

# Lost in Translation: How AI Misinterprets Modern Online Language and Emojis

## Understanding NLP Failures on Culturally Rich and Emoji-Laden Comments

Malachi Randolph

Independent Researcher | Los Angeles, CA

### Background

- NLP tools often fail to interpret internet language—informal, expressive, and filled with slang and emojis.
- This language includes terms like “goat,” “queen,” and “fire” which express sentiment and tone beyond literal meaning.
- Misinterpretations can lead to biased outcomes, especially when language rooted in marginalized communities is treated as “noise” or “foreign.”
- This project examines how sentiment and language models misread culturally influenced internet language.



### Research Questions

- How do NLP models perform on comments that use informal internet language?
- Can human annotation and active learning improve outcomes?
- What are the social implications of these model failures?

### Misclassification Examples

YouTube: Kendrick Lamar's Apple Music Super Bowl Halftime Show					
	Model Sentiment Label	Human Sentiment Label	langid language label	langdetect language label	langua language label
@urbestie1234 Ts so ah 🤡	negative	positive	Khmer	Somali	Sotho
@kdot_fan2025 GOATED HALFTIME	neutral	positive	English	German	English
@the45greatest THEY NOT LIKE US	negative	neutral	English	Vietnamese	English
@strawberriesareyum93	neutral	positive	Hebrew	unknown	unknown

YouTube: Beyoncé & Bruno Mars Crash the Pepsi Super Bowl 50 Halftime Show   NFL					
	Model Sentiment Label	Human Sentiment Label	langid language label	langdetect language label	langua language label
@thequeenbeee That was fenomenol 🤩	negative	positive	English	English	Malay
@natalag3231 So we not gone talk bout Beyoncé all most falling on stage?	negative	neutral	French	English	English
@cheesypleaseformee Daaamn up town funk is soo smooth	negative	positive	English	English	Xhosa

### Internet Language & Cultural Influence

- Much of what is considered Gen Z or internet slang is rooted in African American Vernacular English (AAVE)
- Examples:
  - G.O.A.T. — stands for Greatest of All Time / admiration
  - Fire — approval / excitement
  - Queen — empowerment / admiration
- These phrases carry strong emotional tones, but are often misread by NLP tools not trained on culturally diverse data
- Treating such language as abnormal introduces risk of digital exclusion

### Responsible AI Insight

- Annotations were guided by lived experience and supported by clear labeling guidelines
- Even with limited resources, the project prioritized transparency and fairness
- Highlights need for diverse annotator teams, culturally aware training data, and inclusive evaluation

### Implications & Future Work

- Mislabeling internet language can marginalize voices and distort meaning in social data
- NLP models need to adapt to evolving digital expression
- Recommendations:
  - Dialect-aware open-source datasets
  - Community-involved annotation
  - Responsible model design with the Bender Rule in mind
  - Incorporate AAVE with Standard American English when training models
- Broader goal: ensure NLP systems don't erase or distort culturally rich communication

### Dataset

- 27,000+ comments from Beyoncé's 2016 halftime show
- 134,000+ comments from Kendrick Lamar's 2025 halftime show
- Both performances generated high engagement and culturally rich language, while Kendrick's performance had more emoji-rich comments
- Collected using YouTube API

