

1 A longitudinal analysis of the social information in infants' naturalistic visual experience
2 using automated detections

3

Abstract

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

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

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

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

23 *Keywords:* social cognition; face perception; infancy; head cameras; deep learning

24 Word count: 4392

25 A longitudinal analysis of the social information in infants' naturalistic visual experience
26 using automated detections

27 **Introduction**

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

36 Yet until recently most hypotheses about infants' visual experience have gone untested.
37 Though parents and scientists alike have strong intuitions about what infants see, even the
38 viewpoint of a walking child is hard to intuit (Clerkin, Hart, Rehg, Yu, & Smith, 2017;
39 Franchak, Kretch, Soska, & Adolph, 2011). By equipping infants and toddlers with
40 head-mounted cameras, researchers have begun to document the infant's egocentric
41 perspective on the world (Franchak et al., 2011; Smith, Jayaraman, Clerkin, & Yu, 2018;
42 Smith, Yu, Yoshida, & Fausey, 2015) and the consequences of this changing view for early
43 learning. Using these methods, a growing body of work now demonstrates that the
44 viewpoints of very young infants (less than 4 months of age) are indeed dominated by
45 frequent, persistent views of the faces of their caregivers (Jayaraman, Fausey, & Smith, 2013,
46 2015, 2017; Jayaraman & Smith, 2018; Sugden, Mohamed-Ali, & Moulson, 2014).

47 Beyond these early months, infants' motor and cognitive abilities mature, leading to
48 vastly different perspectives on the world (Iverson, 2010). For example, children see fewer
49 faces and hands when crawling than walking or sitting (Franchak, 2019; Franchak, Kretch, &

50 Adolph, 2017; Kretch, Franchak, & Adolph, 2014; Luo & Franchak, 2020; Sanchez, Long,
51 Kraus, & Frank, 2018; Yamamoto, Sato, & Itakura, 2020) as well as different views of
52 objects (Luo & Franchak, 2020; Smith, Yu, & Pereira, 2011). Further, as infants learn to use
53 their own hands to act on the world, they seem to focus on manual actions taken by their
54 social partners, and their perspective starts to capture views of hands manipulating objects
55 (Fausey et al., 2016a). In turn, caregivers may also start to use their hands with more
56 communicative intent, directing infants' attention by pointing and gesturing to different
57 events and objects during play (Yu & Smith, 2013).

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

74 With hundreds of hours of footage (>40M frames), however, this large dataset
75 necessitates a shift to an automated annotation strategy. Indeed, annotation of the frames

76 extracted from egocentric videos has been prohibitively time-consuming, meaning that most
77 frames are typically not inspected, even in the most comprehensive studies. For example,
78 Fausey et al. (2016a) collected a total of 143 hours of head-mounted camera footage (15.5
79 million frames), of which one frame every five seconds was hand-annotated (by four coders),
80 totalling 103,383 frames (per coder)—an impressive number of annotations but nonetheless
81 only 0.67% of the collected footage. To address this challenge, we use a modern computer
82 vision model of pose detection to automatically detect the presence of hands and faces from
83 the infant egocentric viewpoint. Specifically, we use OpenPose (Cao, Hidalgo, Simon, Wei, &
84 Sheikh, 2018), a model optimized for jointly detecting human face, body, hand, and foot
85 keypoints that operates well on scenes including multiple people, even if they are
86 partially-occluded (see Figure 1). In prior work examining egocentric videos, OpenPose
87 performed comparably to other modern face detection models (Sanchez et al., 2018).

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

96

Method

97 **Dataset**

98 The dataset is described in detail in Sullivan et al. (2021); we summarize these details
99 here. Children wore Veho Muvi miniature cameras mounted on a custom camping headlamp

100 harness (“headcams”) at least twice weekly, for approximately one hour per recording session.
101 One weekly session was on the same day each week at a roughly constant time of day, while
102 the other(s) were chosen arbitrarily at the participating family’s discretion. At the time of
103 the recording, all three children were in single-child households. Videos captured by the
104 headcam were 640x480 pixels, and a fisheye lens was attached to the camera to increase the
105 field of view to approximately 109 degrees horizontal x 70 degrees vertical. Videos¹ with
106 technical errors or that were not taken from the egocentric perspective were excluded from
107 the dataset. We analyze 1745 videos, with a total duration of 391.11 hours (>40 million
108 frames).

109 Detection Method

110 To automatically annotate the millions of frames in SAYCam, we used a pose detector,
111 OpenPose² (Cao et al., 2018; Simon, Joo, Matthews, & Sheikh, 2017), which provided the
112 locations of 18 body parts (ears, nose, wrists, etc.). To do so, a convolutional neural network
113 was used for initial anatomical detection, and part affinity fields were subsequently applied
114 for part association to produce a series of body part candidates. Once these body part
115 candidates were matched to a single individual in the frame, they were finally assembled into
116 a pose. Thus, while we only made use of the outputs of the face and hand detections, the
117 entire set of pose information from an individual was used to determine the presence of a
118 face/hand, making the process more robust to occlusion than methods optimized to detect
119 only faces or hands. Note, however, that these face/hand detections are reliant on the
120 detection of at least a partial pose, so some very up-close views of faces/hands may go
121 undetected.

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

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

¹²² **Detection Validation**

¹²³ To test the validity of OpenPose's hand and face detections, we compared the accuracy
¹²⁴ of these detections relative to human annotations of 24,000 frames selected uniformly at
¹²⁵ random from videos of two children (S and A); 24000 frames sampled from allocentric videos
¹²⁶ were excluded, and these videos were also excluded from the other analyses. Frames were
¹²⁷ jointly annotated for the presence of faces and hands by one author. A second set of coders
¹²⁸ recruited via AMT (Amazon Mechanical Turk) additionally annotated 3150 frames;
¹²⁹ agreement with the primary coder was >95%.

¹³⁰ As has been observed in other studies on automated annotation of headcam data
¹³¹ (e.g. Frank, Simmons, Yurovsky, & Pusiol, 2013; Bambach, Lee, Crandall, & Yu, 2015;
¹³² Sanchez et al., 2018), detection tasks that are easy in third-person video can be quite
¹³³ challenging in egocentric videos, due to difficult angles and sizes as well as substantial
¹³⁴ occlusion. For example, the infant perspective often contains non-canonical viewpoints of
¹³⁵ faces (e.g., looking up at a caregiver's chin) as well as partially-occluded or oblique
¹³⁶ viewpoints of both faces and hands. Further, hand detection tends to be a harder
¹³⁷ computational problem than face detection (Bambach et al., 2015; Simon et al., 2017). We
¹³⁸ thus expected overall performance to be lower in these naturalistic videos than on either
¹³⁹ photos taken from the adult perspective or in egocentric videos in controlled, laboratory
¹⁴⁰ settings (e.g., Sanchez et al., 2018).

