

Optimizing Visual Question Answering Models for Driving: Bridging the Gap Between Human and Machine Attention Patterns

Anonymous CVPR submission

Paper ID *****

Abstract

Visual Question Answering (VQA) models play a critical role in enhancing the perception capabilities of autonomous driving systems by allowing vehicles to analyze visual inputs alongside textual queries, fostering natural interaction and trust between the vehicle and its occupants or other road users. This study investigates the attention patterns of humans compared to a VQA model when answering driving-related questions, revealing disparities in the objects observed. We propose an approach integrating filters to optimize the model's attention mechanisms, prioritizing relevant objects and improving accuracy. Utilizing the LXMERT model for a case study, we compare attention patterns of the pre-trained and Filter Integrated models, alongside human answers using images from the NuImages dataset, gaining insights into feature prioritization. We evaluated the models using a Subjective scoring framework which shows that the integration of the feature encoder filter has enhanced the performance of the VQA model by refining its attention mechanisms.

1. Introduction

Visual Question Answering (VQA) models are integral to autonomous driving systems as they enable vehicles to perceive and understand their surroundings by analyzing visual inputs alongside textual queries, thereby enhancing their perception capabilities. VQA models facilitate natural interaction between the vehicle and its occupants or other road users, fostering trust in autonomous technology. By enabling natural language interaction, VQA models assist in making the autonomous vehicle more transparent and understandable to the driver. When the vehicle can effectively communicate its actions, intentions, and reasoning in a language that humans understand, it fosters a sense of transparency and predictability, which are crucial for building trust.

For instance, if the vehicle encounters a challenging driv-

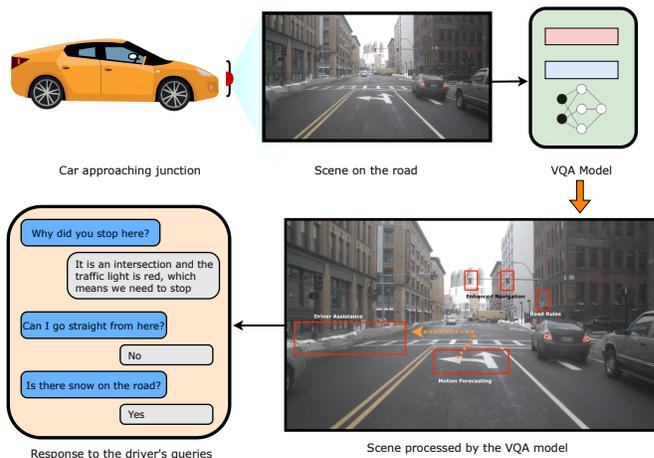


Figure 1. A demo of how VQA models work in a driving scenario. In this scenario, it can better explain its decision-making process to the driver using the VQA model. This allows the driver to better comprehend the situation and feel more confident in the vehicle's capabilities. Moreover, in situations where the driver needs clarification or wants to ask questions about the vehicle's actions or the environment, the VQA model can provide immediate responses, helping to alleviate uncertainties and concerns.

This study is focused on comparing the explanation given for object detection patterns of humans and attention patterns of a Visual Question Answering (VQA) model when answering questions related to driving. Our survey indicated that humans concentrate on objects like road lines, signboards, vehicles in the ego lane, etc when it comes to answering questions related to driving. However, when we looked at the objects observed by a VQA model, it wasn't restricted to only objects related to driving. There were objects like trees, sky, tower, etc which were irrelevant to answer a question like, "How many vehicles are in the ego lane?". The approach here is to streamline the features and objects that the VQA model is taking into consideration by adding a filter when asking a driving-related question. This will optimize the model's attention mechanisms to priori-

059 tize relevant objects and improve its accuracy in answering
060 questions.

061 It also addresses a disparity between human attention
062 patterns and those of the VQA model, aiming to enhance the
063 model's performance in the domain of driving. We are per-
064 forming a case study with a VQA model- LXMERT where
065 we look at how the pretrained model with all its features an-
066 swers a driving question and how a 'filter' integrated model
067 answers the same question while also comparing them with
068 the Human Answers that were provided by human annota-
069 tors. By comparing the attention patterns of the pretrained
070 and streamlined model, we can gain insights into how dif-
071 ferent features and objects are prioritized when answering
072 driving-related questions. This analysis can help in identi-
073 fying the factors that contribute to the models' performance
074 differences.

075 By examining attention mechanisms, we aim to elucidate
076 how VQA models prioritize visual stimuli in their decision-
077 making processes which will help us in the finetuning pro-
078 cess of our experiments.

079 2. Background Study

080 Vision transformers in a Visual Question Answering (VQA)
081 model work by dividing the image into patches and repre-
082 senting them as embeddings [5]. These embeddings, along
083 with the text embeddings of the question, are then fed into a
084 transformer architecture [5]. The transformer processes the
085 embeddings by attending to both visual and textual infor-
086 mation, enabling the model to understand the image and the
087 question simultaneously. Finally, the model generates an
088 answer based on the learned representations from the trans-
089 former layers. In [5], the authors argue that uncertainty
090 in vision is a dominating factor preventing the successful
091 learning of reasoning in vision and language problems. By
092 integrating a filter that focuses on driving-related features,
093 our approach aims to mitigate this uncertainty by providing
094 a VQA model with more relevant visual information tai-
095 lored to the context of driving-related questions.

