

Why We Build Local Large Language Models: An Observational Analysis from 35 Japanese and Multilingual LLMs

Anonymous authors

Paper under double-blind review

Abstract

Why do we build local large language models (LLMs)? What should a local LLM learn from the target language? Which abilities can be transferred from other languages? Do language-specific scaling laws exist? To explore these research questions, we evaluated 35 Japanese, English, and multilingual LLMs on 19 evaluation benchmarks for Japanese and English, taking Japanese as a local language. Adopting an observational approach, we analyzed correlations of benchmark scores, and conducted principal component analysis (PCA) to derive *ability factors*. We found that if LLMs perform well in English on tasks like academic subjects, code generation, arithmetic reasoning, commonsense, and reading comprehension, they also perform well on the same tasks in Japanese. This indicates it is not necessary to specifically train on Japanese text to enhance abilities for solving these tasks. In contrast, training on Japanese text improves question-answering tasks about Japanese knowledge and English-Japanese translation, which indicates that abilities for solving these two tasks can be regarded as *Japanese abilities*. Furthermore, we confirmed that the Japanese abilities scale with the computational budget for Japanese text. Taken together, our findings offer generalizable insights into which tasks benefit from local-language data and what we can expect when building local LLMs.

1 Introduction

Major large language models (LLMs) are English-centric (*English LLMs* hereafter), e.g., Meta Llama 3 (Dubey et al., 2024), Mistral (Jiang et al., 2023), and Phi-3 (Abdin et al., 2024), due to the dominance of English on the internet and the global economy, which results in a limited focus on non-English languages. Several companies and research institutes have been actively developing LLMs targeting non-English languages (*local LLMs* hereafter), e.g., Bllossom (Choi et al., 2024), Chinese-LLaMA (Cui et al., 2024) and openCabrita (Larcher et al., 2023), driven by various motivations. These include advancing research and development in multilingual NLP, mitigating security risks associated with relying on a limited number of foreign companies, and promoting responsible artificial intelligence for their community.

However, the advantages of training LLMs on non-English text remain underexplored—particularly regarding the unique skills or knowledge such LLMs might gain compared to English-centric or Multilingual LLMs. On the one hand, LLMs have demonstrated high multilingual abilities, such as arithmetic reasoning (Shi et al., 2023) and machine translation (Briakou et al., 2023), which casts doubt on the advantage of training on non-English text. On the other hand, training on non-English text has been reported to bring stronger cultural and regional knowledge of the target language (Romanou et al., 2025), although there are mixed findings for other tasks such as commonsense reasoning and reading comprehension (Cui et al., 2024; Choi et al., 2024; Larcher et al., 2023). These two perspectives—multilinguality versus language specificity—suggest that the effectiveness of training on non-English text is inherently task dependent. Indeed, demonstrating an advantage of training on non-English text remains not straightforward. Numerous studies have built non-English LLMs from scratch (Holmström et al., 2023) or via continual pre-training (CPT) over English LLMs (Cui et al., 2024; Choi et al., 2024; Larcher et al., 2023), but their task-specific results are often mixed or contradictly, raising doubts about generalizability (§ 2.1). Because LLM

performances depends on several design choices—such as training from scratch or via CPT, which base model is selected for CPT (Tejaswi et al., 2024a), and how the training data is curated (Penedo et al., 2024; Li et al., 2024)—it is difficult to isolate performance gains specifically attributable to training on non-English text. Given its huge impact, thorough investigation and convincing insights into the advantages of local LLMs are valuable.

To explore what unique skills or knowledge may emerge as the natural consequence of the training on non-English text, we adopt an observational approach (Ruan et al., 2024) for Japanese-centric LLMs (*Japanese LLMs* hereafter), leveraging the exceptionally active development in Japan (e.g., Llama 3.1 Swallow¹ and LLM-jp (LLM-jp et al., 2024)) among non-English initiatives. Specifically, we evaluate 35 publicly available Japanese, English, and multilingual LLMs representing a variety of design choices. We also use 19 comprehensive evaluation benchmarks covering knowledge-based QA, academic subjects, reading comprehension, and more, tasked in Japanese and English. These also includes paired Japanese and English benchmarks so that we can compare the task performance across both languages. Our goal is to derive generalizable insights (i.e., insights that are robust to design choices) by conducting a quantitative analysis.

First, to explore multilinguality versus language specificity, we analyzed score correlations across 19 task benchmarks for 35 LLMs, and applied Principal Component Analysis (PCA) to represent the performance in a low-dimensional *ability space* (Ruan et al., 2024). We found that tasks such as academic subjects, code generation, and arithmetic reasoning exhibited strong cross-lingual correlations on their scores and were associated with the same ability factors across languages. This indicates strong multilingual transferability, suggesting that training in English text would also improve performance on these tasks when tested in Japanese. Conversely, tasks such as QA about Japanese cultural knowledge and English-Japanese translation exhibited weak correlations with other tasks and were strongly associated with an independent ability factor, indicating language-specific abilities.

Second, to investigate the language-specific abilities attributed to training on Japanese text, we examined language-specific scaling laws. Specifically, we defined the language-specific computational budget as the product of the number of parameters and training tokens for each language (Hoffmann et al., 2022), and analyzed the log-linear relationship between these budgets and the ability factors obtained by PCA. We found that the English computational budget showed a strong correlation with the general ability factor but a weak correlation with the Japanese-specific ability factor. In contrast, the Japanese computational budget showed a strong correlation with the Japanese ability factor, suggesting that enhancement of Japanese knowledge and English-Japanese translation skills arise from training on Japanese text itself beyond particular design choice. These knowledge and skill scale with the amount of Japanese training text and are difficult to acquire solely from English text.

2 Related Work

2.1 Effects of Training on Non-English Text

There is a growing number of studies examining the impacts of training local LLMs on target language data: Chinese (Zhao et al., 2024; Cui et al., 2024), Turkish (Toraman, 2024), Portuguese (Larcher et al., 2023), Swedish (Holmström et al., 2023), and Finnish (Luukkonen et al., 2023). Some studies consistently reported gains in reading comprehension (Etxaniz et al., 2024b; Fujii et al., 2024; Dou et al., 2024; Joshi et al., 2025; Vo et al., 2024; Larcher et al., 2023), commonsense reasoning (Etxaniz et al., 2024b; Fujii et al., 2024; Phasook et al., 2024; Dou et al., 2024; Joshi et al., 2025; Vo et al., 2024; Choi et al., 2024; Tejaswi et al., 2024b), and local knowledge QA (Etxaniz et al., 2024b; Fujii et al., 2024; Joshi et al., 2025; Etxaniz et al., 2024a). However, following our survey of 15 previous reports on non-English LLMs (see Table 1 in § A), the evidence remains fragmented for two reasons: 1) Sparse coverage of task types: Prior works evaluated only a small set of benchmarks (an average of 2.5). In particular, machine-translation and coding tasks appear in just 2 and 0 out of 15 studies, respectively. 2) Contradictory results: Some studies drew (self-)contradictory conclusions:

¹<https://swallow-llm.github.io/llama3-swallow.en.html>

e.g., for mathematical reasoning, Etxaniz et al. (2024b) reported positive+neutral effects, whereas Pipatanakul et al. (2023) noted negative+neutral effects; for academic subject, both of Phasook et al. (2024) and Dou et al. (2024) documented positive+neutral effects; and, for summarization, Fujii et al. (2024) observed a negative effect, whereas Joshi et al. (2025) and Tejaswi et al. (2024b) found a positive effect.

2.2 Multilinguality vs Language-specificity

Training on non-English corpora sometimes involve using multilingual corpora. Berend (2022) and Chang et al. (2024a) reported that multilingual training does not always improve performance due to the curse of multilinguality (Conneau et al., 2020). Furthermore, English and multilingual LLMs reportedly show strong multilingual abilities on tasks such as arithmetic and commonsense reasoning (Shi et al., 2023) through cross-language generalization (Zhang et al., 2023). These findings suggest that the benefits of training on non-English text might be limited or task-dependent.

2.3 Correlations between Tasks and Ability Factors

Several prior studies have investigated the correlations between different task benchmarks and associated the task performance with a small number of ability factors (Ruan et al., 2024; Ni et al., 2024; Tiong et al., 2024). These studies have reported strong correlations between knowledge-based QA tasks and identified ability factors specific to arithmetic reasoning and code generation. Additionally, Ruan et al. (2024) observed the log-linear relationship between the computational budget and ability factors. However, these discussions are limited to English monolingual settings, leaving cross-language generalization and scaling laws in multilingual contexts, including Japanese and English as in our study, unexplored.

3 Experimental Settings

3.1 Models

To obtain generalizable insights, we evaluated publicly available 35 Japanese, English, and Multilingual LLMs (see Table 2 in Appendix B.1 for the complete list), which represent diverse design choices, including training data, the number of model parameters, and pre-training approach. The evaluated models include: English LLMs (e.g., Llama 3 (Dubey et al., 2024), Mistral (Jiang et al., 2023), and Mixtral (Jiang et al., 2024)); Japanese LLMs continually pre-trained from English base LLMs on 18–175 billion tokens of Japanese text (e.g., Llama 3 Swallow (Fujii et al., 2024) and Llama 3 Youko (Sawada et al., 2024)); Japanese LLMs pre-trained primarily on 130–1,050 billion tokens of Japanese text from scratch (e.g., LLM-jp (LLM-jp et al., 2024) and Sarashina2; and multilingual LLMs pre-trained on multilingual data including Japanese (e.g., C4AI Command-R² and Qwen2 (Yang et al., 2024)). Notably, all the English LLM families that served as base models for the continually pre-trained Japanese LLMs were evaluated as well. We focused on base models and did not evaluate instruction-tuned models to examine the effect of pre-training and avoid the confounding effects of task-oriented instruction tuning.

To estimate the computational budget for each model, we collected data on the number of model parameters and the number of training tokens for Japanese, English, and total across all languages from official sources such as technical reports, press-release documents, and model cards. Refer to Appendix B.3 for details. For a continually pre-trained model, we calculated the total number of training tokens by summing the tokens used in both initial and continual pre-training stages.

²<https://huggingface.co/CohereForAI/c4ai-command-r-v01>

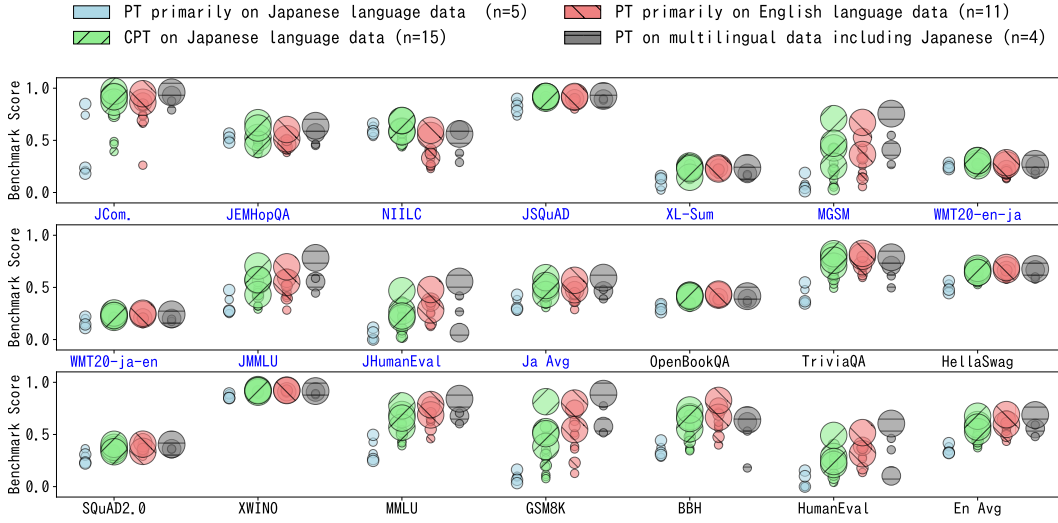


Figure 1: Task performance grouped by primary language of LLMs. Bubble size indicates the number of parameters.

