# Music Identification Using Compression

Authors

Inês Baptista (98384)

João Correia (104360)

Daniel Ferreira (102885)

# 1 Indice

# Contents

# 2    Introduction

Music identification is a crucial task in various applications, such as music recommendation systems, music classification, and music retrieval. Traditional methods for music identification rely on audio features, such as spectrograms, mel-frequency cepstral coefficients, and chroma features. This project aims to develop an Information Theoretical approach to music identification using **compression** techniques.

The main idea is to use the Normalized Compression Distance (NCD) to measure the **similarity** between two audio segments. The NCD is an approximation of the non-computable Normalized Information Distance (NID), which is defined as the maximum of the Kolmogorov complexities of the two segments given each other. To avoid the non-computability of the Kolmogorov complexity, the NCD uses the approximation:

$$\text{NCD}(x, y) = \frac{C(x, y) - \min\{C(x), C(y)\}}{\max\{C(x), C(y)\}} \tag{1}$$

where $C(x)$ denotes the number of bits required by compressor $C$ to represent $x$ and $C(x, y)$ indicates the number of bits needed to compress $x$ and $y$ together (usually, the two strings are concatenated). Distances close to one indicate dissimilarity, while distances near zero indicate similarity.

In this report, we will detail the steps needed to test our approach.

## 2.1    Music identification

Suppose that a certain database contains representations of several (complete) musics, denoted $m_i$, and that we want to identify a segment of music, represented by string $x$. The idea is to compute the values of $\text{NCD}(x, m_i)$ for all $i$, and assign to $x$ the name of the music in the database that yields the smallest value of the NCD.
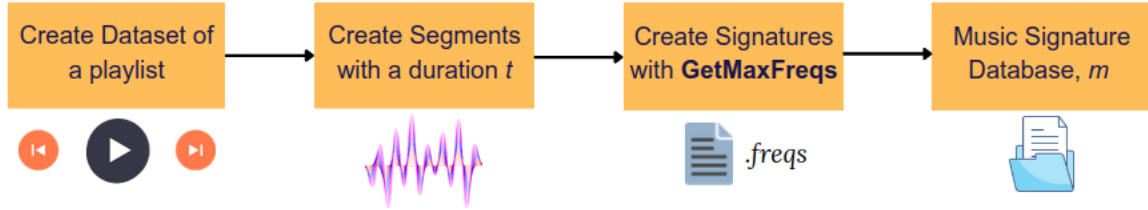
# 3    Methodology



Figure 1: Stage 1: Generate the database $m_i$

## 3.1    Stage 1: Generate the music signature database $m_i$

In the figure 2 below, we see the main steps to generate the database with the signatures that will be used to idenity music within audios.

1. **Playlist Generation**

   The dataset used for this project consists of 25 different musics, downloaded from YouTube. The **create_dataset.py** script automates the process of downloading audio from YouTube videos and playlists, converting the audio to a specified format (`mp3` or `wav`) with desired sample rate, bit depth, and number of channels, and saving it to a designated output directory. It supports batch processing through command-line arguments and input files containing URLs.

2. **Segments Creation**

   Segmenting the original audio files is essential because applications like *Shazam* require only a short duration of audio to accurately identify the music. The **create_audio_segment.py** script is designed to create segments of audio files with specified durations $t$ from given audio files or directories containing audio files. It takes input arguments such as the duration of the segment in seconds, the minimum start time of the segment, and the output path. The script utilizes the `SoX` (Sound eXchange) library to manipulate audio files. The start time for each segment is either randomly chosen or in a specified range.

3. **Signature Generation**

   The audio undergoes a transformation to a representation more suitable for the intended general-purpose compressors, using the provided, by the professors, **GetMaxFreqs** program in `C++`. This process involves segmenting the audio, partially overlapping, and computing the most significant frequencies for each segment. The **GetMaxFreqs** program reads the audio data, converts it to mono, downsamples it, computes the Fast Fourier Transform (FFT) over sliding windows, identifies the top frequencies by power, and saves these frequencies to an output file `.freqs`.

