



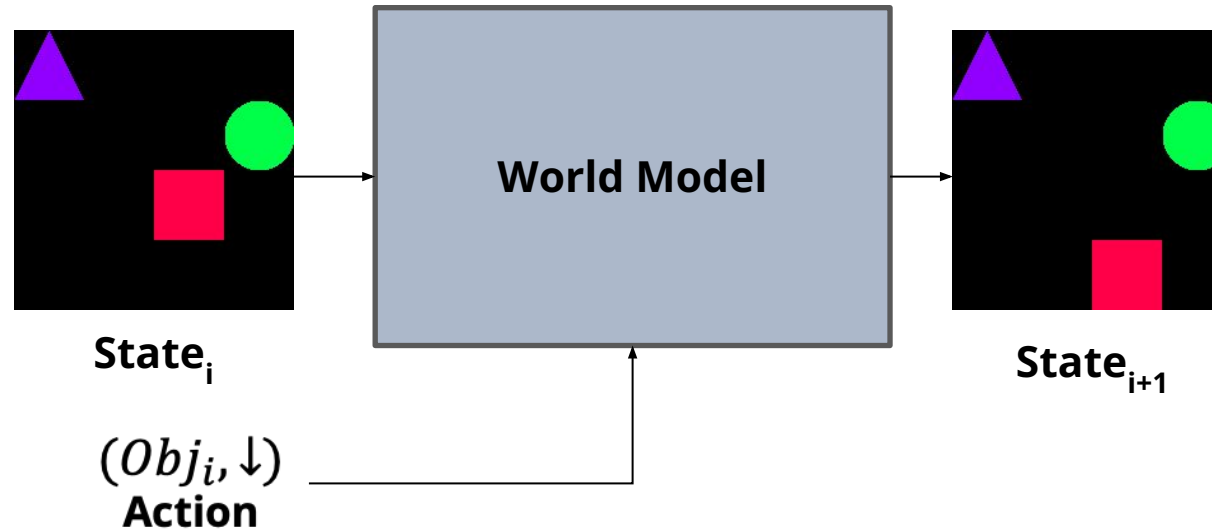
Neurosymbolic Grounding for Object-Oriented Compositional World Models

Atharva Sehgal, Arya Grayeli, Jennifer Sun*, and Swarat Chaudhuri

UT Austin, Caltech*

TL;DW: <https://bit.ly/cosmos-wm>

World Modeling



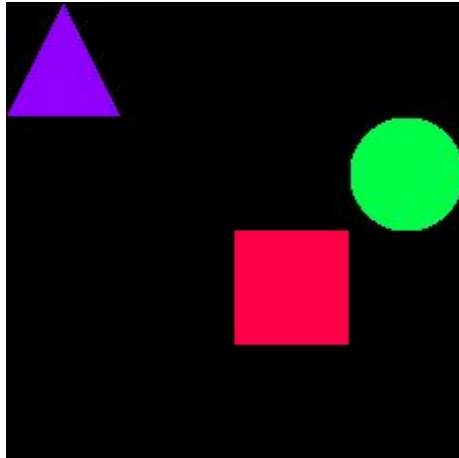
Objective: Implicitly learn the dynamics laws of this domain.

Recurrent world models facilitate policy evaluation. Ha & Schmidhuber, 2018.

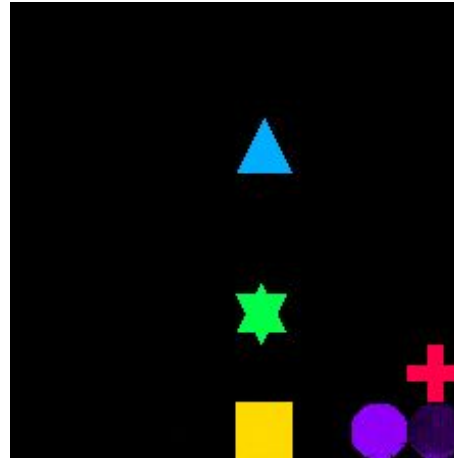
Neural production systems: Learning rule-governed visual dynamics. Goyal, Bengio et al., 2021.

Compositional Generalization

Assuming 5 colors for each object

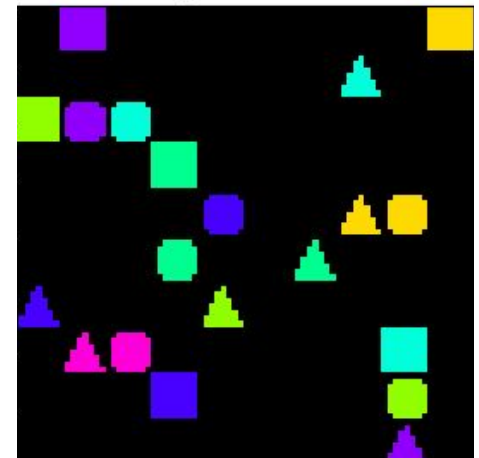


|objects| = $3*5$
#combinations = **455**



|objects| = $5*5$
#combinations = **53130**

...



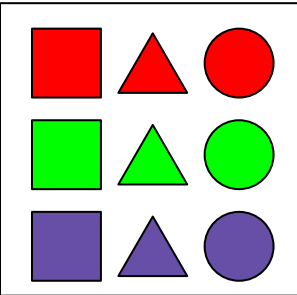
|objects| = $20*5$
#combinations = **5.36×10^{20}**

Compositional Generalization

Entity Composition

Generalization to composition of shapes not seen together during training

Distribution of
Atoms (\mathcal{A})

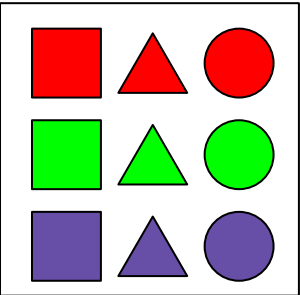


Relational Composition

Entity Composition + Objects with shared attributes have shared dynamics.

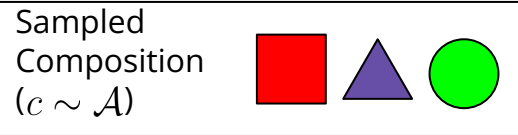
Compositional Generalization

Distribution of
Atoms (\mathcal{A})



Entity Composition

Generalization to composition of shapes not seen
together during training



Relational Composition

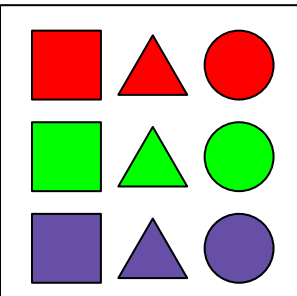
Entity Composition + Objects with shared attributes have shared
dynamics.

Compositional Generalization

Entity Composition

Generalization to composition of shapes not seen together during training

Distribution of Atoms (\mathcal{A})



Sampled Composition ($c \sim \mathcal{A}$)

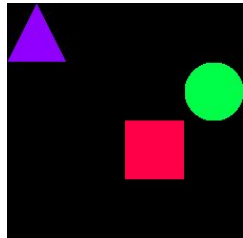
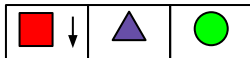


Image ($I^{(t)}$)



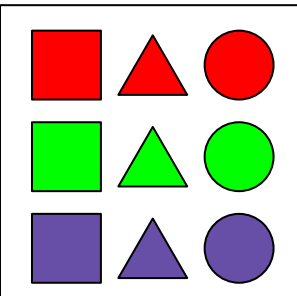
Actions

Relational Composition

Entity Composition + Objects with shared attributes have shared dynamics.

