# Enhancing Journalism with AI: A Study of Contextualized Image Captioning for News Articles using LLMs and LMMs

Aliki Anagnostopoulou<sup>1,2</sup>, Thiago S. Gouvêa<sup>1</sup> and Daniel Sonntag<sup>1,2</sup>

 <sup>1</sup> DFKI German Research Center for Artificial Intelligence
 <sup>2</sup> Applied Artificial Intelligence, Carl von Ossietzky University Oldenburg {firstname.lastname}@dfki.de,

#### Abstract

Large language models (LLMs) and large multi-1 modal models (LMMs) have significantly impacted 2 the AI community, industry, and various economic 3 sectors. In journalism, integrating AI poses unique 4 challenges and opportunities, particularly in en-5 6 hancing the quality and efficiency of news report-7 ing. This study explores how LLMs and LMMs can assist journalistic practice by generating contextu-8 alised captions for images accompanying news ar-9 ticles. We conducted experiments using the Good-10 News dataset to evaluate the ability of LMMs 11 (BLIP-2, GPT-4v, or LLaVA) to incorporate one of 12 two types of context: entire news articles, or ex-13 tracted named entities. In addition, we compared 14 their performance to a two-stage pipeline com-15 posed of a captioning model (BLIP-2, OFA, or ViT-16 GPT2) with post-hoc contextualisation with LLMs 17 (GPT-4 or LLaMA). We assess a diversity of mod-18 19 els, and we find that while the choice of contextu-20 alisation model is a significant factor for the twostage pipelines, this is not the case in the LMMs, 21 where smaller, open-source models perform well 22 compared to proprietary, GPT-powered ones. Ad-23 ditionally, we found that controlling the amount of 24 provided context enhances performance. These re-25 sults highlight the limitations of a fully automated 26 approach and underscore the necessity for an inter-27 active, human-in-the-loop strategy. 28

#### **29 1** Introduction

Large language pre-training [Devlin et al., 2019; Liu et al., 30 2023c] and large vision-language pre-training [Wang et al., 31 2022; Zou et al., 2023], facilitated by advances in deep learn-32 ing and the development of the Transformer [Vaswani et 33 al., 2017] architecture, have significantly impacted research 34 and industry. In journalism, these models offer potential for 35 human-AI collaboration; however, generating news articles 36 with these models is not feasible since these pre-trained mod-37 els lack up-to-date information on current events [Bubeck et 38 al., 2023], among other reasons. Instead, they can assist jour-39 nalists by automating specific tasks, such as captioning im-40 ages that accompany existing news articles. These captions 41



Figure 1: Our proposed architectures. Our two-stage pipeline CIC involves an image captioning system; the generated caption and contextual information are fed into an LLM. We compare this architecture with LLMs, which consider visual and textual input, omitting the need for an additional image captioning component.

should describe the image content *and* provide relevant contextual information that cannot be deduced from the image, including the names of people, locations, and events.

In this study, we investigate the effectiveness of large foun-45 dation models, specifically large language models (LLMs) 46 and large multimodal models for vision-language tasks 47 (LMMs) in performing *contextualised image captioning*. We 48 conduct experiments using the GoodNews dataset, which 49 contains contextualised image captions from news articles. 50 We propose a two-stage pipeline composed of an image cap-51 tioning model with post-hoc contextualisation performed by 52 an LLM. We compare this pipeline, which we denote as CIC, 53 with LLMs (see Figure 1). We evaluate nine configurations, 54 including open-source models such as Llama 3 and LLaVA 55 [Liu et al., 2023b; Liu et al., 2023a; Liu et al., 2024] and 56 closed-source models such as GPT-3 [Liu et al., 2023c] and 57 GPT-4v [OpenAI, 2023]. 58

After briefly presenting related work to AI in journalism and contextualised image captioning (Section 2), we describe our pipelines, the foundation models used, the dataset, and the evaluation metrics in Section 3. We then present and discuss our results in Section 4. Section 5 concludes our work and discusses possible future directions. 64

Image	Caption	Article	NE
	Residents and activists aided a girl who survived amid debris in <b>Aleppo</b> on <b>Sunday</b> after what activists said was an aerial attack that dropped explosive barrels.	ISTANBUL – Syrian government aircraft continued to strike rebel-held areas in Aleppo with makeshift bombs on Sunday, killing at least three dozen people, most of them women and children, antigovernment activists said. [] The government has not commented on the airstrikes other than to mention in the state news media that its forces have killed "terrorists," a blanket term for the opposition.	GPE: Aleppo; DATE: Sunday

Table 1: Example image, caption and relevant context (article and extracted NEs) from the GoodNews dataset. The extracted NEs are also marked in the caption.

