

VISaGE: Understanding Visual Generics and Exceptions

Anonymous ACL submission

Abstract

While Vision Language Models (VLMs) are trained to learn conceptual representations (generalized knowledge across many instances), they are typically *used* to analyze individual instances. When evaluation instances are atypical, this paradigm results in tension between two priors in the model. The first is a *pragmatic prior* that the textual and visual input are both relevant, arising from VLM finetuning on congruent inputs; the second is a *semantic prior* that the conceptual representation is generally true for instances of the category. In order to understand how VLMs trade-off these priors, we introduce a new evaluation dataset, VISaGE, consisting of both typical and *exceptional* images. In carefully balanced experiments, we show that VLMs are typically dominated by the semantic prior, which arises from the language modality, when answering queries about instances. In contrast, conceptual understanding degrades when the assumption of congruency underlying the pragmatic prior is violated with incongruent images.

1 Introduction

Vision-language models (VLMs) are typically used to analyze *instances*: what is going on in a particular image? However, during training they learn a set of *conceptual* representations, generalized knowledge that holds over many instances. While VLMs have been thoroughly tested on their ability to discern minimal differences between image instances (e.g., Johnson et al., 2017; Thrush et al., 2022; Tong et al., 2024), and their conceptual representations, based on exposure to typical instances, have long been analyzed (e.g., Bruni et al., 2014; Silberer et al., 2013; Collell and Moens, 2016), the potential tension between instance and concept representations, as arises in atypical instances, is currently under-explored.

In language, the attributes associated with a conceptual category are often expressed through *gener-*



Typical Cat

Conceptual query:
Do cats have 4 legs?

YES

Instance query:
Does this cat have 4 legs?

YES



Exceptional Cat

Conceptual query:
Do cats have 4 legs?

YES

Instance query:
Does this cat have 4 legs?

NO

Conceptual query with exception name:
Do tripod cats have 4 legs?

NO

Instance query with exception name:
Does this tripod cat have 4 legs?

NO

Figure 1: VISaGE contains both typical category instances and exceptional instances for which generics do not hold. We probe VLMs for conceptual and instance-level understanding, which is congruent in the typical case (top pair) but conflicts in the case of exceptional instances of a category cat (middle pair). However, the same exceptional instance can also be a typical member of the exception category (bottom pair).

ics – generalizations without quantifiers (e.g., cats have four legs). This lack of quantification means that generics remain true regardless of exceptions (tripod cats—cats missing one leg—do not impact the truth of “cats have four legs”). In other words, the attribute is associated as characteristic of the category regardless of how frequent it actually is¹.

Unlike language, which can denote on this generic or conceptual level, as well as refer to a particular instance, VLMs are always grounded in a particular visual instance. Work that has probed for conceptual attributes has used typical instances to stand in for the concept. This conflates instance and conceptual representations. In order to separate the two, visual *exceptions* are required: instances of a category that violate the generic (see Figure 1).

In this vein, we introduce a new evaluation dataset, **VisaGE**: Visual Generics and Exceptions, consisting of conceptual categories with images of

¹This is a substantial simplification of the semantics of generics (cf. Krifka, 1987).

062 both typical and exceptional instances. Specifically,
063 exceptions are always with regard to a particular
064 generic norm, i.e., a typical attribute: a tripod cat
065 is an exception for *cats have four legs*, but is typi-
066 cal for *cats have a long tail*. The category-attribute
067 pairs in VISaGE, along with their exceptions, are
068 extracted from textual generics and carefully manu-
069 ally validated, together with the image instances.

070 Using VisaGE, we investigate two questions:

- 071 1. **(RQ1)** How does conceptual information im-
072 pact VLMs’ ability to recognize *instance at-*
073 *tributes*?
- 074 2. **(RQ2)** How does visual grounding to (po-
075 tentially atypical) instances impact a model’s
076 ability to access *conceptual information*?

077 These research questions examine the effects of two
078 priors in VLMs. The first is a *pragmatic prior*, aris-
079 ing from VLM finetuning, that the textual and vi-
080 sual input are congruent and both relevant; the sec-
081 ond is a *semantic prior* that the category-attribute
082 generic is generally true². In the exceptional image
083 settings we explore with VISaGE, these two priors
084 can conflict: In RQ1, given an atypical instance,
085 the pragmatic prior to focus on the current con-
086 text must overrule the semantic prior of typicality,
087 while for RQ2, the atypical image must be ignored,
088 and the semantic prior should be followed.

