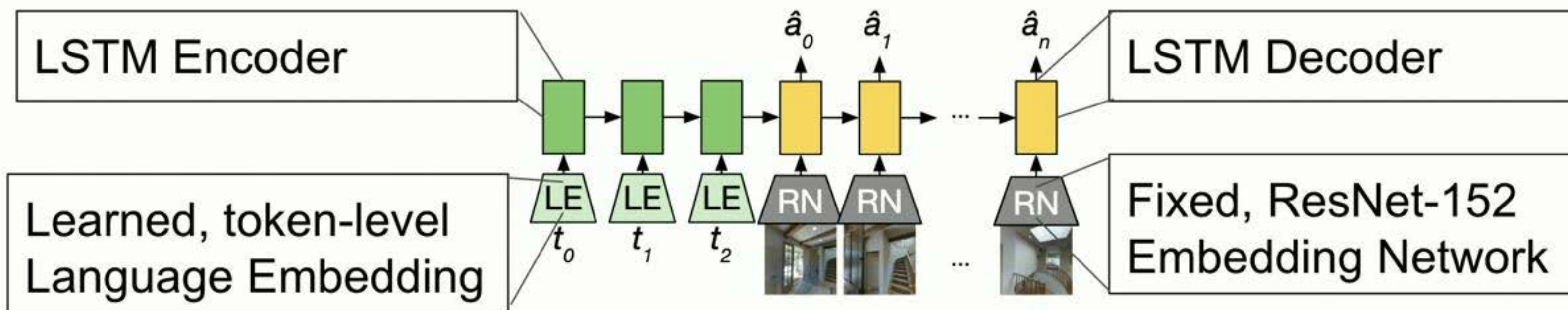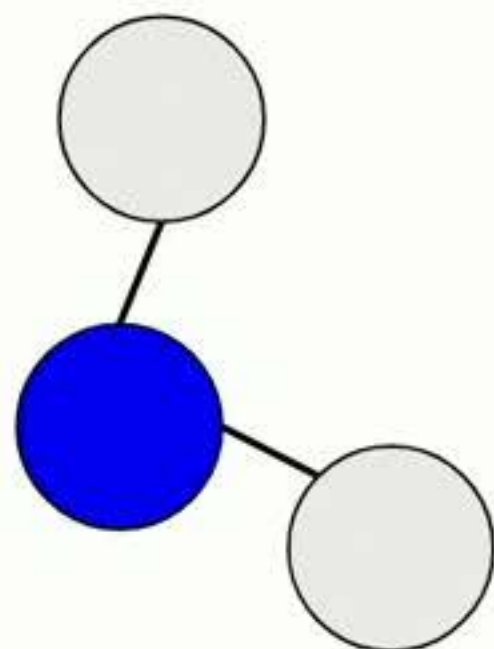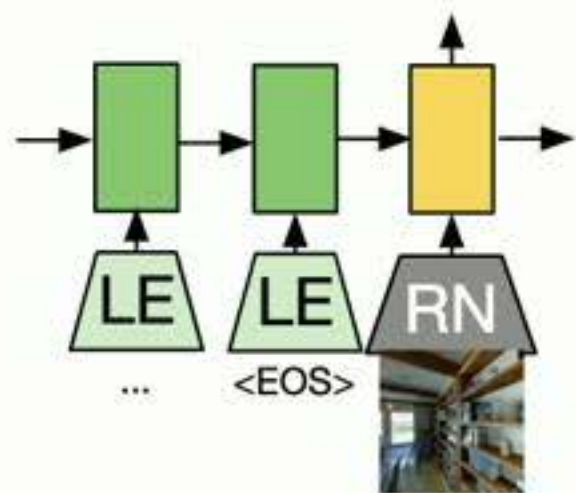
# Sequence-to-Sequence Model

- Encode the language tokens.
- Decode a sequence of actions to take in the environment.
- At every timestep, receive a new visual observation.



Learned, token-level Language Embedding

# Sequence-to-Sequence Model

- Encode the language tokens.

- Decode a sequence of actions to take in the environment.

- At every timestep, receive a new visual observation.



LSTM Encoder

Learned, token-level Language Embedding

# Sequence-to-Sequence Model

- Encode the language tokens.
- Decode a sequence of actions to take in the environment.
- At every timestep, receive a new visual observation.



LSTM Encoder

$\hat{a}_0$   $\hat{a}_1$   $\hat{a}_n$

LSTM Decoder

Learned, token-level Language Embedding

LE   LE   LE   RN   RN   RN

$t_0$   $t_1$   $t_2$

12

# Sequence-to-Sequence Model

- Encode the language tokens.
- Decode a sequence of actions to take in the environment.
- At every timestep, receive a new visual observation.



LSTM Encoder

$\hat{a}_0$  $\hat{a}_1$  $\hat{a}_n$

LSTM Decoder

Learned, token-level Language Embedding

LE  LE  LE  RN  RN  RN

$t_0$  $t_1$  $t_2$

Fixed, ResNet-152 Embedding Network

# Sequence-to-Sequence Model

- Encode the lang___

  We will build on this model.

- Decode a seque___ e in the environment.

- At every timestep, receive a new visual observation.

| LSTM Encoder | | $\hat{a}_0$ $\hat{a}_1$ $\hat{a}_n$ | LSTM Decoder |
|---|---|---|---|
| Learned, token-level Language Embedding | LE LE LE $t_0$ $t_1$ $t_2$ | RN RN ... RN | Fixed, ResNet-152 Embedding Network |

# Sequence-to-Sequence Model

- Train to predict the action a shortest-path planner would take from the current state.

# Sequence-to-Sequence Model

- Train to predict the action a shortest-path planner would take from the current state.

# Sequence-to-Sequence Model

- Train to predict the action a shortest-path planner would take from the current state.

# Language Grounding for QA and Navigation

**Input:**

**Instruction + Frame**

**Question + Frame**

**Question + Frame**



Room-2-Room
[Anderson et al., CVPR'18]

Embodied QA
[Das et al., CVPR'18]

Interactive QA
[Gordon et al., CVPR'18]

**Output:**

**Navigation Actions**

**Navigation Actions +
Answer Action**

**Navigation Actions +
Answer Action**

# Language Grounding for QA and Navigation

**Input:**

**Instruction + Frame**

**Question + Frame**

**Question + Frame**



Room-2-Room
[Anderson et al., CVPR'18]

Embodied QA
[Das et al., CVPR'18]

Interactive QA
[Gordon et al., CVPR'18]

**Output:**

**Navigation Actions**

**Navigation Actions +**
Answer Action

**Navigation Actions +**
Answer Action

# Outline

- Language grounding in visual environments
  - For navigation
  - **Unimodal bias [Thomason et al., NAACL'19]**
- Vision-and-Dialog Navigation [Thomason et al., *in sub*]
  - New dataset - CVDN
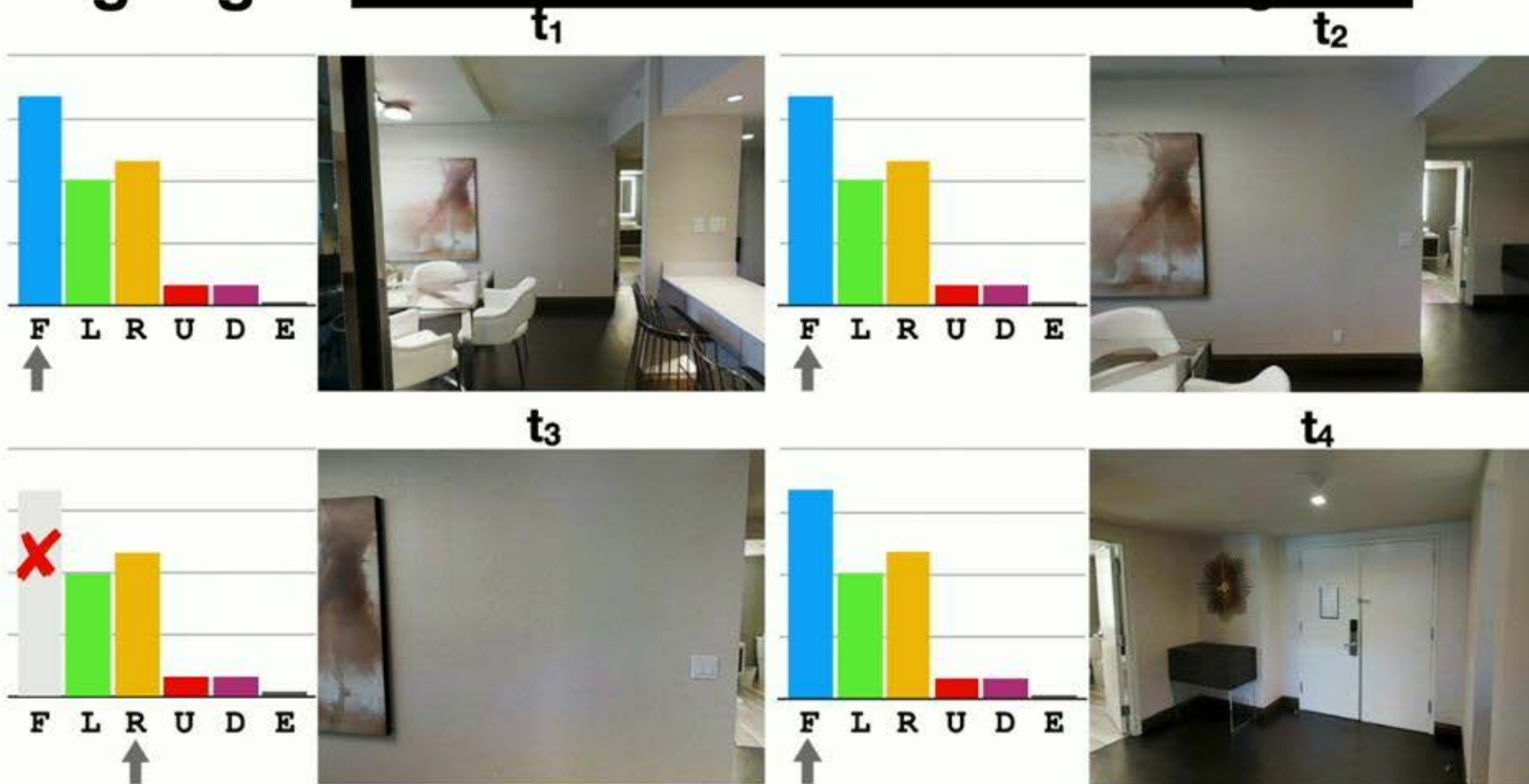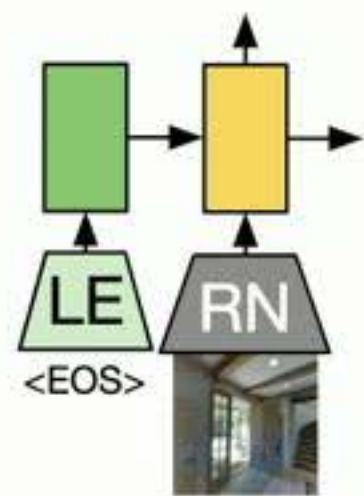  - Navigation from dialog history
- Next steps

