

Intentionomy: a Dataset and Study towards Human Intent Understanding

Menglin Jia^{1,2} Zuxuan Wu^{2,3} Austin Reiter² Claire Cardie¹ Serge Belongie¹ Ser-Nam Lim²
¹Cornell University ²Facebook AI ³Fudan University

Abstract

An image is worth a thousand words, conveying information that goes beyond the mere visual content therein. In this paper, we study the intent behind social media images with an aim to analyze how visual information can facilitate recognition of human intent. Towards this goal, we introduce an intent dataset, Intentionomy, comprising 14K images covering a wide range of everyday scenes. These images are manually annotated with 28 intent categories derived from a social psychology taxonomy. We then systematically study whether, and to what extent, commonly used visual information, i.e., object and context, contribute to human motive understanding. Based on our findings, we conduct further study to quantify the effect of attending to object and context classes as well as textual information in the form of hashtags when training an intent classifier. Our results quantitatively and qualitatively shed light on how visual and textual information can produce observable effects when predicting intent.¹

1. Introduction

Why do we post images on social media platforms like Facebook or Instagram? Are we expressing our feelings to friends and family? Are we seeking to entertain a wide audience? Or is it purely out of habit, or perhaps out of fear of missing out? Images on social media embody more than their explicit visual information, and they tend to be persuasive in commercial ads and even manipulative in the context of political campaigns. Therefore, in the deluge of social media, understanding the intent behind images is critical, especially for tasks like fighting fake news and misinformation [33, 61] on social platforms.

However, understanding human intent behind images from a computer vision point of view is particularly challenging, since it goes beyond standard visual recognition—predicting a set of stuff and thing categories that physically exist in images such as objects [49, 26, 90, 84, 36] and scenes [64, 95, 106]. Additionally, it is a psychological task [72] inherent to human cognition and behavior. It is

¹Intentionomy project page: github.com/kmnp/intentionomy

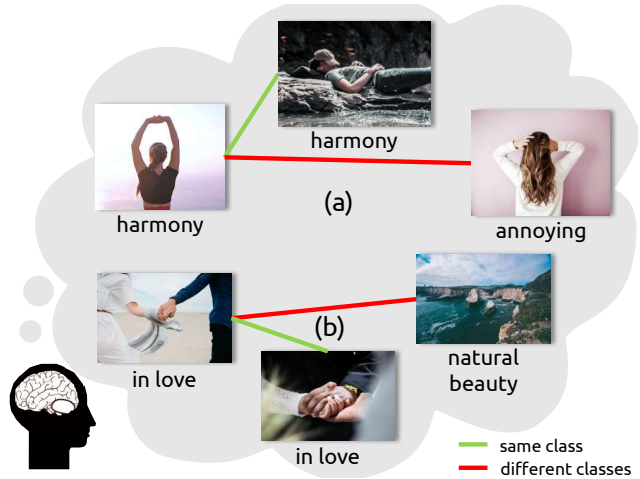


Figure 1. Intent behind images: while (b) shows that the visual motif of holding hands aligns with the common intent of “in love”, (a) illustrates that similarity based on visual appearance alone often would lead to an incorrect match with respect to intent.

similar in spirit to visual commonsense reasoning [100, 74] to derive an answer conditioned on the objects and scenes present in images. In certain cases, intent can be inferred rather directly from representative objects and scenes. For example, a couple holding hands or making a heart symbol clearly have the same motive “in love”(Fig. 1(b)). However, the mapping from visual cue to intent is not always one-to-one. Fig. 1(a) shows that two images with completely different contents (a girl facing the ocean vs. a person relaxing on a rocky surface, with face covered) can represent the same intent (“harmony”). This goes beyond the usual variability (pose, color, illumination, and other nuisances) traditionally addressed in object recognition pipelines [25, 67]. This brings us to the question: *are objects and their image context sufficient for recognizing the intent behind images?*

In this paper, we introduce a human intent dataset, *Intentionomy*, containing 14K images that are manually annotated with 28 intent categories, organized in a hierarchy by psychology experts. To investigate the intangible and subtle connection between visual content and intent, we present a systematic study to evaluate how the perfor-

mance of intent recognition changes as a function of (a) the amount of object/context information; (b) the properties of object/context, including geometry, resolution and texture. Our study suggests that (1) different intent categories rely on different sets of objects and scenes for recognition; (2) however, for some classes that we observed to have large intra-class variations, visual content provides negligible boost to the performance. Furthermore, our study also reveals that attending to relevant object and scene classes brings beneficial effects for recognizing intent.

In light of this, we further study a multimodal framework for intent recognition. In particular, given an intent category, the framework localizes, in a weakly-supervised manner, salient regions in images that are important for recognizing the class-of-interest. These discovered regions are further reinforced during training using a localization loss to guide the network to focus. In addition, we leverage hashtags as a modality complementary to visual information. We demonstrate through extensive evaluations that properly ingesting visual and textual information helps to boost the performance of intent prediction significantly.

Our work makes the following key contributions: (1) A novel dataset of 14,455 high-quality images, each labeled with one or more human *perceived* intent. This dataset, which we call *Intentionomy*, offers a total of 28 intent labels supported by a systematic social psychological taxonomy [72] proposed by experts; (2) A systematic study to show how commonly used object and context information, as well as textual information, contribute to intent recognition; (3) We introduce a framework with the help of weakly-supervised localization and an auxiliary hashtag modality that is able to narrow the gap between human and machine understanding of images.

2. Related Work

Prior work on intent recognition has focused on communicative intents in different contexts. Joo *et al.* [38] define 9 dimensions of persuasive intents of a politician implied through a photo (*e.g.*, trustworthy). Other works [39, 68, 32, 77] also focus on persuasive intents in political images. Additional related work includes image and video advertisement understanding including topics, sentiment and intent [34, 103, 97], or the motivation behind the actions of people from images [60, 88]. Understanding intent is also a key component in persuasive dialogue systems [65, 19, 99, 91]. In this work, we focus on the behavior of the people who post on social media websites. While a large body of work [46, 3, 47, 4, 73, 71] exists that study the motivations behind the usage of social media, relatively much fewer work exists in the area of computer vision.

The most similar work in terms of understanding human motive in social media is from [45], which introduces a multimodal dataset to understand the document intent in Insta-

gram posts. However, we differentiate our work in terms of goals and methods: (1) we emphasize “visual intent” rather than “textual intent”, meaning that we study human motive mainly based on the perceived motives behind images rather than textual data; (2) we systematically analyze how objects and context contribute to the recognition of human motives in the social media domain; (3) our dataset contains more fine-grained categories (28 classes in total) with nine super-categories compared to 8 categories from [45].

Our study on the relationship between intent and content is inspired by [104], which studied the effect of context for object recognition. Other works also proposed context-aware models in various tasks such as object recognition and detection [80, 30, 81, 56, 12, 54, 31, 52, 5, 50, 10], scene classification [96, 11], semantic segmentation [96, 54], scene graph recognition [102], visual question answering [75]. Our work utilizes both object- and scene-level information to distinguish between different intent classes.

3. Intentionomy Dataset

Images Our dataset is built up of free-licensing high-resolution photos from the website Unsplash². We sample images with common keywords that are similar to social media hashtags, including “people”, “happy”, *etc.* The resulting images cover a wide range of everyday life scenes (*e.g.*, from parties, vacations, and work).

Intent taxonomy The selection of intent labels is a non-trivial exercise. The labels must form a representative set of motives³ from social media posts, and it should occur with high enough frequencies in the collection of the dataset. Previous work on motive taxonomies [72] provide a solid foundation for our study. However, not all of the 161 human motives presented in [72] are suitable in the context of social media posts, or can be inferred from single-image inputs. For example, one might need background information about the person inside the image to judge if the intent is “being spontaneous”, “to be efficient”, “to be on time”. Some fine-grained motives in the taxonomy could be merged. For instance, “social group” and “close friends”, “making friends” and “having close friends”. Wherever possible, we further divide certain motives into sub-motives (“in love” and “in love with animal” for instance), for more granularity. Fig. 2 illustrates our resulting ontology in full with hierarchy information and annotated image examples.

Annotation details Amazon Mechanical Turk (MTurk) was recruited to collect labels of perceived intent by employing a similarity comparison task that we call “unsatisfactory substitutes”. We rely on the notion of “mental imagery” [78] – a quasi-perceptual experience that maps example images to a visual representation in one’s mind, along with *games*

²Unsplash Full Dataset 1.1.0

³We use *intent* and *motive* interchangeably

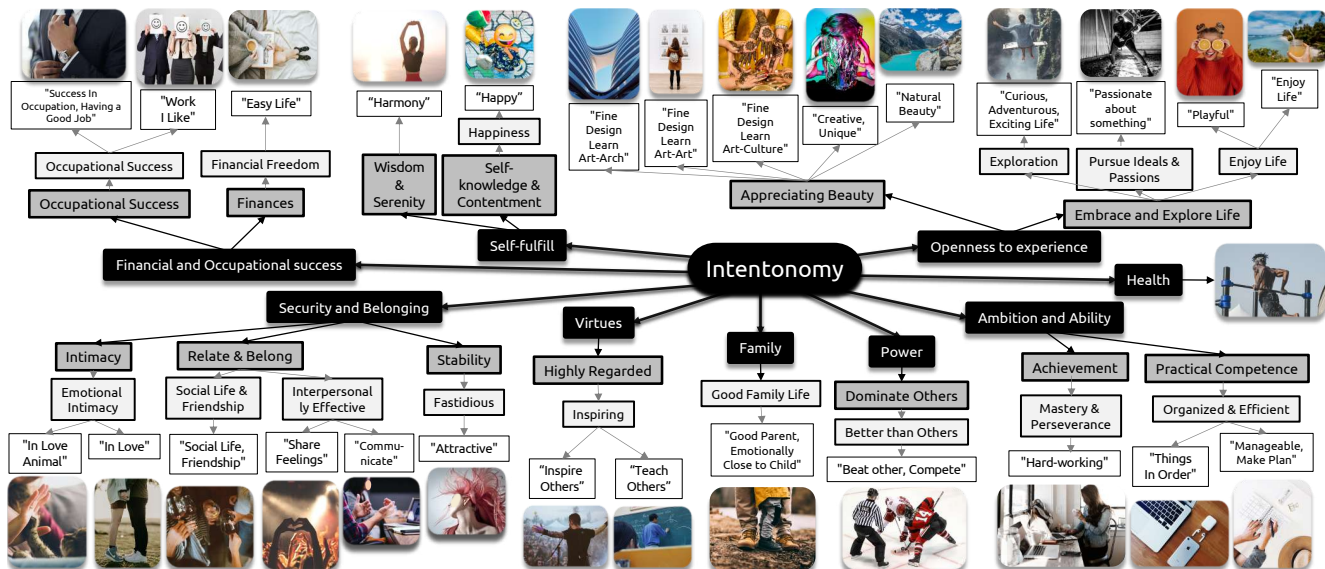


Figure 2. Ontology visualization. We select 28 labels from a general human motive taxonomy used in psychology research [72]. There are 9 super-categories in total (*in black box*), namely “virtues”, “self-fulfill”, “openness to experience”, “security and belonging”, “power”, “health”, “family”, “ambition and ability”, “financial and occupational success”. See the Appendix D for dataset statistics.