## 3.2 Stage 2: Identify a music $x$ given an audio $x$
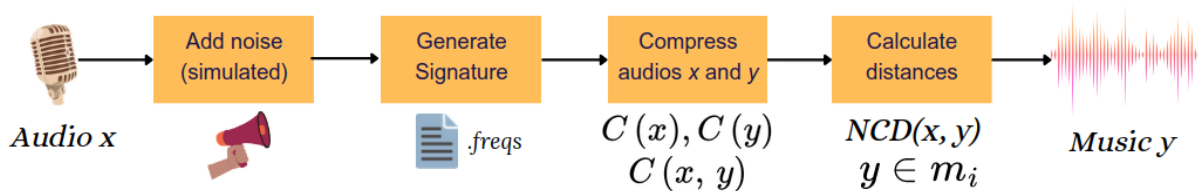


Figure 2: Stage 2: Identify a music $y$ given an audio $x$ containing a music and simulated noise on the background.

1. **Simulating Noise:** To assess the reliability of our music identification system, noise was added to the audio segments used for querying the database. The `create_noise.py` script is utilized to generate different types of noise and add them to the audio segments. Various types of noise can be added, including:

   - **White noise:** Random noise with equal intensity at different frequencies.
   - **Pink noise:** Noise with more energy in lower frequencies, similar to many natural sounds

- **Brown noise:** Noise with even more energy in lower frequencies compared to pink noise.
- **Blue noise:** Noise with more energy in higher frequencies.
- Any **YouTube video**.

The **intensity** of the added noise can be controlled to create different signal-to-noise ratios, allowing for testing the music identification technique under various noise conditions.

2. **Signature Generation**

The audio undergoes a transformation to a representation more suitable for the intended general-purpose compressors, using the provided, by the professors, **GetMaxFreqs** program in `C++`. This process involves segmenting the audio, partially overlapping, and computing the most significant frequencies for each segment. The **GetMaxFreqs** program reads the audio data, converts it to mono, downsamples it, computes the Fast Fourier Transform (FFT) over sliding windows, identifies the top frequencies by power, and saves these frequencies to an output file `.freqs`.

3. **Compreesion**

To compute the Normalized Compression Distance (NCD), we compress both $x$ and $y$ separately to determine their individual compressibility, denoted as $C(x)$ and $C(y)$ respectively. Additionally, we compress their concatenation of $x+y$ to capture shared patterns, represented by $C(x, y)$. These lengths of each compression serve as the basis for computing the NCD, indicating the similarities or dissimilarities between the two entities based on the information extracted through compression, also known as **entropy**. For this task, we used the standard compressors **gzip**, **bz2**, **lzma**, **zstd**, **zlib**, **lz4**, **lz4**, and **snappy**. Later in the analysis, we analyze which one provided the best results.

4. **NCD Calculation:** The database is queried using small samples of the music pieces. For each query, the Normalized Compression Distance (NCD) is calculated, and the music in the database yielding the smallest NCD value is assigned to the query. The NCD computation is performed for each pair of audio segments in the database, serving as a metric to gauge the similarity between the segments.

## 3.3 Automations with pipeline execution

Both stages can be automated through a **pipeline** with custom made `YAML` files. These files define the steps to be taken. For example, a `YAML` file is created with the structure shown in Figure 3.

The pipeline is then executed using the `src/pipeline.py` script. It is important to note that steps within the pipeline can be excluded to avoid repetition. For example, `create_dataset.py` only needs to be run once, as does the creation of original signatures. Furthermore, as depicted in Figure 3, the process includes the calculation of all distances, essentially forming a Cartesian product. However, this might not always be necessary and the `YAML` file can be tailored based to just identify one segment.

