

A First Look: Towards Explainable TextVQA Models via Visual and Textual Explanations

Anonymous NAACL-HLT 2021 submission

Abstract

Explainable deep learning models are advantageous in many situations. Prior works mostly provide unimodal explanations through post-hoc approaches not part of the original system design. Explanation mechanisms also ignore useful textual information present in images. In this paper, we propose MTXNet, an end-to-end trainable multimodal architecture to generate multimodal explanations which focuses on text in the image. We curate a novel dataset TextVQA-X, containing ground truth visual and multi-reference textual explanations that can be leveraged during both training and evaluation. We then quantitatively show that training with multimodal explanations complements model performance and surpasses unimodal baselines by up to 7% in CIDEr scores and 2% in IoU. More importantly, we demonstrate that the multimodal explanations are consistent with human interpretations, help justify the models decision, and provide useful insights to help diagnose an incorrect prediction. Finally, we describe a real-world e-commerce application for using the generated multimodal explanations.

1 Introduction

The ability to explain decisions through voice, text and visual pointing, is inherently human. Deep learning models on the other hand are rather opaque black boxes that don't reveal very much about how they arrived at a specific prediction. Recent research effort, aided by regulatory provisions such as GDPRs "right to explanation" (Goodman and Flaxman, 2017), have focused on peeking beneath the hood of these black boxes and designing systems that inherently enable explanation. Our goal with explainable multimodal architectures is to automate and reduce effort required for compliance and product detail checks in the e-commerce space. Rather than associates manually auditing products in a warehouse, product images can be automatically captured at scale, and passed through models

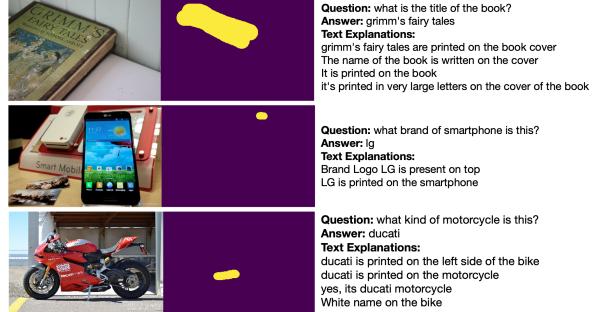


Figure 1: Sample Ground Truth Labels

that detect discrepancies. This can then be provided as evidence to help improve customer and seller partner experiences.

We choose the TextVQA task proposed by Singh et al. (2019) for realizing the system. Our choice is motivated by two reasons. First, the task is multimodal and is naturally suited for generation of multimodal explanations. Second, the task specifically focuses on the text in the image, known to encode essential information for scene understanding and reasoning (Hu et al., 2020), and allows for better quality of explanations including the text recognized. Several approaches have been proposed for the TextVQA task (Singh et al., 2019; Hu et al., 2020; Mishra et al., 2019; Biten et al., 2019; Kant et al., 2020), but they do not include a means for explaining the model decision. In addition to allowing humans to interpret the model's decision, we believe the explanations can also provide valuable insight into what component could be improved.

Most prior explanation approaches (Hendricks et al., 2016, 2018; Li et al., 2018) have been unimodal, generic and do not focus on text in the image. Only recently Huk Park et al. (2018) and Wu and Mooney (2019) generated multimodal explanations for the VQA and Activity Recognition tasks. They curated datasets (VQA-X, ACT-X) consisting of single reference ground truth textual explanations and relied on implicit attention-based visual

001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025

026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042

043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071

explanations without any access to labeled visual ground truth. However, their models cannot read and incorporate text in the image into the explanations. In addition, it is debatable whether attention mechanisms are indeed explanations (Wiegreffe and Pinter, 2019; Jain and Wallace, 2019). Moreover, other works (Das et al., 2017) have shown that current VQA attention models do not seem to look at the same regions as humans, resulting in inconsistent explanations.

The goal of our work is two-fold. First, to collect a multimodal explanations dataset (TextVQA-X) thereby highlighting the need to curate datasets where explanations are not post-hoc but part of the initial interpretable model design. Non post-hoc explanations which may not be faithful to the model decision but are in line with human explanations are still beneficial to end users. Figure 1 provides a representative example. And second, to implement a multimodal explanation system that has the ability to not only read and reason about the text in the image, but more importantly justify its decision with natural language and visually highlight the evidence, useful to even non-experts (Miller et al., 2017). The explanations and model decision must be tightly coupled and mutually influence each other through an end-to-end trainable architecture. In summary, our contributions are as follows:

- We present TextVQA-X, a novel dataset of human annotated multimodal explanations that includes ground truth segmentation maps and multi-reference textual explanations containing text in the image. We will make the dataset publicly available upon acceptance of the paper. (Section 3)
- We propose the first end-to-end trainable MTXNet architecture that produces high quality textual and visual explanations, focusing on the text in the image. (Section 4)
- Qualitative and quantitative results show that textual and visual explanations help justify a models decision and help diagnose the reasons for an incorrect prediction. (Section 5)
- We describe a real-world e-commerce system that can leverage the multimodal explanations and also highlight its challenges. (Section 6)

2 Related Work

VQA / TextVQA. The VQA task (Antol et al., 2015) has received a lot of research attention in

terms of both datasets (Antol et al., 2015; Johnson et al., 2017; Hudson and Manning, 2019) and methods (Anderson et al., 2018; Ben-Younes et al., 2017; Lu et al., 2019). Oftentimes however, these models predict an answer without completely understanding the question and do not change answers across images (Agrawal et al., 2016). Further, they ignore the text in the image and tend to focus on visual components such as objects. To address this limitation, the TextVQA task was proposed by Singh et al. (2019) and has received recent research attention (Kant et al., 2020; Hu et al., 2020; Biten et al., 2019; Mishra et al., 2019). However, not having reliable explanation mechanisms that focus on the text in the image, as part of the system design makes it difficult to diagnose prediction failures. Our work, thus allows for better diagnosis of model failures through explanations in line with human interpretations and focus on text in the image.

