

306 **References**

- 307 [1] Aishwarya Agrawal, Jiasen Lu, Stanislaw Antol, Margaret Mitchell, C. Lawrence Zitnick,
308 Devi Parikh, and Dhruv Batra. Vqa: Visual question answering. *International Journal of*
309 *Computer Vision*, 123:4–31, 2015. 9, 23
- 310 [2] J Alammr. Ecco: An open source library for the explainability of transformer language
311 models. In *Proceedings of the 59th Annual Meeting of the Association for Computational*
312 *Linguistics and the 11th International Joint Conference on Natural Language Processing:*
313 *System Demonstrations*, pages 249–257, Online, August 2021. Association for Computational
314 Linguistics. 8
- 315 [3] Saeed Amizadeh, Hamid Palangi, Alex Polozov, Yichen Huang, and Kazuhito Koishida.
316 Neuro-symbolic visual reasoning: Disentangling. In *International Conference on Machine*
317 *Learning*, pages 279–290. PMLR, 2020. 23
- 318 [4] John R Anderson, C Franklin Boyle, and Brian J Reiser. Intelligent tutoring systems. *Science*,
319 228(4698):456–462, 1985. 2
- 320 [5] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von
321 Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the
322 opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021. 21
- 323 [6] Alexander Bondarenko, Magdalena Wolska, Stefan Heindorf, Lukas Blübaum, Axel-
324 Cyrille Ngonga Ngomo, Benno Stein, Pavel Braslavski, Matthias Hagen, and Martin Potthast.
325 Causalqa: A benchmark for causal question answering. In *ACL*, 2022. 9, 23
- 326 [7] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-
327 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language
328 models are few-shot learners. *Advances in neural information processing systems*, 33:1877–
329 1901, 2020. 2, 6, 21, 23, 24
- 330 [8] Angelos Chatzimparmpas, Rafael M Martins, Kostiantyn Kucher, and Andreas Kerren. Stack-
331 genvis: Alignment of data, algorithms, and models for stacking ensemble learning using perfor-
332 mance metrics. *IEEE Transactions on Visualization and Computer Graphics*, 27(2):1547–1557,
333 2020. 9
- 334 [9] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto,
335 Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating
336 large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021. 23
- 337 [10] Zhuo Chen, Yufen Huang, Jiaoyan Chen, Yuxia Geng, Yin Fang, Jeff Z. Pan, Ningyu Zhang,
338 and Wen Zhang. Lako: Knowledge-driven visual question answering via late knowledge-to-text
339 injection. *ArXiv*, abs/2207.12888, 2022. 9, 24
- 340 [11] Eunsol Choi, Daniel Hewlett, Jakob Uszkoreit, Illia Polosukhin, Alexandre Lacoste, and
341 Jonathan Berant. Coarse-to-fine question answering for long documents. In *ACL*, 2017. 9, 23
- 342 [12] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
343 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:
344 Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022. 2, 9, 24
- 345 [13] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li,
346 Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned
347 language models. *arXiv preprint arXiv:2210.11416*, 2022. 3, 4, 7, 8, 9, 24
- 348 [14] Jonathan H Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly
349 Nikolaev, and Jennimaria Palomaki. Tydi qa: A benchmark for information-seeking question
350 answering in typologically diverse languages. *Transactions of the Association for Computa-*
351 *tional Linguistics*, 2020. 23
- 352 [15] Misha Denil, Alban Demiraj, and Nando De Freitas. Extraction of salient sentences from
353 labelled documents. *arXiv preprint arXiv:1412.6815*, 2014. 8

- 354 [16] Thomas G Dietterich. Ensemble methods in machine learning. In *International workshop on*
355 *multiple classifier systems*, pages 1–15. Springer, 2000. 9, 24
- 356 [17] Virginie Do, Oana-Maria Camburu, Zeynep Akata, and Thomas Lukasiewicz. e-snl-ve:
357 Corrected visual-textual entailment with natural language explanations. *arXiv preprint*
358 *arXiv:2004.03744*, 2020. 2, 18
- 359 [18] Alican Dogan and Derya Birant. A weighted majority voting ensemble approach for classifica-
360 tion. In *2019 4th International Conference on Computer Science and Engineering (UBMK)*,
361 pages 1–6, 2019. 2, 9
- 362 [19] Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli.
363 Eli5: Long form question answering. *arXiv preprint arXiv:1907.09190*, 2019. 23
- 364 [20] Tianyu Gao, Adam Fisch, and Danqi Chen. Making pre-trained language models better
365 few-shot learners. *arXiv preprint arXiv:2012.15723*, 2020. 23
- 366 [21] Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did
367 aristotle use a laptop? a question answering benchmark with implicit reasoning strategies.
368 *Transactions of the Association for Computational Linguistics*, 2021. 23
- 369 [22] Dan Goldwasser and Dan Roth. Learning from natural instructions. *Machine learning*,
370 94(2):205–232, 2014. 23
- 371 [23] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making
372 the v in vqa matter: Elevating the role of image understanding in visual question answering.
373 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages
374 6904–6913, 2017. 2, 18
- 375 [24] Liangke Gui, Borui Wang, Qiuyuan Huang, Alex Hauptmann, Yonatan Bisk, and Jianfeng
376 Gao. Kat: A knowledge augmented transformer for vision-and-language. *arXiv preprint*
377 *arXiv:2112.08614*, 2021. 9, 24
- 378 [25] Jiaxian Guo, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Boyang Li, Dacheng Tao, and
379 Steven CH Hoi. From images to textual prompts: Zero-shot vqa with frozen large language
380 models. *arXiv preprint arXiv:2212.10846*, 2022. 4
- 381 [26] Yaru Hao, Haoyu Song, Li Dong, Shaohan Huang, Zewen Chi, Wenhui Wang, Shum-
382 ing Ma, and Furu Wei. Language models are general-purpose interfaces. *arXiv preprint*
383 *arXiv:2206.06336*, 2022. 4
- 384 [27] Yushi Hu, Hang Hua, Zhengyuan Yang, Weijia Shi, Noah A Smith, and Jiebo Luo. Promptcap:
385 Prompt-guided task-aware image captioning. *arXiv preprint arXiv:2211.09699*, 2022. 4, 19
- 386 [28] Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao
387 Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, et al. Language is not all you need: Aligning
388 perception with language models. *arXiv preprint arXiv:2302.14045*, 2023. 24
- 389 [29] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual
390 reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference*
391 *on computer vision and pattern recognition*, pages 6700–6709, 2019. 9, 23
- 392 [30] Gabriel Ilharco, Mitchell Wortsman, Samir Yitzhak Gadre, Shuran Song, Hannaneh Hajishirzi,
393 Simon Kornblith, Ali Farhadi, and Ludwig Schmidt. Patching open-vocabulary models by
394 interpolating weights. *arXiv preprint arXiv:2208.05592*, 2022. 2, 9
- 395 [31] Anmol Jain, Aishwary Kumar, and Seba Susan. Evaluating deep neural network ensembles
396 by majority voting cum meta-learning scheme. In *Soft Computing and Signal Processing: Proceedings of 3rd ICSCSP 2020, Volume 2*, pages 29–37. Springer, 2022. 9
- 398 [32] Anya Ji, Noriyuki Kojima, Noah Rush, Alane Suhr, Wai Keen Vong, Robert D Hawkins, and
399 Yoav Artzi. Abstract visual reasoning with tangram shapes. In *EMNLP*, 2022. 9

