

A visual Turing test for computer vision systems

Donald Geman^{*}, Stuart Geman[†], Neil Hallonquist^{*} and Laurent Younes^{*}

^{*} Johns Hopkins University, and [†] Brown University

Submitted to Proceedings of the National Academy of Sciences of the United States of America

Today, computer vision systems are tested by their accuracy in detecting and localizing instances of objects. As an alternative, and motivated by the ability of humans to provide far richer descriptions, even tell a story about an image, we construct a “visual Turing test”: an operator-assisted device that produces a stochastic sequence of binary questions from a given test image. The query engine proposes a question; the operator either provides the correct answer or rejects the question as ambiguous; the engine proposes the next question (“just-in-time truing”). The test is then administered to the computer-vision system, one question at a time. After the system’s answer is recorded, the system is provided the correct answer and the next question. Parsing is trivial and deterministic; the system being tested requires no natural language processing. The query engine employs statistical constraints, learned from a training set, to produce questions with essentially unpredictable answers—the answer to a question, given the history of questions and their correct answers, is nearly equally likely to be positive or negative. In this sense, the test is only about vision. The system is designed to produce streams of questions that follow natural story lines, from the instantiation of a unique object, through an exploration of its properties, and onto its relationships with other uniquely instantiated objects.

scene interpretation | computer vision | Turing test | binary questions | unpredictable answers | sparse learning

Going back at least to the mid-twentieth century there has been an active debate about the state of progress in artificial intelligence and how to measure it. Alan Turing [1] proposed that the ultimate test of whether or not a machine could “think,” or think at least as well as a person, was for a human judge to be unable to tell which was which based on natural language conversations in an appropriately cloaked scenario. In a much-discussed variation (sometimes called the “standard interpretation”), the objective is to measure how well a computer can *imitate* a human [2] in some circumscribed task normally associated with intelligent behavior, although the practical utility of “imitation” as a criterion for performance has also been questioned [3]. In fact, the overwhelming focus of the modern AI community has been to assess machine performance more directly by dedicated tests for specific tasks rather than debating about general “thinking” or Turing-like competitions between people and machines.

In this paper we implement a new, query-based test for computer vision, one of the most vibrant areas of modern AI research. Throughout this paper we use “computer vision” more or less synonymously with semantic image interpretation - “images to words.” But of course computer vision encompasses a great many other activities; it includes the theory and practice of image formation (“sensors to images”); image processing (“images to images”); mathematical representations; video processing; metric scene reconstruction; and so forth. In fact, it may not be possible to interpret scenes at a semantic level without taking at least some of these areas into account, especially the geometric relationship between an image and the underlying 3D scene. But our focus is how to evaluate a system, not how to build one.

Besides successful commercial and industrial applications, such as face detectors in digital cameras and flaw detection in manufacturing, there has also been considerable progress in more generic tasks, such as detecting and localizing instances from multiple generic object classes in ordinary indoor and outdoor scenes, in “fine-grained” classification such as identifying plant and animal species, and in recognizing attributes of objects and activities of people. The results

of challenges and competitions (see [4, 5]) suggest that progress has been spurred by major advances in designing more computationally efficient and invariant image representations [6, 7, 8]; in stochastic and hierarchical modeling [9, 10, 11, 12]; in discovering latent structure by training multi-layer networks with large amounts of unsupervised data [13]; and in parts-based statistical learning and modeling techniques [14, 15, 16], especially combining discriminative part detectors with simple models of arrangements of parts [17]. Quite recently, sharp improvements in detecting objects and related tasks have been made by training convolutional neural networks with very large amounts of annotated data [18, 19, 20, 21, 22].

More generally, however, machines lag very far behind humans in “understanding images” in the sense of generating rich semantic annotation. For example, systems that attempt to deal with occlusion, context and unanticipated arrangements, all of which are easily handled by people, typically encounter problems. Consequently, there is no point in designing a “competition” between computer vision and human vision: interpreting real scenes (such as the ones in Figure 1) is virtually “trivial” (at least effortless and nearly instantaneous) for people whereas building a “description machine” that annotates raw image data remains a fundamental challenge.

We seek a quantitative measure of how well a computer vision system can interpret ordinary images of natural scenes. Whereas we focus on urban street scenes, our implementation could easily be extended to other image populations and the basic logic and motivations remain the same. The “score” of our test is based on the responses of a system under evaluation to a series of binary questions about the existence of people and objects, their activities and attributes, and relationships among them, all relative to an image. We have chosen image-based rather than scene-based queries (see Scenes vs. images).

Suppose an image sub-population \mathcal{I} has been specified (“urban street scenes” in Figure 1), together with a “vocabulary” and a corresponding set of binary questions (see Vocabulary and Questions). Our prototype “visual Turing test” (VTT) is illustrated in Figure 2. Questions are posed sequentially to the computer vision system using

Significance

In computer vision, as in other fields of AI, the methods of evaluation largely define the scientific effort. Most current evaluations measure detection accuracy, emphasizing the classification of regions according to objects from a pre-defined library. But detection is not the same as understanding. We present here a different evaluation system, in which a query engine prepares a written test (“visual Turing test”) that uses binary questions to probe a system’s ability to identify attributes and relationships in addition to recognizing objects.

