

Exploring Failure Cases in Multimodal Reasoning About Physical Dynamics

Sadaf Ghaffari, Nikhil Krishnaswamy

Colorado State University, Fort Collins, CO USA
{sadafgh,nkrishna}@colostate.edu

Abstract

In this paper, we present an exploration of LLMs’ abilities to problem solve with physical reasoning in situated environments. We construct a simple simulated environment and demonstrate examples of where, in a zero-shot setting, both text and multimodal LLMs display atomic world knowledge about various objects but fail to compose this knowledge in correct solutions for an object manipulation and placement task. We also use BLIP, a vision-language model trained with more sophisticated cross-modal attention, to identify cases relevant to object physical properties that that model fails to ground. Finally, we present a procedure for discovering the relevant properties of objects in the environment and propose a method to distill this knowledge back into the LLM.

Introduction

Large Language Models (LLMs), exemplified by modern generative systems like GPT-4, display impressive performance on benchmark tasks (Qin et al. 2023). However, these models continue to struggle with questions involving physical reasoning (Goertzel 2023; Wang et al. 2023), and demonstrate an apparent lack of “understanding” of common-sense physical properties and natural laws. Reasoning about properties of simple objects (such as balls being round) appeared to be weak in encoder-style language models such as BERT (Devlin et al. 2019). This could be partially attributed to such knowledge being so common as to be largely uninformative, and sparse in training corpora (Krishnaswamy and Pustejovsky 2022). Newer and larger generative language models appear not to have this problem (Ghaffari and Krishnaswamy 2023), with the autoregressive mechanism able to generate correct output about object properties and affordances. However, this knowledge appears to be atomic, unable to be composed into larger-scale reasoning. Solutions to simple problems that are intuitive to humans appear to not be considered by LLMs despite the level of world knowledge shown about individual items. An example that we will return to throughout this paper is the prompt in Fig. 1 and solutions thereto.

None of these solutions account for the properties of objects and physical laws like inertia and gravity. We argue that

EXAMPLE PROMPT

You need to get to the top of a platform that is 2 meters high. The highest you can jump is 1 meter. You have two blue cubes that are both 1 meter long on all sides, a blue sphere that is 1 meter in diameter, and a blue cylinder whose major axis is 1 meter long. How can you get to the top of the platform?

EXAMPLE RESPONSE (LLAMA 2-7B)

You can use the cube on top of the sphere to get to the desired height. Here’s how:

1. Place the cube on top of the sphere. This will give you a total height of 1.5 meters (1 meter for the cube + 0.5 meters for the sphere).
2. Place the cylinder on top of the cube. This will give you a total height of 2 meters (1 meter for the cube + 1 meter for the cylinder).

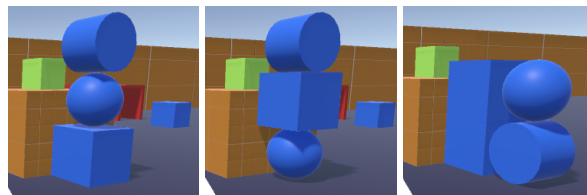


Figure 1: Example of physical reasoning prompt and response, and operationalizations of solutions as given by ChatGPT [L], LLaMA 2-7B (Touvron et al. 2023) [C], and LLaVA (Liu et al. 2023) [R].

this is due to a lack of grounding mechanisms to physical environments and access to non-linguistic modalities beyond imagery.

In this paper, we explore the current abilities of LLMs to perform reasoning about simple objects in dynamic environments. We construct a simulation environment in which to investigate these problems, and develop methods to evaluate LLMs’ solutions to a reasoning problem in terms of correctness of object selection and physical feasibility. We also examine the cross-modal grounding between images from

our environment and key terms in the reasoning task using a cross-attention architecture, and expose a number of shortcomings in this ability relative to the key concepts underlying our example problem. We develop a procedure by which an agent may exploit failure in the zero-shot outputs of LLMs as a trigger to investigate alternative solutions to the problem using object interactions and previously-encoded knowledge of the object semantics. Finally, we propose a method through which knowledge gained from the object interactions can be distilled back into the LLM, and avenues for future research.

Related Work

Recent explorations of generative LLMs have exposed certain classes of problems that they struggle with.

Linguistic and semantic problems Chaturvedi et al. (2022) show that LLM failures in the semantics of tasks like question answering result from features of the training data and regime. Asher et al. (2023) examine the capabilities of LLMs (broadly defined to include both BERT-style and GPT-style models) in semantic problems with a primary focus on quantification. They find that LLMs struggle to learn semantic properties of entailment and consistency, according to formal definitions, and do not generalize beyond the first level of the Borel Hierarchy. Additionally they find that generative models display inconsistency of outputs in this domain, even when temperature is set to 0.

Koçoń et al. (2023) evaluated both ChatGPT and GPT-4 on 25 diverse NLP/NLI tasks and showed that both models approached but did not exceed state of the art performance on multiple tasks, when compared to performance of models specialized for individual or subsets of tasks.

Mathematical and logic problems In what is likely the first independent empirical evaluation of ChatGPT, Shakarian et al. (2023) showed that its performance on math word problems (MWP) was strongly conditioned on a requirement to show its work, *a la* chain-of-thought (CoT) reasoning. However, they also demonstrated that, even given this condition, ChatGPT’s failure rate increases linearly relative to the number of addition and subtraction operations. We also demonstrate failures of LLMs in multistep reasoning problems, albeit in a different domain.

Chen, Zaharia, and Zou (2023) evaluated the performance of ChatGPT (as powered by GPT-3.5 and GPT-4) at different points in time on various problem types, including mathematical and multi-hop reasoning problems. They discovered unpredictable performance and found that the behavior of the same closed service can change substantially over time as models are updated. Increased performance on one type of tasks may have unpredictable effects on other types.

