# **Co-VQA : Answering by Interactive Sub Question Sequence**

#### Anonymous ACL submission

#### Abstract

001 Most existing approaches to Visual Question Answering (VQA) answer questions directly, however, people usually decompose a complex 004 question into a sequence of simple sub questions and finally obtain the answer to the origi-006 nal question after answering the sub question sequence(SQS). By simulating the process, this 007 800 paper proposes a conversation-based VQA (Co-VQA) framework, which consists of three components: Questioner, Oracle, and Answerer. 011 Questioner raises the sub questions using an extending HRED model, and Oracle answers 012 them one-by-one. An Adaptive Chain Visual Reasoning Model (ACVRM) for Answerer is also proposed, where the question-answer pair is used to update the visual representation sequentially. To perform supervised learning 017 for each model, we introduce a well-designed method to build a SQS for each question on 019 VOA 2.0 and VOA-CP v2 datasets. Experimental results show that our method achieves state-of-the-art on VQA-CP v2. Further analyses show that SQSs help build direct seman-023 tic connections between questions and images, provide question-adaptive variable-length reasoning chains, and with explicit interpretability 027 as well as error traceability.

### 1 Introduction

028

041

Visual Question Answering (Agrawal et al., 2015) requires to answer questions about images. It has to process visual and language information simultaneously, which is a basic ability of advanced agents. Therefore, related researches (Anderson et al., 2018; Lu et al., 2016; Goyal et al., 2017b; Agrawal et al., 2018) have attracted more and more attention. The conventional approach (Agrawal et al., 2015) for Visual Question Answering (VQA) is to encode image and question separately and incorporate representation of each modality into a joint representation. Recently, with the proposal of Transformer (Vaswani et al., 2017), based on previous dense co-attention models (Kim et al., 2018;



Figure 1: An illustrative example. After a sequence of four sub questions and their answers  $\{(q1,a1),(q2,a2),(q3,a3),(q4,a4)\}$ , its easier to answer the original question.

Nguyen and Okatani, 2018), some methods (Yu et al., 2019; Gao et al., 2019) further adopt selfattention mechanism to exploit the fine-grained information in both visual and textual modalities. Meanwhile, to enrich indicative information about the image contained in the visual representation, some studies (Cadène et al., 2019; Li et al., 2019) have explored different methods of relational reasoning to capture the relationship between objects.

044

045

047

049

051

058

060

061

062

063

064

065

066

067

068

069

Though above methods have achieved significantly improved performance on real datasets (Agrawal et al., 2015; Goyal et al., 2017b), there are still some issues unsolvable. Most existing approaches answer questions directly, however, it is often difficult, especially to answer complex questions. On the one hand, achieving holistic scene understanding in one round is pretty challenging. On the other hand, performing the whole Q&A process in one round lacks interpretability and is absent to locate errors when the model runs into wrong answers. To address the above difficulties, motivated by theory of mind (Leslie, 1987), as shown in Figure 1, we imagine an internal conversation for answering the original question, where a sub question sequence (SQS, which includes several simple sub questions, we use **SO** to refer to sub question later) is raised and answered one-by-one progressively. Finally, the answer to the original

question is obtained by capturing joint informa-071 tion accumulated in the whole SQS. This way has 072 several significant cognitive advantages: 1) questions with different complexity will decompose SQSs with different lengths, resulting in questionadaptive variable-length reasoning chains, 2) generating SQS gives a clear reasoning path, it therefore 077 provides explicit interpretability and traceability of errors, 3) different questions are likely to contain the same SQ or SQS, these common SQ even SQS help improve the generalization ability of models, 4) SQs are usually more simple and directly related to images, which help further strengthen the semantic connection between language and image. 084

085

086

089

091

094

095

101

102

103

104

105

106

108

109

110

111

113

114

115

116

117

118

119

120

121

To achieve above advantages, we therefore propose a **conversation-based VQA (Co-VQA)** framework which includes an internal conversation for VQA. It consists of three components: **Questioner**, **Oracle** and **Answerer**. As shown in Figure 1, once a question is raised, Questioner asks some SQs, and Oracle provides answers one-byone. Their conversation brings a SQS and the corresponding answer sequence. When there is no more SQ to be generated, the internal conversation is finished. Answerer gives the final answer to the original question.

Questioner employs hierarchical recurrent encoder-decoder architecture (Sordoni et al., 2015), and we adopt a representative VQA model (Anderson et al., 2018) as Oracle. For Answerer, we propose an Adaptive Chain Visual Reasoning Model (ACVRM) to accomplish an explicit progressive reasoning process based on SQS, where SQs are used to guide the update of visual features by a graph attention network (Velickovic et al., 2018) one-by-one. Meanwhile, the answers of SQs are utilized as additional supervision signals to guide the learning process. Further, to provide supervision information for the above three models during training, we propose a well-designed method to construct SQS for each question which is based on linguistic rules and natural language processing technology. VQA-SQS and VQA-CP-SQS datasets are obtained after applying this method to VQA 2.0 (Goyal et al., 2017b) and VQA-CP v2 (Agrawal et al., 2018) datasets.

