

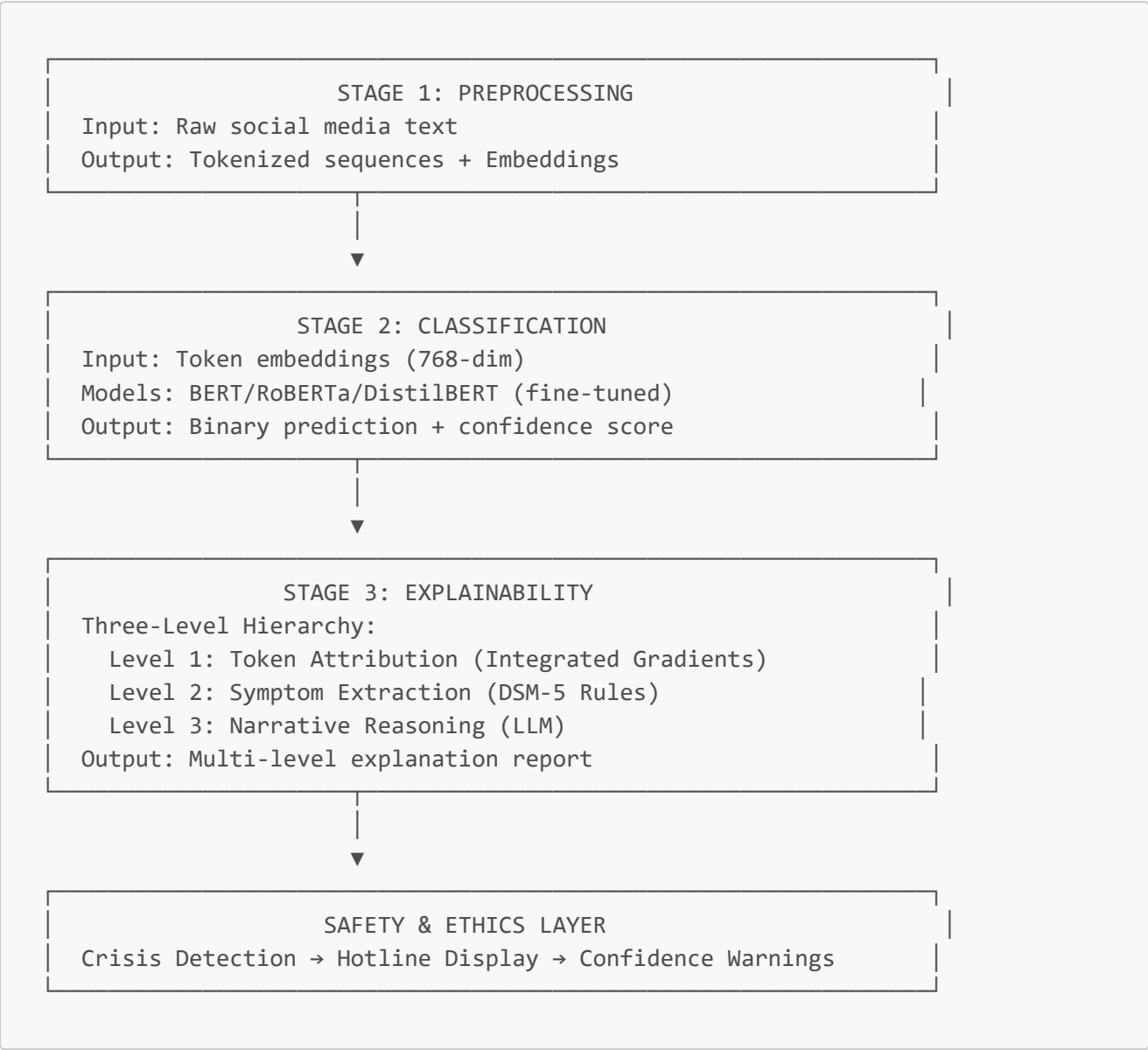
# Methodology

[← Back to Mathematical Modeling](#) | [Next: Experiments →](#)

