# Rapid trial-and-error learning with simulation supports flexible tool use and physical reasoning

Kelsey R. Allen<sup>1,2,\*</sup>, Kevin A. Smith<sup>1,2,\*</sup>, and Joshua B. Tenenbaum<sup>1,2</sup>

<sup>1</sup>Department of Brain and Cognitive Sciences, MIT <sup>2</sup>Center for Brains, Minds and Machines (CBMM) <sup>\*</sup>Indicates equal contribution

Many animals, and an increasing number of artificial agents, display sophisticated capabilities to perceive and manipulate objects. But human beings remain distinctive in their capacity for flexible, creative tool use – using objects in new ways to act on the world, achieve a goal, or solve a problem. To study this type of general physical problem solving, we introduce the Virtual Tools game. In this game, people solve a large range of challenging physical puzzles in just a handful of attempts. We propose that the flexibility of human physical problem solving rests on an ability to imagine the effects of hypothesized actions, while the efficiency of human search arises from rich action priors which are updated via observations of the world. We instantiate these components in the "Sample, Simulate, Update" (SSUP) model and show that it captures human performance across 30 levels of the Virtual Tools game. More broadly, this model provides a mechanism for explaining how people condense general physical knowledge into actionable, task-specific plans to achieve flexible and efficient physical problem-solving.

Keywords: Intuitive Physics, Physical Problem Solving, Tool Use

While trying to set up a tent on a camping trip, you realize that the ground is too hard for the tent stakes, and you have no hammer. What would you do? You might look around for a suitable hammer substitute, passing over objects like pinecones or water bottles in favor of a graspable rock. And if that rock failed to drive in the stakes at first, you might try a different grip, or search for a heavier rock. Most likely, you would only need a handful of attempts before you found an approach that works.

Determining how to pound in tent stakes without a hammer is an example of the flexibility and efficiency of more general physical problem solving. It requires a causal understanding of how the physics of the world works, and sophisticated abilities for inference and learning to construct plans that solve a novel problem. Consider how, when faced with the tent stake challenge, we do not choose an object at random; we choose a rock because we believe we know how we could use it to generate sufficient force on the stake. And

Department of Brain and Cognitive Sciences

if we find that the first rock fails, we again search around for a solution, but use the knowledge of our failures to guide our future search. This style of problem solving is a very structured sort of trial-and-error learning: our search has elements of randomness, but within a plausible solution space, such that the goal can often be reached very quickly.

Here we study the cognitive and computational underpinnings of flexible tool use. While human tool use relies on a number of cognitive systems – for instance, knowing how to grasp and manipulate an object, or understanding how a particular tool is typically used – here we focus on "mechanical reasoning," or the ability to spontaneously repurpose objects in our environment to accomplish a novel goal (Goldenberg & Spatt, 2009; Orban & Caruana, 2014; Osiurak & Badets, 2016).

We target this mechanical reasoning because it is the type of tool use that is quintessentially human. While other animals can manipulate objects to achieve their aims, only a few species of birds and primates have been observed to spontaneously use objects in novel ways, and we often view these activities as some of the most "human-like" forms of animal cognition (e.g., Fig. 1A,B; Shumaker, Walkup, & Beck, 2011). Similarly, while AI systems have become increasingly adept at perceiving and manipulating objects, none perform the sort of rapid mechanical reasoning that people do. Some artificial agents learn to use tools from expert demon-

Corresponding author:

Kelsey R. Allen

Massachusetts Institute of Technology

<sup>77</sup> Massachusetts Avenue, Cambridge, MA 02139

krallen@mit.edu



*Figure 1.* Examples of using objects to achieve a goal. (**A**) Bearded capuchin monkey opening a cashew nut with an appropriately sized stone (photo by Tiago Falótico; Luncz et al., 2016). (**B**) New Caledonian crow using heavy blocks to raise the water level in a tube in order to retrieve food (Jelbert, Taylor, Cheke, Clayton, & Gray, 2014). (**C**) Toddler using a shovel to reach a ball (from youtu.be/hwrNQ93-568?t=198). (**D**) One illustrative trial in the Virtual Tools game (https://sites.google.com/view/virtualtoolsgame). (i) The player must get the red object into the green goal using one of the three tools. (ii) The player chooses a tool and where to place it. (iii) Physics is turned "on" and the tool interacts with other objects. The action results in a near miss.

strations (Xie, Ebert, Levine, & Finn, 2019), which limits their flexibility. Others learn from thousands of years of simulated experience (Baker et al., 2019), which is significantly longer than required for people. Still others can reason about mechanical functions of arbitrary objects but require perfect physical knowledge of the environment (Toussaint, Allen, Smith, & Tenenbaum, 2018), which is unavailable in realworld scenarios. In contrast, even young humans are capable tool users: by the age of four they can quickly choose an appropriate object and determine how to use it to solve a novel task (e.g., picking a hooked rather than straight pipe cleaner to retreive an object from a narrow tube, Fig. 1C; Beck, Apperly, Chappell, Guthrie, & Cutting, 2011).

