

# BabyVLM: Data-Efficient Pretraining of VLMs Inspired by Infant Learning\*

Shengao Wang      Arjun Chandra      Aoming Liu  
Boston University   Boston University   Boston University  
wsashawn@bu.edu      ac25@bu.edu      amliu@bu.edu

Venkatesh Saligrama      Boqing Gong  
Boston University      Boston University  
srv@bu.edu      bgong@bu.edu

## Abstract

*Human infants rapidly develop visual reasoning skills from minimal input, suggesting that developmentally inspired pretraining could significantly enhance the efficiency of vision-language models (VLMs). Although recent efforts have leveraged infant-inspired datasets like SAYCam, existing evaluation benchmarks remain misaligned—they are either too simplistic, narrowly scoped, or tailored for large-scale pretrained models. Additionally, training exclusively on SAYCam overlooks the broader, diverse input from which infants naturally learn. To address these limitations, we propose BabyVLM, a novel framework comprising diverse in-domain evaluation benchmarks and a synthetic training dataset created via child-directed transformations of existing datasets. We demonstrate that VLMs trained with our synthetic dataset achieve superior performance on BabyVLM tasks compared to models trained solely on SAYCam or general-purpose data of the SAYCam size. BabyVLM thus provides a robust, developmentally aligned evaluation tool and illustrates how compact models trained on carefully curated data can generalize effectively, opening pathways toward data-efficient vision-language learning paradigms.*

## 1. Introduction

We propose a novel framework, BabyVLM, for data-efficient pretraining of vision-language models (VLMs). To this end we introduce methods for creating minimal yet naturalistic data—akin to the input human infants receive—as well as diverse in-domain evaluation benchmarks. By carefully curating the training data, we show that our method yields more robust, baby-like representations compared to training on general-purpose corpora, and can further serve

as a template for resource-efficient model training in other specialized domains.

**Challenges in Current VLM Training.** Vision-Language Models have advanced rapidly in recent years [15, 27, 38, 49, 58], but these advancements often rely on massive datasets and prohibitively expensive computational resources. For instance, training large-scale models such as LLaMA [53] or LLaVA [27] can require thousands of GPU hours [14, 46]. Such demands pose fundamental barriers for independent researchers with limited resources, highlighting the need for more accessible pretraining methods.

**Lessons from Infant Learning.** Human infants, by contrast, rapidly acquire complex cognitive and perceptual skills from minimal data and limited environmental exposure [20, 45]. This exceptional efficiency implies that robust representations can be learned from small, carefully curated datasets when these datasets closely mimic natural developmental conditions. Recognizing this, researchers have begun curating datasets such as SAYCam [47], which provides egocentric audiovisual recordings of infants aged 6–32 months. Although our work primarily utilizes SAYCam, other developmentally inspired datasets such as BabyView [28] also support this approach. Our framework capitalizes on these insights, suggesting that intentionally constrained, naturalistic training scenarios can yield efficient, highly generalizable models.

**The Evaluation Gap.** Despite the promise of data-efficient VLM training inspired by infant learning, evaluating such compact models remains a critical challenge. Current benchmarks—such as VQA [2], Winoground [52], and COCO [24]—were designed for large-scale models trained on massive datasets, assessing capabilities that exceed those reasonably achievable by developmentally plausible, compact models. For instance, the Labeled-S benchmark [34], which specifically targets SAYCam data, evaluates only a single classification task and thus cannot comprehensively

\*Project website: [shawnking98.github.io/BabyVLM](https://shawnking98.github.io/BabyVLM)

