# Gendered Abuse Detection in Indic Languages

**Kumar Mrinal**
mrinal22258@iiitd.ac.in

**Piyush Gautam**
piyush21549@iiitd.ac.in

## Abstract

The pervasive problem of online gender-based harassment continues to silence women and marginalized communities across digital platforms, creating hostile environments that restrict free expression. This persistent challenge demands robust technological solutions capable of identifying abusive content across diverse linguistic contexts. Our work tackles this critical issue by developing models trained on annotated datasets in English, Hindi and Tamil. We implement two neural architectures: CNN-BiLSTM with GloVe embeddings serves as our primary baseline (Baseline 1) for Task1, Task2 and Task3, effectively capturing both local patterns and sequential dependencies in textual data. For Task1 and Task3, we additionally deploy GRU-Attention with GloVe embeddings as Baseline 2 to model contextual relationships, while employing Logistic Regression for Task2 (Baseline 2) to provide interpretable classification. Our Code is at, visit Github.

## 1 Introduction

Sexism persists in online communication just as in offline spaces, manifesting through harassment, cyberbullying, and gendered discourse. This problem particularly affects women and marginalized groups in digital platforms, creating hostile environments that limit participation. The challenge is especially acute for Indic language users, where automated detection systems face additional complexities due to linguistic nuances and code-switching patterns.

For this study, we utilize the ICON23 shared task dataset containing annotated Twitter posts in English (7,638 samples), Hindi (7,714 samples), and Tamil (7,914 samples). We implement two baseline approaches: (1) CNN-BiLSTM with GloVe embeddings for all three tasks (Task1: gendered abuse, Task2: targeted attacks, Task3: explicit abuse), and (2) GRU-Attention for Task1/Task3 along with Logistic Regression for Task2. These models address the critical need for effective detection systems while working within the constraints of available Indic language resources.

## 2 Related Work

Research on detecting online gender-based harassment has evolved from traditional machine learning approaches to more sophisticated deep learning techniques. Early methods utilized algorithms like Support Vector Machines (SVMs) and Random Forests for classifying abusive content, though these faced challenges in handling diverse linguistic structures and cultural contexts.

Recent advances employ deep neural networks, with Convolutional Neural Networks (CNNs) proving effective at extracting local textual patterns and Recurrent Neural Networks (RNNs) capturing sequential dependencies. Particularly, hybrid architectures combining CNNs with Bidirectional LSTMs (BiLSTMs) have shown promise by jointly modeling spatial and temporal features in abusive language.

For multilingual detection, approaches leveraging transfer learning from English to Indic languages have been explored to address data scarcity. However, significant challenges remain in adapting models to handle code-switching and cultural nuances in low-resource languages like Hindi and Tamil. The current study builds on these developments while focusing specifically on the ICON23 annotated datasets for English, Hindi, and Tamil.

## 3 Dataset

### 3.1 ICON23 Shared Task Dataset

The **DATASET** provides annotated Twitter posts for gender-based abuse detection across three languages: Indian English, Hindi, and Tamil. The dataset was carefully curated with input from 18 domain experts in gender studies, ensuring high-quality annotations.

## 3.2 Data Collection and Annotation

The dataset comprises:

- 7,638 Indian English posts
- 7,714 Hindi posts
- 7,914 Tamil posts

Each post was annotated for three distinct labels:

1. **Label 1 (Gendered Abuse)**: Identifies abuse targeting non-marginalized genders
2. **Label 2 (Targeted Abuse)**: Detects abuse against gender/sexual minorities
3. **Label 3 (Explicit Hostility)**: Flags overtly hostile content

## 3.3 Annotation Schema

Annotations were collected using the following scheme:

- '1': Agreement with label
- '0': Disagreement with label
- 'NL': Post assigned but not annotated
- 'NaN': Post not assigned to annotator

## 3.4 Dataset Structure

The dataset includes the following key columns:

- `id`: Unique sequence identifier
- `text`: Tweet content
- `language`: Language identifier (English/Hindi/Tamil)
- `key`: Label identifier (question_1–3)
- Annotation columns:
  - `en_a1–en_a6`: English annotators
  - `hi_a1–hi_a5`: Hindi annotators
  - `ta_a1–ta_a6`: Tamil annotators

Table 1: Dataset Statistics by Language

| Language | Post Count |
|---|---|
| Indian English | 7,638 |
| Hindi | 7,714 |
| Tamil | 7,914 |

The dataset's multi-label structure and comprehensive annotation scheme make it particularly valuable for developing robust models for gender-based abuse detection in Indic languages.

