

# Bangla-English Code-Mixing Attacks on Large Language Models: A Comprehensive Security Analysis

- Thesis Overview
- Research Objective
- Key Research Questions
- Methodology
- Major Findings
- ❖ Novel Contributions
- Statistical Validation
- ▀ Ethical Framework
- Future Research Directions
- Technical Specifications
- Academic Impact
- Broader Implications
- Documentation Quality
- >Contact & Repository

# Bangla-English Code-Mixing Attacks on Large Language Models: A Comprehensive Security Analysis

## Thesis Overview

**Title:** Bangla-English Code-Mixing Attacks on Large Language Models: A Comprehensive Security Analysis

**Authors:** Two undergraduate students at SUST-IICT

**Supervisor:** Dr. Ahsan Habib, Associate Professor, IICT, SUST

**Submission Date:** November 2024

**Degree:** Bachelor of Science in Software Engineering

**Institution:** Shahjalal University of Science and Technology (SUST)

## Research Objective

This thesis presents the **first comprehensive study** of Bangla-English code-mixing attacks on Large Language Models (LLMs), investigating vulnerabilities affecting **230 million Bangla speakers worldwide**—the 8th most spoken language globally that has been completely overlooked in adversarial AI safety research.

## Key Research Questions

1. **RQ1:** Does Bangla-English code-mixing with phonetic perturbations bypass LLM safety filters?
2. **RQ2:** Which phonetic and romanization features enable Bangla attacks?
3. **RQ3:** Are all major LLMs vulnerable to Bangla attacks?
4. **RQ4:** Does tokenization disruption explain Bangla attack success?

# Methodology

## Three-Step Attack Pipeline

1. **English Baseline:** 200 harmful prompts across 10 categories (hate speech, illegal activities, etc.)
2. **Code-Mixing (CM):** Convert to Bangla-English mix (optimal 70% Bangla, 30% English)
3. **Phonetic Perturbations (CMP):** Apply misspellings to sensitive English keywords

## Experimental Design

- **Models Tested:** GPT-4o-mini, Llama-3-8B, Mistral-7B (Gemma excluded due to budget)
- **Templates:** 5 jailbreak templates including novel “Sandbox” approach
- **Scale:** 27,000 model responses collected
- **Evaluation:** Automated LLM-as-judge methodology
- **Statistics:** Wilcoxon signed-rank tests, correlation analysis

## Major Findings

### ✓ Core Results (All RQs Answered Affirmatively)

**RQ1 - Attack Effectiveness:** - **40.1% Attack Success Rate** with CMP prompts vs 36.1% English baseline - **Statistically significant** improvement ( $p=0.0070$ ) - Effective across all temperature settings (0.2, 0.6, 1.0)

**RQ2 - Language-Specific Patterns:** - **English word targeting 68% more effective** than Bangla word perturbations - **70:30 Bangla:English ratio** yields optimal attack success - **Vowel substitution and consonant doubling** most effective perturbation types

**RQ3 - Cross-Model Vulnerability:** - **All 3 tested models vulnerable** but with dramatic variation: - Mistral-7B: **81.8% AASR** (critical vulnerability) - Llama-3-8B: **22.7% AASR** (moderate vulnerability) - GPT-4o-mini: **16.0% AASR** (low but non-zero vulnerability)

**RQ4 - Tokenization Mechanism:** - **Strong correlation** between token fragmentation and attack success - Pattern consistent with Hindi-English findings ( $r=0.94$  from prior work) - Phonetic perturbations fragment harmful keywords into semantically inert subword units

## Surprising Discoveries

1. **Jailbreak Templates Reduce Effectiveness** - Contrary to prior research, complex jailbreak templates **decrease** Bangla attack success - Simple “None” template: **46.2% AASR** vs engineered templates: **35.1-42.5% AASR** - Suggests code-mixing provides sufficient obfuscation without additional engineering

**2. English-Centric Safety Filter Design** - Safety filters optimized for English keyword detection - **85% effectiveness gap** between English vs Bangla word targeting - Reveals fundamental architectural limitation in current LLM safety training

## Novel Contributions

### 1. First Bangla Code-Mixing Study

- Systematic evaluation of 230M speaker population previously untested
- Quantitative vulnerability baselines across major LLM architectures
- Independent validation extending beyond Hindi-English research

### 2. English Word Targeting Strategy

- Discovery that perturbing English words within Banglish contexts is dramatically more effective
- Reveals English-centric bias in current safety filter design
- Provides optimization guidance for defensive measures

### 3. Template Ineffectiveness Finding

- Challenges universal applicability of jailbreak templates across languages
- Demonstrates that simpler attacks may be more effective for code-mixing scenarios
- Has implications for both attack optimization and defense prioritization

### 4. Tokenization Mechanism Validation

- Independent confirmation of fragmentation hypothesis for Bangla-English
- Strengthens theoretical understanding beyond language-specific empirics
- Informs development of tokenization-robust defense architectures

### 5. Romanization Variability Analysis

- Identification of non-standardized romanization as unique Bangla vulnerability
- Distinguishes Bangla from languages with established romanization standards
- Suggests differential security implications across linguistic communities

### 6. Scalable Research Framework

- Cost-effective methodology: **\$1.50-2.00 per language** for 50-prompt studies
- Configuration-driven experimentation enabling easy replication
- Directly applicable to **20+ other Indic languages** serving 1B+ speakers

# Statistical Validation

## Significance Testing

- **English → CM:**  $p=0.0209$  (significant)
- **English → CMP:**  $p=0.0070$  (highly significant)
- **CM → CMP:**  $p=0.1291$  (modest, limited by partial data collection)