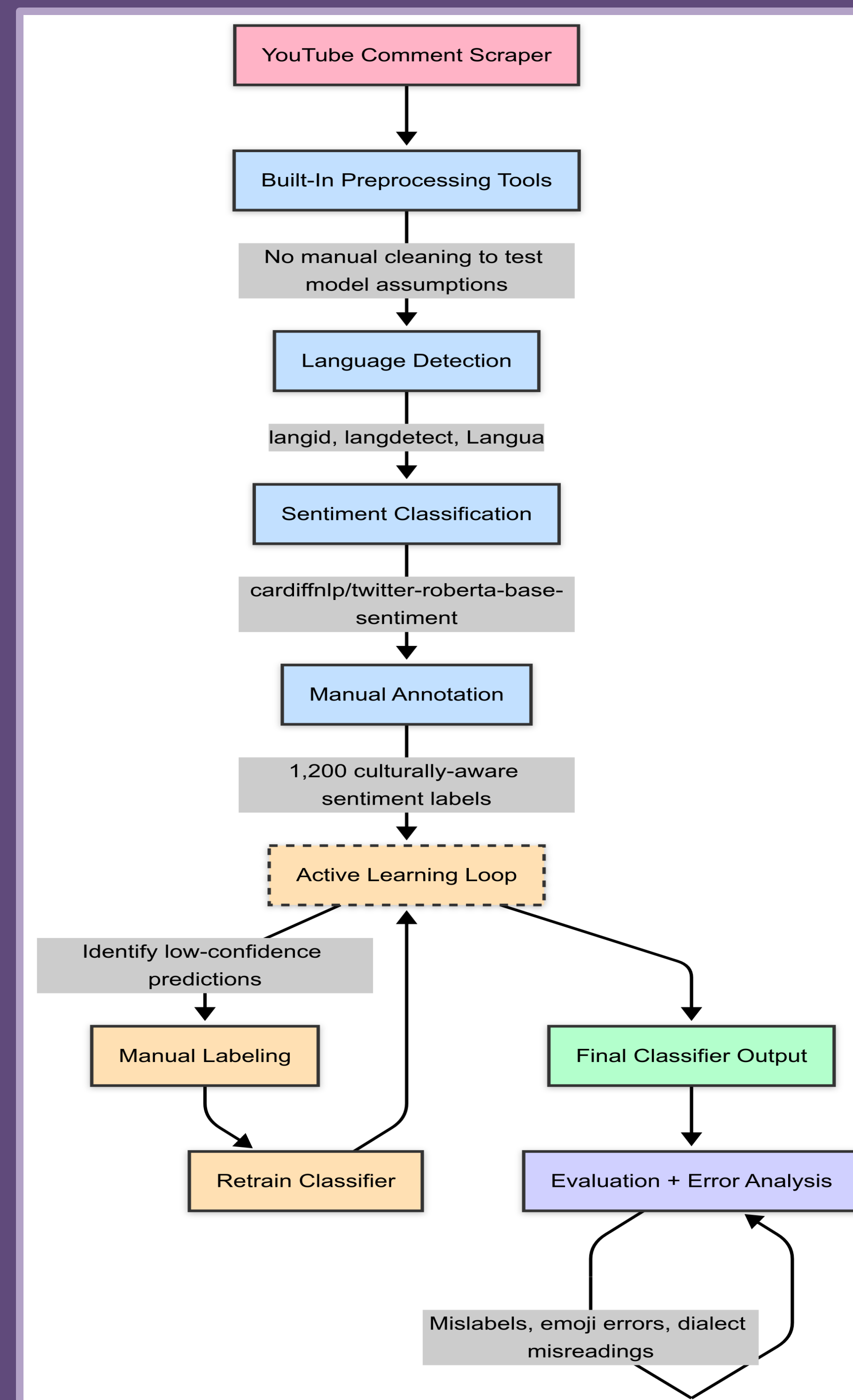
### Annotation Process

- 1,200+ manually labeled comments using culturally-informed guidelines
- Sentiment categories: Positive, Neutral, Negative, Irrelevant

### Modeling

- Used `cardiffnlp/twitter-roberta-base-sentiment` for transformer classification
- Integrated `langid`, `langdetect`, and `langua` for language detection
- Two rounds of active learning (low-confidence sampling) to improve performance

### Methodology Flow Chart



### Results – Language Detection

- The `langdetect` model mislabeled 27% of internet-language-heavy English comments as foreign
- Most errors occurred with short or emoji-filled phrases

### Results – Sentiment Classification

- Transformer model consistently misunderstood culturally coded internet expressions
- Emojis like the skull were read literally, not contextually
- Active learning improved prediction accuracy and confidence



LinkedIn



GitHub repo