Reserved for Publication Footnotes



Fig. 1. Urban street scenes. First row: Athens, Baltimore, Busan, Delhi. Second row: Hong Kong, Miami, Rome, Shanghai.

a “query language” which is defined in terms of an allowable set of predicates. The interpretation of the questions is unambiguous and does not require any natural language processing. The core of the VTT is an automatic “query generator” which is learned from annotated images and produces a sequence of binary questions for any given “test” image $I_0 \in \mathcal{I}$ whose answers are “unpredictable” (see Statistical Formulation). In loose terms, this means that hearing the first $k - 1$ questions and their true answers for I_0 without actually seeing I_0 provides no information about the likely answer to the next question. In order to prepare for the test, designers of the vision systems would be provided with the database used to train the query generator as well as the full vocabulary and set of possible questions, and would have to provide an interface for answering questions. One simple measure of performance is the average number of correct responses over multiple runs with different test images.

Current Evaluation Practice

Numerous datasets have been created to benchmark performance, each designed to assess some vision task (e.g., object detection) on some image domain (e.g., street-scenes). Systems are evaluated by comparing their output on these data to “ground-truth” provided by humans. One well-studied task is classifying an entire image by a general category, either at the object level (“car,” “bike,” “horse,” etc.), where *ImageNet* [5] is a currently popular annotated dataset, or at the scene-level (“beach,” “kitchen,” “forest,” etc.); see for example the SUN dataset [23]. A natural extension of object-level image categorization is detecting and localizing all instances from generic classes in complex scenes containing multiple instances and events; localization refers to providing either a “bounding box” per instance or segmenting the object from the background. Popular datasets for this task include the *Pascal* dataset [4], the *LabelMe* dataset [24], and the *Lotus Hill* dataset [25], all populated by relatively unconstrained natural images, but varying considerably in size and in the level of annotation, ranging from a few keywords to hierarchical representations (*Lotus Hill*). Finally, a few other datasets have been assembled and annotated to evaluate the quality of detected object attributes such as color, orientation and activity; examples are the *Core* dataset [26], with annotated object parts and attributes, and the *Virat* dataset [27] for event detection in videos.

Why not continue to measure progress in more or less the same way with common datasets dedicated to sub-tasks, but using a richer vocabulary? First, as computer vision becomes more ambitious and aims at richer interpretations, it would seem sensible to fold these sub-tasks into a larger endeavor; a system which detects activities and relationships must necessarily solve basic sub-tasks anyway. Then why not simply require competing systems to submit much richer annotation for a set of test images than in previous competitions and then rank systems according to consistency with ground truth supplied by

human annotators? The reason, and the justification for the VTT, is that the current method does not scale with respect to the richness of the representation. Even for the sub-tasks in the competitions mentioned earlier, the evaluation of performance, i.e., comparing the output of the system (e.g., estimated bounding boxes) to the ground-truth is not always straightforward and the quality of matches must be assessed [28]. Moreover, annotating every image submitted for testing at massive levels of detail is not feasible. Hence, objectively scoring the veracity of annotations is not straightforward. As in school, answering specific questions is usually more objective and efficient in measuring “understanding.” Finally, some selection procedure seems unavoidable; indeed, the number of possible binary questions which are both probing and meaningful is virtually infinite. However, selecting a subset of questions (i.e., preparing a test) is not straightforward. We would argue that the only way to ask very detailed questions without having their answers be almost certainly “no” is sequential and adaptive querying—questions which build on each other to uncover semantic structure. In summary, the VTT is one way to “scale up” evaluation.

Proposed Test: Overview

Images of scenes. Our questions are image-centered, but images capture 3D scenes. Whereas we pose our questions succinctly in the form “*Is there a red car?*”, this is understood to mean “*Is there an instance of a red car in the scene partially visible in the image?*”. Similarly, given a designated rectangle of image pixels (see Figure 2 for some examples), the query “*Is there a person in the designated region?*” is understood to mean “*Is there an instance of a person in the scene partially visible in the designated image region?*”. The universal qualifier “partially visible in the image” (or in the designated region) avoids the issue of the scope of the scene and leads naturally to instantiation and story lines.

Estimating uncertainty. The justification for counting all questions the same is the property of unpredictability: at each step k , the likelihood that the true answer for question k is “yes” given the true answers to the previous $k - 1$ questions is approximately one-half. However, generating long strings of “interesting” questions and “story lines” is not straightforward due to “data-fragmentation”: a purely empirical solution based entirely on collecting relative frequencies from an annotated training subset of size n from \mathcal{I} is only feasible if the number of questions posed is approximately $\log_2 n$. Our proposed solution is presented as part of the Statistical Formulation, and in more detail in the Supplemental Information (SI) Appendix; it rests on enlarging the number of images in the dataset which satisfy a given history by making carefully chosen invariance and independence assumptions about objects and their attributes and relationships.

Human in the loop. The operator serves two crucial functions: removing ambiguous questions and providing correct answers. Given

a rich family of questions, some will surely be ambiguous for any specific test image. The solution is “just-in-time truthing”: any question posed by the query generator can be rejected by the operator, in which case the generator supplies another nearly unpredictable one, of which there are generally many. The correct answers may or may not be provided to the system under evaluation at run time. Needless to say, given the state of progress in computer vision, neither of these roles can be served by an automated system. The test can be constructed either offline or “online” (during the evaluation). In either case, the VTT is “written” rather than “oral” since the choice of questions does not depend on the responses from the system under evaluation.