When running the pipeline multiple times, be cautious of calculating results for duplicate files. It is recommended to delete directories that are not overwritten, such as those containing audio segments. There is an option to automate this process or use the `src/clean.py` script. For initial setup, run the `src/bash/install.sh` script to install the necessary libraries.

```
steps:
  - script: src/preprocessing/create_dataset.py
    args:
      file_paths: ["data/playlists.txt"]
      output_path: "data/music"
      audio_format: "wav"
      sample_rate: 44100
      bits_per_sample: 16
      channels: 2
  - script: src/preprocessing/create_segments.py
    args:
      __NO_ARG_NAME__paths: ["data/music"]
      duration: 5
      min_time: 60
      start_time: null
      output_path: "data/segments"
  - script: src/preprocessing/create_noise.py
    args:
      __NO_ARG_NAME__paths: ["data/segments"]
      output_path: "data/noise"
      noise_type: "white"
      intensity: 0.3
  - script: src/preprocessing/create_signatures.py
    args:
      __NO_ARG_NAME__paths: ["data/music"]
      output_path: "data/signatures/original"
      signature-args: "-ws 1024 -sh 256 -nf 8"
      signature_type: "gmf"
  - script: src/preprocessing/create_signatures.py
    args:
      __NO_ARG_NAME__paths: ["data/noise"]
      output_path: "data/signatures/segments"
      signature_type: "gmf"
      signature_args: ""
  - script: src/main/create_distance_results.py
    args:
      __NO_ARG_NAME__paths: ["data/signatures/segments"]
      database_path: "data/signatures/original"
      algorithm: "bz2"
      output_path: "data/distances/{algorithm}/results.csv"
```

Figure 3: Pipeline Configuration YAML file for automation purposes

## 3.4  Enhancing Performance

### 3.4.1  With Multiprocessing

To improve the efficiency and performance of the audio signature generation process, we incorporated the modules `Pool` and `cpu_count` from the `multiprocessing` library in `Python`.

**Creating Signatures in Parallel:**  The `create_gmf_signatures` function sets up a multiprocessing pool to generate signatures in parallel. It uses `Pool` to map the `generate_signature` function to the list of audio file paths:

Listing 1: Python code for the 'create_gmf_signatures' function

```python
def create_gmf_signatures(paths, output_path, args, verbose=False):
    os.makedirs(output_path, exist_ok=True)

    if not check_dependencies() or not compile_get_max_freqs():
        return

```

5

```python
7      with Pool(cpu_count()) as pool:
8          tasks = [(path, output_path, args) for path in paths]
9          results = pool.map(generate_signature, tasks)
10
11     if verbose:
12         for path, result in zip(paths, results):
13             if result:
14                 print(f"Generated signature for {path}")
15             else:
16                 print(f"Failed to generate signature for {path}")
```

We define a function `generate_signature` that generates the signature for a single audio file. This function will be used by the `Pool` object to process files concurrently:

Listing 2: Python code for the 'generate_signature' function

```python
1  def generate_signature(args):
2      path, output_path, signature_args = args
3      output_file = os.path.join(
4          output_path,
5          os.path.basename(path).rsplit('.', 1)[0] + ".freqs")
6      signature_args = signature_args.split()
7      if subprocess.run([
8              "GetMaxFreqs/bin/GetMaxFreqs",
9              "-w",
10             output_file
11             ]
12             + signature_args + [path]
13         ).returncode != 0:
14         print(f"Failed to generate signature for {path}")
15         return False
16     return True
```

### 3.4.2   With Caching

Listing 3: Python code for the 'get_compressed_length' function

```python
1  def get_compressed_length(algorithm, data, cache, cache_key):
2      if cache_key not in cache:
3          cache[cache_key] = len(compress_file(algorithm, data))
4      return cache[cache_key]
```

# 4   Results

## 4.1   Parameter Tuning

The program **GetMaxFreqs** generates signatures with adjustable parameters. The **Window Size (-ws)** determines FFT window size, impacting frequency resolution and computational cost. The **Window Overlap (-sh)** controls sample overlap, enhancing temporal resolution. The **Number of Significant Frequencies (-nf)** limits considered frequencies, balancing computational efficiency with information retention.