Shanahan (2022) maintains a division between the statistical functioning of LLMs and logical structures and methods underlying human reasoning, and mentions how prompt engineering techniques like CoT reasoning effectively reframe the LLM’s problem from a logical one to a statistical one, asking the model to sample from a distribution rather than make an *a priori* inference.

Goertzel (2023) argues that shortcomings in LLMs stem from a lack of cognitive architectures underlying the system, and that incremental (i.e., scale- and data-based) approaches are unlikely to be a path from current LLMs to artificial general intelligence (AGI). Our purview in this paper is not this problem specifically, though we do address how an underlying knowledge base mediated by a simulation can allow a neural system to solve a problem a zero-shot LLM apparently cannot, and propose a method for transferring some of this knowledge into an LLM.

Physical reasoning problems LLMs display often proficient performance on questions involving knowledge of physics (Polverini and Gregorcic 2023), but when it comes to applying common-sense physical reasoning in specific situations, even the most advanced multimodal models (e.g., Say-Can (Ahn et al. 2022) or PALM-E (Driess et al. 2023)) are highly sensitive to changes in the environment and input. Executing tasks in new situations usually involves crafting additional low-level behaviors for the new environment.

The domain in which we conduct our experiments is immediately reminiscent of pick-and-place tasks (Lobbezoo, Qian, and Kwon 2021), wherein an agent, usually a robotic arm, must select a target object and place the object at a specified location and orientation. This has been a popular problem for reinforcement learning (RL) approaches (Li et al. 2020) and has also been adapted to train models to learn about physical intuitions using tasks such as fall prediction (Lerer, Gross, and Fergus 2016). With the advent of generative LLMs, research in the object placement domain has explored how to use LLMs to serve as the back-end planner in situated environments (Ahn et al. 2022; Driess et al. 2023). In these cases, a standard constraint of pick and place tasks is still in force: the locations of objects must be known or deterministic. In addition, grounding objects to the linguistic input follows a standard approach of joint training over image and text distributions, with no multimodal signal besides pixels, and only integrates language with static images, not situated continuous environments.

Physical Reasoning in LLMs

LLMs’ difficulty with physical reasoning is manifest in Fig. 1, where solutions to the problem result in infeasible and unsafe configurations. The LLM is not grounded to the environmental dynamics. To further explore the extent to which an LLM does or does not generate correct outputs regarding relations and physical dynamics between objects, we created a simple scene against which to evaluate LLM outputs in this domain, along with a set of controls on prompts and image inputs. We compare outputs from ChatGPT, LLaMA 2 (Touvron et al. 2023), and LLaVA, an open-weight model that integrates CLIP ViT-L/14 (Radford et al. 2021) with a LLaMA-based generative language model.¹

Methodology

Following previous research (Ghaffari and Krishnaswamy 2022, 2023), we focused on a simple set of geometric ob-

¹We used the 7 billion parameter version of LLaVA-1.5 in these experiments.

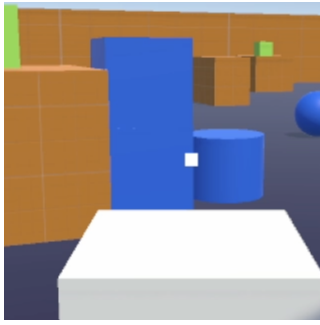


Figure 2: A feasible, physically stable solution to the platform reach problem, that uses the cylinder in an orientation that exploits the stability affordances of its flat ends.

jects inspired by children’s toys, and on the key concepts of flatness and roundness that determine the stability of structures constructed with these objects. We constructed a simple scene in the VoxWorld platform (Krishnaswamy et al. 2022) that reflected the scenario alluded to in Fig. 1. The scene contained 5 interactable objects, each roughly one meter tall/wide—two cubes, two cylinders, and a sphere—and 5 two-meter high platforms. In this paper, we considered only the task of moving objects and stacking them in a configuration to achieve a specified goal.

An agent (represented as a white cube) must navigate the scene and collect objects at the top of each platform. The agent could pick up, carry, rotate, place, and jump on top of any of the interactable objects, with a jump height limited to one meter. The VoxWorld platform for interactive agents takes in linguistic input and operationalizes it in the 3D world by converting the object motions described therein into movements within the environment, which were then executed in real time. Instructions in the output that did not describe moving or placing objects were ignored. Underspecified instructions were handled according to Krishnaswamy (2017), by randomly sampling locations that satisfied logically inferrable constraints.

The prompts we explored are variants of that given in Fig. 1, with changes to control for the output format and type of reasoning we wanted to push the LLM toward. We provided the prompt alone to ChatGPT and LLaMA 2, and to LLaVA along with an image of the scene, including the objects mentioned in the prompt. Prompts were slightly modified for the given model, such as by using “You are in the room shown in the image” with multimodal LLaVA to draw attention to the visual input. We used default temperature values for LLaMA 2 and LLaVA (0.6 and 0.2, respectively).

Use of a fixed problem allows us to assess LLMs’ output relative to a limited set of correct solutions, all of which involve building a staircase-like structure of the kind shown in Fig. 2. We can therefore quickly calculate an intersection-over-union (IoU) of the two sets: the objects that the LLMs included in its solution, and the objects in the correct solution (either two cubes and a cylinder or two cylinders and a cube). This calculation concerns only which objects were mentioned by the LLM in the generated solution, and not

how they were placed. Placement is evaluated using the *stable* metric as discussed below. The outputs were then operationalized in VoxWorld and the simulation run. After applying environmental physics, we scored the stability of the structure as the percentage of objects that remain in place. Therefore a correct solution should have a stability of 100% (all objects remain in place) and an IoU of 1.0 (all correct objects and only correct objects were selected). These two metrics indicate how right or wrong an LLM’s response is in the environment presented. While they may be somewhat abstract, evaluating at this high level allows us to examine the physical feasibility and correctness of the solution provided, without noise introduced in the specific word-level output by known problematic phenomena of LLMs such as hallucinations or problems with counting or size-based reasoning. They are intentionally tuned to the stacking task and provide a template of how to evaluate physical reasoning in tasks such as this, i.e., measuring both the object selection and key task-relevant concept.