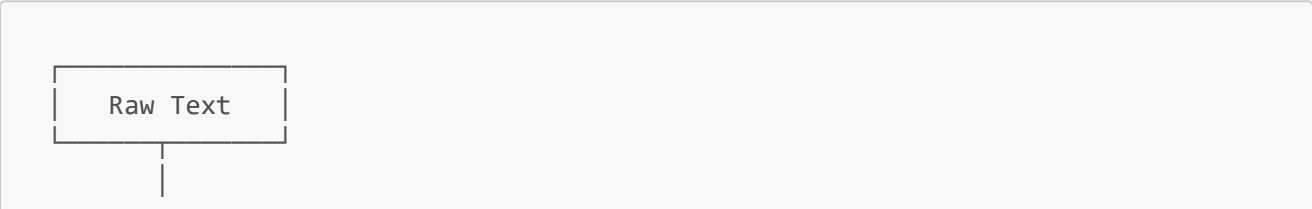
## 1. System Architecture Overview

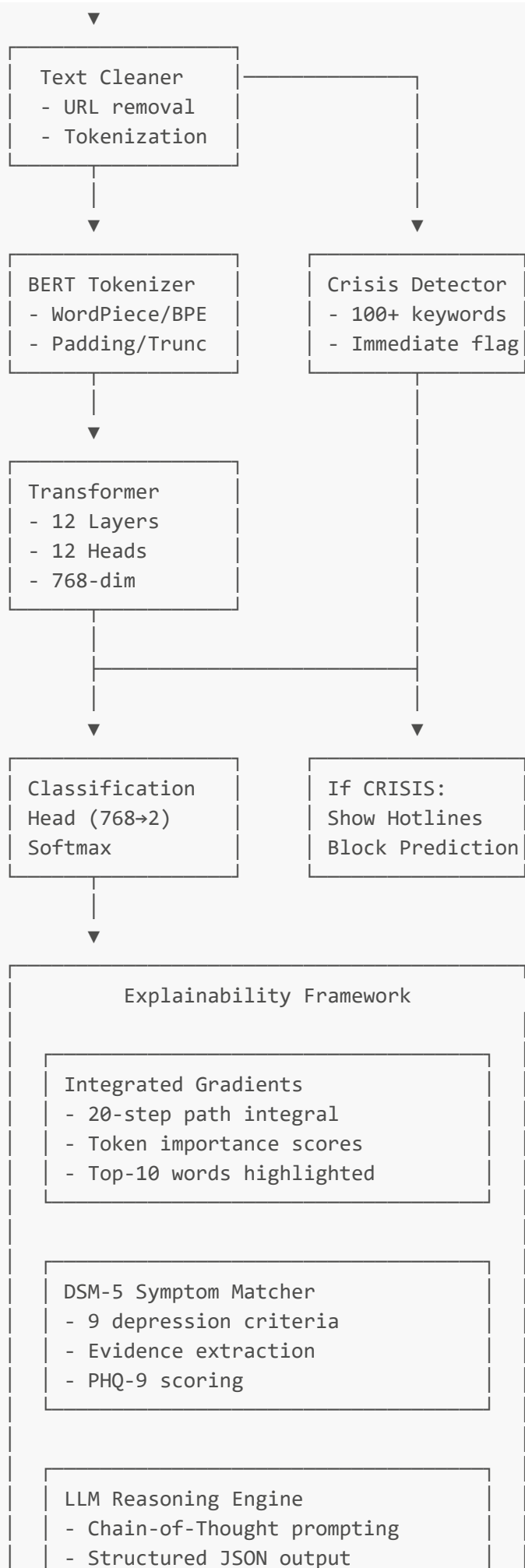
### 1.1 High-Level Pipeline

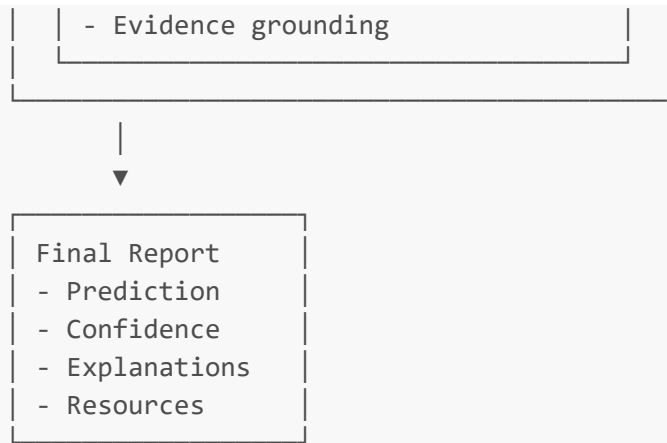
The system follows a three-stage pipeline for explainable depression detection:



### 1.2 Component Interaction Diagram







---

## 2. Stage 1: Preprocessing Module

### 2.1 Text Cleaning Pipeline

#### Implementation:

```
import re
from typing import List

class TextPreprocessor:
    def __init__(self):
        self.url_pattern = r'http\S+|www.\S+'
        self.username_pattern = r'@\w+|u/\w+'
        self.special_char_pattern = r'^\w\s.,!?\'\"-]'

    def clean_text(self, text: str) -> str:
        """Apply all cleaning steps sequentially."""
        # Step 1: Lowercase (optional for BERT)
        # text = text.lower() # Commented out—BERT is case-sensitive

        # Step 2: Remove URLs
        text = re.sub(self.url_pattern, '[URL]', text)

        # Step 3: Remove usernames
        text = re.sub(self.username_pattern, '[USER]', text)

        # Step 4: Remove special characters
        text = re.sub(self.special_char_pattern, '', text)

        # Step 5: Normalize whitespace
        text = re.sub(r'\s+', ' ', text).strip()

        return text

    def batch_clean(self, texts: List[str]) -> List[str]:
```

```
"""Clean multiple texts efficiently."""
return [self.clean_text(text) for text in texts]
```