In the Figure 4, the accuracy is higher when $ws = 2560$, $sh = 0.0625$ and $nf = 3$.
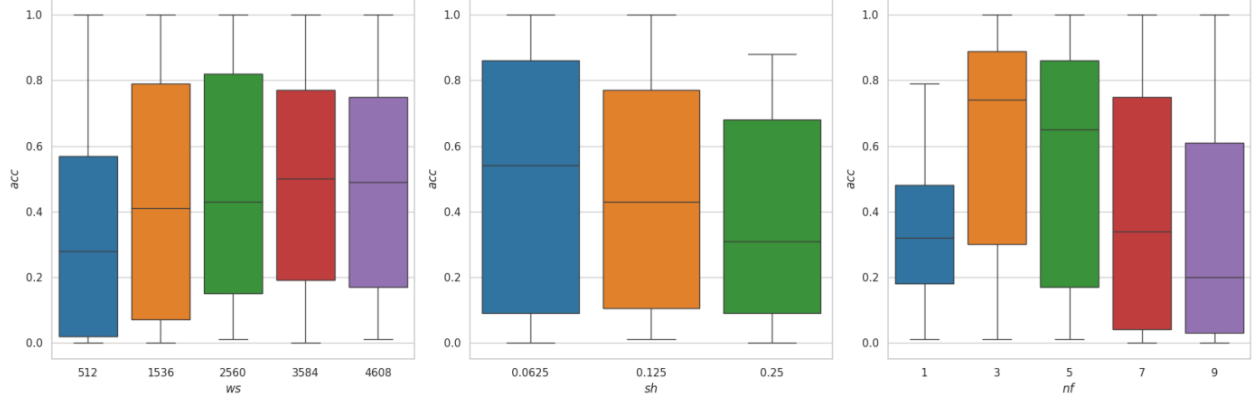
Figure 4: Accuracy results in function of each of the parameters **ws**, **sh** and **nf**.
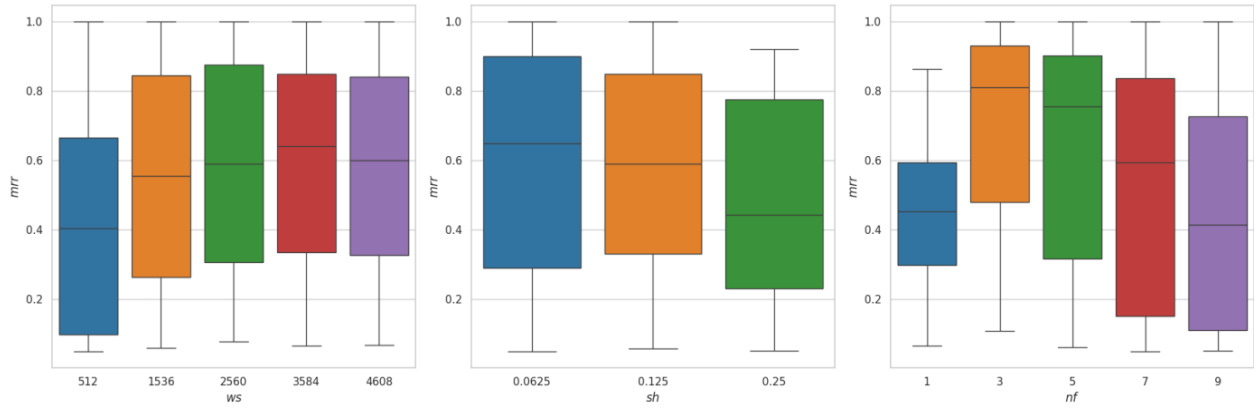


Figure 5: Mean Reciprocal Rank (MRR) results in function of each of the parameters **ws**, **sh** and **nf**.

In the Figure 5, Mean Reciprocal Rank (MRR) is also higher when ws = 2560, sh = 0.0625 and nf = 3. A high MRR (Mean Reciprocal Rank) means that the relevant items are ranked very high in the list of returned results and follows the formula:

$$\text{MRR} = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{\text{rank}_i} \tag{2}$$

### 4.1.1 Performance Evaluation

In the Figure 6, we see the time for different parameter combinations. Each time meaning:

- **t1** is the time to generate the signatures of the complete songs.

- **t2** is the time to generate the signatures of the segments.

