# SOAT: A Scene- and Object-Aware Transformer for Vision-and-Language Navigation

NeurIPS 2021

Abhinav Moudgil[1]

Arjun Majumdar[1]

Harsh Agrawal[1]

Stefan Lee[2]

Dhruv Batra[1]

[1] Georgia Tech

[2] Oregon State University

# Vision-and-Language Navigation (VLN) Task

*Given a natural language instruction, the agent needs to navigate to a goal location by following the given instruction.*

**Input:** *instruction and panoramic observation*
**Output:** *sequence of actions*



Goal: 8.2m

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

# Motivation

## Instruction

*Exit the bedroom and turn left. Continue down the hall and into the room straight ahead and stop before the desk with two green chairs.*

# Motivation

## Instruction

*Exit the <span style="color:red">bedroom</span> and turn left. Continue down the <span style="color:orange">hall</span> and into the room straight ahead and stop before the desk with two green chairs.*

## Scene Descriptions



*bedroom*

*hall*

# Motivation

**Instruction**

**Scene Descriptions**

**Object References**

*Exit the* <span style="color:red">bedroom</span> *and turn left. Continue down the* <span style="color:orange">hall</span> *and into the room straight ahead and stop before the* <span style="color:blue">desk</span> *with two* <span style="color:green">green chairs</span>.
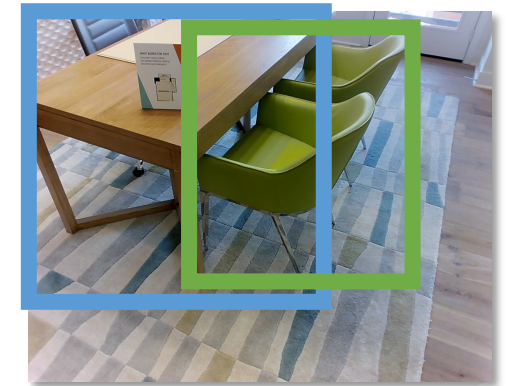
# Motivation

## Instruction

*Exit the <span style="color:red">bedroom</span> and turn left. Continue down the <span style="color:orange">hall</span> and into the room straight ahead and stop before the <span style="color:blue">desk</span> with two <span style="color:green">green chairs</span>.*
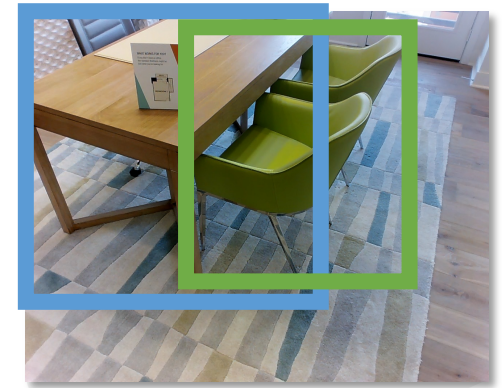
*Most VLN Methods*

## Scene Descriptions

*bedroom*

*hall*



## Object References

*desk*    *green chairs*

# Motivation

**Instruction**

**Scene Descriptions**

**Object References**

*Exit the bedroom and turn left. Continue down the hall and into the room straight ahead and stop before the desk with two green chairs.*

*bedroom*

*hall*

*desk    green chairs*

# Our Approach



*SOAT: Scene- and Object-Aware Transformer*

**Instruction**

*Exit the bedroom and turn left. Continue down the hall and into the room straight ahead and stop before the desk with two green chairs.*
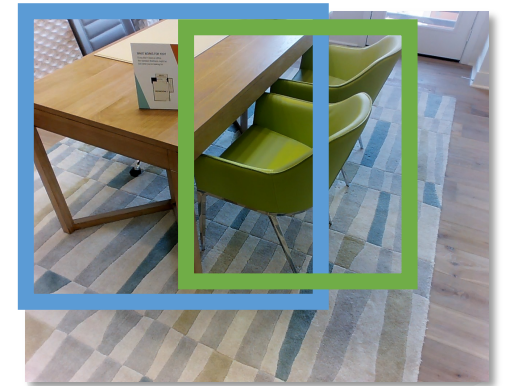
**Scene Descriptions**

*bedroom*

*hall*

**Object References**

*desk*   *green chairs*

# Our Approach

*Exit the <span style="color:red">bedroom</span> and turn left. Continue down the <span style="color:orange">hall</span> and into the room straight ahead and stop before the <span style="color:blue">desk</span> with two <span style="color:green">green chairs</span>.*

**Instruction**

# Our Approach

*Exit the* <span style="color:red">bedroom</span> *and turn left. Continue down the* <span style="color:orange">hall</span> *and into the room straight ahead and stop before the* <span style="color:blue">desk</span> *with two* <span style="color:green">green chairs</span>.