### Example Transformations:

| Raw Text                                  | Cleaned Text               |
|---|----------------------------|
| "Check out https://example.com @user123!" | "Check out [URL] [USER]"   |
| "I feel 😞 so sad... #depression"          | "I feel so sad depression" |
| "u/throwaway123 posted this link"         | "[USER] posted this link"  |

## 2.2 Tokenization Strategy

### BERT (WordPiece):

```
from transformers import BertTokenizer

tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')

text = "I feel worthless and hopeless"
tokens = tokenizer.tokenize(text)
# Output: ['i', 'feel', 'worth', '##less', 'and', 'hope', '##less']

# Convert to IDs
token_ids = tokenizer.convert_tokens_to_ids(tokens)
# Add special tokens
input_ids = tokenizer.encode(text, add_special_tokens=True)
# Output: [101, 1045, 2514, 4276, 3238, 1998, 3246, 3238, 102]
#          [CLS  I   feel worth less and  hope less SEP]
```

### RoBERTa (Byte-Pair Encoding):

```
from transformers import RobertaTokenizer

tokenizer = RobertaTokenizer.from_pretrained('roberta-base')

text = "I feel worthless and hopeless"
tokens = tokenizer.tokenize(text)
# Output: ['I', 'Ġfeel', 'Ġworth', 'less', 'Ġand', 'Ġhope', 'less']
#          (Ġ indicates space before token)

input_ids = tokenizer.encode(text, add_special_tokens=True)
# Output: [0, 100, 619, 966, 1203, 8, 1991, 1203, 2]
#          <s> I   feel worth less and hope less </s>
```

## 2.3 Sequence Padding and Truncation

### Padding Strategy:

```
from transformers import AutoTokenizer
import torch

tokenizer = AutoTokenizer.from_pretrained('roberta-base')

texts = [
    "Short text",
    "This is a much longer text that requires padding to match"
]

# Tokenize with automatic padding
encoded = tokenizer(
    texts,
    max_length=512,
    padding='max_length',      # Pad to max_length
    truncation=True,           # Truncate if exceeds max_length
    return_tensors='pt'        # Return PyTorch tensors
)

print(encoded['input_ids'].shape) # torch.Size([2, 512])
print(encoded['attention_mask'].shape) # torch.Size([2, 512])
```

### Attention Mask:

```
# Example attention mask
# 1 = real token, 0 = padding token
text = "I feel sad"
tokens = ["<s>", "I", "feel", "sad", "</s>", "<pad>", "<pad>", "<pad>"]
attention_mask = [1, 1, 1, 1, 1, 0, 0, 0]
```

**Purpose:** Model ignores padded positions during attention computation.

---

## 3. Stage 2: Classification Module

### 3.1 Model Architecture

#### Base Transformer (BERT/RobERTa):

```
import torch.nn as nn
from transformers import AutoModel

class DepressionClassifier(nn.Module):
```

```

def __init__(self, model_name='roberta-base', num_classes=2, dropout=0.1):
    super().__init__()

    # Load pre-trained transformer
    self.transformer = AutoModel.from_pretrained(model_name)

    # Classification head
    self.dropout = nn.Dropout(dropout)
    self.classifier = nn.Linear(768, num_classes) # 768 → 2

def forward(self, input_ids, attention_mask):
    # Transformer forward pass
    outputs = self.transformer(
        input_ids=input_ids,
        attention_mask=attention_mask,
        output_hidden_states=True,
        output_attentions=True
    )

    # Extract [CLS] token representation
    cls_output = outputs.last_hidden_state[:, 0, :] # Shape: [batch_size,
768]

    # Apply dropout and classification
    cls_output = self.dropout(cls_output)
    logits = self.classifier(cls_output) # Shape: [batch_size, 2]

    return {
        'logits': logits,
        'hidden_states': outputs.hidden_states,
        'attentions': outputs.attentions
    }

```

## 3.2 Fine-Tuning Procedure

### Training Loop:

```

from torch.optim import AdamW
from transformers import get_linear_schedule_with_warmup

# Initialize model
model = DepressionClassifier('roberta-base')
model.to('cuda')

# Loss function with class weights
class_weights = torch.tensor([1.09, 0.93]).to('cuda') # Control, Depression
criterion = nn.CrossEntropyLoss(weight=class_weights)

# Optimizer
optimizer = AdamW(

```

```

    model.parameters(),
    lr=2e-5,
    weight_decay=0.01
)

# Learning rate scheduler
total_steps = len(train_loader) * num_epochs
warmup_steps = 100

scheduler = get_linear_schedule_with_warmup(
    optimizer,
    num_warmup_steps=warmup_steps,
    num_training_steps=total_steps
)

# Training loop
for epoch in range(num_epochs):
    model.train()
    total_loss = 0

    for batch in train_loader:
        input_ids = batch['input_ids'].to('cuda')
        attention_mask = batch['attention_mask'].to('cuda')
        labels = batch['labels'].to('cuda')

        # Forward pass
        outputs = model(input_ids, attention_mask)
        loss = criterion(outputs['logits'], labels)

        # Backward pass
        optimizer.zero_grad()
        loss.backward()

        # Gradient clipping
        torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)

        # Update weights
        optimizer.step()
        scheduler.step()

    total_loss += loss.item()

avg_loss = total_loss / len(train_loader)
print(f"Epoch {epoch+1}, Loss: {avg_loss:.4f}")