- **t3** is the time to calculate the distances.

From the figure, we see that the parameters don't have an impact of **t2** - which makes sense has the segments have all the same duration. Time **t1** increases as **sh** increases and **ws** and **nf**
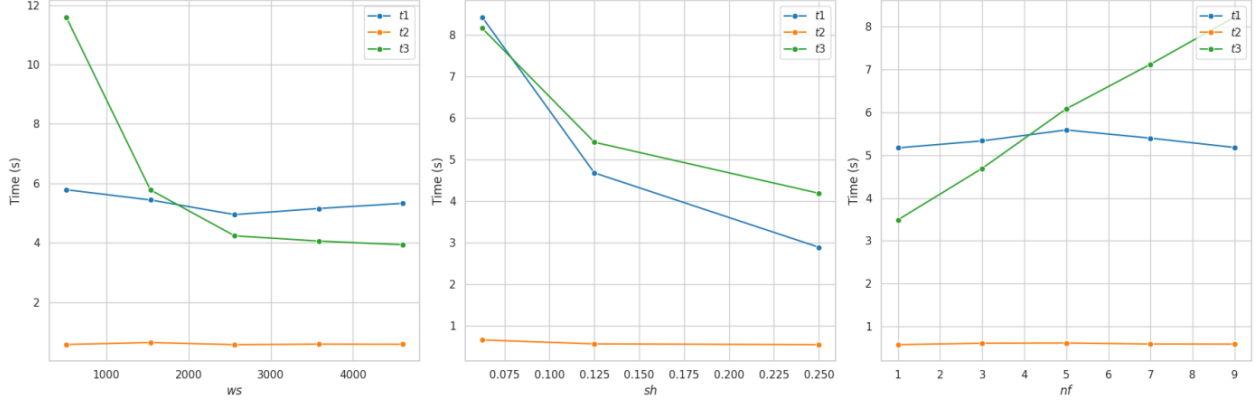
7

Figure 6: Time *t1, t2 and t2* in function of the parameters **ws**, **sh** and **nf**.

also don't have a big impact on this timing. The time to calculate the distances, **t3**, is influenced by all the parameters. It is inversely proportional to **ws** and **sh** (decreases as **ws** and **sh** decrease) and directly proportional to **nf** (increases as **nf** increases).

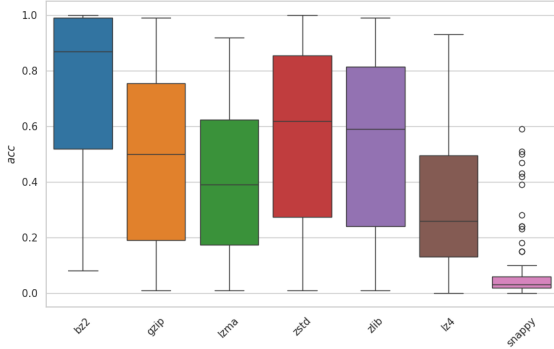## 4.2    Compressors comparison



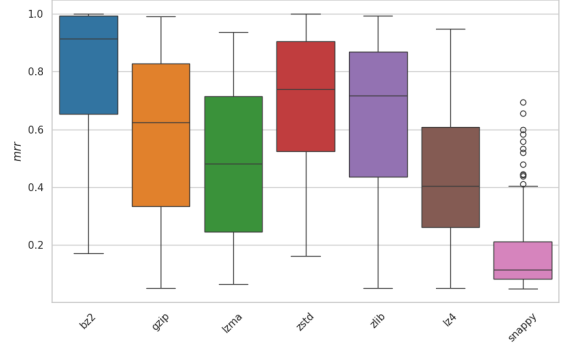Figure 7: Accuracy results in function of each of the compressors.



Figure 8: Mean Reciprocal Rank (MRR) results in function of each of the compressors.

We also compare different compressors - **gzip**, **bz2**, **lzma**, **zstd**, **zlib**, **lz4**, **lz4**, and **snappy** to find the one who retrieves higher accuracy.

