Learning embodied semantics via music and dance semiotic correlations

Francisco Afonso Raposo^{a,c}, David Martins de Matos^{a,c}, Ricardo Ribeiro^{b,c}

^aInstituto Superior Técnico, Universidade de Lisboa, Av. Rovisco Pais, 1049-001 Lisboa, Portugal

^bInstituto Universitário de Lisboa (ISCTE-IUL), Av. das Forças Armadas, 1649-026 Lisboa, Portugal

^cINESC-ID Lisboa, R. Alves Redol 9, 1000-029 Lisboa, Portugal

Abstract

Music semantics is embodied, in the sense that meaning is biologically mediated by and grounded in the human body and brain. This embodied cognition perspective also explains why music structures modulate kinetic and somatosensory perception. We leverage this aspect of cognition, by considering dance as a proxy for music perception, in a statistical computational model that learns semiotic correlations between music audio and dance video. We evaluate the ability of this model to effectively capture underlying semantics in a cross-modal retrieval task. Quantitative results, validated with statistical significance testing, strengthen the body of evidence for embodied cognition in music and show the model can recommend music audio for dance video queries and vice-versa.

1. Introduction

Recent developments in human embodied cognition posit a learning and understanding mechanism called "conceptual metaphor" (Lakoff, 2012), where knowledge is derived from repeated patterns of experience. Neural circuits in the brain are substrates for these metaphors (Lakoff, 2014) and, therefore, are the drivers of semantics. Semantic grounding can be understood as the inferences which are instantiated as activation of these learned neural circuits.

^{*}Corresponding author: Tel.: +351-213-100-313

Email address: francisco.afonso.raposo@tecnico.ulisboa.pt (Francisco Afonso Raposo)

While not using the same abstraction of conceptual metaphor, other theories of embodied cognition also cast semantic memory and inference as encoding and activation of neural circuitry, differing only in terms of which brain areas are the core components of the biological semantic system (Kiefer & Pulvermüller, 2012; Ralph et al., 2017). The common factor between these accounts of embodied cognition is the existence of transmodal knowledge representations, in the sense that circuits are learned in a modality-agnostic way. This means that correlations between sensory, motor, linguistic, and affective embodied experiences create circuits connecting different modalityspecific neuron populations. In other words, the statistical structure of human multimodal experience, which is captured and encoded by the brain, is what defines semantics. Music semantics is no exception, also being embodied and, thus, musical concepts convey meaning in terms of somatosensory and motor concepts (Koelsch et al., 2019; Korsakova-Kreyn, 2018; Leman & Maes, 2014).

The statistical and multimodal imperative for human cognition has also been hinted at, at least in some form, by research across various disciplines, such as in aesthetics (Cook, 2000; Davies, 1994; Kivy, 1980; Kurth, 1991; Scruton, 1997), semiotics (Azcárate, 2011; Bennett, 2008; Blanariu, 2013; Lemke, 1992), psychology (Brown & Jordania, 2011; Dehaene & Cohen, 2007; Eitan & Granot, 2006; Eitan & Rothschild, 2011; Frego, 1999; Krumhansl & Schenck, 1997; Larson, 2004; Roffler & Butler, 1968; Sievers et al., 2013; Phillips-Silver & Trainor, 2007; Styns et al., 2007; Wagner et al., 1981), and neuroscience (Fujioka et al., 2012; Janata et al., 2012; Koelsch et al., 2019; Nakamura et al., 1999; Nozaradan et al., 2011; Penhune et al., 1998; Platel et al., 1997; Spence & Driver, 1997; Stein et al., 1995; Widmann et al., 2004; Zatorre et al., 1994), namely, for natural language, music, and dance. In this work, we are interested in the semantic link between music and dance (movement-based expression). Therefore, we leverage this multimodal aspect of cognition by modeling expected semiotic correlations between these modalities. These correlations are expected because they are mainly surface realizations of cognitive processes following embodied cognition. This framework implies that there is a degree of determinism underlying the relationship between music and dance, that is, dance design and performance are heavily shaped by music. This evident and intuitive relationship is even captured in some natural languages, where words for music and dance are either synonyms or the same (Baily, 1985). In this work, we claim that, just like human semantic cognition is based on multimodal statistical structures,

joint semiotic modeling of music and dance, through statistical computational approaches, is expected to provide some light regarding the semantics of these modalities as well as provide intelligent technological applications in areas such as multimedia production. That is, we can automatically learn the symbols/patterns (semiotics), encoded in the data representing human expression, which correlate across several modalities. Since this correlation defines and is a manifestation of underlying cognitive processes, capturing it effectively uncovers semantic structures for both modalities.

Following the calls for technological applications based on sensorimotor aspects of semantics (Leman, 2010; Matyja, 2016), this work leverages semiotic correlations between music and dance, represented as audio and video, respectively, in order to learn latent cross-modal representations which capture underlying semantics connecting these two modes of communication. These representations are quantitatively evaluated in a cross-modal retrieval task. In particular, we perform experiments on a 592 music audio-dance video pairs dataset, using Multi-view Neural Networks (MVNNs), and report 75% rank accuracy and 57% pair accuracy instance-level retrieval performances and 26% Mean Average Precision (MAP) class-level retrieval performance, which are all statistically very significant effects (p-values < 0.01). We interpret these results as further evidence for embodied cognition-based music semantics. Potential end-user applications include, but are not limited to, the automatic retrieval of a song for a particular dance or choreography video and vice-versa. To the best of our knowledge, this is the first instance of such a joint music-dance computational model, capable of capturing semantics underlying these modalities and providing a connection between machine learning of these multimodal correlations and embodied cognition perspectives.

The rest of this paper is structured as follows: Section 2 reviews related work on embodied cognition, semantics, and semiotics, motivating this approach based on evidence taken from research in several disciplines; Section 3 details the experimental setup, including descriptions of the evaluation task, MVNN model, dataset, features, and preprocessing; Section 4 presents the results; Section 5 discusses the impact of these results; and Section 6 draws conclusions and suggests future work.

