

Validated Image Caption Rating (VICR) Scale, Dataset, and Model

Anonymous ACL submission

Abstract

Assessing the quality of an image caption is a complex task. We propose a new image caption rating system that consists of (1) a robust rating scale that is consistent, teachable, and externally validated, (2) an engaging and scalable data generation approach for the task, (3) a high-quality dataset, and (4) an effective image caption rating predictor. Using contemporary approaches from psychometrics we demonstrate that the proposed scale and rater training routine can support high quality annotation efforts for the task. We introduce two new datasets (one original and another derived) for the task. Our reference-free and multi-level rating predictor performance is on par with state-of-the-art approaches.

1 Introduction

We present a novel image caption rating (ICR) framework that consists of (1) externally validated rating scale, (2) a scalable data generation tool, and (3) high-quality dataset, and (4) an effective ICR prediction model. The problem of image caption quality estimation has received substantial attention in recent years, underscoring the increasing need for reliable solutions (Jiang et al., 2019; Lee et al., 2021; Hessel et al., 2021; Lee et al., 2020; Wang et al., 2018). Existing datasets for the task of image caption rating are generated using the traditional approach of human-driven data annotation efforts, and typically use ad hoc rating scales (Levinboim et al., 2019; Hodosh et al., 2013; Vedantam et al., 2015). All these datasets have been tremendously valuable in advancing the field and have been used extensively (Hessel et al., 2021; Lee et al., 2021). However, several of the datasets suffer

from high skew in ratings and mixed quality annotations. Our work seeks to improve the rigor, quality, and scalability of ICR datasets and data generation process, and provides a robust scoring instrument that is informed by contemporary approaches to measurement – specifically, Item-Response Modeling.

For the problem of image caption rating estimation, the main difference in existing approaches stems from their ability to estimate the rating in the presence or absence of reference caption(s). BLEU (Papineni et al., 2002), METEOR (Denkowski & Lavie, 2014), ROUGE (Lin, 2004), CIDEr (Vedantam et al., 2015) and SPICE (Anderson et al., 2016), BERTscore (Zhang et al., 2019) and ViLBERTscore (Lee et al., 2020) belong to the former category where reference captions are essential, while Visual Semantic Embedding Plus Plus (VSEPP) (Faghri et al., 2017), CLIPScore (Hessel et al., 2021) and approaches proposed by Cui and colleagues (2018) and Levinboim and colleagues (2019) can operate without reference captions. The ability of these approaches to assess caption quality without requiring reference captions has led to rapid progress on this problem. However, the rating granularity employed by these approaches has been restricted to simple binary scale (*good* or *bad* caption). In our work we seek to lift this restriction by employing a 5-level rating scale that can model different aspects of quality in the context of image captions (e.g., correctness, completeness, and inclusion of local and global context), while also retaining the benefits of reference-free rating approach. Although a more detailed scale can offer higher rating capacity, it can also increase the complexity of the rating task; potentially making the task more subjective and tedious. To tackle this downside, we propose a two-pronged solution during data generation: (1) rigorous

81 training procedure with in-built quality control,
82 and (2) gamification.

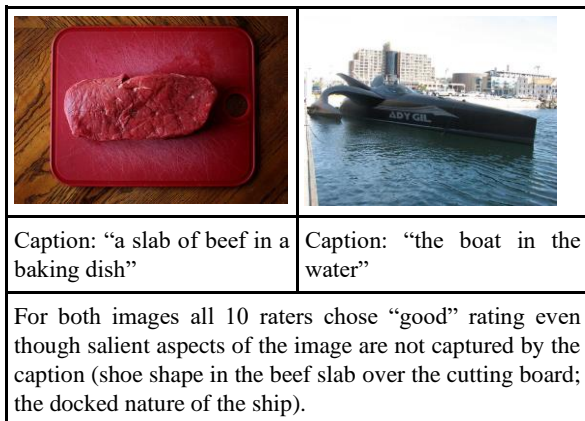
83 Altogether we refer to our work as a Validated
84 Image Caption Rating (VICR) framework; and
85 our specific contributions are introduction of the
86 VICR Scale, VICR Game, VICR Dataset, and
87 VICR Model. The rest of the paper is organized as
88 follows. The next section provides the context and
89 the rich prior work that we build our work on.
90 Section 3 describes the VICR system in detail,
91 followed by Results and Analysis in Section 4,
92 and the conclusions we draw from this work in
93 Section 5.

94 2 Related Work

95 2.1 ICR Scale and Datasets

96 Google Image Caption (GIC) Dataset (Levinboim
97 et al., 2019) and Flickr8k-Expert (Hodosh et al.,
98 2013) are the two widely used large image caption
99 datasets that also include ratings. GIC dataset has
100 600K image-caption ratings. For each
101 image/caption pair, 8-10 binary ratings were
102 collected. The ratio of good ratings to total ratings
103 is used as image caption quality score (range: [0,
104 1]). As is common with binary scales, it does not
105 have the capacity to handle incomplete or partially
106 correct captions. Figure 1 includes two illustrative
107 examples.

108 Similarly, Conceptual Caption Challenge
109 human evaluation studies on the T2 test dataset¹
110 contains 5000 image and caption pairs and the
111 human ratings are collected in the same manner of
112 GIC; each pair has the total rating counts and the
113 good ratings counts.