```

### 3.3 Inference Pipeline

```

def predict(text, model, tokenizer, device='cuda'):
    """Make prediction on single text."""
    model.eval()

```

```

# Preprocess
clean_text = preprocessor.clean_text(text)

# Tokenize
encoding = tokenizer(
    clean_text,
    max_length=512,
    padding='max_length',
    truncation=True,
    return_tensors='pt'
)

input_ids = encoding['input_ids'].to(device)
attention_mask = encoding['attention_mask'].to(device)

# Predict
with torch.no_grad():
    outputs = model(input_ids, attention_mask)
    logits = outputs['logits']
    probs = torch.softmax(logits, dim=1)[0]

# Extract results
prediction = torch.argmax(probs).item()
confidence = probs[prediction].item()

return {
    'prediction': 'depression' if prediction == 1 else 'control',
    'confidence': confidence,
    'probabilities': {
        'control': probs[0].item(),
        'depression': probs[1].item()
    }
}

```

---

## 4. Stage 3: Explainability Framework

### 4.1 Level 1: Token Attribution (Integrated Gradients)

#### Implementation:

```

import torch
from captum.attr import IntegratedGradients

class TokenAttributor:
    def __init__(self, model, tokenizer):
        self.model = model
        self.tokenizer = tokenizer
        self.ig = IntegratedGradients(self._forward_func)

```



```

def _forward_func(self, embeddings):
    """Forward pass with embeddings as input."""
    outputs = self.model.transformer(
        inputs_embeds=embeddings,
        output_hidden_states=False
    )
    cls_output = outputs.last_hidden_state[:, 0, :]
    logits = self.model.classifier(self.model.dropout(cls_output))
    return logits[:, 1] # Depression class logit

def attribute_tokens(self, text, n_steps=20):
    """Compute token attributions using IG."""
    # Tokenize
    encoding = self.tokenizer(
        text,
        max_length=512,
        padding='max_length',
        truncation=True,
        return_tensors='pt'
    )

    input_ids = encoding['input_ids'].to('cuda')

    # Get embeddings
    embeddings =
self.model.transformer.embeddings.word_embeddings(input_ids)

    # Baseline: zero embedding
    baseline = torch.zeros_like(embeddings)

    # Compute attributions
    attributions = self.ig.attribute(
        embeddings,
        baselines=baseline,
        n_steps=n_steps,
        internal_batch_size=1
    )

    # Aggregate over embedding dimension
    attributions = attributions.sum(dim=-1).squeeze(0)

    # Merge subword tokens
    tokens = self.tokenizer.convert_ids_to_tokens(input_ids[0])
    merged_tokens, merged_scores = self._merge_subwords(tokens,
attributions)

    return merged_tokens, merged_scores

def _merge_subwords(self, tokens, scores):
    """Merge WordPiece/BPE subwords back to words."""
    merged_tokens = []
    merged_scores = []

```

```

current_word = ""
current_scores = []

for token, score in zip(tokens, scores):
    # Skip special tokens
    if token in ['<s>', '</s>', '[CLS]', '[SEP]', '<pad>']:
        continue

    # Check if subword continuation
    if token.startswith('##'): # BERT WordPiece
        current_word += token[2:]
        current_scores.append(score.item())
    elif token.startswith('Ġ'): # RoBERTa BPE
        # Save previous word
        if current_word:
            merged_tokens.append(current_word)
            merged_scores.append(sum(current_scores) /
len(current_scores))
        # Start new word
        current_word = token[1:]
        current_scores = [score.item()]
    else:
        # Continuation of current word
        current_word += token
        current_scores.append(score.item())

# Save last word
if current_word:
    merged_tokens.append(current_word)
    merged_scores.append(sum(current_scores) / len(current_scores))

return merged_tokens, merged_scores

```

### Usage Example:

```

attributor = TokenAttributor(model, tokenizer)

text = "I feel worthless and hopeless about everything"
tokens, scores = attributor.attribute_tokens(text)

# Get top-10 important words
top_indices = sorted(range(len(scores)), key=lambda i: scores[i], reverse=True)
[:10]
top_words = [(tokens[i], scores[i]) for i in top_indices]

print("Top Important Words:")
for word, score in top_words:
    print(f" {word}: {score:.3f}")

```

## Output:

```
Top Important Words:
  hopeless: 0.892
  worthless: 0.876
  everything: 0.234
  about: 0.123
  feel: 0.087
  ...
```

