

**Title:** Full interpretation of minimal images.

**Author names and affiliations:** Guy Ben-Yosef<sup>1,2,\*</sup>, Liav Assif<sup>1</sup>, Shimon Ullman<sup>1,2</sup>

1. Department of Computer Science and Applied Mathematics, Weizmann Institute of Science, Rehovot 7610001, Israel
2. Center for Brains, Minds and Machines, Massachusetts Institute of Technology, Cambridge, MA 02139, USA.

\* Present address: Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA 02139, USA.

Corresponding author: shimon.ullman@weizmann.ac.il

Word count for manuscript: 9,486

**Abstract:**

The goal in this work is to model the process of ‘full interpretation’ of object images, which is the ability to identify and localize all semantic features and parts that are recognized by human observers. The task is approached by dividing the interpretation of the complete object to the interpretation of multiple reduced but interpretable local regions. In such reduced regions, interpretation is simpler, since the number of semantic components is small, and the variability of possible configurations is low.

We model the interpretation process by identifying primitive components and relations that play a useful role in local interpretation by humans. To identify useful components and relations used in the interpretation process, we consider the interpretation of ‘minimal configurations’: these are reduced local regions, which are minimal in the sense that further reduction renders them unrecognizable and uninterpretable. We show that such minimal interpretable images have useful properties, which we use to identify informative features and relations used for full interpretation. We describe our interpretation model, and show results of detailed interpretations of minimal configurations, produced automatically by the model. Finally, we discuss implications of full interpretation to difficult visual tasks, such as recognizing human activities or interactions, which are beyond the scope of current models of visual recognition.

**Keywords:**

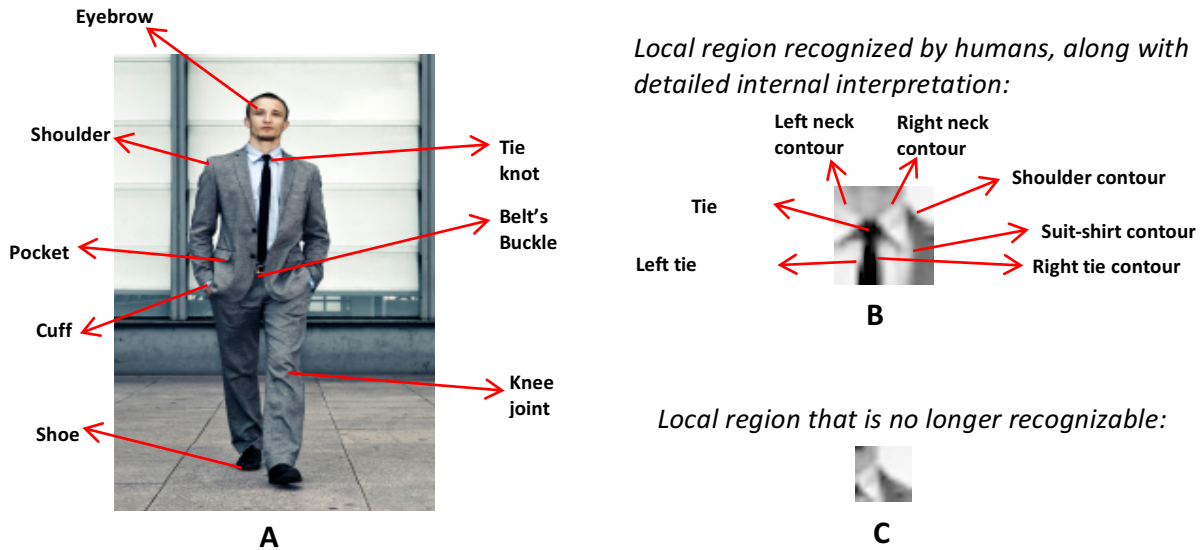
Image interpretation; Minimal images; Parts and relations; Top-down processing;

## **1. Introduction**

Humans can recognize in images not only objects (e.g., a person) and their major parts (e.g., head, torso, limbs), but also multiple semantic components and structures at a fine level of detail (e.g., shirt, collar, zipper, pocket, cuffs etc.), as in Fig. 1A. Identifying detailed components of the objects in the image is an essential part of the visual process, contributing to the understanding of the surrounding scene and its potential meaning to the viewer (Sec. 6.1). Although this capacity is of fundamental importance in human perception and cognition, current understanding of the processes involved in detailed image interpretation is limited.

From the modeling perspective, existing models cannot deal well with the full problem of detailed image interpretation, and, as discussed below, the limitations are of fundamental nature. Computational models of object recognition and categorization have made significant advances in recent years, demonstrating consistently improving results in recognizing thousands of natural object categories in complex natural scenes (Sec. 2). However, existing models cannot provide a detailed interpretation of a scene's components in a way that will approximate human perception. For example, for a given image such as Fig. 1A, existing models can correctly decide if the image contains a person (e.g., Csurka et al., 2004; Simonyan & Zisserman, 2015), and can locate a bounding box around the body (e.g., Dalal & Triggs, 2005; Girshick et al., 2014). At a more refined level, current algorithms can provide an approximate segmentation of the body figure (e.g., Long et al., 2015), and can locate image region containing the main body parts, such as the torso region, the face, or the legs (e.g., Chen et al., 2014; Vedaldi et al., 2014), or keypoints at the joints (e.g., Chen & Yuille, 2014; Wei. et al., 2016). However, existing computational models cannot achieve the accuracy and richness of the local interpretation of image components perceived by a human observer (e.g., as in Fig. 1B).

To clarify the terminology, by the term 'visual interpretation' we refer to a mapping between images and a non-visual domain, such as the domain of objects and object categories, object parts, and other physical entities, which is the semantic domain. For instance, within a face image, a particular contour may correspond to, say, the mouth's upper lip. The contour is an image component, the upper-lip is a semantic component,

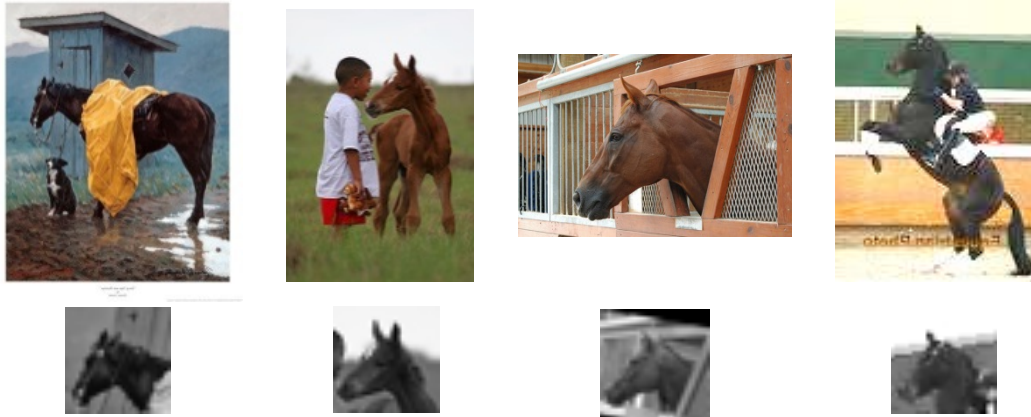


**Figure 1.** (A). Humans can identify a large number of semantic features and parts in an object image. In the image of a walking person, features like the suit's pocket, tie's knot, left shoe, or the right ear, are easily identified by humans, among many others. (B). A detailed interpretation of a small image regions, as identified by human observers. In small local regions, the number of semantic components is significantly smaller than in full images, and variability is reduced. (C). When the local region becomes too limited, human observers can no longer recognize and interpret its content when presented on its own (Ullman et al., 2016).

and the interpretation process maps between the two. The ultimate goal of a 'full interpretation' model is to identify all the semantic components that human observers can identify in an image.

### 1.1. Local image interpretation

Producing a detailed interpretation of an object's image is a challenging task, since a full object may contain a large number of identifiable components in highly variable configurations. We approach this task by decomposing the full object or scene image into smaller, local, regions containing recognizable object components. There are several advantages to perform the interpretation first in local regions, and then combine the results. First, as exemplified in Fig. 1B, in such local regions the task of full interpretation is still possible (Ullman et al., 2016), but it becomes more tractable, since the number of semantic recognizable components is highly reduced. As will be shown (Sec. 5), reducing the number of components plays a key factor in effective interpretation. At the same time, when the interpretation region becomes too limited, observers can no longer interpret or even identify its content, as illustrated in Fig. 1C (Ullman et al., 2016). The goal of the model is therefore to apply the interpretation



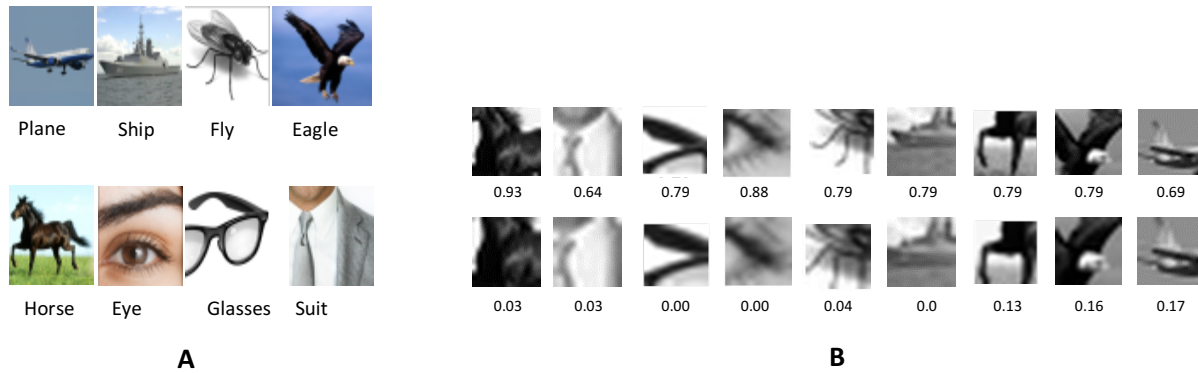
**Figure 2.** Complete horse images taken from ImageNet object recognition benchmark (Deng et al., 2012), and a small recognizable region that is interpretable (similar to Fig. 4A), next to each complete horse image illustrating the reduced variability in small recognizable region vs. the complete object image.

process to local regions that are small, yet interpretable on their own by human observers. A second advantage of applying the interpretation locally is that variability of configurations taken from the same object class, but limited to local regions, is often significantly lower compared with complete object images. For example, the full horse images in Fig. 2 (taken from the ‘horse’ category in ImageNet, Deng et al., 2012), a common benchmark for evaluating object recognition models) are quite different from each other, but can become significantly more similar at the level of local regions. Finally, as will be discussed in the next section, the image of a single object typically contains multiple, partially overlapping regions, where each one can be interpreted on its own. Due to this redundancy, performing the interpretation locally and then combining the results increases the robustness of the full process to local occlusions and distortions.

## 1.2. Minimal configurations

In performing local interpretation, how should an object image be divided into local regions? The approach we take in this study is to develop and test the interpretation model on regions that can be interpreted on their own by human observers, but at the same time are as limited as possible. We used for this purpose a set of local recognizable images derived by a recent study of minimal recognizable images (Ullman et al., 2016). We briefly describe below how these images were obtained, and then explain the reasons for using these local images in developing and testing the interpretation model.

A ‘minimal configuration’ (also termed Minimal Recognizable Configuration, or MIRC) is defined as an image patch that can be reliably recognized by human observers,



**Figure 3.** Minimal configurations adapted from Ullman et al. (2016). (A). The search for minimal images started from different object images (8 shown here), each composed of 50x50 image samples. (B). **Top row:** minimal images discovered by the search. **Bottom row:** sub-minimal configurations, which are slightly reduced versions of the images on top. Numbers below each image show correct recognition rate by 30 human observers. Small changes to the local image at the minimal configuration level can have large effect on recognition. A data set of such pairs is used below for modeling the interpretation of local regions.

which is minimal in the sense that further reduction by either size or resolution makes the patch unrecognizable. To discover minimal configurations, an image patch was presented to observers: if it was recognizable, 5 descendants were generated by either cropping at one corner, or reducing resolution of the original patch. A recognizable patch is identified as a ‘minimal configuration’ if none of its 5 descendants reach recognition criterion (50%). A search started with images from different object classes (Fig. 3A), and identified their minimal configurations over all possible positions, sizes and resolutions. Each subject saw a single patch only from each original image, requiring over 15,000 subjects. Testing was therefore done online using Amazon’s Mechanical Turk platform (MTurk), combined with laboratory controls. At the end of the search, each object class was covered by multiple minimal configurations at different positions and sizes. Minimal configurations were on average about 15 image samples in size; some contained local object parts, others were more global views at a reduced resolution. Examples of identified minimal configurations are shown in the top row of Fig. 3B.

A notable aspect of the results for the purpose of the current study, is the presence of a sharp transition for almost all minimal configurations from a recognizable to a non-recognizable minimal image: a surprisingly small change at the minimal-configuration level can make it unrecognizable. Examples are shown in Fig. 3B, bottom row, together with their respective recognition rates. The small changes between minimal vs. sub-minimal configurations that cause large drop in recognition are used below to identify features and relations used in the interpretation process. It was also found that the large

gap in human recognition rate between minimal and sub-minimal images is not reproduced by current computational models of human object recognition (Serre et al., 2007) and recent deep network models (Krizhevsky et al., 2012; Simonyan & Zisserman, 2015). As shown below (Sec. 5.2), the full interpretation model can provide at least a partial explanation to this sharp drop in recognition.

