

¹ Cross-cultural differences in children's object handling at home

² Marisa Casillas^{1,2} & Mary Elliott^{3,2}

³ ¹ University of Chicago

⁴ ² Max Planck Institute for Psycholinguistics

⁵ ³ University of Texas at Dallas

6

Abstract

7 Object-centered interactions (e.g., child-caregiver toy play) are thought to make significant
8 contributions to children's early word learning. However, it is yet unknown how frequently
9 such interactions occur in children's daily lives. We investigate how often 83 children under
10 age four (and their interactants) handle objects during everyday life in two non-Western
11 communities: one Mayan and one Papuan. Indeed, after infancy, children handle objects
12 relatively frequently (15% or more of their time) and do so in bursts and lulls across the day.
13 Importantly, cultural differences impact how frequently children handle objects across age,
14 perhaps due to differences between these communities in infant carrying practices, young
15 children's daily activities, and the objects available to children. In contrast, we find little
16 evidence to support the idea that object handling by children's interactants might similarly
17 drive early word learning. We discuss the implications of these findings for theoretical and
18 computational models of word learning.

19 *Keywords:* culture, word learning, Mayan, Papuan, daylong recording, egocentric
20 images

21 Word count: X

22 Cross-cultural differences in children's object handling at home

23 **Statement of relevance**

24 Before children go to school, they learn many words at home. Object-centric
25 interactions may propel this early learning: communication is facilitated when the child and
26 an interactant jointly focus on an object (e.g., a toy). However, it is unknown how often
27 children actually engage in such interactions at home, nor whether these interactions are
28 similarly present across diverse cultural contexts. We investigate how often young children
29 and their interactants handle objects during daily life in two Indigenous communities: one
30 Mayan and one Papuan. We find that, while children themselves often handle objects, their
31 interactants rarely handle objects immediately relevant to them. Children also showed
32 different object handling patterns between communities, likely due to differences in infant
33 carrying style and child daily activities. Therefore, focusing on child-led (not parent-led)
34 object interactions and attending to how cultural factors shape everyday activities will be
35 key to illuminating how children learn words at home.

36 **Introduction**

37 As children gain the ability to pick things up, sit on their own, and move around
38 independently, they also experience immense changes in their social and object-centered
39 interactions (Adolph, Karasik, & Tamis-LeMonda, 2010; Franchak, Kretch, Soska, & Adolph,
40 2011; Gaskins, 2000; Kretch, Franchak, & Adolph, 2014; Sanchez, Long, Kraus, & Frank,
41 2018). Object-centric interaction with others, particularly coordinated child-caregiver
42 attention via object handling, has been proposed as a potentially powerful source of
43 information for early word learning (e.g., Yu & Smith, 2013). Egocentric recordings of
44 interaction show decreasing views of faces and increasing views of hands over the first two
45 years of life (Fausey, Jayaraman, & Smith, 2016; Jayaraman, Fausey, & Smith, 2017; but see
46 Long, Kachergis, Agrawal, & Frank, 2020). This focus on active hands can lead to episodes
47 of shared attention in which communication, and thereby word learning, is facilitated (Yu &

⁴⁸ Smith, 2013). And while both children and their caregivers can initiate these interactional
⁴⁹ episodes, children's own holding experiences throw objects into sharp perceptual relief,
⁵⁰ thereby increasing the potential for learning about that object's label, functional affordances,
⁵¹ and more (e.g., Amatuni et al., 2021; Elmlinger, Suanda, Smith, & Yu, 2019; Slone et al.,
⁵² 2018; Soska, Adolph, & Johnson, 2010; see also Clerkin, Hart, Rehg, Yu, & Smith, 2017;
⁵³ Lockman & Kahrs, 2017; Long, Kachergis, Bhatt, & Frank, 2021).

⁵⁴ In order to establish such episodes as plausible and universal drivers of early word
⁵⁵ learning, we must determine how often they occur during children's daily lives and how their
⁵⁶ distribution changes across cultural contexts and child age. If object handling interactions
⁵⁷ play a critical role in word learning, we would expect them to (a) occur frequently, to enable
⁵⁸ the accumulation of evidence across multiple timepoints; and (b) in massed (i.e., bursty)
⁵⁹ distribution, at least for younger children, who may struggle to learn labels when the
⁶⁰ information is spread over longer periods (Clerkin, Hart, Rehg, Yu, & Smith, 2017; Long,
⁶¹ Kachergis, Agrawal, & Frank, 2020; Vlach & Johnson, 2013; Yurovsky & Frank, 2015;
⁶² Yurovsky, Smith, & Yu, 2013; Zhang, Yurovsky, & Yu, 2021). Finally, (c) we would need to
⁶³ see that these distributional properties are robustly maintained across diverse cultural
⁶⁴ contexts.

⁶⁵ We do not yet know how often such opportunities for hand-led intersubjective attention
⁶⁶ arise during children's day-to-day lives around the world. In a free-play setting in the lab,
⁶⁷ US infants and their caregivers handle objects over 90% of the time (Yu & Smith, 2013).
⁶⁸ However, at-home recordings suggest much lower handling rates, with the appearance of
⁶⁹ hands in children's view topping out around 30% (Fausey, Jayaraman, & Smith, 2016; Long,
⁷⁰ Kachergis, Agrawal, & Frank, 2020). While at-home recordings effectively capture a variety
⁷¹ of activity contexts, they have tended to be short and parent-selected in past work, so are
⁷² not representative of children's whole waking days. Further, while prior work has focused on
⁷³ North American children, cross-cultural differences in how children are carried, where they

74 are placed, what kinds of objects typically surround them, and what activities they engage in
75 guarantee wide variability in object handling (Adolph, Karasik, & Tamis-LeMonda, 2010;
76 Gaskins, 2000). For example, children who spend much of their first year carried or tightly
77 bound (e.g., for safety or warmth, Hayashi, 1992; Ishak, Tamis-LeMonda, & Adolph, 2007;
78 LeVine & Lloyd, 1966; Mei, 1994) cannot easily reach out to pick up nearby objects.