In principle, our work is different from existing VQA systems. Our contributions can be concluded into three-fold:

• We introduce a conversation-based VQA (Co-VQA) framework, which consists of three components: Questioner, Oracle and Answerer. 122

123

124

125

126

127

128

129

130

131

132

133

134

135

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

170

- An Adaptive Chain Visual Reasoning Model (ACVRM) for Answerer is proposed, where the question-answer pair is used to update visual representation sequentially.
- Co-VQA achieves state-of-the-art on the challanging VQA-CP v2 dataset. Moreover, SQSs help build direct semantic connections between questions and images, provide questionadaptive variable-length reasoning chains, and with explicit interpretability as well as error traceability.

## 2 Related Work

Visual Question Answering. The current dominant framework for VQA consists of an image encoder, a question encoder, multimodal fusion, and an answer predictor (Agrawal et al., 2015). To avoid the noises caused by global features, methods(Yang et al., 2016; Malinowski et al., 2018) introduce various image attention mechanisms into VQA. Instead of directly using visual features from CNN-based feature extractors, to improve the performance of model, BUTD(Anderson et al., 2018) adopts Faster R-CNN (Ren et al., 2015) to obtain candidate regional features while Pythia(Jiang et al., 2018) integrates the regional feature with grid-level features. Meanwhile, Lu et al. (2016); Nam et al. (2017) put more attention on learning better question representations. To merge information from different modalities sufficiently, MFB(Yu et al., 2017), and MUTAN(Ben-younes et al., 2017) explored higher-order fusion methods. Further, BAN(Kim et al., 2018) and DCN(Nguyen and Okatani, 2018) propose dense co-attention model which directly establish interaction between different modalities with word-level and regional features. Moreover, with the proposal of Transformer (Vaswani et al., 2017), MCAN (Yu et al., 2019) and DFAF (Gao et al., 2019) adopt self-attention mechanism to fully excavate the fine-grained information contained in text and image. Meanwhile, to fully cover the holistic scene in an image, MuREL (Cadène et al., 2019) and ReGAT (Li et al., 2019) explicitly incorporate relations between regions into the interaction process. Selvaraju et al. (2020) also introduces sub questions into their work, the distinctness between us is that it constructs a Sub-VQA dataset for the subset of reasoning questions



Figure 2: Overall illustration and data flow structure diagram of Co-VQA framework.

in VQA dataset to evaluate consistency of VQA models while we adopt sub questions to achieve task-dividing.

Visual Dialogue. Different from VQA, Visual dia-174 logue(VD) is a continuous conversation for images. 175 Several VD tasks (Visual Dialog (Das et al., 2017), 176 GuessWhich (Chattopadhyay et al., 2017), Guess-What?! (de Vries et al., 2017), MMD (Saha et al., 178 2018)) have been proposed. GuessWhat?!, as a 179 goal-directed dialogue task, requires both players to continuously clarify the reference object through 181 dialogue. The Oracle provides the Questioner with relevant information about the target object by con-183 stantly answering yes/no questions raised by the Questioner, and the Guesser generates the final answer based on the historical dialogue. Following 187 the setting, our Co-VQA framework consists of three components, in which Questioner raises SQs, 188 and Oracle answers them one-by-one, finally, An-189 swerer obtains the answer to the original question.

#### 3 Approach

171

172

173

193

195

196

197

200

201

Figure 2 shows the overall illustration and data 192 flow structure diagram of Co-VQA, which consists of the Questioner, the Oracle, and the An-194 swerer. Given an input image I and a question Q, Co-VQA is responsible for predicting the correct answer from the candidate answer set. Specifically, Questioner is responsible for generating 198 a new SQ  $q_t$  for the next round by combining the information in Q, I and the dialogue history  $H_{t-1} = \{(q_1, a_1), \cdots, (q_{t-1}, a_{t-1})\}$ . Then, Oracle produces appropriate answer  $a_t$  for  $q_t$ . After accomplishing the last round of sub questionanswer pair(abbreviated as qa), Answerer utilizes the historical information accumulated throughout the process to obtain the final answer. In this sec-206 tion, we will introduce the three components in Section 3.1-3.3. 208



Figure 3: Overview of the **Questioner** model which is based on extending HRED model. There are three modules: Image Encoder, Hierarchical Encoder, Decoder.

209

210

211

212

213

214

215

216

217

218

219

221

222

223

224

225

227

228

229

230

231

232

234

235

237

#### 3.1 Questioner

At round t, given an image I, a question Q and the dialogue history  $H_{t-1}$ \_  $\{(q_1, a_1), \cdots, (q_{t-1}, a_{t-1})\}$ , Questioner is responsible for generating new sub question  $q_t$ . Generally, we build Questioner based on extending hierarchical recurrent encoder decoder (HRED) (Sordoni et al., 2015). The overall structure of Questioner is depicted in Figure 3.

Image Encoder. Following common practice(Anderson et al., 2018), we extract regional visual features from I in a bottom-up manner using Faster R-CNN model(Ren et al., 2015). Each image will be encoded as a series of M regional visual features  $R \in \mathbb{R}^{M \times 2048}$  with their bounding box  $b = [x, y, w, h] \in \mathbb{R}^{M \times 4}$  ( $M \in [10, 100]$  in our experiments).

Hierarchical Encoder. Embedding matrix Embedder is adopted to map Q and each pair  $(q_i, a_i)$ in  $H_{t-1}$  to  $Q^{emb}$  and  $(q_i^{emb}, a_i^{emb})$  respectively. Then, two question-level encoder GRU,  $GRU_Q$ and  $GRU_q$ , are deployed to obtain corresponding question feature  $Q^{fea}$  and  $q_i^{fea}$  for Q and  $q_i$ .

