

- <sup>1</sup> Detecting social information in a dense database of infants' natural visual experience

### Abstract

The faces and hands of caregivers and other social partners offer a rich source of social and causal information that may be critical for infants' cognitive and linguistic development.

Previous work using manual annotation strategies and cross-sectional data has found systematic changes in the proportion of faces and hands in the egocentric perspective of young infants. Here, we examine the prevalence of faces and hands in a longitudinal collection of more than 1700 headcam videos collected from three children along a span of 6 to 32 months of age—the SAYCam dataset (Sullivan, Mei, Perfors, Wojcik, & Frank, 2021).

To analyze these naturalistic infant egocentric videos, we first validated the use of a modern convolutional neural network of pose detection (OpenPose) for the detection of faces and hands and then applied this model to the entire dataset. First, we found a higher proportion of hands in view than previously reported and a moderate decrease in the proportion of faces in children's view across age. Second, we found variability in the proportion of faces/hands viewed by different children in different locations (e.g., living room vs. kitchen), suggesting that individual activity contexts may shape the social information that infants experience.

Third, we found evidence that children may see closer, larger views of people, hands, and faces earlier in development. These analyses provide new insight into the changes in the social information in view across early development and call for further work that examines their generalizability across populations and their relationship to learning outcomes.

*Keywords:* social cognition; face perception; infancy; head cameras; deep learning  
Word count: 4392

<sup>23</sup> Detecting social information in a dense database of infants' natural visual experience

<sup>24</sup> **Introduction**

<sup>25</sup> Infants are confronted by a blooming, buzzing onslaught of stimuli (James, 1890) which  
<sup>26</sup> they must learn to parse to make sense of the world around them. Yet they do not embark  
<sup>27</sup> on this learning process alone: From as early as 3 months of age, young infants follow overt  
<sup>28</sup> gaze shifts (Gredeback, Theuring, Hauf, & Kenward, 2008), and even newborns prefer to  
<sup>29</sup> look at faces with direct vs. averted gaze (Farroni, Csibra, Simion, & Johnson, 2002), despite  
<sup>30</sup> their limited acuity. As faces are likely to be an important conduit of social information that  
<sup>31</sup> scaffolds cognitive development, psychologists have long hypothesized that faces are  
<sup>32</sup> prevalent in the visual experience of young infants.

<sup>33</sup> Yet until recently most hypotheses about infants' visual experience have gone untested.  
<sup>34</sup> Though parents and scientists alike have strong intuitions about what infants see, even the  
<sup>35</sup> viewpoint of a walking child is hard to intuit (Clerkin, Hart, Rehg, Yu, & Smith, 2017;  
<sup>36</sup> Franchak, Kretch, Soska, & Adolph, 2011). By equipping infants and toddlers with  
<sup>37</sup> head-mounted cameras, researchers have begun to document the infant's egocentric  
<sup>38</sup> perspective on the world. Using these methods, a growing body of work now demonstrates  
<sup>39</sup> that the viewpoints of very young infants (less than 4 months of age) are indeed dominated  
<sup>40</sup> by frequent, persistent views of the faces of their caregivers (Jayaraman, Fausey, & Smith,  
<sup>41</sup> 2015; Jayaraman & Smith, 2018; Sugden, Mohamed-Ali, & Moulson, 2014).

<sup>42</sup> Beyond these early months, infants' motor and cognitive abilities mature, leading to  
<sup>43</sup> vastly different perspectives on the world. For example, crawlers see fewer faces and hands  
<sup>44</sup> than do walking children (Franchak, Kretch, & Adolph, 2017; Kretch, Franchak, & Adolph,  
<sup>45</sup> 2014; Sanchez, Long, Kraus, & Frank, 2018) as well as different views of objects (Smith, Yu,  
<sup>46</sup> & Pereira, 2011). Further, as infants learn to use their own hands to act on the world, they  
<sup>47</sup> seem to focus on manual actions taken by their social partners, and their perspective starts

48 to capture views of hands manipulating objects (Fausey, Jayaraman, & Smith, 2016). In  
49 turn, caregivers may also start to use their hands with more communicative intent, directing  
50 infants' attention by pointing and gesturing to different events and objects during play (Yu  
51 & Smith, 2013).

52 Here, we examine the social information present in the infant visual perspective—the  
53 presence of faces and hands—by analyzing a longitudinal collection of more than 1700  
54 headcam videos collected from three children along a span of 6 to 32 months of age—the  
55 SAYCam dataset (Sullivan et al., 2021). In addition to its size and longitudinal nature, this  
56 dataset is more naturalistic than those previously used in two key ways. First, recordings  
57 were taken under a large variety of activity contexts (Bruner, 1985; Roy, Frank, DeCamp,  
58 Miller, & Roy, 2015) encompassing infants' viewpoints during both activities outside and  
59 inside the home. Even in other naturalistic datasets, the incredible variety in a typical  
60 infant's experience has been largely underrepresented (see examples in Figure 1; e.g., riding  
61 in the car, gardening, watching chickens during a walk, browsing magazines, nursing,  
62 brushing teeth). Second, the head-mounted cameras used in the SAYCam dataset captured a  
63 larger field of view than those typically used, allowing a more complete picture of the infant  
64 perspective. While head-mounted cameras with a more restricted field of view do represent  
65 where infants are foveating most of the time (Smith, Yu, Yoshida, & Fausey, 2015; Yoshida  
66 & Smith, 2008), they may fail to capture short saccades to either faces or hands in the  
67 periphery, as the timescale of head movements is much longer.

68 With hundreds of hours of footage (>40M frames), however, this large dataset  
69 necessitates a shift to an automated annotation strategy. Indeed, annotation of the frames  
70 extracted from egocentric videos has been prohibitively time-consuming, meaning that most  
71 frames are typically not inspected, even in the most comprehensive studies. For example,  
72 Fausey et al. (2016) collected a total of 143 hours of head-mounted camera footage (15.5  
73 million frames), of which one frame every five seconds was hand-annotated (by four coders),

74 totalling 103,383 frames (per coder)—an impressive number of annotations but nonetheless  
75 only 0.67% of the collected footage. To address this challenge, we use a modern computer  
76 vision model of pose detection to automatically detect the presence of hands and faces from  
77 the infant egocentric viewpoint. Specifically, we use OpenPose (Cao, Hidalgo, Simon, Wei, &  
78 Sheikh, 2018), a model optimized for jointly detecting human face, body, hand, and foot  
79 keypoints that operates well on scenes including multiple people, even if they are  
80 partially-occluded (see Figure 1). In prior work examining egocentric videos, OpenPose  
81 performed comparably to other modern face detection models (Sanchez et al., 2018).

82 In this paper, we first describe the dataset and validate the use of this model by  
83 comparing face and hand detections to a human-annotated set of 24,000 frames. Next, we  
84 report how the proportion of faces and hands changes with age in each of the three children  
85 in the dataset. We then investigate sources of variability in our more naturalistic dataset  
86 that may explain differences from prior work, including both the field-of-view of the head  
87 cameras as well as a diversity of locations in which videos were recorded. Finally, making use  
88 of automated annotation of pose bounding boxes, we analyze the size, location, and  
89 variability of detected faces and poses across development.

