

# Aditeya Baral

+1(551)263-8608 | [aditeyabaral@nyu.edu](mailto:aditeyabaral@nyu.edu) | [linkedin.com/in/aditeyabaral](https://www.linkedin.com/in/aditeyabaral) | [aditeyabaral.com](https://aditeyabaral.com) | [Google Scholar](https://scholar.google.com/citations?user=aditeyabaral)

## EDUCATION

**New York University, Courant Institute of Mathematical Sciences**

*Master of Science in Computer Science; GPA – 3.87/4.00*

**Concentration:** Artificial Intelligence

New York City, USA

*Sep 2024 – May 2026*

**PES University**

*Bachelor of Technology in Computer Science & Engineering; GPA – 8.71/10.00*

**Specialization:** Machine Intelligence & Data Science

Bengaluru, India

*Aug 2018 – May 2022*

## EXPERIENCE

**Redis**

*Applied Research Scientist Intern, Redis LangCache; Advisor: Srijith Rajamohan*

San Francisco, USA

*June 2025 – Dec 2025*

- Introduced **Precision–Cache Hit Ratio (P-CHR)** and **Calibration Retention Rate (CRR)**, two *cache-aware metrics* for **Redis LangCache**, reframing semantic-cache model selection as a *calibration problem* rather than a ranking one.
- Curated and open-sourced **LangCache SentencePairs (v1-v3)**, a large-scale dataset family with **1M to 40M** sentence pairs across diverse paraphrase, STS, QA, and adversarial sources for robust semantic-caching fine-tuning.
- Fine-tuned and deployed **LangCache-Embed-v3**, a domain-specific bi-encoder trained with an *ArcFace contrastive objective*, leading all open-source retrievers offline and outperforming larger general-purpose models even without re-ranking.
- Open-sourced the **LangCache ReRanker v1/v2** cross-encoder families with *BCE* and *MNRL* objectives at two training scales, isolating the effect of objective vs. scale on calibration and enabling application-specific behavior.
- Built a comprehensive *evaluation framework* for customer onboarding with a large-scale study of customer data and model baselines across offline and deployment settings, using a two-stage **RedisVL** K-NN retrieval and re-ranking pipeline.
- Demonstrated that the highest-**PR-AUC** models are often the *worst* in deployment, and that re-ranking rarely improves a strong domain retriever, overturning the standard practice in semantic-cache model selection.
- Supported *downstream integration* and development of **LMCache** by *building prototypes* and *benchmarking performance* with **Redis** as an *in-memory KV store*, demonstrating latency and throughput gains.

**Cisco Systems**

*Big Data Engineer, Webex Media Quality Analytics*

Bengaluru, India

*July 2022 – July 2024*

- Developed and deployed streaming jobs in **Scala** and **Flink** to process **1M+ reports/min** and compute **1200+ real-time metrics** from Calls and Meetings.
- Applied *statistical modelling* techniques to investigate and report *media quality insights* to downstream consumers, *reducing errors by 30%* and *analysis time by 15 hrs/week* per team member.
- Led the development of *real-time (<1 min) auditing pipelines* using **Kafka** and **Python** to ensure *per-minute data consistency* between streaming jobs and **Iceberg** and **Pinot** data stores, *reducing manual effort by >80%*.
- Built graphs and dashboards on the **Webex Media Quality Analytics Dashboard** using **Grafana** and **Kibana** to set up alerts and KPIs for **20,000+** clients and customers.

*Big Data Engineering Intern, Webex VideoMesh Analytics*

*Jan 2022 – June 2022*

- Migrated the **Meetings Analytics Engine** from Java and Spark to **Scala** and **Flink** to scale up to **1M+ reports/min** and significantly *improve real-time report generation* by over **40%**.
- Built **VideoMesh Developer APIs** using **Java** and globally rolled them out for **30,000+ enterprises** with **customer-facing applications**.