**Instruction**
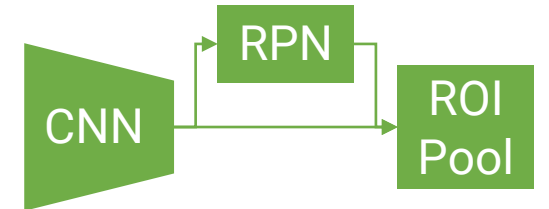


**Panoramic Observation**

# Our Approach



BERT Tokenizer

*Exit the bedroom and turn left. Continue down the hall and into the room straight ahead and stop before the desk with two green chairs.*

**Instruction**

**Panoramic Observation**

# Our Approach

*word tokens*  ☐ ☐ ☐ ☐

BERT Tokenizer



*Exit the bedroom and turn left. Continue down the hall and into the room straight ahead and stop before the desk with two green chairs.*

**Instruction**

**Panoramic Observation**

# Our Approach

*word tokens* □□□□

BERT Tokenizer

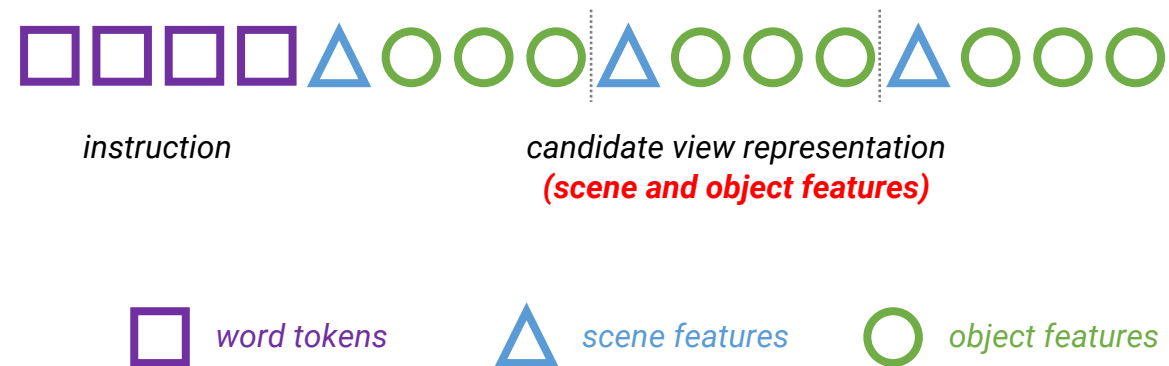*Exit the bedroom and turn left. Continue down the hall and into the room straight ahead and stop before the desk with two green chairs.*