114 Figure 1: Two examples from Google Image Caption Dataset
115 illustrating the limitation of binary scale.

116 Flickr8k-Expert, a subset of Flickr8k dataset
117 (8,000 images and 5 captions per image), has
118 5,822 captions across 1,000 images where each
119 caption has received 3+ ratings from human
120 annotators (21 college students). The rating scale
121 used for Flickr8k dataset consisted of 4 levels
122 (Table 1). The complexity of the ICR task
123 combined with the underspecified rating scale and
124 human error lead to fairly low inter-rater
125 agreement. The rating distribution is also heavily
126 skewed toward levels 1 & 2, indicating overall
127 lower caption quality.

r	Meaning
4	Describes the image without any errors.
3	Describes the image with minor errors.
2	Is somewhat related to the image.
1	Is unrelated to the image.

128 Table 1: Flickr8k-Expert ratings and meanings.

129 The CapEval1K dataset (Lee et al., 2021) is rich
130 for containing fluency, relevance, and
131 descriptiveness rates per caption, but has a rather
132 small size (1,000 captions for 250 images). The
133 PASCAL50s dataset (Vedantam et al., 2015) has
134 50 reference captions per image for 1000 images,
135 but the ratings are not in numeric scale.

136 2.2 Reference-free ICR Estimators

137 VSEPP (Faghri et al., 2017) and CLIPScore
138 (Hessel et al., 2021) are multimodal, visual-
139 linguistic models that use cosine similarity to
140 measure the distance between an image
141 embedding and text embedding in a shared visual-
142 semantic embedding space. Unfortunately, while
143 the cosine similarity does a good job on
144 approximation of the similarity of the vectors in
145 the shared visual-linguistic semantic space, fine
146 tuning or manipulation of the similarity of the
147 image and language features remains difficult.

148 Cui and colleagues (2018) created a deep
149 learning method for determining if a caption for
150 an image was human-written or machine
151 generated. However, this is a binary classifier and
152 is not sufficient for diverse use cases.

153 Levinboim and colleagues (2019) developed an
154 image-caption Quality Estimation (QE) model by
155 training a deep learning model on the GIC dataset.
156 The model inherits the limitations from the dataset
157 discussed in Figure 1.

¹ <https://www.conceptualcaptions.com/winners-and-data>

158 Lee and colleagues (Lee et al., 2021) developed
159 Unreferenced Metric for Image Captioning
160 (UMIC) using UNITER (Chen et al., 2020) via
161 contrastive learning, a process where the model is
162 trained to compare and discriminate the ground-
163 truth captions and diverse synthetic negative
164 samples. Jiang and colleagues (Jiang et al., 2019)
165 developed TIGEr (Text-to-Image Grounding for
166 Image Caption Evaluation) by improving the
167 mapping of the image and the caption pair into
168 carefully grounded vector spaces. These
169 approaches improved consistency with human
170 judgements over prior metrics, but still did not
171 exceed .5 Kendall τ scores on the Flickr8K expert
172 data set.

173 2.3 Integrative Inferential Reasoning (IIR)

174 The importance of a robust rating scale for the
175 ICR task cannot be overstated. Having a
176 theoretical foundation can ensure that a rating
177 scale yields explicit and trainable scoring guides
178 that lead to reliable ratings. Based on industry-
179 accepted image description guidelines, the
180 context in the image must be included in the
181 caption, and thus is an important aspect of caption
182 quality (Rai et al. 2010)². Contextual integration
183 is the backbone of Integrative-Inferential
184 Reasoning (IIR) (Blum et al., 2020). IIR is a
185 cognitive framework that structures context
186 integration in text- and image-based narratives.
187 IIR’s scaled definitions of context and inference
188 offers a roadmap for training humans (and by
189 extension, machines) on how to rate image
190 caption quality based on these characteristics. In
191 its modern form, IIR is a novel approach to
192 capturing combined notions of context and
193 inference; however, the theory stems from older
194 notions of local (e.g., propositional or literal) and
195 global (e.g. schematically or culturally relevant)
196 coherence, which has been investigated in literacy
197 (Graesser et al. 1994; Language and Reading
198 Research Consort...), cognition (Frith and Happé
199 1994; Van der Hallen et al. 2015), neurodiverse
200 populations such as autism (Happé & Frith, 2006;
201 Nuske & Bavin, 2011); and the schema of
202 Question-Answer Relations (Pearson and Johnson
203 1978; Raphael and Au 2005). With its historical
204 theoretical grounding, IIR offers an exciting
205 foundation for developing a new kind of image
206 rating scale.

207 3 Methods

208 The old adage “A picture is worth a thousand
209 words.” perfectly captures the challenge faced by
210 image caption raters (humans and machines). An
211 image can convey layers of nuanced information,
212 while a short textual caption has a very limited
213 information bandwidth. Naturally, assessing the
214 quality of image captions is an inherently tricky
215 task. To tackle this complex problem, we start by
216 unpacking the ICR pipeline. The first source of
217 error is often the rating scale itself. The errors
218 caused by an ill-defined scale propagate
219 downstream and compound. The second source of
220 error is typically humans who are doing the
221 tedious and complicated task of rating the
222 captions. We coalesce these observations to
223 define two key objectives for our work:

224 Objective #1: Design and develop a reliable and
225 scalable data generation approach for the task of
226 image-caption rating. To achieve this objective,
227 we innovate along three areas: (1) Develop a
228 rating scale that accurately captures the nuances
229 and aspects of image caption quality (VICR
230 Scale); (2) Develop an engaging tool (VICR
231 Game) to facilitate high-quality data generation
232 from human raters (VICR Dataset); and (3)
233 Assess the ability of human raters to effectively
234 use this data-generation approach.

