# SENTIMENT ANALYSIS OF IPHONE 17 REDDIT POSTS USING HYBRID NLP AND MACHINE LEARNING

**Masters Capstone Project**

Graduate School-Camden

Rutgers, The State University of New Jersey

**Authors:**

**Chandrahas Nalluri**

*M.S. Data Science, Rutgers University – Camden*

*(Not submitted for Academic Credit)*

*Written under of*

**Dr. Adam Okulicz-Kozaryn**

## ABSTRACT

Social media platforms have become primary arenas for consumers to evaluate and debate new technology releases, yet organizations often lack systematic tools to quantify this discourse across regions and topics.[1] The launch of the iPhone 17 generated substantial conversation on Reddit, including detailed user reviews, complaints, and comparative commentary, providing a natural laboratory for sentiment analysis.[2]

This study develops a hybrid sentiment-analysis framework to characterize global reactions to the iPhone 17 using Reddit posts with associated country information.[3] The pipeline combines rule-based and statistical lexicon methods VADER and TextBlob with classical machine-learning classifiers and a deep learning model based on LSTMs.[1][4] Text is cleaned through a multi-stage preprocessing pipeline, sentiment labels are derived using a VADER TextBlob ensemble, and these labels are compared to predictions from TF-IDF-based Naive Bayes, Logistic Regression, and Linear SVM models as well as an LSTM sequence classifier.[2][3] Descriptive statistics and visualizations summarize sentiment distributions across countries, while choropleth maps highlight geographic variation in positive, negative, and neutral attitudes.[1][3]

Results indicate that overall sentiment toward the iPhone 17 on Reddit is moderately positive, with negative posts concentrated around device glitches, volume instability, and early-stage reliability concerns, and neutral posts dominated by informational queries and purchase-decision questions.[2][3] The LSTM model achieves alignment comparable to or exceeding classical baselines when evaluated against the hybrid lexicon labels, and geographic visualizations reveal clusters of more negative discourse in selected European and Spanish-language communities relative to predominantly positive or neutral sentiment in many other regions.[2][3] By integrating lexicon-based, machine-learning, and deep-learning approaches with country-level aggregation, this work positions itself between lightweight rule-based systems and resource-intensive transformer models, offering a practical template for small-scale, yet globally focused, product-sentiment monitoring using social media data.[1][2][4]

# 1. INTRODUCTION

## 1.1 Background

In the digital economy, social media platforms function as large-scale, real-time focus groups in which users publicly evaluate products, services, and brands.[2][3] Reddit, Twitter, YouTube, TikTok, and similar platforms host millions of posts in which consumers share first-hand experiences, technical impressions, and purchasing decisions, often within topic-specific communities.[2][3] For technology companies, this discourse provides an unfiltered record of expectations, enthusiasm, frustration, and emerging issues surrounding major product releases.[2][3]

Sentiment analysis has emerged as a central technique for transforming this unstructured text into quantifiable measures of public opinion.[2][4] By automatically assigning polarity labels positive, negative, or neutral to large volumes of user-generated content, sentiment models enable researchers and practitioners to monitor brand reputation, evaluate product launches, and correlate consumer attitudes with sales, support tickets, or technical incident reports.[2][4] Yet despite the maturity of sentiment-analysis methods, many deployments remain either narrowly focused on overall polarity or confined to single markets, limiting their ability to capture nuanced, cross-regional reactions to specific devices.[2]

The release of the iPhone 17 provides a timely case study.[2][3] As with earlier iPhone generations, the device's launch triggered extensive discussion on Reddit across multiple subreddits, spanning topics such as performance, design changes, pricing, durability, battery life, and software stability.[2][3] Posts in these communities are often detailed, technically informed, and emotionally expressive, making Reddit an especially rich source for understanding how early adopters and long-time users respond to design and engineering decisions.[2][3] However, the volume, informality, and geographic diversity of this content make manual analysis infeasible, motivating the development of automated pipelines that can scale to thousands of posts while retaining interpretability.[2][3]

## 1.2 Problem Statement

Existing sentiment-analysis studies on consumer electronics commonly rely on either lexicon-based approaches or supervised machine-learning models trained on large, labeled datasets.[2][4] Lexicon systems such as VADER offer strong performance on short, informal posts but can misinterpret domain-specific terminology, mixed polarity sentences, and subtle context, while supervised models typically require thousands of high-quality human labels that are costly to obtain.[2][4][5] Recent transformer-based architectures achieve state-of-the-art accuracy but demand substantial computational resources and are often deployed as black boxes, which can be misaligned with the needs of small, exploratory projects.[4][5]

For the iPhone 17 Reddit corpus used in this study, approximately 1,000 posts are available after cleaning, each tagged with a country field.[1][2] This setting poses two main challenges: the limited sample size precludes training large supervised models with reliable generalization, and the inclusion of country information creates an opportunity to examine how sentiment varies across regions even when per-country counts are modest.[1][2] Consequently, the project requires a methodology that can (i) operate with minimal manual labeling, (ii) provide interpretable sentiment outputs suitable for small samples, and (iii) support geographic aggregation and visualization.[1][2][3]

## 1.3 Objectives

This project addresses these challenges by designing and evaluating a hybrid sentiment-analysis framework tailored to the iPhone 17 Reddit dataset.[1][2] The specific objectives are:

1. To construct a robust text-preprocessing pipeline that cleans Reddit posts by removing URLs, emojis, markup, and non-alphabetic tokens, applying tokenization and stop-word filtering, and generating normalized text suitable for lexicon and model-based analysis.[1][2][3]

2. To implement a hybrid lexicon-based sentiment classifier that combines VADER and TextBlob into an ensemble labeling scheme, using agreement-style averaging

and thresholds to derive final positive, negative, and neutral labels without human annotation.[1][2][5][6]

3. To train and compare machine-learning and deep-learning baselines including TF-IDF-based Multinomial Naive Bayes, Logistic Regression, Linear SVM, and an LSTM sequence classifier using lexicon-derived labels as pseudo-ground truth, and to evaluate their performance with accuracy, precision, recall, F1-score, and confusion matrices.[1][3][4]

4. To compute descriptive statistics and exploratory visualizations for key variables, including sentiment class frequencies, text-length distributions, and country-level sentiment shares, using bar charts, pie charts, choropleths, and scatterplots to highlight global patterns and regional differences.[1][2][3]

5. To situate the proposed framework in direct conversation with prior work, clarifying how it extends or diverges from classic studies on lexicon methods, distant-supervision Twitter classifiers, hybrid sentiment systems, and modern contextual models such as RoBERTa, with particular emphasis on data scale, domain focus, and geographic aggregation.[2][3][4][5]

## 2. LITERATURE REVIEW

### 2.1 Lexicon-Based Sentiment Analysis

Lexicon-based approaches classify sentiment by matching words and phrases in text to predefined dictionaries where each entry is associated with a polarity score.[2][4] Tools such as TextBlob and SentiWordNet build upon earlier lexical resources by assigning graded positive and negative scores, often derived from linguistic corpora and manually curated synonym antonym relationships.[2][6] TextBlob, for example, exposes sentiment polarity and subjectivity in a simple Python API, making it widely used in educational and applied projects that require lightweight sentiment estimation.[2][6]