¹⁴¹ To evaluate OpenPose's performance, we compared its detections to the
¹⁴² manually-annotated gold set of frames, calculating precision (hits / (hits + false alarms)),
¹⁴³ recall (hits / (hits + misses)), and F-score (the harmonic mean of precision and recall). In
¹⁴⁴ our data, for faces, the F-score was 0.64, with a precision of 0.70 and recall of 0.58. For
¹⁴⁵ hands, the F-score was 0.51, with a precision of 0.73 and recall of 0.40. While face and hand
¹⁴⁶ detections showed moderately good precision, face detections were overall slightly more

147 accurate than hand detections. In general, hand detections suffered from fairly low recall,
148 indicating that OpenPose likely underestimated the proportion of hands in the dataset. We
149 also found that restricting our detections to high-confidence face/hand detections (>.5
150 confidence, default threshold for visualization in OpenPose) was not beneficial – improving
151 precision but dramatically impairing recall and thus overall performance: the F-score for
152 high-confidence face detections was 0.41, with a precision of 0.95 and recall of 0.26; for
153 high-confidence hand detections, the F-score was 0.18, with a precision of 0.97 and recall of
154 0.10)

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

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

173 generalizable detectors.

174

Results and Discussion

175 Access to social information across age

176 We analyzed the social information in view across the entire dataset, looking
177 specifically at the proportions of faces and hands detected for each child.³ Data from videos
178 were binned according to the age of the child (in weeks). First, we saw that the proportion
179 of faces in view showed a moderate decrease across this age range (see Figure 2), in keeping
180 with prior findings (Fausey et al., 2016a); in contrast, we did not observe an increase in the
181 proportion of hands in view. These effects were quantified with two separate linear
182 mixed-effect models (see Tables 1 & 2)⁴ After visualizing the data (see Figure 2A), we
183 examined whether linear vs. quadratic terms relating children's age to the proportion of
184 faces/hands detected would provide better fits to the data, and found that this was true in
185 both cases (see Tables 1 & 2), though linear terms also provided a significant fit for faces.
186 Thus, these exploratory results point towards the idea that some children may experience
187 overall more social information in view in the second year of life.

188 However, the most striking result from these analyses is a much overall greater
189 proportion of hands in view than have previously been reported (Fausey et al., 2016a). We
190 found this to be true across all ages, in all three children, and regardless of whether we
191 analyzed human annotations (on the 24K random subset, see dotted lines in Appendix
192 Figure A1) or OpenPose annotations on the entire dataset (see Figure 2A). This is notable

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

⁴Face/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.

193 especially given that OpenPose showed relatively low recall for hands, indicating that this
194 may be an underestimate of the proportion of hands in view. In fact, analysis of the human
195 annotations underscores revealed a much higher proportion of hands relative to faces than
196 the automated annotations.

197 One reason this could be the case is the much larger field of view that was captured by
198 the cameras used in this study: These cameras were outfitted with a fish-eye lens in an
199 attempt to capture as much of the children's field of view as possible, leading to a larger field
200 of view (109 degrees horizontal x 70 degrees vertical) than in many previous studies; for
201 example, in Fausey et al. (2016a) the FOV was 69 x 41 degrees. This larger FOV may have
202 allowed the SAYCam cameras to capture not only the presence of a social partner's hands
203 interacting with objects or gestures, but also the children's own hands, leading to more
204 frequent hand detections.

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

215 As we found that children's hands tended to occur in the lower visual field (see Figure
216 B1), we thus re-analyzed the entire dataset while restricting our analysis to the center field
217 of view, decreasing the proportion of hand detections from 24% to 16%, but only decreased
218 face detections from 20% to 9.90%. This cropping likely removed both the majority of

219 detections of children's own hands but also some detections of adult hands (see Figure B1),
220 especially as OpenPose was biased to miss children's hands when they were in view.
221 Nonetheless, within this modified field of view, we still observed more hand detections than
222 face detections (see dashed lines in Figure 2B). We also still found a higher proportion of
223 hands in view relative to faces when excluding any frames containing child hand's from the
224 human annotated gold sample (see Appendix Figure A1).

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

233 Variability in social information across learning contexts

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

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

241 We did not have access to annotations of activity. Thus, following Roy et al. (2015),
242 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 the dataset, activities varied somewhat predictably by location: for example, eating tended to occur in the kitchen, whereas playtime was the dominant activity in the living room.

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

263 Fine-grained changes in the social information in view

In a third set of analyses, we explored fine-grained changes in the SAYCam infants’ access to social information across development. In these analyses, we capitalize on the fact that OpenPose provides not only face and hand detections but also positional keypoints. In particular, we explored this keypoint dataset with the idea that greater mobility allows older

268 children to be further from their caregivers on average. Thus, younger, less mobile children
269 may tend to see larger faces towards the center of their visual field while older, more mobile
270 children may experience more smaller, more variable views of faces. The same dynamic would
271 be predicted hold for hands as well, as it would be driven by overall differences in distance.

272 Supporting this idea, we found that the averages sizes of the people, faces, and hands
273 in the infant view became smaller over development (Figure 5). This effect was relatively
274 consistent across the three children in the dataset, despite the fact that the three children
275 showed sometimes disparate overall proportions of faces/hands in view. Thus, children may
276 see closer, larger views of people, hands, and faces earlier in development.

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

284 Finally, we also examined whether there were changes in where faces tended to appear
285 in the camera's (and hence, by proxy, the child's) field of view. As expected, faces tended to
286 be located towards the upper field of view, while views of hands were more centrally
287 distributed (see Appendix, Figure C1 for average density distributions). However, we also
288 found evidence that older children tended to see more faces in more variable positions than
289 younger children. Specifically, we examined how variable the horizontal and vertical
290 coordinates were of the faces in the infant view. To do so, we calculated the coefficient of
291 variation of the horizontal (x) and vertical (y) positions of centers of the faces detected by
292 OpenPose (see Figure 6B), and examined changes across age. Faces tended to be more
293 variable in the vertical than their horizontal position (see Figure 6B). We also found that as

294 children got older, they tended to see faces that varied more in their horizontal – but not
295 their vertical position – suggesting that older children might be more likely to see more
296 smaller faces in their periphery (see Figure 6B).

