

**Assessing the alignment between infants' visual and linguistic experience
using multimodal language models**

Alvin W.M. Tan^{*1}, Jane Yang^{*2}, Tarun Sepuri², Khai Loong Aw¹, Robert Z. Sparks¹, Zi Yin³,
Virginia A. Marchman¹, Michael C. Frank¹, and Bria Long²

¹Department of Psychology, Stanford University

²Department of Psychology, University of California, San Diego

³Department of Psychology, Tsinghua University

Author Note

Correspondence at tanawm@stanford.edu

Abstract

Figuring out which objects or concepts words refer to is a central language learning challenge for young children. Most models of this process posit that children learn early object labels from co-occurrences of words and their referents that occur when someone around them talks about an object in the immediate physical environment. But how aligned in time are children's visual and linguistic experiences during everyday learning? To date, answers to this question have been limited by the need for labor-intensive manual annotations of vision–language co-occurrences. Here, we evaluate the use of contrastive language-image pretraining (CLIP) models to automatically characterize vision–language alignment in egocentric videos taken from the infant perspective in home environments. After validating CLIP alignment scores using human alignment judgments, we apply this metric to a large corpus of infant-perspective videos. We show that idealized aligned moments for learning (e.g., “look at the *ball*” with a ball present in the child's view) are relatively rare in children's everyday experiences compared to modern machine learning datasets, and highlight variability in alignment both within and across children. These findings suggest that infrequent alignment is a constraint for models describing early word learning and offer a new method for investigating children's multimodal environment.

Keywords: early word learning; head-mounted cameras; multimodal language models; naturalistic observations

Assessing the alignment between infants' visual and linguistic experience using multimodal language models

Learning that the word “dog” refers to golden retrievers and Chihuahuas, but not cows, requires children to form visual concepts from their everyday experiences. Children’s learning environments provide the basis for these inferences about which visual concepts words refer to, but how often do children actually see clearly labeled exemplars of categories, and how important are these moments for forming visual concepts? Despite the centrality of this question for theories of early word learning, we have limited data about how temporally aligned children’s visual and linguistic experiences are during development. Here, we leverage contrastive language–image pre-training (CLIP) models (Radford et al., 2021) to assess how temporally aligned children’s visual and linguistic experiences are over early development.

Associative models of early word-learning—such as cross-situational word learning—predict that children extract meaning from co-occurrences of frequent words and objects in their environment (Yu & Smith, 2007; Y. Zhang et al., 2021). Egocentric videos taken from the child’s perspective offer an unprecedented window to children’s learning environments during their everyday experiences (Yoshida & Smith, 2008), allowing researchers to quantify how often children actually experience word–object co-occurrences. Initial investigations suggest that the more that children experience a referent (as measured by prevalence in egocentric videos), the more likely they are to learn the word for it in a subsequent looking-while-listening task (Bergelson & Aslin, 2017). Recent computational simulations have further highlighted the importance of clearly labeled examples in building useful category representations: Vong et al. (2024) found that CLIP models can learn rudimentary representations from roughly 60 hours of a single child’s naturalistic head-mounted data, but that directly labeled referents were disproportionately more informative for learning.

Understanding the degree to which children’s visual and linguistic experiences are aligned in everyday learning is thus critical for both building realistic models of early word learning and understanding variability in learning rates across children. Early language learning is highly

variable between children (Frank et al., 2021), and this variability may have downstream consequences for both language and reading success (Marchman & Fernald, 2008). To date, however, our knowledge has been limited by the need for manual annotations (Bergelson & Aslin, 2017; Yoshida & Smith, 2008; Frank et al., 2013), which are unfeasible beyond a few hours of data. Here, we overcome this barrier to progress by using multimodal language models to examine infants' visual and linguistic experience in a large, naturalistic dataset of egocentric visual experience (Long et al., 2024).

Specifically, we leverage advances in multimodal large language models—contrastive language–image pre-training models (Radford et al., 2021)—to quantify the degree to which what children see is temporally aligned with the semantic content in the language that they hear. These multimodal models contain jointly trained visual and language encoders that maximize similarity between images and their captions in large datasets. For example, the language encoder could represent an utterance (i.e., “Do you want to read about the big hungry bear?”) and the vision encoder could represent a concurrent image (i.e., a living room scene with a book in the foreground, see Figure 1); computing cosine similarity of the resulting embeddings from both encoders provides a metric of semantic alignment between visual and linguistic inputs.

However, to date no work has systematically examined whether this metric can capture human judgments about the degree to which utterances and images are aligned. We address this gap by first validating CLIP alignment scores against human performance in a 4-alternative forced choice (4AFC) matching task using a stratified sample of utterances and frames from BabyView, a new open dataset of egocentric videos (Long et al., 2024) with automatic transcriptions and speaker diarization. After establishing that CLIP alignment scores correspond to human judgments, we apply this same metric to concurrent frames that occur during each utterance in the entire dataset, calculating an alignment metric for every utterance.

With these alignment estimates, we quantify consistency and variability in how children's visual and linguistic experiences are aligned across contexts and individuals. Overall, we find that the proportion of children's everyday experiences where the visual and linguistic inputs are highly

aligned is relatively low compared to modern machine learning datasets. We then examine variation in vision–language alignment across developmental age, within individual children, across activity contexts (e.g., reading, eating), and depending on the linguistic content of the utterances. By integrating computational innovations to describe development at scale, we aim to contribute to a broader and more inclusive view of everyday learning contexts, with the ultimate goal of understanding how variability across learning environments influences learning outcomes in this population.