## Effect Sizes

- **Large effects** for GPT-4o-mini and Llama-3-8B (Cohen's  $d > 0.68$ )
- **Negligible effects** for Mistral-7B due to baseline vulnerability saturation

## Confidence Intervals

- **English:** 29.7-35.1% AASR (95% CI)
- **CM:** 38.9-45.3% AASR (95% CI)
- **CMP:** 42.6-49.4% AASR (95% CI)

# Ethical Framework

## Responsible Disclosure Protocol

1. **Pre-publication:** Thesis submission (Nov 2024)
2. **Vendor notification:** 60-90 day patch window (Dec 2024-Jan 2025)
3. **Public disclosure:** After vendor remediation
4. **Dataset protection:** Research-only access, no public release

## Harm Mitigation

- **Limited dataset size** to prevent comprehensive exploitation guide
- **Abstract perturbation descriptions** requiring manual reconstruction effort
- **No automated attack tools** provided
- **Prominent content warnings** throughout documentation

## Community Benefit

- **Advances AI safety** for 230M underserved speakers
- **Enables targeted fixes** through empirical vulnerability documentation
- **Promotes linguistic equity** in global AI deployment
- **Provides scalable framework** for community-driven multilingual safety research

# Future Research Directions

## Immediate Extensions

1. **Scale to 460 prompts** for direct Hinglish comparison

2. **Complete Gemma evaluation** for full 4-model coverage
3. **Human validation study** to calibrate LLM-as-judge reliability
4. **Automated code-mixing** using NMT models for scalability

## Medium-Term Research

1. **Extend to 10+ Indic languages** (Tamil, Telugu, Marathi, Urdu, etc.)
2. **Develop defense mechanisms** including romanization normalization
3. **White-box interpretability** using Integrated Gradients
4. **Multi-turn attack strategies** beyond single-turn evaluation

## Long-Term Vision

1. **Multilingual safety benchmark** across 100+ languages
2. **Tokenization-robust safety architecture** operating at semantic level
3. **Equitable AI safety framework** with mandatory language coverage requirements
4. **Community-driven evaluation** enabling native speaker contributions

# Technical Specifications

## Dataset Characteristics

- **200 total prompts** (scaled from 50-prompt validation)
- **10 harm categories** with balanced distribution
- **3 prompt variants** per scenario (English, CM, CMP)
- **27,000 total responses** collected (75% of planned 36,000)

## Model Coverage

- **GPT-4o-mini:** OpenAI's most deployed model
- **Llama-3-8B:** Meta's open-source benchmark
- **Mistral-7B:** European alternative architecture
- **Gemma-1.1-7B:** Excluded due to budget constraints

## Cost Analysis

- **Total experimental cost:** ~\$1.50 (extremely cost-effective)
- **Per-language replication cost:** \$1.50-2.00 for 50-prompt studies
- **Full-scale cost:** \$15-20 for 460-prompt replication
- **Evaluation cost:** \$0.95 for 27,000 response assessments

# Academic Impact

## Publication Readiness

- **Novel empirical findings** not previously reported
- **Methodological rigor** with statistical validation
- **Practical significance** for 230M speaker community

- **Scalable framework** enabling follow-up research

## Conference Suitability

- **ICML/NeurIPS:** Multilingual AI safety track
- **ACL/EMNLP:** Code-mixing and multilingual NLP
- **IEEE S&P/USENIX:** Adversarial ML and security
- **FAccT:** AI ethics and linguistic equity

## Citation Potential

- **First systematic Bangla study** in adversarial AI
- **Methodological contributions** for multilingual safety research
- **Policy implications** for equitable AI deployment
- **Framework replication** enabling comparative studies

## Broader Implications

### Linguistic Equity in AI

- **Expose systematic bias** in current safety training approaches
- **Advocates for 230M speakers** currently underserved
- **Challenges English-centric** deployment practices
- **Promotes inclusive** AI safety standards

### Security Architecture

- **Reveals tokenization brittleness** as fundamental vulnerability
- **Motivates semantic-level** safety mechanisms
- **Informs defensive strategy** against code-mixing attacks
- **Guides multilingual** red-teaming practices

### Policy Recommendations

1. **Language coverage mandates** for major languages (>100M speakers)
2. **Transparency requirements** for vulnerability disclosure
3. **Proportional safety training** reflecting global speaker distributions
4. **Community engagement protocols** for native speaker involvement

## Documentation Quality

### Comprehensive Coverage

- **9 main chapters** with detailed methodology and results
- **3 appendices** providing replication materials
- **Statistical rigor** with complete test documentation
- **Ethical framework** addressing dual-use concerns

### Reproducibility

- **Complete configuration files** for experimental replication
- **Detailed statistical analysis** enabling verification
- **Clear methodology** with step-by-step procedures
- **Code availability** through modular framework structure

## Academic Standards

- **Proper literature review** situating work in broader context
  - **Transparent limitations** acknowledging scope constraints
  - **Responsible disclosure** following security research best practices
  - **Future work** providing clear research roadmap
- 

## Contact & Repository

**Repository:** sandwipshanto/Thesis

**Institution:** Shahjalal University of Science and Technology

**Department:** Institute of Information and Communication Technology

**Research Area:** Multilingual AI Safety, Adversarial ML, Code-Mixing Attacks

*This thesis represents groundbreaking research in multilingual LLM security, providing the first systematic evaluation of Bangla-English code-mixing vulnerabilities and establishing a foundation for equitable AI safety across linguistic communities worldwide.*