

Sumanth Doddapaneni

PhD Student, IIT Madras

[Website](#) [Github](#) [Google Scholar](#) [Email](#)

Education

Jan 2022 Ongoing	Indian Institute of Technology (IIT), Madras Ph.D, Computer Science & Engineering Advisor: Mitesh M. Khapra	Chennai, India
Aug 2016 May 2020	Indian Institute on Information Technology (IIIT), Sri City B.Tech., Electronics & Communications Engineering	Sri City, India

Experience

Nov 2022 Present	Mila - Quebec AI Institute Visiting Researcher Host: <i>Dr. Rahul Aralikkatte, Dr. Jackie Chi Kit Cheung</i> Exploring pretraining methods to develop better multilingual summarization models.	Remote
Oct 2021 Present	AI4Bharat PhD Researcher Advisors: <i>Dr. Mitesh M. Khapra, Dr. Anoop Kunchukuttan, Dr. Pratyush Kumar</i> Working on Multilingual Language modeling and Machine Translation, with a focus on low-resource Indian languages	Chennai, India
Oct 2020 Sep 2021	Robert Bosch Centre for Data Science and AI Post Baccalaureate Fellow Advisors: <i>Dr. Mitesh M. Khapra, Dr. Anoop Kunchukuttan, Dr. Pratyush Kumar</i> Built SOTA models for Machine Translation (IndicTrans) and Automatic Speech Recognition (IndicWav2Vec) for Indian Languages	Chennai, India

Select Publications

P=Preprints, C=Conference, W=Workshop, J=Journal, *=Equal Contribution

[J.1]	Samanantar: The Largest Publicly Available Parallel Corpora Collection for 11 Indic Languages [🔗] [Code] Gowtham Ramesh*, <u>Sumanth Doddapaneni</u> *, et. al <i>Transactions of the Association for Computational Linguistics</i>	[TACL 2022]
[C.1]	Towards Building ASR Systems For The Next Billion Users [🔗] [Code] Tahir Javed, <u>Sumanth Doddapaneni</u> , Abhigyan Raman, Kaushal Santosh Bhogale, Gowtham Ramesh, Anoop Kunchukuttan, Pratyush Kumar, Mitesh M. Khapra <i>36th AAAI Conference on Artificial Intelligence</i>	[AAAI 2022]
[P.5]	Naamapadam: A Large-Scale Named Entity Annotated Data for Indic Languages [🔗] Arnav Mhaske, Harshit Kedia, <u>Sumanth Doddapaneni</u> , Mitesh M. Khapra, Pratyush Kumar, Rudra Murthy V, Anoop Kunchukuttan <i>Pre-Print</i>	[ArXiv 2022]
[P.4]	IndicXTREME: A Multi-Task Benchmark For Evaluating Indic Languages [🔗] [Code] <u>Sumanth Doddapaneni</u> , Rahul Aralikkatte, Gowtham Ramesh, Shreya Goyal, Mitesh M. Khapra, Anoop Kunchukuttan, Pratyush Kumar <i>Pre-Print</i>	[ArXiv 2022]
[P.3]	Effectiveness of Mining Audio and Text Pairs from Public Data for Improving ASR Systems for Low-Resource Languages [🔗] Kaushal Santosh Bhogale, Abhigyan Raman, Tahir Javed, <u>Sumanth Doddapaneni</u> , Anoop Kunchukuttan, Pratyush Kumar, Mitesh M. Khapra <i>Pre-Print</i>	[ArXiv 2022]
[P.2]	A Survey in Adversarial Defences and Robustness in NLP [🔗] Shreya Goyal, <u>Sumanth Doddapaneni</u> , Mitesh M Khapra, Balaraman Ravindran <i>Pre-Print</i>	[ArXiv 2022]
[P.1]	A Primer on Pretrained Multilingual Language Models [🔗] <u>Sumanth Doddapaneni</u> , Gowtham Ramesh, Mitesh M. Khapra, Anoop Kunchukuttan, Pratyush Kumar <i>Pre-Print</i>	[ArXiv 2021]

Select Research Projects

Multilingual Summarisation

Nov'22 - Present

Dr. Rahul Aralikkatte, Dr. Jackie Chi Kit Cheung

- > Studying the effect of pretraining in multilingual summarisation
- > Pretraining and fine-tuning generative models for abstractive summarisation

Multilingual Language Modeling

Oct'21 - Present

Advisors: Dr. Mitesh M. Khapra, Dr. Anoop Kunchukuttan, Dr. Pratyush Kumar [🔗][Code]

- > Building Large scale monolingual corpora and Language Models for Indian Languages
- > Building strong testsets to evaluate the zero-shot generalisation of the models
- > Studying the various objectives and data regimes that improve zero-shot generalisation of the models
- > Understanding the efficacy of adapters in zero-shot generalisation of the LMs

Speech Recognition

Jun'21 - Feb'22

Advisors: Dr. Mitesh M. Khapra, Dr. Anoop Kunchukuttan, Dr. Pratyush Kumar [Try the Model][Code]

- > Building state-of-the-art models for Automatic Speech Recognition for Indian languages
- > Understanding the various effects of LM on downstream ASR task
- > Accepted as a long paper at AAAI 2022

Machine Translation

Feb'21 - Oct'21

Advisors: Dr. Mitesh M. Khapra, Dr. Anoop Kunchukuttan, Dr. Pratyush Kumar [Try the model][Code]

- > Built state of the art translation model for Indian languages to English and vice versa
- > Created the largest bilingual corpora (50M pairs) across 11 languages for NMT training.
- > Work published in TACL (Volume 10, 2022)

Academic Service

Program Committee MRL @EMNLP'21

Volunteer EACL'21, ICML'21, NeurIPS'21, EMNLP'21, ACL'22, EMNLP'22

Honours and Grants

Google Research Selected to attend the Google Research Week 2023

Naver Labs and Univ. Grenoble Alpes Selected to attend the ALPS Winter School 2022

Robert Bosch Centre Received the Post Baccalaureate Fellowship to work on interdisciplinary AI

TFRC Grant Received TPU Research credits to carry out work on LMs

MSR Travel Grant Received Microsoft Research Travel grant to attend ACL 2022

Volunteering Roles

NLP Reading Group, AI4Bharat Organizer

- > Organizer for NLP Reading Group at AI4Bharat

Volunteer at NLP with Friends Volunteer

- > Help organise talks at NLP with Friends

Invited Talk, Swiggy Data Science Speaker

- > Talk on "Towards Building ASR Systems for the Next Billion Users"

References

- > Dr. Mitesh M. Khapra Associate Professor, IIT Madras, India [🔗]
- > Dr. Anoop Kunchukuttan Senior Applied Researcher, Microsoft AI and Research, India [🔗]
- > Dr. Pratyush Kumar Senior Researcher, Microsoft Research, India [🔗]
- > Dr. Raj Dabre Researcher, NICT, Japan [🔗]
- > Dr. Rahul Aralikkatte Postdoctoral Fellow, MILA & Visiting Researcher, Google [🔗]