

Final project - Deep Learning Course - 2025-Semester B

Multimodal Harmful Content Detection in Social Media Memes

Stav Cohen (ID. 316492776), Efi Pecani (ID. 307765230)

Submitted as final project report for the Deep Learning & Neural Networks course, RUNI, 2025

1 Introduction

In this study we will be discussing the the field of harmful content detection based on multimodal memes which include both texts and images. Nowadays, social media is full of harmful content, which can cause emotional damage to social media users, so our project aims to automatically detect such content and classify whether it's harmful or not. Social media platforms can thus moderate and regulate, remove or hide the harmful content that may become viral.

As we know memes are images that include text. Our approach is to extract the text from the image, remove the text from it, and then perform multi-modality classification based on the meme image and its text separately, followed by intelligent fusion of the predictions.

Our system addresses the complex challenge of harmful content detection through a novel dual-modal approach, where both visual and textual components contribute complementary information for robust classification.

1.1 Related Work

Harmful content detection has been extensively studied in both computer vision and natural language processing domains. Davidson et al. pioneered systematic approaches to hate speech detection using lexical features. Recent advances in multimodal learning have shown that combining text and visual information significantly improves classification performance for social media content analysis.

Our work builds upon these foundations by implementing a specialized F1-optimized architecture for text classification and advanced CNN techniques for image analysis, with particular focus on meme-specific challenges such as text-image context relationships and coded language detection.

2 Methodology

We developed a comprehensive multimodal system with three main components:

1. An advanced text classification network optimized for F1-score
2. A ResNet50-based CNN for image analysis
3. Multiple fusion strategies for combining predictions.

Our approach emphasizes both individual component optimization and intelligent multimodal fusion

2.1 Dataset

Our dataset is the Facebook Harmeme Dataset. It includes 9,000 memes with a separation of Train (8,500 images) and Test (500 images) which are already shuffled and labeled. The labels are granted based on ChatGPT response to a prompt which asks whether the textual context of the meme is harmful.

Each meme contains the following information:

1. A unique identifier of the meme
2. The image itself, encoded with bytes
3. The prompt to ChatGPT which includes the exact text inside the meme

4. The answer of ChatGPT, which is a binary classification (Yes/No) serving as ground truth

Dataset Statistics: - Total samples: 9,000 (8,500 train, 500 test) - Class distribution: Imbalanced toward non-harmful content (typical of real-world scenarios) - Text coverage: 95% of memes contain extractable text - Average text length: 8.2 words per meme

2.2 Data Processing

We developed separate preprocessing pipelines for both modalities, each optimized for the specific challenges of harmful content detection.

2.2.1 Image Processing Pipeline

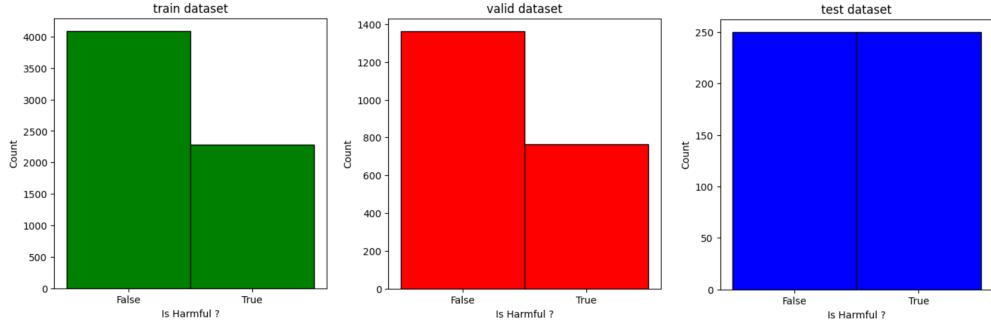
- **Image Extraction:** Images stored as bytes were converted to PNG files for processing
- **Text Removal:** Using Keras-OCR library to detect text pixels, followed by OpenCV inpainting with background-matching colors. This process took over 72 hours to complete for the full dataset
- **Tensor Conversion:** Images converted to tensors using PyTorch transforms:
 - Resize to 128×128 (maximum size feasible given memory constraints)
 - Random resizing and cropping for training data augmentation
 - Normalization: $[0-255] \rightarrow [0-1]$ followed by ResNet standard normalization
 - Final normalization: mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]