**Instruction**



**candidate views**

# Our Approach

*word tokens* ☐☐☐☐

BERT Tokenizer



*Exit the* <span style="color:red">*bedroom*</span> *and turn left. Continue down the* <span style="color:orange">*hall*</span> *and into the room straight ahead and stop before the* <span style="color:blue">*desk*</span> *with two* <span style="color:green">*green chairs*</span>.
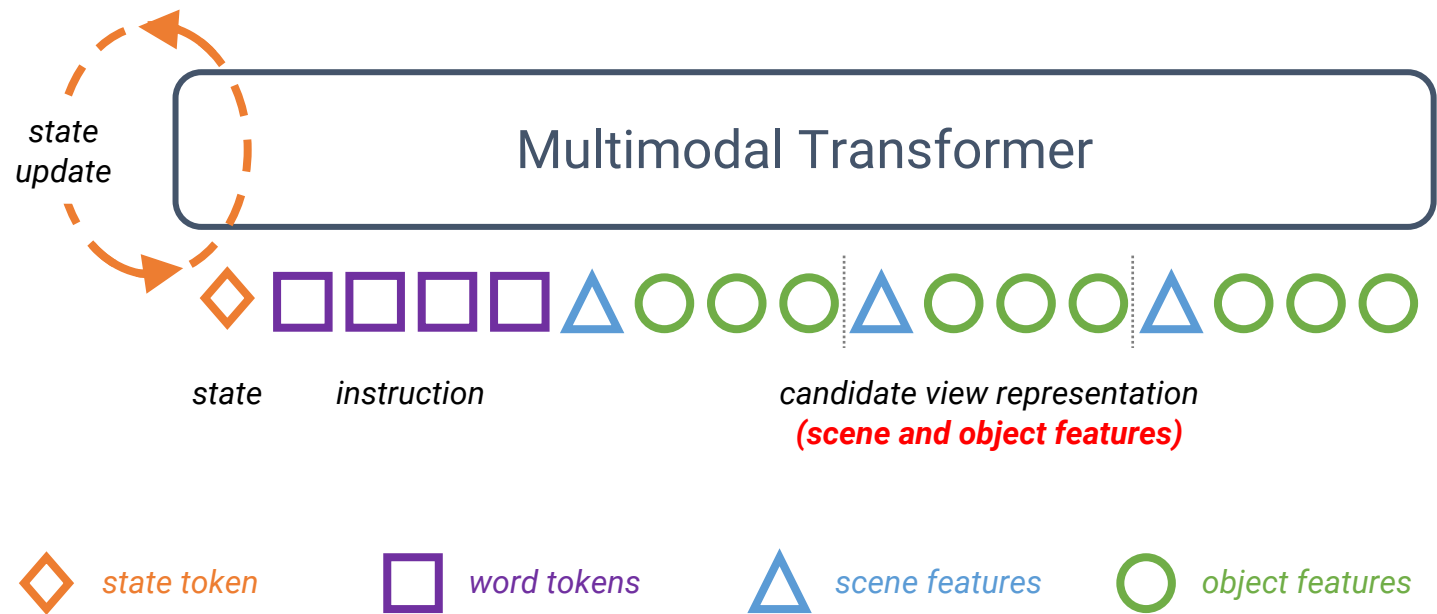
**Instruction**

Scene Classifier (Places)

Object Detector (Visual Genome)

CNN

CNN → RPN → ROI Pool



**candidate views**

# Our Approach

word tokens

candidate view representation
(**scene** and **object** features)

BERT Tokenizer



Scene Classifier
(Places)

CNN

Object Detector
(Visual Genome)

RPN

ROI
Pool

CNN

*Exit the bedroom and turn left. Continue down the hall and into the room straight ahead and stop before the desk with two green chairs.*