Compositional Generalization

Distribution of Atoms (\mathcal{A})



Entity Composition

Generalization to composition of shapes not seen together during training

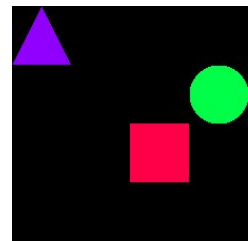
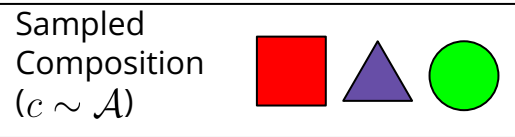
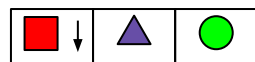
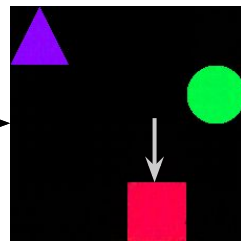


Image ($I^{(t)}$)



Actions



Ground Truth
Next Image
 $I^{(t+1)}$



The object dynamics are invariant to distribution shift!

Relational Composition

Entity Composition + Objects with shared attributes have shared dynamics.

Compositional Generalization

Entity Composition

Generalization to composition of shapes not seen together during training

Sampled Composition
($c \sim \mathcal{A}$)

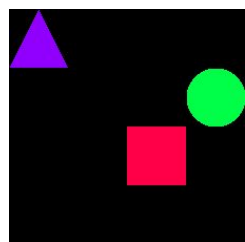
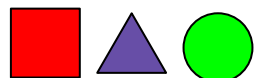
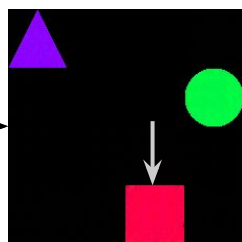


Image ($I^{(t)}$)



Actions



Ground Truth
Next Image
 $I^{(t+1)}$

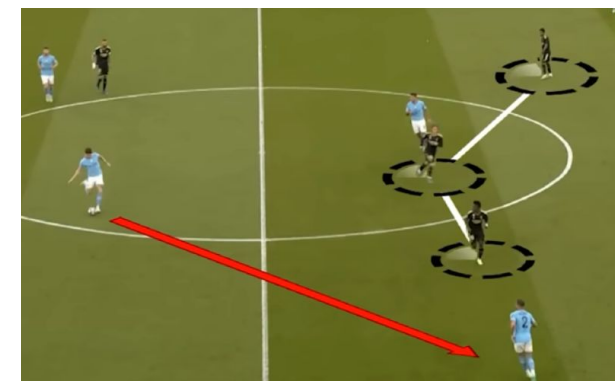


The object dynamics are invariant to distribution shift!

Relational Composition

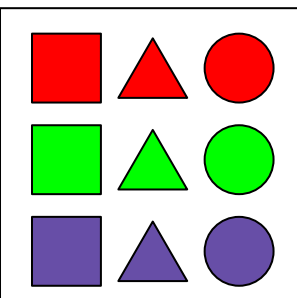
Entity Composition + Objects with shared attributes have shared dynamics.

Higher level rule: same color => cooperation



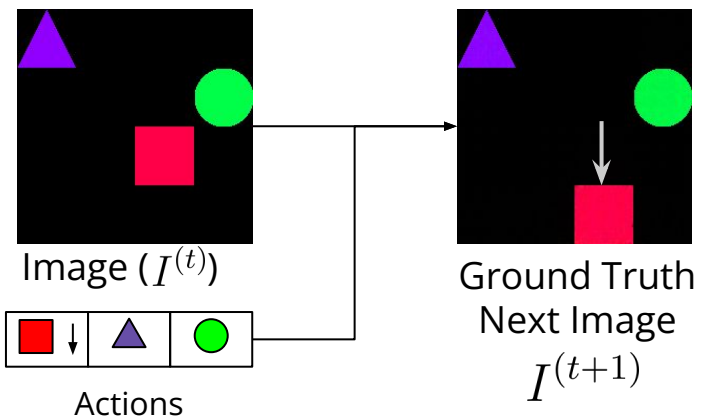
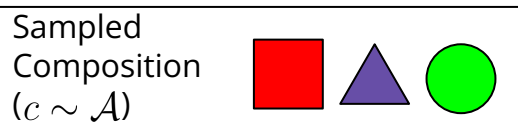
Compositional Generalization

Distribution of Atoms (\mathcal{A})



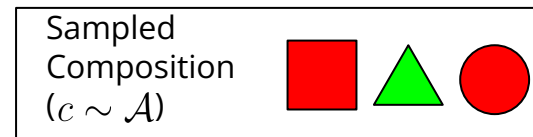
Entity Composition

Generalization to composition of shapes not seen together during training



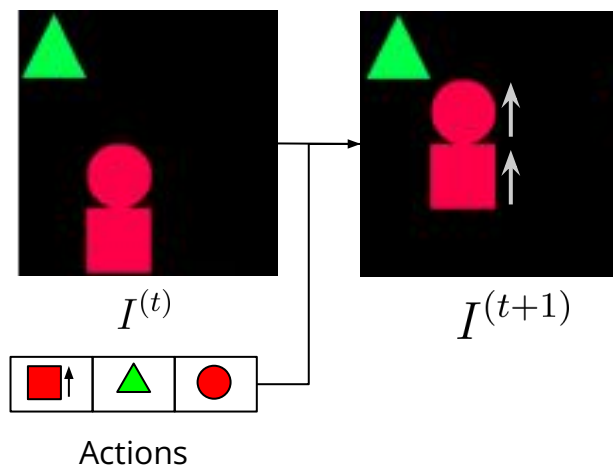
Relational Composition

Entity Composition + Objects with shared attributes have shared dynamics.



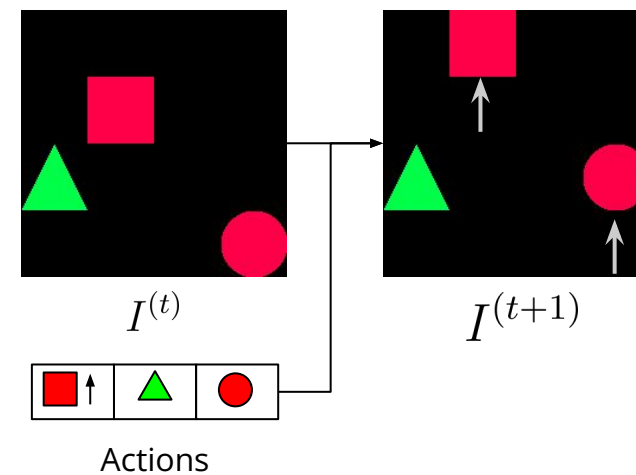
"Sticky"

Shared Attributes = Color + Adjacency



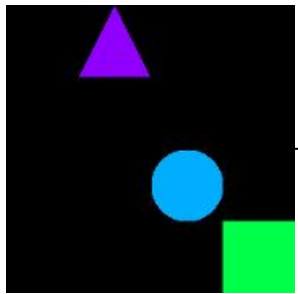
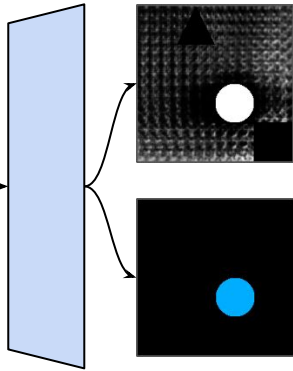
"Team"

Shared Attributes = Color

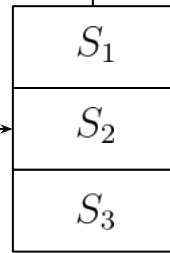
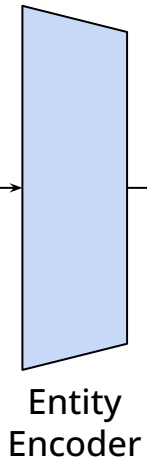