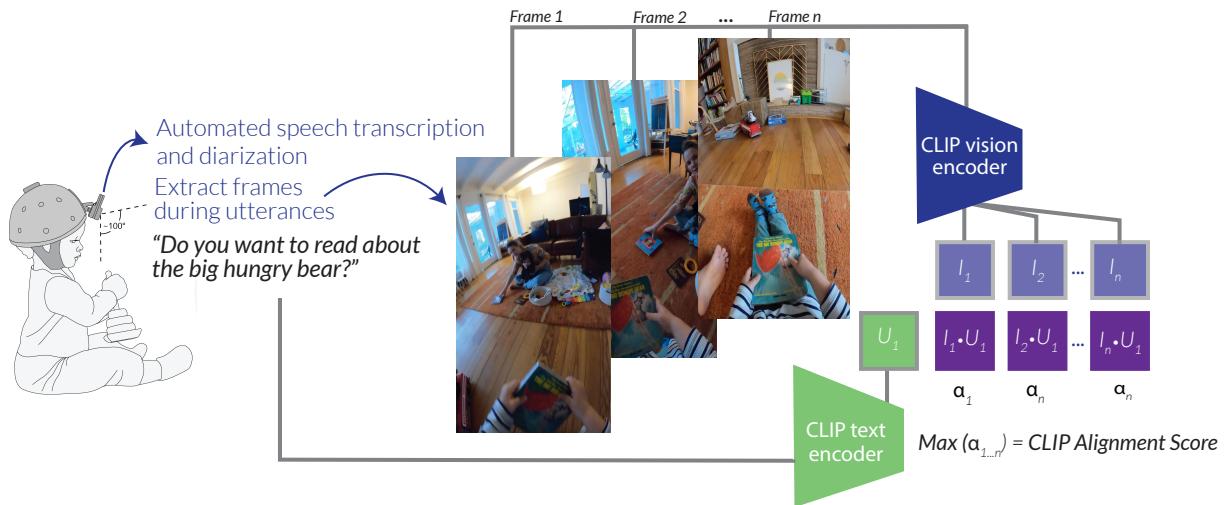


Figure 1

Overview of the pipeline for CLIP alignment scores, showing example frames from a BabyView camera in a naturalistic home environment; CLIP alignment scores are calculated by taking the dot product of the normed embeddings from the CLIP text and vision encoders for each frame associated with an utterance (as transcribed by Distil-Whisper). We then take the maximum alignment score (visualized in purple) for each utterance.

Dataset & Methods

We analyze naturalistic, egocentric video from children’s perspective, collected with the BabyView camera (Long et al., 2024). Compared to prior cameras (Sullivan et al., 2021), the BabyView camera offers substantially better audio quality (increasing the feasibility of automatic transcriptions), and a wider vertical field of view (increasing the likelihood of referents being

visible). We analyze BabyView dataset releases 2025.1 and 2025.2, including only monolingual English-speaking families, yielding 325 hours of speech from the families of 19 children (5–36 months of age).

Automated transcriptions

All videos were transcribed using Distil-Whisper (Radford et al., 2022; Gandhi et al., 2023), specifically the distil-large-v3 model,¹ with a word error rate of .35 (see Long et al. (2024) for further details on the transcripts). We conducted analyses both on all utterances in the videos as well as only those utterances that were identified as being produced by an adult by a speaker diarization model (Lavechin et al., 2020).

CLIP alignment calculation

To calculate vision–language alignment scores, we used an implementation of the OpenAI CLIP ViT/B32 model (Radford et al., 2021; Jina AI, 2023).² For each utterance, we calculated the cosine similarity between the image embedding for each frame (sampled at 1 fps) with the text embeddings for each utterance (“You want the ball”); we refer to this as the *CLIP alignment score*. For each utterance, the CLIP alignment score was the maximum score that occurred in any frame occurring during that utterance. See Figure 1 for an overview of this calculation.

Activity & location annotations

We additionally annotated the frames with activity and location annotations from Sepuri et al. (2025). These annotations were obtained by chunking each video into 10-second clips, and prompting a video question-answering model (VideoLLaMA 3, B. Zhang et al., 2025) to select the appropriate activity and location for each clip from a set of options. These annotations were used to sample a subset of frames for the human validation task described below, ensuring broad coverage of the dataset.

¹ <https://huggingface.co/distil-whisper>

² <https://github.com/jina-ai/clip-as-service>

Human annotations of alignment

We validated CLIP alignment scores through a four-alternative forced choice (4AFC) human annotation task.³ We developed two versions of a forced-choice task that was similar to the contrastive training objective of CLIP models, which evaluated the alignment of the visual and speech information. In one condition, annotators were presented with an utterance and asked to guess which of four frames matched the utterance (*image* condition). In a second condition, annotators were presented with one frame and asked to guess which of four utterances matched the frame (*utterance* condition).

To construct the set of trials shown to annotators, we conducted two types of stratified sampling across utterances. First, we stratified across CLIP alignment scores—we binned each utterance score into one of 5 bins and sampled up to 80 frames in each bin. Next, we stratified across the detected activities and locations (Sepuri et al., 2025), sampling at least 10 frames for each activity and for each location separately. Our goal in using these location and activity locations was to ensure as broad of coverage as possible over this very large dataset while only selecting a small set of utterances and their corresponding frames for analyses. Together, there were 732 frame–utterance pairs sampled for annotation. The three distractors for each trial were randomly sampled from the remaining frames and utterances within the set of sampled pairs.

We recruited 80 English-speaking annotators who completed a web-based version of this task on Prolific. Each annotator saw a mean of 110 test trials in one of the two conditions (~6 annotations per condition per trial), as well as 5 catch trials containing simple vocabulary questions. Annotators were excluded if they answered more than one catch trial incorrectly ($N = 2$). All annotators were blind to the purposes of the study and reported that they were English-speaking adults residing in the United States.

³ Preliminary efforts directly annotating image–utterance pairs for alignment (on a 1–5 Likert scale) suggested that this method was likely to be noisy due to the imprecision of what “alignment” means.

Model Validation

Validation via human annotations

We first assessed the feasibility of using the alignment scores by manually examining a random selection of frames with high alignment scores (CLIP score $\geq .24$, used in prior work (Vong et al., 2024)). We found qualitative evidence that utterances with high alignment tended to refer to a concrete object that was visible (e.g., “Don’t bite your toys”, “There’s your green chair”, “Cornflakes”). In particular, several high alignment frames came from books visible in the frame that caregivers were reading aloud; thus, there was text in the frame that exactly matched caregivers’ utterances. Some mismatches occurred when there was a contextual match (e.g., in a kitchen scene, the *fork* was absent for the utterance “Should I get you a fork?”), such that the alignment score was relatively high despite the referent being absent. However, manual inspection suggested that many utterances that had low alignment tended to refer to absent entities (“when we go back to Ohio, you’ll see lots of cows.”) or were directed towards other adults (“... down by the station later in the evening.”)