089 We test a set of contemporary VLMs and find evi-
090 dence that their conceptual representations do not
091 recognize possible variation in attributes. Specifi-
092 cally, we find evidence that models rely on ex-
093 plicit textual cues to recognize instantiations of ex-
094 ceptional attributes in images, suggesting a strong
095 semantic prior for the generic. Additionally, we
096 observe that the models’ pragmatic prior often in-
097 terferes with conceptual understanding (and the
098 semantic prior) when visual grounding is incongru-
099 ent with the text. This suggests that VLMs’ visually
100 grounded conceptual representations only include
101 typical or generic instances of the category.

102 Our contributions are: 1. a new dataset, VISaGE,
103 consisting of concept-attribute pairs with images
104 of both typical (generic) and exceptional instances;
105 2. experimental evidence that VLM conceptual rep-
106 resentations are visually grounded only in typical
107 or generic instances and do not sufficiently recog-
108 nize within-category variation (exceptions).

²This is analogous to the Gricean maxims of relevance and quality (truthfulness) (Grice, 1975).

2 Background 109

110 Previous work has investigated the semantics of
111 generics with LMs (Ralethe and Buys, 2022; Col-
112 lacciani et al., 2024; Cilleruelo et al., 2025). These
113 studies show LMs often struggle to account for and
114 reason about exceptions in both probing (Allaway
115 et al., 2024) and reasoning (Allaway and McKeown,
116 2025) tasks. However, they have not considered
117 generics in VLMs, particularly how visual ground-
118 ing interacts with generic’s semantics.

119 For evaluating VLMs, most visual benchmarks
120 test situational and configurational instance under-
121 standing (Thrush et al., 2022; Li et al., 2024), some-
122 times with atypical examples (Bitton-Guetta et al.,
123 2023). Although Saleh et al. (2013) create a small
124 dataset of exceptional object images, these are not
125 annotated with semantic attributes, unlike VISaGE.
126 Additionally, our experiments, in which we manip-
127 ulate the image-text congruency, contribute to a
128 line of work investigating the relative importance
129 of different modalities in VLMs (Gat et al., 2021;
130 Frank et al., 2021; Parcalabescu and Frank, 2024).

3 Dataset 131

132 Our dataset VISaGE is constructed by first collect-
133 ing text pairs $(n_{c,a}, e_{c,a})$ where $n_{c,a}$ is a conceptual
134 norm for category c with attribute a and $e_{c,a}$ is an
135 exception to that norm (i.e., a subcategory of c that
136 does not have the attribute a). Then for each pair,
137 we retrieve two sets of images corresponding to
138 cases where the norm applies (generic images V_c)
139 and where it does not (exception images V_e). The
140 resulting dataset then consists of tuples $(n_{c,a}, e_{c,a},$
141 $V_c, V_e)$. Finally, we manually validate and expand
142 the dataset (details in Appendix A).

143 VISaGE contains 1698 exceptional image exam-
144 ples for 441 exception subcategories, derived from
145 972 category-attribute relations (conceptual norms)
146 for 171 categories, balanced with the same number
147 of typical images.³

148 **Norm-Exception Text Pairs** For our initial set of
149 concept-attribute norms, we intersect the category-
150 attribute lists of XCSLB (Devereux et al., 2014;
151 Misra et al., 2022) and the McRae norms (McRae
152 et al., 2005), with the categories in the THINGS
153 object image dataset (Hebart et al., 2019). This re-
154 sults in a robust set of conceptual norms expressed
155 as generics. Finally, for each generic (category-
156 attribute statement) we generate a set of exceptions

³Dataset & code will be released on publication (CC-BY).

$e_{c,a}$ using the LM prompting framework proposed by Allaway et al. (2024). We retain the short exceptions, ideally corresponding to subcategories.

Images We retrieve a large set of images for each exception subcategory using Bing Image Search by querying for the exception name $e_{c,a}$. Subsequent human validation (see below) selects the best images, resulting in a mode of 4 images per exception. A matched number of generic images for each category are taken from the THINGS dataset. These images have been specifically collected to be typical object instances; we further validate the applicability of the generic conceptual norms.