What are the cognitive systems that let us use tools so flexibly, and accomplish our goals so rapidly? It has been suggested that mechanical reasoning relies on mental simulation, which lets us predict how our actions will cause changes in the world (Osiurak & Badets, 2016). This general purpose simulation is a necessary component that supports our ability to reason about objects in novel environments, but by itself cannot explain how we make and update our plans so quickly. We propose that another key to rapid tool use is knowing what sorts of actions to even consider – both from an initial understanding of what actions are useful, and by updating this belief from observing the outcome of our actions, in simulation and in reality.

This paper makes two contributions. First, we introduce

the Virtual Tools game, which presents a suite of physical problem solving challenges, and allows for precise, quantifiable comparisons between human and machine agents. Second, we present a minimal model of flexible tool use, called "Sample, Simulate, Update" (SSUP). This model is built around an efficient albeit noisy simulation engine that allows the model to act *flexibly* across a wide variety of physical tasks. To solve problems rapidly, the SSUP model contains rich knowledge about the world in the form of a structured prior on candidate tools and actions likely to solve the problem, which allows it to limit its simulations to promising candidates. It further learns from its simulations and from observing the outcome of its own actions to update its beliefs about what those promising candidates should be. Across 30 Virtual Tools levels in two experiments, we show that an instantiation of the SSUP model captures the relative difficulties of different levels for human players, the particular actions performed to attempt to solve each level, and how the solution rates for each level evolve.

#### The Virtual Tools game

Inspired by human tool use, as well as mobile physics games (Brain it On, 2015), we propose the Virtual Tools game as a platform for investigating the priors, representations, and planning and learning algorithms used in physical problem solving (https://sites.google.com/view/virtualtoolsgame). This game asks players to place one of several objects ("tools") into a two-dimensional dynamic physical environment in order to achieve a goal: getting a red object into a green region (Fig. 1D). This goal is the same for every level, but what is required to achieve it varies greatly. Once a single tool is placed, the physics of the world is enabled so that players see the effect of the action they took. If the goal is not achieved, players can "reset" the world to its original state and try again; they are limited to a single action on each attempt. We designed 30 levels -20 for the original experiment (Fig. 2) and 10 for a validation experiment (Fig. 7A) - to test concepts such as 'launching', 'blocking', and 'supporting'. Of the first 20 levels, 12 were constructed in six 'matched pairs' which incorporated small differences in the goals or objects in the scene to test whether subtle differences in stimuli would lead to observable differences in behavior.

The Virtual Tools game presents particular challenges that we believe underlie the kinds of reasoning required for rapid physical problem solving more generally. First, there is a **diversity of tasks** that require different strategies and physical concepts to solve, but employ shared physical dynamics that approximate the real world. Second, the game requires **longhorizon causal reasoning**. Since players can only interact with the game by placing a single object, they must be able to reason about the complex cause and effect relationships



*Figure 2.* Twenty levels used in the Virtual Tools game. Players choose one of three tools (shown to the right of each level) to place in the scene in order to get a red object into the green goal area. Black objects are fixed, while blue objects also move; grey regions are prohibited for tool placement. Levels denoted with A/B labels are matched pairs.

of their action long into the future when they can no longer intervene. Finally, the game elicits **rapid trial-and-error learning** in humans. Human players do not generally solve levels on their first attempt, but also generally do not require more than 5-10 attempts in order to succeed. People demonstrate a wide range of problem-solving behaviors, including "a-ha" insights where they suddenly discover the right idea for how to solve a particular task, as well as incremental trialand-error strategy refinement. Figure 3 demonstrates how this occurs in practice, showing four different examples of participants learning rapidly or slowly, and discovering different ways to use the tools across a variety of levels.

## Sample, Simulate, Update Model (SSUP)

We consider the components required to capture both the flexibility and efficiency of human tool use. We propose that people achieve flexibility through an internal mental model that allows them to imagine the effects of actions they may have never tried before ("Simulate"). However, a mental model alone is not sufficient – there are far too many possible actions that could be simulated, many of which are uninformative and unlikely to achieve a specific goal. Some mechanism for guiding an internal search is necessary to focus on useful parts of the hypothesis space. We therefore propose people use structured, object-oriented priors ("Sample") and a rapid belief updating mechanism ("Update") to guide search towards promising hypotheses. We formalize human tool use with these components as the "Sample, Simulate, Update" model (SSUP; Fig. 4A).

SSUP is inspired by the theory of "problem solving as search" (Newell & Simon, 1972), as well as Dyna and other model-based policy optimization methods (Deisenroth, Neumann, Peters, et al., 2013; Sutton, 1991). Crucially, we posit that structured priors and physical simulators must already

be in place in order to solve problems as rapidly as people; thus unlike most model-based policy optimization methods, we do not perform online updates of the dynamics model.

We want to emphasize that we view SSUP as a general modeling framework for physical problem solving, and only present here one instance of that framework: the minimal model (described below, with more detail in SI Appendix, Section S2) that we think is needed to capture basic human behavior in the Virtual Tools game. In the discussion we highlight ways the model will need to be improved in future work, as well as aspects of physical reasoning that rely on a richer set of cognitive systems going beyond the framework presented here.

*Sample: object-based prior.* At a minimum, the actions we should consider to achieve any goal should have the potential to impact our environment. We therefore incorporate an object-based prior for sampling actions. Specifically, the model selects one of the movable objects in the scene, then chooses an x-coordinate in an area that extends slightly beyond the width of the object, and a y-coordinate either above or below that object (Fig. 4B: Prior). For tool choice, we assume participants are equally likely to choose any of the three tools since all tools in the game were designed to be unfamiliar to participants. Samples from this distribution are used to initialize search.

