

Journal Pre-proof

Arm motion symmetry in conversation

Jonathan Windle, Sarah Taylor, David Greenwood, Iain Matthews

PII: S0167-6393(22)00105-4
DOI: <https://doi.org/10.1016/j.specom.2022.08.001>
Reference: SPECOM 2883

To appear in: *Speech Communication*

Received date: 22 February 2022
Revised date: 17 June 2022
Accepted date: 16 August 2022



Please cite this article as: J. Windle, S. Taylor, D. Greenwood et al., Arm motion symmetry in conversation. *Speech Communication* (2022), doi: <https://doi.org/10.1016/j.specom.2022.08.001>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2022 Published by Elsevier B.V.

Arm Motion Symmetry In Conversation

Jonathan Windle, Sarah Taylor, David Greenwood, Iain Matthews

Abstract Data-driven synthesis of human motion during conversational speech is an active research area with applications that include character animation, computer gaming and conversational agents. Natural looking motion is key to both perceived realism and understanding of any synthesised animation. Multi-modal speech and body-motion data is scarce and limited, so it is common to augment real motion data by mirroring the body pose to double the number of training samples. This augmentation is based on the assumption that a person's gesturing is not affected by handedness and that the reflected pose is plausible. In this study, we explore the validity of this assumption by evaluating the reflective symmetry of a speaker's arms during conversational exchanges. We analyse the left and right arm motion of 36 subjects during dyadic conversation and present the per-frame symmetry of the arm gestures. To identify temporal offsets caused by the presence of a leading hand, we compute the time lag between movements of the left and right arms. We perform a nearest neighbour search to test the validity of any mirrored pose. We also consider information theory to examine the information gain from mirroring the data. [We implement a speech-to-gesture generative model to determine the efficacy of lateral mirroring techniques for data augmentation.](#) Our findings suggest that both positional symmetry and left-right motion offsets vary from speaker to speaker. [We conclude that data augmentation by mirroring is valid in certain cases when considering the mirrored pose as a new virtual identity, but that it should be carefully considered as a generic approach if the gesturing style and handedness of the original speaker is to be maintained.](#)

1 Introduction

Co-speech gesturing contributes to language production and perception during conversation. Gesturing provides semantic context, and may be indicative of emotion and emphasis (Kendon, 1994; McNeill, 1985; Studdert-Kennedy, 1994; De Ruiter et al., 2012). Gesturing in conversational speech serves many purposes including contributing to increased understanding, turn taking and listener feedback. Given the multi-modal nature of conversation, it follows that there is a co-dependency between speech and gesture.

Data-driven approaches for automatically driving body motion from speech is an active research area (Alexanderson et al., 2020a,b; Henter et al., 2020; Korzun et al., 2020; Yoon et al., 2020; Ginosar et al., 2019). Applications for these conversational agents include character animation, computer gaming and codec avatars (Bagautdinov et al., 2021). Such systems require multi-modal data comprised of motion captured body pose with a corresponding audio signal. These datasets are typically time-consuming and both financially and computationally expensive to capture, therefore, availability is scarce. A [practised](#) augmentation approach is lateral mirroring (Henter et al., 2020; Alexanderson et al., 2020b; Gong et al., 2021). This is to flip the left and right sided motion with each other.

While lateral mirroring effectively doubles the

amount of training data, we raise the question of how natural and appropriate this augmented data is. [Asymmetry is known to occur in pose from physical body constraints and gesture style types.](#) We present a study of frame-by-frame position and temporal characteristics to investigate if this mirroring produces natural speaker-dependent movement. [This study is not only relevant to gesture generation and data augmentation, it provides an insight into arm symmetry during conversation, providing greater understanding for all relevant fields of research such as gesture recognition and gesture behaviour.](#) Finally we consider the use of this method of analysis as a means to evaluate performance of data-driven synthesised motion.

2 Related Work

We present a review on works relating to speech gesturing, body motion datasets, methods for speech-driven body animation, and techniques for data augmentation used by these methods.

2.1 Arm Gesture and Symmetry

Neither speech nor gesture alone allows a speaker to communicate to their full efficiency. Removing either of these modalities leads to a reduction in semiotic ver-

51 satility (Wagner et al., 2014) and communicative un-
 52 derstanding (Hostetter, 2011). One reason for this is
 53 that each modality represents certain information bet-
 54 ter than the other. For example, hands might better de-
 55 scribe shape or direction by providing visual cues. The
 56 gestures that form these cues may or may not be sym-
 57 metrical, and this may, in part, depend on the particular
 58 shape or direction being described.

59 Environmental conditions contribute a great deal to
 60 the importance of each modality during a conversation.
 61 A small and enclosed space may cause a person to be
 62 conservative with their gesturing, whereas to commu-
 63 nicate the same speech in an expansive, outside envi-
 64 ronment, a person may gesture more actively as they
 65 have more space. Proximity and facing direction of the
 66 conversational partner within the environment will also
 67 effect the extent and type of gesturing. If conversation
 68 is taking place while walking alongside their partner,
 69 this will prompt different behaviour to a static face-to-
 70 face interaction. Similarly, if the partner is far away,
 71 gestures may be emphasised to account for the reduc-
 72 tion in the received audio volume. It has been found
 73 that gesture activity increases during adverse listen-
 74 ing conditions, such as acoustic noise and non-native
 75 speaking conversational partners (Drijvers et al., 2018).

76 Objects surrounding or colliding with the speaker
 77 introduce physical constraints that inhibit or otherwise
 78 affect gesturing. For instance, a wall to one side of
 79 the speaker will limit their available gesture space, con-
 80 strain physical activity and likely increase asymmetry.
 81 Similarly, a speaker's hand might be occupied with an
 82 object such as a glass of water, which would alter ges-
 83 tural behaviour.

84 Individuals exhibit gestural idiosyncrasies. Some
 85 speakers may commonly perform self-adaptor traits
 86 such as self-touching or scratching. Others may have
 87 physiological restrictions, making particular gestures
 88 impossible and affecting the realisation of others. In
 89 each of these cases, asymmetry in the positioning of
 90 the arms is likely.

91 The amount of conversational gesturing that takes
 92 place during an interaction can be linked to a speaker's
 93 personality. It has been found that a speaker's *Big Five*
 94 personality traits (extroversion, neuroticism, conscien-
 95 tiousness, agreeableness and openness to experience)
 96 are correlated with the amount of gesture production
 97 (Hostetter, 2011). In particular, extroversion is posi-
 98 tively correlated with representational gesture produc-
 99 tion, which might be due to extroverted people having
 100 high amounts of energy in social situations and there-
 101 fore gesturing regardless of communicative effect.

102 McNeill defined a gesture space (McNeill, 2011),
 103 stating that the majority of gestures happen in the *cen-*
 104 *tral gesture space* which encompasses the area below
 105 the neck and between the shoulders and elbows. *Pe-*

106 *ripheral gesture space* encapsulates gestures performed
 107 outside of the central gesture space and can be thought
 108 of as the extremes of gesturing. They suggest that the
 109 peripheral gestures aim to capture visual attention.

110 McNeill also defined a classification on the se-
 111 mantic functions of gesture types (McNeill, 2011).
 112 They categorised gestures as either emblematic, iconic
 113 metaphoric, deictic or beat: *Emblematic gestures* bear
 114 a conventionalised meaning; *Iconic gestures* resemble
 115 a certain physical aspect of the conveyed information;
 116 *Metaphoric gesture* is an Iconic gesture resembling ab-
 117 stract content; *Deictic gesture* point out locations in
 118 space; and *Beat gestures* are simple and fast movements
 119 of the hands commonly synchronised with prosodic
 120 events in speech (Pouw et al., 2020). However, in prac-
 121 tice a gesture may perform many semantic functions,
 122 and it has instead been proposed to treat each gesture
 123 category as a dimension on which gestures load to dif-
 124 ferent degrees (McNeill, 2008).

125 A speaker's handedness has been found to impact
 126 gesture production, particularly regarding the position-
 127 ing of the left and right arms. It has been found that
 128 beat-style gestures were more commonly performed
 129 with a speaker's dominant hand, while representational
 130 gestures in right-handed speakers had a right-handed
 131 preference while left-handed speakers did not have a
 132 hand preference (Çatak et al., 2018). There is an as-
 133 sociation between gestural handedness and the emo-
 134 tional dimensions of pleasure and arousal. Kipp and
 135 Martin (Kipp and Martin, 2009) found significant cor-
 136 relation between emotion category and handedness of
 137 the gesture, where speakers consistently used their left
 138 hands to gesture during a relaxed, positive mood and
 139 their right hands to gesture when in a negative, aggres-
 140 sive mood.

141 We have reviewed works that analyse gestural sym-
 142 metry during conversation, however, these works are
 143 limited by the data used. Data is often observed man-
 144 ually from video (McNeill, 2011) or limited to a few
 145 speakers worth of data (Kipp and Martin, 2009). This
 146 reveals a limitation in current studies that we aim to
 147 address.

148 2.2 Body Motion Data and Limitations

149 Conversational body motion data is needed for per-
 150 forming analysis of gestural symmetry, and for training
 151 generative speech-to-body animation models. How-
 152 ever, the availability of such data is scarce and issues
 153 commonly arise during the data collection process re-
 154 sulting in data that is noisy, unnatural or lacking in
 155 quantity. Ideally, motion data is recorded using optical
 156 motion capture systems that track retroreflective mark-
 157 ers on the speaker. The 3D position of each marker is
 158 triangulated between multiple cameras. Issues regard-
 159 ing marker jitter, swapping and occlusion often require

160 motion captured landmarks to be manually cleaned.
 161 Generally, motion capture is both financially and com-
 162 putationally expensive to collect, but can result in high-
 163 quality performance capture. An abundance of body
 164 motion data is available if we use video as a data source.
 165 However, extracting 3D key points from a single video
 166 feed is challenging, often leading to noise and inaccu-
 167 rate depth estimation. This causes a trade off between
 168 data quality and quantity.

169 A dataset that was collected for data-driven syn-
 170 thesis of motion is the Trinity dataset (Ferstl and Mc-
 171 Donnell, 2018). It contains 244 minutes of speech and
 172 motion data that was recorded using 20 Vicon cam-
 173 eras, and the motion data is high quality and accurate.
 174 However, the Trinity dataset contains only one male
 175 speaker producing monologue speech. Gestural motion
 176 and symmetry varies across speakers and therefore it
 177 is difficult to draw conclusions from a single speaker.
 178 Since the speech is monologue, the gesturing that re-
 179 lates to listener understanding and turn taking is also
 180 not captured.

