

Learning to Navigate Unseen Environments: Back Translation with Environmental Dropout

Hao Tan, Licheng Yu, Mohit Bansal

haotan, licheng, mbansal@cs.unc.edu

UNC, Chapel Hill

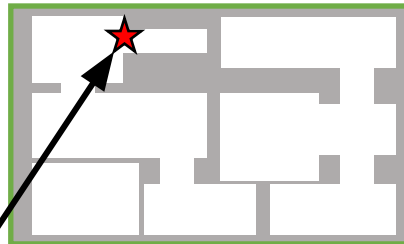


Vision-and-Language Navigation Task

Instruction

Go to the bedroom, and go through the door, continue forward until you can climb three steps to your right...

Bird-View



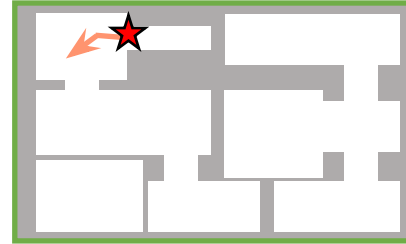
Agent's Start Location

Vision-and-Language Navigation Task

Instruction

Go to the bedroom, and go through the door, continue forward until you can climb three steps to your right...

Bird-View

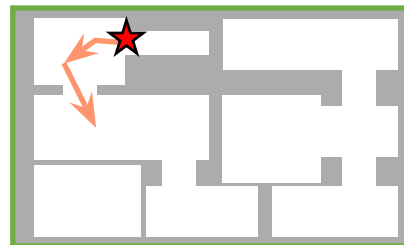


Vision-and-Language Navigation Task

Instruction

Go to the bedroom, and go through the door, continue forward until you can climb three steps to your right...

Bird-View

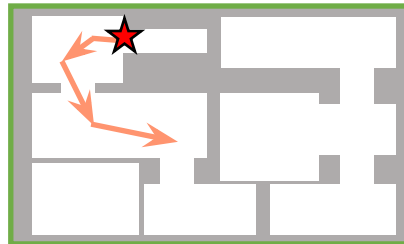


Vision-and-Language Navigation Task

Instruction

Go to the bedroom, and go through the door, continue forward until you can climb three steps to your right...

Bird-View

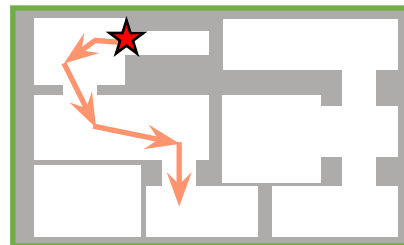


Vision-and-Language Navigation Task

Instruction

Go to the bedroom, and go through the door, continue forward until you can climb three steps to your right...

Bird-View

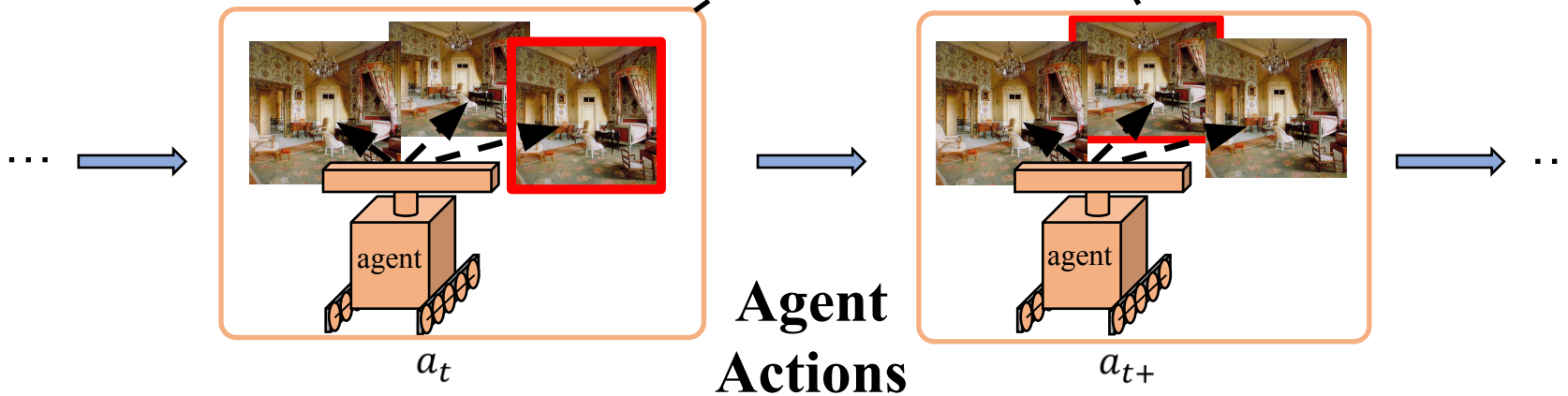
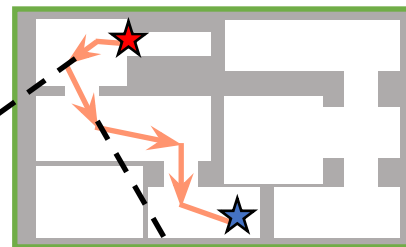


Vision-and-Language Navigation Task

Instruction

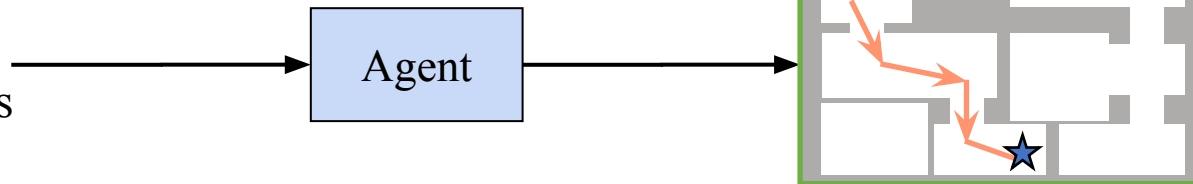
Go to the bedroom, and go through the door, continue forward until you can climb three steps to your right...

Bird-View



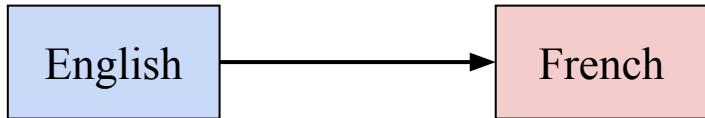
Vision-and-Language Navigation Task

Go to the bedroom, and go through the door, continue forward until you can climb three steps to your right...



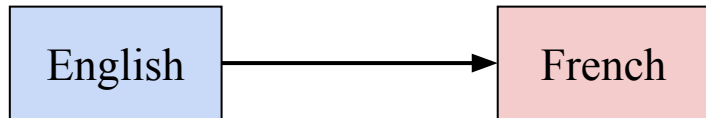
Back Translation

1. Want to learn: En \rightarrow Fr



Back Translation

1. Want to learn: En \rightarrow Fr

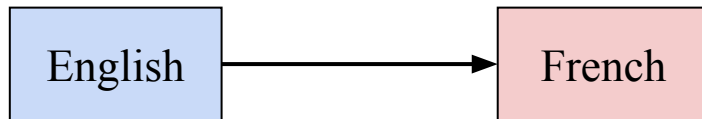


2. Have unpaired Fr corpus



Back Translation

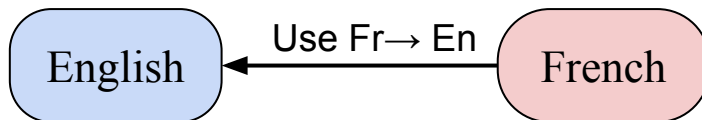
1. Want to learn: En \rightarrow Fr



2. Have unpaired Fr corpus

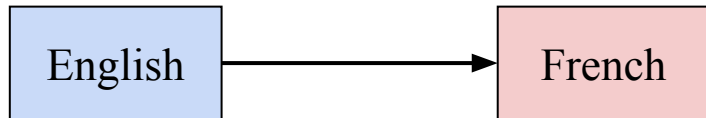


3. Train Fr \rightarrow En and use it to translate unpaired Fr corpus.



Back Translation

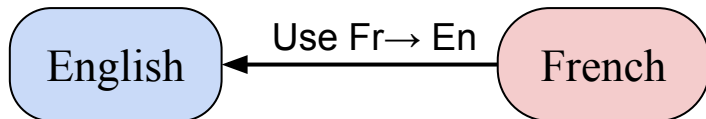
1. Want to learn: $En \rightarrow Fr$.



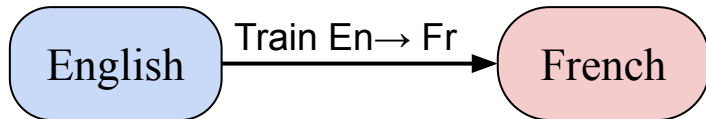
2. Have unpaired Fr corpus.