In Figures 7 and 10, the accuracy and Mean Reciprocal Rank (MRR) is higher for the compressor **bz2**. The second best is the compressor **zstd**.

## 4.3    Best combinations

To identify optimal parameter combinations, we used a composite metric balancing accuracy and computational efficiency. We computed the total processing time (sum of *t1*, *t2*, and *t3*) and normalized it to a range of 0 to 1, using min-max scaling.

The composite metric, a weighted sum of accuracy (60%) and normalized total time (40%), prioritizes high accuracy while penalizing excessive computation time. We segmented the dataset

by compression method, filtered relevant parameters and metrics, and sorted by the composite metric to identify top-performing combinations.

- **Window Sizes (ws)**: {512, 1536, 2560, 3584, 4608}

- **Window Overlap (sh)**: { $\frac{1}{16}$, $\frac{1}{8}$, $\frac{1}{4}$ of ws }

- **Number of Significant Frequencies (nf)**: {1, 3, 5, 7, 9}

- **Compressors**: {'bz2', 'gzip', 'lzma', 'zstd', 'zlib', 'lz4', 'snappy'}

The best combinations for each compressor are shown in Table 1.

Table 1: Combinations and Results

| Combination | $ws$ | $sh$ | $nf$ | Compressor | $acc$ | $MRR$ | $t1$ | $t2$ | $t3$ |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| 1 | 1536 | 192 | 5 | bz2 | 1.0 | 1.0 | 4.8069 | 0.6107 | 5.2554 |
| 2 | 2560 | 160 | 3 | gzip | 0.99 | 0.9909 | 7.2340 | 0.5786 | 2.2616 |
| 3 | 2560 | 160 | 7 | lzma | 0.87 | 0.8997 | 7.3348 | 0.6003 | 22.2095 |
| 4 | 2560 | 160 | 3 | zstd | 1.0 | 1.0 | 7.2340 | 0.5786 | 2.7799 |
| 5 | 2560 | 160 | 3 | zlib | 0.99 | 0.9925 | 7.2340 | 0.5786 | 3.3078 |
| 6 | 4608 | 288 | 5 | lz4 | 0.93 | 0.9482 | 8.4333 | 0.6347 | 1.4662 |
| 7 | 2560 | 160 | 1 | snappy | 0.59 | 0.6942 | 8.1086 | 0.6971 | 1.4702 |

## 4.4 Performance Evaluation

Adding noise to our audios aids to establish if our identifier will perform well identifying musics in the real world as

## 4.5 Robustness Assessment

The robustness of the technique will be assessed by adding noise to the segments used for querying.