181 Social interaction is not limited to conversation.
 182 Joo et al. (Joo et al., 2015) presented a dataset that
 183 contains social interactions during game scenarios, to-
 184 gether with a description of the Panoptic Studio that
 185 was used for the capture. The capture system is com-
 186 prised of a large dome structure containing 480 VGA
 187 cameras for video capture, each with calibrated frame
 188 timers and positions. Using the known positions of the
 189 cameras and 2D pose estimation software, 3D poses are
 190 accurately predicted. While this system produces clean
 191 motion capture, it is both financially and computation-
 192 ally expensive. With 480 cameras, the data-rate is ap-
 193 proximately 29.4 Gbps, requiring a large amount of pro-
 194 cessing power and storage to manage such quantities
 195 of data. While this dataset provides multiple speakers’
 196 motions, the scenarios recorded are not natural conver-
 197 sations but instead social interactions during games,
 198 which will affect the types of gestures that are pro-
 199 duced.

200 There is an abundance of video data available that
 201 contains conversational interaction. This is exploited
 202 by Ginosar et al. who extracted monologue speech and
 203 motion data from videos of talk show hosts, lecturers
 204 and televangelists (Ginosar et al., 2019). The videos are
 205 shot from a single view and therefore only 2D keypoints
 206 were extracted. Further work estimated 3D keypoints
 207 for this dataset (Habibie et al., 2021), however the result
 208 is noisy and includes errors in depth prediction.

209 The main limitations of existing motion captured
 210 data is the number of identities and lack of natural
 211 dyadic conversation. The Talking with Hands dataset
 212 presented by Lee et al. mitigates these limitations and
 213 is selected for our analysis (Lee et al., 2019). This dataset
 214 is described in Section 3.

215 2.3 Speech-driven Body Animation

216 Embodied conversational agents describe both human-
 217 like robots and animations that aim to employ human-
 218 realistic verbal and non-verbal communicative modal-
 219 ities. Data-driven approaches for automatically driv-
 220 ing body motion from speech is an active research area
 221 (Alexanderson et al., 2020a,b; Henter et al., 2020; Korzun
 222 et al., 2020; Yoon et al., 2020; Ginosar et al., 2019). These
 223 approaches aim to estimate a speakers pose, typically
 224 represented by a sparse set of skeleton joints, from their
 225 corresponding speech audio signal.

226 Recent approaches for data-driven motion synthe-
 227 sis typically involve deep learning (Alexanderson et al.,
 228 2020a,b; Henter et al., 2020; Korzun et al., 2020; Yoon
 229 et al., 2020; Ginosar et al., 2019). Their success is highly
 230 dependent on the data used to train them. For instance,
 231 small datasets or those lacking diversity can lead to
 232 models not generalising well or overfitting to training
 233 data (Perez and Wang, 2017). Data quality is also im-
 234 portant as a model can only learn to be as good as the
 235 training data, and inaccurate or poorly labelled data
 236 will cause the model to learn incorrect information. To
 237 mitigate the limited amount of available body motion
 238 data, it is common to augment the dataset. It is key to
 239 ensure that the quality of the data is not compromised
 240 during augmentation, and the focus of our work is to
 241 explore this.

242 2.4 Data Augmentation

243 Data augmentation are techniques used to increase the
 244 amount of data by adding slightly modified copies of
 245 real data or created synthetic data from existing data.
 246 The most common technique for this is through *data*
 247 *warping* defined in (Perez and Wang, 2017) as an ap-
 248 proach to directly augment the input data to the model
 249 in *data space*. Augmentation approaches vary depend-
 250 ing on the data type and the problem domain.

251 When working with image data it is common to ap-
 252 ply simple transformations on each image. These in-
 253 clude flipping, scaling, rotating, translating, noise in-
 254 jection and colour space transformation (Shorten and
 255 Khoshgoftaar, 2019). While flipping, scaling, rotating
 256 and translating are all possible to apply to a 3D skele-
 257 ton representation of body motion data, it is not nec-
 258 essarily appropriate. Scaling the skeleton by a differ-
 259 ent amount in each dimension would alter the iden-
 260 tity. If we scale by the same amount, and if joint an-
 261 gles are used to represent the skeleton pose, this scal-
 262 ing would not provide additional information as the an-
 263 gles would remain identical. Applying a global rotation
 264 to the skeleton might introduce unnatural positioning
 265 (e.g. losing foot contact with the ground). Translating
 266 the skeleton would not effectively augment the data
 267 as the speaker would still move in the same way, but

268 in a different location. Adding noise to the captured
 269 motion would cause unnatural, jittery motion. Flipping
 270 (or laterally mirroring) the skeleton is the only of these
 271 data augmentation approaches that still produces po-
 272 tentially valid human body motion. It is our goal to
 273 determine in what cases this augmentation is a valid
 274 approach.

275 3 Data and Pre-processing

276 This study performs an analysis on the body motion
 277 from the Talking with Hands dataset (Lee et al., 2019).
 278 The dataset consists of 16.2-million frames of motion
 279 at 90 Frames Per Second across 50 different speakers
 280 during dyadic conversation. Unfortunately not all of
 281 this data is currently publicly available and therefore
 282 the available subset of 36 speakers has been used. The
 283 majority of speakers were only captured in conversa-
 284 tion with one other speaker (*shallow* speakers), while
 285 a small number had multiple conversational partners
 286 (*deep* speakers). We removed any non-conversation
 287 segments of the data (e.g. T-Pose sequences) prior to
 288 performing the analysis.

289 The dataset provides a set of 3D skeleton joint key-
 290 points for each frame. Our study focuses on the arm
 291 movements, and considers only the 3D locations of
 292 the left and right shoulder, elbow, forearm and wrist.
 293 The skeleton was translated per frame such that the
 294 mid-point between each shoulder joint was at the ori-
 295 gin. This simplifies the analysis and accounts for large
 296 translations of arms from motion originating from the
 297 spine such as leaning forwards and backwards. This al-
 298 lows us to evaluate translations made by motion gen-
 299 erated from the arms independently of the rest of the
 300 pose. The coordinate system utilised in this paper is as
 301 follows:

- 302 • Y - Height (Up and Down)
- 303 • X - Depth (Back and Forth)
- 304 • Z - Width (Left and Right)

305 We also use a consistent colour scheme through all
 306 figures to represent each forearm. Cyan depicts the
 307 right forearm and Blue depicts the left forearm.

308 4 Mean Pose Symmetry

309 We first evaluate the symmetry of the mean poses for
 310 each speaker, aiming to reveal an impression of the per-
 311 speaker symmetry across all of their motion. Using all
 312 the frames of motion, the per-speaker mean pose is cal-
 313 culated. We then project the right arm to the space of
 314 the left arm by laterally mirroring (along the y-axis).
 315 To evaluate the arm symmetry, the Euclidean distance

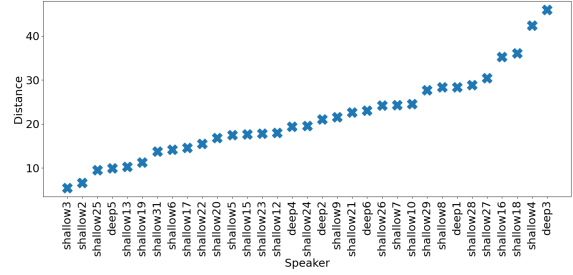


Figure 1: Euclidean distance between mirrored right arm and the left for each speaker.

316 between all joints in the left arm and projected right
 317 arm are calculated. The lower this distance, the closer
 318 the two arms are to each other, which is indicative of a
 319 more symmetrical pose.

320 We show the range of symmetry in Figure 1. We ob-
 321 serve that a person's mean pose is not always symmet-
 322 rical. Shallow3 is found to have the most symmetrical
 323 mean pose, whereas Deep3 has the most asymmetric
 324 pose according to the Euclidean distance.

325 From the 36 speakers we select the two with the
 326 highest and two with the lowest Euclidean distance,
 327 representing the subjects exhibiting the least and most
 328 arm symmetry in their mean pose. We visualise the
 329 level of symmetry by overlaying a perspective projec-
 330 tion of the mirrored right arm onto the left arm. Figure 2
 331 shows this projection from both a frontal and side view
 332 for each of the four speakers. There is clear asymmetry
 333 in the mean arm pose of Deep3 and Shallow4 (columns
 334 one and two). The left arm of Deep3 shows itself an-
 335 gled towards the right side of their body, whereas the
 336 right arm is pointing away from their body, towards the
 337 camera. Shallow4 orients their right wrist away from
 338 their body while their left wrist is pointing towards
 339 their body. At the other extreme, Shallow3 and Shal-
 340 low2 show good symmetry (columns three and four).
 341 In these examples, the mirrored right arm overlaps the
 342 left arm from the shoulder to the elbow with a slight
 343 divergence from the elbow to the wrist.

344 The largest differences between the arm positions is
 345 observed in the side view, whereby each of the left arms
 346 are positioned further forward than the right arms.
 347 While this observation is more prominent on the two
 348 most asymmetric speakers, it holds for each of the
 349 speakers in Figure 2.

350 5 Spatial Symmetry

351 The mean pose analysis in Section 4 provides an indi-
 352 cation of the symmetry of a speaker's most frequent
 353 (or neutral) arm positions. However, it does not explain

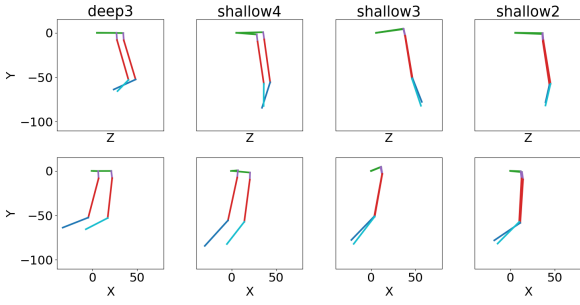


Figure 2: A projection of the mean pose for four speakers. In each case, the right arm (cyan forearm) has been mirrored and overlaid onto the left arm (blue forearm). Top row: front view. Bottom row: side view.

whether the *motion* of the arms is similar or symmetrical. In this section we investigate whether the observed asymmetry is an effect of a speaker's tendency to gesture more on one side than the other, and whether the arms occupy symmetrical gesture spaces. We use 3D keypoints to gather statistics regarding the arm motion of each speaker, discuss the speakers' motion ranges and traits, and define their data-driven gesture spaces.

5.1 Full Arm Motion Range

To reveal whether a similar amount of energy is exerted by the left and right arms, we measure the deviation from the mean pose. We independently compute a frame-wise Euclidean distance from each arm to its respective mean pose. These statistics are calculated over all arm joints.

