

Structure and Interpretation of Neural Codes

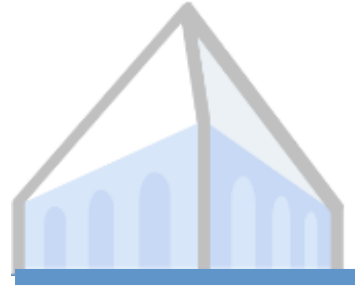


Jacob Andreas

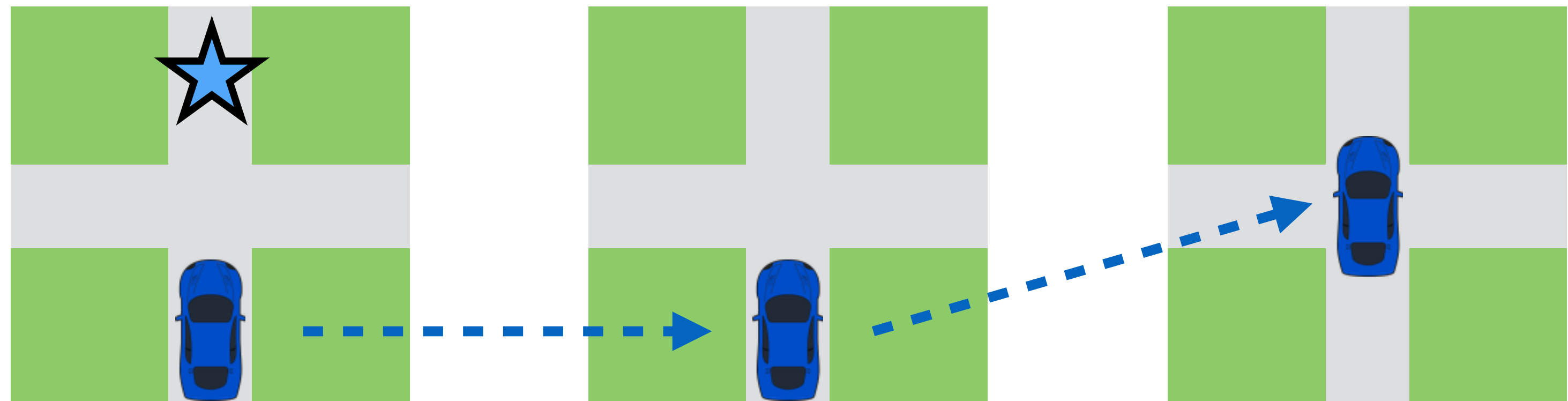
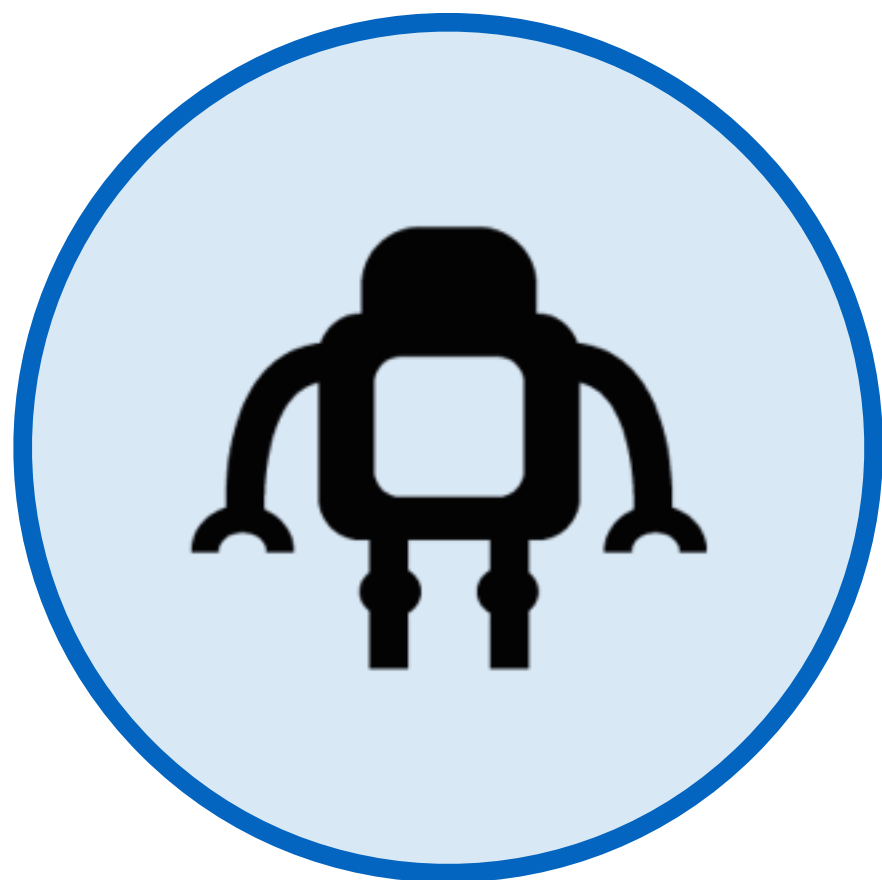
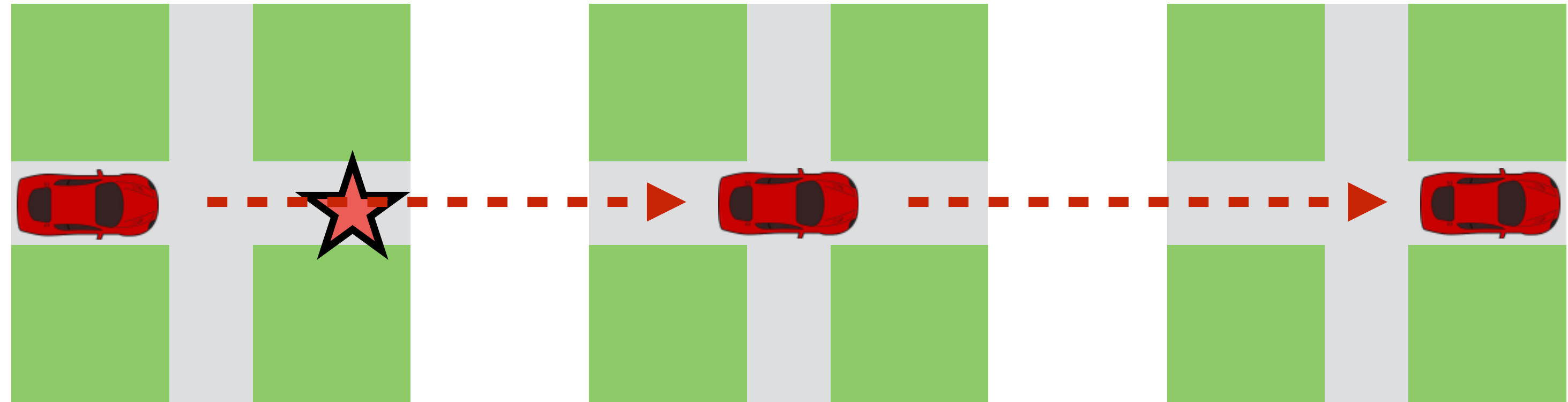
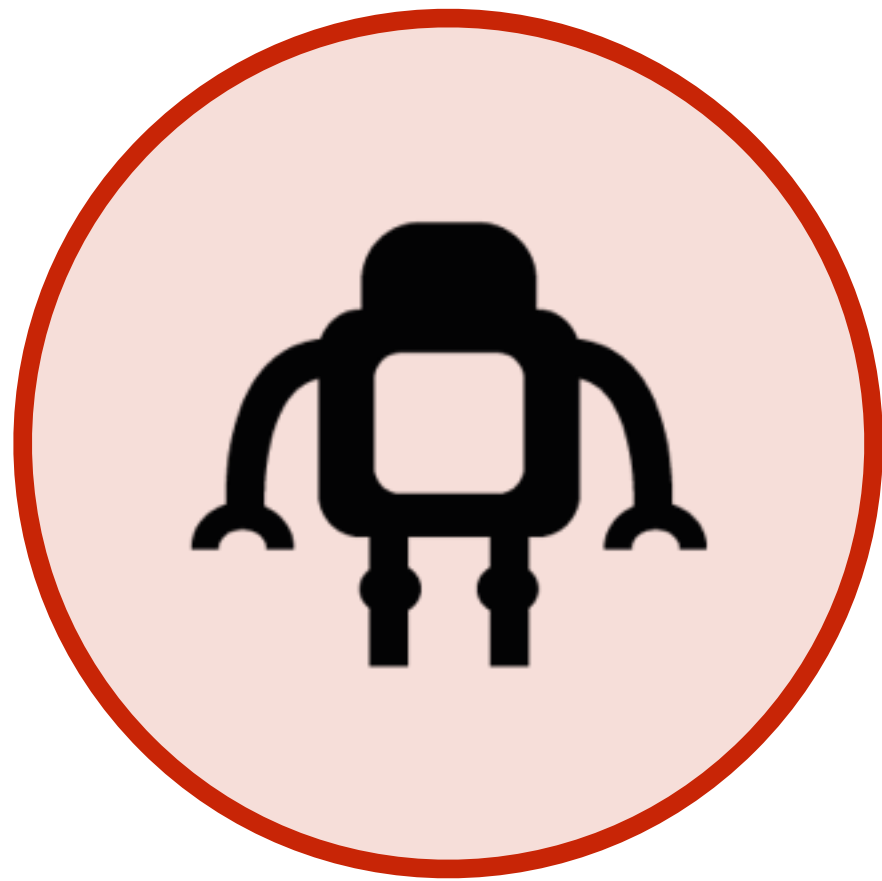
Translating Neuralese

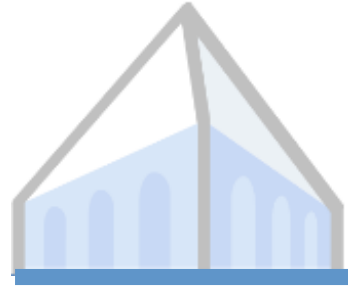


Jacob Andreas, Anca Drăgan and Dan Klein

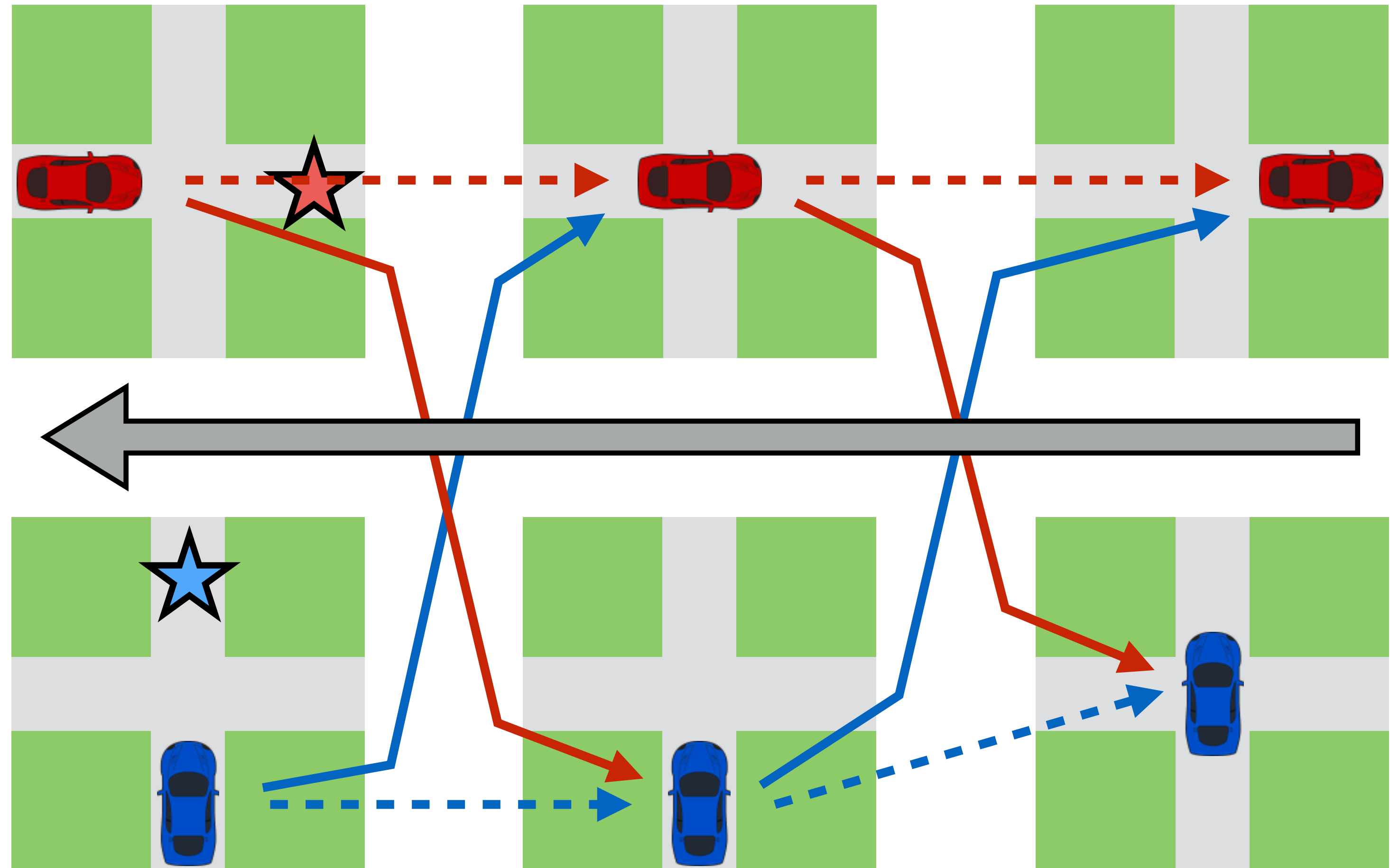
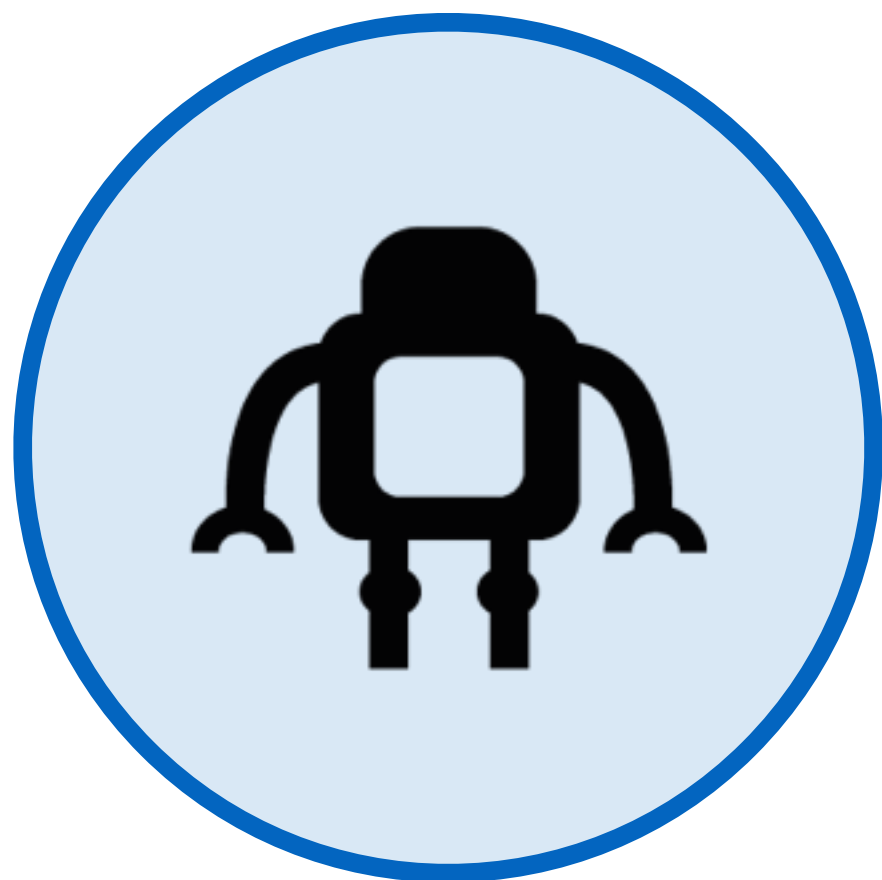
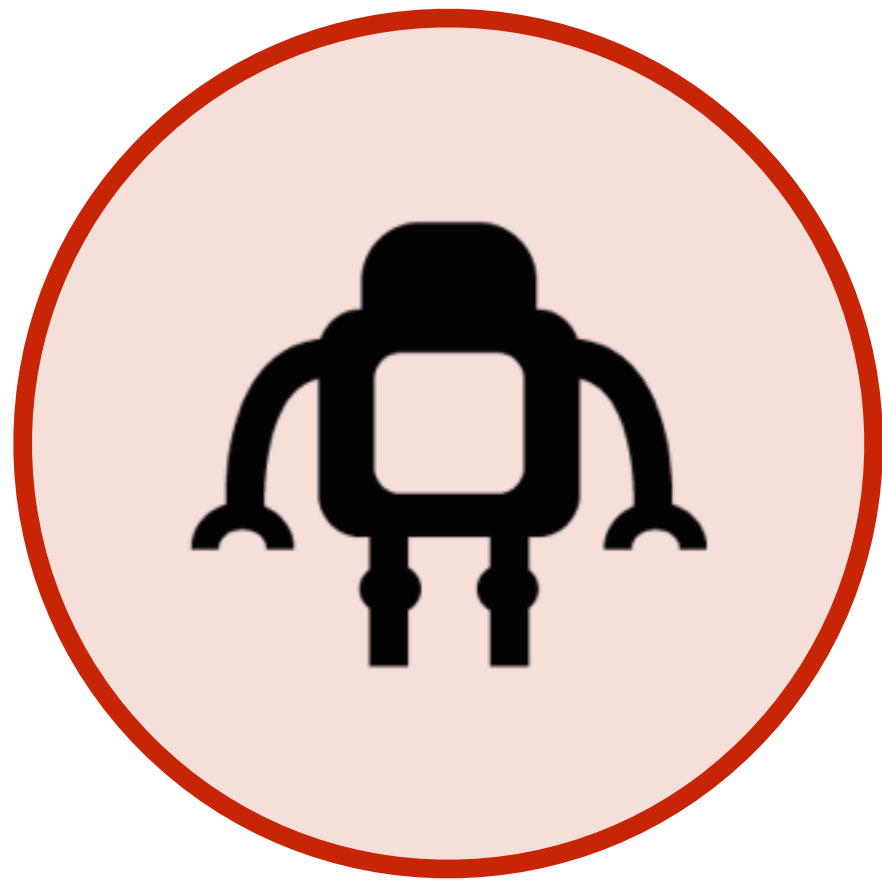


Learning to Communicate



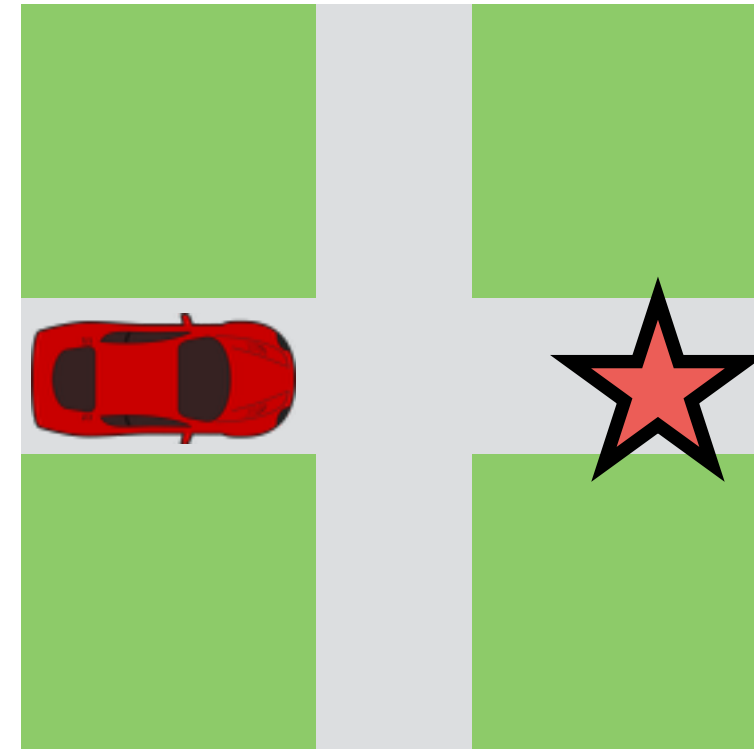
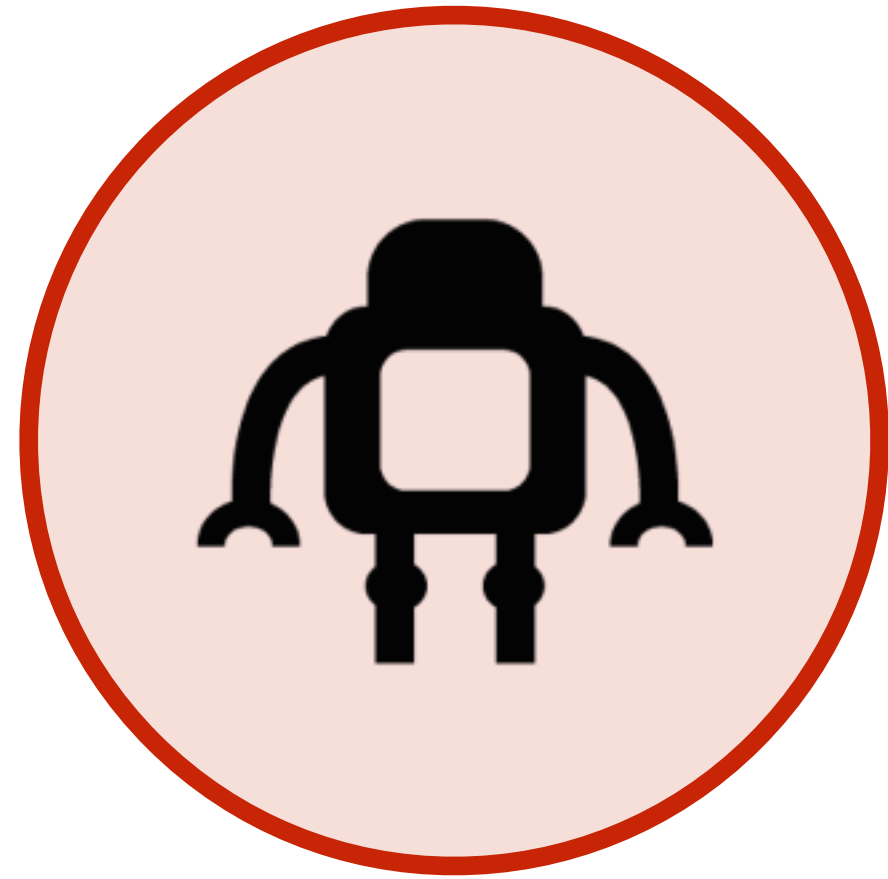


Learning to Communicate

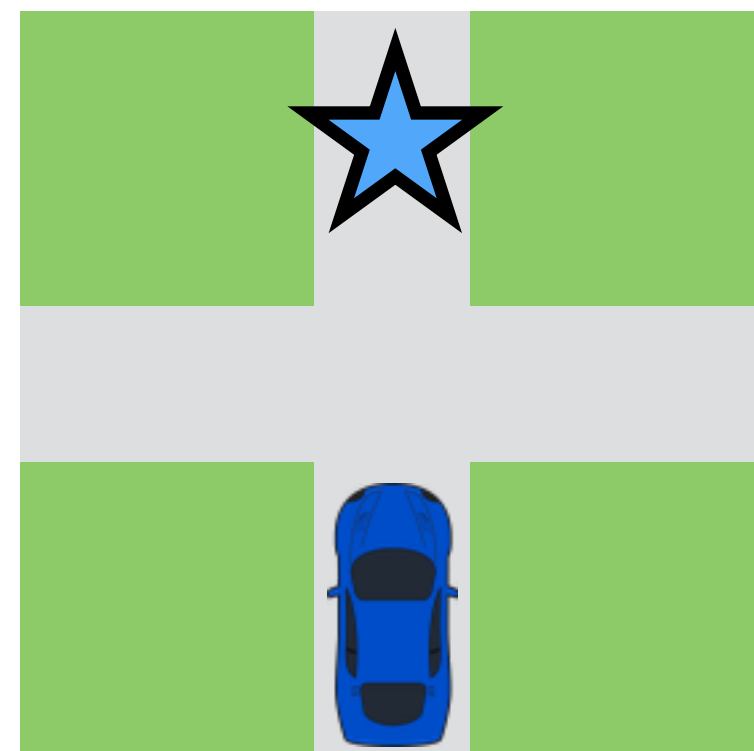
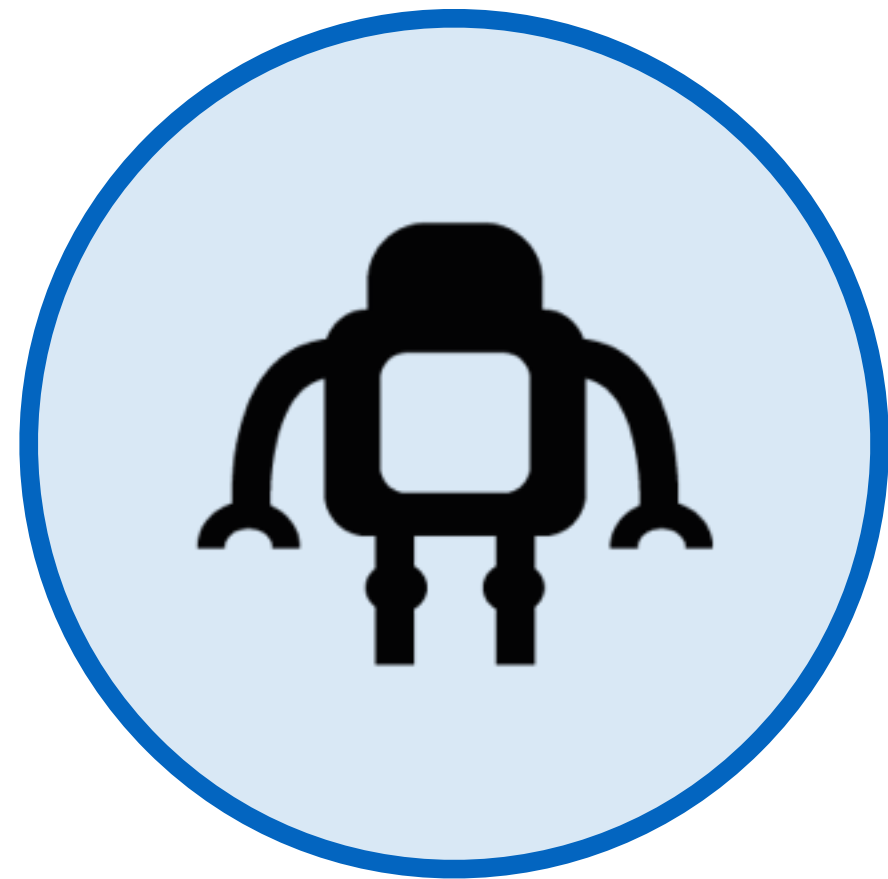
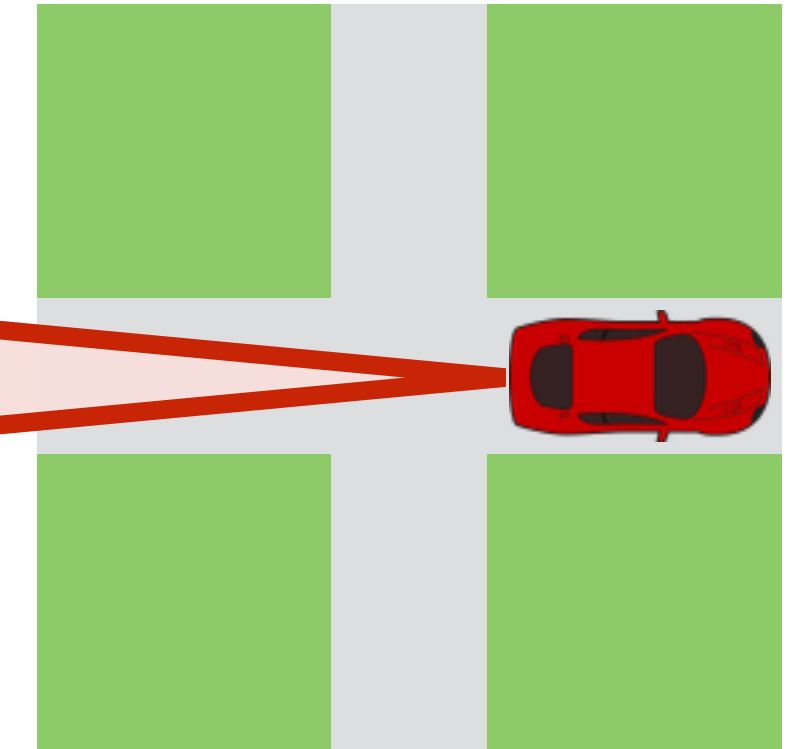




Neurales

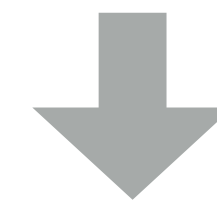
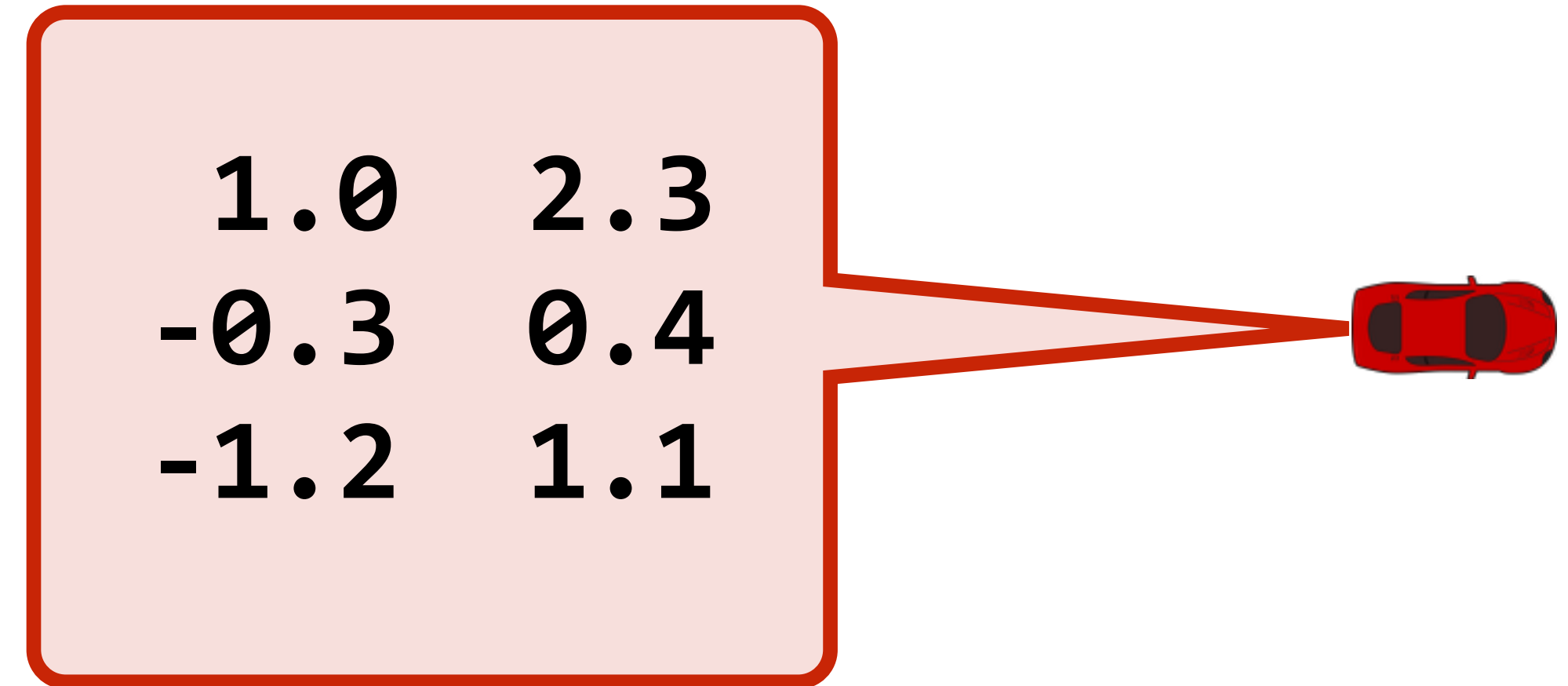
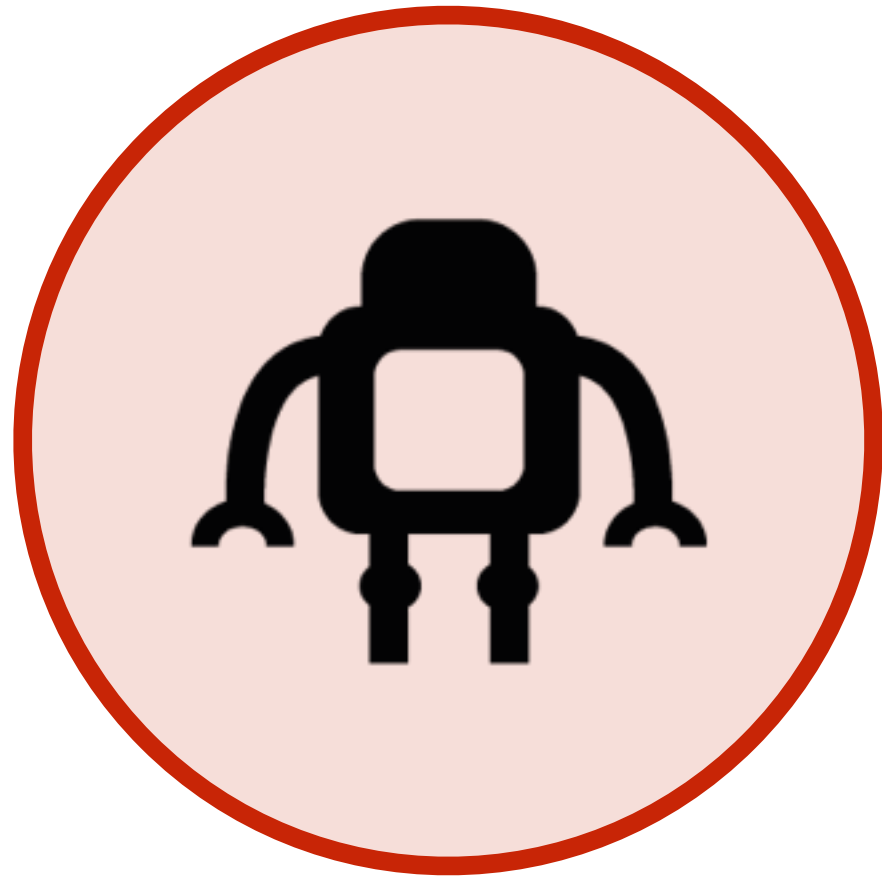


1.0	2.3
-0.3	0.4
-1.2	1.1

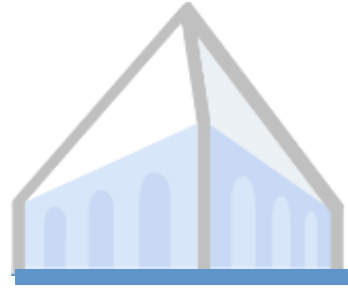




Translating neuralesse

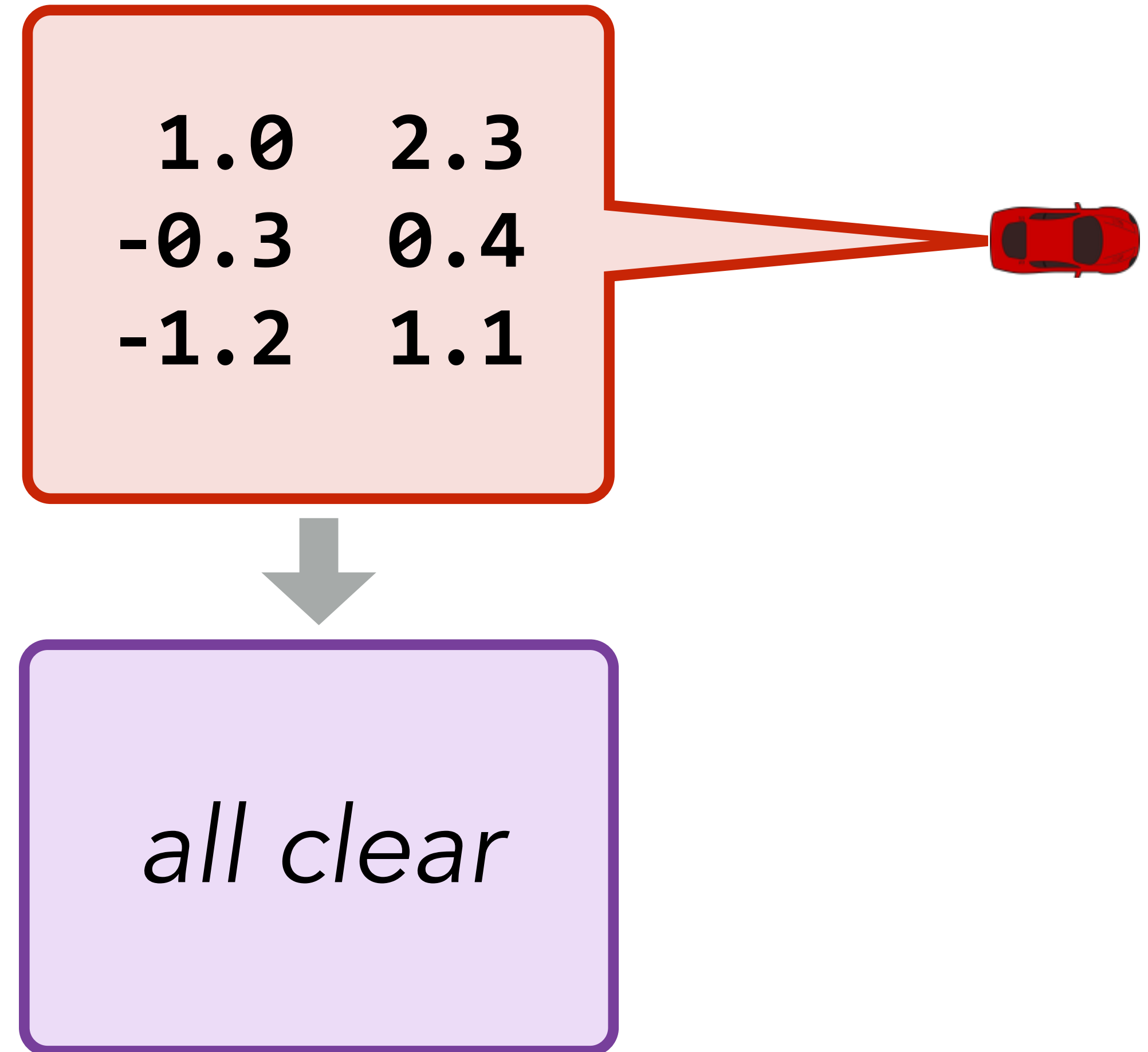


all clear



Translating neuralese

- **Interoperate** with autonomous systems
- **Diagnose** errors
- **Learn** from solutions





Outline

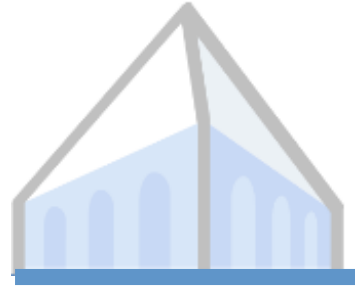
Natural language & neuralese

Statistical machine translation

Semantic machine translation

Implementation details

Evaluation



Outline

Natural language & neuralese

Statistical machine translation

Semantic machine translation

Implementation details

Evaluation



Outline

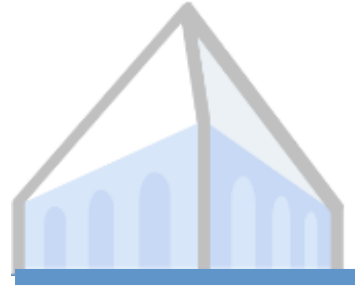
Natural language & neuralese

Statistical machine translation

Semantic machine translation

Implementation details

Evaluation



Outline

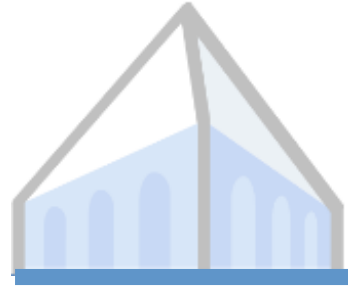
Natural language & neuralese

Statistical machine translation

Semantic machine translation

Implementation details

Evaluation



Outline

Natural language & neuralese

Statistical machine translation

Semantic machine translation

Implementation details

Evaluation



Outline

Natural language & neuralese

Statistical machine translation

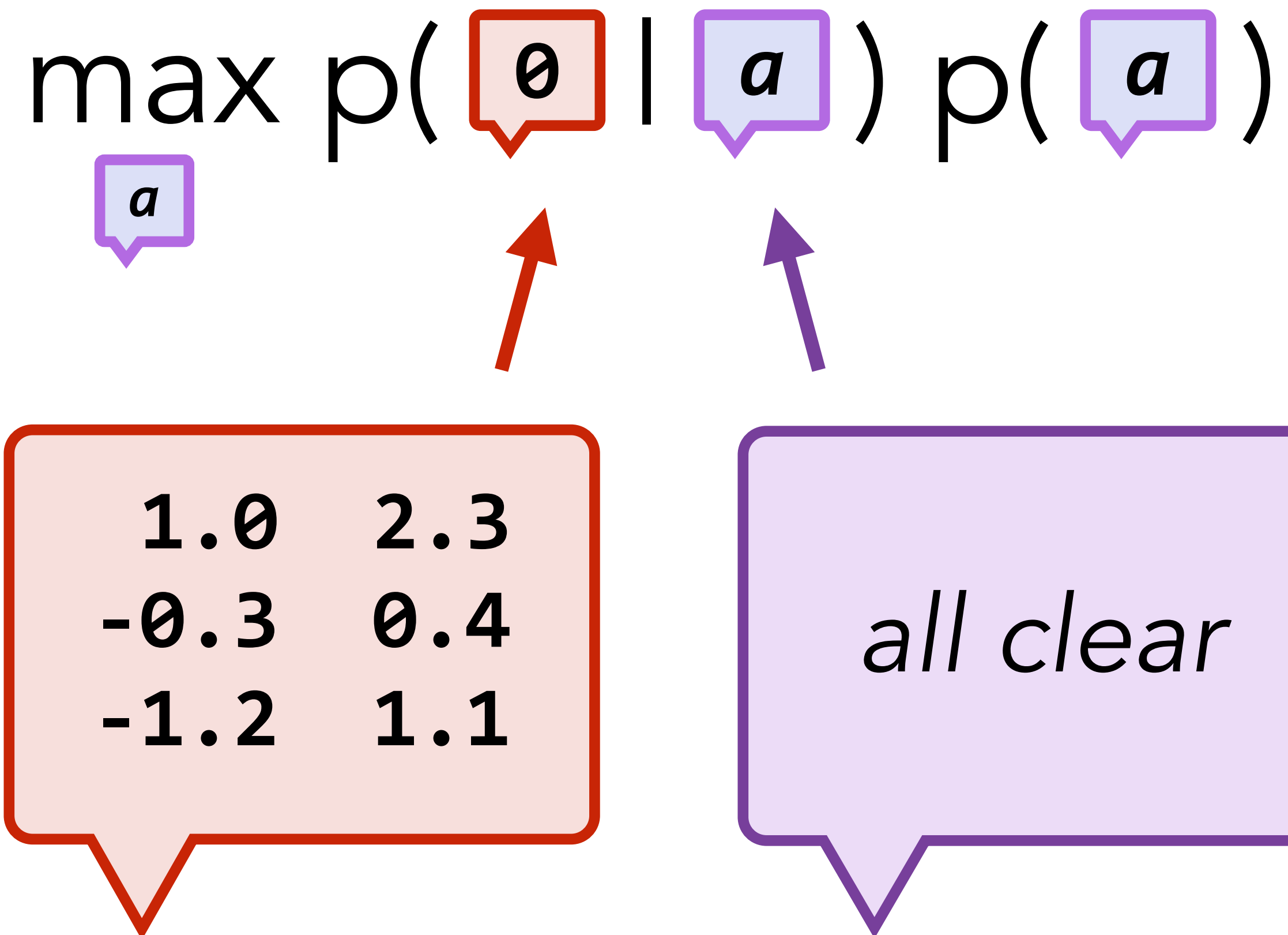
Semantic machine translation

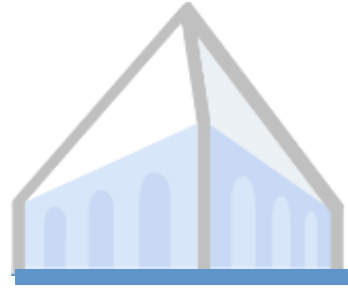
Implementation details

Evaluation



A statistical MT problem





A statistical MT problem



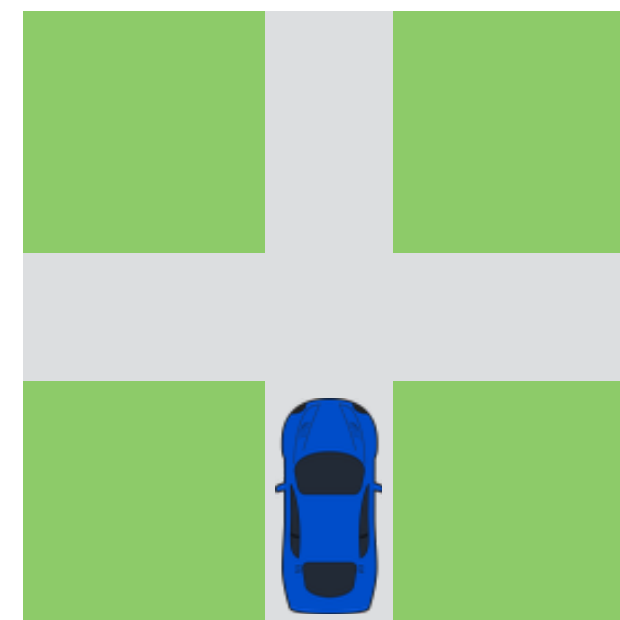
How do we induce a translation model?



