

# Humans choose visual subgoals to reduce cognitive cost

**Felix Binder**

Dept. of Cognitive Science  
UC San Diego  
fbinder@ucsd.edu

**Marcelo Mattar**

Dept. of Psychology  
New York University  
marcelo.mattar@nyu.edu

**David Kirsh**

Dept. of Cognitive Science  
UC San Diego  
kirsh@ucsd.edu

**Judith Fan**

Dept. of Psychology  
Stanford University & UC San Diego  
jefan@stanford.edu

## Abstract

Physical assembly is a difficult planning problem. Humans do it efficiently by breaking large problems into smaller, easier to solve portions. But what governs which portions are chosen? We present a computational model that predicts that humans break physical assembly problems down to minimize cognitive costs. We test this by asking participants to choose which part of a tower they want to build next. Participants reliably choose the easier to solve subgoal out of two otherwise similar options. Beyond the immediate cognitive cost, participants also consider how difficult the rest of the tower will be to solve. A model that takes into account near-future cognitive costs best predicts participants' choices. These findings show that humans can estimate how difficult solving a subgoal will be, and that they choose subgoals to minimize immediate and future cognitive costs. These results help explain how humans make efficient use of cognitive resources to solve complex planning problems.

**Keywords:** planning; problem-solving; physical reasoning; task decomposition; subgoals

## Introduction

Imagine trying to assemble a shelf without the instruction booklet. Hundreds of individual steps might be involved, too many to plan them all at once. Not only are there many different actions to choose from at every turn, combining objects also leads to an explosion of possible world states. These two factors—many possible actions, and their compounding effects—make physical assembly a computationally difficult planning problem. Yet, humans are able to put the shelf together correctly. Given humans' ability to predict the outcome of physical interactions (Battaglia, Hamrick, & Tenenbaum, 2013), planning can be understood as a search over possible sequences of actions and their consequences (Newell & Simon, 1972). How do humans make this computationally explosive problem tractable?

People both think about individual actions (e.g., *tighten this screw*) and break the problem down into smaller parts (e.g., *put together the rear support first*). In many cases, breaking a problem down into subgoals and then solving each subgoal is easier than considering the entire problem at once (Solway et al., 2014; Correa, Ho, Callaway, Daw, & Griffiths, 2022; Newell, Shaw, & Simon, 1958; Maisto, Donnarumma, & Pezzulo, 2015). In the case of physical construction, knowing which subgoals to even consider is not trivial: in task like navigation, potential subgoals are often obvious (namely being in a certain location), but the combinatorial state of the environment in physical problem-solving makes it challenging to know what subgoals to consider. One potential

solution are *visual subgoals*: choosing a subgoal by focussing on a region of the workspace. Rather than imagining the backbone of the shelf as already completed (which would imply knowing how it needs to be completed), one could decide to focus on the area where the back of the shelf is to be constructed, temporarily ignoring the rest of the workspace. This changes the input to the planning system itself.

The notion of changing the input to the planning system to make planning easier is captured by task construals (Ho et al., 2022; Bapst, Sanchez-Gonzalez, Shams, et al., 2019). A task construal is a representation of a problem that has been modified to make solving the problem easier. For example, in the case of constructing the shelf, a task construal might be a representation of only the relevant part of the workspace, tools and parts that are necessary to solve the problem. This reduces the number of actions and resulting states of the world that need to be considered. However, task construals are an optimal theory, meaning they give an account of how the use of cognitive resources is minimized in an ideal way, but it leaves open how people actually propose and select task construals—there are myriad ways to change the construal of a problem. Visual subgoals can be thought of as a particular, constrained form of task construal. Focusing only on a contiguous area of the problem provides a natural way of simplifying planning. The constraints of visual subgoals make this a more tractable approach than task construals generally: proposing and evaluating a task decomposition into areas of the problem is more tractable than proposing and evaluating a general modification to the representation of the problem.

Supposing that people decompose problems into subgoals, what factors govern people's decisions about which subgoal to work on at a given point in time? In particular, to what degree do people represent immediate and future cognitive costs when deciding what part of the problem to work on now? Here, we investigate what governs people's choices between different decompositions of a challenging physical assembly problem. In a set of interlocking experiments, we first ask if the choice of subgoal actually matters for how difficult it is to plan. We find that there are non-trivial differences in planning cost beyond the mere size of a subgoal. Then, we ask if people can tell which subgoal is a good choice when choosing between them. Indeed, participants were sensitive to differences in planning cost even before having planned the subgoal. Finally, do people

consider not only the cost of the next subgoal, but also the cost of the rest of the problem? Overall, we found that participants took some, but not all, future costs into account.

