



Moodplay: Interactive music recommendation based on Artists' mood similarity

Ivana Andjelkovic^{a,*}, Denis Parra^b, John O'Donovan^a

^a University of California - Santa Barbara, Santa Barbara, CA 93106, USA

^b Pontificia Universidad Católica de Chile, Vicuña Mackenna Santiago 4860, Chile

ARTICLE INFO

Keywords:

Recommender systems
Music recommendation
Mood context
Context-aware recommendation
Affective computing
Recommendation interface

ABSTRACT

A large amount of research in recommender systems focuses on algorithmic accuracy and optimization of ranking metrics. However, recent work has unveiled the importance of other aspects of the recommendation process, including explanation, transparency, control and user experience in general. Building on these aspects, this paper introduces *MoodPlay*, an interactive music-artists recommender system which integrates content and mood-based filtering in a novel interface. We show how *MoodPlay* allows the user to explore a music collection by musical mood dimensions, building upon GEMS, a music-specific model of affect, rather than the traditional *Circumplex* model. We describe system architecture, algorithms, interface and interactions followed by use-case and offline evaluations of the system, providing evidence of the benefits of our model based on similarities between the typical moods found in an artist's music, for contextual music recommendation. Finally, we present results of a user study (N = 279) in which four versions of the interface are evaluated with varying degrees of visualization and interaction. Results show that our proposed visualization of items and mood information improves user acceptance and understanding of both the underlying data and the recommendations. Furthermore, our analysis reveals the role of mood in music recommendation, considering both artists' mood and users' self-reported mood in the user study. Our results and discussion highlight the impact of visual and interactive features in music recommendation, as well as associated human-cognitive limitations. This research also aims to inform the design of future interactive recommendation systems.

1. Introduction

Recommender systems are increasingly relied on in many domains for identifying relevant, personalized content from very large information spaces. Well established algorithms, such as Collaborative Filtering (Ekstrand et al., 2011), Content-Based Filtering (Pazzani and Billsus, 2007) and Matrix Factorization (Koren et al., 2009), are used across a variety of domains to recommend digital content or merchandise. Due to its unique consumption characteristics, music falls into a domain where alternative approaches to the traditional recommendation problem can help. These characteristics can be demonstrated by comparing music to two types of content widely offered to users via recommender systems: online movies and merchandise. For example, movies typically require undivided attention for 1–3 h and most often one movie is watched per sitting. On the other hand, we can listen to music throughout the day in almost any situation – while working, exercising, commuting, cooking, socializing and so forth. Similarly, online shopping is usually a focused action that most people engage in for a shorter period of time compared to music listening. While both can depend on a user's

taste, music listening is more often guided by feelings rather than practical reasoning. Overall, compared to other domains, music listening is more context dependent and closely tied to our emotional state. There are several music recommender systems that employ different types of context (daily activity Wang et al., 2012, time of the day (Baltrunas and Amatriain, 2009), music genre Maillet et al., 2009, etc.). However, previous work that integrates affective context for music discovery into a visual and interactive recommendation system is scarce (e.g. Musiccovery)¹. In this paper, we focus on prototyping and evaluating an interactive recommender system that suggests music bands based on artists' mood similarity and user input as an indicator of current preference.

Experimental evidence shows a strong relation between emotion and music (Koelsch, 2009) and previous research in affect-based recommender systems produced improvements over their non-contextual alternatives (Fernández-Tobías et al., 2013; Tkálčič et al., 2010). Previous studies, e.g. (Bostandjiev et al., 2012; Faltings et al., 2004; Gretarsson et al., 2010; Knijnenburg et al., 2012a; Nagulendra and Vassileva, 2014;

¹ <http://www.musiccovery.com>

* Corresponding author.

E-mail addresses: iva.andel@gmail.com (I. Andjelkovic), dparras@uc.cl (D. Parra), jod@cs.ucsb.edu (J. O'Donovan).

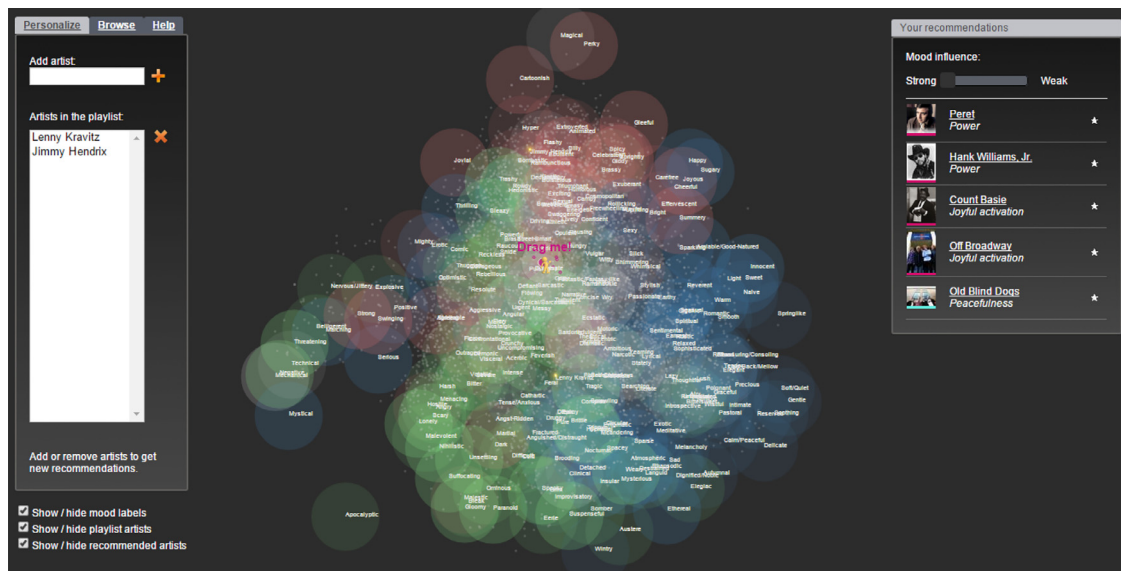


Fig. 1. Screenshot of the MoodPlay interface, divided into three sections: (left) pane for generating user profile by entering artist names, (center) snapshot of the mood space visualization, (right) recommendation list, along with slider for adjusting mood influence. Demonstration video is available at <https://youtu.be/eEdo32oOmcE>.

Parra et al., 2014; Verbert et al., 2013) demonstrate the importance of building interactive recommender interfaces, that go beyond the static-ranked list paradigm to improve user satisfaction. This trend is further supported by results showing that user satisfaction does not depend on recommendation accuracy only, but on factors such as serendipity, novelty, control and transparency as well (Konstan and Riedl, 2012; Mc-Nee et al., 2006). Our goal is to build a prototype recommender system with an interactive interface (Fig. 1) that supports users in discovering unknown, interesting items via interaction in an affective-aware visualization. As a proof of concept, we have designed and implemented MoodPlay, a system that: (a) visually represents affective metadata for a music recommender system, and (b) supports interaction, explanations and control over such visualization. We frame our work around the following research questions (RQ):

- (RQ1) What are the effects of interactive visualizations on the user experience with a recommender system, and what is the right amount of interaction for a music recommender?
- (RQ2) Does affective information improve recommendation accuracy and user experience versus when it is not included?

In our effort to answer the research questions, we have produced the following key contributions:

- *A novel visual interface for recommendation.* A visualization that maps moods and music artists in the same two-dimensional space, supporting item exploration and user control. The space is built using mood tags associated with artists, collected from an established, public database. We extend and visualize the music-specific emotion model - GEMS (defined in Section 2), to better fit a mood-aware music recommendation system.
- *Affect-aware recommendation method.* A novel hybrid recommendation algorithm for mood-aware and audio content-based music recommendation. The algorithm uses both mood tags of artists and audio content of their most popular songs.
- *Enhanced interaction techniques.* We introduce several new interaction mechanisms for hybrid recommendation on a visual mood space. For instance, trail-based and radius-based techniques.
- *Empirical evidence for avoiding high cognitive load.* We present an evaluation of the system through an online experiment ($N = 279$). We discover interesting relations between user interaction, trust, and user perception. We also provide some lessons for interface design in the context of exploratory tasks on recommender systems.

In contrast to our previous work on this system presented in Andjelkovic et al. (2016), graphical design and several interface features have been improved based on the user feedback from the first experiment. The mood space visualization has been updated to show smoother transition between mood categories. We enabled live streaming within the application to provide a more real-world context to the experiment and we display artist information upon clicking on artist nodes in the visualization. Using this improved system, a new experiment was designed, conducted and presented in this paper. Detailed comparison to the previous experiment can be found in 6.

To evaluate our system, we conducted a user study over four different conditions: (1) static recommendations in the form of ranked lists, generated based on a user's selection of seed artists (2) static recommendations, highlighted in a mood space visualization, (3) dynamic recommendation lists generated via user interaction in the mood space visualization, using current user's preference and (4) dynamic, interaction driven, trail-based recommendations.

The rest of this article is structured as follows. In the first two sections we provide important definitions and present related work. Then we describe MoodPlay interface, visual affective model and interactions, followed by a section detailing system architecture and recommendation methods. Next, we describe the user study setup and discuss the key results:

- In general, the system was rated as highly novel and fun to use.
- Visualization of mood information in a visual space significantly improved users' understanding of recommendations.
- Trail-based interaction (example shown in Fig. 5) was considered too confusing.
- Visual conditions (2), (3) and (4) tend to improve system trust, after trust propensity is controlled.

Finally, we share our ideas for future work and the resulting implication for design of future interactive recommender systems.

2. Definitions

Research presented in this article focuses on mood based recommendation. However, throughout the article we use related terms, *emotion* and *affect*, to explain different concepts. Here we provide definitions for each of the terms and other relevant constructs.

Table 1

Structure and description of *MoodPlay* mood hierarchy. Categories and sub-categories marked with * are the expansions from the original GEMS model.

Category	Sub-category	No. of moods	Example moods
Sublimity	Tenderness	24	Delicate, romantic, sweet
	Peacefulness	22	Pastoral, relaxed, soothing
	Wonder	24	Happy, light, springlike
	Nostalgic	9	Dreamy, rustic, yearning
Vitality	Transcendence	10	Atmospheric, spiritual, uplifting
	Power	29	Ambitious, fierce, pulsing, intense
	Joyful activation	32	Animated, fun, playful, exciting
Unease	Tension	32	Nervous, harsh, rowdy, rebellious
	Sadness	18	Austere, bittersweet, gloomy, tragic
	Fear *	10	Spooky, nihilistic, ominous
	Lethargy *	8	Languid, druggy, hypnotic
	Repulsiveness *	10	Greasy, sleazy, trashy, irreverent
Other *	Stylistic *	19	Graceful, slick, elegant, elaborate
	Cerebral *	12	Detached, street-smart, ironic
	Mechanical *	7	Crunchy, complex, knotty

Affect: Colloquial term that covers a broad range of feelings. It encompasses both emotions and moods (George, 1996; Västfjäll, 2002).

Emotions: Intense, short lived feelings, speculated by most researchers to be directed at someone or something (Frijda, 1993; Weiss and Cropanzano, 1996).

Moods: General, low intensity feeling states that often lack a contextual stimulus (Hume, 2012). While the duration of emotions is typically measured in minutes, moods may last several hours or days and cause us to think or brood for a while (Hume, 2012; Paul Ekman, 1994). Emotions become mood states when grouped into positive and negative categories because such grouping allows us to look at emotions more generally instead of in isolation (Hume, 2012). Therefore, emotion models such as Circumplex model of affect (Russell, 1980) are often used to represent moods as well.

Visual Mood Space: We use this term to refer to the 2-D space used to plot artists and moods in *Moodplay*'s visualization.

GEMS: This acronym stands for Geneva Emotional Music Scales. It is a music-specific emotion model proposed and validated by (Zentner and EEROLA, 2011; Zentner et al., 2008), which we used to categorize large number of music related moods. As stated by its authors, "GEMS is the first model and rating instrument specifically designed to capture the richness of musically evoked emotions".² GEMS is a hierarchical model, with three root emotions (Sublimity, Vitality, Unease), 9 corresponding sub-levels, and, in the third level, 45 emotion labels (details in Section 5, Table 1)

3. Related work

The following aspects are the most relevant to our research: affective-aware recommendations, recommendation of music bands, visual approaches to present recommendations beyond a rating list and affective-aware visualizations of music collections.

3.1. Affective computing and recommendations

Research in affective computing has been gaining extensive attention in recent years. Proliferation of mobile and wearable computer devices makes it both necessary and possible to achieve natural and harmonious human-computer interaction. Such devices enable us to track a variety of sources that carry emotional content. For example, different aspects of bodily movement and gestures have been used to recognize emotions: head and hands motion (Glowinski et al., 2011), gait patterns (Karg et al., 2010), body posture (Kleinsmith and Bianchi-Berthouze, 2007), to name a few. In the speech domain, vocal param-

eters such as pitch, speaking rate, formants and modulation of spectral content have also been successfully used to classify emotions in Petrushin (2000), Yu et al. (2001) and Wu et al. (2011). Furthermore, currently the largest data repository of face videos (2 million) owned by Affectiva³ is efficiently used to train computers in detecting emotions from facial expressions in real time.

For recommendation purposes, (Masthoff, 2005) integrated affective state in a group recommender system by modeling satisfaction as mood, while Gonzalez et al. (2007) incorporated the emotional context in a recommender system for a large e-commerce learning guide. More related to our work, Park et al. (2006) developed probably the first context-aware music recommender that exploited mood inferred from context information. And more recently, Tkalcic et al. (2011) and Han et al. (2010) discussed the role of emotions in recommender systems and introduced a framework to identify the stages where emotion can be used for recommendation.

