# Benchmarking Small-Scale Vision-Foundation Models through Developmentally-Inspired Cognitive Assessments
## A New Ecosystem for Academic Pretraining Research
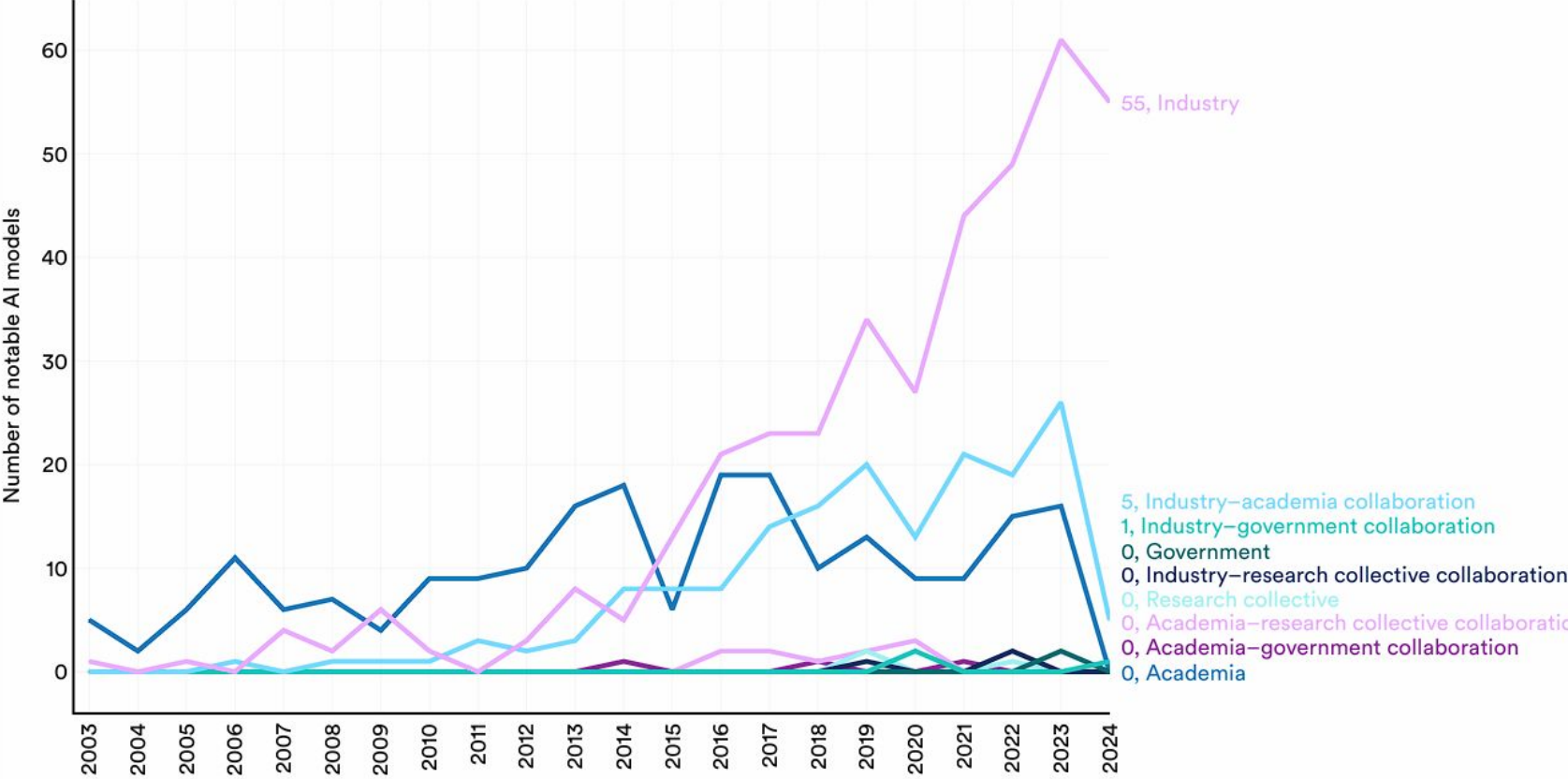
David Li[1, 2], Shawn Li[2], Jeffrey Li[2], Helen Chen[2], Andrew Zhu[2], Andrew Zagula[2], Mimi Zhao[2], Joey Huang[2], Shengao Wang[2], Wenqi Wang[2], Zecheng Wang[2], Michael Wakeham[2], Boqing Gong[2]