## 65 2 Related Work

This section reviews previous research relevant to our study,
 focusing on the application of AI in journalism and the field
 of contextualised captioning.

LLMs and LMMs in journalism [del Barrio and Gática-69 Pérez, 2023] use GPT-3.5 for news frame classification using 70 fine-tuning and prompt engineering. [Bao et al., 2024] de-71 velop a model specialised for question answering and data 72 visualisation in the business and media domain. Following 73 a human-centric approach, [Cheng et al., 2024] propose a 74 model incorporating human input for generating sports news 75 insights. 76

As mentioned in Section 1, LLMs are unsuitable for gener-77 ating news articles due to multiple shortcomings. Besides not 78 being up-to-date, making them prone to hallucinations, LLM 79 news generations can be biased: [Fang et al., 2023] evaluate 80 the gender and racial biases reproduced in LLM-generated 81 content. [Hamilton and Piper, 2022] use GPT-2 to generate 82 counterfactual news articles, finding out that they exhibit a 83 notably more negative attitude towards COVID and a signifi-84 cantly reduced reliance on geopolitical framing. 85

Contextualised image captioning Contextualised image 86 captioning considers additional context to generate an im-87 age caption that describes the image's content and includes 88 relevant external information. The context provided is, in 89 most cases, in textual form. [Biten et al., 2019] and [Tran 90 et al., 2020] use news articles as context; the former uses 91 a template-based architecture, and the latter uses an end-92 to-end architecture, considering additional features such as 93 face and object detection. A modified version of the lat-94 ter is used in [Nguyen et al., 2023] for image captioning on 95 Wikipedia [Srinivasan et al., 2021], while an additional face 96 naming module is also present in [Qu et al., 2023]. [Rajaku-97 mar Kalarani et al., 2023] present a unified architecture for 98 context-assisted image captioning, including contextual vi-99 sual entailment and keyword extraction. 100

## 101 **3 Methods**

This section describes our contextualising captioning experiments, namely the pipelines, type of context used, and evaluation (including datasets and metrics).

#### **TEXT PROMPT - GOODNEWS:**

You are a journalist. Describe this caption in a single sentence so the description suits a news article: [image caption]. Take the following context information into consideration: [context info]

You are a journalist. Describe this image in one sentence so the description suits a news article. Take the following context information into consideration: [context info]

Table 2: Proposed prompt for generating contextualised image captions using the GoodNews news dataset. In the upper row, the prompt for the CIC pipeline is present, and in the lower row the one for LMM.

#### 3.1 Pipelines

We use two approaches to contextualise captions, which we describe below. In the first one, we pair a conventional image captioning (CIC) architecture with an LLM for post-hoc contextualisation. In the second one, we utilise large multimodal models (LMMs), in which the image is directly provided as input, along with the context. Both pipelines are presented in Figure 1.

105

**CIC** As mentioned above, this approach follows two stages. 113 In the base captioning stage, the image captioning architec-114 ture generates a description for a given image. In the con-115 textualising stage, the LLM takes the generated caption and 116 the context as input and generates a caption that includes the 117 context for each image. Such an approach might be benefi-118 cial if access to LLMs is not possible, if privacy issues are 119 present, or if a base caption is needed a priori. A drawback 120 in this case, however, is the *information bottleneck* caused by 121 the contextualising model not having access to all the visual 122 information in the image. We use pre-trained SOTA models 123 for base image captioning: ViT-GPT2<sup>1</sup>, OFA [Wang et al., 124 2022], and BLIP-2 [Li et al., 2023]. For contextualisation, 125 we use GPT-3(.5) [Liu *et al.*, 2023c] and Llama  $3^2$ . 126

LMM LMMs can process visual and text input simultaneously; hence, there is no explicit intermediate caption generation process in this case. BLIP-2, used in the CIC configu-