3. Train $Fr \rightarrow En$ and use it to translate unpaired Fr corpus.

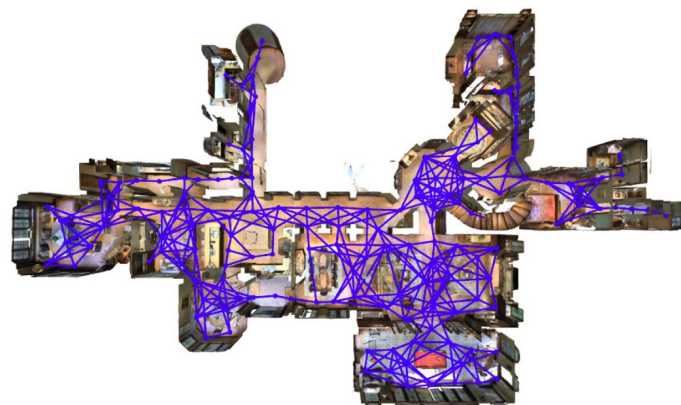


4. Use reversed pairs as additional data for $En \rightarrow Fr$.



Back Translation: Preliminary Setup

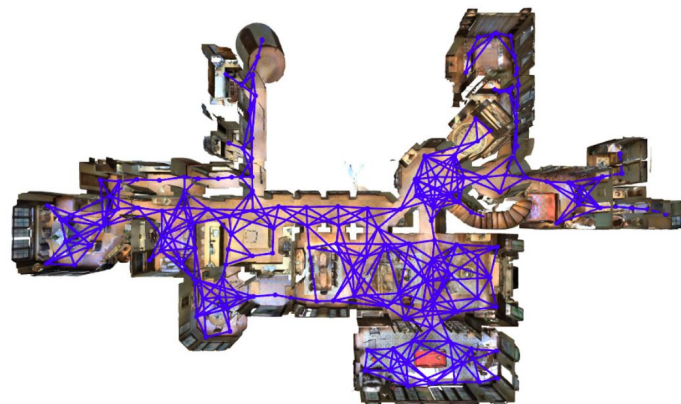
Environment: A set of routes. Some routes have instructions; Some do not.



Back Translation: Preliminary Setup

Environment: A set of routes. Some routes have instructions; Some do not.

Speaker: A pre-trained neural model which generates instructions from routes.

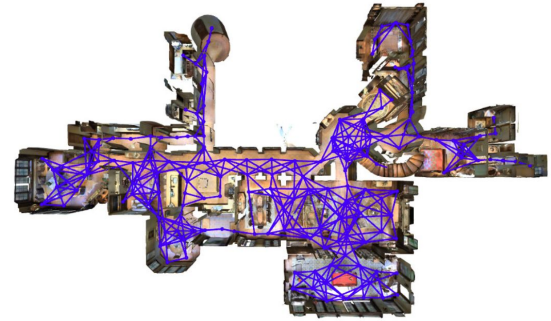


Speaker

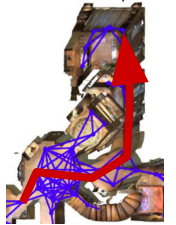
*Walk past
the hall,
turn Left, ...*

Back Translation: Step 1

1. New routes from existing environments.

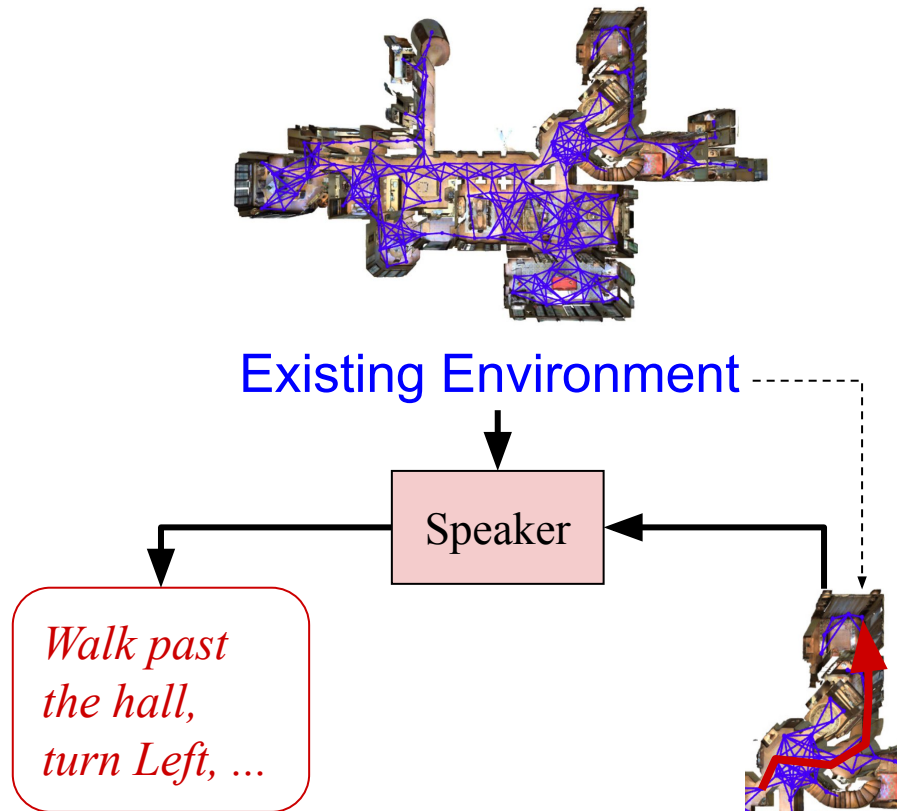


Existing Environment - - -



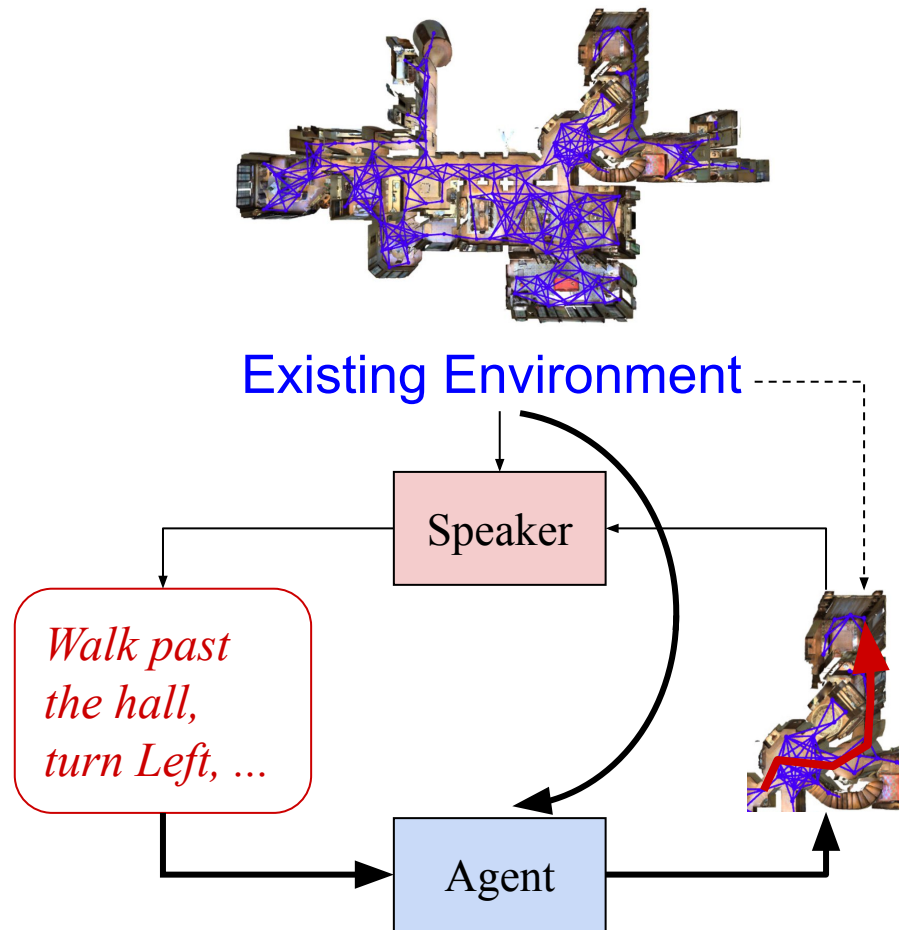
Back Translation: Step 2

1. **New routes** from existing environments.
2. **New instructions** by **pre-trained speaker**.