Explanations. Prior explanation approaches (Shortliffe and Buchanan, 1975; Van Lent et al., 2004; Zeiler and Fergus, 2014; Goyal et al., 2016; Ribeiro et al., 2016; Selvaraju et al., 2017; Das et al., 2017) focus on parts of the input that is relevant to the models decision, but not on explicitly generating explanations as model predictions. Hendricks et al. (2016, 2018) were the first to generate natural language justifications for image classifiers. Unlike our model however, explanations are unimodal and there are no reference human explanations. Closer to our objective Huk Park et al. (2018) generate multimodal explanations and curate a new VQA-X dataset. Wu and Mooney (2019) extend their work to ensure explanations can be traced back to an object ensuring local faithfulness. However, their explanations do not contain the text in the image. They use implicit attention for visual explanations and have no access to visual ground truth during training. Further, they use a single textual explanation reference during training. In contrast, our work incorporates multimodal explanations which focuses on the text in the image.

3 TextVQA-X Dataset

To train and evaluate multimodal explanation models that focus on text in the image, we collect the TextVQA-X dataset by human annotation of a subset of samples from the TextVQA dataset (Singh et al., 2019).

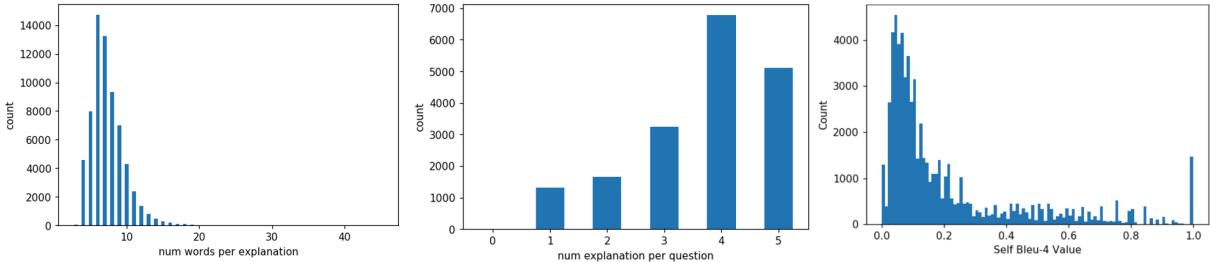


Figure 2: TextVQA-X Dataset Statistics

3.1 Ground Truth Label Collection

We used the Sagemaker Ground Truth (Amazon-AWS, 2018) platform to create a labeling task for gathering visual and textual explanations. Human annotators were asked to provide a single textual explanation that answers the question "Why do you think <answer> is the correct answer for the given question and image pair?". Specific instructions added that annotators should try to incorporate the answer and/or the text in image as part of their explanation. The annotators were also asked to make use of a brush to segment image regions relevant to both the answer and written explanation. Sample annotations are shown in Figure 1. Each image and question pair can have up to 5 distinct human annotators allowing for multi-reference training and evaluation (Zheng et al., 2018). A single segmentation map is obtained by using a threshold of 0.5 obtained as an average over all annotations. Bad actors were identified and most were removed through a combination of heuristics and manual checks. Overall, we collected more than 67K explanations among over 800 unique workers.

3.2 TextVQA Explanation Dataset (TextVQA-X).

Dataset Statistic	Value
Num. Unique Images	11681
Num. Questions	18096
Num. Unique Questions	15374
Num. Visual Explanations	67043
Num. Textual Explanations	67043
Num. Unique Textual Explanations	61458
Avg. Num Textual Explanations per Question	3.71
Avg. Words per Textual Explanation	7.36
Avg. Characters per Textual Explanation	36.92
Textual Explanation Vocab Size	17910

Table 1: TextVQA-X Dataset Summary

In order to obtain a measure of the quality of explanations and to help filter out bad actors, we make use of the Self-BLEU-4 metric (Zhu et al.,

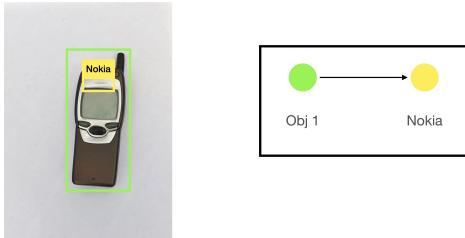
2018). The Self-BLEU score is used to measure how one sentence resembles the rest in a generated collection by regarding one sentence as the hypothesis and the rest as references. A higher Self-BLEU score implies higher similarity of the hypothesis with all the references. A lower Self-BLEU implies higher diversity and lesser overlap. Although we would like to have several diverse textual explanations, we noticed that most good textual explanation annotations have overlap with others. The average Self-BLEU-4 across all annotations was 0.21 indicating consistent overlap and quality.