Symbolic Labelling: Preprocessing

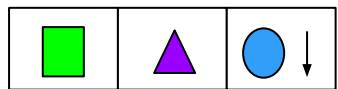
For each slot



Image

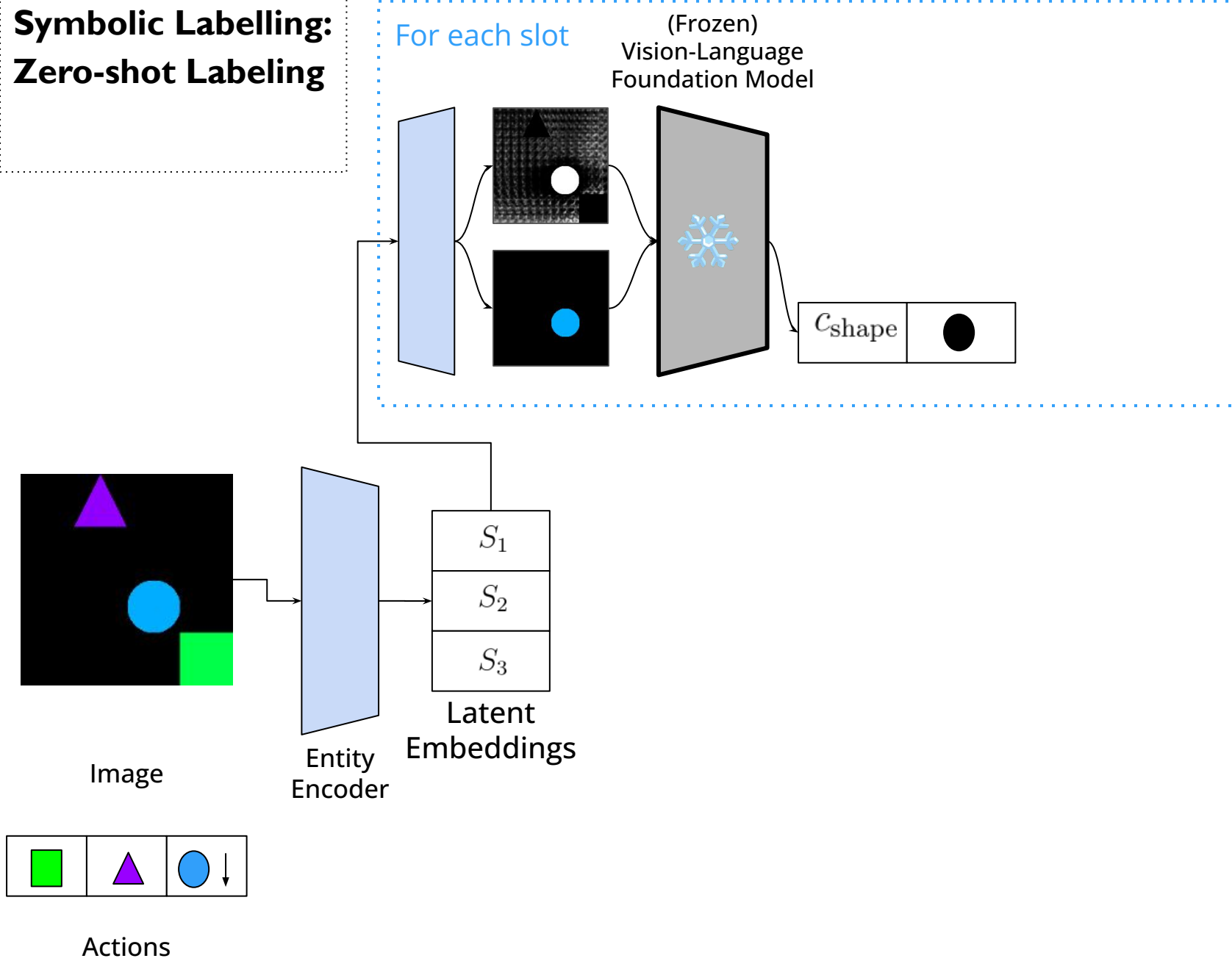


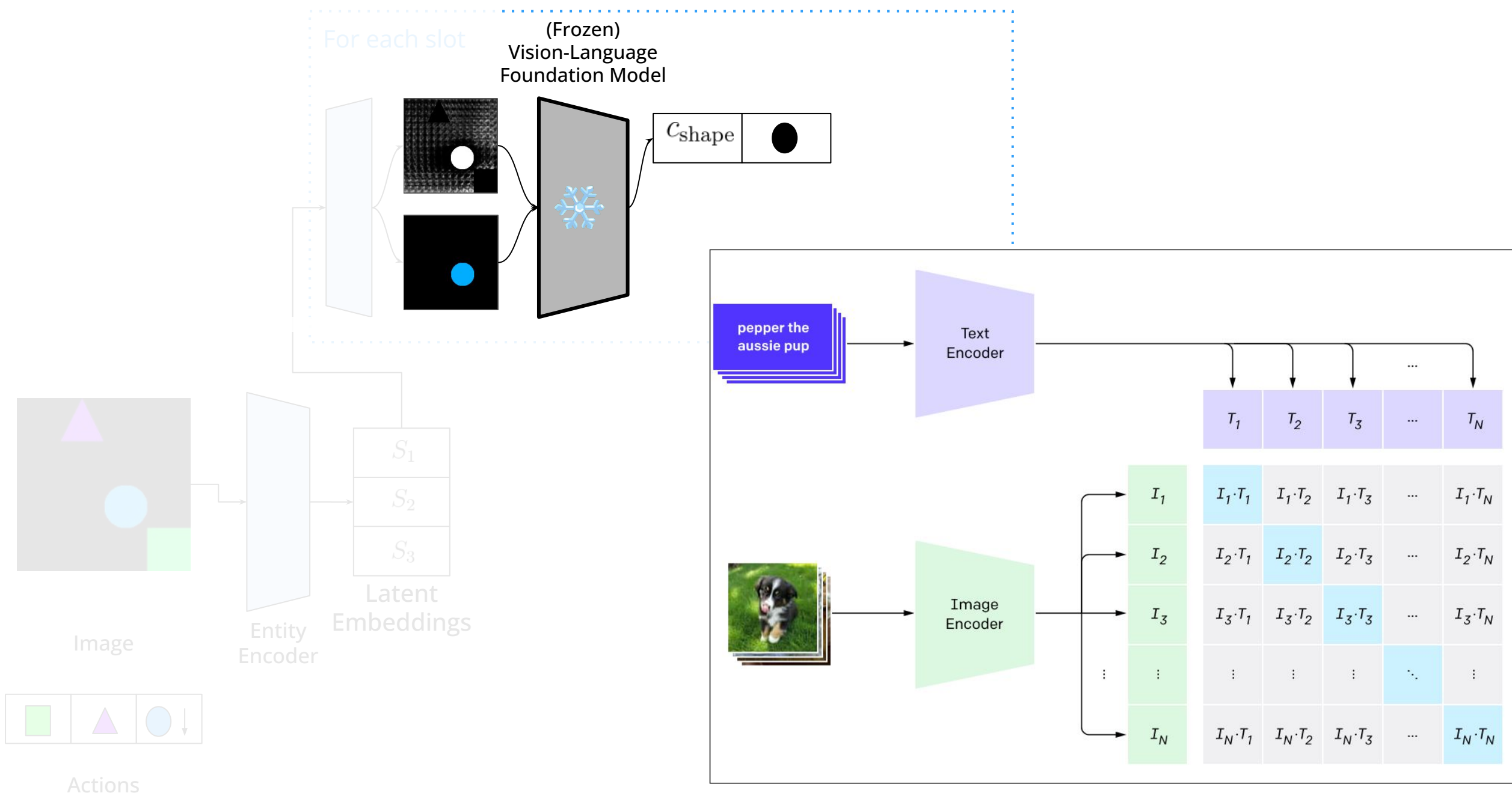
Latent
Embeddings



Actions

