

# Generative Neuro-Symbolic Models of Concept Learning

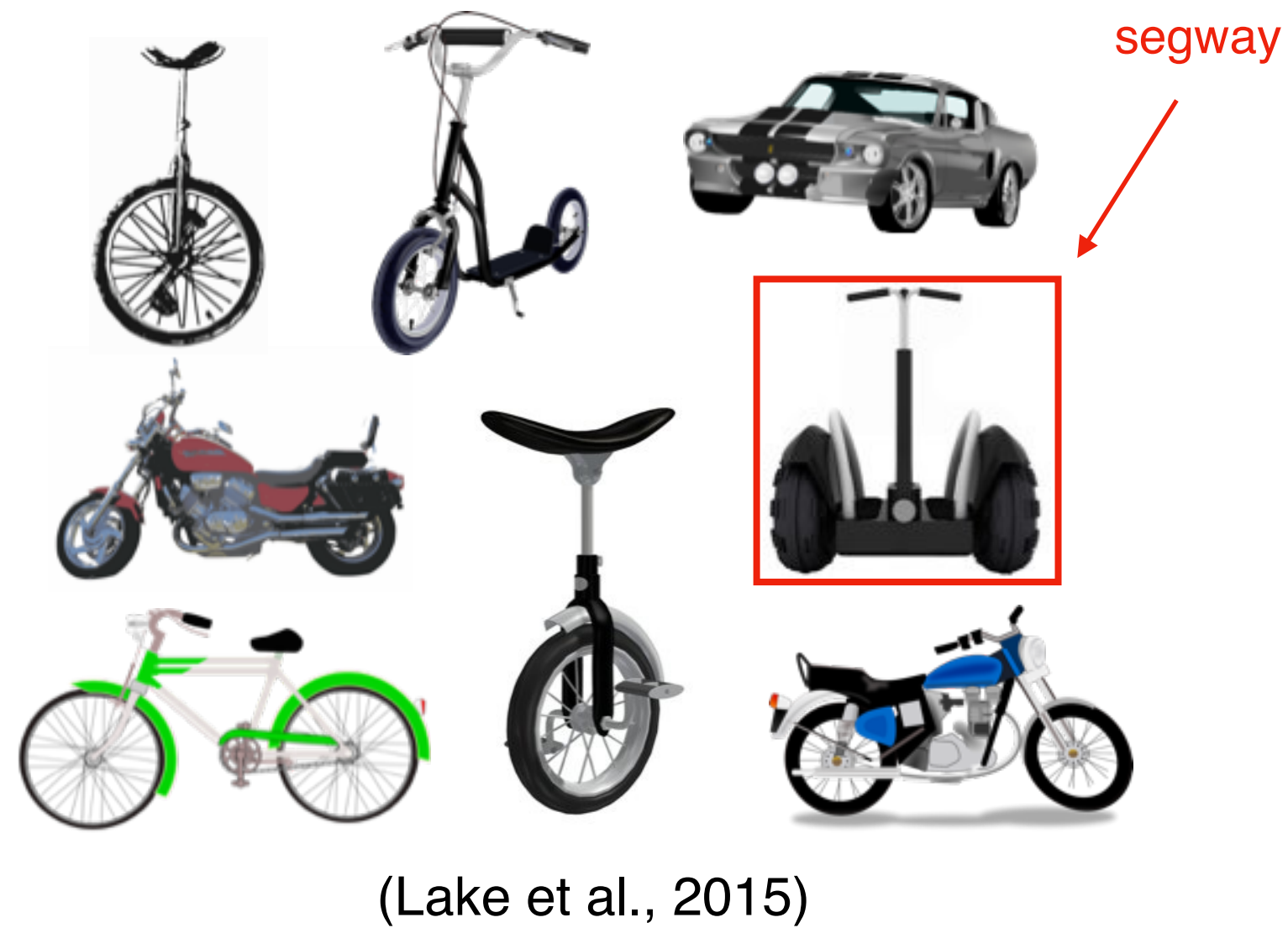
Reuben Feinman

advised by  
Brenden Lake

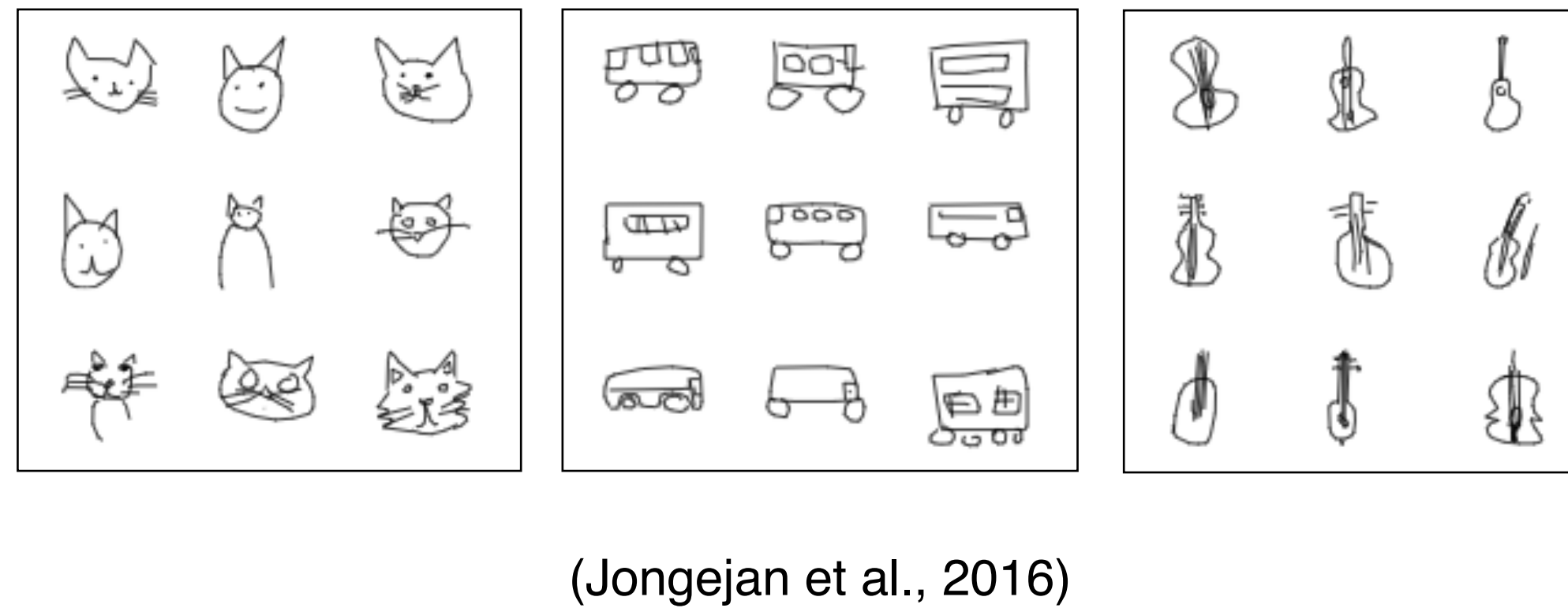


# Human concepts are *task-general*

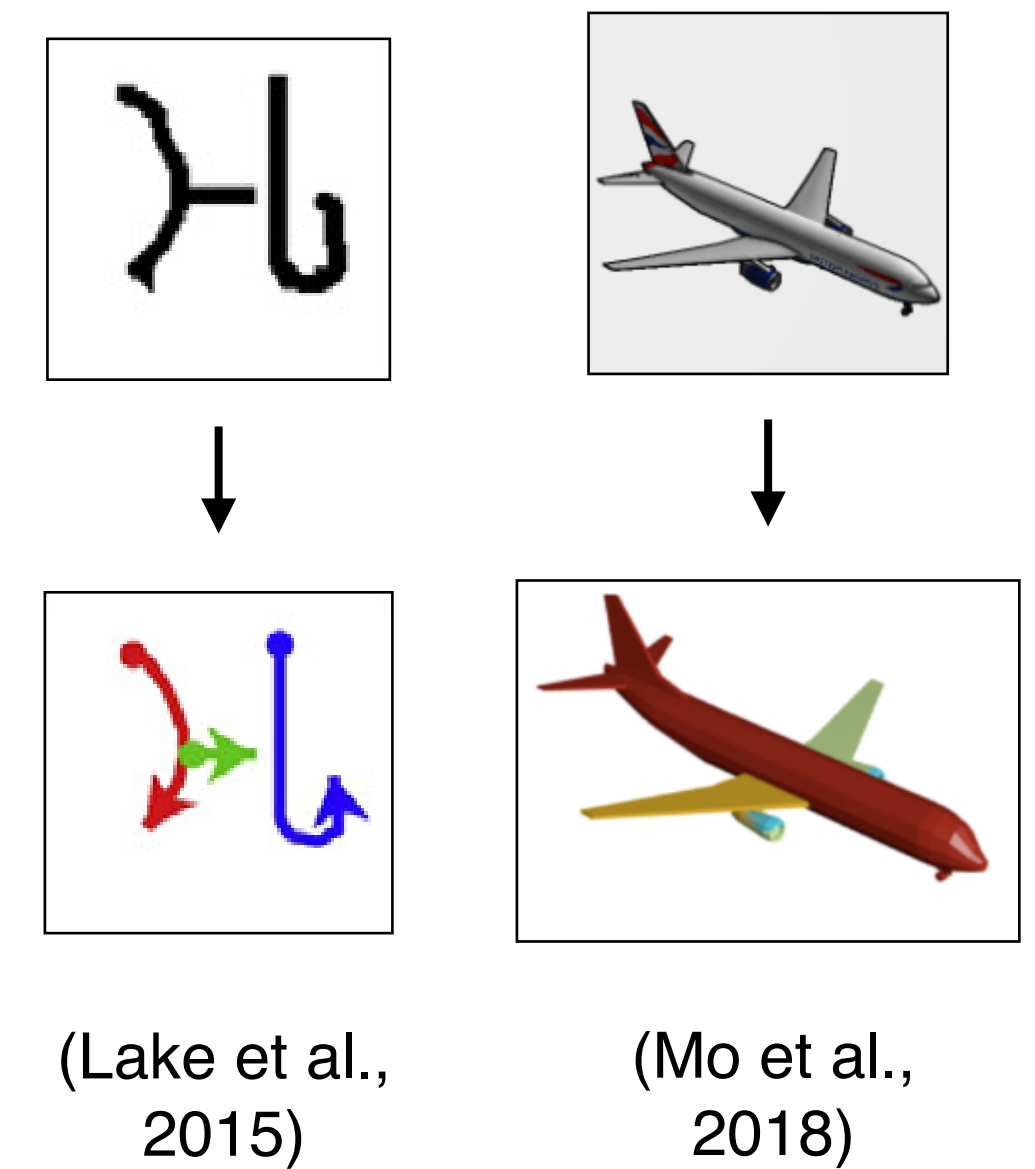
## Recognition



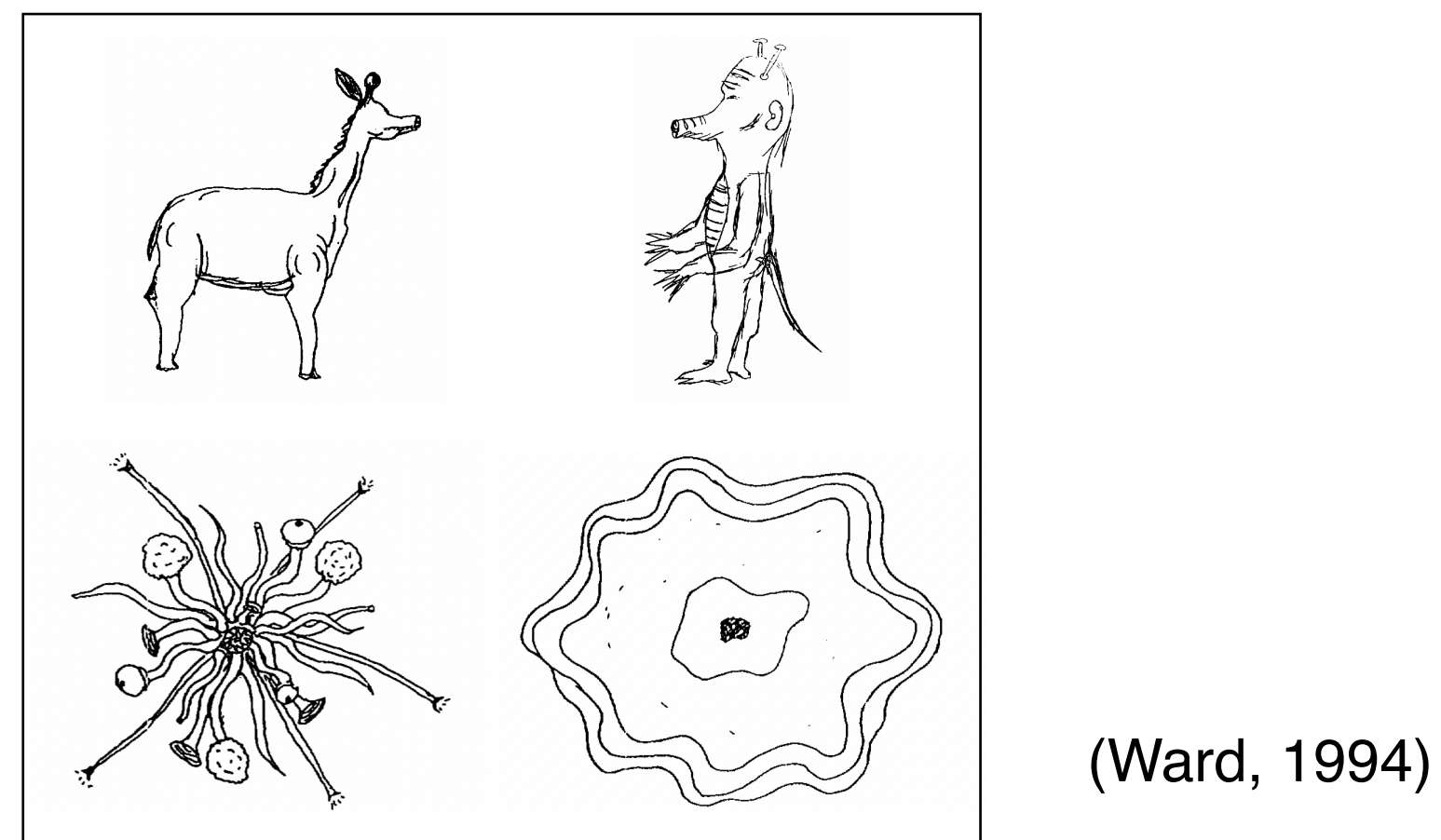
## Generation



## Parsing



## Imagination





# Human concept learning is *fast*

This is a  
“breakfast machine.”



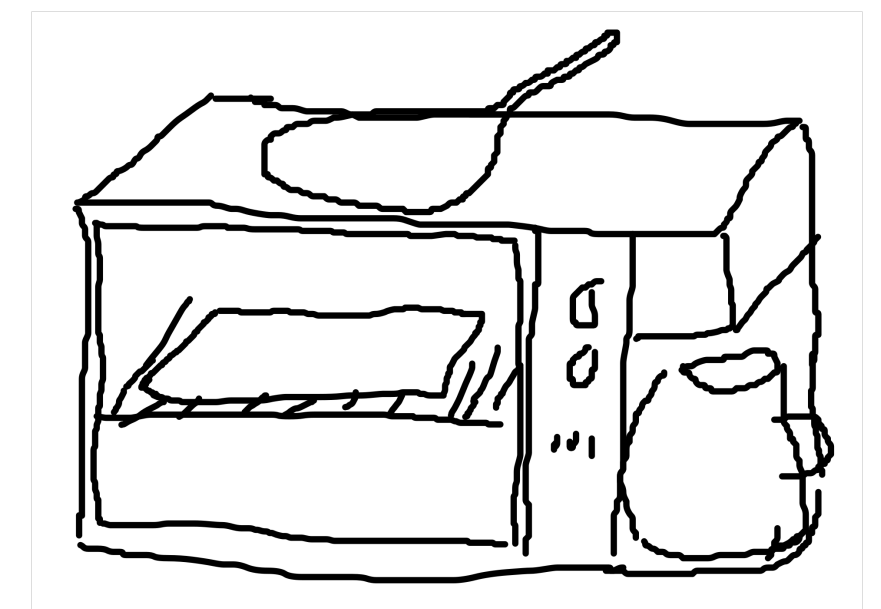
Which is another?



What are its parts?



Create a new one.





# Research questions

- What is the structure of human conceptual representations? How does this structure support a variety of discriminative and generative abilities?
- How do people acquire such rich representations from so little experience?
- How can we understand these abilities in computational terms?

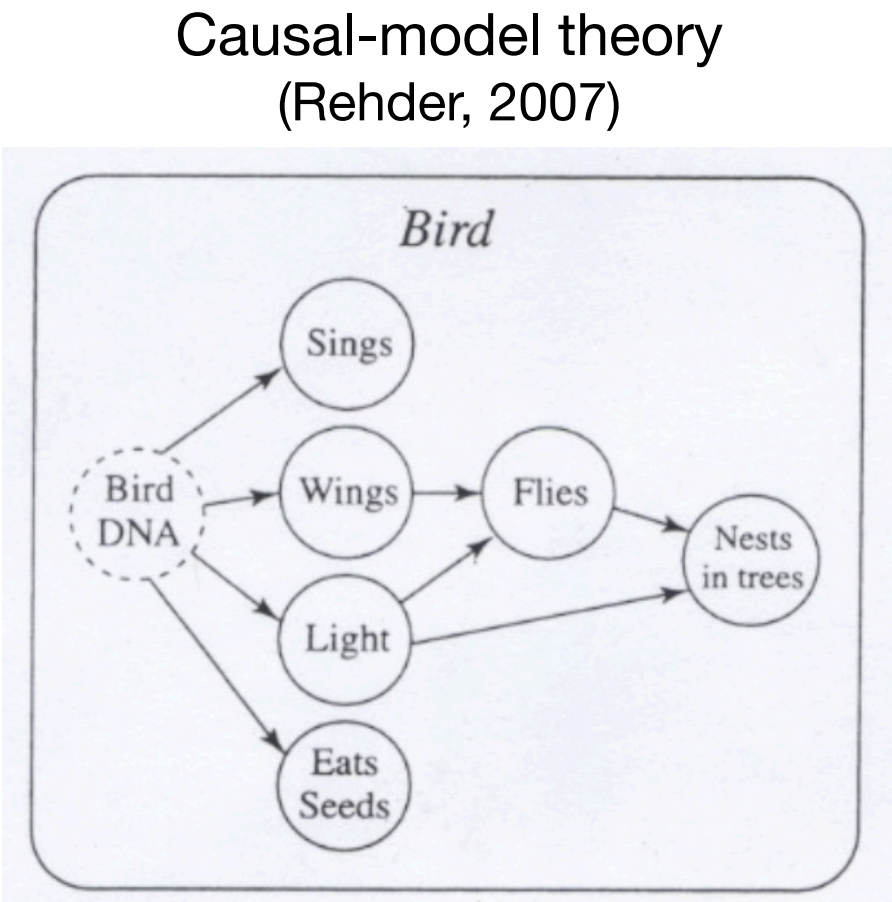


# Modeling Traditions

## Tradition 1: structured knowledge

Intuitive theories  
(Murphy & Medin, 1985)

The "language of thought"  
(Fodor, 1975)

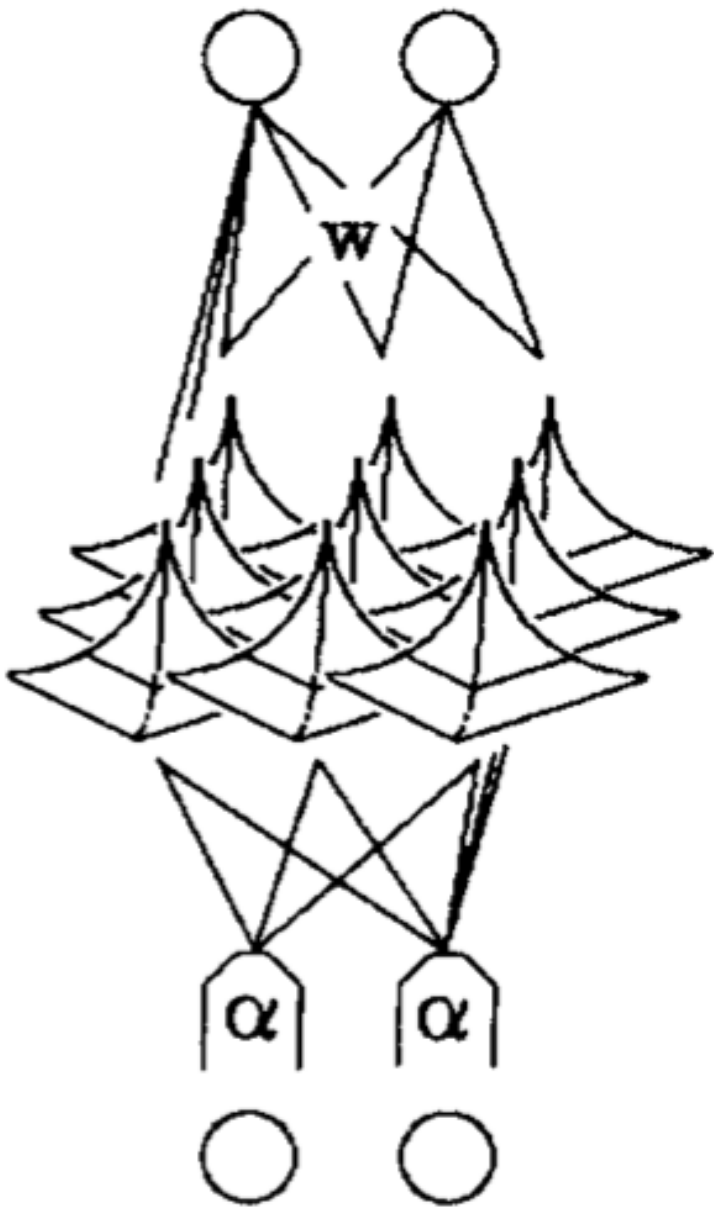


Boolean concepts  
(Feldman, 2000)

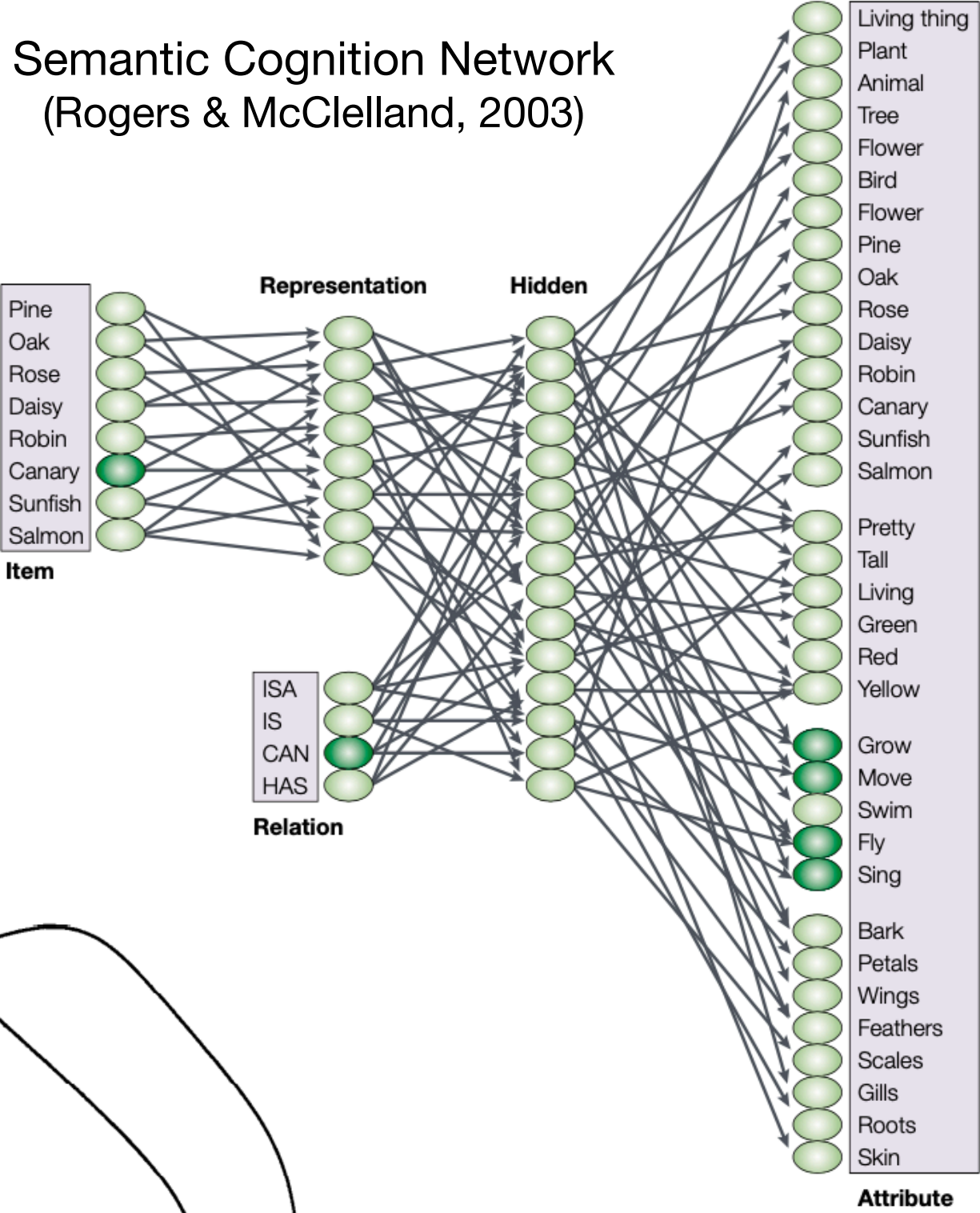
DNF	Minimal formula	Complexity	Illustration
$a'b'c' + a'b'c + a'bc'$	$a'(bc)'$	3	
$a'b'c' + a'b'c + abc'$	$a'b' + abc'$	5	
$a'b'c' + a'bc + ab'c$	$a'(b'c' + bc) + ab'c$	8	

## Tradition 2: statistical knowledge

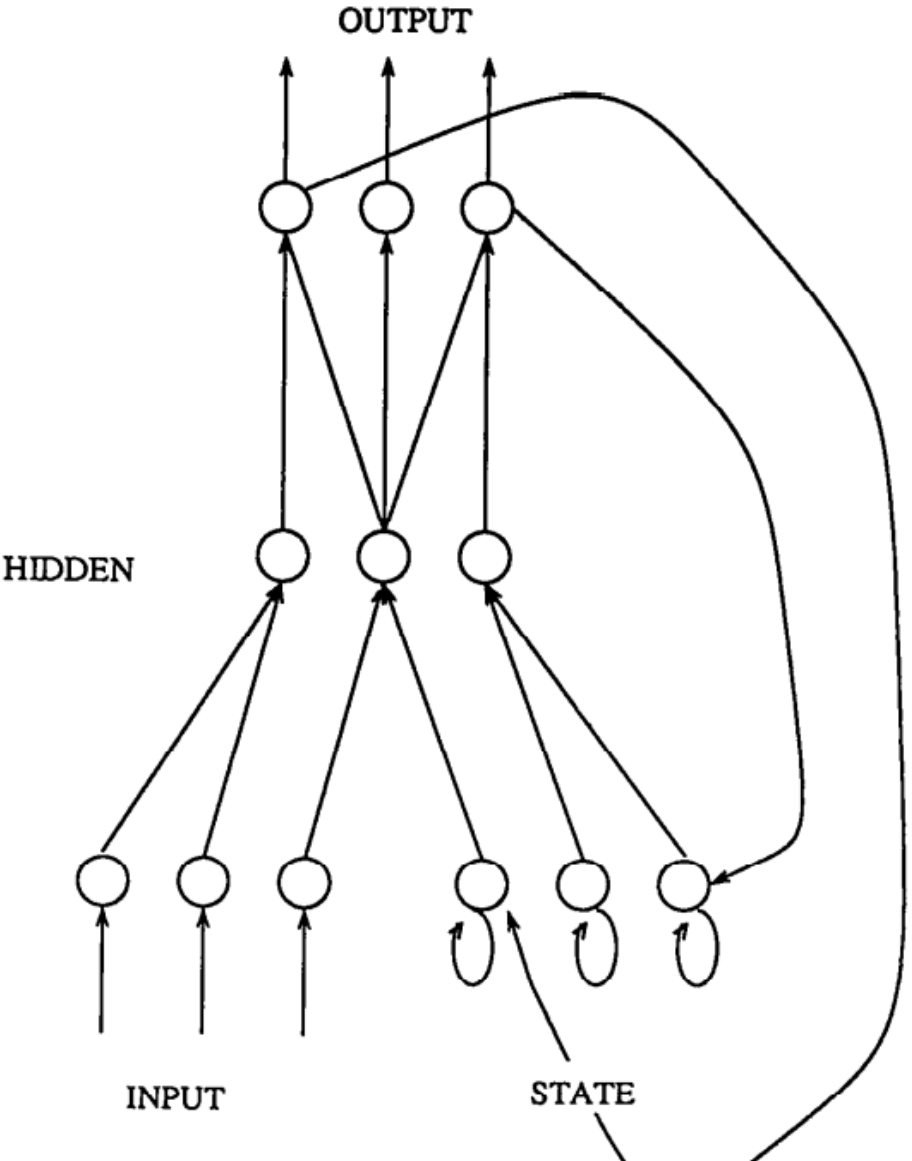
ALCOVE  
(Kruschke, 1992)



Semantic Cognition Network  
(Rogers & McClelland, 2003)



Synthesis?



