

1 A Limitations and Societal Impacts

2 **Limitations** One limitation of our model is its potential for data bias. KOSMOS-1 is trained on a
3 web-scale multimodal corpus, which means that it is likely to be biased towards the data that it was
4 trained on. This could lead to the model generating text that is biased towards certain demographics
5 or viewpoints.

6 Another limitation of KOSMOS-1 is its relatively small size compared to other large language models.
7 This means that the model may not be able to learn as complex relationships between different
8 modalities. This could lead to the model making mistakes when it is asked to perform tasks that
9 require a deep understanding of multiple modalities.

10 Finally, KOSMOS-1 only supports vision modality. This means that the model cannot process other
11 modalities such as speech. This could limit the applications of the model.

12 **Societal Impacts** The broader impact of this paper is that it introduces a new type of large language
13 model that can perceive general modalities, follow instructions, and perform in-context learning. This
14 has the potential to be used for a variety of beneficial applications, such as new educational tools and
15 interactive dialogue assistants in video games. However, there are also potential negative impacts of
16 MLLMs. MLLMs could be used to create fake news articles or social media posts. MLLMs could be
17 used to generate text that reveals private information from web-scale pre-training data.

18 B Hyperparameters

19 B.1 Training

20 We report the detailed model hyperparameter settings of KOSMOS-1 in Table 1 and training hyperpa-
21 rameters in Table 2.

Hyperparameters	
Number of layers	24
Hidden size	2,048
FFN inner hidden size	8,192
Attention heads	32
Dropout	0.1
Attention dropout	0.1
Activation function	GeLU [1]
Vocabulary size	64,007
Soft tokens V size	64
Max length	2,048
Relative position embedding	xPos [2]
Initialization	Magneto [3]

Table 1: Hyperparameters of causal language model of KOSMOS-1

22 B.2 Language-Only Instruction Tuning

23 The detailed instruction tuning hyperparameters are listed in Table 3.

24 C Datasets

25 C.1 Pretraining

26 The models are trained on web-scale multimodal corpora. The training datasets consist of text corpora,
27 image-caption pairs, and interleaved data of images and texts.

Hyperparameters	
Training steps	300,000
Warmup steps	375
Batch size of text corpora	256
Max length of text corpora	2,048
Batch size of image-caption pairs	6,144
Batch size of interleaved data	128
Optimizer	Adam
Learning rate	2e-4
Learning Rate Decay	Linear
Adam ϵ	1e-6
Adam β	(0.9, 0.98)
Weight decay	0.01

Table 2: Training hyperparameters of KOSMOS-1

Hyperparameters	
Training steps	10,000
Warmup steps	375
Batch size of instruction data	256
Batch size of text corpora	32
Batch size of image-caption pairs	768
Batch size of interleaved data	16
Learning rate	2e-5

Table 3: Instruction tuning hyperparameters of KOSMOS-1

28 **Text Corpora** We train our model with The Pile [4] and Common Crawl (CC). The Pile is a
 29 massive English text dataset built for training large-scale language models, which is produced from a
 30 variety of data sources. We exclude data splits from GitHub, arXiv, Stack Exchange, and PubMed
 31 Central. We also include the Common Crawl snapshots (2020-50 and 2021-04) datasets, CC-Stories,
 32 and RealNews datasets [5, 6]. The entire datasets have been purged of duplicate and near-duplicate
 33 documents, as well as filtered to exclude downstream task data.

34 Table 4 provides a full overview of the language datasets that were used in the training of KOSMOS-1
 35 model. These data sources can be divided into the following three categories:

- 36 • **Academic:** NIH Exporter
- 37 • **Internet:** Pile-CC, OpenWebText2, Wikipedia (English), CC-2020-50, CC-2021-04, Realnews
- 38 • **Prose:** BookCorpus2, Books3, Gutenberg [7], CC-Stories

39 **Image-Caption Pairs** The image-caption pairs are constructed from several datasets, including
 40 English LAION-2B [8], LAION-400M [9], COYO-700M [10], and Conceptual Captions [11, 12].
 41 English LAION-2B, LAION-400M, and COYO-700M are collected from web pages of the Common
 42 Crawl web data by extracting image sources and the corresponding alt-text. Conceptual Captions are
 43 also from internet web pages.