# Inputs and Outputs

**Language:** *"Walk past the bar and turn right."*



Actions: **F**orward, turn **L**eft & **R**ight, tilt **U**p & **D**own, **E**nd

# Inputs and Outputs

**Language:** *"Walk past the bar and turn right."*



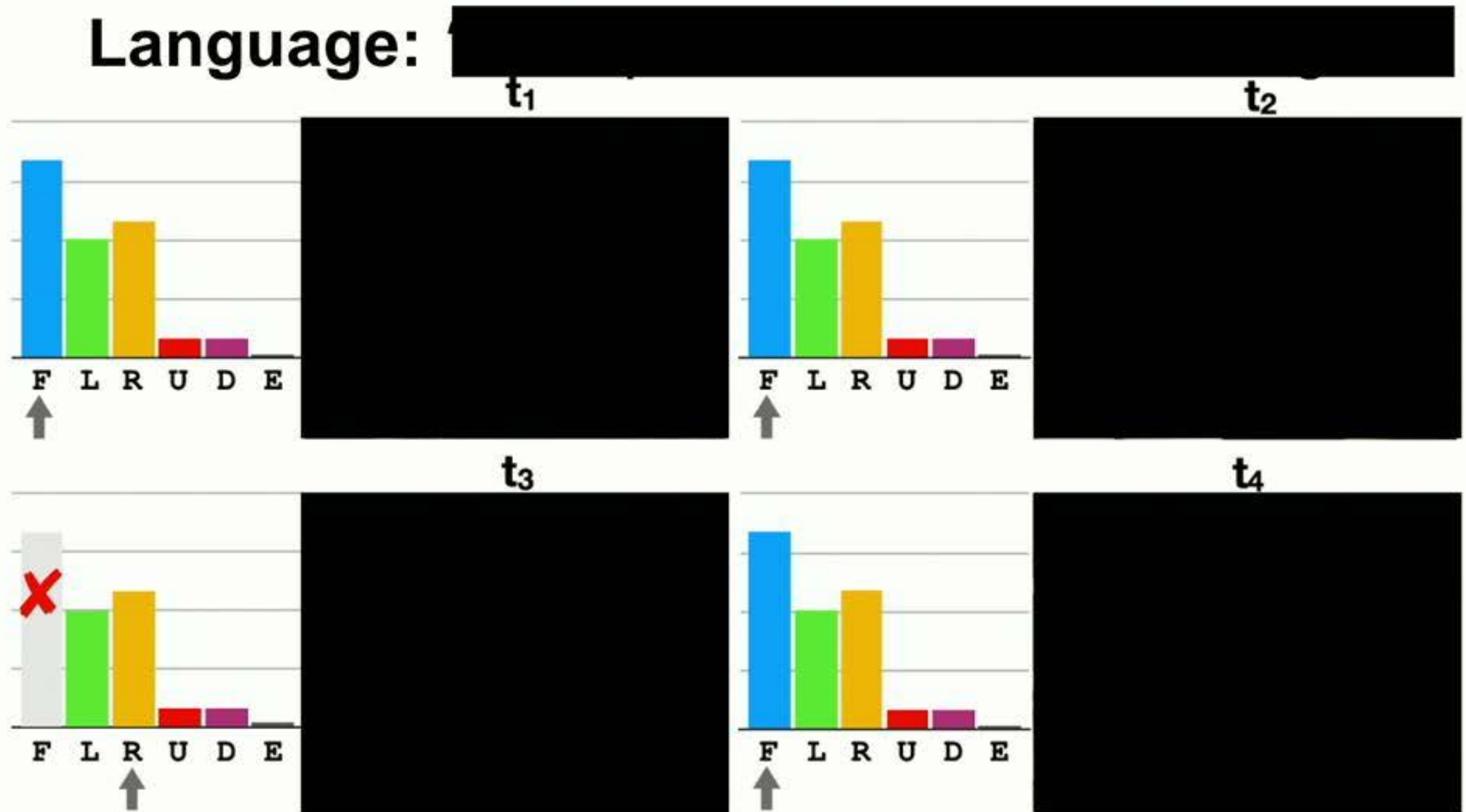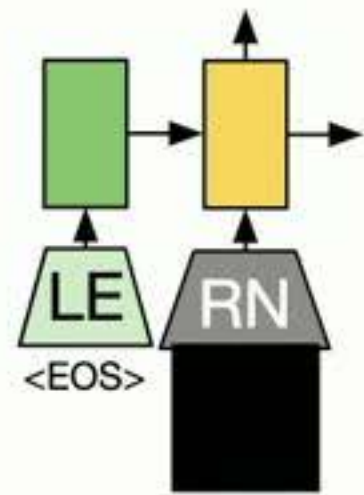Actions: **F**orward, turn **L**eft & **R**ight, tilt **U**p & **D**own, **E**nd

17

# Inputs and Outputs

**Language:** *"Walk past the bar and turn right."*



+Visual Observation
+Available Actions

Actions: **F**orward, turn **L**eft & **R**ight, tilt **U**p & **D**own, **E**nd

# Unimodal Ablation - *Vision Only*

**Language:** ████████████████████████████



Actions: **F**orward, turn **L**eft & **R**ight, tilt **U**p & **D**own, **E**nd

# Unimodal Ablation - *Language Only*

**Language:** *"Walk past the bar and turn right."*



Actions: **F**orward, turn **L**eft & **R**ight, tilt **U**p & **D**own, **E**nd

# Unimodal Ablation - *Action Only*

**Language:**
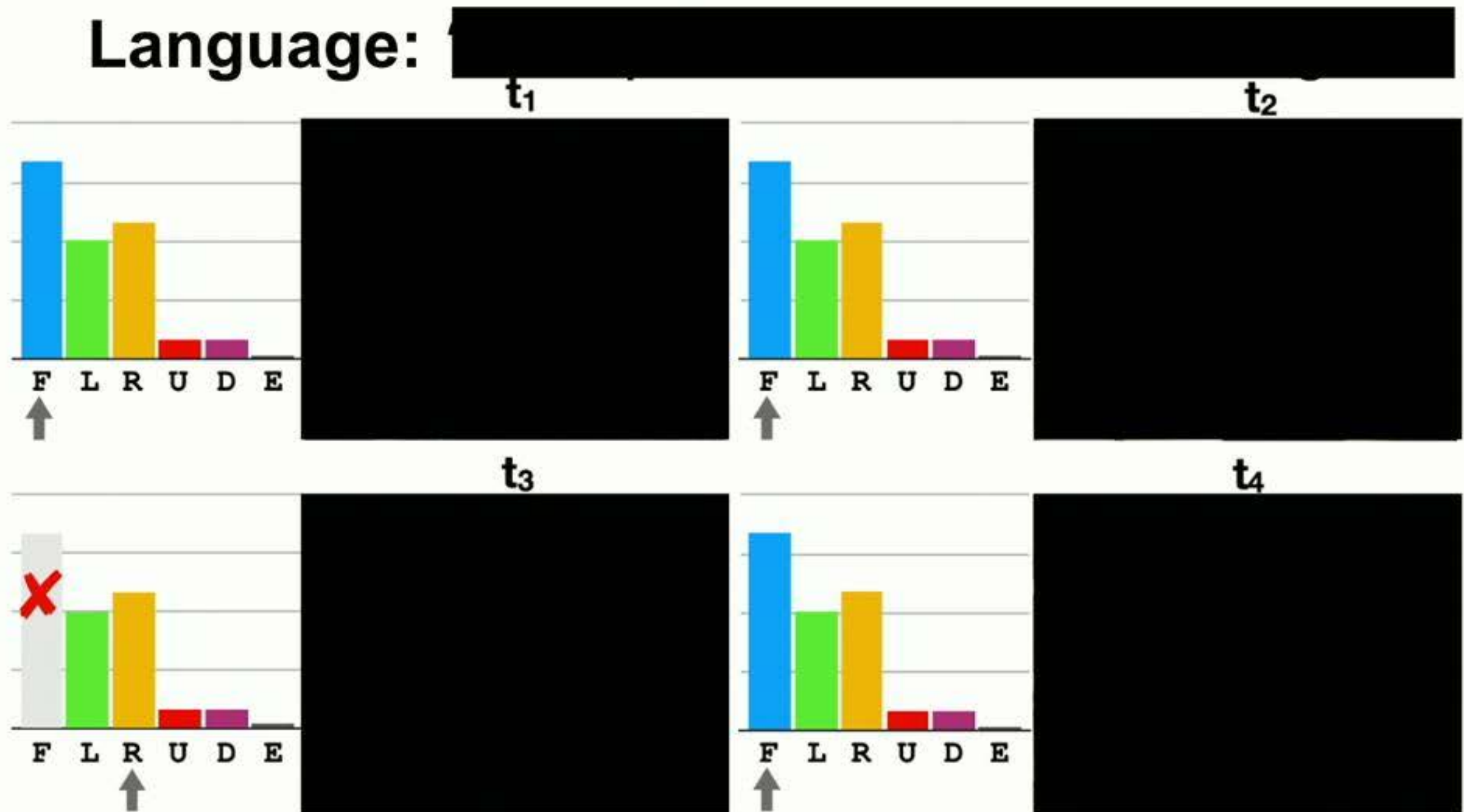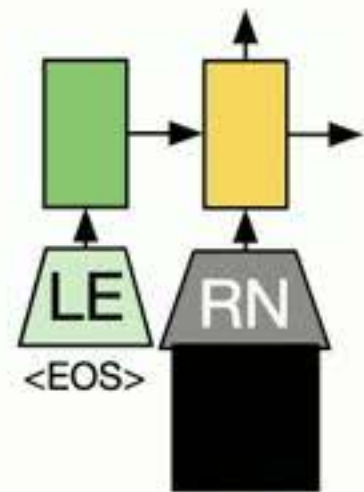


Actions: **F**orward, turn **L**eft & **R**ight, tilt **U**p & **D**own, **E**nd

# Unimodal Model Ablations

# Unimodal Ablation - *Action Only*

**Language:** 



Actions: **F**orward, turn **L**eft & **R**ight, tilt **U**p & **D**own, **E**nd

# Unimodal Model Ablations

# Unimodal Model Ablations



Room-2-Room
[Anderson et al., CVPR'18]

**Beats**
**baseline**

**Vision-, language-, and**
**action-only models.**

# Unimodal Model Ablations



Room-2-Room
[Anderson et al., CVPR'18]

Embodied QA
[Das et al., CVPR'18]