Figure 3. Annotation methods comparison. (a) A standard annotation process: given an image, choose the desired labels from a drop-down list; This approach is time-consuming and highly dependent on the expertise of annotators. (b) Our approach: similarity comparison using “unsatisfactory substitutes” so the annotators can focus on the “swapabilities” of image pairs regarding the intent. The task is to select all the images in the grid that clearly have a different intent than the reference image on the left.

with a purpose [86, 87, 15, 83] as our overall annotation approach. Fig. 3 displays the differences between a standard annotation process and ours. Due to space constraints, we leave other details in the Appendix C.

Although we have implemented strategies to ensure quality (see the Appendix C.3), we acknowledge that there are inevitable inconsistencies in our training data. Different people have different opinions of perceived intent. Prior work [83] shows that there is at least 4% error rate in popular datasets like CUB-200-2011 [89] and ImageNet [14]. Yet these datasets are still effective for computer vision research. Deep learning is robust to label noise in training set [83, 66]. To this end, we create a highly curated test set by enlisting a single domain-specific taxonomic expert to provide the annotations for both validation (val) and test

sets. In our experiments, we regard this expert’s opinions as the “gold standard,” which allows us to focus on self-consistency in val and test sets, but we acknowledge that challenges remain in terms of resolving matters of disagreement among communities of experts. In the end, Intentonomy dataset has 12,740 training, 498 val, 1217 test images. Each image contains one or multiple intent categories.

4. From Visual Content to Human Intent

Our goal is to investigate systematically how visual content within images contributes to the understanding of human intent. To this end, we disentangle the impact of visual content on intent classification by a series of controlled experiments inspired by the methodology in [104]. More specifically, we study the effect of visual content in terms of object (O) and context (C), and focus on the following fundamental aspects: (1) the amount of content information; (2) three different content properties, including geometry, resolution, and low-level texture. We then analyze the relationships between intent classes and specific things and stuff classes. Fig. 4 and 5 provide an overview of our study under different control settings to analyze how visual information affects intent recognition.

More formally, given an image I , we apply a perturbation either to its objects or context to produce a modified image: $I_x^t = f(I, t, x)$, $x \in \mathbf{X}$, $t \in \{O, C\}$, where $f(\cdot)$ indicates a transformation function as will be introduced below and \mathbf{X} is a set of positive integers defining the level of changes. The larger the value of index x , the closer the I_x^t is to the original images. We now introduce different



Figure 4. Example images with full content (far left), and image modifications used for the controlled conditions of our study.

Properties	$ \mathbf{X} $	$I_x^t = f(\cdot)$
Geometry	6	$\text{jigsaw}_{(g \times g)}(t)$, $g = 2^{5-x}, x \in [0, 5]$
Resolution	6	$\text{blur}_\sigma(t)$, $\sigma = 2^{5-x}, x \in [0, 5]$
Low-level texture features	3	$\begin{cases} \text{no t} & x = 0 \\ \mathbb{1}\{\text{texture}(t)\} & x = 1, 2 \end{cases}$

Table 1. Content properties investigation. $t \in \{\mathcal{O}, \mathcal{C}\}$.

transformation functions used to see how intent recognition performance changes based on different visual contents.

Amount of content We control the amount of object or contextual information by expanding (or decreasing for context experiments) the bounding boxes (bbox) of detected objects by e pixels:

$$I_x^t = \begin{cases} \text{bbox}^t & x = 0 \\ \text{bbox}^t \pm e & x \in [1, 7] \quad e = 2^x \\ \text{full image} & x = 8 \end{cases}$$

where $\text{bbox}^{(t \in \mathcal{O})}$ denotes the image area within the bounding box, and $\text{bbox}^{(t \in \mathcal{C})}$ is the area outside the bounding box (see two images in Fig. 4 (the third column from the right) for an example). A total of 9 variations for both objects and context are included. The larger x indicates that the larger the amount of objects or context are presented.

Content properties We also study how visual properties impact intent recognition. We analyze the effect of the following properties of \mathcal{O} and \mathcal{C} , including:

1. **geometry**: regions of objects or context are broken down to $g \times g$ tiles and randomly re-arranged (we call this operation *jigsaw*), while the other content component remains intact;
2. **resolution**: convolving the selected content component with a Gaussian function (zero-mean and various values of standard deviation σ);
3. **low-level texture features**: visual textures are constructed using image statistics [62] for the selected content component.

For all three properties, we only modify the selected regions and paste other intact content components to their original locations. Table 1 describes our method in detail.

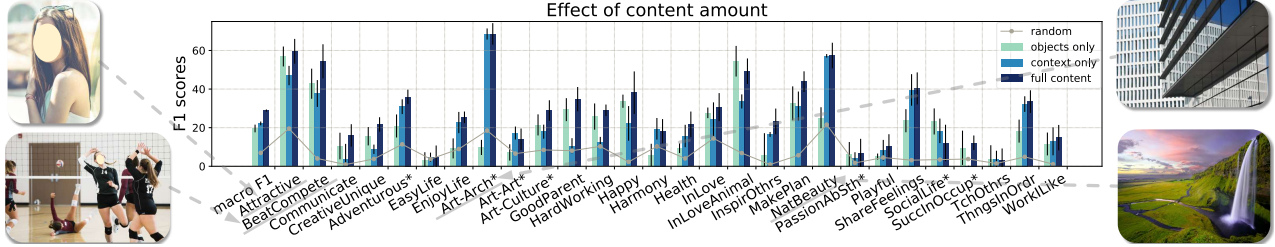
Analysis and discussion Given each transformation f , we finetune a pretrained CNN model and obtain the macro F1 score on the modified validation set. Each model is run multiple times to reduce variance. Fig. 5 shows results of the 4 experiments focusing on content size and three properties. In general, we observe a positive correlation between the amount of content and the macro F1 score. We can see in Fig. 5(a) that recognition F1 score decreases when context/object information is removed, for a majority of motive labels (e.g. “BeatCompete” and “SocialLife*”), confirming that context and objects clues are both important.

Interestingly, there are some exceptions to this trend where either objects or context, on its own, yield comparable results to the original images. For categories like “attractive” or “in love”, object information alone offers comparable F1-scores to full images. In other cases, contextual information achieves decent performance for motives like “appreciate architecture”, “natural beauty”. Such motives are usually associated with representative gestures that provide strong supervisory signals (e.g., see Fig. 1). These signals usually come from single content module, which we further demonstrate in the next subsection.

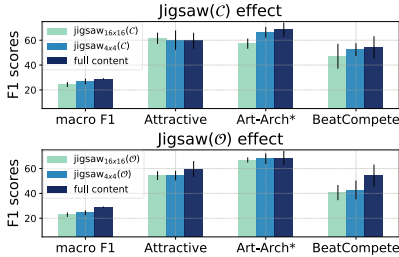
In addition, Figs. 5(b)-5(d) demonstrate how content properties affect intent recognition. We see that geometry, blurred effect, and texture features of the content component decrease the intent recognition performance. See macro F1 score and “beatCompete” in Figs. 5(b)-5(d) for example. Similar to the content size experiment previously, the impact of content properties is different for different classes. The bottom plots of Figs. 5(b)-5(d) show that “Attractive” is sensitive to object manipulation. Motives like “Art-Arch*”, on the other hand, have an opposite trend where context contributes more than objects. The recognition results are robust to object manipulation, yet sensitive to context modulation overall. These observations are further illustrative of the varying importance of objects or contextual information for different classes.

Relationship between intent and object/context classes

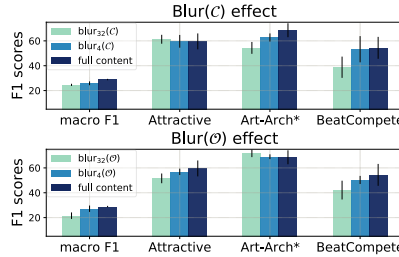
The above analysis demonstrates different intent categories



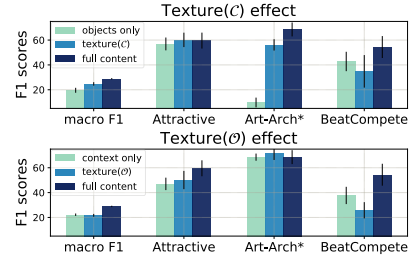
(a) Content size.



(b) Content geometry.



(c) Content resolution.



(d) Content texture.

Figure 5. A study on intent and content. Overall there are three trends among 28 classes, which are presented in Figs. 5(a)-5(d). F1 scores, including average value and standard deviation over 5 runs, and random guess results, for selected classes and selected data variations are displayed. Class names ends with “*” are abbreviated (e.g. “Art-Arch*” is short for “appreciate architecture”).

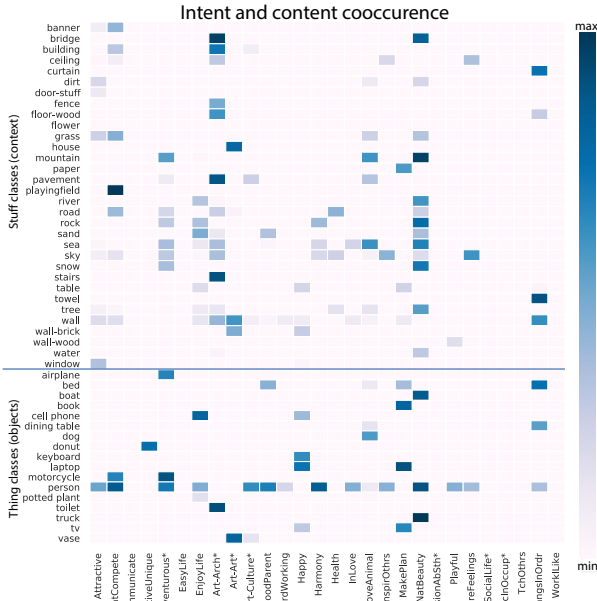


Figure 6. A visualization of Π (Eq. 1), where each entry denotes the correlation between a pair of intent and object/context class.

have different preferences on objects and/or context. We now examine whether there exists relationships between intent categories and *specific* objects/context classes.

More specifically, given an image I with a intent label m and a trained intent recognition model, we use class activation mappings [105], to produce a binary mask $CAM^b(I, m, \tau_{cam})$ (τ_{cam} is a threshold value) to represent the discriminate image regions for class m in the image. We

also feed the image to a segmentation model pretrained on the COCO Panoptic dataset [43] to obtain a binary mask $Pano(I, p, \tau_p)$ (τ_p is a threshold value) for the class p in the COCO dataset. We use the COCO Panoptic dataset [43] because it contains widely used *thing* and *stuff* categories. We then define the correlation between p and m as:

$$\Pi_{p,m} = \frac{CAM^b(I, m, \tau_{cam}) \cap Pano(I, p, \tau_p)}{Pano(I, p, \tau_p)} \quad (1)$$

Here, objects with high scores tend to be semantically meaningful for the corresponding intent categories. Fig. 6 further validates our findings in the content modulation experiments. While there are intent classes requiring both object and context, certain classes are object-oriented while others are context-oriented. Further, it can be observed that certain intent classes are also more dependent on particular object or context classes. For example, “person” is semantically meaningful for intent like “Attractive” and “inHarmony”. It is also consistent that stuff classes like “building”, “bridge” can help discriminate classes like “Architecture”.