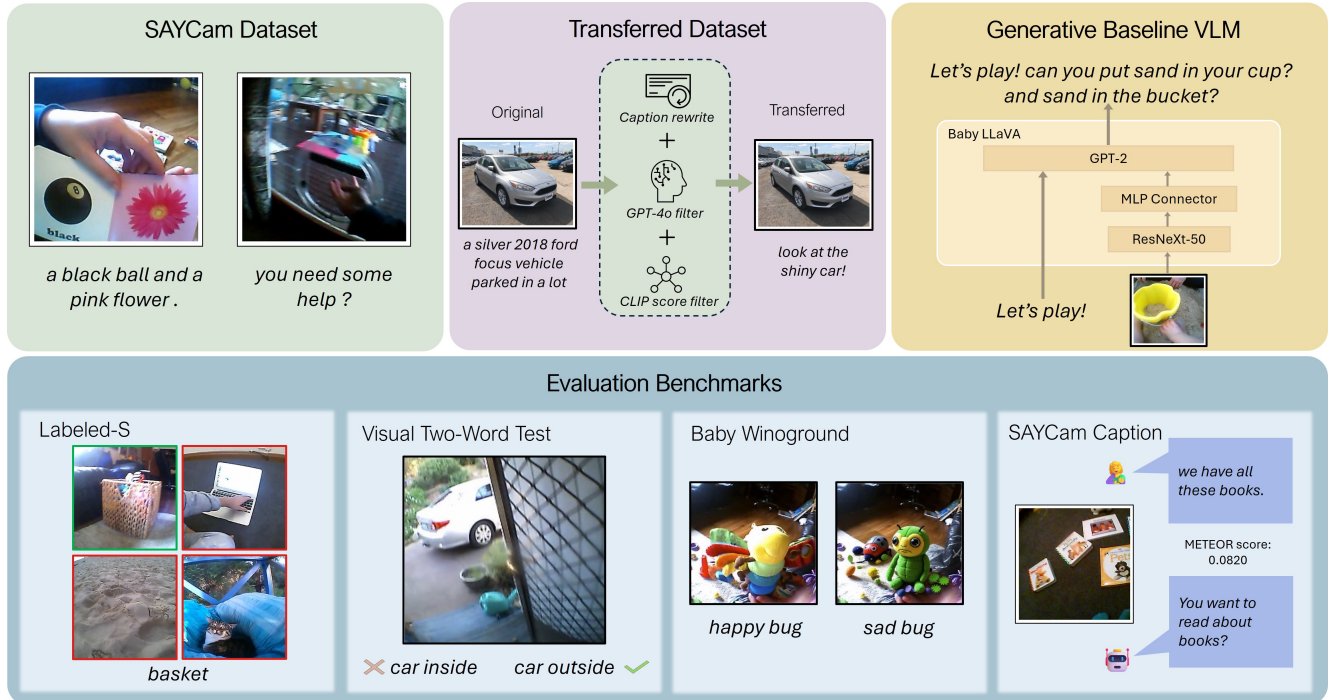


Figure 1. We introduce BabyVLM, a developmentally inspired framework derived from SAYCam, consisting of the original SAYCam dataset [47], a transferred training dataset, a generative baseline VLM, and four evaluation benchmarks.

Benchmark	Task Diversity	Baby-like	In-domain
General purpose (VQA [2], Winoground [52], etc.)	✓	✗	✗
DevBench [48]	✓	✓	✗
Labeled-S [34]	✗	✓	✓
ModelVsBaby [44]	✗	✓	✗
MEWL [17]	✓	✓	✗
BabyVLM	✓	✓	✓

Table 1. Representative features of existing multimodal evaluation benchmarks. **Task Diversity:** The benchmark should include diverse tasks that assess different aspects of a vision-language model’s capability rather than focusing solely on simple tasks (e.g., object classification). **Baby-like:** The benchmark should align with cognitive and linguistic developmental stages observed in human babies. **In-domain:** The testing samples should come from the same data domain as the training dataset, ensuring that evaluation results reflect the model’s ability to generalize within a realistic learning environment.

measure broader vision-language capabilities. Conversely, developmental psychology benchmarks [19, 34, 48] tend to be overly simplistic or not directly relevant to the infant-inspired training data. As summarized in Table 1, this evaluation gap underscores the need for comprehensive, developmentally aligned benchmarks—precisely the gap addressed by our proposed framework, BabyVLM.

To bridge this evaluation gap and realize our goal of data-efficient, developmentally aligned VLM pretraining, we offer three main contributions:

- **In-Domain Evaluation Tasks.** We design three novel evaluation tasks derived from the SAYCam dataset.

These tasks are tailored to reflect the cognitive and perceptual abilities typical of early human development, enabling diverse and meaningful evaluation of compact models trained on developmentally plausible data.

- **Synthetic Data Augmentation.** We introduce a data distillation approach to address the inherent limitations of existing small-scale datasets. By synthesizing simplified, child-directed versions of existing datasets like CC3M [42] using GPT-4o [15], we create training data that more closely mirrors the linguistic and visual complexity encountered by infants.
- **BabyLLaVA: Generative Model Trained from**

**Scratch.** Inspired by recent methods [27, 54], we present BabyLLaVA, the first generative VLM trained entirely on developmentally plausible data. BabyLLaVA demonstrates that compact generative models, when trained on intentionally constrained and naturalistic data, can produce robust, baby-oriented responses from the input of baby viewpoints.

Collectively, these contributions not only demonstrate effective, resource-efficient pretraining within our specific domain but also offer insights that can inform efficient paradigms across diverse applications, thereby lowering barriers to foundational model research.

## 2. Related Work

**Vision-Language Models.** Large vision-language models (VLMs) [4, 15, 27, 29, 38, 49] have significantly advanced multimodal understanding by integrating visual and linguistic data for various tasks, including image captioning, visual question answering, and conversational interaction. Early influential models such as CLIP [38] leveraged contrastive learning paradigms, effectively aligning visual and textual representations within a unified embedding space. More recent generative frameworks, such as LLaVA series [21–23, 25–27, 63], have combined pre-trained visual encoders [38, 60] with large language models [3, 16, 50, 65], enabling more advanced conversational interactions and multimodal generative capabilities. However, these models typically require extensive computational resources and large-scale datasets. In contrast, our approach specifically targets compact generative VLMs trained exclusively on developmentally plausible datasets, providing a framework to improve data efficiency and better align model training with cognitive development processes observed in human infants.

**Developmentally Inspired Learning.** Human infants exhibit remarkable efficiency in acquiring language and visual concepts from limited and naturalistic input, inspiring substantial research into developmentally plausible training paradigms. Early influential datasets like CHILDES [30] facilitated initial explorations into language acquisition through linguistic recordings across diverse languages [1, 8]. Recent initiatives, including the BabyLM Challenge [9, 56], further encouraged the development of models trained on language data scales comparable to those encountered by infants. Extending these ideas into multimodal contexts, datasets such as SAYCam [47] and BabyView [28] have provided egocentric audiovisual data, enabling research that progresses from single-modality learning [33, 34, 36, 43, 55] to visually grounded language acquisition [54, 66]. Our work distinctly builds upon these foundations by explicitly creating synthetic, child-directed multimodal data from general-domain sources, addressing the limitations inherent

in existing infant-inspired datasets and exploring the potential of compact generative VLMs trained in developmentally realistic conditions.

**Multimodal Benchmarks.** Existing multimodal evaluation benchmarks can be broadly classified as general-purpose or developmentally inspired. General-purpose benchmarks, such as Visual Question Answering (VQA) [2, 61] and Winoground [52], evaluate advanced visio-linguistic integration but typically rely on large-scale, non-developmental datasets, rendering them unsuitable for compact models trained on limited developmental data. Conversely, developmentally inspired benchmarks—such as Labeled-S [34], ModelVsBaby [44], DevBench [48], and MEWL [17]—are more aligned with early cognitive processes but often limited to simplistic classification tasks or utilize data not fully reflective of the training domain. Our work explicitly addresses these gaps by proposing cognitively nuanced benchmarks directly aligned with the developmental data domain, thereby enabling accurate and relevant assessments of compact vision-language models trained from minimal, developmentally appropriate multimodal inputs.

