

# Nandan Thakur

thakur-nandan.github.io

Email: [nandan.thakur@uwaterloo.ca](mailto:nandan.thakur@uwaterloo.ca)  
[[Google Scholar](#)] [[Twitter](#)] [[GitHub](#)] [[LinkedIn](#)]

---

- EDUCATION     **University of Waterloo**, Waterloo, ON, Canada     September 2021- Present  
Ph.D. Student, David R. Cheriton of Computer Science  
Advisor: Prof. Jimmy Lin
- Birla Institute of Technology & Science, Pilani**, Goa, India     August 2014 - May 2018  
B.E. (Hons.) in Electronics & Instrumentation, Minor in Finance
- PUBLICATIONS   [1] MIRACL: A Multilingual Retrieval Dataset Covering 18 Diverse Languages.  
Xinyu Zhang\*, **Nandan Thakur\***, Odunayo Ogundepo, Ehsan Kamaloo, David Alfonso-Hermelo,  
Xiaoguang Li, Qun Liu, Mehdi Rezagholizadeh, Jimmy Lin. (\* denotes equal contribution)  
Transactions of the Association for Computational Linguistics (TACL), 2023.
- [2] SPRINT: A Unified Toolkit for Evaluating and Demystifying Zero-shot Neural Sparse Retrieval.  
**Nandan Thakur**, Kexin Wang, Iryna Gurevych, Jimmy Lin.  
SIGIR 2023 - Resource Track.
- [3] Injecting Domain Adaptation with Learning-to-hash for Effective and Efficient Zero-shot Dense  
Retrieval.  
**Nandan Thakur**, Nils Reimers, Jimmy Lin.  
ReNeuIR 2023 Oral Presentation.
- [4] Evaluating Embedding APIs for Information Retrieval.  
Ehsan Kamaloo, Xinyu Zhang, Odunayo Ogundepo, **Nandan Thakur**, David Alfonso-Hermelo,  
Mehdi Rezagholizadeh, Jimmy Lin.  
ACL 2023 - Industry Track.
- [5] Making a MIRACL: Multilingual Information Retrieval Across a Continuum of Languages.  
Xinyu Zhang\*, **Nandan Thakur\***, Odunayo Ogundepo, Ehsan Kamaloo, David Alfonso-Hermelo,  
Xiaoguang Li, Qun Liu, Mehdi Rezagholizadeh, Jimmy Lin. (\* denotes equal contribution)  
WSDM Cup Competition, 2023.
- [6] GPL: Generative Pseudo Labeling for Unsupervised Domain Adaptation of Dense Retrieval.  
Kexin Wang, **Nandan Thakur**, Nils Reimers, Iryna Gurevych.  
NAACL-HLT 2022.
- [7] BEIR: A Heterogenous Benchmark for Zero-shot Evaluation of Information Retrieval Models.  
**Nandan Thakur**, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, Iryna Gurevych.  
NeurIPS 2021 - Datasets and Benchmark Track.
- [8] Augmented SBERT: Data Augmentation Method for Improving Bi-Encoders for Pairwise Sentence  
Scoring Tasks.  
**Nandan Thakur**, Nils Reimers, Johannes Daxenberger, Iryna Gurevych.  
NAACL-HLT 2021.
- PREPRINTS     [9] Leveraging LLMs for Synthesizing Training Data Across Many Languages in Multilingual Dense  
Retrieval.  
**Nandan Thakur**, Jianmo Ni, Gustavo Hernández Ábrego, John Frederick Wieting, Jimmy Lin,  
Daniel Cer  
Arxiv Preprint, 2023.
- [10] HAGRID: A Human-LLM Collaborative Dataset for Generative Information-Seeking with Attribution.  
Ehsan Kamaloo, Aref Jafari, Xinyu Zhang, **Nandan Thakur**, Jimmy Lin.  
Arxiv Preprint, 2023.

- [11] Resources for Brewing BEIR: Reproducible Reference Models and an Official Leaderboard. Ehsan Kamaloo, **Nandan Thakur**, Carlos Lassance, Xueguang Ma, Jheng-Hong Yang, Jimmy Lin. Arxiv Preprint, 2023.
- [12] Simple Yet Effective Neural Ranking and Reranking Baselines for Cross-Lingual Information Retrieval. Jimmy Lin, David Alfonso-Hermelo, Vitor Jeronymo, Ehsan Kamaloo, Carlos Lassance, Rodrigo Nogueira, Odunayo Ogundepo, Mehdi Rezagholizadeh, **Nandan Thakur**, Jheng-Hong Yang, Xinyu Zhang. Arxiv Preprint, 2023.

**RESEARCH EXPERIENCE** **University of Waterloo** Sep 2021 - Present, Waterloo, ON, Canada  
*Ph.D. Student (Supervisor: Prof. Jimmy Lin)*  
 Worked on multilingual information retrieval [1] [12] [9], data and model efficiency [3] and reproducibility [2] [11]. Recently focusing on retrieval-augmented generation with LMs [10] [4].

**Google Research** Sep 2022 - May 2023, California, USA  
*Student Researcher (Mentors: Daniel Cer, Jianmo Ni)*  
 Worked on improving existing multilingual retrieval systems using PaLM 2 generated synthetic data, without expensive human-labeled training data for a wide variety of languages [9].

**UKP Lab, Technical University of Darmstadt** Nov 2019 - Aug 2021, Darmstadt, Germany  
*Research Assistant (Supervisors: Prof. Iryna Gurevych, Nils Reimers)*  
 Worked on developing a benchmark to evaluate zero-shot out-of-domain (OOD) evaluation of retrieval systems [7] and worked on data-augmentation techniques to generate synthetic data for domain adaptation in pairwise sentence [8] and retrieval tasks [6].

