

Development of Cognitive Intelligence in Pre-trained Language Models

Anonymous ACL submission

Abstract

Recent studies show evidence for *emergent cognitive abilities* in Large Pre-trained Language Models (PLMs). The increasing *cognitive alignment* of these models has made them candidates for cognitive science theories. Prior research into the emergent cognitive abilities of PLMs has been *path independent* to model training, i.e. has only looked at the final model weights and not the intermediate steps. However, building plausible models of human cognition using PLMs also requires aligning their performance during training to the developmental trajectories of children’s thinking. Guided by psychometric tests of human intelligence, we choose four task categories to investigate the alignment of ten popular families of PLMs and evaluate each of their available *intermediate and final training steps*: Numerical ability, Linguistic abilities, Conceptual understanding, and Fluid reasoning. We find a striking regularity: regardless of model size, the developmental trajectories of PLMs consistently exhibit a window of maximal alignment to human cognitive development. Before that window, training appears to endow “blank slate” models with the requisite structure to be poised to rapidly learn from experience. After that window, training appears to serve the engineering goal of reducing loss but not the scientific goal of increasing alignment with human cognition.

1 Introduction

Large Pre-trained Language Models (PLMs) like Google’s Gemini (Team et al., 2023), Meta’s LLaMA 2 (Touvron et al., 2023), and OpenAI’s GPT 4 (OpenAI, 2023a) show human-level or even super-human performance on many cognitive performance tasks. This is true in domains such as mathematical reasoning (Shah et al., 2023; Ahn et al., 2024), language comprehension (Warstadt et al., 2020; Ye et al., 2023; Koubaa, 2023), concept understanding (Vemuri et al., 2024), and analogical reasoning (Webb et al., 2023; Hu et al., 2023), con-

tributing to the hype of claims of reaching Artificial General Intelligence (AGI).

Such claims deserve to be scrutinized. Human intelligence is multi-faceted. Furthermore, there is a massive disparity between the training data scale of PLMs and humans. PLMs unintentionally acquire human performance characteristics from the corpora they are trained on, through residues of the values, beliefs, and biases of the authors of the texts (Pellert et al., 2024). We approach the human alignment of PLMs by grounding evaluation in frameworks for *psychometric intelligence*. Psychometric measures of intelligence include multiple subtests spanning a range of abilities, including mathematical thinking, language comprehension, spatial thinking, fluid reasoning, and conceptual understanding (Snow et al., 1984; Carroll, 1993; Sternberg, 2000; McGrew, 2009; Haier, 2023). In this work, we choose representative assessments of different facets of human intelligence, modified for the required textual modality, to evaluate the *cognitive alignment* of PLMs.

A second goal of our work is to move beyond cognitive alignment to also evaluate the *developmental alignment* of PLMs. The claim that the final model state of a PLM approximates adult performance leaves open the question of the path by which it arrived there. Ideally, the model’s performance improvements over training also track the progression of cognitive abilities over development (Elman, 1996; Bengio et al., 2009). This potential parallelism would be stronger evidence for PLMs as cognitive science models. Researchers are increasingly addressing this question by building PLMs trained on a developmentally plausible corpus of child-directed speech, transcribed dialogue, and children’s literature (Huebner et al., 2021; Warstadt et al., 2023; Bhardwaj et al., 2024).

We ask the question of developmental alignment in a theoretically important way: Is the cognitive alignment of PLMs achieved in a *path-independent*

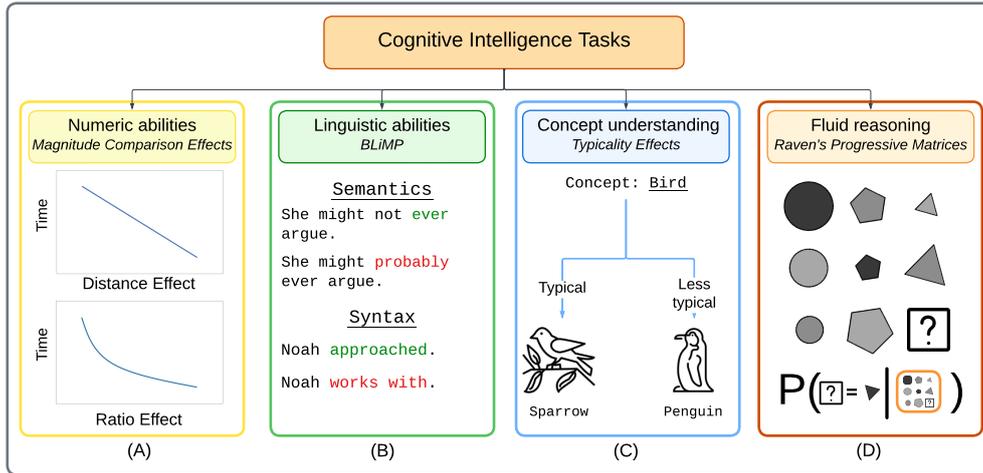


Figure 1: A list of cognitive intelligence tasks under consideration.

084 or *path-dependent* manner? Prior studies focusing
 085 on the cognitive alignment of PLMs have only es-
 086 tablished path independence: that models at the end
 087 of training approximate adult performance across
 088 a range of domains. Here, we also evaluate path
 089 dependence: Do the performance improvements of
 090 PLMs over training track the growth of these abili-
 091 ties in children over development (Holyoak et al.,
 092 1984)? We ask this question for models of different
 093 sizes and track their developmental alignment over
 094 millions and billions of training tokens. If path
 095 independence holds, this opens up applications that
 096 we detail in the conclusion.

097 To summarize, our key contributions are as follows:

- 098 • **Cognitive Modelling using AI:** We test the
 099 appropriateness of PLMs for cognitive model-
 100 ing by evaluating whether their performance
 101 profiles match those of humans.
- 102 • **Developmental trajectories in LLM pre-
 103 training and scaling:** Previous studies have
 104 only evaluated the final training checkpoints
 105 of PLMs for their cognitive plausibility, and
 106 have neglected the question of developmental
 107 trajectories. Here, we also ask: Can PLMs be
 108 used to model human developmental trajec-
 109 tories despite the training data scale mismatch
 110 between PLMs and humans?
- 111 • **Representative tasks:** We choose represen-
 112 tative tasks of human-like psychometric intel-
 113 ligence tests in PLMs. These tasks evaluate
 114 numeric, linguistic, conceptual, and fluid intel-
 115 ligence. We propose these to be a *prerequisite*
 116 to using PLMs for cognitive modeling.

2 Related work 117

2.1 Psychometric theories of intelligence 118

119 Previous intelligence assessments in AI have
 120 looked at singular dimensions, such as numeric
 121 abilities (Zhuang et al., 2023; Fang et al., 2024).
 122 Rather than choose cognitive abilities in a piece-
 123 meal fashion, we look to psychometric theories of
 124 intelligence for guidance (Sternberg, 2000). These
 125 theories distil performance on a large number of
 126 subtests into a small number of latent factors. De-
 127 spite popular attention to “general intelligence” and
 128 the latent factor g , there is a long history of theo-
 129 ries positing that intelligence is composed of multi-
 130 ple domain-specific abilities. An important, early
 131 domain-specific theory of intelligence, (Thurstone,
 132 1938), included seven “primary abilities”. The
 133 most widespread psychometric theory of intelli-
 134 gence today, the Cattell-Horn-Carroll (CHC) theory
 135 (Carroll, 1993; McGrew, 2009), includes among its
 136 “broad” abilities quantitative knowledge, reading
 137 and writing ability, fluid reasoning, and “compre-
 138 hension” knowledge (a subcomponent of which is
 139 conceptual understanding). We evaluate the cogni-
 140 tive and developmental alignment of PLMs along
 141 these four abilities.

2.2 Emergent cognitive abilities in Language Models 142

143 Recently, the performance of language models has
 144 improved as they have increased in size from mil-
 145 lions to billions of parameters, trained on larger cor-
 146 pora, and further tuned in novel ways (instruction
 147 tuned, RLHF). This has led to researchers increas-
 148 ingly advocating for the use of PLMs as cognitive
 149 models (Piantadosi, 2023; Warstadt and Bowman,
 150 2024). Increasing the number of parameters of the
 151

Table 1: Summary of assessments.

Cognitive Domain	Task	Source	License
Numeric Abilities	Magnitude Comparison Effects	(Shah et al., 2023)	(cc by 4.0)
Linguistic Abilities	BLiMP	(Warstadt et al., 2020)	(cc by 4.0)
Concept Understanding	Typicality Effects	(Vemuri et al., 2024; Castro et al., 2021)	(cc by 4.0)
Fluid reasoning	Raven’s Progressive Matrices	(Hu et al., 2023)	(cc by 4.0)

models has given rise to *Emergent Abilities* that cannot be predicted by extrapolating from the performance of smaller models (Wei et al., 2022a). Emergent abilities have been observed in a variety of task types such as multi-task language understanding (Hendrycks et al., 2021), grounded conceptual mapping (Patel and Pavlick, 2022), and truthfulness (Lin et al., 2021). In recent works, Hoffmann et al. (2022); Biderman et al. (2023) show the benefits of training a model for more tokens on problem-solving (Wei et al., 2022b), common-sense reasoning (Sakaguchi et al., 2021), arithmetic abilities (Biderman et al., 2023), and linguistic performance (Paperno et al., 2016). Although the presence of emergent abilities extends to cognitive science domains (Wei et al., 2022b; Goertzel, 2023; Hagen-dorff, 2023), prior studies have been piecemeal in their approach and have failed to (1) consider multiple cognitive abilities as specified by theories of psychometric intelligence and (2) move beyond cognitive alignment to also evaluate the developmental alignment of PLMs over training.

2.3 Pre-trained language model use in developmental modeling

Recently, researchers have begun advocating for the use of PLMs for modeling cognitive development in children (Kosoy et al., 2023; Salewski et al., 2024). For example, Portelance et al. (2023) and Bhardwaj et al. (2024) suggest the use of language models to predict the age of acquisition of words in children. Researchers have also proposed studying second language acquisition and bilingualism by mapping pre-training steps in PLMs to understand the rate of language development (Evanson et al., 2023; Marian, 2023; Sharma et al., 2024). We investigate the assumption that the performance of an intermediate training checkpoint of PLMs maps to the age of child development by looking at the acquisition of human-like psychometric intelligence.

3 A suite of psychometric intelligence tasks

We assemble a suite of cognitively plausible assessments that benchmark PLMs across four abilities of

psychometric intelligence. Table 1 summarizes the tasks along with the licensing details for public use. The details of each assessment and their respective operationalization are given below.¹

3.1 Numeric abilities



Figure 2: Mental Number Line: Organization of magnitude representations in a logarithmically scaled manner.

The question of how humans understand symbolic numbers has been investigated by cognitive scientists for more than half a century. These studies show that people map number symbols to a *mental number line* (MNL, Figure 2) with a log-compressed psychophysical scale (Moyer and Landauer, 1967a).

Prior research on the numerical abilities of PLMs has focused on improving performance on application-driven tasks that require numerical skills in the context of arithmetic equations and word problems (Burns et al., 2021; Amini et al., 2019; Yuan et al., 2023), exact facts (Lin et al., 2020), and measurement estimation (Zhang et al., 2020). However, these tasks fail to directly implicate the key cognitive construct underlying human numerical understanding, the recruitment of a compressed MNL.

In a recent study, Shah et al. (2023) found evidence for a human-like MNL in various PLMs. They show that despite lacking explicit neural circuitry to represent numbers, through experience (i.e., vast amounts of training data), PLMs show human-like performance profiles and learn human-like representations for numerical concepts.

We follow Shah et al. (2023) and look for the two behavioral signatures of a compressed number line representation, the distance effect and the ratio

¹We will add all tasks to a publically available unified language model testing framework, titled *lm-evaluation-harness* (Gao et al., 2023), to support the evaluation of future models on psychometric intelligence assessments.

effect. In humans, these are defined as:

- **Distance effect** (refer to Figure 1A - top): The greater the distance $|x - y|$ between two numbers x and y , the faster they are compared, i.e., the greater (or lesser) number is identified (Moyer and Landauer, 1967b).
- **Ratio effect** (refer to Figure 1A - bottom): The time to compare two numbers x and y is a decreasing function of the ratio of the larger number over the smaller number $\frac{\max(x,y)}{\min(x,y)}$ (Halberda et al., 2008).

These effects can be mapped to language models by adopting the following linking hypothesis: *the greater the cosine similarity of two number representations in a PLM, the longer it takes to discriminate them, i.e., to judge which one is greater (or lesser)*. While we focus on the Distance and Ratio effect, the results for all the effects investigated by Shah et al. (2023) are in Appendix B.1.

Operationalization: We used the same protocol as Shah et al. (2023). For each effect, we test the three formats of number representations of PLMs (mixed-case number words, lower-case number words, and digits). We present the R^2 values for the Distance and Ratio effects, which are averaged across each input representation. The R^2 values for the distance effect in PLMs are obtained by fitting a linear function predicting the cosine similarity of x and y from their distance $|x - y|$. R^2 values for the ratio effect in PLMs are obtained by fitting a negative exponential function predicting the normalized cosine similarity of x and y from their ratio $\frac{\max(x,y)}{\min(x,y)}$. Note: This task requires access to the latent representations of models.

3.2 Linguistic abilities

Language (or verbal) ability is a central component of human cognition and cognitive neuroscience (Hagoort, 2019). At the dawn of the cognitive revolution, it was conceptualized as a largely innate ability, and language acquisition was understood as requiring relatively little learning from experience (Fodor, 1985; Chomsky, 2014). More recently, cognitive developmentalists have shown that infants can learn language through exposure to the statistical regularities of the linguistic environment (Safra et al., 1996; Siegelman, 2020). These findings have been modeled using multi-layer perceptrons

(Elman, 1996) and, more recently, PLMs (Lake and Murphy, 2023).