44 LAION-2B contains about 2B English image-caption pairs, LAION-400M consists of 400M English
 45 image-caption pairs, and COYO-700M has 700M English image-caption pairs. Conceptual Captions
 46 contains 15M English image-caption pairs and consists of two datasets: CC3M and CC12M, which
 47 are also collected from internet webpages using a Flume pipeline. For Conceptual Captions, we
 48 discard pairs whose captions contain special tags such as “<PERSON>”.

49 **Interleaved Image-Text Data** We collect interleaved multimodal data from the Common Crawl
 50 snapshot, which is a publicly available archive of web pages. We use a filtering process to select
 51 about 71M web pages from the original 2B web pages in the snapshot. We then extract the text and

Datasets	Tokens (billion)	Weight (%)	Epochs
OpenWebText2	14.8	21.8%	1.47
CC-2021-04	82.6	17.7%	0.21
Books3	25.7	16.2%	0.63
CC-2020-50	68.7	14.7%	0.21
Pile-CC	49.8	10.6%	0.21
Realnews	21.9	10.2%	0.46
Wikipedia	4.2	5.4%	1.29
BookCorpus2	1.5	1.1%	0.75
Gutenberg (PG-19)	2.7	1.0%	0.38
CC-Stories	5.3	1.0%	0.19
NIH ExPorter	0.3	0.2%	0.75

Table 4: Language datasets used to train the KOSMOS-1 model.

52 images from the HTML of each selected web page. For each document, we limit the number of
53 images to five to reduce noise and redundancy. We also randomly discard half of the documents that
54 only have one image to increase the diversity. By using this corpus, we enable KOSMOS-1 to handle
55 interleaved text and image and improve its few-shot ability.

56 To ensure quality and relevance, we apply several filtering criteria. First, we discard any pages that
57 are not written in English. Second, we discard any pages that do not have images interspersed in
58 the text. Third, we discard any images that have a resolution lower than 64 by 64 pixels or that
59 are single-colored. Fourth, we discard any text that is not meaningful or coherent, such as spam or
60 gibberish. We use some heuristics to identify and remove gibberish text containing emoji symbols,
61 hashtags, and URL links. After applying these filters, we end up with about 71 million documents for
62 training.

63 C.2 Data Format

The training data is organized in the format as follows:

Datasets	Format Examples
Text	<s> KOSMOS-1 can perceive multimodal input, learn in context, and generate output. </s>
Image-Caption	<s> <image> Image Embedding </image> WALL-E giving potted plant to EVE. </s>
Multimodal	<s> <image> Image Embedding </image> This is WALL-E. <image> Image Embedding </image> This is EVE. </s>

Table 5: The examples of the data format to train the KOSMOS-1 model.