3.2 Evaluation Tasks and Benchmarks

We evaluated the models using 19 evaluation benchmarks in both Japanese and English³, which is listed in Table 3 of Appendix B.2. These tasks were selected from the perspective of cross-lingual benchmarking and comprehensiveness for general-purpose LLMs. The evaluation was conducted using zero-shot or few-shot in-context learning settings depending on tasks. Refer to Appendix B.2 for details.

We employed some Japanese benchmarks corresponding to their English counterparts for cross-lingual benchmarking: code generation (JHumanEval (Sato et al., 2024) vs. HumanEval (Chen et al., 2021)), commonsense (JCommonsenseQA (Kurihara et al., 2022) vs. XWINO (Tikhonov & Ryabinin, 2021) and HellaSwag (Zellers et al., 2019)), arithmetic reasoning (MGSM (Shi et al., 2023) vs. GSM8K (Cobbe et al., 2021)), encyclopedic knowledge-based QA (JEMHopQA (Ishii et al., 2023) and NIILC (Sekine, 2003) vs. TriviaQA (Joshi et al., 2017)), reading comprehension (JSQuAD (Kurihara et al., 2022) vs. SQuAD2 (Rajpurkar et al., 2018)), and academic subjects (JMMLU (Yin et al., 2024) vs. MMLU (Hendrycks et al., 2021)). Notably, MGSM, JMMLU, and JHumanEval are translations of GSM8K, MMLU, and HumanEval, respectively. Cross-lingual correlations between these benchmarks provide insights into the multilinguality and language specificity of each task. It is also worth noting that JEMHopQA and NIILC are developed based on Japanese Wikipedia and include instances that assess knowledge specific to Japanese culture, such as history, geography, notable figures and society, making them suitable for evaluating how much LLMs acquire knowledge about Japan.

For comprehensiveness, inspired by the natural language processing taxonomy (Chang et al., 2024b; Guo et al., 2023) and to capture as many ability factors as possible, we included additional task benchmarks beyond cross-lingual benchmarks. Specifically, we employed Japanese automatic summarization (XL-Sum (Hasan et al., 2021)), machine translation between English and Japanese (WMT20-en-ja and ja-en (Barrault et al., 2020)), English question answering (OpenBookQA (Mihaylov et al., 2018)), and logical reasoning (Big-Bench-Hard (Suzgun et al., 2023)). Because we posit that local LLMs serve as foundational models for the target language, our evaluation focused on fundamental knowledge and skills rather than domain-specific tasks (e.g., question answering in financial or medical

³The evaluation scores for each model will be publicly available on Zenodo with a CC-BY Attribute license upon acceptance (for blind review).

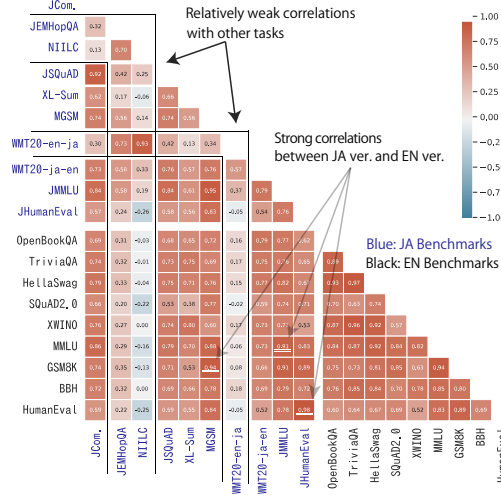


Figure 2: Pearson correlation matrix among task benchmarks ($n = 35$).

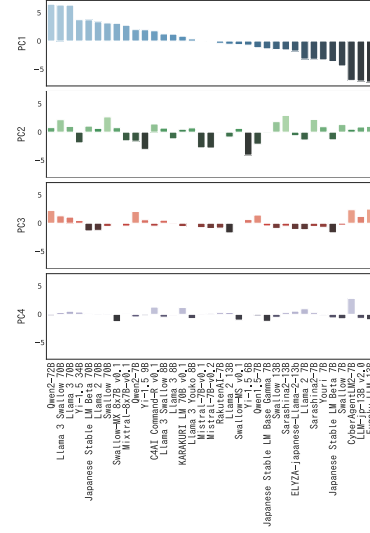


Figure 3: Principal component scores for each LLM.

domains). Furthermore, we excluded safety and bias-related tasks, as these should be addressed in the post-training stage.

3.3 Definition of the Computational Budgets

The Chinchilla scaling laws (Hoffmann et al., 2022) propose an approximation for training FLOPs as $C \approx 6ND$, where C represents the training FLOPs, N is the number of parameters, and D is the number of training tokens. Following this formula, we define ND_l as the computational budget, where D_l is the training tokens for the language l .

3.4 Evaluation Framework and Environment

We evaluated all 35 LLMs on 19 task benchmarks by using a custom implementation⁴ of existing evaluation frameworks such as llm-jp-eval (Han et al., 2024) and the Language Model Evaluation Harness⁵. Refer to Table 4 for the details of implementations used for evaluation. We used NVIDIA A100 GPUs mostly for the evaluations.

4 Experimental Results

Based on the experimental setting explained in the previous section, we obtained a matrix of benchmark scores $X \in \mathbb{R}^{M \times D}$, where M and D are the numbers of LLMs and benchmarks, respectively ($M = 35$ and $D = 19$ in this study) and an element $X_{i,j}$ presents the score of the LLM i on the benchmark j . In this section, we use the benchmark scores matrix X to analyze: 1) the effects of LLM’s primary language on overall performance (§ 4.1), 2) the similarity of benchmarks based on LLM performance (§ 4.2), 3) the ability factors of LLMs (§ 4.3), 4) whether these ability factors align with scaling laws (§ 4.4), and 5) their generalizability to LLMs trained from scratch (§ 4.5).

4.1 Comparison of Benchmark Scores by Pre-trained Languages

Figure 1 presents a bubble chart showing the benchmark score distributions grouped by the primary language of the LLMs: Japanese continually pre-trained (green), Japanese trained

⁴Our implementation has been available on Github, but is hidden here for blind review.

⁵<https://zenodo.org/records/10256836>

PC1 ($r=65, 2$)	0.25	0.13	0.01	0.24	0.21	0.26	0.07	0.24	0.27	0.23	0.25	0.26	0.27	0.22	0.24	0.28	0.26	0.25	0.23
PC2 ($r=15, 4$)	0.06	0.42	0.57	0.13	-0.05	0.05	0.54	0.21	0.08	-0.19	-0.03	-0.03	-0.05	-0.16	-0.01	-0.12	-0.11	-0.04	-0.18
PC3 ($r=7, 0$)	-0.10	0.24	0.03	-0.10	-0.34	0.31	0.03	-0.07	0.17	0.32	-0.21	-0.28	-0.20	0.20	-0.42	0.02	0.27	-0.05	0.33
PC4 ($r=3, 2$)	-0.14	0.12	-0.03	-0.33	-0.54	-0.03	-0.06	0.19	-0.05	-0.23	0.33	0.10	0.22	0.46	0.05	0.01	0.05	0.05	-0.29
	-0.5	-0.0	-0.5																
	Jcom.	JEM-HopQA	NIILC	JSQuAD	XL-Sum	MGSM	WMT20-en-ja	WMT20-en-ja	JMMLU	JHumanEval	OpenBookQA	TriviaQA	HeLa-Swag	SQuAD-2.0	XNLI	MMLU	GSM8K	BBH	HumanEval

Figure 4: Factor Loadings of principal components for each benchmark ($n = 35$; r is the variance explained; blue: Japanese benchmarks; black: English benchmarks).

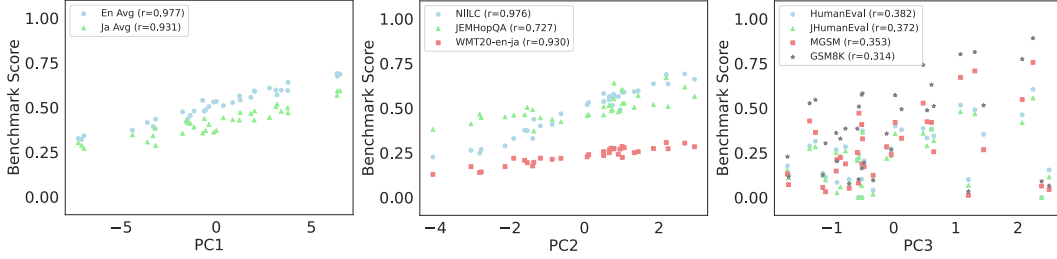


Figure 5: Relationship between principal component scores and raw benchmark scores with significant factor loadings: PC1 vs En/Ja average [left], PC2 vs Japanese knowledge-based QA and En-Ja translation [center], and PC3 vs code-generation and arithmetic reasoning [right] ($n = 35$; r is the pearson correlation coefficient).

195 from scratch (light blue), English (red), and Multilingual (gray). The variable n in each
196 group represents the number of models included.

197 On overall, it is evident that LLMs with larger parameters tend to achieve higher scores in
198 each group. When comparing benchmark scores for smaller models, there is a clear tendency
199 for LLMs continually pre-trained on Japanese text (the green bubbles) to outperform English
200 LLMs (the red bubbles) on Japanese benchmarks (shown in blue) except JHumanEval
201 and MGSM. This indicates the effectiveness of continual pre-training on Japanese text.
202 The advantage is particularly evident in tasks such as Japanese QA (NIILC) and English-
203 Japanese translation (WMT20-en-ja). Refer to Appendix C for detailed discussion. Similarly,
204 Japanese LLMs trained from scratch (the light blue bubbles), despite having relatively few
205 parameters, achieve competitive scores on most Japanese benchmarks, with the exceptions
206 of the arithmetic reasoning (MGSM) and the code-generation (JHumanEval).

207 4.2 Correlation Between Evaluation Benchmarks and Language-Specific Performance

208 To group benchmarks based on the similarities of LLM performance, we computed a Pearson
209 correlation between two benchmarks a and b . More specifically, let the column vectors $X_{:,a}$
210 and $X_{:,b}$ represent the array of two benchmarks a and b , we compute the Pearson correlation
211 $\text{corr}(X_{:,a}, X_{:,b})$. Figure 2 shows the Pearson correlation matrix, revealing two key findings⁶:

212 First, we observed strong cross-lingual correlations on certain tasks: academic subjects
213 (MMLU vs. JMMLU: 0.91), arithmetic reasoning (GSM8K vs. MGSM: 0.94), and code
214 generation (HumanEval vs. JHumanEval: 0.98). In other words, for these tasks, when
215 LLMs perform well on the English benchmarks, they are also likely to perform well on
216 the corresponding Japanese benchmarks. This suggests that multilinguality outweighs
217 language specificity in these tasks, and that LLMs may generalize abilities acquired through
218 training primarily on English text.

219 Second, QA tasks about Japanese knowledge (JEMHopQA, NIILC) and an English-Japanese
220 translation task (WMT20-en-ja) exhibit relatively weak correlations with other tasks re-
221 spectively. In particular, NIILC shows negative correlations with most English tasks, and

⁶We confirmed that using Spearman’s rank correlation produced no significant differences in the findings.