Validation We collect three types of validation annotations for each tuple. First, we validate that the images V_c retrieved from THINGS exhibit the conceptual norms $n_{c,a}$; category-attribute relations that are not visually salient (*birds can sing*) or are not exhibited across images are discarded. Second, we validate that each $e_{c,a}$ is actually an exception to the norm $n_{c,a}$. With this we filter out exception subcategories that are hallucinated (e.g., *strawberry blonde cheetah*) or incorrectly related (e.g., not exceptional or not actually subcategories). Finally, we validate that the retrieved images V_e for each exception are correct. We exclude images that are the wrong category (e.g., images of Ryan Gosling retrieved for the category *gosling*) or that are the wrong style (e.g., not object-centered photographs).

4 Experiments

Using VISaGE, our experiments query VLMs about conceptual and instance attributes across a number of conditions: see Fig. 2 for an overview. Specifically, we vary 1. the type of knowledge being queried (conceptual versus instance); 2. the type of image input (typical versus exceptional images); and, 3. the noun-phrase used to refer to the concept (category-name versus exception-name reference).

Models We test a suite of open-weights VLMs: these are listed in Appendix B. We use the vllm library to wrap our prompts⁴ in the correct model-specific formats.

Evaluation We report the percentage of correct (yes/no) responses for each model, using the first

⁴Conceptual prompt template example: Answer yes or no. Do {concept-pl} have {attribute}?
Instance prompt template example: Answer yes or no. Does this {concept-sg} have {attribute}?

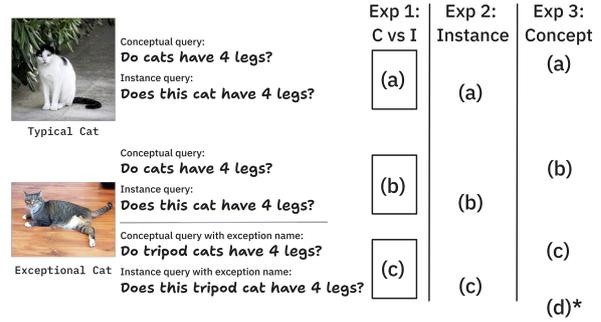


Figure 2: Summary of experiments and conditions: Exp. 1 measures the difference in accuracy between probing at the conceptual level vs. instance level. Exp. 2 tests models’ ability to reason about instances, while Exp 3 tests models’ conceptual understanding. Exp. 3 also includes condition (d), not shown, involving typical images with queries about exception categories.

token of the model output. Note that the correct response depends on the condition: see Figure 1.

4.1 Conceptual vs. Instance Queries

In an initial analysis, we test the ability of VLMs to distinguish between conceptual and instance queries. Specifically, we measure the difference in accuracy between the conceptual and instance queries in three conditions: (1a) typical images with category names, (1b) exceptional images with category names, and (1c) exceptional images with exception names. Condition (1b) is the critical condition, in which the correct prediction is different for conceptual and instance queries (see Fig. 1; numerical results are in App. Table F).

We observe (Figure 3a) that instance queries are in fact harder for VLMs than conceptual queries. When visual input and category name are congruent ((1a) and (1c)), we observe minimal differences (near zero) between the conceptual and instance queries. In contrast, when models are required to consider specific visual features of the input, rather than the semantic information from the category name, as for instance queries in (1b), we observe that most models fail to do this (conceptual accuracy is higher than instance accuracy). The accuracy difference that is visible *only* with incongruent inputs emphasizes the importance of considering how image instances interact with conceptual representations.

4.2 Instance Attribute Recognition

Having shown that VLMs struggle with instance queries requiring visual grounding (§4.1), our