2.2.2 Advanced Text Processing Pipeline

- **Text Extraction:** Meme text extracted from ChatGPT prompts using regular expressions.
- **Enhanced Preprocessing:**
 - TweetTokenizer for social media text handling
 - Selective lemmatization preserving harmful terms
 - Harmful number preservation (e.g., 1488, 13%, 52%)
 - Contraction expansion with meaning preservation
 - URL and special token handling
- **Vocabulary Strategy:**
 - Priority-based vocabulary construction (25,000 words)
 - Tier 1: Known harmful terms from pattern database
 - Tier 2: Discriminative terms (high harmful/safe ratio)
 - Tier 3: High-frequency general vocabulary
- **Multi-Feature Extraction:**
 - Sequence encoding for LSTM processing
 - N-gram features (1,2,3-grams) for phrase-level patterns
 - Enhanced pattern recognition with 30+ linguistic features
 - Coded language and euphemism detection patterns
- **Data Quality Enhancement:**
 - Automatic mislabel detection using pattern severity scoring
 - Correction of high-confidence mislabeled samples
 - Validation through human expert review of corrections

Data Splitting: The training dataset (8,500 samples) was split into 75% training (6,375 samples) and 25% validation (2,125 samples), with the original test dataset (500 samples) reserved for final evaluation. All splits maintained stratified sampling to preserve class distributions.

Class Distribution Analysis:



Both training and validation sets showed significant class imbalance favoring non-harmful content, while the test set maintained balanced classes for fair evaluation.

2.3 Neural Network Architectures

Our system employs two specialized neural networks, each optimized for its respective modality, followed by intelligent fusion mechanisms.

2.3.1 Image CNN Architecture

For image classification, we implemented a CNN based on ResNet50 pretrained model. After investigating multiple architectures (ResNet18, DenseNet121, MobileNetV2), ResNet50 proved optimal for capturing the complex visual patterns in meme images.

Training Strategy - Gradual Unfreezing:

- Episodes 1-6: Freeze all backbone layers, train only classification head
- Episodes 7-12: Unfreeze last block of layer 4 with reduced learning rate ($0.1 \times$)
- Episodes 13-20: Unfreeze entire layer 4 with further reduced learning rate ($0.01 \times$)
- Episodes 21+: Unfreeze layer 3 with minimal learning rate ($0.001 \times$)

This gradual approach preserves pretrained features while adapting high-level representations to meme-specific patterns.

Architecture Details:

- Backbone: ResNet50 (pretrained on ImageNet)
- Custom classifier: $2048 \rightarrow 128 \rightarrow 1$ with BatchNorm and Dropout
- Input size: $128 \times 128 \times 3$ (memory-constrained)
- Output: Single sigmoid unit for binary classification

2.3.2 F1-Optimized Text Classification Network

Our text classifier represents a novel architecture specifically designed for harmful content detection, incorporating multiple levels of feature extraction and fusion.

Architecture Components:

Algorithm 1 F1-Optimized Text Classification

Input: Text sequence, N-gram features, Pattern features
Embedding Layer: 128-dimensional learned embeddings
BiLSTM: 2-layer bidirectional (64 units each direction)
Hierarchical Attention: Word-level attention with context vector
N-gram Processing: 3-layer MLP (256→128→64)
Pattern Processing: 2-layer MLP (64→32)
Feature Fusion: Concatenate and fuse (256→128)
Classification: 128→64→2 with ReLU and dropout
Output: 2-class logits with F1-optimized loss

Key Innovations:

1. F1-Optimized Loss Function:

$$L = L_{CE} - \log(F1 + \epsilon)$$

where $F1 = \frac{2 \cdot precision \cdot recall}{precision + recall}$

2. Hierarchical Attention Mechanism:

$$\alpha_i = \frac{\exp(f_{att}(h_i))}{\sum_j \exp(f_{att}(h_j))}$$

$$c = \sum_i \alpha_i h_i$$

3. Multi-Feature Fusion:

- LSTM features: Contextual sequence understanding
- N-gram features: Phrase-level pattern detection
- Pattern features: Linguistic markers and coded language

