

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 is likely 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 from three children ages 6 to 32 months. To
10 analyze these naturalistic infant egocentric videos, we validated the use of a modern
11 convolutional neural network (OpenPose) for the detection of faces and hands and then
12 applied this model to the entire dataset. First, we found a higher proportion of hands in
13 view than previously reported and a moderate decrease in the proportion of faces in
14 children's view across age. Second, we found substantial variability in the proportion of
15 faces and hands viewed by different children in different locations (e.g., living room
16 vs. kitchen), suggesting that individual activity contexts may shape the social information
17 that infants experience. Third, we found evidence that children may see closer, larger views
18 of people, hands, and faces earlier in development. These analyses provide new insight into
19 the changes in the social information in view across the first few years of life and call for
20 further work that examines their generalizability across populations and their relationship
21 to learning outcomes.

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

23 Word count: 4698

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

26 **Introduction**

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

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

46 Beyond these early months, infants' motor and cognitive abilities mature, leading to
47 vastly different perspectives on the world (Iverson, 2010). For example, children see fewer
48 faces and hands when crawling than walking or sitting (Franchak, 2019; Franchak, Kretch,
49 & Adolph, 2017; Kretch, Franchak, & Adolph, 2014; Luo & Franchak, 2020; Sanchez, Long,

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

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

73 With hundreds of hours of footage (>42M frames), however, this large dataset
74 necessitates a shift to an automated annotation strategy. Indeed, annotation of the frames
75 extracted from egocentric videos has been prohibitively time-consuming, meaning that
76 most frames are typically not inspected, even in the most comprehensive studies. For

example, Fausey et al. (2016a) collected a total of 143 hours of head-mounted camera footage (15.5 million frames), of which one frame every five seconds was hand-annotated (by four coders), totalling 103,383 frames (per coder)—an impressive number of annotations but nonetheless only 0.67% of the collected footage. To address this challenge, we use a modern computer vision model of pose detection to automatically detect the presence of hands and faces from the infant egocentric viewpoint. Specifically, we use OpenPose (Cao, Hidalgo, Simon, Wei, & Sheikh, 2018), a model optimized for jointly detecting human face, body, hand, and foot keypoints that operates well on scenes including multiple people, even if they are partially-occluded (see Figure 1). In prior work examining egocentric videos, OpenPose performed comparably to other modern face detection models (Sanchez et al., 2018).

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

Method

Dataset

The dataset is described in detail in Sullivan et al. (2021); we summarize these details here. Children wore Veho Muvi miniature cameras mounted on a custom camping headlamp harness (“headcams”) at least twice weekly, for approximately one hour per recording session. One weekly session was on the same day each week at a roughly constant

time of day, while the other(s) were chosen arbitrarily at the participating family's discretion. At the time of the recording, all three children were in single-child households. Videos captured by the headcam were 640x480 pixels, and a fisheye lens was attached to the camera to increase the field of view to approximately 109 degrees horizontal x 70 degrees vertical. Videos¹ with technical errors or that were not taken from the egocentric perspective were excluded from the dataset. We analyze 1745 videos, with a total duration of 391.11 hours (>42 million frames).

109 Detection Method

To annotate the millions of frames in SAYCam automatically, we used a pose detector, OpenPose² (Cao et al., 2018; Simon, Joo, Matthews, & Sheikh, 2017). The OpenPose system provides the locations of up to 18 body parts (ears, nose, wrists, etc.) from individual frames. OpenPose relies on a convolutional neural network for initial anatomical detection. It then uses part affinity fields for part association to produce a series of body part candidates. Once these body part candidates are matched to a single individual in the frame, they are finally assembled into a pose. While in this study we only measured face and hand presence, the entire set of pose information from an individual was used to determine the presence of a face/hand, making the process much more robust to occlusion than methods optimized to detect *only* faces or hands. Of course, these face/hand detections are nevertheless reliant on the detection of at least a partial pose, so some very up-close views of faces/hands may still go undetected.

¹ All 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
124 accuracy of these detections relative to human annotations of 24,000 frames selected
125 uniformly at random from videos of two children (S and A); 24000 frames sampled from
126 allocentric videos were excluded, and these videos were also excluded from the other
127 analyses. Frames were jointly annotated for the presence of faces and hands by one author.
128 A second set of coders recruited via AMT (Amazon Mechanical Turk) additionally
129 annotated 3150 frames; 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
146 hand detections showed moderately good precision, face detections were overall slightly
147 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 high-confidence face detections was 0.41, with a precision of 0.95 and recall of 0.26; for high-confidence hand detections, the F-score was 0.18, with a precision of 0.97 and recall of 0.10).

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

Finally, we examined whether the precision, recall, and F-score for hands and faces varied with age or child, and did not find substantial variation. Thus, while OpenPose was trained on photographs from the adult perspective, this model still generalized relatively well to the egocentric infant viewpoint with no fine-tuning or post-processing of the detections. As these detections were imperfect compared to human annotators, fine-tuning these models to better optimize for the infant viewpoint remains an open avenue for future work. Standard computer vision models are rarely trained on the egocentric viewpoint, and we suspect that training these models on more naturalistic data may lead to more robust, 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 the addition of quadratic terms relating children's age to the proportion
184 of faces/hands detected would provide better fit to the data than linear terms alone, and
185 found that this was true in both cases (see Tables 1 & 2), though the linear term was also
186 significant for faces. Thus, these exploratory results point towards the idea that some
187 children may experience overall more social information in view in the second year of life.

188 However, the most striking result from these analyses is a much greater overall
189 proportion of hands in view than has previously been reported (Fausey et al., 2016a). We
190 found this observation to be true across all ages, in all three children, and regardless of
191 whether we analyzed human annotations (on the 24K random subset, see dotted lines in
192 Appendix Figure A1) or OpenPose annotations on the entire dataset (see Figure 2A). This
193 finding is notable especially given that OpenPose showed relatively low recall for hands,
194 indicating that our measurements may in fact be an underestimate of the proportion of

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

⁴ 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.

195 hands in view. In fact, analysis of the human gold standard annotations revealed a much
196 higher proportion of hands relative to faces than the automated annotations.