64

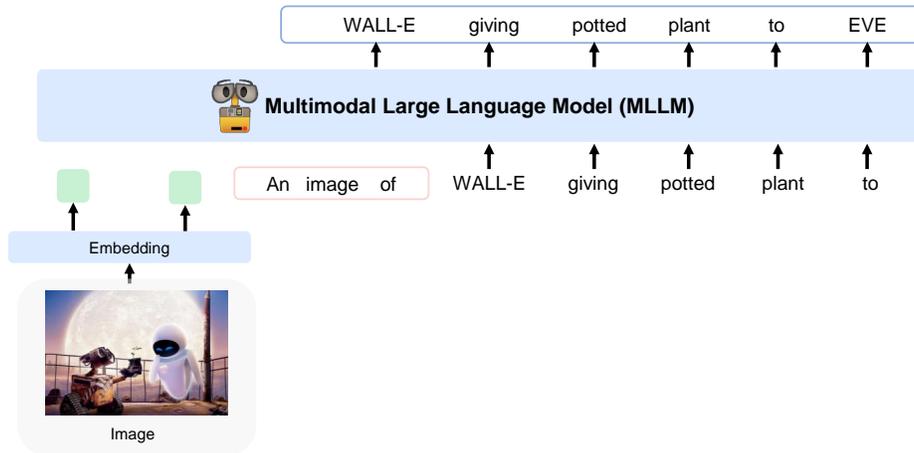
65 D Evaluation

66 D.1 Input Format Used for Perception-Language Tasks

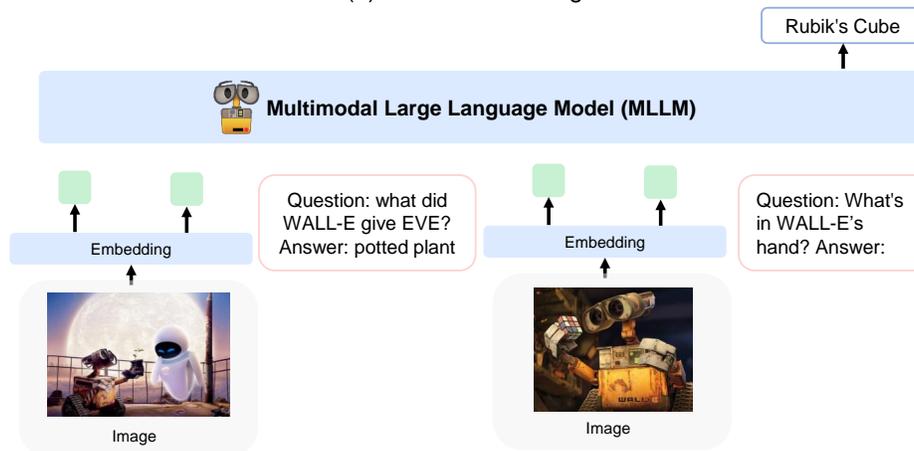
67 Figure 1 shows how we conduct zero-shot and few-shot evaluations on perception-language tasks.

68 D.2 Perception-Language Tasks

69 We evaluate the caption generation on MS COCO Caption [13], and Flickr30k [14]. We use the test
70 set of COCO *Karpathy split* [15], which re-partitions the train2014 and val2014 images [13] into
71 113,287, 5,000, and 5,000 for the training set, validation set, and test set, respectively. We conduct
72 an evaluation on Flickr30k’s *Karpathy split* test set. The image resolution is 224×224. We use
73 beam search to generate the captions, and the beam size is 5. In the few-shot settings, we randomly



(a) Zero-shot learning



(b) Few-shot learning

Figure 1: We evaluate KOSMOS-1 on the perception-language tasks in zero- and few-shot settings. (a) Zero-shot learning, e.g., zero-shot image captioning with language prompts. (b) Few-shot learning, e.g., visual question answering with in-context learning.

74 sample demonstrations from the training set. We use COCOEvalCap¹ to compute CIDEr [16] and
 75 SPICE [17] scores as the evaluation metrics. We prompt KOSMOS-1 with “*An image of*” for zero-shot
 76 and few-shot caption generation experiments.

77 For visual question-answering tasks, we evaluate zero-shot and few-shot results on test-dev set of
 78 VQAv2 [18] and test-dev set of VizWiz [19], respectively. The resolution of images is 224×224. We
 79 use greedy search for the decoding. We follow the normalization rules of the VQAv2 evaluation
 80 code² when computing the VQA accuracy. We evaluate the performance of VQA in an open-ended
 81 setting that KOSMOS-1 generates answers and stops at the </s> (“end of sequence”) token. The
 82 prompt is “*Question: {question} Answer: {answer}*” for visual question answering tasks.

¹<https://github.com/salaniz/pycocoevalcap>

²<https://github.com/GT-Vision-Lab/VQA>

83 D.3 IQ Test Tasks

84 To evaluate the KOSMOS-1 on zero-shot nonverbal reasoning, we construct a dataset of the Raven IQ
85 test. It consists of 50 examples collected from different websites³⁴⁵⁶. Each example has three (i.e.,
86 2×2 matrix), four, or eight (i.e., 3×3 matrix) given images. The goal is to predict the next one. Each
87 instance has six candidate images with a unique correct completion. We measure accuracy scores to
88 evaluate the models. The evaluation dataset is available at <https://aka.ms/kosmos-iq50>.

89 The matrix-style images are flattened and fed into the models one-by-one. To enable the model
90 to better understand the desired task, we also use a textual instruction “*Here are three/four/eight*
91 *images:*”, “*The following image is:*”, and “*Is it correct?*” for conditioning. We append each possible
92 candidate to the context separately and compare the probability that the model outputs “Yes” in a
93 close-ended setting. The candidate that yields the largest probability is regarded as the prediction.

94 D.4 OCR-Free Tasks

95 We evaluate OCR-free language understanding on the Rendered SST-2 [20] test set and Hateful-
96 Memes [21] validation set. We use accuracy as the metric for the Rendered SST-2 and report ROC
97 AUC for the HatefulMemes dataset. We use the prompt “*Question: what is the sentiment of the*
98 *opinion? Answer: {answer}*”, where the answer is either positive or negative for the Rendered SST-2.
99 For the HatefulMemes task, the prompt is “*Question: does this picture contain real hate speech?*
100 *Answer: {answer}*”, where the answer is either yes or no.