Open-world Simulation

In our first setting, the environment was left relatively unconstrained. Objects remained in the locations in which they spawned, and initially various distractors remained in place (e.g., the collectable items on top of the platforms and the player, which were all rendered as cubes, all remained visible). The prompt also terminated in a relatively open-ended question *a la* Fig. 1, resulting in a *free-text response*. A variant on the prompt constrained the response to be a *single sentence* instead of free text (by providing the additional command “Provide your response in one sentence”).

Text-only models (ChatGPT, LLaMA 2) were given the prompt alone *a la* Fig. 1. LLaVA was also provided the complete image of the scene, including the player (white cube, which can be made invisible) and collectible objects (green cubes). Table 1 shows scores for LLM outputs.

Sample responses such as “Place one cube on the ground, stand on it, then stack the other cube and the sphere on top of each other to reach the top of the 2-meter high platform” (ChatGPT) and “The blue cylinder, the blue sphere, and the blue cube on the side that is one meter long” (LLaVA) demonstrate the same problems as seen in Fig. 1: a lack of consideration for the effect of object properties on environmental dynamics (such as what happens when trying to stack another object on a sphere), and also a failure to consider all information in the prompt (such as the fact that both cubes are of the same size).

Controlled Environment

The single-sentence responses from ChatGPT and LLaMA 2 resulted in somewhat better situational descriptions than free-text responses, as did removing the player object distractor in the LLaVA input, so we created a second setting where additional controls could be placed on the environment and visual input.

The second setting involved multimodal input only (hence was only evaluated on LLaVA). We evaluated both free-text responses and responses constrained to be a single sentence. We put a series of controls on the environment to: 1) Remove

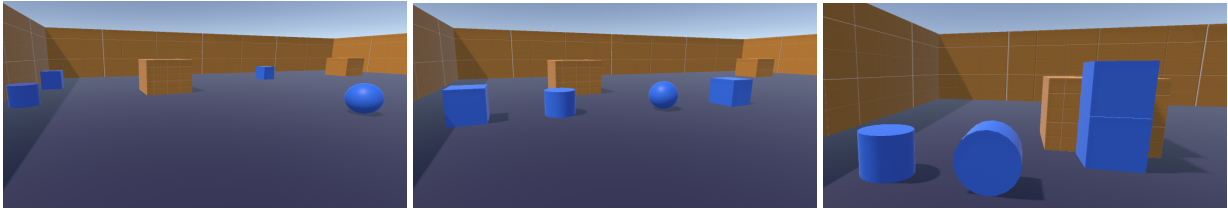


Figure 3: Examples of controls placed on the visual input to LLaVA.

	Free-text		1-sentence	
	Stable	IoU	Stable	IoU
ChatGPT	.33	.50	.67	.75
LLaMA 2-7B	.00	.50	.50	.50
LLaMA 2-70B	.25	.75	1.00 [†]	.33
LLaVA	.50	.75	.50	.75

No player				
	Free-text		1-sentence	
	Stable	IoU	Stable	IoU
LLaVA	.25	.75	.67	.75

Table 1: LLM scores from open-world simulation. The indicated LLaMA 2 models are versions optimized for dialogue (LLaMA 2-chat-*). *No player* denotes variant image inputs (hence only provided to LLaVA), in which the player was not rendered. [†]LLaMA 2-chat-70B’s response in this condition only mentioned one object, which remains stable by default in the simulation.

the collectible objects so that no distractor objects appeared in the visual input (Fig. 3[L]); 2) Place objects directly in the field of view to minimize detection errors due to distance or occlusion (Fig. 3[C]); 3) Provide a partial solution to the problem (such as partial object placements and actions taken up until the last step) and ask the model to supply the objects that the final action should be taken with (Fig. 3[R]). In this setting the prompt also included further controls and specificity, such as giving the correct number of objects in the instructions. Table 2 shows the scores from results in this setting.

Sample responses, such as “*The two cubes and the sphere can be stacked to reach the top of the platform*” (free-text response to the *No distractors* condition) and “*The two cubes and the sphere*” (single-sentence response to the *All objects in FoV* condition) exhibit many of the same problems as above, and additionally appear to be biased toward trying to use every object in the scene, whether or not that object is actually useful for the task (particularly in the free-text condition). In only one case, where the scene was preset to have completed all but the final step, the choice was between two instances of cylinders in different orientations, and the response was constrained to a single sentence, did the LLaVA model generate an output that correctly solved the task (see Fig. 3[R]), to which LLaVA’s response was “*The cylinder on*

the left should be placed in front of the stack of cubes.”).

