

# ALOHa: A New Measure for Hallucination in Captioning Models

Anonymous ACL submission

## Abstract

Despite recent advances in multimodal pre-training for visual description, state-of-the-art models still produce captions containing errors, such as hallucinating objects that are not present in a scene. The existing prominent metric for object hallucination, CHAIR, is limited to a fixed set of MS COCO objects and synonyms. In this work, we propose a modernized open-vocabulary metric, ALOHa, which leverages large language models (LLMs) to measure object hallucinations. Specifically, we use an LLM to extract groundable objects from a candidate caption, measure their semantic similarity to reference objects from captions and/or object detections, and use Hungarian matching to produce a final hallucination score. We show that ALOHa correctly identifies 13.6% more hallucinated objects than CHAIR on HAT, a new gold-standard subset of MS COCO Captions annotated for hallucinations, and 30.8% more on nocaps, where objects extend beyond MS COCO categories.

## 1 Introduction and Background

In recent years, vision-language models have demonstrated remarkable performance. Unfortunately, even state-of-the-art models for visual description still generate captions with object hallucinations – objects or entities that are present in the caption yet are not explicitly supported by visual evidence in the image (Dai et al., 2023). In order to reduce the occurrence of object hallucinations in vision-language models, it is helpful to understand and quantify the problem through *reliable*, *localizable*, and *generalizable* measures of object hallucination. *Reliable* measures are capable of correctly indicating if a given caption contains an object hallucination. *Localizable* measures are capable of indicating which object in a particular caption is hallucinated. *Generalizable* measures are capable of evaluating captions from a wide range of input datasets, across a wide range of object and entity categories.

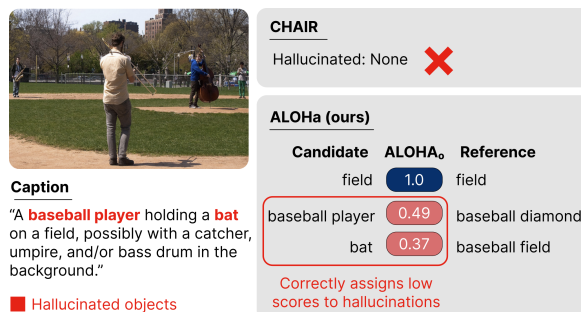


Figure 1: (Top) The SOTA prior object hallucination metric, CHAIR, is limited to MS COCO objects, and fails to detect the hallucinations in this image caption while ALOHa (ours, bottom) correctly assigns low similarity scores to the hallucinations “baseball player” and “bat”. ALOHa does not penalize the caption for “catcher”, “umpire”, and “bass drum”, as the caption indicates uncertainty of their presence.

Recent works that measure object hallucinations in generated text generally fall into two categories: measures that find hallucinations by explicitly matching from a set of objects, and measures that compute distances between latent image and/or text embeddings, indicating a hallucination if the embeddings are too distant. In the first category, CHAIR (Rohrbach et al., 2018) is a measure that explicitly extracts objects from candidate sentences using simple string matching against MS COCO classes and a small set of synonyms. It compares these extracted objects against the ground truth detections, and objects extracted from the ground truth reference captions. CHAIR is both reliable, as string matching on a fixed set of objects is accurate and consistent, and localizable, as individual non-matching strings are identified. However, as seen in Figure 1, CHAIR is not generalizable, as it can only handle a fixed set of predetermined objects. Other uni-modal measures in this category include those for abstractive summarization (Durmus et al., 2020; Kryscinski et al., 2020; Maynez et al., 2020; Son et al., 2022; Sridhar and Visser, 2022; Yuan et al., 2021), dialogue (Huang et al., 2022; Shuster et al., 2021), and structured knowledge (Dhingra et al., 2019). These often generalize poorly to

067 vision-language tasks as they require grounding the  
068 generated text into inputs of the same modality.

069 In the second category, CLIPScore (Hessel  
070 et al., 2021) employs CLIP (Radford et al., 2021)  
071 embeddings to assess image-text matches. While  
072 it is generalizable and reliable, it lacks localization  
073 as it does not pinpoint incorrect spans of text.  
074 CLIPBERTS (Wan and Bansal, 2022) and Ref-  
075 CLIPScore (an extension of CLIPScore accounting  
076 for reference captions) face similar limitations.

077 POPE (Li et al., 2023) evaluates vision-language  
078 models’ likelihood to hallucinate objects with  
079 machine-generated queries consisting of samples  
080 extracted from both reference object detections and  
081 nonexistent objects. The POPE approach, while a  
082 useful object hallucination score for model compar-  
083 ison, addresses a fundamentally different problem  
084 from that which we investigate here – it measures  
085 how often *models* hallucinate rather than localizes  
086 and detects hallucinations within *a single caption*.

087 Inspired by recent successes using LLMs for  
088 evaluation in language-only tasks (Zhang et al.,  
089 2020; Yuan et al., 2021; Bubeck et al., 2023; Chiang  
090 et al., 2023; Zheng et al., 2023), we introduce  
091 Assessment with Language models for Object  
092 Hallucination (ALOHa), a modernized measure  
093 for object hallucination detection that is *reliable*,  
094 *localizable*, and *generalizable*. ALOHa extends the  
095 reliability and localization of CHAIR to new input  
096 domains by leveraging in-context learning of LLMs  
097 combined with semantically-rich text embeddings  
098 for object parsing and matching (Figure 1).

099 For a given image caption, we generate two  
100 measures: **ALOHa<sub>o</sub>**, a numeric score for each  
101 object rating the degree to which that object is a hal-  
102 lucination, and **ALOHa**, an aggregated score rating  
103 the degree to which the whole caption contains a  
104 hallucination. We demonstrate the performance of  
105 ALOHa on a new gold-standard dataset of image  
106 hallucinations, HAT, and show that ALOHa is more  
107 accurate than CLIPScore at detecting object halluci-  
108 nations, and more accurate than CHAIR at correctly  
109 localizing those hallucinations. We conclude by  
110 demonstrating that ALOHa remains reliable and  
111 localizable when generalizing to out of domain data.

## 112 2 ALOHa: Reliable, Localizable, and 113 Generalizable Hallucination Detection

114 ALOHa produces numeric scores rating the degree  
115 of hallucination for each object in a candidate

116 caption as well as an overall caption score, given a  
117 set of ground-truth reference captions and predicted  
118 (or ground truth) image object detections. ALOHa  
119 consists of three stages (Figure 2). (1) Objects  
120 are extracted from the image, reference set, and  
121 candidate caption using a combination of an object  
122 detector and LLM. (2) We filter the object sets  
123 and compute semantic representations of each  
124 object. (3) We compute a maximum-similarity  
125 linear assignment between candidate and reference  
126 objects. The scores from each of the pairs in the  
127 linear assignment, which we call ALOHa<sub>o</sub>, measure  
128 the degree of hallucination for each of the candidate  
129 objects. The minimum similarity in this linear  
130 assignment (the ALOHa score) measures the degree  
131 of hallucination of the caption.

132 **(1) Extracting objects from candidates, refer-**  
133 **ences, and images:** Parsing visually grounded  
134 objects in a caption in an open-domain context is a  
135 surprisingly difficult task. CHAIR (Rohrbach et al.,  
136 2018) relies on a fixed set of MS COCO objects and  
137 synonyms, requiring considerable effort to extend to  
138 other datasets, and sometimes failing at ambiguous  
139 parses (such as mistaking the adjective “orange” for  
140 a noun). SPICE (Anderson et al., 2016) relies on  
141 standard grammar-based object parsing, which can  
142 have similar issues, as purely text-based methods  
143 fall short at identifying which nouns are *visual* – for  
144 instance, avoiding “picture” and “background” in  
145 Figure 2. Captions may also indicate uncertainty  
146 around object presence, such as “a bowl or plate”,  
147 or “a dog biting something, possibly a Frisbee.” We  
148 aim to handle such uncertain objects to avoid unfair  
149 hallucination penalties.

