

Answering visual what-if questions

Wagner, Misha; Basevi, Hector; Shetty, Rakshith ; Li, Wenbin; Malinowski, Mateusz; Fritz, Mario; Leonardis, Ales

DOI:

[10.1007/978-3-030-11009-3_32](https://doi.org/10.1007/978-3-030-11009-3_32)

License:

None: All rights reserved

Document Version

Peer reviewed version

Citation for published version (Harvard):

Wagner, M, Basevi, H, Shetty, R, Li, W, Malinowski, M, Fritz, M & Leonardis, A 2019, Answering visual what-if questions: from actions to predicted scene descriptions. in L Leal-Taixé & S Roth (eds), Computer Vision – ECCV 2018 Workshops: Munich, Germany, September 8-14, 2018, Proceedings. Lecture Notes in Computer Science, vol. 11129, Springer, pp. 521-537, Visual Learning and Embodied Agents in Simulation Environment Workshop at 15th European Conference on Computer Vision (ECCV 2018), Munich, Germany, 9/09/18.
https://doi.org/10.1007/978-3-030-11009-3_32

[Link to publication on Research at Birmingham portal](#)

Publisher Rights Statement:

Checked for eligibility: 21/01/2019

This is a post-peer-review, pre-copyedit version of an article published in Lecture Notes in Computer Science. The final authenticated version is available online at: https://doi.org/10.1007/978-3-030-11009-3_32

General rights

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

Take down policy

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact UBIRA@lists.bham.ac.uk providing details and we will remove access to the work immediately and investigate.

Answering Visual *What-If* Questions: From Actions to Predicted Scene Descriptions

Misha Wagner^{1*}, Hector Basevi^{1*}, Rakshith Shetty², Wenbin Li²,
Mateusz Malinowski², Mario Fritz³, and Aleš Leonardis¹

¹ University of Birmingham, Birmingham, UK

² Max Planck Institute for Informatics, Saarland Informatics Campus, Saarbrücken,
Germany

³ CISP A Helmholtz Center i.G., Saarland Informatics Campus, Saarbrücken,
Germany

Abstract. In-depth scene descriptions and question answering tasks have greatly increased the scope of today’s definition of scene understanding. While such tasks are in principle open ended, current formulations primarily focus on describing only the current state of the scenes under consideration. In contrast, in this paper, we focus on the future states of the scenes which are also conditioned on actions. We posit this as a question answering task, where an answer has to be given about a future scene state, given observations of the current scene, and a question that includes a hypothetical action. Our solution is a hybrid model which integrates a physics engine into a question answering architecture in order to anticipate future scene states resulting from object-object interactions caused by an action. We demonstrate first results on this challenging new problem and compare to baselines, where we outperform fully data-driven end-to-end learning approaches.

Keywords: Scene understanding, Visual Turing Test, Visual question answering, Intuitive physics

1 Introduction

While traditional scene understanding involves deriving bottom-up scene representations such as object bounding boxes and segmentation, in recent years alternative approaches such as *scene captioning* and *question answering* have become increasingly popular. These do not strive for a particular type of representation of the input scene, but rather formulate an alternative task that requires a more holistic scene understanding. Such approaches have been very successful and have shown great advances in extracting the semantic scene content by deriving captions and answers about diverse scene elements.

Beyond the estimation of the “status quo” of a visual scene, recent deep learning approaches have shown improved capabilities of forecasting scenes into the future. This is particularly useful for autonomous agents (e.g., robots or

* These authors contributed equally to this work.

driving assistants) that have to plan ahead and act safely in dynamically changing environments. Recent approaches show extrapolation of complete videos [1], edge information [2] or object trajectories [3, 4]. However, with increasing time horizons and complexity of the scenes, such quantitative predictions become increasingly difficult. In addition, extrapolation of complete image data might be wasteful and overly difficult to achieve.

Furthermore, current work on anticipation and forecasting is typically not interactive, meaning that the agent is acting purely as a passive observer. However, in many real-world applications, an agent is faced with the task of evaluating multiple different potential actions that will cause diverse outcomes. The future is therefore often conditioned on the actions of the agent which is not handled by the state-of-the-art methods.

Therefore, we argue for a qualitative prediction of the future conditioned on an action. We phrase this as the *Answering Visual What-If Questions* task, where the answer is conditioned on an observation and question including a hypothetical action. This formulation allows us to evaluate a model’s forecasting abilities conditioned on a hypothetical action, and at the same time allows for sparse representation of the future where not all details of the scene or object interactions have to be fully modeled.

We provide the first investigation of this challenging problem in a table top scenario where a set of objects is placed in a challenging configuration and different actions can be taken. Dependent on the action, the objects will interact according to the physics of the scene and will cause a certain outcome in terms of object trajectories. The task is to describe the outcome with respect to the action. In order to address this problem we couple the question answering approach with a physics engine – resulting in a hybrid model. While several parts of the method, such as inferring the action from the question and predicting the answer, are data-driven, the core aspect of our model that predicts future outcomes is model-based (using a physics engine).

The contributions of this paper are as follows:

- We define a new challenge called Visual *What-If* Question answering (WIQ) that brings together question answering with anticipation.
- We propose the first dataset for the WIQ task based on table-top interactions between multiple objects called TIWIQ.
- We propose the first hybrid model that uses a physics engine together with a question answering architecture.

