

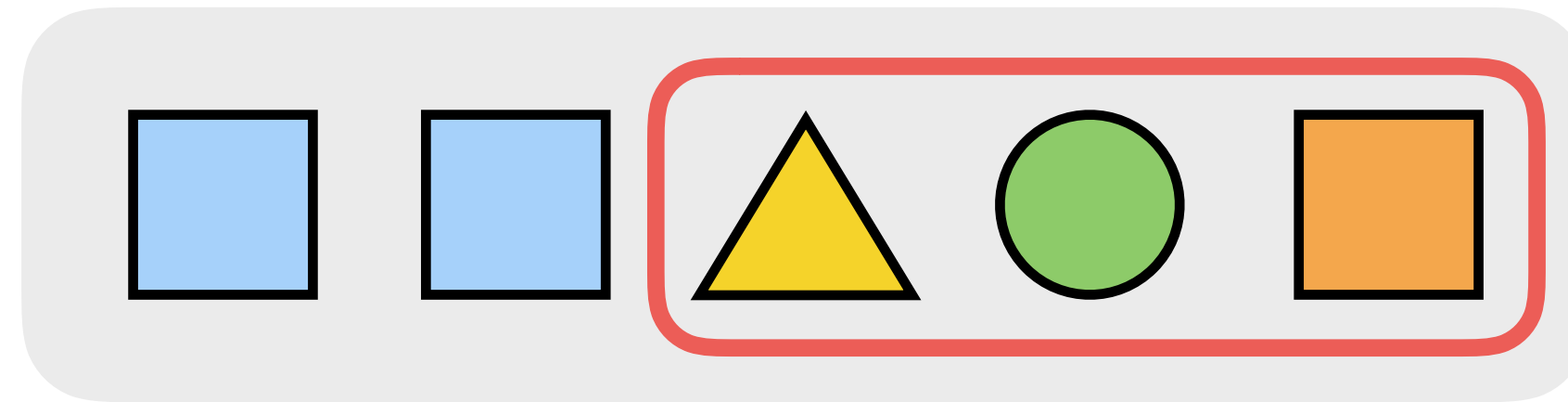
Analogs of Linguistic Structure in Deep Representations



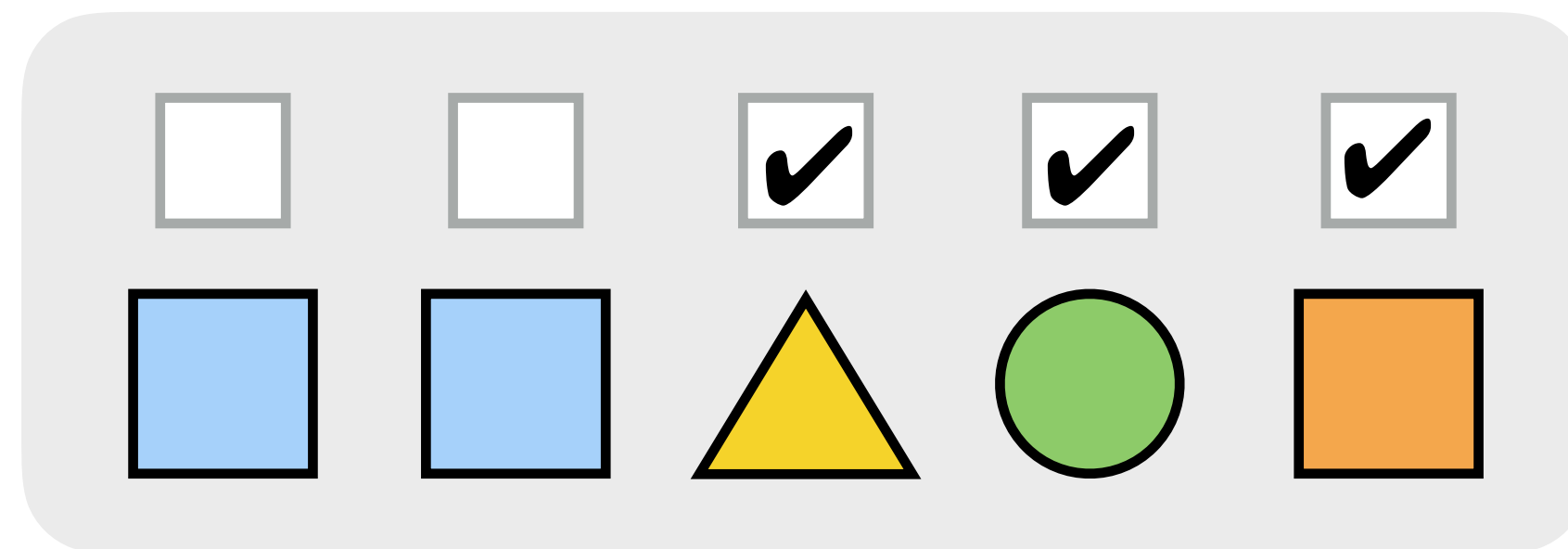
Jacob Andreas and Dan Klein



A game for humans

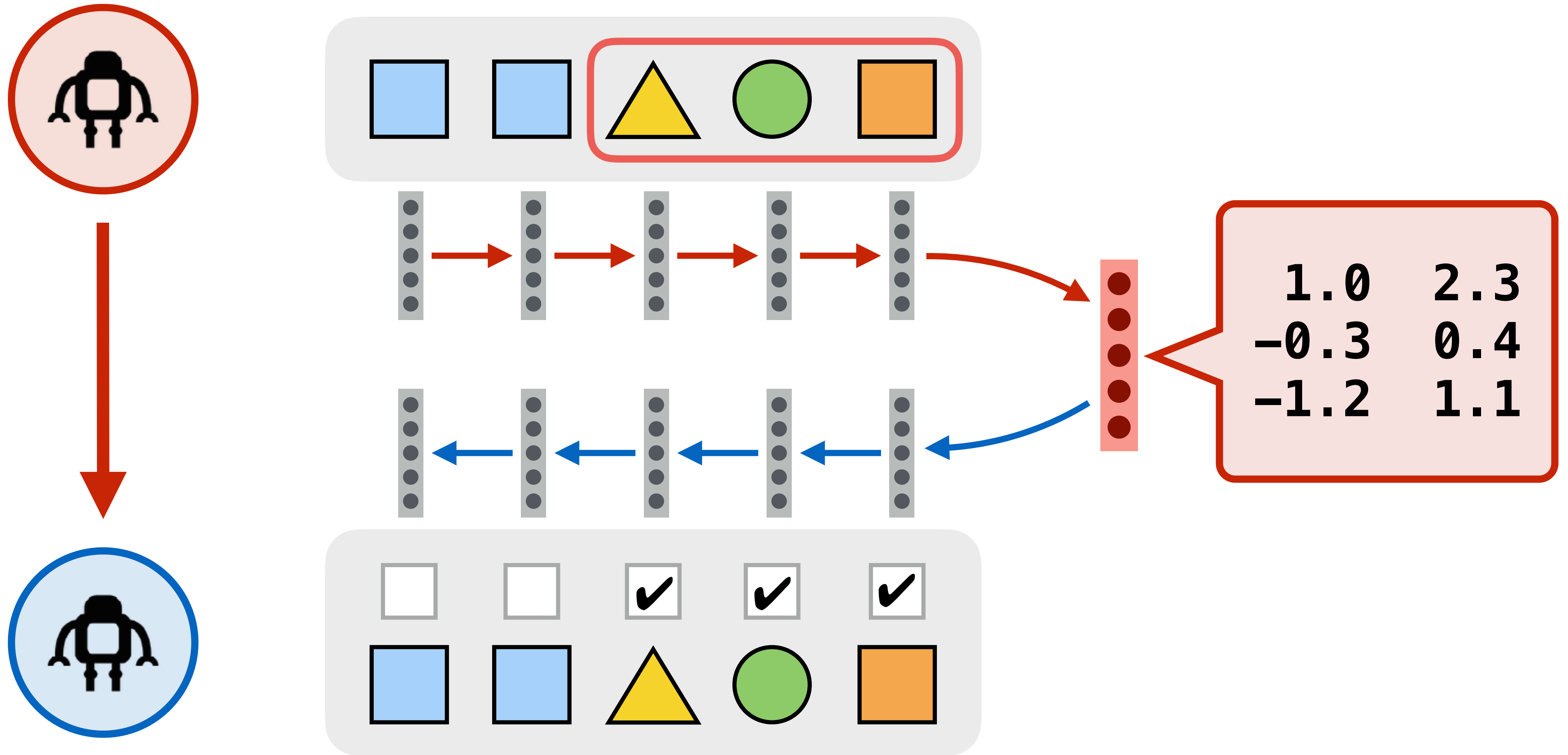


*everything but the blue shapes
orange square and non-squares*





A game for RNNs



[e.g. Lazaridou et al. 2016]



Questions

1. Does the RNN employ a human-like communicative strategy?

*everything but
squares*

?



1.0	2.3
-0.3	0.4
-1.2	1.1



Questions

2. Do RNN representations have interpretable compositional structure?

"not"



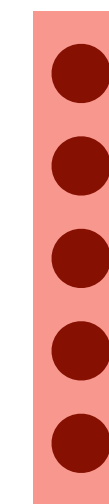
*

"red"



?

≡



1.0	2.3
-0.3	0.4
-1.2	1.1



Computing meaning representations

not the red squares



Computing meaning representations

not the red squares

$\lambda x. \neg(\text{sqr}(x) \wedge \text{red}(x))$



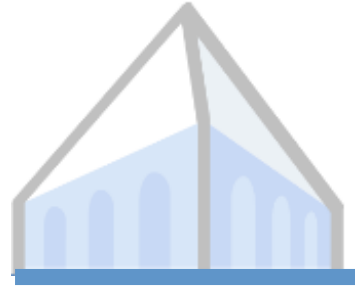
Computing meaning representations

not the red squares

$\lambda x. \neg(\text{sqr}(x) \wedge \text{red}(x))$

not red or not square

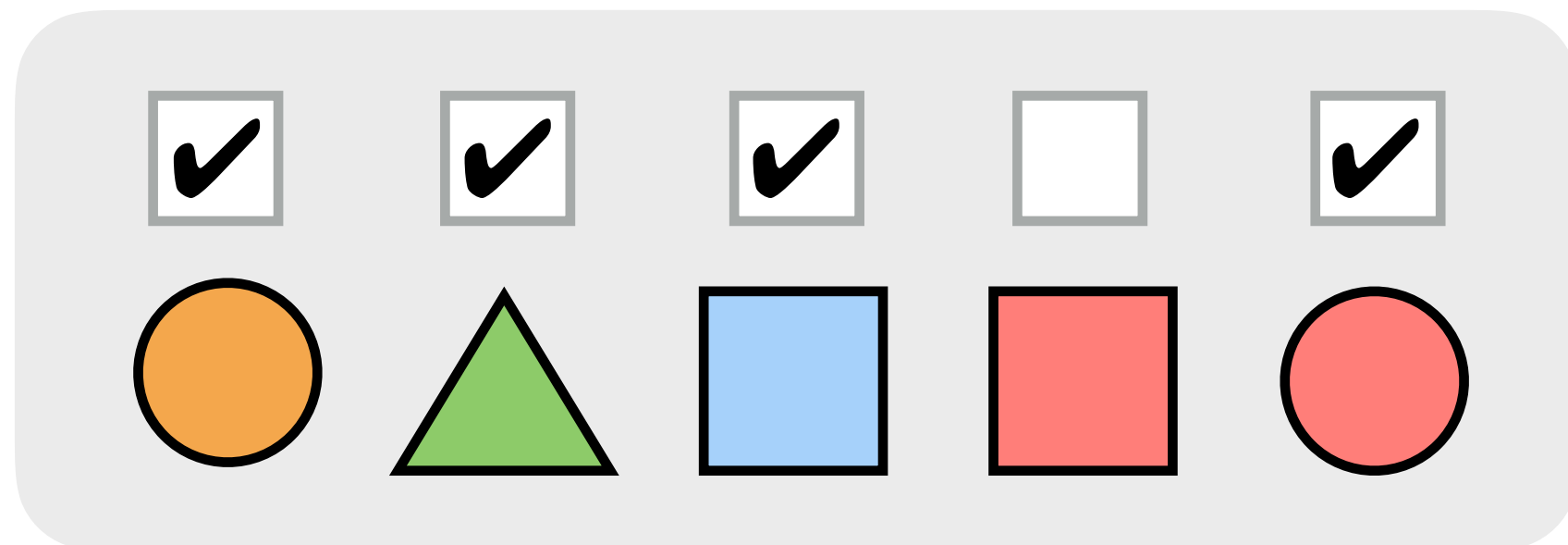
$\lambda x. \neg \text{red}(x) \vee \neg \text{sqr}(x)$



Computing meaning representations

not the red squares

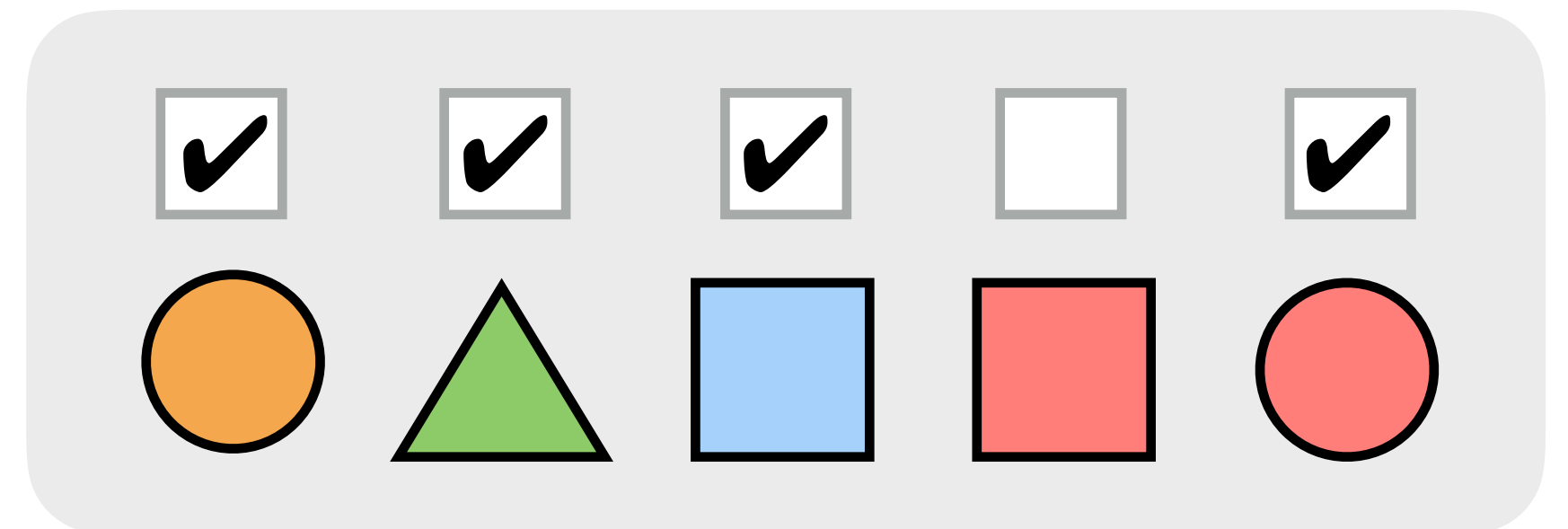
$\lambda x. \neg(\text{sqr}(x) \wedge \text{red}(x))$