150 With the understanding that open-domain  
151 parsing is the primary factor in CHAIR’s lack  
152 of generalization, we leverage the capability of  
153 zero-shot in-context learning in large language  
154 models. Following Brown et al. (2020), we use an  
155 LLM (ChatGPT, OpenAI (2022)) along with the  
156 prompt given in Appendix A to turn the parsing  
157 task into a language completion task easily solvable  
158 by an LLM. We encourage the LLM to extract  
159 visual objects in the scene, consisting primarily of  
160 noun phrases (including any attributes, such as “big  
161 dog” and “purple shirt”), from the candidate and  
162 reference captions. We run the LLM against the  
163 candidate caption to produce the unfiltered object  
164 set  $\mathcal{C}$ , and again for the corresponding reference  
165 captions to produce object set  $\mathcal{R}$ . To extract objects  
166 from the image context, similar to CHAIR, we

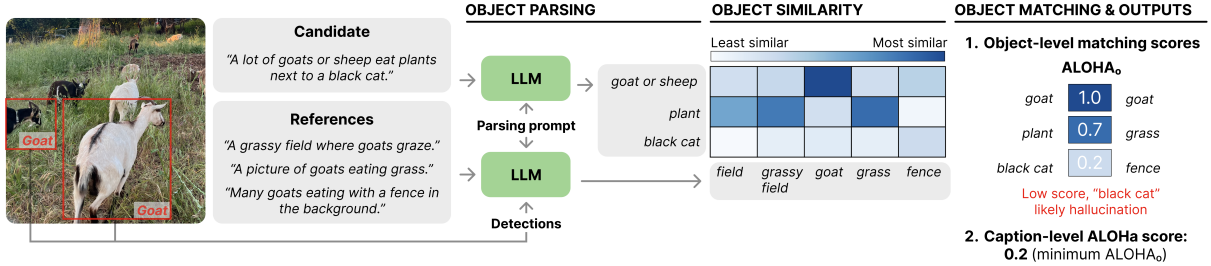


Figure 2: Overview of ALOHa. We prompt an LLM to extract visually-grounded nouns from a candidate machine-generated description and a set of references. We consider uncertain language (e.g., “goat or sheep”), add reference objects with and without modifiers (e.g., both “field” and “grassy field”), and avoid non-visual nouns (e.g., “picture” and “background”). Then, we compute a maximum-similarity linear assignment between candidate and reference object sets, the weights of which form the  $\text{ALOHa}_0$ . Matched pairs with low  $\text{ALOHa}_0$  are likely hallucinations (e.g., “black cat”,  $\text{ALOHa}_0 = 0.2$ ). We additionally output the minimum  $\text{ALOHa}_0$  as a caption-level ALOHa score.

167 augment the set of reference objects with objects  
 168 detected directly from the image using DETR  
 169 (Carion et al., 2020) fine-tuned on MS COCO.

170 **(2) Object filtering:** We further refine candidate  
 171 ( $\mathcal{C}$ ) and reference ( $\mathcal{R}$ ) object sets to better reflect  
 172 specific challenges of object hallucination detection.  
 173 Ideally, hallucination measures should penalize  
 174 specificity when candidate attributes are not  
 175 supported by references (e.g., if “purple shirt”  $\in \mathcal{C}$ ,  
 176 yet “white shirt”  $\in \mathcal{R}$ ), but should not penalize  
 177 generality (e.g., “shirt”  $\in \mathcal{C}$ , yet “white shirt”  $\in \mathcal{R}$ ).  
 178 Thus, we use spaCy (Honnibal et al., 2020a) to  
 179 augment  $\mathcal{R}$  with the root nouns from each *reference*  
 180 noun phrase, but leave the candidates unchanged.

181 Beyond specificity, captions may also express  
 182 uncertainty about the presence of objects in an  
 183 image. For conjunctions (e.g., “fork or knife”),  
 184 we aim to avoid unfair penalties if at least one  
 185 of the objects is grounded. ALOHa considers all  
 186 combinations of selecting a single object from each  
 187 conjunction, denoted as  $\mathcal{C}_{\{1\dots M\}}$  and  $\mathcal{R}_{\{1\dots N\}}$  (e.g.,  
 188 “fork”  $\in \mathcal{R}_0$  and “knife”  $\in \mathcal{R}_1$ ). Additionally, we  
 189 prompt the LLM to indicate uncertain grounding  
 190 by including “possibly” after the object (e.g., “there  
 191 may be a Frisbee” becomes “Frisbee (possibly)”) and  
 192 we remove uncertain objects from  $\mathcal{C}_i$  to  
 193 avoid penalties while maintaining them in  $\mathcal{R}_j$  for  
 194 maximum coverage of more general objects.

195 **(3) Object Matching:** Once we have extracted and  
 196 parsed the candidate and reference object sets, we  
 197 aim to measure the degree of hallucination for each  
 198 candidate object. While we could match objects  
 199 based on string alone (resulting in a binary decision),  
 200 as does CHAIR, often it is useful to understand  
 201 a continuous scale of hallucination – e.g., for a  
 202 reference object “dog”, hallucinating “wolf” should  
 203 be penalized less than “potato.” To capture this scale  
 204 of semantic similarity, for each object text  $o$ , we

205 generate  $o_{\text{emb}} = \phi(o) \in \mathbb{R}^K$ , where  $\phi$  is a semantic  
 206 text embedding model. In our work, we use  
 207 S-BERT (Reimers and Gurevych, 2019). We then  
 208 compute a similarity score for each pair of objects  
 209 (usually the cosine similarity, see Appendix B.3).  
 210 For each ( $\mathcal{C}_i, \mathcal{R}_j$ ) pair, we store these scores in  
 211 a similarity matrix  $\mathcal{S}_{i,j} \in [0, 1]^{|\mathcal{C}_i| \times |\mathcal{R}_j|}$ . We then  
 212 use the Hungarian method (Kuhn, 1955) to find  
 213 an optimal maximum-similarity assignment  $\mathcal{M}_{i,j}$   
 214 between candidate and reference sets of objects.

215 To determine the  $\text{ALOHa}_0$  score for each object,  
 216 we take the maximum score across all possible  
 217 parsings, giving the candidate caption the benefit  
 218 of the doubt, for an object  $c \in \mathcal{C}_i$

$$\text{ALOHa}_0(c) = \max_{i,j} w_{c_i,j} \in \mathcal{M}_{i,j} \quad (1) \quad 219$$

220 While  $0 \leq \text{ALOHa}_0 \leq 1$  indicates the degree of  
 221 hallucination for each object, we also want to  
 222 indicate if an entire caption contains a hallucination.  
 223 We thus define:

$$\text{ALOHa} = \min_{c \in \mathcal{C}} \text{ALOHa}_0(c) \quad (2) \quad 224$$

225 We choose the minimum as the presence of *any*  
 226 hallucinated object indicates that the full caption is  
 227 a hallucination, and even several correct detections  
 228 should not compensate for a hallucination.