In the music recommendation domain, several works infer the users' mood for music recommendation based on movements, temperature and weather (Cunningham et al., 2008) or from the music content (Rho et al., 2009). For example, Griffiths et al. (2013) measure a variety of contextual and physiological indicators (temperature, light, heart activity) in order to detect mood and recommend music by mapping both user's mood and music on the same emotion map. van der Zwaag et al. (2013) take target mood as an input from user and then select songs that direct the user towards the desired mood, while measuring skin conductance to verify the change. Skin temperature (Janssen et al., 2012) and arm gestures (Amelynck et al., 2012) have also been used for inferring mood and querying music collections. Compared to these studies, in our up to date work we use mood information associated with a set of seed artists provided by user to suggest new artists in similar moods. In addition, we propose a rich interface to help users explore mood space and choose music in a desired mood. In the future, this proposed system would be greatly enhanced by incorporating a method for detection of user's mood, using sensors available on wearable devices, social media activity or other contextual information.

3.2. Recommendation of music bands

Recommendations in the music domain is a well-established field within recommender systems, which have shown, among many others, approaches to recommend tracks (Celma and Herrera, 2008; Logan, 2004), albums (Parra and Amatriain, 2011), playlists (Baccigalupo and Plaza, 2006; Hariri et al., 2012; Maillet et al., 2009), music targeted at specific venues (Kaminskas and Ricci, 2011), music targeted

² <http://www.zentnerlab.com/content/musically-evoked-emotions>

³ <http://www.affectiva.com>

at daily activities (Wang et al., 2012), and artists and music bands (Bostandjiev et al., 2012; Hijikata et al., 2012). Since our proposed interface aims at recommending music bands, we focus on presenting related work in this sub-field. Hijikata et al. (2012) used a Naive Bayes recommender to present recommendations of music bands, while Bostandjiev et al. (2012) used a hybrid controllable recommender system with a visual interactive interface, TasteWeights. Compared to these previous approaches, we innovate by using professionally curated mood tags associated with bands to compute similarity, by introducing a user-controllable recommendation interface and by allowing users to explore the music band dataset interactively.

3.3. Visual approaches to recommendation

McNee et al. (2006) highlights the importance of user-centric approaches to evaluating recommender systems, and of developing interfaces and interaction designs, instead of focusing solely on improving recommendation algorithms. Konstan and Riedl (2012), who shows that small improvements in recommender accuracy do not necessarily improve users' satisfaction with a system. However, the development of interfaces that present recommended items in a visual model different from a static ranked list is rather scarce. For example, SFV is (Gou et al., 2011) and Pharos (Zhao et al., 2011) employ visualizations of social, latent communities to recommend new friends and social websites respectively. Other examples include collaborative filtering recommenders with rich user interactions such as PeerChooser (O'Donovan et al., 2008) and SmallWorlds (Gretarsson et al., 2010), and interactive visualizations for recommending conference talks – TalkExplorer (Verbert et al., 2013) and SetFusion (Parra et al., 2014). There is also a range of systems that support dynamic critiquing of an algorithm, such as (Pu et al., 2011) and (Chen and Pu, 2009). Finally, Nagulendra and Vassileva (2014) created an interactive visualization which provides users of social networking sites with awareness of the personalization mechanism. For a detailed review of visual and interactive recommender systems, read He et al. (2016). Although not focused on personalized recommendation, but rather on navigation of musical datasets, Knees et al. (2006) introduced a virtual 3D landscape which allows the user to freely navigate a collection. To the best of our knowledge, Moodplay is the first interactive music recommender system that maps the artists in a latent, navigable, affective visual space based on the recently developed music-specific mood model GEMS (Zentner et al., 2008), further explained in Section 4.2.

3.4. Affective-aware visualizations of music collections

Although affective-based music selection and recommendation are gaining popularity in both research and commercial settings, the development of visual aids for affective information is still scarce. Nearly all existing visualizations are built upon Russell's circumplex model of affect (Russell, 1980). This model is today commonly used to represent emotions and moods as a mixture of two dimensions, valence and arousal, positioning them in the coordinate system. Yang et al. (2008) incorporated it into their music retrieval method, and a commercial application Habu⁴ uses it as a platform for music selection based on mood.

However, many emotions cannot be uniquely characterized by valence and arousal values (Collier, 2007). For example, fear and anger, two distinctive emotions, both have high arousal and negative valence, and are commonly placed close to each other in the circumplex model (Scherer et al., 2003). It is also important to note that models derived from general research in psychology, such as Russell's, may not be suitable for musical emotions. One reason being that music, unlike other life events, possibly induces more contemplative range of emotions (Zentner and EEROLA, 2011). To address this problem, we propose a

novel visual representation of music specific affective dimensions, built upon the GEMS model derived from an extensive psychological study by (Zentner et al., 2008) (see Section 4.2 for details).

4. System overview

The MoodPlay system is accessible via web browser and consists of three sections: input, visualization and recommendation panel. Users construct profiles by entering names of artists via an interactive drop-down list (Fig. 1-left). Based on the mood information associated with profile artists (see Section 4.2 for explanation), the system positions a user avatar in a precomputed visual mood space (Fig. 1-center) and recommends new artists (Fig. 1-right). In this section we provide an overview of the user interface and explain the method for constructing the mood space.

4.1. Interface design

Visualization. Visualization of the mood space along with the artists within it is central to solving the problem of navigation through the music collection and explanations of recommendations. The space consists of 266 moods - similar ones being positioned closer to each other than dissimilar ones (the construction method is detailed in Section 4.2). Furthermore, moods form a hierarchy with three primary categories at the top - vital, sublime and uneasy (see Table 1), portrayed on canvas in different colors. Red, generally associated with passion and high energy (Nijdam, 2005), is used for vital mood category. Blue, more serene color (Nijdam, 2005), denotes sublime category, which includes tender and peaceful moods among others. Lastly, uneasy mood nodes are green - a color that is occasionally associated with sickness (Nijdam, 2005), but here chosen mainly for aesthetic reasons to complement the other two category colors. Mood nodes are semi-transparent, their size is equal and purposefully large enough to cause overlap. This produces an interplay of colors, thus forming the space with gradual transitions between mood categories. Artists from our database are placed within the mood space based on moods associated with their music. Users can stream their music in real-time and see additional artist information by clicking on the nodes.

Our interface design follows Shneidermann's visual information seeking mantra "Overview first, zoom and filter, then details on-demand" (Shneiderman, 1996). As indicated in Table 1 and explained in 4.2, moods form a hierarchy with 3 top categories displayed in distinct colors, and 15 subcategories. User can choose to display individual subcategories via Browse tab on the left interface pane (Fig. 1). Such a hierarchical view of the large number of moods allows the user to explore the mood space by starting from broad terms and then filtering down selection to specific subset of moods in the visualization, while zooming and panning to more closely inspect areas of interest. Furthermore, we implement a dynamic mood labeling algorithm in order to reduce cognitive load in a dense mood space. We rank the moods based on the frequency of their usage to describe different artists, and show only limited number of the most popular moods at a given zoom level.

Recommendations. An ordered list of recommended artists is displayed in the right panel (Fig. 1) and the corresponding artist nodes are highlighted in the mood space. In this way, we aim to provide transparency, trust, efficiency and satisfaction to the user, which are four out of the seven criteria identified by Tintarev and Masthoff (2011) to design explanations in recommender systems (the other three are scrutability, effectiveness and persuasiveness). Items in the recommendation list are linked to audio streams via Rdio⁵ and to Last.fm⁶ profiles of artists. For each recommended artist we also display artist's picture, color of the top mood category (red, blue or green) and name of the sub-category the

⁴ <http://habumusic.com>

⁵ Rdio streaming service was discontinued in December 2015.

⁶ <http://www.last.fm/api/intro>

artist belongs to, with the goal to help users gain some understanding of the music upon visual inspection. Furthermore, because recommended items change as a result of the user's interaction with the system, we display up or down arrows next to artist name if its position changed, a horizontal line if it remained the same or a star if the recommendation is new. Rating of recommended items is enabled only for the purpose of the user study, and is achieved by clicking on one of the five stars below artist names.

Interaction. Adaptivity of music recommenders is particularly important due to the dynamic nature of the listening context (Stober and Nürnberger, 2013). Keeping this in mind, we model the gradual change of a user's preference by enabling the movement of the avatar (Figure 2.1) in the mood space and maintaining the array of trail marks, weighted by distance from the current position (Figure 2.3). As the user navigates away from the initial position, we incorporate the mood information associated with each trail mark into the recommendation algorithm. Removing any of the trail marks is possible by simply clicking on it, and deleting the whole trail is achieved by clicking on the initial position of the avatar.

Finally, fine-tuning of recommendations is further supported by controlling the hybridization of recommendation process. Our recommendation approach accounts for the fact that mood-based similarity between artists does not necessarily match audio based similarity (e.g. techno and punk artists are both energetic, but they do not sound similar). Therefore, we allow users to adjust the mood influence via a slider control which dynamically re-sizes a catchment area around the current avatar position (Figure 2.1 and 2.2). The weaker the mood influence, the more we rely on audio similarity to calculate recommendations, and vice-versa.

4.2. Music specific visual model of affect

A key challenge of this research was showing and explaining inter-relationships between artists and moods in a two-dimensional space. To that end, we analyzed mood metadata for 4927 artists collected from Rovi,⁷ which is to our knowledge the most comprehensive collection of professionally curated mood-artist tags. Each artist in our dataset is characterized by approximately 5 to 20 weighted mood words out of 289 available ones, and represented with a vector $X \in \mathbb{R}^{289}$. Finally, correspondence analysis (Salkind, 2010) was used to reduce data dimensionality, which resulted in a latent variable space containing moods and artists.

Numerous emotion models, both continuous and categorical have been proposed in the psychology field (Izard, 1977; Scherer, 1984; Schimmack and Grob, 2000). For the purpose of identifying potential clusters in our mood space, we explore whether our visual map fits into a hierarchical, and therefore categorical, music-specific emotion model proposed by Zentner et al. (2008). This model, from now on referred to as *Geneva Emotional Music Scales* or GEMS, consists of 3 main categories (vitality, uneasiness, sublimity), 9 sub-categories and 45 music relevant emotion words distributed across different sub-categories. Our hypothesis was that such hierarchy should emerge in the visual mood space built upon professionally curated artist-mood associations. It is important to note that psychology researchers focus on deriving emotion models rather than mood models, and for recommendation purposes, music is generally tagged with mood descriptors. In this paper we use either of the terms depending on the field we address, and an overarching term, affect, in the context of our proposed system.

To perform our hierarchical classification of moods, we employed a WordNet⁸ similarity tool⁹ and calculated similarity scores between 289 Rovi and 45 GEMS mood words. Furthermore, since similarity between

terms in WordNet is based on semantic relatedness and not strictly on synonymity, we evaluated mood classification by subjective observation. For example, the word *volatile* was found to be related more closely to *tender* than *tense* and was placed into sub-category *Tenderness*, rather than *Tension*. Hence, the following steps were taken to reduce observed classification error rate: (1) mood hierarchy is expanded to accommodate moods that do not belong to any of the GEMS categories, (2) 23 of the least frequently used mood words to describe artists in Rovi were discarded, (3) the set of remaining misclassified words are placed into categories that they are more likely to belong to. Table 1 shows the final list of categories and distribution of associated moods.

5. Technical design and implementation

MoodPlay uses diverse data collected from different sources, mostly through public Web APIs. Recommendations have to be computed very quickly, since they are immediately presented in the interface as a result of user interaction. Therefore, the system requires an appropriate architectural design. As depicted in Fig. 3, it has two main components: one for building the library of items with their metadata (*Dataset Construction*) and a second component that generates user recommendations (*Recommendation Framework*). The following subsections describe the architecture design and implementation in detail.

5.1. Dataset

MoodPlay relies on a static dataset of 4927 artists, obtained in several iterations. First, 3275 artists were randomly selected from a subset of the Million Songs Dataset.¹⁰ Artists ranged from very popular to less known, and played music in a variety of genres and over different decades. The pool was then expanded by 2000 most popular artists from the public EchoNest¹¹ database, as measured by proprietary metrics *familiarity* and *hottness*. We complemented the initial set in order to better facilitate an online user study with participants of different ages from different parts of the world. Artists for which we were not able to obtain mood or song data were discarded. Next, mood data for each artist was obtained via Rovi API and the top ten most popular songs for each artist along with corresponding audio analysis data were obtained from EchoNest. Different versions of the same song, having the same title in EchoNest database were discarded. Rdio API was used for music streaming in *MoodPlay*. Finally, artist pictures and links to Last.fm profiles were obtained via Last.fm API.

5.2. Recommendation approaches

Our hybrid cascading recommender (Burke, 2002) operates in two stages as shown in Fig. 4: (1) using the user profile as an input, our system produces a first candidate set of recommendations based on mood similarity, and (2) the output of the first recommender becomes the input to an audio content-based recommender, which re-ranks the artists and produces the final recommendation list. Such layered approach supports our goal to help the user control and understand how recommendations are generated while navigating mood space. The following paragraphs describe the recommendation process in detail.

Offline computation of artist similarity. Artists' pairwise similarity, based on mood and audio content, is calculated offline and stored in two separate data structures. Mood-based similarity between any two artists is a function of their Euclidean distance in the affective space produced by correspondence analysis. To calculate audio-based similarity, we first identify the 10 most popular songs for each artist in our database via the EchoNest API and obtain audio analysis data for the total of 49,270 songs from the same source. We used timbre, tempo, loudness and key

⁷ <http://developer.rovicorp.com/io-docs>

⁸ <https://wordnet.princeton.edu/>

⁹ <http://maraca.d.umn.edu/cgi-bin/similarity/similarity.cgi>

¹⁰ <http://labrosa.ee.columbia.edu/millionsong/pages/getting-dataset#subset>

¹¹ <http://developer.echonest.com>

confidence attributes, which amounted to approximately 10,000 numerical values per song. In order to make the similarity calculations efficient, we represent each song with a vector $v_i \in \mathbb{R}^{515}$ (McFee and Lanckriet, 2011) and build artist data into a KD-tree (Bentley, 1975). Finally, an accelerated approach for nearest-neighbor retrieval that uses maximum-variance KD-tree data structure was used to compute similarity between songs, since it has a good balance of accuracy, scale and efficiency (McFee and Lanckriet, 2011). In this way, time complexity of constructing a similarity matrix was reduced from $\mathcal{O}(n^2)$ to $\mathcal{O}(n \log n)$, while the search for the K nearest neighbors of a given artist is reduced from $\mathcal{O}(K \cdot n)$ to $\mathcal{O}(K \cdot \log n)$. To compute artist similarity, first, for each song we rank all other songs from the dataset from most to least similar. We then calculate average similarity rank of songs per artist (Celma and Herrera, 2008), thus obtaining the artist similarity matrix (Algorithm 1).

