

An Audio-Visual Approach to Music Genre Classification through Affective Color Features

Alexander Schindler¹, Andreas Rauber²

Motivation

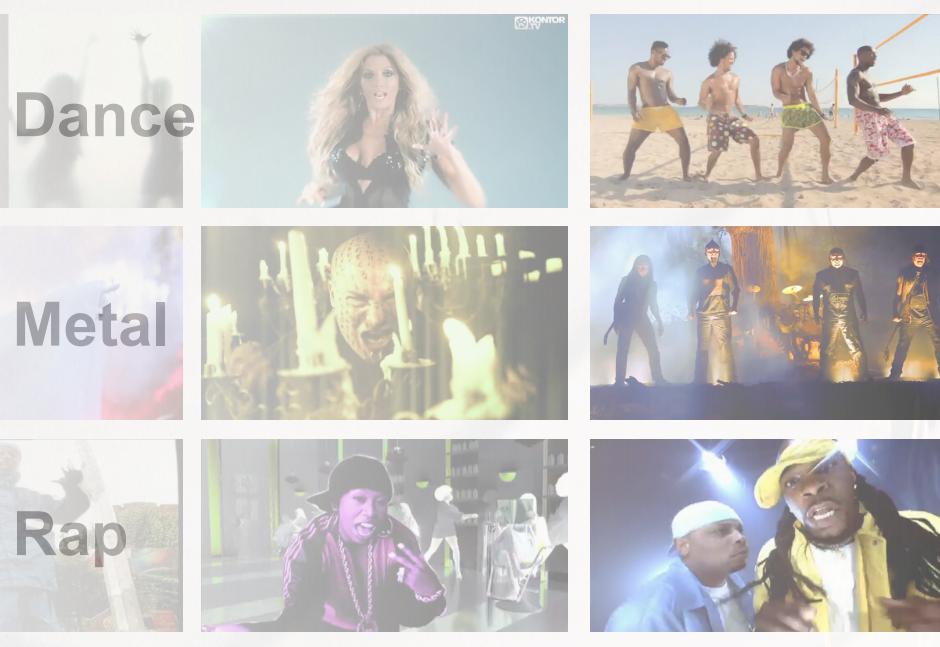
The “look of music” became an important factor in people’s appreciation of music. The rise of music video added an additional semantic layer to a music track. We grew accustomed to a visual vocabulary that is specific for a music style in a way that the genre of a music video can often be predicted by sight only (see Figure 1).

By augmenting music information retrieval technologies with solutions emerging from the video retrieval domain open research challenges could be addressed that are currently problematic to solve through audio content analysis only (e.g., classifying Christmas songs).

Rise in Music Video Importance

In 2011 we conducted a survey among music industry’s decision makers and stakeholders which showed that YouTube is considered to be the leading online music service. Other studies revealed that ‘watching music videos on computer’ was the most mentioned consuming activity of a survey

among 26,644 online consumers in 53 international markets. Music videos are also identified as passive triggers to active music search. Their visual aspects can strongly influence the amount of attention paid to a song.



Music videos contain enough information to estimate their genres by sight.