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

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

215 As we found that children's hands tended to occur in the lower visual field (see
216 Figure Appendix B1), we thus re-analyzed the entire dataset while restricting our analysis
217 to the center field of view, decreasing the proportion of hand detections from 24% to 16%,
218 and 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. To do
226 so, we calculated the number of frames in which infants saw faces and hands together
227 relative to overall proportions of faces/hands that were detected for each child and age
228 range. Faces and hands were jointly present in 11.50 percent of frames (see face
229 hand-occurrences across age in Figure 2C). As shown in Figure 3, all three infants were
230 more likely to see hands independently – without the presence of a face – than they were
231 likely to see faces independently. That is, generally speaking when a face was present, a
232 hand also 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. Some
238 activities seem likely to be characterized by a much higher proportion of faces (e.g., diaper
239 changes) than others (e.g., a car trip). Following this same idea, perhaps other activities
240 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
243 these annotations for a subset of the SAYCam videos. Of the 1745 videos in the dataset,
244 639 were annotated for the location or locations they were filmed in. These location
245 annotations were only available for two children, S and A. Annotated locations mostly

246 consisted of rooms of the house (e.g., “living room”) but also included some other locations
247 (e.g., “car,” “outside”). Of this set, 296 videos were filmed in only a single location (e.g.,
248 the location label did not change within the video), representing 17 percent of the dataset
249 and over 5 million frames. In our viewing of the SAYCam videos and in other annotations
250 available with the dataset, activities varied somewhat predictably by location: for example,
251 eating tended to occur in the kitchen, whereas playtime was the dominant activity in the
252 living room.

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

264 Fine-grained changes in the social information in view

265 In a third set of analyses, we explored fine-grained changes in the SAYCam infants’
266 access to social information across development. In these analyses, we capitalize on the fact
267 that OpenPose provides not only face and hand detections but also positional keypoints. In
268 particular, we explored this keypoint dataset with the idea that greater mobility allows
269 older children to be further from their caregivers on average. Thus, younger, less mobile
270 children may tend to see larger faces towards the center of their visual field while older,
271 more mobile children may experience more smaller, more variable views of faces. The same

272 dynamic would be predicted hold for hands as well, as it would be driven by overall
273 differences in distance.

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

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

286 Finally, we also examined whether there were changes in where faces tended to
287 appear in the camera’s (and hence, by proxy, the child’s) field of view. As expected, faces
288 tended to be located towards the upper field of view, while views of hands were more
289 centrally distributed (see Appendix, Figure C1 for average density distributions). However,
290 we also found evidence that older children tended to see more faces in more variable
291 positions than younger children. Specifically, we examined how variable the horizontal and
292 vertical coordinates were of the faces in the infant view. To do so, we calculated the
293 coefficient of variation of the horizontal (x) and vertical (y) positions of centers of the faces
294 detected by OpenPose (see Figure 6B), and examined changes across age. Faces tended to
295 be more variable in the vertical than their horizontal position (see Figure 6B). We also
296 found that as children got older, they tended to see faces that varied more in their
297 horizontal – but not their vertical position – suggesting that older children might be more

298 likely to see more smaller faces in their periphery (see Figure 6B).

299

General Discussion

300 Here, we analyzed the social information in view in a dense, longitudinal dataset,
301 applying a modern computer-vision model to quantify the hands and faces seen from each
302 of three children's egocentric perspective from 6 to 32 months of age.

303 First, we found a moderate decrease across age in the proportion of faces in view in
304 the videos, in keeping with previous work (Fausey et al., 2016a; Jayaraman et al., 2015).
305 This finding is particularly notable given that, in previous cross-sectional data, this effect
306 seems to be most strongly driven by infants younger than 4 months of age (e.g., Fausey et
307 al., 2016a; Jayaraman et al., 2015; Sugden et al., 2014) who see both more frequent and
308 more persistent faces (Jayaraman & Smith, 2018). We also found this to be true when
309 restricting our analyses to full-field faces, suggesting this effect is not driven by a
310 concurrent shift from more full-view to partial-views of faces.

311 We also found an unexpectedly high proportion of hands in infants' view, even when
312 restricting the field-of-view to the center field-of-view to make the viewpoints comparable
313 to those of headcams used in prior work. Why might this be the case? One idea is that
314 these videos contain the viewpoints of children not only during structured interactions
315 (e.g., play sessions at home or in the lab) but during everyday activities when children may
316 be playing by themselves or simply observing the actions of caregivers and other people in
317 their environment. During these less structured times, caregivers may move about in the
318 vicinity of the child but not interact with them as directly – leading to views where a
319 person and their hands are visible from a distance, but this person's face may be turned
320 away from the infant or occluded (see examples in Figure 1). Indeed, using the same pose
321 detector on videos from in-lab play sessions, Sanchez et al. (2018) found the opposite
322 trend: slightly fewer hand detections than face detections from 8-16 months of age. Work

323 that directly examines the variability in the social information in view across more vs. less
324 structured activity contexts could further test this idea.

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

336 We also used these detailed pose annotations to explore finer-grained changes in how
337 children experience the faces and hands of their caregivers over development. We found
338 that the faces, hands, and people in the infant view tended to become smaller and that
339 faces tended to be farther away and in more variable horizontal positions, in keeping with
340 prior work examining the sizes of faces in the infant view during the first year of life
341 (Jayaraman et al., 2015). Overall, these data support the idea that the social information
342 in view changes across development as infants become increasingly mobile and independent
343 (Fausey et al., 2016b; Franchak et al., 2017). As children explore the world on their own
344 (Xu, 2019), they may experience fewer close-up interactions with their caregivers and more
345 bouts of play where they are exploring the objects in their environment.

346 More broadly, however, these analyses underscore the importance of how, when, from
347 whom, and what data we sample; these choices become central when we attempt to draw
348 conclusions about the regularities of experience. Indeed, while unprecedented in size, this

349 dataset still has many limitations. These videos only represent a small portion of the
350 everyday experience of these three children, all of whom come from relatively privileged
351 households in western societies and thus are not representative in many ways of the global
352 population (Henrich, Heine, & Norenzayan, 2010; Karasik, Tamis-LeMonda, Ossmy, &
353 Adolph, 2018). Any idiosyncrasies in how and when these particular families chose to film
354 these videos also undoubtedly influenced the variability seen here, and may contribute to
355 the individual differences between the three children in this dataset. And without
356 eye-tracking data, we do not know the extent to which children are attending to the social
357 information in view.