**Instruction**



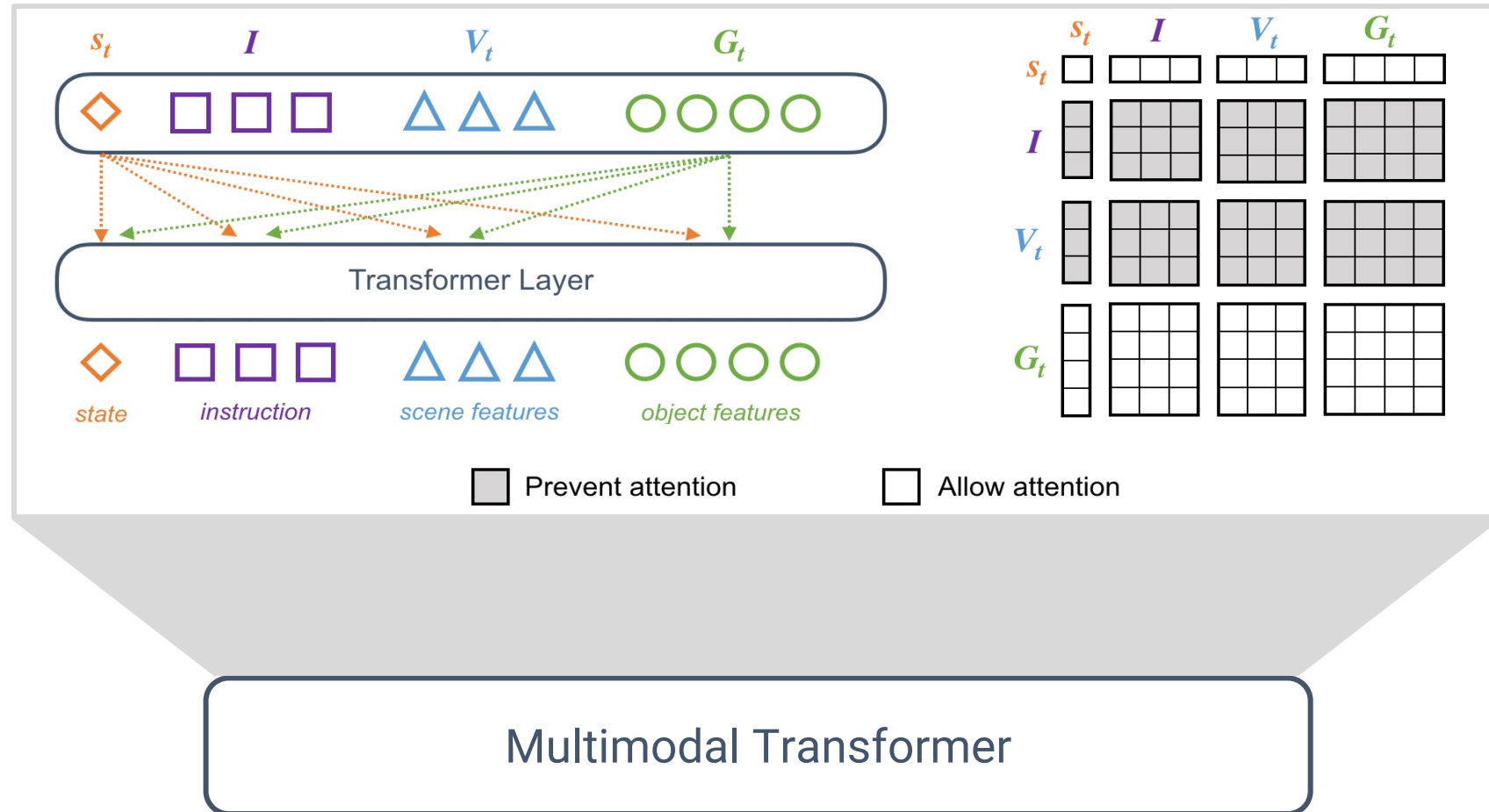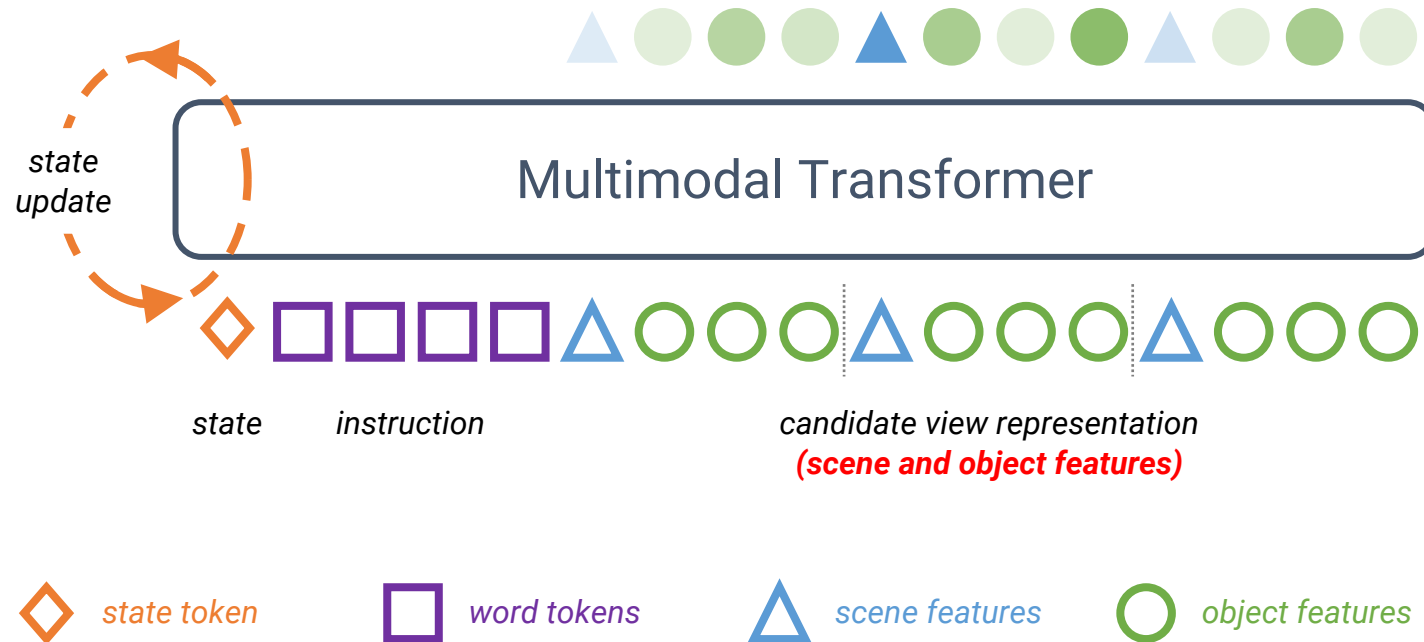candidate
views