Next, we systematically evaluated the correspondence between human and model judgments by examining the human performance for frames on our 4AFC matching task; here, ground truth was the original frame–utterance pairing. We examined how well annotators could match each frame to an utterance (and vice versa) as a function of the alignment score produced by CLIP. As shown in Figure 2, we found that human performance scaled with the alignment score, such that frame–utterance pairs with higher CLIP scores were also more likely to be correctly guessed by humans. This relationship was confirmed by a logistic regression predicting human accuracy as a function of CLIP score, condition, and their interaction. CLIP score emerged as a significant predictor of human accuracy ($b = 18.5 [14.0, 23.1]$, $p < .001$). Neither condition nor the interaction between CLIP score and condition were significant predictors ($p > .8$ for both). This result suggests that the image and utterance conditions returned very similar values. At the utterance level, human accuracies showed a significant correlation between the two conditions ($r = 0.541$, $p < .001$). The predicted intercept (accuracy = 0.5) of the logistic

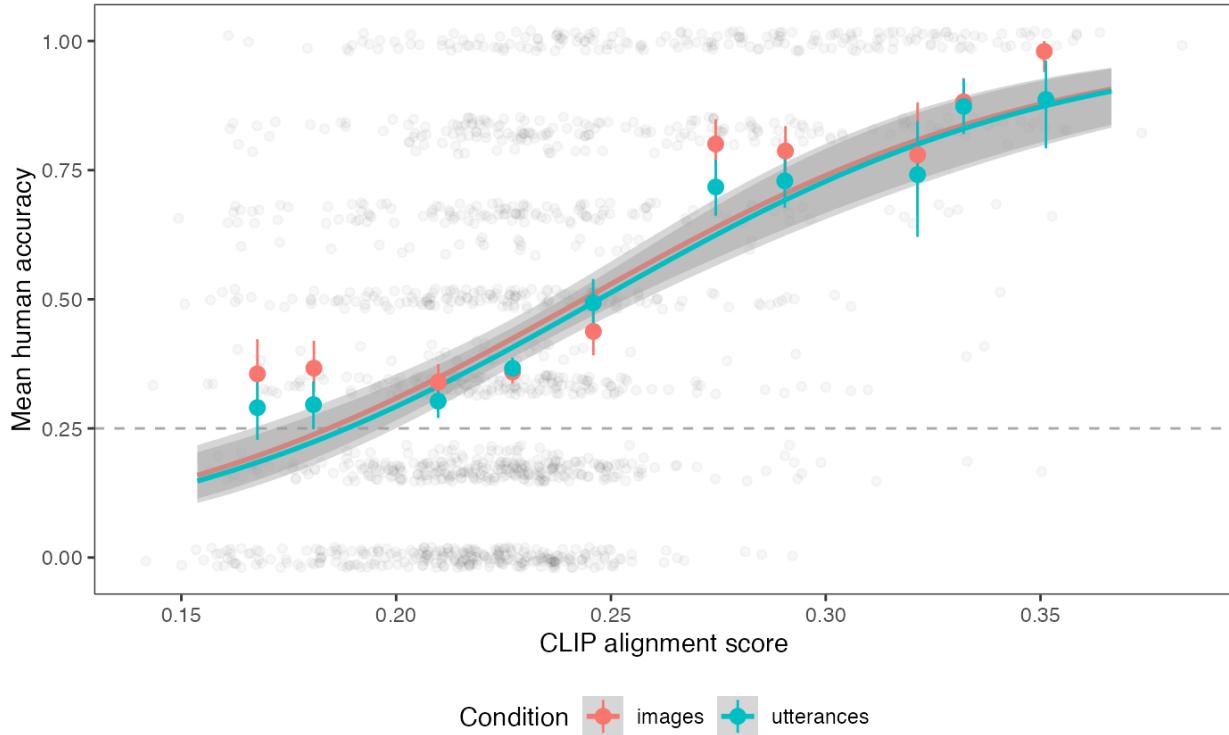


Figure 2

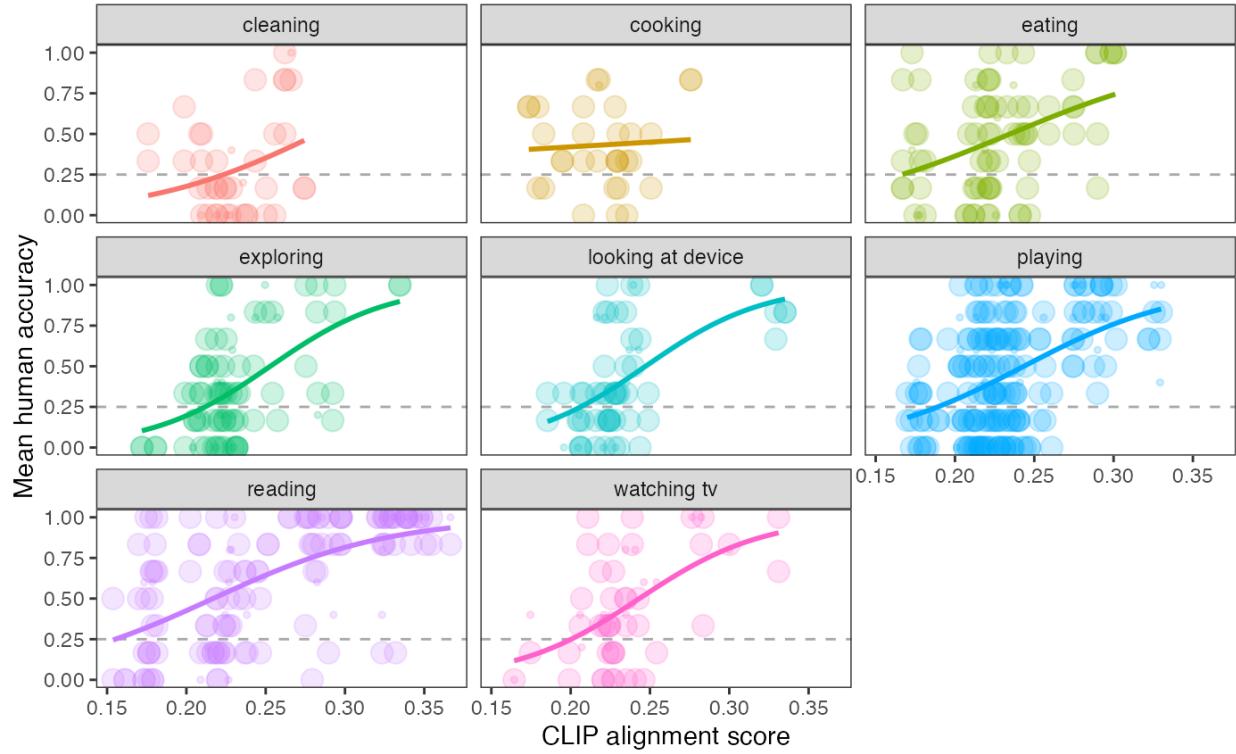
Human 4AFC accuracy by condition as a function of CLIP alignment score. Error bars indicate bootstrapped 95% confidence intervals by CLIP score decile, and lines indicate best-fit logistic curves along with their 95% confidence band. Dashed line indicates chance-level performance.

model occurred at a CLIP score of 0.25; this value was close to the threshold of 0.24 for high alignment used in previous work (Vong & Lake, 2025).

The overall relationship between CLIP alignment score and accuracy also held across different activities, as shown in Figure 3. While there were some observed differences in the intercept of the different logistic curves, these points occurred within the interval (0.24, 0.26).