Feature Extraction

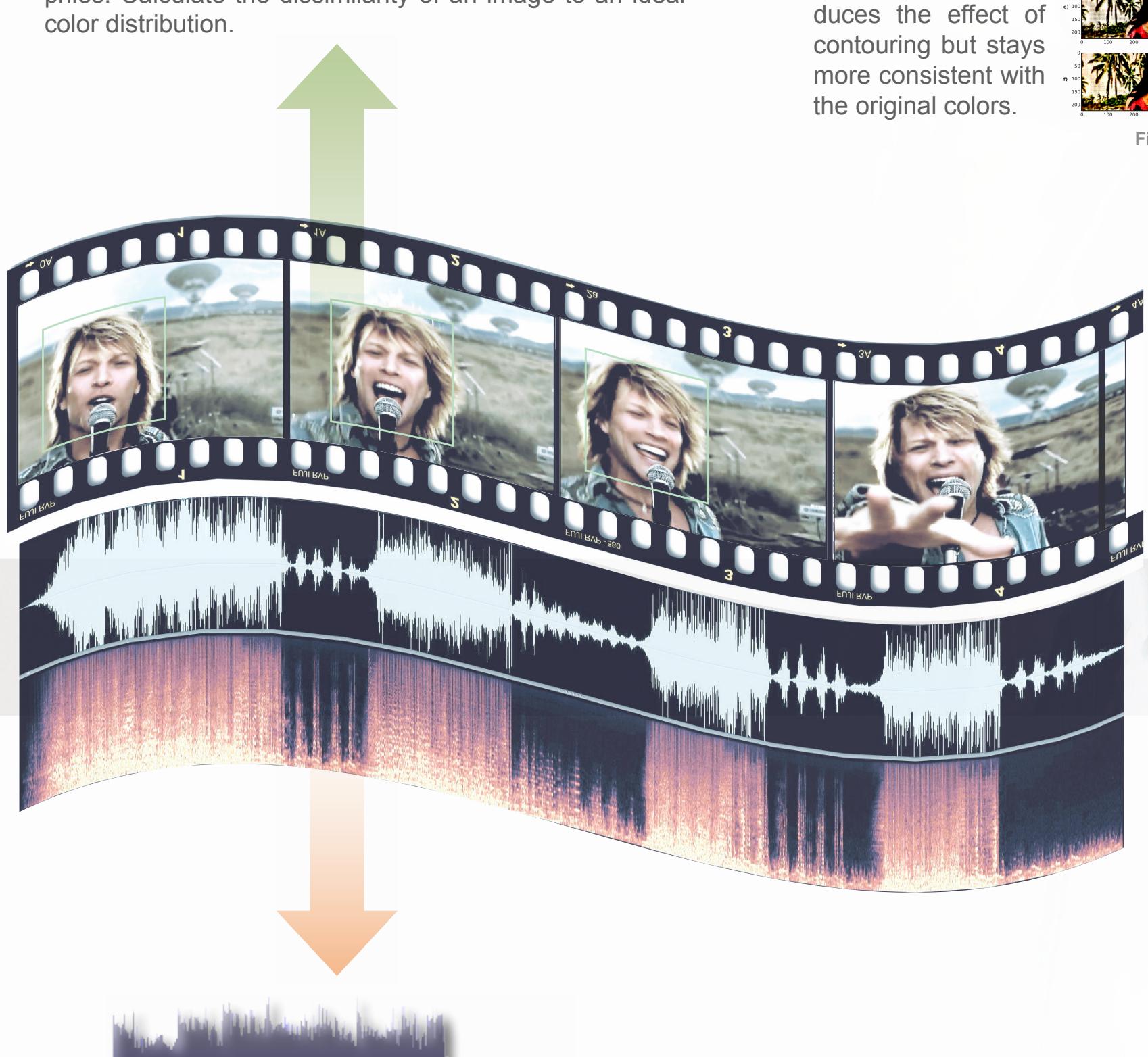
Audio features are extracted from the separated audio channel of the music videos. Visual features are extracted from each frame of a video and aggregated during post-processing by calculating the statistical measures. As a pre-processing step black bars at the borders of video frames, also called Letter-boxing or Pillar-boxing, are removed.

Global Color Statistics

Mean Saturation and Brightness based on the Improved Hue, Luminance and Saturation (IHLS) color space. Hue in IHLS is an angular value. Circular statistics has to be applied to assess mean Hue and its angular deviation. Saturation weighted aggregations are more robust towards weakly saturated colors.

Colorfulness

Used to computationally describe aesthetics in photographies. Calculate the dissimilarity of an image to an ideal color distribution.



MFCC

Mel Frequency Cepstral Coefficients (MFCC) are well known audio features derived from speech recognition.

Chroma

Chroma features project the spectrum onto 12 bins representing the semitones of the musical octave.

Color Names

Color distributions of the reduced Web-safe elementary-color palette: Magenta, Red, Yellow, Green, Cyan, Blue, Black and White. Contrast, brightness and color enhancement using Contrast Limited Adaptive Histogram Equalization (CLAHE). Ordered Dithering was used for color quantization, since it reduces the effect of contouring but stays more consistent with the original colors.

Figure 2: Sequential image enhancement steps.

Psycho-acoustic Music Descriptors

Proposed by Rauber, Lidy et al. are based on a psycho-acoustically modified Sonogram representation that reflects human loudness sensation. **Statistical Spectrum Descriptors (SSD)** subsequently compute statistical moments for the 24 critical bands of hearing. **Rhythm Patterns (RP)** describe fluctuations in modulation frequency which provide a rough interpretation of the rhythmic energy of a song. **Rhythm Histograms (RH)** aggregate the modulation amplitude values of the individual critical bands computed in a RP, providing a lower-dimensional descriptor for general rhythmic characteristics. **Temporal Variants (TSSD, TRH)** describe variations over time through statistical moments calculated from consecutive segments of a track.

