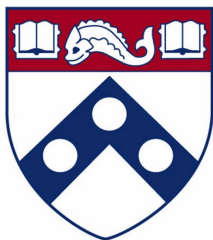


# Image-Grounded Conversations: Multimodal Context for Natural Question and Response Generation

Mostafazadeh et al. (2017)

University of Rochester, Microsoft and University College London



Presenter: Shubham Annadate

# Agenda

- Introduction
- Related Work
- Image Grounded Conversation (IGC)
  - Dataset
  - Task Characteristics
- Models
- Evaluation Setup
- Experimental Results
- Conclusion/Key Contributions
- Thoughts

# Introduction

- Recent work on vision and language → describing or answering questions about image.
- Conversations threads on social media platforms like Twitter.
  - 28% of tweets contain image (as of 2015)
  - Conversation are beyond what is explicitly visible in the image



**User1:** My son is ahead and surprised!

**User2:** Did he end up winning the race?

**User1:** Yes he won, he can't believe it!

Figure 1: A naturally-occurring Image-Grounded Conversation.

# Introduction

- Look at image as a context for interaction rather than an artifact
- Image-Grounded Conversation (IGC)
  - System that generates conversational turns to drive conversation.
  - Falls between chit-chat and goal-oriented dialog systems
  - Combines the threads of language & vision and data-driven conversational modeling



[Goal-oriented dialog systems](#)

# Related Work

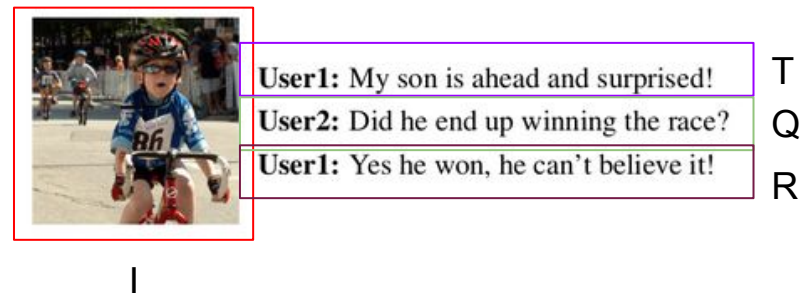
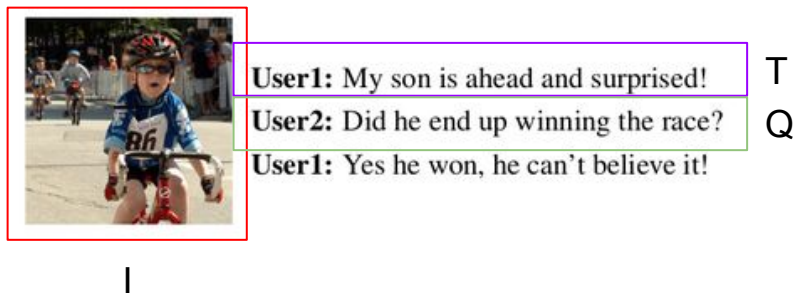
- Vision and Language
  - Visual features combined with language modeling have shown improved performance in image captioning and VQA tasks (2014-15)
  - **Visual Question Generation (VQG) Task** (2016)
- Data-Driven Conversational Modeling
  - Learning conversations from message-response pairs (2011)
  - Context-Sensitive Neural Language Models (2015-16)

# IGC (Task Definition)

Two consecutive conversational steps within the current scope:

Question generation:  $(I, T) \rightarrow Q$

Response Generation:  $(I, T, Q) \rightarrow R$



*I: visual context, T: textual context, Q: question, R: response*

# IGC (Dataset)

- No pre-existing dataset for IGC task
- IGC<sub>Crowd</sub>
  - Sampled *eventful* images from VQG dataset
  - Pair of Amazon MTurk workers have a short conversation about the image
  - For multi-reference evaluation → crowd-sourced 5 additional references per question/response.
  - Used for **validation and testing purpose**
- IGC<sub>Twitter</sub>
  - Used for **training purpose**
  - Sampled 250K quadruples of {I, T, Q, R} tweet threads from Twitter dataset