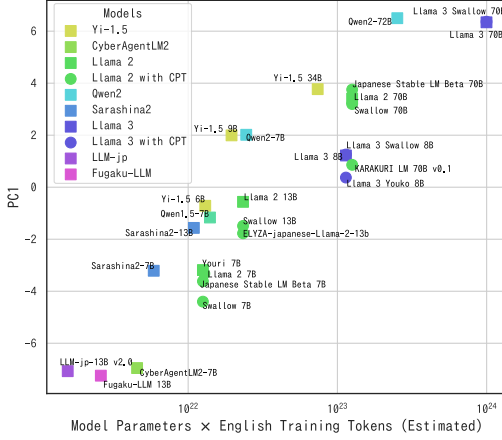


Figure 6: Relationship between the computational budget for English and PC1 scores ($n = 27$).

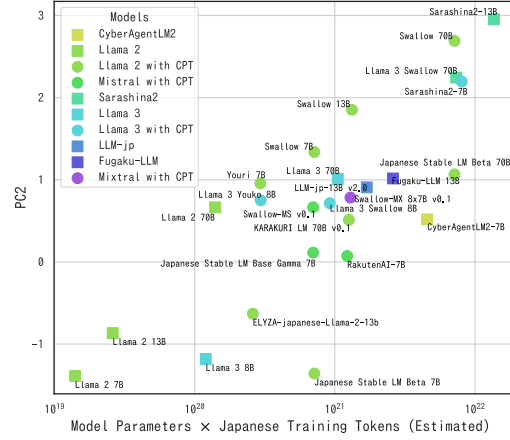


Figure 7: Relationship between the computational budget for Japanese and PC2 scores ($n = 25$).

WMT20-en-ja shows almost no correlations with them. These facts suggest that performance on these tasks may be determined by factors different from those influencing other tasks.

While we observe strong linear correlations between JMMLU, MGSM, and JHumanEval and their English counterparts, given that these are derived from English benchmarks, readers may be concerned that cross-lingual correlations of these benchmarks are overestimated. A straightforward workaround would be to evaluate using random, non-overlapping subsets of instances for each language. Instead of implementing this directly, we approximated the accuracy variation from random splits using the estimated standard error (SE) following [Biderman et al. \(2024\)](#) and confirmed that impact of fluctuation by the SE is negligible on the observed linear trends. For example, MGSM has 250 instances, and the SE for an accuracy of 0.5 is approximately $\sqrt{0.5(1 - 0.5)/250} \approx 0.032$. In contrast, the observed standard deviation of accuracy across LLMs was 0.246, sufficiently larger than the SE.

4.3 Principal Component Analysis (PCA)

We observed benchmark groups from the correlation matrix in the previous subsection. In order to identify ability factors of LLMs, we apply Principal Component Analysis (PCA)⁷ to project the task performance into a low-dimensional ability space.

Formally, we first standardize each column of X to have mean of zero and a standard deviation of one: \hat{X} . Next, we perform eigendecomposition of the correlation matrix as $\hat{X}^\top \hat{X} = U\Lambda U^\top$, where $U = [u_1, u_2, \dots, u_D]$, and $u_j \in \mathbb{R}^D$ is the j -th unit-length eigenvector. We then select the top four principal components (PCs), as their cumulative variance explained (r ; contribution ratio) is 90.8% ($= 65.2\% + 15.4\% + 7.0\% + 3.2\%$ from PC1 to PC4). We define the eigenvectors corresponding to PC1 to PC4, $U_4 = [u_1, u_2, u_3, u_4] \in \mathbb{R}^{D \times 4}$ as the factor loadings and compute corresponding PC scores as $S_4 = \hat{X}U_4$. Given that U is an orthonormal matrix and the total variance explained by PC1–PC4 is about 90%, the original matrix can be approximated as the product of PC scores and factor loadings: $\hat{X} \approx S_4U_4^\top$.

In this way, we decompose standardized benchmark scores \hat{X} into the product of LLM-specific principal component scores (ability factors) $S_4 \in \mathbb{R}^{M \times 4}$ in Figure 3 and benchmark-specific factor loadings $U_4 \in \mathbb{R}^{D \times 4}$ in Figure 4, which represent the associations between the ability factors and task performances⁸.

⁷We used the `sklearn.decomposition.PCA()` method from the `scikit-learn` package.

⁸Since the signs and magnitudes of the PC scores and factor loadings are arbitrary, we adjusted the signs for ease of interpretation and normalized the factor loading vectors to have an L_2 norm of 1.

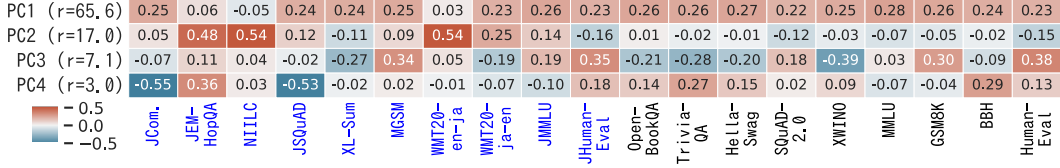


Figure 8: Factor loadings of principal components for each benchmark ($n = 20$: only with models trained from scratch; r is the variance explained; blue: Japanese benchmarks; black: English benchmarks).

The first principal component (PC1) has relatively uniform factor loadings. As shown in Figure 5 left, LLMs with higher PC1 scores tend to have higher average benchmark scores in both English and Japanese, suggesting that PC1 represents a general ability factor. It represents the average performance across most benchmark scores, including commonsense and reading comprehension in Japanese. This indicates that, unlike prior studies (§ 2.1), training on English text is effective and that Japanese-specific training is not necessarily for improving these abilities.

The second principal component (PC2) shows concentrated factor loadings on JEMHopQA, NIILC, and WMT20-en-ja, and relatively small factor loadings on JCommonsenseQA and JSQuAD, indicating the abilities of (encyclopedic) knowledge about Japan and English-Japanese translation. In fact, Figure 3 shows that LLMs pre-trained on Japanese text, such as Swallow and Sarashina2 families, have high PC2 scores, which will be analyzed in detail in § 4.4. Additionally, as shown in Figure 5 center, the higher PC2, the higher benchmark scores on those tasks. For instance, the margin of NIILC accuracy between LLMs with the lowest and highest PC2 scores is approximately 40 points. Considering that PC1 has relatively low factor loadings for these benchmarks, PC2 represents Japanese-specific abilities, such as QA about Japanese knowledge and English-Japanese translation. Given that PC2 strongly associates with Japanese knowledge-based QA tasks, this aligns with previous work (Romanou et al., 2025), which found that multilingual LLMs struggle with cultural questions, especially in languages not included in the pre-training data.

The third principal component (PC3) shows concentrated factor loadings on MGSM, GSM8K, JHumanEval, and HumanEval, representing abilities of multilingualism, language-agnostic arithmetic reasoning, and code generation. As shown in Figure 5 right, there is a moderate trend suggesting that higher PC3 score are associated with higher benchmark scores on code-generation and arithmetic-reasoning.

Finally, the fourth principal component (PC4) shows positive factor loadings for some English benchmarks. However, strong English LLMs, such as Llama-3-70B, do not show higher PC4 scores compared to Japanese LLMs like CyberAgentLM2-7B. In addition, given that the variance explained by PC4 is only 3.2%, PC4 is likely to correspond to residuals that are difficult to interpret in a way tied to specific benchmarks or abilities.

4.4 Scaling Laws between Ability Factors and Computational Budget

In § 4.3, we made two key observations: 1) PC2 represents Japanese ability while PC1 represents a general ability; 2) LLMs pre-trained on Japanese text tend to have higher PC2 scores. Based on these observations, we explore the language-specific scaling laws by examining the log-linear relationship between the computational budgets (§ 3.3) and principal components, which are expected to represent different abilities.

Figure 6 shows the scatter plot with the English computational budget (log scale) and PC1. It reveals that the general ability (PC1) scales with the English computational budget (Pearson’s $\rho = 0.916$)⁹

⁹The correlation with the logarithm of the total computational budget was slightly higher ($\rho = 0.938$). Still, given the weak correlation with the Japanese computational budget, we concluded that it scales more with the English computational budget.

Figure 7 shows the scatter plot with the Japanese computational budget (log scale) and PC2. We can see that the Japanese ability (PC2) moderately scales with the Japanese computational budget ($\rho = 0.779$). We also confirmed that the correlation between PC2 and the English or total computational budget is much weaker ($\rho = 0.164$ and 0.186 , respectively). These findings indicate that PC2 and associated Japanese task performances scale with an increase in Japanese training tokens, thereby supporting our claim in § 4.3 that “PC2 represents Japanese ability.” Furthermore, we argue that the source of Japanese ability lies in the computational budget allocated to Japanese texts.

4.5 PCA for LLMs Trained from Scratch

To verify that our findings are not heavily influenced by the pre-training method, we repeated the analysis after excluding continually pre-trained Japanese LLMs, retaining only 20 LLMs trained from scratch. Figure 8 shows the factor loadings of PCs extracted from the performance of these 20 LLMs, revealing ability factors similar to those identified in the original analysis (§ 4.3). We omit the results of relationships between computational budgets and English and Japanese abilities, but observed the consistent correlations with Figures 6 and 7 (see Figures 13 and 14 in Appendix D.2).

5 Conclusion and Future Work

In this paper, we performed the most comprehensive evaluation to date, testing 35 Japanese, English, and Multilingual LLMs on 19 task benchmarks that assess the abilities in both Japanese and English. This breadth of coverage is one of the key novelties of our study and enables us to extract more generalizable insights than prior work. We then analyzed the cross-task and cross-lingual correlations of benchmark scores, mapped the performance in a low-dimensional ability space, and explored the relationship between ability factors and computational budgets for English and Japanese. The correlation analysis showed strong multilingual abilities in academic subjects, code generation, and arithmetic reasoning tasks. This suggests that, in order to enhance the abilities of these tasks, there is no strong motivation for using Japanese training data.

The low-dimensional factor analysis using PCA identified three ability factors. PC1 represents the general ability and affects nearly all tasks except for QA about Japanese knowledge and English-Japanese translation. PC1 follows a scaling law with the computational budget for English. Complementing PC1, PC2 represents the ability for QA about Japanese knowledge and English-Japanese translation. Interestingly, PC2 follows a scaling law with the computational budget for Japanese data. Although PC3 represents multilingual abilities in arithmetic reasoning and code generation, we have not reached the point of identifying a scaling law that it follows.

From these analyses, we concluded that the advantage of building local LLMs by training on Japanese text is particularly evident in acquiring local knowledge written in Japanese and enhancing the ability to translate from English. This conclusion is likely to characterize Japanese LLMs. Our study is the first broad, unified evaluation across dozens of LLMs and an extensive benchmark suite to reveal which tasks do and do not benefit from target-language training.

We consider two directions as future work. First, we plan to extend the analysis with more LLMs and evaluation tasks to discover additional insights. This includes using LLMs with unique designs, for example, Phi family (Li et al., 2023; Abidin et al., 2024), which were trained on synthetic text. We also want to add evaluation tasks such as Japanese logical reasoning and standardized admission exams. The second direction is to extend our analysis and findings to other languages. We believe that the conclusion of this paper can be generalized to: the advantage of building local LLMs by training in a language is acquiring local knowledge written in the language and enhancing the ability to translate from English to the language. This direction is nontrivial because conducting LLM experiments properly requires a deep understanding of the target languages and cultures. We hope this paper serves as a catalyst for the development and analysis of non-English LLMs.

