

Resolving References in Visually-Grounded Dialogue via Text Generation

Anonymous ACL submission

Abstract

Vision-language models (VLMs) have shown to be effective at image retrieval based on simple text queries, but text-image retrieval based on conversational input remains a challenge. Consequently, if we want to use VLMs for reference resolution in visually-grounded dialogue, the discourse processing capabilities of these models need to be augmented. To address this issue, we propose fine-tuning a causal large language model (LLM) to generate definite descriptions that summarize coreferential information found in the linguistic context of references. We then use a pretrained VLM to identify referents based on the generated descriptions, zero-shot. We evaluate our approach on a manually annotated dataset of visually-grounded dialogues and achieve results that, on average, exceed the performance of the baselines we compare against. Furthermore, we find that using referent descriptions based on larger context windows has the potential to yield higher returns.

1 Introduction

Visually-grounded dialogues are conversations in which participants make references to the visual world. Referring in conversation is understood to be a collaborative process, with shared responsibility for ensuring the successful identification of the referent (Clark and Wilkes-Gibbs, 1986). It is not uncommon for a definite reference to be established over multiple turns, with each separate contribution unlikely to be a minimally distinguishable description of the referent. Taken out of their use context, these referring expressions may be difficult, if not impossible, to resolve. Consider the example dialogue in Figure 1. The underspecified description “*the shiny one*” leads to a clarification question, “*Do you mean that red one?*”. To resolve the expression “*that red one*” to its referent, we need information from earlier in the conversation to understand that “*one*” is a proform of “*apple*”.

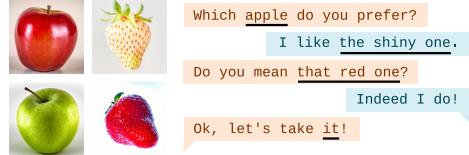


Figure 1: Example dialogue in which two participants discuss fruits. Expressions that denote one or more images are underlined.

Without this linguistic context, the red strawberry and the red apple are equally likely referents.

We can break the problem of reference resolution in visually-grounded dialogue down into three subproblems: (1) mention detection, or finding the expressions that can be grounded in the visual context (“*that red one*”); (2) aggregation of referent-specific information (linking “*apple*”, “*the shiny one*”, and “*that red one*”); and (3) referent identification, or the grounding of language (finding the referent that is best described by the three expressions from among a set of candidate referents). This final step requires bridging the gap between vision and language. For this purpose, we can turn to pretrained vision-language models (VLMs), which have shown to be effective at zero-shot text-image retrieval when given a description of an image (e.g., Radford et al., 2021; Jia et al., 2021; Li et al., 2023). However, current VLMs lack the discourse processing capabilities necessary for reference resolution in visually-grounded dialogue. Although some VLMs may correctly identify the red apple as the referent given the entire dialogue of Figure 1, dialogues are often vastly more complex than this hypothetical exchange. Take, for instance, the dialogue in Appendix A: with multiple mentions of different referents within the same utterance, such a brute-force method would immediately fail. It is clear that if we want VLMs to be effective for this purpose, their discourse processing capabilities need to be augmented.

To this end, we propose fine-tuning a causal large

language model (LLM) for the task of *referent description generation*. Referent description generation can be regarded as a special case of referring expression generation with the goal of always generating the most complete expression possible. For a given mention, the model is trained to generate a definite description that summarizes all information that has been explicitly disclosed about the referent during a conversation. For example, for the mention “*that red one*” in Figure 1 we would want the model to generate the description “*the shiny red apple*”. We will refer to the fine-tuned model as the *conversational referent description generator* (CRDG). The description generated by the CRDG is then used by a pretrained VLM to identify the referent, zero-shot. Our approach can be seen as an exploration of the limits of depending on linguistic context alone for generating referent descriptions, as the discourse processing and eventual grounding of the descriptions are entirely disjoint.

For the experiments presented in this paper we use data from the collaborative image ranking task A Game Of Sorts (AGOS, Willemsen et al., 2022). Referents are represented by separate, but visually similar images from a shared entity category. Due to their largely unrestricted nature and with a focus on the collaborative referential process, the collected dialogues form a challenging test bed for visually-grounded language understanding in conversation. We manually annotate the dialogues by marking mention spans and aligning the spans with the images they denote, and provide ground truth referent descriptions for all marked mentions.

Our main contributions are as follows:

- We present a generative approach to reference resolution in visually-grounded dialogue that frames the discourse processing side of the task as a causal language modeling problem;
- We show that it is possible to fine-tune a causal LLM to generate referent descriptions from dialogue to be used by a pretrained VLM for referent identification, zero-shot;
- We release the discussed materials, including our annotations for A Game Of Sorts (Willemsen et al., 2022)¹.

2 Background

Visually-grounded language understanding is fundamental for conversational agents that engage in

dialogue involving references to the visual world. Researchers have introduced a variety of tasks that provide data for development and frameworks for evaluation of visually-grounded dialogue models. These tasks often take the form of goal-oriented, dyadic interactions but differ in terms of, for example, the visual stimuli used, e.g. abstract figures or realistic photos; the roles assigned to participants, e.g. whether symmetric or asymmetric; constraints on message content, e.g. a fixed vocabulary; and the nature of the task, e.g. navigation, identification, ranking, and multi-turn visual question answering (e.g. Das et al., 2017; De Vries et al., 2017; Shore et al., 2018; Ilinykh et al., 2019; Haber et al., 2019; Udagawa and Aizawa, 2019; Willemsen et al., 2022). It has been noted that the task configuration can significantly impact the extent to which certain dialogue phenomena, such as coreferences and clarification requests, are represented in the collected data, if at all (Agarwal et al., 2020; Haber et al., 2019; Ilinykh et al., 2019; Schlangen, 2019; Willemsen et al., 2022). Tasks that heavily constrain the interactions do not reflect the complex nature of dialogue to the same degree as tasks that have been designed for these phenomena to naturally emerge as part of the discourse, such as A Game Of Sorts (Willemsen et al., 2022), which we use in this paper.

The terms referring expression comprehension (e.g. Yu et al., 2016), referring expression grounding (e.g. Zhang et al., 2018), referring expression recognition (e.g. Cirik et al., 2018), and reference resolution (e.g. Kennington et al., 2015) have been used interchangeably to describe the problem of mapping the language that denotes a referent to a representation of that referent in the visual modality. Prior work noted the importance of referring expressions to conversation, but often modeled the problem independent of the dialogue (e.g. Cirik et al., 2018; Schlangen et al., 2016; Yu et al., 2016; Zhang et al., 2018). The granularity at which grounding occurs may differ between works, as the language may be mapped to bounding boxes of individual objects (Cirik et al., 2018; Schlangen et al., 2016; Yu et al., 2016; Zhang et al., 2018), objects or larger image regions represented by segmentation masks (Liu et al., 2017), or entire images altogether (Haber et al., 2019; Takmaz et al., 2020).