=

not red or not square

$\lambda x. \neg \text{red}(x) \vee \neg \text{sqr}(x)$

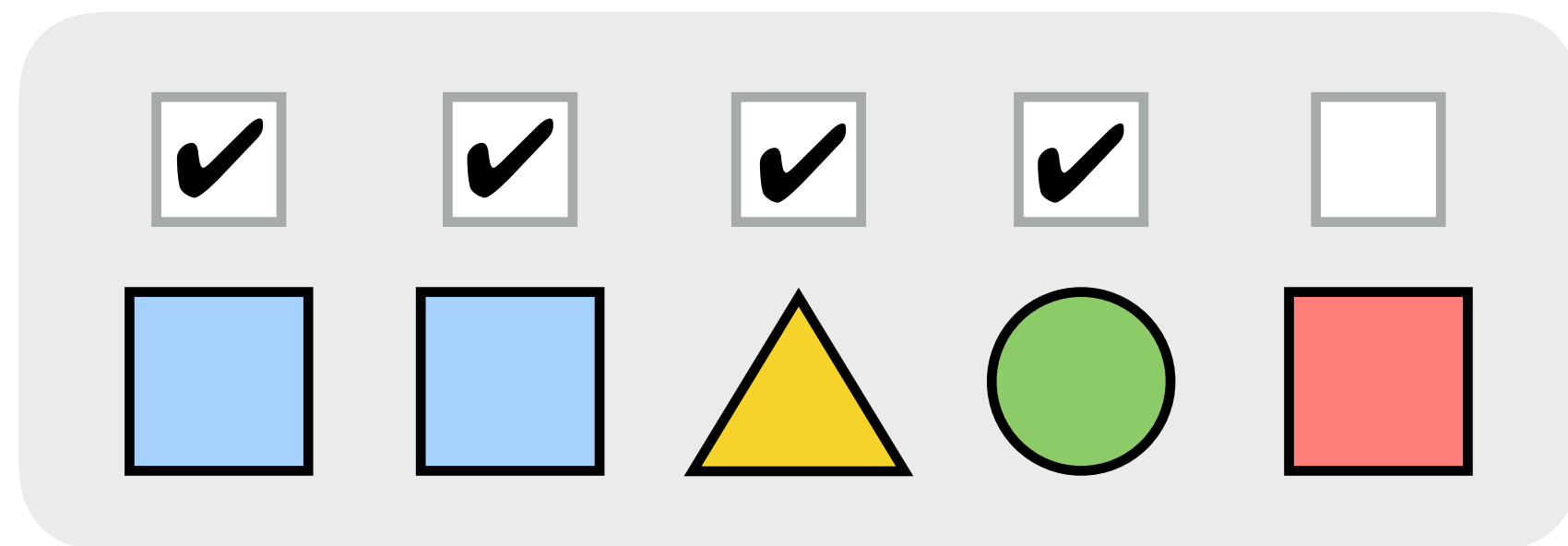




Computing meaning representations

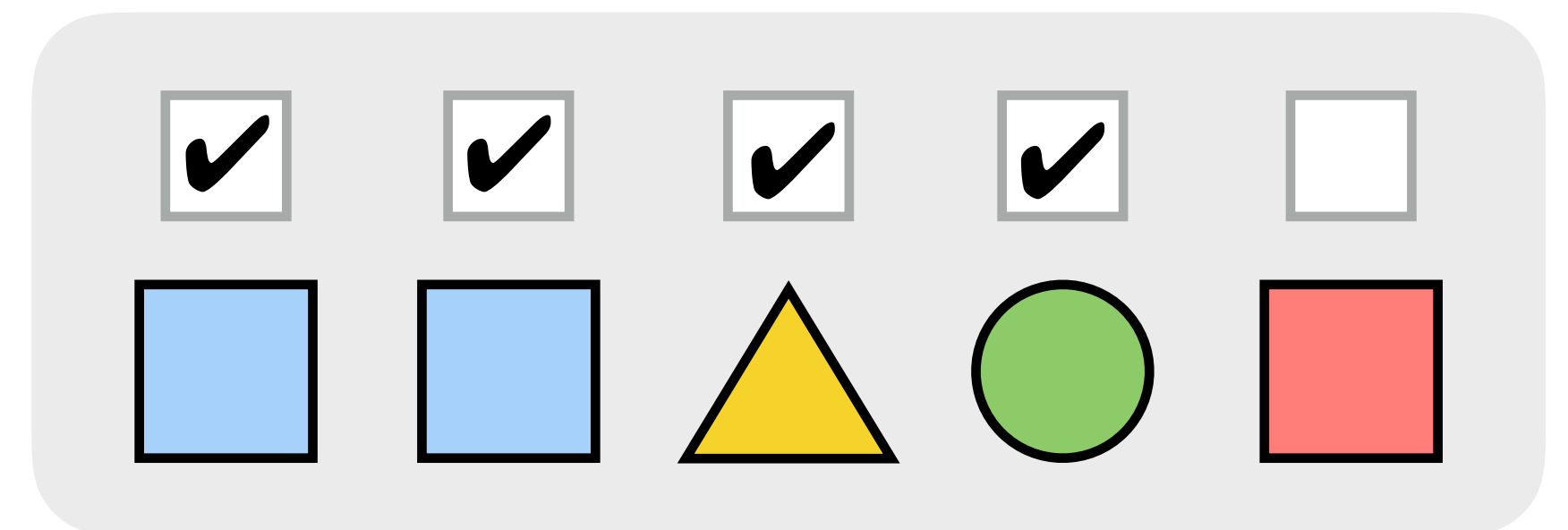
not the red squares

$\lambda x. \neg(\text{sqr}(x) \wedge \text{red}(x))$



not red or not square

$\lambda x. \neg \text{red}(x) \vee \neg \text{sqr}(x)$

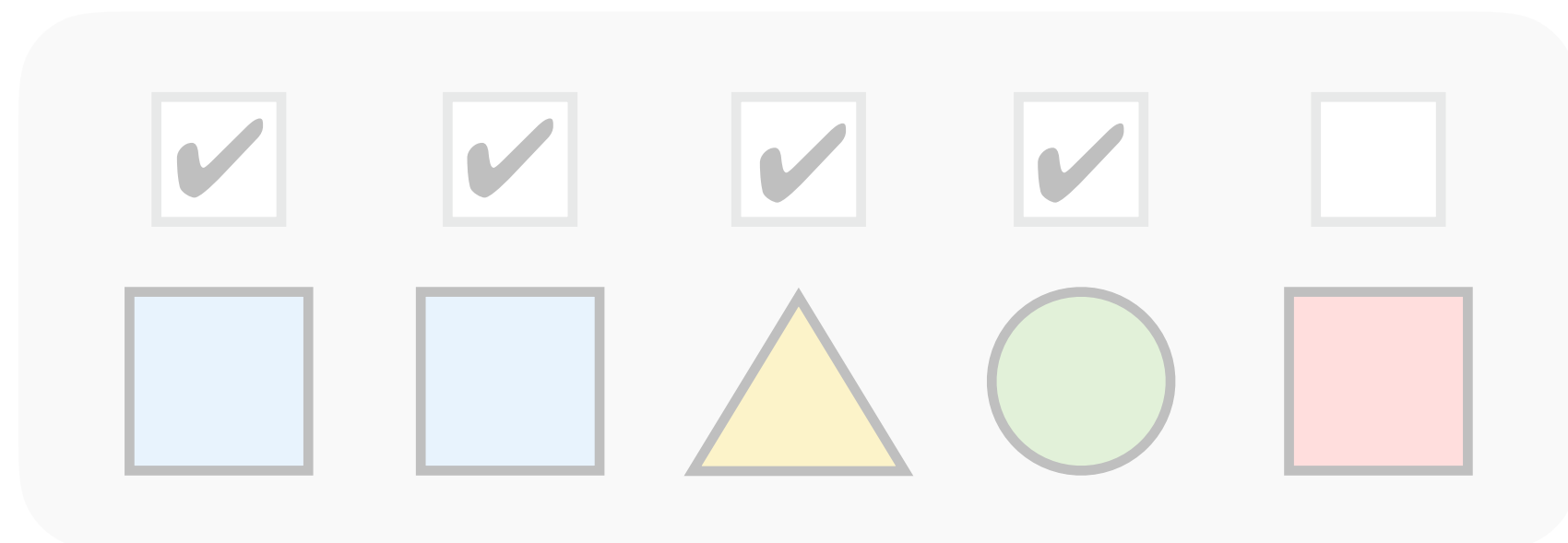




Computing meaning representations

not the red squares

$\lambda x. \neg(\text{sqr}(x) \wedge \text{red}(x))$



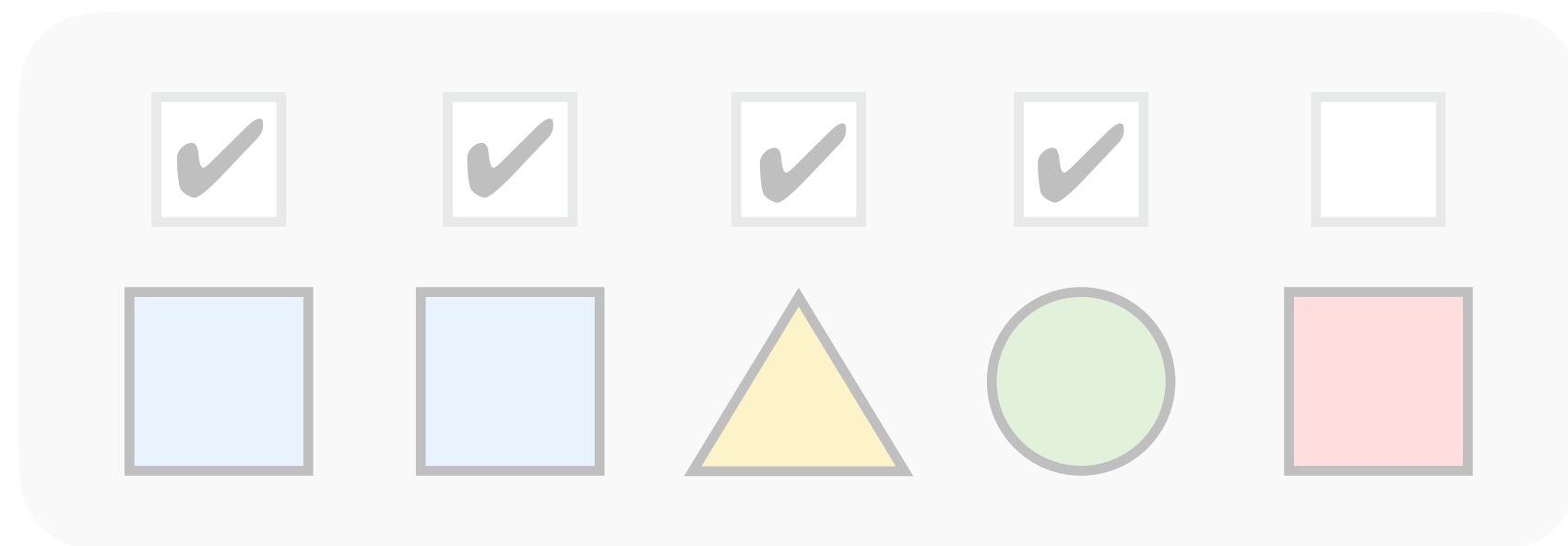
-0.1 1.3 0.5 -0.4



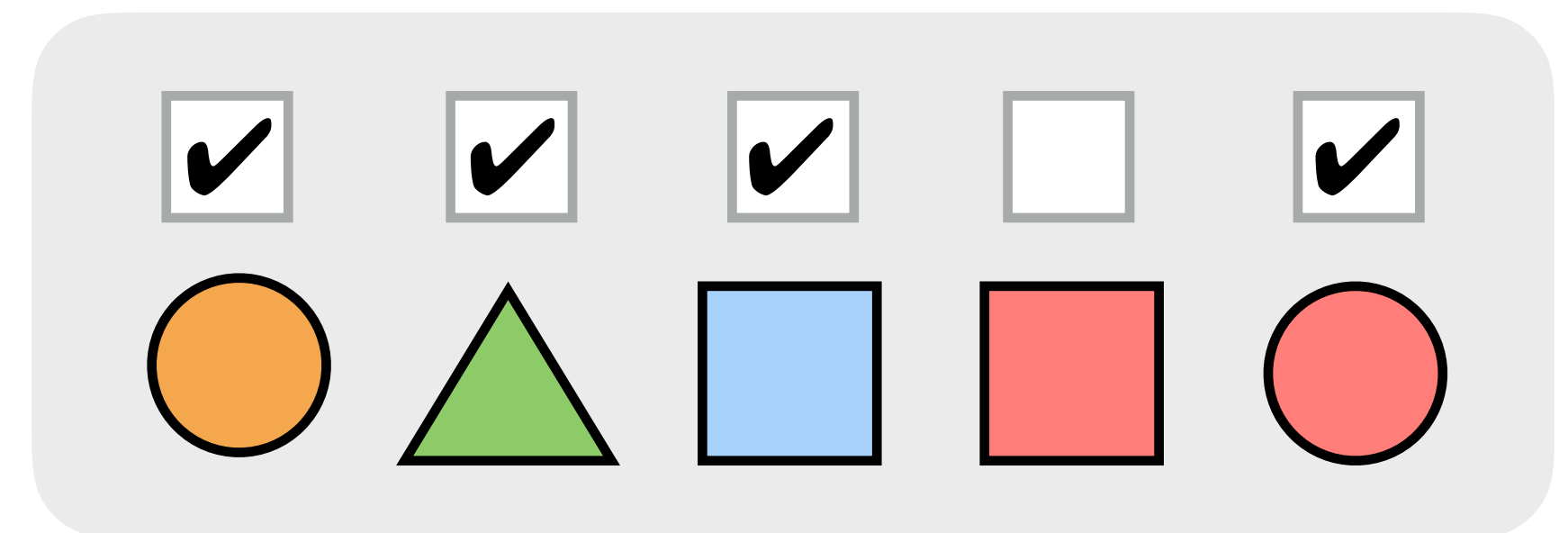
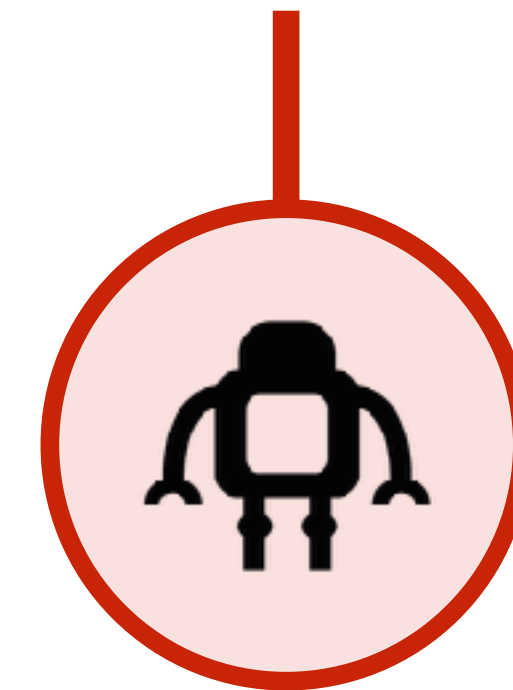
Computing meaning representations

not the red squares

$\lambda x. \neg(\text{sqr}(x) \wedge \text{red}(x))$



-0.1 1.3 0.5 -0.4