Acknowledgments

(Removed for blind review)

Ethics Statement

This study does not evaluate the safety aspects of LLMs, such as harmlessness or honesty (Askell et al., 2021), which are considered to be largely shaped by pre-training data. The same applies when developing local LLMs — they are likely to absorb social group-specific biases (Yanaka et al., 2024), stereotypes, and racism. Consequently, there is a concern that we may be overlooking an inconvenient side effect: it might be unavoidable for local LLMs to reinforce social biases specific to the target language.

Reproducibility Statement

We prioritized reproducibility in our work. As described in § 3.4, all 35 LLMs (Table 2), 19 benchmarks (Table 3), and evaluation frameworks (Table 4) used in our study are publicly accessible. Additionally, evaluation scores for all LLMs, along with models’ metadata—including training data, the number of model parameters, the number of training tokens, and pre-training approach—are available (§ 3.2) to facilitate the reproduction of statistical analyses. Please note that our unified evaluation framework and the results are withheld here to preserve anonymity during the blind review process. For reference, our experiments were primarily conducted on NVIDIA A100 GPUs.

Broader Impacts

We believe our findings will contribute to the development of non-English LLMs. Moreover, this could foster a society in which every country has access to LLMs specialized in its own language and knowledge, thereby reducing the digital divide.

References

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, and Harkirat Behl et al. Phi-3 technical report: A highly capable language model locally on your phone. arXiv:2404.14219, 2024.
01. AI, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, and Jianqun Chen et al. Yi: Open foundation models by 01.AI. arXiv:2403.04652, 2024.
- Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, and Nova DasSarma et al. A general language assistant as a laboratory for alignment. arXiv:2112.00861, 2021.
- Loïc Barrault, Magdalena Biesialska, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, and Eric et al. Joanis. Findings of the 2020 conference on machine translation. In *Proceedings of the Fifth Conference on Machine Translation*, pp. 1–55, 2020. URL <https://aclanthology.org/2020.wmt-1.1>.
- Loubna Ben Allal, Niklas Muennighoff, Logesh Kumar Umapathi, Ben Lipkin, and Leandro von Werra. A framework for the evaluation of code generation models. <https://github.com/bigcode-project/bigcode-evaluation-harness>, 2022.
- Gábor Berend. Combating the curse of multilinguality in cross-lingual WSD by aligning sparse contextualized word representations. In *Proceedings of the 2022 Conference of the*

- 385 *North American Chapter of the Association for Computational Linguistics: Human Language*
 386 *Technologies*, pp. 2459–2471, 2022. URL [https://aclanthology.org/2022.naacl-main.](https://aclanthology.org/2022.naacl-main.176)
 387 [176](https://aclanthology.org/2022.naacl-main.176).
- 388 Stella Biderman, Hailey Schoelkopf, Lintang Sutawika, Leo Gao, Jonathan Tow, Baber
 389 Abbasi, Alham Fikri Aji, Pawan Sasanka Ammanamanchi, Sidney Black, Jordan Clive,
 390 Anthony DiPofi, Julien Etzaniz, Benjamin Fattori, Jessica Zosa Forde, Charles Foster, Jeffrey
 391 Hsu, Mimansa Jaiswal, Wilson Y. Lee, Haonan Li, Charles Lovering, Niklas Muennighoff,
 392 Ellie Pavlick, Jason Phang, Aviya Skowron, Samson Tan, Xiangru Tang, Kevin A. Wang,
 393 Genta Indra Winata, François Yvon, and Andy Zou. Lessons from the trenches on
 394 reproducible evaluation of language models. arXiv:2112.00861, 2024. URL [https://doi.](https://doi.org/10.48550/arXiv.2405.14782)
 395 [org/10.48550/arXiv.2405.14782](https://doi.org/10.48550/arXiv.2405.14782).
- 396 Eleftheria Briakou, Colin Cherry, and George Foster. Searching for needles in a haystack: On
 397 the role of incidental bilingualism in PaLM’s translation capability. In *Proceedings of the*
 398 *61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,
 399 pp. 9432–9452, 2023. URL <https://aclanthology.org/2023.acl-long.524>.
- 400 Tyler A. Chang, Catherine Arnett, Zhuowen Tu, and Ben Bergen. When is multilinguality a
 401 curse? language modeling for 250 high- and low-resource languages. In *Proceedings of the*
 402 *2024 Conference on Empirical Methods in Natural Language Processing*, pp. 4074–4096, 2024a.
 403 URL <https://aclanthology.org/2024.emnlp-main.236/>.
- 404 Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen,
 405 Xiaoyuan Yi, Cunxiang Wang, and Yidong et al. Wang. A survey on evaluation of large
 406 language models. *Association for Computing Machinery Transactions on Intelligent Systems*
 407 *and Technology*, 15(3):1–45, 2024b. URL <https://doi.org/10.1145/3641289>.
- 408 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto,
 409 Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, and Greg Brockman et al.
 410 Evaluating large language models trained on code. arXiv:2107.03374, 2021.
- 411 ChangSu Choi, Yongbin Jeong, Seoyoon Park, Inho Won, HyeonSeok Lim, SangMin Kim,
 412 Yejee Kang, Chanhyuk Yoon, Jaewan Park, and Yiseul et al. Lee. Optimizing language
 413 augmentation for multilingual large language models: A case study on Korean. In
 414 *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language*
 415 *Resources and Evaluation*, pp. 12514–12526, 2024. URL [https://aclanthology.org/2024.](https://aclanthology.org/2024.lrec-main.1095)
 416 [lrec-main.1095](https://aclanthology.org/2024.lrec-main.1095).
- 417 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
 418 Matthias Plappert, Jerry Tworek, Jacob Hilton, and Reiichiro Nakano et al. Training
 419 verifiers to solve math word problems. arXiv:2110.14168, 2021.
- 420 Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume
 421 Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin
 422 Stoyanov. Unsupervised cross-lingual representation learning at scale. In *Proceedings of*
 423 *the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 8440–8451, 2020.
 424 URL <https://aclanthology.org/2020.acl-main.747>.
- 425 Yiming Cui, Ziqing Yang, and Xin Yao. Efficient and effective text encoding for Chinese
 426 LLaMA and Alpaca. arXiv:2304.08177, 2024.
- 427 Longxu Dou, Qian Liu, Guangtao Zeng, Jia Guo, Jiahui Zhou, Xin Mao, Ziqi Jin, Wei Lu,
 428 and Min Lin. Sailor: Open language models for south-East Asia. In *Proceedings of the 2024*
 429 *Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp.
 430 424–435, 2024. URL <https://aclanthology.org/2024.emnlp-demo.45/>.
- 431 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle,
 432 Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, and Angela Fan et al. The Llama
 433 3 herd of models. arXiv:2407.21783, 2024.

- 434 Julen Etxaniz, Gorka Azkune, Aitor Soroa, Oier Lopez de Lacalle, and Mikel
435 Artetxe. Bertaqa: How much do language models know about local culture? In *Advances in Neural Information Processing Systems*, volume 37, pp. 34077–
436 34097, 2024a. URL [https://proceedings.neurips.cc/paper_files/paper/2024/file/](https://proceedings.neurips.cc/paper_files/paper/2024/file/3bb42f6bb1b1ab6809afd6c90865b087-Paper-Datasets_and_Benchmarks_Track.pdf)
437 [3bb42f6bb1b1ab6809afd6c90865b087-Paper-Datasets_and_Benchmarks_Track.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/3bb42f6bb1b1ab6809afd6c90865b087-Paper-Datasets_and_Benchmarks_Track.pdf).
438
- 439 Julen Etxaniz, Oscar Sainz, Naiara Miguel, Itziar Aldabe, German Rigau, Eneko Agirre,
440 Aitor Ormazabal, Mikel Artetxe, and Aitor Soroa. Latxa: An open language model and
441 evaluation suite for Basque. In *Proceedings of the 62nd Annual Meeting of the Association*
442 *for Computational Linguistics (Volume 1: Long Papers)*, pp. 14952–14972, 2024b. URL [https://](https://aclanthology.org/2024.acl-long.799/)
443 aclanthology.org/2024.acl-long.799/.
- 444 Kazuki Fujii, Taishi Nakamura, Mengsay Loem, Hiroki Iida, Masanari Ohi, Kakeru Hattori,
445 Hirai Shota, Sakae Mizuki, Rio Yokota, and Naoaki Okazaki. Continual pre-training for
446 cross-lingual LLM adaptation: Enhancing Japanese language capabilities. In *Proceedings*
447 *of the 1st Conference on Language Modeling*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=TQdd1VhWbe)
448 [TQdd1VhWbe](https://openreview.net/forum?id=TQdd1VhWbe).
- 449 Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Supryadi, Linhao Yu, Yan Liu,
450 Jiaxuan Li, and Bojian Xiong et al. Evaluating large language models: A comprehensive
451 survey. arXiv:2310.19736, 2023.
- 452 Namgi Han, Nobuhiro Ueda, Masatoshi Otake, Satoshi Katsumata, Keisuke Kamata, Hi-
453 rokazu Kiyomaru, Takashi Kodama, Saku Sugahara, Bowen Chen, and Hiroshi Matsuda
454 et al. llm-jp-eval: Automatic evaluation tool for Japanese large language models [llm-jp-
455 eval: 日本語大規模言語モデルの自動評価ツール] (in Japanese). In *Proceedings of the 30th*
456 *Annual Meeting of the Association for Natural Language Processing*, pp. 2085–2089, 2024.
- 457 Tahmid Hasan, Abhik Bhattacharjee, Md. Saiful Islam, Kazi Mubasshir, Yuan-Fang Li, Yong-
458 Bin Kang, M. Sohel Rahman, and Rifat Shahriyar. XL-sum: Large-scale multilingual
459 abstractive summarization for 44 languages. In *Findings of the Association for Computational*
460 *Linguistics*, pp. 4693–4703, 2021. URL [https://aclanthology.org/2021.findings-acl.](https://aclanthology.org/2021.findings-acl.413)
461 [413](https://aclanthology.org/2021.findings-acl.413).
- 462 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
463 Jacob Steinhardt. Measuring massive multitask language understanding. In *Proceedings*
464 *of the Ninth International Conference on Learning Representations*, 2021. URL [https://](https://openreview.net/forum?id=d7KBjmI3GmQ)
465 openreview.net/forum?id=d7KBjmI3GmQ.
- 466 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai,
467 Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, and Aidan
468 et al. Clark. An empirical analysis of compute-optimal large language model training.
469 In *Proceedings of the 36th International Conference on Neural Information Processing Systems*,
470 pp. 1–15, 2022. URL [https://proceedings.neurips.cc/paper_files/paper/2022/file/](https://proceedings.neurips.cc/paper_files/paper/2022/file/c1e2faff6f588870935f114ebe04a3e5-Paper-Conference.pdf)
471 [c1e2faff6f588870935f114ebe04a3e5-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/c1e2faff6f588870935f114ebe04a3e5-Paper-Conference.pdf).
- 472 Oskar Holmström, Jenny Kunz, and Marco Kuhlmann. Bridging the resource gap: Exploring
473 the efficacy of English and multilingual LLMs for Swedish. In *Proceedings of the Second*
474 *Workshop on Resources and Representations for Under-Resourced Languages and Domains*, pp.
475 92–110, 2023. URL <https://aclanthology.org/2023.resourceful-1.13>.
- 476 Ai Ishii, Naoya Inoue, and Satoshi Sekine. Construction of a Japanese multi-hop QA dataset
477 for QA systems capable of explaining the rationale [根拠を説明可能な質問応答システ
478 ムのための日本語マルチホップQAデータセット構築] (in Japanese). In *the 29th Annual*
479 *Meeting of Japanese Association for Natural Language Processing (NLP2023)*, pp. 2088–2093,
480 2023.
- 481 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh
482 Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, and
483 Lucile Saulnier et al. Mistral 7B. arXiv:2310.06825, 2023.