# SOAT: Scene- and Object-Aware Transformer

*instruction*

*candidate view representation*
*(scene and object features)*

□ word tokens    △ scene features    ○ object features

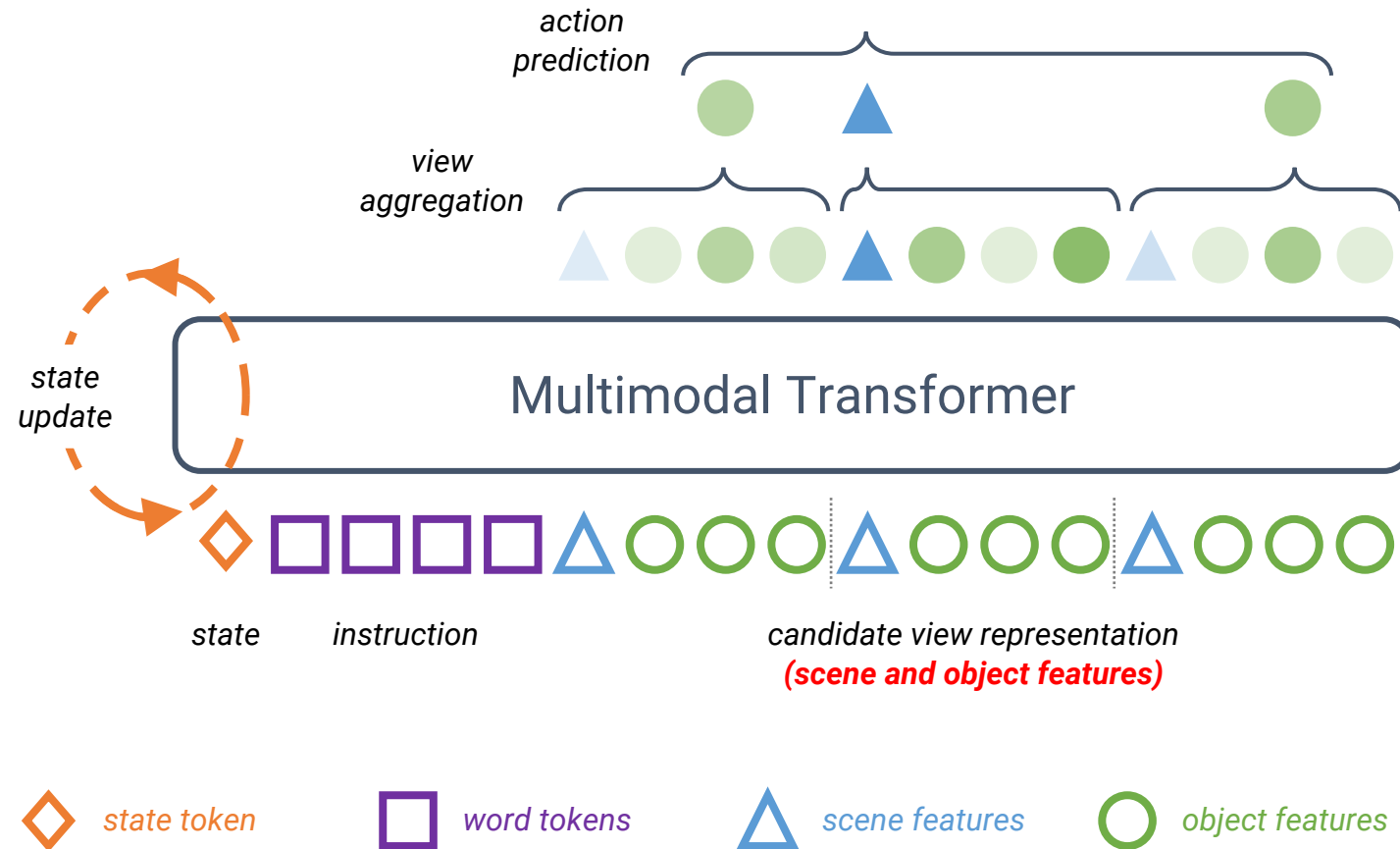# SOAT: Scene- and Object-Aware Transformer