To address the problem computationally, both modalities must in some way be encoded. Engineered visual feature representations and sim-

¹ Available upon publication.

ple language models such as those based on n-grams (e.g. Kennington et al., 2015; Kennington and Schlangen, 2017; Shore and Skantze, 2018) have been mostly replaced with more powerful learned representations that embed the images and text in high-dimensional vector spaces (Haber et al., 2019; Takmaz et al., 2020). This has made it possible to resolve references by computing representational similarity between an encoding of the text that contains a mention and the embeddings of the candidate referents, where the candidate that has the highest matching score is assumed to be the referent (Haber et al., 2019; Takmaz et al., 2020).

Recent work on multimodal representation learning has shown that jointly embedding text and images can work at scale. Trained using a contrastive objective, maximizing representational similarity between true pairings of images and text while simultaneously minimizing similarity of false pairs, vision-language models (VLMs) such as CLIP (Radford et al., 2021), ALIGN (Jia et al., 2021), BLIP (Li et al., 2022), and BLIP-2 (Li et al., 2023), have shown to be effective zero-shot classifiers, outperforming the previous state-of-the-art on various benchmarks without the need for further fine-tuning on specific tasks. However, despite their noteworthy image-text matching performance based on simple text queries, these VLMs lack the discourse processing capabilities required for reference resolution in visually-grounded dialogue. Even a simplified example, such as shown in Figure 1, illustrates a fundamental challenge, namely that of coreference resolution. The interpretation of anaphoric pronouns such as “*it*” is dependent on their antecedents. Without resolving its coreferences first, identifying the referent based on the pronoun alone leads to a random guess.

To improve downstream performance on discourse processing tasks involving coreference, prior work has approached the problem as one of transforming the original input based on linguistic context. This was done either via substitution, such as in Bhattacharjee et al. (2020) where pronouns were substituted for more descriptive mentions of the same referent, or via generation, such as in Quan et al. (2019) where entire utterances were reconstructed in a pragmatically complete manner with coreferences and ellipses resolved. To the best of our knowledge, this approach has not yet been applied to reference resolution in visually-grounded dialogue.

Most contemporary natural language processing (NLP) works use Transformer-based language models (Vaswani et al., 2017). For text generation tasks, it is common to use (unidirectional) autoregressive, or *causal*, language models such as GPT (Radford et al., 2018). While processing sequences, causal language models mask the future, allowing the model to only attend to the current and previous tokens while predicting the next token. A persistent trend has been to scale up language models, both in terms of their parameter count and the size of their training datasets. These increasingly larger models, such as GPT-3 (Brown et al., 2020), OPT (Zhang et al., 2022), PaLM (Chowdhery et al., 2022), and LLaMa (Touvron et al., 2023), have been dubbed *large language models* (LLMs). The current leading paradigm to modeling downstream NLP tasks is based on transfer learning, where a pretrained LLM is fine-tuned for a specific task on a smaller, domain-specific dataset.

3 Method

We treat visually-grounded reference resolution as a text-image retrieval task, where referents are represented by images. We leave finer-grained grounding of words and phrases to image regions or individual entities or parts thereof for future work.

3.1 Proposed Framework

We frame the discourse processing side of the task as a causal language modeling problem. Figure 2 shows a visualization of the proposed framework.

Task Definition We denote the dialogue as $D = (u_1, u_2, \dots, u_n)$, where each u_i represents an utterance. Each utterance consists of an ordered sequence of tokens. An utterance may contain one or more mentions, denoted as M . A mention is an ordered subsequence of tokens from an utterance. A mention has an exophoric referent, denoted as R . A mention is embedded in what we call its linguistic context, denoted as L . As an ordered subsequence of D , the linguistic context of a given mention consists of the utterance in which it is contained and all preceding utterances. The number of preceding utterances, hereafter referred to as the dialogue history, may be capped if a finite size context window is defined. The aim of visually-grounded reference resolution is to resolve a reference to its referent, i.e. to identify R for a given M , from a set of candidate referents, denoted as C , such that $R \subseteq C$.

Referent Description Generation We propose to

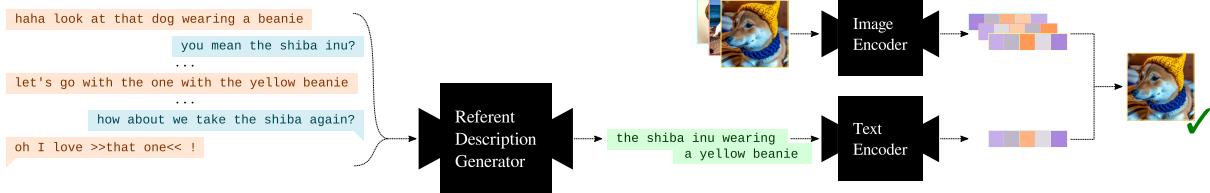


Figure 2: Visualization of the proposed visually-grounded reference resolution framework. With the CRDG we generate a referent description for a marked mention, to be used by a (frozen) pretrained VLM for referent identification.

generate a definite description, denoted as Y , for a given mention M that summarizes all that has been disclosed in L about the referent R . For this purpose, we fine-tune a causal LLM that learns to generate Y conditioned on L . Y is a sequence of tokens expected to be largely constructed from tokens that appear, or are some derivative of tokens that appear, in the coreference chain of R , which is contained in L . We refer to the fine-tuned model as the *conversational referent description generator* (CRDG). For an example of the context dependency of referent description content, see Figure 4 in Appendix B.

LLM Input We mark M in u_i by inserting positional markers as special tokens to indicate the beginning and end of the mention span. We prepend each utterance in L with a speaker token to indicate the source of the contribution. When D is task-oriented, we update L by prepending task instructions, i.e. a special token followed by a sequence of tokens describing the task performed by the dialogue participants. For an example of the input to the LLM, see Figure 5 in Appendix B.

Text-Image Retrieval We use a pretrained VLM to identify R from C based on Y , zero-shot. We use the text encoder of the VLM to encode Y into an n -dimensional feature vector, denoted as \mathbf{v} . We use the image encoder of the VLM to encode each candidate referent of C into an n -dimensional feature vector, which gives a $|C| \times n$ matrix, denoted as \mathbf{A} . We then compute their matrix-vector product. For single-image referents, i.e. when $|R| = 1$, we take the referent to be $R = \text{argmax}(\mathbf{Av})$.

In order to produce accurate referent descriptions, the CRDG must implicitly learn to perform coreference resolution as we do not provide explicit supervision for this subtask. In each sample, only the current mention for which we want the model to generate a description is marked; none of its coreferences are in any way indicated. A principal ad-

vantage of our model is that it can resolve multiple mentions, even when they have different referents, appearing in the same utterance, including nested mentions. Note that for the purpose of this study, we assume mention detection to be solved. As it stands, using this framework in production requires a separate model to propose candidate mentions at the span level.

3.2 Baseline Models

As a lower bound, we report random chance performance. In addition, we compare performance of our approach to baselines based on simple heuristics and a coreference resolution model.

3.2.1 Heuristics

Mention We evaluate the image retrieval performance when the VLMs are presented with just the marked mentions.

