

# Information Retrieval Course Project Proposal

Prateek Kumar Padhy  
Computer Science  
Virginia Tech  
Falls Church, VA, USA  
pprateekkumar@vt.edu

Jayanth Koripalli  
Computer Science  
Virginia Tech  
Falls Church, VA, USA  
jayanth28@vt.edu

Siji Chen  
Computer Science  
Virginia Tech  
Falls Church, VA, USA  
sijic@vt.edu

**Abstract**—Large Language Models (LLMs) have opened the door for revolutionary data analysis methods and insight production in information retrieval. The aim of this project is to investigate and optimize open-source LLMs, particularly those with parameter values under 2 billion, so as to meet the limitations of the available computing power—a 16GB GPU, in particular. This study attempts to find an effective yet potent model that can enable advanced crime data analysis by comparing models like LLaMA-3B, BERT, and RoBERTa. The project is structured around key milestones, including proposal submission, checkpoints, and final presentations, culminating in the application of the selected LLM to three critical areas: identifying crime hotspots, generating public safety reports, and developing a community engagement platform. Through the integration of LLMs with a targeted preprocessing and data handling backend and an intuitive front end, this study aims to provide new perspectives on crime patterns and risk factors, which will benefit the public safety and information science fields.

## I. INTRODUCTION

Crime prevention and prediction is a crucial problem in our society and hence attracted a lot of attention. In academia field, researchers have proposed numerous studies on crime analysis and predictions from various aspects, such as crime hotspot discovery [1], incidence analysis and crime predictions [2]. Meanwhile, the advent of large language models (LLMs) has revolutionized natural language processing, enabling it with the capability to comprehend inference and reasoning tasks. Leveraging this advancement, we will employ the LLM to analyze thousands of crime-related news articles from Timesofindia, delving into the narratives concealed within the headlines.

Our application will furnish users with an intuitive interface to explore crime-related content through a free text queries. For example, users are able to find crime based on location, for which we will use the LLM to extract the location mentioned in the article. In addition to the search interface for information retrieval, we will incorporate a dashboard function to visually display the analyses and patterns unearthed during data mining, such as crime hotspots within the city.

## II. RELATED WORK

### A. GraphBERT

GraphBERT [3] uses an integrated framework to assist the malicious behavior detection on social media platform. It's proposed to bridging the gap between graph and text; which

can detect the behaviors with both semantic and user relational information. This paper doesn't have dataset/code published. They mined twenty million tweets from over one million twitter users; after cleanup, 91,500 tweets from 13,351 users are selected. Tasks include malicious user detection; malicious tweet classification.

### B. Discrimination/ hate crimes across U.S. Cities

[4] City-level analysis. They combined hate crime data from FBI; twitter data collected from 2011 to 2016; census data. Tasks include classifications on discrimination; spatial crime distribution; hate crimes and social media relationship. They found the top 20 features that can predict 25% of the discrimination. Phoenix, AZ had the highest number of race, ethnicity or national-origin based hate crimes over the 6 years.

### C. Crime Analysis and Prediction using Data Mining

[5] Applying data mining techniques to the interface of computer science and criminal justice represents a revolutionary approach to crime research and prevention. This methodology allows the visualization of crime-prone areas and the prediction of regions with a high probability of crime occurrence by focusing on the systematic identification and analysis of crime patterns and trends rather than traditional factors like political enmity or the criminal background of the offender. The introduction of electronic systems has greatly improved the skills of crime data analyzers, enabling law enforcement officials to solve crimes more quickly. Using unstructured data to extract previously undiscovered and valuable information, this method highlights a move toward utilizing technological improvements for more effective and efficient techniques.

### D. Crime Prediction Using Machine Learning and Deep Learning: A Systematic Review and Future Directions

[6] The field of crime prediction has witnessed a significant transformation with the advent of machine learning and deep learning techniques. A comprehensive review which explored over 150 articles, emphasizing the diverse range of algorithms employed in crime prediction. This body of work underlines the critical role of predictive analytics in identifying crime patterns and trends, offering invaluable insights for law enforcement strategies. Key findings from the review highlight the efficacy of various machine and deep learning models in

not only forecasting crime occurrences but also in unveiling underlying factors contributing to criminal activities. The review further sheds light on the datasets commonly used in this research area, providing a foundational resource for future investigations. Acknowledging the potential gaps identified in current methodologies the paper propose future directions aimed at enhancing prediction accuracy, such as integrating socio-economic indicators and adopting more sophisticated neural network architectures. This systematic review serves as a cornerstone for our research, guiding our exploration of novel approaches in crime prediction and contributing to the ongoing dialogue within the academic community on improving public safety through advanced analytics.

### III. PROPOSED APPROACHES

#### A. Selection and Optimization of LLMs

Our aim is to identify the most efficient and effective Large Language Models (LLMs) that are suitable for analyzing crime data within the computational constraints of a 16GB GPU. We plan to evaluate the performance of various LLMs, including LLaMA-3B, BERT, and RoBERTa, on benchmark datasets related to crime. I will implement parameter tuning and model pruning techniques to optimize these models for enhanced performance while ensuring they remain computationally efficient.

#### B. Data Preprocessing and Integration

We intend to develop a robust preprocessing pipeline that can clean, standardize, and enrich crime datasets for more effective analysis. We will design and implement a preprocessing pipeline that incorporates data cleaning, normalization, and the integration of additional data sources for a richer analysis context. We aim to utilize geospatial data processing techniques to aid in the identification and categorization of crime hotspots.