101 D.5 Web Page Tasks

102 We compare the performance on the Web-based Structural Reading Comprehension (WebSRC)
103 dataset [22]. For comparisons, we train a language model (LLM) on the same text corpora with
104 the same training setup as in KOSMOS-1. The LLM takes the text extracted from the web page as
105 input. Its template of the prompt is “*Given the context below from web page, extract the answer from*
106 *the given text like this: Q: Who is the publisher of this book? Answer: Penguin Books Ltd.*
107 *Context: {WebText} Q: {question} A: {answer}*”, where the *{WebText}* presents the text extracted
108 from the web page. Besides using the same prompt, KOSMOS-1 prepends the image before the
109 prompt. Two example images from WebSRC are shown in Appendix D.11. Following the original
110 paper [22], we use exact match (EM) and F1 scores as our evaluation metrics.

111 D.6 Multimodal CoT Tasks

112 We evaluate the ability of multimodal chain-of-thought prompting on the Rendered SST-2. We use the
113 prompt “*Introduce this picture in detail:*” to generate the content in the picture as the rationale. Then,
114 we use the prompt “*{rationale} Question: what is the sentiment of the opinion? Answer: {answer}*”
115 to predict the sentiment, where the answer is either positive or negative.

116 D.7 Zero-shot image classification Tasks

117 Given an input image, we concatenate the image with the prompt “*The photo of the*”. The input
118 is then fed into the model to obtain the category name of the image. We evaluate the model on
119 ImageNet [23], which contains 1.28M training images and 50k validation images in 1k object
120 categories. The prediction is classified as correct if it is exactly the same as the ground-truth category
121 name. The image resolution used for evaluation is 224×224 . We use beam search to generate the
122 category names and the beam size is 2.

123 D.8 Zero-Shot Image Classification with Descriptions

124 Following CUB [24], we construct a bird classification dataset that contains images and natural-
125 language descriptions of categories. The dataset has three groups of binary image classification. Each