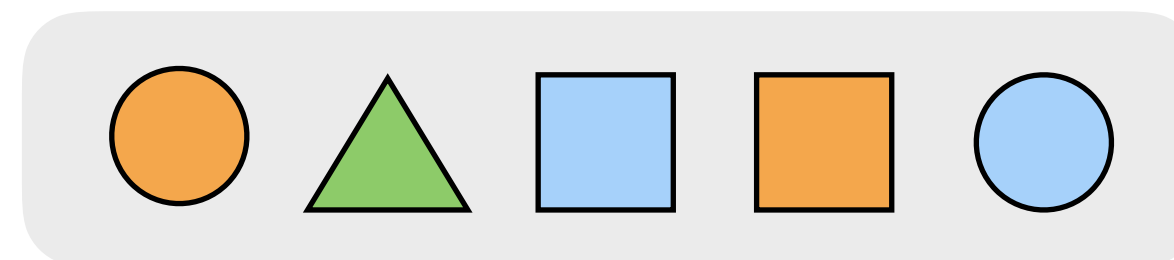
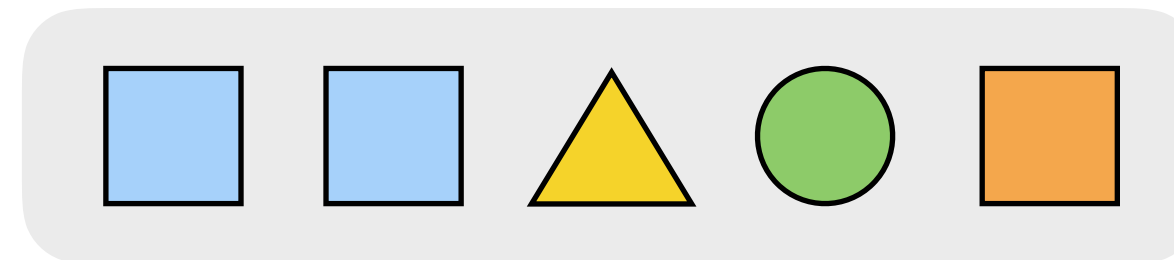
Computing meaning representations

-0.1 1.3
0.5 -0.4
0.2 1.0

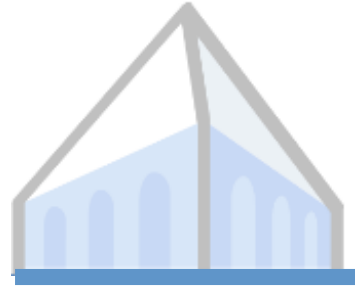


Computing meaning representations

-0.1 1.3
0.5 -0.4
0.2 1.0

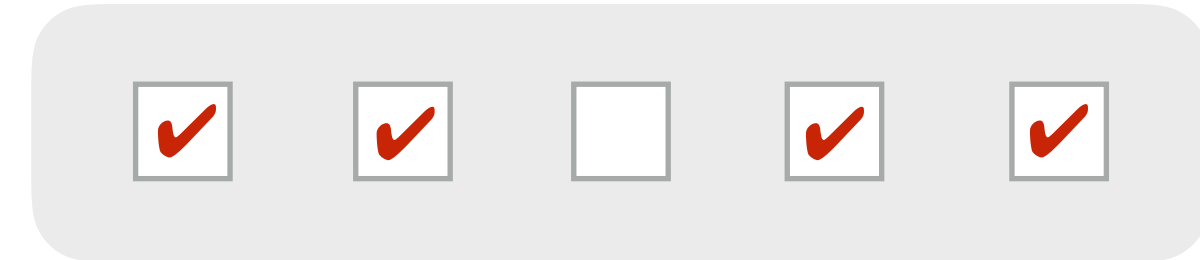
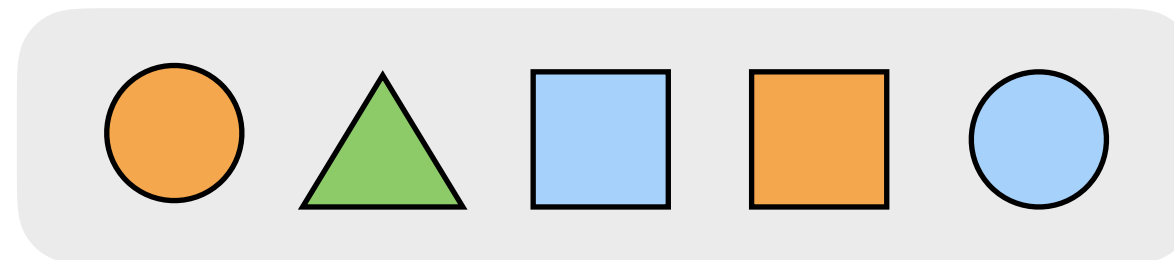
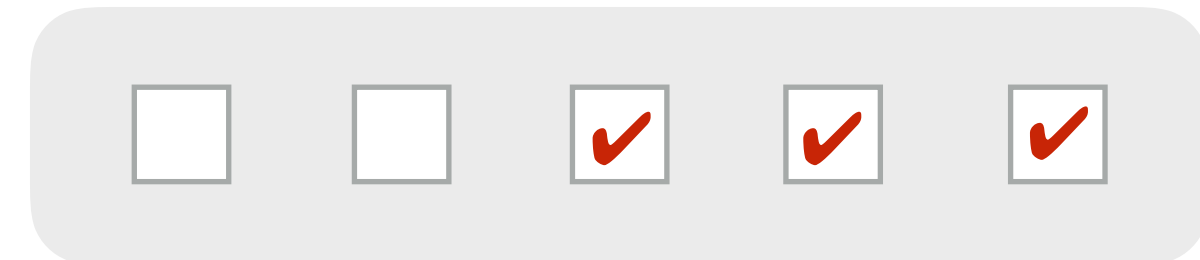
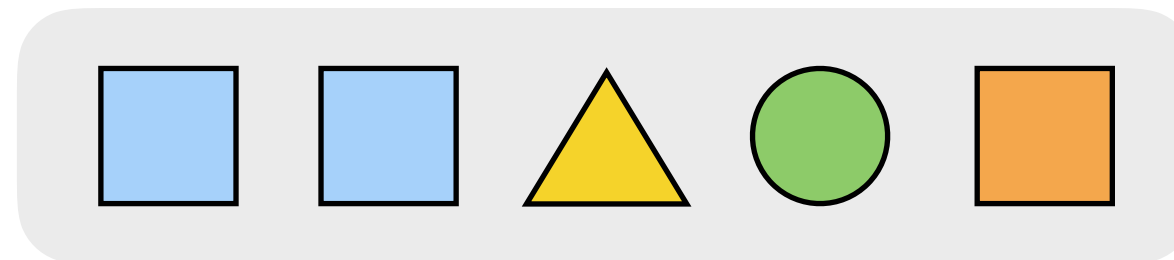


•
•
•



Computing meaning representations

$$\begin{matrix} -0.1 & 1.3 \\ 0.5 & -0.4 \\ 0.2 & 1.0 \end{matrix}$$

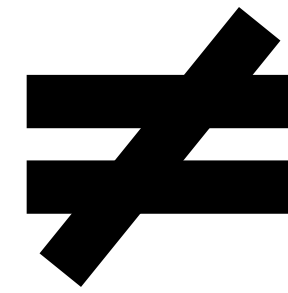


•
•
•

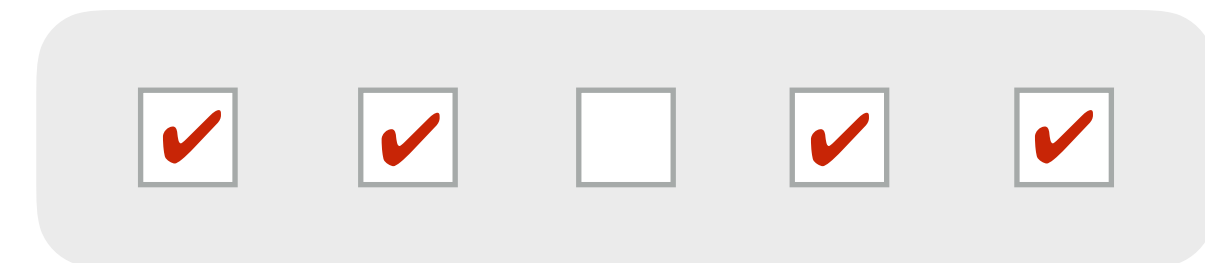
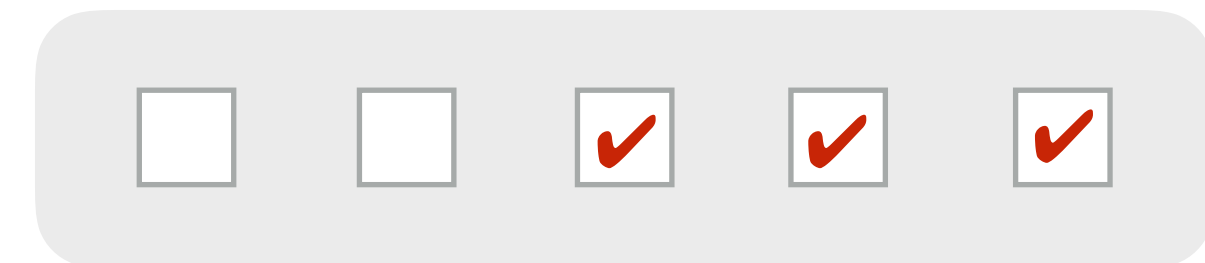
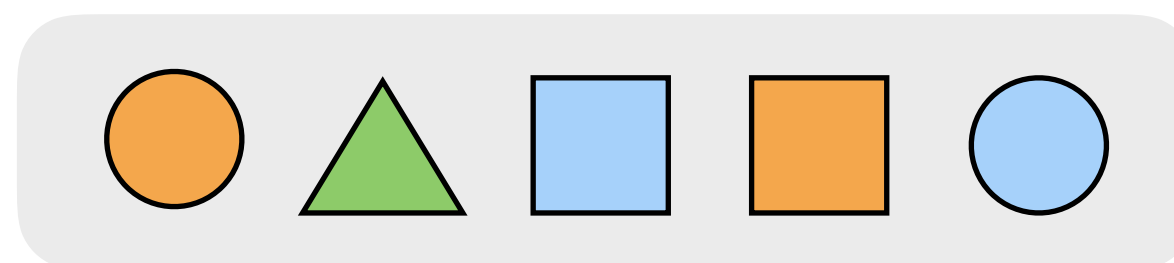
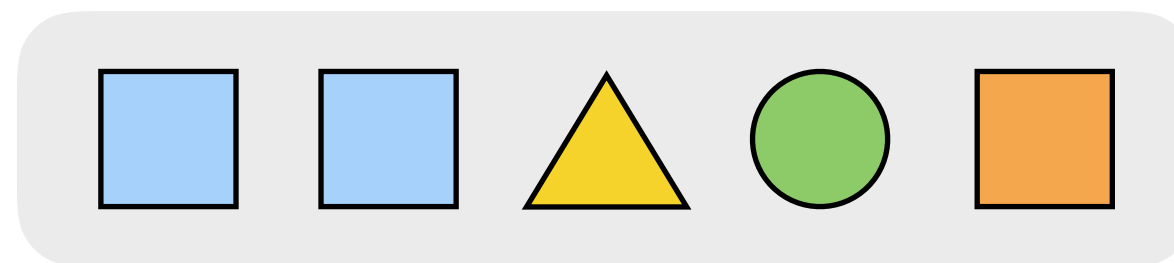
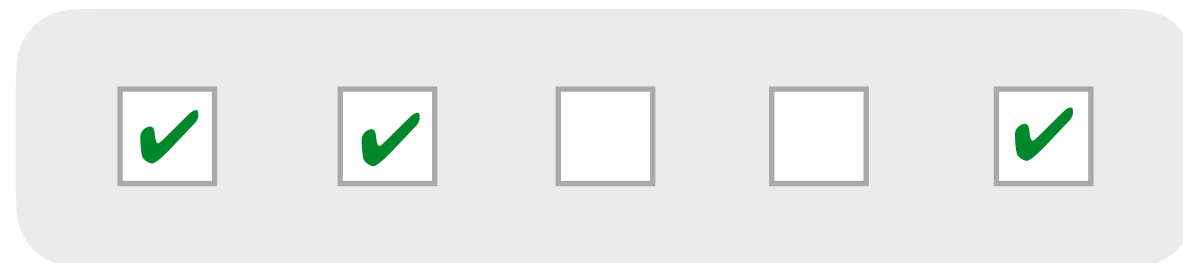
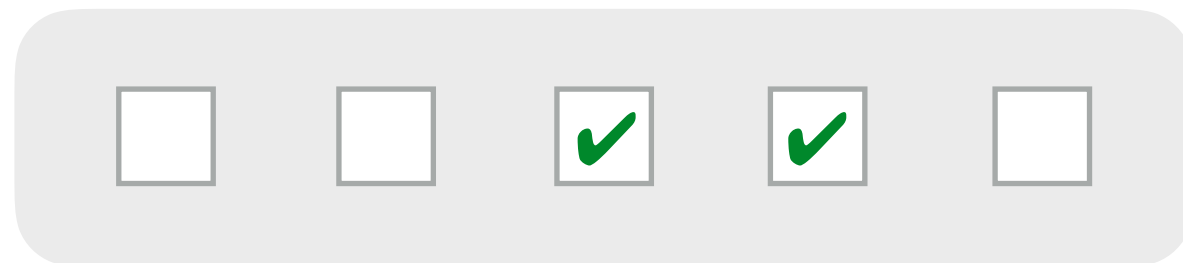


