

ADVISE: Symbolism and External Knowledge for Decoding Advertisements (Supplementary Material)

Keren Ye^[0000–0002–7349–7762] and Adriana Kovashka^[0000–0003–1901–9660]

University of Pittsburgh, Pittsburgh PA 15260, USA
{yekeren,kovashka}@cs.pitt.edu

In this document, we include more information and statistics about our ADVISE model and some implementation decisions. We also provide additional quantitative and qualitative experimental results.

We first describe in more depth the implementation and evaluation setup. We provide details about the implementation and training of our ADVISE model in Sec. 1. In Sec. 2, we describe different ways of adapting Hussain et al. [22]’s method to the ranking task. In Sec. 3, we explain the reason why we choose to rank 50 statements for our main task.

Next, we provide additional quantitative results which demonstrate the contribution of different algorithmic choices that we made. In Sec. 4, we provide details justifying our choice of the value of k' for hard negative mining. In Sec. 5, we demonstrate different strategies and justify our choice of the attention mechanism used. In Sec. 6, we quantitatively demonstrate that it is not sufficient to use any type of label in the domain of interest for the method component described in Sec. 3.4 of the main text; in particular, we show that using the symbol labels as constraints gives more improvement than using topic labels. In Sec. 7, we break down the ranking task evaluation into topics.

Finally, to enable a more intuitive understanding of our model, we provide more qualitative results of the ranking task in Sec. 8, including both statement ranking and hard-statement ranking. In Sec. 9, we show qualitatively that the ADVISE model learns not only the image representation but also meaningful *region* representations.

1 Training the ADVISE model

As shown in Figure 2 in our paper, there are two branches in our model, the main branch (top left), and the knowledge embedding branch (top right). We do not train the branches jointly at the very beginning, instead, we at first train the main branch till the model converges, then we add the knowledge base information. It is beneficial to incorporate knowledge additively using this two-step process. In our experiments, training the knowledge branch requires less time compared to training the main branch. Thus, one could efficiently experiment with multiple extra types of knowledge given that the main branch is trained and the entry-points of symbols are properly set. Another advantage of the two-step training is that the knowledge branch serves a role similar to a residual branch, that is it would not hurt the performance of the existing main branch.

We experimented with using Adagrad, Adam, RMSProp and found Adagrad to give the best results. We use a learning rate of 2.0 and apply no decay strategy on the learning rate. Also, we did not use different gradient multipliers for image and text embedding networks. For both the image embedding network (\mathbf{w} in Eq. 4) and the attention prediction network (\mathbf{w}_a in Eq. 5), we use batch normalization layer, a weight decay of 1e-6, and dropout keep probability of 0.7 for the input Inception V4 feature. For the text embedding network of the statement (Sec. 3.2), DenseCap captions, and symbols (Sec. 3.4), we use a weight decay of 1e-8 and a dropout keep probability of 0.7 for the embedding weights. The DenseCap captions do not share weights with the statements, but they are both initialized from GLOVE word embedding [48]. For the “unknown” words of DenseCap captions, Ads statements and the symbol embedding vectors, they are initialized from the uniform distribution ranging from -0.08 to 0.08. According to our experiments, adding a dropout layer (with keep probability of 0.5) after the pointwise multiplication of image and text embedding of $\|x - y\|_2^2$ is really important. It assures that the model will not overfit on the training set. For the triplet training, we mine the most challenging 32 negative examples in the 128-sized training batch, we weigh 0.1 for the symbol loss and object loss as mentioned in Eq. 8. Based on the settings mentioned above, we train the main branch for 100,000 steps and use Recall@3 as the metric to choose the best model on the validation set.

We build the knowledge branch after getting the checkpoint of the main branch. During this second phase, we freeze all of the parameters we harvest in the first step, saying \mathbf{w} for the image embedding \mathbf{z} , \mathbf{w}_a for attention prediction, \mathbf{t} for Ads statement embedding, \mathbf{c} for DenseCap embedding, and \mathbf{s} for symbol embedding. In the meantime, we also freeze the parameters of the pre-trained symbol classifiers (\mathbf{u}_{symb} as mentioned in Sec. 3.5) since the classifiers are part of our prior knowledge. Therefore, the only parameters of our ADVISE model in the second training phase are the 53 scalar values, that is, α_j^{symb} (for each of $j \in \{1, \dots, C = 53\}$), denoting the importance of the 53 classifiers. If the main branch captures all the information during the training, assigning 0 to all α_j^{symb} would not hurt the performance of the model. In case the main branch

does *not* capture all relevant information, the knowledge branch may provide complementary information to help to improve the final performance. Thus the knowledge branch trained is similar to a residual branch. Please note that the symbol embedding \mathbf{s} is learned in the first phase using symbol constraint, and so these symbol embedding vectors serve as entry points for external knowledge. In order to bound the α_j^{symb} , we apply the *sigmoid* activation on them and multiply them by 2. Thus the 53 confidence scores of the classifiers are ranging from 0 to 2, in which 0 means the associated classifier is not useful. To train the knowledge branch, we use the Adam optimizer, and a learning rate of 0.01. Based on this setting, we train both the branches jointly (with main branch freeze) for 5,000 steps and cross-validate to get the best model.