IGC <sub>Crowd</sub> (val and test sets, split: 40% and 60%)	
# conversations = # images	<b>4,222</b>
total # utterances	<b>25,332</b>
# all workers participated	308
Max # conversations by one worker	20
Average payment per worker (min)	1.8 dollars
Median work time per worker (min)	10.0
IGC <sub>Crowd-multiref</sub> (val and test sets, split: 40% and 60%)	
# additional references per question/response	5
total # multi-reference utterances	<b>42,220</b>

Table 2: Basic Dataset Statistics.

# IGC<sub>Crowd</sub> (Dataset)




<b>Visual Context</b>			
<b>Textual Context Question</b>	This wasn't the way I imagined my day starting. do you think this happened on the highway?	I checked out the protest yesterday. Do you think America can ever overcome its racial divide?	A terrible storm destroyed my house! OH NO, what are you going to do?
<b>Response</b>	Probably not, because I haven't driven anywhere except around town recently.	I can only hope so.	I will go live with my Dad until the insurance company sorts it out.
<b>VQG Question</b>	What caused that tire to go flat?	Where was the protest?	What caused the building to fall over?

Table 1: Example full conversations in our IGC<sub>Crowd</sub> dataset. For comparison, we also include VQG questions in which the image is the only context.



# Task Characteristics (Effectiveness of Multimodal Context)

- IGC task emphasizes modeling both visual and textual context
- Presented human judges with 600 (I,T,Q) triplets from each IGC<sub>Twitter</sub> and IGC<sub>Crowd</sub> and asked to rate *effectiveness* of visual and textual context
  - Whether the question makes sense without the image or text?

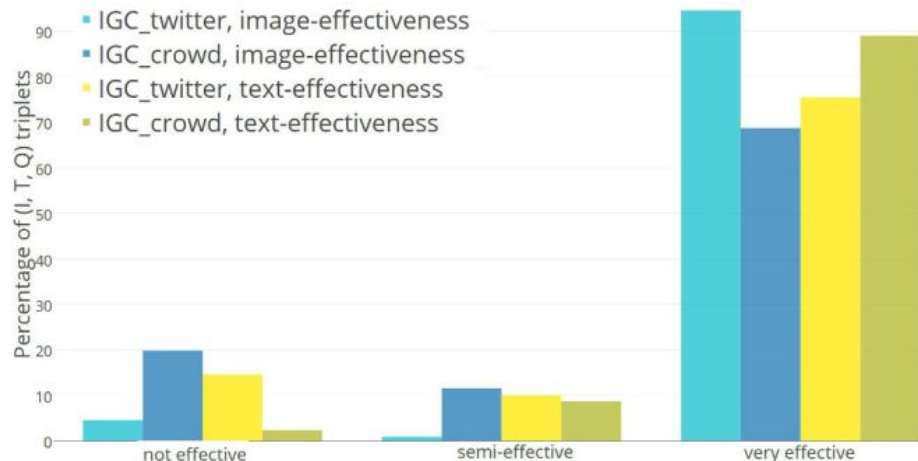


Figure 3: The effectiveness of textual and visual context for asking questions.

# Task Characteristics (Frame Semantic Analysis of Questions)

- Manually annotate a sample of 330 (I,T,Q) triplets in terms of Minsky's Frames
  - Frame: semantic representation of a situation involving participants, props and other conceptual roles
- Annotated the FrameNet frame evoked by I, T and Q.
  - 14%  $I_{FN} = T_{FN}$
  - 32%  $Q_{FN} = I_{FN}$
  - 47%  $Q_{FN} = T_{FN}$


Visual Context	Textual Context	Question
	Look at all this food I ordered!	Where is that from?
FN <i>Food</i>	<i>Request-Entity</i>	<i>Supplier</i>

Table 3: FrameNet (FN) annotation of an example.