## 4.2 Level 2: Symptom Extraction (DSM-5)

### Rule-Based Matcher:

```
import re
from typing import List, Dict

class DSM5SymptomMatcher:
    def __init__(self):
        self.symptom_patterns = {
            'depressed_mood': [
                r'\b(sad|depressed|down|blue|miserable|empty)\b',
                r'\bfeel\s+(bad|terrible|awful|horrible)\b',
                r'\b(crying|tears|weeping)\b'
            ],
            'anhedonia': [
                r'\bno\s+(pleasure|joy|interest|motivation)\b',
                r'\b(nothing|don\'t)\s+(enjoy|care|excites?|interests?)\b',
                r'\b(lost|losing)\s+interest\b',
                r'\bcan\'t\s+enjoy\b'
            ],
            'sleep_disturbance': [
                r'\b(insomnia|can\'t sleep|unable to sleep)\b',
                r'\bsleep\s+(problems?|issues?|difficulties)\b',
                r'\b(hypersomnia|sleeping too much)\b',
                r'\bsleep\s+(\d+)\s+hours\b' # Extract duration
            ],
            'fatigue': [
                r'\b(tired|exhausted|fatigued|drained|no energy)\b',
                r'\black\s+of\s+energy\b',
                r'\bcan\'t\s+get\s+out\s+of\s+bed\b'
            ],
            'worthlessness': [
                r'\b(worthless|useless|failure|burden)\b',
                r'\bno\s+value\b',
                r'\bwaste\s+of\s+space\b'
            ],
            'guilt': [
                r'\b(guilty|ashamed|regret|blame\s+myself)\b',
```

```

        r'\bshould\s+have\b.*\b(better|different)\b'
    ],
    'concentration_difficulty': [
        r'\bcan\'t\s+(focus|concentrate|think)\b',
        r'\b(distracted|foggy|confused)\b',
        r'\btrouble\s+(thinking|deciding|remembering)\b'
    ],
    'psychomotor_changes': [
        r'\b(restless|agitated|fidgety)\b',
        r'\b(slowed|moving slowly|sluggish)\b'
    ],
    'suicidal_ideation': [
        r'\b(suicide|kill myself|end it|not worth living)\b',
        r'\b(death|die|dying)\b.*\b(wish|want|think about)\b',
        r'\bbetter off\s+(dead|without me)\b'
    ]
}

def extract_symptoms(self, text: str) -> List[Dict]:
    """Extract all matching DSM-5 symptoms from text."""
    text_lower = text.lower()
    detected_symptoms = []

    for symptom_name, patterns in self.symptom_patterns.items():
        for pattern in patterns:
            matches = re.finditer(pattern, text_lower, re.IGNORECASE)
            for match in matches:
                # Extract evidence quote
                start = max(0, match.start() - 20)
                end = min(len(text), match.end() + 20)
                evidence = text[start:end].strip()

                detected_symptoms.append({
                    'symptom': symptom_name.replace('_', ' ').title(),
                    'evidence': f"...{evidence}...",
                    'confidence': self._assess_confidence(match.group(),
text_lower)
                })
                break # Only record each symptom once

    return detected_symptoms

def _assess_confidence(self, match_text: str, full_text: str) -> str:
    """Assess confidence based on context."""
    # High confidence: Explicit, strong language
    if any(word in match_text for word in ['suicide', 'hopeless',
'worthless']):
        return 'high'
    # Medium confidence: Clear but less severe
    elif any(word in match_text for word in ['sad', 'tired', 'can\'t']):
        return 'medium'
    # Low confidence: Ambiguous or weak signal
    else:

```

```

        return 'low'

    def compute_phq9_score(self, symptoms: List[Dict]) -> int:
        """Compute PHQ-9 severity score."""
        symptom_scores = {symptom['symptom']: 1 for symptom in symptoms}

        # Adjust for intensity modifiers
        text_combined = ' '.join([s['evidence'] for s in symptoms])
        if re.search(r'\b(always|every day|constantly)\b', text_combined,
re.I):
            multiplier = 3
        elif re.search(r'\b(often|frequently|usually)\b', text_combined, re.I):
            multiplier = 2
        else:
            multiplier = 1

        total_score = sum(symptom_scores.values()) * multiplier
        return min(total_score, 27) # Cap at PHQ-9 maximum

```

### Usage Example:

```

matcher = DSM5SymptomMatcher()

text = "I feel worthless and hopeless. Can't sleep at night. No energy to do
anything."
symptoms = matcher.extract_symptoms(text)
phq9_score = matcher.compute_phq9_score(symptoms)

print(f"Detected Symptoms: {len(symptoms)}")
for symptom in symptoms:
    print(f"  - {symptom['symptom']}: {symptom['evidence']} (Confidence:
{symptom['confidence']})")

print(f"\nPHQ-9 Score: {phq9_score}/27")

```

### Output:

```

Detected Symptoms: 3
  - Worthlessness: ...I feel worthless and... (Confidence: high)
  - Anhedonia: ...hopeless. Can't sleep... (Confidence: high)
  - Sleep Disturbance: ...Can't sleep at night... (Confidence: medium)

PHQ-9 Score: 9/27 (Mild Depression)