**Intel Corporation**

*Applied Research Scientist Intern, Intel VSG; Advisors: Anay Majee, Anbumani Subramanian*

Bengaluru, India

*Aug 2021 – Dec 2021*

- Explored **Few-Shot Learning Object Detection (FSOD)** techniques to reduce *catastrophic forgetting* in constrained and heterogeneous driving environments.
- Investigated and designed novel *representation learning* and *attention mechanisms* to learn *inter/intra-object relationships* using **PyTorch**.
- Outperformed existing approaches at the time on base and novel classes by **0.2 mAP** and **3 mAP** on the *Few-Shot India Driving Dataset*, a benchmark for FSOD.

## SKILLS

**Languages:** Python, Scala, Java, C/C++, Groovy, SQL, L<sup>A</sup>T<sub>E</sub>X

**ML/Stats Libraries:** PyTorch, Tensorflow, HuggingFace, WandB, vLLM, FAISS, pandas, NumPy, scikit-learn, matplotlib

**AI/ML Techniques:** Representation Learning, Mechanistic Interpretability, Transfer Learning, Language Models, RAG

**Big Data/Cloud:** Hadoop, Kafka, Zookeeper, Spark, Flink, Iceberg, Pinot, Redis, ELK

**Frameworks/Tools:** Git, GitHub, Jenkins, Docker, Kubernetes, FastAPI, Grafana, PSQ, MongoDB, AWS, Linux

- [1] **Closing the Calibration Gap in Semantic Caching**  
*Authors: Aditeya Baral\*, Radoslav Ralev\*, Iliya Sotirov Zhechev, Srijith Rajamohan, Jen Agarwal*
- [2] **When ‘+’ Means ‘-’? Probing Arithmetic Circuits Under Symbol Redefinition**  
*Authors: Aditeya Baral, Allen George Ajith, Shauli Ravfogel*
- [3] **Can LLMs *understand* Math? Exploring the Pitfalls in Mathematical Reasoning**  
*Authors: Tiasa Singha Roy\*, Aditeya Baral\*, Ayush Rajesh Jhaveri, Yusuf Baig*
- [4] **CMLFormer – A Dual Decoder Transformer with Switching Point Learning for Code-Mixed Language Modeling**  
*Authors: Aditeya Baral, Allen George Ajith, Roshan Nayak, Mrityunjay Abhijeet Bhanja*
- [5] **Patch and Control – Steering Behavior of Large Vision-Language Models via Latent Activations**  
*Authors: Aditeya Baral, Rijul Dahiya, Dilip Venkatesh*

PAPERS AND PUBLICATIONS

---

- [1] **Training for Compositional Sensitivity Reduces Dense Retrieval Generalization**  
*ICLR 2026, Sci4DL*  
*Authors: Radoslav Ralev, Aditeya Baral, Iliya Sotirov Zhechev, Jen Agarwal, Srijith Rajamohan*
- [2] **ChatBERT – Multi-task approach to Pre-Training for Structured Conversations**  
*Webex AI 2023*  
*Authors: Aditeya Baral (Work done as part of Cisco Webex AI Research)*
- [3] **CalBERT – Code-mixed Adaptive Language Representations using BERT**  
*AAAI 2022, MAKE*  
*Authors: Aditeya Baral, Aronya Baksy, Ansh Sarkar, Deeksha D, Ashwini M Joshi*
- [4] **Information Maximization to Overcome Catastrophic Forgetting in Few-Shot Object Detection**  
*Intel VSG 2021*  
*Authors: Aditeya Baral, Anay Majee, Anbumani Subramanian*
- [5] **MAPLE – MAsking words to generate blackout Poetry using Seq2Seq LERning**  
*ACL 2021, ICNLSP*  
*Authors: Aditeya Baral, Himanshu Jain, Deeksha D, Mamatha H R*
- [6] **Analysis of Kepler Objects of Interest using ML for Exoplanet Identification**  
*IEEE CONIT 2021*  
*Authors: Ameya Rajendra Bhamare, Aditeya Baral, Saarthak Agarwal*