2 Related Work

Learning Physics and Future Predictions: Coping with the physical world by predicting how objects interact with each other using rules of physics is among the pillars of human intelligence. This type of intuitive understanding of physics, often referred to as “intuitive physics” [5], is also becoming of interest to

machine learning researchers. The “NeuroAnimator” is among the first learning-based architectures trained to simulate physical dynamics based on observations of physics-based models [6]. Although the “NeuroAnimator” is mainly motivated by efficiency, others have realized that learning-based architectures may be key components to learn the whole spectrum of physical reasoning that humans possess. For instance, [7] argues that a cognitive mechanism responsible for physical reasoning may resemble an engine that simulates complex interactions between different physical objects, and can be implemented by “graph neural networks” [8, 9] or by an engineered physics engine [10]. In this work, our hybrid model also uses a physics engine, but unlike [10] we are less interested in inferring latent physics properties of two objects from videos, but rather in a forward model of physics for the purpose of answering questions. A complementary line of research has shown that convolutional neural networks (CNN) are capable to some extent of physical reasoning such as stability prediction [3, 11, 12], or future frame synthesis from RGB-D input [13, 1, 2] or even static images [14, 4]. These approaches to physical intelligence focus on testing this understanding either by trying to extrapolate sensory data into the future (predicting video frames) or by inferring individual properties (predicting stability or physics properties). In contrast, we propose to achieve qualitative physical understanding, where we want the model to have general understanding of physical processes but not necessarily the ability to make precise prediction.

Visual Question Answering (Visual QA): This is a recently introduced research area [15, 16] that attempts to build and understand if machines can learn to explain the surrounding environment only based on questions and answers. Since then, the community has seen a proliferation of various datasets [17–23], including the most popular VQA [24], as well as numerous methods [15, 18, 25–29]. Although most of the questions involve static scenes, and are either related to objects, attributes, or activities, there are some that require understanding of physics at the “intuitive level”. Consider even such seemingly simple question as “What is on the table?”. To interpret this question, understanding of “on” is needed, and this involves physical forces such as gravity. In spite of the existence of such questions in the aforementioned datasets, due to lack of interactions, it is hardly possible the learnt models can really understand them, and likely they only rely on visual correlations. In our work, through the interactions, and exploitation of physics, we can train architectures that, we hypothesize, can model physical interactions between objects.

Simulations and Machine Learning: Since it is difficult to generate realistic data that includes complex physical interactions, most approaches either rely on short videos with limited exposition to physics [13, 1, 2] or on synthetically generated data [8, 9, 12, 30]. This problem of lacking good realistic environments with rich physical interactions also governs the research on reinforcement learning [31], where the community often relies on game-like environments [32–35]. Since there is no publicly available realistic environment that has rich enough physical interactions that we are interested in, we build a dataset consisting of 3D scenes, with physical interactions, and with realistically textured objects.

3 Visual *What-If* Questions (WIQ) Task

While Visual QA evaluates the scene understanding of a passive agent, this is not sufficient for an active agent that needs to anticipate the consequences of its actions and communicate about them. To study this aspect of scene understanding, we propose the task of answering “what-if” questions pertaining to a visual scene. The agent is shown an input scene with multiple objects and is given a hypothetical action description. It then has to describe what happens to different objects in the scene, given that it performs the specified action. This amounts to answering questions of the form “If I perform action A , what happens to object X ?”. To answer such questions the agent has to parse the natural language description of the action to infer the action type and target object on which the action is applied, along with the corresponding parameters such as the initial force. Then the agent needs to anticipate the consequences of this action on different objects in the scene, and finally verbalize its predictions. This combines the challenges involved in the standard VQA task [24] with intuitive physics [8] and the future state anticipation tasks [36].

3.1 Table-top Interaction Visual *What-If* Questions (TIWIQ) Dataset

Existing Visual QA datasets [15, 18, 24] focus on static scenes, whereas datasets commonly used in future prediction tasks such as CityScapes [37] involve a passive observer (future states are not conditioned on the agent’s action). Since we are interested in the question answering task involving “physical intelligence”, we collect a new table-top interaction visual *what-if* questions (TIWIQ) dataset. This dataset has synthetic table-top scenes, with pairs of action descriptions and ground-truth descriptions of the outcomes of the specified action. We stick to synthetic scenes and a physics simulation engine to build this dataset as it provides physics, and enables controlled experimentation.

Scenes: To obtain the TIWIQ dataset we instantiate random table-top scenes in a physics engine, simulate actions on these scenes and collect human annotations describing the actions and the consequences. Each training sample in the TIWIQ dataset contains a table-top scene with five objects, each randomly placed upon the table. The five objects are chosen from eight items from the YCB Object Dataset [38]: a foam brick, a cheez-it box, a chocolate pudding box, a mustard bottle, a banana, a softball, a ground coffee can, and a screwdriver.

Actions: A random action is chosen to be performed on a single random object, simulated using the Bullet 3 [39] physics engine. The resulting trajectories are rendered into a video of the interactions. The actions can be one of four: 1. Push an object in a specific direction. 2. Rotate an object clockwise or anti-clockwise. 3. Remove an object from the scene. 4. Drop an object on another object.

Annotation: The objects shown in rendered videos have colored outlines, and when questions are posed to annotators, objects are referred to by their outline color rather than their name. This avoids the questions biasing the annotator’s vocabulary with regard to object names.

Human Baseline: We have also collected a human performance benchmark on the visual what-if question answering task on the TIWIQ dataset. To obtain the human performance baseline, the annotators were shown a static image of the scene and a description of the action to be performed and were asked to describe what happens to different objects in the scene. We compare the performance of the model proposed in section 4 to this human performance benchmark.