79 The present study examines natural patterns of child and interactant object handling
80 at home in two unrelated, non-Western populations. We analyzed the frequency and
81 distribution of object handling in more than 113000 child-perspective photos taken during 83
82 daylong, at-home recordings of children's waking days at home in two communities: one in
83 which children are typically carried on their mother's back for most waking hours during the
84 first year of life (Tseltal; Mayan) and one in which children are typically carried in caregivers'
85 arms during the first year (Rossel; Papuan; Figure 1). Our results suggest that children—but
86 not their interactants—indeed handle objects fairly frequently from toddlerhood onward,
87 with handling episodes distributed in bursts and lulls across the day. That said, object
88 handling patterns across age varied by population, pointing to potentially important effects
89 of cultural context in explaining how object-centric interaction might drive human language
90 development. In what follows we describe each population and explain the methods for data
91 collection, annotation, and analysis before detailing the results. We discuss the implications
92 of the present findings for both theoretical and computational models of early word learning.

93 Method

94 Participating communities

95 We first describe the typical carrying practices, physical and social environments, and
96 patterns of child-directed speech observed in the two communities of study. Neither author is
97 a member of the communities described. Our description is informed by experience talking
98 to and observing families, plus consultation with two researchers who have worked there for
99 several decades (P. Brown and S. C. Levinson). See Supplementary Materials for example



Figure 1. Typical infant holding positions in the studied Tseltal and Rossel communities.

100 images of scenes typical to each context.

101 The Tseltal (Mayan) participants live in a rural swidden horticulturalist community
102 situated within the Chiapas highlands of Southern Mexico. For most of the first year of life,
103 children are carried in a sling on their mother's back during her waking hours and are rarely
104 put on the ground before they walk (Brown, 2011, 2014). When sleeping they are wrapped
105 up completely in the sling, and when waking they are hitched into a sitting position with the
106 head and torso free (Figure 1). We have also observed infants sometimes being carried in the
107 arms, or held on the waist or in a lap. When children *are* placed on the ground, it is usually
108 on a woven mat or blanket, or occasionally in a cardboard box, and always under the
109 surveillance of a nearby caregiver. Infants are occasionally put to sleep in a hammock during
110 the daytime. For these reasons, it is rare to see Tseltal infants crawling. Mothers report that,
111 at some point, infants begin making limb movements while tied up in the sling, indicating
112 their readiness to begin standing and, soon after, taking their first steps. At this point,
113 caregivers may provide standing and walking support with their bodies and household
114 objects until the child is walking on their own. This Tseltal community is situated on a

mountainside such that, when children and their caregivers leave the home (e.g., to visit a relative), they must typically traverse steep dirt paths for part of the way, occasionally walking on a paved road. Thus, while nearly all children between 12 and 18 months walk, they may still be carried for significant parts of the day, especially over longer distances or challenging terrain.

The Rossel (Papuan) participants in our study live in a collection of swidden horticulturalist communities on the northeastern region of Rossel Island, which is the farthest outlying atoll of the Louisiade Archipelago, off the coast of mainland Papua New Guinea. Like Tseltal children, Rossel children are carried for much of their first year. However, they are typically carried in the arms: That is, on the waist, against the chest, and over the shoulder (Figure 1, Brown, 2011; Brown & Casillas, accepted). Children, even young infants, are cared for by a wide network of nearby family members and neighbors (male and female, adult and child), any of whom might carry or hold the infant for sustained periods. Rossel infants, like Tseltal infants, are rarely put directly on the ground. Instead they are more often found laying, sitting, scooting, and crawling on the raised verandas that are constructed under or attached to residences. Similarly, infants are sometimes placed in hammocks for daytime sleep. Once children show signs of interest in walking, caregivers might set out a line of small posts in the ground that the infant can grasp and walk along under a caregiver's supervision. In brief, while infants are rarely placed on the ground in both communities, Tseltal children's movements and ability to grasp nearby objects is more restricted early in development.

Cross-cultural differences in carrying practices less strongly influence children's object handling once they begin to walk, but the social landscape and available objects create persistent differences between communities. While children in both contexts spend significant time interacting with other children (e.g., older siblings and cousins/neighbors), this pattern is especially prominent in the Rossel community, where children join large,

141 independent child playgroups shortly after they start to walk (Brown, 2011, 2014; Brown &
142 Casillas, accepted). These large playgroups sometimes engage in stationary, object-centric,
143 verbally interactive activities (e.g., pretend household play, cracking and eating foraged
144 nuts), but more often facilitate mobile activities in which few objects are relevant and verbal
145 activity is repetitive or routinized (e.g., diving games in the river, chasing games similar to
146 'tag'). Tseltal children tend to participate in smaller child social groups and tend to move
147 within a somewhat more restricted area around their household grounds. Spot observations
148 of Yucatec Mayan children suggest that play, particularly manipulative play with objects and
149 substances, increases in frequency with age, to 40% of the child's time by ages 3–5 (Gaskins,
150 2000, work activities, which may involve object manipulation, also increase in this period).
151 While both communities are remote, the Tseltal community is far more commercially
152 connected to the Western, industrial world—Rossel Island sees only infrequent and irregular
153 boat contact. For this reason, objects that are designed specifically for children's interest or
154 manual manipulation (e.g., toys, crayons, etc.) are even less frequently seen in the Rossel
155 context than the Tseltal one.

156 Recording methods

157 The present data come from a subset of daylong photo-linked audio recordings collected
158 in 2015 (Tseltal) and 2016 (Rossel) of 55+ children under age 5;0 in each site. Children were
159 outfitted with an elastic vest that carried two devices: A lightweight Olympus audio recorder
160 (WS-832 or WS-853) and a wearable camera (Narrative Clip 1) with a miniature fisheye lens
161 (Photojojo Super Fisheye; Figure 2). Some infants (typically those 0;6 and younger) could
162 not wear both devices at once, and so instead wore an infant bodysuit ("onesie") with the
163 audio recorder while their caregiver wore an adult-sized vest with the camera.