90

## Method

### 91 Dataset

92 The dataset is described in detail in Sullivan et al. (2021); we summarize these details  
93 here. Children wore Veho Muvi miniature cameras mounted on a custom camping headlamp  
94 harness (“headcams”) at least twice weekly, for approximately one hour per recording session.  
95 One weekly session was on the same day each week at a roughly constant time of day, while  
96 the other(s) were chosen arbitrarily at the participating family’s discretion. At the time of  
97 the recording, all three children were in single-child households. Videos captured by the

98 headcam were 640x480 pixels, and a fisheye lens was attached to the camera to increase the  
99 field of view to approximately 109 degrees horizontal x 70 degrees vertical. Videos<sup>1</sup> with  
100 technical errors or that were not taken from the egocentric perspective were excluded from  
101 the dataset. We analyze 1745 videos, with a total duration of 391.11 hours (>40 million  
102 frames).

### 103 **Detection Method**

104 To automatically annotate the millions of frames in SAYCam, we used a pose detector,  
105 OpenPose<sup>2</sup> (Cao et al., 2018; Simon, Joo, Matthews, & Sheikh, 2017), which provided the  
106 locations of 18 body parts (ears, nose, wrists, etc.). To do so, a convolutional neural network  
107 was used for initial anatomical detection, and part affinity fields were subsequently applied  
108 for part association to produce a series of body part candidates. Once these body part  
109 candidates were matched to a single individual in the frame, they were finally assembled into  
110 a pose. Thus, while we only made use of the outputs of the face and hand detections, the  
111 entire set of pose information from an individual was used to determine the presence of a  
112 face/hand, making the process more robust to occlusion than methods optimized to detect  
113 only faces or hands. Note, however, that these face/hand detections are reliant on the  
114 detection of at least a partial pose, so some very up-close views of faces/hands may go  
115 undetected.

### 116 **Detection Validation**

117 To test the validity of OpenPose's hand and face detections, we compared the accuracy  
118 of these detections relative to human annotations of 24,000 frames selected uniformly at  
119 random from videos of two children (S and A); 24000 frames sampled from allocentric videos

---

<sup>1</sup>All videos are available at <https://nyu.databrary.org/volume/564>

<sup>2</sup><https://github.com/CMU-Perceptual-Computing-Lab/openpose>

were excluded, and these videos were also excluded from the other analyses. Frames were jointly annotated for the presence of faces and hands by one author. A second set of coders recruited via AMT (Amazon Mechanical Turk) additionally annotated 3150 frames; agreement with the primary coder was >95%.

As has been observed in other studies on automated annotation of headcam data (e.g. Frank, Simmons, Yurovsky, & Pusiol, 2013; Bambach, Lee, Crandall, & Yu, 2015; Sanchez et al., 2018), detection tasks that are easy in third-person video can be quite challenging in egocentric videos, due to difficult angles and sizes as well as substantial occlusion. For example, the infant perspective often contains non-canonical viewpoints of faces (e.g., looking up at a caregiver's chin) as well as partially-occluded or oblique viewpoints of both faces and hands. Further, hand detection tends to be a harder computational problem than face detection (Bambach et al., 2015; Simon et al., 2017). We thus expected overall performance to be lower in these naturalistic videos than on either photos taken from the adult perspective or in egocentric videos in controlled, laboratory settings (e.g., Sanchez et al., 2018).

To evaluate OpenPose's performance, we compared its detections to the manually-annotated gold set of frames, calculating precision (hits / (hits + false alarms)), recall (hits / (hits + misses)), and F-score (the harmonic mean of precision and recall). In our data, for faces, the F-score was 0.64, with a precision of 0.70 and recall of 0.58. For hands, the F-score was 0.51, with a precision of 0.73 and recall of 0.40. While face and hand detections showed moderately good precision, face detections were overall slightly more accurate than hand detections. In general, hand detections suffered from fairly low recall, indicating that OpenPose likely underestimated the proportion of hands in the dataset. We also found that restricting our detections to high-confidence face/hand detections (>.5 confidence, default threshold for visualization in OpenPose) was not beneficial – improving precision but dramatically impairing recall and thus overall performance: the F-score for

<sup>146</sup> high-confidence face detections was 0.41, with a precision of 0.95 and recall of 0.26; for  
<sup>147</sup> high-confidence hand detections, the F-score was 0.18, with a precision of 0.97 and recall of  
<sup>148</sup> 0.10)

<sup>149</sup> We suspected that this was in part because children’s own hands were often in view of  
<sup>150</sup> the camera and unconnected to a pose—a notoriously challenging detection problem  
<sup>151</sup> (Bambach et al., 2015). To assess this possibility, we obtained human annotations for the  
<sup>152</sup> entire subsample of 9051 frames in which a hand was detected; participants (recruited via  
<sup>153</sup> AMT) were asked to draw bounding boxes around children’s and adult’s hands. Overall, we  
<sup>154</sup> found that 43% of missed hand detections were of child hands. When frames with children’s  
<sup>155</sup> hands were removed from the gold set, recall did improve somewhat to 0.57. We also  
<sup>156</sup> observed that children’s hands tended to appear in the lower half of the frames; heatmaps of  
<sup>157</sup> the bounding boxes obtained from these annotations can be seen in Appendix Figure B1.

<sup>158</sup> Finally, we examined whether the precision, recall, and F-score for hands and faces  
<sup>159</sup> varied with age or child, and did not find substantial variation. Thus, while OpenPose was  
<sup>160</sup> trained on photographs from the adult perspective, this model still generalized relatively well  
<sup>161</sup> to the egocentric infant viewpoint with no fine-tuning or post-processing of the detections.  
<sup>162</sup> As these detections were imperfect compared to human annotators, fine-tuning these models  
<sup>163</sup> to better optimize for the infant viewpoint remains an open avenue for future work.  
<sup>164</sup> Standard computer vision models are rarely trained on the egocentric viewpoint, and we  
<sup>165</sup> suspect that training these models on more naturalistic data may lead to more robust,  
<sup>166</sup> generalizable detectors.

167

## Results and Discussion

168 **Access to social information across age**

169 We analyzed the social information in view across the entire dataset, looking  
170 specifically at the proportions of faces and hands detected for each child.<sup>3</sup> Data from videos  
171 were binned according to the age of the child (in weeks). First, we saw that the proportion  
172 of faces in view showed a moderate decrease across this age range (see Figure 2), in keeping  
173 with prior findings (Fausey et al., 2016); in contrast, we did not observe an increase in the  
174 proportion of hands in view. These effects were quantified with two separate linear  
175 mixed-models (see Tables 1 & 2).<sup>4</sup>

176 However, the most striking result from these analyses is a much overall greater  
177 proportion of hands in view than have previously been reported (Fausey et al., 2016). We  
178 found this to be true across all ages, in all three children, and regardless of whether we  
179 analyzed human annotations (on the 24K random subset, see dotted lines in Appendix  
180 Figure A1) or OpenPose annotations on the entire dataset (see Figure 2A). This is notable  
181 especially given that OpenPose showed relatively low recall for hands, indicating that this  
182 may be an underestimate of the proportion of hands in view. In fact, analysis the human  
183 annotations underscores revealed a much higher proportion of hands relative to faces than  
184 the automated annotations.