2. Related work

Conceptual metaphor (Lakoff, 2012) is an abstraction used to explain the relational aspect of human cognition as well as its biological implementation in the brain. Experience is encoded neurally and frequent patterns or correlations encountered across many experiences define conceptual metaphors. That is, a conceptual metaphor is a link established in cognition (often subconsciously) connecting concepts. An instance of such a metaphor implies a shared meaning of the concepts involved. Which metaphors get instantiated depends on the experiences had during a lifetime as well as on genetically inherited biological primitives (which are also learned based on experience, albeit across evolutionary time scales). These metaphors are physically implemented as neural circuits in the brain which are, therefore, also learned based on everyday experience. The learning process at the neuronal level of abstraction is called "Hebbian learning", where "neurons that fire together, wire together" is the motto (Lakoff, 2014). Semantic grounding in this theory, called Neural Theory of Thought and Language (NTTL), which is understood as the set of semantic inferences, manifests in the brain as firing patterns of the circuits encoding such metaphorical inferences. These semantics are, therefore, transmodal: patterns of multimodal experience dictate which circuits are learned. Consequently, semantic grounding triggers multimodal inferences in a natural, often subconscious, way. Central to this theory is the fact that grounding is rooted in primitive concepts, that is, inference triggers the firing of neuron populations responsible for perception and action/coordination of the material body interacting in the material world. These neurons encode concepts like movement, physical forces, and other bodily sensations, which are mainly located in the somatosensory and sensorimotor systems (Desai et al., 2011; Cespedes-Guevara & Eerola, 2018; Koelsch et al., 2019; Lakoff, 2014). Other theories, such as the Controlled Semantic Cognition (CSC) (Ralph et al., 2017), share this core multimodal aspect of cognition but defend that a transmodal hub is located in the Anterior Temporal Lobes (ATLs) instead. Kiefer & Pulvermüller (2012) review and compare several semantic cognition theories and argue in favor of the embodiment views of conceptual representations, which are rooted in transmodal integration of modality-specific (e.g., sensory and motor) features. In the remainder of this section, we review related work providing evidence for the multimodal nature of cognition and the primacy of primitive embodied concepts in music.

Aesthetics suggests that musical structures evoke emotion through isomorphism with human motion (Cook, 2000; Davies, 1994; Kivy, 1980; Scruton, 1997) and that music is a manifestation of a primordial "kinetic energy" and a play of "psychological tensions" (Kurth, 1991). Blanariu (2013) claims that, even though the design of choreographies is influenced by culture, its aesthetics are driven by "pre-reflective" experience, i.e., unconscious processes driving body movement expression. The choreographer interprets the world (e.g., a song), via "kinetic thinking" (von Laban & Ullmann, 1960), which is materialized in dance in such a way that its surface-level features retain this "motivating character" or "invoked potential" (Peirce, 1991), i.e., the conceptual metaphors behind the encoded symbols can still be accessible. The symbols range from highly abstract cultural encodings to more concrete patterns, such as movement patterns in space and time such as those in abstract (e.g., non-choreographed) dance (Blanariu, 2013). Bennett (2008) characterizes movement and dance semantics as being influenced by both physiological, psychological, and social factors and based on space and forces primitives. In music, semantics is encoded symbolically in different dimensions (such as timbral, tonal, and rhythmic) and levels of abstraction (Juslin, 2013; Schlenker, 2017). These accounts of encoding of meaning imply a conceptual semantic system which supports several denotations (Blanariu, 2013), i.e., what was also termed an "underspecified" semantics (Schlenker, 2017). The number of possible denotations for a particular song can be reduced when considering accompanying communication channels, such as dance, video, and lyrics (Schlenker, 2017). Natural language semantics is also underspecified according to this definition, albeit to a much lower degree. Furthermore, Azcárate (2011) emphasizes the concept of "intertextuality" as well as text being a "mediator in the semiotic construction of reality". Intertextuality refers to the context in which a text is interpreted, allowing meaning to be assigned to text (Lemke, 1992). This context includes other supporting texts but also history and culture as conveyed by the whole range of semiotic possibilities, i.e., via other modalities (Lemke, 1992). That is, textual meaning is also derived via multimodal inferences, which improve the efficacy of communication. This "intermediality" is a consequence of human cognitive processes based on relational thinking (conceptual metaphor) that exhibit a multimodal and contextualized inferential nature (Azcárate, 2011). Peirce (1991) termed this capacity to both encode and decode symbols, via semantic inferences, as "abstractive observation", which he considered to be a feature required to learn and interpret by means of experience, i.e., required

for being an "intelligent consciousness".

Human behaviour reflects this fundamental and multimodal aspect of cognition, as shown by psychology research. For instance, Eitan & Rothschild (2011) found several correlations between music dimensions and somatosensoryrelated concepts, such as sharpness, weight, smoothness, moisture, and temperature. People synchronize walking tempo to the music they listen to and this is thought to indicate that the perception of musical pulse is internalized in the locomotion system (Styns et al., 2007). The biological nature of the link between music and movement is also suggested in studies that observed pitch height associations with vertical directionality in 1-year old infants (Wagner et al., 1981) and with perceived spatial elevation in congenitally blind subjects and 4- to 5-year old children who did not verbally make those associations (Roffler & Butler, 1968). Tension ratings performed by subjects independently for either music or a corresponding choreography yielded correlated results, suggesting tension fluctuations are isomorphically manifested in both modalities (Frego, 1999; Krumhansl & Schenck, 1997). Phillips-Silver & Trainor (2007) showed that the perception of "beat" is transferable across music and movement for humans as young as 7 months old. Eitan & Granot (2006) observed a kind of music-kinetic determinism in an experiment where music features were consistently mapped onto kinetic features of visualized human motion. Sievers et al. (2013) found further empirical evidence for a shared dynamic structure between music and movement in a study that leveraged a common feature between these modalities: the capacity to convey affective content. Experimenters had human subjects independently control the shared parameters of a probabilistic model, for generating either piano melodies or bouncing ball animations, according to specified target emotions: angry, happy, peaceful, sad, and scared. Similar emotions were correlated with similar slider configurations across both modalities and different cultures: American and Kreung (in a rural Cambodian village which maintained a high degree of cultural isolation). The authors argue that the isomorphic relationship between these modalities may play an important role in evolutionary fitness and suggest that music processing in the brain "recycles" (Dehaene & Cohen, 2007) other areas evolved for older tasks, such as spatiotemporal perception and action (Sievers et al., 2013). Brown & Jordania (2011) suggest that this capacity to convey affective content is the reason why music and movement are more cross-culturally intelligible than language. A computational model for melodic expectation, which generated melody completions based on tonal movement driven by physical forces

(gravity, inertia, and magnetism), outperformed every human subject, based on intersubject agreement (Larson, 2004), further suggesting semantic inferences between concepts related to music and movement/forces.