### 229 3 Evaluation & Discussion

230 **HAT:** To promote the development of high-quality  
 231 methods for hallucination detection, we collect  
 232 and release HAT (HALLUCINATION TEST), a dataset of  
 233 labeled hallucinations in captions. HAT consists of  
 234 490 samples (90 validation and 400 test) labeled by  
 235 in-domain experts for hallucination on both a word  
 236 level and caption level (See Appendix D). Measures  
 237 are evaluated on two metrics: Average Precision  
 238 (AP) and Localization Accuracy (LA). The AP

Method	LA	AP
Baseline (Majority Vote)	-	33.75
CHAIRs	6.70	36.85
CLIPScore	-	40.10
RefCLIPScore	-	48.40
ALOHa (No Detections)	19.55	48.40
ALOHa (Oracle Detections)	19.55	47.86
ALOHa (DETR Detections)*	<u>20.30</u>	<u>48.62</u>
ALOHa (Oracle+DETR Detections)	<b>21.05</b>	<b>48.78</b>

Table 1: Test set performance for binary hallucination detection on HAT. LA: Localization Accuracy. AP: Average Precision. \* indicates the version of ALOHa used throughout this paper, unless noted otherwise. Oracle detection are human-generated reference detections.

of the method measures reliability, and is defined as how well the measure identifies captions with hallucinations. For CHAIR, decisions are binary, so AP = accuracy. For ALOHa, AP is the weighted mean of precisions across all thresholds. The LA, measured on samples containing hallucinations in HAT, measures localization and is defined as the accuracy of correctly indicating *which* of the specific objects were hallucinated. For CHAIR, a hallucination is correctly localized when at least one detected string mismatch is a hallucination, and for ALOHa when the minimum ALOHa<sub>o</sub> score corresponds to a hallucinated object.

ALOHa’s performance on HAT is shown in Table 1. On AP, ALOHa with DETR detections outperforms both CHAIR and CLIPScore by 11.8% and 8.5% respectively. RefCLIPScore attains a similar AP; however, is not localizable. ALOHa achieves more than twice the LA on HAT CHAIR, a particularly challenging task as HAT includes non-object hallucinations, such as incorrect verbs or relations (see Figure A7). Table 1 further ablates the choice of image detections, and indicates that ALOHa is robust to missing detections.

**FOIL object hallucinations:** To indicate generalizability we evaluate our method on two machine-generated object hallucination datasets. FOIL (Shekhar et al., 2017) contains MS COCO images, where objects are randomly replaced with similar ones (e.g., “bus“ and “car”), and nocaps-FOIL, a similar dataset that we construct on the nocaps dataset (Agrawal et al., 2019) for novel object captioning beyond MS COCO (see Appendix D.1). While both methods are strong on the FOIL dataset, CHAIR fails to transfer to the nocaps-FOIL dataset, as the object set becomes out of scope. CHAIR achieves an AP of only 58.33 (only slightly better than chance) and LA of 14.42, compared to

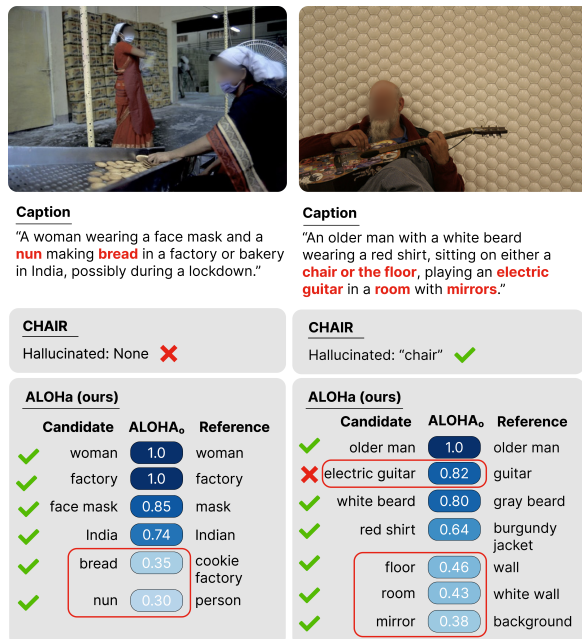


Figure 3: Qualitative Flickr30k examples. (Left) ALOHa correctly assigns low scores to the hallucinated “nun” and “bread”, whereas CHAIR does not detect any hallucinations. (Right) Although ALOHa assigns high similarity between the hallucinated “electric guitar” and reference “(acoustic) guitar”, it assigns low scores to the other 3 hallucinations. CHAIR detects only the hallucination “chair”, missing the others.

ALOHa’s AP of 69.52, and LA of 45.17 (213% relative improvement). See Appendix C.2 for details.

**Flickr30k:** In Figure 3 and Figure A5, we visualize the behavior of CHAIR and ALOHa on several Flickr30k samples (Young et al., 2014), using captions generated by a recent captioning model (Chan et al., 2023) that often produces complex captions with phrases expressing uncertainty.

**Additional results:** As LLMs can hallucinate themselves (Bubeck et al., 2023), we annotate the parsing error rate on HAT in Table A1 and find that GPT-3.5 introduces extraneous objects in 2.97% of samples. In Appendix B we investigate the choice of LLM, similarity measure, and parsing approach.

## 4 Conclusion

This paper introduces ALOHa, a scalable LLM-augmented metric for open-vocabulary object hallucination. ALOHa correctly identifies 13.6% more hallucinated objects on HAT and 31% on nocaps-FOIL than CHAIR. ALOHa represents an important modernization of caption hallucination metrics, and detecting complex hallucinations in actions, quantities, and abstract concepts remains an exciting and challenging task for future exploration.

## 5 Limitations

While ALOHa represents a strong step towards open-domain localized hallucination detection, it comes with several limitations which we discuss in this section.

**Non-determinism** A primary concern with using large language models for an evaluation measure is the natural nondeterminism that comes with them. While in theory language models sampled at a temperature of zero (as we do in this work) are deterministic, it is well documented that small random fluctuations can still occur (OpenAI, 2023). Beyond random fluctuations, the availability of language models long-term can impact the reproducibility of the measure. In this work, we primarily rely on closed source language models, which can change or become unavailable without notice. In Table A1, we demonstrate that ALOHa still functions with open source models such as Koala (Geng et al., 2023), however the performance is significantly degraded due to the parsing capabilities of the model. With time, and more powerful open source LLMs, this will become less of an issue, however relying on a nondeterministic metric for comparative evaluation can easily become a liability.

**Availability of Reference Captions (Reference-Free vs. Reference-Based Measures)** One of the primary limitations of the ALOHa evaluation method is the requirement that reference captions are available for the evaluation dataset (an issue shared by CHAIR). Not only must reference captions be available, but they also must sufficiently cover the salient details in the reference image. When the references are impoverished (as can easily happen with a single reference sentence (Chan et al., 2023)) or when there are no references, and ALOHa must rely entirely on detections, the method under-performs more general methods such as CLIPScore which are reference free, and rely on a large pre-training dataset to encode vision and language correspondences. We strongly believe that the area of reference-free localized hallucination detection is an important area of future research; how can we leverage the tools from large vision and language pre-training in a localized way to understand and interpret where hallucinations lie in hallucinated text? That being said, there is also a place for reference-based measures, as reference-based measures focus on what *humans*

believe to be salient details in the image, whereas reference-free measures always rely on downstream models which *approximate* what humans believe to be important. This means that reference-based measures can often transfer better to new domains than reference-free measures, which often must be trained/fine-tuned in-domain with human-labeled data to achieve strong performance.

**General costs associated with LLMs** The use of large language models for any task incurs significant compute, monetary, environmental, and human costs. ALOHa is a significantly slower evaluation measure than methods like CHAIR (however not that much less efficient than CLIPScore), leading to increased power consumption, and cost during evaluation. In addition, the models that we rely on are generally closed source, and represent a non-trivial monetary expenditure (Experiments in this paper, including ablations, testing, and prototyping required approximately \$120 USD in API fees). Such factors can be limiting to researchers who wish to evaluate large datasets, however we hope that with the advent of larger open source models, and continued investment in hardware and systems research, the cost will decrease significantly. Beyond compute and financial costs, there are environmental and human costs associated with using large language models for evaluation, see Bender et al. (2021) for a detailed discussion of these factors.