### Formal Approach

In order to capture how people choose subgoals, we propose a resource-rational account of planning using subgoals (Gershman, Horvitz, & Tenenbaum, 2015; Callaway et al., 2018). This requires both a way of estimating the planning cost of finding a set of actions for a subgoal, and a strategy of choosing subgoals to reduce the expected planning cost of the overall task.

**Estimating planning costs using search** Classically, planning is thought of as search over a graph of potential actions and their effects (Newell & Simon, 1972; Geffner, 2013). To estimate how difficult a subgoal is to solve, we use the number of states explored by the Best First Search planning algorithm, which has been suggested to best explain human planning (van Opheusden, Galbiati, Bnaya, Li, & Ma, 2017). We operationalize the corresponding planning cost in humans as the time spent planning.

**Choosing subgoals to minimize planning costs** We propose a framework for optimal task decomposition that minimizes the planning cost. Our framework is based on Correa, Ho, Callaway, and Griffiths (2020). Their account predicts that problems are best decomposed into subgoals that minimize the total planning cost of the entire task. The task is fully decomposed into subgoals before any actions are taken. However, in reality, people often decompose tasks incrementally, especially for challenging tasks where it can be difficult to identify all subgoals in advance and to remember a long sequence of subgoals. To address this, our framework allows for determining the next subgoal at a time.

Choosing the next subgoal requires making a trade-off between minimizing the planning cost of the subgoal and maximizing the progress made towards the full goal: do I want to bite off a small piece of the problem that I can solve quickly, or do I want to take on a larger piece of the problem that will take longer to solve but will get me closer to the eventual goal? The value of a particular subgoal is formalized as  $r_i - c_i * \lambda$ , where  $r_i$  is the  $i$ th subgoal as a percentage of the entire goal,  $c_i$  is the planning cost, and  $\lambda$  controls the trade-off between preferring subgoals that are easy to solve over those that make substantial progress.

When choosing a subgoal, the planning cost can be taken into account in three general ways (see Figure 1): A **myopic** strategy for subgoal selection only aims to maximize the value of the current subgoal. While this is the simplest strategy, it can lead to subgoal choices that lead to dead ends. The **lookahead** subgoal selection strategy considers some future costs, but not necessarily all. When planning a sequence of subgoals, this strategy chooses the sequence of length  $d$  (including shorter sequences that complete the

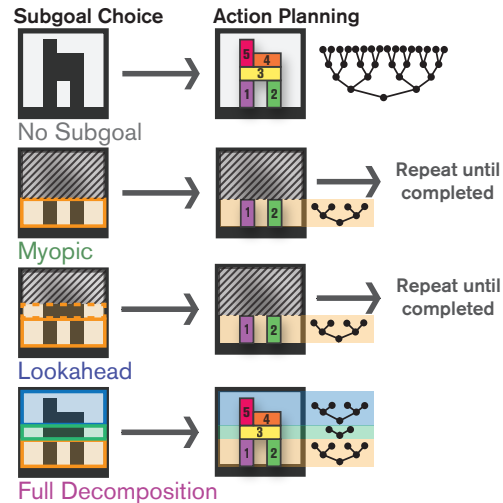


Figure 1: Three different strategies for subgoal selection during physical assembly: **myopic** (only considering the next subgoal), **lookahead** (also considering the next (few) subgoals), and **full decomposition** (breaking down the entire problem in advance). An example action graph is shown on the right. Comparing not using subgoals against full decomposition illustrates how the use of subgoals can reduce the total planning cost.

problem) that maximizes the sum of the utility of the subgoals in the sequence, where the utility of a subgoal is  $c_i * \lambda + r_i$  and then passes the first subgoal in the sequence to the action planner. While this strategy requires more effort when selecting subgoals, it can reduce the risk of making subgoal choices that eventually leads to bad outcomes. Finally, the **full decomposition** strategy considers the entire future planning cost (equivalent to Correa et al. (2020)). The **full decomposition** strategy differs from the **lookahead** strategy with  $d = \infty$  insofar as it decomposes the entire problem once and then sticks to the decomposition, whereas the **lookahead** strategy chooses each subgoal individually, even if it considers future states. The full decomposition is guaranteed to minimize the planning cost of finding actions *given the subgoals*, but it requires considering a potentially large number of potential subgoal combinations, especially for larger problems.