Dataset Statistics: We have generated and annotated 15 batches of data. Each batch has 17 examples of each action, totaling 68 examples per batch. In total, we have 1020 annotated examples. Three batches, totaling 204 examples and 20% of the dataset, are dedicated to testing. For each scene, there are four generated descriptions (one for each object that is not being acted on), therefore there are 4080 ($1020 * 4$) annotated descriptions. However, descriptions relate to movement or interactions between objects only around 25% of the time. This is due to the random placement of objects sometimes resulting in scenes with spatially separated objects, and therefore some actions having no impact on most objects in a scene. This results in approximately 1000 movement and interaction descriptions across the dataset. Only these annotations are used to train the description generation model.

Vocabulary Statistics: The vocabulary of the dataset is explored by counting the number of unique words used across the dataset (1-gram), as well as the number of unique n-grams for values 2 to 5 (2,3,4,5-grams). This is shown in Table 1. These statistics are reported for the action description annotations, the action effect annotations, and the two together. It is worth noting that the vocabulary of the action description dataset is significantly smaller than the vocabulary of the effect description dataset. This is due to the range of actions being specified by the design of the scenario, while the range of effects has no such constraints.

Table 1. The size of the vocabulary for the action and effect descriptions and the whole TIWIQ dataset, including the average sentence length and the number of unique n-grams in each subset of the dataset.

	Descriptions	Length	1-gram	2-grams	3-grams	4-grams	5-grams
Action	9.63	107	323	565	757	867	
Effect	7.663	110	403	724	981	1,075	
All	8.582	152	619	1,171	1,653	1,895	

4 Our Model

Recent advances in Deep Learning architectures have dominated Visual QA and image captioning tasks. The dominant approach is to use end-to-end trainable neural networks which take the inputs, e.g. image and question, and predict

the answer. The whole model is learned purely from the data. Driven by sizable datasets, they have outperformed previous purely model-based approaches e.g. relying on semantic parsing or rule-based systems [15, 25]. Although latest work has also shown early success at applying this end-to-end machine learning paradigm to predicting future states of a scene, the time horizon remains limited to a few frames and rich object interactions are typically not considered, or the scenes are greatly simplified [2, 9, 11]. Therefore, we argue for a hybrid model. We use a physics engine as a backbone in order to simulate outcomes of hypothetical actions, but we embed the simulation into learning-based components that drive the simulation as well as interpret its outcome in terms of a natural language output.

4.1 Model Overview

The proposed hybrid question answering model consists of three distinct components as shown in Fig. 1. There are two inputs to the whole model. The first is a list of object types and their initial pose (position in 3-dimensional space, and a 3x3 rotation matrix) in the scene. We always assume the same table position for every case. The second input is the action description. This was provided by human annotators, and describes some action performed on one of the objects in the scene, for example “The robot pushes the mustard container to the left”.

Both inputs are used by a “parser” (we use a neural network as the parser) to extract parameters of the action to be performed. This includes parsing the action type, object to be acted upon and parameters of the action. This extracted information serves as an input to a physics engine, which then simulates the parsed actions on the input scene to produce trajectories for each object in the scene. While these trajectories encode everything that happened in the simulation, they are not human readable. The description model takes these trajectories as input and produces a natural language summary of the state of each object under the influence of this action. The action parser model and description models are comprised of neural networks and their parameters are learned from the data. The physics engine is model driven and has no trainable parameters. In the following subsections we discuss each of these components and how they interact in more details.

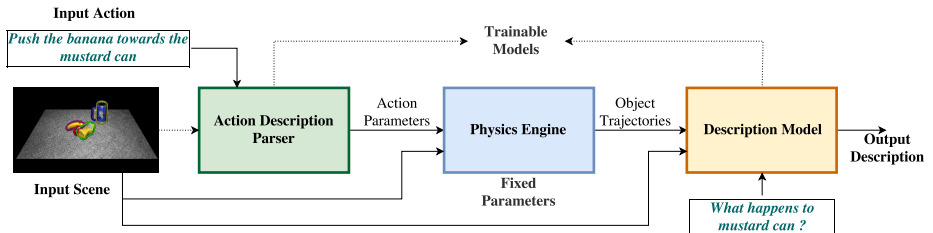


Fig. 1. Overall architecture of the proposed hybrid QA model.

As well as being described in this document, all models are illustrated accurately in the corresponding figures. Each component in the illustrations describes a single layer, whether it is an RNN or a fully connected layer. All details of the layers are given in the supplementary material. This includes layer sizes, dropout, and activation functions.

4.2 Action Description Parser

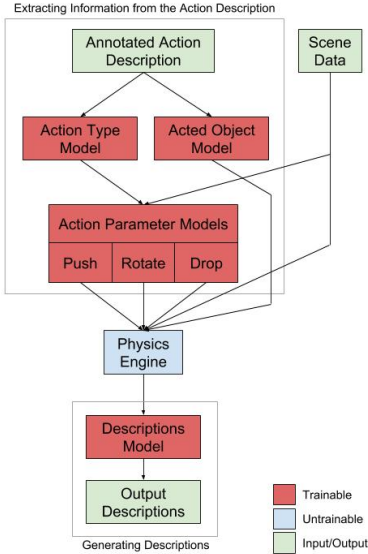


Fig. 2. Illustration of the interaction of subcomponents of the action description parser when inferring parameters of the action to be performed from input sentence.

of *Action Parameter Models*, which infer the exact parameters of the action depending on the action type. Depending on the inferred action type, one of four things happens. If the action type is a push, rotate, or drop action, then the corresponding parameter model is called with the action description and the input scene. If it is a remove action, no parameters need to be inferred as the object is simply removed from the scene.