**Limited Control of Bias** In this work, we do not evaluate the performance of ALOHa on Non-English data, nor do we explicitly control for or measure bias in the creation of HAT (Which is a labeled subset, randomly selected of the MS COCO dataset), or the Nocaps-FOIL dataset (which operates on the same samples as the Nocaps validation dataset). While HAT is a subset of the common MS COCO dataset, we recognize that the creation of such potentially biased datasets has the potential to lead researchers to engineer features and methods which are unintentionally biased against underrepresented groups. We aim to address these shortcomings in the next iteration of HAT, which will not only contain out of domain data for MS COCO trained models, but also aims to better control for bias in the underlying image and caption data. Note that our work, including HAT, is intended for research purposes.

## References

- 398 Harsh Agrawal, Peter Anderson, Karan Desai, Yufei  
399 Wang, Xinlei Chen, Rishabh Jain, Mark Johnson,  
400 Dhruv Batra, Devi Parikh, and Stefan Lee. 2019.  
401 [nocaps: novel object captioning at scale](#). In *2019*  
402 *IEEE/CVF International Conference on Computer*  
403 *Vision, ICCV 2019, Seoul, Korea (South), October*  
404 *27 - November 2, 2019*, pages 8947–8956. IEEE.
- 405 Peter Anderson, Basura Fernando, Mark Johnson, and  
406 Stephen Gould. 2016. Spice: Semantic propositional  
407 image caption evaluation. In *European conference*  
408 *on computer vision*, pages 382–398. Springer.
- 409 Emily M Bender, Timnit Gebru, Angelina McMillan-  
410 Major, and Shmargaret Shmitchell. 2021. On the  
411 dangers of stochastic parrots: Can language models be  
412 too big? In *Proceedings of the 2021 ACM conference*  
413 *on fairness, accountability, and transparency*, pages  
414 610–623.
- 415 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie  
416 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind  
417 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
418 Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
419 Gretchen Krueger, Tom Henighan, Rewon Child,  
420 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,  
421 Clemens Winter, Christopher Hesse, Mark Chen, Eric  
422 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,  
423 Jack Clark, Christopher Berner, Sam McCandlish,  
424 Alec Radford, Ilya Sutskever, and Dario Amodei.  
425 2020. [Language models are few-shot learners](#). In  
426 *Advances in Neural Information Processing Systems*  
427 *33: Annual Conference on Neural Information*  
428 *Processing Systems 2020, NeurIPS 2020, December*  
429 *6-12, 2020, virtual*.
- 430 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan,  
431 Johannes Gehrike, Eric Horvitz, Ece Kamar, Peter Lee,  
432 Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023.  
433 [Sparks of artificial general intelligence: Early exper-](#)  
434 [iments with gpt-4](#). *ArXiv preprint*, abs/2303.12712.
- 435 Nicolas Carion, Francisco Massa, Gabriel Synnaeve,  
436 Nicolas Usunier, Alexander Kirillov, and Sergey  
437 Zagoruyko. 2020. End-to-end object detection with  
438 transformers. In *Computer Vision–ECCV 2020:*  
439 *16th European Conference, Glasgow, UK, August*  
440 *23–28, 2020, Proceedings, Part I 16*, pages 213–229.  
441 Springer.
- 442 David M Chan, Austin Myers, Sudheendra Vijaya-  
443 narasimhan, David A Ross, and John Canny. 2023.  
444 [Ic3: Image captioning by committee consensus](#).  
445 *ArXiv preprint*, abs/2302.01328.
- 446 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,  
447 Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan  
448 Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al.  
449 2023. Vicuna: An open-source chatbot impressing  
450 gpt-4 with 90%\* chatgpt quality. See [https://vicuna.](https://vicuna.lmsys.org)  
451 [lmsys.org](https://vicuna.lmsys.org) (accessed 14 April 2023).
- 452 Wenliang Dai, Zihan Liu, Ziwei Ji, Dan Su, and Pascale  
453 Fung. 2023. [Plausible may not be faithful: Probing](#)  
[object hallucination in vision-language pre-training](#).  
In *Proceedings of the 17th Conference of the Euro-*  
*pean Chapter of the Association for Computational*  
*Linguistics*, pages 2136–2148, Dubrovnik, Croatia.  
Association for Computational Linguistics.
- Bhuwan Dhingra, Manaal Faruqui, Ankur Parikh,  
Ming-Wei Chang, Dipanjan Das, and William Cohen.  
2019. [Handling divergent reference texts when](#)  
[evaluating table-to-text generation](#). In *Proceedings of*  
*the 57th Annual Meeting of the Association for Com-*  
*putational Linguistics*, pages 4884–4895, Florence,  
Italy. Association for Computational Linguistics.
- Esin Durmus, He He, and Mona Diab. 2020. [FEQA: A](#)  
[question answering evaluation framework for faithful-](#)  
[ness assessment in abstractive summarization](#). In *Pro-*  
*ceedings of the 58th Annual Meeting of the Associa-*  
*tion for Computational Linguistics*, pages 5055–5070,  
Online. Association for Computational Linguistics.
- Xinyang Geng, Arnav Gudibande, Hao Liu, Eric  
Wallace, Pieter Abbeel, Sergey Levine, and Dawn  
Song. 2023. [Koala: A dialogue model for academic](#)  
[research](#). Blog post.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan  
Le Bras, and Yejin Choi. 2021. [CLIPScore: A](#)  
[reference-free evaluation metric for image captioning](#).  
In *Proceedings of the 2021 Conference on Empirical*  
*Methods in Natural Language Processing*, pages  
7514–7528, Online and Punta Cana, Dominican  
Republic. Association for Computational Linguistics.
- Matthew Honnibal, Ines Montani, Sofie Van Landeghem,  
and Adriane Boyd. 2020a. [spacy: Industrial-strength](#)  
[natural language processing in python](#).
- Matthew Honnibal, Ines Montani, Sofie Van Landeghem,  
and Adriane Boyd. 2020b. [spacy: Industrial-strength](#)  
[natural language processing in python](#), zenodo, 2020.
- Sicong Huang, Asli Celikyilmaz, and Haoran Li. 2022.  
[Ed-faith: Evaluating dialogue summarization on](#)  
[faithfulness](#). *ArXiv preprint*, abs/2211.08464.
- Wojciech Kryscinski, Bryan McCann, Caiming Xiong,  
and Richard Socher. 2020. [Evaluating the factual](#)  
[consistency of abstractive text summarization](#). In  
*Proceedings of the 2020 Conference on Empir-*  
*ical Methods in Natural Language Processing*  
*(EMNLP)*, pages 9332–9346, Online. Association for  
Computational Linguistics.
- Harold W Kuhn. 1955. The hungarian method for  
the assignment problem. *Naval research logistics*  
*quarterly*, 2(1-2):83–97.
- Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper  
Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab  
Kamali, Stefan Popov, Matteo Mallocci, Alexander  
Kolesnikov, et al. 2020. The open images dataset v4:  
Unified image classification, object detection, and  
visual relationship detection at scale. *International*  
*Journal of Computer Vision*, 128(7):1956–1981.



624 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.  
625 Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with BERT](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

630 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan  
631 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,  
632 Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). *ArXiv preprint*, abs/2306.05685.

635 Deyao Zhu, Jun Chen, Kilichbek Haydarov, Xiaoqian  
636 Shen, Wenxuan Zhang, and Mohamed Elhoseiny.  
637 2023. [Chatgpt asks, blip-2 answers: Automatic questioning towards enriched visual descriptions](#). *ArXiv preprint*, abs/2303.06594.

## 640 Appendix

641 **Appendix A** describes the prompt of the language  
642 model, including the exact language used, the  
643 design choices, and the in-context examples.

644 **Appendix B** describes several explorations of  
645 the hyperparameters, including the language  
646 model chosen for ALOHa, the semantic  
647 embedding method, and the parsing approach.