**(EMBL) European Molecular Biology Laboratory** Jun - Aug 2018, Heidelberg, Germany  
*Research Trainee (Supervisors: Prof. Toby Gibson, Dr. Manjeet Kumar)*  
 Worked on developing a prediction toolkit using machine learning to computationally predict kinase-substrate phosphorylation sites within (CAMK) protein sequences.

**INDUSTRY EXPERIENCE** **KNOLSKAPE** Sep 2018 - Oct 2019, Bengaluru, India  
*Data Scientist (Manager: Chaithanya Yambari)*  
 Worked on developing Krawler.ai, an enterprise product for effectively searching a company's large messy content libraries (pdf, xlsx, docx, etc.) with multimodal search. Implemented search functionality using Elasticsearch and backend data ingestion using Flask, Apache-Airflow and MongoDB.

**Belong.co** Jul - Dec 2017, Bengaluru, India  
*Data Science Intern (Manager: Vinodh K. Ravindranath)*  
 Worked on topic modeling for clustering millions of candidate resumes. Extracted keywords using Flash-Text and automatically clustered candidates using GuidedLDA, a semi-supervised LDA algorithm.

**HONOR AND AWARDS** **University of Waterloo (UW) Graduate Scholarship** 2021 - Present  
**BEIR benchmark** included in CS224U teaching material at Stanford University. 2021  
 Created both the ELLIS NLP 2021 and SustainNLP 2021 workshop websites. 2021  
 Got Selected as a speaker for PyCon Italia in 2020 (Cancelled due to Covid-19) 2020  
 Finalists in Technology Premier League (TPL) held by CIO & Leader, IT Next. 2019  
 Received a **fully-funded Machine Learning (ML) Fellowship** in EMBL Heidelberg 2018

|                        |   |  |
|------------------------|---|--|
| TEACHING<br>EXPERIENCE | <b>Teaching Assistant</b> , University of Waterloo<br><ul style="list-style-type: none"> <li>• CS 135 (Designing Functional Programs) - Fall 2021</li> <li>• CS 136 (Elementary Algorithm Design and Data Abstraction) - Winter 2022, Spring 2023</li> <li>• CS 241 (Foundations of Sequential Programs) - Spring 2022</li> <li>• CS 479/679 (Introduction to Artificial Intelligence) - Winter 2023</li> <li>• CS 370 (Numerical Computation) - Fall 2023</li> </ul>   | 2021 - Present   |
| SERVICES               | <b>Competition Lead Organizer:</b> WSDM Cup 2023.<br><b>Shared-task Lead Organizer:</b> TREC RAG 2024 (Upcoming)<br><b>Reviewer (*CL/NLP conferences):</b> ACL Rolling Review: Oct-Nov (2021), Jan-Apr (2022)<br><b>Reviewer (ML conferences):</b> NeurIPS 2023<br><b>Reviewer (IR conferences):</b> SIGIR 2023, ECIR 2024.   |  |
| INVITED TALKS          | <b>Koç University:</b> Advanced Information Retrieval (Tutorial)<br><b>Stanford University:</b> Heterogenous Benchmarking in IR Research<br><b>OpenNLP Meetup:</b> BEIR, An Open-Source Benchmark for IR Systems  | Virtual, June 2023<br>USA, November 2022<br>Virtual, June 2021                         |
| COURSEWORK             | <b>University of Waterloo:</b> CS 680: Introduction to Machine Learning (Ongoing), CS 889: Data Sources for Emerging Tech, CS 886: Graph Neural Networks, CS 886: Robustness of Machine Learning, CS 679: Neural Networks, CS 848: Information Retrieval, CS 649: Human-Computer Interaction, CS 854: Experimental Performance Evaluation.<br><b>BITS Pilani:</b> Machine Learning, Neural Networks & Fuzzy Logic, Data Structures & Algorithms, Probability & Statistics, Linear Algebra, Econometric Methods, Discrete Mathematics.   | 2021-Present<br>2014-2018  |
| PRESS AND<br>MEDIA     | <b>Making a MIRACL: Multilingual Information Retrieval Across a Continuum of Languages</b> , WSDM Cup Competition 2023<br><b>Domain Adaptation with Generative Pseudo-Labeling (GPL)</b> (Pinecone.ai)<br><b>Extending Neural Retrieval Models to New Domains and Languages</b> (Transformers-at-Work, Zeta Alpha)<br><b>BEIR benchmark as a helpful ML library</b> (ML News by Yannic Kilcher)<br><b>Making the Most of Data: Augmentation with BERT</b> (Pinecone.ai)<br><b>Advance BERT model via transferring knowledge from Cross-Encoders to Bi-Encoders</b> (Towards Data Science) | February 2023<br>August 2022<br>December 2021<br>August 2021<br>March 2021<br>May 2020 |
| COMPETENCES            | <b>Languages</b> Bengali ( <i>native</i> ), English ( <i>fluent</i> , TOEFL 110), Hindi ( <i>fluent</i> ), German ( <i>elementary</i> , A2)<br><b>Programming</b> Python, JavaScript, ReactJS, R, C++, HTML, CSS, Excel, MATLAB, Racket, $\LaTeX$ .<br><b>Libraries and Services</b> Pytorch, JAX, Tensorflow, Flask, Django, SQL, MongoDB, Docker, Elasticsearch, Redis, RabbitMq, Apache-Airflow, Postman.  |  |
| CO-<br>CURRICULAR      | <b>Mime Club Coordinator</b> , BITS Pilani<br>Led a team of 30 student performers in one of the most popular clubs in college. Involved in acting, sound mixing, designing slides and creating stories for more than 10 shows over a span of 4 years.   | 2014-2018  |