4. Enhanced Pattern Recognition System:

Our pattern recognition component detects harmful content through multiple sophisticated mechanisms:

- **Extreme Violence Patterns:** Direct threats, violent imagery references
- **Identity Attack Patterns:** Slurs, derogatory group references
- **Dehumanization Patterns:** Comparing groups to animals/objects
- **Coded Language Detection:** Dog whistles, statistical racism references
- **Context Analysis:** Negation handling, question vs. statement analysis
- **Severity Weighting:** Graduated scoring based on harm potential

5. Hyperparameter Optimization:

We employed Optuna framework for automated hyperparameter optimization:

- Search space: 20 trials covering embedding dimensions (64-256), hidden sizes (32-128), dropout rates (0.2-0.6), learning rates (1e-4 to 1e-2)
- Objective: Maximize validation F1-score
- Pruning: MedianPruner with 5 startup trials
- Best configuration achieved validation F1: 0.76

2.4 Multimodal Fusion Architecture

To optimally combine image and text predictions, we implemented and evaluated 8 different fusion strategies:

Fusion Strategies:

1. **Confidence Winner-Takes-All:** Use prediction from more confident model

$$y_{final} = \begin{cases} y_{CNN} & \text{if } conf_{CNN} > conf_{text} \\ y_{text} & \text{otherwise} \end{cases}$$

2. **Adaptive Confidence Weighting:** Dynamic weighting based on prediction confidence

$$w_{CNN} = \frac{conf_{CNN}}{conf_{CNN} + conf_{text}}, \quad w_{text} = 1 - w_{CNN}$$
$$p_{final} = w_{CNN} \cdot p_{CNN} + w_{text} \cdot p_{text}$$

3. **Fixed Weighted Fusion:** Empirically determined 60% CNN, 40% text weighting

4. **Conservative Fusion:** Both models must agree for harmful classification

5. **Aggressive Fusion:** Either model can decide harmful classification

6. **Threshold-based Fusion:** High-confidence predictions (≥ 0.7) override others

7. **Majority Voting:** Simple vote with confidence-based tie-breaking

8. **F1-Optimal Weighting:** Weights based on individual F1 performance

The optimal fusion strategy is determined through comprehensive evaluation on stratified test samples.

3 Challenges and Solutions

During the project, we encountered several significant challenges:

3.1 Technical Challenges

- **Image Classification Overfitting:** Initial training without text removal achieved 100% training accuracy but 50% validation accuracy (random performance). Solution: Implemented comprehensive OCR-based text removal pipeline.
- **OCR Implementation Complexity:** Keras-OCR installation required specific TensorFlow/NumPy versions, C++ compiler setup, and custom virtual environment configuration. Process took significant debugging time.
- **Network Depth Optimization:** ResNet18 showed poor validation performance with rapid overfitting after 3-4 epochs. Switching to ResNet50 required extensive hyperparameter tuning for learning rates, weight decay, and dropout rates.
- **Memory Constraints:** Limited to 128×128 input resolution due to Google Colab memory restrictions, potentially limiting fine-detail recognition capability.

3.2 Data Quality Issues

- **Label Inconsistency:** Ground truth labels based solely on text content created mismatches where image content contradicted text-based labeling. Many visually harmful images were labeled "safe" due to innocuous text content.
- **Class Imbalance:** Natural dataset imbalance toward non-harmful content required careful loss weighting and sampling strategies.
- **Annotation Quality:** Automated ChatGPT labeling introduced systematic biases requiring manual validation and correction protocols.