*Simulate: a noisy physics engine.* In order to determine which sampled actions are worth trying in the world, we assume people use an "Intuitive Physics Engine" (Battaglia, Hamrick, & Tenenbaum, 2013) to flexibly imagine the effects of their actions. This engine is able to simulate the world forwards in time with approximately correct but stochastic dynamics (Sanborn, Mansinghka, & Griffiths, 2013; Smith & Vul, 2013). Determining the effect of a proposed action therefore involves applying that action to one's mental representation, and using the Intuitive Physics Engine to posit



*Figure 3.* Examples of participants' behavior on three levels, representative of rapid trial-and-error learning: Initial plans are structured around objects, followed by exploring to identify more promising strategies and then refining actions until success. Objects start as shown by light blue/red outlines and follow paths traced out by colored lines. Possible tool choices shown to the right. (A) In the Catapult level, a useful strategy is often identified immediately and rapidly fine-tuned. (B) Other participants first try an unsuccessful strategy but then switch to a more viable strategy and refine it. (C) The Launch (B) level is designed to prevent obvious solutions. This participant may have initially believed the ball would start rolling and attempted to use a tool as a bridge. When this failed they realized they needed to launch the ball, but only discovered after several trials how to use a tool in a non-obvious way to accomplish this, via a hooking motion around the blocking ledge. They then took several more trials to fine-tune this action. (D) In the SeeSaw level, a participant realized on the second attempt they must support the platform for the ball to roll across, then tried different ways of making this happen.



*Figure 4.* (A) The SSUP algorithm. (B) A diagram of the model for the Virtual Tools game. It incorporates an object-based prior, a simulation engine for filtering proposals, and an update module that suggests new proposals based on observations "in the mind" and from actions taken in the world. (C) Illustration of the policy  $\pi'$  evolving while attempting a level. Colored patches represent the Gaussian policy for each tool.

the range of ways that action might cause the world to unfold (Craik, 1943; Dasgupta, Smith, Schulz, Tenenbaum, & Gershman, 2018). Here we implement simulation using a game physics engine with noisy dynamics. People characteristically have noisy predictions of how collisions will resolve (Smith & Vul, 2013), and so for simplicity we assume uncertainty about outcomes is driven only by noise in those collisions (the direction and amount of force that is applied between two colliding objects).<sup>1</sup>

Since the internal model is imperfect, to evaluate an action we produce a small number of stochastic simulations  $(n_{sims}, set here at 4)$  to form a set of hypotheses about the outcome. To formalize how good an outcome is (the reward of a given action), we borrow an idea from the causal reasoning literature for how people conceptualize "almost" succeeding (Gerstenberg, Goodman, Lagnado, & Tenenbaum, 2015). "Almost" succeeding is not a function of the absolute distance an action moved you towards your goal, but instead how much of a *difference* that action made. To capture this, the minimum distance between the green goal area and any of the red goal objects is recorded; these values are averaged across the simulations and normalized by the minimum distance that would have been achieved if no tool had been added. The reward used in SSUP is 1 minus the normalized distance, so that closer objects lead to higher reward.

Once the model finds a good enough action (formalized as the average reward being above some threshold), it takes that action "in the world." Additionally, to model time limits for thinking, if the model considers more than T different action proposals without acting (set here at 5), it takes the best action it has imagined so far. We evaluate the effect of all parameter choices in a sensitivity analysis (see SI Appendix, Fig. S1).

*Update: learning from thoughts and actions.* So far we have described a way of intelligently initializing search to avoid considering actions that will not be useful. But what if the prior still presents an intractably large space of possible actions?

To tackle this, we incorporate an update mechanism that learns from both simulated and real experience to guide future search towards more promising regions of the hypothesis space (Juechems & Summerfield, 2019). This is formally defined as a Gaussian mixture model policy over the three tools and their positions,  $\pi'(s)$ , which represents the model's belief about high value actions for each tool.  $\pi'(s)$  is initialized with samples from the object-oriented prior, and updated using a simple policy gradient algorithm (Williams, 1992). This algorithm will shape the posterior beliefs around areas to place each tool which are expected to move target objects close to the goal, and are therefore likely to contain a solution. Such an update strategy is useful when it finds high value actions that are nearby successful actions, but may also get stuck in local optima where a successful action does not exist. We therefore use a standard technique from reinforcement learning: epsilon-greedy exploration. With epsilon-greedy exploration, potential actions are sampled from the policy  $1 - \epsilon\%$ of the time, and from the prior  $\epsilon\%$  of the time. Note that this exploration is only used for proposing internal simulations; model actions are chosen based on the set of prior simulation outcomes. This is akin to thinking of something new, instead of focusing on an existing strategy.

# Results

We analyze human performance on the first 20 levels of the Virtual Tools game and compare humans to the SSUP model and alternates, including SSUP models with ablations and two alternate learning baselines. We show that the full SSUP model best captures human performance. Access to the game and all data including human and model placements is provided at https://sites.google.com/view/ virtualtoolsgame.