### **1.3. Recognition and interpretation**

With respect to local interpretation, recognition tests of minimal images showed that although the minimal images are ‘atomic’ in the sense that their partial images become unrecognizable, humans can consistently recognize multiple semantic features and parts within them. It was noted (Ullman et al., 2016) that recognition and interpretation of minimal images go hand in hand in the sense that under the tested conditions (unlimited viewing time), when subjects correctly recognized a minimal image, they were also able to provide an internal interpretation of multiple internal components. Since in minimal images all the available information is, by definition, crucial for recognition, we propose in the model below that all the interpreted components of minimal images also contribute to their recognition. As described further below (Sec. 4.3), in the model, the full interpretation process contributes to accurate recognition, since a potential false detection can be rejected if it does not have the expected internal interpretation.

For the purpose of modeling human visual interpretation, our initial focus is on the interpretation of minimal images, for the following reasons. First, they provide a useful test set for the model: since they are interpretable by humans, a theory of human image interpretation should be applicable to such configurations. Second, we use minimal and sub-minimal pairs with a large gap in recognizability and interpretability as a source for inferring useful features for the interpretation of minimal images (Sec. 4). Before describing the model, we briefly describe past work related to visual object interpretation.

## **2. Related work on visual object interpretation**

Visual recognition can take place at different levels of details, from full objects and their main parts, to fine details of objects' structure. In modeling human visual perception, as well as in computer vision, much of the work to date has focused on relatively coarse

levels, rather than full object interpretation considered here. For example, in the Recognition by Components (RBC) model of human object categorization (Biederman, 1987), objects are represented in terms of a small number of 3-D major parts. A leading biological model on the human object recognition system, the HMAX model (Riesenhuber & Poggio, 1999; Serre et al., 2007) produces as its output general category labels of full objects, rather than a detailed interpretation.

A model for human image interpretation (Epshtein et al., 2008) was shown to provide partial image interpretation by a combination of bottom-up with top-down processing. The model uses a hierarchy of informative image patches to represent object parts at multiple levels. The current model also uses a combination of bottom-up and top-down processing, but it provides a significantly richer interpretation, and based on computational and psychophysical considerations, it uses an extended set of elements and relations. A preliminary version of the model was described in Ben-Yosef et al. (2015). The current model extends the early version in the use of minimal images (rather than local image regions), in testing on multiple classes, and in comparisons with human vision.

In computer vision, there has been rapid progress in different aspects of object and scene recognition, based primarily on deep convolutional neural networks and related methods (Hinton, 2007; LeCun et al., 2015; Yamins et al., 2014; Krizhevsky et al., 2012; Simonyan & Zisserman, 2015; He et al., 2016). Such methods have also been adapted successfully for image segmentation, namely the delineation of image regions belonging to different objects. For example, recent algorithms (e.g., Long et al., 2015; Hariharan et al., 2015; Dai et al., 2016) can identify image regions belonging to different objects in the PASCAL (2012) or CoCo benchmarks (2015); however, they do not locate the precise object boundaries, and do not identify the object's semantic components.

A number of studies have begun to address the problem of a fuller object interpretation, including methods for part-based detectors, object parsing, and methods for so-called fine-grained recognition. Recent examples include modeling objects by their main parts, for example an airplane's nose, tail, or wing (Vedaldi et al., 2014), or modeling human-body parts such as the head, shoulder, elbow, or wrist (e.g.,

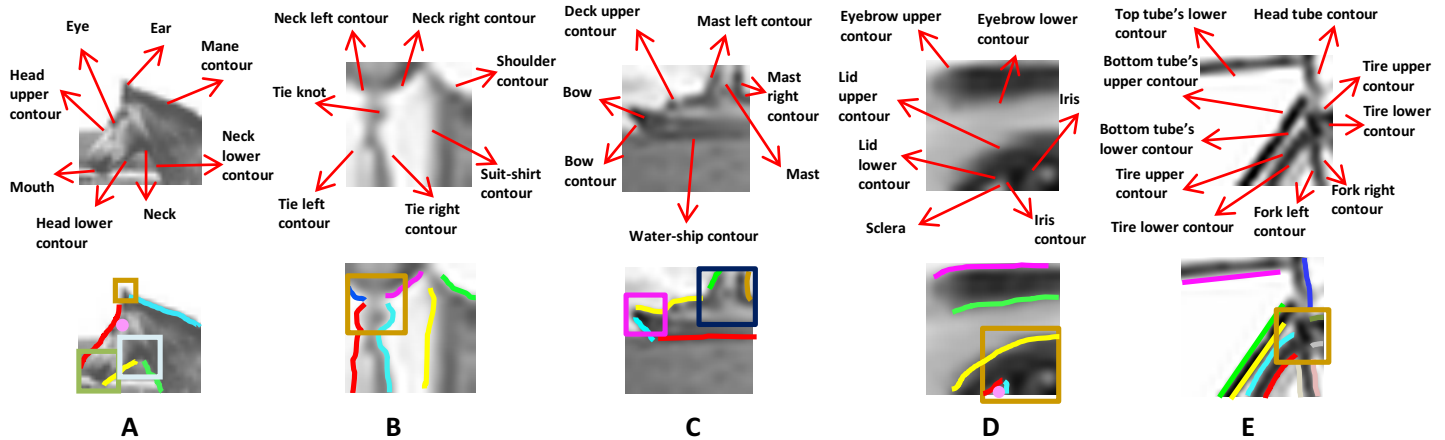
Felzenszwalb et al., 2010; Girshick et al., 2015). Related models provide segmentation at the level of object parts rather than complete objects (applied e.g. to animal body parts such as head, leg, torso, or tail, e.g., Azizpour & Laptev, 2012; Chen et al., 2014).

Another form of interpretation has been the detection of key-points within an object, such as key-points of the human body (e.g., Andriluka et al., 2014; Chen & Yuille, 2014; Tompson et al., 2015) and within the human face (e.g., Yang et al., 2015).

The goal of interpretation models, such as those above, is to produce the semantic structure in an image region. The model is usually given during learning a set of training images together with their interpretation, i.e., a set of semantic elements within each image, and the goal of the model is to identify similar elements in a novel image. In a correct interpretation, the internal components are expected to be arranged in certain consistent configurations, which are often characterized in the model by a set of spatial relations between components. The task of producing the semantic interpretation can therefore be naturally approached in terms of locating within an image region a set of elements (primitives) arranged in a configuration that satisfies relevant relations. The term ‘relations’ also includes properties of single elements (e.g., the curvature, location, or size of a contour), which can be considered as unary relations.

A number of algorithms have been developed and used in the field of machine vision under the general term ‘structured prediction’ to deal with problems related to the learning and discovery of image structures, such as Conditional Random Field (Lafferty et al., 2001), or Structured Support Vector Machine (Joachims et al., 2008). These models are given the set of possible relations to use, and then learn the specific parameters from examples. In terms of properties and relations, in most visual models that deal with image structures, such as the ones above, part properties (unary relations) are limited to local, deep CNN-based features, and binary relations are limited to relative displacements of components (parts or keypoints). As elaborated below, results of the present modeling show that the capacity to provide full interpretations requires the use of features and relations, which go beyond those used in most current recognition models.





**Figure 4.** Human interpretation of minimal configurations. (**Top row**). All components that were identified consistently by human observers (Appendix A). (**Bottom row**). In the interpretation model the components are represented by three types of primitives: points, contours, regions, together with relations between them. For each column, the identified components on the top panel are plotted in different colors on the bottom panel, and by either a point, a contour, or a region (an outlined square).

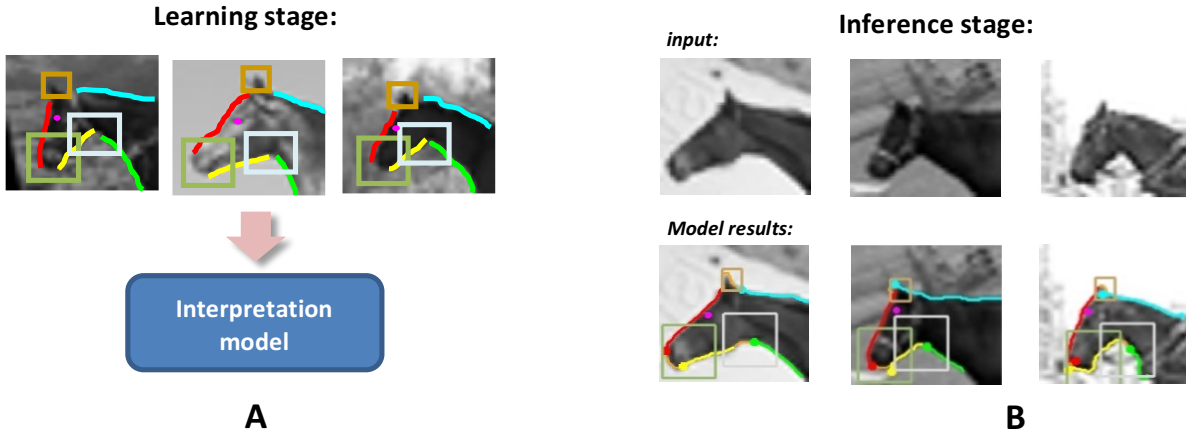
### 3. Model description

Our interpretation scheme has two main components: in the learning stage, it learns the semantic structure of an image region in a supervised manner, and in the interpretation stage, it identifies the learned structure in similar image regions. These two stages are described in the rest of this section.

#### 3.1 Learning setup

The learning stage derives the semantic structure of an object region based on positive examples coming from class images, and negative examples derived by the system from similar but non-class images. We first describe how these training examples are obtained, and then how the region's semantic structure is learned from them.

Positive examples are supplied manually during a preparation stage as a set of image regions with their interpretation, namely, the semantic elements that should be identified and localized. Since the goal is to model humans' ability to obtain a detailed local interpretation, the target set of semantic primitives to identify was collected for different minimal images using human observers. The semantic features to be identified by the model, e.g. 'ear', 'eye', 'tie knot' etc., were features that human observers label consistently in minimal images, verified using a Mechanical Turk procedure (see examples in Fig. 4, top row, and Appendix A for procedure details). The average number of consistently identified elements within a single minimal configuration was 8. To



**Figure 5.** Stages in the interpretation scheme, with horse-head as an example. **(A).** Point, contour, and region primitives that represent the identified parts (cf. Fig. 4) are annotated in training examples (several shown here), and are used to learn an interpretation model, which combines the primitives with relations between them. **(B).** Results of the interpretation model for 3 novel examples of the horse head minimal configuration.

capture the recognized internal components fully as perceived by humans, the primitive elements in the model were divided into three types: two-dimensional (2-D) regions, 1-D contours, and points (0-D). Example sets of primitives for modeling the interpretation of minimal images are shown in Fig. 4, bottom row. For instance, a point-type primitive may describe the eye in the horse head model (Fig. 4A), and a contour-type primitive describes borders such as the borders of the tie in the man-in-suit (Fig. 4A). Larger semantic features marked by observers such as the ship’s ‘bow’ region or the tie’s ‘knot’, were marked as region primitives (outlined squares in Fig. 4, bottom row). The three types of primitives are also supported by psychophysical and physiological studies (e.g., Attenev, 1954; Pasupathy & Connor, 1999).

Given the semantic elements identified by humans in a minimal image of class  $C$  (e.g., a horse-head), we prepared a set of annotated images, in which the semantic components (denoted  $P_C$  below) were marked manually (with automatic refinement). Examples for such annotations are shown in Fig. 5A. The unsupervised learning of components and relations are considered briefly in the final discussion (Sec. 6.2).

Having a set of interpretation examples, the learning process next searches automatically for negative interpretation examples – these are non-class images that are potentially confusable with class images. The procedure for identifying so-called ‘hard negatives’ (e.g., Felzenszwalb et al., 2010; Azizpour et al., 2012; Chen et al., 2014) starts

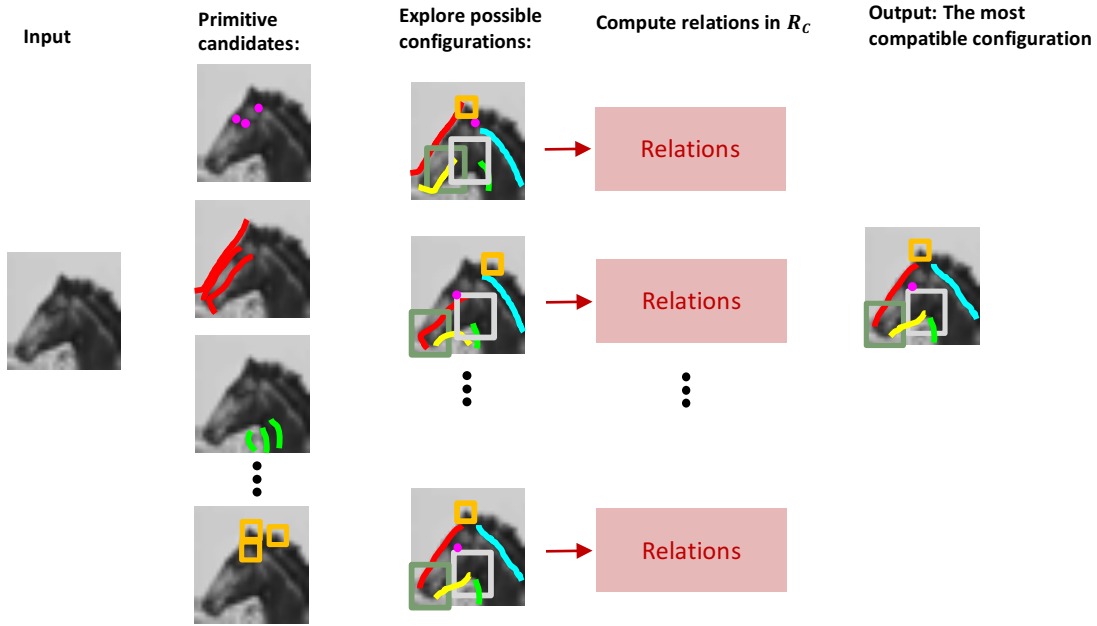
from random non-class examples and then iterates over two steps: finding non-class examples with high interpretation score, then adding them to the training set and re-training the model.