Instantiation. A key mechanism for arriving at semantically interesting questions is instance “instantiation.” A series of positive answers to inquiries about attributes of an object will often imply a single instance, which can then be labeled as “instance k ”. Hence, questions which explicitly address uniqueness are also included, which usually become viable, that is close to unpredictable, after one or two attributes have been established. Once this happens, there is no ambiguity in asking whether “person 1” and “person 2” are *talking* or whether “person 1” is *occluding* “vehicle 2”; see Figure 2. We regard



- | | |
|--|---------|
| 1. Q: Is there a person in the blue region? | A: yes |
| 2. Q: Is there a unique person in the blue region?
(Label this person 1) | A: yes |
| 3. Q: Is person 1 carrying something? | A: yes |
| 4. Q: Is person 1 female? | A: yes |
| 5. Q: Is person 1 walking on a sidewalk? | A: yes |
| 6. Q: Is person 1 interacting with any other object? | A: no |
| ⋮ | |
| 9. Q: Is there a unique vehicle in the yellow region?
(Label this vehicle 1) | A: yes |
| 10. Q: Is vehicle 1 light-colored? | A: yes |
| 11. Q: Is vehicle 1 moving? | A: no |
| 12. Q: Is vehicle 1 parked and a car? | A: yes |
| ⋮ | |
| 14. Q: Does vehicle 1 have exactly one visible tire? | A: no |
| 15. Q: Is vehicle 1 interacting with any other object? | A: no |
| 17. Q: Is there a unique person in the red region? | A: no |
| 18. Q: Is there a unique person that is female in the red region? | A: no |
| 19. Q: Is there a person that is standing still in the red region? | A: yes |
| 20. Q: Is there a unique person standing still in the red region?
(Label this person 2) | A: yes |
| ⋮ | |
| 23. Q: Is person 2 interacting with any other object? | A: yes |
| 24. Q: Is person 1 taller than person 2? | A: amb. |
| 25. Q: Is person 1 closer (to the camera) than person 2? | A: no |
| 26. Q: Is there a person in the red region? | A: yes |
| 27. Q: Is there a unique person in the red region?
(Label this person 3) | A: yes |
| ⋮ | |
| 36. Q: Is there an interaction between person 2 and person 3? | A: yes |
| 37. Q: Are person 2 and person 3 talking? | A: yes |

Fig. 2. A selection of questions extracted from a much longer sequence (one of three shown in §5 of the SI Appendix). Answers, including identifying Q24 as ambiguous, are provided by the operator (see paragraph on “Human in the loop”). Localizing questions include, implicitly, the qualifier “partially visible in the designated region” and instantiation (existence and uniqueness) questions implicitly include “not previously instantiated.” The localizing windows used for each of the four instantiations (vehicle 1, person 1, person 2, and person 3) are indicated by the colored rectangles (blue—thick border, red—thin border, yellow—broken border). The colors are included in the questions for illustration. In the actual test, each question designates a single rectangle through its coordinates, so that “Is there a unique person in the blue region” would actually read “Is there a unique person in the designated region.”

instantiation as identifying the “players” in the scene, allowing for story lines to develop.

Evolving descriptions. The statistical constraints naturally impose a “coarse-to-fine” flow of information, from gist to semantic detail. Due to the unpredictability criterion, the early questions can only inquire about coarse scene properties, such as “*Is there a person in the lefthand side of the image?*” or “*Is there a person wearing a hat?*”, because only these have intermediate probabilities of occurrence in the general population. It is only after objects have been instantiated, i.e., specific instances identified, that the likelihoods of specific relationships among these “players” become appreciably greater than zero.

Vocabulary and Questions

Vocabulary. Our vocabulary \mathcal{V} consists of three components: *types* of objects, \mathcal{T} , type-dependent *attributes* of objects, $\{\mathcal{A}_t, t \in \mathcal{T}\}$, and type-dependent *relationships* between two objects, $\{\mathcal{R}_{t,t'}\}$. For example, for “urban street scenes,” some natural types (or categories) are people, vehicles, buildings, and “parts” such as windows and doors of cars and buildings. Attributes refer to object properties such as clothing and activities of people, or types and colors of vehicles. There may also be attributes based on localizing an object instance within an image, and these provide an efficient method of instantiation (see below). Relationships between two types can be either “ordered,” for instance a person *entering* a car or building, or “un-ordered,” for instance two people *walking* or *talking* together. And some relationship questions may depend on the position of the camera in the underlying 3D scene, such as asking which person or vehicle is closer to the camera. A complete list of objects, attributes, and relationships used in our prototype is included with the SI Appendix.