It is worth mentioning that some intent classes (e.g. “easyLife”, “socialLife”) have no or few correlated thing or stuff classes. Indeed, the F1 score for some motive classes are comparable to random guessing (see Fig. 5(a)). We suspect that visual information only is not enough to represent the inherent visual and semantic diversity in those classes.

5. Multimodal Intent Recognition

The study in Sec. 4 demonstrates that different intent classes have different correlation with context and objects,

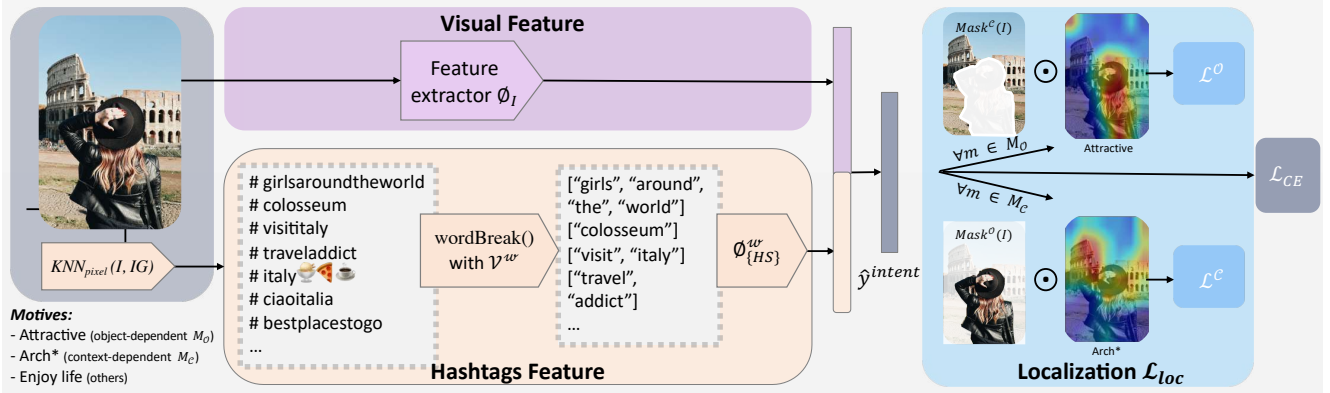


Figure 7. Method overview. Given an image I , we localize important object and context regions for an intent of interest and additionally use hashtags to complement visual information. See texts for more details.

and so using a single “one-size-fits-all” network for intent recognition is sub-optimal. To mitigate this issue, we introduce a localization loss that identifies, for each class, regions in images that are important (Sec. 5.1). In addition, as shown, visual information alone is not sufficient for predicting certain classes of intent. To compensate, we also propose to use an auxiliary channel to provide complementary semantic information (Sec. 5.2). The overall framework is presented in Fig. 7.

5.1. Object/Context localization

Since different intent classes rely on different visual content (either \mathcal{O} or \mathcal{C}), we wish to guide the network to attend to these regions when recognizing a class of interest. In particular, we first split all intent categories into 3 groups based on our study in Sec. 4: object-dependent (M_O), context-dependent (M_C), and others which depends on the entire image. We then use CAM [105] to localize salient regions in a weakly-supervised manner and minimize the overlap area between CAM and the image area that is not a region of interest (Fig. 7).

Formally, given a motive class m and an image sample I , $CAM(I, m)$ denotes the real-valued version of $CAM^b(I, m, \tau_{cam})$ (see Sec.4). Let $Mask^C(I)$ and $Mask^O(I)$ be the aggregated binary masks in image I that represent all detected thing (P_T) and stuff classes (P_S), respectively. See examples of $Mask^O(I)$ and $Mask^C(I)$ in Fig. 7. The localization loss is then defined as:

$$L^O = \sum_{m \in M_O} (CAM(I, m) \odot Mask^C(I)) \quad , \quad (2)$$

$$L^C = \sum_{m \in M_C} (CAM(I, m) \odot Mask^O(I)) \quad , \quad (3)$$

where \odot is element-wise multiplication. The final loss \mathcal{L}_{loc} is the summation of all the entries in L^O and L^C .

Note that our approach is similar to previous work that addresses contextual bias [70]. Both approaches use CAM

as weak annotations to guide training. However, our method does not require a regularization term which grounds CAMs of each category to be closer to the regions from a previously trained model. Therefore, our approach can be trained with a single pass, in an end-to-end fashion.

5.2. Hashtags as an auxiliary modality

Visual information is not sufficient for recognizing certain intent categories (see “EasyLife” in Fig. 6). To further improve intent recognition, we resort to language information as a complementary clue for improved performance. Unfortunately, images from Unsplash are not associated with any text information. We instead leverage visual similarities of the Unsplash images to a larger set of images, which do contain associated metadata that loosely describe the semantics within the images. Instagram (IG) is a social media platform that contains billions of publicly available photos, often with user-provided hashtags. This presents an opportunity to weakly relate images with vastly different visual appearances that contains similar semantic information, by means of hashtags.

In particular, we first compute regional maximum activations of convolutions features [24, 79] from the last activation map of a pretrained Resnext-50 (32x4) model (trained on ImageNet-22k [14]) on 7-days of public photos from IG as well as all the images from our intent dataset. Using these embeddings, we then perform a KNN query for each Unsplash image to retrieve the top k matching IG images for each of the images in our intent dataset. Finally, for each matching IG photo, we collect all of the associated hashtags (additional details are in the supplemental material). The collection of all matched hashtags for a given Unsplash image are represented as an unordered set HT . See Fig. 7 for examples of fetched hashtags. However, directly using hashtags are challenging because: 1) hashtags can be noisy, much like web-scale data tends to be; 2) a hashtag is usually a concatenation of several words, in-

cluding multilingual phrases and emojis (e.g. #coffeeme, #landscapephotography). There are a large amount of out of vocabulary words if one uses a pre-trained word embedding for the entire hashtag. We thus first break the hashtags down using a known dictionary of words (i.e. #coffeeme \rightarrow “coffee” “me”). Subsequently, unusual and noisy tokens/hashtags are automatically filtered out.

Formally, given HT for one image sample and a dictionary \mathcal{V} , we first segment each hashtag hs into a list of tokens based on the given vocabulary: $\text{WordBreak}(hs, \mathcal{V}) = [w]$, $w \in \mathcal{V}$, $hs \in HT$. Separated tokens of one hashtag are mapped to a dense embedding individually, and aggregated into a single representation. Next, all of the resulting hashtag representations are averaged to compute a unified feature for all hashtags associated with a single image. Finally, the hashtag features are concatenated with image features into an integrated representation for classification.

Loss function To capture the different opinions from crowd annotators, we use cross-entropy loss with soft probability, denoted as \mathcal{L}_{CE} , inspired by [51]⁴. More formally, our model computes probabilities \hat{y}^{intent} using a softmax activation, and minimizes the cross-entropy between \hat{y}^{intent} and the target distribution y^{intent} . y^{intent} is a target vector, where each position m contains the number of crowd workers who labeled the associated image to motive class m , normalized by the total number of crowd workers to indicate a probability distribution.

6. Experiments

In this section, we conduct extensive experiments to evaluate the effectiveness of different components of the multimodal framework. More specifically, we report the performance of the following approaches: (1) RANDOM, which is the success rates by random guessing; (2) VISUAL, which finetunes a standard ResNet50 model to classify motives; (3) HASHTAGS (HT), which only uses hashtags to predict intent; (4) VISUAL + HT, which combines visual information and hashtags; (5) VISUAL + \mathcal{L}_{loc} , which augments a visual model with the proposed localization loss; (6) VISUAL + \mathcal{L}_{loc} + HT, which denotes our full model. Among them, (2)-(4) are trained using the standard cross-entropy loss only, \mathcal{L}_{CE} . When the localization loss \mathcal{L}_{loc} is applied, we sum both \mathcal{L}_{loc} and \mathcal{L}_{CE} : $\mathcal{L} = \lambda\mathcal{L}_{loc} + \mathcal{L}_{CE}$. λ is a scalar to determine the contribution of each loss term. Performance are measured using Macro F1, Micro F1, and Samples F1 scores. We repeat each experiment 5 times and report the mean and standard deviation (std).

Table 2 summarizes the results. We can see that the full model achieves a 31.12 macro F1 score, outperforming the VISUAL baseline by +7.76% percent difference, as well as

⁴Similar to [51], we also tried sigmoid cross-entropy loss but obtained worse results.

Method	Macro F1	Micro F1	Samples F1
RANDOM	6.94 \pm 0.09	7.18 \pm 0.10	7.10 \pm 0.10
VISUAL	28.88 \pm 0.56	37.08 \pm 1.07	36.06 \pm 1.51
HT	19.72 \pm 0.88	29.30 \pm 1.62	31.47 \pm 1.64
VISUAL + \mathcal{L}_{loc}	30.37 \pm 0.51 (+1.49)	38.64 \pm 0.95 (+1.56)	37.41 \pm 1.51 (+1.35)
VISUAL + HT	30.32 \pm 0.62 (+1.44)	37.61 \pm 0.85 (+0.53)	38.98 \pm 1.70 (+2.92)
VISUAL + \mathcal{L}_{loc} + HT	31.12 \pm 0.63 (+2.24)	38.49 \pm 0.88 (+1.41)	38.77 \pm 1.74 (+2.71)

Table 2. Experimental results of different approaches for intent recognition measured in Micro F1, Macro F1, Samples F1 scores. (+ ·) indicate the difference comparing to VISUAL. (+ ·) in green denotes that the difference is larger than the std.

the HT baseline by +57.81% percent difference. Furthermore, compared to the VISUAL only approach, adding the localization loss improves macro F1 score by 5.16%. We also observe that visual and text information are complementary, offering 4.99% and 53.75% gain compared to visual and text only, respectively.

To better understand why \mathcal{L}_{loc} and HT improve visual only model, we break down the intent classes into different subsets based on their content dependency, i.e., object-dependent (\mathcal{O} -classes), context-dependent (\mathcal{C} -classes), and *Others* which depends on both foreground and background information; (2) difficulty, which measures how much the VISUAL outperforms achieves than the RANDOM results (“easy”, “medium” and “hard”). More details are given in the Appendix A. Table 3 summarizes the subset results.

The effectiveness of \mathcal{L}_{loc} We see from Table 3 that when adding the localization loss gains are more significant for \mathcal{O} -classes, compared to \mathcal{C} -classes and *Others*. The localization loss depends on the area of either object or context regions in the images, and the objects’ region, which is used in \mathcal{L}^C in Eq. 3, are typically small⁵. As a result, the \mathcal{L}_{loc} has no significant effect on the final score.

We also conduct a qualitative study to understand why the localization loss helps intent recognition. Results are shown in Fig. 8. We can see that the localization loss helps the model to focus on the correct region of interest for both \mathcal{O} - and \mathcal{C} -classes, especially when the image is scattered with multiple objects and scenes. Fig. 8(a) confirms that VISUAL + \mathcal{L}_{loc} works well when both object and context information are presented in the image (bottom 2 examples), or the target region of interest is small (top example). We also note in Fig. 8(b) that for images where the region of interest is located in the center, or is relatively large, both VISUAL and our method give good results.