## 3. Framework

Our proposed framework, BabyVLM, aims to facilitate resource-efficient pretraining and developmentally aligned evaluation of compact VLMs inspired by the minimal yet highly informative learning environments of human infants. To achieve this, BabyVLM comprises: (1) a filtered subset of baby-egocentric audio-visual recording from the SAYCam dataset, (2) a novel synthetic training dataset specifically crafted to reflect infant-directed linguistic and visual experiences, (3) a generative baseline model, BabyLLaVA, trained entirely on this developmentally plausible data, and (4) three novel evaluation benchmarks explicitly tailored to assess multimodal reasoning aligned with early cognitive stages, plus Labeled-S [34], an existing classification benchmark. Please refer to Figure 1 for an overview of our framework.

**Design Principles for the BabyVLM Framework.** A central goal of BabyVLM is to ensure realistic alignment with developmental constraints characteristic of early visual-language learning. To this end, we adopt the following concise guiding principles:

- **Developmentally Appropriate Complexity:** Tasks reflect cognitive capabilities typical of early developmental stages (e.g., basic object and action recognition, simple compositional reasoning), explicitly avoiding tasks requiring more complex reasoning.
- **Limited Generalization Beyond Early Development:** Models should demonstrate intentionally constrained generalization, ignoring performance beyond realistic developmental boundaries.

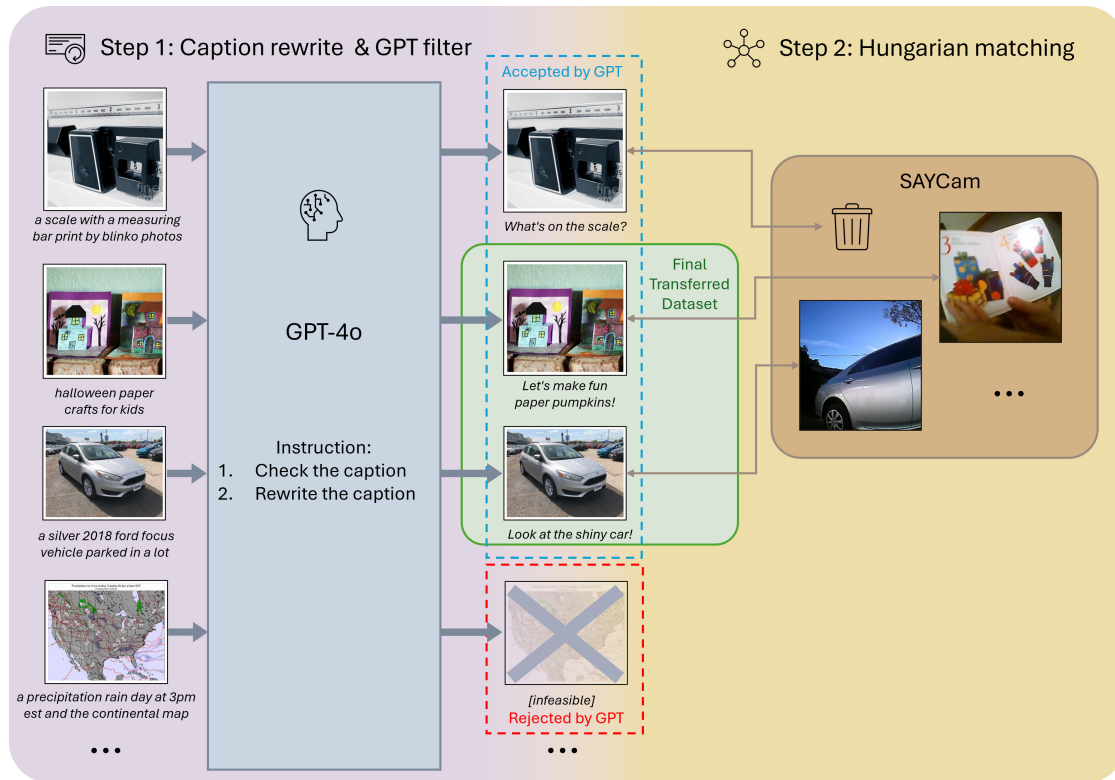


Figure 2. Pipeline for generating the transferred dataset. **Step 1:** We prompt GPT-4o to check whether an input caption is describing something a child would see in daily life and transfer the original image captions into simpler, child-directed utterances. **Step 2:** We use the CLIP similarity score as a metric to represent the distance between two images, and then conduct Hungarian matching to select a small subset of the transferred dataset that is visually aligned with SAYCam images.

- **Linguistic and Visual Simplicity:** Dataset construction explicitly emphasizes simple vocabulary, concrete visual scenes, and straightforward grammatical structures consistent with child-directed interactions.

These principles collectively ensure that the resulting BabyVLMs remain authentic representations of early-stage developmental models, with their effectiveness empirically confirmed in our out-of-domain task evaluations, which are presented in the supplementary material.

### 3.1. Datasets

**Filtered SAYCam Dataset.** The original SAYCam dataset [47] comprises egocentric audiovisual recordings. Following Vong et al. [54], we extract child-directed speech and sample video into image-utterance pairs. To improve data quality, we filter pairs using CLIP similarities [38], retaining only those above a 0.2 threshold to ensure high image-text relevance. This results in 67K pairs. Examples are provided in the supplementary material.