358 Nonetheless, we believe that these advances in datasets and methodologies represent
359 a step in the right direction. The present paper demonstrates the feasibility of using a
360 modern computer vision model to annotate the entirety of a very large dataset (here, >42M
361 million frames) for the presence and size of people, hands, and faces, representing orders of
362 magnitude more data relative to human annotations in prior work. While OpenPose did
363 not provide annotations that were as accurate as those provided by human annotators, we
364 found relatively consistent results with prior literature, suggesting that the sheer scale and
365 density of the annotations provided by this method may overcome some of its limitations.

366 In future work, the adaptation of deep neural networks for the infant egocentric view
367 remains a promising avenue for collaboration between computer vision experts and
368 developmental psychologists. Indeed, this combination has already yielded new insights
369 about the learning mechanisms needed to build visual representations (Orhan, Gupta, &
370 Lake, 2020; Tsutsui, Chandrasekaran, Reza, Crandall, & Yu, 2020; Zhuang, She, Andonian,
371 Mark, & Yamins, 2020). We propose that the use of novel algorithms with large-scale
372 analysis of dense datasets – collected with different fields of view, cameras, and from many
373 different laboratories – will lead to generalizable conclusions about the regularities of infant
374 experience that scaffold learning.

375

Acknowledgements

376 Thanks to the creators of the SAYCam dataset who made this work possible and to
377 Alessandro Sanchez for his contributions to the codebase. This work was funded by a
378 Jacobs Foundation Fellowship to MCF, a John Merck Scholars award to MCF, and NSF
379 #1714726 to BLL.

380

References

- 381 Bambach, S., Lee, S., Crandall, D. J., & Yu, C. (2015). Lending a hand: Detecting hands
382 and recognizing activities in complex egocentric interactions. In *Proc. Of the IEEE*
383 *international conference on computer vision* (pp. 1949–1957).
- 384 Bruner, J. (1985). The role of interaction formats in language acquisition. In *Language and*
385 *social situations* (pp. 31–46). Springer.
- 386 Cao, Z., Hidalgo, G., Simon, T., Wei, S.-E., & Sheikh, Y. (2018). OpenPose: Realtime
387 multi-person 2D pose estimation using Part Affinity Fields. In *ArXiv preprint*
388 *arXiv:1812.08008*.
- 389 Clerkin, E. M., Hart, E., Rehg, J. M., Yu, C., & Smith, L. B. (2017). Real-world visual
390 statistics and infants' first-learned object names. *Phil. Trans. R. Soc. B,*
391 *372*(1711), 20160055.
- 392 Farroni, T., Csibra, G., Simion, F., & Johnson, M. H. (2002). Eye contact detection in
393 humans from birth. *Proceedings of the National Academy of Sciences*, *99*(14),
394 9602–9605.
- 395 Fausey, C. M., Jayaraman, S., & Smith, L. B. (2016a). From faces to hands: Changing
396 visual input in the first two years. *Cognition*, *152*, 101–107.
- 397 Fausey, C. M., Jayaraman, S., & Smith, L. B. (2016b). From faces to hands: Changing
398 visual input in the first two years. *Cognition*, *152*, 101–107.
- 399 Franchak, J. M. (2019). Changing opportunities for learning in everyday life: Infant body
400 position over the first year. *Infancy*, *24*(2), 187–209.
- 401 Franchak, J. M., Kretch, K. S., & Adolph, K. E. (2017). See and be seen: Infant–caregiver
402 social looking during locomotor free play. *Developmental Science*.
- 403 Franchak, J. M., Kretch, K. S., Soska, K. C., & Adolph, K. E. (2011). Head-mounted eye
404 tracking: A new method to describe infant looking. *Child Development*, *82*(6),

- 405 1738–1750.
- 406 Frank, M. C., Simmons, K., Yurovsky, D., & Pusiol, G. (2013). Developmental and
407 postural changes in children's visual access to faces. In *Proceedings of the 35th*
408 *annual meeting of the cognitive science society* (pp. 454–459).
- 409 Gredeback, G., Theuring, C., Hauf, P., & Kenward, B. (2008). The microstructure of
410 infants' gaze as they view adult shifts in overt attention. *Infancy*, 13(5), 533–543.
- 411 Henrich, J., Heine, S. J., & Norenzayan, A. (2010). Most people are not weird. *Nature*,
412 466(7302), 29–29.
- 413 Iverson, J. M. (2010). Developing language in a developing body: The relationship between
414 motor development and language development. *Journal of Child Language*, 37(2),
415 229–261.
- 416 James, W. (1890). *The principles of psychology* (Vol. 1). Macmillan London.
- 417 Jayaraman, S., Fausey, C. M., & Smith, L. B. (2013). Visual statistics of infants' ordered
418 experiences. *Journal of Vision*, 13(9), 735–735.
- 419 Jayaraman, S., Fausey, C. M., & Smith, L. B. (2015). The faces in infant-perspective
420 scenes change over the first year of life. *PLoS One*.
421 <https://doi.org/10.1371/journal.pone.0123780>
- 422 Jayaraman, S., Fausey, C. M., & Smith, L. B. (2017). Why are faces denser in the visual
423 experiences of younger than older infants? *Developmental Psychology*, 53(1), 38.
- 424 Jayaraman, S., & Smith, L. B. (2018). Faces in early visual environments are persistent not
425 just frequent. *Vision Research*.
- 426 Karasik, L. B., Tamis-LeMonda, C. S., Ossmy, O., & Adolph, K. E. (2018). The ties that
427 bind: Cradling in tajikistan. *PloS One*, 13(10), e0204428.
- 428 Kretch, K. S., Franchak, J. M., & Adolph, K. E. (2014). Crawling and walking infants see
429 the world differently. *Child Development*, 85(4), 1503–1518.

