



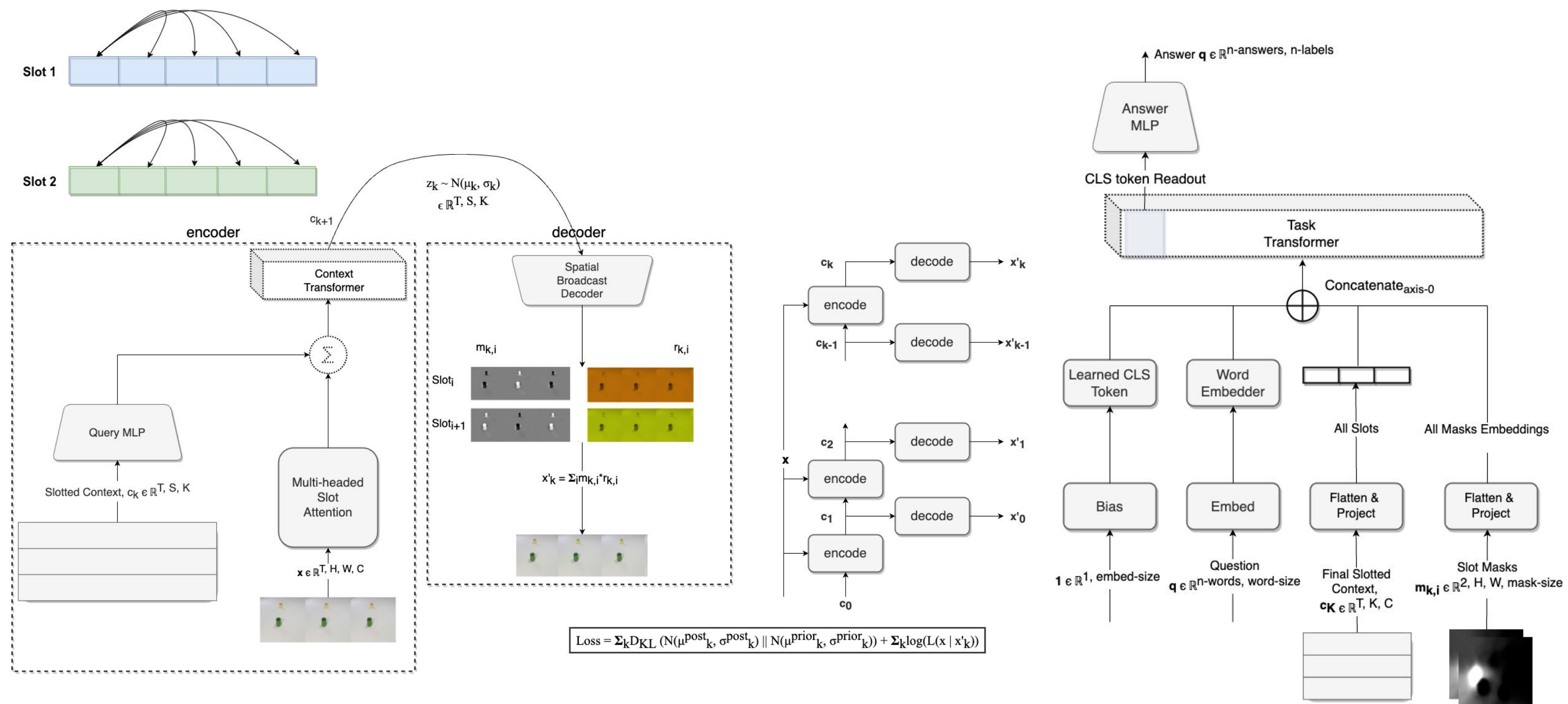
# Solving Reasoning Tasks with a Slot Transformer

Ryan Faulkner, Daniel Zoran

## Introduction & Model Architecture

Our goal is to train a model that can be induced to learn to understand visual scenes that play out over time. Also to generalise across these kinds of domains which may contain complex visual interactions and temporal dynamics.

1. **Slotted Representation** for multi-object scenes.
2. The **Context Transformer** allows relationships among objects to be learned over many timesteps.
3. An **Iterative Model** to enable reasoning across steps.
4. **Generative Model** for improved generalisation over visual reasoning domains.