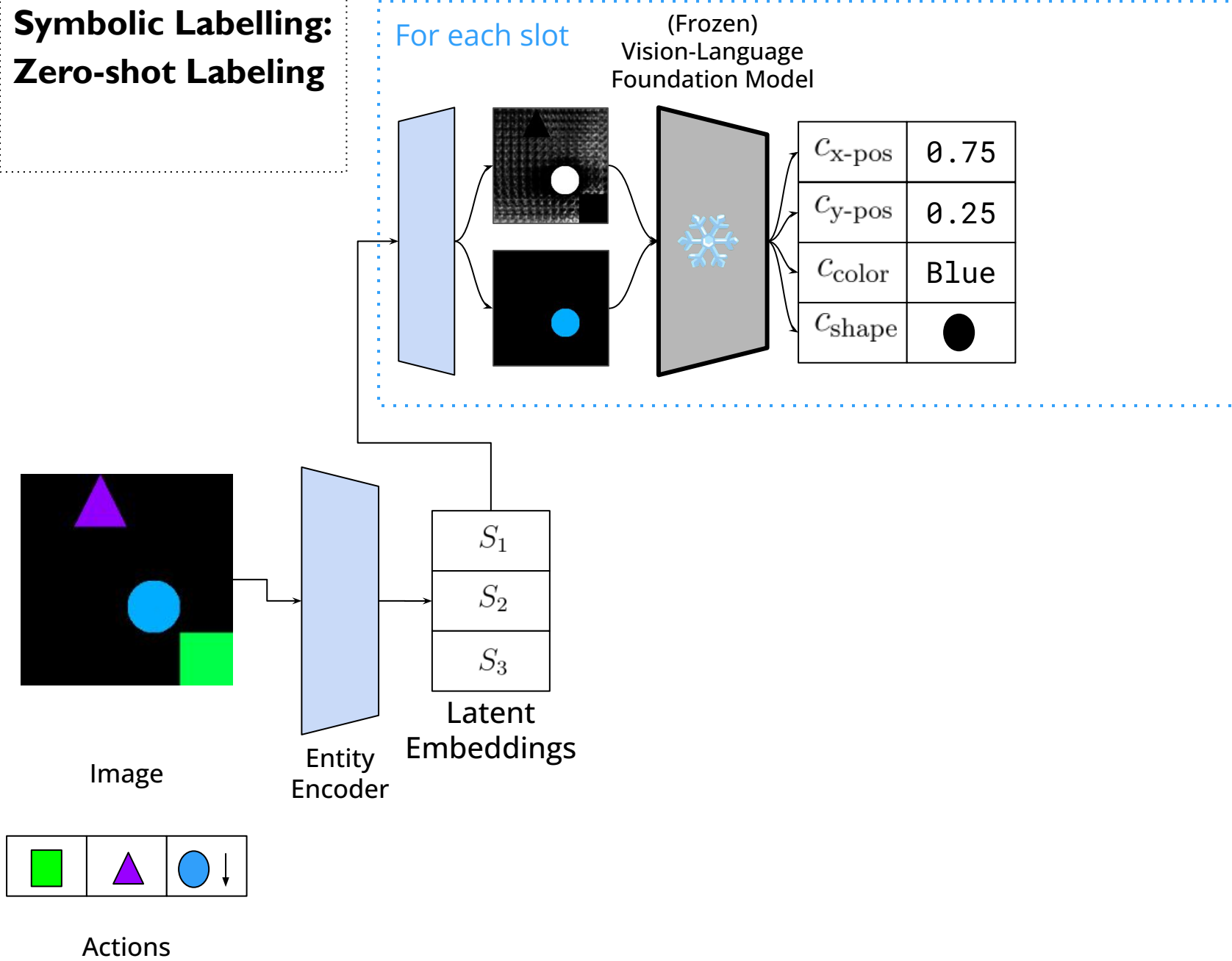
Symbolic Labelling: Zero-shot Labeling



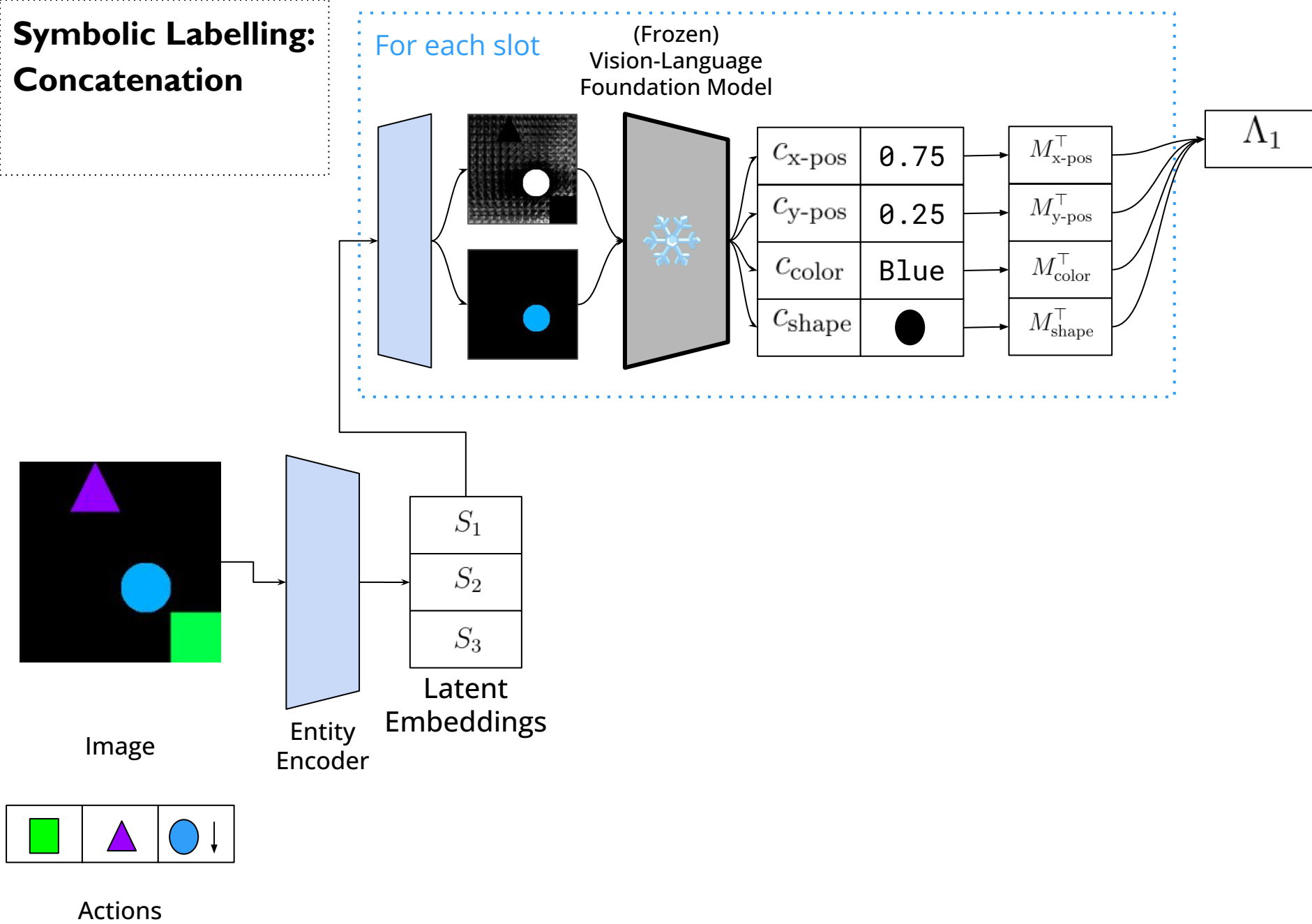


Learning Transferable Visual Models From Natural Language Supervision. Radford et al., 2021.