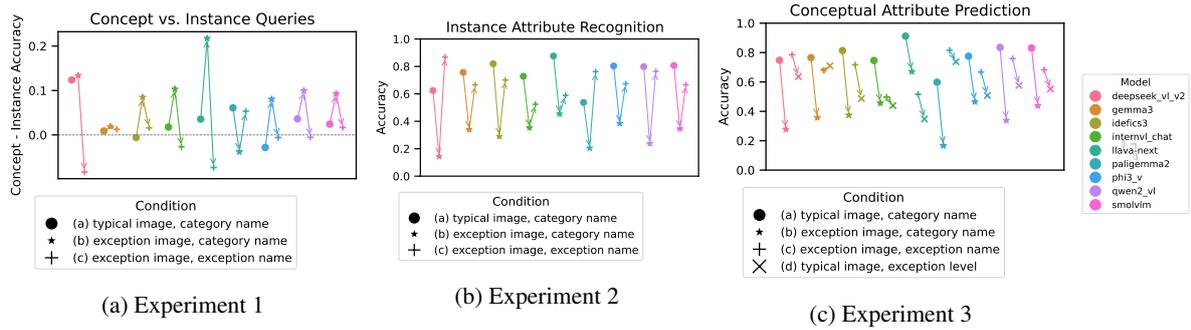


Figure 3: Results: See Fig. 2 for setup. Exp 1: The difference between conceptual and instance accuracy is highest for incongruent pairs (b). Exp 2: Instance attribute recognition declines for exceptional images (b), unless they are named as such (c). Exp 3: Conceptual attribute prediction accuracy decreases for incongruent inputs (b and d).

second set of analyses investigates the role of language-based conceptual activation in misleading models. We compare category-name instance queries (“Does this cat have four legs?”) in two conditions: with (2a) typical and (2b) exceptional images. The third condition (2c), *exception-name* instance queries with exceptional images, provides an explicit language cue to the model about which conceptual representation it should use (the exception rather than the category).

Our results (Figure 3b) show that, despite instance queries directing the model to consider the image, all models still appear to ignore the visual features, relying instead on language-based conceptual cues. In particular, we again see that models have higher accuracy in the conditions ((2a) and (2c)) where the text and image are congruent. When the text and image are *incongruent* (condition (2b)), the models appear to rely on language-based conceptual information. Since the image in (2b) is exceptional for the category, conceptual information activated by the category name does not apply to the image, resulting in a substantial drop in accuracy. This leads to the v-shaped pattern in accuracy. Note that if the models instead prioritized using the visual features of the input, their performance would be relatively stable across conditions.

4.3 Conceptual Attribute Prediction

Since we have shown that visual grounding is often ignored by VLMs when answering instance queries, our final set of experiments studies how it impacts conceptual queries. Specifically, we use two pairs of conditions to investigate the impact of text-image congruency. The first pair of conditions uses the category name in conceptual queries with: (3a) typical (congruent), and (3b) exceptional (in-

congruent) images. The second pair of conditions similarly queries conceptual information about the *exception* subcategory with: (3c) exceptional (congruent), and (3d) typical (incongruent) images.

Our results (Figure 3c) show that VLMs’ ability to answer conceptual questions degrades when the visual grounding is incongruent with the text input. That is, we observe a drop in accuracy in both pairs of conditions when comparing the congruent condition to the incongruent condition ((3a) vs. (3b) and (3c) vs. (3d)). This suggests that the pragmatic prior (considering the image relevant) interferes with the conceptual representation; that is, the image distracts the model from what is actually being asked in the query.

We also observe that incongruency in the input has less impact on accuracy when the queries are about the exception subcategory. One reason for this may be that the generic category is necessarily a well-established concept, since it is derived from a conceptual norm, while the exception subcategory may not be. Models may therefore lack a well-developed multimodal conceptual representation for the exception, resulting in them treating the conceptual query as an instance query. Our results that VLMs rely primarily on language cues for instance queries (§4.2), support this hypothesis.

5 Conclusion

VLMs must balance learned priors with the requirements of the current context. With the use of a new dataset of visual exceptions, VISaGE, we have shown that VLMs have not yet solved this task: Models neither reliably attend to the exception instance, ignoring the conceptual semantic prior, nor can they reliably ignore distractor images to answer generic conceptual queries.

306 Limitations

307 Our categories and attributes are limited to concep-
308 tual norms in American English. This is because
309 the typical images we use for visual grounding (de-
310 rived from THINGS) are based on American En-
311 glish definitions of categories. Conceptual spaces
312 are language-dependent and different languages
313 will make different conceptual distinctions, attend-
314 ing to different attributes. However, we believe the
315 general patterns of results would hold across lan-
316 guages and models, since the distinction between
317 instance-level and conceptual-level reasoning is
318 common across languages.

319 The data collection process focused on quality
320 rather than recall; we may have inadvertently omit-
321 ted particular important exception types. In particu-
322 lar, exceptions that are rare, hard to see, or unlikely
323 to be photographed, are missing (e.g., insomniac
324 owl as an exception for *owls sleep in the day*,
325 cheetah with a broken leg as an exception for
326 *cheetahs are fast*).