 $Q^{fea}$  is utilized as the first step input of sessionlevel encoder GRU, GRUs to grasp global information of original question.  $q_i^{fea}$  and  $a_i^{emb}$  are concatenated as  $qa_i^{fea}$ , which is regard as representation for sub question-answer pair. Meanwhile, it is treated as the i+1-th step input of  $GRU_s$  to



Figure 4: Model architecture of the proposed ACVRM for Answerer. There are four functional modules: Image Encoder, Ouestion Encoder, Sequential Progressive Reasoning and Multimodal Fusion.

(1)

(2)

238

252

256

257

260

261

263

264

#### 3.3 Answerer

3.2

Oracle

method, as our Oracle.

Given a question Q, an image I and a complete dialogue history  $H_T = \{q_1, a_1, ..., q_T, a_T\}$ , the 262 assignment of Answerer is to find out the most accurate  $\hat{a}$  in the candidate answer set, which could be denoted as:

$$\hat{a} = \underset{a \in A}{\operatorname{argmax}} P_{\theta}(a|I,Q,H_T), \qquad (3)$$

obtain context feature  $s_{i+1}$ , which is denoted as:

 $s_{i+1} = GRU_s([q_i^{fea} \mid\mid a_i^{emb}], s_i),$ 

where || represents concatenation. After encoding

 $H_{t-1}$ , we obtain current context representation  $s_t$ .

**Decoder.** At decoding  $q_t$ , we employ an extra one-

layer GRU as decoder, which is initialized by  $s_t$ .

Then a question-guided attention is deployed to

regional features R to obtain the weighted visual

feature  $v_t$ . Further, we fuse  $v_t$  with  $Embedder(q_t^i)$ 

ing, where T is the maximum round of dialogues,

 $L(\theta_Q) = -\sum_{t=1}^{T} log P(q_t | Q, I, H_{t-1}).$ 

The Oracle is responsible for constantly answering

SQs raised by Questioner. Specifically, at round

t, Oracle supplies the answer  $a_t$  for SQ  $q_t$ , based on the image I and  $q_t$ . We regard Oracle as a con-

ventional VQA task and adopt the BUTD (Ander-

son et al., 2018), which is a representative VQA

The negative log-likelihood loss is used for train-

as the input of decoder at every time step *i*.

 $\theta_Q$  is the parameters of Questioner :

where  $\theta$  denotes the parameters of Answerer. To accomplish this task, we propose an Adaptive Chain Visual Reasoning Model (ACVRM), which consists of four components: Image Encoder, Question Encoder, Sequential Progressive Reasoning, and Multimodal Fusion. The overall structure of ACVRM is illustrated as Figure 4.

267

268

270

271

274

275

276

277

278

279

280

282

283

285

286

287

290

291

292

293

294

295

297

#### Image and Question Encoder 3.3.1

Feature extraction modules are shown in the left part of Figure 4. Image encoder is the same as Questioner. For question encoder, we adopt a bidirectional Transformer (Vaswani et al., 2017). Q and each SQ in  $H_T$  will be padded to a maximum length of 14 and be encoded by bidirectional Transformer with random initialization, at last the corresponding question features  $E \in \mathbb{R}^{d_q}, \{SE_i\}_{i=1}^T \in$  $\mathbb{R}^{T \times d_q}$  are obtained after mean pooling. To align the feature dimensions, we linearly map image feature R to  $V_0 \in \mathbb{R}^{M \times d_v}$ . We set  $d_q = d_v = 768$ .

# 3.3.2 Sequential Progressive Reasoning (SPR) **Overall.** To realize progressive visual reasoning under the guidance of SQS, we utilize Graph Visual Reasoning (GVR) module, which will be introduced later, to gradually guide the update of visual features. Specifically, for Q containing T SQs, the t-th step of SPR can be expressed as:

$$V_t^R = GVR(V_{t-1}, SE_t; \theta_G), \qquad (4)$$

where  $V_t^R$  represents the t-th step visual feature, and  $\theta_G$  denotes parameters for GVR. Then, residual connection is deployed in each round to preserve historical information and avoid vanishing of gra-



Figure 5: Flowchart of the **GVR**, including two parts: multimodal fusion based on concatenation and relation reasonsing based graph attention network.

dient. Therefore, the updated visual feature for t-th round can further be depicted as  $V_t = V_{t-1} + V_t^R$ .

Furthermore, each  $q_t$  has a corresponding answer  $a_t$ , which supplies an additional supervision signal for training. For each step t, we adopt a shared two-layer MLP as the sub classifier and then utilize average  $V_t^R$  as input. A cross-entropy loss is used for classification, which is denoted as  $Loss_t^{sub}$ .

301

303

305

306

310

311

312

315

317

318

319

321

327

329

331

**Graph Visual Reasoning.** Inspired by ReGAT (Li et al., 2019), we utilize Graph Attention Network (Velickovic et al., 2018) to learn relations between objects. An overall illustration of GVR is shown in Figure 5. The whole reasoning process is abbreviated as  $V^R = GVR(V,q)$ , which consists of two parts: feature fusion and relational reasoning.