We use BLiMP: The Benchmark of Linguistic Minimal Pairs for English (Warstadt et al., 2020) to evaluate the linguistic abilities of each PLM under consideration. BLiMP consists of 67 datasets of 1000 pairs of minimally different sentences which vary in acceptability and span 12 phenomena at three levels of language: *morphology*, *syntax*, and *semantics*. The 12 phenomena are described in Appendix B.2. Each pair consists of one acceptable sentence and one unacceptable sentence which otherwise differ minimally. BLiMP evaluates the models by measuring if they assign a higher probability to the acceptable vs. unacceptable sentence of each pair. Figure 1B shows two examples of minimal pairs.

Operationalization: We use the LM-eval-harness (Gao et al., 2023) benchmarking suite to test our models on the BLiMP tasks. We evaluate if a model assigns a higher sequential probability to the acceptable sentence. Note: This requires models that can generate probabilities of tokens.

3.3 Concept understanding

On encountering a new stimulus, humans categorize it – assign it to a known concept – in order to make inferences about its unobservable properties (Murphy, 2002). A striking finding is that not all members of a category are equal (Rosch, 1975). Rather, some members (e.g., pigeon) are more typical of a category (e.g., Bird) than other members (e.g., ostrich). This phenomenon, known as the *Typicality Effect*, is a central feature of human categorization (Lakoff, 2008).

Typicality gradients in humans can be measured using the production task, where participants are given a category label (e.g., Bird) and asked to list as many members of the category as they can in a limited time (Battig and Montague, 1969; Van Overschelde et al., 2004; Castro et al., 2021). The typicality of an item is defined as the proportion of participants who produce it.

Language models have shown some evidence of human-like typicality gradients. Heyman and Heyman (2019) used word2vec embeddings to predict the category typicality norms released by De Deyne et al. (2008). More recent work by Misra et al. (2021) and Bhatia and Richie (2022) has looked at correlations of PLMs like BERT, RoBERTa, and GPT-2 to the Rosch (1975) typicality norms for

ten categories. Vemuri et al. (2024) performed the most comprehensive study of the alignment of concept understanding in the latent representations of PLMs. We expand upon their task setup to evaluate human-like concept understanding in the PLMs that are the focus here.

Operationalization: For each model, we calculate the representativeness of a member to its category in three possible ways:

- Closeness judgment problem: Calculate the cosine similarity between the obtained latent representations for the member and the category. This requires models where the latent representations are readily available.
- Surprisal values: For each member in a category, the probability of the sequence a "member" (eg. pigeon) is a "category" (eg. bird). This method requires access to the probability of each token in a sequence.
- Prompting: Prompt the models with the following design: Guidelines, Query, and Options. The Guideline highlights the task of re-ranking the members given in the Options based on appropriateness with the Query. The Query consists of the in-filling task: A ___ is a [category name]. The Options are each of the possible members of the category. Given the complexity of the prompting, usable outputs are only obtained from models that are larger than 30 billion parameters.

For the two in-filling problems (i.e., based on surprisal values and prompting), we also evaluate models on zero to three exemplars as context. The details of the experiments on these different exemplar contexts are given in Appendix B.3.

3.4 Fluid reasoning

Humans can logically parse information and detect patterns in novel stimuli without having to rely on prior experiences or learned information. This ability is called Fluid Reasoning (Cattell, 1963).

We focus on the dominant measure of fluid reasoning, the Ravens Progressive Matrices (RPM) test (Raven, 2003). An example Ravens-like problem is given in Figures 1D and 3. An RPM item consists of a 3x3 matrix of cells with one empty cell. Participants must induce the underlying, abstract patterns that hold across the rows and columns of the matrix, and apply these to infer the image in the empty cell from a given set of options. These images vary in visual attributes like shape and color,

along with more abstract qualities. The RPM is the standard measure of fluid reasoning (Snow et al., 1984) and is highly correlated with analogical reasoning (Goswami, 1986; Webb et al., 2023).

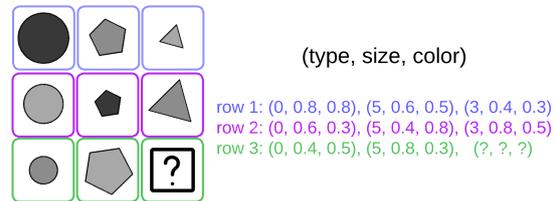


Figure 3: Example adaptation of visual RPM problems to the textual format. Each image is decomposed into tuples of (type, size, color). Type indicates the shape of the image.

Given the visual nature of the RPM, previous work by Hu et al. (2021, 2023) and Webb et al. (2023) mapped the Raven-10000 dataset to a textual format to facilitate the testing of PLMs. The mapping involves reformulating visual elements into text-based numerical tuples to form the I-Raven dataset, representing attributes like shape, size, and color textually, as illustrated in Figure 3. We use their approach with a focus on the “Center Single Alignment” sub-task, which features a single shape per matrix cell. We differ from their work by evaluating a broader set of models.

Operationalization: We determine the model’s preferred answer for a problem by comparing the surprisal values of the whole sequence (instruction, question, candidate tuple) for each of the candidate options, i.e. the probability of each completed digit representation of a matrix. For the example given in Figure 3, this would be checking the probability of this sequence (summation of token probabilities) with the correct answer (3, 0.6, 0.8) to the other candidates. A comprehensive list of the prompts used in this paper is given in Appendix B.4.

4 Models under consideration

We evaluate a wide range of language model families, shown in Table 2. These models are selected based on the following criteria:

Public availability: Open-source models allow us to perform a thorough analysis by accessing the latent representation and the token probability during generation. We follow Holt et al. (2024) while choosing PLMs. Although most models in this study are publicly available and open-source, we use three state-of-art commercial PLMs that are

Table 2: List of language model families under consideration with their statistics.

Models	Source	Latent rep.	Token prob.	Multiple sizes	Intermediate checkpoints	Known training order
Amber	(Liu et al., 2023)	✓	✓	✗	✓	✓
Falcon	(Almazrouei et al., 2023)	✓	✓	✗	✗	✗
Starling	(Zhu et al., 2023)	✓	✓	✗	✗	✗
Llama-2	(Touvron et al., 2023)	✓	✓	✓	✗	✗
Mistral	(Jiang et al., 2023)	✓	✓	✗	✗	✗
Qwen	(Bai et al., 2023)	✓	✓	✓	✗	✗
Pythia	(Biderman et al., 2023)	✓	✓	✓	✓	✓
Gemini	(Team et al., 2023)	✗	✗	✗	✗	✗
GPT-3.5-Turbo	(OpenAI, 2023b)	✗	✓	✗	✗	✗
GPT 4	(OpenAI, 2023a)	✗	✓	✗	✗	✗

Table 3: Performance of Pre-trained Language Models on the tasks. Distance Effect: Averaged R^2 values of different LLMs when fitting a linear function on the cosine-similarity vs. distance plot. Ratio Effect: Averaged R^2 values of different LLMs when fitting a negative exponential function on the cosine-similarity vs. ratio plot. Note: Each value is averaged across all three input types and all model layers to produce one generalizable score. Latent Rep: Average Spearman’s Correlation when using the cosine similarity and latent representation-based approach (Note: * refers to the [prompting approaches](#) for select models which are gated by APIs and not the latent representation-based approach), Zero-Shot: Average Spearman’s Correlation when using the zero-shot surprisal values, BLiMP: The Benchmark of Linguistic Minimal Pairs for English, RPM: Raven’s Progressive Matrices

Model	Numeric Abilities		Linguistic Abilities	Conceptual Understanding		Fluid reasoning
	Distance	Ratio	BLiMP	Latent Rep.	Zero Shot	RPM
	Effect (R^2)	Effect (R^2)	(Acc.)	(Average Spearman’s Correlation)		(Acc.)
Amber-7B	0.913	0.591	0.794	0.083	0.250	0.654
Falcon-7B	0.928	0.838	0.817	-0.116	0.180	0.730
Starling-LM-7B-alpha	0.522	0.187	0.827	-0.003	0.258	0.730
Llama-2-7B	0.670	0.614	0.818	-0.065	0.238	0.752
Llama-2-13B	0.672	0.263	0.793	0.076	0.247	0.756
Mistral-7B	0.641	0.233	0.829	-0.025	0.245	0.756
Mistral-7B-Instruct	0.637	0.543	0.834	0.033	0.255	0.674
Qwen-0.5B	0.833	0.553	0.785	0.072	0.282	0.684
Qwen-1.8B	0.878	0.301	0.792	0.114	0.235	0.746
Qwen-4B	0.881	0.264	0.730	0.001	0.246	0.770
Qwen-7B	0.858	0.616	0.789	0.006	0.229	0.766
Qwen-14B	0.783	0.507	0.792	-0.140	0.249	0.776
Pythia-70M	0.829	0.429	0.723	0.005	0.211	0.194
Pythia-160M	0.947	0.665	0.749	0.067	0.260	0.448
Pythia-410M	0.926	0.679	0.815	0.126	0.284	0.608
Pythia-1B	0.944	0.702	0.806	0.090	0.280	0.674
Pythia-1.4B	0.933	0.764	0.819	0.074	0.283	0.730
Pythia-2.8B	0.961	0.723	0.827	0.221	0.273	0.760
Pythia-6.9B	0.909	0.713	0.809	0.105	0.280	0.716
Pythia-12B	0.846	0.595	0.829	0.184	0.291	0.756
Gemini	NA	NA	NA	0.311 *	NA	NA
GPT-3.5-Turbo	NA	NA	0.825	0.242 *	0.231	0.792
GPT-4	NA	NA	0.849	0.559 *	0.428	0.822

gated behind API calls; GPT-3.5-Turbo (pointing to gpt-3.5-turbo-0613 on the OpenAI platform), GPT-4 (pointing to gpt-4-1106 on the OpenAI platform), and Gemini (also referred to as Gemini-1-Pro at the time of writing). The GPT- x model APIs provide token probabilities of the response, allowing us to calculate surprisal, while Gemini does not.

Availability of multiple sizes: The availability of model sizes for the same architecture and training paradigms allows us to evaluate the emergent cognitive abilities of the models. We have multiple sizes available for the Llama-2, Qwen, and the Pythia family of models.

Availability of intermediate training checkpoints: This allows us to evaluate the effects of pre-training on the model outputs. Together, the availability of multiple model sizes and intermediate training checkpoints allow us to best evaluate the developmental alignment of PLMs. Amber and Pythia’s family of models have available intermediate training checkpoints. While Amber has 360 intermediate checkpoints, the checkpoints are at 4 Billion tokens each and are not at the required granularity.

Pythia Family of models: Pythia (Biderman et al., 2023) is one of the first open-source projects with the goal of scientific and transparent model

415
416
417
418
419
420
421
422
423
424
425
426
427

428
429
430
431
432
433
434
435
436
437
438
439
440

development. It has 8 model sizes ranging from 70 Million to 12 Billion parameters, with each model trained on 286 Billion tokens. The models in the suite are equivalent (in size) to popular decoder architectures like GPT-Neo-(125M, 1.3B, 2.7B) and OPT-(125M, 350M, 1.3B, 2.7B, 6.7B), but with the added benefits of training on a known de-duplicated corpus (Gao et al., 2020), using the same training order for each model size, and having 154 intermediate checkpoints to study the learning trajectories of PLMs. Thus, the Pythia suite of models is ideal for studying the psychometric and developmental alignment of PLMs to humans.

All open-source models are obtained from Huggingface (Wolf et al., 2020), while the gated models are obtained from their respective platforms through API calls. For each model in the Pythia suite, the following intermediate checkpoints are available: [1, 2, 4, 8, ... 512; 1000, 2000, 3000 ... 143000 (exponential increase in checkpoint number until the 512th checkpoint and subsequent progression of 1000 steps until the last checkpoint)], with each checkpoint representing 2 Million tokens seen. Overall, we test 1232 intermediate checkpoints of the Pythia suite of models across all the tasks.

5 Cognitive and developmental alignment of PLMs

The suite of tasks enables comprehensive evaluation of a variety of PLMs on their cognitive alignment to humans across four domains of psychometric intelligence: numeric abilities, linguistic abilities, concept understanding, and fluid reasoning. Table 3 highlights the key results of this evaluation. For the evaluation of conceptual understanding in PLMs, we only report the results for the zero-shot surprisal values and latent representations. This is because we see similar results for zero-shot and few-shot surprisal value-based methods (see comprehensive results in Appendix B.3).

The cognitive alignment of PLMs on psychometrics assessments is summarized below:

- *Numeric abilities:* All PLMs show a human-like distance effect but weakly show a human-like ratio effect. We do not observe any notable changes in alignment with model scaling, indicating the need for the evaluation of future models on this task.
- *Linguistic abilities:* The accuracy of the PLMs on the BLiMP linguistic acceptability tasks

improves upon increasing the number of parameters. Furthermore, we find that all PLMs are substantially more accurate on morphological tasks over syntactic and semantic tasks (*Accuracy: semantic < syntax << morphology*; see Appendix Table 5, Figure 7). Morphological performance develops first followed by syntax and then semantics.

- *Concept understanding:* Prompting methods in commercial models perform substantially better than other methods – closeness judgment and surprisal values – on all open-source models. In the Pythia suite, we observe that larger models outperform smaller counterparts on the same training data.
- *Fluid reasoning:* For all PLM architecture types, larger models outperform their smaller equivalent models.
- Despite differences in PLM architecture type, all models of an approximate size of 7 Billion parameters perform comparably.