096 In [12], they focus on improving the efficiency of visual
097 transformers by removing redundant calculations in trans-
098 former networks. Considering that the attention mechanism
099 in a transformer architecture aggregates different patches
100 layer-by-layer, the authors Yehui Tang et al. present a novel
101 'patch slimming' approach that discards useless patches in
102 a top-down paradigm. Initially, the effective patches in the
103 last layer are identified and then used to guide the patch
104 selection process of previous layers. For each layer, the im-
105 pact of a patch on the final output feature is approximated
106 and patches with less impact will be removed [12]. While
107 this could work for a vision transformer model, it is not
108 necessarily good to implement for a VQA model. Patch
109 slimming aims to improve the efficiency of the model by
110 removing redundant patches throughout the image and the

filter focuses on extracting driving-related features from the
image before passing it through the vision transformer, to
enhance the model's ability to answer driving-related ques-
tions more effectively [12]. The impact of patch slimming
on performance can be more general, affecting the overall
efficiency of the model but potentially risking loss of task-
specific information [12]. However, integrating a filter fo-
cusing on driving-related features directly aims to enhance
performance on driving-related questions by ensuring that
the model receives relevant visual information. Patch slim-
ming is a more general approach that may not adapt specifi-
cally to the requirements of the VQA task, which can result
in a loss of task-specific information. Integrating a filter
specifically designed for driving-related questions ensures
that the model prioritizes relevant features for this task,
leading to improved performance on driving-related ques-
tions while maintaining task specificity.

In [7], a novel object detection framework is proposed
that attempts to extract meaningful and representative fea-
tures across different image scales. The authors do so
by unifying atrous convolutions with a vision transformer
(DIL-ViT). The proposed model uses atrous convolutions
to generate rich multi-scale feature maps and employs a
self-attention mechanism to enrich important backbone fea-
tures [7]. This framework enhances object detection perfor-
mance which could be an excellent feature to add to a VQA
model. However, in our case, the VQA model in question
has to enhance its performance on driving-related questions
by ensuring that the model receives relevant visual infor-
mation. The filter proposed in our work specifically ex-
tracts driving-related features from the input image before
passing it through the vision transformer component of the
VQA model. While both the framework proposed in [7] and
the filter integration approach aim to enhance model perfor-
mance, they differ in their focus, purpose, task applicability,
feature extraction methods, training objectives, and adapta-
tion requirements.

In [14], the authors argue that the existing methods suf-
fer from bias in understanding the image and insufficient
knowledge to solve the problem of VQA. The authors pro-
pose a novel knowledge-based VQA framework (PROOF-
READ) that uses LLM to obtain knowledge explicitly and
the vision language model which can see the image to get
the knowledge answer and a knowledge perceiver that fil-
ters out knowledge that is deemed harmful for getting the
correct final answer [14]. PROOFREAD processes textual
knowledge obtained by a language model, filtering out ir-
relevant or harmful information before combining it with
the visual information [14] whereas our filter focuses on
processing visual information from the image, extracting
driving-related features, and integrating them into the VQA
model's processing pipeline before combining them with
textual information. The framework in [14] is designed to

164 address biases in understanding images and insufficiencies
165 in knowledge to solve VQA problems in general whereas
166 the filter proposed is tailored for improving VQA perfor-
167 mance on driving-related questions specifically, focusing
168 on extracting features relevant to driving scenarios from the
169 image input.

170 3. Proposed Methodology

171 In Visual Question Answering (VQA) models, the model
172 initially encodes the textual question into a numerical rep-
173 resentation to capture its semantic meaning. As the inquiry
174 pertains to a corresponding image, the model extracts vi-
175 sual features using convolutional neural networks (CNNs).
176 Subsequently, the feature extraction mechanism dynam-
177 ically assigns weights to different regions of the image or
178 words in the question based on their relevance to the in-
179 quiry. This weighting enables the model to selectively fo-
180 cus on informative elements while disregarding irrelevant
181 ones. Integrating the weighted features from both the im-
182 age and question encoding, typically through concatenation
183 or element-wise multiplication, the model combines visual
184 and textual information. Finally, the integrated features are
185 fed into a classifier to predict the answer, leveraging the
186 learned associations between input features and correspond-
187 ing answers from training data. Through this process, the
188 attention pattern of the VQA model adapts to the specific
189 question context, facilitating accurate and contextually rel-
190 evant responses across a diverse range of topics.

191 While the architecture of a VQA model aims to repli-
192 cate human cognition and reasoning when responding to in-
193 quires about various scenarios, there exists a gap that re-
194 quires attention. Typically, during driving, humans exhibit
195 focused attention on aspects directly related to driving, of-
196 ten disregarding peripheral details unrelated to the task at
197 hand [13]. When behind the wheel, individuals prioritize
198 observing their immediate surroundings and assessing the
199 next steps in their driving manoeuvres. This selective at-
200 tention ensures optimal performance and safety on the road.
201 For instance, if asked a question ‘Is there snow on the road?’
202 while driving, the driver’s attention would primarily be di-
203 rected towards assessing road conditions. They would ob-
204 serve the road surface for any signs of snow, focusing solely
205 on elements pertinent to their driving task. This focused
206 attention highlights a fundamental distinction between hu-
207 man perception during driving and the holistic scene under-
208 standing performed by VQA models. Therefore, bridging
209 this gap necessitates the creation of a filter that enables the
210 model to prioritize relevant information similar to human
211 attentional patterns, thereby enhancing its ability to discern
212 and respond accurately to questions posed in diverse real-
213 world driving contexts.