### 3.2 Learning the semantic structure:

For a minimal configuration  $C$ , we define its semantic structure  $S_C$  as a pair of two sets: the set of semantic components  $P_C$  mentioned in Sec. 3.1 (also called below the ‘primitives’), and a set of relations between primitives, denoted by  $R_C$ , namely

$$S_C = \langle P_C, R_C \rangle.$$

We include properties of a single primitive as a relation with a single argument. A basic problem at this stage is therefore to learn a set of relations that are useful for identifying configurations, namely, which appear in the positive class examples, and distinguish them from configurations found in the similar but non-class negative examples. The relevant relations for a given image are selected automatically during learning from an initial set (termed ‘relations library’ below) of potentially informative and useful relations to compute (see Sec. 5 on how this set was obtained). For instance, whether the relation ‘containment’ between pairs of primitives should be included in  $R_C$ , all potential pairs of primitives are examined, using the positive and negative examples, to test if one primitive is consistently contained within the other. (See Appendix. B.1 for how the contribution of a relation to the final interpretation was measured.) Each of the relations used in the interpretation scheme is given an index, e.g. the relation ‘containment’ may have the index ‘4’. Following selection, the set of all informative relations identified in a given minimal image  $C$  are represented by the vector  $R_C$ . Each element in  $R_C$  specifies a relation, and its relevant components. For example, the 3<sup>rd</sup> component of  $R_C$  (i.e.,  $R_C^{(3)}$ ) could be the triplet (4, 5, 7). This triplet means that relation 4, which is ‘containment’, holds between components 5 and 7 in the local image model, specifying that component 5 in the local model should be contained inside component 7. Similarly, the element in position 4 in  $R_C$  (i.e.,  $R_C^{(4)}$ ) can be a ‘straightness’ (unary) relation of primitive index 2, etc. Relations in our model could be either binary, e.g., ‘containment’, or represented by a



**Figure 6:** An overview of the interpretation process of a novel image. From left to right, **Input image**. **Detected candidates:** of the primitive components, examples for 3 candidates of each primitive are shown. **Configurations:** examples of possible configurations of detected primitives (denoted by  $\pi$  in Appendix B); the one at the bottom is the optimal one. **Computing relations:** compute the relations in  $R_c$  for each candidate configuration (the vector  $\phi_s(I, \pi)$  in Appendix B). **A compatibility score:** a scoring function ( $g(\phi_s(I, \pi); w)$  in Appendix B) is computed for each configuration. The configuration with highest score is returned as final interpretation.

scalar, e.g., the property ‘location’, specifying the location of a component within the local image.

A detailed description of the learning model and procedure based on ‘structured learning’ framework (e.g. Shalev-Shwartz & Ben-David, 2014) is given in Appendix B. For a novel image, the vector representation  $R_c$  of the image structure is derived as described in the next section, and then used for final interpretation decision.

### 3.3. Interpretation of a novel image

In this section we assume that a local image region has been identified as a likely candidate of a particular object or object part, and the current task is to produce an internal interpretation of the candidate region, and make a final decision about its identity. More details of the algorithm are given in Appendix B, and we also describe later (Sec. 6.3) how the initial detection and full interpretation are integrated together in a combined scheme of a bottom-up stage identifying likely candidates (e.g. by a DNN classifier trained for the task), followed by a top-down interpretation and validation stage.

The interpretation process starts with a candidate region and its proposed classification (e.g., that it contains a horse-head). The process then uses the learned model of the region's structure to identify within the region a structure that best approximates the learned one. This process proceeds in two main stages. The first is a search for local primitives, namely points, contours, and regions in the image, to serve as potential candidates for different components of the expected structure. The second stage searches for a configuration of the components that best matches the learned structure.

To match a given image configuration to the learned structure, we compute the relations in  $R_C$  for this configuration, and then use a compatibility scoring function based on a random forest classifier (Breiman, 2001, Appendix B), which produces a number that evaluates the degree to which the configuration is a correct interpretation of the input image. The interpretation scheme finally selects the highest-scoring configuration. A search among multiple configurations is feasible due to the small number of primitives in the local region. This overall process is illustrated in Fig. 6. A detailed description of the scoring procedure and the optimization part (i.e., finding the most compatible configuration) is given in Appendix B.

#### **4. Useful relations for interpretation**

Producing an interpretation of an image region requires the localization of its participating components, and verifying their correct configuration. The model verifies the structure using inter-elements relations, and a natural question is therefore which relations are useful in modeling local semantic structures. The visual system is known to be sensitive to a range of spatial properties and relations between components such as curvature, straightness, proximity, relative displacement, collinearity, inclusion, bisection, and others, which have been studied both perceptually and physiologically (see review in Sec. 4.1 below). It is unknown, however, which relations play a significant role in the task of visual interpretation. In this section we describe the methods we used to identify informative relations for interpretation, which were then included in the set of interpretation relations.

In contrast with the richness of relations that can be efficiently perceived by the visual system (Sec. 4.1), the majority of models for image recognition and interpretation have been based on a limited number of basic relations. Recognition models based on deep networks obtain high performance in basic categorization, but when the task requires a more detailed interpretation, e.g. identifying keypoints in human pose estimation, performance improves by explicitly incorporating inter-element relations, in particular relative displacement and orientations, using e.g. CRF models (Chen & Yuille, 2014; Wei et al., 2016). We next examined the set of relations which are informative for the full interpretation of local images.

The availability of minimal images allowed us to examine whether basic relations used in previous schemes are sufficient for producing an accurate interpretation by the interpretation model. Minimal configurations are by construction non-redundant visual patterns, and therefore their recognition and interpretation depend on the effective use of all the available visual information. It consequently becomes of interest to examine the performance of a model that uses a limited set of relations when applied to the interpretation of minimal images. We therefore constructed a version of the interpretation scheme, where the set of relations was limited to displacement and proximity relations. Performance for this version proved insufficient compared with human interpretation (see more details in Sec. 5). This limitation motivated the search for additional informative relations, which were shown to improve the interpretation of minimal images. It is worth noting that since minimal images contain small sets of components, it becomes more feasible to use in the model inter-element relations that are more complex and more computationally demanding than used in past models.

We describe in Sections 4.1-4.4 below the process of identifying informative relations for the interpretation process. Previous psychophysical and physiological studies have proposed a number of relations that the visual system is sensitive to. These provided an initial set of candidate relations, and each relation was evaluated by measuring its contribution to the interpretation model applied to a test set of minimal configurations, combined with sub-minimal configurations (Sec. 4.2) and hard-negative examples (Sec. 4.3). We finally describe the relations that were found to be informative for learning

interpretation. We also consider (Sec. 6.2) how a more complete set of informative interpretation relations could be learned and refined over time.

#### **4.1 Relevant visual relations in past literature**

The study of relations between elements in the visual field dates back at least to the Gestalt school and its principles of perceptual organization (Wertheimer, 1923). These principles were based on relations that group visual elements together to be perceived as coherent units, and included proximity, similarity, connectivity, symmetry, and continuity between dots, contours, or regions. Psychophysical experiments since have shown that the human visual system is effortlessly sensitive to a range of spatial properties and relations between visual elements. Such relations include: parallelism and symmetry (e.g., Feldman, 2007; Machilsen et al., 2009), curvature and convexity (e.g., Foster et al., 1993), connectedness of blobs (e.g., Palmer & Rock, 1994), and connectedness of contours (e.g., Jolicoeur et al., 1986), continuity of contours (e.g., Kanizsa, 1979), co-linearity (e.g., Field et al., 1993) and co-circularity (e.g., Parent & Zucker, 1989) of contours, relative length of lines and contours (e.g., Saarela et al., 2009), bisection (e.g., Westheimer et al., 2001), and inclusion (Ullman, 1984).

For many of these relations, it remains unclear whether they are being formed at early stages of visual perception in a bottom-up manner (e.g., Kanizsa, 1979; Field et al., 1993; Parent & Zucker, 1989) or at later stages, applies in a top-down manner to early visual representations (e.g., Ullman, 1984; Jolicoeur et al., 1986; Roelfsema et al., 1998). It is also still unclear which of the relations perceived effortlessly by humans play also a direct role in recognition and interpretation. The computational test described below evaluated directly the contribution of different relations to the interpretation of minimal image. To search for informative relations for interpretation, we started with a list of visual relations identified in past studies listed above, called the ‘candidate relations’, and tested their contribution to the interpretation process applied to minimal and sub-minimal images and hard-negative examples, as discussed next (Sections 4.2-4.3).

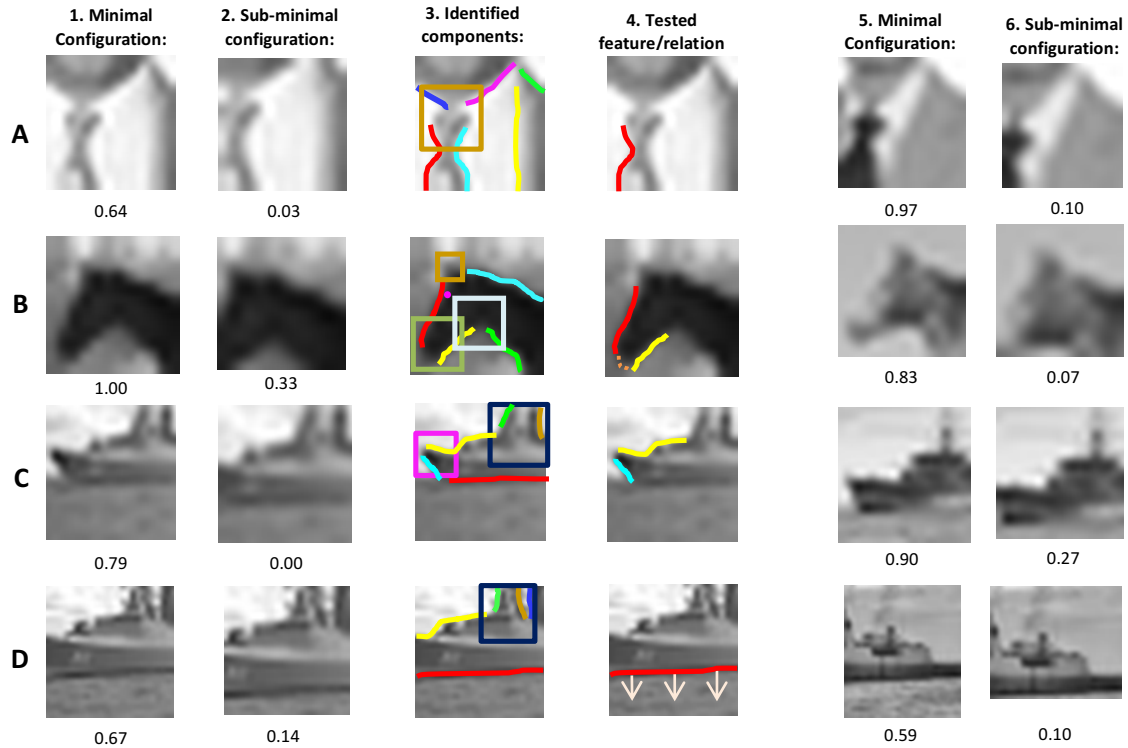
## 4.2 Useful relations from minimal vs. sub-minimal images

The sharp drop in humans' ability to recognize and interpret a minimal configuration when the image is slightly reduced (Ullman et al., 2016), provided a tool for identifying useful relations for modeling human interpretation. A minimal image was compared with its similar, but unrecognized sub-image, to identify either a missing component (e.g., a contour), or a relation (e.g., connected contours that become unconnected), which were present in the minimal image but not in the sub-minimal configuration. Examples are illustrated in Fig. 7, where pairs of minimal vs. sub-minimal configurations are shown (columns 1-2), along with the sets of internal semantic components that were identified by human observers in the minimal images (column 3). By using the human annotations, we found if any components in the minimal image were missing in its sub-minimal image. Using the set of candidate relations, we identified relations that are satisfied in the minimal but not the sub-minimal image. The missing component or relation may not be unique, and in such cases we evaluated a number of alternatives. The examples in Fig. 7 include the existence of the left-side tie contour (7A), connectedness of the two horse muzzle contours (7B), high-curvature meeting of contours (7C), and characteristic texture in the water region (7D).