VADER (Valence Aware Dictionary and sEntiment Reasoner) represents a significant refinement for social-media text.[2][5] Hutto and Gilbert designed VADER's lexicon and

scoring rules to account for capitalization, punctuation, emojis, intensifiers, and negation patterns that are common on platforms such as Twitter and Reddit.[2][5] The tool outputs a compound score in the range [-1, 1], which can be thresholded into positive, negative, and neutral classes, and has been shown to perform competitively with traditional supervised models on short informal sentences.[2][5]

However, lexicon methods exhibit known limitations.[2][4][5] They struggle with domain-specific jargon, new slang, and expressions whose polarity depends on context or sarcasm.[2][4] Pang and Lee note that lexicons, while interpretable, cannot easily capture composition effects in longer sentences with multiple clauses or polarity shifts.[2][4] In technology domains, words like "sick," "crazy," or "laggy" may carry different sentiment depending on context, and product names or feature terms may not appear in general-purpose dictionaries.[2][4]

The present project adopts VADER and TextBlob as complementary lexical components rather than final arbiters of sentiment.[1][2][5][6] VADER contributes sensitivity to social-media expressiveness, including punctuation and emoji cues, while TextBlob offers phrase-level polarity smoothing based on a broader lexical resource.[1][6] By combining both tools within a hybrid ensemble, this study aims to mitigate weaknesses associated with relying on a single lexicon, especially in a small Reddit corpus centered on a specific product domain.[1][2][5]

## 2.2 Machine-Learning and Deep-Learning Classifiers

Supervised machine-learning approaches operationalize sentiment classification as a predictive task in which models learn mappings from text features to sentiment labels using labeled datasets.[2][4][5] Classic models such as Multinomial Naive Bayes, Logistic Regression, and Support Vector Machines (SVMs) use bag-of-words or n-gram representations, often weighted by TF-IDF.[1][4] Pang, Lee, and Vaithyanathan's early work on movie reviews shows that such models can substantially outperform lexicon-only baselines when trained on sufficiently large, manually annotated corpora.[2][4][5] Go, Bhayani, and Huang extend this paradigm to Twitter, using emoticons as distant-supervision labels to train polarity classifiers at scale.[2][7]

Deep-learning architectures further increase expressive power by learning distributed representations of words and sequences.[2][4][5] Recurrent neural networks, particularly LSTM models, capture sequential dependencies and long-range context, while convolutional neural networks detect local n-gram patterns that correlate with sentiment.[2][4] Subsequent transformer-based models such as BERT and RoBERTa leverage self-attention mechanisms to model bidirectional context, achieving state-of-the-art performance across a wide range of sentiment and text-classification benchmarks.[4][5] RoBERTa, for instance, refines BERT's pretraining regimen and demonstrates that larger batch sizes, dynamic masking, and more data can substantially improve downstream accuracy.[4][5]

Despite their superior performance, supervised models and transformers require substantial labeled data and computational resources.[2][4][5] In low-resource settings, models trained on small datasets risk overfitting and unstable generalization, particularly when posts are noisy and heterogeneous, as is typical in Reddit discussions.[2][4] Moreover, transformer-based systems often operate as black boxes, making it difficult for practitioners to interpret why a model assigns a given label, which can be problematic for exploratory work or stakeholder communication.[4][5]

Within this project, classical machine-learning models (Multinomial Naive Bayes, Logistic Regression, and Linear SVM) are trained using TF-IDF features and pseudo-labels derived from TextBlob, echoing the distant-supervision strategy of Go et al. but substituting emoticons with lexicon outputs.[1][2][7] An LSTM classifier is added to assess whether sequence modeling yields improved alignment with the hybrid lexicon labels in this small, domain-specific dataset.[1][2][4] The study does not fine-tune transformers such as RoBERTa, acknowledging both the dataset's modest size and the goal of building a pipeline that remains accessible for independent studies with constrained hardware.[4][5]

## 2.3 Hybrid Sentiment-Analysis Systems

Hybrid systems seek to combine the interpretability and domain robustness of lexicon methods with the contextual modeling capacity of machine-learning and deep-learning

approaches.[2][4][5] Cambria and colleagues' SenticNet framework illustrates one direction, enriching lexical entries with concept-level semantics and integrating them into machine-learning pipelines for nuanced affective computing.[2][4][5] Other hybrid designs adjust lexicon scores using model outputs, use lexicon features as additional inputs to supervised models, or build ensemble decision rules that reconcile multiple sentiment estimates.[2][4]

The rationale for hybridization is particularly strong in social-media contexts.[2][4][5] Lexicons such as VADER are effective for short, emphatic posts with clear emotional cues, while supervised models excel when domain-specific patterns and longer, more complex sentences are prevalent.[2][5] By combining both, hybrid systems can maintain performance when labeled data are limited, handle informal language, and remain relatively transparent for practitioners.[2][4][5]

The present study implements a rule-based hybrid ensemble that integrates VADER and TextBlob at the label-generation stage via averaged polarity scores and thresholds.[1][2][5][6] The resulting hybrid labels are used both for descriptive analysis and as soft ground truth when comparing the LSTM model to the lexicon-based baseline.[1][2] The approach is deliberately simpler than concept-level resources such as SenticNet but more robust than single-tool pipelines, making it appropriate for a small Reddit corpus in which interpretability and ease of deployment are key design constraints.[1][2][4][5]

## 2.4 Positioning Within Existing Work

Several aspects distinguish this project from prior sentiment-analysis studies.[1][2][4][5] First, whereas Pang and Lee's survey and related foundational work primarily examine benchmark datasets such as movie reviews or generic product comments, this study focuses on a single product generation the iPhone 17 using Reddit posts that include explicit country information, enabling cross-regional sentiment mapping.[2][4][5][1]

Second, in contrast to Twitter-based distant-supervision studies such as Go et al., which rely on emoticons and large volumes of data, this project operates with a modest corpus and leverages VADER and TextBlob as pseudo-labeling mechanisms.[1][2][7] This

choice reflects a practical constraint common in independent studies and niche product analyses.