**Beats
baseline**

**Vision-, language-, and
action-only models.**

**Vision-, language-, and
action-only models.**

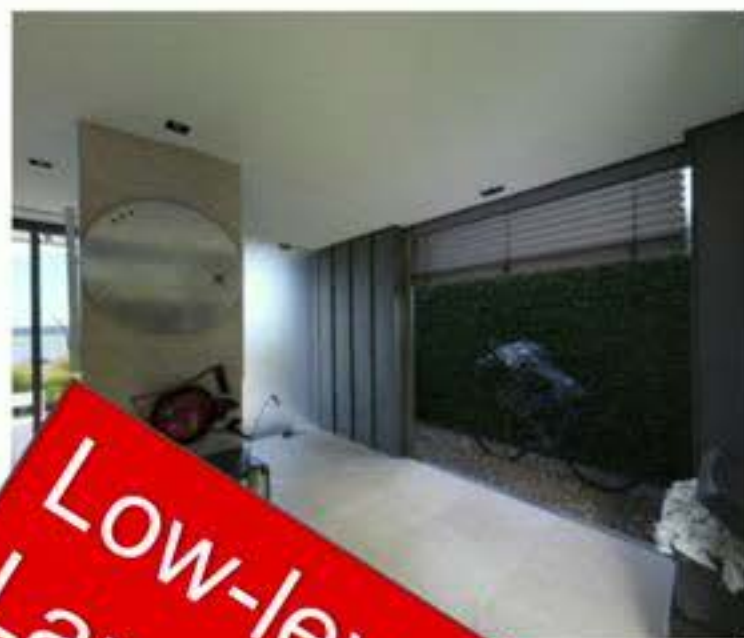# Unimodal Model Ablations



Room-2-Room
[Anderson et al., CVPR'18]



Embodied QA
[Das et al., CVPR'18]



Interactive QA
[Gordon et al., CVPR'18]

**Beats**
**baseline**

**Vision-, language-, and
action-only models.**

**Vision-, language-, and
action-only models.**

**Vision-only model.**

# Unimodal Model Ablations



Room-2-Room
[Anderson et al., CVPR'18]



Embodied QA
[Das et al., CVPR'18]



Interactive QA
[Gordon et al., CVPR'18]

**Beats baseline**

**Vision-, language-, and action-only models.**

**Vision-, language-, and action-only models.**

**Vision-only model.**

**Beats initial model!**

**Language-only model.**

21

# Unimodal Model Ablations



|  | Room-2-Room [Anderson et al., CVPR'18] | Embodied QA [Das et al., CVPR'18] | Interactive QA [Gordon et al., CVPR'18] |
|---|---|---|---|
| **Beats baseline** | Vision-, language-, and action-only models. | Vision-, language-, and action-only models. | Vision-only model. |
| **Beats initial model!** | Language-only model. | Vision-only model. | |

21

# Unimodal Model Ablations



**Low-level Language**

| | [Anders... '18] | Embodied QA [Das et al., CVPR'18] | Interactive QA [Gordon et al., CVPR'18] |
|---|---|---|---|
| **Beats baseline** | Vision-, langu... and action-only models. | Vision-, language-, and action-only models. | Vision-only model. |
| **Beats initial model!** | Language-only model. | Vision-only model. | (none) |

22

# Unimodal Model Ablations



Photorealistic Environments

Low-level Language

[Anders... ...'18]

Embodied QA
[Das et al., CVPR'18]

Interactive QA
[Gordon et al., CVPR'18]

| Beats baseline | Vision-, langu... ...and action-only models. | Vision-, language-, and action-only models. | Vision-only model. |
| Beats initial model! | Language-only model. | Vision-only model. | (none) |

# Unimodal Model Ablations



Photorealistic Environments

Low-level Language

[Anders... ...'18]

Simple Visual Environments

[Das ... ...]

Interactive QA
[Gordon et al., CVPR'18]

| | | |
|---|---|---|
| **Beats baseline** | **Vision-, langu... ...nd action-only models.** | **Vision-, langu... ...nd action-only models.** | **Vision-only model.** |
| **Beats initial model!** | **Language-only model.** | **Vision-only model.** | **(none)** |

22

# Unimodal Model Ablations



Photorealistic Environments

Low-level Language

[Anders... ...'18]

High-level Language

Simple Visual Environments

[Das e... ...]

Interactive QA
[Gordon et al., CVPR'18]

**Beats baseline**
Vision-, langu... ...and action-only models.

Vision-, langu... ...and action-only models.

Vision-only model.

**Beats initial model!**
Language-only model.

Vision-only model.

**(none)**

22

# Unimodal Model Ablations



| | | |
|---|---|---|
| **Beats baseline** | Vision-, language- and action-only models. | Vision-, language- and action-only models. | Vision-only model. |
| **Beats initial model!** | Language-only model. | Vision-only model. | (none) |

# Lessons

# Lessons

- Unimodal baselines expose dataset bias.

# Lessons

- Unimodal baselines expose dataset bias.
  - More appropriate than non-learning baselines.

# Lessons

- Unimodal baselines expose dataset bias.
  - More appropriate than non-learning baselines.
- Rich visual context prevents visual overfitting.

# Lessons

- Unimodal baselines expose dataset bias.
  - More appropriate than non-learning baselines.
- Rich visual context prevents visual overfitting.
- *Underspecified* language context prevents language overfitting.

# Lessons

- Unimodal baselines expose dataset bias.
  - More appropriate than non-learning baselines.
- Rich visual context prevents visual overfitting.
- *Underspecified* language context prevents language overfitting.
  - Also, low-level instructions are somewhat unnatural.

# Lessons

- Unimodal baselines expose dataset bias.
  - More appropriate than non-learning baselines.
- Rich visual context prevents visual overfitting.
- *Underspecified* language context prevents language overfitting.
  - Also, low-level instructions are somewhat unnatural.
- Why not both?

# Outline

- Language grounding in visual environments
  - For navigation
  - Unimodal bias [Thomason et al., NAACL'19]
- Vision-and-Dialog Navigation [Thomason et al., *in sub*]
  - New dataset - CVDN
  - Navigation from dialog history
- Next steps

# "Turn around and exit the library, head down the…"

# "Turn around and exit the library, head down the..."

- *Room-to-Room* uses low-level language.

*"Turn around and exit the library, head down the..."*

- *Room-to-Room* uses low-level language.
- 24% → 80% of human performance since '18.

# "Go to the room with a plant."

- This instruction is *underspecified*.

*"Go to the room with a plant."*

- This instruction is *underspecified*.
- This instruction is *ambiguous*.

# What if we could just… ask?

- *"Should I continue into the living room or go right towards the kitchen?"*

**Visual Fidelity**

Rendered

CLEVR
[Johnson et al., CVPR'17]

Instruction Following
[Chen and Mooney, AAAI'11]

**Visual Context**

Static ← → Dynamic

VQA
[Antol et al., CVPR'15]

Photorealistic

Room-to-Room
[Anderson et al., CVPR'18]

**Visual Fidelity**

Rendered

Static ← → Dynamic

**Visual Context**

A **cylinder** is next to a **yellow object**.

Q1 : What shape is **the object**?

A1 : Sphere

Q2 : And material?

A2 : Metal

Q3 : What about **that cylinder**?

A3 : Rubber

CLEVR-Dialog [Kottur et al., NAACL'19]

**Visual Dialog**

Q: What is the gender of the one in the white shirt ?

A: She is a woman

Q: What is she doing ?

A: Playing a Wii game

Q: Is that a man to her right

A: No, it's a woman

VisDial [Das et al., CVPR'17]

Photorealistic

29

**Language Source**

Templates

**Guidance Abstraction**

Semantic ← → Visual

Humans

Talk the Walk
[de Vries et al., arXiv'18]

Tourist

Guide

Hey, what is near you?

Is that a shop or restaurant?

Your target intersection has three restaurants and a bar

Go to a restaurant corner facing the pub.

Evaluate

Hello, in front of me is a Brook Brothers

It is a clothing shop

Should I go to pub corner? Or one of the restaurant corners?

I'm there

30

**Language Source**

Templates

**Guidance Abstraction**

Semantic ← → Visual

VLNA [Nguyen et al., CVPR'19]

I am lost! Help!

Turn 60 degrees right, go forward, turn left

Hey, what is near you?

Hello, in front of me is a Brook Brothers

Is that a shop or restaurant?

It is a clothing shop

Your target intersection has three restaurants and a bar

Should I go to pub corner? Or one of the restaurant corners?

Go to a restaurant corner facing the pub.

I'm there

Evaluate

Tourist

Guide

**Talk the Walk**
**[de Vries et al., arXiv'18]**

Humans

30

# Outline

- Language grounding in visual environments
  - For navigation
  - Unimodal bias [Thomason et al., NAACL'19]
- Vision-and-Dialog Navigation [Thomason et al., *in sub*]
  - **New dataset**
    - **Cooperative Vision-and-Dialog Navigation (CVDN)**
    - **2k human-human dialogs**
  - Navigation from dialog history
- Next steps

# Outline

- Language grounding in visual environments
  - For navigation
  - Unimodal [L'19]
- Vision-and-D et al., *in sub*]
  - New dataset

    Initial model adapted from R2R does not exhibit unimodal overfitting.

    - Cooperative Vision-and-Dialog Navigation (CVDN)
    - 2k human-human dialogs
  - Navigation from dialog history
- Next steps

**Visible to both Navigator and Oracle**

Hint: The goal room contains a *mat*.    -- this target object is present in at least two rooms, but only one is correct.

**Visible to both Navigator and Oracle**

Hint: The goal room contains a *mat*.

-- this target object is present in at least two rooms, but only one is correct.

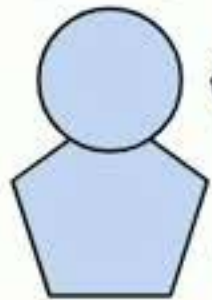**Visible to both Navigator and Oracle**

Hint: The goal room contains a *mat*.

Into the hall or the office?

-- this target object is present in at least two rooms, but only one is correct.

**Visible Only to the Oracle**

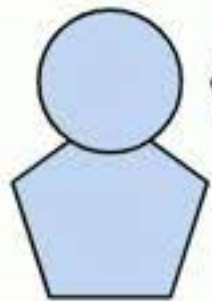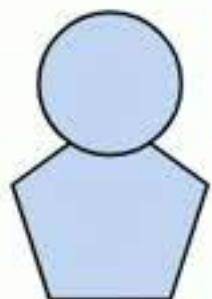-- A shortest-path planner's next steps; up to 5 navigation nodes in the direction of the goal.