There is also neurological evidence for multimodal cognition and, in particular, for an underlying link between music and movement. Certain brain areas, such as the superior colliculus, are thought to integrate visual, auditory, and somatosensory information (Spence & Driver, 1997; Stein et al., 1995). Widmann et al. (2004) observed evoked potentials when an auditory stimulus was presented to subjects together with a visual stimulus that infringed expected spatial inferences based on pitch. The engagement of visuospatial areas of the brain during music-related tasks has also been extensively reported (Nakamura et al., 1999; Penhune et al., 1998; Platel et al., 1997; Zatorre et al., 1994). Furthermore, neural entrainment to beat has been observed as β oscillations across auditory and motor cortices (Fujioka et al., 2012; Nozaradan et al., 2011). Moreover, Janata et al. (2012) found a link between the feeling of "being in the groove" and sensorimotor activity. Korsakova-Kreyn (2018) also explains music semantics from an embodied cognition perspective, where tonal and temporal relationships in music artifacts convey embodied meaning, mainly via modulation of physical tension. These tonal relationships consist of manipulations of tonal tension, a core concept in musicology, in a tonal framework (musical scale). Tonal tension is physically perceived by humans as young as one-day-old babies (Virtala et al., 2013), which further points to the embodiment of music semantics, since tonal perception is mainly biologically driven. The reason for this may be the "principle of least effort", where consonant sounds consisting of more harmonic overtones are more easily processed and compressed by the brain than dissonant sounds, creating a more pleasant experience (Bidelman & Krishnan, 2009, 2011). Leman (2007) also emphasizes the role of kinetic meaning as a translator between structural features of music and semantic labels/expressive intentions, i.e., corporeal articulations are necessary for interpreting music. Semantics are defined by the mediation process when listening to music, i.e., the human body and brain are responsible for mapping from the physical modality (audio) to the experienced modality (Leman, 2010). This mediation process is based on motor patterns which regulate mental representations related to music perception. This theory, termed Embodied Music Cognition (EMC), also supports the idea that semantics is motivated by affordances (action), i.e., music is interpreted in a (kinetic) way that is relevant for functioning in a physical environment. Furthermore, EMC also states that decoding music expressiveness in performance is a sense-giving activity (Leman & Maes, 2014), which falls in line with the learning nature of NTTL. The Predictive Coding (PC) framework of Koelsch et al. (2019) also points to the involvement of transmodal neural circuits in both prediction and prediction error resolution (active inference) of musical content. The groove aspect of music perception entails an active engagement in terms of proprioception and interoception, where sensorimotor predictions are inferenced (by "mental action"), even without actually moving. In this framework, both sensorimotor and autonomic systems can also be involved in resolution of prediction errors.

Recently, Pereira et al. (2018) proposed a method for decoding neural representations into statistically-modeled semantic dimensions of text. This is relevant because it shows statistical computational modeling (in this instance, ridge regression) is able to robustly capture language semantics in the brain, based on functional Magnetic Resonance Imaging (fMRI). This language-brainwaves relationship is an analogue to the music-dance relationship in this work. The main advantage is that, theoretically, brain activity will directly correlate to stimuli, assuming we can perfectly decode it. Dance, however, can be viewed as an indirect representation, a kinetic proxy for the embodied meaning of the music stimulus, which is assumed to be encoded in the brain. This approach provides further insights motivating embodied cognition perspectives, in particular, to its transmodal aspect. fMRI data was recorded for three different text concept presentation paradigms: using it in a sentence, pairing it with a descriptive picture, and pairing it with a word cloud (several related words). The best decoding performance across individual paradigms was obtained with the data recorded in the picture paradigm, illustrating the role of intermediality in natural language semantics and cognition in general. Moreover, an investigation into what voxels were most informative for decoding, revealed that they were from widely distributed brain areas (language 21%, default mode 15%, task-positive 23%, visual 19%, and others 22%), as opposed to being focalized in the language network, further suggesting an integrated semantic system distributed across the whole brain. A limitation of that approach in relation to the one proposed here is that regression is performed for each dimension of the text representation independently, failing to capture how all dimensions jointly covary across both modalities.

3. Experimental setup

As previously stated, multimedia expressions referencing the same object (e.g., audio and dance of a song) tend to display semiotic correlations reflecting embodied cognitive processes. Therefore, we design an experiment to evaluate how correlated these artifact pairs are: we measure the performance of cross-modal retrieval between music audio and dance video. The task consists of retrieving a sorted list of relevant results from one modality, given a query from another modality. We perform experiments in a 4-fold cross-validation setup and report pair and rank accuracy scores (as done by Pereira et al. (2018)) for instance-level evaluation and MAP scores for class-level evaluation. The following sections describe the dataset (Section 3.1), features (Section 3.2), preprocessing (Section 3.3), MVNN model architecture and loss function (Section 3.4), and evaluation details (Section 3.5).

3.1. Dataset

We ran experiments on a subset of the *Let's Dance* dataset of 1000 videos of dances from 10 categories: ballet, breakdance, flamenco, foxtrot, latin, quickstep, square, swing, tango, and waltz (Castro et al., 2018). This dataset was created in the context of dance style classification based on video. Each video is 10s long and has a rate of 30 frames per second. The videos were taken from YouTube at 720p quality and include both dancing performances and practicing. We used only the audio and pose detection data (body joint positions) from this dataset, which was extracted by applying a pose detector (Wei et al., 2016) after detecting bounding boxes in a frame with a real-time person detector (Redmon et al., 2016). After filtering out all instances which did not have all pose detection data for 10s, the final dataset size is 592 pairs.

3.2. Features

The audio features consist of logarithmically scaled Mel-spectrograms extracted from 16,000Hz audio signals. Framing is done by segmenting chunks of 50ms of audio every 25ms. Spectra are computed via Fast Fourier Transform (FFT) with a buffer size of 1024 samples. The number of Mel bins is set to 128, which results in a final matrix of 399 frames by 128 Mel-frequency bins per 10s audio recording. We segment each recording into 1s chunks (50% overlap) to be fed to the MVNN (detailed in Section 3.4), which means that each of the 592 objects contains 19 segments (each containing 39 frames), yielding a dataset of a total of 11,248 samples.