Substitution We improve upon the mention-only baseline by substituting proforms, e.g. pronouns such as “it”, and mentions without descriptive content, e.g. phrases such as “*the one you mentioned*”, with the most recent mention that does not belong to either category. This is expected to be a relatively strong baseline when mentions are specific and anaphora have mostly local antecedents.

3.2.2 Coreference Resolution

We opt for an off-the-shelf² span-based coreference resolution model (**coref**) originally presented in Lee et al. (2018), but that has since been updated to use SpanBERT (Joshi et al., 2020) instead of the original GloVe embeddings (Pennington et al., 2014). For each mention, we use the model to resolve its coreference links and aggregate all coreferential information in its cluster based on the given context window.

²https://github.com/allenai/allennlp-models/tree/main/allennlp_models/coref

348 We experiment with two different representations
349 of the referent descriptions from this model,
350 those being (1) a concatenation of all of the
351 mention’s coreferences and (2) an ordered *set-of-words*
352 representation that contains only the unique lexical
353 items in the cluster. To offset that this model
354 was not specifically trained to handle coreference
355 in conversation, we provide it with the contents of
356 the span of the mention when it does not manage
357 to detect the mention itself and, consequently, not
358 connect it to any of its coreferences. For partial
359 matches, in addition to adding all tokens from the
360 cluster associated with the match, we also add the
361 missing tokens from the span to the description.

362 4 Experiments

363 4.1 Data

364 We use the dialogues from the collaborative image
365 ranking task **A Game Of Sorts** (AGOS, [Willemsen](#)
366 et al., 2022) for our experiments. In AGOS, two
367 players are asked to rank a set of images based on
368 a given sorting criterion. They see the same set of
369 images, but the position of the images on the screen
370 is randomized for each player. Through a largely
371 unrestricted conversation, and without being able
372 to see the perspective of the other player, the play-
373 ers need to agree on how to rank the images given
374 the sorting criterion. Sorting criteria are embedded
375 in scenarios that are intended to create a discussion,
376 leading to mixed-initiative interactions with both
377 parties contributing to the discourse. Each interac-
378 tion takes place over four rounds with the same set
379 of nine images, effectively guaranteeing repeated
380 references. The image sets used for the game cover
381 five different image categories. Each set contains
382 nine images with each image representing an entity
383 from one of these categories as its main subject.
384 [Willemsen et al.](#) (2022) collected three interactions
385 per image set for a total of 15 dialogues.

386 **Ground Truth** Our formulation of the visually-
387 grounded reference resolution problem requires
388 span-based annotations of mentions aligned with
389 the image(s) they denote. These annotations are
390 the basis of our ground truth references for both
391 training and evaluation. We follow [Willemsen](#)
392 et al. (2022) regarding the marking of mentions
393 in AGOS, in that we only annotate those that are
394 either singletons or are part of an identity relation
395 with other mentions that have an exophoric referent
396 that is part of the visual context, i.e. regardless of
397 form, any referring expression that is meant to de-

398 note one or more of the images. During the game,
399 players were asked to provide self-annotations: for
400 each message they sent they were asked to indicate
401 which image(s), if any, they were referring to. We
402 use these self-annotations, post-edited where neces-
403 sary, to manually mark the spans of mentions that
404 can be grounded in the visual context.

405 We create three different representations of the
406 ground truth referent descriptions. Two are auto-
407 matically extracted from the marked mentions and
408 are similar in structure to the labels of the **coref**
409 baseline, i.e. (1) an incremental concatenation of
410 the reference chain and (2) an incremental ordered
411 set of words consisting of the unique lexical items
412 in the cluster. The third are manually constructed
413 labels that summarize reference chains as definite
414 descriptions. For each representation, the context
415 window dictates which references are considered
416 for the label.

417 4.2 Model Specifications

418 For pointers to implementations, we refer the reader
419 to our repository³.

420 4.2.1 LLMs

421 We fine-tune two LLMs, GPT-2 ([Radford et al.](#),
422 2019) and GPT-3 ([Brown et al.](#), 2020), for con-
423 versational referent description generation. For
424 hyperparameters, see our Supplementary Material.

425 **GPT-2** We fine-tune the 1.5 billion parameter GPT-
426 2 model.

427 **GPT-3** We fine-tune the 175 billion parameter
428 [davinci](#) base model using the OpenAI API.

429 4.2.2 VLMs

430 We evaluate the zero-shot text-image retrieval per-
431 formance of several pretrained VLMs for our task,
432 those being CLIP ([Radford et al.](#), 2021), ALIGN
433 ([Jia et al.](#), 2021), BLIP ([Li et al.](#), 2022), and BLIP-2
434 ([Li et al.](#), 2023).

435 **CLIP** We evaluate two variants of CLIP, CLIP ViT-
436 B/16 and CLIP ViT-L/14.

437 **ALIGN** We use the COYO-ALIGN implementa-
438 tion trained from scratch on COYO-700M.

439 **BLIP** We use the BLIP base model.

440 **BLIP-2** We use the BLIP-2 model that was fine-
441 tuned on the [Karpathy and Fei-Fei](#) (2015) training
442 set split of MS COCO ([Lin et al.](#), 2014).

³Available upon publication.

443 4.3 Evaluation

444 We perform (nested) five-fold cross-validation by
445 partitioning the AGOS dataset along the five image
446 sets. To avoid leakage, for each run we use the
447 three dialogues from one image set as the held out
448 test set and train on the twelve dialogues from the
449 four other image sets. To evaluate how dialogue
450 history affects results, we report performance of
451 the different methods for two context windows, **3**
452 and **7**. In addition, we examine whether increasing
453 the size of the context window further would,
454 in principle, lead to greater returns, by assessing
455 ground-truth performance for windows of **13** and
456 the **full** dialogue context. Finally, we conduct an
457 error analysis of the generated descriptions.

458 Note that because our model does not have access
459 to game state information with respect to the
460 visual context, we make a simplifying assumption
461 with regard to the images and reduce the candidate
462 set as the game progresses. A successfully ranked
463 image is no longer considered part of the visual
464 context for that round. Although this does mean
465 that the models will not be able to identify the
466 referent for references to ranked images, as they will
467 not be part of the candidate set, such references
468 are an extremely rare occurrence, as players must
469 discuss the unranked images to progress with the
470 task. For the sake of completeness, we will also
471 report results for the unchanged candidate set.

472 4.3.1 Metrics

473 We measure task success for visually-grounded
474 reference resolution in terms of text-image retrieval
475 performance. In addition, we estimate the quality
476 of the referent descriptions by comparing them to
477 the manually constructed ground truth labels using
478 text similarity metrics.

479 **Text-Image Retrieval** We estimate the image
480 retrieval performance based on accuracy [0, 1], mean
481 reciprocal rank (MRR) [0, 1], and normalized dis-
482 counted cumulative gain (NDCG) [0, 1]. We limit
483 our evaluation to single-image referents. Accuracy
484 is top-1 accuracy.

485 For our random lower bound, we can calculate
486 the expected values for accuracy and MRR. For
487 top-1 accuracy we take 1 over the size of the set
488 of candidate images per item, averaging over all
489 items. For MRR we take 1 over the size of the
490 set of candidate images, divided by two per item,
491 averaging over all items. Calculating an expected
492 value for NDCG of a random model is intractable

493 due to its dependence on relevancy scores.