- 400 [33] Woojeong Jin, Yu Cheng, Yelong Shen, Weizhu Chen, and Xiang Ren. A good prompt is worth
401 millions of parameters: Low-resource prompt-based learning for vision-language models. In
402 *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*
403 *(Volume 1: Long Papers)*, pages 2763–2775, Dublin, Ireland, May 2022. Association for
404 Computational Linguistics. 4
- 405 [34] Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick,
406 and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary
407 visual reasoning. In *Proceedings of the IEEE conference on computer vision and pattern*
408 *recognition*, pages 2901–2910, 2017. 2, 9, 23
- 409 [35] Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale
410 distantly supervised challenge dataset for reading comprehension. In *ACL*, 2017. 23
- 411 [36] Maxime Kayser, Oana-Maria Camburu, Leonard Salewski, Cornelius Emde, Virginie Do,
412 Zeynep Akata, and Thomas Lukasiewicz. e-vil: A dataset and benchmark for natural lan-
413 guage explanations in vision-language tasks. In *Proceedings of the IEEE/CVF International*
414 *Conference on Computer Vision*, pages 1244–1254, 2021. 4
- 415 [37] Daniel Khashabi, Yeganeh Kordi, and Hannaneh Hajishirzi. Unifiedqa-v2: Stronger general-
416 ization via broader cross-format training. *arXiv preprint arXiv:2202.12359*, 2022. 4
- 417 [38] Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without
418 convolution or region supervision. In *International Conference on Machine Learning*, pages
419 5583–5594. PMLR, 2021. 4
- 420 [39] Mike Lewis, Shruti Bhosale, Tim Dettmers, Naman Goyal, and Luke Zettlemoyer. Base
421 layers: Simplifying training of large, sparse models. In *International Conference on Machine*
422 *Learning*, pages 6265–6274. PMLR, 2021. 24
- 423 [40] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed,
424 Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence
425 pre-training for natural language generation, translation, and comprehension. *arXiv preprint*
426 *arXiv:1910.13461*, 2019. 23
- 427 [41] Bo Li, Yifei Shen, Jingkang Yang, Yezhen Wang, Jiawei Ren, Tong Che, Jun Zhang, and
428 Ziwei Liu. Sparse mixture-of-experts are domain generalizable learners. *arXiv preprint*
429 *arXiv:2206.04046*, 2022. 24
- 430 [42] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter:
431 A multi-modal model with in-context instruction tuning, 2023. 2
- 432 [43] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-
433 image pre-training with frozen image encoders and large language models. *arXiv preprint*
434 *arXiv:2301.12597*, 2023. 2, 4
- 435 [44] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-
436 image pre-training for unified vision-language understanding and generation. *arXiv preprint*
437 *arXiv:2201.12086*, 2022. 2, 3, 9, 18, 24
- 438 [45] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A
439 simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*,
440 2019. 4
- 441 [46] Shuang Li, Yilun Du, Joshua B. Tenenbaum, Antonio Torralba, and Igor Mordatch. Composing
442 ensembles of pre-trained models via iterative consensus. *ArXiv*, abs/2210.11522, 2022. 9, 24
- 443 [47] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan,
444 Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In
445 *European conference on computer vision*, pages 740–755. Springer, 2014. 18
- 446 [48] Fangyu Liu, Guy Emerson, and Nigel Collier. Visual spatial reasoning. *arXiv preprint*
447 *arXiv:2205.00363*, 2022. 2, 18