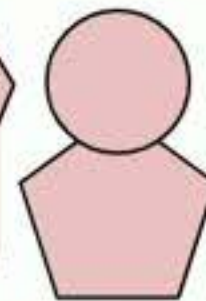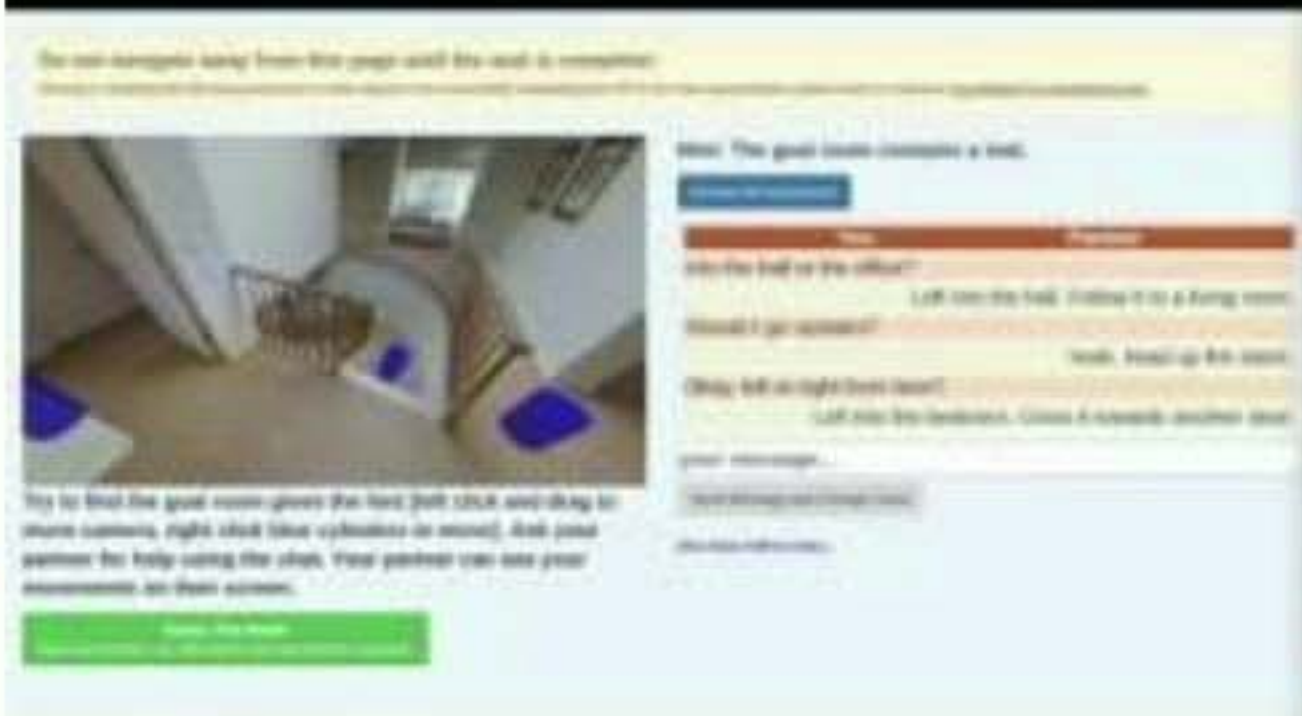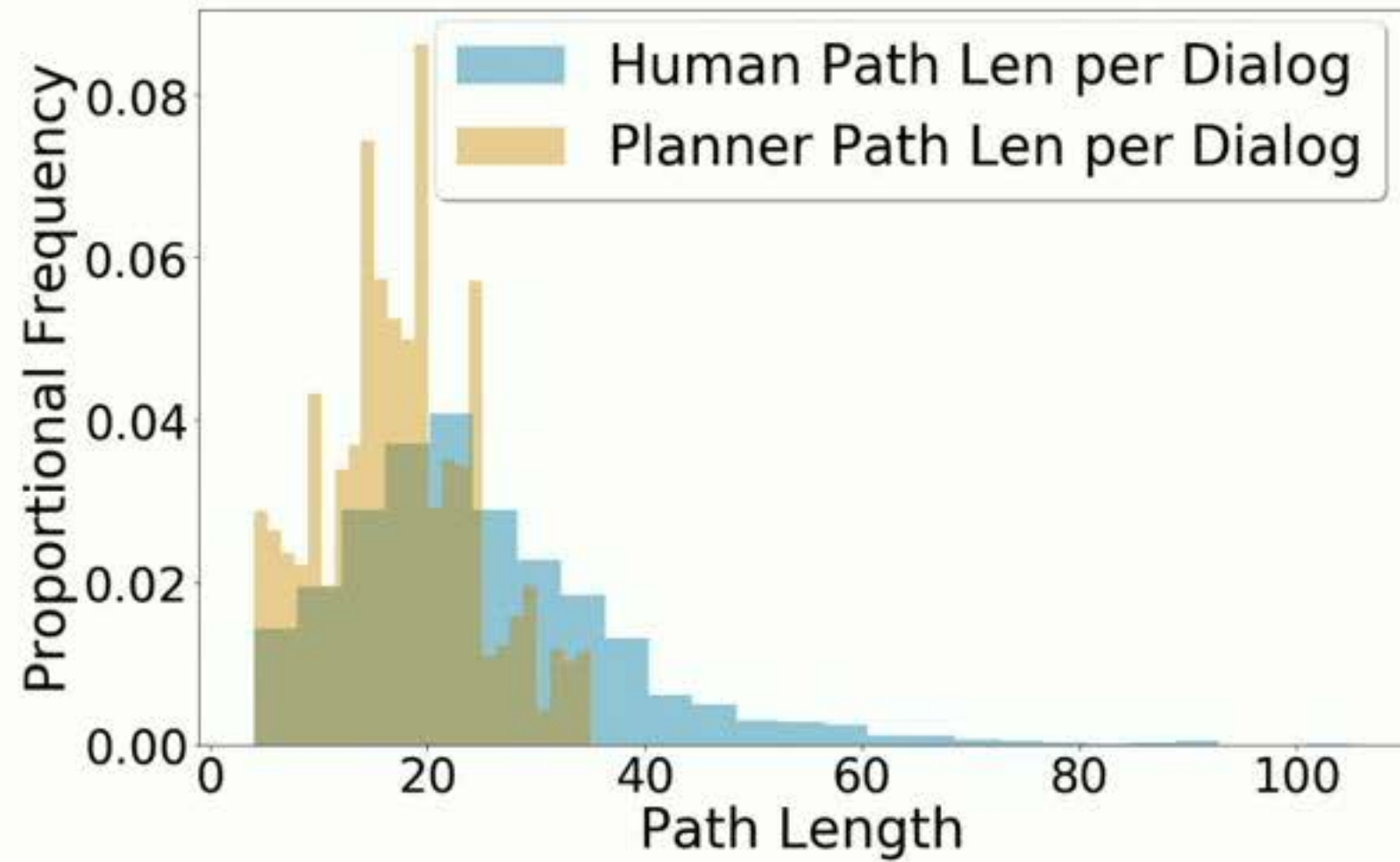**Visible to both Navigator and Oracle**

Hint: The goal room contains a *mat*.

-- this target object is present in at least two rooms, but only one is correct.

Into the hall or the office?

Left into the hall.
Follow it to a living room.

**Visible Only to the Oracle**

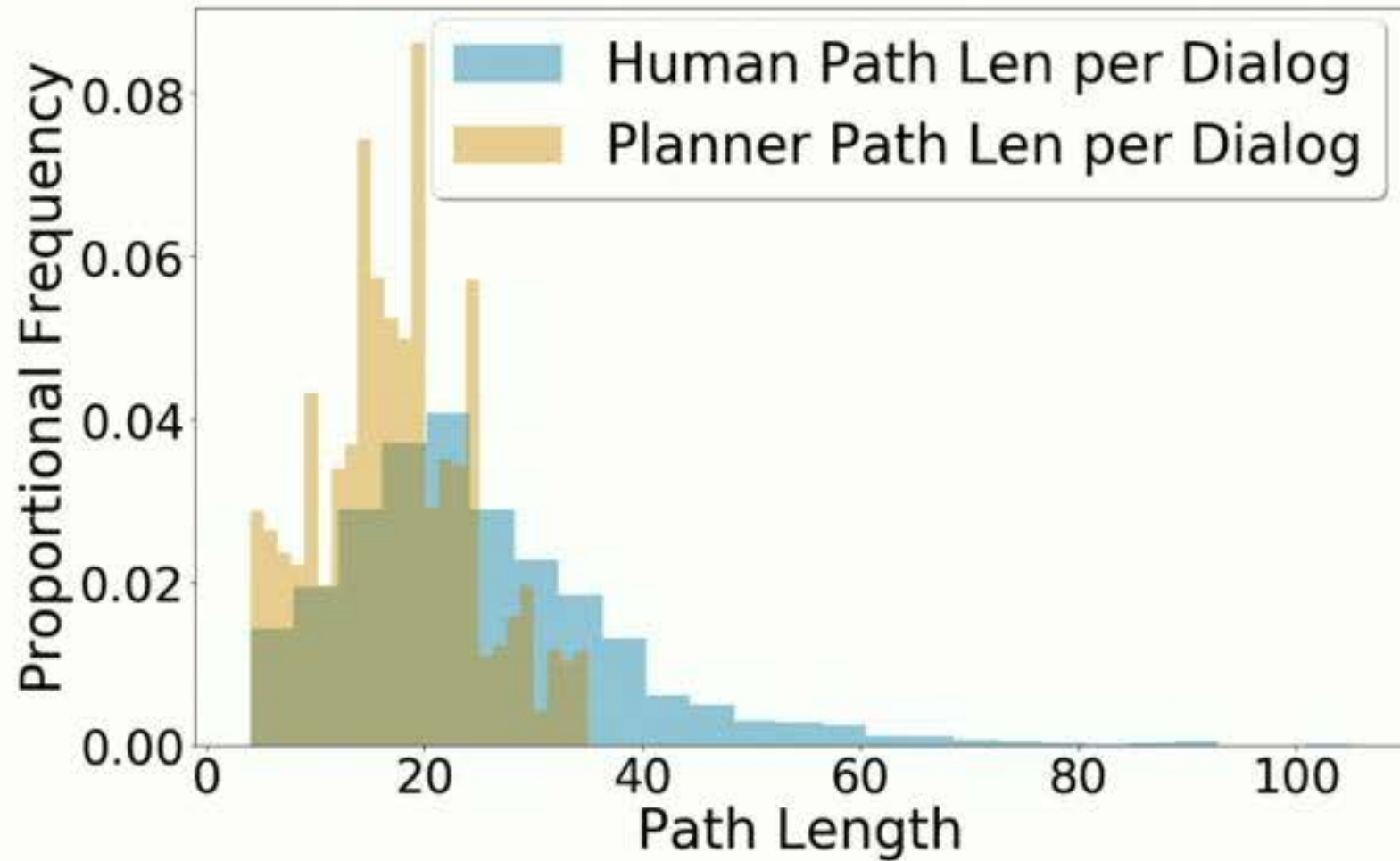-- A shortest-path planner's next steps; up to 5 navigation nodes in the direction of the goal.

**Visible to both Navigator and Oracle**

Hint: The goal room contains a *mat*.

Into the hall or the office?

Left into the hall.
Follow it to a living room.

-- this target object is present in at least two rooms, but only one is correct.

**Visible Only to the Oracle**

-- A shortest-path planner's next steps; up to 5 navigation nodes in the direction of the goal.

**Visible to both Navigator and Oracle**

Hint: The goal room contains a *mat*.

Into the hall or the office?

Left into the hall.
Follow it to a living room.

Should I go upstairs?

-- this target object is present in at least two rooms, but only one is correct.

**Visible Only to the Oracle**

-- A shortest-path planner's next steps; up to 5 navigation nodes in the direction of the goal.