327 Compute limitations restricted the testing of very
328 large VLMs (llama4, pixtral).

329 **Risks** The concepts in our dataset correspond
330 to concrete object categories. However, the diffi-
331 culty of appropriately distinguishing (exceptional)
332 instances vs. conceptual generalizations can also
333 apply to categories that group people, where over-
334 generalization can lead to stereotyping. Under-
335 standing VLM capabilities and limitations is a step
336 towards mitigating these risks.

337 References

338 Emily Allaway, Chandra Bhagavatula, Jena D. Hwang,
339 Kathleen McKeown, and Sarah-Jane Leslie. 2024.
340 [Exceptions, Instantiations, and Overgeneralization:
341 Insights into How Language Models Process Gener-
342 ics](#). *Computational Linguistics*, pages 1291–1355.

343 Emily Allaway and Kathleen McKeown. 2025. [Evaluat-
344 ing defeasible reasoning in LLMs with DEFREAS-
345 ING](#). In *Proceedings of the 2025 Conference of the
346 Nations of the Americas Chapter of the Association
347 for Computational Linguistics: Human Language
348 Technologies (Volume 1: Long Papers)*, pages 10540–
349 10558, Albuquerque, New Mexico. Association for
350 Computational Linguistics.

351 Nitzan Bitton-Guetta, Yonatan Bitton, Jack Hessel, Lud-
352 wig Schmidt, Yuval Elovici, Gabriel Stanovsky, and
353 Roy Schwartz. 2023. [Breaking Common Sense:
354 WHOOPS! A Vision-and-Language Benchmark of
355 Synthetic and Compositional Images](#). In *2023*

*IEEE/CVF International Conference on Computer Vi-
sion (ICCV)*, pages 2616–2627, Paris, France. IEEE.

E. Bruni, N. K. Tran, and M. Baroni. 2014. [Multimodal
Distributional Semantics](#). *Journal of Artificial Intelli-
gence Research*, 49:1–47.

Gustavo Cilleruelo, Emily Allaway, Barry Haddow, and
Alexandra Birch. 2025. [Generics are puzzling. can
language models find the missing piece?](#) In *Proceed-
ings of the 31st International Conference on Compu-
tational Linguistics*, pages 6571–6588, Abu Dhabi,
UAE. Association for Computational Linguistics.

Claudia Collacciani, Giulia Rambelli, and Marianna
Bolognesi. 2024. [Quantifying generalizations: Ex-
ploring the divide between human and LLMs’ sensi-
tivity to quantification](#). In *Proceedings of the 62nd
Annual Meeting of the Association for Computational
Linguistics (Volume 1: Long Papers)*, pages 11811–
11822, Bangkok, Thailand. Association for Compu-
tational Linguistics.

Guillem Collell and Marie-Francine Moens. 2016. Is
an Image Worth More than a Thousand Words? On
the Fine-Grain Semantic Differences between Visual
and Linguistic Representations. In *Coling*.

Barry J. Devereux, Lorraine K. Tyler, Jeroen Geertzen,
and Billi Randall. 2014. [The Centre for Speech,
Language and the Brain \(CSLB\) concept property
norms](#). *Behavior Research Methods*, 46(4):1119.

Stella Frank, Emanuele Bugliarello, and Desmond El-
liott. 2021. [Vision-and-Language or Vision-for-
Language? On Cross-Modal Influence in Multimodal
Transformers](#). In *Proceedings of the 2021 Confer-
ence on Empirical Methods in Natural Language Pro-
cessing*, pages 9847–9857, Online and Punta Cana,
Dominican Republic. Association for Computational
Linguistics.

Itai Gat, Idan Schwartz, and Alex Schwing. 2021. Per-
ceptual Score: What Data Modalities Does Your
Model Perceive? In *Advances in Neural Information
Processing Systems*, volume 34, pages 21630–21643.
Curran Associates, Inc.

Herbert P Grice. 1975. Logic and conversation. In
Speech acts, pages 41–58. Brill.

Martin N. Hebart, Adam H. Dickter, Alexis Kidder,
Wan Y. Kwok, Anna Corriveau, Caitlin Van Wicklin,
and Chris I. Baker. 2019. [THINGS: A database of
1,854 object concepts and more than 26,000 natural-
istic object images](#). *PLOS ONE*, 14(10):e0223792.