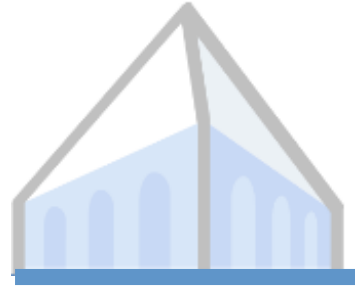
Computing meaning representations

everything but squares



-0.1 1.3
 0.5 -0.4
 0.2 1.0



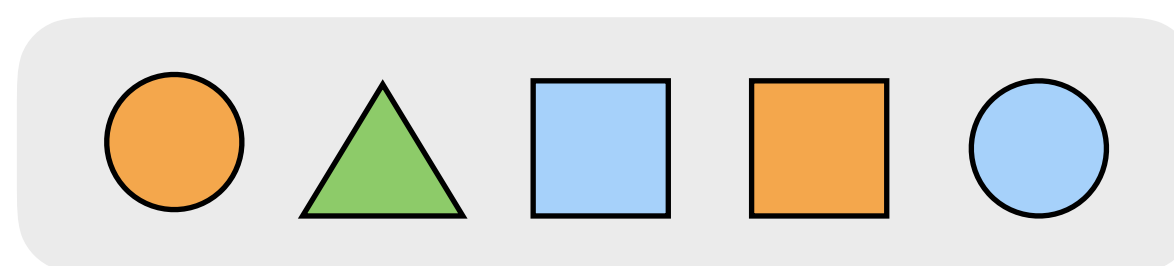
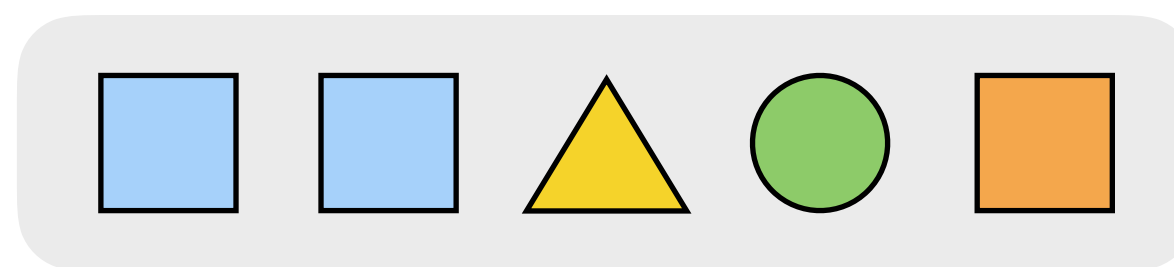
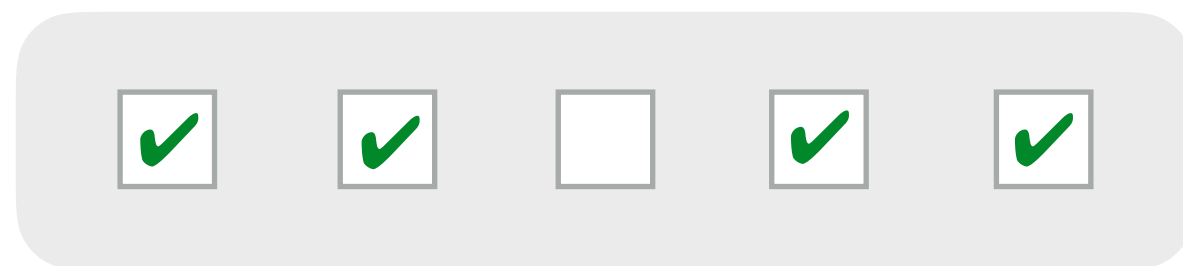
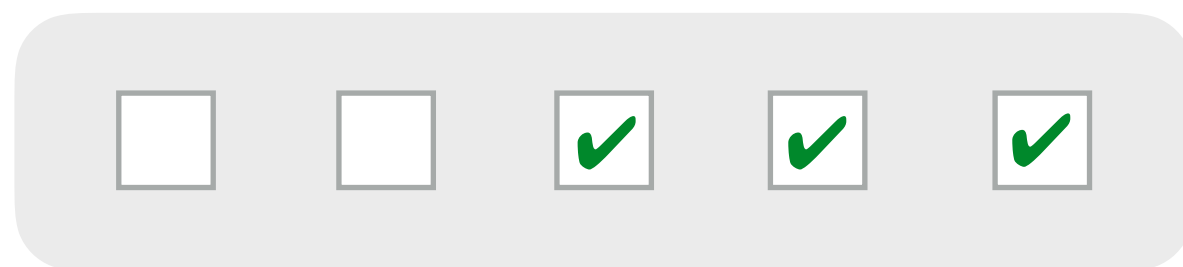


Computing meaning representations

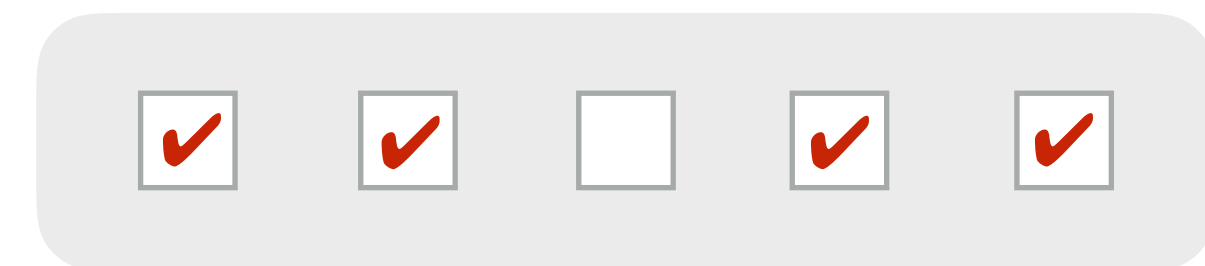
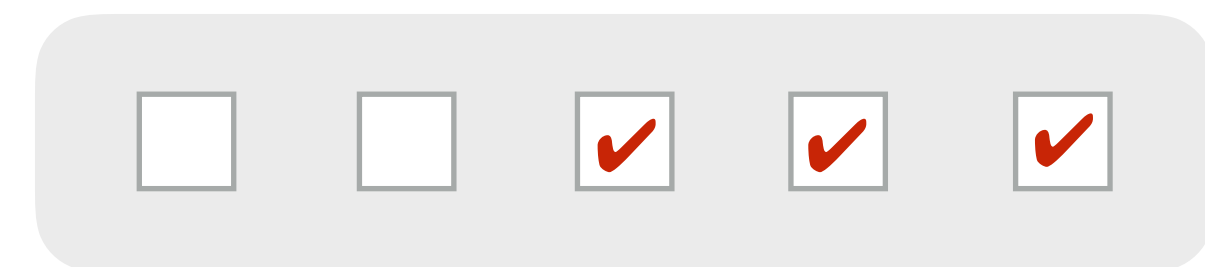
not the blue squares

=

-0.1 1.3
0.5 -0.4
0.2 1.0



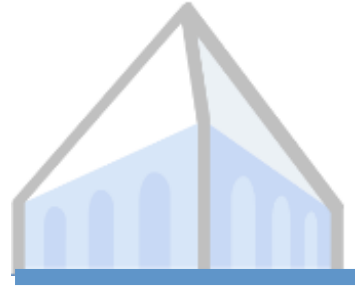
•
•
•





Translating

By comparing denotations from **logical forms** and the **decoder model**, we can find **utterances** and **vectors** with the same meaning.



Questions

1. Does the RNN employ a human-like communicative strategy?
2. Do RNN representations have interpretable compositional structure?

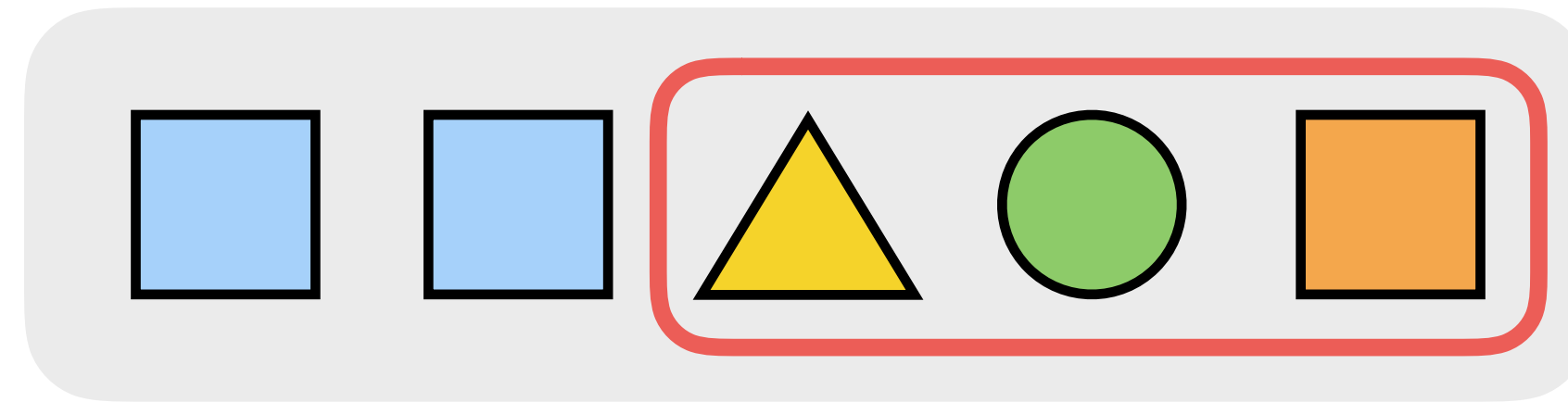


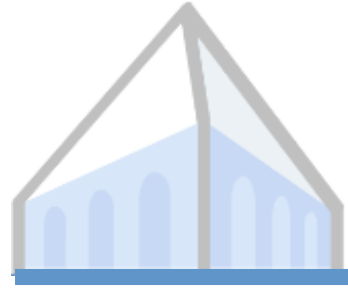
Questions

1. Does the RNN employ a human-like communicative strategy?
2. Do RNN representations have interpretable compositional structure?