At first, the question representation q is concatenated with each of the M visual features  $v_i$ , which we write as  $[v_i || q]$ , then we compute a joint embedding as:

$$v'_i = W([v_i || q]) \quad for \ i = 1, ..., M,$$
 (5)

where  $W \in \mathbb{R}^{d_q \times (d_q + d_v)}$ , and  $v'_i \in \mathbb{R}^{d_q}$  is conducted as initial value of node in the graph G(V, E), where  $e_{ij}$  denotes edges between nodes. Then, to reduce the interference caused by irrelevant information, we design a masked multi-head attention for relational reasoning. Specially, for each head, inspired by Hu et al. (2018), attention weight not only depends on visual-feature weight  $\alpha_{ij}^{h,v}$ , but also bounding-box weight  $\alpha_{ij}^{h,b}$ , we formulate nonnormalized attention weight  $e_{ij}$  as:

$$e_{ij}^h = \alpha_{ij}^{h,v} + \log(\alpha_{ij}^{h,b}), \tag{6}$$

$$\alpha_{ij}^{h,v} = \frac{(W_q^h v_i')^T \cdot W_k^h v_j'}{\sqrt{d_h}},$$
(7)

$$\alpha_{ij}^{h,b} = \max\{0, w \cdot f_b(b_i, b_j)\},$$
(8)

where 
$$d_h = \frac{d_q}{H}$$
,  $H$  denotes the number of head and  
we set  $H = 8$ ,  $W_q^h \in \mathbb{R}^{d_h \times d_q}$ ,  $W_k^h \in \mathbb{R}^{d_h \times d_q}$ 

 $f_b(\cdot, \cdot)$  first computes relative geometry feature  $(log(\frac{|x_i-x_j|}{w_i}), log(\frac{|y_i-y_j|}{h_i}), log(\frac{w_j}{w_i}), log(\frac{h_j}{h_i}))$ , then embeds it into a  $d_h$ -dimensional feature by computing cosine and sine functions of different wavelengths,  $w \in \mathbb{R}^{d_h}$ . Furthermore, according to  $e_{ij}^h$ , to learn a sparse neighbourhood  $N_i^h$ for each node i, we adopt a ranking strategy as  $N_i^h = top_K(e_{ij}^h)$ , where  $top_K$  returns the indices of the K largest values of an input vector, and we set K=15.

334

335

336

337

338

339

340

341

342

343

344

351

352

353

354

356

357

358

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

By employing above mechanism, output features of each head are concatenated, denoted as:

$$v_i^R = ||_{h=1}^H \sigma(\sum_{j \in N_i^h} softmax(e_{ij}^h) \cdot W_v^h v_j'), \quad (9)$$

where  $W_v^h \in \mathbb{R}^{d_h \times d_q}$ .

#### 3.3.3 Fusion Module

Context-aware visual features  $V_T$  are obtained after completing the whole process of SPR. To sufficiently integrate the information of two modalities, we utilize Q to convert  $V_T$  into  $\tilde{V}$  through GVR as  $\tilde{V} = GVR(V_T, E)$ . Then, we employ the same multi-modal fusion strategy as Anderson et al. (2018) to obtain a joint representation H. For Answer Predictor, we adopt a two-layer multi-layer perceptron (MLP) as classifier, with H as the input. Binary cross entropy is used as the loss function. Thus, final loss can be formulated as:

$$Loss = Loss_{BCE} + \sum_{t=1}^{T} Loss_t^{sub}.$$
 (10)

# 4 Experiments

#### 4.1 Datasets

We evaluate our approach on two widely used datasets, including: 1) VQA 2.0 (Goyal et al., 2017b) is composed of real images from MSCOCO (Lin et al., 2014) with the same train/validation/test splits. For each image, an average of 3 questions are generated. These questions are divided into 3 categories: Y/N, Number and Other. 10 answers are collected for each image-question pair from human annotators. The model is trained on the training set, but when testing on the test set, both training and validation set are used for training, and the max-probable answer is selected as the predicted answer. 2) VQA-CP v2 (Agrawal et al., 2018) is a derivation of VQA 2.0. In particular, the distribution of answers with respect to question types differs between training and test splits.

	Test-std				
Model	All	Y/N	Num	Other	All
Bottom-Up	63.37	80.4	43.02	55.96	65.67
BAN	66.04	-	-	-	-
MuREL	65.14	-	-	-	68.41
ReGAT*	67.18	-	-	-	70.58†
DFAF	66.66	-	-	-	70.34†
MCAN	67.2	84.85‡	49.5‡	58.45‡	70.9†
MLIN	66.53	-	-	-	70.28†
Ours	67.26	84.71	50.38	58.44	70.39

Table 1: Performance on VQA 2.0 val split and teststandard splits. "\*" means ensembling result. "†" means training with augmented VQA samples from Visual Genome. "‡" based on our re-implementations.

**Construction of SQS dataset.** To provide the corresponding supervised signal for training Questioner, Oracle, and Answerer, we propose a well-designed method, which is chiefly based on linguistic rules and natural language processing technology. **VQA-SQS** and **VQA-CP-SQS** are obtained by applying this method on VQA 2.0 and VQA-CP v2 datasets. The details of the construction process and the specific statistical information of the two datasets can be found in Appendix.

#### 4.2 Implementation Details

389

391

397

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

**Training and inference.** During training, Questioner, Oracle and Answerer are trained independently. For inference, given a question Q and an image I, SQS is firstly generated through the cooperation between Questioner and Oracle, then Q, I and the complete SQS is combined as the input of Answerer, and obtain the final answer.