297

General Discussion

298 Here, we analyzed the social information in view in a dense, longitudinal dataset,
299 applying a modern computer-vision model to quantify the hands and faces seen from each of
300 three children’s egocentric perspective from 6 to 32 months of age. First, we found a
301 moderate decrease across age in the proportion of faces in view in the videos, in keeping with
302 previous work (Fausey et al., 2016a; Jayaraman et al., 2015). This finding is particularly
303 notable given that, in previous cross-sectional data, this effect seems to be most strongly
304 driven by infants younger than 4 months of age (e.g., Fausey et al., 2016a; Jayaraman et al.,
305 2015; Sugden et al., 2014) who see both more frequent and more persistent faces (Jayaraman
306 & Smith, 2018).

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

309 We also found an unexpectedly high proportion of hands in the view of infants, even
310 when restricting the field-of-view to the center field-of-view to make the viewpoints
311 comparable to those of headcams used in prior work. Why might this be the case? One idea
312 is that these videos contain the viewpoints of children not only during structured
313 interactions (e.g., play sessions at home or in the lab) but during everyday activities when
314 children may be playing by themselves or simply observing the actions of caregivers and
315 other people in their environment. During these less structured times, caregivers may move
316 about in the vicinity of the child but not interact with them as directly—leading to views
317 where a person and their hands are visible from a distance, but this person’s face may be
318 turned away from the infant or occluded (see examples in Figure 1). Indeed, using the same

319 pose detector on videos from in-lab play sessions, Sanchez et al. (2018) found the opposite
320 trend: slightly fewer hand detections than face detections from 8-16 months of age. Work
321 that directly examines the variability in the social information in view across more vs. less
322 structured activity contexts could further test this idea.

323 A coarse analysis based on the location the videos were filmed in further highlights the
324 variability of the social information in view during different activities, showing differences
325 across locations and between individual children. Within a given, well-defined context—e.g.,
326 mealtime in kitchens—S saw more faces than A, and S saw more faces in the kitchen than in
327 other locations. This variability likely stems from the fact that there are at least three ways
328 to feed a young child: 1) sitting in front of the child, facing them as they sit in a high chair;
329 2) sitting behind the child, holding them as they face outward, and 3) sitting side by side.
330 Each of these positions offer the child differing degrees of visual access to faces and hands.
331 While the social information in view may be variable across children in different activity
332 contexts, these analyses suggest they could be stable within a given child’s day-to-day
333 experience.

334 We also used these detailed pose annotations to explore finer-grained changes in how
335 children experience the faces and hands of their caregivers over development. We found that
336 the faces, hands, and people in the infant view tended to become smaller and that faces
337 tended to be farther away and in more variable horizontal positions, in keeping with prior
338 work examining the sizes of faces in the infant view during the first year of life (Jayaraman
339 et al., 2015). Overall, these data support the idea that the social information in view
340 changes across development as infants become increasingly mobile and independent (???:
341 Fausey et al., 2016b). As children explore the world on their own (Xu, 2019), they may
342 experience fewer close-up interactions with the caregivers and more bouts of play where they
343 are exploring the objects and things in the environment around them.

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

whom, and what data we sample; these choices become central when we attempt to draw conclusions about the regularities of experience. Indeed, while unprecedented in size, this dataset still has many limitations. These videos only represent a small portion of the everyday experience of these three children, all of whom come from relatively privileged households in western societies and thus are not representative in many ways of the global population (Henrich, Heine, & Norenzayan, 2010; Karasik, Tamis-LeMonda, Ossmy, & Adolph, 2018). Any idiosyncrasies in how and when these particular families chose to film these videos also undoubtedly influences the variability seen here, and may contribute to the individual differences between the three children in this dataset. And without eye-tracking data, 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 human annotations in prior work. While OpenPose did not provide annotations that were as accurate as those provided by human annotators, we found relatively consistent results with prior literature, suggesting that the sheer scale and density of the annotations provided by this method may overcome some of its limitations.

In future work, the adaptation of deep neural networks for the infant egocentric view remains a promising avenue for collaboration between computer vision experts and developmental psychologists, and indeed has already yielded new insights about the learning mechanisms needed to build visual representations (Orhan, Gupta, & Lake, 2020; Tsutsui, Chandrasekaran, Reza, Crandall, & Yu, 2020; Zhuang, She, Andonian, Mark, & Yamins, 2020). We propose that the use of novel algorithms with 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

³⁷¹ experience that scaffold learning.

³⁷²

Acknowledgements

³⁷³ Thanks to the creators of the SAYCam dataset who made this work possible and to
³⁷⁴ Alessandro Sanchez for his contributions to the codebase. This work was funded by a Jacobs
³⁷⁵ Foundation Fellowship to MCF, a John Mereck Scholars award to MCF, and NSF #1714726
³⁷⁶ to BLL.

377

References

- 378 Bambach, S., Lee, S., Crandall, D. J., & Yu, C. (2015). Lending a hand: Detecting hands
379 and recognizing activities in complex egocentric interactions. In *Proc. Of the IEEE*
380 *international conference on computer vision* (pp. 1949–1957).
- 381 Bruner, J. (1985). The role of interaction formats in language acquisition. In *Language and*
382 *social situations* (pp. 31–46). Springer.
- 383 Cao, Z., Hidalgo, G., Simon, T., Wei, S.-E., & Sheikh, Y. (2018). OpenPose: Realtime
384 multi-person 2D pose estimation using Part Affinity Fields. In *ArXiv preprint*
385 *arXiv:1812.08008*.
- 386 Clerkin, E. M., Hart, E., Rehg, J. M., Yu, C., & Smith, L. B. (2017). Real-world visual
387 statistics and infants' first-learned object names. *Phil. Trans. R. Soc. B*, 372(1711),
388 20160055.
- 389 Farroni, T., Csibra, G., Simion, F., & Johnson, M. H. (2002). Eye contact detection in
390 humans from birth. *Proceedings of the National Academy of Sciences*, 99(14),
391 9602–9605.
- 392 Fausey, C. M., Jayaraman, S., & Smith, L. B. (2016a). From faces to hands: Changing
393 visual input in the first two years. *Cognition*, 152, 101–107.
- 394 Fausey, C. M., Jayaraman, S., & Smith, L. B. (2016b). From faces to hands: Changing
395 visual input in the first two years. *Cognition*, 152, 101–107.
- 396 Franchak, J. M. (2019). Changing opportunities for learning in everyday life: Infant body
397 position over the first year. *Infancy*, 24(2), 187–209.
- 398 Franchak, J. M., Kretch, K. S., & Adolph, K. E. (2017). See and be seen: Infant-caregiver