Interpreting model alignment scores vs. model accuracy

We explored a set of analyses to help interpret CLIP alignment scores relative to human CLIP accuracy scores. We evaluated the same CLIP model (CLIP ViT/B32) on the same 4AFC task presented to our human annotators. We operationalized model choice as the softmax over the alignment scores between an image and the four possible utterances (or vice versa). Accuracy

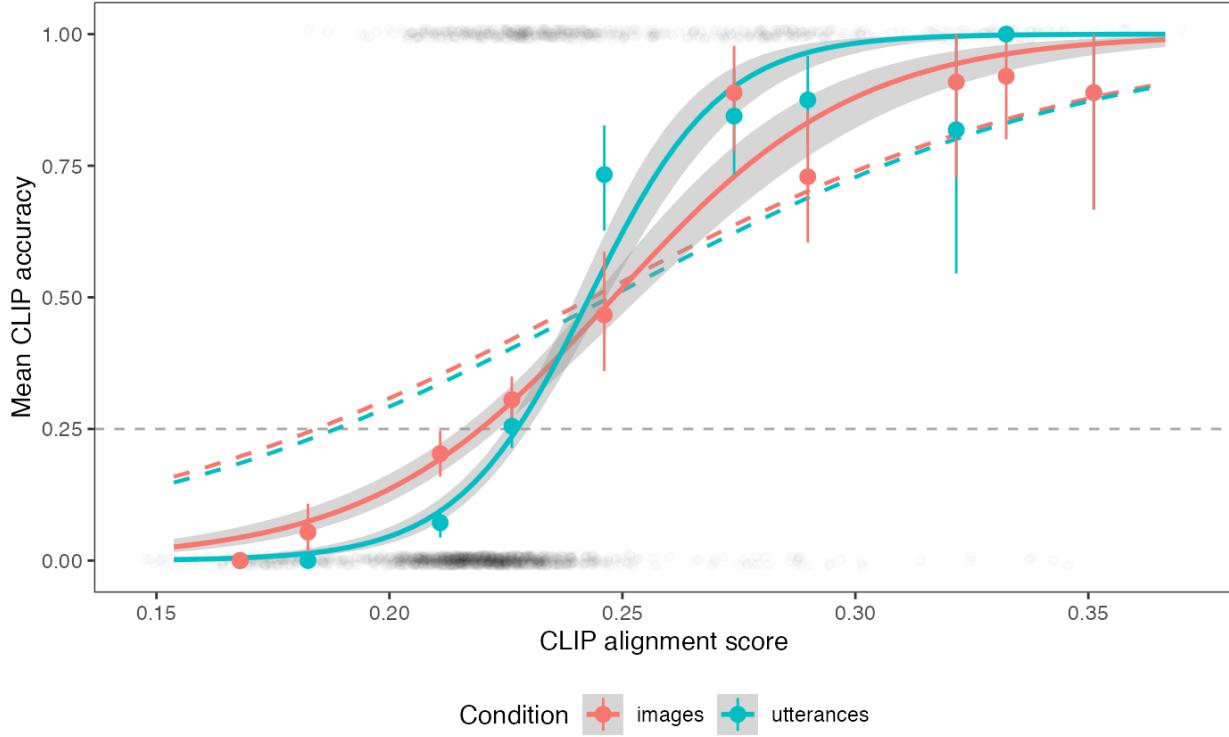
**Figure 3**

Human 4AFC accuracy as a function of CLIP alignment score, for the 8 most frequent activities.

Lines indicate best-fit logistic curves. Dashed line indicates chance-level performance.

was then determined by whether the model chose the correct target. We observed a positive relationship between alignment scores and 4AFC accuracy. This relationship also provides an intuitive interpretation of CLIP score: for a given image—utterance pair with a particular CLIP score and three distractor pairs, it approximates the probability that the CLIP model will select the correct pair.

Additionally, we found an interesting difference between CLIP and humans with respect to the shape of the relationship between CLIP alignment scores and accuracy. At low CLIP alignment scores, the CLIP model is almost at 0% accuracy, followed by a steep increase in accuracy at medium CLIP scores, and close to 100% accuracy at high CLIP scores (see Figure 2). In contrast, humans show approximately chance-level performance at low CLIP scores, and accuracy increases with a shallower slope. Further, humans show above-chance performance at a

**Figure 4**

CLIP 4AFC accuracy as a function of CLIP alignment score. Error bars indicate bootstrapped 95% confidence intervals by CLIP score decile, and lines indicate best-fit logistic curves along with their 95% confidence band. Dashed colored lines reflect human 4AFC accuracies (as in Figure 2). Dashed grey line indicates chance-level performance.

lower threshold than the CLIP model. This difference suggests that CLIP may be wrongly assigning low cosine similarity to frame–utterance pairs on which humans succeeded in the 4AFC task. Nonetheless, the intercepts for both humans and models lie within the interval (0.24, 0.26), further validating these values as a threshold for high alignment.

Descriptive Results & Analyses

Next, we used per-utterance CLIP alignment scores to quantify the consistency and variability in visual–linguistic alignment in a large corpus of egocentric videos taken from the infant perspective. Below, we analyze these alignment scores across the overall dataset, within individual children's data, and with respect to the content of the linguistic utterance.

