

# MPEG Sorter Technical Documentation

Copyright © 2025 Eric Gitonga. All rights reserved.  
This document is licensed under the MIT License.

A Python utility that identifies and sorts media files based on their actual file signatures rather than extensions. Accurately separates MP3 and MP4 files into appropriate directories, correcting mislabeled extensions in the process.

## Implementation Details

### File Structure

`mpeg_sorter.py` - Main Python script

### Dependencies

The implementation relies solely on Python standard library modules: - `os`: For file system operations and CPU core detection - `shutil`: For moving files - `argparse`: For command-line argument parsing - `pathlib`: For path manipulation - `asyncio`: For asynchronous execution framework - `concurrent.futures`: For parallel processing with `ThreadPoolExecutor` - `time`: For performance measurement

### Core Components

1. **File Signature Detection (`get_file_signature`):**
  - Reads the first 12 bytes of each file
  - Compares against known MP3 and MP4 signatures
  - Returns the identified file type ('mp3', 'mp4', or 'unknown')
2. **File Sorting Logic (`sort_files` and `sort_files_async`):**
  - Creates appropriate subdirectories
  - Processes files in the source directory (non-recursively)
  - Uses asynchronous parallel processing with `ThreadPoolExecutor`
  - Identifies file types using signatures
  - Handles extension mismatches
  - Moves files to destination folders
  - Resolves filename conflicts
  - Provides real-time progress updates
3. **Command-line Interface (`main`):**
  - Parses command-line arguments
  - Validates inputs
  - Provides help information
  - Includes worker count configuration option
  - Offers sequential mode for performance benchmarking
  - Creates unknown folder by default (can be disabled with `--no-unknown`)

- Calls the appropriate sorting function based on chosen mode

## Signature Detection

The script uses the following signatures for detection:

### MP3 Signatures

- \xFF\xFB: MPEG-1 Layer 3
- \xFF\xF3: MPEG-2 Layer 3
- \xFF\xF2: MPEG-2.5 Layer 3
- \x49\x44\x33: ID3 tag (common in MP3)

### MP4 Signatures

- \x00\x00\x00\x18\x66\x74\x79\x70: ISO Base Media file (MPEG-4)
- \x00\x00\x00\x20\x66\x74\x79\x70: ISO Base Media file (MPEG-4)
- \x66\x74\x79\x70\x4D\x53\x4E\x56: MPEG-4 video
- \x66\x74\x79\x70\x69\x73\x6F\x6D: ISO Base Media file (MPEG-4)

Additional check for 'ftyp' after the initial 4 bytes is also performed for MP4 container formats.

## Analysis of Strengths and Weaknesses

### Strengths

1. **Signature-based Detection**
  - More reliable than extension-based sorting
  - Can identify mislabeled files with high accuracy
  - Minimal false positives for supported formats
2. **Efficient Implementation**
  - Minimal memory footprint (reads only first 12 bytes of each file)
  - No external dependencies required
  - Files are moved rather than copied, saving disk space
3. **User Experience**
  - Clear feedback for each operation
  - Handles edge cases like duplicate filenames
  - Intuitive command-line interface
4. **Safety Features**
  - Validates source directory existence
  - Creates subdirectories as needed
  - Exception handling for file operations

### Weaknesses

1. **Limited Format Support**
  - Currently only handles MP3 and MP4 formats

- No support for other audio formats (WAV, FLAC, OGG, etc.)
  - No support for other video formats (AVI, MKV, WebM, etc.)
2. **Basic Detection Method**
    - Relies on a small set of signatures for each format
    - May miss some valid files with unusual headers
    - No advanced analysis of file contents beyond headers
  3. **No Recursive Processing**
    - Only processes files in the top level of the source directory
    - Ignores files in subdirectories
  4. **Asynchronous Parallel Processing**
    - Uses ThreadPoolExecutor for concurrent file processing
    - Dynamically scales to available CPU cores for optimal performance
    - Includes progress indication for monitoring large operations
    - Significantly improves processing speed for large collections

## Future Improvements

1. **Expanded Format Support**
  - Add support for additional audio formats (FLAC, WAV, AAC, OGG)
  - Add support for additional video formats (MKV, AVI, WebM, FLV)
  - Create a more extensible signature database
2. **Enhanced Detection**
  - Implement more comprehensive signature detection
  - Add secondary verification methods for ambiguous files
  - Consider using content-based analysis for better accuracy
3. **Performance Enhancements**
  - Add multi-threading support for faster processing
  - Implement progress bars for large operations
  - Add options for batch processing
4. **Additional Features**
  - Add recursive directory processing
  - Add support for custom destination directories
  - Implement dry-run mode to preview changes
  - Add metadata extraction for better organization
  - Create file format reports and statistics
5. **User Interface Improvements**
  - Add a simple GUI option
  - Implement color-coded console output
  - Add interactive mode for uncertain file types
6. **Testing and Validation**
  - Add unit tests for reliability
  - Create test files for various edge cases
  - Implement validation for sorted files

## Implementation Considerations

For future development, consider:

1. **Modularity:** Separate concerns into distinct modules for better maintainability:
  - Signature detection
  - File operations
  - CLI handling
  - Logging/reporting
2. **Configuration:** Add a configuration file for customizable settings:
  - Custom signature definitions
  - Default behavior options
  - Custom destination paths
3. **Performance:** For larger collections, consider:
  - Batch processing to reduce system calls
  - Memory-mapped file access for faster signature detection
  - Worker pools for parallel processing
4. **Integration:** Consider adding hooks for:
  - Media library management systems
  - Backup solutions
  - Automated workflows

By addressing these improvements, MPEG Sorter could evolve into a more comprehensive media management utility suitable for various use cases, from personal collections to professional media libraries.