- 399 social looking during locomotor free play. *Developmental Science*.
- 400 Franchak, J. M., Kretch, K. S., Soska, K. C., & Adolph, K. E. (2011). Head-mounted eye
401 tracking: A new method to describe infant looking. *Child Development*, 82(6),
402 1738–1750.
- 403 Frank, M. C., Simmons, K., Yurovsky, D., & Pusiol, G. (2013). Developmental and postural
404 changes in children's visual access to faces. In *Proceedings of the 35th annual meeting*
405 *of the cognitive science society* (pp. 454–459).
- 406 Gredeback, G., Theuring, C., Hauf, P., & Kenward, B. (2008). The microstructure of infants'
407 gaze as they view adult shifts in overt attention. *Infancy*, 13(5), 533–543.
- 408 Henrich, J., Heine, S. J., & Norenzayan, A. (2010). Most people are not weird. *Nature*,
409 466(7302), 29–29.
- 410 Iverson, J. M. (2010). Developing language in a developing body: The relationship between
411 motor development and language development. *Journal of Child Language*, 37(2),
412 229–261.
- 413 James, W. (1890). *The principles of psychology* (Vol. 1). Macmillan London.
- 414 Jayaraman, S., Fausey, C. M., & Smith, L. B. (2013). Visual statistics of infants' ordered
415 experiences. *Journal of Vision*, 13(9), 735–735.
- 416 Jayaraman, S., Fausey, C. M., & Smith, L. B. (2015). The faces in infant-perspective scenes
417 change over the first year of life. *PLoS One*.
418 <https://doi.org/10.1371/journal.pone.0123780>
- 419 Jayaraman, S., Fausey, C. M., & Smith, L. B. (2017). Why are faces denser in the visual
420 experiences of younger than older infants? *Developmental Psychology*, 53(1), 38.

- 421 Jayaraman, S., & Smith, L. B. (2018). Faces in early visual environments are persistent not
422 just frequent. *Vision Research*.
- 423 Karasik, L. B., Tamis-LeMonda, C. S., Ossmy, O., & Adolph, K. E. (2018). The ties that
424 bind: Cradling in tajikistan. *PloS One*, 13(10), e0204428.
- 425 Kretch, K. S., Franchak, J. M., & Adolph, K. E. (2014). Crawling and walking infants see
426 the world differently. *Child Development*, 85(4), 1503–1518.
- 427 Luo, C., & Franchak, J. M. (2020). Head and body structure infants' visual experiences
428 during mobile, naturalistic play. *Plos One*, 15(11), e0242009.
- 429 Orhan, A. E., Gupta, V. V., & Lake, B. M. (2020). Self-supervised learning through the eyes
430 of a child. *arXiv Preprint arXiv:2007.16189*.
- 431 Roy, B. C., Frank, M. C., DeCamp, P., Miller, M., & Roy, D. (2015). Predicting the birth of
432 a spoken word. *Proc. Of the National Academy of Sciences*, 112(41), 12663–12668.
- 433 Sanchez, A., Long, B., Kraus, A. M., & Frank, M. C. (2018). Postural developments
434 modulate children's visual access to social information. In *Proceedings of the 40th*
435 *annual conference of the cognitive science society*.
- 436 Simon, T., Joo, H., Matthews, I., & Sheikh, Y. (2017). Hand keypoint detection in single
437 images using multiview bootstrapping. In *CVPR*.
- 438 Smith, L. B., Jayaraman, S., Clerkin, E., & Yu, C. (2018). The developing infant creates a
439 curriculum for statistical learning. *Trends in Cognitive Sciences*, 22(4), 325–336.
- 440 Smith, L. B., Yu, C., & Pereira, A. F. (2011). Not your mother's view: The dynamics of
441 toddler visual experience. *Developmental Science*, 14(1), 9–17.
- 442 Smith, L. B., Yu, C., Yoshida, H., & Fausey, C. M. (2015). Contributions of head-mounted

- 443 cameras to studying the visual environments of infants and young children. *Journal*
444 *of Cognition and Development*, 16(3), 407–419.
- 445 Sugden, N. A., Mohamed-Ali, M. I., & Moulson, M. C. (2014). I spy with my little eye:
446 Typical, daily exposure to faces documented from a first-person infant perspective.
447 *Developmental Psychobiology*, 56(2), 249–261.
- 448 Sullivan, J., Mei, M., Perfors, A., Wojcik, E., & Frank, M. C. (2021). SAYCam: A large,
449 longitudinal audiovisual dataset recorded from the infant's perspective. *Open Mind*.
- 450 Tsutsui, S., Chandrasekaran, A., Reza, M. A., Crandall, D., & Yu, C. (2020). A
451 computational model of early word learning from the infant's point of view. *arXiv*
452 *Preprint arXiv:2006.02802*.
- 453 Xu, F. (2019). Towards a rational constructivist theory of cognitive development.
454 *Psychological Review*, 126(6), 841.
- 455 Yamamoto, H., Sato, A., & Itakura, S. (2020). Transition from crawling to walking changes
456 gaze communication space in everyday infant-parent interaction. *Frontiers in*
457 *Psychology*, 10, 2987.
- 458 Yoshida, H., & Smith, L. B. (2008). What's in view for toddlers? Using a head camera to
459 study visual experience. *Infancy*, 13, 229–248.
- 460 Yu, C., & Smith, L. B. (2013). Joint attention without gaze following: Human infants and
461 their parents coordinate visual attention to objects through eye-hand coordination.
462 *PloS One*, 8(11).
- 463 Zhuang, C., She, T., Andonian, A., Mark, M. S., & Yamins, D. (2020). Unsupervised
464 learning from video with deep neural embeddings. In *Proceedings of the ieee/cvpr*
465 *conference on computer vision and pattern recognition* (pp. 9563–9572).