- 448 [49] Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal,
449 and Colin Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context
450 learning. *arXiv preprint arXiv:2205.05638*, 2022. 20, 23
- 451 [50] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic
452 visiolinguistic representations for vision-and-language tasks. *Advances in neural information
453 processing systems*, 32, 2019. 4
- 454 [51] Jiasen Lu, Christopher Clark, Rowan Zellers, Roozbeh Mottaghi, and Aniruddha Kembhavi.
455 Unified-io: A unified model for vision, language, and multi-modal tasks. *arXiv preprint
456 arXiv:2206.08916*, 2022. 4
- 457 [52] Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-
458 Chun Zhu, and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large
459 language models. *arXiv preprint arXiv:2304.09842*, 2023. 2
- 460 [53] Mikołaj Mańkiński and Jacek Mańdziuk. Deep learning methods for abstract visual reasoning:
461 A survey on raven’s progressive matrices. *arXiv preprint arXiv:2201.12382*, 2022. 23
- 462 [54] Mikołaj Mańkiński and Jacek Mańdziuk. A review of emerging research directions in abstract
463 visual reasoning. *arXiv preprint arXiv:2202.10284*, 2022. 2, 23
- 464 [55] Jiayuan Mao, Chuang Gan, Pushmeet Kohli, Joshua B. Tenenbaum, and Jiajun Wu. The neuro-
465 symbolic concept learner: Interpreting scenes words and sentences from natural supervision.
466 *ArXiv*, abs/1904.12584, 2019. 2, 9, 24
- 467 [56] Kenneth Marino, Xinlei Chen, Devi Parikh, Abhinav Gupta, and Marcus Rohrbach. Krisp:
468 Integrating implicit and symbolic knowledge for open-domain knowledge-based vqa. In
469 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
470 14111–14121, 2021. 9, 24
- 471 [57] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A
472 visual question answering benchmark requiring external knowledge. In *Proceedings of the
473 IEEE/cvf conference on computer vision and pattern recognition*, pages 3195–3204, 2019. 2,
474 18
- 475 [58] John McCarthy et al. *Programs with common sense*. RLE and MIT computation center
476 Cambridge, MA, USA, 1960. 23
- 477 [59] Fundamental AI Research Diplomacy Team Meta, Anton Bakhtin, Noam Brown, Emily Dinan,
478 Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu,
479 et al. Human-level play in the game of diplomacy by combining language models with strategic
480 reasoning. *Science*, 2022. 2, 9, 24
- 481 [60] Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. Metaicl: Learning to
482 learn in context. *arXiv preprint arXiv:2110.15943*, 2021. 23
- 483 [61] Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. Cross-task general-
484 ization via natural language crowdsourcing instructions. *arXiv preprint arXiv:2104.08773*,
485 2021. 23
- 486 [62] Hyacinth S Nwana. Intelligent tutoring systems: an overview. *Artificial Intelligence Review*,
487 4(4):251–277, 1990. 2
- 488 [63] OpenAI. Gpt-4 technical report. 2023. 2
- 489 [64] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin,
490 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models
491 to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022. 21, 23
- 492 [65] Jae Sung Park, Chandra Bhagavatula, Roozbeh Mottaghi, Ali Farhadi, and Yejin Choi. Visu-
493 alcomet: Reasoning about the dynamic context of a still image. In *European Conference on
494 Computer Vision*, 2020. 2, 24

- 495 [66] Silviu Pitis, Michael R Zhang, Andrew Wang, and Jimmy Ba. Boosted prompt ensembles for
496 large language models. *arXiv preprint arXiv:2304.05970*, 2023. [2](#)
- 497 [67] Björn Plüster, Jakob Ambsdorf, Lukas Braach, Jae Hee Lee, and Stefan Wermter. Harnessing
498 the power of multi-task pretraining for ground-truth level natural language explanations. *arXiv*
499 *preprint arXiv:2212.04231*, 2022. [4](#), [19](#)
- 500 [68] Chen Qu, Liu Yang, Cen Chen, Minghui Qiu, W Bruce Croft, and Mohit Iyyer. Open-retrieval
501 conversational question answering. In *ACM SIGIR*, 2020. [23](#)
- 502 [69] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
503 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
504 models from natural language supervision. In *International Conference on Machine Learning*,
505 pages 8748–8763. PMLR, 2021. [4](#), [9](#), [24](#)
- 506 [70] Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song,
507 John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language
508 models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*,
509 2021. [24](#)
- 510 [71] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,
511 Yanqi Zhou, Wei Li, Peter J Liu, et al. Exploring the limits of transfer learning with a unified
512 text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67, 2020. [4](#), [9](#), [24](#)
- 513 [72] Nazneen Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. Explain yourself!
514 leveraging language models for commonsense reasoning. In *ACL*, 2019. [23](#)
- 515 [73] Siva Reddy, Danqi Chen, and Christopher D Manning. Coqa: A conversational question
516 answering challenge. *Transactions of the Association for Computational Linguistics*, 2019. [23](#)
- 517 [74] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-
518 networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language*
519 *Processing*. Association for Computational Linguistics, 11 2019. [3](#)
- 520 [75] Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André
521 Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts.
522 *Advances in Neural Information Processing Systems*, 34:8583–8595, 2021. [24](#)
- 523 [76] Omer Sagi and Lior Rokach. Ensemble learning: A survey. *Wiley Interdisciplinary Reviews:*
524 *Data Mining and Knowledge Discovery*, 8(4):e1249, 2018. [9](#), [24](#)
- 525 [77] Shailaja Keyur Sampat, Maitreya Patel, Subhasish Das, Yezhou Yang, and Chitta Baral.
526 Reasoning about actions over visual and linguistic modalities: A survey. *arXiv preprint*
527 *arXiv:2207.07568*, 2022. [9](#), [23](#)
- 528 [78] Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai,
529 Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. Multitask prompted
530 training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207*, 2021. [23](#)
- 531 [79] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettle-
532 moyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach
533 themselves to use tools. *arXiv preprint arXiv:2302.04761*, 2023. [2](#), [9](#)
- 534 [80] Timo Schick and Hinrich Schütze. Exploiting cloze questions for few shot text classification
535 and natural language inference. *arXiv preprint arXiv:2001.07676*, 2020. [23](#)
- 536 [81] Benedikt Schmidt, Reuben Borrison, Andrew Cohen, Marcel Dix, Marco Gärtler, Martin
537 Hollender, Benjamin Klöpper, Sylvia Maczey, and Shunmuga Siddharthan. Industrial virtual
538 assistants: Challenges and opportunities. In *Proceedings of the 2018 ACM International Joint*
539 *Conference and 2018 International Symposium on Pervasive and Ubiquitous Computing and*
540 *Wearable Computers*, pages 794–801, 2018. [2](#)
- 541 [82] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh
542 Mottaghi. A-okvqa: A benchmark for visual question answering using world knowledge.
543 *arXiv preprint arXiv:2206.01718*, 2022. [2](#), [18](#)