Each question is tokenized and padded with 0 to a maximum length of 14. For Questioner and Oracle, each word is embedded using 300-dimensional word embeddings. The dimension of the hidden layer in GRU is set as 1,024(except for  $GRU_Q$  and  $GRU_s$  with 1,324).

Our model is implemented based on Py-Torch(Paszke et al., 2017). In experiments, we use Adamax optimizer for training, with the minibatch size as 256. For choice of learning rate, we employ the warm-up strategy(Goyal et al., 2017a). Specifically, we begin with a learning rate of 5e-4, linearly increasing it at each epoch till it reaches 2e-3 at epoch 4. After 14 epochs, the learning rate is decreased by 0.2 for every 2 epochs up to 18 epochs. We also adopt early stopping strategy. For transformer encoder, we fix the learning rate as 5e-5. Every linear mapping is regularized by weight

Model	All	Y/N	Num	Other
MuREL	39.54	42.85	13.17	45.04
ReGAT*	40.42	-	-	-
MCAN‡	42.35	42.29	14.51	50.02
Ours	42.52	44.42	14.68	49.17

Table 2: State-of-the-art comparison on the VQA-CP v2 dataset. "\*" means ensembling result. "‡" Results based on our re-implementations.

normalization and dropout (p = 0.2 except for the classifier with 0.5).

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

### 4.3 Results

To compare with existing VQA methods, we conduct several experiments to evaluate the performance of our Co-VQA framework, further, to verify the generation quality of the SQs and their impact on the performance of the overall model, Questioner and Oracle are tested additionally. In Table 1, we compare our method with previous work on VQA 2.0 validation and test-standard split. From Table 1, it can be seen that on validation split, Co-VQA achieves the top-tier preformance, our method obtains a accuracy of 67.26, which surpass that of (Yu et al., 2019) by 0.06, and achieves a obvious performance improvement on number questions. On test-standard split, without additional augmented samples from Visual Genome, Our performance still at the third place. We assume the gap between two splits mainly due to the difference in SQS generation quality. To demonstrate the generalizability of Co-VQA, we also conduct experiments on the VQA-CP v2 dataset, where the distributions of the training and test splits are different. Table 2 illustrates the overall performance and from the experimental results, our model gains a significant advantage (2.1) over ReGAT. Compared with MCAN, our model also improved by 0.16. For Questioner and Oracle, we train and evaluate on the train/val split of VQA-SQS dataset.

**Oracle.** The accuracy of Oracle is 93.73 and average F-value is 90.13. On the one hand, the high accuracy is due to SQ itself being simple; On the other hand, decomposition of question leads to many same SQs, strengthening image-language correlation ability at SQ level.

**Questioner.** For Questioner, BLEU score is adopted to measure the quality of the generated SQs. As is shown in Table 3, we attribute the low

<b>BLEU-1</b>	BLEU-2	BLEU-3
26.5	10.4	4.78

Table 3: BLEU evaluation scores of Questioner. We don't report BLEU-4 score because the length of some sub questions is shorter than 4.

Model	All	Y/N	Num	Other
Full	67.26	84.71	50.38	58.44
wo-sub-loss	66.94	84.58	48.95	58.29
wo-SQS	66.55	84.43	46.78	58.18

Table 4: Ablation studies on impact of SQS on VQA2.0 validation set.

BLEU scores to the diversity of syntax details.

#### 4.4 Ablation Study

454

455

456

457

458

459

460

461

462

463

464

465

466

467

We conduct several ablation studies to explore critical factors affecting the performance of Co-VQA.

**The impact of SQS.** In general, as we can observe from Table 4, though there are noises in the answers for SQs, the weak supervision signal provided by them shows a gain of +0.32. Furthermore, the decrease is obvious(-0.71) when we remove total SQS from model, indicating that though the SQS generated from Questioner is not good enough, it still plays a important role in improving performance of model.

Detail Analysis of SQS. To analyze the impact 468 of SQS in detail, we divide the validation split of 469 VQA-SQS into SQS-0 / SQS-1 / SQS-2 / SQS-3&4 470 subsets, where SQS-n represents samples with n 471 SQs. Then, the average accuracy of different mod-472 els on each subset is reported in Table 5. For SQS-1 473 and SQS-2, the additional reasoning brought by 474 SQS achieves a improvement of 1.02 and 0.93 475 respectively. However, for SQS-3&4, the perfor-476 mance decreases compared with wo-SQS, we per-477 form statistics in two aspects to comprehensively 478 explore the causes of this phenomenon. As shown 479 in Table 6, compared with other subsets, SQS-3&4 480 has obviously fewer samples, causing insufficient 481 learning for these samples of long sequence. More-482 483 over, SQs in SQS-34 occur less frequently, thus it is inadequate for model to establish accurate semantic 484 connections between these images and questions. 485

486 **Coherence of SQS.** We also studied the impact 487 of the coherence of SQS on the performance. We

Model	SQS-0	SQS-1	SQS-2	SQS-3&4	All
	(57,411)	(119,285)	(34,226)	(3,432)	(214,254)
Full	69.51	66.70	65.78	63.62	67.26
wo-SQS	69.48	65.68	64.85	64.35	66.55

Table 5: Ablation studies of SQS in detail on VQA 2.0 validation set. SQS-n represents the subset of samples with n SQs in VQA-SQS validation set. We report the average accuracy on each subset.