There are parameter models for three of the four actions: push, rotate, and drop. Each of these models use recurrent networks for embedding the action description specific to the action type.

The *Push Parameter Model* infers the direction of the push by outputting a (x, y) push direction vector in its final layer. The activation for this layer is sigmoid in order to cap both components from -1 to 1. When the physics engine

The first step in the pipeline is parsing the exact nature of the action to be performed from the input sentence description. This model forms part of Fig. 2. It consists of three components with a total of five neural networks, shown in Fig. 3. First component is the *Action Type Model* which infers the type of the action described in the input (push, rotate, drop or remove). This is a recurrent neural network (RNN) that embeds the tokenized action description, iterates over it using a long short-term memory (LSTM) model, and puts the final output of the LSTM through a fully connected layer with softmax activation. These outputs are treated as the probability that each action type was described in the action description.

The second component is the *Acted Object Model*, which predicts the object in the input scene on which the described action is to be performed. The structure is identical to the *Action Type Model*, with the exception that it outputs probabilities of each class being the object to act upon.

Finally, the third component is a set

simulates this push direction, the (x, y) components are converted into an angle, removing the magnitude of the push.

The *Rotate Parameter Model* is a binary classifier which predicts whether the rotation is clockwise or anti-clockwise, using softmax activation for classification.

The *Drop Parameter Model* outputs a classification of which other object the acted object is dropped on. This also has a softmax activation for the classification, running over all possible objects.

Each component model that has text as input requires the text to be embedded. When it provided an improvement in performance, GloVe pre-trained word embeddings[40] were used. Details on which layers used pre-trained embeddings, including the size of the embeddings, is given in the supplementary material.

4.3 Physics Simulation

The Acted Object Model extracts the action type, the object of interest, and the parameters of the action. We use Bullet 3 [39] as a physics engine, with object representations from the YCB Object Dataset [38]. We use object convex decompositions for collision detection which are calculated using VHACD, an algorithm included in the Bullet 3 source code. Pushes and rotations are implemented via impulses.

The physics engine is initialized with initial object poses. The engine is run for one second of simulation time in order for objects to settle in.

The inferred action is then performed on the inferred object, and the simulation is run for a total of five seconds at a sampling rate of 300Hz. Trajectories for each object are extracted from the simulation as a list of translation and rotation pairs, where the translation is a point in 3-dimensional space and the rotation is represented by a 3x3 rotation matrix.

We then run a simple algorithm to check if an object was affected by the action. To do this, we look at the trajectory for a single object, normalize the pose using a standard deviation and mean estimated from the entire training data set, and then calculate the standard deviation of both the translation and rotation, resulting in two floating point values. We say that an object was affected by the action if either of these values exceed a certain threshold. These thresholds were calculated by running a grid

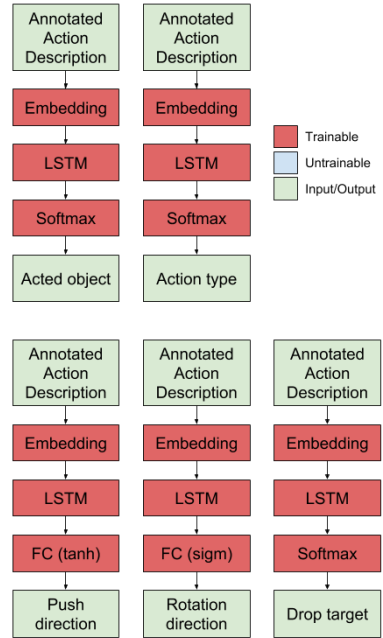


Fig. 3. Illustrations of the action parsing models. Between them, only the final fully connected (FC) layers differ.

search over possible threshold values for either value, and picking the pair that resulted in the best classification accuracy on the training set.

4.4 Generating Descriptions

The *Description Model* shown in Fig. 4 uses the trajectories from the physics simulation to produce a one sentence summary of the effect of the action on each object in the scene. This model is run independently for each generated description with the following inputs: 1. The action description. 2. The object class to describe. 3. The trajectory of the object to describe. 4. A list of other object classes in the scene. 5. A list of other object trajectories in the scene.

The object classes are encoded as a one-hot vector, and the trajectories are encoded as a list of points in 3-dimensional space.

At training time, ground truth trajectories are used, but when the complete hybrid model is being evaluated, predicted trajectories are generated via physics simulation.

The description model works in two stages. First the input trajectories of the target object (whose state is being described) and other objects in the scene are compared and trajectory embeddings are obtained. Then these trajectory embeddings and action description embeddings are input to a decoder LSTM language model, which generates the final description.

To obtain the trajectory embeddings we iterate over each of the other objects in the scene — that is, the ones that are not currently being described. For each object, we compute the difference between its trajectory and the trajectory of the object to be described, at each time step. These difference vectors are then embedded and iterated over using an LSTM. The initial state of this LSTM is provided by embedding both of the object classes and putting them through a fully connected layer. The final hidden state of this LSTM should encode the interactions between the two objects. This output trajectory embedding is concatenated with the object encodings of the two relevant objects. We find that including these embeddings after as well as before the trajectory encoding LSTM improves the overall model’s performance.

The input action description is encoded using an LSTM (as in earlier models, such as the action description model). A fully connected layer is used to transform the concatenated trajectory embedding vector and the encoded input instruction into the right dimensions and is used to initialize the hidden state of the decoder LSTM. The input for the decoder LSTM at time t_0 is the *start of sentence* token, and the input at time t_i is the output from t_{i-1} . At each step the decoder LSTM outputs the next word and this repeats until the *end of sentence* token is predicted. This process is carried out to generate a description for each object in the scene.