235 Objective #2: Develop a novel image-caption
236 rating model (VICR Model) that employs the
237 outcomes from objective #1.

238 Together, these objectives provide a robust,
239 high-quality, and scalable image-caption rating
240 system which is described next.

241 3.1 VICR Scale: Relating IIR to Image Captions

242 Integrative Inferential Reasoning (IIR) is a
243 theoretical construct, developed using the BEAR
244 assessment system (Wilson, 2005). We applied
245 IIR as a theoretical foundation (IC-IIR) to inform
246 the development of VICR Scale. This 5-level
247 scale captures nuances in caption accuracy,
248 completeness, inferential, and contextual
249 information as listed in Table 2.

250 To evaluate the efficacy of the VICR Scale at
251 training raters and at producing consistent ratings
252 across raters, we employed measures of rater
253 competency using the following approach. Rater
254 competency was represented by “items” (the
255 image-caption pairs) and “responses”

² <https://dcmp.org/learn/descriptionkey>

(participants' ratings). We used a 5x5 factorial items design: five images were used, and each image was paired five times, with captions representing each of the five levels of the VICR Scale. Each rater was assigned a *competency score* based on the degree of agreement between their ratings and expert ratings as follows: *Exact Agreement* (participant and expert ratings are equal) received a score of 2; *Adjacent Agreement*, (participant and expert ratings differ by 1) received a score of 1, and *Lack of Agreement* (participant and expert ratings differ by more than 1) received a score of 0. The cumulative score over all 25 image-caption pairs was computed for each rater and analyzed using the Partial Credit Model (PCM) (Masters, 1988; Masters, 2016). The PCM is a Rasch-family measurement model that is used to place items and participants on the same scale and evaluate the quality of an obtained measurement. We also used a Latent Regression (Wilson & De Boeck, 2004) to regress rater competency on their tutorial score obtained during training. Results of these analyses are shared in Section 4.1

<i>r</i>	Meaning
5	Objects, a general scene, and actions are correctly identified if present in the image. The caption describes what is seen and where things are in space.
4	Objects and/or a general scene and/or an action are correctly identified but not every element is completely identified. The caption describes what is seen and where things are in space. There is no interpretation of an event.
3	Relevant objects are correctly identified. The caption describes what is seen but not where objects are in space. There is no description of the overall setting and no interpretation of an event.
2	Objects are partially correctly identified with some errors, but the caption is accurate enough to give an idea of what is happening in the image. The caption identifies most of the objects but might not identify everything. There is no interpretation of what anything means.
1	Objects are incorrectly identified. The caption gives the wrong idea about what is happening in the image.

Table 2: VICR Scale: Ratings and Meanings.

3.2 VICR Dataset Generation: Image Caption Rating Game

To facilitate generation of high-quality and substantially sized data we focus on human rater training and engagement in this phase of the VICR system. Rater training is essential for any data annotation effort, but it is especially important in our project due to the detailed nature of the VICR

Scale. A 5-level scale with each level capturing multiple aspects of caption quality is non-trivial to apply for most humans.

Rater Training: The training is conducted online through a web application that starts by showing the VICR Scale to the human rater. When ready the rater proceeds to a test round where an image and caption pair is displayed, and the rater has to choose the most appropriate rating level from the VICR Scale for the pair. This is repeated for 20 image-caption pairs. The accuracy of the rater's selections is computed using the ground-truth ratings. Raters with accuracy of 0.5 or higher are cleared for data generation, and others are required to redo the training until minimum accuracy is met. The reasoning behind the chosen accuracy threshold is explained in Section 4.1.

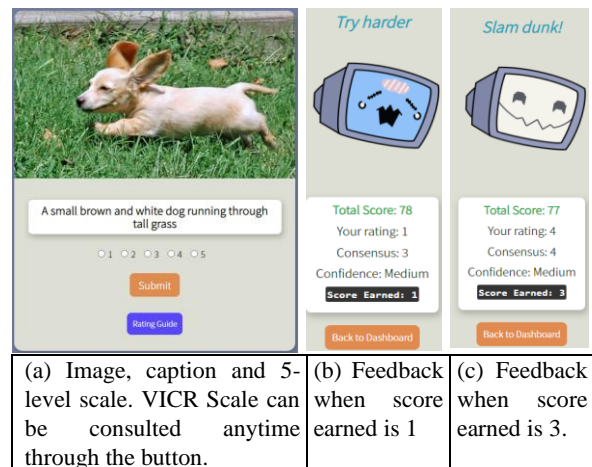


Figure 2: VICR Image Caption Rating Game.