<sup>&</sup>lt;sup>1</sup>https://huggingface.co/nlpconnect/vit-gpt2-image-captioning <sup>2</sup>https://llama.meta.com/llama3/

			Metrics											
			BLEU				ROUGE			METEOR	BERTScore			SBERT
			1	2	3	4	1	2	L		р	r	f1	obbitti
CIC		BLIP-2	0.373	0.182	0.101	0.060	0.356	0.184	0.292	0.344	0.897	0.900	0.898	0.526
	GPT-3	OFA	0.368	0.178	0.097	0.057	0.352	0.182	0.289	0.341	0.895	0.898	0.896	0.507
		ViT-GPT2	0.363	0.175	0.096	0.057	0.345	0.177	0.283	0.334	0.894	0.897	0.895	0.482
		BLIP-2	0.051	0.020	0.009	0.005	0.119	0.051	0.103	0.211	0.742	0.860	0.796	0.632
	Llama	OFA	0.052	0.020	0.010	0.005	0.119	0.050	0.103	0.211	0.745	0.860	0.797	0.630
		ViT-GPT2	0.062	0.024	0.012	0.006	0.123	0.051	0.101	0.220	0.788	0.865	0.824	0.586
LMM		BLIP-2	0.336	0.137	0.062	0.036	0.304	0.143	0.266	0.194	0.882	0.856	0.868	0.402
		GPT-4v	0.283	0.114	0.061	0.035	0.328	0.151	0.248	0.322	0.872	0.888	0.879	0.505
		LLaVA	0.331	0.132	0.072	0.043	0.283	0.129	0.246	0.260	0.885	0.872	0.878	0.399

Table 3: Similarity between ground truth captions (GoodNews dataset) and those generated with named entity context. BERTScore: microaveraged precision (p), recall (r), and F1 score.

130 ration, can also take textual instructions as input, functioning

as an LMM. We additionally consider two additional LMMs,

namely GPT-4v and LLaVA. All GPT models are provided by

OpenAI, while the others are open-source and publicly avail-

able on the Huggingface platform<sup>3</sup>.

**Prompting** Since LLMs and LMMs have different input requirements, we use two prompt versions, slightly modified. The prompts are present in Table 2. The CIC pipeline must include the base image caption and context to generate the appropriate contextualised caption. In contrast, for the LMM pipeline, only the context must be explicitly provided in the text prompt.

# 142 3.2 Evaluation

143 Dataset We use a subset of the GoodNews dataset [Biten *et al.*, 2019]. This dataset contains images and articles, along
145 with contextualised captions. An example is provided in Ta146 ble 1. We only consider images from 1,000 articles, resulting
147 in a total 1,791 images with their respective captions.

Metrics To assess the quality of the generated captions, we 148 use natural language generation metrics: BLEU [Papineni et 149 al., 2002], ROUGE [Lin, 2004], METEOR [Banerjee and 150 Lavie, 2005], and BERTScore [Zhang et al., 2020]. The first 151 three methods are older and widely used to evaluate natural 152 language generation tasks, such as machine translation and 153 image captioning, relying on n-gram overlaps between ref-154 erence and generated text. BERTScore, on the other hand, 155 was introduced more recently. It leverages the pre-trained 156 contextual embeddings from BERT [Devlin et al., 2019] and 157 matches words between reference and generated text by co-158 sine similarity. We additionally measure sentence embedding 159 similarity with a pre-trained SBERT [Reimers and Gurevych, 160 2019] model. 161

### 162 **3.3 Context types**

As seen in Table 5, we consider two kinds of textual context. In the first case, we extract relevant *named entities* (NE), from the target captions, such as person or organisation names, locations, and time. This way, contextualisation can focus on the information appropriate to the caption. In the second case, we consider the whole *article* as context. This provides a larger amount of information to the systems, which, in turn, might not be relevant. Technically speaking, in a larger application, providing the article as context would equal a less controllable but less labour-intensive approach than providing extracted explicit entities.