Note that this is an account of optimal task decomposition, but not necessarily a model of how people decompose tasks. One of the main differences is that this model assumes that the planning cost of a subgoal is known before it is chosen. While proxies for planning cost exist (such as the size of the subgoal), the actual cost of the subgoal is not known until it is solved. This model only aims to minimize the planning cost of finding solutions to the chosen subgoals, but does not itself minimize the cost of choosing subgoals in the first place. For an investigation of the subgoal selection cost of various approaches, see Binder, Mattar, Kirsh, and Fan (2021).

**Hypotheses** For this optimal model to apply to human behavior, the choice of good subgoals need to matter and

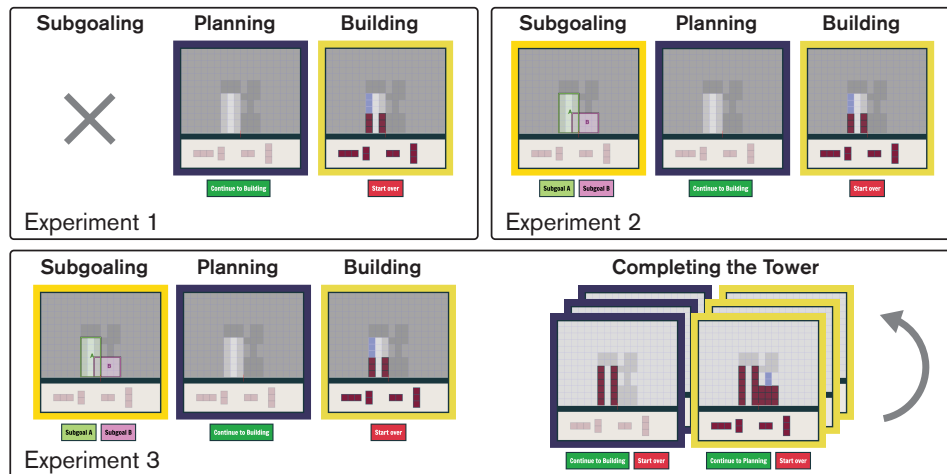


Figure 2: In each experiment, participants plan and construct subgoals on the *Block Tower Reconstruction Task*. During the subgoal selection phase, participants choose between two subgoals (except experiment 1). In the planning phase, they come up with a plan. The plan is executed under time pressure in the building phase. Experiment 3 also requires participants to complete the rest of the tower, where they can switch between planning and building as often as they like.

people need to be able to tell what a good subgoal is. We hypothesize that: (1) different subgoals are easier or harder to plan than others for people—therefore, choosing good subgoals matters for how hard a problem is to solve. (2) People are sensitive to the planning cost of subgoals when choosing between potential subgoals, and (3) when choosing subgoals, people take not only the planning cost of the subgoal itself into account, but also the planning cost of future subgoals. To investigate these hypotheses in turn, we conducted three interlocking experiments.

### Experiment 1: Does the choice of subgoal matter for cognitive cost?

In this experiment, we investigate whether subgoal choice matters for cognitive cost. Intuitively, the larger a subgoal is, the more effort it will take to find a solution. But are there subgoals of the same size that are harder or easier to plan? In other words, should one consider more than just the size of a subgoal when choosing it? To investigate the effect of choosing one subgoal over a similar one, we presented participants with predetermined subgoals and asked them to find solutions to it.

### Methods

**Block Tower Reconstruction Task** We use the *Block Tower Reconstruction Task* adapted from McCarthy, Kirsh, and Fan (2020), in which participants have to construct a physically stable 2D tower of a given shape—see Figure 2. Similar block tower construction tasks have been used to study planning and physical reasoning in artificial agents (Sussman, 1975; Bapst, Sanchez-Gonzalez, Doersch, et al., 2019) and humans (Dietz, Landay, & Gweon, 2019; Cortesa, Jones, Hager, & Khudanpur, 2018).

