# Avisha Das

das.avisha@mayo.edu ⋄ `dasavisha.github.io` ⋄ Google Scholar Profile

## EDUCATION

**Ph.D. in Computer Science** *2014 – 2020*
*University of Houston, Houston, TX*
Dissertation Title: Proactive Defense through Automated Generation of Targeted Attacks

**B.Tech. in Electronics and Communication Engineering** *2010 – 2014*
*West Bengal University of Technology, Kolkata, India*

## EXPERIENCE

**Research Fellow** *November 2023 - Present*
*Arizona Advanced AI & Innovation (A3I) Hub*
*Mayo Clinic, Phoenix, AZ*

**Postdoctoral Research Fellow** *April 2021 - November 2023*
*School of Biomedical Informatics*
*University of Texas Health Science Center at Houston (UTHealth), Houston, TX*

**Data Science-NLP Intern** *May 2019 – August 2019*
*Occidental (Oxy) Petroleum Corporation*
*The Woodlands, TX*

**Summer Research Intern** *June 2018 – August 2018*
*Production Solutions Team, Halliburton Energy Services*
*Houston, TX*

**Data Science Intern** *June 2017 – August 2017*
*2H Offshore Inc.*
*Houston, TX*

**Graduate Research and Teaching Assistant** *August 2014 - December 2020*
*Department of Computer Science,*
*University of Houston, Houston, TX*

## RESEARCH INTERESTS

**Natural Language Processing and Generative Modeling:** Language generation and understanding techniques applied to intelligent dialogue systems and conversational agents, context modeling, and multi-turn dialogue management.

**Security and Large Language Modeling:** Exploring the implications of LLMs in the cyber landscape with a focus on adversarial learning and proactive defensive measures against social engineering threats.

**Knowledge Mining and Retrieval:** Ranking and relevance modeling, and structured semantic retrieval by leveraging semantic annotation, entity recognition, and integration with knowledge graphs.

**Application Areas:** Intelligent automated cognitive therapy, biomedical literature mining, Social engineering threat detection and prevention, adversarial NLP.

## PUBLICATIONS

### Journal Papers

1. Li, Z., Wei, Q., Huang, L.C., Li, J., Hu, Y., Chuang, Y.S., He, J., **Das, A.**, Keloth VK, Yang Y, Diala CS. (2024). **Ensemble pretrained language models to extract biomedical knowledge from literature.** *Journal of the American Medical Informatics Association.*

2. Yang, Y., Zuo, X., **Das, A.**, Xu, H., Zheng, W. Jim (2024). **Representation Learning of Biological Concepts: A Systematic Review**. *Current Bioinformatics.*

3. **Das, A.** & Verma, R. (2020). **Can Machines Tell Stories? A Comprehensive Comparison of Pre-Trained and Fine-Tuned Deep Neural Language Models**. *IEEE Access.*

4. El Aassal, A., Baki, S., **Das, A.**, & Verma, R. (2020). **An In-Depth Benchmarking and Evaluation of Phishing Detection Research for Security Needs.** *IEEE Access.*

5. **Das, A.**, Baki, S., El Aassal, A., Verma, R., & Dunbar, A. (2019). **SoK: A Comprehensive Reexamination of Phishing Research from the Security Perspective.** *IEEE Communications Surveys & Tutorials.*

6. Karimi, S., Moraes, L., **Das, A.**, Shakery, A., & Verma, R. (2018). **Citance-based retrieval and summarization using IR and machine learning.** *Scientometrics.*

### Conference and Workshop Papers

7. **Das, A.**, Tariq, A., Batalini, F., Dhara, B. and Banerjee, I. (2024). **Exposing Vulnerabilities in Clinical LLMs Through Data Poisoning Attacks: Case Study in Breast Cancer** *AMIA Annual Symposium.*

8. **Das, A.**, Li, Z., Wei, Q., Li, J., Huang, L.C., Hu, Y., Li, R., Zheng, W. and Xu, H. (2023). **Extracting Drug-Protein Relation from Literature using Ensembles of Biomedical Transformers.** *The 19th World Ccongress on Medical and Health Informatics (MedInfo).*

9. **Das, A.**, Selek, S., Warner, A., Zuo, X., Hu, Y., Keloth, V., Li, J., Zheng, W., Xu, H. (2022). **Conversational Bots for Psychotherapy: A Study of Generative Transformer Models Using Domain-specific Dialogue.** *BioNLP Workshop (Co-located with Association of Computational Linguistics)*

10. Zeng, V., El Aassal, A., Baki, S., Verma, R., Moraes, L. & **Das, A.** (2020). **Diverse Datasets and a Customizable Benchmarking Framework for Phishing**. *Proceedings of the 5th ACM on International Workshop on Security and Privacy Analytics.*

11. **Das, A.**, & Verma, R. (2019). **Automated email Generation for Targeted Attacks using Natural Language.** *Workshop TA-COS (Co-located with Eleventh International Conference on Language Resources and Evaluation-LREC)*