## Human results

Experiments were approved by the MIT Committee on the Use of Humans as Experimental Subjects under protocol #0812003014. Participants were notified of their rights before the experiment, were free to terminate participation at any time by closing the browser window, and were compensated monetarily for their time.

We recruited 94 participants through Amazon Mechanical Turk and asked each participant to solve 14 levels: all 8 of the unmatched levels, and one variation of each of the 6 matched pairs (randomly selected).

Participants could choose to move on once a problem was solved, or after two minutes had passed. See SI Appendix, Section S1 for further details.

The variation in difficulty between levels of the game was substantial. Participants showed an average solution rate of 81% (sd = 19%), with the range covering 31% for the hardest level to 100% for the easiest. Similarly, participants took an average of 4.5 actions (sd = 2.5) for each level, with a range from 1.5 to 9.4 average attempts. Even within trials, there was a large amount of heterogeneity in the number of actions participants used to solve the level. This would be expected with "rapid trial-and-error" learning: participants who initially tried a promising action would solve the puzzle quickly, while others explored different actions before happening on promising ones (e.g., Fig. 3).

Behavior differed across all six matched level pairs. We study whether these subtle differences do indeed affect behavior, even without feedback on the first action, by asking whether we can identify which level variant each action came

<sup>&</sup>lt;sup>1</sup>We also considered models with additional sources of physics model uncertainty added, but found that the additional parameters did not improve model fit, so we do not analyze those models here.



*Figure 5*. (A) Comparison of average number of human participants' attempts for each level with average number of attempts for the SSUP model. Bars indicate 95% confidence intervals on estimates of the means. (B) Comparison of human participants' accuracy on each trial versus the accuracy of the SSUP model. (C) Comparison of human participants' accuracy to all alternate models. Numbers correspond to the trials in Fig. 2.

from. We find these actions are differentiable across matched levels in 'Shafts', 'Prevention', 'Launch' and 'Table' on the first attempt, but not 'Falling' or 'Towers' (see SI Appendix, Fig. S11 and Sec. S6A for details). However, participants required a different number of actions to solve every level (all ts > 2.7, ps < 0.01). This suggests that people are paying attention to subtle differences in the scene or goal to choose their actions.

#### Model results

We investigate several metrics for comparing the models to human data. First, we look at how quickly and how often each model solves each level, and whether that matches participants. This is measured as the correlation and root mean squared error (RMSE) between the average number of participant attempts for each level and the average number of model attempts for each level, and the correlation and RMSE between human and model solution rates. The SSUP model explains the patterns of human behavior across the different levels well (SI Appendix, Table S2). It uses a similar number of attempts on each level (r = 0.71; 95% CI = [0.62, 0.76]; mean empirical attempts across all levels: 4.48, mean model attempts: 4.24; Fig. 5A) and achieves similar accuracy (r =0.86; 95% CI = [0.76, 0.89]; Fig. 5B).

Across many levels, the SSUP model not only achieves the same overall solution rate as people, but approaches it at the same rate. We measure this by looking at cumulative solution rates – over all participants or model runs, what proportion solved each level within *X* placements – and find that people and the model often demonstrate similar solution profiles (Fig. 6A; see SI Appendix, Section S6B for quantitative comparison).

We can look in more detail how the model accomplishes this by comparing both the first actions that people and the model takes (Fig. 6B), and the actions that both take to solve a level (Fig. 6C). Like our human participants, the model takes significantly different actions on the first attempt between matched level pairs (see SI Appendix, Sec. S6A). More generally, both people and the model will often begin with a variety of plausible actions (e.g., Catapult). In some cases, both will attempt initial actions that have very little impact on the scene (e.g., SeeSaw and Prevention (B)); this could be because people cannot think of any useful actions and so decide to try something, similar to how the model can exceed its simulation threshold. However, in other cases, the model's initial predictions diverge from people, and this leads to a different pattern of search and solutions. For instance, in Falling (A), the model quickly finds that placing an object under the container will reliably tip the ball onto the ground, but people are biased to drop an object from above. Because of this, the model often rapidly solves the level with an object below, whereas a proportion of participants find a way to flip the container from above; this discrepancy can also be seen in the comparison of number of attempts before the solution, where the model finds a solution quickly, while people take a good deal longer (Fig. 5A). For comparisons of the first and last actions across all levels, see SI Appendix, Fig. S11.

**Model comparisons on Virtual Tools.** We compare the full SSUP model against a set of six alternate models. Three models investigate the contribution of each SSUP component by removing the prior, simulation, or updating individually. Two models propose alternate solution methods: learning better world models rather than learning over actions (Parameter Tuning) or replacing the prior and simulator with a learned proposal mechanism (DQN + Updating). The Parameter Tuning alternate model uses inference to learn object densities, frictions and elasticities from observed trajectories. The learned proposal mechanism corresponds to a model-free deep reinforcement learning agent (Mnih et al.,



*Figure 6.* (A) Cumulative solution rate over number of placements for participants vs. the SSUP model. (B) Distribution of model actions (background) versus human actions (points) on the first and last attempts of the level for a selection of four levels. The distribution of model actions is estimated based on fitting a Kernel Density Estimate to the actions taken by the model across 250 simulations. Colors indicate the tool used, with the tools and associated colors shown to the right of each level. In most levels, the SSUP model captures the evolution of participants' solutions well, including the particular actions chosen; in the few cases that it differs, there is no alternative model that systematically explains these differences.