## 4 Results and discussion

Table 3 and Table 4 show the results for generating contex-175 tualised captions given NE context and article context, re-176 spectively. In CIC models, we observe a significantly lower 177 performance when using Llama 3. It is interesting, however, 178 that there is no significant difference between CIC with GPT-179 3 and LMM. In this case, the bottleneck caused by the lack 180 of visual information beyond the text caption is not substan-181 tial - probably because context information contributes more 182 to the meaning of the caption than the content. This is par-183 ticularly interesting in the comparison between GPT-4v and 184 LLaVA: BLEU-3 and -4 scores are higher for LLaVA, and 185 BERTScores between the two models do not differ signifi-186 cantly. This indicates that similar results can be achieved both 187 with closed- and open-source models. 188

Focused context matters Results in Table 3 are in almost 189 all cases significantly higher than in their respective cate-190 gories and metrics in Table 4 - in the case of BLEU-3 and 191 -4, which measures trigram and tetragram overlap, scores are 192 close to or equal to zero. This indicates that less is more; 193 since a model is prone to hallucinating and not correctly fol-194 lowing the instruction in the text prompt, a more controlled 195 approach where only needed information must be explicitly 196 mentioned in the caption is more beneficial. However, this 197 might cause more overhead for the domain expert, as the rel-198 evant context must be identified and integrated manually. In 199 this case, an approach facilitating these processes by inter-200 action and/or integration of additional, controllable modules 201 would prove beneficial. 202

#### 4.1 Ablation study

As an additional ablation experiment, we calculate 204 BERTScore values between our reference ground truth 205

174

203

<sup>&</sup>lt;sup>3</sup>https://huggingface.co/

			Metrics											
			BLEU				ROUGE			METEOP	BERTScore			SBEDT
			1	2	3	4	1	2	L	METEOR	р	r	f1	SDERI
CIC		BLIP-2	0.188	0.034	0.013	0.006	0.171	0.044	0.136	0.164	0.858	0.855	0.857	0.323
	GPT-3	OFA	0.114	0.002	0.000	0.000	0.071	0.002	0.059	0.074	0.831	0.824	0.828	0.314
		ViT-GPT2	0.116	0.002	0.000	0.000	0.072	0.002	0.059	0.075	0.831	0.825	0.828	0.304
		BLIP-2	0.046	0.010	0.004	0.002	0.084	0.021	0.067	0.152	0.788	0.844	0.815	0.567
	Llama	OFA	0.033	0.002	0.000	0.000	0.049	0.003	0.042	0.089	0.779	0.815	0.796	0.563
		ViT-GPT2	0.032	0.002	0.000	0.000	0.049	0.003	0.042	0.089	0.778	0.815	0.796	0.561
LMM		BLIP-2	0.195	0.044	0.018	0.009	0.166	0.046	0.142	0.102	0.853	0.838	0.845	0.270
		GPT-4v	0.105	0.003	0.000	0.000	0.092	0.003	0.074	0.097	0.82	0.824	0.822	0.366
		LLaVA	0.131	0.004	0.000	0.000	0.094	0.004	0.08	0.078	0.834	0.824	0.829	0.276

Table 4: Similarity between ground truth captions (GoodNews dataset) and those generated with article context. BERTScore: micro-averaged precision (p), recall (r), and F1 score.

		N	ΙE	article				
	base	GPT	Llama 3	GPT	Llama 3			
BLIP-2 OFA ViT-GPT2	0.841 0.840 0.854	$egin{array}{c} 0.891\uparrow\ 0.889\uparrow\ 0.887\uparrow \end{array}$	$\begin{array}{c} 0.796 \downarrow \\ 0.797 \downarrow \\ 0.824 \downarrow \end{array}$	$egin{array}{c} 0.857\uparrow \ 0.828\downarrow \ 0.828\downarrow \ 0.828\downarrow \ \end{array}$	$\begin{array}{c} 0.815 \downarrow \\ 0.796 \downarrow \\ 0.796 \downarrow \end{array}$			

Table 5: Ablation study: F1 BERTScores for base captions, compared to contextualised ones (given different contexts and post-hoc contextualisation LLMs).

captions and base captions generated by our three image 206 captioning models of choice. Results for experiments with 207 both contexts are present in Table 5. We expect that the 208 base captions would perform worse than their contextu-209 alised counterparts. However, this is only the case with 210 GPT-3 in combination NE context - compared to the base, 211 non-contextualised captions, scores for the contextualised 212 captions with Llama 3 is lower. This indicates Llama 3's 213 inability to follow the prompt's instructions as expected, 214 which leads to caption generations containing irrelevant 215 information, lowering the automated metric scores. 216

