

# Extracting Coarse Body Movements from Video in Music Performance: A Comparison of Automated Computer Vision Techniques with Motion Capture Data

Kelly Jakubowski<sup>1\*</sup>, Tuomas Eerola<sup>1</sup>, Paolo Alborno<sup>2</sup>, Gualtiero Volpe<sup>2</sup>, Antonio Camurri<sup>2</sup>, Martin Clayton<sup>1</sup>

<sup>1</sup>Music, Durham University, United Kingdom, <sup>2</sup>DIBRIS (Department of Informatics, Bioengineering, Robotics, and Systems Engineering), University of Genova, Italy

*Submitted to Journal:*  
Frontiers in Digital Humanities

*Specialty Section:*  
Digital Musicology

*ISSN:*  
2297-2668

*Article type:*  
Original Research Article

*Received on:*  
08 Jan 2017

*Accepted on:*  
21 Mar 2017

*Provisional PDF published on:*  
21 Mar 2017

*Frontiers website link:*  
[www.frontiersin.org](http://www.frontiersin.org)

*Citation:*  
Jakubowski K, Eerola T, Alborno P, Volpe G, Camurri A and Clayton M(2017) Extracting Coarse Body Movements from Video in Music Performance: A Comparison of Automated Computer Vision Techniques with Motion Capture Data. *Front. Digit. Humanit.* 4:9. doi:10.3389/fdigh.2017.00009

*Copyright statement:*  
© 2017 Jakubowski, Eerola, Alborno, Volpe, Camurri and Clayton. This is an open-access article distributed under the terms of the [Creative Commons Attribution License \(CC BY\)](https://creativecommons.org/licenses/by/4.0/). The use, distribution and reproduction in other forums is permitted, provided the original author(s) or licensor are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

Provisional

1           **Extracting Coarse Body Movements from Video in Music**  
2           **Performance: A Comparison of Automated Computer Vision**  
3           **Techniques with Motion Capture Data**

4  
5   **Kelly Jakubowski<sup>1\*</sup>, Tuomas Eerola<sup>1</sup>, Paolo Albornò<sup>2</sup>, Gualtiero Volpe<sup>2</sup>, Antonio**  
6   **Camurri<sup>2</sup>, Martin Clayton<sup>1</sup>**

7   <sup>1</sup>Department of Music, Durham University, Durham, UK

8   <sup>2</sup>Casa Paganini Research Centre, DIBRIS (Department of Informatics, Bioengineering,  
9   Robotics, and Systems Engineering), University of Genova, Italy

10   **\* Correspondence:**

11   Kelly Jakubowski

12   kelly.jakubowski@durham.ac.uk

13  
14   **Keywords: movement, motion tracking, music performance, musical ensemble**  
15   **coordination, computer vision, video analysis**

16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29

Provisional

30 **Abstract**

31 The measurement and tracking of body movement within musical performances can provide  
32 valuable sources of data for studying interpersonal interaction and coordination between  
33 musicians. The continued development of tools to extract such data from video recordings  
34 will offer new opportunities to research musical movement across a diverse range of settings,  
35 including field research and other ecological contexts in which the implementation of  
36 complex motion capture systems is not feasible or affordable. Such work might also make  
37 use of the multitude of video recordings of musical performances that are already available to  
38 researchers. The present study made use of such existing data, specifically, three video  
39 datasets of ensemble performances from different genres, settings, and instrumentation (a pop  
40 piano duo, three jazz duos, and a string quartet). Three different computer vision techniques  
41 were applied to these video datasets—frame differencing, optical flow, and kernelized  
42 correlation filters (KCF)—with the aim of quantifying and tracking movements of the  
43 individual performers. All three computer vision techniques exhibited high correlations with  
44 motion capture data collected from the same musical performances, with median correlation  
45 (Pearson's  $r$ ) values of .75 to .94. The techniques that track movement in two dimensions  
46 (optical flow and KCF) provided more accurate measures of movement than a technique that  
47 provides a single estimate of overall movement change by frame for each performer (frame  
48 differencing). Measurements of performer's movements were also more accurate when the  
49 computer vision techniques were applied to more narrowly-defined regions of interest (head)  
50 than when the same techniques were applied to larger regions (entire upper body, above the  
51 chest or waist). Some differences in movement tracking accuracy emerged between the three  
52 video datasets, which may have been due to instrument-specific motions that resulted in  
53 occlusions of the body part of interest (e.g. a violinist's right hand occluding the head whilst  
54 tracking head movement). These results indicate that computer vision techniques can be  
55 effective in quantifying body movement from videos of musical performances, while also  
56 highlighting constraints that must be dealt with when applying such techniques in ensemble  
57 coordination research.

58

59

60

61

62

63

64

65

66

## 67 1. Introduction

68 The extraction and quantification of human movement data from musical performances offers  
69 a range of potential uses to researchers of musical interaction. Movement data from  
70 performers can be instrumental to research on interpersonal synchrony and entrainment  
71 between musicians, leader-follower relationships within an ensemble, and musical gestural  
72 analysis, to name just a few examples. Extraction of such data from video recordings can be  
73 particularly useful in situations where more complex or costly motion capture technologies  
74 are not feasible, such as field research and various other ecological performance contexts  
75 (e.g., gigs at nightclubs, rehearsals in music practice rooms, ritual ceremonies and religious  
76 events, etc.). One area that offers a variety of promising techniques for extracting features of  
77 human movement from video is the field of computer vision (Moeslund and Granum, 2001).  
78 The work of computer vision scientists is focussed around developing computational methods  
79 that perform similar tasks to the human visual system using digital images and videos,  
80 including object recognition, event detection, object tracking, and motion estimation (Forsyth  
81 and Ponce, 2002).

82 Researchers have recently begun to test the efficacy of computer vision techniques for  
83 capturing and indexing human body movements during social motor coordination tasks  
84 (Romero et al., 2016) and dance (Solberg and Jensenius, 2016). The work of Romero et al.  
85 (2016) suggests that computer vision methods, as applied to video recordings, can perform  
86 similar tracking of body movements to more expensive techniques, such as motion capture  
87 (MoCap) systems or Microsoft Kinect, under certain conditions. This is advantageous, as  
88 specialised MoCap technologies are not only costly, but can also be invasive in that markers  
89 need to be fixed to a person's body (or for some systems a specialised suit needs to be worn),  
90 time-consuming in terms of set-up and calibration procedures, and difficult to implement in  
91 ecological settings outside of specialised motion capture laboratories. Previous research has  
92 revealed that the conditions under which computer vision methods applied to video most  
93 closely approximate MoCap tracking in terms of body movement quantification include a  
94 fixed video camera angle (e.g., no zooming or panning), stable lighting within the recording  
95 setting, no other movements occurring in the background, and the separation of participants  
96 in space so as to avoid occlusions or the movements of one participant being included in the  
97 analysis space of another (Paxton and Dale, 2012; Romero et al., 2016). However, limitations  
98 of the use of computer vision methods for motion tracking include that these methods have  
99 previously proved more feasible for tracking large-scale, full-body movements than  
100 movements of individual body parts (Paxton and Dale, 2012; Romero et al., 2016) and only  
101 measure movements in two dimensions (cf. MoCap and sensors such as accelerometers,  
102 which measure movements in three dimensions). Additionally, computer vision techniques  
103 are generally applied to data sources with a lower temporal resolution than MoCap  
104 technologies; standard video recordings tend to be recorded at a frame rate of around 25  
105 frames per second (fps), whereas MoCap data is often recorded in the range of 100 to 200  
106 fps.