164 The camera captured images at a fixed interval over the course of the recording day at
165 home, which typically lasted around 9 (Tseltal) or 8 (Rossel) hours. Timestamped images
166 were captured every 30 seconds in the Tseltal data and, following a camera firmware update

167 in late 2015, every 15 seconds in the Rossel data. Participants were able to cover the camera
168 lens at their discretion using a piece of cloth attached to the underside of the camera case
169 (see Casillas, Brown, & Levinson, 2020 for details).



Figure 2. Narrative clip camera with attached fisheye lens.

170 Information about the child's developmental, linguistic, domestic, and demographic
171 profile was collected via interviews between the child's caregivers, the first author, and a
172 research assistant who lives in the community. Dates of birth were collected both verbally
173 and from any available medical documentation. When dates were inconsistent across
174 information sources, we triangulated them using any additional information we could gather
175 (e.g., birth date relative to another child).

176 We here analyze photo data from 40 Tseltal and 43 Rossel daylong recordings (from 38
177 and 42 children, respectively). We focused on children between ages 0;0 and 4;0, which was
178 the target age range during data collection, and hence the most densely sampled. Across
179 sites, the children are balanced as well as possible in terms of age and sex though imbalances
180 remain. Our samples incorporate a majority of the local population in the target age range
181 at the time of data collection.

182 **Annotation**

183 We annotated photos with IMCO (github.com/marisacasillas/ImCo/), an open-source
184 program that allows researchers to efficiently annotate photos using keyboard inputs that are
185 mapped to predefined categories. Each photo dataset typically took a trained research
186 assistant fewer than 10 seconds to annotate. For each photo we annotate: The number of
187 visible adults (i.e., post-pubescent) and children, whether any visible child was crying or
188 breastfeeding, whether the target child was handling an object, whether one or more of the
189 target child's interactants was handling an object directly relevant to the target child (e.g.,
190 while feeding them), and whether the photo was unusable or skipped (e.g., due to
191 overexposure). We only annotated photos between the time the researcher left and the time
192 she returned because caregivers tended to restrict the target child's movements during this
193 period.

194 The present study examines object handling, for which there are four categories: 'C' =
195 the target child is handling an object (e.g. while playing with a toy); 'I' = one or more of the
196 target child's interactants are handling an object relevant to the child (e.g. while washing or
197 feeding the child); 'B' = the target child *and* one or more of the target child's interactants
198 are handling an object (e.g. while playing with a toy together); 'N' = there is no object
199 handling visible.

200 A handled object was defined as something (e.g. a toy, piece of food, or rock) that the
201 target child was holding or manipulating. Large or immovable objects were considered
202 handled if the child was actively engaging with them (e.g., a branch while climbing a tree,
203 but not a table on which a hand was resting). People were not counted as handled objects
204 (e.g. a mother holding her baby or a child holding a breast). Objects near, but not directly
205 in contact with, a child's hands were coded as handled when justifiable (e.g., a hand reaching
206 toward a ball rolling toward it). The chest-worn cameras were also considered held when
207 handled by the participants.#apostrophe

²⁰⁸ **Reliability and data preparation**

²⁰⁹ We analyze 151827 annotations of 113668 photos (41064 from Tseltal recordings and
²¹⁰ 72604 from Rossel recordings) by three annotators. A substantial proportion of photos from
²¹¹ each site were annotated by at least two annotators (Tseltal: 0.44; Rossel: 0.24). One
²¹² annotator contributed the majority of annotations for both sites (Tseltal: 29528; Rossel:
²¹³ 50331 unique photos), while the annotators primarily on one site or the other (Tseltal: 9557
²¹⁴ and Rossel: 39779 vs. Tseltal: 22632 and Rossel: 0).

²¹⁵ We analyze those photos deemed ‘usable’ by at least one annotator (87949 of the
²¹⁶ photos; 77.40% of the full annotated set). The remaining 25719 photos were blocked by the
²¹⁷ provided privacy cover, had too bright/dark lighting to annotate, were taken while the
²¹⁸ researcher was present, or were otherwise uncodeable. On average, this left us with 1,058.95
²¹⁹ photos per recording, but with wide variability between recordings (median: 1044; range:
²²⁰ 13–1737).

²²¹ Annotator agreement was high (85.90%), with comparable scores for Tseltal and Rossel
²²² recordings (Tseltal: 87.80%; Rossel: 83.10%). The primary source of disagreement came
²²³ from whether the target child was handling an object or not ('C' vs. 'N' disagreements; see
²²⁴ Supplementary Materials). Photos with no visible object handling were by far the most
²²⁵ common outcome, so given cases of disagreement, the resulting unweighted Cohen's kappa
²²⁶ score suggests moderate overall agreement ($\kappa = 0.59$).

²²⁷ We derived a single annotation value for each photo as shown in Table 1. We then
²²⁸ derived a single burstiness estimate for child (“C” or “B”) and interactant (“I” or “B”)
²²⁹ object handling for each recording, based on Goh and Barabási (2008)'s B parameter.
²³⁰ Burstiness (B) is calculated as $B = (\sigma_\tau - \mu_\tau)/(\sigma_\tau + \mu_\tau)$ where τ is the distribution of
²³¹ inter-event intervals (IEIs¹). A B score of -1 indicates a completely periodic distribution; 1 a

¹ An “event” here is a photo featuring object handling. Because of the coarse nature of the photo data (i.e.,

Table 1

A single annotation value was derived for each photo from the combination of annotations provided by coders as described below.