648 **Appendix C** contains additional detailed results  
649 and experimental details for experiments in  
650 the paper.

651 **Appendix D** describes the datasets that we col-  
652 lected and constructed, including HAT and  
653 nocaps-FOIL.

## 654 A Prompt

655 The choice of prompt for a large language model  
656 using in-context learning is critical to the perfor-  
657 mance of the model. Each component of the prompt  
658 has some ability to shape the downstream language  
659 distribution. In this work, we use the prompt shown  
660 in [Figure A1](#). This prompt has several rules, which  
661 we discuss here.

662 **Attributes:** We ask that the language model  
663 include all attributes attached to the object if  
664 they are present. By doing so, we can catch  
665 hallucinations such as those shown in [Figure 3](#),  
666 where “electric guitar” appears in the candidate, but  
667 an acoustic guitar is shown in the image. Attributes  
668 are handled differently between reference captions  
669 and candidate captions. For reference captions, we  
670 add both the object with attributes, and the object

You are an assistant that parses visually present objects from an image caption. Given an image caption, you list ALL the objects visually present in the image or photo described by the captions. Strictly abide by the following rules:

- Include all attributes and adjectives that describe the object, if present
- Do not repeat objects
- Do not include objects that are mentioned but have no visual presence in the image, such as light, sound, or emotions
- If the caption is uncertain about an object, YOU MUST include '(possibly)' after the object
- If the caption thinks an object can be one of several things, include 'or' and all the possible objects
- Always give the singular form of the object, even if the caption uses the plural form

Figure A1: The prompt that we use for parsing objects from both captions and sets of reference captions.

without attributes to the set, so the candidate is not penalized for being more general. For the candidate, however, we add only the object with attributes, so if the candidate produces attributes, they must match with something in the reference set.

**Repeated Objects:** In this work, our primary goal is to determine if a particular object is hallucinated, and not focus on the quantity of hallucinations. Thus, we de-duplicate the object set in both the candidate and reference captions, as well as detections coming from the image. By doing this, we focus on if the objects can possibly exist in the image, rather than focus on getting the exact count, which may be incorrect if a candidate caption mentions the same object more than once (and that object is parsed twice).

**Intangible Object:** In many cases, objects mentioned in the candidate or reference set may be intangible, such as color, light, sound, or emotion. To improve the accuracy of the model, we explicitly suggest that such objects should not be included.



Caption: This image shows two pink roses in a tulip-shaped vase on a wooden kitchen counter, next to a microwave and a toaster oven.

Objects:

- pink rose
- tulip-shaped vase
- wooden kitchen counter
- microwave
- toaster oven

Figure A2: An example of a single-caption parsing result.

**Or/Possibly:** Modern captioning methods such as Chat-Captioner (Zhu et al., 2023) and IC3 (Chan et al., 2023) are capable of encoding uncertainty into their approach through the use of words like “possibly” or “maybe”. Additionally, they may make judgments that are uncertain such as “an apple or an orange.” Existing captioning and hallucination detection measures fail to account for this uncertainty, and match both objects, even though the semantics of the caption suggests that the object is uncertain, or may be one of many objects. To account for this, we encourage the LLM to indicate uncertainty in a fixed way, as well as list multiple alternatives on a single line. We then account for this in our matching method, by giving the candidate the benefit of the doubt, scoring only the best match from an alternative set, and ignoring any uncertainty.

**Singularization:** While it is possible to singularize objects using rule-based methods, rule-based methods struggle with challenging nouns, and we found that in general, the LLM was better at performing the singularization set of the post-processing before object matching.

**A.1 In-Context Examples**

In addition to the core prompt text, we provide several contextual samples, which help with in-context learning (Brown et al., 2020). The contextual samples help to align the label space of the model correctly with the target output distribution (Min et al., 2022). An example of such contexts is given in Figure A2 and Figure A3.

**B Hyperparameter Exploration**

In this section, we explore the choices of hyperparameters for ALOHa including the object parsing, semantic embedding, and language model.

Captions:

- Several people riding on a motorcycle with an umbrella open.
- Couples riding motor cycles carrying umbrellas and people sitting at tables.
- A group of people riding scooters while holding umbrellas.
- Some tables and umbrellas sitting next to a building.
- Pedestrians and motorcyclists near an open outdoor market.

Objects:

- person
- couple
- motorcycle
- umbrella
- table
- scooter
- building
- pedestrian
- motorcyclist
- open outdoor market

Figure A3: An example of a multi-caption parsing result.

**B.1 Object Extraction and Semantic Embedding Methods**

In the primary work, we leverage LLMs (OpenAI, 2023) for object extraction, and a BERT-based model (Reimers and Gurevych, 2019) for semantic word embedding. These are, however, not the only choices that can be made for these two components. In Figure A4, we explore the difference in overall performance on HAT’s validation set when using different combinations of object extraction and semantic embedding. Namely, we compare LLM-based extraction to the parse-tree-based noun extraction in SpaCy (Honnibal et al., 2020b), and compare SentenceTransformer (BERT-Based model, (Reimers and Gurevych, 2019)) to Word2Vec (Mikolov et al., 2018), GPT-3 (Ada) embedding, and string matching (strings are case-normalized and lemmatized). In general, we found that combining LLMs with the SentenceTransformer (BERT-Based) model performed better than other methods, and that fuzzy embedding methods often significantly outperform exact string matching when determining hallucination as judged by human raters. This is generally expected: humans have a wide vocabulary that is poorly captured

692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
  
709  
710  
711  
712  
713  
714  
  
715  
  
716  
717  
718  
719  
720  
721  
722  
  
723  
  
724  
725  
726

727  
728  
  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751

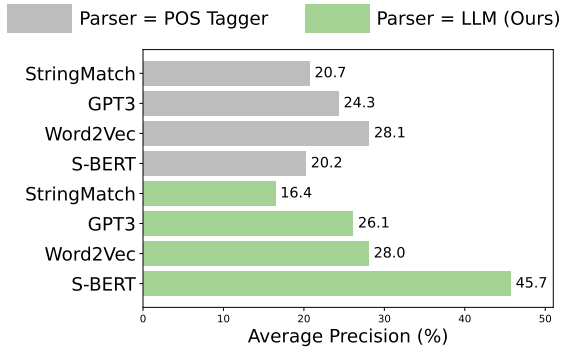


Figure A4: Performance on HAT validation set filtered for hallucinated objects, when comparing embedding methods and object extraction approaches.

by exact string matching. Interestingly, Word2Vec outperforms GPT-3 embeddings. We believe that this is because the GPT-3 embeddings are optimized for sentence-level structures, and may fail to semantically embed single words in a meaningful way.

## B.2 Choice of Large Language Model

The choice of the language model is critical to the overall performance of ALOHa- if the language model does not have sufficient zero-shot parsing capability, it will lead to reduced downstream performance. We investigate the performance of the language model in Table A1 on HAT. In these experiments, we measure the average precision (AP) and LA (see Appendix C.1), as well as the "Parsing error rate" (PER), which is the rate of errors made when parsing objects from reference captions on HAT. We calculate PER (Parsing Error Rate) with manual annotation by taking the fraction of objects output by the LLM that did not exist in the caption (in other words, measuring 1-precision of parsed objects). We additionally annotate and compute the recall - the fraction of objects in the caption that are included in the objects parsed by the LLM. This gives a recall for GPT-3.5 of 98.63%. In these experiments, we find that while Koala (Geng et al., 2023) has strong LA performance on HAT, however ChatGPT (GPT-3.5) (OpenAI, 2023) has both the best average precision, and makes the fewest errors, thus we leverage GPT-3.5 for our primary experiments in the main paper.