## 4 Task Description

This paper addresses the ICON 2023 Shared Task on "Gendered Abuse Detection in Indic Languages." The task focuses on developing computational approaches to identify gender-based online violence, which has significant societal and economic impacts. We explore three distinct subtasks from this challenge:

### 4.1 Subtask 1: Gendered Abuse Classification (Label 1)

This subtask involves developing a classifier for detecting gender-based abuse from the provided dataset. The dataset contains annotations from eighteen domain experts including activists and researchers, offering a solid foundation for building systems capable of identifying gender-based cyber violence across multiple languages.

### 4.2 Subtask 2: Transfer Learning for Abuse Detection

This subtask requires implementing transfer learning approaches by leveraging external open datasets on hate speech and toxic language recognition in Indic languages. We explore how models can benefit from broader patterns and knowledge sources to enhance detection capabilities for subtle forms of gendered abuse.

### 4.3 Subtask 3: Multi-task Classification

The third subtask involves building a multi-task classifier capable of simultaneously detecting both gendered abuse (Label 1) and explicit language (Label 3). This approach examines the interrelationship between different forms of harmful online content and explores the benefits of joint learning across related classification tasks.

## 5 Methodology

### 5.1 Baseline 1: CNN-BiLSTM with GloVe Embeddings

Our first baseline model combines convolutional and recurrent architectures using static word embeddings. The model begins with an input layer processing tokenized sequences of fixed length, followed by a 300-dimensional GloVe embedding layer initialized with pretrained weights that remain frozen during training. A spatial dropout layer (20%) is applied to the embeddings for regularization before feeding them into the convolutional component. The CNN layer employs 64 filters with

kernel size 2 and ReLU activation, capturing local n-gram patterns while preserving sequence length through same padding. The convolutional outputs are then processed by a bidirectional LSTM layer with 128 units in each direction, incorporating 10% dropout on both connections and recurrent states to prevent overfitting. The sequence representations are condensed through global average pooling before passing through a 128-unit dense layer with ReLU activation and 10% dropout. The architecture culminates in a 2-unit softmax output layer for binary classification. This hybrid design effectively combines the CNN's strength in local feature extraction with the BiLSTM's capacity for modeling long-range dependencies, providing a strong baseline for sequence classification tasks.

## 5.2 Baseline 2: GRU-Attention with Glove Embeddings

Our second baseline implements a more sophisticated architecture integrating gated recurrent units with attention mechanisms. Unlike the first baseline, this model initializes its 300-dimensional embedding layer with GloVe vectors but allows them to be fine-tuned during training. The embeddings first pass through a spatial dropout layer (30%) for robust regularization. Two subsequent bidirectional GRU layers process the embeddings - the first with 128 units and the second with 64 units, both employing 20% dropout. The GRU outputs feed into an 8-head self-attention layer with key dimension 16, which learns to dynamically weight important contextual features. We employ skip connections to combine the attention outputs with the original GRU representations, preserving information flow. The combined features undergo parallel global average and max pooling operations to capture different aspects of the sequence. These pooled features then pass through a deep classifier comprising a 256-unit dense layer (with batch normalization and 30% dropout) followed by a 128-unit dense layer (with batch normalization and 20% dropout), both using ReLU activation. The final softmax output layer completes the architecture. This design offers enhanced capacity for modeling complex contextual relationships through its attention mechanism and deep classifier head, while the trainable embeddings allow adaptation to domain-specific language use.

## 5.3 Baseline 2: LR for task2

Used TF-IDF vectorizer to train and validate the data. Imported and fitted the LogisticRegression model with log of odds and L2 regularization to predict labels like hate and non-hate. Printed all evaluation metrics. Created a function to evaluate logistic regression on the test set by cleaning the data, applying label strings, and printing the evaluation metrics similarly.