#### C. Crime Hotspot Identification

Our goal is to use optimized LLMs to analyze crime data effectively and identify potential crime hotspots. We plan to apply natural language processing (NLP) techniques to extract and analyze relevant information from crime reports and related texts. We will integrate geospatial analysis with NLP findings to identify and visualize areas with high crime rates.

#### D. Public Safety Report Generation

We aim to automatically generate comprehensive public safety reports based on the analyzed crime data. We will develop a template for public safety reports that includes essential metrics, trends, and actionable recommendations. Using the optimized LLMs, we plan to populate this template with data-driven insights and analyses.

#### E. Community Engagement Platform Development

We intend to design and implement a web-based platform that is integrated with the backend LLM analysis system. The platform will feature interactive maps, report generation capabilities, and mechanisms for community feedback.

## IV. EXPERIMENT

#### A. Large Language Model (LLM)

We are looking for open-sourced LLM that could be integrated into our project. Based on each type of tasks, we maybe find out the need to use inference only or fine-tune the open source LLM. Hence the GPU memory is a consideration. The GPU Quadro we have is 16GB so need a compact and accurate model which has less than 2 billion parameters. To run the 7B model in full precision, one need  $7 \times 4 = 28GB$  of GPU RAM. One should add `torch_dtype = torch.float16` to use half the memory and fit the model on a T4 which is 16GB. To train, 56GB is needed to store gradients for each parameters. Table I summarized the popular LLM size.

| Model     | Parameter | Memory Full precision |
|-----------|-----------|-----------------------|
| LLaMA-3B  | 3B        | 12 GB                 |
| LlaMA -7B | 7B        | 28 GB                 |
| Vicuna    | 13B       | 52 GB                 |
| BERT      | 345M      | 1.4 GB                |
| RoBERTa   | 355M      | 1.5 GB                |
| MPT-7B    | 7B        | 28 GB                 |
| Flan-T5   | 11B       | 44 GB                 |

TABLE I  
CURRENT OPEN SOURCE LLM PARAMETERS

Based on Table I, we plan to start with the following open source model: 1) BERT based IR+QA Paper 2)Pytorch roBERTa 3)Keras/TensorFlow roBERTa

#### B. Streamlit

Our project will make use of the open-source app framework Streamlit to provide a user-friendly API that links the front end and our backend Large Language Models (LLMs). We will be able to quickly prototype and implement a dynamic interface for our crime data analysis platform thanks to Streamlit's ease of use and effectiveness in developing interactive web apps. With the use of this interface, customers will be able to communicate with our LLMs in real-time and get instant access to functions like community involvement, public safety report generating, and crime hotspot identification. We want to make sure that our platform is not just strong in its analytical skills but also easily accessible and interesting to a diverse user base, which is why we have planned on integrating Streamlit.

## V. DATASET

Crimes in India India is a country of 1.3 billion people, along with it's high population there's tens of thousands of crimes recorded all over the country. This dataset contains links to the articles from times of India website which describes the M.O. of the crimes. This dataset will help us to learn crime trends in different regions of India

Crimes in Washington DC The Washington DC government provided this carefully selected information, which is a crucial representation of specific locations and characteristics of criminal incidences that were reported in 2023. It helps scholars, decision-makers, and the general public comprehend the crime landscape in the nation’s capital by providing a window into the dynamics of the city’s criminal activity. This dataset comprises of 26,091 features along with 28 columns. Key attributes include neighbourhood clusters, offence types, geographical details, timestamps, and associated crime classifications.

## VI. APPENDIX

### A. Timelines

- 1) 2/20 Proposal
- 2) 3/19 check point
- 3) 4/9 presentation
- 4) 4/23, 4/23 final presentation
- 5) 5/7 project due

### B. Tasks

- 1) User interface: webapp/Tableau/Streamlit.
- 2) Backend: model training; model validation and testing; API to UI(maybe not need)
- 3) database: maybe we don’t need inmemory, postgresql and others.
- 4) Testing
- 5) Document writing;

## REFERENCES

- [1] J. E. Eck, S. Chainey, J. G. Cameron, M. Leitner, and R. E. Wilson, “Mapping Crime: Understanding Hot Spots.”
- [2] H. Wang, D. Kifer, C. Graif, and Z. Li, “Crime Rate Inference with Big Data,” in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. San Francisco California USA: ACM, Aug. 2016, pp. 635–644. [Online]. Available: <https://dl.acm.org/doi/10.1145/2939672.2939736>
- [3] J. Wu, C. Zhang, Z. Liu, E. Zhang, S. Wilson, and C. Zhang, “GraphBERT: Bridging Graph and Text for Malicious Behavior Detection on Social Media,” in *2022 IEEE International Conference on Data Mining (ICDM)*. Orlando, FL, USA: IEEE, Nov. 2022, pp. 548–557. [Online]. Available: <https://ieeexplore.ieee.org/document/10027673/>
- [4] K. Relia, Z. Li, S. H. Cook, and R. Chunara, “Race, Ethnicity and National Origin-Based Discrimination in Social Media and Hate Crimes across 100 U.S. Cities,” *Proceedings of the International AAAI Conference on Web and Social Media*, vol. 13, pp. 417–427, Jul. 2019. [Online]. Available: <https://ojs.aaai.org/index.php/ICWSM/article/view/3354>
- [5] S. Sathyadevan, M. Devan, and S. S. Gangadharan, “Crime analysis and prediction using data mining,” in *2014 First international conference on networks & soft computing (ICNSC2014)*. IEEE, 2014, pp. 406–412.
- [6] V. Mandalapu, L. Elluri, P. Vyas, and N. Roy, “Crime prediction using machine learning and deep learning: A systematic review and future directions,” *IEEE Access*, 2023.