VICR Game: To promote rater engagement we frame the annotation task as a single-player asynchronous competitive game; following on the path of image labeling ESP game (Von Ahn & Dabbish, 2004). The web-based VICR Game is designed to provide a similar user experience as the training phase – an image-caption pair is displayed, and the player selects the appropriate rating from the 5-level VICR Scale (Fig. 2a). After rating submission, the player receives feedback that compares their selection with those of the other players so far (Fig. 2b and 2c). Specifically, a *consensus score* (*con*), which is the rounded average of all the previous ratings for that image-caption pair so far, is displayed. Player earns points, *p*, for the rating submission using the following formula:

$$p = \max(4 - [2 |x - con| / c], -1)$$

x = player's rating selection
 con = consensus score
 $c = 1 + (1 + (n - 1) \sigma^2 / V_{max}) / n$

where σ^2 is the variance of the previous ratings and $V_{\max} = 4$, the largest possible variance of the previous ratings, therefore $p \in [-1, 3]$. This formulation models two intuitions: 1. the assigned points should be inversely proportional to the difference between the player’s rating and the consensus score, and 2. the assigned points should be proportional to the degree of agreement among the ratings so far. Together, these intuitions ensure low points for scenarios where agreement among prior ratings is high and the current rating exhibits a large difference from the average. In contrast, if the level of agreement is low, the points decrease only gradually as the difference from the average increases. This is supported by the coefficient c in the formulas above. This coefficient, called confidence, will be between 1 and 2, where 1 represents perfect confidence, and 2 represents the least possible confidence. It is used to modify the distance from the consensus at which various points are awarded.

This formulation provides the ability to penalize ratings that deviate substantially, $p \in [-1, 3]$. We seed the target ratings initially with ratings from VSEPP. For the purposes of calculating mean, variance, and the level of consensus multiplier, we include this initial rating twice, i.e., as two agreeing data points. Once a participant gives their rating for the image-caption pair, this rating replaces one of the two initial ratings, and once a second player has rated the pair, the second initial rating is replaced as well, so that the average and level of consensus are now purely based on the two human ratings, and from then on, the human ratings accumulate as normal.

3.3 VICR Model: Image Caption Rater

We propose a multi-level reference-free image-caption rating predictor, VICR Model (Fig. 3). The rating predictor starts by converting the image and the caption into image and language embeddings, respectively. Preliminary experiments with various image and language embeddings, demonstrated ViLBERT co-fusion embeddings as being the most effective for our model. We use the pooled text and image embeddings of the final hidden layer in ViLBERT and concatenate these into a 2048-dimensional vector as input to our network. For the regressor model, a two-hidden-layer fully-connected neural network with 512 neurons on the first layer, followed by ReLU activation, and 256 neurons on

the second layer, followed by another ReLU activation, with a single neuron with linear activation as the output layer. We used 80% dropout on both hidden layers.

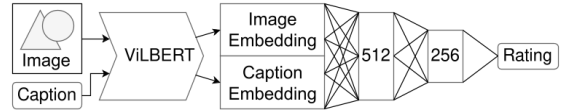


Figure 3: Schematic diagram of the VICR model architecture.

4 Experiments and Results

Our evaluation methodology for the VICR system consists of two user studies for VICR Scale and VICR Game, respectively, a comparative analysis of VICR Dataset, and an empirical evaluation for the VICR Model.

4.1 VICR Scale: Initial Validation of Image-Caption Quality and Rater Consistency

The goal of User Study 1 was to evaluate the efficacy of a VICR Scale and the corresponding training module at generating high-quality ICR data. For this study, 132 fully anonymized participants (college students at a 4-year public university) were recruited. The participants started by undergoing the Rater Training routine (Section 3.2), and the ones who cleared the quality threshold were then prompted to rate 25 items (5 images and 5 captions per image in random order) using the 5-level VICR Scale.

The collected data was then analyzed as per the methodology described in Section 3.1. Specifically, we employed Wright Map – an analytical tool that allows us to place human raters and image caption pairs (i.e., items) visually on the same scale (Embretson 1996; Stachl and Baranger 2020; Blum et al. 2020; Wilson 2005), (Brondfield et al. 2021; Blum 2019) to analyze rater competency when using VICR Scale, and a latent regression analysis understand the strength of an explanatory relationship between our training module and rater competency.

4.1.1 Wright Map

The PCM (section 3.1) uses human rater proficiency and image-caption pair difficulty estimates, and the error associated with them to generate Wright Maps (also known as item-person maps, Fig. 4). The first column displays results from a latent regression (section 4.1.3); the second column shows a histogram of rater

425 competency scores (Section 3.1) in logits (column
426 three).

427 The key observation from this analysis comes
428 from the right side of the Wright Map which
429 reports participants' levels of VICR Competency.
430 Two cumulative thresholds per item (image-
431 caption pair) are represented on the right side (25
432 columns, one for each item along the horizontal
433 axis). The first threshold, marked in yellow,
434 represents where a respondent would be equally
435 likely to score 0 (*LA: Lack of Agreement*) vs. 1 or
436 2 (*AA: Adjacent* or *EA: Exact Agreement*); the
437 second threshold, marked in red, represents where
438 a respondent would be equally likely to score 0 or
439 1 (*LA* or *AA*) vs. 2 (*EA*). This Wright Map shows
440 that in our data most of the items' second
441 thresholds are above most of the items' first
442 thresholds, which represents internal validity of
443 raters' competency in using the VICR Scale
444 (no/minimal confusion).