494 **Text Generation** We evaluate the output from the
495 CRDGs by comparing the generated descriptions
496 to the manually constructed ground truth labels
497 using metrics to quantify similarity. We use the
498 Jaccard index [0, 1] to assess vocabulary overlap.
499 We use BLEU [0, 1] (Papineni et al., 2002) to as-
500 sess similarity based on n-gram overlap (unigrams
501 to four-grams). We use the longest common sub-
502 sequence variant of ROUGE [0, 1] (Lin, 2004), i.e.
503 ROUGE-L, as a further indication of the preser-
504 vation of word order. In addition, we opt for an
505 embedding-based metric as a proxy for semantic
506 equivalence between the predicted and reference
507 sequences. For this purpose, we compute cosine
508 similarity [0, 1] between text embeddings.

509 4.3.2 Human

510 We conduct two different human subject experi-
511 ments to assess human performance for this task.
512 We provide additional details about the experimen-
513 tal setup in the Supplementary Material.

514 **Independent** We conduct an experiment aimed at
515 comparing VLM and human performance on the
516 task where every trial is independent. Participants
517 are given a referent description and are asked to se-
518 lect from a set of candidate images the image they
519 believe is best described by the label. The images
520 and labels are presented to the participants indepen-
521 dent of the dialogue. Note that we evaluate with the
522 reduced candidate set. The referent descriptions are
523 the manually constructed ground truth labels based
524 on the **full** dialogue context. To collect data for
525 all labels, ensuring independence of observations,
526 we recruited 322 participants via crowdsourcing.
527 The crowdworkers were financially compensated
528 for their contributions.

529 **Holistic** We conduct an experiment in which men-
530 tions are shown to participants within the context
531 of the dialogue. For each mention, the participants
532 are presented with the dialogue leading up to and
533 including the message which contains the reference.
534 The start and end of the span of the mention that
535 the participant is asked to resolve are visually in-
536 dicated. For each marked mention, the participant
537 is asked to select which image or images are refe-
538 renced. As they progress with the task, participants
539 will have access to increasingly more of the dia-
540 logue history. For each mention the participants
541 are presented with all images, but with a visual in-
542 dication of their status, i.e. for each image whether
543 the players had managed to successfully rank it

| | Accuracy | | MRR | | NDCG | |
|--------------|----------|-----|-----|-----|------|-----|
| | 3 | 7 | 3 | 7 | 3 | 7 |
| Random | .22 | .22 | .43 | .43 | - | - |
| Mention | .59 | .59 | .73 | .73 | .79 | .79 |
| Substitution | .68 | .68 | .80 | .80 | .85 | .85 |
| coref, chain | .65 | .66 | .78 | .79 | .83 | .84 |
| coref, set | .66 | .66 | .78 | .79 | .84 | .84 |
| GT, chain | .73 | .74 | .83 | .85 | .87 | .88 |
| GT, set | .73 | .75 | .84 | .85 | .87 | .89 |
| GT, manual | .72 | .74 | .83 | .84 | .87 | .88 |
| GPT-2 | .64 | .60 | .74 | .76 | .83 | .80 |
| GPT-3 | .69 | .71 | .81 | .82 | .86 | .86 |

Table 1: Cross-validated image retrieval performance averaged over five folds for single-image referents. *Note.* Scores shown are of VLM that averaged best performance (BLIP-2). Scores are rounded to the nearest hundredth. GT = ground truth.

at that point in the interaction. We recruited 23 participants via crowdsourcing. For each of the 15 AGOS dialogues we collected data from two different participants. Each participant was allowed to provide data for at most one dialogue per image set. The crowdworkers were financially compensated for their contributions.

5 Results

5.1 Text-Image Retrieval

Table 1 shows, for context windows 3 and 7, the zero-shot text-image retrieval performance results for the VLM that averaged best performance over the five folds, namely BLIP-2. For the text-image retrieval accuracy achieved by the other VLMs, performance on the unreduced candidate set, and accuracy per fold for BLIP-2, see Appendix C.

As can be seen from the results presented in Table 1, we achieve best performance with a fine-tuned GPT-3 as the CRDG and BLIP-2 for zero-shot text-image retrieval. In addition to outperforming the baselines, we find that GPT-3 is a more performant discourse processor for this task than GPT-2. This result is consistent between the VLMs.

Results generally show a slight increase in performance when increasing the context window from 3 to 7. Performance on the ground truth reference descriptions for context windows 13 and the full dialogue shows this trend persists, with BLIP-2 achieving approximately 75% and 82% accuracy, respectively. A plot of the performance for the four context windows is shown in Figure 6 in Appendix C. This result suggests that the size of the context window may have a significant impact on performance, with a 10% increase in accuracy

| | GPT-2 | | GPT-3 | |
|---------|-------|-----|-------|-----|
| | 3 | 7 | 3 | 7 |
| BLEU | .55 | .47 | .75 | .70 |
| ROUGE-L | .71 | .65 | .86 | .83 |
| Jaccard | .44 | .35 | .70 | .63 |
| Cosine | .88 | .85 | .96 | .95 |

Table 2: Text generation metrics evaluation results averaged over five folds for single-image referents. *Note.* Scores are rounded to the nearest hundredth.

from 3 to full. Furthermore, the VLMs do not seem overly sensitive to the composition of the referent descriptions, as performance is largely comparable between the automatically generated and the manually constructed ground truth labels.

We find that BLIP-2 is on par with human text-image retrieval performance in terms of top-1 accuracy for the manually constructed ground truth referent descriptions based on the full dialogue history for single-image referents, as our human participants averaged roughly 80% accuracy in the independent setup. However, when we compare these results with the single-image referent text-image retrieval performance in the holistic setup, we see that the upper bound for this task when references are resolved within the combined linguistic and extralinguistic dialogue context is likely considerably higher as our human participants averaged approximately 91% accuracy (average of best performance per dialogue is roughly 93%).

5.2 Text Generation

Table 2 shows the text generation metric results averaged over the five folds, providing an indication of the extent to which the fine-tuned LLMs managed to generate referent descriptions that approximate the manually constructed ground truth labels. We observe that an increase in context window size results in a decrease in scores, which is consistent across metrics. Interestingly, we did not find such a decrease with respect to text-image retrieval performance. We do again find GPT-3 to be more performant than GPT-2, here in terms of approximating the ground truth.

5.3 Error Analysis

Examining the output from the fine-tuned GPT-3 model, we observe a number of recurring errors.