In order to show that our ADVISE really learns the importance of different classifiers, we show the confidence scores of the top-10 most useful and bottom-10 least useful symbol classifiers in Figure 1. The figure accords to our expectation because our ADVISE model tends to weigh more on classifiers such as “smoking”, “animal cruelty” while in the Ads dataset there are only a limited number of training examples for these symbols. Thus incorporating knowledge for these not well-trained symbols is necessary.

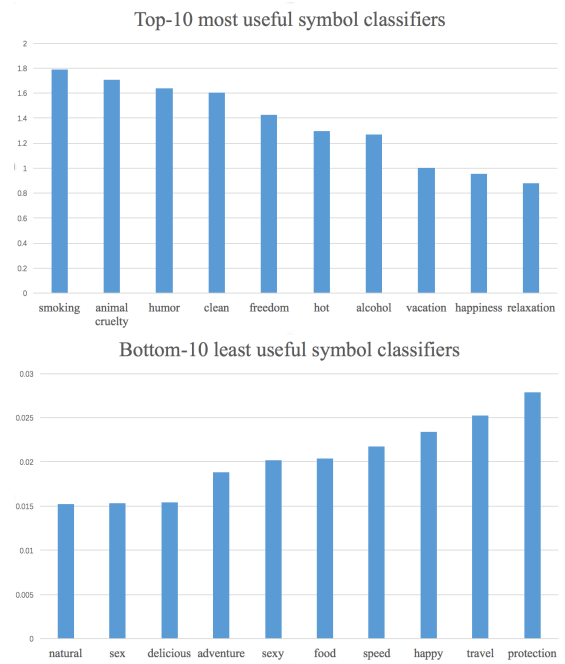


Fig. 1. Confidence scores of the symbol classifiers.

2 Adapting Hussain et al.’s method [22]

Hussain et al. [22] developed the only method we are aware of for understanding the messages of advertisements. In the main text, we show the results of the most promising way of adapting Hussain et al.’s question-answering method for our ranking task. The other ways of adapting Hussain’s method that we tried include retrieving the statement that had the highest similarity between the single-word picked in [22] and the action-reason statements, using a standard GLOVE embedding or the embedding learned in our method. We also tried concatenating the Inception-V4 and symbol features rather than pointwise-adding them, and using the original VGG features used in [22].

3 Reason to choose 50 statements

Our evaluation is similar to [3] which provides 18 choices. The task is challenging, and enlarging the list would make all methods score so poorly (e.g. low Recall@3) that their performance would be hard to compare. Table 1 demonstrates this challenge and shows that the 3 correct statements are internally more similar than a correct and an incorrect statement. We computed D_{within}^{min} , $D_{between}^{min}$, D_{within}^{avg} , $D_{between}^{avg}$ as follows, where \mathbf{x} is an image in Ads dataset \mathbb{D} , $P(\mathbf{x})$ is the set of related statements for \mathbf{x} , $N(\mathbf{x})$ are the randomly sampled statements from images other than \mathbf{x} , and $g(\cdot)$ computes average GLOVE embedding.

$$\begin{aligned}
 D_{within}^{min} &= \text{avg}_{\mathbf{x} \in \mathbb{D}} \min_{a, b \in P(\mathbf{x}), a \neq b} \|g(a) - g(b)\|_2^2 \\
 D_{within}^{avg} &= \text{avg}_{\mathbf{x} \in \mathbb{D}} \text{avg}_{a, b \in P(\mathbf{x}), a \neq b} \|g(a) - g(b)\|_2^2 \\
 D_{between}^{min} &= \text{avg}_{\mathbf{x} \in \mathbb{D}} \min_{a \in P(\mathbf{x}), b \in N(\mathbf{x})} \|g(a) - g(b)\|_2^2 \\
 D_{between}^{avg} &= \text{avg}_{\mathbf{x} \in \mathbb{D}} \text{avg}_{a \in P(\mathbf{x}), b \in N(\mathbf{x})} \|g(a) - g(b)\|_2^2
 \end{aligned}$$

Table 1. The reason to choose 50 statements. We compute L2 distance between statements belonging to same (within) and different images (between). With more candidates, sampled negatives are hard to distinguish: the min “between” distance becomes similar to the “within” distance.

# statements	D_{within}^{min}	D_{within}^{avg}	$D_{between}^{min}$	$D_{between}^{avg}$
10			1.619	1.937
50	1.395	1.551	1.437	1.938
200			1.324	1.938

4 Hard negative mining

In Table 2, we show that negative mining strategy does matter in our task. Using 32 hard negatives is better than using all or just the most challenging example.

Table 2. Hard negative mining with different top-k hyperparam, using batch size 128.