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

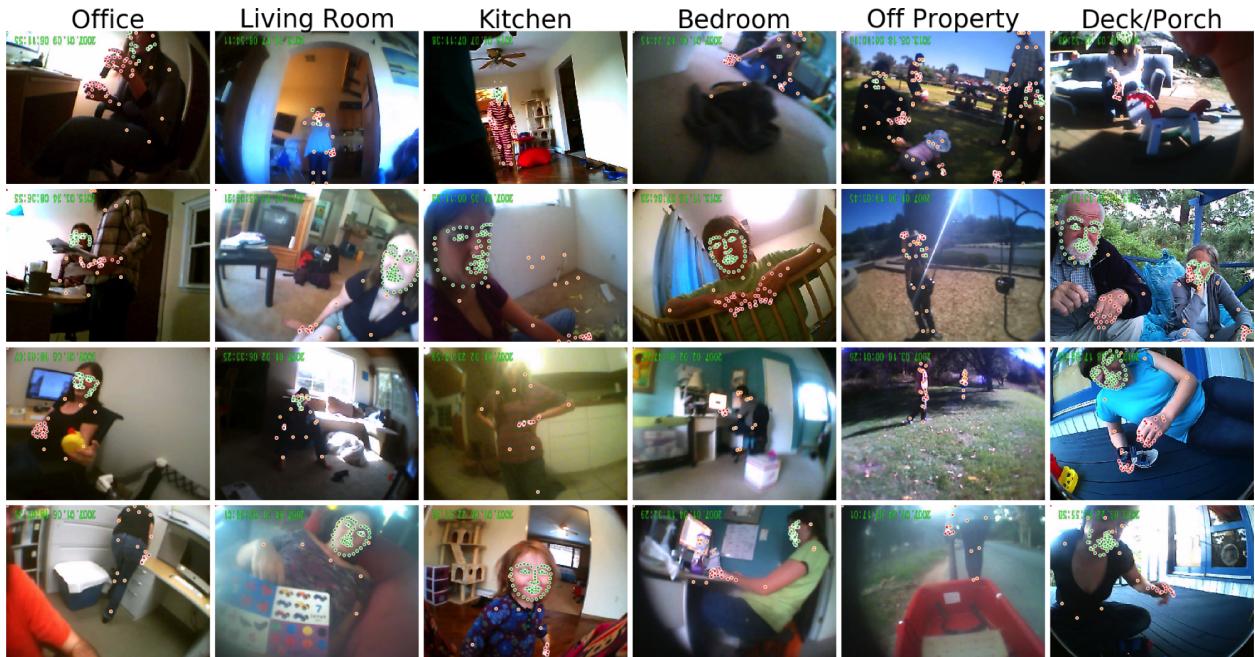


Figure 1. Example frames taken from the dataset, illustrating variability in the infant perspective across different locations. OpenPose detections are shown overlaid on these images (green dots = face, red dots = hands, orange dots = pose).

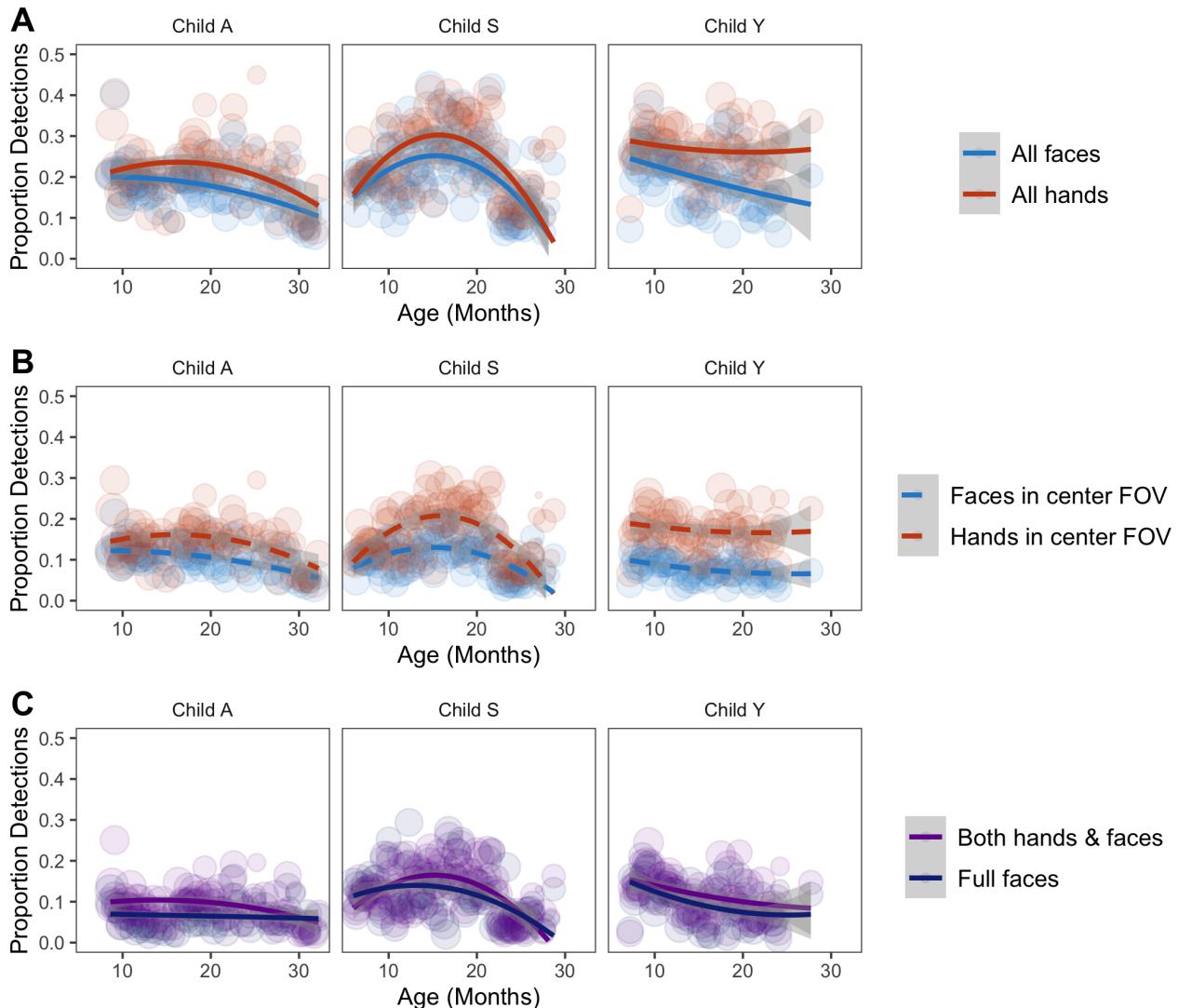


Figure 2. Proportion of frames with (A) All face and hand detections, B. Face/hand detections that fell within the center field-of-view (reducing the contribution of children's own hands) and (C) Face detections that were full faces (e.g., eyes, nose, and mouth all visible) and that co-occurred with 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.