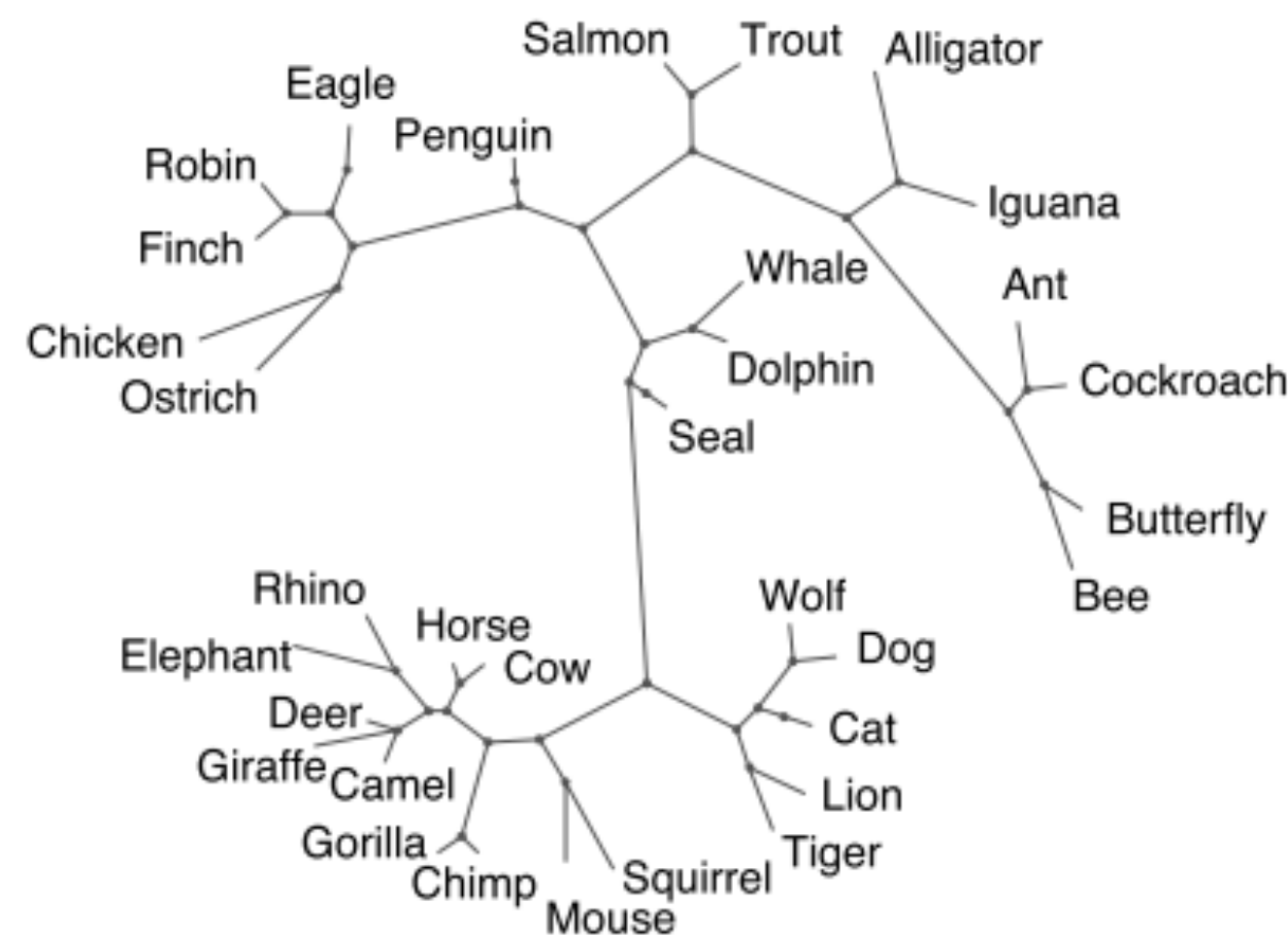
Finding Structure in Time  
(Elman, 1990)



# Prior work: Integrating structure and statistics

Bayes' rule:  $P(\text{structure} \mid \text{data}) \propto P(\text{structure})P(\text{data} \mid \text{structure})$

Structural Forms  
(Kemp & Tenenbaum, 2008)

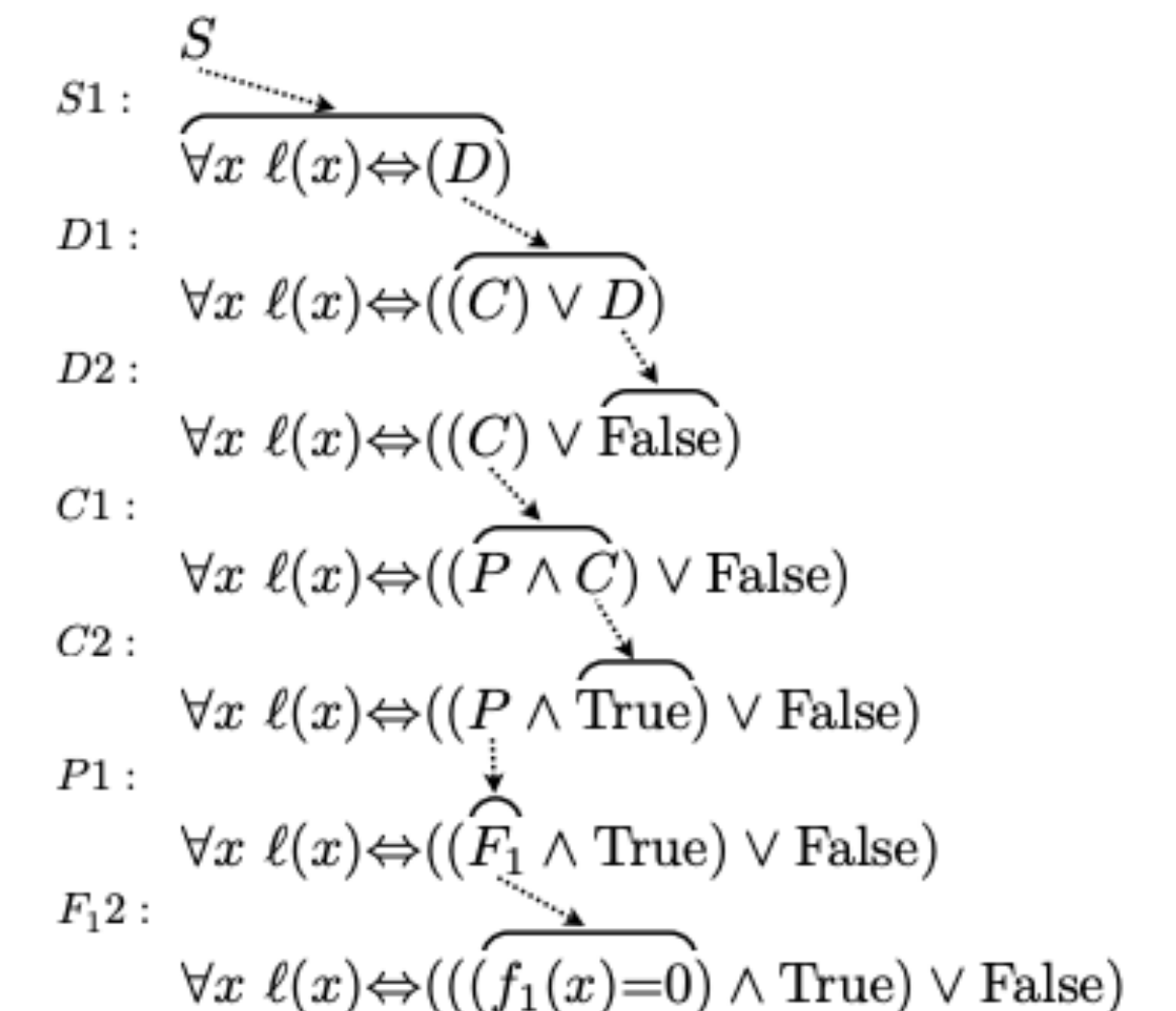


Bayesian Program Learning  
(Lake et al., 2015)

```

procedure GENERATE TYPE
   $\kappa \leftarrow P(\kappa)$  ▷ Sample number of parts
  for  $i = 1 \dots \kappa$  do
     $n_i \leftarrow P(n_i \mid \kappa)$  ▷ Sample number of sub-parts
    for  $j = 1 \dots n_i$  do
       $s_{ij} \leftarrow P(s_{ij} \mid s_{i(j-1)})$  ▷ Sample sub-part sequence
    end for
     $R_i \leftarrow P(R_i \mid S_1, \dots, S_{i-1})$  ▷ Sample relation
  end for
   $\psi \leftarrow \{\kappa, R, S\}$ 
  return @GENERATE TOKEN( $\psi$ ) ▷ Return program
  
```

Rational Rules  
(Goodman et al., 2008)





# Proposal: Generative Neuro-Symbolic (GNS) modeling

---

**procedure** GENERATEEXAMPLE

---

$C \leftarrow 0$

▷ Initialize blank canvas

**for**  $i = 1 \dots, \infty$  **do**

$x_i \leftarrow \text{GENERATEPART}(C)$

▷ Sample part

$r_i \leftarrow \text{GENERATERELATION}(C, x_i)$

▷ Sample relation

$C \leftarrow \text{RENDER}(C, x_i, r_i)$

▷ Render new canvas

**if**  $\text{TERMINATE?}(C)$  **then**

▷ Sample termination (y/n)

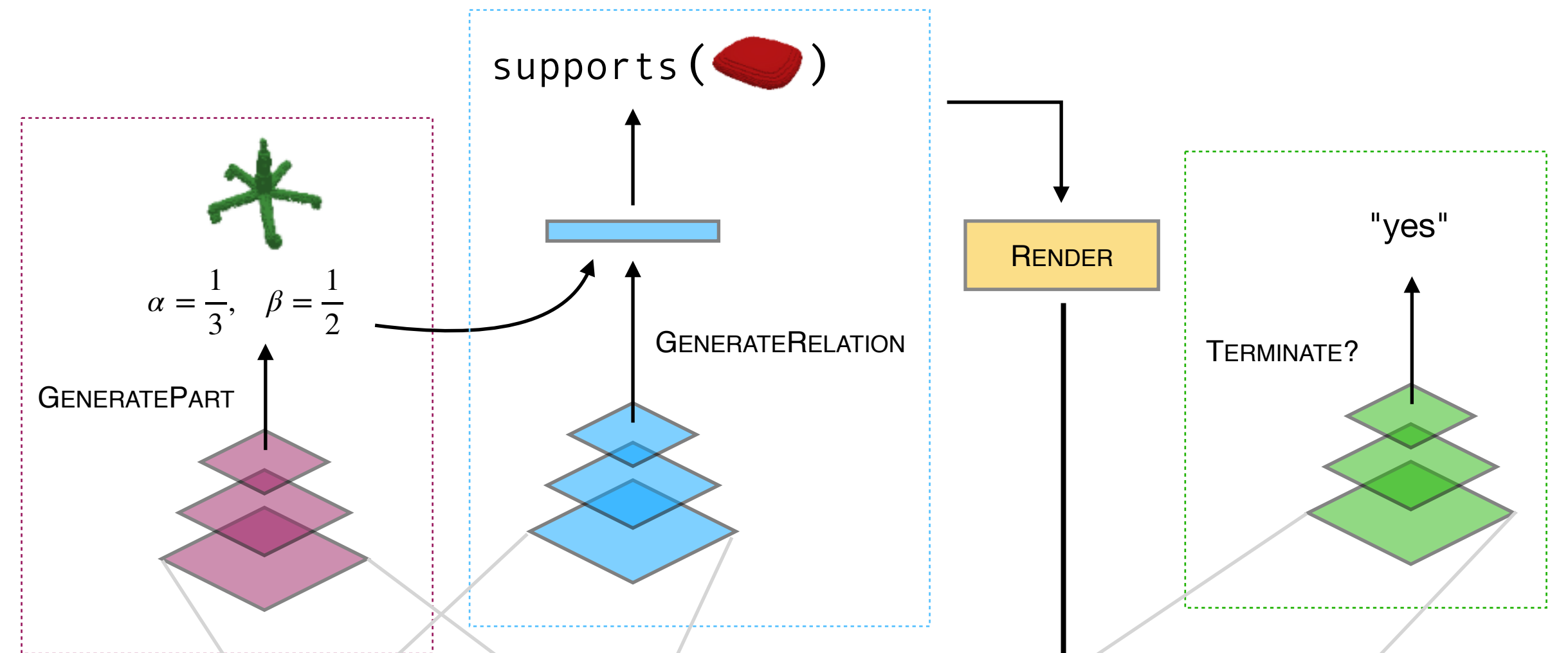
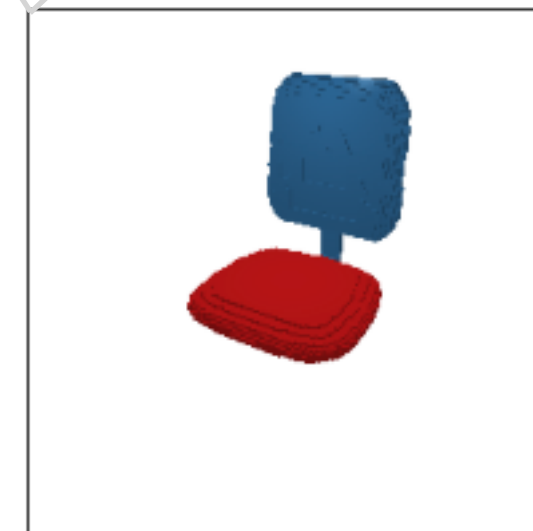
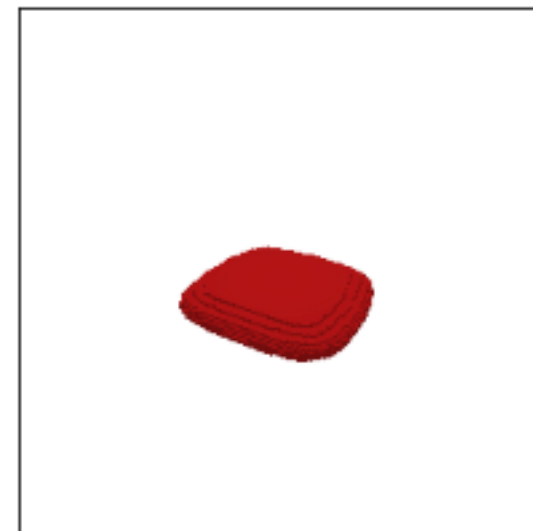
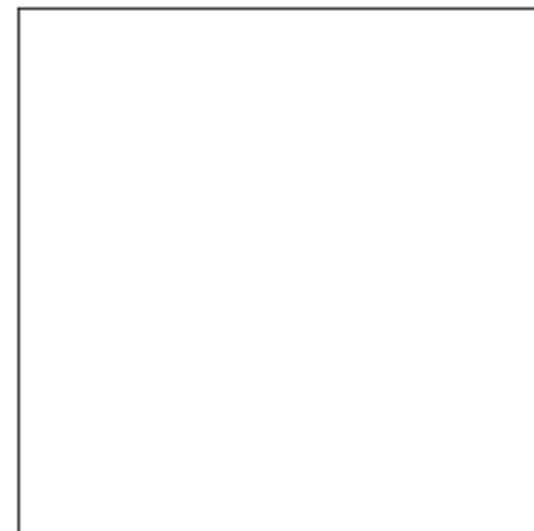
break

**return**  $C$

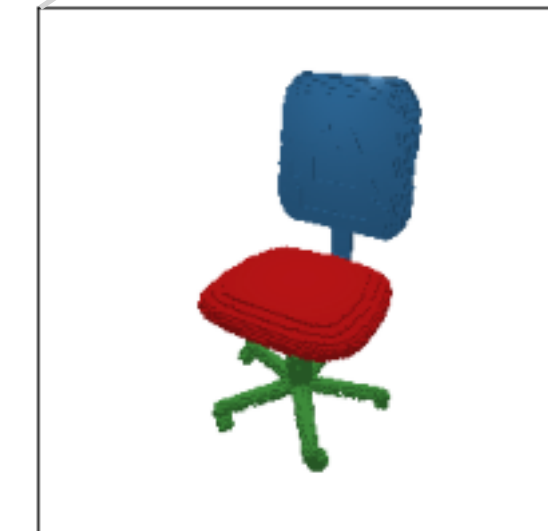
▷ Return example

---

Canvas:  
 $C$



RENDER



New example

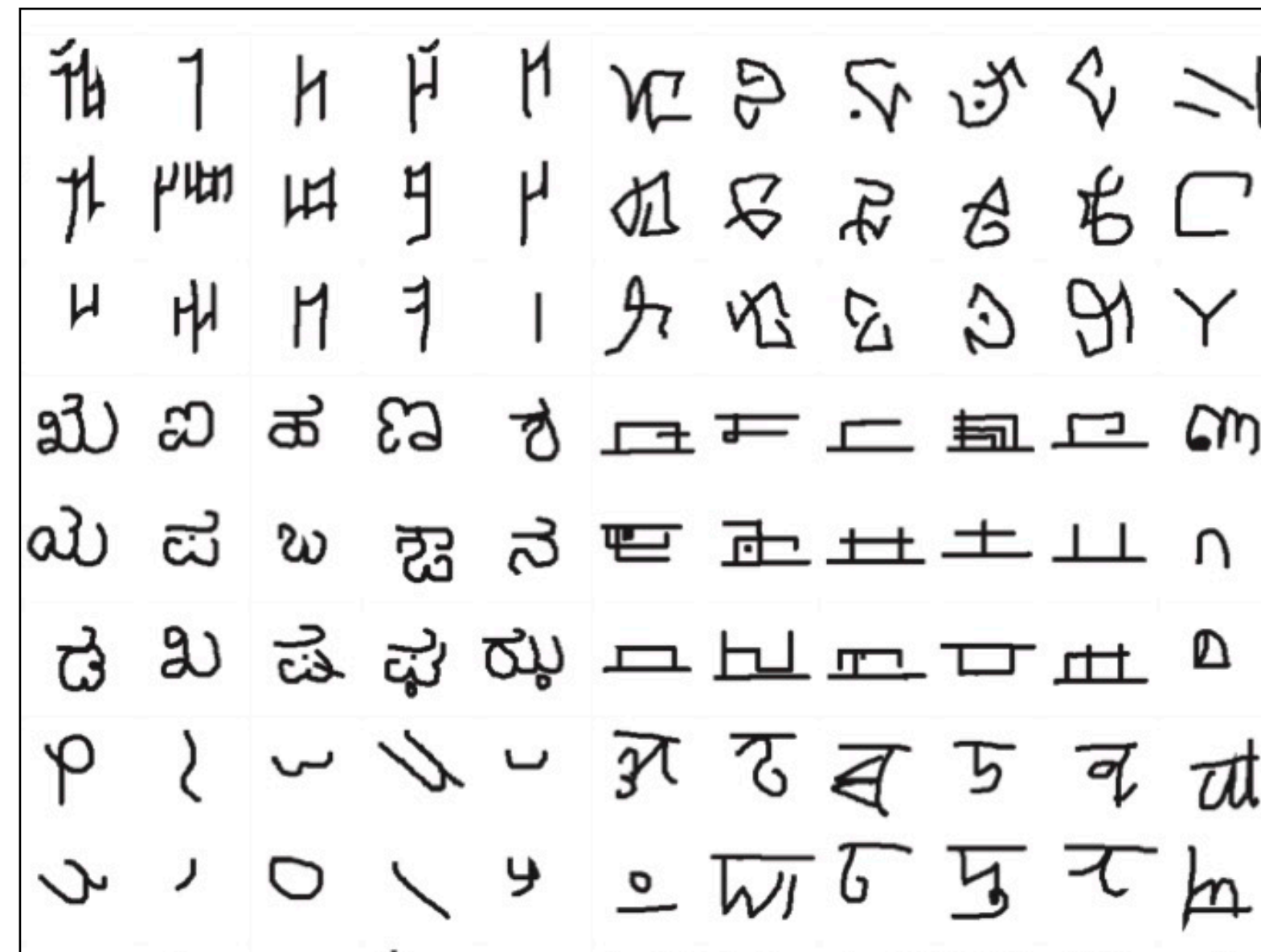


# Agenda

- Case study #1: handwritten characters
- Case study #2: structured visual concepts ("alien figures")
- Additional projects
- Summary & conclusions



# Case study #1: handwritten characters



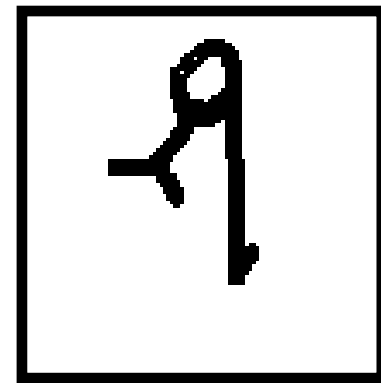
(Lake et al., 2015)



# The Omniglot Challenge

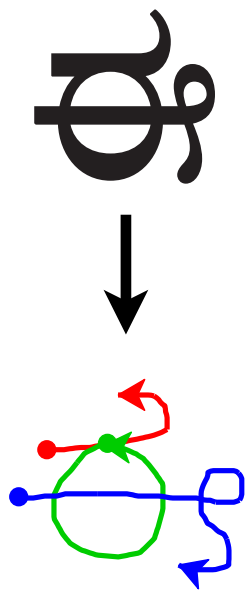
(Lake et al., 2015)

classification



ॐ	ॐ	ॐ	ॐ	ॐ
क	६	१	३	५
८	५	५	८	५
५	५	५	८	५

parsing

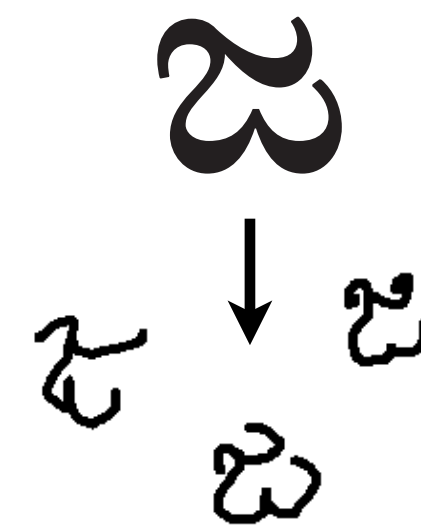


generating  
new concepts

ॐ	ॐ	ॐ	ॐ	ॐ
ॐ	ॐ	ॐ	ॐ	ॐ

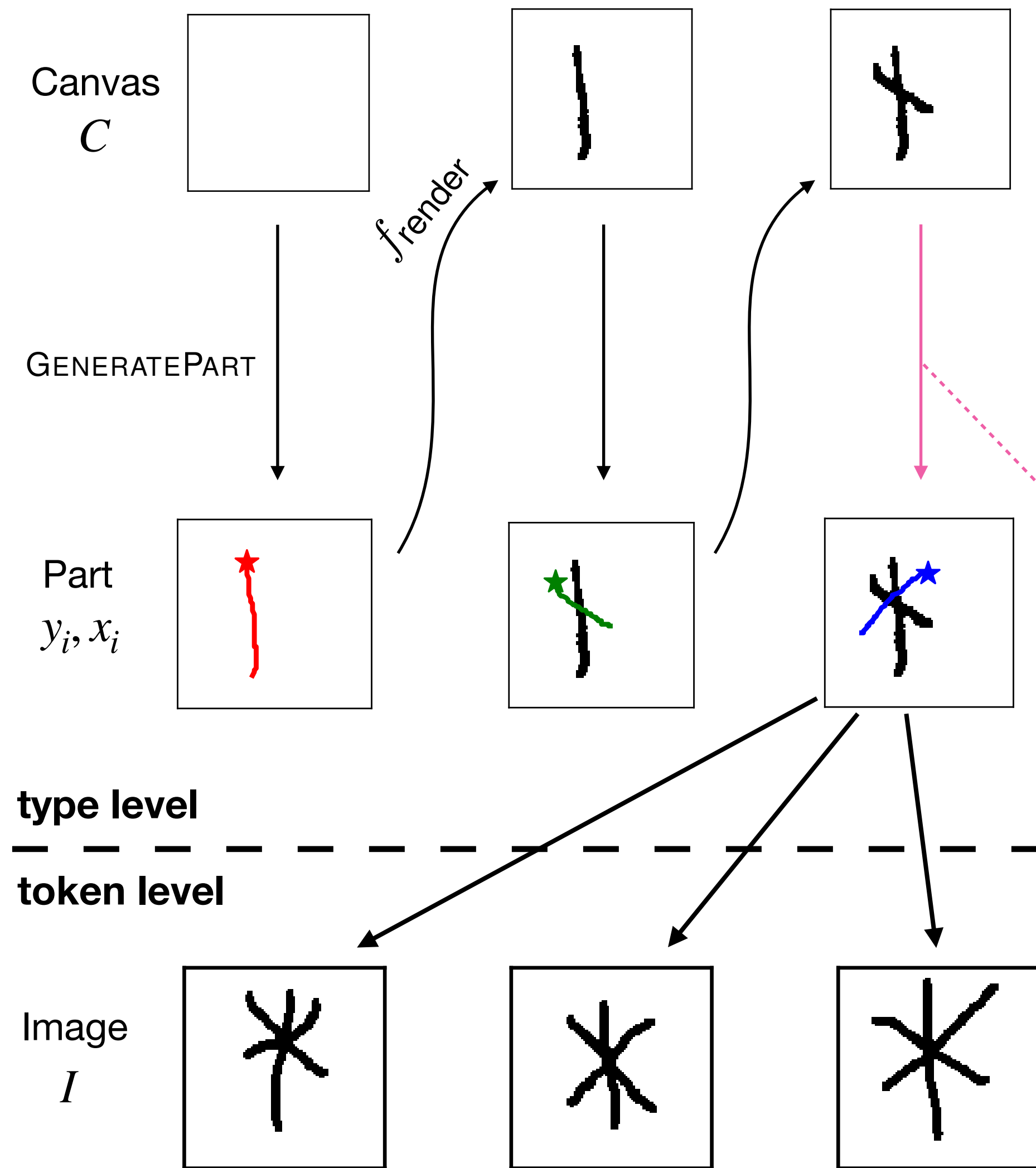
ॐ	ॐ
ॐ	ॐ

generating  
new examples





# GNS model of character concepts



**procedure** GENERATE TYPE

$C \leftarrow 0$

▷ Initialize blank image canvas

**while** *true* **do**

$[y_i, x_i] \leftarrow \text{GENERATEPART}(C)$

▷ Sample part location & parameters

$C \leftarrow f_{\text{render}}(y_i, x_i, C)$

▷ Render part to image canvas

$v_i \sim p(v \mid C)$

▷ Sample termination indicator

**if**  $v_i$  **then**

**break**

▷ Terminate sample

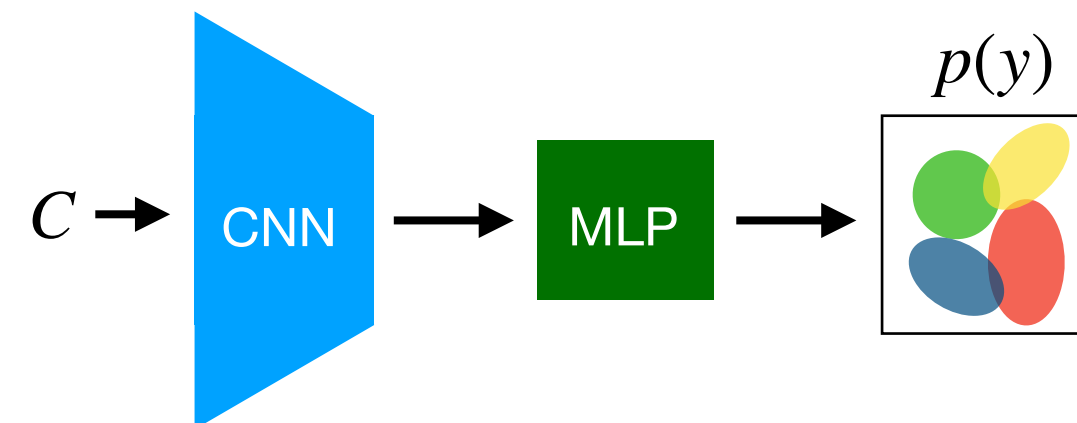
$\psi \leftarrow \{\kappa, y_{1:\kappa}, x_{1:\kappa}\}$