- 484 Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary,
485 Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, and
486 Florian Bressand et al. Mixtral of experts. arXiv:2401.04088, 2024.
- 487 Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. TriviaQA: A large scale
488 distantly supervised challenge dataset for reading comprehension. In *Proceedings of the*
489 *55th Annual Meeting of the Association for Computational Linguistics*, pp. 1601–1611, 2017.
490 URL <https://aclanthology.org/P17-1147>.
- 491 Raviraj Joshi, Kanishk Singla, Anusha Kamath, Raunak Kalani, Rakesh Paul, Utkarsh Vaidya,
492 Sanjay Singh Chauhan, Niranjana Wartikar, and Eileen Long. Adapting multilingual LLMs
493 to low-resource languages using continued pre-training and synthetic corpus: A case
494 study for Hindi LLMs. In *Proceedings of the First Workshop on Natural Language Processing*
495 *for Indo-Aryan and Dravidian Languages*, pp. 50–57, 2025. URL [https://aclanthology.org/](https://aclanthology.org/2025.indonlp-1.6/)
496 [2025.indonlp-1.6/](https://aclanthology.org/2025.indonlp-1.6/).
- 497 KARAKURI Inc. KARAKURI LM 70B v0.1. Hugging Face: karakuri-ai/karakuri-lm-70b-
498 v0.1, 2024.
- 499 Kentaro Kurihara, Daisuke Kawahara, and Tomohide Shibata. JGLUE: Japanese general
500 language understanding evaluation. In *Proceedings of the Thirteenth Language Resources*
501 *and Evaluation Conference*, pp. 2957–2966, 2022. URL [https://aclanthology.org/2022.](https://aclanthology.org/2022.lrec-1.317)
502 [lrec-1.317](https://aclanthology.org/2022.lrec-1.317).
- 503 Celio Larcher, Marcos Piau, Paulo Finardi, Pedro Gengo, Piero Esposito, and Vinicius Caridá.
504 Cabrita: closing the gap for foreign languages. arXiv:2308.11878, 2023.
- 505 Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Gadre, Hritik
506 Bansal, Etash Guha, Sedrick Keh, Kushal Arora, Saurabh Garg, Rui Xin, Niklas Muen-
507 nighoff, Reinhard Heckel, Jean Mercat, Mayee Chen, Suchin Gururangan, Mitchell
508 Wortsman, Alon Albalak, Yonatan Bitton, Marianna Nezhurina, Amro Abbas, Cheng-
509 Yu Hsieh, Dhruva Ghosh, Josh Gardner, Maciej Kilian, Hanlin Zhang, Rulin Shao,
510 Sarah Pratt, Sunny Sanyal, Gabriel Ilharco, Giannis Daras, Kalyani Marathe, Aaron
511 Gokaslan, Jieyu Zhang, Khyathi Chandu, Thao Nguyen, Igor Vasiljevic, Sham Kakade,
512 Shuran Song, Sujay Sanghavi, Fartash Faghri, Sewoong Oh, Luke Zettlemoyer, Kyle
513 Lo, Alaaeldin El-Nouby, Hadi Pouransari, Alexander Toshev, Stephanie Wang, Dirk
514 Groeneveld, Luca Soldaini, Pang Wei Koh, Jenia Jitsev, Thomas Kollar, Alexan-
515 dros G. Dimakis, Yair Carmon, Achal Dave, Ludwig Schmidt, and Vaishaal Shankar.
516 Datacomp-lm: In search of the next generation of training sets for language mod-
517 els. In *Advances in Neural Information Processing Systems*, volume 37, pp. 14200–
518 14282, 2024. URL [https://proceedings.neurips.cc/paper_files/paper/2024/file/](https://proceedings.neurips.cc/paper_files/paper/2024/file/19e4ea30dded58259665db375885e412-Paper-Datasets.and.Benchmarks.Track.pdf)
519 [19e4ea30dded58259665db375885e412-Paper-Datasets.and.Benchmarks.Track.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/19e4ea30dded58259665db375885e412-Paper-Datasets.and.Benchmarks.Track.pdf).
- 520 Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat
521 Lee. Textbooks are all you need II: phi-1.5 technical report. arXiv:2309.05463, 2023.
- 522 LLM-jp, Akiko Aizawa, Eiji Aramaki, Bowen Chen, Fei Cheng, Hiroyuki Deguchi, Rintaro
523 Enomoto, Kazuki Fujii, Kensuke Fukumoto, and Takuya Fukushima et al. LLM-jp: A
524 cross-organizational project for the research and development of fully open Japanese
525 LLMs. arXiv:2407.03963, 2024.
- 526 Risto Luukkainen, Ville Komulainen, Jouni Luoma, Anni Eskelinen, Jenna Kanerva, Hanna-
527 Mari Kupari, Filip Ginter, Veronika Laippala, Niklas Muennighoff, Aleksandra Piktus,
528 and et al. FinGPT: Large generative models for a small language. In *Proceedings of the*
529 *2023 Conference on Empirical Methods in Natural Language Processing*, pp. 2710–2726, 2023.
530 URL <https://aclanthology.org/2023.emnlp-main.164/>.
- 531 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor
532 conduct electricity? a new dataset for open book question answering. In *Proceedings of the*
533 *2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2381–2391, 2018.
534 URL <https://aclanthology.org/D18-1260>.

- 535 Jinjie Ni, Fuzhao Xue, Xiang Yue, Yuntian Deng, Mahir Shah, Kabir Jain, Graham Neubig,
536 and Yang You. Mixeval: Deriving wisdom of the crowd from LLM benchmark mixtures.
537 arXiv:2406.06565, 2024.
- 538 Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Mar-
539 garet Mitchell, Colin Raffel, Leandro Von Werra, and Thomas Wolf. The
540 fineweb datasets: Decanting the web for the finest text data at scale. In *Ad-
541 vances in Neural Information Processing Systems*, volume 37, pp. 30811–30849,
542 2024. URL [https://proceedings.neurips.cc/paper_files/paper/2024/file/
543 370df50ccfd8bde18f8f9c2d9151bda-Paper-Datasets_and_Benchmarks_Track.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/370df50ccfd8bde18f8f9c2d9151bda-Paper-Datasets_and_Benchmarks_Track.pdf).
- 544 Pakawat Phasook, Jessada Pranee, Chananyu Limcharoen, Kittisak Sukhantharat, Anon
545 Saeoueng, Kun Kerdthaisong, Chaianun Damrongrat, and Sarawoot Kongyoung.
546 Thaibkd: Effective of continual pre-training llm in thai language based on knowledge
547 dataset. In *2024 19th International Joint Symposium on Artificial Intelligence and Natural
548 Language Processing (iSAI-NLP)*, pp. 1–7, 2024.
- 549 Kunat Pipatanakul, Phatrasek Jirabovonvisut, Potsawee Manakul, Sittipong Sripaisarn-
550 mongkol, Ruangsak Patomwong, Pathomporn Chokchainant, and Kasima Tharnpipitchai.
551 Typhoon: Thai large language models. arXiv:2312.13951, 2023.
- 552 Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don’t know: Unanswerable
553 questions for SQuAD. In *Proceedings of the 56th Annual Meeting of the Association for
554 Computational Linguistics*, pp. 784–789, 2018. URL <https://aclanthology.org/P18-2124>.
- 555 RakutenGroup, Aaron Levine, Connie Huang, Chenguang Wang, Eduardo Batista, Ewa
556 Szymanska, Hongyi Ding, Hou Wei Chou, Jean-François Pessiot, and Johaness Effendi
557 et al. RakutenAI-7B: Extending large language models for Japanese. arXiv:2403.15484,
558 2024.
- 559 Angelika Romanou, Negar Foroutan, Anna Sotnikova, Sree Harsha Nelaturu, Shivalika
560 Singh, Rishabh Maheshwary, Micol Altomare, Zeming Chen, Mohamed A. Haggag,
561 Snegha A, and et al. INCLUDE: Evaluating multilingual language understanding with
562 regional knowledge. In *Proceedings of the Thirteenth International Conference on Learning
563 Representations*, 2025. URL <https://openreview.net/forum?id=k3gCieTXeY>.
- 564 Yangjun Ruan, Chris J. Maddison, and Tatsunori Hashimoto. Observational scaling laws
565 and the predictability of language model performance. arXiv:2405.10938, 2024.
- 566 Akira Sasaki, Masato Hirakawa, Shintaro Horie, Tomoaki Nakamura, Sam Passaglia,
567 and Daisuke Oba. ELYZA-japanese-Llama-2-13b. [https://huggingface.co/elyza/
568 ELYZA-japanese-Llama-2-13b](https://huggingface.co/elyza/ELYZA-japanese-Llama-2-13b), 2023.
- 569 Yui Sato, Shiho Takano, Teruno Kajiura, and Kimiro Kuramitsu. Do large language models
570 transfer knowledge across languages through additional Japanese training? [llmは日本語
571 追加学習により言語間知識転移を起こすのか?] (in Japanese). In *Proceedings of the 30th
572 Annual Meeting of the Association for Natural Language Processing*, pp. 2897–2901, 2024.
- 573 Kei Sawada, Tianyu Zhao, Makoto Shing, Kentaro Mitsui, Akio Kaga, Yukiya Hono, Toshiaki
574 Wakatsuki, and Koh Mitsuda. Release of pre-trained models for the Japanese language. In
575 *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language
576 Resources and Evaluation (LREC-COLING 2024)*, pp. 13898–13905, 2024. URL [https://
577 aclanthology.org/2024.lrec-main.1213](https://aclanthology.org/2024.lrec-main.1213).
- 578 Satoshi Sekine. Development of a question answering system targeting encyclopedias [百科
579 事典を対象とした質問応答システムの開発] (in Japanese). In *Proceedings of the 9th Annual
580 Meeting of the Association for Natural Language Processing*, pp. 637–640, 2003.
- 581 Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi,
582 Hyung Won Chung, Yi Tay, Sebastian Ruder, and Denny Zhou et al. Language models
583 are multilingual chain-of-thought reasoners. In *Proceedings of the Eleventh International
584 Conference on Learning Representations*, 2023. URL [https://openreview.net/forum?id=
585 fR3wGck-IXp](https://openreview.net/forum?id=fR3wGck-IXp).