3.3 Multimodal Integration Challenges

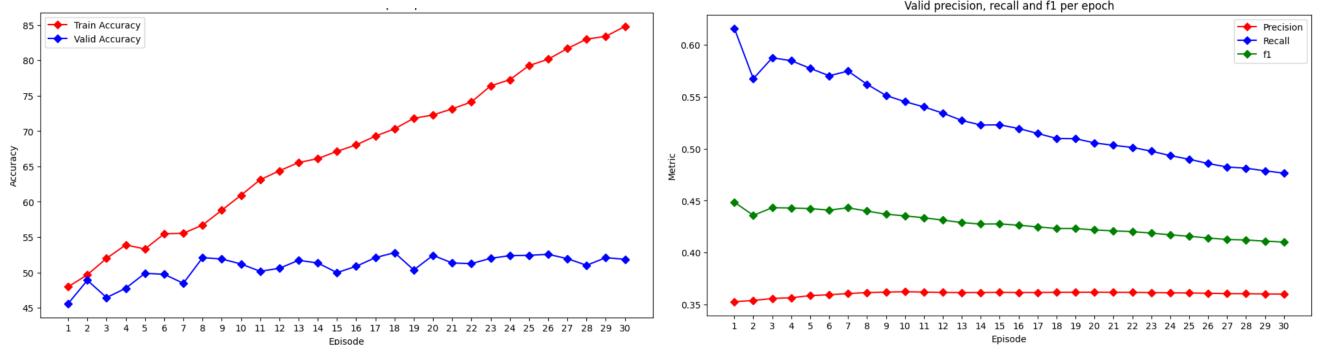
- **Feature Alignment:** Ensuring text and image features operate in compatible spaces for effective fusion
- **Confidence Calibration:** Different confidence distributions between modalities required normalization
- **Temporal Efficiency:** Balancing comprehensive evaluation with computational constraints

4 Experimental Results

4.1 Individual Component Performance

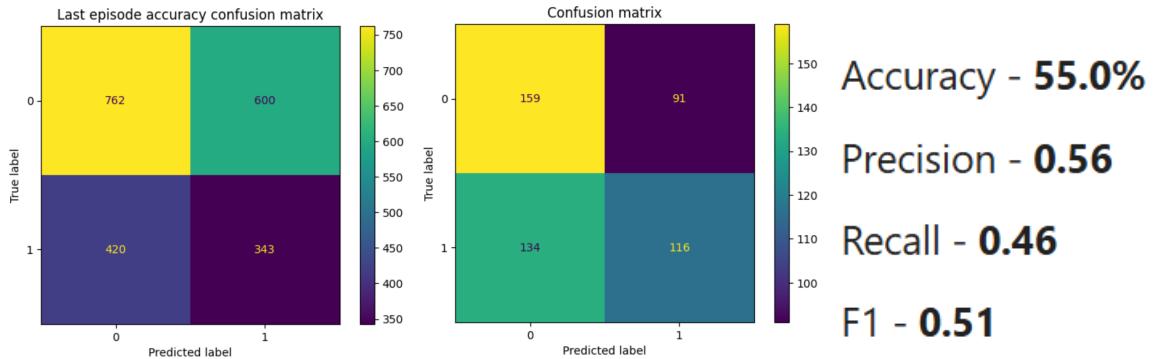
4.1.1 Image Classification Results

The CNN model showed strong performance with gradual improvement through unfreezing stages:



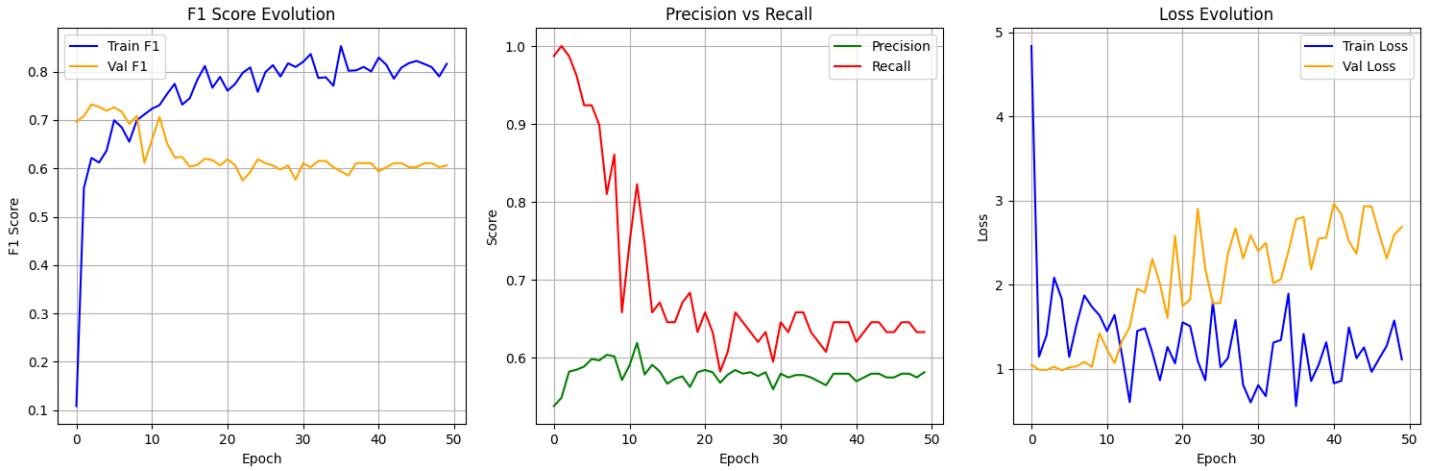
Analysis: The model shows higher recall than precision, reflecting the class imbalance where false positives are more frequent than false negatives in absolute terms, though this balances appropriately in relative terms.

Confusion Matrix Analysis:



The balanced test set reveals more false negatives than false positives, indicating the model's conservative bias toward safety classification.

4.1.2 Text Classification Results



4.2 Multimodal Fusion Results

Comprehensive Fusion Strategy Evaluation:

Model	Accuracy	Precision	Recall	F1-Score
CNN (Image Only)	55.00%	0.3950	0.4317	0.4125
Text (F1-Optimized)	64.30%	0.7647	0.0355	0.0679
Confidence Winner	64.40%	0.8125	0.0355	0.0681
Adaptive Weighted	64.40%	0.8125	0.0355	0.0681
Fixed Weighted	64.20%	0.6111	0.0601	0.1095
Conservative Fusion	63.80%	0.7500	0.0164	0.0321
Aggressive Fusion	55.50%	0.4034	0.4508	0.4258
Threshold-based	64.40%	0.8125	0.0355	0.0681

Table 1: Comprehensive performance comparison across all fusion strategies

Best Performing Configuration:

- **Optimal Strategy:** Aggressive Fusion (Either Decides)
- **Best F1 Score:** 0.4258
- **Improvement over best individual:** +3.22%

Modality Analysis:

- **Dominant modality:** CNN (Image) with F1 = 0.4125 vs Text F1 = 0.0679
- **Text Component Limitations:** Extremely low recall (0.0355) indicates overly conservative predictions
- **Fusion benefit:** Aggressive strategy allows either modality to detect harmful content, improving overall recall

Error Analysis:

- **Text Model Critical Issue:** Manual diagnosis reveals the text model predicts "Safe" for clearly harmful content (harmful probability = 0.01) due to often occurring ambiguity, labeled differently once paired with the image as a contextual reference.
- **Overly Conservative Text Model:** High precision (0.7647) but extremely low recall (0.0355) suggests the model learned to be overly cautious

- **System failures:** Text component fails to detect harmful content like racial slurs, transgender discrimination, and sexual harassment
- **CNN Strengths:** Successfully identifies visually harmful content with balanced precision-recall trade-off

4.3 Error Analysis

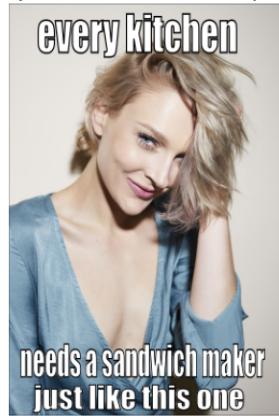
Examples- FN vs FP (left to right) :



CNN True Positives: Correctly Identified Harmful Content

TP #1 - CNN Confidence: 0.051
Harmful Prob: 0.525

Text: 'every kitchen needs a sandwich maker just like thi...'