Method	Rank (Lower ↓ is better)		Recall@3 (Higher ↑ is better)	
	PSA	Product	PSA	Product
1 NEGATIVE	6.033 (± 0.127)	4.647 (± 0.026)	0.853 (± 0.015)	1.091 (± 0.003)
32 NEGATIVES (VSE++)	4.139 (± 0.094)	3.110 (± 0.019)	1.197 (± 0.017)	1.510 (± 0.004)
ALL NEGATIVES (VSE)	4.155 (± 0.091)	3.202 (± 0.019)	1.146 (± 0.017)	1.447 (± 0.004)

5 Region-based v.s. standard attention

We wish to verify that bottom-up region-based attention is more appropriate for our task than standard attention. We refer to the AMC model [6] to implement the baseline standard attention mechanism: we divide the image into 3×3 grids, apply Inception-v4 [58] to get features per cell, obtain regional features (including original image as proposal), then predict attention distribution. Thus we keep the basic feature, number of regions, and resolution the same as for our method. Table 3 shows two ablations of our method and standard attention, for three tasks. Note +ATTENTION from the main text is not our contribution; SYMBOL REGION is. Standard attention is inferior. For PSAs, which is our focus, our symbol-based attention is the strongest attention method overall.

Table 3. Region-based vs. standard image attention. SYMBOL REGION uses mean pooling, REGION-BASED ATT (+ATTENTION in Tab.2 of paper) uses attention pooling over region features, STANDARD ATT applies attention pooling on evenly split 3×3 grids.

Method	Statement				Slogan				Clustering
	Rank ↓		Recall@3 ↑		Rank ↓		Recall@3 ↑		Homogen ↑
	PSA	Product	PSA	Product	PSA	Product	PSA	Product	All
SYMBOL REGION	3.174	2.539	1.442	1.697	3.774	3.344	1.121	1.182	0.331
REGION-BASED ATT	3.258	2.488	1.428	1.726	3.850	3.257	1.155	1.205	0.355
STANDARD ATT	3.382	2.482	1.415	1.720	3.954	3.320	1.073	1.185	0.339

6 Different types of proxies

In Table 4 and Table 5, we show that not any type of label would suffice as constraint. In particular, the Ads Dataset includes 6 times more topic labels that could be used as constraints compared to symbol labels (64,325¹ vs 10,493² images annotated; almost every image is annotated with topics yet only around 15% images are annotated with symbols). Despite this, symbol labels give much greater benefit. Thus, [44]’s proxy approach is not enough; the type of labels must be carefully chosen.

Table 4. Different types of labels as constraints. The baseline method (“No extra components”) uses image attention (Sec. 3.3) but does not have the components from Sec. 3.4-3.5 of the main text.

Method	Rank (Lower ↓ is better)		Recall@3 (Higher ↑ is better)	
	PSA	Product	PSA	Product
No extra components	3.258 (± 0.081)	2.488 (± 0.015)	1.428 (± 0.017)	1.726 (± 0.004)
Symbol labels	3.171 (± 0.081)	2.465 (± 0.015)	1.471 (± 0.017)	1.726 (± 0.004)
Topic labels	3.186 (± 0.079)	2.477 (± 0.015)	1.456 (± 0.017)	1.728 (± 0.004)

Table 5. % improvement for different types of labels as constraints.

Method	PSA		Product	
	Rank ↓	Rec@3 ↑	Rank ↓	Rec@3 ↑
Symbol labels	3	3	1	0
Topic labels	2	2	0	0

¹ The Ads dataset [22] involves 64,832 images with topic annotations. However, the actual topic labels (64,325) we can use as the constraint is the intersection of the topic annotations and the format-filtered (Sec. 4.2) statement annotations.

² There are 13,938 images annotated with symbols in Ads dataset [22], and we use these images to train our region proposal network in Sec. 3.3. However, not all of the free-formed symbol annotations could be transformed into the 53 symbols that [22] used. 10,493 is the number of images that have symbol annotations for the 53 symbols.

7 In-depth quantitative results

Quantitative results with extra measurement. We provide Rank and Recall@3 as the measurement of our model in the paper. In Table 6, we also compute three other metrics: Recall@10, which denotes the number of correct statements ranked in the Top-10; RankAvg, the average ranking value of the averaged-ranked true matching statement; RankMedian, the average ranking value of the median-ranked true matching statement.

The reason that we only use Rank in the paper instead of RankAvg or RankMedian is that we found that in the Ads dataset [22], there are some really noisy annotations that ruin the metrics of RankAvg and RankMedian. This could be seen from Figure 6 where RankAvg is always worse than RankMedian.

Table 6. Our main result with extra measurements. The best method is shown in **bold**. For Recall@3 and Recall@10, higher values (\uparrow) are better. For Rank, RankAvg, and RankMedian, lower values (\downarrow) are better.