- 586 Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won
587 Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, and Denny et al. Zhou. Challenging BIG-
588 bench tasks and whether chain-of-thought can solve them. In *Findings of the Association for*
589 *Computational Linguistics*, pp. 13003–13051, 2023. URL [https://aclanthology.org/2023.](https://aclanthology.org/2023.findings-acl.824)
590 [findings-acl.824](https://aclanthology.org/2023.findings-acl.824).
- 591 Qwen Team. Introducing Qwen1.5, 2024. URL <https://qwenlm.github.io/blog/qwen1.5/>.
- 592 Atula Tejaswi, Nilesh Gupta, and Eunsol Choi. Exploring design choices for building
593 language-specific llms. In *Findings of the Association for Computational Linguistics*, pp.
594 10485–10500, 2024a. URL <https://aclanthology.org/2024.findings-emnlp.614>.
- 595 Atula Tejaswi, Nilesh Gupta, and Eunsol Choi. Exploring design choices for building
596 language-specific LLMs. arXiv:2406.14670, 2024b.
- 597 Alexey Tikhonov and Max Ryabinin. It’s all in the heads: Using attention heads as a baseline
598 for cross-lingual transfer in commonsense reasoning. In *Findings of the Association for*
599 *Computational Linguistics*, pp. 3534–3546, 2021. URL [https://aclanthology.org/2021.](https://aclanthology.org/2021.findings-acl.310)
600 [findings-acl.310](https://aclanthology.org/2021.findings-acl.310).
- 601 Anthony Tiong, Junqi Zhao, Boyang Li, Junnan Li, Steven Hoi, and Caiming Xiong. What
602 are we measuring when we evaluate large vision-language models? an analysis of latent
603 factors and biases. In *Proceedings of the 2024 Conference of the North American Chapter of*
604 *the Association for Computational Linguistics: Human Language Technologies*, pp. 3427–3454,
605 2024.
- 606 Cagri Toraman. Adapting open-source generative large language models for low-resource
607 languages: A case study for Turkish. In *Proceedings of the Fourth Workshop on Multilingual*
608 *Representation Learning (MRL 2024)*, pp. 30–44, 2024. URL [https://aclanthology.org/](https://aclanthology.org/2024.mrl-1.3/)
609 [2024.mrl-1.3/](https://aclanthology.org/2024.mrl-1.3/).
- 610 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
611 Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, and Shruti Bhosale et al. Llama 2:
612 Open foundation and fine-tuned chat models. arXiv:2307.09288, 2023.
- 613 Anh-Dung Vo, Minseong Jung, Wonbeen Lee, and Daewoo Choi. Redwhale: An adapted
614 korean llm through efficient continual pretraining. arXiv:2408.11294, 2024.
- 615 Hitomi Yanaka, Namgi Han, Ryoma Kumon, Jie Lu, Masashi Takeshita, Ryo Sekizawa,
616 Taisei Kato, and Hiromi Arai. Analyzing social biases in Japanese large language models.
617 arXiv:2406.02050, 2024.
- 618 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Cheng-
619 peng Li, Chengyuan Li, Dayiheng Liu, and Fei Huang et al. Qwen2 technical report.
620 arXiv:2407.10671, 2024.
- 621 Ziqi Yin, Hao Wang, Kaito Horio, Daisuke Kawahara, and Satoshi Sekine. Investigating the
622 relationship between prompt politeness and large language model performance [プロン
623 プトの丁寧さと大規模言語モデルの性能の関係検証] (in Japanese). In *Proceedings of the*
624 *30th Annual Meeting of the Association for Natural Language Processing*, pp. 1803–1808, 2024.
- 625 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can
626 a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the*
627 *Association for Computational Linguistics*, pp. 4791–4800, 2019. URL [https://aclanthology.](https://aclanthology.org/P19-1472)
628 [org/P19-1472](https://aclanthology.org/P19-1472).
- 629 Xiang Zhang, Senyu Li, Bradley Hauer, Ning Shi, and Grzegorz Kondrak. Don’t trust
630 ChatGPT when your question is not in English: A study of multilingual abilities and
631 types of LLMs. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language*
632 *Processing*, pp. 7915–7927, 2023. URL <https://aclanthology.org/2023.emnlp-main.491>.
- 633 Jun Zhao, Zhihao Zhang, Luhui Gao, Qi Zhang, Tao Gui, and Xuanjing Huang. Llama
634 beyond English: An empirical study on language capability transfer. arXiv:2401.01055,
635 2024.

A Survey of Prior Work and Comparison with Our Analysis

We systematically surveyed prior works on non-English LLM development in two perspectives: coverage of design choices and effects of training on target languages.

At the first glance on Table 1, we can find that several task types are covered sparsely. Only 0–3 papers address machine translation (in either direction), code generation, or summarization—indicating that these areas remain largely unexplored in the literature.

More importantly, we observed the contradictory evidence for “language-agnostic” tasks. The majority of prior studies actually report gains from target-language training on commonsense reasoning (8 positive, 1 neutral, 0 negative) and academic subject benchmarks (5 positive, 2 neutral, 0 negative). These findings contrast both with our results. Furthermore, there seems no clear consensus on other tasks. For reading comprehension and mathematical reasoning benchmarks, prior work offers mixed or inconclusive evidence regarding the impact of target-language data (6 positive, 3 neutral, 0 negative; 3 positive, 2 neutral, 1 negative, respectively).

Table 1: The impact of training on the target language text. ↗: Positive, ↘: Negative, →: Neutral, –: Not investigated

Reference	Lang	Method	Read-Compr.	Common-Sense Reason.	Math-Reason.	MT to Tgt Lang	MT from Tgt Lang	Acad. Subject	Coding	Local Knowl. QA	Summarization
Ours	JA	PT CPT	→	→	→	↗	↗	→	→	↗	↗
Etzaniz et al. (2024b)	EU	CPT	↗	↗	↗, →	–	–	↗	–	↗	–
Fujii et al. (2024)	JA	CPT	↗	↗	↗	↗	↘	–	–	↗	↘
Phasook et al. (2024)	TH	CPT	→	↗	↗	–	–	↗, →	–	–	–
Dou et al. (2024)	VI TH ID MS LO	CPT	↗	↗	–	–	–	↗, →	–	–	–
Joshi et al. (2025)	HI	CPT	↗	↗	–	–	–	↗	–	↗	↗
Vo et al. (2024)	KO	CPT	↗	↗	–	–	–	–	–	–	–
Choi et al. (2024)	KO	CPT	→	↗	–	–	–	–	–	–	–
Toraman (2024)	TR	CPT	→	→	–	–	–	–	–	–	–
Larcher et al. (2023)	PT TA	CPT	↗	–	–	–	–	–	–	–	–
Tejaswi et al. (2024b)	HI AR TR	CPT	–	↗	–	↗	–	–	–	–	↗
Cui et al. (2024)	ZH	CPT	–	–	–	–	–	↗	–	–	–
Etzaniz et al. (2024a)	EU	CPT	–	–	–	–	–	–	–	↗	–
Holmström et al. (2023)	SV	PT	–	–	↘, →	–	–	–	–	–	–
Luukkonen et al. (2023)	FI	CPT	–	–	–	–	–	–	–	–	–
Pipatanakul et al. (2023)	TH	CPT	–	–	–	–	–	–	–	–	–

B Details of the Experimental Setup

B.1 Evaluated Models

Table 2 shows a list of LLMs evaluated in this study. The table includes the name, the number of active parameters during inference, the base model from which the model was continually pre-trained, the language distribution of the training corpus, the total number of training tokens, the reported or estimated number of training tokens in English and Japanese, and the reference of each model. § B.3 explains the method used to estimate the number of language-specific training tokens. CPT stands for *continual pre-training*.

¹⁰Number of active parameters on inference. The total number of parameters is 47B.

Table 2: List of evaluated LLMs (the number of tokens is in billions [Bil], including estimates).

Model name	Num of params	Construction method	Source of CPT	Corpus	Training tokens	EN tokens	JA tokens	Reference
Yi-1.5 6B	6	PT	—	ZH,EN, Code	3600	2170	—	AI et al. (2024)
CyberAgentLM2-7B	7	PT	—	JA,EN	1300	650	650	cyberagent/calm2-7b
Japanese Stable LM Base Gamma 7B	7	CPT	Mistral-7B-v0.1	JA,EN	—	—	100	stabilityai/japanese-stablelm-base-gamma-7b
Japanese StableLM Beta 7B	7	CPT	Llama2 7B	JA,EN	2100	1794	102	stabilityai/japanese-stablelm-base-beta-7b
Llama 2 7B	7	PT	—	EN	2000	1794	2	Touvron et al. (2023)
Mistral-7B-v0.1	7	PT	—	EN	—	—	—	Jiang et al. (2023)
Mistral-7B-v0.2	7	PT	—	EN	—	—	—	Jiang et al. (2023)
Qwen1.5-7B	7	PT	—	—	4000	2000	—	Team (2024)
Qwen2-7B	7	PT	—	ZH,EN, Code+27	7000	3500	—	Yang et al. (2024)
RakutenAI-7B	7	CPT	Mistral-7B-v0.1	JA,EN	—	—	175	RakutenGroup et al. (2024)
Sarashina2-7B	7	PT	—	JA,EN	2100	840	1050	sbintuitions/sarashina2-7b
Swallow 7B	7	CPT	Llama2 7B	JA,EN	2100	1794	102	Fujii et al. (2024)
Swallow-MS v0.1	7	CPT	Mistral-7B-v0.1	JA,EN, Code	—	—	100	Fujii et al. (2024)
Youri 7B	7	CPT	Llama2 7B	JA,EN	2040	1834	42	Sawada et al. (2024)
Llama 3 8B	8	PT	—	EN	15000	14250	15	Dubey et al. (2024)
Llama 3 Swallow 8B	8	CPT	Llama3 8B	JA,EN, Code	15100	14250	115	Fujii et al. (2024)
Llama 3 Youko 8B	8	CPT	Llama3 8B	JA,EN	15022	14250	37	Sawada et al. (2024)
Yi-1.5 9B	9	PT	—	ZH,EN, Code	3100	2170	—	AI et al. (2024)
ELYZA-japanese-Llama-2-13b	13	CPT	Llama2 13B	JA	2018	1794	20	Sasaki et al. (2023)
Fugaku-LLM 13B	13	PT	—	JA,EN	400	200	200	Fugaku-LLM/Fugaku-LLM-13B
Llama 2 13B	13	PT	—	EN	2000	1794	2	Touvron et al. (2023)
LLM-jp-13B v2.0	13	PT	—	JA,EN, Code	260	120	130	LLM-jp et al. (2024)
Sarashina2-13B	13	PT	—	JA,EN	2100	840	1050	sbintuitions/sarashina2-13b
Swallow 13B	13	CPT	Llama2 13B	JA,EN	2100	1794	102	Fujii et al. (2024)
Yi-1.5 34B	34	PT	—	ZH,EN, Code	3100	2170	—	AI et al. (2024)
C4AI Command-R v0.1	35	PT	—	JA,EN, ZH+8	—	—	—	CohereForAI/c4ai-command-r-v01
Mixtral-8x7B-v0.1	13 ¹⁰	PT	—	EN	—	—	—	Jiang et al. (2024)
Swallow-MX 8x7B v0.1	13 ¹⁰	CPT	Mixtral-8x7B-Instruct-v0.1	JA,EN	—	—	100	Fujii et al. (2024)
Japanese Stable LM Beta 70B	70	CPT	Llama2 70B	JA,EN	2100	1794	102	stabilityai/japanese-stablelm-base-beta-70b
KARAKURI LM 70B v0.1	70	CPT	Llama2 70B	JA,EN	2016	1794	18	KARAKURI Inc. (2024)
Llama 2 70B	70	PT	—	EN	2000	1794	2	Touvron et al. (2023)
Llama 3 70B	70	PT	—	EN	15000	14250	15	Dubey et al. (2024)
Llama 3 Swallow 70B	70	CPT	Llama3 70B	JA,EN, Code	15100	14250	115	Fujii et al. (2024)
Swallow 70B	70	CPT	Llama2 70B	JA,EN	2100	1794	102	Fujii et al. (2024)
Qwen2-72B	72	PT	—	ZH,EN, Code+27	7000	3500	—	Yang et al. (2024)

B.2 Evaluation Tasks and Benchmarks

Table 3 provides an overview of the evaluation benchmarks used in this study. The table includes the benchmark name, a brief description, the language of the task, the metric for scoring the model’s output, the experimental setting (e.g., few-shot, zero-shot, chain-of-thought), and the reference of each benchmark. The scale of evaluation metrics is normalized between 0 and 1, and *EM* means *exact match*.

