# URL CLASSIFICATION

- Suman Kanukollu
- https://www.linkedin.com/in/suman-kanukollu/
- https://www.youtube.com/@sumankanukollu

# INTRODUCTION





- Categorize & filter
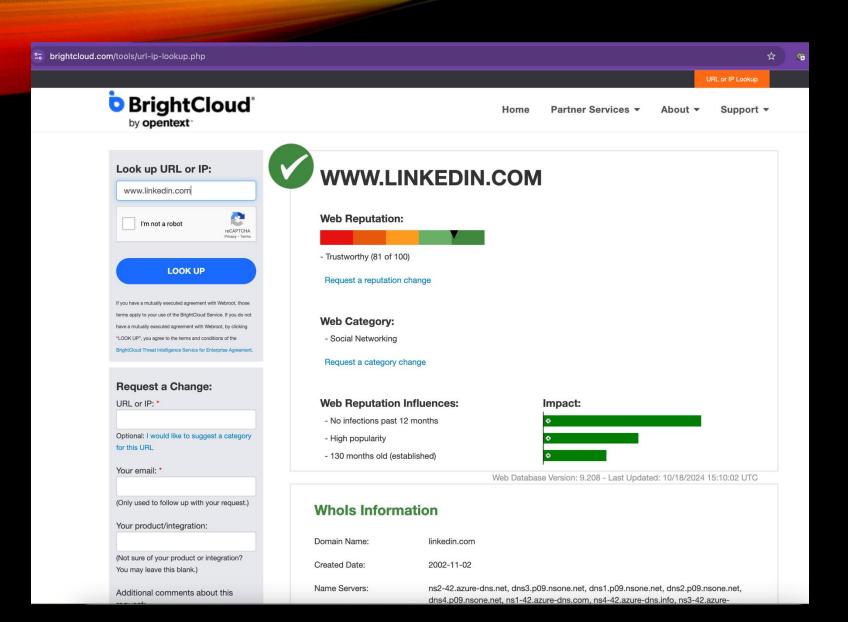Millions of WEB sites in real time?

# SOLUTION OVERVIEW

- Collect data using web scraping (beautiful soup)

- Apply NLTK techniques

- Spacy (large model) extract 18-named Entities

- Sentence-transformers/all-MiniLM-L6-v2 (384-dim embedding vector)

- 786 Vector = 384-Title + 384-Text + 18-NE's

- Encoder only Transformer (with 6-multi-head attention blocks)

- Given a URL, this model processes the content and classifies it into categories like 'Business,' 'Shopping,' or 'News', 'Job Search' ….etc

# DEMO (WITH CHROME EXTENSION)

# REAL WORLD APPLICATIONS

1.  WEB Security

2.  Add targeting

3.  SEO and spam detection

4.  Parental control

5.  Cyber Security

# CONCLUSION

This project showcases the potential of transformers for efficient URL classification, with real world applications in security, marketing, content filtering and much more.