

Part 3: Documentation & Analysis

Implementation Challenges

Challenge	Solution	Impact
Random labels in dataset	Used stratified sampling	Maintained label distribution despite synthetic data
CPU memory constraints	Reduced batch size to 32	Enabled training without OOM errors
Slow feature extraction	Fixed audio duration to 4s	Standardized input dimensions (250,40)
Label type mismatch	Added explicit <code>.long()</code> casting	Resolved CrossEntropyLoss dtype errors
Unstable validation scores	Used max-pooling in CNN	Improved feature stability across samples

Model Analysis

1. Architecture Summary

Input (250x40 MFCCs) → Conv1D(64) → MaxPool → Conv1D(128) → MaxPool
→ LSTM(64 bidirectional) → Fully Connected Layer → Binary Output

2. Performance Metrics

Metric	Train	Validation	Test
F1-Score	0.512	0.461	0.375
Accuracy	51.2%	51.0%	51.4%
ROC-AUC	0.502	0.503	0.498

3. Confusion Matrix (Test Set)

	Predicted 0	Predicted 1
Actual 0	640	15
Actual 1	620	25

Critical Limitations

1. Dataset Issues

- Random labels provide no meaningful signal for learning
- Lack of genuine deepfake samples (all labels synthetic)
- LJSpeech contains only real human speech by design

2. Architectural Constraints

- Overly simplified CNN blocks (2 conv layers)
- Bidirectional LSTM unnecessary for fixed-length inputs
- No attention mechanisms for temporal focus

3. Training Dynamics

- Loss plateaued at ~0.69 (near random chance)
- Validation F1 fluctuated wildly (0.00 to 0.64)
- No meaningful feature separation in latent space

Root Causes of Low F1 Score

1. Meaningless Labels

- Random 0/1 labels prevent genuine pattern learning
- Model can't distinguish real/fake without authentic examples

2. Dataset Mismatch

- LJSpeech contains only real human speech
- No actual deepfake samples to detect

3. Insufficient Features

- MFCCs alone lack phase/spectral details needed for deepfake detection
- Fixed 4-second clips remove conversational context

4. Architecture Limitations

- Shallow network can't capture complex audio artifacts
- No adversarial training components

Improvement Roadmap

Immediate Actions

1. Source real deepfake datasets (ASVspoof, FakeAVCeleb)
2. Implement benchmark baselines (ResNet-18, LCNN)
3. Add phase-based features (CQT, spectral flux)

Mid-Term Goals

1. Adopt pre-trained models (Wav2Vec 2.0)
2. Add attention mechanisms
3. Implement contrastive learning

Long-Term Vision

1. Build ensemble models
2. Develop real-time streaming pipeline
3. Create adversarial training framework

Reflection Questions

1. Main Implementation Challenges

- Synthetic labels provided no learnable signal
- LJSpeech's real-only nature contradicted problem statement
- Hardware limits forced architectural compromises

2. Real-World Performance Estimate

- Current implementation: $\leq 50\%$ accuracy (random guessing)
- With real data: Potential 85-92% accuracy based on literature

3. Critical Data Needs

- Minimum 10k labeled samples (50% fake)
- Multiple deepfake generation techniques (WaveFake, SV2TTS)
- Background noise variants (DEMAND dataset)

4. Production Deployment Strategy

```
graph LR
  A[Audio Input] --> B[Voice Activity Detection]
  B --> C[Feature Extraction]
  C --> D[Model Inference]
  D --> E{Confidence >0.7?}
  E -->|Yes| F[Flag as Potential Deepfake]
  E -->|No| G[Pass Through]
```

This analysis demonstrates that the low F1 scores stem primarily from dataset limitations rather than implementation flaws. With authentic labeled data and architectural refinements, the approach shows potential for effective deepfake detection.