A statistical MT problem

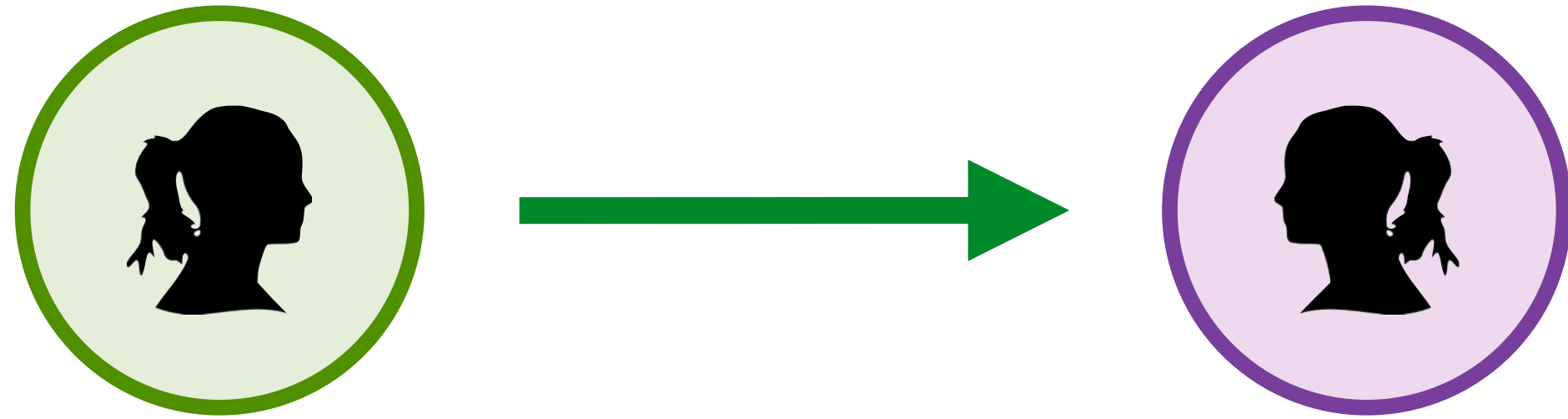
$$\max_a p(\theta | a) p(a)$$

$$\propto \max_a \sum_{\text{map}} p(\theta | \text{map}) p(a | \text{map}) p(\text{map})$$





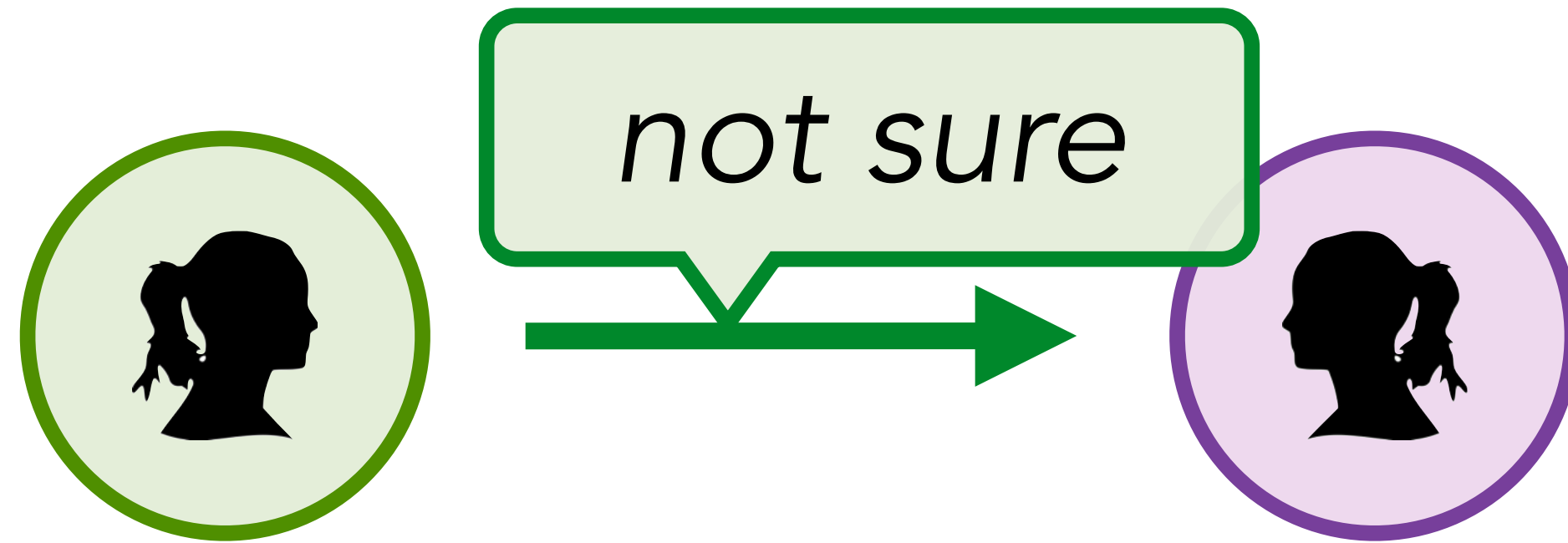
Strategy mismatch



$$\zeta(s) = \frac{1}{\Gamma(s)} \int_0^{\infty} \frac{1}{e^x - 1} x^s \frac{dx}{x}$$



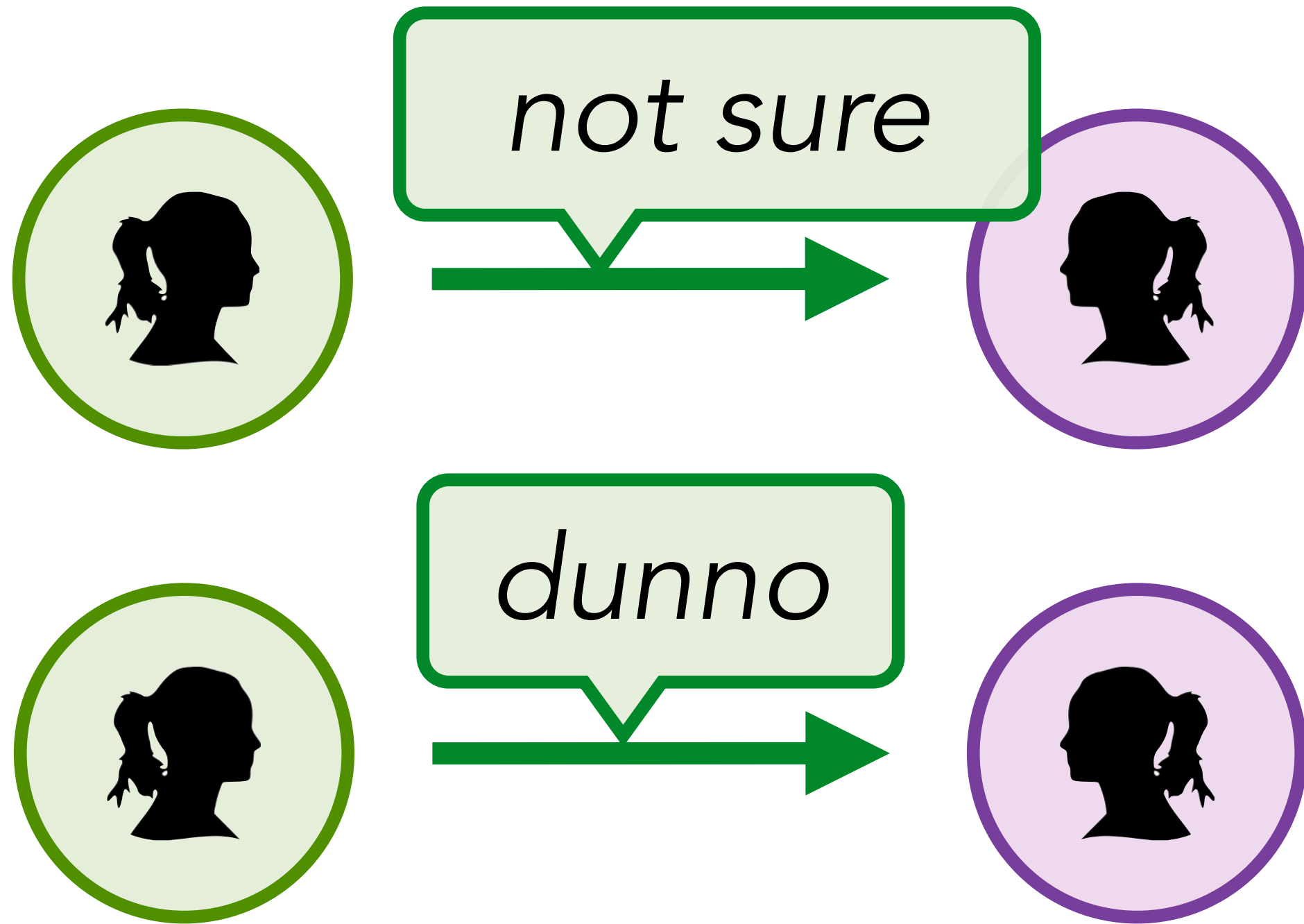
Strategy mismatch



$$\zeta(s) = \frac{1}{\Gamma(s)} \int_0^{\infty} \frac{1}{e^x - 1} x^s \frac{dx}{x}$$

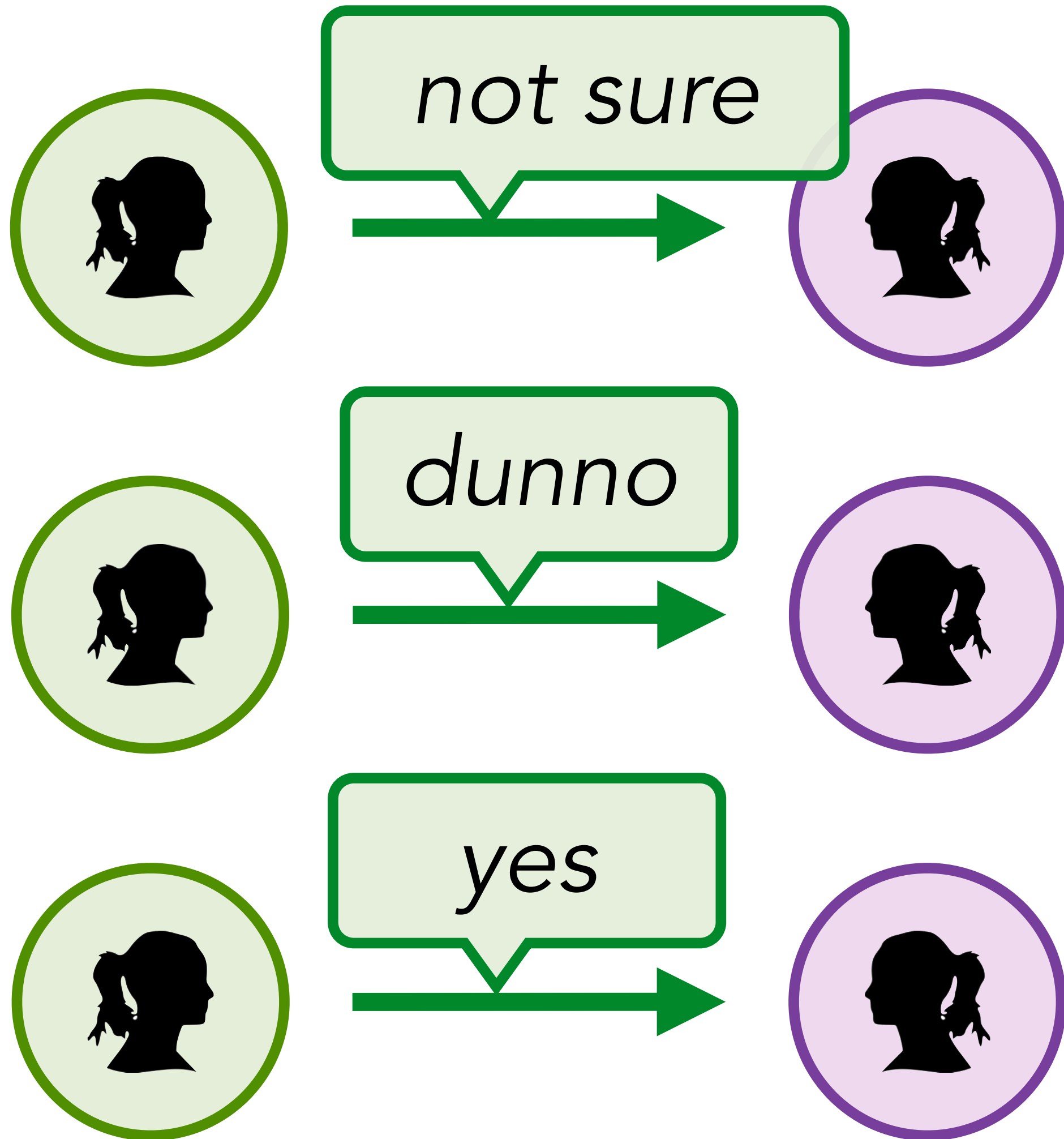


Strategy mismatch



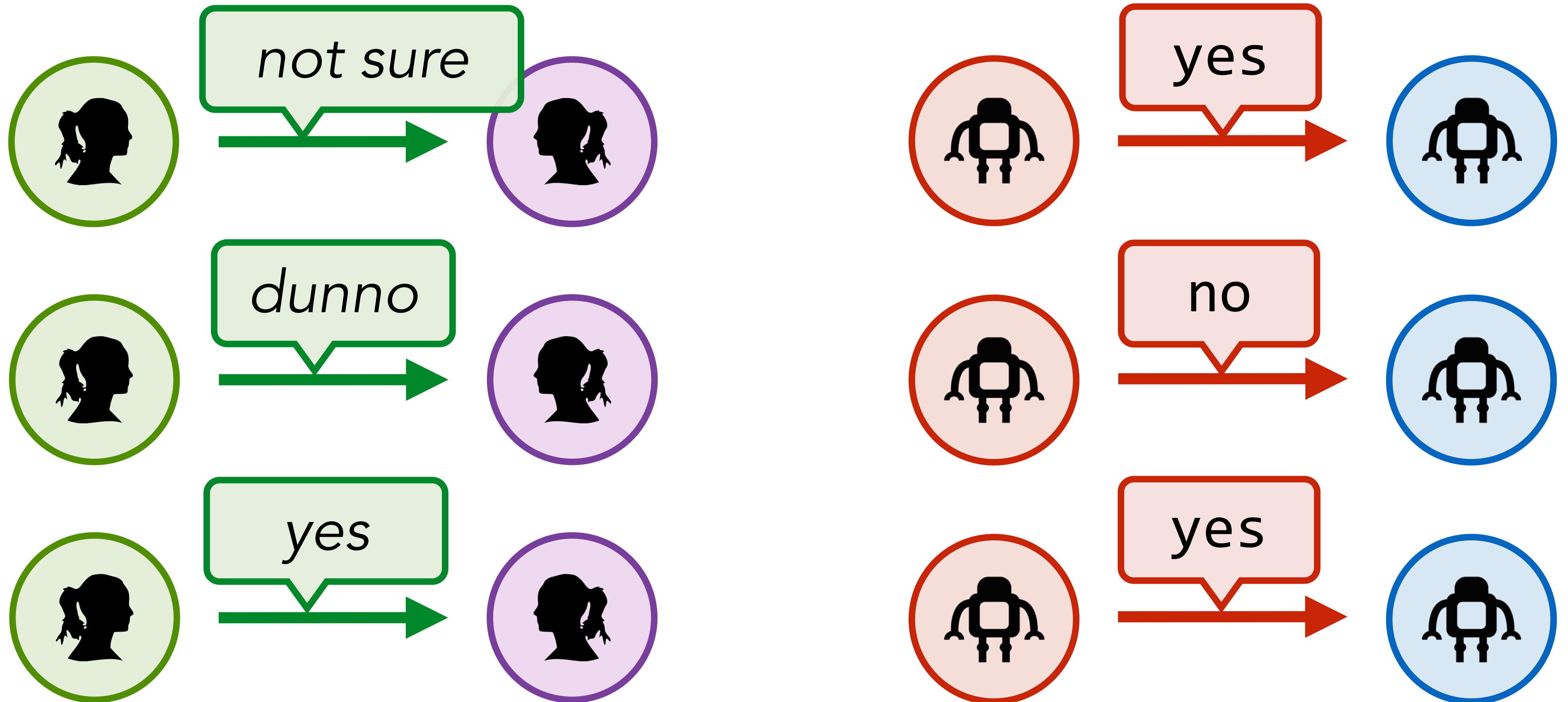


Strategy mismatch





Strategy mismatch

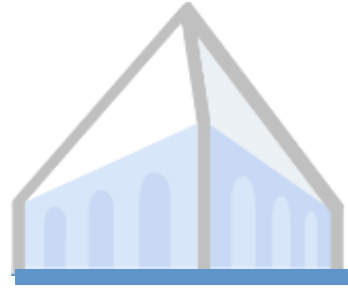




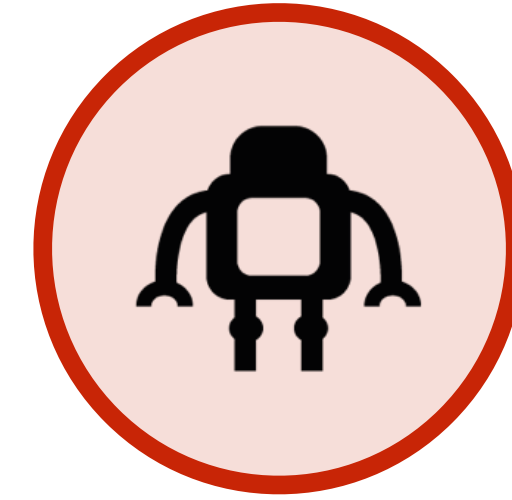
Strategy mismatch



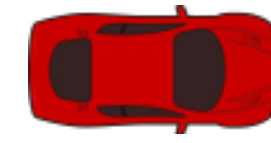
$$\sum p(\theta, \text{map icon} \mid \text{not sure}) p(\text{not sure})$$



Stat MT criterion doesn't capture meaning

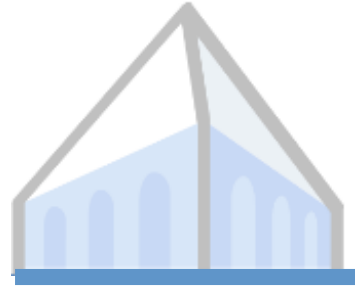


*In the
intersection*



moving
 $(0, 3) \rightarrow (1, 4)$





Outline

Natural language & neuralese

X **Statistical** machine translation

Semantic machine translation

Implementation details

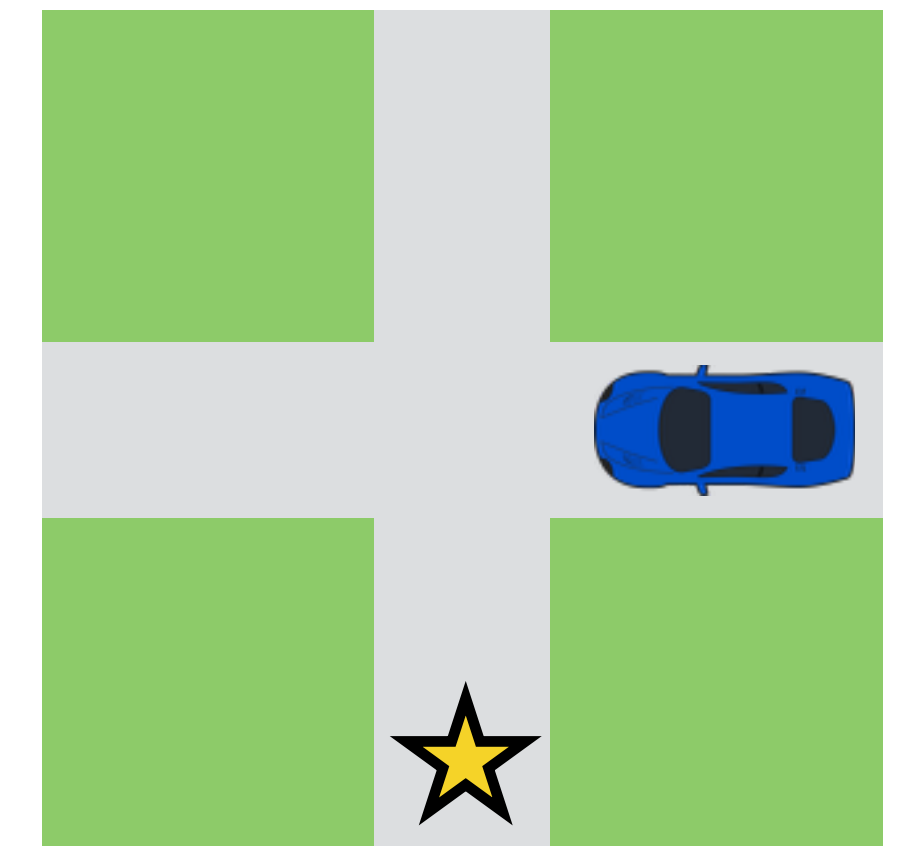
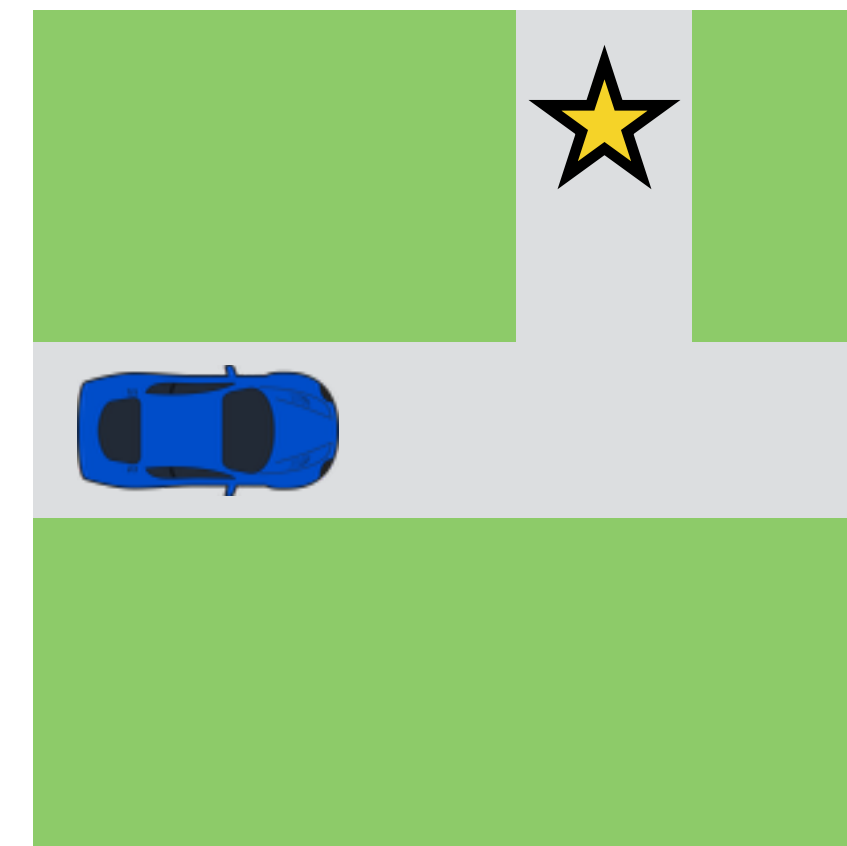
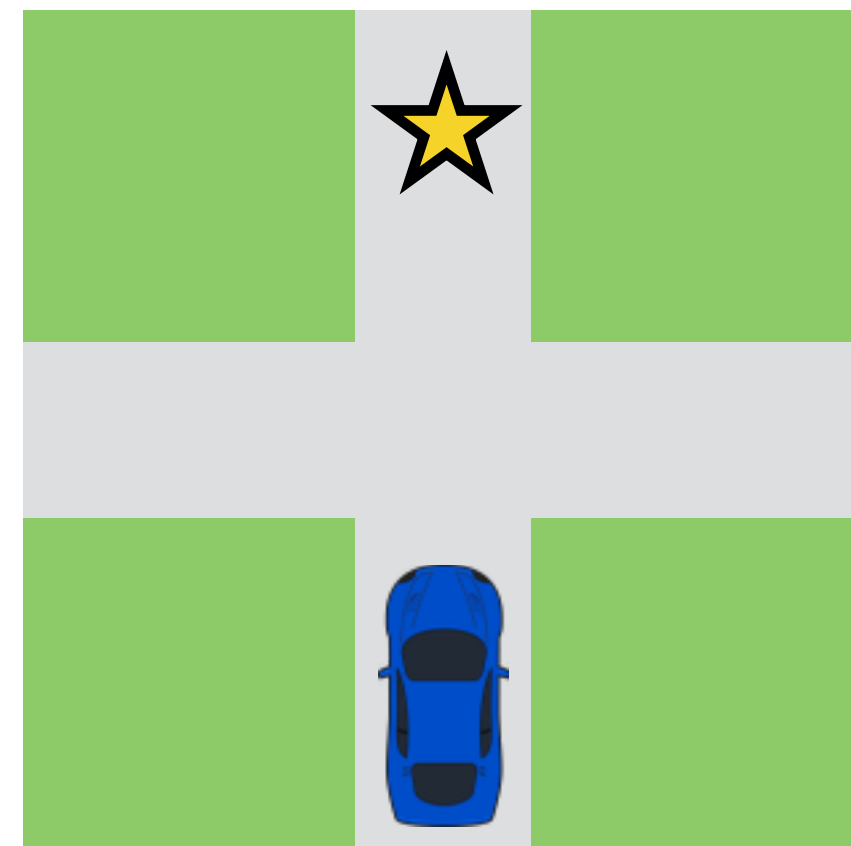
Evaluation



A "semantic MT" problem

The meaning of an utterance is given by its **truth conditions**

I'm going north

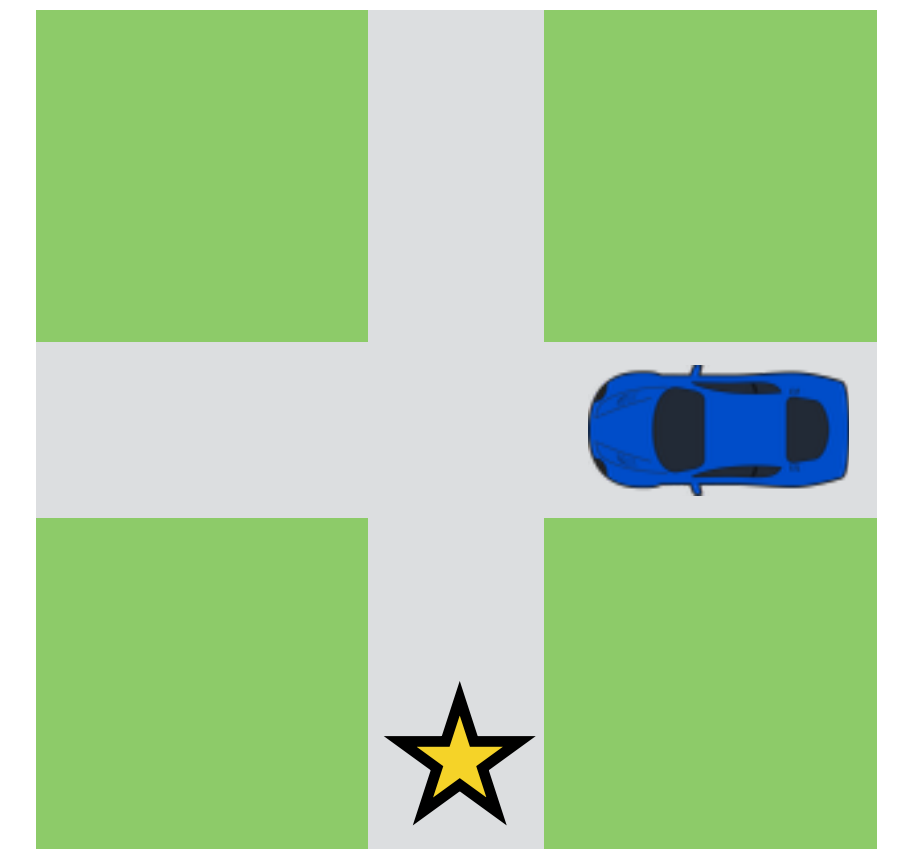
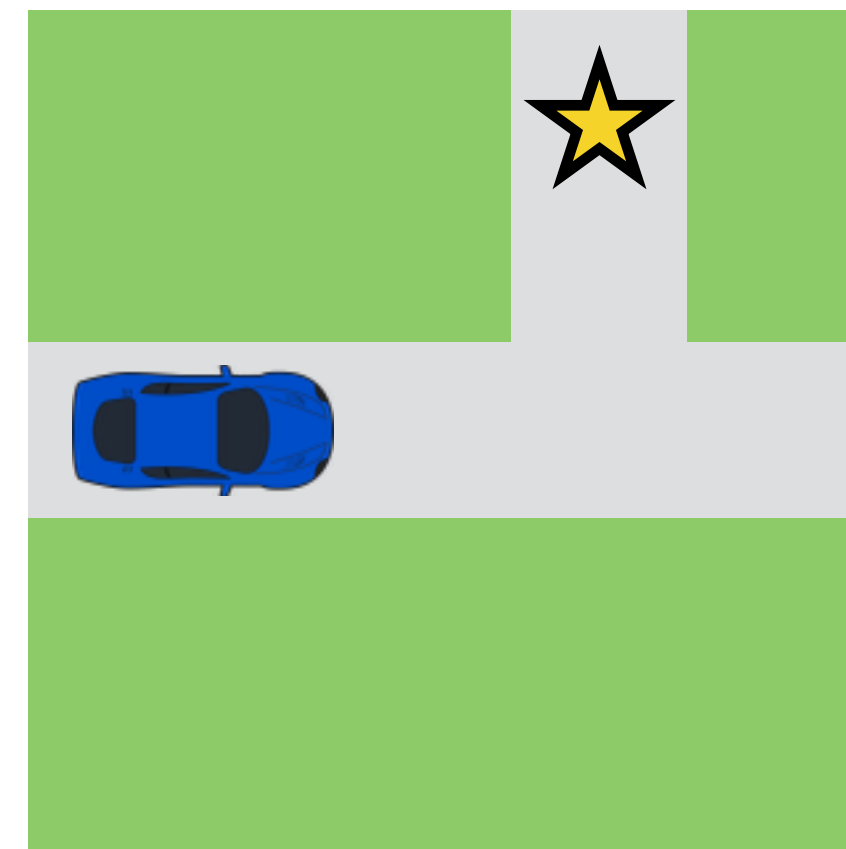
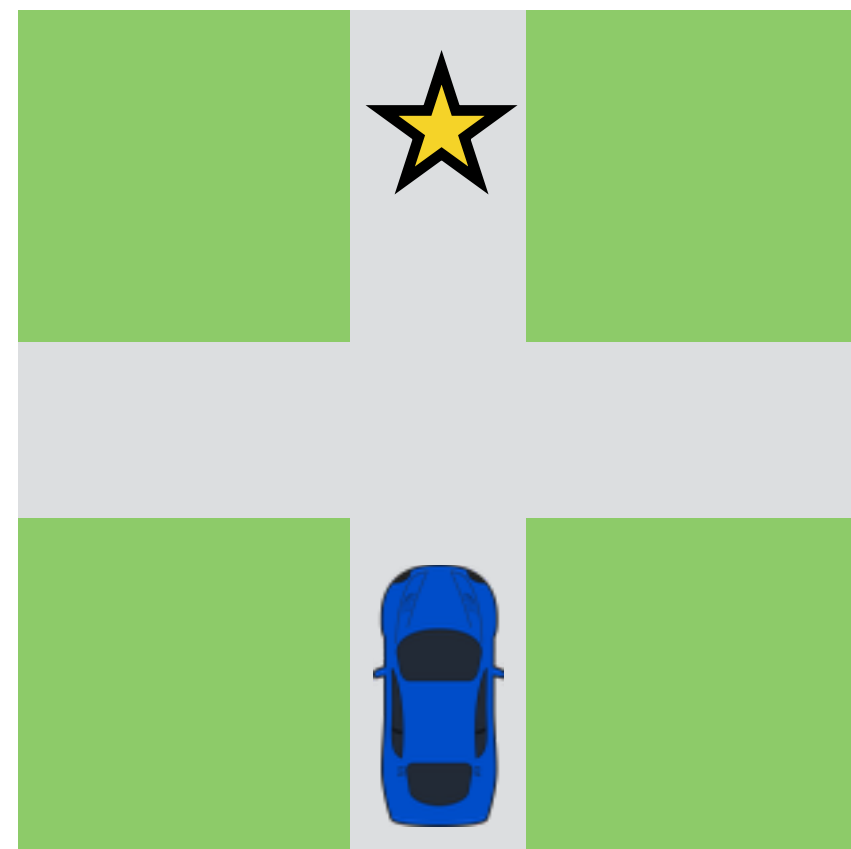


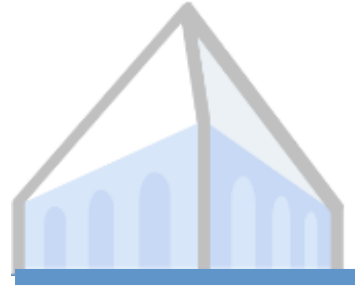


A "semantic MT" problem

The meaning of an utterance is given by its **truth conditions**

I'm going north

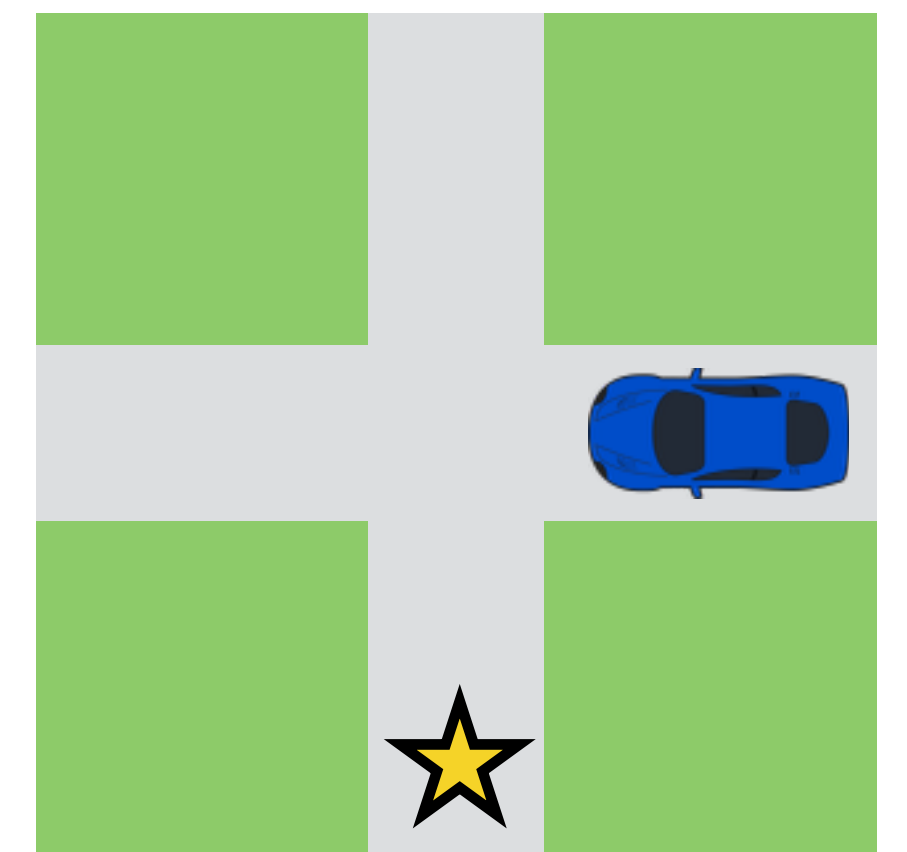
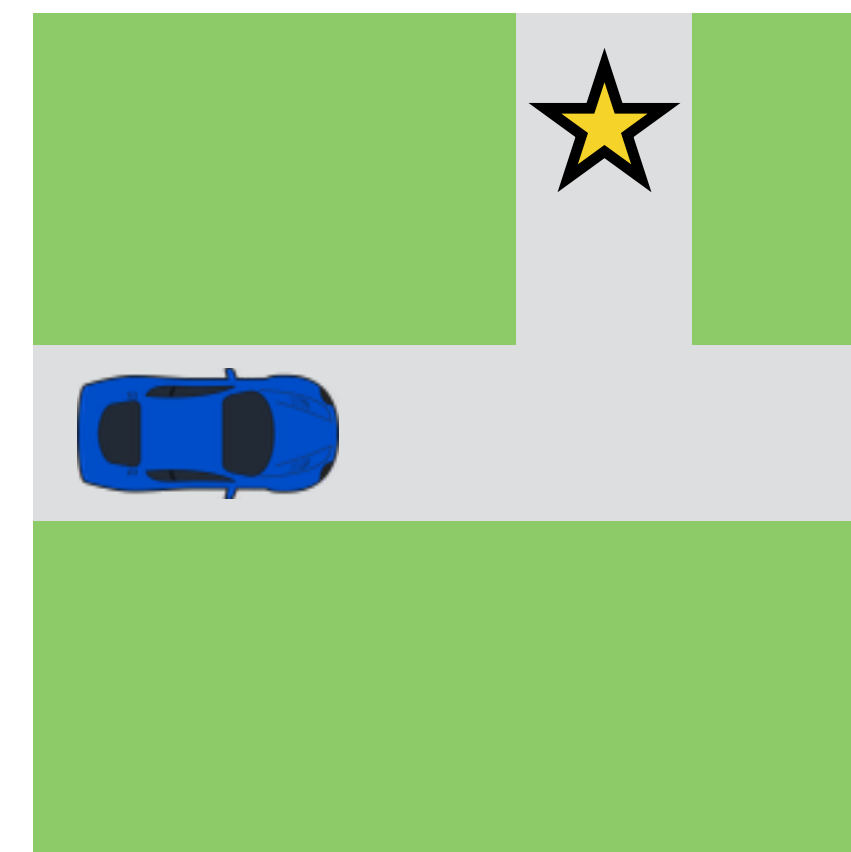
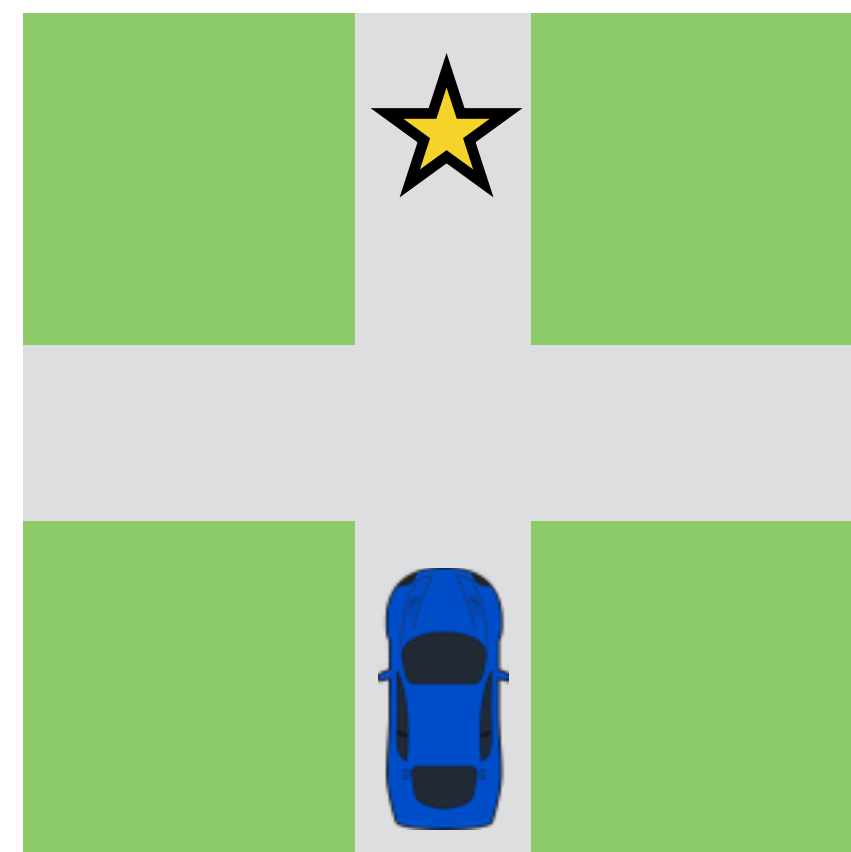




A "semantic MT" problem

The meaning of an utterance is given by its **truth conditions**

I'm going north



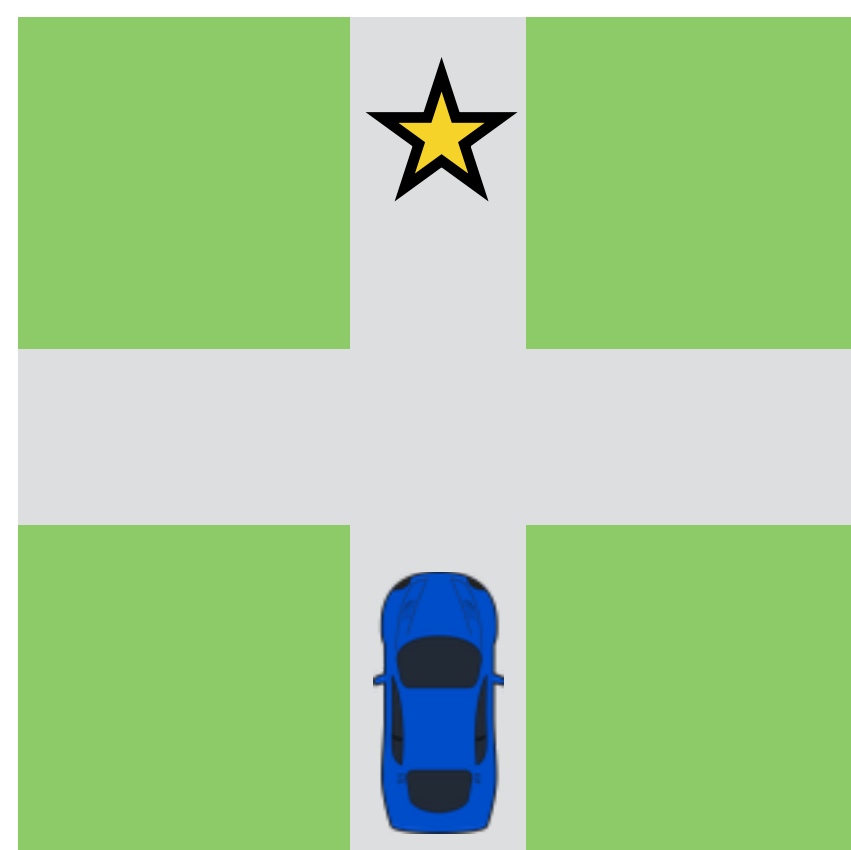
(loc (goal blue) north)



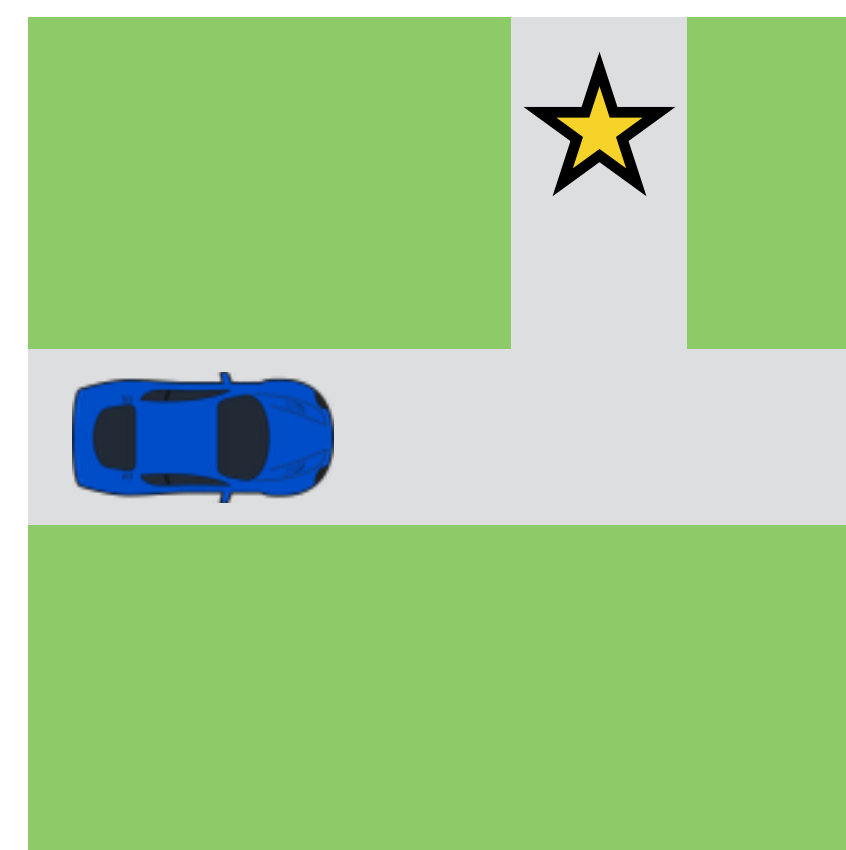
A "semantic MT" problem

The meaning of an utterance is given by its truth conditions
the **distribution over states** in which it is uttered