The most notable errors are those where the model fails to link a mention to (all of) its coreferences that are present in the dialogue segment, or links mentions that denote different referents. For

example, for one mention the ground truth label is “*the sheep dog*”, but the generated label was “*the sheep dog with a leash*”; the model incorrectly attributed the prepositional phrase to the mention as it was actually a descriptor for a different referent. Related, since the CRDGs function at the message level, a mention can have both anaphoric and cataphoric coreferences when there are multiple mentions of the same referent in an utterance. An example of such an utterance is “*Good question. I think the angry one also looks a little wild. So that could be an option as well. I mean the one with white nose and forehead*”, where “*the angry one*”, “*that*”, and “*the one with white nose and forehead*” are all mentions of the same referent with the same ground truth label “*the angry dog with a white nose and forehead*”. The model generates this correctly for the latter two, but not the former one for which only “*the angry dog*” was generated, meaning it correctly substituted the proform but did not link the mention with its cataphoric coreferences.

Finally, some generated referent descriptions differ from the ground truth in terms of lexical choice or syntax, but not in terms of information content. This negatively affects scores of text generation metrics based on overlapping content in particular, but these are otherwise not meaningful errors as there are multiple ways to construct semantically similar descriptions, e.g., “*the big dog which looks scary*” versus “*the big scary-looking dog*”.

6 Discussion

We have presented an approach to visually-grounded reference resolution that frames the discourse processing side of the task as a causal language modeling problem. By fine-tuning an LLM to generate referent descriptions for marked mentions in dialogue segments from the collaborative image ranking task A Game Of Sorts (Willemsen et al., 2022), we demonstrate the possibility of treating referent identification as a zero-shot text-image retrieval problem by using pretrained VLMs for the grounding of the generated labels. As we have not in any way indicated coreferential relations in the fine-tuning training data, our results imply that certain pretrained LLMs, here GPT-3, may learn to resolve coreferences implicitly without the need for explicit supervision for this fundamental subtask.

In this work, we have treated the processing of the discourse as entirely disjoint of the visual modality. As such, it has inherent limitations. The

mentions we find in the dialogues have not been produced void of the extralinguistic context. The dialogue participants could rely on co-observed visual stimuli to help resolve otherwise ambiguous language use. From linguistic context alone, some ambiguities, such as prepositional phrase attachment, may be impossible to resolve. It is, therefore, noteworthy that the downstream zero-shot text-image retrieval performance using the generated descriptions from our unimodal approach far exceeds chance level accuracy, with the potential for results to improve further given access to the full dialogue history, as we found that the ground truth labels based on larger context windows achieve greater text-image retrieval performance. However, the results from our holistic human evaluation support the notion that a multimodal approach should ultimately prove even more effective.

We found that a decrease in text generation metric scores did not necessarily indicate a similar decrease in text-image retrieval performance, suggesting that the generated descriptions captured sufficiently discriminative information about the referents and achieved similar grounding accuracy despite not approximating the ground truth labels to the same extent. It is also important to note that mentions may not have a single, canonical ground truth referent description due to lexical and syntactic variations between referring attempts.

Despite the relatively small size of the dataset collected by Willemsen et al. (2022), we were still able to fine-tune GPT-3 to perform the task with greater accuracy than the baselines, which speaks to the sample efficiency of pretrained LLMs. In comparison, we find that the much smaller GPT-2 is prone to intrusions from the fine-tuning training data and more often fails to resolve the coreferences correctly. Although the complexity of the discourse warrants the use of more powerful models, it is, nevertheless, likely that any LLM used for the task would benefit from a larger fine-tuning dataset. Related, benchmarking performance on other visually-grounded dialogue tasks would provide insights into the generalizability of the method.

In addition to pursuing a multimodal approach, finer-grained grounding, and evaluating our method on other datasets, possible avenues for future work include expanding the annotations to include coreferential relations other than identity relations, addressing multi-image referents, and unifying the method with a mention proposal system.

719 References

- 720 Shubham Agarwal, Trung Bui, Joon-Young Lee, Ioan- 778
721 nis Konstas, and Verena Rieser. 2020. *History for 779
722 Visual Dialog: Do we really need it?* In *Proceedings 780*
723 of the 58th Annual Meeting of the Association for 781
724 Computational Linguistics, pages 8182–8197, On- 782
725 line. Association for Computational Linguistics.
- 726 Santanu Bhattacharjee, Rejwanul Haque, Gideon Mail- 781
727 lette de Buy Wenniger, and Andy Way. 2020. Investi- 782
728 gating Query Expansion and Coreference Resolution 783
729 in Question Answering on BERT. In *Natural Lan- 784
730 guage Processing and Information Systems*, pages 785
731 47–59, Cham. Springer International Publishing.
- 732 Tom Brown, Benjamin Mann, Nick Ryder, Melanie 786
733 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind 787
734 Neelakantan, Pranav Shyam, Girish Sastry, Amanda 788
735 Askell, Sandhini Agarwal, Ariel Herbert-Voss, 789
736 Gretchen Krueger, Tom Henighan, Rewon Child, 790
737 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens 791
738 Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz 792
739 Litwin, Scott Gray, Benjamin Chess, Jack Clark, 793
740 Christopher Berner, Sam McCandlish, Alec Radford, 794
741 Ilya Sutskever, and Dario Amodei. 2020. *Language 795
742 Models are Few-Shot Learners*. In *Advances in Neural 796
743 Information Processing Systems*, volume 33, pages 797
744 1877–1901. Curran Associates, Inc.
- 746 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, 800
747 Maarten Bosma, Gaurav Mishra, Adam Roberts, 801
748 Paul Barham, Hyung Won Chung, Charles Sutton, 802
749 Sebastian Gehrmann, Parker Schuh, Kensen Shi, 803
750 Sasha Tsvyashchenko, Joshua Maynez, Abhishek 804
751 Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vin- 805
752 odkumar Prabhakaran, Emily Reif, Nan Du, Ben 806
753 Hutchinson, Reiner Pope, James Bradbury, Jacob 807
754 Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, 808
755 Toju Duke, Anselm Levskaya, Sanjay Ghemawat, 809
756 Sunipa Dev, Henryk Michalewski, Xavier Garcia, 810
757 Vedant Misra, Kevin Robinson, Liam Fedus, Denny 811
758 Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, 812
759 Barret Zoph, Alexander Spiridonov, Ryan Sepassi, 813
760 David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pil-
761 lai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, 814
762 Rewon Child, Oleksandr Polozov, Katherine Lee, 815
763 Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark 816
764 Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy 817
765 Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, 818
766 and Noah Fiedel. 2022. PaLM: Scaling Language
767 Modeling with Pathways. *_eprint:* 2204.02311.
- 769 Volkan Cirik, Louis-Philippe Morency, and Taylor Berg- 819
770 Kirkpatrick. 2018. *Visual Referring Expression 820
771 Recognition: What Do Systems Actually Learn?* In 821
772 *Proceedings of the 2018 Conference of the North 822
773 American Chapter of the Association for Computational 823
774 Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 781–787, New Orleans, Louisiana. Association for Computational Lin- 824
775 guistics.
- 776 Herbert H. Clark and Deanna Wilkes-Gibbs. 1986. *Re- 781
777 ferring as a collaborative process*. *Cognition*, 22(1):1– 782
778 39.
- 779 Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, 783
780 Deshraj Yadav, José M. F. Moura, Devi Parikh, and 781
781 Dhruv Batra. 2017. *Visual Dialog*. In *2017 IEEE 782
782 Conference on Computer Vision and Pattern Recog- 783
783 nition (CVPR)*, pages 1080–1089.
- 784 Harm De Vries, Florian Strub, Sarah Chandar, Olivier 785
785 Pietquin, Hugo Larochelle, and Aaron Courville. 2017. *GuessWhat?! Visual Object Discovery 786
786 through Multi-modal Dialogue*. In *2017 IEEE 787
787 Conference on Computer Vision and Pattern Recog- 788
788 nition (CVPR)*, pages 4466–4475.
- 789 Janosch Haber, Tim Baumgärtner, Ece Takmaz, Lieke 790
790 Gelderloos, Elia Bruni, and Raquel Fernández. 2019. 791
791 *The PhotoBook Dataset: Building Common Ground 792
792 through Visually-Grounded Dialogue*. In *Proceed- 793
793 ings of the 57th Annual Meeting of the Association for 794
794 Computational Linguistics*, pages 1895–1910, Flo- 795
795 rence, Italy. Association for Computational Linguis- 796
796 tics.
- 797 Nikolai Ilinykh, Sina Zarrieß, and David Schlangen. 800
798 2019. *Meet Up! A Corpus of Joint Activity Dia- 801
799 logues in a Visual Environment*. In *Proceedings of 802
800 the 23rd Workshop on the Semantics and Pragmatics 803
801 of Dialogue - Full Papers*, London, United Kingdom. 804
802 SEMDIAL.
- 803 Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana 805
804 Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, 806
805 Zhen Li, and Tom Duerig. 2021. *Scaling Up Vi- 807
806 sual and Vision-Language Representation Learning 808
807 With Noisy Text Supervision*. In *Proceedings of the 809
808 38th International Conference on Machine Learning*, 810
809 volume 139 of *Proceedings of Machine Learning 811
810 Research*, pages 4904–4916. PMLR.
- 811 Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, 812
812 Luke Zettlemoyer, and Omer Levy. 2020. *Span- 813
813 BERT: Improving Pre-training by Representing and 814
814 Predicting Spans*. *Transactions of the Association 815
815 for Computational Linguistics*, 8:64–77.
- 816 Andrej Karpathy and Li Fei-Fei. 2015. Deep Visual- 817
817 Semantic Alignments for Generating Image Descrip- 818
818 tions. In *Proceedings of the IEEE Conference on 819
819 Computer Vision and Pattern Recognition (CVPR)*.
- 820 Casey Kennington, Livia Dia, and David Schlangen. 821
821 2015. *A Discriminative Model for Perceptually- 822
822 Grounded Incremental Reference Resolution*. In *Pro- 823
823 ceedings of the 11th International Conference on 824
824 Computational Semantics*, pages 195–205, London, 825
825 UK. Association for Computational Linguistics.
- 826 Casey Kennington and David Schlangen. 2017. *A sim- 827
827 ple generative model of incremental reference res- 828
828 olution for situated dialogue*. *Computer Speech & 829
829 Language*, 41:43–67.