- 430 Luo, C., & Franchak, J. M. (2020). Head and body structure infants' visual experiences
431 during mobile, naturalistic play. *Plos One*, 15(11), e0242009.
- 432 Orhan, A. E., Gupta, V. V., & Lake, B. M. (2020). Self-supervised learning through the
433 eyes of a child. *arXiv Preprint arXiv:2007.16189*.
- 434 Roy, B. C., Frank, M. C., DeCamp, P., Miller, M., & Roy, D. (2015). Predicting the birth of
435 a spoken word. *Proc. Of the National Academy of Sciences*, 112(41), 12663–12668.
- 436 Sanchez, A., Long, B., Kraus, A. M., & Frank, M. C. (2018). Postural developments
437 modulate children's visual access to social information. In *Proceedings of the 40th*
438 *annual conference of the cognitive science society*.
- 439 Simon, T., Joo, H., Matthews, I., & Sheikh, Y. (2017). Hand keypoint detection in single
440 images using multiview bootstrapping. In *CVPR*.
- 441 Smith, L. B., Jayaraman, S., Clerkin, E., & Yu, C. (2018). The developing infant creates a
442 curriculum for statistical learning. *Trends in Cognitive Sciences*, 22(4), 325–336.
- 443 Smith, L. B., Yu, C., & Pereira, A. F. (2011). Not your mother's view: The dynamics of
444 toddler visual experience. *Developmental Science*, 14(1), 9–17.
- 445 Smith, L. B., Yu, C., Yoshida, H., & Fausey, C. M. (2015). Contributions of head-mounted
446 cameras to studying the visual environments of infants and young children. *Journal*
447 *of Cognition and Development*, 16(3), 407–419.
- 448 Sugden, N. A., Mohamed-Ali, M. I., & Moulson, M. C. (2014). I spy with my little eye:
449 Typical, daily exposure to faces documented from a first-person infant perspective.
450 *Developmental Psychobiology*, 56(2), 249–261.
- 451 Sullivan, J., Mei, M., Perfors, A., Wojcik, E., & Frank, M. C. (2021). SAYCam: A large,
452 longitudinal audiovisual dataset recorded from the infant's perspective. *Open Mind*.
- 453 Tsutsui, S., Chandrasekaran, A., Reza, M. A., Crandall, D., & Yu, C. (2020). A
454 computational model of early word learning from the infant's point of view. *arXiv*

455 *Preprint arXiv:2006.02802.*

456 Xu, F. (2019). Towards a rational constructivist theory of cognitive development.

457 *Psychological Review, 126*(6), 841.

458 Yamamoto, H., Sato, A., & Itakura, S. (2020). Transition from crawling to walking changes
459 gaze communication space in everyday infant-parent interaction. *Frontiers in*
460 *Psychology, 10*, 2987.

461 Yoshida, H., & Smith, L. B. (2008). What's in view for toddlers? Using a head camera to
462 study visual experience. *Infancy, 13*, 229–248.

463 Yu, C., & Smith, L. B. (2013). Joint attention without gaze following: Human infants and
464 their parents coordinate visual attention to objects through eye-hand coordination.
465 *PloS One, 8*(11).

466 Zhuang, C., She, T., Andonian, A., Mark, M. S., & Yamins, D. (2020). Unsupervised
467 learning from video with deep neural embeddings. In *Proceedings of the ieee/cvpr*
468 *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

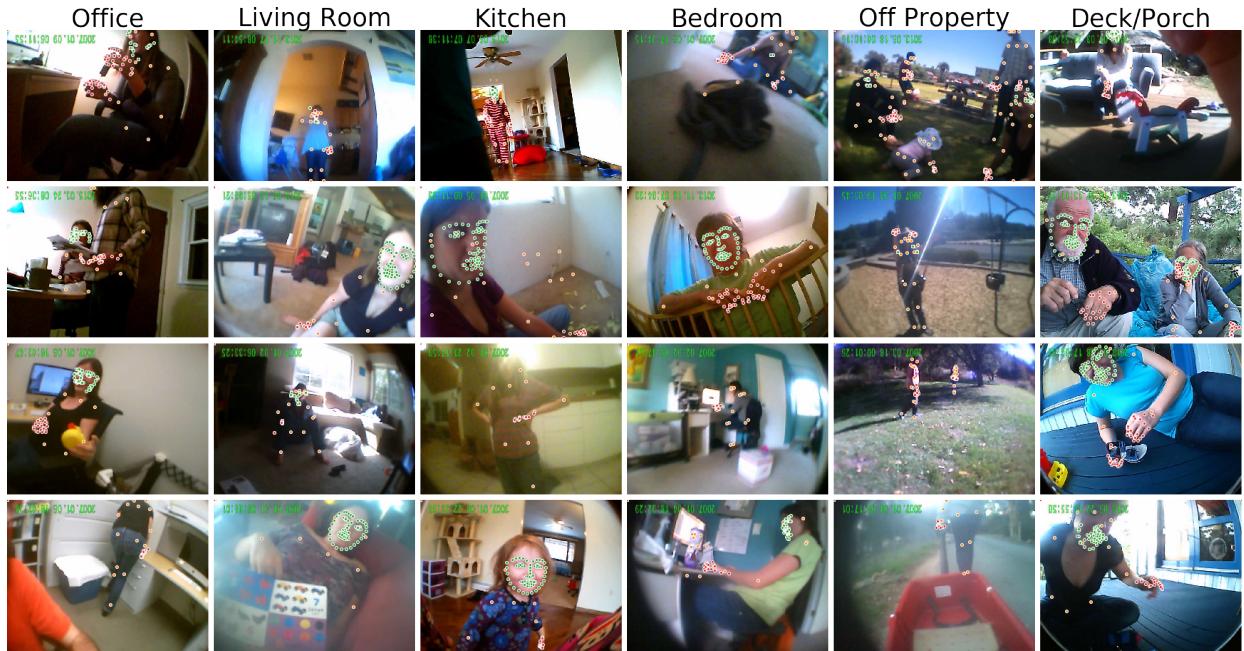


Figure 1. Example frames taken from the dataset, illustrating variability in the infant perspective across different locations. OpenPose detections are shown overlaid on these images (green dots = face, red dots = hands, orange dots = pose).