Subset	SQS-0	SQS-1	SQS-2	SQS-3&4	All
Samples-Num	57,411	119,285	34,226	3,432	214,254
Avg(Freq of SQ)	-	870	851	693	854

Table 6: Data statistics of SQS in detail on VQA 2.0 validation set. The first row shows the number of original questions contained in different SQS sets, and the second row counts the average number of occurrences of the sub questions contained in each subset in the VQA-SQS train split.

ran two different cases : 1) randomly shuffle the SQs in a sequence; 2) remove some SQs in a sequence with 50% probability. As we can observe from Table 7, the declines from original one are not significant, partly due to the coherence of SQS in dataset is not good enough.

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

#### 4.5 Visualization

To better illustrate the effectiveness, explicit interpretability, and traceability of errors of Co-VQA, we visualize and compare the attention maps learned by complete Co-VQA with those learned by wo-SQS. As shown in Figure 6. Column 1 is the original question and ground truth, while Column 2 corresponds to the prediction of model wo-SQS. The middle columns and last column correspond to the generated sub q&a, and the prediction of Co-VQA, respectively. To visualize the attention maps, we use the in-degree of each node as attention value and circle top-2 attended regions with red and blue boxes.

Line 1 shows model wo-SQS only notices one of the dogs and gives a wrong answer "1". However, through SQ "Are there dogs?", Co-VQA focuses on two dogs and gives the correct answer "2". This

Model	All	Y/N	Num	Other
Full	67.26	84.71	50.38	58.44
shuffle	67.15	84.68	49.78	58.42
random	67.08	84.68	49.75	58.28

Table 7: Ablation studies of coherence of SQS on VQA2.0 validation set.



Figure 6: Visualization of attention maps learned by complete Co-VQA with those learned by wo-SQS. The second and last column corresponds to the prediction of wo-SQS and complete Co-VQA respectively. Red and blue bounding boxes shown in each image are the top-2 attended regions.

512 case demonstrates that asking an existence question firstly is beneficial to number question. In Line 2, 513 wo-SQS model focuses on unrelated entities. How-514 ever, Co-VQA attends to the women and the people 515 wearing short sleeves gradually with SQS, and finally, concentrates on the related woman's shirt. Line 3 shows Co-VQA successively attends to cars, 518 blue cars, and the license plate under the guidance 519 of SQS and gets the correct answer. These examples prove that questions with different complexity 521 will correspond to SQS of variable length, and SQ is indeed related to more accurate image attention. 523 Moreover, generating SQ provides not only the 524 logic of reasoning but also additional language interpretation. Thus, compared with previous works 526 that only explain models by attention maps, Co-VQA has significantly better interpretability. 528

> The last line shows Co-VQA gives a wrong answer after adding SQS. However, we can find some possible causes, such as the wrong answer of Q1, Q2 is not related to the question, and the model doesn't attend to relevant entities in the light of

530

531

Q1. It shows that Oracle and Questioner may give wrong answers or generate inappropriate questions, as well as Answerer may establish faulty semantic connections between questions and images, which verifies that Co-VQA has sure traceability for errors and provides guidance for future work.

### 5 Conclusions

We propose a conversation-based VQA (Co-VQA) framework which consists of Questioner, Oracle and Answerer. Through internal conversation based on SQS, our model not only has explicit interpretability and traceability of answer errors, but also can carry out question-adaptive variablelength reasoning chains. Currently, Questioner is relatively simple, and the quality still has a lot of room to improve. Meanwhile, current SQs are only yes/no questions. For future work, we plan to explore how to more effectively generate more diverse and higher quality SQS, and look forward to better model performance.