```

## 4.3 Level 3: LLM Reasoning Engine

### Prompt Engineering:

```

from openai import OpenAI
import json

class LLMReasoner:
    def __init__(self, api_key, model='gpt-4o'):
        self.client = OpenAI(api_key=api_key)
        self.model = model

    def generate_explanation(self, text, prediction, confidence, symptoms):
        """Generate human-readable clinical explanation."""

        # Build Chain-of-Thought prompt
        prompt = f"""You are a clinical psychology assistant with expertise in
depression assessment.

**Input Text:**
"{text}"

**Classification Result:**
- Prediction: {prediction}
- Confidence: {confidence:.1%}
- Detected DSM-5 Symptoms: {len(symptoms)}

**Chain-of-Thought Reasoning:**

Step 1: Identify Primary Emotions
- Analyze the emotional tone: sadness, hopelessness, anger, numbness, etc.
- Rate emotional intensity (low/medium/high)

Step 2: Map Emotions to DSM-5 Symptoms
- Connect expressed emotions to the 9 DSM-5 criteria for Major Depressive
Disorder
- Assess symptom severity based on language intensity

Step 3: Evaluate Duration and Pervasiveness
- Look for temporal indicators (weeks, months, always, every day)
- Assess if symptoms are situational or persistent

Step 4: Check Crisis Risk
- Identify suicidal ideation or self-harm mentions
- Assess urgency level

Step 5: Generate Evidence-Based Conclusion
- Synthesize all findings into a coherent clinical narrative
- Cite specific text evidence for each claim

**Output Format (JSON):**
```json
{{
  "emotion_analysis": {{
    "primary_emotions": ["emotion1", "emotion2"],

```

```

    "emotional_intensity": "low|medium|high"
  }},
  "symptom_mapping": [
    {
      "symptom": "DSM-5 symptom name",
      "evidence": "exact quote from text",
      "severity": "low|medium|high"
    }
  ],
  "duration_assessment": "estimated duration based on text",
  "crisis_risk": true/false,
  "explanation": "One paragraph (100-150 words) clinical summary grounded in
text evidence",
  "confidence_rationale": "Why this confidence level is appropriate"
}

```

Generate ONLY the JSON object. Ensure all evidence quotes are exact substrings from the input text.

"""

```

# Call LLM API
response = self.client.chat.completions.create(
    model=self.model,
    messages=[
        {"role": "system", "content": "You are a clinical psychology expert
providing evidence-based depression assessments."},
        {"role": "user", "content": prompt}
    ],
    response_format={"type": "json_object"},
    temperature=0.3, # Lower temperature for consistency
    max_tokens=800
)

# Parse JSON response
explanation = json.loads(response.choices[0].message.content)

# Validate evidence grounding
self._validate_evidence(text, explanation)

return explanation

def _validate_evidence(self, original_text, explanation):
    """Ensure all evidence quotes exist in original text."""
    for symptom in explanation.get('symptom_mapping', []):
        evidence = symptom.get('evidence', '')
        if evidence and evidence not in original_text:
            raise ValueError(f"Hallucination detected: '{evidence}' not in
original text")

```

**\*\*Usage Example:\*\***

```
```python
reasoner = LLMReasoner(api_key='your-api-key')

text = "I haven't felt joy in months. Sleep is impossible. Nothing matters anymore."
prediction = "depression"
confidence = 0.88
symptoms = matcher.extract_symptoms(text)

explanation = reasoner.generate_explanation(text, prediction, confidence,
symptoms)

print("LLM Clinical Analysis:")
print(json.dumps(explanation, indent=2))
```

### Output:

```
{
  "emotion_analysis": {
    "primary_emotions": ["sadness", "hopelessness", "numbness"],
    "emotional_intensity": "high"
  },
  "symptom_mapping": [
    {
      "symptom": "Anhedonia",
      "evidence": "haven't felt joy in months",
      "severity": "high"
    },
    {
      "symptom": "Sleep Disturbance",
      "evidence": "Sleep is impossible",
      "severity": "high"
    },
    {
      "symptom": "Hopelessness",
      "evidence": "Nothing matters anymore",
      "severity": "high"
    }
  ],
  "duration_assessment": "Multiple months (meets DSM-5 2-week minimum)",
  "crisis_risk": false,
  "explanation": "Text demonstrates 3 core DSM-5 symptoms of Major Depressive Disorder: anhedonia (loss of pleasure for months), insomnia (impossible sleep), and existential hopelessness. The duration of 'months' satisfies the 2-week diagnostic threshold. Language intensity ('impossible', 'nothing matters') indicates moderate-to-severe depression. No explicit suicidal ideation detected, but persistent hopelessness warrants monitoring. Recommend
```



```
professional psychiatric evaluation.",
    "confidence_rationale": "High confidence (88%) justified by multiple explicit
symptoms with clear evidence quotes and extended duration."
}
```