# Selective Attention

# SOAT: Scene- and Object-Aware Transformer



state update

Multimodal Transformer

state

instruction

candidate view representation
**(scene and object features)**

⬦ state token    ☐ word tokens    △ scene features    ◯ object features
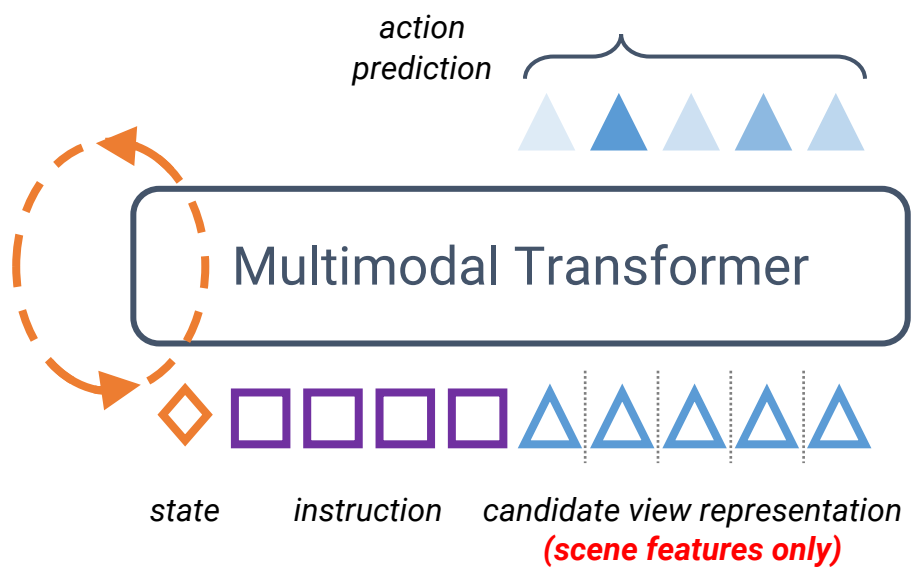
# SOAT: Scene- and Object-Aware Transformer

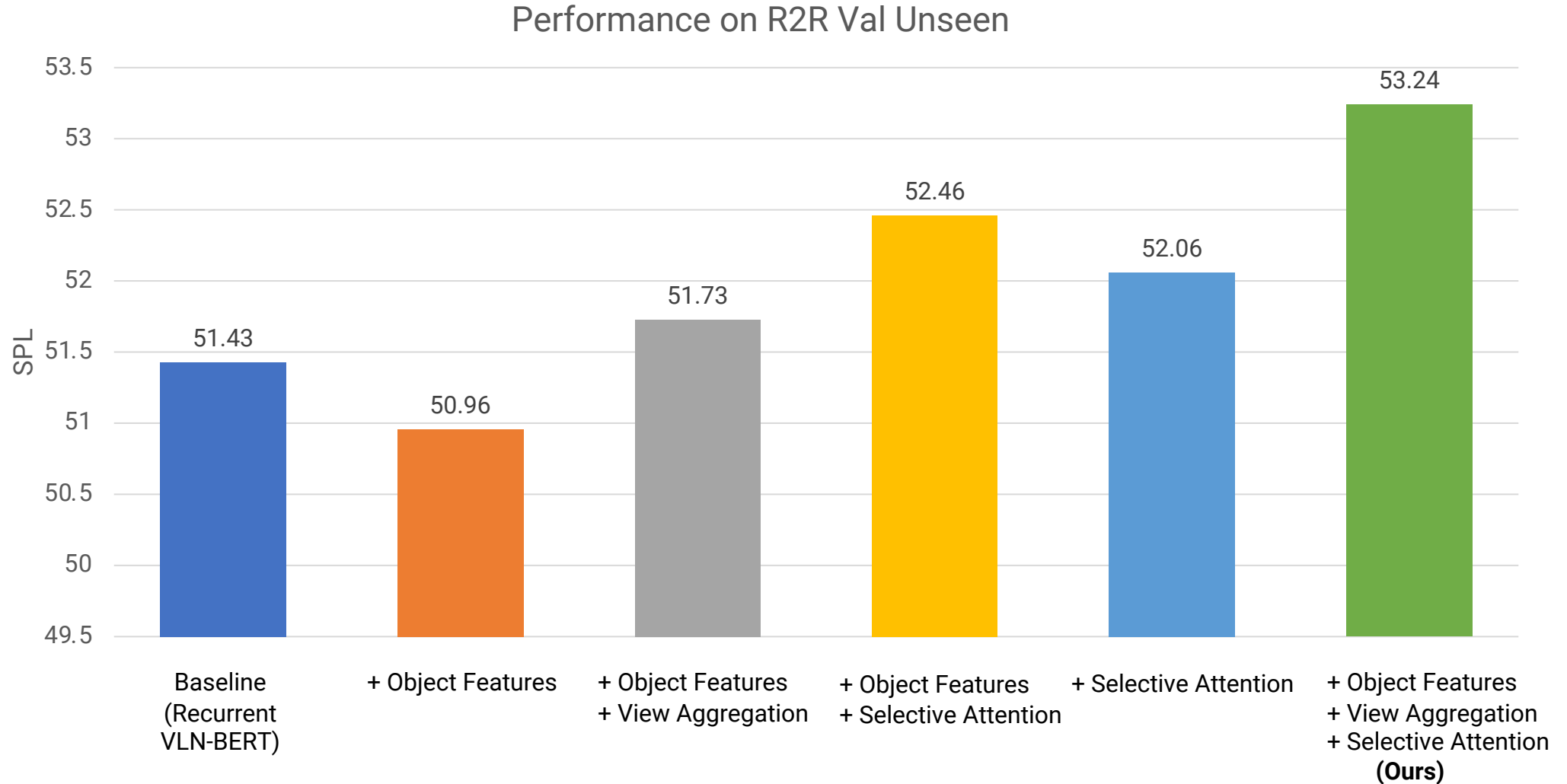# SOAT: Scene- and Object-Aware Transformer