## Experiments & Results

- CLEVRER: QA & Spatial Reasoning
- CATER: QA & Spatial Reasoning
- Kinetics 600: Video Classification

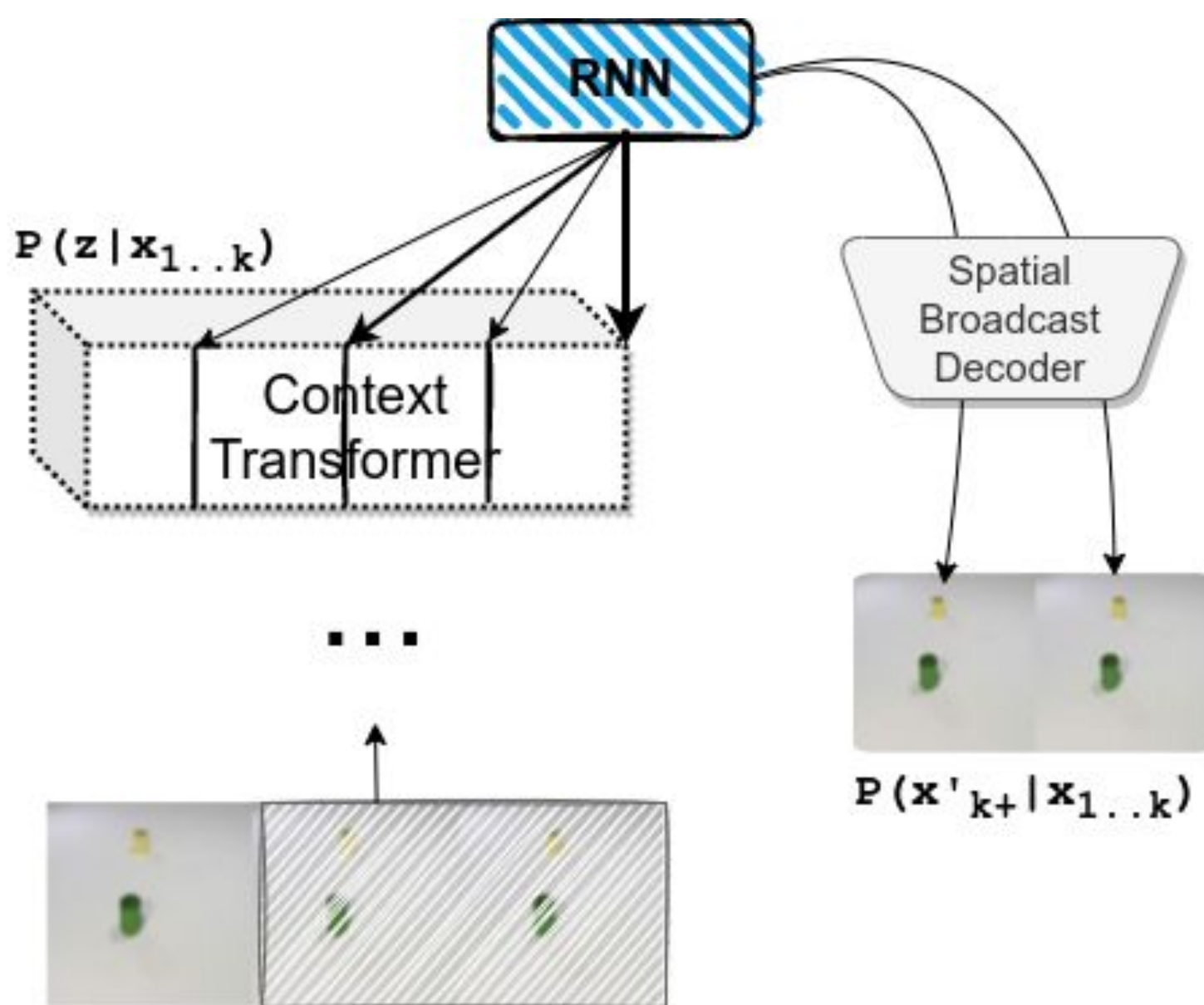
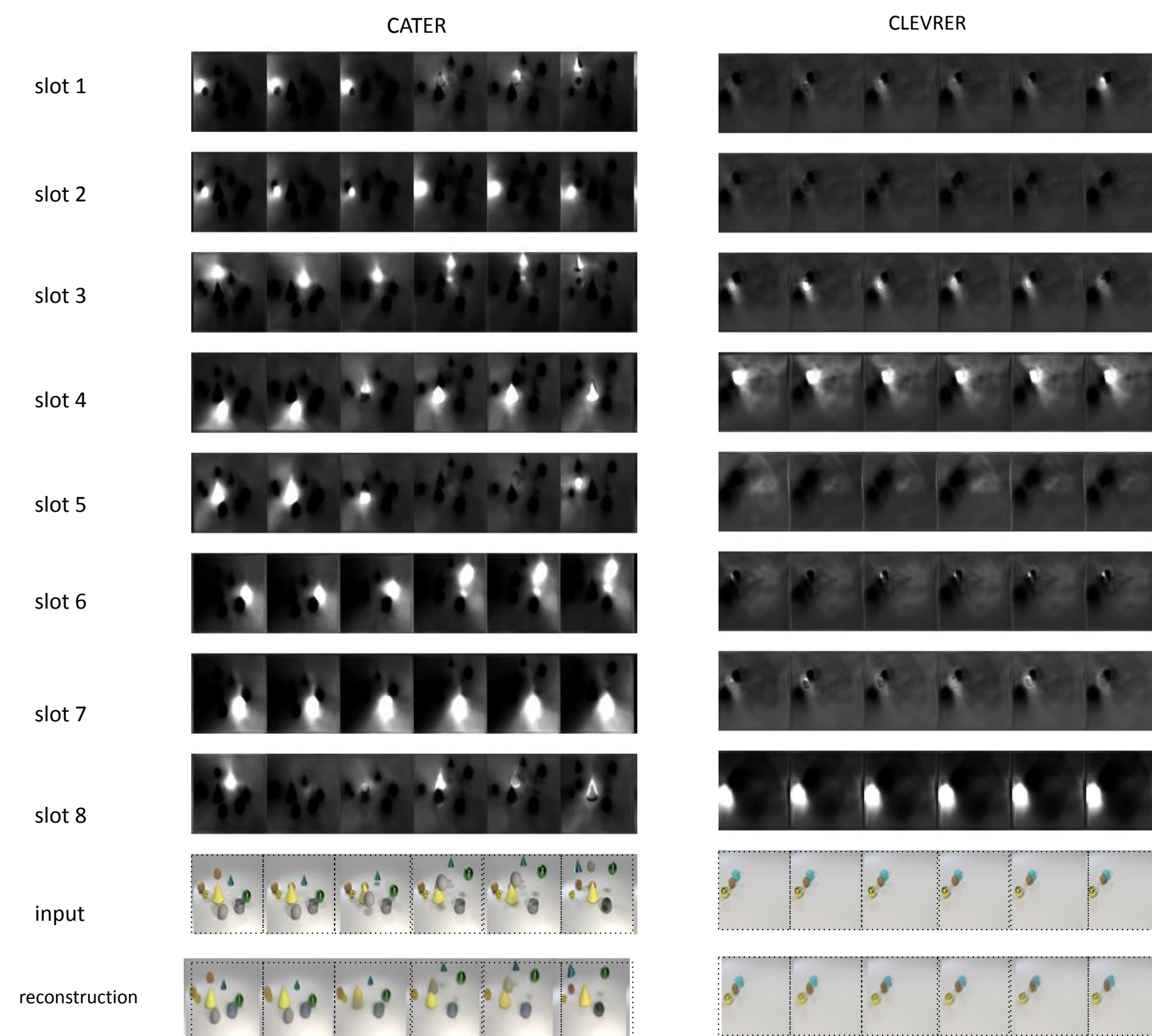