Visual Grounding with Cross-Attention and BLIP

A weakness of the LLaVA model is the lack of a cross-modal attention encoder. Instead, LLaVA uses a linear layer to transform visual token embeddings to word embeddings. In contrast, BLIP (Li et al. 2022) uses a more sophisticated cross-modal component to account for the interaction of text and images, which motivates our use of this model. BLIP is a transformer-based multimodal encoder-decoder, including an image-grounded text encoder which uses cross-attention in every other layer to ground individual words of a text caption to regions of an image. Unlike LLaMA 2/LLaVA/ChatGPT, BLIP does not output long free-text responses to open-ended questions, so it is unsuitable for generating outputs that can be evaluated by operationalizing them in the simulation. However, it is jointly pretrained with three loss functions: image-text contrastive learning, image-text matching, and image-conditioned language modeling, which provides a more sophisticated way of seeing where visual attention is being applied when conditioned on language *input* than LLaVA’s linear projection. In this section we use BLIP as a representative example of an architecture that models the interaction of images and language through cross-attention to examine if and where cross-attention would better capture correct portions of the image containing terms key to physical reasoning problems such as ours. We examine this with a curated set of inputs, and find that in many cases such a model still fails to ground properties/concepts such as flatness/roundness given an image. Examples follow.²

In Fig. 4, the cylinder is on the far left and the sphere is on the far right of the source image. When the word “cylinder” is given, the most heavily weighted object ends up actually being the sphere on the right. This may point to a bias toward round objects like spheres in the object localization components, and perspective distortion effects on the appearance of the cylinder may make it appear more like another object type, such as a cube. When the word “cube” is given, a significant portion of the weight falls on two objects: the cubical platform that is part of the environment, and the cylinder.

In Fig. 5 we see that despite there being two cylinders in the scene, the word “cylinder” is strongly grounded to the upright cylinder (resting on its flat side), even when the text prompt mentions the cylinder on its *round* side. In fact, the model applies more cross-modal attention to the upright cylinder when the word “round” is given than when the word

²The figures omit visual grounding of articles like “the” or “a.”

Free-text						1-sentence					
No distractors		All objects in FoV		Partial solution		No distractors		All objects in FoV		Partial solution	
Stable	IoU	Stable	IoU	Stable	IoU	Stable	IoU	Stable	IoU	Stable	IoU
1.00 [†]	.75	.25	.75	.67	1.00	.67	.75	.67	.75	1.00	1.00

Table 2: LLaVA scores from controlled environment. [†]While technically none of the objects moved from their final positions after physics was applied, LLaVA’s response required placing 2 cubes *on top* of the platform and executing a series of jumps with the other objects without stating how they should be placed. This makes this solution invalid.

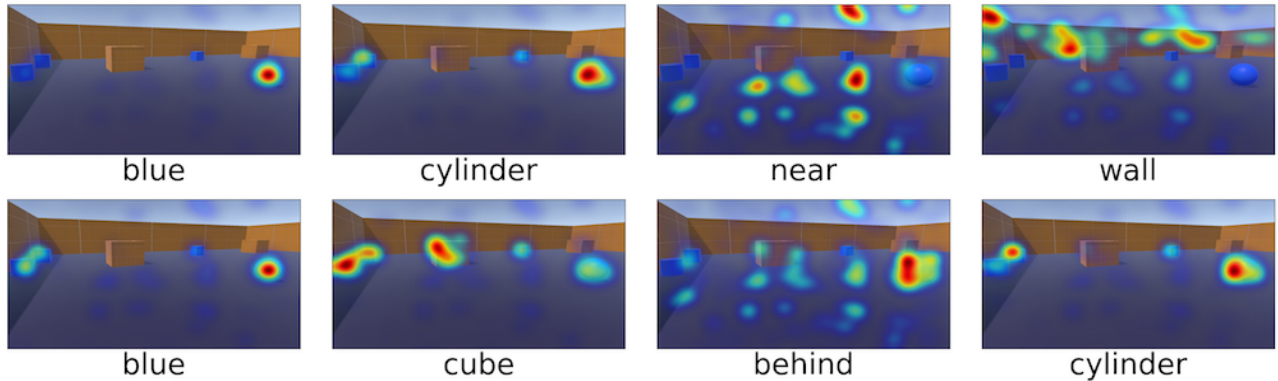


Figure 4: Per-word BLIP visual grounding for “blue cylinder near wall” (top) and “blue cube behind cylinder” (bottom).

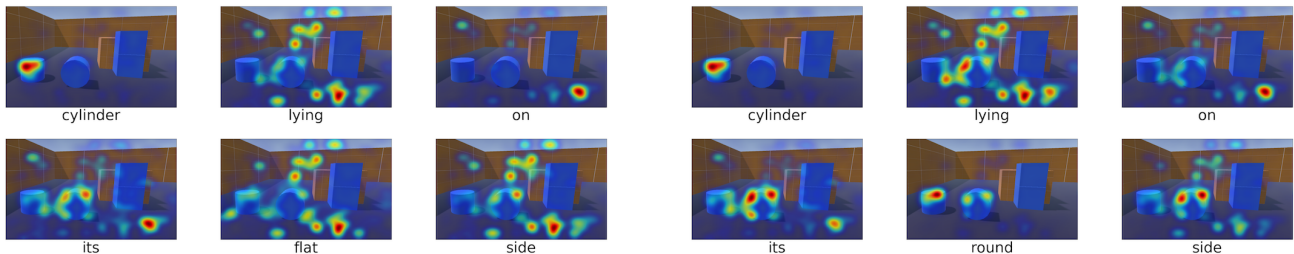


Figure 5: Per-word BLIP visual grounding for “cylinder lying on its flat side” [L] and “cylinder lying on its round side” [R].

“flat” is given. Likewise, in Fig. 7, the highlighted region falls only on the upright cylinder and not the other cylinder. This is likely because cylinders lying on the round side are sparsely or not present in image/text pairs during pre-training. As Barbu et al. (2019) observed, common image datasets often contain biases toward objects in canonical positions. For cylindrical objects, this includes things like soda cans or glasses in the canonical upright orientation.

To examine the extent to which concept words such as *flat* and *round* are visually grounded in a zero-shot manner, we consider various captions (for different images and points of view in our environment) containing these words. Fig. 6 further illustrates that concepts such as “round” and “flat,” critical to our scenario, are not really grounded in this model. The sphere is highlighted for both, indicating that “round” is successfully localized to the sphere likely because of a significant presence in the training data. Despite the contrastive loss used for training BLIP, the antonym “flat” is not localized to a non-round object.