Baseline: VLN↻BERT

# Ablation Study



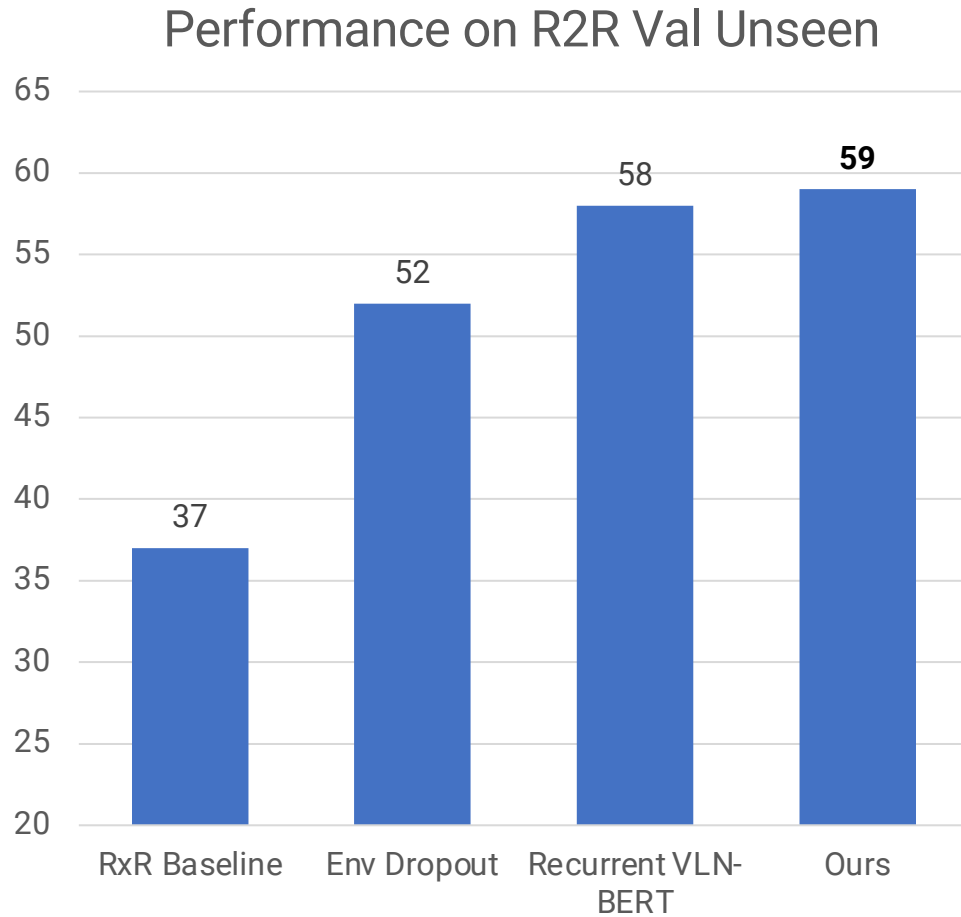Performance on R2R Val Unseen

# Comparison with State-of-the-Art



Performance on R2R Val Unseen

Performance on RxR Val Unseen

# Results on RxR Object-Heavy Instructions Subset

# Results on RxR Object-Heavy Instructions Subset

*You are standing in front of a curio cabinet with lots of dolls in it. You are going turn to your right and enter that doorway. You will see dark wood floors. You are now in a bedroom. It will have a gray and black striped comforter on it. You are going to walk into the bedroom and walk in between the foot of the bed and on your left will be a dresser with a large outdoor painting on it. You are going to stop right there in between those two.*
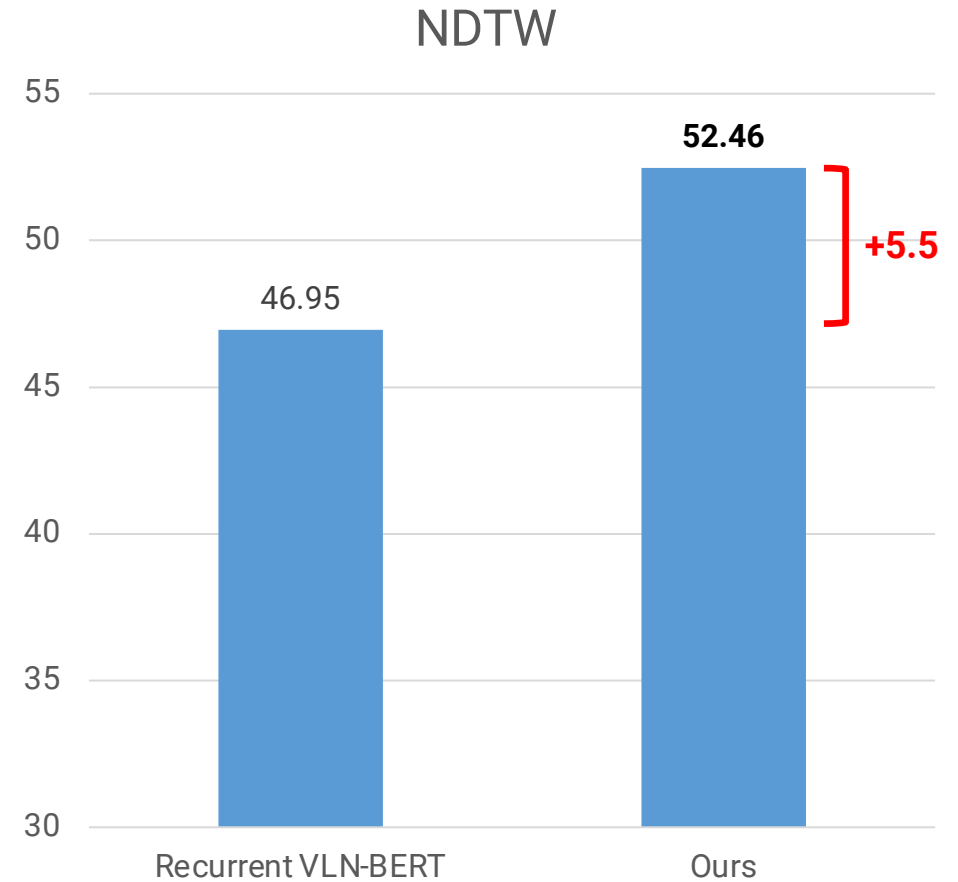
Object-Heavy Instruction Example from RxR Dataset

# Results on RxR Object-Heavy Instructions Subset

*You are standing in front of a curio cabinet[1] with lots of dolls[2] in it. You are going turn to your right and enter that doorway. You will see dark wood floors. You are now in a bedroom. It will have a gray and black striped comforter[3] on it. You are going to walk into the bedroom and walk in between the foot[4] of the bed[5] and on your left will be a dresser[6] with a large outdoor painting[7] on it. You are going to stop right there in between those two.*

Object-Heavy Instruction Example from RxR Dataset

# Results on RxR Object-Heavy Instructions Subset

# Qualitative Example: Baseline VLN↻BERT

*RxR Instruction: "You are going to start facing a front door. Turn to the left and go up the stairs. Once at the top, you are going to keep walking straight, pass the china cabinet on you right and take two more steps, and stop once you are next to the desk and the small dresser with two red glasses on top. Once you are there, you are done."*





— **Recurrent VLN-BERT**
— **GT path**

# Qualitative Example: SOAT (Ours)

*RxR Instruction: "You are going to start facing a front door. Turn to the left and go up the stairs. Once at the top, you are going to keep walking straight, pass the china cabinet on you right and take two more steps, and stop once you are next to the desk and the small dresser with two red glasses on top. Once you are there, you are done."*

# Thank you!