Emotional Values

Pleasure-Arousal-Dominance model based on empirical investigated emotional reactions. The emotional values are linearly related to saturation and brightness in the IHLS color space.

Lightness Fluctuation Patterns

Calculated analogous to the music feature Rhythm Patterns. For each frame a 24 bin histogram of the lightness channel is calculated. Fast Fourier Transform (FFT) is applied to the histogram space of all video frames. This results in a time-invariant representation capturing reoccurring patterns in the video. Based on the observation that light effects, motions and shots are usually beat synchronized in music videos, LFPs can be assumed to express rhythmic structures of music videos.

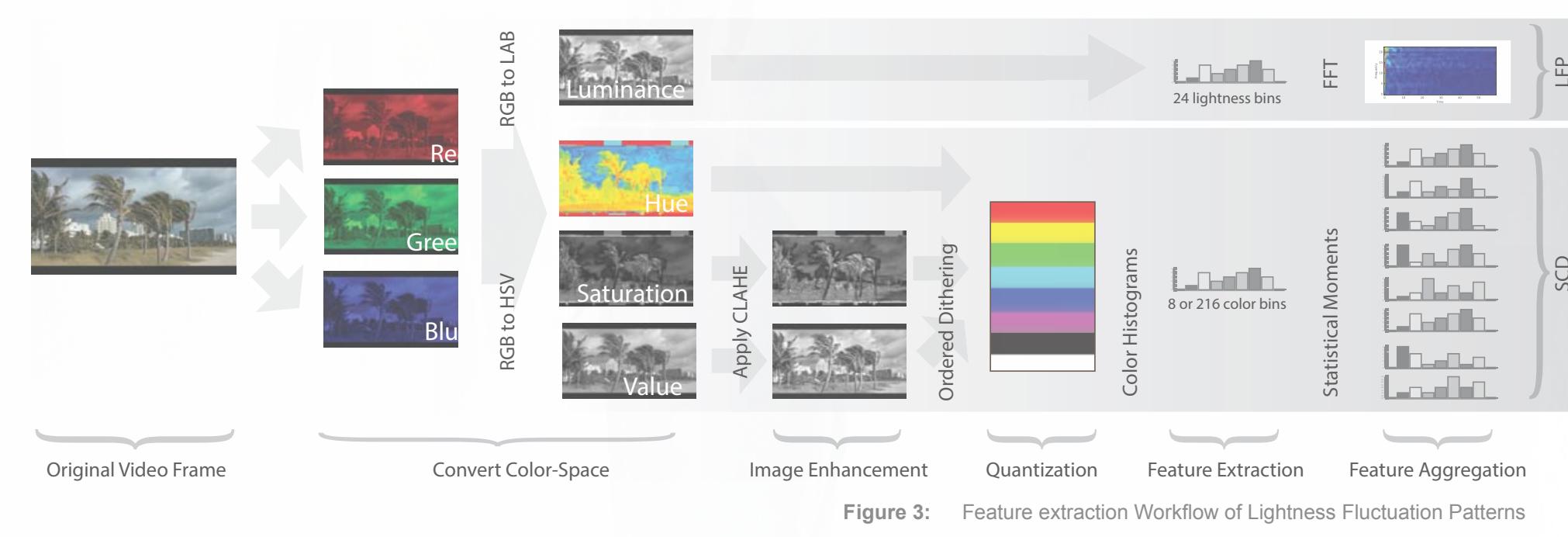


Figure 3: Feature extraction Workflow of Lightness Fluctuation Patterns

Itten Contrasts

A set of art-theory concepts defined for combining colors to induce emotions based on a proportional opponent color model. Contrast calculation is aligned to Wang’s feature extraction. Following contrasts were used: *Contrast of Light and Dark, Contrast of Warm and Cold, Contrast of Saturation, Contrast of Hue*.

Lightness Fluctuation Patterns
Calculated analogous to the music feature Rhythm Patterns. For each frame a 24 bin histogram of the lightness channel is calculated. Fast Fourier Transform (FFT) is applied to the histogram space of all video frames. This results in a time-invariant representation capturing reoccurring patterns in the video. Based on the observation that light effects, motions and shots are usually beat synchronized in music videos, LFPs can be assumed to express rhythmic structures of music videos.