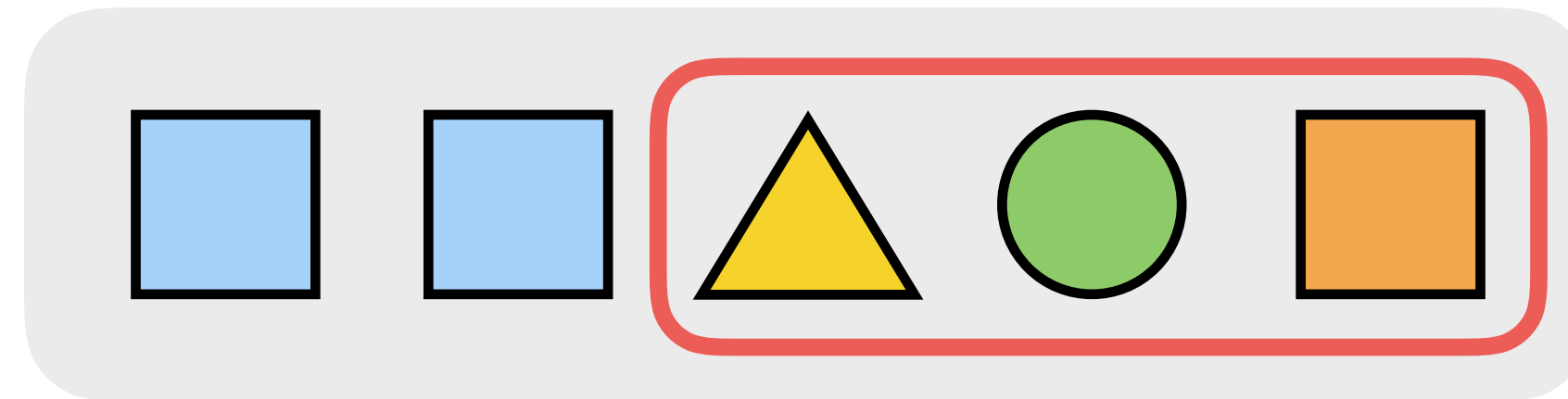
Comparing strategies





Comparing strategies

-0.1	1.3
0.5	-0.4
0.2	1.0

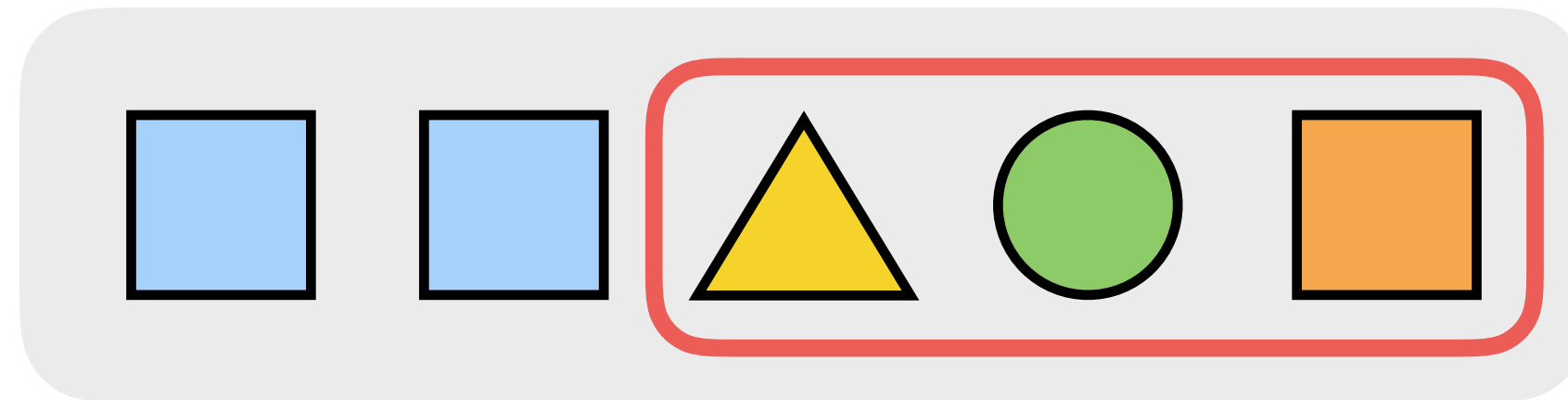


*everything
but squares*

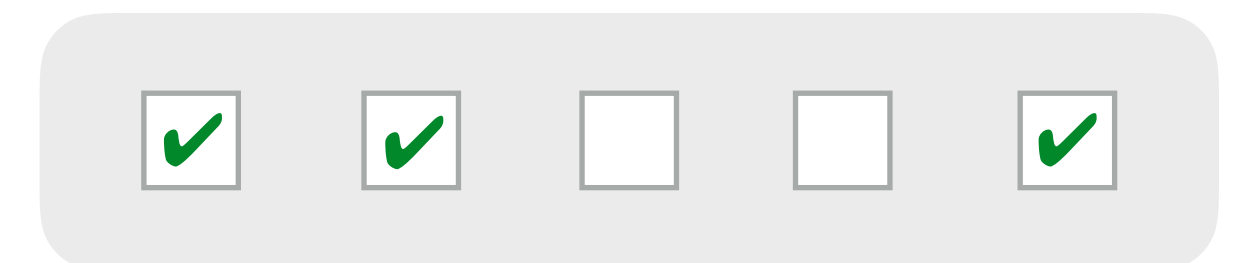
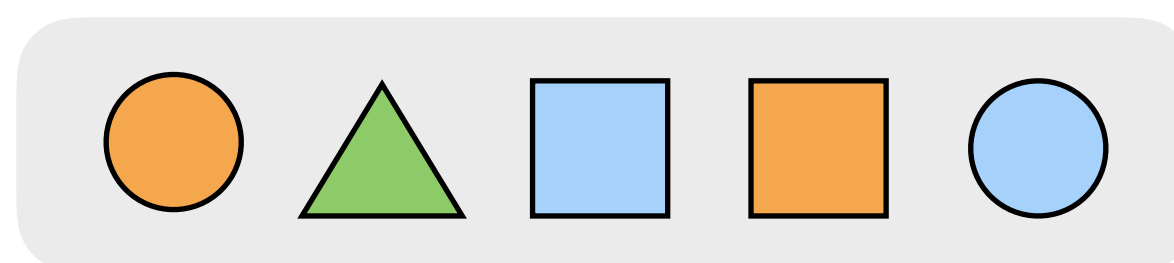
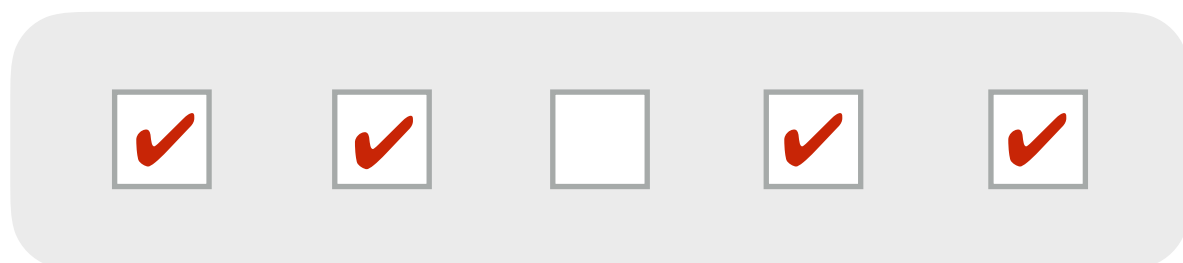
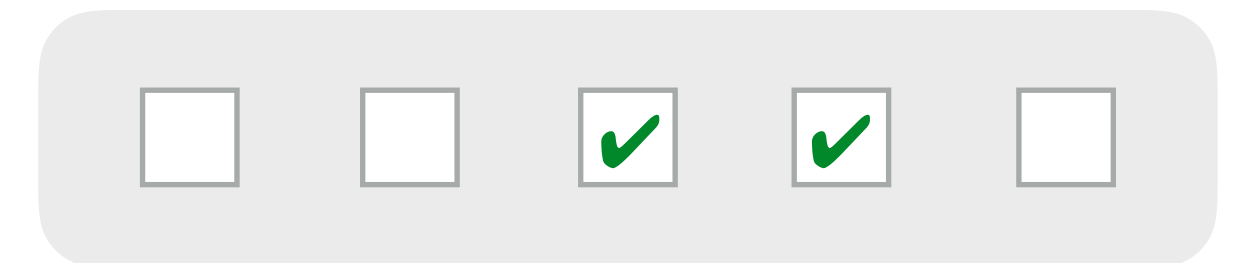
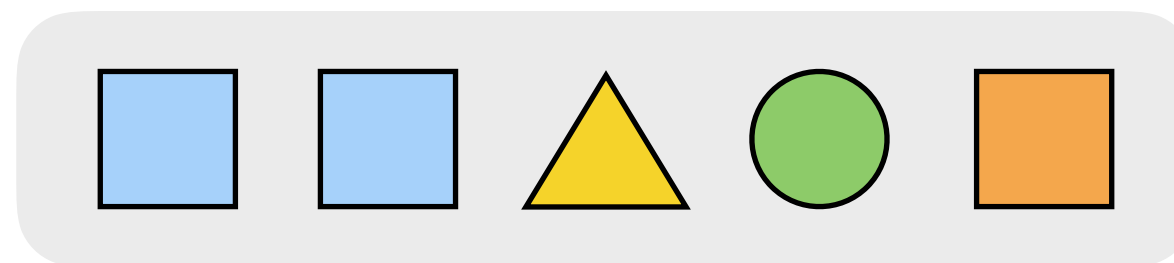
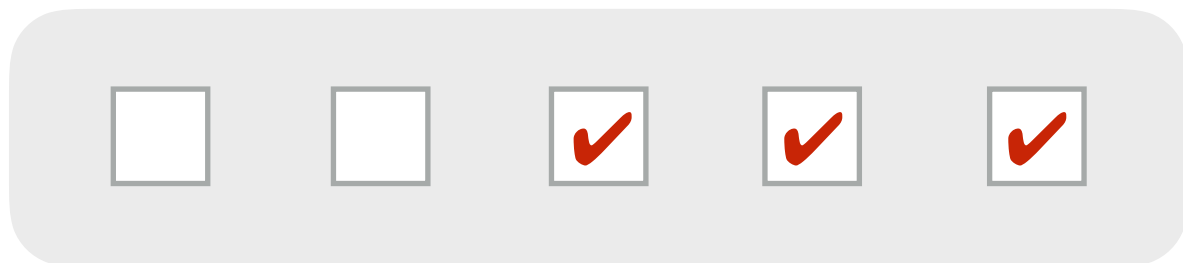


Comparing strategies

-0.1	1.3
0.5	-0.4
0.2	1.0



*everything
but squares*

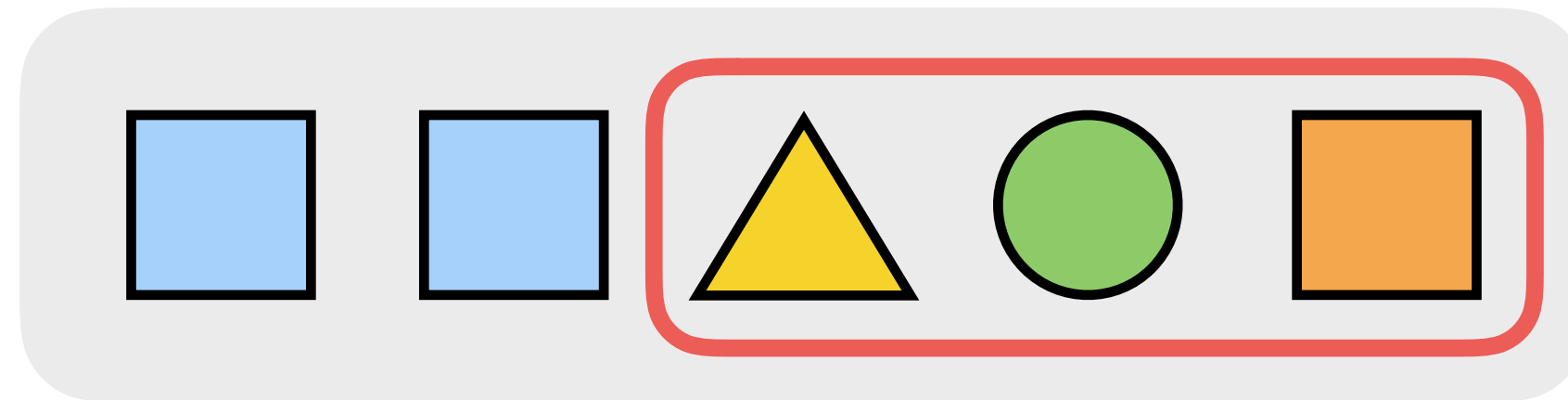


•
•
•

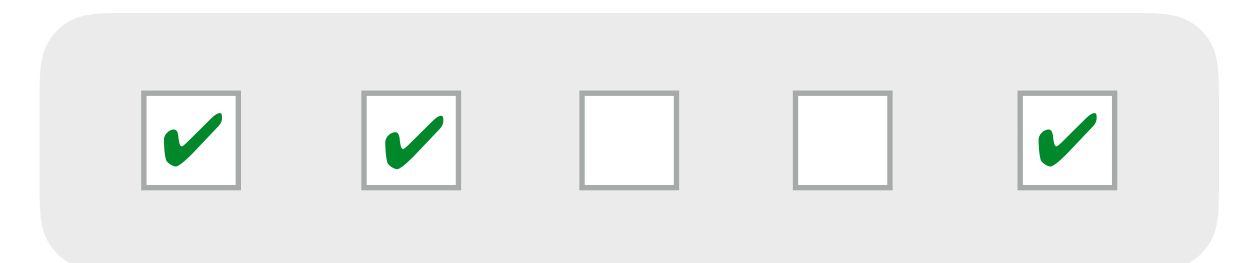
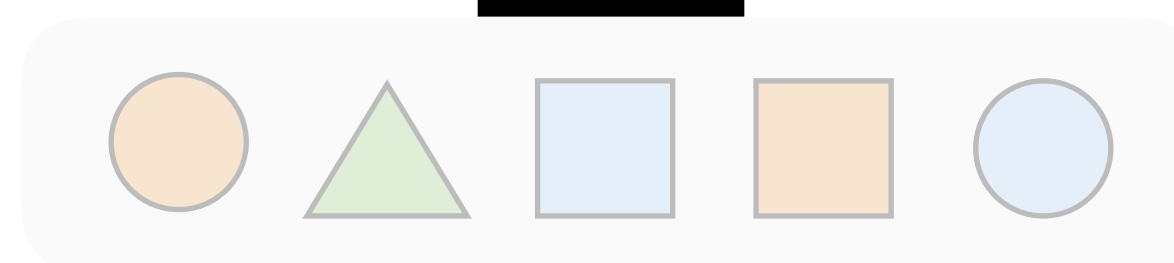
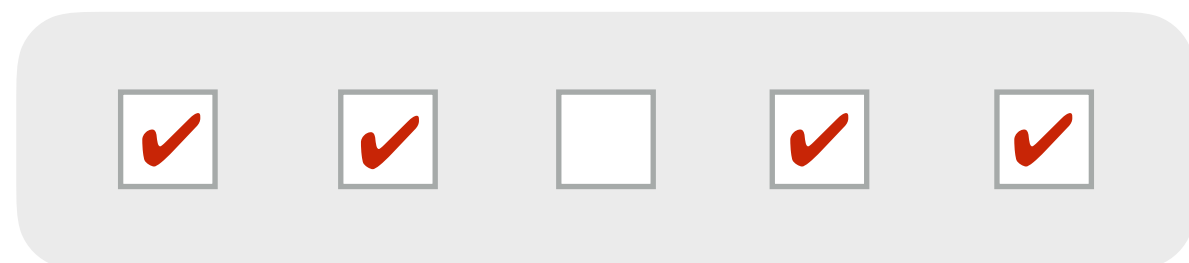
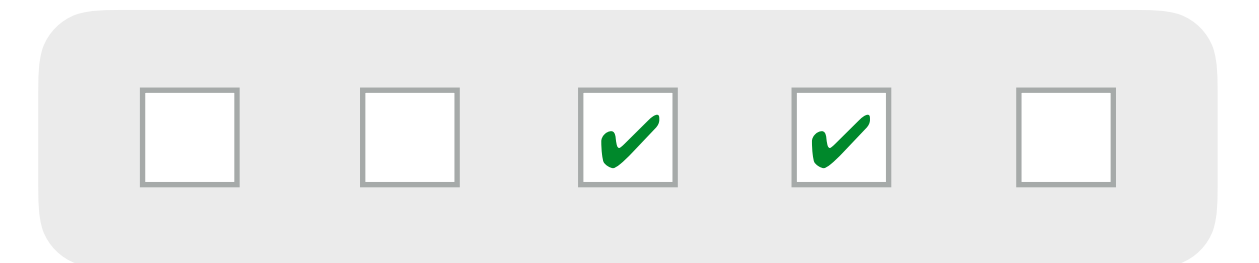
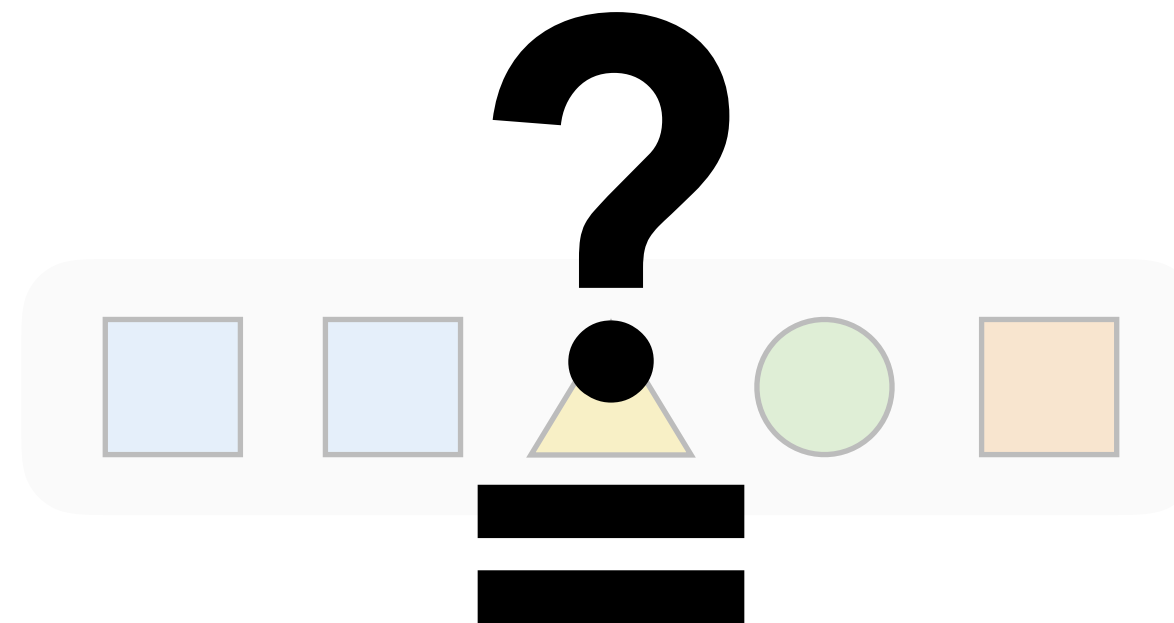
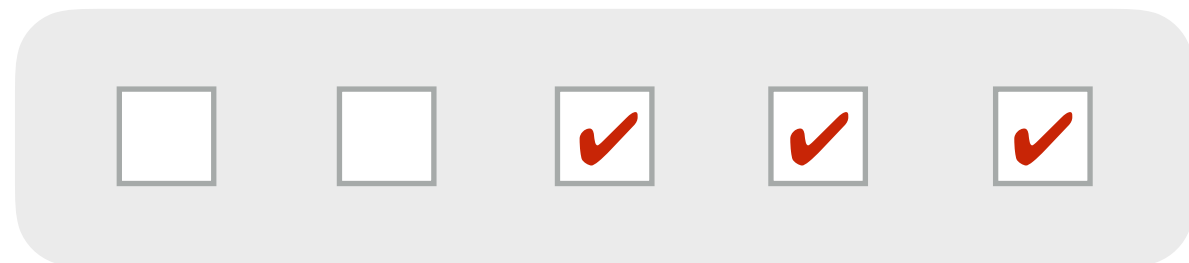