Algorithm 1 Algorithm for computation of audio similarity.

Input:Set of artists: $A = \{a_1, a_2, \dots, a_n\}$ Set of songs for all artists: $S = \{Sa_1 \cup Sa_2 \cup \dots \cup Sa_n\}$ **Output:** Audio similarity ranks: $ARanks = \{a_i \rightarrow \{a_j \rightarrow rank_{ij}\}\}$

```

1: function COMPUTEAUDIOSIMILARITYRANKS
2:   ARanks = {}                                ▷ dictionary of artist similarity ranks
3:   for each artist  $a_i$  in  $A$  do
4:     SRanks = {}                                ▷ dictionary of song similarity ranks
5:     for each song  $s_k$  in  $Sa_i$  do
6:       SRanks[ $s_k$ ] = COMPUTESIMILARITYMAPOF-
       SONGRANKS( $s_k, S$ )
7:     end for
8:     for each artist  $a_j$  in  $A$  do
9:       ARanks[ $a_i$ ][ $a_j$ ] = COMPUTEAVERAGESONGSIMILAR-
       ITY(SRanks,  $Sa_j$ )
10:    end for
11:  end for
12:  return ARanks
13: end function

14: function COMPUTESIMILARITYMAPOFSONGRANKS( $s, S$ )
15:   Rank all songs from  $S$  based on audio similarity to song  $s$ 
16:   for each  $s_j$  in  $S$  do
17:     similarityMapOfSongRanks[ $s_j$ ] =  $rank_j$ 
18:   end for
19:   return similarityMapOfSongRanks
20: end function

21: function COMPUTEAVERAGESONGSIMILARITY(SRanks,  $Sa$ )
22:   average = 0
23:   for each song  $s_i$  in  $Sa$  do
24:     for each song  $s_j$  in  $SRanks.keys$  do
25:       average += SRanks[ $s_j$ ][ $s_i$ ]
26:     end for
27:   end for
28:   average = average / ( $Sa.size * SRanks.size$ )
29:   return average
30: end function

```

Online recommendation. During a user session, MoodPlay recommends new artists similar to the artists the user enters into her profile. First, we determine the overall mood by calculating the centroid $C(u) = (c_x, c_y)$ of profile artist positions, where we then place the user avatar. The coordinates c_x and c_y are calculated as in Eq. (1), where P is the set of artists in the user profile, a_x is the x -axis and a_y is the y -axis coordinate of artist a in profile P .

$$c_x = \sum_{a \in P} \frac{a_x}{|P|} \quad c_y = \sum_{a \in P} \frac{a_y}{|P|} \quad (1)$$

Artists found within the adjustable radius around the centroid are all potential candidates for recommendation because they are considered to reflect the latent moods derived from the user's input. Among the candidate artists, we select the ten most similar to the user profile based on pre-computed audio similarity data, rank them by distance from the user position and display first five as recommended artists (Algorithm 2).

Algorithm 2 Basic algorithm for online music recommendation.

Input:Artists in user profile: $P = \{a_1, \dots, a_n\}$ User position: $u = Centroid(a_1, \dots, a_n)$, $a_i \in P$ or u is a position from user's trail $T = \{u_1, \dots, u_n\}$ Recommendation radius: r Audio similarity ranks: $ARanks = \{a_i \rightarrow \{a_j \rightarrow rank_{ij}\}\}$ Number of recommendations: n_{rec} **Output:**Recommended artists: $R = \{a_1, \dots, a_n\}$

```

1: function RECOMMENDMUSIC( $u$ )
2:    $M = []$                                 ▷ artists within mood radius
3:   for  $a_i$  in  $A - P$  do
4:     if distance( $a_i, u$ ) <  $r$  then  $M[i] = a_i$ 
5:   end if
6: end for
7:  $H = \{\}$                                 ▷ dict. of artists & similarity with  $P$ 
8: for  $a_i$  in  $M$  do
9:    $H[a_i] = AVERAGESIMRANKING(a_i, P)$ 
10: end for
11: sort( $H$ )                                ▷ sort artists by audio similarity
12:  $R = H[1..n_{rec}]$ 
13: return  $R$ 
14: end function

15: function AVERAGESIMRANKING( $a, P$ )
16:   average = 0
17:   for each  $a_i$  in  $P$  do
18:     average += ARanks[ $a$ ][ $a_i$ ]
19:   end for
20:   return average /  $P.size$ 
21: end function

```

Trail-based recommendation. Furthermore, we propose a novel, adaptive recommendation approach that accounts for the preference change in terms of mood, reflected by the repositioning of user's avatar in the affective space. We keep track of each new position and apply a decay function to the preference trail when recommending new artists. Recommendations from the last position in the trail are assigned the greatest weight, because we presume that the most recent mood area of interest is the most relevant to user. The weights further decrease as a function of hop distance from the end of the trail. Pseudocode for the trail based recommendation algorithm is given in Algorithm 3 and here we outline the steps.

At each trail mark, we apply the recommendation algorithm described in the previous sub-sections, which produces an initial set of recommendation candidates. We then calculate adjusted distances d_a between trail marks and surrounding recommendation candidates in two steps. First, we normalize distances between the trail mark and artists because the radius can vary among trail marks. If the distances were not normalized, many relevant artists would be falsely considered irrelevant and would not appear in the final recommendation list. Next, we adjust the normalized distances for each trail mark based on the corresponding weights using the formula $d_a = d_n + \Delta \times (|T| - i)$, where d_n is a normalized distance, Δ is a decay constant, $|T|$ is a total number of trail marks and i is an iterator over the trail marks. After several tests, we found that weight constant Δ performs the best when calculated as: $\Delta = r_{min}/4$, where r_{min} is the minimal recommendation radius.

Algorithm 3 Hybrid recommendation with provenance trails.**Input:**

Trail of user positions: $T = \{u_1, u_2, \dots, u_n\}$, where u_1 is the profile based position and consecutive u_i are positions that user navigated to

Current recommendation radius: r

Minimum recommendation radius: r_{\min}

Number of recommendations: n_{rec}

Output:

Recommended artists: $R = \{a_1, \dots, a_n\}$

```

1: function RECOMMENDMUSICBASEDONTRAIL
2:    $R = \{\}$  ▷ dict. of recommended artists
3:    $\Delta = r_{\min}/4$ 
4:   for  $u_i$  in  $T$  do
5:     for  $a_j$  in RECOMMENDMUSIC( $u_i$ ) do
6:        $d_s = \text{SCALE}(\text{distance}(u_i, a_j), r, r_{\min})$ 
7:        $d_a = d_s + \Delta \times (T.\text{size} - 1 - i)$ 
8:        $R[a_j] = d_a$ 
9:     end for
10:  end for
11:   $\text{sort}(R)$  ▷ sort artists in R by  $d_a$ 
12:  return  $R[1..n_{\text{rec}}]$ 
13: end function

14: function SCALE( $d, r, r_{\min}$ )
15:   $d_c = \text{Convert } d \text{ from range } [0, r] \text{ to } [0, r_{\min}]$ 
16:  return  $d_c$ 
17: end function

```

The larger the value of Δ , the steeper the decay function. Finally, the recommendation candidates are sorted based on adjusted distances, and top five are recommended to the user.

6. Evaluation

Preliminary evaluation of an early version of MoodPlay has been described in Andjelkovic et al. (2016). Compared to the previously published study, here we present modified experiment design and more comprehensive analysis of results. We performed a crowd-sourced study with entirely new set of participants. Number of participants was 378, but 279 remained after filtering out those that we did not deem as valid, i.e. those who incorrectly answered attention check questions or ended the study prematurely.

The focus of the evaluation was to understand the effects of mood-based interactions with a recommendation algorithm and to independently evaluate the influence of the MoodPlay visualization from an explanatory perspective. To improve the previous experiment design, in this study we gave users more freedom to naturally interact with the system and we tracked additional interaction metrics. Furthermore, in the previous study we focused the evaluation on user characteristics, interaction and experience, and placed less attention on ratings-based analysis. Here we report the results of both quantitative and qualitative analysis and address impact of mood based interactions on user experience.

6.1. Mood data in automated recommendation

Before we proceed with our main experiments on interaction with mood data, we first set out to understand the impact that mood information can have on *automated*, non-interactive algorithms. This is an important step that can provide insight into the utility of different inputs about mood during the interactive process that we describe later. In particular, we describe the results of an automated experiment to show quantitatively that using mood as context during recommendation can improve recommendation quality. To do this, we use two versions of

Table 2

Table of experimental conditions and associated features. Conditions increase in complexity, (1) having only two and (4) having all available features.

Feature	(1)	(2)	(3)	(4)
Profile generation	x	x	x	x
Ordered list of recommendations	x	x	x	x
Display of mood space		x	x	x
Navigation in mood space			x	x
Hybridization control			x	x
Trail based recommendations				x
Number of subjects	70	69	70	70

context-aware matrix factorization method from Zheng et al. (2015), trained using a traditional item-rating matrix, along with mood information from our GEMS model, for each item. We compute prediction error to compare against three traditional recommendation algorithms, which are only trained on a matrix of user and item ratings. One limitation of this experiment is that we are predicting over ratings that were gathered through the Moodplay system, which may have influenced the rating in different ways. We plan an additional experiment with ratings gathered from LastFM to verify our results on a separate data set. We run a 5-fold cross validation on a collection of 2548 ratings of 593 items (musical artists) with 5 associated mood tags per item. Data density was 0.29% and the mean item rating was 2.75 on a 1–5 Likert scale, with a standard deviation of 1.31. Fig. 5 shows the results of this experiment for the five methods on the x-axis. The x-axis groupings are MAE and RMSE scores for each algorithm, two popular error metrics used to measure the predictive accuracy of recommendation algorithms Parra and Sahebi (2012). Y-axis shows the value for those metrics. The CAMF or context-aware matrix factorization methods were trained with additional mood information. The CAMF_{CI} algorithm computes mood information over items, while the CAMF_{CU} computes it over individual users. Fig. 5 is sorted from left to right by best performance. It is clear that the three mood based methods outperformed the benchmarks, on both metrics. For the benchmarks, we chose a classic user-KNN collaborative filtering algorithm, along with a simple global average predictor and a pure matrix factorization approach. We note that the poor performance of the latter may be a result of sparsity in the ratings matrix – which enables us to make the point that perhaps mood-aware algorithms can be a good bootstrapping mechanism that helps to combat the sparsity problem for traditional matrix factorization algorithms when such data is available.

6.2. Setup

As in the previous study, we set up four conditions having different features, shown in Table 2. The conditions have increasing visual and interaction complexity (see sub-section *Interaction* in 4.1 for description of more complex features). In each of the conditions users create a profile by entering artist names. The system uses this information to generate recommendations and display them as a list. In condition 1, users see the list of recommendations but do not see the visualization and mood information. In condition 2, mood space and the user's avatar within it are visible, but interaction is not enabled. Condition 3 allows users to navigate the mood space, move the avatar without keeping track of previous positions, and modify the size of catchment area around the avatar. Finally, condition 4 (Fig. 2) tests the full system, including trail based recommendations.

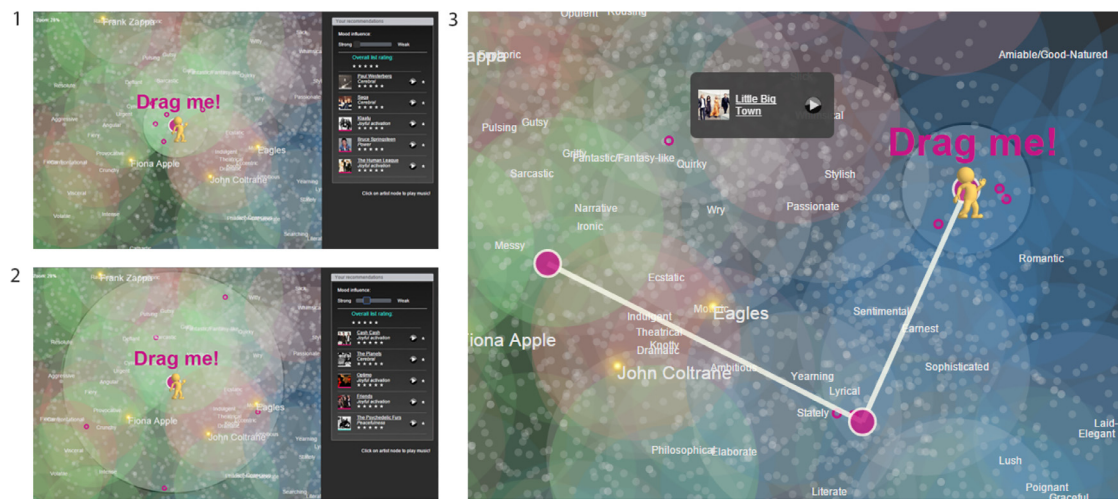


Fig. 2. Screenshots of interactive features in *MoodPlay*: (1) and (2) - the varying recommendation catchment area around user avatar, controlled by a hybridization slider, (3) - trail based interaction, along with a display of artist information box upon clicking on artist node in the visualization.