4.5 Implementation Details

We have implemented the hybrid model and all components in Python using the Keras [41] library with the TensorFlow [42] backend. For the description model, custom layers were introduced into Keras using TensorFlow. Overall runtime of the system is 1.76s, where prediction time of the Action Description Parser and Description Generation is negligible. Almost all time is spent in the simulation part. For reproducibility and to stimulate more work on this challenge, we will release code, models and setup.

5 Results

We evaluate our overall hybrid approach as well as the individual components on the proposed dataset as well as compare to an end-to-end learning and human baseline. We provide example results and analyses that highlight the benefits of our approach.

5.1 Performance of Hybrid Model Components

We separately evaluate the performance of the six components of the hybrid model, using ground truth annotations at these intermediate stages. First, we show the performance of the action description information extraction models (action type model, acted object model, push / rotate / drop parameter model) in table 2. We created simple support vector machine (SVM) baselines in order to benchmark the more powerful neural models. The input to these SVMs is a vector of word counts for every word in the data vocabulary. We find that in all cases bar one, the neural models significantly outperform SVMs, as shown in Table 2. The exception is the rotation parameter model, which is outperformed by 5.8%. The performance for the rotation parameter model is particularly poor due to noisy annotations in the cases of rotation actions. Through looking at a small subset of the rotation action annotations, we have found that 30-40% of the annotations are mislabeled in some way — either giving the wrong rotational direction, or annotated as a push action.

To compare the push parameter model with a classification network, we discretize the angle inferred by the neural model into eight directions (e.g. left, top-left, up). The SVM also classifies to one of those eight directions, allowing us to compare the performance of these two models.

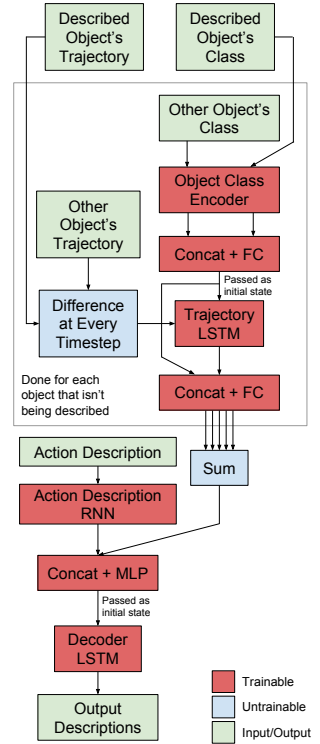


Fig. 4. An illustration of the description model. The section outlined in a gray square is run for every object that isn't the object that is currently being described.

5.2 Quantifying the Hybrid Model Performance

We will now quantify the performance of the proposed hybrid model and baselines on the test set.

Metrics: To measure the description performance we use the standard metrics used in evaluating image captioning models such as BLEU, CIDEr, ROUGE and a custom metric COM (“Correct Objects Mentioned”) that we designed for this specific problem. This metric searches the descriptions for different object names, creating a list of objects mentioned in the text. This is done for both the prediction and the ground truth. The COM metric is computed as the intersection over union of these two sets. The upper-bound for COM is 1, and occurs when all correct objects are mentioned in all predictions for a scene. Image captioning metrics such as BLEU focus on overall n-gram matching of the generated description with the ground truth, regardless of the importance of each word, whereas COM directly measures how well models identify the acting objects in a scene.

Hybrid Model Compared to Baselines: We compare the performance of the hybrid model against three other baselines on the test set. The first is a pure data-driven model, an end-to-end trainable neural network illustrated in figure 5. The inputs to this network are the input action description, the initial scene, and the object to describe. The action description is embedded and then run through an LSTM and the final output of this LSTM is taken. Each class and pose in the initial scene is flattened into a vector, and each of these is put through a fully connected layer and summed together. The object to describe is encoded as a one-hot vector and passed through a fully connected layer. Each of these encodings are concatenated together and treated as input for the decoder LSTM, which generates a description for the specified object class.

Human Baseline and “Upper Bound”: The second set of descriptions is from a human baseline, mentioned in section 3.1. Human annotators were shown the input scene and action description, but not the video of the action taking place. They were asked to describe what happens to each object. This simulates the same task tackled by the hybrid, and pure data-driven models. Finally, the third baseline is obtained by feeding the ground-truth trajectories to the description model. This represents an upper-bound on the hybrid model performance.

Discussion: We find that the hybrid model outperforms the data-driven model in all metrics, with an increase of 15.4% in the BLEU metric, and an increase of 20.5% in the COM metric, as illustrated in table 3. This provides evidence for incorporating a physics engine for solving physics based anticipation problems over a pure data-driven approach. The performance of the hybrid model is close to its upper-bound description model and this gap comes from the cases where

Table 2. Comparing classification accuracy of neural network based and SVM based models on different tasks.

Task	NN	SVM
Action Type	97.5%	97%
Acted Object	94%	90%
Push Parameters	90%	44%
Rotation Parameters	68%	72%
Drop Parameters	90%	36%

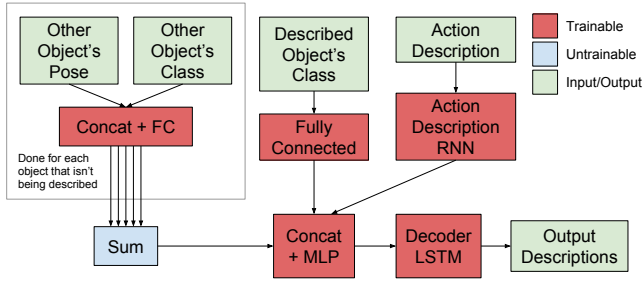


Fig. 5. Illustration of the data-driven model.

the action parsing model failed. However, there is still a gap in performance in terms of the COM metric between the proposed hybrid model and the human benchmark, indicating the scope for future improvements.