Comparing strategies

-0.1	1.3
0.5	-0.4
0.2	1.0

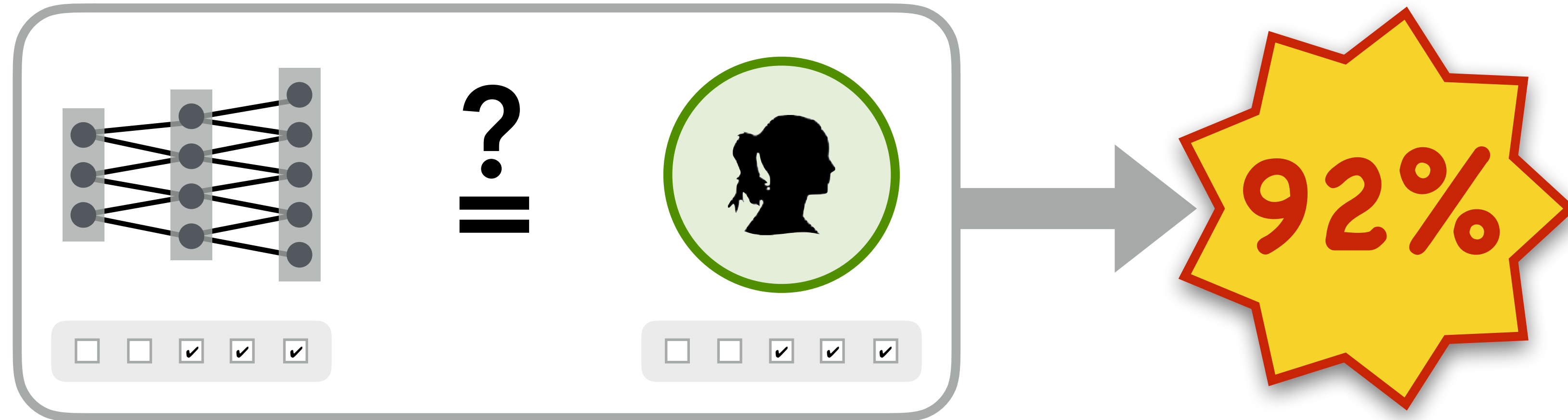


*everything
but squares*



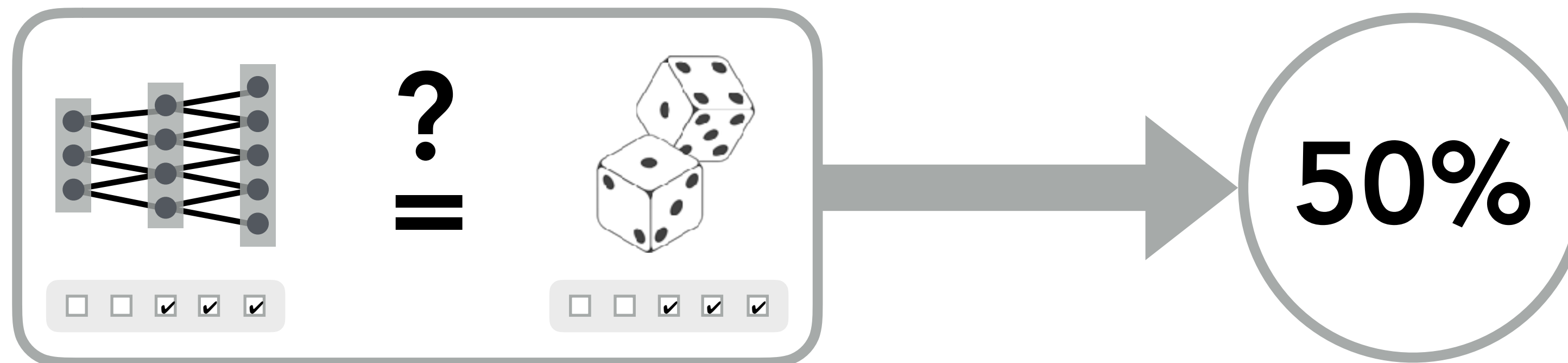
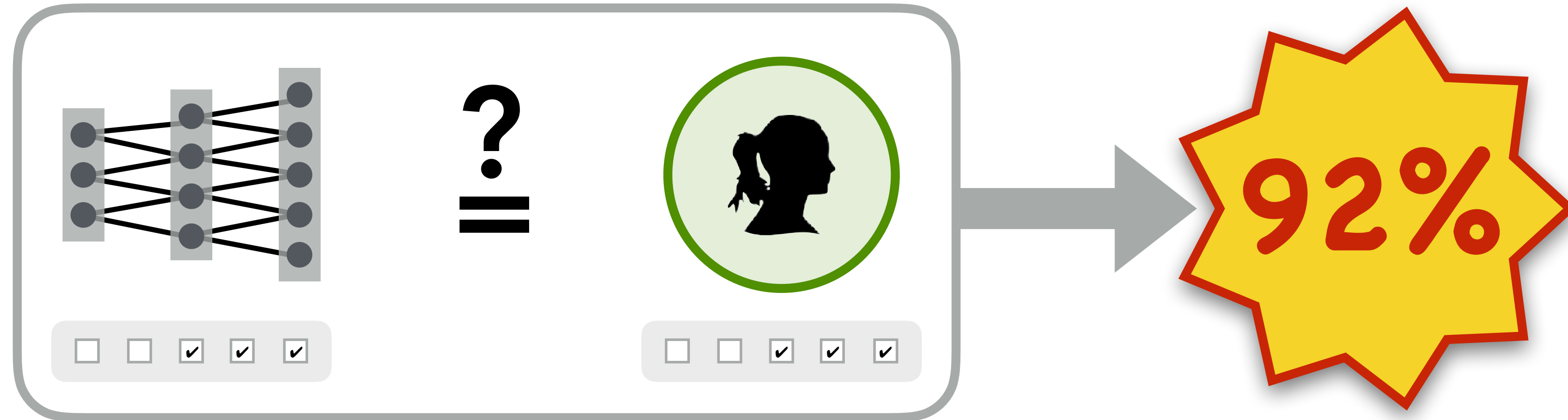


Evaluation: strategies



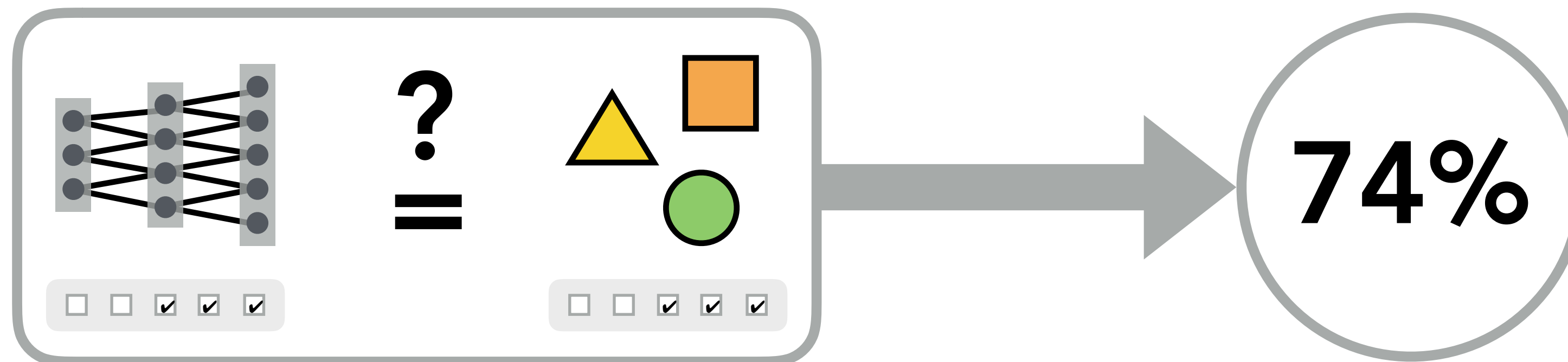
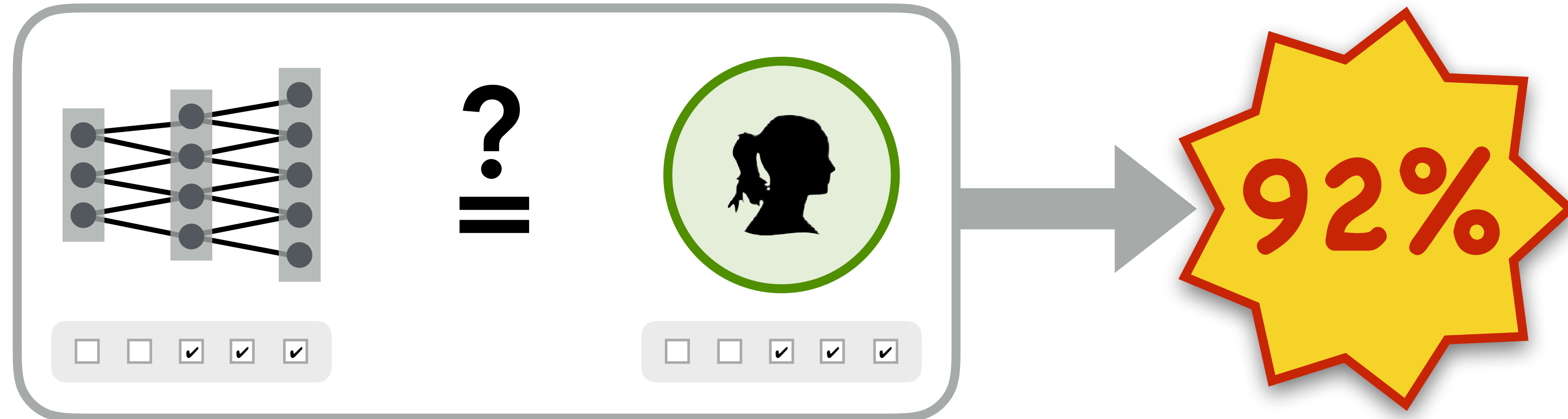


Evaluation: strategies





Evaluation: strategies





Experiments

1. Does the RNN employ a human-like communicative strategy?
2. Do RNN representations have interpretable compositional structure?