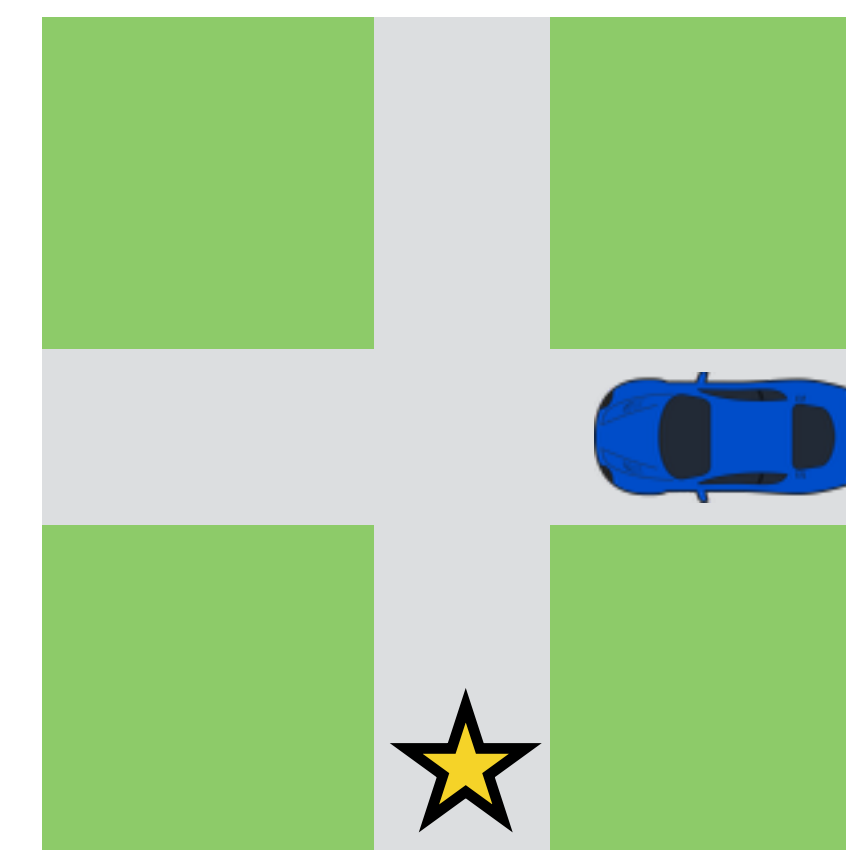
I'm going north



0.4



0.2



0.001



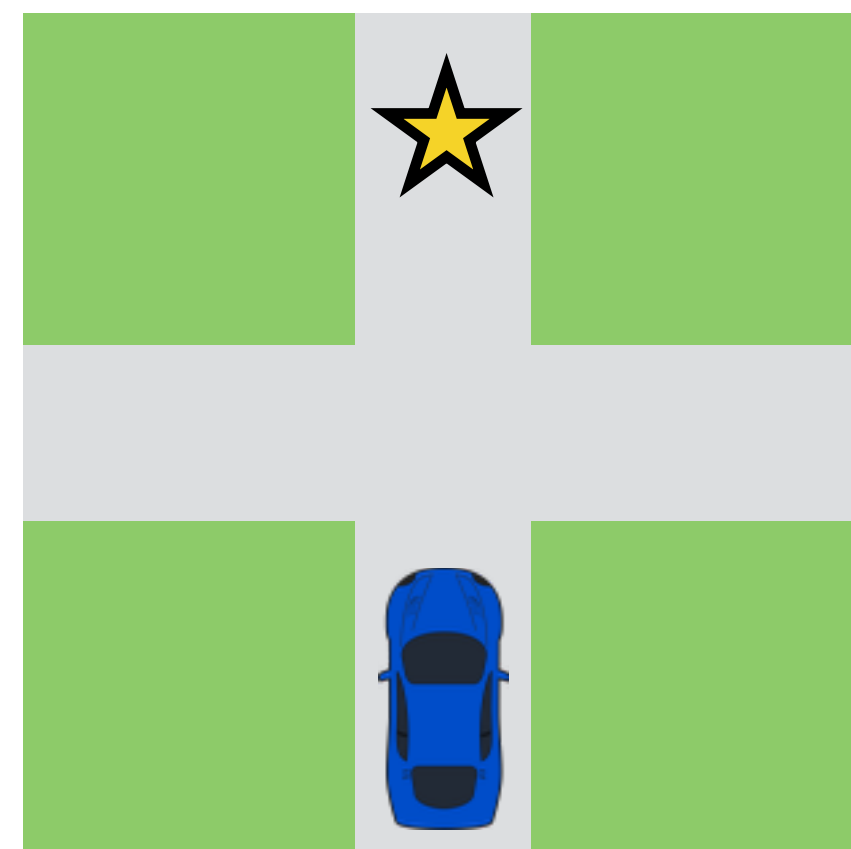
A "semantic MT" problem

The meaning of an utterance is given by its **truth conditions**

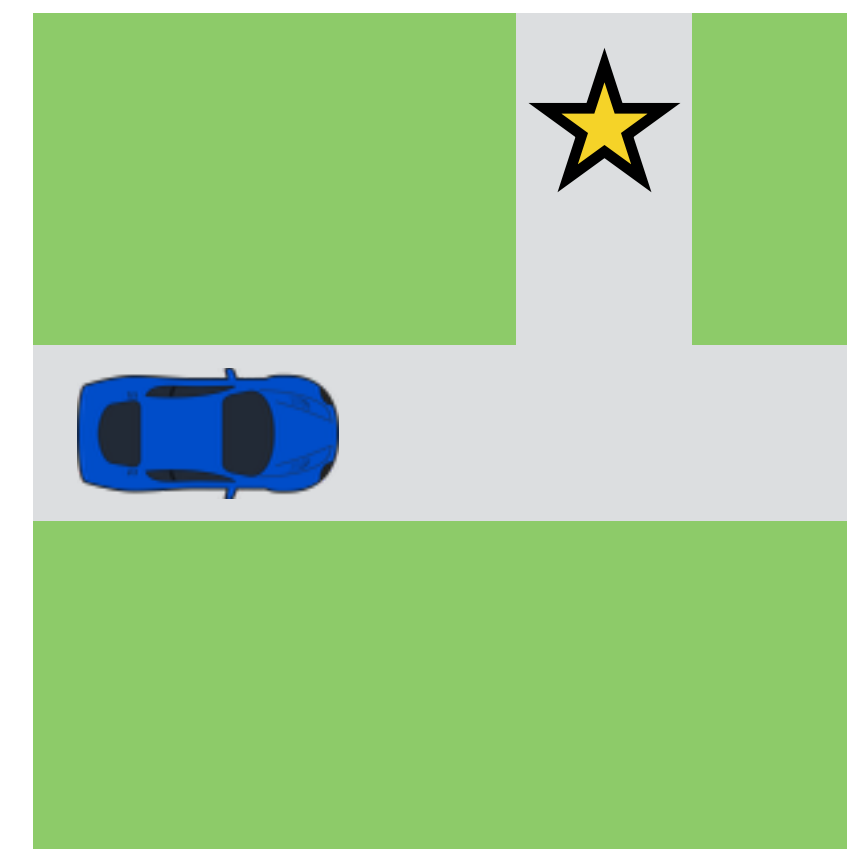
the **distribution over states** in which it is uttered

the **belief** it induces in listeners

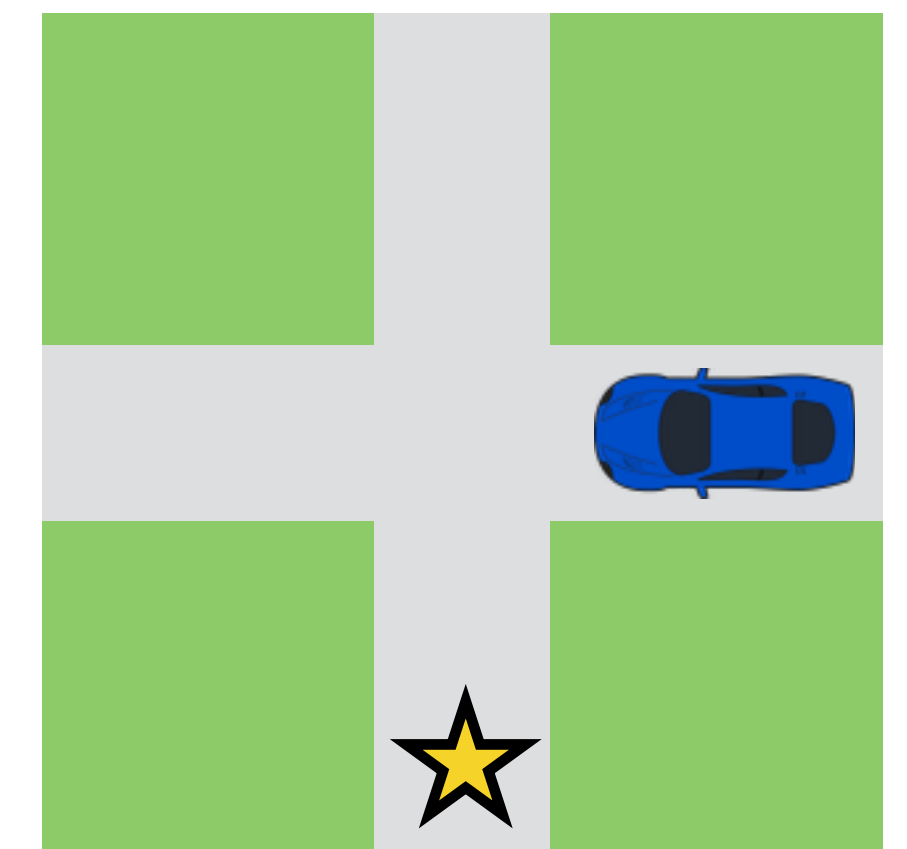
I'm going north



0.4



0.2



0.001



Representing meaning

The meaning of an utterance is given by

the **distribution over states** in which it is uttered

or equivalently, the **belief** it induces in listeners



Representing meaning

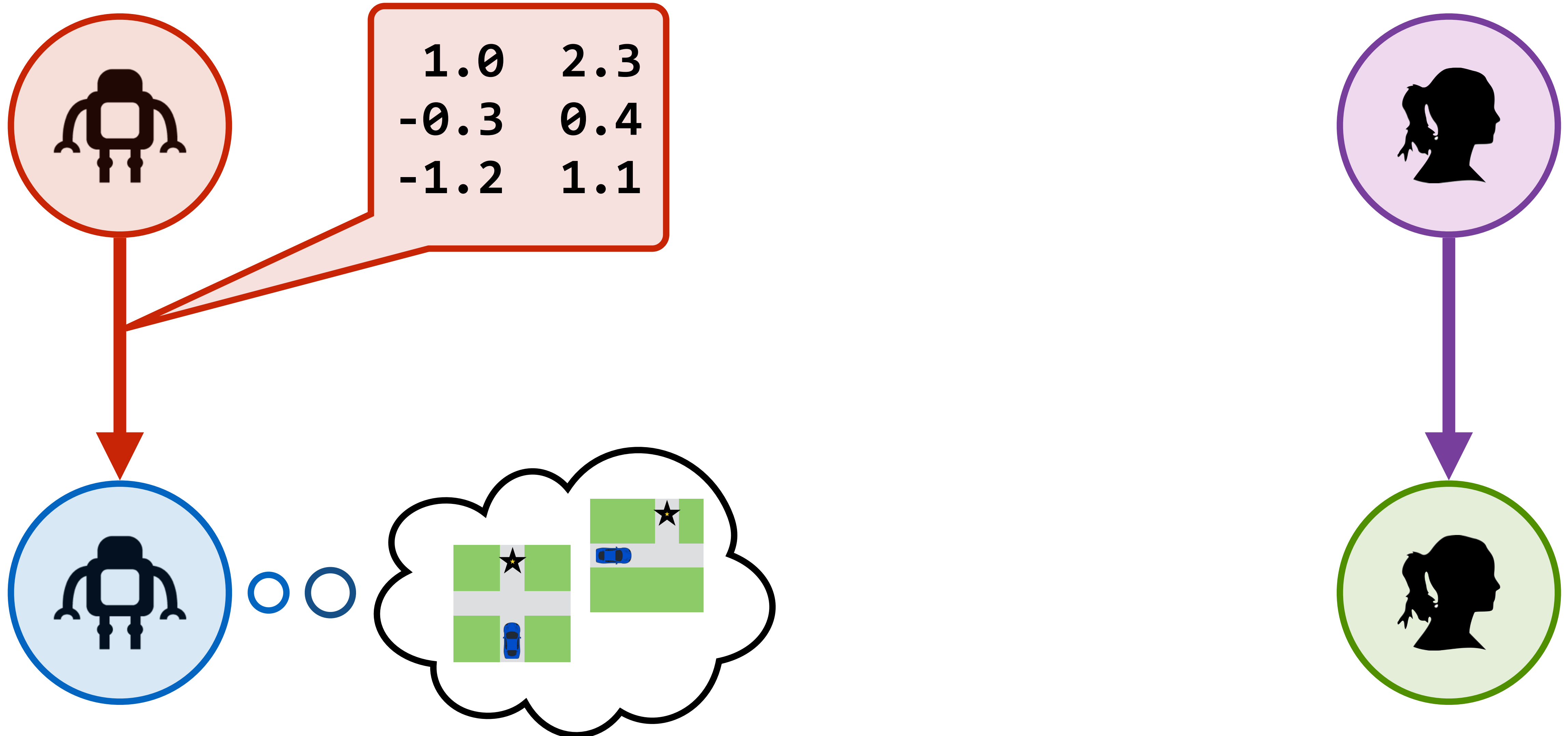
The meaning of an utterance is given by

the **distribution over states** in which it is uttered
or equivalently, the **belief** it induces in listeners

This distribution is well-defined even if the “utterance” is a vector rather than a sequence of tokens.

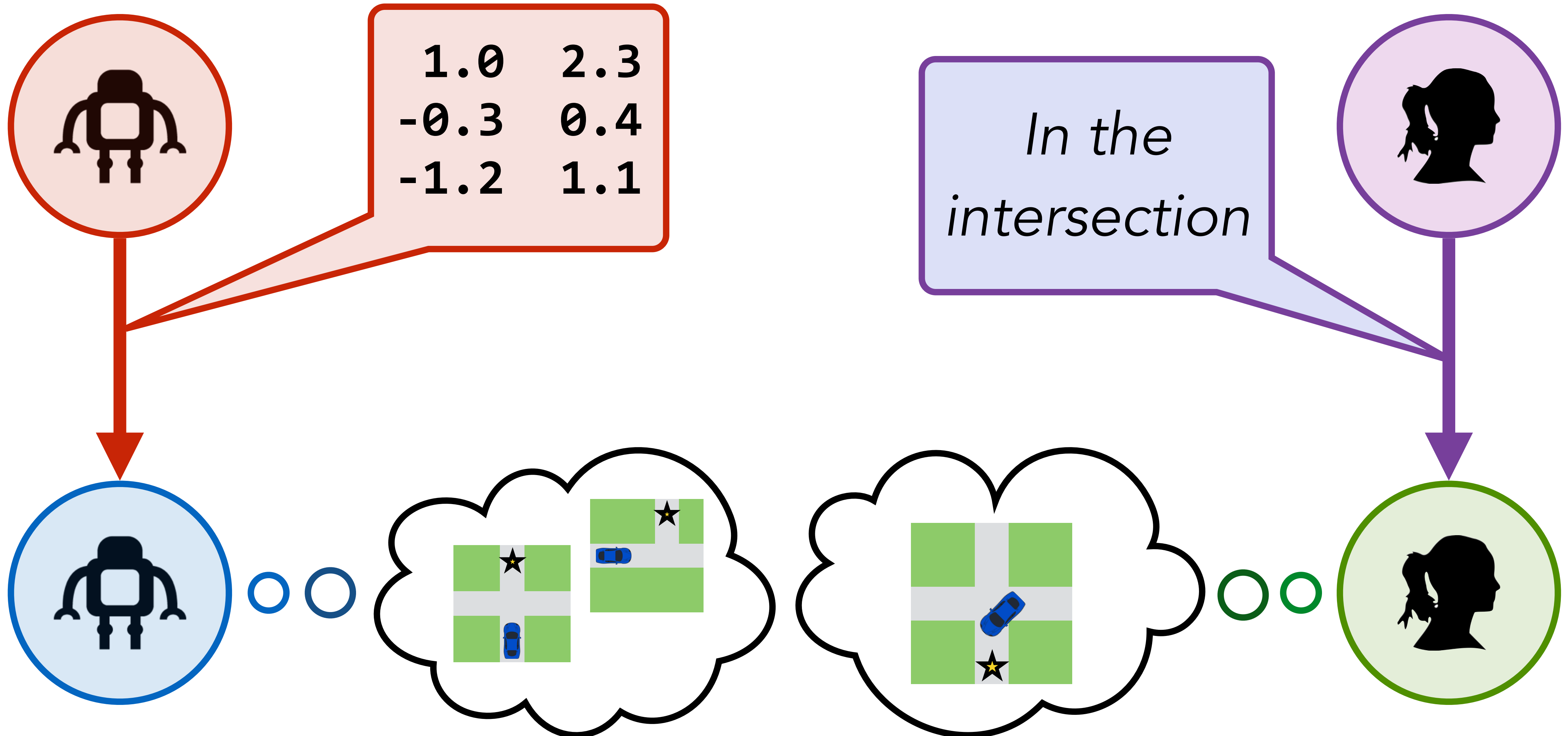


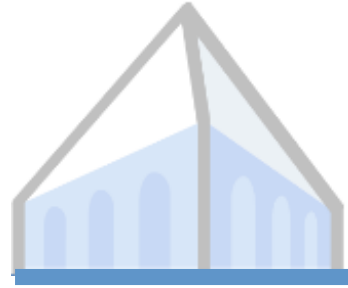
Translating with meaning



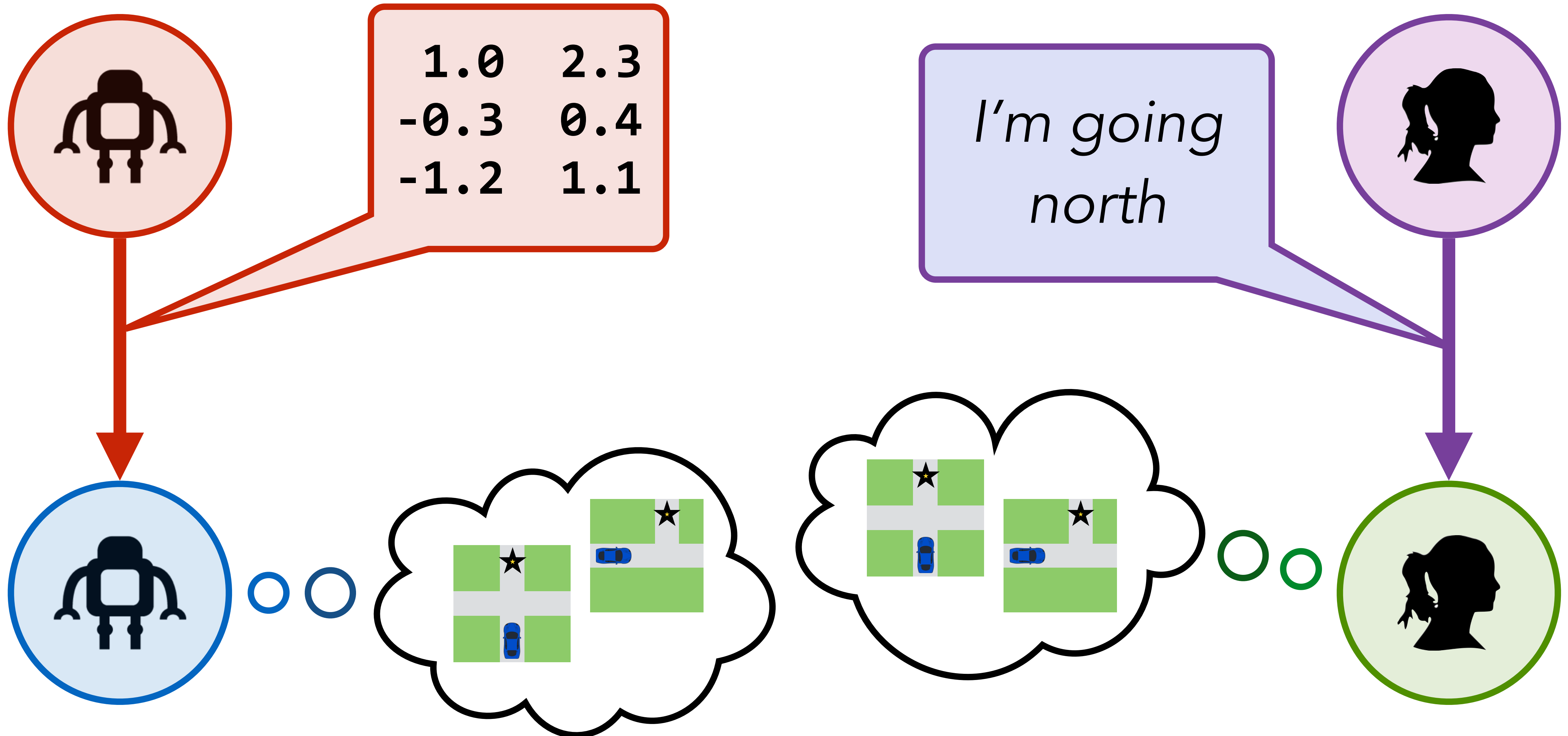


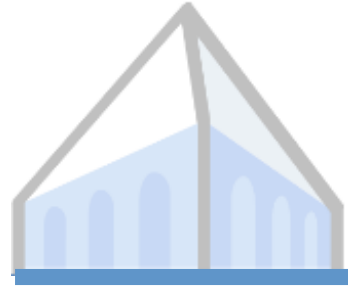
Translating with meaning



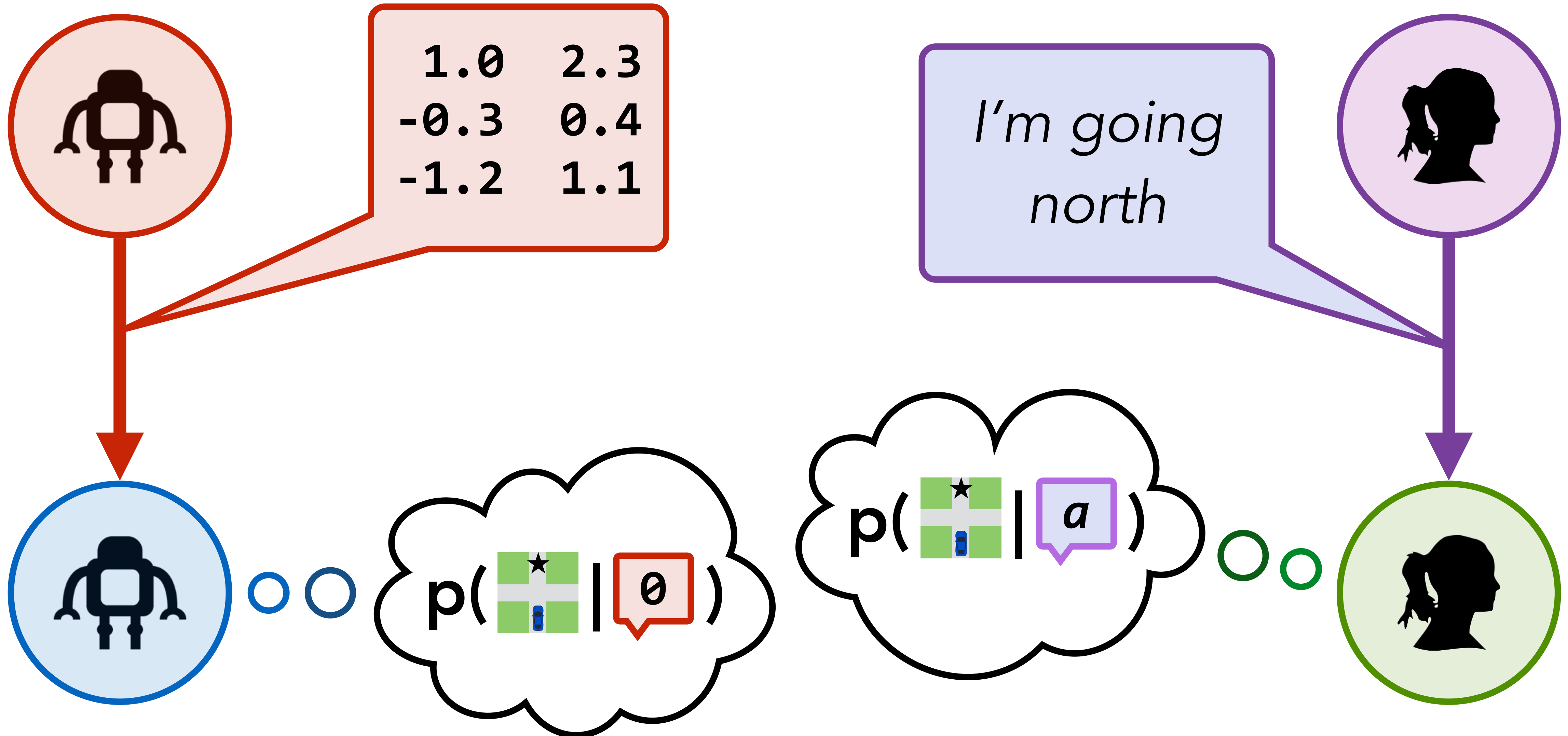


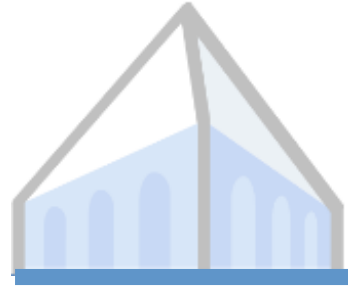
Translating with meaning



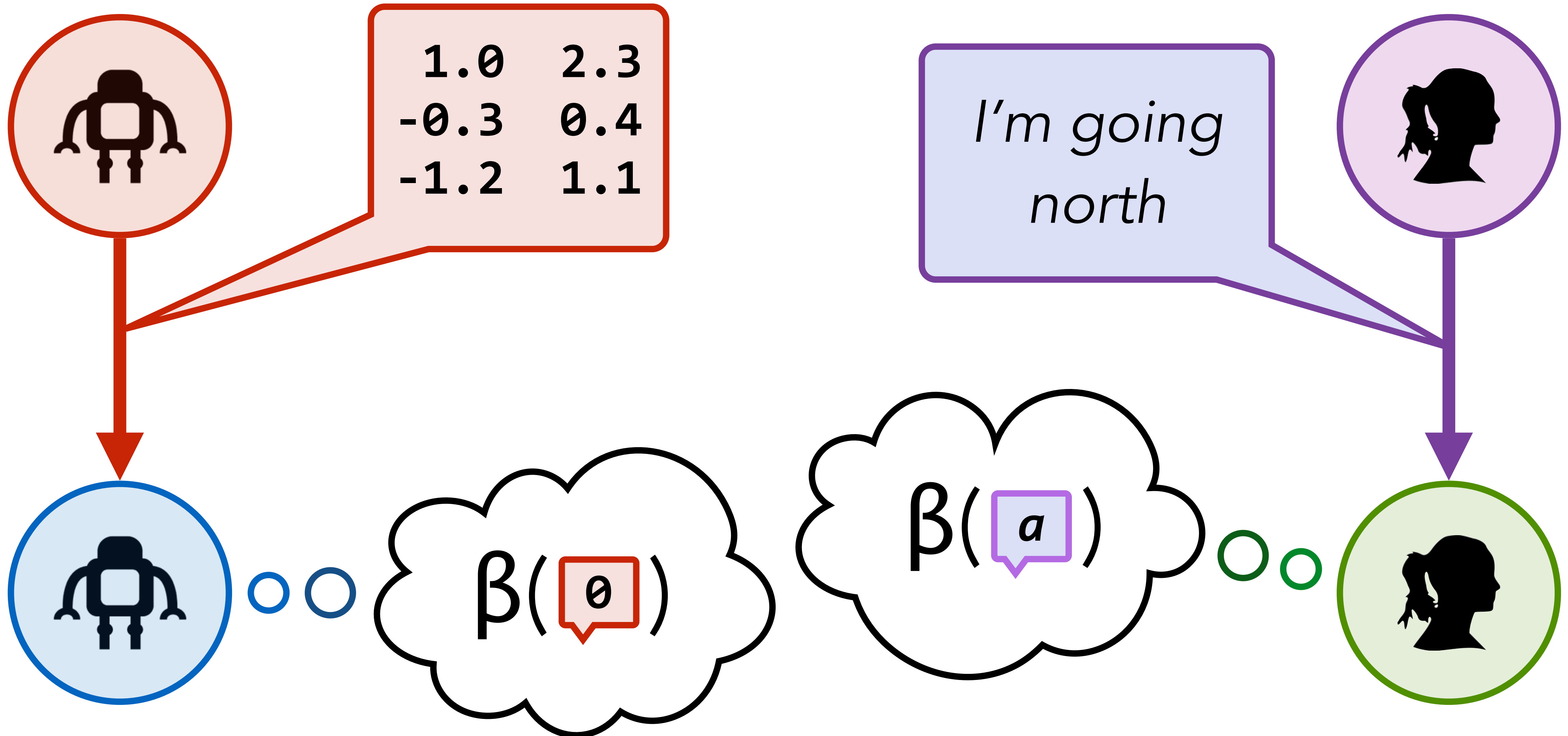


Translating with meaning



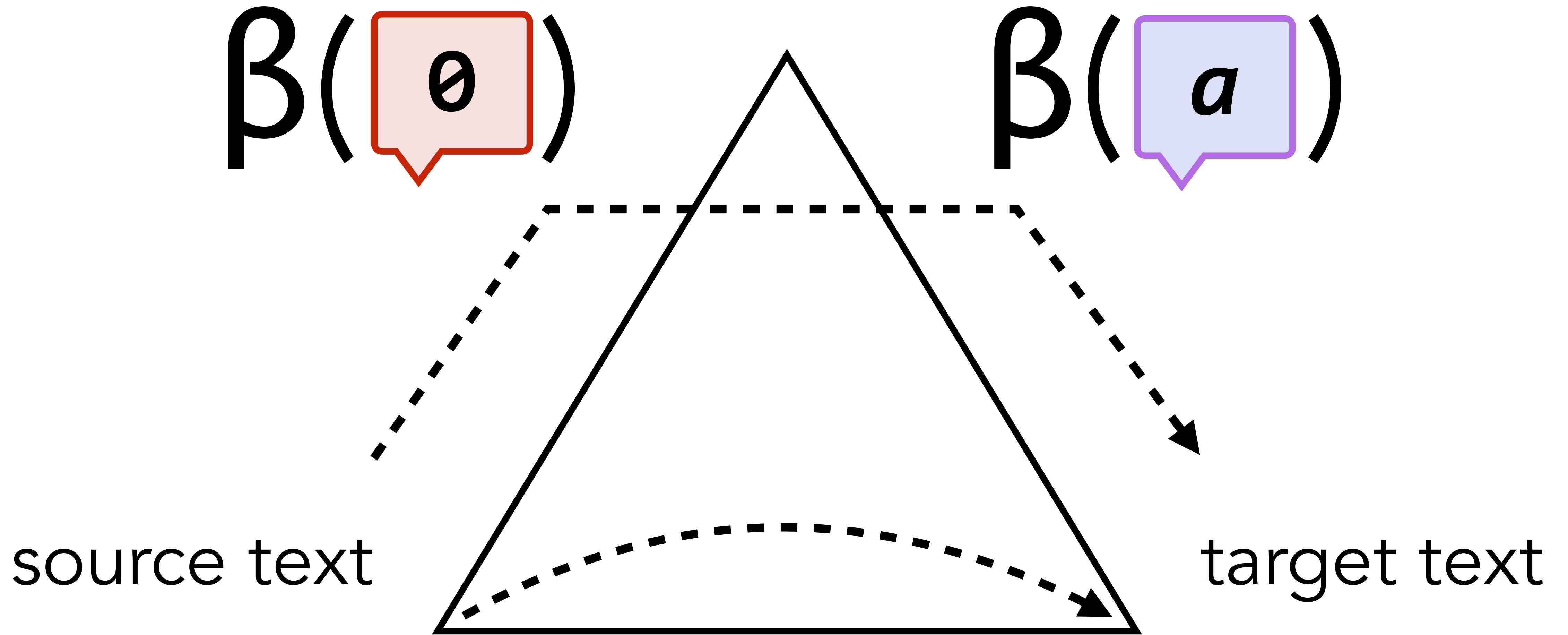


Translating with meaning





Interlingua!





Translation criterion

argmin

a

$$\text{KL}(\beta(\theta) \parallel \beta(a))$$



Translation criterion

argmin

a

$$\text{KL}(\beta(\theta) \parallel \beta(a))$$



Translation criterion

argmin

a

$$\text{KL}(\beta(\theta) \parallel \beta(a))$$



Translation criterion

argmin

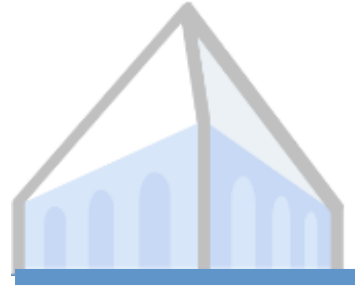
a

$$\text{KL}(\beta(\theta) \parallel \beta(a))$$



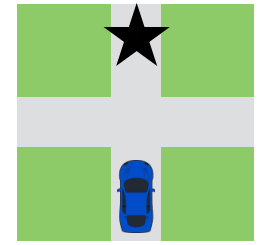


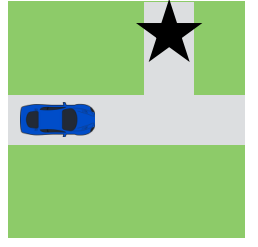
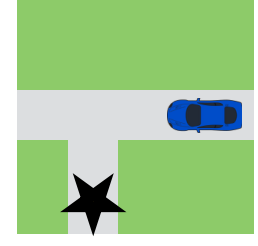
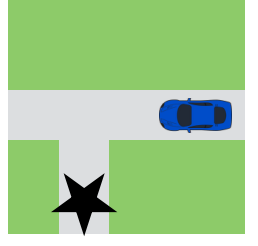
Computing representations

$$\operatorname{argmin}_a \text{KL}(\beta(\theta) \parallel \beta(a))$$



Computing representations: sparsity

$$\operatorname{argmin}_a \text{KL}(\beta(\theta) \parallel \beta(a))$$

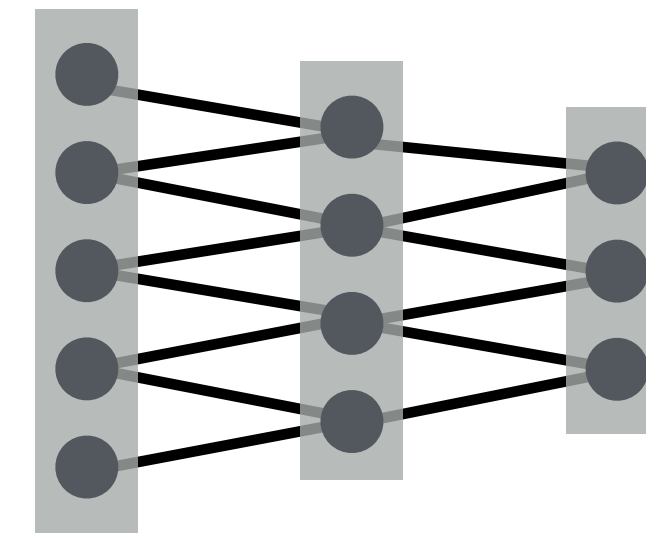
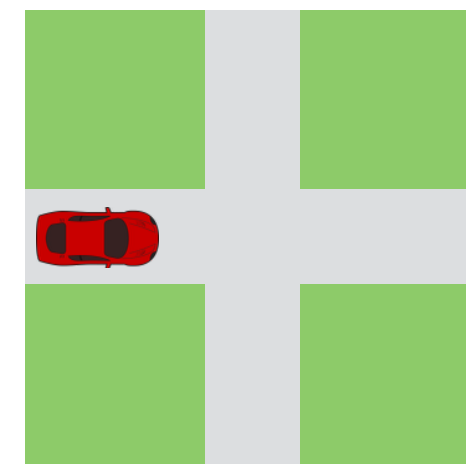
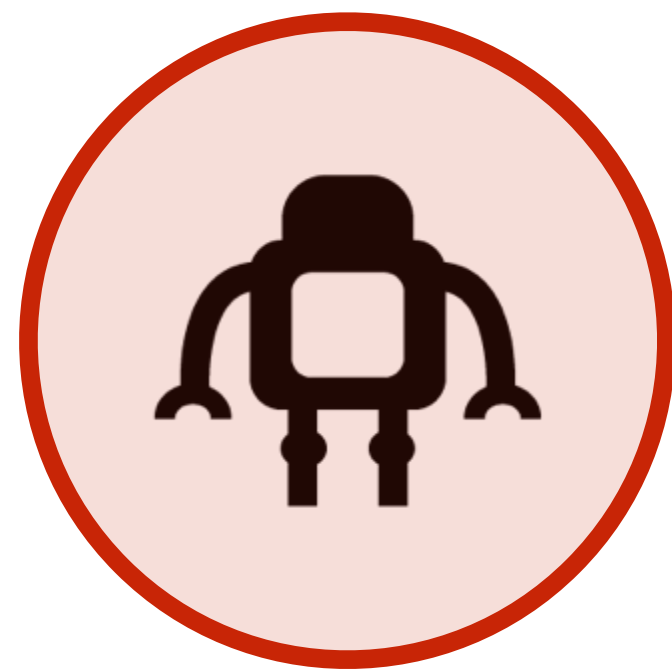
$p(\text{star+car} \mid \theta)$	$p(\text{star+car} \mid a)$
 	
	
	



Computing representations: smoothing

$$\operatorname{argmin}_a \text{KL}(\beta(\theta) \parallel \beta(a))$$

agent
policy



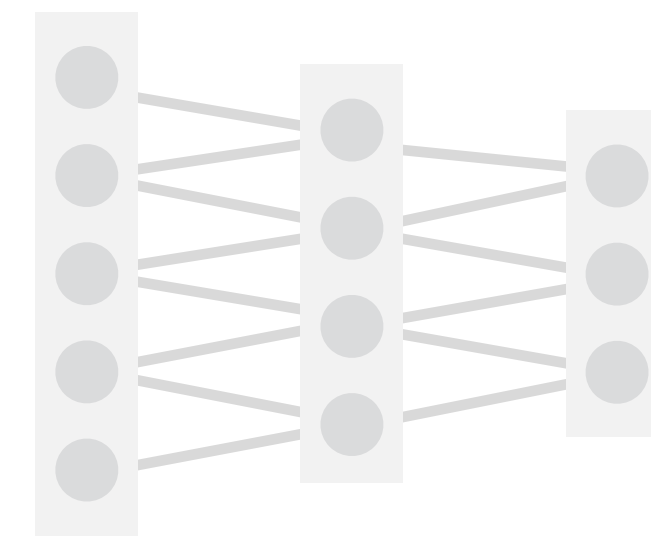
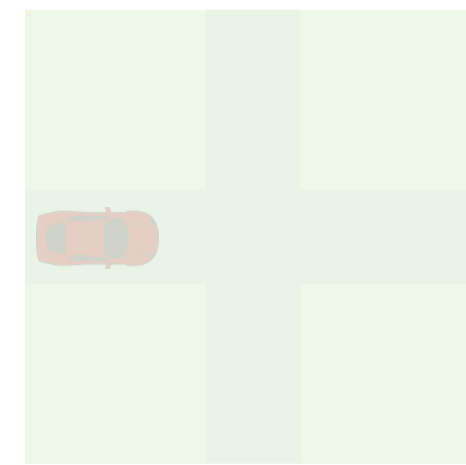
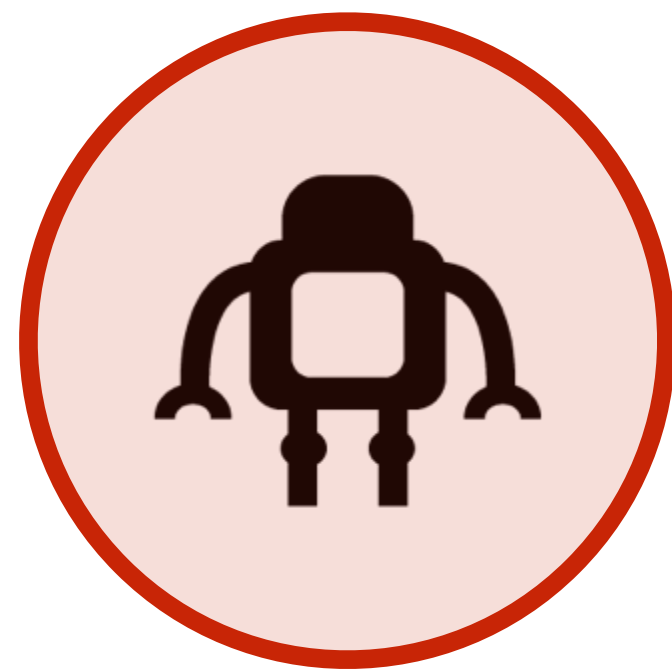
actions &
messages



Computing representations: smoothing

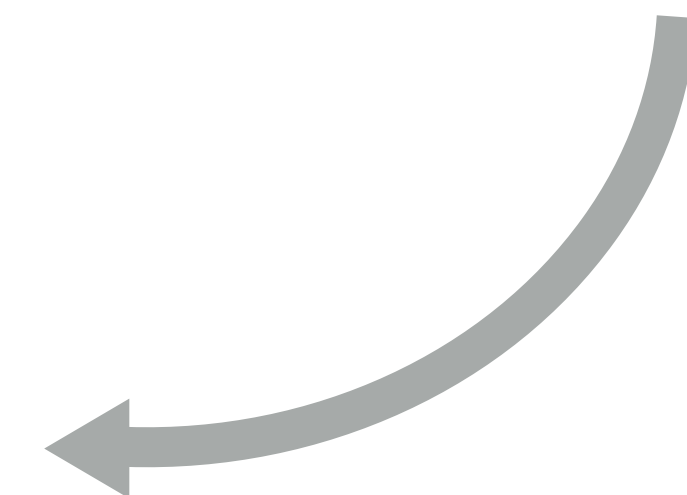
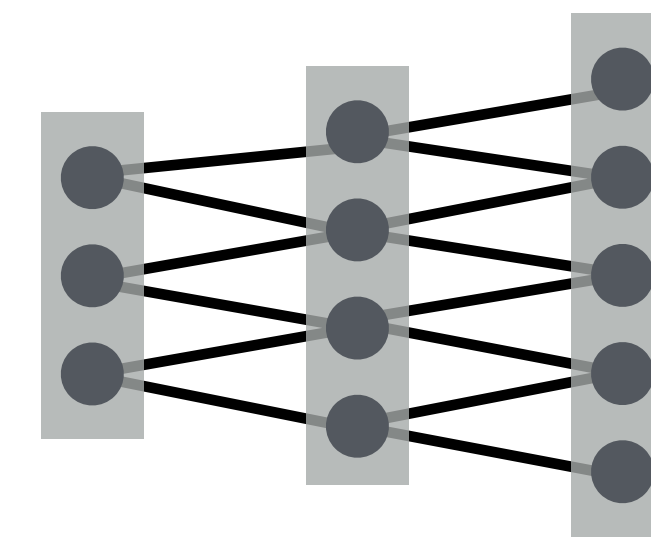
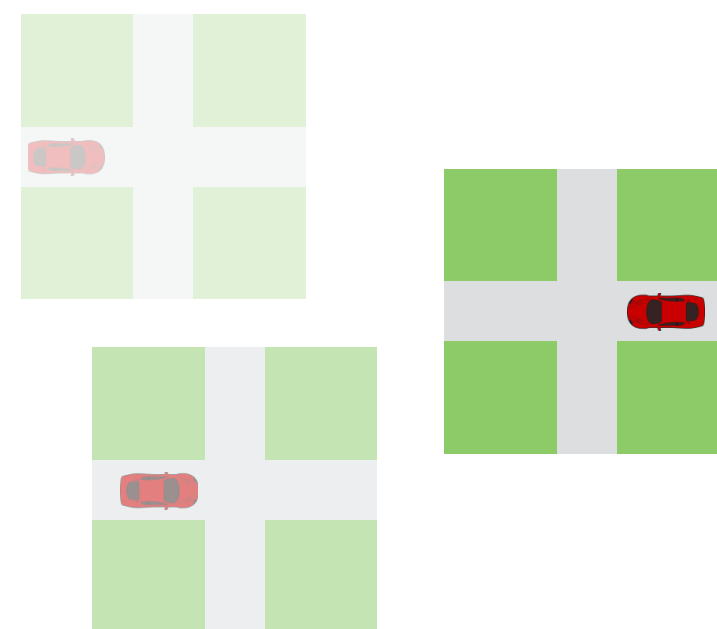
argmin _{a} $KL(\beta(\theta) \parallel \beta(a))$