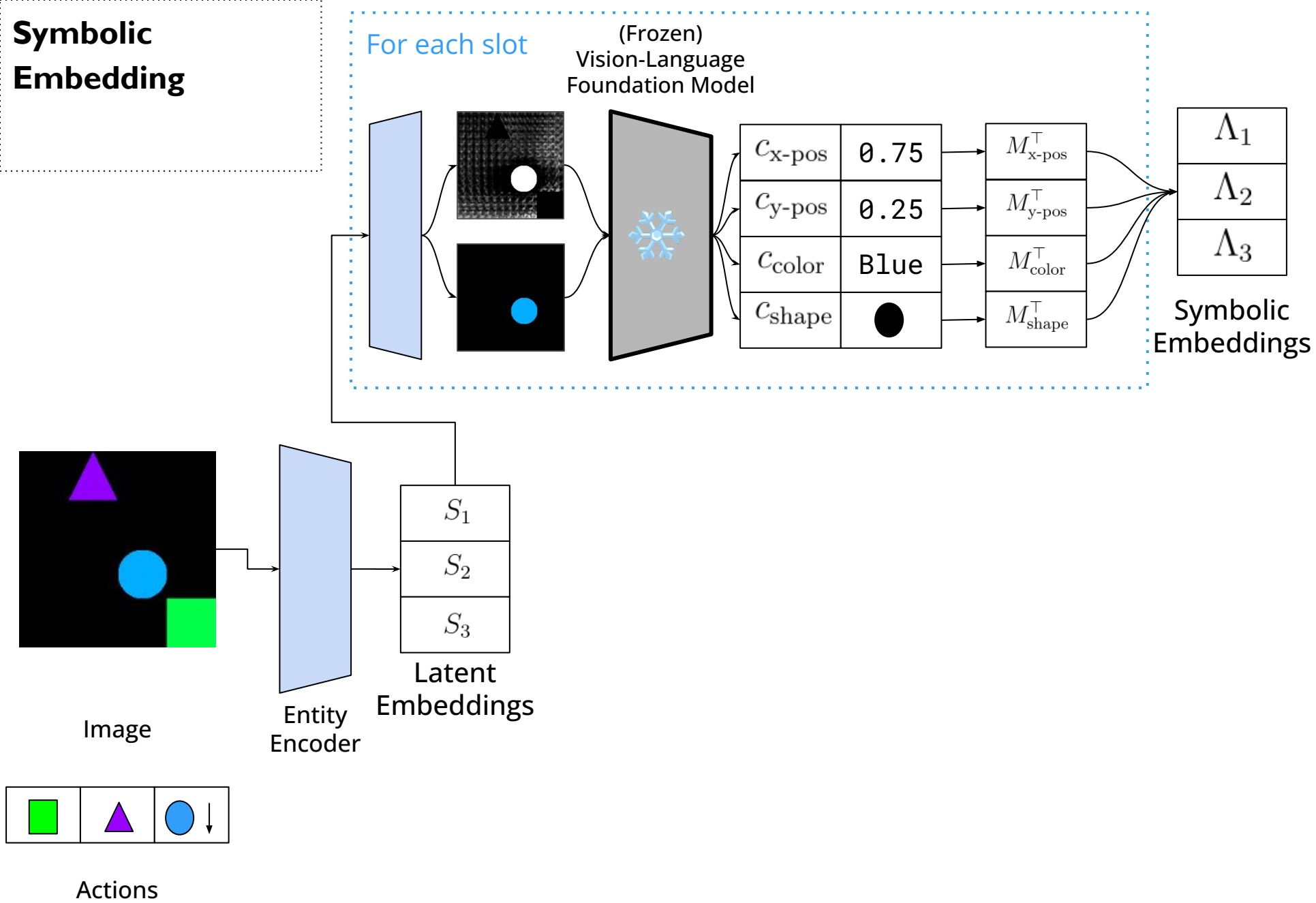
Symbolic Labelling: Zero-shot Labeling



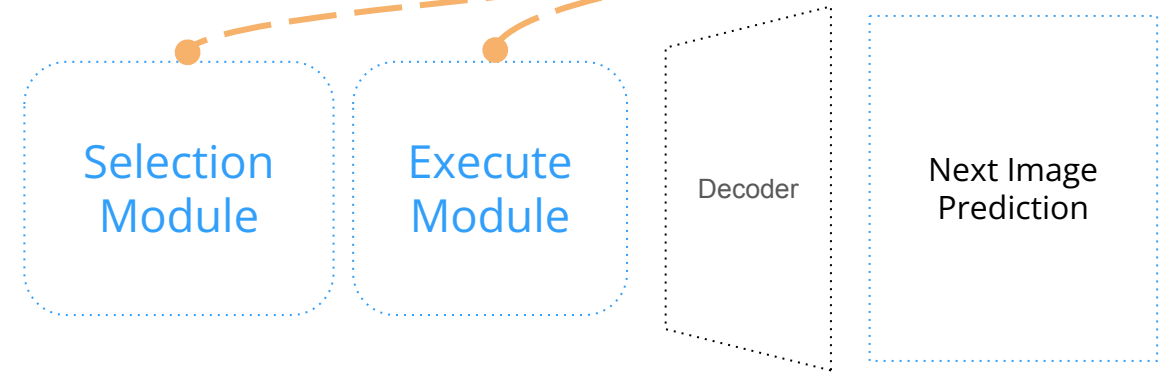
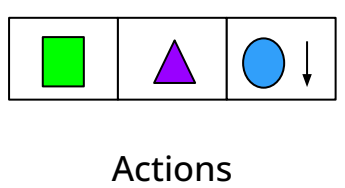
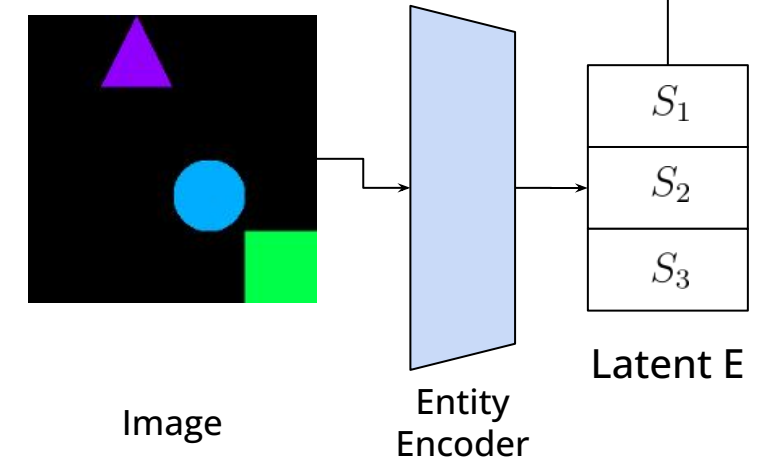
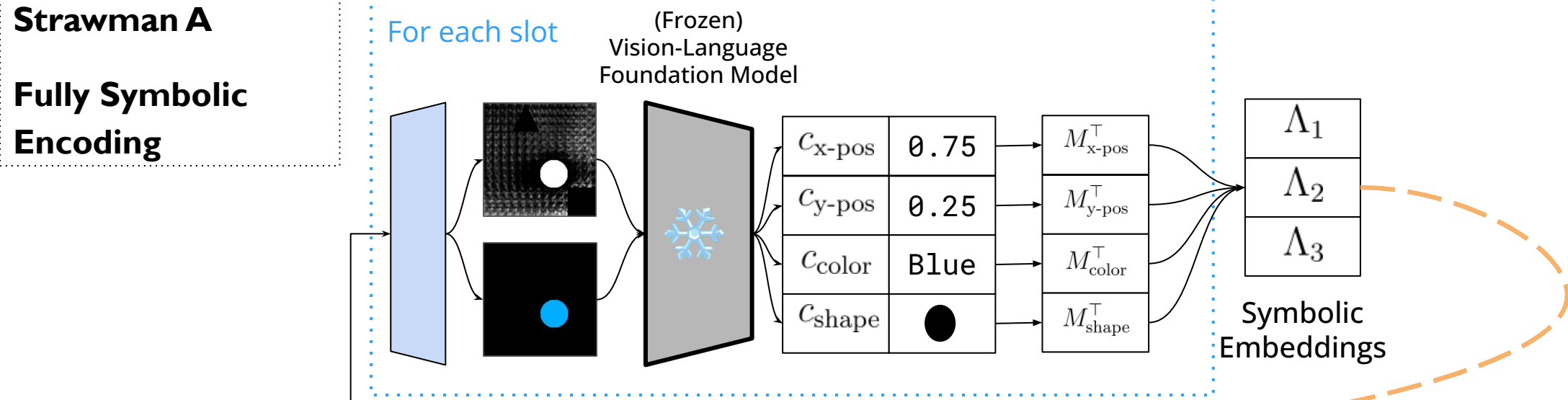
Symbolic Labelling: Concatenation



Symbolic Embedding



Strawman A
Fully Symbolic Encoding



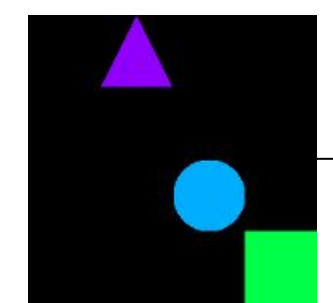
Q: What if we only use Symbolic Embeddings?