		VSE++	HUSSAIN-RANKING	ADVICE (ours)
Product	\uparrow Recall@3	1.510 (\pm 0.004)	1.515 (\pm 0.004)	1.725 (\pm 0.004)
	\uparrow Recall@10	2.379 (\pm 0.003)	2.392 (\pm 0.003)	2.527 (\pm 0.003)
	\downarrow Rank	3.110 (\pm 0.019)	3.093 (\pm 0.019)	2.469 (\pm 0.015)
	\downarrow RankAvg	7.311 (\pm 0.029)	7.122 (\pm 0.028)	6.143 (\pm 0.025)
	\downarrow RankMedian	6.392 (\pm 0.030)	6.297 (\pm 0.029)	5.252 (\pm 0.026)
PSA	\uparrow Recall@3	1.197 (\pm 0.017)	1.258 (\pm 0.017)	1.509 (\pm 0.017)
	\uparrow Recall@10	2.089 (\pm 0.017)	2.151 (\pm 0.017)	2.323 (\pm 0.015)
	\downarrow Rank	4.139 (\pm 0.094)	3.854 (\pm 0.088)	3.013 (\pm 0.075)
	\downarrow RankAvg	9.424 (\pm 0.135)	8.718 (\pm 0.127)	7.553 (\pm 0.119)
	\downarrow RankMedian	8.268 (\pm 0.143)	7.741 (\pm 0.135)	6.394 (\pm 0.121)

Per-topic evaluation. We provide per-topic quantitative results to further compare our ADVICE model and the two strong baselines: VSE++ [11] and HUSSAIN-RANKING [22]. Figures 2 and 3 show the per-topic evaluation results of the statement ranking and hard-statement ranking task respectively. In both figures, we show bar charts of the best performing 10 topics (left) and the worst performing 10 topics (right), in terms of the Rank measurement of our ADVICE model (shown in blue).

Intuitively, the hard-statement ranking task defined in our paper emphasizes more the ability to distinguish the within-topic nuances, yet our main ranking task focus on both the within-topic and between-topic differences. We see from Figures 2 and 3 that ads of different topics perform differently on the two tasks. For example “baby” (baby products) ranks the first in the main ranking task, yet it is also in the worst-10 in terms of Rank in the hard-statement ranking task. Our explanation is that “baby” ads are quite distinguishable from others, yet they do not have sub-categories within themselves. Another example is “clothing”, which has good performance on both tasks. The reason is that “clothing” ads

are distinguishable from others, while they can still be further classified as, e.g., “jean”, “watches”, “outfit”, and so on.



Fig. 2. Bar chart of the statement ranking task. The x-axis denotes the Rank measurement. Lower is better. The error bar shows the standard error, which is defined as δ/\sqrt{n} where δ is the standard deviation of the Rank and n is the number of examples. Our ADVISE model is always better than the other two baselines.