## 6 Experimental Setup

### 6.1 Training Configuration

The models were trained with the following protocol:

- **Optimization**:
  - Adam optimizer with learning rate 0.001 ($\beta_1 = 0.9$, $\beta_2 = 0.999$)
  - Batch size of 32 for all tasks
  - Maximum 15 epochs with early stopping (patience=2)

- **Regularization**:
  - Spatial dropout (20%) on embedding layer
  - Recurrent dropout (10%) in LSTM layers
  - Additional dropout (10%) before classification

- **Training Process**:
  - 80-20 train-validation split (stratified)
  - Model checkpointing based on validation Macro-F1
  - Mixed-precision training (FP16) with NVIDIA GPU acceleration

### 6.2 Evaluation Protocol

Performance was assessed using:

- Primary metric: Macro-F1 score (custom implementation)

- Secondary metrics: Weighted precision and recall

- Full classification reports for each language-task combination

- Confusion matrices visualization

Table 2: Key Training Parameters

| Parameter | Value |
|---|---|
| Base learning rate | 0.001 |
| Batch size | 32 |
| Max epochs | 15 |
| Early stopping patience | 2 |
| Embedding dropout | 20% |
| LSTM dropout | 10% |
| Classifier dropout | 10% |

## 6.3 Implementation Details

The experimental setup featured:

- TensorFlow 2.8 with Keras API

- Custom MacroF1Score metric implementation

- Automatic mixed precision training

- GPU memory growth enabled

- Fixed random seed (42) for reproducibility

- Automated model checkpointing and logging

The training process monitored validation Macro-F1 score for early stopping and model selection, with all models evaluated on held-out test sets for final performance reporting.

## 7 Results

### 7.1 Evaluation Metrics

To comprehensively assess model performance, we employed three key metrics: precision, recall, and macro-F1 score. The macro-F1 score serves as our primary evaluation criterion as it provides a class-balanced perspective by computing the unweighted mean of per-class F1 scores, treating all classes equally regardless of their frequency. This is particularly valuable for our imbalanced dataset where equitable performance across minority classes is crucial.

Precision quantifies the model's ability to correctly identify positive instances (avoiding false positives), while recall measures its capacity to capture all relevant cases (minimizing false negatives). The macro-F1 harmonizes these metrics, offering a single comprehensive measure that accounts for both Type I and Type II errors without being skewed by class distribution. This balanced evaluation approach is especially pertinent for abuse detection tasks, where both over-classification (precision) and under-detection (recall) carry significant consequences.

The results can be seen in Table 3

Table 3: Model Performance Across Languages and Tasks

| Model | Task-Lang | P | R | Macro F1 |
|---|---|---|---|---|
| M1 | T1-EN | 0.84 | 0.80 | 0.47 |
| | T1-HI | 0.58 | 0.76 | 0.43 |
| | T1-TA | 0.58 | 0.57 | 0.55 |
| | T2-EN | 1.00 | 0.99 | 0.50 |
| | T2-HI | 0.57 | 0.57 | 0.57 |
| | T2-TA | 0.62 | 0.58 | 0.55 |
| | T3-EN | 0.63/0.58 | 0.79/0.63 | 0.44/0.40 |
| | T3-HI | 0.58/0.69 | 0.76/0.50 | 0.43/0.54 |
| | T3-TA | 0.53/0.58 | 0.53/0.64 | 0.45/0.49 |
| M2 | T1-EN | 0.79 | 0.81 | 0.68 |
| | T1-HI | 0.74 | 0.72 | 0.63 |
| | T1-TA | 0.77 | 0.76 | 0.76 |
| | T2-EN | 0.84 | 0.84 | 0.66 |
| | T2-HI | 0.36 | 0.48 | 0.35 |
| | T2-TA | 0.74 | 0.74 | 0.53 |
| | T3-EN | 0.70/0.56 | 0.75/0.57 | 0.53/0.53 |
| | T3-HI | 0.66/0.69 | 0.75/0.52 | 0.47/0.50 |
| | T3-TA | 0.62/0.69 | 0.62/0.70 | 0.62/0.64 |

Note: M1=CNN-BiLSTM, M2=GRU-Attention; T1=Task 1 (Gendered Abuse), T2=Task 2 (Targeted Abuse), T3=Task 3 (Explicit Abuse); EN=English, HI=Hindi, TA=Tamil; Task 3 metrics shown as Label 1/Label 3. All results are test set performance.

## 8 Analysis

Based on the performance metrics provided in your table, the GRU-Attention model consistently outperformed the CNN-BiLSTM architecture across all three tasks and languages. Let me analyze why this might be the case. The GRU-Attention model demonstrates significantly better performance particularly in Task 1 (gendered abuse detection), where it achieved macro F1 scores of 0.68, 0.63, and 0.76 for English, Hindi, and Tamil respectively, compared to the CNN-BiLSTM's scores of 0.47, 0.43, and 0.55. This substantial improvement (approximately 20 percentage points on average) suggests that the GRU-Attention architecture is better suited for capturing the nuanced patterns of gendered abusive content.