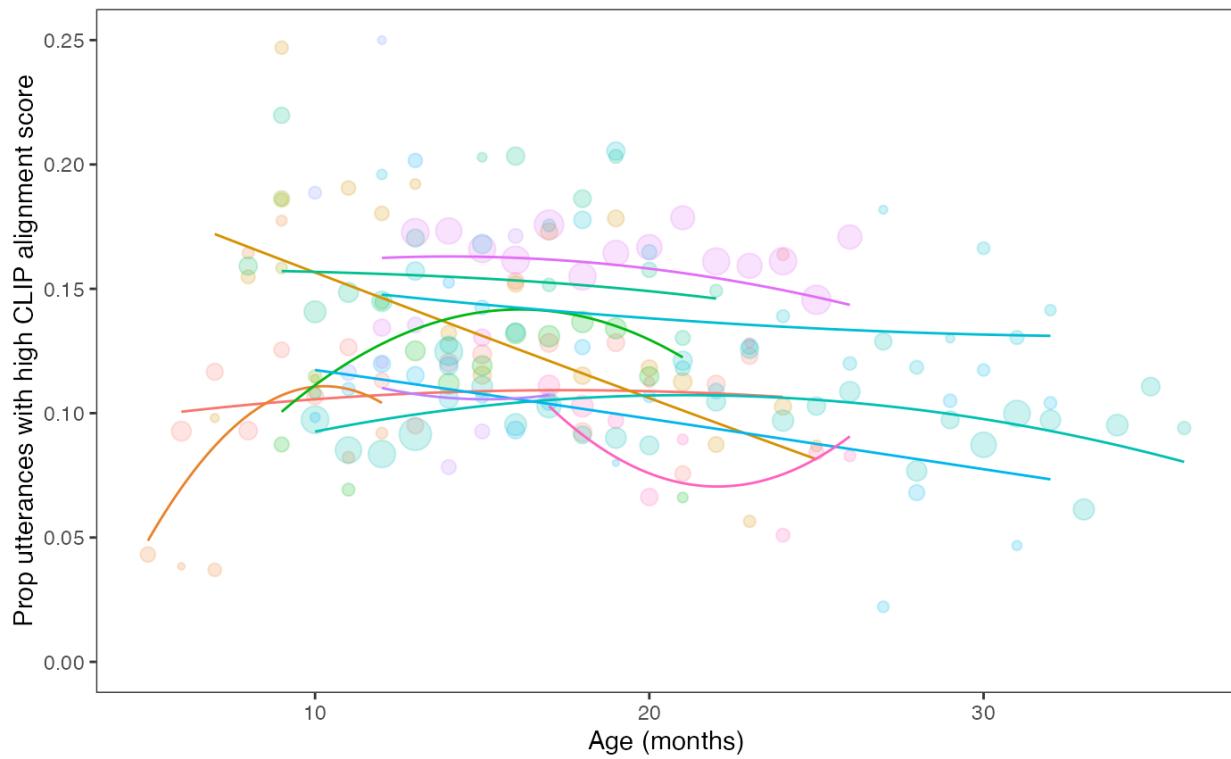
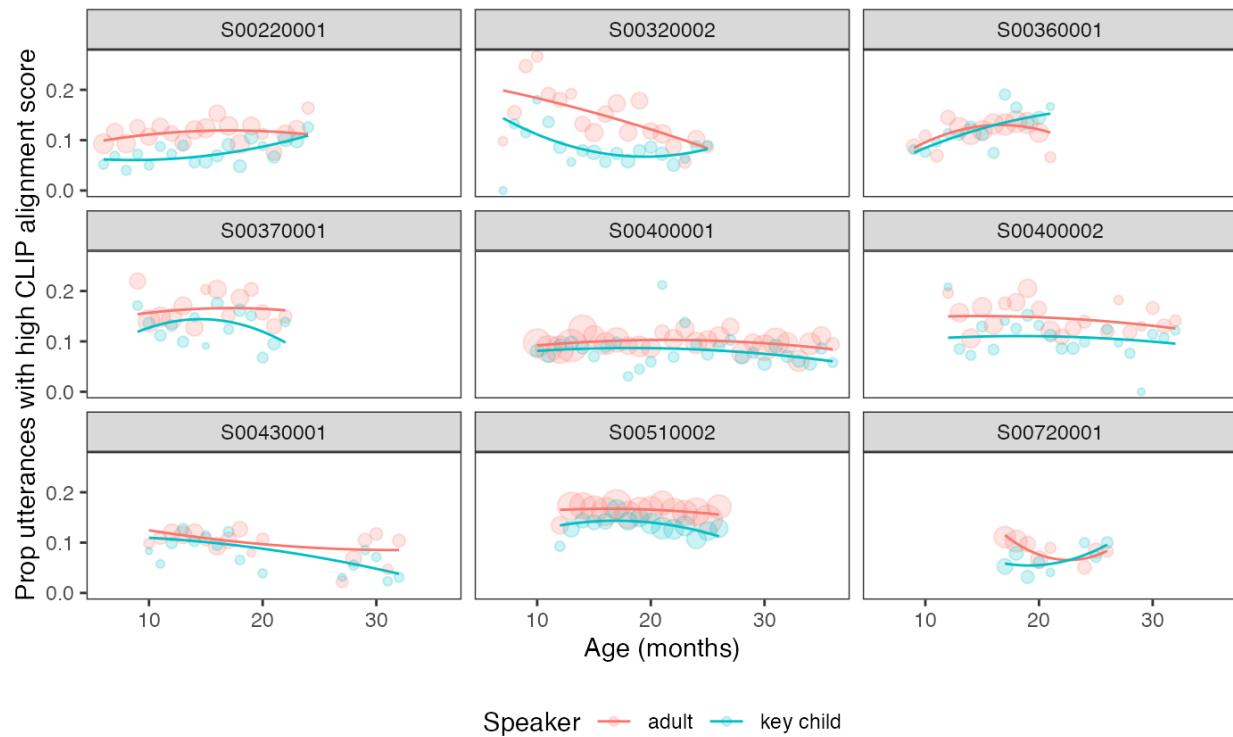


Figure 5

Proportion of highly-aligned adult utterances as a function of the child’s age. Colors reflect individual children, and lines reflect best-fit LOESS curves; the size of the dots scales with the amount of utterances in each age bin.

Vision–language alignment across individuals

Overall, we found that the children’s everyday experiences contained relatively infrequent moments with high visual–linguistic alignment, with considerable variability across individuals. Figure 5 shows the proportion of adult utterances in a given video that had high alignment, as a function of the age of the child at the time of recording; we found the same trend whether we used raw cosine similarity scores or these thresholded scores, as per Vong & Lake (2025). Notably, “highly aligned” frame–utterance pairs occurred less than 20% of the time for most children at almost all ages, and the overall average across our entire dataset was $M = 12.64\%$. However, note that there was also considerable heterogeneity across families, with some families trending around 16% high alignment while others were closer to 8%, a two-fold difference.

**Figure 6**

Proportion of highly aligned utterances as a function of speaker and the child's age, for the 9 families with the most data. Lines represent best-fit LOESS curves.

Higher alignment in adult-produced speech

What drives higher CLIP alignment scores? We examined several possible sources of variance in the alignment scores that we observed. First, we anticipated that utterances produced by adults would have higher alignment than utterances produced by either the child who was recording (i.e., “key child”) or other children in the household. We suspected this would be the case given both higher automated transcription accuracies for adult-produced speech (Long et al., 2024) and the general tendency of caregivers to perhaps try to label or describe a child’s visual environment for them.⁴ Thus, we next examined whether we saw systematic differences in

⁴ The current state-of-the-art diarization algorithm that we employed (Lavechin et al., 2020) is unable to classify whether adult speech is child-directed; however, it is able to distinguish between child-produced and adult-produced speech.

alignment for speech that contained adult- and key child-produced speech, as shown in Figure 6. As expected, we found that adult-produced speech had higher CLIP alignment scores, on average. This finding was supported by a linear mixed effects model predicting the proportion of highly aligned utterances as a function of age, speaker, and their interaction, along with random intercepts and age slopes by child. In this model, speaker type was the only significant predictor ($b = 0.039 [0.018, 0.060]$, $p < .001$). Thus, the higher alignment in adults vs. children was consistent regardless of the age of the child at the time of recording.