#### 217 4.2 Limitations

We identify two significant limitations within our work. The 218 first is the lack of diversity in prompt usage. We experiment 219 with a single prompt as present in Table 2. Especially in cases 220 like CIC with a Llama 3 contextualisation module, experi-221 mentation with prompts might be beneficial and lead to im-222 proved results. The second limitation is related to the metrics 223 we use. The "older" methods (BLEU, ROUGE, METEOR) 224 might not reward the existence of synonyms and paraphrases, 225 as they focus on exact matches and their order. A domain 226 expert, in this case, a journalist, might consider one of the pe-227 nalised captions just as well-formed as its ground truth equiv-228 alent. On the other hand, BERTScore might be rather forgiv-229 ing - hence the higher scores in the tables. When accuracy 230 is required, however, the scores might not capture the differ-231 ence in generated quality. Thus, a user study is necessary to 232 evaluate the presented pipelines' performance. 233

# 5 Conclusion and future work

The appearance of LLMs and LLMs has rendered human-AI 235 synergy more accessible. This work presents a use case for 236 journalism: generating relevant, contextualised image cap-237 tions given different pipelines (called CIC and LMM), in-238 cluding pre-trained image captioning models, LLMs, and 239 LMMs. We evaluate our experiments with automated met-240 rics and conclude that, at least regarding these metrics, the 241 bottleneck caused by using a CIC-like architecture with a 242 textual description of the image rather than the image itself 243 is insignificant. Close-source models such as the GPT fam-244 ily might have an advantage in the CIC configuration. How-245 ever, smaller, open-source models perform similarly well in 246 the LLM configuration. 247

In terms of context, focused information, such as NEs, 248 is more beneficial to the models than the whole article it-249 self. This finding indicated a possible future direction for 250 our work: implementing an interactive system that facilitates 251 journalists' writing captions for their articles. Additionally, 252 we would like to expand our contextualised image captioning 253 experiments to include more datasets, and address the limi-254 tations stated in Section 4, by experimenting with different 255 prompt patterns to increase the efficiency of the proposed ar-256 chitectures and by conducting a user study for a more nu-257 anced evaluation of the quality of the generated captions. 258

## **Ethical Statement**

We have carefully considered the ethical implications of our 260 work and do not foresee any major concerns. The dataset 261 and models we utilize are publicly available. However, we 262 acknowledge the potential risks of disseminating incorrect 263 information, which could lead to the misuse of LLMs and 264 LMMs. This is an important limitation of our research that 265 we strive to mitigate through rigorous evaluation and respon-266 sible usage guidelines. 267

# References

[Banerjee and Lavie, 2005] Satanjeev Banerjee and Alon 269 Lavie. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In 271