Several factors likely contribute to the GRU-Attention model's superior performance. First, the attention mechanism allows the model to focus on the most relevant parts of the text for classification, which is particularly important for detecting subtle forms of abuse that may hinge on specific words or phrases. Gender-based abusive content often contains certain key phrases or contextual cues that might be overlooked in a standard sequential processing approach. The attention mechanism effectively highlights these crucial indicators, giving

them appropriate weight in the final classification decision.

Additionally, GRUs (Gated Recurrent Units) are known for their efficiency in capturing long-range dependencies while mitigating the vanishing gradient problem that can affect LSTM networks. This characteristic is especially valuable when processing complex linguistic structures where abusive content may be embedded within longer sentences or depend on the relationship between words that appear far apart in the text. The simpler gating mechanism in GRUs may have allowed for more effective parameter optimization given the constraints of the dataset size.

In Task 3, the multi-task learning scenario, the performance difference between models is less pronounced, indicating that both architectures face similar challenges when simultaneously optimizing for multiple objectives (gendered abuse and explicit language detection). However, the GRU-Attention model still maintained consistent performance across the different languages. The performance improvement is particularly notable for Hindi and Tamil, suggesting that the GRU-Attention architecture may be better suited for capturing the morphological complexity and syntactic structures of these languages. This is crucial because linguistic features that indicate abuse can manifest differently across languages, and a model that can better adapt to these language-specific characteristics would naturally perform better.

In summary, the GRU-Attention model's superior performance can be attributed to its attention mechanism's ability to focus on relevant parts of text, the GRU's efficient handling of sequential information, and the architecture's adaptability to different linguistic structures across languages. These advantages make it particularly well-suited for detecting subtle forms of abusive content that rely on contextual understanding rather than just the presence of specific keywords.

## 9   Conclusion and Future Work

This paper presented our approach and results for the ICON2023 shared task on identifying gendered abuse in online content. Through systematic experimentation, we demonstrated that GRU-Attention models significantly outperformed CNN-BiLSTM architectures across all three subtasks and languages. The attention mechanism proved particularly valuable for capturing nuanced abusive

language patterns by focusing on relevant textual features, while GRUs effectively handled the sequential nature of linguistic data.

Our performance analysis revealed substantial improvements in the GRU-Attention model's ability to detect gendered abuse in Task 1, with macro F1 scores up to 21 percentage points higher than CNN-BiLSTM baselines. This performance gain was consistent across English, Hindi, and Tamil datasets, suggesting the model's strong adaptability to different linguistic structures. The multi-task learning approach for Task 3 showed that simultaneously predicting gendered abuse and explicit language remains challenging, though our GRU-Attention model maintained more consistent performance.

While our models achieved strong results, challenges remain in handling heavily code-switched content, particularly in Indic languages - an important area for future research. Through this shared task, we have developed performant models addressing the crucial problem of gendered cyber harassment that limits online freedom of expression. Our datasets and model implementations have been open-sourced to facilitate further research and development of more effective countermeasures against online abuse.

## 10   Helping Papers

Sayani Ghosal, Amita Jain, Devendra Kumar Tayal, Varun G Menon, and Akshi Kumar. 2023. Inculcating context for emoji powered bengali hate speech detection using extended fuzzy svm and text embedding models. ACM Transactions on Asian and LowResource Language Information Processing.

Subhajeet Das, Koushikk Bhattacharyya, and Sonali Sarkar. 2023. Performance analysis of logistic regression, naive bayes, knn, decision tree, random forest and svm on hate speech detection from twitter. International Research Journal of Innovations in Engineering and Technology, 7(3):24.

## References

[1] Author(s). (2024). "Title of the paper." *arXiv preprint arXiv:2404.02013*. Available: https://arxiv.org/pdf/2404.02013

[2] Tattle. (2023). "ULI Dataset Repository." Available: https://github.com/tattle-made/uli_dataset

[3] Kaggle. (2023). "Gendered Abuse Detection Shared Task – Overview." Available:

https://www.kaggle.com/competitions/
gendered-abuse-detection-shared-task/
overview

[4] Kaggle. (2023). "Gendered Abuse De-
tection Shared Task – Data." Available:
https://www.kaggle.com/competitions/
gendered-abuse-detection-shared-task/data