Collecting translation data

all the red shapes

blue objects

everything but red

green squares

not green squares



Collecting translation data

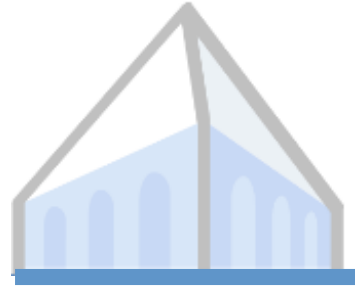
$\lambda x. \text{red}(x)$

$\lambda x. \text{blu}(x)$

$\lambda x. \neg \text{red}(x)$

$\lambda x. \text{grn}(x) \wedge \text{sqr}(x)$

$\lambda x. \neg (\text{grn}(x) \wedge \text{sqr}(x))$



Collecting translation data

$\lambda x. \text{red}(x)$



0.1 -0.3 0.5 1.1

$\lambda x. \text{blu}(x)$



-0.3 0.2 0.1 0.1

$\lambda x. \neg \text{red}(x)$



1.4 -0.3 -0.5 0.8

$\lambda x. \text{grn}(x) \wedge \text{sqr}(x)$



0.2 -0.2 0.5 -0.1

$\lambda x. \neg (\text{grn}(x) \wedge \text{sqr}(x))$



0.3 -1.3 -1.5 0.1



Extracting related pairs

$\lambda x. \text{red}(x)$

0.1 -0.3 0.5 1.1

$\lambda x. \text{grn}(x) \wedge \text{sqr}(x)$

0.2 -0.2 0.5 -0.1

$\lambda x. \neg \text{red}(x)$

1.4 -0.3 -0.5 0.8

$\lambda x. \neg(\text{grn}(x) \wedge \text{sqr}(x))$

0.3 -1.3 -1.5 0.1



Extracting related pairs

$\lambda x. \text{red}(x)$

0.1 -0.3 0.5 1.1

$\lambda x. \neg \text{red}(x)$

1.4 -0.3 -0.5 0.8

$\lambda x. \text{grn}(x) \wedge \text{sqr}(x)$

0.2 -0.2 0.5 -0.1

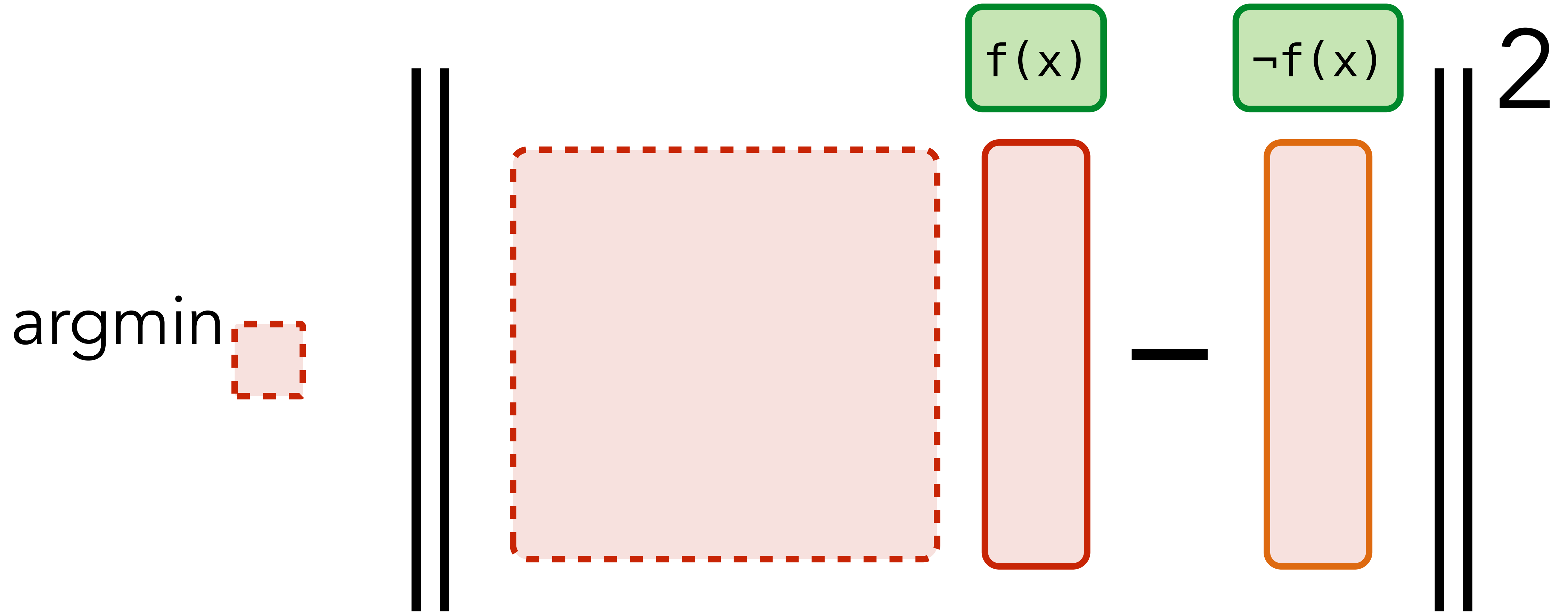
$\lambda x. \neg (\text{grn}(x) \wedge \text{sqr}(x))$

0.3 -1.3 -1.5 0.1





Learning compositional operators





Evaluating learned operators

$\lambda x. \text{red}(x)$

0.1 -0.3 0.5 1.1

$\lambda x. \neg \text{red}(x)$

1.4 -0.3 -0.5 0.8

$\lambda x. \text{grn}(x) \wedge \text{sqr}(x)$

0.2 -0.2 0.5 -0.1

$\lambda x. \neg (\text{grn}(x) \wedge \text{sqr}(x))$

0.3 -1.3 -1.5 0.1

$\lambda x. f(x)$

0.2 -0.2 0.5 -0.1



Evaluating learned operators

$\lambda x. \text{red}(x)$

0.1 -0.3 0.5 1.1

$\lambda x. \text{grn}(x) \wedge \text{sqr}(x)$

0.2 -0.2 0.5 -0.1

$\lambda x. f(x)$

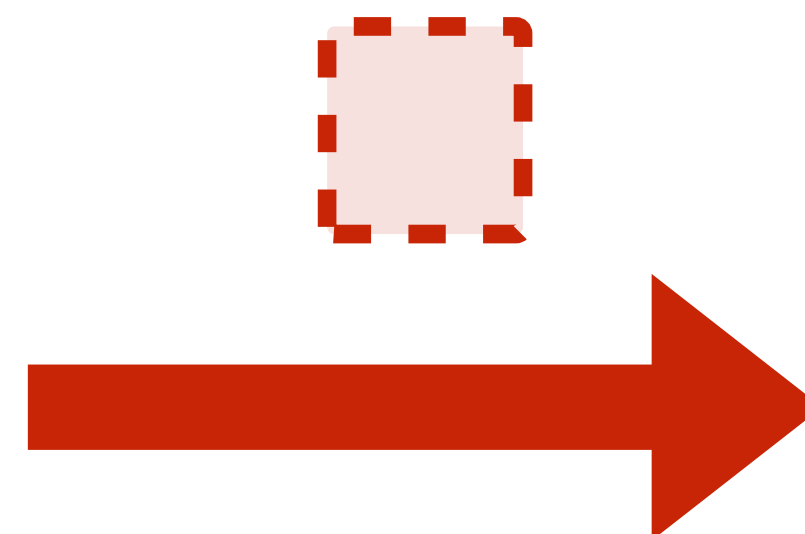
0.2 -0.2 0.5 -0.1

$\lambda x. \neg \text{red}(x)$

1.4 -0.3 -0.5 0.8

$\lambda x. \neg(\text{grn}(x) \wedge \text{sqr}(x))$

0.3 -1.3 -1.5 0.1



-0.2 0.4 -0.3 0.0



Evaluating learned operators

$\lambda x. \text{red}(x)$

0.1 -0.3 0.5 1.1

$\lambda x. \text{grn}(x) \wedge \text{sqr}(x)$

0.2 -0.2 0.5 -0.1

$\lambda x. f(x)$

0.2 -0.2 0.5 -0.1

$\lambda x. \neg \text{red}(x)$

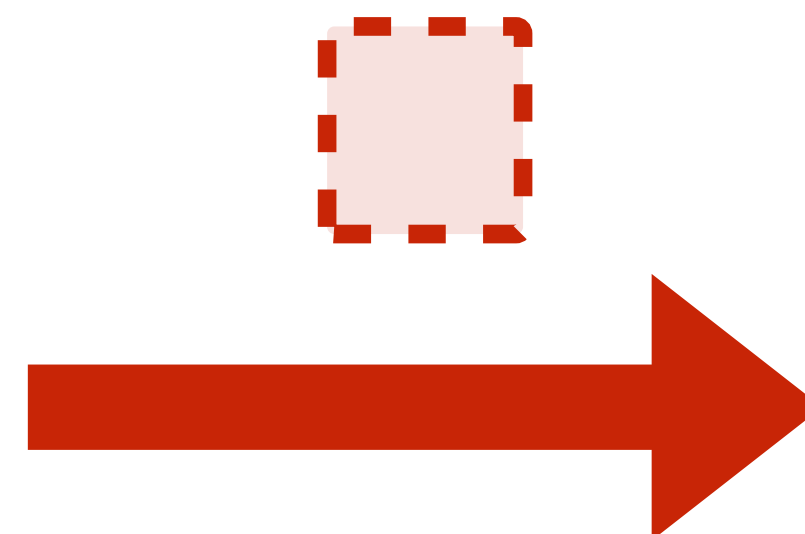
1.4 -0.3 -0.5 0.8

$\lambda x. \neg (\text{grn}(x) \wedge \text{sqr}(x))$

0.3 -1.3 -1.5 0.1

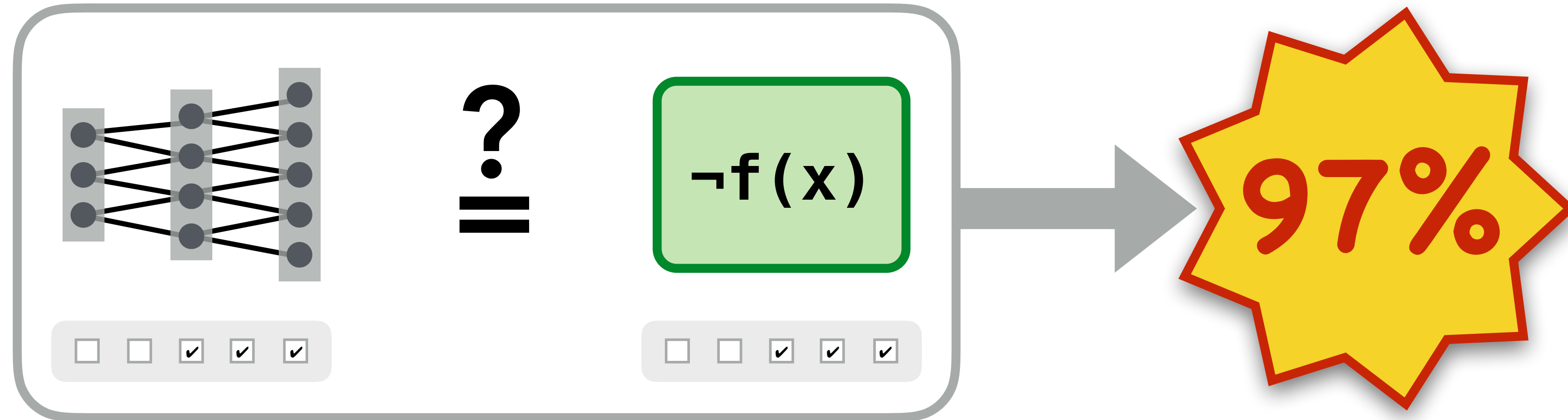
???

-0.2 0.4 -0.3 0.0



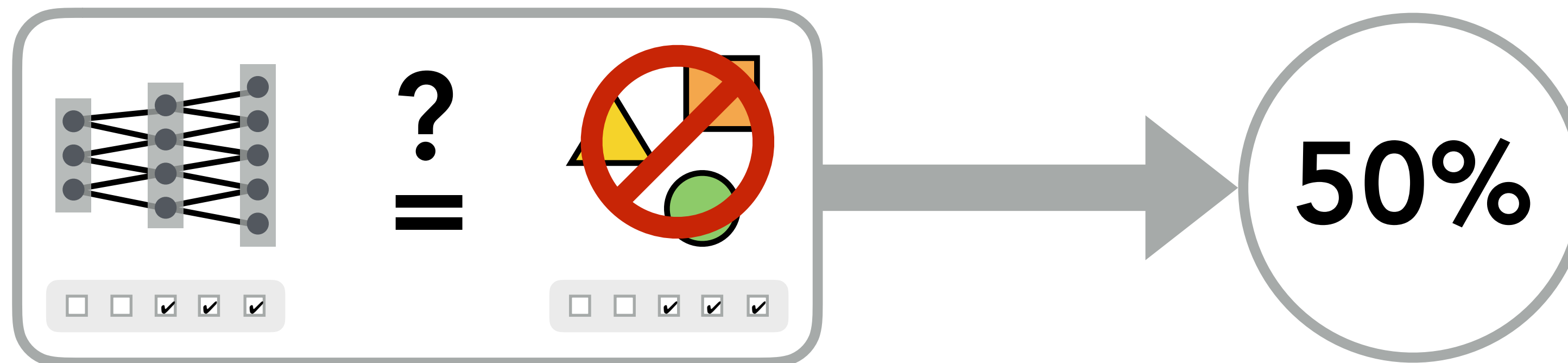
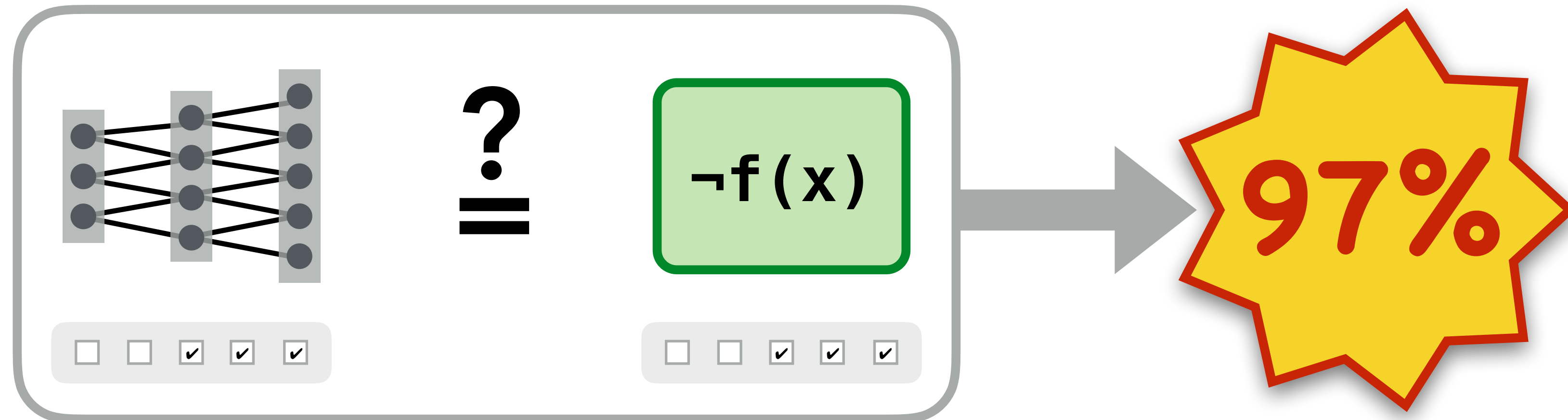