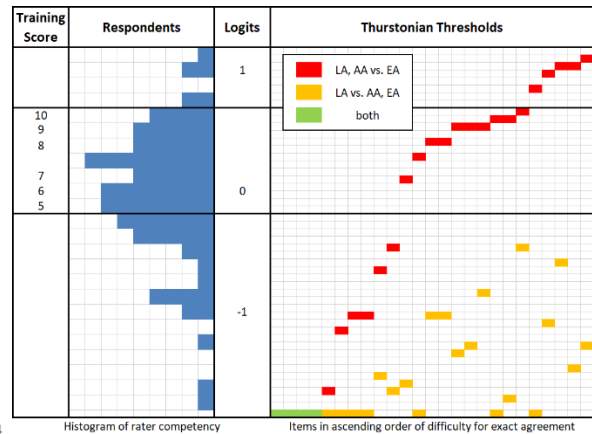
445 4.1.2 Latent Regression

446 Latent regression was used to explain the
447 relationship between VICR training routine and
448 rater competency. The regression coefficient for
449 *training score* was 0.17 (stderr 0.03) which is
450 significant at the .05 level. This means that each
451 additional tutorial item that the respondent rated
452 correctly is associated with a mean increase in
453 rater competency of 0.17 logits. The significance
454 of this output is seen in the leftmost column of the
455 Wright Map, which shows the predicted mean
456 VICR competency score for each possible tutorial
457 score from 5 to 10. At the training score of 5, the
458 predicted mean rater competency (-0.13) is well
459 above all the first thresholds, and above the
460 second threshold for 10 of the 25 items. This
461 suggests that, on average, even a respondent with
462 training score of 5 (weakest rater) is very likely to
463 demonstrate at least *AA* on all items and has more
464 than 50% chance of *EA* on 10 of the items (and
465 less than a 50% chance for *EA* on the other 15
466 items). The significant finding from this analysis
467 is that it provides external validity of raters'
468 competency in using the VICR Scale. This
469 analysis also informs the choice of minimum
470 threshold (training score of 5 = 0.5 accuracy) used
471 for rater selection during training routine (Section
472 3.2).

473 4.1.3 Is the VICR Scale teachable?

474 The short answer is, yes. The positive and
475 significant regression coefficient of 0.17 indicates

476 that respondents who were more successfully
477 trained in using the VICR Scale were better able
478 to reach *EA* on more image-caption pairs. This is
479 also visible in the Wright Map, where raters with
480 an increasing training score of 5 to 10 are more
481 likely to reach *EA* on up to 80% of the items, and
482 most likely to have *AA* on 100% of the items with
483 a training score of 10.



484 Histogram of rater competency
485 Figure 4: Wright Map augmented with predicted means from
486 latent regression

487 4.1.4 Can people use the VICR Scale 488 consistently and well?

489 Based on respondent frequency and locations on
490 the Wright Map, most respondents (103 out of
491 132, i.e., 78%) are above all the first thresholds;
492 as noted above, respondents at that level can
493 reliably achieve high agreement (*EA* or *AA*)
494 (median was 0.01 logits). Respondents who were
495 more successfully trained (higher rater
496 competency score) tend to achieve *EA* on more of
497 the items; respondents with a very high training
498 score (10) tend to achieve *EA* on 80% of the
499 image-caption pairs.

500 4.2 VICR Game and Datasets

501 The goal of User Study 2 was to employ the VICR
502 Game to generate a new dataset for the ICR task.
503 We also created the Combined Dataset that
504 consists of the new VICR Dataset and the
505 Flickr8k-Expert dataset. (The new datasets will be
506 freely available to the research community.)

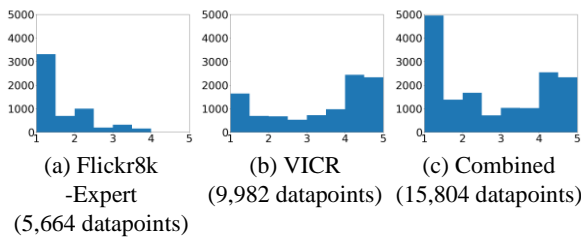
507 **VICR Dataset:** A collection of image-caption
508 pairs was assembled for this user study as follows:
509 8,990 distinct images were chosen at random from
510 MS-COCO 2014 Validation subset (Lin et al.,
511 2014), The captions were selected from 4 sources:
512 (1) the original MS-COCO caption, (2) generated
513 using the Pythia framework (Jiang et al., 2018),
514 (3) generated using the GLACNet model (Kim et
515 al., 2018), and (4) mismatched captions from

516 other images. Leading to 9,982 image-caption
517 pairs.

518 As part of the user study 72 participants (college
519 students) played the VICR Game to rate image-
520 caption pairs from the above collection. On
521 average, the participants played for 102 minutes,
522 earning about \$15 per hour. The participants took
523 about 10 seconds on average to rate an image
524 caption pair. By the end of the user study, a total
525 of 48,174 ratings were collected, so that each of
526 the 9,982 image-caption pairs had at least 4 and at
527 most 7 ratings.