185 One reason this could be the case is the much larger field of view that was captured by  
186 the cameras used in this study: These cameras were outfitted with a fish-eye lens in an

<sup>3</sup>All analyses and preprocessed data files for this paper are available at <https://tinyurl.com/detecting-social-info>

<sup>4</sup>Face/hand detections were binned across each week of filming. Participant's age was converted into months and centered for these analyses. Random slopes for the effect of age by child led to a singular fit and were removed from both analyses; see full model specification in accompanying codebase.

attempt to capture as much of the children's field of view as possible, leading to a larger field of view (109 degrees horizontal x 70 degrees vertical) than in many previous studies; for example, in Fausey et al. (2016) the FOV was 69 x 41 degrees. This larger FOV may have allowed the SAYCam cameras to capture not only the presence of a social partner's hands interacting with objects or gestures, but also the children's own hands, leading to more frequent hand detections.

As we found that children's hands tended to occur in the lower visual field (see Figure B1), we thus re-analyzed the entire dataset while restricting our analysis to the center field of view, decreasing the proportion of hand detections from 24% to 16%, but only decreased face detections from 20% to 9.90%. This cropping likely removed both the majority of detections of children's own hands but also some detections of adult hands (see Figure B1), especially as OpenPose was biased to miss children's hands when they were in view. Nonetheless, within this modified field of view, we still observed more hand detections than face detections (see dashed lines in Figure 2). We also still found a higher proportion of hands in view relative to faces when excluding any frames containing child hand's from the human annotated gold sample (see Appendix FigureA1).

As we found that children's hands tended to occur in the lower visual field (see Figure B1), we thus re-analyzed the entire dataset while restricting our analysis to the center field of view, decreasing the proportion of hand detections from 24% to 16%, but only decreased face detections from 20% to 9.90%. This cropping likely removed both the majority of detections of children's own hands but also some detections of adult hands (see Figure B1), especially as OpenPose was biased to miss children's hands when they were in view. Nonetheless, within this modified field of view, we still observed more hand detections than face detections (see dashed lines in Figure 2B). We also still found a higher proportion of hands in view relative to faces when excluding any frames containing child hand's from the human annotated gold sample (see Appendix Figure A1).

Finally, we analyzed how these two sources of social information co-occurred, finding that faces/hands were jointly present in 11.50 percent of frames (see face hand-occurrences across age in Figure 2C). To do so, we calculated the number of frames in which infants saw faces and hands together relative to overall proportions of faces/hands that were detected for each child and age range. As shown in Figure 3, we found that all three infants were more likely to see hands independently – without the presence of a face – than they were likely to see faces independently. That is, generally speaking when a face was present, a hand also tended to be present.

## 221 Variability in social information across learning contexts

How does the child's context influence the social information in view? Bruner (1985) discussed the role of children's activities in shaping the information present for learning. Following this idea, we investigated whether there were differences in access to faces by the activity that the child was engaged in. This hypothesis seems intuitively appealing.

Some activities seem likely to be characterized by a much higher proportion of faces (e.g., diaper changes) than others (e.g., a car trip). Following this same idea, perhaps other activities involve the presence of more hands in the field of view (e.g., playtime).

We did not have access to annotations of activity. Thus, following Roy et al. (2015), we used spatial location as a proxy for activity context, taking advantage of the presence of these annotations for a subset of the SAYCam videos. Of the 1745 videos in the dataset, 639 were annotated for the location or locations they were filmed in. These location annotations were only available for two children, S and A. Annotated locations mostly consisted of rooms of the house (e.g., “living room”) but also included some other locations (e.g., “car,” “outside”). Of this set, 296 videos were filmed in only a single location (e.g., the location label did not change within the video), representing 17 percent of the dataset and over 5 million frames. In our viewing of the SAYCam videos and in other annotations available with

238 the dataset, activities varied somewhat predictably by location: for example, eating tended  
239 to occur in the kitchen, whereas playtime was the dominant activity in the living room.

240 Figure 4 shows the proportion of faces vs. hands across locations. We found substantial  
241 variation across locations and, to some extent, across children. Separate chi-squared tests for  
242 each child and detection type revealed significant variability in detections by location in each  
243 case, with all  $p < .001$ . For example, while both A and S saw a relatively similar proportion  
244 of faces and hands in the bedroom, the two children saw quite different amounts of faces and  
245 hands from one another in the kitchen. This difference is likely explained by differences in  
246 arrangement of the kitchen in the two children’s households (Sullivan, personal  
247 communication), such that mealtimes in one kitchen resulted in a face-to-face orientation  
248 while it did not in the other). This example illustrates how specifics of the geometry of a  
249 particular context can play an outsize role in the child’s access to social information during  
250 that context.

### 251 **Fine-grained changes in the social information in view**

252 In a third set of analyses, we explored fine-grained changes in the SAYCam infants’  
253 access to social information across development. In these analyses, we capitalize on the fact  
254 that OpenPose provides not only face and hand detections but also positional keypoints. In  
255 particular, we explored this keypoint dataset with the idea that greater mobility allows older  
256 children to be further from their caregivers on average. Thus, younger, less mobile children  
257 may tend to see larger faces towards the center of their visual field while older, more mobile  
258 children may experience more smaller, more variable views of faces. The same dynamic would  
259 be predicted hold for hands as well, as it would be driven by overall differences in distance.

260 Supporting this idea, we found that the averages sizes of the people, faces, and hands  
261 in the infant view became smaller over development (Figure 5). This effect was relatively

262 consistent across the three children in the dataset, despite the fact that the three children  
263 showed sometimes disparate overall proportions of faces/hands in view. Thus, children may  
264 see closer, larger views of people, hands, and faces earlier in development.

265 In keeping with this hypothesis, we also found evidence that faces tended to be farther  
266 away from older children. We restricted our analysis here to faces where both eyes were  
267 detected and computed interpupillary distance as a rough metric of distance, since eyes  
268 should be closer together on average when a face is further from the camera. Figure 6A  
269 shows the average interpupillary distance on faces as a function of each child's age at the  
270 time of recording. There is a trend from larger, closer faces (with a larger interpupillary  
271 distance) to smaller faces that were farther away (with a smaller interpupillary distance).

272 Finally, we also examined whether there were changes in where faces tended to appear  
273 in the camera's (and hence, by proxy, the child's) field of view. As expected, faces tended to  
274 be located towards the upper field of view, while views of hands were more centrally  
275 distributed (see Appendix, Figure C1 for average density distributions). However, we also  
276 found evidence that older children tended to see more faces in more variable positions than  
277 younger children. Specifically, we examined how variable the horizontal and vertical  
278 coordinates were of the faces in the infant view. To do so, we calculated the coefficient of  
279 variation of the horizontal (x) and vertical (y) positions of centers of the faces detected by  
280 OpenPose (see Figure 6B), and examined changes across age. Faces tended to be more  
281 variable in the vertical than their horizontal position (see Figure 6B). We also found that as  
282 children got older, they tended to see faces that varied more in their horizontal – but not  
283 their vertical position – suggesting that older children might be more likely to see more  
284 smaller faces in their periphery (see Figure 6B).

285

## General Discussion

286 Here, we analyzed the social information in view in a dense, longitudinal dataset,  
287 applying a modern computer-vision model to quantify the hands and faces seen from each of  
288 three children's egocentric perspective from 6 to 32 months of age. First, we found a  
289 moderate decrease across age in the proportion of faces in view in the videos, in keeping with  
290 previous work (Fausey et al., 2016). This finding is particularly notable given that, in  
291 previous cross-sectional data, this effect seems to be most strongly driven by infants younger  
292 than 4 months of age (e.g., Fausey et al., 2016; Jayaraman et al., 2015; Sugden et al., 2014)  
293 who see both more frequent and more persistent faces (Jayaraman & Smith, 2018).