Taken together, these results suggest that even a pre-trained vision-language model that contains even a sophisticated cross-attention architecture suffers from problems in grounding descriptive terms to the correct instances of objects, unless supplemented with other, non-linguistic, non-visual information.

Exploiting Plan Failure

Given the results above, it is clear that both text and multimodal language models struggle with correctly modeling the environmental dynamics even in simple situations like ours. This appears to primarily be a shortcoming of using text alone or only text and images as the model inputs, when other modalities besides these are truly informative about object properties and behaviors.

As humans develop object concept representations, they are also learning to individuate objects from the perceptual flow not just based on visual features but also based on experience that includes interacting with them in real

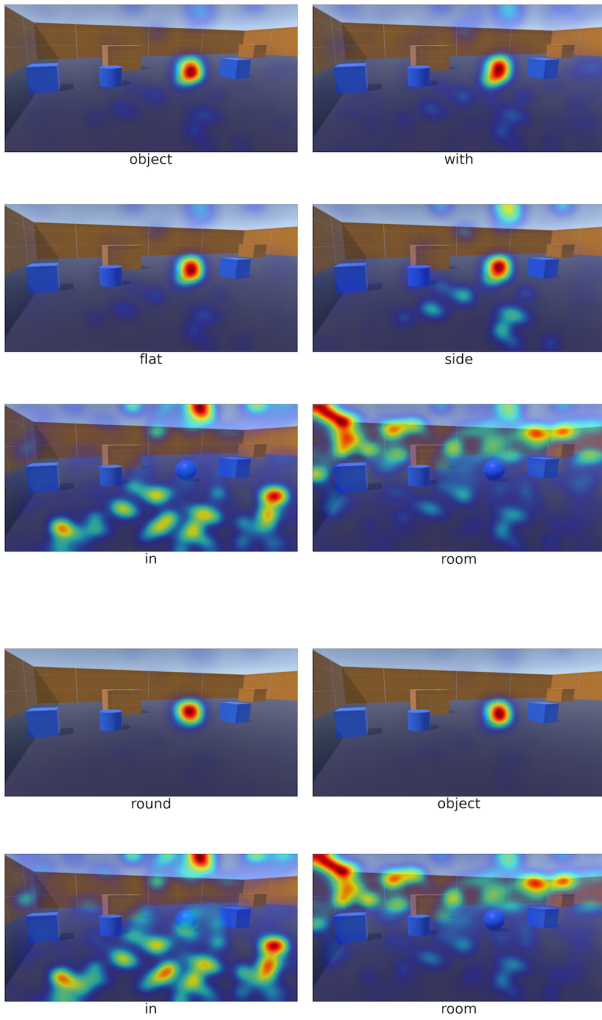


Figure 6: Per-word BLIP visual grounding for “*object with flat side in room*” (top) and “*round object in room*” (bottom).

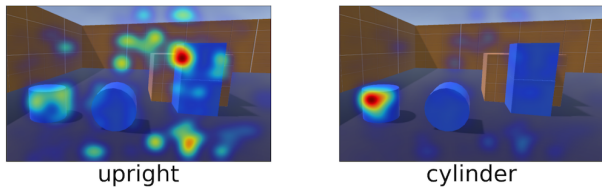


Figure 7: Per-word BLIP visual grounding for “*upright cylinder*”.

time (Spelke 1985, 1990; Baillargeon 1987; Spelke, von Hofsten, and Kestenbaum 1989). We therefore developed a procedure through which an agent embodied in a scene performs a semi-self-guided exploration of objects it encounters and uses that to solve the problem by linking object classifications in the scene to a background knowledge base.

Our knowledge base in this case is a library of *voxemes* in the VoxML modeling language (Pustejovsky and Krishnaswamy 2016). VoxML models object relevant properties

like symmetry, habitats (conditioning environments, as in Pustejovsky (2013)), and affordances, which makes it useful for specifying realistic object behavior in a continuous simulation environment, but our method is in principle friendly to an arbitrary background knowledge base as in Nirenburg, Krishnaswamy, and McShane (2023).

Within the simulation, as objects are interacted with and move, they leave trajectory traces through the environment, governed by the underlying VoxML semantics in interaction with the physics engine. This information includes position, rotation, and movement in response to various actions. For instance, this would take a cylinder’s major axis of symmetry from VoxML (the *Y* axis), and make the cylinder roll along that axis if rolling is an available affordance given the cylinder’s current habitat.

In our procedure, we initially naively follow the step-by-step plan generated by an LLM. At some point, proceeding further in the specified plan becomes impossible. For example, if trying to build the configurations shown in Fig. 1, the agent cannot balance or stack another object on the sphere, which triggers an exploration process.

As our principle concerns in this environment involve successfully stacking objects so they can either support other objects or an agent, we build our exploration strategy on a previous model from Ghaffari and Krishnaswamy (2023). This underlying model uses data gathered in two-object stacking task to classify a set of objects (the same 3 object types as used above, as well as 6 others, which display contrasts between the flatness and roundness, depending on the orientation of the objects). It uses similarity learning to create two high-level clusters in the model’s latent space: one that contains the *flat* objects, or objects that largely remain stable when stacked on a flat surface, and the *round* objects, which do not. Importantly, the model is trained only over objects that are flat on all sides (cube, pyramid, etc.), or entirely round (sphere, egg, etc.). Objects like the cylinder, which display both qualities depending on orientation, are not included in training. This means that the process of grounding the cylinder to the appropriate region involves determining the orientation in which that object is considered *flat* vs. *round*, rather than simply classifying the object.