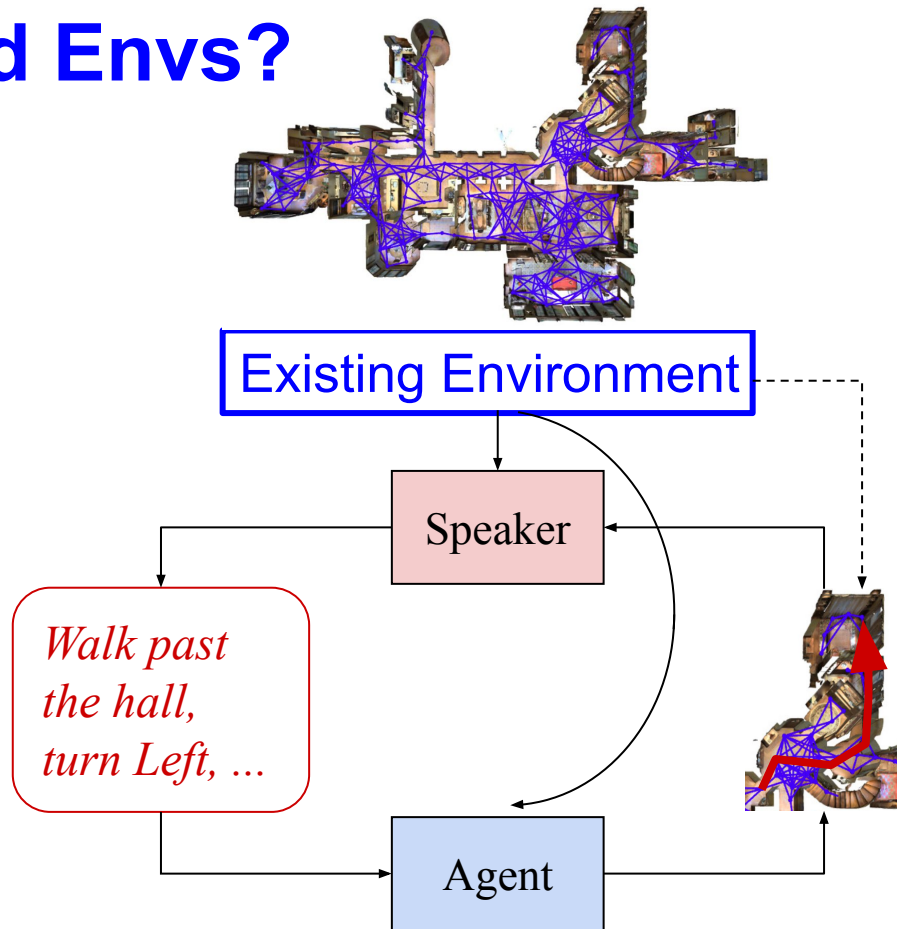
Back Translation: Step 3

1. **New routes** from existing environments.
2. **New instructions** by pre-trained speaker.
3. Train agent on **new routes**, **new instructions**, existing environments.



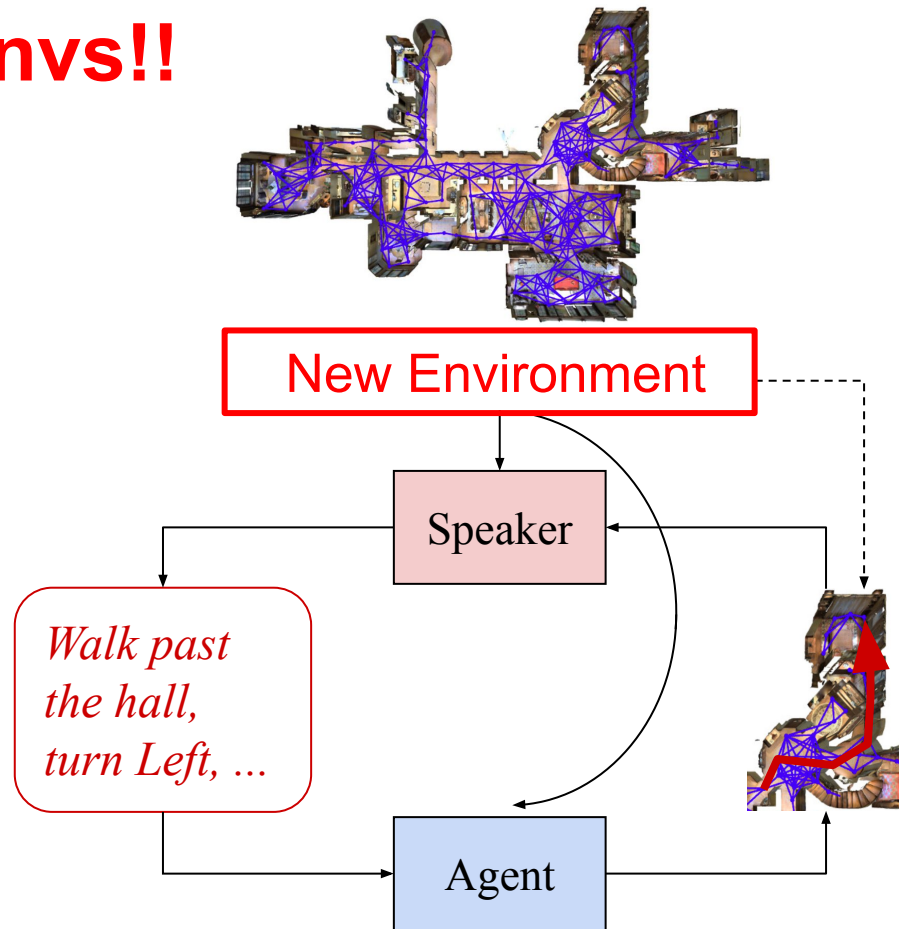
Back Translation: **Limited Envs?**

1. **New routes** from **existing environments.**
2. **New instructions** by pre-trained speaker.
3. Train agent on **new routes, new instructions, existing environments.**



Back Translation: **New Envs!!**

1. **New routes** from **New environments.**
2. **New instructions** by pre-trained speaker.
3. Train agent on **New routes,** **New instructions,** **New environments.**



How to get new environments?

Captured from new houses?

How to get new environments?

Captured from new houses?

Is very expensive...

Matterport Pro2

Made for Professionals

- Pro-grade resolution and accuracy
- Powerful battery for scanning multiple properties in a day
- Derivative assets available (floor plans, print-ready photos, and MatterPak™)
- Compatible with the Matterport Professional and Business subscription plans



\$3,395 USD



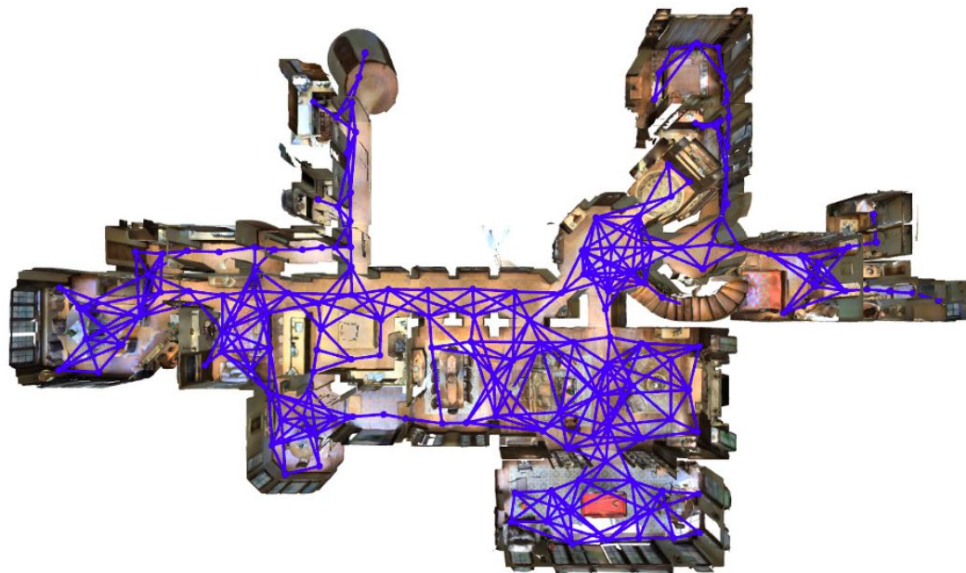
How to get new environments?

Generate new environments?

How to get new environments?

Generate new environments?

Not so easy...



How to get new environments?

Let's modify the existing environments!!