294 We also found this to be true when restricting our analyses to full-field faces, suggesting  
295 this effect is not driven by a concurrent shift from more full-view to partial-views of faces.

296 We also found an unexpectedly high proportion of hands in the view of infants, even  
297 when restricting the field-of-view to the center field of view the videos to make the  
298 viewpoints comparable to those of headcams used in previous work (Fausey et al., 2016).  
299 Why might this be the case? One idea is that these videos contain the viewpoints of children  
300 not only during structured interactions (e.g., play sessions at home or in the lab) but during  
301 everyday activities when children may be playing by themselves or simply observing the  
302 actions of caregivers and other people in their environment. During these less structured  
303 times, caregivers may move about in the vicinity of the child but not interact with them as  
304 directly—leading to views where a person and their hands are visible from a distance, but  
305 this person's face may be turned away from the infant or occluded (see examples in Figure  
306 1). Indeed, using the same pose detector on videos from in-lab play sessions, Sanchez et al.  
307 (2018) found the opposite trend: slightly fewer hand detections than face detections from  
308 8-16 months of age. Work that directly examines the variability in the social information in  
309 view across more vs. less structured activity contexts could further test this idea.

310 A coarse analysis based on the location the videos were filmed in further highlights the

311 variability of the social information in view during different activities, showing differences

312 across locations and between individual children. Within a given, well-defined context—e.g.,

313 mealtime in kitchens—S saw more faces than A, and S saw more faces in the kitchen than in

314 other locations. This variability likely stems from the fact that there are at least three ways

315 to feed a young child: 1) sitting in front of the child, facing them as they sit in a high chair;

316 2) sitting behind the child, holding them as they face outward, and 3) sitting side by side.

317 Each of these positions offer the child differing degrees of visual access to faces and hands.

318 While the social information in view may be variable across children in different activity

319 contexts, these analyses suggest they could be stable within a given child’s day-to-day

320 experience.

321 We also used these detailed pose annotations to explore finer-grained changes in how

322 children experience the faces and hands of their caregivers over development. We found that

323 the faces, hands, and people in the infant view tended to become smaller and that faces

324 tended to be farther away and in more variable horizontal positions. Overall, these data

325 support the idea that the social information in view changes across development as infants

326 become increasingly mobile and independent. As children can navigate the world on their

327 own, they may experience fewer close-up interactions with the caregivers and more bouts of

328 play where they are exploring the objects and things in the environment around them.

329 More broadly, however, these analyses underscore the importance of how, when, from

330 whom, and what data we sample; these choices become central when we attempt to draw

331 conclusions about the regularities of experience. Indeed, while unprecedented in size, this

332 dataset still has many limitations. These videos only represent a small portion of the

333 everyday experience of these three children, all of whom come from relatively privileged

334 households in western societies and thus are not representative in many ways of the global

335 population. Any idiosyncrasies in how and when these particular families chose to film these

<sup>336</sup> videos also undoubtedly influences the variability seen here. And without eye-tracking data,  
<sup>337</sup> we do not know if children are attending to the social information in their visual field.

Nonetheless, we believe that these advances in datasets and methodologies represent a step in the right direction. The present paper demonstrates the feasibility of using a modern computer vision model to annotate the entirety of a very large dataset (here, >40M million frames) for the presence and size of people, hands, and faces, representing orders of magnitude more data relative to prior work. We propose that the large-scale analysis of dense datasets, collected with different fields of view, cameras and from many different laboratories, will lead to generalizable conclusions about the regularities of infant experience that scaffold learning.

346 Acknowledgements

347 Thanks to the creators of the SAYCam dataset who made this work possible and to  
348 Alessandro Sanchez for his contributions to the codebase. This work was funded by a Jacobs  
349 Foundation Fellowship to MCF, a John Mereck Scholars award to MCF, and NSF #1714726  
350 to BLL.

351

## References

- 352 Bambach, S., Lee, S., Crandall, D. J., & Yu, C. (2015). Lending a hand: Detecting hands  
353 and recognizing activities in complex egocentric interactions. In *Proc. Of the IEEE*  
354 *international conference on computer vision* (pp. 1949–1957).
- 355 Bruner, J. (1985). The role of interaction formats in language acquisition. In *Language and*  
356 *social situations* (pp. 31–46). Springer.
- 357 Cao, Z., Hidalgo, G., Simon, T., Wei, S.-E., & Sheikh, Y. (2018). OpenPose: Realtime  
358 multi-person 2D pose estimation using Part Affinity Fields. In *ArXiv preprint*  
359 *arXiv:1812.08008*.
- 360 Clerkin, E. M., Hart, E., Rehg, J. M., Yu, C., & Smith, L. B. (2017). Real-world visual  
361 statistics and infants' first-learned object names. *Phil. Trans. R. Soc. B*, 372(1711),  
362 20160055.
- 363 Farroni, T., Csibra, G., Simion, F., & Johnson, M. H. (2002). Eye contact detection in  
364 humans from birth. *Proceedings of the National Academy of Sciences*, 99(14),  
365 9602–9605.
- 366 Fausey, C. M., Jayaraman, S., & Smith, L. B. (2016). From faces to hands: Changing visual  
367 input in the first two years. *Cognition*, 152, 101–107.
- 368 Franchak, J. M., Kretch, K. S., & Adolph, K. E. (2017). See and be seen: Infant–caregiver  
369 social looking during locomotor free play. *Developmental Science*.
- 370 Franchak, J. M., Kretch, K. S., Soska, K. C., & Adolph, K. E. (2011). Head-mounted eye  
371 tracking: A new method to describe infant looking. *Child Development*, 82(6),  
372 1738–1750.

- 373 Frank, M. C., Simmons, K., Yurovsky, D., & Pusiol, G. (2013). Developmental and postural  
374 changes in children's visual access to faces. In *Proceedings of the 35th annual meeting*  
375 of the cognitive science society (pp. 454–459).
- 376 Gredeback, G., Theuring, C., Hauf, P., & Kenward, B. (2008). The microstructure of infants'  
377 gaze as they view adult shifts in overt attention. *Infancy*, 13(5), 533–543.
- 378 James, W. (1890). *The principles of psychology* (Vol. 1). Macmillan London.
- 379 Jayaraman, S., Fausey, C. M., & Smith, L. B. (2015). The faces in infant-perspective scenes  
380 change over the first year of life. *PLoS One*.  
381 <https://doi.org/10.1371/journal.pone.0123780>
- 382 Jayaraman, S., & Smith, L. B. (2018). Faces in early visual environments are persistent not  
383 just frequent. *Vision Research*.
- 384 Kretch, K. S., Franchak, J. M., & Adolph, K. E. (2014). Crawling and walking infants see  
385 the world differently. *Child Development*, 85(4), 1503–1518.
- 386 Roy, B. C., Frank, M. C., DeCamp, P., Miller, M., & Roy, D. (2015). Predicting the birth of  
387 a spoken word. *Proc. Of the National Academy of Sciences*, 112(41), 12663–12668.
- 388 Sanchez, A., Long, B., Kraus, A. M., & Frank, M. C. (2018). Postural developments  
389 modulate children's visual access to social information. In *Proceedings of the 40th*  
390 *annual conference of the cognitive science society*.
- 391 Simon, T., Joo, H., Matthews, I., & Sheikh, Y. (2017). Hand keypoint detection in single  
392 images using multiview bootstrapping. In *CVPR*.
- 393 Smith, L. B., Yu, C., & Pereira, A. F. (2011). Not your mother's view: The dynamics of  
394 toddler visual experience. *Developmental Science*, 14(1), 9–17.