# Task Characteristics (Event Analysis of Conversations)

Manually annotated 20 conversations with their causal and temporal event structure (CaTeRS Scheme)



User1: My son is ahead and surprised!  
User2: Did he end up winning the race?  
User1: Yes he won, he can't believe it!

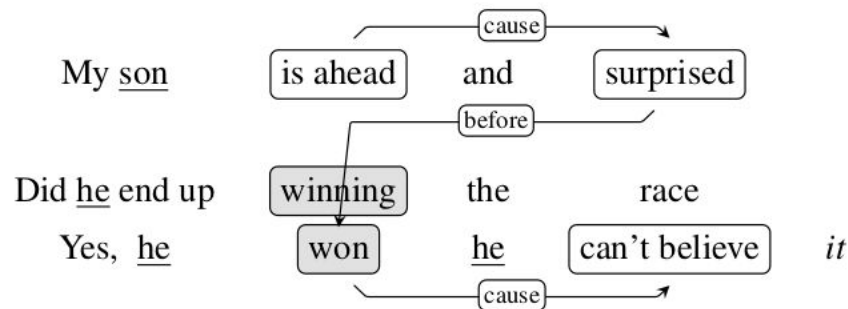


Figure 4: An example causal and temporal (CaTeRS) annotation on the conversation pre-

# Task Characteristics (Event Analysis of Conversations)

- Findings:
  - IGC utterances are rich in events (0.71 event entity mentions)
  - Semantic link annotations reflect common sense relations between event mentions in context of ongoing conversation
  - **Capturing causal and temporal relations between events is necessary for a system to successfully perform IGC task**

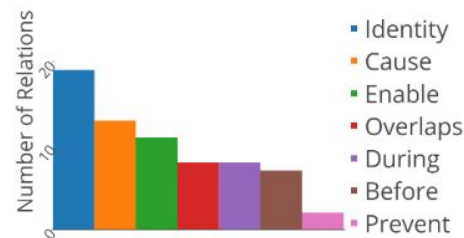
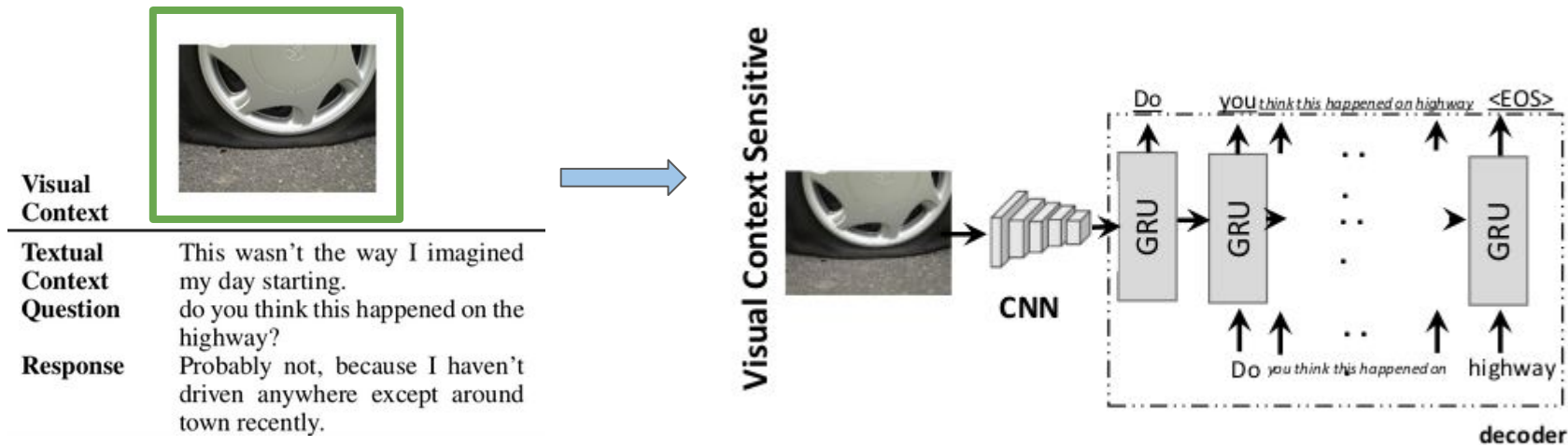


Figure 5: The frequency of event-event semantic links in a random sample of 20 IGC conversations.