528 **Combined Dataset:** We also made a Combined
529 dataset composed of Flickr8k-Expert and VICR to
530 create a bigger data set (15,804 image-caption
531 pairs with ratings). When consolidating the two
532 datasets, we mapped Flickr8k-Expert’s 4-level
533 Scale to the first 4 levels of VICR Scale, since the
534 meanings ratings to 1 to 5 but instead kept them
535 as their original scale of 1 to 4 since their 1 to 4
536 map to our 1 to 4 relatively well with 5 being an
537 extra level in our dataset. The 5th level is
538 essentially not represented in the Flickr8k-Expert
539 rating scale.

540 **Comparative Analysis:** The rating distribution
541 of the Flickr8k-Expert, VICR, and Combined
542 Datasets are illustrated in Fig. 5. For each image-
543 caption pair, the rounded average of all available
544 ratings for that pair is used as the single value
545 rating for the pair.



546 Figure 5: Datasets: Rating Distributions

547 Figure 5 demonstrates that the new VICR
548 Dataset is less skewed in its rating distribution
549 than the Flickr8k-Expert dataset. It does however
550 exhibit bimodal distribution indicating a larger
551 proportion of low- and high-quality captions than
552 average quality captions. The Combined Dataset
553 naturally embodies the properties of both the
554 source datasets.

555 4.3 VICR Model

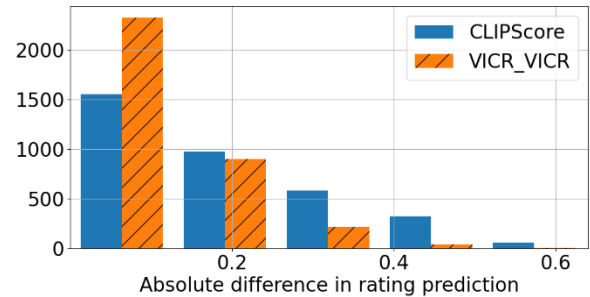
556 We evaluate the effectiveness of our multi-level
557 reference-free image-caption rating predictor,
558 VICR Model, with two empirical experiments.
559 **Experiment 1:** Table 3 provides results for the

560 first experiment where VICR_{VICR} (VICR model
561 trained with VICR Dataset) performance is
562 compared to Reference-based and Reference-free
563 approaches. (We used Adam optimization,
564 minimized on MSE, for 4,000 epochs.)

565 It is not surprising that the Reference-based
566 approaches exhibit higher performance. Within
567 the Reference-free category, CLIPScore provides
568 the highest performance, with VICR_{VICR} being a
569 close second. VICR Model shows good
570 generalizability – despite being trained on VICR
571 Dataset with 5-level rating scale, the predictor
572 provides competitive performance on Flickr8k-
573 Expert dataset with 4-level scale.

Reference-based Approaches	τ_C
BLEU-1	36.3
BLEU-4	33.1
METEOR	43.6
ROUGE	38.1
CIDER	43.7
SPICE	45.9
RefCLIPScore	52.7
ViLBERTScore-F	54.2*
Yi et al.	48.1*
Reference-free Approaches	τ_C
CLIPScore	51.5
UMIC-c	43.1*
TIGER	49.3*
VSEPP	48.1
VisualEntailment	44.6
VICR _{VICR}	50.9

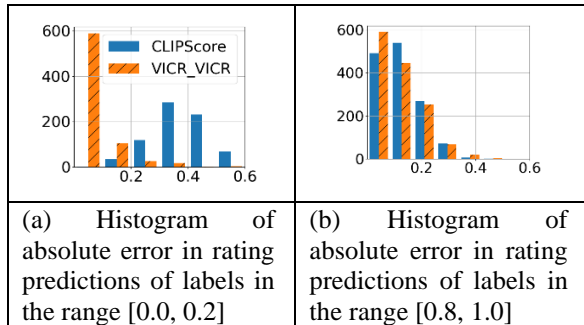
574 Table 3: Kendall τ correlation with ground truth
575 ratings on Flickr8k-Expert dataset for various metrics
576 and predictors. We recreated all the listed results
577 except for the ones with * which are directly from the
578 respective papers. We used “method A” in aggregation
579 (Hessel et al., 2021) and τ_C to be consistent with prior
580 work.



581 Figure 6: Error analysis: Histogram of absolute error in
582 rating predictions.
583

584 Figure 6 provides a deeper analysis of rating
585 predictions by computing the absolute difference
586 between the predicted rating and the ground truth
587 rating for image-caption pairs in the Flickr8k-

588 Expert dataset. The ratings have been normalized
 589 into the range [0, 1]. The x-axis specifies the
 590 absolute error in rating prediction. Notice that the
 591 0-error bar for VICR_{VICR} is substantially higher
 592 than that of CLIPScore. Overall, the histogram
 593 distribution for VICR_{VICR} is heavily skewed to the
 594 left, indicating lower incidence and magnitude of
 595 prediction errors.



596 Figure 7: Error analysis: Histogram of absolute error in
 597 rating predictions on sub-ranges of labels.

598 We further analyzed absolute errors in rating
 599 predictions on sub-ranges of ground truth ratings
 600 (2 shown in Fig. 7), showing higher performance
 601 over CLIPScore in almost all ranges.

