

Structure and Emergence in Human Concepts

Reuben Feinman

advised by
Brenden Lake

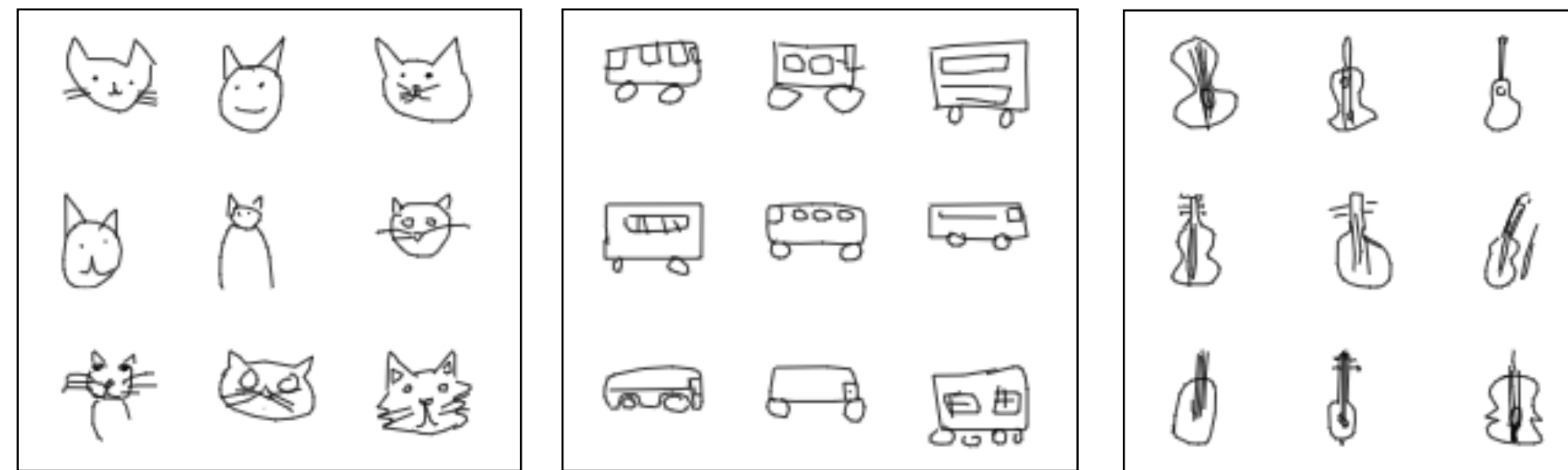
Human Concepts

Recognition



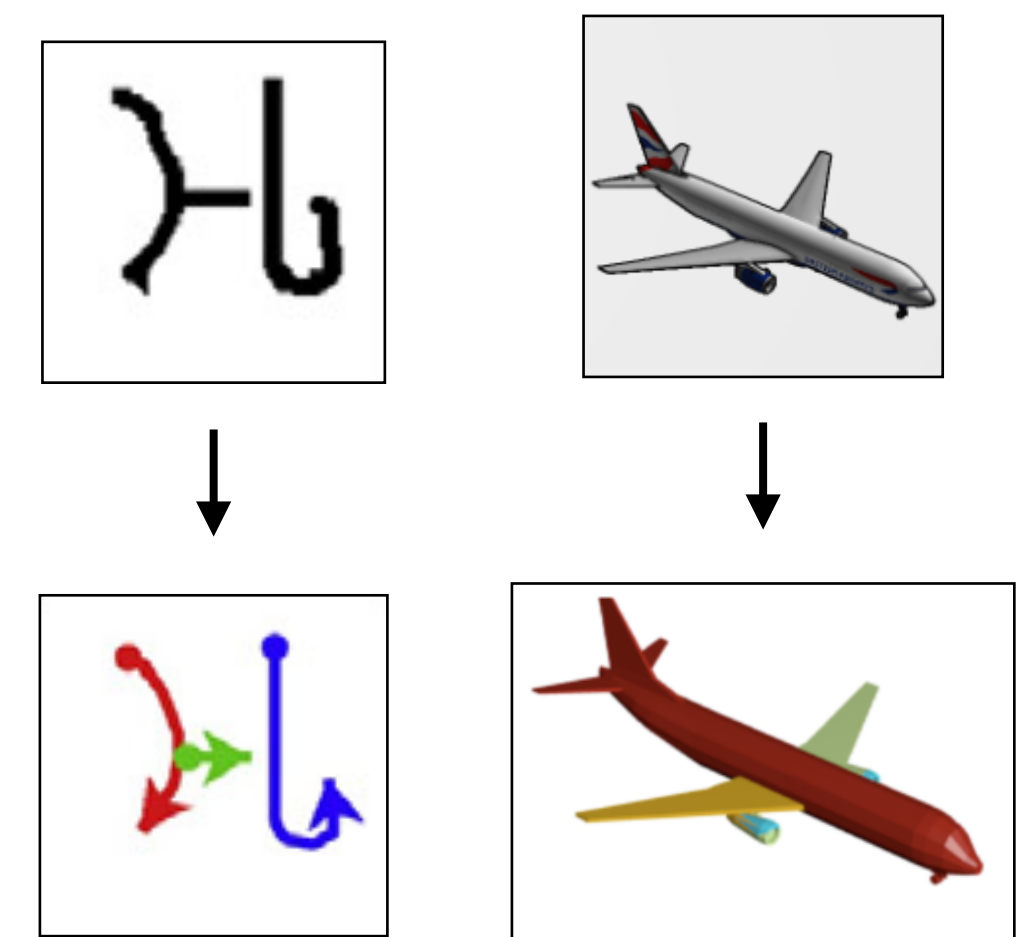
(Lake et al., 2015)

Generation



(Jongejan et al., 2016)

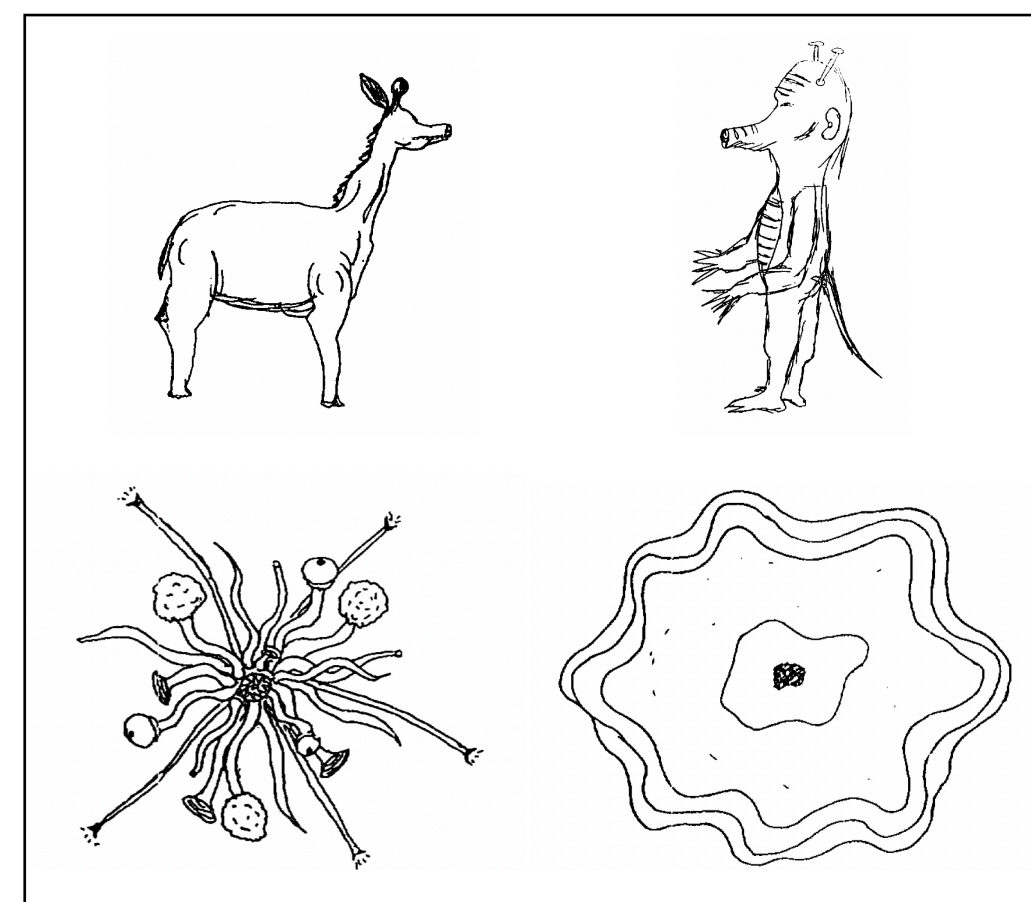
Parsing



(Lake et al., 2015)

(Mo et al., 2018)

Imagination



(Ward, 1994)

Few-shot learning

This is a
“breakfast machine.”