Illustration: Random Removal

Viewpoints
 t
 $t+1$



$O_{t,1}$



$O_{t,2}$



$O_{t+1,1}$



$O_{t+1,2}$

Views

RGB-image
views

Illustration: Random Removal

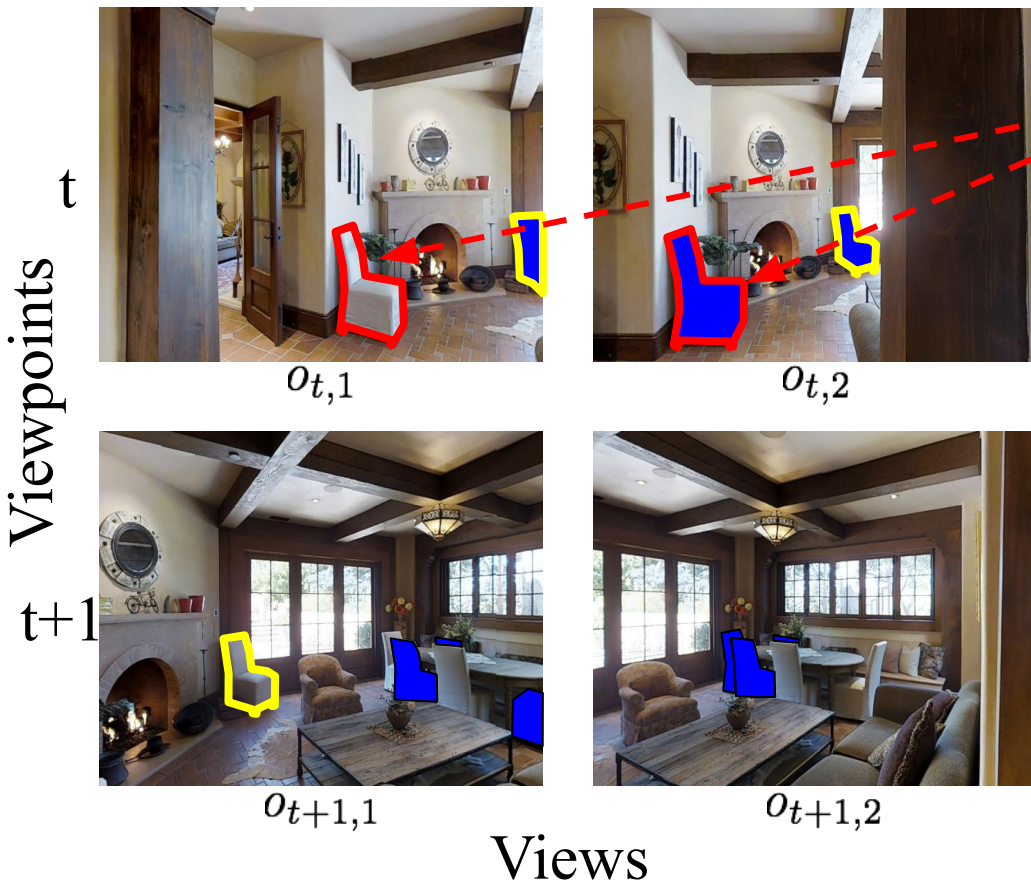
Viewpoints
 t
 $t+1$



Remove objects
(Marked in blue)

Views

Illustration: Random Removal (Two Issues)



**Incomplete
Removal:**
The chair is still
visible from
other views.

Illustration: Random Removal (Two Issues)

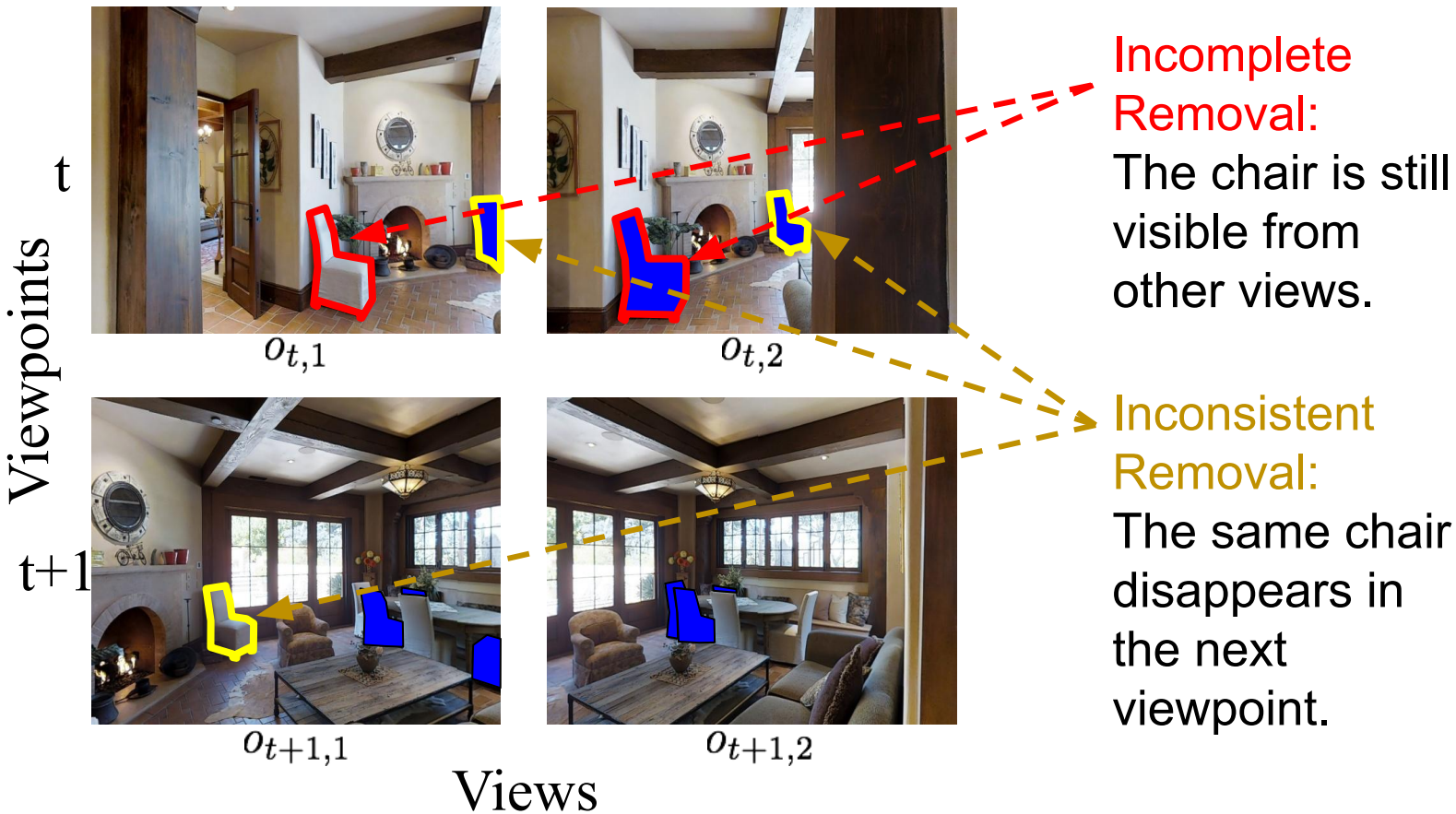
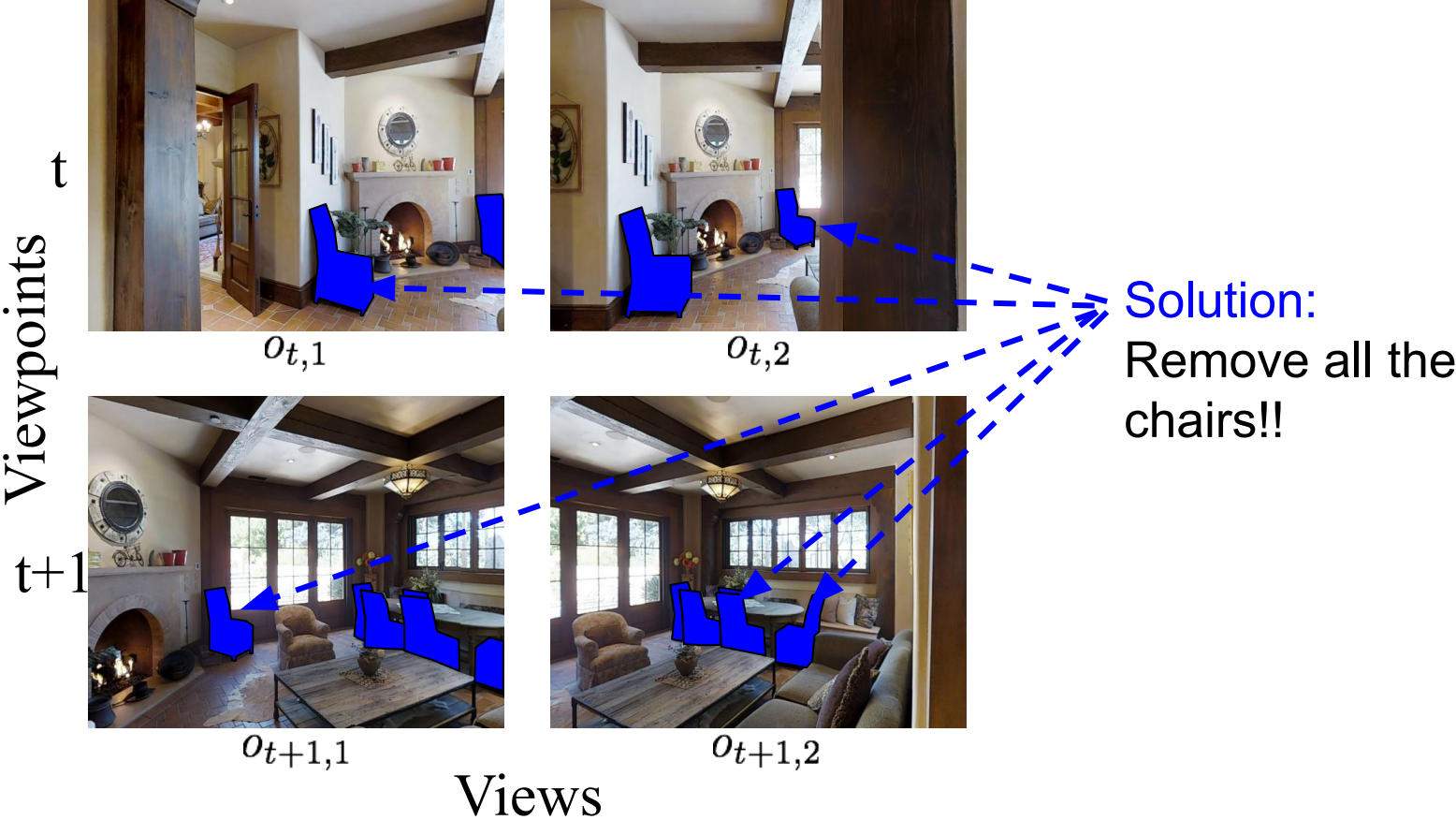