Caption: A group of elephants by some buildings on the water.	Caption: a number of baseball players in a field
CLIPScore: 0.43 VICR _{VICR} : 0.0 human rating avg: 0.0	CLIPScore: 0.44 VICR _{VICR} : 0.06 human rating avg: 0.06
(a)	(b)
Caption: A woman standing on a balcony in front of an elephant float.	Caption: A woman preparing to hit a tennis ball while a man watches.
CLIPScore: 0.69 VICR _{VICR} : 0.81 human rating avg: 0.81	CLIPScore: 0.54 VICR _{VICR} : 0.02 human rating avg: 0.43
(c)	(d)

602 Table 4: Samples of image, caption and metrics.

603 The examples in Table 4 (from VICR Test set)
 604 illuminate this further. For easier comparison, the
 605 ratings are all normalized to lie in the range [0, 1].
 606 There are cases where VICR scores align
 607 perfectly or very closely with human ratings (e.g.,

608 Table 5-a, b, c). There are also cases where
 609 VICR_{VICR} seems even more accurate than human
 610 ratings (e.g., examples d).

611 **Experiment 2:** The second experiment studies
 612 the ability of the three datasets (Flickr8k-Expert,
 613 VICR, and Combined) at training an effective
 614 rating predictor with VICR Model. Each dataset
 615 was split into 64% training, 16% validation, and
 616 20% test for this experiment. Three models,
 617 VICR_{Flickr8k}, VICR_{VICR}, VICR_{Combined}, were trained
 618 on the respective Training sets.

Reference-free Approaches	τ_C Flick8K-Expert	τ_C VICR	τ_C Combined
VICR _{Flickr8k}	52.1*	61.3*	71.2*
VICR _{VICR}	50.6*	66.4*	73.4*
VICR _{Combined}	53.2*	66.0*	75.5*
CLIPScore	(51.2 ¹)	51.5	66.3
VSEPP	48.1	62.3	66.5
VisualEntailment	44.6	54.6	65.0

619 Table 5: Kendall τ correlation with ground truth
 620 ratings for reference-free approaches, *Calculated on
 621 Test set of each dataset. ¹Reported in (Hessel et al.,
 622 2021)

623 The top half of the Table 5 reports performance
 624 of the VICR models with the three Test sets. All
 625 three models perform better on VICR and
 626 Combined Datasets. This trend is also seen with
 627 the other Reference-free approaches (lower half
 628 of Table 5). This suggests that the VICR and
 629 Combined are more reliable ICR datasets than
 630 Flickr8k-Expert.

631 5 Conclusions

632 In this work we introduced an image caption
 633 rating system that consists of a new rating scale,
 634 an engaging data generation approach, a high-
 635 quality dataset, and a rating prediction model. A
 636 multi-level rating scale that captures various
 637 nuances of caption quality can be difficult to
 638 apply. Our user studies suggest that a well-defined
 639 scale along with methodical training and a game-
 640 based data generation setup can provide the right
 641 balance of data quality and quantity. The new
 642 dataset generated by this approach when
 643 employed to train our reference-free rating
 644 predictor provides one of the highest
 645 performances for the image caption rating task.
 646 However, we have not yet explored how to utilize
 647 specific context in rating scale or how robustly it
 648 performs on objects it is not trained on. We also
 649 have not explored potential risks of biases in
 650 image caption ratings.