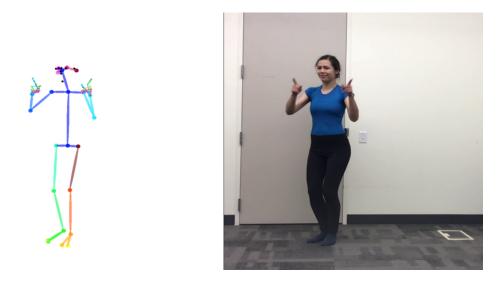


Figure 1: Pose detection illustration taken from (Chan et al., 2018). Skeleton points represent joints.

The pose detection features consist of body joint positions in frame space, i.e., pixel coordinates ranging from 0 (top left corner) to 1280 and 720 for width and height, respectively. The positions for the following key points are extracted: head, neck, shoulder, elbow, wrist, hip, knee, and ankle. There are 2 keypoints, left and right, for each of these except for head and neck, yielding a total of 28 features (14 keypoints with 2 coordinates, x and y, each). Figure 1 illustrates the keypoints. These features are extracted at 30 fps for the whole 10s video duration $(t \in \{t_0...t_{299}\})$, normalized after extraction according to Section 3.3, and then derived features are computed from the normalized data. The position and movement of body joints are used together for expression in dance. Therefore, we compute features that reflect the relative positions of body joints in relation to each other. This translates into computing the euclidean distance between each combination of two joints, yielding 91 derived features and a total of 119 movement features. As for audio, we segment this sequence into 1s segments (50% overlap), each containing 30 frames.

3.3. Preprocessing

We are interested in modeling movement as bodily expression. Therefore, we should focus on the temporal dynamics of joint positions relative to each other in a way that is as viewpoint- and subject-invariant as possible. However, the positions of subjects in frame space varies according to their distance to the camera. Furthermore, limb proportions are also different across subjects. Therefore, we normalize the joint position data in a similar way to Chan et al. (2018), whose purpose was to transform a pose from a source frame space to a target frame space. We select an arbitrary target frame and project every source frame to this space. We start by taking the maximum ankle y coordinate ankl^{clo} (Equation 1) and the maximum ankle y coordinate which is smaller than (spatially above) the median ankle y coordinate ankl^{med} (Equation 2) and about the same distance to it as the distance between it and ankl^{clo} (ankl^{far} in Equation 3). These two keypoints represent the closest and furthest ankle coordinates to the camera, respectively. Formally:

$$ankl = \{ankl_y_t^L\} \cup \{ankl_y_t^R\}$$
$$ankl^{clo} = \max_t(\{y_t : y_t \in ankl\})$$
(1)
$$ankl^{med} = median_t(\{y_t : y_t \in ankl\})$$
(2)

$$ankl^{far} = \max_{t}(\{y_{t} : y_{t} \in ankl \land y_{t} < ankl^{med} \land |y_{t} - ankl^{med}| - \alpha |ankl^{clo} - ankl^{med}| < \epsilon\})$$
(3)

where $ankl_y_t^L$ and $ankl_y_t^R$ are the y coordinates of the left and right ankles at timestep t, respectively. Following (Chan et al., 2018), we set α to 1, and ϵ to 0.7. Then, we computed a scale s (Equation 4) to be applied to the y-axis according to an interpolation between the ratios of the maximum heights between the source and target frames, $heig_{src}^{far}$ and $heig_{tgt}^{far}$, respectively. For each dance instance, frame heights are first clustered according to the distance between corresponding ankle y coordinate and $ankl^{clo}$ and $ankl^{far}$ and then the maximum height values for each cluster are taken (Equations 5 and 6). Formally:

$$s = \frac{\text{heig}_{tgt}^{far}}{\text{heig}_{src}^{far}} + \frac{\text{ankl}_{src}^{avg} - \text{ankl}_{src}^{far}}{\text{ankl}_{src}^{clo} - \text{ankl}_{src}^{far}} \left(\frac{\text{heig}_{tgt}^{clo}}{\text{heig}_{src}^{clo}} - \frac{\text{heig}_{tgt}^{far}}{\text{heig}_{src}^{far}}\right)$$
(4)

$$\begin{split} \text{heig}^{\text{clo}} &= \max_{t}(\{|\text{head}_{-}y_{t}-\text{ankl}_{t}^{\text{LR}}|:|\text{ankl}_{t}^{\text{LR}}-\text{ankl}^{\text{clo}}| < |\text{ankl}_{t}^{\text{LR}}-\text{ankl}^{\text{far}}|\}) \\ \text{heig}^{\text{far}} &= \max_{t}(\{|\text{head}_{-}y_{t}-\text{ankl}_{t}^{\text{LR}}|:|\text{ankl}_{t}^{\text{LR}}-\text{ankl}^{\text{clo}}| > |\text{ankl}_{t}^{\text{LR}}-\text{ankl}^{\text{far}}|\}) \\ & (5) \\ \text{heig}^{\text{far}} &= \max_{t}(\{|\text{head}_{-}y_{t}-\text{ankl}_{t}^{\text{LR}}|:|\text{ankl}_{t}^{\text{LR}}-\text{ankl}^{\text{clo}}| > |\text{ankl}_{t}^{\text{LR}}-\text{ankl}^{\text{far}}|\}) \\ & (6) \\ \text{ankl}_{t}^{\text{LR}} &= \frac{\text{ankl}_{-}y_{t}^{\text{L}}+\text{ankl}_{-}y_{t}^{\text{R}}}{2} \end{split}$$

 $ankl^{avg} = average_t(\{y_t : y_t \in ankl\})$

where $head_y_t$ is the y coordinate of the head at timestep t. After scaling, we also apply a 2D translation so that the position of the ankles of the subject is centered at 0. We do this by subtracting the median coordinates (x and y) of the mean of the (left and right) ankles, i.e., the median of $ankl_t^{LR}$.