2015) which is trained on a set of 4500 randomly generated levels of the game (see SI Appendix, Sec. S5), and then updated online for each of the 20 testing levels using the same mechanism as SSUP. This model has substantially more experience with the environment than other models, and serves as a test of whether model-free methods can make use of this experience to learn generalizable policies that can guide rapid learning. Finally, we compare to a "Guessing" baseline for performance if an agent were to simply place tools randomly. See Fig. 5C and SI Appendix, Table S2 for these comparisons.

Eliminating any of the three SSUP components causes a significant decrease in performance (measured as deviation between empirical and model cumulative solution curves; all bootstrapped ps < 0.0001; see SI Appendix, Sec. S6B, Fig. S6 for further detail). The reduced models typically require more attempts to solve levels because they are either searching in the wrong area of the action space (No Prior), attempting actions that have no chance of being successful (No Simulation), or do not guide search towards more promising areas (No Updating).

DQN + Updating performs worst of all plausible alternate models, using the most actions and solving levels at a rate barely over chance. Because this is equivalent to the No Simulation model with a different prior, its poor performance suggests that generalized action policies cannot easily be learned from repeatedly playing similar levels (see SI Appendix, Sec. S5).

Because the Parameter Tuning model is equivalent to the No Updating model except that the properties of the dynamics model can be learned in Parameter Tuning, comparing those two models allows us to test whether we need to assume that people are learning the dynamics of the world in this game. The fact that both models perform roughly equivalently (see Fig. 5C) suggests that we do not need this assumption here.

Finally, we quantified how well each model captured the particular actions people took. Due to heterogeneity in participants' responses, we were unable to cleanly differentiate models' performance except to find that the DQN + Updating model underperformed the rest (see SI Appendix, Sec. S6C). However, no model reached the theoretical noise ceiling, suggesting components of the SSUP framework could be improved to better explain participants' actions (see the Discussion).

## Validation on novel levels

We conducted a second experiment to test whether the models generalize to novel levels and physical concepts without tuning hyperparameters. For this experiment, we created 10 new levels: 6 novel level types and 4 variants of the originals (Fig 7A), testing an independent sample of 50 partici-

#### ALLEN, SMITH, TENENBAUM



*Figure* 7. Results on 10 additional trials. (A) Trials used for the second experiment. (B) The cumulative solution rate for participants and the SSUP model. (C) Comparison of the number of human and model actions by trial. (D) Comparison of human and model accuracy on each trial.

pants on all levels. The 6 novel level types were designed to test new physical strategies, including balancing, breaking, and removing objects from a ball's path. All other experimental details were identical to the main experiment.

Without tuning any model parameters, we find a good correspondence between human and model solution rates (Fig. 7B), and a strong correlation between the model's performance and human performance across number of placements (Fig. 7C, r = 0.85) and accuracy (Fig. 7D, r = 0.95). Similar to the main experiment, we find a decrement in performance if the prior or simulation are removed, or for the DQN + Updating model (all bootstrapped ps < 0.0001; SI Appendix, Fig. S7). However, while numerically worse, we do not find a reliable difference if the update mechanism is removed (p = 0.055) or swapped for model learning (p = 0.346), suggesting that the particular reward function or update procedure might be less applicable to these levels (see SI Appendix, Sec. S6B).

#### Discussion

We introduce the Virtual Tools game for investigating flexible physical problem solving in humans and machines, and show that human behavior on this challenge expresses a wide variety of trial-and-error problem solving strategies. We also introduce a model for human physical problem solving: "Sample, Simulate, Update." The model presumes that to solve these physics problems, people rely on an internal model of how the world works. Learning in this game therefore involves condensing this vast world knowledge to rapidly learn how to act in each instance, using a structured trial-and-error search.

#### **Model limitations**

Although the SSUP model we used solves many of the levels of the Virtual Tools game in a human-like way, we believe that this is still only a first approximation to the rich set of cognitive processes that people bring to the task. In particular, there are at least two ways in which the model is insufficient: its reliance on very simple priors, and its planning and generalizing only in the forwards direction.

We can see the limits of the object-based prior in the Falling (A) level (Fig. 5B): people are much less likely to consider placing an object underneath the container to tip it over. Instead, many people try to tip it over from above, even though this is more difficult. In this way, people's priors over strategies are *context specific*, which causes them to be slower than the model in this level. In other cases, this context specificity is helpful: for instance, in the hypothetical level shown in Fig. 8A, there is a hole that one of the tools fits suspiciously perfectly into. Many people notice this coincidence quickly, but because the model cannot assess how tools might fit into the environment without running a simulation,

it only succeeds 10% of the time. In future work, a more complex prior could be instantiated in the SSUP framework, but it remains an open question how people might form these context-specific priors, or how they might be shaped over time via experience.

People show greater flexibility than our model in the ability to work backwards from the goal state to find more easily solvable sub-goals (Anderson, 1993). In the hypothetical level in Fig. 8B, the catapult is finicky, which means that most catapulting actions will not make it over the barrier, and therefore will never hit the ball on the left. Instead, the easiest way to increase the objective function is by the incorrect strategy of knocking the ball on the right to get close to the goal, and therefore the model only solves the level 8% of the time. Working backwards to set the first sub-goal of launching the ball over the barrier would prevent getting stuck with knocking the ball as a local minimum. From an engineering standpoint, creating sub-goals is natural with discrete problem spaces (Newell & Simon, 1972), but it is less clear how these might be discovered in the continuous action space of the Virtual Tools game.