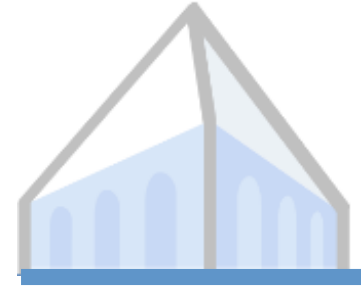
agent
policy



actions &
messages

agent
model

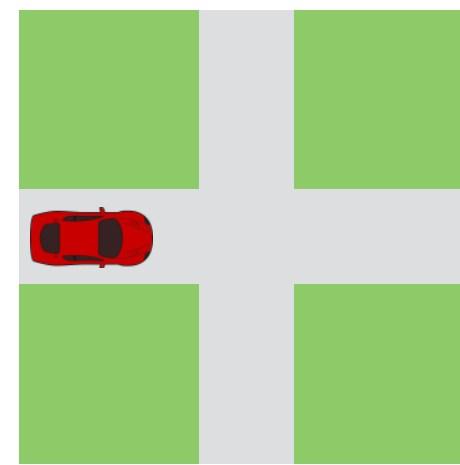




Computing representations: smoothing

argmin _{a} $KL(\beta(\theta) \parallel \beta(a))$

human



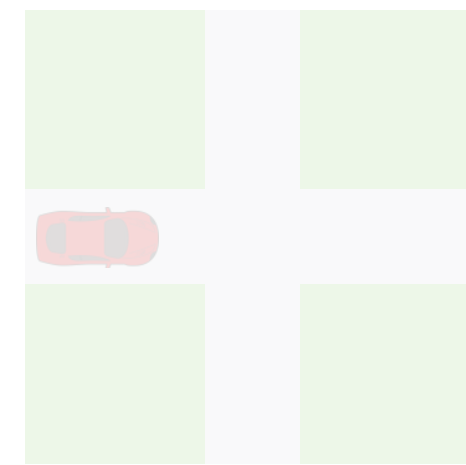
actions & messages



Computing representations: smoothing

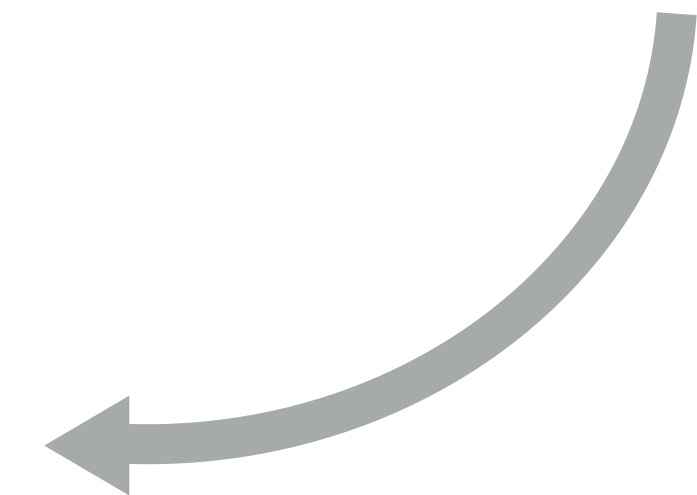
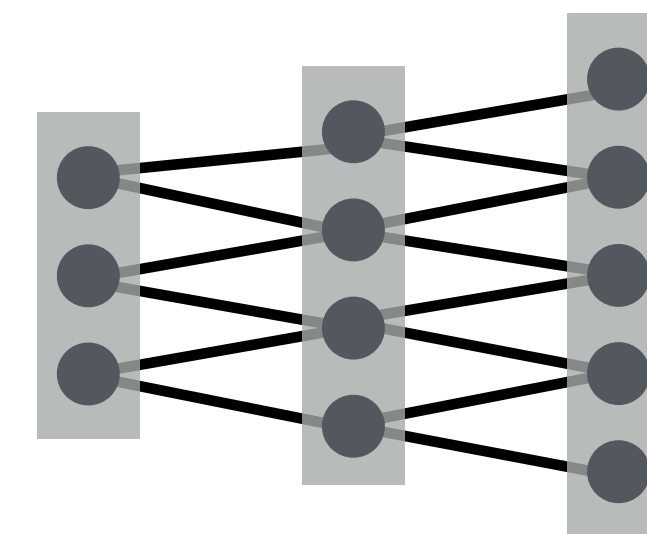
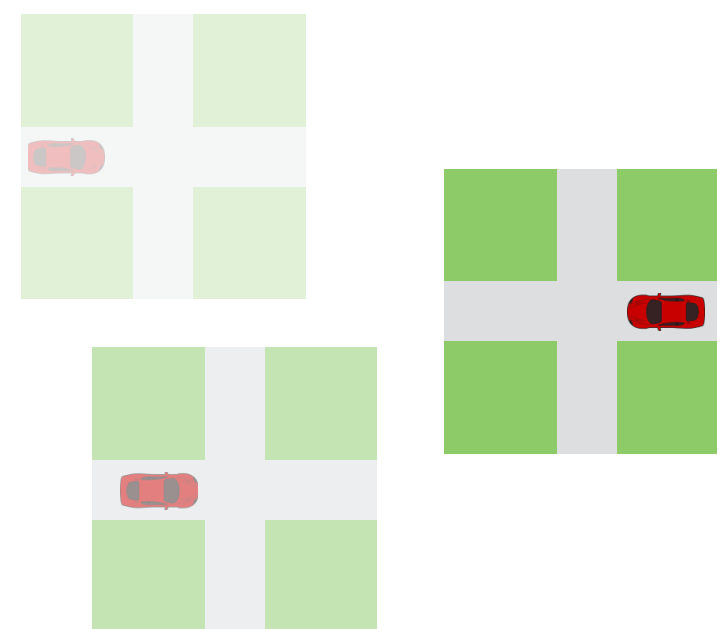
argmin _{a} $KL(\beta(\theta) \parallel \beta(a))$

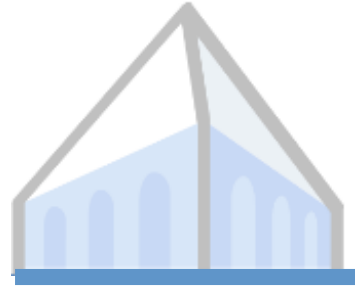
human
policy



actions &
messages

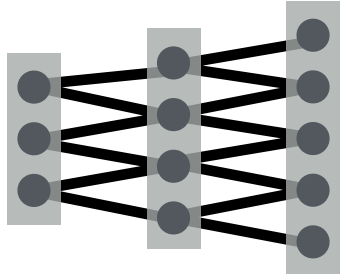
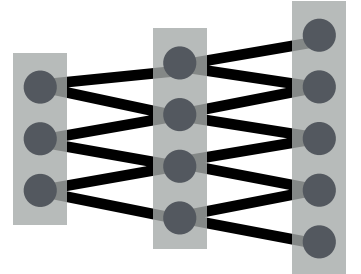
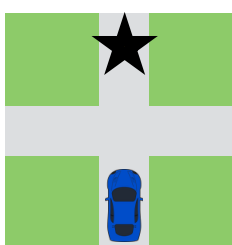
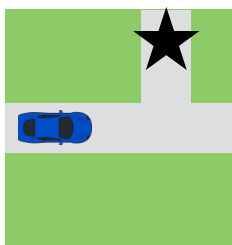
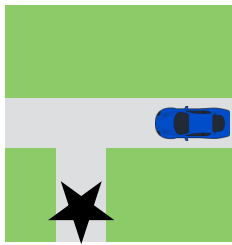
human
model





Computing representations: smoothing

$$\operatorname{argmin}_a \text{KL}(\beta(\theta) \parallel \beta(a))$$

	 θ	 a
	0.10	0.08
	0.05	0.01
	0.13	0.22



Computing KL

$$\operatorname{argmin}_a \text{KL}(\beta(\theta) \parallel \beta(a))$$



Computing KL

$$\operatorname{argmin}_a \operatorname{KL}(\beta(\theta) \parallel \beta(a))$$

$$\operatorname{KL}(p \parallel q) = \mathbf{E}_{\star}^p \frac{p(\star)}{q(\star)}$$



Computing KL: sampling

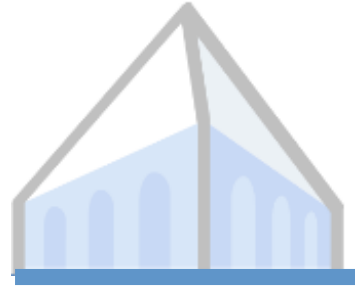
$$\operatorname{argmin}_a \operatorname{KL}(\beta(\theta) \parallel \beta(a))$$

$$\operatorname{KL}(p \parallel q) = \sum_i p(\text{star on car } i) \log \frac{p(\text{star on car } i)}{q(\text{star on car } i)}$$