# Models (Question Generation Models) $(I, T) \rightarrow Q$

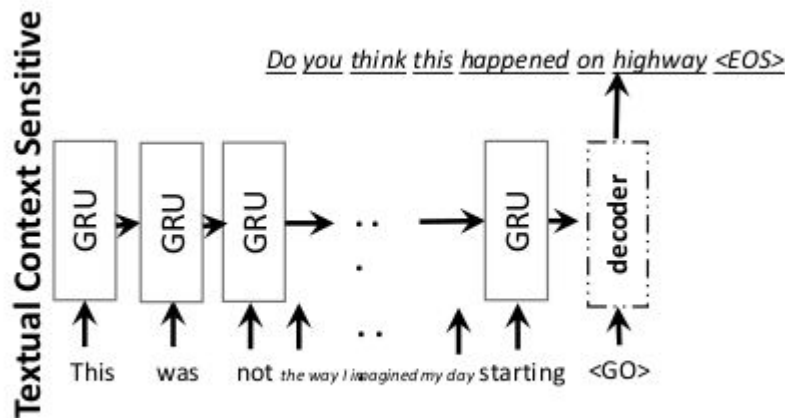
## 1. Visual Context Sensitive Model (V-Gen)



# Models (Question Generation Models) $(I, T) \rightarrow Q$

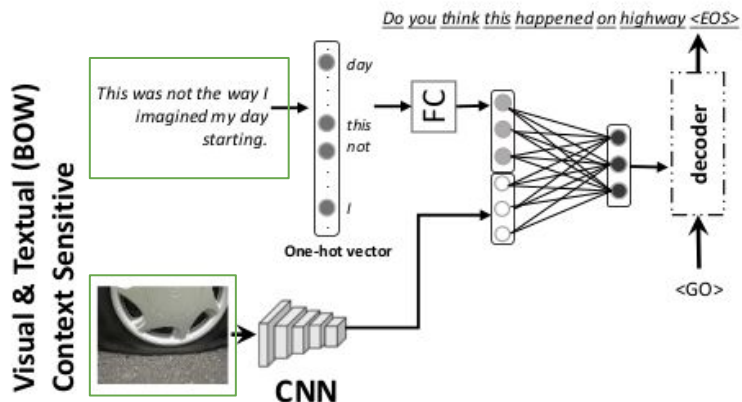
## 2. Textual Context Sensitive Model (T-Gen)

	
<b>Visual Context</b>	
<b>Textual Context</b>	This wasn't the way I imagined my day starting.
<b>Question</b>	do you think this happened on the highway?
<b>Response</b>	Probably not, because I haven't driven anywhere except around town recently.