## **Related cognitive systems**

There is an extensive body of research into the cognitive systems that underlie the use of real-world tools, including understanding how to manipulate them and knowing their typical uses (e.g. Beck et al., 2011; Orban & Caruana, 2014; Osiurak & Badets, 2016; Vaesen, 2012). Here our focus was on "mechanical knowledge" of tools: how to use objects in novel situations. However, in real-world tool use, these systems work together with motor planning and semantic knowledge of tools. Future work can focus on these links, such as how novel tools become familiar, or how our motor limits constrain the plans we might consider.

The Virtual Tools game presents a problem solving task that blends facets of prior work, but encompasses a novel challenge. To rapidly solve these problems requires good prior knowledge of the dynamics – unlike Complex Problem Solving in which the dynamics are learned in an evolving situation (Frensch & Funke, 1995) – and further iteration once a promising solution is considered – unlike the 'a-ha' moment that leads immediately to a solution in Insight Problem Solving (Chu & MacGregor, 2011; Gick & Holyoak, 1980). Unlike in traditional model-based or model-free reinforcement learning, in this task people bring rich models of the world that they can quickly tailor to specific, novel problems.

Distilling rich world knowledge to useful task knowledge is necessary for any agent interacting with a complex world. One proposal for how this occurs is "learning by thinking" (Lombrozo, 2018): translating knowledge from one source (internal models of physics) to another, more specific instantiation (a mapping between actions and outcomes on *this particular* level). We show how SSUP instantiates one example



*Figure 8.* Two problems that demonstrate limitations of the current model. (A) A "suspicious conicidence" that one tool fits perfectly in the hole. (B) Creating a 'sub-goal' to launch the ball onto the other side is useful.

of "learning by thinking": by training a policy with data from an internal model. Evidence for this sort of knowledge transfer has been found in people (Gershman, Markman, & Otto, 2014; Gershman, Zhou, & Kommers, 2017), but has focused on simpler discrete settings in which the model and policy are jointly learned.

## Virtual Tools as an AI Challenge

In preliminary experiments with model-free reinforcement learning approaches (Mnih et al., 2015), we found limited generalization with inefficient learning across almost all of the Virtual Tools levels (see SI Appendix, Section S5) despite significant experience with related levels.

Based on our human experiments, we believe that modelbased approaches will be required to be able to play games like Virtual Tools. Such approaches are becoming increasingly popular in machine learning (Weber et al., 2017), especially when combined with "learning-to-learn" techniques that can learn to adapt quickly to new tasks (Finn, Abbeel, & Levine, 2017; Schmidhuber, Zhao, & Schraudolph, 1998). Learning these models remains challenging, but approaches that incorporate added structure have excelled in recent years (Chang, Ullman, Torralba, & Tenenbaum, 2016; Jayaraman, Ebert, Efros, & Levine, 2018). Within the AI and robotics communities, model-based methods are already popular (Garcia, Prett, & Morari, 1989; Kaelbling & Lozano-Pérez, 2010; Toussaint et al., 2018). Remaining challenges include how to learn accurate enough models that can be used with raw sensor data (Kroemer, Niekum, & Konidaris, 2019), and how to handle dynamic environments.

Virtual Tools adds to a growing set of environments that test artificial agents' abilities to predict and reason using physics, such as the concurrently developed PHYRE benchmark (Bakhtin, van der Maaten, Johnson, Gustafson, & Girshick, 2019) and others (Bapst et al., 2019; Ge, Lee, Renz, & Zhang, 2016; Wenke, Saunders, Qiu, & Fleming, 2019). In contrast, our focus is on providing problems that people find challenging but intuitive, where solutions are non-obvious and do not rely on precise knowledge of world dynamics. By contributing human data to compare artificial and biological intelligence, we hope to provide a test-bed for more human10

like artificial agents.

## **Future empirical directions**

This work provides an initial foray into formalizing the computational and empirical underpinnings of flexible tool use, but there remains much to study. For instance, we do not find evidence that people learn more about the world, perhaps because here there is little benefit to additional precision here. But there are cases where learning the dynamics is clearly helpful (e.g., discovering that an object is abnormally heavy, or glued down), and we *would* expect people to update their physical beliefs in these cases. When and in what ways people update their internal models to support planning is an important area of study.

Children can discover how to use existing objects earlier than they can make novel tools (Beck et al., 2011), suggesting that tool creation is more challenging than tool use. Yet it is the ability to make and then pass on novel tools that is theorized to drive human culture (Tomasello, 1999). It is therefore important to understand not just how people use tools, but also how they develop and transmit them, which we can study by expanding the action space of the Virtual Tools game.

### Conclusion

Understanding how to flexibly use tools to accomplish our goals is a basic and central cognitive capability. In the Virtual Tools game, we find that people efficiently use tools to solve a wide variety of physical problems. We can explain this rapid trial-and-error learning with the three components of the SSUP framework: rich prior world knowledge, simulation of hypothetical actions, and the ability to learn from both simulations and observed actions. We hope this empirical domain and modeling framework can provide the foundations for future research on this quintessentially human trait: using, making, and reasoning about tools, and more generally shaping the physical world to our ends.

