# Image-Chat: Engaging Grounded Conversations

**Kurt Shuster, Samuel Humeau, Antoine Bordes, Jason Weston**

Facebook AI Research

{kshuster,samuelhumeau,abordes,jase}@fb.com

## Abstract

To achieve the long-term goal of machines being able to engage humans in conversation, our models should captivate the interest of their speaking partners. Communication grounded in images, whereby a dialogue is conducted based on a given photo, is a setup naturally appealing to humans (Hu et al., 2014). In this work we study large-scale architectures and datasets for this goal. We test a set of neural architectures using state-of-the-art image and text representations, considering various ways to fuse the components. To test such models, we collect a dataset of grounded human-human conversations, where speakers are asked to play roles given a provided emotional mood or style, as the use of such traits is also a key factor in engagingness (Guo et al., 2019). Our dataset, Image-Chat, consists of 202k dialogues over 202k images using 215 possible style traits. Automatic metrics and human evaluations of engagingness show the efficacy of our approach; in particular, we obtain state-of-the-art performance on the existing IGC task, and our best performing model is almost on par with humans on the Image-Chat test set (preferred 47.7% of the time).

## 1 Introduction

A key way for machines to exhibit intelligence is for them to be able to perceive the world around them – and to be able to communicate with humans in natural language about that world. To speak naturally with humans it is necessary to understand the natural things that humans say about the world they live in, and to respond in kind. This involves understanding what they perceive, e.g. the images they see, what those images mean semantically for humans, and how mood and style shapes the language and conversations derived from these observations.

In this work we take a step towards these goals by considering grounded dialogue involving open-ended discussion of a given image, a setting that is naturally fun for humans (Hu et al., 2014), and study neural conversational models for task. In particular, we explore both generative and retrieval models that handle multimodal dialogue by fusing Transformer architectures (Vaswani et al., 2017) for encoding dialogue history and responses and ResNet architectures (He et al., 2016) for encoding images. We propose ways to fuse those modalities together and perform a detailed study including both automatic evaluations, ablations and human evaluations of our models using crowdworkers.

To train and evaluate such models, we collect a large set of human-human crowdworker conversations, with the aim of training a model to engage a human in a similar fashion, consisting of 202k diverse images and 401k utterances over the images, with 215 different style traits (e.g., optimistic, skeptical or frivolous) to promote engaging conversation. The dataset is made publicly available in ParlAI (Miller et al., 2017) [1].

Our results show that there is a significant gap between state-of-the-art retrieval and generative models on this task. Our best fused retrieval models set a strong baseline, being preferred to human conversationalists 47.7% of the time. We show that both large-scale image and text pre-training, and utilization of style traits, are critical for best results. We then consider transfer to the existing Image Grounded Conversations (IGC) task of Mostafazadeh et al. (2017), where we obtain state-of-the-art results.

## 2 Related Work

The majority of work in dialogue is not grounded in perception, e.g. much recent work explores sequence-to-sequence models or retrieval models for goal-directed (Henderson et al., 2014) or chit-

---

[1] http://parl.ai/projects/image_chat

chat tasks (Vinyals and Le, 2015; Zhang et al., 2018). While these tasks are text-based only, many of the techniques developed can likely be transferred for use in multimodal systems, for example using state-of-the-art Transformer representations for text (Mazare et al., 2018) as a sub-component.

In the area of language and vision, one of the most widely studied areas is image captioning, whereby a single utterance is output given an input image. This typically involves producing a factual, descriptive sentence describing the image, in contrast to producing a conversational utterance as in dialogue. Popular datasets include COCO (Chen et al., 2015) and Flickr30k (Young et al., 2014). Again, a variety of sequence-to-sequence (Vinyals et al., 2015; Xu et al., 2015; Anderson et al., 2018) and retrieval models (Gu et al., 2018; Faghri et al., 2018; Nam et al., 2016) have been applied. These tasks measure the ability of models to understand the content of an image, but not to carry out an engaging conversation grounded in perception. Some works have extended image captioning from being purely factual towards more engaging captions by incorporating style while still being single turn, e.g. (Mathews et al., 2018, 2016; Gan et al., 2017; Guo et al., 2019; Shuster et al., 2019). Our work also applies a style component, but concentrates on image-grounded dialogue, rather than image captioning.

Visual question answering (Antol et al., 2015) and visual dialogue (Das et al., 2017) are another set of tasks which employ vision and language. They require the machine to answer factual questions about the contents of the image, either in single turn or dialogue form. They do not attempt to model natural conversation, but rather assess whether the machine can perform basic perception over the image via a series of questions.

There are some works which directly address dialogue grounded with vision. The work of Pasunuru and Bansal (2018) assesses the ability to execute dialogue given video of computer soccer games. The work of Huber et al. (2018) investigates the use of sentiment-based visual features and facial expressions for emotional image-based dialogue. Perhaps the most related work to ours is Mostafazadeh et al. (2017). Their work considers (visual context, textual context, question, response) tuples, and builds validation and test sets based on 4k eventful images called Image Grounded Conversations (IGC). No training data is provided, but instead the authors use Twitter for that in their experiments. In contrast, we provide training, validation and testing sets over 202k images for our task (that do not overlap with IGC), and consider a general set of images and dialogues, not just events and questions plus responses. In our experiments we also show strong transfer ability of our models to the IGC task.

While there are many ways to measure dialogue quality, human engagement is a popular metric. Engagement itself can be measured in many ways (Bohus and Horvitz, 2009; Yu et al., 2016) but here we adopt the common approach of simply asking humans which speaker they find more engaging, following other works (Li et al., 2019; Dinan et al., 2020).

## 3 Image-Chat

The IMAGE-CHAT dataset is a large collection of (image, style trait for speaker A, style trait for speaker B, dialogue between A & B) tuples that we collected using crowd-workers, Each dialogue consists of consecutive turns by speaker A and B. No particular constraints are placed on the kinds of utterance, only that we ask the speakers to both use the provided style trait, and to respond to the given image and dialogue history *in an engaging way*. The goal is not just to build a diagnostic dataset but a basis for training models that humans actually want to engage with.