Visual subgoals on the task are defined as a rectangular region of the building area and are shown as an overlay. When

building a subgoal, participants were required to perfectly complete the subgoal with no blocks sticking out.

In order to isolate the planning time from the time needed to execute the plan, participants are required to first come up with a plan for the subgoal, then to click a button to advance to the building phase. During the building phase, participants were subject to time pressure in order to prevent them from doing online planning.

**Participants** 86 participants (51 male,  $M_{Age} = 38.5$ ) years were recruited from Prolific and paid a minimum of \$14 per hour. We excluded 6 participants who switched away from the study webpage too often.

**Stimuli** To find subgoals, we procedurally generated 128 stable block tower shapes of varying sizes. For each tower, we generated 4 pairs of subgoals. The subgoals in a pair were matched on their area. The simulated planning cost of solving each subgoal was determined using the Best First Search model. This makes the subgoal with the lower predicted cost the *best* of the pair, the other one the *worst*.

**Procedure** Participants were asked to plan and build 24 predetermined subgoals, corresponding to 12 pairs randomly drawn from the larger set. The overall order of the subgoals was randomized. After completing the subgoal, participants were moved to a different subgoal on a different tower.

### Results

**Overall performance** Overall, participants were able to solve these subgoals in a reasonable amount of time, usually on their first attempt. The average planning time on the first attempt<sup>1</sup> to solve the subgoal was 8.95 seconds

<sup>1</sup>On the block tower construction task, if participants fail to construct the subgoals (due to instability, time pressure or planning errors), they have to start over. Here, we report the planning time on the first attempt at building the subgoal, as this measure is not

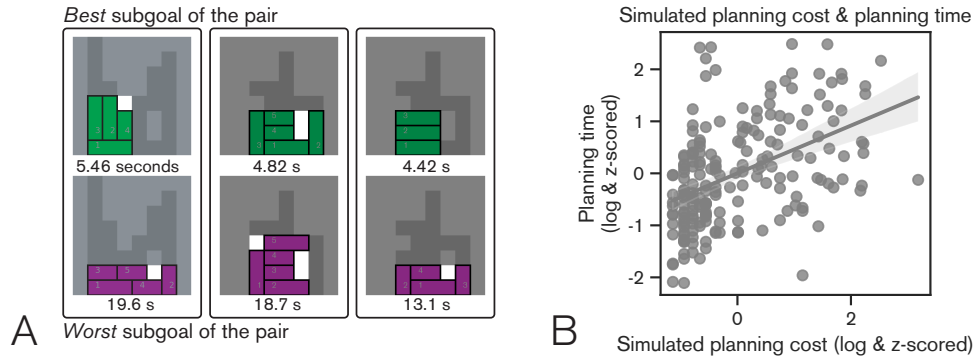


Figure 3: (A) The area-matched subgoals pairs with the highest difference in average planning time (displayed below each subgoal) are shown. The solution found by Best First Search is displayed. (B) The number of costs explored by Best First Search for a given subgoal compared against participants' planning time.

(95% confidence interval (CI): [8.55, 9.34]). Participants required an average of 1.32 attempts to solve the subgoal (95% CI: [1.28, 1.37]), meaning that most of the time participants solved the subgoal on the first try. Moreover, planning time was systematically related to the size of the subgoal (Spearman rank correlation  $\rho(188) = 0.385$ , 95% CI: [0.271, 0.498],  $p < 0.001$ ).

**Difference in planning time between best and worst subgoals** The simulated planning cost was used to generate matched pairs. For the *best* subgoal of the pairs, the participants found a solution in 7.93s (95% CI: [7.41, 8.45]) on average. For the *worst* subgoal, the participants found a solution in 9.96s (95% CI: [9.42, 10.6]) on average. The comparison of the mixed effect model with the fixed effect of subgoal type (best vs. worst) and random effects for participant and subgoal pair against the null model without subgoal type was significant ( $\chi^2(1) = 34.4, p < 0.001$ ). Participants also made fewer errors on the *best* subgoal than on the *worst* subgoal (mixed effect model comparison against null model:  $\chi^2(1) = 21.8, p < 0.001$ ). This difference in planning time between the *best* and *worst* subgoals is consistent with the predictions of the Best First Search model.