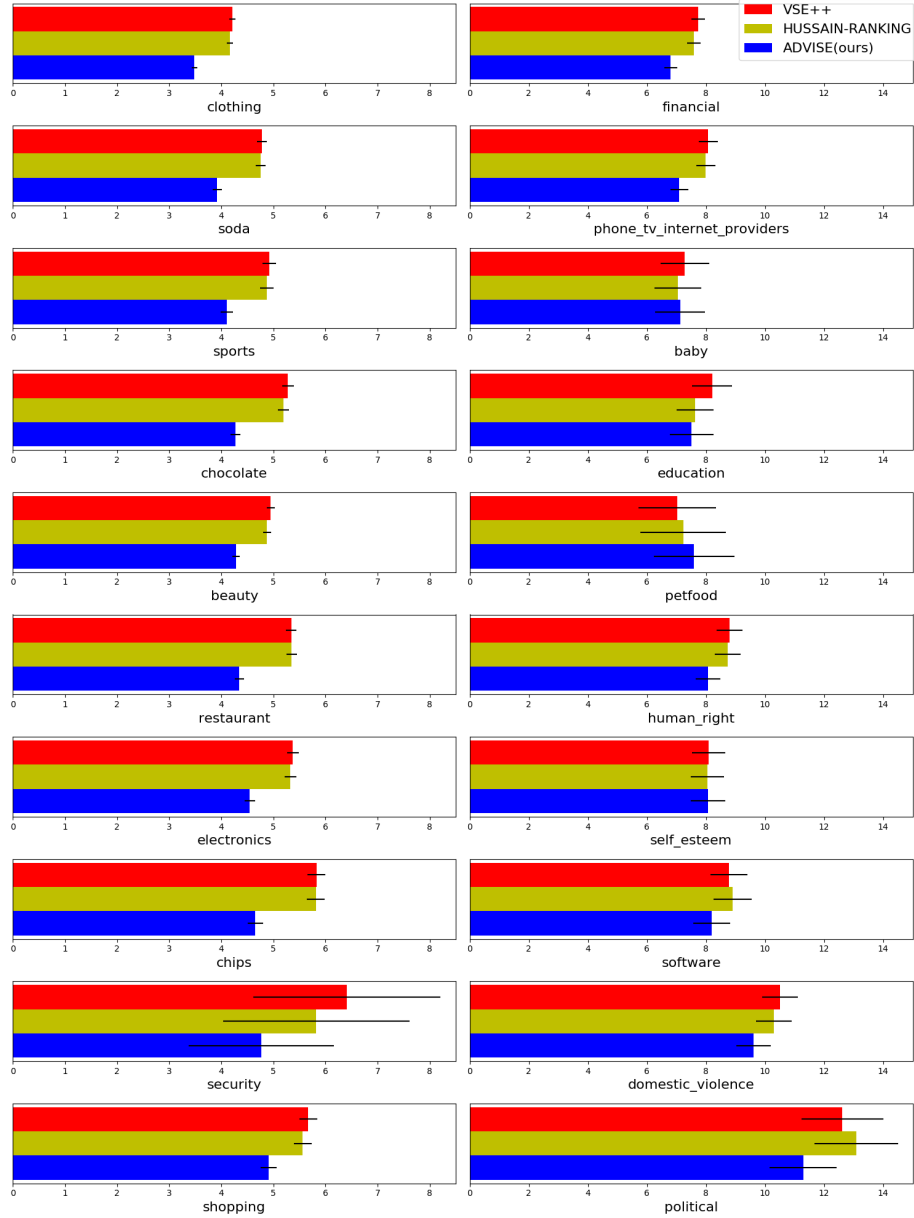


Fig. 3. Bar chart of the hard-statement ranking task. The x-axis denotes the Rank measurement. The error bar shows the standard error, which is defined as δ/\sqrt{n} where δ is the standard deviation of the Rank and n is the number of examples. Our ADVISE model is better than the two baselines in most of the cases.

8 Ranking task - qualitative results

We provide more qualitative examples of both the statement ranking task and the hard-statement ranking task for both PSAs (Figure 4) and product ads (Figure 5). We can see from both figures, that hard-statement ranking task requires the model to have a deep understanding not only about the topic and the purpose of the ads but also the details such as the brand of product, the reasoning of causal relation, etc. We also see from the qualitative examples that failure in the hard-statement ranking task does not always mean we failed to understand the ads. For example, the top-1 hard-statement prediction of row 1 of Figure 4 has already captured all of the information in the ad, “I should not be on my phone while driving because it can cause an accident”, yet this statement is from another similar image which causes the evaluation to not count this prediction as correct.

We see some interesting examples showing that the model understands both the images and statements reasonably well. For example, the result of row 3 in Figure 4, our ADVISE model ranks safety-related statements higher in the results, moreover, the top-ranked statements all involve the keyword “helmet” probably because the model associates the watermelon with head/helmet. Understanding ads is still challenging, and the result of row 4 in Figure 5 shows one obstacle. Our model should recognize beer, yet the ‘pepsi’ mentioned in the second-highest ranked statement is visually quite similar to the bottle in the image thus misleads the model.