Human Baseline Discussion: Under most metrics, the hybrid model outperforms the human baseline. However, this is misleading: the human baseline contains high-quality annotations, and under domain-specific metrics such as COM, is evaluated with a near-perfect score (0.953). Its large error in BLUE, CIDEr, and ROUGE results from the differing vocabularies between the human baseline and the ground truth. For this reason, comparison between the hybrid model and the human baseline is difficult to achieve using these metrics; a similar problem is common in the image captioning domain.

Table 3. Comparison of description generating models. Best values between the hybrid model and data-driven model are highlighted. This shows that the hybrid model exceeds the data-driven model and even the human baseline in some metrics.

Model	BLEU	BLEU-1	BLEU-2	BLEU-3	BLEU-4	CIDEr	ROUGE	COM
Description Model	0.280	0.439	0.345	0.203	0.133	1.118	0.421	0.671
Human Baseline	0.191	0.207	0.192	0.186	0.181	1.849	0.209	0.953
Hybrid Model	0.262	0.407	0.322	0.184	0.134	1.118	0.396	0.640
Data-Driven Model	0.227	0.375	0.282	0.154	0.099	0.896	0.376	0.531

5.3 Qualitative analysis

We provide qualitative examples and analysis in table 5. In these examples, the hybrid model can be seen generating more specific and accurate descriptions of the results of actions compared to the data-driven model. There are three main failure cases of the data-driven model.

The first of these is shown in row 1 of table 5. In this example, the hybrid model correctly predicts the object which hits the foam (in this case the screw-driver) while the data-driven approach predicts that a different object in the scene will hit the foam. Accuracy is lost here due to the data-driven model not being able to reliably infer the object that interacted with the subject object. Our hybrid model performs better in this case, presumably because it was able to use the trajectories from the physics engine to infer the correct object.

The second main failure case is shown in row 2. Both models are correct but the hybrid model gives a more precise description, stating correctly which object hit the mustard container. The data-driven model gives a more vague description by not stating the acting object and just describing the movement.

The third failure case is shown in row 3. Often, the data-driven model produces a description where both the object being acted on and the object affecting it are the same. This could be due to the data-driven model making the best guess it can — if it knows that the class “screw driver” appears in the text but does not know what the other object could be, and it knows that the sentence should reference two objects, then it may choose to mention “screw driver” twice in the sentence. This failure case, although more prevalent in the data-driven model, shows up in the hybrid model too as seen in row 4 of table 5.

There is a failure case unique to the hybrid model. The data-driven model was trained only on cases where the action did have an effect on the object. However, the hybrid model has to infer whether there was an effect. This results in some cases where the hybrid model misclassifies the object as “not moving” and generates the “nothing” description. An example of this case is shown in row 5 of table 5.

5.4 Ablation Analysis

We also analyze the error introduced by the different components within the hybrid model. We do this by introducing, one-by-one, the ground truth values for a particular component instead of the predicted values. The results of this are shown in table 4. We can see that introducing the ground-truth for whether an object moved provides the biggest increase in performance, implying that the hybrid model loses a lot of accuracy when predicting whether an object moved. Conversely, we can also see that the Action Type and Acted Object models introduce relatively small amounts of error, suggesting they correctly model the ground truth.

6 Conclusion

We have proposed a new task that combines scene understanding with anticipation of future scene states. We argue that this type of “physical intelligence” is a key competence of an agent that is interacting with an environment and tries to evaluate different alternatives. In contrast to prior work on quantitative predictions of future states, here, we focus on a qualitative prediction that describes

Table 4. Comparison of how the performance of the hybrid model improves when cumulatively adding truth values for each of the components.

Model	BLEU	CIDEr	ROUGE	COM
All Predictions:	0.262	1.118	0.396	0.640
With True Action Type:	0.264	1.126	0.398	0.644
...and True Acted Object:	0.265	1.129	0.400	0.646
...and True Action Parameters:	0.272	1.153	0.406	0.662
...and True Trajectories:	0.268	1.086	0.401	0.645
...and True Object Acted On:	0.283	1.133	0.424	0.673

the outcome for a certain object with a natural language utterance. Owing to such a formulation, we can train and evaluate our agent on long-term future anticipation, where the model can easily ignore irrelevant details of the scene or the interactions. This contrasts with future frame synthesis where all the details have to be correctly modeled.

Due to the lack of suitable datasets, we introduced the first dataset and an evaluation protocol for this challenging task. Our proposed model is the first that combines a question answering architecture with a physics engine and verbalizes different outcomes dependent on the visual and language input. In particular, our hybrid model outperforms a purely data-driven Deep Learning baseline. We believe that such hybrid models that combine a complex simulation engine with data-driven components represent an exciting avenue for further research as they allow for generalization, scalability and improved performance in challenging scenarios with a high combinatorial complexity.