³<https://en.testometrika.com/intellectual/iq-test/>

⁴<https://en.testometrika.com/intellectual/iq-test-for-kids-7-to-16-year-old/>

⁵<https://iqpro.org/>

⁶<https://iqhaven.com/matrix-g>

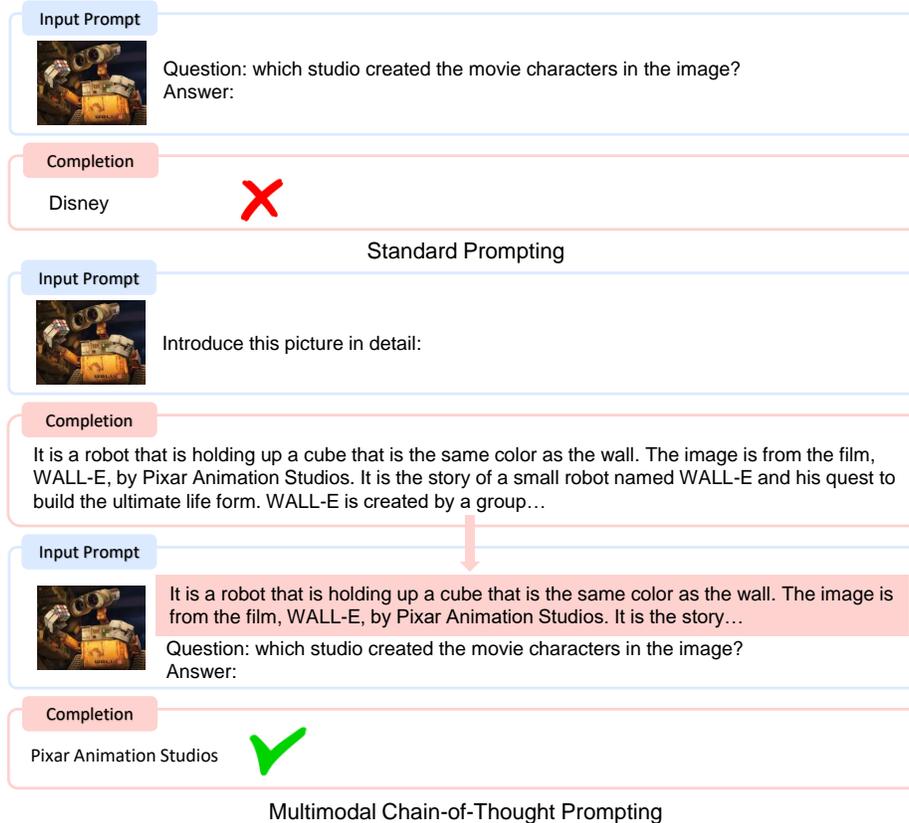


Figure 2: Multimodal Chain-of-Thought prompting enables KOSMOS-1 to generate a rationale first, then to tackle complex question-answering and reasoning tasks.

126 group contains two animal categories with similar appearances. Our goal is to classify images given
 127 the categories’ descriptions. Table 6 presents the data samples. The first group is from [24], while the
 128 other two groups are collected from the website. Each category contains twenty images.

129 The evaluation procedure is illustrated in Figure ???. For the zero-shot setting, we provide detailed
 130 descriptions of two specific categories and use the template “*Question: what is the name of {general*
 131 *category} in the picture? Answer:”* to prompt the model for the specific category name in an open-
 132 ended manner. To evaluate the effect of providing verbal descriptions in context, we also implement
 133 a zero-shot baseline without prompting descriptions. Instead, we provide the corresponding specific
 134 names in the prompt.

135 D.9 Cross-modal Transfer task

136 We compare KOSMOS-1 and the LLM baseline on three object commonsense reasoning datasets,
 137 RELATIVESIZE [25], MEMORYCOLOR [26] and COLORTERMS [27] datasets. Table 7 shows some
 138 examples of object size and color reasoning tasks. RELATIVESIZE contains 486 object pairs from 41
 139 physical objects. The model is required to predict the size relation between two objects in a binary
 140 question-answering format with “Yes”/“No” answers. MEMORYCOLOR and COLORTERMS require
 141 the model to predict the color of objects from a set of 11 color labels in a multiple-choice format. We
 142 use only text as our input and do not include any images. We measure the accuracy of our model on
 143 these three datasets.

144 D.10 Language Tasks

145 We train a language model (LLM) baseline with the same text corpora and training setup. We evaluate
 146 KOSMOS-1 and the LLM baseline on eight language tasks, including cloze and completion tasks (i.e.,

Category 1	Category 2
three toed woodpecker	downy woodpecker
 It has black and white stripes throughout the body and a yellow crown.	 It has white spots on its black wings and some red on its crown.
Gentoo penguin	royal penguin
 It has a black head and white patch above its eyes.	 It has a white face and a yellow crown.
black throated sparrow	fox sparrow
 It has white underparts and a distinctive black bib on the throat.	 It has a reddish-brown plumage and a streaked breast.

Table 6: The detailed descriptions of different categories for in-context image classification.

Task	Example Prompt	Object / Pair	Answer
Object Size Reasoning	<i>Is {Item1} larger than {Item2}? {Answer}</i>	(sofa, cat)	Yes
Object Color Reasoning	<i>The color of {Object} is? {Answer}</i>	the sky	blue

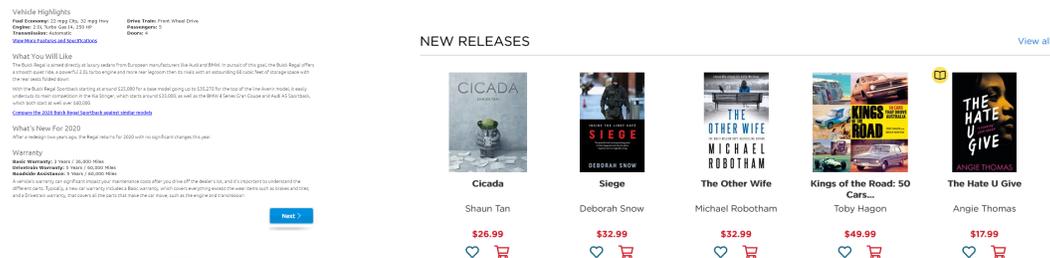
Table 7: Evaluation examples of object size and color reasoning.