**return**  $\psi$

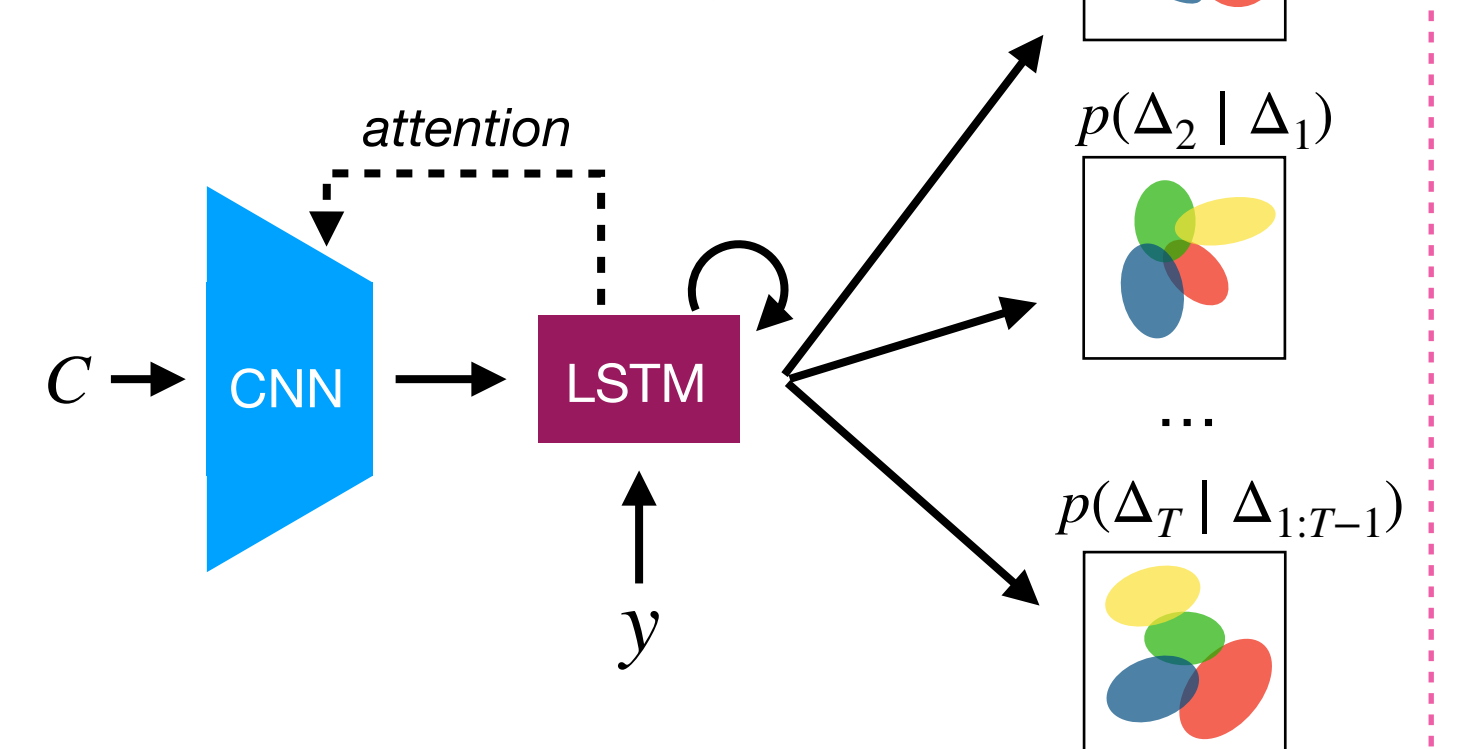
▷ Return concept type

GENERATEPART( $C$ )

location model  $p(y \mid C)$



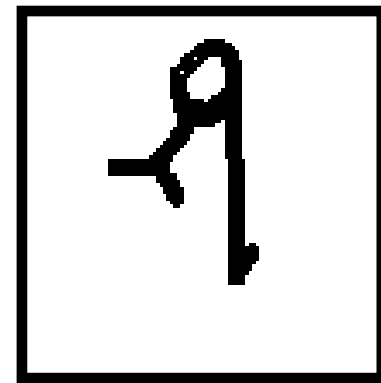
stroke model  $p(x \mid y, C)$





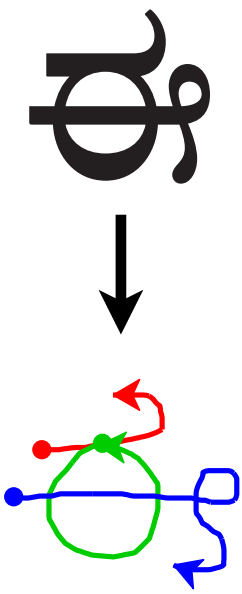
# The Omniglot Challenge

classification

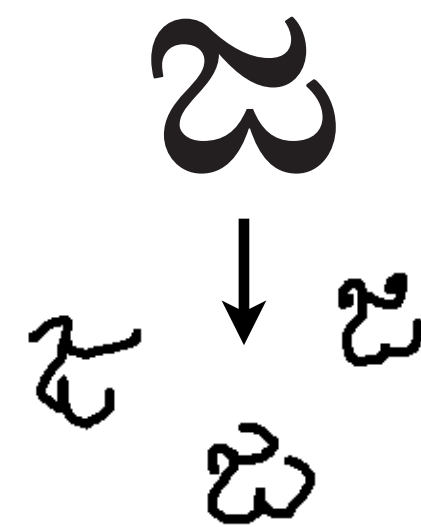


ग	ॠ	ॡ	ॢ	ॣ
क	ख	ग	घ	ङ
च	छ	ज	झ	ञ
ट	ठ	ड	ढ	ण

parsing

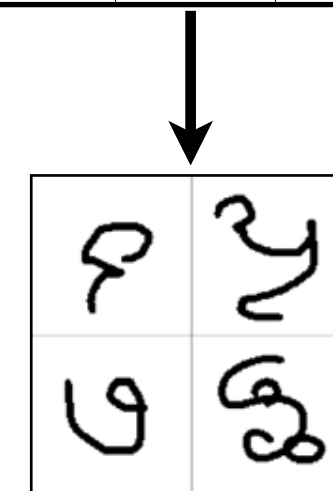


generating  
new examples



generating  
new concepts

ॐ	ॐ	ॐ	ॐ	ॐ
ॐ	ॐ	ॐ	ॐ	ॐ





# Generating new concepts

## 1. Log-likelihoods (LL) of held-out concepts

Test loss per drawing trajectory

GNS	-19.51
H-LSTM	-20.16
LSTM	-19.66

Replicates across different train/test splits

Approximate test LL per pixel image

GNS	-383.67
VHE	-546.84
SG	-861.05

## 2. Model samples

Omniglot

𐄂	𐄃	𐄄	𐄅	𐄆	𐄇
𐄈	𐄉	𐄊	𐄋	𐄌	𐄍
𐄎	𐄏	𐄐	𐄑	𐄒	𐄓
𐄔	𐄕	𐄖	𐄗	𐄘	𐄙
𐄚	𐄛	𐄜	𐄝	𐄞	𐄟
𐄠	𐄡	𐄢	𐄣	𐄤	𐄥

GNS model

𐄂	𐄃	𐄄	𐄅	𐄆	𐄇
𐄈	𐄉	𐄊	𐄋	𐄌	𐄍
𐄎	𐄏	𐄐	𐄑	𐄒	𐄓
𐄔	𐄕	𐄖	𐄗	𐄘	𐄙
𐄚	𐄛	𐄜	𐄝	𐄞	𐄟
𐄠	𐄡	𐄢	𐄣	𐄤	𐄥

fully-symbolic model (BPL)

𐄂	𐄃	𐄄	𐄅	𐄆	𐄇
𐄈	𐄉	𐄊	𐄋	𐄌	𐄍
𐄎	𐄏	𐄐	𐄑	𐄒	𐄓
𐄔	𐄕	𐄖	𐄗	𐄘	𐄙
𐄚	𐄛	𐄜	𐄝	𐄞	𐄟
𐄠	𐄡	𐄢	𐄣	𐄤	𐄥



# Generating new concepts

GNS

Nearest neighbors



GNS samples



Nearest neighbors

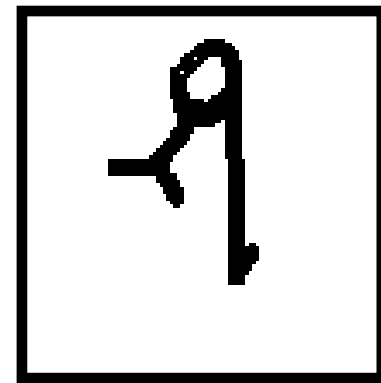


Nearest neighbors are located using the embedding of a convolutional neural network (CNN)



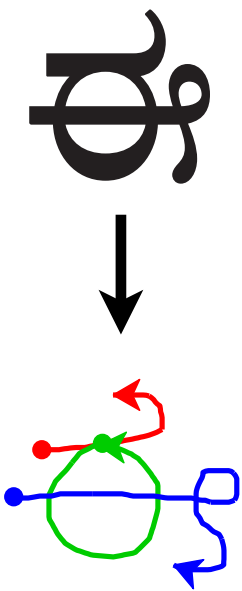
# The Omniglot Challenge

classification



ग	ॠ	ॡ	ॢ	ॣ
क	ख	ग	घ	ङ
च	छ	ज	झ	ञ
ट	ठ	ड	ढ	ण

parsing



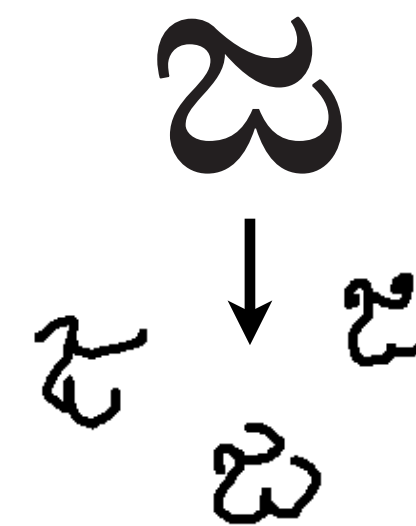
generating  
new concepts

ॐ	ॐ	ॐ	ॐ	ॐ
ॐ	ॐ	ॐ	ॐ	ॐ



ॐ	ॐ
ॐ	ॐ

generating  
new examples



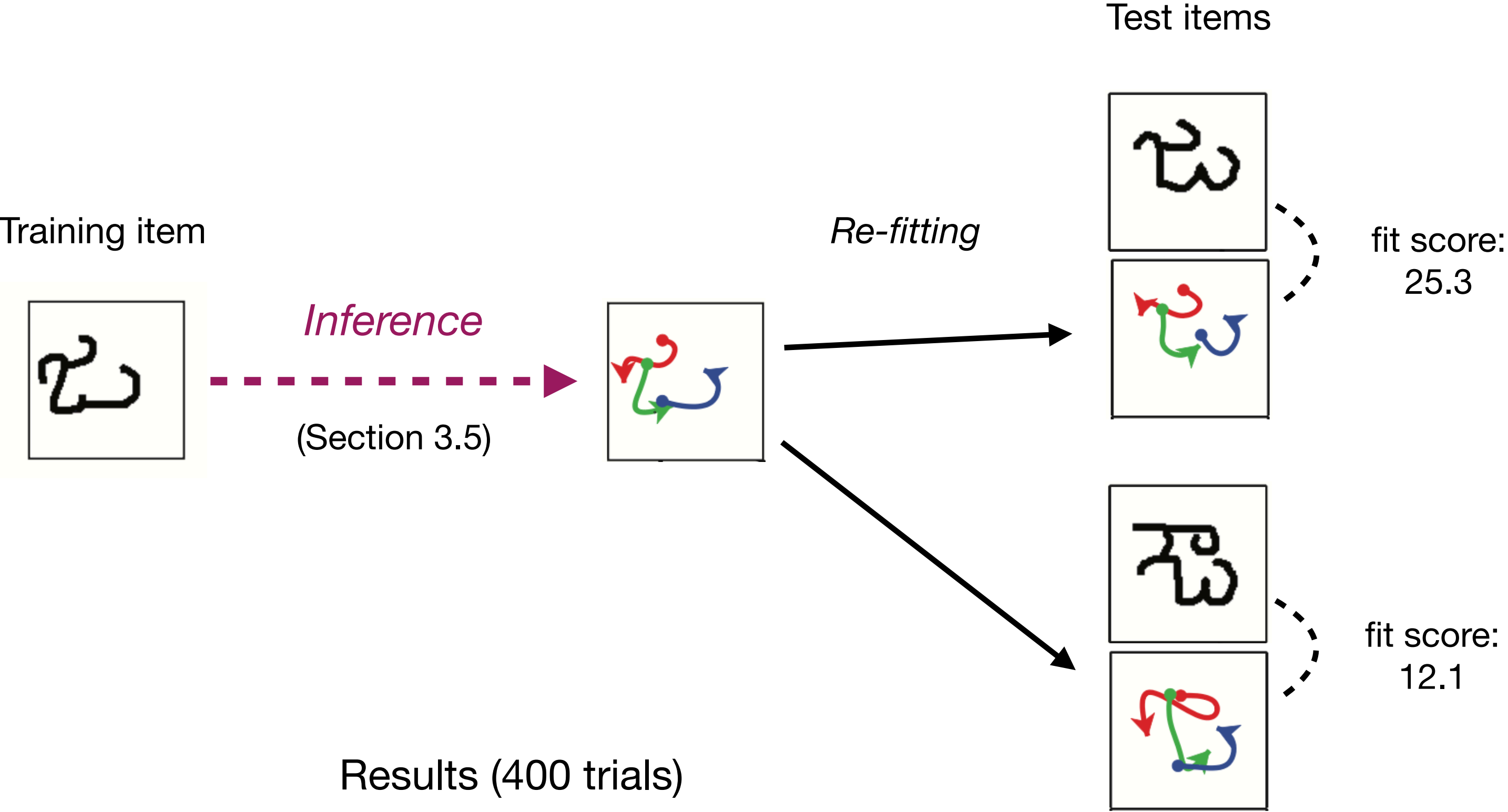


# One-Shot Classification

Target

where is another?

ಅ	ಇ	ಉ	ಎ	ಐ
ಈ	ಖ	ಗ	ಒ	ಝ
ಞ	ತ	ಣ	ತೆ	ದ
ನ	ಯ	ಲ	ಹ	ಳ



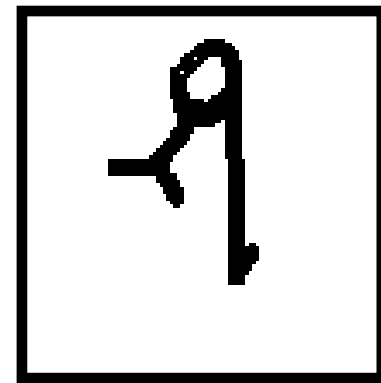
Results (400 trials)

	Accuracy
Humans	95.5%
GNS	94.3%
BPL	96.7%



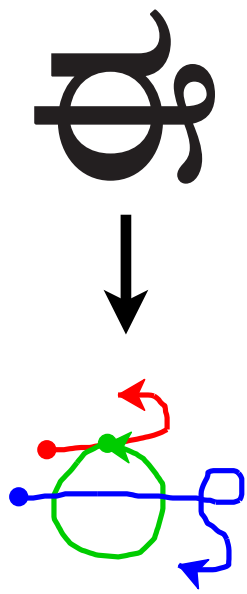
# The Omniglot Challenge

classification



ग	प	म	न	र
क	ख	ग	घ	ङ
च	छ	ज	झ	ञ
ट	ठ	ड	ढ	ण

parsing



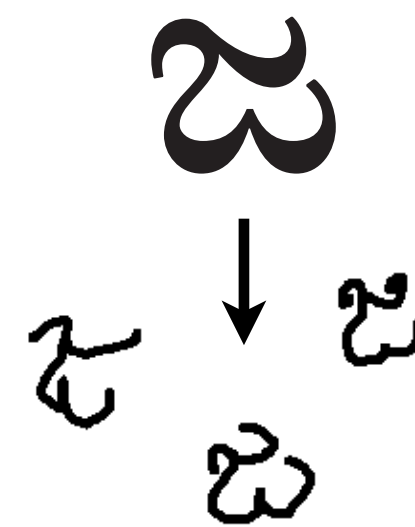
generating  
new concepts

य	३	द	२	७
७	४	३	४	७



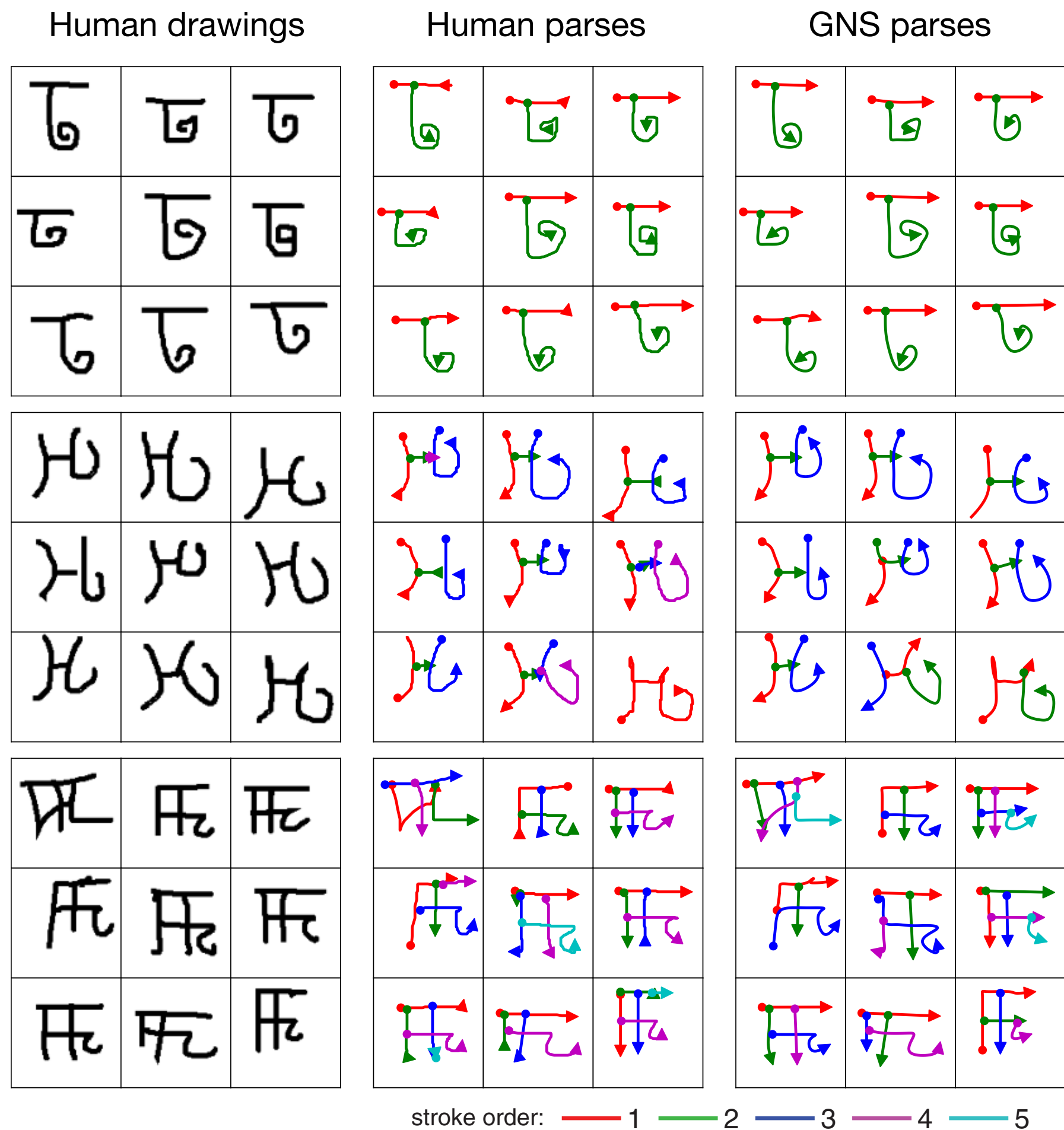
३	२
७	४

generating  
new examples

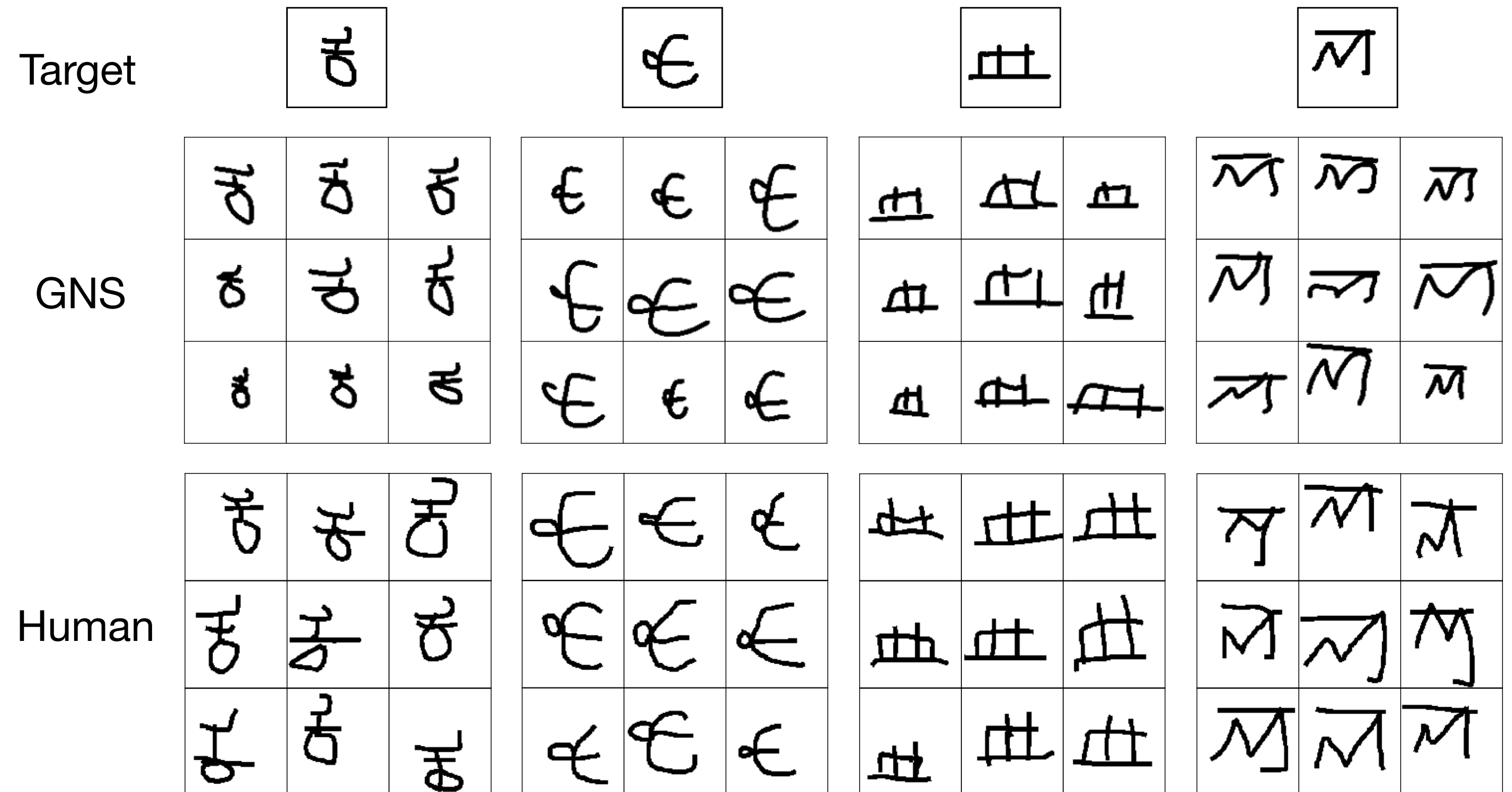




# Parsing



# Generating new exemplars



Many more examples are provided in Section 3.6 and Appendix A.5



# Conclusions: Case study #1

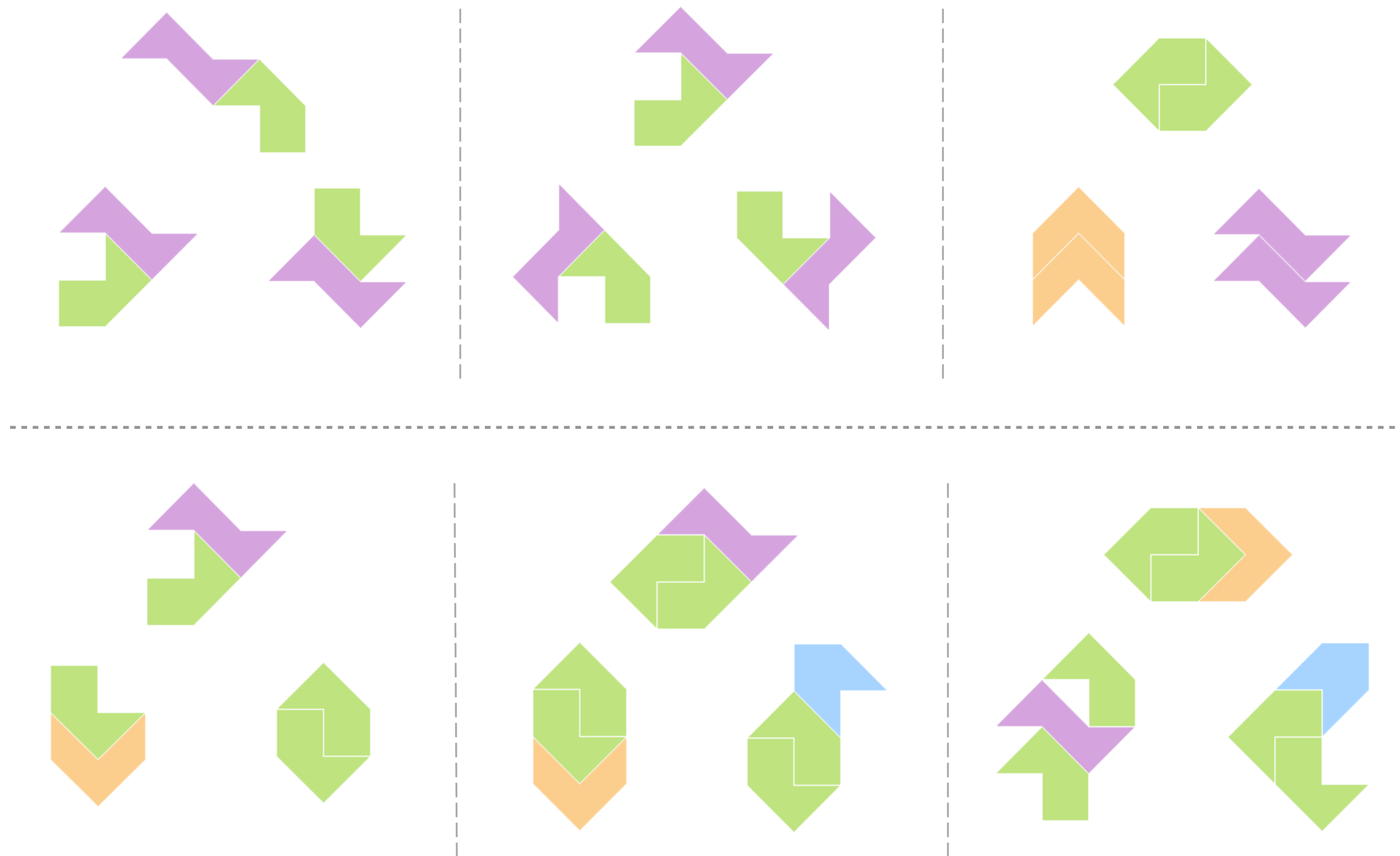
- Humans quickly grasp new concepts and use them in a variety of ways
- Generative Neuro-Symbolic (GNS) models capture the dual structural and statistical components of character concepts and generalize to novel alphabets and a range of tasks
- GNS models offer an account for how previous experience can support the rapid acquisition of new concepts via priors