Wang Emotional Factors

Three factors based on emotional word correlations that are relevant for image retrieval based on emotion semantics.

Feature One: lightness description of a segmented image ranging from very dark to very bright, combined with the hue labels cold and warm.

Feature Two: description of warm or cool regions with respect to different saturations as well as a description of contrast.

Feature Three: combines lightness contrast with an estimation of sharpness. A no-reference perceptual blur measure was used.

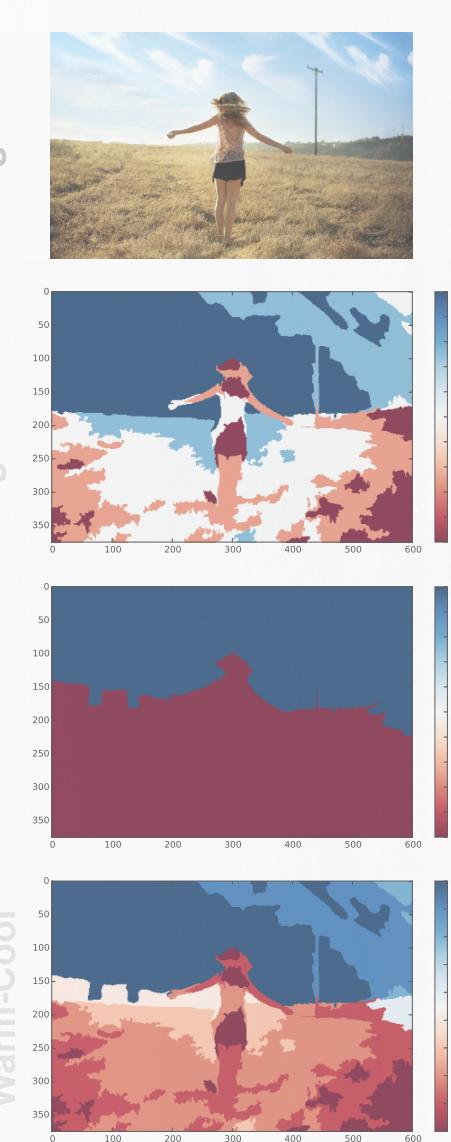


Figure 4: Visualization of contrast features.

Evaluation

The empirical evaluation is based on the Music Video Dataset. Classification experiments using Support Vector Machines (SVM), K-Nearest Neighbors (KNN), Random Forest (RF) and Naive Bayes (NB) classifiers, were conducted. In a second analysis Chi-square feature selection was applied to evaluate the performance of the visual and audio-visual feature-spaces.

Table 2: Overview of all features used in the evaluation. The column # indicates the dimensionality of the corresponding feature set.

Short Name	#	Description
Statistical Spectrum Descriptors (SSD)	168	Statistical description of a psycho-acoustic transformed audio spectrum
Rhythm Patterns (RP)	1024	Description of spectral fluctuations
Temporal Histograms (RH)	60	24 bin histograms of Rhythm Patterns
Temporal SSD and RH	60	Temporal variants of RH (TRH #420), SSD (TSSD #1176)
MFCC	12	Mel Frequency Cepstral Coefficients
Chroma	12	12 distinct semitones of the musical octave
Global Color Statistics	6	mean saturation and brightness, mean angular hue, angular deviation, with/without saturation weighting
Colorfulness	1	colorfulness measure based on Earth Movers Distance
Color Names	8	Magenta, Red, Yellow, Green, Cyan, Blue, Black, White
Pleasure, Arousal, Dominance	3	arousal, emotion values based on brightness and saturation
Itten Contrasts	4	Contrast of Light and Dark, Contrast of Saturation, Contrast of Hue and Contrast of Warm and Cold
Wang Emotional Factors	18	Features for the 3 affective factors by Wang et al. [17]
Lightness Fluctuation Patterns	80	Rhythmic fluctuations in video lightness