## B.3 Semantic Similarity Measure

In ALOHa, we compute the similarity between objects using the cosine distance between embedding vectors generated using the all-MiniLM-L6-v2 S-BERT implementation in the Sentence-

Language Model	LA $\uparrow$	AP $\uparrow$	PER $\downarrow$	PRR $\uparrow$
GPT-3.5	20.30	<b>48.62</b>	<b>2.97</b>	<b>98.63</b>
Claude (Instant)	20.74	41.48	3.31	-
Koala	<b>22.22</b>	38.70	5.07	-

Table A1: Exploration of LLM choice for parsing within ALOHa, on HAT. AP: Average Precision, LA: Localization Accuracy, PER: Parsing Error Rate (%), PRR: Parsing Recall Rate.

Transformers<sup>1</sup> library (Reimers and Gurevych, 2019). While in theory cosine distances should lie in the interval  $[-1,1]$ , in this library, for optimization stability, models are trained with positive samples having similarity 1, and negative samples having similarity 0. This (unintentionally) induces a model which (by optimization) only produces positive cosine similarity scores. ALOHa can still be adapted to negative similarity: our algorithms for maximal assignment and equations 1 and 2 both support negative values (even though they don't appear in this instantiation of the algorithm).

It is further worth noting that S-BERT is not a word similarity measure, and was instead designed to measure distances between sentences. This means that the underlying metric may not be sufficiently optimized at a word level, however, in Figure Figure A4, we give an ablation on the parsing method and examined the effectiveness of different embedding models for semantic similarity. We found among the explored approaches that S-BERT was most effective and that simple word embedding methods such as GLoVe are insufficient. That being said, we acknowledge that leveraging a model trained specifically for semantic similarity between words would be an exciting and powerful extension to the ALOHa framework. With the development of better word embedding models for semantic similarity, we see greater potential in localizing hallucinations with ALOHa.

## C Experimental Details & Additional Experimentation

### C.1 Metrics

We employ several measures in the paper, which we describe in detail here.

**Average Precision** We measure the **Average Precision (AP)** of each hallucination metric to detect sentence-level hallucinations. Specifically,

<sup>1</sup><https://www.sbert.net/>

Method	FOIL				nocaps-FOIL					
	Overall		In-Domain		Near-Domain		Out-Domain		Overall	
	LA	AP	LA	AP	LA	AP	LA	AP	LA	AP
CHAIRs	<b>79.00</b>	<b>92.50</b>	13.47	57.82	17.55	59.14	12.24	58.06	14.42	58.33
CLIPScore	-	76.44	-	<u>71.81</u>	-	<u>70.17</u>	-	<u>78.73</u>	-	<b>73.48</b>
RefCLIPScore	-	<u>80.64</u>	-	<b>79.63</b>	-	<b>78.70</b>	-	<b>85.89</b>	-	<b>81.31</b>
ALOHa	40.00	61.35	<b>47.35</b>	71.80	<b>47.30</b>	66.67	<b>48.84</b>	70.91	<b>45.17</b>	69.52

Table A2: Breakdown of results by domain on nocaps FOIL. AP: Average Precision. LA: Localization Accuracy. Bold and underlined values represent the best and second-best methods respectively.

we label each sample with **1** if it contains a hallucination and **0** otherwise. We then measure AP between those labels and per-sample hallucination measures. For ALOHa, this is:

$$AP = \frac{1}{N} \sum_{i=1}^N \mathbb{I}[\text{label}] \cdot \text{ALOHa}(i) \quad (3)$$

For CHAIR, this is:

$$AP = \frac{1}{N} \sum_{i=1}^N \mathbb{I}[\text{label}] \cdot \mathbb{I}[\text{CHAIR Prediction}] \quad (4)$$

**Localization Accuracy** Localization accuracy (LA) measures the fraction of samples where a metric can correctly identify a hallucinated object, among samples that are known to contain hallucinated objects.

$$LA = \frac{|\{\geq 1 \text{ correctly identified halluc.}\}|}{|\{\geq 1 \text{ halluc.}\}|} \quad (5)$$

A sample receives a LA of 1 if at least one of the predicted hallucinated objects was correct (for CHAIR), or if the object with the minimum matching score was a true hallucination (for ALOHa). We do not measure LA for CLIPScores, as they do not provide hallucination scores per object.

## C.2 FOIL

Table A2 breaks down the results of ALOHa on the FOIL and nocaps-FOIL dataset. The results illustrate a set of subtle results: while ALOHa under-performs CHAIRs in both AP and LA on the original FOIL dataset, this is primarily due to the construction of the dataset itself. FOIL constructs new samples by replacing string-matched COCO objects with a set of hand selected “foil” objects, objects that are visually distinct from the original object, but are near semantic neighbors. This is a best case scenario for CHAIR, as CHAIR relies

on string matching alone, and thus, is easily able to both detect and localized the replaced samples. The inaccuracies in LA and AP come from the synonym set that CHAIR uses for matching, along with parsing errors such as parsing the color “orange” as the object “orange”. In much the way that FOIL is the perfect dataset for CHAIR, FOIL perfectly exploits the strengths of ALOHa. Because of the semantic similarity score, we assign less weight to in-domain hallucinations, making it less likely that the replaced FOIL objects will be detected.

When we move to nocaps-FOIL with non-MS COCO data, however, we see starkly contrasting results. ALOHa significantly outperforms CHAIR, as now the object set that was a strength for in-domain FOIL becomes a liability, and CHAIR is unable to detect any hallucinations at all, due to the restricted string matching. RefCLIPScore, while extremely competitive in the hallucination detection task, cannot perform localization, and thus, serves only as a benchmark for the performance of ALOHa on the FOIL and nocaps-FOIL datasets.

## D Datasets

In this section, we discuss further the data that we use and go into detail on the dataset collection process for HAT (Appendix D.2) and the nocaps-FOIL dataset (Appendix D.1)

### D.1 nocaps-FOIL

The FOIL dataset (Shekhar et al., 2017) is a synthetic hallucination dataset based on samples from the MS-COCO (Xu et al., 2016) dataset. In this dataset, for each candidate-image pair, a “foil” caption is created which swaps one of the objects (in the MS-COCO detection set) in the caption with a different, and closely related neighbor (chosen by hand to closely match, but be visually distinct). While the FOIL dataset provides a useful benchmark for many hallucination detection methods,

it is overly biased towards methods optimized for the MS-COCO dataset. To help evaluate methods that are more general, we introduce a new dataset “nocaps-FOIL” based on the nocaps (Agrawal et al., 2019) dataset. The nocaps dataset consists of images from the OpenImages (Kuznetsova et al., 2020) dataset annotated with image captions in a similar style to MS-COCO. nocaps is split into three sets: an in-domain set, where objects in the images are in the MS-COCO object set, near-domain, where the objects in the image are related to those of MS-COCO, and out-of-domain, where objects in the image are not contained in MS-COCO.

To build the nocaps-FOIL dataset, for each image, we generate the baseline caption by removing a single caption from the reference set. We then generate the foil caption as follows. First, we find any words in the baseline caption that are contained in either the openimages class list (there are 600) or a near neighbor in wordnet. We then randomly select one of these classes to replace. Because there are 600 classes, we do not hand-pick the foil classes, and rather, select a near neighbor class based on sentence embeddings from (Reimers and Gurevych, 2019). We find that in practice, the nearest neighbor is often a synonym, thus, to avoid selecting synonyms, we take the 10th furthest sample, which is often a near neighbor, but is visually distinct. We replace this word in the caption, matching case, and then perform a filter for grammatical correctness using the Ginger<sup>2</sup> API. Any captions which are not grammatically correct are filtered. This leaves us with 2500 image/caption/foil pairs, which we use for evaluation in Table A2.

The OpenImages dataset annotations are under a CC BY 4.0 license, and the images are under a CC BY 2.0 license.