**Case study #2:  
structured visual concepts  
("alien figures")**



# Alien figures



Yanli Zhou



# Human experiments



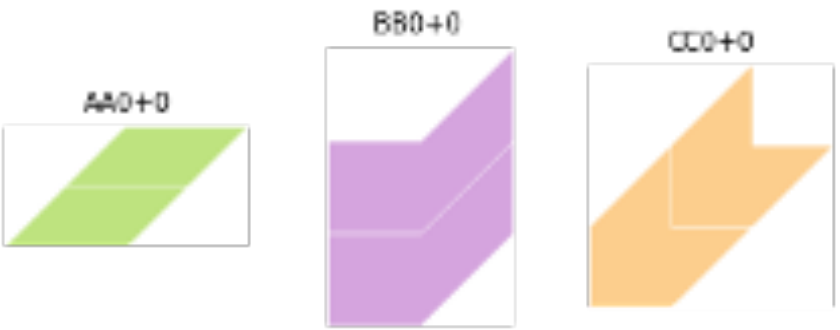
## Categorization

## Generation

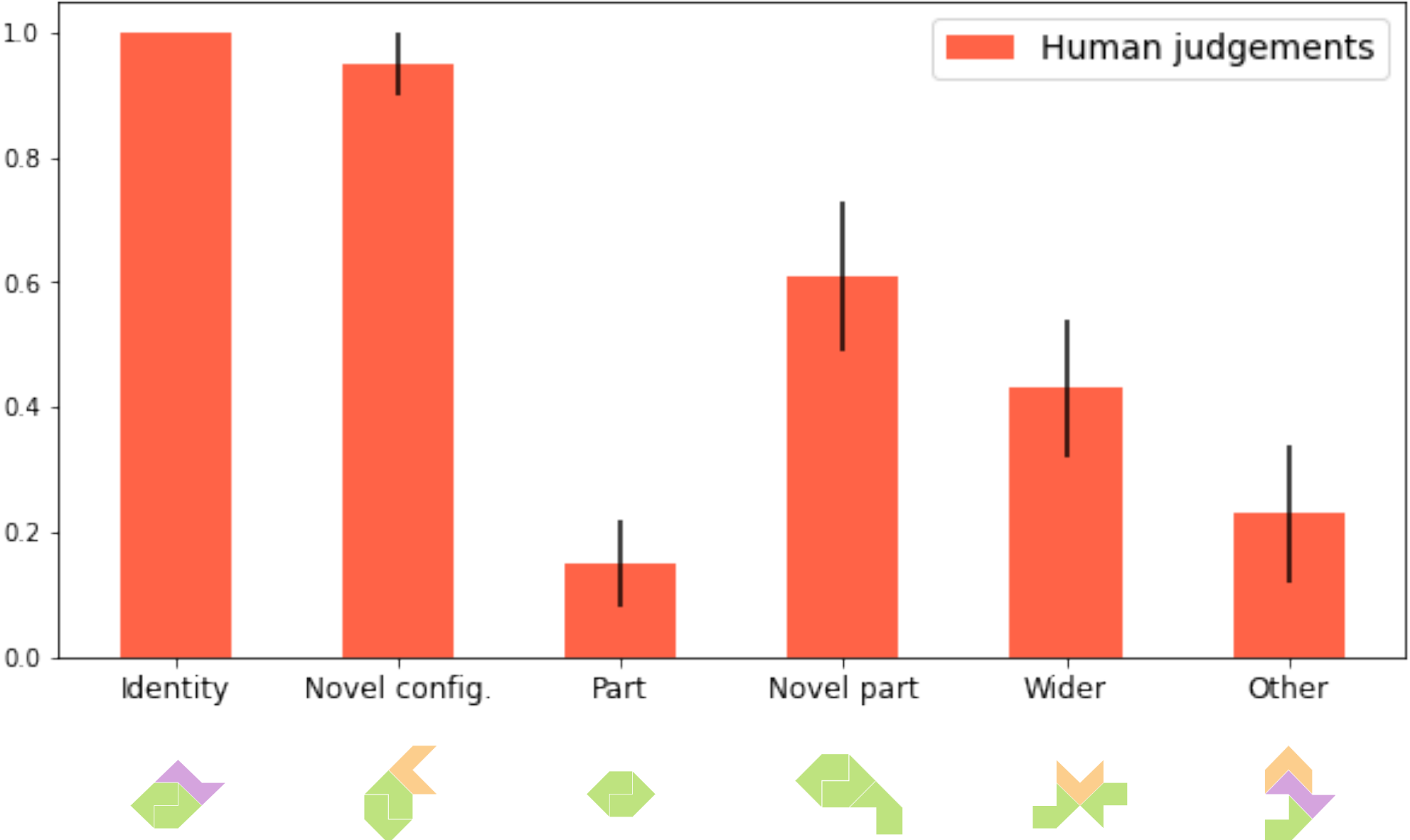
Here are 3 examples of a "wif":



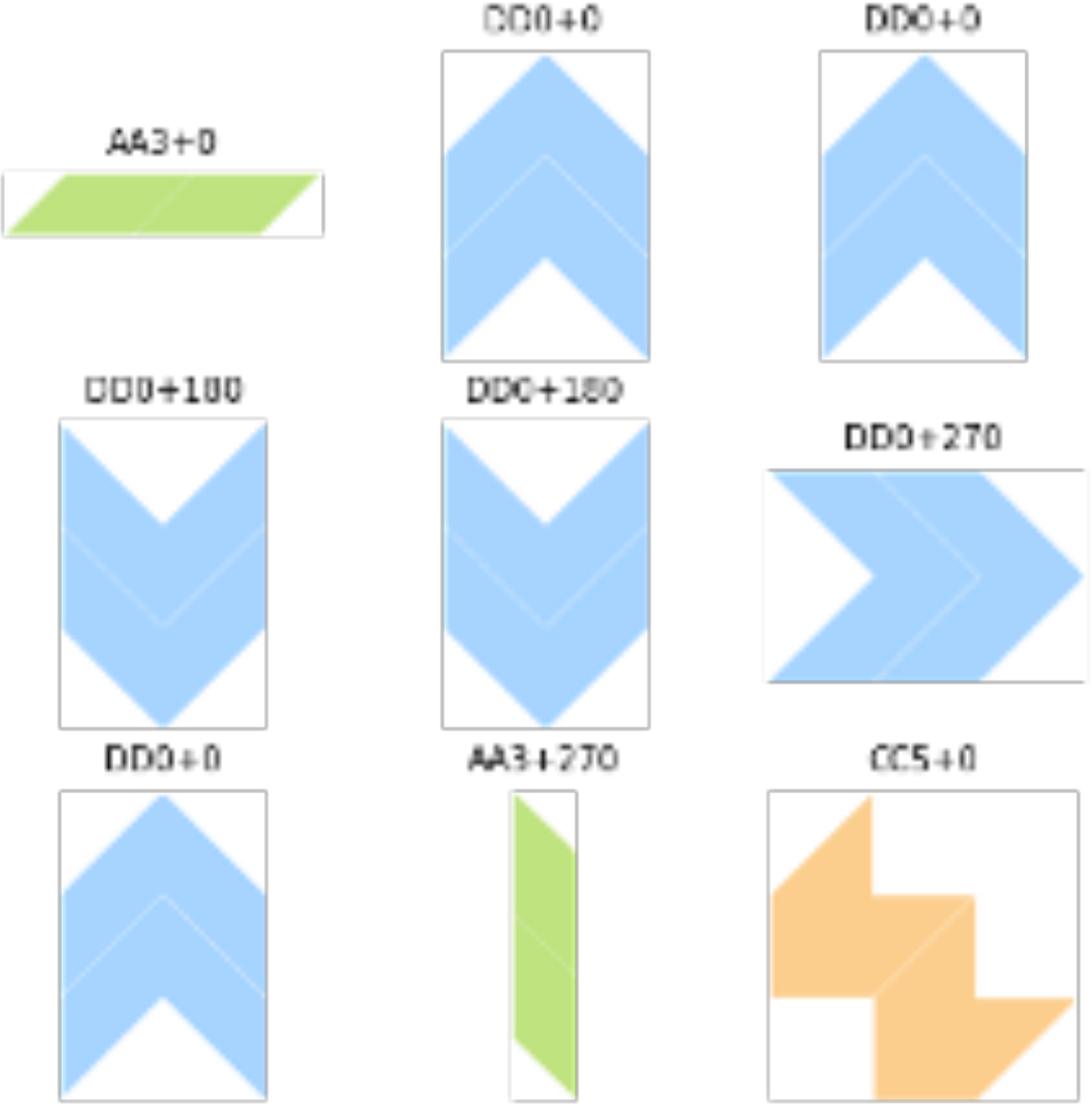
Here are 3 examples of a "dax":



Is this also a "wif"?



Can you make another "dax"?



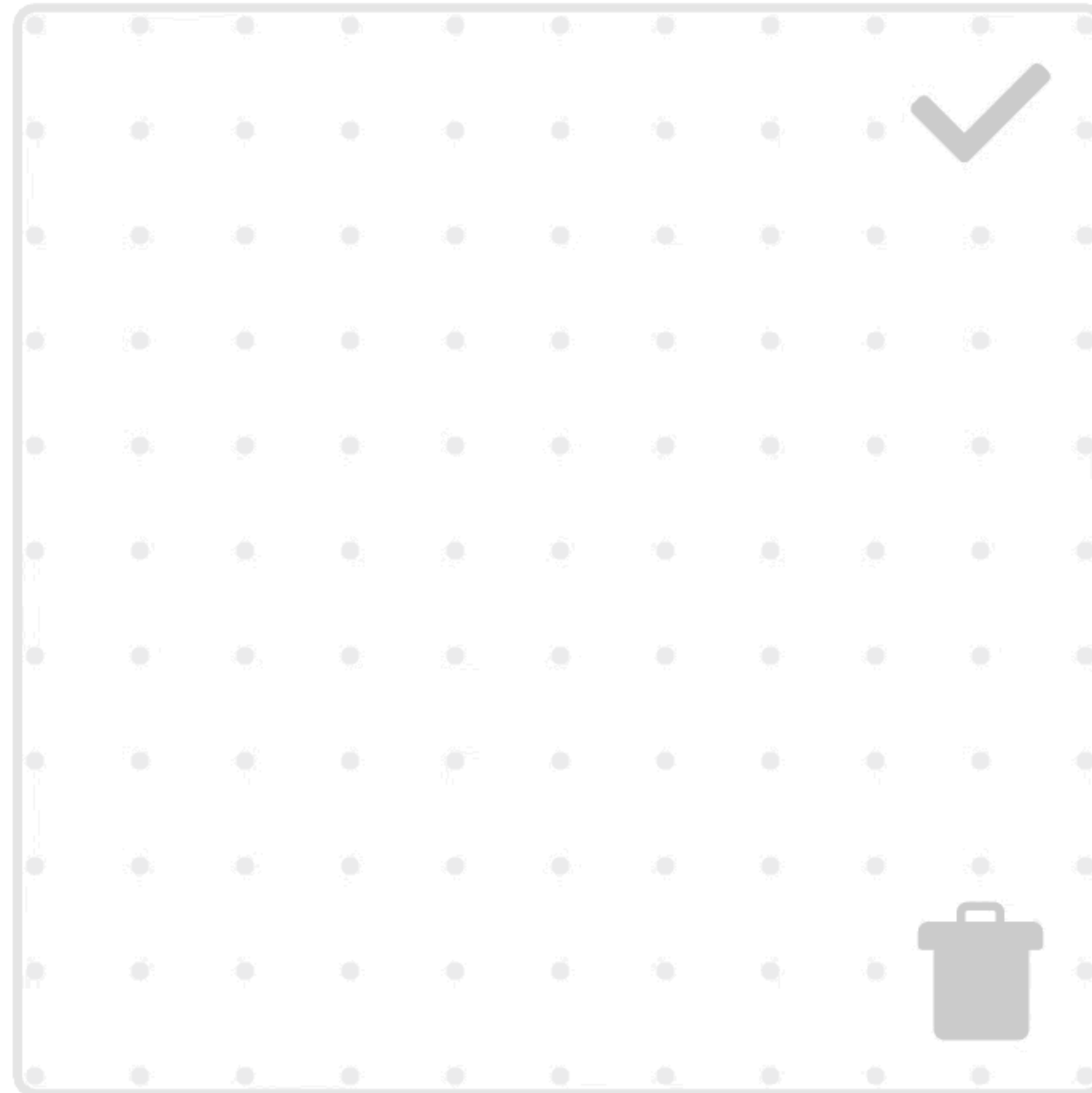
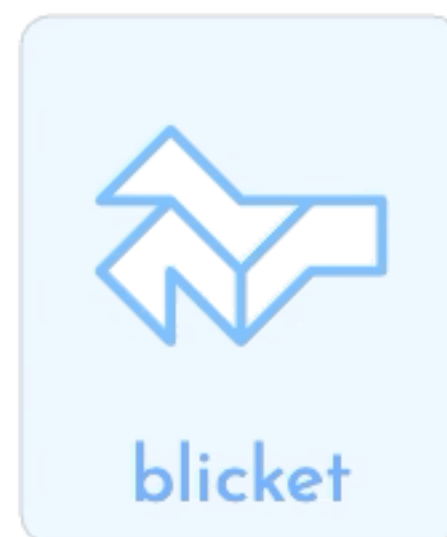
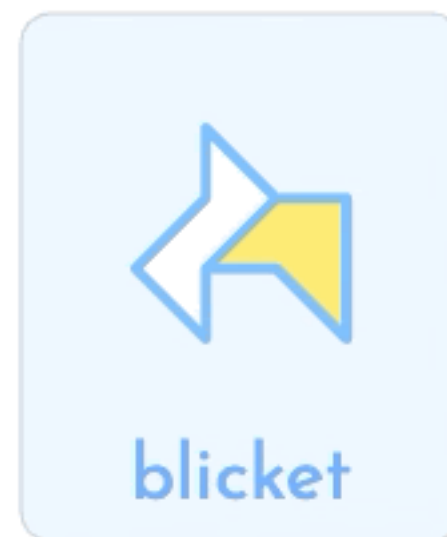
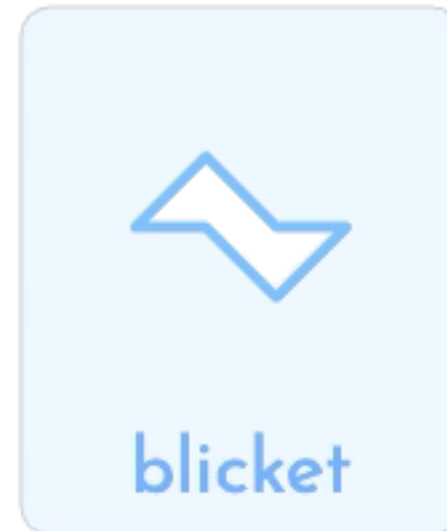
Human generations



# Generation task MTurk interface



Here are 3 examples:



Trial 1/10



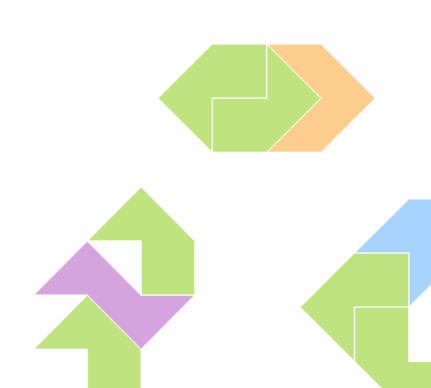
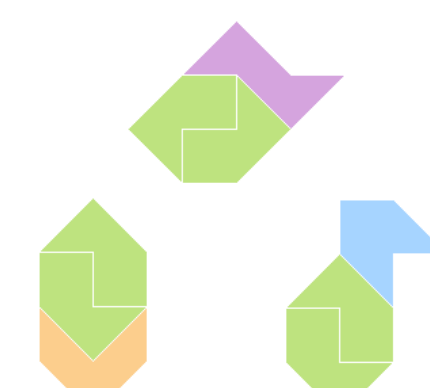
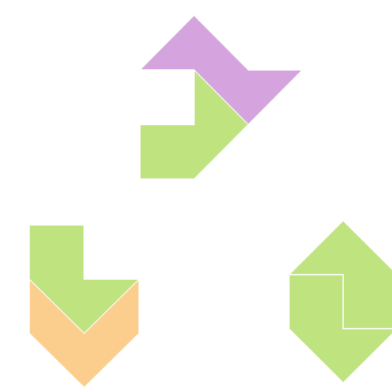
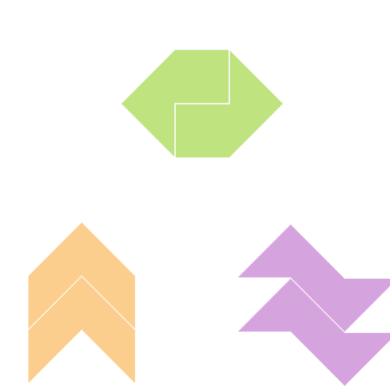
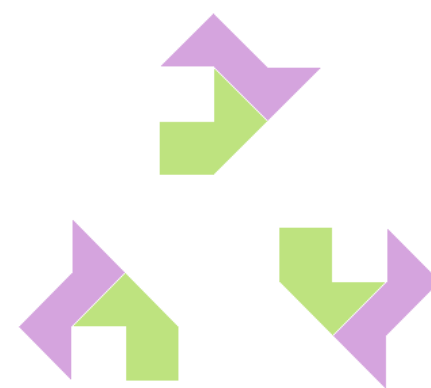
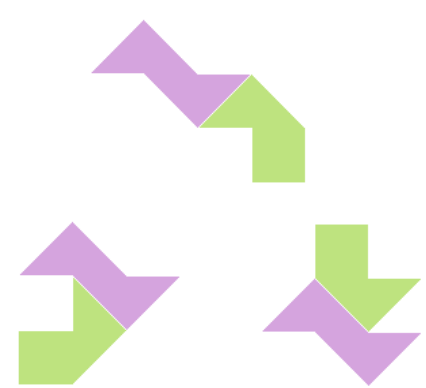
# Symbolic Bayesian model



$$p(h | X) \propto p(h)p(X | h)$$

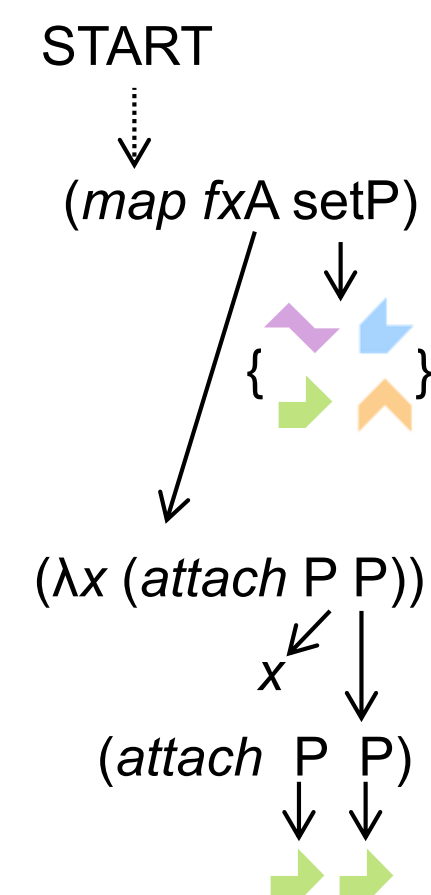
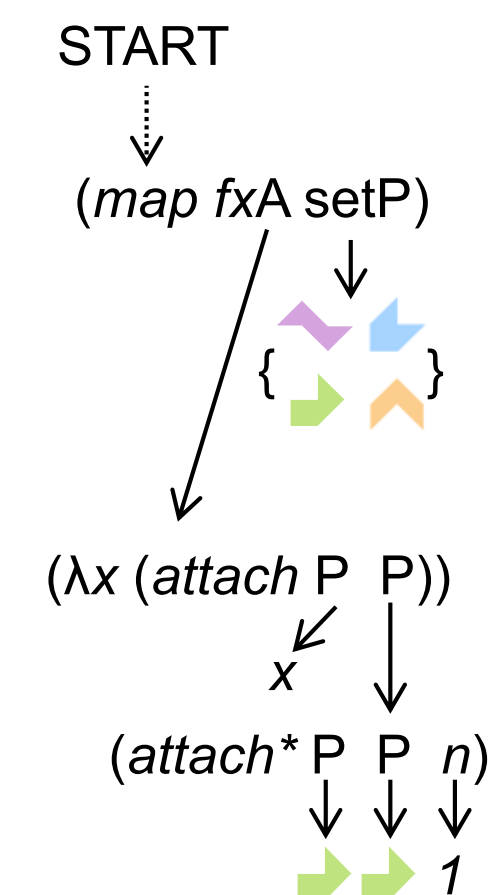
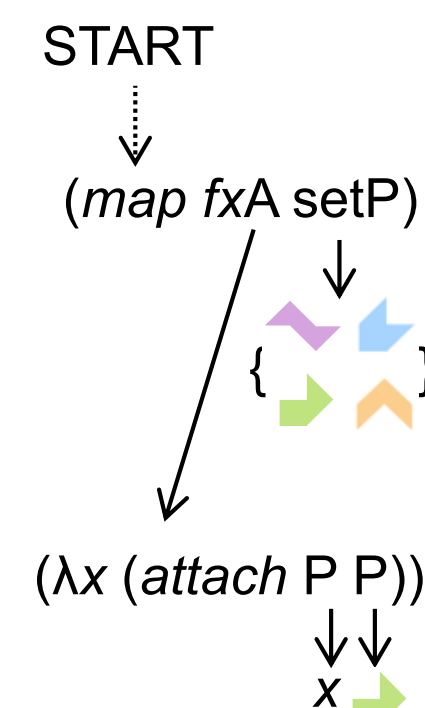
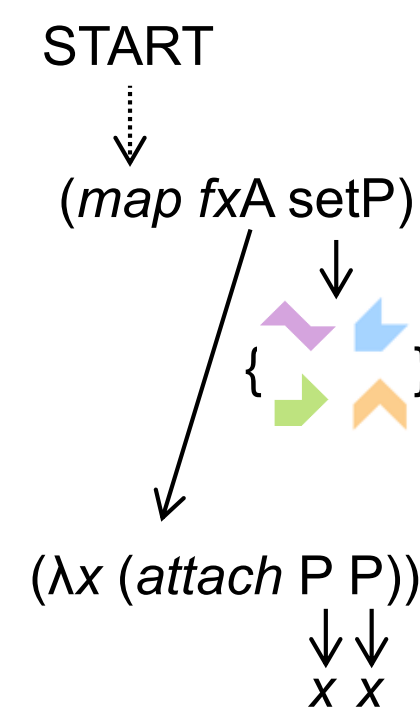
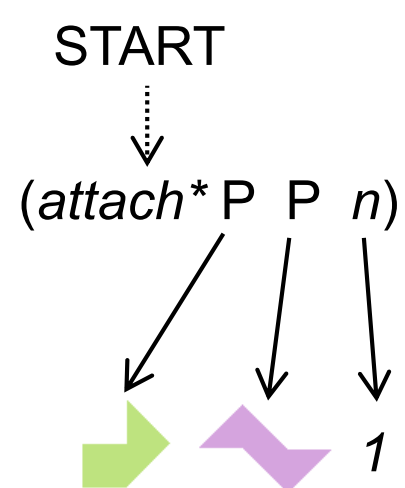
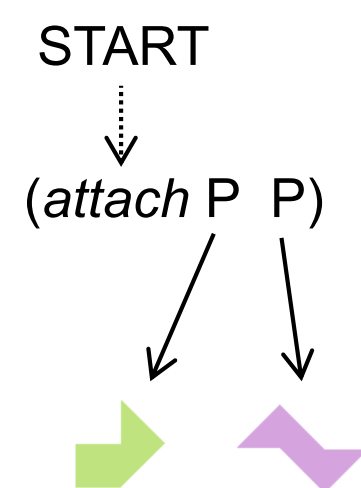
Observation  
("support set")