Mills High School, Millbrae, CA[1]; Boston University, Boston, MA[2]

## INTRODUCTION

The massive resource and data demands of modern Vision-Foundation Models (VFMs) have created a structural divide, pushing fundamental pretraining research out of academia.


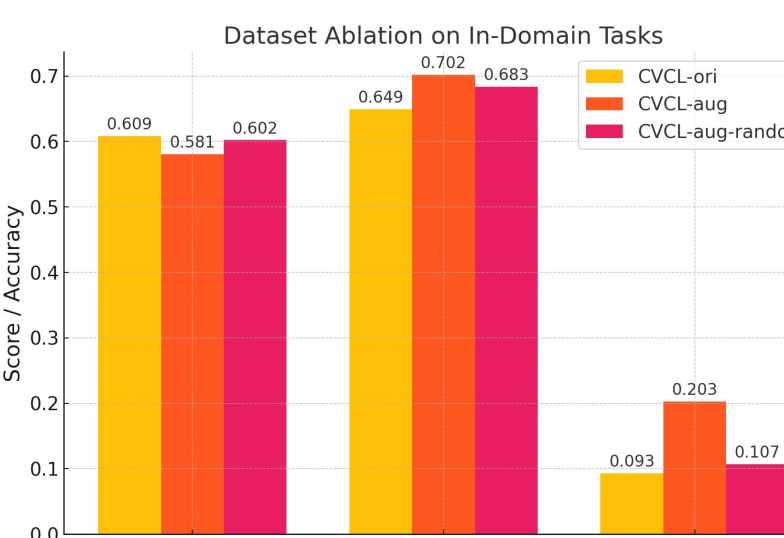**Number of notable AI models by sector, 2003–24**
Source: Epoch AI, 2025 | Chart: 2025 AI Index report

developed 55 notable ML models, academia developed 5.[1]

- Training cutting-edge models such as GPT-4 ($78M) and Gemini Ultra ($191M) vastly exceeds academic funding.[1] In 2024, global private AI investment reached $252.3 billion.[2]



- With academia is unable to match industry's scale, research has shifted to fine tuning corporate models, limiting progress.
- Rather than chasing scale, we aim to use in-domain data and benchmarks grounded in developmental psychology, lowering the barrier to entry for fundamental research.

BabyVLM shows Synthetic child-directed data (CVCL-aug) boosts VTWT and Winoground scores much more than random augmentation.[3]



## METHODS

Data Foundation: SAYCam Dataset[4]
- 472 h of first-person infant video recordings
- Computationally tractable proxy for human development
- Practical academic alternative to petabyte-scale industrial datasets

Benchmark Foundation: NIH Baby Toolbox Tasks[5]
- Visual Delayed Response: attention, memory, executive control
- Picture Vocabulary Test: receptive language (adapted from MacArthur-Bates CDI)[6]
- Memory Task: processing speed and learning efficiency
- Covers developmental age range of 1–42 months for scientific validity

Infrastructure Foundation
- Web-based GUI enabling complex SAYCam metadata queries
- Standardized model wrapper for consistent interfacing
- Automated execution on a private test set & Public leaderboard for transparent, reproducible comparisons

## INFRASTRUCTURE

To bridge large-scale naturalistic data with standardized model evaluation, we developed a comprehensive infrastructure ecosystem. These tools enable researchers to efficiently curate data and benchmark models in a reproducible, scalable manner.

**SAYCam Video & Transcript Navigator**
- Manually navigating the 472-hour SAYCam dataset is prohibitively time-consuming.
- The provided Databrary software is severely lacking. No ability to speed up videos, view/edit their transcripts, categorize them, export them.
- Our browser-based navigator transforms this raw footage into a searchable, interactive database, enabling researchers to explore developmental data without needing data management expertise.