12. El Aassal, A., Moraes, L., Baki, S., **Das, A.**, & Verma, R. (2018). **Anti-Phishing Pilot at ACM IWSPA 2018: Evaluating Performance with New Metrics for Unbalanced Datasets**. *Proceedings of the 1st Anti-Phishing Shared Task Pilot at 4th ACM IWSPA co-located with 8th ACM Conference on Data and Application Security and Privacy (CODASPY 2018)*

13. Verma, R., & **Das, A.** (2017, March). **What's in a URL: Fast feature extraction and malicious URL detection.** *Proceedings of the 3rd ACM on International Workshop on Security and Privacy Analytics.*

14. **Das, A.**, Li, Z., Wei, Q., Li, J., Huang, L.C., Hu, Y., Li, R., Zheng, W.J. and Xu, H. (2021). **UTHealth@BioCreativeVII: Domain-specific Transformer Models for Drug-Protein Relation Extraction.** *BioCreative VII Workshop.*

15. De Moraes, L. F., **Das, A.**, Karimi, S., & Verma, R. (2018). **University of Houston@ CL-SciSumm 2018.** *BIRNDL@ SIGIR.*

16. Karimi, S., Moraes, L. F., **Das, A.**,[1] & Verma, R. (2017). **University of Houston@ CL-SciSumm 2017: Positional language Models, Structural Correspondence Learning and Textual Entailment.** *BIRNDL@ SIGIR.*

**Posters and Abstracts**

17. Tariq, A., Luo, M., Urooj, A., **Das, A.**, Jeong, J., Trivedi, S., Patel, B. and Banerjee, I. (2024). Domain-specific LLM Development and Evaluation–A Case-study for Prostate Cancer. *AMIA Annual Symposium.*

18. **Das, A.**, Anjum, O., Chen, G., Zheng, W., Li, Rongbin (2024). **Efficient Training Corpus Retrieval for Large Language Model Fine Tuning** *AMIA Informatics Summit Paper.*

19. **Das, A.**, Anjum, O., Zheng, W., Diala, C. (2023). **A Multi-faceted Mining Tool for Knowledge and Data Discovery for Cancer Research.** *International Conference on Intelligent Biology and Medicine (ICIBM).*

20. **Das, A.**. (2019) **AskAna: Retrieval Based Virtual Assistant for Digital Operations and Field Development.** *Rice Data Science Conference.*

21. **Das, A.**, & Verma, R. (2017). **What's in a URL: Fast Feature Extraction and Detection of Malicious URLs.** *Women in CyberSecurity (WiCyS) Conference.*

22. **Das, A.**, & Verma, R. (2016). **Analyzing Phishing URLs.** *Poster at Grace Hopper Conference for Celebration of Women.*

23. **Das, A.**, & Verma, R. (2016). **Are Legit and Phishing URLs similar? Hell No! – Lexical characterization and Analysis of URLs.** *Women in CyberSecurity (WiCyS) Conference.*

24. **Das, A.**, & Verma, R. (2016). **Studying Phishing URLs the NLP way.** *Computing Research Association (CRA-W) Grad Cohort Workshop.*

**Book Chapters**

25. Tariq, A., Luo, M., Urooj, A., **Das, A.**, Jeong, J., Trivedi, S., Abdul-Muhsin, H., Ghaffar, U., Yu, N., Patel, B., Banerjee, I. (2024). **Development Of LLM For Prostate Cancer - The Need for Domain-Tailored Training.** *National Cancer Institute.*

**Preprints/Under Review**

26. **Das, A.**, Tariq, A., Batalini, F., Dhara, B., Banerjee, I. (2024). **Framework for Exposing Vulnerabilities of Clinical Large Language Model: A Case Study in Breast Cancer.** *Under Review.*

27. **Das, A.**, Anjum, O., Chen, G., Zheng, W. Jim (2023). **Efficient Training Corpus Retrieval for Large Language Model Fine Tuning**. *Under Review.*

28. **Das, A.**, Keloth, V., Selek, S., Xu, H. (2023). **A Methodological Systematic Review of Deep Learning-based Virtual Assistants for Healthcare**. *Under Review.*

29. Tariq, A., Luo, M., Urooj, A., **Das, A.**, Jeong, J., Trivedi, S., Patel, B. and Banerjee, I. (2024). **Domain-specific LLM Development and Evaluation–A Case-study for Prostate Cancer.** *medRxiv preprint.*

30. **Das, A.** and Verma, R. (2020). **Modeling Coherency in Generated Emails by Leveraging Deep Neural Learners**. *ArXiv preprint.*

# INVITED TALKS

1. **Framework for Exposing Vulnerabilities of Clinical LLMs: Breast Cancer.**
   Stanford MedAI Group Exchange Sessions, Stanford University, 2024.

2. **Large language models and their application in Biomedical Domain.**
   DSICCR Tuesday Seminar Series, UTHealth Houston, 2023.

3. **Domain-specific Transformer Models for Drug-Protein Relation Extraction.**
   CPH Seminar in Precision Medicine, UTHealth Houston, 2022.

4. **Leveraging NLP for Mining Biomedical Data: Named Entity Recognition and Content Recommendation.**
   CPRIT-BIG-TCR Undergraduate Summer Internship Seminar, UTHealth Houston, 2022.