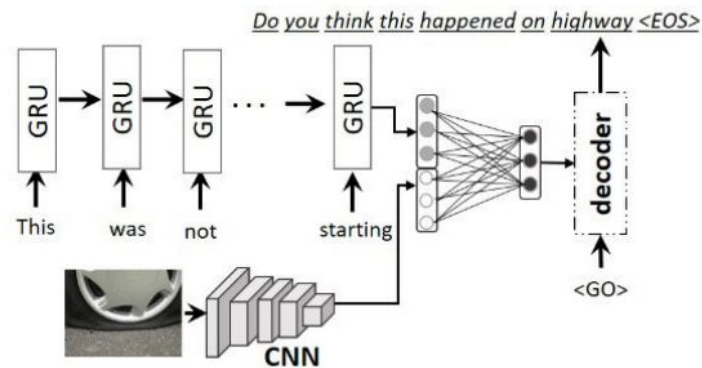


# Models (Question Generation Models) $(I, T) \rightarrow Q$

## 3. Visual & Textual Context Sensitive Model (V&T-Gen)




V&T.BOW-Gen

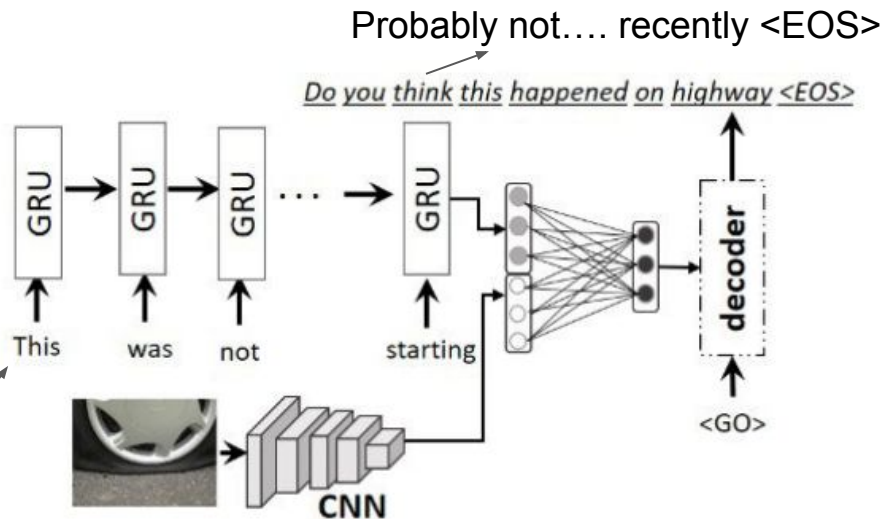


V&T.RNN-Gen

# Models (Response Generation Models) $(I, T, Q) \rightarrow R$

Visual Context	
Textual Context	This wasn't the way I imagined my day starting.
Question	do you think this happened on the highway?
Response	Probably not, because I haven't driven anywhere except around town recently.

Textual context:  
*This was not the way I imagined my day starting <UTT> Do you think this happened on the highway?*





# Models (Decoding and Reranking)

*applies to generation models*

- Greedy Decoding
- Beam Search
  - Generate N-best lists using left-to-right beam search (beam size = 25)
  - Max #tokens = 13
  - Any partial hypothesis that reaches <EOS> → viable for reranking
- Reranking
  - First few hypotheses on top of the N-best list tend to be generic
    - For example, “Where is this?”

$$\log p(h|C) + \lambda \text{idf}(h,D) + \mu |h| + \kappa V(h)$$

Probability of  
hypothesis(h) given  
context (C)

How common the  
hypothesis is across all  
generated N-best lists (D)

# tokens in  
hypothesis

# verb POS in  
hypothesis

# Models (Retrieval Models)

## 1. Visual Context Sensitive Model (V-Ret)

- a. Only uses the given image for retrieval
- b. Finds K nearest training images for the given test image based on cosine similarity of fc7 feature vector  $\rightarrow$  K candidates
- c. Compute textual similarity among the questions in the pool (Smoothed BLEU similarity score)
- d. Emit sentence with highest similarity to rest of the pool.

## 2. Visual & Textual Context Sensitive Model (V&T-Ret)

- a. Uses linear combination of fc7 and word2vec feature vectors
- b. Retrieval process is same as above

# Models (Recap)

- Question/Response generation

- V-Gen
- T-Gen
- V&T-Gen
  - V&T.BOW-Gen (question)
  - V&T.RNN-Gen (response)



- greedy
- beam search (best)
- reranked (best)

- Retrieval

- V-Ret
- V&T Ret

- Train: IGC<sub>Twitter</sub>

- Validation/Test: IGC<sub>Crowd</sub>

# Evaluation Setup

- Human evaluation
  - Crowdsourced on AMT-like system
  - 7 crowd workers rate quality of questions/responses on a scale of 1 to 3 (highest)
  - All system hypotheses are presented at the same time in random order
  - Also present the human gold standard
  - Average the score throughout the test set for each model and the human gold standard
- Automatic Evaluation
  - BLEU with equal weight up to 4 grams at corpus level on the multi-reference  $IGC_{\text{Crowd}}$  test set