Justin Johnson, Bharath Hariharan, Laurens Van
Der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and
Ross Girshick. 2017. [CLEVR: A Diagnostic Dataset
for Compositional Language and Elementary Visual
Reasoning](#). In *2017 IEEE Conference on Computer
Vision and Pattern Recognition (CVPR)*, pages 1988–
1997, Honolulu, HI. IEEE.

| | | | |
|-----|--|--|-----|
| 410 | Manfred Krifka. 1987. An outline of genericity. In | <i>Computer Vision and Pattern Recognition (CVPR)</i> , | 466 |
| 411 | <i>Seminar für natürlich-sprachliche Systeme der Uni-</i> | pages 9568–9578, Seattle, WA, USA. IEEE. | 467 |
| 412 | <i>versität Tübingen.</i> | | |
| 413 | Baiqi Li, Zhiqiu Lin, Wenxuan Peng, Jean de Dieu | | |
| 414 | Nyandwi, Daniel Jiang, Zixian Ma, Simran Khanuja, | | |
| 415 | Ranjay Krishna, Graham Neubig, and Deva Ramanan. | | |
| 416 | 2024. NaturalBench: Evaluating Vision-Language | A Dataset Construction | 468 |
| 417 | Models on Natural Adversarial Samples. <i>Advances</i> | | |
| 418 | <i>in Neural Information Processing Systems</i> , 37:17044– | The McRae norms are conceptual norms elicited | 469 |
| 419 | 17068. | from humans (McRae et al., 2005). Devereux et al. | 470 |
| 420 | Ken McRae, George S. Cree, Mark S. Seidenberg, and | (2014) builds on these in the XCSLB dataset and | 471 |
| 421 | Chris Mcnorgan. 2005. Semantic feature production | then (Misra et al., 2022) further revise them. Each | 472 |
| 422 | norms for a large set of living and nonliving things. | norm can be expressed as a generic. | 473 |
| 423 | <i>Behavior Research Methods</i> , 37(4):547–559. | | |
| 424 | Kanishka Misra, Julia Taylor Rayz, and Allyson Et- | To generate exceptions to the conceptual norms, | 474 |
| 425 | tinger. 2022. A property induction framework for | we use the framework proposed by Allaway et al. | 475 |
| 426 | neural language models. In <i>Proceedings of the 44th</i> | (2024). This framework proposes specific prompt | 476 |
| 427 | <i>Annual Conference of the Cognitive Science Society.</i> | templates for generating exceptions from LLMs, | 477 |
| 428 | Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, | along with a filtering process to ensure the gener- | 478 |
| 429 | Carroll Wainwright, Pamela Mishkin, Chong Zhang, | ated exceptions are true and salient. We use these | 479 |
| 430 | Sandhini Agarwal, Katarina Slama, Alex Ray, and 1 | templates with GPT-3.5 (Ouyang et al., 2022) ⁵ to | 480 |
| 431 | others. 2022. Training language models to follow in- | generate candidate exceptions and remove false | 481 |
| 432 | structions with human feedback. <i>Advances in neural</i> | ones. We keep the top 5 candidates ranked by per- | 482 |
| 433 | <i>information processing systems</i> , 35:27730–27744. | plexity to use in our dataset. | 483 |
| 434 | Letitia Parcalabescu and Anette Frank. 2024. Do Vi- | VISaGE includes substantial human validation, | 484 |
| 435 | sion & Language Decoders use Images and Text | including an iterative process of adding new at- | 485 |
| 436 | equally? How Self-consistent are their Explanations? | tribute norms and exceptions. During validation, | 486 |
| 437 | In <i>The Thirteenth International Conference on Learn-</i> | annotators can revise and expand the dataset by | 487 |
| 438 | <i>ing Representations.</i> | adding additional exceptions and category-attribute | 488 |
| 439 | Sello Ralethe and Jan Buys. 2022. Generic Overgeneral- | relations. Specifically, for valid category-attribute | 489 |
| 440 | ization in Pre-trained Language Models. In <i>Proceed-</i> | relations annotators, can provide an additional ex- | 490 |
| 441 | <i>ings of the 29th International Conference on Com-</i> | ceptional subcategory $\hat{e}_{c,a}$. Additionally, for each | 491 |
| 442 | <i>putational Linguistics</i> , pages 3187–3196, Gyeongju, | exception, annotators can provide a new category- | 492 |
| 443 | Republic of Korea. International Committee on Com- | attribute relation $n_{c,\hat{a}}$ that the exception corre- | 493 |
| 444 | putational Linguistics. | sponds to. This allows us to capture subcategories | 494 |
| 445 | Babak Saleh, Ali Farhadi, and Ahmed Elgammal. 2013. | that are exceptional for the category but not for the | 495 |
| 446 | Object-Centric Anomaly Detection by Attribute- | original attribute a . For example, pixie-bob cats are | 496 |
| 447 | Based Reasoning. In <i>2013 IEEE Conference on Com-</i> | an exception to <i>cats have long tails</i> but not to the | 497 |
| 448 | <i>puter Vision and Pattern Recognition</i> , pages 787–794, | original norm <i>cats have tails</i> . The tuples with the | 498 |
| 449 | Portland, OR, USA. IEEE. | new category-attribute norms $(n_{c,\hat{a}}, e_{c,\hat{a}}, V_c, V_e)$ ⁶ | 499 |
| 450 | Carina Silberer, Vittorio Ferrari, and Mirella Lapata. | are added directly into the dataset while for the new | 500 |
| 451 | 2013. Models of Semantic Representation with Vi- | exceptions $\hat{e}_{c,a}$, new images $V_{\hat{e}}$ are first retrieved | 501 |
| 452 | sual Attributes. In <i>Proceedings of the 51st Annual</i> | and validated before being added to the dataset as | 502 |
| 453 | <i>Meeting of the Association for Computational Lin-</i> | $(n_{c,a}, \hat{e}_{c,a}, V_c, V_{\hat{e}})$. | 503 |
| 454 | <i>guistics.</i> | | |
| 455 | Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet | The annotations were conducted by the authors | 504 |
| 456 | Singh, Adina Williams, Douwe Kiela, and Candace | of this paper. Through the revision and expansion | 505 |
| 457 | Ross. 2022. Winoground: Probing Vision and Lan- | process, we added 121 new tuples of conceptual- | 506 |
| 458 | guage Models for Visio-Linguistic Compositionality. | norm-and-exception (along with their correspond- | 507 |
| 459 | In <i>2022 IEEE/CVF Conference on Computer Vision</i> | ing images). Combined with the added conceptual | 508 |
| 460 | <i>and Pattern Recognition (CVPR)</i> , pages 5228–5238, | norms, we nearly doubled the size of our dataset | 509 |
| 461 | New Orleans, LA, USA. IEEE. | (an increase from 872 tuples to the final 1689 tu- | 510 |
| 462 | Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, | ples). | 511 |
| 463 | Yann LeCun, and Saining Xie. 2024. Eyes Wide | | |
| 464 | Shut? Exploring the Visual Shortcomings of Mul- | | |
| 465 | timodal LLMs. In <i>2024 IEEE/CVF Conference on</i> | | |