		STATEMENT	HARD-STATEMENT
1		<p>[0.809] I should vote because my opinion matters</p> <p>[0.829] I should play soccer because it's godly</p> <p>[0.841] I should not use my phone while driving because it is unsafe and I could hurt myself</p> <p>[0.843] I should shop at Gap because I will look more stylish</p> <p>[0.871] I should order Papa John's pizza because I can get a free 2 liter of soda with my order</p>	<p>[0.713] I should not be on my phone while driving because it can cause an accident</p> <p>[0.732] I should buckle-up for me and my children because that is how we get home safely.</p> <p>[0.746] I should not drink and drive because drinking and driving kills</p> <p>[0.794] I should not drink and drive because it is dangerous</p> <p>[0.802] I should think when biking because thinking leads to wearing a helmet and living</p> <p>[0.637] I should protect the environment because we are all in danger from increasing sea water levels</p> <p>[0.652] I should support Greenpeace because they care about the environment</p> <p>[0.665] I should protect nature because pollution harms the earth</p> <p>[0.666] I should be cautious of global warming before everything dies like the dinosaurs</p> <p>[0.688] I should use this because i like scientific</p>
2		<p>[0.665] I should protect nature because pollution harms the earth</p> <p>[0.747] I should support nature because its beautiful</p> <p>[0.750] I should oppose environmental damage because it will kill the animals</p> <p>[0.760] I should spay/neuter my pet because I love him and want to give him the best life possible.</p> <p>[0.781] I should quit smoking because it's expensive</p> <p>[0.630] I should wear a helmet because it will prevent brain damage</p>	<p>[0.594] I should wear my bike helmet because I want to live</p> <p>[0.630] I should wear a helmet because it will prevent brain damage</p> <p>[0.631] I should look before turning while driving because there may be some one riding a bike</p> <p>[0.690] I should be careful while riding a motorcycle because many of the crash issues are preventable</p> <p>[0.698] I should drive slowly because then I wont get into an accident</p> <p>[0.511] I should follow this announcement in an effort to adopt an animal.</p> <p>[0.521] I should go vegan because it is the only way to protect my family from poisonous meat</p> <p>[0.558] I should become a vegan because it is fashionable.</p> <p>[0.586] I should adopt a vegan lifestyle because vegans are less likely to get diseases</p> <p>[0.590] I should research more about this because I know nothing about it</p> <p>[0.553] I should get help if my partner is abusing me because it is not acceptable for a partner to abuse me</p> <p>[0.578] I should be supporting this message because violence should not be tolerated</p> <p>[0.604] I should support Amnesty International because they are helping to stop domestic violence.</p> <p>[0.607] I should want to not hurt anybody because I dont have time for domestic violence</p> <p>[0.613] I should read flyer because I want to help fight violence against women</p> <p>[0.574] I should not let the color of my skin dictate what I should do for a living because it doesn't define who I am.</p> <p>[0.636] I should speak my truth because women matter</p> <p>[0.660] I should help amnest international because they support victims of domestic violence including their children.</p> <p>[0.660] I should help amnest international because they support victims of domestic violence including their children.</p> <p>[0.729] I should support women's rights because it is the right thing to do</p>
3		<p>[0.741] I should put a helmet on my child because its preventative for head injuries</p> <p>[0.791] I should put a helmet on my child because I don't want my child's head to end up like that melon</p> <p>[0.869] I should but always because it will hold up to leaks</p> <p>[0.898] I should eat Munch Nuts because I will go crazy over them</p>	<p>[0.698] I should drive slowly because then I wont get into an accident</p> <p>[0.511] I should follow this announcement in an effort to adopt an animal.</p> <p>[0.521] I should go vegan because it is the only way to protect my family from poisonous meat</p> <p>[0.558] I should become a vegan because it is fashionable.</p> <p>[0.586] I should adopt a vegan lifestyle because vegans are less likely to get diseases</p> <p>[0.590] I should research more about this because I know nothing about it</p> <p>[0.553] I should get help if my partner is abusing me because it is not acceptable for a partner to abuse me</p> <p>[0.578] I should be supporting this message because violence should not be tolerated</p> <p>[0.604] I should support Amnesty International because they are helping to stop domestic violence.</p> <p>[0.607] I should want to not hurt anybody because I dont have time for domestic violence</p> <p>[0.613] I should read flyer because I want to help fight violence against women</p> <p>[0.574] I should not let the color of my skin dictate what I should do for a living because it doesn't define who I am.</p> <p>[0.636] I should speak my truth because women matter</p> <p>[0.660] I should help amnest international because they support victims of domestic violence including their children.</p> <p>[0.660] I should help amnest international because they support victims of domestic violence including their children.</p> <p>[0.729] I should support women's rights because it is the right thing to do</p>
4		<p>[0.606] I should be more conscious of global warming because it is real and dangerous</p> <p>[0.637] I should buy this scent because it is fancy</p> <p>[0.717] I should focus on conservation because Global warming is real and its happening</p> <p>[0.841] I should shop at ARHAUS furniture because it will make the world more beautiful</p> <p>[0.845] I should know about Dell because it has many features</p>	<p>[0.698] I should report domestic abuse because there is support</p> <p>[0.744] I should write a letter to a soldier overseas because it will help his morale to think of and remember all the little things that made him feel good</p> <p>[0.774] You should get a massage because it will relieve period pain</p> <p>[0.782] I should go to Mount Sinai if I have prostate cancer because I will survive and be with my wife</p> <p>[0.515] I should buy Max Factor makeup because it is what celebrities use</p> <p>[0.729] I should support women's rights because it is the right thing to do</p> <p>[0.783] I should buy Lancome perfume because it will make me feel happy</p> <p>[0.801] I should get the most affordable insurance because esurance could save me money</p> <p>[0.835] I should wear Dior because Sharon Stone wears it</p>
5		<p>[0.698] I should report domestic abuse because there is support</p> <p>[0.744] I should write a letter to a soldier overseas because it will help his morale to think of and remember all the little things that made him feel good</p> <p>[0.774] You should get a massage because it will relieve period pain</p> <p>[0.782] I should go to Mount Sinai if I have prostate cancer because I will survive and be with my wife</p>	<p>[0.515] I should buy Max Factor makeup because it is what celebrities use</p> <p>[0.729] I should support women's rights because it is the right thing to do</p> <p>[0.783] I should buy Lancome perfume because it will make me feel happy</p> <p>[0.801] I should get the most affordable insurance because esurance could save me money</p> <p>[0.835] I should wear Dior because Sharon Stone wears it</p>
6		<p>[0.515] I should buy Max Factor makeup because it is what celebrities use</p> <p>[0.729] I should support women's rights because it is the right thing to do</p> <p>[0.783] I should buy Lancome perfume because it will make me feel happy</p> <p>[0.801] I should get the most affordable insurance because esurance could save me money</p> <p>[0.835] I should wear Dior because Sharon Stone wears it</p>	<p>[0.515] I should buy Max Factor makeup because it is what celebrities use</p> <p>[0.729] I should support women's rights because it is the right thing to do</p> <p>[0.783] I should buy Lancome perfume because it will make me feel happy</p> <p>[0.801] I should get the most affordable insurance because esurance could save me money</p> <p>[0.835] I should wear Dior because Sharon Stone wears it</p>

Fig. 4. Statement/Hard-statement ranking task of PSA ads. We show the top-5 multiple-choice answers ranked by our ADVISE model for both statement and hard-statement ranking task. Statements in **bold** are the correct predictions.