Figure 3 shows the results for the four speakers that were identified as exhibiting the least and most symmetry in their mean pose in Section 4. It is evident that the amount of deviation from the mean pose in the left and right arms is not significantly different if we consider the poses that fall within the whiskers, which represent those within $1.5 \times$ the interquartile range beyond the first and third quartiles. However, the outliers do appear somewhat asymmetrical for speakers Deep3 and Shallow4, each displaying greater divergence from the mean with the right arm compared to the left. Shallow3 and Shallow2 exhibit more symmetrical outliers, indicating that a similar amount of space is encompassed by both arms during these infrequent, larger gestures. The maximum and minimum values for each speaker follow the same trend, with larger maximum values recorded for the right arm in the former two speakers, and similar values for both arms for the latter two.

Figure 4 shows a frontal perspective projection of each speaker's arm pose taken over all of their respective conversations at 1 second intervals. We observe

variability in the gestural symmetry and the amount of gesturing per speaker. Shallow3 appears the most symmetrical with a wide range of positions produced by both arms. Despite having a highly symmetrical mean pose, Shallow2 exhibits a high degree of asymmetry in the peripheral poses whereby the right arm reaches wider poses than the left, but the left arm produces higher gestures than the right. Deep3 and Shallow4 both raise their right hands more frequently than their left, suggesting increased expressiveness in that dominant hand. From these plots it is evident that asymmetry is most apparent in the peripheral gesture space where the extreme gestures are performed. Although relatively infrequent, these extreme gestures capture visual attention and are perceptually significant (McNeill, 2011).

5.2 Gesture Spaces

McNeill defines the *central gesture space* as the area below the neck and between the shoulders and elbows, and the *peripheral gesture space* as any gestures performed outside of the *central gesture space* (McNeill, 2011). Given the variability between the spaces occupied by each speaker's arm and the frequency in which they extend into their respective peripheral spaces, we propose a data-driven approach to defining speaker-specific gesture spaces. We use statistics to define a speaker's *common gesture space* and *extreme gesture space*. The *common gesture space* is the region within a single standard deviation of the respective speaker's mean arm pose. The *extreme gesture space* is the space outside of a single standard deviation of the mean pose, away from the body.

Using our definition, we partition the data into two sections. The *extreme* partition contains all poses with at least one arm in the *extreme gesture space*, and the *common* partition contains the remaining data. We again compute the per-speaker distance from the mean pose for each partition, and the results can be seen in Figure 5. For the majority of the speakers, the distances from the mean for gestures within the common gesture space are similar for both left and right arms (Figure 5, bottom row). An exception is the speaker Deep3 in which the range is larger for the right hand. The greatest differences between the left and right arms are observed in the extreme gesture space (Figure 5, top row), particularly for the asymmetric speakers Deep3 and Shallow4. In each case, one hand diverges further from the mean than the other.

For Deep3, we observe that the left arm is more active in the extreme gesture space than the right, and the reverse is true in the common gesture space. We plot the perspective projection of all poses corresponding to the extreme and common gesture spaces in Figure 6 for each speaker to visualise these differences. The

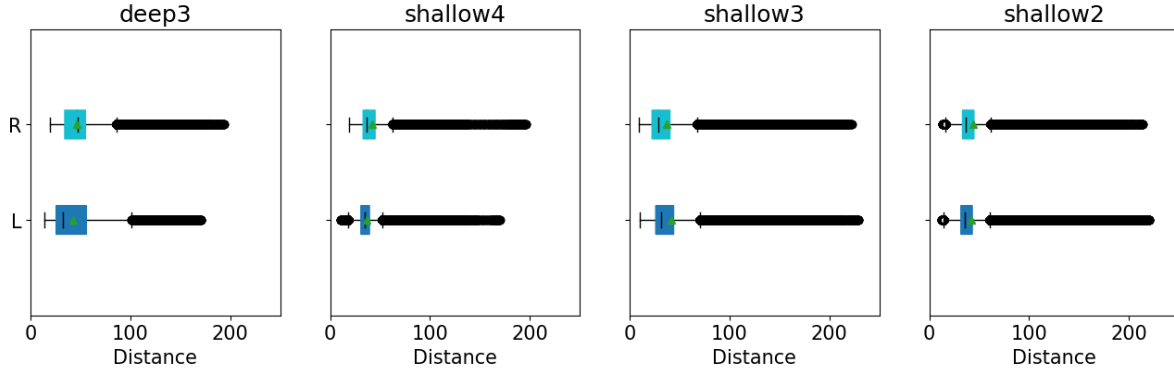


Figure 3: Per-frame Euclidean distance from the mean of each arm for four speakers. L=Left arm, R=Right arm.

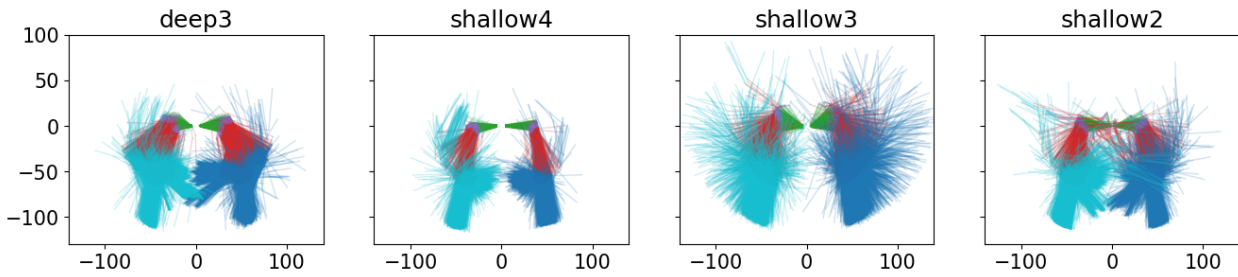


Figure 4: A frontal perspective projection of all poses per speaker, taken at one-second intervals.

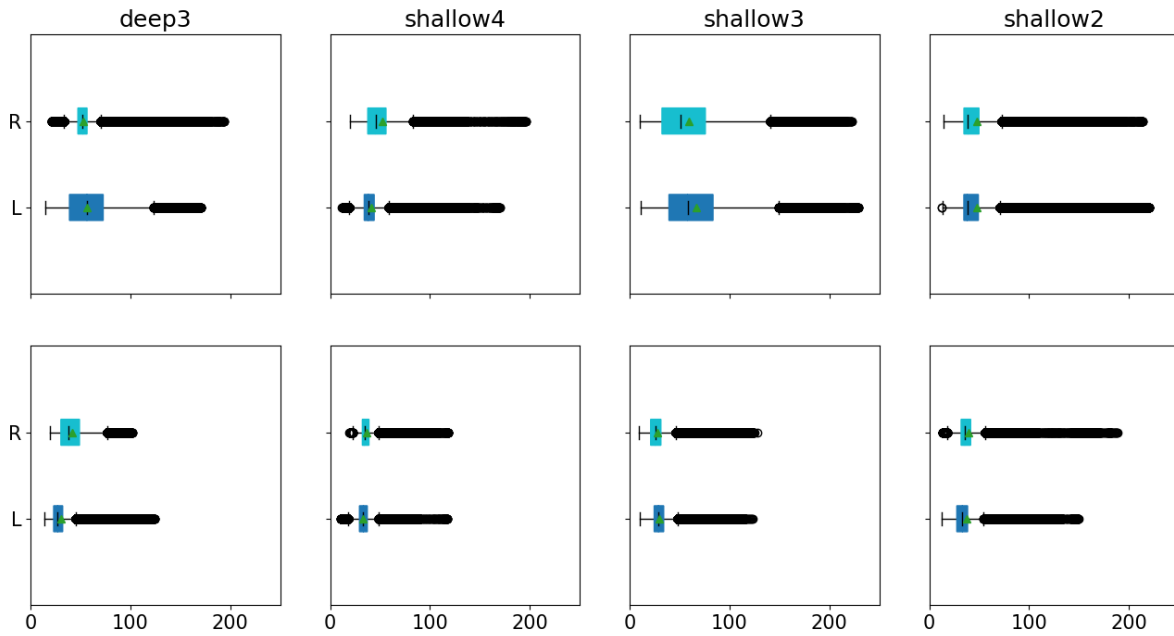


Figure 5: Per-frame Euclidean distance from the mean of each arm, split into *Extreme Gesture Space* (Top) and *Common Gesture Space* (Bottom). L=Left Arm. R=Right Arm.

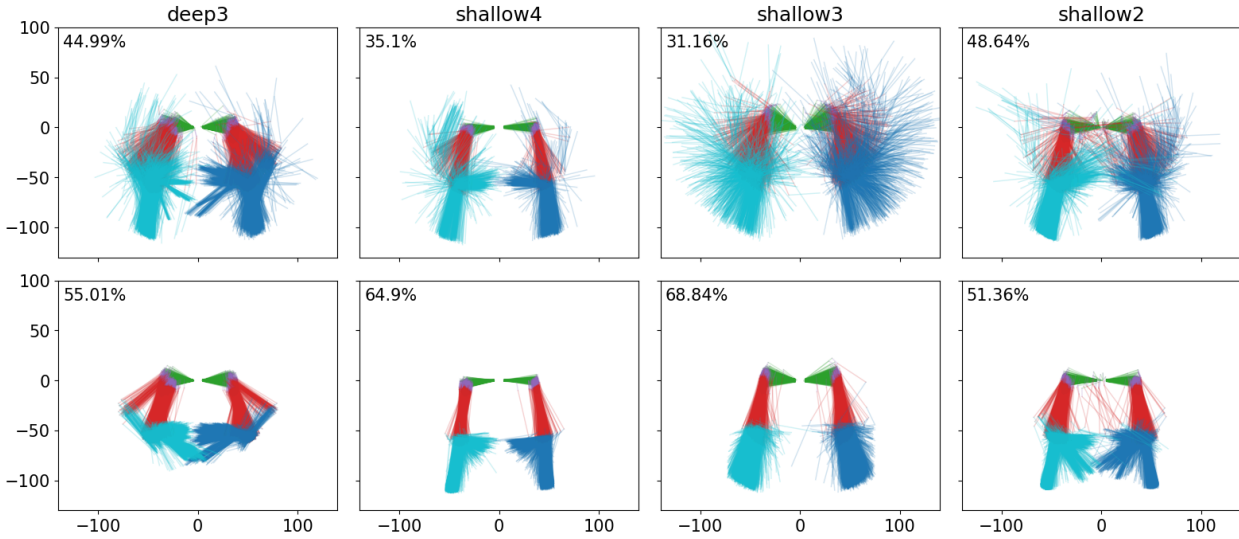


Figure 6: Frontal projections of all poses from four speakers at one-second intervals, split into *Extreme Gesture Space* (Top) and *Common Gesture Space* (Bottom). Percentage in the corner denotes the percentage of poses belonging to the respective gesture space for the respective speaker.

top row reveals that the right arm of Deep3 does contribute to gesturing in the extreme gesture space, but the poses of the left arm are wider, taller and further from the mean pose. In contrast, the bottom row shows more movement in the right arm than the left in the common gesture space, but not significantly.