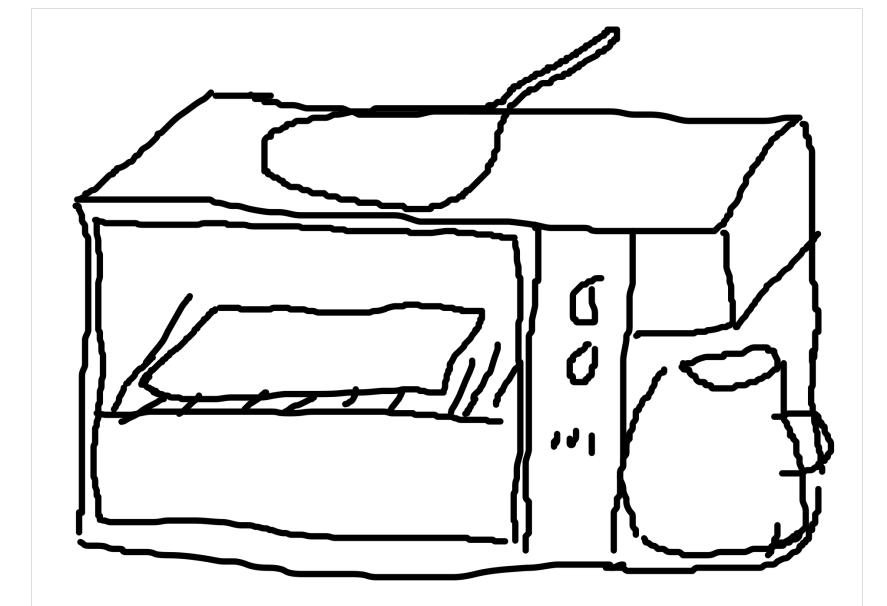
Which is another?



What are its parts?



Create a new one.



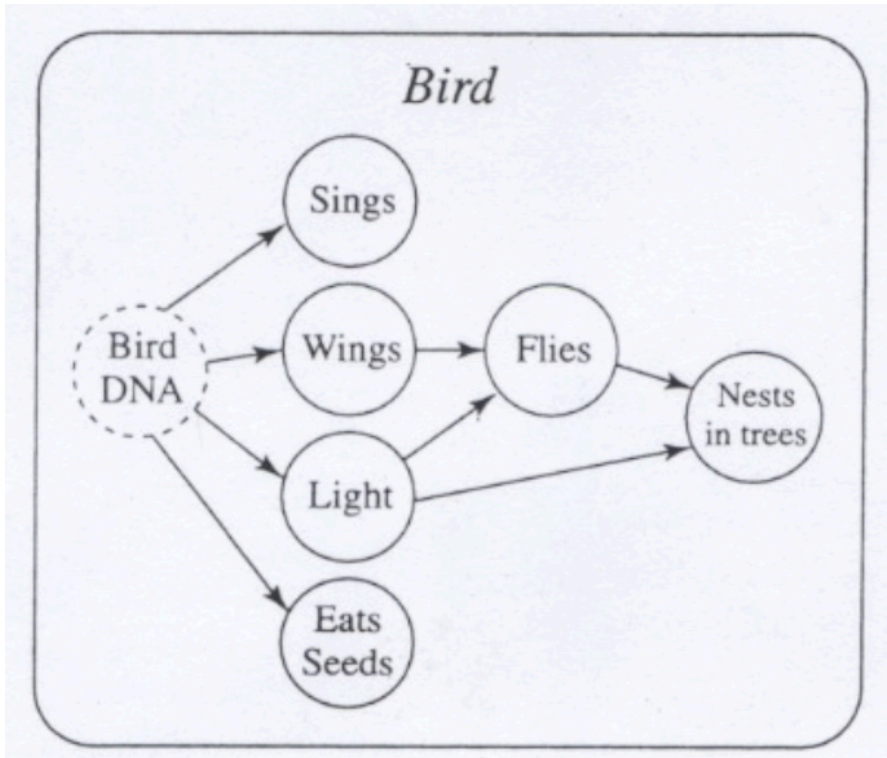
Research Questions

- 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

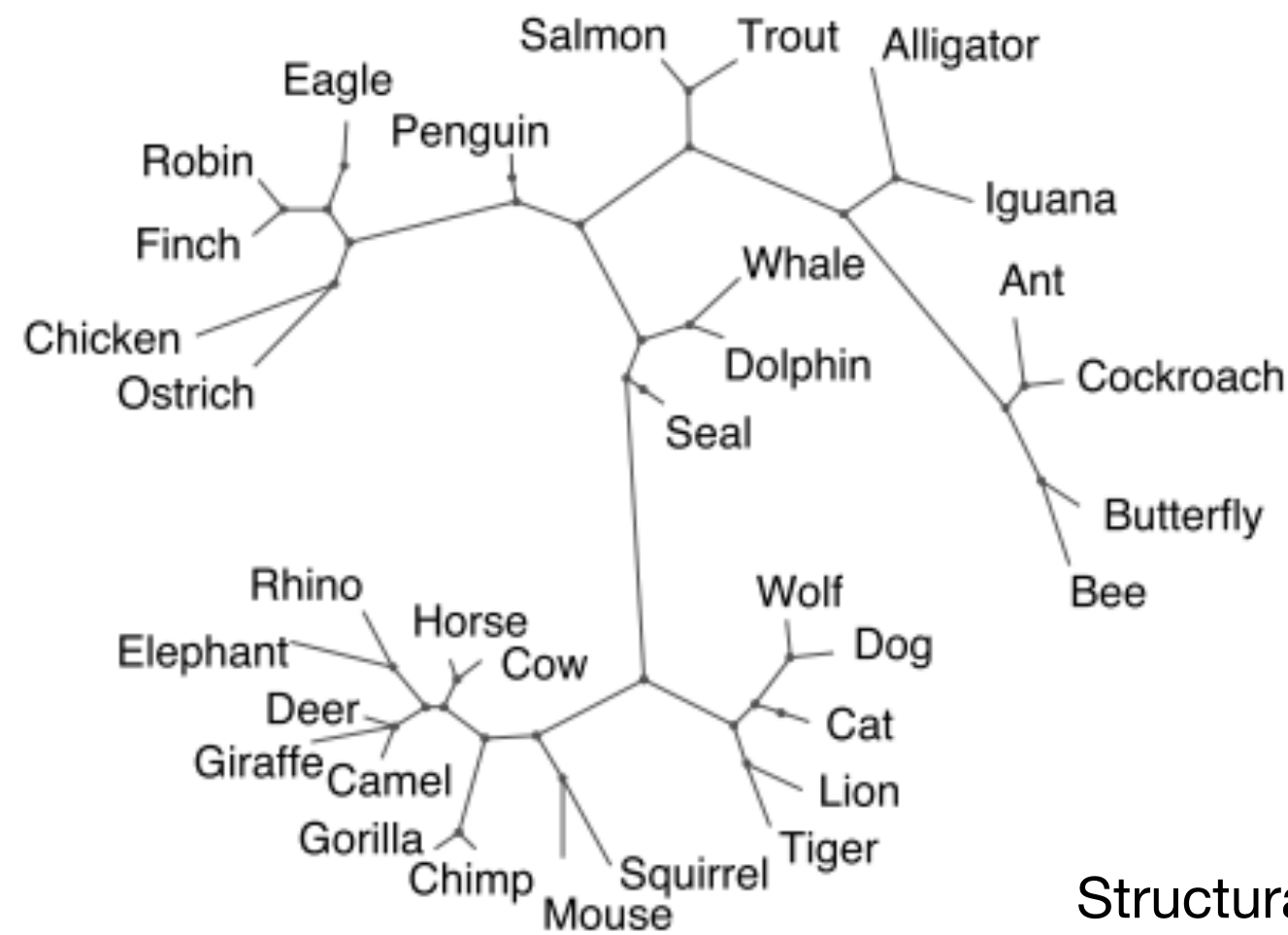
Causal-Model Theory
(Rehder, 2007)