$X$



Most probable formula  
hypothesis inferred by  
the model

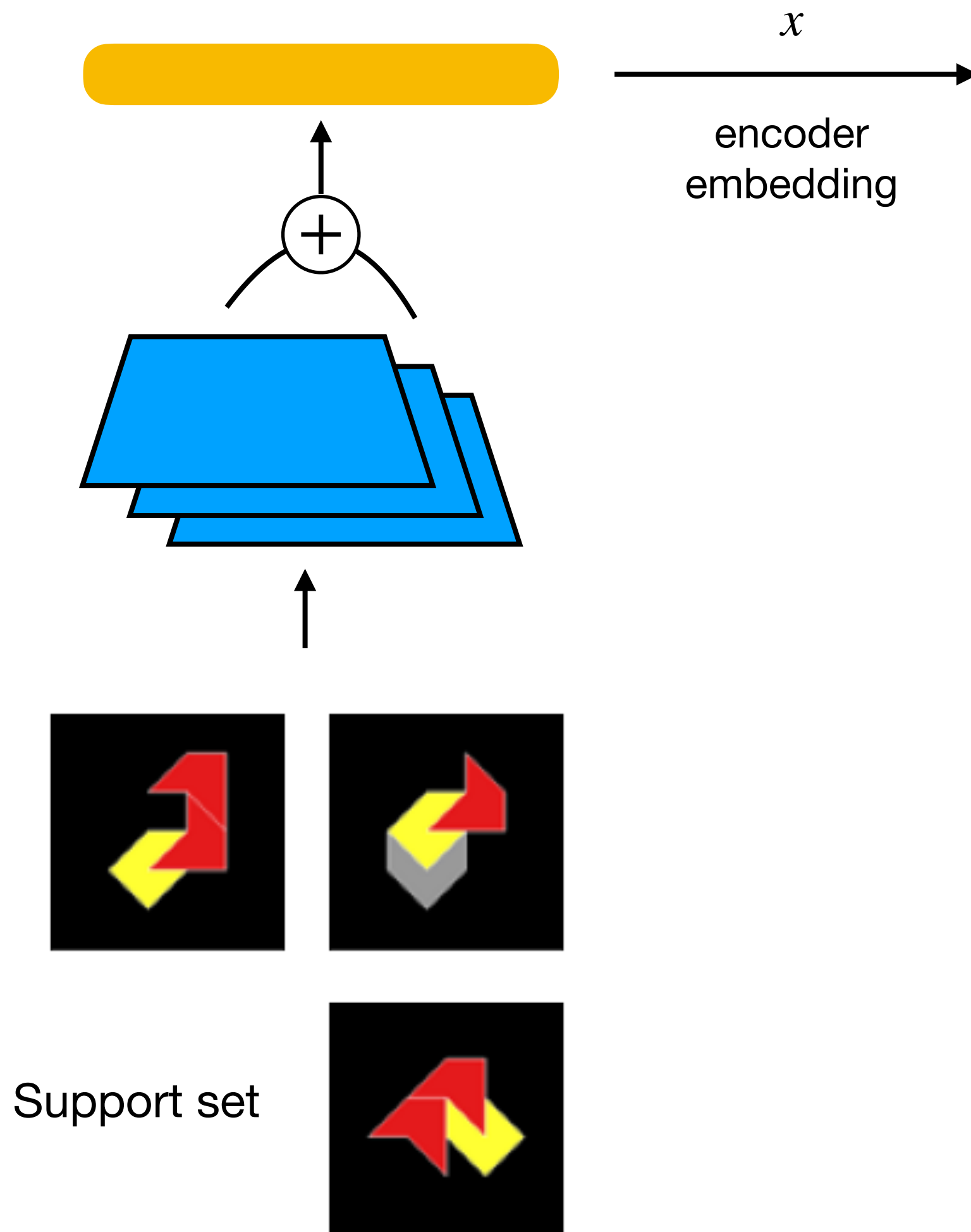
$h$





# Generative neuro-symbolic (GNS) model

## 1. Neural encoder



## 2. Generative neuro-symbolic decoder

**procedure** GENERATE\_TOKEN( $x$ )

$C \leftarrow 0$

**while** True **do**

$c_i \leftarrow \text{GENERATE\_PART}(x, C)$

$r_i \leftarrow \text{GENERATE\_RELATION}(x, C, c_i)$

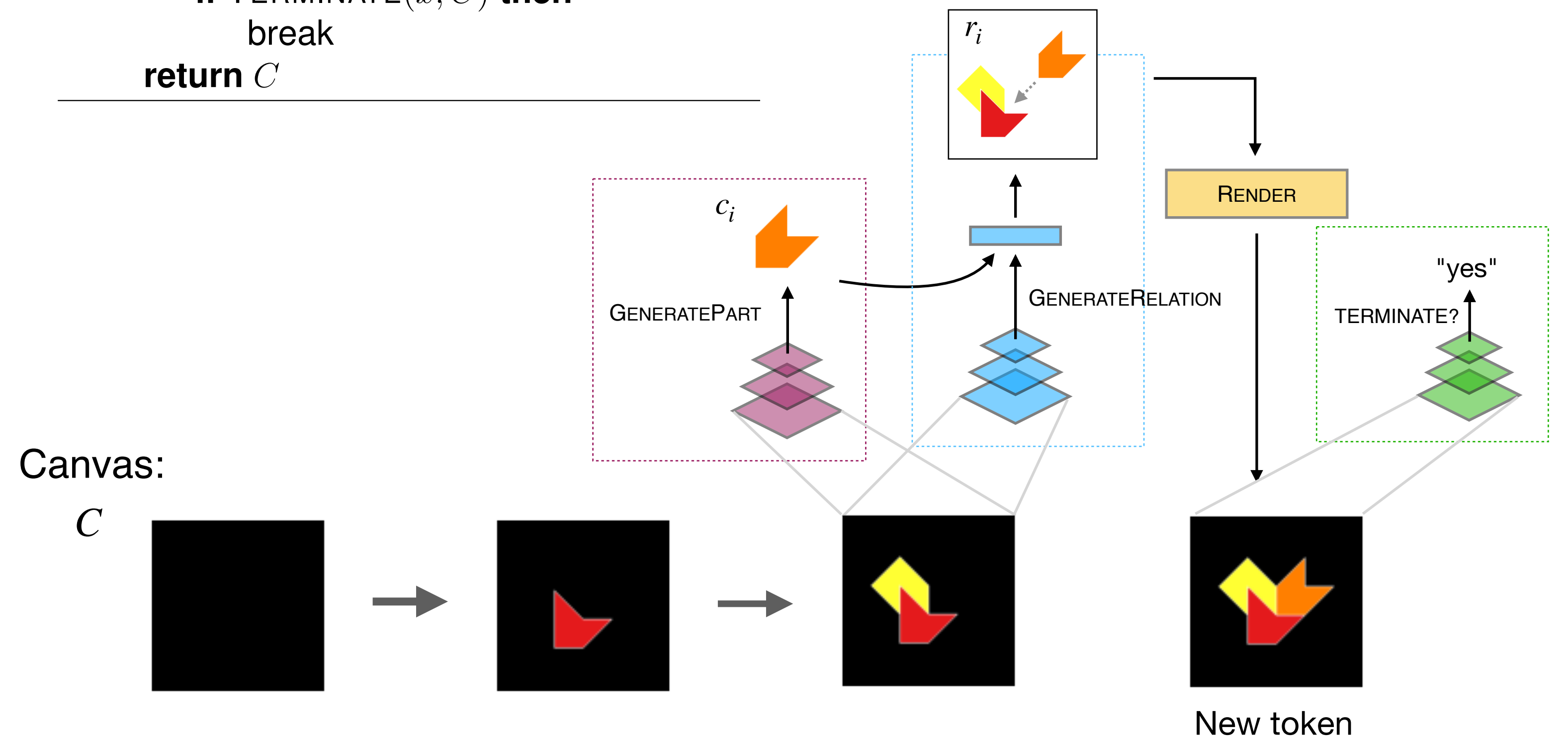
$C \leftarrow \text{RENDER}(C, c_i, r_i)$

**if** TERMINATE( $x, C$ ) **then**  
break

**return**  $C$

Canvas:

$C$

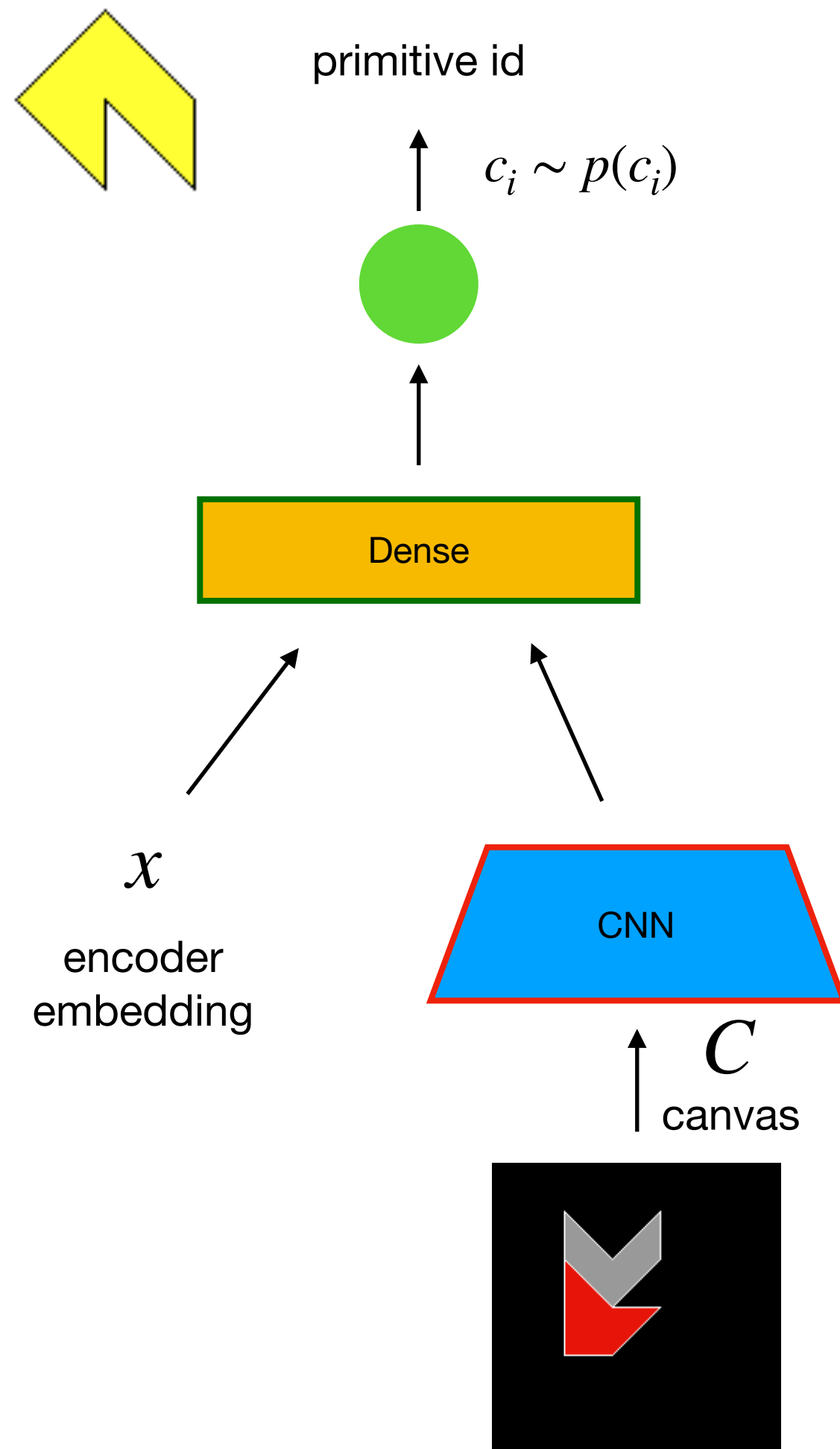




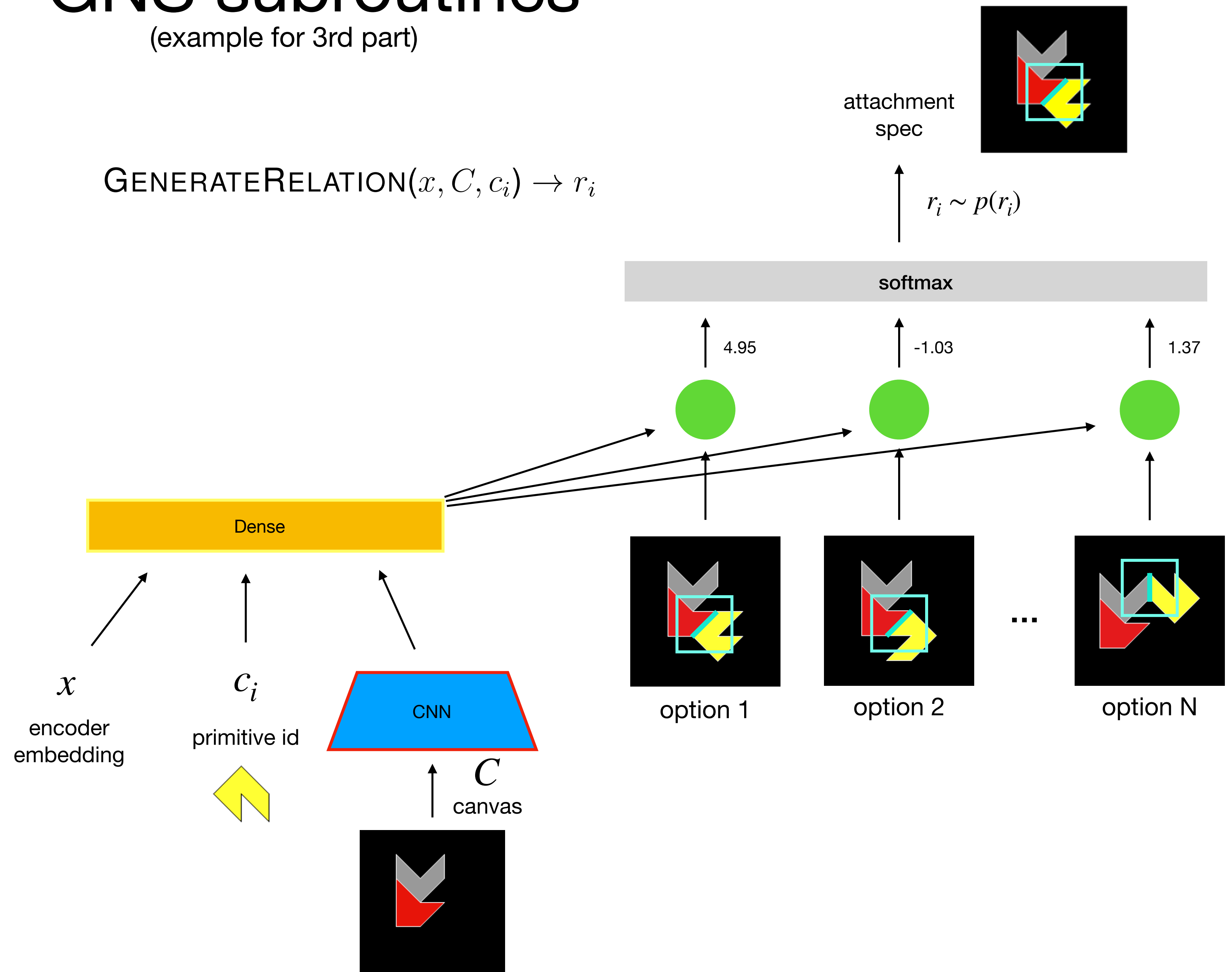
# GNS subroutines

(example for 3rd part)

$$\text{GENERATEPART}(x, C) \rightarrow c_i$$

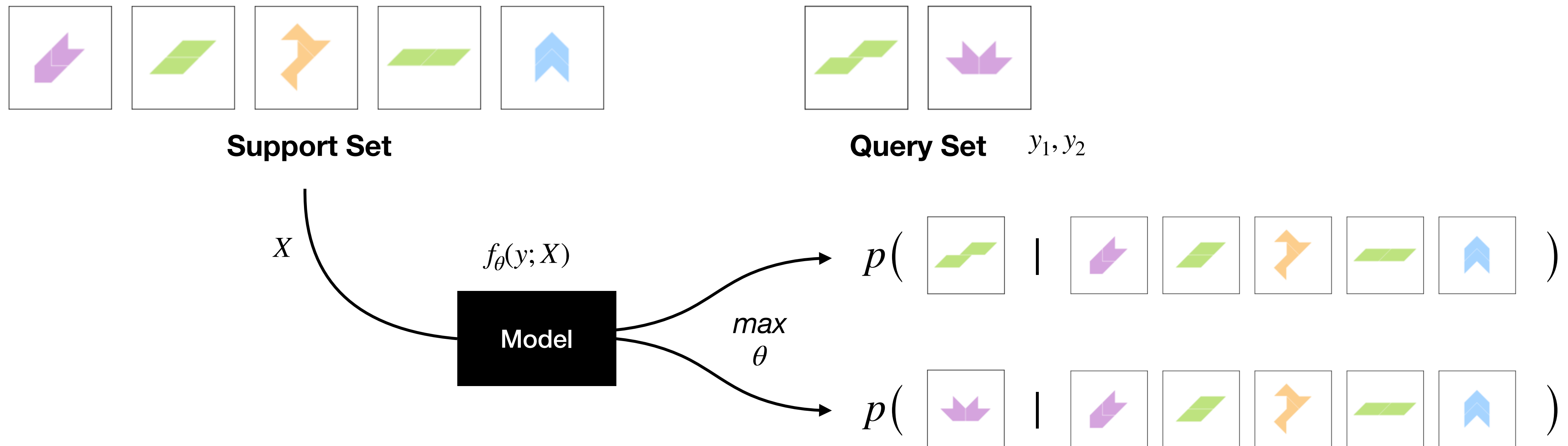


$$\text{GENERATERELATION}(x, C, c_i) \rightarrow r_i$$





# Meta-learning



**Objective:** maximize log-likelihood of query tokens conditioned on the support



# Meta-learning training data

bootstrapping the symbolic Bayesian model

P Synthetic data distribution	<div><div><b>procedure P</b> <math>h \sim p(h)</math> <math>S = x_1, \dots, x_n \sim p(x \mid h)</math> <math>Q = x'_1, \dots, x'_n \sim p(x \mid h)</math> <b>return</b> <math>S, Q</math></div><div><div>▷ Sample formula hypothesis from prior</div><div>▷ Sample support set from formula</div><div>▷ Sample query set from formula</div></div></div>
R <i>Resampled</i> synthetic data distribution	<div><div><b>procedure R</b> <math>S \sim \text{Uniform}(\Phi)</math> <math>h \sim p(h \mid S)</math> <math>Q = x'_1, \dots, x'_n \sim p(x \mid h)</math> <b>return</b> <math>S, Q</math></div><div><div>▷ Sample support set from human trials</div><div>▷ Sample formula hypothesis from posterior</div><div>▷ Sample query set from formula</div></div></div>
H Human distribution	<div><div><b>procedure H</b> <math>S, Q \sim \text{Uniform}(\Phi)</math> <b>return</b> <math>S, Q</math></div><div><div>▷ Sample support &amp; query sets from human trials</div></div></div>
C Bias training distribution	(see Section 4.4 and Appendix B.2)



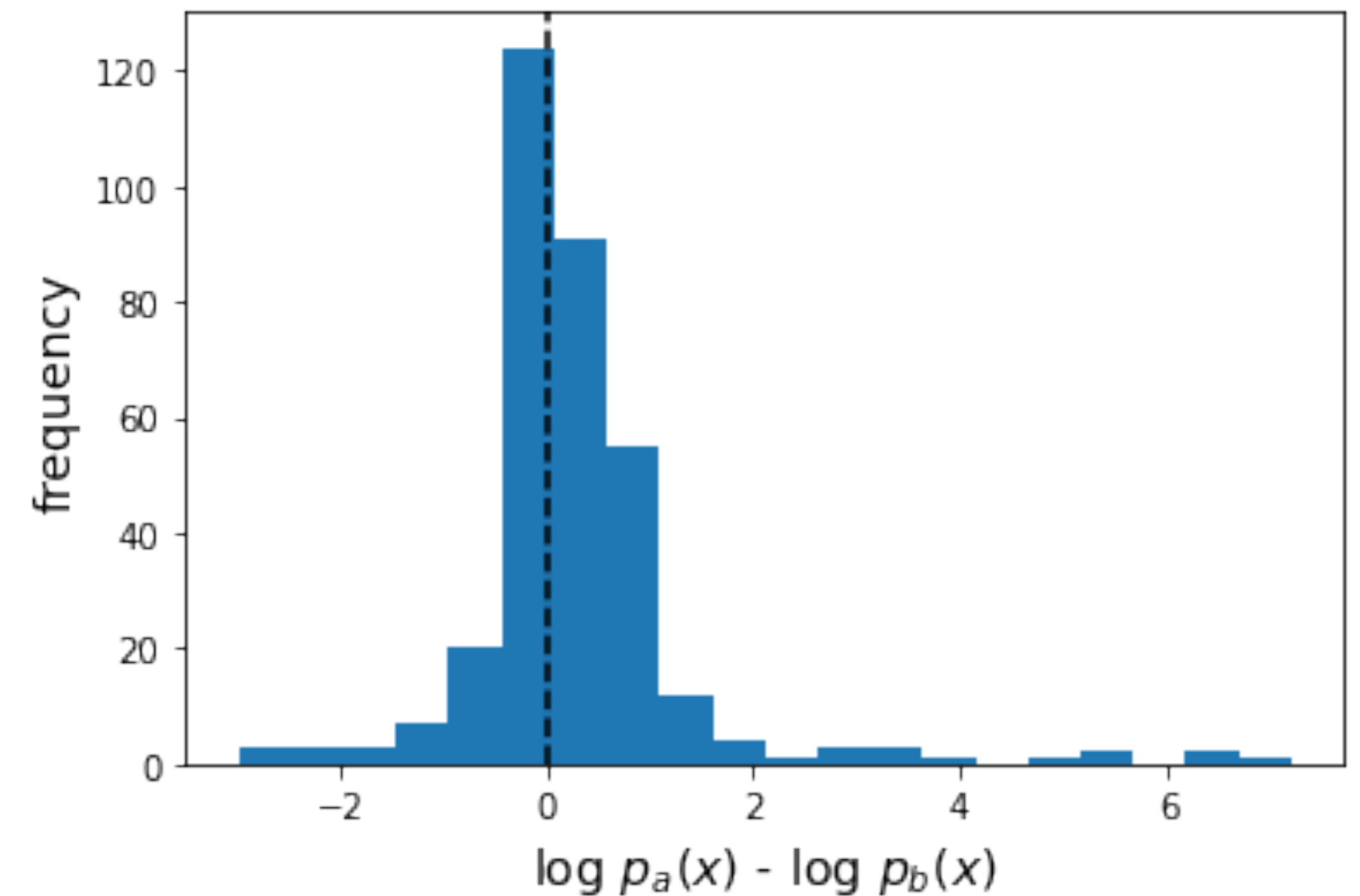
# Log-likelihood evaluations

	Test log-likelihood
Bayesian	-4.741
GNS (P/R/H/C)	<b>-4.444</b>
GNS (P/R/H)	-4.535
GNS (P/R)	-4.645
GNS (P)	-4.930

**Likelihood of held-out human generations.** For each model, the total log-likelihood averaged over the holdout set is reported.

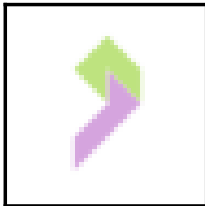
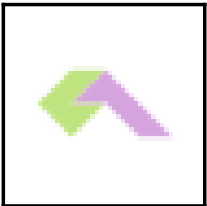
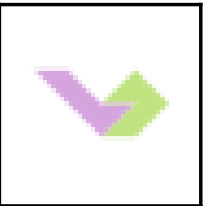
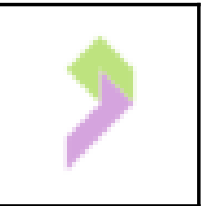
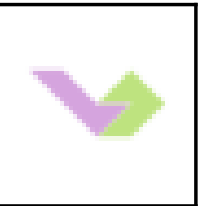
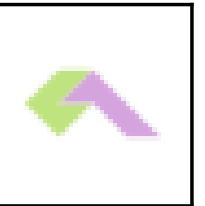
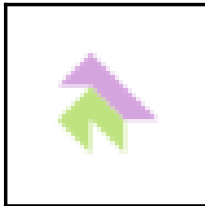
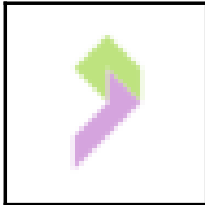
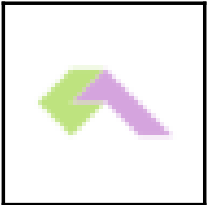
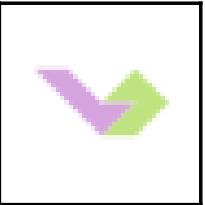
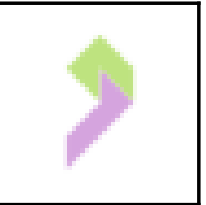
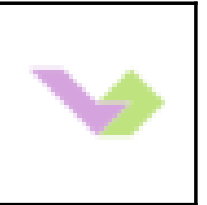
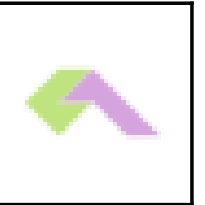
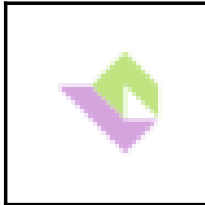
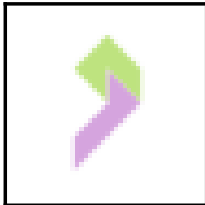
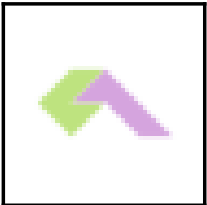
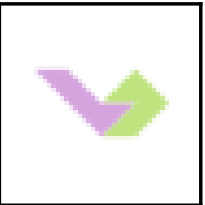
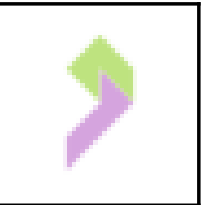
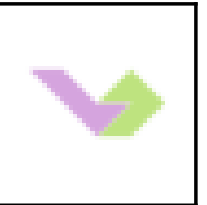
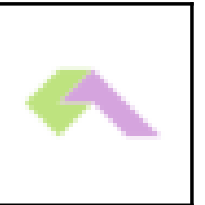
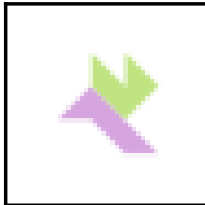
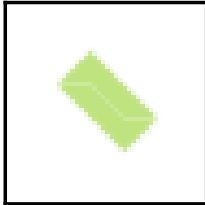
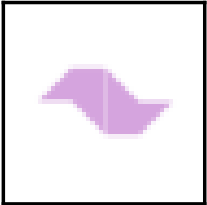
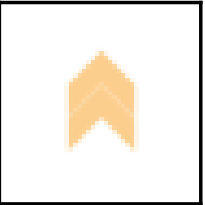
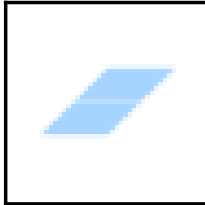
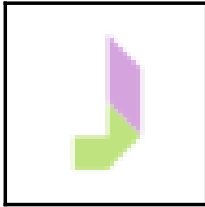
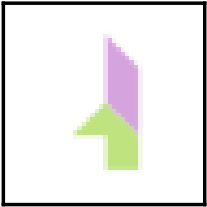
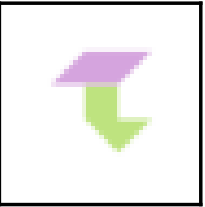
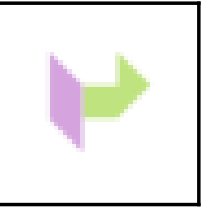
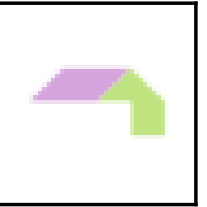
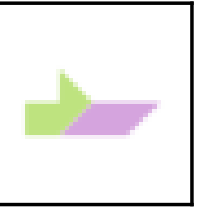