Comparison with VQA-X and VQA-HAT datasets. With respect to textual explanations, the TextVQA-X includes multi-references with an average of 3.71 explanations for each QA pair that can be utilized for both training and testing. In contrast, VQA-X (Huk Park et al., 2018) contains an average of 1.27 explanations with a single textual explanation for QA pairs in the training set and three textual explanations for test/val QA pairs. VQA-HAT (Das et al., 2017) does not include textual explanations. As far as visual explanations are concerned, there are a number of distinctions among these datasets. First, both VQA-X and VQA-HAT are defined on the VQA task which does not require reading text in image. In contrast, the TextVQA-X is specifically designed to focus on text in the image. Second, TextVQA-X includes one ground truth visual explanation for both training and testing (total 67K), whereas VQA-X includes explanations only as part of testing for a small random subset (total 6K). And third, similar to VQA-X, TextVQA-X annotators were asked to directly segment the relevant image region. On the contrary, VQA-HAT annotations were collected by having humans unblur the images and are more likely to introduce noise when irrelevant regions are uncovered.

238 4 Multimodal Text-in-Image 239 Explanation Network (MTXNet)

240 We design our Multimodal Text-in-Image Explanation
241 Network (MTXNet) to allow for end-to-end
242 multitask training of answer prediction, text
243 generation and semantic segmentation extending
244 the M4C model proposed in (Hu et al., 2020). In the
245 subsequent subsections we describe each of the
246 individual components in more detail.

247 4.1 Graph Attention Network (GAT)



248 Figure 4: An example of how to build the graph

249 Many questions in the TextVQA dataset require
250 the model to acknowledge the spatial relationship
251 between objects and OCR tokens. To better encode
252 the relationship between objects and OCR tokens
253 and subsequently generate better quality explana-
254 tions, we leverage graph neural networks. We build
255 the graph using only the visual inputs (object and
256 OCR region bounding boxes). Each object location
257 and OCR token is treated as a node in the graph.
258 Whenever the bounding box associated with node
259 i is contained in the node j , we add an edge from
260 node j to node i . An example is presented in Figure
261 4. We then make use of the Graph Attention Net-
262 work (GAT) (Veličković et al., 2017) to operate on
263 the structured data. Unlike Graph Convolutional
264 Networks (GCN) (Kipf and Welling, 2016) that
265 treat each adjacent node equally, GATs incorpo-
266 rates attention into the layer-wise propagation rule
267 and allows the model to variably weight adjacent
268 nodes based on relevancy.

269 4.2 Multimodal Transformer (MMT)

270 The multimodal transformer operates on three
271 modalities - question words, visual objects and
272 OCR tokens. The feature definition's are identical
273 to that proposed in M4C (Hu et al., 2020) with the
274 addition of textual explanation embeddings whose
275 embedding process resembles that of the question
276 words. We extract features and project them to a

276 common d -dimensional semantic space used for de-
277 coding and prediction. The prediction takes place
278 through a dynamic pointer network (Vinyals et al.,
279 2015) that allows to either predict from a fixed vo-
280 cabulary or from OCR tokens extracted from the
281 image.

282 4.3 Multireferences for Textual Explanations

283 Neural text generation tasks such as machine trans-
284 lation, image captioning and summarization typi-
285 cally only consider a single reference for each ex-
286 ample during training (Zheng et al., 2018). In our
287 case however, considering just a single reference
288 for training is insufficient because of the inherently
289 subjective nature of textual explanations. Thus we
290 leverage the multi-references we have collected in
291 the TextVQA-X dataset during both training and
292 evaluation. We use the *sample one* technique for
293 incorporating multi-references during training. We
294 randomly pick one of the available references in
295 each training epoch.

296 4.4 Visual Explanations through Semantic 297 Segmentation

298 Visual explanations are obtained through a seman-
299 tic segmentation module (Feature Pyramid Net-
300 work - FPN (Kirillov et al., 2017)). They are made
301 an explicit and natural component of end-to-end
302 training by leveraging ground truth label supervi-
303 sion. To incorporate the multimodal embedding
304 from the MMT into the segmentation module, we
305 reshape, pad and concatenate the output with the
306 raw input image along the channel. Thus, the over-
307 all input channels for the segmentation module in-
308 creases to five, with 3 color channels and 2 multi-
309 modal channels. The output of the segmentation
310 model is a continuous mask with higher value im-
311 plying greater relevancy to the inputs. The mask
312 may be binarized through thresholding.

313 4.5 Training

314 The MTXNet architecture is end-to-end trainable
315 with three distinct tasks (1) answer prediction (2)
316 textual explanation generation and (3) visual expla-
317 nation through semantic segmentation. We ensure
318 cross modal feedback between the textual expla-
319 nations and predicted answers by leveraging a phased
320 training process where we randomly choose be-
321 tween one of three choices (1) predict answer then
322 textual explanation (2) predict textual explanation
323 then answer and (3) predict both answer and textual
324 explanation simultaneously. Each task corresponds

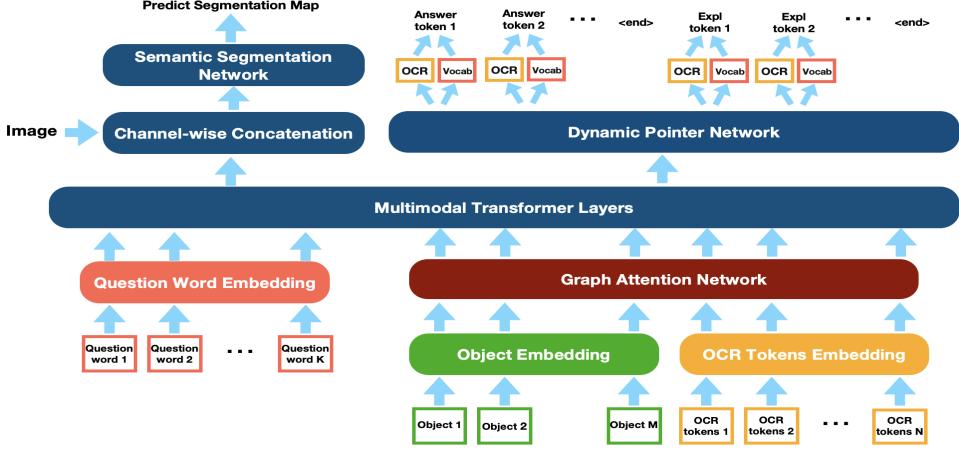


Figure 3: Our Multimodal Text-in-Image Explanation Model (MTXNet) architecture generates multimodal explanations. Explanations and Answers are utilized as part of the iterative autoregressive decoding procedure.