Finding translations

$$\operatorname{argmin}_a \operatorname{KL}(\beta(\theta) \parallel \beta(a))$$



Finding translations: brute force

argmin _{a} $KL(\beta(\theta) \parallel \beta(a))$

going north → 0.5

crossing the intersection → 2.3

I'm done → 0.2

after you → 9.7



Finding translations: brute force

$$\operatorname{argmin}_a \text{KL}(\beta(\theta) \parallel \beta(a))$$

going north → 0.5

crossing the intersection → 2.3

I'm done → **0.2**

after you → 9.7

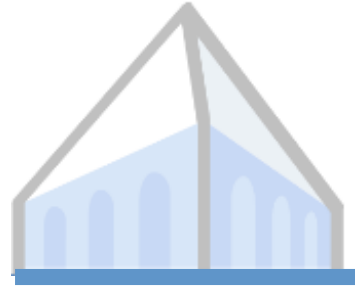


Finding translations

argmin

a

$$\text{KL}(\beta(\theta) \parallel \beta(a))$$



Outline

Natural language & neuralese

Statistical machine translation

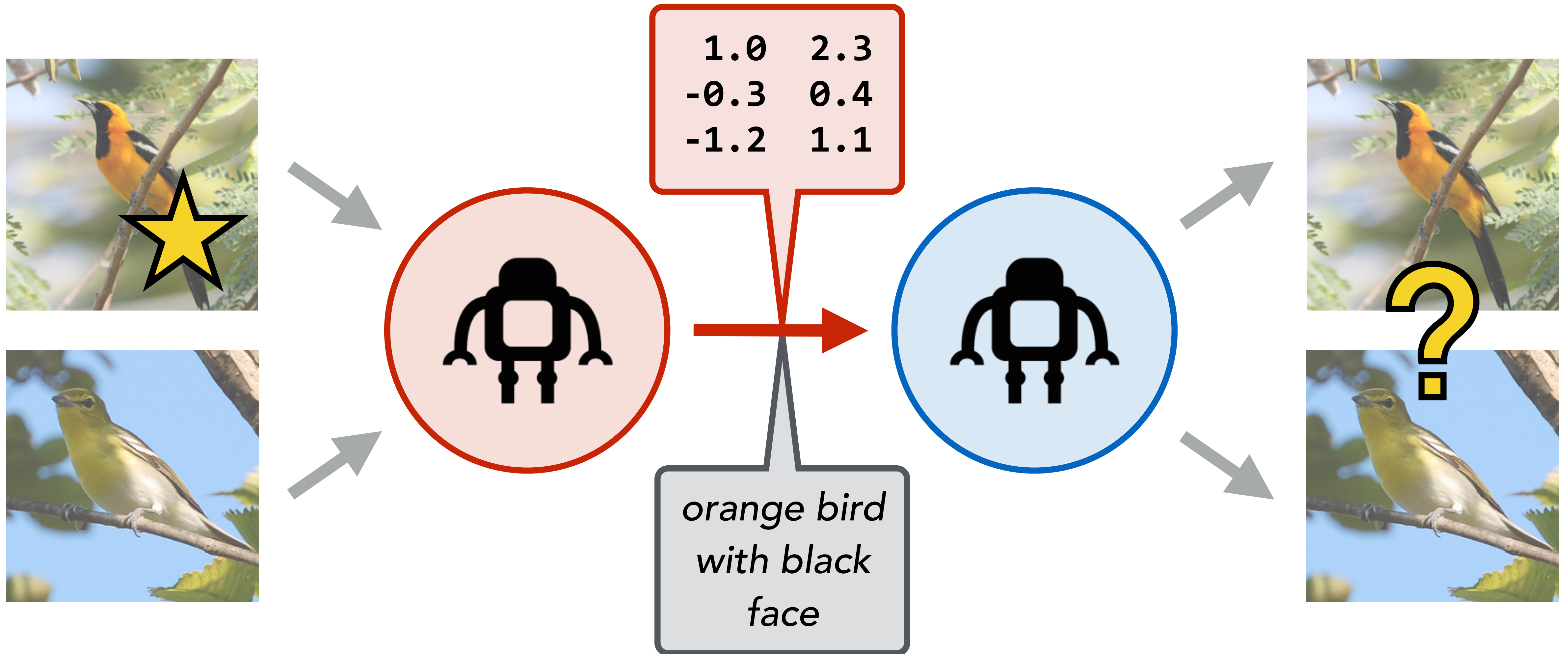
Semantic machine translation

Implementation details

Evaluation

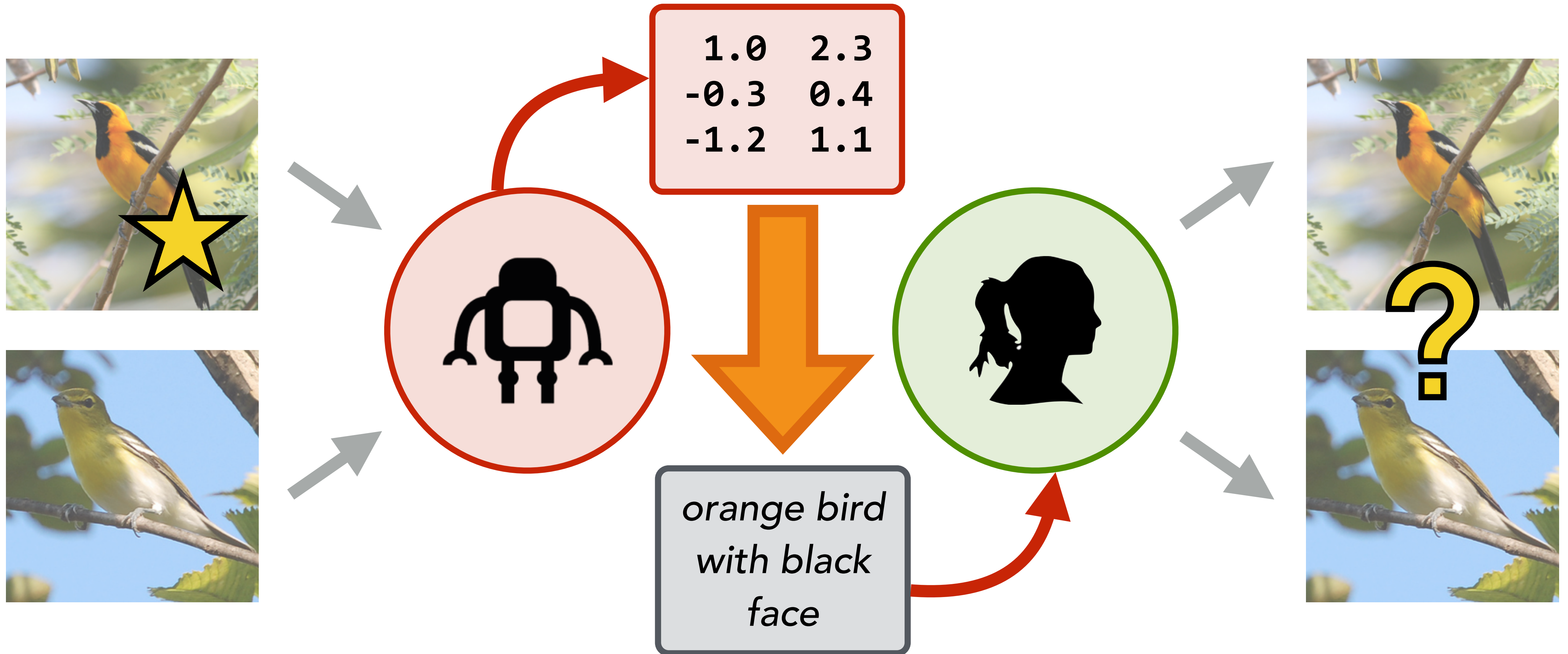


Referring expression games



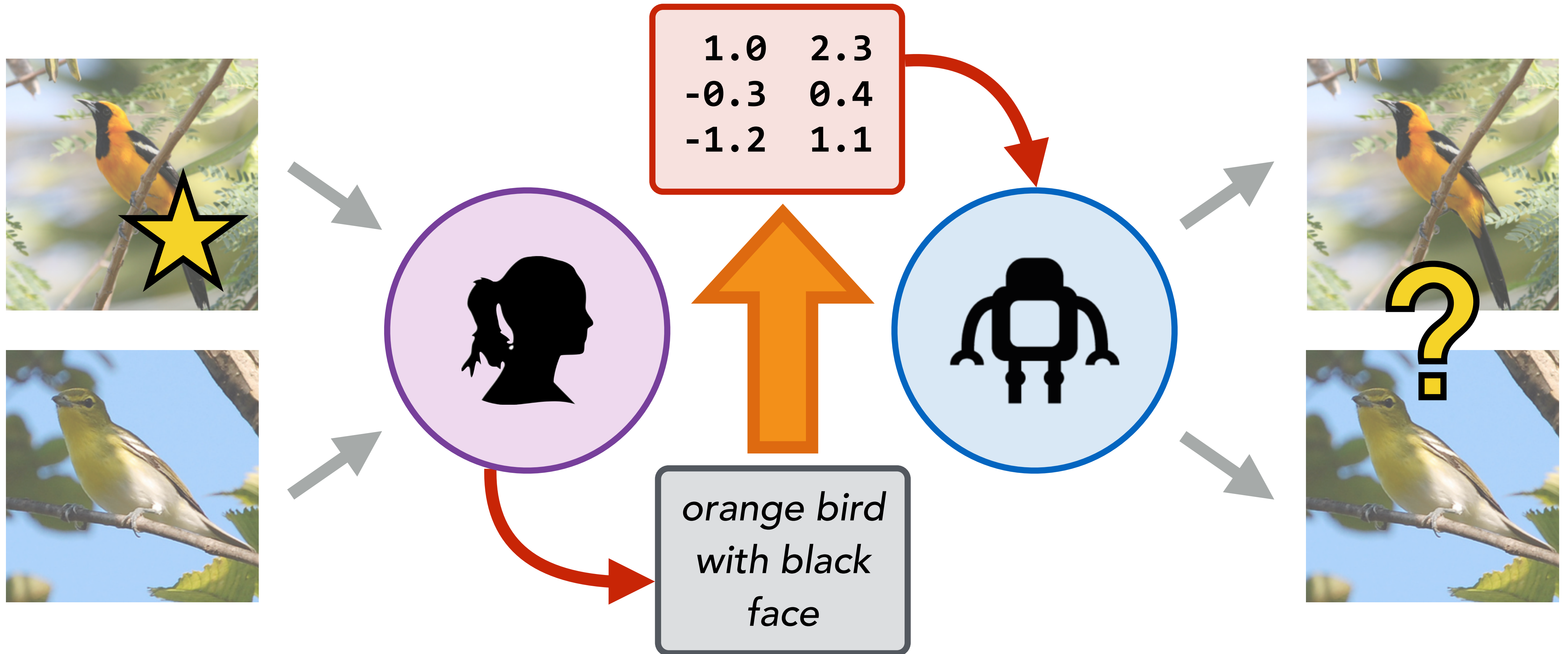


Evaluation: translator-in-the-loop



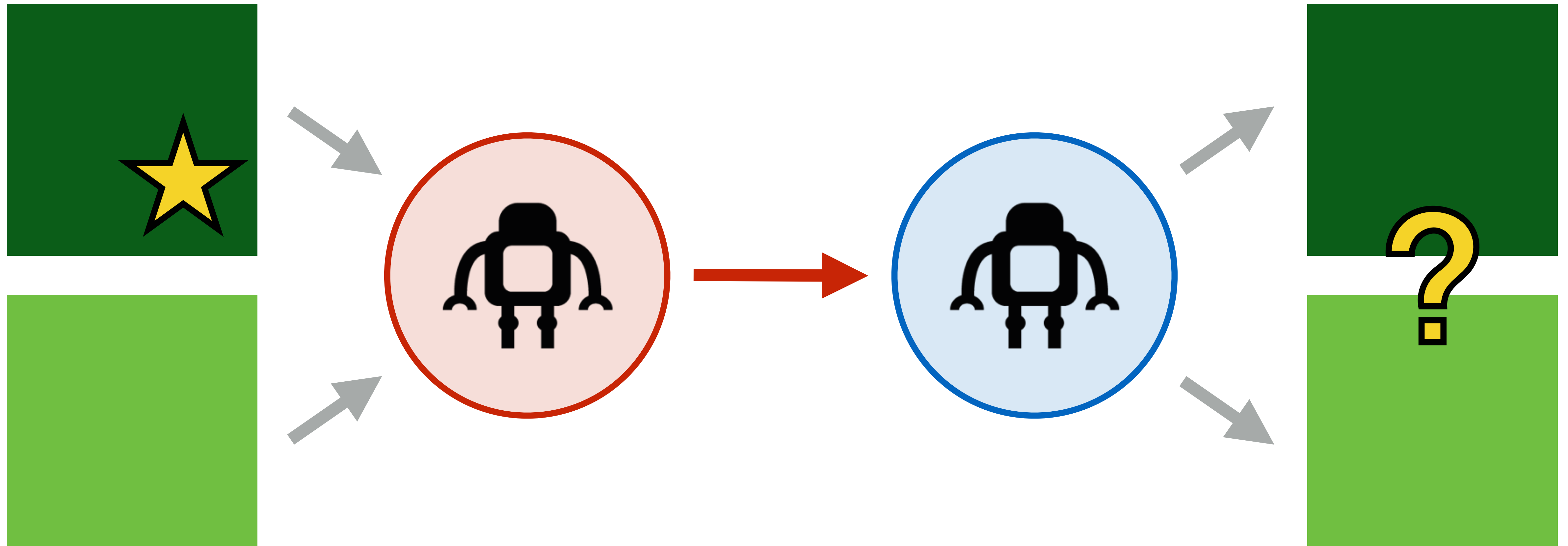


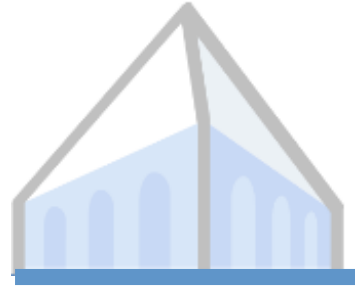
Evaluation: translator-in-the-loop







Experiment: color references





Experiment: color references

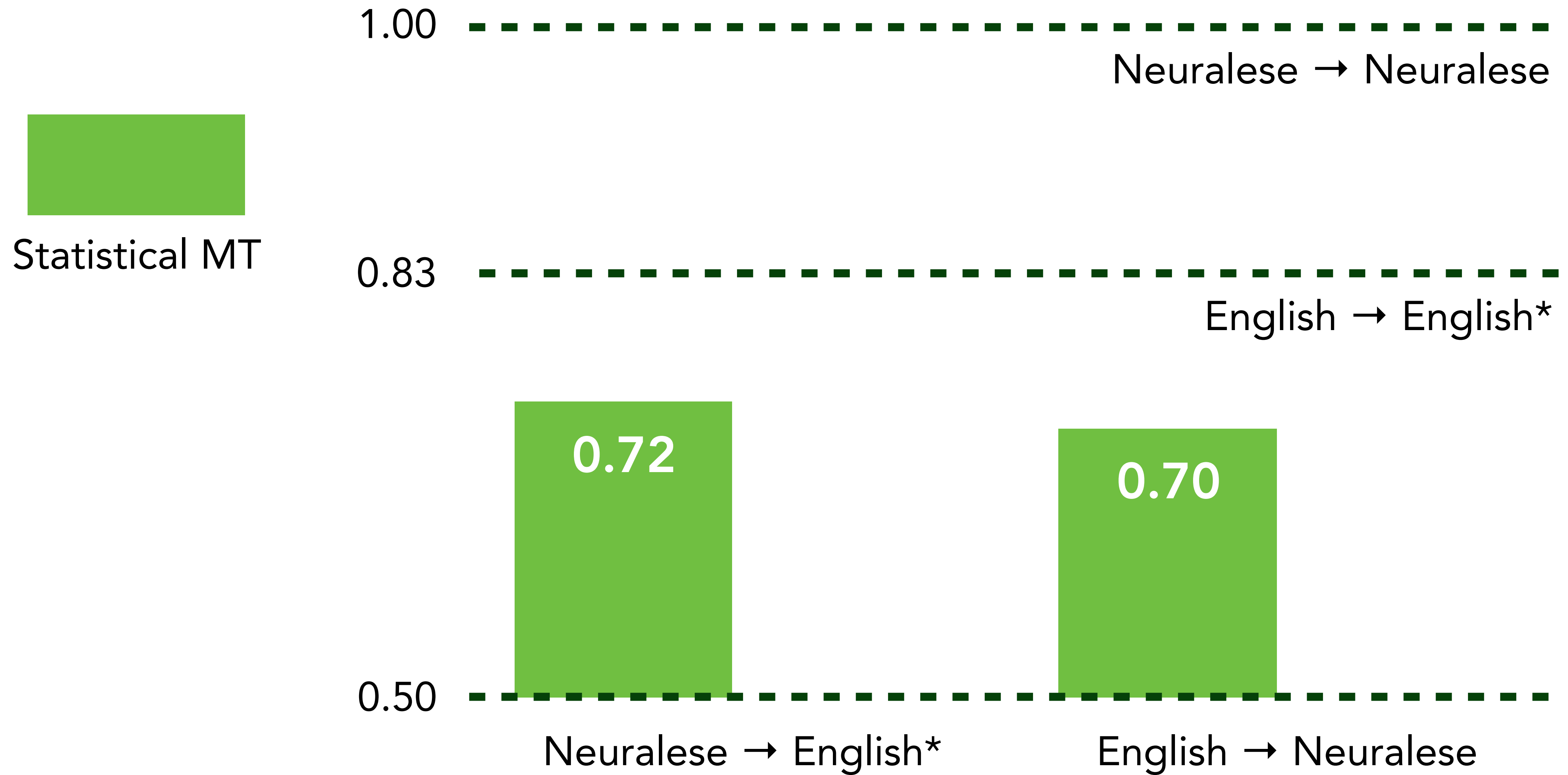
1.00  Neuralese → Neuralese

0.83  English → English*

0.50 

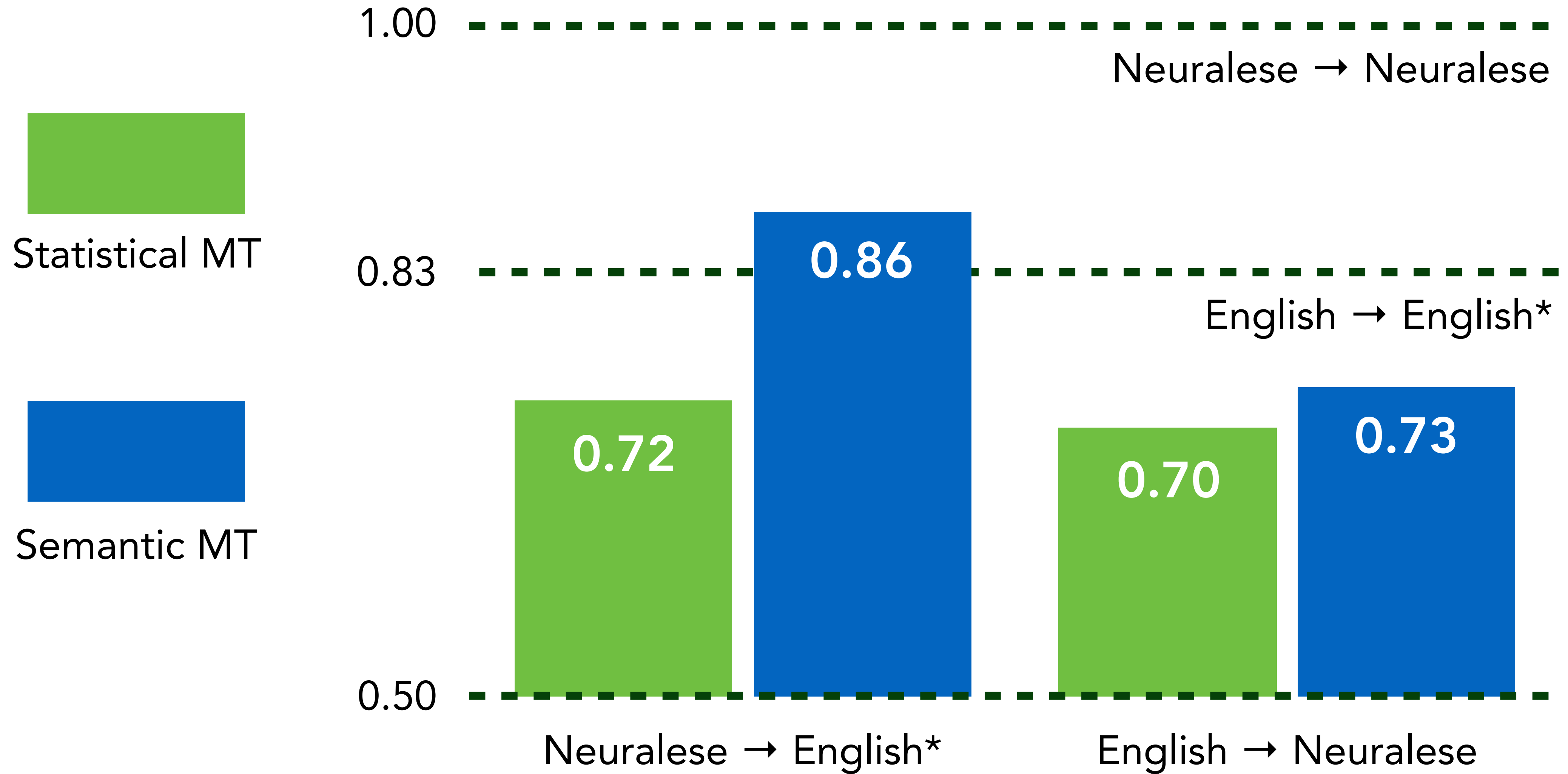


Experiment: color references



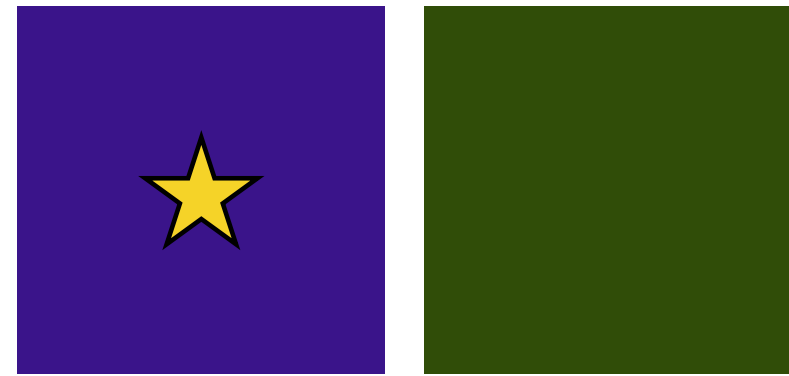


Experiment: color references





Experiment: color references



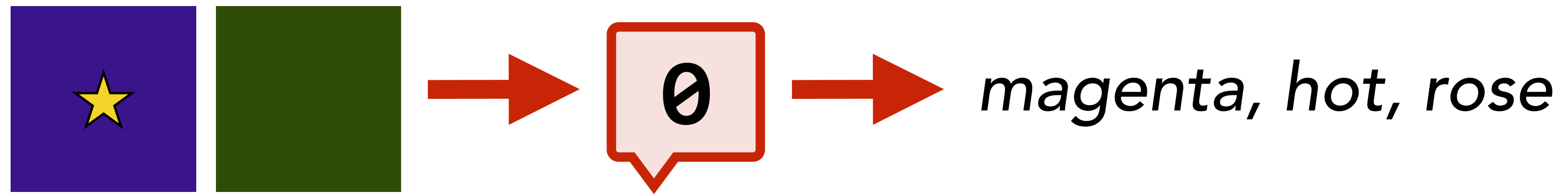


Experiment: color references



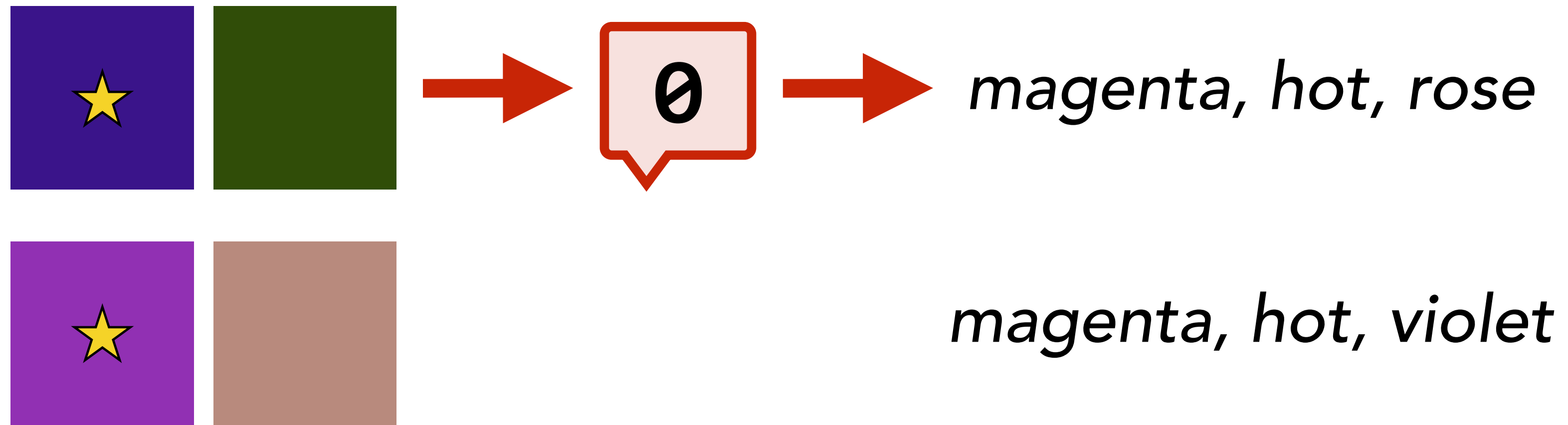


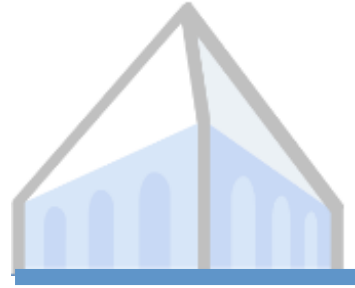
Experiment: color references



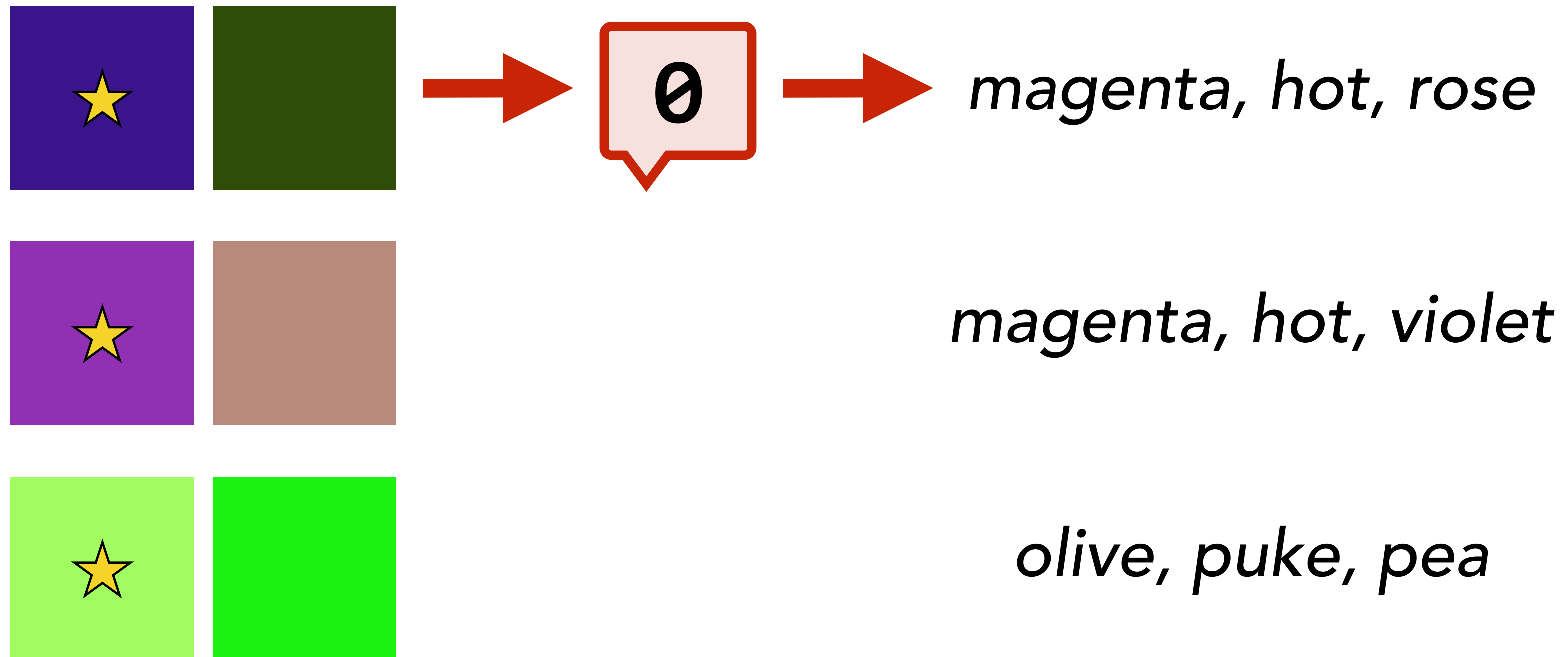


Experiment: color references



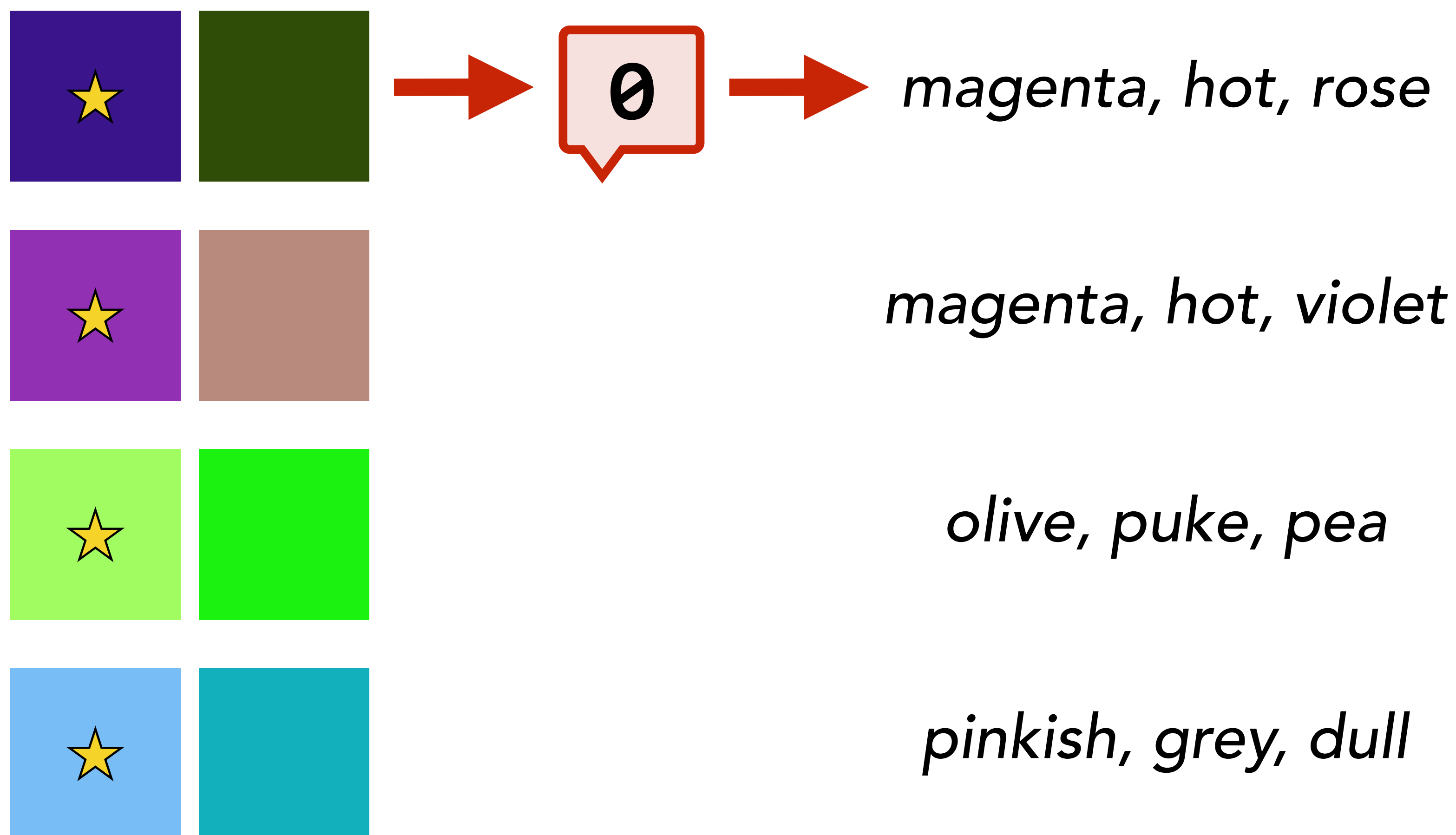


Experiment: color references



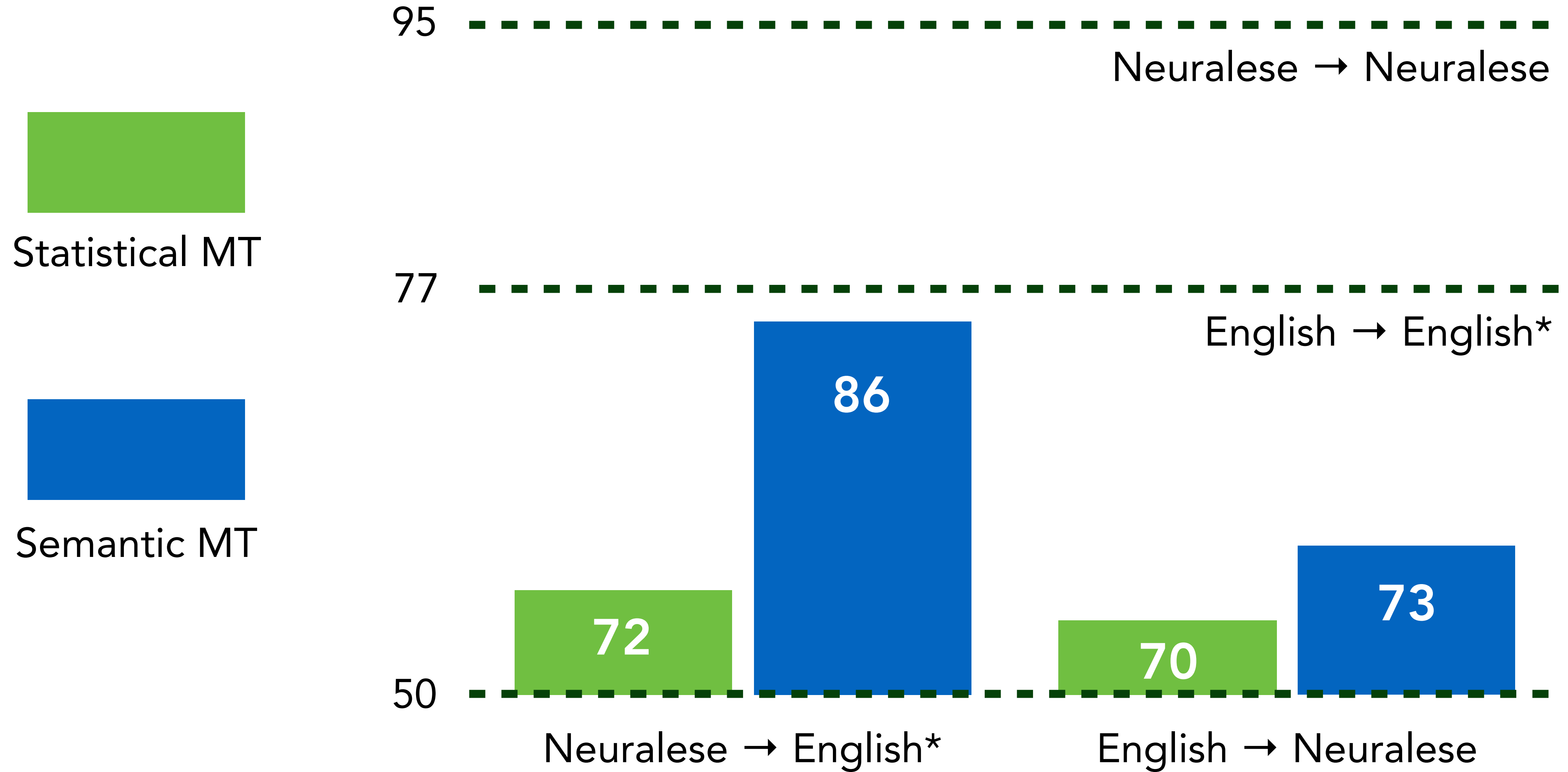


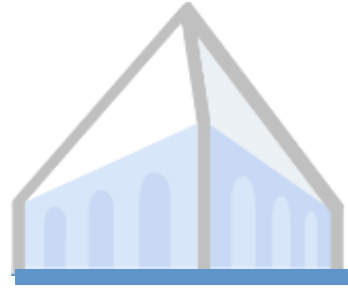
Experiment: color references





Experiment: image references





Experiment: image references



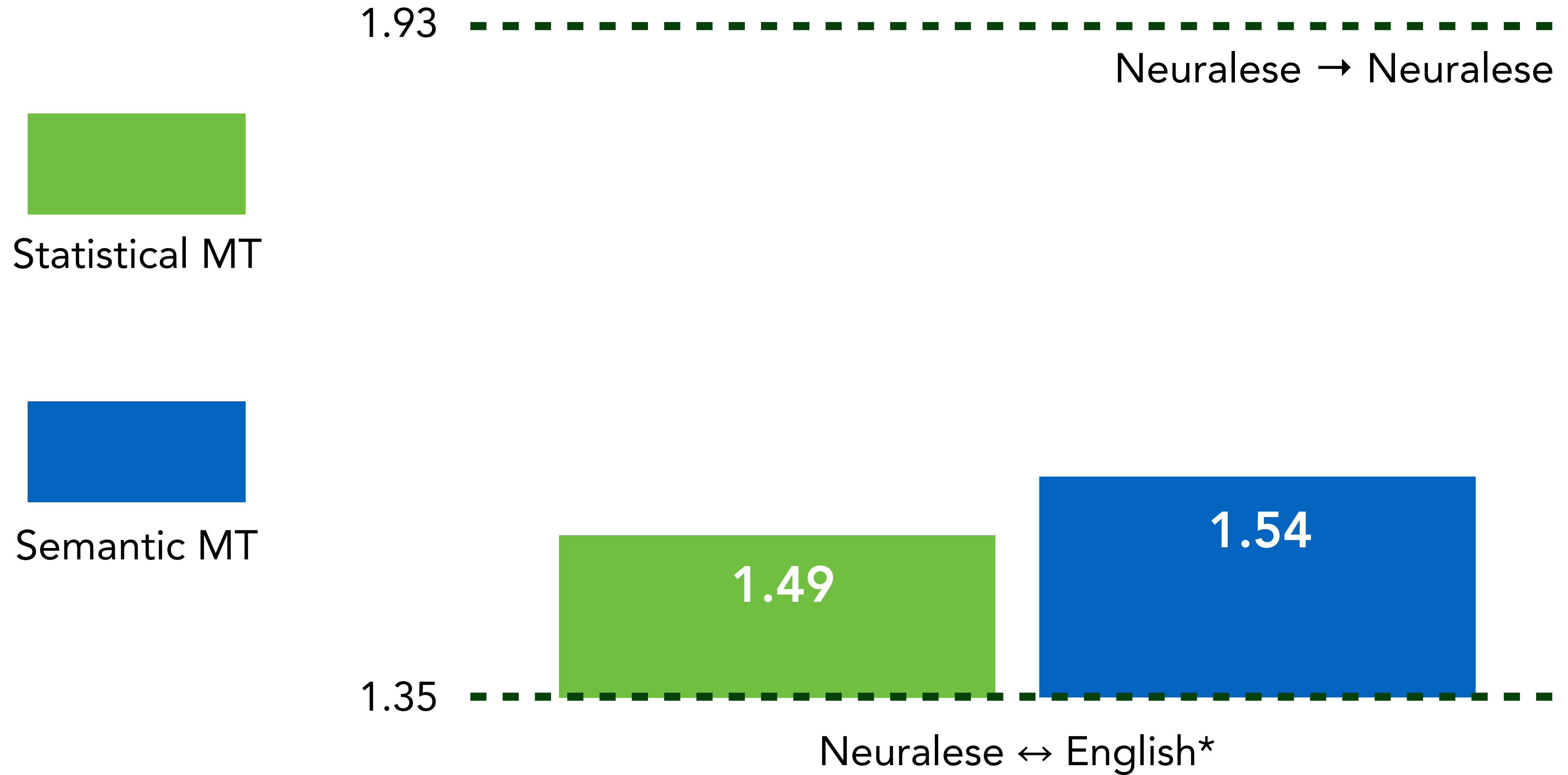
large bird, black wings, black crown



small brown, light brown, dark brown

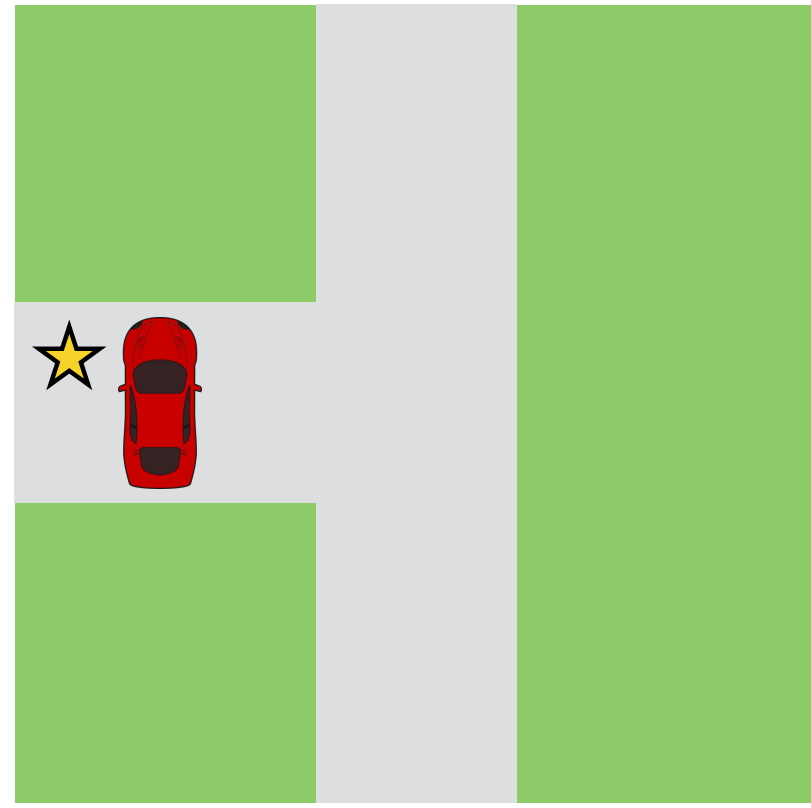


Experiment: driving game

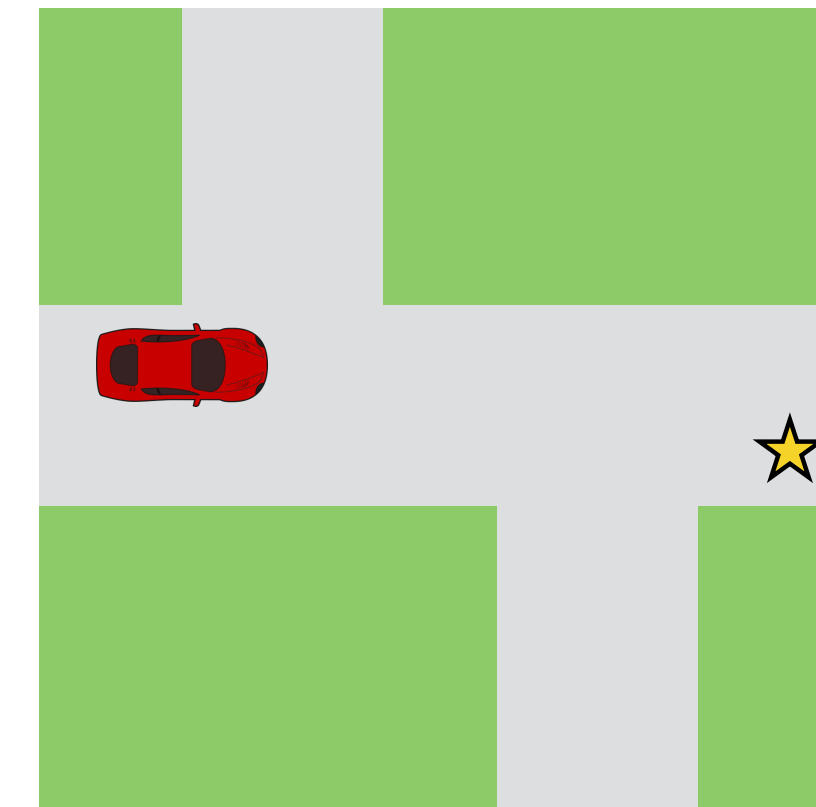




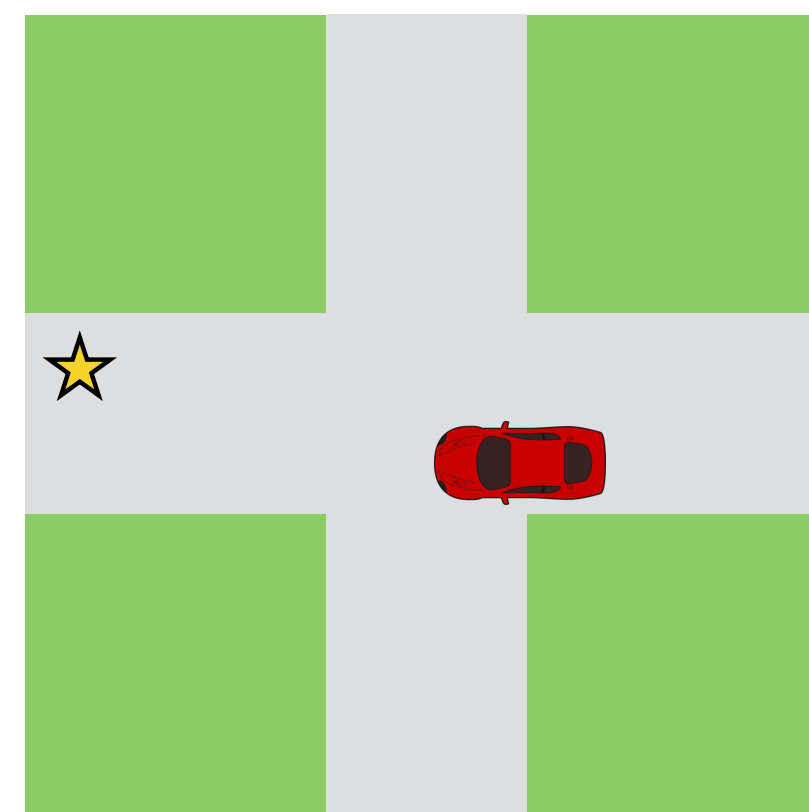
How to translate



*at goal
done
left to top*



*you first
following
going down*



*going in intersection
proceed
going*



Conclusions so far

- Classical notions of “meaning” apply even to un-language-like things (e.g. RNN states)
- These meanings can be compactly represented without logical forms if we have access to world states
- Communicating policies “say” interpretable things!



Conclusions so far

- Classical notions of “meaning” apply even to non-language-like things (e.g. RNN states)
- These meanings can be compactly represented without logical forms if we have access to world states
- Communicating policies “say” interpretable things!



Conclusions so far

- Classical notions of “meaning” apply even to non-language-like things (e.g. RNN states)
- These meanings can be compactly represented without logical forms if we have access to world states
- Communicating policies “say” interpretable things!



Limitations

$$\operatorname{argmin}_a \operatorname{KL}(\beta(\theta) \parallel \beta(a))$$

$$\operatorname{KL}(p \parallel q) = \sum_{i} p(\text{star on green square } i) \log \frac{p(\text{star on green square } i)}{q(\text{star on green square } i)}$$

but what about compositionality?

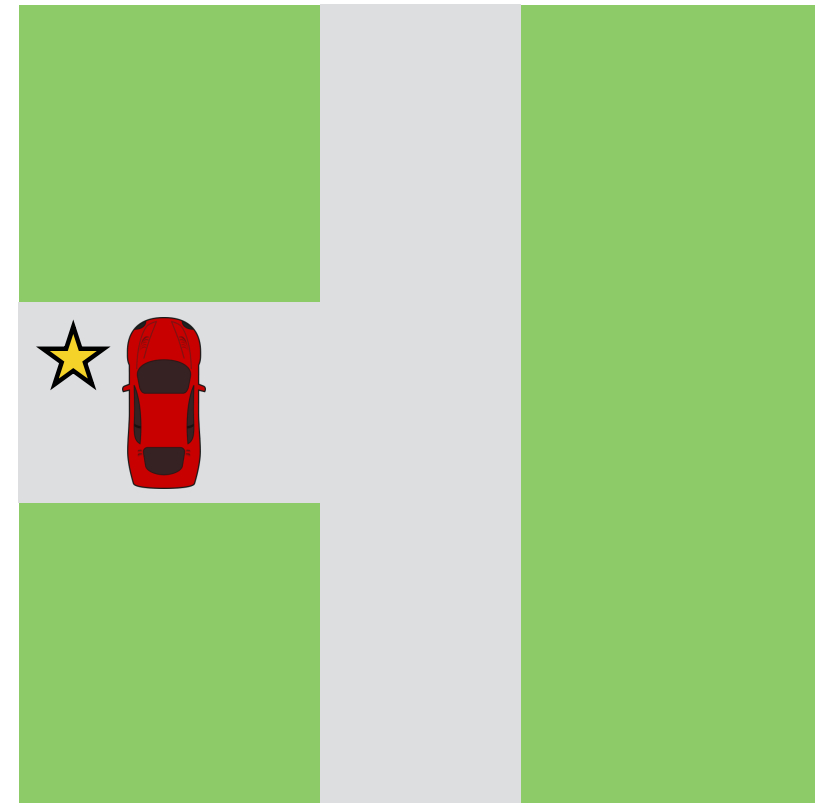
Analogs of linguistic structure in deep representations



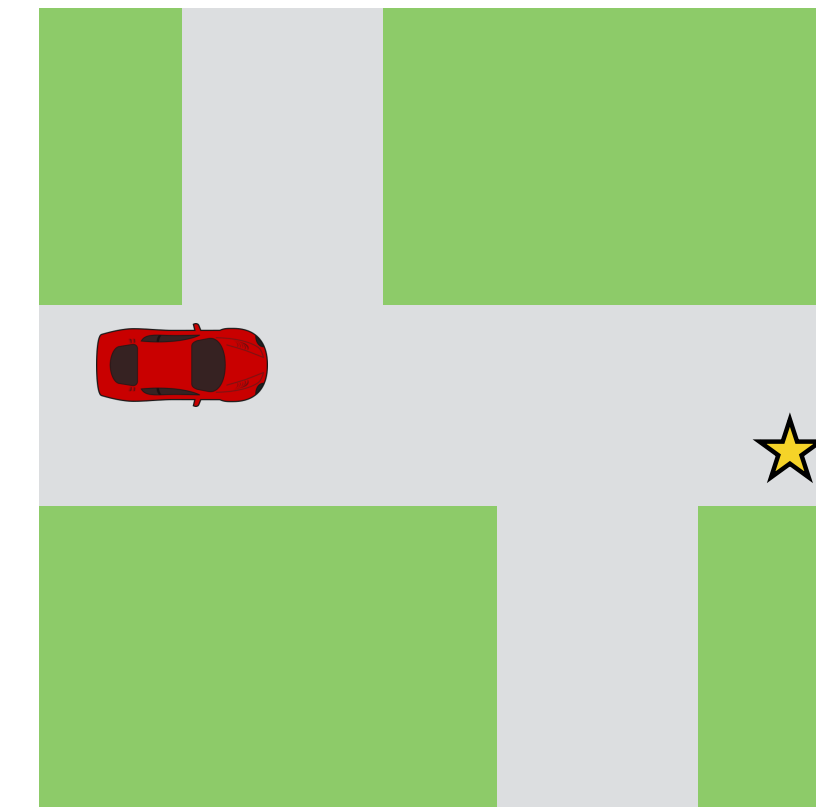
Jacob Andreas and Dan Klein



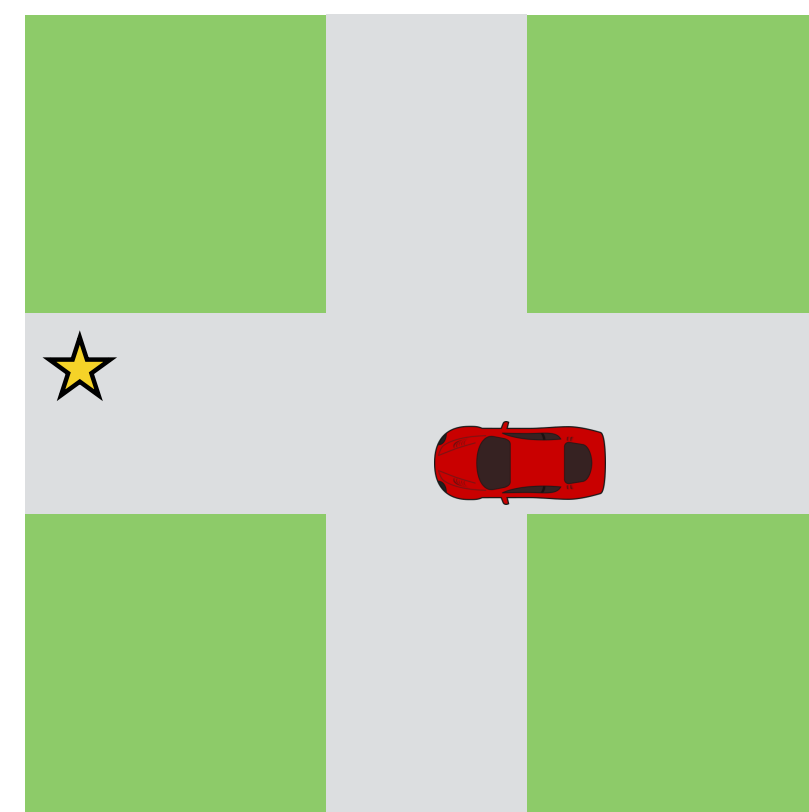
"Flat" semantics



*at goal
done*



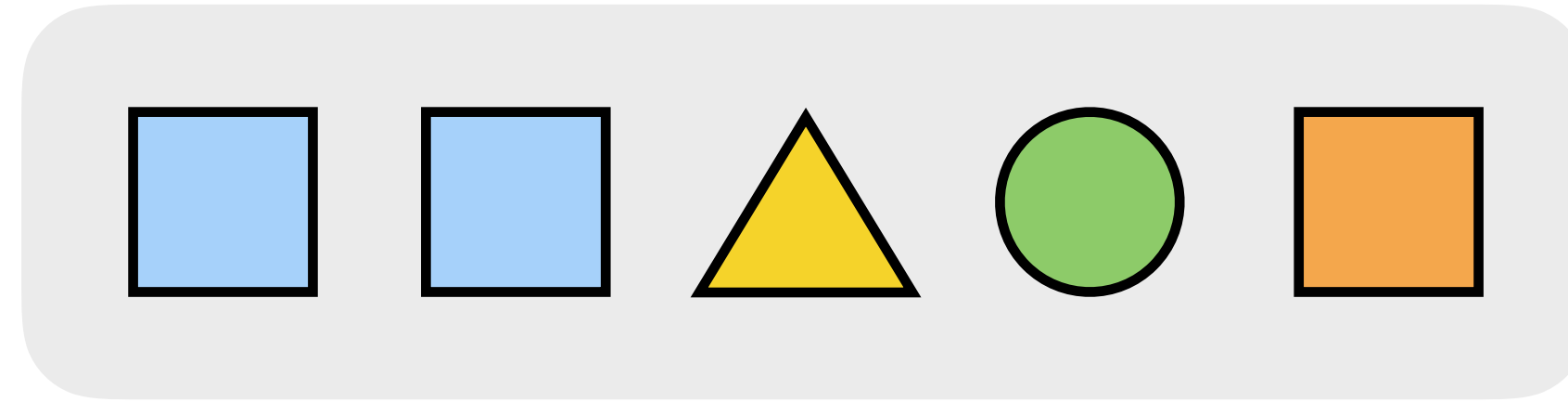
*you first
following*



*going in intersection
proceed
going*

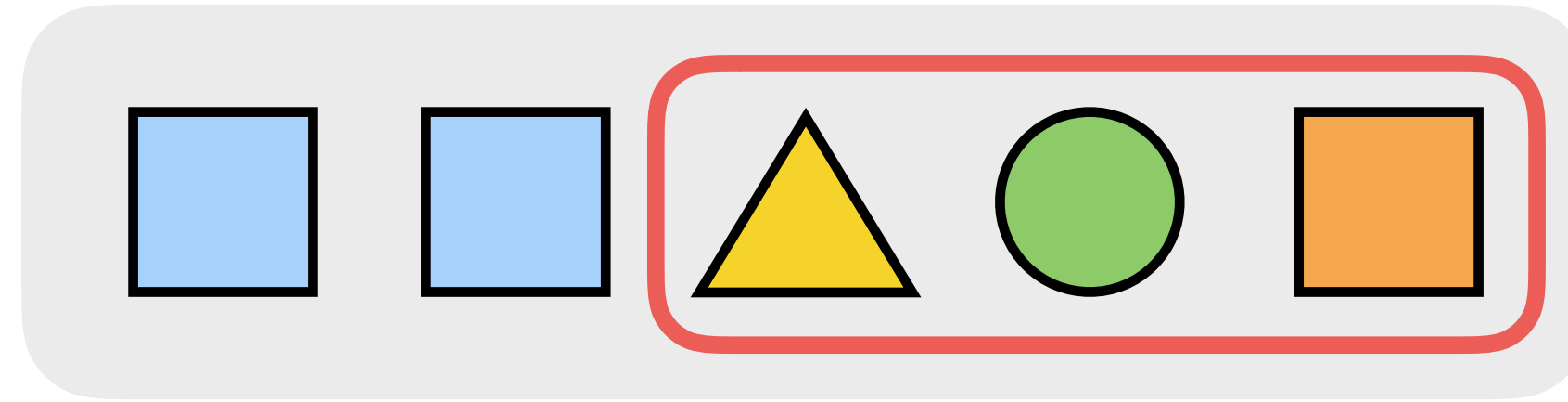


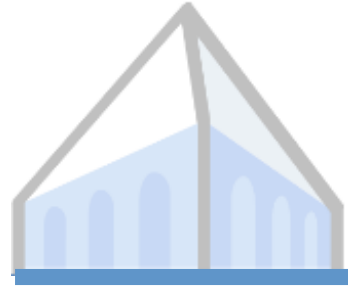
Compositional semantics



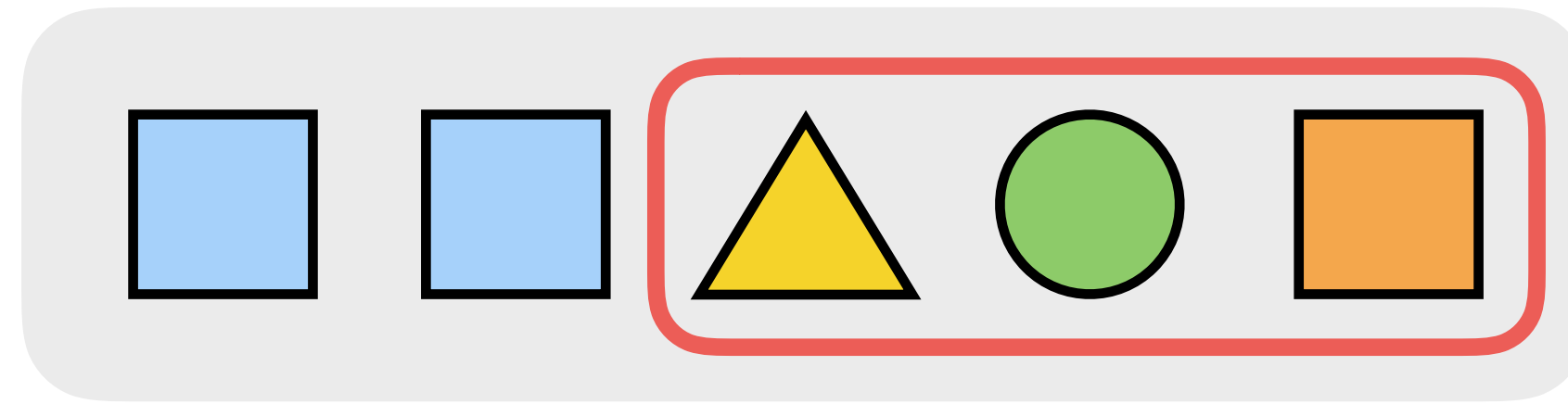


Compositional semantics

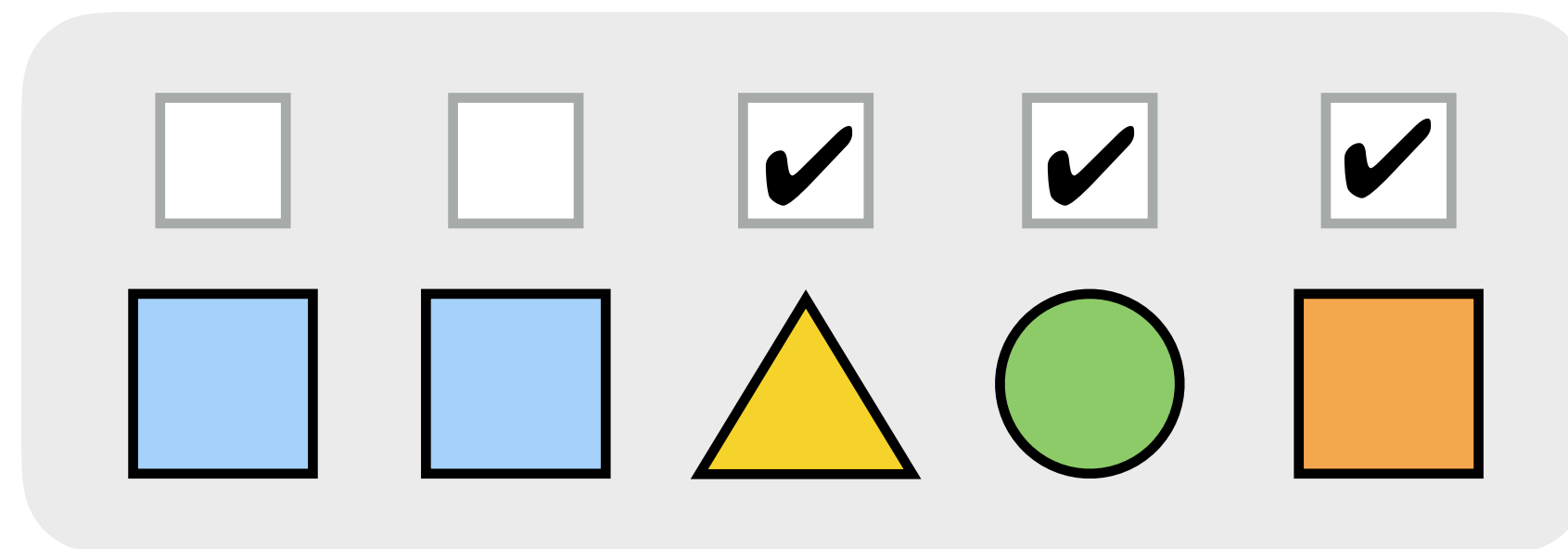


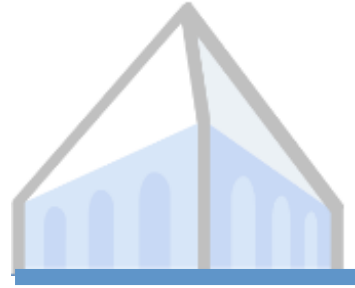


Compositional semantics

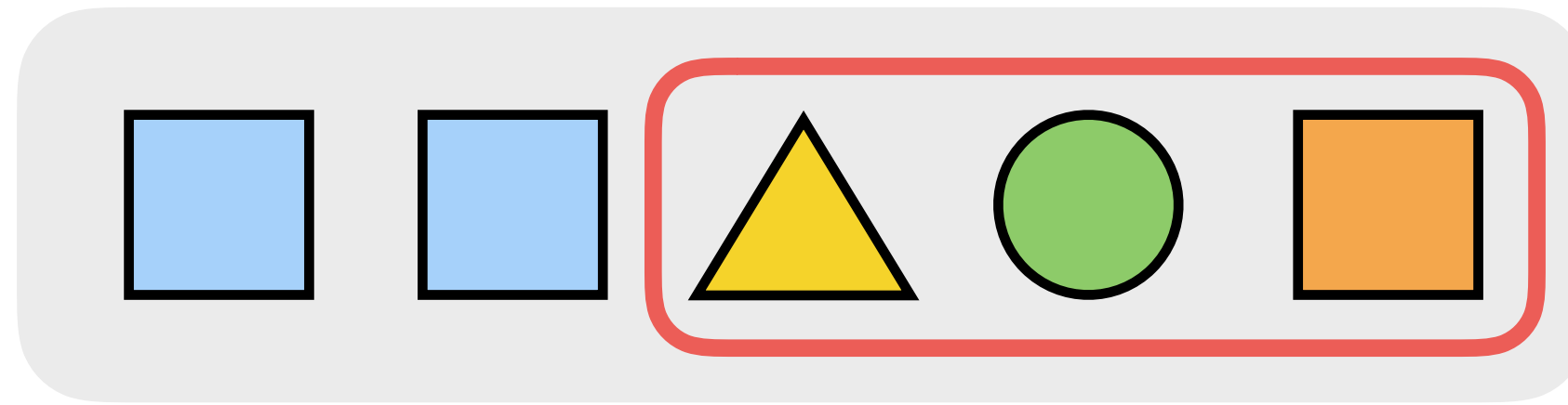


message

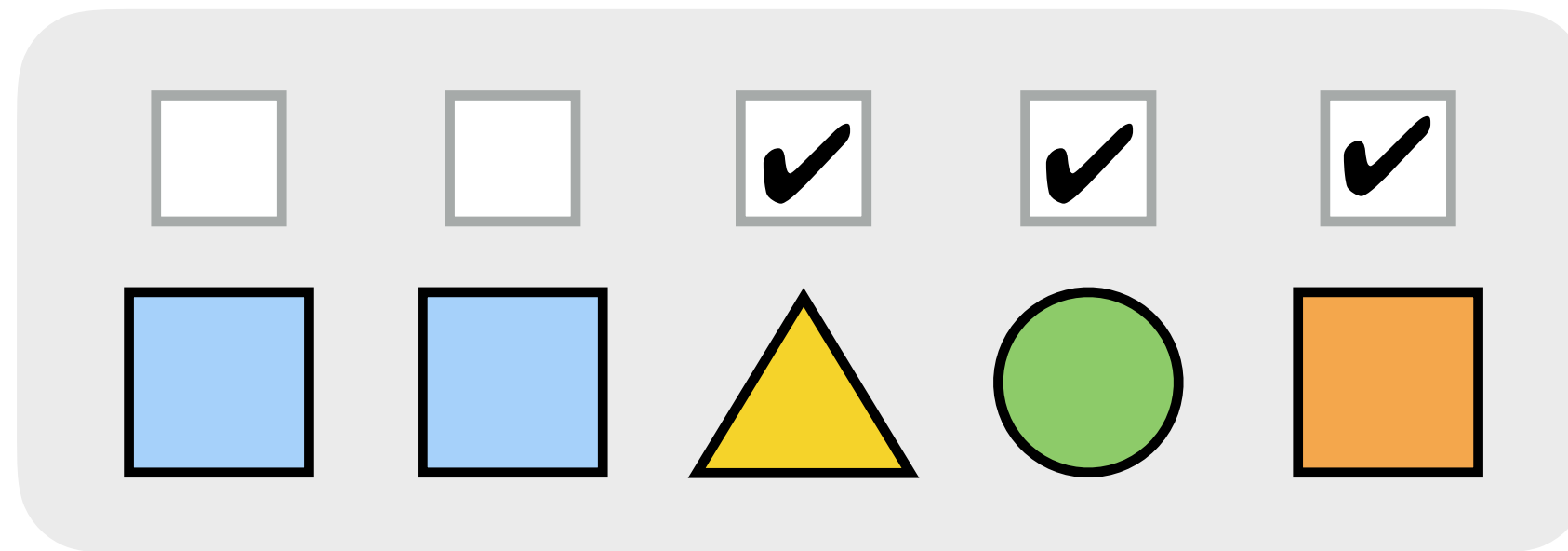




Compositional semantics

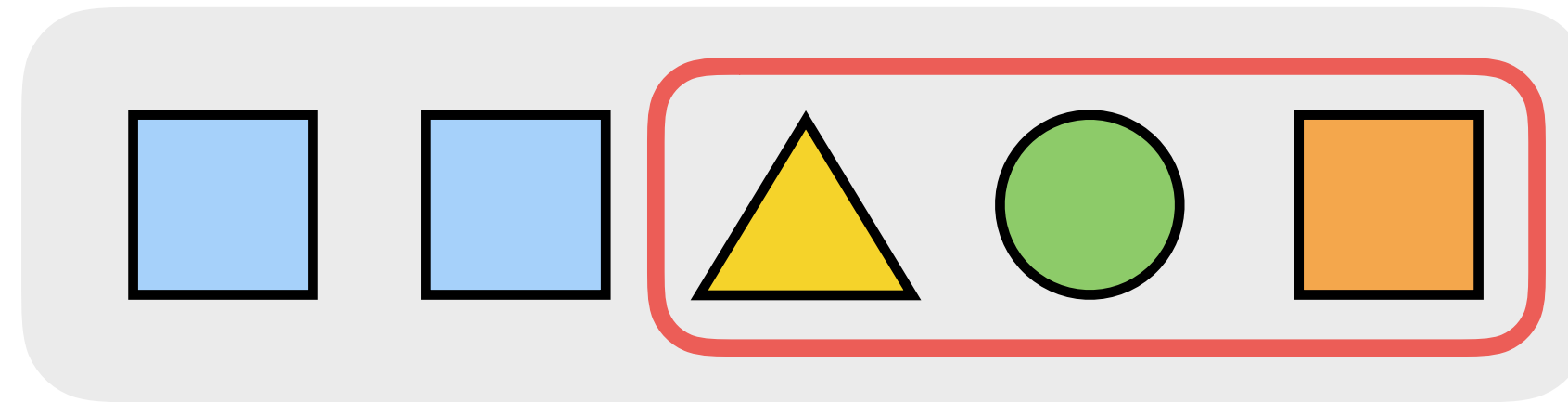


*everything but the blue shapes
orange square and non-squares*



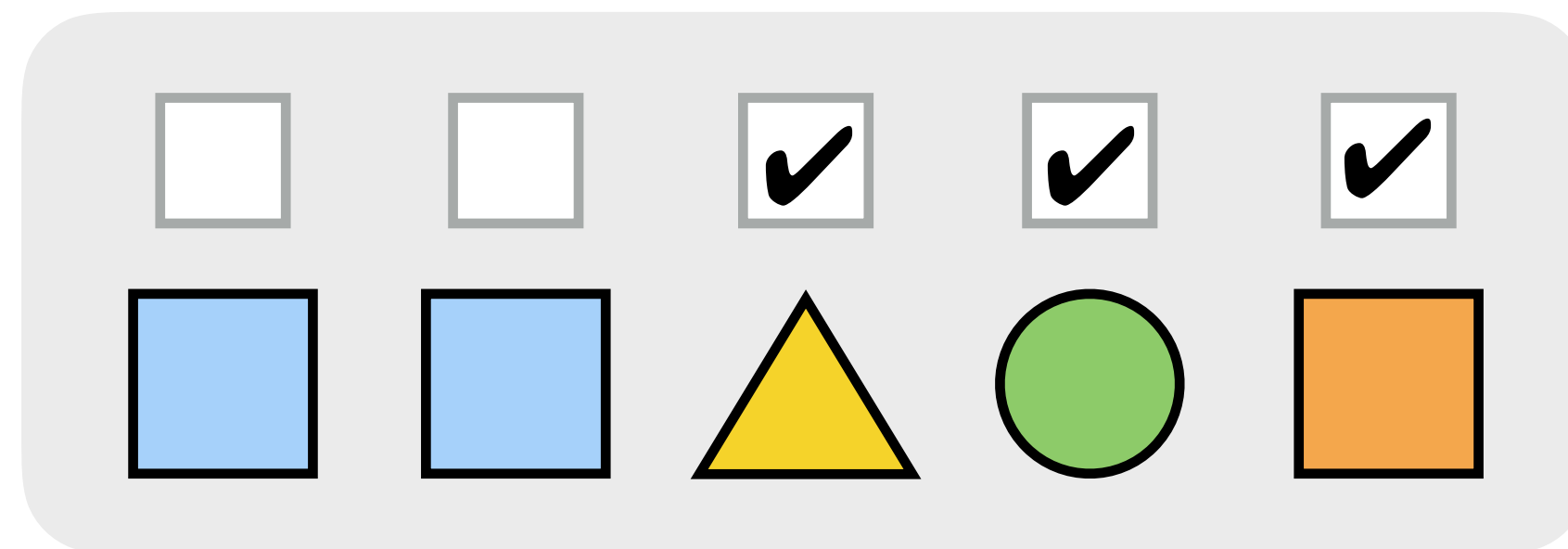


Compositional semantics



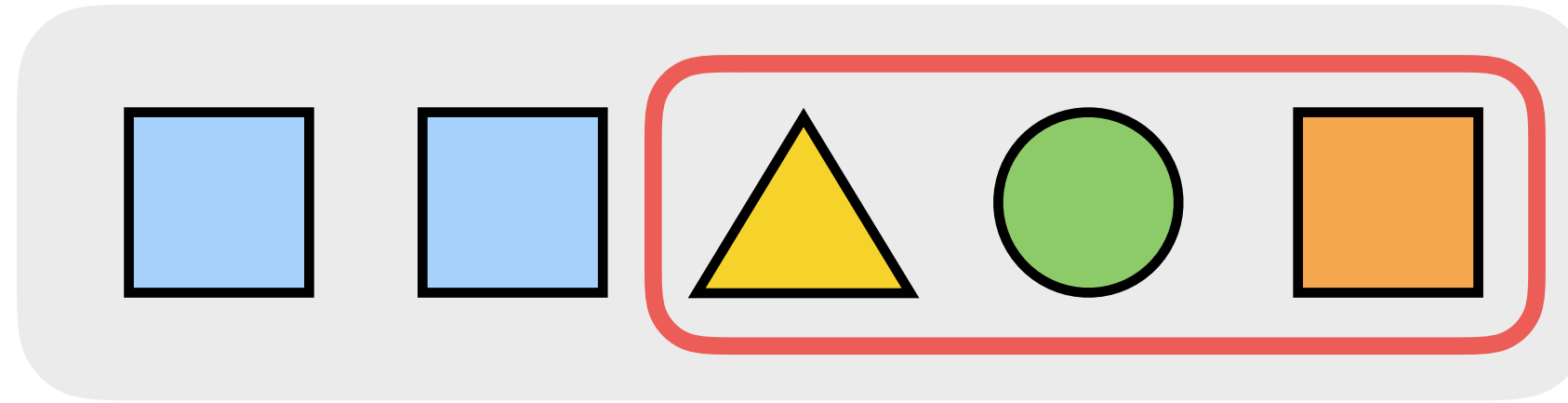
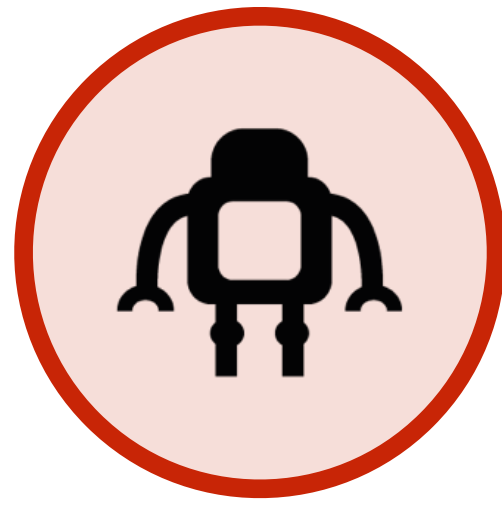
`lambda x: not(blue(x))`

`lambda x: or(orange(x), not(square(x)))`

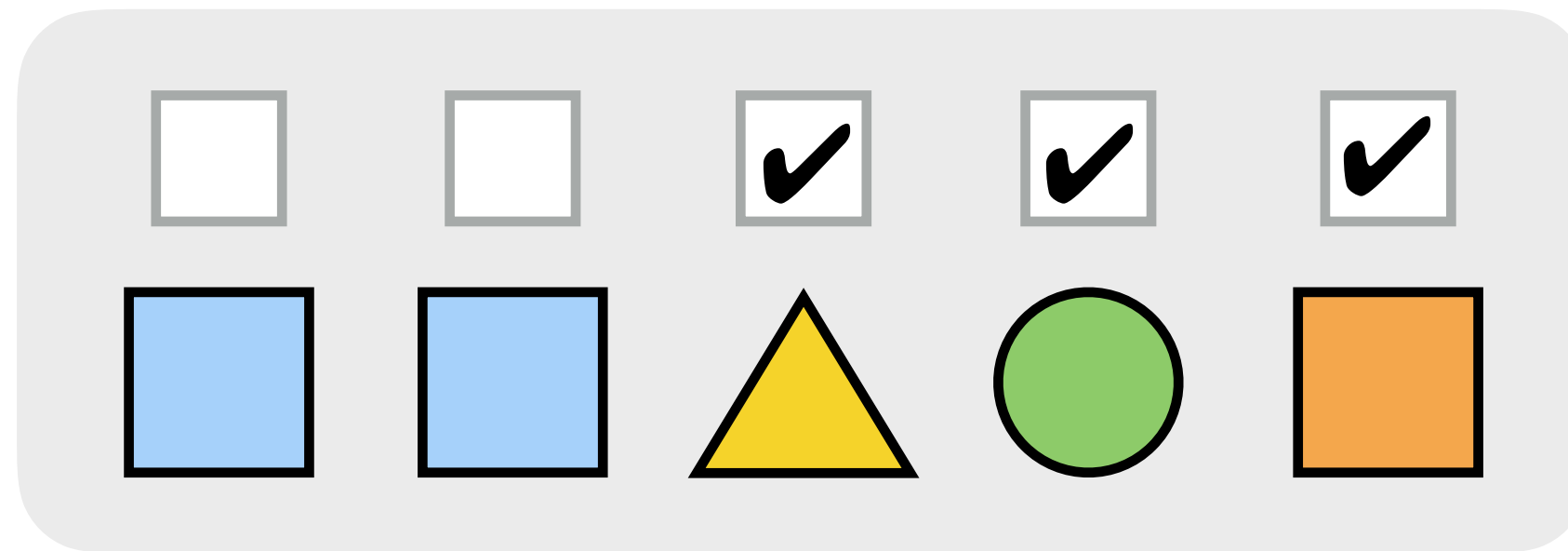
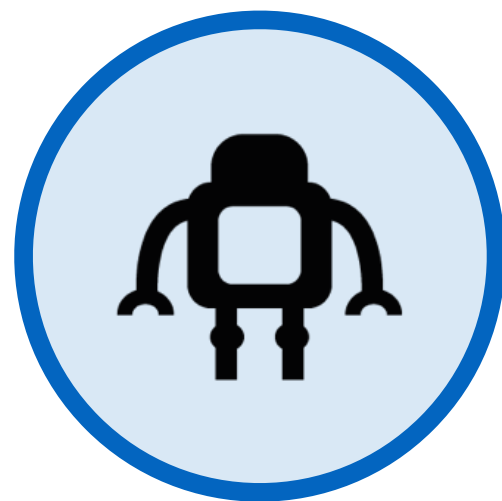




Compositional semantics

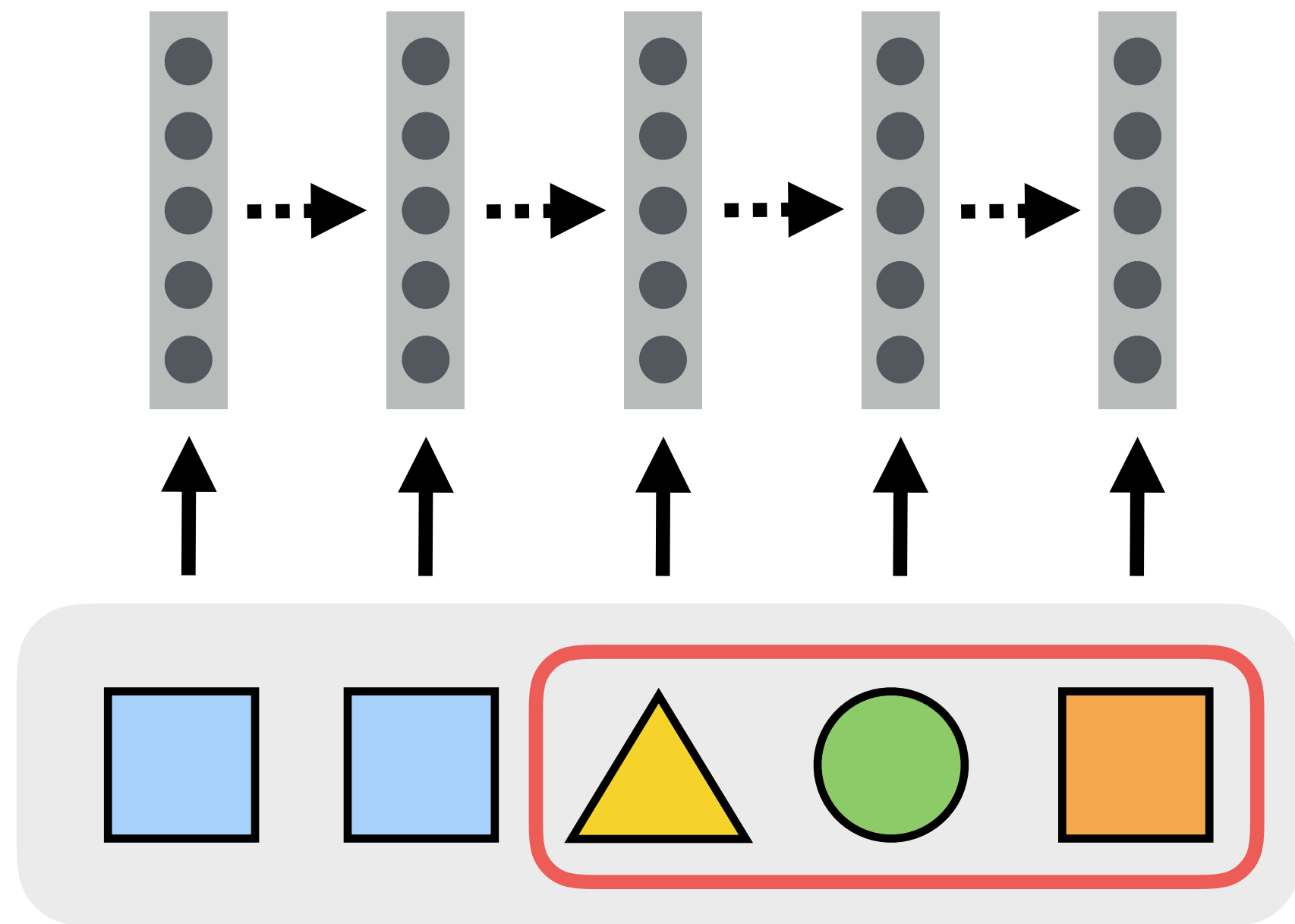


???