107 Music performance serves as another highly relevant case for testing the capabilities of  
108 computer vision techniques, as group music making employs a variety of movement cues to

109 facilitate the coordination of timing and expressivity between performers. This coordination  
110 of timing and expressivity is sometimes referred to as interpersonal entrainment (Clayton et  
111 al., 2005). When producing video recordings of musical performances it is also often possible  
112 to implement solutions to minimise some of the challenges to the application of computer  
113 vision techniques listed above. For instance, the lighting and camera angle may be able to be  
114 fixed to a standardised setting throughout a performance and the performers may be situated  
115 within the performance space such that they do not occlude one another (at least in small  
116 ensembles).

117 Coordination in musical ensembles is achieved through the use and integration of both  
118 auditory (instrumental and vocal sounds) and visual (body movement and eye contact) cues.  
119 The accuracy of temporal coordination in the auditory domain is typically in the order of tens  
120 of milliseconds in expert ensemble performance (e.g., Keller, 2014; Rasch, 1988; Shaffer,  
121 1984). The movements that produce these sounds, such as finger movements of a pianist or  
122 bowing movements of a violinist, often evolve at similarly short timescales. In addition to  
123 these instrumental, sound-producing movements that are required in performance, musicians  
124 also make use of a variety of communicative and sound-facilitating movements that can serve  
125 to coordinate timing and expressive intentions between performers (Jensenius et al., 2010).  
126 These ancillary movements (e.g., head nods, body sway) typically evolve over longer  
127 timescales than instrumental movements (e.g. in the order of seconds; Davidson, 2009;  
128 Wanderley et al., 2005). Importantly, systematic relationships have been observed between  
129 coordination at the level of ancillary body movements and musical sounds (Keller and Appel,  
130 2010; Ragert et al., 2013). Thus, the analysis of such movements can provide information  
131 about the overall level of interpersonal coordination within an ensemble performance. In  
132 contrast to acoustic features and instrumental movements, ancillary body movements tend to  
133 generalise across performers regardless of the instrument played and are also prevalent in  
134 vocal performance. Additionally, the fact that ancillary movements tend to take place across  
135 longer timescales than instrumental movements allows them to be tracked within video  
136 recordings despite its lower temporal resolution in comparison to MoCap. Therefore, it is of  
137 great interest to music researchers to measure and analyse ancillary movements from video  
138 recordings of musical performances.

139 There are a variety of areas within the field of music performance research that may benefit  
140 from the use of computer vision techniques to measure movement data with a view to  
141 quantifying interpersonal coordination. For instance, such techniques could be applied to  
142 study temporal relationships between performers within commercial video recordings of  
143 classical or popular music, or to quantify corporeal interactions between a music or dance  
144 therapist and his/her clients. Ethnomusicologists often make video recordings of musical  
145 performances in ecological settings in which access to sophisticated technologies such as  
146 motion capture is not feasible. Indeed, a large amount of archival material of video recordings  
147 of music performances from across the world already exists. For example, the JVC Video  
148 Anthology of World Music and Dance (JVC, Victor Company of Japan, 1990) comprises  
149 some 30 volumes of field recordings from across the world and the Ethnographic Video for  
150 Instruction & Analysis (EVIA) Digital Archive Project (<http://www.eviada.org/default.cfm>)

151 is a repository of ethnographic videos, including many music performances, which aims to  
152 preserve these materials for the long-term in a digital, online format. As such, if video-based  
153 analysis methods prove to be fruitful in providing new insights about musical interaction, a  
154 large amount of useful research could be done that makes use of such existing video archives  
155 (with the appropriate permissions and taking account of ethical considerations), which could  
156 thereby minimise the costs that are necessarily incurred when collecting new data. The  
157 present study served as a test case in this regard, as it also made use of existing data—in this  
158 case, three existing datasets in which both video and motion capture recordings had been  
159 collected (as reported in Glowinski et al., 2013, Moran et al., 2015, and one previously  
160 unpublished dataset). Our study was therefore able to test whether computer vision  
161 techniques could be used to quantify body movements from video recordings that had  
162 originally been obtained for other research purposes.

163 The computer vision field offers a diverse range of possible techniques for tracking moving  
164 elements and changes in image sequences that were considered for use in the present study.  
165 As the majority of materials in our datasets of musical performances presented a situation in  
166 which only the to-be-tracked targets (the performers) were moving, we first considered  
167 background subtraction techniques. These techniques aim to distinguish an object(s) (in this  
168 case, the performers) in the foreground from a static background and perform further  
169 processing (e.g., tracking or motion detection) on the foreground object. The background  
170 subtraction-based technique that we applied was frame differencing. Frame differencing is  
171 one of the oldest and most widely-used computer vision techniques, which measures the  
172 overall change in pixels within the foreground from one frame to the next (Wren et al., 1997;  
173 see also Jensenius et al., 2005, for an implementation for studying musical gestures). We then  
174 explored two techniques that provide more detailed information on the direction of motion of  
175 each performer. Specifically, we employed a technique based on the variation of the motion  
176 field, known as optical flow (Farnebäck, 2003), and a technique based on pattern similarity  
177 calculation, known as kernelized correlation filters (hereafter referred to as KCF; Henriques  
178 et al., 2015). Optical flow is a technique that has been widely applied within the computer  
179 vision literature (e.g. Fleet and Weiss, 2006; see also Latif et al., 2014, for an application in  
180 studying interpersonal coordination), whereas KCF is a comparatively recently developed  
181 technique. Both of these techniques were used to track the direction of movement of the  
182 performers by providing both horizontal and vertical position data of each performer within  
183 each frame.

184 To summarise, in the present project we applied three automated computer vision techniques  
185 (frame differencing, optical flow, and KCF) to a set of video recordings of musical  
186 performances comprising a variety of performers, performance settings, instrumentations, and  
187 musical styles. The aims were 1) to test the robustness of the computer vision techniques for  
188 capturing body movements across the different performance conditions and 2) to test how  
189 closely these techniques were able to capture the actual motion of performers, as indexed by  
190 motion capture data from the same performances. Finally, as previous studies comparing  
191 motion capture data to computer vision techniques have primarily examined full-body  
192 movements (e.g. Romero et al., 2016), we extended this area of research to include analysis

193 of video data within predefined regions of interest (i.e., head, upper body) to test whether the  
 194 video analysis techniques could also be effective in quantifying movements of specific parts  
 195 of the body. If it was found that computer vision techniques could be effectively applied to  
 196 measure movement in specific body parts such as the head, this would suggest that in some  
 197 cases it may be possible to differentiate sound-producing, instrumental movements from  
 198 sound-facilitating, ancillary movements of musical performers by isolating a part of the body  
 199 that does not play a role in both types of movement (e.g., a guitar or cello player does not  
 200 typically use head movements to produce sounds but rather for communicative purposes).

## 201 2. Methods

### 202 2.1 Materials

203 The project made use of three existing datasets (see Figure 1), in which both video recordings  
 204 and MoCap data of the same musical performances had been collected for other research  
 205 purposes.<sup>1</sup> The first dataset (previously unpublished and hereafter referred to as the “Piano  
 206 Duo”) comprised seven songs performed by singer-songwriters Konstantin Wecker and Jo  
 207 Barnikel. Wecker has been described as one of Germany’s most successful singer-  
 208 songwriters, with a career spanning 40 years at the time of the recording, and Barnikel is a  
 209 leading film and TV composer who had been accompanying Becker on recordings and  
 210 concert tours for over 15 years.

211 The second dataset consisted of three performances by jazz duos, a subset of the Improvising  
 212 Duos corpus described in Moran et al., 2015. In this subset (hereafter referred to as “Mixed  
 213 Instrument Duos”), two duos performed free jazz improvisations and one performed a jazz  
 214 standard (*Autumn Leaves* [J. Kosma, 1945]). Performers in these duos were recruited on the  
 215 basis of public performance experience of around 10 years in their respective styles. Data  
 216 from five of the six performers from this dataset were analysed in respect of performers’  
 217 permissions on data reuse.