to an individual part of the training objective. For answer prediction (\mathcal{L}_{ans}) and textual explanation generation (\mathcal{L}_{text}) we use the *binary cross entropy with logits*¹. For semantic segmentation (\mathcal{L}_{vis}) we use the *dice loss* (Sudre et al., 2017). The naive approach to combine multiple losses is to use a predetermined weighted linear sum of the individual losses. However, the model performance is sensitive to the weights which are hyperparameters and expensive to tune. We thus use a multitask learning loss with homoscedastic uncertainty as proposed by Kendall et al. (2018). The overall objective is present in Equation 1. The weights $\{w_{ans}, w_{text}, w_{vis}\}$ corresponding to the loss terms of the three individual tasks are learned.

$$\mathcal{L} = \sum_i \mathcal{L}_i \exp(-w_i) + w_i, i \in \{ans, text, vis\} \quad (1)$$

5 Experiments

In this section, we detail the experimental setup, present quantitative results with ablations and finally analyze qualitative results.

5.1 Experimental Setup

This subsection discusses the dataset splits, model training, hyperparameter settings and evaluation metrics.

Dataset Splits. We use the TextVQA-X dataset described in Section 3. We choose a random 80/20 split for train and test. The dataset split statistics are present in Table 2. Each question is associated with a single image, one or more textual explanation and

a single visual explanation. The OCR tokens and object regions are already present in the original TextVQA dataset.

Split	#Img.	#Ques.	#Text Expl.	#Vis. Expl.
train	10379	14475	53536	14475
test	3354	3619	13507	3619

Table 2: Train / Test Splits of TextVQA-X Dataset

Preprocessing. The dynamic pointer network is allowed to choose between a fixed 5000 word vocabulary and a maximum of 100 OCR tokens per image. The text explanations and answers are capped to a maximum length of 16 and 12 tokens respectively. For the visual explanations, we use a FPN decoder with ResNeXt50 encoder and $320 \times 320 \times 5$ input feature size.

Model training and hyperparameters. We train the MTXNet model end-to-end in a supervised setting using the Pythia² framework. We use a batch size of 128 and train for a maximum of 8500 epochs using Adam optimizer. The learning rate is set to $1e - 4$ with no weight decay. The best model is selected using test set answer accuracy with an evaluation granularity of every 100 epochs. The entire training task varies from 14-20 hours on 8 Nvidia K80 GPUs.

Evaluation Metrics. Each question in the TextVQA dataset has 10 human annotated answers, and the predicted answer accuracy is measured via a soft voting in accordance with the VQA task evaluation script³. We evaluate the textual explanations

¹<https://pytorch.org/docs/stable/generated/torch.nn.BCEWithLogitsLoss.html>

²<https://github.com/facebookresearch/mmf>

³<https://visualqa.org/evaluation>

380 using the standard BLEU-4 (Papineni et al., 2002),
 381 ROUGE (Lin, 2004), METEOR (Banerjee and
 382 Lavie, 2005) and CIDEr (Vedantam et al., 2015)
 383 metrics computed with the coco-caption⁴
 384 code . All the text generation metrics account for
 385 multi-references by averaging the individual scores.
 386 Finally, we evaluate the visual explanations using
 387 IoU (Intersection over Union) score with a thresh-
 388 old of 0.5.

389 5.2 Ablation Study

390 We ablate MTXNet and compare quantitatively
 391 with a related model on our TextVQA-X dataset
 392 through automatic evaluations for answers and ex-
 393 planations. The results are present in Table 3.

394 **Comparison with existing baselines.** We com-
 395 pute the performance of the baseline model M4C
 396 (Hu et al., 2020) on the TextVQA-X test set (with-
 397 out explanations) and obtain an answer accuracy
 398 of 35.23%. Using the MTXNet architecture and
 399 evaluating on the TextVQA-X test set, we obtain
 400 an answer accuracy of 36.27%. The addition of
 401 explanations thus complements the MTXNet per-
 402 formance.

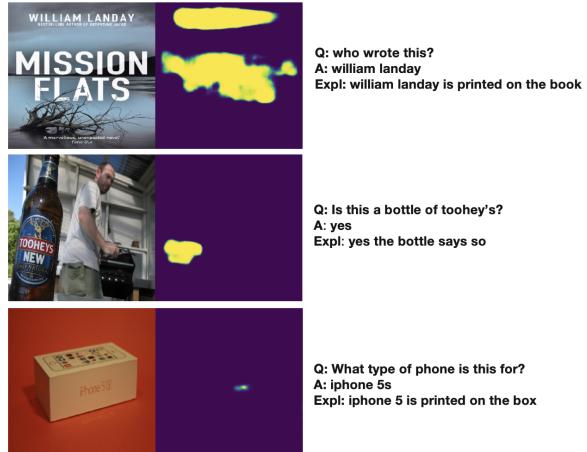
403 **Unimodal vs. Multimodal explanations** We no-
 404 tice that each modality mutually influences the
 405 other as the model learns to jointly optimize for
 406 both modalities of explanations and the answer pre-
 407 diction. Excluding visual explanations results in
 408 the largest drop of up to 7% in CIDEr scores of the
 409 textual explanations. Similarly, the absence of text
 410 explanations results in a 2% drop in IoU of visual
 411 explanations. More importantly, we notice that the
 412 multimodal explanations provide visual and tex-
 413 tual rationale into a models decision. This further
 414 accentuates the value of designing multimodal ex-
 415 planation systems.