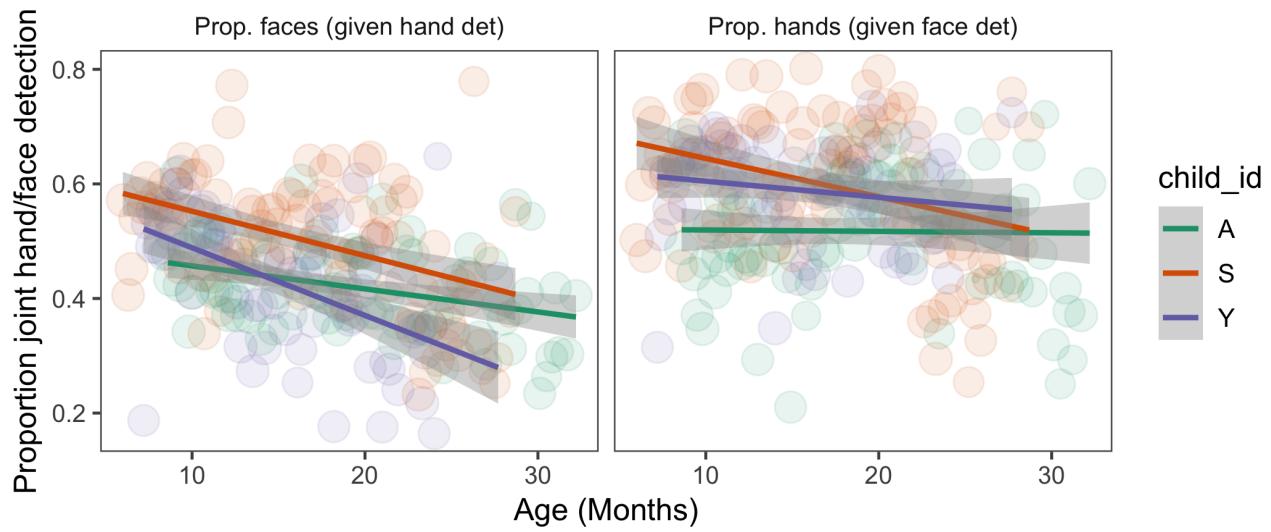


Figure 3. Proportion of joint face and hands detection within frames where hands (left) or faces (right) were detected.

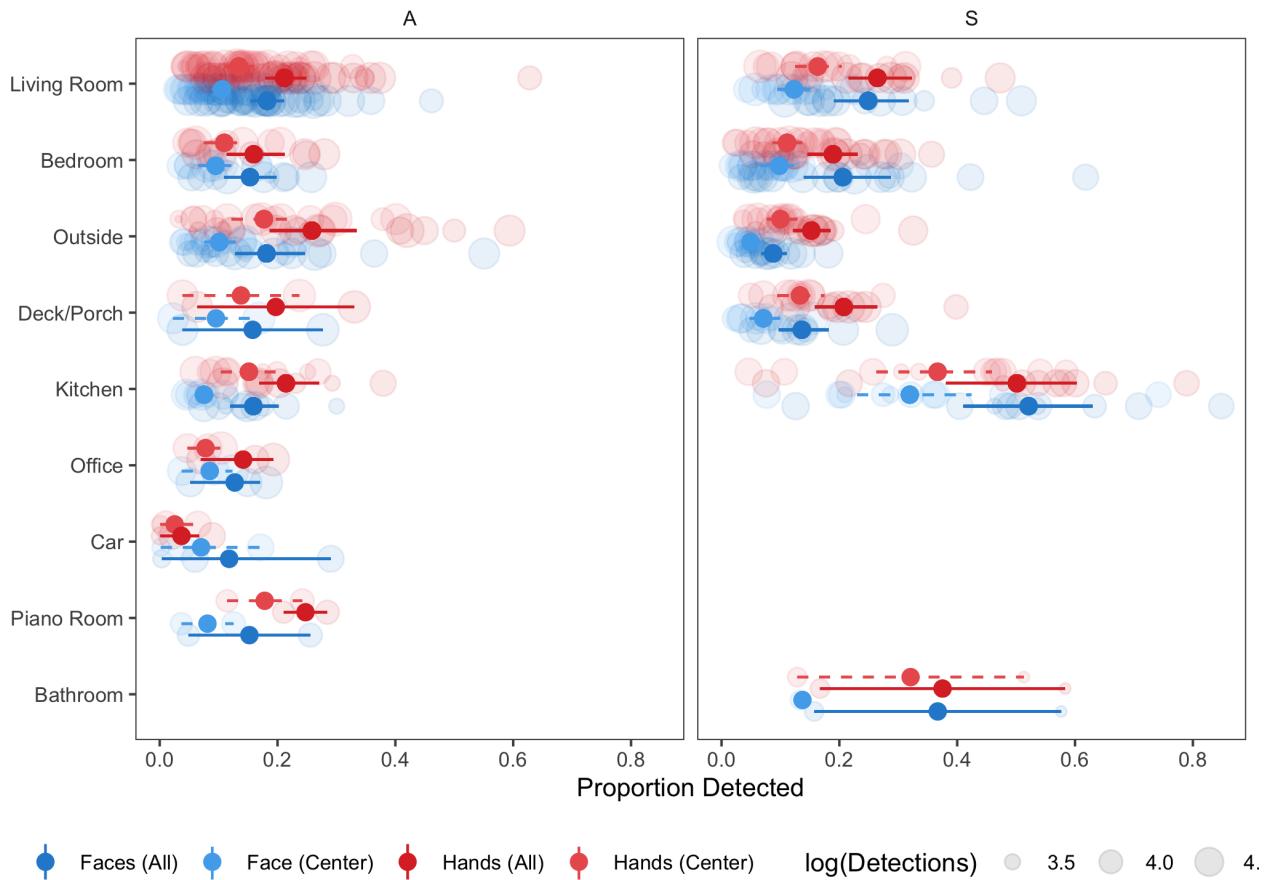


Figure 4. Proportion of faces and hands by location in which egocentric videos were filmed; each panel represents data from an individual child (location annotations were not yet available for Y). Each dot represents data from a week in which videos were filmed and are scaled by the number of frames.