Illustration: Environmental Removal / Dropout

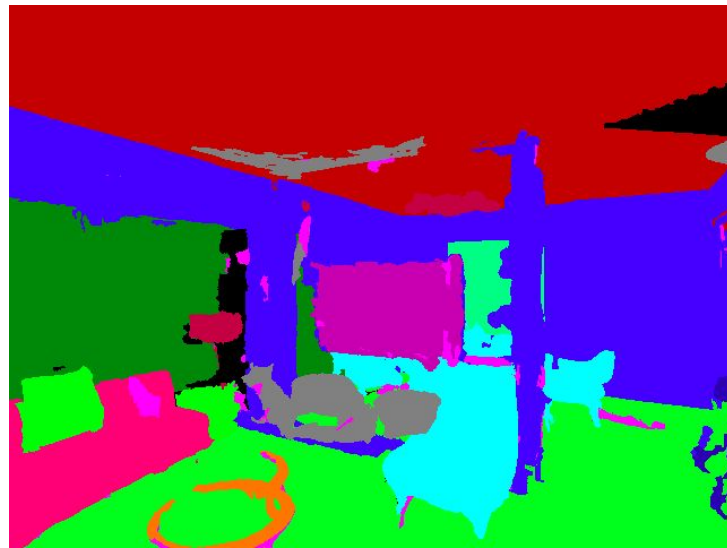


Environmental Dropout: Image-Level Implementation

Object-level annotation is noisy.



RGB Image



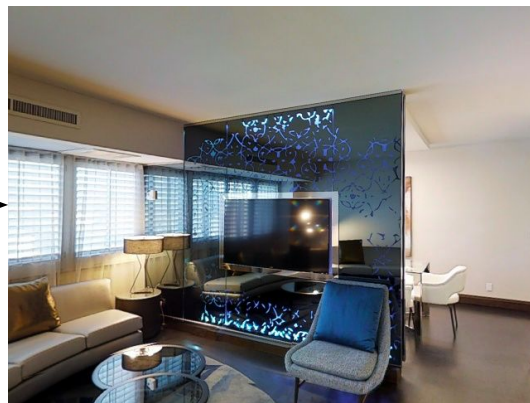
Semantic View

Environmental Dropout: Image-Level

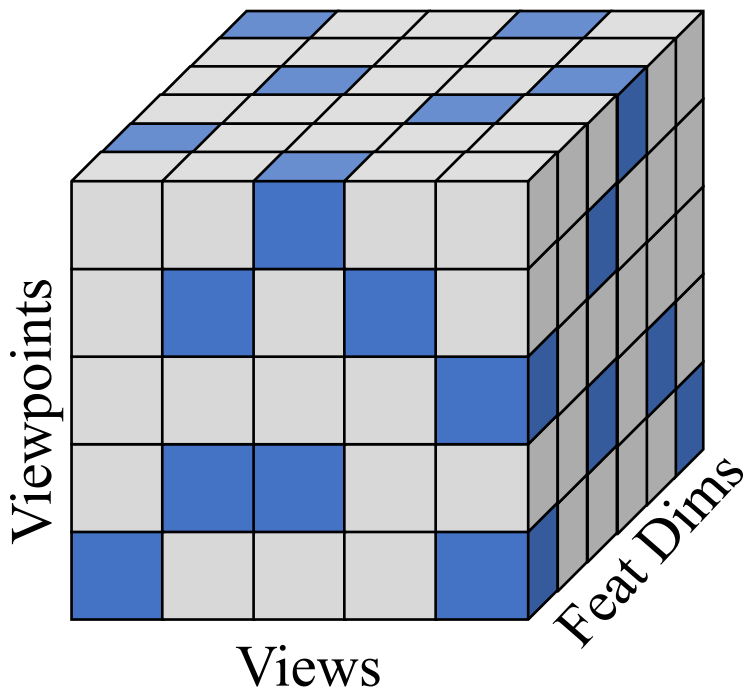
Rendering is slow for training agents.