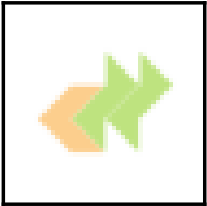

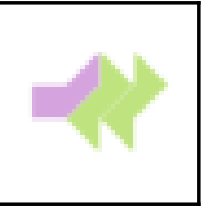
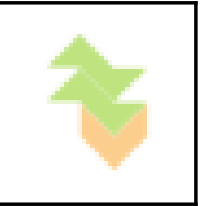


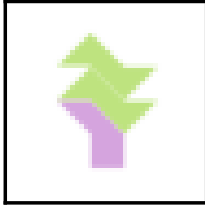
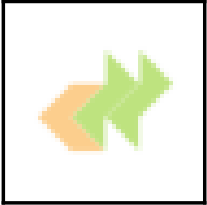

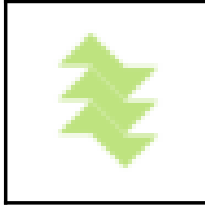




Paired t-test comparing per-example log-likelihood of GNS (P/R/H/C) vs. Bayesian

$$t(336) = 6.197, \quad p < 0.001$$





# Log-likelihood evaluations

	support set						new token	freq.	delta
1								(2)	4.70
2								(1)	3.22
3								(1)	3.19
4								(5)	2.18
5								(1)	1.90
6								(1)	1.69
7								(3)	1.62
8								(1)	1.47



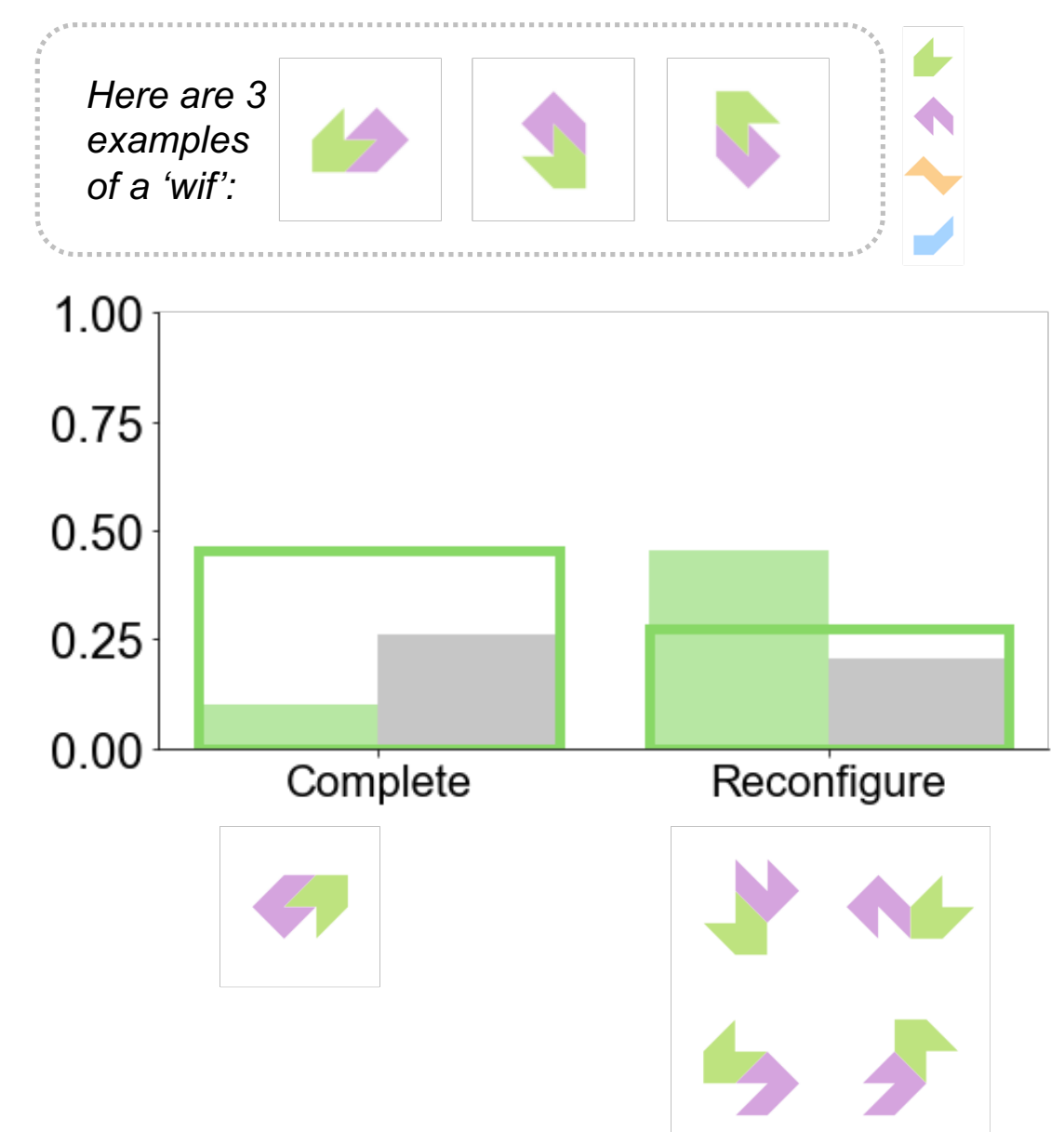
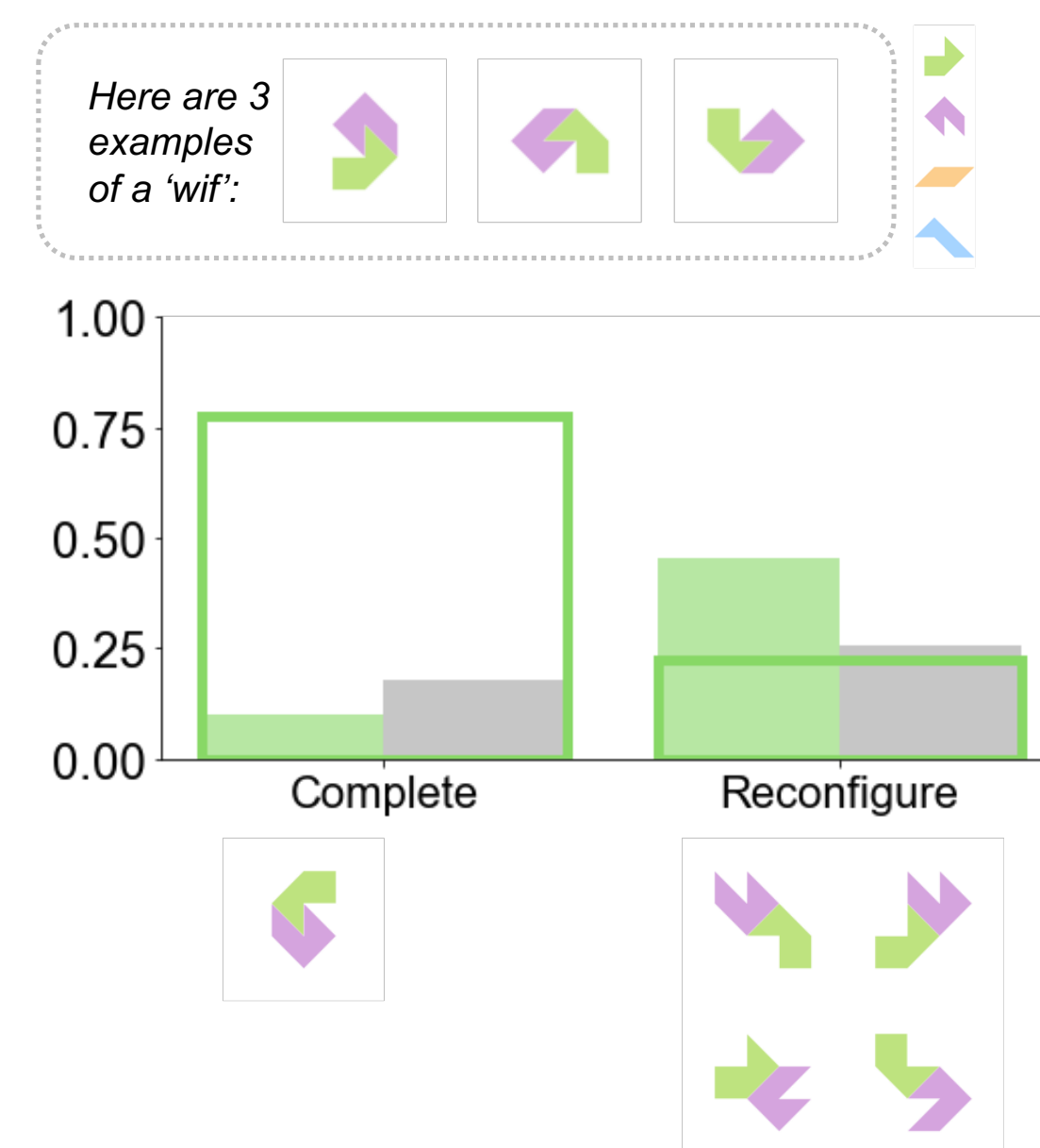
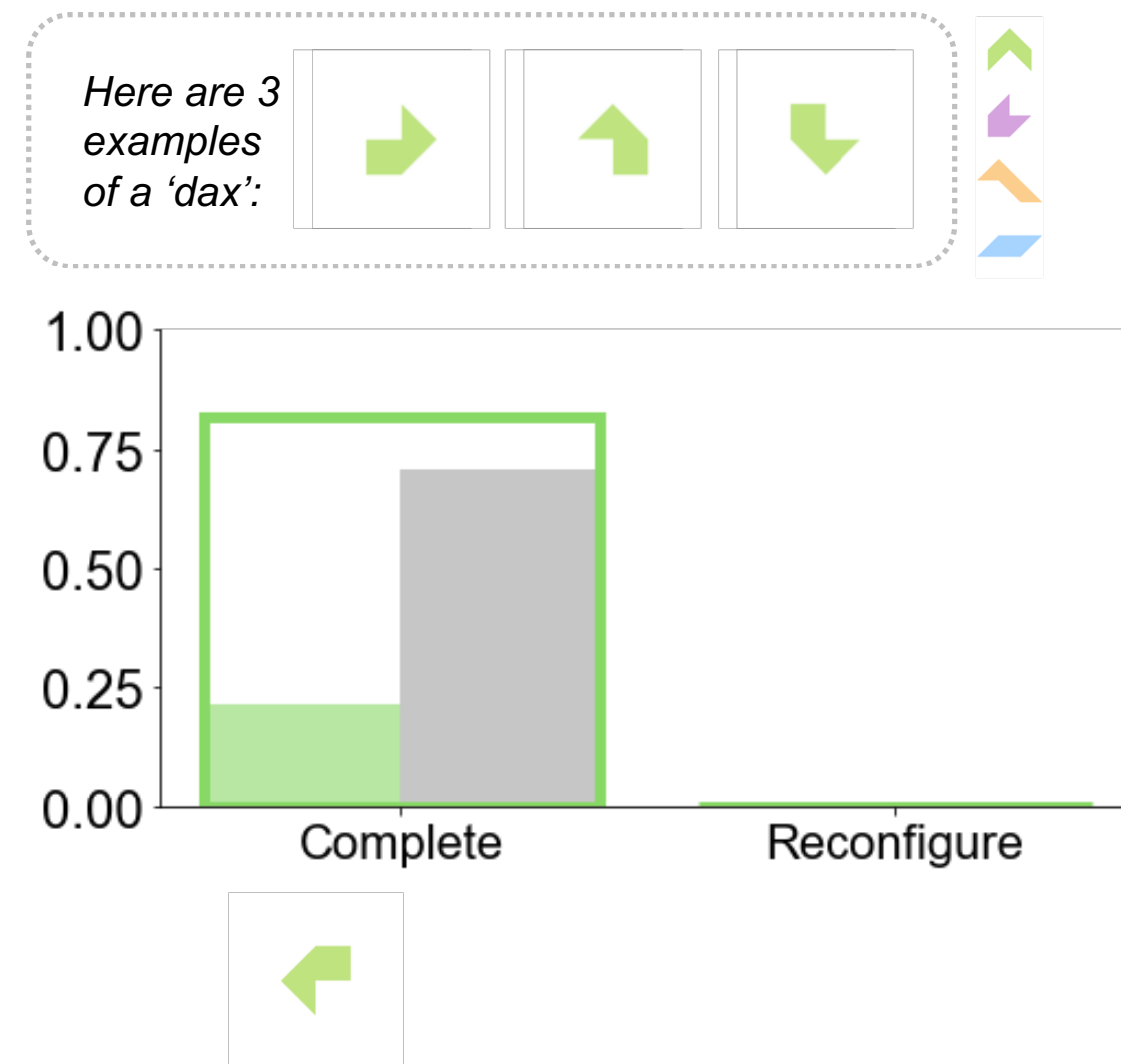
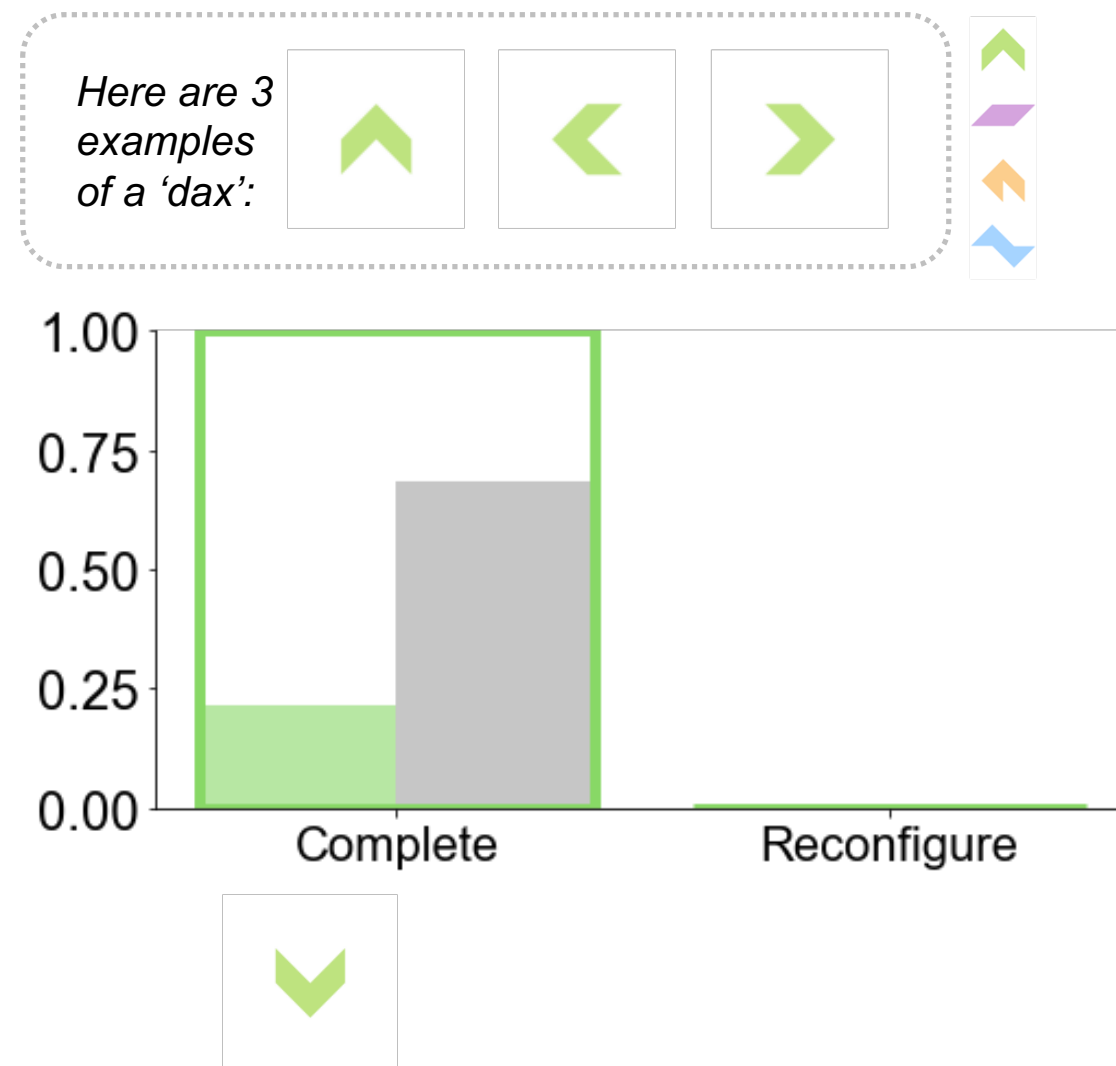
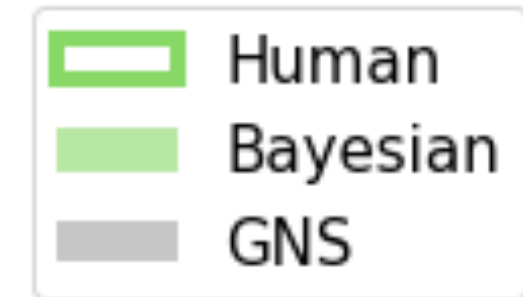
# Accounting for human inductive biases





# Accounting for human inductive biases

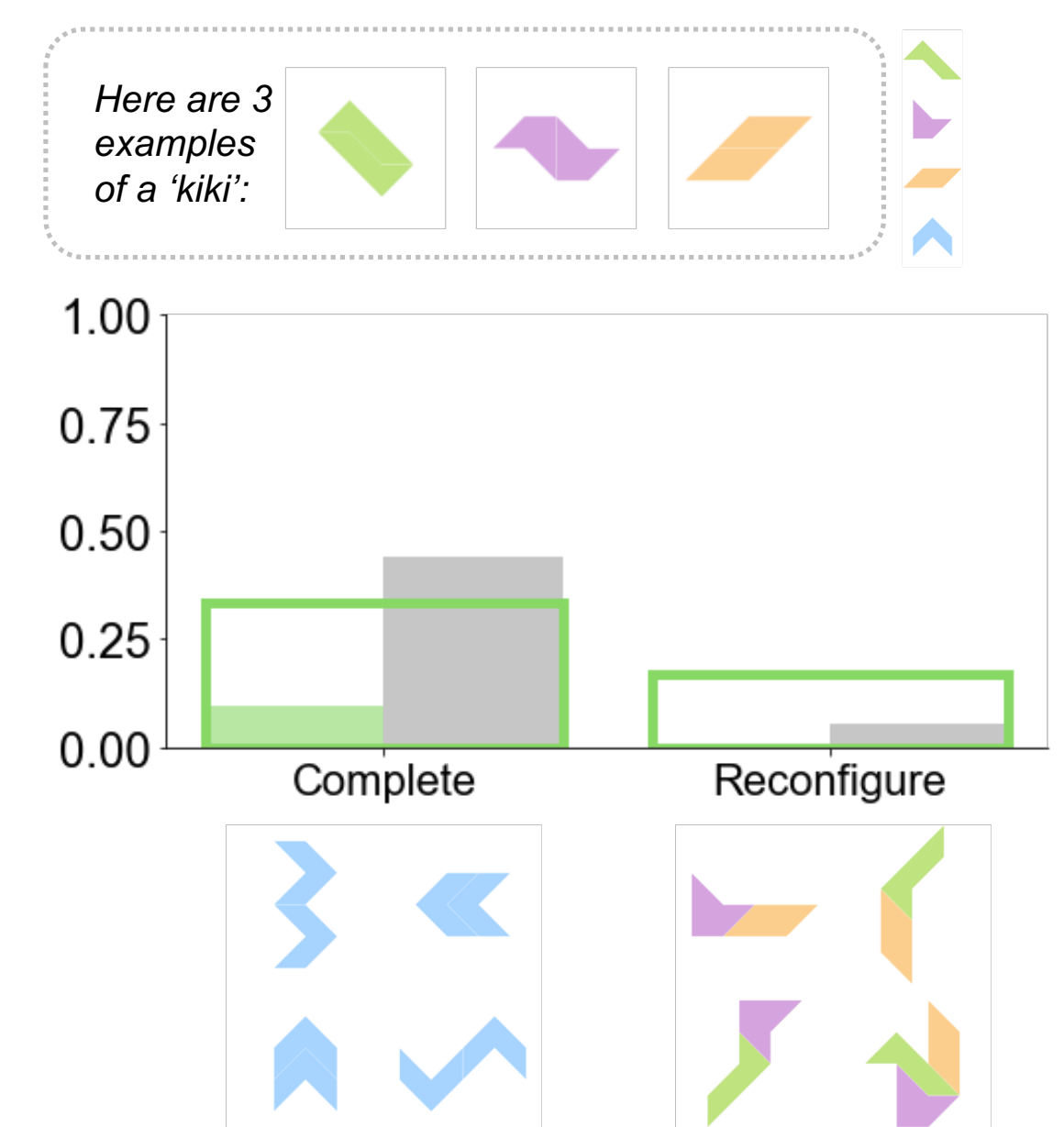
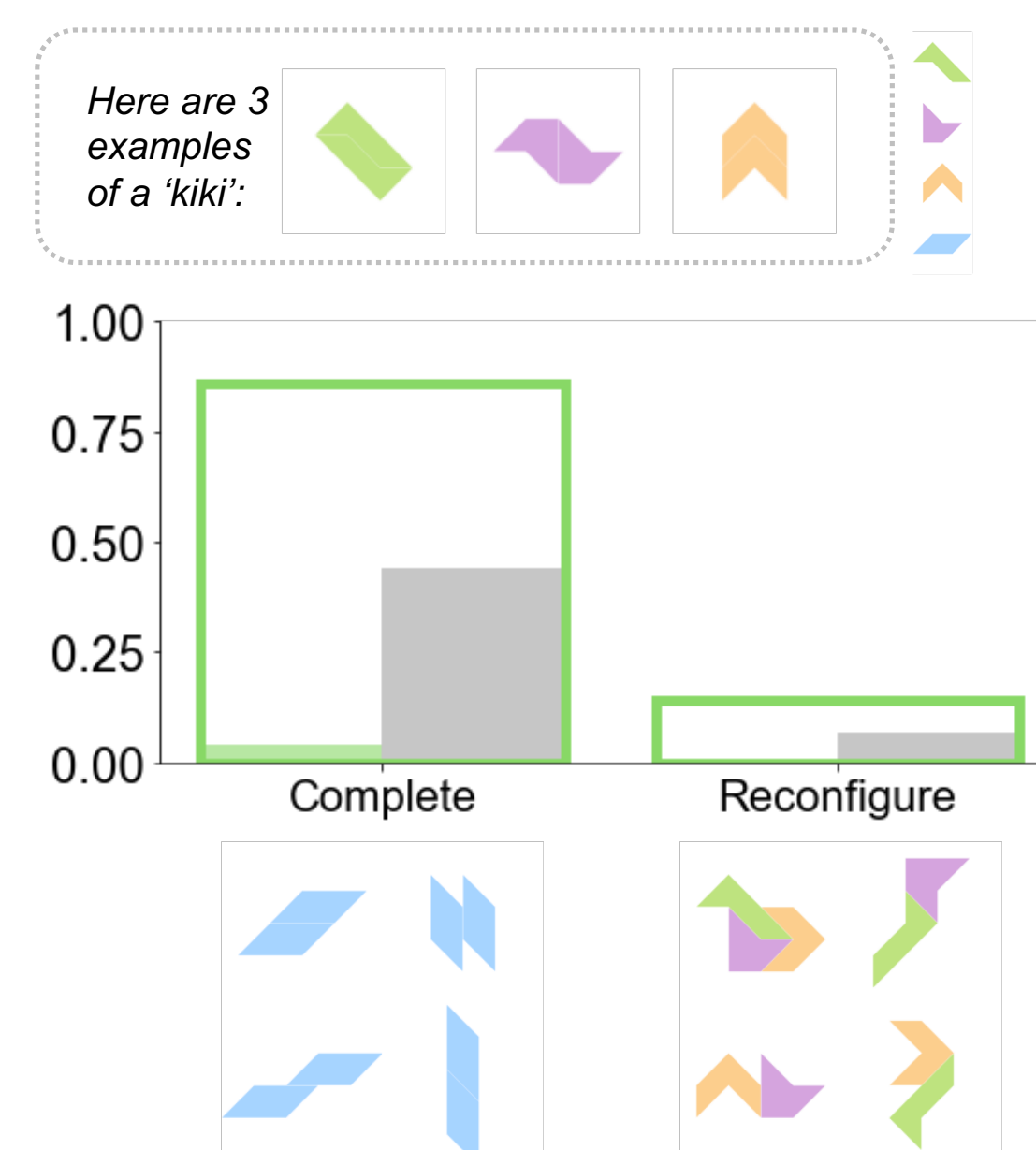
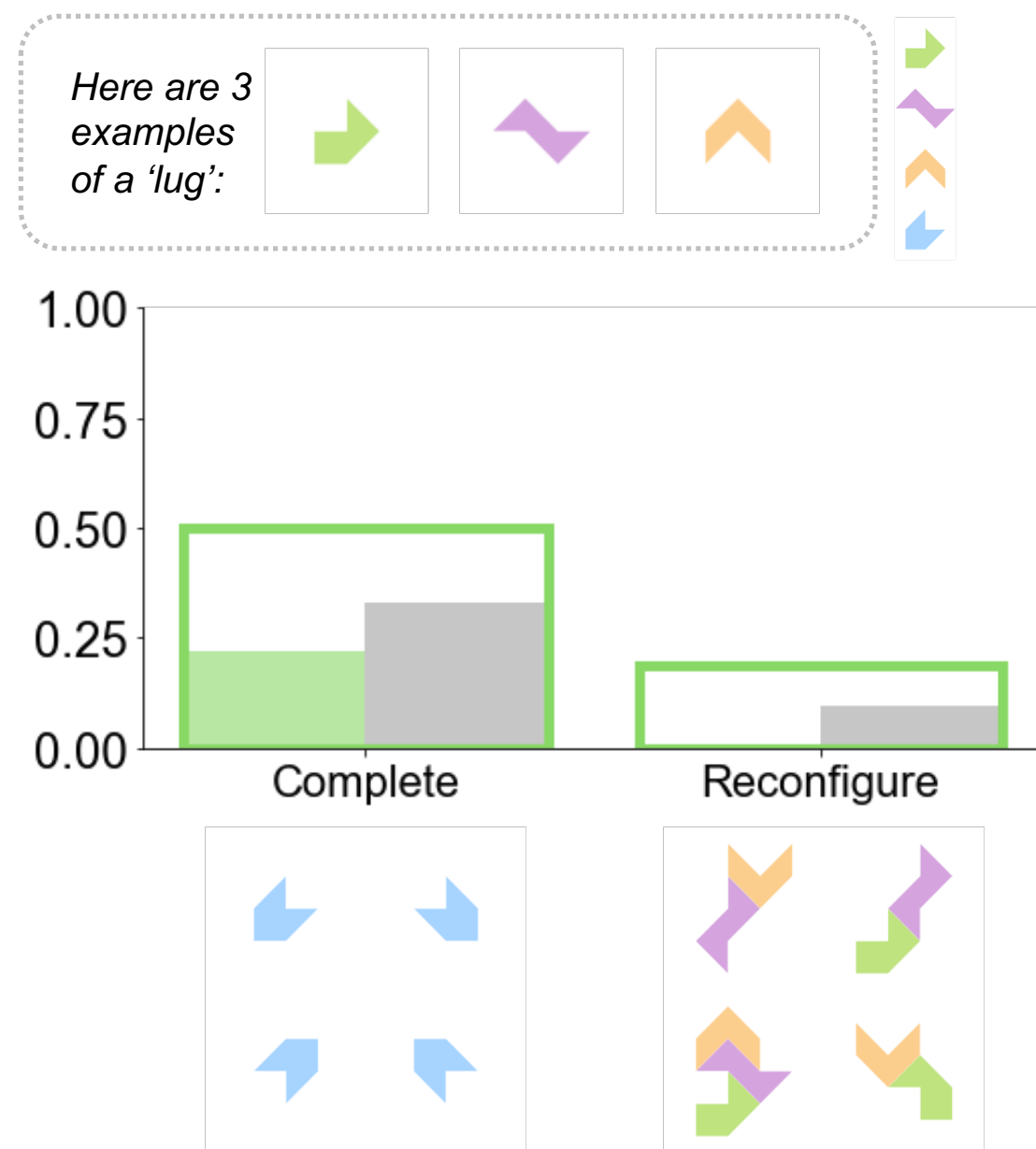
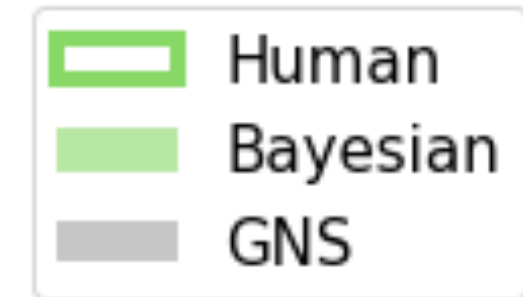
"rotations" trial type