Figure 6 highlights that the positioning of the arms in common gesture space appears to be more symmetrical than in extreme space across all speakers. Each speaker exhibits different types of asymmetry in the extreme gesture space. Shallow4 lowers their left arm and raises the right and shallow2 extends their right arm wider than the left. Shallow3 has highly mobile arms but holds symmetry in both spaces reasonably well, consistent with the findings in Section 5.1. The percentage of poses within each gesture space as shown in Figure 6 impacts the effect of mirroring. Given more symmetry being found in the common gesture space, if a speaker has a lower use of the extreme gesture space, the potential negative impact of mirroring is reduced.

5.3 Self-adaptor Traits

Self-adaptors are movements that occur simultaneously with speech gesturing, and that typically include self-touch, such as scratching of the neck, clapping at an elbow, adjusting hair or interlocking fingers. These traits tend to be realised asymmetrically.

Figure 7 shows the poses of speaker Shallow25 who frequently touches their left hand to their right forearm. The reverse, right hand touching left forearm, is not present in any of the motion. If laterally mirrored,

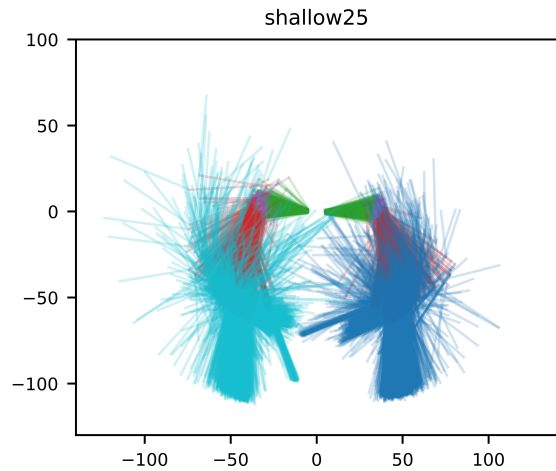


Figure 7: Shallow25 poses taken at one-second intervals. This speaker exhibits self-adaptor movements whereby the left hand frequently touches the right forearm.

474 this self-adaptor movement would not accurately rep-
 475 resent a valid pose from that speaker. The presence and
 476 degree of self-adaptor traits has been found to signifi-
 477 cantly impact the perceived level of neuroticism of a
 478 speaker (Neff et al., 2011), and the effect of reversing the
 479 handedness of the behaviour is not well established.

480 6 Symmetry in Gesture Types

481 When considering the impact of symmetry, the type
 482 of gesture being performed may be important. We re-
 483 viewed a number of speech-motion pairs to determine
 484 what impact may occur from the gesture being mir-
 485 rored. We cannot generalise from these few examples,
 486 but instead should be useful to consider specific aspects
 487 of gesture suitable when mirrored.

488 We observe that beat gestures are often performed
 489 by a single hand. Figure 8 shows a pose plot of a beat
 490 gesture and the values of each wrist position over time.
 491 While the pose plot appears fairly symmetrical with
 492 both arms raised, it is clear that the right arm is moving
 493 up and down, while the left stays fairly static. While we
 494 do not know the dominant hand of this speaker, we ob-
 495 serve some trends similar to those of Çatak et al. (Çatak
 et al., 2018) where one hand is performing the gesture.

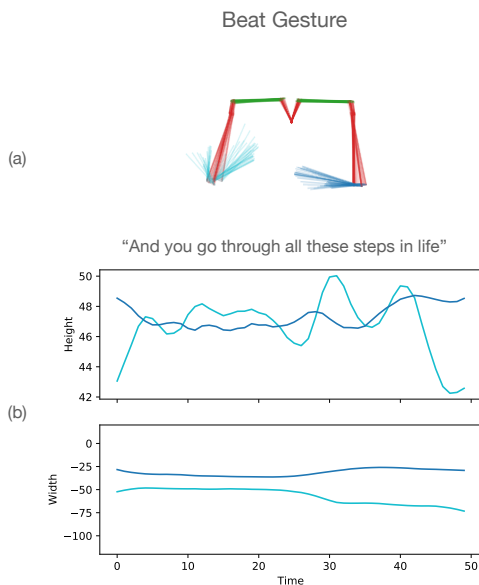


Figure 8: A speaker performing a beat gesture. (a) shows each pose formed over the sequence with the sentence being said below. (b) shows the positions of each wrist in both lateral (left-right) and height (up-down) directions.

496 Çatak et al. (Çatak et al., 2018) suggest that repre-
 497 sentational gestures are performed by a dominant hand

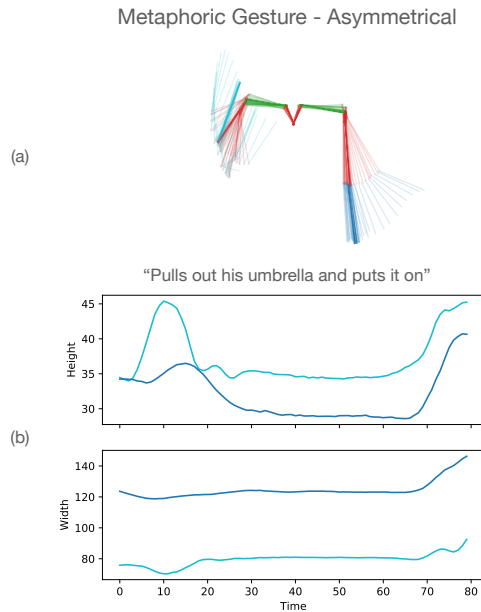


Figure 9: A speaker performing a metaphoric gesture. In this case, the gesture is **asymmetric** due to context. (a) shows each pose formed over the sequence with the sentence being said below. (b) shows the positions of each wrist in both lateral (left-right) and height (up-down) directions.

499 for right-handed speakers but no dominant hand was
 500 found in left-handed speakers. While we cannot com-
 501 pare handedness in this work, we do consider that the
 502 context of the gesture can determine the symmetry of
 503 the gesture performed. Figure 9 shows a metaphoric
 504 gesture being performed, mimicking the use of an um-
 505 brella. It is typical for a person to only use a single hand
 506 while using an umbrella and therefore a single hand is
 507 used to depict this. Should this pose be mirrored, it may
 508 still make logical sense as a single hand will be used
 509 but the handedness of the speaker may not be main-
 510 tained. Figure 10 is a gesture performed by another
 511 speaker, however, they are referring to moving a heavy
 512 object onto a table. Typically moving heavy objects in
 513 the manner outlined in the speech would require two
 514 hands and therefore two hands have been used to de-
 515 pict this. In this instance there are high degrees of sym-
 516 metry between each arm movement, both arms moving
 517 and seemingly at the same or similar time.

518 With regards to directional Deictic gestures, we ob-
 519 served that often the hand closest to the direction was
 520 used. Figure 11 shows a gesture referring to each end of
 521 a building. "That end of the building" is referred to us-
 522 ing the right arm, pointing towards the same direction
 523 to depict an area far away. "this end of the building" is
 524 seemingly the end in which they are stood and a small

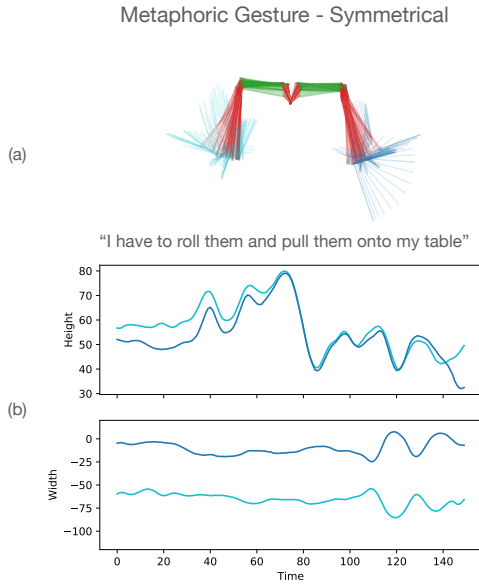


Figure 10: A speaker performing a metaphoric gesture. In this case, the gesture is **symmetric** due to context. (a) shows each pose formed over the sequence with the sentence being said below. (b) shows the positions of each wrist in both lateral (left-right) and height (up-down) directions.

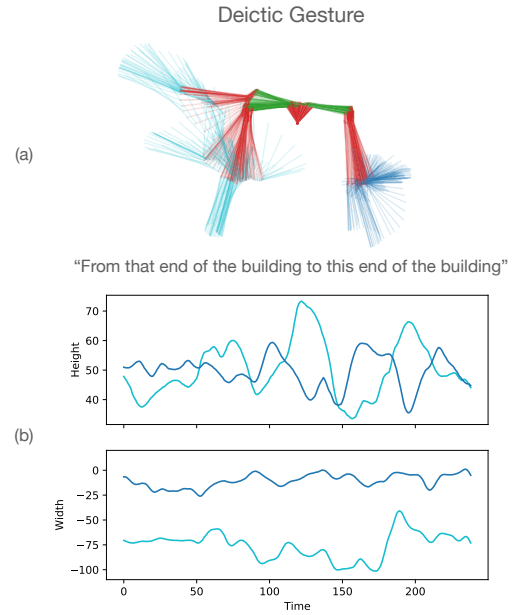


Figure 11: A speaker performing a Deictic gesture. (a) shows each pose formed over the sequence with the sentence being said below. (b) shows the positions of each wrist in both lateral (left-right) and height (up-down) directions.

525 movement of the left arm is used to refer to this. Figure
526 11 time plot shows a clear spike as the right arm moves
527 to the peak directional gesture, the left arm is lowering,
528 suggesting asymmetry.

529 We describe some examples of symmetrical and
530 asymmetrical poses and their associated gesture type.
531 We find that in some cases a mirrored, symmetrical
532 pose may well still portray the same meaning. A good
533 example of this is when a metaphoric action requires
534 the use of both hands to lift something. However, in
535 the example Deictic gesture this would not continue to
536 make sense when performed in the same location.

537 7 Mirrored Pose Validity

538 For some machine learning approaches, the goal of lat-
539 erally mirroring body pose is to generate further, valid
540 examples of the same speaker. In these cases, validity
541 only holds if the mirrored poses fall within the gesture
542 space of the original data belonging to that speaker. In
543 this section we visualise and quantify mirrored pose va-
544 lidity using this definition.