**Graded relationship between planning time and simulated search cost** Beyond the manipulation of *best* and *worst* subgoals, we also tested whether the Best First Search model predicts participants' planning time. The Spearman rank correlation between the predicted costs in number of states explored by Best First Search and the actual planning time incurred by human participants is  $\rho(188) = 0.491$  (95% CI: [0.380, 0.598],  $p < 0.001$ ). While the model leaves some variance unexplained, it is predictive of the time humans incur during action planning.

Taken together, these results suggest that the choice of good subgoals matters for the planning time even beyond the size of the subgoal—therefore, choosing a good subgoal

requires being able to estimate how difficult it will be to find a solution to the subgoal.

## Experiment 2: Do humans choose subgoals that reduce cognitive cost?

The previous experiment showed that there are differences in the planning cost of individual subgoals. Are humans sensitive to these differences when choosing subgoals?

### Methods

**Stimuli** In this study, we draw on the 96 pairs of subgoals from the previous study. From those, we selected 24 pairs of subgoals which had the highest ratio between the *best* and *worst* planning time as measured in the previous experiment.

**Participants** 80 participants (57 male,  $M_{Age} = 34.0$ , 1 excluded) were recruited through Prolific and paid a minimum of \$14 per hour.

**Procedure** Participants were presented with a pair of matched subgoals and were asked to choose one to plan and build (“choose the subgoal that seems easier to complete quickly.”) After completing the subgoal, participants were moved to the next choice.

### Results

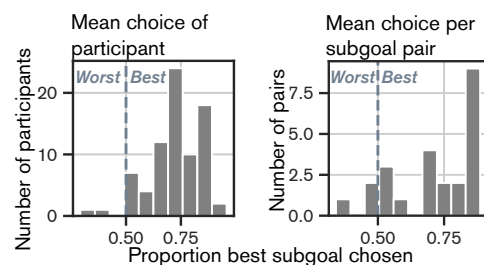


Figure 4: The distribution of participants' choices between the *best* and *worst* subgoal. Left shows the distribution of the average choice that participants made, right shows the distribution of the average choice within each subgoal pair. Bars to the right of the chance line indicate that the *best* subgoals is preferred by that participant/for that subgoal pair.

polluted by knowledge gained on previous attempts and tracks best how difficult participants thought the subgoal is ahead of time.

**Proportion of best subgoal chosen** If the choice of subgoals is driven by the planning cost of the subgoals and participants are sensitive to the expected planning cost, they should choose the *best* subgoal of the pair presented to them. Indeed, in 72.6% (95% CI: [70.6%, 74.7%]) of choices, participants chose the *best* subgoal. This is reliably different from chance ( $t(1895) = 22.1, p < 0.001$ ). This shows that participants are capable of identifying which of the two subgoals is easier to solve before having planned it out. The effect persists when taking into account that for some subgoal pairs, the *best* subgoal can be solved with one fewer block placement than the *worst*. In those pairs, participants chose the *best* subgoal 77.6% of the time (95% CI: [75.1, 80.1]), compared to 65.6% (95% CI: [62.5, 68.7]) in pairs with the same number of placements.

**Outcome of choosing the best subgoal** Choosing the *best* subgoal also leads to better outcomes: participants needed to start over 0.0856 times per subgoal (95% CI: [0.0622, 0.115]) of the *best* subgoals, compared to 0.459 times (95% CI: [0.367, 0.563]) on the *worst* subgoal. A mixed effect linear model with restarts as the dependent variable, subgoal choice as the independent variable, and random intercepts for each participant and tower was compared to a simpler model without subgoal choice. A significant difference was found ( $\chi^2(1) = 79.4, p < 0.001$ ). After choosing the *best* subgoal, participants were also faster to find a solution to it: the average planning time on the first attempt to solve the *best* subgoal was 3.87s (95% CI: [3.68, 4.07]), while the average planning time on the first attempt to solve the *worst* subgoal was 6.72s (95% CI: [6.18, 7.29]); mixed effect model comparison  $\chi^2(1) = 104, p < 0.001$ . However, this could also be explained by underlying skill driving both “better” choices and faster planning.

**Relation of subgoal selection time to choice** Does thinking longer when choosing the subgoal lead to better choices? The longer a participant tends to think when choosing a subgoal, the more likely they were to choose the *best* subgoal ( $r(77) = 0.344, p < 0.001$ ). This indicates that spending more time to break the problem down results in having an easier time to plan down the line.