259

234

268

272 Proceedings of the ACL Workshop on Intrinsic and Extrin-

273 sic Evaluation Measures for Machine Translation and/or

- 274 Summarization, pages 65–72, Ann Arbor, Michigan, June
- 275 2005. ACL.
- [Bao et al., 2024] Yujia Bao, Ankit Parag Shah, Neeru 276 Narang, Jonathan Rivers, Rajeev Maksey, Lan Guan, 277 Louise N. Barrere, Shelley Evenson, Rahul Basole, Con-278 nie Miao, Ankit Mehta, Fabien Boulay, Su Min Park, 279 Natalie E. Pearson, Eldhose Joy, Tiger He, Sumiran 280 Thakur, Koustav Ghosal, Josh On, Phoebe Morrison, 281 Tim Major, Eva Siqi Wang, Gina Escobar, Jiaheng 282 Wei, Tharindu Cyril Weerasooriya, Queena Song, Daria 283 Lashkevich, Clare Chen, Gyuhak Kim, Dengpan Yin, Don 284 Hejna, Mo Nomeli, and Wei Wei. Harnessing business and 285 media insights with large language models. 2024. 286
- [Biten *et al.*, 2019] Ali Furkan Biten, Lluís Gómez, Marçal Rusiñol, and Dimosthenis Karatzas. Good News, Everyone! Context Driven Entity-Aware Captioning for News Images. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 12466–12475. Computer Vision Foundation / IEEE, 2019.
- <sup>294</sup> [Bubeck *et al.*, 2023] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz,
  <sup>296</sup> Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott
  <sup>297</sup> Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint*<sup>299</sup> *arXiv:2303.12712*, 2023.
- [Cheng *et al.*, 2024] Liqi Cheng, Dazhen Deng, Xiao Xie,
   Rihong Qiu, Mingliang Xu, and Yingcai Wu. Snil: Gener ating sports news from insights with large language mod *IEEE transactions on visualization and computer graphics*, PP, 2024.
- [del Barrio and Gática-Pérez, 2023] David Alonso del Barrio and Daniel Gática-Pérez. Framing the news: From human perception to large language model inferences. *Proceedings of the 2023 ACM International Conference on*
- 309 *Multimedia Retrieval*, 2023.
- [Devlin et al., 2019] Jacob Devlin, Ming-Wei Chang, Ken-310 ton Lee, and Kristina Toutanova. BERT: Pre-training of 311 deep bidirectional transformers for language understand-312 ing. In Jill Burstein, Christy Doran, and Thamar Solorio, 313 editors, Proceedings of the 2019 Conference of the North 314 American Chapter of the Association for Computational 315 Linguistics: Human Language Technologies, Volume 1 316 (Long and Short Papers), pages 4171-4186, Minneapo-317 lis, Minnesota, June 2019. Association for Computational 318 Linguistics. 319
- [Fang *et al.*, 2023] Xiao Fang, Shangkun Che, Minjia Mao,
   Hongzhe Zhang, Ming Zhao, and Xiaohang Zhao. Bias of
   AI-generated content: an examination of news produced
- by large language models. *Scientific Reports*, 14, 2023.
- [Hamilton and Piper, 2022] Silvia Durianova Hamilton and
   Andrew Piper. The covid that wasn't: Counterfactual journalism using gpt. *ArXiv*, abs/2210.06644, 2022.

- [Li et al., 2023] Junnan Li, Dongxu Li, Silvio Savarese, and 327 Steven C. H. Hoi. BLIP-2: Bootstrapping Language-328 Image Pre-training with Frozen Image Encoders and Large 329 Language Models. In Andreas Krause, Emma Brunskill, 330 Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and 331 Jonathan Scarlett, editors, International Conference on 332 Machine Learning, ICML 2023, 23-29 July 2023, Hon-333 olulu, Hawaii, USA, volume 202 of Proceedings of Ma-334 chine Learning Research, pages 19730-19742. PMLR, 335 2023. 336
- [Lin, 2004] Chin-Yew Lin. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization* 338
   *Branches Out*, pages 74–81, Barcelona, Spain, July 2004. 339
   ACL. 340
- [Liu *et al.*, 2023a] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023. 343
- [Liu *et al.*, 2023b] Haotian Liu, Chunyuan Li, Qingyang 344 Wu, and Yong Jae Lee. Visual instruction tuning, 2023. 345
- [Liu *et al.*, 2023c] Xin Liu, Daniel McDuff, Geza Kovacs, 346
  Isaac R. Galatzer-Levy, Jacob E. Sunshine, Jiening Zhan, 347
  Ming-Zher Poh, Shun Liao, Paolo Di Achille, and Shwetak N. Patel. Large language models are few-shot health learners. *CoRR*, abs/2305.15525, 2023. 350
- [Liu *et al.*, 2024] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024. 354
- [Nguyen et al., 2023] Khanh Nguyen, Ali Furkan Biten, 355 Andrés Mafla, Lluís Gómez, and Dimosthenis Karatzas. 356 Show, Interpret and Tell: Entity-Aware Contextualised Im-357 age Captioning in Wikipedia. In Brian Williams, Yiling 358 Chen, and Jennifer Neville, editors, Thirty-Seventh AAAI 359 Conference on Artificial Intelligence, AAAI 2023, Thirty-360 Fifth Conference on Innovative Applications of Artificial 361 Intelligence, IAAI 2023, Thirteenth Symposium on Edu-362 cational Advances in Artificial Intelligence, EAAI 2023, 363 Washington, DC, USA, February 7-14, 2023, pages 1940-364 1948. AAAI Press, 2023. 365
- [OpenAI, 2023] OpenAI. GPT-4 Technical Report, 2023. 366 \_eprint: 2303.08774. 367
- [Papineni et al., 2002] Kishore Papineni, Salim Roukos, 368
  Todd Ward, and Wei-Jing Zhu. Bleu: a Method for Automatic Evaluation of Machine Translation. In Proceedings 370
  of the 40th Annual Meeting of the ACL, pages 311–318, 371
  Philadelphia, Pennsylvania, USA, July 2002. ACL. 372
- [Qu *et al.*, 2023] Tingyu Qu, Tinne Tuytelaars, and Marie-Francine Moens. Visually-aware context modeling for news image captioning. *CoRR*, abs/2308.08325, 2023. 375
- [Rajakumar Kalarani *et al.*, 2023] Abisek Rajaku mar Kalarani, Pushpak Bhattacharyya, Niyati Chhaya, and
   Sumit Shekhar. "let's not quote out of context": Unified
   vision-language pretraining for context assisted image
   captioning. In Sunayana Sitaram, Beata Beigman Kle banov, and Jason D Williams, editors, *Proceedings of* 381