- 395 Smith, L. B., Yu, C., Yoshida, H., & Fausey, C. M. (2015). Contributions of head-mounted  
396 cameras to studying the visual environments of infants and young children. *Journal  
397 of Cognition and Development*, 16(3), 407–419.
- 398 Sugden, N. A., Mohamed-Ali, M. I., & Moulson, M. C. (2014). I spy with my little eye:  
399 Typical, daily exposure to faces documented from a first-person infant perspective.  
400 *Developmental Psychobiology*, 56(2), 249–261.
- 401 Sullivan, J., Mei, M., Perfors, A., Wojcik, E., & Frank, M. C. (2021). SAYCam: A large,  
402 longitudinal audiovisual dataset recorded from the infant's perspective. *Open Mind*.
- 403 Yoshida, H., & Smith, L. B. (2008). What's in view for toddlers? Using a head camera to  
404 study visual experience. *Infancy*, 13, 229–248.
- 405 Yu, C., & Smith, L. B. (2013). Joint attention without gaze following: Human infants and  
406 their parents coordinate visual attention to objects through eye-hand coordination.  
407 *PloS One*, 8(11).

Table 1

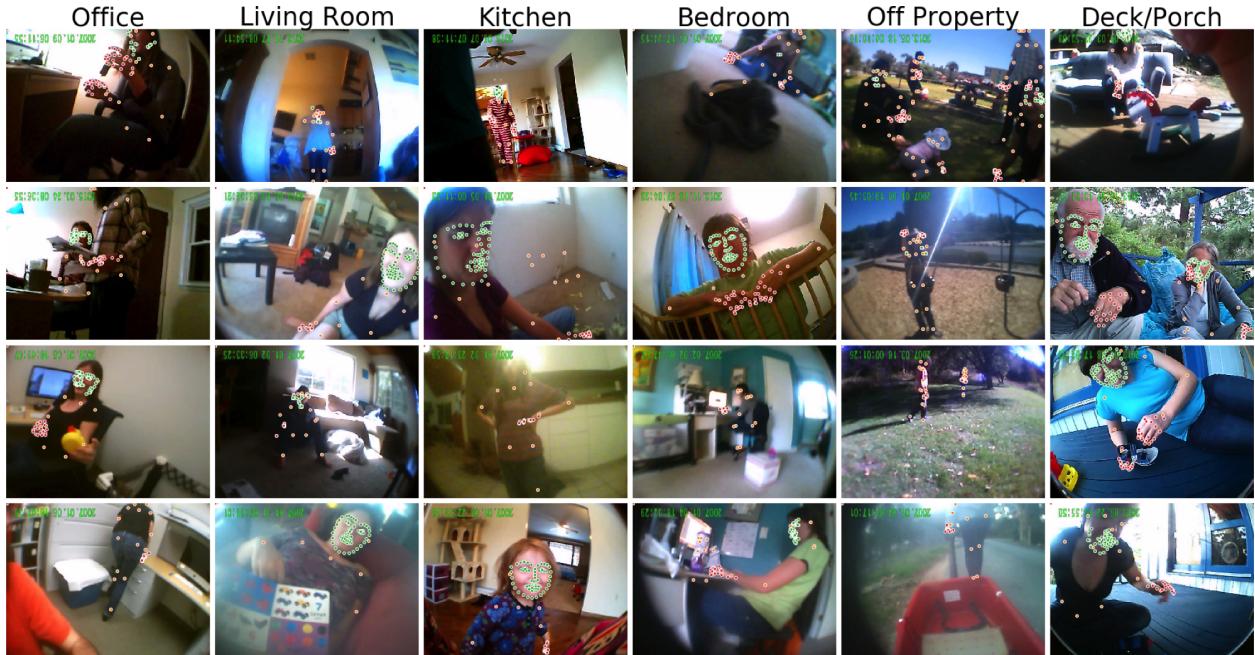
*Coefficients from a mixed-effects regression predicting the proportion of faces seen by infants in the center FOV.*

	Estimate	Std. Error	df	t value	Pr(> t )
(Intercept)	0.098	0.011	1.953	8.850	0.013
Age	-0.195	0.060	429.926	-3.257	0.001
Age**2	-0.160	0.059	429.032	-2.708	0.007