**Relation of subgoal selection time to planning time** The average time to choose a subgoal is 6.67s (95% CI: [6.32, 7.07]), which is less than the time taken to plan those subgoals themselves in the previous study (10.3s, 95% CI: [10.3, 10.3]). This rules out that participants might plan out both subgoals before choosing one of them and therefore have access to how difficult it is to solve.

### Experiment 3: Do humans choose subgoals that reduce future cognitive cost?

Experiment 2 shows that participants take the cost of solving a subgoal into account when choosing between subgoals. So

far, the cost of the subgoal itself was the only relevant factor. In real life, however, choosing one subgoal will affect how the rest of the problem can be broken down. Do participants take the cost of the rest of the problem into account when choosing the next subgoal?

## Methods

**Participants** 200 participants (133 male,  $M_{\text{Age}} = 36.4$ ) were recruited from Prolific and paid about \$14/hour. 21 were excluded and 8 were unable to finish the study due to technical issues.

**Stimuli** The aim is to isolate the effect that future planning costs have on the subgoal choice. Therefore, we want to select pairs of subgoals that are matched on both planning cost and progress, but vary with respect to predicted future costs under different subgoal selection strategies.

To estimate how attractive each initial subgoal is under different strategies, we generated predictions from the **myopic**, **lookahead** and **full decomposition** strategies marginalized over the dynamic range of  $\lambda$ . We chose 54 subgoal pairs where the preferences of the different subgoal selection strategies differed maximally. To be sure that we can distinguish random behavior from the **myopic** strategy, we included 18 pairs for which the initial subgoals were allowed to vary in size and progress.

**Procedure** Participants were asked to choose between two initial subgoals, then to plan and build the chosen subgoal for 12 pairs of subgoals, randomly drawn from the larger set of 72. Since this study investigates sensitivity to future costs, participants were required to complete building the rest of the tower. After the initial subgoal, subgoals were not provided. Rather, participants moved between the planning and time-pressured construction phase as often as they wanted. This allowed for the collection of data on the time spent planning the rest of the tower.

## Results

**Simple cost of the rest of the tower & subgoal choice** Are participants sensitive to the cost of the rest of the tower? The simplest way of measuring this is to ask whether the ratio in the predicted cost of the rest of the tower after completing the subgoal predicts the choice of the subgoal. We found that the mere cost of planning the rest of the tower without the use of subgoals predicts the subgoal choice poorly ( $r(70) = 0.0938$ , 95% CI bootstrapped by resampling participants:  $[-0.0379, 0.199], p = 0.469$ ). We found that the model based on the mere cost of planning the rest of the tower is a poor explanation

**Subgoal selection strategies & subgoal choice** The basic planning cost of the rest of the tower does not predict the subgoal choice—what about models of subgoal choice that take future decompositions into account? To answer this question, we compared the predictions of the three strategies to the choice proportion for each pair across all



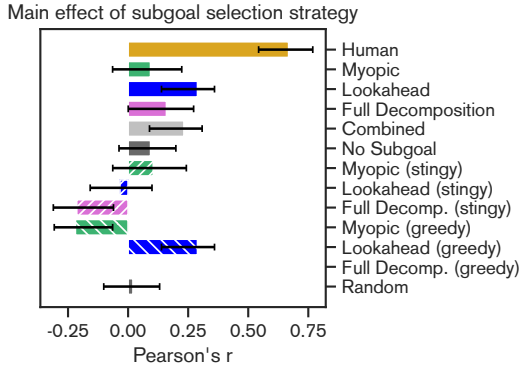


Figure 5: The correlation coefficient for each fitted model of subgoal selection. Error bars are the 95% bootstrapped confidence intervals. The greedy and stingy models are lesions that are only sensitive to progress and cost, respectively. The bar for “Human” shows the corrected split-half correlation across samples of the participants.

participants. The value for  $\lambda$  and the softmax temperature  $T$  (which translates the predicted value of two subgoals into the proportion for one) were fitted for each planner to maximize the cosine similarity of predicted and actual choice proportions across the subgoal pairs. The **lookahead** strategy predicts participants behavior, with a Pearson  $r(70) = 0.289$  (95% CI: [0.139, 0.360],  $p = 0.013815$ ) between the predicted choice proportion and the human choice proportions across the subgoal pairs. However, the **myopic** strategy and the **full decomposition** strategy do not predict participants choices well ( $r(70) = 0.0935$ , 95% CI: [-0.0645, 0.223],  $p = 0.435$  and  $r(70) = 0.160$ , 95% CI: [0.000508, 0.273],  $p = 0.180$ ). Simply averaging over the predicted choice proportion of the three strategies yields a Pearson correlation of  $r(70) = 0.233$  (95% CI: [0.0891, 0.308],  $p = 0.0493$ )—the strategies taken together predict participants behavior reasonably well.