---

## 5. Safety & Ethics Module

### 5.1 Crisis Detection System

#### Implementation:

```
import re
from typing import Dict, List

class CrisisDetector:
    def __init__(self):
        # High-risk keywords (weight=1.0)
        self.high_risk_keywords = [
            'suicide', 'kill myself', 'end my life', 'not worth living',
            'better off dead', 'want to die', 'plan to die', 'goodbye world'
        ]

        # Medium-risk keywords (weight=0.6)
        self.medium_risk_keywords = [
            'hopeless', 'no point', 'give up', 'can\'t go on', 'burden to
everyone',
            'wish I was dead', 'think about death', 'nothing to live for'
        ]

        # Intent indicators
        self.intent_patterns = [
            (r'\bI will\b', 0.9),      # Strong intent
            (r'\bI\'m going to\b', 0.9),
            (r'\bI plan to\b', 0.8),
            (r'\bI want to\b', 0.6),
            (r'\bI think about\b', 0.4) # Weak intent
        ]

        # Plan specificity indicators
        self.plan_indicators = [
            r'\b(gun|pills|rope|jump|bridge|overdose)\b',
            r'\b(tonight|tomorrow|today|soon)\b',
            r'\b(note|goodbye|last)\b'
        ]

    def detect_crisis(self, text: str) -> Dict:
        """Comprehensive crisis assessment."""
        text_lower = text.lower()
```

```

# Initialize risk score
risk_score = 0.0
detected_keywords = []

# Check high-risk keywords
for keyword in self.high_risk_keywords:
    if keyword in text_lower:
        risk_score += 1.0
        detected_keywords.append((keyword, 'high'))

# Check medium-risk keywords
for keyword in self.medium_risk_keywords:
    if keyword in text_lower:
        risk_score += 0.6
        detected_keywords.append((keyword, 'medium'))

# Assess intent
intent_score = 0.0
for pattern, score in self.intent_patterns:
    if re.search(pattern, text, re.IGNORECASE):
        intent_score = max(intent_score, score)

# Assess plan specificity
plan_score = 0.0
for pattern in self.plan_indicators:
    if re.search(pattern, text_lower):
        plan_score += 0.3
plan_score = min(plan_score, 1.0)

# Combined risk score
total_risk = 0.5 * risk_score + 0.3 * intent_score + 0.2 * plan_score

# Crisis threshold
is_crisis = total_risk > 0.8 or any(severity == 'high' for _, severity
in detected_keywords)

return {
    'is_crisis': is_crisis,
    'risk_score': total_risk,
    'detected_keywords': detected_keywords,
    'intent_level': intent_score,
    'plan_specificity': plan_score,
    'urgency': self._assess_urgency(intent_score, plan_score)
}

def _assess_urgency(self, intent, plan):
    """Categorize urgency level."""
    if intent > 0.8 and plan > 0.5:
        return 'IMMEDIATE' # High intent + specific plan
    elif intent > 0.6 or plan > 0.5:
        return 'HIGH' # Strong intent or plan
    elif intent > 0.4 or plan > 0.3:
        return 'MODERATE'

```

```
else:
    return 'LOW'
```

## 5.2 Hotline Resource Display

```
class CrisisResourceProvider:
    def __init__(self):
        self.hotlines = {
            'US': {
                'name': 'National Suicide Prevention Lifeline',
                'phone': '988',
                'text': 'Text "HELLO" to 741741',
                'url': 'https://988lifeline.org'
            },
            'India': {
                'name': 'AASRA',
                'phone': '91-22-2754-6669',
                'hours': '24/7',
                'url': 'http://www.aasra.info'
            },
            'International': {
                'name': 'Befrienders Worldwide',
                'url': 'https://www.befrienders.org',
                'description': 'Find helplines in 30+ countries'
            }
        }

    def get_resources(self, country='US'):
        """Retrieve crisis hotlines for specific country."""
        return self.hotlines.get(country, self.hotlines['International'])

    def format_crisis_message(self):
        """Generate crisis intervention message."""
        return """
🚨 **CRISIS LANGUAGE DETECTED** 🚨

If you are in immediate danger, please contact emergency services:

**United States:**
- 📞 National Suicide Prevention Lifeline: 988
- 💬 Crisis Text Line: Text "HELLO" to 741741
- 🌐 Online Chat: https://988lifeline.org

**India:**
- 📞 AASRA: 91-22-2754-6669 (24/7)
- 📞 iCall: 91-22-2556-3291
- 🌐 http://www.aasra.info

**International:**
- 🌐 Find a helpline: https://findahelpline.com
```

- 🌐 Befrienders Worldwide: <https://www.befrienders.org>

```
**You are not alone. Help is available.**  
"""
```

