

- [74] Y.-F. Ma, L. Lu, H.-J. Zhang, and M. Li, “A user attention model for video summarization,” in *Proceedings of the tenth ACM international conference on Multimedia*, 2002, pp. 533–542.
- [75] Y.-F. Ma, X.-S. Hua, L. Lu, and H.-J. Zhang, “A generic framework of user attention model and its application in video summarization,” *IEEE transactions on multimedia*, vol. 7, no. 5, pp. 907–919, 2005.
- [76] R. Krishna, K. Hata, F. Ren, L. Fei-Fei, and J. Carlos Niebles, “Dense-captioning events in videos,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 706–715.
- [77] J. Wang, W. Jiang, L. Ma, W. Liu, and Y. Xu, “Bidirectional attentive fusion with context gating for dense video captioning,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7190–7198.
- [78] C. Sun, A. Myers, C. Vondrick, K. Murphy, and C. Schmid, “Videobert: A joint model for video and language representation learning,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 7464–7473.
- [79] D. Zhao, Z. Xing, C. Chen, X. Xu, L. Zhu, G. Li, and J. Wang, “Seenomaly: Vision-based linting of gui animation effects against design-don’t guidelines,” in *2020 IEEE/ACM 42nd International Conference on Software Engineering (ICSE)*. IEEE, 2020, pp. 1286–1297.
- [80] S. Feng, M. Xie, and C. Chen, “Efficiency matters: Speeding up automated testing with gui rendering inference,” *arXiv preprint arXiv:2212.05203*, 2022.
- [81] L. Gomez, I. Neamtiu, T. Azim, and T. Millstein, “Reran: Timing-and touch-sensitive record and replay for android,” in *2013 35th International Conference on Software Engineering (ICSE)*. IEEE, 2013, pp. 72–81.
- [82] P. Krieter and A. Breiter, “Analyzing mobile application usage: generating log files from mobile screen recordings,” in *Proceedings of the 20th international conference on human-computer interaction with mobile devices and services*, 2018, pp. 1–10.
- [83] D. Nurmuradov and R. Bryce, “Caret-hm: recording and replaying android user sessions with heat map generation using ui state clustering,” in *Proceedings of the 26th ACM SIGSOFT International Symposium on Software Testing and Analysis*, 2017, pp. 400–403.
- [84] L. Zhu and Y. Yang, “Actbert: Learning global-local video-text representations,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 8746–8755.
- [85] Z. Zhang, Y. Shi, C. Yuan, B. Li, P. Wang, W. Hu, and Z.-J. Zha, “Object relational graph with teacher-recommended learning for video captioning,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 13 278–13 288.