5. **Natural Language Understanding and Generation**
   Advanced Natural Language Processing Course, University of Houston, 2022.

## MEDIA COVERAGE

**Automated Email Generation for Targeted Attacks**. AD-Tech, DataSkeptic Podcast, 2022 Oct 31. Link.

## TEACHING EXPERIENCE

**Teaching Assistant,** *University of Houston*
- Artificial Intelligence (COSC 6368) [Summer'20]
- Software Design (COSC 4353/6353) [Spring'20]
- Machine Learning (COSC 6342) [Fall'19]
- Computer Organization and Architecture (COSC 6323) [Fall'18]
- Security Analytics (COSC 4397/COSC 6346) [Spring'18, Spring'19]
- Software Design (COSC 4353/6353) [Fall'17]
- Data Structures and Algorithms (COSC 3320) [Fall'16, Spring'17]

## AWARDS, HONORS AND OTHERS

**Awards and Honors**
1. **CPRIT BIG-TCR Postdoctoral Training Program Fellowship,**[1] 2022-2024.
   Cancer Prevention and Research Institute of Texas, UTHealth Houston.
2. **Second place, Litcoin NLP Challenge,**[2] March 2022.
   National Center for Advancing Translational Sciences (NCAT), UTHealth Houston.
3. **Cullen Graduate Success Fellowship**, Fall 2020.
   UH Alumni Association, University of Houston.
4. **Govt. of India Merit-based Scholarship for Undergraduate Education**, 2010 -2014.
   Ministry of Human Resources-India (MHRD), India.

**Travel Grants**
1. Annual Meeting of the Association for Computational Linguistics (ACL), 2020, 2022
2. Grace Hopper Conference for Women in Computing (GHC), 2015, 2016, 2018
3. International Workshop on Security and Privacy Analytics (IWSPA), 2017, 2018
4. Empirical Methods in Natural Language Processing Conference (EMNLP), 2016
5. Women in CyberSecurity Conference (WiCyS), 2016, 2017
6. Computing Research Association for Women (CRA-W), 2015

**Other**
1. First Place (Winner), CodeRED Discovery (2018), University of Houston
2. Third Place, CodeRED Exploration (2017). University of Houston.
3. Winner, Social Track at HackRice 7 (2017), Rice University.

## PROFESSIONAL/ACADEMIC SERVICE

**Journal Club**

· Organizer, MedAI Group Exchange Sessions, Stanford University-Mayo Clinic Arizona.

**Editorial Services**

· Review Editor, Text-mining and Literature-based Discovery, Frontiers in Research Metrics and Analytics Journal.

---

[1]https://www.uth.edu/big-tcr/people/trainees.htm
[2]Part of the UTHealth-SBMI Team (Result)

**Reviewing Services**

· **Journals**
  1. Artificial Intelligence in Medicine Journal (IF: 7.011)
  2. Journal of Biomedical Informatics (JBI) (IF: 8.0)
  3. Computers & Security Journal (IF: 5.105)
  4. Journal of Information Security and Applications (IF: 4.96)
  5. IEEE Open Access Journal (IF: 3.475)
  6. Neural Computing and Applications (NCAA) (IF: 5.102)
  7. PLOS Digital Health (IF:4.01)
· **Conferences**
  1. Association for the Advancement of Artificial Intelligence (AAAI), 2024
  2. Empirical Methods in Natural Language Processing (EMNLP), 2021, 2022, 2023
  3. Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (AACL), 2021, 2022, 2023
  4. International Joint Conference on Natural Language Processing (IJCNLP), 2022, 2023
  5. International Conference on Bioinformatics and Biomedicine (BIBM), 2022
  6. Annual Meeting of the Association for Computational Linguistics (ACL), 2019, 2018
  7. ACM International Workshop on Security and Privacy Analytics (Co-located with CODASPY), 2018, 2019

**Program and Organizing Committee**

· Program committee member for Workshop on Multimodal4Health 2024 (co-located with ICHI)
· Program committee member for Workshop on Natural Language Processing for Bangla 2023 (co-located with EMNLP)
· Program committee member for EMNLP 2022 (Tracks include Language Modeling and Analysis of Language Models, Natural Language Generation, and Summarization tracks)
· Program committee member for AACL-IJCNLP 2022, AACL-IJCNLP 2023
· Chair of Organizing Committee for the First Security and Privacy Analytics Anti-Phishing Shared Task 2018 (co-located with CODASPY 2018)

**Mentoring**

· **Mentor**, Machine Learning for Health (ML4H) Workshop (Co-located with NeurIPS 2022).
· **Graduate**

  1. Rongbin Li (Ph.D. candidate), UTHealth, Houston.

  2. Ayman El Aassal (Ph.D. candidate), University of Houston, Houston.
· **Undergraduate**

  1. Boddhisattwa Dhara, BITS-Pilani (Hyderabad Campus), India.

  2. Gal Egozi, University of Houston, Houston.

## ACADEMIC REFERENCES

Available on Request