Questions. Each question $q \in \mathcal{Q}$ belongs to one of four categories: existence questions, $\mathcal{Q}_{\text{exist}}$, uniqueness questions, $\mathcal{Q}_{\text{uniq}}$, attribute questions, \mathcal{Q}_{att} , or relationship questions, \mathcal{Q}_{rel} . The goal of the existence and uniqueness questions is to instantiate objects, which are then labeled (“person 1,” “vehicle 3,” ...) and subsequently available, by reference to the label, in attribute and relationship questions (“Is person 1 partially occluding vehicle 3?”). Consequently, questions in \mathcal{Q}_{att} and \mathcal{Q}_{rel} refer only to previously instantiated objects. See Figure 2 for examples drawn from $\mathcal{Q}_{\text{exist}}$ (e.g., 1, 19, 26), $\mathcal{Q}_{\text{uniq}}$ (e.g., 2, 9, 17), \mathcal{Q}_{att} (e.g., 3, 10, 23), and \mathcal{Q}_{rel} (e.g., 25, 36, 37). (Summarizing, the full set of questions is $\mathcal{Q} = \mathcal{Q}_{\text{exist}} \cup \mathcal{Q}_{\text{uniq}} \cup \mathcal{Q}_{\text{att}} \cup \mathcal{Q}_{\text{rel}}$.)

As already mentioned, we use “in the designated region” as shorthand for “in the scene that is partially visible in the designated region of the image.” Similarly, so as to avoid repeated discovery of the same objects, all existence and uniqueness questions include the additional qualifier “not previously instantiated,” which is always implied rather than explicit. So “Is there a person in the designated region wearing a hat?” actually means “Is there a person in the scene partially visible in the designated region of the image, wearing a hat and not previously instantiated?”

We assume the answers are unambiguous *for humans* in nearly all cases. However, there is no need to identify *all* ambiguous questions for any image. Filtering is “as needed”: given $I_0 \in \mathcal{I}$, any question q which is elicited by the query generator but is in fact ambiguous *for* I_0 will be rejected by the human operator during the construction of the VTT. (Examples include question 24 in the partial stream shown in Figure 1 and two others in the complete streams shown in §5 of the SI Appendix.)

Statistical Formulation

Selecting questions whose answers are unpredictable is only meaningful in a statistical framework in which answers are random variables relative to an image population \mathcal{I} , which serves as the underlying sample space, together with a probability distribution P on \mathcal{I} .

Query generator. Given an image $I \in \mathcal{I}$, the query generator interacts with an oracle (human being) to produce a sequence of questions and correct answers. The human either rejects a question as ambigu-

ous or provides an answer, in which case the answer is assumed to be a (deterministic) function of I . The process is recursive: given a *history* of binary questions and their answers, $H = ((q_1, x_1), \dots, (q_k, x_k))$, $q_i \in \mathcal{Q}$ and $x_i \in \{0, 1\}$, the query generator either stops, for lack of additional unpredictable questions, or proposes a next question q , which is either rejected as ambiguous or added to the history along with its correct answer x :

$$H \rightarrow [H, (q, x)] \triangleq ((q_1, x_1), \dots, (q_k, x_k), (q, x)), \quad x \in \{0, 1\}$$

Not all sequences of questions and answers make sense. In particular, attribute and relationship questions (\mathcal{Q}_{att} and \mathcal{Q}_{rel}) always refer to previously instantiated objects, restricting the set of meaningful histories, which we shall denote by \mathbb{H} . A key property of histories $H = ((q_1, x_1), \dots, (q_k, x_k)) \in \mathbb{H}$ produced by the query generator is that each question q_i , given the history $((q_1, x_1), \dots, (q_{i-1}, x_{i-1}))$, is “unpredictable,” a concept which we will now make precise.

Given a history H , only some of the questions $q \in \mathcal{Q}$ are good candidates for followup. As already noted, references to labeled objects cannot precede the corresponding instantiation questions, and furthermore there is a general ordering to the questions designed to promote natural story lines. For a given query generator, we will write \mathcal{Q}_H to indicate the set of possible followup questions defined by these *non-statistical* constraints. Typically, \mathcal{Q}_H contains many candidates, most of which are highly predictable given the history H , and therefore unsuitable.

The set of histories, \mathbb{H} , can be viewed as a set of binary random variables: $H = H(I) = 1$ if $H = ((q_1, x_1), \dots, (q_k, x_k)) \in \mathbb{H}$ and if the sequence of questions (q_1, \dots, q_k) produces the sequence of unambiguous answers (x_1, \dots, x_k) for the image I , and $H = 0$ otherwise. We will write P_H for the *conditional* probability on \mathcal{I} given that $H(I) = 1$.

Consider now the probability under P_H that a question $q \in \mathcal{Q}_H$ elicits the (unambiguous) response $X_q = X_q(I) \in \{0, 1\}$, for a given history $H \in \mathbb{H}$:

$$P_H(X_q = x) \triangleq \frac{P\{I : H(I) = 1, X_q(I) = x\}}{P\{I : H(I) = 1\}} \quad [1]$$

For simplicity, we have represented the set $\{I : [H, (q, x)](I) = 1\}$ in the numerator of [1] with the more intuitive expression $\{I : H(I) = 1, X_q(I) = x\}$, though this is decidedly an abuse of notation since the function $X_q(I)$ is not defined in the absence of the history H . Still, under P_H , X_q is a binary random variable which may or may not be “unpredictable.” To make this precise, we define the *predictability* of $q \in \mathcal{Q}_H$, given the history $H \in \mathbb{H}$, by $\rho_H(q) = |P_H(X_q = 1) - 0.5|$. Evidently, $\rho = 0$ indicates q is totally unpredictable and $\rho = 0.5$ indicates q is totally predictable.