# Accounting for human inductive biases

"primitives" trial type





# Conclusions: Case study #2

- GNS models are an effective way to understand and simulate human few-shot learning of structured visual concepts
- Compared to a strong symbolic baseline model, GNS provides an improved likelihood account of human few-shot generation
- GNS can account for human inductive biases that are not well-explained by alternatives



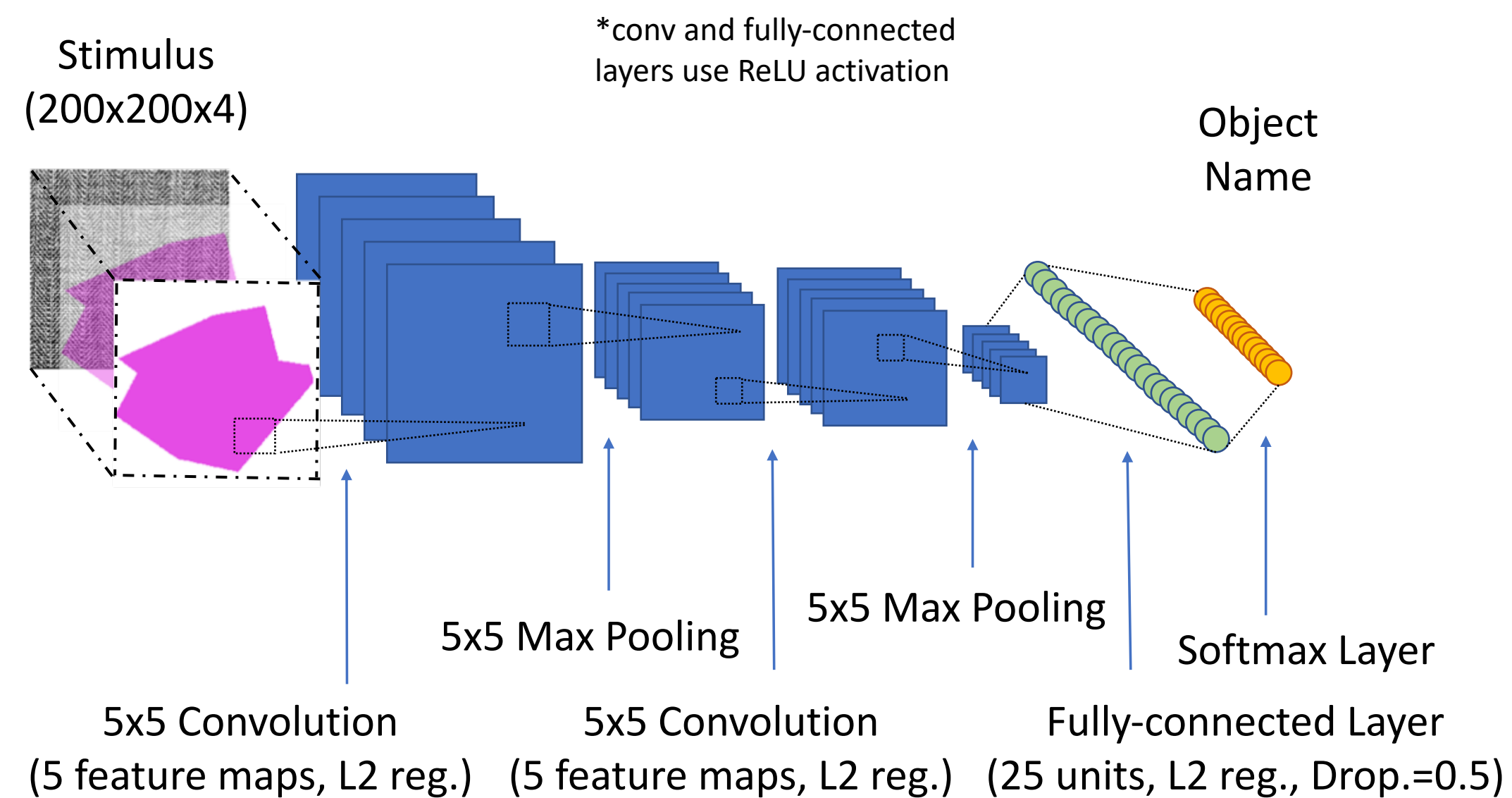
**Additional projects**



# Learning inductive biases with simple neural networks

(Feinman & Lake, 2018)

## Convolutional neural network (CNN) architecture



## Shape bias test

This is a “dax.”



(Smith et al., 2002)

Where is the other “dax?”

1



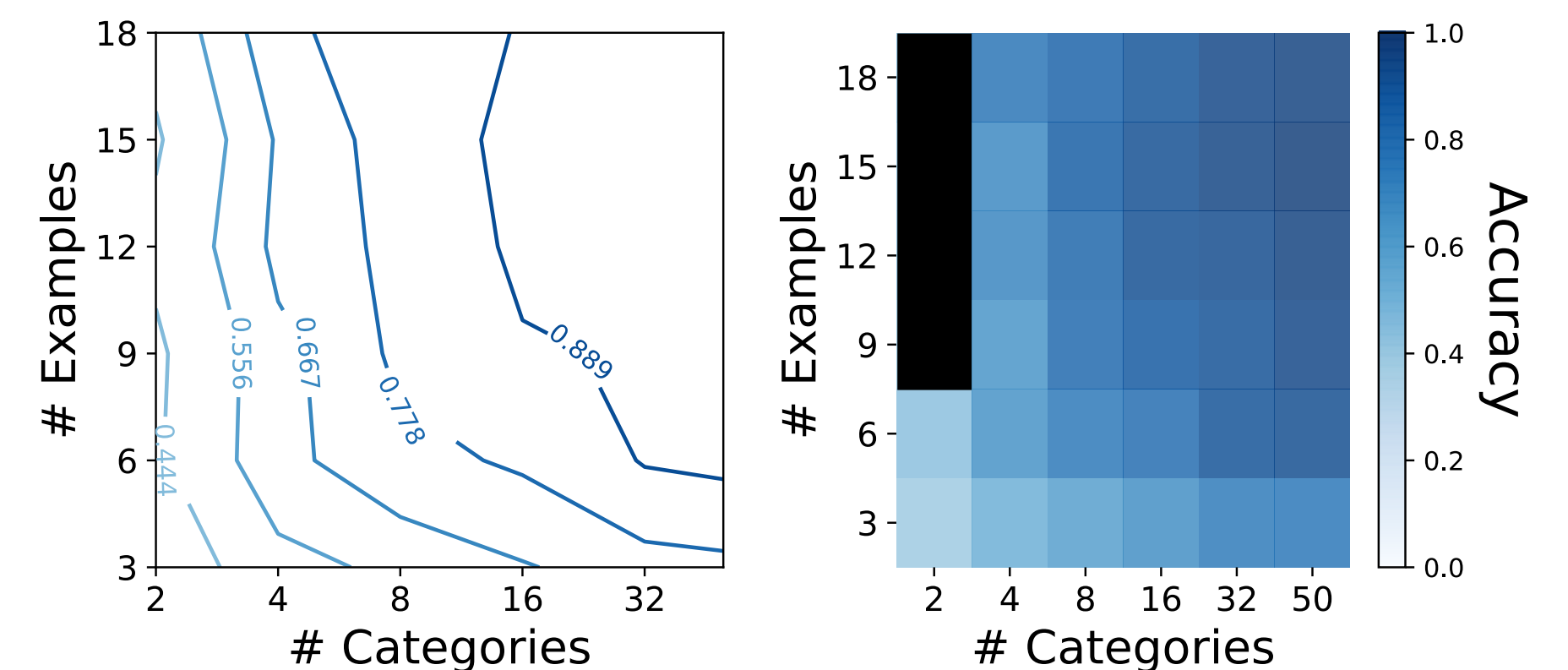
2



3



## CNN shape bias strength vs. dataset size





# Learning a smooth kernel regularizer for convolutional neural networks

(Feinman & Lake, 2019)

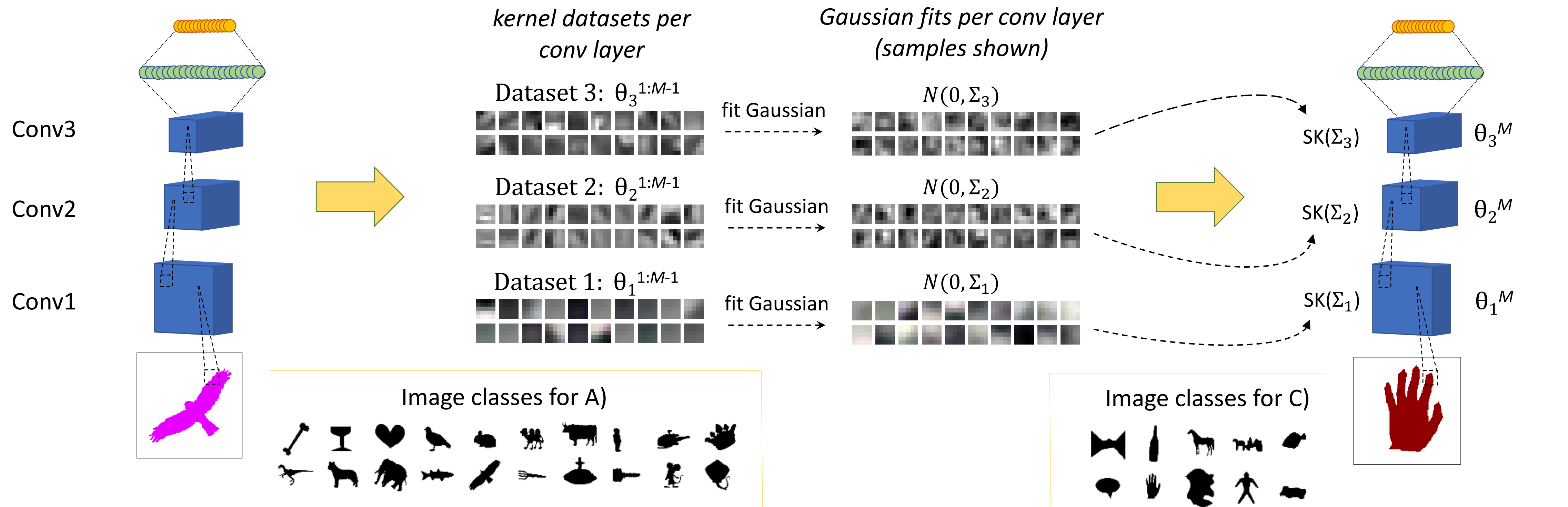
## PHASE 1

## PHASE 2

### A) Train CNN (repeat 20x)

### B) Extract kernel statistics

### C) Apply SK-reg to new task





# Summary & Conclusions



# introduced Generative Neuro-Symbolic (GNS) modeling

---

**procedure** GENERATEEXAMPLE

---

$C \leftarrow 0$

▷ Initialize blank canvas

**for**  $i = 1 \dots \infty$  **do**

$x_i \leftarrow \text{GENERATEPART}(C)$

▷ Sample part

$r_i \leftarrow \text{GENERATERELATION}(C, x_i)$

▷ Sample relation

$C \leftarrow \text{RENDER}(C, x_i, r_i)$

▷ Render new canvas

**if**  $\text{TERMINATE?}(C)$  **then**

▷ Sample termination (y/n)

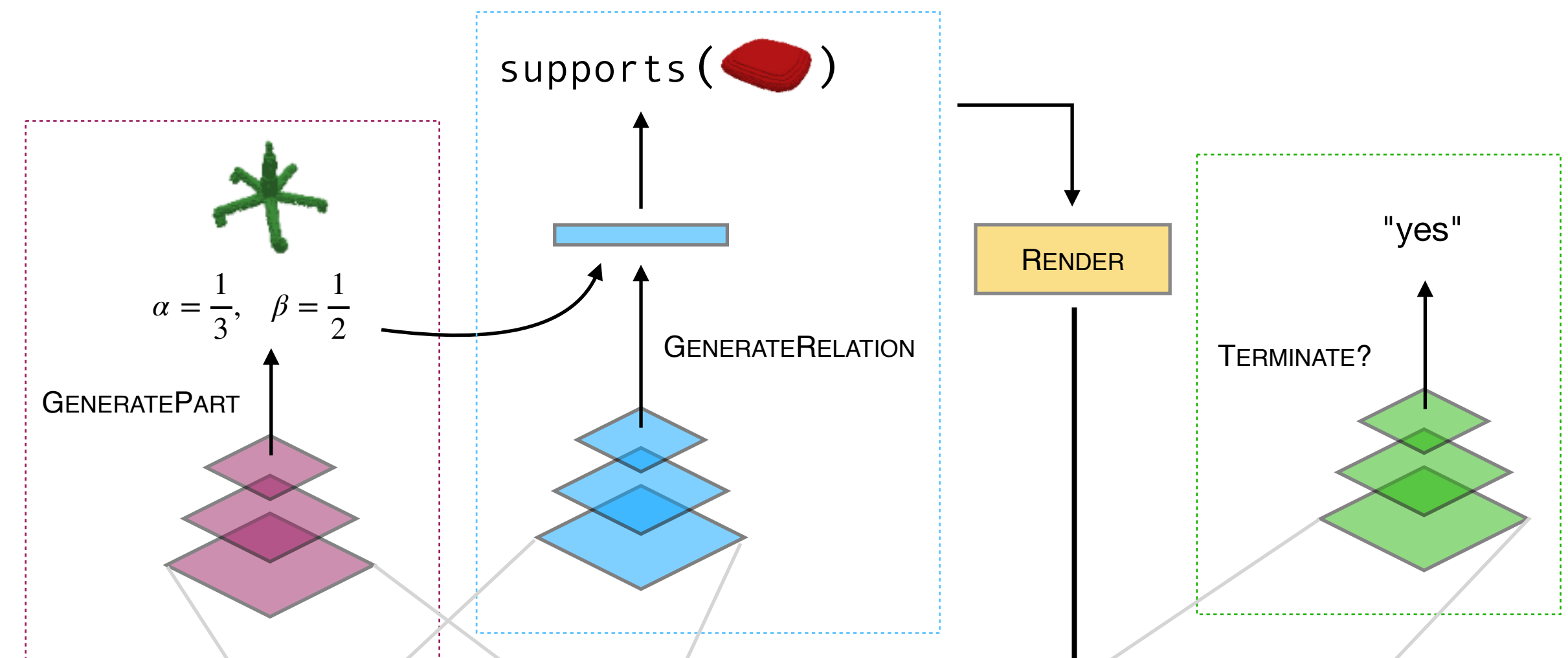
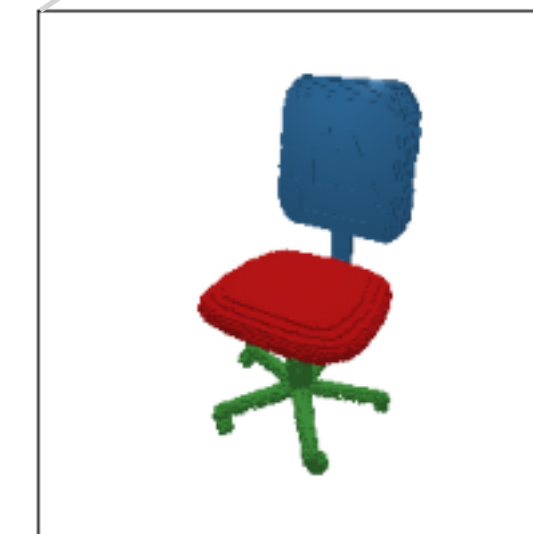
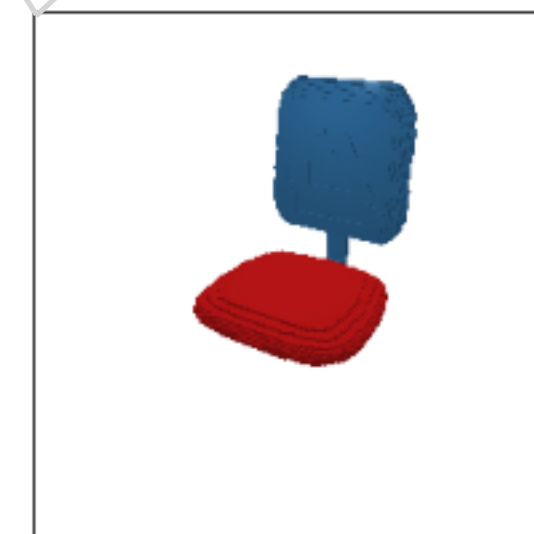
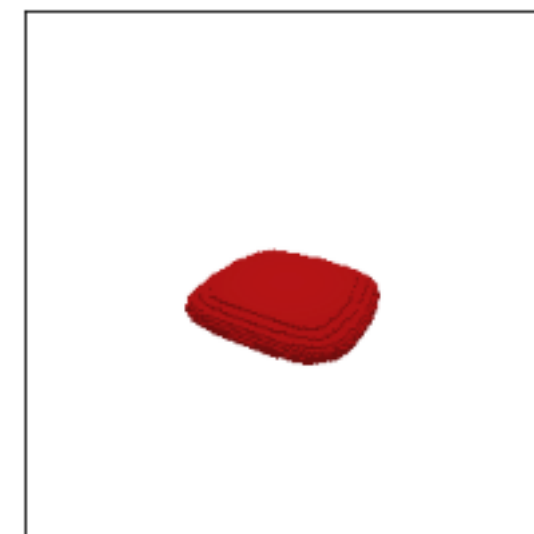
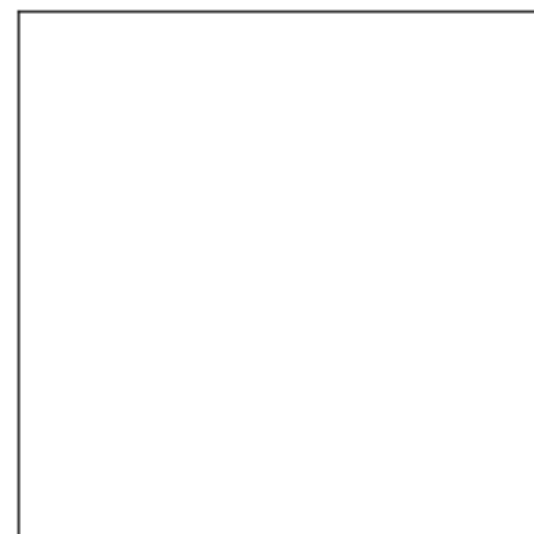
break

**return**  $C$

▷ Return example

---

Canvas:  
 $C$





# GNS model of handwritten character concepts

---

**procedure** GENERATECHARACTER

$C \leftarrow 0$

**for**  $i = 1 \dots \infty$  **do**

$r_i \leftarrow \text{GENERATERELATION}(C)$

$x_i \leftarrow \text{GENERATEPART}(C, r_i)$

$C \leftarrow \text{RENDER}(C, x_i, r_i)$

$v_i \leftarrow \text{TERMINATE?}(C)$

**if**  $v_i$  **then**

break

**return**  $C$

---

▷ Initialize blank canvas

▷ Sample relation

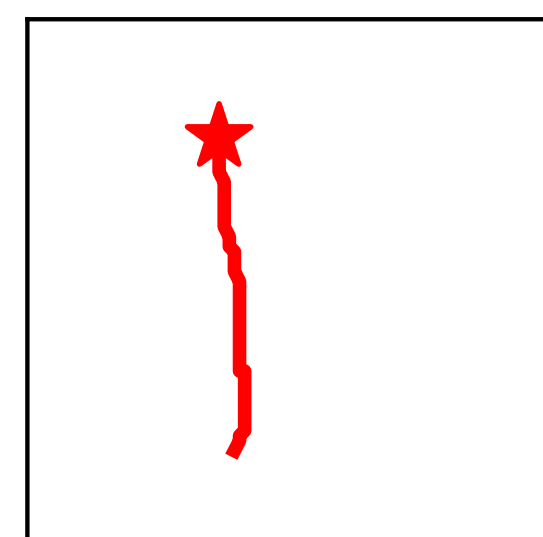
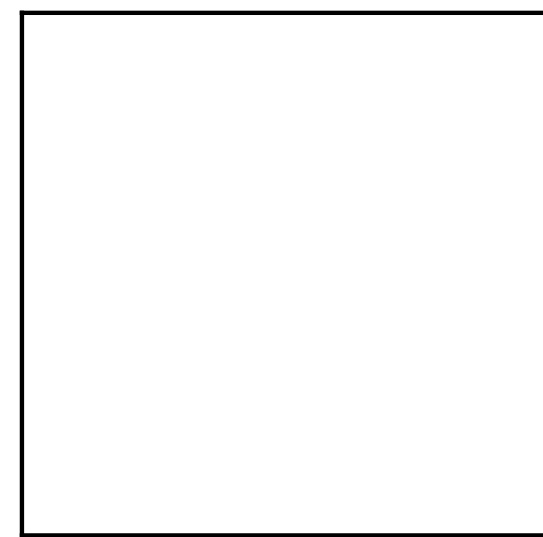
▷ Sample part

▷ Render to canvas

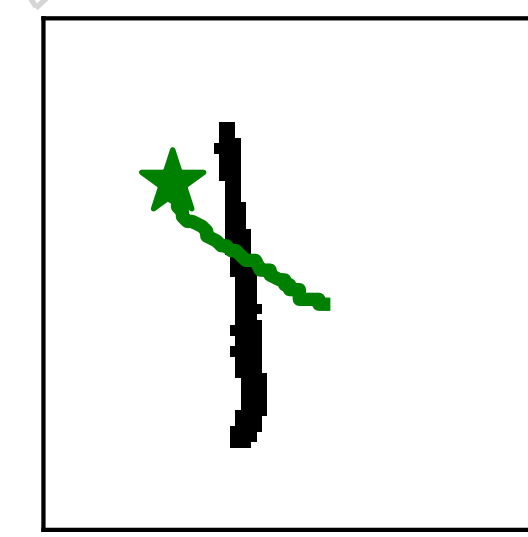
▷ Sample termination indicator

▷ Terminate sample

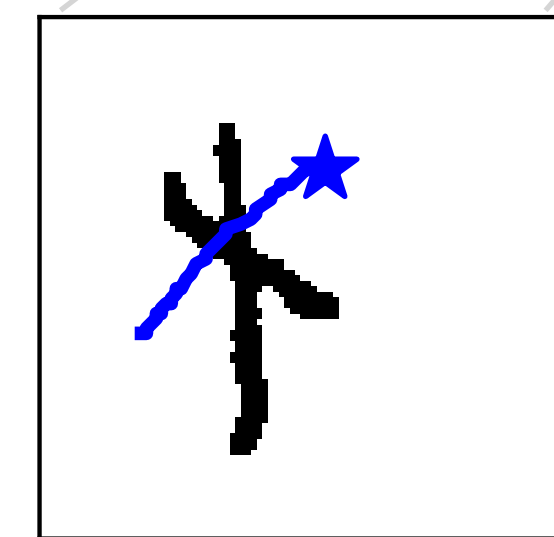
Canvas:  
 $C$



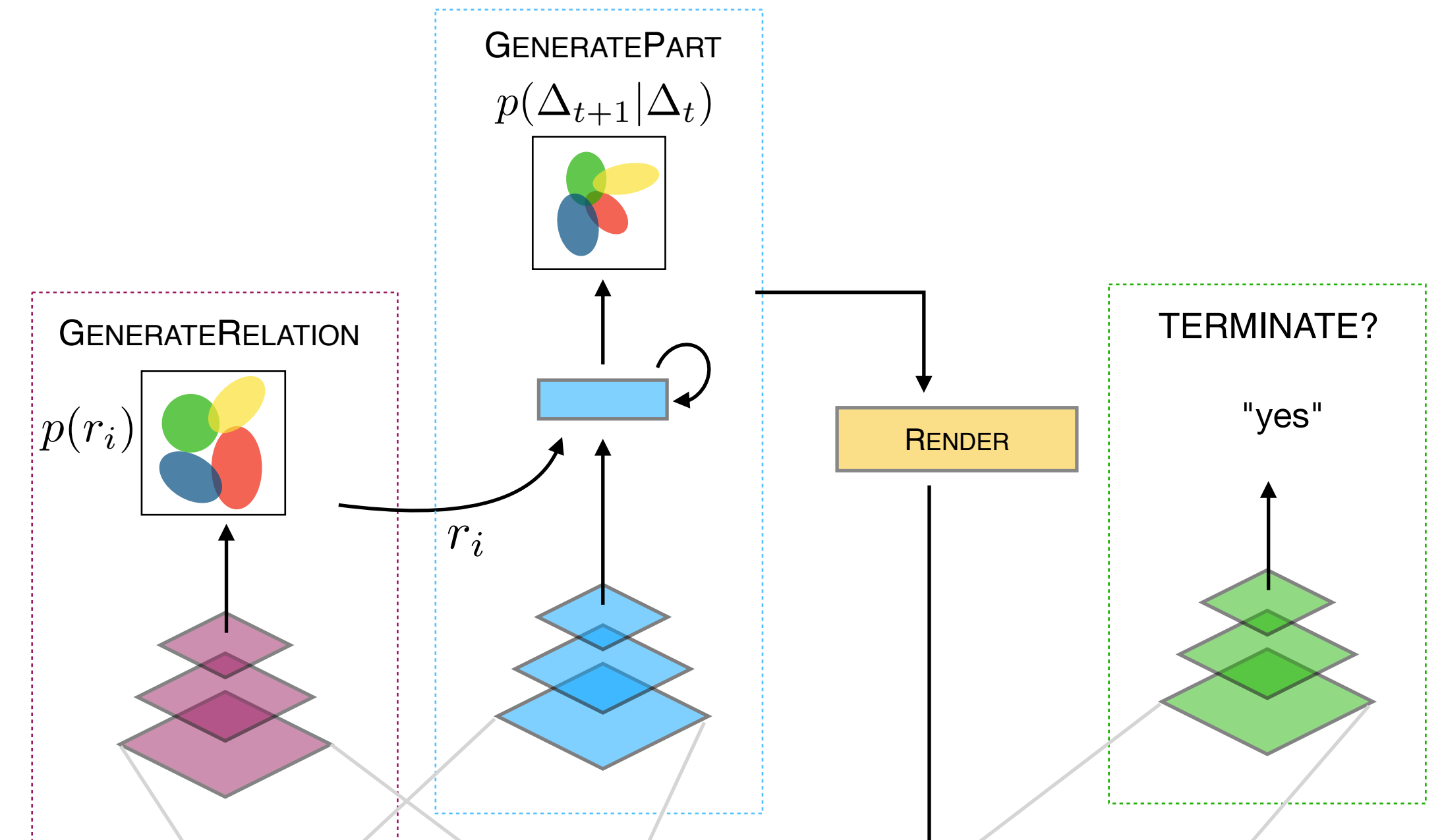
$i = 1$



$i = 2$

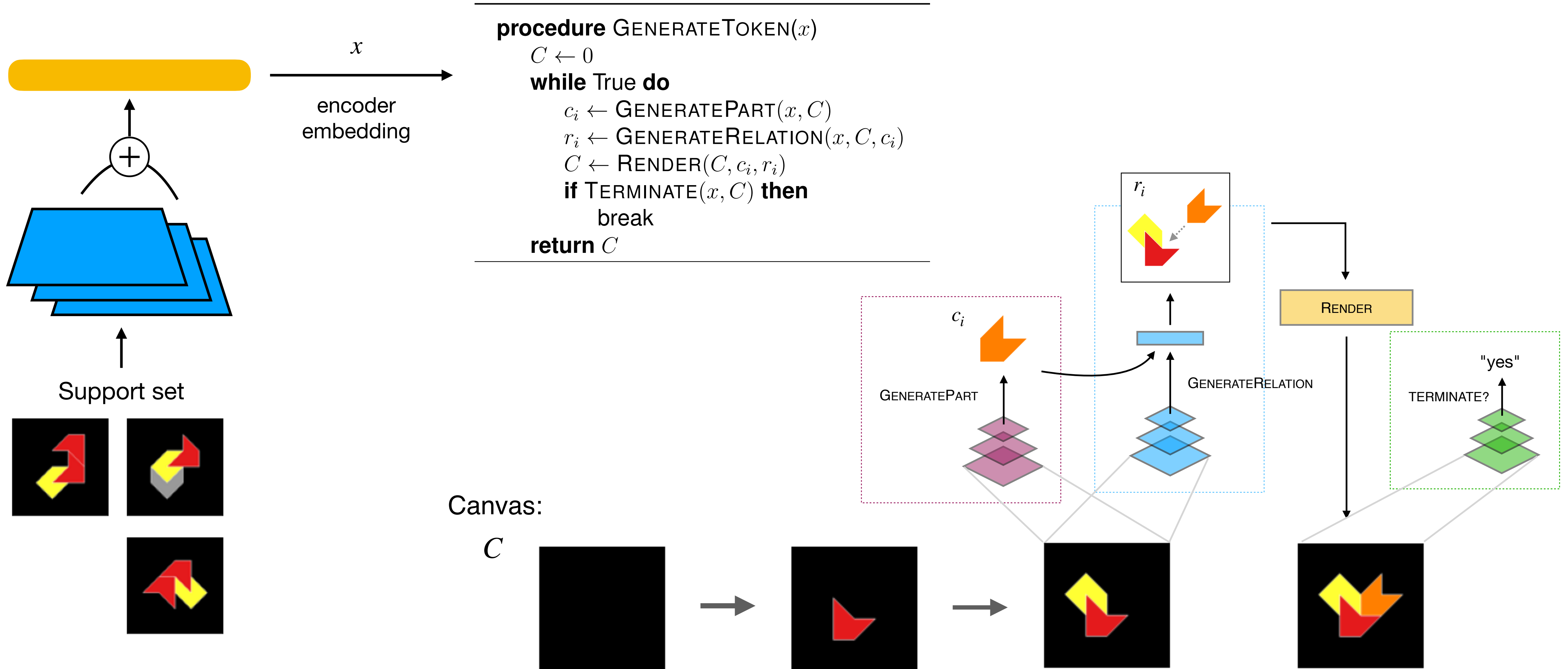


$i = 3$





# GNS model of synthetic part-based concepts ("alien figures")





# General conclusions

- Generative neuro-symbolic (GNS) modeling provides a novel synthesis of ideas from the structured and statistical modeling traditions
- By combining these ingredients in a computational model, we can account for human concept learning in ways that purely- symbolic and neural models fall short
- GNS models can help us understand the dual structural and statistical natures of human knowledge and direct us toward a more accurate representation of concepts