Table 3: List of benchmarks used for evaluation.

Name	Description	Lang.	Eval. metric ^{9,10}	Exp. setup	Reference
JcommonsenseQA (JCom.)	Multiple-choice questions with 5 options based on a knowledge base	JA	Acc.	4-shot	Kurihara et al. (2022)
JEMHopQA	Free-form question answering to evaluate knowledge and reasoning ability	JA	Char F1	4-shot	Ishii et al. (2023)
NIILC	Free-form question answering where answers can be obtained from an encyclopedia	JA	Char F1	4-shot	Sekine (2003)
JSQuAD	Free-form question answering on Wikipedia articles	JA	Char F1	4-shot	Kurihara et al. (2022)
XL-Sum	Generating summaries from BBC articles	JA	ROUGE-2	1-shot	Hasan et al. (2021)
MGSM	Japanese translation of the primary school math word problem dataset (GSM8K)	JA	Acc. (EM)	4-shot	Shi et al. (2023)
WMT20(en-ja)	English-Japanese translation of news articles	JA	BLEU	4-shot	Barrault et al. (2020)
WMT20(ja-en)	Japanese-to-English translation of news articles	JA	BLEU	4-shot	Barrault et al. (2020)
JMMLU	Japanese translation of the multiple-choice benchmark MMLU (53 subjects)	JA	Acc.	5-shot	Yin et al. (2024)
JHumanEval	Japanese translation of HumanEval	JA	pass@1	0-shot 10 trials	Sato et al. (2024)
OpenBookQA	Multiple-choice questions based on scientific knowledge and common sense	EN	Acc.	4-shot	Mihaylov et al. (2018)
TriviaQA	Free-form question answering based on trivia knowledge	EN	Acc. (EM)	4-shot	Joshi et al. (2017)
HellaSwag	Multiple-choice questions to predict the next event	EN	Acc.	4-shot	Zellers et al. (2019)
SQuAD2	Free-form question answering based on a supporting document	EN	Acc. (EM)	4-shot	Rajpurkar et al. (2018)
XWINO	Binary-choice questions to identify the antecedent of a pronoun in a sentence	EN	Acc.	4-shot	Tikhonov & Ryabinin (2021)
MMLU	Multiple-choice questions across 57 subjects	EN	Acc.	5-shot	Hendrycks et al. (2021)
GSM8K	Primary school math word problem dataset	EN	Acc. (EM)	4-shot	Cobbe et al. (2021)
BBH	23 challenging tasks from the BIG-Bench dataset	EN	Acc. (EM)	3-shot CoT	Suzgun et al. (2023)
HumanEval	Evaluation of code generation ability via unit tests	EN	pass@1	0-shot 10 trials	Chen et al. (2021)

B.3 Estimating the Number of Training Tokens

The numbers of language-specific training tokens (in billions) were either obtained from or calculated based on official sources such as technical reports, release documents, or model cards. When an exact number was unavailable in the source, we used the following estimates:

- Ratio of Japanese training tokens:
 - Llama 2, Llama 3: 0.1%
 - Mistral, Mixtral: 0%
 - Full-scratch Japanese LLMs: 50%
 - Japanese LLMs with CPT: 100%
- Ratio of English training tokens:
 - Qwen1.5, Qwen2: 50%
 - Yi-1.5: 70%
 - Llama 2: 89.7%
 - Llama 3: 95%

Table 4: List of evaluation frameworks.

Name	Description	Reference
LLM-jp eval (1.3.0)	Automatic evaluation tool for Japanese LLMs	Han et al. (2024)
JP Language Model Evaluation Harness (commit #9b42d41)	An evaluation framework for Japanese LLMs	zenodo.10256836
Language Model Evaluation Harness (0.4.2)	An evaluation framework for LLMs	zenodo.10256836
Code Generation LM Evaluation Harness (commit #0261c52)	An evaluation framework for code generation task	Ben Allal et al. (2022)

Table 5: Breakdown of LLM groups used in Figure 1.

Category	Models
Japanese LLMs pre-trained from scratch	CyberAgentLM2-7B, Sarashina2-7B, Sarashina2-13B, Fugaku-LLM 13B, LLM-jp-13B v2.0
LLMs continually pre-trained on Japanese text	Japanese Stable LM Base Gamma 7B, Japanese Stable LM Beta 7B, RakutenAI-7B, Swallow 7B, Swallow-MS v0.1, Youri 7B, Llama 3 Swallow 8B, Llama 3 Youko 8B, ELYZA-japanese-Llama-2-13b, Swallow 13B, Swallow-MX 8x7B v0.1, Japanese Stable LM Beta 70B, KARAKURI LM 70B v0.1, Llama 3 Swallow 70B, Swallow 70B
English LLMs	Yi-1.5 6B, Llama 2 7B, Mistral-7B-v0.1, Mistral-7B-v0.2, Llama 3 8B, Yi-1.5 9B, Llama 2 13B, Yi-1.5 34B, Mixtral-8x7B-v0.1, Llama 2 70B, Llama 3 70B
Multilingual LLMs	C4AI Command-R v0.1, Qwen1.5-7B, Qwen2-7B, Qwen2-72B

A symbol ‘-’ in Table 2 indicates that the number could not be obtained or estimated despite our best efforts. We excluded these LLMs from the analysis of the scaling laws in § 4.4.

B.4 Evaluation Framework

Table 4 reports a list of evaluation frameworks used in this study. The table shows the framework name, a brief description, and the reference of the framework. We slightly customized these evaluation frameworks to cover benchmarks that are not officially supported and to implement workarounds for LLMs; for example, some LLMs require special tokens or line breaks in the prompt to generate valid outputs. We will release the customized implementation upon acceptance.

B.5 Details of LLM Grouping

Table 5 shows the breakdown of LLM groups used in Figure 1.

C Analysis of the Evaluation Results

This section presents detailed observations that complement the explanation in § 4.1.

C.1 Performance Difference between the Pre-trained Languages

Figure 1 reveals a notable observation: the scores of Japanese LLMs pre-trained from scratch (the blue box) are consistently lower than those of continually pre-trained models. This

may be due to the relatively small number of parameters of the LLMs in this category (e.g. CyberAgentLM2-7B, Sarashina2-7B, Fugaku-LLM 13B), as well as the limited training budget (i.e., number of training tokens) available for developing LLMs from scratch. This highlights a challenge in developing local LLMs in Japan.

Additionally, compared to other groups, multilingual LLMs (the black box) performed significantly better in arithmetic reasoning (MGSM and GSM8K) and code generation (JHumanEval and HumanEval) tasks. However, we believe that this does not reflect the overall strength of multilingual LLMs, but rather the strengths of Qwen family (Yang et al., 2024), which represents three out of four LLMs in this group.

C.2 Variations in Task Scores

Figure 1 highlights tasks with both high and low score variances. Tasks with low score variances can be grouped into two categories:

1. Benchmarks evaluated with n-gram based metrics (e.g. WMT20-ja-en and WMT20-en-ja with BLEU, and XL-Sum with ROUGE-2).
2. Tasks requiring essential skills (e.g. JSQuAD and SQuAD2.0 (reading comprehension), and OpenBookQA and XWINO (commonsense)).

In contrast, tasks with high score variances can be grouped into two categories:

1. Tasks requiring specific capabilities (e.g. MGSM, GSM8K (arithmetic reasoning), JHumanEval and HumanEval (code generation))
2. Knowledge-intensive tasks (e.g. NIILC, JMMLU, MMLU, and TriviaQA)

The scores for these tasks heavily depend on whether a model possesses the necessary capabilities or specialized knowledge, which leads to a greater variance.

D Robustness Check of Findings Obtained from Experimental Results

To test the robustness of the findings presented in § 4, we conducted two additional analyses using different methods and settings: the use of maximum likelihood estimation and Promax rotation¹¹ instead of PCA (in § 4.3); and exclusion of continually pre-trained models to focus on models trained from scratch. Moreover, we performed leave-one-out cross-validation to confirm that our insights derived from observational approach are robust to statistical errors.

D.1 Maximum Likelihood Estimation and Promax Rotation

Figure 10 presents factor loadings with Promax rotation applied. This figure reveals two factors similar to those identified in § 4.3: ability factor for arithmetic reasoning and code generation (Factor 2 for PC3), and ability factor Japanese (Factor 3 for PC2). In contrast, the first factor (Factor 1) seems to represent English ability, not the general ability (PC1), since the loading scores are strongly positive on the English task benchmarks such as OpenBookQA, TriviaQA, HellaSwag, and XWINO.

Additionally, the fourth factor (Factor 4) seems to be a distinct ability factor for Japanese at first glance since the loading scores are strongly positive on two Japanese task benchmarks (JCom. and JSQuAD). However, the correlation coefficient with the logarithm of the computational budget for Japanese is as small as 0.241, much lower than that of the computational budget for English (0.788). Figure 9 shows small Factor 4 scores on Japanese LLMs, such as Llama 3 Youko 8B, Japanese Stable LM Beta 7B, CyberAgentLM2-7B, LLM-jp-13B v2.0 and Fugaku-LLM 13B. Even strong Japanese LLMs (e.g., Llama 3 Swallow 70B, Japanese Stable LM Base Gamma 7B) do not show high scores compared to non-Japanese LLMs. Therefore,

¹¹We used the `factor_analyzer.FactorAnalyzer()` and `factor_analyzer.Rotator()` method from the `factor_analyzer` package.

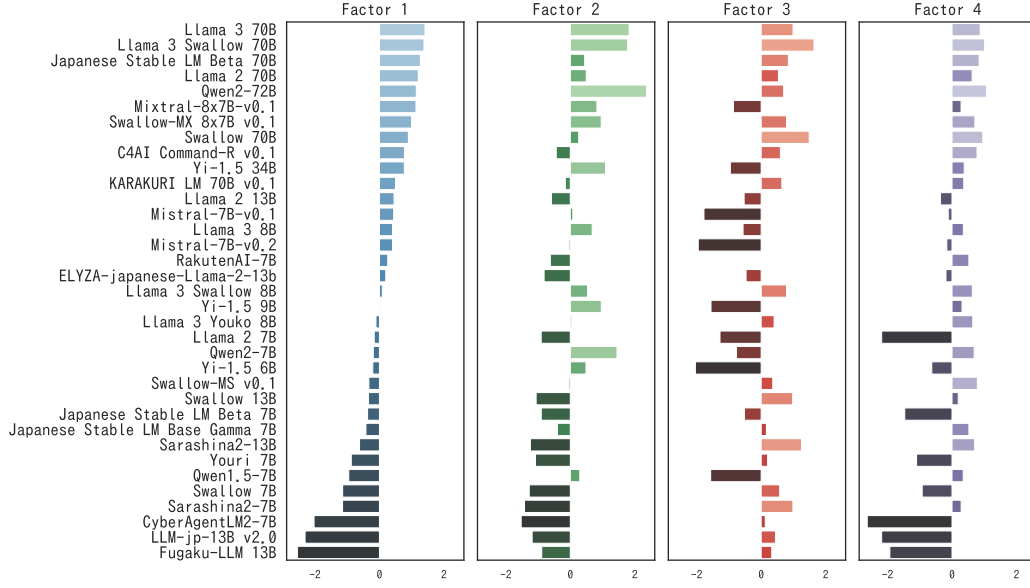


Figure 9: Factor scores for each model with Promax rotation applied.