3.1. Object Perception and Cognitive Processes 214

215 When presented with a question about a driving scenario,
216 humans instinctively assess various factors to formulate a
217 response. They consider the context, including details like
218 location, weather, and traffic conditions, while also identi-
219 fying potential hazards such as other vehicles, pedestrians,
220 or adverse road conditions [13]. Drawing on their knowl-
221 edge of traffic rules and regulations, they analyze the sce-
222 nario through the lens of right-of-way, speed limits, and
223 relevant guidelines [9]. Decision-making involves weigh-
224 ing the available options against safety, efficiency, and legal
225 considerations, with a keen spatial awareness guiding their
226 understanding of distances and relative speeds. Through-
227 out this process, safety remains the most important concern,
228 leading to actions aimed at minimizing risks and promoting
229 responsible driving behaviour [13].

230 This complicated process has to be kept in mind while
231 designing an autonomous driving system. These learnings
232 also need to be incorporated into a VQA model if we want
233 it to answer all our questions related to driving. Achiev-
234 ing this requires a deep understanding of human attention
235 patterns, which can then be mirrored in the attention mech-
236 anisms of VQA models. We discuss in the following sec-
237 tions how aligning these attention patterns can improve the
238 effectiveness of VQA models in handling driving-related in-
239 quires.

3.1.1 Human Answer Explanation Patterns 240

241 To gain an insight into the factors humans consider when
242 given a driving scenario and posed with a question, we sur-
243 veyed ten individuals with a minimum of five years of driv-
244 ing experience. Participants were asked to provide answers
245 to questions depicted in Figures Tab. 1 and Sec. 5. The re-
246 sponses with the highest number of votes were selected as
247 the definitive answers.

248 The features observed via answers to these questions
249 were all cumulated together by asking the humans about
250 the features observed using the same questionnaire. This
251 explanation of features observed while making the decision
252 to answer the given question helped us understand the recur-
253 ring attention patterns in human observation and also com-
254 pile a list of features that are commonly useful in answering
255 driving-related questions.

3.1.2 Attention Patterns of VQA models 256

257 The attention mechanism in a Visual Question Answering
258 (VQA) model typically shows the focus or weight assigned
259 to different regions of an image. Specifically, it indicates
260 which parts of the input (such as image features or words in
261 the question) are deemed most relevant or informative for
262 answering the given question. By visualizing the attention

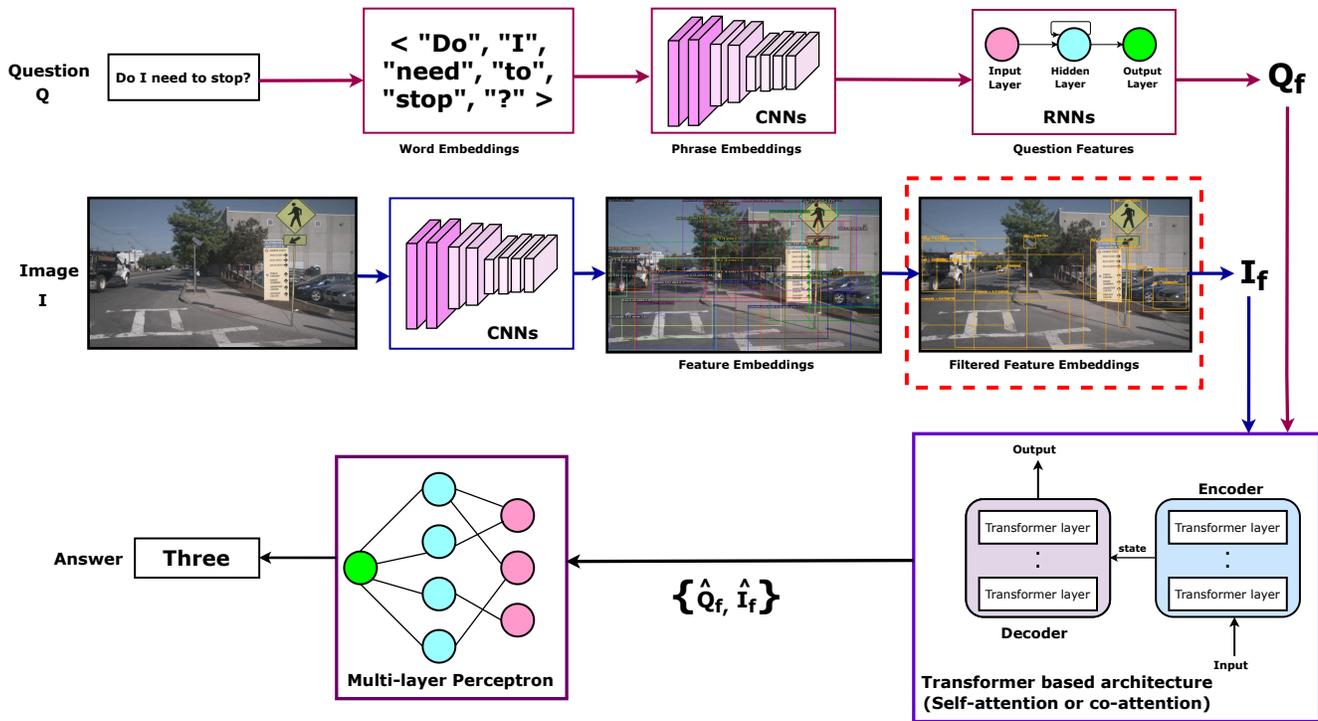


Figure 2. Refining VQA architecture: Integration of the filter into a general VQA architecture

weights, we can discern which areas of the image or words in the question the model prioritizes in its decision-making process. This helps in understanding the reasoning behind the model’s responses and provides insights into how it processes and interprets visual and textual information to generate answers.