		STATEMENT	HARD-STATEMENT
1		<p>[0.350] I should buy Revlon makeup because they are pretty and natural</p> <p>[0.355] I should use Revlons lip balms and mascara because it will enhance the look of my lips and lashes</p> <p>[0.392] I should buy Revlon makeup because it will enhance my features</p> <p>[0.444] I should use Heinz because it does not have unnatural things in it</p> <p>[0.614] I should drink this bacardi because it makes the world seem different</p>	<p>[0.330] I should use Men Degree because it feels great</p> <p>[0.346] I should buy this product because it's sexy.</p> <p>[0.346] I should buy this lipstick because it has cool shades</p> <p>[0.349] I should buy this product because its use keeps you alert</p> <p>[0.350] I should buy Revlon makeup because they are pretty and natural</p> <p>[0.449] I should wear Chanel because it's sophisticated</p>
2		<p>[0.468] I should buy Polo clothes because it will make me handsome</p> <p>[0.477] I should buy these clothes because handsome preppies wear them</p> <p>[0.547] I should wear Polo by Ralph Lauren because it will make me attractive</p> <p>[0.595] I should eat Pringles because Pringles is down with the holidays.</p> <p>[0.609] I should purchase these security cameras from this distributor because they offer a variety of products</p>	<p>[0.468] I should buy Polo clothes because it will make me handsome</p> <p>[0.477] I should buy these clothes because handsome preppies wear them</p> <p>[0.530] I should buy a Givenchy suit because it will look cool</p> <p>[0.532] I should wear jean shorts because it's unexpected</p>
3		<p>[0.578] I should shop on Google because there is a lot of information</p> <p>[0.634] I should shop online with my phone because it is more convenient</p> <p>[0.711] I should sign up for this website to get black Friday deals because it helps you to save money when you're not paying regular price</p> <p>[0.793] I should buy westinghouse appliances because they're easy to clean</p> <p>[0.853] I should use this because it has local inventory and related items</p>	<p>[0.550] I should use a Window's Phone because I get free music for a year</p> <p>[0.553] I should download this app because I want to hear about updates on Applebees</p> <p>[0.578] I should shop on Google because there is a lot of information</p> <p>[0.600] I should consider using this software because it can allow me to create and edit ads from my smartphone</p> <p>[0.625] I should use admob because it is on my phone</p>
4		<p>[0.348] I should drink Bud Light because I have good taste</p> <p>[0.516] I should drink more Pepsi because I cannot get enough</p> <p>[0.555] I should buy this beer because I can get a boyfriend</p> <p>[0.684] I should drink Budweiser because it will make sexy women appear</p> <p>[0.690] I should want to drink this because its good for women</p>	<p>[0.348] I should drink Bud Light because I have good taste</p> <p>[0.384] I should drink this beer because it tastes good</p> <p>[0.430] I should drink Pilsner Beer because it will associate me with standards</p> <p>[0.440] I should drink Schaefer beer because it is a fine tasting beer</p> <p>[0.462] I should buy coors light because it is the worlds most refreshing can</p>
5		<p>[0.591] I should buy a volkswagen because it can get you through the toughest terrain</p> <p>[0.659] I should drive this car because its tough</p> <p>[0.686] I should buy volkswagen because it can weather the storm</p> <p>[0.723] I should be drinking Coca-Cola because I am a star</p> <p>[0.785] I should get an Ericsson because it takes pictures that look real</p>	<p>[0.536] I should use Shell because they are efficient</p> <p>[0.536] I should purchase a boles-aero because travel in style</p> <p>[0.591] I should buy a volkswagen because it can get you through the toughest terrain</p> <p>[0.621] I should get an oldsmobile because it is a summer classic</p> <p>[0.659] I should drive this car because its tough</p>
6		<p>[0.382] I should buy a cycle ops bicycle because it will get me out of the house</p> <p>[0.516] I should buy cycle oops because its healthy</p> <p>[0.779] I shouldn't wear fur because fur is unattractive on a person</p> <p>[0.882] I should buy this golfing equipment because I want to use the best and not something that will fall apart</p> <p>[0.882] I should give gifts because I can give a gift to someone religious</p>	<p>[0.312] I should visit the Museum of Science because there's a piece of history for everyone</p> <p>[0.332] I should love to surf because the ad is telling me to</p> <p>[0.382] I should buy a cycle ops bicycle because it will get me out of the house</p> <p>[0.504] I should join this website So my daughter can have a better life too</p> <p>[0.516] I should buy cycle oops because its healthy</p>

Fig. 5. Statement/Hard-statement ranking task of product ads. We show the top-5 multiple-choice answers ranked by our ADVISE model for both statement and hard-statement ranking task. Statements in **bold** are the correct predictions.

9 k-NN retrieval on image regions

After embedding each image region and weighing them using the attention mechanism, our ADVISE model learns a good image-level representation that could be used to distinguish among statements. Since the image representation is a weighted sum of region representations, regions serve as visual words in our model and the image is analogous to a sentence. In order to know if the model also has the ability to assign concept to these visual words (image regions), we made the following qualitative experiments.

Product words as query. We choose 11 discriminative words from the vocabulary of product ads and use k-NN to retrieve the 10 most related image regions from the test images. The retrieval results are shown in Figure 6. Though we manually choose the 11 words, the k-NN results are entirely generated by the ADVISE model. Please note that we have no such labels associated with regions at training time, that is, we only use the label-agnostic bounding box location annotations (we ignore the symbol categorical labels and semantic meaning of the box region), and the image-level statement annotations, to train the model. However, the model itself successfully learns to associate the concept with a specific region.