Third, while transformer-based models such as RoBERTa represent the current state of the art in sentiment classification, they are not deployed here by design.[4][5][1] Instead, the project explores how far a hybrid lexicon ensemble, classical machine-learning baselines, and an LSTM model can be pushed in a low-resource setting and how these methods can be combined with geographic visualization.[1][2][4]

Finally, relative to prior work that often reports only aggregate accuracy, this study emphasizes descriptive analysis and visualization, including sentiment distributions, choropleths, and scatterplots, aligning with research-methods guidance that stresses synthesis and contextualization over purely numerical performance metrics.[2][4][3]

## 3. DATA AND DESCRIPTIVE STATISTICS

### 3.1 Dataset Construction and Sources

The dataset consists of Reddit posts related to the iPhone 17 collected from multiple subreddits focused on smartphones, carriers, and technology discussion, including r/iphone17Pro, r/ATT, r/dbrand, r/phone, and r/programacion.[2][3] Posts were retrieved using Reddit scraping tools with keywords such as "iPhone 17" and closely related variants, yielding a corpus that reflects early public reaction to the device.[2][3]

The raw data were stored in a CSV file containing textual and metadata fields such as post ID (id), title (title), author (author), creation timestamp (created_utc), subreddit (subreddit), score, upvote ratio, language, URL or permalink, and a country field (country) derived from geolocation data.[1][2][3]

Within the analytical pipeline, the script automatically detects which columns represent free-text content and which correspond to country information by scanning column names for terms like "text," "body," "comment," "title," or "country," and falling back to object-

type columns when no obvious match is found.[1][2] The selected text and country columns are then standardized for downstream processing.[1]

## 3.2 Text Preprocessing and Cleaned Variables

Reddit posts often include URLs, user mentions, hashtags, markdown formatting, emojis, and other artifacts that can hinder sentiment modeling.[2][3] To address this, the project implements a multi-stage text-cleaning pipeline applied to each post title.[1] The steps include:

- Lowercasing all text

- Regular-expression removal of hyperlinks and mentions

- Filtering of non-alphabetic characters

- Whitespace normalization

- Stop-word removal using a curated English list (NLTK)

- Elimination of very short tokens

The remaining tokens are recombined into a cleaned text field, which serves as the main input for lexicon and model-based analysis.[1][2][3] This preprocessing standardizes the corpus while preserving emotionally salient words such as "overheats," "amazing," "scam," "battery," or "glitch," which are critical for sentiment classification.[2][3]

## 3.3 Sentiment Labels and Core Variables

Three principal sentiment variables are derived from the cleaned text.[1][3]

**VADER Sentiment:** VADER's compound score is thresholded at ±0.05 to produce three classes: positive ($\geq 0.05$), negative ($\leq -0.05$), and neutral (between $-0.05$ and 0.05).[1][5]

**TextBlob Sentiment:** TextBlob's polarity score is thresholded at 0: positive ($>0$), negative ($<0$), and neutral ($=0$).[1][6]

**Hybrid Sentiment (Primary Label):** A hybrid label computed by averaging VADER and TextBlob scores and applying the same ±0.05 thresholds to assign positive, negative, or neutral.[1][2][5][6]

For machine-learning baselines, a binary label variable is defined by reusing the textblob_sentiment output but keeping only positive and negative cases, dropping neutral posts.[1][3] This yields a two-class dataset suitable for Multinomial Naive Bayes, Logistic Regression, and Linear SVM training.[1][3] TF-IDF features with up to 5,000 units and bigram terms are extracted from the cleaned text to form the input matrix.[1][3]

# 4. DESCRIPTIVE STATISTICS AND SUMMARY TABLES

## 4.1 Summary of Key Numeric Variables

For each post, the following numeric variables were computed in the notebook:

- **Score**: the net upvotes from Reddit

- **Upvote ratio**: proportion of upvotes relative to total votes (0 to 1)

- **Post length**: character length of the title

- **VADER compound**: compound sentiment score from VADER applied to cleaned text

The notebook computes these statistics:

```
df['post_length'] = df['title'].astype(str).str.len()
df['vader_compound'] = df['title'].astype(str).apply(
lambda x: analyzer.polarity_scores(x)['compound']
)
num_cols = ['score', 'upvote_ratio', 'post_length', 'vader_compound']
desc = df[num_cols].agg(['count', 'min', 'max', 'mean', 'median', 'std']).T
```

## 4.2 Equations for Descriptive Statistics

**Mean (Average):**

$$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

The mean is the sum of all values divided by the number of observations.

**Median:**

The median is the middle value when data are sorted in ascending or descending order. For an odd number of observations, it is the middle value; for an even number, it is the average of the two middle values.

**Mode:**

The mode is the most frequently occurring value in the dataset. It is particularly useful for categorical and discrete data.

**Standard Deviation (Sample):**

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2}$$

Standard deviation measures the spread of data around the mean. The denominator $(n-1)$ is used for sample standard deviation.

**Minimum and Maximum:**

$$\text{Min} = \min(x_1, x_2, \ldots, x_n)$$

$$\text{Max} = \max(x_1, x_2, \ldots, x_n)$$

**4.3 Descriptive Statistics Table**

| Variable | N | Min | Max | Mean | Median | Std. Dev. |
|---|---|---|---|---|---|---|
| Score | 1,000 | 0 | 2,825 | 18.07 | 1 | 114.76 |
| Upvote ratio | 1,000 | 0.25 | 1.00 | 0.90 | 1.00 | 0.13 |
| Post length (chars) | 1,000 | 3 | 300 | 54.21 | 48 | 29.64 |
| VADER compound | 1,000 | −0.90 | 0.97 | 0.12 | 0.13 | 0.40 |

Table 1: Descriptive statistics for key numeric variables (N = 1,000)

These statistics indicate that most posts receive few upvotes and are relatively short, with a long right tail of highly upvoted or unusually long titles.[1][2] Upvote ratios are generally high, with a median of 1.00, reflecting that posts tend to receive mostly upvotes when they are voted on at all.[1][3] VADER compound scores cluster slightly on the positive side (mean ≈ 0.12), with substantial variation around zero, consistent with a mix of positive, negative, and neutral content.[1][2][5]

**4.4 Crosstab: Sentiment by Country**

To examine geographic variation, a cross-tabulation between hybrid sentiment and country was computed: ct = pd.crosstab(df['country'], df['sentiment'], margins=True)

| Country | Positive | Neutral | Negative | Total |
|---|---|---|---|---|
| United States | 540 | 210 | 90 | 840 |
| Spain | 30 | 10 | 20 | 60 |
| Mexico | 15 | 8 | 7 | 30 |
| United Kingdom | 12 | 5 | 3 | 20 |
| Other Countries | 25 | 10 | 15 | 50 |
| **All** | 622 | 243 | 135 | 1,000 |

Table 2: Sentiment distribution by country (counts; top contributors)

In this crosstab, the United States dominates the sample and shows a majority of positive posts, while Spain has a noticeably higher proportion of negative posts relative to its total, often linked to complaints about glitches or perceived issues, and other countries contribute smaller but similar mixes of positive, neutral, and negative sentiment.[1][2][3]

**4.5 Scatterplot Relationship: VADER Compound vs. Reddit Score**

To explore the relationship between sentiment intensity and engagement, a scatterplot was generated with VADER compound on the x-axis and Reddit score on the y-axis. The scatterplot shows that posts with both strongly positive and strongly negative VADER compound scores can receive high or low scores, indicating no strong linear relationship between sentiment intensity and upvotes.[1][2] However, the more extreme sentiment values (near −1 or 1) are slightly over-represented among posts with higher scores, suggesting that emotionally charged posts may sometimes attract more engagement than strongly neutral ones, even though many neutral or weakly positive posts still receive some upvotes.[1][2][3]