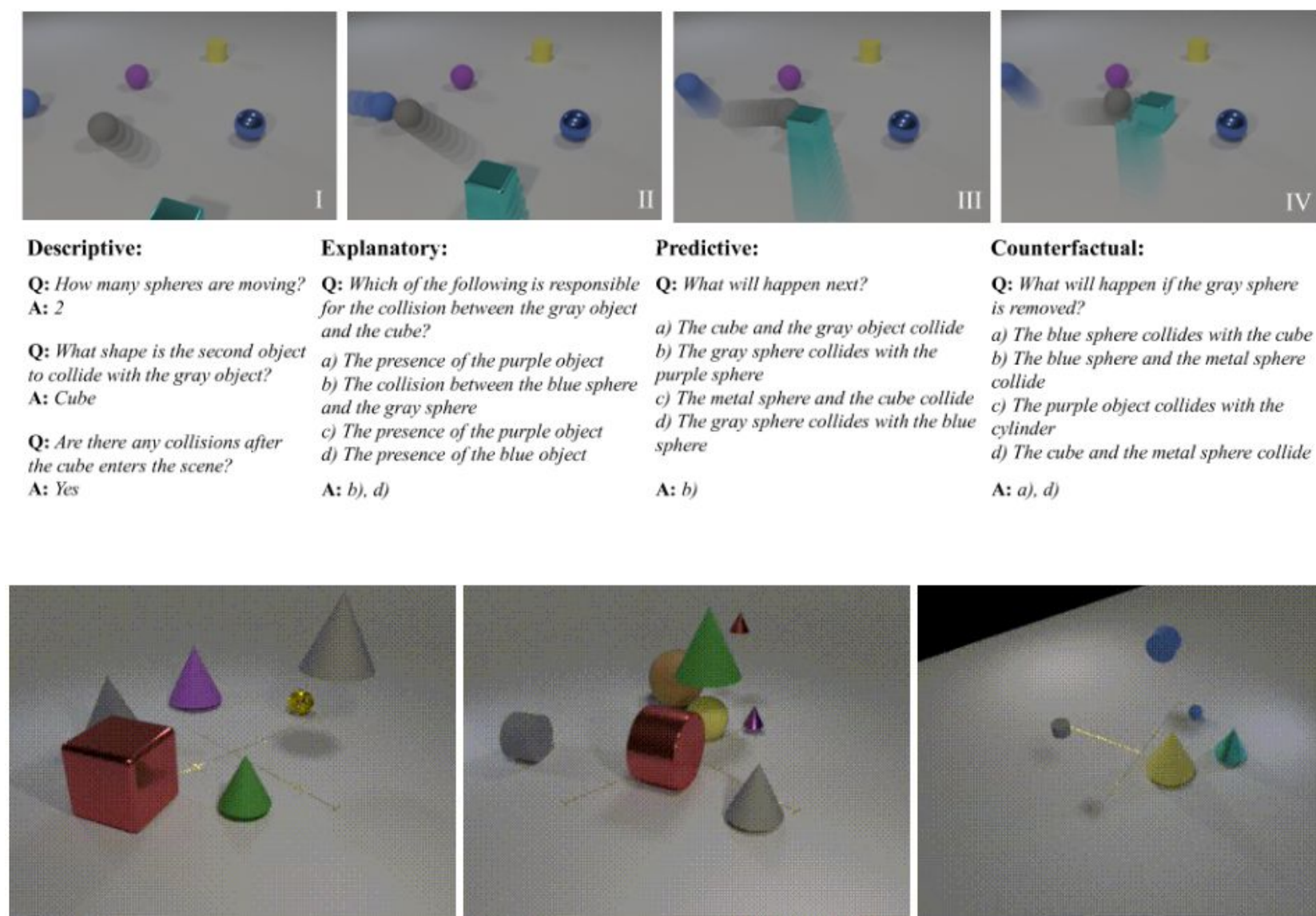
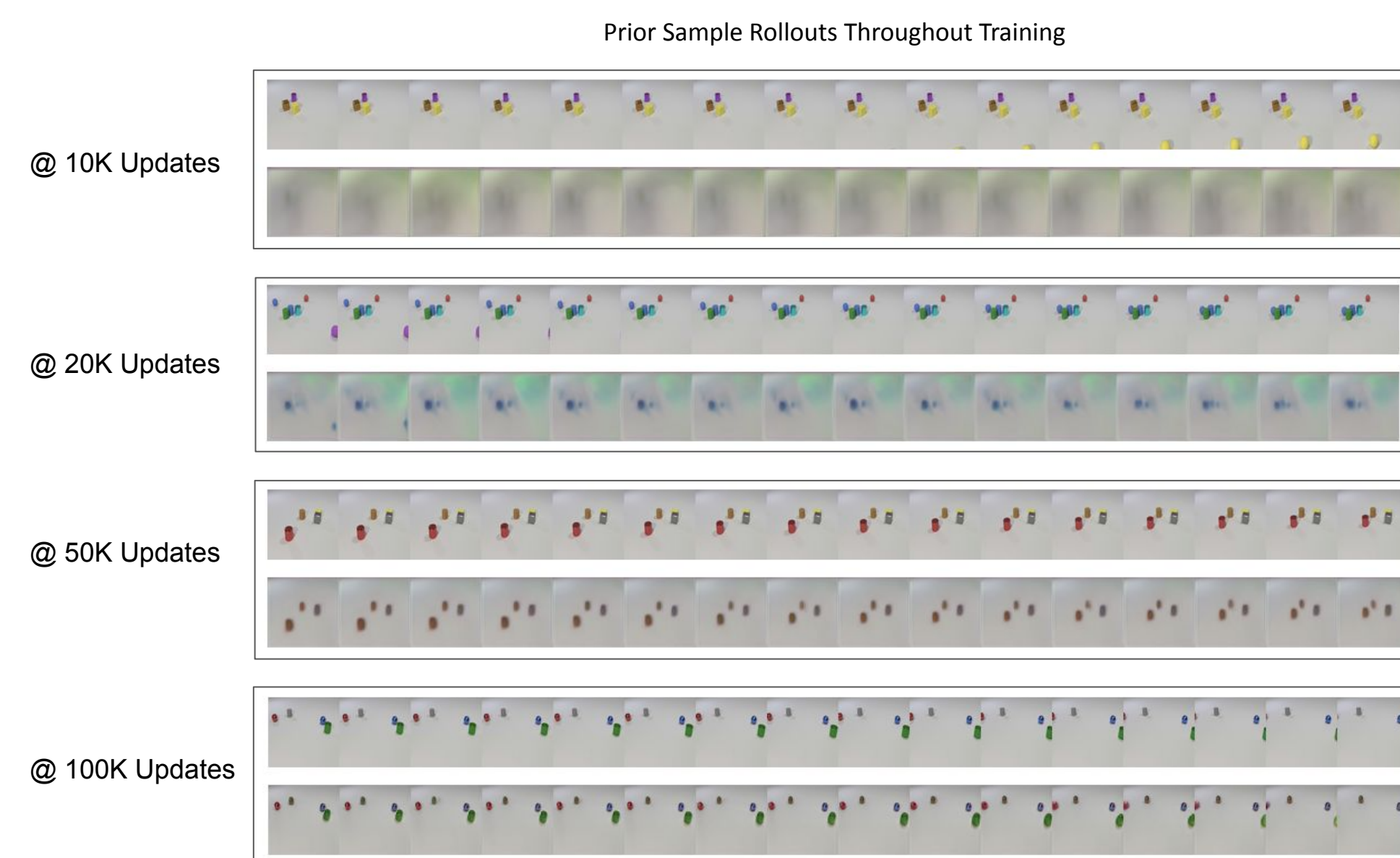