When the failure of the initial LLM-generated plan triggers exploration, the agent traverses the environment, selects objects, and attempts to stack each one on top of itself. Because the agent is represented as a cube, the exploration process replicates the original stacking task the model was trained on. The agent has a goal to seek out objects that can be grounded to the *flat* region in its underlying model.

As the cylinder is unseen by the underlying model, the agent makes its judgments through analogizing the behavior of samples of the new object to samples of previously-encountered objects using vector similarity in the latent space. It then indexes the conditions under which those behaviors occur, which links the raw trajectory features to the embedding vector and its nearest neighbors. A schematic overview of this procedure is shown in Fig. 9.

In our method, the agent samples an object, determines the configuration, or habitat, the object requires to satisfy the stackability property, and places it, resulting in a sequence of

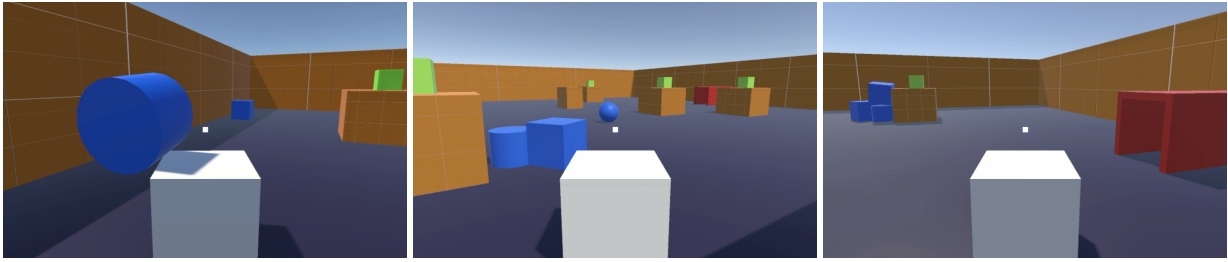


Figure 8: Agent executing plan with our method. L: Agent samples cylinder. C: Agent places cylinder flat side down, then samples cube and places it beside cylinder. R: Agent samples cube and places it on top of cylinder.

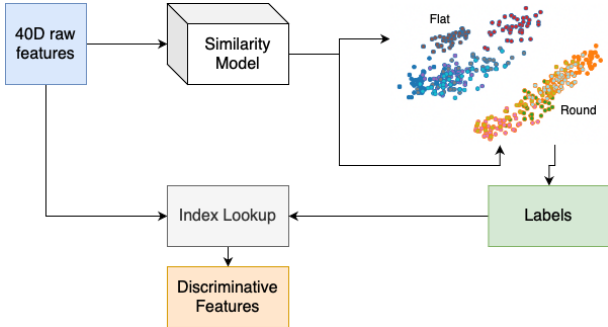


Figure 9: Schematic overview of exploration procedure.

events like that shown in Fig. 8. The agent first samples the cylinder and determines how to place it so it can participate in a successful stacking relation (flat side down). It then samples a cube and places it beside the cylinder. It then samples another cube and places it on top of the previously-placed cylinder, creating an approximation of a 2-step staircase that can be used to climb the platform.

A Proposal for Grounding Language Models to Physical Laws

We have demonstrated that current LLMs struggle with tasks of this nature, and have proposed a method to use exploration of objects under interaction to achieve a successful solution in this class of problems. The technical challenge is to retain LLMs’ open-domain, grammatical generation capabilities while accounting for environmental dynamics and object properties *in situ*.

We frame this as a knowledge distillation problem (Hinton, Vinyals, and Dean 2015) with the added challenge of needing to transfer knowledge into an LLM from a model specialized to object classification over images and trajectory data. Even a relatively small LLM like LLaVA-7B is still significantly larger than a candidate object classification model like those we have used in this line of research so far, and a standard soft logit distribution over the object class labels is not likely to provide sufficient information for a generative LLM with thousands of potential token outputs.

In addition, there is a substantial distribution mismatch between the autoregressive language model and the models in which we are encoding the information we want to

source, such as relations between objects, their properties, and their afforded behaviors. This information has been bootstrapped through a knowledge base or modeling language like VoxML, and needs to be first vectorized into a subsymbolic form, and those vectors need to be aligned with the distribution of the generative LLM.

To address this challenge, we propose leveraging the information that can be directly extracted from a simulation environment like ours. We know the location and extents of different objects in space, and the camera position from which images are captured. All these are expressed in Cartesian coordinates and a quaternion for rotations, compressed into a single 4×4 transformation matrix. Therefore, we propose projecting the object locations from 3D world space into pixel space to target patches in the images where we know attention should be paid if the correct object is to be extracted from the image. The spatial trajectory data, their localization in pixel space (e.g., as bounding boxes), along with the images, would be passed into a transformer encoder. Self-attention would be trained to detect objects in the image, with an additional object localization signal from the bounding box with the object label attached.

Attention to the correct region of the image features should also condition tokens in the outputs that describe object-relevant properties or actions. In Ghaffari and Krishnaswamy (2023) we already showed that by grounding object terms from a language model to objects from a trajectory-based classifier, we get information about related terms “for free” (e.g., grounding terms for flat objects to the flat cluster also grounded terms like “stack”).

We propose to encode the raw object movement and/or visual features in a self-supervised fashion through an attention encoder, to tease out correlations in object features relevant to object identification. For example, the major axis of a cylinder, i.e., as encoded in VoxML, is a strong correlate to its “stackability” or flat surface in one orientation, and its lack thereof in another. Since this information is implicitly encoded in the object trajectory features from the simulation, we want to enable the language model to learn correlations between them and task-relevant tokens.