Combination	Result	# Photos
B only	B	566
C + C/B	C	14011
I + I/B/C	I	982
N + anything	N	72390

²³² maximally bursty distribution; and 0 a random distribution. For reference, multimodal event
²³³ distributions in adult matcher-director games have been found to typically fall around $B =$
²³⁴ 0.15–0.2. (Abney, Dale, Louwerse, & Kello, 2018).

²³⁵ Results

²³⁶ All analyses and figures were generated in R (Aust & Barth, 2020; R Core Team, 2020;
²³⁷ Wickham et al., 2019). The anonymized data and analysis code are available at
²³⁸ github.com/marisacasillas/daylong-photos. Full tabular regression outputs can be found in
²³⁹ the Supplementary Materials.

²⁴⁰ Frequency and distribution of object handling.

²⁴¹ Children.

²⁴² Across both daylong photo datasets, children 4;0 and under handled objects in an
²⁴³ average of 15.80% of photos (median = 12.90%, range = 0%–56.73%). Prevalence of object

samples every ~15 or ~30 seconds) we cannot precisely say when handling events started and stopped. IEI is therefore the time in seconds between two consecutive photos featuring object handling of the desired type (here “C”/“B” or “I”/“B” for the child and interactant analysis, respectively).

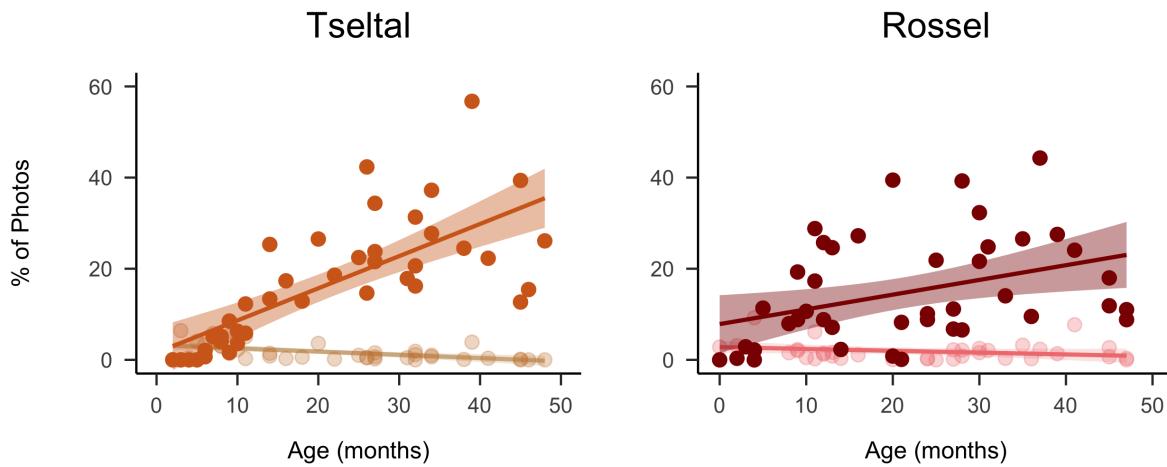


Figure 3. Proportion of photos in which children (dark) and their interactants (light) are seen handling an object during the Tseltal (left) and Rossel (right) daylong recordings.

244 handling was overall similar between sites, with an average of 16.91% of photos showing
 245 object handling for the Tseltal dataset (median = 15.82%, range = 0%–56.73%) and 14.77%
 246 of photos for the Rossel dataset (median = 11.02%, range = 0%–44.29%). However,
 247 inspection of object handling by age across sites reveals a pattern whereby Tseltal infants
 248 handle objects less frequently than Rossel children early on, but then do so more frequently
 249 than Rossel children later (Figure 3).

250 A linear regression modeling the effects of age, cultural group, and their interaction on
 251 the prevalence of child object handling (i.e., the proportion of photos in which the child was
 252 handling an object; $N = 83$; $AIC = -136.73$)² revealed impacts of both age and cultural
 253 group. Specifically, children handled objects more frequently with age overall ($b = 0.003$,
 254 95% CI [0.002, 0.004], $t(79) = 2.802$, $p = 0.006$), and children in the Tseltal data based
 255 showed a larger increase in object handling with age compared to Rossel children ($b = 0.004$,
 256 95% CI [0.002, 0.005], $t(79) = 2.322$, $p = 0.023$).

² $\text{glm}(\text{Proportion of photos showing object handling} \sim \text{Child age (months; numeric)} * \text{Site (Tseltal/Rossel; factorial)})$

257 Children's object handling was bursty in all 77 recordings with two or more instances
 258 of child object handling (Figure 4); $B_{mean} = 0.45$, $B_{sd} = 0.17$, $B_{range} = -0.05\text{--}0.81$). A linear
 259 regression³ ($N = 77$; AIC = -48.63) revealed no effects of age or cultural context on
 260 burstiness.

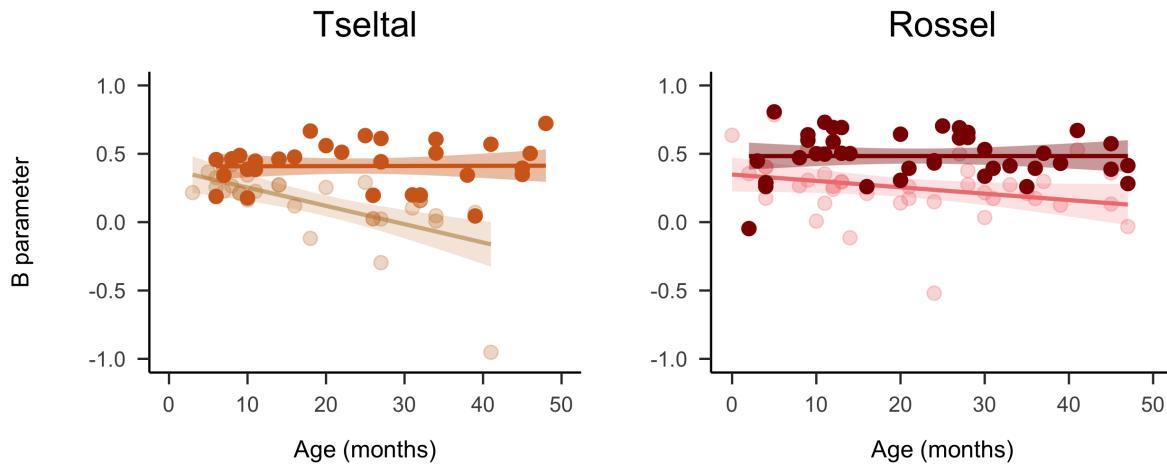


Figure 4. Burstiness values for child (dark) and their interactants' (light) object handling during the Tseltal (left) and Rossel (right) daylong recordings.