Randomization. In general, many questions have answers with low predictability at each step k . Rather than select the *most* unpredictable question at step k , we make a random selection from the set of *almost unpredictable* questions, defined as those for which $\rho_H(q) \leq \epsilon$, where H is the history preceding the k 'th question. (In practice we choose $\epsilon = 0.15$, and we designate all such questions “unpredictable.”) In this way, we can generate many query streams for a given test image I , and develop multiple story lines within a query stream. In doing so, a path to instantiation might be $\{X_{ta} = 1, X_{tb} = 1, X_{ut\{a,b\}} = 1\}$, meaning that once there are instances of object type t with attribute a and also instances with attribute b , then the likelihood of having a unique (‘ u ’) instance with *both* attributes may rise to approximately one-half. Commonly, a designated region serves as an important instantiating attribute, as in the chain $\{X_{ta} = 1, X_{uta} = 0, X_{t\{a,b\}} = 1, X_{ut\{a,b\}} = 1\}$, where a is the designated region. Here, for example, t might refer to a person, of which several are partially visible in region a , but only one pos-

sesses the additional attribute b (e.g., “sitting”, “female”, or “wearing a hat”). There are more examples in Figure 2, and two complete sequences of questions in §5 of the SI Appendix.

Story lines and the simplicity preference. We impose constraints on the set of questions allowed at each step—the set of available followup questions given the history H , which we have denoted by \mathcal{Q}_H , is a small subset of the set of all possible questions, \mathcal{Q} . The main purpose is to encourage natural sequences, but these constraints also serve to limit the number of conditional likelihoods that must be estimated.

The loop structure of the query engine enforces a general question flow that begins with existence and uniqueness questions (\mathcal{Q}_{exist} , \mathcal{Q}_{uniq}), with the goal of instantiating objects. As objects are instantiated, the vision system is interrogated about their properties, meaning their attributes, and then their relationships to the already-instantiated objects. After these “story lines” are exhausted, the outer loops are revisited in search of new instantiations. The query engine halts when there are no more unpredictable existence or uniqueness questions. As already mentioned, all loops include randomization, meaning that the next query is randomly selected from the questions in \mathcal{Q}_H that are found to be unpredictable.

The pose attribute is especially useful to an efficient search for uniquely characterized objects, i.e. instantiation. Once the existence of an object that is partially visible within a region w is established, ensuing existence and uniqueness queries are restricted to w or its sub-regions. As these regions are explored, the unpredictability constraint then favors questions about the same object type, but annotated with additional attributes. Eventually, either an object partially visible in a sub-region of w is instantiated or the collection of unpredictable questions about such an object is exhausted. In the latter case the query engine returns to the outer loop and begins a new line of questions; in the former, it explores the attributes and relationships of the newly instantiated object. (All regions are rectangular and the full set, \mathcal{W} , is specified in the SI Appendix.)

Finally, there is a simplicity constraint that further promotes a natural line of questions. This can be summarized, roughly, as “one new thing at a time.” An existence, uniqueness, or attribute question, q , is considered simpler than an alternative question of the same type, q' , if q contains fewer attributes than q' . Given the unpredictable subset of \mathcal{Q}_H , simpler questions are favored over more complex questions, and questions of equal complexity are chosen from with equal likelihood. Further detail and pseudocode—see Algorithm—can be found in the SI Appendix.

Estimating predictability. The conditional likelihoods, $P_H(X_q = 1)$, are estimated from a training set \mathbb{T} in which all answers (or equivalent information—see Figure 3) are provided for each of n images from \mathcal{I} . The methods used to gather and annotate the training images are discussed in the next section, on the prototype VTT. The objects, people and vehicles, are located with bounding boxes and labelled with their attributes, and pairs of objects are labelled with their relationships.

The task of estimating conditional likelihoods, and therefore predictability, is guided in part by the ordering of questions built into the query engine, which, as already noted, begins with a search for an instantiated object, immediately followed by questions to determine its attributes, and then finally by an exploration of its relationships with any previously instantiated objects.

For instantiation questions, $q \in \mathcal{Q}_{inst} \triangleq \mathcal{Q}_{exist} \cup \mathcal{Q}_{uniq}$, the natural estimator $\hat{P}_H(X_q = 1)$ is the relative frequency (maximum likelihood) estimator

$$\frac{\#\{I \in \mathbb{T} : H(I) = 1, X_q(I) = x\}}{\#\{I \in \mathbb{T} : H(I) = 1\}} \quad [2]$$

Observe, though, that the number of images in the training set which satisfy the history H (i.e., for which $H(I) = 1$) is cut approximately in half at each step, and hence after about $\log_2 n$ steps direct estima-

tion is no longer possible. Consequently, to generate tests with more than ten or so questions, we are obliged to make “invariance” assumptions to allow for data pooling so as to expand the number of images from which these relative frequencies are computed. Specifically, if we assume that $X_q, q \in \mathcal{Q}_{inst}$, given the history $H \in \mathbb{H}$, depends only on a subsequence, H'_q of H , then the distribution on X_q is invariant to the questions and answers in H that were dropped, and the estimator [2] can be modified by substituting the condition $H(I) = 1$ by $H'_q(I) = 1$.