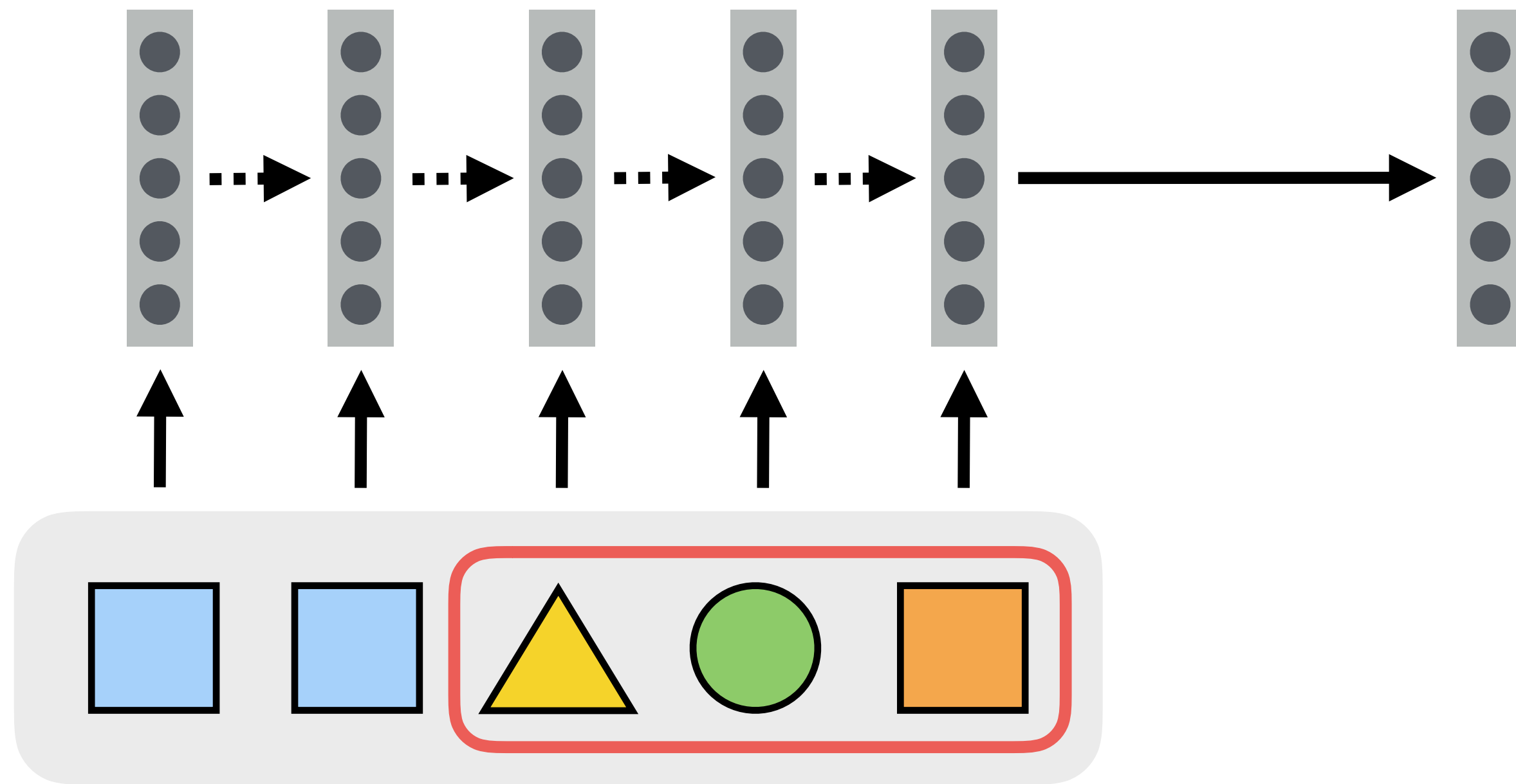


Model architecture



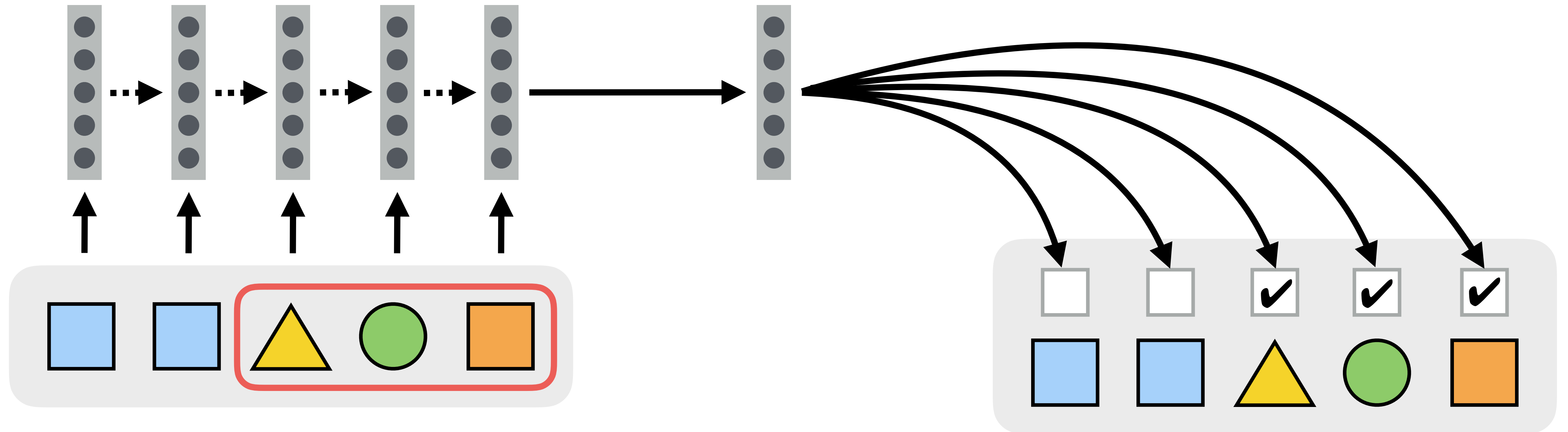


Model architecture



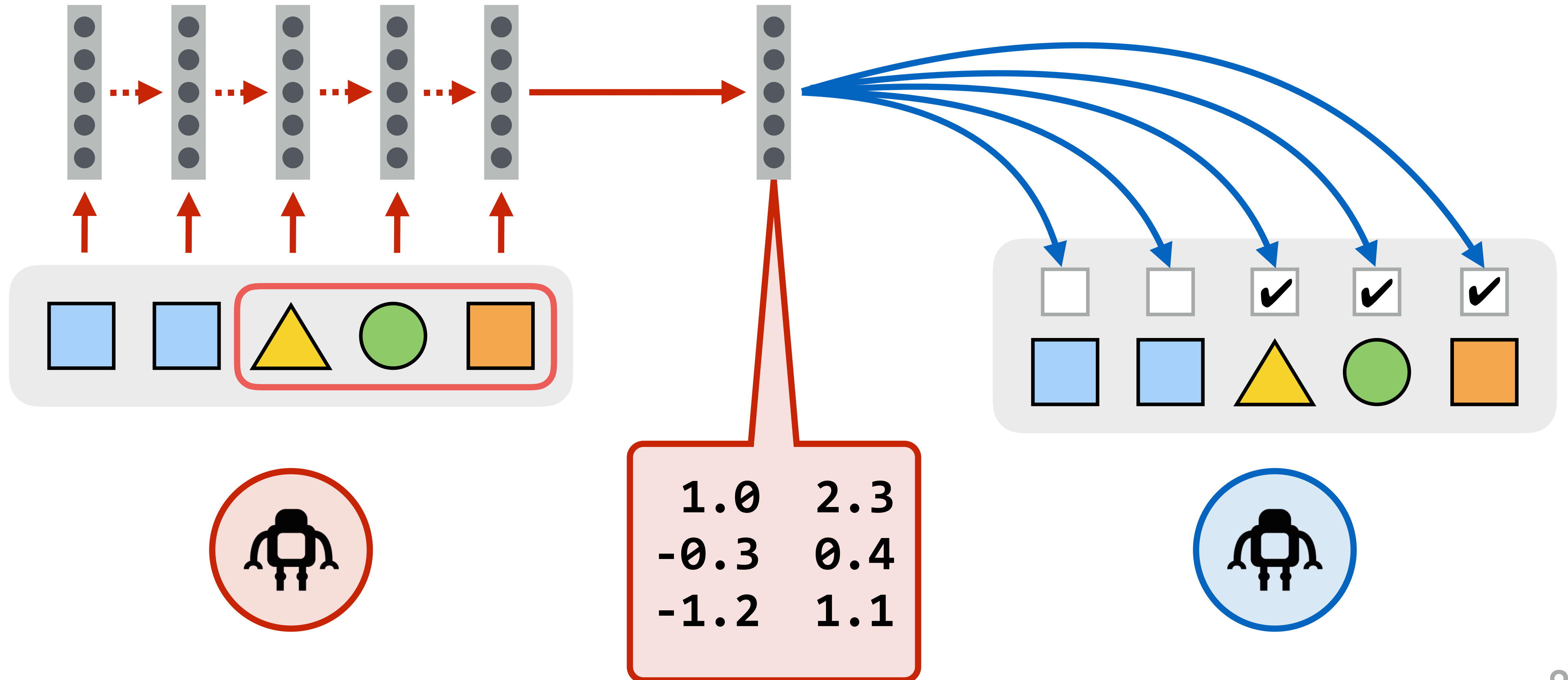


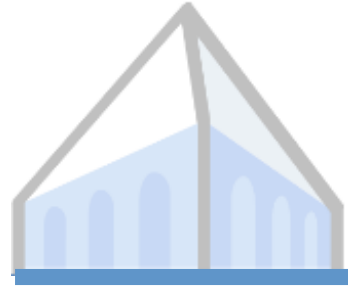
Model architecture





Model architecture

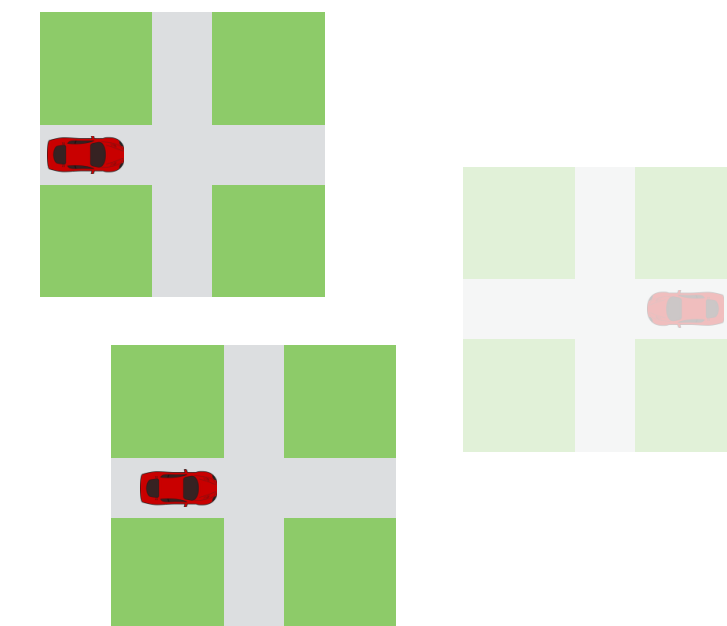
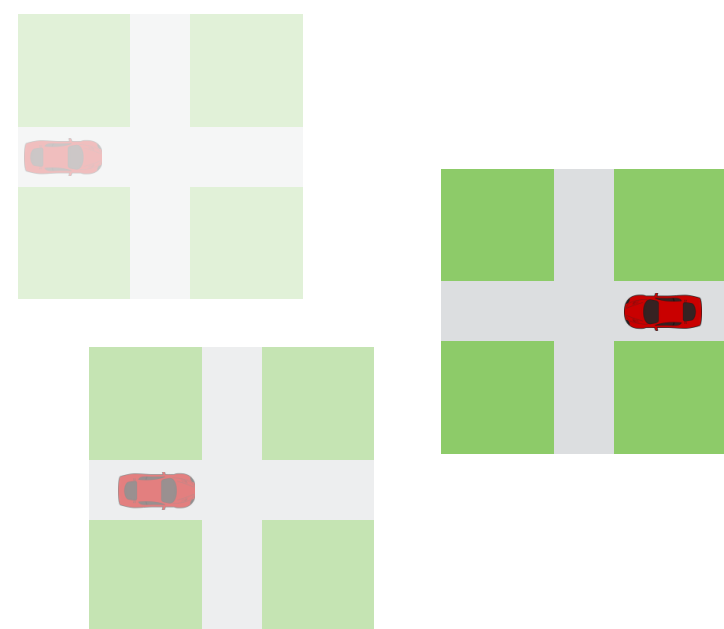


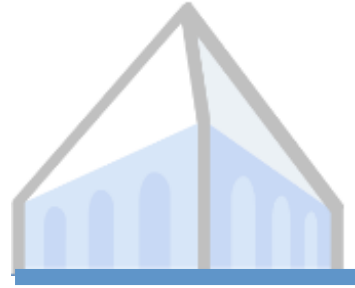


Computing meaning representations

1.0 2.3
-0.3 0.4
-1.2 1.1

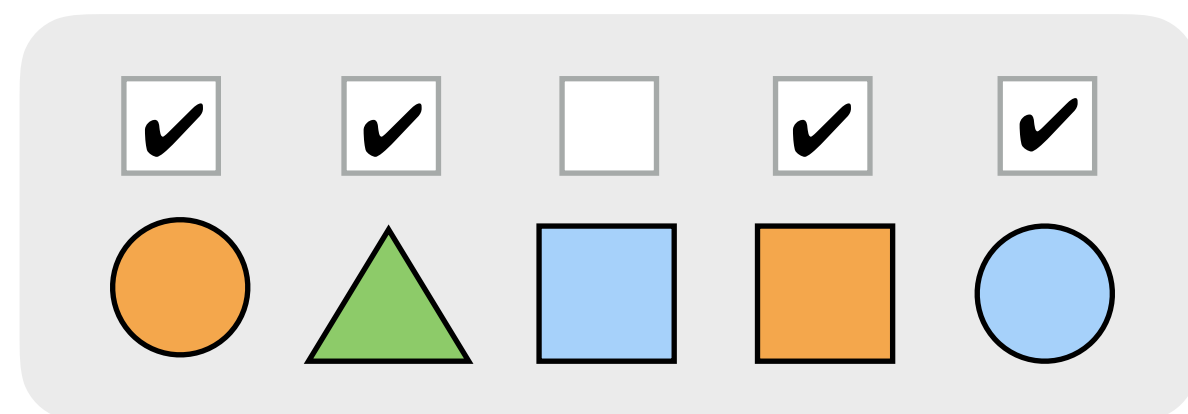
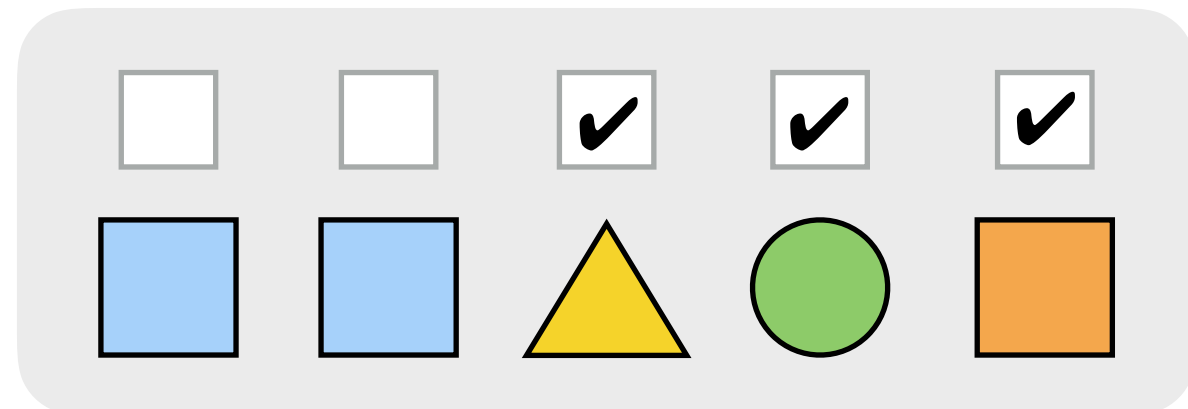
on the left



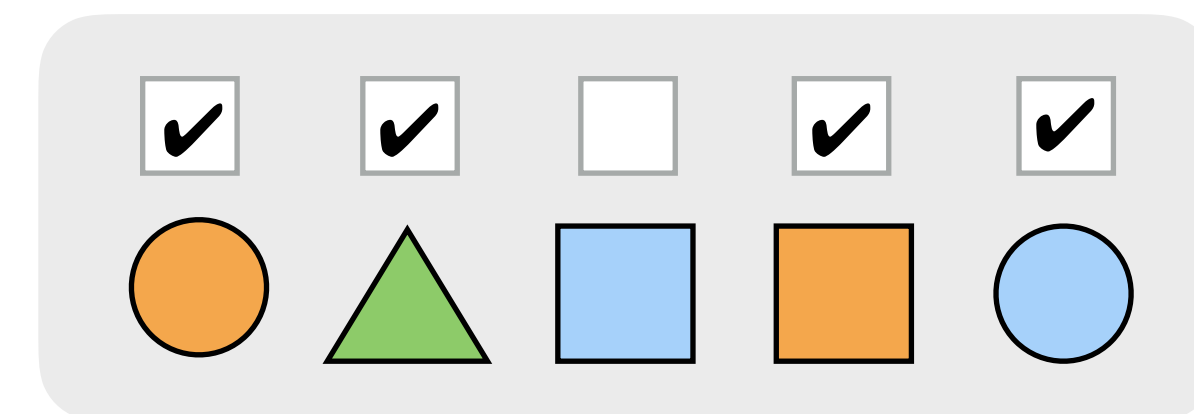
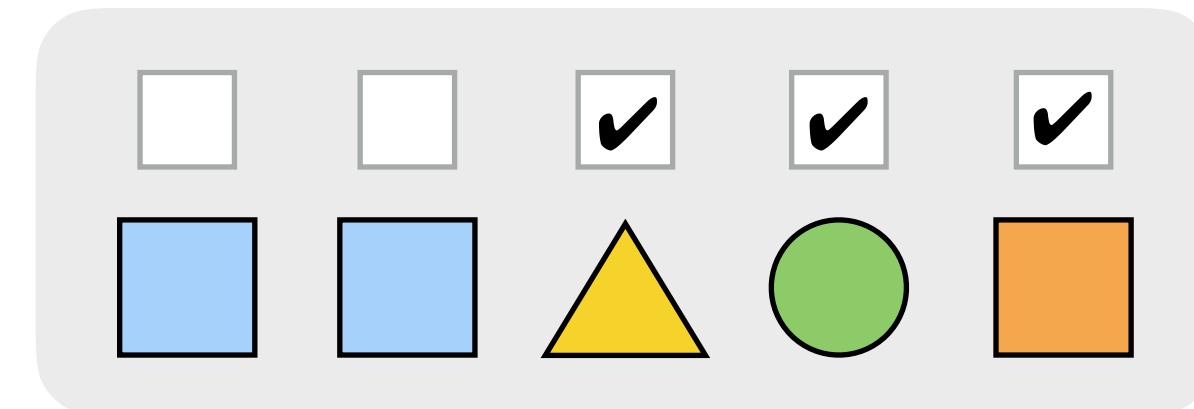


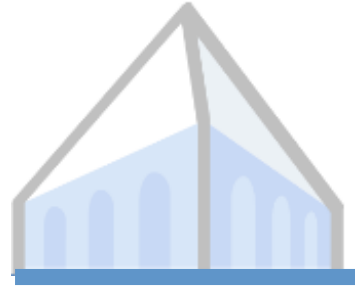
Computing meaning representations

-0.1 1.3
0.5 -0.4
0.2 1.0



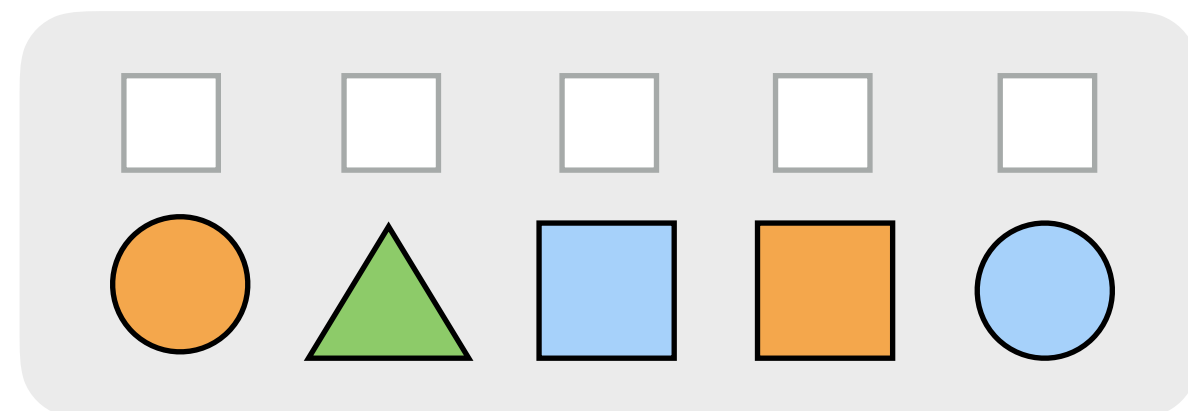
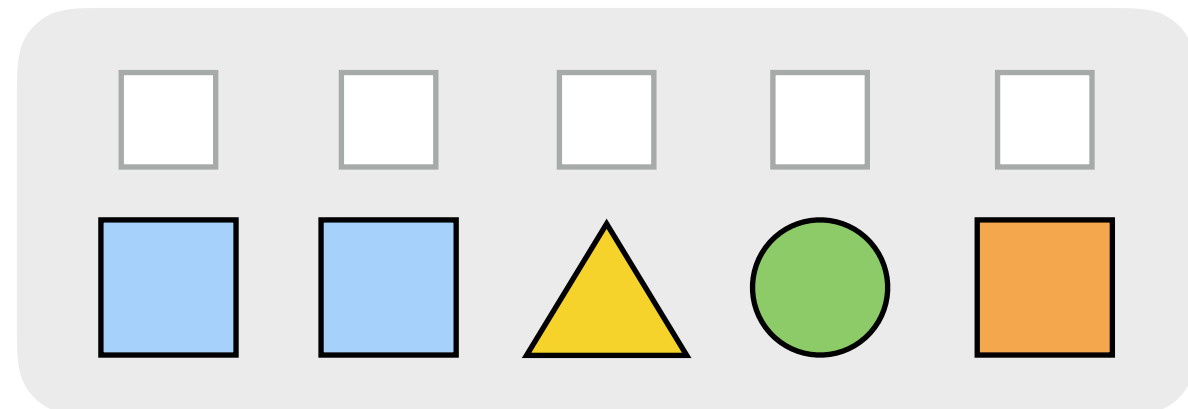
*everything but
squares*





Computing meaning representations

-0.1 1.3
0.5 -0.4
0.2 1.0

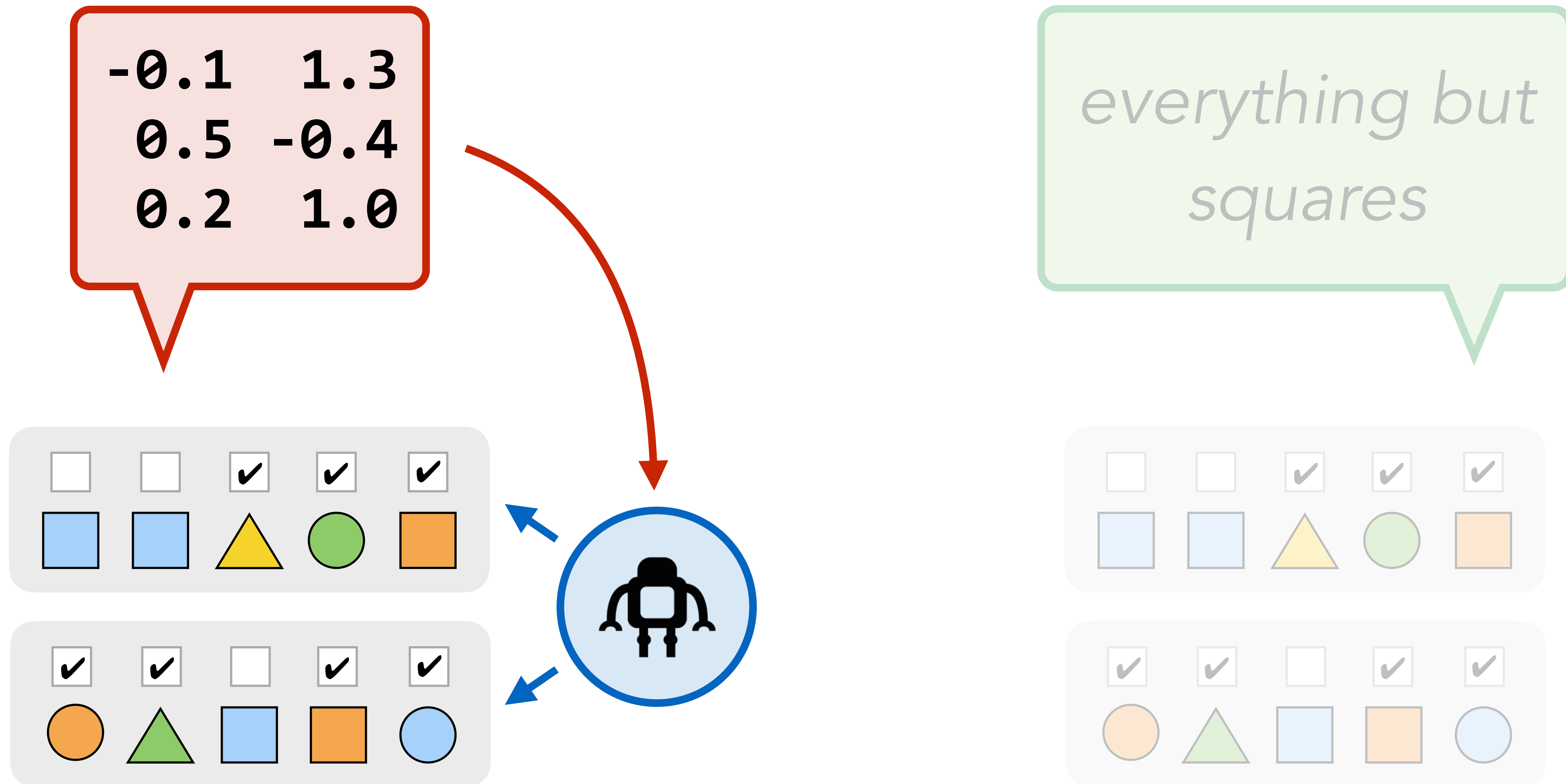


*everything but
squares*



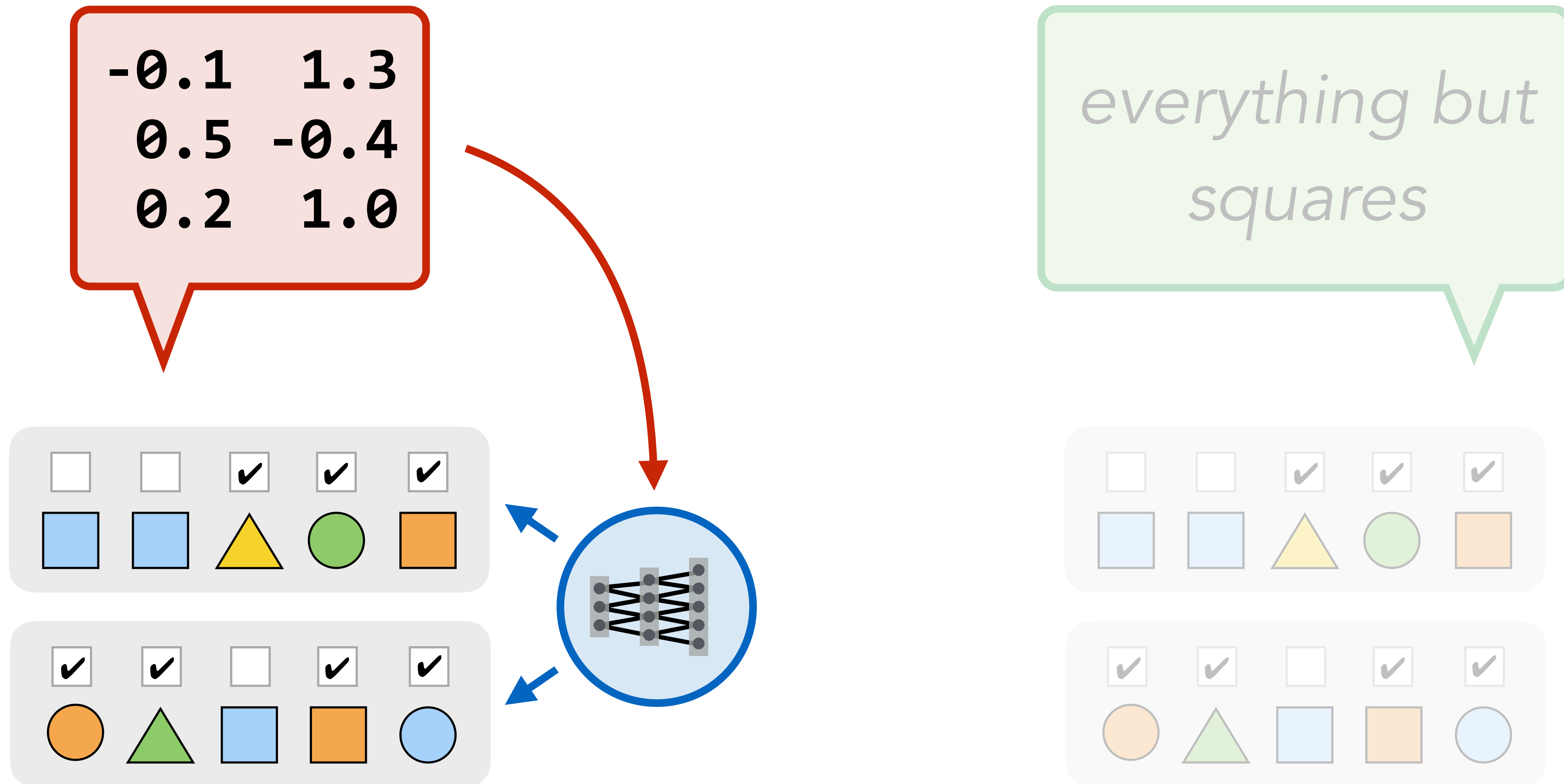


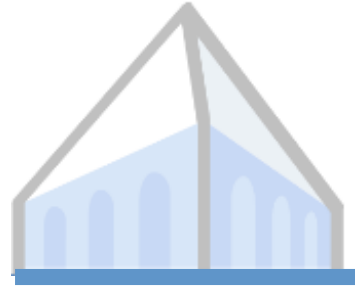
Computing meaning representations





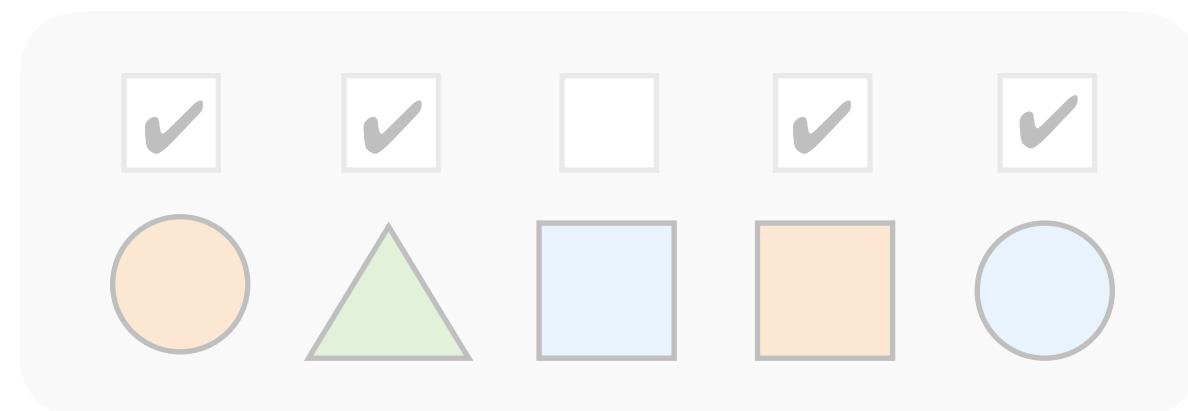
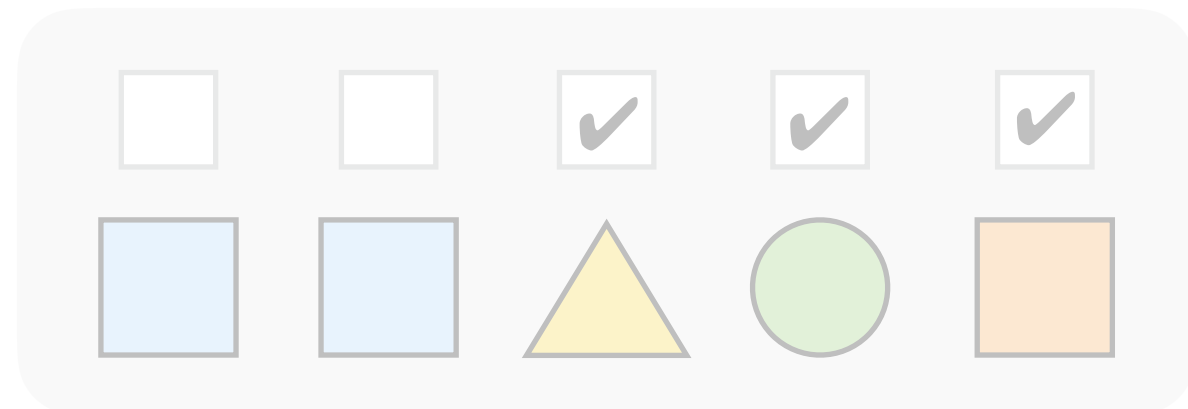
Computing meaning representations



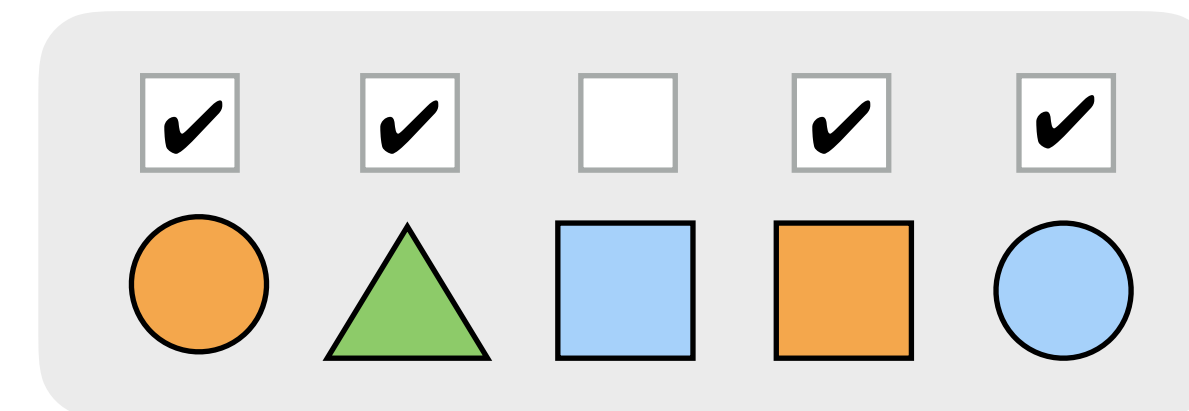
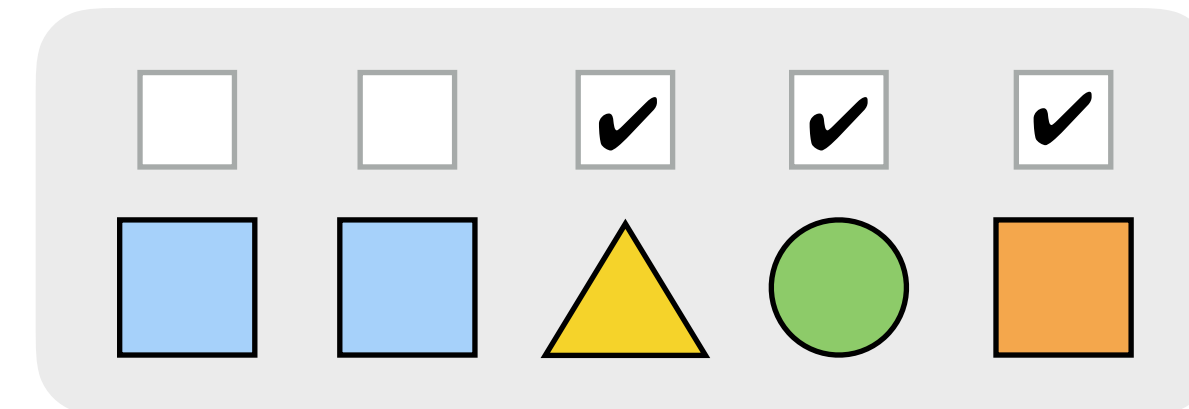


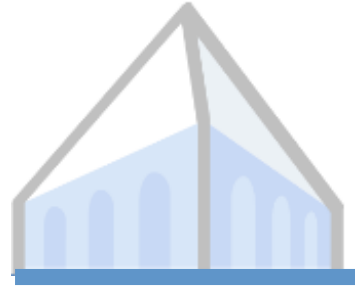
Computing meaning representations

-0.1 1.3
0.5 -0.4
0.2 1.0

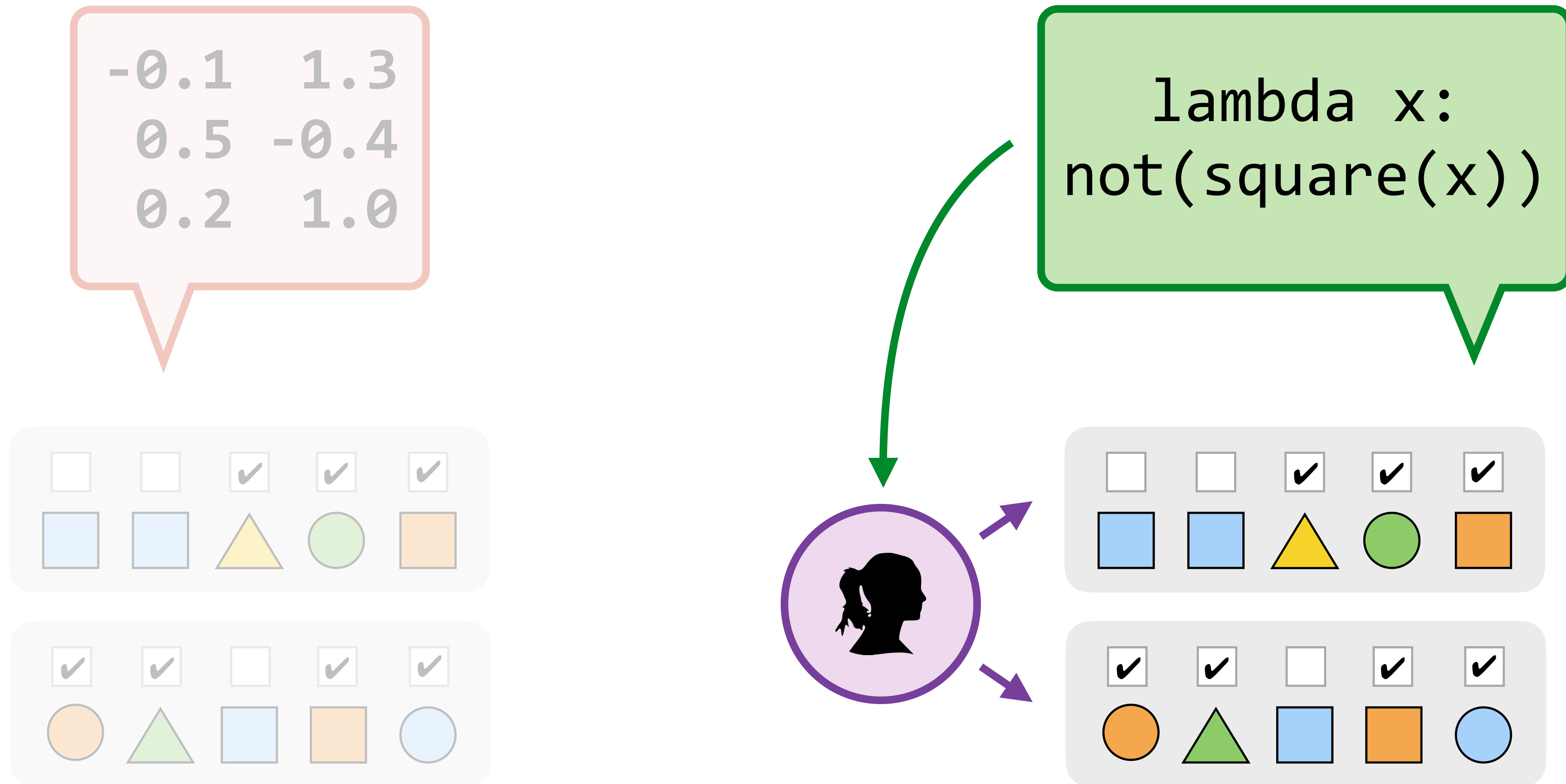


*everything but
squares*



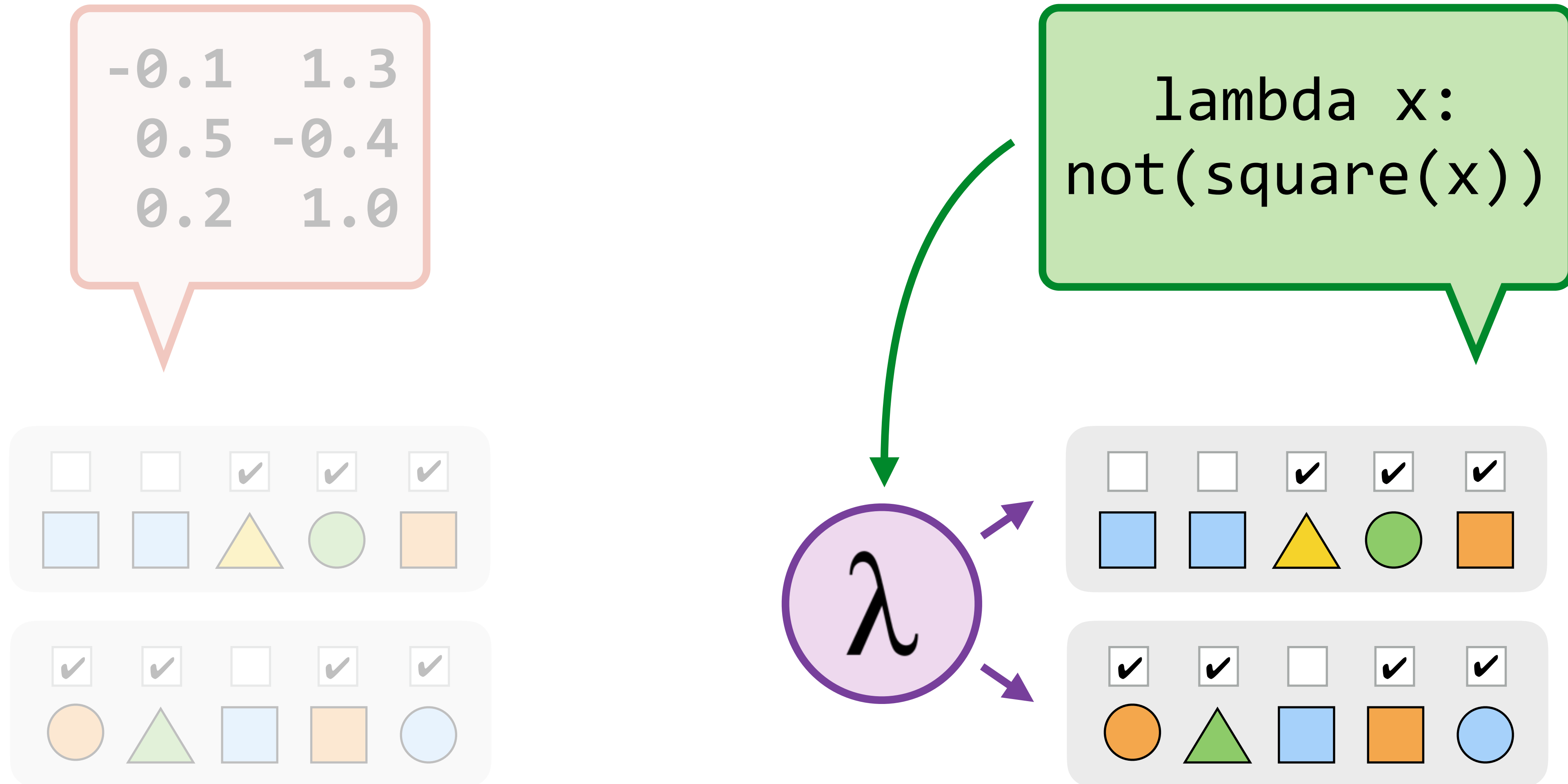


Computing meaning representations





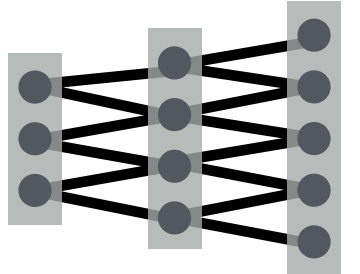
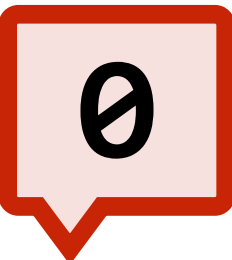
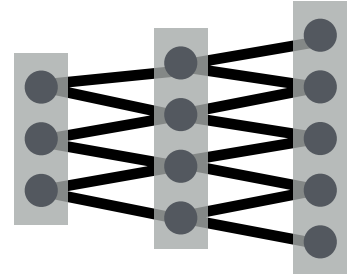
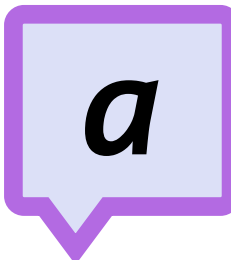
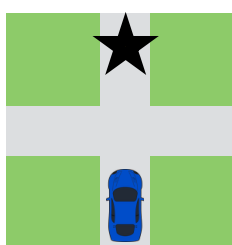
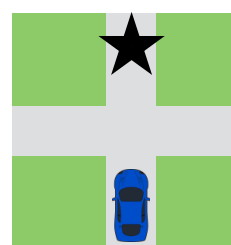
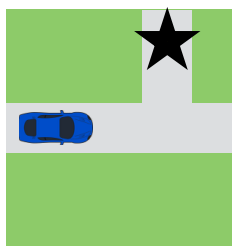
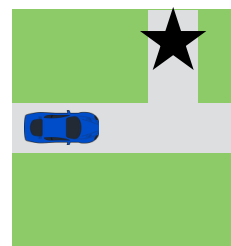
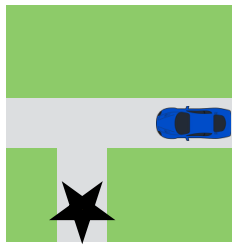
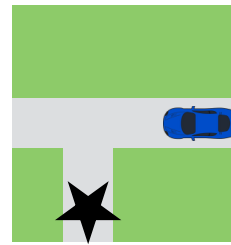
Computing meaning representations