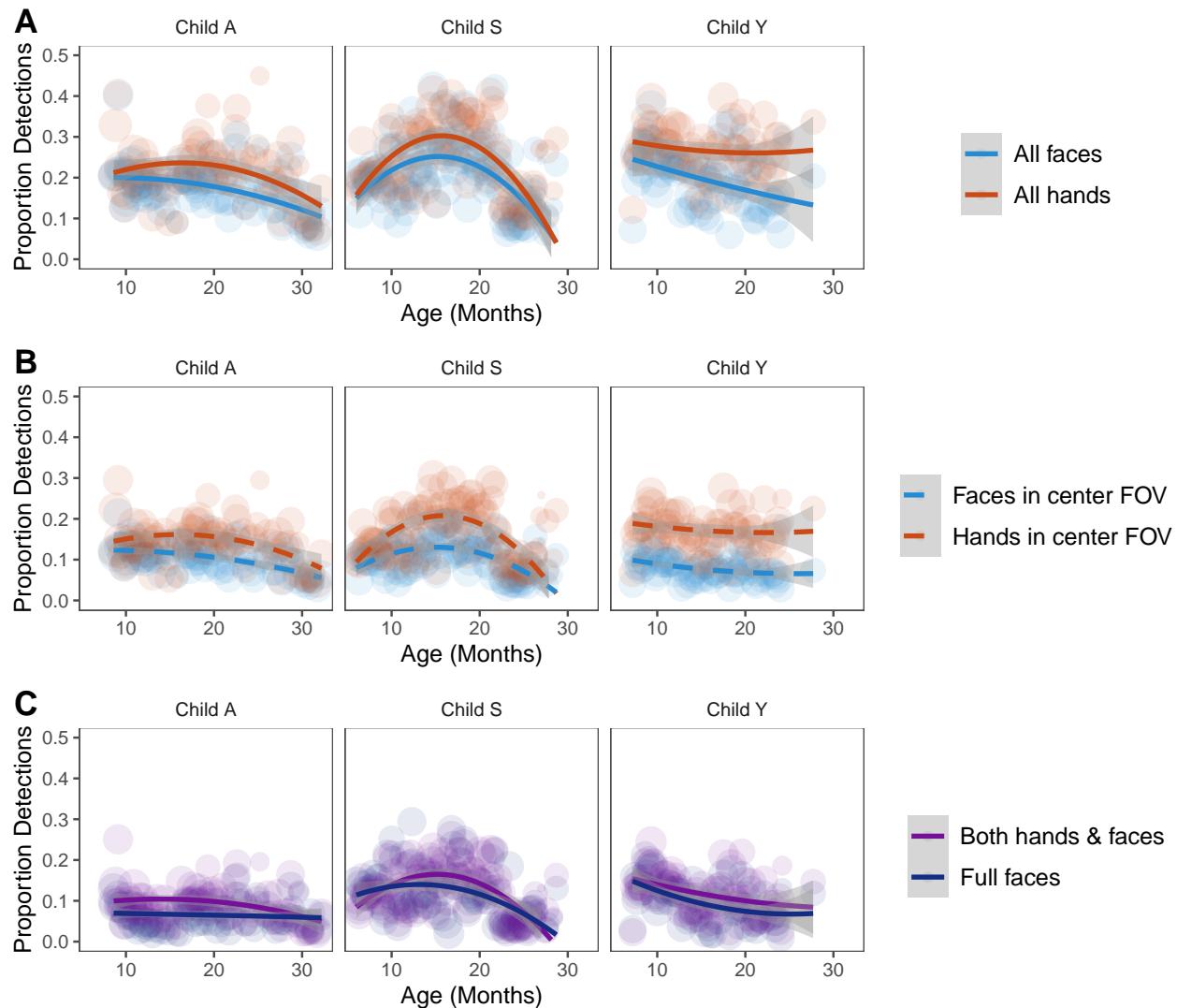


Figure 2. Proportion of frames with (A) All face and hand detections, (B) Face/hand detections that fell within the center field-of-view (reducing the contribution of children's own hands) and (C) Face detections that were full faces (e.g., eyes, nose, and mouth all visible) and that co-occurred with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

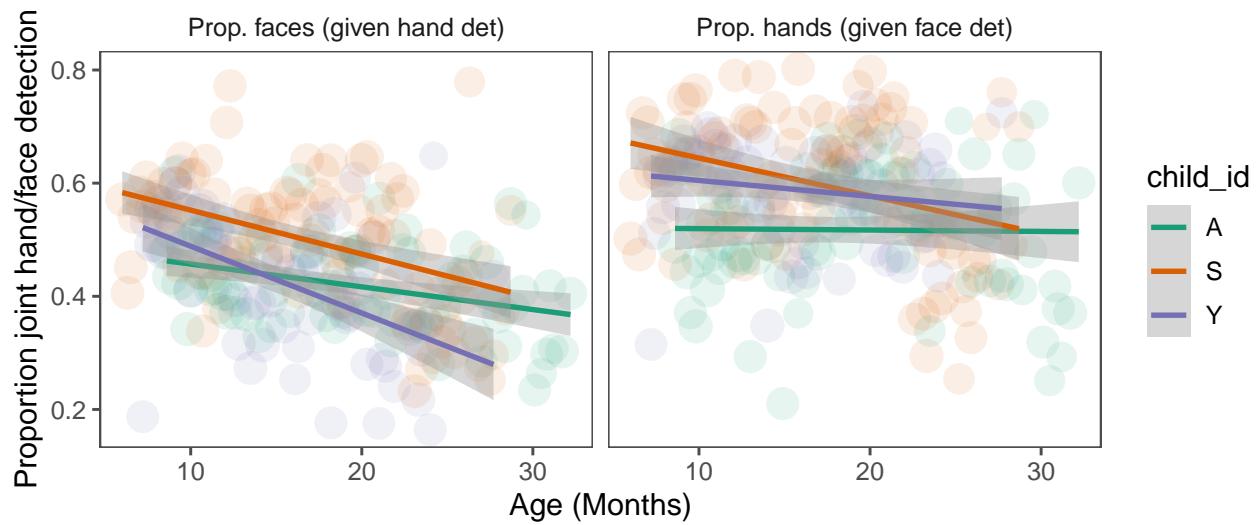


Figure 3. Proportion of joint face and hands detection within frames where hands (left) or faces (right) were detected.

