

Question asking as program generation

computation and cognition lab // new york university

Anselm Rothe¹, Brenden Lake^{1,2}, & Todd Gureckis¹

¹Department of Psychology, ²Center for Data Science, New York University

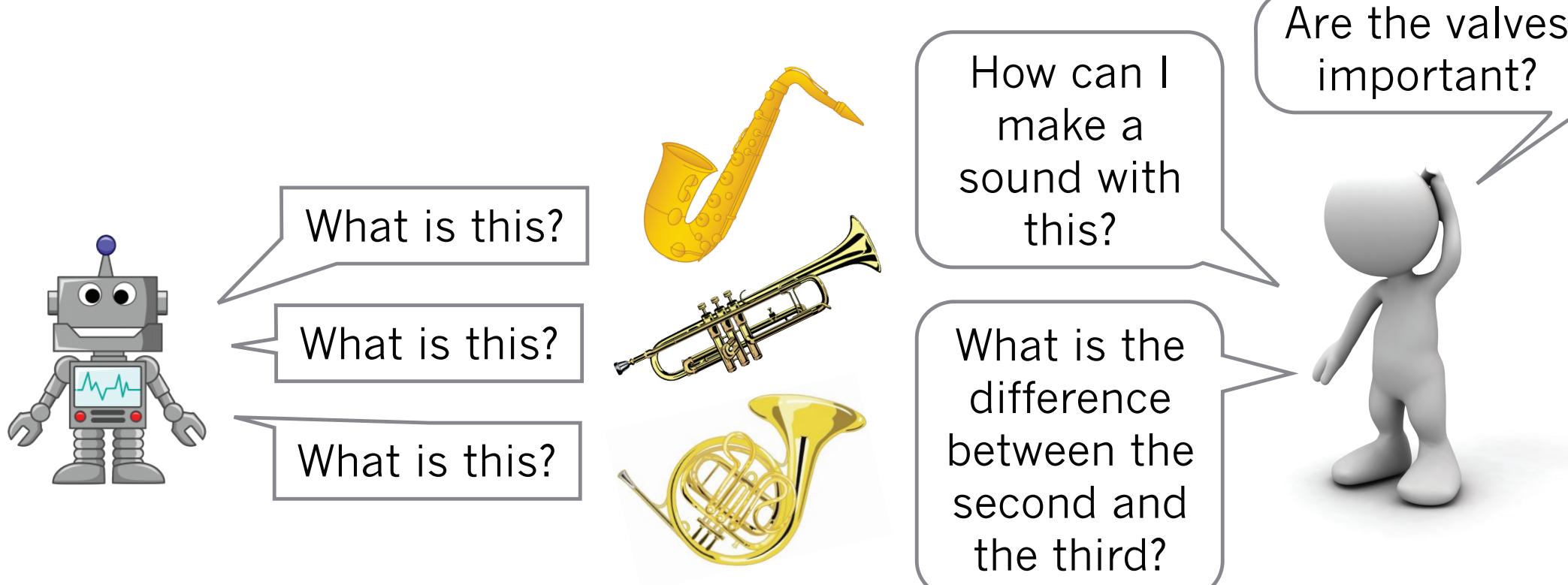


IF YOU COULD ASK ANYTHING, WHAT WOULD YOU ASK?

A question is an expression used to make a **request for information**.

Nearly all progress in Active Learning has been made with focus on a simple type of questions (label queries).

However, people use a much richer set of questions to obtain information in everyday life.