Evaluation: negation





Evaluation: negation



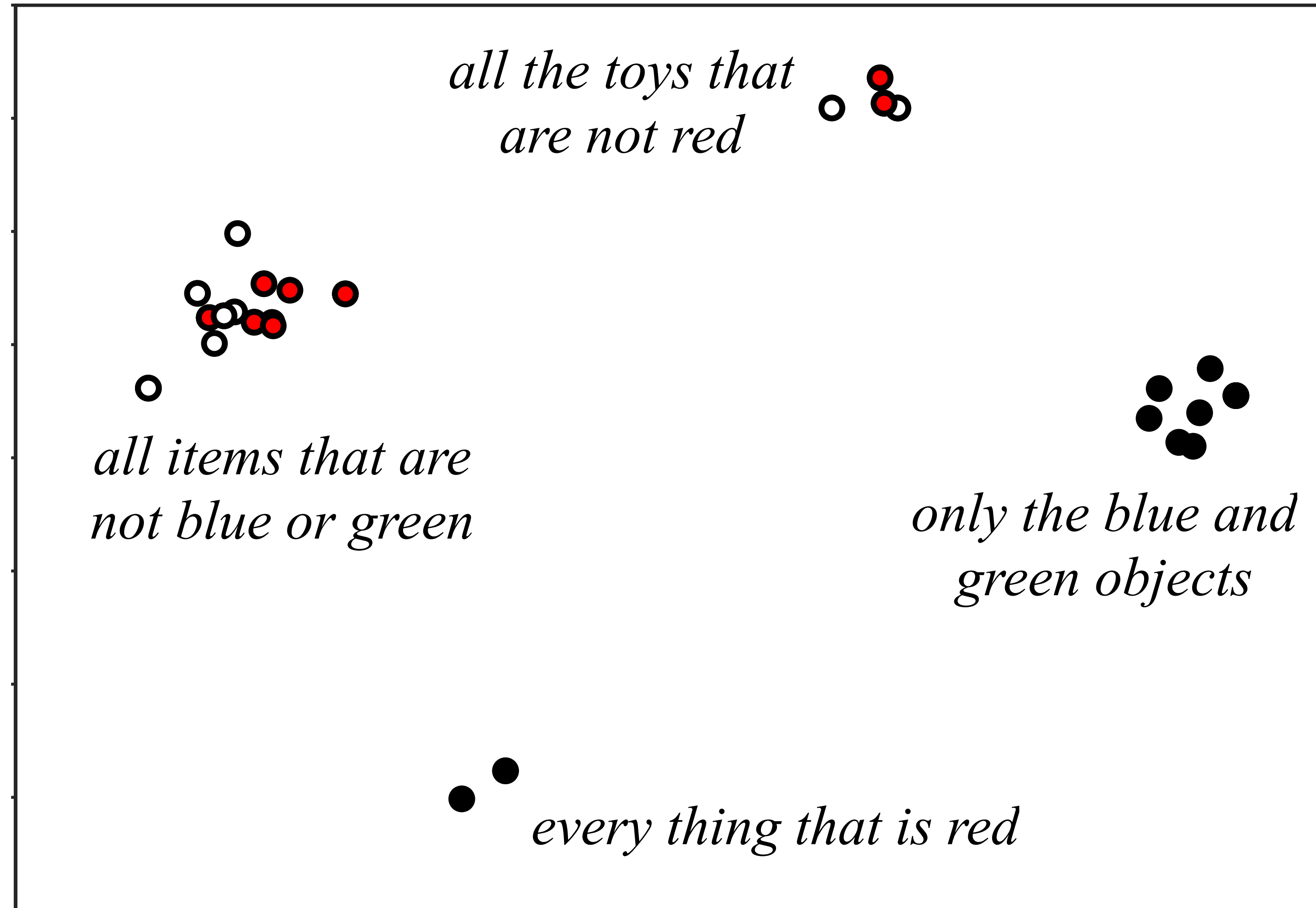


Visualizing negation

●
Input

●
Predicted

○
True



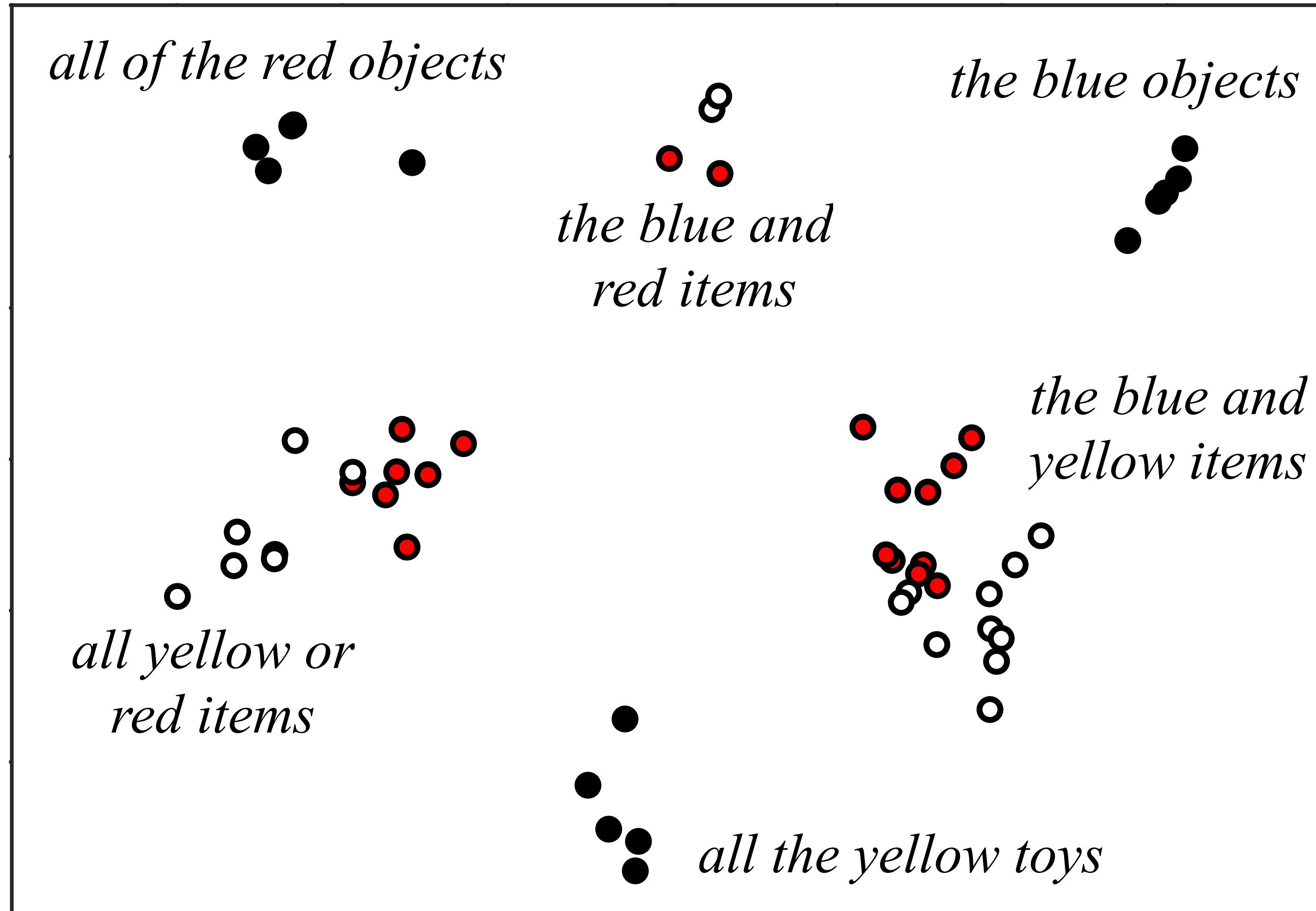


Visualizing disjunction

●
Input

●
Predicted

○
True





Conclusions

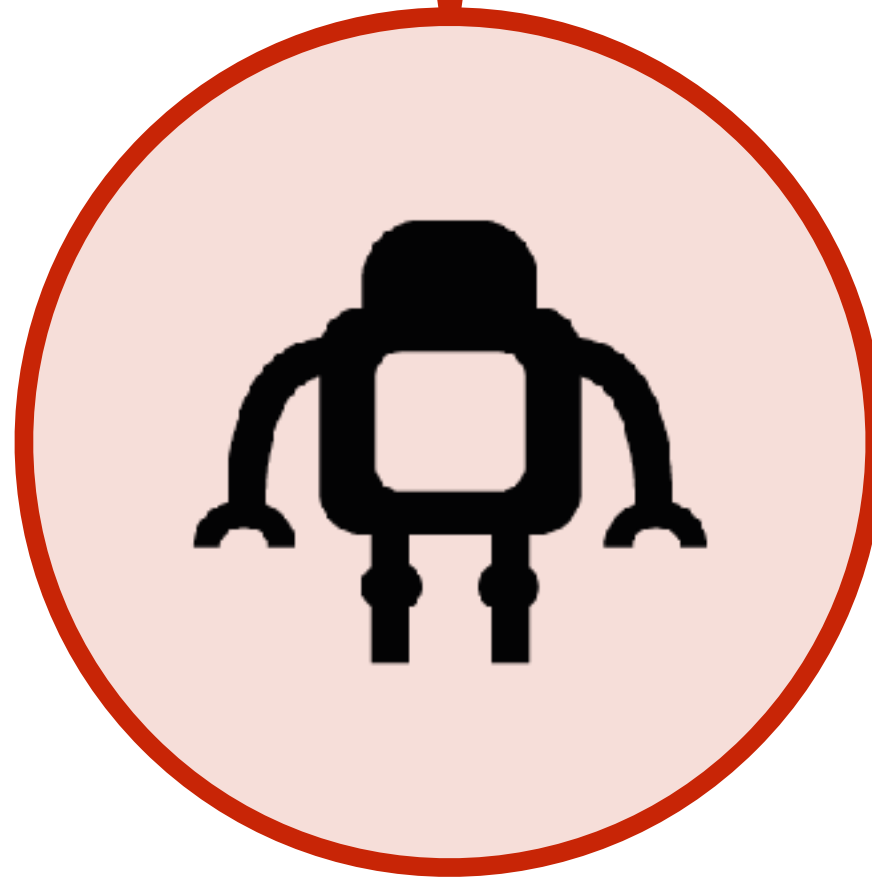
- Under the right conditions, RNN reprs exhibit interpretable pragmatics & compositional structure
- Not just communication games—language might be a good general-purpose tool for interpreting deep reprs.



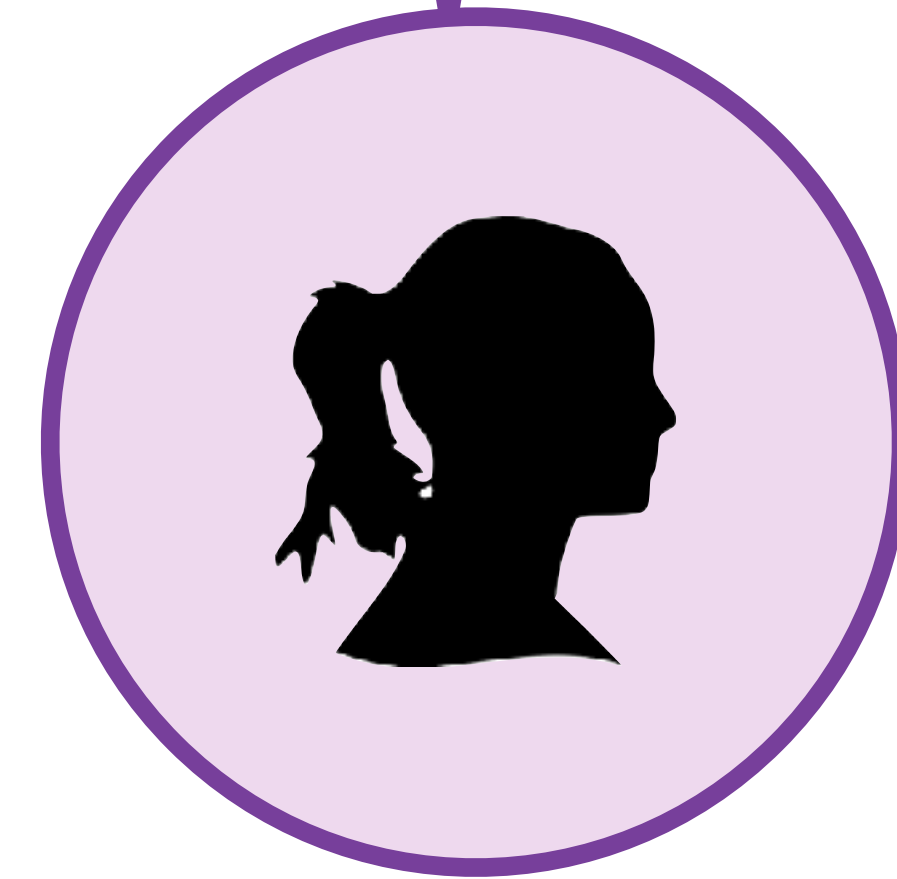
Conclusions

- Under the right conditions, RNN reprs exhibit interpretable pragmatics & compositional structure
- Not just communication games—language might be a good general-purpose tool for interpreting deep reprs.

1.0	2.3
-0.3	0.4
-1.2	1.1



Thank you!



<http://github.com/jacobandreas/rnn-syn>