**Full decomposition strategy & actual planning time** The **full decomposition** strategy operates on an estimation of the planning cost of the entire tower. While it does not predict participants choices, it is predictive of the actual planning time that participants incur when completing the tower? Indeed, the choice proportion of the **full decomposition** strategy predicts the ratio of the total planning time for the rest of the tower ( $r(70) = -0.272$ , 95% CI: [-0.454, -0.0821],  $p = 0.0209$ ).<sup>2</sup> The **myopic** and **lookahead** strategies do not predict the planning time well ( $r(70) = -0.0482$ , 95% CI: [-0.283, 0.180],  $p = 0.688$  and  $r(70) = 0.214$ , 95% CI: [-0.0165, 0.404],  $p = 0.0711$ ).

**Lesioned strategies** The strategies presented so far involve a trade-off between the cost of the subgoal and the cost of the rest of the tower. Lesioned versions of the strategies that do not take into account the cost of

<sup>2</sup>The negative correlation means that participants preferred the easier subgoal of the pair.

subgoals, only their progress—“greedy”—are not predictive of human behavior for the **myopic** ( $r(70) = -0.219$ , 95% CI: [-0.307, -0.0643],  $p = 0.0640$ ), but are for the **lookahead** strategy ( $r(70) = 0.289$ , 95% CI: [0.139, 0.360],  $p = 0.0138$ ).<sup>3</sup> Conversely, versions that only care about avoiding costs—“stingy” strategies—are not predictive of human behavior for the **myopic**, **lookahead**, and **full decomposition** strategies ( $r(70) = 0.105$ , 95% CI: [-0.0639, 0.242],  $p = 0.381$ ;  $r(70) = -0.0340$ , 95% CI: [-0.158, 0.0994],  $p = 0.777$ ;  $r(70) = -0.214$ , 95% CI: [-0.311, -0.0608],  $p = 0.0713$ ).

## Discussion

Here we presented a model of subgoal selection in physical assembly that incorporates both current and future planning costs. To test the degree to which this model describes patterns in human subgoal selection, we conducted an experiment in which people made decisions between different portions of the tower—i.e., visual subgoals—they would prefer to build next. We found that participants’ decisions were sensitive to differences in how much planning time was required to figure out how to complete each subgoal. While none of the strategies were strongly predictive of participants’ subgoal choices, they are best described by the **lookahead** strategy, which is consistent with participants taking some, but not all the future costs into account when choosing a subgoal. This is consistent with a resource-rational model of subgoal decomposition (choosing subgoals to minimize cognitive cost). It is also consistent with humans decomposing problems into subparts as they go along. Taken together, results contribute to our understanding of how humans solve complex physical construction problems. However, they raise the question of how humans are capable of the metacognitive feat of estimating the planning costs of potential subgoals.

Notably, the subgoal selection strategies as outlined are optimal theories that describe what the best choice would be. One assumption of the theory is that the actual planning cost of a subgoal is known ahead of time. Another is the value for a subgoal is derived from the cheapest possible sequence of subgoals that follow it. Both are required to determine the optimal subgoal choice, but neither is necessarily true of humans. Beyond that, the analyses as presented combine the behavior of all participants, but the strategies that participants follow may vary systematically across individuals as well as across problem contexts. Measuring the likelihood of human choice behavior under different strategies is a promising avenue for converging insight.

While the block tower reconstruction task is a richer task than the simple grid world tasks often used to study planning, it is still a greatly simplified version of the real world. Future work aims to extend this work to more realistic settings with richer action spaces, uncertainty and imperfect information.

<sup>3</sup>The greedy **full decomposition** strategy rates all full decompositions as equally good.

## Acknowledgments

The authors would like to thank Will McCarthy, as well as the other members of the Cognitive Tools Lab at UC San Diego for their comments and support. This work was supported by an ONR Science of Autonomy and a NSF CAREER Award #2047191 to J.E.F.