---

## 6. Multi-Model Comparison Framework

### 6.1 Ensemble Strategy

```
from typing import List, Dict  
import numpy as np  
  
class ModelEnsemble:  
    def __init__(self, models: List[DepressionClassifier], names: List[str]):  
        self.models = models  
        self.names = names  
  
    def predict_ensemble(self, text: str, tokenizer) -> Dict:  
        """Run all models and aggregate predictions."""  
        predictions = []  
  
        for model, name in zip(self.models, self.names):  
            result = predict(text, model, tokenizer)  
            predictions.append({  
                'model': name,  
                'prediction': result['prediction'],  
                'confidence': result['confidence'],  
                'probabilities': result['probabilities']  
            })  
  
        # Aggregate results  
        depression_votes = sum(1 for p in predictions if p['prediction'] ==  
        'depression')  
        control_votes = len(predictions) - depression_votes  
  
        # Weighted average of probabilities  
        avg_depression_prob = np.mean([p['probabilities']['depression'] for p  
in predictions])  
        avg_control_prob = np.mean([p['probabilities']['control'] for p in  
predictions])  
  
        # Consensus decision  
        consensus = 'depression' if depression_votes > control_votes else  
        'control'  
        agreement_pct = max(depression_votes, control_votes) / len(predictions)  
        * 100  
  
        return {  
            'consensus': consensus,
```

```

        'agreement_percentage': agreement_pct,
        'individual_predictions': predictions,
        'ensemble_probabilities': {
            'depression': avg_depression_prob,
            'control': avg_control_prob
        }
    }
}

```

## 7. Complete Workflow Integration

### 7.1 End-to-End Pipeline

```

class ExplainableDepressionDetector:
    def __init__(self, model_path, tokenizer_name, api_key):
        # Initialize components
        self.preprocessor = TextPreprocessor()
        self.tokenizer = AutoTokenizer.from_pretrained(tokenizer_name)
        self.model = DepressionClassifier.from_pretrained(model_path)
        self.attributor = TokenAttributor(self.model, self.tokenizer)
        self.symptom_matcher = DSM5SymptomMatcher()
        self.llm_reasoner = LLMReasoner(api_key)
        self.crisis_detector = CrisisDetector()
        self.resource_provider = CrisisResourceProvider()

    def analyze(self, text: str) -> Dict:
        """Complete analysis pipeline."""
        # Stage 0: Crisis check (priority)
        crisis_result = self.crisis_detector.detect_crisis(text)

        if crisis_result['is_crisis']:
            return {
                'status': 'CRISIS_DETECTED',
                'crisis_info': crisis_result,
                'resources': self.resource_provider.format_crisis_message(),
                'prediction_blocked': True
            }

        # Stage 1: Preprocessing
        clean_text = self.preprocessor.clean_text(text)

        # Stage 2: Classification
        prediction_result = predict(clean_text, self.model, self.tokenizer)

        # Stage 3: Explainability
        # Level 1: Token attribution
        tokens, scores = self.attributor.attribute_tokens(clean_text)
        top_tokens = sorted(zip(tokens, scores), key=lambda x: x[1],
                             reverse=True)[:10]

```

```

# Level 2: Symptom extraction
symptoms = self.symptom_matcher.extract_symptoms(clean_text)
phq9_score = self.symptom_matcher.compute_phq9_score(symptoms)

# Level 3: LLM reasoning
llm_explanation = self.llm_reasoner.generate_explanation(
    text=clean_text,
    prediction=prediction_result['prediction'],
    confidence=prediction_result['confidence'],
    symptoms=symptoms
)

# Compile final report
return {
    'status': 'SUCCESS',
    'prediction': prediction_result,
    'explainability': {
        'token_attribution': top_tokens,
        'symptoms': symptoms,
        'phq9_score': phq9_score,
        'llm_reasoning': llm_explanation
    },
    'safety': {
        'crisis_detected': False,
        'low_confidence_warning': prediction_result['confidence'] < 0.7
    }
}

```

### Usage:

```

detector = ExplainableDepressionDetector(
    model_path='models/trained/roberta-base',
    tokenizer_name='roberta-base',
    api_key='your-openai-api-key'
)

text = "I feel hopeless and worthless. Can't sleep, no energy. What's the point?"
result = detector.analyze(text)

print(f"Prediction: {result['prediction']['prediction']}")
print(f"Confidence: {result['prediction']['confidence']:.1%}")
print(f"\nTop Important Words:")
for word, score in result['explainability']['token_attribution'][:5]:
    print(f"  {word}: {score:.3f}")
print(f"\nPHQ-9 Score: {result['explainability']['phq9_score']}/27")
print(f"\nLLM Explanation: {result['explainability']['llm_reasoning']['explanation']}")

```

## 8. Summary

This methodology integrates:

1. ☒ **Preprocessing**: Text cleaning + tokenization
2. ☒ **Classification**: Fine-tuned transformers (88% accuracy)
3. ☒ **Explainability**: 3-level hierarchy (IG + DSM-5 + LLM)
4. ☒ **Safety**: Crisis detection + hotlines
5. ☒ **Ethics**: Non-diagnostic language + confidence warnings

**Key Innovation:** First mental health NLP system combining Integrated Gradients, clinical symptom extraction, and LLM reasoning in a unified framework.

---

[← Back to Mathematical Modeling](#) | [Next: Experiments →](#)