3.4. Multi-view neural network architecture

The MVNN model used in this work is composed by two branches, each modeling its own view. Even though the final embeddings define a shared and correlated space, according to the loss function, the branches can be arbitrarily different from each other. The loss function is Deep Canonical Correlation Analysis (DCCA) (Andrew et al., 2013), a non-linear extension of Canonical Correlation Analysis (CCA) (Hotelling, 1936), which has also been successfully applied to music by Kelkar et al. (2018) and Yu et al. (2019). CCA linearly projects two distinct view spaces into a shared correlated space and was suggested to be a general case of parametric tests of statistical significance (Knapp, 1978). Formally, DCCA solves:

$$\left(w_{\mathbf{x}}^{*}, w_{\mathbf{y}}^{*}, \varphi_{\mathbf{x}}^{*}, \varphi_{\mathbf{y}}^{*}\right) = \operatorname*{argmax}_{\left(w_{\mathbf{x}}, w_{\mathbf{y}}, \varphi_{\mathbf{x}}, \varphi_{\mathbf{y}}\right)} \operatorname{corr}\left(w_{\mathbf{x}}^{\mathbf{T}} \varphi_{\mathbf{x}}\left(\mathbf{x}\right), w_{\mathbf{y}}^{\mathbf{T}} \varphi_{\mathbf{y}}\left(\mathbf{y}\right)\right)$$
(7)

where $\mathbf{x} \in \mathbb{R}^{\mathbf{m}}$ and $\mathbf{y} \in \mathbb{R}^{\mathbf{n}}$ are the zero-mean observations for each view. $\varphi_{\mathbf{x}}$ and $\varphi_{\mathbf{y}}$ are non-linear mappings for each view, and $w_{\mathbf{x}}$ and $w_{\mathbf{y}}$ are the canonical weights for each view. We use backpropagation and minimize:

$$-\sqrt{\mathrm{tr}\left(\left(C_{XX}^{-1/2}C_{XY}C_{YY}^{-1/2}\right)^{\mathrm{T}}\left(C_{XX}^{-1/2}C_{XY}C_{YY}^{-1/2}\right)\right)}$$
(8)

$$C_{XX}^{-1/2} = Q_{XX} \Lambda_{XX}^{-1/2} Q_{XX}^{\mathbf{T}}$$
(9)

where X and Y are the non-linear projections for each view, i.e., $\varphi_{\mathbf{x}}(\mathbf{x})$ and $\varphi_{\mathbf{y}}(\mathbf{y})$, respectively. C_{XX} and C_{YY} are the regularized, zero-centered covariances while C_{XY} is the zero-centered cross-covariance. Q_{XX} are the eigenvectors of C_{XX} and Λ_{XX} are the eigenvalues of C_{XX} . $C_{YY}^{-1/2}$ can be computed analogously. We finish training by computing a forward pass with the training data and fitting a linear CCA model on those non-linear mappings. The canonical components of these deep non-linear mappings implement our semantic embeddings space to be evaluated in a cross-modal retrieval task. Functions $\varphi_{\mathbf{x}}$ and $\varphi_{\mathbf{y}}$, i.e., the audio and movement projections are implemented as branches of typical neural networks, described in Tables 1 and 2. We use *tanh* activation functions after each convolution layer. Note that other loss functions, such as ones based on pairwise distances (Hermann & Blunsom, 2014; He et al., 2017), can theoretically also be used for the same task. The neural network models were all implemented using TensorFlow (Abadi et al., 2015).

Table 1: Audio Neural Network Branch								
layer type	dimensions					# params		
input	39	X	128	×	1	0		
batch norm	39	\times	128	\times	1	4		
2D conv	39	×	128	Х	8	200		
2D avg pool	13	×	16	X	8	0		
batch norm	13	×	16	×	8	32		
2D conv	13	×	16	×	16	2064		
2D avg pool	3	×	4	×	16	0		
batch norm	3	Х	4	\times	16	64		
2D conv	3	×	4	×	32	6176		
2D avg pool	1	\times	1	×	32	0		
batch norm	1	Х	1	×	32	128		
2D conv	1	×	1	\times	128	4224		
Total params						12892		

Table 1. Audio Neural Network Branch

Table 2: Movement Neural Network Branch						
layer type	dimensions			# of params		
input	30	×	119	0		
batch norm	30	×	119	476		
gru	1	×	32	14688		
Total params				15164		

3.5. Cross-modal retrieval evaluation

In this work, cross-modal retrieval consists of retrieving a sorted list of videos given an audio query and vice-versa. We perform cross-modal retrieval

on full objects even though the MVNN is modeling semiotic correlation between segments. In order to do this, we compute object representations as the average of the CCA projections of its segments (for both modalities) and compute the cosine similarity between these cross-modal embeddings. We evaluate the ability of the model to capture semantics and generalize semiotic correlations between both modalities by assessing if relevant crossmodal documents for a query are ranked on top of the retrieved documents list. We define relevant documents in two ways: instance- and class-level. Instance-level evaluation considers the ground truth pairing of cross-modal objects as criterion for relevance, (i.e., the only relevant audio document for a dance video is the one that corresponds to the song that played in that video). Class-level evaluation considers that any cross-modal object sharing some semantic label is relevant (e.g., relevant audio documents for a dance video of a particular dance style are the ones that correspond to songs that played in videos of the same dance style). We perform experiments in a 4-fold cross-validation setup, where each fold partitioning is such that the distribution of classes is similar for each fold. We also run the experiments 10 runs for each fold and report the average performance across runs.

We compute pair and rank accuracies for instance-level evaluation (similar to Pereira et al. (2018)). Pair accuracy evaluates ranking performance in the following way: for each query from modality X, we consider every possible pairing of the relevant object (corresponding cross-modal pair) and non-relevant objects from modality Y. We compute the similarities between the query and each of the two cross-modal objects, as well as the similarities between both cross-modal objects and the corresponding non-relevant object form modality X. If the corresponding cross-modal objects are more similar than the alternative, the retrieval trial is successful. We report the average values over queries and non-relevant objects. We also compute a statistical significance test in order to show that the model indeed captures semantics underlying the artifacts. We can think of each trial as a binomial outcome, aggregating two binomial outcomes, where the probability of success for a random model is $0.5 \times 0.5 = 0.25$. Therefore, we can perform a binomial test and compute its p-value. Even though there are 144×143 trials, we consider a more conservative value for the trials 144 (the number of independent queries). If the p-value is lower than 0.05, then we can reject the null hypothesis that the results of our model are due to chance. Rank accuracy is the (linearly) normalized rank of the relevant document in the retrieval list: ra = 1 - (r-1)/(L-1), where r is the rank of the relevant cross-modal

object in the list with L elements. This is similar to the pair accuracy evaluation, except that we only consider the query from modality X and the objects from modality Y, i.e., each trial consists of one binomial outcome, where the probability of success for a random model is 0.5. We also consider a conservative binomial test number of trials of 144 for this metric.