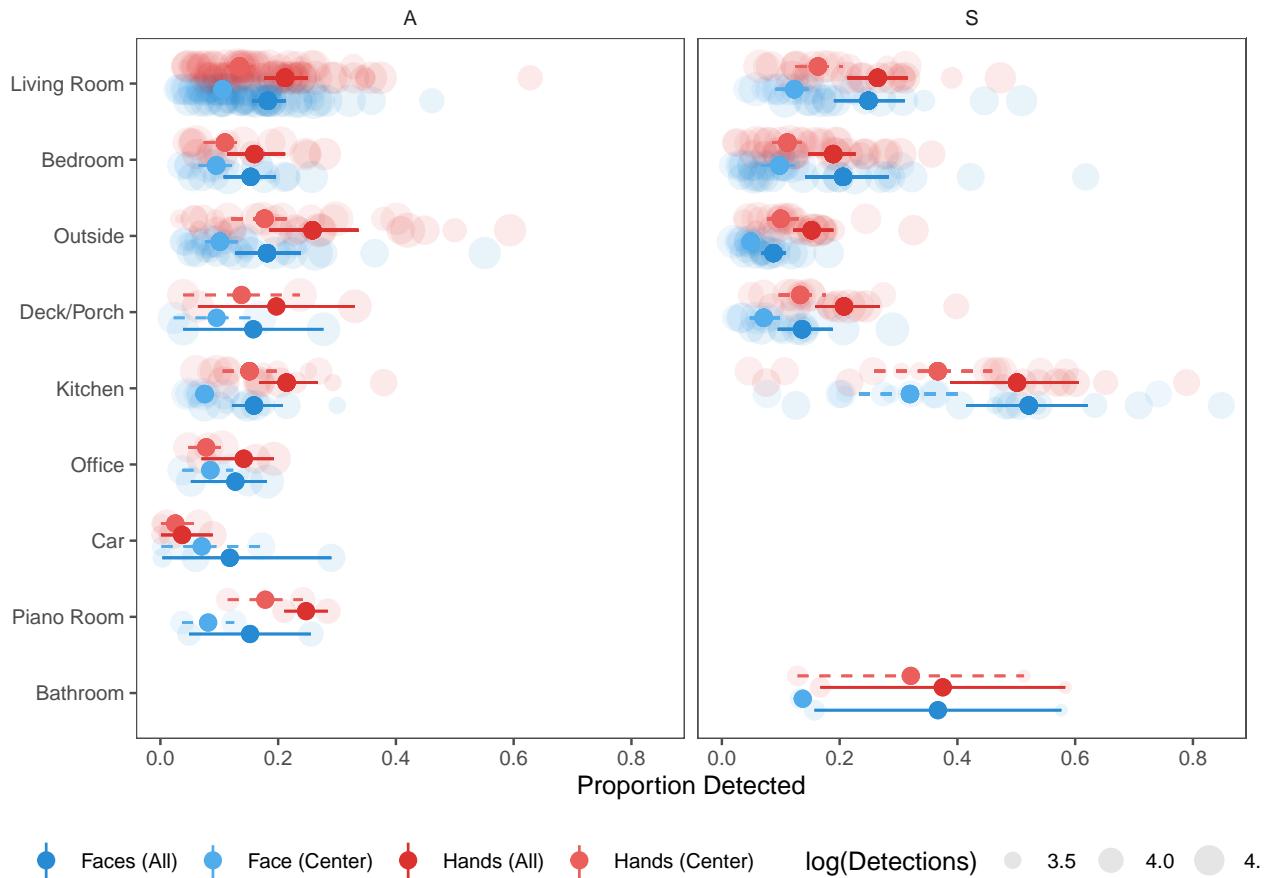


Figure 4. Proportion of faces and hands by location in which egocentric videos were filmed; each panel represents data from an individual child (location annotations were not yet available for Y). Each dot represents data from a week in which videos were filmed and are scaled by the number of frames.

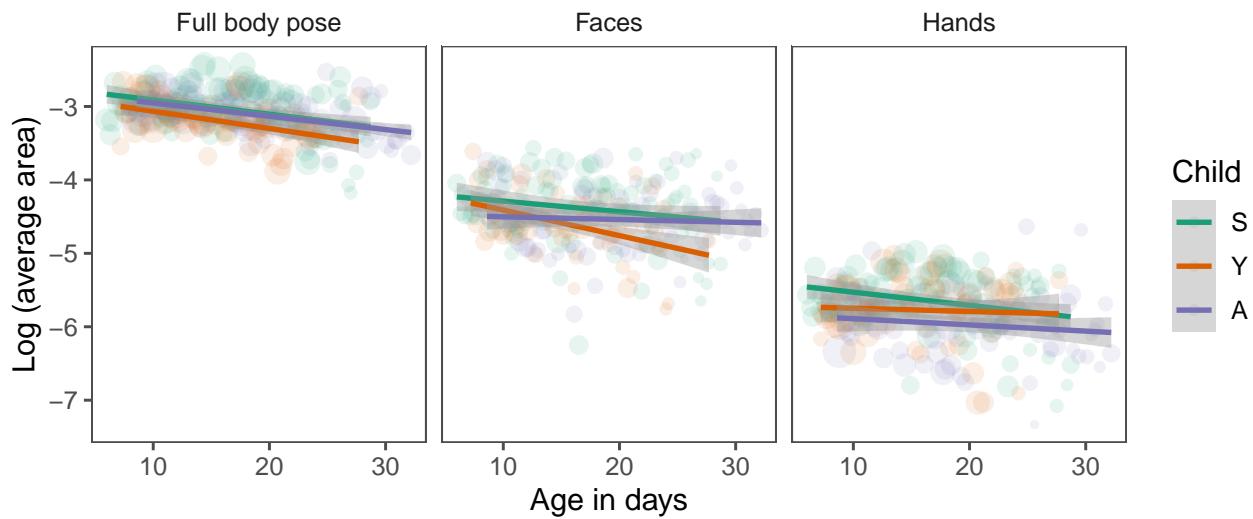


Figure 5. Average size of poses, faces, and hands detected in the dataset between eyes in faces detected as a function of age for each child in the dataset (each color = different child). Data are binned by each week that the videos were filmed and scaled by the number of frames in that age range.