833 Kenton Lee, Luheng He, and Luke Zettlemoyer. 2018.
834 **Higher-Order Coreference Resolution with Coarse-**
835 **to-Fine Inference.** In *Proceedings of the 2018 Con-*
836 *ference of the North American Chapter of the Asso-*
837 *ciation for Computational Linguistics: Human Lan-*
838 *guage Technologies, Volume 2 (Short Papers)*, pages
839 687–692, New Orleans, Louisiana. Association for
840 Computational Linguistics.

841 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi.
842 2023. BLIP-2: Bootstrapping Language-Image Pre-
843 training with Frozen Image Encoders and Large Lan-
844 guage Models. *_eprint:* 2301.12597.

845 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi.
846 2022. **BLIP: Bootstrapping Language-Image Pre-**
847 **training for Unified Vision-Language Understanding**
848 **and Generation.** In *Proceedings of the 39th Inter-*
849 *national Conference on Machine Learning*, volume
850 162 of *Proceedings of Machine Learning Research*,
851 pages 12888–12900. PMLR.

852 Chin-Yew Lin. 2004. **ROUGE: A Package for Auto-**
853 **matic Evaluation of Summaries.** In *Text Summariza-*
854 *tion Branches Out*, pages 74–81, Barcelona, Spain.
855 Association for Computational Linguistics.

856 Tsung-Yi Lin, Michael Maire, Serge Belongie, James
857 Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,
858 and C. Lawrence Zitnick. 2014. **Microsoft COCO:**
859 **Common Objects in Context.** In *Computer Vision –*
860 *ECCV 2014*, pages 740–755, Cham. Springer Interna-
861 tional Publishing.

862 Chenxi Liu, Zhe Lin, Xiaohui Shen, Jimei Yang, Xin
863 Lu, and Alan Yuille. 2017. **Recurrent Multimodal**
864 **Interaction for Referring Image Segmentation.** In
865 *2017 IEEE International Conference on Computer*
866 *Vision (ICCV)*, pages 1280–1289.

867 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-
868 Jing Zhu. 2002. **Bleu: a Method for Automatic Eval-**
869 **uation of Machine Translation.** In *Proceedings of the*
870 *40th Annual Meeting of the Association for Compu-*
871 *tational Linguistics*, pages 311–318, Philadelphia,
872 Pennsylvania, USA. Association for Computational
873 Linguistics.

874 Jeffrey Pennington, Richard Socher, and Christopher
875 Manning. 2014. **GloVe: Global Vectors for Word**
876 **Representation.** In *Proceedings of the 2014 Confer-*
877 *ence on Empirical Methods in Natural Language Pro-*
878 *cessing (EMNLP)*, pages 1532–1543, Doha, Qatar.
879 Association for Computational Linguistics.

880 Jun Quan, Deyi Xiong, Bonnie Webber, and Changjian
881 Hu. 2019. **GECOR: An End-to-End Generative El-**
882 **lipsis and Co-reference Resolution Model for Task-**
883 **Oriented Dialogue.** In *Proceedings of the 2019 Con-*
884 *ference on Empirical Methods in Natural Language*
885 *Processing and the 9th International Joint Confer-*
886 *ence on Natural Language Processing (EMNLP-IJCNLP)*, pages 4547–4557, Hong Kong, China. As-
887 sociation for Computational Linguistics.

888 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
889 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-
890 tray, Amanda Askell, Pamela Mishkin, Jack Clark,
891 Gretchen Krueger, and Ilya Sutskever. 2021. **Learn-**
892 **ing Transferable Visual Models From Natural Lan-**
893 **guage Supervision.** In *Proceedings of the 38th Inter-*
894 *national Conference on Machine Learning*, volume
895 139 of *Proceedings of Machine Learning Research*,
896 pages 8748–8763. PMLR.

897 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya
898 Sutskever, and others. 2018. Improving language
899 understanding by generative pre-training. Publisher:
900 OpenAI.

901 Alec Radford, Jeffrey Wu, Rewon Child, David Luan,
902 Dario Amodei, and Ilya Sutskever. 2019. Language
903 Models are Unsupervised Multitask Learners.