**Transferred Synthetic Training Dataset.** While the SAYCam dataset provides naturally curated developmental data, there are inherent limitations of relying on this

dataset exclusively. First, videos in SAYCam were recorded in 60 to 80 minute sessions twice a week, documenting only a small subset of each child’s developmental experience. Moreover, due to practical constraints, SAYCam videos were often recorded at fixed times each week, reducing variation in the infant’s recorded environment and further limiting our ability to capture the diverse multi-modal input streams from which babies learn [59]. To address these limitations, we created a synthetic auxiliary training corpus by adapting general-purpose multimodal datasets—CC3M [42], LAION [41], and SBU [32]—to match infant learning conditions.

Our approach comprises two steps, as illustrated in Figure 2. In the first step, we prompt GPT-4o to rewrite captions into concise, child-directed utterances, emulating speech typically used with two-year-olds. GPT-4o also flags which image-caption pairs are misaligned with an infant’s daily experience for exclusion. By emphasizing everyday vocabulary, simple grammar, and concrete objects and actions, we ensured linguistic alignment with early-stage learners.

In the second step, to further maintain visual consis-

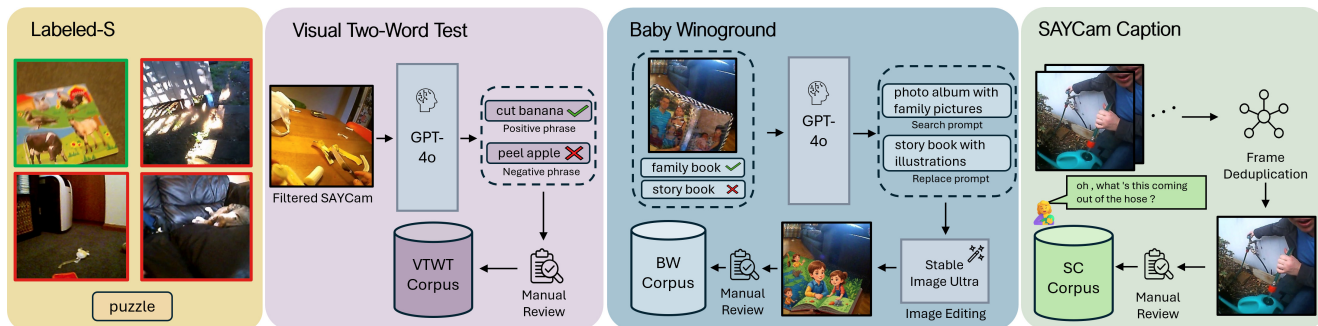


Figure 3. Illustrations of in-domain evaluation benchmarks in the BabyVLM framework. **Labeled-S:** The category label must be matched to the target referent among 4 candidates. **Visual Two-Word Test:** The positive phrase must be matched to the image. Positive and negative phrases are generated by GPT-4o. **Baby Winoground:** The positive and negative phrases must be matched with their corresponding images. Negative images are generated by Stable Diffusion [11], with prompts enhanced by GPT-4o. **SAYCam Caption:** The generated image caption must match the ground truth image caption. All image-caption pairs come from a de-duplicated subset of the SAYCam test split.

tency, we apply the Hungarian algorithm [18] and utilize the CLIP similarity as a distance metric to select a subset of images resembling SAYCam. The number of selected samples matches that of the filtered SAYCam training set, resulting in a dataset that maintains visual alignment while balancing diversity and domain relevance. More details of the transformation guidelines and examples of rewritten captions are included in the supplementary material.

### 3.2. BabyLLaVA: Generative VLM Baseline

We then train a compact VLM, called BabyLLaVA, using the compiled dataset. Inspired by LLaVA [27], BabyLLaVA integrates a compact language model (GPT-2 [37], 7M parameters) and vision encoder (ResNeXt-50 [57], 23M parameters) through a lightweight multilayer perceptron connector. To examine the impact of model capacity, we also provide a larger variant (Llama-1.1B [62] + ViT-L [10]). Consistent with our guiding principles, BabyLLaVA’s compact size and simplified architecture explicitly limit the model’s complexity, aligning closely with the realistic developmental constraints of early-stage learners. Training exclusively on our provided developmental dataset, BabyLLaVA provides a suitable baseline model to evaluate the effectiveness of developmentally aligned multimodal learning. Additional details can be found in the supplementary material.

### 3.3. Evaluation Tasks

To rigorously evaluate multimodal reasoning within the intended developmental scope, we introduce three novel benchmarks explicitly designed around cognitive milestones typical of early learners, in addition to the existing Labeled-S task [34]. These tasks deliberately embody simplicity and developmental appropriateness, ensuring alignment with our guiding principles (Figure 3).

**Labeled-S.** The Labeled-S testing dataset was first introduced by Orhan et al. [34] and has since been used in studies such as [33, 54]. Leveraging SAYCam annotations, Orhan et al. manually curated a dataset comprising approximately 58K labeled frames across 26 categories. Following Vong et al. [54] and aligning with standard child testing procedures [31], we use a subset of Labeled-S and evaluate models by presenting a target category label alongside four candidate images, requiring the model to identify the correct match.

**Visual Two-Word Test (VTWT).** Inspired by the linguistic milestone known as the “two-word” stage (typically 18–24 months) [6, 7, 35, 39], this task assesses compositional semantic reasoning. Models must correctly match SAYCam images with appropriate two-word phrases (e.g., “wash cup” vs. “fill cup”). Starting with a sub-sampled test split from SAYCam, we generate 5117 phrase pairs using GPT-4o. These are manually reviewed for linguistic and visual appropriateness, yielding 967 final pairs. Table 2 summarizes the distribution of phrase types tested. Detailed annotation guidelines and procedures are provided in the supplementary material.

Types of Differences	Proportions (%)	Examples
Verb	27.2	“wash cup” vs. “fill cup”
Adjective	21.4	“happy faces” vs. “sad faces”
Noun	17.0	“car outside” vs. “bike outside”
Verb + Noun	21.4	“spread jam” vs. “cut bread”
Adjective + Noun	11.5	“yellow flower” vs. “green tree”
Verb + Adjective	1.5	“small frown” vs. “big smile”

Table 2. Proportions and examples of each type of differences between positive and negative phrases in the Visual Two-Word Test.