545 We perform a nearest neighbour search of each mir-
546 rored pose in the original motion data per speaker. The
547 distance metric used is the Euclidean distance which is
548 computed over the joint locations in both arms. We fo-

549 cus on the poses that fall within the extreme gesture
550 space, defined as any pose outside of one standard de-
551 viation away from the mean pose (Section 5.2). We first
552 present a visualisation of the nearest neighbours in Fig-
553 ure 12. In this plot the top row shows a subset of the
554 mirrored poses for each speaker, and the bottom row
555 shows the nearest neighbours from the original motion
556 data. It is evident from this figure that it is not possible
557 to cover the full range of motion found in the mirrored
558 poses in the original data. For each speaker there are
559 areas in world space for which the arm does not reach
560 in the original data.

561 In the rightmost column of Figure 12 we observe
562 that, with speaker Shallow2, for the left arm to reach
563 out as wide as it does in the mirrored poses, in the orig-
564 inal data, the right arm also has to extend. This sug-
565 gests that in the original data, it is characteristic for
566 either both arms to move to a wide position together,
567 or for the right arm to move out wide independently.
568 It is uncharacteristic for the left arm to reach out in-
569 dependently from the right arm. For both Deep3 and
570 Shallow4 (leftmost columns), when the mirrored poses
571 are at their most extreme poses (i.e. the arms elevated
572 to their highest and widest positions), it is not possible
573 to match these in the original data.

574 Figure 13 shows mean distances between the mir-
575 rored poses and the closest match in the original data.

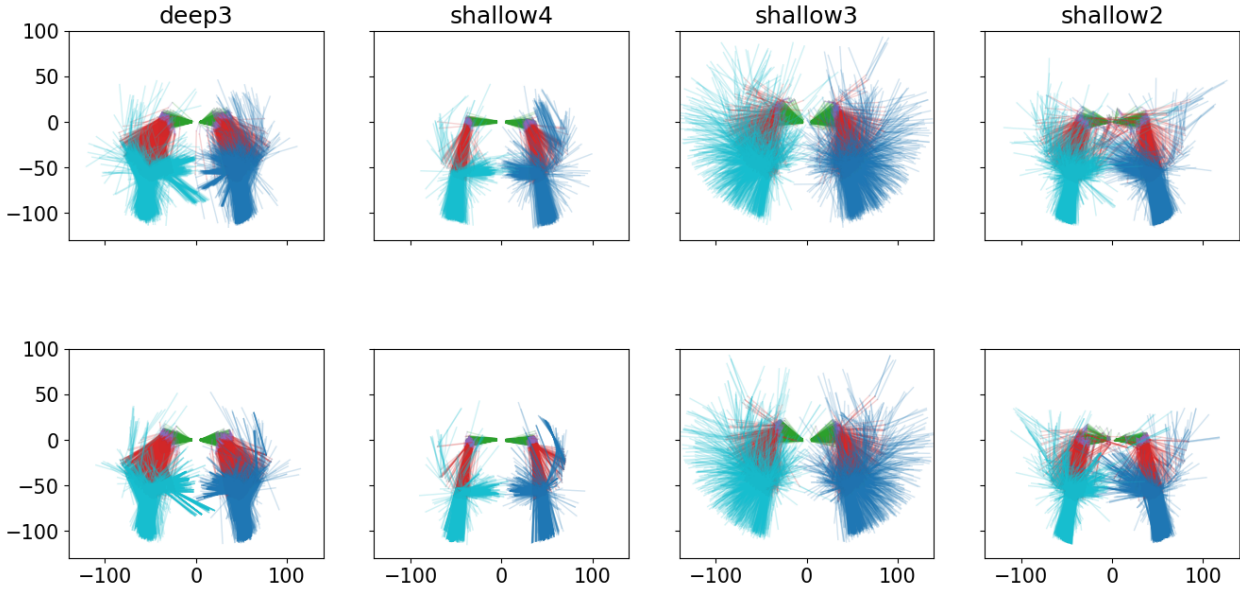


Figure 12: The frontal 2D projections of mirrored poses that are at least 1 standard deviation away from their mean pose (top) and the closest respective mean poses from the original data (bottom).

Although Deep3 was associated with the least symmetrical mean pose from the dataset (Section 5), we observe that, in the extreme gesture space, they produce similar gestures with both left and right hands.

8 Temporal Symmetry

Our analysis so far has considered only frame-wise statistics, which does not account for differences in the dynamics of each arm. Lateral mirroring for body data augmentation swaps the positions of the arms on a frame-by-frame basis, so the dynamics of the respective arms are inherently swapped. In practice, there may exist an asynchrony, or a temporal shift, between the motion of the two arms, particularly if the speaker gestures with a dominant hand. In this section we perform a cross-correlation analysis to reveal any temporal lag between left and right hands.

Correlation between the left and right hand positions is computed over a 401-frame window ($\approx 4.5s$), centred at frame t . For each windowed frame in the left hand data, $t = 0, \dots, T$, we slide the window over the right arm data from frames $t - 200$ to $t + 200$ and compute the correlation coefficient between the segments. A larger window size was not used since we observed that a lag longer than 2.2 seconds was more commonly due to a rhythmic motion than an asynchrony caused by a leading hand. The cross-correlation analysis is performed for each motion sequence on a per-speaker basis. We independently run the analysis on each direc-

tional axis and the Euclidean distance to the mean pose of each hand, and the results can be seen in Figure 14.

Although Shallow2 has a relatively symmetrical gesture space (Figure 4), Figure 14 clearly shows a dominant hand in the temporal domain. This indicates that this speaker leads with their right hand with a mean offset of 28 frames ($\approx 0.31s$) when considering the distance from the mean pose. If we consider the individual axes, we observe that the right hand leads in all cases, and in the X and Y axes the offset is greater than 0.5s. This suggests that, although a symmetrical pose is formed, there is a temporal offset between hands achieving this pose.

It is evident that other speakers' motions are more symmetrical and very small temporal offsets were found. Shallow3 in Figure 14 is an example where the mean offset does not exceed a mean of 17 frames (0.19s) in any axis.

9 Mutual Information

In this Section we explore mirroring for data augmentation from an information theory perspective. Specifically, whilst mirroring effectively doubles the amount of data, how much additional *information* does it introduce? We compute the mutual information between the original data and its mirrored counterpart to reveal the dependence between the two distributions.

We measure Normalised Mutual Information (NMI) (Strehl and Ghosh, 2002) on a per-speaker, per-axis ba-

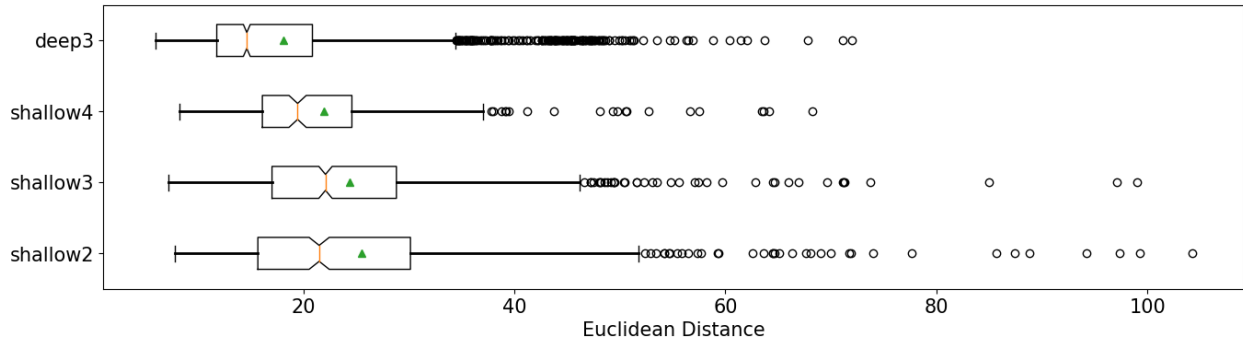


Figure 13: Euclidean distance between mirrored arm position and the closest pose from the original data for poses in the extreme gesture space.

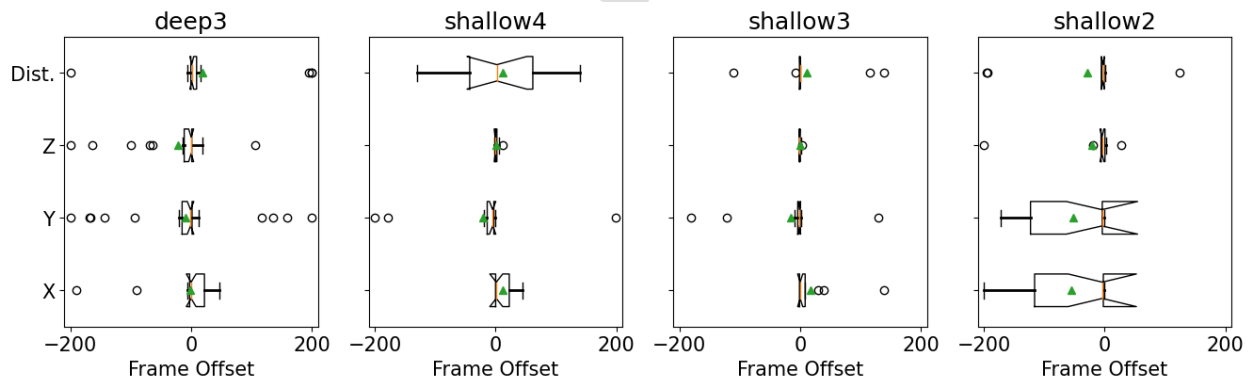


Figure 14: Cross correlation analysis between left and right hand position for each directional axis and Euclidean distance from the mean.

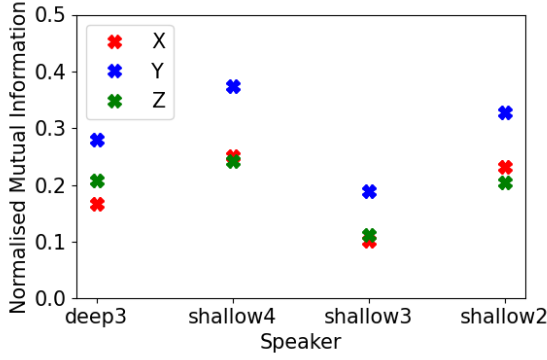


Figure 15: Normalised Mutual Information per-speaker, per-axis measured between the original and mirrored wrist joints. Lower values represent a higher degree of independence.

sis at the wrist joint. NMI is computed using the following:

$$NMI(X, \tilde{X}) = \frac{I(X, \tilde{X})}{\sqrt{H(X)H(\tilde{X})}} \quad (1)$$

where $I(X, \tilde{X})$ is the mutual information between the original and mirrored data, and $H(X)$ and $H(\tilde{X})$ is the entropy of the original and mirrored data respectively. The entropy is calculated using the nearest neighbour approach (Kozachenko and Leonenko, 1987).