Bayesian Program Learning
(Lake et al., 2015)

```

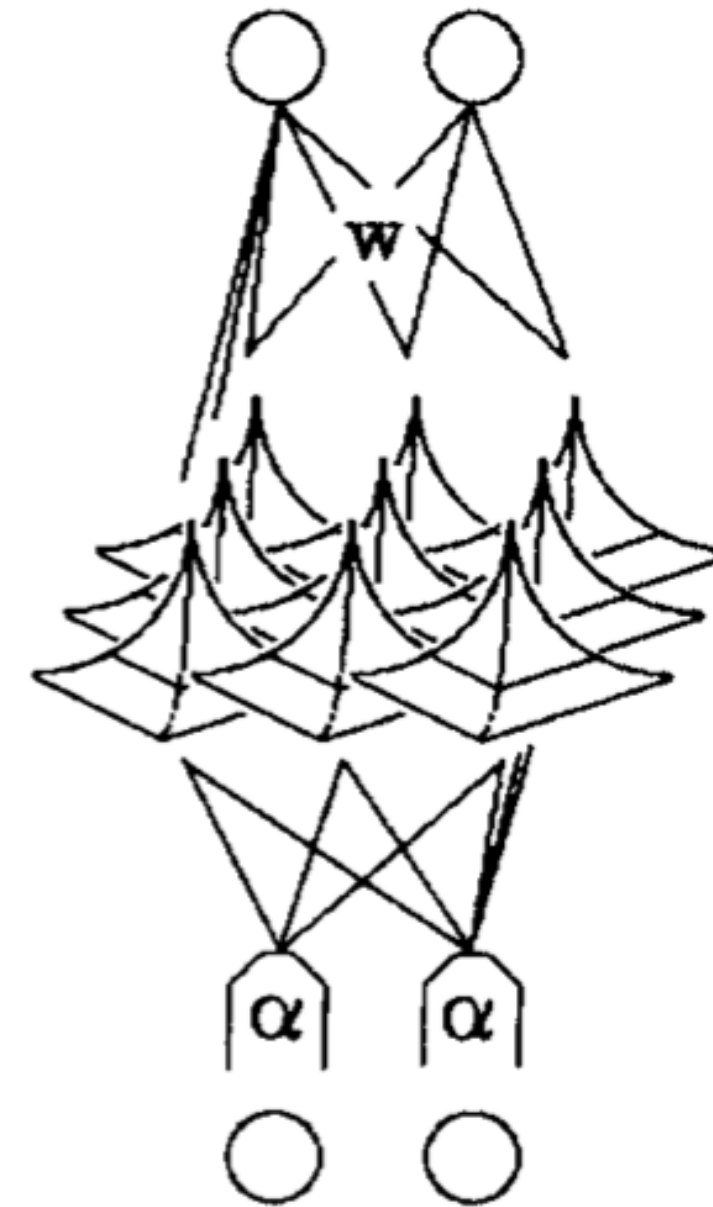
procedure GENERATETYPE
   $\kappa \leftarrow P(\kappa)$  ▷ Sample number of parts
  for  $i = 1 \dots \kappa$  do
     $n_i \leftarrow P(n_i|\kappa)$  ▷ Sample number of sub-parts
    for  $j = 1 \dots n_i$  do
       $s_{ij} \leftarrow P(s_{ij}|s_{i(j-1)})$  ▷ Sample sub-part sequence
    end for
 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



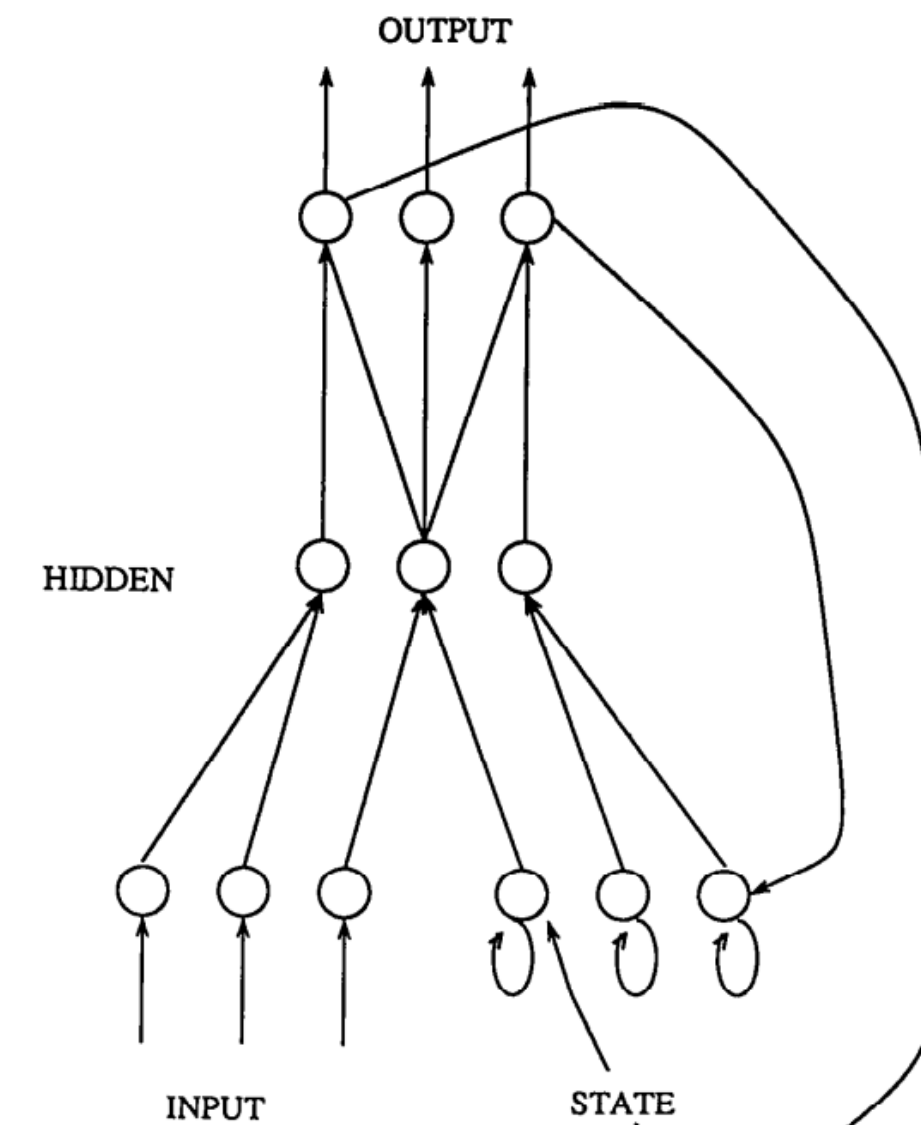
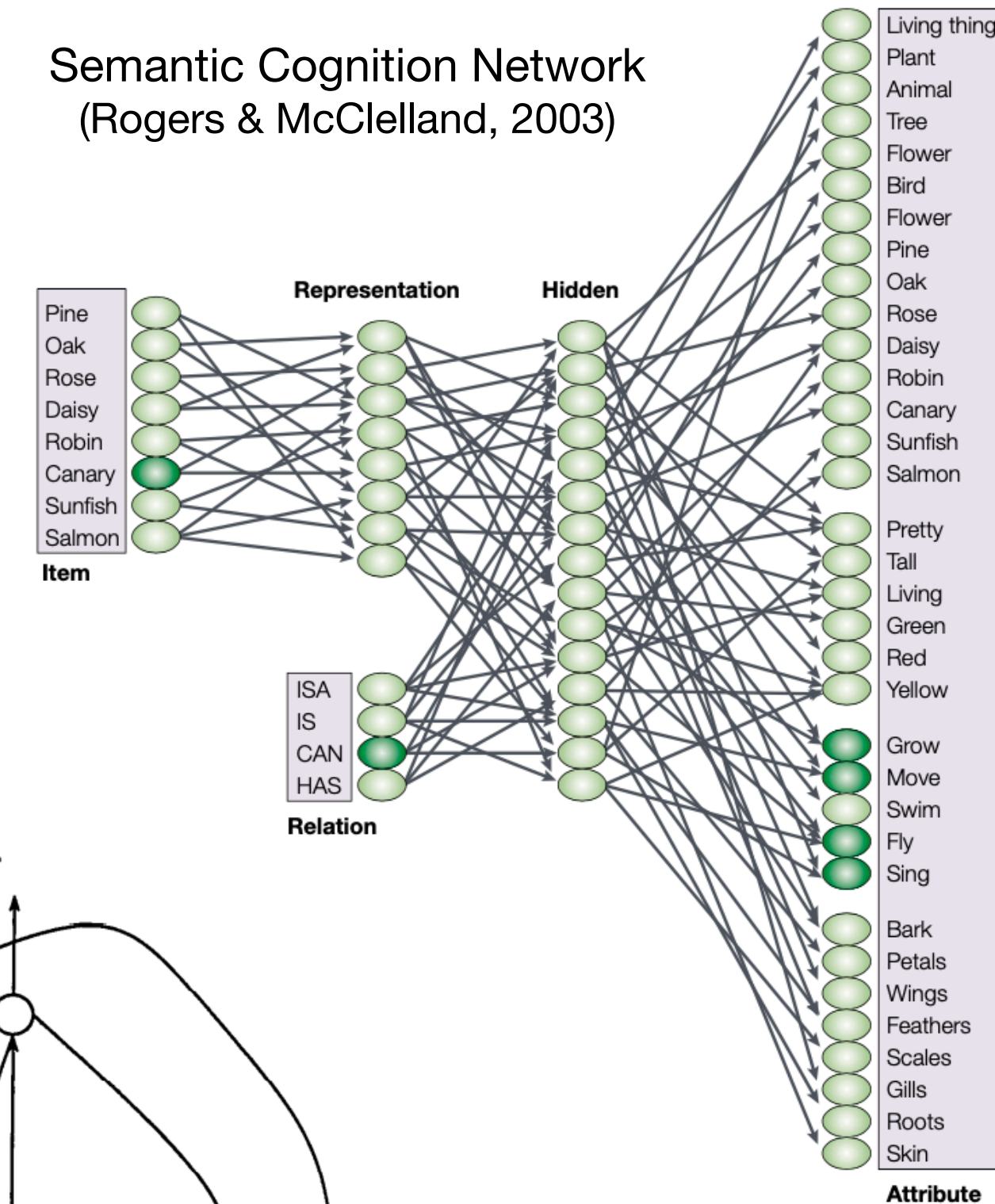
Structural Forms
(Kemp & Tenenbaum, 2008)

Tradition 2: emergent "statistical" knowledge

ALCOVE
(Kruschke, 1992)