The developmental alignment of the PLMs on the tasks is shown in Figure 4. We make the following key observations:

- *Training endows the “blank slate” with requisite structure:* In each assessment, the model “warm-ups” in training on a few million/ billion tokens, moving from a “blank slate” to possessing the requisite structure. This structure can be thought of as the child’s endowment at birth. Development of the four abilities begins only after reaching this state.
- *Training shows a region of development:* For all four tasks, we see a window of monotonic development, in which all models gain the respective cognitive abilities.
- *After development, training appears to serve an engineering goal:* After the window of development, training appears to only serve the engineering goal of loss reduction. This observation is especially pronounced for numeric abilities and conceptual understanding.
- *Assessments for Fluid Reasoning and Linguistic Abilities show significant gains in scaling and greater pre-training:* For the Fluid Reasoning and Linguistic Abilities assessments, we see that the alignment score continues to increase as the PLMs are trained on a greater number of tokens. Furthermore, for these abilities, models also show scaling effects, with larger models outperforming smaller ones.

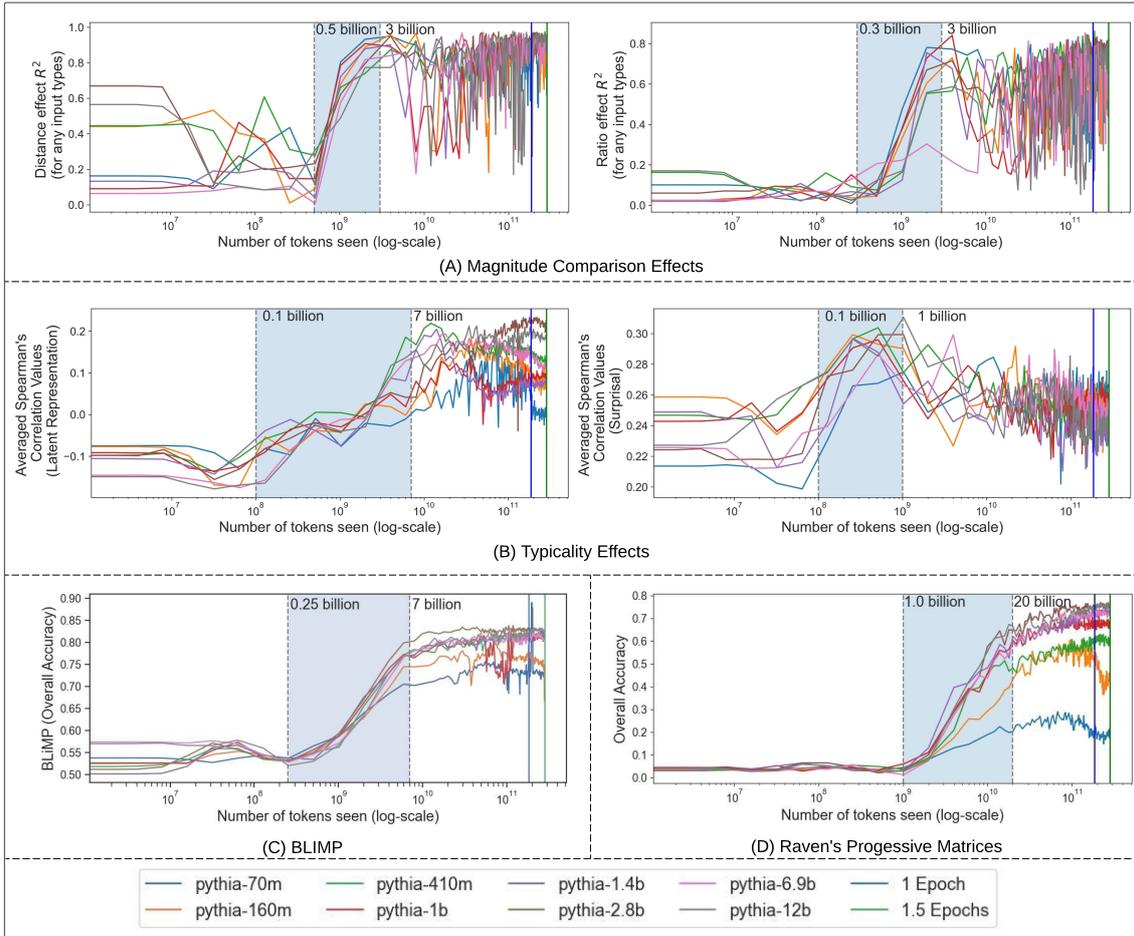


Figure 4: Developmental trajectory of the Pythia suite of models on the psychometric intelligence tasks as a function number of tokens seen. We display the x-axis in a log-scaled manner as maximal development occurs in the range of 100 Million to 20 Billion tokens seen for all tasks. The windows of maximal development are illustrated by the blue shading.

- *The windows weakly align with human ages of development:* Variation in the onsets of windows replicates what is known of cognitive development. For example, children acquire language early (i.e., during the preschool years), whereas the onset of improving fluid reasoning is later, when children enter elementary school, and continues for longer, throughout adolescence. Correspondingly, the models significantly develop linguistic abilities while training on 250 Million to 7 Billion tokens, whereas they acquire fluid reasoning abilities later, while training on 1 to 20 Billion tokens.

6 Conclusions

This paper investigates the evidence appropriateness of using PLMs for human cognitive and developmental modeling with the help of adapted psychometric intelligence assessments. It uses repre-

sentative assessments of four facets of human intelligence: numeric abilities, linguistic abilities, conceptual understanding, and fluid reasoning. Our experiments show that PLMs develop cognitive abilities purely through their experience in the world, indicating that cognitive abilities in humans may not be innate, but rather learned similarly through the world. Most significantly, we find a window of monotonic development in which all models improve approximately linearly on the four cognitive abilities. Before that window, we interpret training as endowing “blank slate” models with the requisite structure for rapid learning. Also notable is the finding of PLM scaling effects for the assessments of linguistic abilities and fluid reasoning. We propose evaluation against these tasks as a *prerequisite* before treating PLMs as models of human cognition and its development.

7 Limitations

Some limitations of the work are as follows: (1) We use an aggregation of psychometric tests for PLMs. The limitations of each test are inherited in the suite of tasks. (2) The alignment scores may be wrongly interpreted when evaluating PLMs with these tasks. Alignment scores show the similarity of PLM outputs to human outputs on psychometric tests and indicate that PLMs do not need explicit neural circuitry for these intelligence tests. We do not suggest these models as proxies for humans in any manner and recommend further testing before use. (3) The developmental alignment of the models points towards the acquisition of human-like performance on the four psychometric assessments in the range of 100 Million to 20 Billion training tokens. This conclusion has two limitations: Pythia is the only suite of models with available intermediate checkpoints and, while unlikely, the observed developmental trajectories might be artifacts of the pre-training order. (4) The psychometric assessments for PLMs are adapted from similar human psychometric tests. Different ways of adaptation may lead to different results. Furthermore, while representative, these assessments are not exhaustive tests of human intelligence. Future work can expand to other tests like spatial and commonsense reasoning. (5) Some open source models like Llama-2 have larger 70 Billion parameter variants but we lack the compute resources to evaluate them. Large open-source models would lead to appropriate comparisons of performance with commercial models like GPT-4. (6) While our work evaluates changes in cognitive alignment with an increase in model size and the number of pre-training tokens, we do not control for different tuning methodologies like instruction tuning and reinforcement learning with human or artificial intelligence feedback. Accounting for different tuning methods is computationally intensive for the 1200+ model checkpoints across 10 architectures.

8 Ethical Considerations

All tasks and corresponding datasets have low ethical risks and none expose sensitive information. Additionally, we obtain approval from the authors of each dataset for their use and release. There are no major risks associated with conducting this research beyond those associated with working with PLMs. There may be risks in misinterpreting the alignment scores when evaluating with the tests.

The psychometric analysis of this study is one-way: we look for human performance characteristics and behaviors in PLMs. PLMs are experimental technologies and future work using this research should proceed with caution. Assessment of the tasks indicates PLM alignment – or the lack thereof – to human cognitive behavior. Indications of higher human alignment do not indicate an absolute proxy for humans. The goal of tasks in this work is a pre-cursor assessment of PLMs on their ability to act as cognitive models. Therefore, researchers and users should perform more tests before use. .

References

- Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. 2024. Large language models for mathematical reasoning: Progresses and challenges. *arXiv preprint arXiv:2402.00157*.
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Hestlow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. Falcon-40B: an open large language model with state-of-the-art performance.
- Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. *MathQA: Towards interpretable math word problem solving with operation-based formalisms*. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2357–2367, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Sheng-guang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- W. F. Battig and W. E. Montague. 1969. Category norms of verbal items in 56 categories a replication and extension of the connecticut category norms. *Journal of Experimental Psychology Monographs*, 80:1–46.
- Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In

681	<i>Proceedings of the 26th annual international conference on machine learning</i> , pages 41–48.	Cody S. Ding. 2018. <i>Fundamentals of Applied Multidimensional Scaling for Educational and Psychological Research</i> . Springer International Publishing.	733
682			734
683	Khushi Bhardwaj, Raj Sanjay Shah, and Sashank Varma.		735
684	2024. Pre-training llms using human-like development data corpus .	Jeffrey L Elman. 1996. <i>Rethinking innateness: A connectionist perspective on development</i> , volume 10. MIT press.	736
685			737
686	S. Bhatia and R. Richie. 2022. Transformer networks of human conceptual knowledge. <i>Psychological Review</i> .		738
687		Linnea Evanson, Yair Lakretz, and Jean-Rémi King. 2023. Language acquisition: do children and language models follow similar learning stages? <i>arXiv preprint arXiv:2306.03586</i> .	739
688			740
689	Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. 2023. Pythia: A suite for analyzing large language models across training and scaling .		741
690			742
691		Qixiang Fang, Daniel L Oberski, and Dong Nguyen. 2024. Patch–psychometrics-assisted benchmarking of large language models: A case study of mathematics proficiency. <i>arXiv preprint arXiv:2404.01799</i> .	743
692			744
693			745
694			746
695			
696	I. Borg and P.J.F. Groenen. 2005. <i>Modern Multidimensional Scaling: Theory and Applications</i> . Springer.	Gustav Theodor Fechner. 1860. <i>Elements of psychophysics</i> . 1.	747
697			748
698	Gregory C. Burgess, Jeremy R. Gray, Andrew R. A. Conway, and Todd Samuel Braver. 2011. Journal of experimental psychology : General neural mechanisms of interference control underlie the relationship between fluid intelligence and working memory span .	Jerry A Fodor. 1985. <i>Precis of the modularity of mind. Behavioral and brain sciences</i> , 8(1):1–5.	749
699			750
700		Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. The pile: An 800gb dataset of diverse text for language modeling .	751
701			752
702			753
703	Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the MATH dataset . <i>CoRR</i> , abs/2103.03874.	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2023. A framework for few-shot language model evaluation .	754
704			755
705			756
706			757
707	John B Carroll. 1993. <i>Human cognitive abilities: A survey of factor-analytic studies</i> . 1. Cambridge University Press.		758
708			759
709			760
710	Nichol Castro, Taylor Curley, and Christopher Hertzog. 2021. Category norms with a cross-sectional sample of adults in the united states: Consideration of cohort, age, and historical effects on semantic categories . <i>Behavior research methods</i> , 53(2):898–917.		761
711			762
712			763
713			764
714		Ben Goertzel. 2023. Generative ai vs. agi: The cognitive strengths and weaknesses of modern llms . <i>arXiv preprint arXiv:2309.10371</i> .	765
715	Raymond B Cattell. 1963. Theory of fluid and crystallized intelligence: A critical experiment. <i>Journal of educational psychology</i> , 54(1):1.		766
716			767
717		Usha Goswami. 1986. Children's use of analogy in learning to read: A developmental study . <i>Journal of Experimental Child Psychology</i> , 42(1):73–83.	768
718	Raymond Bernard Cattell. 1987. <i>Intelligence: Its structure, growth and action</i> . Elsevier.		769
719			770
720	Noam Chomsky. 2014. <i>Aspects of the Theory of Syntax</i> . 11. MIT press.	Thilo Hagendorff. 2023. Machine psychology: Investigating emergent capabilities and behavior in large language models using psychological methods. <i>arXiv preprint arXiv:2303.13988</i> .	771
721			772
722	Andrew R. A. Conway, Nelson Cowan, Michael F. Bunting, David J. Theriault, and Scott R. B. Minkoff. 2002. A latent variable analysis of working memory capacity, short-term memory capacity, processing speed, and general fluid intelligence . <i>Intelligence</i> , 30:163–183.		773
723			774
724		Peter Hagoort. 2019. <i>Human language: From genes and brains to behavior</i> . MIT Press.	775
725			776
726		Richard J Haier. 2023. <i>The neuroscience of intelligence</i> . Cambridge University Press.	777
727			778
728	S. De Deyne, S. Verheyen, E. Ameel, W. Vanpaemel, M. J. Dry, W. Voorspoels, and G. Storms. 2008. Exemplar by feature applicability matrices and other dutch normative data for semantic concepts. <i>Behavior research methods</i> , 40:1030–1048.	Justin Halberda, Michèle M. M. Mazzocco, and Lisa Feigenson. 2008. Individual differences in non-verbal number acuity correlate with maths achievement . <i>Nature</i> , 455(7213):665–668.	779
729			780
730			781
731			782
732			