32

Navigator View

Oracle View

ON AVERAGE, PATHS IN VDN ARE OVER
THREE TIMES LONGER THAN R2R PATHS.

# Dialog Leads to Long Paths and Rich Language

# Dialog Leads to Long Paths and Rich Language



- Path Length Average:
  - Human (25.0); Planner (17.4)
  - Room-to-Room (6.0)

# Dialog Leads to Long Paths and Rich Language



- **Path Length Average:**
  - Human (25.0); Planner (17.4)
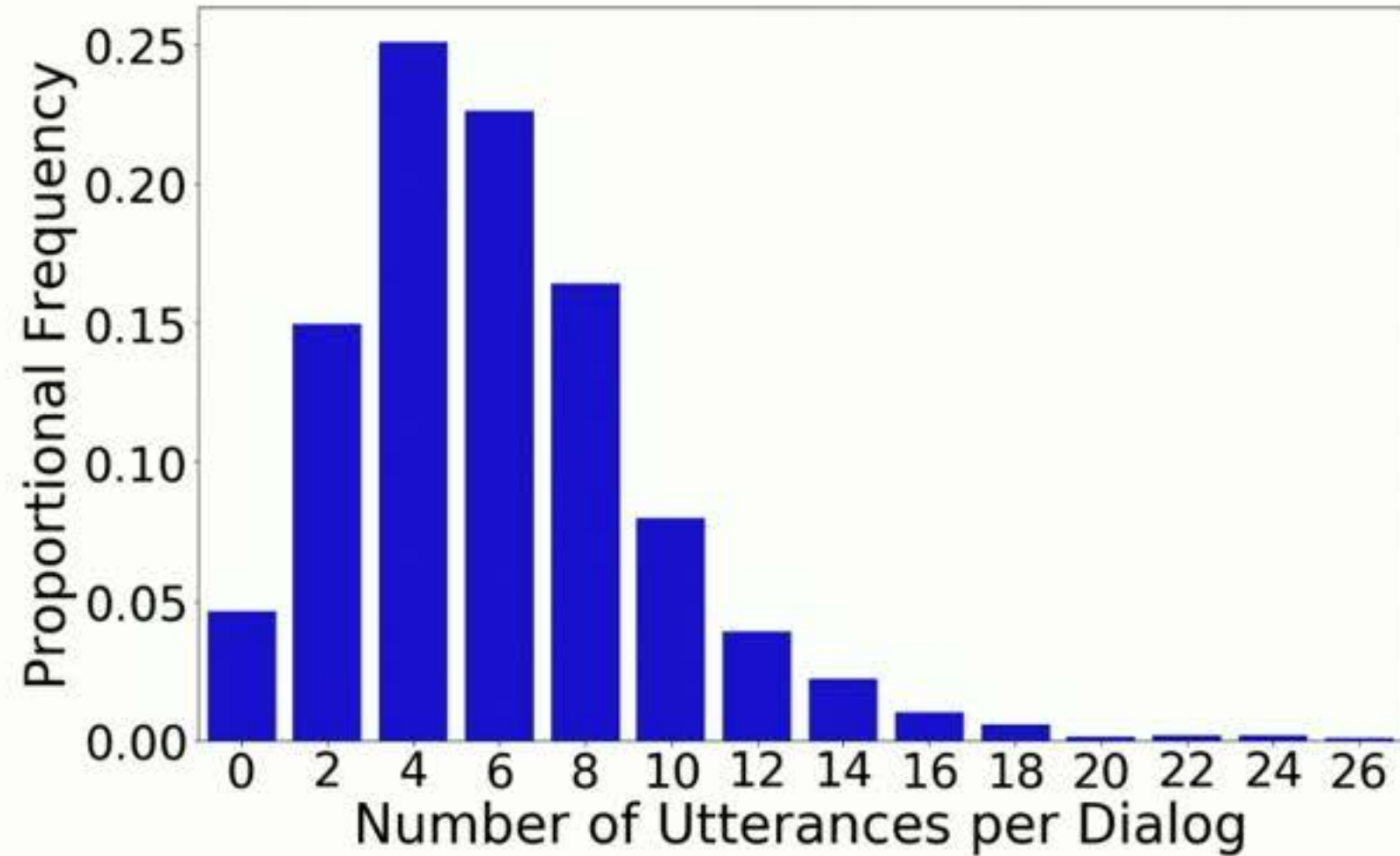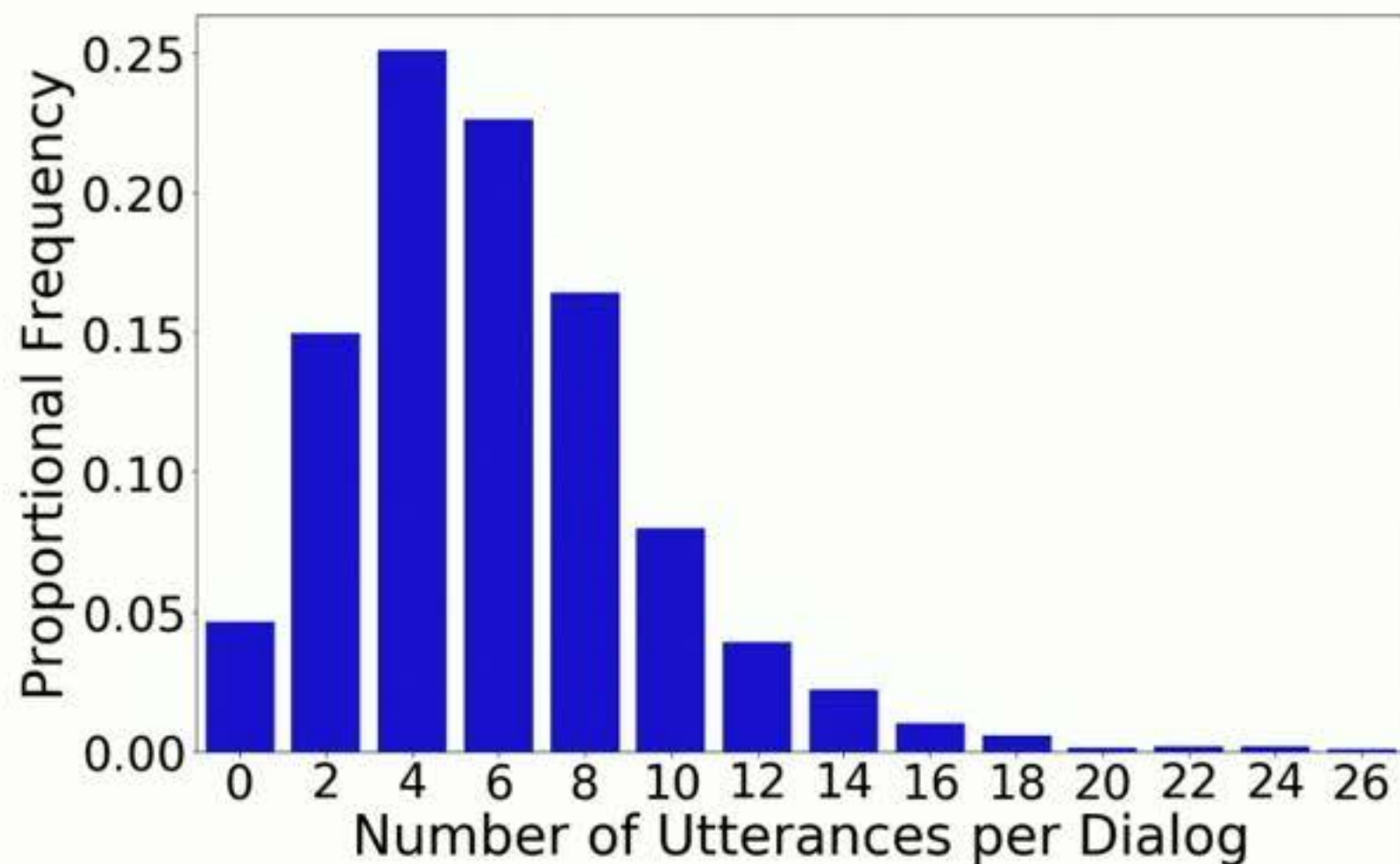  - Room-to-Room (6.0)

# Dialog Leads to Long Paths and Rich Language



- Path Length Average:
  - Human (25.0); Planner (17.4)
  - Room-to-Room (6.0)
- Average total words:
  - CVDN (82)
  - Room-to-Room (29)

# Dialog Leads to Long Paths and Rich Language

# Dialog Leads to Long Paths and Rich Language



| | |
|---|---|
| **The goal room contains a *rug*.** | ➡ 3 steps |
| ***Navigator***: Should I go to the left or right? | |
| ***Oracle***: Go left and turn right after the bathroom. | ➡ 4 steps |
| ***Navigator***: Do I need to go in the room with the run or keep on going right? | |
| ***Oracle***: Turn right and take the tiny hallway on the right. You will ascend the stairs you find on the right. | ➡ 7 steps |
| ***Navigator***: Should I go into the kitchen or to the right? | |
| ***Oracle***: Turn toward the front door and go up the stairs you see on the right. | ➡ 6 steps |
| ***Navigator***: Do I go left or right? | |
| ***Oracle***: Go along the railing to the right. Stop at the room with a brown chair. | ➡ 4s 🛑 |

35

# Dialog Leads to Long Paths and Rich Language



**Shared Visual Context → Egocentric Expressions**

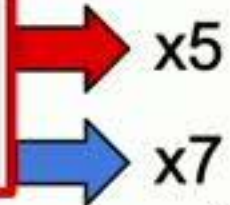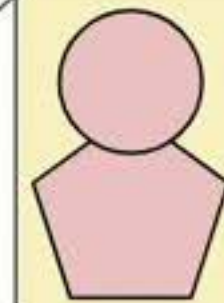| | |
|---|---|
| The goal room contains a *rug*. | ➡ 3 steps |
| *Navigator*: Should I go to the left or right? | |
| *Oracle*: Go left and turn right after the bathroom. | ➡ 4 steps |
| *Navigator*: Do I need to go in the room with the run or keep on going right? | |
| *Oracle*: Turn right and take the tiny hallway on the right. You will ascend the stairs you find on the right. | ➡ 7 steps |
| *Navigator*: Should I go into the kitchen or to the right? | |
| *Oracle*: Turn toward the front door and go up the stairs you see on the right. | ➡ 6 steps |
| *Navigator*: Do I go left or right? | |
| *Oracle*: Go along the railing to the right. Stop at the room with a brown chair. | ➡ 4s ⬣ |

36

# Outline

- Language grounding in visual environments
  - For navigation
  - Unimodal bias [Thomason et al., NAACL'19]
- Vision-and-Dialog Navigation [Thomason et al., *in sub*]
  - New dataset - CVDN
  - **Navigation from dialog history**
    - **7k dialog-based navigation inputs**
- **Next steps**

# Navigation from Dialog History

| | |
|---|---|
| **The goal room contains a *rug*.** | → 3 steps |
| *Navigator*: Should I go to the left or right? | |
| *Oracle*: Go left and turn right after the bathroom. | → 4 steps |
| *Navigator*: Do I need to go in the room with the run or keep on going right? | |
| *Oracle*: Turn right and take the tiny hallway on the right. You will ascend the stairs you find on the right. | → 7 steps |
| *Navigator*: Should I go into the kitchen or to the right? | |
| *Oracle*: Turn toward the front door and go up the stairs you see on the right. | → 6 steps |
| *Navigator*: Do I go left or right? | |
| *Oracle*: Go along the railing to the right. Stop at the room with a brown chair. | → 4s 🛑 |

# Navigation from Dialog History

| |
|---|
| **The goal room contains a *rug*.** |
| *Navigator*: Should I go to the left or right? |
| *Oracle*: Go left and turn right after the bathroom. |
| *Navigator*: Do I need to go in the room with the run or keep on going right? |
| *Oracle*: Turn right and take the tiny hallway on the right. You will ascend the stairs you find on the right. |

**x5**

**x7**

**Oracle**



**Navigator**

# Navigation from Dialog History

The goal room contains a *rug*.