416 **GAT better captures structural dependencies.**
 417 The removal of GAT from the MTXNet architec-
 418 ture adversely impacts the quality of explanations
 419 and answers. The greatest drop of 7% is observed
 420 for the CIDEr metric. We believe the GAT helps
 421 better encode the relationship between objects and
 422 OCR tokens enhancing the relationship reasoning
 423 ability. The image region corresponding to the text
 424 is also highlighted better as seen in the 2% increase
 425 in IoU when GAT is included in MTXNet.

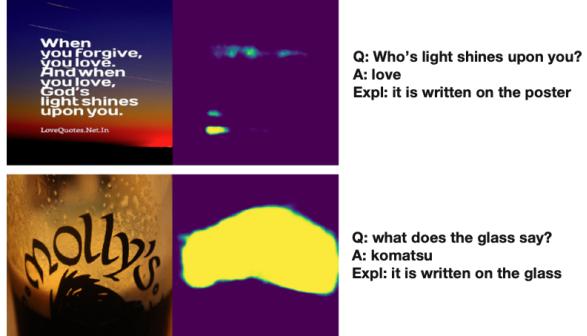
426 **Multi-reference training improves text genera-
 427 tion.** Training with multi-references significantly
 428 outperforms training with a single randomly chosen

429 sample fixed for all epochs. The largest increase
 430 of up to 25% was noticed in CIDEr score, with the
 431 increase being consistent across all text generation
 432 metrics. This underscores the benefits of having
 433 multi-references for both training and evaluation
 434 and designing systems that utilize this effectively.

435 5.3 Qualitative Samples



436 Figure 5: Examples where the MTXNet model pro-
 437 duces high quality explanations.



438 Figure 6: Examples where the MTXNet model fails.

439 As can be seen in Figure 5, the MTXNet is able
 440 to accurately answer the given question while also
 441 justifying its decision through textual and visual
 442 explanations. In certain cases, the OCR engine
 443 could be inaccurate and lead to wrong tokens being
 444 predicted, but the overall answer and explanations
 445 are correct. Figure 6 depicts two failure cases. The
 446 upper subimage indicates this could be due to incor-
 447 rect visual localization while the lower subimage in-
 448 dicates a potential OCR prediction error, although
 449 the visual explanation is correct. Despite being
 446 generic and dull the textual explanations are cor-
 447 rect. In other cases, the model fails due to incorrect
 448 visual localization as seen in Figure 7.

⁴<https://github.com/tylin/coco-caption>

Ablation	Approach	Visual Explanation		Textual Explanation			
		IoU		B	R	M	C
No visual explanation (VE)	MTXNet (GAT + MR + TE)	-		25.16	47.63	21.76	88.43
No textual explanation (TE)	MTXNet (GAT + MR + VE)	16.10		-	-	-	-
No graph attention (GAT)	MTXNet (MR + TE + VE)	16.55		27.87	49.28	21.61	88.57
No multireferences (MR)	MTXNet (GAT + TE + VE)	17.52		5.92	28.05	11.65	70.60
Consolidated architecture	MTXNet (GAT + MR + TE + VE)	18.86		31.07	53.87	22.06	95.07

Table 3: Quantitative Evaluation of Answer and Explanations. All metrics are in %. VE: visual explanation, TE: textual explanation, GAT: graph attention network, MR: multi-references. Evaluated automatic metrics: Intersection over Union (IoU), BLEU-4 (B), METEOR (M), ROUGE (R), CIDEr (C).



Figure 7: Example where the explanation is consistent with an incorrect prediction.

Explanations help explain incorrect decisions of model. In Figure 7, we see that the right answer to the question is “target”. However, the model predicts “dollar tree”. From the visual and textual explanations we see that the image region localized is incorrect and the model fails to grasp the meaning of “fading”. This potentially results in it focusing on the more prominent “dollar tree” text. Such an analysis provides insights into the component of the system that is failing and deserves further attention.

6 Applications to E-Commerce Businesses

E-commerce businesses need to comply with industry-wide, and country-specific regulations, to provide accurate and useful information of products to improve customer experience that leads to more business. Our long-term goal with explainable multimodal architectures is to automate and reduce manual effort required for compliance and product detail checks. This will enable businesses to scale compliance and customer experience improvement efficiently without linear increases in cost. Further, these architectures help validate if models are performing as intended and used for the right purposes.

A potential customer experience issue arises when the physical product in a warehouse is different from that uploaded by a seller on the product details page. A possible reason could be that the seller or manufacturer labeled the product erroneously when they packaged it. Many sellers taking advantage of lower cost manufacturing in a global supply chain may not be able to audit every batch of product leaving the factory. Such discrepancies will almost certainly lead to product returns, because the customer didn’t get what they wanted and increases costs. Such discrepancies may also be due to more nefarious reasons, such as opportunistic bad actors taking advantage of sellers that have successful products by introducing poorer quality or mismatched offers at a lower price to unsuspecting customers. Examples of compliance issues include detecting products that contain batteries and chemicals to comply with transportation and logistics regulations, as well as identifying products that require additional safety documentation and checks, such as products that may have unintended use by children (e.g. toys and products that may end up as toys should not have heavy metals or other poisons that cause illness or death when accidentally ingested). While not all answers can be obtained with product images alone, manual investigation processes utilize these images to identify potential risks that warrant additional steps in the process (e.g. lab testing).