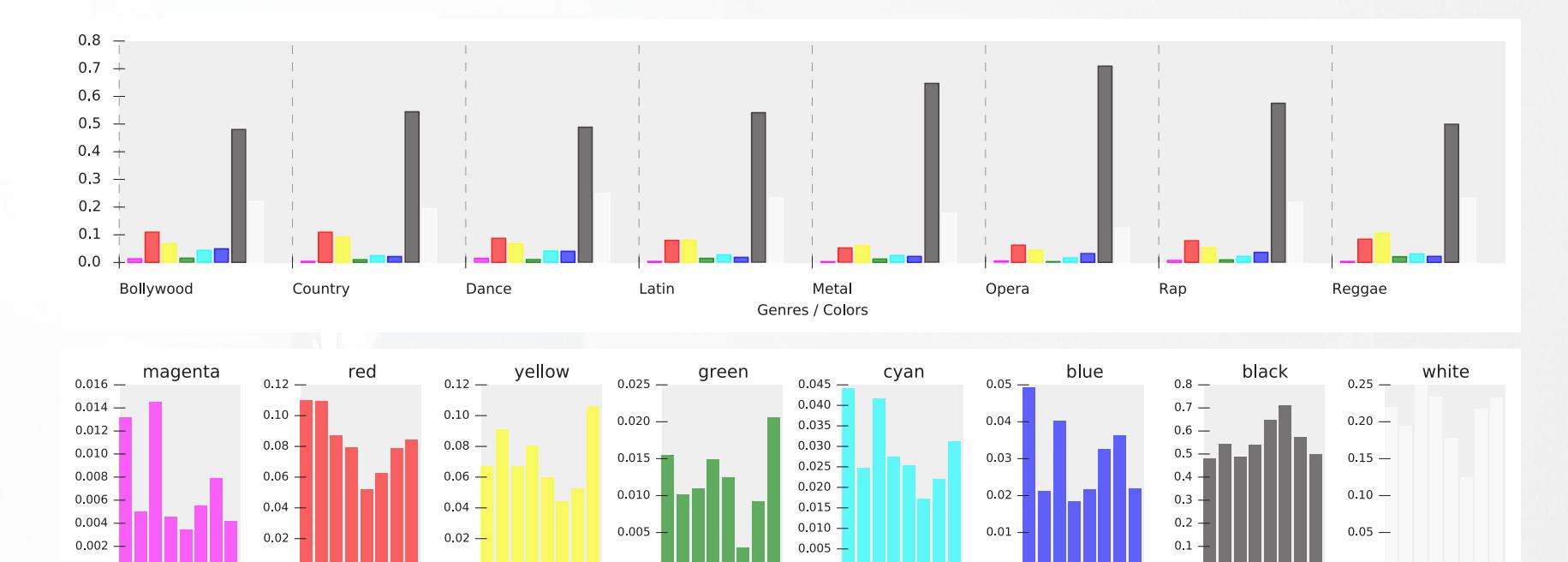


Figure 5: Distribution of named colors in the MV-MIX dataset

Visual

Using visual features only an accuracy of 50.13% could be reached for Support Vector Machines (SVM) for the MVD-VIS set. Accuracies for other sets or classifiers range from 17.89% to 39.38%. Because all classes equal in size these results are above a baseline of 12.5% or 6.25% respectively. Yet, the performance of the visual features alone is not representative.

Audio-Visual

The audio-visual results show interesting effects. Generally, there is no improvement of the performance over the top performing audio features. The results show that combining the visual features with chroma and rhythm descriptors has a positive effect on the accuracy while it is negative with spectral and timbral features.

Feature Selection

Applying ranked Chi-square attribute selection on the visual features shows, that affective features as well as the frequencies of black and white pixels have highest values. Further, more information is provided by variance and min/max aggregated values than by mean values.



Figure 6: Chi Square Feature Evaluation in descending order from left to right. Dark blue areas correspond with high Chi Square values.

Table 1: The Music Video Dataset - Detailed Overview of structure including class descriptions

Genre	Videos	Artists	BPM
MVD-VIS	100	32	122 (21)
Country	100	84	128 (1)
Dance	100	72	131 (5)
Latin	100	79	105 (1)
Metal	100	NA	104 (25)
Opera	100	81	104 (25)
Rap	100	75	114 (36)
MVD-MM			
80ies	100	75	135 (10)
Dubstep	100	66	116 (28)
Folk	100	69	129 (27)
Hard Rock	100	64	135 (26)
Indie	100	68	120 (25)
Pop Rock	100	69	128 (25)
Post-Rock	100	67	114 (25)
MVD-MIX			
16 Genres	1600	1040	