Even though the proposed model and loss function do not explicitly optimize class separation, we expect it to still learn embeddings which capture some aspects of the dance genres in the dataset. This is because different instances of the same class are expected to share semantic structures. Therefore, we perform class-level evaluation, in order to further validate that our model captures semantics underlying both modalities. We compute and report MAP scores for each class, separately, and perform a permutation test on these scores against random model performance (whose MAP scores are computed according to Bestgen (2015)), so that we can show these results are statistically significant and not due to chance. Formally:

$$MAP_{C} = \frac{1}{|Q_{C}|} \sum_{q \in Q_{C}} AP_{C}(q)$$
(10)

$$AP_C(q) = \frac{\sum_{j=1}^{|R|} \operatorname{pr}(j) \operatorname{rel}_C(r_j)}{|R_C|}$$
(11)

where C is the class, Q_C is the set of queries belonging to class C, $AP_C(q)$ is the Average Precision (AP) for query q, R is the list of retrieved objects, R_C is the set of retrieved objects belonging to class C, pr (j) is the precision at cutoff j of the retrieved objects list, and rel_C(r) evaluates whether retrieved object r is relevant or not, i.e., whether it belongs to class C or not. Note that the retrieved objects list always contains the whole (train or test) set of data from modality Y and that its size is equal to the total number of (train or test) evaluated queries from modality X. MAP measures the quality of the sorting of retrieved items lists for a particular definition of relevance (dance style in this work).

4. Results

Instance-level evaluation results are reported in Tables 3 and 4 for pair and rank accuracies, respectively, for each fold. Values shown in the X / Y format correspond to results when using audio / video queries, respectively. The model was able to achieve 57% and 75% for pair and rank accuracies, respectively, which are statistically significantly better (p-values < 0.01) than the random baseline performances of 25% and 50%, respectively.

Table 3: Instance-level Pair Accuracy

Fold 0	Fold 1	Fold 2	Fold 3	Average	Baseline
$0.57 \ / \ 0.57$	$0.57 \ / \ 0.56$	0.60 / 0.59	0.55 / 0.56	0.57 / 0.57	0.25

Table 4: Instance-level Rank Accuracy

Fold 0	Fold 1	Fold 2	Fold 3	Average	Baseline
0.75 / 0.75	$0.75 \ / \ 0.75$	0.77 / 0.76	$0.74 \ / \ 0.74$	0.75 / 0.75	0.50

Class-level evaluation results (MAP scores) are reported in Table 5 for each class and fold. The model achieved 26%, which is statistically significantly better (p-value < 0.01) than the random baseline performance of 13%.

5. Discussion

Our proposed model successfully captured semantics for music and dance, as evidenced by the quantitative evaluation results, which are validated by statistical significance testing, for both instance- and class-level scenarios. Instance-level evaluation confirms that our proposed model is able to generalize the cross-modal features which connect both modalities. This means the model effectively learned how people can move according to the sound of music, as well as how music can sound according to the movement of human bodies. Class-level evaluation further strengthens this conclusion by showing the same effect from a style-based perspective, i.e., the model learned how people can move according to the music style of a song, as well as how music can sound according to the dance style of the movement of human bodies. This result is particularly interesting because the design of both the model and experiments does not explicitly address style, that is, there is no stylebased supervision. Since semantic labels are inferenced by humans based on

Table 5: Class-level MAP

Style	Fold 0	Fold 1	Fold 2	Fold 3	Average	Baseline
Ballet	0.43 / 0.40	0.33 / 0.31	0.51 / 0.41	0.37 / 0.32	0.41 / 0.36	0.10
Breakdance	0.18 / 0.17	0.18 / 0.14	0.18 / 0.14	0.23 / 0.22	0.19 / 0.17	0.09
Flamenco	0.20 / 0.18	0.16 / 0.19	0.15 / 0.16	0.16 / 0.17	0.17 / 0.17	0.12
Foxtrot	0.22 / 0.24	0.23 / 0.24	0.21 / 0.21	0.16 / 0.18	0.20 / 0.22	0.12
Latin	0.23 / 0.23	0.19 / 0.20	0.21 / 0.22	0.20 / 0.19	0.21 / 0.21	0.14
Quickstep	0.21 / 0.20	0.14 / 0.12	0.19 / 0.19	$0.21 \ / \ 0.16$	0.19 / 0.17	0.09
Square	0.28 / 0.26	$0.34 \ / \ 0.29$	0.30 / 0.26	0.30 / 0.29	0.30 / 0.27	0.16
Swing	0.22 / 0.21	$0.22 \ / \ 0.22$	0.22 / 0.23	$0.24 \ / \ 0.26$	0.23 / 0.23	0.15
Tango	0.28 / 0.29	$0.39 \ / \ 0.37$	$0.34 \ / \ 0.38$	$0.31 \ / \ 0.33$	0.33 / 0.34	0.17
Waltz	0.52 / 0.51	$0.35 \ / \ 0.35$	0.38 / 0.31	0.48 / 0.41	0.43 / 0.40	0.15
Average	0.28 / 0.27	$0.25 \ / \ 0.24$	$0.27 \ / \ 0.25$	$0.27 \ / \ 0.25$	0.26 / 0.25	0.13
Overall	0.28 / 0.27	$0.27 \ / \ 0.26$	$0.27 \ / \ 0.26$	$0.28 \ / \ 0.26$	0.28 / 0.26	0.14

semiotic aspects, this implies that some of the latent semiotic aspects learned by our model are also relevant for these semantic labels, i.e., these aspects are semantically rich. Therefore, modeling semiotic correlations, between audio and dance, effectively uncovers semantic aspects.

The results show a link between musical meaning and kinetic meaning, providing further evidence for embodied cognition semantics in music. This is because embodied semantics ultimately defends that meaning in music is grounded in motor and somatosensory concepts, i.e., movement, physical forces, and physical tension. By observing that dance, a body expression proxy for how those concepts correlate to the musical experience, is semiotically correlated to music artifacts, we show that music semantics is kinetically and biologically grounded. Furthermore, our quantitative results also demonstrate an effective technique for cross-modal retrieval between music audio and dance video, providing the basis for an automatic music video creation tool. This basis consists of a model that can recommend the song that best fits a particular dance video and the dance video that best fits a particular song. The class-level evaluation also validates the whole ranking of results, which means that the model can actually recommend several songs or videos that best fit the dual modality.