Semantic Cognition Network
(Rogers & McClelland, 2003)



Finding Structure in Time
(Elman, 1990)

Synthesis?

Proposal:

Generative Neuro-Symbolic (GNS) Modeling

- Goal: model the *compositional* and *causal* structure in how concepts are formed, while simultaneously modeling nonparametric statistical relationships
- Proposal: probabilistic programs with neural network sub-routines

- probabilistic program representation facilitates explicit causal, compositional structure
- individual parts, and correlations between parts, are represented implicitly by neural networks

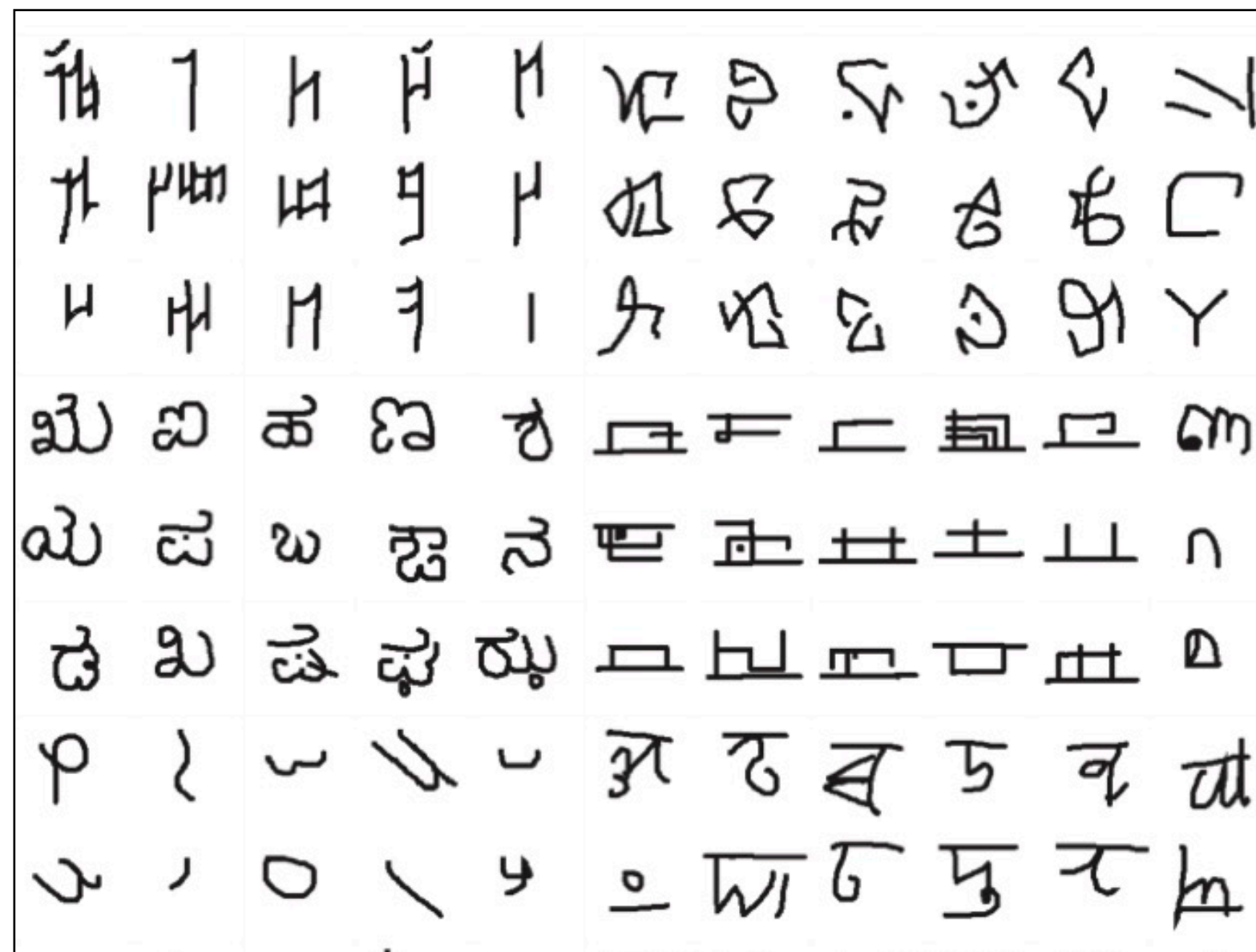
procedure GENERATECONCEPT