147 StoryCloze, HellaSwag), Winograd-style tasks (i.e., Winograd, Winogrande), commonsense reasoning
 148 (i.e., PIQA), and three datasets BoolQ, CB, and COPA from the SuperGLUE benchmark [28]. The
 149 detailed descriptions of these datasets are provided in Appendix D.10. We conduct experiments under
 150 zero-shot and few-shot settings. We evaluate each test example by randomly sampling examples from
 151 the training set as demonstrations. We set the number of shots to 0, 1, and 4 in our experiments.

152 We conduct experiments on language tasks in four categories:

- 153 • Cloze and completion tasks: StoryCloze [29], HellaSwag [30]
- 154 • Winograd-style tasks: Winograd [31], Winogrande [32]
- 155 • Commonsense reasoning: PIQA [33]
- 156 • Three datasets from SuperGLUE benchmark [28]: BoolQ [34], CB [35], COPA [36]

157 D.11 WebSRC Task Examples



The screenshot shows a 'NEW RELEASES' section with five book covers: 'Cicada' by Shaun Tan (\$26.99), 'Siege' by Deborah Snow (\$32.99), 'The Other Wife' by Michael Robotham (\$32.99), 'Kings of the Road: 50 Cars...' by Toby Hagon (\$49.99), and 'The Hate U Give' by Angie Thomas (\$17.99). Each book has a heart icon and a shopping cart icon below its price.

(a) Question is “What is the type of this drive?”

(b) Question is “Who is the author of “Cicada”?”

Figure 3: Examples from WebSRC [22].