3.2. Human-Guided Feature Filter

We cumulated the features that are being observed by humans when answering driving-related questions and incorporated them in the construction of a filter aimed at capturing pertinent visual information. The objects like roads, lines, curbs, sidewalks, crosswalks, bikes, cars, trucks, etc., are recurring features in any given driving scenario which were used when creating the filter. This filter is designed to be integrated before the vision transformer component of the VQA model, ensuring that it focuses solely on relevant driving-related features as shown in Figure 2. This approach mimics human attention patterns, thereby enhancing the model’s ability to effectively answer questions about driving scenarios by prioritizing the most relevant visual cues.

Filtering out irrelevant visual data reduces computational complexity and memory requirements, making the model more efficient and faster in processing information. This filter aligns the model’s attention with human observation patterns, and the reasoning behind its predictions becomes more interpretable and aligned with human intuition. By

emphasizing commonly observed features, it is observed in Case Studies (Section 4) that a VQA model can generalize better to new or unseen driving scenarios, enhancing its robustness and applicability in real-world settings. Prioritizing relevant visual cues related to driving can improve the safety and reliability of autonomous driving systems, ensuring they focus on critical information for making informed decisions on the road.

3.3. Filter: Algorithm and Need

The holistic approach typically employed by VQA models to capture and utilize intricate data patterns appears ineffective when narrowing the focus solely to driving-related questions. Thus, a filter is necessary to prevent the VQA model from expending computational resources on irrelevant learnings.

Integrating a filter before the vision transformer component of the VQA model helps to improve the model’s performance when asked driving-related questions, as detailed further in Section 4. The advantages of this filter are listed as follows:

- **Feature Relevance:** By incorporating a filter specifically designed to capture driving-related features, the model can prioritize and emphasize information relevant to driving tasks. This can help the model to better focus on important visual cues such as road signs, vehicles, lanes, traffic lights, and road conditions, which are crucial for understanding and answering driving-related questions.

- 317 • **Reduced Noise:** Filtering out irrelevant visual informa-
318 tion can help reduce noise in the input data, providing the
319 model with cleaner and more focused inputs. This can
320 prevent the model from being distracted by non-driving-
321 related elements in the image, leading to more accurate
322 predictions for driving-related questions.
- 323 • **Improved Attention Mechanism:** By pre-processing the
324 input with a filter targeting driving-related features, the at-
325 tention mechanism within the vision transformer compo-
326 nent can be guided to attend more effectively to relevant
327 regions of the image. This can enhance the model’s abil-
328 ity to extract and leverage important visual information
329 when generating answers to driving-related questions.
- 330 • **Enhanced Generalization:** Focusing the model’s atten-
331 tion on driving-related features during pre-processing can
332 help improve its generalization capabilities, allowing it to
333 handle better variations in driving scenarios, lighting con-
334 ditions, and camera perspectives. This can lead to more
335 robust performance across different driving-related ques-
336 tion types and real-world conditions.

337 3.3.1 Algorithm

338 The filter proposed is shown in Algorithm ???. It filters out
339 irrelevant predictions based on predefined classes, extracts
340 relevant information from the filtered predictions, converts
341 the data to suitable types, and returns the filtered features.

Algorithm 1: Feature filter for Vision Transformer
block in a VQA model

- 1 **Input:** Extract predicted classes, scores, bounding
boxes, normalized bounding boxes, and ROI
features from outputs tensor;
 - 2 **Output:** Filtered features for VQA;
 - 3 Initialize empty lists for filtered boxes, classes,
labels, indices, normalized bounding boxes, and
ROI features;
 - 4 **if** *predicted class is in a predefined list of classes*
then
 - 5 Append the box, class, label, index, normalized
 bounding box, and ROI feature to the
 corresponding lists;
 - 6 Convert filtered boxes, normalized bounding boxes,
and ROI features to suitable data types;
 - 7 **Return:** filtered boxes, classes, labels, indices,
normalized bounding boxes, and ROI features;
-

342 This process helps in focusing the model’s attention on
343 the most relevant visual features for answering questions,
344 thereby improving the overall performance of the VQA
345 model.

4. Case Study 346

347 We perform a case study by incorporating the filter into a
348 VQA model and observing the different answers before and
349 after the filter is integrated. By comparing the model’s re-
350 sponses before and after the filter’s integration, we gain a
351 clear understanding of the enhancements brought about by
352 focusing on relevant driving-related features. We examine
353 how the model’s attention patterns evolve post-filter integra-
354 tion and can discern whether they align more closely with
355 human observation patterns in driving scenarios. It is ob-
356 served that this alignment enhances the model’s ability to
357 answer driving-related questions accurately along with their
358 interpretability and generalization capabilities.

4.1. Dataset 359

360 The images collected to test the filter’s performance are
361 from the nuImages dataset. nuImages is a dataset of 93000
362 2d annotated images from a larger pool of data (nuScenes
363 dataset). The images we used are randomly selected sam-
364 ple images from nuImages. We chose two images per cam-
365 era as it allows us to evaluate the VQA model’s ability to
366 comprehend changes in perspective resulting from differ-
367 ent camera angles. This approach ensures a diverse range
368 of viewpoints, enabling a comprehensive assessment of the
369 model’s performance across various perspectives.