```
 $M \leftarrow 0$  ▷ Initialize memory state  
while True do  
   $x_i, r_i \sim p(x, r \mid M)$  ▷ Sample part and relation from neural net  
   $M \leftarrow f_{render}(x_i, r_i, M)$  ▷ Render part to memory (differentiable)  
   $v_i \sim p(v \mid M)$  ▷ Sample termination indicator  
  if  $v_i$  then  
    break  
return  $\{X, R\}$  ▷ Return concept type
```

GNS program to generate a concept "type," a prototype for a new conceptual class

**Case study:
handwritten characters**

Omniglot dataset

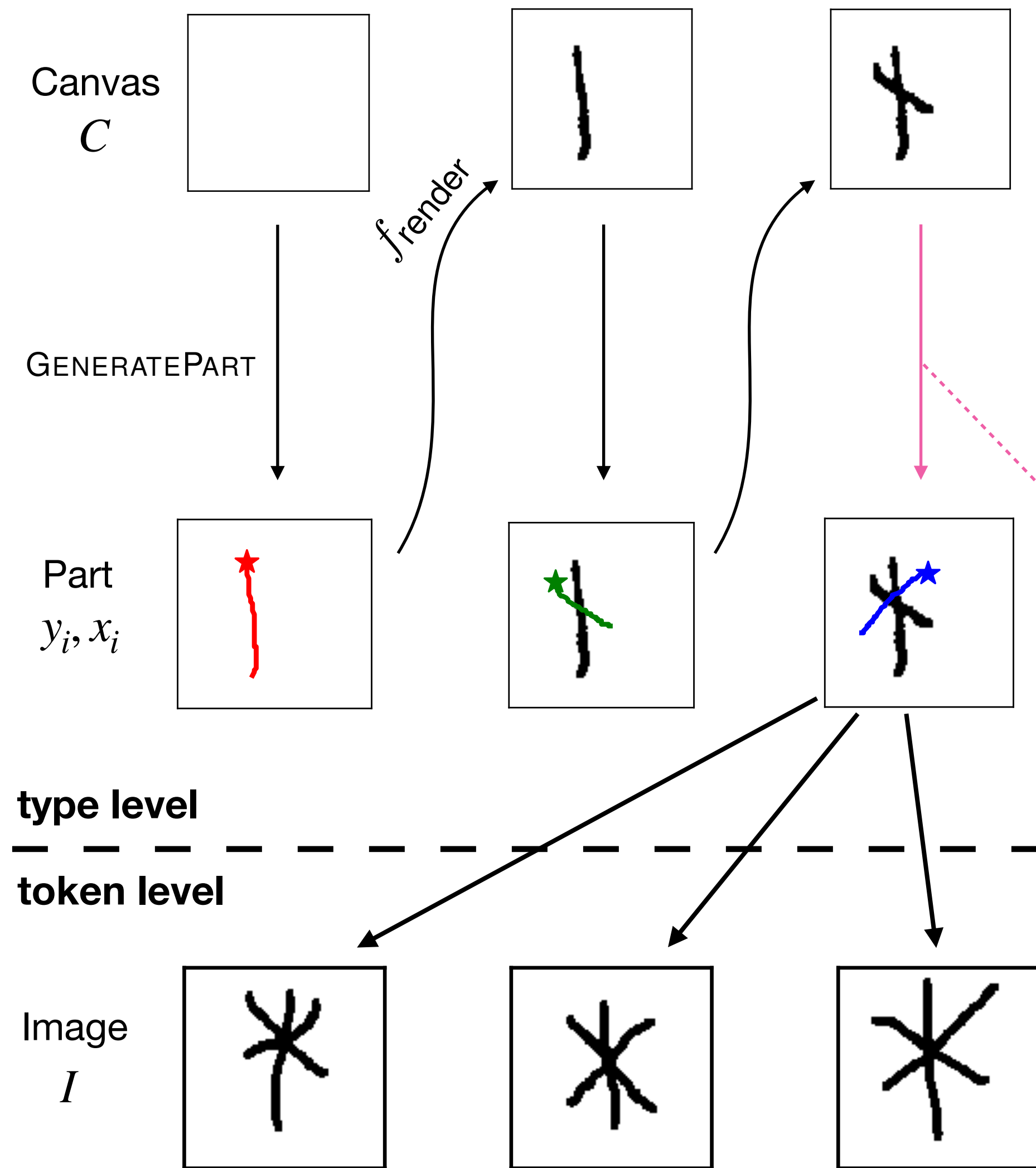


(Lake et al., 2015)

Objectives:

1. Train a GNS model to learn background knowledge of characters using a "background set" of character classes as proxy for human experience
2. Evaluate the model in a series of few-shot concept learning tasks with novel, unseen character classes (from new alphabets) and compare to human behaviors

GNS model of character concepts



procedure GENERATE TYPE

$C \leftarrow 0$

while *true* **do**

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

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

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

if v_i **then**

break

