

FINAL PROJECT: NEWSBOT INTELLIGENCE SYSTEM 2.0

Reflection Journal

Course Information:

NLP ITAI 2373

Submitted by:

Faiza Abdullah

([https://github.com/AbdullahFaiza/NLP-](https://github.com/AbdullahFaiza/NLP-ITAI2373/tree/main/Final%20Project%3A%20NewsBot%20Intelligence%20System%202.0)

ITAI2373/tree/main/Final%20Project%3A%20NewsBot%20Intelligence%20System%202.0)

Professor:

Patricia Mcmanus

Key Insights on NewsBot 2.0 Development

The development of NewsBot 2.0 was a transformative experience, revealing the complexity and power of integrating advanced NLP techniques into a cohesive system. One profound insight was the critical role of system integration in ensuring that individual components—like topic modeling, transformer fine-tuning, and multilingual processing—work seamlessly together. The project's modular architecture taught me how to orchestrate diverse NLP tasks, from preprocessing 39 multilingual articles to generating trend predictions (e.g., `neutral: 53.85%`, `negative_media: 46.15%`). This process underscored that a successful NLP system is not just about individual algorithms but about creating robust workflows that handle real-world data challenges, such as inconsistent formats or missing labels.

Another key takeaway was the importance of iterative debugging and validation. The project required resolving multiple errors (e.g., `eval_strategy`, `insights parsing`, `metadata.widgets`), which highlighted how NLP systems demand meticulous attention to detail. For instance, fixing the Enhancement success rate from 0% to 100% by ensuring valid `insight_category` outputs showed me that even small oversights in data structures can cascade across the pipeline. This experience deepened my appreciation for validation as a tool to catch and correct issues early, ensuring reliability for production-ready systems like NewsBot 2.0.

Challenges Encountered

The project presented several technical challenges that tested my problem-solving skills. One significant issue was the Enhancement success rate of 0%, caused by the `ContentEnhancementEngine` failing to produce valid insights with an `insight_category`. Initially, I assumed the issue was in the summarization step, but debugging revealed that insights were being saved as empty dictionaries. Fixing this required modifying the `enhance_content` method to ensure consistent dictionary outputs (e.g., `{'insight_category':`

'neutral'})), which took several iterations and taught me the importance of logging intermediate outputs (e.g., Sample insights in outputs/logs.txt).

Another challenge was the transformer fine-tuning error in (Bonus part) `TrainingArguments.__init__()` got an unexpected keyword argument 'evaluation_strategy'. This arose from an outdated transformers library argument, which I resolved by updating to `eval_strategy`, compatible with `transformers>=4.28.0`. The process of identifying and fixing this error, which occurred around 8:07 PM on August 7, 2025, emphasized the need to stay updated with library APIs and verify dependency versions (e.g., `!pip show transformers`).

Another bonus area of trend prediction error ('str' object has no attribute 'get') was particularly frustrating, as it stemmed from stringified insights in `enhanced_articles.csv`. I initially struggled to understand why `insights.get('insight_category')` failed, but inspecting `df_enhanced['insights'].head()` revealed JSON-like strings (e.g., `'{\'insight_category\': \'neutral\'}'`). Adding `json.loads` to parse these strings fixed the issue, producing correct trend scores (neutral: 53.85%, negative_media: 46.15%). This challenge highlighted the importance of data serialization in pipelines.

The final hurdle was a Git error ('state' key is missing from 'metadata.widgets') when saving the notebook. This was caused by incomplete widget metadata from interactive outputs (e.g., matplotlib plots). Running a cleaning script to remove `metadata.widgets` and setting `matplotlib.use('Agg')` in Cell 1 resolved it, but it required learning about Jupyter's metadata structure, a process that took significant time and caused immense frustration leaving me not much time to work on bonus segment of Website development.

Connections to Real-World Applications

The project's diverse components connected directly to real-world NLP applications. The multilingual processing, using M2M100 to translate 39 articles (English, Spanish, French), mirrors how global news platforms like BBC or Reuters analyze cross-language coverage.

Preserving cultural nuances during translation, as seen in `translated_articles.csv`, is critical for understanding regional perspectives, such as differing views on climate change.

The transformer fine-tuning, enhanced with Wandb logging, relates to industry practices at companies like Hugging Face, where models are fine-tuned for specific tasks (e.g., news categorization). Integrating the Wandb API (`caf3b27fcb1507b0a0cfd31a9dab2ec98670ad5`) to track metrics like training loss felt like a real-world MLOps workflow, reinforcing the value of experiment tracking in production systems.

Trend prediction has clear applications in social media analytics, where platforms like X use topic modeling to identify trending topics. The visualization (`trend_prediction.png`) showing neutral (53.85%) and negative_media (46.15%) trends could inform media strategies, such as targeting positive tech news. The conversational interface aligns with chatbot development, enabling queries like “Latest news on climate change,” which companies like Google use for news aggregation.

Comparisons Between Approaches

Using `nlk` for tokenization and stop word removal was lightweight and fast, similar to Lab 2, but less robust for complex tasks. In contrast, sentence-transformers for semantic search provided richer embeddings, enabling accurate query matching (e.g., top-5 articles with similarity >0.1). This trade-off between simplicity (NLTK) and accuracy (sentence-transformers) mirrors Lab 2’s NLTK vs. `spaCy` comparison, confirming that task-specific tools are critical.

For topic modeling, LDA implementation was effective for identifying themes (e.g., neutral, negative_media) but required careful preprocessing to avoid noise. Compared to potential alternatives like NMF, LDA was chosen for its interpretability, though NMF could offer better sparsity for larger datasets. This decision-making process echoed Lab 2’s stemming vs. lemmatization debate, where task goals (speed vs. accuracy) drive choices.

The use of transformers versus simpler models like TF-IDF highlighted a trade-off between computational cost and performance. Fine-tuning `distilbert-base-uncased` was resource-intensive (20–30 minutes on CPU), but its accuracy in categorizing articles outperformed TF-IDF-based methods. This reinforced that advanced models require more setup but yield superior results for complex tasks.

Future Applications

This project has inspired several ideas for future NLP applications. One is a real-time news monitoring system, extending multilingual processing to integrate APIs like NewsAPI for live article feeds. This could power a dashboard for journalists, using trend prediction to highlight emerging topics. Incorporating real labels (e.g., via NewsAPI) would improve fine-tuning accuracy, making it suitable for industry-grade classification.

Another application is enhancing the conversational interface into a full chatbot for news exploration, similar to customer service bots. Using spaCy's POS tagging could improve query parsing, enabling complex inputs like "Compare tech news from 2024 and 2025." The project's validation framework could be adapted to monitor chatbot performance, ensuring high success rates.

Applying fact-checking by integrating APIs like Google Fact Check Tools, building on the mock implementation, is another experimental area. This would create a robust bias detection system for media companies, leveraging sentiment analysis and external verification to ensure credible reporting. These applications feel achievable, thanks to the project's comprehensive pipeline and debugging lessons.

Conclusion

Developing NewsBot 2.0 was a challenging yet rewarding journey, transforming my understanding of NLP from basic preprocessing to a sophisticated, integrated system. Overcoming errors like `Enhancement 0%`, `eval_strategy`, `insights parsing`, and

metadata.widgets taught me resilience and the importance of systematic debugging. The project's real-world connections to news analytics and chatbots made it highly relevant, while comparisons between tools (e.g., NLTK vs. transformers) deepened my technical decision-making skills. I'm excited to apply these insights to future projects, confident that NewsBot 2.0 has equipped me with the skills to tackle complex NLP challenges.