Supervision can then be sourced from attention heads over the object encodings to those of the language model. Given the previously-uncovered correlations between object and behavior terms in language models, an attention loss (Eq. 1, for up to h attention heads, where L denotes the language

model and O denotes the object model) should optimize the model to apply attention from a cylinder representation to related action and property tokens, like “stack” if the cylinder is vertical, or “roll” if it is horizontal. a_i denotes the i -th attention head. Obj denotes the representation obtained by language model (L) or object model (O). This loss function minimizes the Euclidean distance between attention from the object and language models.

$$\mathcal{L}_{att} = \sum_{i=1}^h \|Obj_{a_i}^L - Obj_{a_i}^O\|_2^2 \quad (1)$$

An embedding loss (Eq. 2) then would minimize the Euclidean distance between object embeddings drawn from the final hidden state (s) of the object classifier and those of object denoting tokens. Because the two models are trained over different data with different initializations and regimes, the resulting embedding spaces are not directly comparable, so a learned linear matrix, denoted $W_{V \rightarrow L}$, projects the visual embeddings into the language model’s space.

$$\mathcal{L}_{emb} = \|Obj_s^L - Obj_s^V W_{V \rightarrow L}\|_2^2 \quad (2)$$

Fig. 10 shows a design schematic of the proposed model.

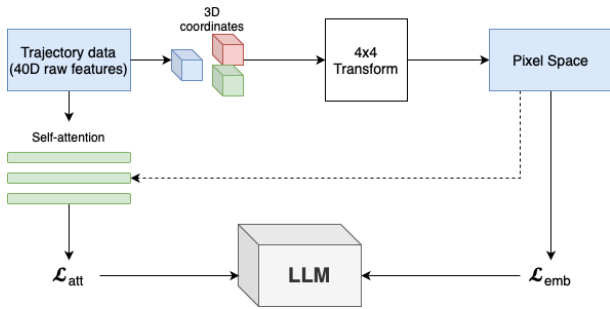


Figure 10: Proposed architecture for distilling object property and motion information to LLM representations.

Within the LLM, a contrastive method *a la* Yang et al. (2023) provides a method for generating both positive and negative samples of LLM responses for scoring and training a preference model, which addresses a profound challenge in standard reinforcement learning from human feedback (RLHF) approaches (Christiano et al. 2017): sourcing a sufficient distribution of scored responses is expensive and time consuming. However, here we are only concerning ourselves with generating physically sensible and plausible responses that account for object properties and the operation of environmental dynamics over them. Therefore, we propose to source the scoring directly from the simulation itself. For example, if we are concerned with optimizing the LLM toward generating correct responses of physical reasoning about properties relevant to object stacking, then a “good” response to prompts of this nature would be those that, when operationalized, generate stable configurations. “Bad” responses would be those that either cannot be logically satisfied (e.g., due to placing multiple objects in the same place), or results in unstable configurations. This would enable us to rapidly source a preference signal in a task-focused fashion.

The preference model is then trained using the sum of a contrastive loss, the attention loss, and the embedding loss with weighting terms Λ applied to each function and tuned over a validation set to minimize the average loss. Eq. 3 provides the combined loss function.

$$\mathcal{L} = \lambda_1 \mathcal{L}_{cont} + \lambda_2 \mathcal{L}_{att} + \lambda_3 \mathcal{L}_{emb} \quad (3)$$

The outputs of the preference model flow into the RL training of the LLM, using existing methods such as proximal policy optimization (PPO).

Conclusion

Despite the problems in LLMs’ reasoning that we have demonstrated, and that others have also studied in different domains, the way that LLMs are discussed in both scientific literature and popular media (Shanahan 2022) creates an impression of human-like abilities. Sometimes these apparent abilities do manifest in certain domains and conditions, but this should not be taken as an indicator of general human-like reasoning abilities, especially in the very physically- and situationally-grounded type of problems we examine.

In this paper we presented evidence of the difficulty LLMs have in physical reasoning problems, by mediating the LLMs’s solutions through a simulation environment to examine the effects of physics and environmental dynamics on the presented solutions. Problems manifested particularly when executing a multi-step plan where completing future steps depends on successfully exploiting the relationships created between object during past steps. We focused on multimodal models, specifically LLaVA due to its open-weight nature, and compared model outputs from this and text-only models, in different conditions, including allowing free-form or restricted responses, and putting controls on the visual input from the environment to clearly select for the information needed to solve the problem without distractors. Only in one condition were we able to coax the correct solution from the LLM: when all but the last step of the presented problem were already given, and the LLM simply had to choose the object to be used in the final step. The results point to a weakness in LLMs at reasoning *causally*, such as successfully predicting what will happen to an object configuration due to the application of consistent environmental physics after it is created. We also investigated BLIP and identified cases where it still cannot ground concepts and properties inherent to the objects in our environment in a zero-shot manner despite large-scale image/text pretraining.

This led to the second part of our contribution: an alternative method for determining the right object to be used for the task through interaction and exploration, and a proposal to use the information extracted from such explorations in the LLM to direct outputs toward better (in this case, more physically feasible) solutions by distilling grounding signals from the simulation environment back into the model. This amounts to a process of making what is implicit in linguistic input more explicit (Krishnaswamy 2017), and opens up many more avenues in future work toward grounding LLMs in a realistic understanding of causality and natural laws.