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

return ψ

▷ Initialize blank image canvas

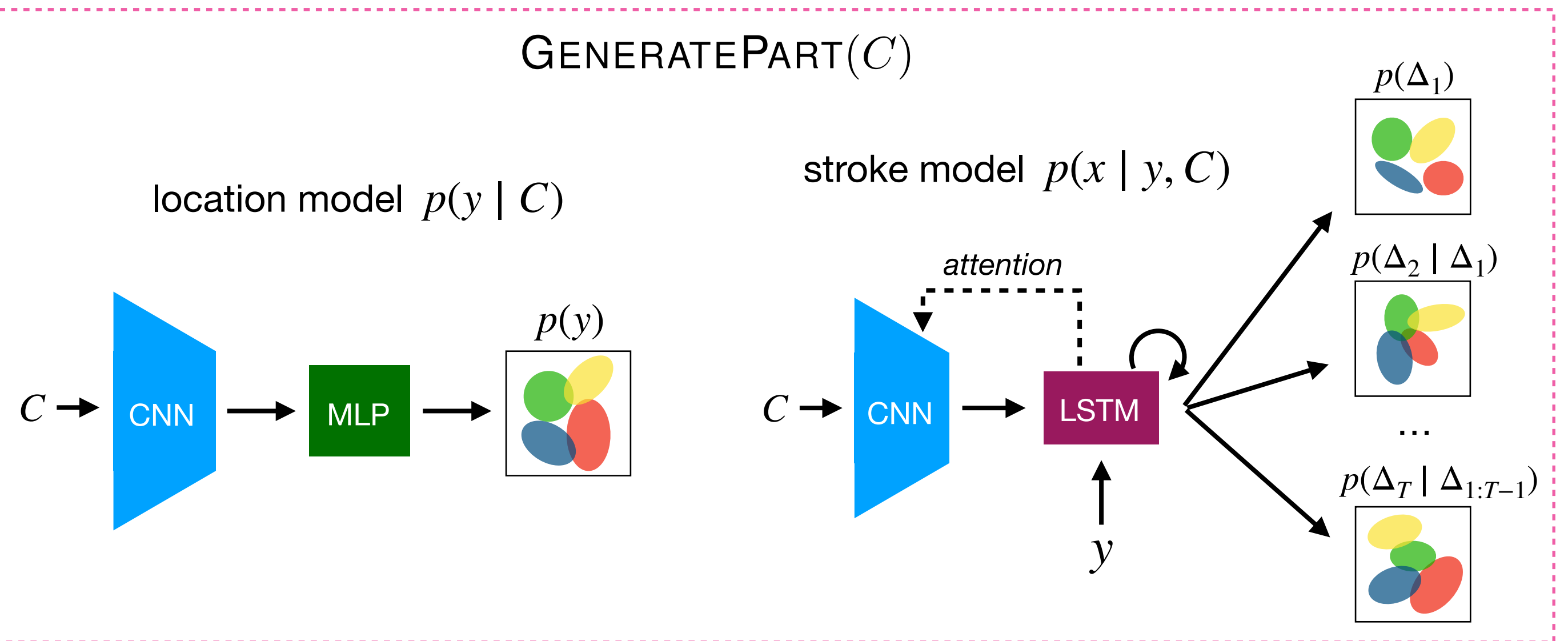
▷ Sample part location & parameters

▷ Render part to image canvas

▷ Sample termination indicator

▷ Terminate sample

▷ Return concept type



Type prior

1. Evaluations on held-out concepts

Test losses

Neuro-Sym	19.51
H-LSTM	20.16
Baseline	19.66

Replicates across different train/test splits

2. Generating new concepts

Humans

5.	l	す	y	c	才
⊙	:	ち	o	m	ま
λ	l	π	ω	γ	l
ψ	γ	ま	r	ね	⊙
⊙	o	か	μ	H	H
π	μ	⊙	H	μ	⊙.

GNS model

≡	Z	Δ	α	γ	ρ
⊙	P	↑	⊕	▷	⊙
⊙	H	⊕	⊕	⊕	4
⊙	H	3	⊕	⊕	▷
R	⊙	F	⊕	⊕	⊕
7	3l	≡	⊙	⊕	⊙

fully-symbolic model (BPL)

π	γ	⊕	⊙	⊕	⊙
⊕	⊕	⊕	⊕	⊕	⊕
≡	2	U	⊕	⊕	⊕
⊕	⊕	⊕	⊕	⊕	⊕
⊕	⊕	⊕	⊕	⊕	⊕
⊕	⊕	⊕	⊕	⊕	⊕
⊕	⊕	⊕	⊕	⊕	⊕

Concept learning tasks

(Lake et al., 2015)

(a) One-shot classification

Where is another?

ಗ	೧	೨	೩	೪
೫	೬	೭	೮	೯
೧೦	೧೧	೧೨	೧೩	೧೪
೧೫	೧೬	೧೭	೧೮	೧೯

Where is another?

೨೦	೨೧	೨೨	೨೩	೨೪
೨೫	೨೬	೨೭	೨೮	೨೯
೩೦	೩೧	೩೨	೩೩	೩೪
೩೫	೩೬	೩೭	೩೮	೩೯

(b) Parsing

೪	೫	೬
೭	೮	೯
೧೦	೧೧	೧೨

Human parses

೪	೫	೬
೭	೮	೯
೧೦	೧೧	೧೨

Machine parses

೪	೫	೬
೭	೮	೯
೧೦	೧೧	೧೨

stroke order 1 2 3 4 5

(c) Generating new exemplars

Human or Machine?

೧	೨	೩
೪	೫	೬
೭	೮	೯

೧	೨	೩