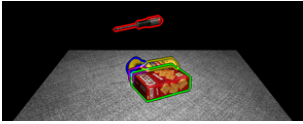
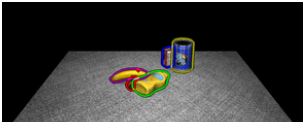

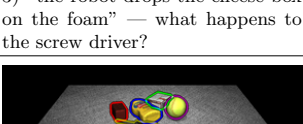
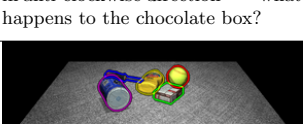
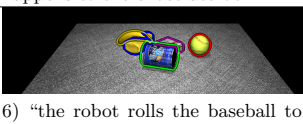
Acknowledgements

We acknowledge MoD/Dstl and EPSRC for providing the grant to support the UK academics involvement in a Department of Defense funded MURI project through EPSRC grant EP/N019415/1.

References

1. Mathieu, M., Couprie, C., LeCun, Y.: Deep multi-scale video prediction beyond mean square error. In: International Conference on Learning Representations. (2016)
2. Bhattacharyya, A., Malinowski, M., Schiele, B., Fritz, M.: Long-term image boundary extrapolation. In: Association for the Advancement of Artificial Intelligence. (2018)
3. Lerer, A., Gross, S., Fergus, R.: Learning physical intuition of block towers by example. In: International conference on machine learning. (2016)
4. Mottaghi, R., Rastegari, M., Gupta, A., Farhadi, A.: “What happens if...” Learning to predict the effect of forces in images. In: European Conference on Computer Vision, Springer (2016) 269–285

Table 5. Examples of different scenes in the data set with annotations from ground truth, human baseline, the predictions from the hybrid model, and predictions from the baseline data-driven model. The hybrid can be seen giving more precise descriptions than the data-driven model.

Input Scene And Action	Output Description
 <p>1) "the robot drops the screw driver on the foam" — what happens to the foam?</p>	<p>Ground Truth "the foam is pushed a little by the screw driver"</p> <p>Human Baseline "the foam is pushed because the screw driver drops on it"</p> <p>Hybrid Model Prediction "the foam is pushed by the screw driver"</p> <p>Data-Driven Model Prediction "the foam is pushed by the mustard container"</p>
 <p>2) "the robot spins the screw driver in anti-clockwise direction" — what happens to the mustard container?</p>	<p>Ground Truth "the screw driver pushes the mustard container"</p> <p>Human Baseline "the mustard container moves a little due to the impact of spinning screw driver"</p> <p>Hybrid Model Prediction "the screw driver pushes the mustard container"</p> <p>Data-Driven Prediction "the mustard container shakes a little from the impact"</p>
 <p>3) "the robot drops the cheese box on the foam" — what happens to the screw driver?</p>	<p>Ground Truth "the screw driver is pushed by the cheese box"</p> <p>Human Baseline "the screw driver is pushed by the cheese box"</p> <p>Hybrid Model Prediction "the screw driver is pushed by the cheese box"</p> <p>Data-Driven Prediction "the screw driver is pushed by the screw driver"</p>
 <p>4) "the robot spins the screw driver in anti-clockwise direction" — what happens to the chocolate box?</p>	<p>Ground Truth "the chocolate box is pushed by the screw driver"</p> <p>Human Baseline "the chocolate box is pushed a slightly by the screw driver"</p> <p>Hybrid Model Prediction "the chocolate box pushes the chocolate box"</p> <p>Data-Driven Prediction "the chocolate box is pushed by the screw driver"</p>
 <p>5) "the robot pushes the baseball to the middle of the table" — what happens to the chocolate box?</p>	<p>Ground Truth "the chocolate box is pushed by the baseball"</p> <p>Human Baseline "the chocolate box is pushed by the baseball"</p> <p>Hybrid Model Prediction "nothing"</p> <p>Data-Driven Prediction "the chocolate box is pushed by the chocolate box"</p>
 <p>6) "the robot rolls the baseball to the north-west side of the table and it drops off" — what happens to the banana?</p>	<p>Ground Truth "nothing"</p> <p>Human Baseline "nothing"</p> <p>Hybrid Model Prediction "nothing"</p> <p>Data-Driven Prediction N/A</p>