261 **Interactants.**

262 Overall, children's interactants handled child-relevant objects in an average of 1.82% of
 263 photos (median = 1.09%, range = 0%–11.01%). Prevalence of interactant object handling
 264 was also similar between sites, with an average of 1.70% of photos for the Tseltal dataset
 265 (median = 0.92%, range = 0%–6.41%) and 1.94% for the Rossel dataset (median = 1.36%,
 266 range = 0%–11.01%). Interactant object handling, unlike child object handling, appears
 267 descriptively similar across age between sites (Figure 3). Using the same model structure, we
 268 found no evidence to support impacts of age, cultural group, or their interaction on
 269 interactant object handling ($N = 83$; AIC = -405.52; all $|t| < 0.54$).

270 Interactants' object handling was bursty in 63 of the 69 recordings with two or more
 271 instances of interactant object handling (0.91; Figure 4; $B_{mean} = 0.20$, $B_{sd} = 0.24$, $B_{range} =$

³ $\text{glm}(B \sim \text{Child age (months; numeric)} * \text{Site (Tseltal/Rossel; factorial)})$

272 -0.95–0.79). A linear regression ($N = 69$; AIC = -15.22) revealed a significant age-by-site
273 interaction whereby burstiness decreased more across age for Tseltal interactants than Rossel
274 ones ($b = -0.009$, 95% CI [-0.013, -0.005], $t(65) = -2.113$, $p = 0.038$). There was no evidence
275 for further effects of age or cultural context.

276

Discussion

277 The present data suggest that object handling episodes have the greatest potential for
278 word learning after age one, though this may vary across cultural contexts; the data do *not*
279 support interactant object handling as critical for early word learning. Specifically, children's
280 object handling increased with age, reaching 15% of the time or higher for children over age
281 1;0, but showed differing developmental trajectories between sites—Tseltal children initially
282 handled objects less often than Rossel children did, but then later handled objects more
283 often. Interactants' child-relevant object handling was very rare—around 1–2% of the
284 time—and was not impacted by age or cultural context. With little exception, handling
285 events occurred in bursty distribution across the day by both children and interactants, and
286 across age and cultural contexts.

287 **Cross-cultural differences**

288 The difference in developmental trajectories between these two sites is likely driven by
289 multiple factors, including carrying practices and the social and economic organization of
290 daily life. Carrying practices may drive the differences that appear before children begin to
291 walk; Tseltal infants under 1;0, who are carried in a sling, seldom handled objects, while
292 Rossel infants, who are typically carried in the arms, steadily increased object handling over
293 the first year. Social organization impacts the number and composition of interactants
294 present, as well as the types of activities children engage in; in our dataset, children on
295 Rossel Island typically spent long stretches of the day with large, independent play groups
296 that frequently engaged in non-object-centric activities (e.g., playing tag, diving in the river)
297 whereas Tseltal children tended to spend more of their day doing domestic themed,

298 object-centric activities (e.g., (playing at) preparing foods, drawing and/or toy play) in
299 smaller groups, and stayed nearer to the immediate area around the family home. More
300 market integration in the Tseltal community with its surrounding industrialized economies
301 also translates to a much higher preponderance of toys and other small manmade objects
302 compared to Rossel Island (e.g., cups, pens, plastic bottles, etc.), though we note these
303 objects are still less prevalent than they are in, e.g., the EuroAmerican homes where much of
304 the prior work on object-centric interaction has taken place.

305 **How much is enough?**

306 Establishing naturalistic home rates of object handling in interaction will helpfully
307 constrain the input assumptions made by theoretical and computational models of word
308 learning in the first few years of life. In-lab investigations of US object-centric interactions
309 show objects being held by children or their interactants over 90% of the time (Yu & Smith,
310 2013), which is not sustainable during everyday life at home, including housework activities,
311 adults socializing, daytime sleep periods, and so on, all of which are captured in home
312 daylong recordings. The present results place the likely frequency of object handling closer
313 to indirect estimates based on the general prevalence of hands during short, at-home
314 recordings—a maximum of around 25–30% of the time for children over 1;0 (Fausey,
315 Jayaraman, & Smith, 2016; Long, Kachergis, Agrawal, & Frank, 2020).

316 Importantly, we do not know how much object handling is “enough” to support early
317 word learning. If this type of learning episode is very potent or frequently accompanied by
318 actual object labels, even low rates of object handling may result in robust word learning.
319 Given the relatively low rates of child-directed talk during at-home recordings in these and
320 Western communities (Bunce et al., 2020; Casillas, Brown, & Levinson, 2020, 2021) we
321 anticipate that object handling episodes that are actually accompanied by talk about the
322 object are very infrequent. Object handling activity is, at least, bursty; repeated or sustained
323 handling may provide a boost to label learning, compensating for the low overall frequency

324 of these events.