೪	೫	೬
೭	೮	೯

(d) Generating new concepts (from type)

Alphabet

೧	೨	೩	೪	೫
೬	೭	೮	೯	೧೦

Human or Machine?

೧	೨
೩	೪

೫	೬
೭	೮

(e) Generating new concepts (unconstrained)

Human or Machine?


೧	೨
೩	೪

೫	೬
೭	೮

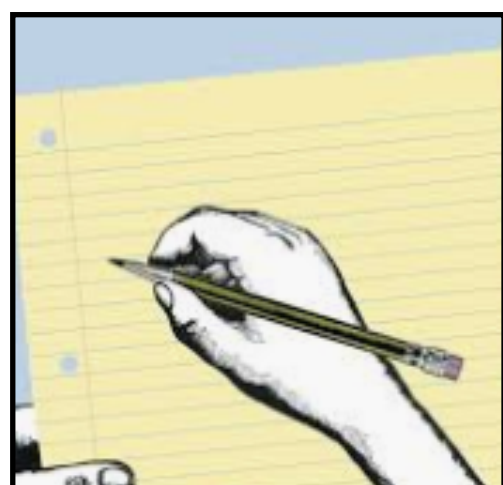
೧	೨
೩	೪

೫	೬
೭	೮

future work

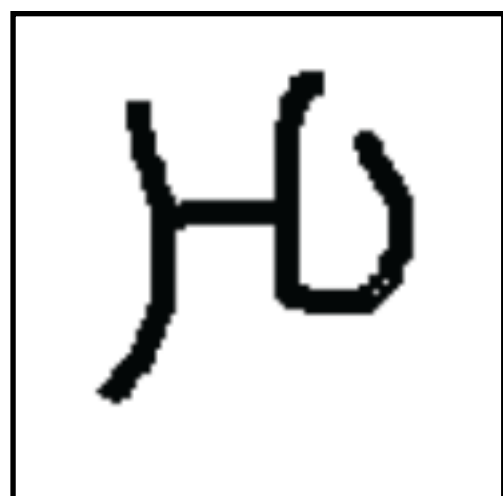
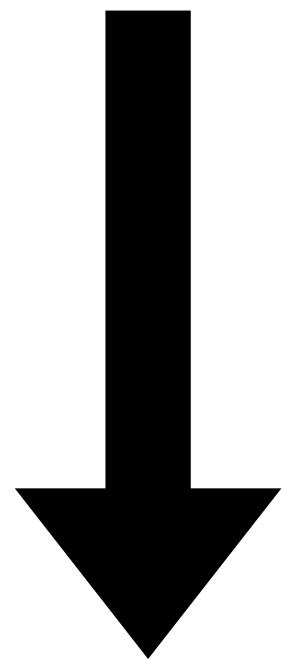


Probabilistic inference



θ

Latent program



I

Image
(observation)

Inference

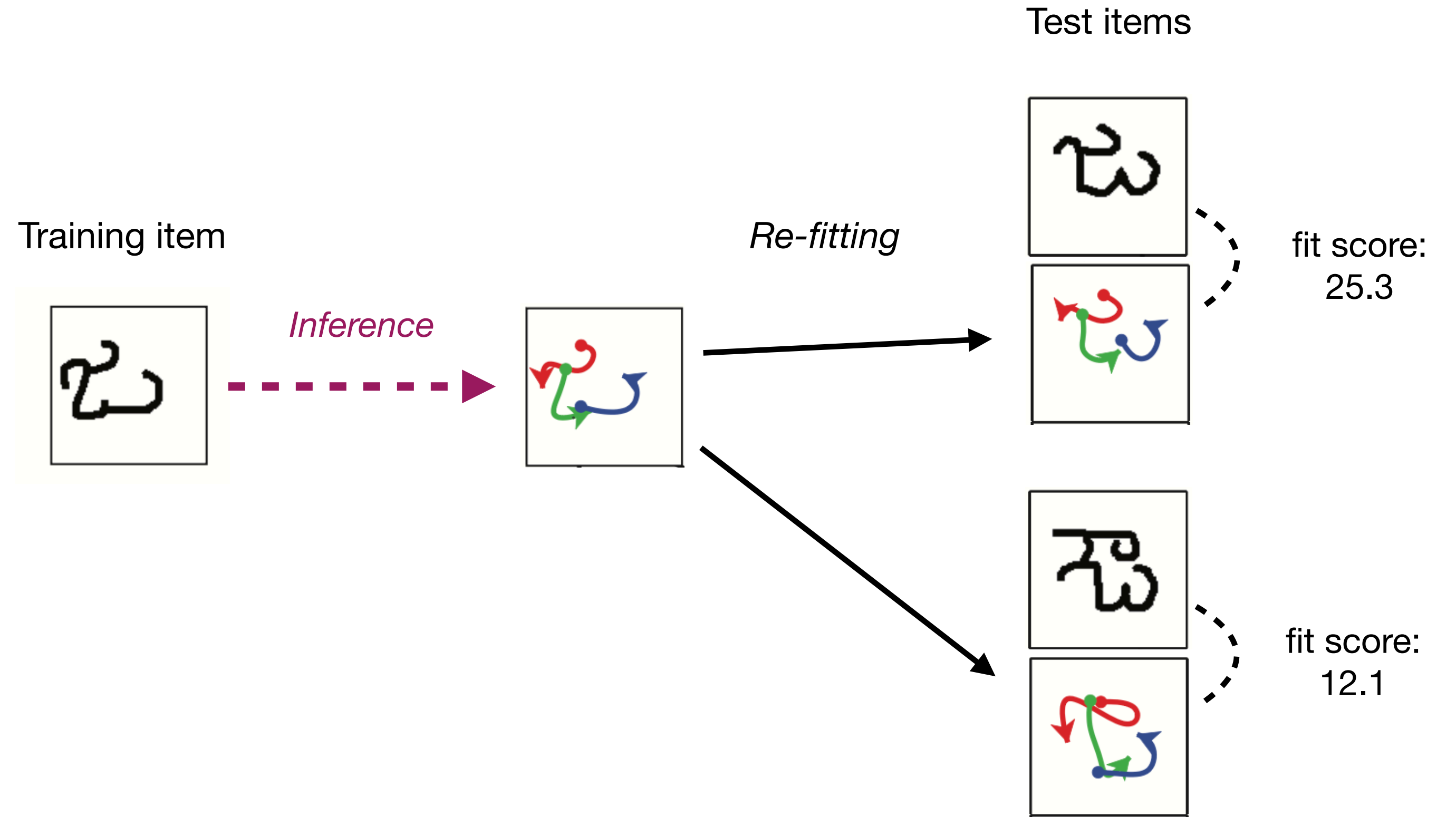
Bayes' rule:

$$P(\theta | I) = \frac{P(I | \theta)P(\theta)}{P(I)}$$

One-Shot Classification



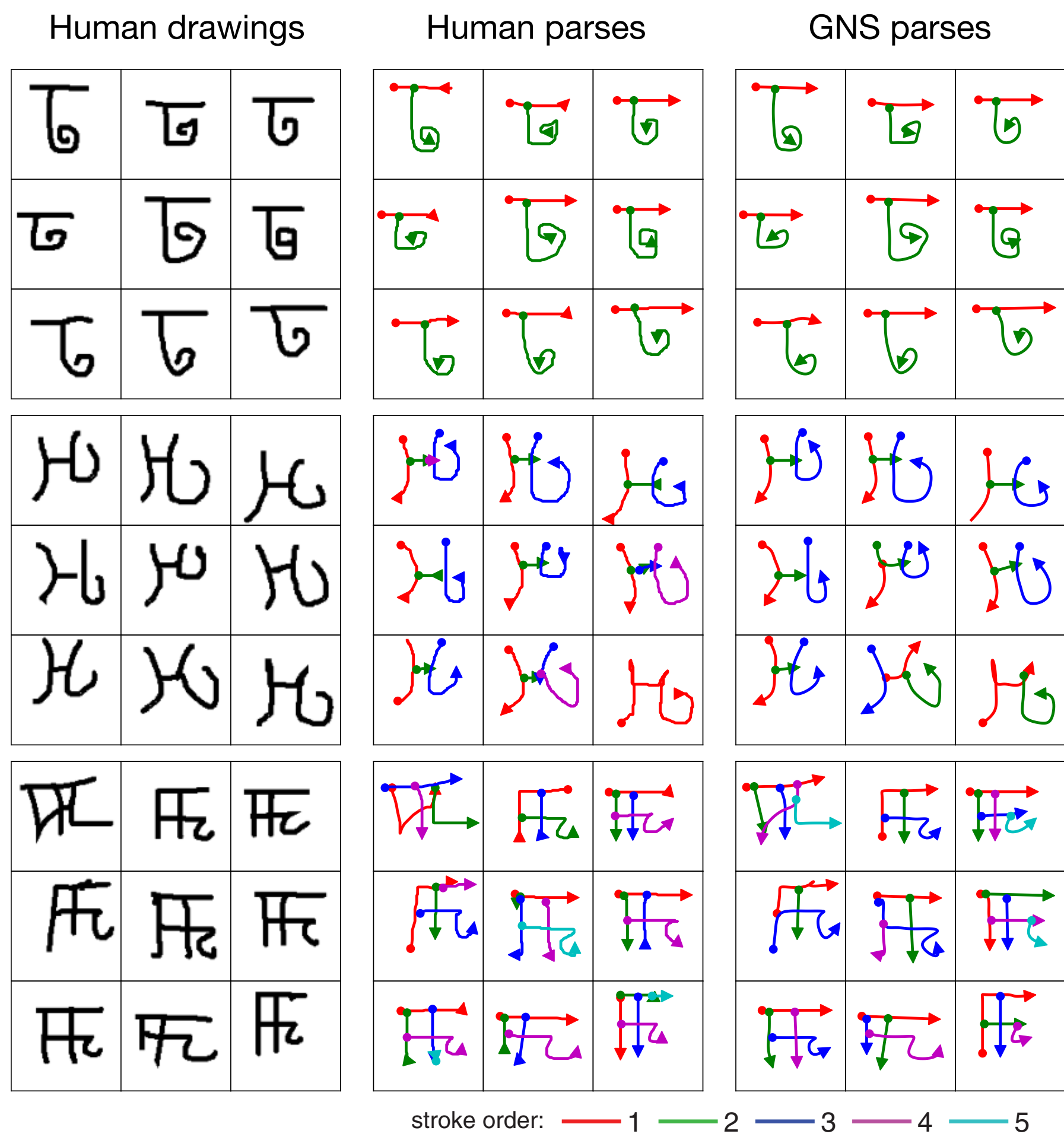
ఓ	ఇ	లు	ఎ	ఔ
ఉ	బి	గా	ఓ	ఝ
షె	త	ణ	త్ర	ద
న	య	ల	రా	ళ



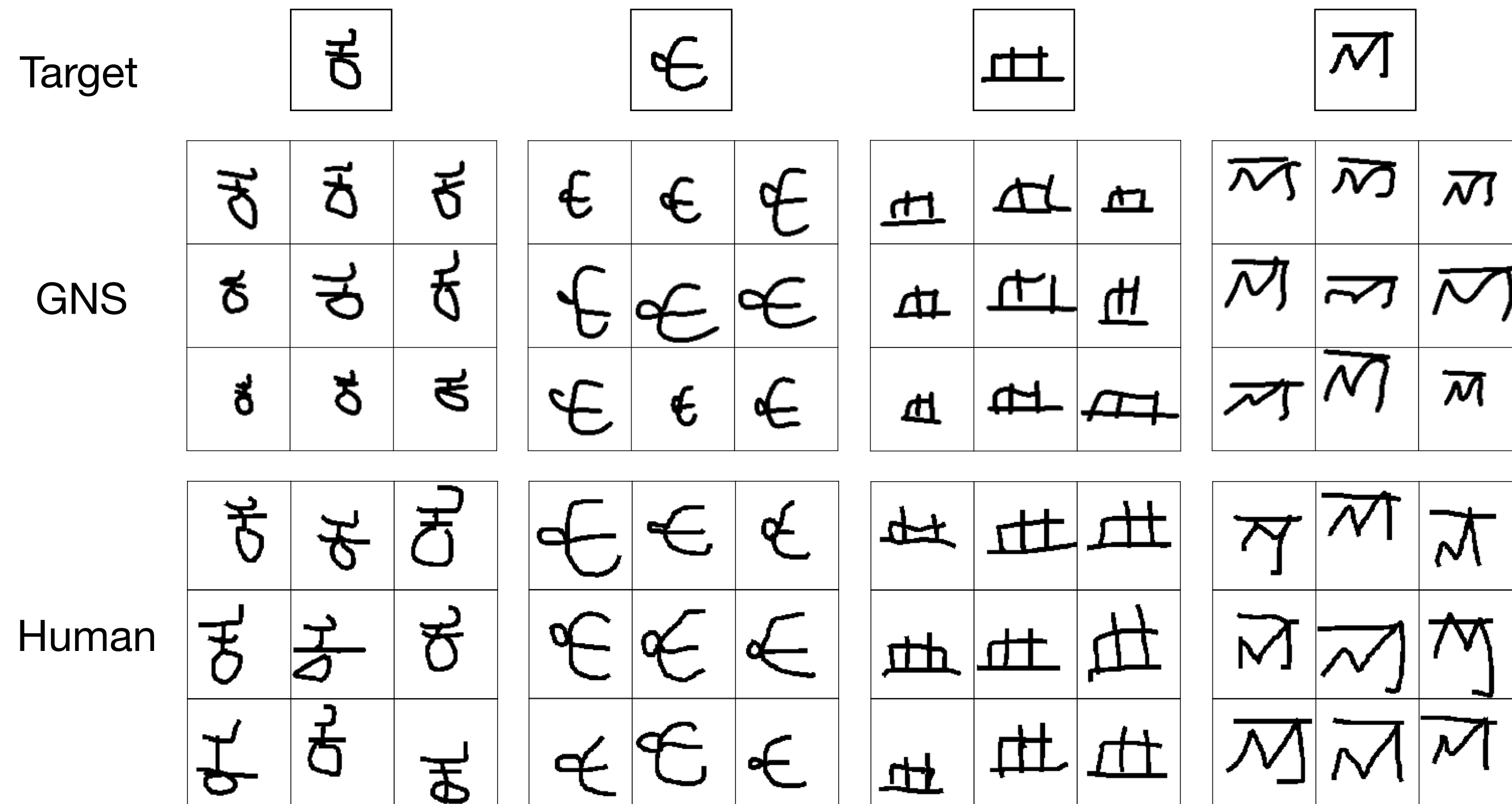
Results (400 trials)

GNS	94.3%
Humans	95.5%

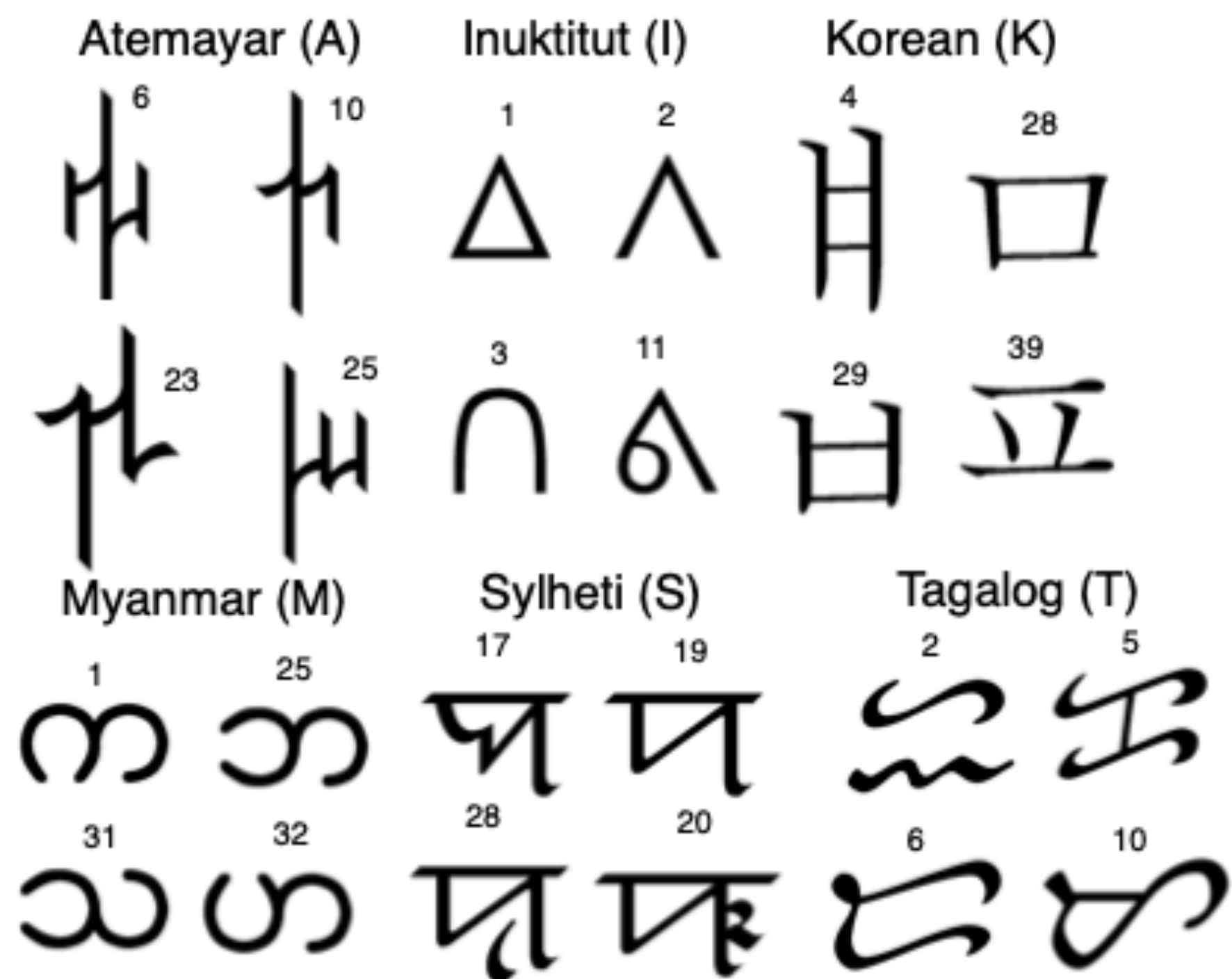
Parsing



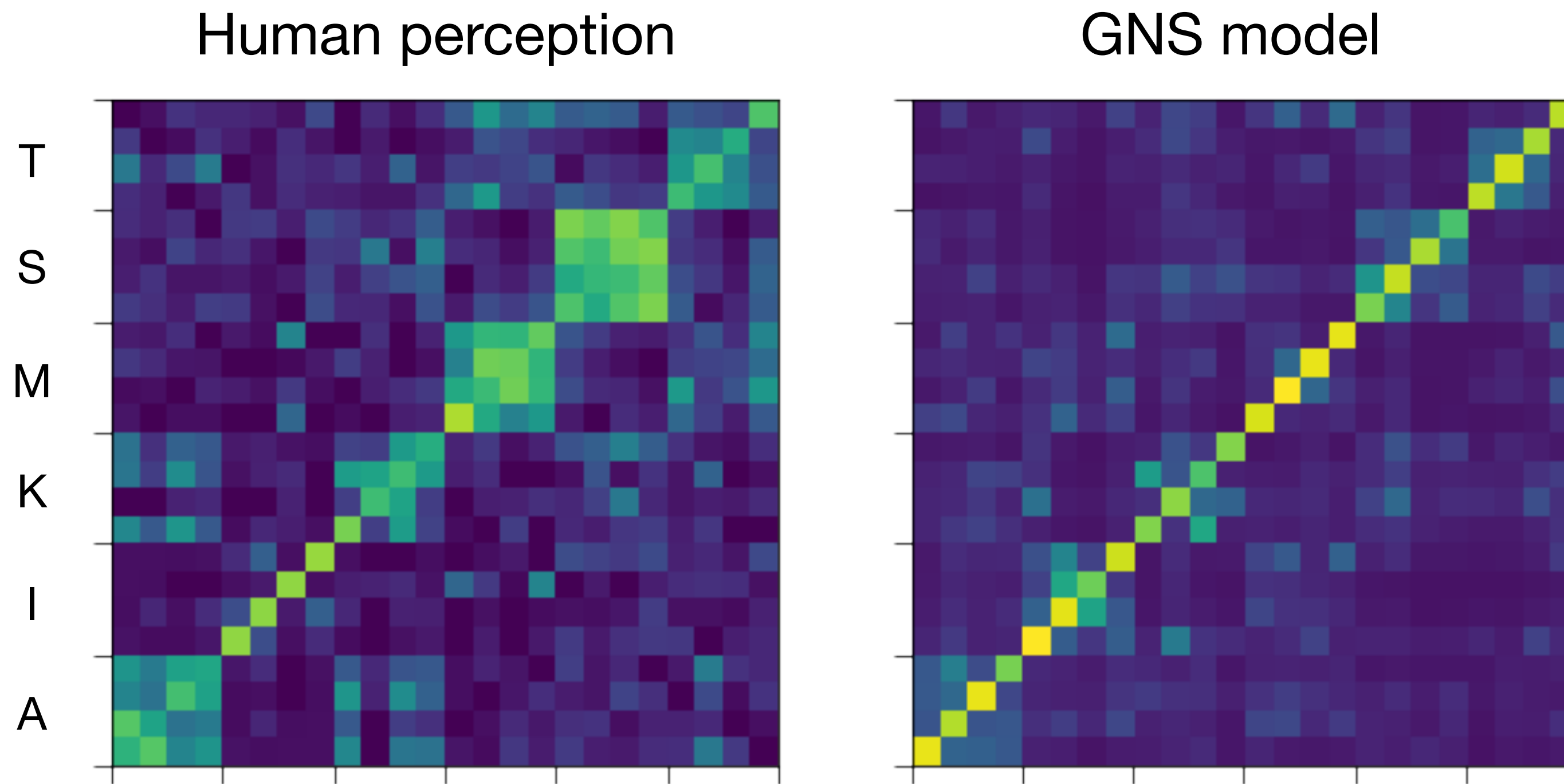
Generating new exemplars



Fit to human perceptual discrimination



(Lake et al., 2011)



$$r(576) = 0.650; \quad p < 0.001$$

Conclusions

- Human concepts go far beyond classification: they enable a variety of discriminative and generative abilities