**Baby Winoground.** Extending VTWT, Baby Winoground tests more advanced visio-linguistic compositional reason-

ing. Inspired by Winoground [52], this task presents two images and two corresponding phrases (one positive, one negative). Negative images were created by modifying specific visual elements of original images through targeted prompts provided to Stability AI’s Stable Image Ultra model [11]. All samples undergo a manual review to ensure minimal domain gaps and precise visio-linguistic mappings, resulting in 365 high-quality test samples. Full prompt-engineering details and validation methods are in the supplementary material.

For evaluation, we adopt a modified group score to better understand model performance under distributional shifts. Each example consists of two image-phrase pairs: one positive pair (a real SAYCam frame and a matching phrase) and one negative pair (a modified phrase and a synthetic image). The standard group score requires the model to correctly identify both the positive and negative pairs over four pairwise comparisons. To gain finer insight, we break this down into two context-conditioned variants:

- **Positive Context Score:** Measures whether the model correctly identifies the matching pair when using the original SAYCam image or phrase as context.
- **Negative Context Score:** Measures the same, but when using the synthetic image or modified phrase as context.

**SAYCam Caption.** The SAYCam Caption benchmark evaluates generative captioning skills by requiring models to generate accurate, contextually relevant descriptions for SAYCam images. Captions are sourced from the test split of SAYCam, using child-directed utterances as ground truth. To refine the dataset, we deduplicate frames with identical utterances, retaining only those with the highest CLIP image-caption similarities. This yields 1598 distinct image-caption pairs, which are then manually verified, resulting in 294 final test samples. Evaluation is performed using the METEOR metric [5]. This task measures a model’s ability to generate coherent, semantically appropriate child-directed descriptions. Examples of test samples are provided in the supplementary material.

## 4. Experiments

Our experimental evaluation aims to compare VLM architectures and training paradigms within a developmentally plausible setting, investigate the effectiveness of our synthetic child-directed dataset, and perform a fine-grained analysis of compositional reasoning. To validate that our baby models align with the cognitive and linguistic limitations of early-stage learners, we also assess baby models on tasks that exceed typical infant-level developmental capacities, presented in the supplementary materials.

### 4.1. In-Domain Benchmark Results

We begin by evaluating multiple models, including baby models trained purely on SAYCam (BabyLLaVA, CVCL [54]) and larger upper bound models that are either directly used out of the box (LLaVA-v1.5-7B [25], CLIP-large [38]) or further fine-tuned on our SAYCam data (LLaVA-v1.5-7B-ft). These models are assessed on four in-domain benchmarks: Labeled-S, Visual Two-Word Test (VTWT), Baby Winoground, and SAYCam Caption. Table 3 summarizes these results.

Notably, CVCL—a contrastive model—consistently outperforms the generative BabyLLaVA model across most tasks. This observation aligns with existing literature [12, 13, 51, 64], suggesting that contrastive models may be better suited to discriminative tasks, possibly due to their direct objective of learning joint visual-textual alignment. However, generative models like BabyLLaVA demonstrate reasonable performance on simpler compositional tasks such as VTWT, indicating substantial potential for improvement on more sophisticated compositional tasks like Baby Winoground. Interestingly, the larger BabyLLaVA-Llama variant performs similarly or even worse than BabyLLaVA-GPT2, despite being about 50 times larger—suggesting overfitting to the limited data; as such, we default to using the GPT2 version for all the remaining discussion. In particular, Baby Winoground reveals a stark asymmetry: baby models perform above chance when reasoning from in-distribution (positive) context, but below chance from out-of-distribution (negative) context, highlighting a systematic failure under distribution shift. Moreover, generative captioning, measured by SAYCam Caption scores, remains challenging for all models, emphasizing the additional complexity inherent in generating full linguistic descriptions from minimal data.

### 4.2. Transferred Dataset Ablation

We next perform an ablation study (Table 4) comparing models trained on several different dataset settings: using our filtered SAYCam dataset only (*filtered-only*), SAYCam plus our transferred, child-directed dataset (*filtered-aug*), transferred dataset only (*aug-only*), SAYCam plus randomly selected general-domain dataset of the same size as the transferred dataset (*filtered-random*), random general-domain dataset only (*random-only*), and a bigger filtered SAYCam dataset whose number of samples is doubled by relaxing the threshold (*filtered-double*). Additional analysis regarding data efficiency can be found in the supplementary material.

We observe clear performance improvements in CVCL and BabyLLaVA when using our carefully curated dataset compared to random augmentation, particularly on compositional reasoning tasks such as VTWT and Baby Winoground. These results indicate that explicitly adapt-

Category	Model	Labeled-S	VTWT	Baby Winoground			SAYCam Caption
				Overall	Pos. Ctx	Neg. Ctx	
<b>Upper Bound Models</b>	LLaVA-v1.5-7B	0.7400	0.7851	0.4274	0.6575	0.6301	0.1657
	LLaVA-v1.5-7B-ft	0.6591	0.7038	0.3205	0.5644	0.6027	0.1798
	CLIP-large	0.7100	0.8625	0.6740	0.7315	0.8603	N/A
<b>Baby Models</b>	BabyLLaVA-GPT2	0.4195	0.6252	0.0658	0.3890	0.2301	<b>0.1379</b>
	BabyLLaVA-Llama	0.4200	0.6029	0.0521	0.3945	<b>0.2466</b>	0.1287
	CVCL	<b>0.6086</b>	<b>0.6494</b>	<b>0.0932</b>	<b>0.5068</b>	0.2246	N/A
<b>Random Guess</b>	-	0.25	0.5	0.1667	0.25	0.25	N/A

Table 3. Evaluation results of in-domain tasks, where a higher score indicates better performance. For Labeled-S, we use the same target and foil testing samples as [34] and report accuracy. For the Visual Two-Word Test (VTWT), we report accuracy. For Baby Winoground, we report the group score for different contexts. For SAYCam Caption, we report the METEOR score.