**Style Traits** A number of works have shown that style traits for image captioning help provide creative captions (Mathews et al., 2018, 2016; Gan et al., 2017; Shuster et al., 2019). We apply that same principle to image grounded dialogue, considering a set of 215 possible style traits, using an existing set from Shuster et al. (2019). The traits are categorized into three classes: positive (e.g., sweet, happy, eloquent, humble, witty), neutral (e.g., old-fashioned, skeptical, solemn, questioning) and negative (e.g., anxious, childish, critical, fickle, frivolous). We apply these to both speakers A and B, who will be assigned different style traits for each given conversation.

**Images** The images used in our task are randomly selected from the YFCC100M Dataset[2] (Thomee et al., 2016).

**Dialogue** For each image, we pick at random two style traits, one for speaker A and one for speaker

*A: Peaceful    B: Absentminded*

A: I'm so thankful for this delicious food.

B: What is it called again?

A: Not sure but fried goodness.

*A: Fearful    B: Miserable*

A: I just heard something out there and I have no idea what it was.

B: It was probably a Wolf coming to eat us because you talk too much.

A: I would never go camping in the woods for this very reason.

*A: Erratic    B: Skeptical*

A: What is the difference between the forest and the trees? Oh look, dry pavement.

B: I doubt that's even a forest, it looks like a line of trees.

A: There's probably more lame pavement on the other side!

Figure 1: Some samples from the IMAGE-CHAT training set. For each sample we asked humans to engage in a conversation about the given image, where the two speakers, A and B, each have a given provided style.

B, and collect the dialogue using crowdworkers who are asked to both assume those roles, and to be engaging to the other speaker while doing so. It was emphasized in the data collection instructions that the style trait describes a trait of the speaker, not properties of the content of the image they are discussing. Some examples from the training set are given in Figure 1.

**Data Quality**  During data collection crowdsourcers were manually monitored, checking to ensure they were following the instructions. Poor performers were banned, with comments discarded. A verification process was also conducted on a subset of the data, where separate annotators were asked to choose whether the utterance fit the image, style, or both, and found that 92.8% of the time it clearly fit the image, and 83.1% the style, and 80.5% both. Note, given that not all utterances should directly reference an image property or invoke the style, we do not expect 100%.

**Overall Dataset**  The overall dataset statistics are given in Table 1. This is a fairly large dialogue dataset compared to other existing publicly available datasets. For example, PersonaChat (Zhang et al., 2018) (which is not grounded in images) consists of 162k utterances, while IGC (Mostafazadeh et al., 2017) (grounded in images) consists of 4k of validation and test set examples only, compared to over 400k utterances in IMAGE-CHAT.

| Split | train | valid | test |
|---|---|---|---|
| Number of Images | 186,782 | 5,000 | 9,997 |
| Number of Dialogues | 186,782 | 5,000 | 9,997 |
| Number of Utterances | 355,862 | 15,000 | 29,991 |
| Style Types | 215 | 215 | 215 |
| Vocabulary Size | 46,371 | 9,561 | 13,550 |
| Tokens per Utterance | 12.3 | 12.4 | 12.4 |

Table 1: IMAGE-CHAT dataset statistics.

## 4 Models

We consider two major types of dialogue model: retrieval and generative. Both approaches make use of the same components as building blocks. We use three sub-networks for the three modalities of input: (i) an image encoder, (ii) a dialogue history encoder; and (iii) a style encoder. In the retrieval model these are then fed into a combiner module for combining the three modalities. Finally, there is a response encoder for considering candidate responses and this is scored against the combined input representations. An overview of the retrieval architecture is shown in Figure 2. For the generative model, the three encoders are used as input, and a further decoder Transformer is used for outputting a token sequence; beam search is applied.

**Image Encoder**  We build our models on top of pretrained image features, and compare the performance of two types of image encoders. The first is a residual network with 152 layers described in He et al. (2016) trained on ImageNet (Russakovsky et al., 2015) to classify images among 1000 classes, which we refer to in the rest of the pa-
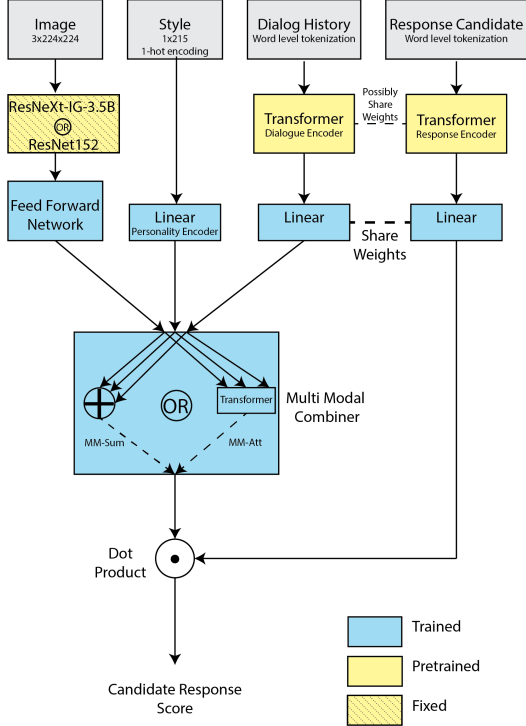
Figure 2: The TRANSRESNET$_{RET}$ multimodal architecture for grounded dialogue. There are several options: different image encoders (ResNet152 or ResNeXt-IG-3.5B), text encoders (shared or separate Transformers for history and response), and different multimodal combiners (sum or attention-based).

per as *ResNet152* features. We used the implementation provided in the torchvision project (Marcel and Rodriguez, 2010). The second is a ResNeXt $32 \times 48d$ (Xie et al., 2017) trained on 3.5 billion Instagram pictures following the procedure described by Mahajan et al. (2018), which we refer to in the rest of the paper as *ResNeXt-IG-3.5B*. The representation $r_I$ of an image $I$ is obtained by using the 2048-dimensional output of the image encoder as input to a feed-forward network: a multi-layer perceptron with ReLU activation units and a final layer of 500 dimensions in the retrieval case, and a linear layer in the generative case.