## Acknowledgments

The authors would like to thank Leslie Kaelbling, Roger Levy, Eric Schulz, Jessica Hamrick, and Tom Silver for helpful comments, and Mario Belledonne for help with the Parameter Tuning model. This work was supported by NSF STC award CCF-1231216, ONR MURI N00014-13-1-0333, and research grants from ONR, Honda, and Mitsubishi Electric.

#### References

Anderson, J. R. (1993). Problem Solving and Learning. Am Psychol, 10.

- Baker, B., Kanitscheider, I., Markov, T., Wu, Y., Powell, G., Mc-Grew, B., & Mordatch, I. (2019). Emergent tool use from multi-agent autocurricula. *ICLR*.
- Bakhtin, A., van der Maaten, L., Johnson, J., Gustafson, L., & Girshick, R. (2019). Phyre: A new benchmark for physical reasoning. *NeurIPS*.
- Bapst, V., Sanchez-Gonzalez, A., Doersch, C., Stachenfeld, K. L., Kohli, P., Battaglia, P. W., & Hamrick, J. B. (2019, April). Structured agents for physical construction. *ICML*. Retrieved 2019-04-08, from http://arxiv.org/abs/1904 .03177
- Battaglia, P. W., Hamrick, J. B., & Tenenbaum, J. B. (2013, November). Simulation as an engine of physical scene understanding. *PNAS*, 110(45), 18327–18332. Retrieved 2019-03-07, from http://www.pnas.org/cgi/doi/10.1073/ pnas.1306572110 doi: 10.1073/pnas.1306572110
- Beck, S. R., Apperly, I. A., Chappell, J., Guthrie, C., & Cutting, N. (2011, May). Making tools isn't child's play. *Cognition*, 119(2), 301–306. Retrieved 2019-05-28, from https://linkinghub.elsevier.com/retrieve/ pii/S0010027711000163 doi: 10.1016/j.cognition.2011 .01.003
- Brain it on. (2015). https://brainitongame.com/.
- Chang, M. B., Ullman, T., Torralba, A., & Tenenbaum, J. B. (2016, December). A Compositional Object-Based Approach to Learning Physical Dynamics. *International Conference on Learning Representations*. Retrieved 2019-05-20, from http://arxiv.org/abs/1612.00341 (arXiv: 1612.00341)
- Chu, Y., & MacGregor, J. N. (2011, February). Human Performance on Insight Problem Solving: A Review. *The Journal of Problem Solving*, 3(2). Retrieved 2019-05-28, from https:// docs.lib.purdue.edu/jps/vol3/iss2/6 doi: 10 .7771/1932-6246.1094
- Craik, K. J. W. (1943). *The Nature of Explanation*. CUP Archive. (Google-Books-ID: wT04AAAAIAAJ)
- Dasgupta, I., Smith, K. A., Schulz, E., Tenenbaum, J. B., & Gershman, S. J. (2018). Learning to act by integrating mental simulations and physical experiments. In *Proceedings of the* 40th Annual Meeting of the Cognitive Science Society. Retrieved 2019-03-07, from http://biorxiv.org/lookup/ doi/10.1101/321497 doi: 10.1101/321497
- Deisenroth, M., Neumann, G., Peters, J., et al. (2013). A survey on policy search for robotics. *Foundations and Trends in Robotics*, 2(1–2), 1–142.
- Finn, C., Abbeel, P., & Levine, S. (2017). Model-agnostic metalearning for fast adaptation of deep networks. In *Proceedings* of the 34th international conference on machine learningvolume 70 (pp. 1126–1135).
- Frensch, P. A., & Funke, J. (1995). Complex problem solving: The european perspective. Psychology Press.
- Garcia, C. E., Prett, D. M., & Morari, M. (1989). Model predictive control: theory and practice—a survey. *Automatica*, 25(3), 335–348.
- Ge, X., Lee, J. H., Renz, J., & Zhang, P. (2016). Hole in one: Using qualitative reasoning for solving hard physical puzzle problems. In *Proceedings of the twenty-second european conference on artificial intelligence* (pp. 1762–1763).

