

JUDE KHOUJA +1 240 200 4022 jude@latynt.com

SUMMARY

NLP: LLM IFT, DPO, LoRA, Representation learning (contrastive and triplet loss), Semantic Clustering, Summarization

DL/ML: Transformers, EncDec, RNNs, KNN, SVM, PCA, Clustering, LDA **Tools:** Python, Pytorch, Deepspeed, AWS Sagemaker, hydra, SQL

Management: Establishing ML teams and functions, Scoping ML team's strategy and OKRs, Recruiting and Interviewing, Mentoring ML Scientists and Interns and establishing data and annotation partnerships.

HONORS Fulbright Scholarship, U.S. Department of State 2010.

SELECTED PUBLICATIONS

Hajij et al. "TopoX: A Suite of Python Packages for Machine Learning on Topological Domains". (In review)

Khouja, J. Stance Prediction and Claim Verification: An Arabic Perspective. Proceedings of the Third Workshop on Fact Extraction and VERification (FEVER) workshop at ACL 2020

Knowledge and Theme Discovery across Very Large Biological Data Sets Using Distributed Queries: A Prototype Combining Unstructured and Structured Data. PLOS 2013

Mr. LDA: A Flexible Large Scale Topic Modeling Package using Variational Inference in MapReduce. ACM International Conference on World Wide Web, 2012.

HIGHLIGHTED EXPERIENCE

PRINCIPAL APPLIED SCIENTIST, FORETHOUGHT.AI Dec 22 - Oct 23

Leading work in Large Language Modeling instruction-finetuning and evaluation. Built FLAN2-style instruction templates of 12 tasks using production, synthetic and human annotations.

SR. STAFF APPLIED SCIENTIST, FORETHOUGHT.AI Feb 21 - Nov 22

Built and deployed a ML pipeline for semantic clustering of support tickets which included training new T5-based abstractive summarizer and a custom encoder model trained with triplet loss from synthetic and human labels which improved offline clustering coverage by 46% relative and 8% in online metrics for beta customers while reducing duplicates. The work resulted in the release of a new product (discover).

Designed and implemented company's first ML experimentation repo and infrastructure using pytorch, hydra, mlflow and sagemaker. Grew and managed the ML team and established ML infra functions. Setup hiring and interviewing processes and advised leadership team on ML initiatives, OKRs and objectives.

FOUNDER, LATYNT LLC Sep 19 - Present

ML advising and consulting for the private sector. Released Arabic News Stance Corpus (ANS) one of the first datasets to study misinformation detection in the Arabic language with BERT-based baselines.

SR. PRINCIPAL ML SCIENTIST, SAGE INTACCT Apr 20 – Jan 21

Lead the company's first AI product by revamping all ML models for time tracking which improved relative F1-score 60-80%. Helped build the ML engineering team and hire ML and Data engineers.

SENIOR APPLIED SCIENTIST, MICROSOFT May 18 – Sep 19

Contributed to scaling up Language Model distributed training to tens of billions of words and explored the use of sub-word representations (BPE, Wordpiece) in speech services's language modeling team.

SR. /PRINCIPAL DATA SCIENTIST, SALESFORCE Sep 15 – Feb 18

Drove the team's NLP DL foundation and automated business processes by building text classification models for the customer support processes.

Developed reusable machine learning pipeline and feature engineering libraries.

DATA SCIENCE LEAD, CRITTERCISM Jul 14 – Nov 14

Applied clustering techniques and text matching for cardinality reduction of mobile device types.

Built an internet facing analytics portal for tracking live mobile performance metrics worldwide (Patented).

Built AWS on demand analytics infrastructure using S3 and EMR.

DATA ANALYTICS LEAD, IREX Nov 13 – Jun 14

Lead the organizational technical projects, processes and team for documenting human rights violations from crowdsourced content.

Oversaw the design and implementation of video and image based annotation systems.

BIG DATA SPECIALIST, ORACLE Mar 13 – Oct 13

Prototyped machine learning proof of concepts in the public sector using technologies including Hadoop, Hive, Pig, R Enterprise, Mahout and other proprietary and open source tools.

DATA SCIENTIST, ORACLE Jun 12 – Aug 12

Lead Data Scientist in the Oracle/NCI partnership project that won the "2012 Best Government Big Data Solution" Award.

Developed MapReduce programs in Java and Python for generating synthetic Gene data of 900 million patients.

GRADUATE RESEARCH ASSISTANT, UMIACS Oct 11 – May 12

Evaluated distributed Topic Modeling (LDA) algorithms and applied them for unsupervised lexicon expansion.

Developed an Arabic version of the Word Count tool (LIWC) for sentiment analysis and honor dictionary validation.

EDUCATION & Training

University of Oxford - 2023 - present (part-time)

Ph.D. in Information, Communication and the Social Sciences

Stanford University - 2017

Graduate coursework in NLP with Deep Learning (CS224n)

University of Maryland College Park - 2012

Masters in Information Management

Damascus University - 2007

B.E. Computer Science, Focus: Artificial Intelligence