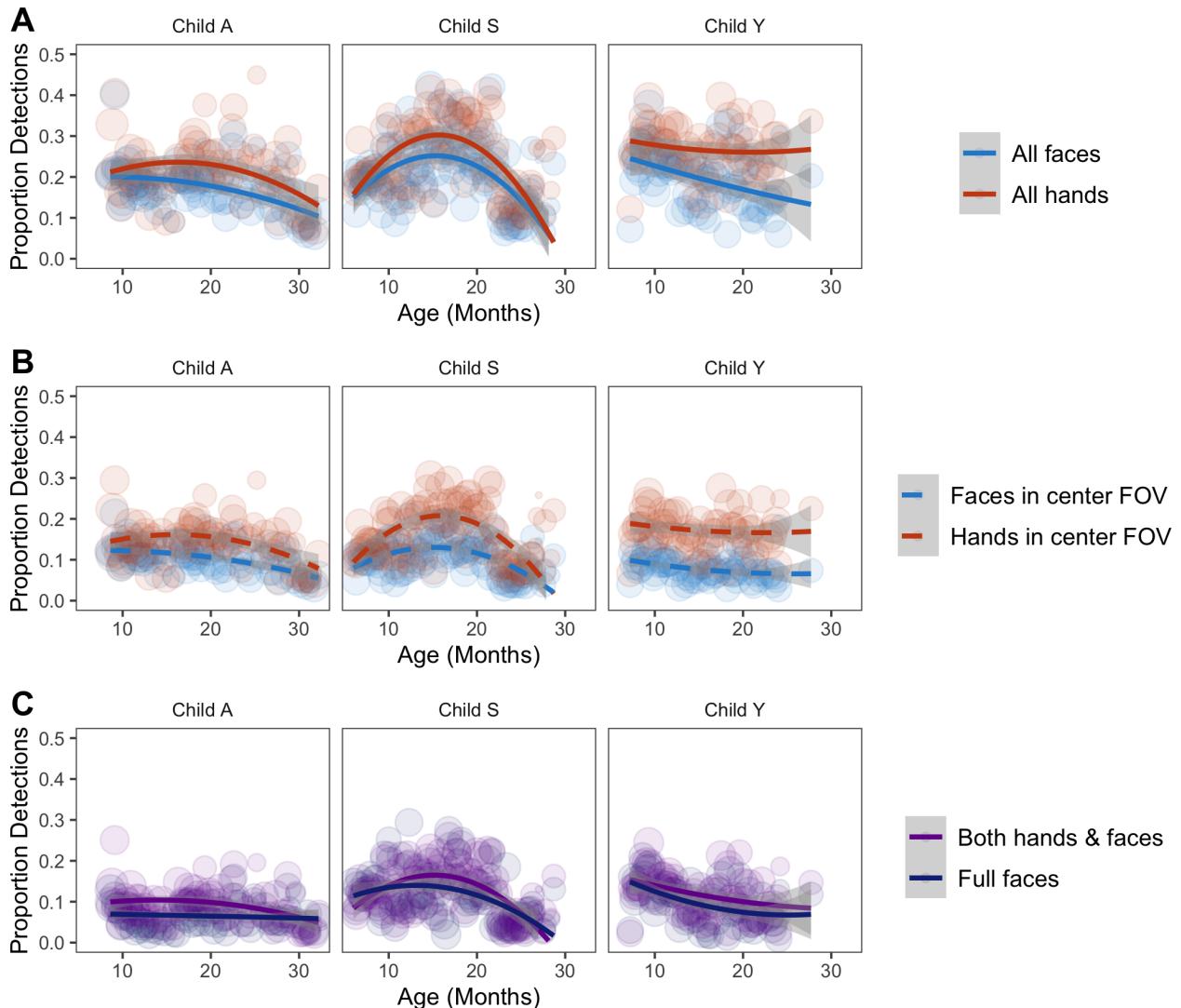
Table 2

*Coefficients from a mixed-effects regression predicting the proportion of hands seen by infants in the center FOV.*

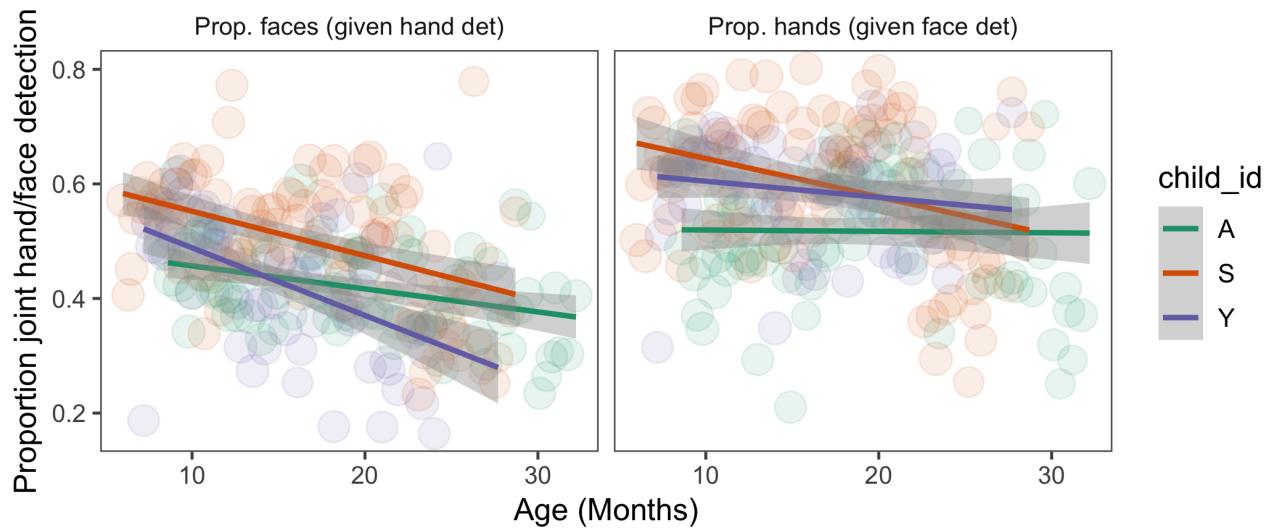
	Estimate	Std. Error	df	t value	Pr(> t )
(Intercept)	0.161	0.007	1.828	21.906	0.003
Age	-0.145	0.078	422.334	-1.855	0.064
Age**2	-0.319	0.077	429.968	-4.134	<.001



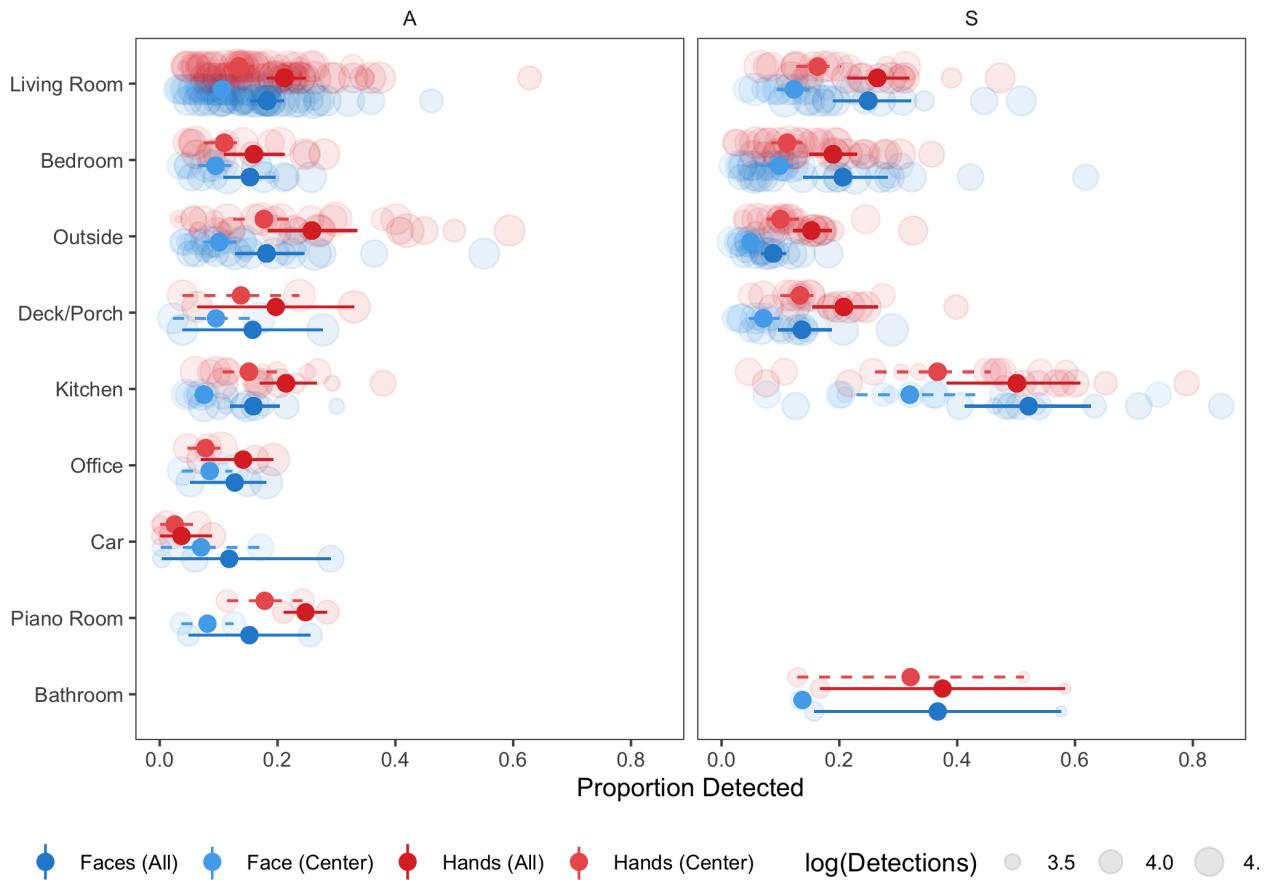
*Figure 1.* Example frames taken from the dataset, illustrating variability in the infant perspective across different locations. OpenPose detections are shown overlaid on these images (green dots = face, red dots = hands, orange dots = pose).