In-depth search: transcript, tags, etc. Metadata storage:

Speed controller and Frame sampler - play or export at sampling FPS



Star and Tags (comma-separated)

Video Information & notes

Clickable, auto-scrolling transcript

Timeline Transcript Editing: draggable / resizable words

https://github.com/eleusinianexpositor/saycam_visualizer

To ensure benchmark validity, we built a scalable labeling pipeline with Label Studio, automating clip ingestion into verifiable annotation tasks. Custom interfaces for tasks like Visual Delayed Response were rapidly developed using Gemini, streamlining iteration.



We designed a **Model Evaluation Framework** that defines a standardized interface for model submission.
- Consistent, fair evaluation under identical conditions
- Streamlined integration accelerates computational research
- Automated contest foundation underpins a scalable model

**The Blueprint**
BaseVLMWrapper (Standardized Abstract Contract)
```
__init__(...)
select_choice(...)
generate_text(...)
```

**Implementation**
Researchers implement methods to wrap their model's specific logic for loading weights and running inference.

**Engine**
Evaluation Backend
```
model = VLM_Wrapper(path)
prediction = model.select_choice(instance)
score(prediction, ground_truth)
```

The final element is a **Public Contest Website**. This site will provide:
- open access to all benchmark documentation and starter code
- manage model submissions via our wrapper framework
- display results on a live, transparent leaderboard

## BENCHMARKS

**Visual Delayed Response (VDR)**
*An object appears in screen and then moves off-screen, left, right, up, or down. The VFM must identify direction.* Extracted 2200 usable clips from SAYCam, with diversity in direction and object type, yielding robust benchmark.

**Memory Task (MT)**
- *Learning: View overlapping image pairs.*
- *Testing: Given one old and one new image; pick the old.*
81 varying baby-level animal / object candidates found from SAYCam, can be combined to create adequate pairs.

**Picture Vocabulary (PV)**
*The VFM sees an image with 2–4 labeled objects and must pick the correct bounding box based on a prompt.* Generated 409 questions and 1227 total semantic, phonological, and categorical distractors.

## DISCUSSION

- Grounded in developmental science, our benchmarks level the playing field.
- Open-source tools make exploration and contribution to this new research accessible.
- Task suite will be broadened to include additional tasks, from NIH Baby Toolbox and other resources like Mullen Scales of Early Learning.
- We also plan to scale up our evaluation engine to support external model submissions and to refine our annotation tools with semi-automated methods.
- We release these resources publicly to catalyze a effort around data-efficient, developmentally-inspired pretraining.

## REFERENCES

(1) AI Index Steering Committee. The 2025 AI Index Report; Stanford Institute for Human-Centered Artificial Intelligence: Stanford, CA, 2025.
(2) Cottier, B.; Rahman, R.; Fattorini, L.; Maslej, N.; Besiroglu, T.; Owen, D. The Rising Costs of Training Frontier AI Models; 2025.
(3) Wang, S.; Chandra, A.; Liu, A.; Saligrama, V.; Gong, B. BabyVLM: Data-Efficient Pretraining of VLMs Inspired by Infant Learning; 2025.
(4) Sullivan, J.; Mei, M.; Perfors, A.; Wojcik, E.; Frank, M. C. SAYCam: A Large, Longitudinal Audiovisual Dataset Recorded From the Infant's Perspective. Open Mind 2021, 5, 20–29. https://doi.org/10.1162/opmi_a_00039.
(5) Gershon, R., Novack, M. A., & Kaat, A. J. (2024). The NIH Infant and Toddler Toolbox: A new standardized tool for assessing neurodevelopment in children ages 1–42 months. Child development, 95(6), 2252–2254.
(6) Marchman, V. A.; Dale, P. S. The MacArthur-Bates Communicative Development Inventories: Updates from the CDI Advisory Board. Frontiers in Psychology 2023, 14. https://doi.org/10.3389/fpsyg.2023.1170303.

## ACKNOWLEDGEMENTS