## D.2 HAT

HAT is based on MS-COCO and aims to be a gold-standard benchmark for the evaluation of hallucination in image captioning methods. While it is relatively small, it is densely annotated by in-domain experts for several types of hallucination including object hallucination, action hallucination, and numeric hallucination among others. HAT consists of 90 validation samples, and 400 test samples, each containing a machine candidate caption generated by one of BLIP (Li et al., 2022), OFA (Wang et al., 2022), IC3

<sup>2</sup><https://www.gingersoftware.com/>

(Chan et al., 2023) or Chat-Captioner (Zhu et al., 2023), and annotations which mark which word in the captions are hallucinated (See Figure A8 for exact instructions given to annotators). An image/caption pair is considered a hallucination if at least one of the words in the caption is hallucinated.

Screenshots of the interface for data collection are given in Figure A8. While initial versions of the dataset were collected using AMT workers, we found that the quality of annotations was not sufficiently high, and thus, trained experts explicitly in hallucination detection, and leveraged expert ratings for the samples in the test dataset.

MS-COCO is under a Creative Commons Attribution 4.0 License.

## E Qualitative Examples

We provide additional qualitative examples from the following scenarios:

### E.1 Flickr30k Examples

Figure A5 shows several examples on the Flickr-30k dataset Young et al. (2014) with captions generated by IC3 (Chan et al., 2023), a modern image captioning model that often generates longer, more complex captions including uncertain language such as “possibly.” We highlight objects with  $ALOH_{a_0} \leq 0.5$  as likely hallucinations. For samples going from left to right:

1. The caption hallucinates the word “mother”, as there is no visual evidence that the woman is specifically a mother. CHAIR does not capture this, as “mother” is mapped to a synonym for “person”, which it counts as a grounded (non-hallucinated) object. ALOHa matches “mother” to the reference “person”, assigning a borderline  $ALOH_{a_0}$  of 0.5.
2. The image does not contain a hallucination. CHAIR flags “table” as hallucinated, yet the caption expressed uncertainty with a conjunction: “chair or table.” ALOHa successfully parses this conjunction and selects “cloth” with  $ALOH_{a_0} = 1.0$  to the exact reference match.
3. CHAIR does not detect the hallucinated “bridge”, which is successfully assigned a low  $ALOH_{a_0} = 0.35$ .

- 987 4. The caption hallucinates the word “father”.  
 988 In most cases, the specific relationship of  
 989 “father” is unlikely to be grounded (similar  
 990 to “mother” in sample 1); yet, in this image,  
 991 it is even more clear as there are only children  
 992 present. CHAIR maps “father” as another  
 993 synonym for “person” and does not consider  
 994 it a hallucination, whereas “father” has a low  
 995  $\text{ALOH}_{a_0} = 0.34$ .

## 996 E.2 HAT Examples

997 We present 4 random samples from HAT each  
 998 for cases without hallucinations (Figure A6) and  
 999 with hallucinations (Figure A7). Because these  
 1000 examples contain more nuance that we discuss  
 1001 below, we do not indicate binary hallucination  
 1002 decisions as in Appendix E.1.

1003 Starting with Figure A6), samples with captions  
 1004 that were labeled as correct, from left to right:

- 1005 1. Both CHAIR and ALOHa successfully do  
 1006 not find any hallucinations.
- 1007 2. CHAIR does not flag any hallucinations.  
 1008 ALOHa assigns a low  $\text{ALOH}_{a_0} = 0.36$  for  
 1009 “sun“, an incorrect parse from the phrase  
 1010 “sunny day”. However, the other objects are  
 1011 successfully matched. Interestingly, ALOHa  
 1012 adds “snowboard” as an object, inferring that  
 1013 the physical item would need to be present  
 1014 given the verb “snowboarding”.
- 1015 3. CHAIR again does not flag any hallucina-  
 1016 tions.  $\text{ALOH}_{a_0}$  for “tall building” is the  
 1017 mid-range 0.59, matched with the reference  
 1018 “building”, indicating a somewhat uncertain  
 1019 attribute. This may be reasonable given the  
 1020 point of view in the image.
- 1021 4. CHAIR finds no hallucinations. “Cloudy sky”  
 1022 receives a somewhat low  $\text{ALOH}_{a_0} = 0.45$ .  
 1023 Although this phrase is accurate given the  
 1024 image, this is a failure case in which the  
 1025 references are incomplete.

1026 Next, we discuss Figure A7, showing samples  
 1027 that were labeled to contain a hallucination. Recall  
 1028 that labels capture *all* types of caption errors, includ-  
 1029 ing those other than object hallucinations, to serve  
 1030 as a valuable source for research around general  
 1031 caption correctness. As a result, there exists non-  
 1032 object hallucinations in HAT that are impossible for  
 1033 CHAIR or ALOHa to localize. From left to right:

- 1034 1. The attribute “tall” is labeled as a hallucination,  
 1035 as the building next to the bus is only one story.  
 1036 Similar to sample 3 in Figure A6,  $\text{ALOH}_{a_0}$  for  
 1037 “tall building” is somewhat uncertain at 0.59.  
 1038 Other objects are correctly grounded.
- 1039 2. The object “table” is a hallucinated, misclas-  
 1040 sified object; e.g., one reference opts for the  
 1041 more general “wooden surface.” However, the  
 1042 reference mentions a “table” that it is placed  
 1043 on, leading CHAIR to avoid considering it  
 1044 as a hallucination. For ALOHa, this example  
 1045 shows one of the 2.97% of cases (Table A1)  
 1046 where ALOHa hallucinates a reference object,  
 1047 “dining table”. The candidate “round wooden  
 1048 table” is matched to it, with an erroneously  
 1049 high  $\text{ALOH}_{a_0}$  of 0.74.
- 1050 3. This sample contains a complex error, in which  
 1051 the arrow is not, in fact, “pointing in different  
 1052 directions.” This non-object hallucination  
 1053 is impossible for the object-specific CHAIR  
 1054 and ALOHa to localize correctly. However,  
 1055 it demonstrates ALOHa’s capability to extract  
 1056 more complex attributes such as “red street  
 1057 sign” and “orange detour sign.”
- 1058 4. The cat’s location “on top of a small chair”  
 1059 is labeled as an error. CHAIR does not flag  
 1060 any hallucinations.  $\text{ALOH}_{a_0}$  for “small chair”  
 1061 is 0.59, yet both metrics cannot capture the  
 1062 specific relation.