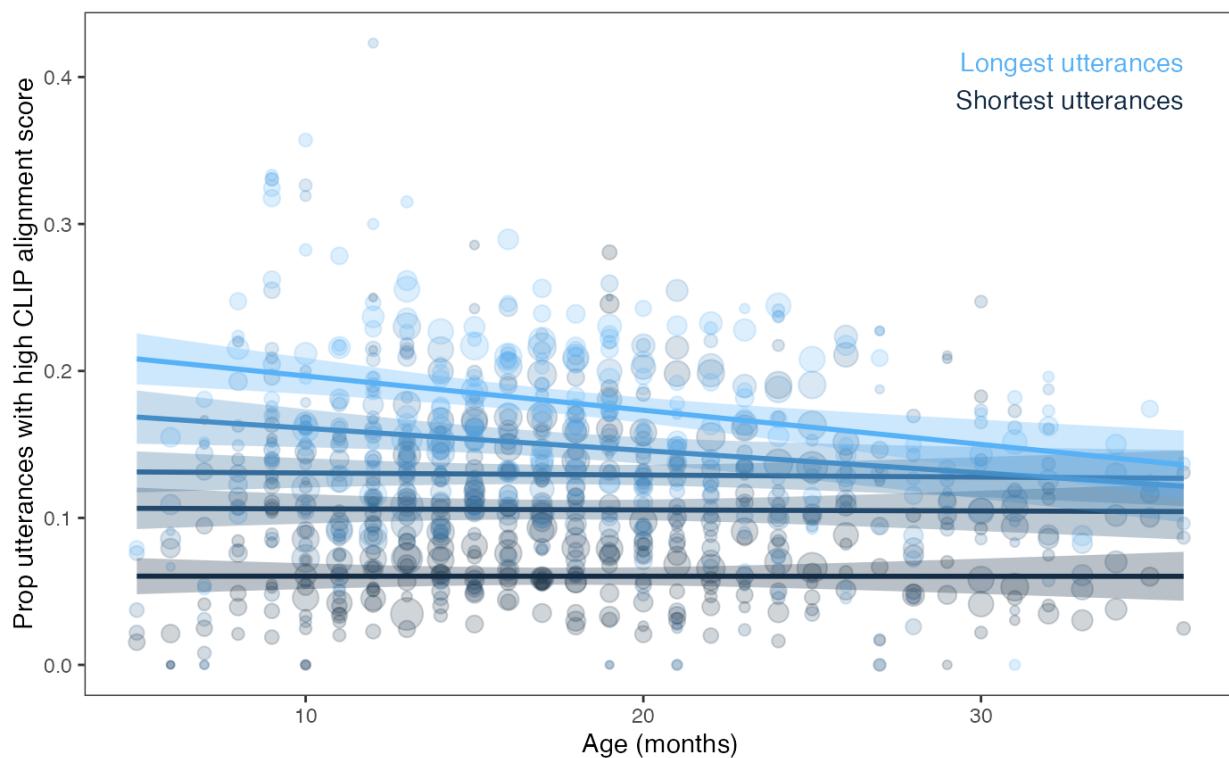


Figure 7

Proportion of utterances with high vision–language alignment scores as a function of children's age during recording and utterance duration; here, we restrict our analyses to adult-produced speech. Utterance duration is binned into quintiles; lines reflect best-fit linear trend lines for each utterance length decile; dot size reflects the number of utterances in each age/utterance duration bin. Lighter lines/dots refer to utterances of longer duration.

Higher alignment for longer utterances

We next examined the possible role of utterance duration in driving alignment scores, as shown in Figure 7. Analyzing adult utterances only, we found that utterances that spanned a longer duration had higher CLIP alignment scores than shorter utterances, on average. This relationship was confirmed using a linear mixed effects model predicting CLIP alignment score as a function of utterance length, age, and their interaction, along with random slopes for utterance length and intercepts by child (random slopes for age failed to converge). In this model, there was a significant effect of utterance duration ($b = 9.460 \times 10^{-4}$ [7.920×10^{-4} , 1.100×10^{-3}], $p < .001$). That is, longer utterances (e.g., “Here are your blue shorts with the dinosaurs”) tended to have higher vision–language alignment than shorter utterances (e.g., “Hi baby”, “Your bottle”). In addition, vision–language alignment decreased slightly across age for *longer* but not shorter utterances, which had overall lower alignment, as evidenced by an interaction in this same model between utterance duration and age ($b = -1.072 \times 10^{-5}$ [-6.567×10^{-6} , -1.487×10^{-5}], $p < .001$). Thus, these exploratory analyses suggest that adult speech may decrease slightly in its overall visual–linguistic alignment across age as children become older—perhaps reflecting changes in care routines or in the degree to which adults are referencing absent or future objects and events.

Higher alignment for utterances containing more frequent and concrete lemmas

We also investigated the possible role of the content of the utterances in driving alignment, focusing on properties of individual words in the utterances. We estimated the CLIP alignment score per word by lemmatizing all utterances using UDPipe (Straka et al., 2016). Then, for each lemma, we calculated the mean CLIP alignment score for all utterances containing that lemma. We then filtered down to lemmas that appeared in at least 10 different utterances to remove idiosyncratic low-frequency items, retaining 1213 lemmas (~20%). We then merged these lemmas with several psycholinguistic predictors: *log frequency* of occurrence of the lemmas in the BabyView dataset, *concreteness* (Brysbaert et al., 2014), *imageability* (Cortese & Fugett, 2004; Schock et al., 2011), and *sensorimotor* and *action strength* (Lynott et al., 2020). These

psycholinguistic predictors were chosen as they may relate to the extent to which the meaning of words can be represented with elements in the visual field, and that have been shown to affect word processing in adults (Khanna & Cortese, 2021).

Because the psycholinguistic norms did not cover all the lemmas present in our dataset, we used multiple imputation (with 5 imputations) using predictive mean matching to handle missing data. We then fit a linear model predicting lemma mean CLIP alignment score, with the fixed effects of log frequency, concreteness, imageability, sensorimotor strength, and action strength. Only two predictors emerged as significant. More concrete lemmas tended to occur in utterances with higher CLIP alignment scores than less concrete lemmas ($b = 0.0016 [0.0007, 0.0025]$, $p < .001$), and lemmas that were more frequent occurred in utterances with higher CLIP alignment scores than less frequent lemmas ($b = 0.0007 [0.0003, 0.0010]$, $p < .001$). We do note, however, that lemmas seem to follow a Zipfian distribution in frequency, thus there are a few highly frequent lemmas that may be driving most of the effect of log frequency.