- 544 [83] Zhenwei Shao, Zhou Yu, Meng Wang, and Jun Yu. Prompting large language models
545 with answer heuristics for knowledge-based visual question answering. *arXiv preprint*
546 *arXiv:2303.01903*, 2023. 4
- 547 [84] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarczyk, Andy Davis, Quoc Le, Geoffrey Hinton,
548 and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts
549 layer. *arXiv preprint arXiv:1701.06538*, 2017. 24
- 550 [85] Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory
551 cost. In *International Conference on Machine Learning*, pages 4596–4604. PMLR, 2018. 18
- 552 [86] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang.
553 Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface. *arXiv preprint*
554 *arXiv:2303.17580*, 2023. 2, 9
- 555 [87] Derek Sleeman and John Seely Brown. *Intelligent tutoring systems*. London: Academic Press,
556 1982. 2
- 557 [88] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec
558 Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback.
559 *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020. 22
- 560 [89] Md Arafat Sultan, Shubham Chandel, Ramón Fernández Astudillo, and Vittorio Castelli. On
561 the importance of diversity in question generation for qa. In *ACL*, 2020. 23
- 562 [90] Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution
563 for reasoning. *arXiv preprint arXiv:2303.08128*, 2023. 2
- 564 [91] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from
565 transformers. *arXiv preprint arXiv:1908.07490*, 2019. 4
- 566 [92] Anthony Meng Huat Tiong, Junnan Li, Boyang Li, Silvio Savarese, and Steven CH Hoi.
567 Plug-and-play vqa: Zero-shot vqa by conjoining large pretrained models with zero training.
568 *arXiv preprint arXiv:2210.08773*, 2022. 4
- 569 [93] Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris, Alessandro Sordoni, Philip Bachman,
570 and Kaheer Suleman. Newsqa: A machine comprehension dataset. In *Rep4NLP@ACL*, 2016.
571 9, 23
- 572 [94] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: Lessons
573 learned from the 2015 mscoco image captioning challenge. *IEEE transactions on pattern*
574 *analysis and machine intelligence*, 39(4):652–663, 2016. 2
- 575 [95] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou,
576 Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a
577 simple sequence-to-sequence learning framework. In *International Conference on Machine*
578 *Learning*, pages 23318–23340. PMLR, 2022. 2, 3, 9, 18, 24
- 579 [96] Wenhui Wang, Hangbo Bao, Li Dong, and Furu Wei. Vlm0: Unified vision-language pre-
580 training with mixture-of-modality-experts. *ArXiv*, abs/2111.02358, 2021. 9, 24
- 581 [97] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. Rationale-
582 augmented ensembles in language models. *arXiv preprint arXiv:2207.00747*, 2022. 2
- 583 [98] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. Self-
584 consistency improves chain of thought reasoning in language models. *arXiv preprint*
585 *arXiv:2203.11171*, 2022. 2
- 586 [99] Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan
587 Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv*
588 *preprint arXiv:2109.01652*, 2021. 7, 8, 23
- 589 [100] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny
590 Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint*
591 *arXiv:2201.11903*, 2022. 2

- 592 [101] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony
593 Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer,
594 Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain
595 Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art
596 natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in
597 Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020.
598 Association for Computational Linguistics. [3](#)
- 599 [102] Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-
600 Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and
601 Ludwig Schmidt. Model soups: averaging weights of multiple fine-tuned models improves
602 accuracy without increasing inference time. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song,
603 Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International
604 Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*,
605 pages 23965–23998. PMLR, 17–23 Jul 2022. [2](#), [9](#)
- 606 [103] Ning Xie, Farley Lai, Derek Doran, and Asim Kadav. Visual entailment: A novel task for
607 fine-grained image understanding. *ArXiv*, abs/1901.06706, 2019. [9](#), [18](#), [23](#)
- 608 [104] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of
609 in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021. [5](#)
- 610 [105] Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Yumao Lu, Zicheng Liu, and Lijuan
611 Wang. An empirical study of gpt-3 for few-shot knowledge-based vqa. In *AAAI*, 2022. [2](#), [24](#)
- 612 [106] Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and
613 Joshua B. Tenenbaum. Clevrer: Collision events for video representation and reasoning. *ArXiv*,
614 abs/1910.01442, 2019. [2](#), [24](#)
- 615 [107] Kexin Yi, Jiajun Wu, Chuang Gan, Antonio Torralba, Pushmeet Kohli, and Joshua B. Tenen-
616 baum. Neural-symbolic vqa: Disentangling reasoning from vision and language understanding.
617 *ArXiv*, abs/1810.02338, 2018. [2](#), [9](#), [24](#)
- 618 [108] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. Deep modular co-attention networks
619 for visual question answering. In *Proceedings of the IEEE/CVF conference on computer vision
620 and pattern recognition*, pages 6281–6290, 2019. [4](#)
- 621 [109] Munazza Zaib, Wei Emma Zhang, Quan Z Sheng, Adnan Mahmood, and Yang Zhang. Con-
622 versational question answering: A survey. *Knowledge and Information Systems*, 2022. [9](#),
623 [23](#)
- 624 [110] Rufai Yusuf Zakari, Jim Wilson Owusu, Hailin Wang, Ke Qin, Zaharaddeen Karami Lawal,
625 and Yuezhou Dong. Vqa and visual reasoning: An overview of recent datasets, methods and
626 challenges. *arXiv preprint arXiv:2212.13296*, 2022. [2](#), [9](#), [23](#)
- 627 [111] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition:
628 Visual commonsense reasoning. *2019 IEEE/CVF Conference on Computer Vision and Pattern
629 Recognition (CVPR)*, pages 6713–6724, 2018. [2](#), [9](#), [24](#)
- 630 [112] Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi,
631 and Yejin Choi. Merlot: Multimodal neural script knowledge models. In *Neural Information
632 Processing Systems*, 2021. [9](#), [24](#)
- 633 [113] Andy Zeng, Adrian S. Wong, Stefan Welker, Krzysztof Choromanski, Federico Tombari,
634 Aveek Purohit, Michael S. Ryoo, Vikas Sindhwani, Johnny Lee, Vincent Vanhoucke, and
635 Peter R. Florence. Socratic models: Composing zero-shot multimodal reasoning with language.
636 *ArXiv*, abs/2204.00598, 2022. [2](#), [9](#), [24](#)
- 637 [114] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen,
638 Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained
639 transformer language models. *arXiv preprint arXiv:2205.01068*, 2022. [4](#), [6](#)