✓: Trivially generalizes to attribute compositions. Robust selection module!

✗: Symbols bottleneck expressivity. The circle is *Cyan*, not *Blue*! Erroneous image reconstruction.

Strawman B

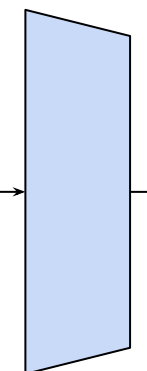
Fully Neural Encoding



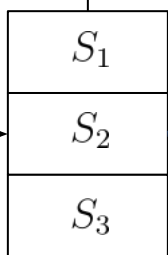
Image



Actions



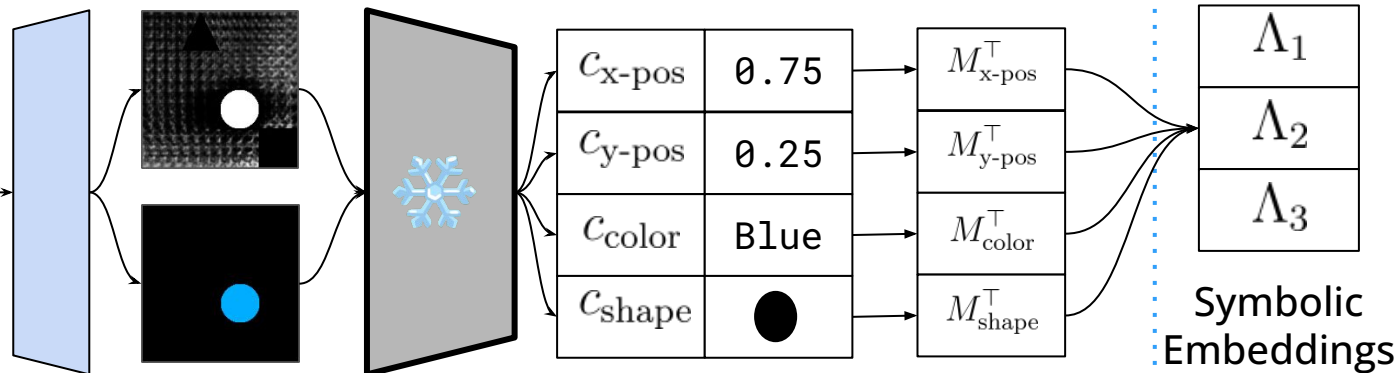
Entity Encoder



Latent Embeddings

For each slot

(Frozen)
Vision-Language
Foundation Model



Selection Module

Execute Module

Decoder

Next Image Prediction

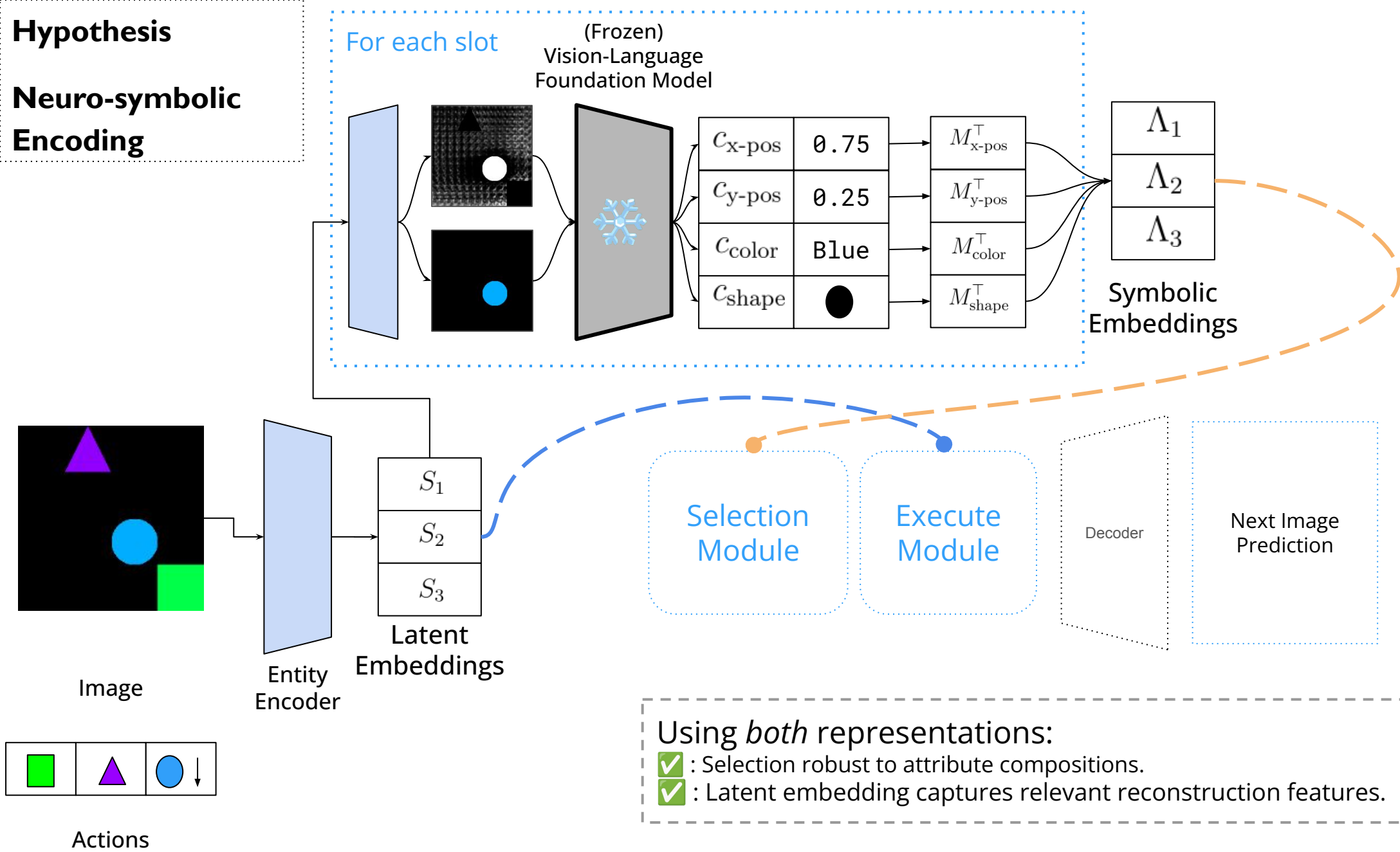
Q: What if we only use Latent Embeddings?

: Tendency to overfit to attribute compositions seen during training.