We next evaluated for each of the missing components or relations, how consistent it is among other examples of minimal images, and how informative it is for the interpretation process, using our full data set of training examples. We start by testing for consistency in the set of minimal and sub-minimal pairs of the same class namely, finding additional pairs separated by the same component or relation (Fig. 7, columns 5-6). As an initial filtering stage, components or relations playing a role in at least 3 additional pairs were kept for the next stage, in which they were tested by their contribution to the performance of the interpretation algorithm. Each relation (similarly for candidate components) was tested by adding it to the set of relations (namely, to the relations  $R_C$ ), training a new interpretation algorithm, and measuring the difference in interpretation performance with and without this relation.

In more details, to test how informative is a given relation to the interpretation process, we have trained and compared two alternative versions of the interpretation model. The first version, (termed 'basic'), included a limited set of relations commonly





**Figure 7.** *Inferring relations between internal components with large contribution to recognition and interpretation. Minimal and sub-minimal pairs (columns 1,2, recognition rate shown below the images), are shown with internal components recognized by humans in the minimal images (column 3). To identify useful components and relations for interpretation, we compared the minimal and sub-minimal images. Using the identified components, we found if any component in (1) are missing in (2). Using the set of candidate relations, we identified relations that are satisfied in (1) but not in (2). The contribution of each missing component or relation was then evaluated using training examples (see text). When necessary, several alternatives were evaluated. Examples of informative components and relations are shown in column 4. Examples of additional MIRC / sub-MIRC pairs in the training set with the same missing component or relation, with its effect on recognition, are shown in columns 5,6. Inferred components and relations illustrated in the figure are: missing contour element (in A), connectedness of two contours (B), contours meet at high curvature (C), and characteristic texture in a region (bounded by the red contour and image border) in (D).*

used in the visual structure modeling literature (Sec. 2), namely, unary relations based on local texture and shape appearance, and binary ones based on the relative displacement of components. The basic model is then compared to a second (termed ‘augmented’) interpretation model, where the basic set of relations is augmented with the relation we wish to test. Performance of both models was evaluated on a data set, which included for each of the minimal images in Fig. 4, a set of 120 positive examples, and 8000 negative (non-class) examples, split between training and validation sets. Performance of the two models was compared by classification by the random forest classifier (using the Out-Of-Bag test for strength of random forest features, Brieman, 2001, Appendix B), to assess the contribution of each new relation. Relations that improved random forest

classification average precision by 1% or more (found in pilot experiments to be significant) were incorporated in our set of relations. The final set was subsequently used in the overall evaluation of the model applied to the interpretation and recognition of minimal and sub-minimal images (Sec. 5). Fig. 7 illustrates the process for example relations which were found to be informative for interpreting the corresponding minimal images.

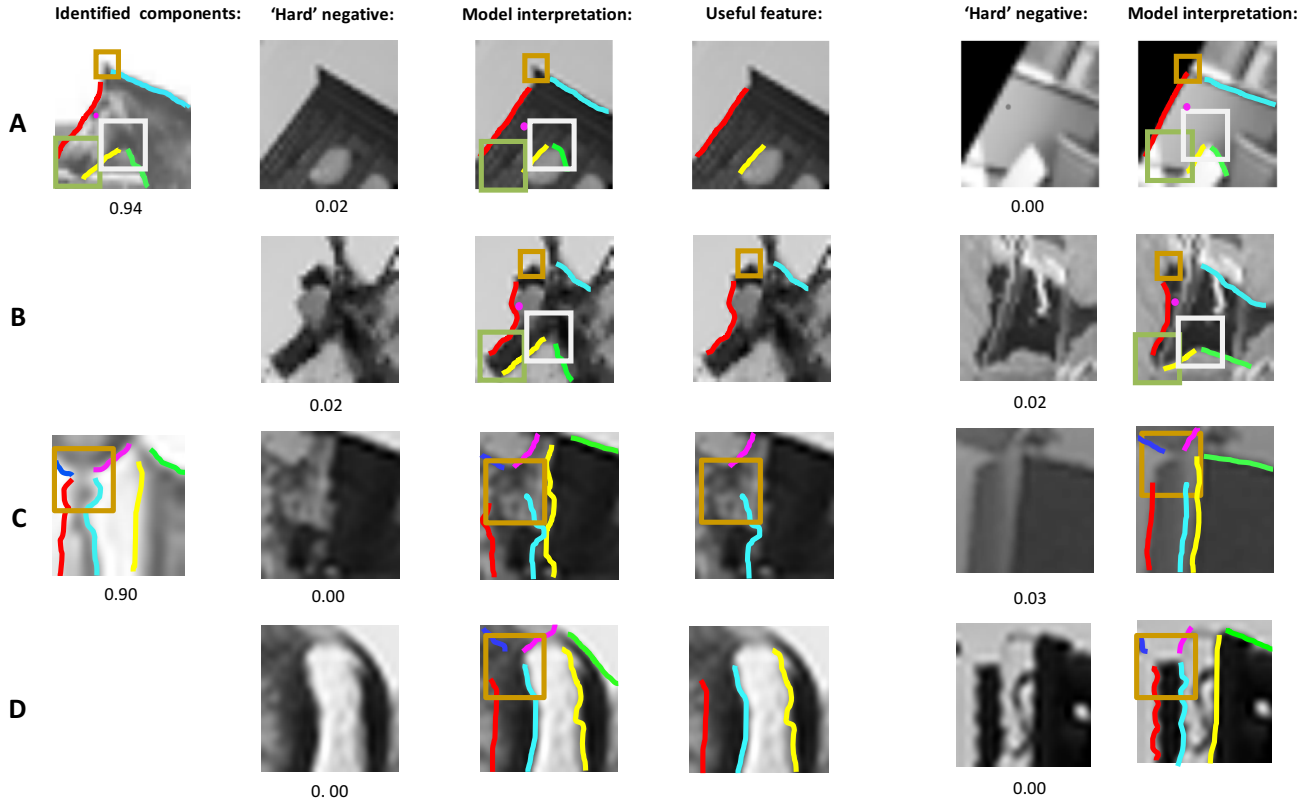
### **4.3 Useful relations from ‘hard’ negative examples**

In addition to the sub-minimal images test discussed above, which compared images from the same class, a complementary source for identifying useful relations for full interpretation is a comparison of minimal configurations with ‘hard’ non-class examples, which are difficult in the sense that they are confusable with true class examples by current computational models (a deep net model, Simonyan & Zisserman, 2015, and a human recognition model, Serre et al., 2007). Such a comparison can identify components and relations that are informative for human recognition and interpretation, but are missing from current models. We describe next how hard-negative examples were generated and how they were used to identify useful relations for interpretation.

To identify hard negative examples for a given minimal image, we trained a deep CNN-based classifier (Simonyan & Zisserman, 2015) using 120 examples of the minimal image (details in Sec. 5), and a large set of negative examples (200,000 local regions cropped and rescaled from various non-class images). We then applied the classifier on a validation set (equal in size to the training set), and we retained the 4000 non-class image regions with the highest detection scores. These are the hard-negative examples, used in the search for informative relations. Similar to the use of sub-minimal images described above, the search proceeds along the following steps.

We start with the ‘basic’ interpretation model as defined in Sec. 4.2 and iterate over the following procedure:

- i) Keep the k hard-negative images that received the highest interpretation score (since images later required MTurk tests, we used the limit k=40).
- ii) Confirm (using MTurk testing) that these negative examples are not confusable for human observers. (Examples that were also difficult for humans were



**Figure 8:** Useful relations for interpretation extracted from ‘hard’ negative examples. Columns show (left to right): minimal images with their human interpretation, non-class examples with high detection score with their human recognition rate, interpretation applied to the negative example by the model. Differences in components or relations are identified and evaluated, see text. Column 4 shows relations found to be informative for the interpretation model. They include: high straightness of two contours, typical of man-made objects (in A), connectedness of two contours through the ear region (in B), connectedness of two contours through a tie knot region (in C), coherent texture between the two shirt parts, see text (in D). The identified relations were used to reject hard negatives, examples in the last two columns.

removed from the set in practice, no more than 2 examples were removed at this stage).

- iii) Compare the interpretation produced by the model for the images collected in (i), to human annotations of the corresponding minimal image examples. As in Sec. 4.2, identify components or relations (from the list of candidate relations) present in the positive examples but not in the hard negatives.
- iv) For each such a relation, test its contribution to the interpretation model by the difference in random forest classification with and without this addition, as in Sec. 4.2.
- v) Once relations from all hard negative images were tested, and the contributing subset was added to the relations set, train a new version of interpretation model

and repeat the search from step (i), to discover additional relations from hard negatives to the new version.

We iterated this procedure until no new contributing relations were found (at most 3 iterations were needed per class).

Fig. 8 illustrates examples of hard negatives discovered and used to identify informative relations, and the process of finding these relations. Examples include ‘highly-straight’ contours (typical for man-made objects) in the horse head (e.g., the red and yellow contours in Fig. 8A), the connectedness of horse head contours through the ear region (red and cyan contours in Fig. 8B), sharp corners at the tie knot’s (cyan and magenta contours, connected inside the brown square in Fig. 8C), and coherent visual texture (or intensity level) between the two shirt parts (the area that is left to the red contour and the area that is bounded by the contours in cyan and yellow in Fig. 8D).

	<i>Relation Operands</i>	<i>Description</i>		<i>Relation Operands</i>	<i>Description</i>
1	All primitives	<b>Location and relative location:</b> for all primitives, and for all pairs of primitives in the structure.	8	Contour, Contour	<b>Length ratio</b> between two contours
2	Point	<b>Strength of intensity maxima/minima,</b> center-surround filter responses at a point location.	9	Contour, Contour	<b>Parallelism</b> between two contours
3	Contour	<b>Deviation from line/circular arc:</b> in particular for man-made objects.	10	Region, Region	<b>Coherent visual appearance</b> similar appearance/texture features in region i and in region j
4	Contour	<b>Visual appearance along contour</b> distribution of visual appearance/texture features along contour.	11	Contour, Point	<b>Cover of a point by a contour:</b> if a contour i covers a point j. For ‘cover’ refer to appendix C.
5	Region	<b>Visual appearance inside a region</b> distribution of visual appearance/texture features in a region	12	Contour, Region	<b>Contour ends in a region:</b> if a contour i ends in a region j.
6	Contour, Contour	<b>Relative location of contour endings:</b> between endings of two different contours	13	Point, Region	<b>Containment:</b> if point i is inside region j
7	Contour, Contour	<b>Continuity:</b> smooth continuation between two given contour endings.	14	Contour, Contour, Region	<b>Contour Bridging:</b> Testing whether two disconnected contour elements can be bridged (linked in the edge map).

**Table 1.** Relations that were found informative for the learning process, by the method and criterion discussed in Sec. 4.2 and 4.3. See implementation details for relation procedures in Appendix C.

#### 4.4 The final set of relations

The final set of relations, obtained by comparing MIRCs to both sub-MIRCs and hard negatives, includes unary relations (properties), binary relations, and relations among three or more primitives. Relations in the set are composed of basic relations as listed in Sec. 4.2, augmented with candidate relations which proved to contribute to the recognition and interpretation accuracy by the computational experiments in Sec. 4.2 and 4.3. Relations in the library range from low-complexity ones such as computing relative location between primitives, to higher complexity procedures such as computing the continuity, bridging, or parallelism of contours. Table 1 lists relations with the highest contribution, as measured in Sec. 4.2 and 4.3. Technical details for implementing the relation procedures are discussed in Appendix C.