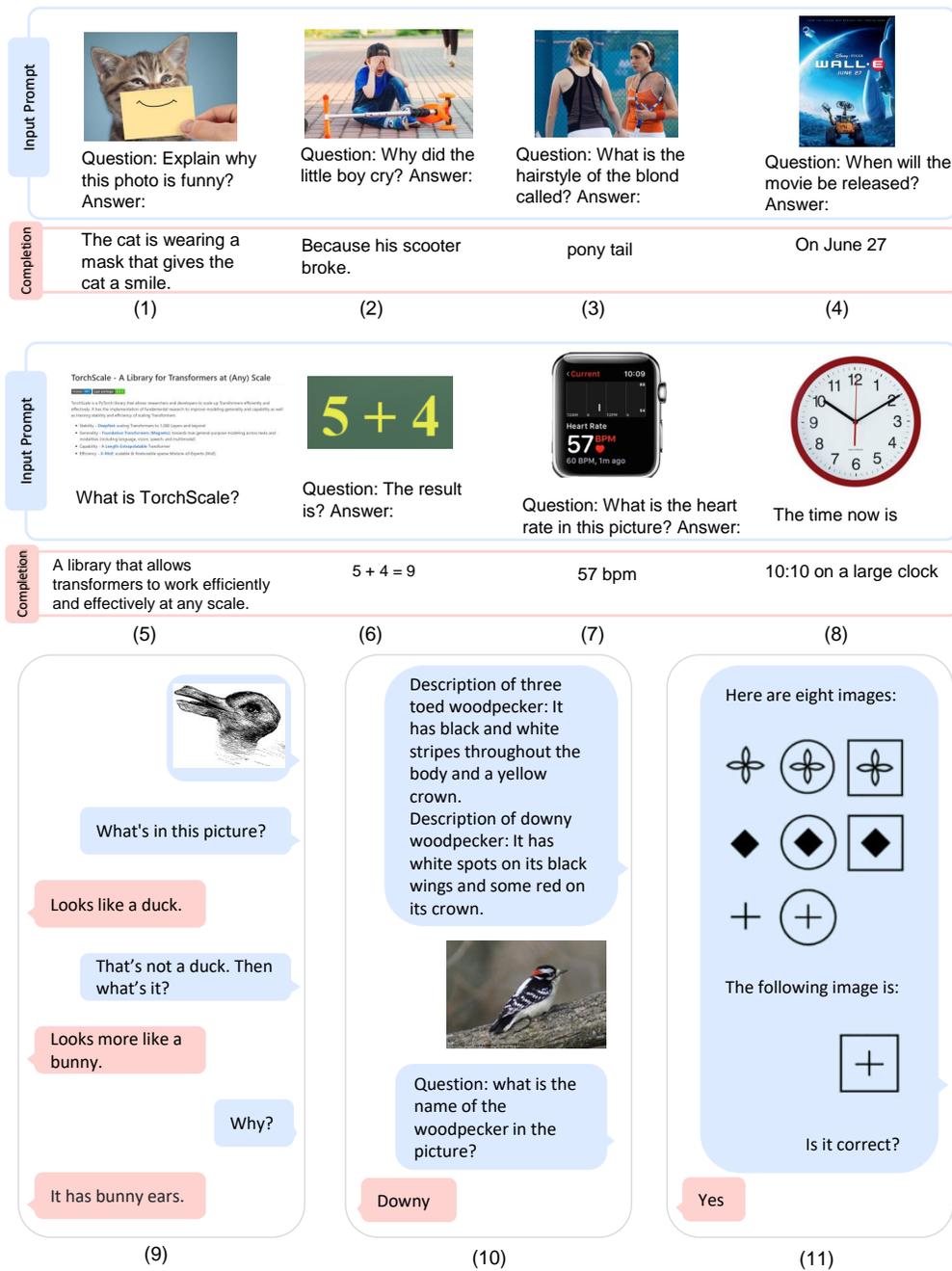


Figure 4: Selected examples generated from KOSMOS-1. Blue boxes are input prompt and pink boxes are KOSMOS-1 output. The examples include (1)-(2) visual explanation, (3)-(4) visual question answering, (5) web page question answering, (6) simple math equation, and (7)-(8) number recognition, and (9)-(11) visual dialogue.



Figure 5: Selected examples generated from KOSMOS-1. Blue boxes are input prompt and pink boxes are KOSMOS-1 output. The examples include (1)-(2) image captioning, (3)-(6) visual question answering, (7)-(8) OCR, and (9)-(11) visual dialogue.

References

- 159
- 160 [1] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (GELUs). *arXiv preprint*
161 *arXiv:1606.08415*, 2016.
- 162 [2] Yutao Sun, Li Dong, Barun Patra, Shuming Ma, Shaohan Huang, Alon Benhaim, Vishrav
163 Chaudhary, Xia Song, and Furu Wei. A length-extrapolatable transformer. *arXiv preprint*
164 *arXiv:2212.10554*, 2022.
- 165 [3] Hongyu Wang, Shuming Ma, Shaohan Huang, Li Dong, Wenhui Wang, Zhiliang Peng, Yu Wu,
166 Payal Bajaj, Saksham Singhal, Alon Benhaim, Barun Patra, Zhun Liu, Vishrav Chaudhary, Xia
167 Song, and Furu Wei. Foundation transformers. *CoRR*, abs/2210.06423, 2022.
- 168 [4] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason
169 Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse
170 text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 171 [5] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan
172 Catanzaro. Megatron-lm: Training multi-billion parameter language models using model
173 parallelism. *arXiv preprint arXiv:1909.08053*, 2019.
- 174 [6] Shaden Smith, Mostofa Patwary, Brandon Norick, Patrick LeGresley, Samyam Rajbhandari,
175 Jared Casper, Zhun Liu, Shrimai Prabhunoye, George Zerveas, Vijay Korthikanti, Elton Zhang,
176 Rewon Child, Reza Yazdani Aminabadi, Julie Bernauer, Xia Song, Mohammad Shoeybi,
177 Yuxiong He, Michael Houston, Saurabh Tiwary, and Bryan Catanzaro. Using DeepSpeed and
178 Megatron to train Megatron-Turing NLG 530B, a large-scale generative language model, 2022.
- 179 [7] Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, Chloe Hillier, and Timothy P. Lillicrap.
180 Compressive transformers for long-range sequence modelling. In *ICLR*, 2020.
- 181 [8] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman,
182 Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-
183 5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint*
184 *arXiv:2210.08402*, 2022.
- 185 [9] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton
186 Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open
187 dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- 188 [10] Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon
189 Kim. Coyo-700m: Image-text pair dataset, 2022.
- 190 [11] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A
191 cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings*
192 *of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018,*
193 *Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, pages 2556–2565. Association
194 for Computational Linguistics, 2018.
- 195 [12] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing
196 web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the*
197 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3558–3568, 2021.
- 198 [13] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan,
199 Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*,
200 pages 740–755, 2014.
- 201 [14] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to
202 visual denotations: New similarity metrics for semantic inference over event descriptions. *TACL*,
203 2:67–78, 2014.
- 204 [15] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image
205 descriptions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(4):664–676,
206 2017.