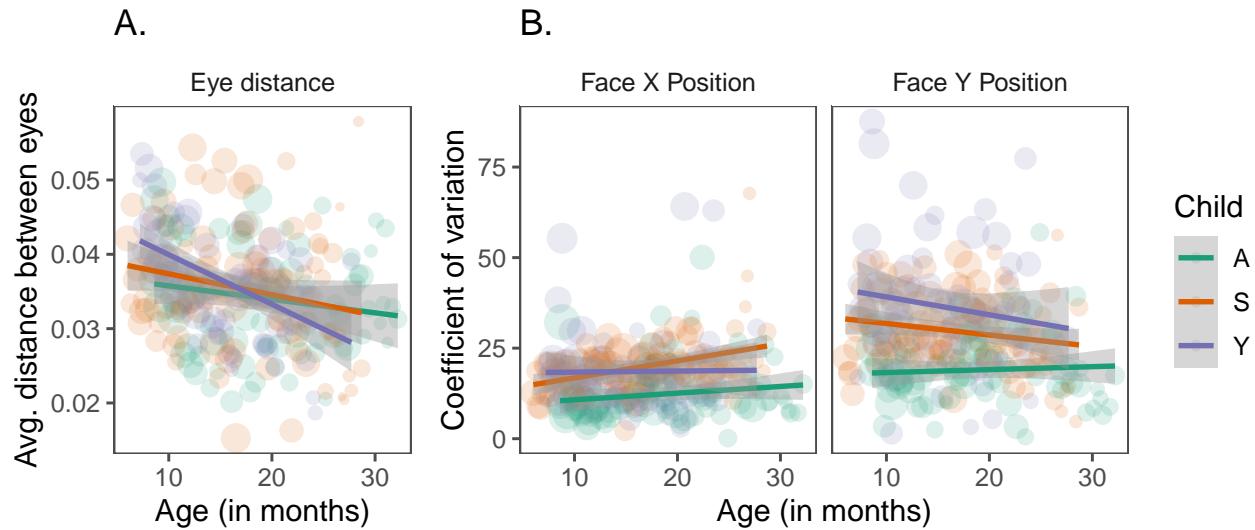


Figure 6. (A) Average distance between eyes and (B) average coefficient of variation for the x and y position of faces detected by OpenPose as a function of each child's age at the time of filming. Data in (A) are restricted to faces where both eyes were detected. Data are binned by each week that the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

Appendix A

Face/hand detections relative to human annotations

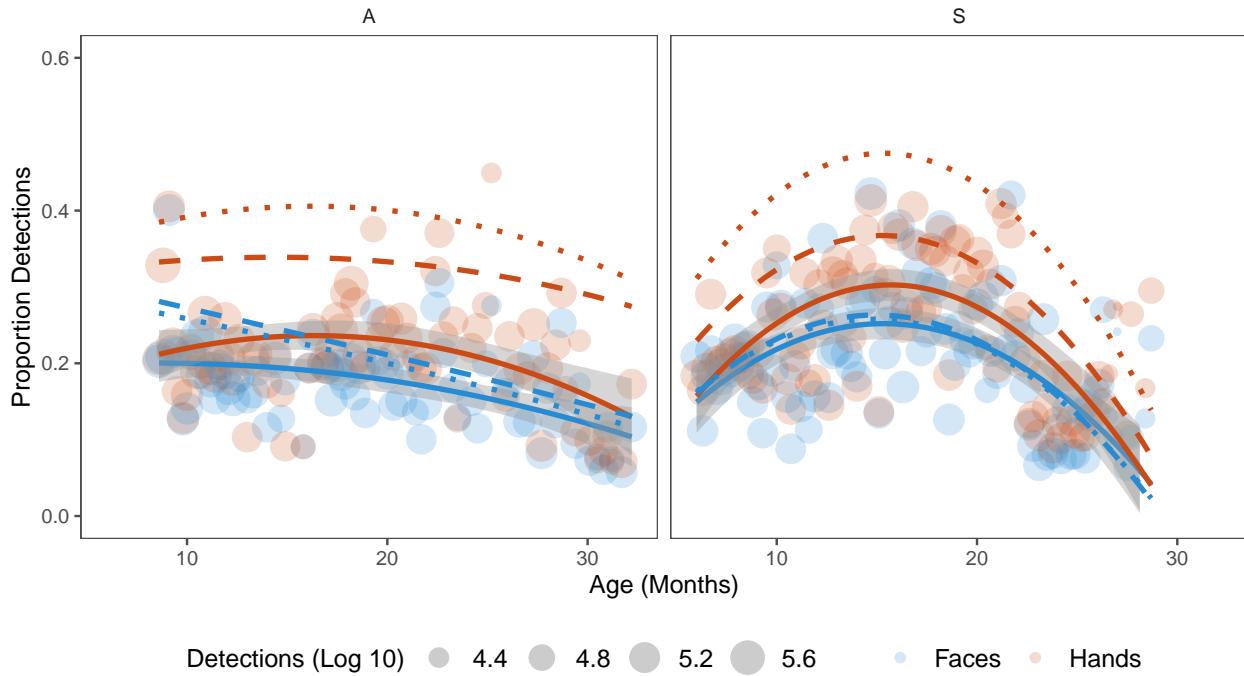
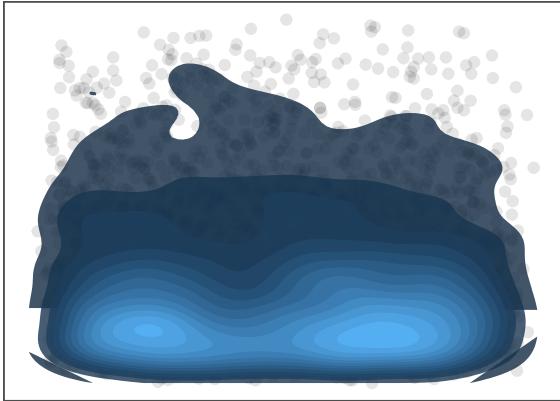