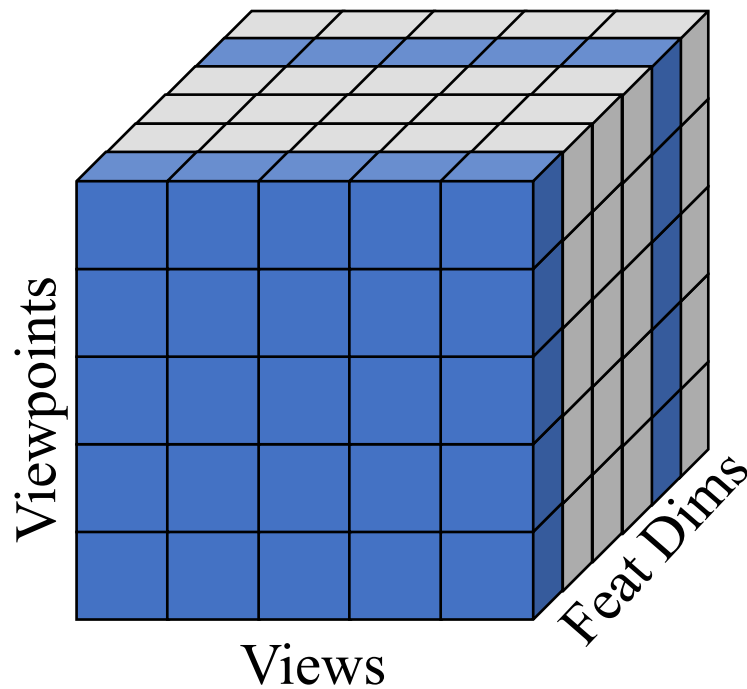
Render



Environmental Dropout: Feature-Level

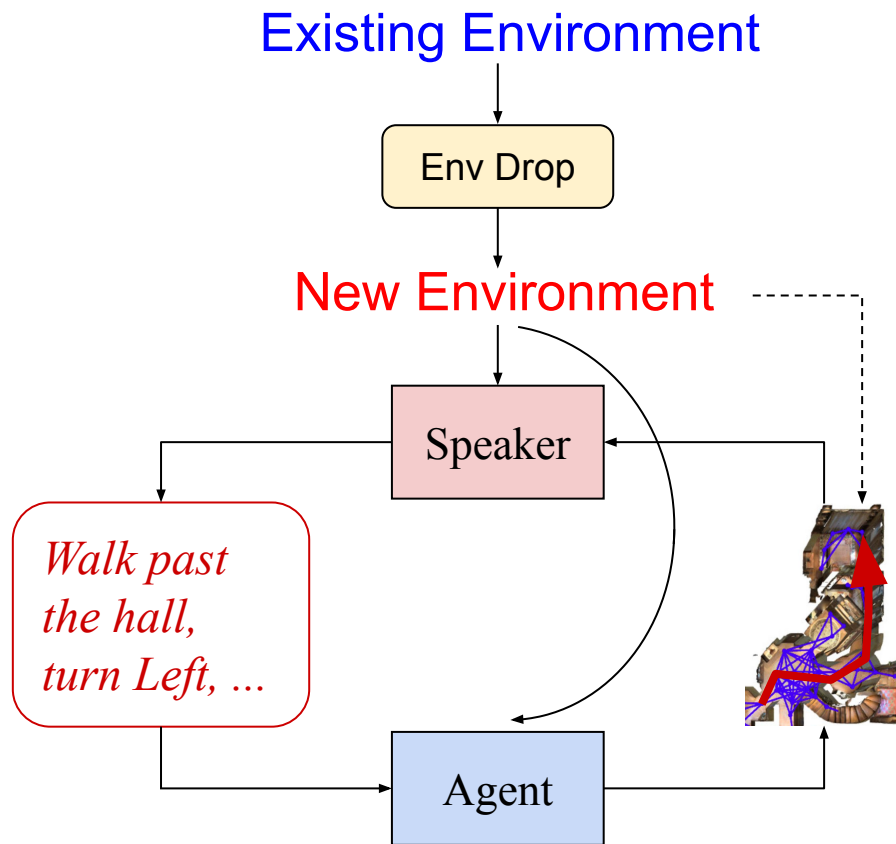


Random Feature Dropout



Environmental Dropout

Environmental Dropout: Full Pipeline

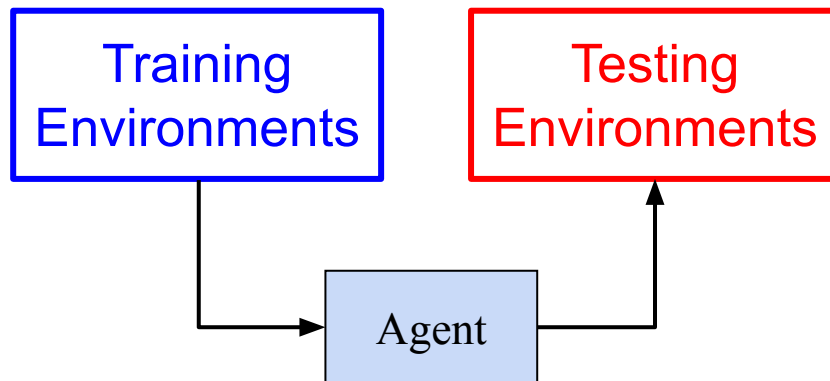


Results Comparison

Metric: Success Rate

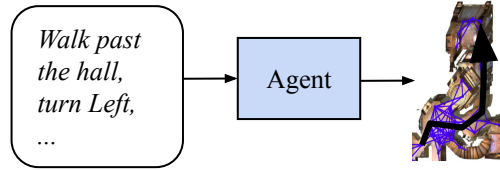


Evaluated in Unseen
Environments



Results Comparison

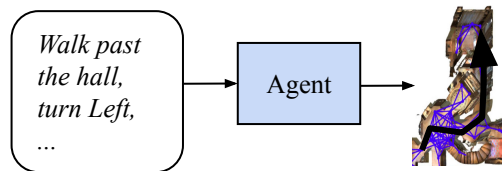
Agent Training



46.5%

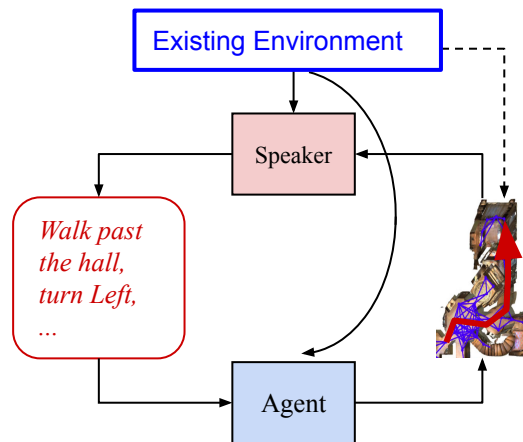
Results Comparison

Agent Training



46.5%

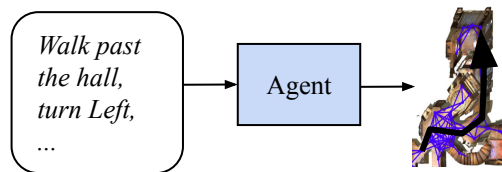
Back Translation



48.2% (+1.7%)

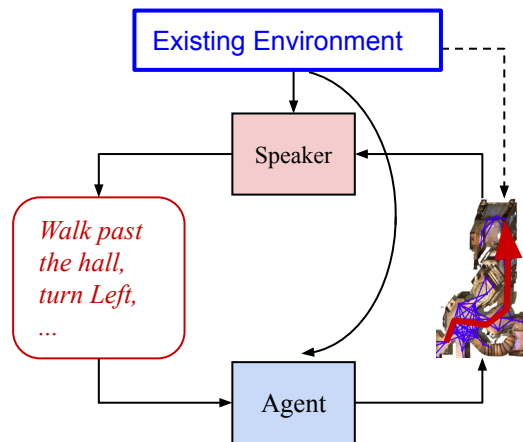
Results Comparison

Agent Training



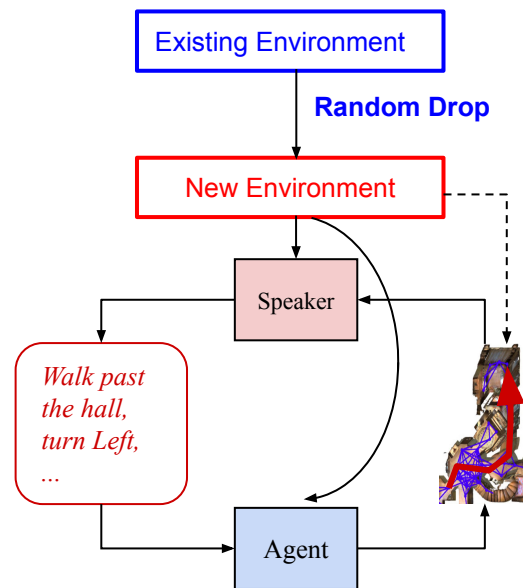
46.5%

Back Translation



48.2% (+1.7%)

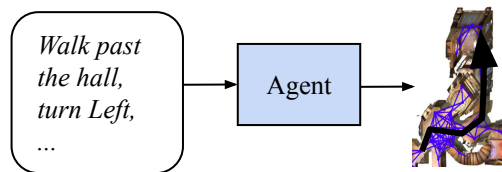
Back Translation w/ Random Dropout



48.4% (+1.9%)

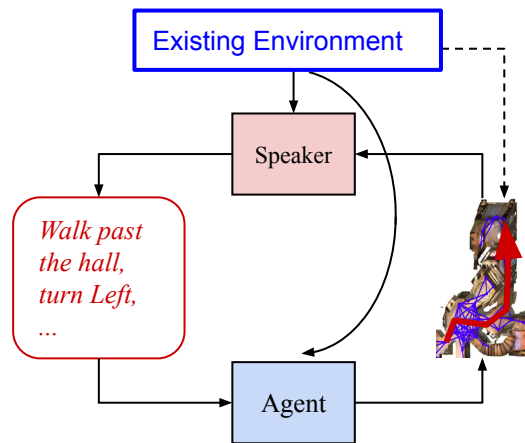
Results Comparison

Agent Training



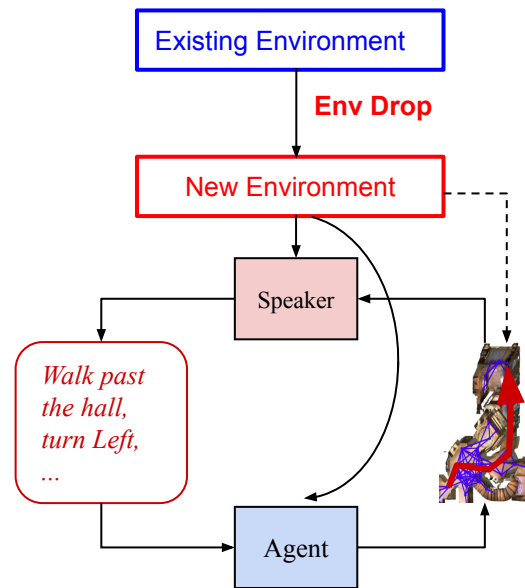
46.5%

Back Translation



48.2% (+1.7%)

Back Translation w/ Env Dropout



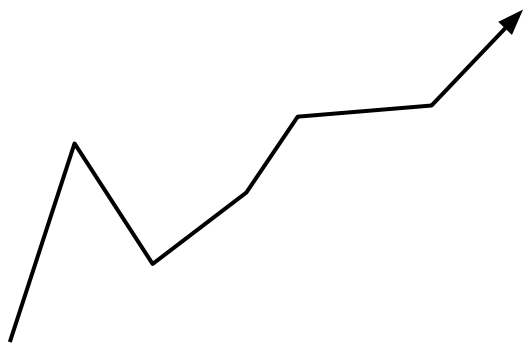
52.2% (+5.7%)

Leaderboard Results

Self-Monitoring Navigation Agent via Auxiliary Progress Estimation,
Ma et.al., 2019

Reinforced Cross-Modal Matching and Self-Supervised Imitation
Learning for Vision-Language Navigation, Wang et.al., 2019

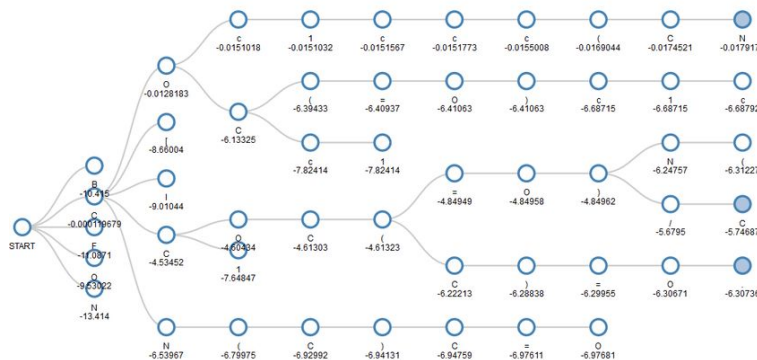
Greedy Decoding



Previous Best: 48.0%

Ours: **51.5% (+3.5%)**

Beam Search



Previous Best: 63.0%

Ours: **68.9% (+6.9%)**

Sufficient

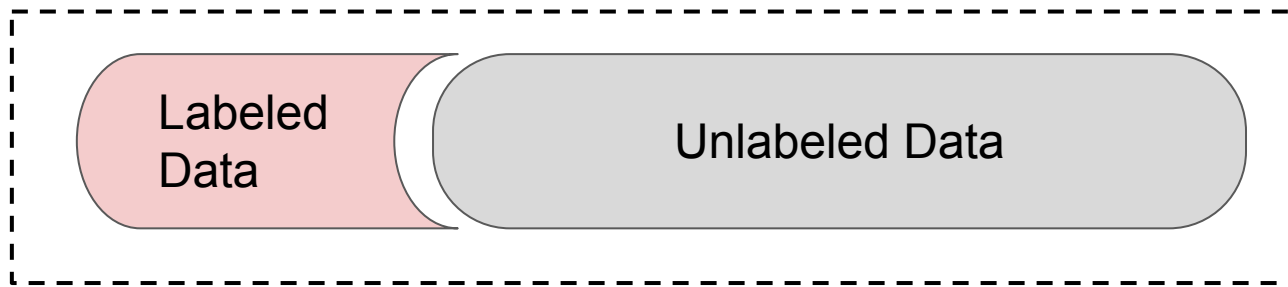
If we use “new” environment,
the result is better.