Model	Labeled-S	VTWT	Baby Winoground			SAYCam Caption
			Overall	Pos. Ctx	Neg. Ctx	
CVCL-filtered	0.6086	0.6494	0.0932	0.5068	0.2246	N/A
CVCL-filtered-aug	0.5805	0.7021	0.2027	0.4657	0.4493	N/A
CVCL-filtered-random	0.6023	0.6835	0.1068	0.4739	0.2958	N/A
BabyLLaVA-filtered	0.4195	0.6252	0.0658	0.3890	0.2301	0.1379
BabyLLaVA-filtered-aug	0.5364	0.6933	0.0822	0.3726	0.3096	0.1592
BabyLLaVA-filtered-random	0.5155	0.6553	0.0877	0.3616	0.2712	0.1778
BabyLLaVA-filtered-double	0.4659	0.6383	0.0658	0.2411	0.2301	0.1799
BabyLLaVA-aug-only	0.5000	0.6239	0.0630	0.3370	0.3644	0.0615
BabyLLaVA-random-only	0.4400	0.5098	0.0548	0.2904	0.3836	0.0615

Table 4. Ablation study of our transferred dataset on in-domain tasks. **XX-filtered**: Only SAYCam. **XX-filtered-aug**: SAYCam + child-directed transferred data. **XX-filtered-random**: SAYCam + random general-domain data. **XX-filtered-double**: SAYCam with double-numbered samples. **XX-aug-only**: Only random general-domain data. **XX-random-only**: Only random general-domain data.

ing general-domain data to reflect the linguistic simplicity and visual content of infant environments significantly enhances the data efficiency and overall alignment of the resulting models. Using only the transferred data degrades the model performance, which shows the importance of having the original SAYCam dataset as an anchor. Notably, for Baby Winoground, training with the transferred dataset substantially improves performance in the negative context setting, despite a slight drop in the positive context score; while the randomly selected dataset also improves the negative context score, the gains are smaller, indicating that our transferred dataset is more effective in helping baby models generalize to broader domains. In contrast, improvements in generative captioning remain modest, suggesting that further refinements, such as enriching the linguistic variety or introducing narrative structures, could improve generative performance.

### 4.3. Assessing Language Bias in VTWT

To confirm the robustness of our VTWT benchmark, we conducted an experiment removing visual context entirely (Table 5). The resulting performance drop from around 78% accuracy (with image) to approximately random chance (53% without image) demonstrates that the task cannot be solved through language biases alone. This validates VTWT as a rigorous evaluation of genuine multi-modal compositional reasoning rather than simple linguistic pattern-matching, confirming the robustness and appropriateness of our benchmark.

Model	VTWT (w/ image)	VTWT (w/o image)
LLaVA-v1.5-7B	0.7851	0.5307
BabyLLaVA	0.6252	0.5360

Table 5. Ablation study of language-only bias on VTWT.

## 4.4. Investigating Compositional Reasoning

We further dissect the VTWT performance by examining model accuracy on different types of compositional differences (noun, verb, adjective, or combinations thereof) in Table 6.

Type of Difference	CVCL-filtered	CVCL-filtered-aug	CVCL-filtered-random
Verb	0.6221	0.6564	0.7404
Adjective	0.5507	0.6086	0.5797
Noun	0.7317	0.7682	0.7073
Verb + Noun	0.7087	0.7572	0.6699
Adjective + Noun	0.6936	0.7927	0.7297
Verb + Adjective	0.4285	0.6428	0.7857

Table 6. Performance breakdown on VTWT by part-of-speech differences.

We observe that all three model variants perform worse on adjective differences than adjective + noun differences. We suspect this is because adjective differences alone are often less visually explicit in an image, and the presence of an additional noun difference helps the models disambiguate these cases.

Additionally, models trained exclusively on developmentally plausible data (CVCL-filtered and CVCL-filtered-aug) exhibit distinct performance patterns. For single-component differences (the first three rows of Table 6), both models achieve their highest performance on noun differences and perform worse on verb and adjective differences (e.g., 76% vs. 65% and 60% respectively for CVCL-filtered-aug). This result aligns with linguistic development findings from [6, 40], which suggest that early-stage language learners use nouns at least twice as often as verbs and are also slower to acquire adjectives. However, similar phenomenon is not observed from CVCL-filtered-random.

The alignment between our empirical results and developmental psychology literature reinforces that our targeted synthetic data transformations effectively facilitate more robust baby-like representations — a central objective identified in our introduction.

## 4.5. Discussion

Overall, our experiments reinforce several key insights central to our initial narrative. First, child-directed transformations of general-domain datasets provide substantial gains within our developmentally plausible domain, validating our approach. However, generative models face heightened difficulties in compositional reasoning and full-sentence generative tasks, highlighting significant room for further development. Lastly, while the specialized infant-oriented approach offers promising efficiencies, it inherently limits performance in broader contexts. These findings suggest fruitful future directions, including expanding dataset richness, exploring hybrid generative-discriminative training methods, and generalizing our approach to other spe-

cialized domains, as envisioned in our introduction.

## 5. Conclusion and Future Work

In this work, we introduced BabyVLM, a framework for data-efficient pretraining and evaluation of compact vision-language models (VLMs) inspired by the developmental learning conditions of human infants. Our approach is grounded in explicitly enforcing developmental constraints on both data and model design, ensuring that baby models operate within a realistic cognitive scope. To achieve this, we curated a filtered subset of the SAYCam dataset, constructed a novel synthetic training dataset that aligns with child-directed language and visual experiences, and introduced three evaluation benchmarks designed to test multimodal reasoning at early developmental stages.