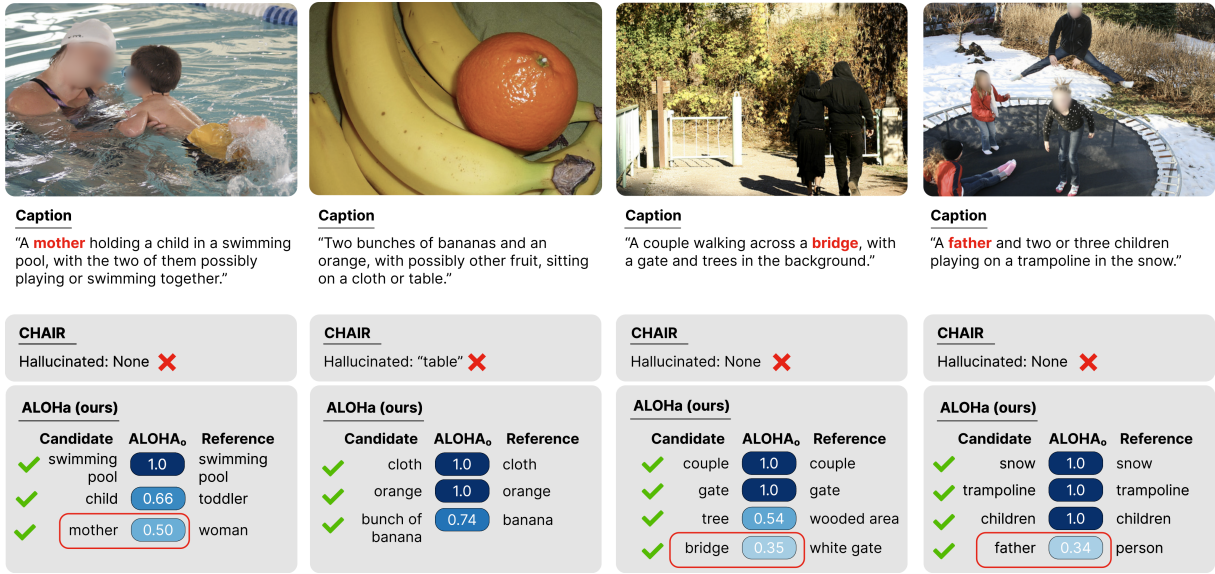


Figure A5: Qualitative samples of ALOHa evaluated on the Flickr-30k dataset, with candidate captions generated by IC3 (Chan et al., 2023). Hallucinated objects in the caption text are red and bolded. See Appendix E.1 for discussion.

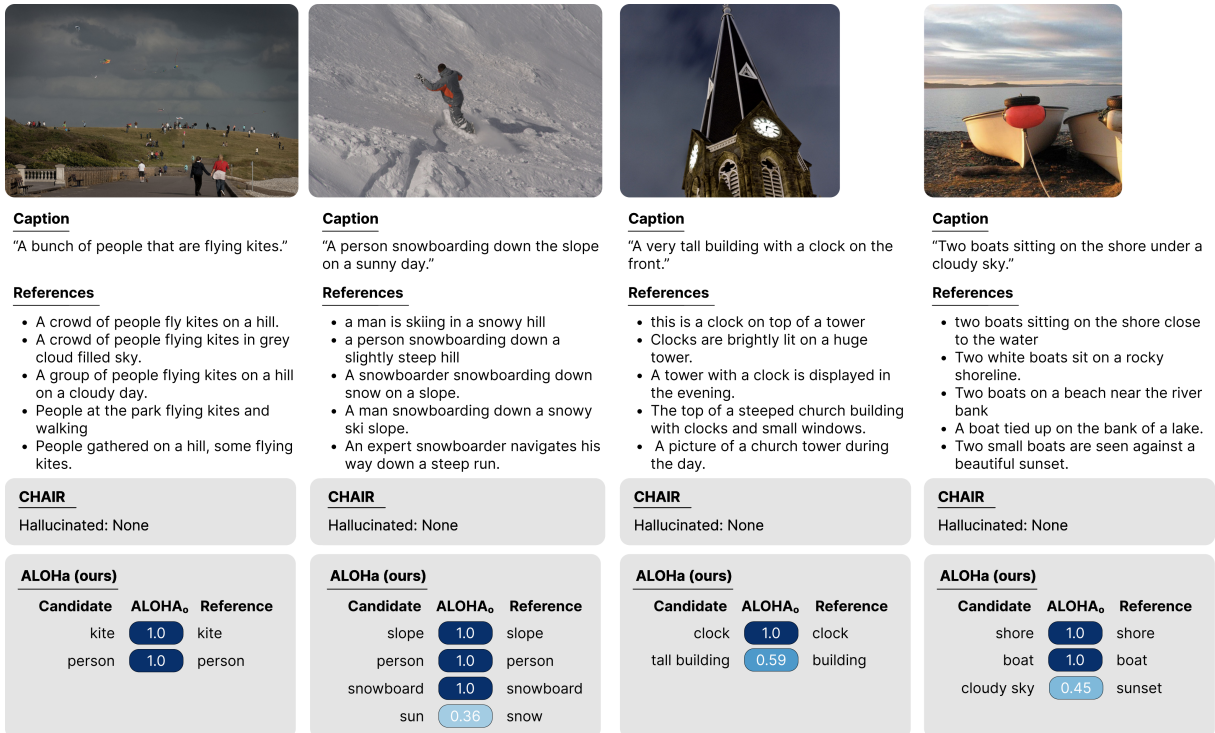


Figure A6: Randomly selected qualitative examples of ALOHa evaluated on the HAT dataset when there is no hallucination in the ground truth. See Appendix E.2 for discussion.



**Caption**

"A bus driving down a street next to a tall building."

**References**

- A large long bus on a city street.
- A city bus on the street in front of buildings.
- A blue bus traveling down an incline of a busy street.
- A city bus with full side advertisement in front of a building.
- a public transit bus on a city street

**CHAIR**

Hallucinated: None

**ALOHa (ours)**

Candidate	ALOHa <sub>o</sub>	Reference
bus	1.0	bus
street	1.0	street
tall building	0.59	building



**Caption**

"A round wooden table with a small pizza."

**References**

- A platter with a baked good on it
- A plain piece of bread resting on a wooden plate.
- A whole cheese pizza sitting on a wood pan on a table.
- a close up of a pizza on a wooden surface on a table
- A white cracker looking pizza is on a cutting board.

**CHAIR**

Hallucinated: None

**ALOHa (ours)**

Candidate	ALOHa <sub>o</sub>	Reference
round wooden table	0.74	dining table
small pizza	0.69	pizza



**Caption**

"A street sign with a detour pointing in different directions."

**References**

- An orange detour sign hanging from a metal pole under a cloudy sky.
- Red street sign with black letters sitting on metal post.
- A street pole with an orange detour sign.
- a close up of a street sign with a sky background
- A red detour sign that is on a pole.

**CHAIR**

Hallucinated: None

**ALOHa (ours)**

Candidate	ALOHa <sub>o</sub>	Reference
street sign	0.83	red street sign
detour	0.59	orange detour sign



**Caption**

"A cat stands on top of a small chair."

**References**

- A cat perched on top of a dresser. A cat walks along the top of a bedroom dresser.
- a cat sits on a dresser next to a rocking chair
- Black cat standing on a blue dresser next to a chair.
- A cat laying on top of a blue dresser near a chair.

**CHAIR**

Hallucinated: None

**ALOHa (ours)**

Candidate	ALOHa <sub>o</sub>	Reference
cat	1.0	cat
small chair	0.59	chair

Figure A7: Randomly selected qualitative examples of ALOHa evaluated on the HAT dataset when there is a hallucination in the ground truth. These hallucinations are generally challenging to detect. See Appendix E.2 for discussion.

## Description Rating Tool

**Instructions:** Review the image and text caption of that image, then click on any content words (nouns, adjectives, verbs, and numbers) in the caption which are not necessarily supported by the image content. Do not click on words like "The", "A", or "An".

For example, if the caption says "The cat is sleeping on the rug," yet there is nothing on the rug, click on the words "cat" and "sleeping". If the caption says "The vase contains three red roses," but there are only two roses in the image, click on the word "three".

If the caption uses an incorrect verb to describe an action in the image, click on that word. For example, if the caption reads "The woman is swimming in the ocean," but the image shows the woman walking on the beach, click on the word "swimming."

If a word is a compound word, such as "sofa chair," select either both words or neither word.

If it is impossible to tell whether a word is supported by the image or not, select that word anyways. For example, if the caption says "The child is smiling" and the image only shows the back of the child, it may be difficult to tell the child's facial expression. In this case, select the word "smiling" even if it's unclear whether or not it is accurate.

If no words are incorrect, select "Caption is correct". If either the caption or the image is not visible, press the "Not Visible" button.

HIT Tasks Completed: 100



Caption: A man holding a tennis racquet on a tennis court.

Select any incorrect words:

A man holding a tennis racquet on a tennis court.

Caption is correct

Image/Captions Not Visible

Submit

Figure A8: The hallucination dataset collection interface.