Sufficient **and Necessary**

If we use “new” environment,
the result is better.

If we **do not** use “new” environment,
the result **would not be** better.

Upper Bound of Back Translation (on Existing Envs)

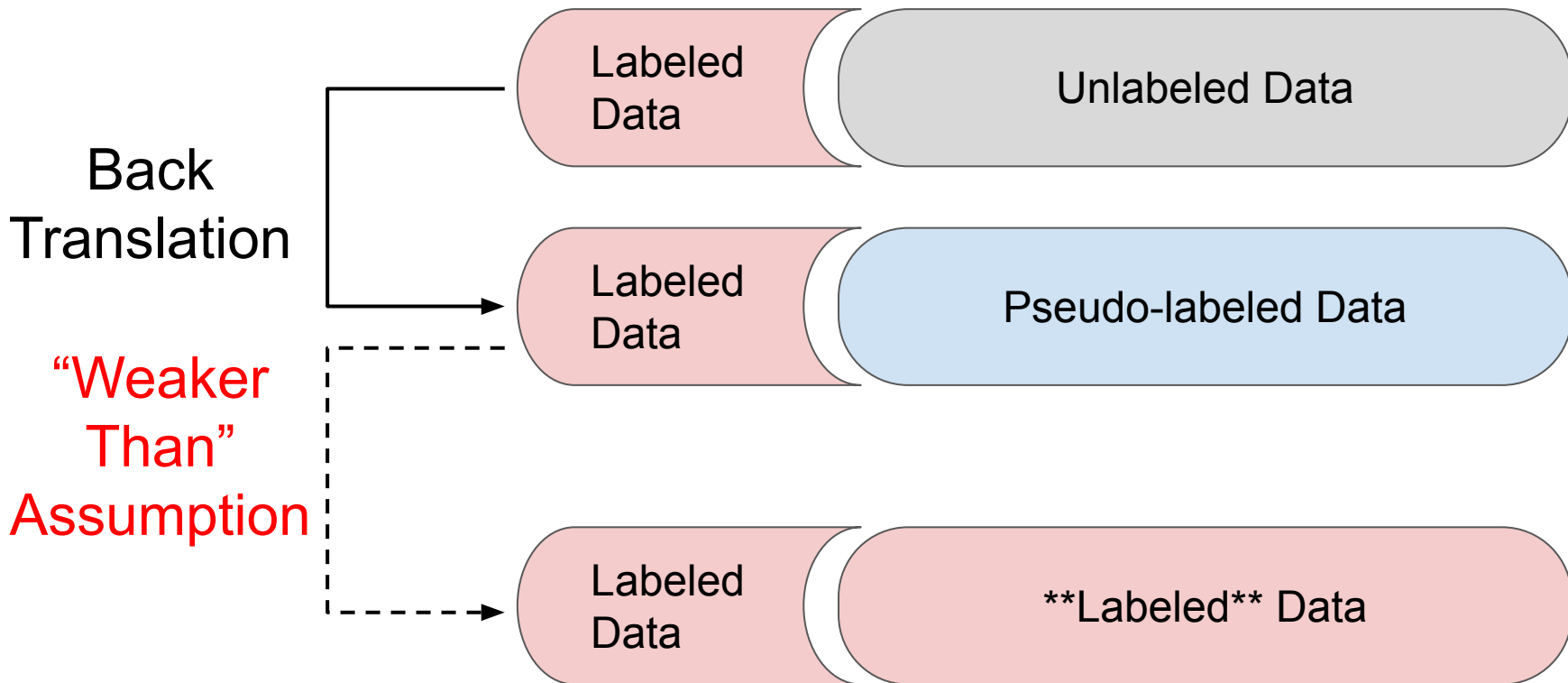


Existing Environments

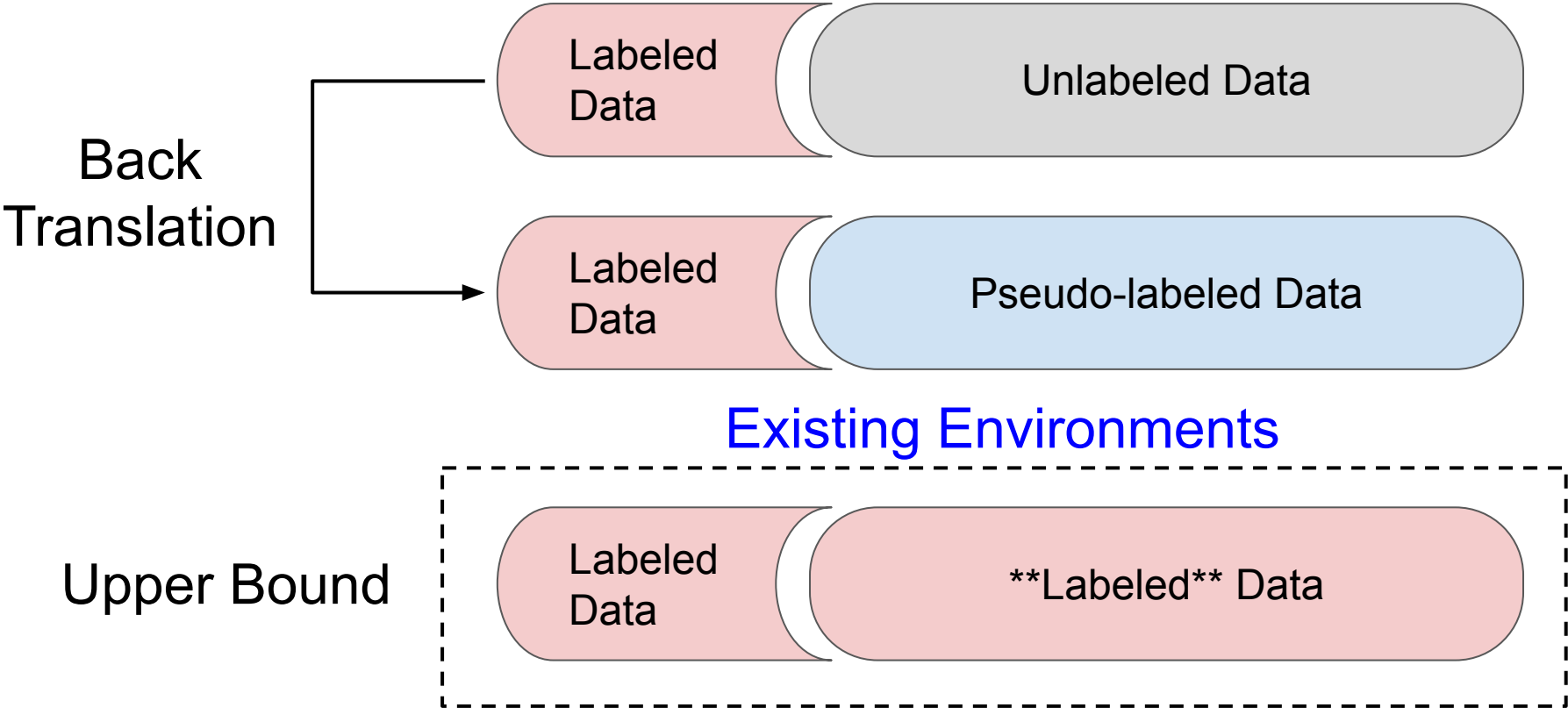
Upper Bound of Back Translation (on Existing Envs)



Upper Bound of Back Translation (on Existing Envs)



Upper Bound of Back Translation (on Existing Envs)

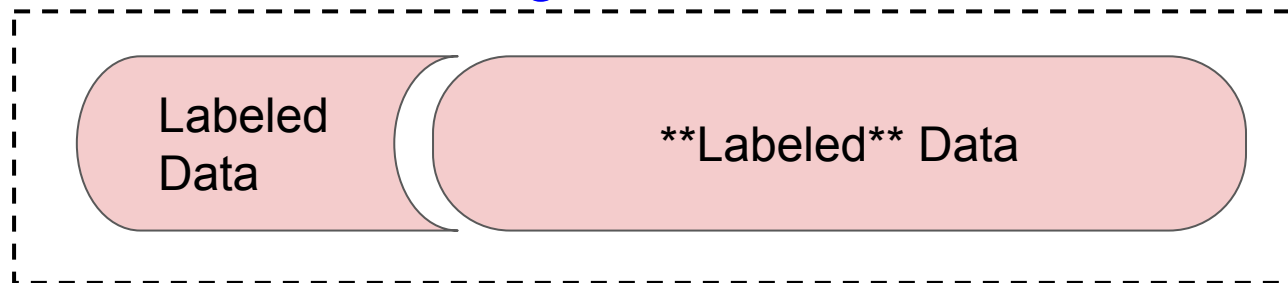


Upper Bound of Back Translation (on Existing Envs)

How to calculate (approximate)
this upper bound?