Let $w \in \mathcal{W}$ be the localizing region, possibly the entire image, referenced in the instantiation question q . H'_q is derived from H by assuming that the event $X_q = x$ is independent of all attribute and relationship questions in H , and all existence and uniqueness questions which involve localizations $w' \in \mathcal{W}$ which are disjoint from w , with the important exception of uniqueness questions that answered positive ($q' \in \mathcal{Q}_{uniq}, X_{q'} = 1$) and therefore instantiated a new object. In other words, the approximation is that, conditioned on the history, the distribution of an instantiation question depends only on the uniqueness questions that instantiated objects, and the existence and uniqueness questions that are localized to regions intersecting w . By preserving the instantiating questions in H , which addresses the potential complications introduced by the implied qualifier “not previously instantiated,” we guarantee that $H(I) = 1 \Rightarrow H'_q(I) = 1$ for all $I \in \mathbb{T}$, so that the population of images used to estimate $P_H(X_q = 1)$ with $H'_q(I)$ is no smaller than the one with $H(I)$ and typically far larger. More discussion, and a further invariance assumption leading to further improvement in population size, are included with the SI Appendix.

As for attribute questions, $q \in \mathcal{Q}_{att}$, which are always about the most recently instantiated object and always precede any relational information, the natural (relative frequency) estimator for $P_H(X_q = 1)$ is in terms of the population of labelled objects found in the training images, rather than the images themselves. Given a history H , consider a question of the form $q = o_t a$: “Does object o_t have attribute a ?” where o_t is an object of type $t \in \{person, vehicle\}$ and $a \in \mathcal{A}_t$. The history, H , defines a (possibly empty) set of attributes, denoted A , that are already known to belong to o_t . Let \mathcal{O}_T be the set of all annotated objects in the training set, and, for each $o \in \mathcal{O}_T$, let $\mathcal{T}_T(o)$ be the type of o and $\mathcal{A}_T(o)$ be the set of attributes belonging to o , e.g., $\mathcal{T}_T(o) = \{person\}$ and $\mathcal{A}_T(o) = \{female, adult, standing\}$ for the right-most object in Figure 3. The relative frequency estimator for $P_H(X_q = 1)$, using the population of annotated objects, is

$$\frac{\#\{o \in \mathcal{O}_T : \mathcal{T}_T(o) = t, A \cup \{a\} \subseteq \mathcal{A}_T(o)\}}{\#\{o \in \mathcal{O}_T : \mathcal{T}_T(o) = t, A \subseteq \mathcal{A}_T(o)\}} \quad [3]$$

There is again the sparsity problem, which we address in the same way—through invariance assumptions that effectively increase the number of objects. The set of attributes for objects of type t can be partitioned into subsets that can be reasonably approximated as independent conditioned on belonging to a particular object o_t . As an example, if $t = person$ then *crossing a street* is not independent of *standing still*, but both are approximately independent of gender, $\{male, female\}$, and of *child* versus *adult*, as well as whether or not o_t is *carrying something* or *wearing a hat*. These conditional independence assumptions decrease the size of the set A in [3], thereby increasing the set of $o \in \mathcal{O}_T$ used to estimate $P_H(X_q = 1)$.

The approach to relationship questions, $q \in \mathcal{Q}_{rel}$, is essentially the same as the approach to attribute questions, except that the training population is the set of *pairs* of objects in the training images, rather than the individual objects. The independence (invariance) assumptions include relationships that are independent of the attributes of the related objects (e.g., the relationship *driving/riding* a vehicle is assumed to be independent of both the gender of the person driving or riding, as well as whether the vehicle is dark or light colored, or whether or not its tires are visible) and relationships that are independent of each other (e.g., whether one vehicle is closer to the camera

than another vehicle is assumed to be independent of which vehicle is larger). A systematic accounting of the independence assumptions used in our prototype VTT, for both attribute and relationship questions, can be found in the SI Appendix and its accompanying tables.

A Prototype VTT

The data collection and annotation was performed by undergraduate workers at Johns Hopkins University. Unlike “crowd-sourcing”, this allowed for more customized instructions. Our dataset has 2,591 images, collected online using search engines such as Google street view and required to meet certain basic criteria: portray a standard city street scene; be obtained during daytime; have a camera height from roughly head-level to several feet above; contain clearly visible objects, attributes, and relationships from our vocabulary. The images are from large cities from many countries.

For annotation, we can rule out directly answering each binary question $q \in \mathcal{Q}$, since the questions only make sense in the context of a history— \mathcal{Q}_{att} and \mathcal{Q}_{rel} always refer to instantiated objects, and \mathcal{Q}_{exist} and \mathcal{Q}_{uniq} always include the not-previously-instantiated qualification. As discussed, a history itself can be viewed as a binary function of the image, but there are far too many for an exhaustive annotation. Instead, an essentially equivalent, but more compact and less redundant, representation was used. For example, once a “bounding box” is provided for every object (see the examples in Figure 3), there is a high likelihood that the correct answer to a localization question “Is x partially visible in region w ?” is determined by whether or not the bounding box of x intersects the region w . Thus bounding boxes were drawn around every instance of an object for which the annotator had no uncertainty about its category. For partially occluded objects, the bounding box was placed over the region of the image that the annotator expected the object would occupy had the object not been partially occluded. Attributes were annotated only for objects in which all the attributes were unambiguous, which alleviated the annotation of distant objects. Relationships were only annotated between pairs of objects with bounding boxes and for which at least one relationship from the type-dependent list was present.