**Navigator**: Should I go to the left or right?

**Oracle**: Go left and turn right after the bathroom.

**Navigator**: Do I need to go in the room with the run or keep on going right?

**Oracle**: Turn right and take the tiny hallway on the right. You will ascend the stairs you find on the right.

**Navigator**: Should I go into the kitchen or to the right?

**Oracle**: Turn toward the front door and go up the stairs you see on the right.

**Navigator**: Do I go left or right?

**Oracle**: Go along the railing to the right. Stop at the room with a brown chair.

# Navigation from Dialog History

| |
|---|
| **The goal room contains a *rug*.** |
| ***Navigator***: Should I go to the left or right? |
| ***Oracle***: Go left and turn right after the bathroom. |
| ***Navigator***: Do I need to go in the room with the run or keep on going right? |
| ***Oracle***: Turn right and take the tiny hallway on the right. You will ascend the stairs you find on the right. |
| ***Navigator***: Should I go into the kitchen or to the right? |
| ***Oracle***: Turn toward the front door and go up the stairs you see on the right. |
| ***Navigator***: Do I go left or right? |
| ***Oracle***: Go along the railing to the right. Stop at the room with a brown chair. |

➡️ 3 steps

# Navigation from Dialog History

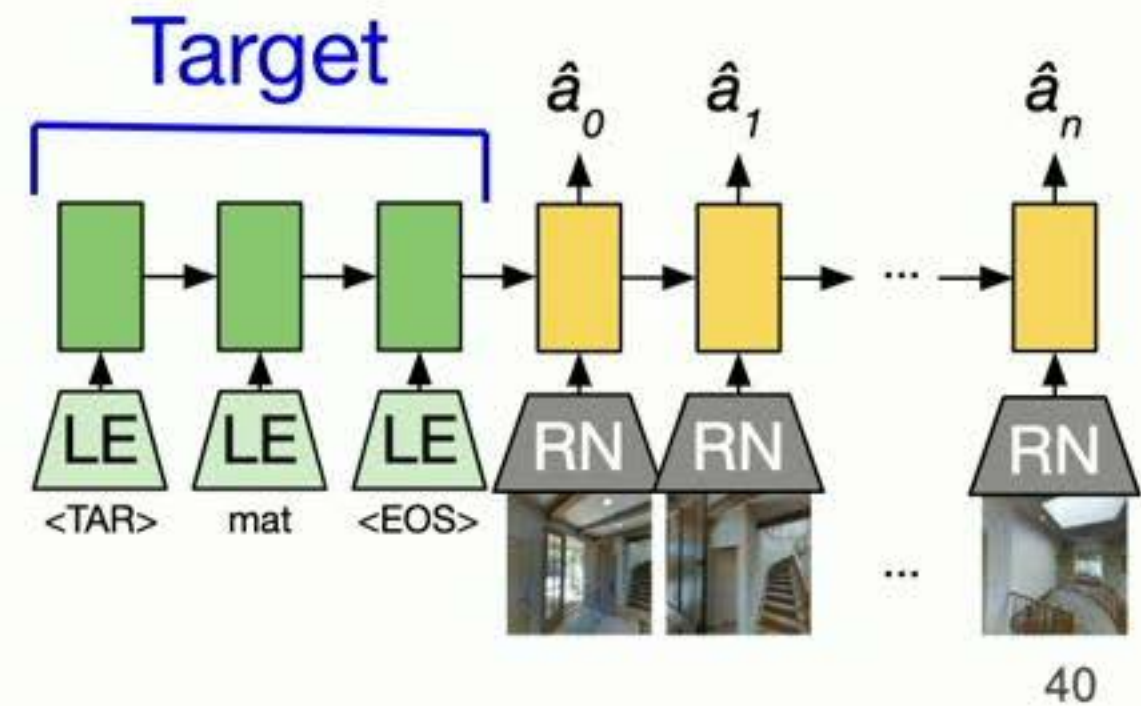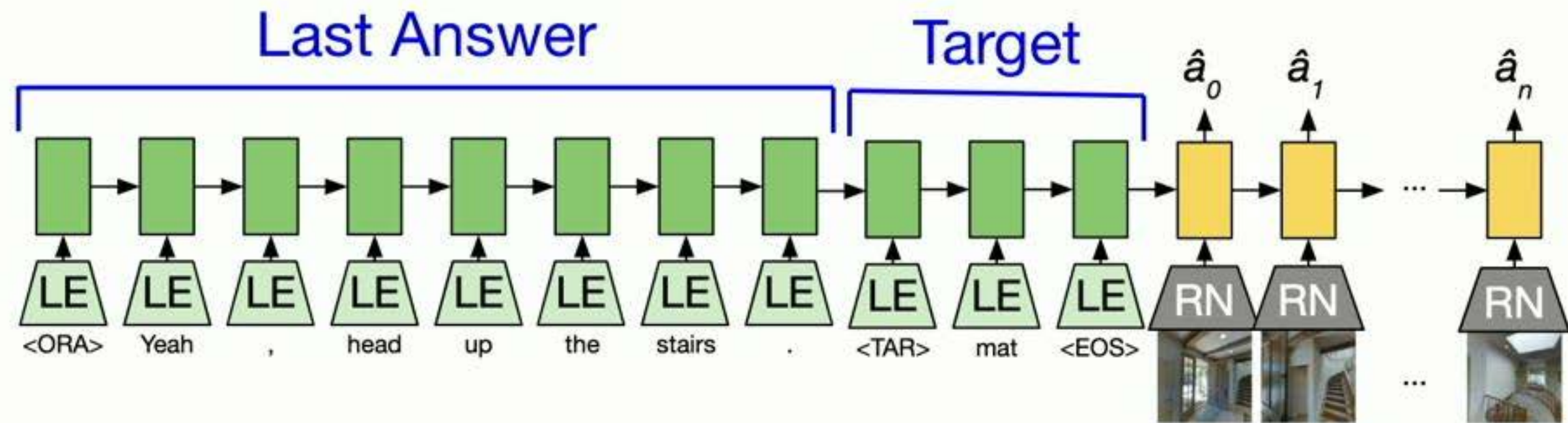| |
|---|
| **The goal room contains a *rug*.** |
| ***Navigator***: Should I go to the left or right? |
| ***Oracle***: Go left and turn right after the bathroom. |
| ***Navigator***: Do I need to go in the room with the run or keep on going right? |
| ***Oracle***: Turn right and take the tiny hallway on the right. You will ascend the stairs you find on the right. |
| ***Navigator***: Should I go into the kitchen or to the right? |
| ***Oracle***: Turn toward the front door and go up the stairs you see on the right. |
| ***Navigator***: Do I go left or right? |
| ***Oracle***: Go along the railing to the right. Stop at the room with a brown chair. |

➡ 3 steps

➡ 4 steps

# Navigation from Dialog History

| |
|---|
| **The goal room contains a *rug*.** |
| *Navigator*: Should I go to the left or right? |
| *Oracle*: Go left and turn right after the bathroom. |
| *Navigator*: Do I need to go in the room with the run or keep on going right? |
| *Oracle*: Turn right and take the tiny hallway on the right. You will ascend the stairs you find on the right. |
| *Navigator*: Should I go into the kitchen or to the right? |
| *Oracle*: Turn toward the front door and go up the stairs you see on the right. |
| *Navigator*: Do I go left or right? |
| *Oracle*: Go along the railing to the right. Stop at the room with a brown chair. |

⇨ 3 steps

⇨ 4 steps

- **Input:** History so far + visual frame per timestep.

# Navigation from Dialog History

| |
|---|
| **The goal room contains a *rug*.** |
| *Navigator*: Should I go to the left or right? |
| *Oracle*: Go left and turn right after the bathroom. |
| *Navigator*: Do I need to go in the room with the run or keep on going right? |
| *Oracle*: Turn right and take the tiny hallway on the right. You will ascend the stairs you find on the right. |
| *Navigator*: Should I go into the kitchen or to the right? |
| *Oracle*: Turn toward the front door and go up the stairs you see on the right. |
| *Navigator*: Do I go left or right? |
| *Oracle*: Go along the railing to the right. Stop at the room with a brown chair. |

➡ 3 steps

➡ 4 steps

- **Input:** History so far + visual frame per timestep.
- **Output**: Navigation action per timestep.

# Navigation from Dialog History

| |
|---|
| The goal room contains a *rug*. |
| *Navigator*: Should I go to the left or right? |
| *Oracle*: Go left and turn right after the bathroom. |
| *Navigator*: Do I need to go in the room with the run or keep on going right? |
| *Oracle*: Turn right and take the tiny hallway on the right. You will ascend the stairs you find on the right. |
| *Navigator*: Should I go into the kitchen or to the right? |
| *Oracle*: Turn toward the front door and go up the stairs you see on the right. |
| *Navigator*: Do I go left or right? |
| *Oracle*: Go along the railing to the right. Stop at the room with a brown chair. |

➡ 3 steps

➡ 4 steps

- **Input:** History so far + visual frame per timestep.
- **Output**: Navigation action per timestep.
- **Goal**: Get closer to the target object room.

39

# Navigation from Dialog History

| |
|---|
| **The goal room contains a *rug*.** |
| *Navigator*: Should I go to the left or right? |
| *Oracle*: Go left and turn right after the bathroom. |
| *Navigator*: Do I need to go in the room with the run or keep on going right? |
| *Oracle*: Turn right and take the tiny hallway on the right. You will ascend the stairs you find on the right. |
| *Navigator*: Should I go into the kitchen or to the right? |
| *Oracle*: Turn toward the front door and go up the stairs you see on the right. |
| *Navigator*: Do I go left or right? |
| *Oracle*: Go along the railing to the right. Stop at the room with a brown chair. |

⟹ 3 steps

⟹ 4 steps

- **Input:** History so far + visual frame per timestep.
- **Output**: Navigation action per timestep.
- **Goal**: Get closer to the target object room.
- 2k dialogs → 7k histories.