The effectiveness of hashtags From Table 3, we observe

⁵Note that \mathcal{L}^C minimizes the overlap region between object area and the salient region (CAM).

Method	Content			Difficulty		
	\mathcal{O} -classes	\mathcal{C} -classes	Others	Easy	Medium	Hard
RANDOM	7.75 \pm 5.47	12.53 \pm 5.96	6.05 \pm 5.23	19.86 \pm 1.28	7.11 \pm 3.40	2.81 \pm 1.80
VISUAL	34.92 \pm 3.63	41.27 \pm 3.53	25.34 \pm 1.13	61.84 \pm 4.90	33.71 \pm 2.24	11.73 \pm 1.74
HT	26.96 \pm 0.80	35.15 \pm 4.18	15.43 \pm 0.87	63.58 \pm 1.79	19.68 \pm 1.70	6.63 \pm 1.43
VISUAL + \mathcal{L}_{loc}	38.82 \pm 1.95 (+3.9)	43.14 \pm 3.00 (+1.87)	25.90 \pm 1.35 (+0.56)	63.67 \pm 1.47 (+0.09)	34.72 \pm 1.26 (+1.01)	13.83 \pm 1.13 (+2.10)
VISUAL + HT	37.71 \pm 2.70 (+2.79)	42.17 \pm 3.62 (+0.90)	26.36 \pm 1.17 (+1.02)	66.67 \pm 2.12 (+4.83)	32.93 \pm 1.57 (-0.78)	15.52 \pm 0.98 (+3.79)
VISUAL + \mathcal{L}_{loc} + HT	39.82 \pm 1.56 (+4.90)	42.09 \pm 2.57 (+0.82)	26.77 \pm 1.13 (+1.43)	66.18 \pm 4.56 (+4.34)	33.86 \pm 1.08 (+0.15)	16.50 \pm 1.80 (+4.77)

Table 3. Experimental results in terms of how much object/context information intent categories need (content), and how difficult it is for VISUAL to outperforms the RANDOM results (difficulty). (+ ·) indicates the difference comparing to VISUAL. (+ ·) in green denotes that the difference is larger than the std.

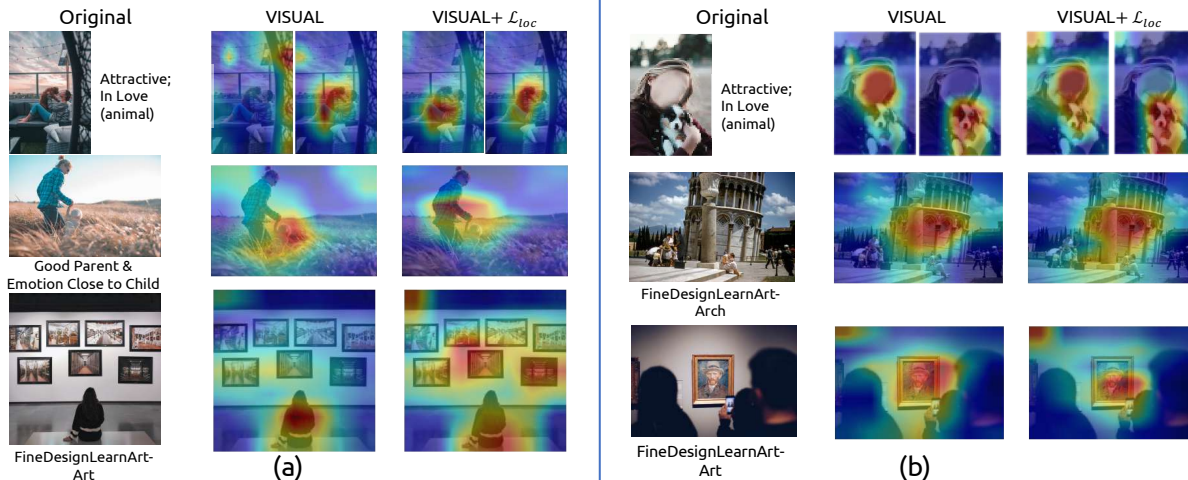


Figure 8. Analysis of the proposed localization loss. (a) VISUAL + \mathcal{L}_{loc} approach learns to isolate appropriate regions of interest, comparing to VISUAL. For example, our method learn to focus on the dog and girl respectively for “In love (animal)” and “Attractive” respectively, which are \mathcal{O} -classes. (b) Examples for which both VISUAL + \mathcal{L}_{loc} and VISUAL produce similar visualizations. Both methods focus on the correct region, which are located in center and account for a larger area of the image.

that the model using both images and hashtags outperforms the uni-modal approaches over “easy” and “hard” classes, without hurting the “medium” classes. This suggests there is value in the auxiliary information to help close the semantic gap. Therefore, our results suggest that images and hashtags do in fact complement each other in the motive recognition task. For example, #love is directly indicative of the intent label “in love”, as is #workout of “health” (see Fig. 7 for more hashtag examples). Interestingly, for “easy” classes, HT model outperforms the VISUAL model by 8.2%, however it struggles with the “medium” and “hard” classes (Table 3). This suggests that hashtags provided by users, while noisy, do still contain information about intent to some extent.

It is perhaps counter-intuitive that hashtags do not outperform visual signals entirely. While hashtags seem to capture the essence of human intents (see examples in Fig. 7), careful inspection of the fetched hashtags shows that not all hashtags are useful in practice. Obscurity and ambiguity exist, including typos, slang, inside jokes, and irrelevant information. More effective modeling of hashtags remains an open research problem.

7. Conclusion

In this work, we studied the problem of modeling human motives in social media posts. We introduced a new dataset that taps into mental imagery in a novel annotation game with a purpose to acquire labels from MTurk, and collected a rich image dataset with 28 human motives supported by a social psychological taxonomy. We conducted rigorous studies to explore the connections between content and intent. Our results show that there is still much room for improvement (for context-dependent, and hard classes for example). We therefore hope that the new Intentionomy dataset will facilitate future research to better understand the cognitive aspects of images.

Acknowledgement We thank Luke Chesser and Timothy Carbone from Unsplash for providing the images, Kimberly Wilber and Bor-chun Chen for tips and suggestions about the annotation interface and annotator management, Kevin Musgrave for the general discussion, and anonymous reviewers for their valuable feedback. This work is supported by a Facebook AI research grant awarded to Cornell University.

Appendix

The aim of our work is to investigate the complex psycho-emotional landscape hidden behind social media posts, and to lay the groundwork for the research in this domain. Such research can foster the development of systems to identify harmful posts and to reduce social media abuse and misinformation. In our work we proposed to explore human intent understanding by introducing a new image dataset along with a new annotation process. We conduct an extensive analysis on the relationship between *content* and *intent*. We also presented a framework with two complementary modules for the task. In the supplemental material, we provide the following items that shed further insight on these contributions:

- Details for reproduce our results (A);
- An extended discussion of hashtag experiments (B);
- Information about data collection process (C);
- Intentionomy data analysis (D);
- A datasheet for our motive taxonomy (E);
- Additional related work (F) and other questions regarding our work (G).

A. Experiment Details

A.1. Experimental setup

Training details To extract visual information, we use a ResNet50 [29] model which is pretrained on ImageNet [14]) as the backbone of our framework. We use Pytorch [57] to implement and train all the models on a single NVIDIA V100 GPU. We adopt standard image augmentation strategy during the training (randomly resize crop to 224×224 , horizontal flip). We use stochastic gradient descent with 0.9 momentum with batch size as 128. The learning rate is warmed up linearly from 0 to base learning rate ($1e-3$ for image only models, $5e-4$ for the rest) during the first five epochs. Since the dataset is not balanced, we follow [48, 13] to stabilize the training processing by initializing the the bias for the last linear classification layer with $b = -\log((1 - \pi) / \pi)$, where the prior probability π is set to 0.01.

Localization loss For the \mathcal{L}_{loc} , we conducted grid search for λ with the range $\{0.5, 0.1, 0.01, 0.001\}$. We set $\lambda = 0.1$ in the end, which is also consistent with the parameter used in previous work [70].

Hashtags To obtain hashtags, we index the Unsplash photos using KNN [37], and retrieved a total of 661,505 Instagram images with associated hashtags. We experiment with a range of k for the nearest neighbor search: further details

are shown in Sec. B and D. We also compare four different word embeddings [7, 59, 9, 17], which all utilize wiki data for pretraining. The hashtag features are followed by a 2 layer MLP [1024, 2048], with a ReLU activation using a dropout of 0.25, before concatenated with image feature.

Intent vs. content study To obtain $Mask^O(I)$, we use a pretrained mask-RCNN (X101 32x8d FPN 3x) model⁶ [28] trained on COCO dataset [49] to obtain objects’ segmentation masks with a threshold of 0.6. Multiple objects are merged together. $Mask^C(I)$ is defined as the pixel area in an image I that does not belong to $Mask^O(I)$. A ResNet50 [29] model, pretrained on ImageNet [14] and fine-tuned on each variation of the dataset. All images are resized to longest side of 1280 before processing.

To analyze the relations between content disruption levels and intent recognition scores, we fit a line $\alpha\mathbf{X} + \beta = \mathbf{y}_{F1}$, and define the correlation $\rho(\mathbf{X}, \mathbf{y}_{F1}) \in \{\text{positive, neutral, negative}\}$ based on the normalized slope values ($\bar{\alpha} = \alpha / |\mathbf{X}| \times 10$). The value of $\bar{\alpha}$ and $\rho(\mathbf{X}, \mathbf{y}_{F1})$ are used to group intent classes as described in Sec. A.2.

To investigate the relationship between intent and specific thing and stuff classes, we use a pretrained panoptic FPN segmentation model [42] trained on COCO panoptic dataset and obtain masks for both thing and stuff classes in the images (with a threshold of $\tau_p = 0.7$, p with area less than 10% of the whole image are ignored). The CAM heatmaps are averaged over all five trained model results with $\tau_{cam} = 0.4$. All images are resized to longest side of 1280 before processing.

A.2. Identifying intent classes

To quantify and analyze the experimental results, we group 28 classes into subsets based on two different criteria, *i.e.*, content and difficulty. Table 4 shows a summary.

By content Intent categories are grouped into object-dependent (O -classes), context-dependent (C -classes), and *Others* which depends on both foreground and background information.

By difficulty Based on random guessing and standard classification results using full content information, we categorize classes based on how far the CNN model achieves than the random results. Formally, given a random guessing score r and model result s for a class m , the information gain is defined as $D(m) = r \log(s/r)$. $D(m)$ takes both the value of r and the relative gain from s to r into considerations. The larger D is, the easier the class m is for a standard CNN model to learn.

⁶detectron2 model zoo https://github.com/facebookresearch/detectron2/blob/master/MODEL_ZOO.md

	Classes	Frequency	Definition
Content	\mathcal{O} -classes	7 24.9%	$\bar{\alpha}_{\mathcal{O}} > \bar{\alpha}_{\mathcal{C}}, \rho_{\mathcal{C}} \neq \text{positive}$
	\mathcal{C} -classes	2 11.1%	$\bar{\alpha}_{\mathcal{O}} < \bar{\alpha}_{\mathcal{C}}, \rho_{\mathcal{O}} \neq \text{positive}$
	Others	19 64.0%	o.w.
Difficulty	Easy	3 23.8%	$D \leq 5$
	Medium	15 51.6%	$D \in (5, 15]$
	Hard	10 24.6%	$D > 15$

Table 4. Intent classes categorization. We propose to group 28 classes based on two criteria and report the definition, frequency (in the forms of [number of classes | training image percentage]). See text for definition of D .

