

analyses such as VALUEFLOW with safeguards, transparency mechanisms, and human oversight when applied beyond controlled research settings.

## References

- Shuaihang Chen, Yuanxing Liu, Wei Han, Weinan Zhang, and Ting Liu. 2025. **A survey on llm-based multi-agent system: Recent advances and new frontiers in application.** *Preprint*, arXiv:2412.17481.
- Christian Haerpfer, Ronald Inglehart, Alejandro Moreno, Christian Welzel, K Kizilova, Jaime Diez-Medrano, Marta Lagos, Pippa Norris, Eduard Ponarin, Bi Puranen, and 1 others. 2020. World values survey: Round seven-country-pooled datafile. madrid, spain & vienna, austria: Jd systems institute & wvs secretariat. *Version: http://www.worldvaluessurvey.org/WVSDocumentationWV7.jsp*.
- Mengkang Hu, Pu Zhao, Can Xu, Qingfeng Sun, Jianguang Lou, Qingwei Lin, Ping Luo, and Saravanan Rajmohan. 2025. **Agentgen: Enhancing planning abilities for large language model based agent via environment and task generation.** *Preprint*, arXiv:2408.00764.
- Yoichi Ishibashi and Yoshimasa Nishimura. 2024. **Self-organized agents: A llm multi-agent framework toward ultra large-scale code generation and optimization.** *Preprint*, arXiv:2404.02183.
- Han Jiang, Xiaoyuan Yi, Zhihua Wei, Ziang Xiao, Shu Wang, and Xing Xie. 2024a. Raising the bar: Investigating the values of large language models via generative evolving testing. *arXiv preprint arXiv:2406.14230*.
- Liwei Jiang, Sydney Levine, and Yejin Choi. 2024b. **Can language models reason about individualistic human values and preferences?** In *Pluralistic Alignment Workshop at NeurIPS 2024*.
- Liwei Jiang, Taylor Sorensen, Sydney Levine, and Yejin Choi. 2025. **Can language models reason about individualistic human values and preferences?** *Preprint*, arXiv:2410.03868.
- Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. 2023. **Dspy: Compiling declarative language model calls into self-improving pipelines.** *Preprint*, arXiv:2310.03714.
- Hannah Rose Kirk, Bertie Vidgen, Paul Röttger, and Scott A Hale. 2024. The benefits, risks and bounds of personalizing the alignment of large language models to individuals. *Nature Machine Intelligence*, pages 1–10.
- Yunxuan Li, Yibing Du, Jiageng Zhang, Le Hou, Peter Grabowski, Yeqing Li, and Eugene Ie. 2024. **Improving multi-agent debate with sparse communication topology.** *Preprint*, arXiv:2406.11776.
- Siyang Liu, Trisha Maturi, Bowen Yi, Siqi Shen, and Rada Mihalcea. 2024. The generation gap: Exploring age bias in the value systems of large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 19617–19634.
- Yuanyi Ren, Haoran Ye, Hanjun Fang, Xin Zhang, and Guojie Song. 2024. **Valuebench: Towards comprehensively evaluating value orientations and understanding of large language models.** *Preprint*, arXiv:2406.04214.
- Shalom H Schwartz. 1992. Universals in the content and structure of values: Theoretical advances and empirical tests in 20 countries. In *Advances in experimental social psychology*, volume 25, pages 1–65. Elsevier.
- Shalom H Schwartz. 1994. Are there universal aspects in the structure and contents of human values? *Journal of social issues*, 50(4):19–45.
- Shalom H Schwartz. 2012. An overview of the schwartz theory of basic values. *Online readings in Psychology and Culture*, 2(1):11.
- Shalom H Schwartz, Jan Cieciuch, Michele Vecchione, Eldad Davidov, Ronald Fischer, Constanze Beierlein, Alice Ramos, Markku Verkasalo, Jan-Erik Lönnqvist, Kursad Demirütke, and 1 others. 2012. Refining the theory of basic individual values. *Journal of personality and social psychology*, 103(4):663.
- Hua Shen, Nicholas Clark, and Tanushree Mitra. 2025. **Mind the value-action gap: Do llms act in alignment with their values?** *Preprint*, arXiv:2501.15463.
- Hua Shen, Chieh-Yang Huang, Tongshuang Wu, and Ting-Hao Kenneth Huang. 2023. Convxai: Delivering heterogeneous ai explanations via conversations to support human-ai scientific writing. In *Companion Publication of the 2023 Conference on Computer Supported Cooperative Work and Social Computing*, pages 384–387.
- Hua Shen, Tiffany Kneare, Reshma Ghosh, Kenan Alkiek, Kundan Krishna, Yachuan Liu, Ziqiao Ma, Savvas Petridis, Yi-Hao Peng, Li Qiwei, and 1 others. 2024a. Towards bidirectional human-ai alignment: A systematic review for clarifications, framework, and future directions. *arXiv preprint arXiv:2406.09264*.
- Hua Shen, Tiffany Kneare, Reshma Ghosh, Yu-Ju Yang, Tanushree Mitra, and Yun Huang. 2024b. **Valuecompass: A framework of fundamental values for human-ai alignment.** *arXiv preprint arXiv:2409.09586*.
- Hua Shen, Yuguang Yang, Guoli Sun, Ryan Langman, Eunjung Han, Jasha Droppo, and Andreas Stolcke. 2022. Improving fairness in speaker verification via group-adapted fusion network. In *ICASSP 2022–2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7077–7081. IEEE.

Taylor Sorensen, Jared Moore, Jillian Fisher, Mitchell Gordon, Niloofar Mireshghallah, Christopher Michael Rytting, Andre Ye, Liwei Jiang, Ximing Lu, Nouha Dziri, and 1 others. 2024. A roadmap to pluralistic alignment. *arXiv:2402.05070*.

Qiaosi Wang, Michael Madaio, Shaun Kane, Shivani Kapania, Michael Terry, and Lauren Wilcox. 2023. Designing responsible ai: Adaptations of ux practice to meet responsible ai challenges. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pages 1–16.

Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao Ge, Furu Wei, and Heng Ji. 2024. Unleashing the emergent cognitive synergy in large language models: A task-solving agent through multi-persona self-collaboration. *Preprint*, arXiv:2307.05300.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models. *Preprint*, arXiv:2201.11903.

Zihao Yi, Jiarui Ouyang, Zhe Xu, Yuwen Liu, Tianhao Liao, Haohao Luo, and Ying Shen. 2025. A survey on recent advances in llm-based multi-turn dialogue systems. *Preprint*, arXiv:2402.18013.

Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, and Zhi Jin. 2024. Codeagent: Enhancing code generation with tool-integrated agent systems for real-world repo-level coding challenges. *Preprint*, arXiv:2401.07339.

Xinyang Zhang, Ningfei Wang, Hua Shen, Shouling Ji, Xiapu Luo, and Ting Wang. 2020. Interpretable deep learning under fire. In *29th {USENIX} security symposium ({USENIX} security 20)*.

Han Zhou, Xingchen Wan, Ruoxi Sun, Hamid Palangi, Shariq Iqbal, Ivan Vulić, Anna Korhonen, and Serkan Ö. Arik. 2025. Multi-agent design: Optimizing agents with better prompts and topologies. *Preprint*, arXiv:2502.02533.

Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbulin, and Jürgen Schmidhuber. 2024. Language agents as optimizable graphs. *Preprint*, arXiv:2402.16823.