Discussion

Can machine learning models be used to explore the alignment of linguistic and visual experience in children's everyday experiences? To answer this question, we first validated a computational technique for assessing alignment in videos of children's experience. We found that the similarity of the linguistic and visual embeddings in a contrastive language–image pre-training model (CLIP) was related to empirical human judgments of vision–language alignment, suggesting that this metric can be used as a proxy for vision–language alignment in other datasets.

Using this validated approach, we then asked: How aligned in time is what children see with what they hear? While idealized moments of joint attention between children and their caregivers may produce some highly aligned moments, we found that such aligned moments occurred with considerable variability across individuals and time—at most every one in five utterances in any given age bin, and far less often in many cases. Instead, children often heard speech between two caregivers talking with each other, from another child, or referring to an

object or situation that was not in the “here and now” (e.g., “Let’s go to the park”). Notably, this lower level of alignment contrasts with the amount of visual–linguistic alignment in curated datasets used to train vision–language models (e.g., COCO (Lin et al., 2015) or Flickr (Young et al., 2014)). Furthermore, the considerable heterogeneity across families in alignment rates highlights that young children’s word learning mechanisms need to be relatively robust to fairly large differences in the alignment of the visual and linguistic information.

A few different features significantly predicted CLIP alignment scores. First, adult-produced speech was significantly more aligned than child-produced speech. There are several possible reasons for this finding—for example, adult speech may have been more contextually contingent, or the adults may have been engaging in more pedagogical talk that was directly tied to the immediate context. It is also possible that adult utterances were simply more well-structured or better-transcribed compared to child utterances, and thus, more similar to the training data of CLIP. A related finding is that longer utterances exhibited higher alignment—and adults, given their better command of language, typically produce longer utterances than children. Future work could potentially disentangle these hypotheses by manipulating adult and child utterances in various ways to determine the impact of register, sentence structure, and context contingency on alignment.

Additionally, some types of lemmas occurred more often in highly aligned utterances than other lemmas. In particular, more frequent lemmas occur in more aligned utterances. This effect may be specific to child-directed speech, which may be more grounded in the here and now than general adult speech, such that more frequently occurring lemmas are also more grounded. Additionally, concrete lemmas occur in more aligned utterances, suggesting that these words are more likely to occur in contexts in which the referred objects are also in the visual scene. The causal direction of these effects remains an open question—for example, it could be that utterances containing concrete words result in children’s correct attending to the referred object, or it could be that having concrete objects in the environment drives adult to talk about those objects. In-lab interventional experiments (e.g., varying the number and concreteness of objects

in the environment) could help to determine the direction of causality.

We anticipated that alignment would, on average, decline with age as speech becomes more abstract (“Are you feeling happy?”) and adults are more likely to refer to absent referents (e.g., “Can you go get froggy?”). However, we did not find strong evidence for this trend in our dataset, as only our exploratory analyses revealed this effect for the longest utterances. However, a key avenue for future work is quantifying additional sources of variability in at-home egocentric recordings (e.g., adult- vs. child-directed speech, episodes of joint attention with caregivers) that may influence alignment.

We suspect that assessing vision–language alignment within more flexible time windows (e.g., 2 seconds before or after an utterance) may lead to the recovery of some highly aligned moments that were not captured by the present analyses. More flexible time windows may help to identify the episodes that are likely to be most valuable for learning and that may be underestimated in our analyses. Analyzing longer time windows would also help to determine the role of alignment dynamics over the course of back-and-forth exchange, such as whether the caregiver or the child initiates a bid for joint attention (see Bianco et al., 2025). Averaging across the embeddings of neighboring utterances and frames during such exchanges could lead to increased vision–language alignment as well (He et al., 2025).

One additional direction for future work relates to the fact that CLIP conducts a center crop on the image before encoding it; as such, information that lies outside of the center crop (including, notably, text on a book during shared book reading that may appear at the bottom of the visual field) would not be used to determine the CLIP alignment score. An initial exploration of alternative image processing steps, such as padding and squashing, to include information outside of the center crop did produce moderate differences in 4AFC accuracy. In future work, analyzing CLIP alignment scores between an utterance and different segments of the visual field will also allow us to understand the distribution of aligned visual information across the visual field.

More broadly, this work focused on at-home recordings of monolingual English-speaking

families in the United States. We suspect that visual–linguistic alignment properties may differ depending on context—for example, due to cultural differences in how caregivers draw infants’ attention to objects (Senzaki et al., 2016). Collection and analyses of data from broader and more diverse populations would help to determine the generalizability of these findings.

Overall, our results suggest that CLIP alignment score is a meaningful metric to estimate the vision–language alignment in naturalistic egocentric videos, but that the alignment in young children’s everyday environments is relatively infrequent compared to alignment in modern machine learning datasets. Nonetheless, this alignment varies in systematic ways based on speaker, utterance properties, and features of the lemmas in those utterances. Future work exploring the dynamics, distribution, and predictors of variation in alignment will help to better capture the complex landscape of children’s early visual and linguistic input, and therefore inform our theories and models of language learning in young children.

Acknowledgments

We gratefully acknowledge the families who participated in the BabyView Dataset. This work was funded by an NIH R00HD108386 grant to B.L., by a grant from Schmidt Futures, by a gift from Meta, by the Stanford Center for the Study of Language and Information John Crosby Olney Fund, and by the Stanford Human-Centered AI Initiative (HAI) Hoffman-Yee grant program.