# Thank You

Brenden Lake

Yanli Zhou

Guy Davidson

Emin Orhan

Wai Keen Vong

+ HMLL lab

Tuan-Anh Le

Maxwell Nye

Joshua Tenenbaum

Lucas Tian

+ CoCoSci lab

Nikhil Parthasarathy

NYU Neuroscience cohort

Andy Feinman

Mary Van Hoomissen

Nick Feinman

Charlotte Walmsley



# Questions?

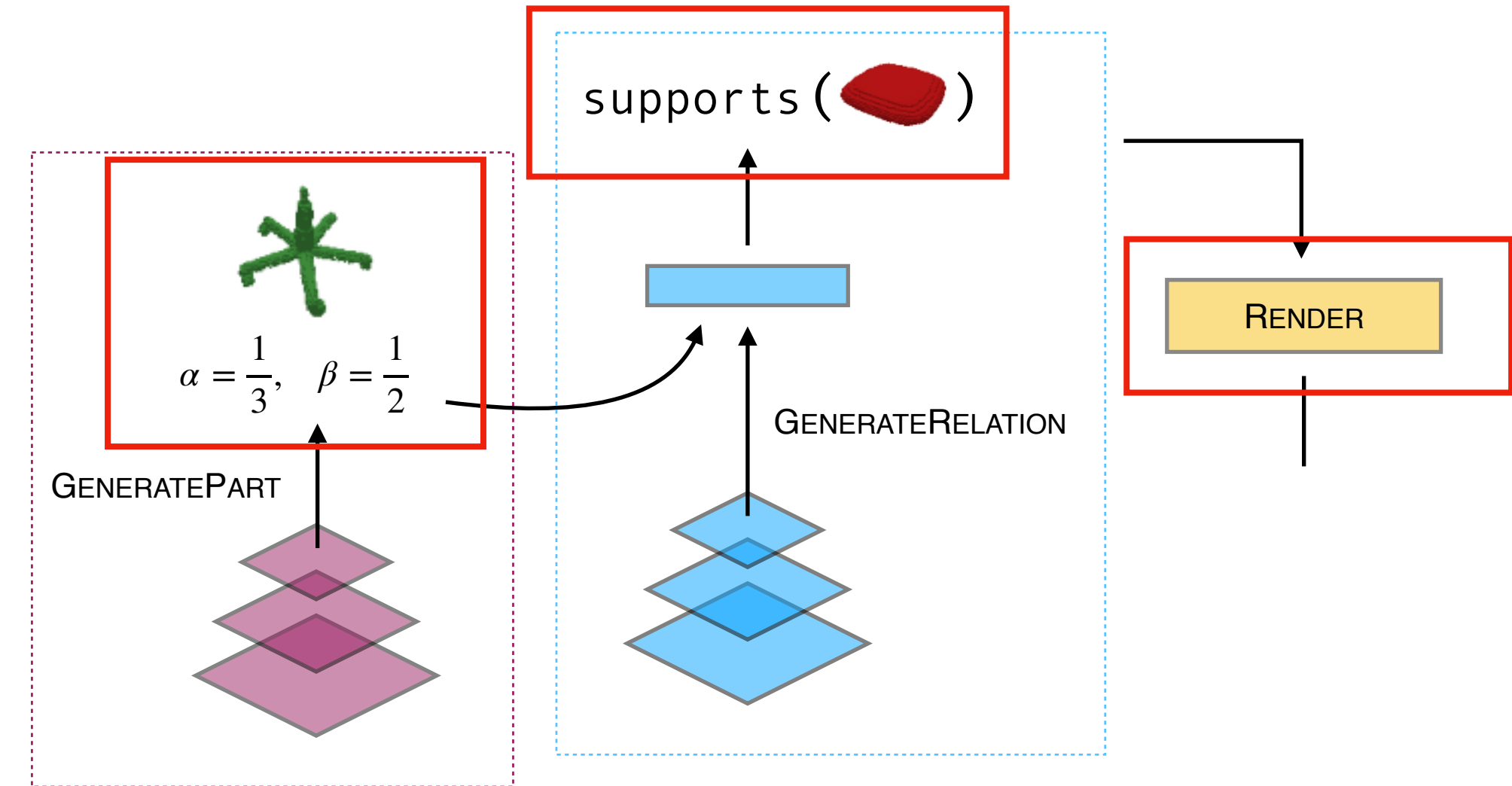
*"What I cannot create, I do not understand."*

—Richard Feynman

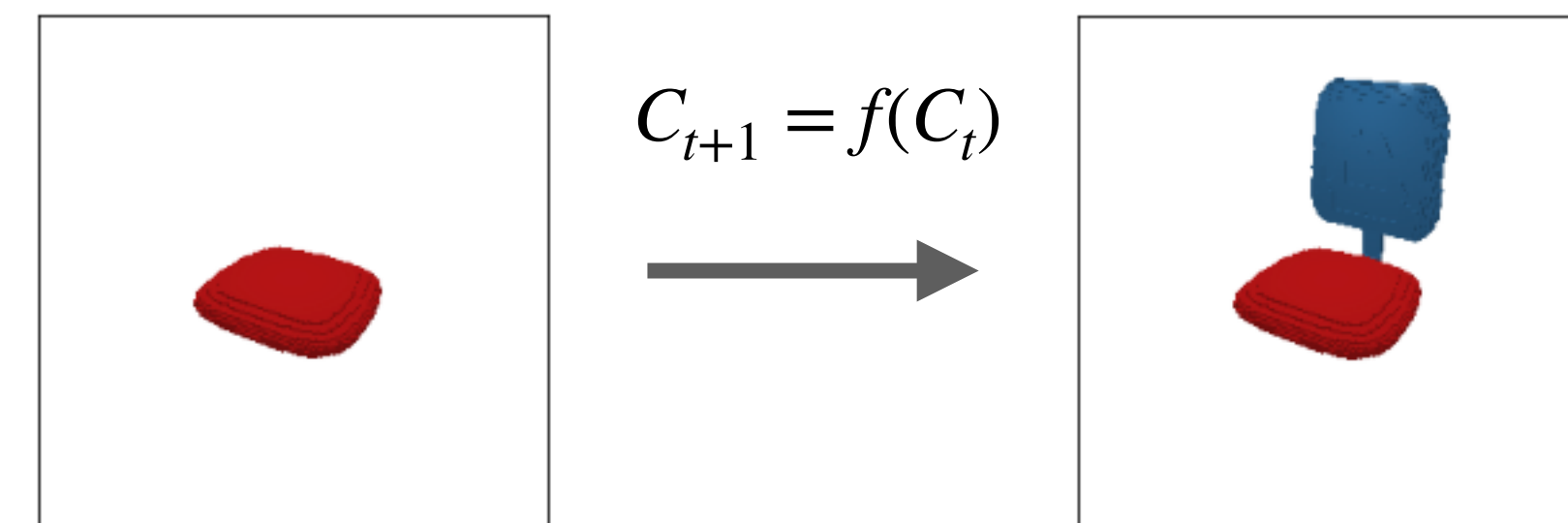


# Inductive biases of GNS architecture

Explicit notion of *causality* via symbolic primitives and renderer

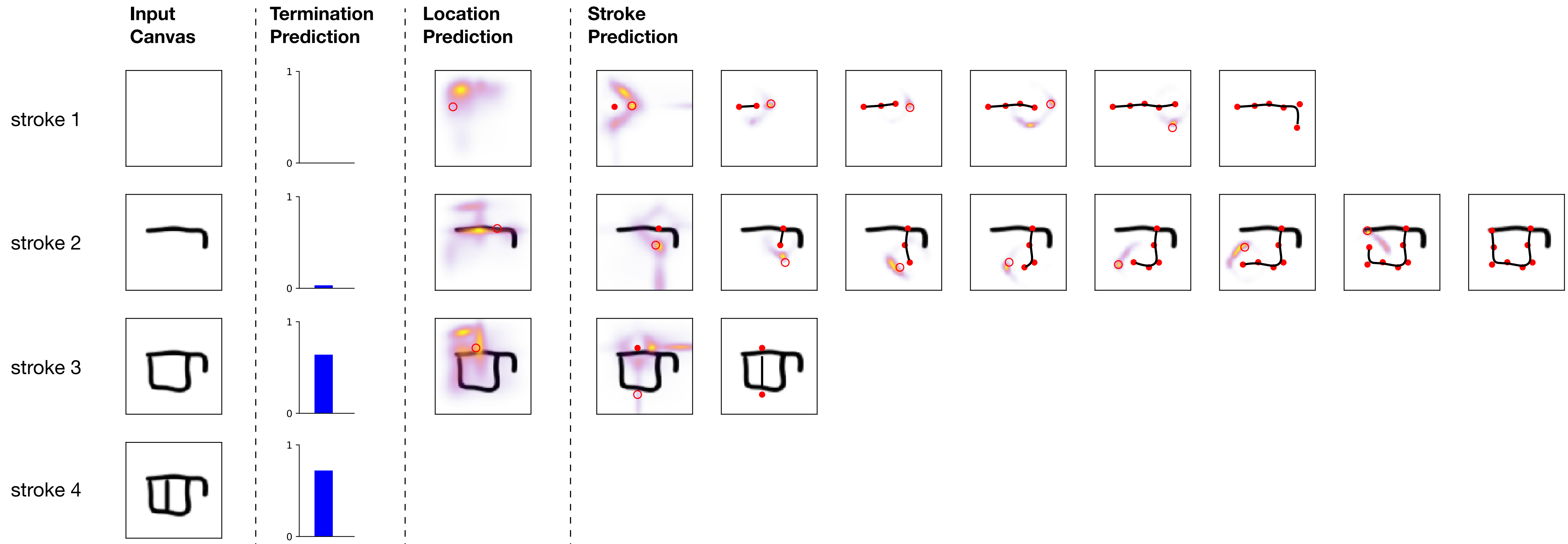


Compositional representation via modular subroutines and controlled memory state



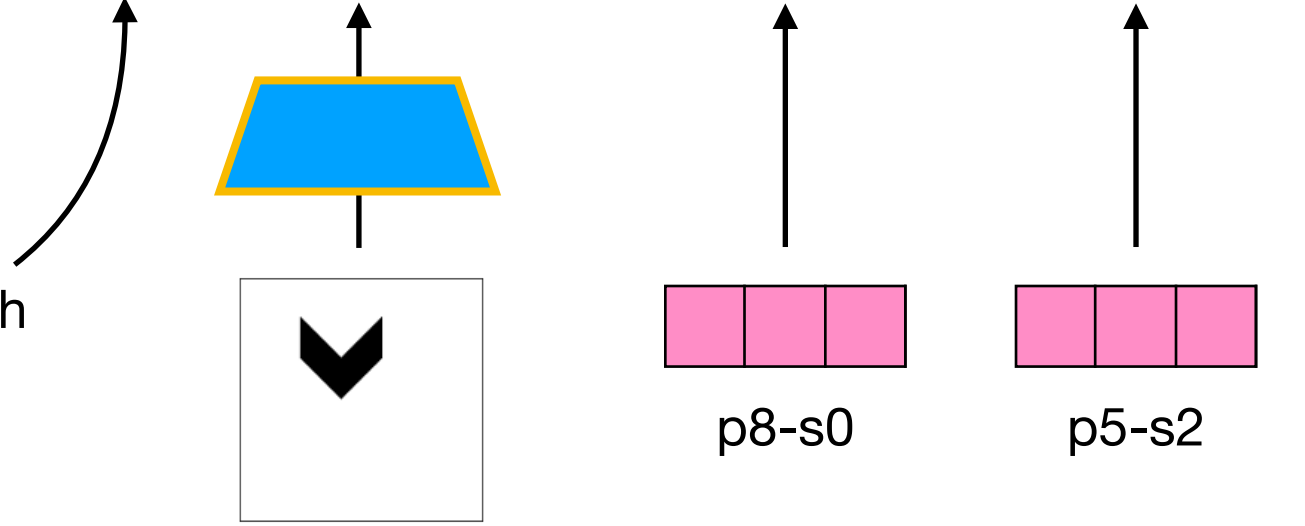
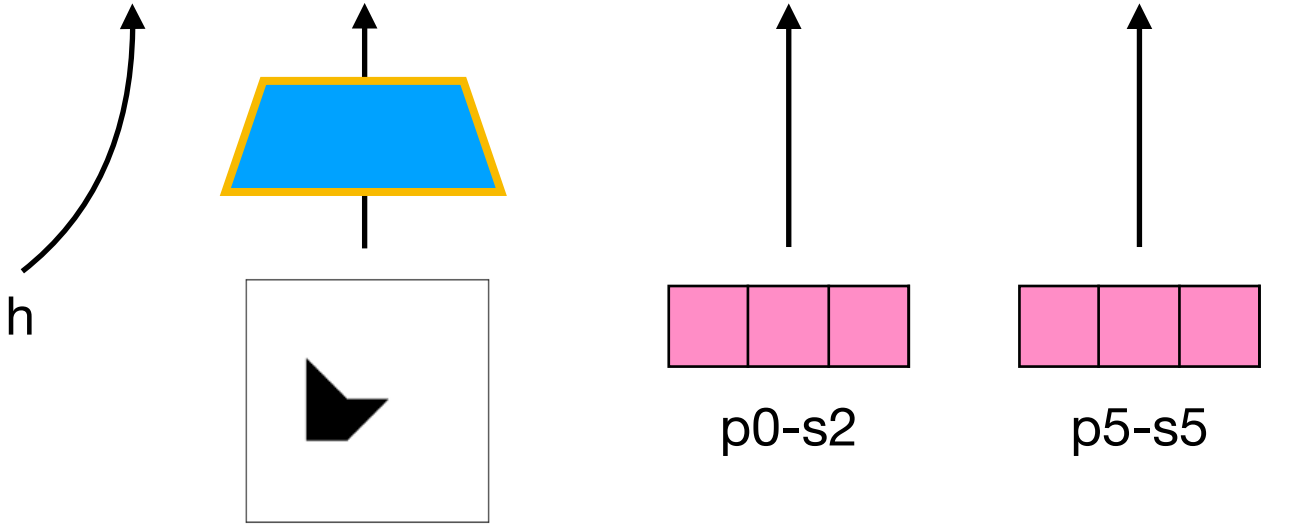
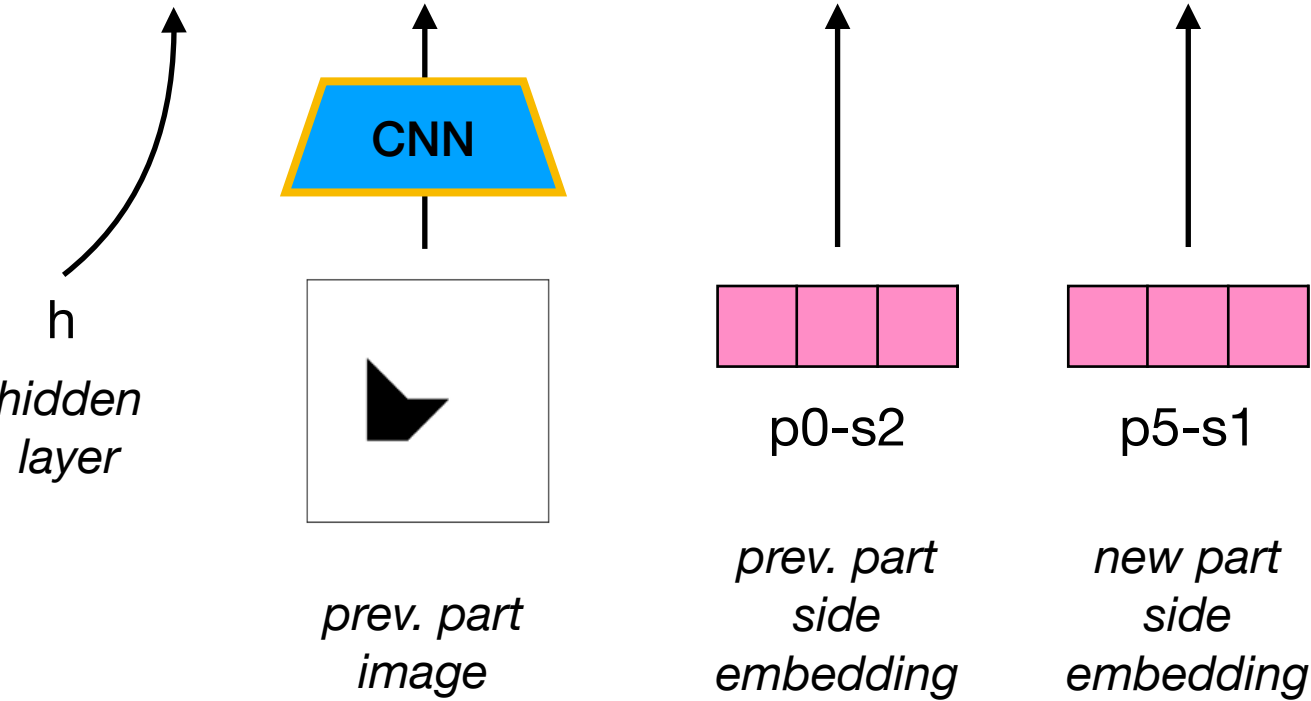
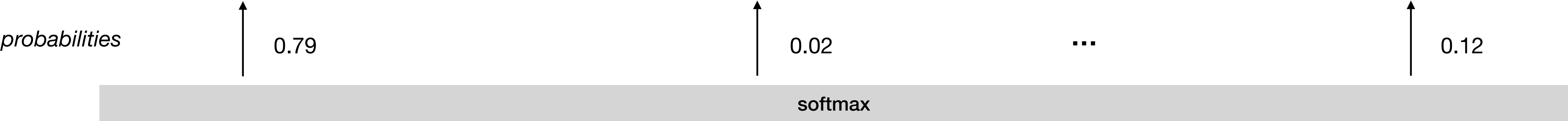


# Forward model in action

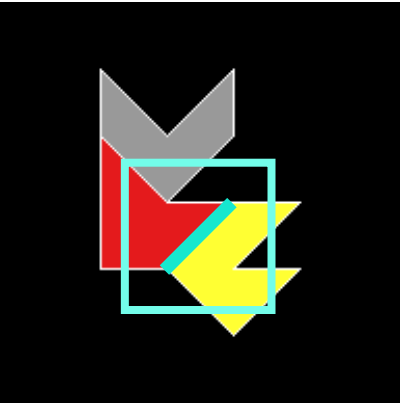




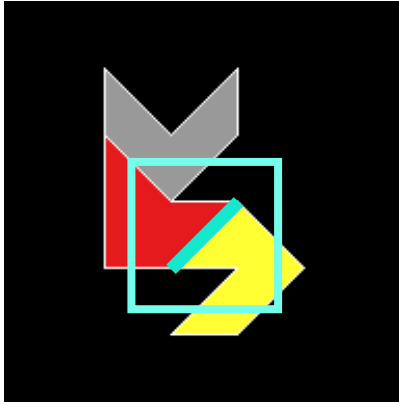
$\text{GENERATE\_RELATION}(x, C, c_i) \rightarrow r_i$



option 1

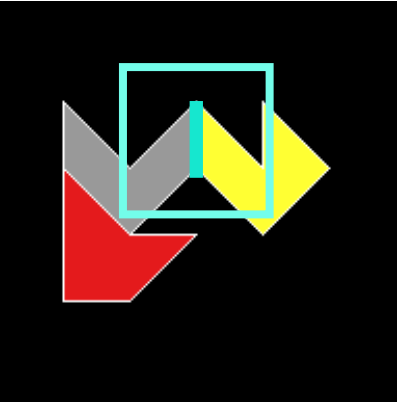


option 2



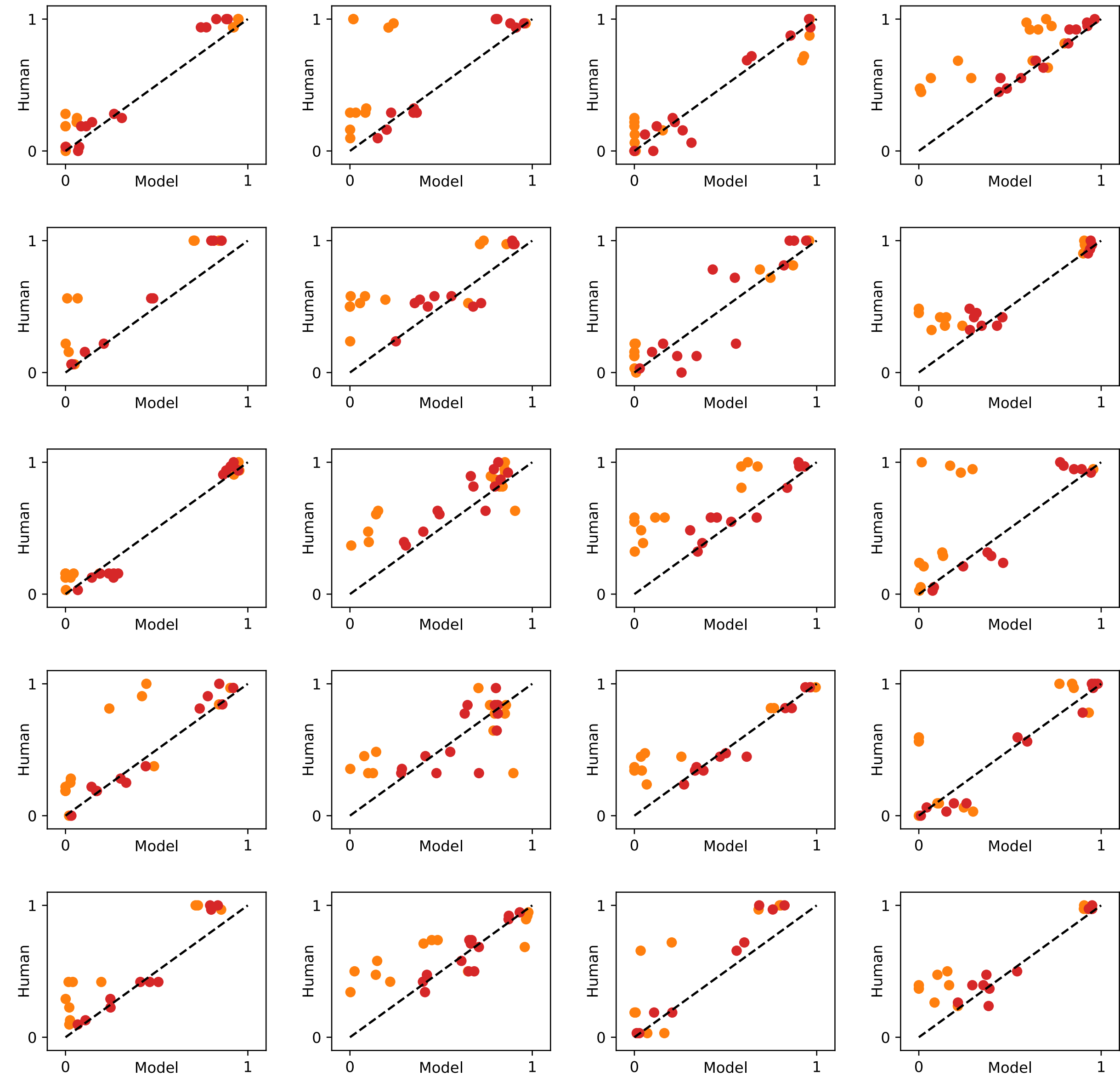
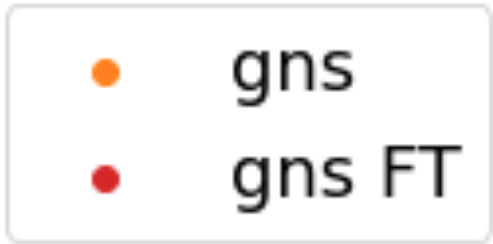
...

option N





# Alien figures: categorization task



GNS: Best-performing GNS model from the generation task (experiment 1), evaluated without any modification

GNS FT: A *finetuned* variant of the GNS model from generation. The model is initialized with the generation parameters and further optimized using (a subset of) human categorization data

	Pearson r	Spearman r
GNS	0.761	0.637
GNS FT	0.953	0.881

**Correlation with human judgements.** Correlation coefficients are computed for each concept type, and the average coefficient across types is reported.