Rather than manually auditing products in a warehouse, product images can be automatically captured at scale, and passed through models that detect such discrepancies. With the help of subject matter experts, attributes such as quantity, color and brand names, and other common misleading attributes are identified apriori. Relevant questions that target these attributes are formulated. The image and question are then inputs to a multimodal explainable system (such as MTXNet) that can provide an answer and justify its prediction through

517 multimodal explanations. Answers can then be
518 compared against the information extracted from
519 the product detail pages on the website. Any dis-
520 crepancies found can be noted and a selling partner
521 can be provided evidence through the multimodal
522 explanations to take corrective steps.

523 An example use-case is as follows. Given a large
524 container of cereal, with smaller boxes within, poten-
525 tial questions are: "How many cereal boxes are
526 within the container?". This information is usu-
527 ally written on the larger container present in the
528 warehouse and can be answered based on reading
529 the text in the image. If there is any discrepancy
530 encountered in the number of boxes of cereal in the
531 warehouse and that listed on the website, appropri-
532 ate action can be taken. Other similar questions
533 include: "How heavy is the product?", "Is the chair
534 red?", "Does the item contain allergens?", and "Did
535 the product pass the lead test?".

536 As with most explainability techniques, the chal-
537 lenges with the use of such explainable systems are
538 two-fold. First, since there can be multiple stake-
539 holders with diverse expertise and expectations, we
540 need to clearly define the level of abstraction at
541 which they interact with the system. For instance,
542 while a scientist can use the explanations to im-
543 prove the model, a business operations associate
544 may use the explanations to identify and audit prod-
545 uct discrepancies. Second, we need fine grained
546 evaluation methodologies and metrics that take into
547 account the stakeholders as well.

548 References

- 549 Aishwarya Agrawal, Dhruv Batra, and Devi Parikh.
550 2016. Analyzing the behavior of visual question
551 answering models. *arXiv preprint arXiv:1606.07356*.
- 552 Jon Almazán, Albert Gordo, Alicia Fornés, and Ernest
553 Valveny. 2014. Word spotting and recognition with
554 embedded attributes. *IEEE transactions on pattern
555 analysis and machine intelligence*, 36(12):2552–
556 2566.
- 557 Amazon-AWS. 2018. SageMaker Ground Truth.
558 [https://aws.amazon.com/sagemaker/
559 groundtruth/](https://aws.amazon.com/sagemaker/groundtruth/).
- 560 Peter Anderson, Xiaodong He, Chris Buehler, Damien
561 Teney, Mark Johnson, Stephen Gould, and Lei
562 Zhang. 2018. Bottom-up and top-down attention for
563 image captioning and visual question answering. In
564 *Proceedings of the IEEE conference on computer vi-
565 sion and pattern recognition*, pages 6077–6086.
- 566 Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Mar-
567 garet Mitchell, Dhruv Batra, C Lawrence Zitnick,

568 and Devi Parikh. 2015. Vqa: Visual question an-
569 swering. In *Proceedings of the IEEE international
570 conference on computer vision*, pages 2425–2433.

571 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hin-
572 ton. 2016. Layer normalization. *arXiv preprint
573 arXiv:1607.06450*.

574 Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An
575 automatic metric for mt evaluation with improved
576 correlation with human judgments. In *Proceedings
577 of the acl workshop on intrinsic and extrinsic eval-
578 uation measures for machine translation and/or sum-
579 marization*, pages 65–72.

580 Hedi Ben-Younes, Rémi Cadene, Matthieu Cord, and
581 Nicolas Thome. 2017. Mutan: Multimodal tucker
582 fusion for visual question answering. In *Proceed-
583 ings of the IEEE international conference on com-
584 puter vision*, pages 2612–2620.

585 Ali Furkan Biten, Ruben Tito, Andres Mafla, Lluis
586 Gomez, Marçal Rusinol, Ernest Valveny, CV Jawahar,
587 and Dimosthenis Karatzas. 2019. Scene text
588 visual question answering. In *Proceedings of the
589 IEEE International Conference on Computer Vision*,
590 pages 4291–4301.

591 Piotr Bojanowski, Edouard Grave, Armand Joulin, and
592 Tomas Mikolov. 2017. Enriching word vectors with
593 subword information. *Transactions of the Associa-
594 tion for Computational Linguistics*, 5:135–146.

595 Mark G Core, H Chad Lane, Michael Van Lent, Dave
596 Gomboc, Steve Solomon, and Milton Rosenberg.
597 2006. Building explainable artificial intelligence
598 systems. In *AAAI*, pages 1766–1773.

599 Abhishek Das, Harsh Agrawal, Larry Zitnick, Devi
600 Parikh, and Dhruv Batra. 2017. Human attention in
601 visual question answering: Do humans and deep net-
602 works look at the same regions? *Computer Vision
603 and Image Understanding*, 163:90–100.

604 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
605 Kristina Toutanova. 2018. Bert: Pre-training of deep
606 bidirectional transformers for language understand-
607 ing. *arXiv preprint arXiv:1810.04805*.