References

- Ahn, M.; Brohan, A.; Brown, N.; Chebotar, Y.; Cortes, O.; David, B.; Finn, C.; Fu, C.; Gopalakrishnan, K.; Hausman, K.; et al. 2022. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*.
- Asher, N.; Bhar, S.; Chaturvedi, A.; Hunter, J.; and Paul, S. 2023. Limits for Learning with Language Models. *arXiv preprint arXiv:2306.12213*.
- Baillargeon, R. 1987. Object permanence in $3\frac{1}{2}$ - and $4\frac{1}{2}$ -month-old infants. *Developmental psychology*, 23(5): 655.
- Barbu, A.; Mayo, D.; Alverio, J.; Luo, W.; Wang, C.; Gutfreund, D.; Tenenbaum, J.; and Katz, B. 2019. Objectnet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. *Advances in neural information processing systems*, 32.
- Chaturvedi, A.; Bhar, S.; Saha, S.; Garain, U.; and Asher, N. 2022. Analyzing Semantic Faithfulness of Language Models via Input Intervention on Conversational Question Answering. *arXiv preprint arXiv:2212.10696*.
- Chen, L.; Zaharia, M.; and Zou, J. 2023. How is ChatGPT’s behavior changing over time? *arXiv preprint arXiv:2307.09009*.
- Christiano, P. F.; Leike, J.; Brown, T.; Martic, M.; Legg, S.; and Amodei, D. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Burstein, J.; Doran, C.; and Solorio, T., eds., *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 4171–4186. Minneapolis, Minnesota: Association for Computational Linguistics.
- Driess, D.; Xia, F.; Sajjadi, M. S.; Lynch, C.; Chowdhery, A.; Ichter, B.; Wahid, A.; Tompson, J.; Vuong, Q.; Yu, T.; et al. 2023. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*.
- Ghaffari, S.; and Krishnaswamy, N. 2022. Detecting and accommodating novel types and concepts in an embodied simulation environment. In *Proceedings of the 10th Annual Conference on Advanced in Cognitive Systems*.
- Ghaffari, S.; and Krishnaswamy, N. 2023. Grounding and distinguishing conceptual vocabulary through similarity learning in embodied simulations. In *Proceedings of the 15th International Conference on Computational Semantics*.
- Goertzel, B. 2023. Generative AI vs. AGI: The Cognitive Strengths and Weaknesses of Modern LLMs. *arXiv preprint arXiv:2309.10371*.
- Hinton, G.; Vinyals, O.; and Dean, J. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- Kocoń, J.; Cichecki, I.; Kaszyca, O.; Kochanek, M.; Szydło, D.; Baran, J.; Bielaniec, J.; Gruza, M.; Janz, A.; Kancierz, K.; et al. 2023. ChatGPT: Jack of all trades, master of none. *Information Fusion*, 101861.
- Krishnaswamy, N. 2017. *Monte Carlo Simulation Generation Through Operationalization of Spatial Primitives*. Brandeis University.
- Krishnaswamy, N.; Pickard, W.; Cates, B.; Blanchard, N.; and Pustejovsky, J. 2022. The VoxWorld platform for multimodal embodied agents. In *LREC proceedings*, volume 13.
- Krishnaswamy, N.; and Pustejovsky, J. 2022. Affordance embeddings for situated language understanding. *Frontiers in Artificial Intelligence*, 5: 774752.
- Lerer, A.; Gross, S.; and Fergus, R. 2016. Learning physical intuition of block towers by example. In *International conference on machine learning*, 430–438. PMLR.
- Li, J.; Li, D.; Xiong, C.; and Hoi, S. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, 12888–12900. PMLR.
- Li, R.; Jabri, A.; Darrell, T.; and Agrawal, P. 2020. Towards practical multi-object manipulation using relational reinforcement learning. In *2020 IEEE international conference on robotics and automation (ICRA)*, 4051–4058. IEEE.
- Liu, H.; Li, C.; Wu, Q.; and Lee, Y. J. 2023. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*.
- Lobbezoo, A.; Qian, Y.; and Kwon, H.-J. 2021. Reinforcement learning for pick and place operations in robotics: A survey. *Robotics*, 10(3): 105.
- Nirenburg, S.; Krishnaswamy, N.; and McShane, M. 2023. Hybrid Machine Learning/Knowledge Base Systems Learning through Natural Language Dialogue with Deep Learning Models. In *AAAI Spring Symposium: Challenges Requiring the Combination of Machine Learning and Knowledge Engineering*.
- Polverini, G.; and Gregorcic, B. 2023. How understanding large language models can inform the use of ChatGPT in physics education. *European Journal of Physics*.
- Pustejovsky, J. 2013. Dynamic event structure and habitat theory. In *Proceedings of the 6th International Conference on Generative Approaches to the Lexicon (GL2013)*, 1–10.
- Pustejovsky, J.; and Krishnaswamy, N. 2016. VoxML: A Visualization Modeling Language. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, 4606–4613.
- Qin, C.; Zhang, A.; Zhang, Z.; Chen, J.; Yasunaga, M.; and Yang, D. 2023. Is ChatGPT a general-purpose natural language processing task solver? *arXiv preprint arXiv:2302.06476*.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PMLR.
- Shakarian, P.; Koyyalamudi, A.; Ngu, N.; and Mareedu, L. 2023. An Independent Evaluation of ChatGPT on Mathematical Word Problems (MWP). *arXiv preprint arXiv:2302.13814*.
- Shanahan, M. 2022. Talking about large language models. *arXiv preprint arXiv:2212.03551*.

- Spelke, E. S. 1985. Perception of unity, persistence, and identity: Thoughts on infants' conceptions of objects.
- Spelke, E. S. 1990. Principles of object perception. *Cognitive science*, 14(1): 29–56.
- Spelke, E. S.; von Hofsten, C.; and Kestenbaum, R. 1989. Object perception in infancy: Interaction of spatial and kinetic information for object boundaries. *Developmental Psychology*, 25(2): 185.
- Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Wang, Y.; Duan, J.; Fox, D.; and Srinivasa, S. 2023. NEWTON: Are Large Language Models Capable of Physical Reasoning? In *Findings of the Association for Computational Linguistics: EMNLP 2023*, 9743–9758.
- Yang, K.; Klein, D.; Celikyilmaz, A.; Peng, N.; and Tian, Y. 2023. Rlcd: Reinforcement learning from contrast distillation for language model alignment. *arXiv preprint arXiv:2307.12950*.