904 David Schlangen. 2019. **Grounded Agreement Games:**
905 **Emphasizing Conversational Grounding in Visual**
906 **Dialogue Settings.** *CoRR*, abs/1908.11279. ArXiv:
907 1908.11279.

908 David Schlangen, Sina Zarrieß, and Casey Kenning-
909 ton. 2016. **Resolving References to Objects in Photo-**
910 **graphs using the Words-As-Classifiers Model.** In
911 *Proceedings of the 54th Annual Meeting of the Asso-*
912 *ciation for Computational Linguistics (Volume 1:*
913 *Long Papers)*, pages 1213–1223, Berlin, Germany.
914 Association for Computational Linguistics.

915 Todd Shore, Theofronia Androulakaki, and Gabriel
916 Skantze. 2018. **KTH Tangrams: A Dataset for Re-**
917 **search on Alignment and Conceptual Pacts in Task-**
918 **Oriented Dialogue.** In *Proceedings of the Eleventh*
919 *International Conference on Language Resources and*
920 *Evaluation (LREC 2018)*, Miyazaki, Japan. Eu-
921 *ropean Language Resources Association (ELRA).*

922 Todd Shore and Gabriel Skantze. 2018. **Using Lexical**
923 **Alignment and Referring Ability to Address Data**
924 **Sparsity in Situated Dialog Reference Resolution.**
925 In *Proceedings of the 2018 Conference on Empiri-*
926 *cal Methods in Natural Language Processing*, pages
927 2288–2297, Brussels, Belgium. Association for Com-
928 putational Linguistics.

929 Ece Takmaz, Mario Giulianelli, Sandro Pezzelle, Ara-
930 bella Sinclair, and Raquel Fernández. 2020. **Refer,**
931 **Reuse, Reduce: Generating Subsequent References**
932 **in Visual and Conversational Contexts.** In *Proced-*
933 *ings of the 2020 Conference on Empirical Methods in*
934 *Natural Language Processing (EMNLP)*, pages
935 4350–4368, Online. Association for Computational
936 Linguistics.

937 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier
938 Martinet, Marie-Anne Lachaux, Timothée Lacroix,
939 Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal
940 Azhar, Aurelien Rodriguez, Armand Joulin, Edouard
941 Grave, and Guillaume Lample. 2023. **LLaMA: Open**
942 **and Efficient Foundation Language Models.** *_eprint:*
943 2302.13971.

944

945 Takuma Udagawa and Akiko Aizawa. 2019. A Natural
946 Language Corpus of Common Grounding under
947 Continuous and Partially-Observable Context. In
948 *Proceedings of the Thirty-Third AAAI Conference on*
949 *Artificial Intelligence and Thirty-First Innovative Ap-*
950 *plications of Artificial Intelligence Conference and*
951 *Ninth AAAI Symposium on Educational Advances in*
952 *Artificial Intelligence, AAAI'19/IAAI'19/EAAI'19.*
953 AAAI Press. Event-place: Honolulu, Hawaii, USA.

954 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
955 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
956 Kaiser, and Illia Polosukhin. 2017. Attention is All
957 you Need. In *Advances in Neural Information Pro-*
958 *cessing Systems*, volume 30. Curran Associates, Inc.

959 Bram Willemsen, Dmytro Kalpakchi, and Gabriel
960 Skantze. 2022. Collecting Visually-Grounded Di-
961 alogue with A Game Of Sorts. In *Proceedings of*
962 *the Thirteenth Language Resources and Evaluation*
963 *Conference*, pages 2257–2268, Marseille, France. Eu-
964 ropean Language Resources Association.

965 Licheng Yu, Patrick Poirson, Shan Yang, Alexander C.
966 Berg, and Tamara L. Berg. 2016. Modeling Con-
967 text in Referring Expressions. In *Computer Vision –*
968 *ECCV 2016*, pages 69–85, Cham. Springer Interna-
969 tional Publishing.

970 Hanwang Zhang, Yulei Niu, and Shih-Fu Chang. 2018.
971 *Grounding Referring Expressions in Images by Vari-*
972 *ational Context*. In *2018 IEEE/CVF Conference*
973 *on Computer Vision and Pattern Recognition*, pages
974 4158–4166.

975 Susan Zhang, Stephen Roller, Naman Goyal, Mikel
976 Artetxe, Moya Chen, Shuhui Chen, Christopher De-
977 wan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mi-
978 haylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel
979 Simig, Punit Singh Koura, Anjali Sridhar, Tianlu
980 Wang, and Luke Zettlemoyer. 2022. OPT: Open
981 Pre-trained Transformer Language Models. *_eprint:*
982 2205.01068.

Appendices

A Dialogue Excerpt

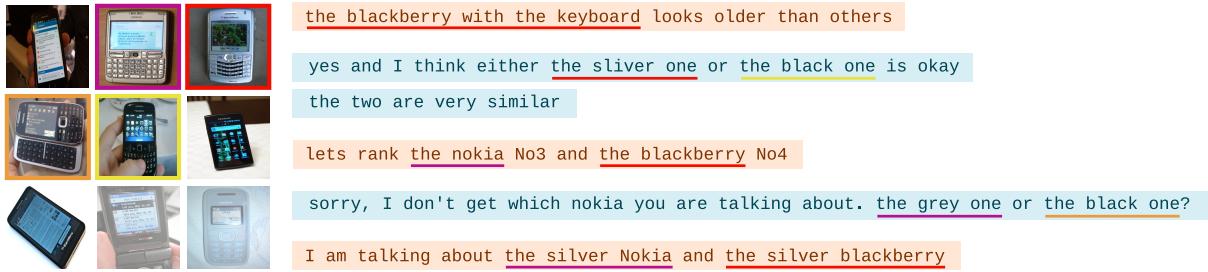


Figure 3: Excerpt of an AGOS dialogue with references to single-image referents underlined; the color indicates the referent. *Note.* The two images that have been ranked successfully at this point in the interaction have a faded appearance.

B Model Input

```

A: Which one do you think is the smallest?
B: the black ear beagle! -> the beagle dog with black ears
B: looking to the left -> the beagle dog with black ears looking to the left
A: The one ears longer than the head? -> the beagle dog with black ears longer than its head looking to the left
B: yes

```

Figure 4: Excerpt of an AGOS dialogue with messages paired with manually constructed ground truth referent descriptions. Mentions are in purple and made bold for illustrative purposes. Ground truth labels for the referent denoted by the mention in green.

```

"Speaker" token + task instructions <N> You are looking to hang a picture on your wall,
                                         but you have no hammer at your disposal to put the nail in the wall.
                                         Which of these phones would you consider most suitable to use as an impromptu hammer and why?
Speaker token + utterance <A> both of them seem quite hard to use. Yeah let's choose bigger one first
Speaker token + utterance <B> I heard Nokia is pretty solid
Speaker token + utterance with a marked mention <A> yeah I was thinking the same
Inference token <B> Maybe >>the one with rubbers<< ? Easier to grab
Ground truth referent description for the marked mention ->
                                         the phone with rubber END

```

Figure 5: Sample input to LLM, deconstructed for demonstration purposes (the sample is otherwise a flat sequence of tokens). Left (text in purple): explanation of input; right (text in black): input. *Note.* The ground truth is only available to the model during training, not during inference.