218 The third dataset (“String Quartet”) comprised eight recordings by the Quartetto di Cremona  
 219 string quartet performing the first movement of Schubert’s *String Quartet No. 14* (“Death and  
 220 the Maiden”; Glowinski et al., 2013). Two of these recordings featured only the first violinist  
 221 performing his part alone. For the other six recordings, two of the four performers were  
 222 selected for whom the least occlusions were observed (i.e. another player was not moving in  
 223 front of him/her regularly). In total, the three datasets allowed for the analysis of 33 cases of  
 224 10 different performers playing six different instruments (see Table 1).

225

226 -INSERT FIGURE 1 ABOUT HERE-

227

---

<sup>1</sup> In all instances the primary focus of the original research was on the collection of MoCap data, thus the performance settings were optimised for MoCap data collection and video was collected as a secondary measure for reference purposes only.

228 For each of the three datasets, the recordings were made in the same room under similar  
229 performance conditions (e.g. all string quartet recordings were made with performers situated  
230 in a similar position on the same stage using the same video camera and MoCap system). The  
231 Piano Duo and Mixed Instrument Duos were both recorded at the Max Planck Institute in  
232 Leipzig, Germany, using a Vicon Nexus 1.6.1 optical motion capture system with ten  
233 cameras and a sampling rate of 200 Hz. A SONY HDR-HC9 camera was used to make the  
234 video recordings. The video files were recorded in AVI format at a frame rate of 25 fps and  
235 frame size of 720 x 576 pixels. The String Quartet was recorded at Casa Paganini Research  
236 Centre (University of Genova, Italy), using a Qualisys Oqus300 motion capture system with  
237 eleven cameras and a sampling rate of 100 Hz. A JVC GY-HD-251 camera was used to  
238 capture video of the performances. The video files were recorded in AVI format at a frame  
239 rate of 25 fps and frame size of 720 x 576 pixels.

240

241 -INSERT TABLE 1 ABOUT HERE-

242

## 243 2.2 Analysis

### 244 2.2.1 Motion capture data

245 All MoCap data were processed using the MoCap Toolbox (Burger and Toiviainen, 2013) in  
246 Matlab. Each dataset was first rotated in order to orient the MoCap data to the same  
247 perspective as the camera angle of the video recording. This was done manually by  
248 inspecting animations generated from the MoCap data in comparison to the video recording  
249 (see Figure 2). Once the optimal rotation was achieved, a subset of markers was selected  
250 from each performer, comprising one marker from the head and one from the torso or each  
251 shoulder (if a torso marker was not present, as was the case for the String Quartet). If  
252 multiple markers were present for a specific body part (e.g. four head markers), the marker  
253 for which the least amount of data points were missing was selected. Markers were also  
254 selected in consideration of the camera angle of the video. For instance, if only the back of  
255 the head of a performer was visible in the video, a marker from the back of the head was  
256 selected. The three-dimensional coordinates from each selected marker were saved for further  
257 analysis. The horizontal and vertical coordinates of the MoCap data are subsequently referred  
258 to as the x- and y-dimensions respectively, which were compared to the two-dimensional data  
259 that were derived from the video recordings by the computer vision techniques.

260

### 261 2.2.2 Video data

262 The computer vision techniques (frame differencing, optical flow, and KCF) were  
263 implemented in EyesWeb XMI 5.6.2.0 ([http://www.infomus.org/eyesweb\\_ita.php](http://www.infomus.org/eyesweb_ita.php)). The first  
264 step when applying each technique was to manually define relevant regions of interest (ROIs)  
265 on which to apply the technique to each video. A rectangular ROI was selected around each

266 performer whilst ensuring that only that individual performer was serving as the main source  
267 of motion in the ROI (see Figure 2). This was generally achieved to a high standard, although  
268 there were a few cases in which the hands or bows of another performer occasionally moved  
269 into the ROI in the Piano Duo and String Quartet. Two sets of ROIs were defined for each  
270 performer in each video—a larger ROI that comprised the upper body (from the mid-chest or  
271 the waist up to the top of the head, depending on how much of the performer could be seen in  
272 the video<sup>2</sup>) and a smaller ROI around the head only. Frame differencing and optical flow  
273 were both applied using the same sets of upper body and head ROIs for each video. A slightly  
274 different set of upper body and head ROIs were defined for KCF, due to the way this  
275 technique is implemented. In typical implementations of KCF, the entire ROI moves  
276 dynamically throughout the process of tracking the performer. Conversely, frame  
277 differencing and optical flow were applied on static ROIs that do not move during the  
278 analysis process. As such, larger ROIs were needed that could encompass the whole range of  
279 movement of a performer for frame differencing and optical flow, whereas KCF is more  
280 suited to smaller ROIs since the ROI shifts from frame to frame.

281 In frame differencing, the foreground, i.e. the moving element(s) of interest (in this case, the  
282 performers), is separated from the background and further processing is performed on the  
283 foreground. In the present study, frame differencing was implemented using the Pfinder  
284 algorithm of Wren et al. (1997). A version of this algorithm has previously been implemented  
285 in EyesWeb for studying interpersonal musical coordination in Indian duos (Alborno et al.,  
286 2015). The Pfinder algorithm uses adaptive background subtraction, in which the background  
287 model that is subtracted from the foreground is constantly updated throughout the analysis  
288 process. The speed at which the background model is updated is determined by the alpha  
289 constant, which was set in the present study to 0.4, following an optimisation process in  
290 which this parameter was manually adjusted to a range of values and tested on a subset of the  
291 present videos. The analysis that was performed on the foreground elements measures the  
292 overall Quantity of Motion (QoM) in each ROI for each frame, which is computed based on  
293 the number of pixels that change in the foreground from one frame to the next. This analysis  
294 produces one column of output values for each performer.

295 Optical flow is the distribution of apparent velocities of movement of brightness patterns in  
296 an image. In optical flow, characteristics such as edges or angles are identified within each  
297 section of the video frame. In the next frame, such characteristics are sought again. A speed is  
298 then associated to each pixel in the frame; the movement is determined by the ratio between  
299 the distance in pixels of the displacement of the characteristic in question and the time  
300 between one frame and another. The version of optical flow that was implemented in the  
301 present study is known as dense optical flow<sup>3</sup> and is based on the algorithm of Farnebäck

---

<sup>2</sup> In some cases the waist of a performer could not be seen, as it was behind their instrument (e.g. for some pianists).

<sup>3</sup> Traditional optical flow methods (e.g. as implemented by Lucas and Kanade (1981)) compute optical flow for a sparse feature set, i.e. using only specific parts of the image, such as detected corners. Dense optical flow, as implemented by Farnebäck (2003), performs optical flow computation on all pixels in the image for each frame. The use of dense optical flow can increase the accuracy of the optical flow results, with a tradeoff of slower computation speed.

302 (2003). This technique has previously been implemented in EyesWeb in work of Alborno et  
303 al. (2015) on Indian music duos, as well as to develop a “virtual binocular” installation in  
304 which users' movements are tracked and estimated by computation of optical flow on the face  
305 (Camurri et al., 2010). A similar optimisation procedure was followed to that used for frame  
306 differencing in which the “pyramid layers” parameter was adjusted to a range of values and  
307 tested on a subset of the present videos. This parameter allows for the tracking of points at  
308 multiple levels of resolution; increasing the number of pyramid layers allows for the  
309 measurement of larger displacements of points between frames but also increases the number  
310 of necessary computations. The optimal value that was selected for this parameter was 12.  
311 The resulting output that was provided by the optical flow analysis was two columns of data  
312 per performer, which represent movement of the barycentre of the ROI along the x-  
313 (horizontal) and y- (vertical) axes. The barycentre of the ROI is computed based on pixel  
314 intensities. The video image is converted to greyscale and the barycentre coordinates are  
315 calculated as a weighted mean of the pixel intensities within the ROI; this is done separately  
316 for the x- and y-dimensions.