Existing Environments

Upper Bound



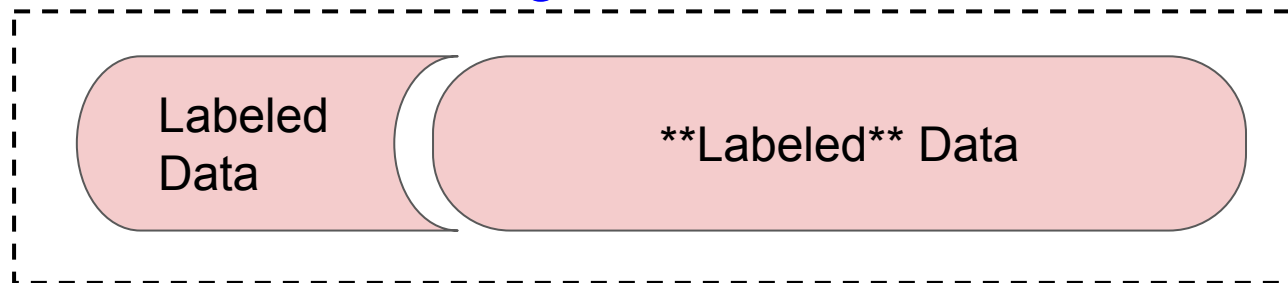
Upper Bound of Back Translation (on Existing Envs)

How to calculate (approximate)
this upper bound?

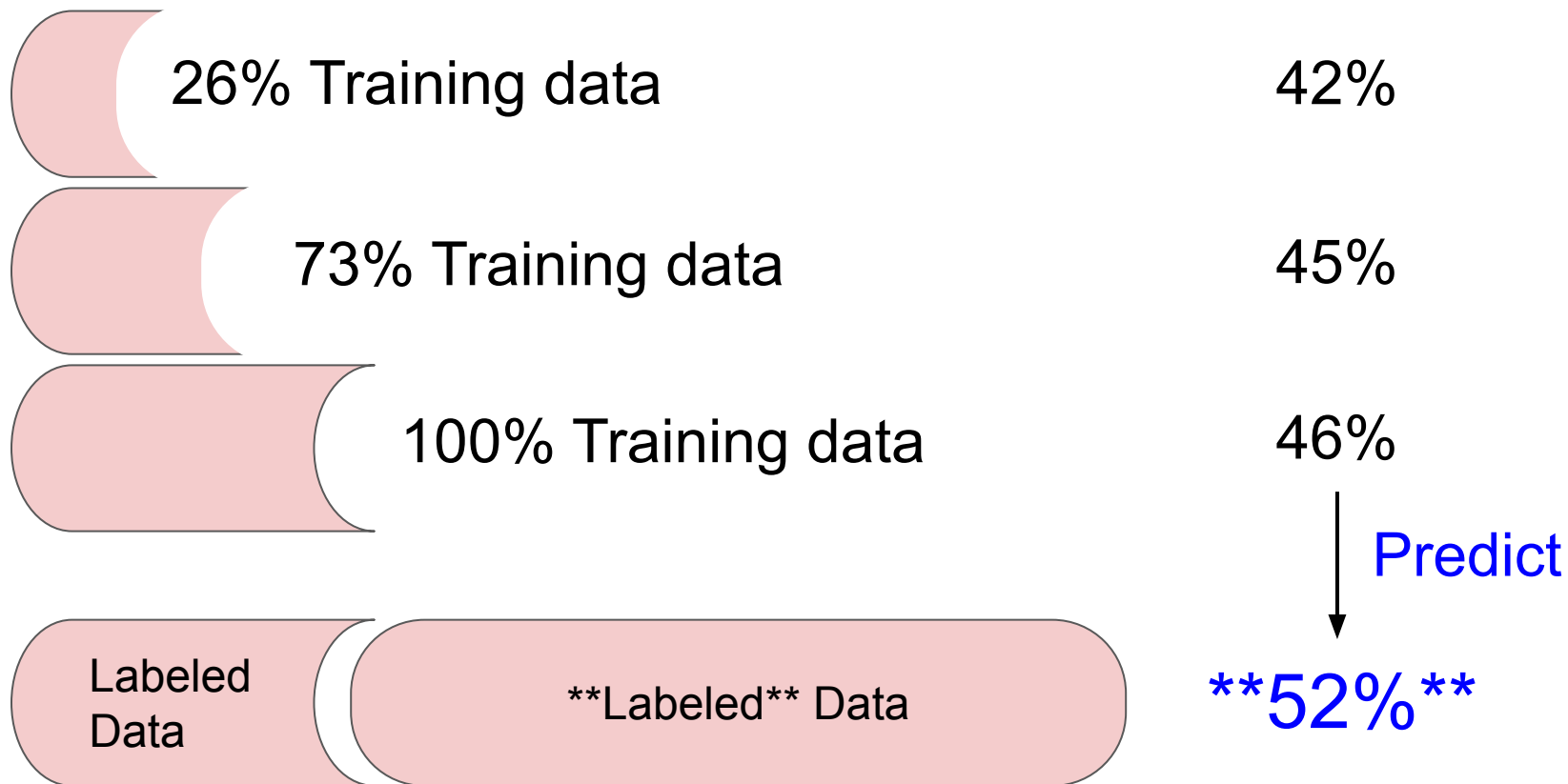
“Result Extrapolation” Approximation

Existing Environments

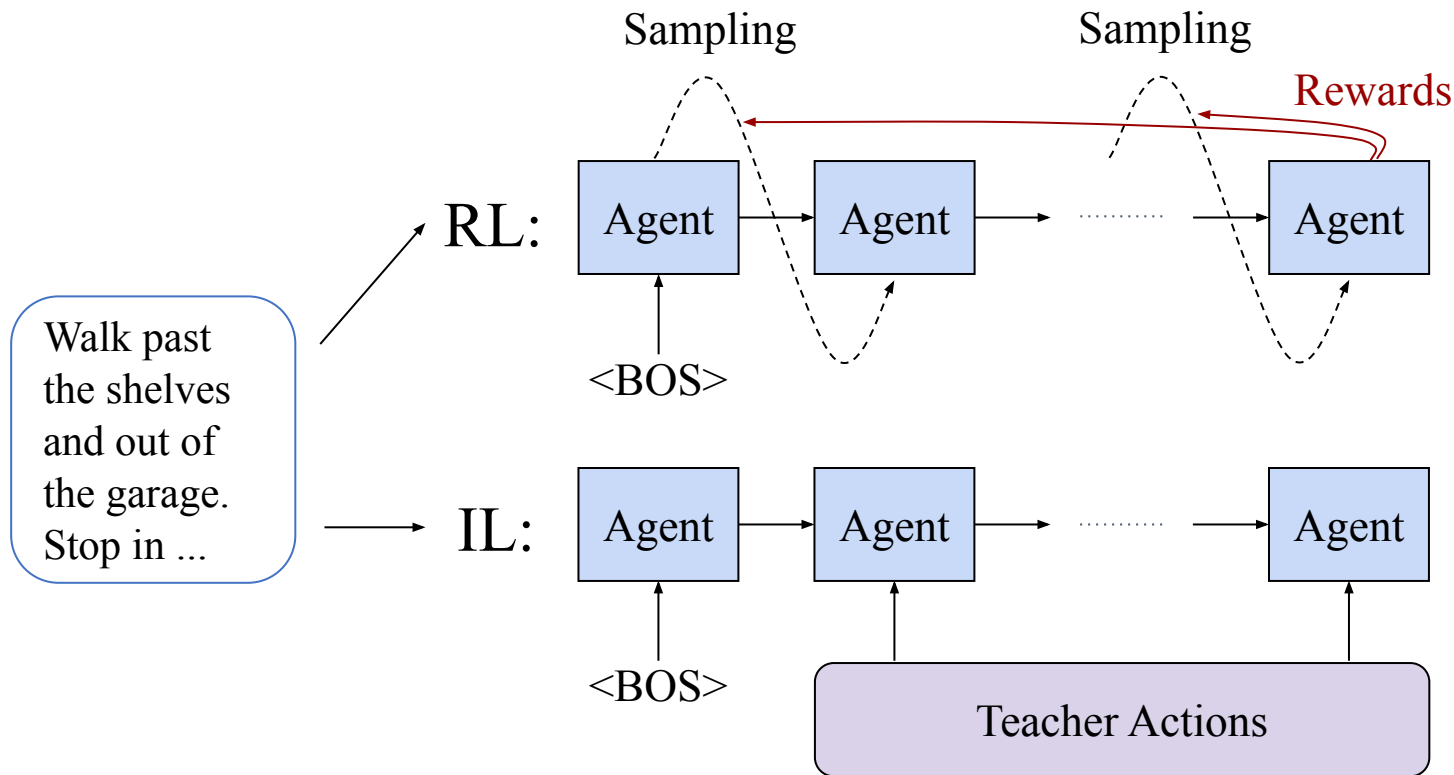
Upper Bound



“Result Extrapolation” Approximation



Reinforcement Learning + Imitation Learning



Thank you!

Hao Tan, Licheng Yu, Mohit Bansal

Code released at:

<https://github.com/airsplay/R2R-EnvDrop>

UNC Chapel Hill



Supported by ARO-YIP, ONR, Google, Facebook, Adobe, Baidu, and Salesforce.