

Lokesh Madasu

+91-9133832696 Website maadasulokesh@gmail.com GitHub LinkedIn

Education

International Institute of Information Technology <i>Master of Science by Research in Computer Science - CGPA: 8.83</i>	Hyderabad, India 2021 - ongoing
Rajiv Gandhi University of Knowledge Technologies <i>Bachelor of Technology in Computer Science and Engineering - CGPA: 8.43</i>	Nuzvid, India 2015 - 2019

Experience

IIIT Hyderabad <i>Research Assistant under <u>Dr. Manish Shrivastava</u> - NLP Lab</i> <ul style="list-style-type: none">Working on text generation problems, creating high-quality datasets using crowdsourcing, with a primary focus on low-resource Indian languages. Exploring the performance of state-of-the-art pre-trained models/LLMs on Indic language datasets, specifically in tasks such as summarization and headline generation.	Hyderabad, India Dec 2021 - ongoing
IIIT Hyderabad <i>Research Intern under <u>Dr. Manish Shrivastava</u> - NLP Lab</i> <ul style="list-style-type: none">Developed site-specific web scrapers to collect high-quality data from a variety of news websites in 8 Indian languages, resulting in 4.8 million news article-headline pairs for headline generation task.Created the first largest human-annotated Abstractive Summarization dataset for Telugu, comprising 20329 article-summary pairs.	Hyderabad, India May 2020 - Dec 2021
ChatGen.ai <i>Data Engineer Intern</i> <ul style="list-style-type: none">Worked on data gathering, cleaning, and extracting useful information from diverse web sources.Contributed to the development of conversational chatbots.	Mumbai, India Oct 2017 - Dec 2017

Publications

- Lokesh Madasu, Gopichand Kanumolu, Nirmal Surange, Manish Shrivastava. *Mukhyansh: A Headline Generation Dataset for Indic Languages*. 37th Pacific Asia Conference on Language, Information and Computation, PACLIC-2023 (Accepted as a long paper, to be published).
- Gopichand Kanumolu, Lokesh Madasu, Pavan Baswani, Ananya Mukherjee, Manish Shrivastava. *Unsupervised Approach to Evaluate Sentence-Level Fluency: Do We Really Need Reference?*. South East Asian Language Processing Workshop, AACL-2023 (**Won Best Paper Award**)
- Gopichand Kanumolu, Lokesh Madasu, Nirmal Surange, Manish Shrivastava. *TeClass: A Human-Annotated Relevance-based Headline Classification and Generation Dataset for Telugu*. (Under review).

Research Projects

A Semantic Textual Relatedness Benchmark <ul style="list-style-type: none">Collaborating with researchers from NRC-Canada, Cardiff University to develop a benchmark dataset for the Semantic Relatedness task, focusing on the Telugu, Hindi, and Marathi languages.The project involved creating 4100 sentence pairs, each demonstrating varying degrees of relatedness.It is part of the ongoing shared task in SemEval-2024, where participants are required to develop robust baselines for the dataset.	Aug 2023 - Current
News Headline Generation for 8 Indic languages <ul style="list-style-type: none">Created an extensive multilingual dataset comprising 3.4 million article-headline pairs for headline generation task.Experimented with state-of-the-art pre-trained models and achieved an average R-L score of 31.43 across all eight languages.	May 2022 - Aug 2023
Relevance-based News Headline Classification <ul style="list-style-type: none">Created the largest Telugu news headline classification dataset, comprising 26,178 article-headline pairs.Given a news article and headline as input, the model classifies the headline into three classes based on its relevance to the article.	Feb 2022 - Oct 2023
Measure Text Fluency <ul style="list-style-type: none">Implemented a language model based metric to measure the fluency in 9 Indian languages, and Sinhala.The results of baseline models show a strong correlation with human judgments.	Feb 2022 - July 2022

Skills

- **Languages:** Python, C
- **Databases:** SQL, MongoDB
- **Tools & Frameworks:** PyTorch, Keras, Tensorflow, NLTK, SpaCY, Scikit-Learn, LLMs, Huggingface Transformers
- **Development:** Git, VS Code

Course Projects

- Abstractive Summarization for Telugu** Python, Keras, Transformers
- Given a news article as input, the model generates a concise summary. Implemented various state-of-the-art transformer models on TeSum dataset and achieved ROUGE-L score of 34.0.
- CodeMixed Machine Translation** Python, PyTorch, Transformers
- Implemented a code-mixed machine translation system that translates English sentences into Hinglish sentences(a combination of words from Hindi and English).
- Clickbait Intensity & Style Analysis** Python, Keras, Transformers
- The goal is to predict and reduce clickbait intensity in headlines. Various regression algorithms and pre-trained models were used to predict intensity, along with a paraphrase model for reduction.
- Wikipedia Search Engine** Python, Streamlit
- Developed a scalable search engine for a 90GB English Wikipedia dump, utilizing inverted indexing, searching and relevance ranking. Achieved query response times of under 10 seconds, handling both plain and field queries, while providing accurate and timely results in the form of relevant Wikipedia page titles.

Relevant Courses

- Advanced NLP, Introduction to NLP, Statistical Methods in AI, Information Retrieval & Extraction, Data Analytics

Volunteer Experience

- Mentored over 400+ student interns in the NLP course. Designed and assessed quizzes, assignments, and projects, as well as resolved queries related to assignments, projects, and NLP concepts.
- Volunteered at the Advanced Summer School on NLP (IASNLP-2022) organized by Language Technologies Research Center, IIIT-Hyderabad.

Reference

Dr. Manish Shrivastava (Assistant Professor, MT-NLP Lab, IIIT-Hyderabad)
m.shrivastava@iiit.ac.in