317 KCF is a relatively recently developed tracking technique (Bolme et al., 2009), based on  
318 older correlation filter methods (Hester, 1980), that works using pattern similarity  
319 calculations on a frame-by-frame basis. KCF was implemented in EyesWeb<sup>4</sup> in the present  
320 study using the OpenCV C++ implementation<sup>5</sup> of the algorithm of Henriques et al. (2015).  
321 When the KCF algorithm is initialised, a visual tracker is placed at the centre pixel of the pre-  
322 defined ROI for the first frame of the video. In the second frame, similarity and classification  
323 computations are performed by searching for the set of pixels with the maximum correlation  
324 to the initial tracker position in terms of its multi-channel RGB colour attributes, and so on  
325 for each subsequent frame. In effect, this allows the technique to track the movement of the  
326 performers across the ROI. Similarly to optical flow, the output of the KCF analysis is two  
327 columns of data per performer, which represent movement of the barycentre of the ROI along  
328 the x- and y-axes. In this case, since the ROI moves dynamically with the performer, the  
329 barycentre that is used is the geometric barycentre at the intersection of the two diagonals of  
330 the rectangular ROI.

### 331 2.2.3 Motion capture and video comparison

332 As video data collection was not the primary focus of the original studies, the video and  
333 MoCap data were not synchronised with an external timecode. As such, these two data  
334 sources were aligned in the present study using automated cross-correlational methods. Each  
335 video analysis output from EyesWeb was cross-correlated with its corresponding MoCap  
336 target (e.g. the x-coordinate of the head from the optical flow analysis within the head ROI  
337 was cross-correlated with the x-coordinate of the MoCap head marker). This allowed us to  
338 determine the optimal lag time for each trial, which was defined as the lag at which the  
339 maximum correlation value between the video and MoCap data was reached. The median  
340 optimal lag time from all cross-correlational analyses from the same video (taking account of

---

<sup>4</sup> The KCF block has recently been released within the Image Processing Library of EyesWeb.

<sup>5</sup> [http://docs.opencv.org/trunk/d2/dff/classcv\\_1\\_1TrackerKCF.html](http://docs.opencv.org/trunk/d2/dff/classcv_1_1TrackerKCF.html)

341 analysis of all position data from both performers in each video) was taken as the optimal lag  
 342 time for that particular video. The median optimal lag time across all video and MoCap  
 343 pairings in the dataset was 0.05 seconds (range = -0.10 to 0.42 seconds). Before computing  
 344 any statistical comparisons between the video and MoCap data, the MoCap data were down-  
 345 sampled to match the lower sampling rate of the videos at 25 fps, and all video and MoCap  
 346 data outputs were de-trended and normalised. Figure 2 depicts the data preparation and  
 347 extraction process for video and MoCap for one example performance from the Mixed  
 348 Instrument Duos.

349

350 -INSERT FIGURE 2 ABOUT HERE-

351

### 352 3. Results

353 The main focus of the subsequent data analysis was to compare the efficacy of the three  
 354 computer vision techniques (frame differencing, optical flow, and KCF) for measuring body  
 355 movements of musical performers across the three different datasets ( Piano Duo, Mixed  
 356 Instrument Duos, and String Quartet)<sup>6</sup> and two sets of ROIs (upper body and head). For the  
 357 upper body ROI, we compared the outputs of the computer vision analyses to the coordinates  
 358 of the torso marker from the MoCap data (or the right shoulder marker, in the case of the  
 359 String Quartet<sup>7</sup>) for each trial. For the head ROI, we compared the computer vision data to  
 360 the coordinates of the MoCap head marker.

361 Since frame differencing provides a single, overall estimate of movement of each performer  
 362 (rather than two-dimensional tracking), the optical flow, KCF, and corresponding MoCap  
 363 data were converted from Cartesian (x and y) to polar (radial and angular) coordinates. We  
 364 then computed the absolute change of the radial coordinate on a frame-by-frame basis for  
 365 each trial; this absolute change measure was used in subsequent comparisons to the one-  
 366 dimensional frame differencing results. Both the resultant absolute change data and the QoM  
 367 data from frame differencing were kernel smoothed in R using the Nadaraya–Watson kernel  
 368 regression estimate with a bandwidth of 1.<sup>8</sup> The video and MoCap data for each trial were  
 369 then compared using correlations (Pearson’s  $r$ ); a summary of these comparisons is reported,  
 370 by dataset, in Table 2.<sup>9</sup> These descriptive statistics suggest that the two-dimensional tracking  
 371 methods (optical flow and KCF) tend to perform more accurately than the more coarse-

---

<sup>6</sup> Although the primary research question is focused on evaluating and comparing the three computer vision techniques within the two ROIs, “dataset” is also included as an independent variable in subsequent analyses to take account of the fact that the three datasets vary on a number of parameters, including setting, recording session, lighting, camera angle, and instrumentation.

<sup>7</sup> This analysis was also tested with the left shoulder marker and the average of the left and right shoulder markers, however these analyses revealed similar patterns of results and did not increase the overall correlations.

<sup>8</sup> This smoothing procedure was applied because both the video and MoCap data contained small random fluctuations, which were smoothed without tampering with the overall shape of the trajectories. Filtering had a minor positive effect on the overall results (mean increase in video/MoCap correlation values of 0.07).

<sup>9</sup> Median values (rather than means) are reported as descriptive statistics throughout this paper due to some non-normal data distributions and the relative robustness of the median to the presence of statistical outliers.

372 grained method (frame differencing) and that performance of all three computer vision  
 373 techniques is improved when concentrated on a smaller ROI (head, as compared to upper  
 374 body).

375

376 -INSERT TABLE 2 ABOUT HERE-

377

378 For the data using the upper body ROI, a 3x3 mixed ANOVA was conducted to test the  
 379 effects of computer vision technique (frame differencing, optical flow, KCF) and dataset  
 380 (Piano Duo, Mixed Instrument Duos, String Quartet) on accuracy of overall movement  
 381 measurement (as indexed by the correlation of each video analysis output with the MoCap  
 382 data; see Table 2). Prior to entering the correlation values as the dependent variable in the  
 383 ANOVA, these values were subjected to a Fisher z-transformation to normalise the  
 384 distribution. The ANOVA revealed significant main effects of computer vision technique  
 385 ( $F(2, 60) = 16.51, p < .001, \eta_p^2 = .355$ ) and dataset ( $F(2, 30) = 15.41, p < .001, \eta_p^2 = .507$ ), as  
 386 well as a significant technique by dataset interaction ( $F(4, 60) = 18.82, p < .001, \eta_p^2 = .557$ ).  
 387 Bonferroni-corrected, paired-samples t-tests indicated that optical flow provided a more  
 388 accurate measure of performers' movements than both frame differencing ( $t(32) = 3.67, p =$   
 389  $.003$ ) and KCF ( $t(32) = 3.38, p = .006$ ); no significant difference was found between the  
 390 frame differencing and KCF techniques. Tukey HSD tests revealed that overall movement  
 391 measurements were more accurate for the Piano Duo than both the Mixed Instrument Duos  
 392 (mean difference = 0.528, SE = 0.152,  $p = .004$ ) and the String Quartet (mean difference =  
 393 0.583, SE = 0.110,  $p < .001$ ); no significant difference was found between the Mixed  
 394 Instrument Duos and the String Quartet. Bonferroni-corrected, independent-samples t-tests  
 395 indicated that the optical flow technique exhibited more accurate performance for the Piano  
 396 Duo than the Mixed Instrument Duos ( $t(17) = 4.06, p = .009$ ) and the String Quartet ( $t(26) =$   
 397  $6.80, p < .001$ ). The KCF technique also achieved more accurate performance for the Piano  
 398 Duo than the String Quartet ( $t(26) = 3.39, p = .018$ ). All other pairwise comparisons of the  
 399 three datasets by computer vision technique failed to reach statistical significance.