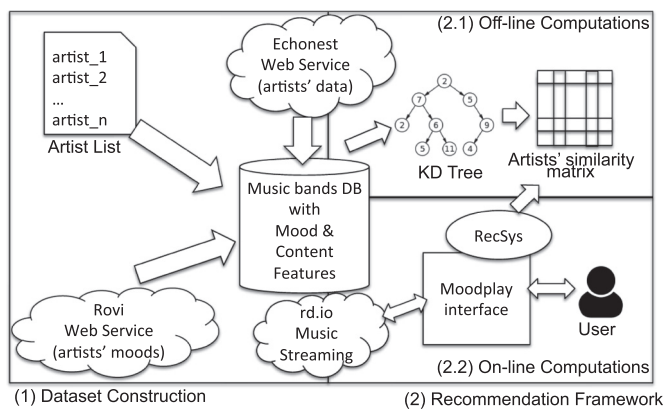


Fig. 3. MoodPlay system architecture indicating the modules for: (1) dataset construction and (2) recommendation framework.

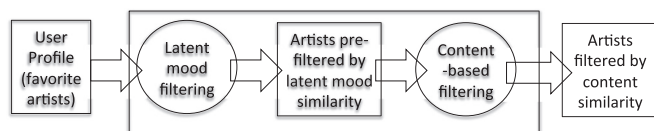


Fig. 4. Schematic representing our hybrid cascading recommender which pre-filters based on mood similarity and then post-filters based on content similarity.

6.3. Study procedure

Participants accepted the study on Mechanical Turk¹² and were redirected to a Qualtrics¹³ pre-study survey with demographic and propensity related questions. Among the questions in this survey, we collected users' current mood by asking them to choose one of the following options: (a) Sublime (e.g.: joyful, warm and tender moods), (b) Vital (e.g.: stimulating moods such as "lively", "energetic" or "fierce") or (c) Uneasy (e.g.: negative moods such as "sad", "tense" or "fearful").

Following this, they were assigned to a random condition and performed the main task. Finally, participants gave qualitative feedback in a post-study survey, also administered through the Qualtrics platform.

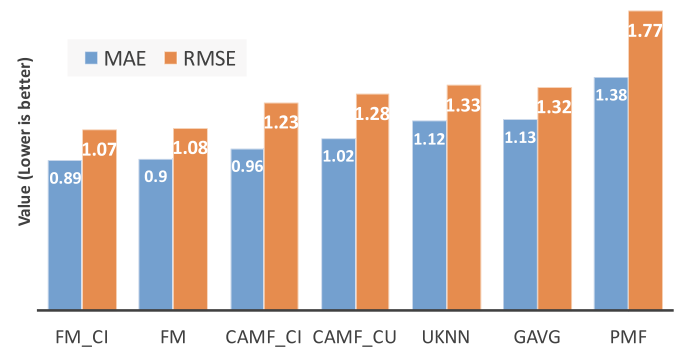


Fig. 5. Predictive error (MAE and RMSE) for five recommendation algorithms: PMF (Probabilistic Matrix Factorization), GAVG (Global average rating), UKNN (User-based K Nearest Neighbors), CAMF (Context-aware Matrix Factorization), FM (Factorization Machines). CAMF_CI, CAMF_CU and FM_CI are tuned with mood data at the item (CI) and user levels (CU), respectively, while the other three are benchmarks and are only tuned on a traditional user-item ratings matrix.

During the main task, participants were given step by step instructions in the form of interactive MoodPlay system tutorial. They were asked to enter at least three profile items (music bands) from a drop-down list, shown on the left in Fig. 1. In all conditions, this profile was used to generate a list of 5 recommendations, that were shown on the right side of the screen. Ratings were collected for the recommendation list as a whole and 5 individual items in the list. Participants were then allowed to interact freely with the system and generate as many intermediate recommendation lists as they wished. Once satisfied, they again rated the full list of items prior to finishing the MoodPlay interaction task. In our study, ratings were an indicator of users' perceived recommendation quality, or simply how well they liked suggested artists. Although music is subjective, and users may have different criteria for rating (e.g. expectations at a given moment, taste, current mood, similarity of suggestions to profile items), by comparing ratings across different conditions we can evaluate the impact of MoodPlay's features on the perceived quality of recommendations. To ensure that users spent sufficient time in the experiment, we displayed a non-numerical timer and gave users the opportunity to proceed to the post-study after at least 1.5 min of interaction.

¹² <https://www.mturk.com>

¹³ <http://www.qualtrics.com>

6.4. Participants

The 279 valid participants were equally distributed across all 4 conditions: 70, 69, 70 and 70. Studies lasted an average of 20 minutes and participants were paid an amount of \$1.00 per study. Age ranges of participants were reported from 18 to over 65, with an average range of 25–30. 52% were female. 13% did not finish college, 40% had a four-year college degree and 47% had a graduate degree. 74% were familiar with data visualization; 66% used a mouse for the interactive study and 34% had a trackpad. When asked about music tastes, 89% said they listen to music frequently. Reported use of streaming services such as Pandora was normally distributed. 71% of participants reported that they preferred a mix of popular and esoteric music. Participants were asked to rate the statement *I am a trusting person* on a scale of 1 to 5, in order to evaluate whether their trust in the recommendation system stems from their trust propensity or interaction with the system. The results were approximately evenly distributed across low, medium and high trust bins. During the design stage of this experiment, approximately 10 informal lab-based studies were also conducted and participants were interviewed to gauge their experiences with the system. Among the questions in the pre-study survey, we collected users' current mood by asking them "Which of the following best describes your current mood?". They choose one of the following options: (a) Sublime (e.g.: joyful, warm and tender moods), (b) Vital (e.g.: stimulating moods such as "lively", "energetic" or "fierce") or (c) Unease (e.g.: negative moods such as "sad", "tense" or "fearful"). The results indicated that most people felt to be in a sublime mood (56%), followed by vital (29%), and the fewest, unease(15%).

7. Results

We present our results in five subsections. We first provide details on how subjects interacted with the interface in the different conditions, in order to understand how the design decisions affected the user behavior. Then, in the second subsection, we present results in terms of the diversity and accuracy of the system, comparing artist streaming activity, rating and nDCG¹⁴ differences among conditions. In the third subsection we analyze the effect of mood in the results – both self-reported user mood prior to the study and artists' mood category. Next, in the fourth section we present qualitative results to understand subjective aspects of user behavior. Finally, in section five, we summarize the results by combining both quantitative and qualitative data into an integral analysis. This allows us to explain how visual and interactive aspects of each condition affect the results of objective and subjective metrics.

7.1. User behavior from log analysis

We recorded the amount of time users spent on the interface, but we found no significant differences among conditions. We also logged several user interactions with the system, most of which were clicks on different interface components as shown in Fig. 6. While some of these actions were available across all conditions since they were recorded on the user profile and recommendation panels (adding and removing artists in the profile list, playing music by clicking on artists in the recommendation list), other actions were available only in the *visually-enhanced* conditions (clicking on artist nodes and playing music through artist nodes in the visualization). Finally, two interactions were available only in the most advanced condition, where users could actually draw a trail when moving the avatar (creating and removing trail marks). We highlight two results from analyzing these actions and detailed statistics can be found in Table 3.

Preference elicitation. In MoodPlay, music artists were recommended based on their mood similarity to the artists in the user profile. We found

Table 3

Statistics describing user interactions with the interface in different conditions. Superscript numbers indicate conditions over which the significant difference was found. Significance is obtained via multiple *t*-test with Bonferroni correction, except for *# users who removed artists*, where it was obtained via multiple proportion test.

Statistic	Condition			
	(1)	(2)	(3)	(4)
# users	70	69	70	70
avg. artists added/user	4.47 ^{2, 3, 4}	3.78 ^{3, 4}	3.03	3.09
avg. artists removed/user	0.37 ^{3, 4}	0.21	0	0
# users who removed artists	16 ^{3, 4}	11 ^{3, 4}	0	1
avg. total interactions/user	18.11	18.91	24.27 ^{1, 2}	23.25 ^{1, 2}
avg. recommended lists/user	1.87 ²	1.53	2.31 ²	2.24 ^{1, 2}

that users added significantly more artists in conditions 1 ($M=4.47$, $SE=0.23$), $p < .002$, and 2 ($M=3.78$, $SE=0.13$), $p < .003$, than in condition 3 ($M=3.03$, $SE=0.02$) and 4 ($M=3.09$, $SE=0.09$). Furthermore, while 16 and 11 users removed artists from their profile in conditions 1 and 2 respectively, only one user removed an artist in condition 4 and no user did it in condition 3. In conditions 1 and 2 the only way that users could update their list of recommendations was by adding or removing artists in their profile. In conditions 3 and 4, users could update the recommendation list simply by moving the avatar in the interactive visualization. On average they generated more recommendation lists than users in conditions 1 and 2. Despite the ease to get a new set of recommendations in conditions 3 and 4 compared to 1 and 2, the results still show that users in all conditions were interested in exploring recommendations beyond the first list – either because they were curious, they enjoyed the system or were not content with the initial recommendations list.

Diversity. One of the most interesting results of our study is that the right amount of interaction functionality in a visual interface can promote diversity among the consumed items. Promoting diversity in recommender systems is one of the most important topics in the area (Ziegler et al., 2005), particularly helpful in preventing the creation of filter bubbles (Pariser, 2011). We measured this effect by comparing the number of unique artists rated and played per user in each condition. With respect to artists played, we considered "playing activity" in any component of the interface (visualization and recommendation panel) and in the recommendation list only, to make a fair comparison against condition 1. Plots in Fig. 7 show these distributions. Significant differences were assessed with Wilcoxon signed-rank tests since data departs from normality. We accounted for multiple comparisons with Bonferroni correction. The most important result is that condition 3 significantly outperforms all the other conditions in the three aforementioned metrics: amount of rated items ($M=10.59$, $S.E.=0.41$), $p < .001$, number of artists played anywhere ($M=10.71$, $S.E.=0.81$), $p = .002$, and artists played on the recommendation panel only ($M=10.61$, $S.E.=0.82$), $p < .003$. Also notable, condition 1 shows significantly more diversity than condition 2 in terms of unique artists rated ($M=8.56$, $S.E.=0.28$), $p < .001$, and played in the recommendation list ($M=7.63$, $S.E.=0.43$), $p < .02$.

Ranking. In addition to differences in diversity of explored artists, we analyzed differences in ranking among the different conditions. During the study, users had to rate an initial and a final list of recommendation, and they were free to rate more lists in between. We used the metric nDCG (Manning et al., 2008), since it is a common metric used in recommender systems (Parra and Sahebi, 2013). nDCG measures the gain of a recommendation discounted by the logarithm of its position in the list. This accumulated gain is high when relevant items (rated 4 or 5) appear at the top of the list and the non-relevant elements (rated 1, 2 or 3) are placed at the bottom. Table 4 shows the average nDCG of the first and last lists at each condition. To analyze the differences in nDCG ranking between conditions, we conducted multiple pairwise *t*-tests with Bonferroni correction. We found no differences in nDCGs at the initial

¹⁴ Normalized Discounted Cumulative Gain – a measure of ranking quality in terms of usefulness of recommended item based on its position in the recommendation list.

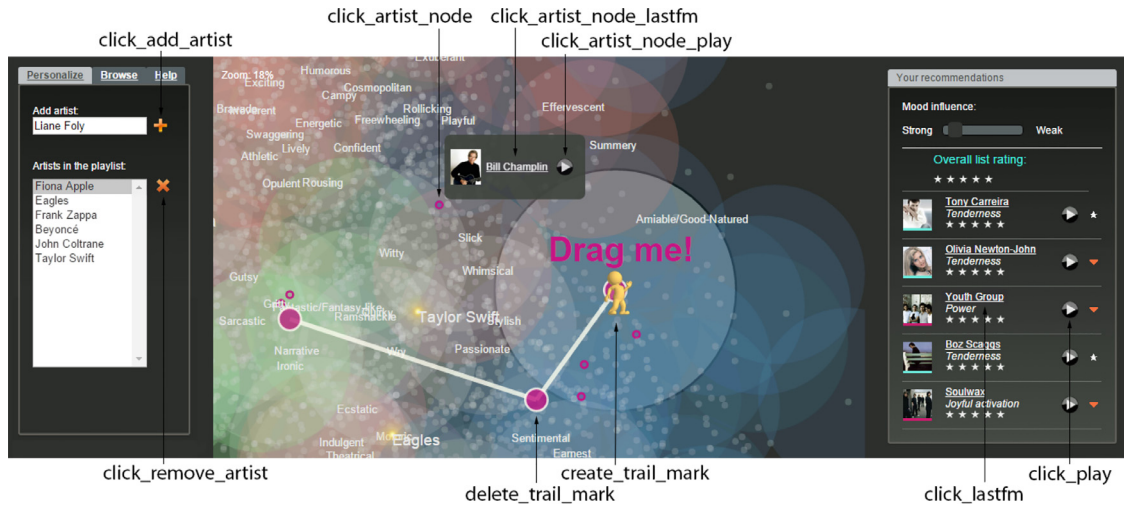


Fig. 6. User actions logged during the user study, contextualized on the interface of condition 4, which includes all system's features.

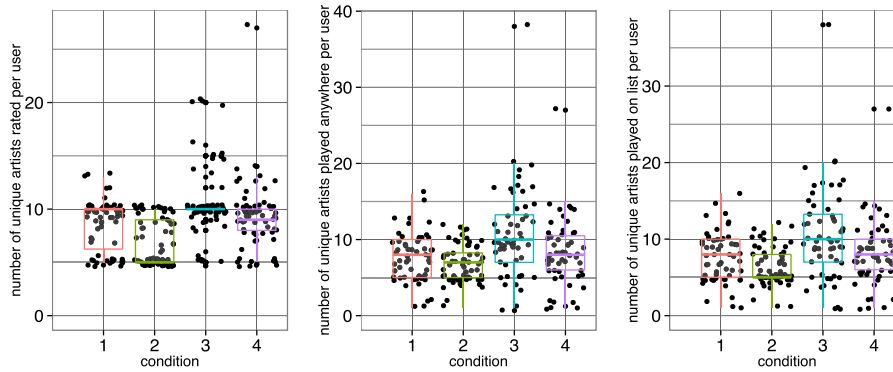


Fig. 7. Consumption of unique items per user: rating-based (left), and played-based interactions all over the interface (center) and on the recommendation list only (right).