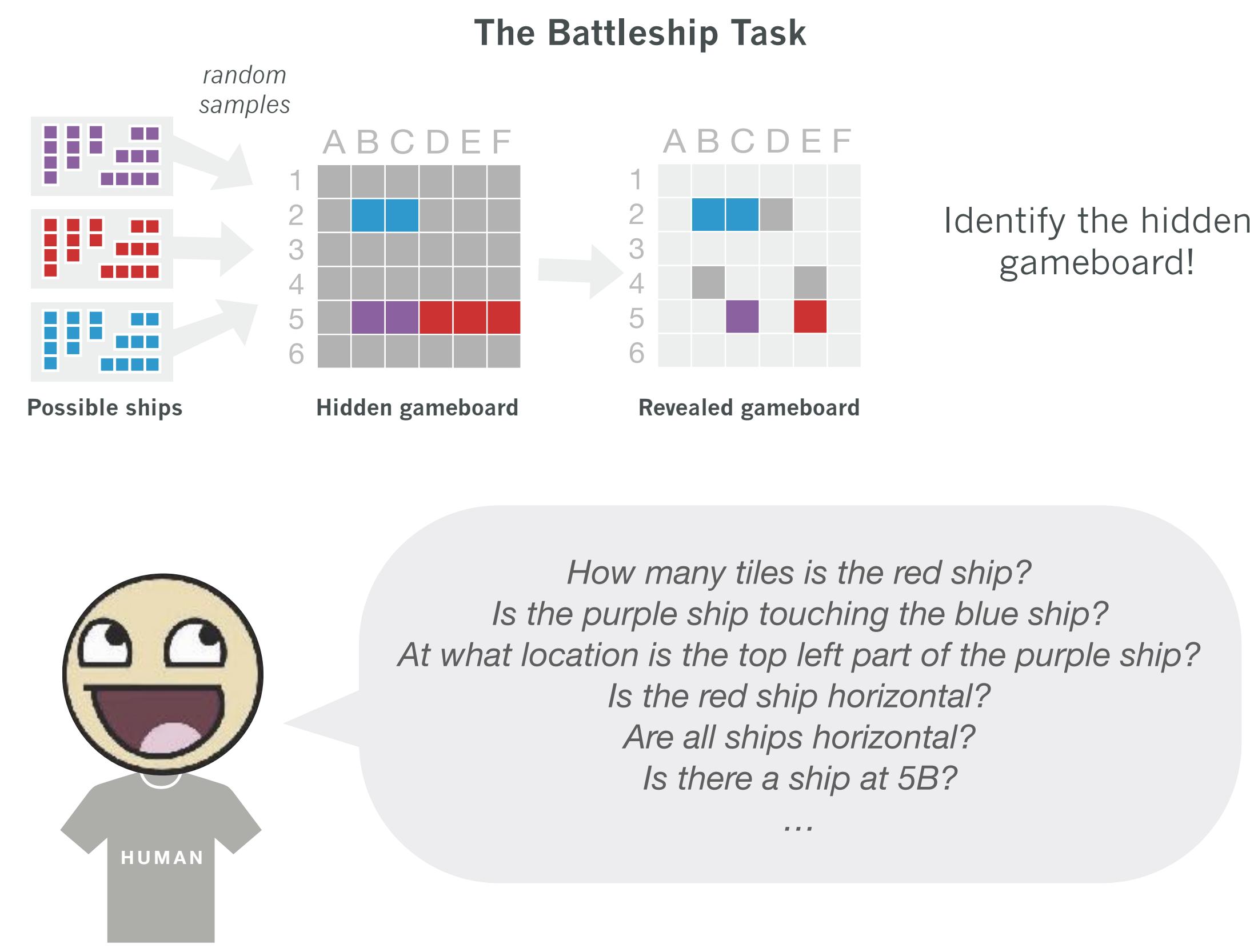


We propose a **computational framework** of how people **generate** rich questions, treating question asking as program synthesis.

A long term goal is to develop algorithms with a human-like capacity to learn by asking rich questions. A second goal is to understand more about the computational aspects of human question asking.

QUESTION DATA SET

We used the human question data set from Rothe, Lake, and Gureckis (2016) with 605 questions across 18 game contexts.