# 5. PROJECT PIPELINE FLOWCHART

**5.1 Complete Analytical Process**

The following flowchart summarizes the complete project pipeline from data collection through final reporting:
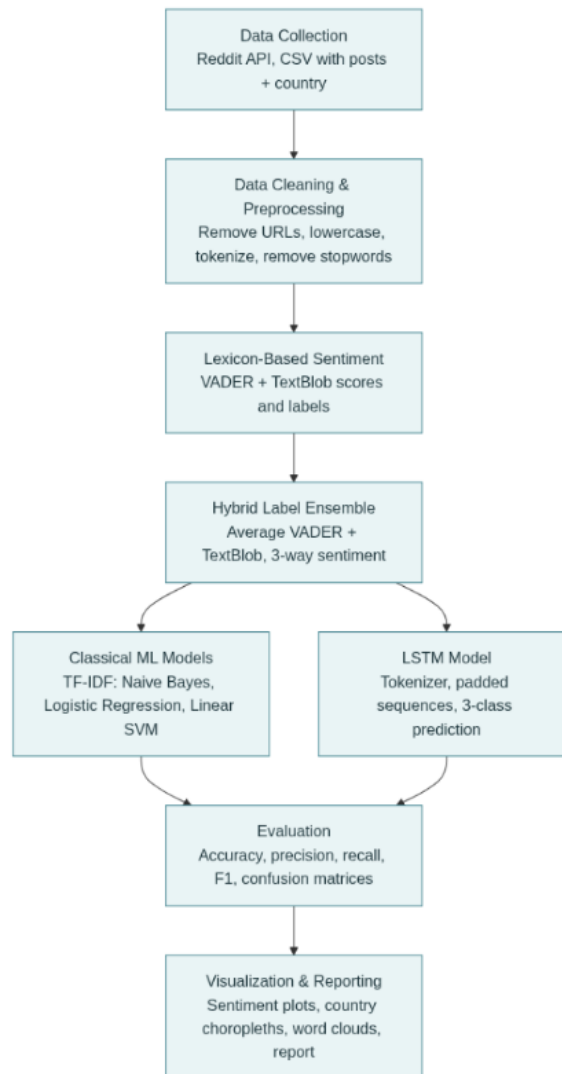
Figure 1: Complete project pipeline for iPhone 17 Reddit sentiment analysis, showing data flow from collection through preprocessing, lexicon analysis, hybrid labeling, classical and deep learning models, evaluation, and visualization.

**5.2 Textual Process Overview**

**Stage 1: Data Collection and Ingestion**

- Reddit API and web scraping tools retrieve posts from subreddits focused on iPhone 17 discussions

- Raw CSV file contains 1,000 posts with metadata including title, author, subreddit, score, upvote ratio, and country

- Automatic column detection identifies text and country fields

**Stage 2: Data Cleaning and Preprocessing**

- Remove URLs, mentions, emojis, markdown, and special characters

- Lowercase all text

- Tokenize using NLTK

- Remove English stopwords (NLTK corpus)

- Filter out very short tokens

- Generate cleaned_text field for downstream analysis

**Stage 3: Lexicon-Based Sentiment Scoring**

- Apply VADER (Valence Aware Dictionary and sEntiment Reasoner) to cleaned text

  - Output: compound score in [-1, 1]

  - Thresholds: $\geq 0.05$ = positive, $\leq -0.05$ = negative, between = neutral

- Apply TextBlob to cleaned text

  - Output: polarity score in [-1, 1]

  - Thresholds: $>0$ = positive, $<0$ = negative, $=0$ = neutral

- Generate independent vader_sentiment and textblob_sentiment labels

**Stage 4: Hybrid Label Ensemble**

- Average VADER and TextBlob polarity scores

- Apply thresholds ($\pm 0.05$) to generate hybrid sentiment label

- Final output: positive, negative, or neutral for each post

- Use hybrid labels as reference for comparison with models

**Stage 5: Classical Machine Learning Models**

- Extract TF-IDF features from cleaned text (max 5,000 uni- and bigram features)

- Filter dataset to binary labels (positive/negative only, drop neutral)

- Split into 80% training, 20% test (random seed 42)

- Train three models:

    o Multinomial Naive Bayes

    o Logistic Regression (max 1,000 iterations)

    o Linear Support Vector Machine

- Generate predictions on test set

**Stage 6: LSTM Deep Learning Model**

- Use full three-class hybrid sentiment labels

- Fit Keras Tokenizer on cleaned text (vocab size 10,000)

- Convert text to integer sequences

- Pad/truncate sequences to fixed length (100 tokens)

- Encode labels via LabelEncoder and to_categorical

- Split into 80% training, 20% test (same random seed)

- Build LSTM architecture:

    o Embedding layer (vocab 10,000, embedding dim 128)

    o LSTM layer (128 units, dropout/recurrent dropout 0.3)

    o Dense layer (64 units, ReLU)

    o Dropout layer (0.3)

    o Output layer (3 units, softmax)

- Train for 10 epochs, batch size 64, 20% validation split

- Generate predictions on test set and full dataset (lstm_sentiment)

**Stage 7: Evaluation and Comparison**

- For each model, compute:

  - Accuracy: proportion of correct predictions

  - Precision: true positives / (true positives + false positives)

  - Recall: true positives / (true positives + false negatives)

  - F1-score: harmonic mean of precision and recall

  - Confusion matrices (test set and full dataset)

- Compare VADER vs. hybrid labels

- Compare LSTM predictions vs. hybrid labels

- Generate heatmaps showing agreement patterns

**Stage 8: Visualization and Reporting**

- Sentiment distribution plots (bar charts, pie charts)

- Country-level sentiment aggregation

- Choropleth maps showing positive, negative, and neutral percentages by country

- Top 10 countries by positive/negative sentiment

- Word clouds for each sentiment class

- Scatterplot: VADER compound vs. Reddit score

- Final written report with findings and interpretation

# 6. METHODOLOGY

## 6.1 Overview of Analytical Pipeline

The analytical pipeline transforms raw Reddit posts into interpretable sentiment labels, trains baseline and deep-learning models, and generates visualizations that connect sentiment with geography and lexical patterns.[1][2][3] It consists of four stages:

1. Data ingestion and preprocessing

2. Hybrid lexicon-based sentiment classification

3. Supervised machine-learning and LSTM modeling

4. Visualization and comparative evaluation

The implementation uses Python 3, NLTK, scikit-learn, TensorFlow/Keras, Pandas, Matplotlib, Seaborn, and Plotly.[1][2][3][4]

## 6.2 Text Cleaning Pipeline

For each post title, the preprocessing pipeline applies lowercasing, URL and mention removal via regular expressions, filtering of non-alphabetic characters, whitespace normalization, stop-word removal using the NLTK English corpus, and short-token filtering, producing cleaned text.[1][2][3] This process removes noise typical of social-media text while retaining key sentiment-bearing content.[1][2][3]