- 640 [115] Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew M Dai,
641 Quoc V Le, James Laudon, et al. Mixture-of-experts with expert choice routing. *Advances in*
642 *Neural Information Processing Systems*, 35:7103–7114, 2022. 24
- 643 [116] Fengbin Zhu, Wenqiang Lei, Chao Wang, Jianming Zheng, Soujanya Poria, and Tat-Seng
644 Chua. Retrieving and reading: A comprehensive survey on open-domain question answering.
645 *arXiv preprint arXiv:2101.00774*, 2021. 23

646 A Experimental Details

647 In this section, we elaborate on our training and evaluation details, prompt templates, and more
648 qualitative examples for analysis.

649 A.1 Datasets

650 Our experiments are conducted on a challenging suite of three diverse visual reasoning tasks, including
651 outside knowledge VQA, visual entailment, and visual spatial reasoning. For each task, we select the
652 following dataset respectively.

653 **Visual Question Answering v2** [23] (VQA v2) is a large-scale benchmark containing over 1
654 million images from the COCO dataset and more than 250,000 human-generated question-answer
655 pairs. The dataset is designed to test the ability of machine learning models to understand both the
656 visual content of an image and the meaning behind natural language questions. The questions in VQA
657 v2 cover a wide range of topics and are often open-ended, requiring models to reason and generalize
658 about the world. VQA v2 has been widely used to evaluate the performance of state-of-the-art models
659 in the field of computer vision and natural language processing.

660 **Augmented Outside Knowledge VQA** [82] (A-OKVQA) contains about 25k questions paired
661 with both multiple choice (MC) answer options. Unlike most existing VQA datasets, the questions in
662 A-OKVQA cannot often be answered by querying the knowledge base, but rather involve some type
663 of commonsense reasoning and outside knowledge about the situation portrayed in the image.

664 **Outside Knowledge VQA** [57] (OK-VQA) includes more than 14,000 questions that require
665 external knowledge to answer. The answers are provided in free-text direct answer form. Both
666 A-OKVQA and OK-VQA sample images from the COCO dataset, with no overlapping.

667 **e-SNLI-VE** [17] dataset is an extended version of SNLI-VE dataset [103], which contains about
668 190k question pairs and human-annotated natural language explanations for the ground-truth labels.
669 The text premise provides a statement about the contents of the image. The task is to determine
670 whether the statement is true or false based on the image content.

671 **Visual Spatial Reasoning** [48] (VSR) consists of 65 spatial relations (*e.g.*, under, in front of, facing,
672 *etc.*) of instances in images. VSR has more than 10k question pairs, associated with 6940 images
673 from MS COCO [47].

674 A.2 Finetuning Details

675 We adopt pretrained BLIP [44]² and OFA [95]³ as VLMs, and freeze their parameters without
676 updating. The finetuning only happens on the language model part. The training set of each dataset is
677 used for finetuning. We use the whole training set unless otherwise specified in low-data finetuning
678 discussion.

679 We use an AdaFactor optimizer [85] at the learning rate of 1e-4 for all Cola-FT experiments. The
680 batch size is by default set to 16, though we find Cola-FT insensitive to batch size. We finetune and
681 evaluate the models on NVIDIA V100 or A100 GPUs. The finetuning ranges from 1 hour to about
682 15 hours, varying by the dataset.

683 Following the common experiment protocols, we employ a teacher forcing and greedy decoding
684 strategy for finetuning.

685 A.3 Evaluation Details

686 As specified, we use the validation or test set multiple choice accuracy as the evaluation metric. In
687 A-OKVQA, we report `val/test` accuracy, and `val` accuracy in e-SNLI-VE, `test` (zero-shot split)
688 accuracy in VSR. For simplicity and consistency, we evaluate ablation experiments on A-OKVQA

²BLIP: <https://github.com/salesforce/BLIP>

³OFA: <https://huggingface.co/OFA-Sys>

689 validation set. Following the common experiment protocols [27, 67], we report the single run results
690 for performance comparison.

691 The exemplars at the inference of Cola-Zero are randomly sampled from the training set, i.e.,
692 supposedly help the LM learn the input data distribution and output format but do not leak relevant
693 information to the evaluation question.

694 A.4 A-OKVQA Direct Answer Results

695 In addition to MC accuracy, we present the direct answer (DA) accuracy of models on the A-OKVQA
696 validation set in Tables 6 and 7.

	FLAN-T5-Small	FLAN-T5-Base	FLAN-T5-XL	FLAN-T5-XXL
Cola-FT	56.5	60.6	64.1	65.4
Cola-Zero (2-shot)	30.3	34.6	57.6	61.0
Cola-Zero (0-shot)	28.6	36.0	55.0	59.3

Table 6: A-OKVQA validation set DA performance. Extension of Figure 5.

	1-shot	2-shot	3-shot	4-shot
Cola-Zero	60.2	61.0	60.7	59.2

Table 7: Cola-Zero in-context few-shot learning DA performance on A-OKVQA validation set. Extension of Figure 6.

697 A.5 Qualitative Examples

698 In this section, we provide more qualitative examples on A-OKVQA (Figure 8), e-SNLI-VE (Fig-
699 ure 9), and VSR (Figure 10) datasets.