400 An analogous 3x3 mixed ANOVA was conducted for the data using the head ROIs. A  
 401 significant effect of computer vision technique was found ( $F(2, 60) = 24.23, p < .001, \eta_p^2 =$   
 402  $.447$ ), with no significant effect of dataset ( $F(2, 30) = 3.14, p = .058, \eta_p^2 = .173$ ). The  
 403 technique by dataset interaction term was statistically significant ( $F(4, 60) = 5.59, p = .001,$   
 404  $\eta_p^2 = .272$ ). Bonferroni-corrected, paired-samples t-tests revealed that optical flow and KCF  
 405 both provided more accurate measures of performers' movements than frame differencing  
 406 ( $t(32) = 3.88, p = .001$  and  $t(32) = 8.38, p < .001$ , respectively) and KCF provided a more  
 407 accurate measure than optical flow ( $t(32) = 2.77, p = .027$ ). Bonferroni-corrected,  
 408 independent-samples t-tests indicated that the optical flow technique achieved more accurate  
 409 performance for the Piano Duo than the String Quartet ( $t(26) = 4.42, p = .001$ ). All other  
 410 pairwise comparisons of the three datasets by computer vision technique failed to reach  
 411 statistical significance.

412 Finally, we compared performance of the computer vision techniques between the upper  
 413 body ROI versus the head ROI. A paired-samples t-test indicated that movement  
 414 measurement was more accurate overall when restricted to a smaller ROI (the head) than a  
 415 larger ROI (upper body),  $t(98) = 2.54, p = .013$ .

416 We next looked in more detail at tracking in the horizontal versus vertical dimensions for  
 417 both optical flow and KCF, as compared to the MoCap data. These results are displayed in  
 418 Table 3, broken down by tracking dimension. Paired-samples t-tests for both the optical flow  
 419 (upper body ROI:  $t(32) = 6.22, p < .001$ ; head ROI:  $t(32) = 5.21, p < .001$ ) and KCF data  
 420 (upper body ROI:  $t(32) = 6.82, p < .001$ ; head ROI:  $t(32) = 5.77, p < .001$ ) indicated that  
 421 tracking by the computer vision techniques was significantly more accurate in the horizontal  
 422 than the vertical dimension. To probe this difference further, we explored whether the overall  
 423 lower performance in vertical movement tracking might be due to the computer vision  
 424 techniques also picking up on the missing, third dimension (depth) in which movement can  
 425 be made, in addition to the vertical dimension. It is plausible that this might especially be the  
 426 case when a performer is orthogonal to the video camera, and thus movement forward and  
 427 backward appears in the video as increases or decreases in the size of the performer. We  
 428 conducted two sets of regression analyses in which 1) the vertical dimension of the MoCap  
 429 data was used as a predictor of the vertical dimension of the video data and 2) the vertical  
 430 dimension of the MoCap data *and* the depth dimension of the MoCap data were used as  
 431 predictors of the vertical dimension of the video data. We then computed the change in  
 432 adjusted  $R^2$  values between the two regression analyses. For optical flow analysis, the  
 433 adjusted  $R^2$  values for the Mixed Instrument Duos and String Quartet only increased on  
 434 average by 0.03 and 0.06 respectively when taking the third MoCap dimension into account.  
 435 In both of these datasets the performers were viewed from the side or were situated  
 436 diagonally with respect to the camera (see Figure 1). However, in the Piano Duo, where the  
 437 performers were seated orthogonally to the camera (see Figure 1), the  $R^2$  values of the  
 438 regression models increased on average by 0.22 when the depth dimension of the MoCap  
 439 data was added as a predictor in addition to the vertical dimension. Although all of the  
 440 increases in adjusted  $R^2$  values were statistically significant (Mixed Instrument Duos:  $t(9) =$   
 441  $2.59, p = .029$ ; String Quartet:  $t(27) = 3.92, p = .001$ ; Piano Duo:  $t(27) = 3.99, p < .001$ ), the  
 442 raw adjusted  $R^2$  values indicate that the inclusion of the depth dimension made the most  
 443 substantial contribution to explaining the previously unaccounted variance in the Piano Duo.  
 444 A similar pattern emerged for the KCF data (adjusted  $R^2$  change values: Piano Duo = 0.13,  
 445 Mixed Instrument Duos = 0.03, String Quartet = 0.06). This change was statistically  
 446 significant within the Piano Duo ( $t(27) = 3.95, p = .001$ ) and String Quartet datasets ( $t(27) =$   
 447  $3.44, p = .002$ ) but not the Mixed Instrument Duos ( $t(9) = 2.12, p = .063$ ).

448

449 -INSERT TABLE 3 ABOUT HERE-

450

451 **4. Discussion**

452 The results of the present study indicate that the quantification of movement of musical  
453 performers from video using computer vision techniques closely approximates measurements  
454 from more sophisticated and costly technologies such as motion capture systems under  
455 certain conditions. Specifically, frame differencing, optical flow, and KCF techniques all  
456 achieved generally high correlations with MoCap data collected from the same musical  
457 performances, with median correlation values of .75 to .94, depending on the ROI, dataset,  
458 and computer vision technique. These results are in line with the work of Romero et al.  
459 (2016), who found specifically that frame differencing methods could provide a close  
460 approximation to MoCap data when tracking movement during social coordination tasks  
461 involving tapping, pointing, and clapping. It should also be noted that the promising results of  
462 the present study were obtained despite the fact that the video datasets were originally  
463 collected as a secondary measure to MoCap and the performance settings were not optimised  
464 with video data collection or computer vision analysis in mind. This suggests that the  
465 performance of these computer vision techniques might improve even further when working  
466 with video data that is optimised for the present research purposes, but also that existing  
467 video corpora that have been compiled for other aims could still provide promising data  
468 sources for subsequent research in which quantification of movement from video is required.

469 Our results also extend previous research (e.g., Paxton and Dale, 2012; Romero et al., 2016)  
470 by suggesting that the more recently developed, two-dimensional tracking techniques (optical  
471 flow and KCF) tend to outperform the older method of frame differencing. In addition,  
472 tracking of the head within the head ROI was more accurate overall than tracking of the torso  
473 within the upper body ROI. The KCF technique in particular displayed marked performance  
474 improvements in comparison to the other two techniques when constrained to the head ROI  
475 as compared to the upper body. A plausible explanation for the improved performance within  
476 the head ROIs is that the larger ROIs set around the upper body contain a variety of sources  
477 of movement, including not just torso movement but head movement and, in some cases,  
478 hands, bows of stringed instruments, etc., thereby resulting in decreased tracking accuracy of  
479 the torso. Researchers aiming to make use of larger ROIs (such as the upper body ROI from  
480 our study) to address particular research questions in the future might note that we were still  
481 able to provide a reasonable approximation of overall movement of musical performers as  
482 compared to MoCap data. However, it can be difficult to differentiate between various  
483 sources of movement within a large ROI, for example, sound-producing/instrument-specific  
484 movements (e.g., movement of the violin bow or shifting of the left hand up and down the  
485 neck of a cello) versus sound-facilitating/ancillary gestures (e.g., head nods or swaying  
486 together in time). Thus, ROI size should be taken into account in future research when the  
487 objective is to track movement from specific body parts or to measure only specific types of  
488 movement. On the other hand, if the objective is to provide an overall estimate of a  
489 performer's movement and there is no need to clarify the body part from which the  
490 movement originates or its expressive/functional purpose a larger ROI could still be suitable.