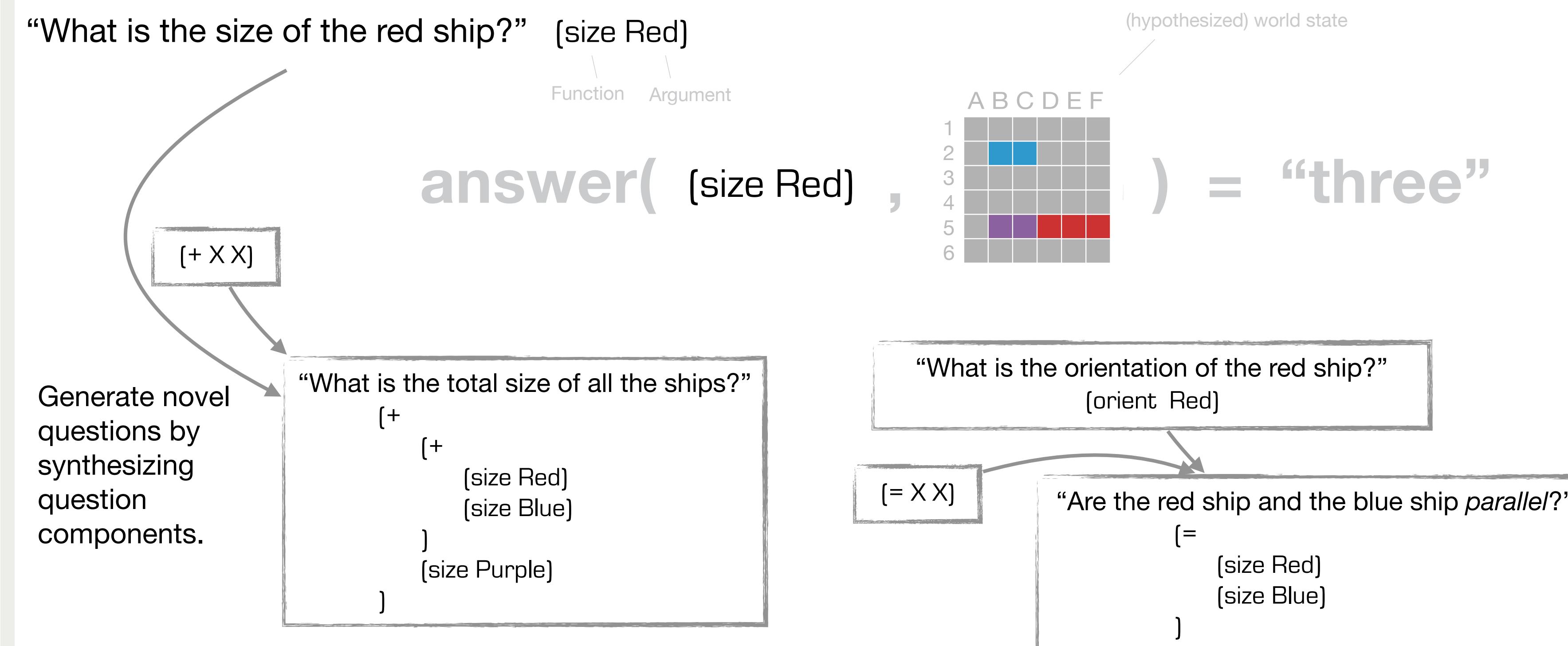
Only questions with a one-word answer were allowed and no combination of questions.

This research was supported by NSF grant BCS-1255538, the John Templeton Foundation Varieties of Understanding project, a John S. McDonnell Foundation Scholar Award to TG, and the Moore-Sloan Data Science Environment at NYU.



QUESTIONS AS PROGRAMS

We view questions as programs that, when executed on the state of the world output an answer.