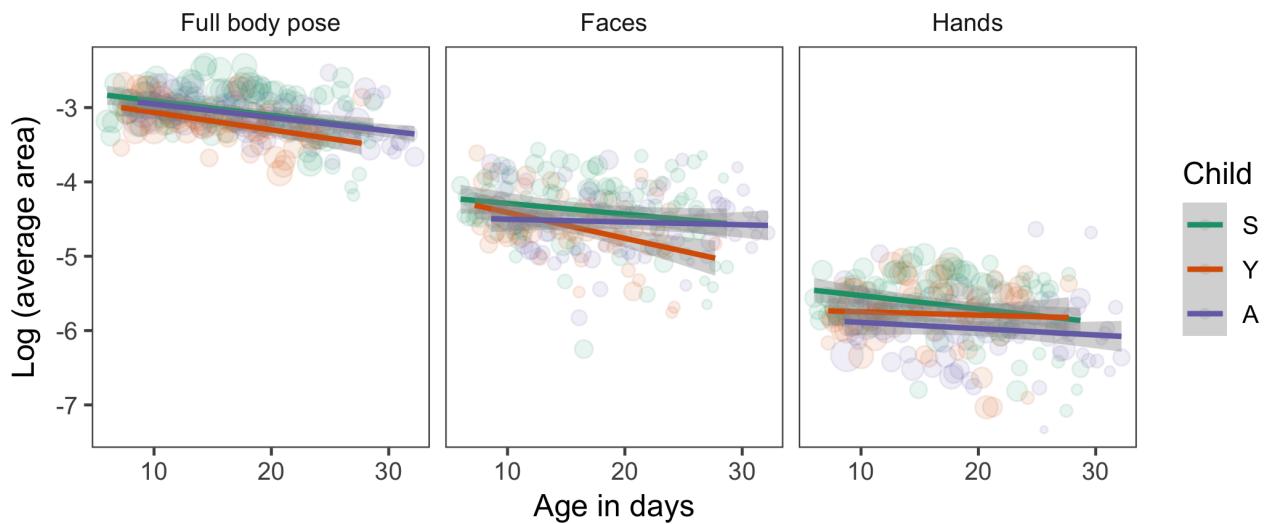


Figure 5. Average size of poses, faces, and hands detected in the dataset between eyes in faces detected as a function of age for each child in the dataset (each color = different child). Data are binned by each week that the videos were filmed and scaled by the number of frames in that age range.

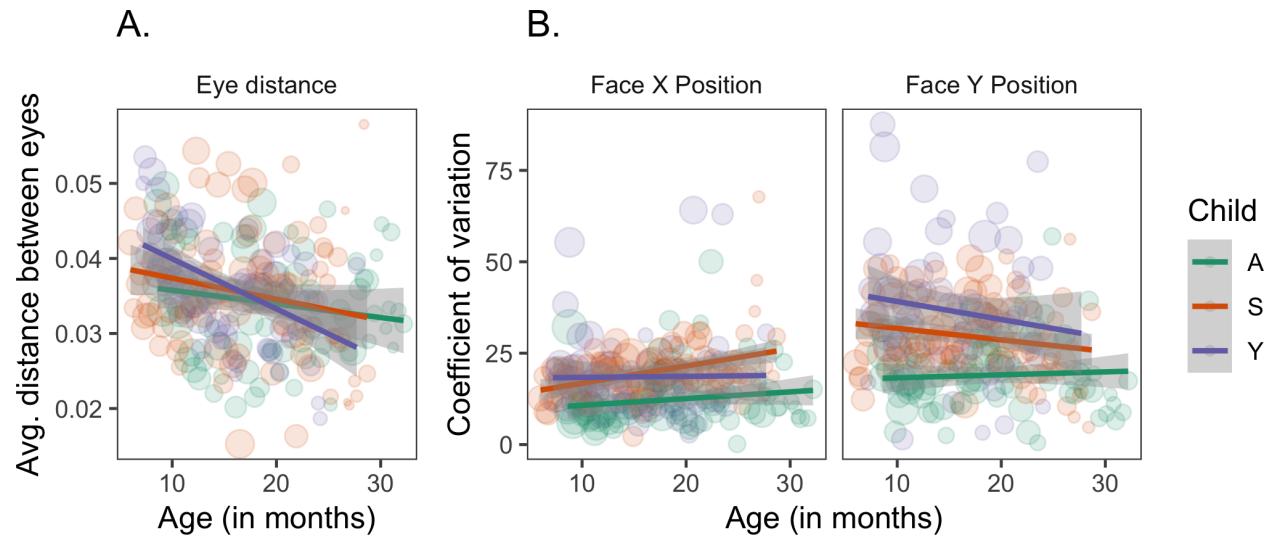


Figure 6. (A) Average distance between eyes and (B) average coefficient of variation for the x and y position of faces detected by OpenPose as a function of each child's age at the time of filming. Data in (A) are restricted to faces where both eyes were detected. Data are binned by each week that the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

Appendix A

Face/hand detections relative to human annotations

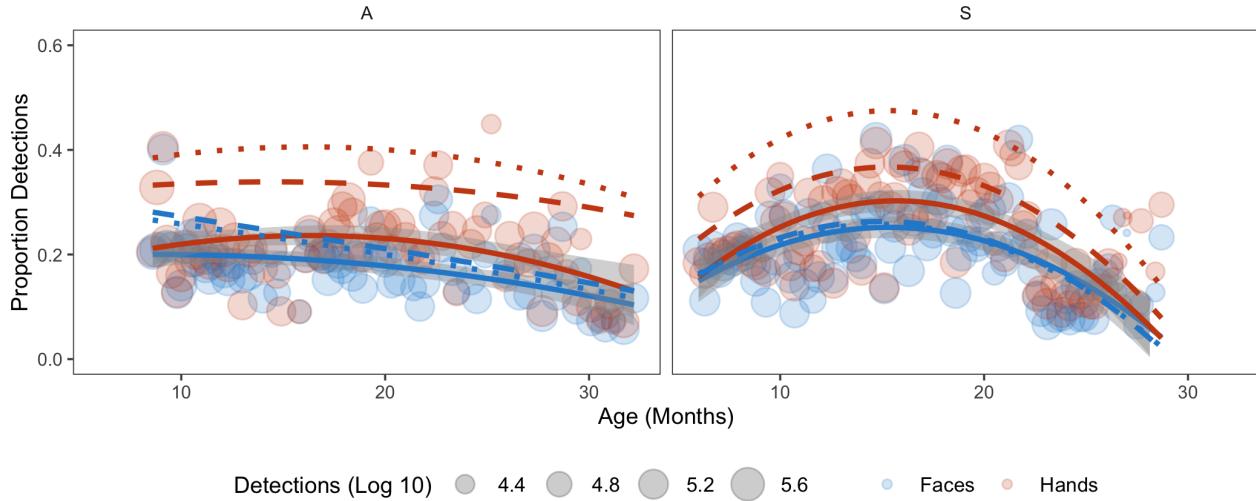
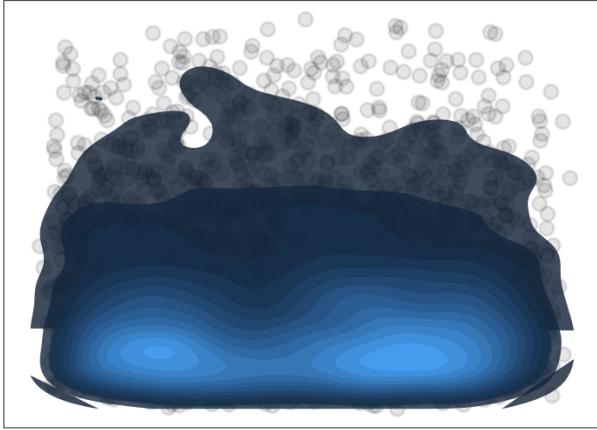