| Short Name | HF Model Name |
|----------------|--------------------------------------|
| deepseek_v1_v2 | deepseek/deepseek-v1.2-tiny |
| gemma3 | google/gemma-3-4b-it |
| idefics3 | HuggingFaceM4/Idfics3-8B-Llama3 |
| internvl_chat | OpenGVLab/InternVL2-2B |
| llava-next | llava-hf/llava-v1.6-mistral-7b-hf |
| paligemma2 | google/paligemma2-3b-ft-docci-448 |
| phi3_v | microsoft/Phi-3.5-vision-instruct |
| qwen2_v1 | Qwen/Qwen-VL |
| smolvlm | HuggingFaceTB/SmolVLM2-2.2B-Instruct |

Table 1: Models used in experiments.

B Models

See Table 1 for the details of the models used. Models are downloaded from HuggingFace; model details can be found at https://huggingface.co/MODEL_NAME.

C Compute

Experiments were performed using either Nvidia A100 or A4500 GPUs. On average, each evaluation (single model, condition) took approximately 15m, including model loading.

D Experimental Details

We used the `vllm`⁷ package, version `0.8.5.post1` with `transformers v4.52.0.dev0` and `torch v2.6.0`. Models were evaluated with default settings, apart from limiting the model’s output size in order to deal with memory limitations. We only evaluated the first output token.