PSAs words as query. Aligning abstract words from PSA ads to the image regions is a more challenging task. We show in our paper quantitatively that our ADVISE model performs worse in PSAs than that in product ads. To understand the challenge, we show several qualitative examples in Figure 7. Similar to the previous visualization of product ads, we choose multiple discriminative words and retrieve image regions using k-NN. We retrieve the top-20 image regions for each query, merge some queries (such as “kill”, “kills”, and “killing”), and manually select (since the results are not as good as that of product ads) 10 typical examples to visualize.

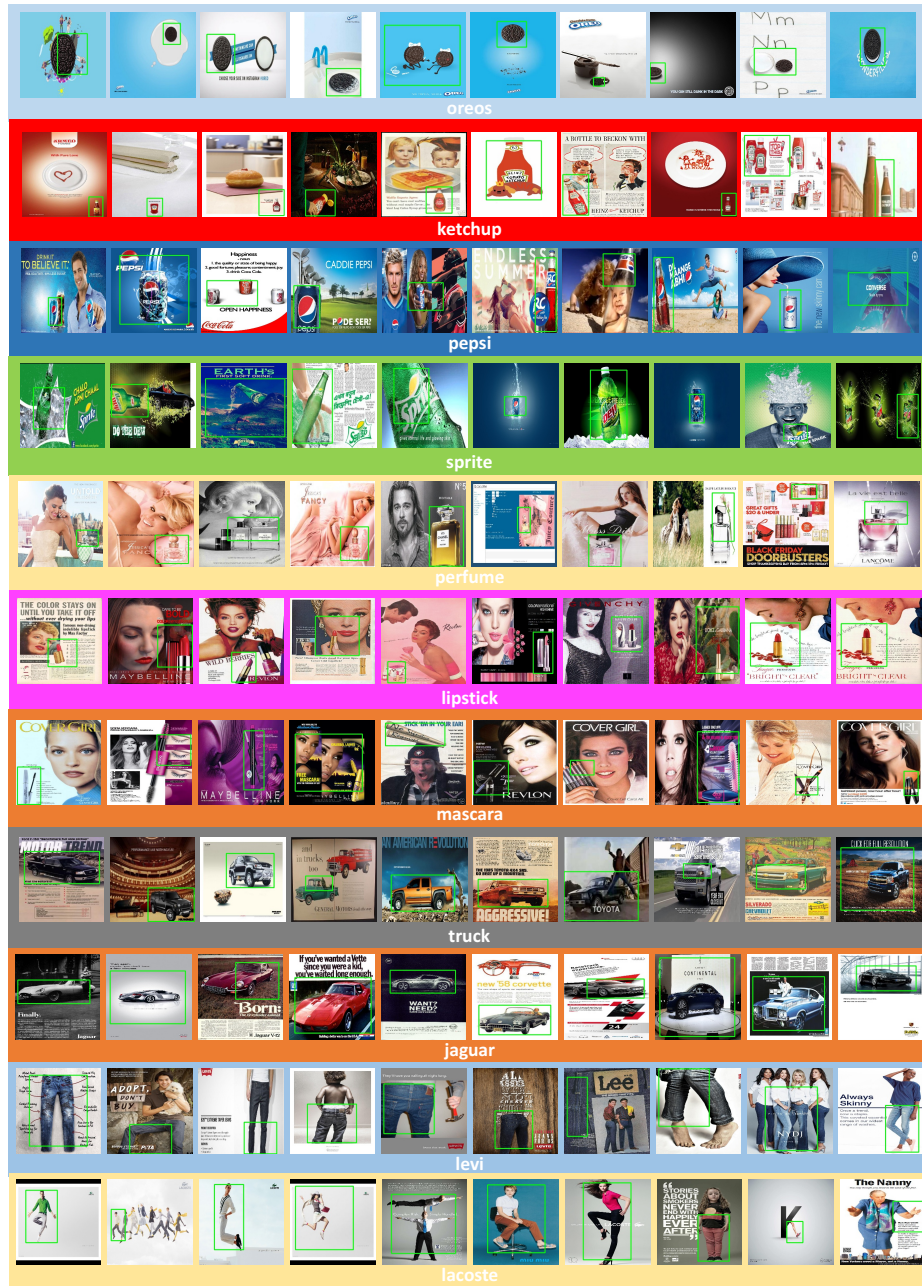


Fig. 6. Product words and the retrieved image regions. We select confusing word pairs such as “pepsi” and “sprite”, “lipstick” and “mascara”, “truck” and “jaguar”. We see that our ADVISE model makes mistakes occasionally such as retrieving “pepsi” for query “sprite”. However, the model knows the nuances in general.



Fig. 7. PSA words and the retrieved image regions. It is a more challenging task to associate PSA words with image regions since the words in PSAs tend to be more abstract than that in product ads. The qualitative examples such as “warming” and “litter” remind us that the embedding of image region and the attention mechanism may also depend on the other regions in the same image (e.g. in case that a beautiful woman exists, polar bear does not symbolize “warming” any longer), which is an interesting research direction for our future work.