783	Joshua K. Hartshorne and Laura T. Germine. 2015.	and weaknesses of lamda responses. <i>arXiv preprint</i>	838
784	When does cognitive functioning peak? the asyn-	<i>arXiv:2305.11243</i> .	839
785	chronous rise and fall of different cognitive abilities		
786	across the life span. <i>Psychological Science</i> , 26:433 –	Anis Koubaa. 2023. GPT-4 vs. GPT-3.5: A concise	840
787	443.	showdown.	841
788	Iris Hendrickx, Su Nam Kim, Zornitsa Kozareva,	Brenden M Lake and Gregory L Murphy. 2023. Word	842
789	Preslav Nakov, Diarmuid O Séaghdha, Sebastian	meaning in minds and machines. <i>Psychological re-</i>	843
790	Padó, Marco Pennacchiotti, Lorenza Romano, and	<i>view</i> , 130(2):401.	844
791	Stan Szpakowicz. 2019. Semeval-2010 task 8 - multi-	George Lakoff. 2008. <i>Women, fire, and dangerous</i>	845
792	way classification of semantic relations between pairs	<i>things: What categories reveal about the mind.</i> Uni-	846
793	of nominals. <i>arXiv preprint arXiv:1911.10422</i> .	versity of Chicago press.	847
794	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,	Bill Yuchen Lin, Seyeon Lee, Rahul Khanna, and Xiang	848
795	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	Ren. 2020. Birds have four legs?! NumerSense:	849
796	2021. Measuring massive multitask language under-	Probing Numerical Commonsense Knowledge of Pre-	850
797	standing.	Trained Language Models. In <i>Proceedings of the</i>	851
798	Tom Heyman and Gert Heyman. 2019. Can prediction-	<i>2020 Conference on Empirical Methods in Natural</i>	852
799	based distributional semantic models predict typical-	<i>Language Processing (EMNLP)</i> , pages 6862–6868,	853
800	ity? <i>Quarterly Journal of Experimental Psychology</i> ,	Online. Association for Computational Linguistics.	854
801	72:2084–2109.	Stephanie Lin, Jacob Hilton, and Owain Evans. 2021.	855
802	Jordan Hoffmann, Sebastian Borgeaud, Arthur Men-	Truthfulqa: Measuring how models mimic human	856
803	sch, Elena Buchatskaya, Trevor Cai, Eliza Ruther-	falsehoods. <i>arXiv preprint arXiv:2109.07958</i> .	857
804	ford, Diego de Las Casas, Lisa Anne Hendricks,	Zhengzhong Liu, Aurick Qiao, Willie Neiswanger,	858
805	Johannes Welbl, Aidan Clark, et al. 2022. Train-	Hongyi Wang, Bowen Tan, Tianhua Tao, Junbo Li,	859
806	ing compute-optimal large language models. <i>arXiv</i>	Yuqi Wang, Suqi Sun, Omkar Pangarkar, Richard Fan,	860
807	<i>preprint arXiv:2203.15556</i> .	Yi Gu, Victor Miller, Yonghao Zhuang, Guowei He,	861
808	Faye Holt, William Held, and Diyi Yang. 2024. Per-	Haonan Li, Fajri Koto, Liping Tang, Nikhil Ranjan,	862
809	ceptions of language technology failures from south	Zhiqiang Shen, Xuguang Ren, Roberto Iriondo, Cun	863
810	asian english speakers.	Mu, Zhiting Hu, Mark Schulze, Preslav Nakov, Tim	864
811	Keith J Holyoak, Ellen N Junn, and Dorrit O Billman.	Baldwin, and Eric P. Xing. 2023. Llm360: Towards	865
812	1984. Development of analogical problem-solving	fully transparent open-source llms.	866
813	skill. <i>Child development</i> , pages 2042–2055.	Viorica Marian. 2023. Studying second language ac-	867
814	Sheng Hu, Yuqing Ma, Xianglong Liu, Yanlu Wei, and	quisition in the age of large language models: Unlock-	868
815	Shihao Bai. 2021. Stratified rule-aware network for	ing the mysteries of language and learning, a commen-	869
816	abstract visual reasoning. In <i>Proceedings of the</i>	tary on “age effects in second language acquisition:	870
817	<i>AAAI Conference on Artificial Intelligence (AAAI)</i> ,	Expanding the emergentist account” by catherine l.	871
818	volume 35, pages 1567–1574.	caldwell-harris and brian macwhinney. <i>Brain and</i>	872
819	Xiaoyang Hu, Shane Storcks, Richard L Lewis, and	<i>language</i> , 246.	873
820	Joyce Chai. 2023. In-context analogical reasoning	Kevin S McGrew. 2009. Chc theory and the human	874
821	with pre-trained language models. <i>arXiv preprint</i>	cognitive abilities project: Standing on the shoulders	875
822	<i>arXiv:2305.17626</i> .	of the giants of psychometric intelligence research.	876
823	Philip A. Huebner, Elior Sulem, Fisher Cynthia, and	K. Misra, A. Ettinger, and J. T. Rayz. 2021. Do lan-	877
824	Dan Roth. 2021. BabyBERTa: Learning more gram-	guage models learn typicality judgments from text?	878
825	mar with small-scale child-directed language. In <i>Pro-</i>	<i>arXiv preprint arXiv:2105.02987</i> .	879
826	<i>ceedings of the 25th Conference on Computational</i>	Robert S. Moyer and Thomas K. Landauer. 1967a. Time	880
827	<i>Natural Language Learning</i> , pages 624–646, Online.	required for judgements of numerical inequality. <i>Nat-</i>	881
828	Association for Computational Linguistics.	<i>ure</i> , 215(5109):1519–1520.	882
829	Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-	Robert S. Moyer and Thomas K. Landauer. 1967b. Time	883
830	sch, Chris Bamford, Devendra Singh Chaplot, Diego	required for judgements of numerical inequality. <i>Nat-</i>	884
831	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	<i>ure</i> , 215(5109):1519–1520.	885
832	laume Lample, Lucile Saulnier, et al. 2023. Mistral	G. Murphy. 2002. <i>The Big Book of Concepts</i> . MIT	886
833	7b. <i>arXiv preprint arXiv:2310.06825</i> .	press.	887
834	Eliza Kosoy, Emily Rose Reagan, Leslie Lai, Alison	OpenAI. 2023a. Gpt-4 technical report.	888
835	Gopnik, and Danielle Krettek Cobb. 2023. Com-	OpenAI. 2023b. New and improved embedding model.	889
836	paring machines and children: Using developmen-	Accessed: 2023-08-14.	890
837	tal psychology experiments to assess the strengths		

891	Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. 2016. The lambda dataset: Word prediction requiring a broad discourse context. <i>arXiv preprint arXiv:1606.06031</i> .	944
892		945
893		946
894		
895		
896		
897	John M. Parkman. 1971. Temporal aspects of digit and letter inequality judgments . <i>Journal of Experimental Psychology</i> , 91(2):191–205.	
898		
899		
900	Roma Patel and Ellie Pavlick. 2022. Mapping language models to grounded conceptual spaces . In <i>International Conference on Learning Representations</i> .	
901		
902		
903	Inc. Pearson. 2021. Miller’s analogy test preparation .	
904	Max Pellert, Clemens M. Lechner, Claudia Wagner, Beatrice Rammstedt, and Markus Strohmaier. 2024. Ai psychometrics: Assessing the psychological profiles of large language models through psychometric inventories . <i>Perspectives on Psychological Science</i> , 0(0):17456916231214460. PMID: 38165766.	
905		
906		
907		
908		
909		
910	Steven Piantadosi. 2023. Modern language models refute chomsky’s approach to language. <i>Lingbuzz Preprint</i> , lingbuzz, 7180.	
911		
912		
913	Eva Portelance, Yuguang Duan, Michael C. Frank, and Gary Lupyan. 2023. Predicting age of acquisition for children’s early vocabulary in five languages using language model surprisal . <i>Cognitive science</i> , 47 9:e13334.	
914		
915		
916		
917		
918	Jean Raven. 2003. Raven progressive matrices. In <i>Handbook of nonverbal assessment</i> , pages 223–237. Springer.	
919		
920		
921	Nils Reimers and Iryna Gurevych. 2019. Sentence transformers: Multilingual sentence embeddings using bert / roberta / xlm-roberta & co. with pytorch . Accessed: 2023-08-14.	
922		
923		
924		
925	Eleanor Rosch. 1975. Cognitive representations of semantic categories. <i>Journal of Experimental Psychology: General</i> , 104(3):192.	
926		
927		
928	Jenny R Saffran, Richard N Aslin, and Elissa L Newport. 1996. Statistical learning by 8-month-old infants. <i>Science</i> , 274(5294):1926–1928.	
929		
930		
931	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavathula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. <i>Communications of the ACM</i> , 64(9):99–106.	
932		
933		
934		
935	Leonard Salewski, Stephan Alaniz, Isabel Rio-Torto, Eric Schulz, and Zeynep Akata. 2024. In-context impersonation reveals large language models’ strengths and biases. <i>Advances in Neural Information Processing Systems</i> , 36.	
936		
937		
938		
939		
940	Raj Sanjay Shah, Vijay Marupudi, Reba Koenen, Khushi Bhardwaj, and Sashank Varma. 2023. Numeric magnitude comparison effects in large language models .	
941		
942		
943		
	Mihir Sharma, Ryan Ding, Raj Sanjay Shah, and Sashank Varma. 2024. Monolingual and bilingual language acquisition in language models.	944
		945
		946
	Noam Siegelman. 2020. Statistical learning abilities and their relation to language. <i>Language and Linguistics Compass</i> , 14(3):e12365.	947
		948
		949
	Richard E Snow, Patrick C Kyllonen, Brachia Marshalek, et al. 1984. The topography of ability and learning correlations. <i>Advances in the psychology of human intelligence</i> , 2(S 47):103.	950
		951
		952
		953
	Robert J Sternberg. 2000. <i>Handbook of intelligence</i> . Cambridge University Press.	954
		955
	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Slav Petrov, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul R. Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, Alexandre Frechette, Charlotte Smith, Laura Culp, Lev Proleev, Yi Luan, Xi Chen, James Lottes, Nathan Schucher, Federico Lebron, Alban Rrustemi, Natalie Clay, Phil Crone, Tomas Kocisky, Jeffrey Zhao, Bartek Perz, Dian Yu, Heidi Howard, Adam Bloniarz, Jack W. Rae, Han Lu, Laurent Sifre, Marcello Maggioni, Fred Alcober, Dan Garrette, Megan Barnes, Shantanu Thakoor, Jacob Austin, Gabriel Barth-Maron, William Wong, Rishabh Joshi, Rahma Chaabouni, Deeni Fatiha, Arun Ahuja, Ruibo Liu, Yunxuan Li, Sarah Cogan, Jeremy Chen, Chao Jia, Chenjie Gu, Qiao Zhang, Jordan Grimstad, Ale Jakse Hartman, Martin Chadwick, Gaurav Singh Tomar, Xavier Garcia, Evan Senter, Emanuel Taropa, Thanumalayan Sankaranarayanan Pillai, Jacob Devlin, Michael Laskin, Diego de Las Casas, Dasha Valter, Connie Tao, Lorenzo Blanco, Adrià Puigdomènech Badia, David Reitter, Mianna Chen, Jenny Brennan, Clara Rivera, Sergey Brin, Shariq Iqbal, Gabriela Surita, Jane Labanowski, Abhi Rao, Stephanie Winkler, Emilio Parisotto, Yiming Gu, Kate Olszewska, Yujing Zhang, Ravi Ad-danki, Antoine Miech, Annie Louis, Laurent El Shafey, Denis Teplyashin, Geoff Brown, Elliot Catt, Nithya Attaluri, Jan Balaguer, Jackie Xiang, Pidong Wang, Zoe Ashwood, Anton Briukhov, Albert Webson, Sanjay Ganapathy, Smit Sanghavi, Ajay Kannan, Ming-Wei Chang, Axel Stjerngren, Josip Djolonga, Yuting Sun, Ankur Bapna, Matthew Aitchison, Pedram Pejman, Henryk Michalewski, Tianhe Yu, Cindy Wang, Juliette Love, Junwhan Ahn, Dawn Bloxwich, Kehang Han, Peter Humphreys,	956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000
		1001
		1002
		1003