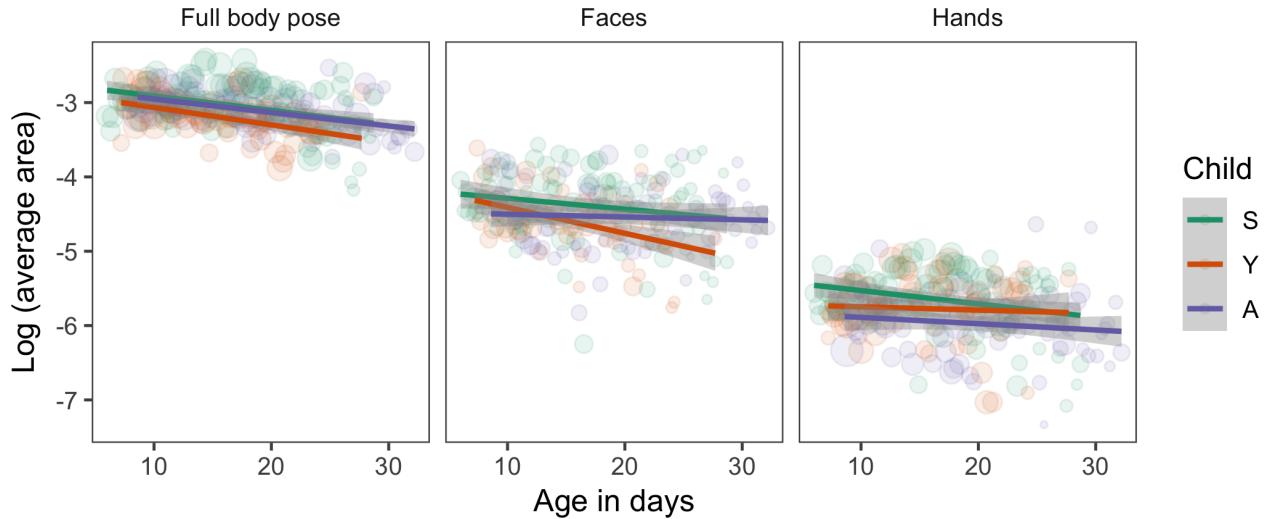
*Figure 2.* Proportion of frames with (A).All face and hand detections, B.Face/hand detections that fell within the center field-of-view (reducing the contribution of children's own hands) and (C) Face detections that were full faces (e.g., eyes, nose, and mouth all visible) and that co-occurred with hands, plotted as a function of age for each child (A, S, and Y). Data are binned by each week that the videos were filmed and scaled by the number of frames in that age range.



*Figure 3.* Proportion of joint face and hands detection within frames where hands (left) or faces (right) were detected.



*Figure 4.* Proportion of faces and hands by location in which egocentric videos were filmed; each panel represents data from an individual child (location annotations were not yet available for Y). Each dot represents data from a week in which videos were filmed and are scaled by the number of frames.



*Figure 5.* Average size of poses, faces, and hands detected in the dataset between eyes in faces detected as a function of age for each child in the dataset (each color = different child). Data are binned by each week that the videos were filmed and scaled by the number of frames in that age range.

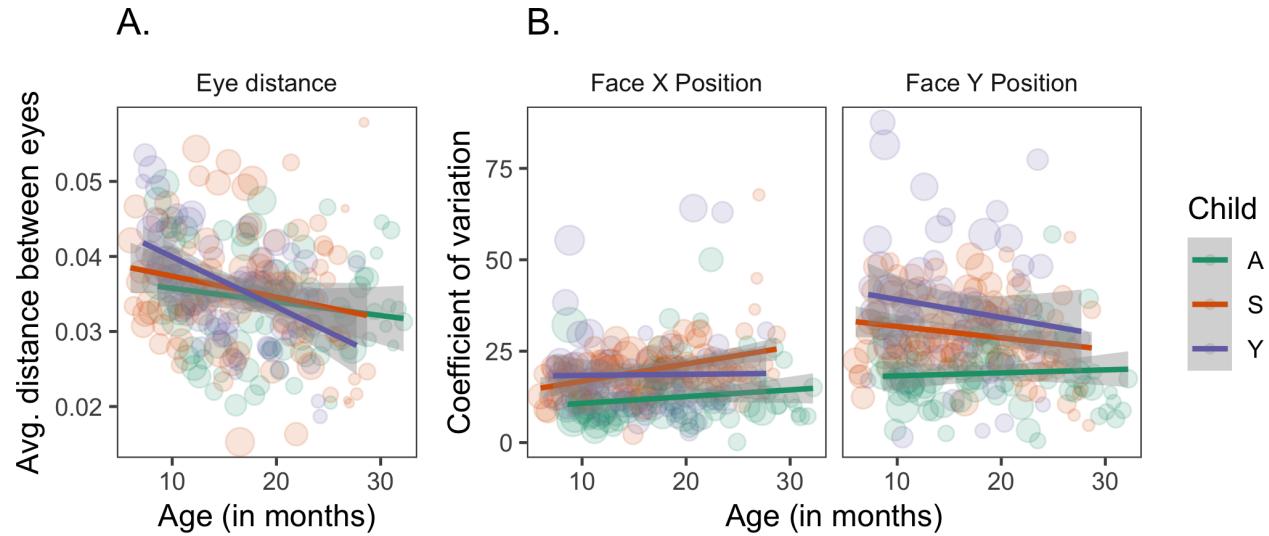
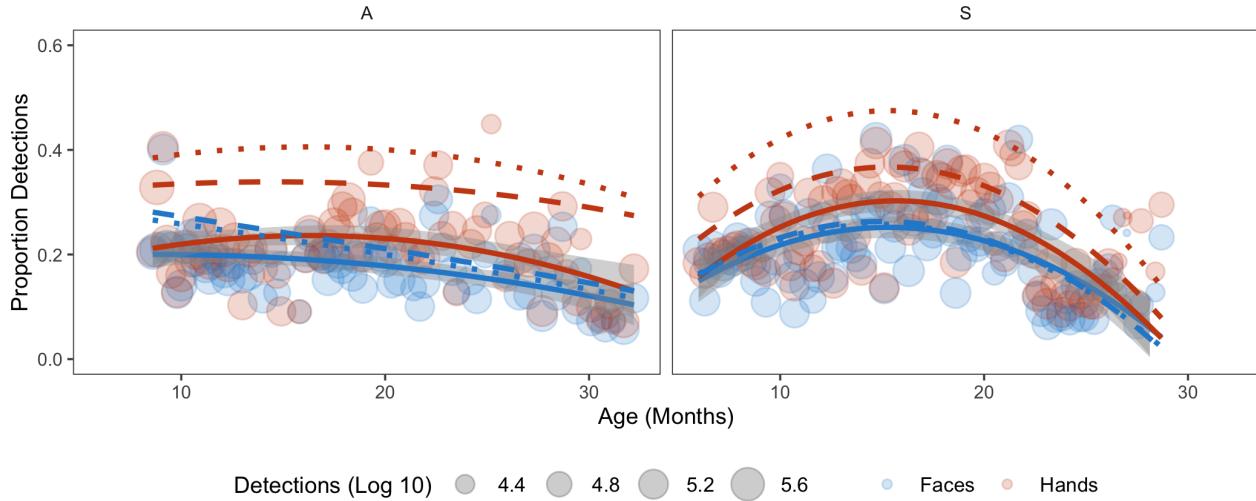


Figure 6. (A) Average distance between eyes and (B) average coefficient of variation for the x and y position of faces detected by OpenPose as a function of each child's age at the time of filming. Data in (A) are restricted to faces where both eyes were detected. Data are binned by each week that the videos were filmed and scaled by the number of face detections in that age range.