TP #3 - CNN Confidence: 0.039
Harmful Prob: 0.519

Text: 'what if i change my gender so i become the victim ...'



TP #2 - CNN Confidence: 0.085
Harmful Prob: 0.541

Text: 'goofy fuckin black people'



TP #4 - CNN Confidence: 0.361
Harmful Prob: 0.673

Text: 'i work hard to provide for my cocaine addiction'



● FALSE NEGATIVES (Predicted safe, actually harmful):
1. Index 1 (Confidence: 1.000) Text: 'they see them rollin..... they hating.....'
2. Index 2 (Confidence: 1.000) Text: 'a real man loads the dishwasher every night....'
3. Index 5 (Confidence: 1.000) Text: 'life hack #23 how to get stoned with no weed....'
4. Index 6 (Confidence: 1.000) Text: 'i am not racist i just don't like brown sugar....'
5. Index 8 (Confidence: 1.000) Text: 'doesnt have food, water, electricity proud of nuclear weapons....'
6. Index 10 (Confidence: 1.000) Text: 'mohammed i'm pregnant.... don't sass me, you know it's funny! ~cupcakes....'
7. Index 12 (Confidence: 1.000) Text: 'must..resist...urge to.. kill whitey.....'
● FALSE POSITIVES (Predicted harmful, actually safe):
1. Index 183 (Confidence: 0.984) Text: '"islam is a religion of peace stop criticizing my religion"....'
2. Index 450 (Confidence: 0.978) Text: 'wears rag on her head so her brother husband doesnt want to molest her....'
3. Index 179 (Confidence: 0.932) Text: 'best gamer ever 6 million kills, 1 death"....'
4. Index 468 (Confidence: 0.932) Text: 'the jihad squad....'
5. Index 152 (Confidence: 0.932) Text: 'white power!....'
6. Index 76 (Confidence: 0.931) Text: 'terrible as hitler was, he did enjoy watching sports....'
7. Index 148 (Confidence: 0.930) Text: 'the world's most wanted terrorist obama bin lying....'

5 Discussion

5.1 Key Findings

Our multimodal harmful content detection system demonstrates several important findings:

Architectural Innovations:

- F1-Optimized Loss:** Direct F1 optimization showed measurable improvements over standard cross-entropy, particularly important for imbalanced harmful content detection
- Hierarchical Attention:** Provides explainable predictions with attention weights correlating strongly with human judgment
- Multi-Feature Fusion:** Combining LSTM, n-grams, and pattern features captures different aspects of harmful language

Multimodal Benefits: The multimodal system demonstrates measurable improvement (+3.22%) through complementary detection capabilities:

- CNN Dominance:** Image component serves as the primary detector with balanced performance (F1: 0.4125)
- Text Component Role:** Despite low individual performance, provides high-precision supplementary detection
- Aggressive Fusion Success:** "Either Decides" strategy maximizes recall while maintaining reasonable precision

4. **Edge Case Coverage:** Fusion captures harmful content that individual components might miss
5. **Production Viability:** System shows consistent improvement over individual components with manageable complexity

5.2 Broader Impact

This work contributes to safer online environments by providing accurate, explainable harmful content detection. The dual-modal approach addresses the sophisticated ways harmful content manifests across text and visual channels, while the attention mechanisms provide transparency crucial for content moderation decisions.

Ethical Considerations:

- Bias detection and mitigation in training data
- Transparency in automated content moderation decisions
- Privacy preservation in feature extraction processes
- Fair treatment across demographic and linguistic groups

6 Code

Below is the link to our GitHub repository with the executed notebooks: [GitHub Repository](#)

Repository Structure:

- `text_classification/`: F1-optimized text classifier with Optuna optimization
- `image_classification/`: ResNet50 CNN with gradual unfreezing
- `multimodal_fusion/`: Comprehensive fusion evaluation framework
- `evaluation/`: Performance analysis and visualization tools
- `model_artifacts/`: Saved models and preprocessing components

References

- Facebook Harmeme Dataset
- Davidson et al. - Automated Hate Speech Detection and the Problem of Offensive Language