700 Due to the large span of the three figures, for better visibility, we put the detailed description directly
701 in each figure’s caption part. We illustrate how Cola-FT and Cola-Zero process the VLMs answers in
702 each example. Overall, in these examples, we can observe that even if BLIP and OFA provide wrong
703 answers, Cola can still present the correct answer based on the captions provided by OFA and BLIP,
704 as well as the choice set. This may illustrate how Cola amazingly accomplishes visual reasoning
705 tasks via coordinating BLIP and OFA.

706 A.6 Failure Cases

707 In Figure 11, we provide a few failed cases to analyze the specific behavior of Cola.

708 The leftmost example’s correct answer is *kayaking*, but there are no hints from OFA and BLIP’s
709 answers and captions. Therefore Cola-Zero incorrectly provides the answer *OFA* without sufficient
710 information as hints, while surprisingly Cola-FT answered correctly from OFA’s *boating* answer.

711 The left example again has insufficient information from captions. While BLIP answers *no* and OFA
712 answers *yes*, Cola-FT chooses to answer *maybe*, which looks natural but unfortunately picks the
713 wrong choice.

714 The right example’s captions contain enough information this time. But both Cola-FT and Cola-
715 Zero are misled by BLIP’s wrong answer *no parking*.

716 The rightmost example also has insufficient information from captions. In this situation, Cola has no
717 choice but to believe either BLIP or OFA’s answer, but it mistakenly prefers BLIP’s wrong answer.

718 A.7 Prompt Templates

719 Across three datasets, the prompt template is roughly the same, with minor differences mainly in
720 the format of the questions and choices. We list the prompt templates adopted in A-OKVQA and
721 e-SNLI-VE/VSR in Table 8 and Table 9, respectively.



Question	Why might people sit here?	The room can be described as what?	In what type of location are they playing with the body board?	What is in front of the monitor?
OFA caption	colorful umbrellas on the riverwalk	living room layout and decor medium size how to decorate a small living room dining combo mant	person, left, and person look at a painting of a great white shark.	a desk with a computer, a lamp, a laptop, and a plant.
BLIP caption	a colorful umbrella umbrella with colorful umbrellas	a dining room table with a glass table and chairs	a man holding a surfboard while another man is standing next to him	a desk with a computer and a lamp
Choices	[<u>to testify</u> , <u>to rest</u> , 'to shop', 'get tattoo']	[<u>tidy</u> , 'messy', 'on fire', 'destroyed']	[<u>room</u> , 'beach', 'park', 'store']	[<u>keyboard</u> , 'phone', 'mouse', 'headphones']
OFA answer	to eat	living room	bedroom	a keyboard
BLIP answer	yes	dining room	beach	monitor
Cola-Zero answer	to rest	tidy	beach	keyboard
Cola-FT answer	to rest	tidy	room	keyboard

Figure 8: **A-OKVQA qualitative examples.** Leftmost: LM doesn't use BLIP and OFA's answers, but may observe from captions to derive the correct final answer. Left: As shown on the left, LM does not follow the wrong answers from OFA and BLIP but gets the correct answers from captions. Right: With both OFA and BLIP answering incorrectly, LM derives the correct one from both VLMs' captions and answers. Rightmost: After assessing the questions, answers, and captions, LM goes with OFA's answer and rewrites it to match the expression in the choices. The correct choices are underlined. Cola-Zero answers are given in zero-shot settings.

VQA Prompt Template

Answer the following multiple-choice question by OFA and BLIP's description and their answers to the visual question. OFA and BLIP are two different vision-language models to provide clues.

OFA's description: <OFA caption>

BLIP's description: <BLIP caption>

Q: <Question>

OFA's answer: <OFA answer>

BLIP's answer: <BLIP answer>

Choices: <Choices to the question>

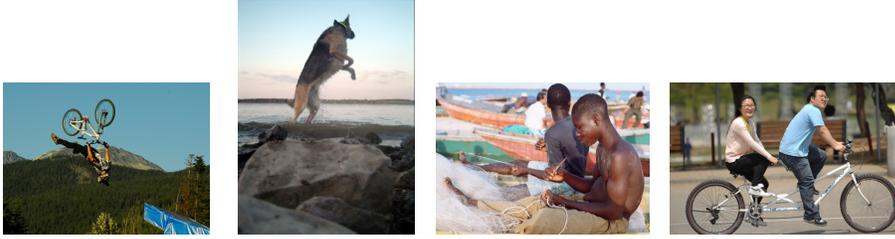
A:

Table 8: **VQA prompt template for the LM, for VQA v2 / OK-VQA / A-OKVQA.** The LM is instructed to coordinate VLMs. Each question set defines *visual context*, *question with choices*, and *plausible answers*.

722 A.8 Parameter-efficient Finetuning

723 To further reduce the computation cost in model adaptation, we explored parameter-efficient finetuning
 724 (PEFT) techniques to reduce finetuning parameter counts. Specifically, we use $(IA)^3$ [49], which
 725 finetunes an overhead of 1 million parameters, equivalent to 0.01% of the full parameters of FLAN-
 726 T5-XXL.

727 Compared to full finetuning, $(IA)^3$ requires more iterations to converge. The performance of a
 728 $(IA)^3$ finetuned FLAN-T5-XXL model is on par with a fully finetuned FLAN-T5-Small (80 million



Question	Does the image describe "A professional daredevil"?	Does the image describe "the dog is a shitz"?	Does this image describe "Two twenty-somethings prepare to catch salmon while other older men catch catfish"?	Does this image describe "A little girl gets hit by a woman riding a bike"?
OFA caption	person doing a flip on a mountain bike	a dog jumping out of the water.	men repairing fishing nets on the beach in zanzibar, tanzania	a man and a woman on a tandem bike
BLIP caption	a man doing a trick on a bike in the air	a dog jumping over rocks in the water	a man sitting on a boat with a fishing net net	a man and woman riding a bicycle in a parking lot
Choices	['yes', ' <u>maybe</u> ', 'no']	['yes', ' <u>maybe</u> ', 'no']	['yes', ' <u>maybe</u> ', 'no']	['yes', 'maybe', ' <u>no</u> ']
OFA answer	yes	no	yes	yes
BLIP answer	yes	no	yes	no
Cola-Zero answer	yes	no	no	no
Cola-FT answer	maybe	maybe	maybe	no