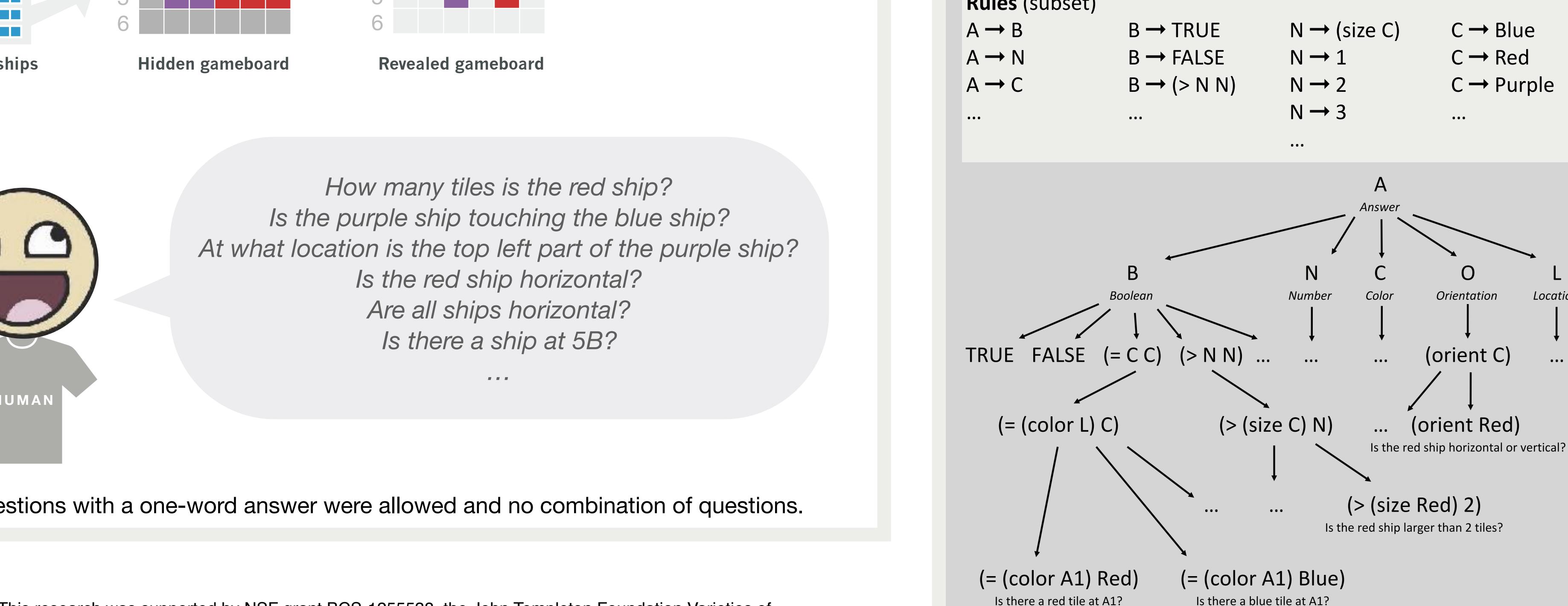
PROBABILISTIC GENERATIVE MODEL

The model should be capable of asking **novel questions** in **new contexts**.

The model aims to predict **which questions people will ask**.

We used **energy-based learning** to fit the relative importance of question features (using the human question data set as training data).

← The space of questions X is defined by a grammar.



RESULTS

QUANTITATIVE RESULTS

A We fit 4 **model variants** to all but one context and let it predict the remaining one. The log-likelihoods of the human questions were averaged across held out contexts.

Models that had one key feature lesioned achieved a lower log-likelihood than the full model using all features.

B The predictions of the **full model** showed strong alignment with the question frequencies in the data set for some contexts and more modest alignment for others (average correlation $\rho = .64$).

QUALITATIVE RESULTS

The full model was also able to generate **novel, "human-like" questions** that no human had asked.

We removed duplicate questions that were equivalent to the human questions (determined by the mutual information of their answer vectors, as well as functional equivalence up to a swapping of the arguments, e.g. (size Blue) is form equivalent to (size Red)).

Context	Question (manual translation)	Program (samples from model)
1	What is the row of the bottom right red tile? What is the row of the top left red tile? How many tiles have the same color as tile 2F? Is tile 5F a water tile? What is the column of the top left of the tiles that have the color of the bottom right corner of the board?	[rowL (bottomright [coloredTiles Red])) (rowL (topleft [coloredTiles Red])) (setSize [coloredTiles (color 2F)]) (isSubset [coloredTiles Water] [coloredTiles (color 5F)]) (colL (topleft [coloredTiles (color (bottomright [set 1A ... 6F]))]))
2	What is the column of the bottom right water tile? What is the column of the top left water tile? How many tiles have the same color as the bottom right tile of the board? What is the bottom right tile that has the same color as the tile 1A? Are all the ships oriented horizontally?	(colL (bottomright [coloredTiles Water])) (colL (topleft [coloredTiles Water])) (setSize [coloredTiles (color (bottomright [set 1A ... 6F]))]) (bottomright [coloredTiles (color (topleft [set 1A ... 6F]))]) (all [map [lambda x == H [orient x]] (set Blue Red Purple)])
3		
4		
5		
6		
7		
8		
9		
10		
11		
12		
13		
14		
15		
16		
17		
18		
19		
20		

TAKE HOME POINTS

Our model can produce **interesting and informative questions** that people could have but did not ask in the question data set.

Our model predicts to a certain degree what questions people will ask in a given game context.

The **compositionality** of our approach is important as about 15% of the human questions did only appear in a single game context. Any model unable to synthesize novel questions would be guaranteed to fail at these 15%.