**Style Encoder** To condition on a given style trait, we embed each trait to an $N$-dimensional vector to obtain its representation $r_S$. We used $N = 500$ for retrieval and $N = 300$ for generation.

**Dialogue Encoder** The entire dialogue history $D$ is encoded into a fixed size vector $r_D$ using a Transformer architecture (Vaswani et al., 2017), followed by a linear layer. Such Transformers have been shown to perform strongly on a variety of dia-

logue tasks previously (Yang et al., 2018; Mazare et al., 2018). We use a Transformer with 4 layers, 300 hidden units, and 6 attention heads. The outputs are pooled (mean) to give a final vectorial encoding.

We pretrain the entire encoder following the setup described in Mazare et al. (2018): we train two encoders on a next-utterance retrieval task on a Reddit dataset of dialogues containing 1.7 billion pairs of utterances, where one encodes the context and another the candidates for the next utterance; their dot product indicates the degree of match, and they are trained with negative log-likelihood and $k$-negative sampling. We then initialize our system using the weights of the candidate encoder only, and then train on our task in either generative or retrieval mode.

### 4.1 Retrieval Models

**Multimodal combiner module** We consider two possible combiner modules for the inputs:

*Multimodal sum combiner (MM-sum)*: Given an input image, style trait and dialogue $(I, S, D)$, together with a candidate response $C$, the score of the final combination is computed as $s(I, S, D, C) = (r_I + r_S + r_D) \cdot r_C$.

*Multimodal attention combiner (MM-att)*: A more sophisticated approach is to use an attention mechanism to choose which modalities are most relevant for each example by stacking Transformers. We concatenate the three representation vectors $r_I$, $r_S$ and $r_D$ and feed them to a second Transformer (4 attention heads, 2 layers, 500 hidden units) which performs self-attention over them. The three modalities are thus reweighted by the corresponding attention weights to give the final input representation vector $r_T$, which is used to compute the score for a given candidate using $r_T \cdot r_C$.

**Response encoder** We employ the same Transformer architecture as in the dialogue encoder for encoding candidate responses. We tried two variants: either sharing or not sharing the weights with the input dialogue encoder.

**Training and Inference** Given a tuple $I, S, D$, and a set of candidates $(c_1, .., c_N)$, at inference time the predicted utterance is the candidate $c_i$ that maximizes the score $s(I, S, D, c_i)$. At training time we pass a set of scores through a softmax and train to maximize the log-likelihood of the correct responses. We use mini-batches of 500 training

examples; for each example, we use the gold responses of the other examples of the batch as negatives. During final human evaluation all candidates from the training set are considered to produce a response (356k candidates in our experiments).

## 4.2 Generative Models

**Dialogue Decoder** The encoding from the image encoder has a final linear layer of dimension 2048 × 300. This projects it to the same size of the token encoding of the dialogue decoder. We thus add it as an extra token at the end of the Transformers encoder output. For style, we simply prepend the style to the beginning of the dialogue history, and it is thus encoded in the dialogue encoder. We then treat this as a standard seq2seq Transformer in order to generate dialogue responses.

**Training and Inference** We train with a batch size of 32 and learning rate of .0001 using adam, and apply beam search with a beam of size 2 and trigram blocking at inference time. Hyperparameters are chosen on the validation set.

## 5 Experiments

We test our models on the IMAGE-CHAT and IGC datasets using automatic metrics and human evaluations. We analyze the performance of the different module and architecture choices, as well as ablation studies to determine the importance of each of the model's inputs.

### 5.1 Automatic Evaluation on IMAGE-CHAT

**Module Choices** We first compare various module configurations of our TRANSRESNET$_{RET}$ model, and additionally show the results for a simple information retrieval baseline, in which the candidates are ranked according to their weighted word overlap to the input message. We measure recall at 1 and 5 (R@1/100 and R@5/100) retrieval metrics, where for each sample there are 100 candidates to rank: 99 random candidates chosen from the test set, and the true label. Note that in human evaluations we use all the train set candidates.

The results are shown in Table 2. We report the average metrics for the total task, as well as the breakdown of the performance on each turn of dialogue (turns 1, 2 and 3). The average metrics indicate that using the ResNeXt-IG-3.5B image encoder features improves performance significantly across the whole task, as we obtain 50.3% R@1 for our best ResNeXt-IG-3.5B model and only 40.6%

for our best ResNet152 model. When broken down by turn, it appears that the ResNeXt-IG-3.5B features are particularly important in the first round of dialogue, in which only the image and style are considered, as the difference between their best models increases from 9.7% in the full task to 19.5% in the first turn. Our baseline multimodal sum combiner (MM-Sum) outperforms the more sophisticated self-attention (MM-Att) combiner, with the latter scoring 49.3% on the full task. Having separate candidate and dialogue history text encoders also works better than sharing weights.

In subsequent experiments we use the best performing system for our retrieval model. As ResNeXt-IG-3.5B performs best we use that for our generative model going forward as well.

**Full & Ablation Study** We now perform experiments for both retrieval and generative models for the full system, and additionally we remove modalities (image, style, and dialogue history). For the generative models we report the ROUGE-L metric. The results are shown in Table 3, which we now analyze.

*Turn 1:* In the first round of dialogue the models produce utterances given the image and style only, as there is no dialogue history yet. For both models, image is more important than style, but using both together helps.

*Turn 2:* In the second turn, in which a model produces a response to a first utterance, the models perform similarly when using only the image or only the dialogue history, while performing poorly with just the style. Any combination of two modalities improves the results, with the style + dialogue combination performing slightly higher than the other two. Using all modalities works best.