1004	Thibault Sellam, James Bradbury, Varun Godbole,	Zach Gleicher, Thi Avrahami, Anudhyan Boral,	1067
1005	Sina Samangoeei, Bogdan Damoc, Alex Kaskasoli,	Hansa Srinivasan, Vittorio Selo, Rhys May, Kon-	1068
1006	Sébastien M. R. Arnold, Vijay Vasudevan, Shubham	stantinos Aisopos, Léonard Hussenot, Livio Baldini	1069
1007	Agrawal, Jason Riesa, Dmitry Lepikhin, Richard Tan-	Soares, Kate Baumli, Michael B. Chang, Adrià Rec-	1070
1008	burn, Srivatsan Srinivasan, Hyeontaek Lim, Sarah	casens, Ben Caine, Alexander Pritzel, Filip Pavetic,	1071
1009	Hodkinson, Pranav Shyam, Johan Ferret, Steven	Fabio Pardo, Anita Gergely, Justin Frye, Vinay	1072
1010	Hand, Ankush Garg, Tom Le Paine, Jian Li, Yu-	Ramasesh, Dan Horgan, Kartikeya Badola, Nora	1073
1011	jia Li, Minh Giang, Alexander Neitz, Zaheer Abbas,	Kassner, Subhrajit Roy, Ethan Dyer, Víctor Cam-	1074
1012	Sarah York, Machel Reid, Elizabeth Cole, Aakanksha	pos, Alex Tomala, Yunhao Tang, Dalia El Badawy,	1075
1013	Chowdhery, Dipanjan Das, Dominika Rogozińska,	Elsbeth White, Basil Mustafa, Oran Lang, Ab-	1076
1014	Vitaly Nikolaev, Pablo Sprechmann, Zachary Nado,	hishek Jindal, Sharad Vikram, Zhitao Gong, Sergi	1077
1015	Lukas Zilka, Flavien Prost, Luheng He, Marianne	Caelles, Ross Hemsley, Gregory Thornton, Fangxi-	1078
1016	Monteiro, Gaurav Mishra, Chris Welty, Josh Newlan,	aoyu Feng, Wojciech Stokowiec, Ce Zheng, Phoebe	1079
1017	Dawei Jia, Miltiadis Allamanis, Clara Huiyi Hu,	Thacker, Çağlar Ünlü, Zhishuai Zhang, Moham-	1080
1018	Raoul de Liedekerke, Justin Gilmer, Carl Saroufim,	mad Saleh, James Svensson, Max Bileschi, Piyush	1081
1019	Shruti Rijhwani, Shaobo Hou, Disha Shrivastava,	Patil, Ankesh Anand, Roman Ring, Katerina Tsihlas,	1082
1020	Anirudh Baddepudi, Alex Goldin, Adnan Ozturel,	Arpi Vezer, Marco Selvi, Toby Shevlane, Mikel Ro-	1083
1021	Albin Cassirer, Yunhan Xu, Daniel Sohn, Deven-	driguez, Tom Kwiatkowski, Samira Daruki, Keran	1084
1022	dra Sachan, Reinald Kim Amplayo, Craig Swans-	Rong, Allan Dafoe, Nicholas FitzGerald, Keren	1085
1023	on, Dessie Petrova, Shashi Narayan, Arthur Guez,	Gu-Lemberg, Mina Khan, Lisa Anne Hendricks,	1086
1024	Siddhartha Brahma, Jessica Landon, Miteyan Patel,	Marie Pellat, Vladimir Feinberg, James Cobon-	1087
1025	Ruizhe Zhao, Kevin Villela, Luyu Wang, Wenhao	Kerr, Tara Sainath, Maribeth Rauh, Sayed Hadi	1088
1026	Jia, Matthew Rahtz, Mai Giménez, Legg Yeung,	Hashemi, Richard Ives, Yana Hasson, YaGuang	1089
1027	Hanzhao Lin, James Keeling, Petko Georgiev, Di-	Li, Eric Noland, Yuan Cao, Nathan Byrd, Le Hou,	1090
1028	ana Mincu, Boxi Wu, Salem Haykal, Rachel Sapu-	Qingze Wang, Thibault Sottiaux, Michela Paganini,	1091
1029	tro, Kiran Vodrahalli, James Qin, Zeynep Cankara,	Jean-Baptiste Lespiau, Alexandre Moufarek, Samer	1092
1030	Abhanshu Sharma, Nick Fernando, Will Hawkins,	Hassan, Kaushik Shivakumar, Joost van Amers-	1093
1031	Behnam Neyshabur, Solomon Kim, Adrian Hut-	foort, Amol Mandhane, Pratik Joshi, Anirudh	1094
1032	ter, Priyanka Agrawal, Alex Castro-Ros, George	Goyal, Matthew Tung, Andrew Brock, Hannah Shea-	1095
1033	van den Driessche, Tao Wang, Fan Yang, Shuo yiin	han, Vedant Misra, Cheng Li, Nemanja Rakićević,	1096
1034	Chang, Paul Komarek, Ross McIlroy, Mario Lučić,	Mostafa Dehghani, Fangyu Liu, Sid Mittal, Junhyuk	1097
1035	Guodong Zhang, Wael Farhan, Michael Sharman,	Oh, Seb Noury, Eren Sezener, Fantine Huot, Matthew	1098
1036	Paul Natsev, Paul Michel, Yong Cheng, Yamini	Lamm, Nicola De Cao, Charlie Chen, Gamaleldin	1099
1037	Bansal, Siyuan Qiao, Kris Cao, Siamak Shakeri,	Elsayed, Ed Chi, Mahdis Mahdieh, Ian Tenney, Nan	1100
1038	Christina Butterfield, Justin Chung, Paul Kishan	Hua, Ivan Petrychenko, Patrick Kane, Dylan Scand-	1101
1039	Rubenstein, Shivani Agrawal, Arthur Mensch, Kedar	inaro, Rishub Jain, Jonathan Uesato, Romina Datta,	1102
1040	Soparkar, Karel Lenc, Timothy Chung, Aedan Pope,	Adam Sadovsky, Oskar Bunyan, Dominik Rabiej,	1103
1041	Loren Maggiore, Jackie Kay, Priya Jhakra, Shibo	Shimu Wu, John Zhang, Gautam Vasudevan, Edouard	1104
1042	Wang, Joshua Maynez, Mary Phuong, Taylor Tobin,	Leurent, Mahmoud Alnahlawi, Ionut Georgescu, Nan	1105
1043	Andrea Tacchetti, Maja Trebacz, Kevin Robinson,	Wei, Ivy Zheng, Betty Chan, Pam G Rabinovitch,	1106
1044	Yash Katariya, Sebastian Riedel, Paige Bailey, Ke-	Piotr Stanczyk, Ye Zhang, David Steiner, Subhajt	1107
1045	fan Xiao, Nimesh Ghelani, Lora Aroyo, Ambrose	Naskar, Michael Azzam, Matthew Johnson, Adam	1108
1046	Slone, Neil Houlsby, Xuehan Xiong, Zhen Yang,	Paszke, Chung-Cheng Chiu, Jaume Sanchez Elias,	1109
1047	Elena Gribovskaya, Jonas Adler, Mateo Wirth, Lisa	Afroz Mohiuddin, Faizan Muhammad, Jin Miao,	1110
1048	Lee, Music Li, Thais Kagohara, Jay Pavagadhi, So-	Andrew Lee, Nino Vieillard, Sahitya Potluri, Jane	1111
1049	phie Bridgers, Anna Bortsova, Sanjay Ghemawat,	Park, Elnaz Davoodi, Jiageng Zhang, Jeff Stanway,	1112
1050	Zafarali Ahmed, Tianqi Liu, Richard Powell, Vijay	Drew Garmon, Abhijit Karmarkar, Zhe Dong, Jong	1113
1051	Bolina, Mariko Inuma, Polina Zablotskaia, James	Lee, Aviral Kumar, Luowei Zhou, Jonathan Evens,	1114
1052	Besley, Da-Woon Chung, Timothy Dozat, Ramona	William Isaac, Zhe Chen, Johnson Jia, Anselm	1115
1053	Comanescu, Xiance Si, Jeremy Greer, Guolong Su,	Levskaya, Zhenkai Zhu, Chris Gorgolewski, Peter	1116
1054	Martin Polacek, Raphaël Lopez Kaufman, Simon	Grabowski, Yu Mao, Alberto Magni, Kaisheng Yao,	1117
1055	Tokumine, Hexiang Hu, Elena Buchatskaya, Yingjie	Javier Snaider, Norman Casagrande, Paul Suga-	1118
1056	Miao, Mohamed Elhawaty, Aditya Siddhant, Nenad	nthan, Evan Palmer, Geoffrey Irving, Edward Loper,	1119
1057	Tomasev, Jinwei Xing, Christina Greer, Helen Miller,	Manaal Faruqui, Isha Arkatkar, Nanxin Chen, Izhak	1120
1058	Shereen Ashraf, Aurko Roy, Zizhao Zhang, Ada Ma,	Shafraan, Michael Fink, Alfonso Castaño, Irene Gian-	1121
1059	Angelos Filos, Milos Besta, Rory Blevins, Ted Kli-	nomis, Wooyeol Kim, Mikołaj Rybiński, Ashwin	1122
1060	menko, Chih-Kuan Yeh, Soravit Changpinyo, Jiaqi	Sreevatsa, Jennifer Prendki, David Soergel, Adrian	1123
1061	Mu, Oscar Chang, Mantas Pajarskas, Carrie Muir,	Goedeckemeyer, Willi Gierke, Mohsen Jafari, Meenu	1124
1062	Vered Cohen, Charline Le Lan, Krishna Haridasan,	Gaba, Jeremy Wiesner, Diana Gage Wright, Yawen	1125
1063	Amit Marathe, Steven Hansen, Sholto Douglas, Ra-	Wei, Harsha Vashisht, Yana Kulizhskaya, Jay Hoover,	1126
1064	jkumar Samuel, Mingqiu Wang, Sophia Austin,	Maigo Le, Lu Li, Chimezie Iwuanyanwu, Lu Liu,	1127
1065	Chang Lan, Jiepu Jiang, Justin Chiu, Jaime Alonso	Kevin Ramirez, Andrey Khorlin, Albert Cui, Tian	1128
1066	Lorenzo, Lars Lowe Sjösund, Sébastien Cevey,	LIN, Marin Georgiev, Marcus Wu, Ricardo Aguilar,	1129

1130	Keith Pallo, Abhishek Chakladar, Alena Repina, Xi-	Zheng, Francesco Pongetti, Mukarram Tariq, Yan-	1193
1131	hui Wu, Tom van der Weide, Priya Ponnappalli, Car-	hua Sun, Lucian Ionita, Mojtaba Seyedhosseini,	1194
1132	oline Kaplan, Jiri Simsa, Shuangfeng Li, Olivier	Pouya Tafti, Ragha Kotikalapudi, Zhiyu Liu, An-	1195
1133	Dousse, Fan Yang, Jeff Piper, Nathan Ie, Minnie	mol Gulati, Jasmine Liu, Xinyu Ye, Bart Chrzaszcz,	1196
1134	Lui, Rama Pasumarthi, Nathan Lintz, Anitha Vi-	Lily Wang, Nikhil Sethi, Tianrun Li, Ben Brown,	1197
1135	jayakumar, Lam Nguyen Thiet, Daniel Andor, Pedro	Shreya Singh, Wei Fan, Aaron Parisi, Joe Stanton,	1198
1136	Valenzuela, Cosmin Padurarur, Daiyi Peng, Kather-	Chenkai Kuang, Vinod Koverkathu, Christopher A.	1199
1137	ine Lee, Shuyuan Zhang, Somer Greene, Duc Dung	Choquette-Choo, Yunjie Li, TJ Lu, Abe Ittycheriah,	1200
1138	Nguyen, Paula Kurylowicz, Sarmishta Velury, Se-	Prakash Shroff, Pei Sun, Mani Varadarajan, Sanaz Ba-	1201
1139	bastian Krause, Cassidy Hardin, Lucas Dixon, Lili	hargam, Rob Willoughby, David Gaddy, Ishita Das-	1202
1140	Janzer, Kiam Choo, Ziqiang Feng, Biao Zhang,	gupta, Guillaume Desjardins, Marco Cornero, Brona	1203
1141	Achintya Singhal, Tejasi Latkar, Mingyang Zhang,	Robenek, Bhavishya Mittal, Ben Albrecht, Ashish	1204
1142	Quoc Le, Elena Allica Abellan, Dayou Du, Dan McK-	Shenoy, Fedor Moiseev, Henrik Jacobsson, Alireza	1205
1143	innon, Natasha Antropova, Tolga Bolukbasi, Orgad	Ghaffarkhah, Morgane Rivière, Alanna Walton, Clé-	1206
1144	Keller, David Reid, Daniel Finchelstein, Maria Abi	ment Crepy, Alicia Parrish, Yuan Liu, Zongwei	1207
1145	Raad, Remi Crocker, Peter Hawkins, Robert Dadashi,	Zhou, Clement Farabet, Carey Radebaugh, Praveen	1208
1146	Colin Gaffney, Sid Lall, Ken Franko, Egor Filonov,	Srinivasan, Claudia van der Salm, Andreas Fidje-	1209
1147	Anna Bulanova, Rémi Leblond, Vikas Yadav, Shirley	land, Salvatore Scellato, Eri Latorre-Chimoto, Hanna	1210
1148	Chung, Harry Askham, Luis C. Cobo, Kelvin Xu,	Klimczak-Plucińska, David Bridson, Dario de Ce-	1211
1149	Felix Fischer, Jun Xu, Christina Sorokin, Chris Al-	sare, Tom Hudson, Piermaria Mendolicchio, Lexi	1212
1150	berti, Chu-Cheng Lin, Colin Evans, Hao Zhou, Alek	Walker, Alex Morris, Ivo Penchev, Matthew Mauger,	1213
1151	Dimitriev, Hannah Forbes, Dylan Banarse, Zora	Alexey Guseynov, Alison Reid, Seth Odoom, Lucia	1214
1152	Tung, Jeremiah Liu, Mark Omernick, Colton Bishop,	Loher, Victor Cotruta, Madhavi Yenugula, Dominik	1215
1153	Chintu Kumar, Rachel Sterneck, Ryan Foley, Rohan	Grewe, Anastasia Petrushkina, Tom Duerig, Antonio	1216
1154	Jain, Swaroop Mishra, Jiawei Xia, Taylor Bos, Ge-	Sanchez, Steve Yadlowsky, Amy Shen, Amir Globerson,	1217
1155	offrey Cideron, Ehsan Amid, Francesco Piccinno,	Adam Kurzrok, Lynette Webb, Sahil Dua, Dong	1218
1156	Xingyu Wang, Praseem Banzal, Petru Gurita, Hila	Li, Preethi Lahoti, Surya Bhupatiraju, Dan Hurt, Ha-	1219
1157	Noga, Premal Shah, Daniel J. Mankowitz, Alex	roon Qureshi, Ananth Agarwal, Tomer Shani, Matan	1220
1158	Polozov, Nate Kushman, Victoria Krakovna, Sasha	Eyal, Anuj Khare, Shreyas Rammohan Belle, Lei	1221
1159	Brown, MohammadHossein Bateni, Dennis Duan,	Wang, Chetan Tekur, Mihir Sanjay Kale, Jinliang	1222
1160	Vlad Firoiu, Meghana Thotakuri, Tom Natan, An-	Wei, Ruoxin Sang, Brennan Saeta, Tyler Liechty,	1223
1161	had Mohananey, Matthieu Geist, Sidharth Mudgal,	Yi Sun, Yao Zhao, Stephan Lee, Pandu Nayak, Doug	1224
1162	Sertan Girgin, Hui Li, Jiayu Ye, Ofir Roval, Reiko	Fritz, Manish Reddy Vuyyuru, John Aslanides, Nidhi	1225
1163	Tojo, Michael Kwong, James Lee-Thorp, Christo-	Vyas, Martin Wicke, Xiao Ma, Taylan Bilal, Ev-	1226
1164	pher Yew, Quan Yuan, Sumit Bagri, Danila Sinopal-	genii Eltyshev, Daniel Balle, Nina Martin, Hardie	1227
1165	nikov, Sabela Ramos, John Mellor, Abhishek Sharma,	Cate, James Manyika, Keyvan Amiri, Yelin Kim,	1228
1166	Aliaksei Severyn, Jonathan Lai, Kathy Wu, Heng-	Xi Xiong, Kai Kang, Florian Luisier, Nilesh Tripu-	1229
1167	Tze Cheng, David Miller, Nicolas Sonnerat, Denis	raneni, David Madras, Mandy Guo, Austin Waters,	1230
1168	Vnukov, Rory Greig, Jennifer Beattie, Emily Cave-	Oliver Wang, Joshua Ainslie, Jason Baldrige, Han	1231
1169	ness, Libin Bai, Julian Eisenschlos, Alex Korchem-	Zhang, Garima Pruthi, Jakob Bauer, Feng Yang, Ri-	1232
1170	niy, Tomy Tsai, Mimi Jasarevic, Weize Kong, Phuong	ham Mansour, Jason Gelman, Yang Xu, George	1233
1171	Dao, Zeyu Zheng, Frederick Liu, Fan Yang, Rui	Polovets, Ji Liu, Honglong Cai, Warren Chen, Xi-	1234
1172	Zhu, Mark Geller, Tian Huey Teh, Jason Sanmiya,	angHai Sheng, Emily Xue, Sherjil Ozair, Adams Yu,	1235
1173	Evgeny Gladchenko, Nejc Trdin, Andrei Sozanschi,	Christof Angermueller, Xiaowei Li, Weiren Wang, Ju-	1236
1174	Daniel Toyama, Evan Rosen, Sasan Tavakkol, Lint-	lia Wiesinger, Emmanouil Koukoumidis, Yuan Tian,	1237
1175	ing Xue, Chen Elkind, Oliver Woodman, John Car-	Anand Iyer, Madhu Gurusurthy, Mark Goldenson,	1238
1176	penter, George Papamakarios, Rupert Kemp, Sushant	Parashar Shah, MK Blake, Hongkun Yu, Anthony	1239
1177	Kafle, Tanya Grunina, Rishika Sinha, Alice Tal-	Urbanowicz, Jennimaria Palomaki, Chrisantha Fer-	1240
1178	bert, Abhimanyu Goyal, Diane Wu, Denese Owusu-	nando, Kevin Brooks, Ken Durden, Harsh Mehta,	1241
1179	Afriyie, Cosmo Du, Chloe Thornton, Jordi Pont-	Nikola Momchev, Elahe Rahimtoroghi, Maria Geor-	1242
1180	Tuset, Pradyumna Narayana, Jing Li, Sabaer Fatehi,	gaki, Amit Raul, Sebastian Ruder, Morgan Red-	1243
1181	John Wieting, Omar Ajmeri, Benigno Uria, Tao Zhu,	shaw, Jinhyuk Lee, Komal Jalan, Dinghua Li, Ginger	1244
1182	Yeongil Ko, Laura Knight, Amélie Héliou, Ning	Perng, Blake Hechtman, Parker Schuh, Milad Nasr,	1245
1183	Niu, Shane Gu, Chenxi Pang, Dustin Tran, Yeqing	Mia Chen, Kieran Milan, Vladimir Mikulik, Trevor	1246
1184	Li, Nir Levine, Ariel Stolovich, Norbert Kalb, Re-	Strohman, Juliana Franco, Tim Green, Demis Has-	1247
1185	becca Santamaria-Fernandez, Sonam Goenka, Wenny	sabis, Koray Kavukcuoglu, Jeffrey Dean, and Oriol	1248
1186	Yustalim, Robin Strudel, Ali Elqursh, Balaji Laksh-	Vinyals. 2023. Gemini: A family of highly capable	1249
1187	minarayanan, Charlie Deck, Shyam Upadhyay, Hyo	multimodal models .	1250
1188	Lee, Mike Dusenberry, Zonglin Li, Xuezhi Wang,	Louis Leon Thurstone. 1938. Primary mental abilities.	1251
1189	Kyle Levin, Raphael Hoffmann, Dan Holtmann-	<i>Psychometric monographs</i> .	1252
1190	Rice, Olivier Bachem, Summer Yue, Sho Arora,	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	1253
1191	Eric Malmi, Daniil Mirylenka, Qijun Tan, Christy	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	1254
1192	Koh, Soheil Hassas Yeganeh, Siim Põder, Steven		