4.2. VQA Model: LXMERT 370

371 LXMERT (Learning Cross-Modality Encoder Represent-
372 ations from Transformers) is a large-scale Transformer
373 model that consists of three encoders: an object relationship
374 encoder, a language encoder, and a cross-modality encoder
375 [11]. The model uses the Adam optimizer with a linear-
376 decayed learning rate schedule and a peak learning rate at
377 $1e - 4$. The model is trained for 20 epochs which is roughly
378 670K4 optimization steps with a batch size of 256. The
379 pretraining of VQA tasks, however, is only for the last 10
380 epochs because this task converges faster and empirically
381 needs a smaller learning rate [11]. An illustration of the
382 networks in LXMERT is shown in Figure 4.

383 The VQA architecture in LXMERT facilitates compre-
384 hensive question-answering by integrating language and vi-
385 sual inputs. Using transformer layers of self-attention and
386 cross-attention respectively, the model encodes contextual
387 information from both textual queries and holistic visual
388 features extracted from images. Through the collaborative
389 operation of these components, with Lxmert Visual Feature
390 Encoder and Lxmert Encoder, the model achieves a holistic
391 understanding of the interplay between language and visual
392 information to generate answers. However, the holistic ap-
393 proach of visual features is not necessarily a great idea when
394 we want the model to only answer driving-related queries
395 (examples in Supplementary Material).

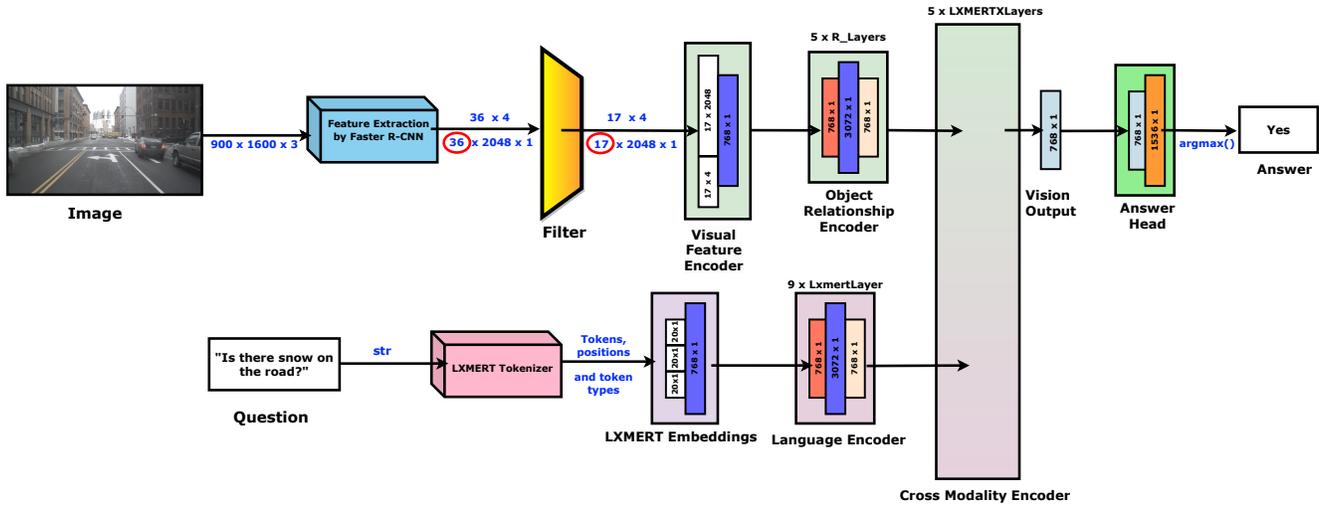


Figure 3. Visualizing the Functionality of a LXMERT with the filter integrated: An Illustrative Approach

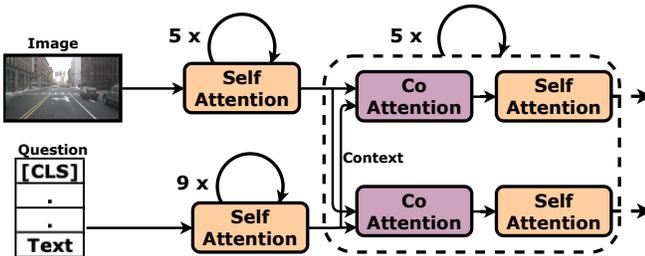


Figure 4. An Illustration of the architecture of LXMERT: self-attention with co-attention encoder