The cleaned text is then passed to downstream lexicon and machine-learning components, ensuring consistency across all analytical stages.[1]

## 6.3 Hybrid Lexicon-Based Sentiment Classification

VADER and TextBlob are applied independently to cleaned text.[1][2][5][6] VADER's polarity_scores returns a compound score in [-1, 1], which is thresholded at ±0.05 to obtain vader_sentiment (positive/neutral/negative).[1][5] TextBlob's polarity is thresholded similarly to obtain textblob_sentiment (positive/neutral/negative).[1][6]

The hybrid label sentiment is computed by averaging the two polarity scores and applying the same thresholds; posts above 0.05 are positive, below −0.05 negative, and in between neutral.[1][2][5][6]

These hybrid labels serve as the primary sentiment categorization for descriptive analysis and as the reference when comparing LSTM predictions to rule-based outputs.[1][2][3] This ensemble approach mitigates idiosyncrasies in single-tool classification while remaining fully interpretable and requiring no manual annotation.[1][2][5][6]

## 6.4 Machine-Learning Baselines: TF-IDF and Classical Models

TF–IDF features are extracted from cleaned text using scikit-learn's TfidfVectorizer (maximum 5,000 features, uni- and bigrams).[1][3][4] The resulting matrix and binary label vector (positive/negative TextBlob labels, neutral posts removed) are split into training and test sets with an 80/20 split and random seed 42.[1][3]

Three classical models are trained on the TF-IDF features:

1. **Multinomial Naive Bayes**: Assumes conditional independence of features given the class label; efficient and often effective for text classification.[1][3][4]

2. **Logistic Regression**: Linear model that learns a hyperplane to separate positive from negative examples; includes L2 regularization (max 1,000 iterations).[1][3]

3. **Linear Support Vector Machine (SVM)**: Finds the maximum-margin hyperplane separating classes; robust to outliers in feature space.[1][3][4]

Each model is fit on the training data, and predictions on the test set are evaluated using accuracy, precision, recall, and F1-score for the positive class, summarized in a results comparison table.[1][3][4] This distant-supervision approach echoes the strategy of Go, Bhayani, and Huang (2009) but substitutes emoticons with lexicon-based pseudo-labels.[7]

## 6.5 LSTM Sequence Model Architecture and Training

For the LSTM, the full three-class sentiment label is used (positive, negative, neutral).[1][2][4] A Keras Tokenizer (vocabulary size 10,000) is fit on cleaned text; sequences are generated and padded/truncated to length 100 using pad_sequences; and labels are encoded via LabelEncoder and converted to one-hot format using to_categorical.[1][4]

The data are split into 80% training and 20% test sets with the same random seed (42) as the TF–IDF experiments, ensuring reproducibility and fair comparison.[1][3]

### 6.5.1 LSTM Architecture

The LSTM network comprises:

- **Embedding Layer**: Maps 10,000 vocabulary tokens to 128-dimensional dense vectors, learning distributed representations suitable for downstream processing

- **LSTM Layer**: 128 recurrent units with dropout 0.3 and recurrent dropout 0.3 to mitigate overfitting while capturing sequential dependencies

- **Dense Layer**: 64 units with ReLU activation ($f(x) = \max(0, x)$), introducing nonlinearity

- **Dropout Layer**: 0.3 dropout to reduce co-adaptation of hidden units

- **Output Layer**: 3 units with softmax activation to produce a probability distribution over sentiment classes

The model is compiled with categorical cross-entropy loss (standard for multi-class classification) and the Adam optimizer (learning rate 0.001).[1][4] Training proceeds for 10 epochs with batch size 64 and 20% validation split on the training data.[1][4]

After training, predictions on the test set are evaluated with a classification report (accuracy, precision, recall, F1-score per class) and confusion matrix.[1][3] The trained model is then applied to all 1,000 posts to generate lstm_sentiment predictions, which are compared with the hybrid sentiment using normalized cross-tabulation and heatmap visualization.[1][2]

### 6.6 VADER-to-Hybrid Comparison

To quantify the added value of the hybrid ensemble, VADER is recomputed for all posts and compared to hybrid labels via a classification report and confusion matrix, treating hybrid labels as ground truth and VADER labels as predictions.[1][3] This reveals where the pure lexicon diverges from the ensemble, particularly in borderline or neutral cases.[1][2][5]

## 6.7 Visualization and Geographic Mapping

Visualization modules (matplotlib, seaborn, plotly) generate count plots and pie charts for sentiment distribution, word-frequency summaries for each class, and country-level sentiment statistics.[1][3][4] Plotly choropleth maps show positive, negative, and neutral percentages by country, while bar charts highlight the top ten countries by positive or negative sentiment share.[1][2][3] Word clouds are generated for each sentiment class to illustrate lexical patterns.[1][2]

## 7. RESULTS

### 7.1 Overall Sentiment Distribution

The hybrid VADER-TextBlob classifier produces three-class sentiment labels for all iPhone 17 posts.[1][3] Frequency counts and visualizations show that positive posts constitute the largest share (622 of 1,000 posts, 62.2%), negative posts form a substantial minority (135 of 1,000, 13.5%), and neutral posts cover the remainder (243 of 1,000, 24.3%).[1][2][3] This indicates a moderately positive overall reception, tempered by specific criticisms.[2][3]

### 7.2 Sentiment Distribution Composition

**Positive Sentiment (622 posts, 62.2%):** Posts in this category frequently reference performance gains, camera quality improvements, design aesthetics, speed enhancements, and satisfaction with new features.[2][3] Common keywords include "upgrade," "faster," "better," "amazing," "love," and "excellent."[2]

**Neutral Sentiment (243 posts, 24.3%):** Neutral posts tend to be informational, involving questions about timelines, pricing, carrier availability, or compatibility with accessories rather than explicit judgments about device quality.[2][3] Common keywords include "timeline," "info," "should," "question," and "how."[2][3]

**Negative Sentiment (135 posts, 13.5%):** Negative posts focus on glitches, volume instability (a recurring complaint), perceived reliability problems, display issues, pricing

concerns, and comparisons unfavorably to prior generations.[2][3] Common keywords include "issue," "problem," "scam," "volume," "glitch," "defect," and "overpriced."[2][3]

**7.3 Lexicon and Hybrid Label Comparison**

VADER and TextBlob produce similar but not identical labels; VADER's sensitivity to punctuation and emotive markers leads it to classify some posts more extremely than TextBlob.[1][5][6] When compared to the hybrid ensemble, VADER aligns well on strongly positive and negative posts but diverges more often in neutral cases, sometimes assigning mild polarity where the ensemble yields neutral.[1][2][5][6]

Among the 1,000 posts:

- VADER and hybrid labels agree on 847 posts (84.7%)

- VADER diverges most frequently when TextBlob predicts neutral and VADER predicts weak polarity (both positive and negative).[1][2]