Factor 1 ($r=23.2$)	0.12	0.03	0.07	0.15	0.26	-0.10	-0.09	0.26	0.09	0.02	0.92	1.08	0.89	0.04	1.00	0.29	0.06	0.58	-0.02
Factor 2 ($r=21.1$)	-0.02	0.03	-0.07	0.17	0.02	0.63	-0.01	0.02	0.39	1.13	0.09	0.10	0.03	0.21	-0.20	0.27	0.63	0.32	1.11
Factor 3 ($r=15.4$)	-0.02	0.44	1.01	0.10	-0.01	0.15	0.90	0.09	0.09	0.01	0.00	0.04	-0.06	0.02	-0.22	-0.03	-0.12	0.12	-0.07
Factor 4 ($r=10.8$)	0.91	-0.05	0.06	0.73	-0.02	-0.08	-0.01	0.03	0.20	-0.00	-0.08	-0.04	0.01	0.06	0.14	0.20	-0.02	-0.09	0.04
	JCom.	JEMHopQA	NIILC	JSQuAD	XL-Sum	MGSM	WMT20-en-ja	WMT20-ja-en	JMMLU	JHumanEval	OpenBookQA	TriviaQA	HeLLaSwag	SQuAD2.0	XWINO	MMLU	GSM8K	BBH	HumanEval

Figure 10: Factor loadings by task with Promax rotation applied ($n = 35$; r represents a contribution; blue and black colors correspond to Japanese and English task benchmarks, respectively).

the fourth factor should be considered as a residual that is difficult to interpret; therefore, commonsense tasks and reading comprehension do not determine Japanese abilities.

To sum, these results confirm two similar factors to those identified in § 4.3 (an ability factor for arithmetic reasoning and code generation, and a Japanese ability factor) and two unique factors (an English ability factor and a residual factor).

D.2 Analysis with only Full-scratch Models

We removed continually pre-trained LLMs, which are categorized as *LLMs continually pre-trained on Japanese text* in Table 5 and conducted the same analysis as in § 4.2 to § 4.4.

Figure 15 shows the Pearson correlation matrix of benchmark scores. The figure reveals that JEMHopQA, NIILC (QA about Japanese knowledge) and WMT20-en-ja (English-Japanese translation) are weakly correlated with other tasks. In addition, the figure shows strong correlations across languages in benchmarks of arithmetic reasoning (GSM8K vs. MGSM), academic subjects (MMLU vs. JMMLU), and code generation (HumanEval vs. JHumanEval). These findings are consistent with those identified with continually pre-trained LLMs in § 4.2.

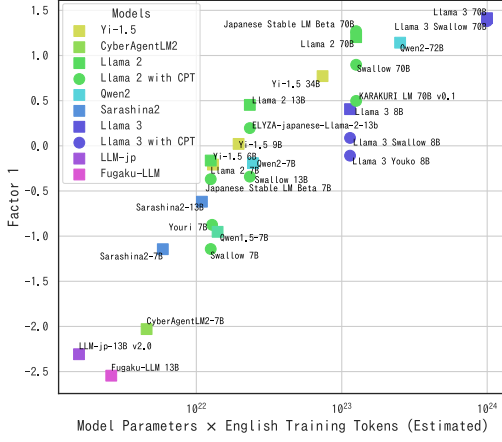


Figure 11: Relationship between the computational budget for English and Factor 1 ($n = 27$).

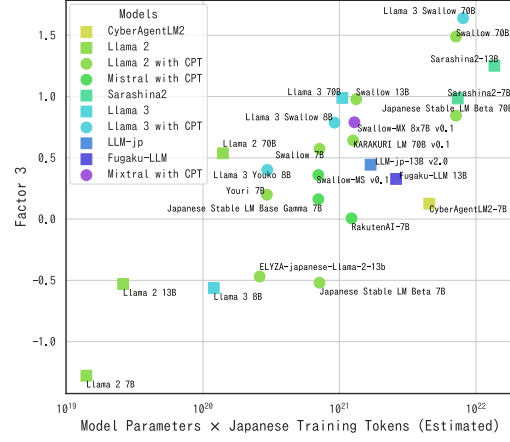


Figure 12: Relationship between the computational budget for Japanese and Factor 3 ($n = 27$).

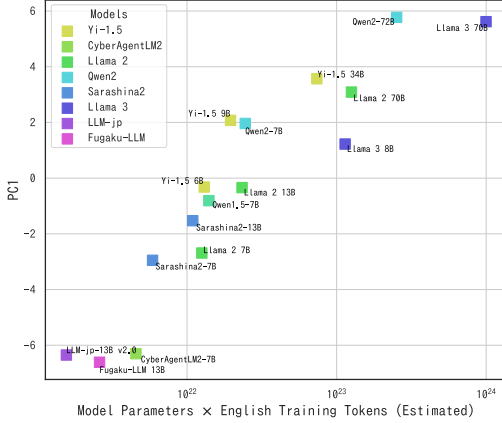


Figure 13: Relationship between the computational budget for English and PC1 ($n = 16$; only with models trained from scratch).

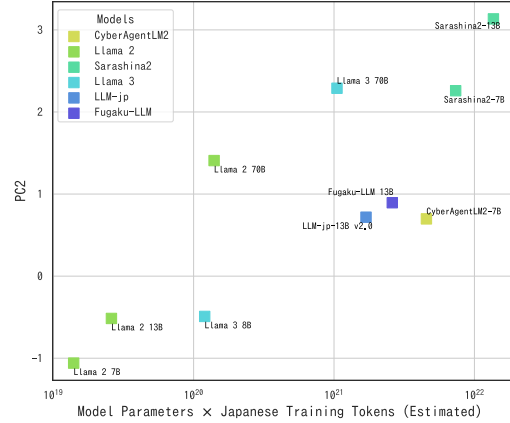


Figure 14: Relationship between the computational budget for Japanese and PC2 ($n = 10$; only with models trained from scratch).

Figure 16 shows the factor loadings for each task benchmark. The figure highlights four factors: a general ability factor with uniform scores on each benchmark (PC1); a Japanese ability factor with high scores on JEMHopQA, NIILC, and WMT20-en-ja (PC2); an ability factor for arithmetic reasoning and code generation with high scores on HumanEval, JHumanEval, MSGM, and GSM8K (PC3); and a residual factor that is difficult to interpret (PC4). These observations are consistent with those obtained with continually pre-trained LLMs in § 4.3.

Lastly, we examined the relationship between the computational budget for English and PC1 (Figure 13) and the one between the computational budget for Japanese and PC2 (Figure 14). Figure 13 exhibits a strong positive correlation between PC1 (general ability) and computational budget for English ($\rho = 0.923$), and Figure 14 indicates a moderate positive correlation between PC2 (Japanese ability) and computation budget for Japanese ($\rho = 0.779$). These relationships are the same as those confirmed with continually pre-trained LLMs in § 4.4.

In this way, we could confirm the findings observed in § 4.2 to § 4.4 even with the LLMs built from scratch, which indicates the robustness of the findings against the construction methods of LLMs.

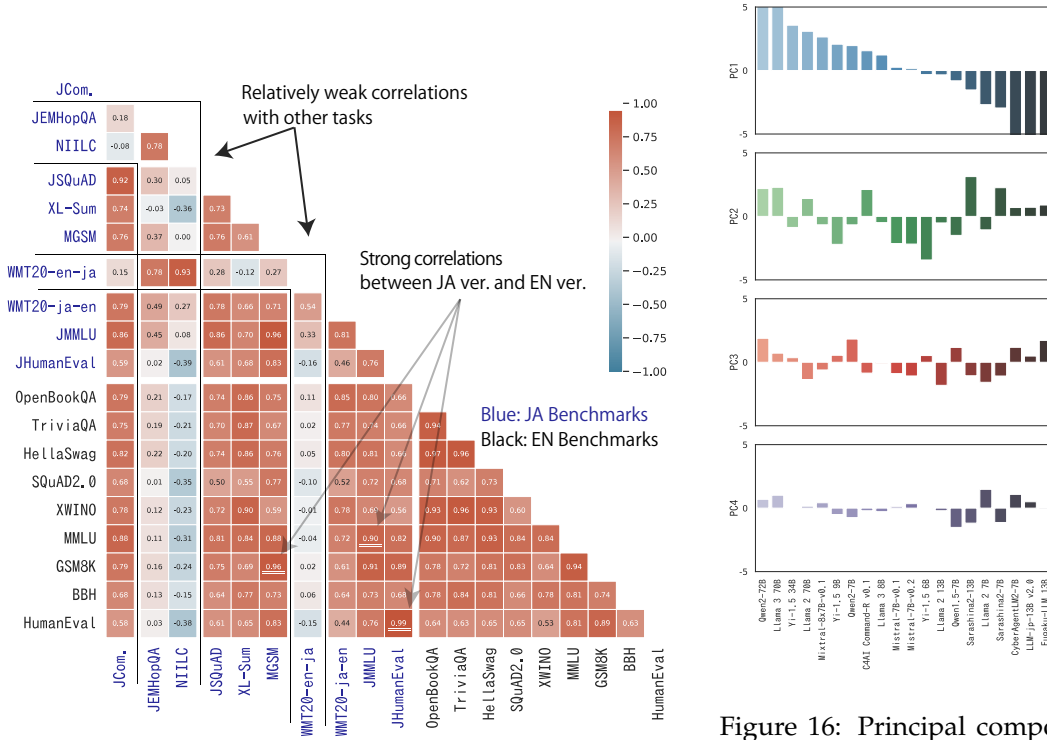


Figure 15: Pearson correlation matrix among benchmark scores ($n = 20$; only with models trained from scratch).

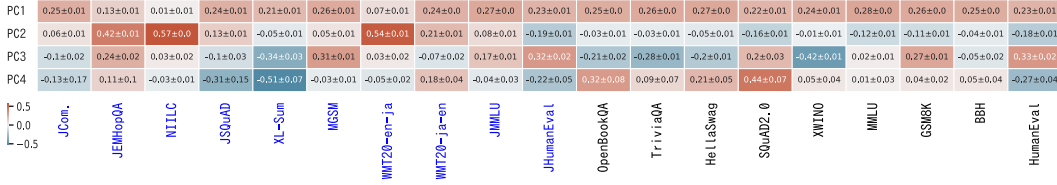


Figure 17: Leave-One-Out CV statistics: mean and standard deviations of the factor loadings ($n = 35$, blue: Japanese benchmarks, black: English benchmarks).

771 D.3 Leave-One-Out Cross-Validation

772 We assessed the statistical error of factor loadings using leave-one-out cross-validation on
 773 the analyzed LLMs (see Figure 17) and confirmed that the standard deviations were small
 774 relative to the absolute values.

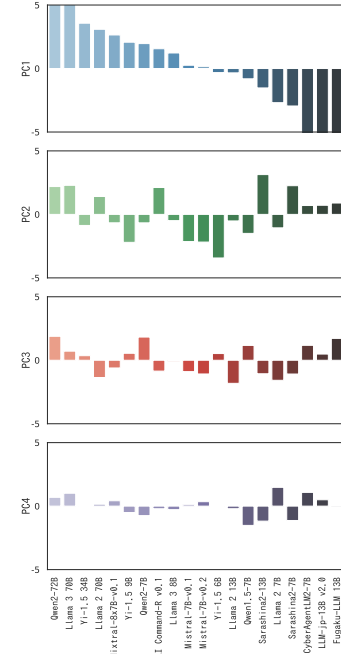


Figure 16: Principal component scores for each model ($n = 20$; only with models trained from scratch).