All code and materials available at:  
[https://github.com/cogtoolslab/  
tools\\_block\\_construction\\_human\\_experiments](https://github.com/cogtoolslab/tools_block_construction_human_experiments)

## References

- Bapst, V., Sanchez-Gonzalez, A., Doersch, C., Stachenfeld, K. L., Kohli, P., Battaglia, P. W., & Hamrick, J. B. (2019, May). Structured agents for physical construction. *arXiv:1904.03177 [cs]*.
- Bapst, V., Sanchez-Gonzalez, A., Shams, O., Stachenfeld, K., Battaglia, P. W., Singh, S., & Hamrick, J. B. (2019, October). Object-oriented state editing for HRL. *arXiv:1910.14361 [cs, stat]*.
- Battaglia, P. W., Hamrick, J. B., & Tenenbaum, J. B. (2013, November). Simulation as an engine of physical scene understanding. *Proceedings of the National Academy of Sciences*, 110(45), 18327–18332. doi: 10.1073/pnas.1306572110
- Binder, F., Mattar, M., Kirsh, D., & Fan, J. (2021). Visual scoping operations for physical assembly. *Proceedings of the 43th Annual Conference of the Cognitive Science Society*, 7.
- Callaway, F., Lieder, F., Das, P., Gul, S., Krueger, P. M., & Griffiths, T. L. (2018). A resource-rational analysis of human planning. , 6.
- Correa, C. G., Ho, M. K., Callaway, F., Daw, N. D., & Griffiths, T. L. (2022, November). *Humans decompose tasks by trading off utility and computational cost* (No. arXiv:2211.03890). arXiv.
- Correa, C. G., Ho, M. K., Callaway, F., & Griffiths, T. L. (2020, July). Resource-rational Task Decomposition to Minimize Planning Costs. *arXiv:2007.13862 [cs]*.
- Cortesa, C. S., Jones, J. D., Hager, G. D., & Khudanpur, S. (2018). Constraints and Development in Children's Block Construction. , 6.
- Dietz, G., Landay, J. A., & Gweon, H. (2019). Building blocks of computational thinking: Young children's developing capacities for problem decomposition. , 7.
- Geffner, H. (2013, July). Computational models of planning. *Wiley Interdisciplinary Reviews: Cognitive Science*, 4(4), 341–356. doi: 10.1002/wcs.1233
- Gershman, S. J., Horvitz, E. J., & Tenenbaum, J. B. (2015, July). Computational rationality: A converging paradigm for intelligence in brains, minds, and machines. *Science*, 349(6245), 273–278. doi: 10.1126/science.aac6076
- Ho, M. K., Abel, D., Correa, C. G., Littman, M. L., Cohen, J. D., & Griffiths, T. L. (2022, June). People construct simplified mental representations to plan. *Nature*, 606(7912), 129–136. doi: 10.1038/s41586-022-04743-9
- Maisto, D., Donnarumma, F., & Pezzulo, G. (2015, March). Divide et impera: Subgoalng reduces the complexity of probabilistic inference and problem solving. *Journal of The Royal Society Interface*, 12(104), 20141335. doi: 10.1098/rsif.2014.1335
- McCarthy, W., Kirsh, D., & Fan, J. (2020). Learning to build physical structures better over time. *Proceedings of the 42nd Annual Conference of the Cognitive Science Society*.
- Newell, A., Shaw, J. C., & Simon, H. A. (1958). Elements of a theory of human problem solving. *Psychological Review*, 65(3), 151–166. doi: 10.1037/h0048495
- Newell, A., & Simon, H. A. (1972). *Human problem solving* (Vol. 104). Prentice-Hall Englewood Cliffs, NJ.
- Solway, A., Diuk, C., Córdoba, N., Yee, D., Barto, A. G., Niv, Y., & Botvinick, M. M. (2014, August). Optimal Behavioral Hierarchy. *PLoS Computational Biology*, 10(8), e1003779. doi: 10.1371/journal.pcbi.1003779
- Sussman, G. J. (1975). *A computer model of skill acquisition*. New York; London: American Elsevier ; Elsevier.
- van Opheusden, B., Galbiati, G., Bnaya, Z., Li, Y., & Ma, W. J. (2017). A computational model for decision tree search. , 6.