Method		Macro F1	
WordBreak	Embeddings	All	Hard
	fastText [7]	19.92 \pm 0.86	6.47 \pm 0.93
✓	BERT [17]	6.58 \pm 0.13	0.0 \pm 0.0
✓	fastText [7]	20.04 \pm 0.53	6.63 \pm 1.45
✓	GloVe [59]	21.37 \pm 0.19	6.64 \pm 0.83
✓	static BERT [9]	18.97 \pm 0.23	7.47 \pm 0.86

Table 5. Model performance with HT feature only on val set. Static BERT with our proposed WordBreak method gives best result.

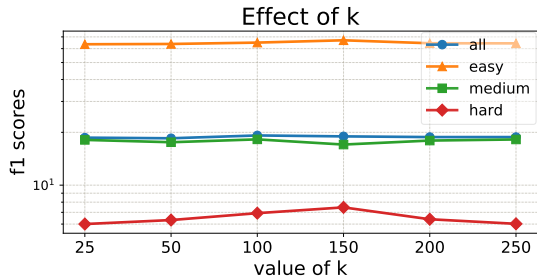


Figure 9. Effect of k for HT features on val set. In general, F1 score peaks at $k = 150$, $k \in \{25, 50, 100, 150, 200, 250\}$. Y-axis is in log scale.

B. Additional Hashtags Results

Separating hashtags benefits hard classes In Table. 5, we report performance using HT only and compare different hashtag representation methods. Hashtags, despite not constituting a natural language, are compact by definition. We observe that separating hashtags into phrases outperforms subword-level embedding for the whole hashtag. FastText embedding [7] utilizes sub-word information and usually works well with rare words. Yet separating hashtags is able to achieve a 15.5% gain on *hard* classes, and 7% on overall macro F1 score. We use static-BERT [9] for all the other experiments since improving hard classes is the reason why we propose using multiple modalities. Note that BERT [17] yields results comparable to random guessing. A possible reason is that the average token length for a hashtag is 4.7 (std = 3.5), which suggests a low level of contextual information within any given hashtag.

Hashtags from k nearest neighbours How does the noise in collected Instagram hashtags impact classification results? We collect hashtags by fetching pixel-level similar Instagram posts using KNN. Thus the collected hashtags are less and less relevant to the image, as k increases. As pointed out by [51, 53, 35] and mentioned above, hashtags are prone to noise: one may include irrelevant hashtags for the post (e.g. #likesforlikes, #igers). We study the performance of resulting hashtag features by varying the number of top nearest neighbors for each sample i . Fig. 9 shows that F1 score for “hard” and “easy” classes peak at $k = 150$. “Medium” classes are less sensitive to the value of k and peak at 100. We use $k = 150$ for all the other experiments.

C. Dataset Creation Details

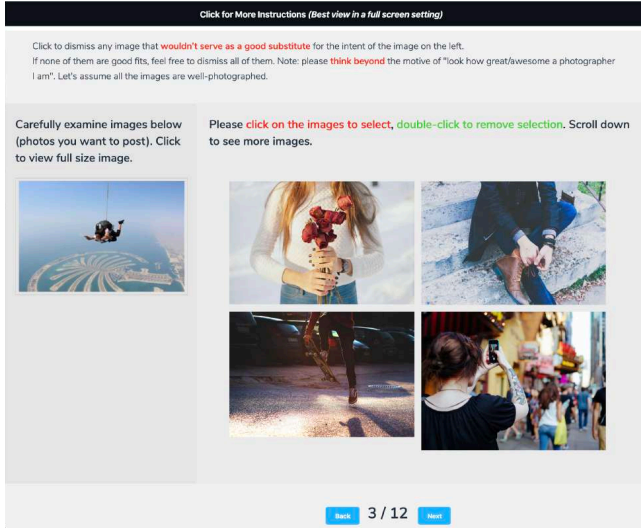
Given the inherent abstract nature of intent understanding, one challenge we are facing is that how to collect reasonable labels in an effective manner. A standard annotation process for image classification task is to ask qualified annotators to select from a list of labels given one image. Annotators become qualified after a series training sessions for the label information [36]. This approach would have been time-consuming and highly dependent on the expertise of our annotators. We instead adopt a *game with a purpose* approach to keep annotators engaged and let them focus on the “swapabilities” of image pairs regarding the intent. We use relative similarity comparison in batch using grid format following [93]. The annotation task is to select all the images in the grid that clearly have a different intent than the reference image on the left. Note that the resulting labels represent the *perceived* intent: the viewer’s opinion of the intent of the image. This section provide more details on the dataset acquisition process.

C.1. Annotation interface

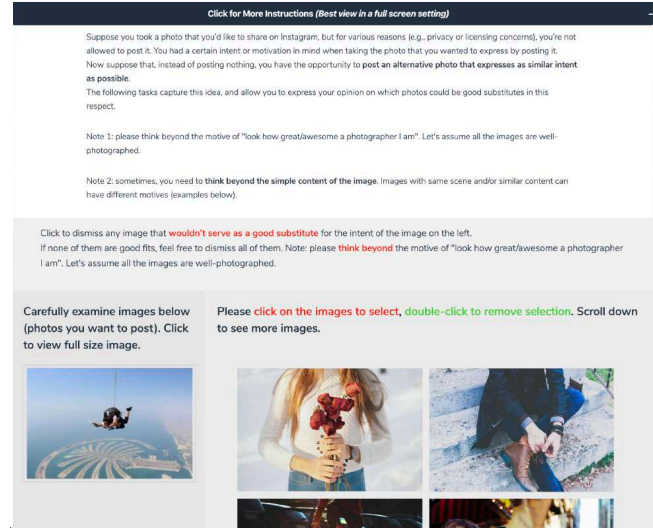
As noted in [83], *games with a purpose* annotation approach, like the ESP Game [86], reCAPTCHA [87] and BubbleBank [15], require some artistry to design tools that keep user engaged. Keeping this principle in mind, we design an interface⁷ that displays a probe image and a 2×2 images grid side by side. Amazon Mechanical Turk workers are asked to select all the images in the grid that clearly have a different motive than the reference image on the left. A welcome splash page is shown at the beginning of each annotation task, to briefly introduce or remind the annotators.

Fig. 10 shows the main annotation interface. There is a collapsible section on top of the interface that display instructions. Images inside the grid are sorted dynamically ac-

⁷Interface is modified based on simpleamt, which use Jinja2 as backend. UI design was adapted from *Snapshot by TEMPLATED, templated.co @templatedco*.



(a)



(b)

Figure 10. Annotation interface. We present a story to the workers to put them into the mindset of the imagined user who want to post the image presented. (a) Main annotation page, with probe image and 2×2 image grid displayed side-by-side. (b) Collapsible instruction on the top of the interface.

ording to the height-to-width ratios, so the interface looks nicer (inspired by [6]). The probe image on the left is always kept shown on the screen throughout scrolling up and down the page.

Since human motives are inherently abstract to understand, we provide a narrative, which is shown below, for the annotators so they could focus on the swapability of images. The narrative presents a story for the workers, which bring them to the scenario of the imagined user who want to post the image presented on the left. We also provided example selections inside the collapsible instructions and the welcome splash page (see Fig. 10(b)).

Annotation narratives: Suppose you took a photo that you'd like to share on Instagram, but for various reasons (e.g., privacy or licensing concerns), you're not allowed to post it. You had a certain intent or motivation in mind when taking the photo that you wanted to express by posting it. Now suppose that, instead of posting nothing, you have the opportunity to *post an alternative photo that expresses as similar intent as possible*. The following tasks capture this idea, and allow you to express your opinion on which photos could be good substitutes in this respect.

We used 4 images per grid, 12 grids per HITs, including 1 catch trials. We only use annotation results that pass the catch trials. In order to get a richer similarity representations, and to examine the quality of the annotators, we also use 3 annotators for the same HIT.

Annotators' feedbacks of our interface and general annotation system include: "I've been enjoying doing these

Keywords	# Instagram post	# Unsplash photos sampled
"people"	39,174,751	8,000
"travel"	479,354,358	4,500
"happy"	564,642,361	5,500
"business"	60,129,975	2,000

Table 6. Keywords and hashtags mapping

hits", "I enjoy these tasks, so I would like to keep doing them", "I truly enjoy these hits and always appreciate the feedback!" "I hope to see more from you guys soon i love doing these!" "Thanks for Your HITs, i really enjoyed working on them and i hope i did good." "I enjoy the HITs and am glad to be able to contribute."

C.2. Images selection

Candidate images Our goal is to fetch photos from Unsplash, that is similar to images uploaded to social medias like Instagram. All the images are visually and aesthetically pleasing content generated by users. Each photo of Unsplash has a list of associated keywords, produced by an online deep-learning based API. We use these keywords to query photos from Unsplash. Criterias for the chosen keywords are: 1) it should be reasonable and possible to appear in Instagram; 2) it should cover a wide range of scenarios in everyday life. With such requirement in mind, we chose four keywords by browsing Unsplash website and using common sense: "people", "travel / vacation", "happy",

and “business”, which were selected according to the popular hashtags on Instagram. Table 6 summarizes the keywords and related number of public Instagram posts as of 2020/3/20. A total of 20,000 images were fetched using these four keywords. During annotation process, our annotators found that around 5K images do not have any intent labels, so we discard those in the analysis and experiments.

Probe images We carefully chose probe images that cover a reasonably large range of scenes and objects [20], including both cluttered and relative uniform scenes, and diverse range of objects, colors, textures and shapes. In order to reduce possible ambiguity during annotations, the probe image also uniquely represents one human motive only. The probe image are manually inspected by all the authors.

C.3. Annotators management

To ensure quality, we restrict access to MTurks who pass our qualification task. And we constantly check the performance and send feedback to MTurks. After first 100 annotation tasks (HIT) we launched at MTurk, we limit the annotation task to the top annotators

Each annotator needs to take a qualification test in order to get access to our annotation task. The purposes of qualification test are two folds: firstly, to help us to select qualified workers who understand that we are annotating motives; secondly, to help workers get familiar with our designed narratives in the annotation. A total of four questions are presented to the potential annotators. Aside from the requirement of having an Instagram account, we provide three questions that serve as an introductory training and qualification task. Three image triplets (a probe image, and two substitute options) were carefully curated, and each triplet was presented as three images side by side. We specifically selected images that either has similar content but different motives with the probe image, or similar motive but different motives.

Periodically, we check the annotation progress and send messages to workers to inform how many catch trials they failed. We received positive responses from annotators about such feedback system. One annotator commented that “It’s always nice when workers receive feedback from requesters on MTurk about the quality of the work being done, and it was reassuring to receive emails (even if they were more-or-less automated) from your team to let me know I was doing well.”

C.4. Annotation methods comparison

Instead of selecting from a list motive labels given each image, we adopted image comparison approach, using “unsatisfactory substitutes” and mental imagery. The average annotation time for one annotation task is 20.60 (± 9.65) minutes. Each annotation task contains 48 images. There-

fore, the annotators spend 25.75 second per image on average.