325 **Lack of interactant object handling.** Interactant object handling was rare in the
326 present data, suggesting that in these communities it plays, at most, a minor role in early
327 word learning. Our estimates stand in stark contrast to those from in-lab object-centric
328 child-caregiver interaction in the US, in which adult caregivers and young one-year-olds
329 handle objects with approximately equal prevalence (each 25% of the time on their own and
330 43% of the time jointly, Yu & Smith, 2013). Informally, we note that interactant-held objects
331 were typically limited to a few prototypical items associated with basic care, e.g., food,
332 clothing, and daily hygiene—not toys, books, or other items that would more often be found
333 in recordings with middle-class Western families. If present data are closer to children's true
334 at-home experiences, in Western homes or elsewhere, theories of word learning via hand-led
335 intersubjective episodes should focus on cases when the child is doing the object handling,
336 and not their interactant(s).

337 **Looking ahead**

338 The present work provides preliminary benchmarks for computational modeling and
339 future comparative work on word learning during multimodal interaction. TThe findings also
340 form a basis for further theory development around how object-handling events and early
341 word learning are influenced by the child's cultural and socioeconomic milieu in concert with
342 their motor development. Beyond word learning, naturalistic object handling patterns have
343 implications for children's development of visual, tactile, and event-based representations of
344 the world (e.g., the canonical form and function of familiar objects and the expected form
345 and function of novel objects). Continued work with cross-cultural collections of highly
346 naturalistic egocentric data will be key to understanding how children come to learn about
347 the world, and how their learning process is continually reshaped by features of both their
348 home environment and their own developmental gains.

349

Contributions

350 MC developed the study concept, collected the data, and contributed funding. ME
351 annotated the data and trained other annotators. Both contributed to the writing and
352 approved the final version of this manuscript.

353

Acknowledgements

354 We are indebted to Rebeca Guzmán López, Humbertina Gómez Pérez, Taakêmê
355 Ñamono, Y:aaw:aa Pikuwa, and the participating families and community leaders in each
356 population sampled. We also thank the PNG National Research Institute, the
357 Administration of Milne Bay Province, and the Centro de Investigaciones y Estudios
358 Superiores en Antropología Social (CIESAS) Sureste. Our work builds on that of Penelope
359 Brown and Stephen C. Levinson who were invaluable consultants in conducting this study.
360 We are grateful to Shawn C. Tice for his technical contributions to the photo-annotation
361 tool, and to Ine Alvarez van Tussenbroek and Cielke Hendriks for their help with annotation.
362 This work is supported by a NWO Veni Innovational Scheme grant (275-89-033) to MC and
363 by fieldwork funding from the Max Planck Institute for Psycholinguistics.

364

References

365

Abney, D. H., Dale, R., Louwerse, M. M., & Kello, C. T. (2018). The bursts and lulls of multimodal interaction: Temporal distributions of behavior reveal differences between verbal and non-verbal communication. *Cognitive Science*, 42(4), 1297–1316.

366

367

Adolph, K. E., Karasik, L. B., & Tamis-LeMonda, C. S. (2010). Motor skill. In M. H. Bornstein (Ed.), *Handbook of cultural developmental science* (pp. 61–88). Psychology Press: New York, NY.

368

369

Amatuni, A., Schroer, S., Zhang, Y., Peters, R., Reza, Md. A., Crandall, D., & Yu, C. (2021). Characterizing the object categories two children see and interact with in a dense dataset of naturalistic visual experience. In *Proceedings of the 43rd annual conference of the cognitive science society* (pp. 265–271).

370

371

Aust, F., & Barth, M. (2020). *papaja: Create APA manuscripts with R Markdown*. Retrieved from <https://github.com/crsh/papaja>

372

373

Brown, P. (2011). The cultural organization of attention. In A. Duranti, E. Ochs, & and Bambi B Schieffelin (Eds.), *Handbook of Language Socialization* (pp. 29–55). Malden, MA: Wiley-Blackwell.

374

375

Brown, P. (2014). The interactional context of language learning in Tzeltal. In I. Arnon, M. Casillas, C. Kurumada, & B. Estigarribia (Eds.), *Language in interaction: Studies in honor of Eve V. Clark* (pp. 51–82). Amsterdam, NL: John Benjamins.

376

377

Brown, P., & Casillas, M. (accepted). Childrearing through social interaction on Rossel Island, PNG. In A. J. Fentiman & M. Goody (Eds.), *Esther Goody revisited: Exploring the legacy of an original inter-disciplinarian* (pp. XX–XX).

- 388 New York, NY: Berghahn.
- 389 Bunce, J., Soderstrom, M., Bergelson, E., Rosemberg, C., Stein, A., Alam, F., ...
- 390 Casillas, M. (2020). *A cross-cultural examination of young children's everyday*
- 391 *language experiences.*
- 392 Casillas, M., Brown, P., & Levinson, S. C. (2020). Early language experience in a
- 393 Tseltal Mayan village. *Child Development, 91*(5), 1819–1835.
- 394 Casillas, M., Brown, P., & Levinson, S. C. (2021). Early language experience in a
- 395 papuan community. *Journal of Child Language, 48*(4), 792–814.
- 396 Clerkin, E. M., Hart, E., Rehg, J. M., Yu, C., & Smith, L. B. (2017). Real-world
- 397 visual statistics and infants' first-learned object names. *Philosophical Transactions*
- 398 *of the Royal Society B: Biological Sciences, 372*(1711), 20160055.
- 399 Elmlinger, S. L., Suanda, S. H., Smith, L. B., & Yu, C. (2019). Toddlers' hands
- 400 organize parent-toddler attention across different social contexts. In *2019 joint*
- 401 *IEEE 9th international conference on development and learning and epigenetic*
- 402 *robotics (ICDL-EpiRob)* (pp. 296–301). IEEE.
- 403 Fausey, C. M., Jayaraman, S., & Smith, L. B. (2016). From faces to hands: Changing
- 404 visual input in the first two years. *Cognition, 152*, 101–107.
- 405 Franchak, J. M., Kretch, K. S., Soska, K. C., & Adolph, K. E. (2011). Head-mounted
- 406 eye tracking: A new method to describe infant looking. *Child Development, 82*(6),
- 407 1738–1750.
- 408 Gaskins, S. (2000). Children's daily activities in a Mayan village: A culturally
- 409 grounded description. *Cross-Cultural Research, 34*(4), 375–389.