*Turn 3:* By the third turn of dialogue, the conversation history proves to be by far the most important in isolation compared to the other two modalities in isolation. Conditioning on the style+dialogue is the most effective of any combination of two modalities. Again, using all modalities still proves best.

### 5.2 Human Evaluations on IMAGE-CHAT

We test our final models using human evaluation.

**Evaluation Setup** We use a set of 500 images from YFCC-100M that are not present in IMAGE-CHAT to build a set of three-round dialogues pairing humans with models in conversation. We then

| Model | Combiner | Text Encoders R@1 | Image Encoder R@1 | Turn 1 R@1 | Turn 2 R@1 | Turn 3 R@1 | All R@1 | All R@5 |
|---|---|---|---|---|---|---|---|---|
| IR Baseline | n/a | n/a | n/a | - | - | - | 2.15 | 5.86 |
| TRANSRESNET$_{RET}$ | MM-Att | Separate | ResNet152 | 35.7 | 44.5 | 40.5 | 40.2 | 67.0 |
| TRANSRESNET$_{RET}$ | MM-Sum | Separate | ResNet152 | 34.5 | 46.0 | 41.3 | 40.6 | 67.2 |
| TRANSRESNET$_{RET}$ | MM-Sum | Shared | ResNeXt-IG-3.5B | 53.6 | 47.0 | 41.3 | 47.3 | 73.1 |
| TRANSRESNET$_{RET}$ | MM-Att | Shared | ResNeXt-IG-3.5B | **54.4** | 49.0 | 43.3 | 48.9 | 74.2 |
| TRANSRESNET$_{RET}$ | MM-Att | Separate | ResNeXt-IG-3.5B | 53.5 | 50.5 | 43.8 | 49.3 | 74.7 |
| TRANSRESNET$_{RET}$ | MM-Sum | Separate | ResNeXt-IG-3.5B | 54.0 | **51.9** | **44.8** | **50.3** | **75.4** |

Table 2: Module choices on IMAGE-CHAT. We compare different module variations for TRANSRESNET$_{RET}$.

| | TRANSRESNET$_{RET}$ (R@1/100 ) | | | | TRANSRESNET$_{GEN}$ (ROUGE-L) | | | |
|---|---|---|---|---|---|---|---|---|
| Modules | Turn 1 | Turn 2 | Turn 3 | All | Turn 1 | Turn 2 | Turn 3 | All |
| Image Only | 37.6 | 28.1 | 20.7 | 28.7 | 21.1 | 21.9 | 22.4 | 21.8 |
| Style Only | 18.3 | 15.3 | 17.0 | 16.9 | 20.2 | 20.9 | 22.0 | 21.0 |
| Dialogue History Only | 1.0 | 33.7 | 32.3 | 22.3 | 18.9 | 22.7 | 23.7 | 21.8 |
| Style + Dialogue (*no image*) | 18.3 | 45.4 | 43.1 | 35.4 | 20.4 | 24.1 | 24.8 | 23.1 |
| Image + Dialogue (*no style*) | 37.6 | 39.4 | 32.6 | 36.5 | 21.3 | 22.8 | 23.6 | 22.6 |
| Image + Style (*no dialogue*) | **54.0** | 41.1 | 35.2 | 43.4 | **23.7** | 23.2 | 23.8 | 23.5 |
| Style + Dialogue + Image (*full model*) | **54.0** | 51.9 | 44.8 | 50.3 | **23.7** | 24.2 | 24.9 | 24.3 |

Table 3: Ablations on IMAGE-CHAT. We compare variants of our best TRANSRESNET generative and retrieval models (ResNeXt-IG-3.5B image encoder, and MM-Sum + separate text encoders for retrieval) where we remove modalities: image, dialogue history and style conditioning, reporting R@1/100 for retrieval and ROUGE-L for generation for dialogue turns 1, 2 and 3 independently, as well as the average over all turns.

conduct evaluations at each round of dialogue for each example in the evaluation set; we have a separate set of human evaluators look at the provided conversation turns, and ask them to compare two possible utterances for the next turn of conversation, given the image, dialogue history and relevant style (which is the same for both human author and model, so there is no advantage). We ask the evaluators in a blind test to choose the "more engaging" of the two possible utterances: one from a human, and the other from a model.

**Human annotation vs. TRANSRESNET model** We compare human-authored utterances to those produced by our models. The human conversations are collected in the same fashion as in IMAGE-CHAT but on test images. As for humans, the model outputs are conditioned on the image, style and previous dialogue history. TRANSRESNET$_{GEN}$ simply generates a response, whereas TRANSRESNET$_{RET}$ retrieves candidate utterances from the IMAGE-CHAT training set. The latter is given a separate set of candidates corresponding to the round of dialogue – e.g. when producing a response to turn 1, the model retrieves from all possible round 1 utterances from the train set (in that case 186,858 possible choices).

The results are shown in Fig. 4, comparing all models on the first round (left): TRANSRESNET$_{GEN}$ and TRANSRESNET$_{RET}$ us-

ing ResNeXt-IG-3.5B, and TRANSRESNET$_{RET}$ using ResNet152 features. As in automatic evaluations, ResNet152 features performed more poorly. The retrieval model outperformed the generative model, a result that has been observed in other (text-only) dialogue tasks (Dinan et al., 2019; Zhang et al., 2018). In turn 1, TRANSRESNET$_{RET}$ (ResNeXt-IG-3.5B) has a win rate against humans of 49.4% (difference not significant using a binomial two-tailed test, $p > 0.5$), while both other models are significantly outperformed by humans ($p < 2 \times 10^{-7}$ compared to ResNet152 features), showing the importance of our retrieval architecture and image feature choices. We thus compare only TRANSRESNET$_{RET}$ (ResNeXt-IG-3.5B) to humans in all three turns (Fig. 4, right). That model performs well, with an overall win rate against humans of 47.7% (difference is significant, $p < 7 \times 10^{-5}$). Example predictions of TRANSRESNET$_{RET}$ (ResNeXt-IG-3.5B) are given in Figure 3.