Table 1: CLEVRER per Question Accuracies (Chen et al. [2021])

Methods	Extra Labels Attr.	Descriptive Prog.	Predictive	Explanatory	Counterfactual
CNN+MLP			48.4	18.3	13.2
CNN+LSTM			51.8	17.5	31.6
Memory	No	No	54.7	13.9	33.1
HCRN			55.7	21.0	21.0
MAC (V)			85.6	12.5	16.5
<b>Slot Transformer (Ours)</b>			<b>87.4</b>	<b>48.3</b>	<b>65.3</b>
TVQA+ MAC (V+)	Yes	No	72.0	<b>23.7</b>	<b>48.9</b>
			<b>86.4</b>	22.3	42.9
IEP (V)			52.8	14.5	9.7
TbD-net (V)	No	Yes	79.5	3.8	6.5
DCL			<b>90.7</b>	<b>82.8</b>	<b>82.0</b>
NS-DR			88.1	79.6	68.7
NS-DR (NE)	Yes	Yes	85.8	74.3	54.1
DCL-Oracle			<b>91.4</b>	<b>82.0</b>	<b>82.1</b>

Table 2: CATER Results on Task 3: "Localization"

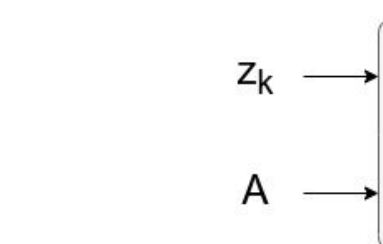
Methods	Top-1	Top-5	Static Camera	Static Camera	Moving Camera	Moving Camera
			Top-1	Top-5	Top-1	Top-5
TimeSformer (+pre-train)	<b>82.2</b>	<b>95.6</b>				
ISD (+pre-train)	71.7	90.4				
SimCLR (+pre-train)	51.6	-	2.8	13.8	3.9	-
ImageNet (+pre-train)	54.7	-	60.2	81.8	1.2	28.6
CVRL (+pre-train)	72.9	-	46.2	69.9	1.5	38.6
Aloe Ding et al. [2020]			<b>74.0</b>	<b>94.0</b>	<b>0.44</b>	<b>59.7</b>
<b>Slot Transformer (Ours)</b>	<b>68.2</b>	<b>87.6</b>	<b>62.9</b>	<b>84.7</b>	<b>0.86</b>	<b>35.8</b>



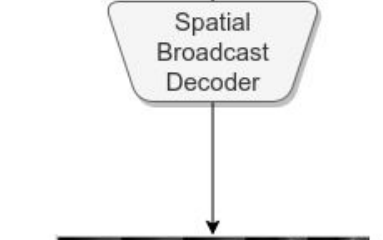
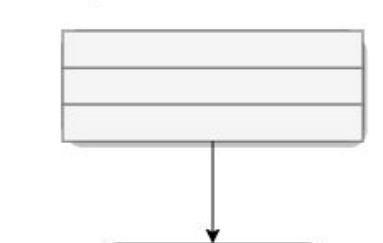
## Discussion & Conclusions

- Scene and sequence representation.
  - Masks + slot ablation show reasonable segmentation.
- More Iteration has a positive effect on the model's ability to generalize
- Auxiliary Losses help with both CLEVRER & CATER test performance