To compare two annotation approaches, two authors of this study annotated 57 random sampled images from our dataset using standard image tagging annotation method (82.2 second) per image. Our image comparison method using “unsatisfactory substitutes” requires less annotation time per image.

C.5. Human-in-the-loop

We adopted a hybrid human-in-the-loop strategy to incrementally learn a motive classifier in the annotation process. Starting from a set of randomly selected images, the dataset is enlarged by an iterative process that utilizes a trained classifier to recommend relevant images to annotators. At each iteration and for each motive label, we train a deep learning classifier using 90% of the labeled data. 10% of the held-out data is always added to the test set. The trained classifier is applied to the rest of unlabeled data, and images with a score larger than 0.35 are sent back to annotators for verification. We applied this method until there is no positive image left in the unlabeled set for each label. See Fig. 11 for examples of our dataset images.

C.6. Inter annotator agreement

As explained in the main text, each image was inspected by three annotators. We use Fleiss’ kappa score [21] to measure inter annotator agreements per annotation task. The average score is 59.84%, indicating “moderate” agreement [27]. The inter-annotator agreement score demonstrates the complexity of the annotation task, and the inherently abstract nature of human intent understanding.

C.7. Test set annotation

We ask one author, as chief executive to annotate the validation and test set. The annotation process took three weeks. We found that there are more images per label in the resulting annotation, comparing to the MTurk result. This further demonstrate the MTurks are able to identify correct motive labels using our *game with a purpose* approach. Yet in general MTurks tend to miss some of the labels.

D. Dataset Analysis

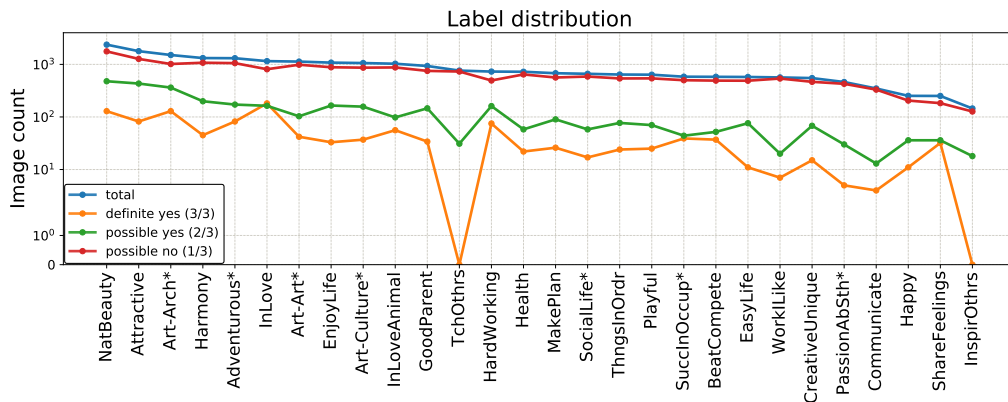
In this section, we analyze the properties of the dataset in more detail, and examine the inter-annotator agreement.

Dataset statistics Fig. 12 shows the label distribution of whole training data, over 28 classes, 9 super categories, 3 content classes, and 3 difficulty classes. It shows there is class imbalance in our dataset, which is the property of datasets in the real world [82].

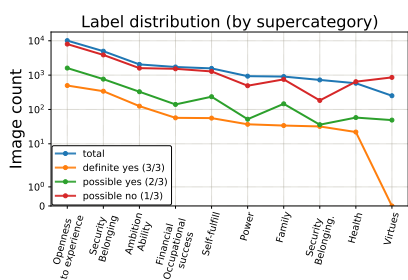
Hashtags We also fetched hashtags from Instagram, with the hope of further capturing the semantics of human in-

Label	Definite yes (3/3)	Possible yes (2/3)	Possible no (1/3)
beat and compete			
enjoy life			
manageable, making plans			
natural beauty			
things in order			

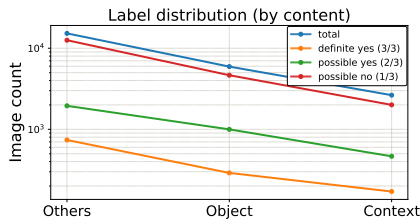
Figure 11. Sample motive labels, and images that are respectively marked as definite yes (3 out of 3 annotators agree), possible yes (2 out of 3 agree), and possible no (1 out of 3 agree). Images that belong to “definite yes” and “possible yes” can have completely different objects, scenes. This further illustrate the high intra-class variance nature of intent classification.



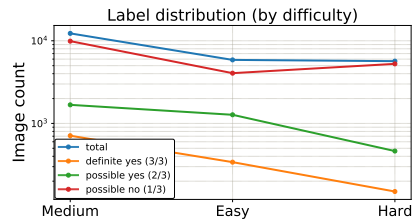
(a) Per-class distribution.



(b) Distribution by supercategories.



(c) By content groups.



(d) By difficulty groups.

Figure 12. Training data distribution. Class names ends with “*” are abbreviated.

Intent classes	Top hashtags
Attractive	#portrait #fashionblogger #womenempowerment #makeup
SocialLife Friendship	#family #sun #sea #beach
NaturalBeauty	#mountains #landscape #sunrise #sunset #naturelovers
Playful	#travel #guitar #lifestyle #puppy #livemusic
Happy	#smile #newbornphotography #mood #headshot #vibes
WorkILike	#entrepreneur #smallbusiness #motivation #marketing

Table 7. Common hashtags for six intent classes.

Dataset	Intent type	# Intent classes	# Images
MDID [45]	Textual	8	1299
Intentionomy	Visual	28	14455

Table 8. Comparison with prior work.

tents. In total, we fetched 1,700,915 unique hashtags. Each Unsplash photo has an average of 457.6 (± 317.512) hashtags. As noted in [85], hashtags serve as a medium of self-expression that not limited to objective descriptions of image content. Table 7 lists a collections of top hashtags for selected intent classes

Lexical statistics We fetch the accompanying text description with the images found on the website. These descriptions are generated by a deep-learning based API and verified by human. We report the lexical (word-level) statistic of the dataset. Specifically, the top words occurred in the descriptions of validation images are presented. Table 9 shows top 10 frequent non-stopping words per class, shedding light on the properties of the images. Although the descriptions can be heavily biased, Table 9 illustrates that, as they should, the occurrences of image objects and properties are relatively balanced across all the classes, indicating that most of the frequent words are not necessarily directly predictive of the intent label. However, we do admit that there are exceptions. Certain words can be correlated to certain human intents. For example, “face” occurs frequently in the class “CreativeUnique”. “Smiling” is one of the top 10 frequent words in “Playful”, “Happy”, and “InspirOthrs”.

Note that there are 985 “man” and 1714 “woman” in total in the test set, indicating the existence of gender bias in our dataset, which is a common issue in nowadays machine learning systems [8, 23, 16, 22, 63]. “Woman” occurs 74% more than “man”. We observe that female gender word tend to associates with classes like “attractive”, “happy”, “enjoy life”. Male gender, on the other hand, associates with “exciting life”, “health”, “beat and compete”. As pointed out

in [23], such gender-specific associations, even with subjectively positive words such as the intent labels presented, are *benevolent sexism*. We would like to raise the awareness of such phenomenon. Any machine learning down-streaming tasks should always apply fairness into consideration during algorithm development.

E. Intent Taxonomy

Table 10 lists the detailed taxonomy and explanation for each intent class. Note that there are similarities between emotions and motives. For example, the category “happy/joy” appears in both emotion recognition [58, 44, 34, 92] and intent recognition [38]. Indeed, the common Latin root word of “emotion” and “motivation” is “movere” (to move) [69]. Young [98] argues both emotion and motivation influence human behavior, and that emotion arises from the interplay (e.g. conflict, frustration, satisfaction) of motives. Emotions can also be viewed as a reward or punishment for a specific motivated behavior [76].

F. Additional Related Work

Comparison with prior work Table 8 summarizes the differences between Intentionomy and prior work that focuses on social media intent. Other discussions can be found in the main text (Sec. 2).

Subjective attributes Recently, there are some progress in building datasets describing the subjective perspective of images [1]. For example, [101] studies visual common-sense reasoning, requiring computational model to answer challenging questions about an image and provide a rationale justification. Some prior work studies visual rhetoric from different perspectives: 1) protest activity from social media images [94]; 2) memorability [41]; 3) personality [55]; 4) evoked emotions and sentiment [58, 44, 34, 2, 40, 92]. Our work focuses on the *perceived* intent recognition, which is another psychological feature⁸.

G. Other Concerns and Thoughts

Comparison with human performance A proper human experiment involves careful experimental design accounting for variables including demographic information, life experience, and cultural background. At present, such an effort is out of the scope of our study. We believe, however, that our project provides a starting point for future studies with human subjects.

⁸The definition of “motives” according to Merriam Webster [18] is that something (such as a need or desire) that causes a person to act.

Class	Top words
Attractive	woman (257), wearing (100), white (56), standing (55), black (52), photography (41), man (37), near (36), top (35), holding (33)
BeatCompete	man (48), person (29), woman (25), daytime (24), black (21), white (21), holding (16), photography (14), standing (14), riding (13)
Communicate	woman (29), black (16), sitting (13), photography (13), person (11), brown (10), holding (10), man (10), white (10), two (9)
CreativeUnique	woman (36), man (20), holding (16), person (14), photography (14), white (14), face (12), black (12), green (11), blue (11)
CuriousAdventurousExcitingLife	man (62), person (57), woman (55), daytime (51), standing (42), white (39), photography (36), near (34), wearing (29), gray (27)
EasyLife	woman (51), white (20), person (18), daytime (18), sitting (18), photography (17), man (14), standing (13), near (13), wearing (12)
EnjoyLife	woman (86), daytime (45), standing (35), near (35), person (34), man (33), holding (30), sitting (30), white (29), water (28)
FineDesignLearnArt-Arch	photography (54), white (47), building (46), woman (43), daytime (41), near (39), photo (36), concrete (32), people (30), brown (27)
FineDesignLearnArt-Art	woman (61), person (40), daytime (40), man (39), white (38), black (31), holding (29), photography (29), standing (28), near (28)
FineDesignLearnArt-Culture	woman (89), man (47), wearing (37), standing (37), daytime (29), holding (27), white (26), black (25), near (24), photography (23)
GoodParentEmoCloseChild	woman (71), man (42), wearing (37), daytime (33), white (32), holding (28), black (27), near (26), standing (26), photography (23)
Happy	woman (94), wearing (48), standing (28), man (27), smiling (23), black (20), shirt (17), white (16), brown (14), photography (13)
HardWorking	macbook (14), person (13), book (10), holding (9), woman (8), white (8), man (7), brown (7), using (6), sitting (6)
Harmony	woman (94), standing (62), man (52), person (50), near (47), daytime (43), white (33), sitting (33), photo (30), photography (30)
Health	man (44), woman (32), person (20), daytime (18), people (18), white (17), photography (16), body (15), near (15), water (15)
InLove	woman (97), man (68), wearing (36), standing (35), near (34), daytime (33), white (29), person (28), photography (28), sitting (26)
InLoveAnimal	woman (53), white (45), man (36), standing (33), person (32), photography (28), near (26), black (26), daytime (25), brown (24)
InspirOthrs	man (11), person (9), standing (8), holding (7), woman (7), black (5), stage (5), playing (4), wearing (3), smiling (3)
ManagableMakePlan	white (37), black (28), person (22), near (20), woman (20), brown (15), photo (15), holding (13), book (13), macbook (12)
NatBeauty	woman (107), standing (98), man (96), daytime (76), mountain (72), near (65), person (64), photography (64), water (57), white (56)
PassionAbSmthing	woman (27), wearing (17), man (16), white (15), standing (13), black (13), daytime (12), near (12), photography (12), brown (11)
Playful	woman (69), wearing (29), man (26), black (23), holding (19), white (16), smiling (14), standing (14), near (14), daytime (13)
ShareFeelings	people (16), man (10), person (8), group (8), holding (7), woman (7), black (7), smartphone (5), focus (4), photography (4)
SocialLifeFriendship	woman (46), photography (22), wearing (22), man (20), black (20), person (16), standing (16), people (15), daytime (15), sitting (14)
SuccInOccupHavGdJob	woman (43), man (31), black (24), white (23), wearing (22), standing (18), photo (13), person (13), holding (12), near (11)
TchOthrs	woman (50), man (37), person (23), white (22), black (21), photography (21), near (21), wearing (21), standing (19), daytime (18)
ThngsInOrdr	white (25), woman (23), brown (19), black (19), standing (18), man (18), top (15), person (12), near (12), daytime (12)
WorkILike	woman (49), man (34), person (25), black (21), wearing (20), sitting (18), near (18), holding (18), daytime (18), white (15)