491 Within the present study the two-dimensional computer vision techniques exhibited greater  
492 precision in tracking horizontal than vertical movement. This seems to be at least partially  
493 explained by the missing dimension (depth) that cannot be precisely tracked by video

494 analysis methods in the same way as afforded by MoCap. The implication of this finding is  
495 that studies which aim to track precise directionality of vertical movement such as head nods  
496 might encounter a certain degree of measurement error, whereas horizontal movements such  
497 as side-to-side swaying can be tracked with a greater degree of spatial precision. However,  
498 combining these two tracking dimensions into polar coordinates (as in Table 2) tends to  
499 provide a good approximation of the overall movement of a performer, with median  
500 correlations above .80 for the upper body ROI and above .90 for the head ROI in both optical  
501 flow and KCF. Another possible avenue for future research would be to record video of  
502 musical performances using multiple camera angles in an attempt to recover the missing third  
503 dimension that cannot be measured from the present video data.

504 Some differences between the three datasets emerged, particularly in regard to the upper body  
505 ROI. In general, measurements of performers' movements were more accurate for the Piano  
506 Duo than the String Quartet and, in some cases, the Mixed Instrument Duos. This may be  
507 due, at least in part, to the fact that within the String Quartet dataset and certain examples  
508 from the Mixed Instrument Duos (cellist and double bassist), the bows of the violinist/violist  
509 and the left hands of the cellist/double bassist often entered the ROIs and created an extra  
510 source of motion that could be picked up by the computer vision techniques. This was the  
511 case even when the ROI was focused around the head, as the bow or left hand sometimes  
512 occluded the face. These cases provide examples of a discrepancy in differentiating the  
513 sound-producing, instrumental movements of a performer from ancillary movements of the  
514 head, and highlight that the specific demands and idiosyncrasies of performing on certain  
515 instruments should be taken into account when conducting research that aims to quantify  
516 musicians' movements from video. In the case of the string quartet, a different camera angle  
517 could be considered to avoid occlusions within the ROI. Or, depending on the research  
518 question of interest, other body parts could be tracked that do not present this occlusion  
519 problem, for instance, the tapping of performers' feet in time to the music.

520 It should also be noted that some of the differences in movement tracking/quantification  
521 accuracy between the three datasets could have arisen from differences in the video source  
522 material, such as lighting, camera angle, and distance of the performers from the camera.  
523 Future research should aim to test the independent contributions of each of these factors.  
524 Additionally, there may have been fundamental differences between the *types* of ancillary  
525 movements that performers in the different datasets made, which could be affected both by  
526 the instrument being played and the musical style itself (e.g. free jazz improvisation and  
527 notated string quartets might require different types of communicative gestures for different  
528 purposes). Although classifying movement types is beyond the scope of the present study,  
529 future research could also test whether certain classes of body movements are more  
530 accurately tracked than others.

531 These results open new avenues for researchers of musical movement. In our own future  
532 research we aim to apply some of these computer vision techniques to examine how the  
533 relationships between the movements of co-performers stabilise or change over time and how  
534 these corporeal relationships affect audience appraisals of a performance. We also aim to  
535 conduct cross-cultural comparisons of what it means to "play in time together" within

536 different musical traditions, using music that is performed for a variety of different functions  
 537 (e.g., rituals, dance, concert performance, etc.; Clayton, 2013). Additional possible  
 538 applications of these computer vision techniques for future research could include the study  
 539 of leader-follower relationships, the relationship between visual movement coordination and  
 540 synchrony/asynchrony in the auditory modality, and studies of movement coordination  
 541 differences between expert versus novice performers.

542

## 543 **References**

- 544 Alborno, P., Volpe, G., Camurri, A., Clayton, M., and Keller, P. (2015). Automated video  
 545 analysis of interpersonal entrainment in Indian music performance. In *7th International*  
 546 *Conference on Intelligent Technologies for Interactive Entertainment (INTETAIN)* (pp. 57-  
 547 63). IEEE. doi: 10.4108/icst.intetain.2015.259521
- 548 Bolme, D. S., Draper, B. A., and Beveridge, J. R. (2009). Average of synthetic exact filters.  
 549 Proceedings from *Computer Vision and Pattern Recognition* (pp. 2105-2112). IEEE. doi:  
 550 10.1109/CVPR.2009.5206701
- 551 Burger, B. and Toiviainen, P. (2013). MoCap Toolbox-A Matlab toolbox for computational  
 552 analysis of movement data. In *Proceedings of the Sound and Music Computing Conference*  
 553 *2013*. Berlin: Logos Verlag Berlin. ISBN 978-3-8325-3472-1
- 554 Camurri, A., Canepa, C., Coletta, P., Cavallero, F., Ghisio, S., Glowinski, D., and Volpe, G.  
 555 (2010). Active experience of audiovisual cultural content: The virtual binocular interface. In  
 556 *Proceedings of the Second Workshop on eHeritage and Digital Art Preservation* (pp. 37-42).  
 557 Firenze, Italy. doi: 10.1145/1877922.1877934
- 558 Clayton, M (2013). Entrainment, ethnography and musical interaction. In M. Clayton, B.  
 559 Dueck, & L. Leante (Eds.), *Experience and Meaning in Music Performance* (pp. 17-39).  
 560 Oxford: Oxford University Press.
- 561 Clayton, M., Sager, R., and Will, U. (2005). In time with the music: The concept of  
 562 entrainment and its significance for ethnomusicology. In *European Meetings in*  
 563 *Ethnomusicology* 11: 1-82.
- 564 Davidson, J. W. (2009). Movement and collaboration in musical performance. In S. Hallam,  
 565 I. Cross, & M. Thaut (Eds.), *The Oxford Handbook of Music Psychology* (pp. 364-376).  
 566 Oxford: Oxford University Press.
- 567 Farneback, G. (2003). Two-frame motion estimation based on polynomial expansion. In J.  
 568 Bigun and T. Gustavsson (Eds.), Proceedings from *13th Scandinavian Conference on Image*  
 569 *Analysis* (pp. 363-370). Heidelberg: Springer Berlin.
- 570 Fleet, D. and Weiss, Y. (2006). Optical flow estimation. In N. Paragios, Y. Chen, & D.  
 571 Olivier (Eds.), *Handbook of mathematical models in computer vision* (pp. 237-257). Springer  
 572 US.