References

- Bergelson, E., & Aslin, R. N. (2017). Nature and origins of the lexicon in 6-mo-olds. *Proceedings of the National Academy of Sciences*, 114(49), 12916–12921.
- Bianco, C., Pang, J., & Yu, C. (2025, September). Contingent behavior during caregiver-child interaction improves the quality of word learning opportunities. In *2025 IEEE International Conference on Development and Learning (ICDL)* (pp. 1–6). doi: 10.1109/ICDL63968.2025.11204427
- Brysbaert, M., Warriner, A. B., & Kuperman, V. (2014, October). Concreteness ratings for 40 thousand generally known English word lemmas. *Behavior Research Methods*, 46(3), 904–911. Retrieved from <http://dx.doi.org/10.3758/s13428-013-0403-5> doi: 10.3758/s13428-013-0403-5
- Cortese, M. J., & Fuggetta, A. (2004, August). Imageability ratings for 3,000 monosyllabic words. *Behavior Research Methods, Instruments, & Computers*, 36(3), 384–387. Retrieved from <http://dx.doi.org/10.3758/BF03195585> doi: 10.3758/bf03195585
- Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2021). *Variability and consistency in early language learning: The wordbank project*. The MIT Press. Retrieved from <http://dx.doi.org/10.7551/mitpress/11577.001.0001> doi: 10.7551/mitpress/11577.001.0001
- Frank, M. C., Tenenbaum, J. B., & Fernald, A. (2013, January). Social and discourse contributions to the determination of reference in cross-situational word learning. *Language Learning and Development*, 9(1), 1–24. Retrieved from <http://dx.doi.org/10.1080/15475441.2012.707101> doi: 10.1080/15475441.2012.707101
- Gandhi, S., von Platen, P., & Rush, A. M. (2023). *Distil-whisper: Robust knowledge distillation via large-scale pseudo labelling*. Retrieved from <https://arxiv.org/abs/2311.00430>

He, Z. W., Trott, S., & Khosla, M. (2025). Seeing through words, speaking through pixels: Deep representational alignment between vision and language models. In *Proceedings of the 2025 conference on empirical methods in natural language processing* (pp. 35645–35660).

Jina AI. (2023). *Clip-as-service*. Retrieved from <https://clip-as-service.jina.ai/>

Khanna, M. M., & Cortese, M. J. (2021, May). How well imageability, concreteness, perceptual strength, and action strength predict recognition memory, lexical decision, and reading aloud performance. *Memory*, 29(5), 622–636. Retrieved from
<http://dx.doi.org/10.1080/09658211.2021.1924789> doi:
10.1080/09658211.2021.1924789

Lavechin, M., Bousbib, R., Bredin, H., Dupoux, E., & Cristia, A. (2020, October). An Open-Source Voice Type Classifier for Child-Centered Daylong Recordings. In *Proceedings of Interspeech 2020* (pp. 3072–3076). Shanghai, China. doi: 10.21437/Interspeech.2020-1690

Lin, T.-Y., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., ... Dollár, P. (2015). *Microsoft coco: Common objects in context*. Retrieved from
<https://arxiv.org/abs/1405.0312>

Long, B., Sparks, R. Z., Xiang, V., Stojanov, S., Yin, Z., Keene, G. E., ... others (2024). The babyview dataset: High-resolution egocentric videos of infants' and young children's everyday experiences. *arXiv preprint arXiv:2406.10447*.

Lynott, D., Connell, L., Brysbaert, M., Brand, J., & Carney, J. (2020, December). The lancaster sensorimotor norms: multidimensional measures of perceptual and action strength for 40, 000 english words. *Behavior Research Methods*, 52(3), 1271–1291. Retrieved from
<http://dx.doi.org/10.3758/s13428-019-01316-z> doi: 10.3758/s13428-019-01316-z

Marchman, V. A., & Fernald, A. (2008, May). Speed of word recognition and vocabulary knowledge in infancy predict cognitive and language outcomes in later childhood.

- Developmental Science*, 11(3). Retrieved from
<http://dx.doi.org/10.1111/j.1467-7687.2008.00671.x> doi:
10.1111/j.1467-7687.2008.00671.x
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., ... others (2021). Learning transferable visual models from natural language supervision. In *International conference on machine learning* (pp. 8748–8763).
- Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., & Sutskever, I. (2022). *Robust speech recognition via large-scale weak supervision*. Retrieved from
<https://arxiv.org/abs/2212.04356>
- Schock, J., Cortese, M. J., & Khanna, M. M. (2011, October). Imageability estimates for 3,000 disyllabic words. *Behavior Research Methods*, 44(2), 374–379. Retrieved from
<http://dx.doi.org/10.3758/s13428-011-0162-0> doi: 10.3758/s13428-011-0162-0
- Senzaki, S., Masuda, T., Takada, A., & Okada, H. (2016, January). The communication of culturally dominant modes of attention from parents to children: A comparison of canadian and japanese parent-child conversations during a joint scene description task. *PLOS ONE*, 11(1), e0147199. Retrieved from <http://dx.doi.org/10.1371/journal.pone.0147199> doi: 10.1371/journal.pone.0147199
- Sepuri, T., Aw, K. L., Tan, A. W. M., Sparks, R. Z., Marchman, V. A., Frank, M. C., & Long, B. (2025, October 10). Characterizing young children's everyday activities using video question-answering models. *PsyArXiv preprint*. Retrieved from
https://doi.org/10.31234/osf.io/gndy9_v1 doi: 10.31234/osf.io/gndy9_v1
- Straka, M., Hajič, J., & Straková, J. (2016, May). UDPipe: Trainable pipeline for processing CoNLL-U files performing tokenization, morphological analysis, POS tagging and parsing. In N. Calzolari et al. (Eds.), *Proceedings of the tenth international conference on language*

- resources and evaluation (LREC'16)* (pp. 4290–4297). Portorož, Slovenia: European Language Resources Association (ELRA). Retrieved from <https://aclanthology.org/L16-1680/>
- Sullivan, J., Mei, M., Perfors, A., Wojcik, E., & Frank, M. C. (2021, May). SAYCam: A Large, Longitudinal Audiovisual Dataset Recorded From the Infant's Perspective. *Open Mind: Discoveries in Cognitive Science*, 5, 20–29. doi: 10.1162/opmi_a_00039
- Vong, W. K., & Lake, B. M. (2025). On the robustness of modeling grounded word learning through a child's egocentric input. *arXiv preprint arXiv:2507.14749*.
- Vong, W. K., Wang, W., Orhan, A. E., & Lake, B. M. (2024). Grounded language acquisition through the eyes and ears of a single child. *Science*, 383(6682), 504–511.
- Yoshida, H., & Smith, L. B. (2008). What's in view for toddlers? using a head camera to study visual experience. *Infancy*, 13(3), 229–248.
- Young, P., Lai, A., Hodosh, M., & Hockenmaier, J. (2014). From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *TACL*, 2, 67–78.
- Yu, C., & Smith, L. B. (2007). Rapid word learning under uncertainty via cross-situational statistics. *Psychological science*, 18(5), 414–420.
- Zhang, B., Li, K., Cheng, Z., Hu, Z., Yuan, Y., Chen, G., ... Zhao, D. (2025). VideoLLaMA 3: Frontier multimodal foundation models for image and video understanding. *arXiv preprint arXiv:2501.13106*. Retrieved from <https://arxiv.org/abs/2501.13106>
- Zhang, Y., Yurovsky, D., & Yu, C. (2021). Cross-situational learning from ambiguous egocentric input is a continuous process: Evidence using the human simulation paradigm. *Cognitive science*, 45(7), e13010.