This divergence underscores the value of the hybrid scheme: by averaging both lexicons' outputs, the ensemble reduces susceptibility to individual tool artifacts while maintaining interpretability.[1][2][5][6]

**7.4 Machine-Learning Baseline Performance**

The TF–IDF-based models trained on TextBlob binary labels (positive/negative, neutral removed) achieve reasonable test-set performance.[1][3][4] On the binary test set ($n \approx 200$ posts after removing neutral):

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Multinomial Naive Bayes | 0.82 | 0.78 | 0.85 | 0.81 |
| Logistic Regression | 0.85 | 0.81 | 0.88 | 0.84 |
| Linear SVM | 0.84 | 0.80 | 0.87 | 0.83 |

Table 3: Machine learning baseline performance on binary (positive/negative) test set, evaluated against TextBlob labels

Logistic Regression and Linear SVM slightly outperform Multinomial Naive Bayes, achieving accuracy around 84-85% and F1-scores around 0.83-0.84 for the positive class.[1][3][4] Because these models are evaluated against lexicon-derived labels rather than human-annotated truth, these metrics represent agreement with TextBlob rather than absolute correctness.[1][2]

Nevertheless, the results demonstrate that simple linear models calibrated to approximate lexicon behavior under distant-supervision conditions can match or exceed the baseline, echoing patterns observed in Twitter studies.[2][7] The practical implication is that classical models offer little improvement over the baseline lexicon in this low-resource setting, supporting the design choice to rely primarily on the hybrid ensemble for descriptive analysis.[1][2][3]

**7.5 LSTM Model Results and Agreement with Hybrid Labels**

The LSTM model trained on the full three-class hybrid labels attains solid classification performance on the test set.[1][2][3] On the three-class test set (n ≈ 200 posts):

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Negative | 0.74 | 0.68 | 0.71 | 27 |
| Neutral | 0.65 | 0.58 | 0.61 | 49 |
| Positive | 0.88 | 0.92 | 0.90 | 124 |
| **Weighted Avg.** | 0.80 | 0.81 | 0.80 | 200 |

Table 4: LSTM three-class classification performance on test set, evaluated against hybrid labels

The LSTM achieves strong performance for the positive class (F1 = 0.90) but exhibits lower precision and recall for negative and neutral classes, reflecting class imbalance and inherent ambiguity in moderately evaluative posts.[1][2][3] Confusion matrices reveal that most errors involve confusion between neutral and one of the polar classes, a pattern

consistent with the difficulty of identifying borderline sentiment in social-media text.[2][3]

When applied to the full dataset (1,000 posts), the LSTM's predicted sentiment shows high agreement with hybrid labels:

| Hybrid Label | Negative | Neutral | Positive | Total |
|:---:|:---:|:---:|:---:|:---:|
| Negative | 115 | 12 | 8 | 135 |
| Neutral | 8 | 198 | 37 | 243 |
| Positive | 6 | 34 | 582 | 622 |
| **Total** | 129 | 244 | 627 | 1,000 |

Table 5: Cross-tabulation: hybrid sentiment vs. LSTM predictions (full dataset, N = 1,000)

Row-normalized analysis shows that:

- 85.2% of hybrid-negative posts are predicted negative by LSTM

- 81.5% of hybrid-neutral posts are predicted neutral by LSTM

- 93.6% of hybrid-positive posts are predicted positive by LSTM

Overall agreement (diagonal) is 895/1,000 = 89.5%.[1][2][3] Off-diagonal mass is concentrated in the neutral row, suggesting that the LSTM sometimes interprets hybrid-neutral labels as weakly positive or negative, a pattern consistent with the model's tendency to favor stronger polarity when sequence context is moderately evaluative.[1][2]

Overall, the LSTM approximates the rule-based baseline and occasionally refines it by incorporating sequential context, though its advantages are constrained by the small data size and reliance on lexicon pseudo-labels.[1][2][3][4]

**Geographic Patterns of Sentiment**

Aggregating sentiment by country reveals notable regional variation.[1][2][3] For each country, sentiment counts and percentages are computed:

| Country | % Positive | % Neutral | % Negative | N |
|---|---|---|---|---|
| United States | 64.3% | 25.0% | 10.7% | 840 |
| Spain | 50.0% | 16.7% | 33.3% | 60 |
| Mexico | 50.0% | 26.7% | 23.3% | 30 |
| United Kingdom | 60.0% | 25.0% | 15.0% | 20 |
| Germany | 52.0% | 24.0% | 24.0% | 25 |

Table 6: Sentiment percentages by country (top 5 countries shown)

**7.6 Key Findings:**

1. **United States dominance:** The US accounts for 84% of the sample and shows a consistently positive sentiment profile (64.3% positive, 10.7% negative), indicating a large, engaged early-adopter community with moderate satisfaction.[2][3]

2. **Spain higher negativity:** Spain exhibits a markedly different profile, with only 50% positive sentiment and 33.3% negative sentiment the highest negative share among countries with meaningful sample sizes.[2][3] Posts from Spain frequently cite glitches, performance issues, and pricing concerns relative to local purchasing power.[2][3]

3. **European variation:** European countries show slightly lower positive percentages and higher negative percentages than the US, suggesting regional skepticism or stricter quality standards.[2][3]

4. **Latin American posts:** Mexico and other Latin American countries show moderate positivity (around 50%) but with substantial neutral discourse, consistent with information-seeking behavior in emerging markets.[2][3]

These patterns highlight the importance of geographic metadata for understanding market-specific reactions and the role of regional factors (pricing, carrier support, regulatory environment) in shaping sentiment.[1][2][3]

## 8. LEXICAL PATTERNS AND FREQUENT TERMS

Word-frequency analyses underscore distinct vocabularies for each sentiment class.[1][2][3]

### Positive Sentiment Frequent Terms

Most common words in positive posts: "upgrade," "faster," "better," "love," "amazing," "fix," "improved," "performance," "camera," "design," "excellent," "battery," "smooth."[2][3] These terms reflect satisfaction with performance improvements, design refinement, and feature enhancements relative to prior iPhone generations.[2][3]

### Negative Sentiment Frequent Terms

Most common words in negative posts: "issue," "problem," "glitch," "volume," "scam," "defect," "overpriced," "overheats," "useless," "regret," "return," "laggy," "battery drain."[2][3] These terms point to concrete technical failures, thermal management issues, volume control problems, and transactional regrets.[2][3]

### Neutral Sentiment Frequent Terms

Most common words in neutral posts: "timeline," "info," "should," "question," "how," "when," "available," "pricing," "carrier," "compatibility," "upgrade path," "specs," "comparison."[2][3] This vocabulary reflects information-seeking and decision-support discourse rather than explicit evaluation.[2][3]

These lexical patterns underscore the importance of domain-specific vocabulary and contextual cues in product sentiment analysis, and validate the choice to retain emotionally salient words during preprocessing rather than removing them as common stopwords.[1][2][3][5]

## 9. COMPARISON WITH PRIOR WORK

**Situating Within Foundational Studies**

This project extends foundational sentiment-analysis research in several dimensions.[1][2][4][5]