1255	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models .	1309
1256		1310
1257		1311
1258		1312
1259	Peter D Turney. 2005. Measuring semantic similarity by latent relational analysis. <i>arXiv preprint cs/0508053</i> .	1313
1260		1314
1261	Peter D Turney and Patrick Pantel. 2010. From frequency to meaning: Vector space models of semantics. <i>Journal of artificial intelligence research</i> , 37:141–188.	1315
1262		1316
1263		1317
1264		1318
1265	J. P. Van Overschelde, K. A. Rawson, and J. Dunlosky. 2004. Category norms: An updated and expanded version of the norms. <i>Journal of Memory and Language</i> , 50:289–335.	1319
1266		1320
1267		1321
1268		1322
1269	Siddhartha Vemuri, Raj Sanjay Shah, and Sashank Varma. 2024. Evaluating typicality in combined language and vision model concept representations. In <i>Under review</i> .	1323
1270		1324
1271		1325
1272		1326
1273	Alex Warstadt and Samuel R. Bowman. 2024. What artificial neural networks can tell us about human language acquisition .	1327
1274		1328
1275		1329
1276	Alex Warstadt, Aaron Mueller, Leshem Choshen, Ethan Wilcox, Chengxu Zhuang, Juan Ciro, Rafael Mosquera, Bhargavi Paranjabe, Adina Williams, Tal Linzen, and Ryan Cotterell, editors. 2023. <i>Proceedings of the BabyLM Challenge at the 27th Conference on Computational Natural Language Learning</i> . Association for Computational Linguistics, Singapore.	1330
1277		1331
1278		1332
1279		1333
1280		1334
1281		1335
1282		1336
1283	Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananeey, Wei Peng, Sheng-Fu Wang, and Samuel R Bowman. 2020. Blimp: The benchmark of linguistic minimal pairs for english. <i>Transactions of the Association for Computational Linguistics</i> , 8:377–392.	1337
1284		1338
1285		1339
1286		1340
1287		1341
1288	Taylor Webb, Keith J Holyoak, and Hongjing Lu. 2023. Emergent analogical reasoning in large language models. <i>Nature Human Behaviour</i> , 7(9):1526–1541.	1342
1289		1343
1290		1344
1291	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022a. Emergent abilities of large language models. <i>arXiv preprint arXiv:2206.07682</i> .	1345
1292		1346
1293		1347
1294		1348
1295		1349
1296		1350
1297		1351
1298	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022b. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in Neural Information Processing Systems</i> , 35:24824–24837.	1352
1299		1353
1300		1354
1301		1355
1302		1356
1303	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	1357
1304		1358
1305		1359
1306		1360
1307		1361
1308		1362
	Junjie Ye, Xuanting Chen, Nuo Xu, Can Zu, Zekai Shao, Shichun Liu, Yuhan Cui, Zeyang Zhou, Chao Gong, Yang Shen, Jie Zhou, Siming Chen, Tao Gui, Qi Zhang, and Xuanjing Huang. 2023. A comprehensive capability analysis of GPT-3 and GPT-3.5 series models. <i>arXiv preprint arXiv:2303.10420</i> .	1363
	Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, and Songfang Huang. 2023. How well do large language models perform in arithmetic tasks?	1364
	Xikun Zhang, Deepak Ramachandran, Ian Tenney, Yanai Elazar, and Dan Roth. 2020. Do language embeddings capture scales? In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 4889–4896, Online. Association for Computational Linguistics.	1365
	Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, and Jiantao Jiao. 2023. Starling-7b: Improving llm helpfulness and harmlessness with rlaiif.	1366
	Yan Zhuang, Qi Liu, Yuting Ning, Weizhe Huang, Rui Lv, Zhenya Huang, Guanhao Zhao, Zheng Zhang, Qingyang Mao, Shijin Wang, et al. 2023. Efficiently measuring the cognitive ability of llms: An adaptive testing perspective. <i>arXiv preprint arXiv:2306.10512</i> .	1367
		1368
		1369
		1370
		1371
		1372
		1373
		1374
		1375
		1376
		1377
		1378
		1379
		1380
		1381
		1382
		1383
		1384
		1385
		1386
		1387
		1388
		1389
		1390
		1391
		1392
		1393
		1394
		1395
		1396
		1397
		1398
		1399
		1400
		1401
		1402
		1403
		1404
		1405
		1406
		1407
		1408
		1409
		1410
		1411
		1412
		1413
		1414
		1415
		1416
		1417
		1418
		1419
		1420
		1421
		1422
		1423
		1424
		1425
		1426
		1427
		1428
		1429
		1430
		1431
		1432
		1433
		1434
		1435
		1436
		1437
		1438
		1439
		1440
		1441
		1442
		1443
		1444
		1445
		1446
		1447
		1448
		1449
		1450
		1451
		1452
		1453
		1454
		1455
		1456
		1457
		1458
		1459
		1460
		1461
		1462
		1463
		1464
		1465
		1466
		1467
		1468
		1469
		1470
		1471
		1472
		1473
		1474
		1475
		1476
		1477
		1478
		1479
		1480
		1481
		1482
		1483
		1484
		1485
		1486
		1487
		1488
		1489
		1490
		1491
		1492
		1493
		1494
		1495
		1496
		1497
		1498
		1499
		1500
		1501
		1502
		1503
		1504
		1505
		1506
		1507
		1508
		1509
		1510
		1511
		1512
		1513
		1514
		1515
		1516
		1517
		1518
		1519
		1520
		1521
		1522
		1523
		1524
		1525
		1526
		1527
		1528
		1529
		1530
		1531
		1532
		1533
		1534
		1535
		1536
		1537
		1538
		1539
		1540
		1541
		1542
		1543
		1544
		1545
		1546
		1547
		1548
		1549
		1550
		1551
		1552
		1553
		1554
		1555
		1556
		1557
		1558
		1559
		1560
		1561
		1562
		1563
		1564
		1565
		1566
		1567
		1568
		1569
		1570
		1571
		1572
		1573
		1574
		1575
		1576
		1577
		1578
		1579
		1580
		1581
		1582
		1583
		1584
		1585
		1586
		1587
		1588
		1589
		1590
		1591
		1592
		1593
		1594
		1595
		1596
		1597
		1598
		1599
		1600
		1601
		1602
		1603
		1604
		1605
		1606
		1607
		1608
		1609
		1610
		1611
		1612
		1613
		1614
		1615
		1616
		1617
		1618
		1619
		1620
		1621
		1622
		1623
		1624
		1625
		1626
		1627
		1628
		1629
		1630
		1631
		1632
		1633
		1634
		1635
		1636
		1637
		1638
		1639
		1640
		1641
		1642
		1643
		1644
		1645
		1646
		1647
		1648
		1649
		1650
		1651
		1652
		1653
		1654
		1655
		1656
		1657
		1658
		1659
		1660
		1661
		1662
		1663
		1664
		1665
		1666
		1667
		1668
		1669
		1670
		1671
		1672
		1673
		1674
		1675
		1676
		1677
		1678
		1679
		1680
		1681
		1682
		1683
		1684
		1685
		1686
		1687
		1688
		1689
		1690
		1691
		1692
		1693
		1694
		1695
		1696
		1697
		1698
		1699
		1700
		1701
		1702
		1703
		1704
		1705
		1706
		1707
		1708

Table 4: Magnitude Comparison effects. Distance Effect: Averaged R^2 values of different LLMs when fitting a linear function on the cosine-similarity vs distance plot. Size Effect: Averaged R^2 values of different LLMs when fitting a linear function on the cosine-similarity vs size-difference plot. Ratio Effect: Averaged R^2 values of different LLMs when fitting a negative exponential function on the cosine-similarity vs ratio plot. Note: Each value is averaged across all three input types and all model layers to produce one generalizable score. MDS Stress: The stress value is a measure of how well the distances between the points in the multidimensional space represent the dissimilarities of the original data points (lower is better). MDS Correlation: Correlation between the MDS solutions and the expected values of human MNL. Range (Sim): This indicates the range of the cosine-similarities. Max (sim): This indicates the maximum similarity between any two numbers. Range and Max (sim) describe the y-axis.

Model	Distance Effect	Ratio Effect	Size Effect	MDS Stress	MDS Correlation	Range (Sim)	Max (Sim)
Amber-7B	0.913	0.591	0.607	0.157	0.572	0.008	0.995
Falcon-7B	0.928	0.838	0.725	0.183	0.655	0.286	0.779
Starling-LM-7B-alpha	0.522	0.187	0.494	0.320	0.305	0.001	0.995
Llama-2-7B	0.670	0.614	0.535	0.122	0.547	0.016	0.983
Llama-2-13B	0.672	0.263	0.421	0.234	0.372	0.002	0.999
Mistral-7B	0.641	0.233	0.244	0.287	0.425	0.001	0.996
Mistral-7B-Instruct	0.637	0.543	0.182	0.317	0.512	0.001	0.992
Qwen-0.5B	0.833	0.553	0.215	0.246	0.679	0.064	0.911
Qwen-1.8B	0.878	0.301	0.330	0.198	0.328	0.107	0.902
Qwen-4B	0.881	0.264	0.330	0.215	0.581	0.160	0.763
Qwen-7B	0.858	0.616	0.257	0.153	0.636	0.129	0.734
Qwen-14B	0.783	0.507	0.206	0.248	0.369	0.138	0.710
Pythia-70M	0.829	0.429	0.418	0.204	0.463	0.060	0.949
Pythia-160M	0.947	0.665	0.382	0.231	0.715	0.042	0.970
Pythia-410M	0.926	0.679	0.393	0.210	0.710	0.041	0.972
Pythia-1B	0.944	0.702	0.470	0.196	0.725	0.037	0.973
Pythia-1.4B	0.933	0.764	0.600	0.203	0.658	0.022	0.983
Pythia-2.8B	0.961	0.723	0.459	0.256	0.737	0.009	0.993
Pythia-6.9B	0.909	0.713	0.535	0.195	0.663	0.013	0.990
Pythia-12B	0.846	0.595	0.540	0.189	0.620	0.007	0.993

- **Size effect:** Given two comparisons of the same distance (i.e., of the same value for $|x - y|$), the smaller the numbers, the faster the comparison (Parkman, 1971).
- **Ratio effect** (refer to figure 1 (A) bottom): The time taken by humans to compare two numbers (x,y) is a decreasing function of the ratio of the larger number over the smaller number $\frac{\max(x,y)}{\min(x,y)}$ (Halberda et al., 2008).
- **Multidimensional scaling:** Along with the three effects, we investigate the consistency of the latent number representations of PLMs with the human MNL using multidimensional scaling (Borg and Groenen, 2005; Ding, 2018). MDS recovers the latent representation from the cosine (dis)similarities between the vector representations of all pairs of numbers (for a given LLM, layer, and number format). This is evaluated by the correlation between the positions of the numbers 1 to 9 in the MDS solution and the expected values ($\log(1)$ to $\log(9)$) of the human MNL (refer to the correlation value in table 4).