Our experiments validate the effectiveness of this approach. In-domain evaluations demonstrate that baby models can learn meaningful multimodal associations from developmentally appropriate data, while out-of-domain evaluations confirm their intentional constraints, preventing over-generalization beyond early cognitive capabilities. Notably, we find that the observed performance gaps between baby models and larger models arise from multiple factors—model capacity, task complexity, and data alignment—rather than capacity alone. This underscores the importance of dataset and task design in modeling early-stage learning.

Moving forward, our work opens several avenues for further research. First, expanding the dataset to incorporate additional multimodal learning signals—such as temporal context or richer object interactions—could further refine developmental modeling. Second, investigating hybrid models that balance generative and contrastive training may provide insights into optimizing learning efficiency in data-limited regimes. Lastly, our benchmarks and methodology can serve as a foundation for broader inquiries into how developmental constraints shape representation learning in neural models.

By establishing a principled framework for modeling early-stage multimodal learning, BabyVLM provides a meaningful step toward understanding and replicating data-efficient learning in artificial systems, with potential implications for both machine learning and cognitive science.

## References

- [1] Raquel G. Alhama, Caroline Rowland, and Evan Kidd. Evaluating word embeddings for language acquisition. In *Proceedings of the Workshop on Cognitive Modeling and Computational Linguistics*, pages 38–42, Online, 2020. Association for Computational Linguistics. 3
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi

- Parikh. VQA: Visual Question Answering. In *International Conference on Computer Vision (ICCV)*, 2015. 1, 2, 3
- [3] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. 3
- [4] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 3
- [5] Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72, 2005. 6
- [6] Stephanie Berk and Diane Lillo-Martin. The two-word stage: Motivated by linguistic or cognitive constraints? *Cognitive psychology*, 65(1):118–140, 2012. 5, 8
- [7] Amanda C Brandone, Sara J Salkind, Roberta Michnick Golinkoff, and Kathy Hirsh-Pasek. Language development. 2006. 5
- [8] Michael R. Brent and Timothy A. Cartwright. Distributional regularity and phonotactic constraints are useful for segmentation. *Cognition*, 61(1-2):93–125, 1996. 3
- [9] Leshem Choshen, Ryan Cotterell, Michael Y Hu, Tal Linzen, Aaron Mueller, Candace Ross, Alex Warstadt, Ethan Wilcox, Adina Williams, and Chengxu Zhuang. [call for papers] the 2nd babyLM challenge: Sample-efficient pretraining on a developmentally plausible corpus. *arXiv preprint arXiv:2404.06214*, 2024. 3
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 5
- [11] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024. 5, 6
- [12] Gregor Geigle, Radu Timofte, and Goran Glavaš. African or european swallow? benchmarking large vision-language models for fine-grained object classification. *arXiv preprint arXiv:2406.14496*, 2024. 6
- [13] Hulingxiao He, Geng Li, Zijun Geng, Jinglin Xu, and Yuxin Peng. Analyzing and boosting the power of fine-grained visual recognition for multi-modal large language models. *arXiv preprint arXiv:2501.15140*, 2025. 6
- [14] Jordan Hoffmann, Sebastian Borgeaud, Andreas Mensch, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022. 1
- [15] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 1, 2, 3
- [16] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. 3
- [17] Guangyuan Jiang, Manjie Xu, Shiji Xin, Wei Liang, Yujia Peng, Chi Zhang, and Yixin Zhu. Mewl: Few-shot multimodal word learning with referential uncertainty. In *International Conference on Machine Learning*, pages 15144–15169. PMLR, 2023. 2, 3
- [18] Roy Jonker and Ton Volgenant. A shortest augmenting path algorithm for dense and sparse linear assignment problems. In *DGORN/SOR: Papers of the 16th Annual Meeting of DGOR in Cooperation with NSOR/Vorträge der 16. Jahrestagung der DGOR zusammen mit der NSOR*, pages 622–622. Springer, 1988. 5
- [19] Talia Konkle, Timothy F Brady, George A Alvarez, and Aude Oliva. Conceptual distinctiveness supports detailed visual long-term memory for real-world objects. *Journal of experimental Psychology: general*, 139(3):558, 2010. 2
- [20] Alexander LaTourrette, Dana Michelle Chan, and Sandra R Waxman. A principled link between object naming and representation is available to infants by seven months of age. *Scientific reports*, 13(1):14328, 2023. 1
- [21] Bo Li, Hao Zhang, Kaichen Zhang, Dong Guo, Yuanhan Zhang, Renrui Zhang, Feng Li, Ziwei Liu, and Chunyuan Li. Llava-next: What else influences visual instruction tuning beyond data?, 2024. 3
- [22] Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. Llava-next: Stronger llms supercharge multimodal capabilities in the wild, 2024.
- [23] Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv preprint arXiv:2407.07895*, 2024. 3
- [24] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision—ECCV 2014: 13th European conference, Zurich, Switzerland, September 6–12, 2014, proceedings, part v 13*, pages 740–755. Springer, 2014. 1
- [25] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024. 3, 6
- [26] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024.
- [27] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 1, 3, 5
- [28] Bria Long, Violet Xiang, Stefan Stojanov, Robert Z Sparks, Zi Yin, Grace E Keene, Alvin WM Tan, Steven Y Feng, Chengxu Zhuang, Virginia A Marchman, et al. The