## 5. Experimental evaluation

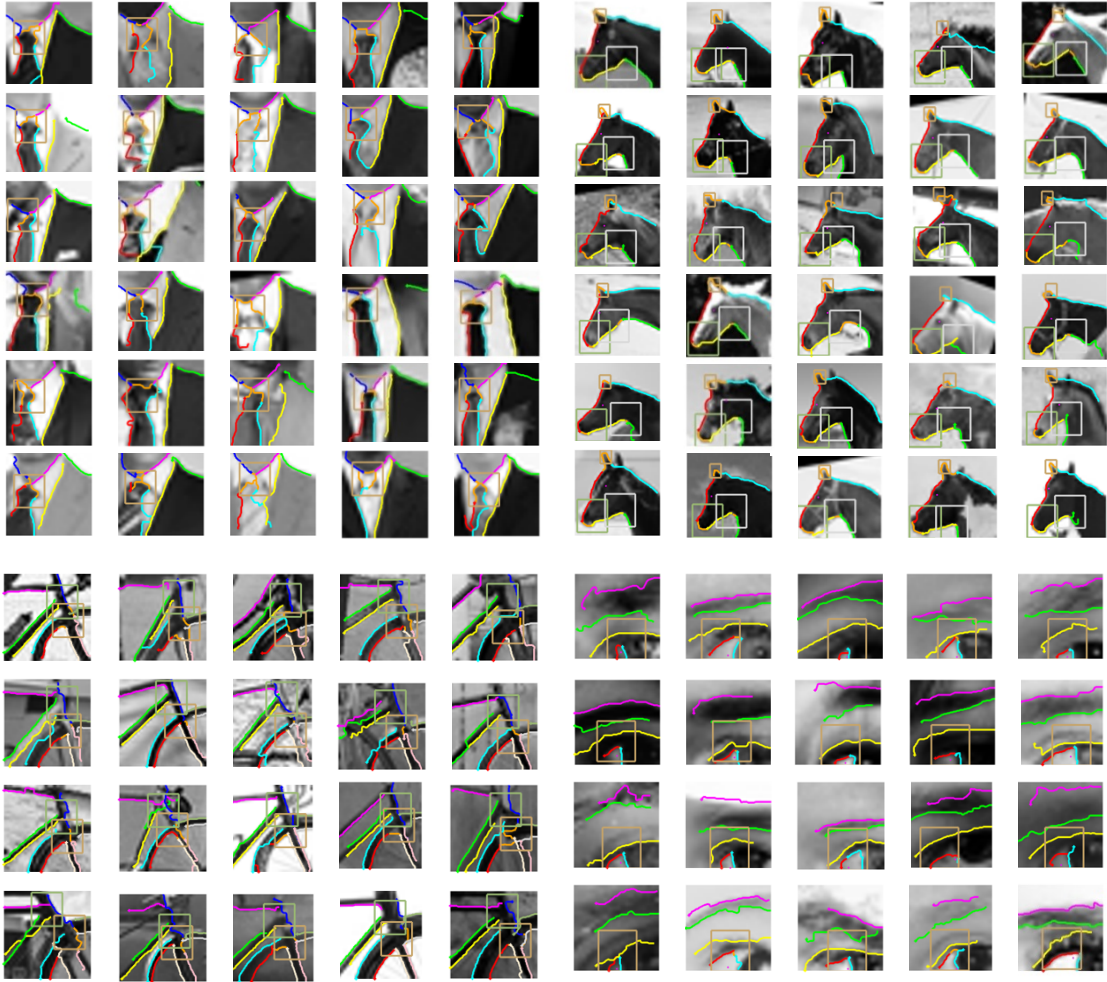
So far we have identified the useful components and relations when tested individually. We next combined all of them in the full interpretation model (as described in Sec. 3) and tested its performance. The full set of relations for the trained model was composed of the relations listed in Table 1. To evaluate the full interpretation model, we performed experiments to assess (i) the interpretation correctness on novel images, (ii) the ability of the interpretation model to predict human recognition at the level of minimal image, and (iii) the contribution of informative relations included in the model to human recognition, using modified minimal images.

Training of the model was obtained as described in Sec. 3, with annotated examples of minimal images, and non-class (negative) examples. To get positive class examples for the minimal image we wanted to model (e.g., 'horse-head'), we collected full-object images from known data sets (Flicker, Google images, ImageNet), and manually extracted from each image a local region at the position and size similar to the discovered minimal image (Ullman et al., 2016). The minimal image examples used for training were in slightly higher resolution than the minimal images found in Ullman et al., 2016 (image resolution was increased by factor of 1.5), since we found that using this scale during training improved the model results when applied to novel images.

To have ground truth for the interpretation, two human subjects provided annotation of the set of primitives (e.g., Fig. 4A) for all examples (one annotator used for ground truth, the other for measuring consistency, details in Appendix A). Negative (non-class) examples for training were collected automatically from cropped windows in non-class images at similar size to the minimal image. To get hard negative examples, we trained a deep CNN classifier (Simonyan & Zisserman, 2015), as described in Sec. 4.3, and collected images that received high recognition scores. We next turn to describe our three testing procedures, in Sec. 5.1-5.3 below.

### **5.1 Comparing model output to human interpretation**

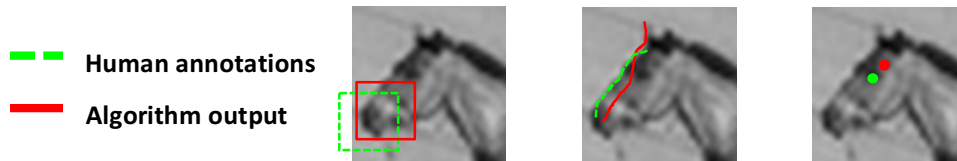
The interpretations produced by the model were compared with the ground truth annotations supplied by the human annotators. Since the model is novel in terms of



**Figure 9:** Interpretation results for minimal images belonging to (clockwise) a horse-head, a man in a suit, an eye, and a bike. (cf. Fig. 4).

producing full interpretation, it cannot be compared directly with any existing alternative models. However, we made our set of annotations publicly available, and the current model provides a baseline to also evaluate future results. To assess the role of the compound relations derived in Sec. 4.2 and 4.3, we compared results from two versions of our model, which differed in the relations included in the model: one using only relations based on local appearance, location, and displacement (termed the *basic set* below, indexed 1,4,5 in Table 1), and a second, using the full set of relations in Table 1 (termed the *compound set* below).

Fig. 9 shows examples of the interpretations produced by the model for novel test images. To assess the interpretations, we matched the model output to human annotations for multiple examples. Our training set contained 120 positive examples, and 25,000



**Figure 10:** *Quantitative evaluation of the model interpretation results. We compared interpretation results to human annotations based on the Jaccard measure similarity criteria: for regions, contours, and points (see Appendix D for details).*

negative examples for each interpretation model. Our test set contained 480 examples for the horse head minimal image (Fig. 4A), 330 examples for the man-in-suit minimal image (Fig. 4B), and 120 of the eye (Fig. 4D) and the bike (Fig. 4E) minimal images. We automatically matched the ground truth annotated primitives to the interpretation output by the so-called Jaccard index, (Tan et. al., 2006), which is a commonly used similarity measure for comparing automatic detection results (high Jaccard means similar interpretations). This index compares the similarity of two regions, by the area of the regions' intersection divided by area of their union, and was adapted to compare the accuracy of detecting region, contour, and point primitives, as illustrated in Fig. 10, and explained in more details in Appendix D. Table 2 shows results for the basic and full relation sets, as well as agreement between different human annotators, which can serve as an upper bound for comparing interpretation performance. Interpretation using the compound relations was significantly closer to the 'ground truth' human interpretation compared with the use of basic set of relations ( $P < 4.99 \times 10^{-11}$  for all primitives in 4 classes,  $n=33$ , one-tailed paired  $t$  test). However, the agreement between the model and ground truth interpretations was still lower than the agreement between different human interpretations ( $P < 1.14 \times 10^{-13}$  for all primitives in 4 classes,  $n=33$ , one-tailed paired  $t$  test).

## 5.2 Interpretation for predicting minimal and sub-minimal images

The link between interpretation and recognition, as discussed in Sec. 1.3, suggests that the interpretation score (which is a part of the model output) may be used as a part of the human recognition process at the minimal image level. In particular, it

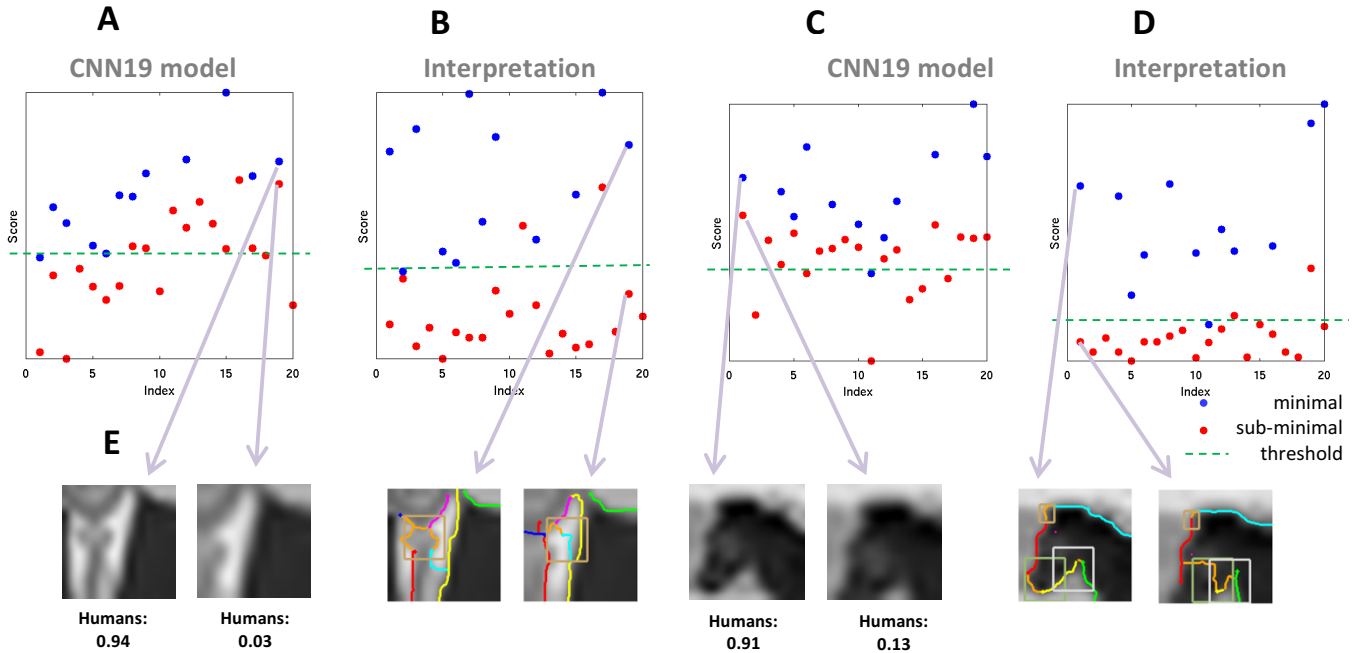
	<i>Basic</i>	<i>Compound</i>	<i>Humans</i>		<i>Basic</i>	<i>Compound</i>	<i>Humans</i>
<b>Horse-head</b>				<b>Man-In-Suit</b>			
Ear region	0.11	0.37	0.60	Knot region	0.62	0.66	0.74
Mouth region	0.69	0.76	0.85	Left tie contour	0.48	0.55	0.72
Neck region	0.55	0.68	0.74	Right tie contour	0.47	0.53	0.72
Upper head contour	0.44	0.69	0.84	Suit-shirt contour	0.64	0.73	0.83
Mane contour	0.34	0.61	0.79	Shoulder contour	0.50	0.63	0.66
Lower head contour	0.46	0.66	0.79	Left neck contour	0.49	0.65	0.84
Lower neck contour	0.32	0.63	0.74	Right neck contour	0.39	0.49	0.77
Eye point	0.29	0.49	0.60	<b>All primitives</b>	<b>0.51</b>	<b>0.61</b>	<b>0.75</b>
<b>All primitives</b>	<b>0.40</b>	<b>0.61</b>	<b>0.75</b>				
<b>Eye</b>				<b>Bike</b>			
Iris region	0.39	0.56	0.79	Fork region	0.72	0.73	0.80
Lower lid contour	0.47	0.62	0.73	Tire lower contour (left side)	0.68	0.75	0.86
Cornea contour	0.33	0.60	0.81	Tire lower contour (right side)	0.62	0.75	0.90
Upper lid contour	0.41	0.64	0.74	Bottom tube's upper contour	0.59	0.74	0.86
Lower eyebrow contour	0.51	0.64	0.83	Bottom tube's lower contour	0.54	0.70	0.87
Upper eyebrow contour	0.45	0.51	0.81	Top tube's lower contour	0.36	0.43	0.84
Sclera point	0.56	0.54	0.79	Head tube contour	0.49	0.60	0.85
<b>All primitives</b>	<b>0.44</b>	<b>0.59</b>	<b>0.78</b>	Tire upper contour(right side)	0.50	0.62	0.81
				Tire lower contour(right side)	0.53	0.57	0.78
				Fork left contour	0.60	0.68	0.82
				Fork right contour	0.59	0.71	0.83
				<b>All primitives</b>	<b>0.56</b>	<b>0.66</b>	<b>0.84</b>

**Table 2.** Accuracy of the interpretation results, comparing the basic model, compound model, and human annotators. Accuracy is measured by the average Jaccard index between the model interpretation and ground truth supplied by human annotations. For comparison, human accuracy is measured by the agreement, measured by the Jaccard index, between the human annotators.

is interesting to compare the interpretation scores for minimal and sub-minimal images, to assess the usefulness of interpretation for recognition.