Translation criterion

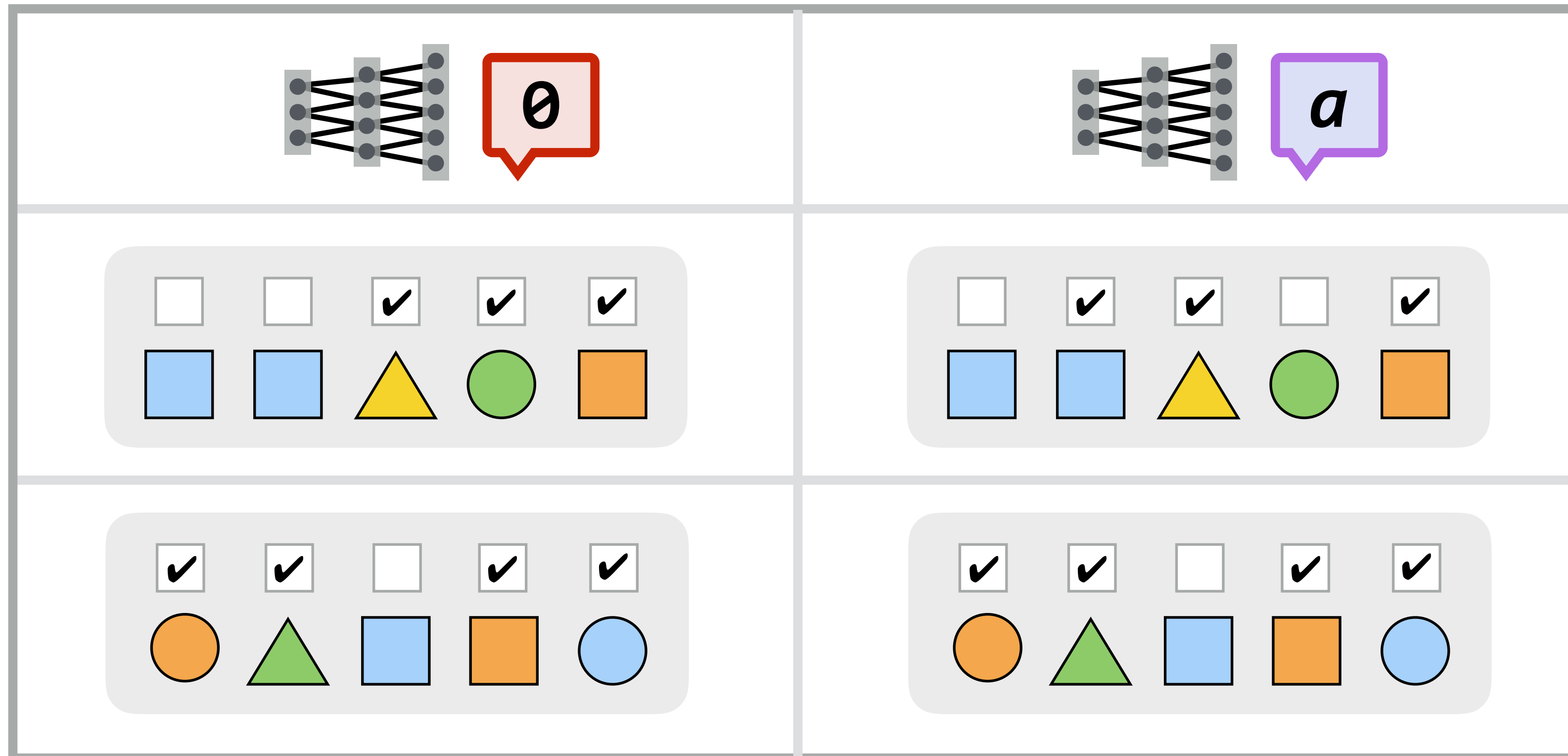
$$q(\theta, a) = \text{KL}(\beta(\theta) \parallel \beta(a))$$

 	 
 0.10	 0.08
 0.05	 0.01
 0.13	 0.22



Translation criterion

$$q(\theta, a) = \mathbf{E}[\beta(\theta) = \beta(a)]$$





Experiments

“High-level” communicative behavior

“Low-level” message structure



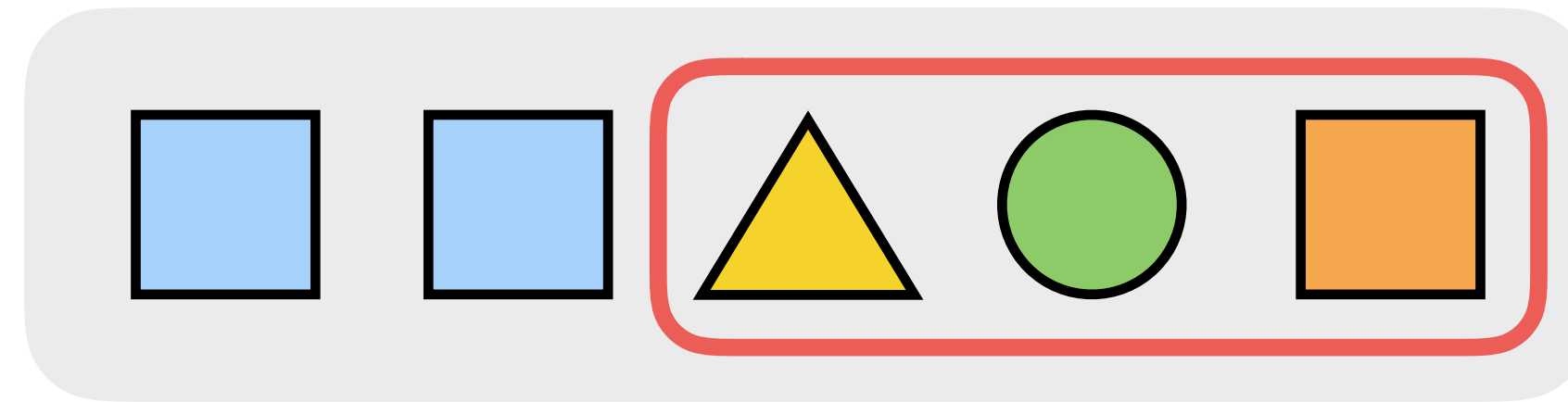
Experiments

“High-level” communicative behavior

“Low-level” message structure



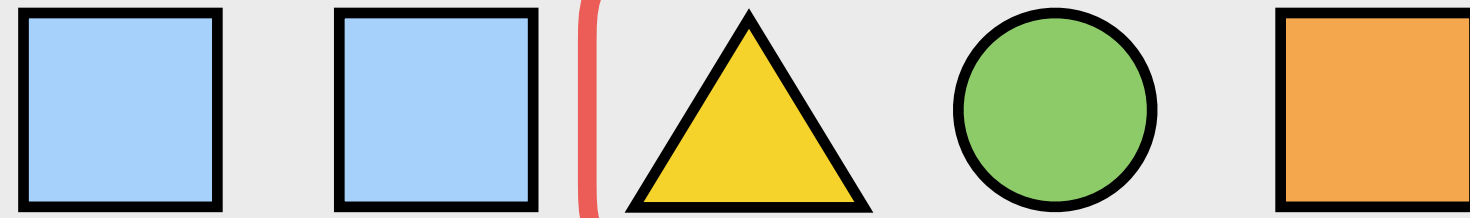
Comparing strategies



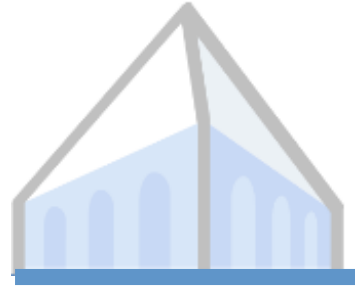


Comparing strategies

-0.1 1.3
 0.5 -0.4
 0.2 1.0

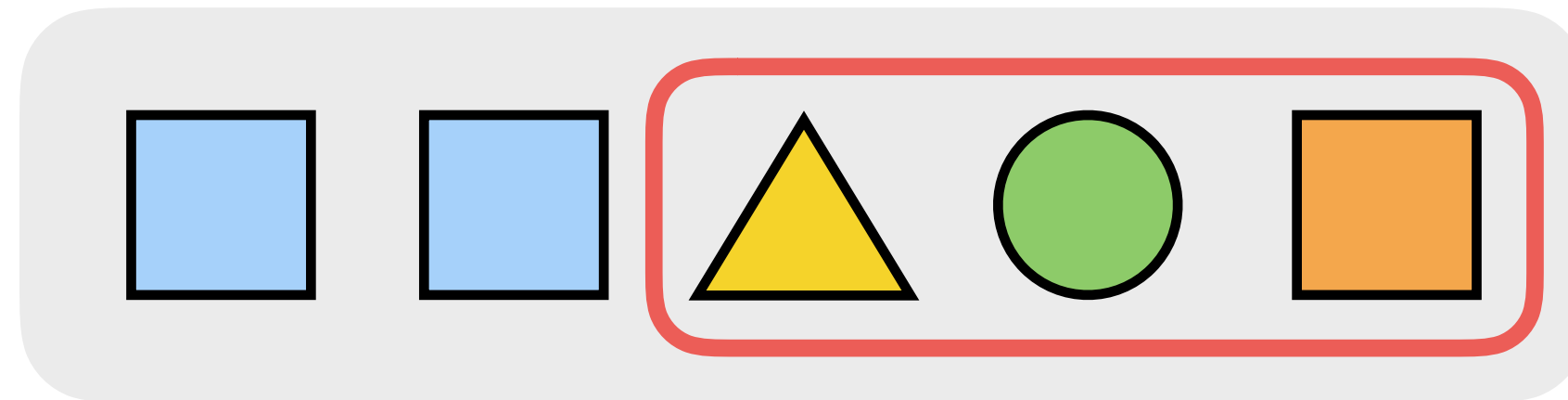


*everything
but squares*

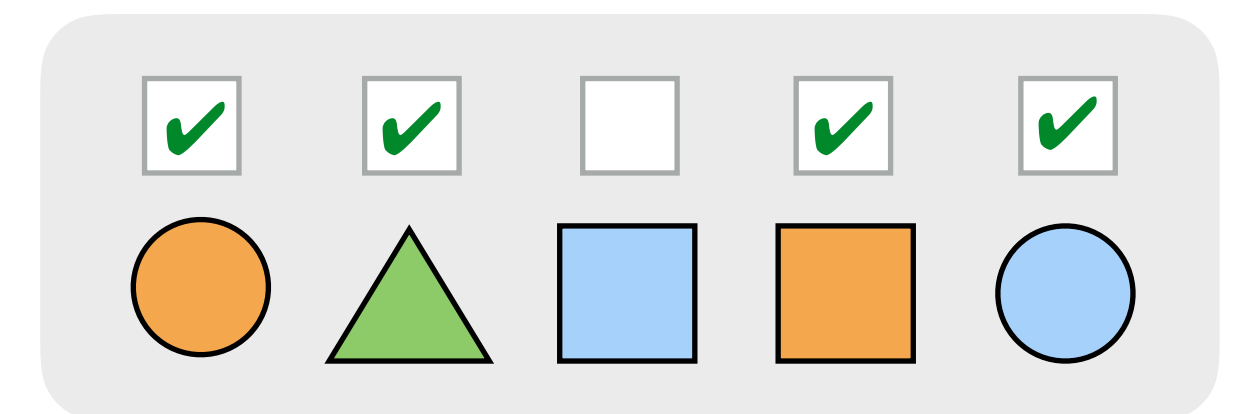
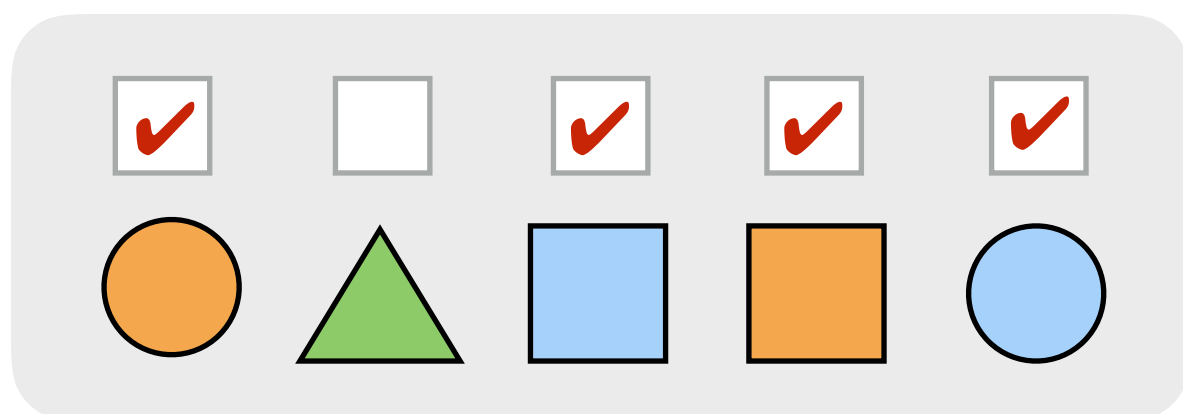
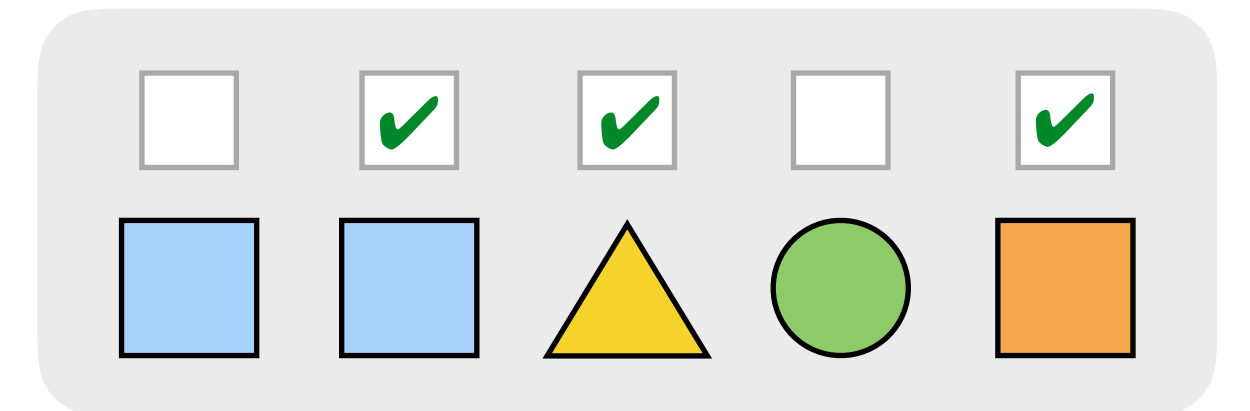
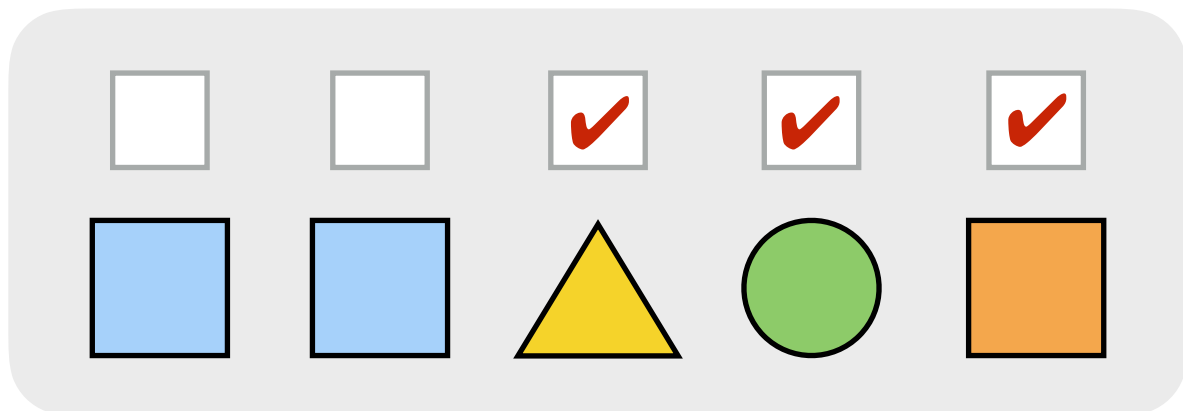


Comparing strategies

-0.1	1.3
0.5	-0.4
0.2	1.0



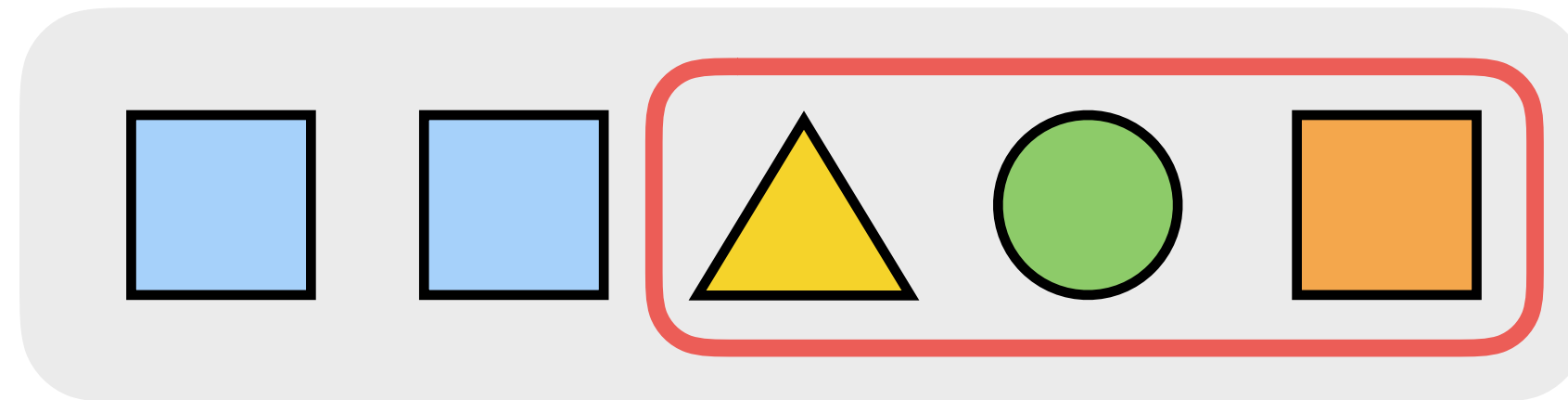
*everything
but squares*



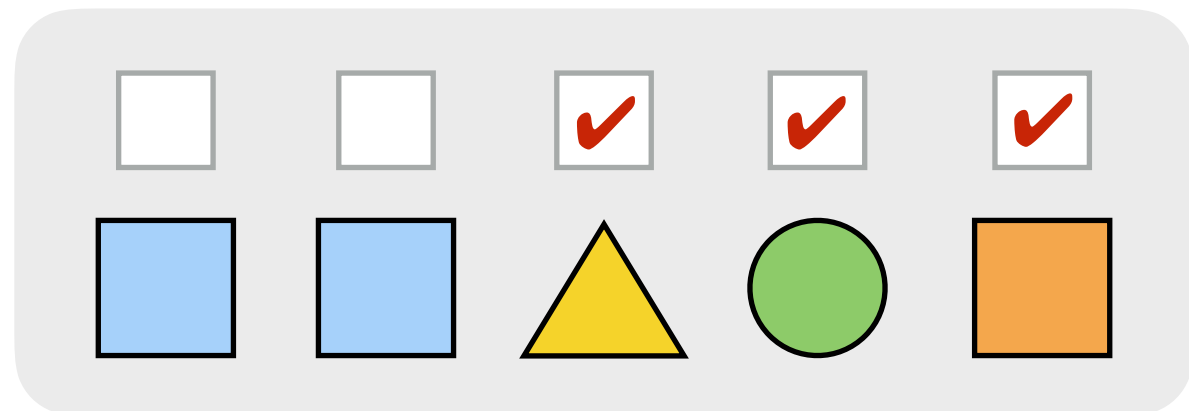


Comparing strategies

-0.1	1.3
0.5	-0.4
0.2	1.0

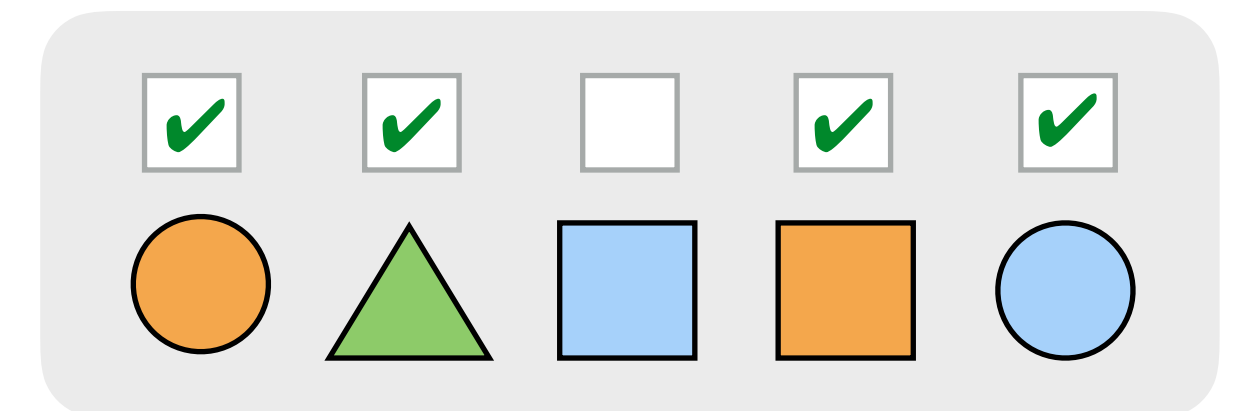
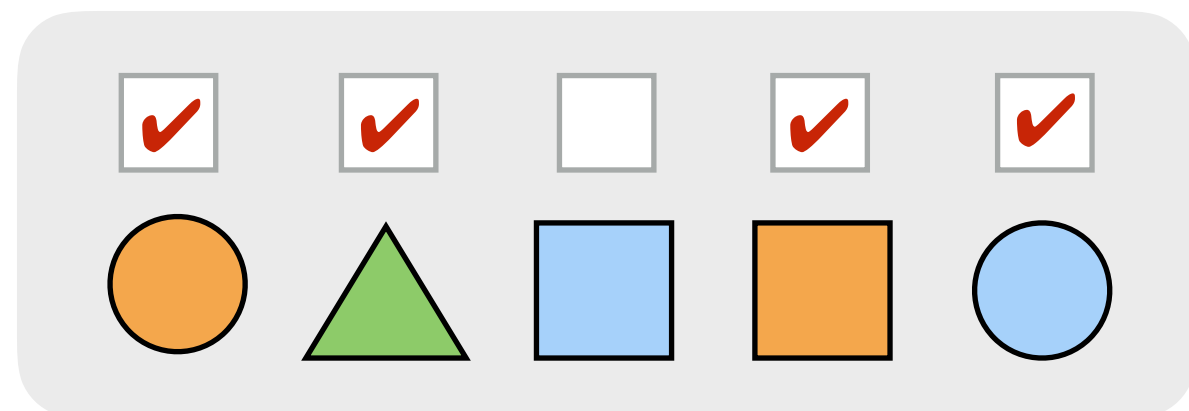
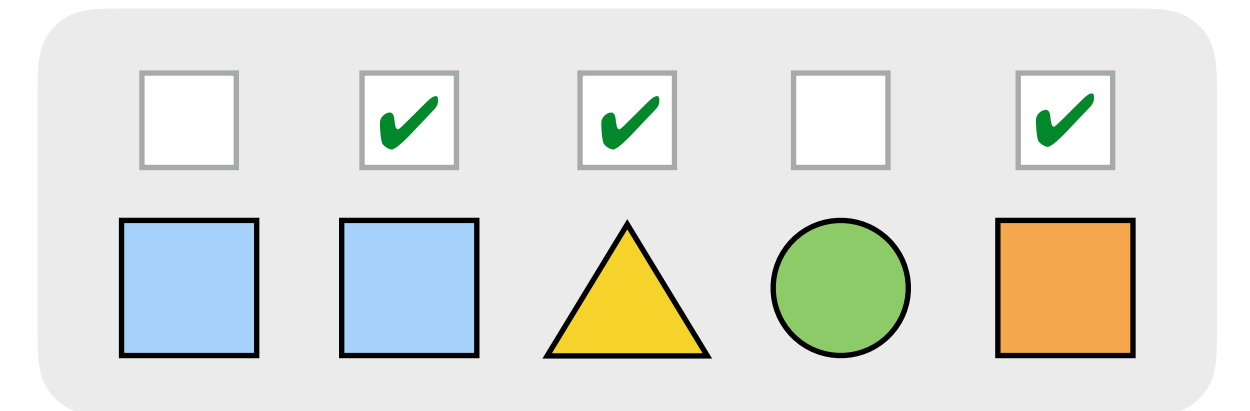


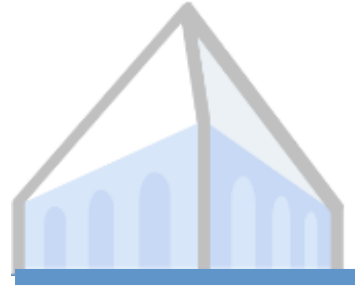
*everything
but squares*



?

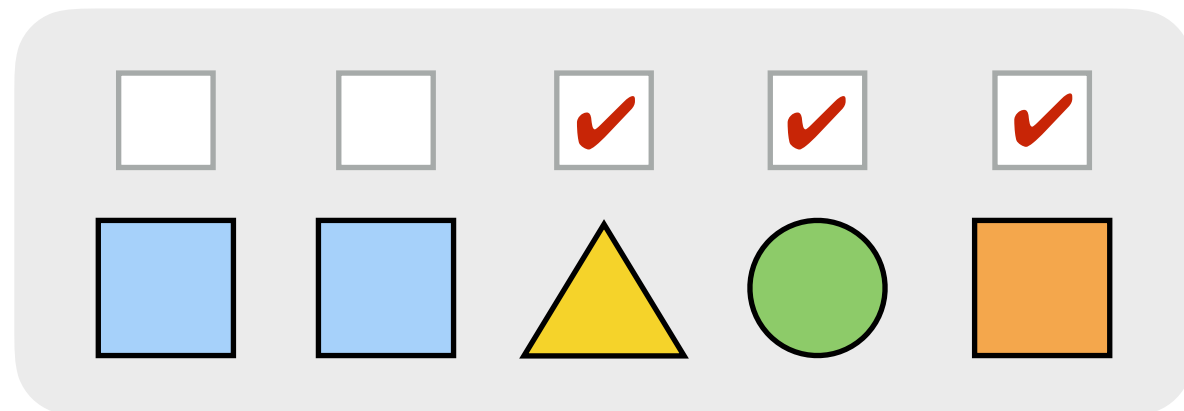
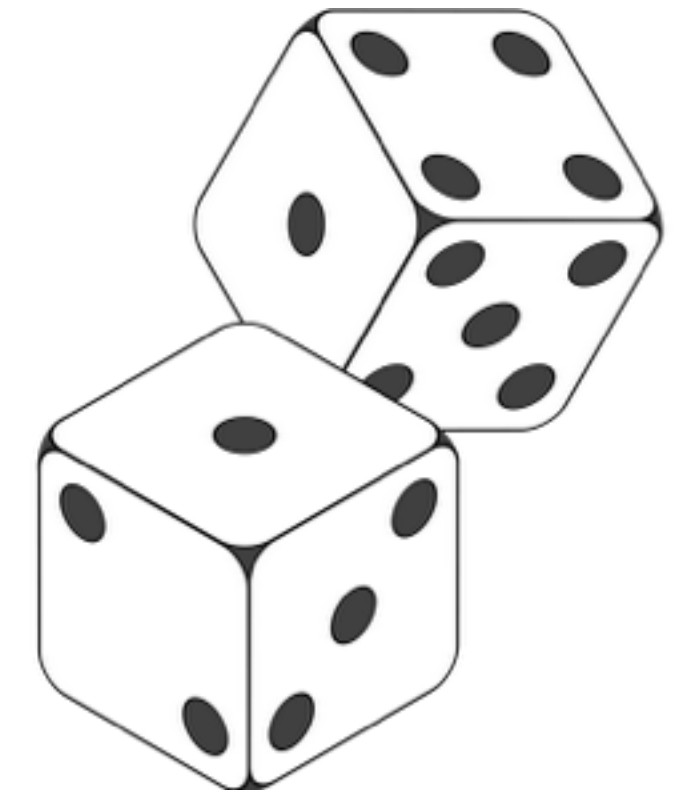
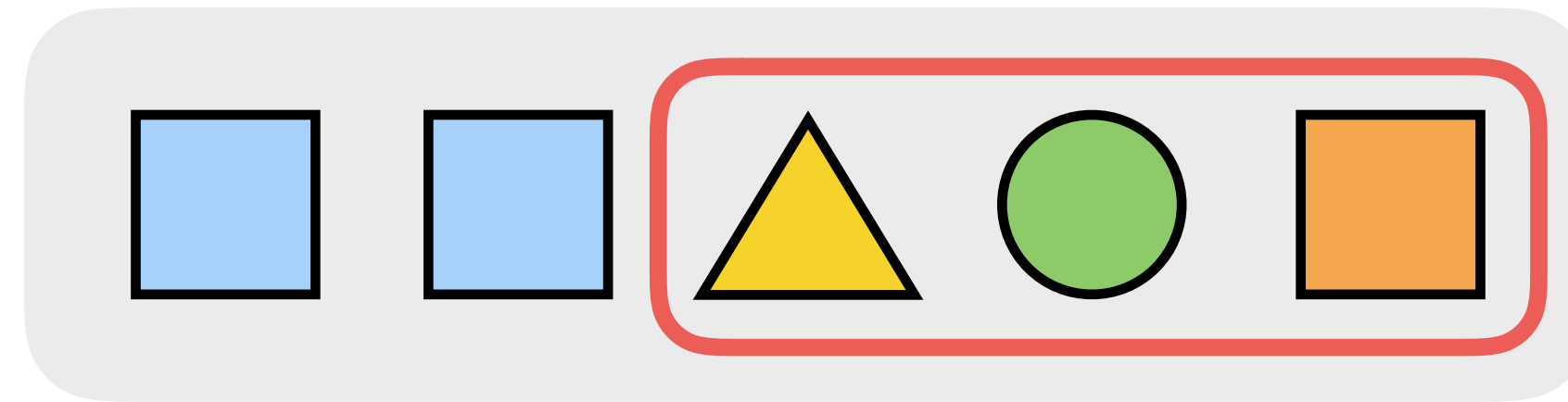
==





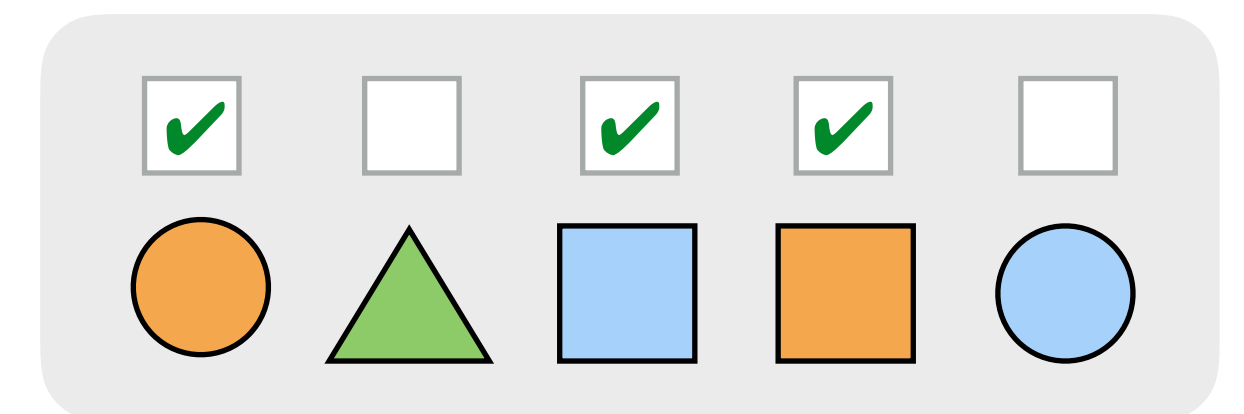
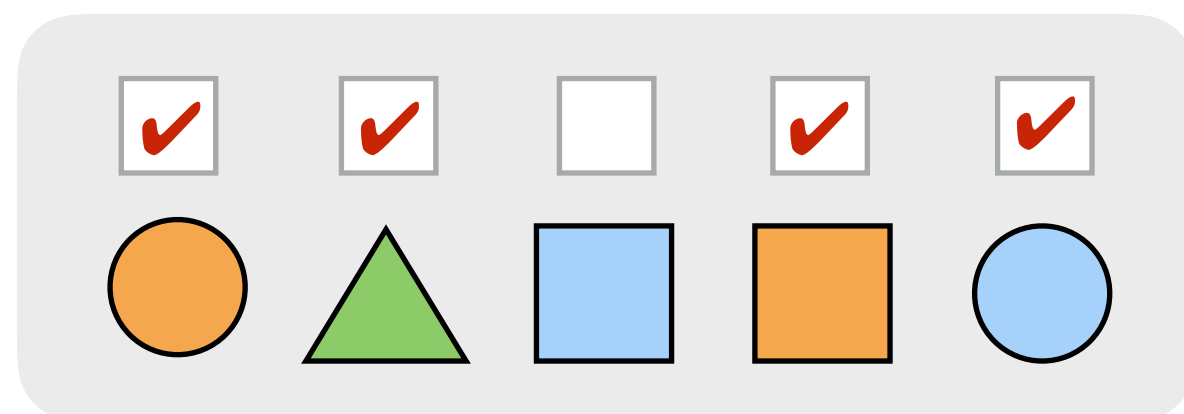
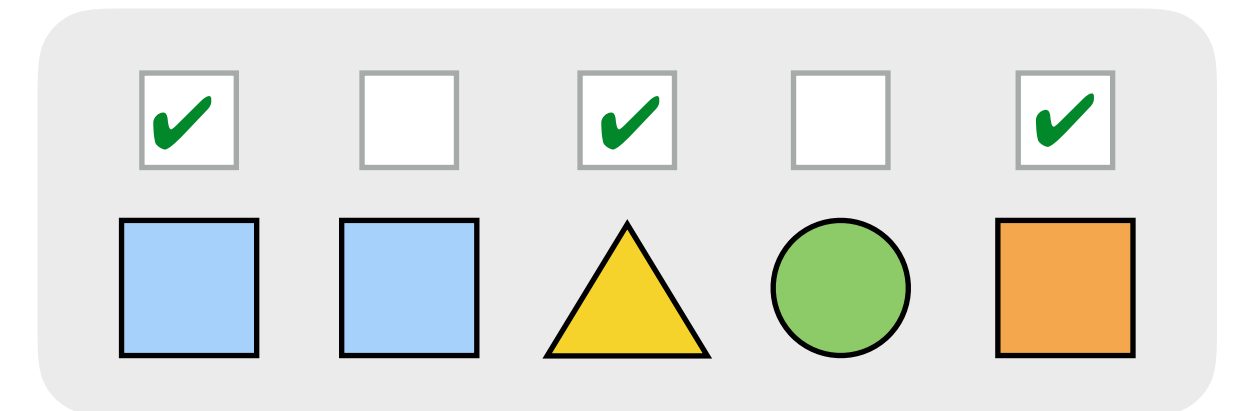
Theories of model behavior: random

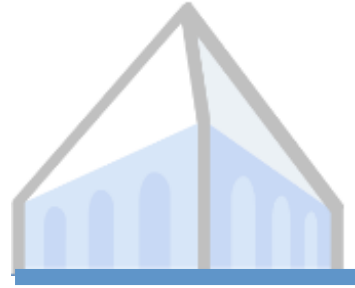
-0.1	1.3
0.5	-0.4
0.2	1.0



?