Table 4

Normalized Discounted Cumulative Gain (NDCG) and standard error for the first and last rated recommendation list, per condition.

	NDCG and standard error per condition			
	1	2	3	4
First recommended list	0.58 ± 0.02	0.54 ± 0.02	0.58 ± 0.02	0.55 ± 0.02
Last recommended list	0.54 ± 0.03	0.47 ± 0.03	0.58 ± 0.02	0.53 ± 0.03

lists, but by comparing the nDCGs at the end of the study, we found that condition 3 had a significantly larger nDCG ($M=0.58, SE=0.02$) than condition 2 ($M=0.47, SE=0.02$), $p=0.048$. Since the recommendation algorithm was the same in all four conditions, only the visualization and interaction could explain the observed differences among conditions. Condition 2 provides a visualization which allows users to explore the dataset (artists) in terms of mood, but unlike condition 3, it does not allow them to update the recommendation list through the visual mood space. This might explain the better ranking quality that was observed in condition 3.

7.2. Mood analysis

We designed *Moodplay* to investigate the effect of music exploration based on mood categories upon user satisfaction. In this investigation, we also explored whether there is a connection between users' self-reported and the affective profile of the music that they listened to. During the pre-study, we collected users' self-reported mood as described in Section 6.3: Sublime, Vital or Uneasy.

For the sake of understanding the context, the distribution of artists' primary mood categories and users' self-reported mood are shown in Fig. 8. The axes of users' self-reported mood distribution are flipped to facilitate comparison with Table 5, which shows the average weights of each primary mood category for artists (Sublimity, Uneasy, Vitality, Style) versus users' self reported mood (Sublime, Unease, Vital). In addition, Table 5 compares two groups of data: artists that users added to their profile, and artists that were rated. We observe the following main trends in this analysis:

- In Fig. 8, we observe that sublimity is the most frequent mood in both: (a) artists' primary mood category, and (b) users' self-reported mood. This is followed by vitality and uneasy. In the case of artists, the least frequent primary mood is style, the one we created in our research to expand the GEMS model.
- Artists have weights in several mood categories. In the first row of Table 5, we show the average weight of each category over the whole artist dataset. The mood category with largest average weight is Sub-

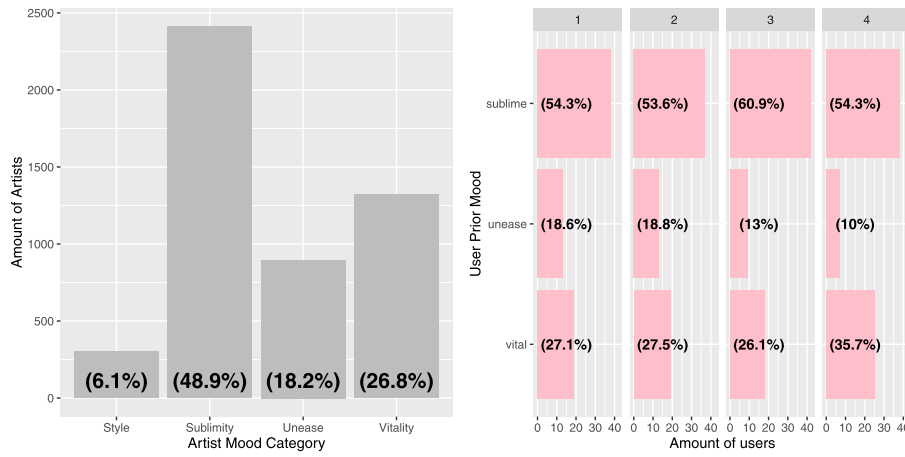


Fig. 8. Distributions of primary artists' mood (left) and users' self-reported moods in each condition, prior to the beginning of user study (right).

Table 5

Average weight of each artist mood (columns) in the whole dataset (first row), among artists added to user profile (rows 2–4), and among artists rated by users (rows 5–6). Statistical tests were performed column-wise. Within each column, we found statistically significant weights only among artists rated, showing that usage of the MoodPlay system actually changed the consumption depending on user's self-reported mood.

		Average Artists' Moods Weights			
	User-Reported Mood	Art. Sublimity	Art. Uneasy	Art. Vitality	Art. Style
Artists in User Profile	(overall dataset)	0.38	0.21	0.27	0.14
	Sublime	0.37	0.22	0.27	0.14
	Unease	0.36	0.24	0.27	0.13
	Vital	0.37	0.20	0.29	0.14
Artists Rated	Sublime	0.40	0.16	0.29^U	0.15
	Unease	0.38	0.20^{S, V}	0.26	0.16
	Vital	0.38	0.16	0.30^U	0.15

limity (0.38), followed by Vitality (0.27), Uneasy (0.21), and finally Style (0.14).

- In rows 2–4 of Table 5 we split the data based on users' self-reported mood and consider the artists added to user profiles (965 total). When comparing each artist primary mood category (columns) across the three potential users' reported moods (rows), we found no statistically significant differences. This means that, on average, users added artists with similar primary mood distribution to their profiles, independent of their self-reported mood.
- The last three rows in Table 5 show the average weights in each primary mood, for each users' self reported mood, but this time considering the artists rated by users (2,704 ratings in total). We found a couple of statistically significant differences. For example, users who had Unease as their reported mood were more likely to listen to/rate music with high uneasy mood (0.2) than users in Sublime (0.16) or Vital (0.16) self-reported mood. Additionally, users in either Sublime (0.29) or Vital (0.3) reported mood, rated/listened to artists with significantly higher Vitality than users who reported feeling unease (0.26).

7.3. Analysis of post-study survey

In the post-study survey, we analyzed user perception of the system (Table 6). As expected, perceived ease of use drops-off with higher interface complexity and confusion increases, with condition 1 being significantly less confusing than all the rest, and also easier to use than conditions 2 and 4. Interestingly, the difference is not significant when compared to condition 3, which is visually as complex as condition 2, but offers more controllable functionality (e.g. draggable user avatar).

Furthermore, we did not see clear differences in average ratings per condition, but the perception of accuracy in condition 3 is significantly higher than in condition 1. This result is very interesting, because we are using exactly the same recommendation algorithm in both conditions, but the perception of accuracy changes with the addition of visualization and draggable avatar. As expected, people also felt that condition 3 allowed them significantly more control than condition 1.

Several aspects were not perceived significantly different among conditions, such as trust (rows 1 and 2), diversity of recommendations and helpfulness of interface to compare moods of different artists. Nonetheless, users perceived condition 3 more helpful in understanding how recommendations were generated than condition 1.

All of these results indicate that the visual layout of moods and artists, with the addition of ability to re-position the avatar in the mood space and control the hybrid recommendation algorithm, gradually improve user experience. However, introduction of trails in condition 4 has a negative effect, most likely because of the cognitive overload. In addition, we suspect that trails may be perceived as limiting for the exploration. It is possible that users expected to receive recommendations based on the most recent trail point, whereas the system accounts for all previous trail points.

7.4. Qualitative analysis

In each condition, participants were asked in the post-study survey to leave feedback on their experience and give suggestions for improving the system. Table 7 lists representative comments, grouped by condition and sentiment. On the positive side, many users had fun using MoodPlay and enjoyed discovering new artists in different moods in conditions 3 and 4. Participants reported in all four conditions that the artist database

Table 6

Summary of the most relevant variables in the post-study survey. Numbers indicate average user agreement (on a scale from 1 to 100) with mean \pm S.E. Values in bold indicate significant difference over a condition indicated by the superscript number. Multiple comparisons were adjusted using Bonferroni correction.

Statement	Mean agreement and standard error per condition			
	1	2	3	4
I trusted recommendations from the system	37.1 \pm 3.6	44.6 \pm 3.5	48.8 \pm 3.5	38.4 \pm 3.6
Interaction with the interface increased my trust in the recommendations	43.4 \pm 3.6	47.1 \pm 3.8	49.4 \pm 3.8	39.2 \pm 3.7
The recommendations were diverse	60.9 \pm 3.3	65.3 \pm 3.3	68.6 \pm 3.3	59.9 \pm 3.3
The interface helped me understand and compare moods of different artists	49.4 \pm 3.5	55.7 \pm 3.5	55.7 \pm 3.3	46.3 \pm 3.4
The interface helped me understand how recommendations were generated	42.8 \pm 3.6	54.4 \pm 3.9	58.6¹ \pm 3.8	50.3 \pm 3.7
The interface allowed me to control the recommendations	42.7 \pm 3.6	53.8 \pm 3.4	63.8¹ \pm 3.9	52.7 \pm 3.5
The interface was confusing	23.1 \pm 3.1	45.3¹ \pm 4.1	46¹ \pm 3.9	52.6¹ \pm 3.9
Overall, the recommendations were accurate	36.2 \pm 3.6	40.7 \pm 3.6	49.8¹ \pm 3.7	38.7 \pm 3.5
The system was easy to use	73.9^{2,4} \pm 3.6	58.3 \pm 3.8	63.8 \pm 3.6	53.2 \pm 4
The interface was slow	22.2 \pm 3.6	32.2 \pm 3.84	31.9 \pm 3.6	37.8¹ \pm 3.2
The tutorial explained the system reasonably well	72.6⁴ \pm 3.3	62.9 \pm 3.4	65.8 \pm 3.2	57.5 \pm 3.6
By the end of the session I was satisfied with the recommendations	42.2 \pm 4.1	44.2 \pm 4	49.3 \pm 3.9	38 \pm 3.6

Table 7

Selected positive and negative user feedback grouped by experimental condition.

Cond.	Positive comments	Negative comments
1	All good. It was really fun. I enjoyed using this!	Add more bands/artists to the search- for example, neither Silversun Pickups nor Smashing Pumpkins were found to add to my list. The recommendations didn't seem to match the artists I chose. Show more information on how the mood of a song/artist is determined.
2	I think this could be a great tool. Good luck with the progress I am anxious to give it a try when it is finished I really liked this, it is a new concept that I've never seen. It helped introduce me to artists in different genres that I had never heard before and were very good. The mood cloud is awesome, and I didn't know there could be so many different music moods, that was great, but not being able to explore the artists within each specific mood circle causes some frustration. Making the cloud more dynamic to dragging and clicking would enhance the tool.	I put in 3 rappers and it gave me like oldies and pop songs. Genre plays roles in certain moods. It runs a little slow, should improve optimization for older computers. I really didn't understand it.
3	Really good player, i would change nothing it actually made me listen to a couple of artists i did not know about and liked their music. An interesting concept. I use Pandora a lot, and my stations are usually based off of my mood that day. This tool would be useful for randomization of choices of music. This is really cool, I do not listen to much music and I think this would help me find some new artists or even be used as a therapy tool.	Make the interface simpler and more concise. Speed up loading times It was slow and laggy and some of the recommendations didn't have a play button. I'd like the option to buy a track if I heard one I really liked, or to save a playlist if I really enjoyed it. Larger music selection, possibly change the strong week slider, to broad or specific to the particular mood you are feeling.
4	its a cool design Neat program! If I could practice with it more I think I would really enjoy it. It was excellent! Thanks to the developers for developing wonderful tool.	Some of recommended artists didn't relate to my mood close enough. There is a lot of text on the page and it's a little overwhelming. Instead of starting off with so many "moods," maybe just have 20 initially listed. Make the interface faster and smoother. There was too much choppiness when I was using the visualization tool.

was small, compared to commercial systems, and also had mixing of genres in the recommendation lists. In addition, visualization rendering was sluggish for some users. These problems can be addressed in the future by considering genre in the recommendation algorithm and by optimizing the visual design for even larger artist database.

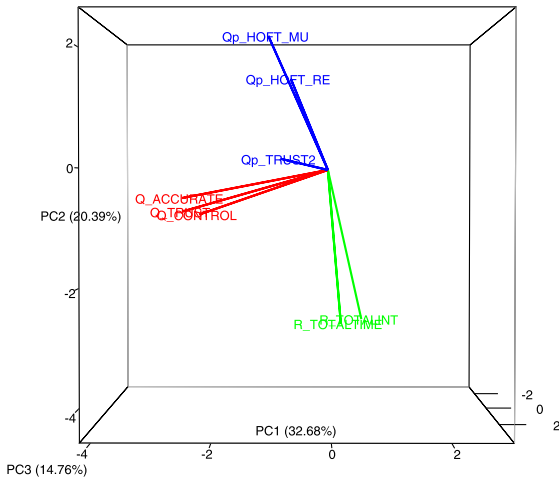
7.5. Connecting behavioral and perception measures

In order to explore the relationships between quantitative and qualitative experimental results collected during the user study, we performed Principal Component Analysis (PCA) (Kroonenberg, 2008), a technique for dimensionality reduction, over the variables that have shown significant effects in previous studies (Bostandjiev et al., 2012; Knijnenburg et al., 2012a; 2012b; Parra and Brusilovsky, 2015). Our analysis focused on the following variables:

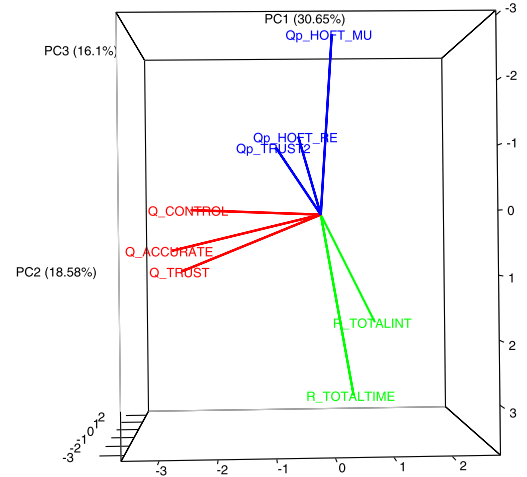
- Qp_HOFT_MU (pre-study question): *How often do you listen to music online?*
- Qp_HOFT_RE (pre-study question): *How often do you use recommender systems?*

- Qp_TRUST2 (pre-study question): *Are you a trusting person?*
- Q_ACCURATE (post-study question): *How accurate do you think the recommendation were?*
- Q_CONTROL (post-study question): *Did you feel in control of the interface?*
- Q_TRUST (post-study question): *How much do you trust the recommendations suggested during the study?*
- R_TOTALINT: Number of user interactions with the system (clicks, music plays, ratings, etc.)
- R_TOTALTIME: Duration of the user study.