In human perception, there is a sharp drop in recognition rates at the minimal image level: a small change to the image can have drastic effects on recognition rate (Sec. 1.2, above, Ullman et al., 2016). This sharp drop was not reproduced by computational models of recognition, and it therefore becomes of interest to examine whether the internal interpretation of minimal image may provide a basis for this perceptual sensitivity. It is possible, for example, that even small changes to a minimal image could disrupt the presence of key elements and their relations. To test this possibility, we measured the gap between human recognition rates for minimal and sub-minimal images (via MTurk search on new image examples) and compared it to the gap predicted by two models: the current interpretation model, and a classifier based on deep convolutional networks (very-deep CNN, Simonyan & Zisserman, 2015), trained on minimal image examples. Our test set included 12 examples of minimal images and 20 examples of sub-minimal images for each of two minimal image categories: the horse-head (Fig. 4A) and

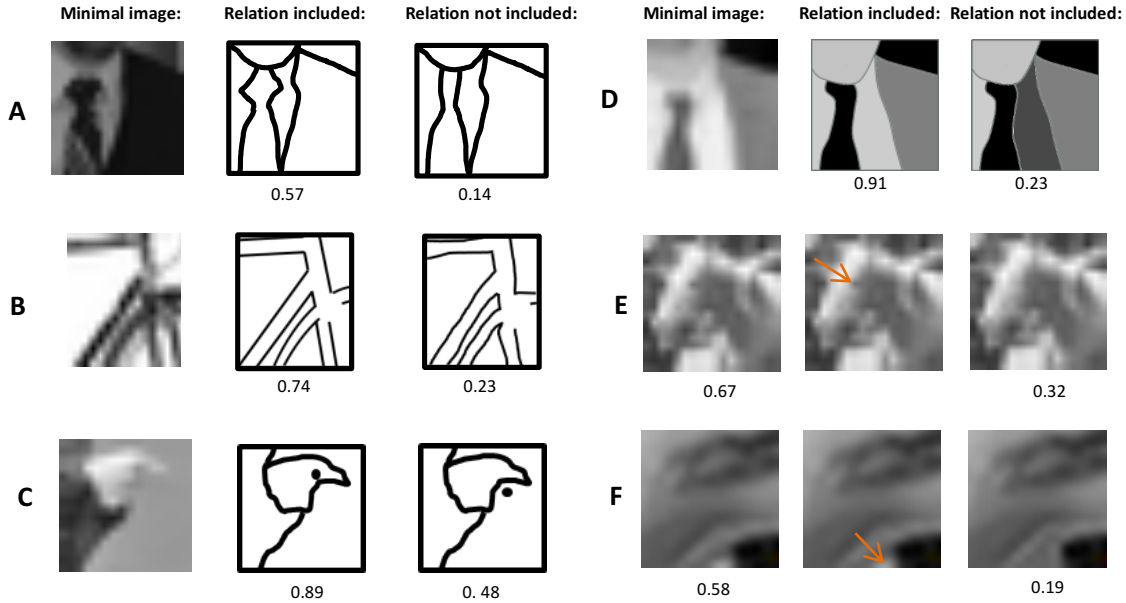




**Figure 11:** Recognition of minimal and sub-minimal images. The recognition scores of two models were compared against human recognition: a 19-layer feed-forward CNN classifier, and the interpretation model. Columns show minimal and sub-minimal pairs (a single minimal image can have more than one sub-minimal image) of the man-in-suit (in **A,B**), and the horse head (in **C,D**). The upper plots show the recognition scores of the CNN and interpretation models. Green dashed line represents the human recognition rate threshold. The separation of minimal and sub-minimal images by interpretation increases, and is closer to human recognition gap than separation by the deep CNN (**E**). Examples of minimal and sub-minimal pairs and their interpretation. The interpretation of the minimal images is more accurate compared with the sub-minimal ones. The gray arrows point to the corresponding score of each image by the interpretation and CNN models. Notice the scores for the minimal and sub-minimal images are similar by CNN, and different by interpretation.

man-in-suit (Fig. 4B). The average gap measured between human recognition rates for minimal images and for sub-minimal images was 0.75 for the horse head, 0.74 for man-in-suit. This sharp gap in human recognition at the minimal image level was compared with the computational models as described next.

To compute the recall gap of models, the model’s classification score was compared against an acceptance threshold, and scores above threshold were considered true detections. For each model, we set the acceptance threshold to match the human recognition rate. For example, for the man-in-suit, the average human recognition across all 12 examples was 0.88, and the model threshold was set so that 11/12 examples will be accepted (see Fig. 11A). Recognition rate for the sub-minimal images was then derived from the fraction of sub-minimal images exceeding the threshold, and the difference in recognition rates defines the model’s recognition gap.



**Figure 12:** 'Intervention': Testing informative relations via transformed minimal images. (A-C). Rendering sketches from images (D). Creating  $k$ -color cartoons. (E,F). Re-coloring a small set of pixels ( $\leq 4$ , pointed by the red arrow) at the same color of their neighboring pixels. In a transformed image, a relation is removed to test its predicted role in human perception. Relations tested: sharp curvature in the tie contour (in A), high contour straightness (in B), containment of a point in bounded contours (in C), coherent color/texture in the two parts (in D), minimum intensity (in E), and maximum intensity (in F).

The CNN model was trained prior to testing on 120 examples of minimal images for each class, and 200,000 non-class examples. The model's scores for the minimal and sub-minimal images on the test sets are shown in Fig. 11A,C. The gaps computed for the horse-head and man-in-suit were 0.20, and 0.37, respectively, both considerably smaller than the human recognition gap. The second model tested the interpretation trained as in Sec. 5.1. Interpretation scores are shown in Fig. 11B,D, along with the interpretation examples of minimal and sub-minimal pair from each category. The average interpretation gap was 0.75 for the horse-head and 0.76 for the man-in-suit, closely similar to the gaps measured for humans. The differences in recognition gap between the CNN and interpretation models were highly significant ( $P < 2.44 \times 10^{-4}$  for horse-head,  $P < 5.7 \times 10^{-3}$  for man-in-suit,  $n=20$ , Fisher's exact test). The difference is likely to arise because the interpretation model incorporates class-specific properties and relations that are not included in the CNN model. We discuss this difference further in Sec. 6.3, 6.4 below.

### 5.3 Testing predicted relations via intervention on minimal images

The interpretation model includes informative relations between components, which were identified using the data sets of sub-minimal images and hard negative. The model predicts that disrupting these relations should reduce the ability of human observers to recognize and interpret minimal images. To further verify the role of these relations, we used direct intervention (Pearl, 2009) on minimal images, testing whether removing specific relations from the minimal image will decrease human recognition. For this purpose, we created transformed versions of the minimal images, in which specific relations were selectively manipulated. The transformed versions were then tested psychophysically via the MTurk.

The transformations applied to minimal images included rendering sketches, including rendering  $k$ -color cartoons ( $k \leq 5$ ), and re-coloring a small set of pixels (number of re-colored pixels  $\leq 4$ ), examples in Fig 12. To create sketches, we traced contours of the original MIRC image, either manually as in Fig. 12 A, B (right column), and C, or semi-manually using straight lines, as in Fig. 12 B, middle column. Cartoon sketches are similar, but using a small number of grey-levels ( $\leq 5$ ) for the regions (12D). Re-coloring images were done with interactive graphics design tools (Irfan, Photoshop). For all sketches, we kept all contours or segments in the minimal image that are used as primitives in the interpretation model, and verified that the sketched images were still recognizable (e.g., Fig. A-D, middle column).

In the sketch images, a specific contour or a region can be selectively modified, with minimal or no change to other image parts. We created a modified version for each sketch, where selected contours or regions were changed based on the tested property or relation (e.g., Fig. A-D, right column). Since we know how a relation is computed in the model, we can change contours or regions such that this relation will no longer be detected. We then tested whether the specific disruption of a single relation will cause a significant drop in MIRC recognition as predicted by the model.

The tested relations were taken from the set of the most informative relations in the relations set (Table 1). For each tested relation, we first applied a manipulation which

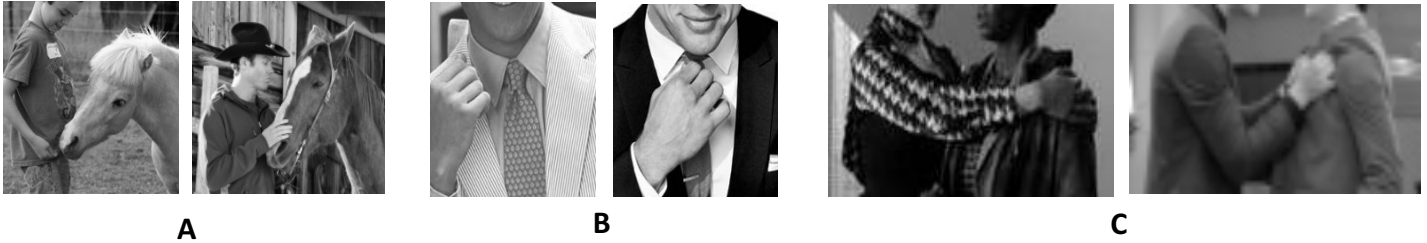
removes the relation from the model relations vector (the computed  $R_C$ ) while keeping the rest of the relations intact (the model can provide interpretation for both natural and sketched minimal images).

Each relation was tested using five different pairs of manipulated and non-manipulated versions, and the average human recognition drop for each relation was measured. Example results are shown in Fig. 12. Fig. 12A-D used sketches from minimal images. The sketched versions eliminated specific relations in the representations: sharp curvature (12A, the tie knot, cf. Fig. 8C), high straightness measure (12B, bike contours, cf. Fig. 8A), containment of a point in region (12C, bird's eye) and the coherent appearance (in intensity or texture) between two regions (12D, cf. Fig. 8D). In Fig. 12E-F, a local change was introduced to disrupt the model property of minimal (12E) or maximal (12F) local intensity. The change was induced by re-coloring 3-4 pixels, to match the average intensity of their neighboring pixels.

For all tested relations in Fig 12, the manipulation resulted in a significant drop in human recognition rate. (For example, Fig. 12A, 5 image pairs, average drop = 0.41,  $P < 2.46 \times 10^{-4}$ ,  $n=5$ , one-tailed paired  $t$  test. In similar one-tailed paired  $t$  tests for Fig. B-F,  $P < 0.0052$  for all cases). In summary, the results show a sharp drop in recognition following intervention to eliminate a relation predicted by the model to be highly informative for the interpretation of the relevant minimal image. This agreement between the model and human recognition supports the proposed role of the tested relations in human recognition and interpretation of minimal images.

## **6. Discussion and implications**

In this work, we described a model for local image interpretation, applied to minimal recognizable images. The ultimate goal of full image interpretation is to recognize meaningful semantic components anywhere in the image, but we used minimal images for development and testing of the model for two reasons. First, local interpretation reduces the number of components and the complexity of the model, and second, using a data set of minimal and sub-minimal images is useful for identifying informative components and relations which play a part in the interpretation process.



**Figure 13.** Examples of fine interpretation in recognizing human actions and interactions. **(A).** Recognizing petting vs. feeding a horse depends on the exact location of the human hand on the horse muzzle. **(B).** Whether the hand is touching the knot or not, determines the action of ‘fixing a tie’. **(C).** The hands contact locations provide important cues for recognizing a ‘hug’ interaction between the agents.

The interpretation model was shown to produce reliable interpretation of local image regions. It also helps to explain the sharp drop in recognition between minimal and sub-minimal images, which is characteristic of human observers, but not reproduced by current bottom-up computational models. It will be interesting to further test in the future the agreement between human recognition errors of difficult images and errors made by recognition models, with and without an interpretation stage.

Similar to other cognitive and computational models, interpretation is defined in the model in terms of a local structure, composed of components, properties, and relations. Our empirical testing of properties and relations proposed in past studies, showed that a number of them contributed to the performance of the model (Table 1). In comparison, restricting the relations to relative displacements between components (relations 1, 4 and 5 in Table 1), which are commonly used in computational models, proved insufficient for reliable interpretation (Sec. 5.1). Consistent with this computational evidence, a subset of the relations used by the model were found to directly affect human recognition, as human recognition of modified minimal images, to exclude selected relations, dropped significantly. Taken together, the role of the components and relations incorporated in the interpretation model is supported by three complementary sources of evidence: their contribution to correct interpretation by the model, the effect they have on the sharp difference in recognition between MIRC and sub-MIRC, and the effects of their selective elimination from minimal images on human recognition of these images.

Future work in modeling the interpretation process should go beyond the interpretation of local regions discussed in this study, towards the interpretation of full, natural images. The interpretation of full images is likely to be goal-directed, namely,

providing detailed interpretation of regions of interest, rather than uniformly across the image. Minimal images, at multiple scales, can provide a natural starting point for the fuller interpretation process, because they can be reliably recognized and interpreted on their own, independent of the surrounding context, and can subsequently help in further disambiguation and interpretation of nearby regions.

### **6.1 Detailed interpretation for complex visual tasks**

Full interpretation of semantic components at the level produced by the current model can play a useful role for extracting meaning from complex configurations, arising in tasks such as recognizing actions or social interactions between agents. The reason is that the exact meaning of an image may depend on fine localization of object parts and the relations between relevant parts, as illustrated in Fig. 13. Relatively little work has been done to date on modeling the recognition of complex interactions between agents and objects, or between agents. It will be of interest to extend in the future the current work, to study the role of detailed image interpretation in the recognition of complex actions and agents' interactions.

### **6.2 Learning relations**

In the current model, relations between components of the local interpretation are used to identify the correct structure. There are two main questions regarding the relations used for the purpose of interpretation. The first is the full set of relations that are useful for the task, and the second is identifying informative relations for a particular local structure (e.g., horse-head). Since the set of so-called 'basic' relations proved insufficient, we evaluated a larger set of relations, using minimal, sub-minimal, and difficult non-class images. The resulting set is not necessarily complete, and future studies may identify additional relevant relations. In terms of the human visual systems, such relations could be in part pre-existing in the visual system, and in part learned from visual experience. Regarding the identification of informative relations for a novel class of images, the approach in the model was to examine the full set of possible relations, and identify the informative ones using positive and negative examples, where the negative examples came from high-scoring non-class examples. It will be of interest to examine in the future the possibility of replacing this search by network learning models, based on

positive and negative examples, but without using an explicit set of possible relations. The issue of unsupervised learning of semantic components is left for future studies, we only note that some components may be learned based on their independent motion within the image (e.g. an eye or mouth within a face), or based on points of contact between an agent and an object (such as a cup-handle or door-knob).