Table 9. Lexical statistics of the image descriptions in the validation set. Top 10 most frequent non-stopping words per class. The numbers next to each word is the count within that specific class

Class	Descriptions
Attractive	Being good looking, attractive.
BeatCompete	Beat people in a competition.
Communicate	To communicate or express myself.
CreativeUnique	Being creative (e.g., artistically, scientifically, intellectually). Being unique or different.
CuriousAdventurousExcitingLife	Exploration - Being curious and adventurous. Having an exciting, stimulating life.
EasyLife	Having an easy and comfortable life.
EnjoyLife	Enjoying life
FineDesignLearnArt-Arch	Appreciating fine design (man-made wonders like architectures)
FineDesignLearnArt-Art	Appreciating fine design (artwork)
FineDesignLearnArt-Culture	Appreciating other cultures
GoodParentEmoCloseChild	Being a good parent (teaching, transmitting values). Being emotionally close to my children.
Happy	Being happy and content. Feeling satisfied with one's life. Feeling good about myself.
HardWorking	Being ambitious, hard-working.
Harmony	Achieving harmony and oneness (with self and the universe).
Health	Being physically active, fit, healthy, e.g. maintaining a healthy weight, eating nutritious foods. To be physically able to do my daily/routine activities. Having athletic ability.
InLove	Being in love.
InLoveAnimal	Being in love with animal
InspirOthers	Inspiring others, Influencing, persuading others.
ManagableMakePlan	To keep things manageable. To make plans
NatBeauty	Experiencing natural beauty.
PassionAbSmthing	Being really passionate about something.
Playful	Being playful, carefree, lighthearted.
ShareFeelings	Sharing my feelings with others.
SocialLifeFriendship	Being part of a social group. Having people to do things with. Having close friends, others to rely on. Making friends, drawing others near.
SuccInOccupHavGdJob	Being successful in my occupation. Having a good job.
TeachOthers	Teaching others.
ThngsInOrdr	Keeping things in order (my desk, office, house, etc.).
WorkILike	Having work I really like.

Table 10. The taxonomy for our Intentionomy dataset.

References

- [1] Xavier Alameda-Pineda, Miriam Redi, Nicu Sebe, and Shih-Fu Chang. *Understanding Subjective Attributes of Data*, 2019 (accessed October 07, 2020). 14
- [2] Xavier Alameda-Pineda, Elisa Ricci, Yan Yan, and Nicu Sebe. Recognizing emotions from abstract paintings using non-linear matrix completion. In *CVPR*, pages 5240–5248, 2016. 14
- [3] Tamar Ashuri, Shira Dvir-Gvishman, and Ruth Halperin. Watching me watching you: How observational learning affects self-disclosure on social network sites? *Journal of Computer-Mediated Communication*, 23(1):34–68, 2018. 2
- [4] Saeideh Bakhshi, David A Shamma, Lyndon Kennedy, and Eric Gilbert. Why we filter our photos and how it impacts engagement. In *Ninth International AAAI Conference on Web and Social Media*, 2015. 2
- [5] Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss, editors, *ECCV*, 2018. 2
- [6] Sean Bell, Paul Upchurch, Noah Snavely, and Kavita Bala. OpenSurfaces: A richly annotated catalog of surface appearance. *ACM Trans. on Graphics (SIGGRAPH)*, 32(4), 2013. 11
- [7] Piotr Bojanowski, Edouard Grave, Armand Joulin, and

- Tomas Mikolov. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146, 2017. 9, 10
- [8] Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in neural information processing systems*, pages 4349–4357, 2016. 14
- [9] Rishi Bommasani, Kelly Davis, and Claire Cardie. Interpreting Pretrained Contextualized Representations via Reductions to Static Embeddings. In *ACL*, pages 4758–4781, Online, July 2020. Association for Computational Linguistics. 9, 10
- [10] Wieland Brendel and Matthias Bethge. Approximating cnns with bag-of-local-features models works surprisingly well on imagenet. *ICLR*, 2019. 2
- [11] Xinlei Chen, Li-Jia Li, Li Fei-Fei, and Abhinav Gupta. Iterative visual reasoning beyond convolutions. In *CVPR*, June 2018. 2
- [12] Myung Jin Choi, Antonio Torralba, and Alan S Willsky. Context models and out-of-context objects. *Pattern Recognition Letters*, 33(7):853–862, 2012. 2
- [13] Y. Cui, M. Jia, T.Y. Lin, Y. Song, and S. Belongie. Class-balanced loss based on effective number of samples. In *CVPR*, 2019. 9
- [14] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR*, 2009. 3, 6, 9
- [15] Jia Deng, Jonathan Krause, and Li Fei-Fei. Fine-grained crowdsourcing for fine-grained recognition. In *CVPR*, pages 580–587, 2013. 3, 10
- [16] Sunipa Dev and Jeff Phillips. Attenuating bias in word vectors. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 879–887, 2019. 14
- [17] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NACCL*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. 9, 10
- [18] Merriam-Webster Dictionary. Merriam-webster. *On-line at <http://www.mw.com/home.htm>*, 2002. 14
- [19] Arie Dijkstra. The psychology of tailoring-ingredients in computer-tailored persuasion. *Social and personality psychology compass*, 2(2):765–784, 2008. 2
- [20] Marin Ferecatu and Donald Geman. A statistical framework for image category search from a mental picture. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(6):1087–1101, 2008. 12
- [21] Joseph L Fleiss, Bruce Levin, and Myunghee Cho Paik. *Statistical methods for rates and proportions*. john wiley & sons, 2013. 12
- [22] Nikhil Garg, Londa Schiebinger, Dan Jurafsky, and James Zou. Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences*, 115(16):E3635–E3644, 2018. 14
- [23] Peter Glick and Susan T Fiske. The ambivalent sexism inventory: Differentiating hostile and benevolent sexism. In *Social Cognition*, pages 116–160. Routledge, 2018. 14
- [24] Albert Gordo, Jon Almazán, Jérôme Revaud, and Diane Larlus. Deep image retrieval: Learning global representations for image search. *CoRR*, abs/1604.01325, 2016. 6
- [25] Ralph Gross, Iain Matthews, Jeffrey Cohn, Takeo Kanade, and Simon Baker. Multi-pie. *Image and Vision Computing*, 28(5):807–813, 2010. 1
- [26] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *CVPR*, 2019. 1
- [27] Lisa Hartling, Michele Hamm, Andrea Milne, Ben Vandermeer, P Lina Santaguida, Mohammed Ansari, Alexander Tsertsvadze, Susanne Hempel, Paul Shekelle, and Donna M Dryden. *Validity and inter-rater reliability testing of quality assessment instruments*. Agency for Healthcare Research and Quality (US), 2012. 12
- [28] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. 9
- [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 9
- [30] Derek Hoiem, Alexei A Efros, and Martial Hebert. Geometric context from a single image. In *ICCV*, 2005. 2
- [31] Hexiang Hu, Guang-Tong Zhou, Zhiwei Deng, Zicheng Liao, and Greg Mori. Learning structured inference neural networks with label relations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. 2
- [32] Xinyue Huang and Adriana Kovashka. Inferring visual persuasion via body language, setting, and deep features. In *CVPRW*, pages 73–79, 2016. 2
- [33] Minyoung Huh, Andrew Liu, Andrew Owens, and Alexei A Efros. Fighting fake news: Image splice detection via learned self-consistency. In *ECCV*, 2018. 1
- [34] Zaeem Hussain, Mingda Zhang, Xiaozhong Zhang, Keren Ye, Christopher Thomas, Zuha Agha, Nathan Ong, and Adriana Kovashka. Automatic understanding of image and video advertisements. In *CVPR*, pages 1705–1715, 2017. 2, 14
- [35] Hamid Izadinia, Bryan C. Russell, Ali Farhadi, Matthew D. Hoffman, and Aaron Hertzmann. Deep classifiers from image tags in the wild. In *Proceedings of the 2015 Workshop on Community-Organized Multimodal Mining: Opportunities for Novel Solutions*, MMCommons ’15, page 13–18, New York, NY, USA, 2015. Association for Computing Machinery. 10
- [36] Menglin Jia, Mengyun Shi, Mikhail Sirotenko, Yin Cui, Claire Cardie, Bharath Hariharan, Hartwig Adam, and Serge Belongie. Fashionpedia: Ontology, segmentation, and an attribute localization dataset. In *ECCV*, 2020. 1, 10
- [37] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. *arXiv preprint arXiv:1702.08734*, 2017. 9
- [38] Jungseock Joo, Weixin Li, Francis F Steen, and Song-Chun Zhu. Visual persuasion: Inferring communicative intents of images. In *CVPR*, pages 216–223, 2014. 2, 14