- 573 Forsyth, D. A. and Ponce, J. (2002). *Computer vision: A modern approach*. Englewood  
574 Cliffs, NJ: Prentice Hall Professional Technical Reference.
- 575 Glowinski, D., Gnecco, G., Piano, S. and Camurri, A. (2013). Expressive non-verbal  
576 interaction in string quartet. In Proceedings of *Conference on Affective Computing and*  
577 *Intelligent Interaction (ACII 2013)*. Geneva, Switzerland.
- 578 Henriques, J. F., Caseiro, R., Martins, P., and Batista, J. (2015). High-speed tracking with  
579 kernelized correlation filters. *IEEE Transactions on Pattern Analysis and Machine*  
580 *Intelligence* 37: 583-596. doi: 10.1109/TPAMI.2014.2345390
- 581 Hester, C. F. and Casasent, D. (1980). Multivariant technique for multiclass pattern  
582 recognition. *Applied Optics* 19: 1758-1761.
- 583 Jensenius, A. R., Godøy, R. I., and Wanderley, M. M. (2005). Developing tools for studying  
584 musical gestures within the Max/MSP/Jitter environment. In *Proceedings of the International*  
585 *Computer Music Conference* (pp. 282-285). ISSN 2223-3881
- 586 Jensenius, A. R., Wanderley, M. M., Godøy, R. I., and Leman, M. (2010). Concepts and  
587 methods in research on music-related gestures. In R.I. Godøy & M. Leman (Eds.), *Musical*  
588 *Gestures: Sound, Movement, and Meaning* (pp. 12–35). New York: Routledge.
- 589 JVC Video Anthology of World Music and Dance (1990). Tokyo: JVC, Victor Company of  
590 Japan.
- 591 Keller, P.E. (2014). Ensemble performance: Interpersonal alignment of musical expression.  
592 In D. Fabian, R. Timmers, & E. Schubert (Eds.), *Expressiveness in Music Performance:*  
593 *Empirical Approaches Across Styles and Cultures* (pp. 260-282). Oxford: Oxford University  
594 Press.
- 595 Keller, P.E. and Appel, M. (2010). Individual differences, auditory imagery, and the  
596 coordination of body movements and sounds in musical ensembles. *Music Perception* 28: 27-  
597 46. doi: 10.1525/mp.2010.28.1.27
- 598 Latif, N., Barbosa, A. V., Vatiokiotis-Bateson, E., Castelhana, M. S. and Munhall, K. G.  
599 (2014). Movement coordination during conversation. *PloS One* 9: e105036.
- 600 Moeslund, T. B. and Granum, E. (2001). A survey of computer vision-based human motion  
601 capture. *Computer Vision and Image Understanding* 81: 231-268. doi:  
602 10.1006/cviu.2000.0897
- 603 Moran, N., Hadley, L. V., Bader, M. and Keller, P. E. (2015). Perception of ‘back-  
604 channeling’ nonverbal feedback in musical duo improvisation. *PloS One* 10: e0130070.
- 605 Paxton, A. and Dale, R. (2013). Frame-differencing methods for measuring bodily synchrony  
606 in conversation. *Behavior Research Methods* 45: 329-343. doi: 10.3758/s13428-012-0249-2

- 607 Ragert, M., Schroeder, T., and Keller, P.E. (2013). Knowing too little or too much: The  
 608 effects of familiarity with a co-performer's part on interpersonal coordination in musical  
 609 ensembles. *Frontiers in Auditory Cognitive Neuroscience* 4: 368. doi:  
 610 10.3389/fpsyg.2013.00368
- 611 Rasch, R. A. (1988). Timing and synchronization in ensemble performance. In J. Sloboda  
 612 (Ed.), *Generative Processes in Music: The Psychology of Performance, Improvisation, and*  
 613 *Composition* (pp. 70-90). Oxford: Oxford University Press.
- 614 Romero, V., Amaral, J., Fitzpatrick, P., Schmidt, R. C., Duncan, A. W., and Richardson, M.  
 615 J. (2016). Can low-cost motion-tracking systems substitute a Polhemus system when  
 616 researching social motor coordination in children? *Behavior Research Methods*. doi:  
 617 10.3758/s13428-016-0733-1
- 618 Shaffer, L. H. (1984). Timing in solo and duet piano performances. *The Quarterly Journal of*  
 619 *Experimental Psychology* 36: 577-595. doi: 10.1080/14640748408402180
- 620 Solberg, R. T. and Jensenius, A. R. (2016). Optical or inertial? Evaluation of two motion  
 621 capture systems for studies of dancing to electronic dance music. Proceedings from *Sound*  
 622 *and Music Computing*, Hamburg, Germany. ISSN 2518-3672
- 623 Wanderley, M. M., Vines, B. W., Middleton, N., McKay, C., and Hatch, W. (2005). The  
 624 musical significance of clarinetists' ancillary gestures: An exploration of the field. *Journal of*  
 625 *New Music Research* 34: 97-113. doi: 10.1080/09298210500124208
- 626 Wren, C. R., Azarbayejani, A., Darrell, T., and Pentland, A. P. (1997). Pfinder: Real-time  
 627 tracking of the human body. *IEEE Transactions on Pattern Analysis and Machine*  
 628 *Intelligence* 19: 780-785. doi: 10.1109/34.598236

629

### 630 **Conflict of Interest Statement**

631 This research was conducted in the absence of any commercial or financial relationships that  
 632 could be construed as a potential conflict of interest.

### 633 **Funding Statement**

634 This work was supported by the Arts and Humanities Research Council [grant number  
 635 AH/N00308X/1].

### 636 **Acknowledgments**

637 The authors would like to thank Peter Keller and Nikki Moran for sharing video and motion  
 638 capture data for use in this project and providing feedback on an earlier version of this paper.  
 639 Thanks to Simone Tarsitani for help with the data preparation and storage. Recording credit  
 640 for the Piano Duo goes to Marie Ragert, Kerstin Traeger, Maria Bader, and Jan Bergmann,  
 641 and recording credit for the Mixed Instrument Duos goes to Kerstin Traeger, Maria Bader,

642 and Jan Bergmann. Recording credit for the String Quartet goes to Corrado Canepa, Paolo  
643 Coletta, Nicola Ferrari, Simone Ghisio, Donald Glowinski, Giorgio Gnecco, Maurizio  
644 Mancini, Stefano Piana, Roberto Sagoleo, and Giovanna Varni. Thanks also to the musicians  
645 of the Piano Duo, Mixed Instrument Duos, and the Quartetto di Cremona string quartet, for  
646 their crucial contribution to the creation of the video and motion capture datasets.

#### 647 **Figure captions**

648 Figure 1. Screenshots of one example video from each of the three datasets.

649 Figure 2. An example of the data preparation and extraction process from the Mixed  
650 Instrument Duos. The left panel shows the selection of ROIs for the head of each performer  
651 and a corresponding head marker from the MoCap data. The right panel shows the KCF data  
652 and MoCap trajectories for the x- and y-coordinates of each performer's head as time series.

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672 Table 1. Summary of Performance Details for Each Dataset

Dataset	No. of video recordings	No. of different performers	No. of trials analysed*	Instrumentation	Mean Duration in seconds (SD)
Piano Duo	7	2	14	two pianists/ vocalists	119.68 (1.78)
Mixed Instrument Duos	3	5	5	cellist, soprano saxophonist, double bassist, two pianists	76.08 (43.14)
String Quartet	8	3	14	violinist, violist, cellist	125.74 (22.92)
<b>Total</b>	<b>18</b>	<b>10</b>	<b>33</b>	<b>6 instruments</b>	<b>115.11 (27.70)</b>

673 \*Note: A trial was defined as one video-recorded performance by one performer.