C Additional VLM Results

| | CLIP-B | | CLIP-L | | ALIGN | | BLIP | |
|--------------|---------------|----------|---------------|----------|--------------|----------|-------------|----------|
| | 3 | 7 | 3 | 7 | 3 | 7 | 3 | 7 |
| Random | .11 | .11 | .11 | .11 | .11 | .11 | .11 | .11 |
| Mention | .36 | .36 | .44 | .44 | .44 | .44 | .40 | .40 |
| Substitution | .42 | .42 | .51 | .51 | .52 | .52 | .50 | .50 |
| coref, chain | .41 | .42 | .49 | .49 | .47 | .46 | .47 | .46 |
| coref, set | .42 | .41 | .48 | .48 | .49 | .48 | .47 | .47 |
| GT, chain | .45 | .47 | .54 | .56 | .53 | .53 | .52 | .54 |
| GT, set | .46 | .48 | .54 | .56 | .54 | .54 | .53 | .55 |
| GT, manual | .47 | .48 | .53 | .55 | .58 | .59 | .55 | .57 |
| GPT-2 | .41 | .38 | .46 | .43 | .49 | .44 | .47 | .43 |
| GPT-3 | .44 | .45 | .52 | .52 | .54 | .55 | .52 | .52 |

Table 3: Cross-validated image retrieval accuracy averaged over five folds for single-image referents (candidate set not reduced). *Note.* Scores are rounded to the nearest hundredth. GT = ground truth; CLIP-B = CLIP ViT-B/16; CLIP-L = CLIP ViT-L/14.

| | CLIP-B | | CLIP-L | | ALIGN | | BLIP | |
|--------------|--------|-----|--------|-----|-------|-----|------|-----|
| | 3 | 7 | 3 | 7 | 3 | 7 | 3 | 7 |
| Random | .22 | .22 | .22 | .22 | .22 | .22 | .22 | .22 |
| Mention | .49 | .49 | .55 | .55 | .56 | .56 | .54 | .54 |
| Substitution | .55 | .55 | .62 | .62 | .64 | .64 | .64 | .64 |
| coref, chain | .54 | .54 | .61 | .61 | .60 | .60 | .61 | .61 |
| coref, set | .54 | .53 | .61 | .60 | .61 | .61 | .61 | .61 |
| GT, chain | .58 | .59 | .66 | .67 | .66 | .67 | .66 | .68 |
| GT, set | .58 | .60 | .66 | .68 | .67 | .67 | .66 | .69 |
| GT, manual | .59 | .60 | .64 | .66 | .69 | .70 | .69 | .70 |
| GPT-2 | .53 | .49 | .58 | .54 | .61 | .58 | .60 | .58 |
| GPT-3 | .57 | .58 | .63 | .63 | .66 | .66 | .67 | .68 |

Table 4: Cross-validated image retrieval accuracy averaged over five folds for single-image referents (candidate set reduced). *Note.* Scores are rounded to the nearest hundredth. GT = ground truth; CLIP-B = CLIP ViT-B/16; CLIP-L = CLIP ViT-L/14.

| | Cars | | Dogs | | Paintings | | Pastries | | Phones | |
|--------------|------|-----|------|-----|-----------|-----|----------|-----|--------|-----|
| | 3 | 7 | 3 | 7 | 3 | 7 | 3 | 7 | 3 | 7 |
| Random | .22 | .22 | .22 | .22 | .22 | .22 | .22 | .22 | .22 | .22 |
| Mention | .52 | .52 | .62 | .62 | .60 | .60 | .61 | .61 | .58 | .58 |
| Substitution | .63 | .63 | .70 | .70 | .70 | .70 | .68 | .68 | .67 | .67 |
| coref, chain | .59 | .60 | .69 | .69 | .66 | .67 | .67 | .68 | .63 | .63 |
| coref, set | .60 | .57 | .68 | .68 | .69 | .68 | .69 | .70 | .62 | .62 |
| GT, chain | .66 | .66 | .76 | .78 | .72 | .74 | .75 | .78 | .71 | .69 |
| GT, set | .66 | .65 | .74 | .77 | .73 | .78 | .76 | .80 | .73 | .73 |
| GT, manual | .64 | .63 | .75 | .78 | .77 | .80 | .70 | .72 | .74 | .74 |
| GPT-2 | .62 | .62 | .67 | .62 | .67 | .62 | .63 | .61 | .57 | .50 |
| GPT-3 | .63 | .63 | .75 | .78 | .70 | .70 | .68 | .72 | .70 | .69 |

Table 5: Cross-validated image retrieval accuracy per fold for single-image referents (candidate set not reduced). *Note.* Scores shown are of VLM that averaged best performance (BLIP-2). Scores are rounded to the nearest hundredth. GT = ground truth; CLIP-B = CLIP ViT-B/16; CLIP-L = CLIP ViT-L/14.

| | Accuracy | | MRR | | NDCG | |
|--------------|----------|-----|-----|-----|------|-----|
| | 3 | 7 | 3 | 7 | 3 | 7 |
| Random | .11 | .11 | .22 | .22 | - | - |
| Mention | .47 | .47 | .63 | .63 | .72 | .72 |
| Substitution | .55 | .55 | .71 | .71 | .78 | .78 |
| coref, chain | .53 | .51 | .69 | .68 | .76 | .76 |
| coref, set | .53 | .51 | .69 | .68 | .77 | .76 |
| GT, chain | .60 | .61 | .75 | .76 | .81 | .82 |
| GT, set | .60 | .62 | .75 | .77 | .81 | .83 |
| GT, manual | .63 | .64 | .76 | .78 | .82 | .83 |
| GPT-2 | .54 | .48 | .69 | .65 | .77 | .73 |
| GPT-3 | .60 | .60 | .74 | .74 | .80 | .81 |

Table 6: Cross-validated image retrieval performance averaged over five folds for single-image referents (candidate set not reduced). *Note.* Scores shown are of VLM that averaged best performance (BLIP-2). Scores are rounded to the nearest hundredth. GT = ground truth.

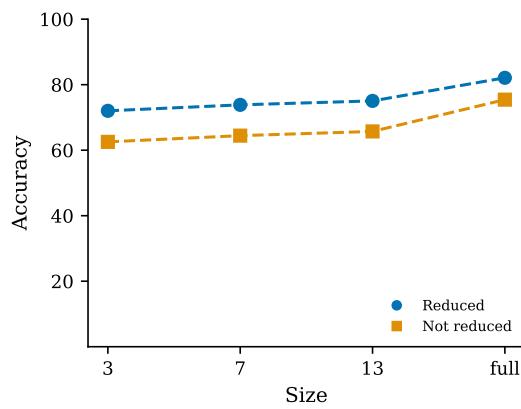


Figure 6: Text-image retrieval accuracy as a function of the size of the context window. Results are shown for BLIP-2 on the manually constructed ground truth referent descriptions based on their respective windows. We show results for both the reduced candidate set and the not reduced candidate set.