# Experimental Results (Human Evaluation)

	Human	Generation (Greedy)			Generation (Beam, best)				Generation (Reranked, best)			Retrieval	
	Gold	Textual	Visual	V & T	Textual	Visual	V & T	VQG	Textual	Visual	V & T	Visual	V & T
Question	<u>2.68</u>	1.46	1.58	1.86	1.07	1.86	<b>2.28</b>	2.24	1.03	2.06	2.13	1.59	1.54
Response	<u>2.75</u>	1.24	–	1.40	1.12	–	<b>1.49</b>	–	1.04	–	1.44	–	1.48

## Key Takeaways:

- Multimodal V&T outperforms Textual and Visual
- Top generation in N-best list is preferred over reranked
- Human gold standards are favoured throughout the table
- IGC<sub>Crowd Test Set</sub> robust test set for benchmarking progress on this task

# Experimental Results (Automatic Evaluation)

	Textual	Generation			Retrieval	
		Visual	V & T	VQG	Visual	V & T
Question	1.71	3.23	4.41	<b>8.61</b>	0.76	1.16
Response	1.34	–	<b>1.57</b>	–	–	0.66

## Key Takeaways:

- BLEU scores are low
- Multimodal V&T outperforms all other models except VQG
- For both automatic and human evaluation → performance on question generation is better than response generation

# Experimental Results (Examples)




Visual Context				
Question Generation	<b>Textual Context</b>	The weather was amazing at this baseball game.	I got in a car wreck today!	My cousins at the family reunion.
	<b>Gold Question</b>	Nice, which team won?	Did you get hurt?	What is the name of your cousin in the blue shirt?
	<b>V&amp;T-Ret</b>	U at the game? or did someone take that pic for you?	<b>You driving that today?</b>	<b>U had fun?</b>
	<b>V-Gen</b> <b>V&amp;T-Gen</b>	Where are you? <b>Who's winning?</b>	Who's is that? <b>What happened?</b>	Who's that guy? <b>Where's my invite?</b>
Response Generation	<b>Textual Context</b>	The weather was amazing at this baseball game. <UTT> Nice, which team won?	I got in a car wreck today! <UTT> Did you get hurt?	My cousins at the family reunion. <UTT> What is the name of your cousin in the blue shirt?
	<b>Gold Response</b>	My team won this game.	No it wasn't too bad of a bang up.	His name is Eric.
	<b>V&amp;T-Ret</b>	10 for me and 28 for my dad.	<b>Yes.</b>	lords cricket ground . beautiful.
	<b>V&amp;T-Gen</b>	ding ding ding!	<b>Nah, I'm at home now.</b>	<b>He's not mine!</b>

Table 4: Example question and response generations on IGC<sub>Crowd</sub> test set. All the generation models use beam search with reranking. In the textual context, <UTT> separates different utterances. The generations in bold are acceptable utterances given the underlying context.

# Conclusion/Contributions

- Introduced a new task on multimodal image-grounded conversation for formulating questions and responses around images.
- Released to research community a crowdsourced dataset of 4222 high quality multi-turn conversations around eventful images and multiple references.
- Experiments suggest that capturing multimodal context improves the quality of question and response generation



# Thoughts

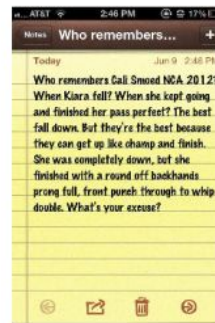
- First step into combining threads of language & vision and conversations
- Including other kinds of grounding
  - Temporal, geolocation ...
- Attention based mechanism
- Performance across different automatic evaluation metrics
- More to a conversation than just question and response
- Better quality training dataset compared to IGC<sub>Twitter</sub>

# IGC<sub>Twitter</sub> Training Dataset Problems

- Conversation is not grounded in image and textual context but rather in participant's established relation or prior history
  - Around 46% are like this!
- Abundance of screenshots



Smile.  
Why are you so obsessed  
with me?  
Oh pls



What's your excuse?  
Nca nationals? which day?  
  