# Navigation from Dialog History

| | |
|---|---|
| **The goal room contains a** *rug*. | ➡ 3 steps |
| *Navigator*: Should I go to the left or right? | |
| *Oracle*: Go left and turn right after the bathroom. | ➡ 4 steps |
| *Navigator*: Do I need to go in the room with the run or keep on going right? | |
| *Oracle*: Turn right and take the tiny hallway on the right. You will ascend the stairs you find on the right. | ➡ 7 steps |
| *Navigator*: Should I go into the kitchen or to the right? | |
| *Oracle*: Turn toward the front door and go up the stairs you see on the right. | ➡ 6 steps |
| *Navigator*: Do I go left or right? | |
| *Oracle*: Go along the railing to the right. Stop at the room with a brown chair. | ➡ 4 steps |

- **Input:** History so far + visual frame per timestep.
- **Output**: Navigation action per timestep.
- **Goal**: Get closer to the target object room.
- 2k dialogs → 7k histories.

39

# Initial, Sequence-to-Sequence Model

# Initial, Sequence-to-Sequence Model



Last Answer      Target

$\hat{a}_0$   $\hat{a}_1$   $\hat{a}_n$

LE LE LE LE LE LE LE LE   LE LE LE   RN RN   ...   RN

&lt;ORA&gt; Yeah , head up the stairs .   &lt;TAR&gt; mat &lt;EOS&gt;

40

# Initial, Sequence-to-Sequence Model

*Navigator*: Should I turn left down the hallway ahead?
*Oracle*: ya

# Initial, Sequence-to-Sequence Model

*Navigator*: Should I turn left down the hallway ahead?
*Oracle*: ya

# Initial, Sequence-to-Sequence Model

*Oracle*: Through the lobby. So go through the door next to the green towel. Go to the left door next to the two yellow lights. Walk straight to the end of the hallway and stop

…

*Navigator*: Are these the yellow lights you were talking about?

# Initial, Sequence-to-Sequence Model



**All Previous Questions and Answers**

LE &lt;NAV&gt; | LE Into | LE the | LE hall | LE or | LE the | LE office | LE ? | LE &lt;ORA&gt; | LE Left | LE into | LE the | LE hall | LE . | LE Follow | LE it | LE to | LE a | LE living | LE room | LE .

**Last Question**      **Last Answer**      **Target**

$\hat{a}_0$   $\hat{a}_1$   $\hat{a}_n$

LE &lt;NAV&gt; | LE Should | LE I | LE go | LE upstairs | LE ? | LE &lt;ORA&gt; | LE Yeah | LE , | LE head | LE up | LE the | LE stairs | LE . | LE &lt;TAR&gt; | LE mat | LE &lt;EOS&gt; | RN | RN | ... | RN

# Evaluation - Validation (best perf over 200 epochs)



- **Seen Environments:**
  - Novel dialogs.
  - Houses seen at training time.

# Initial, Sequence-to-Sequence Model

**All Previous Questions and Answers**



<NAV> Into the hall or the office ? <ORA> Left into the hall . Follow it to a living room .

**Last Question**      **Last Answer**      **Target**

$\hat{a}_0$    $\hat{a}_1$    $\hat{a}_n$



<NAV> Should I go upstairs ? <ORA> Yeah , head up the stairs . <TAR> mat <EOS>

# Evaluation - Validation (best perf over 200 epochs)



- **Seen Environments:**
  - Novel dialogs.
  - Houses seen at training time.

# Evaluation - Validation (best perf over 200 epochs)



- **Seen Environments:**
  - Novel dialogs.
  - Houses seen at training time.

41

# Evaluation - Validation (best perf over 200 epochs)



- **Seen Environments:**
  - Novel dialogs.
  - Houses seen at training time.

# Evaluation - Validation (best perf over 200 epochs)



Bar chart — Goal Progress (m) vs Dialog History Encoded:
- target object: 5.71
- + last answer: 6.04
- + last question: 6.16
- + all

- **Seen Environments:**
  - Novel dialogs.
  - Houses seen at training time.

# Evaluation - Validation (best perf over 200 epochs)

- **Seen Environments:**
  - Novel dialogs.
  - Houses seen at training time.

41

# Evaluation - Validation (best perf over 200 epochs)



- **Seen Environments:**
  - Novel dialogs.
  - Houses seen at training time.

# Evaluation - Test (epoch of best Val Unseen)



- **Unseen Envs:**
  - Novel dialogs.
  - Novel houses not seen during training.

# Evaluation - Test (epoch of best Val Unseen)



- **Unseen Envs:**
  - Novel dialogs.
  - Novel houses not seen during training.

# Evaluation - Test (epoch of best Val Unseen)



- **Unseen Envs:**
  - Novel dialogs.
  - Novel houses not seen during training.

42

# Evaluation - Test (epoch of best Val Unseen)



- **Unseen Envs:**
  - Novel dialogs.
  - Novel houses not seen during training.

# Evaluation - Test (epoch of best Val Unseen)



- **Unseen Envs:**
  - Novel dialogs.
  - Novel houses not seen during training.

# Evaluation - Test (epoch of best Val Unseen)



- **Unseen Envs:**
  - Novel dialogs.
  - Novel houses not seen during training.

# Evaluation - Unimodal Baselines

# Evaluation - Unimodal Baselines

# Evaluation - Unimodal Baselines

# Evaluation - Unimodal Baselines

Action-only

# Evaluation - Unimodal Baselines

Action-
only

Vis-
only



43

# Evaluation - Unimodal Baselines



43

# Evaluation - Unimodal Baselines



43

# Evaluation - Unimodal Baselines

Action-
only

Vis-
only

Lang-
only



43

# Evaluation - Unimodal Baselines



43

# Evaluation - Unimodal Baselines

Action-
only

Vis-
only

Lang-
only

Goal Progress (m) vs Model

Legend:
- Validation (Seen)
- Test (Unseen)

| Model | Validation (Seen) | Test (Unseen) |
|---|---|---|
| Action-only | 0.91 | 0.52 |
| Vis-only | 5.72 | 1.74 |
| Lang-only | 1.58 | 1.40 |
| Lang+Vis | 5.92 | 2.35 |
| Best Possible | | |

# Evaluation - Unimodal Baselines



43

# Evaluation - Unimodal Baselines



43

# Evaluation Lessons

# Evaluation Lessons

- Dialog history:

# Evaluation Lessons

- Dialog history:
    - Longer context leads to better performance.
    - Particularly helpful in unseen environments.

# Evaluation Lessons

- Dialog history:
  - Longer context leads to better performance.
  - Particularly helpful in unseen environments.
- Unimodal baselines:
  - Initial model makes use of multimodal information

# Evaluation Lessons

- Dialog history:
  - Longer context leads to better performance.
  - Particularly helpful in unseen environments.
- Unimodal baselines:
  - Initial model makes use of multimodal information

# Evaluation Lessons

- Dialog history:
  - Longer context leads to better performance.
  - Particularly helpful in unseen environments.
- Unimodal baselines:
  - Initial model makes use of multimodal information
  - Multimodal most helpful in unseen environments.

# Evaluation Lessons

- Dialog history:
  - Longer context leads to better performance.
  - Particularly helpful in unseen environments.
- Unimodal baselines:
  - Initial model makes use of multimodal information
  - Multimodal most helpful in unseen environments.
- Headroom remains for more nuanced models.

# Outline

- Language grounding in visual environments
  - For navigation
  - Unimodal bias [Thomason et al., NAACL'19]
- Vision-and-Dialog Navigation [Thomason et al., *in sub*]
  - New dataset - CVDN
  - Navigation from dialog history
- **Next steps**

# Incorporating Navigation History

# Incorporating Navigation History

*Oracle*: You were there briefly but left. There is a turntable behind you a bit. Enter the bedroom next to it.



<NAV> Into the hall or the office ? <ORA> Left into the hall . Follow it to a living room .

$\hat{a}_0$   $\hat{a}_1$   $\hat{a}_n$

<NAV> Should I go upstairs ? <ORA> Yeah , head up the stairs . <TAR> mat <EOS>

# Incorporating Navigation History

# Incorporating Navigation History

Visible to both Navigator and Oracle

Hint: The goal room contains a *mat*.

Into the hall or the office?
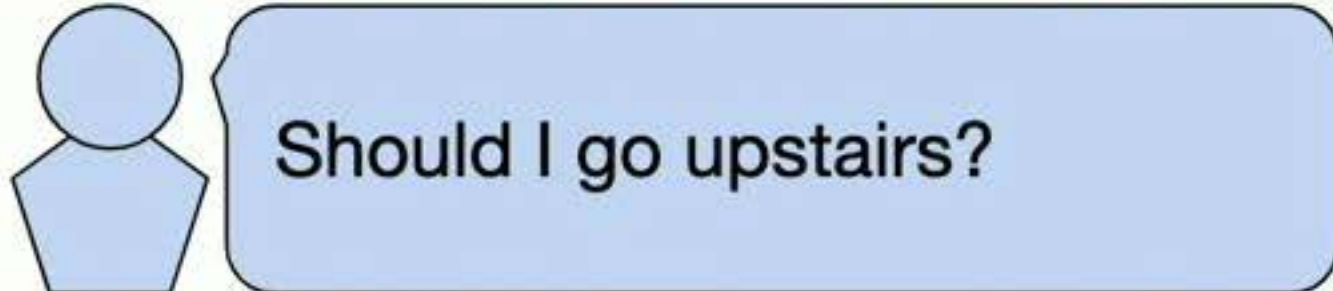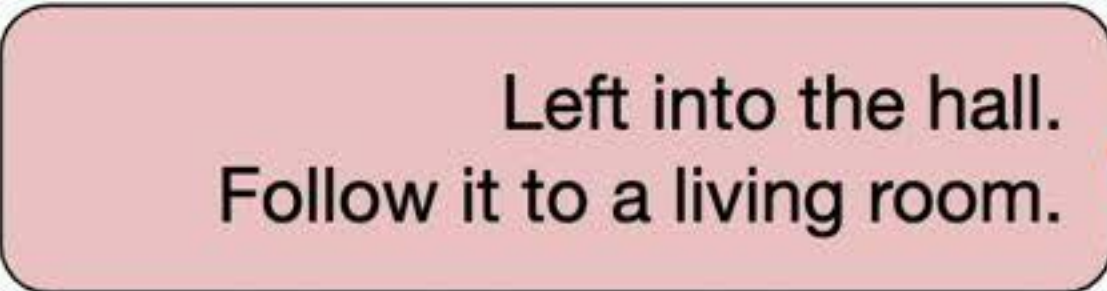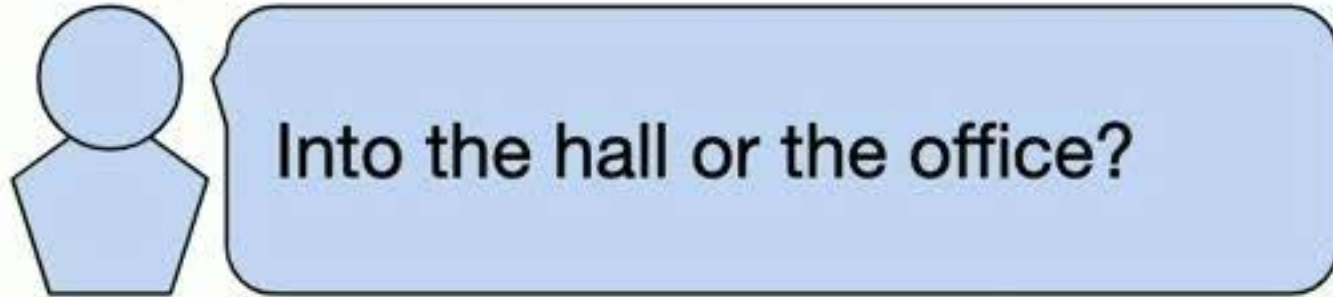
Left into the hall.
Follow it to a living room.