# 5 Conclusion

The project aims to develop a novel approach to music identification using compression techniques. The results show that the NCD-based approach is effective in identifying musics. The robustness of the technique is also assessed by adding noise to the segments used for querying.
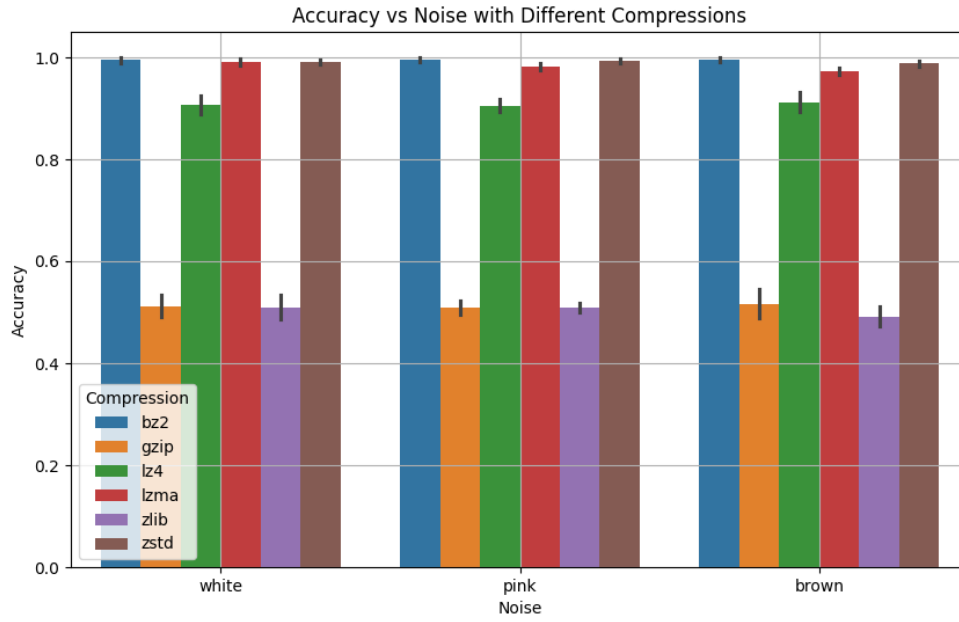
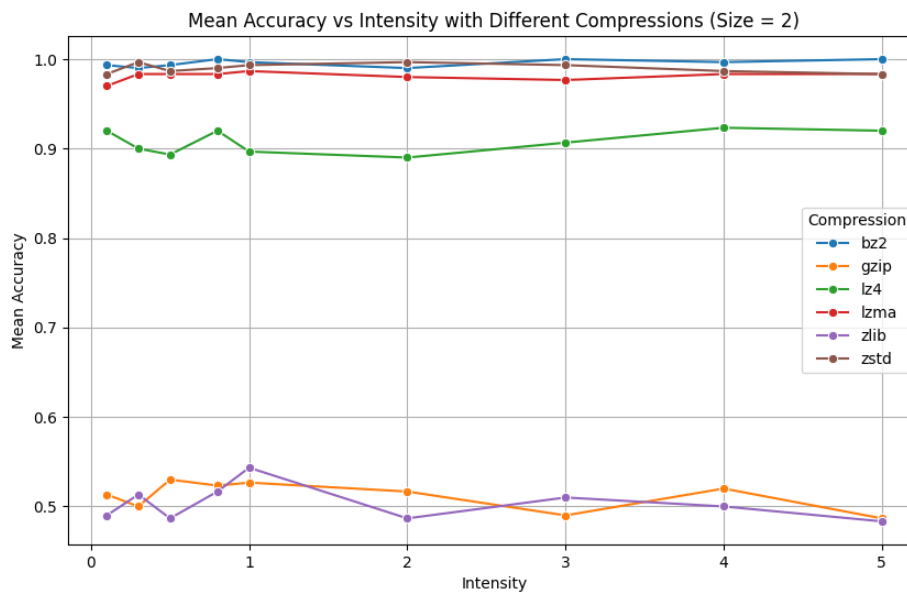Figure 9: Mean Reciprocal Rank (MRR) results in function of each of the compressors.



Figure 10: Mean Reciprocal Rank (MRR) results in function of each of the compressors.

# Annex

```yaml
steps:
  - script: src/preprocessing/create_dataset.py
    args:
      file_paths: ["data/playlists.txt"]
      output_path: "data/music"
      audio_format: "wav"
      sample_rate: 44100
      bits_per_sample: 16
      channels: 2
  - script: src/preprocessing/create_segments.py
    args:
      __NO_ARG_NAME__paths: ["data/music"]
      duration: 5
      min_time: 60
      start_time: null
      output_path: "data/segments"
  - script: src/preprocessing/create_noise.py
    args:
      __NO_ARG_NAME__paths: ["data/segments"]
      output_path: "data/noise"
      noise_type: "white"
      intensity: 0.3
  - script: src/preprocessing/create_signatures.py
    args:
      __NO_ARG_NAME__paths: ["data/music"]
      output_path: "data/signatures/original"
      signature_type: "gmf"
  - script: src/preprocessing/create_signatures.py
    args:
      __NO_ARG_NAME__paths: ["data/noise"]
      output_path: "data/signatures/segments"
      signature_type: "gmf"
      signature_args: ""
  - script: src/main/create_distance_results.py
    args:
      __NO_ARG_NAME__paths: ["data/signatures/segments"]
      database_path: "data/signatures/original"
      algorithm: "bz2"
      output_path: "data/distances/{algorithm}/results.csv"
```