- babyview dataset: High-resolution egocentric videos of infants' and young children's everyday experiences. *arXiv preprint arXiv:2406.10447*, 2024. 1, 3
- [29] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, et al. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*, 2024. 3
- [30] Brian MacWhinney and Catherine Snow. The child language data exchange system: An update. *Journal of child language*, 17(2):457–472, 1990. 3
- [31] Dana McDaniel, Cecile McKee, and Helen Smith Cairns. *Methods for assessing children's syntax*. Mit Press, 1998. 5
- [32] Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million captioned photographs. *Advances in neural information processing systems*, 24, 2011. 4
- [33] A Emin Orhan and Brenden M Lake. Learning high-level visual representations from a child's perspective without strong inductive biases. *Nature Machine Intelligence*, 6(3):271–283, 2024. 3, 5
- [34] Emin Orhan, Vaibhav Gupta, and Brenden M Lake. Self-supervised learning through the eyes of a child. *Advances in Neural Information Processing Systems*, 33:9960–9971, 2020. 1, 2, 3, 5, 7
- [35] William O'Grady and Sook Whan Cho. First language acquisition. *Contemporary linguistics: An introduction*, pages 409–448, 2001. 5
- [36] Yulu Qin, Wentao Wang, and Brenden M Lake. A systematic investigation of learnability from single child linguistic input. *arXiv preprint arXiv:2402.07899*, 2024. 3
- [37] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 5
- [38] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 1, 3, 4, 6
- [39] Nicholas Ricciardi, Xuan Yang, and Rutvik H Desai. The two word test as a semantic benchmark for large language models. *Scientific Reports*, 14(1):21593, 2024. 5
- [40] C. Sandhofer and L. B. Smith. Learning adjectives in the real world: How learning nouns impedes learning adjectives. *Language Learning and Development*, 3(3):233–267, 2007. 8
- [41] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in neural information processing systems*, 35:25278–25294, 2022. 4
- [42] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of ACL*, 2018. 2, 4
- [43] Saber Sheybani, Himanshu Hansaria, Justin Wood, Linda Smith, and Zoran Tiganj. Curriculum learning with infant egocentric videos. *Advances in Neural Information Processing Systems*, 36:54199–54212, 2023. 3
- [44] Saber Sheybani, LB Smith, Z Tiganj, SS Maini, and A Dendukuri. Modelvsbaby: A developmentally motivated benchmark of out-of-distribution object recognition. *PsyArXiv*. <https://doi.org/10.31234/osf.io/83gae>, 2024. 2, 3
- [45] Linda B Smith and Lauren K Slone. A developmental approach to machine learning? *Frontiers in psychology*, 8: 296143, 2017. 1
- [46] Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for deep learning in nlp. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3645–3650, 2019. 1
- [47] Jessica Sullivan, Michelle Mei, Andrew Perfors, Erica Wojcik, and Michael C Frank. Saycam: A large, longitudinal audiovisual dataset recorded from the infant's perspective. *Open mind*, 5:20–29, 2021. 1, 2, 3, 4
- [48] Alvin Wei Ming Tan, Sunny Yu, Bria Long, Wanjing Anya Ma, Tonya Murray, Rebecca D Silverman, Jason D Yeatman, and Michael C Frank. Devbench: A multimodal developmental benchmark for language learning. *arXiv preprint arXiv:2406.10215*, 2024. 2, 3
- [49] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. 1, 3
- [50] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024. 3
- [51] Piotr Teterwak, Ximeng Sun, Bryan A Plummer, Kate Saenko, and Ser-Nam Lim. Clamp: contrastive language model prompt-tuning. *arXiv preprint arXiv:2312.01629*, 2023. 6
- [52] Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. Winoground: Probing vision and language models for visiolinguistic compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5238–5248, 2022. 1, 2, 3, 6
- [53] Hugo Touvron et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 1
- [54] Wai Keen Vong, Wentao Wang, A Emin Orhan, and Brenden M Lake. Grounded language acquisition through the eyes and ears of a single child. *Science*, 383(6682):504–511, 2024. 3, 4, 5, 6
- [55] Wentao Wang, Wai Keen Vong, Najoung Kim, and Brenden M Lake. Finding structure in one child's linguistic experience. *Cognitive science*, 47(6):e13305, 2023. 3
- [56] Alex Warstadt, Leshem Choshen, Aaron Mueller, Adina Williams, Ethan Wilcox, and Chengxu Zhuang. Call for papers—the babylm challenge: Sample-efficient pretraining

- on a developmentally plausible corpus. *arXiv preprint arXiv:2301.11796*, 2023. 3
- [57] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017. 5
- [58] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024. 1
- [59] Lorijn Zaadnoordijk, Tarek R Besold, and Rhodri Cusack. The next big thing (s) in unsupervised machine learning: Five lessons from infant learning. *arXiv preprint arXiv:2009.08497*, 2020. 4
- [60] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11975–11986, 2023. 3
- [61] Peng Zhang, Yash Goyal, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Yin and yang: Balancing and answering binary visual questions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5014–5022, 2016. 3
- [62] Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tynyllama: An open-source small language model, 2024. 5
- [63] Yuanhan Zhang, Bo Li, haotian Liu, Yong jae Lee, Liangke Gui, Di Fu, Jiashi Feng, Ziwei Liu, and Chunyuan Li. Llava-next: A strong zero-shot video understanding model, 2024. 3
- [64] Yuhui Zhang, Alyssa Unell, Xiaohan Wang, Dhruva Ghosh, Yuchang Su, Ludwig Schmidt, and Serena Yeung-Levy. Why are visually-grounded language models bad at image classification? *arXiv preprint arXiv:2405.18415*, 2024. 6
- [65] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena.” arxiv. *arXiv preprint cs.CL/2306.05685*, 2023. 3
- [66] Cheng Zhuang, Shuang Yan, Arash Nayebi, Mark Schrimpf, Michael C Frank, James J DiCarlo, and Daniel L.K. Yamins. Unsupervised neural network models of the ventral visual stream. *Proceedings of the National Academy of Sciences of the United States of America*, 118(3):e2014196118, 2021. 3