- Gershman, S. J., Markman, A. B., & Otto, A. R. (2014, February). Retrospective revaluation in sequential decision making: A tale of two systems. J Exp Psychol Gen, 143(1), 182–194. Retrieved 2019-07-21, from https://search.proquest.com/docview/ 1237023272/abstract/53A55607AF4A47F7PQ/1 doi: http://dx.doi.org/10.1037/a0030844
- Gershman, S. J., Zhou, J., & Kommers, C. (2017, July). Imaginative Reinforcement Learning: Computational Principles and Neural Mechanisms. J Cognitive Neurosci, 29(12), 2103–2113. Retrieved 2019-07-18, from https://doi.org/10.1162/jocn\_a\_01170 doi: 10.1162/jocn\_a\_01170
- Gerstenberg, T., Goodman, N. D., Lagnado, D. A., & Tenenbaum, J. B. (2015). How, whether, why: Causal judgments as counterfactual contrasts. In 37th Annual Meeting of the Cognitive Science Society (p. 6).
- Gick, M. L., & Holyoak, K. J. (1980, July). Analogical problem solving. *Cognitive Psychol*, 12(3), 306–355. Retrieved 2019-05-28, from http://www.sciencedirect .com/science/article/pii/0010028580900134 doi: 10.1016/0010-0285(80)90013-4
- Goldenberg, G., & Spatt, J. (2009). The neural basis of tool use. Brain, 132(6), 1645-1655. doi: 10.1093/brain/awp080
- Jayaraman, D., Ebert, F., Efros, A. A., & Levine, S. (2018). Timeagnostic prediction: Predicting predictable video frames. *In*ternational Conference on Learning Representations.
- Jelbert, S. A., Taylor, A. H., Cheke, L. G., Clayton, N. S., & Gray, R. D. (2014). Using the aesop's fable paradigm to investigate causal understanding of water displacement by New Caledonian crows. *PloS One*, 9(3), e92895.
- Juechems, K., & Summerfield, C. (2019). *Where does value come from?* (Preprint). PsyArXiv. doi: 10.31234/osf.io/rxf7e
- Kaelbling, L. P., & Lozano-Pérez, T. (2010). Hierarchical planning in the now. In Workshops at the twenty-fourth aaai conference on artificial intelligence.
- Kroemer, O., Niekum, S., & Konidaris, G. (2019). A review of robot learning for manipulation: Challenges, representations, and algorithms. arXiv preprint arXiv:1907.03146.
- Lombrozo, T. (2018). "Learning by thinking" in science and in everyday life. In *The Scientific Imagination*. New York, NY: Oxford University Press.
- Luncz, L. V., Falótico, T., Pascual-Garrido, A., Corat, C., Mosley, H., & Haslam, M. (2016). Wild capuchin monkeys adjust stone tools according to changing nut properties. *Sci Rep*, 6, 33089.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., ... others (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540), 529.
- Newell, A., & Simon, H. A. (1972). Human problem solving. Oxford, England: Prentice-Hall.
- Orban, G. A., & Caruana, F. (2014). The neural basis of human tool use. *Front Psychol*, 5. doi: 10.3389/fpsyg.2014.00310

- Osiurak, F., & Badets, A. (2016). Tool use and affordance: Manipulation-based versus reasoning-based approaches. *Psychological Review*, 123(5), 534–568. Retrieved 2019-05-21, from http://doi.apa.org/getdoi.cfm?doi=10 .1037/rev0000027 doi: 10.1037/rev0000027
- Sanborn, A. N., Mansinghka, V. K., & Griffiths, T. L. (2013, April). Reconciling intuitive physics and Newtonian mechanics for colliding objects. *Psychol Rev*, 120(2), 411–437. Retrieved 2019-05-31, from https://search.proquest.com/docview/ 1314700727/abstract/E7521331AA754DDBPQ/1 doi: http://dx.doi.org/10.1037/a0031912
- Schmidhuber, J., Zhao, J., & Schraudolph, N. N. (1998). Reinforcement learning with self-modifying policies. In *Learning to learn* (pp. 293–309). Springer.
- Shumaker, R. W., Walkup, K. R., & Beck, B. B. (2011). Animal Tool Behavior: The Use and Manufacture of Tools by Animals. JHU Press.
- Smith, K. A., & Vul, E. (2013, January). Sources of Uncertainty in Intuitive Physics. *TopiCS*, 5(1), 185–199. Retrieved 2019-03-07, from http://doi.wiley.com/10 .1111/tops.12009 doi: 10.1111/tops.12009
- Sutton, R. S. (1991, July). Dyna, an Integrated Architecture for Learning, Planning, and Reacting. *SIGART Bull.*, 2(4), 160– 163. Retrieved 2019-05-30, from http://doi.acm.org/ 10.1145/122344.122377 doi: 10.1145/122344.122377
- Tomasello, M. (1999). The Cultural Origins of Human Cognition. Harvard University Press. (Google-Books-ID: hRFfCgAAQBAJ)
- Toussaint, M., Allen, K. R., Smith, K. A., & Tenenbaum, J. B. (2018, June). Differentiable Physics and Stable Modes for Tool-Use and Manipulation Planning. In *Robotics: Science and systems*. Retrieved 2019-03-07, from http:// www.roboticsproceedings.org/rss14/p44.pdf doi: 10.15607/RSS.2018.XIV.044
- Vaesen, K. (2012, August). The cognitive bases of human tool use. Behav Brain Sci, 35(4), 203–218. Retrieved 2019-05-21, from https://www.cambridge.org/core/ product/identifier/S0140525X11001452/type/ journal\_article doi: 10.1017/S0140525X11001452
- Weber, T., Racanière, S., Reichert, D. P., Buesing, L., Guez, A., Rezende, D. J., ... others (2017). Imagination-augmented agents for deep reinforcement learning. *NeurIPS*.
- Wenke, S., Saunders, D., Qiu, M., & Fleming, J. (2019). Reasoning and generalization in rl: A tool use perspective. arXiv preprint arXiv:1907.02050.
- Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4), 229–256.
- Xie, A., Ebert, F., Levine, S., & Finn, C. (2019). Improvisation through physical understanding: Using novel objects as tools with visual foresight. *Robotics: Science and Systems*.