- 410 Goh, K.-I., & Barabási, A.-L. (2008). Burstiness and memory in complex systems.
411 *EPL (Europhysics Letters)*, 81(4), 48002.
- 412 Hayashi, K. (1992). The influence of clothes and bedclothes on infants' gross motor
413 development. *Developmental Medicine & Child Neurology*, 34(6), 557–558.
- 414 Ishak, S., Tamis-LeMonda, C. S., & Adolph, K. E. (2007). Ensuring safety and
415 providing challenge: Mothers' and fathers' expectations and choices about infant
416 locomotion. *Parenting: Science and Practice*, 7(1), 57–68.
- 417 Jayaraman, S., Fausey, C. M., & Smith, L. B. (2017). Why are faces denser in the
418 visual experiences of younger than older infants? *Developmental Psychology*,
419 53(1), 38.
- 420 Kretch, K. S., Franchak, J. M., & Adolph, K. E. (2014). Crawling and walking
421 infants see the world differently. *Child Development*, 85(4), 1503–1518.
- 422 LeVine, R. A., & Lloyd, B. B. (1966). *Nyansongo: A Gusii community in Kenya*.
423 Wiley.
- 424 Lockman, J. J., & Kahrs, B. A. (2017). New insights into the development of human
425 tool use. *Current Directions in Psychological Science*, 26(4), 330–334.
- 426 Long, B., Kachergis, G., Agrawal, K., & Frank, M. C. (2020). *Detecting social
427 information in a dense database of infants' natural visual experience*.
- 428 Long, B., Kachergis, G., Bhatt, N., & Frank, M. C. (2021). Characterizing the object
429 categories two children see and interact with in a dense dataset of naturalistic
430 visual experience. In *Proceedings of the 43rd annual conference of the cognitive
431 science society* (pp. 279–285).

- 432 Mei, J. (1994). The northern chinese custom of rearing babies in sandbags:
433 Implications for motor and intellectual development. In J. H. A. van Rossum & J.
434 I. Laszlo (Eds.), *Motor development: Aspects of normal and delayed development*.
435 VU Uitgeverij: Amsterdam, NL.
- 436 R Core Team. (2020). *R: A language and environment for statistical computing*.
437 Vienna, Austria: R Foundation for Statistical Computing. Retrieved from
438 <https://www.R-project.org/>
- 439 Sanchez, A., Long, B., Kraus, A. M., & Frank, M. C. (2018). Postural developments
440 modulate children's visual access to social information. In *Proceedings of the 40th*
441 *annual conference of the cognitive science society* (pp. 2412–2417).
- 442 Slone, L. K., Abney, D. H., Borjon, J. I., Chen, C., Franchak, J. M., Pearcy, D., ...
443 Yu, C. (2018). Gaze in action: Head-mounted eye tracking of children's dynamic
444 visual attention during naturalistic behavior. *Journal of Visualized Experiments*,
445 (141).
- 446 Soska, K. C., Adolph, K. E., & Johnson, S. P. (2010). Systems in development:
447 Motor skill acquisition facilitates three-dimensional object completion.
448 *Developmental Psychology*, 46(1), 129–138. <https://doi.org/10.1037/a0014618>
- 449 Vlach, H. A., & Johnson, S. P. (2013). Memory constraints on infants'
450 cross-situational statistical learning. *Cognition*, 127(3), 375–382.
- 451 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., ...
452 Yutani, H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*,
453 4(43), 1686. <https://doi.org/10.21105/joss.01686>
- 454 Yu, C., & Smith, L. B. (2013). Joint attention without gaze following: Human infants
455 and their parents coordinate visual attention to objects through eye-hand

- 456 coordination. *PloS One*, 8(11), e79659.
- 457 Yurovsky, D., & Frank, M. C. (2015). An integrative account of constraints on
458 cross-situational learning. *Cognition*, 145, 53–62.
- 459 Yurovsky, D., Smith, L. B., & Yu, C. (2013). Statistical word learning at scale: The
460 baby's view is better. *Developmental Science*, 16(6), 959–966.
- 461 Zhang, Y., Yurovsky, D., & Yu, C. (2021). Cross-situational learning from ambiguous
462 egocentric input is a continuous process: Evidence using the human simulation
463 paradigm. *Cognitive Science*, 45(7), e13010.

¹ Supplementary Materials: Cross-cultural differences in children's object handling at home

² Marisa Casillas^{1,2} & Mary Elliott^{3,2}

³ ¹ University of Chicago

⁴ ² Max Planck Institute for Psycholinguistics

⁵ ³ University of Texas at Dallas

- 6 Supplementary Materials: Cross-cultural differences in children's object handling at home

7 **Developmental context**



Figure 1. Four images illustrating aspects of the Tseltal and Rossel environments. Left: Tseltal target infant in the mother's lap at a family morning meal in the kitchen hut. Left-center: Tseltal mother walking with a man up a main path (target child is in sling). Right-center: Rossel children using stones to crack open foraged nuts under the stairs entering a household structure. Right: Rossel child playing in the net hammock under their house with a caregiver laying nearby on the veranda.