## 5.3 Transfer to the IGC Task

To test the strength of our task and models we consider transfer to the IGC of task of Mostafazadeh et al. (2017). In particular, we focus on their response task, which provides an image and a dialogue history of two utterances: a context utterance, followed by a question. The task is to then pro-

| Image | Style | Conversation |
|---|---|---|
| | **Turn 1 examples** | |
| | A: Artful<br>A: Monstrous<br>A: Earnest | *Model predictions:*<br>This looks like a painting. Something out of a Norman Rockwell.<br>Some wretched town!.<br>Yeah, we have finally arrived at Grandpa's old barnhouse! I can't wait<br>to get started with the restoration! |
| | A: Wishful<br>A: Opinionated<br>A: Imaginative | *Model predictions:*<br>I hope one day to have a dog this majestic.<br>This puppy looks cold get him a blanket.<br>Puppies are just the universe's way of telling us everything will be okay. |
| | A: Respectful<br><br>A: Shy<br><br>A: Confused | *Model predictions:*<br>What an honor to have beautiful places like these to contemplate<br>natures rocks at their best.<br>While the area looks nice, I'd stay away, fear of tsunamis leaves me<br>paralyzed at the beach.<br>I don't understand why would you put rocks there? |
| | **Turn 2 examples** | |
| | A: Empathetic<br><br>B: Aloof<br>B: Frivolous | I understand that an open field full of beautiful foliage is a testament to<br>our environment.<br>*Model predictions:*<br>It's just a regular old field.<br>So is every grass in the country |
| | A: Playful<br><br>B: Questioning<br>B: Amusing | i would jump inside that waterfall<br>*Model predictions:*<br>Do you really think it is feasible?<br>Well I would try to top you with a backflip. |
| | A: Maternal<br><br>B: Whimsical<br>B: Blunt | It must be very dangerous if children play there.<br>*Model predictions:*<br>Nonsense, this would be a great place to play!<br>I would rather find some lawn. Those rocks look uncomfortable |
| | **Turn 3 examples** | |
| | A: Happy<br>B: Anxious<br><br>A: Happy | This looks like something fun and exciting to do on vacation!<br>oh dear what happens if he gets taken with the wind<br>*Model prediction:*<br>Ooh, the wind looks perfect to fly this awesome kite! |
| | A: Pompous<br><br>B: Offhand<br><br>A: Pompous | This archway is so dull and boring. I could have done a much better<br>job with the garden in my own back yard.<br>Its about as boring as you!!<br>*Model prediction:*<br>Please. I could show you much prettier places in my own back yard. |
| | A: Sweet<br><br>B: Spontaneous<br><br>A: Sweet | This is so beautiful, I love nature and would love to see this place and<br>meet the people who live there!<br>Then lets go there right now! Just you and me.<br>*Model prediction:*<br>I wish you could take me there, I would love to go with you! |

Figure 3: Example predictions from our TRANSRESNET$_{RET}$ (MM-Sum) model on the evaluation set using all candidates for turns 1–3 . Two speakers A & B with given style traits discuss a photo. The dialogue context before the model prediction is completed by humans, followed by one or more possible model responses, given different style conditioning. The model clearly uses the image, given style and dialogue history in formulating its response.
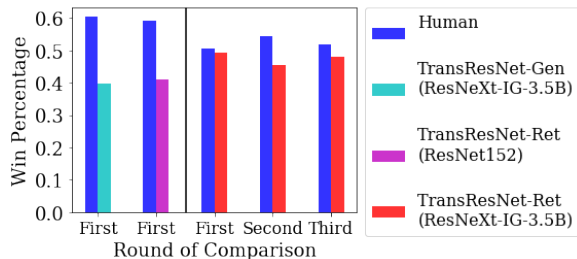
Figure 4: Human evaluations on IMAGE-CHAT. Engagingness win rates of pairwise comparisons between human utterances and TRANSRESNET$_{RET}$ (ResNet152 or ResNeXt-IG-3.5B) or TRANSRESNET$_{GEN}$, comparing over the rounds of dialogue.



Figure 5: IGC Evaluations. The best model from Mostafazadeh et al. (2017) is compared to our best TRANSRESNET$_{RET}$ and TRASNRESNET$_{GEN}$ models. On the left, annotator's ratings of responses from the models are shown as a percentage of the annotator's ratings of human responses. On the right, BLEU-4 scores on the response task are shown.

duce a response. This is clearly related to our task, except it focuses on answering questions, which our task does not. Our task is more varied as it was collected in an unconstrained way, unlike in IGC where they were asked to write a question. Nevertheless, assuming a question contains a *?* or starts with *who*, *what*, *when*, *where*, *why* or *how*, our dataset contains 40,076 training utterances that are questions (11.3% of the data) and so it could be possible to produce responses to them. Without any fine-tuning at all, we thus simply took exactly the same best trained models and used them for their question response task as well.

Unfortunately, after contacting the authors of Mostafazadeh et al. (2017) they no longer have the predictions of their model available, nor have they made available the code for their human evaluation setup. However, the test set is available. We therefore attempted to reproduce the same setup as in their experiments, which we will also make publicly available upon acceptance.

**Automatic Evaluation**  We measure our best TRANSRESNET$_{GEN}$ model's performance on the IGC test set in terms of BLEU-4. The results are shown in Fig. 5 (right). We find that our model outperforms the model from Mostafazadeh et al. (2017), achieving a score of 2.30 compared to 1.49.