5. McCloskey, M.: Intuitive physics. *Scientific American* **248**(4) (1983) 122–131
6. Grzeszczuk, R., Terzopoulos, D., Hinton, G.: Neuroanimator: Fast neural network emulation and control of physics-based models. In: *Proceedings of the 25th annual conference on Computer graphics and interactive techniques*, ACM (1998) 9–20
7. Battaglia, P.W., Hamrick, J.B., Tenenbaum, J.B.: Simulation as an engine of physical scene understanding. *Proceedings of the National Academy of Sciences* **110**(45) (2013) 18327–18332
8. Battaglia, P., Pascanu, R., Lai, M., Rezende, D.J., et al.: Interaction networks for learning about objects, relations and physics. In: *Advances in neural information processing systems*. (2016) 4502–4510
9. Watters, N., Tacchetti, A., Weber, T., Pascanu, R., Battaglia, P., Zoran, D.: Visual interaction networks. In: *Advances in neural information processing systems*. (2017)
10. Wu, J., Yildirim, I., Lim, J.J., Freeman, B., Tenenbaum, J.: Galileo: Perceiving physical object properties by integrating a physics engine with deep learning. In: *Advances in neural information processing systems*. (2015) 127–135
11. Li, W., Leonardis, A., Fritz, M.: Visual stability prediction for robotic manipulation. *Proceedings of the IEEE International Conference on Robotics and Automation* (2017)
12. Li, W., Azimi, S., Leonardis, A., Fritz, M.: To fall or not to fall: A visual approach to physical stability prediction. *CoRR* **abs/1604.00066** (2016)
13. Ranzato, M., Szlam, A., Bruna, J., Mathieu, M., Collobert, R., Chopra, S.: Video (language) modeling: a baseline for generative models of natural videos. *CoRR* **abs/1412.6604** (2014)
14. Mottaghi, R., Bagherinezhad, H., Rastegari, M., Farhadi, A.: Newtonian scene understanding: Unfolding the dynamics of objects in static images. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2016) 3521–3529
15. Malinowski, M., Fritz, M.: A multi-world approach to question answering about real-world scenes based on uncertain input. In: *Advances in neural information processing systems*. (2014) 1682–1690
16. Malinowski, M., Fritz, M.: Towards a visual turing challenge. *CoRR* **abs/1410.8027** (2014)
17. Geman, D., Geman, S., Hallonquist, N., Younes, L.: Visual Turing test for computer vision systems. *Proceedings of the National Academy of Sciences* **112**(12) (2015) 3618–3623
18. Ren, M., Kiros, R., Zemel, R.: Exploring models and data for image question answering. In: *Advances in neural information processing systems*. (2015)
19. Yu, L., Park, E., Berg, A.C., Berg, T.L.: Visual madlibs: Fill in the blank description generation and question answering. In: *Proceedings of the IEEE International Conference on Computer Vision*, IEEE (2015) 2461–2469
20. Zhu, Y., Groth, O., Bernstein, M., Fei-Fei, L.: Visual7W: Grounded question answering in images. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2016) 4995–5004
21. Tapaswi, M., Zhu, Y., Stiefelhagen, R., Torralba, A., Urtasun, R., Fidler, S.: Movieqa: Understanding stories in movies through question-answering. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2016) 4631–4640
22. Agrawal, A., Batra, D., Parikh, D., Kembhavi, A.: Don’t just assume; look and answer: Overcoming priors for visual question answering. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2018)

23. Kaffe, K., Cohen, S., Price, B., Kanan, C.: Dvqa: Understanding data visualizations via question answering. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2018)
24. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Lawrence Zitnick, C., Parikh, D.: VQA: Visual question answering. In: Proceedings of the IEEE International Conference on Computer Vision. (2015) 2425–2433
25. Malinowski, M., Rohrbach, M., Fritz, M.: Ask your neurons: A deep learning approach to visual question answering. *International Journal of Computer Vision* **125**(1-3) (2017) 110–135
26. Yang, Z., He, X., Gao, J., Deng, L., Smola, A.: Stacked attention networks for image question answering. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2016) 21–29
27. Fukui, A., Park, D.H., Yang, D., Rohrbach, A., Darrell, T., Rohrbach, M.: Multi-modal compact bilinear pooling for visual question answering and visual grounding. *CoRR* **abs/1606.01847** (2016)
28. Hu, R., Andreas, J., Rohrbach, M., Darrell, T., Saenko, K.: Learning to reason: End-to-end module networks for visual question answering. In: Proceedings of the IEEE International Conference on Computer Vision. (2017)
29. Santoro, A., Raposo, D., Barrett, D.G., Malinowski, M., Pascanu, R., Battaglia, P., Lillicrap, T.: A simple neural network module for relational reasoning. In: Advances in neural information processing systems. (2017) 4974–4983
30. Ehrhardt, S., Monszpart, A., Mitra, N.J., Vedaldi, A.: Taking visual motion prediction to new heightfields. *CoRR* **abs/1712.09448** (2017)
31. Sutton, R.S., Barto, A.G.: Reinforcement learning: An introduction. Volume 1. MIT press Cambridge (1998)
32. Beattie, C., Leibo, J.Z., Teplyashin, D., Ward, T., Wainwright, M., Küttler, H., Lefrancq, A., Green, S., Valdés, V., Sadik, A., Schrittwieser, J., Anderson, K., York, S., Cant, M., Cain, A., Bolton, A., Gaffney, S., King, H., Hassabis, D., Legg, S., Petersen, S.: Deepmind lab. *CoRR* **abs/1612.03801** (2016)
33. Kempka, M., Wydmuch, M., Runc, G., Toczek, J., Jaśkowski, W.: Vizdoom: A doom-based ai research platform for visual reinforcement learning. In: IEEE Conference on Computational Intelligence and Games, IEEE (2016) 1–8
34. Shah, S., Dey, D., Lovett, C., Kapoor, A.: Airsim: High-fidelity visual and physical simulation for autonomous vehicles. In: Field and service robotics, Springer (2018) 621–635
35. Wu, Y., Wu, Y., Gkioxari, G., Tian, Y.: Building generalizable agents with a realistic and rich 3d environment. *CoRR* **abs/1801.02209** (2018)
36. Finn, C., Goodfellow, I., Levine, S.: Unsupervised learning for physical interaction through video prediction. In: Advances in neural information processing systems. (2016) 64–72
37. Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2016)
38. Çalli, B., Walsman, A., Singh, A., Srinivasa, S., Abbeel, P., Dollar, A.M.: Benchmarking in manipulation research: The YCB object and model set and benchmarking protocols. *CoRR* **abs/1502.03143** (2015)
39. Coumans, E.: Bullet 3. <https://github.com/bulletphysics/bullet3> (2018)
40. Pennington, J., Socher, R., Manning, C.D.: Glove: Global vectors for word representation. In: Empirical Methods in Natural Language Processing (EMNLP). (2014) 1532–1543

41. Chollet, F., et al.: Keras. <https://github.com/keras-team/keras> (2015)
42. Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G.S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., Zheng, X.: TensorFlow: Large-scale machine learning on heterogeneous systems (2015) Software available from [tensorflow.org](https://www.tensorflow.org).