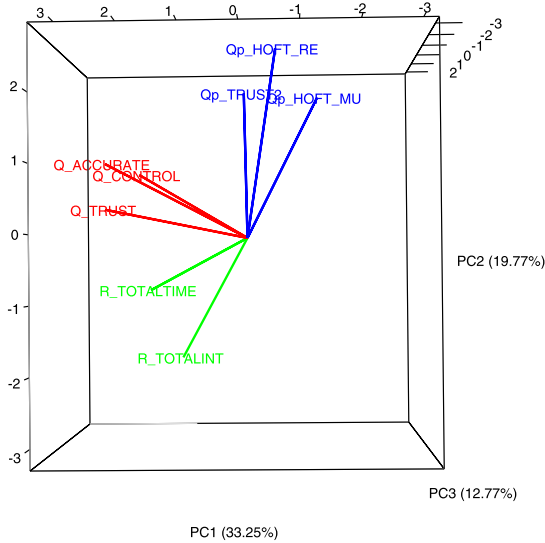
Fig. 9 shows biplots drawn from the output of PCA for each condition separately, with arrows denoting each of the variables in the above list. For interpreting PCA plots we used the guidelines described in Kroonenberg (2008): (a) the length of a vector represents the variance of that variable (within the principal components used in the biplot), (b) the cosine of the angle between a vector and an axis indicates the importance of the contribution of the variable to the corresponding principal component, (c) the cosine of the angle between pairs of variables indi-



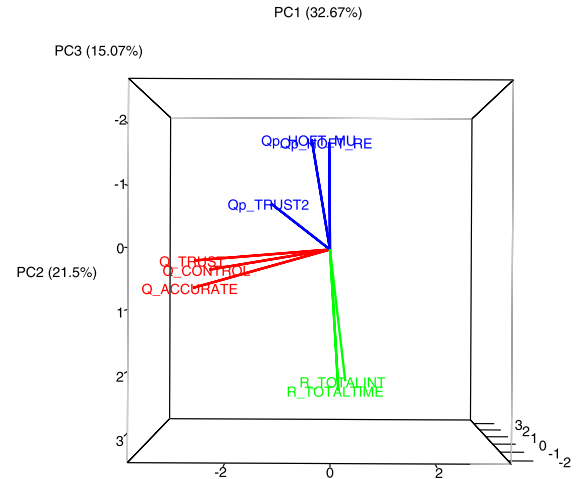
(a) Condition 1: no visualization, static list of recommendations.



(b) Condition 2: visualization without interactive update of recommendations.



(c) Condition 3: visualization with avatar and interactive recommendations.



(d) Condition 4: visualization with avatar, trails and interactive recommendations.

Fig. 9. 3D Biplots for Principal Component Analysis of the experiment variables: (i) pre-study survey (blue), (ii) post-study survey (red) and (iii) user interaction (green). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

cate how correlated they are, and (d) uncorrelated variables are at right angles with respect to each other.

Based on these guidelines, we plan to compare the influence of user interaction data such as time and clicks, on the results (accuracy, control, trust) versus the influence of user characteristics (trust propensity) and prior experience (familiarity with recommenders and music).

Condition 1. We observe that the red post-study variables (Q_CONTROL, Q_TRUST, Q_ACCURACY) are strongly correlated with each other. On the other side, they are almost orthogonal to variables R_TOTALINT and R_TOTALTIME, indicating that the amount of time spent and interaction with the interface had small or no effect on variables which have been shown to influence the final user satisfaction with the system (Knijnenburg et al., 2012a). Furthermore, blue-colored

post-study variables are loading in the same direction as users' pre-existing level of trust (Qp_TRUST2), but the short length of this vector tells us that its total variability is not well explained by principal components PC1, PC2 and PC3. Finally, the familiarity of users with music (Qp_HOFT_MU) and how often they listen to music online (Qp_HOFT_RE) are strongly correlated between each other, but they do not explain the variability of red post-study variables.

Condition 2. Similar to condition 1, the prior levels of trust (Qp_TRUST2) and familiarity with recommendation systems (Qp_HOFT_RE) load in the same direction as perceived control, trust and accuracy (Q_CONTROL, Q_TRUST, Q_ACCURACY) in PC1. However, Q_CONTROL departs from Q_ACCURATE and Q_TRUST on the projection over PC2, which implies that the user perception of

accuracy and trust diverted from the perception of control, compared to condition 1. This observation may explain some previous negative results in condition 2, such as the low number of unique artists rated and played compared to other interfaces. Since users were able to explore the mood space visualization, but they could not update the list of recommendations by interacting with it, their perception of control diverted from the perception of trust and accuracy.

Condition 3. The preference of users for this condition, shown in the previous analyses, can be explained holistically with the PCA plot. This is the only condition where Q_TOTALTIME and Q_TOTALINT load in the same left direction as Q_CONTROL, Q_ACCURATE and Q_TRUST in PC1 – the PC which explains most of the data variance in this condition. Notably, the acute angle between the red post-study variables and R_TOTALTIME shows that the amount of time that users spent on the interface explains the post-study variables, and especially Q_TRUST, more than in any other condition. This is an important result since it might indicate that both visualization and interaction combined help to increase the trust that users have in the recommendation system. By looking at the angles between variables in PC1, we can also see that Q_TRUST (final perceived trust) is more correlated to R_TOTALTIME than to Qp_TRUST2 (the initial level of user trust).

Condition 4. Similar to conditions 1 and 2, and unlike 3, this condition shows a disconnection between the amount of time and interaction on one hand, and the user perception in terms of Q_ACCURATE, Q_CONTROL and Q_TRUST on the other. The acute angle between Qp_TRUST2 and Q_TRUST shows that user's trust in the system is more likely to be determined by the inherent user trust than by the time spent interacting with the system. The plot also shows an opposite relation between the prior user experience with recommender and the amount of user interaction (since vectors form roughly a 180° angle), which may suggest that the user was confused and not fully taking advantage of advanced interface features. This could explain the drop in user satisfaction of this interface compared to condition 3.

8. Discussion

There is a fertile ground for expansions and branching of this research in several directions. The overarching idea is to build a system that recommends music according to user's musical taste, and guides the user from her current mood to the desired (target) mood. One caveat about the discussion, and our results, is that our experiment is based on single-session interactions with the system. Ideally, a longitudinal study in a real-world music listening context should be performed, and the authors are exploring possible ways to achieve this. We now discuss the research questions in light of the study results and summarize the development process and features of MoodPlay system. In the next section, we follow with limitations and avenues for future work.

(RQ1) What are the effects of interactive visualizations on the user experience with a recommender system, and what is the right amount of interaction for a music recommender?

The user study results clearly showed that the interface design and a certain combination of interactive features improve objective and perceived recommendation accuracy, as well as self-reported user satisfaction. We have shown that introduction of hybridization control for recommendation algorithm and the ability to move a user avatar, yielded positive effects across a variety of examined metrics. However, tracking of user mood states in the form of locations on the provenance trail introduced undesirable effects. First, we suspect that this increased system complexity beyond a comfortable threshold and caused cognitive overload, although another potential reason is that the trail did not match the users' mental model, preventing them from navigating the collection of artists effectively for the purpose of identifying relevant artists based on mood. Second, users who are unfamiliar with the system and participate in short listening sessions may be more inclined to rapidly investigate the mood space than those who are familiar with it and use the system in a more natural setting. Thus, the trail may have been per-

ceived as a limitation during relatively short experiment sessions. Nevertheless, modeling of changing mood preference is a fruitful research endeavor and our future work can address trails that follow smoother mood transitions, that are optional and are used during longer listening sessions.

(RQ2) Does affective information improve recommendation accuracy and user experience versus when it is not included?

Our analyses in Section 7.2 show evidence that both user mood prior to the study and the primary mood associated with the artists have an effect on the distribution of ratings. This difference in distributions is more pronounced between the interface with no visualization (condition 1) versus the other interfaces which have visualization (2, 3, 4). This result shows that making people aware of the mood of the artists combined with appropriate interactivity in a music recommender, can change the way they perceive the accuracy of the same music algorithm. In particular, when users' self-reported mood was *Unease* (associated with anger, sadness, depression) their overall rating decreased compared to users with different self-reported mood (*Sublime* and *Vital*) only at the conditions with visualization of moods. Moreover, when users conducted the study while in *Sublime* mood, they were more likely to provide higher ratings to artists with mood *Style*, while the interfaces with visualization received higher ratings for artists with primary mood *Unease*. Overall, there is a need for further research to establish a causal link between users' self-reported mood, artists' mood, and interface, but the results of our study hint towards a direction where all these variables play an important role in building recommender system interfaces.

Visualization of affective metadata for a music recommender system

Mood information has been visually represented in several preceding works, with the goal of enabling user selection of artists in desired moods. Typically, users choose a mood point in the visual space and the system plays music associated with the selected mood. To our knowledge, all up to date visualizations of moods for this purpose are based on a circumplex model of affect, which represents moods along valence and arousal dimensions (see Section 3.4). We argue that there exists a need to use a music specific mood model for the purpose of music recommendation, and propose an approach to fulfill it. Specifically, a dimensionality reduction method was applied to high-dimensional data containing mood-artist associations. It was then shown that a mood model, previously developed in music psychology research, emerges in the obtained two-dimensional, visual mood space.

To use this space during recommendation process, and help users to get a better understanding of it, several design aspects were addressed when incorporating it into an interactive system. Though not all explicitly tested in the user study: choice of colors, item sizes and transparency, dynamic labeling of mood nodes and node filtering based on mood categories, they all aim to explain the mood space and support the recommendation.

Supporting interaction, explanations and control over such visualization

The interaction with the system ranges from zooming and panning the visualization to explore the moods and artists, to controlling the hybridization level of the recommendation algorithm. Both user profile items and recommended artists are highlighted in the visualization, which helps users understand how those two sets are related based on moods. The explanation and exploration are further supported by providing links to external artist profiles, music streaming on demand, and displaying mood categories for recommended items. Moreover, a user avatar is positioned within the mood space at the centroid of user profiles items. The ability to move the avatar and form a trail of mood markers serves as a mechanism for modeling the change in user preference. This way users can control the recommendations, which are regenerated whenever the position is changed. A second way to influence the recommendation algorithm is by setting the ratio between mood and audio based filtering. This is achieved via simple slider control and visually explained to the user by re-sizing the catchment area around the avatar when the slider is moved. As the area increases, the recommendation

results depend more on the audio similarity to the profile items and less on the mood metadata.

9. Limitations and future work

Visual mood space. In this research, a key goal was to explore the use of a (visual) hierarchical mood model, capable of handling different granularity in the way moods are represented. Our model encompasses moods that are difficult to represent on the traditional valence arousal scale (Russell's Circumplex model). Our aim was also to enable location and exploration of moods via hierarchy, and therefore bring more diversity to research in mood scales. In the future work, however, we would like to compare our results with the traditional model as a benchmark.

Recommendation algorithm and interface. In terms of accuracy, users perceived the system as lacking prediction power, since the final survey resulted in an average evaluation among 36.2 - 49.8 out of a maximum of 100. We acknowledge this weakness of our implementation, and highlight the following aspects for improvement. First, building the mood space using a larger artist database could improve the mood-based component of the recommendation algorithm. Next, the MoodPlay system accounts for audio similarity when recommending music, but audio content analysis does not always accurately distinguish between music genres. Therefore, the recommendation algorithm can be improved by incorporating genre information. In addition, the system uses an audio similarity method that previously yielded satisfactory results but further investigation and comparison of algorithms could produce a better outcome. Finally, although the off-line algorithm evaluation found strong results when using state-of-the-art methods such as factorization machines and context-aware matrix factorization, for the user study in Mechanical Turk we used a simpler approach. Since we faced a cold-start problem (we had no previous preference information of the users), we relied on a version of our hybrid content and mood-based algorithm which led to less accurate predictions in some cases, but could be easily tuned with increased training data. Importantly, the comparison between interfaces and the effect of mood, which are the main aspects of this research are not affected by this.

Identifying user mood and musical preference. In the current implementation of MoodPlay, users build their profile by manually selecting several artists and we make recommendations based on the overall mood derived from that profile. We argue that it is acceptable to use mood data at the artist level, rather than on song level, because multiple moods associated with each artist in our database describe that artist's repertoire of songs. Nevertheless, using individual songs as an input and recommending tracks accordingly could perhaps yield greater precision. Another important consideration is that MoodPlay was introduced to users as a platform where they can create a list of favorite artists and be recommended new artists in similar moods. Hence, user's profile in MoodPlay reflects musical taste, but possibly it reflects the combination of both user's taste and mood. Although the core of our study was to explore how different interactive features affect the user experience with a mood aware recommendation system, and not to build a taste profile or auto-detect user's mood, it is important to note that taste and current preference based on mood can be treated as separate, but related parameters. Both can be determined by explicitly asking a user to provide the information, or implicitly, based on relevant data that has been collected automatically. In the following paragraphs, we focus on ideas for determining current preference as reflected by current mood, which could also capture the dynamics of a user's musical taste if tracked over longer period of time.