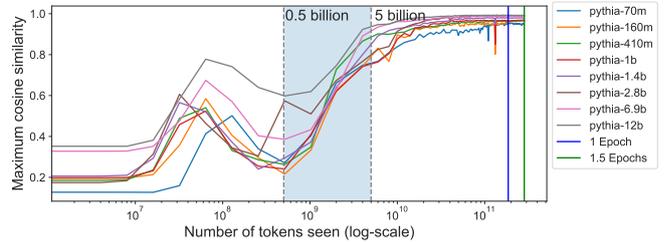


Figure 5: Development of the idea of "numbers" in Pythia. The y-axis indicates the maximum cosine similarity between the latent representations of any two number words/ digits.

Beyond these effects, we investigate the development of the latent understanding of the concept of "numbers" in the PLMs. As PLMs see more data, the average values of the similarity become larger, indicating that models learn the distinctions among numbers better (refer to figure 5). This is further substantiated by figure 6, where the similarities between number words develop to be greater than the similarity between (number, non-number) words and (non-number, non-number) words.

B.2 Linguistic Abilities

The 12 phenomena tested by BLiMP are as follows:

Table 5: Accuracy of different language models on the BLiMP linguistic acceptability tasks.

Model	BLiMP	Syntax	Semantic	Morphology
Amber-7B	0.794 (\pm 0.174)	0.779 (\pm 0.011)	0.736 (\pm 0.011)	0.888 (\pm 0.009)
Falcon-7B	0.817 (\pm 0.173)	0.797 (\pm 0.011)	0.758 (\pm 0.011)	0.917 (\pm 0.008)
Starling-LM-7B-alpha	0.827 (\pm 0.161)	0.799 (\pm 0.011)	0.788 (\pm 0.011)	0.938 (\pm 0.007)
Llama-2-7B	0.818 (\pm 0.165)	0.792 (\pm 0.011)	0.782 (\pm 0.011)	0.917 (\pm 0.008)
Llama-2-13B	0.793 (\pm 0.184)	0.757 (\pm 0.011)	0.767 (\pm 0.011)	0.898 (\pm 0.008)
Mistral-7B	0.829 (\pm 0.174)	0.801 (\pm 0.011)	0.780 (\pm 0.010)	0.940 (\pm 0.007)
Mistral-7B-Instruct	0.834 (\pm 0.149)	0.808 (\pm 0.011)	0.788 (\pm 0.011)	0.931 (\pm 0.008)
Qwen-0.5B	0.785 (\pm 0.176)	0.759 (\pm 0.012)	0.718 (\pm 0.012)	0.907 (\pm 0.008)
Qwen-1.8B	0.792 (\pm 0.162)	0.777 (\pm 0.012)	0.764 (\pm 0.011)	0.875 (\pm 0.010)
Qwen-4B	0.730 (\pm 0.154)	0.694 (\pm 0.013)	0.728 (\pm 0.013)	0.814 (\pm 0.012)
Qwen-7B	0.789 (\pm 0.156)	0.769 (\pm 0.012)	0.736 (\pm 0.012)	0.885 (\pm 0.010)
Qwen-14B	0.792 (\pm 0.144)	0.775 (\pm 0.012)	0.747 (\pm 0.012)	0.881 (\pm 0.010)
Pythia-70M	0.723 (\pm 0.210)	0.701 (\pm 0.012)	0.628 (\pm 0.012)	0.872 (\pm 0.010)
Pythia-160M	0.749 (\pm 0.207)	0.717 (\pm 0.012)	0.718 (\pm 0.011)	0.864 (\pm 0.010)
Pythia-410M	0.815 (\pm 0.169)	0.785 (\pm 0.011)	0.752 (\pm 0.011)	0.935 (\pm 0.007)
Pythia-1B	0.806 (\pm 0.198)	0.782 (\pm 0.011)	0.728 (\pm 0.011)	0.935 (\pm 0.007)
Pythia-1.4B	0.819 (\pm 0.173)	0.792 (\pm 0.011)	0.768 (\pm 0.011)	0.931 (\pm 0.008)
Pythia-2.8B	0.827 (\pm 0.156)	0.800 (\pm 0.011)	0.782 (\pm 0.011)	0.925 (\pm 0.007)
Pythia-6.9B	0.809 (\pm 0.179)	0.792 (\pm 0.011)	0.750 (\pm 0.011)	0.913 (\pm 0.008)
Pythia-12B	0.829 (\pm 0.158)	0.804 (\pm 0.011)	0.778 (\pm 0.011)	0.932 (\pm 0.007)
Gemini	NA	NA	NA	NA
GPT-3.5-Turbo	0.825 (\pm 0.166)	0.818 (\pm 0.010)	0.781 (\pm 0.011)	0.931 (\pm 0.007)
GPT-4	0.849 (\pm 0.120)	0.797 (\pm 0.010)	0.801 (\pm 0.009)	0.941 (\pm 0.007)

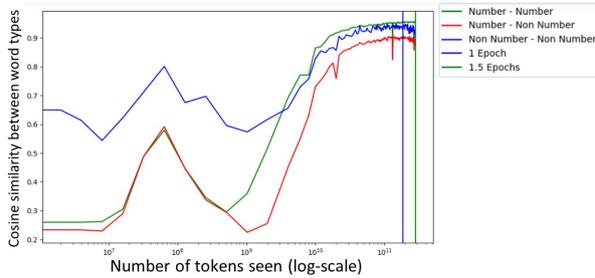


Figure 6: Development of the idea of "numbers" in Pythia. The y-axis shows the cosine similarity between word types. The cosine similarity values are averaged over all input types, all model layers, and all model sizes.

- Anaphor agreement (morphology): This linguistic phenomenon tests if an anaphor (pronoun) adheres to the antecedent (noun or phrase it refers to) in terms of gender, number, or person.
- Argument Structure (syntax): The argument structure tests the relationship between a verb and its arguments (such as nouns or noun

phrases).

- Binding (syntax, semantics): This tests the structural relationship between an anaphor (pronoun) and its antecedent (noun or phrase it refers to).
- Control/ Raising (syntax, semantics): These structures test how semantics differ by syntactical variations of subjects/verbs in subordinate and main clauses.
- Determiner-noun agreement (morphology): This tests the agreements of the determiners with the corresponding nouns in number (singular or plural) and sometimes gender (e.g., "his" for masculine nouns, "her" for feminine nouns).
- Ellipsis (syntax): This refers to the omission of words from a sentence that can be understood from the context.
- Filler-gap (syntax): This tests the syntactic structure of sentences that include phrasal movements (wh-questions, relative clauses).

- Irregular forms (morphology): Forms in language that do not follow regular patterns and may need to be memorized. For example, the superlative of good is better, best, and not gooder, goodest.
- Island effects (syntax): These test the constraints on syntactic environments where the gap in a filler-gap dependency can occur.
- NPI licensing (semantics): This phenomenon tests the constrained situations where negative polarity items like any and ever are limited to the scope of negation.
- Quantifiers (semantics): This phenomenon tests the constraints regarding the placement of quantifiers. Specifically, BLiMP looks at superlative quantifiers (such as "at least") that cannot occur within negation, and definite quantifiers and determiners cannot function as subjects in existential "there" constructions.
- Subject-verb agreement (morphology): The subject and tense forms of the verb must agree on the number, for example, singular vs plural.

Table 5 shows that the PLMs are more accurate in morphology than in language syntax and semantics. Most models also perform better on syntactic language features than semantic language features.

B.3 Conceptual Understanding

Table 7 shows the human alignment of PLMs on their concept understanding for different operationalization methods. We see that Gemini, GPT-3.5-Turbo, and GPT-4 perform better than other models. Furthermore, Surprisal and Prompting-based methods are stronger techniques for evaluating conceptual understanding of models than representation-based methods. Given the higher performance of Prompting methods on three API-based models, we only show the category-wise results for those models. The final prompt design is given in section B.3.1 and table 11. Tables 8, 9, and 10 show Spearman’s correlation on the categories along with the standard deviation, the minimum correlation, and the maximum correlation. We perform the same infilling tasks 50 times for each category to account for variations in generations. We note that the models often failed to return all the options in the in-filling task. We discard such situations in our analysis.

Note: Under the closeness judgment protocol, our experiments fail to match up to the performance of the models used by Vemuri et al. (2024). This is because our choice of open-source models only provides token representations, on which we later perform an aggregation operation. This aggregation operation leads to a loss of information. In contrast, Vemuri et al. (2024) use sentence-transformer models (Reimers and Gurevych, 2019), which provide singular latent representation for longer text. This variation in experimentation leads to the difference in alignment scores.

Table 6: Typicality effects: Comparing Average Spearman’s correlation score across categories from tables 8, 9, and 10.

Categories	GPT 3.5	GPT 4	Gemini
bird	0.183	0.536	0.353
carpenters tool	0.418	0.679	0.610
clothing	0.022	0.594	0.155
color	-0.016	0.882	0.569
dwelling	0.208	0.335	0.340
earth formation	0.251	0.496	0.155
fabric	0.48	0.708	0.504
fish	0.183	0.643	0.247
flower	0.48	0.772	0.515
flying thing	0.07	0.249	0.184
footwear	0.118	0.521	0.218
four-legged animal	0.435	0.818	0.537
fruit	0.465	0.726	0.508
furniture	0.069	0.525	0.147
gardeners tool	0.355	0.557	0.507
green thing	0.196	0.572	0.335
insect	0.18	0.629	0.286
instrument	0.194	0.709	0.450
kitchen utensil	0.384	0.624	0.252
ship	0.104	0.233	-0.078
snake	0.177	0.419	0.328
toy	0.299	0.480	0.169
tree	0.333	0.557	0.445
vegetable	0.096	0.783	0.121
vehicle	0.17	0.381	0.033
weapon	0.348	0.421	0.239
weather	0.333	0.255	0.274
Average	0.242	0.559	0.311

Table 7: Results for the typicality effects using the three methods

Model	Latent Representations	Surprisal Values				Prompting
		Zero-shot	One-shot	Two-shot	Three-shot	
Amber-7B	0.083	0.250	0.227	0.261	0.247	NA
Falcon-7B	-0.116	0.180	0.215	0.242	0.200	NA
Starling-LM-7B-alpha	-0.003	0.258	0.211	0.215	0.235	NA
Llama-2-7B	-0.065	0.238	0.213	0.202	0.207	NA
Llama-2-13B	0.076	0.247	0.163	0.183	0.170	NA
Mistral-7B	-0.025	0.245	0.219	0.261	0.257	NA
Mistral-7B-Instruct	0.033	0.255	0.192	0.204	0.235	NA
Qwen-0.5B	0.072	0.282	0.264	0.288	0.250	NA
Qwen-1.8B	0.114	0.235	0.246	0.251	0.215	NA
Qwen-4B	0.001	0.246	0.217	0.252	0.193	NA
Qwen-7B	0.006	0.229	0.203	0.220	0.220	NA
Qwen-14B	-0.140	0.249	0.224	0.207	0.199	NA
Pythia-70M	0.005	0.211	0.266	0.291	0.285	NA
Pythia-160M	0.067	0.260	0.263	0.276	0.264	NA
Pythia-410M	0.126	0.284	0.235	0.282	0.242	NA
Pythia-1B	0.090	0.280	0.309	0.287	0.264	NA
Pythia-1.4B	0.074	0.283	0.249	0.267	0.235	NA
Pythia-2.8B	0.221	0.273	0.286	0.267	0.236	NA
Pythia-6.9B	0.105	0.280	0.264	0.250	0.220	NA
Pythia-12B	0.184	0.291	0.248	0.274	0.270	NA
Gemini	NA	NA	NA	NA	NA	0.311
GPT-3.5-Turbo	NA	0.231	0.248	0.299	0.270	0.242
GPT-4	NA	0.428	0.471	0.399	0.402	0.559

Table 8: Average Spearman’s correlation score for each category on 50 runs of each in-filling experiment on the Gemini-Pro model.