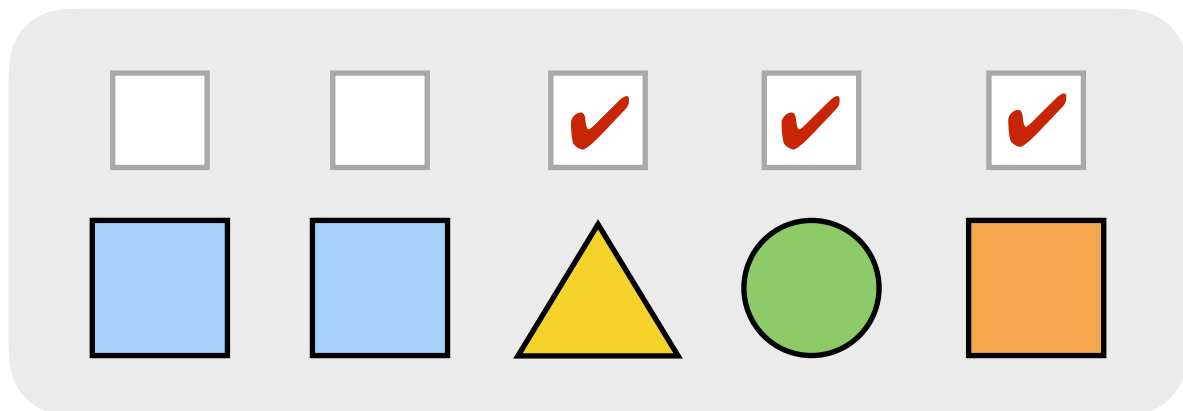
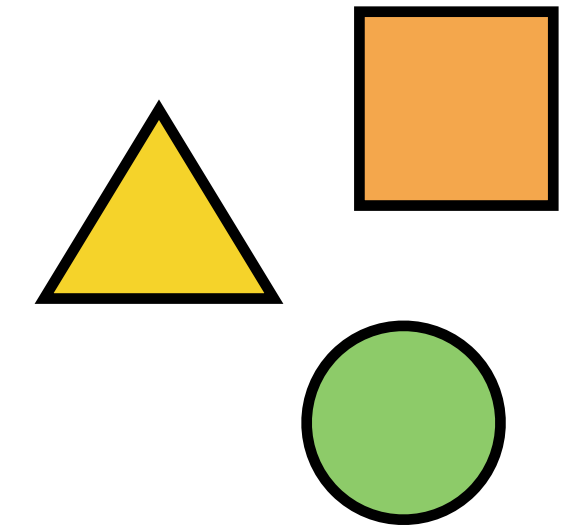
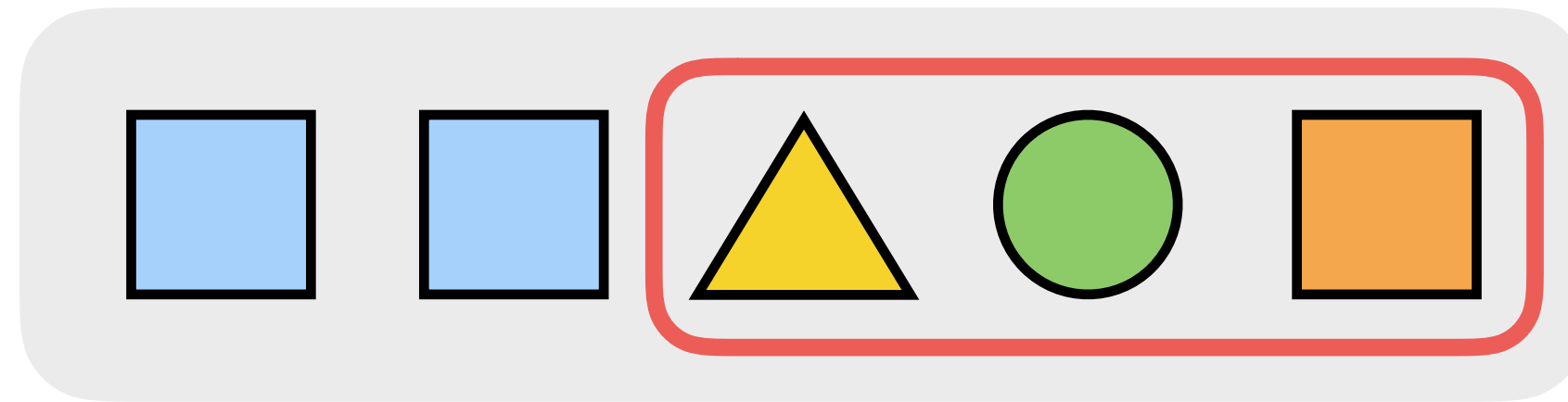
==





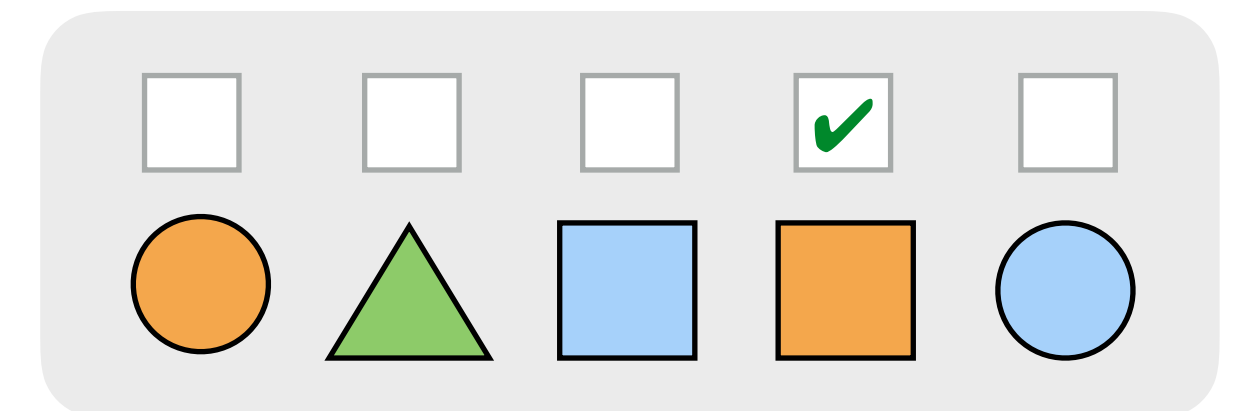
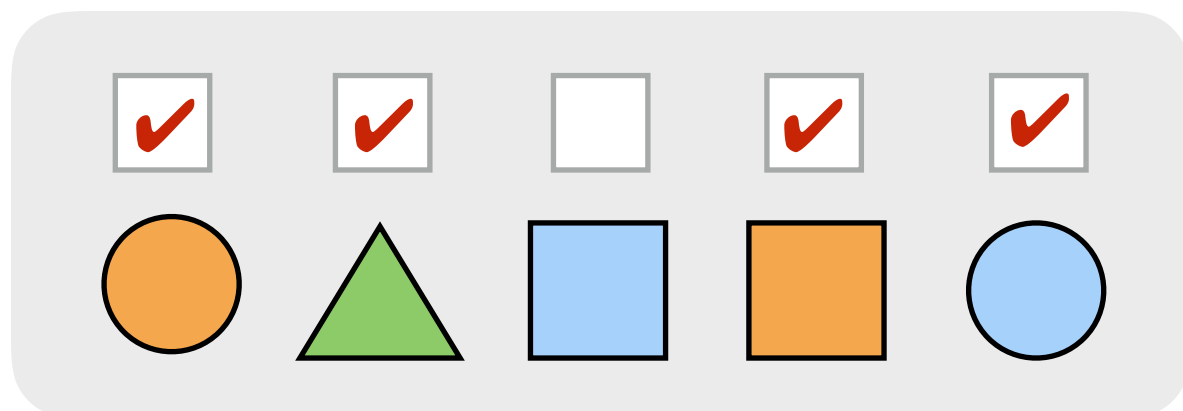
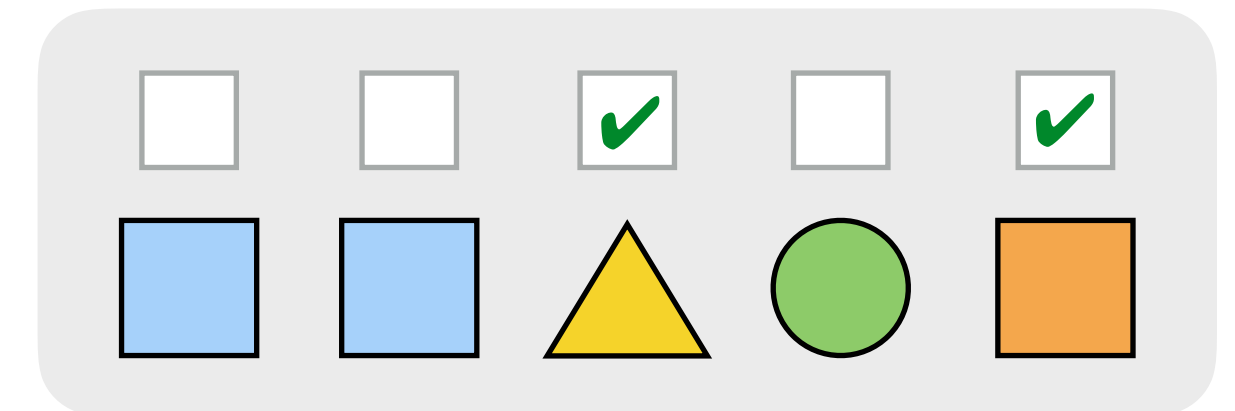
Theories of model behavior: literal

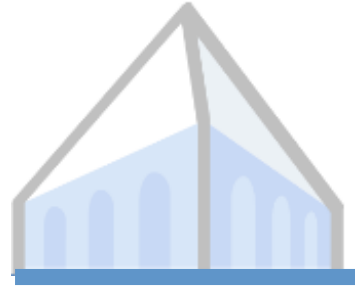
-0.1	1.3
0.5	-0.4
0.2	1.0



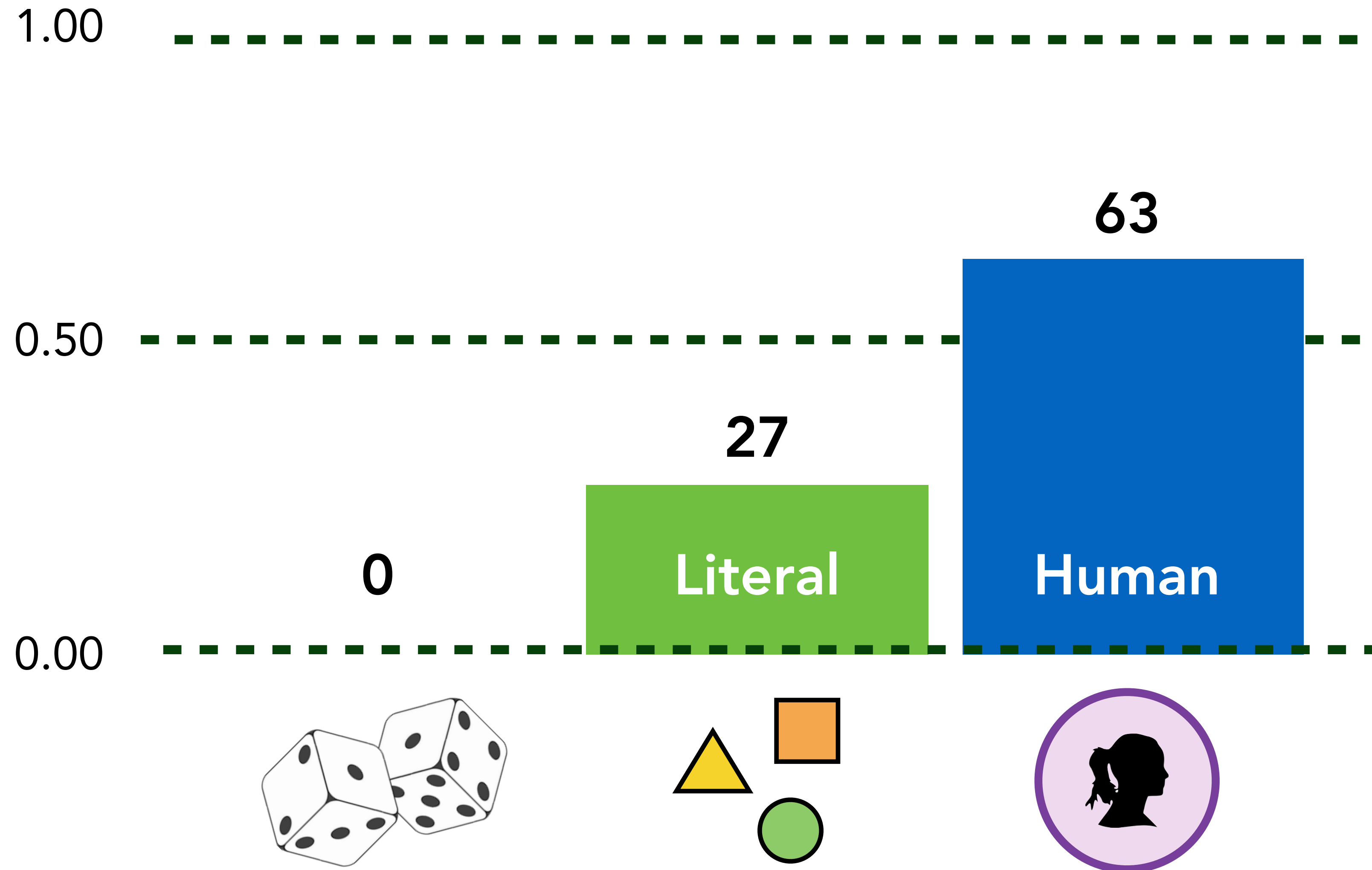
?

==



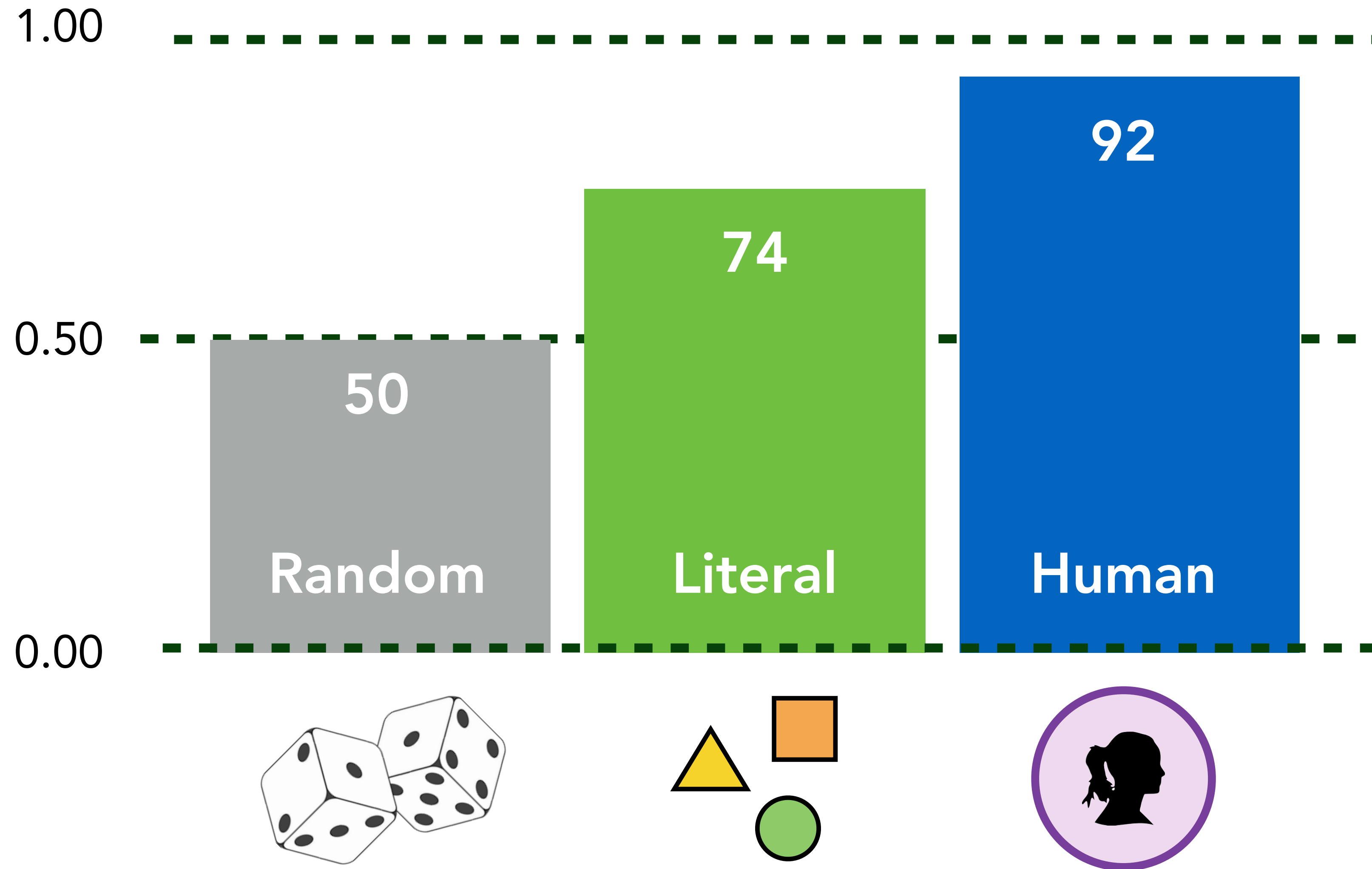


Evaluation: high-level scene agreement





Evaluation: high-level object agreement

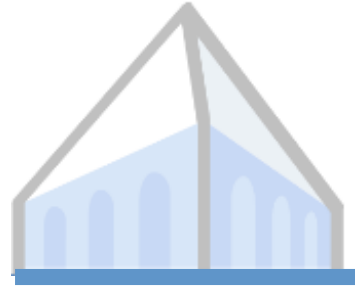




Experiments

“High-level” communicative behavior

“Low-level” message structure



Collecting translation data

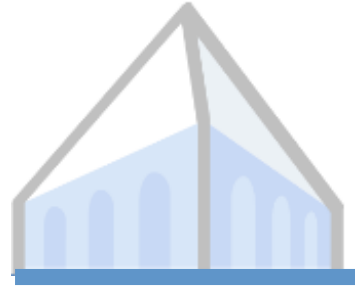
all the red shapes

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

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

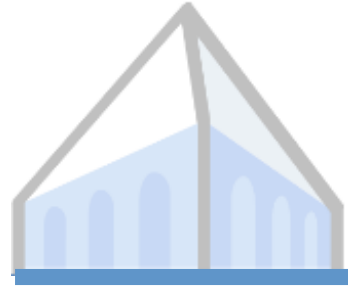


0.2 -0.2 0.5 -0.1

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



0.3 -1.3 -1.5 0.1



Extracting related pairs

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

0.1 -0.3 0.5 1.1

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

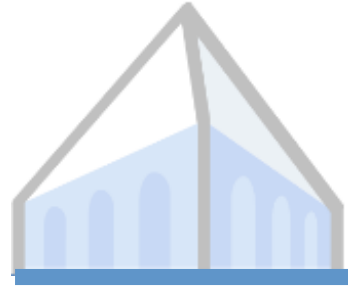
0.2 -0.2 0.5 -0.1

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

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

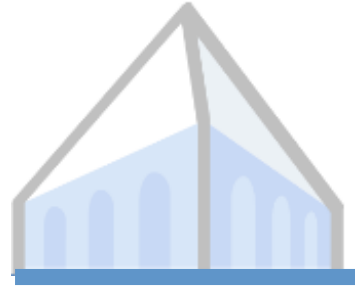
$\lambda x. \text{grn}(x) \wedge \text{sqr}(x)$

0.2 -0.2 0.5 -0.1

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

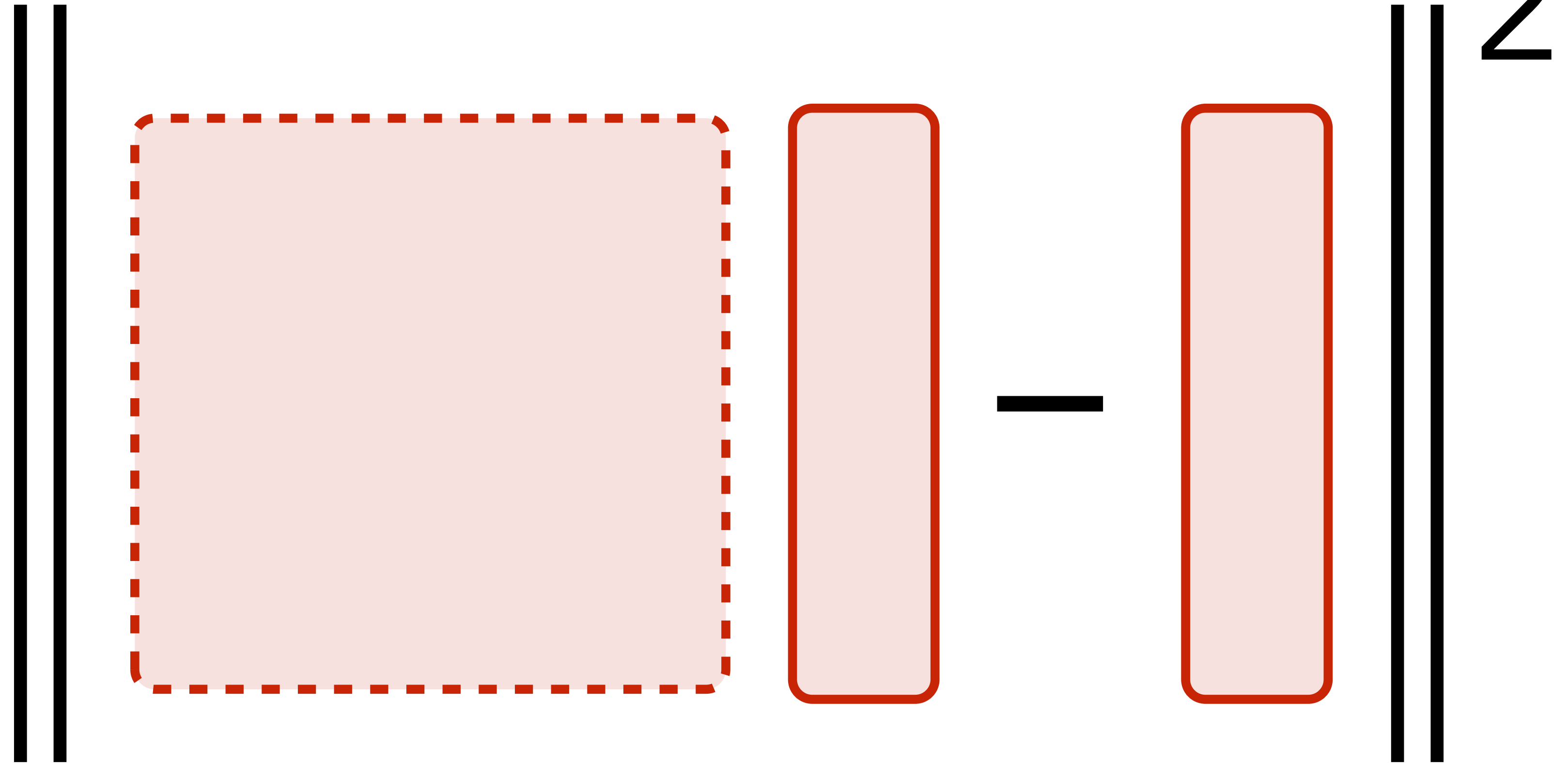
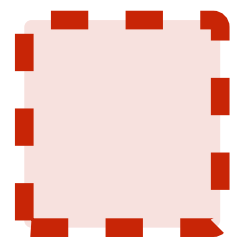
0.3 -1.3 -1.5 0.1





Learning compositional operators

argmin





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

$\lambda x. f(x)$

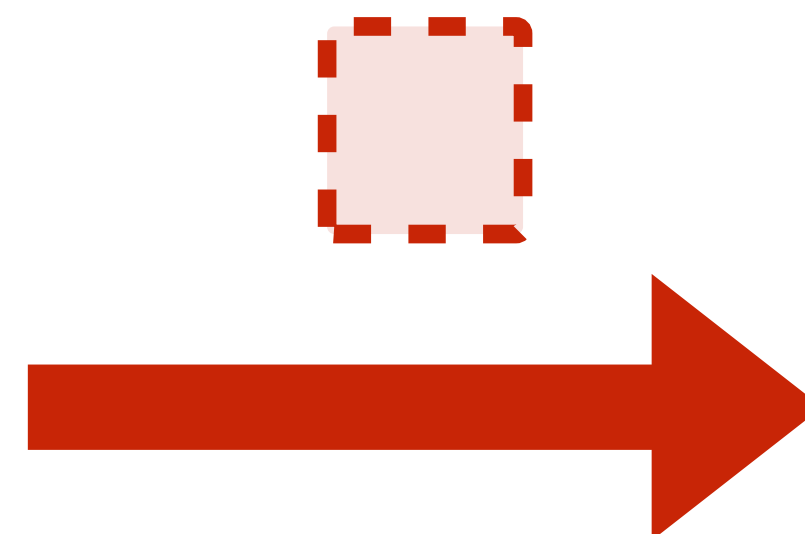
0.2 -0.2 0.5 -0.1

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

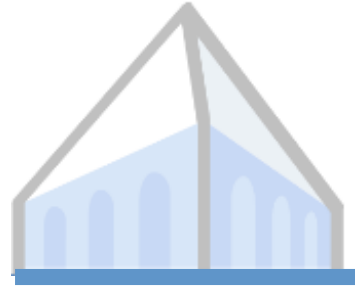
1.4 -0.3 -0.5 0.8

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

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

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

1.4 -0.3 -0.5 0.8

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

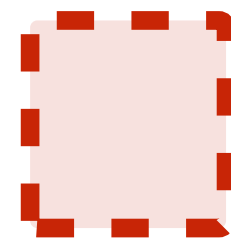
0.2 -0.2 0.5 -0.1

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

0.3 -1.3 -1.5 0.1

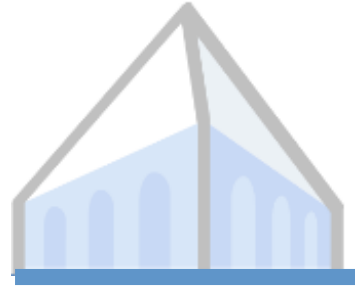
$\lambda x. f(x)$

0.2 -0.2 0.5 -0.1

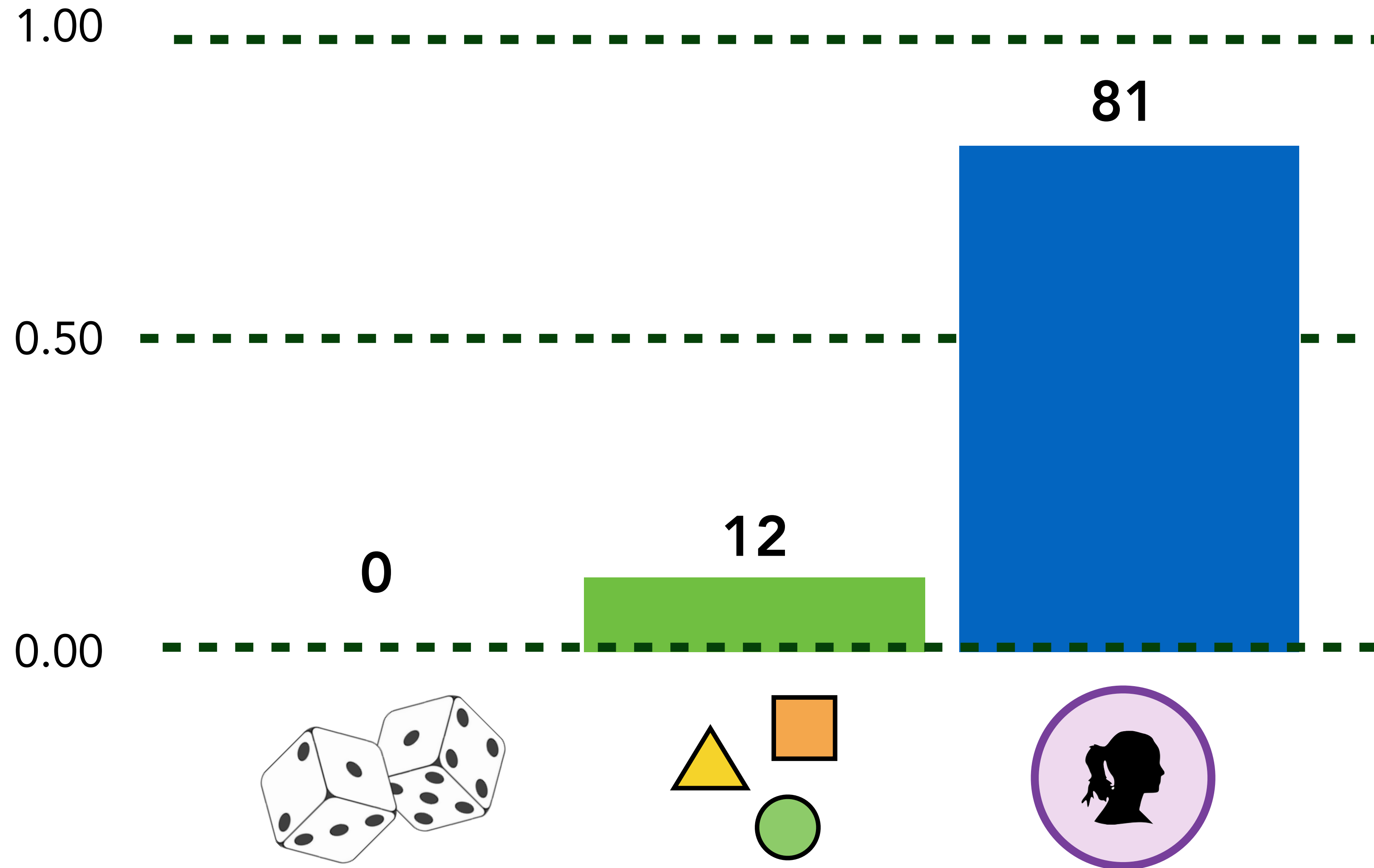


???

-0.2 0.4 -0.3 0.0



Evaluation: scene agreement for negation



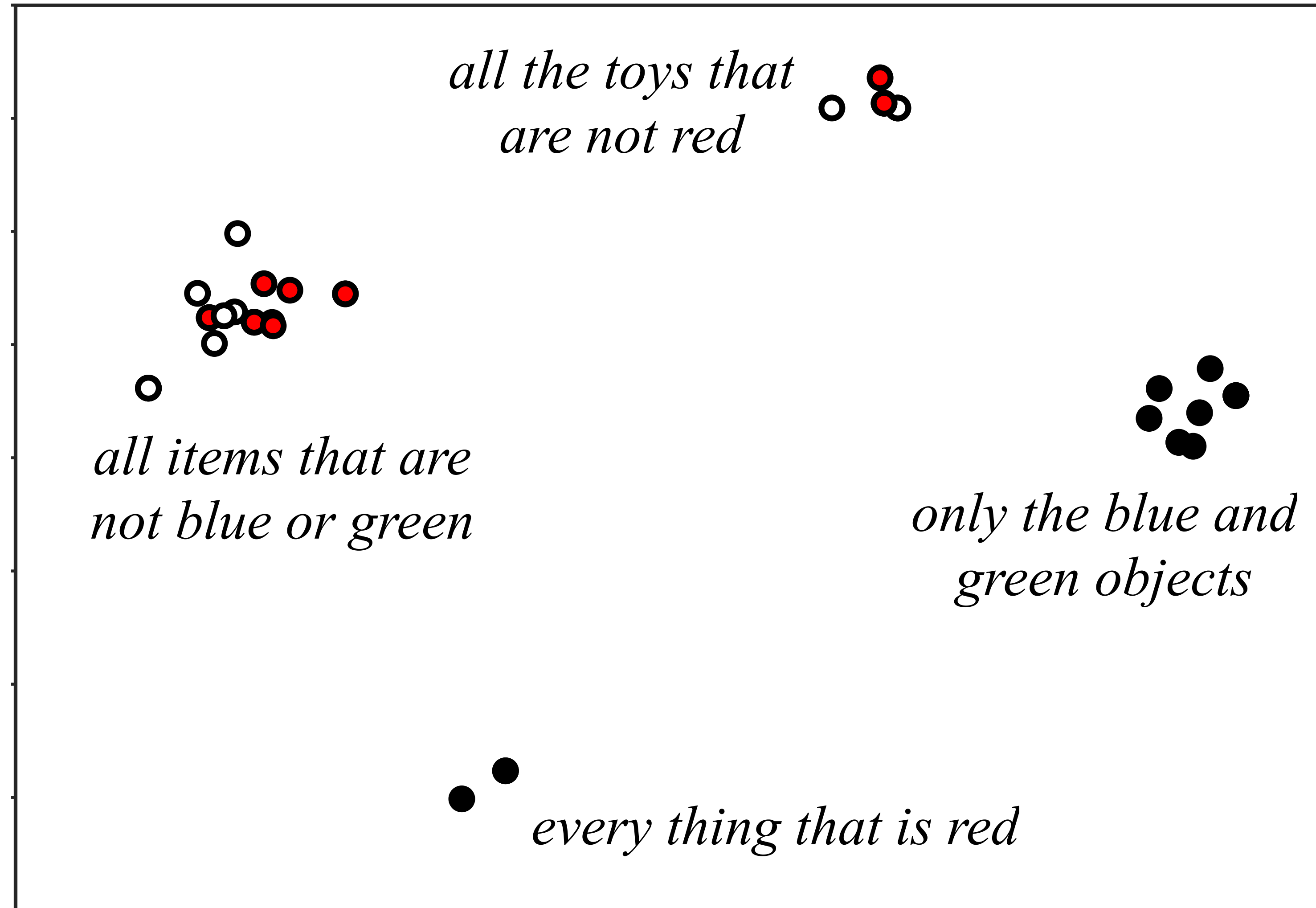


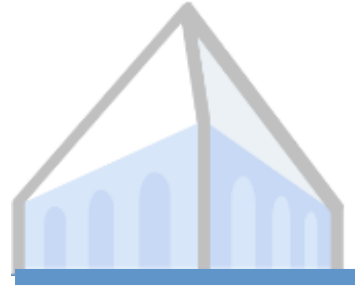
Visualizing negation

●
Input

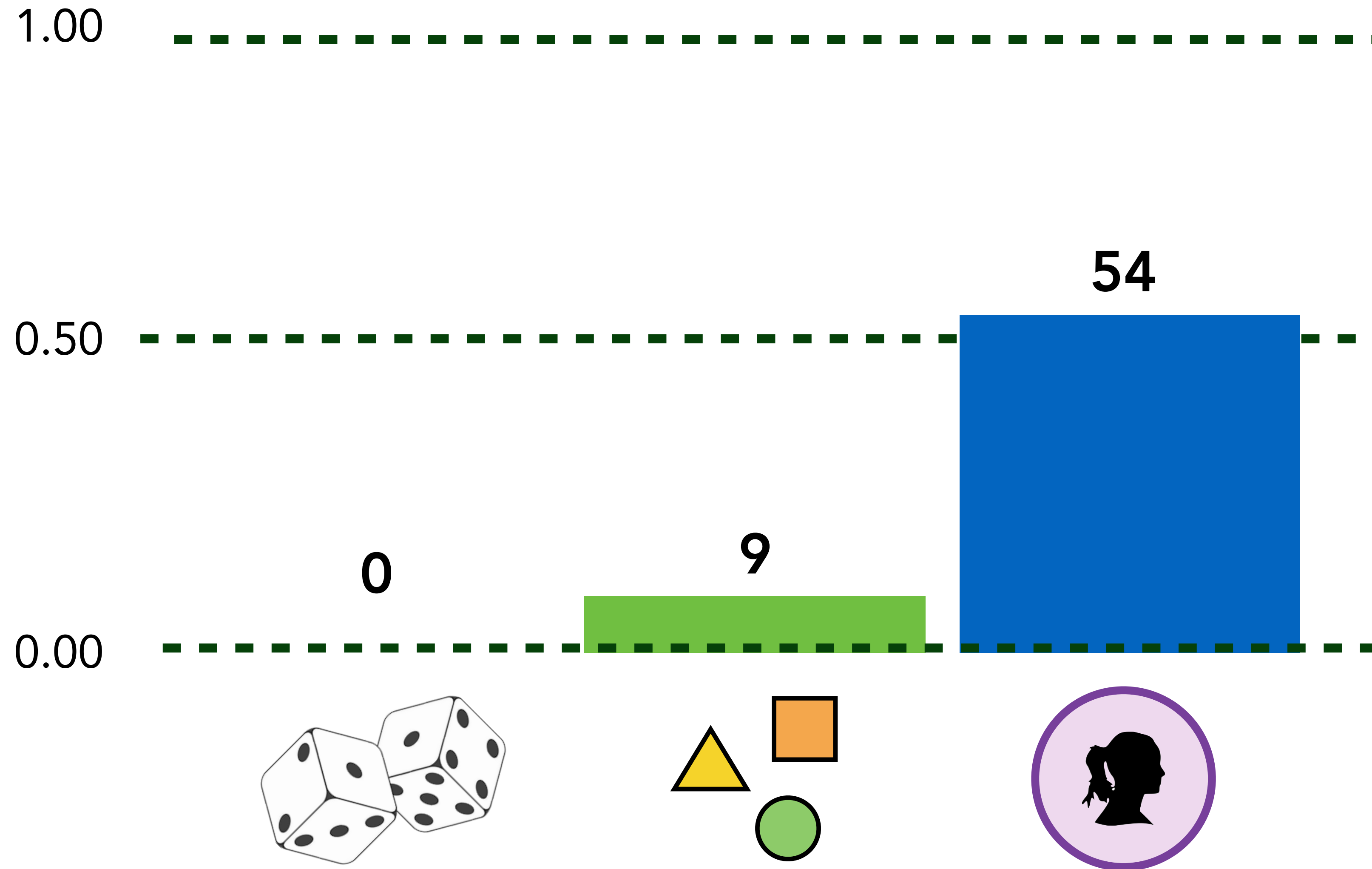
●
Predicted

○
True





Evaluation: scene agreement for disjunction



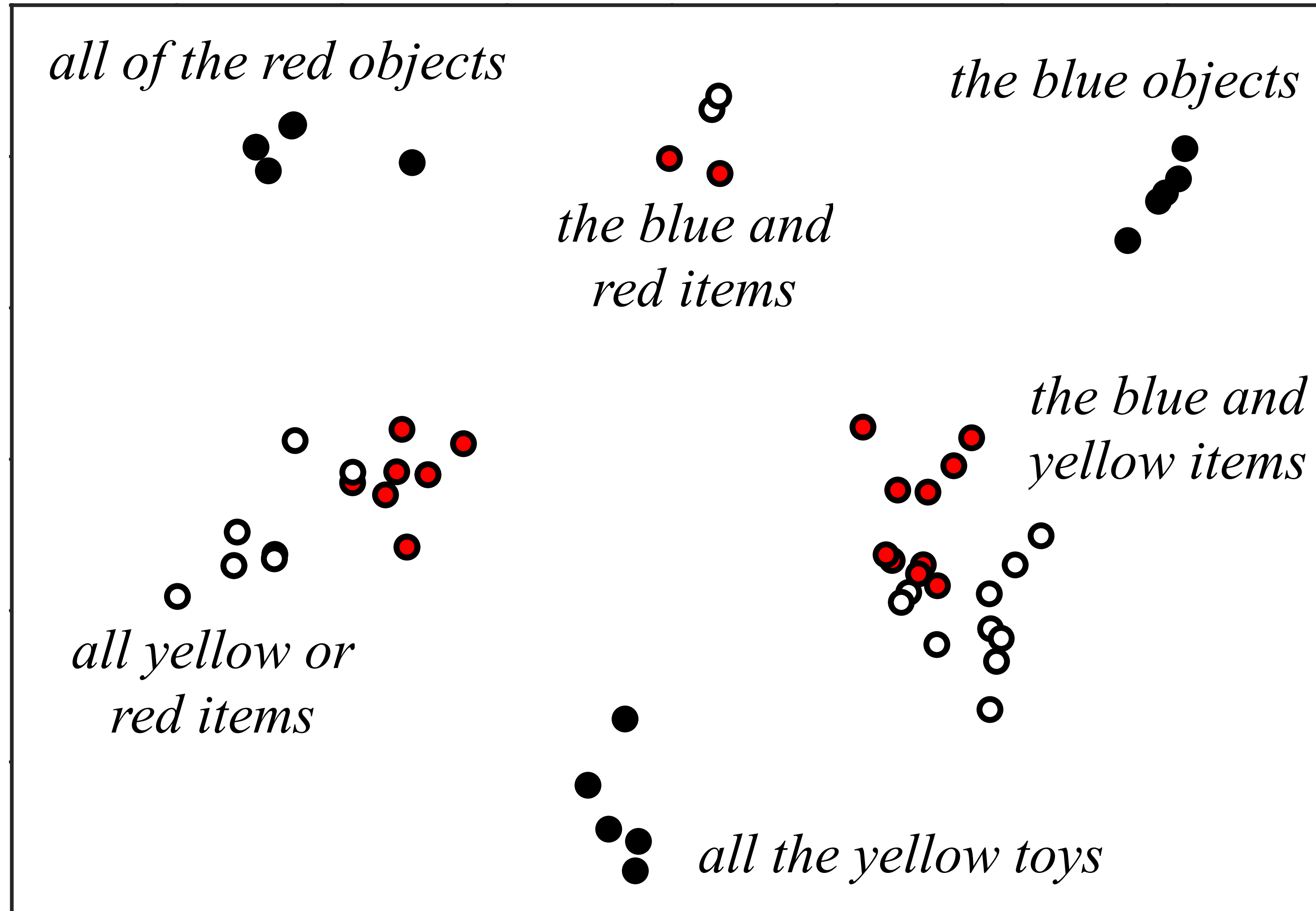


Visualizing disjunction

●
Input

●
Predicted

○
True





Conclusions

- We can translate between neuralese and natural lang. by grounding in distributions over world states
- Under the right conditions, neuralese exhibits interpretable pragmatics & compositional structure
- Not just communication games—language might be a good general-purpose tool for interpreting deep reprs.



Conclusions

- We can translate between neuralesse and natural lang. by grounding in distributions over world states
- Under the right conditions, neuralesse exhibits interpretable pragmatics & compositional structure
- Not just communication games—language might be a good general-purpose tool for interpreting deep reprs.

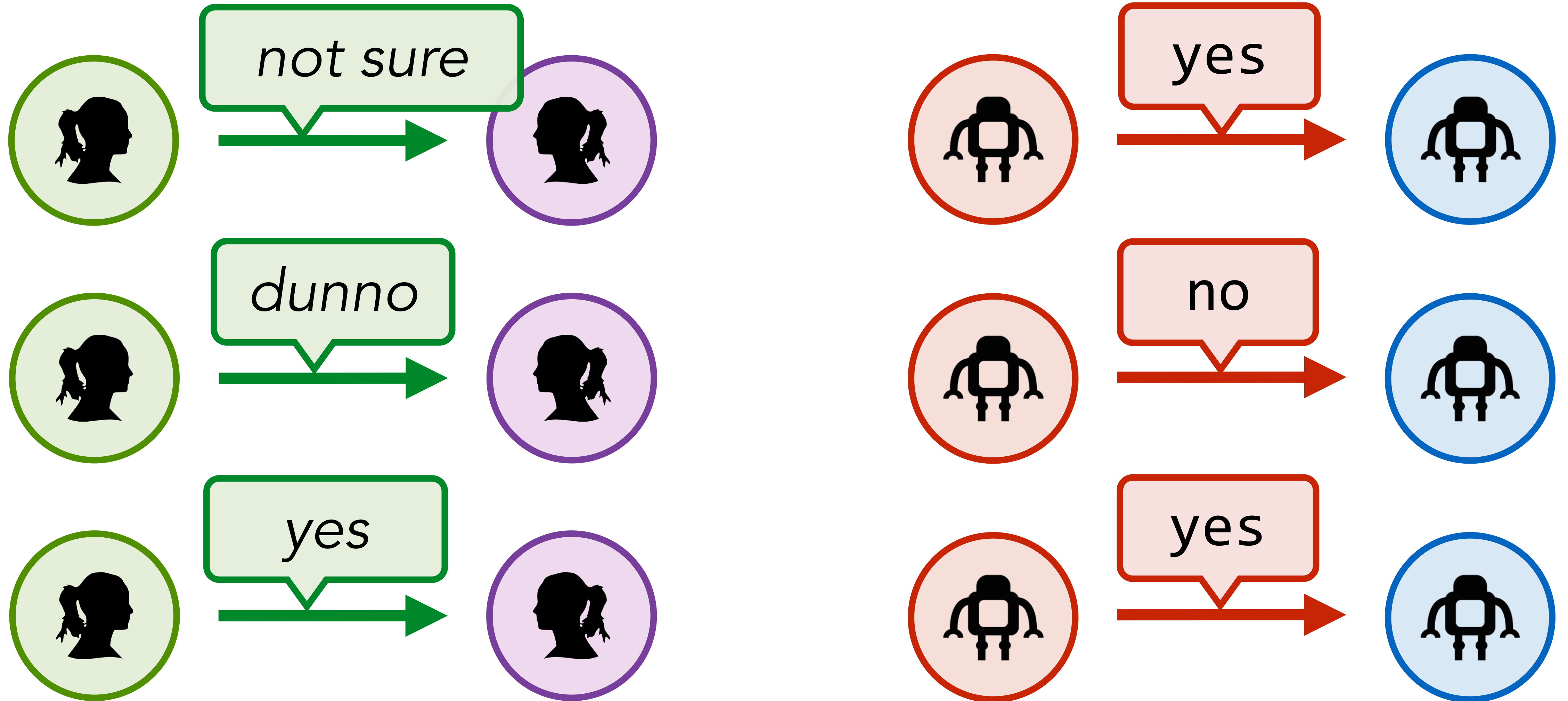


Conclusions

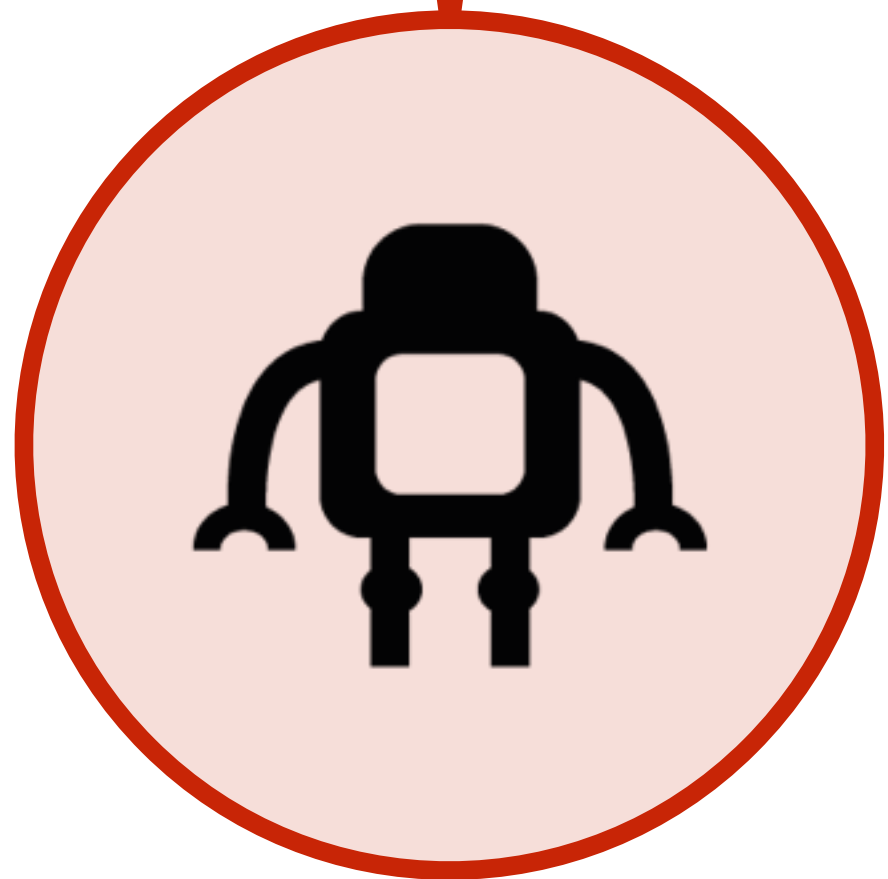
- We can translate between neuralese and natural lang. by grounding in distributions over world states
- Under the right conditions, neuralese exhibits interpretable pragmatics & compositional structure
- Not just communication games—language might be a good general-purpose tool for interpreting deep reprs.



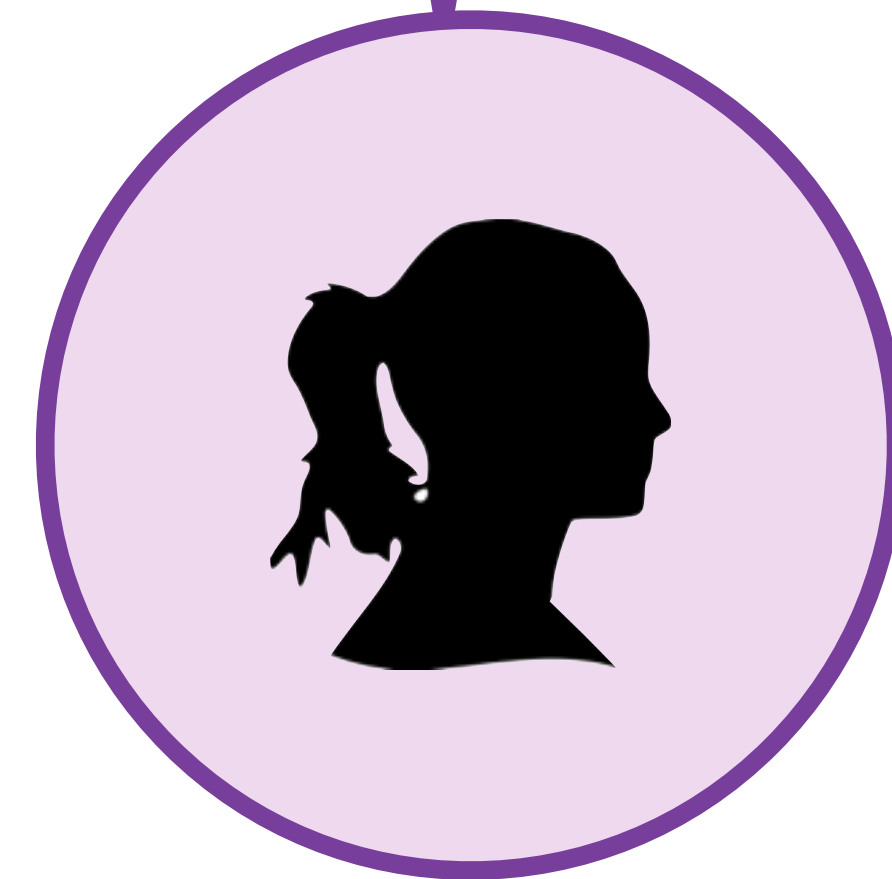
Conclusions



1.0	2.3
-0.3	0.4
-1.2	1.1



Thank you!



<http://github.com/jacobandreas/{neuralese,rnn-syn}>