**Human Evaluation**  We compare the provided human response (from the test set) with 7 variants of our TRANSRESNET$_{RET}$ model (mimicking their setup), whereby we have our model condition on 7 styles for which it performed well on evaluations in section 5.2. Annotators rated the quality of responses on a scale from 1 to 3, where 3 is the highest, reporting the mean over ∼2k questions. We then scale that by the score of human authored
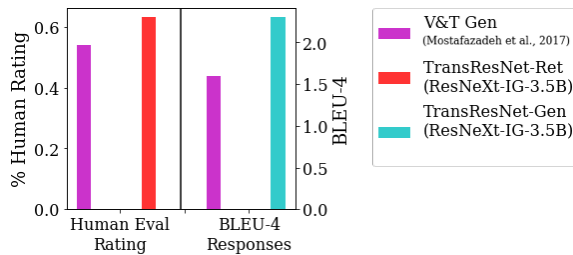
responses, to give a percentage. The results are shown in Fig. 5 (left). Our model narrows the gap between human and model performance, yielding a higher percentage of the human score (62.9% vs. 54.2%). More detailed results and example predictions of our model can be found in Appendices E and F, including examples of highly rated and poorly rated outputs from our model.

## 6  Conclusion

This paper presents an approach for improving the way machines can generate grounded conversations that humans find engaging. Focusing on the case of chit-chatting about a given image, a naturally useful application for end-users of social dialogue agents, this work shows that our best proposed model can generate grounded dialogues that humans prefer over dialogues with other fellow humans almost half of the time (47.7%). This result is made possible by the creation of a new dataset IMAGE-CHAT[3].

Our work shows that we are close to having models that humans can relate to in chit-chat conversations, which could set new ground for social dialogue agents. However, our retrieval models outperformed their generative versions; closing that gap is an important challenge for the community. While our human evaluations were on short conversations, initial investigations indicate the model as is can extend to longer chats, see Appendix G, which should be studied in future work. The next challenge will also be to combine this engagingness with other skills, such as world knowledge (Antol et al., 2015) relation to personal interests (Zhang et al., 2018), and task proficiency.

---

[3]http://parl.ai/projects/image_chat

## References

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 vqa. *CVPR*.

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433.

Dan Bohus and Eric Horvitz. 2009. Models for multi-party engagement in open-world dialog. In *Proceedings of the SIGDIAL 2009 Conference: The 10th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 225–234. Association for Computational Linguistics.

Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. 2015. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*.

Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José MF Moura, Devi Parikh, and Dhruv Batra. 2017. Visual dialog. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.

Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander Miller, Kurt Shuster, Jack Urbanek, Douwe Kiela, Arthur Szlam, Iulian Serban, Ryan Lowe, et al. 2020. The second conversational intelligence challenge (convai2). In *The NeurIPS'18 Competition*, pages 187–208. Springer.

Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2019. Wizard of Wikipedia: Knowledge-powered conversational agents. In *Proceedings of the International Conference on Learning Representations (ICLR)*.

Fartash Faghri, David J Fleet, Jamie Ryan Kiros, and Sanja Fidler. 2018. Vse++: Improving visual-semantic embeddings with hard negatives.

Chuang Gan, Zhe Gan, Xiaodong He, Jianfeng Gao, and Li Deng. 2017. Stylenet: Generating attractive visual captions with styles. In *Proc IEEE Conf on Computer Vision and Pattern Recognition*, pages 3137–3146.

J. Gu, J. Cai, S. Joty, L. Niu, and G. Wang. 2018. Look, imagine and match: Improving textual-visual cross-modal retrieval with generative models. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7181–7189.

Longteng Guo, Jing Liu, Peng Yao, Jiangwei Li, and Hanqing Lu. 2019. Mscap: Multi-style image captioning with unpaired stylized text. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4204–4213.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.

Matthew Henderson, Blaise Thomson, and Jason D Williams. 2014. The second dialog state tracking challenge. In *Proceedings of the 15th Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL)*, pages 263–272.

Yuheng Hu, Lydia Manikonda, and Subbarao Kambhampati. 2014. What we instagram: A first analysis of instagram photo content and user types. In *Eighth International AAAI Conference on Weblogs and Social Media*.

Bernd Huber, Daniel McDuff, Chris Brockett, Michel Galley, and Bill Dolan. 2018. Emotional dialogue generation using image-grounded language models. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, page 277. ACM.

Margaret Li, Jason Weston, and Stephen Roller. 2019. Acute-eval: Improved dialogue evaluation with optimized questions and multi-turn comparisons. *arXiv preprint arXiv:1909.03087*.

Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens van der Maaten. 2018. Exploring the limits of weakly supervised pretraining. In *Computer Vision – ECCV 2018*, pages 185–201, Cham. Springer International Publishing.

Sébastien Marcel and Yann Rodriguez. 2010. Torchvision the machine-vision package of torch. In *Proceedings of the 18th ACM International Conference on Multimedia*, MM '10, pages 1485–1488. ACM.

Alexander Mathews, Lexing Xie, and Xuming He. 2018. Semstyle: Learning to generate stylised image captions using unaligned text. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8591–8600.

Alexander Patrick Mathews, Lexing Xie, and Xuming He. 2016. Senticap: Generating image descriptions with sentiments. In *AAAI*, pages 3574–3580.

Pierre-Emmanuel Mazare, Samuel Humeau, Martin Raison, and Antoine Bordes. 2018. Training millions of personalized dialogue agents. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2775–2779, Brussels, Belgium. Association for Computational Linguistics.

A. H. Miller, W. Feng, A. Fisch, J. Lu, D. Batra, A. Bordes, D. Parikh, and J. Weston. 2017. Parlai: A dialog research software platform. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 79–84.

Nasrin Mostafazadeh, Chris Brockett, Bill Dolan, Michel Galley, Jianfeng Gao, Georgios Spithourakis, and Lucy Vanderwende. 2017. Image-grounded conversations: Multimodal context for natural question and response generation. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 462–472, Taipei, Taiwan. Asian Federation of Natural Language Processing.

Hyeonseob Nam, Jung-Woo Ha, and Jeonghee Kim. 2016. Dual attention networks for multimodal reasoning and matching. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2156–2164.

Ramakanth Pasunuru and Mohit Bansal. 2018. Game-based video-context dialogue. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 125–136, Brussels, Belgium. Association for Computational Linguistics.

Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. 2015. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252.

Kurt Shuster, Samuel Humeau, Hexiang Hu, Antoine Bordes, and Jason Weston. 2019. Engaging image captioning via personality. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Bart Thomee, David A. Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. 2016. Yfcc100m: The new data in multimedia research. *Commun. ACM*, 59(2):64–73.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008.

Oriol Vinyals and Quoc Le. 2015. A neural conversational model. In *Proceedings of the 31st International Conference on Machine Learning, Deep Learning Workshop*, Lille, France.

Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.

S. Xie, R. Girshick, P. Dollr, Z. Tu, and K. He. 2017. Aggregated residual transformations for deep neural networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5987–5995.

Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057.

Yinfei Yang, Steve Yuan, Daniel Cer, Sheng-yi Kong, Noah Constant, Petr Pilar, Heming Ge, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. Learning semantic textual similarity from conversations. In *Proceedings of The Third Workshop on Representation Learning for NLP*, pages 164–174, Melbourne, Australia. Association for Computational Linguistics.

Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78.

Zhou Yu, Leah Nicolich-Henkin, Alan W Black, and Alexander Rudnicky. 2016. A wizard-of-oz study on a non-task-oriented dialog systems that reacts to user engagement. In *Proceedings of the 17th annual meeting of the Special Interest Group on Discourse and Dialogue*, pages 55–63.

Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pages 2204–2213, Melbourne, Australia. Association for Computational Linguistics.

# A  More Details of IGC Evaluations

In this section we describe a few choices we made and implementation details regarding the IGC human evaluation in the section regarding Transfer to the IGC Task.

**Multiple Traits**  In the IGC human evaluation setup from (Mostafazadeh et al., 2017), human annotators were shown eight choices when rating the quality of responses to questions: seven responses from various models, and one human response. To mirror this setup as closely as possible, we chose seven of our highest performing style traits to condition on to display in addition to the human response. We show the results of each trait in Table 4.

**Automatic Evaluation**  In (Mostafazadeh et al., 2017), the authors provide BLEU scores for their models in an attempt to evaluate their effectiveness via automated metrics. The authors note that the scores are very low, "as is characteristic for tasks with intrinsically diverse outputs." Additionally, it has been shown in (Shuster et al., 2019) that BLEU scores for image captioning retrieval models are generally far lower than those of generative models (as retrieval models do not optimize for such a metric), and yet human evaluations can show the complete opposite results. In fact, in that work retrieval models were shown to be superior to generative models in human evaluations, which is why we adopted them here. For these reasons we omit BLEU scores of our retrieval models on the IGC test set as uninteresting. We do however compare BLEU scores with our generative model in the main paper.

**Test Set Size**  The IGC test set provides the urls to all 2591 images for which (context, question, response) tuples were collected. We were only able to recover 2195 images from this initial set, as some of the urls provided are no longer associated with the corresponding images. Thus, our human evaluations are conducted on this subset.

| Style | Score |
|---|---|
| Neutral | 1.55 |
| Charming | 1.55 |
| Extravagant | 1.55 |
| Calm | 1.57 |
| Sweet | 1.58 |
| Spirited | 1.60 |
| Enthusiastic | 1.61 |
| Human | 2.55 |

Table 4: IGC Human Evaluation on responses from our TRANSRESNET MM-SUM model conditioned on various personalities. Responses were rated on a quality scale from 1 to 3, where 3 is the highest.

# B IMAGE-CHAT Human Annotation Setup



## Respond to a Comment on an Image
### Description

In this task, you will be shown 5 images, each of which has a comment about the image. The goal of this task is to write an engaging response to this comment as if you were continuing a dialog about the image.

### STEP 1

With each new photo, you will be given a **personality trait** that you will try to emulate in your response to the comment on the image. For example, you might be given "**snarky**" or "**sentimental**". The personality describes *YOU*, not the picture. It is *you* who is snarky or sentimental, not the contents of the image nor the original comment about the image.

### STEP 2

You will then be shown an image and a comment that goes with the image, for which you will write a response *in the context of your given personality trait*. Please make sure your response has at least **three words**. Note that these are *responses* to the comments on the image, and not simply image captions.

*Reminder - please do not write anything that involves any level of discrimination, racism, sexism and offensive religious/politics comments, otherwise the submission will be rejected.*

## Image

Someone wrote the following comment on this image:

**Peace and tranquility should be more abundant. This greenery evokes those feelings for me and I'm very thankful.**

Write your response as if you were: **Profound**

Figure 6: Instructions pane for crowdworkers when collecting the second round of dialogue.

## Continue a Dialog on an Image
### Description

In this task, you will imagine that you are speaking with your friend about 5 separate images. For each image, you will be shown "your" initial comment on the image, and your friend's response to the comment. The goal of this task is to write an engaging response to your friend as if you were continuing a dialog about the image.

### STEP 1

With each new photo, you will be given a **personality trait** that you will try to emulate in your response. For example, you might be given "**adventurous**". The personality describes *YOU*, not the picture. It is *you* who is adventurous, not the contents of the image.

### STEP 2

You will then be shown an image, "your" initial comment that goes with the image, and your friend's response. You will continue the dialog by responding to your friend's response *in the context of your given personality trait*. Please make sure your response has at least **three words**. Note that these are not simply image captions, but **engaging** responses.

*Reminder - please do not write anything that involves any level of discrimination, racism, sexism and offensive religious/politics comments, otherwise the submission will be rejected.*

## Image

YOU wrote the following comment on this image:
**I would be worried about getting cold out there.**

YOUR FRIEND responded:
**It's nice to just sit out in the snow and watch it fall. It's like being a whole different world.**

Write your response as if you were: **Maternal (Mother-like)**

Figure 7: Instructions pane for crowdworkers when collecting the third round of dialogue.

## C  IMAGE-CHAT Human Evaluation Setup



Figure 8: Instructions pane for crowdworkers when collecting the IMAGE-CHAT Evaluations.