608 David K Duvenaud, Dougal Maclaurin, Jorge Ipar-
609 raguirre, Rafael Bombarell, Timothy Hirzel, Alán
610 Aspuru-Guzik, and Ryan P Adams. 2015. Convolu-
611 tional networks on graphs for learning molecular
612 fingerprints. In *Advances in neural information pro-
613 cessing systems*, pages 2224–2232.

614 Wenqi Fan, Yao Ma, Qing Li, Yuan He, Eric Zhao,
615 Jiliang Tang, and Dawei Yin. 2019. Graph neural
616 networks for social recommendation. In *The World
617 Wide Web Conference*, pages 417–426.

618 Akira Fukui, Dong Huk Park, Daylen Yang, Anna
619 Rohrbach, Trevor Darrell, and Marcus Rohrbach.
620 2016. Multimodal compact bilinear pooling for
621 visual question answering and visual grounding.
622 *arXiv preprint arXiv:1606.01847*.

- 623 Chenyu Gao, Qi Zhu, Peng Wang, Hui Li, Yuliang Liu, 678
 624 Anton van den Hengel, and Qi Wu. 2020. Structured 679
 625 multimodal attentions for textvqa. *arXiv preprint*
 626 *arXiv:2006.00753*.
- 627 Bryce Goodman and Seth Flaxman. 2017. European 680
 628 union regulations on algorithmic decision-making 681
 629 and a “right to explanation”. *AI magazine*, 38(3):50– 682
 630 57.
- 631 Yash Goyal, Tejas Khot, Douglas Summers-Stay, 683
 632 Dhruv Batra, and Devi Parikh. 2017. Making the 684
 633 v in vqa matter: Elevating the role of image under-
 634 standing in visual question answering. In *Proceed- 685
 635 ings of the IEEE Conference on Computer Vision 686
 636 and Pattern Recognition*, pages 6904–6913.
- 637 Yash Goyal, Akrit Mohapatra, Devi Parikh, and Dhruv 687
 638 Batra. 2016. Towards transparent ai systems: Interpreting 688
 639 visual question answering models. *arXiv preprint arXiv:1608.08974*. 689
- 641 Lisa Anne Hendricks, Zeynep Akata, Marcus 690
 642 Rohrbach, Jeff Donahue, Bernt Schiele, and Trevor 691
 643 Darrell. 2016. Generating visual explanations. In 692
 644 *European Conference on Computer Vision*, pages 693
 645 3–19. Springer.
- 646 Lisa Anne Hendricks, Ronghang Hu, Trevor Darrell, 694
 647 and Zeynep Akata. 2018. Grounding visual explana- 695
 648 tions. In *Proceedings of the European Conference 696
 649 on Computer Vision (ECCV)*, pages 264–279.
- 650 Ronghang Hu, Amanpreet Singh, Trevor Darrell, and 697
 651 Marcus Rohrbach. 2020. Iterative answer prediction 698
 652 with pointer-augmented multimodal transform- 699
 653 ers for textvqa. In *Proceedings of the IEEE/CVF 700
 654 Conference on Computer Vision and Pattern Recog- 701
 655 nition*, pages 9992–10002.
- 656 Drew A Hudson and Christopher D Manning. 2019. 702
 657 Gqa: a new dataset for compositional question 703
 658 answering over real-world images. *arXiv preprint 704
 659 arXiv:1902.09506*, 3(8).
- 660 Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, 705
 661 Anna Rohrbach, Bernt Schiele, Trevor Darrell, and 706
 662 Marcus Rohrbach. 2018. Multimodal explanations: 707
 663 Justifying decisions and pointing to the evidence. In 708
 664 *Proceedings of the IEEE Conference on Computer 709
 665 Vision and Pattern Recognition*, pages 8779–8788.
- 666 Sarthak Jain and Byron C Wallace. 2019. Attention is 710
 667 not explanation. *arXiv preprint arXiv:1902.10186*.
- 668 Justin Johnson, Bharath Hariharan, Laurens van der 711
 669 Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross 712
 670 Girshick. 2017. Clevr: A diagnostic dataset for 713
 671 compositional language and elementary visual reasoning. 714
 672 In *Proceedings of the IEEE Conference on Compu- 715
 673 ter Vision and Pattern Recognition*, pages 2901– 716
 674 2910.
- 675 Yash Kant, Dhruv Batra, Peter Anderson, Alex 717
 676 Schwing, Devi Parikh, Jiasen Lu, and Harsh 718
 677 Agrawal. 2020. Spatially aware multimodal 719
 678 transformers for textvqa. *arXiv preprint arXiv:2007.12146*.
- 680 Andrej Karpathy and Li Fei-Fei. 2015. Deep visual- 720
 681 semantic alignments for generating image descrip- 721
 682 tions. In *Proceedings of the IEEE conference on 722
 683 computer vision and pattern recognition*, pages 723
 684 3128–3137.
- 685 Alex Kendall, Yarin Gal, and Roberto Cipolla. 2018. 724
 686 Multi-task learning using uncertainty to weigh 725
 687 losses for scene geometry and semantics. In *Pro- 726
 688 ceedings of the IEEE conference on computer vision 727
 689 and pattern recognition*, pages 7482–7491.
- 690 Thomas N Kipf and Max Welling. 2016. Semi- 728
 691 supervised classification with graph convolutional 729
 692 networks. *arXiv preprint arXiv:1609.02907*.
- 693 Alexander Kirillov, Kaiming He, Ross Girshick, and Pi- 730
 694otr Dollár. 2017. A unified architecture for instance 731
 695 and semantic segmentation.
- 696 H Chad Lane, Mark G Core, Michael Van Lent, Steve 732
 697 Solomon, and Dave Gomboc. 2005. Explainable arti- 733
 698 ficial intelligence for training and tutoring. Technical 734
 699 report, UNIVERSITY OF SOUTHERN CALI- 735
 700 FORNIA MARINA DEL REY CA INST FOR CRE- 736
 701 ATIVE
- 702 Qing Li, Jianlong Fu, Dongfei Yu, Tao Mei, and Jiebo 737
 703 Luo. 2018. Tell-and-answer: Towards explainable 738
 704 visual question answering using attributes and cap- 739
 705 tions. *arXiv preprint arXiv:1801.09041*.
- 706 Chin-Yew Lin. 2004. Rouge: A package for automatic 740
 707 evaluation of summaries. In *Text summarization 741
 708 branches out*, pages 74–81.
- 709 Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan 742
 710 Lee. 2019. Vilbert: Pretraining task-agnostic visi- 743
 711 olinguistic representations for vision-and-language 744
 712 tasks. In *Advances in Neural Information Process- 745
 713 ing Systems*, pages 13–23.
- 714 Tim Miller, Piers Howe, and Liz Sonenberg. 2017. Ex- 746
 715 plainable ai: Beware of inmates running the asyl- 747
 716 um or: How i learnt to stop worrying and love 748
 717 the social and behavioural sciences. *arXiv preprint 749
 718 arXiv:1712.00547*.
- 719 Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, 750
 720 and Anirban Chakraborty. 2019. Ocr-vqa: Visual 751
 721 question answering by reading text in images. In 752
 722 *2019 International Conference on Document Analy- 753
 723 sis and Recognition (ICDAR)*, pages 947–952. IEEE.
- 724 Kishore Papineni, Salim Roukos, Todd Ward, and Wei- 754
 725 Jing Zhu. 2002. Bleu: a method for automatic eval- 755
 726 uation of machine translation. In *Proceedings of the 756
 727 40th annual meeting of the Association for Compu- 757
 728 tational Linguistics*, pages 311–318.
- 729 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine 758
 730 Lee, Sharan Narang, Michael Matena, Yanqi Zhou, 759
 731 Wei Li, and Peter J Liu. 2019. Exploring the limits