- 207 [16] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image
208 description evaluation. In *CVPR*, pages 4566–4575, 2015.
- 209 [17] Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Spice: Semantic
210 propositional image caption evaluation. In *ECCV*, pages 382–398, 2016.
- 211 [18] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the
212 v in vqa matter: Elevating the role of image understanding in visual question answering. In
213 *CVPR*, pages 6325–6334, 2017.
- 214 [19] Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo,
215 and Jeffrey P Bigham. Vizwiz grand challenge: Answering visual questions from blind people.
216 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages
217 3608–3617, 2018.
- 218 [20] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
219 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
220 models from natural language supervision. In *International conference on machine learning*,
221 pages 8748–8763. PMLR, 2021.
- 222 [21] Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik
223 Ringshia, and Davide Testuggine. The hateful memes challenge: Detecting hate speech in
224 multimodal memes. In *Advances in Neural Information Processing Systems*, volume 33, pages
225 2611–2624, 2020.
- 226 [22] Xingyu Chen, Zihan Zhao, Lu Chen, JiaBao Ji, Danyang Zhang, Ao Luo, Yuxuan Xiong, and
227 Kai Yu. WebSRC: A dataset for web-based structural reading comprehension. In *Proceedings of
228 the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4173–4185,
229 Online and Punta Cana, Dominican Republic, November 2021. Association for Computational
230 Linguistics. doi: 10.18653/v1/2021.emnlp-main.343. URL [https://aclanthology.org/
231 2021.emnlp-main.343](https://aclanthology.org/2021.emnlp-main.343).
- 232 [23] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
233 hierarchical image database. In *2009 IEEE Computer Society Conference on Computer Vision
234 and Pattern Recognition (CVPR 2009), 20-25 June 2009, Miami, Florida, USA*, pages 248–255.
235 IEEE Computer Society, 2009.
- 236 [24] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge J. Belongie. The
237 caltech-ucsd birds-200-2011 dataset. 2011.
- 238 [25] Hessam Bagherinezhad, Hannaneh Hajishirzi, Yejin Choi, and Ali Farhadi. Are elephants bigger
239 than butterflies? reasoning about sizes of objects. *ArXiv*, abs/1602.00753, 2016.
- 240 [26] Tobias Norlund, Lovisa Hagström, and Richard Johansson. Transferring knowledge from vision
241 to language: How to achieve it and how to measure it? *ArXiv*, abs/2109.11321, 2021.
- 242 [27] Elia Bruni, Gemma Boleda, Marco Baroni, and Nam Khanh Tran. Distributional semantics in
243 technicolor. In *ACL*, 2012.
- 244 [28] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill,
245 Omer Levy, and Samuel R Bowman. SuperGLUE: A stickier benchmark for general-purpose
246 language understanding systems. *arXiv preprint arXiv:1905.00537*, 2019.
- 247 [29] Nasrin Mostafazadeh, Michael Roth, Annie Louis, Nathanael Chambers, and James Allen.
248 Lsdsem 2017 shared task: The story cloze test. In *Proceedings of the 2nd Workshop on Linking
249 Models of Lexical, Sentential and Discourse-level Semantics*, pages 46–51, 2017.
- 250 [30] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can
251 a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the
252 Association for Computational Linguistics*, 2019.
- 253 [31] Hector J. Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. In
254 *Principles of Knowledge Representation and Reasoning*, 2012.

- 255 [32] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. WinoGrande: An
256 adversarial winograd schema challenge at scale. In *AAAI*, pages 8732–8740, 2020.
- 257 [33] Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning
258 about physical commonsense in natural language. In *Thirty-Fourth AAAI Conference on*
259 *Artificial Intelligence*, 2020.
- 260 [34] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and
261 Kristina Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions.
262 In *Proceedings of the 2019 Conference of the North American Chapter of the Association*
263 *for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short*
264 *Papers)*, pages 2924–2936, Minneapolis, Minnesota, June 2019. Association for Computational
265 Linguistics. doi: 10.18653/v1/N19-1300. URL <https://aclanthology.org/N19-1300>.
- 266 [35] Marie-Catherine de Marneffe, Mandy Simons, and Judith Tonhauser. The CommitmentBank:
267 Investigating projection in naturally occurring discourse. *Proceedings of Sinn und Bedeutung*,
268 23(2):107–124, Jul. 2019.
- 269 [36] Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S. Gordon. Choice of plausible
270 alternatives: An evaluation of commonsense causal reasoning. In *AAAI Spring Symposium*,
271 2011.