The complete vocabulary is given in first two tables of the SI Appendix. The prototype includes only two types of objects: $\mathcal{T} = \{people, vehicles\}$. However, we also consider a few “parts”—things carried by people, and tires of vehicles, folded into the attribute categories. There is also an “attribute” for every element of a multi-scale collection \mathcal{W} of rectangular subsets of pixels referred to as “regions” (see SI Appendix for the complete collection); the “attribute” corresponding to $w \in \mathcal{W}$ is that the object instance is partially visible within w . The set of attributes also includes properties which are independent of positioning in the underlying scene, such as *female, child, wearing a hat, carrying something* for type *people*, and *car, truck, motorcycle, bicycle, light colored* for type *vehicle*. Still others refer to pose and context—e.g., *sitting, crossing a street, walking on a sidewalk, entering/exiting a building* for *people*, and *moving, stopped, parked, one tire visible, two tires visible* for *vehicles*. Ad-

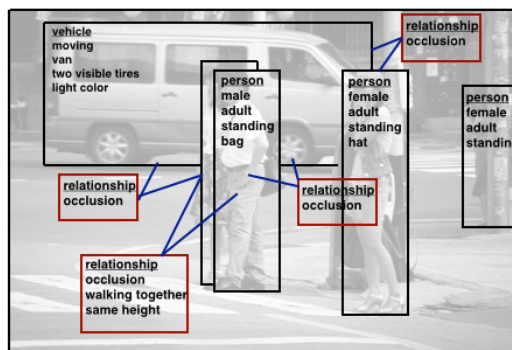


Fig. 3. Annotation provided by human workers.

ditionally, for both people and vehicles, the attribute *interacting with something* refers to any of a specific collection of relationships: for a person, *talking, walking together, holding hands* with another person, or *driving/riding, exiting, entering* a vehicle, and for a vehicle, *immediately behind, immediately in front* of another vehicle.

Relationships could be ordered or unordered. The unordered relationships between people are *talking, walking together, holding hands* and the ordered ones are about which person is *taller, closer to the camera*, and possibly *occluding* the other. The ordered relationships for two vehicles are the same, with *taller* replaced by *larger* and the addition of *immediately behind, immediately in front*; there are no unordered relationships between vehicles. Finally, a person may be *driving/riding, exiting, entering, occluding, or occluded by* a vehicle.

Level of difficulty. The vocabulary was selected to avoid query streams that would be considered hopelessly difficult by today’s computer-vision standards. Nevertheless, there are plenty of subtleties to challenge, and likely defeat, the best existing systems, e.g. the third stream in the SI Appendix (§5.3), which includes an example of extreme occlusion, two examples which require inferring that bicycles are moving, rather than stopped, and another occlusion that rests on the interpretation of a small number of pixels. A few additions to the vocabulary would dial up the difficulty, considerably, say adding the relationship “playing catch” or other objects like windows, signs, and tables and chairs, which are often nearly impossible to identify without context, especially when partially occluded.

Discussion

In the decades following Alan Turing computer vision became one of the most active areas of AI. The challenge of making computers “see” has attracted researchers from across science and engineering and resulted in a highly diverse set of proposals for formulating the “vision problem” in mathematical terms, each with its ardent advocates. The varying popularity of competing strategies can be traced in the proceedings of conferences.

Debates persist about what actually works and how to measure success. Until fairly recently, each new method was “validated” on homegrown data and with homegrown metrics. Recently, the computer vision community has accepted testing on large common datasets, as reviewed above, and various well-organized “challenges” have been accepted by many research groups. Many believe that adopting uniform metrics has made it easier to sort out what works appreciably better than before and accelerated progress.

But these metrics, such as false positive and false negative rates for sub-tasks such as detecting and localizing people, do not yet apply to the richer descriptions that human beings can provide, for example in applying contextual reasoning to decide whether or not a car is “parked” or is “larger” than another, or a person is “leaving” a building or “observing” something, or two people are “walking and talking together.” If annotating ordinary scenes with such precision is accepted as a benchmark for vision, then we have argued for raising the bar and proceeding directly to metrics for full-scale scene interpretation. We have proposed a “written” VTT as a step in this direction.

Many design decisions were made, some more compelling than others. “Story lines” approximate natural sequences of questions and are well handled by the loop structure of the algorithm. On the other hand, whereas conditional independence assumptions are probably a necessary approach to the data sparsity problem, the prototype lacks a unified implementation. Scaling to substantially larger vocabularies and more complex relationships, and deeper part/whole hierarchies, would be difficult to manage by simply enlarging the existing brute-force tabulation of dependency relationships (see SI Appendix). Possibly, the right approach is to build full-blown generative scene models, at least for the placements of parts and objects, and object groupings, from which predictability could be estimated via sampling or inferred by direct calculation.