### 6.3 Interpretation and Top-Down processing

Our model suggests that the relations required for a detailed interpretation are in part considerably more complex than spatial relations used in current recognition models (Sec. 2). Furthermore, the experimental results show that the relations used for interpretation are often class-specific, in the sense that the most informative relations for

<i>Horse-head</i>	<i>Man-in-suit</i>	<i>Eye</i>	<i>Bike</i>
<b>Intensity minimum</b> (at the eye point)	<b>Contour appearance</b> (along the tie)	<b>Deviation from circular</b> (lid upper contour)	<b>Parallelism</b> (tube contours)
<b>Contour Bridging</b> of the mane and mouth upper contours	<b>Region appearance</b> (suit region)	<b>Cover of point by contour</b> (sclera by lid contour)	<b>Continuity</b> (tire upper contours)
<b>Contour Bridging</b> (at the mouth)	<b>Contour ending in region</b> (tie contour in knot region)	<b>Relative contour endings</b> (lower lid and the iris contours)	<b>Region appearance</b> (wheel region)

**Table 3.** Top 3 informative relations found for the different class models of minimal images

the interpretation of a given class often depend on the class. This is illustrated in Table 3, which shows the most informative relations found by the model for the interpretation of 4 different classes of minimal images. Since the subsets of informative relations are class-dependent, it will be computationally efficient to compute the more complex relations selectively, in a class-specific manner, rather than computing all possible relations for all candidate classes. In such a scheme, the interpretation process will be naturally divided into two main stages. The first is a bottom-up recognition stage, similar to current feed-forward models. This stage will lead to the activation of one or several objects classes, but without detailed object interpretation. The activated classes will then trigger a top-down process for the computation of further class-specific components and relations required for a detailed interpretation. The interpretation will be also used for validation of the activated classes in the first stage, by rejecting bottom-up detections which do not have the expected interpretation. Future studies could explore this two-stage proposal further by psychophysical and physiological methods. For example, since the accurate recognition of minimal images depends in the model on its internal interpretation, the

top-down component predicts that the reliable recognition and interpretation of minimal images will be a relatively slow process compared with a single feed-forward pass.

#### **6.4 Interpretation by network models**

Recognition models based on deep convolutional networks have shown to produce high-accuracy results in object classification and promising results in related tasks, such as segmentation (e.g., Long et al., 2015). The current model combines network algorithms with other methods to extract complex relations and identify the final structure. Similar combinations have been used recently by other models that extract complex structures (e.g. human pose, Chen and Yuille, 2014, combining CNN with a subsequent conditional random field stage; Lake et al., 2015, in the domain of written characters). We found that existing feed-forward network models have limited accuracy when applied to the interpretation of minimal images. Our evaluation trained a recent semantic segmentation network (Long et al., 2015) to identify interpretation components of minimal images. The accuracy of the resulting interpretation was closer to the ‘basic’ model compared with the full version of the current model (Sec. 5.1). It is plausible, however, that extended network models, such as models using recurrence and memory, will cope more successfully with local interpretation. It will be of interest to develop such models in future work, and compare network structures that prove successful for local interpretation, with aspects of cortical circuitry in the visual system, e.g. in terms of using recurrent and feedback connectivity.

**Acknowledgements:** We thank Daniel Harari for sharing psychophysics data and for help with data collection. This work was supported by ERC Advanced Grant “Digital Baby”, the EU’s Horizon 2020 research and innovation program under grant agreement No. 720270, Israeli Science Foundation grant 320/16, and the Center for Brains, Minds and Machines, funded by NSF Science and Technology Centers Award CCF-1231216.



## **Appendix A. Psychophysics experimental methods**

### **A.1. Labeling all semantic components in a minimal image:**

This experiment was used for identifying semantic elements which humans can consistently identify in minimal images. Subjects ( $n=30$ ) were presented with a minimal image in which a red arrow pointed to a location in the image (e.g., the horse eye, or the center of the mouth region), and were asked to name the indicated location. Similarly, a contour was marked in red on the image, and subjects produced two labels for the two sides of the contours (e.g., tie and shirt). In both cases subjects were asked to also name the object they saw in the image (without the markings). To map the scope of ‘full’ human-level interpretation, we put the red arrows and contours at multiple image locations, and tested their consistent labeling. We considered a recognized component if more than 50% of human tags were consistent. Presentation time was unlimited, and the subjects responded by typing the labels. All experiments and procedures were approved by the institutional review boards of the Weizmann Institute of Science, Rehovot, Israel. All participants gave informed consent before starting the experiments.

### **A.2. Annotating point, contour, and region components in minimal image examples:**

Subjects ( $N=2$ ) were presented with examples of the semantic components found for a given minimal image by the experiment in Appendix A.1 (annotated by points, contours, and regions, as in Fig. 4B), and were asked to produce similar annotations in novel examples. Annotators were given partially overlapping sets of examples from each class, which together covered the complete training and testing sets. At least 50 examples from each class were annotated by two different subjects, and were used to test consistency in human annotations (see Table 2). The annotated images served as the ‘ground truth’ in evaluating the performance of the interpretation model (Sec. 5.1, and Table 2).

## **Appendix B. The learning model and procedure**

### **B.1. A structured learning model based on random forest**

The problem of local interpretation can be viewed as an instance of so-called ‘structured learning’ (e.g. Shalev-Shwartz & Ben-David, 2014). As described in Sec. 3.2, given a structure  $S_C$  consisting of a set of primitives  $P_C$ , and a vector  $R_C$  of relations between

them, we wish to learn an interpretation function  $f_S$  that finds the structure  $S_C$  (denoted  $S$  below for simplicity) in an image  $I$

$$f_S(I) = \pi$$

where  $I$  is the object image, and  $\pi$  is not just a class label, but a full assignment, which is in our case a mapping between components in the structure  $S$  and points, contours, and regions in the image  $I$ . We refer to  $\pi$  as an ‘assignment’, since it assigns to any primitive in the model  $S$ , a counterpart in the image, identified by  $\pi_i$ .  $\pi$  is then a vector  $\pi = [\pi_1, \pi_2, \dots, \pi_N]$ , where  $N$  is the number of primitives in the model  $S$ . For example, if the minimal image is the horse head, and the primitive set in  $S$  includes, among others, the horse eye (primitive index = 1, type = point), and the horse mane contour (primitive index = 5, type = contour), then,  $\pi_1$  is a point in  $I$  assigned to the horse’s eye, and  $\pi_5$  is a contour in  $I$  assigned to the horse’s mane.

It is common to express the function  $f_S$  using a (learnable) *scoring* function  $g(I, \pi; w)$ , which measures the compatibility between the model structure  $S$ , and the corresponding structure identified in the image. The additional variables  $w$  are parameters of the interpretation function, described below.  $f_S(I)$  then takes the form:

$$1) f_S(I; w) = \underset{\pi}{\operatorname{argmax}} \{ g(I, \pi; w) \},$$

namely, given an image  $I$  (with parameters  $w$  already fixed), find the assignment  $\pi$  into  $I$  that has the highest compatibility with the model structure  $S$ . The goal of the function  $f_S$  is then to find the configuration of elements within the image  $I$ , which is as compatible as possible with the model structure  $S$ .

The function  $g$  in our interpretation measures the compatibility between properties and relations specified by the structure  $S$  of the model, and the same properties and relations computed for the corresponding image elements, identified by the assignment  $\pi$ . This compatibility is computed as follows. Given an assignment  $\pi$  of the model primitives to the image  $I$ , we denote the results of measuring all the model relations in the specific image  $I$  by the vector  $\phi_S(\pi, I)$ . Following the example in Sec. 3.2, position 3 in the vector  $\phi_S(\pi, I)$  could be ‘true’ (or 1), indicating that primitive 5 is

contained in primitive 7, and position 4 could be 0.9 indicating the degree of straightness for primitive 2.

The relations vector  $\phi_S(I, \pi)$  is then used to measure the compatibility of the image structure with the model structure. This is obtained in our model by a random forest algorithm (Amit & Geman, 1996; Breiman, 2001), which is learned from training examples. A random forest is a non-linear model composed of a set of classification trees:

$$\{t_1, t_2, \dots, t_j, \dots\},$$

where  $t_j$  is the  $j$ -th tree in a forest. The parameters  $w$  in this model (in the definition of  $f_S$  and  $g$ ) are the queries in the tree nodes, and a standard learning procedure for random forests (Breiman, 2001) is used to set these parameters based on training examples. Each tree is applied to the relations vector  $\phi_S(I, \pi)$  to produce a decision whether the given assignment, represented by  $\pi$ , is consistent with a class structure or not (i.e., the relations vector  $\phi_S(I, \pi)$  was classified as 1 or 0). Finally, the function  $g$  returns the average of all tree votes:

$$2) \quad g(I, \pi, w) = \frac{1}{K} \sum_{j=1}^K t_j(\phi_S(I, \pi)),$$

where  $K$  is the number of trees in the forest. The assignment we seek is the one that maximizes the value of this expression, an effective optimization search is described in Appendix B.2 below.

The random forest algorithm also provides a method for evaluating the individual contribution of each of the relations in the model to the learning process. This is obtained by removing a single relation in  $\phi_S(I, \pi)$  in all vectors in our data, and measuring the interpretation correctness by the random forest with and without this relation. (Referred to as the ‘Out of bag estimate’ for strength of random forest features, Brieman, 2001). We used this method in Sec. 4 to derive a set of relations, which are useful for the interpretation process. ‘Informative’ relations in Sec. 4 are measured by the difference in the performance of the model (the interpretation correctness) with and without the relation in question.

## B.2. Detecting primitive candidates and an effective optimization search

We describe below how we implemented the calculation of  $f_S$  (Eq. 2), namely, derive the best assignment  $\pi$  for a given image  $I$ . Our implementation includes two stages: (i) finding  $k$  ( $k = 10$ ) candidates for each primitive in  $S$ , and (ii) seeking the candidate combination that forms the best assignment. In more details, the two stages are

- i. *Primitive candidates*: For primitives of type 'point' and 'region' we find candidates in a bottom-up manner: for 'point', we consider all pixels in the minimal image, and for 'region' we take all image windows of the region size in a 'sliding window' search. For type 'contour' we find the candidates in a top-down manner, as follows: We project ground truth annotated contours on an edge map (Arbelaez et al., 2011), to get edge contour fragments similar in their location and shape to the ground truth ones. We then used connected pairs of fragments (by Kovese's edge linking toolbox, 2000) as candidates for the contour primitive. We rank all candidates of point, contour, and region types by their unary relations in  $R_C$ , and keep the top  $k$  for each primitive. Unary relations used for ranking include visual appearance of regions and contours (relations 4 and 5 in Table 1), and intensity minima/maxima of points (relation 2 in Table 1).
- ii. *Finding the best assignment*: Given an image  $I$ , a trained model  $w$ , and a set of candidates for each primitive in  $P_C$ , we run over different configurations of candidates in a coordinate descent manner (Bertsekas, 1999). We start with a random configuration, and then optimize successively one candidate at a time. Specifically, the procedure is:
  - 1) Start with a random configuration of primitive candidates  $\pi = [\pi_1, \pi_2, \dots, \pi_i, \dots, \pi_N]$ .
  - 2) Repeat until  $g$  converges:  
For each primitive  $i$ , go over all candidates  $\pi_i'$  and update:
    - $\pi' = [\pi_1, \pi_2, \dots, \pi_i', \dots, \pi_N]$
    - $\pi \leftarrow \operatorname{argmax}\{g(I, \pi, w), g(I, \pi', w)\}$
  - 3) Return  $\pi$ .

Such a procedure is guaranteed to converge to a local optimum (Bertsekas, 1999; a

similar optimization search was used for Hopfield networks, Hopfield, 1982).

Experimentally, because the search space in minimal images is limited due to small number of primitives, 3 initiations of the procedure were usually sufficient to get good convergence.

### **Appendix C. Details of computing relation**

Table 1 in Sec. 4.3 contains the set of relations used in our models. In this appendix we add technical details about the computational procedures for computing the different relations. For all procedures described here,  $x, y$  represent the coordinates of the image plane. All procedures were implemented in MATLAB, code is available from the authors.

**Containment:** Given a pixel point  $[x, y]$  and a set of pixels comprising a region  $R$ , we return true if the point is in the region, i.e.,  $[x, y] \in R$ .  $R$  can be either a single region primitive, or a region bounded by two (or more) contour primitives.

**Contour ends in a region:** Given an end point pixel  $[x_1^c, y_1^c]$  of a contour  $C$ , and a set of region pixels  $R$ , we return ‘true’ if the end point is in the region, i.e.,  $[x_1^c, y_1^c] \in R$ .

**Parallelism:** Given two contours,  $C_a$  and  $C_b$ , we compute a binary mask  $M$ :

$$M(x, y) = 1 \quad \text{if } [x, y] \in C_a \text{ or } [x, y] \in C_b$$

$$M(x, y) = 0 \quad \text{otherwise.}$$

We then compute the distance transform map (Maurer et al., 2003) for  $M$ , denoted  $DT\{M\}$ , followed by a non-maxima suppression to get the ridges  $R$  of  $DT\{M\}$ . The ridges  $R$  is the set of pixels that are at equal distance from both contours. The two contours are considered parallel if the variance of  $R$  is close to zero. We exclude cases where the size of  $R$  is small. We thus return ‘true’ if  $Var[R] < \varepsilon$ , where  $\varepsilon$  is a threshold close to zero (we chose empirically  $\varepsilon = 0.2$ ).