6. Conclusions and future work

We proposed a computational approach to model music embodied semantics via dance proxies, capable of recommending music audio for dance video and vice-versa. Quantitative evaluation shows this model to be effective for this cross-modal retrieval task and further validates claims about music semantics being defined by embodied cognition. Future work includes correlating audio with 3D motion capture data instead of dance videos in order to verify whether important spatial information is lost in 2D representations, incorporating Laban movement analysis features and other audio features in order to have fine-grained control over which aspects of both music and movement are examined, test the learned semantic spaces in transfer learning settings, and explore the use of generative models (such as Generative Adversarial Networks (GANs)) to generate and visualize human skeleton dance videos for a given audio input.

References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., & Zheng, X. (2015). TensorFlow: Large-scale Machine Learning on Heterogeneous Systems, .
- Andrew, G., Arora, R., Bilmes, J., & Livescu, K. (2013). Deep Canonical Correlation Analysis. In Proceedings of the 30th International Conference on Machine Learning (pp. 1247–1255).
- Azcárate, A. L.-V. (2011). Intertextuality and Intermediality as Crosscultural Communication Tools: A Critical Inquiry. *Cultura. International Journal of Philosophy of Culture and Axiology*, 8, 7–22.
- Baily, J. (1985). Music Structure and Human Movement. In P. Howell, I. Cross, & R. West (Eds.), *Musical Structure and Cognition*. London: Academic Press.

- Bennett, K. (2008). The Language of Dance. Textos & Pretextos, (pp. 56–67).
- Bestgen, Y. (2015). Exact Expected Average Precision of the Random Baseline for System Evaluation. The Prague Bulletin of Mathematical Linguistics, (pp. 131–138).
- Bidelman, G. M., & Krishnan, A. (2009). Neural Correlates of Consonance, Dissonance, and the Hierarchy of Musical Pitch in the Human Brainstem. *Journal of Neuroscience*, 29, 13165–13171.
- Bidelman, G. M., & Krishnan, A. (2011). Brainstem Correlates of Behavioral and Compositional Preferences of Musical Harmony. *Neuroreport*, 22, 212–216.
- Blanariu, N. P. (2013). Towards a Framework of a Semiotics of Dance. CLCWeb: Comparative Literature and Culture, 15.
- Brown, S., & Jordania, J. (2011). Universals in the World's Musics. *Psychology of Music*, 41, 229–248.
- Castro, D., Hickson, S., Sangkloy, P., Mittal, B., Dai, S., Hays, J., & Essa, I. A. (2018). Let's Dance: Learning from Online Dance Videos. CoRR, abs/1801.07388.
- Cespedes-Guevara, J., & Eerola, T. (2018). Music Communicates Affects, Not Basic Emotions - A Constructionist Account of Attribution of Emotional Meanings to Music. Frontiers in Psychology, 9.
- Chan, C., Ginosar, S., Zhou, T., & Efros, A. A. (2018). Everybody Dance Now. CoRR, abs/1808.07371.
- Cook, N. (2000). Analysing Musical Multimedia. Oxford University Press.
- Davies, S. (1994). Musical Meaning and Expression. Cornell University Press.
- Dehaene, S., & Cohen, L. (2007). Cultural Recycling of Cortical Maps. Neuron, 56, 384–398.
- Desai, R. H., Binder, J. R., Conant, L. L., Mano, Q. R., & Seidenberg, M. S. (2011). The Neural Career of Sensory-motor Metaphors. *Journal of Cognitive Neuroscience*, 23, 2376–2386.

- Eitan, Z., & Granot, R. Y. (2006). How Music Moves: Music Parameters and Listeners' Images of Motion. *Music Perception*, 23, 221–248.
- Eitan, Z., & Rothschild, I. (2011). How Music Touches: Musical Parameters and Listeners' Audio-tactile Metaphorical Mappings. *Music Perception*, 39, 449–467.
- Frego, R. J. D. (1999). Effects of Aural and Visual Conditions on Response to Perceived Artistic Tension in Music and Dance. *Journal of Research in Music Education*, 47, 31–43.
- Fujioka, T., Trainor, L. J., Large, E. W., & Ross, B. (2012). Internalized Timing of Isochronous Sounds is Represented in Neuromagnetic Beta Oscillations. *Journal of Neuroscience*, 32, 1791–1802.
- He, W., Wang, W., & Livescu, K. (2017). Multi-view Recurrent Neural Acoustic Word Embeddings. In Proceedings of the 5th International Conference on Learning Representations.
- Hermann, K. M., & Blunsom, P. (2014). Multilingual Distributed Representations Without Word Alignment. In Proceedings of the 2nd International Conference on Learning Representations.
- Hotelling, H. (1936). Relations Between Two Sets of Variates. Biometrika, 28, 321–377.
- Janata, P., Tomic, S. T., & Haberman, J. M. (2012). Sensorimotor Coupling in Music and the Psychology of the Groove. *Journal of Experimental Psychology: General*, 141, 54–75.
- Juslin, P. N. (2013). What does Music Express? Basic Emotions and Beyond. Frontiers in Psychology, 4.
- Kelkar, T., Roy, U., & Jensenius, A. R. (2018). Evaluating a Collection of Sound-tracing Data of Melodic Phrases. In *Proceedings of the 19th International Society for Music Information Retrieval* (pp. 74–81).
- Kiefer, M., & Pulvermüller, F. (2012). Conceptual Representations in Mind and Brain: Theoretical Developments, Current Evidence and Future Directions. *Cortex*, 48, 805–825.

