# India Speaks Preprocessing Pipeline - Design Report

## 1. Modular Architecture

- Decision: Separated concerns into distinct modules (text normalization, audio validation, quality filtering, CLI)

- Rationale: Enables independent testing, easier maintenance, and extensibility for new languages/features

- Trade-off: Slightly more complex structure vs. single monolithic file

## 2. Stub-Based Audio Validation

- Decision: Implemented audio validation stubs without actual audio processing

- Rationale: Sample data contained only S3 URLs without actual audio files for validation

- Impact: Pipeline validates metadata but cannot detect actual audio corruption or quality issues

## 3. Conservative Quality Filtering

- **Decision**: Aggressive filtering with language path mismatch detection as primary filter

- **Rationale**: Prioritizes high-quality training data over quantity (6% acceptance rate)

- **Trade-off**: Rejects many potentially valid samples to ensure data consistency

## 4. Multi-Language Text Normalization

- **Decision**: Single normalizer class supporting 12+ languages with language-specific rules

- **Rationale**: Centralized logic for consistent processing across all Indic languages

- **Limitation**: Uses basic heuristics rather than ML-based language detection

### 5. CLI-First Design

- **Decision**: Both Python API and command-line interface with simple wrapper script

- **Rationale**: Supports both programmatic integration and standalone usage

- **Implementation**: Fallback to Python import method due to Windows entry point issues

## Known Limitations:

## 1. Audio Processing Constraints

- Limitation: No actual audio file validation (stubs only)

- Impact: Cannot detect corrupted audio, actual duration mismatches, or quality issues

- Mitigation: Relies on CSV metadata for duration/quality validation

## 2. Language Detection Accuracy

- Limitation: Basic script validation using Unicode ranges rather than sophisticated language detection

- Impact: May incorrectly classify mixed-script text or transliterated content

- Example: Hindi text in Latin script would be misclassified

## 3. Processing Performance

- Limitation: Single-threaded processing (~2,300 samples/sec)

- Impact: Large datasets (>100K samples) may require significant processing time

- Potential: Could be parallelized for 3-5x speed improvement

## 4. Data Quality Assessment

- Limitation: Simple quality scoring based on word count, repetition, and basic heuristics

- Impact: May miss subtle quality issues or domain-specific requirements

- Example: Technical terminology might be incorrectly flagged as low quality

## 5. Memory Usage

- Limitation: Loads entire dataset into memory for processing

- Impact: Large CSV files (>1GB) may cause memory issues

- Constraint: Not optimized for streaming/batch processing

## Current Performance Metrics

- **Processing Speed**: ~2,300 samples/second

- **Memory Usage**: ~50MB per 10K samples

- **Acceptance Rate**: 6% (highly conservative filtering)

- **Primary Rejection**: Language path mismatch (90% of rejections)

## Recommended Improvements

- **Audio Integration**: Implement actual audio validation with torchaudio

- **Parallel Processing**: Add multiprocessing for 3-5x speed improvement

- **Streaming Support**: Process large datasets in chunks to reduce memory usage

- **ML Language Detection**: Replace heuristic-based language validation

- **Quality Scoring**: Implement more sophisticated text quality assessment

## Design Philosophy

The pipeline prioritizes data quality over quantity, using conservative filtering to ensure high-quality training data for ASR/TTS models. The modular design enables future enhancements while maintaining production stability.