the 61st Annual Meeting of the Association for Compu-

tational Linguistics (Volume 5: Industry Track), pages
 695–706, Toronto, Canada, July 2023. Association for

Computational Linguistics.
[Reimers and Gurevych, 2019] Nils Reimers and Iryna

Gurevych. Sentence-BERT: Sentence embeddings using 387 Siamese BERT-networks. In Kentaro Inui, Jing Jiang, 388 Vincent Ng, and Xiaojun Wan, editors, Proceedings 389 of the 2019 Conference on Empirical Methods in Nat-390 ural Language Processing and the 9th International 391 Joint Conference on Natural Language Processing 392 (EMNLP-IJCNLP), pages 3982–3992, Hong Kong, 393 China, November 2019. Association for Computational 394 Linguistics. 395

[Srinivasan *et al.*, 2021] Krishna Srinivasan, Karthik Raman, Jiecao Chen, Michael Bendersky, and Marc Najork.
WIT: Wikipedia-based Image Text Dataset for Multimodal
Multilingual Machine Learning. In Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, and
Tetsuya Sakai, editors, *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in*

403 Information Retrieval, Virtual Event, Canada, July 11-15,

<sup>404</sup> 2021, pages 2443–2449. ACM, 2021.

[Tran *et al.*, 2020] Alasdair Tran, Alexander Patrick Mathews, and Lexing Xie. Transform and Tell: Entity-Aware
News Image Captioning. In 2020 IEEE/CVF Conference

408 on Computer Vision and Pattern Recognition, CVPR 2020,
 409 Seattle, WA, USA, June 13-19, 2020, pages 13032–13042.

410 Computer Vision Foundation / IEEE, 2020.

[Vaswani *et al.*, 2017] Ashish Vaswani, Noam Shazeer, Niki
Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
Łukasz Kaiser, and Illia Polosukhin. Attention is all you

414 need. In Proceedings of the 31st International Confer-

ence on Neural Information Processing Systems, NIPS'17,

page 6000–6010, Red Hook, NY, USA, 2017. Curran Associates Inc.

[Wang *et al.*, 2022] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. OFA: Unifying Architectures, Tasks, and Modalities Through a Simple Sequenceto-Sequence Learning Framework. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang

- Niu, and Sivan Sabato, editors, International Conference
   on Machine Learning, ICML 2022, 17-23 July 2022, Bal-
- *timore, Maryland, USA*, volume 162 of *Proceedings of*

Machine Learning Research, pages 23318–23340. PMLR, 2022.

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

435 [Zou et al., 2023] Xueyan Zou, Zi-Yi Dou, Jianwei Yang,

- 436 Zhe Gan, Linjie Li, Chunyuan Li, Xiyang Dai, Harki-
- rat Behl, Jianfeng Wang, Lu Yuan, Nanyun Peng, Lijuan

438 Wang, Yong Jae Lee, and Jianfeng Gao. Generalized

decoding for pixel, image, and language. In *IEEE/CVF* 439 *Conference on Computer Vision and Pattern Recognition*, 440 *CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, 441 pages 15116–15127. IEEE, 2023. 442