Figure A1. Proportion of faces and hands seen as a function of age for each child in the dataset. Data are binned by each week that the videos were filmed and scaled by the number of frames in that age range. Dashed lines show estimated trend lines from proportion of faces/hands in view when analyzing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

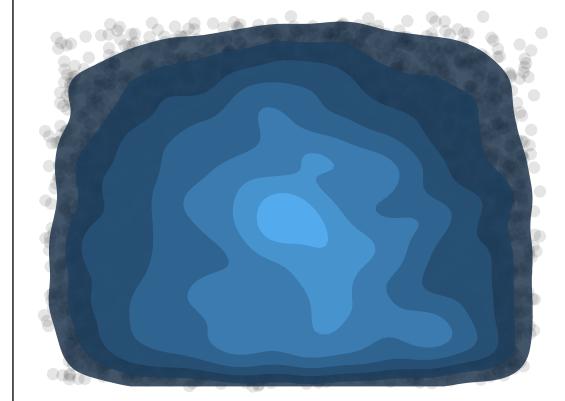


Figure B1. Density estimates for the child (left) and adult (right) hands that were detected in the 24K frame random gold set; each dot represents the center of a bounding box made by an adult participant. Brighter values indicate more detections.

Appendix C

Distribution of faces and hands in the visual field

⁴⁶⁹ We explored where in the visual field children tended to see faces and hands, suspecting
⁴⁷⁰ that these distributions might become wider as children grow older and learn to locomote
⁴⁷¹ on their own, following preliminary analyses from Frank (2012). As expected, faces tended
⁴⁷² to appear in the upper visual field in contrast to hands, which tended to be more centrally
⁴⁷³ located (see Figure C1). However, we found little evidence for any changes in the positions
⁴⁷⁴ of faces and hands across age, suggesting that this is a relatively stable property of infants'
⁴⁷⁵ visual environment from 6 months of age.

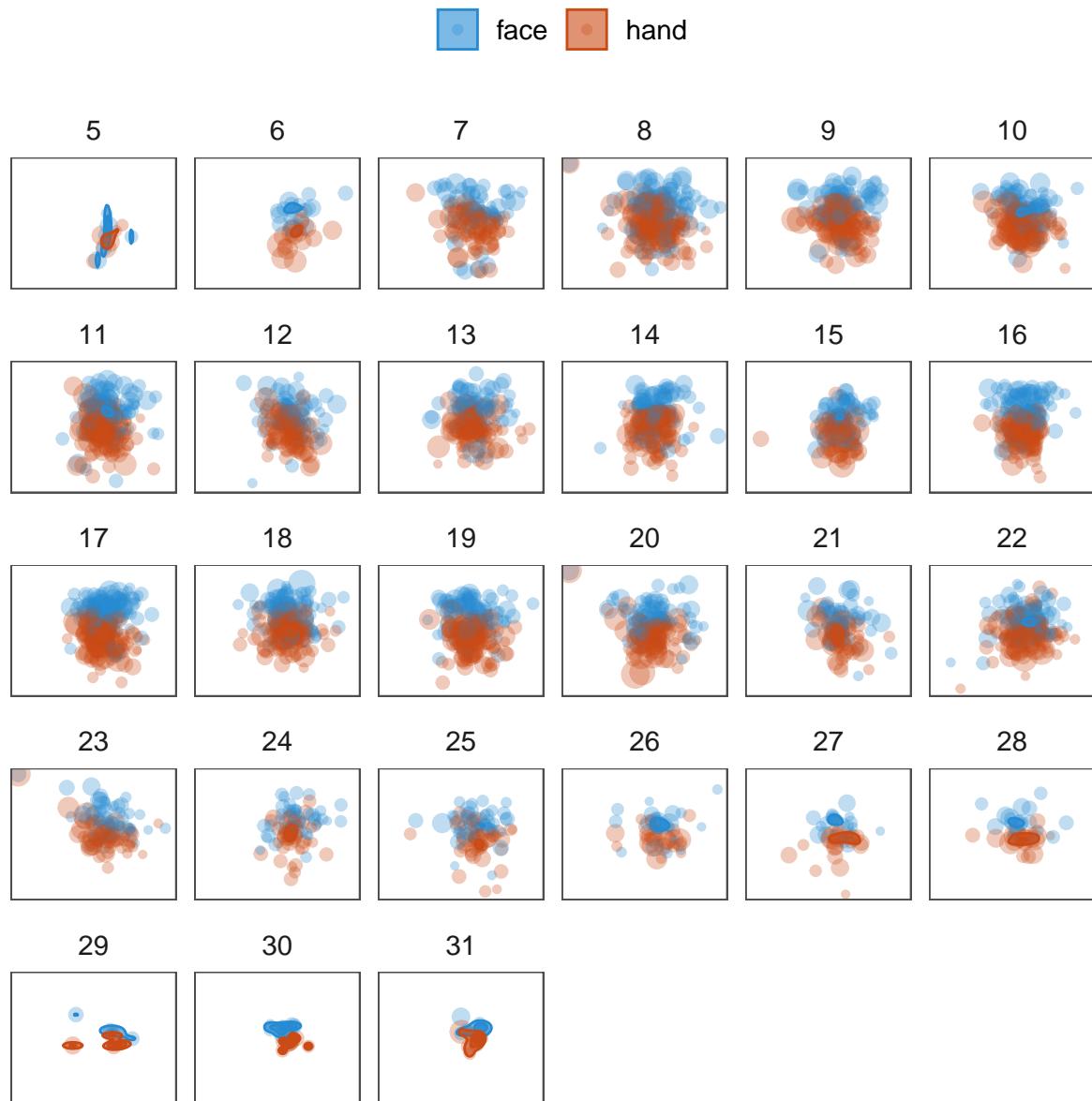


Figure C1. Each panel shows the average position of faces and hands in the visual field; each dot represents the average position from one video within a given age range.