An approach used by some commercial recommendation systems (e.g. Spotify¹⁵, GooglePlay¹⁶) is to let users type in a mood or select it from a predefined list. This is not always an efficient method nor an

easy task for users given the large number of available mood tags. In particular, mood data in our system is very detailed and attempts to capture nuances that characterize different artists (e.g. rowdy, playful, graceful, elegant), whereas typical users may employ a smaller vocabulary and less specific words to describe mood of the music (e.g. happy, sad, energetic, calm). An alternative approach is to use the mood hierarchy embedded in MoodPlay to pick a top mood category from a list, followed by a subcategory and finally select specific mood. In this manner we can mitigate matching issues that might arise from granularity or choice of language for mood.

Implicitly determining user mood in an automated fashion on a granular level is even more challenging. However, an implicit approach can be effective if used with less specificity because it can entirely free the user from interaction. If greater granularity is desired, it can be improved by asking for some minimal input. Extensive research in affective computing, and discussed in Section 3.1 considers multiple mechanisms for improving data collection in MoodPlay. For example, a user's mood and current preference could be determined from contextual data such as: social media statuses, time of the day, weather, activity automatically inferred from GPS location or proximity of friends in the network, facial expression captured by mobile device or bodily functions measured by wearable devices. Another key benefit that arises from rich passive profiling data, is that mood can be inferred through behavior, and can serve to inform the system in a better way than a direct self-report from the user. In addition, there are indirect ways to measure users' mood, by asking them to choose certain colors, images or sound clips, which reflect how they are feeling at a certain time. For instance, the system AMARA (Affective Museum of Art Resource Agent) allows the users to explore art collections by asking them simple questions about their current feelings and interests in artwork (Park et al., 2012). This is especially useful when the user is in some state of denial about a current mood, or has any other metacognitive issue with reporting current mood.

Identifying target mood. It is not always desirable to play music that directly matches the listener's current mood. Instead, listeners may be interested in hearing music that changes how they feel. For example, happy music can uplift a listener who is feeling sad. Conversely, some people enjoy bitter-sweet music when sad, while at other times they might prefer springlike or playful songs. Target mood largely depends on a personal preference and current conditions, and therefore the recommender requires complex input or a highly advanced sensing algorithm to determine it.

Depending on the listening context and preference, the recommender can either suggest music in the target mood or find and follow a path from current to target mood. Commercial recommendation systems already offer playlists for different moods and activities (e.g. mellow, music for work or gym), which are effective for short term, action-based listening. However, to the best of our knowledge, there are no recommenders that allow transitions from one state to another or adapt to changes in how user feels or changes in listening context.

Adaptive recommendation systems have been an active area of research in recent years. Looking beyond their applications in entertainment, adaptive music recommenders can be of particular value in music therapy. Recent studies show positive effects of music on recovery of movement (e.g. in patients with stroke or Parkinson's disease) and speech (Thaut and McIntosh, 2014). Music therapy with the goal to modulate emotions has been studied less extensively, but its benefits to pain and mood management have been documented (Juslin and Laukka, 2004; Siedliecki and Good, 2006). The current version of MoodPlay has attracted interest from music therapists because its engaging interface can aid choosing music during therapy sessions for hospitalized children and elderly people with dementia. However, in a broad sense, adaptive recommendation systems can help to create a profound impact on a listener's well being, outside of formal therapeutic settings. By being able to continuously monitor feedback about a user's state and context, and adapt to changes, the therapeutic benefits of music can be improved. Our future work will look at these monitoring mechanisms with a view

¹⁵ <http://www.spotify.com>

¹⁶ <http://play.google.com/music>

to tuning MoodPlay to adapt readily to observed changes in patients' physical and emotional contexts.

Path from one mood state to another. The trail algorithm in MoodPlay can be viewed as a crude way to create a trail (path) from one mood state to another and generate recommendations accordingly. Through the evaluations we have performed with the system, we observed via numerous metrics that users preferred recommendations obtained by navigating the music collection freely, over the recommendations given by a trail based algorithm. We do not assume that this means that path-based computation of music recommendations are a bad thing, but rather that we need to improve our visual and interaction design for this aspect. User feedback comments and interviews lead us to believe that level of control of the path-based algorithm is key factor in user satisfaction. For example, users could be given a choice whether to use the system in an exploratory mode and freely navigate, or in preference modeling mode where they build the trail. Depending on a user's activity, available time and listening context, she could choose to engage more or less with the system. In cases when the user chooses to build a trail, recommending items along the trail (in between the trail marks) could provide more gradual change in the recommendations and possibly offer a more enjoyable listening experience during long sessions. Such a recommendation method would require evaluation in a more natural setting, and over a longer period of time.

Scalability. MoodPlay was developed on a database of 5000 artists. In comparison, online streaming services offer access to tens of millions of artists. In order to maximally scale the system, extensive work is needed in several areas. Even though there are efficient ways to perform dimensionality reduction of millions of data points, visualization design has to be adapted to accommodate such a large number. One simple way to achieve this is to show only limited number of artists on different zoom levels, according to some criteria such as popularity or relevance to a user based on preference data. A challenge in such a filtering method is to determine what artists the user is interested in seeing, and to show popular artists but also encourage discovery by introducing less known artists.

In a survey of dimensionality reduction methods, Fodor (2002) argues for the usefulness of dimensionality reduction for high dimensional data. The argument is that not all of the variables are "important" for understanding the phenomenon of interest in a high-dimensional data set. Our multidimensional data, where many moods are associated with each artist, poses a similar challenge. We apply a correspondence analysis, but note that other methods (for example genetic approaches, factor analysis or multi-dimensional scaling) may yield different layout of moods, and therefore different recommendations. However, an exhaustive comparison of these methods is beyond the scope this work. Furthermore, dimensionality reduction can introduce noise, creating clusters which did not exist in the original high-dimensional data. In future work, we believe it will be revealing to inspect the results of several dimensionality reduction techniques, with a tool such as the one introduced in Stahnke et al. (2016), and make advances in dynamic interactive labeling of the reduced space to help inform users of the underlying semantics of the space.

10. Conclusion

This paper presented and evaluated *MoodPlay* –a hybrid recommender system for musical artists which introduces a novel interactive visualization of moods and artists. The system supports explanation and control of a recommender system via manipulation of an avatar within the visualization. Design and implementation of an online experiment (N=279) was presented to evaluate the system through four conditions with varying degrees of visualization, interaction and control. Our key results have shown that interface design and a certain combination of interactive features improve objective and perceived recommendation accuracy, as well as self-reported user satisfaction with the recommender system (RQ1), and that making people aware of the typical mood of

an artist's music, combined with appropriate interactivity in a music recommender, can change the way users perceive the accuracy of the recommendation algorithm (RQ2).

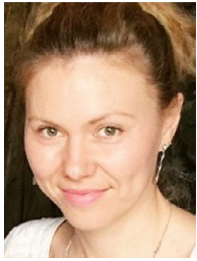
11. Funding Sources

This work was partially supported by the U.S. Army Research Laboratory under Cooperative Agreement No. W911NF-09-2-0053; The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of ARL, NSF, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation here on. The author Denis Parra was funded by the Chilean research agency CONICYT, FONDECYT grant number 11150783.

References

- Amelynck, D., Grachten, M., van Noorden, L., Leman, M., 2012. Toward e-motion-based music retrieval a study of affective gesture recognition. *T. Affect. Comput.* 3 (2), 250–259.
- Andjelkovic, I., Parra, D., O'Donovan, J., 2016. Moodplay: interactive mood-based music discovery and recommendation. In: *Proceedings of the 2016 Conference on User Modeling Adaptation and Personalization*. ACM, New York, NY, USA, pp. 275–279. doi:10.1145/2930238.2930280.
- Baccigalupo, C., Plaza, E., 2006. Case-based sequential ordering of songs for playlist recommendation. In: *Advances in Case-Based Reasoning*. Springer, pp. 286–300.
- Baltrunas, L., Amatriain, X., 2009. Towards time-dependant recommendation based on implicit feedback. *Workshop on Context-Aware Recommender Systems (CARS'09)*.
- Bentley, J.L., 1975. Multidimensional binary search trees used for associative searching. *Commun. ACM* 18 (9), 509–517. doi:10.1145/361002.361007.
- Bostandjiev, S., O'Donovan, J., Höllerer, T., 2012. Tasteweights: a visual interactive hybrid recommender system. In: *Proceedings of the Sixth ACM Conference on Recommender Systems*. ACM, pp. 35–42.
- Burke, R., 2002. Hybrid recommender systems: survey and experiments. *User Model User-adapt Interact* 12 (4), 331–370. doi:10.1023/A:1021240730564.
- Celma, O., Herrera, P., 2008. A new approach to evaluating novel recommendations. In: *Proceedings of the 2008 ACM Conference on Recommender Systems*. ACM, New York, NY, USA, pp. 179–186. doi:10.1145/1454008.1454038. <http://doi.acm.org.ezproxy.puc.cl/10.1145/1454008.1454038>
- Chen, L., Pu, P., 2009. Interaction design guidelines on critiquing-based recommender systems. *User Model User-adapt Interact* 19 (3), 167–206. doi:10.1007/s11257-008-9057-x.
- Collier, G.L., 2007. Beyond valence and activity in the emotional connotations of music. *Psychol. Music* 35 (1), 110–131. doi:10.1177/0305735607068890.
- Cunningham, S., Caulder, S., Grout, V., 2008. Saturday night or fever? context-aware music playlists. *Proc. Audio Mostly* <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.489.7774&rep=rep1&type=pdf>.
- Ekstrand, M.D., Riedl, J.T., Konstan, J.A., 2011. Collaborative filtering recommender systems. *Found. Trends Hum.-Comput. Interact.* 4 (2), 81–173. doi:10.1561/1100000009.
- Faltings, B., Pu, P., Torrens, M., Viappiani, P., 2004. Designing example-critiquing interaction. In: *Proceedings of the 9th International Conference on Intelligent User Interfaces*. ACM, pp. 22–29.
- Fernández-Tobías, I., Cantador, I., Plaza, L., 2013. An emotion dimensional model based on social tags: crossing folksonomies and enhancing recommendations. In: *E-Commerce and Web Technologies*. Springer, pp. 88–100.
- Fodor, I., 2002. A Survey of Dimension Reduction Techniques. Technical Report. Lawrence Livermore National Lab., CA (US).
- Frijda, N.H., 1993. Moods, Emotion Episodes and Emotions. In: Lewis, M., Haviland, J.M. (Eds.), *Handbook of Emotions*. New York: Guilford Press, pp. 381–403.
- George, J.M., 1996. Individual Differences and Behavior in Organizations. San Francisco: Jossey-Bass, p. 145.
- Glowinski, D., Dael, N., Camurri, A., Volpe, G., Mortillaro, M., Scherer, K., 2011. Toward a minimal representation of affective gestures. *IEEE Trans. Affect. Comput.* 2 (2), 106–118.
- Gonzalez, G., De La Rosa, J.L., Montaner, M., Delfin, S., 2007. Embedding emotional context in recommender systems. In: *Proceedings of the 2007 IEEE 23rd International Conference on Data Engineering Workshop*. IEEE, pp. 845–852.
- Gou, L., You, F., Guo, J., Wu, L., Zhang, X.L., 2011. Sfviz: interest-based friends exploration and recommendation in social networks. In: *Proceedings of the 2011 Visual Information Communication-International Symposium*. ACM, p. 15.
- Gretarsson, B., O'Donovan, J., Bostandjiev, S., Hall, C., Höllerer, T., 2010. Smallworlds: Visualizing social recommendations. In: *Computer Graphics Forum*, 29. Wiley Online Library, pp. 833–842.
- Griffiths, D., Cunningham, S., Weinle, J., 2013. A discussion of musical features for automatic music playlist generation using affective technologies. In: Delsing, K., Liljedahl, M. (Eds.), *Audio Mostly Conference*. ACM, pp. 13:1–13:4.
- Han, B.-j., Rho, S., Jun, S., Hwang, E., 2010. Music emotion classification and context-based music recommendation. *Multimed Tools Appl* 47 (3), 433–460.