Figure 9: **e-SNLI-VE qualitative examples**. Leftmost: As the connection to *daredevil* is not obvious in BLIP and OFA’s captions, although Cola-Zero is misled, Cola-FT correctly answers *maybe*. Left: Similar to the left example, Cola-FT answer correctly as no obvious connections are seen from the captions to this question. Right: Similar to the left example, the fact of *catch catfish* is not reasonable from the captions, Cola-FT picks the correct answer *maybe*. Rightmost: As *girl gets hit* is not obvious in BLIP and OFA’s captions and answers, Cola-Zero and Cola-FT both follow BLIP to choose the correct answer *no*. The correct choices are underlined. Cola-Zero answers are given in zero-shot settings.

729 parameters) counterpart (Figure 5). Notably, the former is associated with more computation and
 730 memory footprint as a consequence of more parameters in the forward pass.

731 A.9 Extended Ablation Studies

732 **Do caption labels offer useful information to LLM? How would more prompt variations affect**
 733 **the performance of Cola?** We tested Cola-Zero with and without caption labels on A-OKVQA
 734 validation set, observing a slight decrease in performance when without them (70.39% w/t vs. 69.97%
 735 w/o). More ablative experiments showed that removing the VLM’s answer labels led to a substantial
 736 drop in performance (70.39% w/t vs. 67.62% w/o). Removing the model characteristic descriptions
 737 also led to a decrease (70.39% w/t vs. 68.37% w/o).

738 **Do longer image captions improve reasoning performance?** On A-OKVQA validation set, we
 739 tested longer image descriptions (>50 tokens) but found no gain compared to Cola or single VLMs.
 740 Longer captions decreased FLAN-T5+OFA’s accuracy by 0.61% and FLAN-T5 with BLIP by 0.69%
 741 on the A-OKVQA validation set. Cola (captions <30 tokens) reached 77.73%, outperforming
 742 individual VLMs. Longer captions lacked meaningful visual context, possibly due to short text and
 743 image pairs in their training datasets. This experiment reaffirms Cola’s effectiveness in aggregating
 744 individual VLM functionalities.

745 B Extended Related Works

746 B.1 Finetuning Large Language Models

747 Large language models [7, 64, 5] pretrained on massive amounts of unstructured data have gradually
 748 demonstrated great performance by finetuning on additional task-specific instances. Finetuning

				
Question	Does this image describe "The truck contains the elephant" ?	Does this image describe "The bed is under the handbag" ?	Does this image describe "The couch is behind the hot dog" ?	Does this image describe "The bowl contains the banana" ?
OFA caption	an elephant being transported on a truck in sri lanka	a black and white tuxedo cat with a white nose, yellow eyes, and white	person enjoying a meal by the fire	bananas and mangoes in a bowl
BLIP caption	a truck with a large elephant in the back of it	a black cat laying on a bed with a pillow	a man sitting on a couch with a plate of food	a bowl of fruit is shown in this bowl
Choices	[' <u>yes</u> ', 'no']	['yes', ' <u>no</u> ']	['yes', ' <u>no</u> ']	[' <u>yes</u> ', 'no']
OFA answer	yes	no	yes	yes
BLIP answer	no	no	yes	no
Cola-Zero answer	no	no	no	yes
Cola-FT answer	yes	no	no	yes

Figure 10: **VSR qualitative examples.** Leftmost: As OFA caption mentioned *elephant being transported* and OFA provides the correct answer, Cola-FT follows OFA's choice. Left: As OFA and BLIP provide the same answer, Cola-Zero and Cola-FT follow the choice. Right: As the captions do not provide obvious information, even BLIP and OFA provide the same answer, Cola-Zero and Cola-FT are not misled to the wrong choice. Rightmost: As the captions provide strong clue *bananas in a bowl*, although BLIP's answer is incorrect, Cola-Zero and Cola-FT still choose the correct answer. The correct choices are underlined. Cola-Zero answers are given in zero-shot settings.

				
Question	What are the people doing in the water?	Does the image describe "The man is making a vase"?	What kind of zone is this bike parked in?	Does this image describe "The motorcycle is beside the truck" ?
OFA caption	black and white photo of a man on a bike looking at a canoe in the river	person on the potter's wheel	a city made by people bucharest	men walking past a truck in kabul, afghanistan.
BLIP caption	a man and woman on a bike in a park	a man is sitting on a chair and is using a wheel	a bicycle parked next to a pedestrian crossing sign	a man walking down the street in a city
Choices	['surfing', 'fishing', ' <u>kayaking</u> ', 'swimming']	['yes', 'maybe', ' <u>no</u> ']	['temporary', ' <u>pedestrian</u> ', 'no parking', 'handicap']	[' <u>yes</u> ', 'no']
OFA answer	boating	yes	pedestrian	yes
BLIP answer	swimming	no	no parking	no
Cola-Zero answer	OFA	no	no parking	no
Cola-FT answer	kayaking	maybe	no parking	no

Figure 11: **Failed cases.** The correct choices are underlined. Cola-Zero answers are given in zero-shot settings.

749 a large language model can be considerably more sample efficient than re-training from scratch,
 750 although acceptable performance may still require a considerable quantity of data [88]. Recent

e-SNLI-VE / VSR Prompt Template

Answer the following multiple-choice question by OFA and BLIP’s description and their answers to the visual question. OFA and BLIP are two different vision-language models to provide clues.

OFA’s description: <OFA caption>

BLIP’s description: <BLIP caption>

Q: does the image describe <hypothesis> ?