Visible Only to the Oracle

Should I go upstairs?

48

**Visible to both Navigator and Oracle**

Hint: The goal room contains a *mat*.

Into the hall or the office?
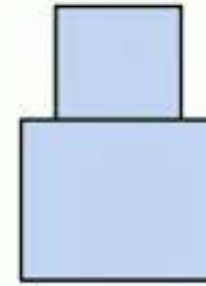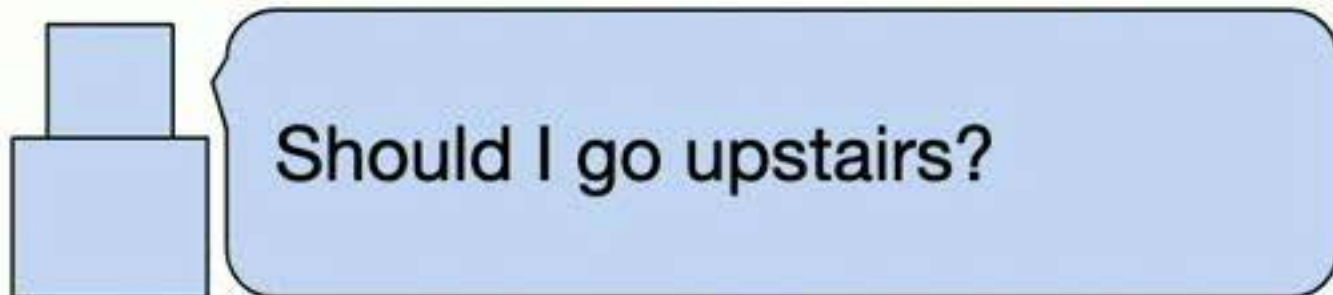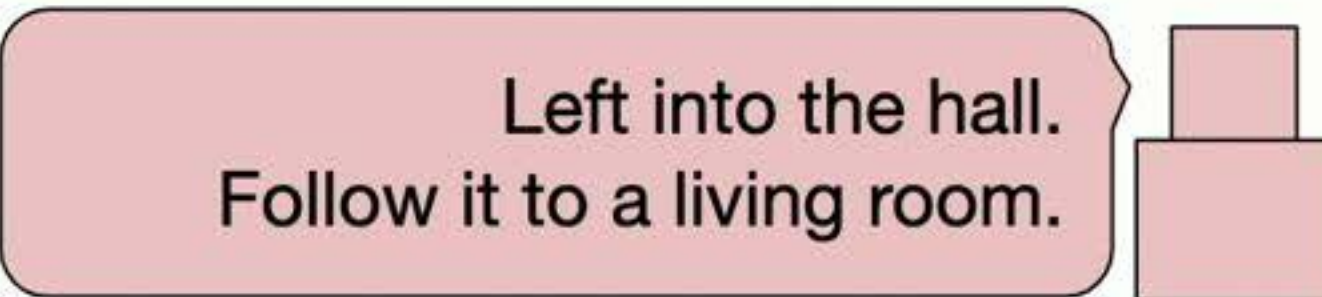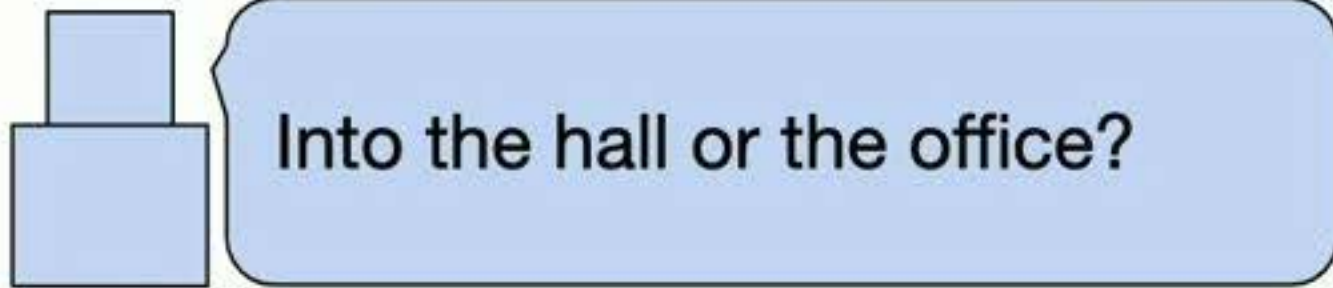
Left into the hall.
Follow it to a living room.

Should I go upstairs?

Navigation

Navigation
Question
Generation

**Visible Only to the Oracle**

Question
Answering

49

Hint: The goal room contains a *mat*.

Navigation

Question Generation

Question Answering

Hint: The goal room contains a *mat*.

Navigation

Question
Generation

● Environment exploration

Question
Answering

Hint: The goal room contains a *mat*.

Navigation

Question
Generation

- Environment exploration
- Self-play
- "Language" evolution

Question
Answering

Hint: The goal room contains a *mat*.

Navigation

Question Generation

- Environment exploration
- Self-play
- "Language" evolution
- Reinforcement Learning

Question Answering

# Bringing Robots from Industrial to Human Spaces

Industrial





50

# Bringing Robots from Industrial to Human Spaces

Industrial

Human

50

# Bringing Robots from Industrial to Human Spaces



50

# Bringing Robots from Industrial to Human Spaces



Industrial

Human

Natural Language →

Navigation →

Additional Safety →

Robust Perception →

⋮

50

# Takeaways

# Takeaways

- For vision-and-language navigation:

# Takeaways

- For vision-and-language navigation:
  - Unimodal ablations expose dataset bias.

# Takeaways

- For vision-and-language navigation:
  - Unimodal ablations expose dataset bias.
- For vision-and-dialog navigation:

# Takeaways

- For vision-and-language navigation:
  - Unimodal ablations expose dataset bias.
- For vision-and-dialog navigation:
  - Cooperative dialog facilitates a mix of high- and low-level language.

# Takeaways

- For vision-and-language navigation:
  - Unimodal ablations expose dataset bias.
- For vision-and-dialog navigation:
  - Cooperative dialog facilitates a mix of high- and low-level language.
  - Dialog context helps agents infer better navigation actions.

# Vision-and-Dialog Navigation

# Vision-and-Dialog Navigation

- Jesse Thomason, Daniel Gordon, and Yonatan Bisk.

  "Shifting the Baseline: Single Modality Performance on Visual Navigation and QA", NAACL'19 [https://arxiv.org/abs/1811.00613]

# Vision-and-Dialog Navigation

- Jesse Thomason, Daniel Gordon, and Yonatan Bisk.

  "Shifting the Baseline: Single Modality Performance on Visual Navigation and QA", NAACL'19 [https://arxiv.org/abs/1811.00613]

- Jesse Thomason, Michael Murray, Maya Cakmak, and Luke Zettlemoyer.

  "Vision-and-Dialog Navigation", *in submission* [https://arxiv.org/abs/1907.04957]

# Vision-and-Dialog Navigation

- Jesse Thomason, Daniel Gordon, and Yonatan Bisk.

  "Shifting the Baseline: Single Modality Performance on Visual Navigation and QA", NAACL'19 [https://arxiv.org/abs/1811.00613]

- Jesse Thomason, Michael Murray, Maya Cakmak, and Luke Zettlemoyer.

  "Vision-and-Dialog Navigation", *in submission* [https://arxiv.org/abs/1907.04957]

- Cooperative Vision-and-Dialog Navigation dataset

# Vision-and-Dialog Navigation

- Jesse Thomason, Daniel Gordon, and Yonatan Bisk.

  "Shifting the Baseline: Single Modality Performance on Visual Navigation and QA", NAACL'19 [https://arxiv.org/abs/1811.00613]

- Jesse Thomason, Michael Murray, Maya Cakmak, and Luke Zettlemoyer.

  "Vision-and-Dialog Navigation", *in submission* [https://arxiv.org/abs/1907.04957]

- Cooperative Vision-and-Dialog Navigation dataset
  - Data+code: https://github.com/mmurray/cvdn/

# Vision-and-Dialog Navigation

- Jesse Thomason, Daniel Gordon, and Yonatan Bisk.
  "Shifting the Baseline: Single Modality Performance on Visual Navigation and QA", NAACL'19 [https://arxiv.org/abs/1811.00613]

- Jesse Thomason, Michael Murray, Maya Cakmak, and Luke Zettlemoyer.
  "Vision-and-Dialog Navigation", *in submission* [https://arxiv.org/abs/1907.04957]

- **Cooperative Vision-and-Dialog Navigation dataset**
  - Data+code: https://github.com/mmurray/cvdn/
  - Live demo: https://cvdn.dev/

# Vision-and-Dialog Navigation

- Jesse Thomason, Daniel Gordon, and Yonatan Bisk.

  "Shifting the Baseline: Single Modality Performance on Visual Navigation and QA", NAACL'19 [https://arxiv.org/abs/1811.00613]

- Jesse Thomason, Michael Murray, Maya Cakmak, and Luke Zettlemoyer.

  "Vision-and-Dialog Navigation", *in submission* [https://arxiv.org/abs/1907.04957]

- Cooperative Vision-and-Dialog Navigation dataset
  - Data+code: https://github.com/mmurray/cvdn/
  - Live demo: https://cvdn.dev/

- Navigation from Dialog History task

# Vision-and-Dialog Navigation

- Jesse Thomason, Daniel Gordon, and Yonatan Bisk.

  "Shifting the Baseline: Single Modality Performance on Visual Navigation and QA", NAACL'19 [https://arxiv.org/abs/1811.00613]

- Jesse Thomason, Michael Murray, Maya Cakmak, and Luke Zettlemoyer.

  "Vision-and-Dialog Navigation", *in submission* [https://arxiv.org/abs/1907.04957]

- ## Cooperative Vision-and-Dialog Navigation dataset
  - ### Data+code: https://github.com/mmurray/cvdn/
  - ### Live demo: https://cvdn.dev/

- ## Navigation from Dialog History task
  - ### Leaderboard coming soon!

# Future Work: More Expressive Simulator



"*But a slice of bread in the microwave.*"

- **High-level instructions.**
  - ○ **With optional, accompanying low-level.**
- **Navigation + manipulation + interaction.**
- **No dialog to support clarifications.**