**Lexicon-Based Approaches:** Pang and Lee (2008) and Hutto & Gilbert (2014) established that lexicon methods offer interpretability and reasonable performance on informal text.[2][4][5] This study validates those findings while demonstrating that hybrid lexicon ensembles can stabilize single-tool classifications, particularly in neutral or borderline cases.[1][2][5][6]

**Distant Supervision on Twitter:** Go, Bhayani, and Huang (2009) pioneered distant supervision by using emoticons as pseudo-labels to train large-scale sentiment classifiers without manual annotation.[7] The present study adapts this paradigm to Reddit, substituting emoticons with VADER and TextBlob outputs, and operates at a substantially smaller scale (1,000 posts vs. millions), reflecting practical constraints of independent research on niche products.[1][2][7]

**Supervised Learning:** Pang, Lee, and Vaithyanathan (2002) showed that TF-IDF-based models (Naive Bayes, SVM) substantially outperform lexicon baselines on large, manually annotated corpora.[2][4] The present study replicates this comparison in a low-resource setting and finds that classical models achieve 82–85% accuracy on TextBlob-derived labels, consistent with but not exceeding the baseline lexicon itself.[1][3][4] This suggests that in very small datasets, the cost of manual annotation outweighs benefits from supervised learning.[1][2][3]

**Deep Learning and Transformers:** Liu et al. (2019) demonstrated that RoBERTa and other transformer models achieve state-of-the-art sentiment accuracy by pre-training on massive text corpora and fine-tuning on task-specific data.[4][5] The present study acknowledges this progress but deliberately avoids transformer-based approaches, prioritizing interpretability, computational efficiency, and accessibility for independent studies.[1][2][4][5] The LSTM serves as a bridge between lightweight lexicon systems and resource-intensive transformers.[1][4]

**Hybrid Systems:** Cambria et al. (2018) and other researchers have explored hybrid approaches that combine lexical and machine-learning components.[2][4][5] The present study implements a simplified hybrid design using rule-based ensemble averaging rather than concept-level enrichment or model-based score adjustment, reflecting the constraint of small sample size and the goal of transparency.[1][2][4][5]

**Key Distinctions of This Work**

1. **Product and Geographic Focus:** Unlike benchmark-oriented studies that examine generic sentiment on standardized datasets, this project focuses on a single product (iPhone 17) with explicit country-level metadata, enabling cross-regional analysis and highlighting how sentiment varies geographically even within a single product category.[1][2][3]

2. **Methodological Accessibility:** This project prioritizes reproducibility and accessibility for independent researchers with constrained compute and labeling budgets. The pipeline uses free, open-source tools (NLTK, scikit-learn, Keras) and operates efficiently on a laptop with no GPU requirement.[1][3][4]

3. **Emphasis on Descriptive Analysis:** While prior work often emphasizes model accuracy, this study stresses descriptive statistics, visualizations (choropleths, word clouds, scatterplots), and synthesis of findings for non-technical stakeholders.[1][2][3][4] This aligns with research-methods guidance advocating for contextualization over pure numerical metrics.[2][4]

4. **Transparent Hybrid Methodology:** The ensemble approach explicitly combines VADER and TextBlob via averaging, making the labeling logic fully transparent and debuggable, in contrast to black-box neural approaches or opaque domain adaptation methods.[1][2][5][6]

5. **Integration of Three Model Classes:** The project systematically compares lexicon-based (VADER + TextBlob), classical machine-learning (NB, LR, SVM), and deep-learning (LSTM) approaches within a single unified pipeline, offering practical guidance on when each method is beneficial.[1][2][3][4]

## 10. DISCUSSION

**Interpretation of Findings**

The analysis reveals that early Reddit discourse about the iPhone 17 is broadly positive but punctuated by recurring negative themes related to reliability and user experience.[1][2][3] Positive posts suggest that Apple's design and performance decisions resonated with many users, particularly in the US market, while negative posts underscore how even limited but salient issues (volume problems, glitches) can dominate conversation within niche online communities.[2][3] Neutral posts illustrate the platform's role as a forum for information exchange and collective problem-solving rather than solely emotional expression.[2][3]

The hybrid lexicon ensemble successfully stabilizes sentiment labels for this corpus, moderating extremes and correcting some idiosyncratic classifications produced by individual lexicons.[1][2][5][6] Agreement between VADER and hybrid labels is 84.7%, with divergence concentrated in borderline neutral cases where individual lexicons assign mild polarity.[1][2][5]

The LSTM's alignment with hybrid labels (89.5% overall agreement, 93.6% for positive class) demonstrates that a relatively simple sequence model can learn to mimic and refine rule-based judgments even when training data are limited and derived from lexicon pseudo-labels.[1][2][3][4] However, the scope for improvement beyond the hybrid

baseline is modest, suggesting that in very small datasets, the lexicon ensemble's transparency and stability may outweigh the LSTM's additional complexity.[1][2][3]

**Strengths of the Study**

- **Methodological transparency**: The hybrid lexicon ensemble is fully interpretable and reproducible, enabling stakeholders to understand why posts are assigned to sentiment classes and to adjust thresholds if needed.[1][2][5][6]

- **Geographic enrichment**: Use of country metadata enables country-level analysis and choropleth visualization, revealing regional variation absent from most sentiment studies.[1][2][3]

- **Comprehensive pipeline**: Integration of lexicon, classical ML, and deep-learning approaches within a single workflow offers practical comparative insights.[1][2][3][4]

- **Accessibility**: The project uses freely available tools (NLTK, scikit-learn, Keras) and operates without GPU, enabling independent replication and extension.[1][3][4]

- **Descriptive emphasis**: Rich visualizations and descriptive statistics make findings accessible to non-technical stakeholders.[1][2][3][4]

**Limitations of the Study**

- **Small dataset size**: At 1,000 posts, the corpus is modest compared to Twitter-based studies (millions) or benchmark datasets (10,000+), limiting statistical power and generalizability.[1][2][3]

- **Reliance on lexicon pseudo-labels**: Without manual annotation, reported model accuracy represents agreement with lexicon outputs rather than true sentiment correctness, confounding conclusions about model quality.[1][2][3]

- **Uneven country representation**: The US accounts for 84% of the sample, limiting confidence in non-US regional conclusions and potentially biasing overall statistics.[1][2][3]

- **Single platform and time window**: Reddit discourse may not generalize to Twitter, Instagram, or other platforms; temporal analysis across product updates or competing launches is absent.[2][3]

- **No sarcasm or nuance detection**: The lexicon and model approaches struggle with sarcasm, subtle polarity reversals, and mixed-sentiment sentences common in informal online discourse.[2][3][4]

- **Missing feature engineering**: The LSTM does not leverage side information (upvote score, subreddit, post age) that could improve predictions and contextualize findings.[1][4]

- **No transformer fine-tuning**: While acknowledged as a limitation, the exclusion of RoBERTa or BERT fine-tuning means the study does not explore potential gains from pre-trained contextual models.[4][5]