Normalising the Mutual Information allows for easy comparison between speakers and axis, producing a value between 0-1. This NMI value describes the dependence of the two variables. At zero NMI, the variables are completely independent, and as the NMI increases to 1, it indicates a reduction in uncertainty and largely dependent variables.

The NMI for each speaker is shown in Figure 15. This shows that the amount of mutual information in the wrists is speaker-dependent. However, when considering the relative mutual information between axes, the Y-axis (movement of the wrist in the vertical axis) consistently has higher values. Therefore, our analysis suggests that more information will be gained in the movement along the X-axis (forward-back) and the Z-axis (left-right) from augmenting the dataset with mirrored poses. Information symmetry is revealed from NMI. Low levels of NMI and therefore, low information symmetry indicates the importance of both wrists to predictive models. This is particularly important when regarding motion datasets gathered from video. As occlusion is common, arms are often interpolated or missing from the data. By removing or including potentially incorrect arm movement on one side, you are losing important information or introducing large amounts of

uncharacteristic information.

10 Generative Modelling

To further support our findings, we train a Long Short-Term Memory (LSTM) model on different splits of data and use various augmentation settings to map from speech to body pose. We aim to determine the impact of including the potentially uncharacteristic mirrored motion for a speaker and whether including the mirrored speaker as a new *virtual identity* improve results.

10.1 Motion Representation

Of the 36 speakers released, only 18 have both audio and motion capture available and therefore we use this subset. Mocap was down-sampled to 30fps to ensure realistic motion was maintained, but training time was reduced. A test sequence is randomly held out for each speaker and the remaining data, 20% is held out for validation and 80% is used for training. The global position for each speaker is inconsistent and therefore, the respective mean global root position is removed from each frame on a per-sequence basis. 3D positions in world space are the target values which are standardised by subtracting the mean pose and dividing by the standard deviation computed over all speakers across all training sequences.

10.2 Audio Representation

Mel Spectrograms or Mel Frequency Cepstral Coefficients (MFCCs) are often used in speech-to-motion pipelines (Habibie et al., 2021; Alexanderson et al., 2020a; Taylor et al., 2021). We instead use a model trained using a multi-task learning framework that is comprised of 12 regression tasks. (PASE+) (Ravanelli et al., 2020) features encode an audio waveform and should implicitly encode MFCCs and other speech-related information, including prosody and speech content. Speech is downsampled using a band-sinc filtering method from 44.1KHz to 16KHz.

10.3 Generative Model

Using an LSTM-based model, we train using a single motion frame's worth of audio (33ms) to predict a frame of motion. To ensure motion is speaker-specific, we condition the speech using a learned feature vector that encodes a speaker's identity. This learned feature vector should adequately associate the speaker and their gesturing style. With this learned feature vector, it should allow us to introduce a speaker's potentially uncharacteristic mirrored motion to the model, without affecting the gesturing style of the speaker.

The LSTM model contains 4 bi-directional layers, each with 1024 hidden units and a 40% dropout followed by a ReLU non-linearity layer and a fully connected layer. The output from the fully connected layer is the estimated (standardised) body pose at that frame.

10.4 Training Procedure

Models are trained using the Adam optimiser with a learning rate of 0.0001 and batch size of 256. Not all sequences contain hand motion, where this is the case, we compute the loss against all joints in the body except the hands. We use 30-frame long sequences to train, with a 25-second overlap on each window.

We use a multi-term loss function. We minimise the position values as an L_2 loss on joint positions and also an L_2 loss on joint velocity and acceleration. Introducing the velocity and acceleration allows the model to produce smoother and more realistic transitions. On observation of some bone stretching artefacts due to positions not having any constraint on distance apart, we include an L_1 loss on bone length. The final loss L_c is computed as:

$$\begin{aligned} L_p &= L_2(y, \hat{y}) \\ L_v &= L_2(f'(y), f'(\hat{y})) \\ L_a &= L_2(f''(y), f''(\hat{y})) \\ L_b &= L_1(y_{lengths}, \hat{y}_{lengths}) \\ L_c &= L_p + L_v + L_a + L_b \end{aligned} \quad (2)$$

where y and \hat{y} is the ground truth and predicted motion, and $y_{lengths}$ and $\hat{y}_{lengths}$ are Euclidean distances between each joint and its parent in the skeleton hierarchy for the ground truth and predicted motion respectively. The term L_p is representative of positional accuracy, L_v velocity accuracy, L_a acceleration accuracy, L_b bone length accuracy and L_c is the combined loss. L_1 and L_2 represent Mean Absolute Error and Mean Squared Error respectively.

10.5 Experimental Setup

We train the same model architecture on each of the settings defined as follows:

All Data. We form a baseline using all available training data with no augmentation.

Half Data. A random subsample of the training data reduces the number of samples by approximately 50%. We train a model using this reduced data to enable us to compare the effect of doubling the size of the training set by augmentation versus adding additional ground truth data.

Mirrored Same Identity. We augment the *Half Data* training set by laterally mirroring the pose at each

frame. Mirrored data is assigned the **same** identity label as the original speaker. This allows us to determine the impact of introducing uncharacteristic motion for a specific speaker.

Mirrored Virtual Identity. We augment the *Half Data* training set by laterally mirroring the pose at each frame. During training, we assign a **new** virtual identity label to the mirrored data. This allows us to determine if adding motion that could be considered characteristic for a different speaker aids or hinders performance.

All Data Mirrored Virtual Identity We additionally train our model on all available training data plus the laterally mirrored augmentation. As in the *Mirrored Virtual Identity* setting, the augmented sequences are assigned new virtual identity labels. This represents our optimal setting.

10.6 Results

We continue to use motion characteristics to evaluate performance. These include positional pose plots, distances from the mean pose and temporal handedness. Our analysis should provide an indication of how characteristic the predicted motion is and whether the introduction of motion has had an impact on performance. We follow the same procedure as in Section 3 and translate per frame so that the midpoint of the left and right shoulders and centred on the origin.

10.6.1 Using the same identity

We observe two key findings; the mirrored data produced far more muted and symmetrical motion than desired.

We found the movement generated to be positionally symmetrical over the whole pose but particularly with arm movements. Figure 16a shows each of the arms consistently raising simultaneously when using mirrored data as the same identity. While using just half of the data and no mirror augmentation, there are more asymmetrical poses which are closer to the characteristics performed in the ground truth.

Figure 16b indicates the amount of time and distance away from the mean pose. It is a common trend across speakers that the distance from the mean pose was lower in the mirrored with the same identity split when compared to motion generated from half of the data and the ground truth. This is indicative of the muted motion observed, producing slow and small movements.

Temporal symmetry is notably present when using the same identity. When the left-hand moves, the right hand also moves at the same time producing unnatural motion. Figure 16 shows a strong correlation between the left and right wrists moving at a temporal lag offset

of ± 1 frame. When compared to the ground truth, this high temporal symmetry is very uncharacteristic of the speaker.

10.6.2 Augmenting With a Virtual Identity

With a detrimental effect of including mirrored data under the same identity, we examine the effects of including mirrored data under a virtual identity (*Mirrored Virtual Identity*).

We identified improvements in generated motion quality varied between speakers, however, we did not find a negative impact on performance. Mirroring with a virtual identity was found to be competitive with a model trained with all of the available data, often improving positioning, adding some more movement that closely resembles the ground truth and generating motion from all of the data.

An example of improvement from including lateral mirrored data is shown in Figure 17. The distribution of distances from the mean pose shown in Figure 17b decreases from half of the data and half mirrored as a virtual identity. We also note the poses in Figure 17a appear closer to the predictions using all of the data and ground truth. By seemingly lowering the arms more often than the generated motion using half of the data, this supports the hypothesis that the addition of mirrored data as a virtual identity can be competitive with a model including all data.

11 Discussion

We discuss our findings on arm symmetry during dyadic conversation and its impact on lateral mirroring for body motion data augmentation. We present the potential issues that could arise, and when it would and would not be a suitable data augmentation approach.

If lateral mirroring is used for body data augmentation, caution should be taken if gesturing style and handedness of the speaker are to be preserved. From our analysis it is clear that mirroring can result in both valid poses and dynamics for certain speakers who move with a high degree of arm symmetry. Statistical analysis can be performed on a per-speaker basis to ensure that this is the case. However, for these highly symmetrical speakers, the information gained from mirroring the arm motion might be minimal. In the majority of cases, the speakers did not move symmetrically, and the mirrored data would not reflect the true characteristics of a speaker's gesturing style. While mirroring could produce a physically valid pose for a speaker, it may not fit with their motion style or handedness.

From our generative modelling, a naive mirroring implementation did not predict characteristic or plau-

sible motion and was found to be detrimental to model performance. We instead suggest the use of a new *virtual identity* for the mirrored poses. We found that the amount of improvement was speaker-dependent. We speculate this may be due to the non-uniform distribution of data across the speakers. As the dataset used has *shallow* and *deep* speakers, the amount of data available per speaker varies. Although the models appeared to capture the speaker identities well, there is a chance that with small amounts of data for some speakers, the motion characteristics required to describe this speaker's motion are simply not present in the training data. We speculate the improvement may be due to an increase in generalised characteristics common across all speakers. If the aim is to preserve the gesturing style and handedness of the original speaker, lateral mirroring should instead be used to increase the number of speakers in a dataset by treating the mirrored data as its own *virtual identity*. Care must still be taken to account for directional cues in the training data speech that could lead to a multi-modal disparity.

Shallow₂₅ in Figure 7 is an example of an asymmetrical self-adaptor trait that is characteristic to that speaker. The left arm touching the right arm is common in their data, but the right arm does not appear to touch the left arm in the same manner. If this stylistic motion was to be maintained, simply mirroring the body pose would not suffice.

Mirroring the data has the potential to cancel out temporal offset characteristics. We have observed that certain speakers gesture with a leading hand. We found a generative model that has been trained on both the original and augmented motion data with the same identity removes any temporal offsets and produces temporally symmetrical motion. This synthetic motion would not be faithful to the original speaker.

Given the speaker-dependent nature of the amount of symmetry, we expect the inclusion of a symmetry statistic to aid in numerous tasks. We discuss the use of statistics for synthetic motion evaluation in Section 11.1, however, we also suggest considering the use of these statistics for identity classification. Motion symmetry could be important to the classification of speaker identity. We expect that a discriminatory model (i.e. "Does this motion resemble the expected speaker?") could be successful when classifying using symmetry motion characteristics. More work is required to determine what degree of success classifying a speaker's identity using motion symmetry alone could provide.