534

# 540 541

542

543

544

545

546

547

548

549

550

551

552

#### References

554

555

559 560

565

566

571

572

573

574

575

577

579

581

583

584

587

591

592

593

594

595

596

598

- Aishwarya Agrawal, Dhruv Batra, Devi Parikh, and Aniruddha Kembhavi. 2018. Don't just assume; look and answer: Overcoming priors for visual question answering. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 4971– 4980.
- Aishwarya Agrawal, Jiasen Lu, Stanislaw Antol, Margaret Mitchell, C. Lawrence Zitnick, Devi Parikh, and Dhruv Batra. 2015. Vqa: Visual question answering. *International Journal of Computer Vision*, 123:4–31.
- Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang.
  2018. Bottom-up and top-down attention for image captioning and visual question answering. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 6077–6086.
- Hedi Ben-younes, Rémi Cadène, Matthieu Cord, and Nicolas Thome. 2017. Mutan: Multimodal tucker fusion for visual question answering. 2017 IEEE International Conference on Computer Vision (ICCV), pages 2631–2639.
- Rémi Cadène, Hedi Ben-younes, Matthieu Cord, and Nicolas Thome. 2019. Murel: Multimodal relational reasoning for visual question answering. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 1989–1998.
- Prithvijit Chattopadhyay, Deshraj Yadav, Viraj Prabhu, Arjun Chandrasekaran, Abhishek Das, Stefan Lee, Dhruv Batra, and Devi Parikh. 2017. Evaluating visual conversational agents via cooperative humanai games. In *HCOMP*.
- Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José M. F. Moura, Devi Parikh, and Dhruv Batra. 2017. Visual dialog. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 1080–1089.
- Harm de Vries, Florian Strub, A. P. Sarath Chandar, Olivier Pietquin, H. Larochelle, and Aaron C. Courville. 2017. Guesswhat?! visual object discovery through multi-modal dialogue. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 4466–4475.
- Peng Gao, Hongsheng Li, Haoxuan You, Zhengkai Jiang, Pan Lu, Steven C. H. Hoi, and Xiaogang Wang. 2019. Dynamic fusion with intra- and inter-modality attention flow for visual question answering. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 6632–6641.
- Priya Goyal, Piotr Dollár, Ross B. Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. 2017a. Accurate, large minibatch sgd: Training imagenet in 1 hour. ArXiv, abs/1706.02677.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv 609 Batra, and Devi Parikh. 2017b. Making the v in vga 610 matter: Elevating the role of image understanding in 611 visual question answering. 2017 IEEE Conference 612 on Computer Vision and Pattern Recognition (CVPR), 613 pages 6325–6334. 614 Han Hu, Jiayuan Gu, Zheng Zhang, Jifeng Dai, and 615 Yichen Wei. 2018. Relation networks for object de-616 tection. 2018 IEEE/CVF Conference on Computer 617 Vision and Pattern Recognition, pages 3588–3597. 618 Yu Jiang, Vivek Natarajan, Xinlei Chen, Marcus 619 Rohrbach, Dhruv Batra, and Devi Parikh. 2018. 620 Pythia v0.1: the winning entry to the vqa challenge 621 2018. ArXiv, abs/1807.09956. 622 Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. 623 2018. Bilinear attention networks. In NeurIPS. 624 Alan M. Leslie. 1987. Pretense and representation: The 625 origins of "theory of mind.". Psychological Review, 626 94:412-426. 627 Linjie Li, Zhe Gan, Yu Cheng, and Jingjing Liu. 628 2019. Relation-aware graph attention network for 629 visual question answering. 2019 IEEE/CVF Interna-630 tional Conference on Computer Vision (ICCV), pages 631 10312-10321. 632 Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James 633 Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, 634 and C. Lawrence Zitnick. 2014. Microsoft coco: 635 Common objects in context. In ECCV. 636 Edward Loper and Steven Bird. 2002. Nltk: The natural 637 language toolkit. CoRR, cs.CL/0205028. 638 Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. 639 2016. Hierarchical question-image co-attention for 640 visual question answering. In NIPS. 641 Mateusz Malinowski, Carl Doersch, Adam Santoro, and 642 Peter W. Battaglia. 2018. Learning visual question 643 answering by bootstrapping hard attention. In ECCV. 644 Hyeonseob Nam, Jung-Woo Ha, and Jeonghee Kim. 645 2017. Dual attention networks for multimodal rea-646 soning and matching. 2017 IEEE Conference on 647 Computer Vision and Pattern Recognition (CVPR), 648 pages 2156-2164. 649 Duy-Kien Nguyen and Takayuki Okatani. 2018. Im-650 proved fusion of visual and language representations 651 by dense symmetric co-attention for visual question 652 answering. 2018 IEEE/CVF Conference on Com-653 puter Vision and Pattern Recognition, pages 6087-654 6096. 655 Adam Paszke, Sam Gross, Soumith Chintala, Gregory 656 Chanan, Edward Yang, Zach DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 658 2017. Automatic differentiation in pytorch. 659

Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39:1137–1149.

663

664

665

666

668

673

674

675

676

677 678

679

682

684

692

695

698

701

- Amrita Saha, Mitesh M. Khapra, and Karthik Sankaranarayanan. 2018. Towards building large scale multimodal domain-aware conversation systems. In AAAI.
- Ramprasaath R. Selvaraju, Purva Tendulkar, Devi Parikh, Eric Horvitz, Marco Túlio Ribeiro, Besmira Nushi, and Ece Kamar. 2020. Squinting at vqa models: Introspecting vqa models with sub-questions. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 10000–10008.
- Alessandro Sordoni, Yoshua Bengio, Hossein Vahabi, Christina Lioma, Jakob Grue Simonsen, and Jianyun Nie. 2015. A hierarchical recurrent encoder-decoder for generative context-aware query suggestion. Proceedings of the 24th ACM International on Conference on Information and Knowledge Management.
- Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *ArXiv*, abs/1706.03762.
- Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio', and Yoshua Bengio. 2018. Graph attention networks. *ArXiv*, abs/1710.10903.
- Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alex Smola. 2016. Stacked attention networks for image question answering. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 21–29.
- Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. 2019. Deep modular co-attention networks for visual question answering. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 6274–6283.
- Zhou Yu, Jun Yu, Jianping Fan, and Dacheng Tao. 2017. Multi-modal factorized bilinear pooling with co-attention learning for visual question answering. 2017 IEEE International Conference on Computer Vision (ICCV), pages 1839–1848.

703

Α

Appendix

A.1 Data source

VQA-CP v2 datasets.

question sequence.

following process.

A.2 Construction principle

order-1 SQs and order-2 SQs.

A.3 Construction method

Here we introduce our method for construcing SQS

We construct SQS dataset based on VQA 2.0 and

The principle of data construction is based on an

intuitive idea: high-order questions can adopt cor-

responding low-order questions as their sub ques-

tions, then these sub questions are arranged accord-

ing to the order from low to high to form a sub

We determine the order of questions according to

the templates in the Table 8. For order-0 and order-

1 questions, there is no corresponding SQ, order-2

questions can construct corresponding order-1 SQs,

while the order-3 questions can construct multiple

we construct the sub questions according to the

1) For each question, we first adopt Spacy  $^{1}$  and

NLTK toolkit (Loper and Bird, 2002) to identify

all noun blocks in the question and filter out some

noun blocks based on the predefined phrase list.