396 When we look deeper into the architecture of LXMERT,
 397 the input dimensions start with the image input, sized at 900
 398 x 1600 pixels with 3 channels (RGB). The feature extraction
 399 of the image inputted to LXMERT is done using Faster R-
 400 CNN [2], where features are taken in 36 x 2048 x 1 and the
 401 boxes are taken in 36 x 4 dimensions. The filter we pro-
 402 pose takes outputs from Faster R-CNN used in LXMERT
 403 along with parameters like device and detection threshold.
 404 It processes these outputs to extract relevant information
 405 such as predicted classes, scores, bounding boxes, normal-
 406 ized bounding boxes, and region of interest (ROI) features.
 407 It filters out predictions based on a predefined set of classes
 408 (e.g., signs, curbs, people, vehicles, etc.) using a detec-
 409 tion threshold (circled in Red), which is 17 x 2048 x 1 di-
 410 mensions. The function then returns the filtered informa-
 411 tion including filtered bounding boxes, classes, labels, in-
 412 dices, normalized bounding boxes, and ROI features which
 413 becomes the input for the Lxmert Visual Feature Encoder
 414 as shown in the Figure 3. These dimensions undergo trans-
 415 formations through convolutional layers and pooling lay-
 416 ers resulting in higher-level feature representations (eg: 17
 417 x 2048 and 3072 x 1) while reducing spatial dimensions
 418 to 768 x 1 as shown in Visual Feature Encoder and Ob-
 419 ject Relationship Encoder. This approach allows for the

420 model to learn complex and abstract representations in the
 421 intermediate layers with 3072 features, potentially captur-
 422 ing more nuanced information or patterns. Then, by reduc-
 423 ing the dimensionality back to 768 in the subsequent lay-
 424 ers (R_Layers), the model can consolidate and distil this in-
 425 formation into a more compact representation suitable for
 426 further processing or downstream tasks. Therefore, even
 427 though the input to the Cross modality Encoder has fewer
 428 features (768), the attention mechanism can still effectively
 429 capture relationships and dependencies across the input se-
 430 quence.

431 Meanwhile, the question input is initially represented as
 432 word embeddings by taking tokens, positions, and token
 433 types as the input. It undergoes text processing in the Lan-
 434 guage Encoder block to capture the semantic information
 435 of the question. This process transforms the input ques-
 436 tion into a fixed-length vector representation. After sepa-
 437 rate processing of the image and question inputs, their fea-
 438 tures are combined in the Cross Modality Encoder enabling
 439 the model to leverage both visual and textual information.
 440 This joint representation retains relevant information from
 441 both modalities, facilitating the capture of complex patterns
 442 in the data. Subsequent layers, including the Vision Out-
 443 put layer (768 x 1), further process these combined features
 444 to capture intricate relationships between visual and textual
 445 cues. Finally, a probability distribution over possible an-
 446 swers, with dimensions corresponding to the number of an-
 447 swer classes in the dataset (1536) is processed in the An-
 448 swer Head block. The final output answer with the most
 449 probability is chosen using $\text{argmax}()$.

450 Table 1 is intended to show the difference in answers be-
 451 tween an LXMERT model with and without the filter for
 452 better readability. The column features correspond to the
 453 features observed by the model when it generated that re-
 454 spective answer.

Table 1. Analyzing Feature Detection: LXMERT Pretrained Model with and without Filter vs. Human Observations

Camera	Image	Questions	Human Answers	Features	Pretrained Answers	Features	Filter added	Features
Back Camera		How many vehicles are there?	3	cars or bikes	1	tree, building, clouds, truck, road, pole, crosswalk, lines, sidewalk, scene	0	truck, road, crosswalk, lines, sidewalk
		Which camera is this image from?	Back Camera	car behind	front		unknown	
		Are there any vehicles in ego lane?	No	vehicles in the lane	Yes	scene, sky, street, light, headlights, road	No	road, car, line
		Is it safe to initiate a lane change?	unable to tell	not enough information	Yes		No	
Back Left Camera		Which camera is this image from?	unable to tell	Road Signs	front	sky, tree, pole, building, road, shadow, line, sidewalk, person	top	road, pole, line, person
		Is it okay to initiate a lane change?	No	continuous white line	Yes		No	
		Are there any vehicles coming behind?	unable to tell	not enough information	Yes	tree, sky, pole, leaves, sign, building, grass, road, vehicle, line, sidewalk, bottle	No	road, pole, sign, line
		Which street is this?	Summer Street	Name plate	unknown		unknown	
Back Right Camera		Which camera is this image from?	Back Right	road edges and markings	unknown	sky, tree, building, sign, street, sidewalk, crosswalk, vent, tire, street, car, road, pole	unknown	road, truck, crosswalk, sign, pole, car
		Do I need to stop?	No	already halfway in the turn	No		No	
		Are there any pedestrians on the sidewalk?	No	sidewalk and pedestrians	No	stret, sidewalk, building, tree, road, bike, bus, door	No	road
		Can I park on the right?	No	No parking slots	No		No	
Front Camera		Which camera is this image from?	Front Camera	road markings and vehicles in the front	front	building, water, road, car, lines, street, city	top	road, lines, water, car
		Is there snow on the road?	Yes	Road and kerb	No		No	
		Are there any pedestrians?	Yes	pedestrians	No	building, tree, street, road, sidewalk, line, van	No	road, lines, line
		Can I go right in this lane?	Yes	junction road to the right	No		No	
Front Left Camera		Which camera is this image from?	Front Left	kerb and pedestrians	front	ceiling, tree, building, pole, people, ground, sidewalk, line, man, window, person	unknown	pole, person, man, line
		Can I take a left from here or should I go straight?	Don't know	not enough information	Yes		Yes	
		How many pedestrians are there?	Ten	Pedestrians	0	sky, tree, pole, building, sign, woman, crosswalk, road, median, line, shirt, pants, man, people, person	0	road, crosswalk, person, sign, line, man
		Do I need to stop till pedestrians cross to turn left?	Yes	vehicle orientation and pedestrians	No		No	
Front Right Camera		Can I park here?	No	parking slots	No	building, bus, van, circle, car, sign, tire, ceiling	No	car, sign
		Why can't I park here?	No space	empty slots not available	parking		No	
		Which camera is this image from?	Front Right	sign boards and directions	front	grass, road, curb, man, tree, sign, sky, building	unknown	road, sign, man, curb, pole, person
		Which direction can I drive in?	Only straight	road markings and kerb	right		right	