- Hariri, N., Mobasher, B., Burke, R., 2012. Context-aware music recommendation based on latent topic sequential patterns. In: *Proceedings of the Sixth ACM Conference on Recommender Systems*. ACM, New York, NY, USA, pp. 131–138. doi:10.1145/2365952.2365979.
- He, C., Parra, D., Verbert, K., 2016. Interactive recommender systems: a survey of the state of the art and future research challenges and opportunities. *Expert Syst Appl* doi:10.1016/j.eswa.2016.02.013.
- Hijikata, Y., Kai, Y., Nishida, S., 2012. The relation between user intervention and user satisfaction for information recommendation. In: *Proceedings of the 27th Annual ACM Symposium on Applied Computing*. ACM, pp. 2002–2007.
- Hume, D., 2012. Emotions and moods. *Organ. Behav.* 258–297.
- Izard, C., 1977. *Human Emotions. Emotions, Personality, and Psychotherapy*. Springer. <https://books.google.com/books?id=A8K1G9FXMsEC>
- Janssen, J.H., van den Broek, E.L., Westerink, J.H.D.M., 2012. Tune in to your emotions: a robust personalized affective music player. *User Model. User-Adapt. Interact.* 22 (3), 255–279.
- Juslin, P.N., Laukka, P., 2004. Expression, perception, and induction of musical emotions: a review and a questionnaire study of everyday listening. *J. New Music Res.* 33 (3), 217–238. doi:10.1080/0929821042000317813.
- Kaminskas, M., Ricci, F., 2011. Location-adapted music recommendation using tags. In: *User Modeling, Adaption and Personalization*. Springer, pp. 183–194.
- Karg, M., Kühnlenz, K., Buss, M., 2010. Recognition of affect based on gait patterns. *IEEE Trans. Syst. Man Cybern. Part B* 40 (4), 1050–1061.
- Kleinsmith, A., Bianchi-Berthouze, N., 2007. Recognizing affective dimensions from body posture. In: *Proceedings of the 2Nd International Conference on Affective Computing and Intelligent Interaction*. Springer-Verlag, Berlin, Heidelberg, pp. 48–58.
- Knees, P., Schedl, M., Pohle, T., Widmer, G., 2006. An innovative three-dimensional user interface for exploring music collections enriched. In: *Proceedings of the 14th ACM International Conference on Multimedia*. ACM, New York, NY, USA, pp. 17–24. doi:10.1145/1180639.1180652.
- Knijnenburg, B.P., Bostandjiev, S., O'Donovan, J., Kobsa, A., 2012a. Inspectability and control in social recommenders. In: *Proceedings of the Sixth ACM Conference on Recommender Systems*. ACM, pp. 43–50.
- Knijnenburg, B.P., Willemsen, M.C., Gantner, Z., Soncu, H., Newell, C., 2012b. Explaining the user experience of recommender systems. *User Model User-adapt Interact* 22 (4–5), 441–504.
- Koelsch, S., 2009. A neuroscientific perspective on music therapy. *Ann. N. Y. Acad. Sci.* 1169 (1), 374–384.
- Konstan, J.A., Riedl, J., 2012. Recommender systems: from algorithms to user experience. *User Model User-adapt Interact* 22 (1–2), 101–123.
- Koren, Y., Bell, R., Volinsky, C., 2009. Matrix factorization techniques for recommender systems. *Comput. (Long Beach Calif)* 42 (8), 30–37. doi:10.1109/MC.2009.263.
- Kroonenberg, P.M., 2008. *Applied Multiway Data Analysis*, Vol. 702. John Wiley & Sons.
- Logan, B., 2004. Music recommendation from song sets. In: *Proceedings of the ISMIR*.
- Maillet, F., Eck, D., Desjardins, G., Lamere, P., et al., 2009. Steerable playlist generation by learning song similarity from radio station playlists. In: *Proceedings of the ISMIR*, pp. 345–350.
- Manning, C.D., Raghavan, P., Shtetz, H., 2008. *Introduction to Information Retrieval*. Cambridge University Press, New York, NY, USA.
- Masthoff, J., 2005. The pursuit of satisfaction: affective state in group recommender systems. In: *User Modeling 2005*. Springer, pp. 297–306.
- McFee, B., Lanckriet, G.R.G., 2011. Large-scale music similarity search with spatial trees. In: Klapuri, A., Leider, C. (Eds.), *ISMIR*. University of Miami, pp. 55–60. <http://dblp.uni-trier.de/db/conf/ismir/ismir2011.html#1185McFeeL11>
- McNee, S.M., Riedl, J., Konstan, J.A., 2006. Being accurate is not enough: how accuracy metrics have hurt recommender systems. In: *CHI'06 Extended Abstracts on Human Factors in Computing Systems*. ACM, pp. 1097–1101.
- Nagulendra, S., Vassileva, J., 2014. Understanding and controlling the filter bubble through interactive visualization: a user study. In: *Proceedings of the 25th ACM Conference on Hypertext and Social Media*. ACM, New York, NY, USA, pp. 107–115. doi:10.1145/2631775.2631811.
- Nijdam, N.A., 2005. Mapping emotion to color URL: <https://pdfs.semanticscholar.org/5f0d/e6e7bc1d5443243f9f42f2379db9639a933d.pdf>.
- O'Donovan, J., Smyth, B., Gretarsson, B., Bostandjiev, S., Höllerer, T., 2008. Peerchooser: visual interactive recommendation. In: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, pp. 1085–1088.
- Pariser, E., 2011. *The Filter Bubble: What the Internet is Hiding from you*. Penguin, UK.
- Park, H.-S., Yoo, J.-O., Cho, S.-B., 2006. A context-aware music recommendation system using fuzzy Bayesian networks with utility theory. In: *Fuzzy Systems and Knowledge Discovery*. Springer, pp. 970–979.
- Park, S.J., Chae, G., MacDonald, C., Stein, R., Wiedenbeck, S., Kim, J., 2012. Amara: the affective museum of art resource agent. In: *CHI '12 Extended Abstracts on Human Factors in Computing Systems*. ACM, pp. 1035–1038. doi:10.1145/2212776.2212379.
- Parra, D., Amatriain, X., 2011. Walk the talk: Analyzing the relation between implicit and explicit feedback for preference elicitation. In: *Proceedings of the 19th International Conference on User Modeling, Adaption, and Personalization*. Springer-Verlag, Berlin, Heidelberg, pp. 255–268. <http://dl.acm.org/citation.cfm?id=2021855.2021878>
- Parra, D., Brusilovsky, P., 2015. User-controllable personalization. *Int. J. Hum.-Comput. Stud.* 78 (C), 43–67. doi:10.1016/j.ijhcs.2015.01.007.
- Parra, D., Brusilovsky, P., Trattner, C., 2014. See what you want to see: visual user-driven approach for hybrid recommendation. In: *Proceedings of the 19th International Conference on Intelligent User Interfaces*. ACM, New York, NY, USA, pp. 235–240. doi:10.1145/2557500.2557542.
- Parra, D., Sahebi, S., 2012. *Recommender Systems: Sources of Knowledge and Evaluation Metrics*. In: *Advanced Techniques in Web Intelligence-2*. Springer Berlin/Heidelberg, pp. 149–175.
- Parra, D., Sahebi, S., 2013. Recommender systems: sources of knowledge and evaluation metrics. In: et al. (Eds.), J.V. (Ed.), *Advanced Techniques in Web Intelligence-2: Web User Browsing Behaviour and Preference Analysis*. Springer-Verlag, Berlin Heidelberg, pp. 149–175. 7
- Paul Ekman, R.J.D., 1994. *The Nature of Emotion: Fundamental Questions*. Oxford University Press.
- Pazzani, M.J., Billsus, D., 2007. *Content-based Recommendation Systems, The Adaptive Web*. Springer-Verlag, Berlin, Heidelberg, pp. 325–341.
- Petrushin, V.A., 2000. Emotion recognition in speech signal: experimental study, development, and application. In: *Proceedings of the ICSLP*, pp. 222–225.
- Pu, P., Faltungs, B., Chen, L., Zhang, J., Viappiani, P., 2011. Usability guidelines for product recommenders based on example critiquing research. In: Ricci, F., Rokach, L., Shapira, B., Kantor, P.B. (Eds.), *Recommender Systems Handbook*. Springer US, pp. 511–545. doi:10.1007/978-0-387-85820-3.16.
- Rho, S., Han, B.-j., Hwang, E., 2009. Svr-based music mood classification and context-based music recommendation. In: *Proceedings of the 17th ACM international conference on Multimedia*. ACM, pp. 713–716.
- Russell, J., 1980. A circumplex model of affect. *J. Pers. Soc. Psychol.* 39 (6), 1161–1178.
- Salkind, N.J., 2010. *Encyclopedia of Research Design*. SAGE Publications, Inc. 0 edition. doi:10.4135/9781412961288.
- Scherer, K.R., 1984. *On the Nature and Function of Emotion: A Component Process Approach*. Lawrence Erlbaum, Hillsdale, NJ, pp. 293–317.
- Scherer, K.R., Johnstone, T., Klasmeyer, G., 2003. Vocal expression of emotion. In: *Handbook of Affective Sciences*, pp. 433–456.
- Schimmack, U., Grob, A., 2000. Dimensional models of core affect: a quantitative comparison by means of structural equation modeling. *Eur. J. Pers.* 14 (4), 325–345.
- Shneiderman, B., 1996. The eyes have it: a task by data type taxonomy for information visualizations. In: *IEEE Symposium on Visual Languages*, 1996. *Proceedings. IEEE*, pp. 336–343.
- Siedlecki, S.L., Good, M., 2006. Effect of music on power, pain, depression and disability. *J. Adv. Nurs.* 54 (5), 553–562. doi:10.1111/j.1365-2648.2006.03860.x.
- Stahnke, J., Dörk, M., Müller, B., Thom, A., 2016. Probing projections: interaction techniques for interpreting arrangements and errors of dimensionality reductions. *IEEE Trans. Vis. Comput. Graph* 22 (1), 629–638.
- Stober, S., Nürnberger, A., 2013. Adaptive music retrieval—a state of the art. *Multimedia Tools Appl.* 65 (3), 467–494. doi:10.1007/s11042-012-1042-z.
- Thaut, M.H., McIntosh, G.C., 2014. Neurologic music therapy in stroke rehabilitation. *Curr. Phys. Med. Rehabil. Rep.* 2 (2), 106–113. doi:10.1007/s40141-014-0049-y.
- Tintarev, N., Masthoff, J., 2011. Designing and evaluating explanations for recommender systems. In: *Recommender Systems Handbook*. Springer, pp. 479–510.
- Tkalčič, M., Burnik, U., Košir, A., 2010. Using affective parameters in a content-based recommender system for images. *User Model User-adapt Interact* 20 (4), 279–311.
- Tkalcic, M., Kosir, A., Tasic, J., 2011. Affective recommender systems: the role of emotions in recommender systems. In: *Proceedings of the RecSys 2011 Workshop on Human Decision Making in Recommender Systems*. Citeseer, pp. 9–13.
- Västfjäll, D., 2002. Emotion induction through music: a review of the musical mood induction procedure. *Musicae Scientiae* 5 (1 suppl), 173–211.
- Verbert, K., Parra, D., Brusilovsky, P., Duval, E., 2013. Visualizing recommendations to support exploration, transparency and controllability. In: *Proceedings of the 2013 International Conference on Intelligent User Interfaces*. ACM, New York, NY, USA, pp. 351–362. doi:10.1145/2449396.2449442.
- Wang, X., Rosenblum, D., Wang, Y., 2012. Context-aware mobile music recommendation for daily activities. In: *Proceedings of the 20th ACM International Conference on Multimedia*. ACM, pp. 99–108.
- Weiss, H.M., Cropanzano, R., 1996. Affective Events Theory: A theoretical discussion of the structure, causes and consequences of affective experiences at work. In: Staw, B.M., Cummings, L.L. (Eds.), *In: Research in organizational behavior: An annual series of analytical essays and critical reviews*, 18. Elsevier Science/JAI Press, US, pp. 1–74.
- Wu, S., Falk, T.H., Chan, W.-Y., 2011. Automatic speech emotion recognition using modulation spectral features. *Speech Commun.* 53 (5), 768–785.
- Yang, Y.-H., Lin, Y.-C., Cheng, H.T., Chen, H.H., 2008. Mr. emo: music retrieval in the emotion plane. In: El-Saddik, A., Vuong, S., Griwodz, C., Bimbo, A.D., Candan, K.S., Jaimas, A. (Eds.), *ACM Multimedia*. ACM, pp. 1003–1004. <http://dblp.uni-trier.de/db/conf/mm/mm2008.html#YangLCC08>
- Yu, F., Chang, E., qing Xu, Y., yeung Shum, H., 2001. Emotion detection from speech to enrich multimedia content. In: *Proceedings of the Second IEEE Pacific-Rim Conference on Multimedia*, pp. 550–557.
- Zentner, M., EEROLA, T., 2011. Self-report measures and models. *Handbook of Music and Emotion: Theory, Research, Applications*.
- Zentner, M., Grandjean, D., Scherer, K.R., 2008. Emotions evoked by the sound of music: characterization, classification, and measurement. *Emotion* 8 (4), 494–521.
- Zhao, S., Zhou, M.X., Zhang, X., Yuan, Q., Zheng, W., Fu, R., 2011. Who is doing what and when: social map-based recommendation for content-centric social web sites. *ACM Trans. Intell. Syst. Technol. (TIST)* 3 (1), 5.
- Zheng, Y., Mobasher, B., Burke, R., 2015. Carskit: A java-based context-aware recommendation engine. In: *Proceedings of the 2015 IEEE International Conference on Data Mining Workshop (ICDMW)*, pp. 1668–1671. doi:10.1109/ICDMW.2015.222.
- Ziegler, C.-N., McNee, S.M., Konstan, J.A., Lausen, G., 2005. Improving recommendation lists through topic diversification. In: *Proceedings of the 14th International Conference on World Wide Web*. ACM, pp. 22–32.
- van der Zwaag, M.D., Janssen, J.H., Westerink, J.H.D.M., 2013. Directing physiology and mood through music: validation of an affective music player. *T. Affect. Comput.* 4 (1), 57–68.



Ivana Andjelkovic works in music and technology, approaching the field through visualization of music related data, music recommendation and audio signal processing. She holds B.S. and M.S. degrees in Computer Science and a Ph.D. in Media Arts and Technology from University of California, Santa Barbara. She currently works as a Senior Audio Software Engineer, developing algorithms for audio manipulation.



Denis Parra is Assistant Professor at the Department of Computer Science, School of Engineering at PUC Chile. Professor Parra holds a B.Eng. from Universidad Austral de Chile and a Ph.D. from the University of Pittsburgh, PA, USA. His main research interests are Recommender Systems, applications of Data Mining and Machine Learning, Information Visualization and Social Networks. Professor Parra has published in important journals in the area such as IJHCS, ComCom, ESWA and ACM TiiS, as well as in IUI, UMAP, Hypertext, RecSys, CSCW, and WWW conferences. He leads the SocVis Laboratory at PUC Chile.



John O'Donovan is an Associate Researcher at the Department of Computer Science, University of California Santa Barbara, where he co-directs the FourEyes laboratory. His research background is in AI and HCI, with a focus on recommender systems. He has published more than 60 research papers in peer reviewed conferences and journals in this area. His research on human computer interaction and intelligent interfaces has won multiple best paper awards. He received a highest impact paper award from ACM IUI in 2017. He served as general co-chair of ACM IUI conference for 2016, and program co-chair of ACM RecSys for 2018.