Categories	Average SpearmanR	Minimum Values	Maximum Values	Std Dev
bird	0.353	-0.156	0.582	0.144
carpenters tool	0.610	0.417	0.885	0.104
clothing	0.155	-0.104	0.523	0.141
color	0.569	-0.147	0.916	0.260
dwelling	0.340	0.140	0.499	0.086
earth formation	0.155	-0.449	0.494	0.191
fabric	0.504	0.125	0.811	0.168
fish	0.247	-0.505	0.611	0.265
flower	0.515	-0.183	0.779	0.208
flying thing	0.184	-0.068	0.602	0.193
footwear	0.218	-0.340	0.569	0.215
four-legged animal	0.537	0.225	0.689	0.099
fruit	0.508	-0.019	0.802	0.222
furniture	0.147	-0.479	0.663	0.310
gardeners tool	0.507	0.025	0.771	0.151
green thing	0.335	0.037	0.535	0.117
insect	0.286	-0.121	0.635	0.193
instrument	0.450	0.092	0.832	0.175
kitchen utensil	0.252	-0.164	0.691	0.243
ship	-0.078	-0.414	0.277	0.179
snake	0.328	-0.156	0.596	0.147
toy	0.169	-0.203	0.526	0.174
tree	0.445	0.257	0.585	0.073
vegetable	0.121	-0.322	0.596	0.184
vehicle	0.033	-0.053	0.236	0.055
weapon	0.239	-0.173	0.577	0.193
weather	0.274	-0.029	0.591	0.147

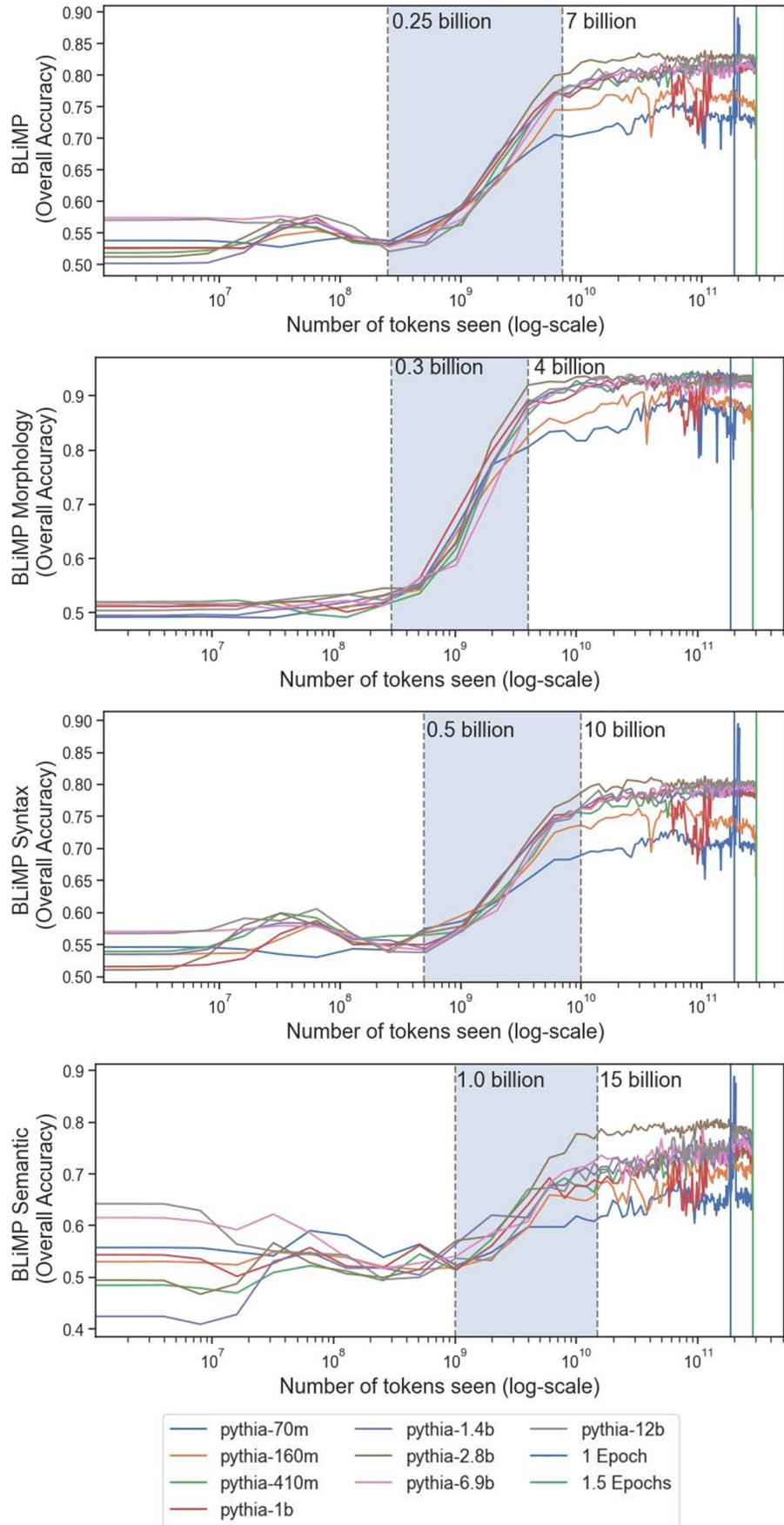


Figure 7: Developmental trajectory of the Pythia suite of models on the BLiMP linguistic acceptability tasks.

Table 9: Average Spearman’s correlation score for each category on 50 runs of each in-filling experiment on the GPT-3.5-Turbo model.

Categories	Average SpearmanR	Minimum Values	Maximum Values	Std Dev
bird	0.183	-0.209	0.552	0.209
carpenters tool	0.418	-0.162	0.858	0.282
clothing	0.022	-0.321	0.540	0.192
color	-0.016	-0.596	0.564	0.261
dwelling	0.208	-0.053	0.400	0.123
earth formation	0.251	-0.296	0.562	0.217
fabric	0.480	-0.044	0.767	0.233
fish	0.183	-0.326	0.690	0.280
flower	0.480	-0.301	0.800	0.269
flying thing	0.070	-0.181	0.377	0.149
footwear	0.118	-0.439	0.581	0.241
four-legged animal	0.435	-0.264	0.869	0.292
fruit	0.465	-0.006	0.868	0.241
furniture	0.069	-0.325	0.447	0.195
gardeners tool	0.355	-0.311	0.796	0.294
green thing	0.196	-0.337	0.572	0.211
insect	0.180	-0.248	0.503	0.201
instrument	0.194	-0.242	0.466	0.191
kitchen utensil	0.384	-0.610	0.797	0.334
ship	0.104	-0.314	0.599	0.250
snake	0.177	-0.244	0.591	0.196
toy	0.299	-0.210	0.603	0.180
tree	0.333	-0.199	0.731	0.289
vegetable	0.096	-0.191	0.542	0.172
vehicle	0.170	-0.381	0.381	0.201
weapon	0.348	-0.058	0.609	0.156
weather	0.333	-0.425	0.662	0.236

Table 10: Average Spearman’s correlation score for each category on 50 runs of each in-filling experiment on the GPT-4 model.

Categories	Average SpearmanR	Minimum Values	Maximum Values	Std Dev
bird	0.536	0.355	0.756	0.098
carpenters tool	0.679	0.549	0.843	0.078
clothing	0.594	0.350	0.751	0.100
color	0.882	0.813	0.952	0.035
dwelling	0.335	0.183	0.497	0.070
earth formation	0.496	0.373	0.628	0.061
fabric	0.708	0.583	0.801	0.052
fish	0.643	-0.237	0.817	0.218
flower	0.772	0.629	0.869	0.057
flying thing	0.249	-0.118	0.704	0.221
footwear	0.521	0.191	0.721	0.112
four-legged animal	0.818	0.634	0.906	0.056
fruit	0.726	0.567	0.868	0.069
furniture	0.525	0.381	0.605	0.055
gardeners tool	0.557	0.314	0.757	0.098
green thing	0.572	0.444	0.709	0.050
insect	0.629	0.451	0.871	0.103
instrument	0.709	0.585	0.885	0.064
kitchen utensil	0.624	0.358	0.750	0.075
ship	0.233	-0.346	0.618	0.232
snake	0.419	0.002	0.575	0.108
toy	0.480	0.277	0.675	0.111
tree	0.557	0.300	0.781	0.106
vegetable	0.783	0.413	0.892	0.102
vehicle	0.381	0.166	0.699	0.119
weapon	0.421	0.268	0.650	0.082
weather	0.255	0.122	0.357	0.061

Table 11: Prompt design for evaluating typicality effects in models bigger than 30 billion parameters.

Prompt region	Description	Actual prompt
Guidelines	Describe the overall idea of typicality to the model and the task guidelines	Appendix B.3.1
Query	This is the actual fill-in-the-blanks task	The ___ is a "Category-Name"
Options	List of items in a randomized order and separated by a new line	—

B.3.1 Conceptual Understanding - Final Prompt

1483

1484

Typicality effects refer to the influence of the typicality or prototypicality of an object or category on various cognitive processes, including perception, categorization, and memory. The concept of typicality stems from the prototype theory, which suggests that our mental representations of categories are based on prototypes or typical examples.

In the context of perception, typicality effects can influence how we perceive and recognize objects. Objects that are more prototypical or representative of a category are typically perceived more quickly and accurately than atypical objects. For example, when shown a series of pictures of birds, a typical bird like a robin would be recognized faster than a less typical bird like a penguin.

In categorization tasks, typicality effects can influence how we classify objects into categories. Prototypical or highly typical objects are more likely to be assigned to their corresponding category than atypical objects. For instance, when asked to categorize fruits, an apple, being a highly typical fruit, is more likely to be classified as a fruit compared to a less typical fruit like a durian.

Typicality effects also impact memory processes. Prototypical objects are typically better remembered than atypical objects. When asked to recall a list of animals, participants are more likely to remember prototypical animals such as dogs or cats compared to less typical animals like lemurs or armadillos.

Overall, typicality effects demonstrate how the typicality or prototypicality of objects within a category influences our perception, categorization, and memory processes, highlighting the role of prototypes in cognitive functioning.

Based on the typicality effect definitions, give rankings for filling the blank task without any description from the following options.

Make sure to include all the items from the options. Please return items in the following manner:

1. item1
2. item2
3. item3

Also make sure to use the same items as given in the options.

Query:

A ____ is a [Category Name]

Options:

[A]

[B]

1485 B.4 Fluid Reasoning

1486 Humans cannot completely operate without relying on prior experience. The pervasive role of prior
1487 knowledge in shaping cognition is a foundational tenet of the cognitive revolution. However, “Fluid
1488 intelligence” is the ability to solve novel and abstract problems (Raven, 2003). It is a core cognitive
1489 ability, closely related to other domain-general cognitive abilities like working memory, and executive
1490 function, both correlationally (Conway et al., 2002) and in terms of the underlying neural correlates (i.e.,
1491 in the prefrontal cortex) (Burgess et al., 2011). It is distinguished from crystallized intelligence, which is
1492 composed of the domain-specific knowledge and skills one acquires through one’s lifetime (Hartshorne
1493 and Germine, 2015). This distinction is a classic one in psychology (Carroll, 1993).

1494 B.4.1 Scholastic Assessment Test analogy questions

1495 Previous work has shown that fluid reasoning correlates with analogical reasoning (Goswami, 1986; Snow
1496 et al., 1984; Cattell, 1987). AI, ML, and NLP research has focused on analogical reasoning because this
1497 requires many componential abilities: syntactic parsing, semantic understanding, categorization, inductive
1498 reasoning, mathematical reasoning, and so on (Pearson, 2021). Research on the cognitive alignment of
1499 PLMs has focused on performance on the 374 Scholastic Assessment Tests (SAT) analogy questions
1500 by Turney (2005). Despite being broadly used in literature (Turney, 2005; Turney and Pantel, 2010;
1501 Hendrickx et al., 2019; Webb et al., 2023), our pilot experiments show that PLMs like GPT-3.5-Turbo,
1502 GPT-4, and Gemini perform nearly at ceiling on this test, while other open source models perform poorly
1503 on the same test. This hints that the set of questions in the test may be part of the GPT-X/ Gemini training
1504 or tuning data.

1505 **Operationalization:** Each problem is of the form A:B::?, with answer choices containing candidates
1506 for C:D. We evaluate the performance of models in three ways:

- 1507 • Closeness judgment problem: Calculate the cosine similarity between the obtained latent representa-
1508 tions for the member and the category. This requires models where the latent representations are
1509 readily available. These cosine similarities are calculated in different ways:
 - 1510 – 3-cos-add: $\cos(\text{vector}(D), \text{vector}(C) - \text{vector}(A) + \text{vector}(B))$
 - 1511 – 3-cos-mul: $\cos(\text{vector}(D), \text{vector}(B)) * \cos(\text{vector}(D), \text{vector}(C)) / (\cos(\text{vector}(D), \text{vector}(A)) + e)$;
1512 e is a small constant to prevent overflow.
 - 1513 – Concat-cos: $\cos([\text{vector}(A) \parallel \text{vector}(B)], [\text{vector}(C) \parallel \text{vector}(D)])$
- 1514 • Surprisal values: Calculating the summation of probabilities for each token with the as=to relation-
1515 ship; forming the sequence A is to B as C is to D.
- 1516 • Prompting: Prompt the models with the following design: Guidelines, Query, and Options. The
1517 Guideline highlights the task of solving the analogy problem. The Query consists of A:B. The
1518 options are the candidate pairs C:D.

B.4.2 Raven's Progressive Matrices - list of prompts used in experiments

1519

1520

1. "Solve the following Raven Progressive Matrix problem by identifying the pattern in the sequences. Select the correct choice for the missing element."
2. "Identify the correct option to complete the Raven Progressive Matrix. Consider the patterns in numeric and fractional values across the rows to solve the problem."
3. "Solve the Raven Progressive Matrix problem"
4. "Solve the Raven Progressive Matrix problem. Select the correct choice for the missing element in row 3."
5. "Complete the pattern in the Raven Progressive Matrices problem"
6. "Apply abstract reasoning to solve the following Raven Progressive Matrices problem:"
7. "Solve the Raven Progressive Matrices by identifying patterns and drawing analogies. Select the correct choice for the missing element in row 3."
8. "Select the correct choice for row 3, using the patterns and analogies from rows 1 and 2."

row1: (2,0.5,100), (4,0.5,100), (3,0.5,100)

row2: (3,0.7,50), (2,0.7,50), (4,0.7,50)

row3: (4,0.2,70), (3,0.2,70), ?

Choices: (1,0.2,70), (5,0.2,30), (5,0.2,70), (2,0.2,70), (5,0.2,110), (4,0.2,70), (3,0.5,70), (2,0.2,90)"