- Generative neuro-symbolic (GNS) models can capture the dual structural and statistical characteristics of human concepts that enable flexible generalization to a range of tasks
- GNS models offer an account for how previous experience can support the rapid acquisition of new concepts through priors

Thank You

Brenden Lake (NYU)

Joshua Tenenbaum (MIT)

Tuan-Anh Le (MIT)

Maxwell Nye (MIT)

Lucas Tian (Rockefeller U.)

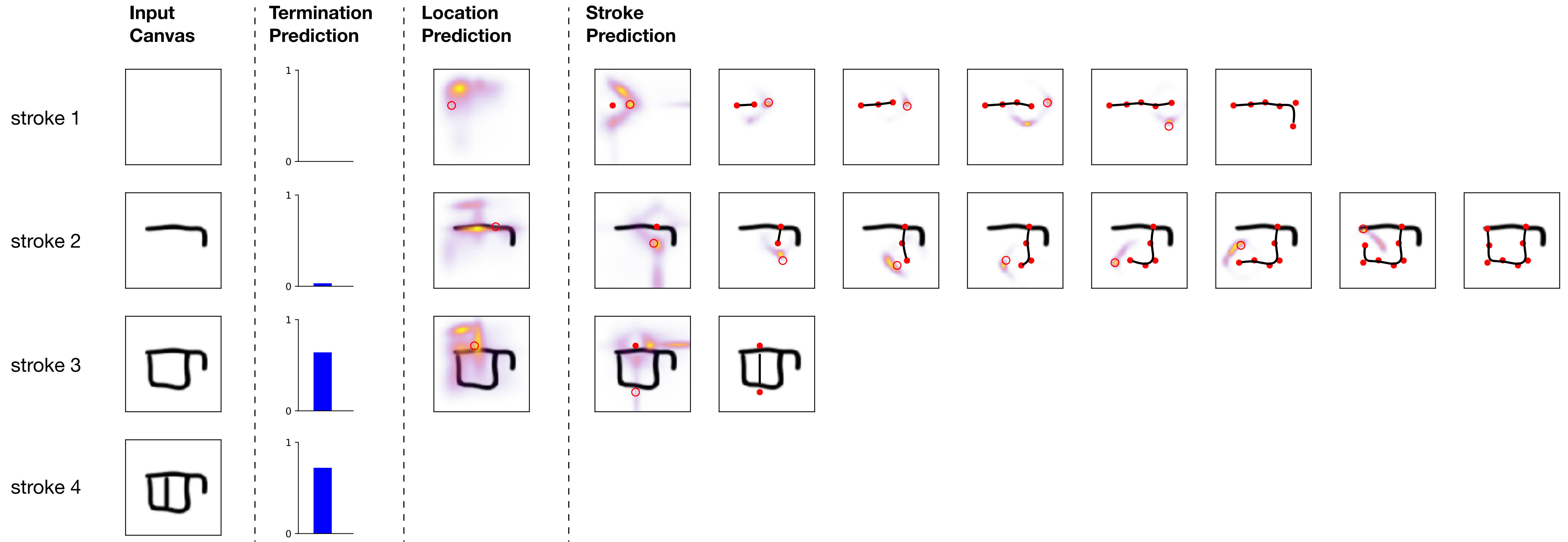
Stéphane Deny (Facebook)

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

—Richard Feynman

Extras

GNS Type Prior



Novelty of character samples

GNS Samples

୧	୨	୩	୪	୫	୬	୭	୮	୯	୧୦
୧୧	୧୨	୧୩	୧୪	୧୫	୧୬	୧୭	୧୮	୧୯	୨୦
୨୧	୨୨	୨୩	୨୪	୨୫	୨୬	୨୭	୨୮	୨୯	୩୦
୩୧	୩୨	୩୩	୩୪	୩୫	୩୬	୩୭	୩୮	୩୯	୪୦
୪୧	୪୨	୪୩	୪୪	୪୫	୪୬	୪୭	୪୮	୪୯	୫୦
୫୧	୫୨	୫୩	୫୪	୫୫	୫୬	୫୭	୫୮	୫୯	୬୦
୬୧	୬୨	୬୩	୬୪	୬୫	୬୬	୬୭	୬୮	୬୯	୭୦
୭୧	୭୨	୭୩	୭୪	୭୫	୭୬	୭୭	୭୮	୭୯	୮୦
୮୧	୮୨	୮୩	୮୪	୮୫	୮୬	୮୭	୮୮	୮୯	୯୦
୯୧	୯୨	୯୩	୯୪	୯୫	୯୬	୯୭	୯୮	୯୯	୧୦୦

Nearest training neighbors

୧	୨	୩	୪	୫	୬	୭	୮	୯	୧୦
୧୧	୧୨	୧୩	୧୪	୧୫	୧୬	୧୭	୧୮	୧୯	୨୦
୨୧	୨୨	୨୩	୨୪	୨୫	୨୬	୨୭	୨୮	୨୯	୩୦
୩୧	୩୨	୩୩	୩୪	୩୫	୩୬	୩୭	୩୮	୩୯	୪୦
୪୧	୪୨	୪୩	୪୪	୪୫	୪୬	୪୭	୪୮	୪୯	୫୦
୫୧	୫୨	୫୩	୫୪	୫୫	୫୬	୫୭	୫୮	୫୯	୬୦
୬୧	୬୨	୬୩	୬୪	୬୫	୬୬	୬୭	୬୮	୬୯	୭୦
୭୧	୭୨	୭୩	୭୪	୭୫	୭୬	୭୭	୭୮	୭୯	୮୦
୮୧	୮୨	୮୩	୮୪	୮୫	୮୬	୮୭	୮୮	୮୯	୯୦
୯୧	୯୨	୯୩	୯୪	୯୫	୯୬	୯୭	୯୮	୯୯	୧୦୦