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695 Table 2. Median Correlations between the Computer Vision and Motion Capture Data

Region of Interest	Dataset	Number of trials	FD: Median correlation (SD)	OF: Median correlation (SD)	KCF: Median correlation (SD)
Upper Body	Piano Duo	14	.80 (0.14)	.98 (0.07)	.89 (0.09)
	Mixed	5	.71 (0.16)	.85 (0.06)	.80 (0.07)
	Instrument Duos				
	String Quartet	14	.77 (0.08)	.75 (0.26)	.77 (0.25)
	<b>All Datasets</b>	<b>33</b>	<b>.75 (0.13)</b>	<b>.87 (0.22)</b>	<b>.84 (0.20)</b>
Head	Piano Duo	14	.73 (0.18)	.94 (0.03)	.95 (0.06)
	Mixed	5	.72 (0.17)	.92 (0.13)	.92 (0.18)
	Instrument Duos				
	String Quartet	14	.83 (0.13)	.80 (0.21)	.94 (0.07)
	<b>All Datasets</b>	<b>33</b>	<b>.79 (0.16)</b>	<b>.91 (0.17)</b>	<b>.94 (0.10)</b>

696 Note: FD = Frame Differencing, OF = Optical Flow, KCF = Kernelized Correlation Filters; x- and y-coordinates  
697 are combined into polar coordinates for motion capture, OF, and KCF data

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716 Table 3. Median Correlations between the OF/ KCF and Motion Capture Data, by Dimension

Technique	Region of Interest	Dataset	Number of trials	Median correlation, x-dimension (SD)	Median correlation, y-dimension (SD)		
OF	Upper Body	Piano Duo	14	.98 (0.06)	.93 (0.45)		
		Mixed Instrument Duos	5	.85 (0.05)	.66 (0.55)		
		String Quartet	14	.74 (0.25)	.39 (0.28)		
		<b>All Datasets</b>	<b>33</b>	<b>.87 (0.22)</b>	<b>.65 (0.41)</b>		
	Head	Piano Duo	14	.94 (0.03)	.82 (0.22)		
		Mixed Instrument Duos	5	.91 (0.11)	.85 (0.05)		
		String Quartet	14	.81 (0.20)	.57 (0.17)		
		<b>All Datasets</b>	<b>33</b>	<b>.92 (0.16)</b>	<b>.75 (0.21)</b>		
		KCF	Upper Body	Piano Duo	14	.86 (0.10)	.62 (0.34)
				Mixed Instrument Duos	5	.79 (0.07)	.45 (0.51)
String Quartet	14			.75 (0.24)	.33 (0.41)		
<b>All Datasets</b>	<b>33</b>			<b>.80 (0.19)</b>	<b>.55 (0.41)</b>		
Head	Piano Duo		14	.96 (0.04)	.60 (0.28)		
	Mixed Instrument Duos		5	.91 (0.17)	.78 (0.09)		
	String Quartet		14	.93 (0.07)	.80 (0.17)		
	<b>All Datasets</b>		<b>33</b>	<b>.94 (0.09)</b>	<b>.78 (0.22)</b>		

717 Note: OF = Optical Flow, KCF = Kernelized Correlation Filters

718



Piano Duo



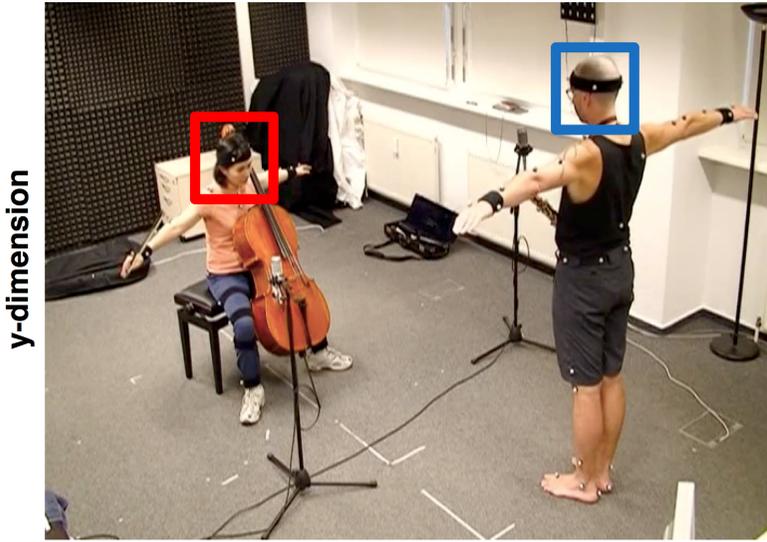
Mixed Instrument Duos



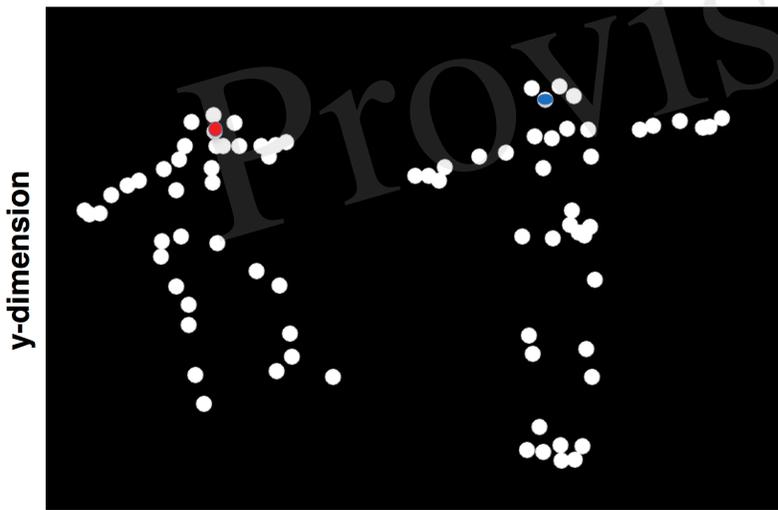
String Quartet

Provisional

**Performer 1**      **Performer 2**

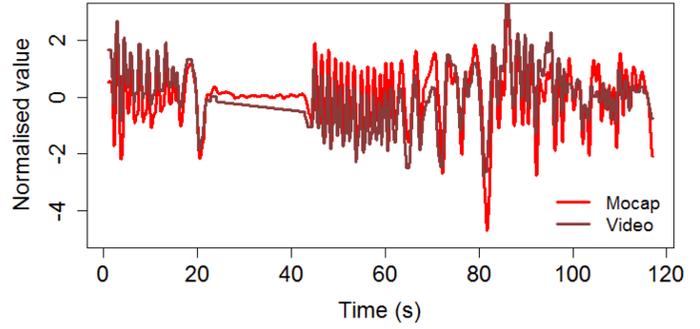


**Video data**

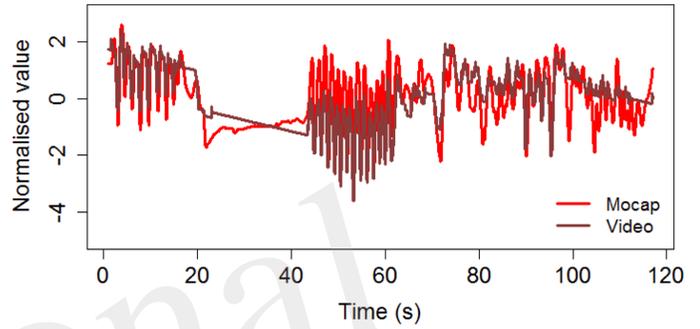


**MoCap data**

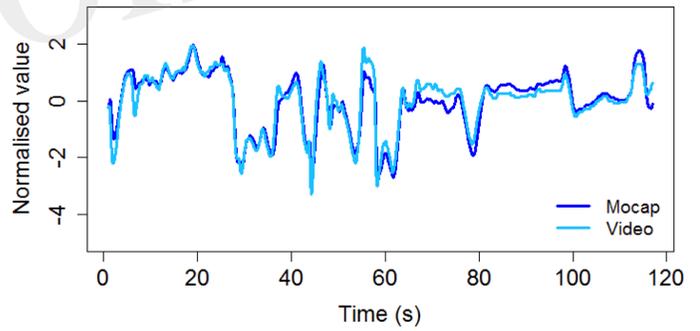
**Performer 1: x-dimension (r = .69)**



**Performer 1: y-dimension (r = .71)**



**Performer 2: x-dimension (r = .94)**



**Performer 2: y-dimension (r = .91)**