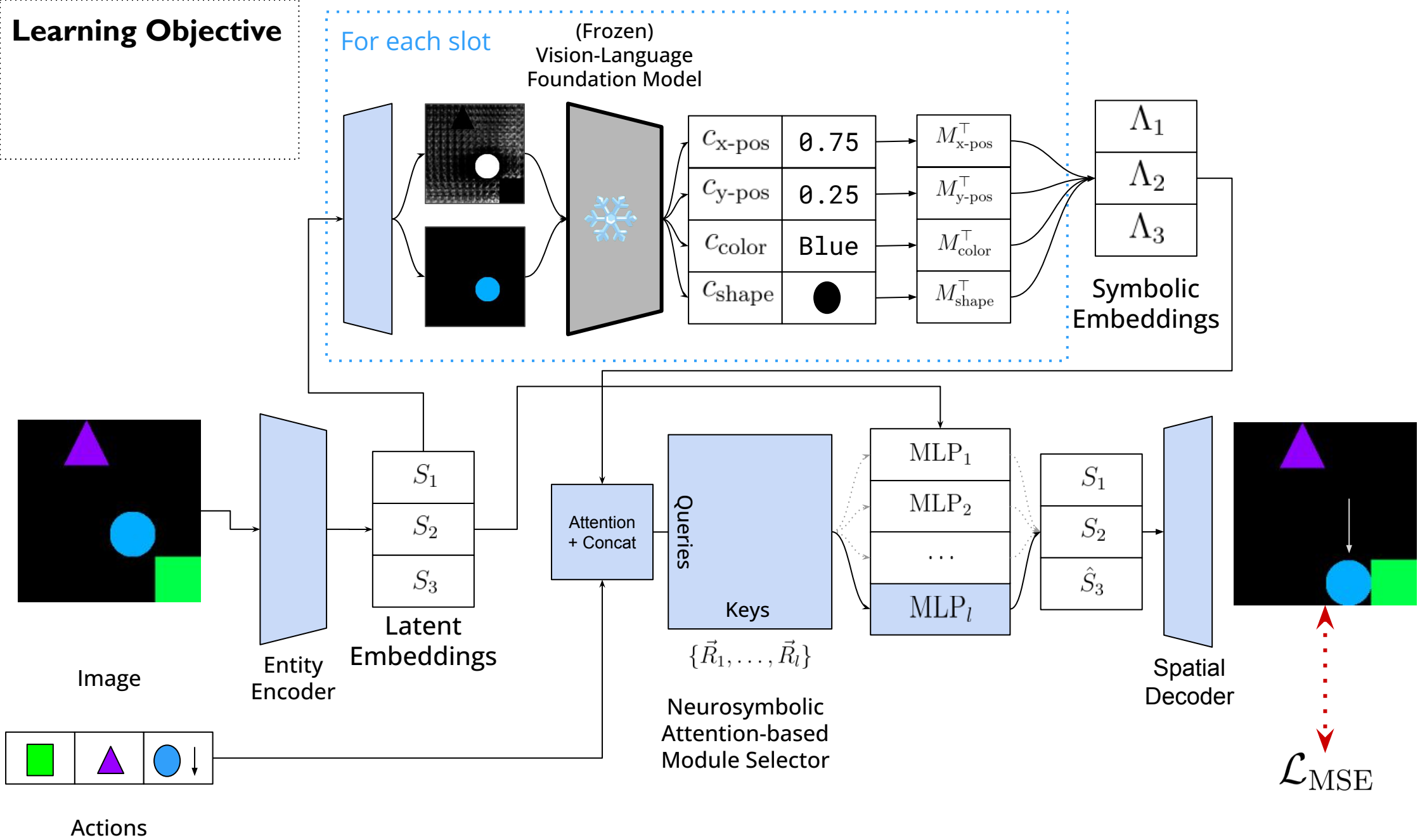
: Latent embedding contains very expressive features. Good image reconstructions!