- Kivy, P. (1980). The Corded Shell: Reflections on Musical Expression. Princeton University Press.
- Knapp, T. R. (1978). Canonical Correlation Analysis: A General Parametric Significance-testing System. Psychological Bulletin, 85, 410–416.
- Koelsch, S., Vuust, P., & Friston, K. (2019). Predictive Processes and the Peculiar Case of Music. Trends in Cognitive Sciences, 23, 63–77.
- Korsakova-Kreyn, M. (2018). Two-level Model of Embodied Cognition in Music. Psychomusicology: Music, Mind, and Brain, 28, 240–259.
- Krumhansl, C. L., & Schenck, D. L. (1997). Can Dance Reflect the Structural and Expressive Qualities of Music? A Perceptual Experiment on Balanchine's Choreography of Mozart's Divertimento No. 15. *Musicae Scientiae*, 1, 63–85.
- Kurth, E. (1991). Ernst Kurth: Selected Writings. Cambridge University Press.
- von Laban, R., & Ullmann, L. (1960). *The Mastery of Movement*. London: MacDonald & Evans.
- Lakoff, G. (2012). Explaining Embodied Cognition Results. Topics in Cognitive Science, 4, 773–785.
- Lakoff, G. (2014). Mapping the Brain's Metaphor Circuitry: Metaphorical Thought in Everyday Reason. *Frontiers in Human Neuroscience*, 8.
- Larson, S. (2004). Musical Forces and Melodic Expectations: Comparing Computer Models and Experimental Results. *Music Perception*, 21, 457– 498.
- Leman, M. (2007). Embodied Music Cognition and Mediation Technology. MIT Press.
- Leman, M. (2010). An Embodied Approach to Music Semantics. Musicae Scientiae, 14, 43–67.
- Leman, M., & Maes, P.-J. (2014). The Role of Embodiment in the Perception of Music. *Empirical Musicology Review*, 9, 236–246.

- Lemke, J. L. (1992). Intertextuality and Educational Research. Linguistics and Education, 4, 257–267.
- Matyja, J. R. (2016). Embodied Music Cognition: Trouble Ahead, Trouble Behind. *Frontiers in Psychology*, 7.
- Nakamura, S., Oohashi, N. S. T., Nishina, E., Fuwamoto, Y., & Yonekura, Y. (1999). Analysis of Music-brain Interaction with Simultaneous Measurement of Regional Cerebral Blood Flow and Electroencephalogram Beta Rhythm in Human Subjects. *Neuroscience Letters*, 275, 222–226.
- Nozaradan, S., Peretz, I., Missal, M., & Mouraux, A. (2011). Tagging the Neuronal Entrainment to Beat and Meter. *Journal of Neuroscience*, 31, 10234–10240.
- Peirce, C. S. (1991). On Signs: Writings on Semiotic. University of North Carolina Press.
- Penhune, V. B., Zatorre, R. J., & Evans, A. C. (1998). Cerebellar Contributions to Motor Timing: A PET Study of Auditory and Visual Rhythm Reproduction. *Journal of Cognitive Neuroscience*, 10, 752–765.
- Pereira, F., Lou, B., Pritchett, B., Ritter, S., Gershman, S. J., Kanwisher, N., Botvinick, M., & Fedorenko, E. (2018). Toward a Universal Decoder of Linguistic Meaning from Brain Activation. *Nature Communications*, 9.
- Phillips-Silver, J., & Trainor, L. J. (2007). Hearing What the Body Feels: Auditory Encoding of Rhythmic Movement. *Cognition*, 105, 533–546.
- Platel, H., Price, C., Baron, J.-C., Wise, R., Lambert, J., Frackowiak, R. S. J., Lechevalier, B., & Eustache, F. (1997). The Structural Components of Music Perception: A Functional Anatomical Study. *Brain*, 120, 229– 243.
- Ralph, M. A. L., Jefferies, E., Patterson, K., & Rogers, T. T. (2017). The Neural and Computational Bases of Semantic Cognition. *Nature Reviews Neuroscience*, 18, 42–55.
- Redmon, J., Divvala, S. K., Girshick, R. B., & Farhadi, A. (2016). You Only Look Once: Unified, Real-time Object Detection. In *Proceedings of the* 29th IEEE Conference on Computer Vision and Pattern Recognition (pp. 779–788).

- Roffler, S. K., & Butler, R. A. (1968). Localization of Tonal Stimuli in the Vertical Plane. Journal of the Acoustical Society of America, 43, 1260– 1265.
- Schlenker, P. (2017). Outline of Music Semantics. Music Perception, 35, 3–37.
- Scruton, R. (1997). The Aesthetics of Music. Oxford University Press.
- Sievers, B., Polansky, L., Casey, M., & Wheatley, T. (2013). Music and Movement Share a Dynamic Structure that Supports Universal Expressions of Emotion. Proceedings of the National Academy of Sciences of the United States of America, 110, 70–75.
- Spence, C., & Driver, J. (1997). Audiovisual Links in Exogenous Covert Spatial Orienting. Perception & Psychophysics, 59, 1–22.
- Stein, B. E., Wallace, M. T., & Meredith, M. A. (1995). Neural Mechanisms Mediating Attention and Orientation to Multisensory Cues. In M. S. Gazzaniga (Ed.), *The Cognitive Neurosciences* (pp. 683–702). MIT Press.
- Styns, F., van Noorden, L., Moelants, D., & Leman, M. (2007). Walking on Music. Human Movement Science, 26, 769–785.
- Virtala, P., Huotilainen, M., Partanen, E., Fellman, V., & Tervaniemi, M. (2013). Newborn Infants' Auditory System is Sensitive to Western Music Chord Categories. *Frontiers in Psychology*, 4.
- Wagner, S., Winner, E., Cicchetti, D., & Gardner, H. (1981). "Metaphorical" Mapping in Human Infants. *Child Development*, 52, 728–731.
- Wei, S.-E., Ramakrishna, V., Kanade, T., & Sheikh, Y. (2016). Convolutional Pose Machines. In Proceedings of the 29th IEEE Conference on Computer Vision and Pattern Recognition (pp. 4724–4732).
- Widmann, A., Kujala, T., Tervaniemi, M., Kujala, A., & Schröger, E. (2004). From Symbols to Sounds: Visual Symbolic Information Activates Sound Representations. *Psychophysiology*, 41, 709–715.
- Yu, Y., Tang, S., Raposo, F., & Chen, L. (2019). Deep Cross-modal Correlation Learning for Audio and Lyrics in Music Retrieval. ACM Transactions on Multimedia Computing, Communications, and Applications, 15.

Zatorre, R. J., Evans, A. C., & Meyer, E. (1994). Neural Mechanisms Underlying Melodic Perception and Memory for Pitch. *Journal of Neuroscience*, 14, 1908–1919.