(#fig:context.imgs)

8

Reliability

9 Confusion matrices

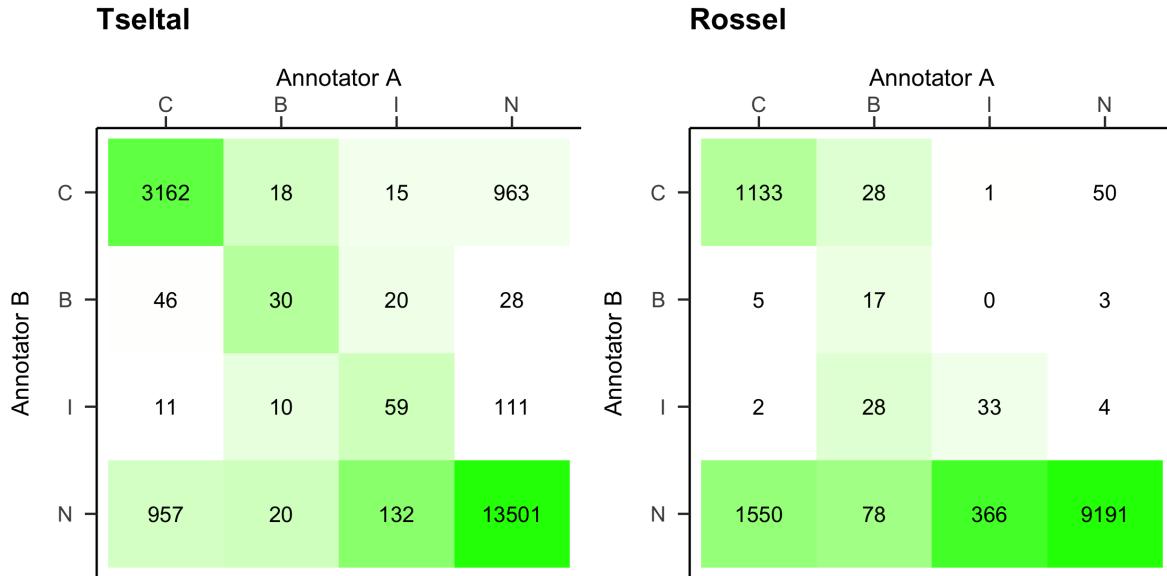


Figure 2. Confusion matrices of possible photo annotation values by site: Tseltal (left) and Rossel (right). The value in each cell indicates the number of photos with the given combination of annotations from person A and person B. The shading indicates proportion of the label category for Annotator A distributed over Annotator B's possible responses.

(#fig:confmat.fig)

Table 1

Linear regression output for the analysis of annotation agreement by child age, population, and their interaction.

	B	SE	t	p
(Intercept)	0.90	0.04	22.55	0.00
Age	0.00	0.00	-1.84	0.08
SiteMayan	0.02	0.05	0.44	0.66
Age:SiteMayan	0.00	0.00	0.20	0.84

Effects of age and cultural context

To test whether annotation was more reliable at different ages or in different sites, we conducted a linear regression of mean agreement value by child age (in months; numeric), site (Tseltal vs. Papuan; factorial), and their interaction ($\text{Mean_agmt} \sim \text{Age} * \text{Site}$; $N = 29$; $\text{AIC} = -70.13$). We found no evidence for significant impacts of age or cultural context on annotation agreement scores.

Full regression output of object handling analyses

Below are the full model outputs for each regression reported in the main text as well as equivalent models with a log-transformed dependent variable (which often yielded more normally distributed residual estimates, but showed no qualitative differences from the results reported in the main text).

Table 2

Linear regression output for the analysis of proportion photos showing child object handling by child age, population, and their interaction.

	B	SE	t	p
(Intercept)	0.08	0.03	2.69	0.01
Age	0.00	0.00	2.80	0.01
SiteMayan	-0.06	0.04	-1.49	0.14
Age:SiteMayan	0.00	0.00	2.32	0.02

Table 3

Linear regression output for the analysis of log-transformed proportion photos showing child object handling by child age, population, and their interaction.

	B	SE	t	p
(Intercept)	-3.01	0.24	-12.52	0.00
Age	0.04	0.01	3.81	0.00
SiteMayan	-0.53	0.35	-1.53	0.13
Age:SiteMayan	0.03	0.01	2.02	0.05

Table 4

Linear regression output for the analysis of child object handling burstiness by child age, population, and their interaction.

	B	SE	t	p
(Intercept)	0.48	0.05	9.53	0.00
Age	0.00	0.00	0.00	1.00
SiteMayan	-0.08	0.08	-0.97	0.34
Age:SiteMayan	0.00	0.00	0.03	0.97

Table 5

Linear regression output for the analysis of log-transformed child object handling burstiness by child age, population, and their interaction.

	B	SE	t	p
(Intercept)	-0.92	0.17	-5.28	0.00
Age	0.01	0.01	1.27	0.21
SiteMayan	0.13	0.27	0.51	0.61
Age:SiteMayan	-0.01	0.01	-1.09	0.28

Table 6

Linear regression output for the analysis of proportion photos showing child-relevant interactant object handling by child age, population, and their interaction.

	B	SE	t	p
(Intercept)	0.03	0.01	4.83	0.00
Age	0.00	0.00	-1.76	0.08
SiteMayan	0.00	0.01	0.54	0.59
Age:SiteMayan	0.00	0.00	-0.96	0.34

Table 7

Linear regression output for the analysis of log-transformed proportion photos showing child-relevant interactant object handling by child age, population, and their interaction.

	B	SE	t	p
(Intercept)	-3.52	0.16	-21.40	0.00
Age	-0.01	0.01	-1.67	0.10
SiteMayan	0.26	0.24	1.12	0.27
Age:SiteMayan	-0.02	0.01	-1.63	0.11

Table 8

Linear regression output for the analysis of child-relevant interactant object handling burstiness by child age, population, and their interaction.

	B	SE	t	p
(Intercept)	0.35	0.06	5.77	0.00
Age	0.00	0.00	-1.91	0.06
SiteMayan	0.04	0.10	0.40	0.69
Age:SiteMayan	-0.01	0.00	-2.11	0.04

Table 9

Linear regression output for the analysis of log-transformed child-relevant interactant object handling burstiness by child age, population, and their interaction.

	B	SE	t	p
(Intercept)	0.25	0.16	1.58	0.12
Age	0.00	0.01	-0.59	0.56
SiteMayan	0.33	0.25	1.31	0.20
Age:SiteMayan	-0.03	0.01	-2.65	0.01