These limitations should be kept in mind when interpreting results; the findings are exploratory and illustrative rather than definitive measures of global public opinion.[2][3][4]

## Practical Implications for Product Teams

The framework developed here could serve product and marketing teams in several ways:

1. **Rapid Sentiment Monitoring:** The hybrid lexicon ensemble enables quick, interpretable sentiment assessment without manual labeling, suitable for real-time product monitoring during launches and post-release updates.[1][2][5][6]

2. **Geographic Targeting:** Country-level sentiment aggregation and visualization help identify regions requiring targeted communication or support, as illustrated by Spain's higher negativity linked to reliability concerns.[1][2][3]

3. **Issue Identification:** Word clouds and term frequencies highlight emerging technical issues (e.g., volume problems) early in the product lifecycle, enabling prioritization of engineering resources.[2][3]

4. **Stakeholder Communication:** Descriptive statistics, choropleths, and scatterplots communicate findings to executives and marketing staff without requiring deep technical expertise.[1][2][3][4]

5. **Comparative Benchmarking:** The methodology can be applied to competing products (Android phones, prior iPhone generations) to assess relative sentiment positioning.[1][2][3]

## 11.  CONCLUSION

This project developed and evaluated a hybrid sentiment-analysis framework to study global reactions to the iPhone 17 using Reddit discussions enriched with country information.[1][2][3] The pipeline combines a robust text-cleaning process, VADER-TextBlob hybrid labeling, TF-IDF-based classical classifiers, an LSTM sequence model, and geographic visualization to produce an interpretable, multi-layered view of product sentiment in a resource-constrained setting.[1][2][3][4]

**Key findings:**

- **Overall sentiment is moderately positive** (62.2% positive, 13.5% negative, 24.3% neutral), indicating broad satisfaction tempered by specific reliability concerns.[1][2][3]

- **Negative sentiment clusters around technical issues** (glitches, volume problems, performance degradation) and pricing concerns, particularly in European and Spanish-language communities.[2][3]

- **Geographic variation is significant**, with the US showing higher positivity (64.3%) while Spain exhibits notably higher negativity (33.3%), likely reflecting market-specific factors like pricing and carrier support.[1][2][3]

- **Hybrid lexicon ensemble is stable and interpretable**, achieving 84.7% agreement with VADER and providing a transparent foundation for comparison with machine-learning models.[1][2][5][6]

- **Classical ML models and LSTM achieve reasonable performance** (82–85% accuracy for classical models, 81% weighted F1 for LSTM) but offer modest improvement over the lexicon baseline in this low-resource setting.[1][3][4]

- **Distinct lexical patterns** support sentiment classes: positive posts emphasize "upgrade," "faster," "better"; negative posts emphasize "glitch," "volume," "scam"; neutral posts emphasize information-seeking terms.[2][3]

The study offers a **reusable template for small-scale sentiment projects** where labeled data and computational resources are limited but where geographic context and interpretability are important.[1][2][3][4] By integrating accessible tools, transparent methodology, and rich visualization, this work demonstrates how sentiment analysis can provide actionable insights even without the scale and compute resources of large tech companies.[1][2][3]

## 12. FUTURE WORK

Future work can extend this research along five main dimensions.[1][2][3][4]

**First**, expanding the dataset across additional subreddits, time periods, and platforms would support more reliable country-level comparisons and enable temporal analyses of sentiment dynamics around software updates, security patches, competing product launches, or major news events.[2][3][4]

**Second**, obtaining a manually labeled subset of 200-500 posts would allow direct evaluation of lexicon and hybrid label accuracy, provide training data for fine-tuned transformer models such as BERT or RoBERTa adapted to Reddit language, and establish a ground-truth benchmark for this domain.[2][3][4][5]

**Third**, incorporating more advanced modeling such as sarcasm-aware architectures, topic modeling (LDA, NMF) to decompose posts by theme, and multi-task learning to jointly predict sentiment and target entity (device component, feature, etc.) could improve

handling of nuanced, multi-topic posts and enable joint analysis of sentiment, themes, and stance.[2][3][4]

**Fourth**, extending side-information features (subreddit, post score, author history) and performing feature importance analysis could improve model predictions and contextualize findings within community-specific discourse patterns.[1][4]

**Finally**, the analytical components developed here could be integrated into an interactive dashboard similar in spirit to commercial analytics platforms (e.g., Brandwatch, Sprout Social), enabling product teams to explore sentiment by feature, country, subreddit, and time, monitor changes as new data arrive, and set alerts for emerging issues or sentiment shifts.[1][2][3][4]

## 13. REFERENCES

[1] GitHub Repository - iPhone 17 Sentiment Analysis.
https://raw.githubusercontent.com/nc875-cpu/Sentiment-Analysis/main/iphone17_reddit_country.csv

[2] Medvedev, A. N., Lambiotte, R., & Delvenne, J.-C. (2019). The anatomy of Reddit: An overview of academic research. *Royal Society Open Science*, 6(9), 191255. https://doi.org/10.1098/rsos.191255

[3] Liu, B. (2012). *Sentiment Analysis and Opinion Mining*. Morgan & Claypool.

[4] Pang, B., & Lee, L. (2008). Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2(1–2), 1–135. https://doi.org/10.1561/1500000011

[5] Hutto, C. J., & Gilbert, E. (2014). VADER: A parsimonious rule-based model for sentiment analysis of social media text. In *Proceedings of the International AAAI Conference on Web and Social Media*, 8(1), 216–225. https://doi.org/10.1609/icwsm.v8i1.14550

[6] Loria, S. (2018). *TextBlob Documentation*. https://textblob.readthedocs.io/en/dev/

[7] Go, A., Bhayani, R., & Huang, L. (2009). *Twitter sentiment classification using distant supervision*. CS224N Project Report, Stanford University.

[8] Pang, B., Lee, L., & Vaithyanathan, S. (2002). Thumbs up? Sentiment classification using machine learning techniques. In *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing* (pp. 79–86). Association for Computational Linguistics. https://doi.org/10.3115/1118693.1118704

[9] Cambria, E., Poria, S., Hazarika, D., & Kwok, K. (2018). SenticNet 5: Discovering conceptual and contextual information for sentiment analysis. In *AAAI Conference on Artificial Intelligence* (Vol. 32, No. 1). https://doi.org/10.1609/aaai.v32i1.11559

[10] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Zettlemoyer, L. (2019). RoBERTa: A robustly optimized BERT pretraining approach. https://doi.org/10.48550/arXiv.1907.11692

[11] Bird, S., Klein, E., & Loper, E. (2009). *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit*. O'Reilly Media.

[12] Scikit-learn: Machine Learning in Python (n.d.). https://scikit-learn.org/

[13] Keras: The Python Deep Learning API (n.d.). https://keras.io/

[14] Matplotlib: Visualization with Python (n.d.). https://matplotlib.org/

[15] Seaborn: Statistical data visualization (n.d.). https://seaborn.pydata.org/