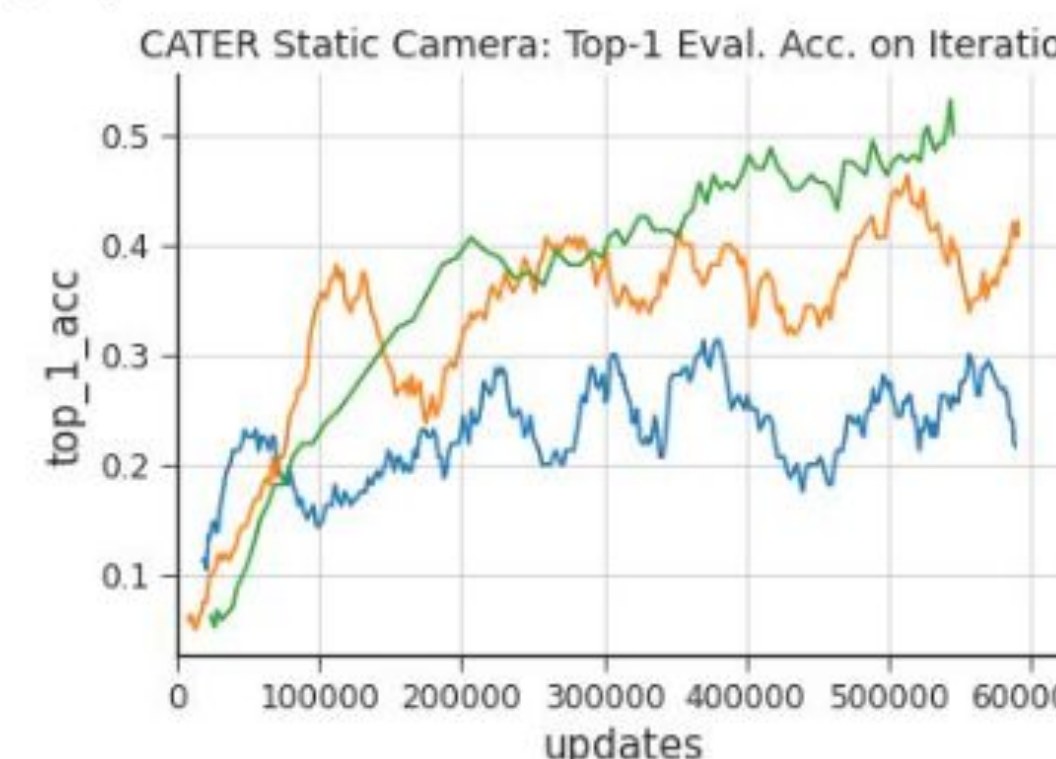
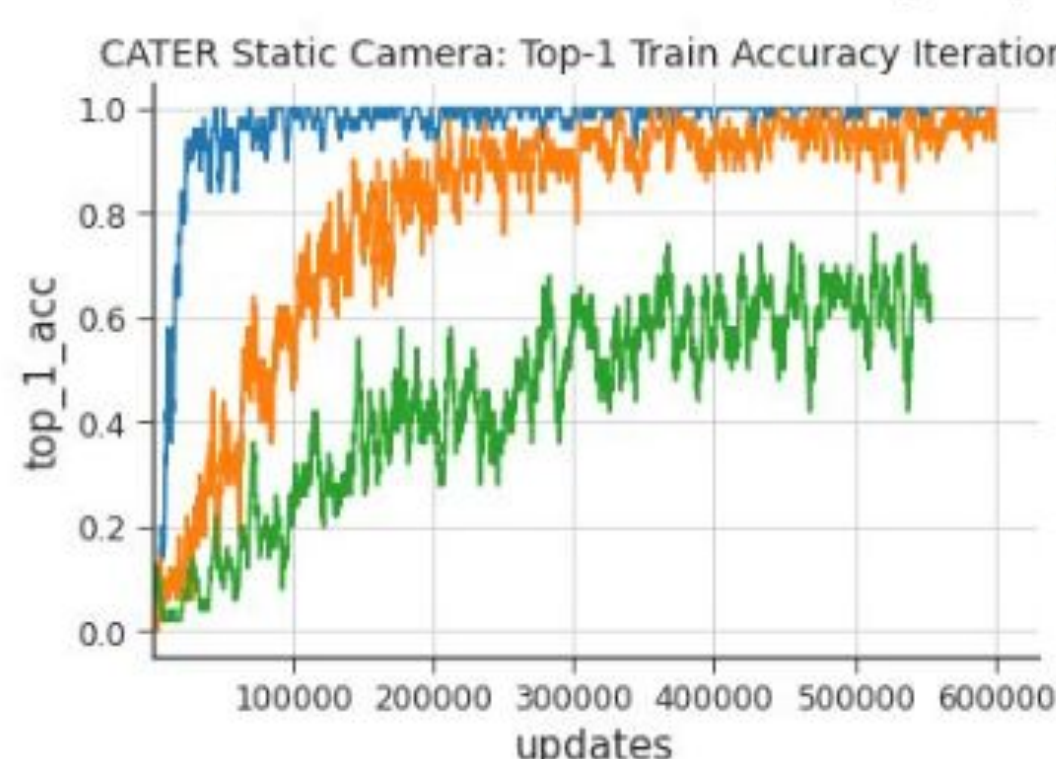
Question Prediction



Object Prediction



Train (Left) vs. Test (Right) over iterations:



Ablate Attention