Figure A1. Proportion of faces and hands seen as a function of age for each child in the dataset. Data are binned by each week that the videos were filmed and scaled by the number of frames in that age range. Dashed lines show estimated trend lines from proportion of faces/hands in view when analyzing 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

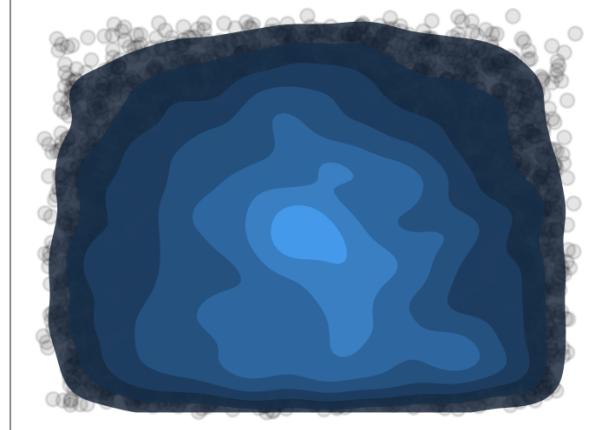


Figure B1. Density estimates for the child (left) and adult (right) hands that were detected in the 24K frame random gold set; each dot represents the center of a bounding box made by an adult participant. Brighter values indicate more detections.

Appendix C

Distribution of faces and hands in the visual field

⁴⁶⁶ We explored where in the visual field children tended to see faces and hands, suspecting that
⁴⁶⁷ these distributions might become wider as children grow older and learn to locomote on their
⁴⁶⁸ own, following preliminary analyses from Frank (2012). As expected, faces tended to appear
⁴⁶⁹ in the upper visual field in contrast to hands, which tended to be more centrally
⁴⁷⁰ located. However, we found little evidence for any changes in the positions of faces and hands
⁴⁷¹ across age, suggesting that this is a relatively stable property of infants' visual environment
⁴⁷² from 6 months of age.

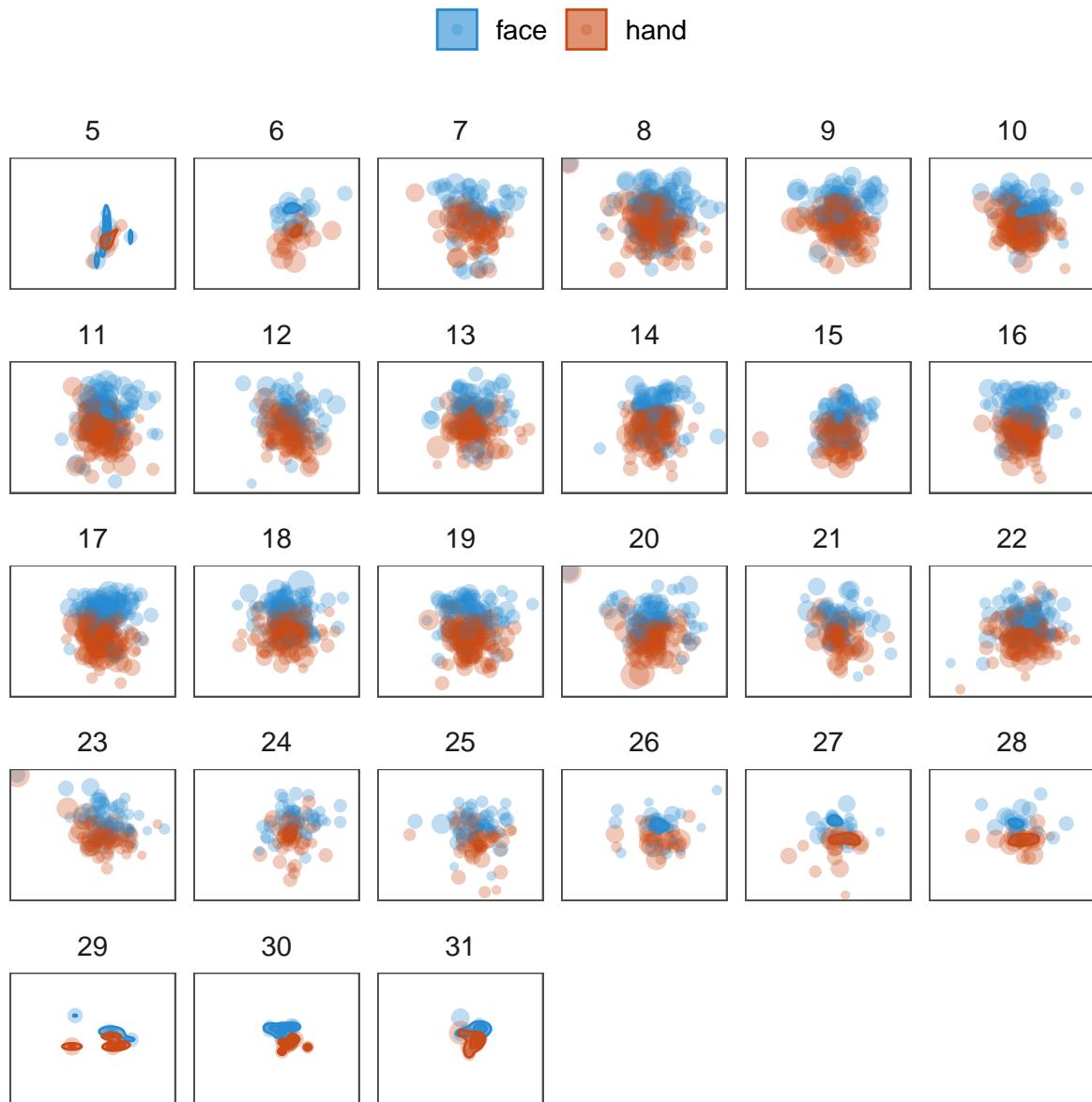


Figure C1. Each panel shows the average position of faces and hands in the visual field; each dot represents the average position from one video within a given age range.