Finally, coming back to a “conversation” with a machine, another possibility is a more free-form, open-ended “oral test”: the operator formulates and delivers a query to the system under evaluation, awaits an answer, and then chooses the next query, presumably based on the history of queries and system answers. As before, the operator may or may not provide the correct answer. This has the advantage that the operator can “probe” the system capacities with the singular efficiency of a human, for example detect and focus on liabilities and ask “confirmatory” questions. But the oral test has the disadvantage of being subjective and requiring rapid, basically real-time, responses from the system. On balance, the written test seems to be more practical, at least for the time being.

ACKNOWLEDGMENTS. We received helpful input from Mark Johnson of Macquarie University and Vincent Velten of the Air Force Research Laboratory. The work was partially supported by the Office of Naval Research under Contract ONR N000141010933, the Defense Advanced Research Projects Agency under Contract FA8650-11-1-7151, the National Science Foundation under Grant 0964416.

1. Turing AM (1950) Computing machinery and intelligence. *Mind* pp 433–460.
2. Saygin AP, Cicekli I, Akman V (2003) in *The Turing Test* (Springer), pp 23–78.
3. Russell SJ, Norvig P (2003) *Probabilistic Reasoning* (Prentice Hall).
4. Everingham M, Van Gool L, Williams CK, Winn J, Zisserman A (2010) The pascal visual object classes (voc) challenge. *International journal of computer vision* 88:303–338.
5. Deng J, et al. (2009) Imagenet: A large-scale hierarchical image database (IEEE), pp 248–255.
6. Lowe DG (2004) Distinctive image features from scale-invariant keypoints. *International journal of computer vision* 60:91–110.
7. Zhu Q, Yeh MC, Cheng KT, Avidan S (2006) Fast human detection using a cascade of histograms of oriented gradients (IEEE), Vol. 2, pp 1491–1498.
8. Yu G, Morel JM (2009) A fully affine invariant image comparison method (IEEE), pp 1597–1600.
9. Zhu SC, Mumford D (2007) *A stochastic grammar of images* (Now Publishers Inc).
10. Ommer B, Sauter M, Buhmann JM (2006) Learning top-down grouping of compositional hierarchies for recognition (IEEE), pp 194–194.
11. Chang LB, Jin Y, Zhang W, Borenstein E, Geman S (2011) Context, computation, and optimal roc performance in hierarchical models. *International journal of computer vision* 93:117–140.
12. Lu W, Lian X, Yuille A (2014) Parsing semantic parts of cars using graphical models and segment appearance consistency. *arXiv preprint arXiv:1406.2375*.
13. Hinton G, Osindero S, Teh YW (2006) A fast learning algorithm for deep belief nets. *Neural computation* 18:1527–1554.
14. Fei-Fei L, Fergus R, Perona P (2003) A Bayesian approach to unsupervised one-shot learning of object categories (IEEE), pp 1134–1141.
15. Amit Y, Trouné A (2007) Pop: Patchwork of parts models for object recognition. *International Journal of Computer Vision* 75:267–282.
16. Felzenszwalb PF, Huttenlocher DP (2005) Pictorial structures for object recognition. *International Journal of Computer Vision* 61:55–79.
17. Felzenszwalb PF, Girshick RB, McAllester D, Ramanan D (2010) Object detection with discriminatively trained part-based models. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 32:1627–1645.
18. Krizhevsky A, Sutskever I, Hinton GE (2012) Imagenet classification with deep convolutional neural networks pp 1097–1105.
19. Girshick R, Donahue J, Darrell T, Malik J (2013) Rich feature hierarchies for accurate object detection and semantic segmentation. *arXiv preprint arXiv:1311.2524*.
20. Oquab M, Bottou L, Laptev I, Sivic J, et al. (2014) Learning and transferring mid-level image representations using convolutional neural networks.
21. Zhang N, Paluri M, Ranzato M, Darrell T, Bourdev L (2013) Panda: Pose aligned networks for deep attribute modeling. *arXiv preprint arXiv:1311.5591*.
22. Hariharan B, Arbeláez P, Girshick R, Malik J (2014) in *Computer Vision—ECCV 2014* (Springer), pp 297–312.
23. Xiao J, Hays J, Ehinger KA, Oliva A, Torralba A (2010) Sun database: Large-scale scene recognition from abbey to zoo (IEEE), pp 3485–3492.
24. Russell BC, Torralba A, Murphy KP, Freeman WT (2008) Labelme: a database and web-based tool for image annotation. *International journal of computer vision* 77:157–173.
25. Yao B, Yang X, Zhu SC (2007) Introduction to a large-scale general purpose ground truth database: methodology, annotation tool and benchmarks (Springer), pp 169–183.
26. Endres I, Farhadi A, Hoiem D, Forsyth DA (2010) The benefits and challenges of collecting richer object annotations (IEEE), pp 1–8.
27. Oh S, et al. (2011) A large-scale benchmark dataset for event recognition in surveillance video (IEEE), pp 3153–3160.
28. Özdemir B, Aksoy S, Eckert S, Pesaresi M, Ehrlich D (2010) Performance measures for object detection evaluation. *Pattern Recognition Letters* 31:1128–1137.