E AI Agent Use

We used coding agents (copilot) to assist with code development. We did not use any AI agents for writing.

F Full Results

Numerical results for all experiments and conditions are in Table 4.

G Annotation Tool

See Figure 5.

⁵`gpt-3.5-turbo-0613`

⁶Note that $e_{c,\hat{a}} = e_{c,a}$; the changed index is for clarity.

⁷<https://docs.vllm.ai>

| Conds | prompt | image type | name | deepseek | idefics3 | qwen2_v1 | phi3_v | paligemma2 | gemma3 | llava-next | internvl_chat | smolvlm |
|----------|----------|------------|-----------|----------|----------|----------|--------|------------|--------|------------|---------------|---------|
| (1a, 3a) | concept | generic | category | 0.75 | 0.81 | 0.83 | 0.78 | 0.60 | 0.77 | 0.91 | 0.75 | 0.83 |
| (1a, 2a) | instance | generic | category | 0.62 | 0.82 | 0.80 | 0.80 | 0.54 | 0.76 | 0.88 | 0.73 | 0.81 |
| (1b, 3b) | concept | exception | category | 0.28 | 0.37 | 0.34 | 0.47 | 0.17 | 0.36 | 0.67 | 0.46 | 0.44 |
| (1b, 2b) | instance | exception | category | 0.14 | 0.29 | 0.24 | 0.38 | 0.20 | 0.34 | 0.45 | 0.35 | 0.35 |
| (1c, 3c) | concept | exception | exception | 0.79 | 0.72 | 0.76 | 0.67 | 0.82 | 0.68 | 0.52 | 0.50 | 0.68 |
| (1c, 2c) | instance | exception | exception | 0.87 | 0.70 | 0.76 | 0.67 | 0.76 | 0.67 | 0.59 | 0.52 | 0.67 |
| (3d) | concept | generic | exception | 0.64 | 0.49 | 0.58 | 0.51 | 0.74 | 0.71 | 0.35 | 0.44 | 0.55 |
| - | instance | generic | exception | 0.34 | 0.61 | 0.53 | 0.60 | 0.36 | 0.37 | 0.68 | 0.54 | 0.53 |

Figure 4: Accuracy results for all experiment conditions.

At concept 6 of 25

alligator: generic images (THINGS)



alligator: generic features that should apply to THINGS images. Check the correct attributes.

alligator: alligators have scales

How exception for the attribute: alligators have scales

EXCEPTIONS

alligator: soft toy alligators is an exception to the rule "alligators have scales"

How generic the exception (soft toy alligator) is for

Exception: soft toy alligator

| | | | |
|--|---|--|---|
| <input type="checkbox"/> keep c3738fa6c24ff985e6d878542555f-c1.jpg  | <input type="checkbox"/> keep 64579e8b3-7983-462b-acc5-720237056a2.jpg  | <input type="checkbox"/> keep Alligator Toy Figures for Kids-5x6-.jpg  | <input type="checkbox"/> keep 8d85e7223a43a11a1d3d0748a0d4e-af.jpg  |
| <input type="checkbox"/> keep 54d81c3ba3b6d07e6b385429f77-4a.jpg  | <input type="checkbox"/> keep 6c08892371c875a275e4d75af9a3-af.jpg  | <input type="checkbox"/> keep c22d32b-c2bd-4-f6-9b-b-953af4e4922.jpg  | <input type="checkbox"/> keep 42a2833ab3acc2b288de3a209f4c.jpg  |
| <input type="checkbox"/> keep 5a191e-Alligator Stuffed Animal-Cr.jpg  | <input type="checkbox"/> keep web6a6b-5cc4-45ab-4d87-c238d3923d.jpg  | <input type="checkbox"/> keep 7323ba4b7c0f5ee756462b49372b-af.jpg  | <input type="checkbox"/> keep 75493c58c3d81c45466ca229c6a0.wmp  |
| <input type="checkbox"/> keep 808180a-24-Alligator Stuffed Animal.jpg  | <input type="checkbox"/> keep 19a01c9f1c2819f9a254732ca43d8mshdng.jpg  | <input type="checkbox"/> keep 71800779e4c.jpg  | <input type="checkbox"/> keep 83a0c3c3c4c-82200.jpg  |

Update Concept Annotations

Figure 5: Annotation interface for dataset validation and expansion