**Continuity of contours:** Given a contour  $C_a$  with one of its endings:  $[x_1^{C_a}, y_1^{C_a}]$ , and a contour  $C_b$  with one of its ending:  $[x_1^{C_b}, y_1^{C_b}]$ , we estimate the local orientations at the endings, namely  $\theta_1^{C_a}$  and  $\theta_1^{C_b}$ , and use them to compute the completion path between

$[x_1^{C_a}, y_1^{C_a}, \theta_1^{C_a}]$  and  $[x_1^{C_b}, y_1^{C_b}, \theta_1^{C_b}]$  (Ben-Yosef & Ben-Shahar, 2012). We consider ‘good continuation’ between the two contours if the completed path does not contain inflection points. We return ‘true’ if the number of inflection points in the path equals to zero.

**Bridging contours:** Given a contour  $C_a$  with one of its endings  $[x_1^{C_a}, y_1^{C_a}]$ , a contour  $C_b$  with one of its endings  $[x_1^{C_b}, y_1^{C_b}]$ , and the image  $I$  from which the two contours are extracted, we test for an image contour connecting them. We compute the UCM map (an edge map, Arbelaez et al., 2011) for  $I$  and define a graph  $G = \langle V, E \rangle$ , where  $V$  is the set all pixels in the UCM map, namely

$$v_i \in V : UCM(v_i) > \tau,$$

$\tau$  is a UCM threshold ( $\tau=0.1$ ), and  $E$  is a set of weighted edges. An edge  $e \in E$  is put for each pair of pixels in  $V$  that are immediate image neighbors. The weight of an edge  $e = \{v_i, v_j\}$  is defined as the difference in UCM levels between pixels:

$$w(e) = UCM(v_j) - UCM(v_i)$$

(The graph  $G$  is computed in a pre-process stage.) We return the shortest weighted path in  $G$  (if exists) between  $[x_1^{C_a}, y_1^{C_a}]$  and  $[x_1^{C_b}, y_1^{C_b}]$ .

The bridging procedure was also extended in two versions: (i) finding a path in  $G$  that is the most consistent with the ways contours  $C_a$  and  $C_b$  are connected in positive train images, and (ii) finding a path in  $G$  that is constrained to pass through region primitive.

**Visual appearance inside regions or along contours:** Given a candidate image region  $R_a$  for a primitive  $R$  in the model, we match the distribution of the visual appearance features in  $R_a$  and in the training examples of  $R$ . Visual appearance features were ‘visual words’ features (Arandjelovic & Zisserman, 2013), and deep neural network features (Long et al., 2015). For a contour candidate, we used a similar match of visual appearance features, this time along a thin region surrounding the contour.

**Coherent visual appearance:** Given two candidate image regions  $R_a$  and  $R_b$ , we match the distribution of the visual appearance features in these two regions. Visual appearance

features were ‘visual words’ features (Arandjelovic & Zisserman, 2013), and deep neural network features (Long et al., 2015).  $R_a$  or  $R_b$  could be either a single region primitive, or a region bounded by two (or more) contour primitives.

**Cover of a point by a contour:** Given a pixel point  $[x, y]$  and a contour  $C$ , we project  $C$  on the X axis of the image plane, and return ‘true’ if  $x$  is within the range of projection. We composed procedures for different directions of cover, namely for a contour covers a point from top or from bottom. Similar ‘cover’ procedures were also for the Y axis.

#### **Appendix D. Evaluating similarity between elements: points, contours, and regions**

This process was used for evaluating the correctness of the interpretation produced by the model (Sec. 5.1). For two regions, A and B, the standard Jaccard measure ( $|A \cap B|/|A \cup B|$ , Tan et al., 2006) was used. For two points, we construct a small square region around each point (size of 12% of the minimal image), and then evaluate the Jaccard index of these regions. For two contours, we used a simple extension of the Jaccard index to contours, by extending the contours into tube shaped regions (tube width was 4% of the minimal image) and measure the Jaccard index between these regions.

#### **References**

1. Attneave, F. (1954). Some informational aspects of visual perception. *Psychological review*, 61(3), 183-193.
2. Amit, Y., & Geman, D. (1997). Shape quantization and recognition with randomized trees. *Neural computation*, 9(7), 1545-1588.
3. Azizpour, H., & Laptev, I. (2012). Object detection using strongly-supervised deformable part models. *Proceedings of the European Conference on Computer Vision*, 836-849.
4. Andriluka, M., Pishchulin, L., Gehler, P., & Schiele, B. (2014). 2d human pose estimation: New benchmark and state of the art analysis. *Proceedings of the IEEE Conference on computer Vision and Pattern Recognition*, 3686-3693.
5. Arandjelovic, R., & Zisserman, A. (2013). All about VLAD. *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 1578-1585.
6. Arbelaez, P., Maire, M., Fowlkes, C., & Malik, J. (2011). Contour detection and hierarchical image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 33(5), 898-916.
7. Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5-32.
8. Biederman, I. (1987). Recognition-by-components: a theory of human image understanding. *Psychological review*, 94(2), 115-147.
9. Bertsekas, D. P. (1999). *Nonlinear programming* (pp. 1-60). Belmont: Athena scientific.
10. Ben-Yosef, G., & Ben-Shahar, O. (2012). A tangent bundle theory for visual curve completion. *IEEE transactions on pattern analysis and machine intelligence*, 34(7), 1263-1280.
11. Ben-Yosef, G., Assif, L., Harari, D., & Ullman, S. (2015). A model for full local image interpretation. *Proceedings of the annual meeting of the Cognitive Science Society*, 220-225.

12. Chen, X., & Yuille, A. L. (2014). Articulated pose estimation by a graphical model with image dependent pairwise relations. *Advances in Neural Information Processing Systems*, 1736-1744.
13. Csurka, G., Dance, C., Fan, L., Willamowski, J., & Bray, C. (2004). Visual categorization with bags of keypoints. *Proceedings of the workshop on statistical learning in computer vision, European Conference on Computer Vision* 1(1), 1-2.
14. Dalal, N., & Triggs, B. (2005). Histograms of oriented gradients for human detection. *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 886-893.
15. Epshtein, B., Lifshitz, I., & Ullman, S. (2008). Image interpretation by a single bottom-up top-down cycle. *Proceedings of the National Academy of Sciences*, 105(38), 14298-14303.
16. Everingham, M., Van Gool, L., Williams, C. K., Winn, J., & Zisserman, A. (2010). The Pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2), 303-338.
17. Felzenszwalb, P. F., Girshick, R. B., McAllester, D., & Ramanan, D. (2010). Object detection with discriminatively trained part-based models. *IEEE transactions on pattern analysis and machine intelligence*, 32(9), 1627-1645.
18. Feldman, J. (2007). Formation of visual “objects” in the early computation of spatial relations. *Perception & Psychophysics*, 69(5), 816-827.
19. Foster, D. H., Simmons, D. R., & Cook, M. J. (1993). The cue for contour-curvature discrimination. *Vision research*, 33(3), 329-341.
20. Field, D. J., Hayes, A., & Hess, R. F. (1993). Contour integration by the human visual system: evidence for a local “association field”. *Vision research*, 33(2), 173-193.
21. Girshick, R., Donahue, J., Darrell, T., & Malik, J. (2014). Rich feature hierarchies for accurate object detection and semantic segmentation. *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 580-587.
22. Girshick, R., Iandola, F., Darrell, T., & Malik, J. (2015). Deformable part models are convolutional neural networks. *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 437-446.
23. Hinton, G. E. (2007). Learning multiple layers of representation. *Trends in cognitive sciences*, 11(10), 428-434.
24. Deng, J., Berg, A., Satheesh, S., Su, H., Khosla, A., & Fei-Fei, L. (2012). *Imagenet large scale visual recognition competition 2012 (ILSVRC2012)*. net.org/challenges/LSVRC/2012/.
25. Joachims, T., Hofmann, T., Yue, Y., & Yu, C. N. (2009). Predicting structured objects with support vector machines. *Communications of the ACM*, 52(11), 97-104.
26. Jolicoeur, P., Ullman, S., & Mackay, M. (1986). Curve tracing: A possible basic operation in the perception of spatial relations. *Memory & Cognition*, 14(2), 129-140.
27. Kanizsa, G. (1979). *Organization in vision: Essays on Gestalt perception*. Praeger Publishers.
28. Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 1097-1105.
29. Kovesi, P. D. (2000). MATLAB and Octave functions for computer vision and image processing. Online: <http://www.csse.uwa.edu.au/~pk/Research/MatlabFns/#match>.
30. Lake, B. M., Salakhutdinov, R., & Tenenbaum, J. B. (2015). Human-level concept learning through probabilistic program induction. *Science*, 350(6266), 1332-1338.
31. LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436-444.
32. Lafferty, J., McCallum, A., & Pereira, F. (2001). Conditional random fields: Probabilistic models for segmenting and labeling sequence data. *Proceedings of the International Conference on Machine Learning*, 282-289.
33. Long, J., Shelhamer, E., & Darrell, T. (2015). Fully convolutional networks for semantic segmentation. *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 3431-3440.
34. Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., ... & Zitnick, C. L. (2014). Microsoft coco: Common objects in context. *Proceedings of the European Conference on Computer Vision*, 740-755.
35. Machilsen, B., Pauwels, M., & Wagemans, J. (2009). The role of vertical mirror symmetry in visual shape detection. *Journal of Vision*, 9(12), 11-11.
36. Maurer, C. R., Qi, R., & Raghavan, V. (2003). A linear time algorithm for computing exact Euclidean distance transforms of binary images in arbitrary dimensions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(2), 265-270.



37. Palmer, S., & Rock, I. (1994). Rethinking perceptual organization: The role of uniform connectedness. *Psychonomic bulletin & review*, 1(1), 29-55.
38. Pasupathy, A., & Connor, C. E. (1999). Responses to contour features in macaque area V4. *Journal of Neurophysiology*, 82(5), 2490-2502.
39. Pearl, J. (2009). *Causality*. Cambridge university press.
40. Parent, P., & Zucker, S. W. (1989). Trace inference, curvature consistency, and curve detection. *IEEE Transactions on pattern analysis and machine intelligence*, 11(8), 823-839.
41. Roelfsema, P. R., Lamme, V. A., & Spekreijse, H. (1998). Object-based attention in the primary visual cortex of the macaque monkey. *Nature*, 395(6700), 376-381.
42. Riesenhuber, M., & Poggio, T. (1999). Hierarchical models of object recognition in cortex. *Nature neuroscience*, 2(11), 1019-1025.
43. Saarela, T. P., Sayim, B., Westheimer, G., & Herzog, M. H. (2009). Global stimulus configuration modulates crowding. *Journal of Vision*, 9(2), 1-11.
44. Serre, T., Oliva, A., & Poggio, T. (2007). A feedforward architecture accounts for rapid categorization. *Proceedings of the National Academy of Sciences*, 104(15), 6424-6429.
45. Simonyan K, Zisserman A (2015) Very deep convolutional networks for large-scale image recognition. *International Conference on Learning Representations*.
46. Tan, P. N., Steinbach, M., & Kumar, V. (2006). *Introduction to data mining* (Vol. 1). Boston: Pearson Addison Wesley.
47. Tompson, J. J., Jain, A., LeCun, Y., & Bregler, C. (2014). Joint training of a convolutional network and a graphical model for human pose estimation. *Advances in neural information processing systems*, 1799-1807.
48. Ullman, S., Assif, L., Fetaya, E., & Harari, D. (2016). Atoms of recognition in human and computer vision. *Proceedings of the National Academy of Sciences*, 113(10), 2744-2749.
49. Ullman, S. (1984). Visual routines. *Cognition*, 18(1-3), 97-159.
50. Yamins, D. L., Hong, H., Cadieu, C. F., Solomon, E. A., Seibert, D., & DiCarlo, J. J. (2014). Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the National Academy of Sciences*, 111(23), 8619-8624.
51. Vedaldi, A., Fulkerson, B.: VLFeat: An open and portable library of computer vision algorithms. <http://www.vlfeat.org/> (2008)
52. Vedaldi, A., Mahendran, S., Tsogkas, S., Maji, S., Girshick, R., Kannala, J., ... & Taskar, B. (2014). Understanding objects in detail with fine-grained attributes. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3622-3629.
53. Wertheimer, M. (1923). Untersuchungen zur Lehre von der Gestalt. II. *Psychological Research*, 4(1), 301-350.
54. Westheimer, G., Crist, R. E., Gorski, L., & Gilbert, C. D. (2001). Configuration specificity in bisection acuity. *Vision research*, 41(9), 1133-1138.
55. Wei, S. E., Ramakrishna, V., Kanade, T., & Sheikh, Y. (2016). Convolutional Pose Machines. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 4724-4732.
56. Yang, S., Luo, P., Loy, C. C., & Tang, X. (2015). From facial parts responses to face detection: A deep learning approach. *Proceedings of the IEEE International Conference on Computer Vision*, 3676-3684.