Day 2 i believe ! if you go  
on youtube it should be the  
first one !

Table 8: Example Twitter conversations that add noise to the dataset.

Questions/Comments?

# Supplementary Material

# IGC<sub>Twitter</sub> Training Dataset Example

					
<b>Visual</b>	<b>Con-</b>				
<b>text</b>					
<b>Textual</b>	<b>Con-</b>	Oh my gosh, i'm so buying this shirt.	I found a cawaii bird.	Stocking up!!	Only reason I come to carnival.
<b>text</b>		Where did you see this for sale?	Are you going to collect some feathers?	Ayee! what the prices looking like?	Oh my God. How the hell do you even eat that?
<b>Question</b>					
<b>Response</b>		Midwest sports	There are so many crows here I'd be surprised if I never found one.	Only like 10-20% off..I think I'm gonna wait a little longer.	They are the greatest things ever chan. I could eat 5!

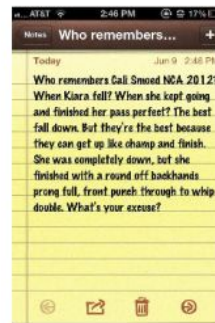
Table 7: Example conversations in the IGC<sub>Twitter</sub> dataset.

# IGC<sub>Twitter</sub> Training Dataset Problems

- Conversation is not grounded in image and textual context but rather in participant's established relation or prior history
  - Around 46% are like this!
- Abundance of screenshots



Smile.  
Why are you so obsessed  
with me?  
Oh pls



What's your excuse?  
Nca nationals? which day?  
  
Day 2 i believe ! if you go  
on youtube it should be the  
first one !

Table 8: Example Twitter conversations that add noise to the dataset.

# Data Analysis (length of sentences)

On average, IGC<sub>Twitter</sub> has longer sentences



Figure 8: Distribution of the number of tokens across datasets.

# Data Analysis (diversity in questions)



Figure 9: Distributions of n-gram sequences in questions in VQG, IGC<sub>Twitter</sub>, and IGC<sub>Crowd</sub>.

IGC<sub>Twitter</sub> is most diverse, with light-colored part of the circle indicating sequences with less than 0.1% representation in the data.



# Data Analysis (IGC questions characteristics)

- IGC<sub>Twitter</sub> has largest vocabulary size → challenging for training
- IGC<sub>Twitter</sub> and IGC<sub>Crowd</sub> have highest ratio of concrete to abstract nouns
- Contextually grounded questions of IGC<sub>Crowd</sub> are competitive with VQG in inter-annotation similarity

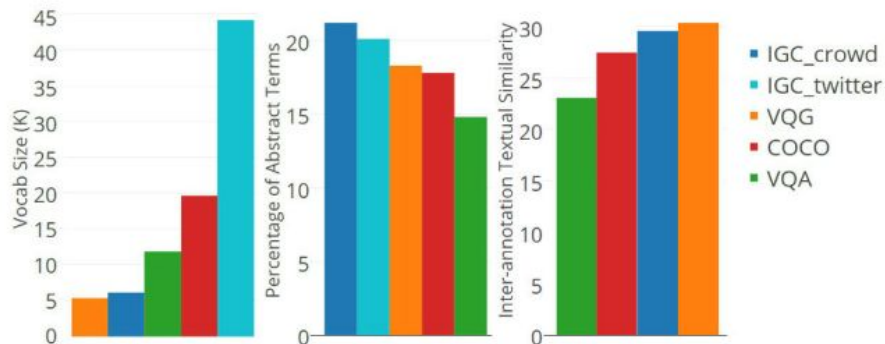


Figure 10: Comparison of V&L datasets.