455 It can be seen from the ‘Camera’ column that we tried to
 456 keep diverse driving scenarios in mind while designing the
 457 case study. The answers received from LXMERT, both pre-
 458 trained and when the filter has been integrated, have been
 459 listed along with the features extracted in each case (to the
 460 right of the corresponding column). It provides a visual rep-
 461 resentation of the model’s performance in addressing the
 462 posed questions, allowing for an assessment of their effec-
 463 tiveness based on the Human Answers. The reason for com-
 464 paring the outputs of three VQA models with human an-
 465 swers, using colour coding (green for correct, red for wrong,
 466 yellow for partially correct), is to visually emphasize per-
 467 formance and discrepancies between the models and human
 468 responses. This visual representation allows for a quick and
 469 intuitive understanding of the accuracy and effectiveness of
 470 the models in comparison to human performance. Further

discussion of this rationale is considered in the paper [10] 471
 and the results in the table are discussed in 5. 472

5. Results and Discussion 473

We use the subjective scoring framework for VQA mod- 474
 els [1] in autonomous driving to gauge the improvement 475
 of LXMERT after the filter has been added. This scoring 476
 system analyses the answers provided by the VQA model 477
 using multiple types of natural language processing mod- 478
 els (BERT-base-uncased, NLI-distilBERT-base, all-mpnet- 479
 base-v2 and GPT-2) [4] and sentence similarity benchmark 480
 metrics (Cosine Similarity) [6]. The results are shown in 481
 the Table 2. 482

It can be observed from Figure 5 that there is a notice- 483
 able enhancement in the model’s performance after the in- 484
 tegration of the filter as the MAE (Figure 5a) and RMSE 485

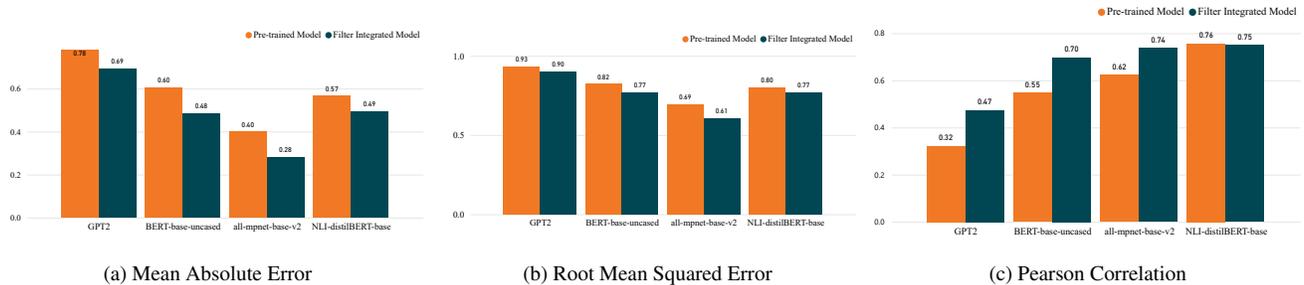


Figure 5. Assessment of LXMERT using the Subjective Scoring Framework

Table 2. Evaluation of LXMERT Performance using Subjective Scoring Framework Metrics: MAE, RMSE, and Pearson Correlation

LXMERT model	Mean Absolute Error		Root Mean Squared Error		Pearson Correlation	
	Pretrained	Filtered	Pretrained	Filtered	Pretrained	Filtered
NLI-distilBERT-base	0.5660	0.4942	0.7998	0.7712	0.7558	0.7504
all-mpnet-base-v2	0.3989	0.2802	0.6932	0.6068	0.6231	0.7370
BERT-base-uncased	0.6042	0.4840	0.8229	0.7675	0.5480	0.6968
GPT2	0.7778	0.6931	0.9344	0.9010	0.3220	0.4737

(Figure 5b) scores have lowered when compared to the pre-trained model. The increase in Pearson correlation (Figure 5c) scores shows that the answers given by the filter-integrated model are closer to the human answers which is ultimately the goal for any VQA model. However, it has to be acknowledged that there are erroneous responses despite this enhancement. These inaccuracies are due to the inherent limitation of the VQA model, as it was not originally designed or trained specifically for driving-related queries.

To address this discrepancy, fine-tuning the model with a driving dataset is a viable solution. This process of fine-tuning will equip the model with the necessary contextual knowledge to interpret questions from a driving perspective accurately, consequently refining its responses accordingly.

After integrating the filter, it's evident from the observed features (Figure 1) that the model has begun to emulate human attention patterns to a remarkable extent. This enhancement is significant for its ability to focus on relevant information. By aligning more closely with human attention patterns, the model becomes more adept at understanding nuanced context, discerning subtle cues, and prioritizing relevant data points. This heightened cognitive alignment improves the model's interpretability and enhances its adaptability in driving scenarios. We further show examples of a few cases in the Supplementary Material where we observe in the figures the differences in object detection and the model's answers due to different filter weights at the Feature extraction stage.