## Appendix A

## Face/hand detections relative to human annotations

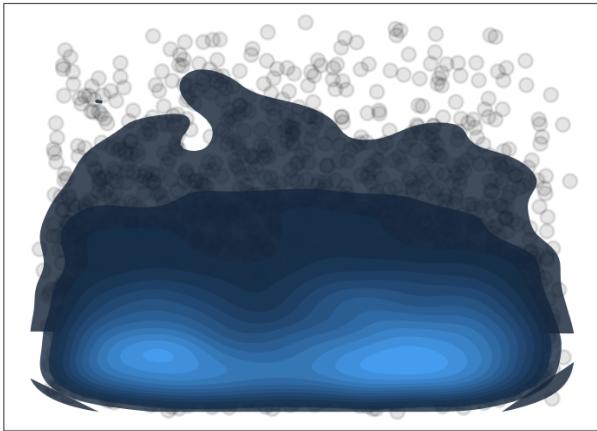


*Figure A1.* Proportion of faces and hands seen as a function of age for each child in the dataset. Data are binned by each week that the videos were filmed and scaled by the number of frames in that age range. Dashed lines show estimated trend lines from proportion of faces/hands in view when analyzing the gold set of frames made by human annotators. Dotted lines show trend lines from the goldset when frames when children's own hand were detected.

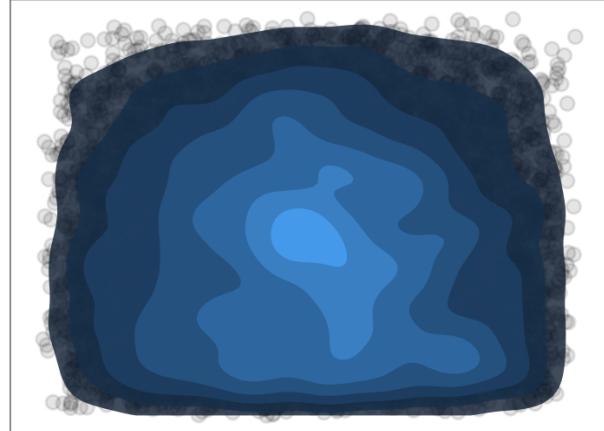
## Appendix B

Density of child vs. adults hands in the visual field

A. Child hand density



B. Adult hand density



*Figure B1.* Density estimates for the child (left) and adult (right) hands that were detected in the 24K frame random gold set; each dot represents the center of a bounding box made by an adult participant. Brighter values indicate more detections.

## Appendix C

## Distribution of faces and hands in the visual field

408 We explored where in the visual field children tended to see faces and hands, suspecting that  
409 these distributions might become wider as children grow older and learn to locomote on their  
410 own, following preliminary analyses from Frank (2012). As expected, faces tended to appear  
411 in the upper visual field in contrast to hands, which tended to be more centrally  
412 located. However, we found little evidence for any changes in the positions of faces and hands  
413 across age, suggesting that this is a relatively stable property of infants' visual environment  
414 from 6 months of age.

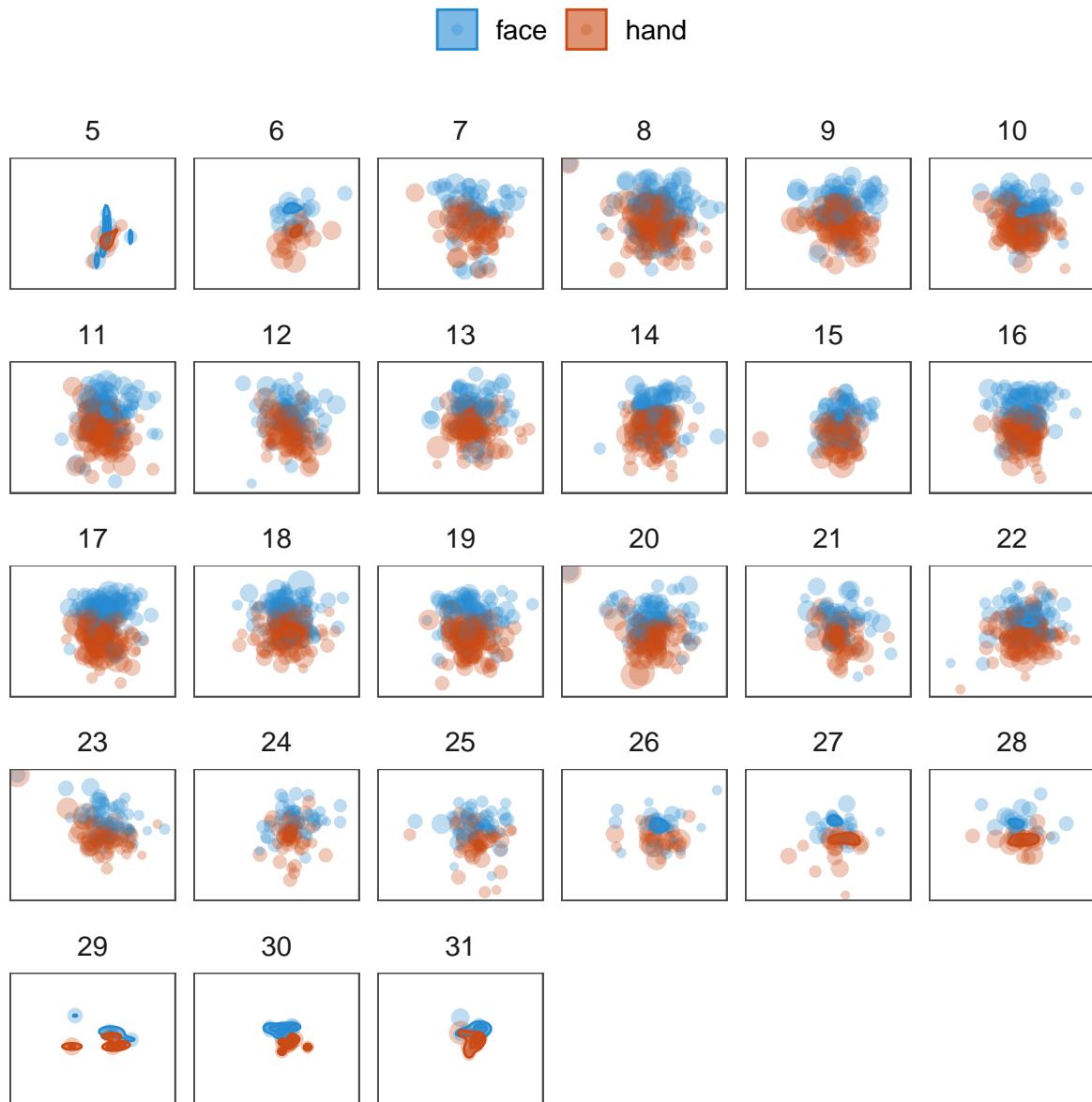


Figure C1. Each panel shows the average position of faces and hands in the visual field; each dot represents the average position from one video within a given age range.