The mutual dependence between the mirrored poses and original is speaker dependent, and we observe that some information is gained through lateral mirroring. *More information* may be enticing, however, this measure does not inform on appropriateness, and

the added information may introduce uncharacteristic motions.

Previous work by Çatak et al. (Çatak et al., 2018) has considered the impact of handedness on beat and representational gestures. They found that beat gestures had a preference for the dominant hand of the speaker, whereas representational gestures varied. In left-handed speakers, there was no preference, but in right-handed speakers, there was a right-handed preference. This suggests that, although arm positions could be reflectively similar, the types of gesturing could be varied. When training a generative body motion model using mirrored motion, there is a risk that both hands will produce beat gestures in the synthesised animation, which may reduce realism or even understanding.

We analysed a few gesture types and their relationship to symmetry. While we cannot generalise from this small analysis, it would be sensible to consider when certain gesture types could be adequately mirrored. It is essential that handedness is maintained during directional or positional gestures, such as pointing to communicate a direction. If a speaker uses a gesture to signify to the left and the augmented version points to the right with no adaptation of the corresponding audio speech, this would lead to a disparity in the multimodal context. When building gesture-generation systems, it would be beneficial to keep the handedness of gestures produced consistent.

Further study is required to determine the impact of modifying positional and temporal symmetry on realism and understanding. However, our findings suggest that care should be taken when augmenting data using lateral mirroring. There is a risk that with this augmented data the motion could lose speaker-dependent characteristics.

11.1 Evaluating Synthetic Motion

A significant challenge in data-driven synthesis of embodied agents is how to evaluate the synthesised body animation. It is common to evaluate performance of generative models by means of a user study (Alexanderson et al., 2020a). Assuming the synthesised data is to represent that of a particular speaker, the analysis from this study could also be considered as a performance evaluation method.

If the goal is to generate animated body motion that is faithful to the style of a particular speaker, we would expect the animation to possess the same positional and temporal characteristics as the speaker's ground truth motion. We propose that statistical analyses based on the work presented in this paper would provide good indicators of these qualities. The per-speaker percentage of time spent in the extreme gesture space, degree of spatial symmetry and temporal lag of

the animated result compared to the ground truth motion would be indicative of the similarities in both gesturing style and handedness.

12 Conclusion

We have studied four subjects from the *Talking with Hands* dataset to examine the symmetry in arm motion during dyadic conversation. We found that motion symmetry is highly speaker dependent. We derived a data-driven approach for defining a per-speaker gesture space, and found that the arms exhibited more lateral symmetry when in the common gesture space (closer to the mean pose) than when in the extreme space (further from the mean). We discovered that some speakers gesture with a leading hand, and others maintain left-right temporal alignment. We used information theory to find there is a large amount of information to be gained from both wrists. We employed a speech-to-motion model to support our findings.

Using these findings we have determined the efficacy of lateral mirroring for data augmentation and the considerations that should be made. If the goal is to maintain a speaker's gesturing style and handedness, mirroring for generating further examples of that speaker can only be used in certain cases, and is not suitable as a generic data augmentation approach. However, we suggest it is suitable for increasing the number of speakers in the training set by treating the mirrored data as a new *virtual identity*.

Finally, we propose our statistical analysis for evaluating the performance of speech-driven conversational agents to ensure that speaker characteristics have been retained in the synthesised motion.

Acknowledgements

Sarah Taylor was supported by the Engineering and Physical Research Council (Grant number EP/S001816/1).

References

- Alexanderson, Simon, Gustav Eje Henter, Taras Kucherenko, and Jonas Beskow. 2020a. Style-controllable speech-driven gesture synthesis using normalising flows. In *Computer Graphics Forum*, volume 39, pages 487–496. Wiley Online Library.
- Alexanderson, Simon, Éva Székely, Gustav Eje Henter, Taras Kucherenko, and Jonas Beskow. 2020b. Generating coherent spontaneous speech and gesture from text. In *Proceedings of the 20th ACM International Conference on Intelligent Virtual Agents*, pages 1–3.

- 1013 Bagautdinov, Timur, Chenglei Wu, Tomas Simon,
1014 Fabián Prada, Takaaki Shiratori, Shih-En Wei,
1015 Weipeng Xu, Yaser Sheikh, and Jason Saragih. 2021.
1016 Driving-signal aware full-body avatars. *ACM Trans.*
1017 *Graph.*, 40(4).
- 1018 Çatak, Esra Nur, Alper Açık, and Tilbe Göksun. 2018.
1019 The relationship between handedness and valence:
1020 A gesture study. *Quarterly Journal of Experimental*
1021 *Psychology*, 71(12):2615–2626.
- 1022 De Ruiter, Jan P, Adrian Bangerter, and Paula Dings.
1023 2012. The interplay between gesture and speech in
1024 the production of referring expressions: Investigating
1025 the tradeoff hypothesis. *Topics in Cognitive Science*,
1026 4(2):232–248.
- 1027 Drijvers, Linda, Asli Özyürek, and Ole Jensen. 2018.
1028 Hearing and seeing meaning in noise: Alpha, beta,
1029 and gamma oscillations predict gestural enhance-
1030 ment of degraded speech comprehension. *Human*
1031 *Brain Mapping*, 39(5):2075–2087.
- 1032 Ferstl, Ylva and Rachel McDonnell. 2018. Investigat-
1033 ing the use of recurrent motion modelling for speech
1034 gesture generation. In *Proceedings of the 18th Interna-*
1035 *tional Conference on Intelligent Virtual Agents*, pages
1036 93–98.
- 1037 Ginosar, Shiry, Amir Bar, Gefen Kohavi, Caroline Chan,
1038 Andrew Owens, and Jitendra Malik. 2019. Learning
1039 individual styles of conversational gesture. In *Pro-*
1040 *ceedings of the IEEE/CVF Conference on Computer Vi-*
1041 *sion and Pattern Recognition*, pages 3497–3506.
- 1042 Gong, Kehong, Jianfeng Zhang, and Jiashi Feng. 2021.
1043 Poseaug: A differentiable pose augmentation frame-
1044 work for 3d human pose estimation. In *Proceedings*
1045 *of the IEEE/CVF Conference on Computer Vision and*
1046 *Pattern Recognition*, pages 8575–8584.
- 1047 Habibie, Ikhsanul, Weipeng Xu, Dushyant Mehta,
1048 Lingjie Liu, Hans-Peter Seidel, Gerard Pons-Moll,
1049 Mohamed Elgharib, and Christian Theobalt. 2021.
1050 Learning speech-driven 3d conversational gestures
1051 from video. In *ACM International Conference on In-*
1052 *telligent Virtual Agents (IVA)*.
- 1053 Henter, Gustav Eje, Simon Alexanderson, and Jonas
1054 Beskow. 2020. Moglow: Probabilistic and control-
1055 lable motion synthesis using normalising flows. *ACM*
1056 *Transactions on Graphics (TOG)*, 39(6):1–14.
- 1057 Hostetter, Autumn B. 2011. When do gestures com-
1058 municate? a meta-analysis. *Psychological bulletin*,
1059 137(2):297.
- 1060 Joo, Hanbyul, Hao Liu, Lei Tan, Lin Gui, Bart Nabbe,
1061 Iain Matthews, Takeo Kanade, Shohei Nobuhara, and
1062 Yaser Sheikh. 2015. Panoptic studio: A massively
1063 multiview system for social motion capture. In *The*
1064 *IEEE International Conference on Computer Vision*
1065 *(ICCV)*.
- 1066 Kendon, Adam. 1994. Do gestures communicate? a
1067 review. *Research on language and social interaction*,
1068 27(3):175–200.
- 1069 Kipp, Michael and Jean-Claude Martin. 2009. Gesture
1070 and emotion: Can basic gestural form features dis-
1071 criminate emotions? In *2009 3rd international confer-*
1072 *ence on affective computing and intelligent interaction*
1073 *and workshops*, pages 1–8. IEEE.
- 1074 Korzun, Vladislav, Ilya Dimov, and Andrey Zharkov.
1075 2020. The finemotion entry to the genea challenge
1076 2020.
- 1077 Kozachenko, LF and Nikolai N Leonenko. 1987. Sample
1078 estimate of the entropy of a random vector. *Problemy*
1079 *Peredachi Informatsii*, 23(2):9–16.
- 1080 Lee, Gilwoo, Zhiwei Deng, Shugao Ma, Takaaki Shira-
1081 tori, Siddhartha S Srinivasa, and Yaser Sheikh. 2019.
1082 Talking with hands 16.2 m: A large-scale dataset of
1083 synchronized body-finger motion and audio for con-
1084 versational motion analysis and synthesis. In *Pro-*
1085 *ceedings of the IEEE/CVF International Conference on*
1086 *Computer Vision*, pages 763–772.
- 1087 McNeill, David. 1985. So you think gestures are non-
1088 verbal? *Psychological review*, 92(3):350.
- 1089 McNeill, David. 2008. *Gesture and thought*. University
1090 of Chicago press.
- 1091 McNeill, David. 2011. *Hand and mind*. De Gruyter Mou-
1092 ton.
- 1093 Neff, Michael, Nicholas Toothman, Robeson Bowmani,
1094 Jean E Fox Tree, and Marilyn A Walker. 2011. Don't
1095 scratch! self-adaptors reflect emotional stability. In
1096 *International Workshop on Intelligent Virtual Agents*,
1097 pages 398–411. Springer.
- 1098 Perez, Luis and Jason Wang. 2017. The effectiveness of
1099 data augmentation in image classification using deep
1100 learning. *arXiv preprint arXiv:1712.04621*.
- 1101 Pouw, Wim, Steven J Harrison, Núria Esteve-Gibert,
1102 and James A Dixon. 2020. Energy flows in gesture-
1103 speech physics: The respiratory-vocal system and
1104 its coupling with hand gestures. *The Journal of the*
1105 *Acoustical Society of America*, 148(3):1231–1247.
- 1106 Ravanelli, Mirco, Jianyuan Zhong, Santiago Pascual,
1107 Pawel Swietojanski, Joao Monteiro, Jan Trmal, and
1108 Yoshua Bengio. 2020. Multi-task self-supervised
1109 learning for robust speech recognition. In *ICASSP*