## D  IGC Human Evaluation Setup



Figure 9: Instructions pane for crowdworkers when collecting the IGC Evaluations.

# E   Highly Rated Examples from IGC

| Image | IGC Round | Output |
|---|---|---|
|  | **Context** | These flowers are growing in my back yard. |
| | **Question** | What type of flowers are they? |
| | **Human Response** | I don't know but they are so pretty. |
| | **TransResNet MM-Sum** | I don't know but these flowers are gorgeous and look so bright! |
|  | **Context** | I bought this at a flea market. |
| | **Question** | What is this for? |
| | **Human Response** | I think it's a mortar and pestle. |
| | **TransResNet MM-Sum** | I'm not sure, but you could sell it for some cash! |
|  | **Context** | I'd love to visit this city. |
| | **Question** | Is this some city in Europe? |
| | **Human Response** | It definitely looks like it. |
| | **TransResNet MM-Sum** | Oh my goodness, yes! I love Amsterdam! |
|  | **Context** | These just flew over my house! |
| | **Question** | Wow, aren't those Blue Angels? |
| | **Human Response** | Yeah, they practice in this area for their shows. |
| | **TransResNet MM-Sum** | These are the Blue Angels, they are truly one of a kind |
|  | **Context** | My nephew's choir did so well yesterday! |
| | **Question** | That is great. You must be proud |
| | **Human Response** | I am so proud, because he's been working really hard. |
| | **TransResNet MM-Sum** | I am so proud of all that they accomplished. |
|  | **Context** | This would be a great place to have a picnic. |
| | **Question** | yes it would it seems very relaxing doesnt it? |
| | **Human Response** | Yes, and I would quite enjoy spending a day here. |
| | **TransResNet MM-Sum** | I think a picnic is more than a great possibility! |

Table 5: Highly rated examples from the IGC dataset test split where TRANSRESNET$_{RET}$ MM-Sum responses were rated the highest (score of 3) by human evaluators.
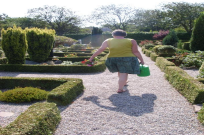
# F Poorly Rated Examples from IGC

| Image | IGC Round | Output |
|---|---|---|
|  | **Context** <br> **Question** <br><br> **Human Response** <br> **TransResNet MM-Sum** | These are Easter eggs we made for art class. <br> How did you make them? <br><br> We started with a glass base and went from there. <br> OHMYGOSH EASTER EGG HUNT! I want to hunt for Easter eggs too! |
|  | **Context** <br> **Question** <br><br> **Human Response** <br> **TransResNet MM-Sum** | I'm thinking of getting these shocks on my bike. <br> Are the shocks black in color? <br><br> Yes, except for the silver portion at the end. <br> Just got my new wheels! |
|  | **Context** <br> **Question** <br><br> **Human Response** <br> **TransResNet MM-Sum** | I cannot decide if this is milk or something else. <br> It looks like milk, but it could also be kefir. <br><br> What is kefir? <br> You can get all the protein you want form the sugar. |
|  | **Context** <br> **Question** <br><br> **Human Response** <br> **TransResNet MM-Sum** | This woman is visiting the grave of a loved one. <br> Where is the cemetery? <br><br> The cemetary is at the foot a beautiful quartz mountain. <br> Well she's just a walking, talking garden! |

Table 6: Low rated examples from the IGC dataset test split where TRANSRESNET$_{RET}$ MM-Sum responses were rated the lowest (score of 1) by human evaluators.

## G  Longer Chat Example



A: *TransResNet$_{RET}$ MM-Sum (Extraordinary)*     B: *Human (no style)*

A: It's the most beautiful mountain in the world!

B: It looks really cool! I wonder if you can ski on it

A: They are a mountain range I've climed many times!

B: You have climbed that? That seems really hard

A: I did it; i conquered this climb in only 7 hours and 10 minutes!

B: Wow, that's fast. I would ski down that but I would need a lift to take me up

A: But could you? Could you truly climb this?

B: I really don't think I could

A: Climbing a mountain can give one a special strength, you need to experience it

B: Maybe one day on a smaller mountain

A: It would take hard will and determination to scale that mighty peak

Figure 10: Long-form conversation with the model. The model is given a style here, while the human is not.

## H  Additional Ablation Results

| Modules | TRANSRESNET$_{GEN}$ (F1) | | | | TRANSRESNET$_{GEN}$ (BLEU-4) | | | |
|---|---|---|---|---|---|---|---|---|
| | Turn 1 | Turn 2 | Turn 3 | All | Turn 1 | Turn 2 | Turn 3 | All |
| Image Only | 10.8 | 11.0 | 11.2 | 11.0 | 1.1 | 1.3 | 1.2 | 1.2 |
| Style Only | 10.4 | 9.8 | 10.4 | 10.2 | 1.4 | 1.5 | 1.4 | 1.4 |
| Dialogue History Only | 9.9 | 11.4 | 12.2 | 11.2 | 1.0 | 1.9 | 1.8 | 1.6 |
| Style + Dialogue *(no image)* | 9.6 | **12.5** | **13.1** | 11.7 | 1.5 | **2.1** | 2.0 | 1.9 |
| Image + Dialogue *(no style)* | 10.7 | 11.1 | 11.7 | 11.2 | 1.1 | 1.7 | 1.6 | 1.5 |
| Image + Style *(no dialogue)* | 12.1 | 11.6 | 11.6 | 11.8 | 1.6 | 1.5 | 1.5 | 1.6 |
| Style + Dialogue + Image *(full model)* | **12.3** | **12.5** | **13.1** | **12.6** | **1.7** | **2.1** | **2.0** | **1.9** |

Table 7: Ablations on IMAGE-CHAT. We compare variants of our best TRANSRESNET generative model (ResNeXt-IG-3.5B image encoder) where we remove modalities: image, dialogue history and style conditioning, reporting F1 and BLEU-4 for generation for dialogue turns 1, 2 and 3 independently, as well as the average over all turns.