6. Conclusion and Future Work

In conclusion, this study has introduced a novel filter designed to enhance the performance of VQA models specifically in driving-related tasks. Through our case study, we have demonstrated the efficiency of the filter in mimicking

human attention patterns to a significant extent, thereby laying the groundwork for improved VQA capabilities. The limitation of this approach is that we assume that the human is telling what they are actually observing which is leaving a scope for subjectivity in data. For future experiments, we would like to use the eye tribe tracker that delivers real-time data of where a person is looking on a screen similar to [8]. This would potentially improve the accuracy and reliability of the observations in future experiments. However, it's essential to acknowledge that VQA models are not inherently trained for driving tasks, highlighting the need for further optimization and adaptation. Our future work will focus on fine-tuning at least three VQA models using an exclusive driving dataset such as Nuscenes MQA [3], tailored to the complexities of driving environments. By training VQA models on annotated driving scenes and questions, we aim to bolster their performance and adaptability in addressing driving-related queries. Additionally, we plan to conduct a thorough analysis of the fine-tuned models' performances to gain insights into the effectiveness of model adaptation. We also intend to explore the integration of a layer capable of understanding camera information into VQA models. This enhancement will enable the models to perform spatial reasoning tasks more effectively, analyze object positioning within the camera frame, and provide dynamic and adaptive responses to queries about the driving environment. By configuring the filter so that it is capable of leveraging camera information, we aim to bridge the gap between human and machine attention patterns, thereby advancing the capabilities of VQA models in driving scenarios.

References

- [1] The authors. Subjective scoring framework for vqa models in autonomous driving. *IEEE access (In Review)*, 2024. 7

- 552 [2] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE inter-*
553 *national conference on computer vision*, pages 1440–1448,
554 2015. 6
- 555 [3] Yuichi Inoue, Yuki Yada, Kotaro Tanahashi, and Yu Yam-
- 556 aguchi. Nuscenes-mqa: Integrated evaluation of captions
557 and qa for autonomous driving datasets using markup anno-
- 558 tations. In *Proceedings of the IEEE/CVF Winter Conference*
559 *on Applications of Computer Vision*, pages 930–938, 2024.
560 8
- 561 [4] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina
- 562 Toutanova. Bert: Pre-training of deep bidirectional trans-
- 563 formers for language understanding. In *Proceedings of*
564 *naacL-HLT*, page 2, 2019. 7
- 565 [5] Corentin Kervadec, Theo Jaunet, Grigory Antipov, Moez
- 566 Baccouche, Romain Vuillemot, and Christian Wolf. How
- 567 transferable are reasoning patterns in vqa? In *Proceedings*
568 *of the IEEE/CVF Conference on Computer Vision and Pat-*
569 *tern Recognition*, pages 4207–4216, 2021. 2
- 570 [6] Alfirna Rizqi Lahitani, Adhistya Erna Permanasari, and
- 571 Noor Akhmad Setiawan. Cosine similarity to determine sim-
- 572 ilarity measure: Study case in online essay assessment. In
- 573 *2016 4th International Conference on Cyber and IT Service*
574 *Management*, pages 1–6. IEEE, 2016. 7
- 575 [7] Arthur Lam, J Lim, Ricky Sutopo, and Vishnu Monn
- 576 Baskaran. Paying attention to varying receptive fields: ob-
- 577 ject detection with atrous filters and vision transformers. In
- 578 *Proceedings of the 32nd British Machine Vision Conference*,
579 2021. 2
- 580 [8] Kristien Ooms, Lien Dupont, Lieselot Lapon, and Stanislav
- 581 Popelka. Accuracy and precision of fixation locations
- 582 recorded with the low-cost eye tribe tracker in different ex-
- 583 perimental setups. *Journal of eye movement research*, 8(1),
584 2015. 8
- 585 [9] Michael I Posner and Steven E Petersen. The attention sys-
- 586 tem of the human brain. *Annual review of neuroscience*, 13
- 587 (1):25–42, 1990. 3
- 588 [10] Kaavya Rekanar, Ciaran Eising, Ganesh Sistu, and Martin
- 589 Hayes. Towards a performance analysis on pre-trained vi-
- 590 sual question answering models for autonomous driving. In
- 591 *Proceedings of the Irish Machine Vision and Image Process-*
592 *ing Conference*, 2023. 7
- 593 [11] Hao Tan and Mohit Bansal. Lxmert: Learning cross-
- 594 modality encoder representations from transformers. In *Pro-*
595 *ceedings of the 2019 Conference on Empirical Methods in*
596 *Natural Language Processing*, 2019. 5
- 597 [12] Yehui Tang, Kai Han, Yunhe Wang, Chang Xu, Jianyuan
- 598 Guo, Chao Xu, and Dacheng Tao. Patch slimming for ef-
- 599 ficient vision transformers. In *Proceedings of the IEEE/CVF*
600 *Conference on Computer Vision and Pattern Recognition*,
601 pages 12165–12174, 2022. 2
- 602 [13] Geoffrey Underwood, Peter Chapman, Neil Brocklehurst,
- 603 Jean Underwood, and David Crundall. Visual attention while
- 604 driving: sequences of eye fixations made by experienced and
- 605 novice drivers. *Ergonomics*, 46(6):629–646, 2003. 3
- 606 [14] Yang Zhou, Pengfei Cao, Yubo Chen, Kang Liu, and Jun
- 607 Zhao. Prompting vision language model with knowledge
- 608 from large language model for knowledge-based vqa. *arXiv*
609 *preprint arXiv:2308.15851*, 2023. 2