The phrase list mainly includes meaningless quan-

tifiers, pronouns and abstract nouns, such as lots,

someone, something, you, they, it, the day, the pic-

2) After the first step of filtering, for questions con-

taining noun blocks, according to the dependency

relation between the extracted noun blocks, some

noun blocks may be used as prepositional phrases.

For remaining noun blocks, we use Part-of-Speech

Tagging of Spacy to find out corresponding nouns,

adjectives and quantifiers. For nouns, we save

them separately, while for adjectives, quantifiers

and prepositional phrases, we save these modifiers

with the noun in a form of 2-tuple, such as (noun,

3) After the first step of filtering, for questions

without noun blocks, considering there may be

omissions in the process of extraction, we perform

pattern matching through Spacy based on the pre-

defined matching template to determine their cate-

ture, a body, emotion, this, type etc.

and the statistical information of datasets.

- 712
- 710

711

713 714

715

716

717

719

720 721

722

723

724 725

727 728

730 731

733

737 738

736

740

741 742

744

745

746

747

749

gories. Table 9 illustrates partial matching patterns <sup>1</sup>https://spacy.io/

modifier).

order	question template
0	no entity
1	single entity
2	entity & attribute
3	comparison between different entities

Table 8: Templates for question of different order.

for different type of questions. Specially, for existence questions, no additional processing is required, while for other type of questions, we save the nouns that are exist in the questions.

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

4) We further filter the nouns and tuples saved in 2) and 3). We aim to filter out abstract nouns, nonsubstantial nouns, and 2-tuple corresponding to these nouns. The following are some cases to be filtered:

a) Abstract Noun: direction, design, surface, area, emotion, skill etc.

b) Non Substantive Noun: mode, base, day, love, name, print, piece etc.

5) For the remaining nouns and their corresponding 2-tuple, we use the pre-defined question template to construct the corresponding sub questions. To facilitate the process of construction and model training, we design all sub questions as yes / no questions and reveal the matching pattern for each type of sub questions in Table 11. The construction process of ground-truth for sub questions can be illustrated as following:

Existence SQ and Attribute SQ we first extract the label and attribute information of the entity by using the detection model, and then combine these information to produce the answer.

Prep SQ and Position SQ the location information obtained by the detection model is utilized to judge the relationship of overlapping and orientation between entities, we use the obtained relationship to generate the corresponding answer.

Number SQ we first make a rough quantity estimation based on the image, and then make manual correction.

6) Considering there may be wrong answers, incoherent sequences and nonstandard question grammar in the process of automatic construction, and to increase the diversity of SQs, we invite ten students in our laboratory to further manually correct some samples(about 5K samples).

The SOS datasets obtained by performing the above operations on VQA 2.0 and VQA-CP V2 datasets are called VQA-SQS and VQA-CP-SQS respec-

Question Type	Matching Pattern
Existence	(do you see)?[DET   PRON   ADP]* [NOUN   PROPN]* NOUN?
Verb	(do you see)? [DET   PRON   ADP]* [NOUN   PROPN]* NOUN?[VBG   VBN]?
Attribute	BE [DET   PRON   ADP]* [NOUN   PROPN]* NOUN? ADJ?
Num	BE [DET   PRON   ADP]* NUM NOUN NOUN* ?
Prep	BE [DET   PRON   ADP]* [NOUN   PROPN]* NOUN VERB? ADP DET NOUN NOUN* ?

Table 9: Matching patterns for different type of questions

Dataset	Split	#Images	#Q&A	#Non-empty SQS	Avg(SQ)
VQA-SQS	Train	82,783	443,757	328,140	0.94
VQA-SQS	Val	40,504	214,354	156,943	0.925
VQA-CP-SQS	Train	120,932	438,183	322,200	0.93
VQA-CP-SQS	Test	98,226	219,928	162,883	0.946

Table 10: Dataset statistics of VQA-SQS and VQA-CP-SQS.



Figure 7: Dataset distribution of VQA-SQS.



Figure 8: Dataset distribution of VQA-CP-SQS.



Figure 9: Some samples of VQA-SQS, including existence SQ, attribute SQ, prep SQ and number SQ.

SQ Type	Matching Pattern		
	Is there any [entity]?		
Existence	Is there any [color] [entity]?		
	Are there [entites]?		
	Is the [entity] [color]?		
Attribute	Is any [entity]?		
	Are these [entites] in similar size?		
Drop	Is there any [entity] on the [entity2]?		
Flep	Is there any [entity] in the [entity2]?		
Number	Are there [number] [entites]?		
Nulliber	Is there only one [entity]?		
Position	Is the [entity] on the left?		
	Is the [entity] on the right?		
	' Is the [entity] in the middle?		

Table 11: Sub question generation template for different SQ types.

#### tively.

### A.4 Dataset statistics

Table 10 shows general statistical information of the two SQS datasets, then, Figure 7 and Figure 8 respectively reveal three fine-grained distribution of two datasets including number distribution of SQ (7-a / 8-a), type distribution of SQ (7-b / 8b) and answer distribution of SQ (7-c / 8-c). To display more convenient, in (7-a / 8-a) and (7-b / 8-b), the ordinate axis adopts logarithmic scale.

Figure 9 displays four samples of VQA-SQS dataset.