- 1110 2020-2020 IEEE International Conference on Acoustics,
1111 Speech and Signal Processing (ICASSP), pages 6989–
1112 6993. IEEE.
- 1113 Shorten, Connor and Taghi M Khoshgoftaar. 2019. A
1114 survey on image data augmentation for deep learn-
1115 ing. *Journal of Big Data*, 6(1):1–48.
- 1116 Strehl, Alexander and Joydeep Ghosh. 2002. Cluster
1117 ensembles—a knowledge reuse framework for com-
1118 bining multiple partitions. *Journal of machine learn-
1119 ing research*, 3(Dec):583–617.
- 1120 Studdert-Kennedy, Michael. 1994. Hand and mind:
1121 What gestures reveal about thought. *Language and
1122 Speech*, 37(2):203–209.
- 1123 Taylor, Sarah, Jonathan Windle, David Greenwood, and
1124 Iain Matthews. 2021. Speech-driven conversational
1125 agents using conditional flow-vaes. In *European Con-
1126 ference on Visual Media Production*, pages 1–9.
- 1127 Wagner, Petra, Zofia Malisz, and Stefan Kopp. 2014.
1128 Gesture and speech in interaction: An overview.
- 1129 Yoon, Youngwoo, Bok Cha, Joo-Haeng Lee, Minsu Jang,
1130 Jaeyeon Lee, Jaehong Kim, and Geehyuk Lee. 2020.
1131 Speech gesture generation from the trimodal context
1132 of text, audio, and speaker identity. *ACM Transac-
1133 tions on Graphics (TOG)*, 39(6):1–16.
-

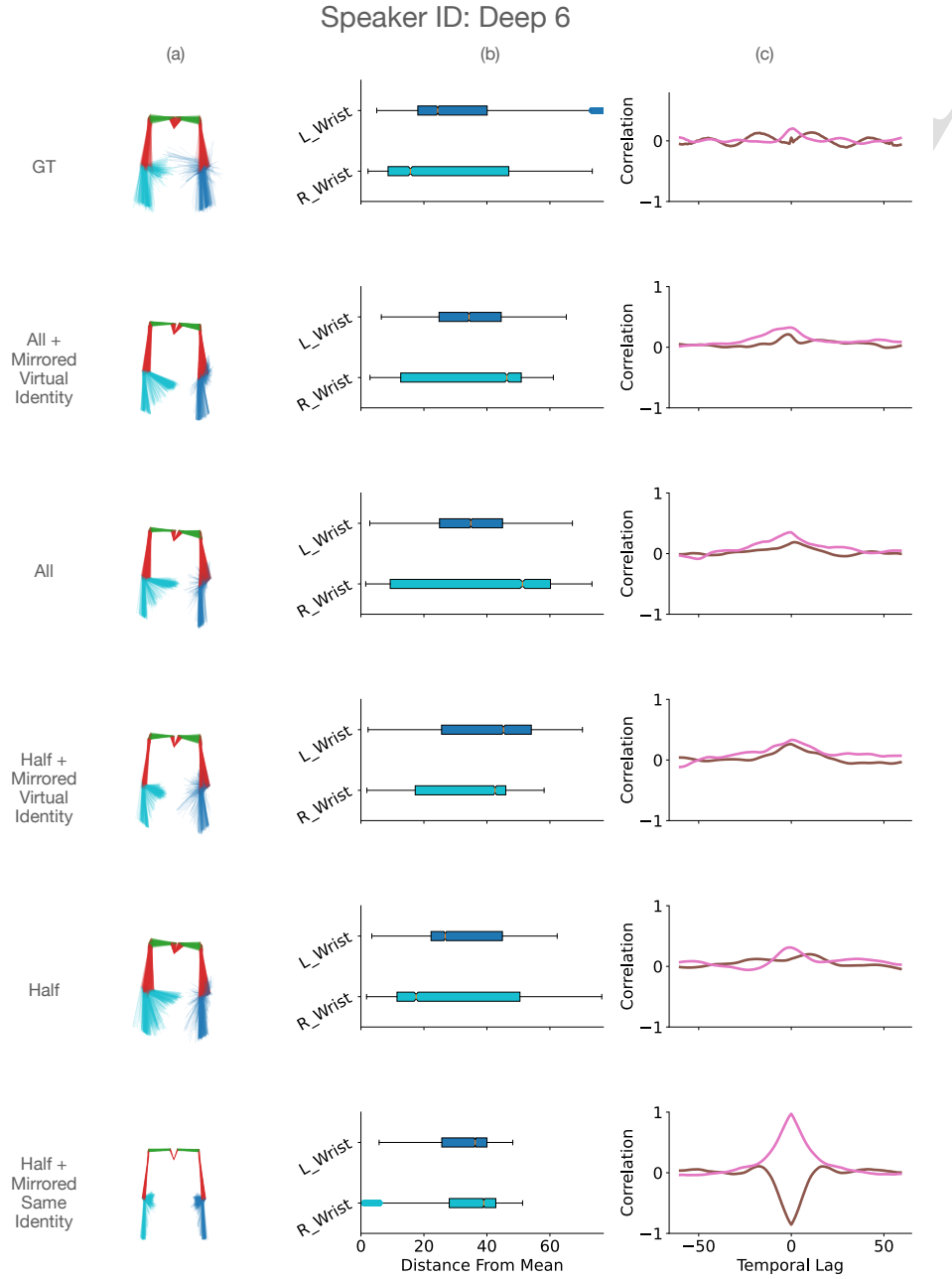


Figure 16: A comparison for a single speaker's generated motion showing detrimental impact of including mirrored motion under the same identity. Each row corresponds to a different data split used. Column (a) contains the orthographic projection of a pose at every second in the sequence. Column (b) shows the distribution of distances from the mean arm pose. Column (c) shows the cross correlation lags between the onset of left wrist motion given right wrist motion in the Z (left-right) and Y (up-down) shown in **brown** and **pink** respectively

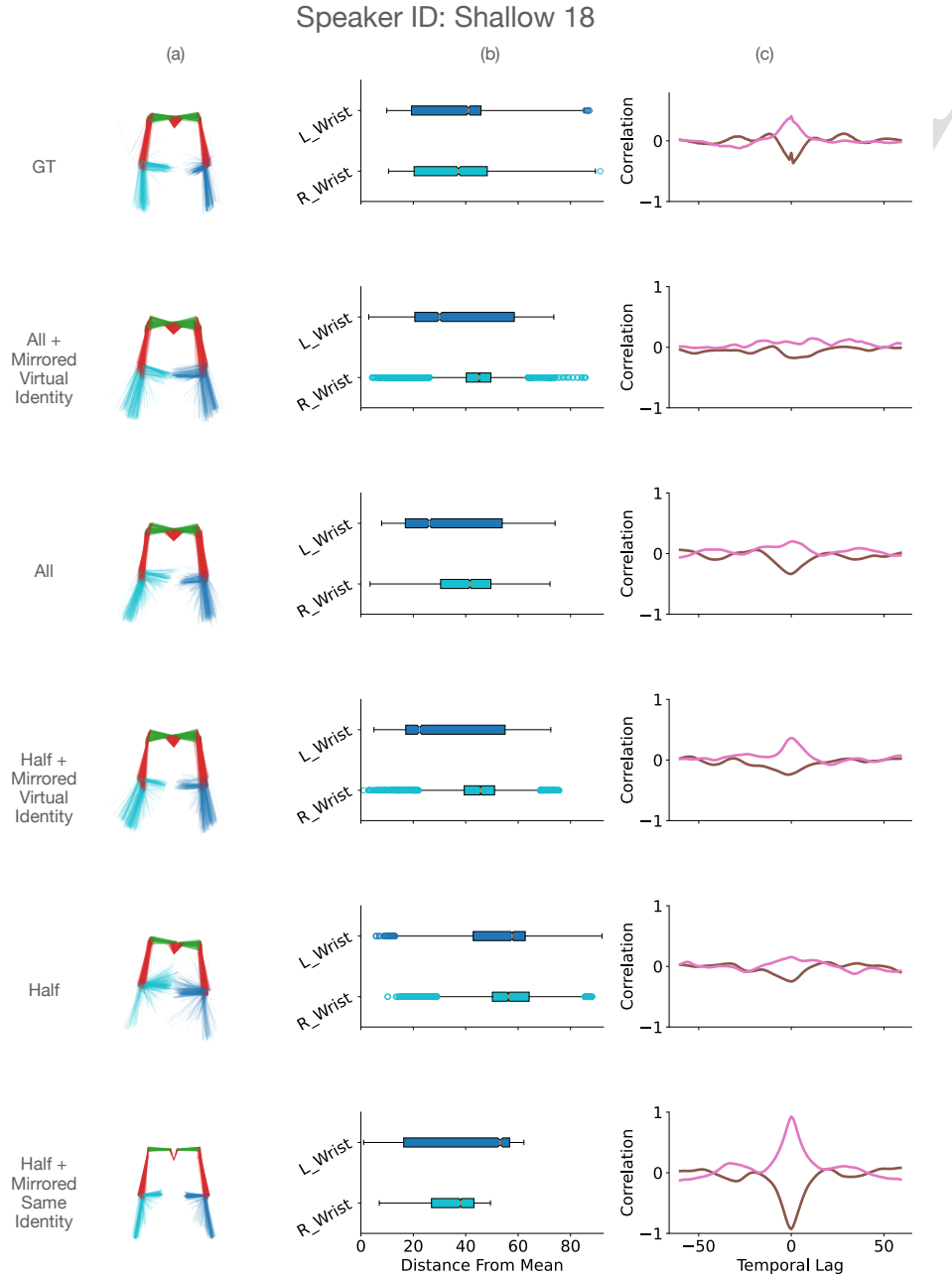


Figure 17: A comparison for a single speaker's generated motion showing detrimental impact of including mirrored motion under the same identity. Each row corresponds to a different data split used. Column (a) contains the orthographic projection of a pose at every second in the sequence. Column (b) shows the distribution of distances from the mean arm pose. Column (c) shows the cross correlation lags between the onset of left wrist motion given right wrist motion in the Z (left-right) and Y (up-down) shown in **brown** and **pink** respectively

Highlights

- Review the motion symmetry of multiple speakers during dyadic conversation, analysing positional, temporal and informational symmetry.
- Discuss the efficacy of lateral mirroring of the human body as a means of data augmentation.
- Conclude lateral mirroring is only applicable in certain cases and is not suited as a generic approach.
- Suggest lateral mirroring is suitable for increasing the number of identities in a data set, including the mirrored data as a new speaker.
- Propose our statistical analysis for evaluating performance of speech-driven conversational agents.

Jonathan Windle: Conceptualisation, Methodology, Software, Formal analysis, Writing- Original Draft. **Sarah Taylor:** Conceptualisation, Methodology, Writing- Reviewing and Editing. **David Greenwood:** Conceptualisation, Methodology, Writing- Reviewing and Editing. **Iain Matthews:** Conceptualisation, Methodology, Writing- Reviewing and Editing

Journal Pre-proof

Declaration of interests

☐ The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

☒ The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Iain Matthews reports a relationship with Epic Games that includes: employment.