Hypothesis

Neuro-symbolic Encoding

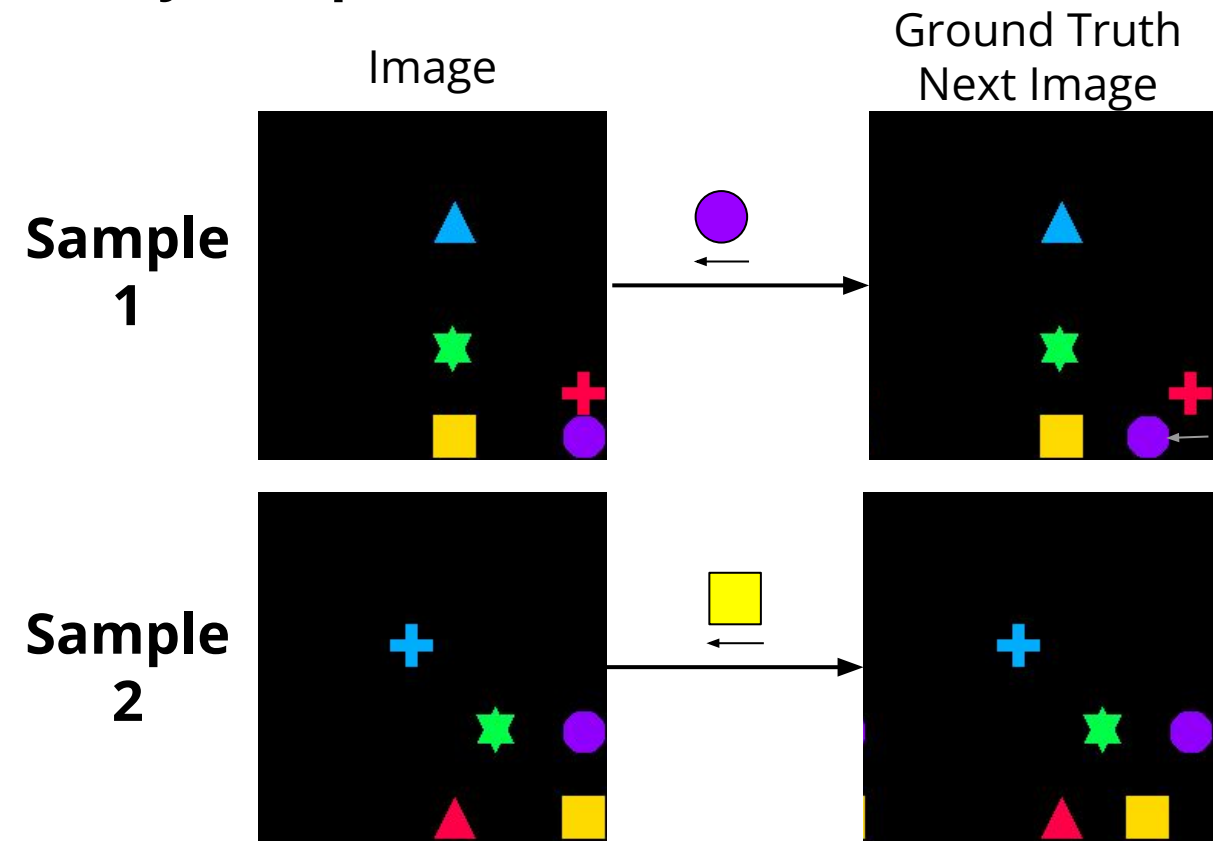


Learning Objective



Results

Entity Composition

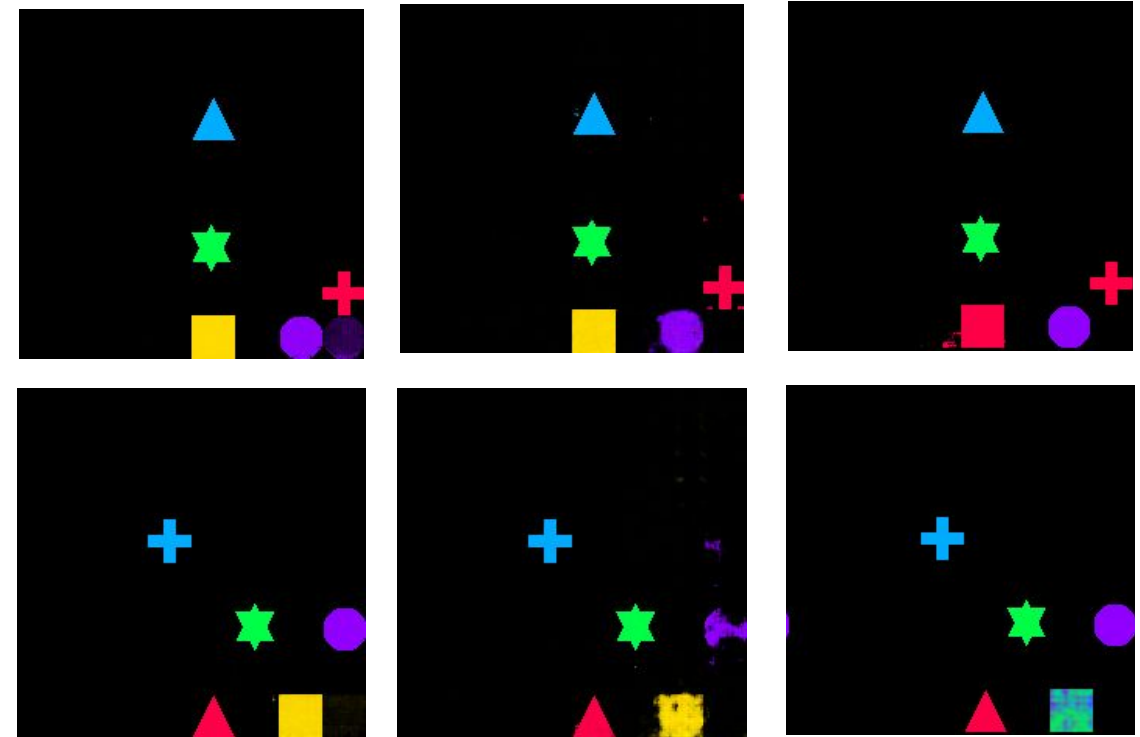


Predicted Next Image

COSMOS

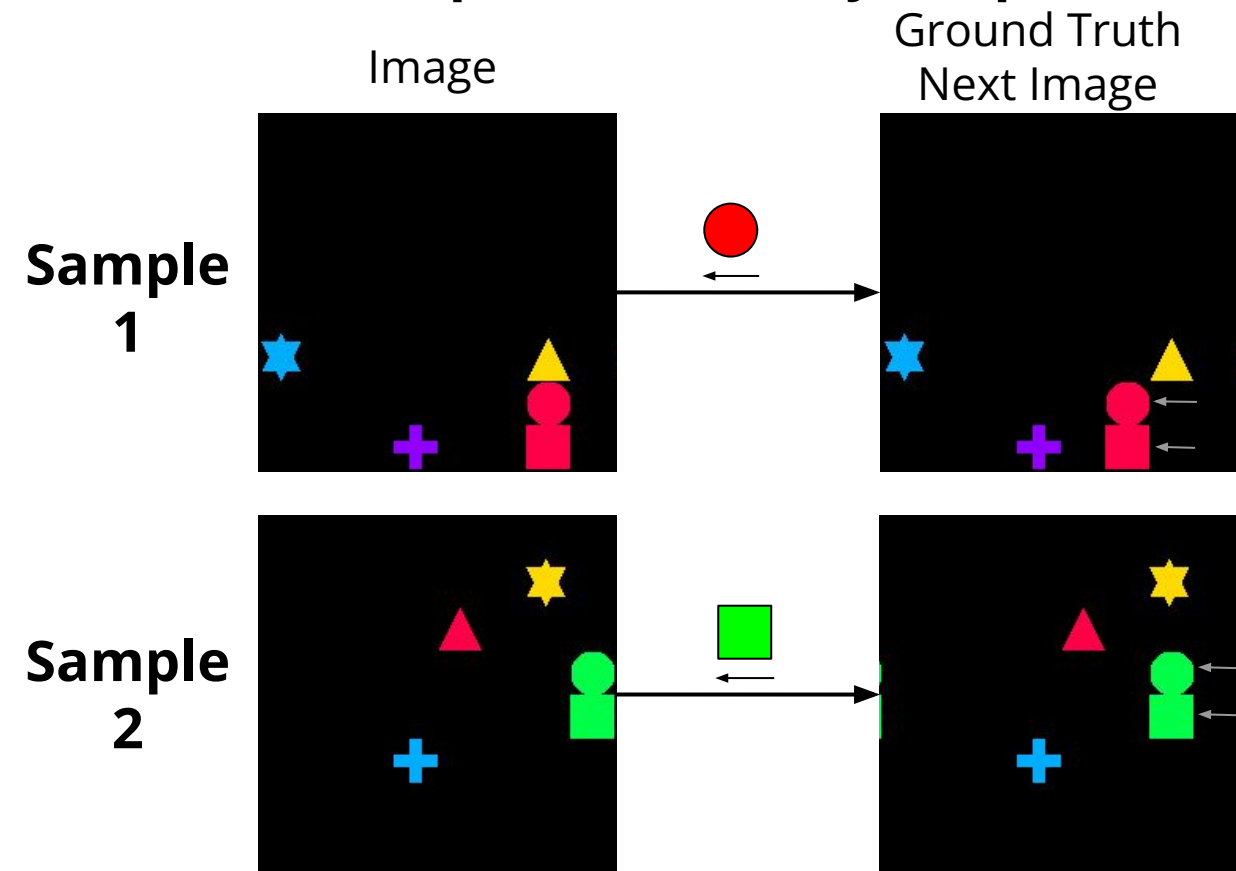
AlignedNPS

GNN



Results

Relational Composition (Sticky Shapeworld)

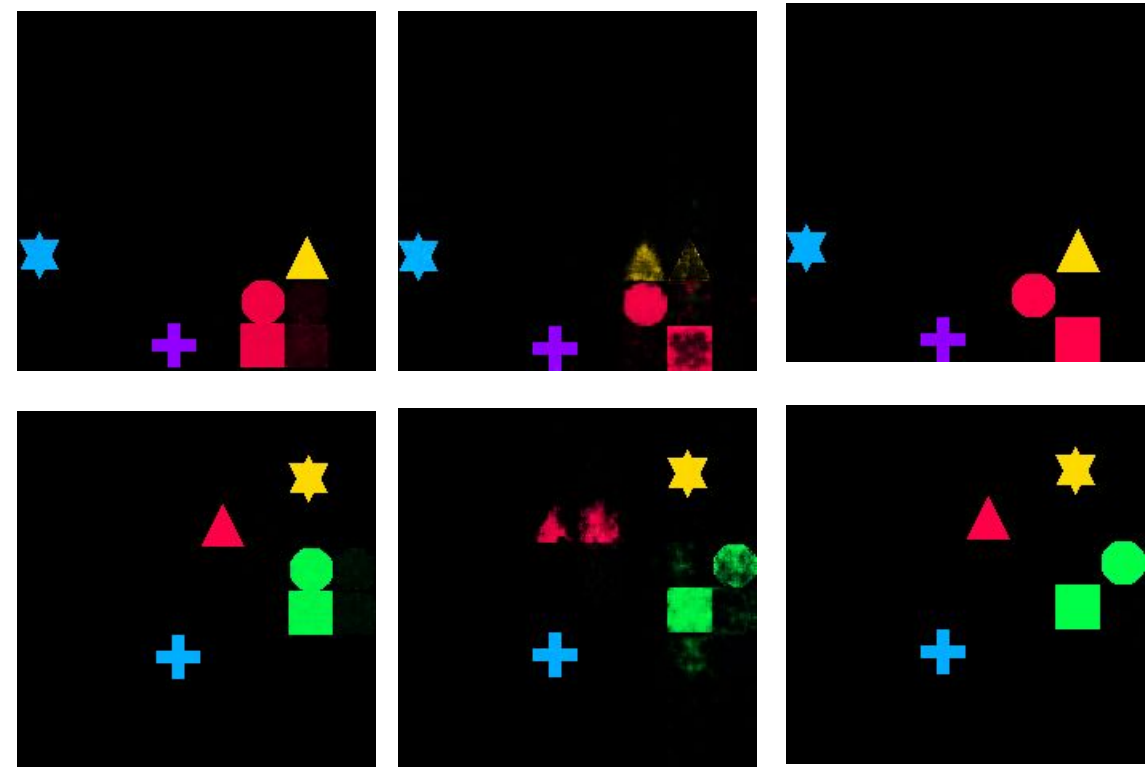


Predicted Next Image

COSMOS

AlignedNPS

GNN



Takeaways

- **Explicit symbolic knowledge** helps with compositionality
- **Extend**, rather than replace, deep representations
- **Foundation models** over language (and code) give symbols for free

More Information:

- <https://bit.ly/cosmos-wm>