References

- Anderson, P., Fernando, B., Johnson, M., & Gould, S. (2016, October). Spice: Semantic propositional image caption evaluation. In *European conference on computer vision* (pp. 382-398). Springer, Cham.
- Blum, A. M. (2019). *Deficit or Difference? Assessing Narrative Comprehension in Autistic and Typically Developing Individuals: Comic vs. Text*. University of California, Berkeley.
- Blum, A. M., Mason, J. M., Kim, J., & Pearson, P. D. (2020). Modeling question-answer relations: the development of the integrative inferential reasoning comic assessment. *Reading and Writing*. <https://doi.org/10.1007/s11145-020-10026-4>.
- Bronfield, S., Seol, A., Hyland, K., Teherani, A., & Hsu, G. (2019). Integrating Concept Maps into a Medical Student Oncology Curriculum. *Journal of Cancer Education*, 1-7.
- Chen, Yen-Chun et al. "UNITER: UNiversal Image-TExt Representation Learning." Lecture Notes in Computer Science (2020): 104–120. Crossref. Web.
- Cui, Y., Yang, G., Veit, A., Huang, X., & Belongie, S. (2018). Learning to evaluate image captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 5804-5812).
- Denkowski, M., & Lavie, A. (2014, June). Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the ninth workshop on statistical machine translation* (pp. 376-380).
- Embretson, S. E. (1996). Item response theory models and spurious interaction effects in factorial ANOVA designs. *Applied Psychological Measurement*, 20(3), 201-212.
- Faghri, F., Fleet, D. J., Kiros, J. R., & Fidler, S. (2017). Vse++: Improving visual-semantic embeddings with hard negatives. *arXiv preprint arXiv:1707.05612*.
- Frith, U., & Happé, F. (1994). Autism: beyond "theory of mind." *Cognition*, 50(1), 115–132.
- Graesser, A. C., Singer, M., & Trabasso, T. (1994). Constructing inferences during narrative text comprehension. *Psychological Review*, 101(3), 371–395.
- Happé, F., & Frith, U. (2006). The weak coherence account: detail-focused cognitive style in autism spectrum disorders. *Journal of Autism and Developmental Disorders*, 36(1), 5–25.
- Hessel, J., Holtzman, A., Forbes, M., Bras, R. L., & Choi, Y. (2021). CLIPScore: A Reference-free Evaluation Metric for Image Captioning. *arXiv preprint arXiv:2104.08718*.
- Hodosh, M., Young, P., & Hockenmaier, J. (2013). Framing image description as a ranking task: Data, models and evaluation metrics. *Journal of Artificial Intelligence Research*, 47, 853-899.
- Jiang, M., Huang, Q., Zhang, L., Wang, X., Zhang, P., Gan, Z., ... & Gao, J. (2019). TIGER: text-to-image grounding for image caption evaluation. *arXiv preprint arXiv:1909.02050*.
- Jiang, Y., Natarajan, V., Chen, X., Rohrbach, M., Batra, D., & Parikh, D. (2018). Pythia v0.1: the winning entry to the vqa challenge 2018. *arXiv preprint arXiv:1807.09956*.
- Kim, T., Heo, M. O., Son, S., Park, K. W., & Zhang, B. T. (2018). Glac net: Glocal attention cascading networks for multi-image cued story generation. *arXiv preprint arXiv:1805.10973*.
- Language and Reading Research Consortium, & Muijselaar, M. M. L. (2018). The dimensionality of inference making: Are local and global inferences distinguishable? *Scientific Studies of Reading: The Official Journal of the Society for the Scientific Study of Reading*, 22(2), 117–136.
- Lee, H., Yoon, S., Dernoncourt, F., Kim, D. S., Bui, T., & Jung, K. (2020, November). ViLBERTScore: Evaluating Image Caption Using Vision-and-Language BERT. In *Proceedings of the First Workshop on Evaluation and Comparison of NLP Systems* (pp. 34-39).
- Lee, H., Yoon, S., Dernoncourt, F., Bui, T., & Jung, K. (2021). UMIC: An unreferenced metric for image captioning via contrastive learning. *arXiv preprint arXiv:2106.14019*.
- Levinboim, T., Thapliyal, A., Sharma, P., & Soricut, R. (2019). Quality estimation for image captions based on large-scale human evaluations. *arXiv preprint arXiv:1909.03396*.
- Lin, C. Y. (2004, July). Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out* (pp. 74-81).
- Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., ... & Zitnick, C. L. (2014, September). Microsoft coco: Common objects in context. In *European conference on computer vision* (pp. 740-755). Springer, Cham.
- Masters, G. N. (1988). The Analysis of Partial Credit Scoring. *Applied Measurement in Education*, 1(4), 279–297.
- Masters, G. N. (2016). Partial credit model. In W. J. van der Linden (Ed.), *Handbook of Item Response Theory* (Vol. 1, pp. 109–126). Taylor and Francis.
- Nuske, H. J., & Bavin, E. L. (2011). Narrative comprehension in 4-7-year-old children with autism: Testing the Weak Central Coherence account. *International Journal of Language & Communication Disorders*, 46(1), 108–119.
- Papineni, K., Roukos, S., Ward, T., & Zhu, W. J. (2002, July). Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics* (pp. 311-318).
- Pearson, P. D., & Johnson, D. D. (1978). Questions. In *Teaching reading comprehension* (pp. 153–178). Holt, Rinehart and Winston.
- Raphael, T. E., & Au, K. H. (2005). QAR: Enhancing comprehension and test taking across grades and content areas. *The Reading Teacher*, 59(3), 206–221.
- Stachl, C. N., & Baranger, A. M. (2020). Sense of belonging within the graduate community of a research-focused STEM department: Quantitative assessment using a visual narrative and item response theory. *PloS one*, 15(5), e0233431.
- Van der Hallen, R., Evers, K., Brewaeys, K., Van den Noortgate, W., & Wagemans, J. (2015). Global processing takes time: A meta-analysis on local-global visual processing in ASD. *Psychological Bulletin*, 141(3), 549–573.
- Vedantam, R., Lawrence Zitnick, C., & Parikh, D. (2015). Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 4566-4575).
- Von Ahn, L., & Dabbish, L. (2004, April). Labeling images with a computer game. In *Proceedings of the SIGCHI conference on Human factors in computing systems* (pp. 319-326).

779 Wang, X., Chen, W., Wang, Y. F., & Wang, W. Y. (2018).
780 No metrics are perfect: Adversarial reward learning for
781 visual storytelling. *arXiv preprint arXiv:1804.09160*.
782 Wilson, M. (2005). *Constructing measures: An item response*
783 *modeling approach*. Lawrence Erlbaum.
784 Wilson, M., & De Boeck, P. (2004). Descriptive and
785 explanatory item response models. In P. De Boeck & M.
786 Wilson (Eds.), *Explanatory Item Response Models: A*
787 *Generalized Linear and Nonlinear Approach* (pp. 43–74).
788 Springer.
789 Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q., & Artzi,
790 Y. (2019). Bertscore: Evaluating text generation with
791 bert. *arXiv preprint arXiv:1904.09675*.