- [39] J. Joo, F. F. Steen, and S. Zhu. Automated facial trait judgment and election outcome prediction: Social dimensions of face. In *ICCV*, pages 3712–3720, 2015. 2
- [40] Brendan Jou, Tao Chen, Nikolaos Pappas, Miriam Redi, Mercan Topkara, and Shih-Fu Chang. Visual affect around the world: A large-scale multilingual visual sentiment ontology. In *Proceedings of the 23rd ACM international conference on Multimedia*, pages 159–168, 2015. 14
- [41] Aditya Khosla, Akhil S Raju, Antonio Torralba, and Aude Oliva. Understanding and predicting image memorability at a large scale. In *ICCV*, pages 2390–2398, 2015. 14
- [42] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In *CVPR*, pages 6399–6408, 2019. 9
- [43] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *CVPR*, pages 9404–9413, 2019. 5
- [44] Ronak Kosti, Jose M Alvarez, Adria Recasens, and Agata Lapedriza. Emotion recognition in context. In *CVPR*, pages 1667–1675, 2017. 14
- [45] Julia Kruk, Jonah Lubin, Karan Sikka, Xiao Lin, Dan Jurafsky, and Ajay Divakaran. Integrating text and image: Determining multimodal document intent in instagram posts. In *EMNLP*, pages 4614–4624, 2019. 2, 14
- [46] Chih-Hui Lai. Motivations, usage, and perceived social networks within and beyond social media. *Journal of Computer-Mediated Communication*, 24(3):126–145, 2019. 2
- [47] So-Hyun Lee and Hee-Woong Kim. Why people post benevolent and malicious comments online. *Communications of the ACM*, 58(11):74–79, 2015. 2
- [48] T.Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *CVPR*, 2018. 9
- [49] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 1, 9
- [50] Yong Liu, Ruiping Wang, Shiguang Shan, and Xilin Chen. Structure inference net: Object detection using scene-level context and instance-level relationships. In *CVPR*, pages 6985–6994, 2018. 2
- [51] D.K. Mahajan, R.B. Girshick, V. Ramanathan, K. He, M. Paluri, Y. Li, A. Bharambe, and L. van der Maaten. Exploring the limits of weakly supervised pretraining. In *ECCV*, 2018. 7, 10
- [52] Ishan Misra, Abhinav Gupta, and Martial Hebert. From Red Wine to Red Tomato: Composition with Context. In *CVPR*, 2017. 2
- [53] Ishan Misra, C. Lawrence Zitnick, Margaret Mitchell, and Ross Girshick. Seeing through the Human Reporting Bias: Visual Classifiers from Noisy Human-Centric Labels. In *CVPR*, 2016. 10
- [54] R. Mottaghi, X. Chen, X. Liu, N. Cho, S. Lee, S. Fidler, R. Urtasun, and A. Yuille. The role of context for object detection and semantic segmentation in the wild. In *CVPR*, 2014. 2
- [55] Nils Murrugarra-Llerena and Adriana Kovashka. Cross-modality personalization for retrieval. In *CVPR*, pages 6429–6438, 2019. 14
- [56] Dennis Park, Deva Ramanan, and Charless Fowlkes. Multiresolution models for object detection. In *ECCV*, pages 241–254. Springer, 2010. 2
- [57] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017. 9
- [58] Kuan-Chuan Peng, Tsuhan Chen, Amir Sadovnik, and Andrew C Gallagher. A mixed bag of emotions: Model, predict, and transfer emotion distributions. In *CVPR*, pages 860–868, 2015. 14
- [59] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014. 9, 10
- [60] Hamed Pirsiavash, Carl Vondrick, and Antonio Torralba. Inferring the why in images. Technical report, MASSACHUSETTS INST OF TECH CAMBRIDGE, 2014. 2
- [61] Kashyap Papat, Subhabrata Mukherjee, Andrew Yates, and Gerhard Weikum. DeClarE: Debunking fake news and false claims using evidence-aware deep learning. In *EMNLP*, pages 22–32, Brussels, Belgium, Oct.-Nov. 2018. Association for Computational Linguistics. 1
- [62] Javier Portilla and Eero P Simoncelli. A parametric texture model based on joint statistics of complex wavelet coefficients. *International journal of computer vision*, 40(1):49–70, 2000. 4
- [63] Marcelo OR Prates, Pedro H Avelar, and Luís C Lamb. Assessing gender bias in machine translation: a case study with google translate. *Neural Computing and Applications*, pages 1–19, 2019. 14
- [64] Ariadna Quattoni and Antonio Torralba. Recognizing indoor scenes. In *CVPR*, pages 413–420. IEEE, 2009. 1
- [65] Barbara K Rimer and Matthew W Kreuter. Advancing tailored health communication: A persuasion and message effects perspective. *Journal of communication*, 56:S184–S201, 2006. 2
- [66] David Rolnick, Andreas Veit, Serge Belongie, and Nir Shavit. Deep learning is robust to massive label noise. *arXiv preprint arXiv:1705.10694*, 2017. 3
- [67] Florian Schroff, Tali Treibitz, David Kriegman, and Serge Belongie. Pose, illumination and expression invariant pairwise face-similarity measure via doppelgänger list comparison. In *ICCV*, Barcelona, 2011. 1
- [68] Behjat Siddiquie, Dave Chisholm, and Ajay Divakaran. Exploiting multimodal affect and semantics to identify politically persuasive web videos. In *Proceedings of the 2015 ACM on International Conference on Multimodal Interaction*, pages 203–210, 2015. 2
- [69] Sarah M Sincero. *Motivation and Emotion*, 2012 (accessed October 07, 2020). 14
- [70] Krishna Kumar Singh, Dhruv Mahajan, Kristen Grauman, Yong Jae Lee, Matt Feiszli, and Deepti Ghadiyaram. Don’t judge an object by its context: Learning to overcome contextual bias. In *CVPR*, 2020. 6, 9

- [71] Flávio Souza, Diego de Las Casas, Vinícius Flores, Sunbum Youn, Meeyoung Cha, Daniele Quercia, and Virgílio Almeida. Dawn of the selfie era: The whos, wheres, and hows of selfies on instagram. In *Proceedings of the 2015 ACM on conference on online social networks*, pages 221–231, 2015. [2](#)
- [72] Jennifer R Talevich, Stephen J Read, David A Walsh, Ravi Iyer, and Gurveen Chopra. Toward a comprehensive taxonomy of human motives. *PloS one*, 12(2):e0172279, 2017. [1](#), [2](#), [3](#)
- [73] Jennifer R Talevich, Stephen J Read, David A Walsh, Ravi Iyer, and Gurveen Chopra. Toward a comprehensive taxonomy of human motives. *PloS one*, 12(2):e0172279, 2017. [2](#)
- [74] Makarand Tapaswi, Yukun Zhu, Rainer Stiefelhagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. MovieQA: Understanding Stories in Movies through Question-Answering. In *CVPR*, 2016. [1](#)
- [75] Damien Teney, Lingqiao Liu, and Anton van Den Hengel. Graph-structured representations for visual question answering. In *CVPR*, pages 1–9, 2017. [2](#)
- [76] Julian F Thayer and Richard D Lane. A model of neurovisceral integration in emotion regulation and dysregulation. *Journal of affective disorders*, 61(3):201–216, 2000. [14](#)
- [77] Christopher Thomas and Adriana Kovashka. Predicting the politics of an image using webly supervised data. In *Advances in Neural Information Processing Systems*, pages 3625–3637, 2019. [2](#)
- [78] Nigel J.T. Thomas. Mental imagery. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, summer 2019 edition, 2019. [2](#)
- [79] Giorgos Toliás, Ronan Sifre, and Hervé Jégou. Particular object retrieval with integral max-pooling of cnn activations. In *Proceedings of the International Conference on Learning Representations*, 2016. [6](#)
- [80] Antonio Torralba. Contextual priming for object detection. *International journal of computer vision*, 53(2):169–191, 2003. [2](#)
- [81] Antonio Torralba, Kevin P Murphy, and William T Freeman. Using the forest to see the trees: exploiting context for visual object detection and localization. *Communications of the ACM*, 53(3):107–114, 2010. [2](#)
- [82] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *CVPR*, Salt Lake City, UT, 2018. [12](#)
- [83] Grant Van Horn, Steve Branson, Ryan Farrell, Scott Haber, Jessie Barry, Panos Ipeirotis, Pietro Perona, and Serge Belongie. Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection. In *CVPR*, Boston, MA, 2015. [3](#), [10](#)
- [84] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *CVPR*, 2018. [1](#)
- [85] Andreas Veit, Maximilian Nickel, Serge Belongie, and Laurens van der Maaten. Separating self-expression and visual content in hashtag supervision. In *CVPR*, Salt Lake City, UT, 2018. [14](#)
- [86] Luis Von Ahn. Games with a purpose. *Computer*, 39(6):92–94, 2006. [3](#), [10](#)
- [87] Luis Von Ahn, Benjamin Maurer, Colin McMillen, David Abraham, and Manuel Blum. recaptcha: Human-based character recognition via web security measures. *Science*, 321(5895):1465–1468, 2008. [3](#), [10](#)
- [88] Carl Vondrick, Deniz Oktay, Hamed Pirsiavash, and Antonio Torralba. Predicting motivations of actions by leveraging text. In *CVPR*, pages 2997–3005, 2016. [2](#)
- [89] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011. [3](#)
- [90] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011. [1](#)
- [91] Xuewei Wang, Weiyan Shi, Richard Kim, Yoojung Oh, Sijia Yang, Jingwen Zhang, and Zhou Yu. Persuasion for good: Towards a personalized persuasive dialogue system for social good. In *ACL*, pages 5635–5649, 2019. [2](#)
- [92] Zijun Wei, Jianming Zhang, Zhe Lin, Joon-Young Lee, Niranjan Balasubramanian, Minh Hoai, and Dimitris Samaras. Learning visual emotion representations from web data. In *CVPR*, pages 13106–13115, 2020. [14](#)
- [93] Michael Wilber, Sam Kwak, and Serge Belongie. Cost-effective hits for relative similarity comparisons. In *Human Computation and Crowdsourcing (HCOMP)*, Pittsburgh, 2014. [10](#)
- [94] Donghyeon Won, Zachary C Steinert-Threlkeld, and Jungseock Joo. Protest activity detection and perceived violence estimation from social media images. In *Proceedings of the 25th ACM international conference on Multimedia*, pages 786–794, 2017. [14](#)
- [95] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *CVPR*, pages 3485–3492. IEEE, 2010. [1](#)
- [96] Jian Yao, Sanja Fidler, and Raquel Urtasun. Describing the scene as a whole: Joint object detection, scene classification and semantic segmentation. In *CVPR*, pages 702–709. IEEE, 2012. [2](#)
- [97] Keren Ye, Narges Honarvar Nazari, James Hahn, Zaeem Hussain, Mingda Zhang, and Adriana Kovashka. Interpreting the rhetoric of visual advertisements. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. [2](#)
- [98] Paul Thomas Young. *Emotion in man and animal; its nature and relation to attitude and motive*. Wiley, 1943. [14](#)
- [99] Dian Yu and Zhou Yu. Midas: A dialog act annotation scheme for open domain human machine spoken conversations. *arXiv preprint arXiv:1908.10023*, 2019. [2](#)

- [100] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *CVPR*, June 2019. [1](#)
- [101] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *CVPR*, pages 6720–6731, 2019. [14](#)
- [102] Rowan Zellers, Mark Yatskar, Sam Thomson, and Yejin Choi. Neural motifs: Scene graph parsing with global context. In *CVPR*, June 2018. [2](#)
- [103] Mingda Zhang, Rebecca Hwa, and Adriana Kovashka. Equal but not the same: Understanding the implicit relationship between persuasive images and text. *arXiv preprint arXiv:1807.08205*, 2018. [2](#)
- [104] Mengmi Zhang, Claire Tseng, and Gabriel Kreiman. Putting visual object recognition in context. In *CVPR*, pages 12985–12994, 2020. [2](#), [3](#)
- [105] B. Zhou, A. Khosla, Lapedriza. A., A. Oliva, and A. Torralba. Learning Deep Features for Discriminative Localization. *CVPR*, 2016. [5](#), [6](#)
- [106] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017. [1](#)