- 732 of transfer learning with a unified text-to-text trans- 788
 733 former. *arXiv preprint arXiv:1910.10683*.
 734 Marco Tulio Ribeiro, Sameer Singh, and Carlos 789
 735 Guestrin. 2016. "why should i trust you?" explain- 790
 736 ing the predictions of any classifier. In *Proceed- 791
 737 ings of the 22nd ACM SIGKDD international con- 792
 738 ference on knowledge discovery and data mining*, 793
 739 pages 1135–1144.
 740 Ramprasaath R Selvaraju, Michael Cogswell, Ab- 794
 741 hishek Das, Ramakrishna Vedantam, Devi Parikh, 795
 742 and Dhruv Batra. 2017. Grad-cam: Visual explana- 796
 743 tions from deep networks via gradient-based local- 797
 744 ization. In *Proceedings of the IEEE international 798
 745 conference on computer vision*, pages 618–626.
 746 Edward H Shortliffe and Bruce G Buchanan. 1975. A 801
 747 model of inexact reasoning in medicine. *Mathemat- 802
 748 ical biosciences*, 23(3-4):351–379.
 749 Amanpreet Singh, Vivek Natarajan, Meet Shah, 803
 750 Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, 804
 751 and Marcus Rohrbach. 2019. Towards vqa models 805
 752 that can read. In *Proceedings of the IEEE Confer- 806
 753 ence on Computer Vision and Pattern Recognition*,
 754 pages 8317–8326.
 755 Carole H Sudre, Wenqi Li, Tom Vercauteren, Sebastien
 756 Ourselin, and M Jorge Cardoso. 2017. Generalised
 757 dice overlap as a deep learning loss function for
 758 highly unbalanced segmentations. In *Deep learn- 761
 759 ing in medical image analysis and multimodal learn- 760
 760 ing for clinical decision support*, pages 240–248.
 761 Springer.
 762 Michael Van Lent, William Fisher, and Michael Man-
 763 cuso. 2004. An explainable artificial intelligence
 764 system for small-unit tactical behavior. In *Proced- 767
 765 ings of the national conference on artificial intelli- 766
 766 gence*, pages 900–907. Menlo Park, CA; Cambridge,
 767 MA; London; AAAI Press; MIT Press; 1999.
 768 Ramakrishna Vedantam, C Lawrence Zitnick, and Devi
 769 Parikh. 2015. Cider: Consensus-based image de-
 770 scription evaluation. In *Proceedings of the IEEE*
 771 *conference on computer vision and pattern recogni-*
 772 *tion*, pages 4566–4575.
 773 Petar Veličković, Guillem Cucurull, Arantxa Casanova,
 774 Adriana Romero, Pietro Lio, and Yoshua Bengio.
 775 2017. Graph attention networks. *arXiv preprint*
 776 *arXiv:1710.10903*.
 777 Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly.
 778 2015. Pointer networks. In *Advances in neural in- 779
 779 formation processing systems*, pages 2692–2700.
 780 Sarah Wiegreffe and Yuval Pinter. 2019. Atten-
 781 tion is not explanation. *arXiv preprint*
 782 *arXiv:1908.04626*.
 783 Jialin Wu and Raymond Mooney. 2019. Faithful mul-
 784 timodal explanation for visual question answering.
 785 In *Proceedings of the 2019 ACL Workshop Black- 786
 786 boxNLP: Analyzing and Interpreting Neural Net- 787
 787 works for NLP*, pages 103–112.