OFA’s answer: <OFA answer>

BLIP’s answer: <BLIP answer>

e-SNLI-VE Choices: [yes, no, maybe]

VSR Choices: [yes, no]

A:

Table 9: **e-SNLI-VE/VSR prompt template for the LM.** The LM is instructed to coordinate VLMs. Each question set defines *visual context*, *hypothesis*, and *plausible answers*.

	Accuracy	# Finetuning Params
Finetuning	77.73	11B (100%)
PEFT, (IA) ³	63.76	1M (0.01%)

Table 10: (IA)³ [49] **parameter-efficient tuning (PEFT) performance.** We finetune a FLAN-T5-XXL model on the A-OKVQA training set and evaluate it on the A-OKVQA validation set.

751 works have finetuned task-specific models that demonstrate amazing capabilities in many real-world
752 applications, such as Copilot for program synthesis [9].

753 B.2 Instruction-based Learning

754 Recent advances in the capabilities of language models have piqued researchers’ curiosity in the field
755 of instruction-based learning [22, 58, 80, 20]. The core of instruction-based learning is to explore
756 the knowledge of the language model itself. In contrast to prompt learning to stimulate the language
757 model’s ability to complete blanks, instruction tuning more focuses on activating the language model’s
758 comprehension by giving obvious instructions to models and expecting correct feedback. Earlier
759 work [61] finetune BART [40] using instructions and few-shot exemplars in question answering,
760 text classification, and text modification. Their findings suggest that few-shot instruction tuning
761 improves performance on unseen tasks. [60] finetunes GPT-2 Large and also observes that few-shot
762 exemplar instruction tuning could improve performance. [78] finetunes T5-11B with more diverse
763 instruction templates and observe similar improvements in zero-shot learning. More recent work [99]
764 performs large-scale experiments with a 137B FLAN-T5 model and instruction-tune it on over 60
765 datasets verbalized via instruction templates. They observe FLAN-T5 substantially improves over
766 zero-shot GPT-3 (175B) on 20 of 25 evaluation datasets. OpenAI also releases InstructGPT [64]
767 based on GPT-3 [7], it makes use of human annotations to steer desired model behavior through both
768 instruction tuning and reinforcement learning of human feedback. They discover that InstructGPT is
769 favored by humans over unmodified GPT-3.

770 B.3 Visual Reasoning

771 Beyond the uni-modal reasoning tasks such as question answering (QA) [93, 35, 11, 73, 72, 19,
772 68, 14, 89, 21, 116, 109, 6], visual reasoning requires models to not only understand and interpret
773 visual information but also to apply high-level cognition to derive rational solutions [34, 29, 3,
774 53, 54, 77, 110]. Several tasks have been introduced to address visual reasoning, such as visual
775 question answering (VQA) [1], in which models are expected to provide answers to questions
776 related to an image and visual entailment (VE) [103], where the model is required to determine

777 the similarity or relationship between a given image and a description. Classic visual reasoning
778 methods have employed an image encoder and a text encoder, along with a reasoning block that
779 utilizes attention mechanisms [111, 65, 112, 96], neuro-symbolic methods [107, 55, 106], or external
780 knowledge [56, 24, 10] to perform reasoning.

781 Recent progress in large pre-trained models has led to the development of language models (LMs)
782 that possess exceptional commonsense reasoning capabilities [71, 13, 12, 70]. These models can
783 potentially replace the reasoning block in visual reasoning tasks, and LMs' lack of perception can
784 be compensated by incorporating multiple vision-language models (VLMs) trained on different
785 domains [69, 95, 44]. For example, PICa [105] converts the image into captions that GPT-3 [7]
786 can understand, and adapts GPT-3 to solve the VQA task in a few-shot manner by providing a few
787 in-context VQA examples. However, there is still a lack of research on how to harness the collective
788 power of these complementary VLMs for visual reasoning tasks.

789 **B.4 Model Ensembling**

790 Model ensembling is a powerful machine learning technique that combines the predictions of multiple
791 models to improve the overall performance of a given task [16]. Classic model ensembling methods
792 include simple averaging, weighting the predictions based on model performance, and stacking the
793 models. By combining the predictions of multiple models, ensembling can reduce the variance and
794 bias of the final predictions, resulting in a more robust and accurate model [76]. Ensemble methods
795 have been shown to perform well in a wide range of tasks, including image classification, natural
796 language processing, and time series forecasting. However, when it turns to multimodal tasks such as
797 visual reasoning, a simple combination is not applicable to heterogeneous models as their inputs and
798 outputs vary.

799 The Mixture-of-Experts (MoE) [84, 75, 115, 39, 41] can be conceptualized as a model ensemble
800 strategy implemented at the level of network architecture. MoE-based multi-modal models [28] excel
801 in leveraging the specific strengths of each expert, thereby delivering the performance that often
802 outstrips that of any individual expert. In these networks, the credibility of each expert's output is
803 dynamically weighted, facilitating a comprehensive and nuanced response to multimodal tasks.

804 However, even within this sophisticated framework, challenges can arise, particularly when managing
805 heterogeneous pre-trained multimodal models. To address this problem, an innovative approach
806 known as Socratic Models (SMs) [113] has been proposed. SMs employ prompt engineering to guide
807 these diverse models through multimodal discussions, effectively combining their varied knowledge.
808 This method promotes a more harmonious and effective integration of different models, enhancing
809 the ensemble's ability to handle complex tasks.

810 With a similar goal, [46] proposes a closed-loop iterative consensus optimization method to utilize
811 the strengths of individual models. However, previous methods do not fully explore the potential of a
812 centralized solution or adapt to the separate functionalities of different models, particularly in the
813 visual reasoning scenario. Recent studies, such as CICERO [59], have shown that language models
814 possess strong capabilities in coordinating multiple agents, which inspires us to reorganize pre-trained
815 multimodal models with a focus on the language models.

816 **Broader Impact**

817 This study inherits ethical risks of biases from pretrained VLMs and LMs, depending on their training
818 data. We suggest the users consider the possible biases in reasoning and prompt the model to interpret
819 its predictions in natural languages when necessary.